

Dan Tufiş, Florin Gh. Filip (coordonatori)

**Limba Română
în
Societatea Informațională - Societatea Cunoașterii**



ACADEMIA ROMÂNĂ

Secția de Știința și Tehnologia Informației
Institutul de Cercetări pentru Inteligență Artificială

Limba Română în Societatea Informațională - Societatea Cunoașterii

Coordonatori: Dan TUFIȘ, Florin Gh. FILIP

Colecția
Societatea Informațională
Coordonator:
Prof. dr. ing. Doina BANCIU

ACADEMIA ROMÂNĂ

Secția de Știința și Tehnologia Informației
Institutul de Cercetări pentru Inteligență Artificială

Limba Romana

111

Societatea Informațională Societatea Cunoașterii

Coordonatori:
Dan TUFIȘ, Florin Gh. FILIP

ACADEMIA ROMÂNĂ

București, decembrie 2002

Volumul de față a fost produs de
Institutul de Cercetări pentru Inteligență Artificială (RACAI) al Academiei Române
în cadrul Proiectului "Strategii și soluții pentru Societatea Informațională -
Societatea Cunoașterii în România"
din Programul Național de Cercetare-Dezvoltare "INFOSOC",
condus de Institutul Național pentru Cercetare-Dezvoltare în informatică (ICI),
București

EDITURA

DEDICAȚIE

*Acest volum este dedicat Academicianului Mihai Drăgănescu, Profesor și mentorul unei întregi generații de specialiști în știința și tehnologia informației în general și al problemelor societății informaționale și a cunoașterii în special. Marea majoritate a contribuțiilor din acest volum aparțin unor experți ce fac parte din Comisia de Informatizare a Limbii Române, comisie a Academiei Române la carei naștere un rol esențial l-a avut Profesorul Drăgănescu, președintele Secției Știința și Tehnologia Informației. Savantul Mihai Drăgănescu are numeroase contribuții în știința contemporană, binecunoscute atât în țară cât și în străinătate. Pentru cine îl cunoaște pare incredibilă puterea sa de muncă, debordantă creativitate și neostoita căutare a noului. Profesorul Drăgănescu este indiscutabil port-drapelul conceptului de societate informațională-societate a cunoașterii în România. În lucrările sale din urmă cu peste 25-30 de ani se regăsesc cu claritate multe concepte foarte actuale în zilele noastre, previziuni curajoase atunci, acum realități cotidiene. În lucrările domniei sale din ultima vreme, apare un nou concept ce avem convingerea că se va impune: Societatea Conștiinței, o treaptă superioară a societății cunoașterii. Nu este de mirare deci că în contextul societății informaționale și a cunoașterii profesorul Drăgănescu a susținut cu consecvență a afirmat cu claritate rolul Inteligenței Artificiale în devenirea noilor societăți a cunoașterii. Între domeniile Inteligenței Artificiale un loc de frunte în promovarea principiilor societății cunoașterii îi revine Tehnologiei Limbajului Natural. Profesorul Drăgănescu a fost unul dintre pușinii oameni de știință români care au înțeles și au sprijin total aceste direcții. Cu aproape douăzeci de ani în urmă (1983), Profesorul Drăgănescu edita (împreună cu Adrian Davidoviciu și Ioan Georgescu) volumul "Inteligența Artificială și Robotica" pentru ca trei ani mai târziu (împreună cu Corneliu Burileanu) să editeze un alt volum de referință "Analiza și sinteza semnalului vocal". Astăzi, cercetările mondiale în domeniul tehnologiilor lingvistice au atins un nivel de maturitate ce permit sinergizarea eforturilor lingviștilor, informaticienilor, matematicienilor și a altor specialiști din sectorul academic sau industrial, să abordeze proiecte mari, interdisciplinare având ca obiectiv prelucrarea automată, în mediile de comunicare electronică, a din ce în ce mai multe*limbi naturale. Printre acestea, limba română își face loc încet dar sigur. Volumul de față este o mărturie în acest sens. În același timp, volumul constituie într-o nouă confirmare a realităților pe care Profesorul Mihai Drăgănescu le prefigura cu mulți ani în urmă.*

Coediție



EDITURA
Expert

București, România

Editor și coordonare editorială: **Valeriu IOAN-FRANC**

Redactori: **Mircea FAȚĂ, Paula NEACȘU, Irina STĂNESCU**

Concepția grafică, machetare și tehnoredactare: **Lumița LOGIN**

Coperta: **Nicolae LOGIN**

Toate drepturile asupra acestei ediții aparțin Academiei Române. Reproducerea în totalitate și parțială și pe orice suport, este interzisă fără acordul prealabil al editorului, fiind supusă prevederilor legii drepturilor de autor.

ISBN 973-8177-83-9

Apărut 2002—

Dr. Dan Tufiș, m.c.A.R, Acad. Florin Gh. Fi

CUPRINS

INTRODUCERE.....	9
------------------	---

SECȚIUNEA I:

LINGVISTICĂ TEORETICĂ ȘI FORMALĂ; TERMINOLOGIE

Resurse lingvistice pentru limba română elaborate la Institutul de Lingvistică "Iorgu Iordan" - <i>Ioana Vintilă-Rădulescu</i>	19
Contribuția lingvisticii la studiul terminologiilor științifice - <i>Angela Bidu-Vrănceanu</i>	33
Gramaticile generative nontransformaționale - <i>Emil Ionescu</i>	39
Către o teorie X-bar funcțională - <i>Neculai Curteanu</i>	51
Teoria HPSG. Studiu de caz: acordul încrucișat - <i>Ana-Maria Barbu</i>	87
După 10 ani de experiență terminografică: noul model de date terminologice al TermRom - <i>Dan Matei</i>	109
Probleme de reprezentare a datelor terminografice într-o bază de date relațională - <i>Sorin Ghețaru</i>	121

SECȚIUNEA II:

TEHNOLOGII ALE LIMBAJULUI SCRIS

RO - B A L K A N E T-ontologie lexicalizată, în context multilingv, pentru limba română - <i>Dan Tufiș, Dan Cristea</i>	137
Algoritmi de segmentare a textului în unități de tip clauzal - <i>Dan Gălea, Neculai Curteanu, Cristian Linteș</i>	165
O metodă automată pentru inserarea diacriticelor în texte în limba română - <i>Rada F. Mihalcea, Vivi A. Năstase</i>	191
Contribuții privind structura statistică de cuvinte în limba română scrisă - <i>Adriana Vlad, Adrian Mitrea</i>	207
Dezambiguizarea automată a cuvintelor din corpusuri paralele folosind echivalenții de traducere - <i>Dan Tufiș</i>	235

Referențialitate și cursivitate în relație cu structura de discurs - <i>Dan Cristea</i>	269
DLIR - un sistem de căutare documentară multilingv - <i>Amalia Todirașcu</i>	303
Mediu hermenofor pentru asistarea învățării unor concepte dintr-o limbă străină - <i>Ștefan Trăușan-Matu</i>	317

SECȚIUNEA III: TEHNOLOGII ALE LIMBAJULUI VORBIT

Experimente în vederea recunoașterii vorbitorului - <i>Corneliu Burileanu</i> , <i>Luigi Bojan</i>	3 3 5
Prelucrarea inițială a textului de intrare în cadrul unui sistem de sinteză a vorbirii pornind de la text în limba română - <i>Dragoș Burileanu</i>	359
Utilizarea tehnicilor nuanțate (fuzzy) și de dinamică neliniară pentru sinteza adaptivă a vorbirii - <i>Horia-Nicolai L. Teodorescu</i>	381
Dicționarele multimedia ale limbii române. Secvențe de implementări și experimentări - <i>Dumitru Todoroi, Diana Micusa, Zinaida Todoroi</i> , <i>Ion Lingă, Ion Covalenco, Nicolae Objeleanu, Ștefan Spătaru, Stela</i> <i>Lungu, Virginia Țurcanu, Elena Cozlov, Nadejda Ambrozii, Victor</i> <i>Slobodeanu, Igor Coșeru, Cătălina Suruceanu</i>	401
Mediu pentru editarea transcrierilor fonetice în limba română. Realizarea atlasului lingvistic român pe regiuni - <i>Silviu Bejinariu, Vasile Apopei, Mariana Roman</i>	423

SECȚIUNEA IV: DEZBATERI ȘI DISCUȚII

Asupra a doi vectori funcționali ai societății cunoașterii: managementul cunoașterii și învățarea electronică. Cultura și societatea cunoașterii - <i>Mihai Drăgănescu</i>	441
între lingvistica matematică și cea computațională - <i>Solomon Marcus</i>	471
între lingvistica matematică și cea computațională: o altă perspectivă - <i>Dan Tufiș</i>	481

INTRODUCERE

Programul de cercetare aplicativă "Strategii și soluții pentru Societatea Informațională - Societatea Cunoașterii în România (SI-SC), din subprogramul strategic, al Programului Național INFOSOC a avut ca principale obiective stabilirea unui program de veghe conceptuală pentru menținerea pe linia tendințelor mondiale ale avansului SI-SC, sensibilizarea factorilor de decizie și a publicului larg, crearea unui cadru de reflecție prospectivă pe teme prioritare ale SI-SC: economie, sociale, culturale, tehnologice, ambientale, precum și operaționalizarea unor soluții de interes prioritar pe plan național. În cadrul acestui proiect a fost elaborat volumul "Societatea Informațională - Societatea Cunoașterii. Concepte, soluții și strategii pentru România" (publicat la Ed. Expert în anul 2000), realizat sub coordonarea Academicianului Florin Gheorghe Filip. Acest volum avea ca scop construirea unei viziuni și conținea o serie de studii și cercetări care au aprofundat rezultatele programului prioritar al Academiei Române privind *Societatea Informațională - Societatea Cunoașterii* și au identificat o serie de orientări strategice cerute de susținerea unei dezvoltări de tip "salt" a SI-SC în România. Prin prisma obiectivelor proiectului, au fost analizate principalele aspecte conceptuale ale SI-SC, probleme legate de infrastructurile informatice și de comunicații ale SI-SC, formarea profesională și pregătirea generală a populației în și pentru SI-SC, rolul științei în cercetării și inovării, aspecte sociale și juridice, instituțiile statului și relația lor cu cetățeanul, dezvoltarea economiei și afacerilor, dimensiunea culturală a SI-SC, actorii sociali ai creării și difuzării tehnologiei informației și comunicațiilor în contextul SI-SC. Studiile tematice, ancheta Delphi pentru consultarea opiniei experților privind tendințele globale și opțiunile posibile de raportare la ele, scenariile de evoluție elaborate au susținut funcția prospectivă a proiectului.

Funcția operativă a acestui proiect, respectiv identificarea de soluții tehnice privind rezolvarea principalelor priorități identificate în faza analizei prospectivă, urma să se manifeste în perioada imediat următoare, printr-o dintr-o serie de cercetări/dezvoltări tehnologice ce vor trata pe larg problematica specifică fiecăruia dintre direcțiile amintite anterior. Această serie este deschisă pentru prezentul volum ce înglobează contribuții ale unor specialiști români reprezentativi în domeniul prelucrării automate a limbajului natural și a resurselor lingvistice necesare utilizării limbii române în mediile de comunicare electronică.

În [1] este definit conceptul de "Societate Informațională - Societatea Cunoașterii" (SI-SC) precum și principalii săi vectori tehnologici și funcționali. În acest context "internetul dezvoltat" (ca vector tehnologic) și "managementul

utilizării morale a cunoașterii la nivel global" (ca vector funcțional) sunt prezentați ca factori motrici esențiali ai Societății Cunoașterii, și în perspectivă, a Societății Conștiinței. "Din momentul în care intervine Internetul cu marile avantaje pe care acesta le aduce (e-mail, comerț electronic și tranzacții electronice, piața Internet, distribuția de 'conținut') prin cuprinderea în sfera informației electronice a unui număr cât mai mare de cetățeni se trece la societatea informațională. Cunoașterea este informație cu înțeles și informație care acționează. De aceea societatea cunoașterii nu este posibilă decât grefată pe societatea informațională și nu poate fi separată de aceasta. În același timp, ea este mai mult decât societatea informațională prin rolul major care revine informației-cunoaștere în societate." [1]

În 1984, William Gibson, un dizident cognitiv - după cum se auto-caracterizează, publică volumul SF *"Neuromancef"* (Ace Book, July 1984, ISBN: 0-441-56959-5), carte care pe lângă o mulțime de premii literare i-a adus notorietatea și pentru crearea termenului "cyberspace": "the total interconnectedness of human beings through computers and telecommunication without regard to physical geography... A consensual hallucination experienced daily by billions of legitimate operators, in every nation, by children learning mathematical concepts...a graphical representation of data abstracted from the banks of every computer in the human system. Unthinkable complexity. Lines of light ranged in the non-space of the mind. Clusters and constellations of data. Like city lights receding..." (op. cit).

Termenul a făcut carieră, actualmente fiind o noțiune care din punct de vedere tehnic subsuma conceptul "Internet"(scris cu majusculă): "cyberspace: The impression of space and community formed by computers, computer networks, and their users; the virtual "world" that Internet users inhabit when they are online

The term internet (spelled with a lower case T) is distinguished from the Internet (spelled with the "I" capitalized). The Internet refers to a specific, historic, ubiquitous worldwide digital communication network." (cf. Glossary of Telecommunications, American National Standard T1.523-2001, www.atis.org/ta2k/cyberspace.html, 05.08.2002).

Dimensiunea tehnică (evocată mai sus) a noțiunii de *"ciberspațiu"* este complementată de dimensiunea socio-culturală și din această perspectivă de problemele "satului global" previzionat de Societatea Informațională - Societatea Cunoașterii. Ideea atenuării schismei dintre specialiștii din domeniul tehnic și cei din zona științelor umaniste în contextul SI-SC este susținută puternic și de M. Derouzos [5], cel care a propus conceptul de "piața informațională", pe care îl consideră mai realist decât cel de "ciberspațiu". De altfel, dimensiunile socio-culturale ale SI-SC au fost evocate în capitolele 2, 3, 4 și 6 ale volumului *"Societatea Informațională - Societatea Cunoașterii. Concepte, soluții și strategii pentru România"*

Printre componentele socio-culturale ale SI-SC, utilizarea limbii materne în mediile de lucru și comunicare electronice a accesului universal la ciberspațiu [2, 3, 4] constituie priorități absolute.

În contextul actual, al comunicării mediate de tehnologia informației și telecomunicații, limba devine obiect al investigației tehnice. Tehnologia limbajului impune metodologii specifice de cercetare/dezvoltare, dezvoltarea sau adaptarea resurselor lingvistice fundamentale cum ar fi dicționarele, tezaurele, corpusurile gramaticile computerizate, în conformitate cu standardele sau recomandările existente. În funcție de resursele lingvistice disponibile, de volumul și calitatea de compatibilitate a codificării lor în raport cu recomandările și standardele internaționale etc., se poate vorbi de *nivelul de tehnologizare* al unei limbi naturale. Nivelul de tehnologizare al unei limbi naturale este în corespondență directă cu statutul de limbă de *circulație electronică*. Această sintagmă, o parafrază a expresiei *limbă de circulație internațională*, încearcă să elimine antinomia, pe care o cunoscută pe atât de goală în conținut spiritual și cultural, "limbi mari/limbi mici". Conceptul de "limbă de circulație electronică", pe lângă semnificația lui directă, are și profunde implicații culturale, sociale și nu în ultimul rând economice implică dreptul fiecărui cetățean de a avea acces în propria limbă la cunoștințele informaționale și serviciile ciberspațiului.

Promovarea limbii române în SI-SC presupune informatizarea limbii române ca factor infrastructural fundamental (vector funcțional) precum și stimularea utilizării curente (prin vectori tehnologici) a limbii române în utilizarea tehnologiilor și serviciilor informatice. Acest obiectiv presupune eforturi umane și materii substanțiale și de dimensionarea lor se leagă orizontul de timp al realizării sale.

Volumul de față reunește lucrări ce tratează aspecte specifice prelucrării limbajului natural, în marea lor majoritate cu aplecare directă asupra limbii române. Înerent, volumul de față nu poate acoperi întreaga arie problematică a domeniului după cum nici reprezentarea specialiștilor români în domeniul tehnologiei limbajului nu este completă, dar cititorul va găsi un larg evantai de direcții de cercetare, care specialiștii români au obținut rezultate importante.

Volumul este structurat în patru părți (aspecte teoretice și probleme de terminologie, prelucrarea limbajului scris, prelucrarea limbajului vorbit, dezbateri și discuții) care pot fi citite în mod independent, în funcție de interesul specific al cititorului.

Prima parte "Lingvistică teoretică și formală; terminologie" cuprinde lucrări din domeniul lexicografiei, sintaxei și terminologiei.

În lucrarea "Resurse lingvistice elaborate la Institutul de Lingvistică «Iordan»" Ioana Vintilă Rădulescu face o trecere în revistă a celor mai importante resurse lingvistice realizate în cei peste 50 de ani de activitate la Institutul de Lingvistică «Iordan».

Angela Bidu-Vrânceanu prezintă în lucrarea "Contribuția lingvisticii la studiul terminologiilor științifice" concluziile a trei contracte de cercetare științifice

având ca obiect studiul terminologic al limbajului folosit în diverse domenii (matematică, filozofie, mineralogie, arte plastice).

Articolul "Gramaticile nontransformaționale" al lui Emil Ionescu face o prezentare generală a gramaticilor bazate pe unificare și constrângeri precum și a principalelor realizări, în contextul acestei paradigme, în cercetarea lingvistică din România.

Neculai Curteanu propune în lucrarea "Către o teorie X-bar funcțională" o reconsiderare a teoriei clasice X-bar prin perspectiva modelului propriu SCD (Segmentare-Coeziune-Dependență).

Ana-Maria Barbu prezintă în lucrarea sa "Teoria HPSG: studiu de caz: acordul încrucișat" principalele caracteristici ale teoriei HPSG și discută în acest context un caz de dependență încrucișată specific limbii române, respectiv clauzele relative în care pronumele relativ este precedat de articolul genitival.

O serie de probleme legate de terminologia computațională sunt prezentate în ultimele două lucrări ale primei secțiuni. În articolul "După 10 ani de experiență terminologică: noul model de date terminologice al TERMROM" Dan Matei prezintă modelul dezvoltat în conformitate cu noile tendințe și standarde în domeniu și adoptat de Asociația Română de Terminologie - TERMROM.

Lucrarea lui Sorin Gețaru "Probleme de reprezentare a datelor terminografice într-o bază de date relațională" aduce în discuție aspecte specifice reprezentărilor standardizate necesare realizării dezideratului de interschimb și interoperabilitate între diverse tezaure terminologice și discută elementele distinctive ale standardului ISO-12200 MARTIF (Machine-Readable Terminology Interchange Format).

Secțiunea a doua a volumului ("Tehnologii ale limbajului scris") este deschisă de lucrarea lui Dan Tufiș și Dan Cristea "RO-BALKANET - ontologie lexicalizată în context multilingv pentru limba română" care descrie stadiul dezvoltării unui dicționar, pentru limba română, structurat ca o rețea semantică, de tip EuroWordNet, rezultat al unui program european ce-și propune extensia EuroWordnet (în prezent implementat pentru 10 limbi europene) cu încă 5 limbi.

Articolul lui Dan Gâlea, Neculai Curteanu și Cristian Linteș "Algoritmi de segmentare a textului în unități de tip clauzal" tratează o problemă delicată a prelucrării limbajului natural, respectiv cea a identificării, în raport cu un anumit criteriu funcțional, a structurilor "clauzale" și prezintă contrastiv doi algoritmi diferiți (unul dintre ei aparținând autorilor), atât prin prisma modelării lingvistice cât și al performanței computaționale.

Rada Mihalcea și Vivi Năstase prezintă în articolul lor o metodă de inserare automată a caracterelor diacritice în texte scrise (cu studiu de caz pentru

limba română) fără diacritice și comentează rezultatele proprii în comparație cu cele ale altor metode dezvoltate pentru rezolvarea aceleiași probleme.

Adriana Vlad și Adrian Mitrea prezintă în lucrarea lor "Contribuții privind structura statistică de cuvinte în limba română scrisă" rezultate recente privind caracterizarea statistică a limbii române scrise, prin aproximarea ei ca un Markov ergotic multiplu cu ordin de multiplicitate mai mare decât 30, rezultate obținute prin analiza riguroasă a unui corpus foarte mare de texte.

Articolul "Dezambiguizarea semantică automată în corpusuri paralele" al lui Dan Tufiș prezintă o alternativă la spinoasa problemă a dezambiguizării cuvintelor polisemantice, bazându-se pe extragerea cunoștințele implicite existente într-un corpus multilingv (creat de traducători profesioniști) și apelând la tehnici euristici ale lingvisticii corpusului.

Dan Cristea prezintă în articolul "Referențialitate și cursivitate în structura discursului" elementele definitorii ale teoriei sale asupra structurii discursivității textelor (teoria nervurilor) și își exemplifică argumentația prin analiza dihotomă a structură-referențialitate și structură-coerență.

În lucrarea "DLIR - un sistem de căutare documentară multilingv" Amalia Todirașcu prezintă o abordare bazată pe logici terminologice, ontologii și tehnici de prelucrare a corpusurilor în implementarea unui sistem de regăsire documentară bilingv (română și franceză).

Partea a doua a volumului se încheie cu articolul lui Ștefan Trăușan-Matei "Medii hermenofor pentru asistarea învățării unor concepte într-o limbă străină" care după o prezentare a noțiunilor cu care operează în lucrare, descrie un model de prelucrare a metaforelor utilizate în limbaje specializate (studiu de caz: limbaj financiar) incorporat într-un sistem de instruire inteligentă în învățarea conceptelor într-o limbă străină, sistem distribuit dezvoltat în cadrul unui proiect european.

Secțiunea a treia a volumului este dedicată problemelor de prelucrare a vorbirii. Corneliu Burileanu și Luigi Bojan se opresc asupra tehnicilor de recunoaștere a vorbitorului ca etapă distinctă și strict necesară pentru recunoașterea automată a vorbirii și prezintă o parte a rezultatelor obținute de către autori.

Lucrarea lui Dragoș Burileanu "Prelucrarea inițială a textului de intrare în cadrul unui sistem de sinteză a vorbirii pornind de la text în limba română" abordează problemele sintezei limbajului vorbit pornind de la un text în formă electronică și detaliază etapa de preprocesare a textului ca etapă primară a procesului transformării sale în semnal vocal inteligibil și coerent.

Tot în domeniul sintezei vorbirii se plasează și lucrarea lui Horia Nicolae Teodorescu "Utilizarea tehnicilor nuanțate (fuzzy) și de dinamică neliniară per

sinteza adaptivă a vorbirii" ce subliniază rolul esențial al prozodiei și al modelării sale algoritmice în realizarea unor sinteze vocale de calitate, purtătoare de informație emoțională.

Un proiect de anvergură, este prezentat de Dumitru Todoroi, Diana Micusa, Zinaida Todoroi, Ion Lingă, Ion Covalenco, Nicolae Objeleanu, Ștefan Spătaru, Stela Lungu, Virginia Turcanu, Elana Cozlov, Nadejda Ambrozii, Victor Slobodeanu, Igor Coșeru și Cătălina Suruceanu în lucrarea "Dicționarele multimedia ale limbii române. Secvențe de implementări și experimentări".

Secțiunea a treia a volumului se încheie cu lucrarea elaborată de Silviu Bejinariu, Vasile Apopei și Mariana Roman "Mediu pentru editarea transcrierilor fonetice în Limba Română. Realizarea Atlasului Lingvistic Român pe Regiuni" ce prezintă un instrument ce permite realizarea facilă a transcrierilor fonetice într-un limbaj standardizat (IPA), oferă extensii specifice de adnotare fonetică (realizate până acum manual) și prefigurează realizarea variantei computerizate a atlaselor lingvistice românești.

Ultima secțiune a volumului (Dezbateri și discuții) conține trei contribuții. Prima dintre ele, elaborată de Mihai Drăgănescu, "Asupra a doi vectori funcționali ai Societății Cunoașterii: Managementul Cunoașterii și învățarea Electronică. Cultura și Societatea Cunoașterii" reprezintă liantul dintre volumul precedent (*Societatea Informațională - Societatea Cunoașterii. Concepte, soluții și strategii pentru România*, coordonator Fl. Gh. Filip) și volumul de față, rafinând clasificarea din lucrarea anterioară și adâncind o serie de probleme ridicate în [1].

Ultimele două contribuții reprezintă două puncte de vedere asupra problematicii prelucrării limbajului natural, prima poziție "între lingvistica matematică și cea computațională" fiind susținută de Solomon Marcus, iar cea de a doua "între lingvistica matematică și cea computațională: o altă perspectivă" fiind prezentată de Dan Tufiș.

Mulțumiri

Coordonatorii acestui volum, mulțumesc tuturor celor care au participat la realizarea proiectului "Strategii și soluții pentru societatea informațională-societatea cunoașterii în România" derulat cadrul programului național INFOSOC. Mulțumiri speciale se cuvin directorului programului INFOSOC, Profesor Doina Banciu, care a susținut și a manifestat un interes deosebit față de desfășurarea acestui proiect.

Referințe bibliografice

- [1] M. Drăgănescu "Societatea informațională și a cunoașterii. Vectorii societății cunoașterii" în *F.G. Filip (coord.) Societatea Informațională - Societatea Cunoașterii. Concepte, soluții și strategii pentru România*. Academia Română, Editura Expert, ISBN 973-8177-42-1, 2001, pp. 43-112
- [2] *** The Multilingual Information Society, Report of Commission of the European Communities, COM(95) 486/final, Brussels, November 1995.
- [3] *** Multilingualism in an Information Society, International Symposium organized by EC/DGXIII, UNESCO and Ministry of Foreign Affairs of the French Government, Paris 4-6 December 1997.
- [4] *** Promotion and Use of Multilingualism and Universal Access to Cyberspace, UNESCO 31st session, November 2001.
- [5] M. Dertouzos. "What It will Be". Harper Edge. New York, 1997 (trad. în română "Ce va fi", Ed. Tehnică, București, 2000).

Secțiunea I

**LINGVISTICA TEORETICA
SI FORMALĂ;
TERMINOLOGIE**

Resurse lingvistice pentru limba română elaborate la Institutul de Lingvistică "Iorgu Iordan"

Ioana VINTILĂ-RĂDULESCU
Institutul de Lingvistică "Iorgu Iordan - Al. Rosetti"
București, Calea 13 Septembrie 13
e-mail: ioanar@fx.ro

1. Considerații generale

Numind *resursă* în general o "rezervă sau sursă de mijloace (materiale sau spirituale) susceptibile de a fi valorificate într-o împrejurare dată", înțelegem prin *resurse lingvistice* pentru limba română izvoarele fundamentale de informații cu privire la aceasta, stocate convenabil (chiar dacă încă preponderent în manieră tradițională) și care, în calitate de componente ale *culturii* în sensul cel mai larg, sunt susceptibile de a fi valorificate pentru studierea limbii române, precum și în diverse scopuri conexe, inclusiv aplicative, în cadrul *societății informatice* actuale.

Cât privește Institutul de Lingvistică "Iorgu Iordan"², acesta nu mai există formal ca atare, deoarece la începutul anului 2002, printr-o hotărâre de guvern adoptată la propunerea conducerii Academiei Române, s-a produs re-unirea sa și a Institutului de Fonetice și Dialectologie "Al. Rosetti". (Spunem reunire întrucât cercetările de fonetică și de dialectologie formaseră inițial obiectul unui sector, respectiv al unei secții a Institutului de Lingvistică din București al Academiei Române (înființat în 1949), devenită din 1961 centru și apoi institut independent.) Întrucât în 1998 fusese oficializată, tot prin hotărâre de guvern, propunerea celor două institute, aprobată de Prezidiul Academiei, de a-și adăuga fiecare în titulatură numele fostului său director, institutul în cadrul căruia cele două nuclee care au fuzionat acum își continuă de fapt activitatea poartă numele dublu de Institutul de Lingvistică "Iorgu Iordan - Al. Rosetti".

Fără îndoială, cele mai numeroase și mai importante resurse lingvistice pentru limba română s-au realizat la acum fostul Institut de Lingvistică "Iorgu

*** (1975). Dicționarul limbii române (DLR). Serie nouă. Tomul IX, Litera R, *București*, s.v.

² Pentru o imagine de ansamblu asupra activității acestui institut și a istoriei sale v. Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu (coordonatori) (1999). Institutul de Lingvistică "Iorgu Iordan". 50 de ani de existență (1949-1999), *București*.

lordan", înglobând, până în 1961 direct și apoi numai indirect, și contribuția colegilor foneticieni și dialectologi³, precum și, în unele cazuri, în colaborare cu alte institute de specialitate din țară ale Academiei - Institutul de Lingvistică și Istorie Literară "Sextil Pușcariu" din Cluj și Institutul de Filologie Română "Alexandru Philippide" din Iași - și cu cadre didactice de la facultățile de profil mai ales ale Universității din București. Această activitate este continuată și în noul cadru organizatoric de sectoarele fostului institut, pe care în cele ce urmează îl vom numi, pe scurt, *Institutul*.

2. Resurse lexicografice

Dintre resursele lingvistice tradiționale dezvoltate până în prezent de Institut, cele mai importante din punctul de vedere care interesează aici sunt cele **lexicografice - dicționarele** (mono- și bilingve) -, activitatea lexicografică din Institut, începută încă de la înființarea sa, desfășurându-se din 1959 în cadrul unui sector specializat cu acest profil, condus până în 1985 de Mircea Seche, iar de atunci încoace de Ion Dănăilă⁴.

2.1. Dicționare monolingve

2.1.1. Dintre dicționarele românești monolingve se distinge, prin anumite trăsături ale sale, dicționarul "explicativ general academic" intitulat pur și simplu **Dicționarul limbii române** - dar mai cunoscut ca "Dicționarul Academiei" - a cărui realizare se apropie de sfârșit și care va cuprinde o mare parte a "tezaurului" lexical al limbii române - fără a putea și nici a intenționa să includă însă ansamblul cuvintelor românești folosite în toate epocile, în toate regiunile și în toate domeniile⁵. În ciuda marilor sale calități, care sunt bine cunoscute și asupra cărora nu credem deci că mai este nevoie să insistăm aici, acest dicționar prezintă un dezavantaj major din punctul de vedere al utilizării sale ca resursă de bază (pe lângă faptul că nu se prezintă și sub forma unei variante electronice, care nici nu putea fi imaginată până nu de mult) și anume caracterul său fatalmente neunitar,

datorat faptului că a fost elaborat pe parcursul a aproape un secol⁶, de unde mai deosebiri dintre cele două părți ale sale: cea publicată între 1907 și 1949 și conducerea marelui lingvist Sextil Pușcariu și cea care a început să apară în 1965 și a cărei publicare se apropie, în fine, de sfârșit. "Seria veche" a dicționarului academic, desemnat de aceea prin sigla DA, cuprinde literele A-C (inclusiv puținele neologisme scrise acum cu *k*-, iar în DA cu *c/?*-) și F-J complet, iar literele D și L parțial (până la cuvântul de, respectiv *lojniță*), totalizând 3.142 de pagini tipar, format mare, dintre ele lipsind în întregime, după cum se observă, litera R. Această primă jumătate a dicționarului se distinge prin lista de cuvinte, bogată și ales sub aspectul fondului tradițional, prin tratarea amănunțită a semantismului bazată pe numeroase citate, prin dimensiunile și valoarea comentariilor etimologice, precum și prin traducerea sensurilor în limba franceză⁷. Desigur, acestea aveau cum figura în aceste prime volume numeroasele neologisme încetățenite în românește după elaborarea lor, ilustrarea sensurilor prin utilizarea lor de către autori mai noi și în general toate aspectele care sunt rodul evoluției ulterioare a limbii române, al cercetărilor dialectale, etimologice, filologice etc. mai recente și dezvoltării lingvisticii și metodelor ei, în general. Din 1965 dicționarul și-a reînceput apariția, în format asemănător, ca *Serie nouă* (de data aceasta sub o siglă diferită menționată în titlu, DLR), cu litera M, sub conducerea, la început, a lui lordan, Alexandru Graur și Ion Coteanu, iar actualmente a lui Gh. Mihăilă și Marius Sala. Noua serie păstrează, în mare, principiile lui Sextil Pușcariu, dar beneficiază de toate avantajele elaborării sale mai aproape de zilele noastre: ea include modificări și amplificări reflectând evoluția limbii române, a lexicografiei române și a studiului limbii române în ansamblu, precum și a lingvisticii în general, dar mai cuprinde, în schimb, traducerea sensurilor (în anii '60 nefiind considerat oportun acest lucru, deși era util mai ales pentru cunoașterea limbii române către străini, fără a fi, este drept, uzual într-un dicționar monolingv explicativ), secțiunea etimologică a fost redusă, dicționarul păstrându-și însă caracterul istoric (sensurile sunt date în ordinea atestării lor în texte și în alte surse)⁸. Institutul din bucureștean a redactat literele M, N, P, S și Z⁹ și este pe cale de a încheia

³ V., printre altele, Marius Sala (1999). Institutul de Lingvistică "Iorgu Iordan" la 50 de ani. Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit, p. 35-37.

⁷ Pentru o descriere amănunțită a DA v. Mircea Seche (1969). Activitatea lexicografică a Sextil Pușcariu, în Schiță de istorie a lexicografiei române, voi. II, De la 1880 până astăzi, București, p. 42-72.

⁸ V. și Mircea Seche (1969). Seria nouă a Dicționarului academic general în Schiță de istorie a lexicografiei române, voi. II, De la 1880 până astăzi, București, p. 72-79. Iorgu Iordan, Al. Graur, I. Coteanu (red. resp.) et al. (1965-2000). *Dicționarul limbii române (DLR). Serie nouă*, București: T. VI, *Litera M*, 1965-1968 (apărut inițial în fascicule); Partea 1, *Litera N*, 1971; Partea a 2-a, *Litera O*, 1969] VIII, *Litera P*. Partea 1, P-PĂZĂ, 1972; Partea a 2-a, PE-PÎNAR, 1974; Partea a 3-a, PÎNĂ-POGRIBANIE, 1977; Partea a 4-a, POGRIJENIE-PRESIMȚIRE, 1980; Partea a 5-a, PRESIN-PUZZOLANĂ, 1984; *Litera R*, 1975; X. *Litera S*. Partea 1, S-SCLĂBUC, 1986; Partea a 2-a, SCLĂBUC-SEMÎNȚĂRIE, 1987; Partea a 3-a, SEMN-SÎVEICĂ, 1990; Partea a 4-a, SCLĂBUC-SLĂBUC, 1991.

³ Aceștia au produs mai ales "resurse" de un tip specializat, concretizate în principal în atlase lingvistice și în arhiva fonogramică a limbii române, de care nu ne vom ocupa în mod direct aici, dar care, ca și contribuțiile similare ale altor institute, au avut și un aport indirect la resursele fundamentale despre care vorbim, printre izvoarele cărora s-au numărat

⁴ Pentru detalii cu privire la lucrările acestuia v. Ion Dănăilă (1999). Sectorul de lexicologie și lexicografie, în Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit, p. 98-113.

⁵ Ideea, relativ utopică și controversată, a înregistrării și chiar a descrierii semantice a întregului inventar lexical al limbii române (ILEX) din toate timpurile, incluzând atât numele comune, cât și cele proprii (v. Ion Dănăilă (1993). Pentru un inventar general al limbii române, în "Limba română" XLII, nr. 2, p. 61-68), nici nu a început a fi pusă în practică.

reluarea și terminarea literei *D* absentă din prima parte (trei volume); numai primele patru litere elaborate la București însumează 51.847 de cuvinte și variante, totalizând 5.839 p. Institutului din Cluj i-au revenit literele *O*, *R*, *T*, *Ț* (totalizând 2.044 de pagini de tipar), *U* (aflată sub tipar) și, din prima parte, reluarea și terminarea unei părți din litera *L*, iar celui din Iași - literele *Ș*, *V* (*Ș* și prima parte din cele trei ale literei *V* - singura dintre acestea apărută până acum - totalizând 599 de pagini de tipar), *IV*, *X*, *Y*, precum și, din prima parte, elaborarea literei lipsă *£* și reluarea și terminarea unei părți din litera *L* pentru etimologii au fost consultați specialiști din mai multe centre universitare. Majoritatea literelor au apărut, unele pe sărite (*M* între 1965 și 1968, *N* în 1971, *O* în 1969, *P* între 1972 și 1984, *R* în 1975, *S* între 1986 și 1994, *Ș* în 1978, *T* în 1983, *T* în 1994, prima parte din *V* (până la a *veni*) în 1997 și *Z* în 2000) - în total 20 de volume -, cu excepția literelor *D*, *E*, *K*, *L*, *U*, a puținelor cuvinte începând cu litera *Q* și a ultimelor părți ale literei *V* (începând cu *venin*), la care se adaugă literele *W*, *X* și *Y*. Deosebirea cea mai importantă constă în tipurile de cuvinte reprezentate în cele două serii: la majoritatea primelor litere ale alfabetului (cu excepția celor care s-au redactat abia acum), neologismele sunt slab reprezentate, nu numai din cauza faptului că foarte multe nici nu se încetățeniseră încă în limba română la vremea elaborării volumelor respective, dar și din cauza reticenței lui Pușcariu cu privire la acest sector al vocabularului; într-o situație asemănătoare se află termenii regionali, deoarece cercetările dialectale se aflau în acea vreme abia la început. Prima parte prezintă în schimb avantajul de a putea servi ca bază pentru o prelucrare bilingvă, întrucât includea și traducerea sensurilor în limba franceză, la care a trebuit să se renunțe în perioada comunistă. Reluarea și completarea acestui dicționar, absolut necesară, nu ni se mai pare astăzi recomandabil și nici posibil de realizat prin mijloace tradiționale (fișe etc), ci exclusiv pe baze informatizate. Ea ar trebui să valorifice, printre altele, și banca de texte și cea de inovații a limbii române, despre care va fi vorba mai departe. Ar fi necesar ca partea publicată înainte de 1949 să fie reluată și adusă la zi, cu atât mai mult cu cât puține persoane și chiar biblioteci posedă dicționarul în întregime (chiar în cazul seriei noi, tirajele diverselor litere au fost diferite și în continuă scădere), iar îmbătrânirea hârtiei în cazul seriei vechi o face fragilă și greu de consultat. Având în vedere că pentru noua serie a dicționarului s-au adunat, manual, peste șase milioane de fișe cu extrase și atestări (dintre acestea, în DLR au fost incluse cea 3.200.000 de citate¹⁰, reprezentând aproximativ 88% din totalul textului), este de sperat că la reluarea, într-un viitor mai mult sau mai puțin apropiat, a primei serii se va putea uza de

avantajele elaborării computerizate, valorificându-se băncile de date în curs de elaborare în institut.

Având în vedere diferențele semnalate (dintre care unele se regăsesc între primele și ultimele litere din seria nouă), este foarte binevenită ideea actuali responsabili ai DLR de a se publica, pentru operativitate, un *Supliment* - care poate realiza relativ mai lesne - "care să înregistreze neologismele adoptate în limba literară de la începutul secolului" 20 "până în prezent, precum și o serie de cuvinte regionale incluse în atlasele lingvistice și în culegeri de pe teren și termeni vechi extrași din documente ale secolelor al XVI-lea - al XVIII-lea, editate în ultimele decenii"¹¹.

2.1.2. Din motivele expuse mai sus, la care se adaugă și faptul că DA/DL este accesibil mai ales specialiștilor și mai puțin publicului larg, institutul bucureștean pregătește între timp, la sugestia conducerii Academiei Române, sinteză a marelui dicționar academic, fără citate și izvoare și cu un sistem foarte economic de prezentare a informațiilor lexicografice. Acest *Mic dicționar academic* (MDA)¹² (care va avea totuși patru volume), inclus, alături de DL printre lucrările fundamentale ale Academiei Române, va avea cea 175 000 de intrări (cea 125.000 de cuvinte și cea 50.000 de variante); primul volum (A-C) a fost publicat în anul 2001 de editura Univers Enciclopedic. Proiectul *Micului dicționar academic*, numit astfel în opoziție cu "marele" dicționar academic, și-a propus să reducă decalajul dintre cele două serii ale acestuia, îmbogățind primele litere cu baza unor surse lexicografice mai noi. La rândul său, acest nou dicționar prezintă însă dezavantajul de a fi fost obligat, prin dimensiuni, să renunțe la citate și ilustrații, ceea ce limitează posibilitatea utilizării lui ca sursă de informații morfologice, gramaticale și stilistice; numărul neobișnuit de mare de abrevieri ne transparente, utilizate din același motiv de economie, constituie un argument suplimentar în favoarea realizării unei versiuni electronice a MDA care să permită regăsirea automată a informațiilor.

2.1.3. Spre deosebire de DA/DLR, o reflectare în general unitară a vocabularului limbii române oferă *Dicționarul explicativ al limbii române* despre a cărui siglă, DEX, se afirmă, pe drept cuvânt, că a devenit un apelat denumirea, care ar fi trebuit protejată prin înregistrare, a fost preluată abuziv în *Noul dicționar explicativ al limbii române* publicat pe CD-Rom de firmele Litera și sigla NODEX, sugerând că ar fi "un nou DEX". Prima ediție, un volum de 1.049 pagini, cuprinzând 56.569 de cuvinte și variante, a fost urmată de un *Supliment*

¹⁰ Marius Sala, G. Mihăilă (2000). Cuvânt înainte, în *Dicționarul limbii române* (DLR). Seria nouă. Tomul XIV. Litera Z, București, p. VI.

¹¹ V. I. Dănăilă (1994). De ce este nevoie de un MDA?, în "*Limba română*" XLIII, 9-10, pp. 397-406 și Marius Sala (2001). Prefața, în *Micul dicționar academic* (MDA), voi. I, A-B, București.

¹² I. Coteanu, Luiza Seche, M. Seche (conducătorii lucrării) et al. (1975). *Dicționarul explicativ al limbii române* (DEX), București.

¹ SPONGHIOS, 1992; Partea a 5-a, SPONGIAR-SWING, 1994; XI Partea 1, Litera Ș, 1978; Partea a 2-a, Litera T, T-TOCĂLIȚĂ, 1982; Partea a 3-a, TOCĂNA-TWIST, 1983; XII, Partea 1, Litera Ț, 1994; XIII, Partea 1, Litera V, V-VENI, 1997; XIV, Litera Z, 2000.

² în legătură cu reflectarea noilor norme ortografice ale limbii române în volumele DLR elaborate după 1993, semnalăm faptul că forma sânt, reflectând un fonetism real, vechi și popular, este păstrată în citatele în care nu era folosită.

Dicționarul explicativ al limbii române (DEX-S)¹⁶. Ediția a doua a DEX¹⁵ totalizează 1.204 pagini; această ediție, care se publică în continuare în tiraje succesive, totalizase numai în primii patru ani de la apariție 65.000 de exemplare vândute, după un calcul sumar rezultând că la 42 de locuitori ai României revenea un DEX. Actualmente, se poate într-adevăr afirma că, prin DEX, *best-sellerul* lingvisticii românești, Institutul a intrat în marea majoritate a caselor din România. Se preconizează ca DEX să fie realizat, în fine, într-un viitor relativ apropiat, și în format electronic. El a fost deja supus, de către Centrul de Cercetări Avansate în învățarea Automată, Prelucrarea Limbajului Natural și Modelarea Conceptuală al Academiei Române, codificării conform TEI¹⁶. Se estimează că ediția a III-a a DEX, concepută sub conducerea lui Ion Dănilă, va avea în plus față de precedenta cea 30.000 de cuvinte. Sub conducerea lui Ion Coteanu și Ion Dănilă, la sectorul de specialitate al Institutului a fost conceput și un **Nou dicționar explicativ al limbii române** (NEX), cu caracteristici diferite de cele ale DEX: inventar de cea 100.000 de cuvinte și variante (deci aproape de două ori mai multe decât prima ediție a DEX), definiții mai concise, prin eliminarea sinonimelor și - din păcate!-, neincluzarea etimologiei cuvintelor; revizuit de cei doi responsabili, el așteaptă introducerea în calculator, în vederea efectuării corelațiilor semantice definiționale și sinonimice.

2.1.4. DEX a scos practic din circulație dicționarele explicative mai vechi, limitate la limba română literară, DLRLC și DM¹⁷. Prima siglă reprezintă **Dicționarul limbii române literare contemporane**¹⁸, elaborat de institutele din București și Cluj pornind de la "baza manuscrisă" a DA și apărut între 1955 și 1957 în patru volume. El se mai folosește și astăzi - deși din el lipsesc cuvintele, sensurile și citatele neconforme cu ideologia vremii - pentru citatele cu care, spre deosebire de dicționarele de dimensiuni comparabile mai noi, sunt ilustrate sensurile cuvintelor (chiar dacă, pentru unele neologisme, citatele provin, așa cum era obligatoriu în epocă, din traducerile "operelor clasice" marxism-leninismului!). El mai merită deci atenție în virtutea faptului că, spre deosebire de DEX și de MDA,

¹⁶ Ion Coteanu, Ion Dănilă, Nicoleta Tiugan (conducătorii lucrării) et al. (1988). Supliment la Dicționarul explicativ al limbii române (DEX-S) București.

¹⁷ Ion Coteanu, Lucreția Mareș (sub conducerea) et al. (1996), Dicționarul explicativ al limbii române (DEX), ediția a II-a, București.

¹⁸ Dan Tufiș (2000). Cercetare și colaborare internațională în ingineria lingvistică la RACAI, în "Terminologia în România și în Republica Moldova", Cluj-Napoca, p. 34-36 și Recherche et collaboration internationale en industries de la langue à l'Académie Roumaine, în "Terminometre Hors-serie n° 4. La terminologie en Roumanie et en République de Moldova", p. 38-40.

¹⁹ Pentru detalii cu privire la aceste două dicționare v. Mircea Seche (1969). Dicționarele explicative ale limbii române literare, în Schiță de istorie a lexicografiei române, voi. II, De la 1880 până astăzi, București, p. 135-147.

²⁰ D. Macrea, E. Petrovici (sub direcția) et al. (1955-1957). Dicționarul limbii române literare contemporane (DLRLC), Editura Academiei, București, voi. I, A-C; II, D-L, 1956' III M-R, 1957; IV, S-Z, 1957.

include citate ilustrative, care din păcate au fost eliminate din dicționarele următoare.

2.1.5. O versiune prescurtată a acestui dicționar, cu un inventar puțin mărit și cu adăugarea etimologiei cuvintelor, dar cu eliminarea citatelor, a fost publicat de Institutul din București în 1958 sub titlul **Dicționarul limbii române moderne** (abreviat DM).

2.1.6. Un dicționar de un tip special, cu o utilitate mult mai largă decât aceea care i se recunoaște de obicei, elaborat de data aceasta de colectivul de gramatică al Institutului (condus până de curând de Mioara Avram²⁰), este **Dicționarul ortografic, ortoepic și morfologic al limbii române** (DOOM) \ Este singurul dicționar al limbii române (mai bogat decât DEX) care conține ample informații cu privire la formele flexionare ale cuvintelor variabile incluse, putând servi astfel (chiar dacă aceste informații nu sunt exhaustive) ca sursă pentru studii și aplicații de morfologie. Institutul are în prezent în lucru, sub conducerea subsemnatei, o a doua ediție, parțial revăzută și adăugită, a DOOM (care va cuprinde și cuvinte neînregistrate în nici un dicționar românesc până în prezent). Aceasta va apărea în anul 2003, inclusiv pe CD-Rom, și va trebui să servească drept bază unui nou corector ortografic și morfologic, care să țină seamă de modificarea unor recomandări oficiale în raport cu cele încă în vigoare.

2.1.7. În fine, un dicționar mai puțin obișnuit, **Dicționarul invers**²¹, în care cuvintele sunt ordonate alfabetic pornind dinspre sfârșitul lor, este deosebit de util specialiștilor pentru studierea terminațiilor, a desinențelor și a sufixelor, dar și poezilor, fiind utilizabil și ca dicționar de rime. Această lucrare - care, spunând "legenda", a valorificat experiența din copilărie a uneia dintre autoare, care folosește în joacă o *păsărească* de acest fel - ar merita și ea o nouă elaborare, pe baza unui inventar mai bogat și actualizat de cuvinte și a unui program care să permită "răsturnarea" lor automată.

2.1.8. Institutul a publicat, încă din 1968, un dicționar al lexicului unor autori, primul ales neputând fi altul decât Eminescu - **Dicționarul limbii poetice a lui Eminescu**²², care însă, la acea vreme, nu se putea baza, evident, pe stabilirea concordanțelor așa cum se realizează ea în zilele noastre.

2.1.9. Institutul a elaborat de asemenea o serie de dicționare ale limbii române pe epoci sau pe probleme, cum sunt **Dicționarul limbii române literare**

¹⁹ D. Macrea (sub direcția) (1958). Dicționarul limbii române moderne, București.

²⁰ Pentru activitatea acestuia v. Mioara Avram (1999). Colectivul de gramatică, în Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit., p. 113-125.

²¹ Mioara Avram (red. resp.) et al. (1982). Dicționarul ortografic, ortoepic și morfologic al limbii române (DOGM), București, 1982.

²² *** (fsgjy Dicționar invers, București. V. și Mircea Seche (1969). Schiță de istorie a lexicografiei române, voi. II, De la 1880 până astăzi, București, p. 254-255.

²³ Tudor Vianu (sub redacția) et al. (1968). Dicționarul limbii poetice a lui Eminescu, București.

vechi²⁵ și *Dicționarul împrumuturilor latino-romanice în limba română veche*²⁶, publicate de sectorul de limbă literară, filologie și poetică²⁶, condus de Ion Gheție,²⁷ iar în prezent de Alexandru Mareș - și *Dicționarul elementelor românești din documentele slavo-române*²⁷, elaborat la sectorul de slavistică²⁸ - dicționare destinate în primul rând specialiștilor.

2.1. 10. Un cercetător din institut, Constant Măneca, a publicat, împreună cu Florin Marcu, un extrem de util, cu toate criticile care i s-au adus, *Dicționar de neologisme*²⁹, reluat și dezvoltat, după moartea celui dintâi, de Florin Marcu, în numeroase variante, de diverse dimensiuni, la diferite edituri, inclusiv pe CD-Rom.

2.1.11. Se află în lucru și *Dicționarul etimologic al limbii române* (DELR) - coordonator: Marius Sala -, altă lucrare fundamentală a Academiei Române, la care colaborează cercetători din toate sectoarele Institutului, cercetători din Cluj și Timișoara și cadre didactice de la universitățile din București, Cluj și Timișoara.

2.1.12. Pe lângă resursele privitoare la numele comune, Institutul a elaborat și importante lucrări consacrate numelor proprii³⁰.

Astfel, în domeniul toponimiei, după clasică lucrare a lui Iorgu Iordan³¹, s-a realizat în Institut *Dicționarul toponimic al României*, partea I, *Oltenia*³², elaborat sub conducerea lui Gh. Bolocan în colaborare cu cadre didactice de la Universitatea din Craiova, din care au apărut în perioada 1993-2001 primele trei volume, precum și al doilea dicționar din serie, consacrat *Munteniei* și aflat în curs de definitivare.

În domeniul onomasticii, de asemenea urmând altei lucrări clasice a lui Iordan³³, Institutul colaborează și la proiectul internațional PatRom, care realizează un dicționar istoric de antroponomie romanică, în care este reprezentată și limba română, și din care până acum a fost publicat un prim volum de prezentare³⁴.

2.2. Dicționare bilingve și multilingve

2.2.1. Pe lângă dicționarele monolingve ale limbii române, Institutul a realizat și unele din cele mai importante dicționare bilingve³⁵ (englez-român, german-român³⁷, rus-român³⁸, ceh-român³⁹ și sârb-român⁴⁰ - perechea sârb-român este un dicționar român-sârb, fiind în curs de redactare; un dicționar francez-român rămas nepublicat) și frazeologice (spaniol-român, sub tipar, și român-spaniol, în curs de elaborare), cărora li se adaugă dicționare bilingve⁴¹ - care au început să fie transpuse și pe CD-Rom - și dicționare frazeologice românești⁴² și bilingve românești⁴³ elaborate de unii membri ai Institutului; *Dicționarul elen-român*, lucrare colectivă, este aproape și el de sfârșit.

2.2.2. Institutul a colaborat și la mai multe dicționare multilingve⁴⁴, din care se distinge în mod deosebit un lexicon multilingv de un tip special - adevărată premieră internațională - *Dicționarul elementelor latinești savante din limbile romanice*, elaborat la sectorul de romanistică (condus inițial de mar-

³³ Iorgu Iordan (1983). Dicționar al numelor de familie românești, București.

³⁴ *** (fQQj) pictionnaire historique d'anthroponymie romane (PatRom). Presentation de la collection, 2001, Institutul de Lingvistică și Literatură, proiect, Tübingen.

³⁵ V. și Ilinca Constantinescu. (1999). Fostul sector de germanistică, în Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit., p. 174-179.

³⁶ L. Levițchi (red. resp.) et al. (1974). Dicționar englez-român, București. Suplimentul la acest dicționar, care nu a mai apărut, coordonat de Ilinca Constantinescu, va fi inclus în nouă ediție, mult mărită, a dicționarului, aflată sub tipar și care va reprezenta cel mai bogat dicționar englez-român.

³⁷ M. Isbășescu, Măria Iliescu (coord. și revizie) et al. (1966, 1988). Dicționar german-român, București, 1966; ediția a II-a revăzută și îmbogățită, București, 1988.

³⁸ Gheorghe Bolocan (redactor responsabil) (1964). Dicționar rus-român, București.

³⁹ S. Stați (red. resp.) et al. (1967). Dicționar ceh-român, București.

⁴⁰ M. Tomici (1998-2000). Dicționar sârb-român, 3 voi, Timișoara.

⁴¹ Gh. Bolocan (1972). Dicționar bulgar-român, București - Sofia; Gh. Bolocan et al. (1983). Dicționar român-rus, București - Moscova; Al. Calciu, C. Duhăneanu, D. Munteanu (1979). Dicționar român-spaniol, București; Ana Canarache (coord.) (1967, 1978). Dicționar român-francez, București, 1978; M. Isbășescu (red. resp.) (1963), Dicționar român-german, București; Valeria Neagu (2001). Dicționar român-spaniol (cu transpuneri pe CD-Rom), București.

⁴² V. Breban et al. (1969). Dicționar de expresii și locuțiuni românești București.

⁴³ Gh. Bolocan et al. (1968). Dicționar frazeologic rus-român, București; H. Mantsch et al. (1979). Dicționar frazeologic român-german, București.

⁴⁴ *** fIQQij Dictionnaire de la presse écrite et audiovisuelle. Espagnol-français-italien-portugais-roumain, Paris; *** (2001). Usage Dictionary of Anglicisms in Selected European Languages (UDASEL) Oxford ș.a.

²⁴ Mariana Costinescu, Magdalena Georgescu, Florentina Zgraon (1987). *Dicționarul limbii române literare vechi (1640-1780). Termeni regionali* București.

²⁵ Gh. Chivu, Emanuela Buză, Alexandra Roman Moraru (1992). *Dicționarul împrumuturilor latino-romanice în limba română veche (1421-1760)* București.

²⁶ V. Ion Gheție (1999). Colectivul de limbă literară și filologie, în Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit., p. 132-143.

²⁷ Gheorghe Bolocan (redactor responsabil) et al. (1981). *Dicționarul elementelor românești din documentele slavo-române. 1374-1600*, București.

²⁸ *Cu privire la care v. Virgil Nestorescu (1999). Sectorul de lexicografie bilingvă. Fostul sector de slavistică, în Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit. p. 165-174.*

²⁹ F. Marcu, C. Măneca (1961-1978). *Dicționar de neologisme*, București, 1961; ed. II revăzută și adăugită, 1966; 1978. V. și Mircea Seche (1969). *Schită de istorie a Jexicografiei române*, voi. II, *De la 1880 până astăzi*, București, p. 154-159.

Pentru activitatea în acest domeniu v. Gheorghe Bolocan, Ecaterina Mihăilă (1999). Colectivul de onomastică și Domnița Tomescu (1999). Grupul de lucru PatRom, în Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit., p. 125-132.

³¹ Iorgu Iordan (1952-1963). *Nume de locuri românești în Republica Populară Română* București, 1952; *Toponimia românească*, București, 1963.

³² Gh. Bolocan (sub redacția) et al. (1993-2001). *Dicționarul toponimic al României. Oltenia* (DTRO), voi. I-III, Craiova, Editura Universitară.

romanist Iorgu Iordan, apoi de Marius Sala și în prezent de subsemnata)⁴⁵, în colaborare cu cadre didactice de la Facultatea de Limbi și Literaturi Străine a Universității din București și în coordonarea prof. dr. Sanda Reinheimer Rîpeanu, decanul Facultății. Negăsindu-și un editor "clasic" din cauza costurilor prea ridicate, acest dicționar va fi publicat direct pe Internet, sub auspiciile Universității din București.

3. Bănci de date

3.1. Institutul a avut în proiect încă din anii 1978-80 realizarea primei bănci computerizate de date lingvistice din România (**Banca de date fonomorfosemantice a limbii române - BANDASEM**)⁴⁶, cel dintâi modul fiind cel de semantică, proiectat pentru un *Dicționar confruntativ de sinonime, de analogii și de asociații al limbii române* (DCSAAs). Redactarea acestuia, care a ajuns la litera S, s-a făcut însă cu mijloace tradiționale, deși prin colaborarea cu Centrul de Calcul al Universității din București se elaborase un modul de program în sistemul Socrate pentru recunoașterea și selectarea, ca probă, a analogiilor și a asociațiilor cuvântului *blitz*. Elaborarea DCSAAs a fost întreruptă pentru un timp în favoarea lucrărilor prioritare al Academiei, iar reluarea lui se va putea face, sperăm, cu mijloacele informatice disponibile actualmente⁴⁷.

3.2. O minibancă inițiată în cadrul sectorului de gramatică al Institutului, a cărei alimentare a fost din păcate întreruptă în favoarea concentrării forțelor pentru realizarea ediției a doua a "Gramaticii Academiei", este **Banca de inovații a limbii române**, bazată pe monitorizarea presei scrise și audiovizuale actuale.

3.3. Având în vedere că în DOOM informația este atomizată, în folosul cititorului neprofesionist, în cadrul fiecărui cuvânt-titlu în parte, dar este greu de sistematizat de către specialist, Institutul are în proiect, începând din 2003, realizarea unui baze de date care să permită nu numai elaborarea unui **Nou dicționar ortografic, ortoepic și morfologic al limbii române** și a unor dicționare specializate de un tip asemănător, precum și aducerea lor permanentă la zi, ci și gruparea cuvintelor în clase în funcție de caracteristicile lor fonetice, grafice și morfologice⁴⁸.

Cu privire la activitatea acestuia v. Marius Sala (1999). Sectorul de limbi romanice și clasice, în Mioara Avram, Marius Sala, Ioana Vintilă-Rădulescu, op. cit., p. 147-164.

⁴⁶ Ion Dănăilă (2000). Proiecte de prelucrare electronică a vocabularului limbii române, în *"Terminologia în România și în Republica Moldova"*, Cluj-Napoca, p. 36-37.

⁴⁷ Partea de fonetică/grafematică și de morfologie a BANDASEM a fost cedată institutului omolog din Cluj, pentru care v. Felicia Șerban et al. (2000). Baza de date a limbii române, în *"Terminologia în România și în Republica Moldova"*, Cluj-Napoca, p. 37-38 și La base de données de la langue roumaine, în *"Terminometre Hors-serie n° 4. La terminologie en Roumanie et en République de Moldova"*, p. 40-42.

⁴⁸ Clasificarea cuvintelor românești conform modului lor de flexiune, realizată de Alf Lombard, Constantin Gâdei (1981). Dictionnaire morphologique de la langue roumaine,

3.4. Institutul are în proiect și elaborarea sau definitivarea terminologice⁴⁹ (dicționare terminologice bi- și multilingve, valoare elaborate în cadrul proiectului PRACTEAST din cadrul programului C O al Comisiei Europene⁵⁰ și un dicționar al termenilor oficiali); de al membri ai Institutului au colaborat la realizarea **Băncii de date** (BDT) multilingve a Asociației Române TermRom⁵¹, care, cu sprijin terminologie și inginerie lingvistică a Uniunii Latine, este accesibil TermRom găzduit de CIMEC (<http://www.cimec.ro/tr>) și, de curând, și Rețeaua României (prin subsemnata) în Rețeaua Panlatină de (Realiter)⁵² și în Rețeaua Francofonă de Amenajare Lingvistică constitui desigur un sprijin în dezvoltarea resurselor terminologice română în conformitate cu normele și recomandările internaționale.

Lund - București, bazată pe inventarul DEX, prezintă unele inexactități insuficientei cunoașteri de către autori a limbii române actuale; ea co bazele realizării, în Republica Moldova, a unui pachet de programe destin nivel morfologic, pentru care v. Elena Boian et al. (2000). Instrumentar lingvistice, în "Terminologia în România și în Republica Moldova", Cluj-Napoca, p. 13-14. Instruments pour applications linguistiques, în "Terminometre Hors-serie terminologie en Roumanie et en République de Moldova", p. 42-44; o gr unui număr limitat de cuvinte ale limbii române a fost realizată de Flora Șoșa (1999) în îndreptar ortografic și morfologic; București.

⁴⁹ V. Ioana Vintilă-Rădulescu (1999). Institutul de Lingvistică "Iorgu Iordan" din Cluj-Napoca, în *"Terminologia în România și în Republica Moldova"*, Cluj-Napoca, p. 13-14. Linguistică Iorgu Iordan de Bucarest, în *"Terminometre Hors-serie n° 4. en Roumanie et en République de Moldova"*, p. 22-23.

⁵⁰ Nicoleta Petuhov. (2000). Colaborarea românească la proiectul PRACTEAST, în *"Terminologia în România și în Republica Moldova"*, Cluj-Napoca, p. 64-66. collaboration roumaine au projet Practeast, în *"Terminometre Hors-serie terminologie en Roumanie et en République de Moldova"*, p. 64-66.

⁵¹ Dan Matei (2000). Banca de date terminologice a TermRom, în *"Terminologia și în Republica Moldova"*, Cluj-Napoca, p. 29-30 și La banque de données de TermRom, în *"Terminometre Hors-serie n° 4, La terminologie en Roumanie et en République de Moldova"*, p. 32-33.

⁵² Dan Matei (2000). Prezența românească în rețeaua panlatină de terminologie, în *"Terminologia în România și în Republica Moldova"*, Cluj-Napoca, p. 51-52 și La coopération dans le cadre de l'Agence Intergouvernementale de la Francophonie), în *"Terminometre Hors-serie terminologie en Roumanie et en République de Moldova"*, p. 57-58.

4. Corpusuri

O altă categorie importantă de resurse lingvistice o constituie corpusurile, la Institut fiind în curs de realizare o **Bancă de texte românești**, care cuprinde texte din secolele al XVI-lea - al XVII-lea, introduse integral în calculator, și în care se prevede introducerea câtorva sute de texte din toate epocile. Inițiată de directorul institutului, acad. Marius Sala, Banca a fost deja valorificată în elaborarea unor teze de doctorat, printre altele la aceea a Janei Balacciu-Matei. Pentru exploatarea ei deplină în vederea identificării primelor atestări ale cuvintelor limbii române din fondul vechi, necesare MDA și *Dicționarului etimologic al limbii române* (DELR), a îmbogățirii dicționarului limbii române în general și a dezvoltării studiilor privind istoria limbii române literare și a limbii noastre în ansamblu este necesară achiziționarea unor programe de ultimă oră, precum și specializarea unor persoane pentru utilizarea lor eficientă. Sperăm de asemenea că într-un viitor nu prea îndepărtat se va realiza și dorita joncțiune cu Banca de texte din faza modernă și contemporană a limbii române, proiectată a se realiza la Centrul de Studii Românești de pe lângă Universitatea din Anvers, inaugurat în primăvara anului 2000 sub conducerea cunoscutei romaniste și romaniste Liliane Tasmowski.

5. Resurse bibliografice

Amintim pe scurt și principalele resurse bibliografice privitoare la limba română elaborate de Institut sau de membri ai acestuia⁵⁴. *Bibliografia limbii române*, inițiată de Al. Rosetti și definitivată de Aurel Nicolescu, a rămas nepublicată. *Bibliografia românească de lingvistică (BRL)* referitoare la lucrările de lingvistică apărute în țară începând din 1944 apare anual în revista "Limba română"; în 1999, ea totalizase deja 64.340 de titluri, în peste 3.300 de pagini de tipar; se preconizează introducerea în calculator a tuturor numerelor din BRL în vederea publicării unui volum cu itemurile ordonate pe autori și pe domenii (descrise și separate mai amănunțit decât în forma apărută, cronologic, cu indice de domenii, materii, cuvinte, autori etc).

Pentru domeniul terminologiei s-au realizat bibliografiile ale dicționarului terminologic, respectiv ale studiilor de terminologie⁵⁵ și ale standardelor românești de/cu terminologie⁵⁶, precum și un repertoriu bio-bibliografic al terminologilor

⁵⁴ I. Coteanu, I. Dănăilă (1970). Introducere în lingvistica și filologia românească Probleme. Bibliografie, București; T. Vianu (red. resp.) et al. (1972). Bibliografia analitică a limbii române literare. 1780-1866, București; Gh. Chivu, Mariana Costinescu (1974). Bibliografia filologică românească. Secolul al XVI-lea, București.

⁵⁵ Anca Fezi et al. (2000). Bibliografia lucrărilor de terminologie (1990-1999). România, în "Terminologia în România și în Republica Moldova", Cluj-Napoca, p. 103-113 și pe discheta anexată revistei Terminometro Hors-serie n° 4. La terminologie en Roumanie et en Republique de Moldova".

⁵⁶ Aurora Peșan, EdySăvescu (2000). Standarde românești de/cu terminologie (1990-1999). România, în "Terminologia în România și în Republica Moldova", Cluj-Napoca, 2000, p.

din România⁵⁷, inclus în repertoriul internațional al terminologilor din domeniul neolatin pregătit de Uniunea Latină și accesibil pe Internet.

6. Concluzii

Nu ne vom referi aici la alte tipuri de lucrări (gramatici⁵⁸, tratate⁵⁹, enciclopedii⁶⁰ etc.) elaborate de Institut sau de cercetători ai acestuia ori la alte tipuri de resurse care ar merita să fie elaborate de noul institut, pentru a înlocui lucrări mai vechi și a valorifica posibilitățile oferite culturii de societatea informațională, de exemplu un nou dicționar de frecvență al limbii române ș.a.

Deși dicționarele pe CD-Rom și cele pe Internet sunt solicitate de tot mai mulți utilizatori din țară și din străinătate, care cer tot mai des informații cu privire la eventuale dicționare românești on-line, până în prezent a existat la noi o anumită reticență a editurilor proprietare ale drepturilor asupra edițiilor pe suport tradițional de hârtie față de acest nou mod de difuzare. Nu trebuie însă să existe temerea că folosirea și a noilor suporturi ar diminua vânzarea cărților, în condițiile în care, în ciuda tuturor eforturilor, un procent încă infim din populația României are acces la PC-uri. De altfel, practica altor țări a arătat că, în mod neașteptat, difuzarea și în format electronic chiar a sporit desfacerea cărților, cărora le-a făcut în felul acesta reclamă și care prezintă, la rândul lor, alte avantaje în utilizare în raport cu CD-Romurile, cele două tipuri specializându-se și în funcție de necesități. Astfel, având în vedere culegerea lor computerizată, atât DEX, cât și MDA și DOOM ar putea fi primele dicționare ale Institutului difuzate în viitor și pe CD-Rom.

Credem că și diverse lucrări valoroase ale Institutului, care, exclusiv din motive financiare, nu-și găsesc editori de ani de zile, nici în țară, nici în străinătate (ca *Bibliografia limbii române*, *Dicționarul spaniolei americane* ș.a.), ar putea

117-126 și pe discheta anexată revistei Terminometro Hors-serie n° 4. La terminologie en Roumanie et en Republique de Moldova".

⁵⁷ Adriana Marinescu (2000). Repertoriul bibliografic al terminologilor. România, în "Terminologia în România și în Republica Moldova", Cluj-Napoca, 2000, p. 128-139 și pe discheta anexată revistei Terminometro Hors-serie n° 4. La terminologie en Roumanie et en Republique de Moldova".

⁵⁸ *** (1954, 1963). Gramatica limbii române, București, ed. I, 1954; ed. a II-a, revăzută și adăugită, 1963; Mioara Avram (1986, 1997, 2001). Gramatica pentru toți, București, 1986, 1997, 2001.

⁵⁹ Al. Rosetti (redactor responsabil) et al. (1965, 1969). Istoria limbii române. București, vol. I. Limba latină, voi. al II-lea; Al. Graur, Mioara Avram (1970-1989). Formarea cuvintelor limbii române, București: I. Fulvia Ciobanu, Finuța Hasan (1970). Compunerea; II. Mioara Avram et al. (1978). Prefixele, 1978; III. Laura Vasiliu (1989). Sufixe, 7. Derivare verbală etc.

⁶⁰ Marius Sala, Ioana Vintilă-Rădulescu (1981). Limbile lumii. Mică enciclopedie, București (1984). Les langues du monde. Petite encyclopedie, București - Paris; Marius Sala (coord.) et al. (1989). Enciclopedia limbilor romanice, București; (2001), Enciclopedia limbii române, București.

valorificare prin aducerea lor la cunoștința celor interesați pe această cale, tot mai utilizată în societatea informațională actuală. O condiție pentru viitor este realizarea din capul locului a lucrărilor institutului pe calculator, care a devenit posibilă prin tot mai buna dotare tehnică a Institutului, realizată prin eforturile directorului său, precum și prin însușirea, de către un număr tot mai mare de cercetători din Institut, în special din generațiile tânără și mijlocie, a cunoștințelor de operare pe calculator, inclusiv, în unele cazuri, a lucrului cu baze de date.

Prin realizarea proiectelor de editare pe CD-Rom și pe Internet vom recupera relativa întârziere în acest domeniu față de difuzarea în România, de către Grupului Editorial Litera din Republica Moldova și firma Litera Internațional, cu sediul în București, a unor CD-Romuri cuprinzând, în diverse combinații, mai multe titluri⁶¹. Sperăm că CD-Romurile consacrate unor dicționare ale Institutului vor fi, deși tot protejate, mai ușor de instalat decât cele de la Litera și că vor oferi mai multe facilități în utilizare decât acestea, care nu sunt foarte practice, mai ales pentru cercetători, în ciuda structurii lor modulare și a interfeței lor comune, despre care în reclamă se spune că permit activarea simultană a tuturor dicționarelor.

Pentru progresul cercetărilor și dezvoltarea și prelucrarea resurselor la nivelul exigențelor pe plan mondial, credem că în viitor se impune o mai bună colaborare, în interes reciproc, între lingviști și informaticienii preocupați de probleme asemănătoare.

Contribuția lingvisticii la studiul terminologiilor științifice

Angela BIDU-VRÂNCEANU

Universitatea din București, Edgar Quinet nr. 5-7

vrancean@gpsnet.ro

1. Se admite "laicizarea" științelor [1] sau importanța lor socio-culturală economică și pedagogică tot mai mare în societățile moderne. Aceasta înseamnă că *limbajele specializate* și *terminologiile* lor nu mai reprezintă coduri tot mai inaccesibile vorbitorilor obișnuiți, nespecializați sau de altă specialitate. În direcția deschiderii, chiar și parțiale a codurilor științifice, *dicționarele generale* [2], care includ un număr destul de mare de termeni științifici joacă un rol deosebit pentru a asigura accesul la sensul specializat oricărui vorbitor insuficient informat, pentru a ajuta să rezolve ambiguitățile de diferite tipuri și chiar să utilizeze adecvat terminologie. Permanenta raportare la dicționarele generale ca forme instituționalizate de reglare a uzului nu numai al cuvintelor din limba comună, ci și al termenilor specializați constituie premisa de la care pornim pentru a susține importanța lingvisticii în descrierea terminologiilor științifice, în receptarea și utilizarea lor adecvată chiar și de către nespecialiști.

Pe aceste poziții s-a situat activitatea în cadrul a trei contracte de cercetare științifică pe anii 1997, 1999 și 2000, finanțate de CNCSIS (Consiliul Național de Cercetare Științifică). Au fost studiate limbajul **filozofic**, terminologiile **matematică**, **mineralogică** și din **artele plastice** și, dintr-o perspectivă mai limitată **medicină**, **lingvistica** și **științele politice**. Rezultatele cercetărilor au fost publicate în două volume: *Lexic comun*, *lexic specializat* [3], care conține studii cu caracter monografic și *Lexic științific interdisciplinar* [4], reprezentând o sinteză lexicografică generală și specializată pentru termenii din fiecare dintre domeniile studiate care apar mai mult decât într-o terminologie științifică.

În toate cercetările întreprinse s-a urmărit adoptarea *unei grile metodologice comune* atât pentru clase de cuvinte din limba comună (*abstractele*), cât și pentru termenii specializați din orice domeniu. S-a obținut atât caracterizarea fiecărei terminologii studiate în parte, cât și desprinderea unor trăsături generale ale terminologiilor științifice, relevante din punct de vedere lingvistic. S-au avut în vedere aspecte *paradigmatice* privind diferitele modalități de definire a sensului și relațiile semantice (*monosemie/polisemie*, *hiponimie*, *sinonimie*) din perspectiva necesității ca termenii științifici să fie monoreferențiali, univoci din punct de vedere

Corectorul electronic ORTO 2001 ROM SP, Dicționarul ortografic al limbii române, Gramatica uzuală a limbii române, Noul dicționar explicativ al limbii române⁶¹ Marele dicționar de neologisme de Florin Marcu, Dicționarul de dublete etimologice ale limbii române de Marcu Gabinschi și un Dicționar de termeni de afaceri englez-român.

semantic și să nu aibă sinonime. Analiza **sintagmatică** a gradului de non-determinare contextuală ca o condiție de exprimare a sensului specializat a individualizat terminologiile științifice studiate, de la o *libertate contextuală* mai mare (terminologia **matematică**, **mineralogică**) sau relativă (terminologia **filozofică**) până la o *strictă determinare contextuală* (terminologia **politică** și din **artele plastice**). Acolo unde independența contextuală e mai mare, determinările contextuale exprimă în mod similar în diferite terminologii (**matematică**, **filozofică**, **lingvistică**) subcategoriile științifice care dezambiguizează lexicul științific interdisciplinar. Caracterizarea termenilor științifici prin mărci diastratice în dicționarele generale și enciclopedice ca tipuri de informații sintagmatice reprezintă un aspect foarte important pentru uzajul adecvat de către specialiști, aspect deficitar, inegal rezolvat.

De pe poziția receptorului nespecializat care decodează sensul total sau parțial, un rol important îl are *definiția lexicografică* care, spre deosebire de cea *terminologică* trebuie să fie mai mult sau mai puțin *naturală* și prin aceasta accesibilă. Existența celor două tipuri de definiții ale termenilor specializați este în general admisă și compararea lor este favorizată de prezentarea sintetică, sinoptică propusă de noi [4]. Chiar și în cazul definițiilor strict terminologice, Em. Vasiliu [5] a susținut și demonstrat prin diferite exemple relevanța diferită a unor componente de sens pentru vorbitorul specialist sau non-specialist. Pornind de la aceste constatări de principiu, ar fi justificat ca termenii științifici să aibă *definiții alternative*, științifice și pre-științifice [6], condiționate atât de o interpretare semantică, cât și de una pragmatică. Din această perspectivă, definițiile termenilor științifici în dicționarele generale ar trebui să difere de cele din dicționarele specializate pentru a facilita deschiderea codurilor științifice și pentru a dezambiguiza lexicul științific interdisciplinar (din principiu, de interes mai larg) sau tangențele cu limba comună. Din păcate, cu mici excepții (**matematica**) selecția termenilor științifici și definirea lor nu diferă aproape deloc în dicționarele generale și în cele specializate.

2. Din perspectivă lingvistică, terminologiile investigate prezintă o serie de particularități:

Matematica se caracterizează prin cel mai mare grad de abstractizare și de ermetism la nivelul sensurilor și definițiilor lor. Compararea definițiilor specializate cu cele din dicționarele generale arată că acestea din urmă definesc diferit și mai accesibil termenii, fără a afecta precizia lor semantică. Sensurile univoce, fără sinonime nu sunt condiționate contextual; sintagmele mai mult sau mai puțin fixe diferențiază subcategoriile conceptuale (de ex. *sistem de ecuații*, ~ *de curbe*, ~ *de numerație*, ~ *de referință*) și nu afectează independența semantică a acestora. Această terminologie dispune de cea mai bună marcă diastratică în DEX, chiar dacă există numeroase situații în care apartenența la matematică rezultă numai din definiție (manieră de caracterizare practică sistematic și nu întotdeauna convenabil de DEX în cazul altor terminologii). **Matematica** are cel

mai bogat lexic științific interdisciplinar, cei mai numeroși termeni comuni fiind cu **fizica**, **filozofia**, **logici**, dar și cu **lingvistica**, **biologia**, **arhitectura** ș.a.; termenii interdisciplinari își păstrează aproape neschimbat sensul, indiferent de domeniul în care se utilizează. Dacă în unele cazuri (relația cu **fizica**, **logica**, **filozofia**) punctul de plecare pentru lexicul interdisciplinar nu se poate stabili cu certitudine, în destule alte situații, **matematica** este sursa "împrumutului" făcut de alte științe (**artele plastice**, **arhitectură**, **lingvistică** ș.a.)

Mineralogia reprezintă și ea un grad mare de ermetism sau închidere a codului, majoritatea termenilor fiind univoci semantic, monoreferențiali și implicit, independenți contextual. Determinările contextuale reprezintă subtipuri, ca și în alte terminologii (**matematică**, **filozofie** de ex.; *acvamarin brazilian*, ~ *sintetic*, ~ *siamez*, etc.) Are un număr mai limitat de termeni comuni cu alte științe (**chimia**, **artele plastice**, **simbolistica**) și, cel puțin pentru ultimele două, **mineralogia** este punctul de origine al termenilor interdisciplinari. În ciuda caracterului strict specializat al acestei terminologii, marcarea diastratică din dicționarele generale este deficitară.

Terminologia **filozofică** se caracterizează printr-un grad oarecare de ambiguitate, determinat de variații de interpretare în funcție de curente și tipuri de texte, dar și de contactele cu alte științe sau cu limba comună. De aceea definițiile termenilor **filozofici** nu se pot limita la dicționare, fiind necesară analiza strategiilor argumentative și a figurilor textuale; Invers proporțional cu această necesitate de dezambiguizare, DEX-ul prezintă o marcă diastratică deficitară atât pentru termenii filozofici, cât și pentru celelalte terminologii cu care se stabilesc interdisciplinarități, cum ar fi **matematica**, **lingvistica** și alte domenii **umaniste**. O bună parte a lexicului științific interdisciplinar are ca punct de plecare **filozofia**, al cărei sens se păstrează ca o medie semantică în majoritatea disciplinelor. Ca și în alte științe, determinarea contextuală exprimă în general subtipuri (de ex. *sistem al științelor*, ~ *axiomatic*, ~ *filozofic*).

Terminologia **artelor plastice** prezintă aspecte paradoxale. Maniera de înregistrare și de definire echivocă, imprecisă a acestor termeni în dicționarele generale dă impresia unui nespecialist de falsă accesibilitate, interpretare contrazisă categoric de definițiile precise, riguroase din dicționarele și textele specializate. Dependența contextuală strictă a numeroși termeni din **artele plastice**, al căror sens specializat e condiționat de sintagmele fixe în care apare (de ex. *acord cromatic*, *compoziție de gen*, *semn plastic*) reprezintă o altă caracteristică a acestei terminologii. **Artele plastice** au un lexic științific interdisciplinar bogat, în care se remarcă faptul că sunt preluați cu unele modificări semantice (privind interesul pentru acest domeniu) termeni din alte științe, cum ar fi **chimia**, **mineralogia**, **matematica**, **fizica**. DEX-ul nu utilizează decât mărcile diastratice (pictură), (sculptură) dispuse nesistematic și rar, ceea ce contribuie la o tratare deficitară a acestei terminologii.

Lexicul științelor politice prezintă, din perspectiva analizei întreprinse de noi, o serie de particularități (unele asemănătoare cu artele plastice). Se remarcă dependența contextuală strictă a acestei terminologii, nici unul dintre termeni nefiind total liber contextual. Sensul specializat în științele politice se exprimă, deci, aproape exclusiv pe cale sintagmatică, în contexte mai mult (*celulă de criză, agregare de interese*, de ex.) sau mai puțin fixe (diverse combinații cu adjectivul politic în sintagme nominale: *capital politic, cartel ~ algoritm ~, contract ~, dialog ~ alternanță politică*). Preia (fără să fie niciodată punct de plecare termeni din numeroase și variate științe: economia, filozofia, dreptul, dar și lingvistica, biologia, medicina, geografia, fizica, psihologia, sportul. În majoritatea acestor cazuri nu există o motivare de conținut strictă (dincolo de întrebuintarea metaforică), ceea ce determină, în mare parte, mai curând un lexic științific interferent decât unul interdisciplinar. Poate și din cauza modificărilor continue și rapide din domeniul politicii, DEX-ul înregistrează în mică măsură termeni și sensuri din acest domeniu diastratic, ceea ce constituie un dezavantaj în impunerea acestei terminologii.

3. Analiza lingvistică a limbajelor științifice (care ar putea fi extinsă) permite caracterizarea unor terminologii ca "*puternice*" (matematica, mineralogia de ex.), iar a altora mai "*slabe*" în diferite forme și grade (de ex. științele politice, artele plastice), cu dificultăți mai mari de deschidere a codurilor în cazul primei categorii.

Delimitarea componentelor de sens relevante diferit în funcție de vorbitori specializați și nespecializați ar putea constitui o bază obiectivă pentru rezolvarea mai eficientă a definițiilor alternative în dicționarele generale, foarte importante în "laicizarea" științelor necesară în grade diferite în epoca actuală. Expriarea sensului specializat condiționat de dependențele contextuale mai mici (pentru terminologiile "puternice") sau mai mari (pentru terminologiile "slabe") constituie o caracterizare lingvistică relevantă. În schimb, în unele cazuri (ca pentru terminologia politică), determinările contextuale sunt mai favorabile, "transparentei" semantice sau deschiderii codurilor specializate.

Analiza lexicului științific interdisciplinar (LSI) poate contribui și ea la determinarea specificului unor terminologii. Științele care constituie sursa, punctul de plecare pentru o parte a LSI își susțin, și pe această cale, statutul de terminologie "puternică" (de ex. matematica, fizica și, din acest punct de vedere filozofia). Dimpotrivă, atunci când punctul de plecare nu se poate stabili aproape niciodată la nivelul unor terminologii (științele politice, artele plastice), aceasta constituie o modalitate de determinare specifică. Diferențierea interdisciplinarităților (cu o motivare de conținut determinată de considerarea referentului din diferite puncte de vedere sau de un transfer conceptual) de simplele interferențe (mai puțin sau deloc motivate, cu modificări de sens ale termenilor, multe metafore) se bazează pe aprecierea distanței semantice, verificată obiectiv.

Dat fiind rolul dicționarelor generale în impunerea și extinderea terminologiilor științifice, de interes pentru diferite categorii de vorbitori, carențele constatate în tratarea sensului și în marcarea lor diastratică riguroasă conduc la concluzia necesității unei reconsiderări și remedieri a manierei de tratare din perspectiva "laicizării" științelor.

Referințe bibliografice

- [1] F. Rastier (1995) Le terme; entre ontologie et linguistique. *Banque des mots* 1995/7, p. 35-65.
- [2] DEX - Dicționar explicativ al limbii române, (1996) ed.a 2-a sub coord. acad- I. Coteanu și Dr. Lucreția Mareș, Ed. Univers Enciclopedic, București 1996.
- [3] A. Bidu-Vrânceanu - coordonator (2000). *Lexic comun, lexic specializat*, Editura Universității din București, 2000, cu colaboratorii: Alice Toma (matematică), Silvia Săvulescu (mineralogie), Claudia Ene (filozofie), Alexandra Vrânceanu (arte plastice).
- [4] A. Bidu-Vrânceanu - coordonator (2001). *Lexic științific interdisciplinar*, Editura Universității din București, 2001, cu colaboratorii: Silvia Săvulescu (științe politice și mineralogie), Alice Toma (matematică), Claudia Ene (filozofie), Alexandra Vrânceanu (arte plastice).
- [5] Em. Vasiliu (1980). Sens și definiție lexicografică "Studii și cercetări lingvistice", an XXXI, 465, 1980.
- [6] Em. Vasiliu (1982/1983). Adevăr analitic și definiție lexicografică "Analele științifice ale Universității "Al. I. Cuza" din Iași", secțiunea III, tom XXVIII/XXIX, 1982/1983.

Gramaticile generative nontransformationale

Emil IONESCU
Universitatea București, Facultatea de Litere
Str. Edgar Quinet nr. 5-7,
Email: eionescu@racai.ro

Acest articol este o prezentare generală a gramaticilor generative nontransformationale (GNT) și a prezenței lor în cercetarea lingvistică din România. În prima secțiune a articolului este descrisă geneza acestei direcții, iar în secțiunile a doua, a treia și a patra se menționează principalele realizări științifice și existența instituțională ale curentului. Partea a cincea este consacrată descrierii și pașilor care au dus la pătrunderea acestor gramatici în mediile academice și științifice noi. Concluziile articolului se vor a fi o pledoarie în sprijinul eforturilor în această direcție în cultura științifică românească.

1. Gramaticile generative nontransformationale: apariție

Gramaticile generative nontransformationale reprezintă o direcție în lingvistica formală contemporană, o direcție extrem de influentă și dinamică. Istoria acestei direcții este, desigur, mai recentă decât cea a gramaticii generative în general, dar face parte dintr-o tradiție generativismului din care face parte. Este însă o istorie deja bogată și diversă. Printre altele, diversitatea se exprimă și prin faptul că suntem obștinși să vorbim despre *gramatici* și nu despre o gramatică nontransformațională, pu-

Putem plasa începuturile acestei istorii la cumpăna dintre anii 1950 și 1960. Sunt anii când programul gramaticii universale al lui Noam Chomsky a ajuns la punctul să depășească starea de impas atinsă prin faza denuțită de mișcării "teoria standard". Privită din perspectiva prezentului, lucrarea lui Chomsky ("Lectures on Government and Binding") tocmai aceiași lucru înseamnă: depășirea crizei prin propunerea unui model nou de gramatică generativă.

Punctele în care gramatica universală este reformulată în ceea ce privește "Government and Binding" (GB) nu sunt puține și nici neînsemnate. Una dintre cele mai importante modificări a fost operată într-una din componentele care au fost cele mai mari speranțe: componenta transformărilor. Formulată su-

conceptului de transformare în cadrul modelului GB înseamnă două lucruri: simplificare și îngrădire. Simplificare, deoarece marea varietate de transformări se reduce acum la o singură operație: deplasarea unui constituent oarecare a. Și îngrădire, pentru că deplasarea nu se poate produce oricum, ci numai în condițiile în care anumite reguli foarte generale, numite *principii*, sunt respectate.

Nu toți adepții generativismului au fost însă mulțumiți cu noua propunere. Ceea ce s-a reproșat a fost că transformările rămăneau mai departe mecanisme prea puternice - în ciuda îngrădirilor și a simplificărilor - deoarece ele operau pe un domeniu prea larg: cel al structurilor sintactice. O altă obiecție viza temeiurile mentale ale operației de deplasare: în ciuda plauzibilității aparente a acestei ipoteze, nu există dovezi - susțineau criticii - că mintea implicată în utilizarea limbajului ar face uz de o astfel de operație. În sfârșit, existau cercetători care considerau că noul model de gramatică universală era greoi din punct de vedere computațional, tocmai din cauza operației de deplasare: anume, pentru fiecare deplasare de constituenți, este necesară o verificare a compatibilității dintre principii și deplasarea constituentului.

În ansamblu, divergențele legate de conceptul de transformare au pregătit cea mai mare ruptură pe care a cunoscut-o în istoria sa curentul gramaticii universale. Criticii radicali ai conceptului de transformare au propus renunțarea la acest mecanism, propunere pe care Chomsky și cei ce l-au urmat nu au acceptat-o niciodată. Începând cu anul 1981, ruptura se oficializează. Apar pe rând Gramatica Lexico-Funcțională (LFG - Bresnan și Kaplan), Gramatica Sintagmatică Generalizată (GPSG - Gazdar, Klein Pullum și Sag), Gramatica Arborilor Adăugați (TAG - Joshi), Gramatica Centrilor de Sintagmă (HPSG - Pollard și Sag), Gramaticile Categoriale de Unificare (CUG- Uzokreit)

2. Caracteristicile GNT

Dincolo de varietatea lor, gramaticile nontransformationale au un set de trăsături comune:

- Exploatează în mod generalizat reprezentările în termeni de trăsături
- Fac recurs la mecanismul unificării
- Se bazează pe constrângeri
- Sunt gramatici lexicaliste
- Au adecvare computațională

2.1. Reprezentări: structurile de trăsături

Reprezentările în termeni de trăsături sunt bine cunoscute în lingvistica modernă, datorită fonologiei și semanticii structurale. GNT au meritul de a fi generalizat această notație la scara întregii teorii lingvistice. Prin perechea

trăsătură (atribut)-valoare, orice fel de informație lingvistică morfologică, sintactică semantică, pragmatică - își găsește adecvată. Câteva exemple: notația [P(arte de)V(orbire): nume] desemnează o anumită entitate lingvistică este un nume. Reprezentarea [F(gerunziu)] precizează că avem a face cu un verb la gerunziu. [RAM(ură): v(aloare)n(on)v(idă)] spune că obiectul lingvistic în structură internă și este prin urmare o sintagmă. Este ușor de remarcat că atribut-valoare aplică principiul general al funcțiilor: unui atribut corespunde o anumită valoare, întocmai cum unui argument dat unei funcții corespunde o anumită valoare, datorită unei legi specifice de corespondență. Reprezentările care se face uz în GNT sunt denumite *structuri de trăsături*.

2.2 Unificarea

GNT se mai numesc și gramatici de unificare. Unificarea este operația de structurare a structurilor de trăsături. Unificarea a două structuri de trăsături A și B este structura minimală de trăsături care cuprinde în același timp structurile A și B. Dacă o astfel de structură nu există, unificarea "eșuează" (ceea ce înseamnă că unificarea verifică așadar compatibilitatea dintre două structuri de trăsături și produce o structură rezultantă care conține toată informația din structurile de unificării, iată câteva exemple:

(1) [CAT: det] u [CAT: nume] = ± (eșec)

CAT: det

(2) [CAT: det] u [ACORD: [NUM: singl=

ACORD: [NUM

(3) [CAT: nume] u $\left| \begin{array}{l} \text{ACORD:} \\ \text{GEN: mase} \\ \text{NUM: sing} \end{array} \right| \left| \begin{array}{l} \text{CAT: nume} \\ \text{ACORD:} \end{array} \right|$

Operația de unificare din primul exemplu eșuează pe motiv că structura rezultantă ar trebui să conțină atributul CAT cu două valori diferite (nume și det). Unificarea se realizează normal în (2) și (3), și produce structuri de trăsături complexe.

Se poate remarca faptul că rolul unificării este acela de a organiza informația care este corect în variate compartimente de limbă. Dacă are loc unificarea de informații fonologice, aceasta explică un aspect al corectitudinii fonologice în limbă dată. O unificare de informații morfologice dă seama de corectitudinea morfologică, ș.am.d. Nu e însă exclusă nici unificarea de informații diferite, de exemplu, semantice și morfologice, semantice și sintactice.

2.3. Constrângeri

În exemplul (1) din paragraful precedent, unificarea eșuează deoarece nici o structură de trăsături nu poate avea valori diferite pentru același atribut. Aceasta este o "lege" inerentă unificării, tot astfel cum în logica bivalentă o "lege inerentă" este terțiul exclus. Se poate spune că (1) definește o limită a unificării și implică o constrângere asupra acestei operații. Constrângerea este de natură formală, pentru că derivă din natura însăși a unificării. Dar pentru scopurile unei teorii lingvistice, astfel de constrângeri nu pot fi suficiente. Polona, de pildă, face la verbele de persoana I deosebirea între verbele folosite de un bărbat și cele folosite de o femeie. Verbul are așadar gen în polonă, dar nu și în română. Pentru a face această diferență între cele două limbi trebuie să se admită că unificarea informației de gen cu cea de verb se poate face în polonă dar nu se poate face și în română. Numai că de această dată constrângerea privind unificările nu mai are temei formal. Nu se poate spune că în mod necesar verbul are sau nu gen. Unificările acestor informații sunt prin urmare "contingente", sau cu un alt termen, "empirice", tocmai pentru că ele nu derivă din natura însăși a operației. Gramatica unei limbi se descrie mai ales în termenii unificărilor "contingente".

2.4 Lexicalism

În teoriile contemporane ale gramaticii, lexicalismul este o opțiune privitoare la modul în care este concepută structura cuvintelor în relația lor cu sintaxa. Există teorii, precum GB, care consideră că procesul de constituire morfologică a cuvintelor are loc în sintaxă. În acest sens, GB este o morfosintaxă deoarece generalizează operația de deplasare la nivelul morfologiei înseși, prin mecanismul numit "deplasare centru-centru" (engl. "Head to Head Movement"). Gramaticile de unificare adoptă o strategie distinctă: ele consideră că procesele de constituire morfologică a cuvintelor sunt independente de sintaxă. În această perspectivă, rezultatul proceselor morfologice furnizează sintaxei inputul necesar: cuvintele gata formate. Modularizarea celor două componente ale gramaticii se dovedește preferabilă mai ales în cazul limbilor cu morfologie bogată.

Un alt aspect al lexicalismului asumat de GNT este ilustrat de modul în care sunt construite explicațiile de gramaticalitate. Explicațiile în GNT se sprijină în măsura posibilului (dar într-o măsură mult mai mare decât în alte teorii) pe proprietățile cuvintelor. În istoria generativismului, pasivul, de pildă, a fost considerat multă vreme o structură explicabilă *sintactic*, adică o construcție rezultată din transformări ale unei alte structuri sintactice. GNT afirmă însă că nu e nevoie să se recurgă la structuri sintactice anumite, deoarece toate elementele de care e nevoie pentru a explica o construcție pasivă pot fi codificate la nivelul cuvintelor'. Un tratament asemănător

Preferința aceasta pentru un compartiment de limbă în defavoarea altui compartiment, atunci când se pune problema mecanismelor care justifică o anumită construcție nu e înțeleasă încă nici azi de unii lingviști. Este vorba de aceia care cred că a avansa o explicație lexicalistă atunci când există deja una sintactică pentru un fenomen oarecare

poate fi observat în cazul dependențelor la distanță, sau în cel al ridicării (engl. "raising"), unde rolul unităților lexicale în determinarea acestor construcții este de asemenea semnificativ.

2.5 Adecvare computațională

În lingvistică, o teorie este considerată adecvată, dacă este aplicabilă în domeniul de fapte pentru care este construită ca o explicație. O morfologie de limbă, de pildă, este adecvată dacă prin regulile propuse dă seamă de morfologie corecte ale limbii supuse analizei.

Acest principiu foarte general a fost nuanțat de către Chomsky. Este deja celebră: pornind de la ideea că utilizarea limbajului este o activitate minții omenești, Chomsky a susținut că o teorie trebuie socotită adecvată pentru că produce explicații ale cazurilor de corectitudine, și nu pentru că mecanismele utilizate sunt dovedite (sau cel puțin presupuse) a fi în concordanță cu mintea omenească. Quine afirmase că dacă avem două gramatici care descriu diferite explică aceeași realitate lingvistică, nu există criterii simple pentru alegere a uneia dintre ele. Chomsky a replicat că un astfel de criteriu ar fi el fiind măsura în care fiecare dintre aceste gramatici se folosește în mod cunoscut ca aparținând minții în procesele ei cognitive.

Criteriul suplimentar formulat de Chomsky în evaluarea teoriei este că ea este apropiată de cea a psihologilor și psiholingviștilor. Criteriul este a apropiat comunitatea generatiștilor de cea a psihologilor și psiholingviștilor. Cercetările de psiholingvistică. S-au obținut rezultate interesante și au susținut ipoteze neașteptate. De pildă, regulile de constituenți sînt socotite adecvate operațiuni cu mare probabilitate de a fi folosite de inteligența omenească. Recursivitatea este și ea considerată a fi o proprietate de care inteligența omenească face uz în utilizarea limbajului.

Criteriul lui Chomsky a condus însă și la cercetări cu rezultate neașteptate. Judecat. De pildă, despre realitatea psihologică a *urmelor*, concepute în teoria GB, s-a argumentat și pro și contra, și este foarte dificil ca o teorie să poată lua o poziție.

Un lucru este cert totuși în evoluția raporturilor dintre teorie și realitate ei psihologică: comparativ cu faza de început, interesul psiholingviștilor față de ipotezele venite din comunitatea "chomsky" a crescut semnificativ. A crescut însă interesul psiholingviștilor pentru ipotezele psihologice ale lumii inteligenței artificiale. Este celebră în acest sens ipoteza de adecvare a cunoștințelor lexicale a lui Quillian, care a atras atenția în mod special

înseamnă doar a propune variațiuni pe aceeași temă. Diferențele sunt în măsura în care și privesc mecanismele cognitive angajate în utilizarea limbajului. Este interesant că procesarea unităților lexicale este mai ușor de efectuat decât unele procesuri structurale sintactice. Acest fapt oferă un criteriu valoros de judecare a adecvării gramaticii privite din unghi cognitiv.

de psihologi și de psiholingviști. Un al treilea factor intra astfel în joc, rezultatul fiind că unele teorii lingvistice au devenit atente la operațiile și mecanismele utilizate de inteligența artificială. Erau exact teoriile generative netrtransformaționale. Consecința principală a acestei deplasări de interes a fost că teoriile în cauză au devenit accesibile utilizării automate. Cu alte cuvinte - și spre deosebire de gramaticile lui Chomsky - ele pot fi implementate computațional.

Vom numi adecvarea unei teorii la domeniul de fapte pe care îl abordează *adecvare lingvistică*. Măsura în care o teorie lingvistică aparține (sau poate fi presupusă a aparține) minții omenești definește *adecvarea ei psihologică*. Iar gradul în care ea este livrabilă inteligenței artificiale indică *adecvarea ei computațională*. Direcția actuală a curentului de idei pare să fie următoarea: legăturile și dialogul dintre psihologia cognitivă și inteligența artificială sunt într-o continuă creștere, astfel încât adecvarea computațională a unei teorii lingvistice are șanse mari să-i confere și adecvare psihologică. Pe această direcție sunt plasate gramaticile generative netrtransformaționale.

3. Realizări

Una dintre cele mai importante realizări ale gramaticilor nontransformationale îl reprezintă numărul mare de aplicații. O enumerare a limbilor supuse analizelor nu este posibilă aici, dar se poate preciza că aproximativ doua treimi din familiile de limbi (considerate în eșantioanele lor reprezentative) au fost analizate din perspectiva netrtransformațională. Este caracteristic acestor analize faptul că refuză deosebirea chomskyană centru-periferie ("core-periphery"). Ele se concentrează asupra varietății de date oferite de corpusuri.

Ceea ce este însă cel mai important sub aspectul realizărilor este faptul că GNT au reușit să producă replici viabile la analizele paradigmei dominante, cea chomskyană. O serie de fenomene gramaticale - privite de obicei ca fiind de la sine caracterizabile prin mecanismul deplasării constituenților - au primit în cadrul GNT analize alternative. Așa s-a întâmplat cu construcțiile pasive, cu fenomenul de ridicare (și mai general cu fenomenele de depedență limitată), cu construcțiile nonlocale (precum topicalizările, structurile relative și interogative). În această privință, GNT au continuat tradiția firească, inaugurată de structuralism, tradiție constând în regândirea fenomenelor de limbă odată cu fiecare nouă școală lingvistică.

4. Forme instituționale de susținere

GNT sunt bine reprezentate instituțional. Ele și-au făcut loc în primul rând în programele curriculare ale unor universități de prestigiu, precum Universitatea

Stanford, Universitatea Statului Ohio (Columbus), Universitatea Tuebingen, Universitatea Saarbruecken, Universitatea Groningen, King's College din Londra, Universitatea Edinburgh, Universitatea Paris 7. Extensiile acestor programe curriculare sunt școlile de vară. O prestigioasă școală de acest fel ("European Summer School in Logic Language and Information" - ESSLLI) este organizată anual din 1989, cu rolul de diseminare a evoluțiilor și curentelor formate în interiorul gramaticilor netrtransformaționale. Este apoi de semnalat, în aceeași linie a "didacticii" gramaticilor nontransformationale, nou înființata școală de vară de la Konstanz (Germania).

În planul congreselor științifice, HPSG și LFG au de multă vreme propriile lor conferințe anuale. Iar un congres ținut o dată la doi ani - cel de gramaticile formale - urmărește să adune sub același acoperiș toate școlile aceleiași familii.

Până de curând, gramaticile nontransformationale nu au avut o revistă proprie. Lucrările însă au fost și sunt publicate în reviste de prestigiu, precum "Computational Linguistics", "Natural Language and Linguistic Theory", "Journal of Linguistics", "Language" sau "Langages". O revistă orientată explicit spre aceste gramatici este editată de puțină vreme la cunoscuta editură olandeză Kluwer. Este vorba despre revista "Grammars". De asemenea, pe lângă Centrul de Studii asupra Limbajului și Informației de la Universitatea Stanford există de mai multă vreme o deja celebră editură care publică lucrările esențiale ale domeniului.

5. Gramaticile nontransformationale în România

Prezența GNT în România poate fi discutată având în vedere două coordonate: cea a contribuțiilor științifice și cea a programelor curriculare.

Din primul punct de vedere, întâia contribuție (după cunoștința noastră, ce puțin) a venit din partea Adrianei Costăchescu ([14]). Adriana Costăchescu este autorul unui studiu, din perspectiva GPSG (teorie care a precedat și inspirat HPSG), asupra relației dintre coordonarea adversativă și subordonarea concesivă. Studiul a fost elaborat în 1993 și publicat în 1996.

Lucrări de prezentare generală a diferitelor forme de GNT sau, dimpotrivă, de prezentare a trunchiului comun - unificarea - au fost publicate în ultimii șase ani de Adrian Atanasiu, Verginica Barbu, Ana-Maria Barbu, Florentina Hristea, Emilia Ionescu și Rodica Tătar.

Printre "pionierii" aplicațiilor acestor gramatici la limba română trebuie menționați Liviu Ciortuz și cercetătoarea italiană Paola Monachesi. Amândoi au folosit teoria HPSG. Rolul lui Monachesi în stimularea aplicațiilor de acest tip la limba română trebuie în mod special subliniat. Studiile sale asupra criticelor pronominale din română au determinat o "mobilizare" a energiilor câtorva

cercetători români. Este vorba despre Ana-Maria Barbu, Emil Ionescu și Amalia Todirașcu.

Ana-Maria Barbu a aplicat HPSG în analiza elementelor gravitând în jurul verbului - adverbul de negație, semiadverbele, auxiliarele - și a ajuns la concluzia că acestea sunt mai apropiate de afixe decât de cuvinte. Concluzia analizei se întâlnește cu concluzia exprimată în lucrarea Valeriei Guțu Romalo, "Morfologie structurală a limbii române", în care formele compuse ale verbelor sunt considerate forme cu afix mobil.

O alta contribuție a Anei-Maria Barbu privește ordinea constituenților în grupul nominal. Valorificând sugestiile de analiză ale lui Valerio Allegranza', Ana-Maria Barbu a propus o clasificare a constituenților grupului nominal, care este relevantă pentru problema ordinii acestora. Analiza produce astfel soluții clare și eficiente într-o problemă complicată de gramatică a limbii române.

Semnalând unele neajunsuri în analiza GB a fenomenului de anticipare clitică a complementului direct nominal în română, Verginica Barbu și Emil Ionescu propun o abordare alternativă HPSG. Analiza poate fi extinsă și la alte limbi care prezintă fenomenul în cauză. Analiza susține că pronumele neaccentuate nu au un comportament uniform, proprietățile lor depinzând de faptul dacă participă sau nu la structuri de dublare. Noutatea abordării vine din faptul că fenomenul anticipării obiectului direct este în mod ultim justificat prin proprietățile lexicale ale verbului tranzitiv.

Un fenomen care, în aparență cel puțin, implică recursul la mecanismul deplasării - este vorba de prezența pronomelor neaccentuate în acuzativ în contexte în care ele nu sunt subordonate față de vreun element din acel context - este tratat într-un alt studiu asupra cliticelor pronomiale românești² (). Studiul arată că ipoteza deplasării constituenților nu este necesară în analiza fenomenului. Este propusă în alternativă o analiză fără deplasări care captează toate proprietățile fenomenului.

O analiză HPSG este propusă de asemenea pentru fenomenul negației duble și multiple în română [23]. În sfârșit, Amalia Todirașcu abordează într-unui din studiile sale asupra limbii române, o categorie de dependențe limitate (așanumitele *tough-constructions*), din aceeași perspectivă HPSG.

În aceeași linie a contribuțiilor științifice, merită amintită o inițiativă instituțională: acreditarea de către CNCSIS, în anul 2001, a Centrului de Lingvistică Computațională de pe lângă Facultatea de Litere. Centrul este perechea universitară a Centrului de Studii Avansate în Inteligență Artificială. Apariția sa a fost semnalată în buletinul european ELSNEWS. Unul dintre programele de cercetare pe anul 2002 ale centrului are în vedere dezvoltarea aplicațiilor de gramatici netrăsformationale la limba română.

În engleză, fenomenul este cunoscut sub numele de clitic climbing", și este ilustrat în română de structuri de tipul Nu-l pot suferi pe Ion.

În planul programelor curriculare, GNT și-au făcut loc mai gîntâmpinate uneori nu doar cu neîncredere, ci și cu ostilitate. A fost o fericire un sprijin substanțial și constant al factorilor de decizie. Decanul Facultății de Litere, acad. prof. Dan Horia Mazilu, la rectorul București, prof. dr. Ioan Mihăilescu, la prorectorul aceleiași instituții Pânzaru, și la acad. Dan Ioan Tufiș, directorul Centrului de Studii Inteligență Artificială al Academiei Române, cărora autorul acestor lucrări exprimă via și profunda sa grațitudine, pentru susținerea pe care au acordat în inițiativele sale. Mulțumită acestui sprijin, au devenit realitate unele lucrări care pot fi considerate succese:

- În programa cursurilor opționale de limbă pentru anul 2000, la Facultății de Litere a fost introdus în 1996 un curs de gramatică HPSG, iar din 1997 pînă în 2001 s-a ținut un curs de gramatică de unificare cu referire specială la HPSG.
- Din 1999, se predă la Facultatea de Matematică a Universității București un curs opțional de un an de prelucrare a limbii naturale, în care un loc important îl ocupă gramaticile de unificare.
- Din 1997 pînă în prezent masteratul de lingvistică teoretică la Facultatea de Litere din cadrul aceleiași universități găzduiește un curs de semestru de teorie HPSG aplicată la limba română.
- Din 1999, același masterat oferă un seminar de gramatică de unificare implementare computațională.
- În anul 2000, un proiect de dezvoltare a componentei de unificare computațională în cadrul masteratului de lingvistică teoretică a beneficiat de sprijin de finanțare din partea Băncii Mondiale și a Uniunii Europene. României, sprijin care a făcut posibil printre altele organizarea unor cicluri de conferințe pe teme de GNT (în special HPSG) la Facultatea de Litere a Universității București. Au conferențiat Ivan Sagot (Facultatea Stanford), Anne Abeille și Daniele Godard (Universitatea Paris 7), Ștefan Muller (Universitatea din Jena), Robert Malouf (Universitatea Groningen), Howard Gregory (King's College, Londra), și Hans-Joachim Lenz (Universitatea Tübingen), toți fiind personalități renumite în domeniul. Mulțumită aceluiași program, cercetătorii din România au petrecut stagii de specializare la universitățile din Lille și au putut participa la manifestări reprezentative, cum ar fi conferința UNESCO asupra spațiilor virtuale și multilingvismului (aprilie 2001), colocviul de gramatici bazate pe unificare la Trondheim (august 2001), sau congresul de prelucrare a limbilor naturale de la Tokyo, (noiembrie, 2001). Cea mai importantă realizare legată de acest program, a constat însă în promovarea mobilității studențești, concretizate în vizitele de studiu la masteratul de lingvistică teoretică, la universitățile din Tübingen, Paris 7 și Siena.

6. Concluzii

Deși GNT au pătruns în mediile științifice din România mai târziu decât în alte țări, faptul că ele sunt prezente la noi este un lucru încurajator. Există tentația de a privi aceste eforturi de sincronizare cu mișcarea de idei din domeniul lingvisticii formale drept tentative mimetice și superficiale. Este o greșeală gravă. Diversele comunități de lingviști pot desigur ignora un curent, precum cel prezentat mai sus, dar aceasta este o atitudine, pentru a spune așa, pe proprie răspundere. GNT și teoria lingvistică pe care ele au inspirat-o și-au făcut deja loc în lingvistica zilelor noastre și au devenit una din paradigmele majore. În plus, dubla deschidere a acestor gramatici către psihologia cognitivă, pe de-o parte, și către inteligența artificială, pe de altă parte, recomandă această paradigmă drept cadrul *privilegiat* de dialog interdisciplinar din științele umaniste ale contemporaneității. Din acest triunghi, sunt așteptate să apară noi aplicații - unele au și apărut deja - care vor extinde într-un mod neașteptat conceptul de lingvistică aplicată. Pentru toate aceste motive, tentativele de a păstra un contact viu și de perspectivă cu comunitatea științifică internațională a GNT reprezintă o investiție sigură pe termen lung.

Bibliografie

- [1] Abeille, A. *Les nouvelles syntaxes. Grammaires d'unification et analyse du français*, Armând Colin, Paris, 1993
- [2] Atanasiu, A. *Curs de lingvistică matematică*, Editura Universității București, 1998
- [3] Barbu, A.M. *Gramatici categoriale. Studiu comparativ cu gramaticile de constituenți*, "Limba Română", XLVI, 4-6, p 239-252, Ed. Academiei, 1997
- [4] Idem, *Complexul verbal*, "Studii și Cercetări Lingvistice", Ed. Academiei, sub tipar.
- [5] Idem, *Romanian Determiners: Order and Classification*, "Revue Roumaine de Linguistique", Ed. Academiei, sub tipar
- [6] Idem, *Funcțiile sintactice în Teoria X-Bară*, "Studii și Cercetări Lingvistice", Ed. Academiei, sub tipar Barbu, A.M. și E. Ionescu *Teorii gramaticale contemporane: Gramatica Centrilor de Sintagmă*, "Limba Română" 1 1996 31-55
- [7] Idem, *Accusative Clitic Doubling in Romanian*, Liviu Ciortuz, Paola Monachesi, Hans Uszkoreit (editori) "Informai Proceedings of the GE&GL Workshop: Grammar Engineering and Grammar Learning", Tușnad, România, 1997

- [8] Barbu, V. *Despre gramaticile de unificare*, Analele Universității București, seria limbă și literatură română, 2001, p. 45-52
- [9] Barbu, V. și E. Ionescu *Anticiparea complementului direct în limba română în perspectiva HPSG*, Lucrările colocviului "Perspective moderne asupra limbii române", București, Editura Universității din București, (sub tipar)
- [10] Borsley, R. *Syntactic Theory: A Unified Approach*, Edward Arnold, London, 1991
- [11] Bresnan, J (editor) *The Mental Representation of Grammatical Relations*, MIT, Press, Ca. Mass, 1982
- [12] Ciortuz, L. *An HPSG Kernel for Romanian*, manuscris, 1996
- [13] Ciortuz, L, P. Monachesi, și H. Uszkoreit (editori; *Informai Proceedings of the GE&GL Workshop: Grammar Engineering and Grammar Learning*, Tușnad, România, 1997
- [14] Costăchescu, A. *"Coordination" adversative et "subordination" concessive*, Iliescu, M. și S. Sora, (editori), Rumänisch: Typologie, Klassifikation, Sprachcharakteristik, Mtinchen, 1996, p. 121-134
- [15] Gazdar, G, E. Klein, G. Pullum și I. Sag, *Generalized Phrase Structure Grammar*, Cambridge, Harvard University Press, 1985
- [16] Gerlach, B. și J. Grijzenhout (editori) *Clitics in Phonology, Morphology and Syntax*, John Benjamins Publishing Company, Amsterdam / Philadelphia, 2000
- [17] Hristea, F. *Introducere în procesarea limbajului natural cu aplicații în PROLOG*, Editura Universității București, București, 2000
- [18] Iliescu, M. și S. Sora, (editori), *Rumänisch: Typologie, Klassifikation, Sprachcharakteristik*, Mtinchen, 1996, p. 121-134
- [19] Ionescu, E. *A Type of SOV Construction in Romanian*, "Cahiers de Linguistique Theorique et Appliquee", tomes XXXII-XXXIII, 1995-1996, 19-39
- [20] Idem, *Accusative Weak Pronouns in Romanian*, "Cahiers de Linguistique Theorique et Appliquee", tomes XXXII-XXXIII, 1995-1996, 19-39
- [21] Idem, *Accusative Clitic Doubling in Romanian*, "Cahiers de Linguistique Theorique et Appliquee" tomes XXXII-XXXIII, 1995-1996, 53-73
- [22] Idem, *Accusative Clitic Climbing in Romanian*, "Cahiers de Linguistique Theorique et Appliquee", tomes XXXII-XXXIII, 1995-1996, 74-87
- [23] Idem, *A Quantification-based Approach to Negative Concord in Romanian* in Geert-Jan M. Kruijff and Richard T. Oehrle (editori), *Proceedings of Formal Grammar Conference Utrecht*, 1999, p. 25-36
- [24] Idem, *pro-Drop: An HPSG Account without Lexical Rules*, "Bucharest Working Papers in Linguistics", voi. I, nr.1, 1999, 117-124

- [25] Idem, *On the Status of PE in the Direct Object Construction in Romanian*, Romanian Journal of Information Science and Technology, volume 4, numbers 3-4, 2001, p. 293-310
- [26] Joshi, A. Introduction to Tree Adjoining Grammar, *Manaster Ramer, A. (editor)* The Mathematics of Language, John Benjamins, Amsterdam, 1987, p. 87-114
- [27] Kruijff, G.-J. M. and R. T. Oehrle (editors), Proceedings of Formal Grammar Conference, Utrecht, 1999
- [28] Manaster Ramer, A. (ed.) *The Mathematics of Language*, John Benjamins Publishing Company, Amsterdam, 1987
- [29] Monachesi, P. *Clitic Placement in the Romanian Verbal Complex*, Gerlach and Grijzenhout (2000), p. 255-294.
- [30] Pollard, C. și I. A. Sag, *Information-based Syntax and Semantics*, CSLI, University of Chicago Press 1987
- [31] Idem, *Head-driven Phrase Structure Grammar*, The University of Chicago Press, Chicago, 1994
- [32] Shieber, St. *An Introduction to Unification-based Theories of Grammar*, CSLI, University of Chicago Press, 1986
- [33] Tătar, D. *Inteligență artificială*, Editura Albastră, Cluj, 2001
- [34] Todirașcu, A. *Romanian Tough-Constructions*, Ciortuz, L, P. Monachesi, și H. Uszkoreit (editors; *Informai Proceedings of the GE&GL Workshop: Grammar Engineering and Grammar Learning*, Tușnad, România, 1997
- [35] Wood, M. McGee, *Categorial Grammars*, Routledge London and New York, 1993

Către o teorie X-bar funcțională

Neculai CURTEANU

Institutul de Informatică Teoretică, Academia Română, Filiala Iași
curteanu@iit.tuiasi.ro

1. Teorii X-bar mai vechi și mai noi

Scopul prezentei lucrări este dublu: (a) de a propune o nouă X-bar schemă, numită X-bar schemă *funcțională* și *recursivă* (pe scurt, FX-bar schemă), mai generală și mai adecvată decât cele existente, care să satisfacă cerințele unei abordări funcționale a limbajului natural (LN), în particular, ale strategiei lingvistice SCD (Segmentare-Coeziune-Dependență) [1], [2], și (b) de a pune în evidență faptul că teoria FX-bar propusă poate reprezenta o posibilă (și necesară) soluție la următoarea problemă ridicată de Noam Chomsky în teoria Minimalist Program [3]: în două capitole diferite, Chomsky afirmă (în două abordări diferite, aparent contradictorii, asupra structurii sintactice a LN) atât importanța crescândă a teoriei X-bar cât și posibilitatea ca teoria X-bar standard să fie "*largely eliminated in favor of bare essentials*" (vezi secțiunea 5).

1.1. Teoria X-bar clasică

Printre (sub)teoriile care reprezintă substanța majoră pentru câteva teorii formale importante asupra sintaxei (LN), un rol fundamental este jucat de către așa-numita teorie X-bar. X-bar schemele propuse sunt de obicei însoțite de definiții, ipoteze, restricții, principii și alte (sub)teorii gramaticale care specifică într-o cât mai mare măsură modul concret în care X-bar schemele sunt utilizate pentru a construi structurile sintactice de bază ale LN. În general, teoria X-bar stabilește categoriile gramaticale principale, proiecțiile lor lingvistice (minimale și maximele), relațiile de dominare dintre categorii în cadrul acestor proiecții, sub-, co-, sau supra-ordonarea lor. Toate aceste aspecte asigură numai coloana vertebrală (infrastructura) consistentă a structurii sintactice în reprezentarea LN. Un capitol de o importanță deosebită este relația dintre teoria X-bar și alte sub(teorii) sintactice și semantice care formează întregul corpus al unei anumite teorii lingvistice.

Prima formă a X-bar teoriei este propusă de către Noam Chomsky în lucrarea *Remarks on Nominalizations* (1970) [4]. Chomsky scoate în evidență diferențele reale existente în următoarele sintagme nominale:

(1.1) *John's criticism of the book*

(1.2) *John's criticizing the book;*

În special datorită șablonului verbal (similar cu al verbului "criticize") rezultat din gerunziul nominal (pentru engleză) "criticizing", în comparație cu forma nominală derivată "criticism".

Teoria X-bar originală propusă de Chomsky identifică trei categorii lexicale primitive, N [Eng: *noun*], V [Eng: *verb*] și A [Eng: *adjective*], fiecare dintre ele cu câte două categorii sintagmatice corespunzătoare. Mai exact, utilizând notația X = N, V, A, categoria gramaticală X se întâlnește ca *nucleu* [Eng: *head*] într-o categorie intermediară X' (sau X1, sau X'), tradițional numită X-bar, precum și într-o categorie maximală X" (sau X2, sau X²), tradițional numită XP, reprezentând *proiecția maximală* a categoriei gramaticale X (lexicală sau nelexicală). Categoria X este numită *nucleul sintagmelor* X' (sau X1) și X" (sau X2) care o conțin. Să mai notăm că prescurtarea pentru categoria *prepozițională* este P.

Ulterior au fost considerate *patru* categorii lexicale, bazate pe următoarele combinații ale celor două trăsături N și V (considerate ca fiind generice pentru categoriile lexicale):

N	este o categorie X cu trăsăturile [+N, -V];
V	este o categorie X cu trăsăturile [-N,+V];
A	este o categorie X cu trăsăturile [+N, +V];
P	este o categorie X cu trăsăturile [-N, -V].

Teoria X-bar poate fi înțeleasă și ca o specificare a modalității în care unele categorii gramaticale sunt *dominate* de către altele, deci ca o teorie a dominanței gramaticale (sau, așa cum spune Chomsky, a "*guvernării*"), care arată cum un nucleu (sau o categorie lingvistică) X se proiectează (se extinde) către categoriile mai complexe (structurile sintagmatice) X' (sau X1) și X" (sau X2, sau XP). Structurile sintactice X1 și X2 devin categorii gramaticale esențiale ale organizării și reprezentării textului în LN.

Deci, X-bar teoria clasică consideră că X, împreună cu o secvență de *complemente* (sau *argumente*, notate Argj) este imediat dominată de X1, în timp ce X1 împreună cu o secvență de *specificatori* (notată Spec) este imediat dominată de către X2 (sau. XP). Utilizând binecunoscutele notații din domeniul teoriilor lingvistice formale, (X' = X1, X" = X2 = XP), categoriile lexicale și gramaticale ale teoriei X-bar clasice a lui Chomsky sunt următoarele:

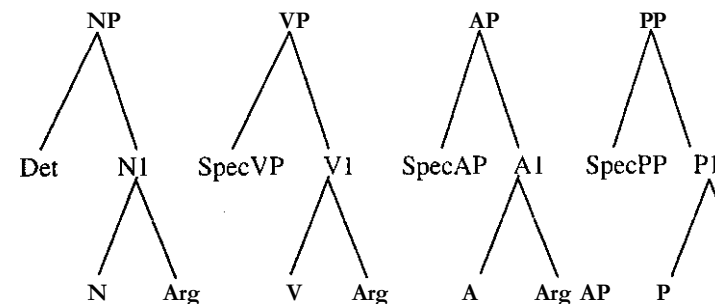


Figura 1.1. Proiecțiile categoriilor lexicale din teoria X-bar clasică

1.2. Extinderea teoriei X-bar la categorii non-lex

Stowell [5] propune ca teoria X-bar clasică să fie extinsă la categorii nelexicale sau *funcționale*. În particular, categoria gramaticală S (sau S1, în română: *frază*), care corespunde uneia sau mai multor propoziții gramaticale, este văzută ca I2 sau IP, deci ca proiecția maximală a categoriei nelexicale I (sau I1, sau I2). Nucleul nelexical I (INFL) reprezintă mulțimea de flexionare atribuite nucleului lexical al clauzei-matrice (propoziția sau clauza-matrice, sau chiar una regentă) dintr-o frază, așa cum sunt timpul, aspectul etc. în fraze. Remarcăm *categoria* S, care introduce un anumit grad de complexitate în analiza gramaticală, atât în engleză cât și în română. Termenul de *realitate lingvistică* codificată de categoria S ar trebui să fie acoperit de categoria "gramaticală" pentru engleză [Eng: (*grammatical*) *clause*], și de categoria "gramaticală" pentru limba română, cu două sorturi principale de clauze: clauza prescurtată CLF sau mai simplu CL, și *clauză infinită*, prescurtată CI.

Astfel în extensia nelexicală a teoriei X-bar, S este proiecția maximală a categoriei virtuale (nelexicale) I, în timp ce S1 este văzută ca CP sau CP, unde nucleul C este un *complementizator*, o categorie gramaticală care corespunde unei expresii (unui delimitator) sau unei sintagme care introduce o clauză subordonată, e.g. pronume relativ, conjuncție, locuțiune conjuncțională. Teoria X-bar extinsă acreditează următoarele structuri:

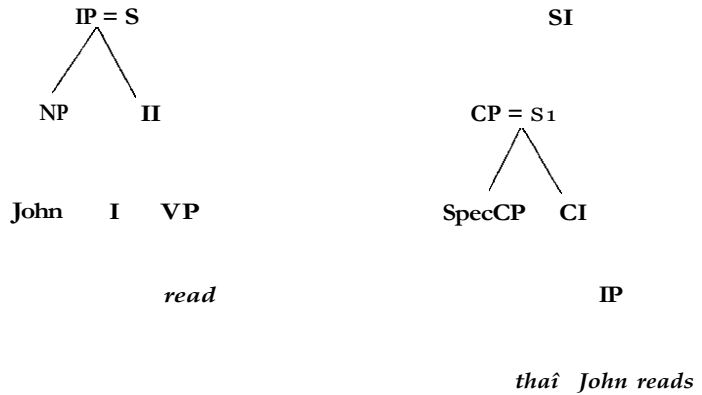


Figura 1.2. Teoria X-bar extinsă la categorii nelexicale

Sunt necesare câteva remarci:

(a) Teoria X-bar extinsă utilizează terminologia de "categorii nelexicale (sau funcționale)", prin care Stowell, Chomsky și alți lingviști definesc *noile nuclee* ale structurilor sintactice considerate. Categoria virtuală "I" este, desigur, una nelexicală, și susține o anumită funcționalitate depinzând de categoria lexicală căreia îi este atribuită. Categoria C nu este, de obicei, nelexicală (exceptând situația, posibilă, când ea lipsește) deoarece C corespunde unor categorii gramaticale lexical nevide. În ceea ce privește funcționalitatea lui C, suntem de acord că C corespunde într-adevăr unor funcții și relații sintactice și semantice importante pe care le numim *marcheri de propoziție* [1], [2], [6], uneori incluși în clase mai largi cum sunt cea a *marcherilor de discurs* [7], reprezentând în același timp și un element (deci o relație) de co-referință în cadrul fenomenului de *legare*, și/sau o "barieră" [8] în cadrul *teoriei limitării* [9]. Aceste aspecte multi-funcționale ale categoriei C nu sunt contradictorii ci doar complementare, întregind un tablou complex al funcționalității lexical-semantice pentru o categorie lingvistică atât de specială cum este C.

(b) A doua observație este dedicată rolului unor categorii nelexicale în cadrul X-bar schemelor extinse. Din Fig. 1.2. reiese că subiectul NP are rolul (nesigur) al unui specificator pentru S = IP, în timp ce VP reprezintă complementul categoriei virtuale I. De asemenea, S1 = CP se consideră a fi proiecția maximală a categoriei C, în timp ce complementul sintagmei CP este IP. Admițând că în engleză, din punct de vedere sintactic, această supoziție are sens deoarece categoria C reprezintă nucleul acestor sintagme, în alte limbaje, inclusiv româna, acest lucru este nedecis, în special din perspective semantice și funcționale. Unele abordări funcționale ale acestor probleme sunt discutate în mai multe lucrări, dar

ne vom restrânge să menționăm aici soluțiile oferite de către *teoria gramaticii funcționale* [10] și *strategia lingvistică SCD* [1], [2], [6]. Un interes special prezintă abordarea *lexicală* (inclusiv *funcțională*) a teoriei X-bar ca subteorie de bază în cadrul teoriei sintactice HPSG [Eng: *Head-driven Phrase Structure Grammar*] [11]. O analiză comparativă cu *FX-bar schema* propusă în această lucrare va fi făcută într-o lucrare viitoare.

1.3. X-bar schemele din teoria GB

X-bar schemele propuse de teoria *Government and Binding* (GB) a lui Chomsky [5] sunt următoarele:

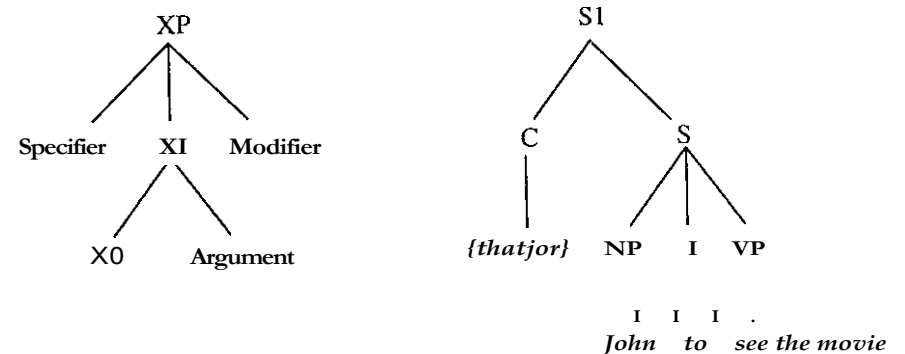


Figura 1.3. X-bar schema generală din GB, X = N, V, A, P, S

în teoria GB există următoarele X-bar *echivalențe* pentru proiecțiile categoriilor gramaticale (lexicale și nelexicale).

Tabelul 1.3

X	XI	X2
N	NI	NP
V	VI	VP
A	AI	AP
P	PI	PP
I	S	SI

în lucrările GB [5] și cele care urmează, Chomsky consideră categoria I ca fiind nucleul lui S, iar complementizatorul C ca fiind nucleul lui S1. În subsecțiunea următoare teoria sintactică GPSG a lui G. Gazdar [12] face un important pas

înainte către lexicalitate și către utilizarea explicită a trăsăturilor lingvistice atribuite categoriilor gramaticale.

1.4. Teoria X-bar în GPSG

În teoria lingvistică GPSG [Eng: *Generalized Phrase Structure Grammar*] [12], [13] etc, (sub)teoria X-bar joacă de asemenea un rol central, o sintagmă a LN fiind definită ca proiecția *trăsăturilor lingvistice* atribuite *nucleului* [Eng: *head*] acelei sintagme. Informația cuprinsă în trăsăturile nucleului determină caracteristicile principale ale comportamentului sintactic al sintagmelor LN. Reamintim că o categorie sintactică în GPSG se reprezintă ca o mulțime de perechi <trăsătură, valoare>. De exemplu, eticheta NP [Eng: *noun phrase*] (sau N2), prin care se notează o sintagmă nominală, reprezintă o abreviere pentru mulțimea {<N, +>, <V, ->, <BAR, 2>}, unde BAR este *numele trăsăturii* ce codifică *nivelul de proiecție* a categoriei sintactice N = {<N, +>, <V, ->}. Trăsătura BAR poate lua valorile 0, 1, 2. Teoria GPSG consideră N, V, A și P ca fiind *categoriile sintactice majore*. Toate celelalte sunt considerate de GPSG ca fiind *categoriile minore*: determinatori, complementizatori, marcheri, cuantificatori, alte particule etc. Categoriile majore sunt considerate de către teoria GPSG ca având întotdeauna o *valoare* pentru *trăsătura* BAR. Valoarea BAR pentru categoriile minore nu este definită niciodată în GPSG.

Teoria sintactică a GPSG aduce câteva elemente noi și interesante comparativ cu teoria GB: **(a)** X-bar schemele au, ca și în GB, trei nivele de proiecție (valorile trăsăturii BAR); **(b)** Pentru economia reprezentării, GPSG propune ca în X-bar schemele de bază, nivelul proiecției lingvistice să fie conservat când se trece de la nucleu către expresiile subcategorizate, mai puțin în cazul în care acest lucru se face prin (alte) reguli explicite; **(c)** Printr-un mecanism de *moștenire implicită*, nivelele BAR de proiecție a nodului-rădăcină și ale nodurilor-fiice rămân aceleași, mai puțin în cazul în care există o indicație contrară expresă.

O altă caracteristică este aceea că în GPSG nu se întâlnesc categorii abstracte, non-lexicale, cum ar fi "I" (INFL) din GB. Acest lucru este posibil deoarece în GPSG, pentru aceste categorii nelexicale, nu există un nivel de proiecție pe care ele să fie reprezentate (sub nivelul lexical BAR = 0). Consecința este aceea că, în GPSG, S este proiecția unei categorii V. Mai exact, proiecțiile maximale ale lui V sunt VP, S, și S1, depinzând de următoarele valori luate de către trăsăturile SUBJ și COMP (= complementizator = C):

V[BAR 2][SUBJ-][COMP NIL] = VP;

V[BAR 2][SUBJ +][COMP NIL] = S;

V[BAR 2][SUBJ +][COMP a] = S1; unde a e {*that, for, whether, if*}.

În sfârșit, trebuie să remarcăm că GPSG trebuie să rezolve problemele întâlnite în mod obișnuit în formalismele gramaticale bazate pe unificarea lingvistică (și/sau logică), de exemplu PATR-II [14], HPSG [15], [16] etc. O astfel

de problemă este, în particular, transmiterea informației despre *timpul verbului* între forma flexionară codificată de verb și nodul S. Pentru teoriile lingvistice care permit inserarea în arborele de derivare a cuvintelor flexionate, așa cum este cazul cu GPSG, HPSG etc, informația despre forma flexionară trebuie să poată fi mutată în ambele direcții pe nivelele X-bar schemei. Din aceasta derivă, în GPSG, condiția ca V să fie *nucleul structurii clauzale* care corespunde categoriei S. Pe de altă parte, în GB, informația asupra timpului unui verb poate fi transmisă dinspre nodul I către proiecția sa în S înainte ca I să fie combinat cu forma flexionată a verbului din S. Această situație poate produce potențiale dificultăți procedurale și de reprezentare.

Este important de menționat că proiecțiile categoriilor din Tabelul 1.3 rămân aceleași pentru GPSG și LFG [Eng: *Lexical Funcțional Grammar*] (vezi de exemplu [13]), cu diferența notabilă că prima celulă din ultima linie a Tabelului 1.3 este goală, deoarece în aceste două teorii lingvistice (ca și în altele), categoria virtuală I lipsește.

1.5. O formulare recursivă a X-bar schemelor din teoria Tbarr

Vom propune în această subsecțiune o *formulare recursivă* a teoriei X-bar avându-și originea în *teoria barielor* (Tbarr) [8], [17] și fiind compatibilă cu teoria sintactică a *Programului Minimalist* (MinP) [3] și cu modelul său gramatical din *Principii și Parametri* (P&P) [3]. În conformitate cu MinP și P&P, gramaticile concrete ale limbajelor naturale (LN) reale pot fi modelate de mulțimi de parametri și valorile lor, care specifică principii și teorii lingvistice universale valabile. Pentru o asemenea setare (asignare) a valorilor parametrilor, relațiile de precedență (de ordonare liniară) dintre categoriile gramaticale sunt obținute din proprietăți ca marcarea cazuală, atribuirea de roluri tematice ((0 - roluri și 8 - marcheri), împreună cu alte relații și marcheri ce se aplică la nivelul sintagmelor, clauzelor, și unităților de discurs. Din acest motiv, relațiile de precedență pentru X-bar schemele propuse pot fi utilizate independent pe arborii sintactici considerați, informația de ordonare (liniară) a categoriilor fiind dată de următorii *parametri de precedență*.

(OrdPar) *Un anumit parametru (depinzând de limbaj) precizează dacă secvența de specificatori precede sau succede nucleul, iar un alt parametru (depinzând de limbaj) precizează când secvența complementelor precede sau succede nucleul din X-bar schemă.*

De exemplu, în *engleză*, specificatorii preced de obicei nucleele lor nominale, în timp ce în *română*, în mod normal, ei succed nucleelor. În general, complementele (argumentele) succed nucleele lor și în *engleză* și în *română*. Un caz special al argumentului este *subiectul* (sintactic). Această exprimare a (OrdPar) poate fi încă particularizată în funcție de categoriile lexicale concrete, din LN concrete. De exemplu, atât în *română* cât și în *engleză*, când o sintagmă

adjectivală (adverbială) este *predicațional activă*, fiind urmată de anumite argumente (complemente sau adjuncți), atunci este obligatoriu ca ea să succedă propriului nucleu și nu să îl precedă.

Consecința principală a parametrizării dependentă de limbaj a precedentei categoriilor lingvistice este că în exprimarea teoriilor lingvistice se pot utiliza arbori neordonați, iar principiile propuse de teoria X-bar primesc un puternic caracter de independență relativ la regulile de dominare ale structurilor sintagmatice. Este important faptul ca X-bar schemele obținute în cadrul teoriei X-bar considerate să asigure proiecții adecvate ale categoriilor lexicale, permițând inserarea *adjuncțiilor*, obținerea categoriilor de proiecție maximală, și acceptarea faptului că unele proiecții minimale sau maxime din *structura de adâncime* pot fi *vide* (deci noduri care să domine *categorii vide*), conform [9], [8], [17].

Fiind stabilit *principiul* (OrdPar), teoriile GB și Tbarr consideră următoarele trei nivele ale proiecției din teoria X-bar, sintetizate de următoarele *reguli* (*principii*) și de *X-bar schemele* corespunzătoare:

(PX0) Fiecare nod XO dintr-o schemă X-bar este fie *vid*, neavând nici o trăsătură, fie este nodul-mamă al unui element lexical a cărei categorie gramaticală și trăsături sunt specificate la nivelul lexiconului.

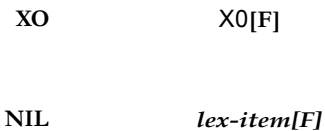


Figura 1.5.1. Nodul XO în TBarr

(PX1) Fiecare nod X1 (X' sau X') având trăsăturile lexicale F este fie nodul-rădăcină al exact unui nod X (care este *nucleu*) cu trăsăturile F și al unei secvențe de noduri XP (care sunt *complemente*, sau *argumente*), fie este rădăcina unui nod identic X1 împreună cu exact un nod XP (care este *adjunct*).

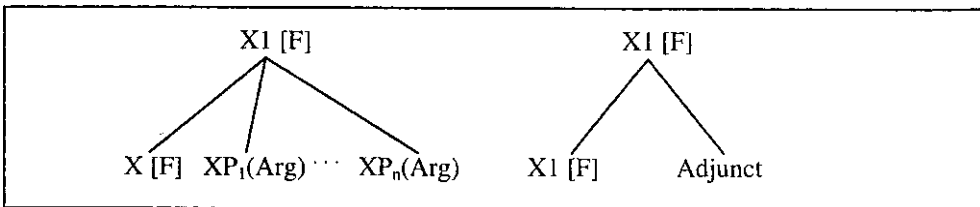


Figura 1.5.2. Nodul X1 în TBarr

(PX2) Fiecare nod XP care are trăsăturile lexicale F trebuie să satisfacă una și numai una din următoarele condiții: **(i)** XP este un nod-frunză (nu mai are nici un nod-fiică) și mulțimea F este *vidă*; **(ii)** XP este rădăcina unei secvențe de

XP's (*specificatori*) și a exact unui nod X1 moștenind trăsăturile F; **(iii)** XP este rădăcina unei secvențe de XP's (*complemente*, sau *argumente*) și a exact unui nod X cu trăsăturile F; **(iv)** XP este rădăcina unui alt nod XP moștenind trăsăturile F și a exact unui nod XP.

O *observație importantă* este aceea că unele dintre *secvențele* XP specificate în regulile (PX1) și (PX2) pot fi *vide*.

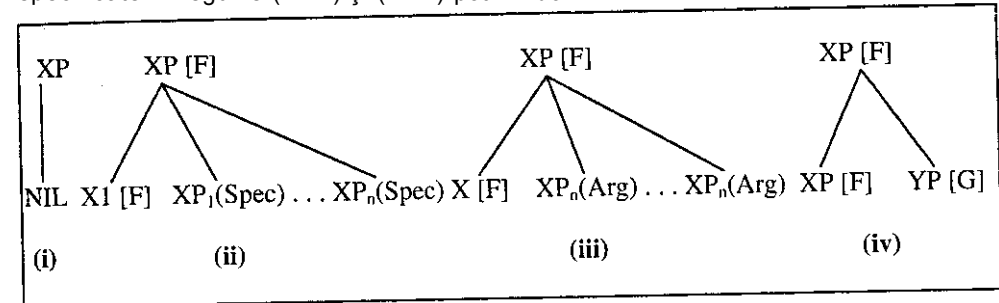


Figura 1.5.3. Nodul X2 în teoria TBarr

Combinând recursiv X-bar schemele rezultate din *regulile* (X_o)-(X_n)-(X_p2) se pot obține toate structurile sintactice întâlnite în *X-bar teoria clasică* și *extinsă*

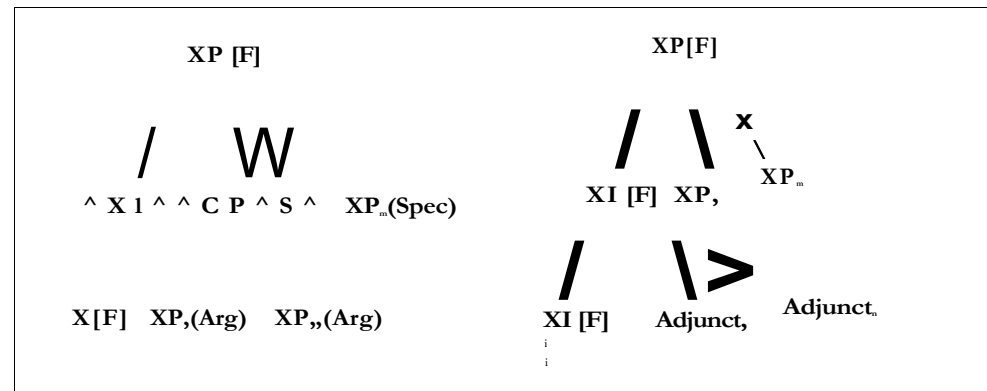


Figura 1.5.4. Formele generale (și recursive) ale X-bar schemelor din TBarr

2. X-bar teoria din modelul P&P al teoriei MinP

2.1 Sistemul Chomskyan al gramaticii universale

Această subsecțiune conturează câteva aspecte implicate de către teoria X-bar în cadrul teoriilor MinP (*Minimalist Program*) și P&P (*Principles and Parameters*) [3]. Pentru a înțelege contextul, este necesar să schițăm teoria lui Chomsky a gramaticii universale UG [Eng: *Universal Grammar*] și a relațiilor sale cu abordarea MinP bazată pe P&P. Sunt introduse următoarele concepte de UG.

Capacitatea utilizării și înțelegerii LN se bazează în esență pe proceduri care pot genera obiecte numite *descrieri structurale* (SDs). SDs sunt *expresii* de limbaj. Teoria unui LN particular constituie gramatica acestuia, în timp ce teoria tuturor limbajelor și a expresiilor pe care le generează ele reprezintă *Gramatica Universală* (UG).

Se consideră că UG specifică anumite *nivele lingvistice*, sau sisteme de reprezentare a informației lingvistice. UG a lui Chomsky [3] presupune că fiecare SD este o secvență (8, a, n, X) de patru reprezentări pe următoarele nivele, respectiv: *structură de adâncime* (D-structură), *structură de suprafață* (S-structură), *formă fonetică* (PF) și *formă logică* (LF). O ipoteză constructivă pentru UG este aceea că limbajul este scufundat în *sisteme de performanță* care permit ca exprimări în LN să fie folosite pentru articulare, interpretare, referire, interogare, reflecție și alte acțiuni, în timp ce SDs devin un complex de instrucțiuni pentru aceste sisteme de performanță.

O altă ipoteză standard pentru construcția UG este aceea că un LN este format din două componente: un *lexicon* și un *sistem computațional*. Această construcție este o inovație esențială comparativ cu teoria GB, care pretinde independența sa față de orice aspecte computaționale sau de implementare. Lexiconul specifică elementele de intrare pentru sistemul computațional, în timp ce acesta folosește intrările de lexicon pentru a genera derivări și SDs. Derivarea unei exprimări lingvistice particulare implică alegerea elementelor din lexicon și evaluarea, construind perechea pe două nivele de performanță, numite și *reprezentări de interfață*. Una din ipotezele de bază ale teoriei lui Chomsky *Minimalist Program* este aceea că în construcția SD, utilizând lexiconul și sistemul de evaluare, sunt luate în considerare *numai două* nivele de interfață, corespunzând lui PF (formă fonetică) și lui LF (formă logică), împreună cu mulțimile de perechi (n, X) rezultate din cele două forme.

În abordarea P&P a teoriei lingvistice MinP, UG asigură un *sistem de principii* fixat, asociat cu un tablou finit de *parametri evaluați* (pe un număr finit de valori). Regulile pentru un LN particular se reduc la alegerea valorilor pentru acești parametri. Noțiunea de construcție gramaticală este eliminată, împreună cu regulile particulare de construcție, specifice gramaticilor generative. Construcții ca

VP, clauză relativă, pasivul etc. devin doar elemente ale unei taxonomii generale, sau colecții de fenomene explicate prin interacțiunea principiilor de UG, legate (setate) cu anumite valori fixate ale parametrilor.

În sistemul computațional al UG există un set de *principii invariante*, fiecare cu un domeniu de *opțiuni* restrânse la elementele funcționale și proprietățile generale ale lexiconului. O *selecție* Z printre aceste opțiuni determină LN concret. În schimb, un limbaj determină o mulțime infinită de SDs lingvistice, fiecare pereche (n, X) fiind obținută din nivelele de interfață (PF, LF), respectiv. *Achiziția de limbaj* implică fixarea mulțimii 2, în timp ce *gramatica* limbajului se reduce la specificarea lui 2. În fine, un *sistem de parsare* care este invariant și neantrenat (cum adesea se presupune) poate fi văzut ca o transformare a perechii (I, TI) într-o schemă structurată similară cu o SD. Condițiile asupra reprezentărilor LN impuse pentru diferite principii și (sub)teorii, cum ar fi teoria *legării*, teoria *cazurilor*, \wedge -teoria etc., sunt satisfăcute pe nivelele de interfață ale sistemelor de performanță. Toate aceste ipoteze fac parte din teoria MinP a lui Chomsky și din construcția sa pentru UG.

2.2 (Sub)teoria X-bar în contextul teoriei MinP

Sistemul computațional al unui LN concret preia reprezentările unei forme date și le modifică, în timp ce UG trebuie să furnizeze mijloacele de a reprezenta o mulțime de elemente din *lexicon* într-o formă care să poată fi accesată și procesată de către sistemul computațional. Forma sub care este accesat lexiconul de către sistemul computațional poate fi considerată ca fiind o anumită *versiune* a teoriei X-bar. Schemele X-bar pot fi asociate în mod natural cu *structuri de trăsături lingvistice* [18], ca un *tip de date lingvistice* standard și invariant pentru a reprezenta și a procesa LN eficient. În strategia SCD, *schemele X-bar augmentate* [19] considerate până acum nu sunt doar tipuri de reprezentare a datelor la nivelul lexiconului ci ele pot asigura structurile invariante fundamentale pentru a reprezenta și a procesa textul în LN la nivel sintactic [1], [2], [6].

În teoria *Minimalist Program* și modelarea P&P a UG, proprietățile și relațiile esențiale sunt formulate în termenii simpli și elementari ai *teoriei X-bar*. Astfel, o *structură X-bar* este compusă din *proiecțiile* lingvistice ale *nucleelor* selectate din lexicon. În *schema X-bar* a teoriei MinP reprezentată în Fig. 2.2.1. sunt prezente două relații locale: relația *Specificator-Nucleu* de la ZP la X, și relația *Nucleu-Complement* de la X și YP (ordinea categoriilor nu este esențială, fiind stabilită de către parametri P&P adecvați de ordonare). Relația Nucleu-Complement (Nucleu-Argument) nu este numai "locală" ci și fundamentală deoarece este asociată (8) *relațiilor tematice*.

Dacă, pentru moment, nu este luată în considerare *relația de adjuncție* sau *adjuncții* se consideră a se afla printre argumentele-complemente, X-bar structurile pot fi reduse la X-bar schema din Fig. 2.2.1, cu următoarele specificări: (a) Sunt considerate numai relațiile locale (deci nici o relație de proiecție între X și vreo sintagmă inclusă în proiecțiile maximale YP sau ZP); (b) Relația *Nucleu-Complement* reprezintă *relația locală de nucleu* [Eng: *core relation*]; (c) O relație locală *admisibilă* a schemei X-bar din MinP este cea *Nucleu-Nucleu*. De exemplu, relația unui verb predicativ cu nucleul predicțional (deverbal) al unei sintagme nominale pe care o subcategorizează; (d) O altă relație în X-bar schema din MinP este *legătura de lanț* [Eng: *chain /ln/c*], corespunzând unui *lanț de dominare* sau de *governare*.

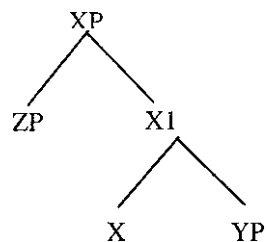


Figura 2.2.1. Schema X-bar din teoria MinP

Guvernarea realizată de nucleu joacă un rol central în toate componentele teoriei MinP asupra UG. Una dintre problemele-cheie este asignarea corectă a trăsăturilor nucleului. În HPSG și SCD, de exemplu, acest lucru este realizat la nivel lexical (BAR = 0), după aplicarea flexionării, cât și la nivel de lexicon (nivel de proiecție notat convențional cu BAR = -1) pentru clasa categoriilor lingvistice cu *proprietăți funcționale* (predicative, relaționale), fie ele verbe, substantive, adjective, markeri de sintagmă, markeri de discurs etc. care antrenează un comportament sintactic funcțional [2], [6]. În particular, pentru teoria MinP, subteorii ca *S-governarea* și *governarea de caz*, corespunzând *6-marcării* și *Caz-marcării*, sunt cele mai importante forme de dominare. Un studiu comparativ al guvernării categoriilor (dependență, dominare), relație prezentă firesc în cele mai importante teorii sintactice formale existente în acest moment, este inclus în [20].

Structurile propuse de teoria X-bar trebuie "animate" de către (sub)teoriile (de asemenea complementare) conținute în MinP și P&P, și care explicitează fenomenele de *governare*, *legare*, *limitare* etc. ce s-au dovedit a fi importante pentru orice teorie lingvistică deoarece ele asigură reguli pentru organizarea *lexiconului* și a *sistemului computațional* care generează și recunoaște SDs.

De exemplu, în funcționarea *teoriei cazurilor* în contextul schemelor X-bar din MinP, ipoteza standard din MinP este aceea că, într-o frază (propoziție), relația *Specificator-Nucleu* atrage după sine *cazul structural* pentru *poziția de subiect*, în timp ce *poziția de obiect* primește cazul sub guvernarea nucleului V, incluzând construcții în care obiectul marcat cazual de către un verb nu este complementul său ci doar un adjunct (așa-numita *marcare de caz excepțională*).

În continuare este prezentată structura X-bar de bază a *clauzei* în teoria MinP, cu următoarele notații uzuale: C = COMP = Complementizator, T = Timpul, Agr_s = acordul subiectului; Agr_o = acordul obiectului etc.

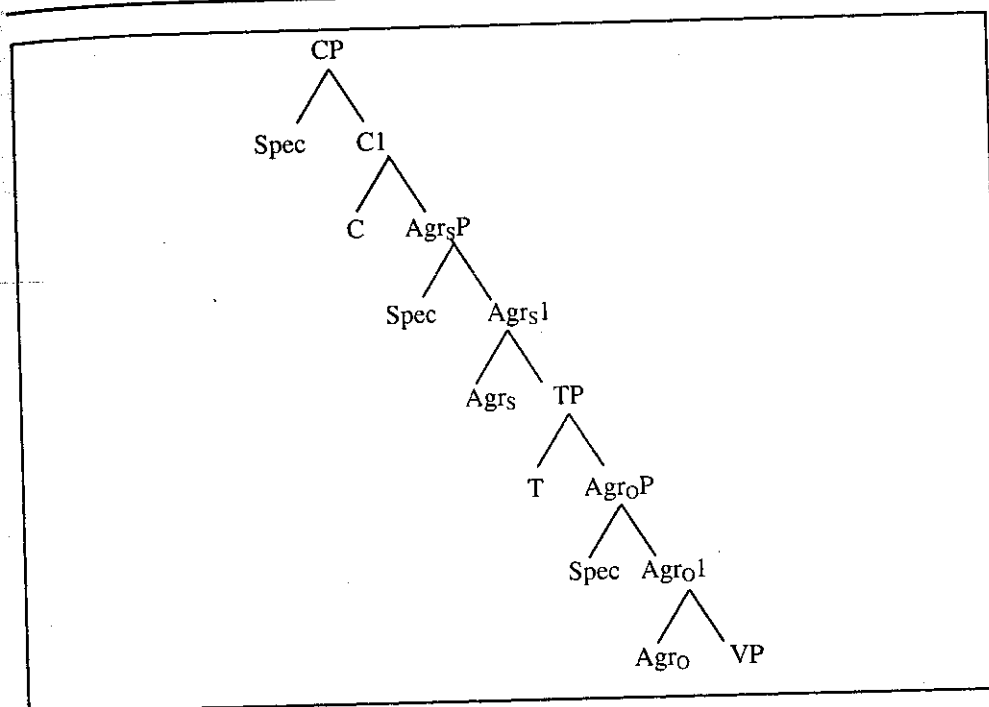


Figura 2.2.2. X-bar structura clauzei în teoria MinP

Schemele X-bar clauzale clasice din Fig. 1.2. și Fig. 1.3. sunt expandate în Fig. 2.2.2., cu următoarea posibilă interpretare funcțională: X-bar schema MinP are ca nucleu VP, care își selectează sintagma-Obiect (sau argument, mai general) prin acord și marcare, afectată apoi de Specificator. Un timp finit T aplicat sintagmei Verb-Obiect generează sintagma TP [Eng: *tensed phrase*], căreia i se aplică apoi aceleași funcții de selecție a subiectului (acord, marcare, specificare), generând sintagma Verb-Obiect-Subiect, care este de fapt clauza finită simplă (notată S). În fine, prin aplicarea asupra lui S (văzută ca sintagmă Agr_sP) a unui complementizator C (sau marker clauzal, marker de discurs etc.) se obține o clauză "completă" ce poate, prin recursie, să genereze orice frază [Eng: *sentence*].

Alte exemple de X-bar scheme bazate pe MinP și P&P, ce pot fi discutate în contextul mai general al fenomenelor de guvernare sunt date de Fig. 2.2.3. care urmează.

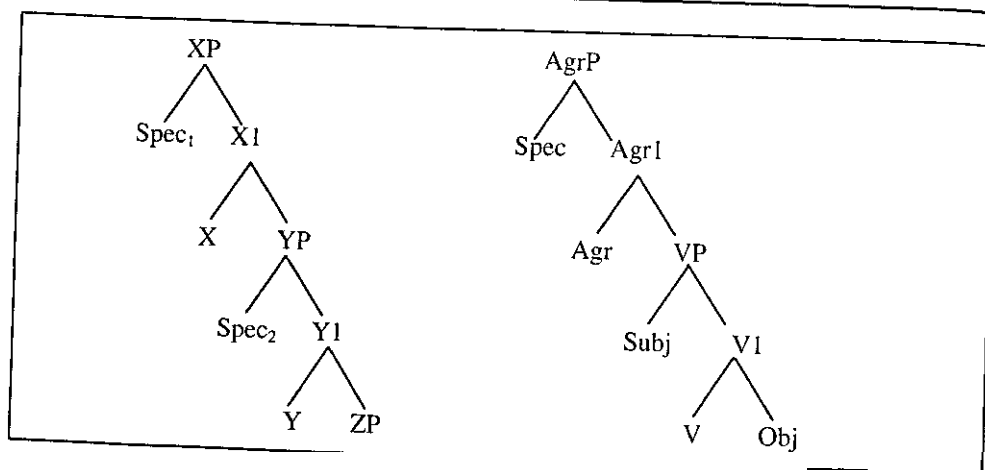


Figura 2.2.3. X-bar scheme în fenomene de "ridicare"
la nivel de Spec în MinP

Concluzia este aceea că teoria X-bar din MinP sintetizează relațiile fundamentale de dependență, descrise de X-bar schemele propuse, și implicate în procesele de organizare a lexiconului și a sistemului computațional din UG. X-bar teoria în abordarea MinP reflectă în principal *aspectele statice* întâlnite în fenomenele de *guvernare* (c-comandă, m-comandă, bariere, categorii de blocare etc), în *teoria legării* și în procesele de *referință-coreferință*, în stabilirea dependențelor la mare distanță (extra-clauzale) etc Nu vrem să intrăm în detalii și să explicităm mecanismele de lucru ale X-bar schemelor considerate, ci mai curând să atragem atenția asupra *teoriei X-bar* ca o *componentă fundamentală* a unei teorii lingvistice noi și elaborate cum este MinP și modelul său P&P [3].

Teoriile MinP și P&P nu reprezintă un punct-terminus pentru evoluția teoriei X-bar. Dimpotrivă, asigură o bază de pornire pentru o strategie radical diferită în care Chomsky examinează cele mai serioase argumente pentru a abandona (!) teoria X-bar [3; Cap. *Categorii și transformări*]. Această alternativă și consecințele sale sunt discutate în secțiunea 5, și ar trebui să reprezinte una dintre cele mai importante provocări prezente pentru domeniul analizei și proiectării teoriilor lingvistice [21].

Unul dintre principalele scopuri ale secțiunii care urmează este de a introduce propunerea noastră de *scheme X-bar funcționale* (*scheme FX-bar*) în cadrul strategiei lingvistice SCD. Propunerea noastră o considerăm a fi o poziție pragmatică și echilibrată în *direcția* teoriei X-bar, atrăgând atenția asupra adevăratului său rol și oportunităților computaționale din lingvistica reală, înțelegerea corectă a aspectelor *statice* și *dinamice* ale acestei versiuni a teoriei X-bar ar trebui să fie de asemenea o consecință a unei priviri cuprinzătoare a întregului context al teoriilor lingvistice care stabilesc principiile de dependență,

pasele de markeri, categoriile și ierarhiile, regulile de referire și structurare, în Strânsă relație cu formele și regulile de construcție ale (sub) schemelor FX-bar.

3. Scheme X-bar funcționale și strategia lingvistică SCD

în [19], în contextul *strategiei lingvistice SCD* (*Segmentare-Coeziune-Dependență*) [22], [19], [1], [2], [6], este definită o clasă de *scheme X-bar augmentate* (*scheme AX-bar*), scheme destinate a reprezenta *invarianti sintactici generali* de reprezentare și operare cu structurile gramaticale ale LN, în particular pentru limba română, ca soluție la problemele de analiză și generare automată a LN. *Schemele FX-bar* (*funcționale*) propuse aici completează și extind *schemele AX-bar* [19], și pot fi interpretate în mai multe moduri: (1) din punct de vedere *static*, schemele FX-bar pot furniza câteva de tipuri fundamentale de date pentru reprezentarea informației lingvistice în structuri de trăsături lingvistice, standardizate și tipizate; (2) din punct de vedere *dinamic*, schemele FX-bar pot codifica informația lingvistică în formă procedurală ca *funcții și relații standard* ce sunt (recursiv) apelate în cadrul proceselor de analiză și generare a LN; (3) schema FX-bar generală poate fi de asemenea interpretată și utilizată ca un automat pe baza căruia să se realizeze o analiză *on-line* a textului unei fraze, cuvânt cu cuvânt.

3.1. Câteva preliminarii asupra SCD

Sunt necesare unele precizări asupra noțiunilor și notațiilor cu care lucrează *strategia lingvistică* SCD. Unul dintre elementele importante este că nivelul 2 (BAR = 1) în X-bar schema clasică joacă un *rol-cheie* în SCD pentru construcția structurilor sintactice, și este utilizat sub numele de *grup nominal* (NG), *grup verbal* (VG), *grup adjectival-adverbial* (AG), în general XG, pentru X = N, V, A. *Grupul XG* corespunde *proiecției lexicale* X1, cu X = N, V, A, și *clauzei minimale* CLO, în X-bar schema fundamentală propusă în Fig. 3.2.1.

Să menționăm că orice XG (X1) este un XP (X2), dar nu și invers, deoarece proiecția categoriei X în cazul XG lucrează numai pentru nivelul BAR ^ 1. SCD face de asemenea distincție între câteva tipuri de NGs (NGs elementare, Predicaționale, non-predicaționale, etc), VGs (VGs la un timp finit și la un timp non-finit) etc

O altă trăsătură esențială și specifică a SCD este un tratament adecvat al *Proprietăților funcționale* ale categoriilor lingvistice, ca și al tuturor *categoriilor Naționale* și sintagmelor (expresiilor) de *discurs*. Mecanismul utilizat pentru a obține acest lucru se bazează pe *clase de markeri lingvistici și ierarhiile* lor [1, 12] și [6], [7], [45]. Câteva observații se impun:

(a) *Marcherii* din SCD, numiți *marcheri de structuri sintagmatice* (PS-Ms) [Eng: *phrase-structure markers*], sunt cu totul diferiți de ceea ce teoria lui Chomsky numește formal "*marcheri de sintagmă*" [Eng: *phrase-markers*] în [17], sau *T(ree)-marcheri* în [3]. *Marcherii Chomsky* sunt definiți ca "*tăieturi orizontale*" (sau "factorizări") în cadrul unui arbore de derivare, sau ca fiind arborele însuși. Mult mai apropiați de ceea ce sunt PS-Ms în HPSG [16], *marcherii de structuri sintagmatice* (PS-Ms) din SCD sunt acele categorii lexicale și nelexicale care se aplică cuvintelor și structurilor sintagmatice (PSs) cu scopul de *evidențiere*, de a *marca*, anumite funcții și relații sintactice și semantice pe care PSs respective le joacă în cadrul unei exprimări. Punerea în evidență a anumitor funcții care se aplică PSs se referă la (cel puțin) câteva elemente: *tipul funcției* (sintactic, semantic, relațional, logic, pragmatic, discursiv etc), *locul*, în text, *unde începe aplicarea funcției sau relației*, și *domeniul (domeniile, conexe sau nu)* de aplicare a funcției sau relației (limitele textuale între care se aplică).

Exemple tipice de PS-Ms din SCD sunt: (a) trăsăturile *predicative* generate de către *categoriile predicative* (de fapt, verbe, substantive, adjective și adverbe predicative); (b) acele *mijloace gramaticale* prin care sunt introduse *noi* NGs (grupuri nominale în limbajul SCD), VGs, AGs (Caz-marcarea, acordul, gradele de comparație, etc); (c) acele *categorii și expresii* (numite și *marcheri de discurs*) care introduc *noi* clauze; (d) PS-Ms care introduc *proprietăți relaționale* asupra PSs și clauzale (de exemplu de *marcheri* de tip logic cum sunt structurile *dacă-atunci-altfel*, *deoarece*, etc, dar și *marcheri* de tip sintactico-semantic cum sunt aceia care introduc *categorii și clauze subordonate* etc)

(b) SCD se aseamănă din unele puncte de vedere cu abordarea [16] a HPSG și, parțial, cu [15], care exploatează, pentru prima oară în clasa teoriilor lingvistice bazate pe *gramatici* de PSs (PS-Gs), într-o mult mai mare măsură, categoria lingvistică a *marcherilor* PS-Ms. În [16], Pollard & Sag "postulează o nouă parte a *marcherilor* de discurs,... ce se remarcă ... printr-un nou atribut al categoriilor (în plus față de NUCLEU și SUBCAT) numită MARKING, cu valori din sortul *marking*". Teoria HPSG enunță PRINCIPIUL MARCĂRII [16, p. 400] după cum urmează:

"într-o sintagmă cu nucleu, valoarea trăsăturii MARKING este lexical-identică cu cea a trăsăturii MARKER-DAUGHTER dacă aceasta există, și cu cea a trăsăturii HEAD-DAUGHTER în caz contrar.

Modul în care HPSG [16] pune la lucru PS-Ms reprezintă un bun și esențial pas înainte, deși credem că nu exploatează îndeajuns potențialul funcțional și relațional al diferitelor clase de *marcheri* și ierarhiile acestora (așa cum face strategia SCD, vezi și [7], [45]).

(c) Continuând și extinzând construcția limbajului, ca o expresie de convergență între gramatica categorială și *Minimalist Program*, Chomsky [3] consideră *transformările generalizate* (GTs) și concepe un demers de înlocuire a

X-bar teoriei, ce explică în Programul Minimalist structura constiuenților (sintagmatici) complecși, prin GT *Merge* care construiește obiecte sintactice pornind de la obiecte sintactice simple (de exemplu, "*speaks*" și "*French*" sunt "reunite" într-un nou obiect sintactic "*speaks French*" etc). Mai multe formalizări ale acestui nou curent al ideilor lui Chomsky pot fi găsite în cadrul gramaticilor logice multi-modale și de tipuri categoriale, e.g. [21], [23], [24] etc (vezi și secțiunea 5).

(d) Dintr-o perspectivă diferită dar oarecum similară, *gramatica funcțională* (FG) [25] a lui Simon Dik, orientată funcțional și semantic, încearcă să facă aceleași lucruri. Ca și în SCD, FG găsește *patru tipuri ierarhice* de bază ale categoriilor relaționale, aceste tipuri corespunzând într-o bună măsură cu *clasele de marcheri* PS-Ms și *ierarhiile* lor stabilite în SCD [7], [2], [6]. PS-Ms reprezintă acele mijloace lingvistice de "suprafață" pe care le utilizează un limbaj natural pentru a organiza sintactic și semantic structurile codificate în construcții gramaticale. Se impune în viitor o analiză comparativă între cele *patru nivele* sau "*straturi*" din organizarea formală și semantică furnizată de FG [25], și cele *patru nivele* de proiecție lingvistică, împreună cu clasele de *marcheri* corespunzătoare, din SCD: (1) *cuvântul* (lexical); (2) *sintagma* XG (X = N, V, A) subclauzală; (3) *clauza* (finită și infinită); (4) *unitatea de discurs* (una sau mai multe fraze, care să formeze un *segment* de discurs).

(e) În fine, privitor la utilizarea intensivă a *caracterului predicativ* pe care categoriile lexicale majore (N, V, A) îl poartă (proprietate moștenită sau dobândită apoi de alte categorii gramaticale), *strategia lingvistică* SCD este comparabilă în special cu FG, cu accentul particular pe ierarhiile de delimitare și marcarea aplicate structurilor sintactico-semantic. SCD pornește de la *lexicon* și stabilește la acest nivel o *taxonomie predicativă inițială* pentru categoriile lexicale majore. Un exemplu simplu al acestei taxonomii predicative este dat de către cele două categorii importante de substantive comune: *substantive existențiale* sau *obiectuale*, a căror *predicativitate* (trăsătură PRED) este EXIST (e.g. [Eng: *student*, *table*; Rom: *elev-student*, *masă*]) și a căror reprezentare funcțională reflectă categorii individuale sau personale, de exemplu predicatul uni-variabil *student(X)*, *masă(X)* etc, și substantive de tip-predicativ, a căror predicativitate (trăsătură PRED) are valoarea ACT, e.g. [Rom: *întâlnire*, *invidie*, *marcare* etc], și ale căror reprezentări funcționale depind de mai multe variabile, de exemplu *întâlnire(X, Y, ...)*, *invidie(X, Y, ...)*, *marcare(X, Y)* etc. Substantivele proprii și/sau personificările sunt codificate prin constante ale variabilelor din predicatele de mai sus. Câteva din remarcile anterioare vor fi aprofundate în concluziile finale ale lucrării.

Schemele FX-bar, ca și precursorile lor *schemele AX-bar* [19], reflectă Pentru SCD faptul că un XPG (grupul sintagmatic de nucleu X), sau mai simplu XG, conține un *nucleu*, reprezentat printr-o categorie lexicală (nevidă) sau printr-o categorie virtuală (vidă), înconjurat (prin relații de *coeziune*) de specificatori și/sau Codificatori de tipul A (adjectival-adverbial). Este esențial să facem următoarea specificare: un XG din SCD nu include nici un complement (argument obligatoriu)

sau adjunct. Complementele și adjuncții, împreună cu *nucleele* de nivel BAR = 1 formează nivelul BAR = 2 în FX-bar schema propusă în Fig. 3.2.1. Pentru un anumit nivel de specificare semantică, FX-bar schemele nu fac o distincție clară între complementele (argumente obligatorii) și adjuncții, considerând toate structurile subcategorizate ca fiind argumente *sintactice*] clasificări ulterioare (suplimentare) sunt făcute pe baza șabloanelor verbale și *restricțiilor* sintactice, semantice, și pragmatice asupra componentelor șablonului, la nivel de lexicon.

O problemă a cărei soluție poate influența în mod special și teoria X-bar este aceea a asignării corecte a complementelor și adjuncțiilor, în particular, a stabilirii corecte a dependențelor dintre grupurile nominale (NGs). Soluția acestei probleme nu se poate obține la nivel sintactic, iar o soluție completă nu se poate obține uneori nici chiar în contextul unui nivel semantic minimal (vezi [26], [27]). Chomsky remarcă realitatea că "... *the distinction between modifiers and arguments is notoriously difficult in certain cases*" [9, p. 44]. Exemple simple ilustrează această problemă: în TBarr [8], sintagma "*the students of physics*" este văzută ca un NP cu un *argument* PP, în timp ce sintagma "*the students in the yard*" este considerată a fi un NP cu un *adjunct modifier* PP. De fapt, în numeroase LNs, inclusiv engleză, se pot aduce multiple argumente serioase pentru ca cele două sintagme să poată fi la fel de bine interpretate fie într-un fel, fie în celălalt.

Soluția SCD pentru acest exemplu foarte particular este următoarea (schițând și soluția problemei generale): substantivul "*students*" este obiectual, adică *nu are* o natură *predicațională* prin el însuși, astfel că ambele sintagme nominale care îl succed sunt considerate de către SCD ca fiind *modificatori* pentru NG "*students*". Natura acestor modificatori poate fi diferită deoarece "*physics*" este introdus de markerul de caz (genitiv) "of, în timp ce "*the yard*" este introdus de markerul *prepozițional* "in". În general, când nucleul lui NG posedă o trăsătură *predicațională*, atunci NG care urmează nucleului predicațional asigură o distribuție sintactică ce satisface un anumit șablon (verbal) al predicatului (verbului) corespunzător.

Clasele din PS-Ms și ierarhiile lor din SCD [7], [45] sunt responsabile pentru delimitarea structurilor sintagmatice propuse de schemele FX-bar, și pentru stabilirea dependențelor sintactico-semantice. Diferitele tipuri de markeri sunt adesea aplicate simultan (deci multiplu) asupra acelorași categorii gramaticale, în cadrul anumitor nivele de structurare (proiecții pe BAR-nivel). Similar cu unele teorii lingvistice (LFG, FG, și parțial HPSG) dar contrar altora (GB, GPSCB etc), SCD nu consideră *prepoziția* (X = P) ca fiind o categorie lexicală majoră. În SCD, P primește rolul unui marker (funcțional), având atât proprietăți de marker de caz cât și de complementizator. Categoriile HPSG PP[+PRD] sau PP[-PRD] (vezi [16]) sunt irelevante pentru SCD deoarece trăsătura +PRD în HPSG este atribuită numai lui PP subcategorizat de un V, în timp ce trăsătura (*predicațională*) PRED din SCD poate fi în mod egal atribuită lui V, N, sau A (la nivelul lexiconului, cel puțin) dar nu și lui P.

În SCD proprietățile de subcategorizare sunt exploatate *ab initio*, la nivelul de organizare a lexiconului, pe baza *trăsăturii funcționale* PRED de *predicaționalitate*, asignată sau nu, unora din categoriile sintactice majore N, V, A. Observații lingvistice

Empirice ne-au convins, încă de la începuturile cristalizării SCD [22], funcțională și predicativă adecvată ar trebui să reprezinte punctul oricărei teorii lingvistice, atât din motive teoretice cât și computaționale din abordările actuale (cum ar fi [27]-[32]) aduc o susținere puternică pentru ideile esențiale din SCD, în special folosirea intensivă a proprietății funcționalității descrierilor lexical-semantice ale categoriilor lingvistice și procesarea automată a LN cât și în organizarea bazelor de cunoștințe

[19] propune următoarea specificare a *Principiului Proiecției Maximală* (PMP) [Eng: *Principle of Maximal Projection*], ca un pas important în realizarea unei teorii lingvistice intensivă a trăsăturilor predicaționale (funcționale) ale categoriilor lingvistice în SCD. Propunem aici

O specificare a P M P (formă actualizată):

Proprietățile de subcategorizare ale categoriilor sintactice majore N, V, A depind de trăsătura lor lexical-semantică PRED(icativity), cu valori

EXIST, și de trăsătura lor morfo-semantică TENS(e), cu valorile ACT și INFI(nite).

Trăsătura PRED, atribuită categoriilor majore N, V, A la nivel de lexicon primește două valori: valoarea ACT, pentru acele categorii care sunt *predicaționale* (în literatură este folosit adesea termenul "*deverbalitate*"); EXIST, pentru acele categorii N, V, A cu caracter *existențial*, care sunt *predicaționale*. Trăsătura TENS primește valorile FINI(te) pentru atributele *existențiale* ale categoriei V care posedă un timp sau aspect finit, personal, și valori pentru toate celelalte categorii și situații. Exemple:

[Eng: *boy, pencil* \ Rom: *băiat, pix*] PRED:= EXIST; și TENS:= ACT

[Eng: *attempt, showing, proved*; Rom: *încercare, arătând, demonstrat*] PRED:= ACT; și TENS:= INFI

[Eng: *are*; Rom: *sunt*] PRED:= EXIST; și TENS:= ACT

[Eng: *gives*; Rom: *dă*] PRED:= ACT; și TENS:= INFI

Într-un grup verbal VG reprezentând un compus la un timp verbal "pozitiv" de trăsături, cum sunt ACT sau FINI sunt *moștenite* de la trăsăturile VG, sau pot fi obținute *cumulativ* prin trăsăturile sintactice.

Specificarea PMP de mai sus este o funcție *proiecției maxime* a trăsăturii PRED din SCD deoarece în multe LNs, inclusiv în română, *calitatea devine funcțională*, deci funcțională) a categoriilor lexicale tradiționale *non-verbale*. A trebuie descoperită cât mai devreme posibil și asignată la nivelul lexiconului. Exemplu, în engleză, deși pentru substantivele care *verbalizează* valoarea trăsăturii lor TENS este INFI, aceste substantive posedă, în SCD, aceeași valoare ACT sau EXIST pe care o au verbele

substantivele (sau gerunziile) în "-/ng", și astfel posedă *aceleași* proprietăți de subcategorizare ca ale verbului de origine.

3.2. Ipoteze de lucru și aspecte caracteristice ale FX-bar schemei

Continuând ideile de bază ale schemelor AX-bar din [19], propunem, pentru SCD, FX-bar schema generală din Fig. 3.2.1. Muchiile din stânga conțin *noduri* cu rol *funcțional* sau *relațional*: marcheri, cuantificatori, specificatori, modificatori (eventual adjuncți). Pentru a obține reprezentări sintactice și semantice corecte, nodurile funcționale se aplică (recursiv) *nucleelor* X_k și CL_k , $k = 0, 1, 2$, iar nucleeele, cu rol funcțional (predicațional, X_1) sau relațional (eventual X_2), au ca argumente clauze infinite (complemente, X_1) sau finite (X_2). Precizăm că la acest nivel nu se poate face distincția dintre complemente COMPLi (argumente obligatorii) și adjuncți ADJCTj (argumente opționale). În mod normal, în Fig. 3.2.1., ADJCTj sunt "amestecați" printre ARGj, la nivel sintactic nefiind discernabili de complementele obligatorii ale unui nucleu predicațional. Poziția funcțională (la stânga nodului X_1) a nodurilor ADJCT poate rezulta doar în urma unor calcule semantice și pragmatice suplimentare, din care se obține rolul tematic al argumentelor ARGs ale lui X_1 .

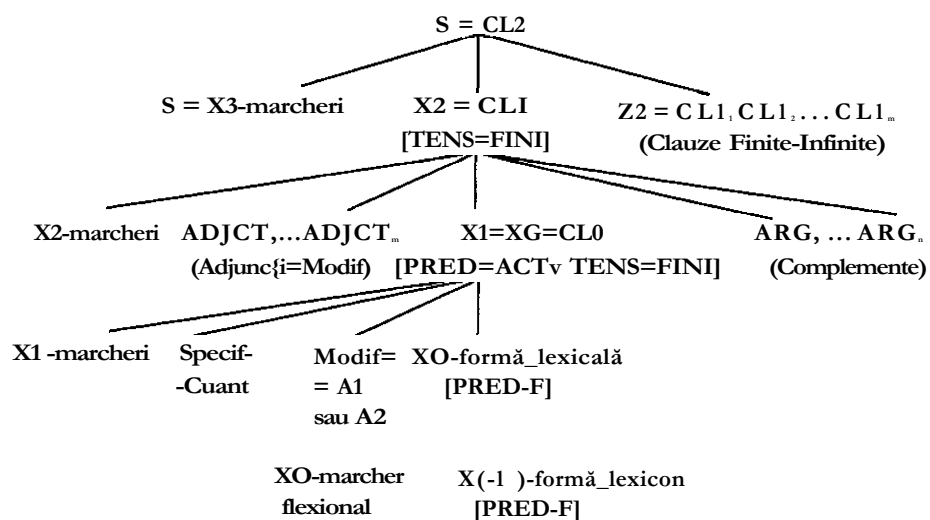


Figura 3.2.1. Schema (funcțională) FX-bar generală

(*) **Aspecte specifice ale schemei FX-bar propuse:** (*1) Sunt permise un număr arbitrar de *argumente* (sau *sateliți* în sensul [10], [31]), toate notate cu ARGs. În SCD, ARGs sunt formate din *complemente* obligatorii (COMPLs) și din *adjuncți*

(ADJCTs), sau *complemente opționale*. ADJCTs pot fi reprezentați la nivel sintactic tot ca argumente ale nucleului, însă la nivel semantic ADJCTs au rol de *modificatori* ai nucleului. Notația "A-poziție" din teoriile Chomskyene, care înseamnă ARG-poziție, nu trebuie confundată cu notația noastră pentru categoria A = adjectiv-adverb. În teoriile și notația lui Chomsky, COMPLs sunt în A-poziție (ARG-poziție), în timp ce ADJCTs nu. SCD se situează pe o poziție sintactică similară cu HPSG [16], care utilizează lista SUBCAT pentru a codifica toate sintagmele pe care le subcategorizează un nucleu semantic, adică atât COMPLs cât și ADJCTs (sau ARGs din SCD). (*2) Sintagmele AG = A0 sau A1, sau AP = A2 sunt *postulate* de către SCD ca fiind de tipul categoriei *funcționale* Modif, manifestate prin categoriile A (de nivel XO, și aplicabile la nivel XO), ADJCTs (de nivel X1, și aplicabile la nivel XO și X1), și clauza relativă (de nivel X2, și aplicabilă la nivel XO și X1). (*3) Categoria generică Specif (sau Spec), în care intră cuvintele și sintagmele ce desemnează *cuantificatori* de toate tipurile (generalizați), determinatori (în particular), este postulată de către SCD ca fiind o *categorie funcțională* ce poartă trăsături de natură *cuantificațională* la nivel lexical (în particular, *negația* la nivelul X1), inclusiv articularea (hotărâtă sau nu), suprapunându-se deci uneori peste X1-marcheri de trăsături funcționale cum este *acordul*. Relațiile (funcționale) de *acord* sunt esențiale pentru *coeziunea locală* și *globală* în cadrul strategiei SCD: acordul dintre XO-Modif și XO-Specif cu nucleul XO (la nivel X1), acordul Nucleu-Subj (sau chiar Nucleu-COMPL) și acordul COMPL-Pron_{em}fat (Pronume emfatic) (la nivel X2), o anumită *corespondență* a timpurilor evenimentelor într-o clauză și între clauze. Aceste tipuri de relații de acord, referință și coreferință, coeziune, coerență, etc. sunt responsabile pentru o largă clasă de dependențe locale și globale, inclusiv dependențe la distanță mare și în extra-poziție. Accentul în componenta de *coeziune* a strategiei SCD (Segmentare-Coeziune-Dependență) cade pe mijloacele *sintactice* și de "*suprafață*", mai curând decât pe cele semantice, încercând să găsim, să extragem, și să utilizăm într-o măsură maximală informații de ordin superior, cum ar fi informația de discurs [34], pragmatică, semantică etc. (*4) Sintagma tradițională PP din teoriile lingvistice clasice, iar în SCD, grupul prepozițional PG (format dintr-un grup nominal NG care este precedat de o prepoziție sau o locuțiune prepozițională) este întotdeauna considerată un ARG (COMPL sau ADJCT) în FX-bar schemele al căror nucleu (lexical nevid sau vid) este N, V, A. Această ipoteză de bază asupra PG este justificată de SCD prin faptul că P *nu* este considerată o categorie majoră, adică o categorie de nivel X1 în schema FX-bar din Fig. 3.2.1. ci doar o categorie de nivel XO. Proprietățile de subcategorizare ale N, V, A (dar nu și P) pot fi asignate *ab initio*, 'a nivel de lexicon, începând cu trăsătura lexicală PRED a categoriilor predicaționale. Categoria P poate primi proprietăți funcționale, cel mai adesea ca *marcher de caz*, uneori proprietăți *relaționale* (de exemplu [Eng: *on*; Rom: *asupra*]), dar nu și proprietăți de subcategorizare. (*5) Subiectul (Subj) în SCD, lexical nevid sau vid (^oRO), este considerat ca un argument special al proiecțiilor maxime ale categoriilor X = N, V, A într-o clauză finită (de nivel X2) sau infinită (de nivel X1). (*6) În ipotezele (*5) și (*2) de mai sus, categoria lingvistică tradițională VP este dizolvată într-un grup verbal VG (finit sau infinit), înconjurat (de cele mai multe ori

urmat) ca nucleu de ARGs și formând o clauză finită, respectiv infinită. (*7) *Teoria limitării* și multe probleme majore legate de TBarr [8], [9], [17] sunt explicitate și rezolvate în cadrul realizat de SCD și schemele FX-bar, în principal datorită delimitării clare a funcțiilor și relațiilor care se aplică cuvintelor și sintagmelor, a reprezentării lor lexicale prin clasele de PS-Ms, și a specificării domeniului lor de aplicare. Acest rol este realizat explicit în cadrul claselor și ierarhiilor de markeri propuse și utilizate de SCD [2], [6], [7], [45]. Trebuie să remarcăm că în lucrările sale cele mai recente [34], [35], Chomsky adoptă o tehnică similară de "limitare" a operațiilor de *construire* [Eng: merge] și *transformare* [Eng: move] doar la "domeniul" sintactic al unei "faze" [Eng: phase], o unitate textuală (care în general coincide cu clauza!) în care Chomsky propune următorul *principiu de impenetrabilitate* "într-o fază (clauză n.n.) F cu nucleul H, domeniul lui H nu este accesibil la operații în exteriorul lui F, ci este accesibil numai H și muchia sa (nodul său ascendent)" [34]. Exact așa este construită și funcționează schema FX-bar! De asemenea, fenomene de *teoria legării* [9], [8], [3], [16], *legăturile* [Eng: linking] din [27], mecanisme de *coeziune* (locală și globală) și *discurs* întâlnite în [36], [31], [33], etc. sunt mai ușor de pus în evidență și de rezolvat în cadrul oferit de strategia lingvistică SCD și teoria FX-bar.

(*) Observații asupra ipotezelor de lucru pentru schema FX-bar din Fig.

3.2.1.: (41) Schema FX-bar este proiectată să lucreze în asociere cu un parser care este capabil să recunoască clasele de PS-Ms și structurile sintagmatice considerate de strategia lingvistică SCD. Schema FX-bar este organizată pe *patru nivele* de proiecție BAR = (H-3 (deasupra nivelului de lexicon, notat convențional BAR = -1); *trei nivele X0-X1-X2* corespund proiecției dintre nivelul lexical (BAR = 0) și nivelul *clauzal*, al structurilor *uni-eveniment*, alte *trei nivele CL0(=X1)-CL1(=X2)-CL2* corespund proiecției dintre nivelul *clauzal minimal* CLO = X1 și nivelul *frazei*, al structurilor *multi-eveniment*. Nivelele uni-eveniment X0-X1-X2 exprimă predicția clauzei (propoziției) simple în care sunt distribuite categoriile lexicale de bază și sintagmele pe care le generează, în timp ce nivelele CL0-CL1-CL2 exprimă relațiile logice și predicționale (de ordinul doi) dintre clauzele simple. Schema FX-bar lucrează într-o manieră recursivă (top-down sau bottom-up), atât în situațiile de analiză cât și în cele de generare în care este antrenat parserul asociat, în strânsă cooperare cu strategia lingvistică SCD, cu clasele de PS-Ms și ierarhiile lor și, mai ales, pe baza *meta-algoritmilor* SCD de analiză-generare [1], [2], [6], [7]. Să mai observăm că FX-bar schema din Fig. 3.2.1. poate fi utilizată independent de așa numita *ordine canonică* (sau *sistemică*) a cuvintelor și sintagmelor dintr-o clauză, specifică fiecărui LN [37], [38]. (^2) Valoarea ACT de trăsătură (funcțională) pentru categoriile N și A (și implicit V) este atribuită acestor categorii la nivel de lexicon atunci când ele corespund unor evenimente cu actanți și/sau stări multiple. Valoarea EXIST este implicit sau explicit introdusă de formele și înțelesurile verbelor existențiale (a fi), modale (a trebui), etc. (^3) Trăsătura (funcțională) TENS este similară cu categoriile virtuale I (INFL) și T (Tense) din teoriile GB și TBarr ale lui Chomsky și din schemele S-bar corespunzătoare (Fig. 1.3. și Fig. 2.2.3.). Pentru un VG finit (TENS = FINI), structura V2 corespunzătoare devine clauza finită clasică.

Dacă sintagma XG (X1) este un grup a cărei categorie-nucleu X posedă valorile de trăsături PRED = ACT și TENS = INFI, atunci XG devine noul nucleu al unei clauze infinite ce face parte dintr-o structură de nivel X2 (XP). (^4) Poziția specială a *subiectului sintactic* (Subj) este considerată de către SCD atât o ARG-poziție (asemănătoare, de fapt, cu o COMPL-poziție) cât și o Caz-poziție. În concordanță cu TBarr [8] și cu HPSG [16], Subj primește poziția specială a *primului element* din lista SUBCAT [16]. Aceasta este în esență o poziție sintactică, iar Subj poate primi o funcție tematică (Opoziție) autentică doar ca rezultat al unor calcule sintactice și semantice suplimentare! (^5) Așa cum rezultă din schema FX-bar din Fig. 3.2.1., sintagmele AP și PP din teoriile lingvistice clasice sunt segmentate de către markerii SCD [7], [45] în sintagme mai mici XG, X = N, V, A. Așa cum am precizat deja, SCD atribuie noilor sintagme următoarele roluri: AG = Modif, cu rol funcțional la nivelul de proiecție X1, și PG = ARG (COMPL sau ADJCT), ADJCT purtând de asemenea rol de Modif al nucleului de nivel X2. PG devine deci un NG P-marcant, iar orice categorie A are de la început reprezentarea (nesaturată) A(X), unde X = N, V, A este nucleul ([existent, viitor, sau lipsind pur și simplu] al sintagmei de nivel X1 în care Modif = A. În mod similar, orice categorie Specif (determinator, cuantificator, etc.) joacă un rol similar, schema FX-bar impunând reprezentarea funcțională Specif(X), unde X este nucleul sintagmei. (^6) În ciuda anumitor asemănări (inerente) între schemele FX-bar și versiunea MinP a teoriei X-bar, există diferențe de bază în ce privește organizarea și funcționarea constructivă dintre schemele (F)X-bar din Fig. 3.2.1. și Fig. 2.2.1. De exemplu, în schema FX-bar, fiecare element lexical se proiectează într-o categorie obiectuală sau funcțională (relațională), aceasta este (coeziv și ^recursiv) înconjurată de către Specif și/sau Modif, iar dacă valoarea ACT a trăsăturii HPRED a nucleului este prezentă, atunci această valoare ACT este moștenită de către întreaga sintagmă al cărei nucleu a fost specificat sau modificat. Această sintagmă cu nucleu predicțional își subcategorizează complementele (argumentele obligatorii COMPLs) și adjuncții ADJCTs (care modifică sintagma-nucleu). În schema X-bar din Fig. 2.2.1., se întâmplă tocmai invers deoarece "*The Head-Complement relation is the "most local" relation of an XP to a terminal Head Y, all other relations within YP being Head-Specifier (apari from adjunction, ...)*" [3: p. 53]. (^7) Deși schema FX-bar generală a fost proiectată având în vedere în primul rând limba română, ea poate fi aplicată pentru a reprezenta, grafic și logic, structuri sintactico-semantice ale LNs cu valori ale parametrilor gramaticali foarte diferite, cum ar fi engleză-germană sau franceză-germană. Distribuția complementelor (argumentelor) în română (engleză, franceză) poate fi foarte diferită de cea din germană; de exemplu, într-o clauză al cărei verb principal din compusul său verbal VG se află în poziție finală, sau pentru o categorie A (adjectiv-adverb) având valoarea de trăsătură PRED = ACT.

Ex. 3.2.2.R. /Paharul /spart //de Ion/ cu mingea /de fotbal/

Ex. 3.2.2.E. /The glass /broken //by Ion/ with / the football /

Ex. 3.2.2.G. /Das/von Ion/mit/dem Fußball //zerbrochene //Glass/

După cum am remarcat în (t1), schema FX-bar poate fi utilizată independent de regulile structurilor sintagmatice și ordinea lor (din română sau

germană), aceasta deoarece principiile rămân aceleași și diferă numai anumiți parametri și valorile lor pentru LNs distincte: în română (și engleză, franceză) argumentele succed o categorie A ce reprezintă un nucleu predicțional, în timp ce în germană ele îl pot precede. Dacă un nucleu V al unei clauze are valorile de trăsături PRED = ACT și TENS = FINI, atunci distribuția ARGs este similară cu cea din română, cu posibile (și probabile) diferențe impuse de *ordinea sistemică*, strict dependentă de LN, a ARGs (a se vedea [37] dar și [27]).

Dacă se încearcă utilizarea formei FX-bar ca "schelet" pentru un automat (sau gramatică formală) de analiză și generare a LN, un asemenea automat ar trebui să mimeze atât forma generală a schemei FX-bar cât și regulile gramaticale de analiză-generare. Partea din automat care reflectă cele patru nivele de organizare a structurilor LN în schema FX-bar ar trebui să fie independentă de limbaj (cel puțin pentru o largă clasă de limbaje europene), în timp ce (sub)partea constituantă care recunoaște structurile lingvistice pe fiecare nivel individual X_k (k = 1, 2, 3) trebuie să fie dependentă de limbaj (acest fapt este binecunoscut și parametrizat). Reprezentarea schemei FX-bar pentru Ex.3.2.2.G. este aceeași cu reprezentările FX-bar pentru Ex.3.2.2.R.-E., și similară cu figura pentru Ex.4.1.2.R.-E.

4. Exemple de aplicare a schemelor FX-bar

Vom expune câteva exemple de aplicare a schemelor FX-bar la reprezentarea sintagmelor, clauzelor și frazelor. În exemplele prezentate, categoriile gramaticale pentru care PRED = ACT sau TENS = FINI vor fi subliniate, iar PS-Ms care se aplică sintagmelor X_k (k = 0, 1, 2) sunt reprezentați grafic în text prin apariția unuia sau mai multor semne 'slash' /. Să notăm că schemele (augmentate) AX-bar din [19], deși oarecum asemănătoare în spirit sunt efectiv scufundate în schema FX-bar generală, diferențele substanțiale constând în forma unitară a FX-bar schemei și în criteriile sintactice și logico-semantice mai clare, pe baza cărora clasele de PS-Ms și ierarhiile lor sunt explicit propuse și aplicate în funcționarea schemei FX-bar.

Care este relația dintre exemplele de FX-bar scheme și formulele logice atașate după reprezentarea grafică? *Prima formulă* este o reprezentare uzuală a LN, care folosește limbajul logicii predicatelor, reprezentare mai apropiată de exprimarea în LN, conținând toate variabilele ce codifică referințele-coreferințele, dar (pentru simplitate) fără cuantificatorii corespunzători. *A doua formulă* este traducerea mai completă a primei formule în limbajul de programare logică Prolog, folosind tehnici clasice de reprezentare a cunoștințelor de LN în Prolog. Pe o scală ascendentă a măsurii în care sintagmele LN ar fi analizate, *schema* FX-bar poate fi văzută ca un prim rezultat al procesului de parsare (analiză), *prima formulă* ar urma procesului de parsare, incorporând fenomenele de referință (și coreferință, rezoluție a anaferei, etc), iar *a doua formulă* ar reprezenta o rafinare a primei formule. Formulele de tipul doi reprezintă de asemenea atât un stadiu final al procesului de analiză a frazei cât și punctul de pornire în procesul de generare a

...i fraze (conform cu abordarea [39], [6] a generării automate a LN, de [31], de exemplu).

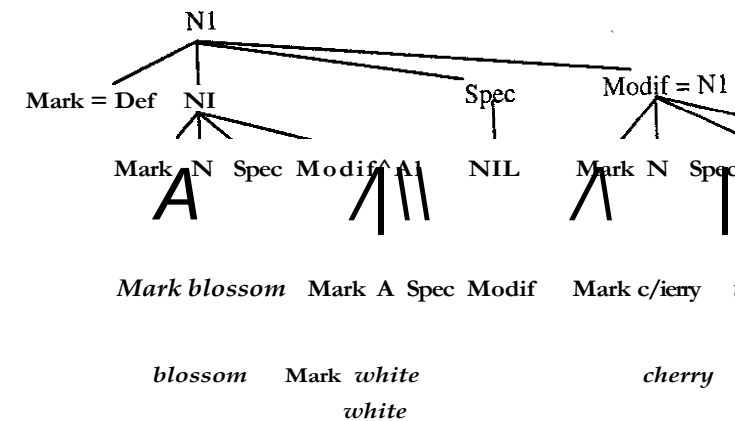
Este important să remarcăm că schema FX-bar propusă principal, relațiile de *dependență* dintre diferitele categorii, sintagme dintr-o frază, împreună cu markerii corespunzători care controlează, comportamentul lor distribuțional. Deoarece am văzut în ce măsură argumentelor este (parametric) dependentă de limbaj în schemele FX-bar pot codifica nu numai situații în care argumentele succed (situația obișnuită care ele preced nucleul lor semantic (Ex.3.2.2.), dar și în care aceluiasi nucleu sunt interschimbabile. Deci aceleiasi scheme FX-bar atribui mai multe formule logice corespunzătoare "echivalente".

4.1. De la text la scheme FX-bar

Strategia SCD propune următoarele scheme FX-bar pentru ... mai jos. Deși muciile ale căror noduri sunt Modif sau Specif su dreapta nucleului corespunzător (pentru conveniențe grafice), ele tre ca având rol funcțional (situate la stânga și aplicându-se nucleului X1 unii adjuncți, la nivel X2. Diferențele dintre codificarea formei pentru e pentru română sunt nesemnificative (cu excepția unor aspecte sup acord, care sunt puse în evidență). Forma codificată a textului engleză este un argument suplimentar pentru versatilitatea schemelor propuse.

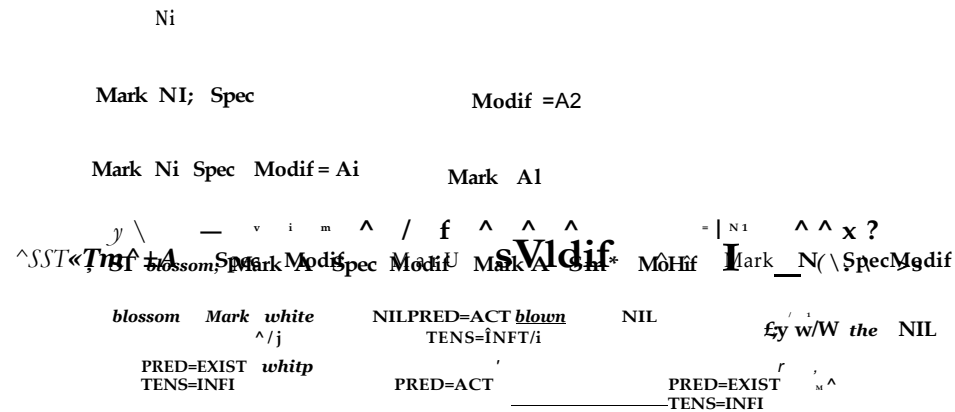
~Ex. 4.1.1.R. //floare albă/de cireș /

Ex. 4.1.1.E. / the cherry / white blossom /



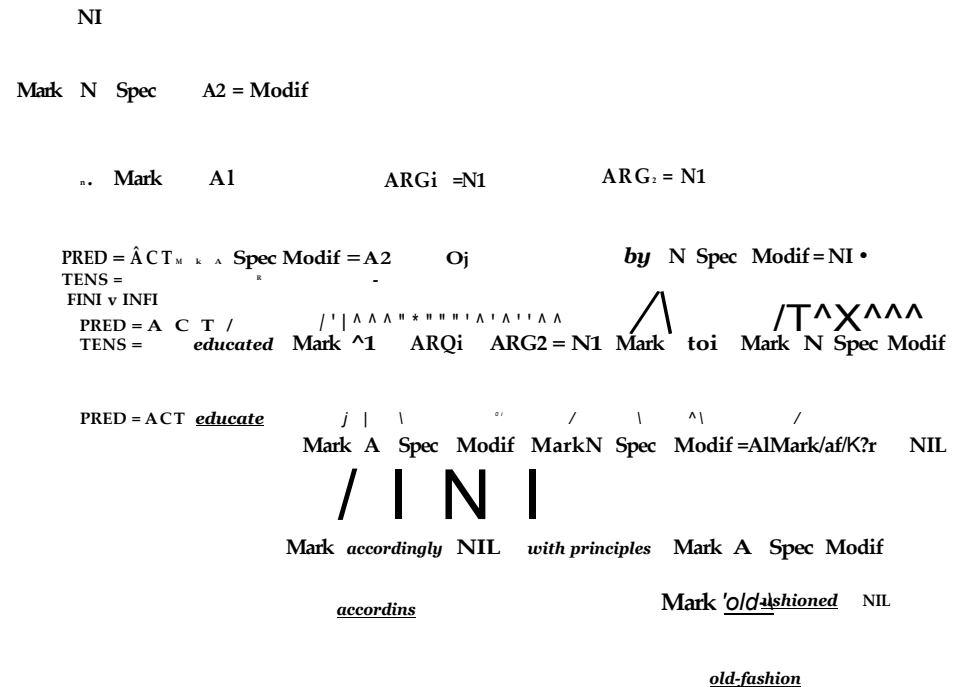
LR.4.1.1.R. de(cireş)(albă floare(X));
 LR.4.1.1.E. quant(indef, X, white(blossom(X)), cherry(X)).

Ex. 4.1.2.R. //floare albă/ //bătută//de vânt/
 Ex. 4.1.2.E./tf7e white blossom/ //blown//bv the wind/
 object, = O,; eventj = e,



LR.4.1.2.R. albă(floare(X)) A bătută(de(vânt(Y)), X);
 LR.4.1.2.E. quant(indef, X, white(blossom(X)),
 quant(indef, Y, by(the(wind(Y))), blown(Y, X))).

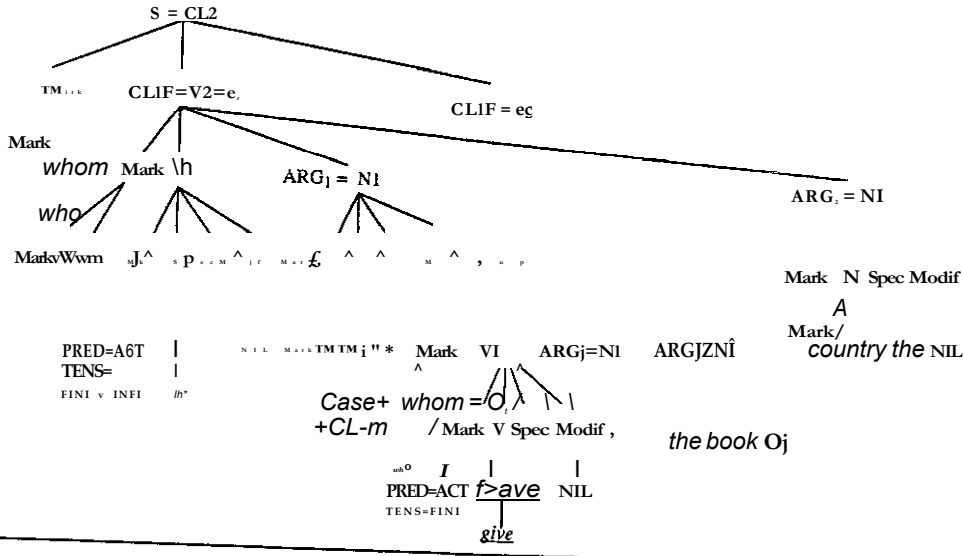
Ex. 4.1.3.R. //educat//[de tatăl său//corespunzător//cu vechile principii//
 Ex. 4.1.3.E. //educated//[bv his father//accordinalv// with old-fashioned principles/



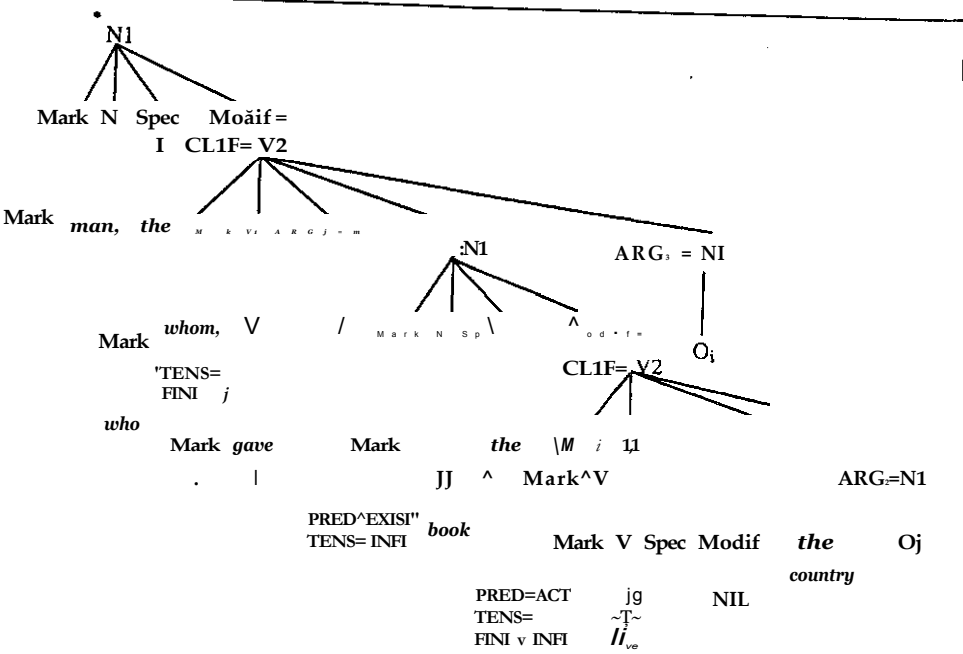
LF.4.1.3.R. corespunzător X, cu(vechile(principii(Y)))
 (educaț(X, de(său(tatăl(Z))));
 LF.4.1.3.E. quant(indef, X, educated(X, by(his(father(Z))))), quant(indef, Y, with(
 old(principles(Y))), accordinglv(X, Y))).

Ex.4.1.4.R. // Omul, // căruia, // PRO, / h -am dat // cartea // PROj a părăsit // țara.//
 Ex.4.1.4.E. // The man\ //whom\ / /pave//the book//PRO\ left //the country.//

4.1.4.E. Reading 1 (left = pastjense(feave))



4.1.4.E. Reading 2 (left = past_participle(leave))

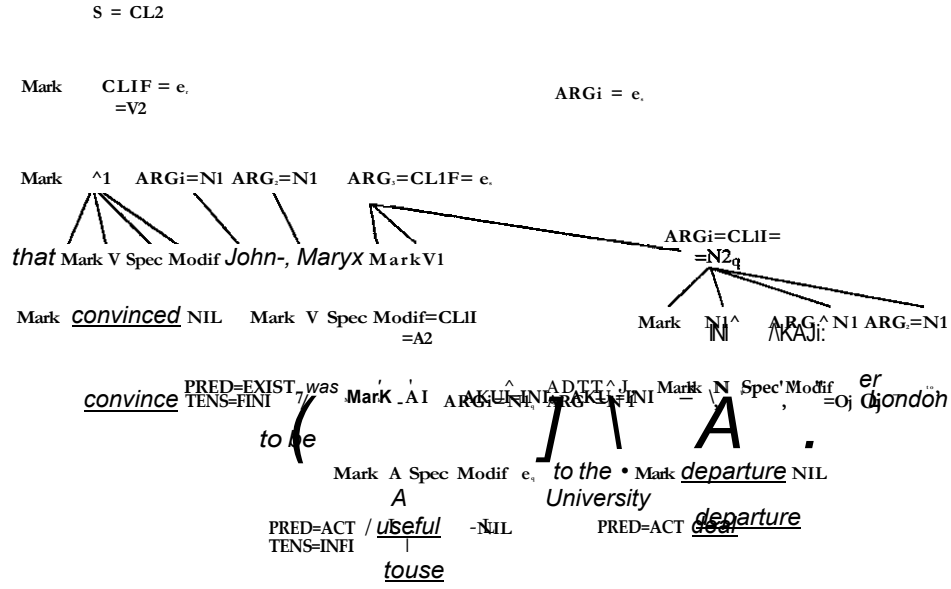


4.1.4.E. Reading 3 ([Eng: left] - [Rom: stanga]).

LF.4.1.4.R. a-părăsit(omul(X) A am-da(Y, cartea(Z), X), țara(T);
 LF.4.1.4.E. quant(def, X, and(man(X), quant(def, Y, l(Y), quant(def, Z, book(Z),
 gavef Y, Z, X))), quant(def, T, country(T), jeft(X, T))).

Ex. 4.1.5.R. //lon\ //a convins^ //pe Mariaj //că //deplasarea* //e/ //la Lodra /
 /a fost utilă //e, / Universității. //
 Ex. 4.1.5.E. //John\//convincd //Mary^ //that //her^ //departure* // to London
 //was useful//e, / to the University. //

objectj = Oj; eventj = ej



LF.4.1.5.R. a-convins(ion, -o(pe(mariaj)), că(a-fost-utilă(
 deplasarea(ei(Xj), la(londra)), universității(Y))));
 LF.4.1.5.E. convinced(john = X, mary = Y, quant(def, X, her(X), departuref X, to(
 london)) = E, quant(def, Z, university(Z), was-useful(E, Z)))).

4.2. Observații generale

(*1) NU este scopul prezentei lucrării să arate cum sunt obținute
 ^prezentările FX-bar ale structurilor LN (într-o manieră mai mult sau mai puțin
 3'goritică), ci doar să propună schema FX-bar generală ca un mecanism esențial
 ^representare a informației lingvistice, să sugereze cum lucrează, și să explice

rațiunile introducerii acestui mecanism. Teoria FX-bar este integrată ca⁴ o componentă importantă a strategiei lingvistice SCD, însă ea poate fi utilizată și în alte contexte computaționale, cu condiția de a include ingredientele necesare, și anume, clasele de PS-Ms, ierarhiile acestor clase, o taxonomie funcțională (predicațională) și relațională a categoriilor majore și a marcherilor, un algoritm (în particular, algoritmi SCD) de obținere a structurilor de dependență, etc. Aspecte mai detaliate ale SCD au fost prezentate în [1], [2], [6], [7], [45]. (*2) Funcționarea corectă a schemelor FX-bar expuse arată clar cât de necesară este utilizarea (intensivă) a trăsăturilor *predicative* și *funcțional-relaționale* pentru fiecare categorie lexicală. Din experiența noastră în ce privește analiza și generarea automată a limbii române [6], considerăm că accentul pus pe trăsăturile funcționale ale categoriilor gramaticale, cuplat cu punerea în evidență a PS-Ms, reprezintă elemente-cheie în utilizarea cu succes a teoriilor X-bar curente în procesarea automată a LN și în cadrul unor teorii lingvistice moderne (UG, FG, HPSG, etc). (A3) Punerea în valoare a trăsăturilor *funcționale* (în particular, predicaționale) ale categoriilor majore N, V, A, și a celor *relaționale ale claselor de marcheri* (marcheri numiți în literatură și "*cue phrases*" [Rom: *sintagme indicatoare*] [28], [31], sau *conective* [29], [30] etc), deși esențiale, nu poate rezolva toate problemele. De exemplu, asignarea dependențelor corecte în juxtapunerea de NGs este o problemă binecunoscut de dificilă, imposibil de rezolvat complet doar la nivel sintactic. Există însă în prezent un puternic curent către acest tip de abordări, aceasta deoarece ele reflectă mult mai adecvat structura reală a textului de LN (cel puțin pentru o clasă largă de LNs europene). Aceste abordări pot diferi substanțial în instrumentele și tehnicile de parsare, însă principiile rămân foarte similare (de exemplu, [19], [29], [31], [33], etc (*4) PS-Ms (marcherii de structuri sintagmatice) joacă un rol fundamental în delimitarea structurilor sintactice și semantice, și stabilirea dependențelor corecte între aceste structuri, SCD a pus accentul încă de la începuturi pe acest aspect [22]. Se remarcă în prezent o întreagă mișcare către reconsiderarea rolului esențial al marcherilor, în special la nivel de *discurs* și în analize complexe ale marilor unități textuale (regăsirea informației, rezumare automată, planificare și generare automată de text, etc). Strategia SCD, cu componenta ei de teorie FX-bar, încearcă să pună la lucru întreaga paletă de PS-Ms, de la nivel *lexical* și de *coeziune* (locală), până la nivel de *discurs* (*coeziune* și *coerență* globală), punând accentul pe sintaxă (nivelul de "*suprafață*", [Eng: *shallow*]) și pe un nivel minimal de semantism. În funcție de problema de LN ce trebuie rezolvată, acest nivel poate fi amplificat în mod corespunzător, (A5) Cuplarea schemelor FX-bar cu: (a) clasele de marcheri SCD și cu ierarhia lor ce corespunde celor *patru* nivele de proiecție lingvistică din FX-bar [7], [45]; (b) o taxonomie bazată pe predicaționalitate a categoriilor majore N, V, A; (c) exploatarea maximală a trăsăturilor *funcționale* (predicaționale) și *relaționale* ale tuturor categoriilor lexicale și nelexicale (deci și ale PS-Ms); (d) o schemă X-bar simplă și unică, apelată recursiv pe cele patru nivele ale sale, pornind de la *lexicon* (convențional, BAR = -1) și până la nivelul de *discurs* al frazei multi-eveniment (BAR = 3), aceste aspecte reprezintă principalele diferențe (și noutăți) dintre teoria FX-bar și teoriile X-bar precedente, (A6) Schema FX-bar poate fi de asemenea

asociată cu un automat dependent de limbaj (pentru o largă clasă de LNs), care începe să lucreze pentru fiecare frază, primește *on-lirie* cuvânt cu cuvânt, și se oprește odată cu semnul de punctuație final al frazei. Pentru valori adecvate ale parametrilor de LN cum sunt *ordinea cuvintelor (argumentelor)* și *direcția proiecției lingvistice* pentru categoriile majore și pentru marcheri, schema FX-bar poate reprezenta corect dependențele structurilor lingvistice (inclusiv pentru Ex.3.2.2.G).

5. Problema X-bar teoriei actuale

Mai este necesară X-bar teoria sau nu? Este teoria X-bar pe moarte sau nu? Care este valoarea teoretică și, mai ales, practică a (sub)teoriei X-bar în cadrul teoriilor lingvistice și al tehnologiilor actuale ale LN? Cum trebuie să percepem în mod corect X-bar teoria atunci când, în aceeași carte a lui Chomsky, găsim următoarele două pasaje:

(Chomsky1): "*The concepts of X-bar theory are therefore fundamental. In a minimalist theory, the crucial properties and relations will be stated in the simple and elementary terms of X-bar theory.*" [3, p. 172],

(Chomsky2): "*Standard X-bar theory is thus largely eliminated in favor of bare essentials.*" [3; p. 246].

Subliniem că aceste citate nu sunt extrase din text astfel încât să nu aibă relevanță în context, cu intenția de a provoca confuzie. Dimpotrivă! De asemenea, scopul nostru nu este de a căuta o posibilă incoerență ci de a pune în evidență noua poziție a lui Noam Chomsky, între 1992 și 1995. Încercăm să deschidem o discuție pe această temă deoarece considerăm că există o problemă, și că ea este de o reală importanță.

În această secțiune urmărim cinci obiective: **(A)** Să enunțăm problema X-bar teoriei. **(B)** Să rezumăm soluțiile existente în momentul de față. **(C)** Să stabilim rolul X-bar teoriei în interiorul contextului teoriilor lingvistice și să sugerăm posibile dezvoltări. **(D)** Să specificăm poziția FX-bar schemelor propuse privitor la dilema eliminării complete a X-bar teoriei și, în special, relația noii FX-bar teorii conturate în contextul strategiei lingvistice SCD. **(E)** Câteva concluzii și perspective.

(A) Să considerăm următoarea *problemă*: reflectă *teoria X-bar* o *realitate lingvistică* a LNs, și dacă da, prin ce *mijloace* această realitate lingvistică ar putea fi *cel mai bine* reflectată? Proiecția categoriilor lingvistice este un fapt lingvistic de netăgăduit. Chomsky și alți distinși lingviști au fost în completă eroare în ultimii 25-30 de ani? Credem că nu. Problema este dacă teoria X-bar poate încă să mai fie un *bun model*, sau vehicul, care să exprime acest *fapt*, și cu ce preț de utilitate. *Principiul Proiecției Extinse* [3, p. 55] și *Principiul Proiecției Maximale* (propus în [19] și secțiunea 3.1.) au ca scop să stabilească forma și marginile cele mai probabile ale unităților textuale obținute în cadrul procesului de proiecție a categoriilor lingvistice.

(B) Ipoteza (Chomsky1) de mai sus dă un răspuns afirmativ la această întrebare în timp ce (Chomsky2) reprezintă, aparent, opusul acestui răspuns. Abordarea din [3, Cap. *Categories and Transformations*] pentru ipoteza (Chomsky2) este că disoluția schemelor X-bar, deci a proiecției categoriilor lingvistice, poate fi înlocuită cu succes prin folosirea proprietăților de funcționalitate, predicativitate, tipologie și transformare intrinseci acestor categorii, deși aceste proprietăți sunt reprezentate în [3] cu același aparat X-bar pe care îl combat! În cadrul unei teorii a "*structurii sintagmatice pure*", operațiile unui sistem computațional al LN "*construiesc recursiv obiecte sintactice*", iar "*categoriile sunt construcții elementare rezultate din proprietățile elementelor lexicale*", cu condiția "*să nu fie adăugate obiecte noi în cursul procesării, înafară de rearanjări ale proprietăților lexicale*" [3]. Rezultatul pare să fie spectacular: dispar nivelele de proiecție (în sensul teoriei X-bar), astfel spus, nu se face nici o deosebire între elementele lexicale și nucleele proiectate din ele, în timp ce "*teoria structurilor sintagmatice poate fi eliminată în întregime, se pare, pe baza celor mai elementare ipoteze*" [3, p. 294].

Nu ar fi pentru întâia oară când teoria lingvistică încearcă să renunțe la (sub)teoria X-bar. Chomsky sugerează că nivelele de proiecție lingvistică pot fi înlocuite de către "*proprietățile (funcționale n.n.) ale elementelor lexicale*". Acesta este chiar cazul *gramaticii funcționale* (FG) [25] în care, formal, lipsește teoria X-bar. Dar chiar și în gramatica funcțională a lui Dik, conținutul ascuns al teoriei X-bar este scufundat de fapt în cele *patru nivele* de structuri ierarhice ale *functorilor* și *operatorilor* ce se aplică pe categoriile și structurile cu care FG lucrează la fiecare nivel sintactic. O situație specială avem în SCD, unde nivelele de proiecție a categoriilor lingvistice sunt recuperate pe baza unei funcționalități ierarhice a elementelor lexicale, iar FX-bar schema propusă poate fi utilizată (recursiv) ca un invariant sintactic constructiv al structurilor sintagmatice în cadrul proceselor de analiză și generare automată a LN (limbii române).

Schema FX-bar propusă (Fig. 3.2.1.) poate fi considerată ca un compromis, o negociere, între (Chomsky1) și (Chomsky2), deoarece (Chomsky2) se prezintă fără mecanisme concrete pentru a-și susține ipoteza: în timp ce teoriile X-bar clasice nu mai pot fi utilizate ca instrumente operaționale pentru a reflecta o viziune exclusiv funcțională (și relațională) asupra sintaxei, teoria FX-bar propusă poate face acest lucru.

(C) Poziția noastră privind problema (A) asupra teoriei X-bar poate fi rezumată astfel: (C1) Proiecția categoriilor gramaticale este un fapt lingvistic. (C2) Acest fapt poate fi corect reflectat prin "*nuclee*" și "*nivele (bar) cle proiecție*" în interiorul schemelor X-bar, dar și prin proprietățile funcționale "*intrinsec?*" ale categoriilor lexicale și gramaticale. (C3) Teoria X-bar include deci o componentă de adevărată construcție lingvistică, iar ingredientul său de bază este confecționat din relațiile funcționale stabilite între elementele lexicale (și nelexicale) conținute în cadrul schemelor X-bar. (C4) Atunci când proprietățile funcționale ale categoriilor lexicale nu sunt evaluate și exploatate corespunzător, teoria X-bar este inconsistentă și produce dificultăți de calcul și rezultate incorecte. (C5) Acestea

sunt consecințele unui aspect mult mai general, și anume că teoria X-bar nu trebuie să fie văzută ca o teorie gramaticală singulară, construită pentru sine, ci ca un dispozitiv component al unui *mecanism lingvistic teoretic și computațional mai generat*, ale cărui principii să guverneze teoria X-bar. Axiomatica (bazele constructive ale) teoriei X-bar trebuie să fie un *rezultat* al bunei ei funcționări, pe fenomenele concrete de limbaj, și nu invers! (C6) *Ad limitum*, se poate concepe că mecanismul lingvistic teoretic menționat mai înainte poate funcționa și fără includerea dispozitivului reprezentat de teoria X-bar, așa cum încearcă teoria MinP să propună în [3, Cap. *Categories and Transformations*] (dar folosindu-se în explicare tot de aparatul de reprezentare al teoriei X-bar), precum și în cazul FG [25].

(D) Considerăm că schemele (funcționale) FX-bar propuse furnizează un (sub)sistem necesar și folositor în cadrul oricărei teorii sintactice asupra LN, inclusiv (și în special) pentru strategia lingvistică SCD. O condiție esențială pentru schemele FX-bar este ca ele să reflecte corespunzător proprietățile *funcționale* și *relaționale* ale tuturor categoriilor *lexicale* și *gramaticale*. Exemplele 4.1.1.-4.1.5 arată cum sunt construite schemele FX-bar, cum se obțin (prin apel recursiv pe nivele) structurile sintagmatice complexe ale LN, și cum acestea rămân închise la operatorul de compunere (adjuncție) pe baza principiilor și regulilor SCD.

Schimbând perspectiva, prin definirea teoriei FX-bar ca o componentă a strategiei lingvistice SCD, și parafrazând formalismul bine-cunoscut al gramaticilor TAG [Eng: *tree adjoining grammar*], strategia SCD poate fi văzută și ca o teorie a evaluării și adjuncției de FX-bar scheme. Este doar o mostră a rolului important pe care teoria X-bar îl poate încă juca în cadrul teoriei și tehnologiei LN.

(E) Un *element original* propus de schemele FX-bar în peisajul teoriilor X-bar cunoscute este rolul lor dublu ce îl pot juca în cadrul strategiei SCD (și nu numai): Schemele FX-bar pentru X = N, V, A, CL (CL = clauză) trebuie conceput ca un set de invarianți sintactici (dinamici) ce pot fi folosiți (1) la *reprezentarea* informației lingvistice la nivel de lexicon (în mod similar cu structurile de trăsătură lingvistice [18], dar într-o manieră mai simplă și mai regulată), și (2) la *procesarea (analizarea și generarea) automată* de text în LN (inclusiv, și mai ales, pentru limbă română), de la structurile sintagmatice simple până la cele de discurs.

Derivarea de *automate* și *gramatici formale bazate* pe schema FX-bar pentru analiza LN, ar fi o consecință normală și o provocare a prezentei propunerii. Modul *recursiv*, *ascendent* și *incremental* (prin apelul de funcții și relații cu nivel lingvistic multiplu), dar și *descendent* (bazat pe sateliții nucleelor semantice) utilizarea la maximum a *contextualității* marcherilor de toate tipurile poate reprezenta o motivație naturală pentru cercetarea relației dintre strategia SCD (cu componenta ei de *teorie* FX-bar), și modelele generative oferite de către *gramaticile contextuale* Marcus [41], [42], un formalism *context-dependent* puternic, destinat parsării dar și analizei semantice și de discurs (articularea *topic-focus* [37]) a LN. *Gramatici contextuale* Marcus aparțin unei serii de formalisme care includ *gramatici* TAG [4], *dramatici orientate-nucleu* [15], [16], *gramatici indexate*, *gramatici X-bar*, *gramatici context-free marcate* [44] etc, formalisme ce realizează o modelare mai realistă a comportamentului sintactic, semantic și discursiv al LN.

Referințe bibliografice

- [1] N. Curteanu (1990). *A Marker-Hierarchy-based Approach Supporting the SCD Parsing Strategy*. Research Report no. 18, Institute of Technical Cybernetics, Bratislava.
- [2] N. Curteanu (1994). *From Morphology to Discourse Through Marker Structures in the SCD Parsing Strategy. A Marker-Hierarchy Based Approach*. Language and Cybernetics, Akademia Libroservo, Prague, 61-73.
- [3] Noam Chomsky (1995). *The Minimalist Program*. The MIT Press, Cambridge, Massachusetts.
- [4] N. Chomsky (1970). *Remarks on Nominalizations*. In R. Jacobs and P. Rosenbaum (eds.), *Readings in Transformational Grammar*, Ginn and Co., Boston, 184-221.
- [5] T. Stowell (1981). *Origins of Phrase Structure*. Ph.D. Dissertation, Dept. of Linguistics and Philosophy, MIT, Cambridge.
- [6] N. Curteanu, G. Holban (1996). *Strategia lingvistică SCD aplicată la analiza și generarea limbii române*. *Limbaj și Tehnologie* (Dan Tufiș, Ed.), Academia Română, București, p. 169-176.
- [7] N. Curteanu, C. Linteș (2002). *Segmentation Algorithms for Clause-Type Textual Units*, Research Report, Institute of Theoretical Informatics, Romanian Academy.
- [8] Noam Chomsky (1986). *Barriers*. The MIT Press, Cambridge.
- [9] Noam Chomsky (1981). *Lectures on Government and Binding*. Foris, Dordrecht.
- [10] Simon C. Dik (1989). *The Theory of Functional Grammar*. Foris Publishers, Dordrecht.
- [11] Cari Pollard, Ivan Sag (1994). *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago & London.
- [12] Gerald Gazdar, E. Klein, G. Pullum, I. Sag (1985). *Generalized Phrase Structure Grammar*. Harvard University Press, Cambridge, Massachusetts.
- [13] Peter Sells (1985). *Lectures on Contemporary Syntactic Theories*. CSLI, Stanford, California.
- [14] Stuart Shieber (1986). *An Introduction to Unification-Based Approaches to Grammar*. CSLI, Stanford, California.
- [15] Cari Pollard, Ivan Sag (1987). *Information-based Syntax and Semantics*. CSLI, Stanford, California.
- [16] Cari Pollard, Ivan Sag (1994). *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago & London.
- [17] E.P. Stabler Jr. (1992). *The Logical Approach to Syntax: Foundations, Specifications and Implementations of Theories of Government and Binding*. The MIT Press, Cambridge, Massachusetts.
- [18] N. Curteanu, G. Holban (2000). *A Set-Theoretic Approach to Linguistic Feature Structures and Unification Algorithms* (I, II). *Computer Science Journal of Moldova*, 8(2): 116-149, 8(3): 223-246.
- [19] Neculai Curteanu (1988). *Augmented X-bar Schemes*. COLING'88 Proceedings, Budapest, 130-132.
- [20] Neculai Curteanu, A. Todirașcu, G. Holban (1997). *Teorii sintactice ale limbajului natural*. Raport de cercetare, Institutul de Informatică Teoretică, Academia Română, Iași, 66 p.
- [21] Alain Lecomte (1998). *Multimodal Logic for Syntax*. *Logica Trianguli*, 2: 49-72.
- [22] Neculai Curteanu (1983). *Algoritmi de analiză sintactică a frazei și propoziției românești*. *INFO-IAȘI'83*, p. 533-549.
- [23] M. Moortgat (1997). *Categorial Type Logics*. *Handbook of Logic and Language*, Elsevier.
- [24] E.P. Stabler Jr. (1997). *Derivational Minimalism*. *Logical Aspects of Computational Linguistics*, LNCS no. 1328, Springer-Verlag, Berlin.
- [25] Simon Dik (1989). *The Theory of Functional Grammar*. Foris Publishers, Dordrecht.
- [26] Robert Kasper (1993). *Adjuncts in the Mittelfeld*. în "German Grammar in HPSG" (J. Nerbonne et al., Eds.), CSLI, Stanford, California.
- [27] Denis Bouchard (1995). *The Semantics of Syntax. A Minimalist Approach to Grammar*. The Univ. of Chicago Press, Chicago & London.
- [28] Julia Hirschberg, D. Litman (1993). *Empirical Studies on the Disambiguation of Cue Phrases*. *Computational Linguistics* 19(3): 501-530.
- [29] Jacques Jayez, C. Rossari (1999). *Pragmatic Connectives as Predicates. The Case of Inferential Connectives*. în "Predicative Forms in Natural Language and in Lexical Knowledge Bases" (P. Saint-Dizier, Ed.), Kluwer Academic Publishers, Dordrecht.
- [30] Patrick Saint-Dizier (Ed.) (1999). *Predicative Forms in Natural Language and in Lexical Knowledge Bases*. Kluwer Academic Publishers, Dordrecht.
- [31] Daniel Marcu (2000). *The Theory and Practice of Discourse Parsing and Summarization*. The MIT Press, Cambridge.
- [32] V. Raskin, S. Nirenburg (1999). *Lexical Rules for Deverbal Adjectives*. "Breadth and Depth of Semantic Lexicons", Kluwer Academic Publishers, Dordrecht.

- [33] O. Popârda, N. Curteanu (2002). *L'evolution du discours juridique frangais analyse par la strategie linguistique SCD*. In "Representation du Sens Linguistique" (D. Bouchard, Ed.), LINCOS Studies in Theoretical Linguistics, LINCOS EUROPA, Munchen.
- [34] Noam Chomsky (2000). *Minimalist inquiries: the framework*. în R. Martin *et al.* (Eds) "Step by step. Essays on Minimalist Syntax in Honor of Howard Lasnik", MIT Press, Cambridge, p. 89-155.
- [35] Noam Chomsky (2001). *Derivation by phase*: în M. Kenstowicz (Ed.) "Ken Hale: a life in language", MIT Press, Cambridge, p. 1-52.
- [36] Jane Morris; G. Hirst (1991). *Lexical Cohesion Computed by Thesaural Relations as an Indicator of the Structure of Text*. Computational Linguistics 17(1): 21-48.
- [37] Eva Hajicova, H. Skoumalova, P. Sgall (1995). *An Automatic Procedure for Topic-Focus Identification*. Computational Linguistics, 21(1): 81-94.
- [38] P. Sgall, E. Hajicova, J. Panevova (1986). *The Meaning of the Sentence in Its Semantic and Pragmatic Aspects*. Kluwer Academic Publishers, Dordrecht.
- [39] S. Shieber, F. Pereira, G. Van Noord, R. Moore (1990). *Semantic Head-Driven Generation*. Computational Linguistics 16(1): 30-41.
- [40] Șteven Abney (1996). *Part-Of-Speech Tagging and Parțial Parsing*. în "Corpus-Based Methods in Language and Speech", (K. Church *et al.*, Eds.), Kluwer Acad. Publishers, Dordrecht.
- [41] Solomon Marcus (1997). *Contextual Grammars and Natural Language*. în Cap. 5 (Voi. 2) din "The Handbook of Formal Languages", G. Rozenberg, A. Salomaa, Eds., Springer-Verlag, Berlin, 215-235.
- [42] Gheorghe Păun (1997). *Marcus Contextual Grammars*. Kluwer Academic Publishers, Dordrecht.
- [43] Michele Abrusci, Christophe Fouquere, Jacqueline Vauzeille (1999). *Tree Adjoining Grammars in a Fragment of the Lambek Calculus*. Computational Linguistics, 25(2): 209-236.
- [44] Philip Miller (1999). *Strong Generative Capacity. The Semantics of Linguistic Formalism*. CSLI Publications, Stanford, California.
- [45] D. Gâlea, N. Curteanu, C. Linteș (2002). *Algoritmi de segmentare a textului în unități de tip causal*. (în prezentul volum)

Teoria HPSG. Studiu de caz: acordul încrucișat

Ana-Maria BARBU

RACAI, Calea 13 Septembrie nr.13, București

abarbu@racai.ro

1. Introducere

Oricât ar fi de mare entuziasmul creat de performanțele realizate de calculatorul, care cuprinde deopotrivă și domeniul prelucrării limbajului natural, rezultate temeinice nu se pot obține dacă acestea nu sunt fundamentate pe îndelungi și profunde analize teoretice. Nu putem aspira la obiective majore în ingineria lingvistică, precum analizarea și generarea de texte, construirea de verificatoare ortografice și gramaticale sau chiar de traducătoare automate, dacă se ignoră particularitățile inerente ale obiectului în studiu, anume ale limbajului natural în general, și a limbii de aplicație, în special. Or aceste particularități sunt oferite, sub un aspect sau altul, tocmai de teoriile gramaticale. Experiența a dovedit că eșecurile din ingineria lingvistică au avut ca posibile surse eșecurile descrierea corespunzătoare a fenomenelor de limbă, dar și succesele, la rândul lor, s-au datorat în parte acurateței, exactității, și nu în ultimul rând caracteristicilor computaționale ale unui model gramatical teoretic.

Iată de ce alegerea unei teorii lingvistice adecvate, cu scopul de a scrie pe baza acesteia o gramatică computațională a unei limbi particulare, în speță a limbii române, este un act de primă însemnătate.

După anii primelor dezvoltări ale gramaticii generative, sintaxa formală este, de aproape două decenii, repusă în discuție ca obiect de studiu autonom și distinct în aceeași timp de cel al lexicului și cel al sensului. Mai multe curente teoretice, cunoscute sub numele generic de "gramatici de unificare" sau "gramatici bazate pe constrângeri", s-au născut din această reconsiderare a sintaxei. E vorba de modele recente (cele mai vechi datând de la începutul anilor '70) dezvoltate în cea mai mare parte în Statele Unite, și în general aproape necunoscute publicului român. Aceste modele se pretează scrierii de gramatici pentru calculator, dar ambiția lor este mai întâi de a constitui teorii lingvistice sine stătătoare. Autorii lor se înscriu pe linia programului gramaticii generative chomskyene din 1957, de la care preiau grija pentru o formalizare operatorie a sintaxei, dar se disting suficient de modelul actual al Școlii de la Cambridge (născut din Government and Binding) pentru a prezenta teorii alternative. Printre puncte

comune ale gramaticilor de unificare, se află pe de o parte atenția acordată unei articulări mai explicită a lexicului, sintaxei și semanticii, pe de altă parte accentul pus pe descrierile lingvistice și recurgerea la un stil de analiză sintactică mai "concret", care limitează recurgerea la elemente "vide" (nerealizate concret) și care restrânge numărul etapelor intermediare în producerea unui enunț.

În acest articol vom prezenta pe scurt una dintre teoriile lingvistice amintite, anume "Gramatica sintagmatică ghidată de centru", denumită abreviat HPSG după numele său din engleză "Head-driven Phrase Structure Grammar". Apoi vom ilustra modul în care poate fi aplicată această teorie în reprezentarea unui fenomen mai special de limbă română prin aceea că presupune dependențe încrucișate de acord. Este vorba de structuri relative de tipul *băiatul a cărui soră cântă* unde articolul genitival a se acordă cu substantivul *soră*, iar pronumele relativ *cărui* se acordă cu substantivul *băiatul*.

2. Teoria lingvistică HPSG

2.1. Scurt istoric

Modelul gramaticii sintagmatică ghidate de centru (engl. *Head-driven Phrase Structure Grammar*, sau HPSG) a fost conceput la începutul anilor '80 de Cari Pollard și Ivan Sag cu scopul de a permite o integrare mai explicită a diferitelor nivele de analiză lingvistică: fonetic, sintactic și semantic. El a luat naștere în principal din Gramatica Sintagmatică Generalizată (GPSG) și din lucrările lui C. Pollard despre *Head Grammar* [1], dar autorii lor s-au inspirat deopotrivă din numeroase alte teorii. Ei au preluat de la modelul chomskyan al Guvernării și Anaforicității (GB) noțiunea de modularitate și recurgerea la principii foarte generale (Principiul anaforicității, al controlului etc). De la Gramatica Funcțională de Unificare (FUG) [2] au împrumutat reprezentarea uniformă a elementelor lexicale, a sintagmelor și regulilor gramaticale sub formă de structuri de trăsături. S-au inspirat de la Gramatica Lexical Funcțională (LFG) pentru îmbogățirea cadrelor de subcategorizare și a noțiunii de regulă lexicală. Au luat de la gramaticile categoriale ideea de saturare progresivă a predicatelor și recurgerea la o ierarhie de funcții gramaticale (cf. [3]). S-au inspirat, în sfârșit, dintr-un punct de vedere mai formal, din lucrări de logică și informatică asupra tipurilor și moștenirii.

Teoria este prezentată în cele două lucrări ale lui C. Pollard și Ivan Sag: [4] și [5]. Majoritatea exemplurilor privesc limba engleză și tratează fenomene variate: fenomene de acord, construcții infinitivale, anafore, construcții relative și comparative. Fenomenele de control sunt totodată dezvoltate în [6], iar o analiză a anaforelor este propusă în [7]. Primele lucrări au conferit de la bun început o dimensiune multilinguală acestei teorii prin abordări privind germana ([8], [9]), catalana ([10]), japoneza ([11]), dar și coreana ([12]), franceza ([13]) și italiana ([14]).

C. Pollard și I. Sag preiau din modelul GPSG noțiunea de gramatică sintagmatică, cu distincția între o componentă ierarhică (scheme DI -de dominanță imediată) și o componentă liniară (principii de precedență liniară), precum și recurgerea la principii foarte generale de partaj și de propagare a trăsăturilor. Totuși ei se separă de modelul original în câteva puncte. Structurile sintagmatică sunt în întregime exprimate în termeni de structuri de trăsături, cu introducerea unui atribut Ramuri. Structurile de trăsături sunt la rândul lor organizate în ierarhii de tipuri, comportând fiecare trăsături predefinite. Modelul HPSG oferă astfel anumite simplificări în raport cu GPSG: întregul arsenal de reguli DI este redus la șase scheme de bază; metareguliile sunt eliminate în favoarea regulilor lexicale. S-a urmărit deosebirea clară între ceea ce ține de domeniul constrângerilor universale și ceea ce ține de descrierea unei limbi particulare. Principiile de coocurență a trăsăturilor din GPSG, care amestecă constrângerile universale și cele specifice unei limbi date, au fost suprimate.

2.2 Organizarea generală a HPSG

2.2.1 Caracteristici specifice gramaticilor de unificare

Se poate considera că gramaticile de unificare, sau gramaticile bazate pe constrângeri, reprezintă noile teorii sintactice ale anilor '80. Este vorba de modele care urmăresc o articulare explicită între lexic, sintaxă și semantică. Proprietățile lingvistice corespunzătoare sunt concepute ca "informații" asociate morfemelor sintagmelor sau construcțiilor, combinate prin operații variate, dintre care unificarea ocupă un rol central. Această concepție "integratoare" este unul dintr-aturile lor pentru tratarea automată a limbajelor naturale. Un alt avantaj este că ele se bazează pe modele logice sau matematice (gramatici de constituenți structuri de trăsături), pentru care au fost definite metode de programare. Ele sunt în general rezultatul unui compromis între expresivitatea lingvistică (grija de a facilita exprimarea diferitor principii lingvistice adăugându-se variante notaționale sau operatori) și eficacitate (notații concentrate, puține operații).

Aici, ne vom rezuma să punctăm trăsăturile lor comune cele mai pregnante, dintre care:

- reabilitarea descrierilor de suprafață;
- reînnoirea descrierilor sintactice prin definirea de trăsături complexe;
- definirea de principii generale de bună formare a enunțurilor;
- integrarea lexicului, sintaxei și semanticii.

Gramaticile de unificare îmbogățesc aparatul formal al gramaticilor de constituenți cu un număr de noțiuni importante. În acest capitol ne vom limita la prezentarea principalelor noțiuni utilizate pe parcursul lucrării, pentru detalii putând fi consultate S. Shieber [21] sau H. Uszkoreit [38].

2.2.1.1 Structuri de trăsături

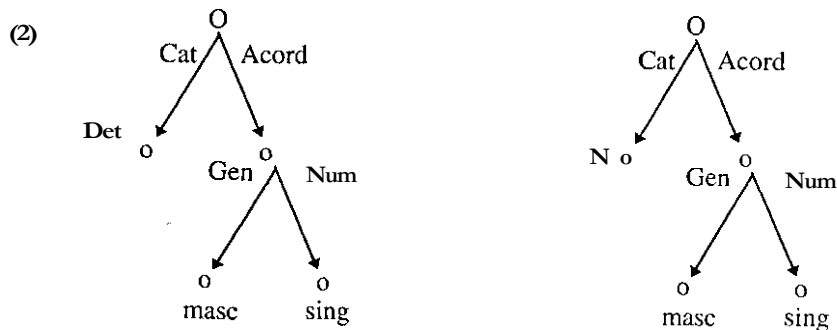
Structurile de trăsături (engl. *feature structure*) sunt primitive ale teoriilor sintactice bazate pe unificare și reprezintă ansambluri de trăsături, numite și complexe de trăsături (engl. *feature complexes* sau *feature bundles*), care pot fi reprezentate sub formă de matrice. O **trăsătură** este o pereche atribut-valoare, valorile putând fi simboluri atomice sau trăsături. Trăsăturile cu valoare non atomică conduc la structuri de trăsături care prezintă îmbricări.

Spre exemplu, cuvintelor *acest* și *câine* li se asociază o trăsătură *Cat* cu valoare atomică (pentru categorie) și o trăsătură complexă *Acord* care ia ca valoare conjuncția a două trăsături *Num* (pentru număr) și *Gen*:

(1)	acest	câine
	Cat = Det	Cat = N
	Acord = Gen = mase Num = sing	Acord Gen = mase Num = sing

O structură este rău formată când conține de două ori același atribut (la același nivel de imbricare) cu o valoare diferită.

Și alte reprezentări de structuri de trăsături (sau structuri atribut-valoare) sunt posibile, fiind echivalente formal. Cele mai utile, pentru implementarea informatică, sunt cele care utilizează grafuri orientate: arcuri care poartă nume de trăsături și punctează spre noduri care sunt etichetate cu valoarea trăsăturii (dacă e vorba de trăsături cu valoare atomică) sau sunt puncte de plecare pentru alte arce (pentru trăsături cu valoare non atomică). De pildă, pentru exemplele de mai sus vom avea următoarele reprezentări:



În termeni de grafuri, echivalentul interdicției ca un același atribut să apară de două ori la același nivel cu valori diferite este interdicția ca două arcuri care poartă aceeași etichetă să puncteze, plecând din același nod, către două noduri

diferite (ceea ce e o restricție generală asupra grafurilor ce corespund automatelor deterministe).

Structurile de grafuri pot fi ciclice sau non ciclice. Acestea din urmă se numesc **grafuri aciclice orientate** (engl. *Directed Acyclic Graph* sau DAG), denumire adesea folosită pentru a desemna structurile de trăsături.

În lucrul cu structuri de trăsături complexe se impun unele distincții, de pildă, între structurile identice și structurile cu valori partajate (sau reentrante). Cele din urmă sunt identice și vor rămâne astfel indiferent de modificările survenite ulterior, ceea ce nu se întâmplă cu primele. În exemplul ce urmează structura de trăsături A comportă două atribute cu valori identice *Acord* și *Num*. În structura B, cele două atribute *Acord* sunt coindexate (prin indicele 1), ceea ce face ca ele să partajeze în mod egal trăsătura [*Num* = sing].

(3).	A:	B:
	Det = [Acord=[Num = sing]]	Det = [Acord = 1] [Num = sing]]
	Nume = [Acord = [Num = sing]]	Nume = [Acord = 1]

Dacă se unifică fiecare din aceste structuri cu structura C de mai jos rezultatul nu va fi același:

C: [Det = [Acord = [Gen = mase]]

(4)	C ^ A:
	Det = [Acord = [Num = sing, Gen = mase]]
	Nume = [Acord = [Num = sing]]

(5)	C ^ B:
	Det = [Acord = 11] [Num = sing, Gen = mase]]
	Nume = [Acord = 1]

După unificare, trăsătura *Acord* îmbricată sub atributul *Nume* va avea și o trăsătură *Gen* specificată în cazul lui C u B, dar nu și în cazul C u A.

În termeni de grafuri, reprezentarea unei structuri reentrante ca B este următoarea:

(6)

B:

Nume	Det
Acord	Acord

Num
I, o
sing

2.2.1.2 Extensiune și unificare

Se definește o relație de **extensiune** între structuri de trăsături după cum urmează:

O structură de trăsături A este o **extensiune** a unei structuri de trăsături B (notându-se $A \supseteq B$) dacă și numai dacă:

- toate trăsăturile cu valoare atomică prezente în B sunt prezente și în A cu aceeași valoare,

- pentru orice trăsătură $\langle \rangle$ cu valoare non atomică, valoarea lui $\langle \rangle$ în A este o extensiune a valorii lui $\langle \rangle$ în B.

De exemplu, structura de trăsături asociată cuvântului *câine* în (1), reluată în (7) stânga, este o extensiune a structurii din (7), dreapta, dar reciproca nu este adevărată pentru că structura din dreapta nu are trăsătura [Num = sing] prezentă în cea a cuvântului *câine*:

	Cat = N	Cat = N
(7)	Acord = Gen = mase Num = sing	Acord \Rightarrow [Gen = mase]

Dacă numărul de atribute nu este limitat se poate obține o infinitate de structuri care sunt extensii ale unei structuri date. Relația inversă a extensiei se numește **subsumare**, A subsuma B dacă și numai dacă B este o extensie a lui A.

Pe baza acestei relații de ordine parțială putem defini o structură de latice, cu o limită superioară și o limită inferioară. Este de notat că aici nu există o relație de ordine strictă pentru că orice structură este o extensie a ei înseși (ADA). Structura care le subsumează pe toate celelalte (pentru care toate celelalte sunt extensii) este structura vidă (notată T), pe care o putem interpreta ca disjuncția tuturor cuplurilor atribut-valoare ale gramaticii. Dacă dorim să plasăm o limită inferioară, structura care va fi o extensie a tuturor celorlalte (care este subsumată de toate celelalte) va fi cea care conține conjuncția tuturor cuplurilor atribut-valoare posibile (notată 1) adică o structură "falsă" sau rău formată.

Această relație de ordine parțială e folosită pentru a defini unificarea. Această operație a luat naștere din cercetările în logică și informatică (limbajul

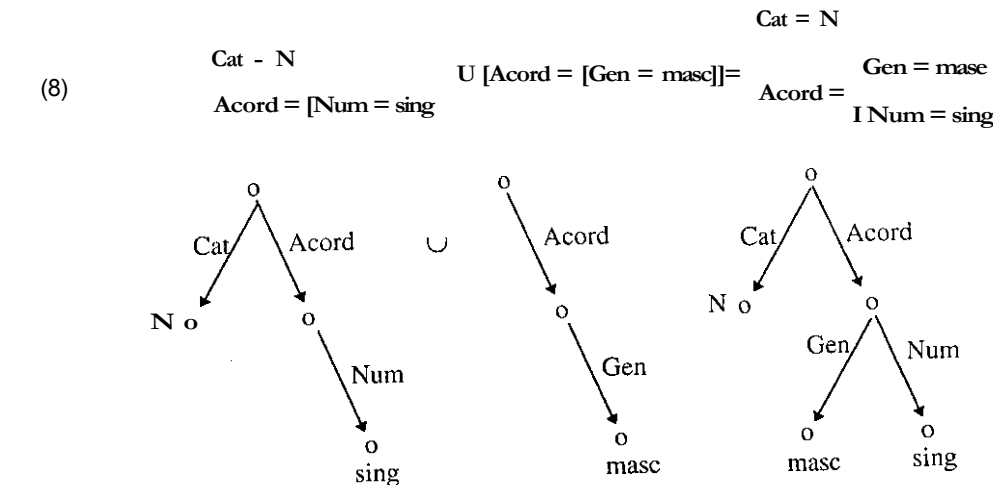
Prolog). Definită la început ca procedură de rezolvare pentru logica predicatelor de ordinul întâi, cf. [15], ea a fost introdusă în lingvistică de A. Colmerauer, [16], apoi de M. Kay, [17], pentru a testa, fuziona și propaga trăsături sintactice. Ea este definită în felul următor:

Unificarea a două structuri de trăsături A și B (notată $A \cup B$) este structura minimală care este în același timp o extensiune a lui A și a lui B. Dacă o astfel de structură nu există, unificarea "eșuează" (ceea ce e notat cu 1).

Altfel spus, unificarea verifică compatibilitatea dintre două structuri de trăsături și produce o structură rezultantă care este cea mai mică structură care conține toată informația din prima structură și toată informația din a doua structură.

Unificarea este o operație idempotentă ($A \cup A = A$), comutativă ($A \cup B = B \cup A$) și asociativă ($A \cup (B \cup C) = (A \cup B) \cup C$), spunem de asemenea că este declarativă (dacă $A = A'$ și $B = B'$ atunci $A \cup B = A' \cup B'$) și monotonă ($A \cup B \supseteq A$ și $A \cup B \supseteq B$; dacă $A \supseteq C$ și $B \supseteq D$ atunci $A \cup B \supseteq C \cup D$), ceea ce vrea să spună că relațiile de extensiune sunt conservate prin unificare. Colocvial spus, unificarea adaugă informație, fără să o scadă.

În termeni de grafuri, echivalentul operației de unificare este fuziunea definită pentru automatele cu număr finit de stări. Pentru exemplul din (8a) se obține reprezentarea grafică din (8b):



Anumiți operatori pot fi adăugați structurilor de trăsături, cei mai utili fiind negația (notată \sim sau $*$ pentru trăsături cu valoare atomică) și disjuncția (notată prin acolade sau semnul \vee). Folosirea negației permite să se renunțe la anumite valori ale disjuncției. Există de exemplu echivalență între următoarele două ecuații, dacă considerăm că atributul Mod are 8 valori posibile în română (indicativ, conjuncție, imperativ, prezumtiv, infinitiv, gerunziu, supin, participiu):

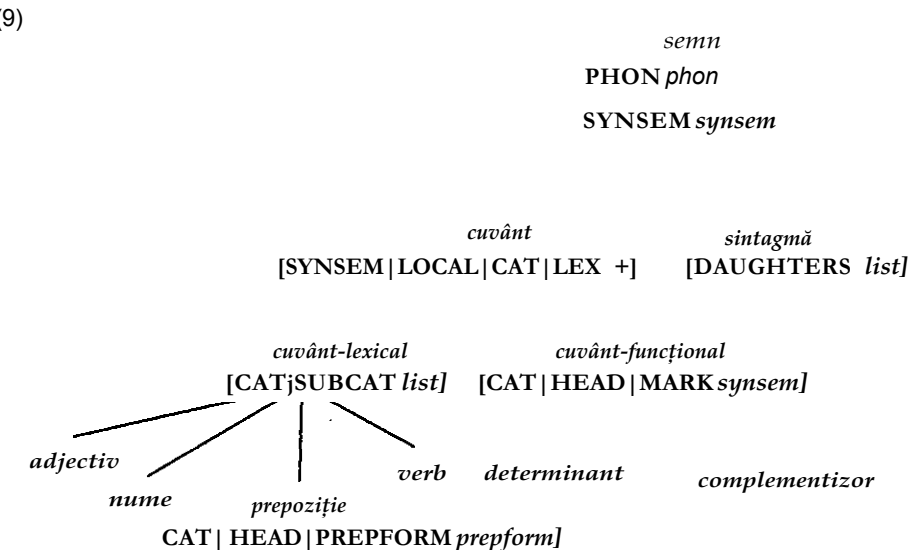
[Mod * inf] [Mod = ind/conj7prez/imp/ger/sup/part].

În secțiunea următoare vom trece la descrierea caracteristicilor specifice ale teoriei HPSG care o fac distinctă de toate celelalte teorii bazate pe unificare. Trebuie spus de la bun început că autorii modelului HPSG au preluat o mulțime de caracteristici ale teoriilor apărute anterior, inclusiv de la gramatica generativă, tocmai din dorința de a aduna într-un singur formalism tot ce e mai adecvat pentru reprezentarea lingvistică în general. Pentru o paralelă detaliată între HPSG și alte teorii bazate pe constrângeri a se vedea [18].

2.2 Caracteristici specifice HPSG

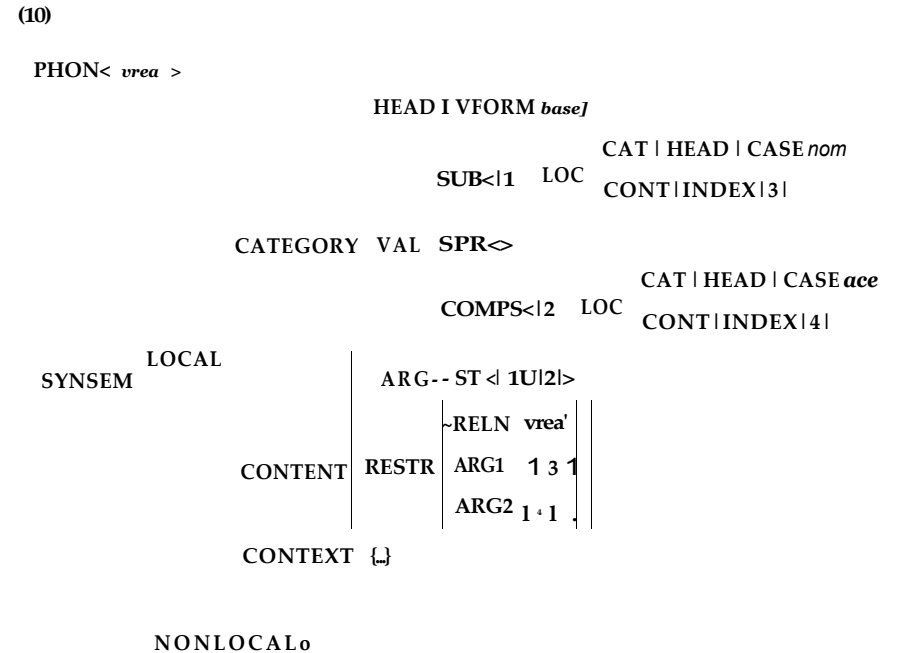
În HPSG, structurile de trăsături, utilizate în LFG pentru reprezentarea funcțiilor gramaticale, iar în GPSG pentru reprezentarea categoriilor, sunt sistematizate pentru a include atât structurile de constituenți cât și regulile gramaticale. Ele corespund la ceea ce se numește un semn lingvistic, adică un cuvânt, o sintagmă sau o regulă, conținând informații fonetice, sintactice, semantice și discursive. Structurile de trăsături sunt cât se poate de adecvate pentru organizarea într-o notație comună a informațiilor lingvistice eterogene.

Spre deosebire de celelalte teorii lingvistice bazate pe unificare, HPSG utilizează ierarhizarea tipologică. Fiecare structură de trăsături este încadrată într-un anumit tip pentru care sunt predefinite anumite constrângeri și care își are locul într-o ierarhie de tipuri. În cadrul ierarhiei funcționează relația de moștenire a constrângerilor tipurilor superioare asupra descendenților lor. Un exemplu de ierarhie de tipuri este dată în (9).



Pentru fiecare tip sunt definite anumite trăsături specifice (sau anumite constrângeri) care se adaugă constrângerilor moștenite de la tipurile din care descind. Trebuie adăugat că într-o ierarhie de tipuri sunt permise moșteniri multiple, adică sunt permise tipuri care au mai mulți părinți.

Cel mai general tip în HPSG este "semnul" (în engleză *sign*). El conține informație fonologică (prin trăsătura PHON) și informație sintactico-semantică (prin trăsătura SYNSEM). Semnul, la rândul lui, poate fi un cuvânt sau o sintagmă, după cum se vede în (9), mai sus. Sintagma are spre deosebire de cuvânt o trăsătură în plus, numită DAUGHTERS (adică ramuri-surori) care, în plus, are o listă de semne combinate în sintagmă. Un exemplu de semn lexical împreună cu descrierea trăsăturilor specifice acestuia este dată în (10) pentru verbul a vrea.



Combinarea cuvintelor în sintagme se face pe baza unor reguli exprimate la rândul lor sub formă de structuri de trăsături tipizate, purtând numele de scheme de Dominanță Imediată (scheme DI). Asupra regulilor acționează suplimentar principiile, care, la rândul lor, sunt exprimate prin constrângeri asupra anumitor trăsături. În cele ce urmează vom prezenta pe scurt principiile și schemele DI de bază. Dintre principii, ne rezumăm prezentarea la următoarele:

- a. Principiul Trăsăturilor Centrale
- b. Principiul de Subcategorizare

c. Principiul Semantic

a. Principiul Trăsăturilor Centrale

Pentru majoritatea sintagmelor se definește un atribut HEAD ("centru"), inclus în trăsătura CATEGORY (CAT), a cărei valoare trebuie să fie partajată cu cea a atributului HEAD din semnul ramurii-centru HEAD-DTR a sintagmei. Principiul Trăsăturilor Centrale poate fi exprimat prin descrierea următoare (notând valoarea partajată prin indicele [1]):

- (11) "SYNSEM | CAT | HEAD [1]
DAUGHTERS | HEAD-DTR | SYNSEM | CAT | HEAD [1]

Semnul HEAD-DTR poate fi sintagmatic sau lexical,

b. Principiul de Subcategorizare

Atributul SUBCAT are ca valoare o listă care este actualizată progresiv, pe măsură ce sintagma se "saturează", în sensul că atunci când complementele sunt realizate, ele sunt eliminate din lista SUBCAT a sintagmei respective. O sintagmă se numește saturată (sau completă) când valoarea listei SUBCAT este vidă. Principiul de Subcategorizare poate fi enunțat astfel:

Valoarea listei SUBCAT a ramurii HEAD-DTR a unei sintagme trebuie să corespundă concatenării listei L1 ca valoare a atributului SUBCAT al sintagmei și a listei L2 a semnelor ce aparțin ramurii de componente COMPS-DTR (sau, mai precis, nu lista semnelor, ci a trăsăturilor SYNSEM a acestor semne).

Acesta poate fi reprezentat prin structura de trăsături următoare (notând prin simbolul © concatenarea listelor):

- (12) ["SYNSEM | CATEGORY | SUBCAT L1
DAUGHTERS HEAD-DTR | SYNSEM | CAT | SUBCAT L1©L2
COMPS - DTR L2

două scheme DI: "completa" și "subcategorizare" pot fi descrise următoarele

- (13) $cam \wedge L \wedge \wedge \wedge$ s. sintagma saturată cu ramură Comp. emente: **head-**
SYNSEM | CATEGORY | SUBCAT <>
DAUGHTERS HEAD-DTR | SYNSEM | CAT | SUBCAT < X >
COMPS - DTR < X >

, 2. Schema DI pentru o sintagmă non saturată cu ramură Complemente: **head-compl**

- (14) SYNSEMj CATEGORY | SUBCAT < X >
DAUGHTERS HEAD-DTR | SYNSEM | CAT | SUBCAT < X,Y1,Y2...Yn >
COMPS-DTR < Y1,Y2...Yn >

3. Schema DI pentru o sintagmă cu ramură Adjunct: **head-adjunct**

Modificatorii (adjective atributive, adverbe, componente circumstanțiale) sunt introduși într-o ramură specială numită ramura Adjunct (sau ADJCT-DTR). Modificatorii selecționează categoria pe care o modifică (N' pentru adjective, V sau GV pentru adverbe). Această selecție se face printr-un atribut MODIF, care are ca valoare o structură de trăsături SYNSEM. Pentru o sintagmă centru-adjunct bine formată trebuie să aibă loc unificarea valorii trăsăturii MODIF a adjunctului cu valoarea trăsăturii SYNSEM a centrului. Astfel adjectivele pot selecționa numele pentru care sunt atribute, iar adverbele pot selecționa verbele respective, adică se poate preciza în intrarea lor lexicală trăsăturile Categorie, Conținut, Index etc. ale numelui sau verbului așteptat. Descrierea unei sintagme cu Adjunct este următoarea:

- (15) DAUGHTERS HEAD-DTR | SYNSEM | I
ADJCT - DTR | SYNSEM | CAT | HEAD | MODIF 111

c. Principiul Semantic

Principiul semantic reglementează propagarea trăsăturilor semantice, adică cele două trăsături CONTENT și CONTEXT. Se urmărește pe de o parte ca sintagmele să partajeze valoarea trăsăturii CONTENT din ramura centrului cu trăsătura proprie CONTENT, iar pe de altă parte să determine "ridicarea" la nivelul sintagmelor superioare a eventualelor cuantificatori și a variabilelor care le pot corespunde.

HPSG face apel la noțiunea de centru semantic, acesta fiind identic cu centrul sintactic, în afara cazului sintagmelor cu adjunct. În acest caz, centrul Sintactic este categoria modificată, dar centrul semantic este modificatorul (care joacă rolul de predicat semantic). Principiul Semantic poate fi exprimat astfel:

Valoarea atributului CONTENT a categoriei dominante este identică cu valoarea atributului CONTENT a categoriei care este centru semantic (ramura Adjunct sau, implicit, ramura HEAD).

O altă schemă DI, *head-functor*, propusă de Allegranza în [19], reprezintă o modificare a schemei *head-adjunct* cu scopul de a satisface exigențele de

reprezentare a determinantilor într-un grup nominal. Determinatorii sunt tratați ca funcții aplicați centrului. Ei selectează centrul prin atributul ARG-SLOT și marchează sintagma rezultată cu anumite trăsături specifice determinantului respectiv prin partajarea valorii atributului MARKER între ramura Functor și nodul mamă. Descrierea acestei scheme este dată mai jos.

4. Schema DI pentru o sintagmă cu ramură Functor: **head-functor**

SYNSEM | LOCAL | CAT; MARKER III

(16)

FUN-DTR	SYNSEM	LOCAL	CAT	HEAD	MARKER] 1
DAUGHTERS					ARG-SLOT 12 i
	HEAD-DTR	SYNSEM	2		

Cu aparatul formal oferit de HPSG, în secțiunea care urmează, dăm spre exemplificare analizarea unei structuri concrete din limba română. Structura propusă conține un centru nominal modificat de o propoziție relativă al cărei element de relație este în cazul genitiv precedat de articolul genitival. Această structură este interesantă prin faptul că prezintă un fenomen, acela de acord încrucișat, care pare să scape reprezentărilor gramaticilor independente de context. Avantajul teoriei lingvistice discutate aici, însă, oferă o soluție pe cât de unitară, pe atât de elegantă, după cum sperăm să reiasă din cele ce urmează.

3. Structuri relative cu acord încrucișat

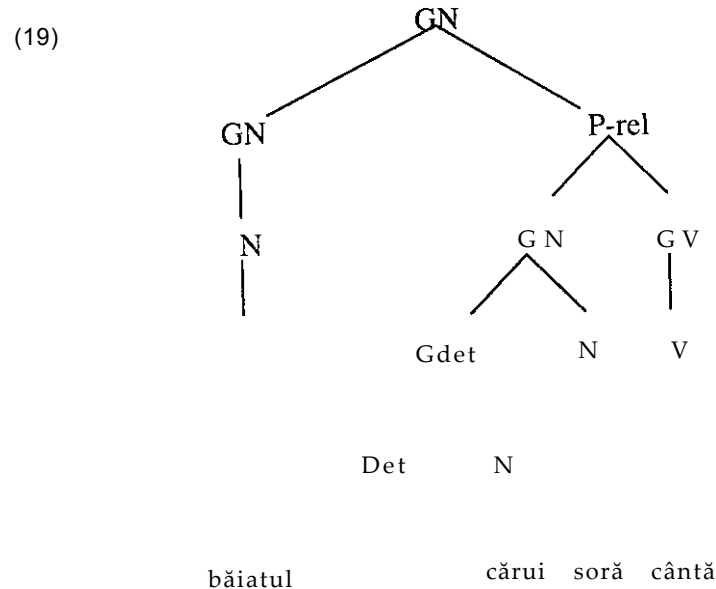
În limba română structurile care prezintă acord încrucișat sunt propozițiile relative în care pronumele relativ este precedat de articolul genitival, ca în exemplul de mai jos.

(17) băiatul a cărui soră cântă

Acordul este încrucișat prin aceea că pronumele relativ propriu-zis se acordă cu substantivul determinat de propoziția relativă, *băiatul*, iar articolul genitival *a*/se acordă cu subiectul relativei, *soră*, după următoarea schemă:

(18) băiatul a cărui soră cântă

Structura internă a acestui grup nominal este reprezentată în arborele de mai jos.

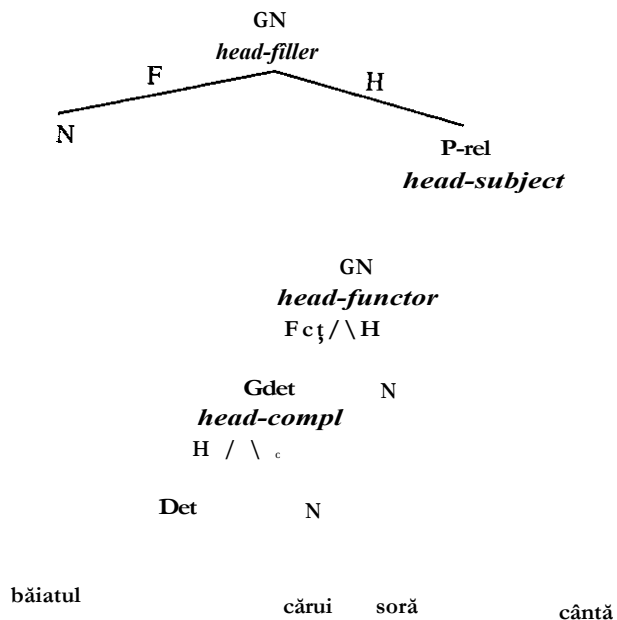


După cum se vede în acest arbore, exemplul din (17) este format dintr-un substantiv centru, *băiatul*, modificat de o propoziție relativă al cărei subiect, a *cărui soră*, cuprinde elementul de relație care face legătura dintre numele amintit și propoziția relativă.

Dacă ne-am limita descrierea la regulile independente de context sugerate în arbore, nu am putea da seamă de fenomenul de acord încrucișat pe care-l discutăm aici. Acest lucru este însă posibil dacă folosim o gramatică HPSG, beneficiind de avantajele oferite de mecanismul unificării și de reprezentările prin structuri de trăsături.

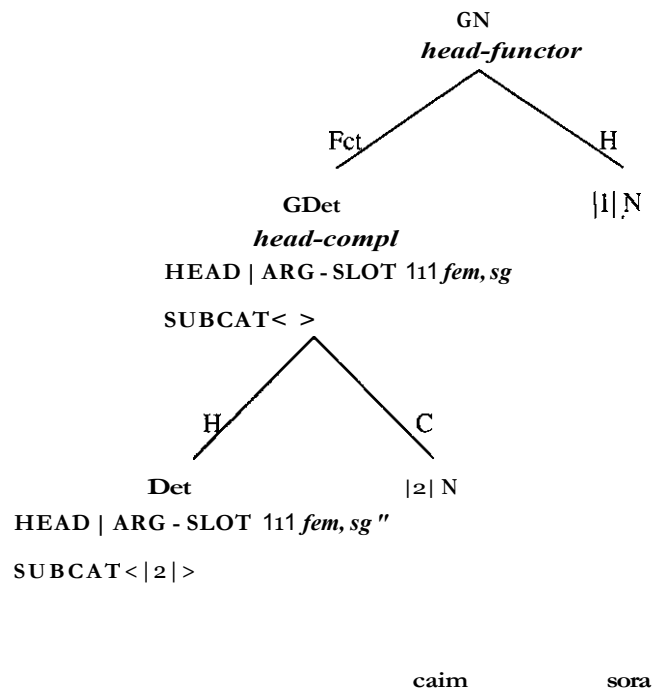
Aplicând schemele de dominanță imediată și principiile specifice teoriei HPSG, arborele de mai sus poate fi adnotat cu regulile HPSG aplicate, în felul următor (unde am folosit ca notații funcționale H=centrul sintagmei, C=complement, Fct=functor, F=filler).

(20)



Fenomenul de acord încrucișat presupune pe de o parte acordul determinantului a cu substantivul *soră*, iar pe de altă parte acordul pronumelui relativ *cărui* cu substantivul *băiatul*. Primul acord amintit se face relativ banal. Intrarea lexicală a determinantului a, în calitate sa de functor, specifică în valoarea atributului său central ARG-SLOT ce trăsături de acord trebuie să aibă substantivul pe care urmează să-l modifice. Când detrminatorul a se combină cu complementul său *cărui*, principiul trăsăturilor centrale face ca această informație să fie percolată la nodul mamă GDet. Mai departe, schema DI head-functor verifică dacă trăsăturile de acord ale GDet unifică cu cele ale centrului său nominal. Acest mecanism este ilustrat în arborele de mai jos.

(21)

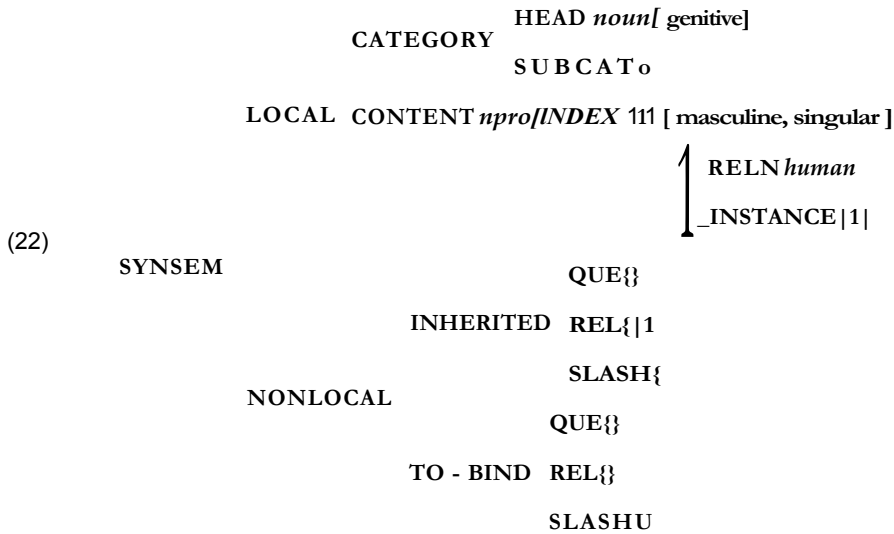


Al doilea tip de acord, în schimb, ridică anumite dificultăți pentru că se realizează într-un arbore local, adică nu se realizează între ramurile unui același nod. Prin urmare, trăsăturile de acord ale pronumelui relativ sunt percolate până la nivelul nodului P-rel (din (19)) pentru a putea fi utilizate în unificarea de regula head-filler cu trăsăturile de acord ale substantivului determinat.

Mecanismul din teoria HPSG care dă seama de propagarea anumitor trăsături se numește mecanismul dependențelor la distanță și este aplicabil fenomenelor de limbă precum interogațiile, topicalizările și, cum este cazul, construcțiile relative. Aici ne vom ocupa numai de tratarea relației de dependență dintre celelalte fenomene a se vedea [5].

Ideea principală a acestui mecanism este că pronumele relative conțin în intrările lor lexicale informații despre numele la care se referă. Într-o intrare Pronumelui relativ din exemplul nostru va conține, prin urmare, informația (22).

PHONOLOGY < *cărui* >



Valoarea trăsăturii NONLOCAL | INHERITED indică acele trăsături care vor fi supuse Principiului Trăsăturilor Nonlocale. Aceste trăsături pot fi specifice elementelor interogative, definite prin atributul QUE, elementelor dislocate, date de atributul SLASH sau pot fi specifice elementelor relative indicate prin atributul REL. După cum se observă în (22), acest ultim atribut are în cazul de față valoare non-vidă, coindexată cu conținutul semantic de masculin-singular al pronumelui.

Potrivit Principiului Trăsăturilor Nonlocale, formulat în (23), valoarea atributului nonlocal INHERITED ("moștenit") este trecută din nod în nod spre vârful arborelui până va întâlni o ramură soră ale cărei trăsături locale unifică cu cele moștenite.

(23) **Principiului Trăsăturilor Nonlocale**

Pentru fiecare trăsătură nonlocală, valoarea atributului INHERITED a nodului mamă este egală cu reuniunea valorilor atributului INHERITED ale ramurilor fiice mai puțin valoarea atributului TO-BIND a ramurii centru.

Atributul TO-BIND, practic, oprește propagarea trăsăturilor moștenite în momentul în care se realizează elementul căutat, adică elementul care a făcut necesară această propagare. De exemplu, trăsăturile de acord ale pronumelui relativ, în exemplul nostru *cărui*, se propagă la nivelul propoziției relative până când este realizat substantivul la care se referă acest pronume, adică *băiatul*.

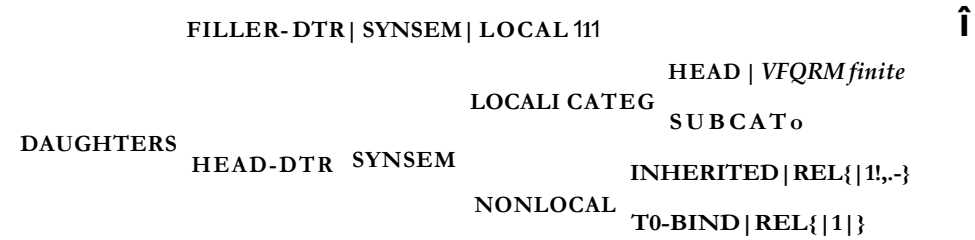
Regula care asignează o valoare atributului TO-BIND în momentul în care are loc unificarea trăsăturilor locale ale unui element cu trăsăturile moștenite pe

emă de dominanță imediată numită *head-filler* (*filler* ar

Tal*.î " d ? e p » -ceTa ce vine sâ completeze o lipsă") si es,e descnsa ,n

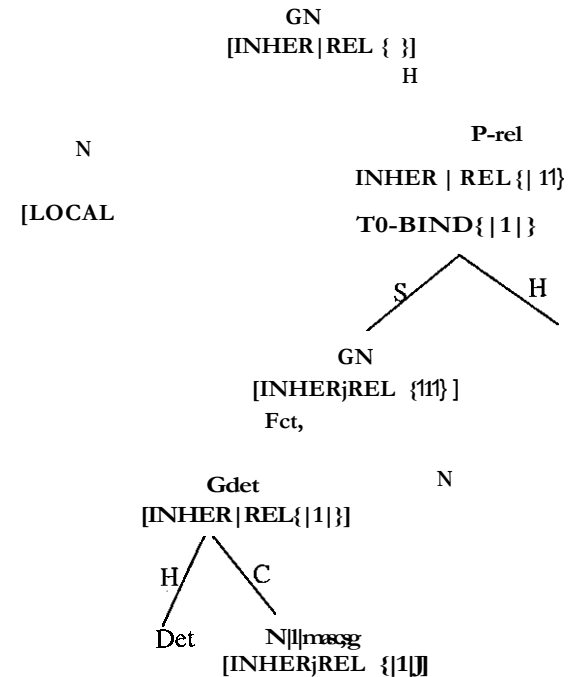
(24).

(24) Schema DI *head-filler*



în trăsăturit dacă aplicăm Principiul Trăsăturilor Nonlocale și schema DI *head-filler* care if-vem în vedere se realizează în man.era ilustrată în arborele de mai jos.

(25)



băiatul a cărui sora cântă

În concluzie, acordul încrucișat avut în vedere presupune, pe de o parte, un acord local, cel dintre articolul genitival și substantivul determinat, în cazul nostru subiectul propoziției relative, iar pe de altă parte un acord la distanță, cel dintre pronumele relativ și substantivul determinat, exterior propoziției relative. Primul tip de acord se face pe baza Principiului Trăsăturilor Centrale și a acordului banal dintre functor și centrul său, pe când cel de al doilea tip de acord face uz de Principiul Trăsăturilor Nonlocale și de schema de Dominanță Imediată *head-filler*.

4. Concluzii

Analiza oferită aici pune în lumină faptul că un fenomen dificil precum acordul încrucișat poate fi tratat într-o manieră relativ simplă și elegantă cu ajutorul unei teorii lingvistice adecvate, cum este teoria Head-driven Phrase Structure Grammar.

Prin aparatul formal și adecvarea lingvistică pe care le oferă această teorie, descrierea fenomenelor limbii române devine incontestabil mai unitară, mai explicită și mult mai riguroasă. Acestor avantaje li se adaugă încă unul, extrem de important, acela al adecvării teoriei pentru implementarea informatică. Este deschis astfel drumul pentru construirea de gramatici computaționale ale limbii române și dezvoltarea componentei informatizate a acesteia.

Aplicațiile informatice ale teoriei HPSG sunt, de altfel, în plină dezvoltare și nu am dori să încheiem înainte de a aminti câteva aspecte în acest sens. ^

Modelul HPSG a făcut parte încă de la origine dintr-un sistem de tratare automată a englezei dezvoltat în laboratoarele de cercetare Hewlett Packard din Palo Alto ([20]). Apoi, au fost propuse diferite implementări, unele bazate pe sistemul PATR ([21]), altele realizate direct în Prolog ([22], [23]). Dintre implementările de sisteme de gestiune a structurilor de trăsături tipologizate și cu moștenire, se poate cita sistemul *Typed Feature Structure* (TFS) al lui M.Emele și R. Zajac [24] și sistemul ALE al lui B. Carpenter [25].

Teoria HPSG a inspirat deopotrivă noul formalism european ALEP, a cărui implementare (în Prolog) presupune un mecanism de gestionare de gramatici și lexicoane, un analizor, un generator și un modul de transfer pentru traduceri automate. Este de altfel utilizat în mai multe centre de cercetare universitară (precum DFKI la Saarbrücken, *Center for Cognitive Science* în statul Ohio, CSLI la Stanford) sau industriale, în special la ATR în Japonia (pentru traducerea automată englezo-japoneză pentru stabilirea de întâlniri prin telefon).

O altă aplicație informatică a acestei teorii, pe cât de recentă, pe atât de importantă este cea cuprinsă în proiectul Verbmobil, [26], care s-a ocupat cu traducerea bidirecțională, în timp real, a textelor vorbite în trei limbi (germană, engleză și japoneză).

Head-driven Phrase Structure Grammar este o teorie care s-a impus incontestabil în lingvistica modernă atât prin numeroasele sale aplicații informatice, cât și prin "generalitatea" aparatului său care o face adecvată pentru numeroase limbi ale lumii, așa cum se poate vedea din impresionanta bibliografie electronică HPSG oferită de pagina www.dfki.de/lt/HPSG. Nu trebuie trecute cu vederea lucrările de limba română dezvoltate în acest cadru, dintre care le amintim pe cele ale lui Ionescu ([27]-[33]), Monachesi ([34]-[36]) și Barbu ([37]) la care s-ar cuveni să se adauge multe altele spre afirmarea limbii române în lingvistica internațională.

Referințe bibliografice

- [1] Pollard, C. - *Generalized Context-Free Grammars, Head Grammars and Natural Language*. Teză de doctorat. Universitatea din Stanford, 1984.
- [2] Kay, Martin - "Funcțional Grammars", *Actes 5° annual meeting of the Berkeley Linguistics Society*, Berkeley, 1979, pp. 142-158.
- [3] Oehrle, Richard; Bach, Emmon; Wheeler, Deirdre (eds.) - "Categorial Grammars and Natural Language Structures", Dordrecht: Reidel, 1988.
- [4] Pollard, C; Sag, I. - *Information-based Syntax and Semantics*, CSLI, University of Chicago Press, 1987.
- [5] Pollard, C; Sag, I. - *Head-driven Phrase Structure Grammar*, CSLI, University of Chicago Press, 1994.
- [6] Sag, I; Pollard, C. - "An integrated theory of complement control", *Language*, 67:1, 1991, pp. 63-113.
- [7] Pollard, C; Sag, I. - "Anaphors in English and the scope of binding theory", *Linguistic Inquiry*, 23:2, 1992, pp. 261-303.
- [8] Pollard, C. - "On head non-movement", *Actele Colocviului Discontinuous constituency*, Tilburg, 1990.
- [9] Nerbonne, J.; Netter, K.; Pollard, C. (eds.) - "German grammar in HPSG", CSLI, University of Chicago Press, 1993.
- [10] Balari, S. - "Feature structures, linguistic information and grammatical theory", Teză de doctorat, Universitatea Autonomă din Barcelona, 1993.
- [11] Gunji, T. - *Japanese Phrase Structure Grammar*, Reidel, 1987.
- [12] Chung, C. - "Korean auxiliary verb constructions without VP modes", *Harvard Workshop on Korean Linguistics*, V; în C. Pollard, I. Sag (eds.), *Readings in HPSG*, 1993.

- [13] Miller, P.; Sag, I. - *French clitic movement without clitics or movement*, LSA Meeting, Los Arigeles, 1993.
- [14] Monachesi, P. - "Object clitics and clitic climbing in Italian HPSG grammar", *Actes 6° European ACL*, Utrecht, 1993, pp. 431-437.
- [15] Robinson, J. - "A machine-oriented logic based on the resolution principle", *Journal of the ACM*, 12, 1965, pp.23-44.
- [16] Colmerauer, A. - "Les grammaires de metamorphose", Universite d'Aix Marseille, 1975, reluat în L. Bolc (ed.) *Natural Language Communication with computers*, Springer, Verlag, 1978.
- [17] Kay, M. - "Funcțional grammars", *Actes 5° annual meeting of the Berkeley Linguistics Society*, Berkeley, 1979, pp. 142-158.
- [18] Abeille, A. - *Les nouvelles syntaxes. Grammaires d'unification et analyse du français*, Armând Colin, Paris, 1993.
- [19] Allegranza, V. - "Determiners as Functors: NP Structure in Italian" în S. Balari & L. Dini (eds.) *Românçe in HPSG*, CSLI, Stanford, 1998.
- [20] Proudian, D.; Pollard, C. - "Parsing Head-driven Phrase Structure Grammar", *Actes 23°ACL*, Chicago, 1985, pp. 167-171.
- [21] Shieber, S. - *An Introduction to unification-based theories of grammar*, CSLI, University of Chicago Press, 1986.
- [22] Oliva, K. - "Simple parser for an HPSG-style grammar implemented in Prolog", *Actes13°COLING*, Helsinki, vol.3,1990, pp.434-436.
- [23] Carpenter, B. - "The generative power of Categorical grammars and Head-driven Phrase Structure grammar with lexical rules", *Computational Linguistics*, 17:3, 1991, pp. 301-314.
- [24] Emele, M.; Zajac, R. - "Typed-unification grammars", *Actes 13° COLING*, Helsinki, vol.3, 1990, pp. 293-298.
- [25] Carpenter, B. - "The Logic of typed Feature Structures with Applications to Unification Grammars, Logic Programs and Constraint Resolution", Cambridge University Press [Implementarea sistemului ALE], 1992.
- [26] Wahlster, W. (ed.) - *Verbobil: Foundations of Speech-to-Speech Translation*, Springer, Berlin, 2000.
- [27] Ionescu, E. - "A Type of SOV Construction in Romanian", *Cahiers de Linguistique Theorique et Appliquee*, tomes XXXII-XXXIII, 1995-1996, 19-39.
- [28] Ionescu, E. - "Accusative Weak Pronouns in Romanian", *Cahiers de Linguistique Theorique et Appliquee*, tomes XXXII-XXXIII, 1995-1996, 40-52.
- [29] Ionescu, E. - "Accusative Clitic Doubling in Romanian", *Cahiers de Linguistique Theorique et Appliquee* tomes XXXII-XXXIII, 1995-1996, 53-73.
- [30] Ionescu, E. - "Accusative Clitic Climbing in Romanian", *Cahiers de Linguistique Theorique et Appliquee*, tomes XXXII-XXXIII, 1995-1996, 74-87.
- [31] Ionescu, E. - "A Quantification-based Approach to Negative Concord in Romanian" in Geert-Jan M. Kruijff and Richard T. Oehrle (editori), *Proceedings of Formal Grammar Conference Utrecht,1999*, p. 25-36.
- [32] Ionescu, E. - *Pro-Drop: An HPSG Account without Lexical Rules*, "Bucharest Working Papers in Linguistics", voi. I, nr.1, 1999, 117-124.
- [33] Ionescu, E. - *On the Status of PE in the Direct Object Construction in Romanian*, Romanian Journal of Information Science and Technology, volume 4, numbers 3-4, 2001, p. 293-310.
- [34] Monachesi, P. - "The morphosyntax of Romanian cliticization" în P.-A. Coppen, H. van Halteren, & L. Teunissen, eds., *Proceeding of Computational Linguistics in The Netherlands 1997*, pp. 99-118, Amsterdam-Atlanta:Rodopi.
- [35] Monachesi, P. - "Linearization properties of the Romanian verbal complex" în *Proceedings of WECOL 98*, Tempe, 1999.
- [36] Monachesi, P. - "Clitic Placement in the Romanian verbal complex", în B. Gerlach and J. Grijzenhout (eds.) *Clitics in Phonology, Morphology and Syntax*, LA 36, Amsterdam: John Benjamins Publishing Company, 2000.
- [37] Barbu, A.M. - "Romanian determiners:order and classification" în *Revue Roumaine de Linguistique*, XLIII, nr.5-6, pp.299-315, București, 1998.
- [38] Uszkoreit, H. - "From Feature Bundles to Abstract Data Types: New Directions in the Representation of Linguistic Knowledge, in H. Blaser *Natural Language at the Computer*, Berlin: Springer, 1989.

După 10 ani de experiență terminografică: noul model de date terminologice al TermRom

Dan MATEI, Institutul de Memorie Culturală
Piața Presei Libere, nr. 1, C.P. 33-90, 713411, București
dan@cimec.ro

A. Preambul

Din 1991 – când a fost înființată — Asociația Română de Terminologie (TermRom) a desfășurat o activitate terminografică materializată într-o bază de date proprie (accesibilă, în parte, pe web la www.cimec.ro/tr/) și într-o serie de publicații specifice. Formatul terminografic utilizat — descris în [1] –, derivat din formatul standard MicroMATER (ISO 6156), se bazează pe un model de date (relativ) complex, serializat pe două nivele: nivelul conceptului și nivelul termenului. Practica terminografică (ce se traduce prin prelucrarea unei mari diversități de date terminologice) ne-a revelat o tensiune între complexitatea datelor reale și insuficiența complexitate a modelului folosit. În plus, necesitatea transferului de date între aplicații diverse a scos la iveală utilitatea consemnării cu o granularitate sporită a elementelor înregistrării terminologice. Mai mult, "entuziasmul" cu care ISO revizuește standardele terminologice în ultimii ani¹, cu alte cuvinte, relativa instabilitate a standardelor din acest domeniu, îndeamnă la o și mai fină granularitate, pentru a spori șansele de compatibilitate cu normele de transfer viitoare. Pe de altă parte, pe măsura acumulării experienței, era din ce în ce mai limpede că modelul de date folosit ar trebui să acomodeze o mai mare diversitate și complexitate de metadata bibliografice, ca și o fină și flexibilă tratare a metadatelor "administrative", de gestionare a colecției terminologice (vezi și [2]).

Aceste considerente au dus la elaborarea unui model de date obiectual, care, pe lângă cerințele expuse mai sus, să fie și suficient de abstract ca să permită o serializare convenabilă (pentru transfer de date), — probabil bazată pe XML, de exemplu în formatul MARTIF [3] — și să nu ceară elaborarea de aplicații informatice de o complexitate excesivă.

¹ *Atât ISO 12200 cât și ISO 12620 sunt în revizie (deși ambele datează doar din 1999), iar ISO 16642, este încă nedefinitivat. Desigur, această stare a lucrurilor probează și faptul că domeniul nu este încă bine "așezat".*

B. Modelul

Clasă

Generalizare

Dependență

Asociere

Convențional, modelul este împărțit în secțiuni ("pachete" [packages], în terminologia UML). La nivelul cel mai de sus, se disting secțiunea (așa zis) funcțională și secțiunea administrativă.

B.1. Secțiunea funcțională

În fig. 1 se prezintă clasele funcționale esențiale și asocierile lor. Practic, orice element al modelului este o 'înregistrare'. Cu alte cuvinte, 'înregistrare' este clasa generică. Existența unei clase generice oferă — pe lângă gruparea proprietăților comune tuturor elementelor — și posibilitatea de a avea un identificator unic pentru fiecare înregistrare din baza de date ce implementează acest model.

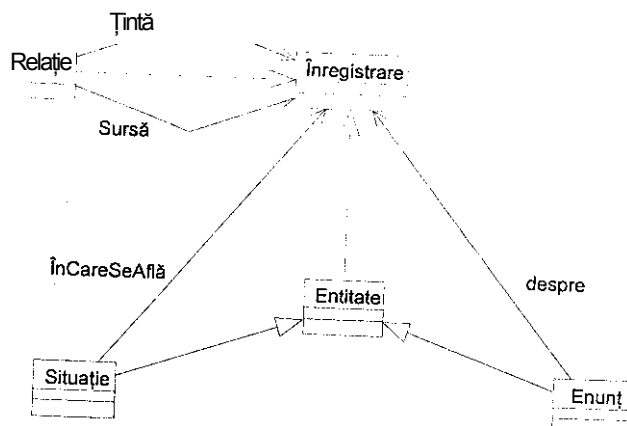


Figura 1 - Secțiunea funcțională (generică)

Clasa 'înregistrare' are două subclase: 'entitate' (care grupează ce au o existență autonomă) și 'relație' (care grupează asocierile înregistrări). Se observă că sunt acceptabile chiar și relațiile binare lucru folositor și în practică.

Reificarea relațiilor binare între înregistrări simplifică modelul și constituie o manieră flexibilă de a consemna o mare varietate de elemente ale modelului. O relație R poate avea două caracteristici utile în cadrul modelului:

- simetria: dacă x este în relația R cu y, y este în relația R cu x.
- tranzitivitatea: dacă x este în relația R cu y și y este în relația R cu z, atunci x este în relația R cu z.

Pentru fiecare instanță a clasei 'relație', aceste caracteristici se consemnează ca un atribut al tipului respectiv de relație (nereprezentat în model)². Consemnarea acestor proprietăți ale relațiilor poate fi foarte utilă pentru programele care ar exploata baza de date.

Pentru a se rezolva (relativ) simplu și flexibil asocierile înregistrări, s-a introdus subclasa 'situație' a clasei 'entitate'. După o figură, o instanță (sau mai multe) a clasei 'situație' se asociază cu o instanță a clasei 'înregistrare', iar obiectul 'situație' este conectat cu oricâte alte instanțe banale ale clasei 'relație'. În practică, cele mai frecvente obiecte de acest tip sunt ca încarnări de contexte și evenimente. În plus, o altă subclasă a clasei 'entitate' este 'enunț'. Acest tip de obiect consemnează atribute ale unei înregistrări care n-au fost aprioric prevăzute, cu alte cuvinte el găzduiește mențiuni pentru care se dorește un set de simple note, și anume care se doresc a fi colocabile și/sau indexabile.

În continuare se prezintă doar subsecțiunile secțiunii funcționale de interes în contextul acestui volum.

B.1.1. Secțiunea terminologică

Fig. 2 prezintă entitățile (i.e. subclasele clasei 'entitate') din secțiunea terminologică.

O categorie de relații importantă în terminologie — este cea a relațiilor tranzitive și asimetriche.

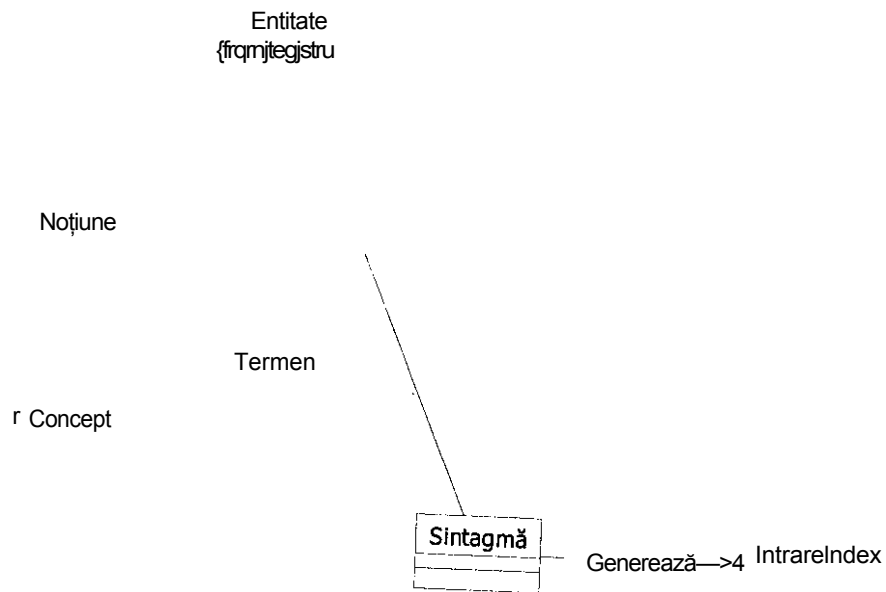


Figura 2 - Secțiunea terminologică

Principala clasă a acestei subsecțiuni este 'noțiune'. Instanțele ei consemnează noțiunile vehiculate în baza de date terminologică, independent de limbă. Din rațiuni practice, și anume din necesitatea de a cuprinde în baze de date terminologice și materialul organizat de obicei în tezaure terminologice, s-a decis să se cuprindă în modelul de date nu doar conceptele pure, ci și unități semantice mai largi, precum cele desemnate de termenii compuși într-un tezaur (sau ceea ce ISO 12620 numește 'unități frazeologice' [A.2.1.18]). Clasa acestor unități conceptuale care cuprinde conceptele și unitățile semantice mai largi este clasa 'noțiune'. Distincția fină între 'noțiune' și 'concept' este formulată în logică astfel [5]:

Noțiune: formă logică fundamentală care reflectă însușirile caracteristice necesare și generale ale unei clase de obiecte.

Concept: noțiune care reflectă însușirile esențiale ale unei clase de obiecte.

Asadar, o noțiune care nu e concept cuprinde mai mulți factori semantici, deci poate fi factorizată.

A doua subclasă a acestei secțiuni este 'termen'. Instanțează doar "denumirile" conceptelor (A.1. în ISO 12620). Consemnează ceea ce au în comun o familie de expresii lingvistice pentru un concept¹. Expresiile lingvistice propriu-zise sunt consemnate în "sintagmă"². Din pricina faptului că un termen poate fi exprimat prin expresii lingvistice (flexiuni, variante ortografice etc), s-a precizat în "termenului" de expresiile sale lingvistice, în felul acesta nu ne referim la definiția pentru 'termen', din ISO 12620 (A.1): "a designation of a concept in a special language by a linguistic expression".

Se poate observa în figură faptul că sintagmele generează intrări de index. În fapt, o sintagmă poate genera — prin inversare/permutare — mai multe intrări de index, dacă terminograful decide că asta ar fi în folosul utilizatorilor. Sintagma generează intrări de index la fiecare "factor" semnificativ. Exemple:

Sintagma	Intrări de index
efect Doppler	efect Doppler Doppler, efect
pseudofonetism	pseudofonetism fonetism, pseudo-
completivă indirectă anticipată	completivă indirectă anticipată indirectă anticipată, completivă indirectă anticipată, completivă indirectă

Clasa 'relație' este vitală pentru consemnarea asocierii conceptelor în modelului. Pentru a ilustra modul în care se consemnează informații de relație esențială, în fig. 3 s-au reprezentat tipurile de relații esențiale care pot exista între concepte, o parte, conceptele cu termenii care le desemnează, iar pe de altă parte, sintagmele care-i exprimă. De asemenea, se vede cum o "situație" poate caracteriza această ilustrare — implică (cel puțin) un loc, o perioadă de timp care caracterizează designarea.

O regulă simplă, pragmatică de a distinge o noțiune care este concept de o noțiune care este concept și-ar găsi locul într-un dicționar terminologic, nu.

Exemple de "familie de expresii lingvistice" sunt: a) clădire, clădiri; b) expresiv, expresive.

În acest context, 'sintagmă' desemnează — printr-un abuz de limbaj — o expresie de cuvinte.

¹ d e n * , M r e n u s u n i c o n c e p t e ; b l o n a * i m w t o t _

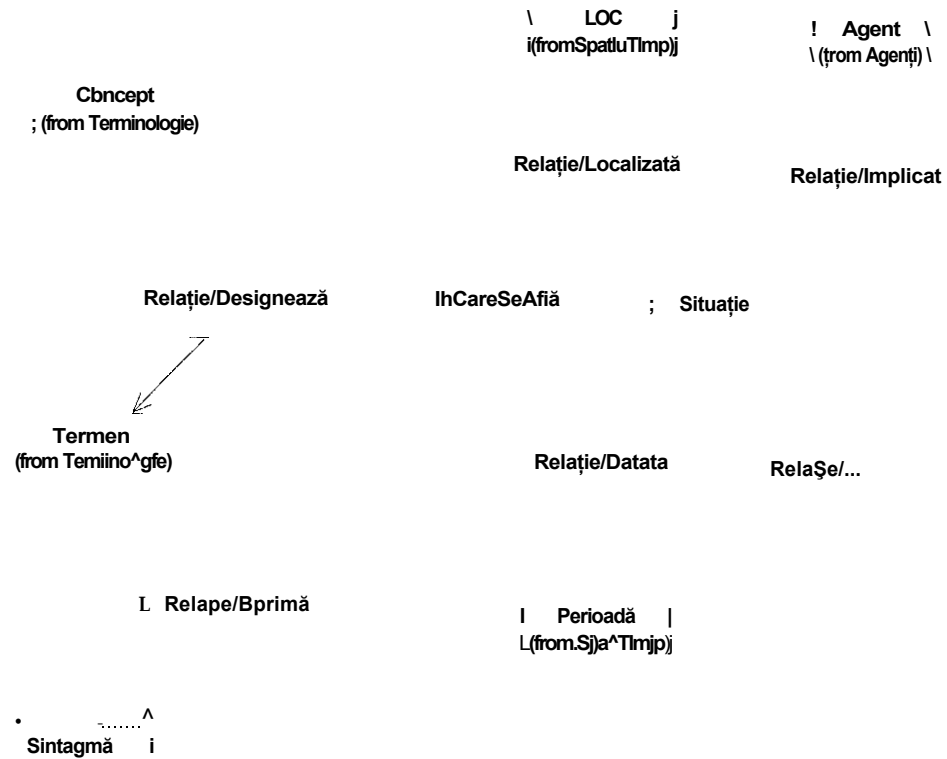


Figura 3 - Ilustrare a reprezentării informației terminologice

Într-o astfel de schemă, se pot reprezenta cu acuratețe cazuri precum:

a) Concept: *mic arbust cu flori roșietice din familia ericaceae ...*

- **Relație/designează:**
 Termen (științific) [latină]:
 Relație/exprimă:
 Sintagmă: *Kalmia latifolia*
- **Relație/designează:**
 Situație/context:
 Relație/localizează:
 Loc: *nordul Statelor Unite*
 Termen [engleză]:

- **Relație/exprimă:**
 Sintagmă: *mountain laurel*
- **Relație/designează:**
 Situație/context:
 Relație/localizează:
 Loc: *sudul Statelor Unite*
 Termen [engleză]:
 Relație/exprimă:
 Sintagmă: *calico bush*
- **Relație/designează:**
 Situație/context:
 Relație/localizează:
 Loc: *sudul Statelor Unite*
 Termen [engleză]:
 Relație/exprimă:
 Sintagmă: *sheep's bane*
- **Relație/designează:**
 Termen [română]:
 Relație/exprimă:
 Sintagmă [s.m.sg.]: *laur de munte*
 Relație/exprimă:
 Sintagmă [s.m.pl.]: *lauri de munte*
- b) Concept: comandant de călărimie
 - **Relație/designează:**
 Situație/context:
 Relație/localizează:
 Loc: *Moldova*
 Relație/localizează:
 Loc: *Țara Românească*
 Relație/datează:
 Perioadă: *sec. XVII-XVIII*
 Termen [română]:
 Relație/exprimă:
 Sintagmă [s.m.sg.]: *serdar*
 Relație/exprimă:
 Sintagmă [s.m.pl.]: *serdari*
- c) Concept: boier de rang mijlociu

- Relație/definează
 Situație/context:
 Relație/datează:
 Perioadă: *sec. XVIII-XIX*
 Termen [română]:
 Relație/exprimă:
 Sintagmă: *serdar* [s.m.sg.]
 Relație/exprimă:
 Sintagmă: *serdari* [s.m.pl.]

- Situație/etimologie:
 Relație/provine din:
 Termen [grecă]:
 Relație/exprimă:
 Sintagmă:
 Relație/provine din:
 Termen [latină]:
 Relație/exprimă:
 Sintagmă:

Tot ca o ilustrare, în fig. 4 se prezintă modul cum se consemnează etimologia unui termen, cu ajutorul clasei 'situație': o situație de tip 'etimologie' se asociază cu termenul de bază, iar termenii din care acesta provine sunt asociați cu situația prin intermediul unor relații de tip 'provineDin'.

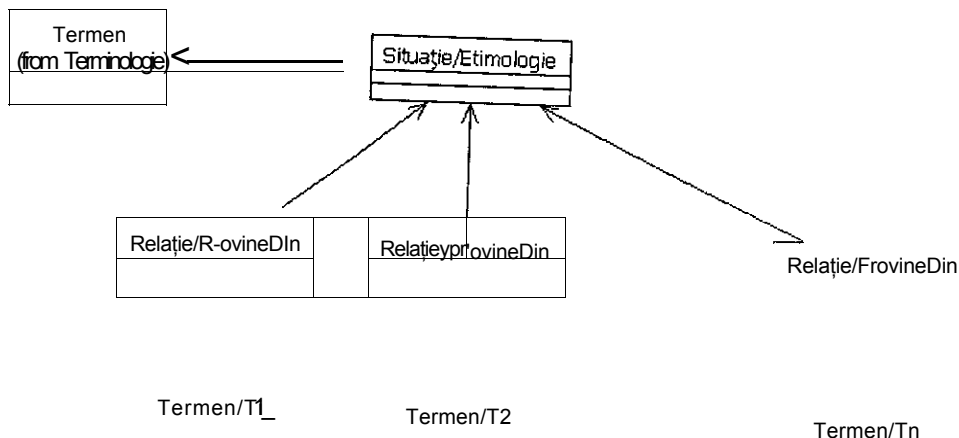


Figura 4 - Ilustrare a reprezentării etimologiei

De pildă:

- Concept: fixat la vârf
- Relație/definează:
 Termen [română]:
 Relație/exprimă:
 Sintagmă: *acrofix*

6.7.2. Secțiunea bibliografică
 Fig. 5 prezintă entitățile (i.e. subclasele clasei 'entitate bibliografică', cu alte cuvinte este o secțiune de metadate. Secțiunea deoarece o bună parte din multitudinea de date bibliografice sunt ajutate de relații. Clasa esențială este 'ediție'; cea care conține informații bibliografice a unei ediții citate.

Entitatea 'lucrare' consemnează metadatele specifice unei lucrări textuale, în cazul nostru), i.e. "abstractizează" ceea ce au în comun toate lucrările. Utilitatea ei imediată este colocarea tuturor manifestărilor indiferent de limbă sau ediție. O subclasă importantă a clasei 'lucrare' este 'serial'. Aici se consemnează și periodicele, adică entitățile ce aparțin clasei 'NumărPeriodic', cu alte cuvinte publicațiile-gazdă ale articolelor. Asupra acestor clase și a relațiilor între ele depășește cadrul acestui capitol.

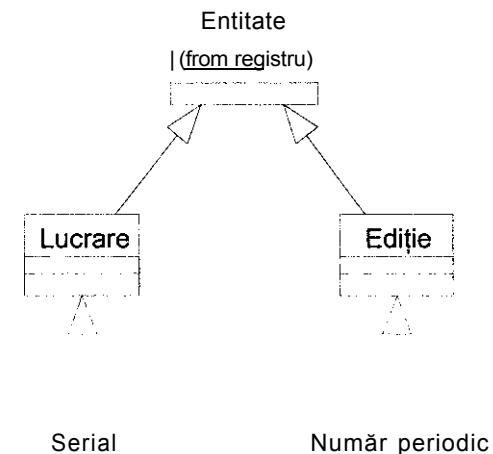


Figura 5 - Secțiunea bibliografică

B.2. Secțiunea administrativă

În fig. 6 se prezintă clasele de natură administrativă și relațiile esențiale între ele. Rolul acestor clase este de a consemna modificările survenite în baza de date, în succesiunea lor. În acest fel se poate urmări geneza înregistrărilor și se pot identifica responsabilitățile. În plus, deoarece se prevede și stocarea datelor modificate, se creează premisele revenirii la stări anterioare ale bazei de date. În instanțele clasei 'intervenție' se consemnează fiecare modificare operată asupra unei înregistrări. Fiecare asemenea instanță este asociată — prin intermediul instanțelor clasei 'contribuție' — cu agentul (i.e. operatorul) care a produs-o. În plus o intervenție este asociată și cu sursele ei documentare. Se observă cum clasa 'referință' poate avea ca instanțe atât referințe bibliografice (citând o ediție), cât și referințe personale (citând o comunicare personală).

Clasa 'înregistrareArhivă' este foarte importantă, instanțele ei fiind chiar versiunile "desuete" (i.e. cele dinainte de modificări) ale atributelor înregistrărilor.

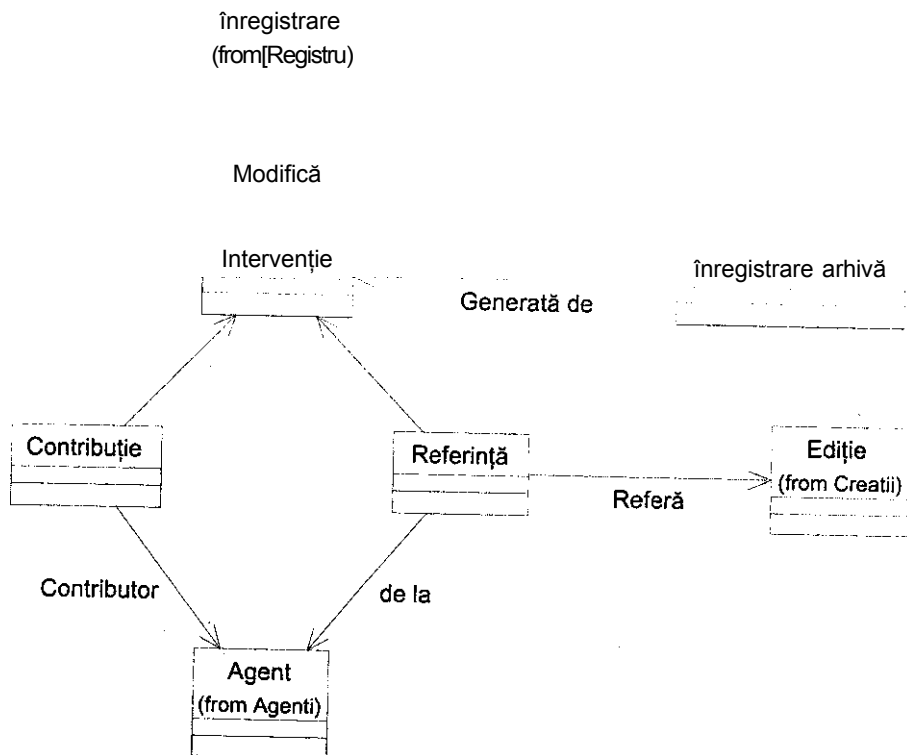


Figura 6 - Secțiunea administrativă

C. Remarci finale

Modelul prezentat pare suficient de flexibil pentru a satisface cerințele funcționale atât ale unei baze de date terminologice, cât și a uneia lexicografice (mai ales datorită distincției între termeni și expresiile lor lingvistice). El este și suficient de abstract pentru ca schema unei baze de date ce l-ar folosi ca fundament să fie relativ comodă la implementare.

TermRom are în curs un proiect de elaborare a unei astfel de baze de date terminologice. După finalizarea acesteia, este de așteptat un proces traumatic de convertire a bazei de date curente. Sporul de funcționalitate obținut va compensa însă efortul.

D. Referințe

- [1] Matei, Dan. *Banca de date terminologice a TermRom și problemele ei neologice*, în *Limbaj și Tehnologie* / Dan Tufiș - editor. - București: Editura Academiei Române, 1996'
- [2] ISO/CD 16642:1999, *Computer applications in terminology - Metamodel for representing terminological data collections*
- [3] ISO 12200:1999, *Computer applications in terminology - Machine-readable terminology interchange format (MARTIF) - Negotiate interchange*
- [4] ISO 12620:1999, *Computer applications in terminology - Data categories*
- [5] Chețan, Octavian, Radu Sommer. *Dicționar de filozofie / Coordonare științifică Octavian Chețan, Radu Sommer.* — București: Editura Politică, 1978

Probleme de reprezentare a date terminografice într-o bază de date rel

Sorin GHEȚARII

TERMOROM, Str. Meșterul Manole nr. 3

gsorin@fx.ro

Oriunde și oricând se creează, comunică, înregistrează, stochează, transformă sau refolosește informație sau cunoștințe este implicată într-un fel sau altul și terminologia. Comunicarea în domeniu a devenit un discurs specializat cu texte de specialitate în nenumărate forme. Atunci când se definește terminologia ca o mulțime de concepte și denumirile lor într-un anumit domeniu, ea poate fi considerată infrastructura cunoașterii de specialitate. Scrierea textelor tehnice de tehnică devin astfel imposibile fără o utilizare corectă a unor resurse. Deoarece producerea textelor tehnice implică frecvent multă terminologie multilingve de înaltă calitate au devenit bunuri mult căutate găsit pe înfloritoare piață a industriilor limbajelor și cunoașterii.

Există numeroase baze de date terminologice disponibile on-line sau pe CD-ROM (TERMIUM, EURODICAUTOM), pe disc sau prin intermediul unor dicționare electronice sau ca baze de date personale realizate de ingineri, specialiști în calculatoare, chimiști care lucrează în traducătorii, autori de texte tehnice. Aceste baze de date sunt utilizate în:

- traducere asistată de calculator;
- scrierea de texte tehnice și științifice asistată de calculatoare;
- sisteme informatice (administrarea componentelor etc.);
- cercetări terminologice în lingvistică, filozofia științei și tehnologiei etc.

Pentru asemenea obiective au fost dezvoltate aplicații (programe de management al bazelor de date terminologice), unele pe piața terminologică internațională, altele ca prototipuri în cadrul cercetare academice.

MARTIF este formatul standardizat pentru managementul terminologic. Posibilitatea organizării terminologiei în baze de date diferite face nerealistă presupunerea ca s-ar putea cădea de acord asupra unui anumit format de bază de date relațională, așa cum este SQL, pentru schimburile terminologice. De aceea s-a mers pe linia

format la dispoziția publică fără obligații materiale și care să fie independent de platforma de lucru. Rezultatul este MARTIF (Machine-Readable Terminology Interchange Format cunoscut și ca ISO 12200. In ISO 12620 sunt descrise 150 de categorii de date, un număr imens care nu urmărește decât să le arate pe cele posibile și modul în care acestea pot fi structurate. Categoriile MARTIF sunt împărțite în 10 secțiuni grupate în 4 clase. Acestea sunt:

- termen: cuprinde categoria de date termen (1);
- informație în legătură cu termenii: conține informația legată de termeni (2) și informația privind gradul de echivalență;
- informație descriptivă: relație cu domeniul (4), descrierea conceptului (5), relații între concepte (6), categorii de date care leagă un concept de poziția sa în sistemul de concepte (7), note (8);
- informație administrativă: categorii de date care leagă un concept de un element al unui tezaur sau de o altă formă de documentare (9), categorii de date care cuprind informații administrative.

Un avantaj major al faptului că MARTIF este scris folosind cod SGML este acela că, deși se poate aprecia că lectura codului nu este facilă, ea este totuși posibilă ca urmare a faptului că nu face apel decât la caracterele ASCII. Un alt avantaj al sistemului MARTIF este acela că el acceptă referințe către alte documente chiar din interiorul documentului. Inițial MARTIF presupune că înainte de implementarea produselor software pentru importul sau exportul datelor programatorii sunt obligați să examineze sursele implicate. Pentru a asigura un acces așa numit "orb" care să permită oricui să transfere baze de date terminologice din orice sistem spre sau dinspre MARTIF este necesară o standardizare suplimentară a categoriilor de date, domeniilor specifice etc.

Tabela ce urmează enumera cea parte a "elementelor" MARTIF care sunt de cea mai mare importanță pentru realizarea unei resurse terminologice Multilingve.

<termEntry>	<p>ivnrimTf ⁶ U m C o e o . atetermin o'ogice _{pentru un concept} "</p> <p>S ^ n S S S ^ ? ^ ? administrative codate lor sau, în cazul unei abordări bilingve sau multilingve, două sau mai multe date descriptive și administrative asociate lor</p> <p>Atributele includ:</p> <p>type, care clasifică setul de date terminologice conform categoriile de date specificate de ISO 12620</p>
<langSet>	<p>Limba, în caarul unui element <termEntry> va fi folosit pentru a grupa mai multe <tig> și <ntig> asociate unei singure limbi</p> <p>S S S S T * . an 9 8546 Ob.igatorie! ? n afara - u ' S care</p>

HI **<tig>**

Grup de informații terminologice; în cadrul unui element <termEntry>, va conține elemente de informații asociate cu un singur termen, fiecare dintre acestea funcționând la rând cu alte cuvinte nu este permisă imbricarea între elemente subordonate unui <tig>.

Prezența atributului lang este obligatorie, în afara cazului în care el este moștenit.

<ntig>

Grup încuibat de informații terminologice; va fi folosit în cadrul unui element <termEntry> dacă anumite elemente interne sunt asociate mai curând cu elemente interne, decât cu elementele <tig>.

<term>

Următoarele elemente vor fi folosite în cadrul <ntig> pentru a găzdui alte date terminologice: <termGrp>, <termNote>, <descripGrp> și <adminGrp>.

Prezența atributului lang este obligatorie, în afara cazului în care el este moștenit.

[<termGrp>

Va conține un termen format dintr-un singur cuvânt sau din multe cuvinte, sau o desemnare simbolică privită ca termen.

"Va conține un element <term> și posibil, cel puțin un element încuibat în plus față de termen."

<termNote>

"Va conține informații legate de termen.

Atributele includ:

type care clasifică <termNote> conform categoriilor de date terminologice; în ISO 12200.

[<termNoteGrp>

Va conține un element <termNote> și posibil cel puțin un element încuibat în plus față de informația legată de termen pentru a găzdui un nivel suplimentar de imbricare în cadrul elementului <termGrp>

<descrip>

Va conține informații descriptive precum definiția, explicații descriind concepte și termeni.

Atributele includ:

type, care clasifică <descrip> potrivit categoriilor de date terminologice în ISO 12200.

<descripGrp>

[<admin>

"Va conține date administrative.

Atributele includ:

type care clasifică <admin> în funcție de categoriile de date terminologice în ISO 12200.

<adminGrp>

<date> Va conține o singură dată de formatul YYYY-MM-DD, cu opțiunea notării dată-timp YYYY-MM-DDhh:mm:ss. Atributele includ:
type, care clasifică <date> după categoriile specificate în ISO 12200.

<note> Va conține o notă sau o adnotare drept comentariu legat fie de un întreg <termEntry>, un întreg <tig> sau <ntig> ori de unul din elementele <...Grp>.

<descripNote> Va fi folosit în cazul informațiilor de tipul <note> folosite în cadrul <descripGrp> când conținutul notei este legat de o listă de opțiuni.

<adminNote> Va fi folosit în cazul informațiilor de tipul <note> folosite în cadrul <adminGrp> când conținutul notei este legat de o listă de opțiuni

<ptr>¹ Va consta dintr-un indicator către o altă locație din documentul curent.
 (Atributele includ:
 type, care clasifică <ptr> conform Anexei A, A. 12
 target, care precizează destinația referirii, ca unul sau mai mulți identificatori SGML.

<ref> Va defini o referire către o altă locație din documentul curent, în termeni de unul sau mai multe elemente identificabile. <ref>GI este asociat cu text suplimentar drept conținut al elementului, deci constă dintr-o etichetă-start cu o țintă integrată, urmată de textul asociat și închisă de o etichetă-sfârșit.
 Atributele includ:
 type, care clasifică <ref> conform Anexei A.
 target, care precizează destinația referirii ca unul sau mai mulți identificatori SGML.

<xref>² Va defini o referință la un grafic, ilustrație, figură, tabel sau alt document extern sau fișier folosind o notație indicativă extinsă ca valoare a atributului țintă a <xref>, de ex. <xref target='documentIdentifier'>, unde valoarea 'documentIdentifier' este un cod de identificare pentru documentul țintă. Utilizatorul va documenta notația indicativă extinsă care este folosită incluzând un comentariu adecvat în elementul <encodingDesc> ale header DTD.
 Atributele includ:
 type, care clasifică <xref> conform Anexei A.
 target, care precizează destinația referirii ca unul sau mai mulți identificatori SGML

¹ Mmcăt ZTJdoar "To I t S i Z IT <xref> sunt toate c o n ^

or drept co n t i n u t * amantului, into p t a – Elementele <ptr>, <ref> și

* Z ~ u r ^ de <TMf> **>»> să fi accesibile Sllui-țintă pentru

>> 3

SUI

<foreign>

<refObjectList>

<refObject>⁴

Va fi folosit pentru a marca un cuvânt sau o frază ca grafic în contrast cu textul înconjurător.
 Atributele includ:
 type, care clasifică <ref> conform Anexei A.
 target, care precizează destinația referirii ca unul sau mai mulți identificatori SGML **sau o fraza ca apărând al**
 cea a textului înconjurător.
 Atributele includ:
 lang, care identifică limba cuvântului sau frazei marcate
 lang, unit? monioa: m m / M, - * «7r7irn il o o n m

Va fi folosit în back-matter și va conține unul sau mai multe obiecte back-matter, mai ales resurse comune ca: bibliografice, date de responsabilitate, identificatori jnamespace (URL-uri și FPI-uri), material textual la referiri dese, liste de locații geografice, fișiere externe, asemenea.
 (Atributele includ:
 itype, care clasifică <refObjectList> după categoriile specificate în ISO 12620 Anexa A, A. 11.4.1
 Va conține o data constând în general dintr-o resursă ca: date bibliografice, date de responsabilitate, iden inamespace (URL-uri și FPI-uri), material textual la referiri dese, liste de locații geografice, fișiere externe, asemenea. Datele bibliografice ar trebui să rezide sau într-un document extern (caz în care se va face datele bibliografice din back matter folosind elem
 Atributele includ:
 type, care clasifică <refObject> după categoriile d specificate în ISO 12620 Anexa A, A. 11.4.2. Dacă altfel, tipul <refObject> este moștenit de la <refOB [respectiv.

Notă - în managementul terminologiei o utilizare frecventă a <hi> se face în termeni necesari, adică termeni folosiți într-o definiție, notă sau alt material sunt definiți altundeva în resursa terminologică. Vezi de asemenea Anexa A.11.4.1.
Notă - Unele documente terminologice cuprind date bibliografice ca referință la surse externe. Această practică încurajează redundanța și efortul mărit pentru îngrijirea informațiilor ar trebui convertite în obiecte back matter (informații bibliografice) însoțite de un back-matter posibil.

<itemSet>	Va fi folosit în back matter și va conține unul sau mai multe obiecte individuale care în mod tradițional sunt grupate împreună, de ex. obiectele numele autorului și prenumele autorului vor fi grupate împreună într-un <itemSet> de tip=autor Atributele includ: type care clasifică <itemSet> în principal conform categoriilor de date listate în ISO 12620 Anexa B. Totuși acest Standard Internațional nu specifică întregul spectru al categoriilor de date care pot fi folosite cu <itemSet>
<item>	Va conține un exemplu individual de informație back matter Atributele includ: type, care clasifică <itemSet> în principal conform categoriilor de date listate în ISO 12620 Anexa B pentru informații bibliografice Totuși acest Standard Internațional nu specifică întregul spectru al categoriilor de date care pot fi folosite cu <item>
<itemGrp>	Va conține unul sau mai multe <item> împreună cu <otr> <ref> sau <note>. Atributele includ: type, care clasifică <item> în principal conform categoriilor de date listate în ISO 12620 Anexa B pentru informații bibliografice Totuși acest Standard Internațional nu specifică întregul spectru al categoriilor de date care pot fi folosite cu <itemSet>

Din acest tabel au mai fost eliminate elementele (aproape la fel de numeroase) specifice informațiilor bibliografice. Instanțierea elementelor enumerat mai sus se face prin intermediul "categoriilor de date" standardizate de ISO 12620. Numărul acestora este de aproximativ 200. În cea mai amplă resursă terminologică (EURODICAUTOM) sunt în prezent prezente mai puțin de 20 astfel de categorii de date.

Uniunea Europeană în activitatea sa este unul dintre utilizatorii majori ai procedurilor de traducere a textelor și terminologiei. Aceasta se datorește parțial faptului că legislația sa este direct aplicabilă în statele membre și de aceea ea trebuie să fie disponibilă în toate limbile de lucru oficiale. Ca rezultat, traducătorii Comisiei Europene produc mai mult de 1 milion de pagini pe an și au de-a face cu cel puțin 6-7 milioane de termeni (în medie sunt 8 sau 9 termeni care ridică probleme pe fiecare pagină).

Unitatea pentru Terminologie a Comisiei Europene este destinată asigurării suportului lingvistic pentru toate limbile oficiale ale Uniunii Europene. Au fost elaborate glosare de specialitate, multe dintre ele în nouă limbi. Domeniile acoperite sunt tratatele importante cum ar fi cele de la Maastricht și Roma, cele economice și administrative (Taxa pe Valoarea Adăugată, buget) dar și unele legate de subiectele centrale sau puternic inovatoare ale științei și tehnologiei (fizica plasmei, biotehnologie, minerit). Deosebit de rolul lor de resurse

terminologice și de surse terminologice pentru domeniile de inovare, acest glosare documentează ceea ce se numește "Eurolect", adică frazele și cuvintele care își au origina în cadrul Uniunii Europene și pentru care nu există echivalențe naționale.

Monitorizând toate modificările apărute ca urmare a unei evoluții permanente a bazei de date EURODICAUTOM am constatat că, recent, a avut loc schimbarea suportului hardware și odată cu aceasta pot fi observate următoarele:

- Indicarea mult mai frecventă a referinței la documentul sursă al termenului;
- Indicarea frecventă a referinței la documentul sursă al definiției acestuia;
- Indicarea documentului sursă și pentru sinonime și abrevieri;
- Utilizarea mai frecventă a notelor pentru adăugarea unor informații suplimentare asupra termenilor, acestea putând fi grupate astfel:
 - o {NTE} explicații și informații generale asupra termenilor;
 - o {TXT} contextul (de cele mai multe ori un exemplu de utilizare a termenului respectiv);
 - o {GRM} informații gramaticale (gen, număr);
 - o {USG} indicarea mediului în care este utilizat termenul: "*technical jargon*";
 - o {REG} notă asupra unor utilizări locale speciale sau asupra regionalismelor;
 - o {DOM} indicarea unui domeniu sau subdomeniu care completează clasificarea obișnuită folosită anterior și care a rămas încă prezentă.

De asemenea se prevede ca în cel mai scurt timp să fie implementate următoarele:

- afișarea tuturor caracterelor și diacriticelor (ca și a informațiilor lingvistice, dacă se cere);
- îmbunătățirea sistemului de clasificare a domeniilor;
- introducerea link-urilor interne și externe.

Modelele de date terminologice orientate în exclusivitate către terminologie au avantajul de a fi relativ intuitive pentru terminolog. Transcrierea directă a elementelor și relațiilor dintre acestea într-o bază de date este din ce în ce mai dificilă și mai riscantă.

Există încercări meritorii de realizare a unor interfețe "cuprinzătoare" pentru consultarea resurselor terminologice. Exemplele următoare sunt edificatoare în acest sens.

Primul exemplu ar putea provoca comentarii legate de complexitatea reală a înregistrării referințelor bibliografice cele mai obișnuite.

Copyright Cycorn Limited 2002 <http://www.cy.com.co.uk/>

These details identify the source of some Text appearing within one of the term entries.

• • x

```

Identificator ISO 1087-1:2000
Author family name TC B7/SC 1
Article title -
Page number 1
Book title OS/400
ISBN
Book edition Draft
Publication date 1999-04-22
Publisher

```

numeroaselor scrisuri, formate de date și limbi existente. În același timp trebuie adaptată și interfața utilizator potrivit locului și culturii căreia îi aparține acest printr-un proces nu mai puțin important de "localizare"

Multă vreme, prelucrarea automată a datelor a fost considerată satisfăcătoare realizabilă prin utilizarea setului ASCII de caractere. În prezent este însă absolut necesar ca:

- Utilizatorul calculatorului să poată tasta caractere și simboluri (vestimentare europene, est-europene, grecești și cirilice, cel puțin) folosind claviatură standard.
- Aplicația să prelucreze și să afișeze sau să imprime șiruri de caractere formate corect folosind seturi de caractere specifice fiecărei limbi.

Aceste cerințe pot fi realizate prin valorificarea calităților standardului Unicode de codificare prin utilizarea unor coduri de 16 biți pentru reprezentarea tuturor caracterelor pentru calculatoarele moderne care includ simbolurile tehnice și semnele speciale necesare imprimării textelor.

^ Explore the tabs below to set in numerous properties of the term. It is OK to leave many properties blank (undefined)

^ Explore the tabs below to set in numerous properties of the term. It is OK to leave many properties blank (undefined)

; Main | Gialli | Sage | Sound | Status | Descriptions

```

Terminology
Term type
Term ID
j Antonym term
j False friend
1 Short form of another term
Abbreviated form of another term
s Generate unique term ID
Target term
Target term
Target term
Target term

```

More term type

Commit changes Commit changes and close Rollback changes and close

Al doilea, ne determină să luăm în considerare următoarele:

La nivelul Uniunii Europene numărul limbilor pentru care este necesar suport terminologic este atât de mare (și sperăm încă în creștere) încât nu mai este posibilă multiplicarea tabelor bazelor de date potrivit numărului de limbi de lucru. Din fericire, "balizarea" documentelor permite identificarea și prelucrarea corect dependentă de limba în care au fost concepute acestea. Se vine astfel în sprijinul "globalizării" aplicațiilor informatice care sunt suport al resurselor terminologice multilingve dând posibilitatea acceptării, prelucrării și prezentării

Cu alte cuvinte la nivelul seturilor de semne necesare unei resurse terminologice multilingve se poate conta pe serviciile standardului Unicode și pe cele ale oricărei baze de date relaționale care acceptă Unicode.

Pentru indicarea formatelor de prezentare (fonte, punere în pagină, seturi de caractere) și a limbii utilizate se face apel la balizare astfel încât la nivelul câmpului vom găsi șiruri de caractere Unicode balizate.

Înscriserea datelor terminologice este facilitată de înscrierea lor în "categorii de date" bine definite (vezi ISO 12620). Dar numărul mare al acestor categorii și mai ales incidența ridicată a aparițiilor neprevăzute dinainte a unora noi face imposibilă alocarea unui câmp de date fiecărei categorii de date. Aceași observație poate fi făcută și asupra relațiilor dintre diferitele categorii de date care reflectă direct relațiile dintre elementele MARTIF. O soluție este o abstractizare suplimentară a datelor terminologice după încadrarea lor succesivă în șiruri de caractere balizate, categorii de date, elemente MARTIF.

În centrul modelului de date se află un set de 13 entități (atomi):

<u>Entitate</u>	<u>Descriere</u>
data category	o anumită clasă de informații terminologice (de exemplu: <u>term, part of speech</u>)
jdata category name	un nume agreat de utilizator (user-friendly), dependent de limbă, al unei anumite categorii de date (de exemplu, în română, " <u>termen</u> " pentru <u>term</u>)
data category index type	o strategie de indexare corespunzătoare unei anumite categorii de date (ISO 12620) (de exemplu: nu se indexează, se indexează ca valoare unică, se indexează <u>cuvânt cu cuvânt</u>)
	o anumită limbă, care dispune de o schema de codare uniformă care utilizează un singur set de caractere (de exemplu: <u>French, German, Italian</u>)
picklist	o combinație unică de caractere care poate fi utilizată pentru reprezentarea unei singure sau mai multor limbi (de exemplu: <u>ISO 8879-1, ISO 8859-2</u>)
	o mulțime de valori posibile ale unor date terminologice aparținând unei anumite categorii de date (ISO 12620) (de exemplu, pentru categoria "parte de vorbire": <u>noun, verb, adjective</u>)
	<u>o dată terminologică unică</u>
	<u>o dată (time stamp) care constituie valoarea unui element</u> <u>un număr care constituie valoarea unui element</u> un membru al unei liste care reprezintă valoarea unui element
	<u>șir de caractere care constituie valoarea unui element</u> în șir de caractere <u>r.prⁿ r onⁿ - : - ^ x—</u>

Primele 6 "articole" sunt "meta-entități"; ele sunt create și tabelele corespunzătoare sunt completate cu informații înainte de încărcarea oricărei date terminologice în baza de date. Prin completarea acestor table se conturează și se activează chiar modelul de date al bazei de date terminologice. Cu alte cuvinte ansamblul "meta-tabelor" definește structura care impune condiții și unifică datele terminologice de nivel molecular. Ele pot fi considerate atomi catalizatori a reacțiilor necesare combinării altor atomi în interacțiuni moleculare.

Celelalte 7 entități se încarcă direct prin proceduri de introducere a datelor sau prin import și cuprind datele terminologice vizibile pentru utilizatorul bazei de date. Informațiile conținute de aceste entități pot fi validate la nivel molecular folosind interogări SQL standard. Majoritatea interogărilor formulate de utilizator bazei de date se concentrează aproape în întregime asupra informațiilor încărcate în aceste entități.

Elementul central al aplicației pentru întreținerea unei astfel de baze de date este componenta de tip *parser* pentru crearea, validarea și prelucrarea documentelor MARTIF în particular (fără a ignora documentele SGML, HTML, XML). În mod obișnuit un *parser* este un modul software care examinează un document SGML prin confruntarea acestuia cu DTD-ul corespunzător. Rezultatul acestei examinări este de cele mai multe ori simplu: 'da' în situația în care documentul reprezintă o instanțiere validă a DTD-ului și 'nu' în cazul contrar. De cele mai multe ori *parser-ul* este capabil să 'normalizeze' documentul validat (aducându-l la o 'formă canonică') astfel încât facilitează formatarea, editarea și încărcarea documentului în baza de date.

Alături de *parser* și legat de acesta se află un *editor structurat*. Pornind de la DTD acesta propune utilizatorului pas cu pas opțiunile de compunere, sau modificare a unui document în conformitate cu definiția tipului corespunzător documentului. În cazul în care obiectivul este compunerea unui document SGML acesta poate asigura completarea teg-urilor necesare.

De cele mai multe ori sistemele de management al bazelor de date orientate spre text folosesc *fișiere inversate de indexare* a conținutului acestor baze pentru regăsirea informațiilor. Căutarea poate urmări apariția unui anume cuvânt sau a unui model oarecare într-un document sau în o parte a acestuia. Identificarea subdiviziunilor documentului se poate face folosind tocmai tag-urile care acesta este marcat, respectiv modul în care acestea au fost transcrise în relații dintre tabelele bazei de date.

În fine, o componentă deosebit de importantă este aceea care realizează funcțiile de *import-export* ale datelor terminologice spre și dinspre baza de date.

Terminologia calității

Realizarea unor resurse terminologice multilingve este de mai multă vreme în centrul preocupărilor Asociației Române pentru Terminologie (TERMROM). Începând de anul trecut pe lista temelor având aceeași orientare se înscrie proiectul "Terminologie armonizată cu prevederile EURODICAUTOM în domeniul calitate și standardizare". Proiectul a fost inițiat de Ministerul Educației și Cercetării și este finanțat în cadrul Programului CALIST.

Obiectivele principale ale acestui subprogram sunt:

- Asigurarea flexibilității necesare pentru a răspunde operativ la cerințele concrete de rezolvare a unor teme de cercetare care decurg din prioritățile stabilite prin strategiile guvernamentale adoptate pe domenii specifice, în procesul integrării României în U.E.
- Asigurarea condițiilor de dezvoltare și armonizare a sistemului de standarde naționale în conformitate cu cerințele organismelor de standardizare europene și internaționale;
- Asigurarea unei baze terminologice științifice pentru elaborarea standardelor de calitate românești, precum și în ceea ce privește condițiile de aplicabilitate a prevederilor standardelor internaționale și europene adaptate ca standarde românești;
- Clarificarea condițiilor pe care trebuie să le îndeplinească produsele românești în vederea pătrunderii lor pe piața unică a Uniunii Europene și produsele introduse în România.

Pentru realizarea obiectivelor proiectului au fost prevăzute următoarele activități:

- Întocmirea unui Proiect Terminologic pentru definirea și înregistrarea terminologiei domeniilor calitate și standardizare utilizate în documentele oficiale ale Uniunii Europene, conform prevederilor EURODICAUTOM și standardelor internaționale;
- Extragerea, traducerea și structurarea terminologiei domeniilor calitate și standardizare;
- Proiectarea, programarea și implementarea unei Baze de date conform Proiectului Terminologic capabilă să gestioneze toate domeniile EURODICAUTOM;
- Înregistrarea în baza de date a terminologiei domeniilor calitate și standardizare;
- Elaborarea unei aplicații informatice de administrare a bazei de date terminologice și de transfer de date terminologice conform formatului standard ISO pentru lucrul în rețea;

Realizarea unui site web pentru promovarea Bazei de date terminologice și punerea acesteia la dispoziția publicului.

A fost avizat Proiectul Terminologic, au fost stabilite cerințele pe care să le satisfacă suportul informatic, s-a constituit un fond de termeni specifici extrași din EURODICAUTOM și din Tezaurul rațional al CEI și au fost demarate activitățile pentru realizarea unei baze de date relaționale EUROCAST pentru înregistrarea acestora.

Bibliografie

- [1] **ISO 639:1988**
Code for the representation of names of languages
- [2] **ISO 639-2:1998**
Code for the representation of names of languages - Part 2: Alpha-3 code
- [3] **ISO 704:2000**
Terminology work - Principles and methods
- [4] **ISO 860:1996**
Terminology work - Harmonization of concepts and terms
- [5] **ISO 1087-1:2000**
Terminology work - Vocabulary - Part 1: Theory and application
- [6] **ISO 1087-2:2000**
Terminology work - Vocabulary - Part 2: Computer applications
- [7] **ISO 1951:1997**
Lexicographical symbols particularly for use in classified defining vocabularies
- [8] **ISO 6156:1987**
Magnetic tape exchange format for terminological/lexicographical records (MATER)
- [9] **ISO 10241:1992**
Preparation and layout of international terminology standards
- [10] **ISO 12199:2000(E)**
Alphabetical ordering of multilingual terminological and lexicographical data represented in the Latin alphabet
- [11] **ISO 12200:1999**
Computer applications in terminology - Machine-readable terminology interchange format (MARTIF) - Negotiated interchange
- [12] **ISO/TR 12618:1994**
Computer aids in terminology - Creation and use of terminological databases and text corpora
- [13] **ISO 12620:1999**
Computer applications in terminology - Data categories
- [14] **ISO 15188:2001**
Project management guidelines for terminology standardization

SECȚIUNEA i1

TEHNOLOGII ALE LIMBAJULUI SCRIS

Ro-Balkanet - ontologie lexicalizată, în context multilingv, pentru limba română

Dan TUFIS, Institutul de Cercetări pentru Inteligența Artificială,
Academia Română, București
Calea 13 Septembrie nr. 13, 74311, sector 5

tufis@racai.ro

Dan CRISTEA, Facultatea de Informatică, Universitatea A.I.Cuza, Iași
Str. General Berthelot, nr. 16

dcristea@infoiasi.ro

Rezumat

Cerințele creării unei ontologii multilingve de tipul EuroWordNet sunt frecvent contradictorii și dacă problemele de compatibilitate nu sunt considerate în etapele timpurii ale construcției, o armonizare tardivă se poate dovedi dificilă sau imposibilă. Mai exact, există două probleme majore de compatibilitate care trebuie avute în vedere și anume: acoperirea conceptuală - în sensul că fiecare lexicon monolingv ar trebui să conțină lexicalizări ale aceluiași fond conceptual și coeziunea interpretativă - în sensul că interpretarea relațiilor folosite în fiecare din ontologiile cuprinse în ontologia multilingvă trebuie să fie identică. În lucrare sunt discutate ambele aspecte și prezentate soluțiile adoptate în vederea satisfacerii criteriilor de consistență și coerență multilinguală a wordnet-ului pentru limba română.

1. Limbă, resurse lingvistice și comunicare electronică

Cercetarea în domeniul tehnologiilor limbajului este un domeniu ce are deja istorie în știința calculatoarelor, dar, actualmente, motivațiile sale depășesc sfera interesului pur științific sau comercial. Păstrarea identității limbilor și culturilor naționale în cadrul globalizant al societății informaționale și a cunoașterii readuce în actualitate avertismentul lui Alain Danzin [1]: **"în era electronică, este esențial pentru supraviețuirea unei limbi ca ea să fie folosită în sistemele de informare electronică."** Avansul științific și tehnologic obținut în cei 10 ani scurși de la raportul prezentat de Danzin Comisiei Europene a condus la maturizarea unor teorii, tehnologii, metode și la dezvoltarea altora noi, dar mai ales a permis

definirea unor standarde pentru realizarea unitară a ceea ce generic se numește *resurse lingvistice fundamentale* ale unei limbi. Caracterul multilingual al societății cunoașterii, în care conceptul de "unitate prin diversitate" se referă în primul rând la prezervarea limbilor și culturilor actuale, a generat o deosebită efervescentă, puternic stimulată de organismele internaționale - în primul rând de Comisia Europeană - asupra cercetării în domeniul resurselor multilingve. Metodologic, tehnologia limbajului natural creează o distincție netă între prelucrări și date, între "mașinăria software de prelucrare a limbajului" numită și *lingware* și cunoștințele lingvistice, numite cum arătam *resurse lingvistice*, necesare funcționării acestei mașinării. Dihotomia *lingware* - *resurse lingvistice*, susținută de standardele de reprezentare și codificare a cunoștințelor lingvistice permite dezvoltarea independentă a celor două componente ale unui sistem de prelucrare a limbajului. *Lingware*-ul este independent de limbă și intră tot mai pregnant în zona ingineriei software. El poate fi dezvoltat de specialiști de oriunde fără ca aceștia să fie preocupați de limba pentru care va fi folosit. Resursele lingvistice însă sunt de competența specialiștilor vorbitori nativi ai limbii respective. În condițiile în care aceste resurse lingvistice sunt realizate în conformitate cu standardele sau practicile internaționale, ele pot fi integrate în sistemele de comunicare electronică, nu doar pentru prelucrare monolingvă ci mai ales pentru prelucrări multilingve. Beneficiile alinierii la standardele internaționale în realizarea resurselor lingvistice sunt enorme, și putem considera un exemplu foarte simplu. Să presupunem că suntem interesați de un anumit subiect și, folosind imensul ocean informațional ce este Internet-ul, apelăm la un așa numit "motor de căutare", un program a cărui funcționalitate asigură identificarea documentelor electronice ce conțin informații potențial relevante pentru subiectul nostru de interes. Acest gen de serviciu informațional este asigurat de "motoare de căutare" precum Google, Altavista, Excite și multe altele. Documentele interesante din punctul nostru de vedere ar putea să fie scrise în limba engleză, franceză, germană, română sau orice altă limbă. Dar pentru a le regăsi pe toate, indiferent în ce limbă am formulat cererea noastră de regăsire, motorului general de căutare îi sunt necesare resursele lingvistice specifice limbilor în care documentele ar putea exista. Dacă aceste resurse lingvistice există pentru engleză, franceză, germană, italiană etc. și ele sunt reprezentate în același format standardizat, rezultatul cercetării noastre documentare va fi o colecție de documente tratând subiectul de interes în oricare dintre aceste limbi. Un astfel de serviciu, numit regăsire documentară multilingvă este o realitate pentru toate limbile "mari", o calificare ce nu are acoperire în substratul cultural ci doar în ceea ce se numește "nivelul de informatizare al limbii". Procesul de informatizare a unei limbi naturale permite potențarea și diseminarea ei prin mijloacele tehnologice ale societății informaționale.

2. Lexicalizarea abordărilor în tehnologia limbajului și conceptul "wordnet"

Lexicul este fără îndoială cea mai importantă resursă lingvistică a unei limbi. Marea majoritate a cercetării actuale, atât în lingvistica formală cât mai ales în tehnologia limbajului, plasează componenta lexicală în centrul modelelor de limbă, sub influența a ceea ce a fost numită abordarea *lexicalizată* sau *lexicalistă* a studiului limbii. Nu este de mirare, deci, enormul interes pentru dezvoltarea de resurse lexicale multilingve. Studiul computațional al dicționarelor electronice natura informației ce trebuie inclusă în ele și tipul de prelucrări pe care le poate facilita o anumită structurare a unui mare voium lexical a fost, fără îndoială fundamental influențat de proiectul WordNet, lansat în urmă cu mai mult de 25 de ani la Universitatea din Princeton sub conducerea reputatului psiholingvist George Miller. WordNet, resursă publică, este o uriașă rețea semantică lexicală în care peste 100.000 de *înțelesuri* lexicalizate în limba engleză prin mai mult de 130.000 de cuvinte sunt asociate între ele prin relații semantice și/sau lexicale [2]. Fondul lexical este distribuit în 4 rețele semantice corespunzând categoriilor gramaticale deschise: substantive, verbe, adjective și adverbe. Noțiunea de *înțeles* (*meaning*) este în WordNet echivalată cu cea de concept și este reprezentată printr-o serie sinonimică în care fiecare cuvânt al seriei are asociat un număr ce identifică sensul în care cuvântul respectiv are înțelesul asociat conceptului. Seria sinonimică ce identifică un înțeles se numește *sinset*. Relațiile existente între sinseturi sunt de diferite tipuri, depinzând de categoria gramaticală a cuvintelor ce alcătuiesc un anumit sinset (antonimie/sinonimie, hiponimie/hiperonimie, holonimie/meronimie, troponimie etc). Influența proiectului WordNet a fost enormă în domeniul tehnologiei limbajului (exprimată poate și prin faptul că acum, în limbajul tehnic ce puțin, cuvintele "wordnet" și "synset" au devenit substantive comune, importate pr calchiere în mai toate limbile), iar beneficiile acestui concept sunt atât de evidente încât Comisia Europeană, între 1996 și 1998, a finanțat un proiect similar de mare anvergură numit EuroWordNet [3]. Acest proiect, extrem de ambițios și-a propus nu numai realizarea concertată de wordneturi monolingve pentru limbile europene de circulație internațională (engleză, franceză, germană, italiană, olandeză, spaniolă) dar a introdus o cerință fundamental nouă, anume corelarea multilingua a celor 6 rețele semantice lexicale, astfel încât dintr-un sinset al unei limbi să poată ajunge în echivalentul de traducere al oricăreia dintre celelalte 5 limbi. Fa de relațiile originale din WordNet, EuroWordNet propune un inventar mult bogat (90) de relații cum ar fi cele tematice de tip cazual (Agent, Patient, Instrument, Location, Direction) sau cele corelând sensurile derivaților lexic (XPOS-SYNONYMY: a adora - adorație).

Soluția tehnică pentru corelarea multilinguală a rețelelor semantice monolingve a fost definirea unui index interlingual (ILI), independent de limbă, conținând reprezentări conceptuale ale înțelesurilor lexicalizabile în limb

proiectului. Fiecare înțeles din oricare din limbile reprezentate în rețeaua semantică multilingvă este pus în corespondență, în general, cu un singur concept al indexului interlingual. Aceste corespondențe se realizează prin intermediul a 20 de tipuri distincte de relații binare. Sinseturile (seriile sinonimice) din două sau mai multe limbi care sunt puse în corespondență cu același concept din ILI sunt considerate echivalenți de traducere, natura echivalenței de traducere fiind definită de tipul relațiilor ce definesc corespondența dintre sinseturile respective și conceptul comun.

Inițial, indexul multilingual a fost constituit ca o mulțime nestructurată a tuturor înțelesurilor lexicalizate în WordNet (cu alte cuvinte în engleză). Ulterior, prin dezvoltarea wordneturilor monolingve, ILI a fost îmbogățit și cu reprezentări conceptuale cu lexicalizări ce nu se regăsesc în engleză.

O altă inovație a proiectului EuroWordNet a fost adoptarea unei mulțimi de primitive semantice, independente de limbaj, în termenii cărora așa-numitele *concepte de bază* din ILI au fost asociate cu descrieri *ontologice*. Prin importul acestor descrieri la nivelul lexicalizărilor prin echivalenți de traducere (și, prin moștenire, la hiponimii acestora) în fiecare dintre wordneturile monolingve, în EuroWordNet se poate vorbi de o ontologie lexicală multilingvă. O prezentare în detaliu a proiectului EuroWordNet se poate găsi în [4].

După 3 ani, proiectul EuroWordNet inițial a fost extins pentru o perioadă de încă doi ani (EuroWordNet II) și a încorporat încă 4 limbi: bască, catalană, cehă și estoniană. Proiectul EuroWordNet II s-a încheiat în anul 2000 cu realizarea unor nuclee a căror extensie a rămas în exercițiul financiar al autorităților naționale.

3. Limba română în contextul proiectului BALKANET, extensie a EuroWordNet

În septembrie 2001 a fost lansat proiectul european BALKANET (IST - 2000 - 29388), o continuare firească a proiectului EuroWordNet II care aduce alături de cele 10 limbi europene alte 5 limbi din zona balcanică: bulgară, greacă, română, sârbo-croată, turcă [5]. Ca și în EuroWordNet, ontologiile lexicele monolingve sunt corelate printr-o mulțime de concepte interlinguale, corespondențele fiind stabilite cu ajutorul unor relații de echivalență complexe (*eq-synonymy*, *eq-near-synonymy*, *eq-has-hyperonym*, *eq-has-hypernym* etc).

Reprezentanții din România în acest proiect, care va dura trei ani, sunt Institutul Academiei Române de Cercetări pentru Inteligență Artificială din București (coordonator Dan Tufiș) și Facultatea de Informatică a Universității A.I.Cuza din Iași (coordonator Dan Cristea) și în realizarea obiectivelor proiectului sunt implicați numeroși specialiști, atât informaticieni cât și lingviști. Desigur, participarea românească în acest proiect și angajarea față de obiectivele

proiectului nu s-au bazat numai pe entuziasm ci pe activități și rezultate anterioare importante, pe *surse lingvistice* primare [6] de referință ale limbii române, implementate ca *resurse lingvistice* [6] în format standardizat și pe o multitudine de programe de prelucrare dezvoltate de-a lungul a mulți ani de cercetare, în cea mai mare parte prin finanțare internațională.

3.1. Corpusuri

În cadrul proiectelor europene Multext-East și TELRI [7], [8], [9], [10], [11] a fost creat un corpus paralel în 7 limbi, foarte detaliat adnotat, bazat pe romanul "1984" al lui Orwell și un alt corpus paralel în 25 de limbi, bazat pe "Republica" lui Platon. Adnotarea folosită inițial a fost conformă cu standardul TEI (<http://www.tei-c.org/>), dar ulterior, odată cu cristalizarea standardului CES [12], corpusurile au fost re-adnotate (automat) în conformitate cu CES. Acestea sunt două corpusuri relativ mici (câte aproximativ 110.000 cuvinte în fiecare limbă) dar, datorită acurateței proceselor de etichetare și de aliniere (validate manual), au fost extrem de folositoare pentru diverse aplicații, de la construirea modelelor lingvistice pentru etichetare morfo-sintactică [13], clasificare a documentelor [14], extragere de echivalenți de traducere [15], până la discriminarea automată a sensurilor [16]. Pe lângă corpusurile multilingve s-au construit alte două corpusuri monolingve mult mai mari: un corpus literar bazat pe diverse romane (conținând aproximativ 1.500.000 cuvinte) și un corpus jurnalistic (conținând peste 100.000.000 cuvinte). Ambele corpusuri au fost segmentate, etichetate și lematizate automat¹.

3.2. Dicționare explicative: WEB-LEX și XML-LEX

Principalul dicționar pe care l-am folosit în analiza noastră este Dicționarul Explicativ al Limbii Române [17], referința lexicografică pentru limba română contemporană, dicționar realizat de Institutul de Lingvistică "Iorgu Iordan"² a Academiei Române. În urma analizelor statistice de frecvență în corpusurile menționate, au fost selectate și introduse în format electronic cele mai frecvente 23.000 de cuvinte titlu din DEX. Acest nucleu DEX a fost convertit într-o bază de date lexicală în cadrul proiectului european CONCEDE (*CONortium for Central European Dictionary Encoding*) [11] și al proiectului prioritar al Academiei WEB-LEX [18]. Ulterior, îmbogățit continuu prin culegere manuală din alte câteva dicționare explicative (DEX'84, DOOM, DLRM), la inițiativa unor tineri entuziași atât din țară cât și din diasporă (vezi de pildă: <http://dex.franco.com>), WEB-LEX a fost corectat sub aspect sintactic-structural și codificat într-un format standardizat respectând convențiile lexicografice utilizate de DEX și, în măsura posibilului, conținutul său textual. Uneori, din considerente legate de consistența structurală

¹ Multe dintre aceste resurse pot fi găsite pe situl Consorțiului de Informatizare pentru Limba Română (ConslLR) la adresa <http://consilr.info.uaic.ro>

² Noua sa denumire este Institutul de Lingvistică "Iorgu Iordan-Al. Rosetti

s-au operat o serie de modificări asupra conținutului. De asemenea, o serie de erori evidente în sursa primară au fost corectate de specialiști avizați. Deși mai bogat (în prezent WEB-LEX conține aproape 70.000 de intrări, față de cele circa 56.000 de intrări din DEX'96), influența DEX a fost fundamentală în dezvoltarea WEB-LEX. Pe de altă parte, eventualele critici asupra conținutului, acolo unde ne-am despărțit de DEX, în nici un caz nu trebuie puse în seama Institutul de Lingvistică "Iorgu Iordan-AI. Rosetti" ci a noastră. Din acest motiv, preferăm să ne referim la WEB-LEX ca la un dicționar **de tip** DEX și nu ca variantă computațională a DEX-ului.

Codificarea conținutului WEB-LEX s-a realizat folosind limbajul de anotare XML. Implementarea, ce explicitează toate convențiile tipografice precum și informațiile implicite, a condus la un volum textual de date de circa 8-10 ori mai mare față de conținutul textual echivalent al DEX-ului. Anotarea XML a fost realizată automat, cu ajutorul compilatorului DIC [18]. Compilatorul a fost generat automat folosind JavaCC®, pe baza unei gramatici LL(7) ce descrie structura formală a intrărilor în DEX. DIC poate fi folosit pentru a genera documente XML (conform cu DTD-ul CONCEDE) pentru orice dicționar ce folosește convențiile tipografice adoptate în DEX. În [19] sunt prezentate o multitudine de dicționare realizate sau aflate în curs de realizare la Institutul de Lingvistică "Iorgu Iordan-AI. Rosetti" și presupunând că ele urmăresc convențiile tipografice și lexicografice adoptate în DEX, toate aceste surse lingvistice de referință pentru limba română ar putea fi transformate, cu efort minim, în resurse computaționale fundamentale pentru prelucrarea automată.

Varianta codificată a dicționarului nostru este numită XML-LEX iar structura sa este descrisă de DTD-ul (*Document Type Definition*) pe care îl reproducem în figura 1, dezvoltat în cadrul proiectului CONCEDE.

```

<!-- CONCEDE project - Deliverable DR2.1: concede.dtd -->
<!-- copyright CONCEDE project consortium, 1999 -->
<!-- ENTITY DECLARATIONS -->
<!ENTITY % a.global'
    id      ID      #IMPLIED
    n       CDATA   #IMPLIED
    lang    IDREF   #IMPLIED'
<!ENTITY % a.text'
    %a.global;
    rend    CDATA   #IMPLIED
    wsd     CDATA   #IMPLIED'
<!ENTITY % basetags'

```

```

(orth|pron|hyph|syll|stress|pos|gen|case|number|gram|tns|
mood|q|source|gloss|usg|def|per|aspect|degree|voice|eg|
etym|xr|trans|itype|subc)' >
<!ENTITY % dictbase.seq '#PCDATA | na' >
<!-- STRUCTURAL ELEMENTS -->
<!ELEMENT dictionary (body) >
<!ATTLIST dictionary %a.global;
    type      CDATA      #IMPLIED
    version   CDATA      #REQUIRED
    xmkspace  (default | preserve) 'preserve' >
<!ELEMENT body (entry+) >
<!ATTLIST body %a.global; type CDATA #IMPLIED >
<!ELEMENT entry
    (hw, (%basetags;|struc|alt|brack)*) >
<!ATTLIST entry %a.global; type CDATA #IMPLIED >
<!ELEMENT struc (%basetags; j struc | alt | brack)* >
<!ATTLIST struc %a.global; type CDATA #IMPLIED >
<!ELEMENT trans (%basetags; | struc | alt | brack)* >
<!ATTLIST trans %a.global; type CDATA #IMPLIED >
<!ELEMENT alt (%basetags; | brack)* >
<!ATTLIST alt %a.global; type CDATA #IMPLIED >
<!ELEMENT brack (%basetags;)* >
<!ATTLIST brack %a.global; type CDATA #IMPLIED s
<!-- CONTENT ELEMENTS -->
<!ELEMENT voice (%dictbase.seq;)* >
<!ATTLIST voice %a.text; >
<!ELEMENT tns (%dictbase.seq;)* >
<!ATTLIST tns %a.text; >
' <!ELEMENT syll (%dictbase.seq;)* >
<!ATTLIST syll %a.text; >
<!ELEMENT subc (%dictbase.seq;)* >
<!ATTLIST subc %a.text; >

```



```

<!ELEMENT stress (%dictbase.seq;)* >
<!ATTLIST stress %a.text; >
<!ELEMENT source (%dictbase.seq;)* >
<!ATTLIST source %a.text; >
<!ELEMENT pos (%dictbase.seq;)* >
<!ATTLIST pos %a.text; >
<!ELEMENT per (%dictbase.seq;)* >
<!ATTLIST per %a.text; >
<!ELEMENT number (%dictbase.seq;)* >
<!ATTLIST number %a.text; >
<!ELEMENT na (#PCDATA) >
<!ATTLIST na %a.text; >
<!ELEMENT mood (%dictbase.seq;)* >
<!ATTLIST mood %a.text; >
<!ELEMENT m (%dictbase.seq;)* >
<!ATTLIST m %a.text; >
<!ELEMENT lang (%dictbase.seq;)* >
<!ATTLIST lang %a.text; >
<!ELEMENT itype (%dictbase.seq;)* >
<!ATTLIST itype %a.text; >
<!ELEMENT hw (%dictbase.seq;)* >
<!ATTLIST hw %a.text; >
<!ELEMENT gram (%dictbase.seq;)* >
<!ATTLIST gram %a.text; >
<!ELEMENT gen (%dictbase.seq;)* >
<!ATTLIST gen %a.text; >
<!ELEMENT degree (%dictbase.seq;)* >
<!ATTLIST degree %a.text; >
<!ELEMENT case (%dictbase.seq;)* >
<!ATTLIST case %a.text; >
<!ELEMENT aspect (%dictbase.seq;)* >
<!ATTLIST aspect %a.text; >

```

```

<!ELEMENT hyph (%dictbase.seq;)* >
<!ATTLIST hyph %a.text; >
<!ELEMENT eg (source | q | gloss)* >
<!ATTLIST eg %a.global; >
<!ELEMENT pron (%dictbase.seq;)* >
<!ATTLIST pron %a.text; type CDATA #IMPLIED >
<!ELEMENT q
(%dictbase.seq; | gloss |ptr [xptr |oref])* >
<!ATTLIST q %a.text; type CDATA #IMPLIED >
• <!ELEMENT etym
(%dictbase.seq; | gloss | lang | m |ptr [xptr |oref])* >
<!ATTLIST etym %a.text; type CDATA #IMPLIED >
<!ELEMENT xr (%dictbase.seq; | ptr [xptr])* >
<!ATTLISTxr %a.text; type CDATA #IMPLIED >
<!ELEMENT def (%dictbase.seq; | ptr [xptr |oref |usg])* >
<!ATTLIST def %a.text; type CDATA #IMPLIED >
I <!ELEMENT gloss (%dictbase.seq; | ptr [xptr |oref])* >
<!ATTLIST gloss %a.text; type CDATA #IMPLIED >
<!ELEMENT orth (%dictbase.seq; | ptr [xptr |oref |usg])* >
<!ATTLIST orth %a.text;
expansion NMTOKEN #IMPLIED
extent (full | pref | suff | part) "full"
type CDATA #IMPLIED >
<!ELEMENT usg (%dictbase.seq;)* >
<!ATTLISTusg %a.text;
type (syn|hyper|colloc|comp|plev|acc|lang|gram|obj|
subj|verb|hint|geo|domjregister|time|style|
| hyponym | antonym | other) "other" >
| <!ELEMENT oref EMPTY >
j <!ATTLIST oref %a.text;
| target IDREF #IMPLIED
| fullform NMTOKEN #IMPLIED >

```

```

<!ELEMENT ptr EMPTY
<!ATTLIST ptr %a.text;
  corresp IDREFS #IMPLIED
  next IDREF #IMPLIED
  prev IDREF #IMPLIED
  type CDATA #IMPLIED
  resp CDATA #IMPLIED
  crdate CDATA #IMPLIED
  targType NMTOKEN #IMPLIED
  targOrder (y | n | u) v
  evaluate (all | one | none) #IMPLIED
  target IDREFS #REQUIRED >
<!ELEMENT xptr EMPTY >
<!ATTLIST xptr %a.text;
  corresp IDREFS #IMPLIED
  next IDREF #IMPLIED
  prev IDREF #IMPLIED
  type CDATA #IMPLIED
  resp CDATA #IMPLIED
  crdate CDATA #IMPLIED
  targType NMTOKEN #IMPLIED
  targOrder (y | n | u) "u"
  evaluate (all | one | none) #IMPLIED
  target NMTOKEN #REQUIRED

```

Figura 1: DTD-ul Concede, utilizat la implementarea XML-LEX

Această structură de codificare a fost adoptată în implementarea unui număr mare de dicționare, reprezentând un standard "de facto" în lexicografia computațională actuală [20]. Detalii suplimentare privind semantica entităților folosite în codificare și a atributelor acestora pot fi găsite în documentația tehnică a proiectului la adresa www.itri.bton.ac.uk/projects/concede/. În tabelul de mai jos, sunt exemplificate reprezentarea tipografică (de tip DEX) și reprezentarea codificată în XML.

DEX

ZA[^], zale, s.f. 1. Fiecare dintre ochiurile unui lanț; p. gener. (la pl.) lanț. • Lănțișor de metal întrebuințat uneori ca podoabă. • Cusătură în formă de lănțișor, executată de obicei la broderii. 2. (La pl.) împletitură executată din inele mici de fier legate unul de altul; p. ext. armură făcută din această împletitură, cu care se îmbrăcau oștenii în antichitate și în evul mediu, spre a se apăra de loviturile dușmanilor. [Var.: (reg.) zălă.zea s.f.] - Cf. ngr. zăva.

XML-LEX

```

<entry type="homonym" id="ZĂ.2">
  <hw>ZA</hw>
  <alt>
    <brack>
      <gram>nominativ_feminin__singularjndefinit</gram>
        <orth>ZA</orth>
      </brack>
      <gram>nominativ_feminin_pluraljndefinit</gram>
        <orth>zale</orth>
      </brack>
    </alt>
    <pos>substantiv</pos>
    <gen>feminin</gen>
    <struc n="1">
      <alt>
        <def>Fiecare dintre ochiurile unui lanț</def>
      <brack>
        <usg type="hyper">prin generalizare </usg>
        <usg>la pl.</usg>
        <def>lanț.</def>
      </brack>
    </alt>
    <struc type="Sec">
      <def>Lănțișor de metal întrebuințat uneori ca podoabă. </def>
    </struc>
    <struc type="Sec">
      <def>Cusătură în formă de lănțișor, executată de obicei la broderii.</def>
    </struc>
    </struc>
    <struc n="2">
      <usg>La pl.</usg>
    <alt>
      <def>împletitură executată din inele mici de fier legate unul de altul</def>
    <brack>
      <usg type="hyper">prin extensiune</usg>
      <def>armură făcută din această împletitură, cu care se îmbrăcau oștenii în antichitate și în evul mediu, spre a se apăra de loviturile dușmanilor.</def>
    </brack>
  </entry>

```

```

</alt>
</struc>
<struc type="Varianta">
<alt>
<brack>
<orth> zală</orth>
<stress> zălă</stress>
<usg>reg.</usg>
</brack>
<orth> zea</orth>
</alt>
<pos>substantiv</pos>
<gen>feminin</gen>
</struc>
<etym>
Cf.
<lang>ngr.</lang>
zâva.
</etym>
</entry>

```

Figura 2: Conținut primar și codificarea echivalentă în XML (cf. CONCEDE.dtd)

În tabelul din Figura 2, sunt exemplificate reprezentarea tipografică (de tip DEX) și reprezentarea codificată în XML. Menționăm că reprezentarea tipografică din coloană stângă a Figurii 2 s-a obținut automat, folosind un convertor XML de format, proiectat astfel încât rezultatul generării (interpretarea marcajului XML) să fie cât mai apropiat de aspectul dicționarului tipărit. Structura de dicționar, definită mai jos, este suficient de generală pentru a permite implementarea diferitelor tipuri de dicționare. În fapt, DTD-ul CONCEDE a fost utilizat pentru codificarea a două dicționare bilingve: un dicționar Sloven-Englez și un dicționar Român-Francez.

Adnotarea XML fiind independentă atât de convențiile tipografice cât și de limba dicționarului, este posibilă căutarea multi-criterială a informației în unul, două sau mai multe dicționare explicative ale unor limbi diferite. De pildă, o căutare multi-criterială ar putea fi parafrazată astfel:

Găsește și afișează toate intrările ce corespund substantivelor feminine, de origine neo-greacă și al căror cuvinte titlu încep cu secvența de litere ZA. O astfel de căutare va avea ca rezultat tipărirea cel puțin a intrării corespunzătoare cuvântului titlu ZA²:

Z A², zale, s.f. 1. Fiecare dintre ochiurile unui lanț; *p. gener.* (la pl.) lanț. • Lănțisor de metal întrebuițat uneori ca podoabă. • Cusătură în formă de

lănțisor, executată de obicei la broderii. 2. (La pl.) împletitură executată din inel mici de fier legate unul de altul; *p. ext.* armură făcută din această împletitură, care se îmbrăcau oștenii în antichitate și în evul mediu, spre a se apăra de loviturile dușmanilor. [Var.: (reg.) zălă, zea s.f.] - Cf. ngr. zâva.

33. Alte dicționare, lexicoane; indexul interlingual

Unul dintre rezultatele proiectului Multext-East îl constituie un lexicon de forme ocurență (LFO), cu peste 450.000 de intrări, care conține triplete de tip <cuvânt, lernă, cod_morfo-sintactic>. Acest lexicon va fi completat cu formele flexionare (generate automat) ale lemelor din XML-LEX nereprezentate în LFO. Codificarea folosită este compatibilă cu recomandările Eagles (<http://www.ilc.pcnr.it/EAGLES/home.html>) pentru adnotarea morfo-sintactică și este documentată pe larg în [10].

O altă resursă lexicală esențială a fost Dicționarul de Sinonime al Limbii Române - DSLR [21], care a fost transpus în formă electronică la Facultatea de Informatică a Universității "A.I.Cuza" din Iași. Forma electronică a DSLR a fost convertită în format XML astfel încât aceeași interfață ce a fost dezvoltată pentru XML-LEX funcționează și cu XML-DSLR.

Din corpusurile paralele menționate mai sus și folosind programul conceput să implementeze metodologia noastră de extragere a echivalențelor de traducere [22], [23], [24] s-a construit un dicționar bilingv Român - Englez (de asemenea transpus în format XML). Acest lexicon bilingv a fost validat manual și îmbogățit cu noi intrări din diverse surse publice.

În sfârșit, o resursă extrem de valoroasă a fost și Indexul Interlingual (IL) al EuroWordNet, exportat în format XML cu editorul VisDic produs la Universitatea Masaryk din Brno [25].

3.4. Alegerea nucleului lexical

Vom da câteva definiții ale unor noțiuni pe care le vom folosi în cele ce urmează.

Când ne plasăm într-un context monolingv, vorbim despre *sensuri înțelesuri* și *sinseturi*. Un cuvânt are unul sau mai multe *sensuri*. Un sens referă un *înțeles*. În EuroWordNet sensurile unui cuvânt sunt numerotate în funcție de frecvența lor, iar sensul unei leme este denotat adăugând numărul sensului la forma ortografică a acesteia. O mulțime de sensuri astfel specificate (ex. *action, activity, activeness*) care referă același înțeles este numit *sinset* și constituie în suși denotația înțelesului sensurilor din sinset. Cu alte cuvinte, un sinset reprezintă *lexicalizarea unui înțeles* în contextul monolingv curent.

Dacă abstractizăm noțiunea de *înțeles*, definită ca mai sus, astfel încât să nu mai facem referirea la un anumit context monolingv, vom vorbi despre *concepte*

care sunt referite de *înțelesurile* lexicalizate în diferitele limbi. Așadar, putem vorbi despre concepte care au sau nu realizare lingvistică într-o limbă sau alta. Un concept este un construct cognitiv, independent de limbă, care în EuroWordNet este totdeauna lexicalizat cel puțin într-una dintre limbi. Un concept este mai departe rafinat în termeni de distincții semantice elementare (trăsături semantice), deci putem vorbi despre gruparea conceptelor în funcție de trăsăturile lor semantice.

În EuroWordNet și deci și în BALKANET, ILI este definit ca o colecție nestructurată de intrări de forma: <ILI-index><descriere ontologică><glosă> {domeniu}. Indexul interlingual inițial a fost construit plecând de la versiunea 1.5 a Wordnet-ului și deci glosele pentru fiecare concept au fost importate direct din sinsetul englezesc care se referă la înțelesul conceptualizat în ILI.

Pentru a facilita o cât mai bună intercorelare a wordneturilor monolingve din cadrul proiectului și pentru a înlesni extensia lor ulterioară, consorțiul proiectului a decis ca procesul implementărilor paralele să fie centrat pe concepte (independente de limbă) selectate de comun acord, la momente succesive de timp.

O primă selecție a constituit-o mulțimea așa-numitelor "concepte de bază" definite în EuroWordNet ca fiind acele concepte din ILI lexicalizate în limba engleză (în WORDNET) prin sinseturi plasate pe un nivel ierarhic cât mai sus și, în plus, care au un număr mare de hiponimi direcți (tot în WORDNET). Rațiunea acestei decizii a constat în faptul că, aceste concepte fiind foarte generale și totodată productive în definirea unor concepte mai particulare, este foarte probabil ca ele să fie lexicalizate în majoritatea limbilor de interes. Acest lucru a fost probat atât în EuroWordNet cât și în BALKANET. Mulțimea *conceptelor de bază* (o motivație mai detaliată a selecției lor este prezentată în [4] în raport cu obiectivele EuroWordNet) conține 1.310 concepte, fiecăruia dintre ele fiindu-i atașată o glosă explicativă și o *descriere ontologică* (vezi [26]).

După implementarea, în toate cele 5 limbi ale proiectului, a nucleelor de ontologii lexicale corespunzând conceptelor de bază, s-a făcut o nouă selecție, de data aceasta, conținând 4.000 de noi concepte interlinguale.

Selecția a avut în vedere, pe de o parte maximizarea compatibilității cu EuroWordNet, iar pe de altă parte relevanța stocului lexical pentru fiecare limbă din perspectivă monolingvă. Primul criteriu a fost operaționalizat alegându-se acele concepte lexicalizate în cele mai multe limbi din EuroWordNet. Limita inferioară a numărului de limbi a fost fixată la 5, astfel încât după implementarea acestor concepte în BALKANET ele să fie lexicalizate în cel puțin 10 limbi.

Criteriul relevanței monolingve a condus la propunerea mai multor mulțimi candidate de concepte. Pentru fiecare limbă a proiectului au fost efectuate analize cantitative în context strict monolingv. Metodele de analiză au diferit de la partener

la partener, în raport cu datele și instrumentele disponibile pentru limbă. După analiza acestor mulțimi, au fost incluse în mulțimea finală acele concepte care au apărut în cel puțin două propuneri. Mulțimea finală a conceptelor a fost ordonată după numărul de limbi din EuroWordNet ce le lexicalizează și după numărul de limbi din BALKANET care le-au propus. Primele 4000 de concepte din această listă au fost de comun acord alese ca țintă comună pentru următoarea etapă a proiectului.

În continuare prezentăm metodologia folosită pentru limba română în selecția fondului lexical în cadrul BALKANET. Analiza cantitativă a fost realizată asupra unui corpus foarte mare, format din mai multe romane și dintr-un număr de texte jurnalistice culese de pe web. Corpusul (conținând mai mult de 10 milioane de cuvinte) a fost supus unor prelucrări statistice, fiind lematizat automat, iar cuvintele care prezentau interes (substanțive, verbe, adjective și adverbe) au fost sortate în funcție de frecvența lor. Am extras în acest fel o listă de mai mult de 30.000 de leme. În funcție de frecvența acestora în textele din corpus, această listă a fost împărțită în trei clase corespunzând celor mai frecvente 10.000 de leme (I), următoarele 10.000 de leme (II) și restul (III). Frecvența dintr-un corpus este o măsură foarte mult subiectivă. Printre cele mai puternice criterii pentru a măsura numărul și reprezentativitatea textelor incluse în corpusul folosit au fost cantitativă. Luând în calcul faptul că din ce în ce mai multe texte sunt culese pe web, mărimea corpusului nu mai reprezintă o problemă semnificativă. Reprezentativitatea rămâne în continuare un punct slab. Definirea unei mulțimi de texte care trebuie incluse într-o analiză cantitativă face obiectul unor discuții polemice și nu vom insista asupra ei. Având în vedere că datele noastre sunt aproape în întregime din texte jurnalistice, problema reprezentativității este în continuare ridicată. Dicționarul de Frecvențe al Cuvintelor Române (DFC) publicat cu mult timp în urmă, bazat pe un *corpus balansat* de 500.000 de cuvinte (teatru, nuvele și scurte povestiri, eseuri memorii și corespondențe, articole jurnalistice, literatură tehnică) conține cele mai frecvente 5.000 de cuvinte. Dacă este foarte controversat, FDLW este încă folosit de mulți lingviști ca referință. Comparația pe care am făcut-o a arătat că mai toată lista de leme inventariate de FDLW se găsesc și în lista obținută de noi, chiar dacă aceleași scoruri de frecvență. Pe lângă frecvența în corpus am avut în vedere două criterii mai puțin controversate și care au putut fi operaționalizate: numărul resurselor lingvistice disponibile și instrumentele noastre de analiză. Primul este numărul de sensuri pe care un cuvânt (împreună cu sinonimele și expresiile în care participă) îl are într-un dicționar. Al doilea este numărul de definiții de dicționar în care apare un anumit cuvânt. Al treilea criteriu este încă în analiză, ar putea fi numărul de derivate lexicale ale unui cuvânt. Pentru o pertinentă analiză din acest punct de vedere, o excelentă lucrare este

În această fază a proiectului BALKANET, ne-am concentrat atenția asupra substantivelor din limba română, iar datele experimentale raportate mai jos se referă doar la acestea. Având însă în vedere că procedurile tehnice nu depind de categoria gramaticală, metodologia și procedura vor fi aceleași și pentru verbe, adjective și adverbe. Luând în calcul numai primele două clase de frecvență descrise mai sus (primele 20.000 cele mai frecvente din corpusul jurnalistic) am extras din XML-LEX mai mult de 8.000 de intrări de substantive și substantive compuse (care însumează aproximativ 35.000 de sensuri) astfel încât productivitatea definițională PD (numărul de definiții în care participă un substantiv) să fie cel puțin 3. Lista a fost sortată în funcție de productivitatea definițională și numărul de sensuri ale fiecărui cuvânt titlu.

Substantiv	Productivitate definițională	Număr de sensuri	FRECV _{range}
acțiune	2279	13	I
persoană	1979	9	I
parte	1882	94	I
formă	1286	21	I
obiect	1204	16	I
fapt	1044	11	I
rasism	3	1	II

Figura 3: Ordonarea candidaților

Pentru toate aceste substantive am extras traduceri englezești din dicționarul de echivalenți de traducere. Procedurile pentru extragerea automată a echivalenților de traducere din corpusuri paralele ca și procedura de discriminare a sensurilor sunt descrise pe larg în [22], [23], [29], [15], [16]. Fiecare substantiv din limba română a fost pus în corespondență cu lista tuturor conceptelor din ILI corespunzătoare traducerilor sale în engleză. Conceptele astfel identificate au fost sortate după rangul corelat al substantivelor românești de la care s-a pornit.

Interesant de remarcat este că dintre cele 4000 de concepte selectate în final prin armonizarea propunerilor tuturor partenerilor, circa 2600 s-au regăsit și în primele 4000 de concepte ale ierarhiei noastre. Toate cele 4000 de concepte selectate de consorțiu se regăsesc printre primele 6000 de concepte ale ierarhiei noastre.

Toate substantivele reprezentând potențiale lexicalizări ale celor 4000 de concepte din cea de a doua selecție au fost automat puse în corespondență cu toate definițiile lor din XML-LEX. De asemenea, ele au fost corelate cu lexicalizările din limba engleză ale celor 4.000 de concepte. Prin intermediul dicționarului de echivalenți de traducere englez-român, fiecare concept a fost asociat cu

lexicalizarea din limba engleză (extrasă din WORDNET) și cu potențialele lexicalizări în limba română.

Dicționarul de Sinonime al Limbii Române (DSLRL), digitizat și codificat în XML, a fost folosit pentru a extrage seriile sinonimice pentru cuvintele românești selectate. În XML-DSLRL unii membri ai seriilor sinonimice sunt arhaisme sau regionalisme. Discuțiile preliminare au condus către ideea de a elimina toate cuvintele care fac parte din aceste clase (ne-am bazat pe cerința de a construi un nucleu lexical de uz general în limba română contemporană). Totuși, pentru eventualitatea în care aceste cuvinte filtrate (împreună cu informațiile despre uz) vor fi necesare mai târziu, s-a asigurat recuperabilitatea lor. Seriile sinonimice românești au fost considerate ca posibile sinseturi și adăugate la asociațiile descrise mai sus.

4. Instrumente software dezvoltate pentru proiectul BALKANET

Materialul lingvistic de bază descris în secțiunea anterioară, a fost asamblat prin intermediul unor programe unitare, astfel încât toată această informație este disponibilă într-o interfață "prietenosă", prin care lexicograful alege echivalențele corecte de sens dintre cele potențiale. Această interfață este generată și "personalizată" automat în funcție de mulțimea conceptelor interlinguale furnizată ca parametru de intrare unui generator de interfețe. Printr-un astfel de model arhitectural, a fost posibil ca sarcina construirii wordnet-ului pentru limba română să fie distribuită între membrii celor două colective românești participante la proiect și judicios controlată. Pentru fiecare dintre aceștia s-a generat o interfață personalizată pentru o submulțime distinctă de concepte dintr-una din cele agreate de consorțiu proiectului. Utilizatorul acestei interfețe, pe care generăm îl numim în continuare *lexicograf*, va lucra în mod independent de ceilalți construind, ca urmare a interacțiunii, fragmente ale wordnetului pentru limba română. La un moment dat, lexicograful alege un concept din mulțimea ce i-a fost repartizată căruia dorește să-i atașeze un sinset românesc. El are la dispoziție simultan sinsetul ce lexicalizează în limba engleză conceptul respectiv și, pentru fiecare cuvânt englezesc din acest sinset, toate potențialele lui traduceri în limba română, aceste traduceri având atașate toate definițiile conținute în XML-LEX. În plus, fiecare cuvânt românesc are atașate toate seriile sinonimice din XML-DSLRL în care el este prezent. Ceea ce trebuie să decidă lexicograful este (vezi figura 4)

- care este cuvântul românesc a cărui definiție este cea mai apropiată de definiția conceptului lexicalizat în limba engleză;
- care este cea mai bună serie sinonimică a acestui cuvânt;

- c. care dintre definițiile atașate cuvintelor dintr-o serie sinonimică este cea mai adecvată pentru a fi aplicabilă tuturor cuvintelor din seria respectiva.

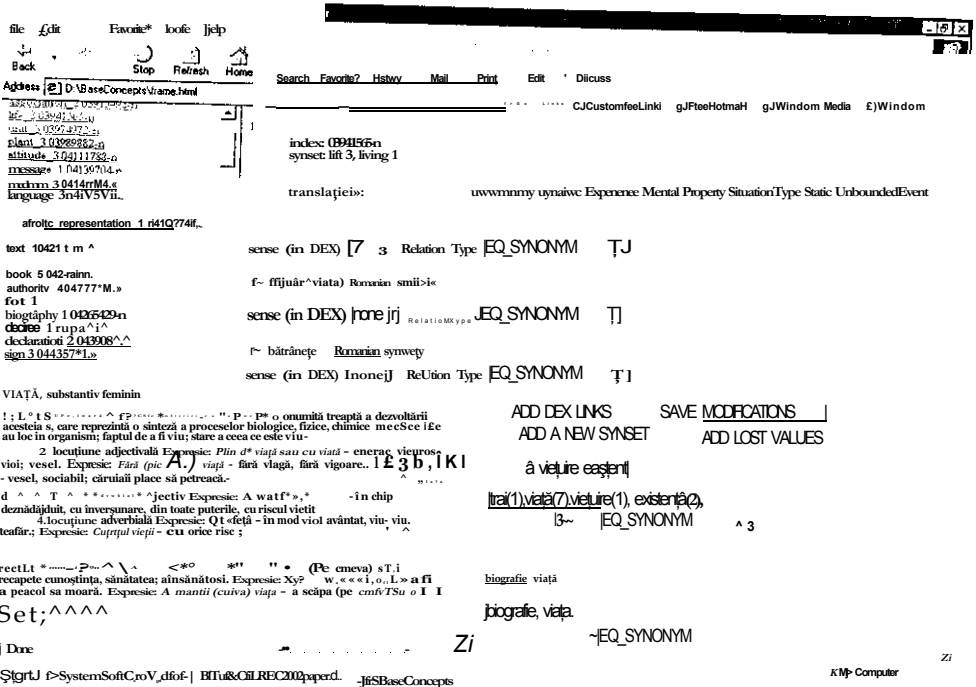


Figura 4: Editorul pentru construirea sinseturilor

În majoritatea cazurilor, definițiile extrase din XML-LEX corespunzând sinonimelor dintr-un sinset nu sunt identice, lexicografii alegând p S a mai apropiata de definiția conceptului corespunzător (vezi figura 5)

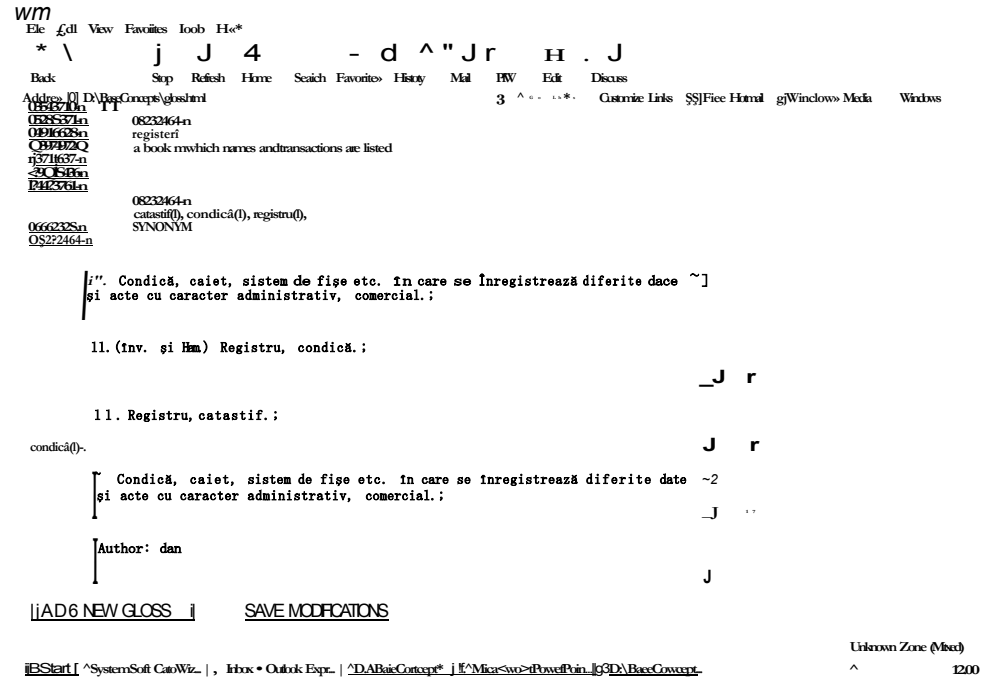


Figura 5: Editorul pentru asignarea gloselor

Merită menționat că în faza asocierii gloselor a devenit evidentă incorectitudinea alcătuirii unor sinseturi, ele fiind modificate. În alte cazuri Dicționarul Explicativ al Limbii Române include în aceeași definiție două sensuri care sunt demarcate în ILI ca două concepte diferite. În astfel de situații strategia generală a fost să se despartă definiția românească și să se atașeze ca glosă partea relevantă.

Fragmente create de fiecare lexicograf sunt agregate în mod incremental în structuri din ce în ce mai complexe și mai acoperitoare din punct de vedere lexical. Acest proces de agregare se realizează în mod centralizat, astfel încât corectitudinea structurilor rezultate să poată fi controlată și, în cazul conflictelor, să se poată identifica și corecta sursele de conflict (de exemplu: același sens pus în corespondență cu concepte diferite, sensuri diferite ale aceluiași cuvânt puse în corespondență cu același concept, literalii fără identificatori de sens etc). Corectarea unor conflicte între două porțiuni ale structurii agregate poate să genereze conflicte între alte părți ale sale. Pentru evitarea acestui pericol au fost proiectate mecanisme de control centralizat al unificării subseturilor de wordnet ce gestionează efectul global al oricăror modificări locale.

4.1. Importul relațiilor taxonomice; vizualizare sincronizată a mai multor wordneturi

Construcția sinseturilor și punerea lor în corespondență cu conceptele interlinguale reprezintă doar una din cele două dimensiuni fundamentale ale procesului de construire a unei rețele semantice lexicale pusă în corespondență cu indexul interlingual, respectiv cea de implementare a nodurilor și echivalarea acestora cu conceptele interlinguale. Cea de a doua dimensiune a procesului construcției rețelei o constituie definirea relațiilor (intrainguale) între nodurile create și echivalate în prima fază. Deosebit de importante sunt relațiile taxonomice care stabilesc o ierarhie de tip generic-specific între sinseturile unui wordnet.

Stabilirea relațiilor taxonomice între sinseturile wordnetului pentru limba română s-a făcut automat (urmată de validarea umană) în baza principiului "echivalenței ierarhice interlinguale" [30]. În esență, acest principiu afirmă că:

1. dacă sinsetul S_{1LA} din limba LA și sinsetul S_{1LB} din limba LB sunt echivalate cu același concept C_1 din ILI și
2. dacă sinsetul S_{2LA} din limba LA și sinsetul S_{2LB} din limba LB sunt echivalate cu același concept C_2 din ILI și
3. dacă în limba A sinseturile S_{1LA} și S_{2LA} sunt într-o relație ierarhică H' (H' denotă compunerea de un număr de ori cel puțin egal cu 1 a relației H , în cazul nostru: has-as-hypernym), atunci:

în limba B sinseturile S_{1LB} și S_{2LB} sunt într-o relație ierarhică similară H' (deși lanțurile de relații H pot fi de lungimi diferite în cele două limbi).

Principiul explicitează necesitatea ca interpretarea relațiilor folosite în ontologia multilingvă să fie similară, așadar definește coeziunea interpretativă a relațiilor ontologice în toate limbile participante la proiect. Acest principiu este reprezentat schematic în figura 6:

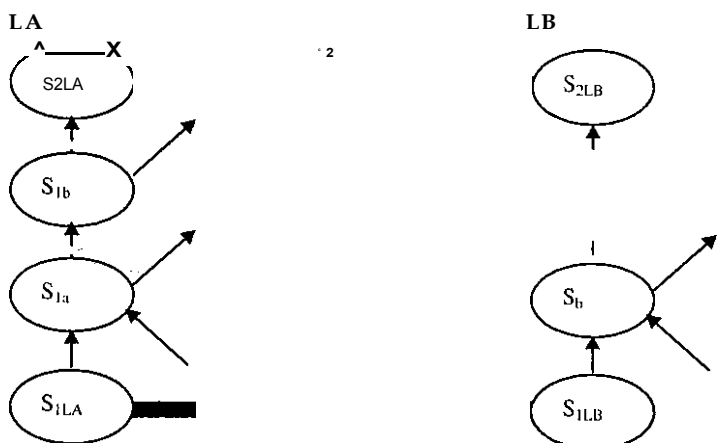


Figura 6: $(S_{1LA} \text{ EQ-SYN } S_{1LB}) \& (S_{2LA} \text{ EQ-SYN } S_{2LB}) \& (S_{1LA} \text{ H' } S_{2LA}) \quad (S_{1LB} \text{ H' } S_{2LB})$

În secțiunea următoare vom arăta pe un caz concret cum poate fi exploatat acest principiu pentru a importa (și eventual valida/corecta manual) relațiile dintr-un wordnet în care structurile ierarhice au fost stabilite, într-un wordnet pentru care au fost stabilite doar relațiile de echivalență translatională cu indexul interlingual (ILI).

Ultima etapă a construirii unui grup de sinseturi este transformarea rezultatelor interacțiunii lexicografului cu interfața descrisă anterior într-un format independent de limbă (codificare XML) și specific editorului multilingual de ontologii lexicale numit VisDic [25]. Odată generat acest format, el poate fi încărcat în VisDic, iar wordnetul pentru limba română poate fi vizualizat în mod sincron cu toate celelalte wordneturi încărcate. În figura de mai jos este ilustrată afișarea în mod sincron a sinsetului românesc (*ființă_1*, *formă de viață*, *viețuitoare_1*, *vietateji*) și a celui englezesc (*being_1* *life form*, *living thing*, *organism_1*) și a arborilor lor de hiponimi. Cele două sinseturi sunt aliniate via ILI, ambele fiind echivalate independent cu conceptul interlingual cu identificatorul 00002728-n.

Figura 7: Vizualizarea sincronizată a două ontologii lexicale cu ajutorul VisDic

Editorul de ontologii multilingve, VisDic, a fost dezvoltat în cadrul proiectului BALKANET pentru a substitui funcționalitatea asigurată în cadrul EuroWordNet de editorul Poiaris, dezvoltat de firma Lernout & Hauspie. Implementat inițial pentru ca rezultatele proiectului BALKANET să poată fi utilizate în regim liber de restricții comerciale (Poiaris poate fi utilizat doar contra cost), VisDic este constant îmbunătățit cu facilități noi a căror necesitate apare pe măsura evoluției proiectului BALKANET, fiind deja unul dintre cele mai puternice instrumente existente pentru gestiunea ontologiilor multilinguale.

5. Principiul conservării trans-linguale a ierarhiei lexicale. Studiu de caz: Condimente, mirodenii, sosuri și alte ingrediente

Vom considera fragmentele din RO-WordNet și WordNet 1.5 arătate în figura 8. Săgețile reprezintă relațiile taxonomice (de la hiponime spre hipernime) în cele două wordneturi. Liniile groase reprezintă relațiile de echivalență de traducere (EQ-SYN) dintre sinseturile celor două limbi, aceasta însemnând că sinseturile respective sunt puse în corespondență cu același concept din ILI. Linia groasă întreruptă reprezintă o relație EQ-SYN identificată ca nerespectând principiul conservării trans-linguale a ierarhiilor lexicale din cele două wordneturi. Inconsistența este semnalată deoarece în română relațiile ierarhice (de hiponimie) dintre *mirodenie*(RO) și *condiment*(RO) ca și dintre *ketchup*(RO) și *sos*(RO) nu sunt verificate de echivalenții lor în limba engleză: *sp/ce*(EN) este frate cu *condiment*(EN) și respectiv *ketchup*(EN) este frate cu *sauce*(EN). Dacă structura variantei 1.5 a WordNet este considerată cea corectă, acest exemplu arată că principiul păstrării ierarhiei nu este irefutabil. Pe de altă parte, dacă ar fi rezonabil să considerăm că WN 1.5 este amendabil (de exemplu făcând *mustard*(EH) și *ketchup*(EH) hiponimii direcți ai lui *sauce*(en)) ca în figura 9, atunci principiul păstrării ierarhiei ar putea fi o puternică probă a consistenței!

În urma restructurării ierarhice și de echivalare translațională, necesare pentru respectarea principiului conservării trans-linguale a ierarhiei lexicale (arătate în figura 9), interesant este faptul că a dispărut relația de echivalență între cuvântul românesc *condiment* și cuvântul englezesc *condiment*.

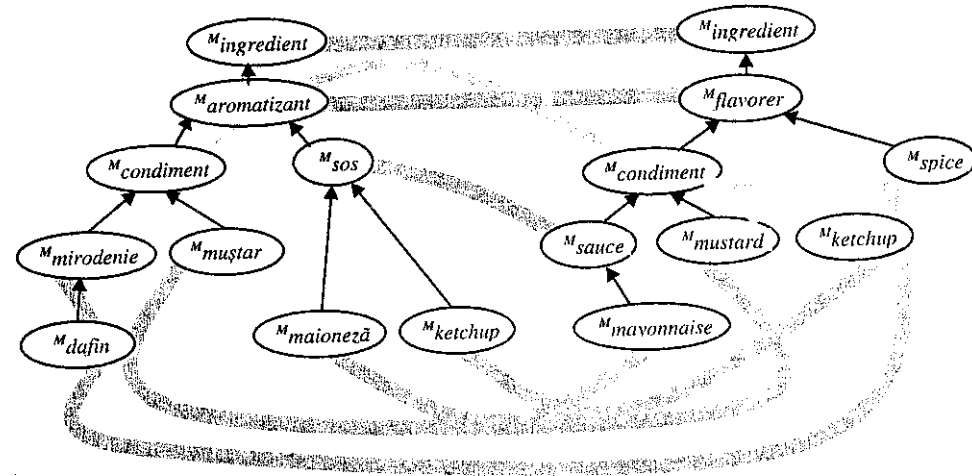


Figura 8: Nerespectarea principiului conservării trans-linguale a ierarhiei lexicale

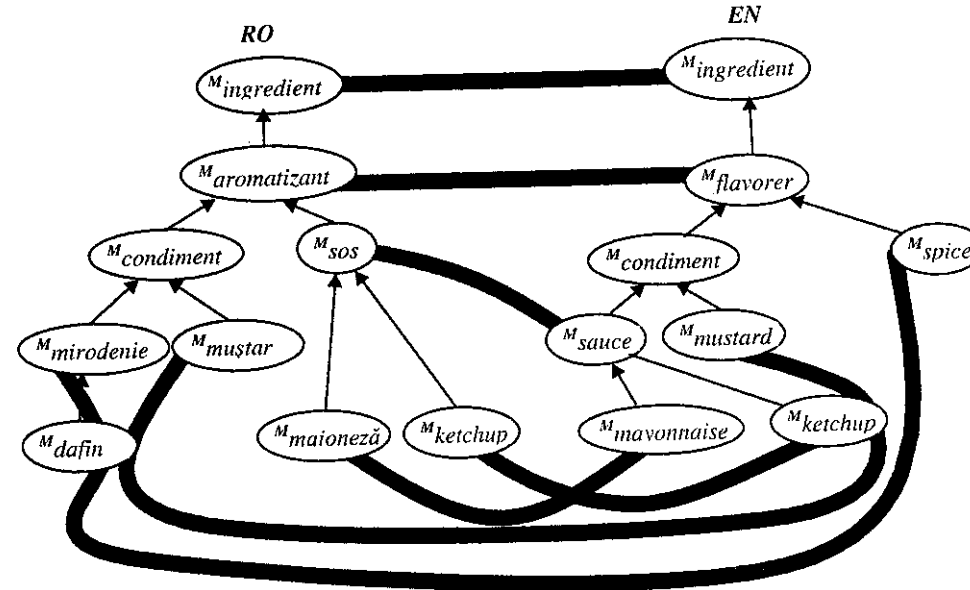


Figura 9: Reprezentare conformă cu principiul conservării trans-linguale a ierarhiei lexicale

Consultată recent asupra acestei probleme, Christiane Felbaum a confirmat existența unei erori în ierarhia WN1.5, probată, de altfel, și de glosa /u/ketchup (thick spicy sauce made from tomatoes).

Pentru ca această echivalență să fie posibilă, în condițiile principiului conservării trans-linguale a ierarhiei lexicale, ar trebui ori ca în limba engleză *spice* să fie un hiponim al lui *condiment* iar *sauce* să nu fie un hiponim al lui *condiment* ci frate, ori în limba română *sos* să fie un hiponim al lui *condiment* iar *mirodenie* să nu fie un hiponim al lui *condiment* ci frate. Ambele variante au fost respinse de experții consultați, lexicografi și vorbitori nativi ai limbii engleze și respectiv române. Singura concluzie posibilă este că în română și engleză cuvântul *condiment* nu reprezintă exact același lucru.

6. Concluzii

Realizarea ontologiei lexicale pentru limba română, în contextul multilingual definit de proiecte de tipul EuroWordNet, Balkanet și GlobalWordnet (www.globalwordnet.org), este esențială pentru procesul de informatizare a limbii române. Experiența internațională arată că un astfel de proiect nu este niciodată închis, reclamând actualizare și întreținere continuă, apărând mereu noi idei de îmbunătățire a performanțelor și noi cerințe de exploatare. Specialiștii de la Princeton au anunțat deja versiunea 1.7.1 a Wordnet, mult îmbunătățită. În variantele ce vor urma, pe lângă extensia în continuare a fondului lexical, toate cuvintele non-funcționale apărând în definiții vor conține referințe spre sinsetul corespunzător contextului de utilizare. Cu alte cuvinte, Wordnet va deveni simultan și un dicționar și un corpus adnotat la nivelul sensului. O altă dezvoltare semnificativă o va reprezenta traducerea definițiilor din Wordnet în formule logice, adecvate prelucrărilor inferențiale. Acest proiect, coordonat de Dan Moldovan și Sanda Harabagiu se află în derulare la Universitatea Texas din Dallas [31], [32].

Astfel de extensii vor trebui considerate în viitor și în wordnetul pentru limba română aflat deocamdată în fază incipientă. Obiectivul final prevăzut pentru cei trei ani de derulare ai proiectului BALKANET (septembrie 2004) este realizarea unui nucleu de câte 8.000 de sinseturi în fiecare din limbile proiectului.

În acest moment, la mai puțin de un an de la începerea proiectului, wordnetul românesc se află cu mult înaintea graficului prevăzut, având deja create peste 6.000 de sinseturi. Se poate estima că, în condiții normale, în cei peste doi ani care au mai rămas wordnetul românesc va ajunge la peste 20.000 de sinseturi, acoperind peste 40.000 de literali. Atingerea unui volum lexical similar cu al altor wordneturi necesită însă continuarea proiectului și după anul 2004, atragerea unor noi colective de specialiști în această întreprindere și desigur găsirea surselor de finanțare, în principal interne, care să permită dezvoltarea și întreținerea wordnetului românesc. Operaționalizarea acestui obiectiv poate fi facilitată de contextul organizatoric creat de curând prin înființarea la Academia Română a Comisiei de Informatizare pentru Limba Română (CILR) precum și a Consorțiului

de Informatizare pentru Limba Română (ConsILR: <http://www.consilr.info.uaic.ro/>), for executiv al CILR.

A fost construită o platformă software de dezvoltare incrementală a rețelei semantice ce permite implementarea independentă de regiuni ale rețelei și integrarea ulterioară a acestora. Viabilitatea acestui concept arhitectural și a demersului de dezvoltare distribuită a wordnetului au fost validate prin implicarea în procesul de construire a 10 specialiști, cărora li s-au adăugat încă 12 studenți masteranzi de la Facultatea de Litere a Universității București și Facultatea de Informatică a Universității "A.I. Cuza" (cele două facultăți ce au programe de Mașter în domeniul prelucrării limbajului natural și al lingvisticii computaționale). Rezultatele produse în mod independent au fost agregate fără nici o dificultate. Mediul lingware de dezvoltare conține un modul special de verificare a corectitudinii deciziilor lingvistice la crearea sinseturilor românești sau la punerea lor în corespondență cu conceptele indexului interlingual. După cum era de așteptat, procesul de integrare a rezultatelor parțiale furnizate de fiecare membru al celor două echipe de realizare a evidențiat o serie de inconsistențe cu explicații diverse:

- neatenție în asignarea sensurilor, generată de oboseala expertului decident uman;
- granularitate semantică diferită între sensurile explicitate în XML-LEX și sensurile conceptelor din ILI;
- absența lexicalizării în limba română a unor concepte existente în ILI și introducerea unor forme perifrastice cu definiții ad-hoc;
- erori sau incompletitudini existente în sursele lingvistice primare folosite în implementare.

Inconsistențele depistate, atât de natură structurală, dar mai ales cele de natură semantică au fost înregistrate, analizate și unele dintre ele corectate. Altele, necesită o analiză mai profundă și rezolvarea lor a fost amânată pentru o etapă ulterioară a proiectului. Aceasta cu atât mai mult cu cât, prin analiza similară pe care am efectuat-o asupra wordneturilor pentru celelalte limbi din proiect, am constatat că există multe similități ale acestor genuri de inconsistențe. Sunt puse astfel în evidență o serie de concepte din ILI pentru care diferența semantică dintre lexicalizările lor este prea mică pentru a fi sesizată ușor chiar și de către un vorbitor nativ al limbii respective. Distincții atât de rafinate au, din perspectiva prelucrării automate și mai ales a traducerii automate, o utilitate limitată iar în context multilingv pot fi chiar surse de eroare. Pericolul micșorării distanței semantice (am putea numi acest fenomen pulverizarea conceptuală) între conceptele din ILI este amplificat de adăugarea unor concepte ce au lexicalizări într-o singură limbă sau într-un număr mic de limbi. O soluție pentru evitarea idiosincraziilor lexicale într-un context multilingv și a disparităților de traducere este

gruparea conceptelor foarte apropiate semantic în ceea ce s-ar putea numi *concepte agregat* Lexicalizările înțelesurilor din două sau mai multe limbi, puse în corespondență cu aceleași concepte din ILI sau cu concepte membre ale unui agregat, vor putea fi folosite ca echivalenți de traducere în pofida unor diferențieri semantice specifice unei limbi sau alteia (*ciorbă, sarmale, pepper pot, porcupine ball* etc; vezi și exemplele din secțiunea precedentă). Analiza inconsistentelor interumane în echivalarea înțelesurilor dintr-o limbă cu conceptele interlinguale din ILI, precum și identificarea conceptelor distincte puse în corespondență cu echivalenți de traducere (extrași automat din corpusuri paralele sau găsiți într-un dicționar bilingv clasic) pot furniza informații calitative mult mai interesante (cel puțin din perspectiva psiho-lingvisticii) și mai demne de încredere decât o analiză statistică. Aceasta este o promițătoare direcție de cercetare ce se dezvoltă în paralel cu activitatea principală de construcție a wordnetului pentru limba română.

Referințe bibliografice

- [1] Danzin, A. - Towards a European Language Infrastructure" raport al Comisiei Europene, 1992.
- [2] Fellbaum, Ch. (ed.) - *WordNet: An Electronic Lexical Database*, MIT Press, 1998, 423 p.
- [3] Bloksma, L., Diez-Orzas and Vossen, P. - The User Requirements and Funcțional Specification of the EuroWordNet-project *EWN-deliverable D.001*, LE-4003, 1996.
- [4] Vossen, P. (ed.) - "A Multilingual Database with Lexical Semantic Networks", Kluwer Academic Publishers, Dordrecht, 1998.
- [5] Stamou, S., Oflazer, K., Pala, K., Christoudoulakis, D., Cristea, D., Tufiș, D., Koeva S., Totkov, G., Dutoit, D., Grigoriadou, M. - BALKANET A Multilingual Semantic Network for the Balkan Languages, in *Proceedings of the International Wordnet Conference*, Mysore, India, 1997.
- [6] Tufiș, D. - "Promovarea Limbii Române în SI-SC", în "Societatea Informațională - Societatea cunoașterii: concepte, soluții și strategii pentru România", Florin Gh. Filip (coord.), Ed. Expert, București, ISBN973-8177-42-1, 2001, pp. 131-142.
- [7] Erjavec, T., Ide, N., Tufiș, D. - Encoding and Parallel Alignment of Linguistic Corpora in Six Central and Eastern European Languages" in Michael Levison (ed) *Proceedings of the Joint ACH/ALL Conference* Queen's University, Kingston, Ontario (also on <http://www.qucis.queensu.ca/achallc97>), 1997.
- [8] Dimitrova, L., Erjavec, T., Ide, N., Kaalep, H.J., Petkevic, V., Tufiș, D. - Multext-East: Parallel and Comparable Corpora and Lexicons for Six Central and Eastern European Languages. *Proceedings of COLING*, Montreal, Canada,
- [9] Tufiș, D., Bruda, Șt. - Structure Markup in CES and Preliminary Statistics on Romanian Translation of Plato's "Republica", *Proceedings of International Seminar on Encoding*, Ljubliana, February, 1997, also in *TELRI News*, nr. 5.
- [10] Tufiș, D., Barbu, A.M., Pătrașcu, V., Rotariu, G., Popescu, C. - Corpora and Corpus-Based Morpho-Lexical Processing, in Tufiș D., P. Andersen (eds.) *Recent Advances in Romanian Language Technology*, Editura Academiei, 1997, pp. 115-128.
- [11] Tufiș, D., Rotariu, G., Barbu, A.M. - TEI-Encoding of a Core Explanatory Dictionary of Romanian. In Kiefer, F. and Pajzs J. (eds.) *Papers in Computational Lexicography*, Hungarian Academy of Sciences, 1999, pp. 219-228.
- [12] Ide, N. - *Corpus Encoding Standard: SGML Guidelines for Encoding Linguistic Corpora* First International Language Resources and Evaluation Conference, Granada, Spain, 1998, See also <http://www.cs.vassar.edu/CES/>.
- [13] Tufiș, D. - Tiered Tagging and Combined Classifiers in F. Jelinek, E. Noth (eds) *Text, Speech and Dialogue*, Lecture Notes in Artificial Intelligence 1692, Springer, 1999.
- [14] Tufiș, D., Popescu, C, Roșu, R - Automatic classification of documents by random sampling in *Proceeding of the Romanian Academy*, Series A, voi. 1, no. 2, 2000, pp. 18-28.
- [15] Tufiș, D. - "A cheap and fast way to build useful translation lexicons" in *Proceedings of the 19th International Conference on Computational Linguistics, COLING2002*, Taipei, 25-30 August, 2002, pp. 1030-1036.
- [16] Ide, N., Erjavec, T., Tufiș, D. - "Sense Discrimination with Parallel Corpora" in *Proceedings of the SIGLEX Workshop on Word Sense Disambiguation: Recent Successes and Future Directions. ACL2002*, Philadelphia, 2002, pp. 54-60.
- [17] DEX - Coteanu, I, Seche, L., Seche, M. (coord.). *Dicționarul Explicativ al Limbii Române*, ediția a II-a, *Univers Enciclopedic*, 1996, București.
- [18] Tufiș, D. - Blurring the distinction between machine readable dictionaries and lexical databases. *Research Report, RACAI-RR56, 2000, p. 56.*
- [19] Vintilă-Rădulescu, I. - "Resurse lingvistice pentru limba română elaborate la Institutul de Lingvistică «Iorgu Iordan»", în acest volum, 2002.
- [20] Erjavec, T., Evans, R., Ide, N., Kilgariff, A. - The CONCEDE Model for Lexical Databases. *Proceedings of the Second Language Resources and Evaluation Conference (LREC)*, Athens, Greece, 2000, pp. 355-362.
- [21] Seche, L., Seche, M. - *Dicționarul de sinonime al limbii române*. Univers Enciclopedic, București, 1997.
- [22] Tufiș, D., Barbu, A.M. - *Computational Bilingual Lexicography: Automatic Extraction of Translation Dictionaries*, in *International Journal on Science and*

Technology of Information, Romanian Academy, ISSN 1453-8245, Voi. 4, No. 3-4, 2001, pp. 325-352.

- [23] Tufiş, D., Barbu, A.M. - *Extracting multilingual lexicons from parallel corpora*, in Proceedings of the ACH-ALLC conference, New York, 12-17 June, 2001, 4 p.
- [24] Tufiş, D., Barbu, A.M. - "Lexical token alignment: experiments, results and applications" In Proceedings of LREC2002, Las Palmas, Spain, 2002, pp. 458-465.
- [25] Pavelek, T., Pala, K. - *VisDic: A new Tool for WordNet Editing* in Proceedings of the 1st International Wordnet Conference, Mysore, 2002.
- [26] Rodriguez, H., Climent, S., Vossen, P., Bloksma, L., Peters, W., Alonge, A., Bertagna, F., Roventini, A. - The Top-Down Strategy for Building EuroWordNet: Vocabulary Coverage, Base Concepts and Top Ontology. In Piek Vossen (ed.) *EuroWordNet: A Multilingual database with lexical semantic networks*, Computers and Humanities, Voi. 32, Nos. 2-3, 1998.
- [27] Juilland, A., Edwards, P.M.G., Juilland, I. - The Frequency Dictionary of Rumanian Words. *Mouton & CC*, London-The Hague-Paris, 1965.
- [28] Dinu, M. - Personalitatea limbii române, Editura ALL, 1996, 368 p.
- [29] Erjavec, T., Ide, N., Tufiş, D. - *Automatic Sense Tagging Using Parallel Corpora*, in Proceedings of the 6th Natural Language Processing Pacific Rim Symposium, Tokyo, Japan, 27-29 November, 2001, pp. 212-219.
- [30] Tufiş, D., Cristea, D. - "Methodological issues in building the Romanian Wordnet and consistency checks in Balkanet", In Proceedings of LREC2002, Las Palmas, Spain, May, 2002, pp. 35-41.

[31] Moldovan, D. - "Question Answering Systems in Knowledge Management", IEEE Intelligent Systems, voi. 16, nr. 6, 2001, pp. 90-92.

[32] Harabagiu, S., Miller, G., Moldovan, D. - "WordNet 2 - A Morphologically and Semantically Enhanced Resource", in *Proceedings of SIGLEX-99*, Univ. of Maryland, 1999, pp. 1-8.

Algoritmi de segmentare a textului în unități de tip clauzal

D. GÂLEA,

Institutul de informatică Teoretică, Academia Română, Filiala Iași,
dgatea@iit.tuiasi.ro

N. CURTEANU,

Institutul de Informatică Teoretică, Academia Română, Filiala Iași,
curteanu@iit.tuiasi.ro

C. LINTEȘ,

Institutul de Informatică Teoretică, Academia Română, Filiala Iași

1, Introducere

Scopul lucrării este dublu: (a) Să prezinte și să compare doi algoritmi de segmentare a frazei (românești) în unități de tip clauzal. (b) Să întrezească și să susțină două componente de bază ale *strategiei lingvistice* SCD (Segmentare-Coeziune-Dependență) [1], [2] de analiză a limbajului natural (LN): procesul de segmentare a textului de LN, și *teoria* FX-bar [3], [23]. Segmentarea textului poate continua sau interfera cu stabilirea arborilor de dependență între unitățile clauzale și subclauzale (unități sintagmatice) ale textului. Unitățile de tip-clauzal corespund, în general, *relațiilor retorice* dintre unitățile minimale ale discursului, astfel încât algoritmi de segmentare pot fi (și chiar sunt) utilizați în aplicații ce țin de *teoria și procesarea discursului*. Primul algoritm este o aplicare la limba română a segmentării frazei în unități de tip-clauzal, algoritm dezvoltat de Daniel Marcu în [4], [5] (și prescurtat în cele se urmează "*algoritm-Marcu*", sau "*algoritm-M*"). Al doilea algoritm reprezintă o rafinare a segmentării în clauze și grupuri sintagmatice din cadrul strategiei lingvistice SCD [1], [2], [3], [6] (prescurtat în cele se urmează prin "*algoritm SCD*"). Acești algoritmi sunt implementați într-un mediu specializat de procesare (dezvoltat sub C++), și este realizată o comparație computațională a execuției segmentării de tip-clauzal pe un set consistent de fraze românești [7].

Segmentarea textului LN a devenit în ultimii ani un subiect intens cercetat și cu multiple aplicații. O atenție specială a primit segmentarea textului de LN în unități de discurs, în particular, segmentarea frazei în *unități minimale* de discurs de multe ori și pe bună dreptate asociate cu unități de tip clauzal în numeroase

teorii sintactice, semantice, și de discurs. Unitățile textuale de tip-clauzal obținute (sau proiectate) prin mijloace orientate sintactic, către analiza de suprafață [Eng: *shallow*], s-au dovedit a fi esențiale în numeroase tipuri de procesare a LN: parsare, traducere automată, generare de LN, interpretare de discurs, extragere de date lingvistice, regăsirea informației, rezumare automată, rezoluția anaforei etc. Un caz special de segmentare a textului este ceea ce numim '*chunking*', un proces dedicat obținerii unor tipuri de "*segmente*" [Eng: *chunks*] dominate de anumite categorii (verb, substantiv, adjectiv-adverb, clauză). În continuare, vom folosi doar termenul de *segmentare* a textului de LN, considerând *chunking-ul* drept un caz particular al procesului de segmentare a LN.

În analiza și implementarea celor doi algoritmi de segmentare, *algoritmul-M* și *algoritmul-SCD*, cel puțin două aspecte le considerăm a fi importante: **(a)** Se demonstrează că *algoritmul-M* de segmentare este scufundat în *algoritmul-SCD*, ceea ce înseamnă că primul dintre cei doi algoritmi poate fi obținut ca un caz particular al claselor de marcheri, ierarhiei acestor clase, și a segmentării (dependențelor) obținute de cel de-al doilea algoritm, **(b)** *Algoritmul-SCD* de segmentare poate fi conceput ca un bun punct de start în *proiectarea unui cadru general* pentru *algoritmii de segmentare* a textului de LN. Un asemenea cadru ar fi compus din: (b1) mai multe *sisteme de transformare* aplicate în cascadă, fiecare sistem component fiind format din *seturi și subseturi specifice de etichete*, (b2) o *ierarhie* stabilită între câteva dintre cele mai importante *clase* ale acestor *etichete*, și (b3) o *gramatică formală* (sau un automat finit) pentru recunoașterea (sub)secvențelor și arborilor de etichete (în concordanță cu ierarhia claselor de etichete). În abordarea prezentă, aceste aspecte sunt exemplificate de către o implementare C++ a celor doi algoritmi într-un mediu specializat, o bază de date a marcherilor (de discurs) românești, și o ierarhie specifică a claselor de marcheri lingvistici. Cei doi algoritmi de segmentare considerați sunt executați și comparați pentru un set consistent de fraze românești. Posibile dezvoltări și aplicații sunt menționate în [21], [22].

Importanța segmentării de tip clauzal a frazei în procesul de parsare a textului a fost scoasă în evidență încă de la începutul anilor '80, iar studii teoretice datează mult mai devreme. În România, primele lucrări științifice și contracte de lingvistică computațională au conținut, printre alte realizări meritorii, și primele încercări de realizare a segmentării automate a frazei în clauze finite (și non-finite) [8], [9], [10], [11]. În pofida unor mijloace formale neadecvate (gramatici formale) și de programare (rețele ATN) disponibile la acel timp, ideile principale pe care se bazează abordările menționate nu numai că au reprezentat premiere pentru acele timpuri, dar multe din ideile de atunci își păstrează încă o surprinzătoare actualitate, aceste fenomene de *come-back* ciclic fiind frecvente (și perfect explicabile de evoluția tehnologică) în momentul de față. Trebuie menționate aici folosirea intensivă a *marcherilor de discurs* (*cue phrases*, *connectives*), întâlnită și

în [4], [5], [12], [13], a *predicativității* (aparitia categoriilor *deverbale*) [14], [15], a utilizării automatelor finite în analiza LN etc.

De fapt, o versiune a gramaticii formale preluată din [8] este folosită în *Pasul 6* al *algoritmului SCD-2002* de segmentare, în concatenarea marcherilor de nivel M3 și M2 (vezi Secțiunea 4), în timp ce rudimente ale unor reguli similare din aceeași gramatică se regăsesc în algoritmul-M de segmentare, la compunerea acțiunilor care lucrează cu apariția multiplă a marcherilor (de discurs) [4], [5] (vezi Secțiunea 3).

2. Segmentarea de tip clauzal cu algoritmul M-1997

Prescurtat în continuare ca "*algoritmul de segmentare M-1997*", sau simplu "*algoritmul M-1997*", algoritmul de *segmentare-Marcu* a frazei în unități de tip-clauzal [4], [5] funcționează ca un automat finit, sau ca o rețea de tranziție, bazat pe un set de *stări și acțiuni*. În [4] se face o analiză de corpus a potențialilor *marcheri de discurs*, numiți și "*sintagme indicatoare*" [Eng: *cive phrases*] și "*conective*", cu scopul de a evalua contribuția potențială a diferiților marcheri la determinarea (delimitarea) unităților textuale elementare pe care sunt definite *relațiile retorice*, în cadrul unității textuale standard care este fraza. În încercarea de a stabili principalele tipuri de funcții ale marcherilor, și anume de tip clauzal, frazai, de discurs, sau pragmatic, algoritmul de segmentare M-1997 consideră mai întâi următoarele *trei clase* de marcheri:

- (Mari) În *prima clasă* sunt cuprinși marcherii (sintagmele indicatoare) care joacă un rol în cadrul discursului pentru majoritatea fragmentelor de text ale corpusului analizat. Elementele din (Mari) vor fi numite în cele ce urmează "*marcheri de discurs*", iar specifici acestei prime clase sunt marcheri ca "*deși*" [Eng: *although*], "*pe lângă*" [Eng: *besides*], "*dacă*" [Eng: *if*], "*atunci*" [Eng: *then*] etc.
- (Mar2) Marcherii din a *doua clasă*, numiți "*marcherii de frază/clauză*", joacă în discurs, pentru majoritatea fragmentelor de text în care apar, rolul de *adiacenți* la alți marcheri de discurs sau clauzali. Un membru specific al clasei (Mar2) este considerat a fi "*și*" [Eng: *and*], deoarece are rol clauzal de fiecare dată când apare înaintea altui marcher de discurs sau clauzal, cu toate că poate avea atât rol de discurs cât și clauzal atunci când apare izolat.
- (Mar3) A *treia clasă* conține marcheri care s-au dovedit că joacă un rol de *delimitare a clauzelor* în majoritatea fragmentelor de text investigate în [4]; ei vor fi referiți, simplu, ca "*marcheri clauzali*" (sau "de clauză"). (Mar3) include, de asemenea, acei marcheri pentru care analiza de corpus nu a putut distinge între funcția lor de discurs și cea clauzală. "*După*" [Eng: *after*] este un astfel de element reprezentativ al (Mar3).

Marcu [4] a selectat mai mult de 450 de *marcheri* (pentru engleză) în cadrul analizei sale de corpus pentru marcherii de discurs și de frază/clauză. Marcherii sunt stocați și procesați într-o *bază de date* ale cărei înregistrări conțin următoarele câmpuri:

- a. Câmpul denumit *Example* conține un *fragment de text* din care a fost extras marcherul.
- b. Câmpul *Marker* codifică *marcherul* însuși, împreună cu marcherii de punctuație contextuali și, atunci când este necesar, ceilalți *marcheri adiacenți*.
- c. Câmpul *Usage* furnizează unul sau mai multe dintre *rolurile funcționale* ale marcherului:
 - (d) *Frază/clauzal (S)*, atunci când marcherul nu îndeplinește *nici o funcție* în structurarea discursului;
 - (c2) *De discurs (D)*, când marcherul evidențiază o *relație de discurs* între două unități textuale;
 - (c3) *Pragmatic (P)*, dacă există o *relație* între o construcție lingvistică (sau non-lingvistică) care conține marcherul, și convingerile, planurile, intențiile și/sau scopurile de comunicare ale vorbitorului.
- d. Câmpul *Break_action* (acțiune de oprire) conține un nume de *acțiune* din mulțimea acțiunilor ce vor fi executate în cadrul procesului de segmentare. Acest proces este controlat de către un set de semnalizatori (*flaguri*). Execuția unei acțiuni din mulțimea {`NOTHING`, `NORMAL`, `COMMA`, `NORMALJ`"HEN_c6mMA", `END`, `MATCH_PAREN`, `COMMA_PAREN`, `MATCH_DASH`, `SET_AND`, `SET_OR`, `DUAL`} are unul dintre următoarele efecte:
 - (d1) creează o *margină* pentru unitatea textuală elementară în *string-ul* de intrare;
 - (d2) setează un semnalizator (*flag*).
- e. Câmpul *Position* specifică poziția marcherului de discurs în cadrul unității textuale căreia îi aparține. Valorile acestui câmp sunt B, M și E, după cum marcherul este situat *la început* (B), *în mijlocul* (M) sau, respectiv, *la sfârșitul* (E) unității textuale.

3. Algoritm de segmentare M-1997

Algoritm M-1997 primește în intrare o frază S și masivul *markers[n]* al marcherilor potențiali de discurs și clauzali din fraza S. Masivul *markers[n]* conține marcherii recunoscuți în S. Fiecare element al acestui masiv este caracterizat de către următoarea *structură de trăsături*:

- *Acțiunea* asociată aceluia marcher;

- *Poziția* marcherului în cadrul unității textuale elementare (B, M sau E);
- *Semnalizatorul hasjdiscoursejunction* care inițial este setat la valoarea "no".

Câteva dintre variabilele importante cu care lucrează algoritmul M-1997 sunt: *"status"*, *"parentheticaC"* și *"clauses"*.

Algoritm M-1997 pentru identificarea unităților de tip-clauzal din cadrul unei fraze are *două* părți principale:

(1) Când variabila *"status"* este NIL, algoritmul M-1997 execută acțiuni care pot introduce margini ale unității textuale sau pot modifica variabila, influențând procesarea marcherilor ulteriori. Pentru partea (1) a algoritmului M-1997, atunci când variabila *"status"* la valoarea NIL, sunt considerate următoarele situații:

- (1a) Dacă tipul de marcher este DUAL, determinarea marginilor unității textuale depinde de marcherul adiacent care precede marcherul curent analizat. În această situație, algoritmul M-1997 setează variabila *"status"* la aceeași valoare ca și în cazul unui marcher de tip COMMA.
- (1b) Dacă marcherul analizat curent nu este adiacent cu marcherul imediat precedent, atunci este identificată o margine a unității textuale.
- (1c) Cel mai frecvent tip de marcher (și de acțiune) este NORMAL, marcher care identifică o nouă unitate de tip clauzal a cărei margine dreapta este dată de marcherul curent analizat.
- (1d) Când marcherul de tip COMMA este precedat de un marcher de discurs, *sau*
- (1e) Tipul marcherului este NORMAL__THEN_COMMA, *atunci* algoritmul M-1997 identifică o nouă unitate de tip-clauzal ca și în cazul marcherului de tip NORMAL.
În oricare dintre cazurile (1c), (1d), (1e), variabila *"status"* este actualizată astfel încât o margine a unității textuale să fie identificată la prima apariție a unei virgule (COMMA).
- (1f) Pentru marcherul de tip NOTHING, singura acțiune constă în a atribui marcherului o utilizare specifică discursului.
- (1g) Marcherii care introduc posibile apariții de unități textuale parantetice (texte între paranteze) au doar efectul de a actualiza variabila *"status"*, ca și în cazul apariției marcherilor "sf și *"sau"*.

(2) Atunci când variabila *"status"* nu este NIL, algoritmul M-1997 execută acțiuni specifice pentru a realiza:

- (2a) Tratatrea informației din paranteze. O dată identificată o paranteză deschisă, o linie-de-despărțire [Eng: *dash*] (între două asemenea linii se introduce de obicei o apozitie sau un text explicativ), sau un



marcher de discours a cărei acțiune asociată este COMMA_PAREN, algoritmul M-1997 caută prima paranteză închisă, linie-de-despărțire, sau virgulă, ignorând toți ceilalți marcheri întâlniți pe parcurs. Acest tratament atrage după sine faptul că informației parantetizate *nu* îi este atribuită nici o stare pentru unitățile textuale elementare. Totuși, algoritmul M-1997 evită stabilirea de margini parantetizate în cazurile în care prima virgulă care urmează după un marcher COMMA_PAREN este imediat urmată de un marcher "și" ori "sau". De menționat este, de asemenea, că tratamentul aplicat informației dintre paranteze în algoritmul M-1997 poate conduce la rezultate eronate, ca în exemplul "*I-am dat lui Ion o rachetă de tenis, care i-a plăcut și o minge de plastic, care nu i-a plăcut*". Acest tip de erori poate fi evitat printr-o tratare mult mai adecvată în cadrul algoritmului de segmentare SCD.

(2b) Dacă variabila "status" conține acțiunea COMMA, apariția primei virgule care nu este adiacentă unui marcher "și" ori "sau" determină identificarea unei noi unități elementare de discours. Algoritmul M-1997 nu este, capabil, în general, să distingă suficient de precis între rolurile de discours și frazale/clauzale ale marcherilor "și" și "sau". Anumite situații sunt totuși recunoscute ca introducând funcții de discours, ca de exemplu apariția unui marcher de discours imediat după un "și" ori "o" în care valoarea semnalizatorului *has_discourse_function* este stabilită la "yes".

Forma originală a algoritmului M-1997 [4], [5] este extinsă și îmbunătățită în implementarea noastră pentru limba română (subsecțiunea 5.3) cu o analiză mai detaliată la nivelul ei superior, pentru apariții multiple și corelate ale marcherilor de discours/clauză [7].

4. Algoritmul de segmentare SCD-2002

Această secțiune prezintă partea de segmentare și dependență, în principal la nivel de clauză, desprinsă din strategia lingvistică SCD (*Segmentare-Coeziune-Dependență*) [1], [2], [3], [6], [23]. Forma actuală a algoritmului, referită în restul articolului prin prescurtarea SCD-1994, este foarte apropiată de versiunea publicată în [1], [2]. Noutatea principală a algoritmului SCD-2002 față de SCD-1994 constă într-o rafinare a claselor de marcheri, o nouă ierarhie a acestora, și în noul algoritm de stabilire a segmentării și dependenței (structurării) clauzelor și grupurilor sintagmatice. Vom pune în evidență relația dintre algoritmul M-1997 și algoritmi SCD-1994 și SCD-2002, arătând că primul este scufundat în ceilalți doi [7].

Rezultatele obținute prin execuția algoritmilor de segmentare M-1997 și SCD-2002 pe aceleași fraze conduc la aceeași concluzie: SCD-2002 are o granularitate (mult) mai fină a claselor de marcheri în comparație cu cea a claselor algoritmului M-1997, iar rafinarea acțiunilor implicate în SCD-2002 conduce la

delimitarea de unități textuale de tip-clauzal mai precise (de fapt mai corecte și mai adecvate) decât cele obținute de către algoritmul M-1997, prețul computațional ce trebuie plătit pentru acest fapt rămânând să fie analizat.

Este de menționat că segmentarea clauzală practică de SCD-2002 este doar un aspect particular al segmentării textului, deoarece se obțin și alte "*bucăți*" mici de text dominate de nuclee semantice de tip N (Substantiv), V (Verb), A (Adjectiv-Adverb). Segmentarea rezultată din clasele de marcheri SCD-2002 se află într-o strânsă relație cu noua teorie *X-bar funcțională* (FX-bar) [3], [23], o altă componentă importantă a strategiei lingvistice generale SCD.

Din schema generală FX-bar propusă în [3] se detașează următoarele nivele de proiecție la nivel lexical și gramatical:

Tabelul 4.1.

Nivele de proiecție ale schemei FX-bar (vezi [3], [23])

Marcheri	Nivelul de Proiecție	Structura gramaticală	Exemple
trăsătura PRED sau EXIST (OBJECT)	nivel de lexicon; prin convenție, (BAR = -1)	[forma de dicționar a cuvântului; X = N, V, A, Pron,	<i>la ploua</i> ^ <i>conducere</i> (trăsătura ; PRED) { <i>clădire</i> (trăsătura EXIST, înțelesul obiectual) <i>clădire</i> (PRED, pentru înțelesul acțiunii) <u><i>creion</i></u> (EXIST)
MO-marcher reprezintă aplicarea inflexiunii	XO (BAR = 0)	forma lexicală (de text) a cuvântului; X=N, V, A, ...	<i>filouă</i>
M1-marcher se aplică nucleului XO M1(X0)=X1	X1 = CL0; (BAR=1) poate fi identificat și cu nivelul 0 de proiecție a clauzei, BAR-CL = 0	sintagme XG (X=N, V, A), i.e. grupuri nominale, verbale, adjectivale <u>adverbiale</u>	<i>orice steag alb</i> <i>ploua</i> <i>aleargă repede</i> <i>nu aleargă deloc</i> <i>foarte bine studiat</i> <i>Măria i-a dat un măr</i> <i>fiicei sale.</i> <i>O femeie dăruind un măr unui bărbat conține o</i> <u><i>clauză infinită.</i></u>
M2-marcher se aplică proiecției X1 M2(X1)=X2=CL1 M2 se aplică unei <u>singure clauze</u>	proiecția X2 = C BAR = 2 și BAR-CL = 1	relații de discours între clauze finite	<i>Dacă plouă atunci plec</i> <i>mai devreme și îmi iau și</i> <u><i>umbrela.</i></u>
M3(CL1, CL1)=CL2 marcheri de discours; M3 se aplică la două sau mai multe clauze	nivelul de proiecție X3 = CL2; BAR = 3 si BAR-CL = 2		

4.1. Clasele de marcheri pentru algoritmul SCD-2002

Pentru algoritmul de segmentare SCD-2002 propunem o anumită rafinare a claselor de marcheri și a ierarhiilor acestor clase din [1], [2], [7], schimbări ce constau în următorul set de marcheri, în concordanță cu Tabelul 4.1. de mai sus:

M3 = { marcheri (de discurs) inter-clauzali}.

Clasa de marcheri M3 este formată din *funcții* sau *relații* (atunci când marcherii sunt corelați), având ca argumente două sau mai multe clauze finite (unele dintre ele pot fi infinite). Acești marcheri sunt ceea ce [4], [5], precum și alte abordări numesc "*marcheri de discurs*", și se aplică proiecțiilor sintactice de nucleu $X2 = CL1$ (și de nivel $X3$), de tip clauzal în teoria FX-bar (vezi Tabelul 4.1.).

M3 poate fi partiționată în următoarele subclase (în ordinea descrescătoare a *priorității* de definire a relațiilor de dependență - vezi Fig. 4.1.1.):

M33 = {marcheri (de discurs) inter-clauzali care introduc o dependență (neambiguă) de *supra-ordonare strictă*}. *Supra-ordonarea strictă* înseamnă *ridicarea efectivă a (cel puțin) unui nivel* de dependență clauzală, și este reprezentată de marcheri precum "*atuncl*", "*altfef*" etc.

M32 = {marcheri (de discurs) inter-clauzali care introduc dependență de *supra-ordonare*, incluzând *semnele de punctuație* precum două puncte, punct-și-virgulă, paranteză închisă, linie-de-despărțire etc.}. *Supra-ordonarea* presupune *ridicarea* unuia sau mai multor nivele de dependență clauzală, sau rămânerea pe același nivel de dependență în cadrul unei dependențe de tip-coordonare. Exemple tipice de marcheri din clasa M32 sunt "*da!*", "*așadaf*", "*chiar*", "*lajeljde*", "*în_comparațieJcu*" etc.

M31 = {marcheri (de discurs) inter-clauzali care introduc unul sau mai multe nivele de dependență de *sub-ordonare*, incluzând *semne de punctuație* ca paranteza deschisă, linia-de-despărțire etc.} Aceasta este o clasă largă de marcheri de discurs formată din numeroase tipuri de relații între clauze: logice, sintactice, semantice, pragmatice etc.

Așa cum a fost menționat mai sus, fiecare dintre clasele M33, M32 și M31 poate, la rândul ei, să fie partiționată în subclase care conțin marcheri de tip relațional (exprimați prin corelație), ce stabilesc relații între clauze, sau ca funcții de clauze (cu cel puțin două argumente).

M2 = { marcheri care introduc o clauză (finită sau infinită), sau un grup sintagmatic al cărui nucleu semantic este una din categoriile sintactice predicative N, V, A }. Compusul sintactic (sau *grupul sintagmatic* în termenii [3]) XG , $X = N, V, A$, poate fi asimilat unei clauze degenerate, infinite (vezi Tabelul 4.1) în cazul $X = N, A$.

M2 este divizată în următoarele subclase (în ordinea descrescătoare a *priorității* de introducere a relațiilor de dependență):

M25 = {marcheri care introduc *clauza relativă*}.

Explicația constă în faptul că o clauză relativă reprezintă cea mai complexă unitate sintagmatică ce joacă rol de modificador, și care se aplică nucleului NG al clauzei relative:

M24 = {aparitia unui *grup verbal finit* (FVG) sau, echivalent, apariția valorii FINITE pentru trăsătura TENS atribuită unui verb, introducând deci o *clauză finită*}.

Întregul grup verbal poate moșteni valoarea trăsăturii FINITE dacă nucleul său V sau altă componentă importantă din VG poartă această valoare a trăsăturii TENS (de exemplu, auxiliarul din VG).

M23 = {aparitia unei sintagme *predicative* XG (sau $X1$), $X=V, N, A$, al cărei nucleu semantic este o *categorie predicativă*, purtând valoarea $PRED = ACT$ (posibil încă la nivel de lexicon), și introducând astfel o *clauză infinită*}.

Clasele de marcheri M24 și M23 introduc structuri de nivel- $X2$, și anume clauze finite sau infinite, formate dintr-o sintagmă $X1$ (sau grup XG , $X = N, V, A$) care reprezintă *nucleul semantic, finit* ($TENS = FINITE$) sau *predicativ* ($PRED = ACT$ ional), al structurii de nivel- $X2$, urmată de sateliți (argumente și/sau conjuncții) corespunzători de tip NG (inclusiv NG-uri prefixate de o prepoziție, deci clasicele sintagme PP). Unele dintre argumente, cum este cazul *subiectului gramatical*, pot preceda nucleul semantic de tip $X1$ al clauzei căreia îi aparțin [3]. Să mai precizăm că există o *ordine sistemică (canonică)* [18], [19], a sateliților, sau "*actanților*" (argumente și conjuncții) dintr-o clauză (finită sau infinită): $ACT(or)$, $PAT(ient)$, $ADDR(essee)$, $ORIG(ine)$, $LOC(ation)$ etc. Ordinea canonică este specifică fiecărui nivel NG, și se poate obține în urma unei cercetări statistice și lingvistice foarte atente [18].'

Putem găsi recent un *principiu de predicativitate* similar cu cel folosit în *strategia lingvistică SCD*, și aplicat la sintagmele nominale din limba italiană [14], sau la adjectivele "*deverbele*" [14], [16]. În timp ce predicativitatea verbelor este frecventă și naturală, trăsătura de *nepredicativitate* [17, p. 22] (de fapt nepredicativitate) a verbelor de tip *existențial* este și ea la fel de frecventă (formele lui "*a fi*"), valoarea lor FINITE, dublată sau nu de valoarea trăsăturii $PRED = ACT$, anunțând totuși apariția unei clauze finite.

M22 = {marcheri care introduc relații de tip-JOIN, adică *conjuncții* de tip V . "*sau*", "*lajel_ca(și)*", "*împreunăjcu*"}

M21 = { $COMMA$ (sau $VIRGULA$) }.

Clasele M22 și M21 cuprind marcheri cu un grad important de ambiguitate deoarece pot introduce orice structură de tip $X1$ (grupuri XG , $X = N, V, A$) sau $X2$ (clauze finite sau infinite).

M1 = { marcheri care delimitează (introduc) structuri XG }.

Conform strategiei SCD și teoriei FX-bar [3], [23], clasa de marcheri M1 constă în *marcheri de r?/Ve/-X1*, X = N, V, A, adică marcheri care se aplică construcțiilor sintactice de *nivel-X1* (denotat și XG, și numit X-grup). Aceste sintagme constau, de fapt, dintr-un *nucleu semantic* înconjurat de *modificatori* (adjective sau adverbe) și/sau *specificatori* (sau *cuantificatori*, unii generalizați, printre cuantificatori incluzându-se determinatorii, negația etc).

Așa cum există o ordine sistemică a sateliților unui nucleu semantic într-o clauză (sintagmă de nivel-X2), în mod similar există o "*ordine structurală*", dată de "*distanța*" modificatorilor, cuantificatorilor, prepozițiilor etc. față de nucleul XO, pentru constituenții unei sintagme de nivel-X1. Astfel, în limba română (franceză, engleză), cel mai "apropiat" față de nucleul XO trebuie să fie *modificatorul* (adjectivul sau adverbul), urmează apoi *cuantificatorul* (care ocupă locul modificatorului dacă acesta lipsește), apoi prepoziția (ad-poziția, în general) etc. De exemplu, nu este sintactic corectă sintagma *frumos orice copil* sau *orice frumos pe copil*. Nucleul XO înconjurat de modificatori și/sau specificatori (cuantificatori) poate fi marcat funcțional prin *pre-poziții* (în cazul grupului nominal NG din română, engleză, franceză), dar și prin *post-poziții* (în cazul NG sau VG din engleză sau germană). Marcarea clitic-funcțională (prin particule pre- sau *post-poziționale*) poate exprima *cazul* (pentru NG), sau *timpul*, *semantica* (pentru VG) etc. Principalele elemente componente ale unei structuri XG corespund și subclaselor de marcheri ai clasei M1.

M1 poate fi divizată în subclase de marcheri, subclase utile în delimitarea substructurilor XG (X1), X = N, V, A, în conformitate cu un criteriu cum este *distanța* dintre nucleul semantic XO și elementele funcționale care îl "*înconjoară*", un asemenea nucleu este, în ultimă instanță, un substantiv comun obiectual (numit și *autosemanticin* [19]), un nume propriu, sau un substantiv personalizat (dar fără nume propriu, denominalizat).

M14 = {aparitia unui substantiv comun *obiectual* (nepredicațional, autosemantic), a unui nume propriu, sau a unui substantiv personalizat denominalizat}

M13 = {aparitia unui *modificator* (adjectiv, adverb, adjectiv pronominal)}

M12 = {aparitia unui *cuantificator* (generalizat)}

M11 = {pre-poziții sau *post-poziții* exprimând *cazul* (pentru N), *timpul* sau *semantismul* (pentru V) etc }

Ultima clasă de marcheri, notată MO (sau M00 pentru uniformitate), și ai cărei marcheri se aplică *formei de dicționar* a cuvântului, este reprezentată de *rolul funcțional* al *flexionării*.

Recapitulând, clasele de marcheri considerate de strategia lingvistică SCD, în particular de *algoritmul de segmentare* SCD-2002, pot fi reprezentate grafic de următoarea ierarhie [7]:

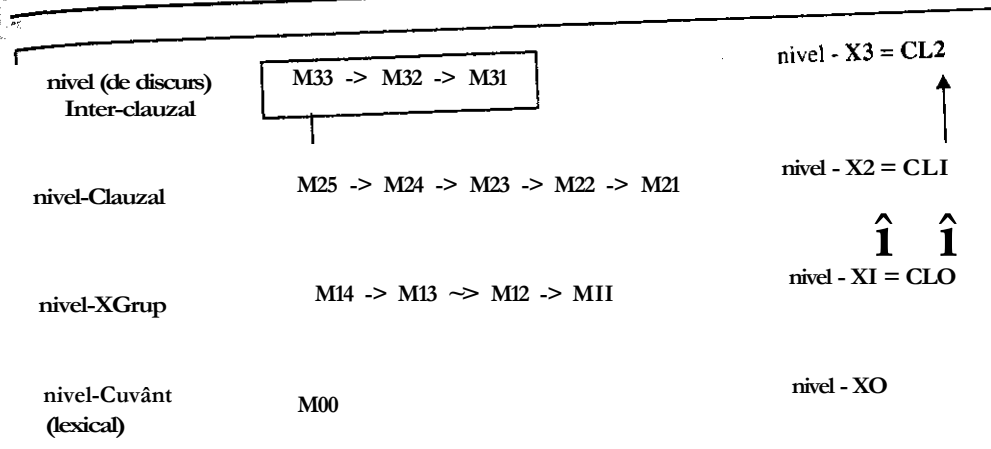


Figura 4.1.1. Clasele de marcheri SCD-2002 și ierarhia lor

Orientarea arcelor din Fig. 4.1.1., stabilite între clasele și subclasele de marcheri, provine dintr-o *ordine de prioritate descrescătoare* între marcherii considerați, și este reprezentată mai jos prin relația ">" dintre clasele și subclasele de marcheri. Această ierarhie este o *ipoteză de bază* impusă în *strategia lingvistică* SCD și, prin consecință, și în *algoritmul de segmentare SCD* [7].

$$(4.1.2) \quad VQ = 1+4) \quad M(k+1)(i+1) * M(k+1)j \quad (k = 0+2);$$

$$(4.1.3) \quad \vee (k = 0+2) \quad M(k+1)i > Mkj \quad 0 = 1+5), Q = 0+5).$$

Aceste inegalități ne spun că marcherii din subclasa $M(k+1)(j+1)$ sunt de *prioritate mai mare* în comparație cu marcherii din subclasa $M(k+1)j$, ($k = 0+2$), ($j = 1+4$), în cadrul aceleiași clase $M(k+1)$ de marcheri aflată pe *aceleiași nivel* de proiecție lingvistică, iar marcherii din aceeași clasă $M(k+1)$ au o *prioritate mai mare* față de marcherii din clasa Mk de pe nivelul *inferior* de proiecție lingvistică.

Această ierarhie a marcherilor și claselor de marcheri este considerată de noi ca fiind validă pentru *limba română*. Probabil că anumite modificări vor fi necesare când se trece de la un LN la altul. Dacă ne situăm în domeniul mai restrâns al limbajelor indo-europene (cum sunt franceza, engleza, germana, italiana, spaniola, posibil rusa), se poate aprecia că structurile și clasele de marcheri propuse în Tabelul 4.1. și Fig. 4.1.1. rămân aceleași sau foarte asemănătoare, cu anumite modificări parametrizate în funcție de limbaj.

4.2. Algoritmul SCD-2002 de segmentare și stabilire a dependențelor

Urmând algoritmi de segmentare și dependență (numiți și *meta-algoritmi* SCD) propuși în [1] și [2] (denotați în continuare SCD-1994), rafinați cu clasele de markeri considerate în subsecția precedentă, se obține forma prezentă a *algoritmului de segmentare SCD* (denotată SCD-2002), conform [7]. Dezvoltăm aici forma *secvențial-liniară* a acestui algoritm, însă în [1] sunt expuse și o formă *secvențial-recursivă*, ca și o versiune *paralelă* a algoritmului. O formă "*inversată*" (pentru care în intrare avem un arbore de derivare sau o formulă logică, iar în ieșire - ca și în intrarea în algoritmul standard - avem o frază) poate fi folosită pentru a ghida procesul de generare a unei fraze de LN [2], schimbând operația de recunoaștere a markerilor cu cea de generare a lor, și analiza (parsarea) compuşilor sintactici cu generarea lor.

În *descrierea algoritmului* de segmentare SCD-2002 sunt folosite câteva *operații* al căror înțeles este bine să fie precizat de la început.

- (4.2.a) *Recunoașterea markerilor* înseamnă inserarea în text a unor etichete adecvate, ce corespund markerilor care realizează delimitarea unităților textuale sintactice, semantice, și de discurs.
- (4.2.b) *Verificarea markerului* înseamnă preluarea, din baza de date a markerilor, a celor mai importante valori din structura de trăsături a acelui marker.
- (4.2.c) *Segmentarea* implică o analiză liniară (parsare) a secvenței de etichete de markeri, și *recunoașterea* unei subsecvente (eventual discontinuă) care face parte din secvența originală de etichete de markeri.
- (4.2.d) *Recunoașterea structurii sintactice* înseamnă *segmentarea* și *recunoașterea* structurilor sintactice elementare cum sunt NG, VG, AG, clauza infinită, și clauza finită.
- (4.2.e) *Compunerea structurilor (de dependență)* constă în stabilirea dependențelor (sub-ordonare, co-ordonare, supra-ordonare), succesive dintre structurile sintactice recunoscute, pe baza rolului funcțional specific al markerilor care delimitează aceste structuri, și utilizând ierarhia corespunzătoare dintre clasele cărora le aparțin acești markeri (vezi Fig. 4.1.1. și relațiile 4.1.2.-4.1.3.).

Algoritmul de segmentare SCD-2002

Step01. Recunoașterea pe text a markerilor din clasa M3;

Step01. Recunoașterea pe text a markerilor din clasa M2;

Step03. Verificarea contextuală și recunoașterea apariției corelate a markerilor de tip M3 și M2⁽¹⁾;

Step04. Segmentarea frazei în clauze finite;

Step05. Segmentarea (chunking), dacă este necesar, a clauzelor finite în clauze infinite;

[Stop: Dacă scopul procesării este de a obține o structură liniară a clauzelor finite și/sau infinite din frază].

Step06. Verificarea markerilor M3 și stabilirea relațiilor de dependență inter-clauzală⁽²⁾;

[Stop: Dacă scopul procesării este doar de a obține arborele de dependență a clauzelor finite (și infinite) din frază].

Step07. Recunoașterea pe text a markerilor din clasa M1;

Step08. Verificarea contextuală și recunoașterea (eventualei apariții corelate) a markerilor M1⁽³⁾;

Step09. Recunoașterea structurilor XG (X = N, V, A)⁽⁴⁾;

Step10. Verificarea markerilor M24 și M23, și stabilirea relațiilor de dependență dintre structurile infinite, intra-clauzale de tip XG⁽⁵⁾;

[Stop].

Indicii superiori (**n**) care apar în algoritmul de mai sus corespund următoarelor *remarci*:

- (1) Markerii corelați pot fi reprezentați ca upluri ordonate (liste) de markeri.
- (2) Relațiile de dependență clauzală pot fi stabilite (ca în [8, Anexa 9, p. 108], de exemplu) prin utilizarea unei *gramatici formale* (ambigue) definită pe secvențe de markeri din (sub)clasele M3, M25, M22, și M21.
- (3) Markerii complecși pot fi sintagme sau expresii de tipul gradelor de comparație a adjectivelor, diferiți cuantificatori generalizați etc.
- (4) În execuția acestui pas se realizează parsarea sintagmelor XG dintr-o clauză finită și infinită.
- (5) Dependențele dintre structurile de tip XG sunt stabilite în principal prin utilizarea trăsăturilor și valorilor de trăsături TENS = FINITE sau INFINITE, și PRED = AC sau EXIST, pe care le posedă nucleele semantice ale sintagmelor XG, X = N, V, A (a se vedea [3], [23]). Aceste valori pot fi *moștenite* din reprezentarea olexicon a cuvintelor care poartă aceste trăsături și care formează XG, sau pot fi *dobândite* de către nucleul semantic al XG în procesul de recunoaștere (parsare) a structurii.

5. Compararea algoritmilor de segmentare

5.1. Algoritmii de segmentare SCD-1994 și SCD-2002

Algoritmii SCD-1994 expuși în [1], [2] se bazează pe *patru* (sub)clase principale de marcheri, denotate acolo prin (clasele de) "1-marcheri" până la "4-marcheri". Aceste subclase de marcheri din SCD-1994 corespund următoarelor (sub)clase de marcheri din prezentul algoritm SCD-2002 [7]:

(5.1.1) 1-marcheri = M3 u M25 u M22 ;

2-marcheri = M24;

3-marcheri = M23;

4-marcheri = M21 u M1

Prezentăm în continuare algoritmul de segmentare SCD-1994 (în forma secvențial-recursivă), așa cum a fost expusă în [1, p.68-69], având ca scop parsarea LN. Algoritmul SCD-1994 (în forma secvențial-liniară) și destinat sarcinii de *generare* a LN este prezentat în [2, p. 172-173].

Algoritmul de segmentare SCD-1994 în formă secvențial-recursivă (SR)

Step01. Recunoașterea marcherilor de clauză.

Step02. Recunoașterea sintagmelor VG (grupuri verbale) finite și infinite.

Step03. Verificarea contextuală a marcherilor.

Step04. Segmentarea clauzală.

Step05. Segmentarea sub-clauzală.

Step06. Recunoașterea 1-marcherului;

Recunoașterea 1-structurii:

Wait-until 1-structura este completă.

Step07. Recunoașterea 2-marcherului;

Recunoașterea 2-structurii:

Wait-until structura de nivel-X2 este completă*.

Step08. Recunoașterea 3-marcherului;

Recunoașterea 3-structurii.

Step09. Recunoașterea 4-marcherului;

Procesarea 4-structurii.

Step10. 3-structură completă?

Nu: *Go-to* Step08.

Da: Compune 3-structuri; *Go-to* Step1.

Step11. 2-structură completă ?

Nu: *Go-to* Step07.

Da: Compune 2-structuri; *Go-to* Step12.

Step12. 1-structură completă ?

Nu: *Go-to* Step06.

Da: Compune 1-structuri; *Go-to* Stop.

Stop.

* *Structuri* AX-bar (în original, în [1]), înțelegând structuri sintactice derivate din *schemele* X-bar *augmentate*, definite în [20] și extinse în [3]. Scopul acestui pas al algoritmului este de a completa clauza finit introdusă printr-un grup verbal finit.

Principala problemă cu algoritmul de segmentare și dependență SCD-1994 (forma SR) este că sunt necesare "*multiple nivele de recursie pentru a completa și compune structurile*" [1, p.69].

5.2. Algoritmii de segmentare M-1997 și SCD-2002

În această subsecțiune vom arăta că algoritmul de segmentare M-1997 este *scufundat* în algoritmul SCD-2002 (de fapt, și în SCD-1994) [7].

M-1997 este un algoritm de "*suprafață*" destinat segmentării discursului în unități textuale de tip-clauzal. În timp ce, pentru acest scop, M-1997 folosește numai *marcheri de discurs* {"*cue phrases*" sau *conective*), algoritmul SCD-2002 utilizează un set de clase de marcheri mai larg și în același timp mai rafinat, se care include clasele de marcheri din M-1997 ca un caz particular. Mai precis, relațiile dintre clasele de marcheri Mari, Mar2, și Mar3 (vezi Secțiunea 2) utilizat pentru M-1997, și clasele de marcheri Mkj ale algoritmului SCD-2002 sunt următoarele:

(5.2.1) Mari u Mar2 u Mar3 c M3 u M25 u M22 u M21

sau, posibil, mai precis:

(5.2.2) Mari u Mar3 c M3 u M25 și Mar2 c M22 u M21

Diferența dintre algoritmii M-1997 și SCD-2002 nu constă doar în faptul că al doilea algoritm are un număr mai mare de clase, care sunt mai fine (mai precise), ci, mai important este faptul că aceste clase formează un *sistem ierarhic* (expus în Fig. 4.1.1.) ce este utilizat în procesele de *segmentare* și de *stabilire a dependențelor*. SCD-2002 furnizează noi clase de marcheri, cum sunt M23 și M24 (aparitia categoriilor predicative și/sau având un timp finit), precum și clasa M

cu subclasele sale (aparitia unor componente ale sintagmei XG, X = N, V, A). Acesta este un *prim argument* din care rezultă că M-1997 este *scufundat* în SCD-2002. "Scufundarea" este un termen care reflectă, de fapt, un proces de rafinare și de creștere a preciziei în calculul marginilor (limitelor) unităților textuale și a dependențelor dintre ele, pentru SCD-2002 în comparație cu M-1997.

Al doilea argument important care susține validitatea relației afirmate între cei doi algoritmi este următorul: fiecare *acțiune* din M-1997 are un corespondent într-o *operație* (sau o mulțime de *operații*) din algoritmul SCD-2002 (subsecțiunea 4.2).

Pentru segmentare, algoritmul M-1997 asociază fiecărui marker, în baza de date a markerilor, o anumită *acțiune* ce este statistic determinată de către analiza de corpus efectuată în [4]. Corespondența dintre operațiile algoritmului-SCD, și o *acțiune* din algoritmul-M, se face în felul următor:

(5.2.a) *Acțiunea* (și marcherii) NORMAL din algoritmul-M are același efect cu operațiile de procesare a markerilor de discurs din clasa M3 a algoritmului-SCD. Când este întâlnit un asemenea marker, aceasta înseamnă că o clauză (în SCD-2002) sau o unitate de tip-clauzal (în M-1997) este pe cale de a se încheia și o altă clauză, respectiv unitate de tip-clauzal, este probabil că va începe.

(5.2.b) *Acțiunile* COMMA, SET_AND, și SET_OR din algoritmul-M sunt folosite pentru a dezambiguiza rolul unor marcheri din M3 pentru care nu se poate aplica întotdeauna regula generală (*acțiunea* NORMAL). Acești marcheri sunt următorii pentru limba română: ",", [Eng: *comma*], "și" și "sau". Rolul acestor marcheri este ambiguu deoarece comportamentul lor nu este uniform în cadrul delimitării unităților textuale. SCD-2002 rezolvă aceste cazuri cu ajutorul utilizării unei gramatici formale de marcheri care descrie principalele reguli de delimitare și dependență a clauzelor (în limba română). Această gramatică (vezi indicele superior **(2)** din SCD-2002 și remarca corespunzătoare) are ca scop să recunoască secvențele cele mai frecvente de marcheri din clasele M3 și M2. Numai *câteva* dintre aceste reguli sunt incorporate în mod explicit în algoritmul M-1997 original.

(5.2.c) O unitate de tip-clauzal din M-1997 nu este în mod necesar o clauză finită în sensul gramatical al noțiunii, așa cum este folosit în algoritmul-SCD. O asemenea unitate de tip-clauzal, în sens M-1997, poate fi o întregă frază, formată din mai multe clauze finite. M-1997 folosește, de fapt, pentru segmentarea liniară a frazei în unități de tip-clauzal numai *trei reguli* din cele folosite de SCD-2002, iar aceste reguli sunt sintetizate de către *acțiunile* COMMA, SET_AND, SET_OR.

(5.2.d) *Acțiunile* MATCH_PAREN, MATCH_DASH, COMMA_PAREN sunt utilizate de către M-1997 pentru a delimita acele întinderi de text care pot fi omise atunci când fraza este segmentată în unități de tip-clauzal. Aceste părți "explicative" din text, considerate a nu fi importante, sunt, în text, puse între *paranteze*, (perechi de) *liniuțe de-despărțire*, sau (perechi de) *virgule*. Algoritmul M-1997 nu tratează aceste întinderi "parantetizate" de text ca fiind unități de tip-clauzal propriu-zise, ci le consideră doar ca fiind scufundate în unitatea de tip-clauzal de care aparțin. Pentru SCD-2002, aceste *acțiuni* M-1997 nu au un corespondent specific deoarece *paranteza* (închisă și deschisă), *virgula*, și *liniuța-de-despărțire* sunt tratate ca marcheri de discurs (M3), și fac parte din *gramatica de marcheri compuși* (concatenați) care este asociată cu algoritmul SCD-2002 de segmentare și dependență a clauzelor dintr-o frază.

(5.2.e) Din același motiv ca cel menționat mai sus, în (5.2.d), *acțiunile* DUAL, NORMAL_THEN_COMMĂ din M-1997 nu au, nici ele, un corespondent explicit în SCD-2002; aceste două acțiuni sunt de asemenea înglobate în *gramatica formală de secvențe de marcheri de discurs*, care se dovedește a fi, în mod clar, mai generală, ușor de extins (sau de restrâns), este dependentă de LN specific analizat, și modelează comportamentul marcherilor simpli și compuși (concatenați) de tip M3 și M2.

Relațiile (5.2.1-2) și observațiile (5.2.a-e) demonstrează că algoritmul de segmentare M-1997 este (chiar strict) scufundat în algoritmul-SCD (atât SCD-2002 cât și SCD-1994) [7]. Acest fapt, stabilit teoretic aici, este confirmat de către rezultatele empirice ale implementărilor, prezentate în subsecțiunea care urmează.

5.3. Execuția segmentării pentru algoritmi M-1997 și SCD-2002

Actuala etapă de implementare a algoritmilor de segmentare [7] este prezentată în exemplele care urmează. *Step06* din SCD-2002, și *Step12* din SCD-1994 stabilesc relațiile de *dependență inter-clauzală*, folosind o gramatică formală pentru marcherii de discurs, simpli și compuși (concatenați), din clasele M3 și M2. Această fază a algoritmului nu este încă implementată, în prezent. Să menționăm că stabilirea *dependențelor intra-clauzale* este (parțial) implementată prin utilizarea, pentru moment, (numai) a subclaselor M2 și M1 de marcheri. Marginile inter-clauzale din text sunt reprezentate prin *paranteze pătrate*, în timp ce pentru marginile și dependențele intra-clauzale sunt folosite *parantezele rotunde* (obișnuite). Indicii inferiori ai parantezelor pătrate arată numărul curent al unităților textuale de tip-clauzal din algoritmul M-1997, respectiv numărul curent al clauzelor obținute din algoritmul SCD-2002.

Exemplul 5.3.1.

Ex.5.3.1.Tag. (Etichetarea morfologică realizată cu mediul *TexTag* - vezi Fig. 5.4.1. și Fig. 5.4.2.)

<NSRY,23,0>Câmpul</NSRY,23,0> <V3,24,0>era verde</V3,24,0>
 <CR,25,0>și</CR,25,0> <NSRY,26,0>vița</NSRY,26,0>
 <S,27,0>de</S,27,0> <NSRN,28,0>vie</NSRN,28,0>
 <PXA,29,0>se</PXA,29,0> <V3,30,0>acoperise</A/3,30,0>
 <S,31,0>cu</S,31,0> <NPN,32,0>lăstari</NPN,32,0>
 <APN,33,0>verzi</APN,33,0><COMMA,34,0>,</COMMA,34,0>
 <NPRY,35,0>copacii</NPRY,35,0> <S,36,0>de pe</S,36,0>
 <NSRY,37,0>marginea</NSRY,37,0>
 <NSOY,38,0>șoselei</NSOY,38,0>
 <V3,39,0>înfrunziseră</V3,39,0> <CR,40,0>și</CR,40,0>
 <NSRY,41,0>briza</NSRY,41,0> <V3,42,0>sufila</V3,42,0>
 <S,43,0>dinspre</S,43,0>
 <NSRN,44,0>mare</NSRN,44,0><POINT,45,0>.</POINT,45,0>

Ex.5.3.1.Mar. (Rezultatul segmentării (fără dependențe), obținut prin aplicarea algoritmului M-1997 în cadrul mediului *ClauSEGM* - vezi Fig. 5.4.3.)

[Câmpul era verde și vița de vie se acoperise cu lăstari verzi, copacii de pe marginea șoselei înfrunziseră și briza sufla dinspre mare.]i

Ex.5.3.1.SCD. (Rezultatul segmentării (fără dependențe), obținut prin aplicarea algoritmului SCD-2002 în cadrul mediului *ClauSEGM* - vezi Fig. 5.4.4.)

[(Câmpul) era verde ^ și[(vița) (de (vie)) se acoperise (cu (lăstari) (verzi)) h ,[(copacii) (de pe (marginea (șoselei))) înfrunziseră], și[(briza) sufla (dinspre (mare)).].

Exemplul 5.3.2.**Ex.5.3.2.Tag.**

<S,1,0>în</S,1,0> <NSN,2,0>întuneric</NSN,2,0> <V2,3,0> ai fi
 zis</V2,3,0> <C,4,0>că</C,4,0> <V3,5,0>fulgeră </A/3,5,0> <R,6,0>ca</R,6,0>
 <NSRY,7,0>vara</NSRY,7,0> <COMMA,8,0>,</COMMA,8,0> <C,9,0>dar</C,9,0>
 <NPRY, 10,0>noaptea</NPRY,10,0> <V3,11,0>erau reci</A/3,11,0>
 <CR,12,0>și</CR,12,0> <QZ,13,0>nu</QZ,13,0> <PPSD, 14,0>ți</PPSD,14,0>
 <PXA,15,0>se</PXA,15,0> <V3,16,0> părea</V3,16,0> <R,17,0>deloc</R,17,0>
 <C,18,0>că</C, 18,0> <PXA,19,0>se</PXA,19,0> <V3,20,0>apropie</
 V3,20,0><NSRY,21,0>furtuna</NSRY,21,0><POINT,22,0>.</POINT,22,0>

Ex.5.3.2.Mar.

[În întuneric ai fi zis]i [că fulgeră ca vara,], [dar nopțile erau
 părea deloc], [că se apropie furtuna.].

EX.5.3.2.SCD.

[(în (întuneric)) ai fi zis]i [că fulgeră (ca (vara))], [dar (nop
), și [nu (ți) se părea (deloc)]. [că se apropie (furtuna)].

Exemplul 5.3.3.**Ex.5.3.3.Tag.**

<NSRY,46,0>Poarta</NSRY,46,0> <V3,47,0>era deschis
 <COMMA,48,0>,</COMMA,48,0> <TSR,49,0>un</TSR,49,0> <NS
 </NSN,50,0> <V3,51,0>ședeau</A/3,51,0> <S,52,0>la</S,52,0> <NSN
 </NSN,53,0> <S,54,0>pe</S,54,0> <TSR,55,0> o</TSR,55,0>
 bancă</NSRN,56,0><COMMA, 57,0>,</COMMA,57,0> <TSR,58,0>
 <NSRN, 59,0>ambulanță</NSRN,59,0> <V3,60,0>aștepta</A/3,60,0>
 </S,61,0> <NSRY, 62,0>ușa</NSRY,62,0> <S, 63,0>de</S,63,0>
 serviciu</NSN,64,0> <CR, 65,0>și</CR,65,0> <VG,66,0>intr
 <V1,67,0> am simțit</V1,67,0> <NSRY,68,0>mirosul</NSRY,68,0>
 pardoselii</NSOY,69,0> <S,70,0>de</S,70,0> <NSRN,71,0>marmură
 <S,72,0>și</S,72,0> <S,73,0>de</S,73,0> <NSN,74,0>spital</NSN
 75,0>.</POINT,75,0>

Ex.5.3.3.Mar. (întindere de text între paranteze acolade {...})

[Poarta era deschisă, {un soldat ședeau la soare pe o bancă
 aștepta la ușa de serviciu și intrând am simțit mirosul pardoselii de
 spital.

Ex.5.3.3.SCD.

[(Poarta) era deschisă } , [(un (soldat)) ședeau (la (soare)) (p
 k , [(o (ambulanță)) aștepta (la (ușa) (de (serviciu)))], și[intrând am
 (pardoselii)) (de (marmură) (și (de (spital))))].

Exemplul 5.3.4.**Ex.5.3.4.Tag.**

<NPRY,1,0>Trupele</NPRY,1,0> <V3,2,0>treceau</A/3,2,0>
 lângă</S,3,0><NSRN,4,0>casă</NSRN,4,0><COIv1MA,5,0>,</COM
 6,0>pe</S,6,0> <NSRN,7,0>șosea</NSRN,7,0><COMMA,8,0>,
 <CR,9,0>și</CR,9,0> <NSRY,10,0>praful</NSRY,10,0> <RELO,1
 RELO.H.OxZ.^O^/Z.^OxPPSA.IS.OH^PPSA.IS^<V
 V3,14,0> <PXA,15,0>se</PXA,15,0> <V3,16,0>astemă</A/3,16,0>

</S,17,0> <NPRY,18,0>frunzele</NPRY,18,0> <NPOY,19,0>copacilor</NPOY,19,0><POINT,20,0>.</POINT,20,0>

Ex.5.3.4.Mar.

[Trupele treceau pe lângă casă, pe șosea, și praful[^] [pe care-l ridicau se așternea pe frunzele copacilor.]₂

Ex.5.3.4.SCD. (clauză relativă - atributivă)

[(Trupele) treceau (pe lângă (casă)) , (pe (șosea))]t ,[și (praful) [pe care-(l) ridicau se așternea (pe (frunzele (copacilor)))].]₂

Exemplul 5.3.5.

Ex.5.3.5.Tag.

<QZ,76,0>Nu</QZ,76,0> <PPSA,77,0>m</PPSA,77,0><Z,78,0>-</Z,78,0>
<V3,79,0>a văzut<A/3,79,0> <CR,80,0>și</CR,80,0> <QZ,81,0>n</QZ,81,0><Z,
82,0>T,</Z,82,0><V1,83,0>am știut</V1,83,0> <C,84,0>dacă</C,84,0> <V3,85,0>
e</V3,85,0> <NSRY,86,0>cazului</NSRY,86,0> <C,87,0>să</C,87,0> <PPSA,
88,0>mă</PPSA,88,0> <V1,89,0>duc<A/1,89,0> <S,90,0>la</S,90,0> <PPS,91,0>
el</PPS,91,0> <C,92,0>să</C,92,0><Z,93,0>-</Z,93,0><PPSA,94,0>i</PPSA,
94,0> <V1,95,0>raportez<A/1,95,0> <C,96,0>că</C,96,0> <V1,97,0>am sosit</
V1,97,0> <C,98,0>sau dacă</C,98,0> <QZ,99,0>nu</QZ,99,0> <V3,100,0>e mai
bine</V3,100,0> <C,101,0>să</C,101,0> <PPSA,102,0>mă</PPSA,102,0> <V1,
103,0>duc<A/1,103,0> <C,104,0>să</C,104,0> <PPSA,105,0>mă</PPSA,105,0>
<V1,106,0>aranjez<A/1,106,0> <R,107,0>putin</R,107,0><POINT,108,0>.</POINT,
108,0>

Ex.5.3.5.Marc.

puțin.]₂

Ex.5.3.5.SCD

[Nu (m)-a văzut]i și[n-am știut]₂ [dacă e (cazul)]₂ [să (mă) duc (la (el))]₂
[să-(i) raportez]₂ [că am sosit]₂ [sau dacă nu e mai bine]₂ [să (mă) duc]₂ [să (mă)
aranjez (puțin)]₂

Exemplul 5.3.6.

Ex.5.3.6.Tag.

<NSRY,109,0>Fereastra</NSRY,109,0> <V3,110,0>era deschisă<A/3,
110,0><COMMA,111,0>,</COMMA,111,0> <NSRY,112,0>patu</NSRY,112,0>
<PSS,113,0>meu</PSS,113,0> <V3,114,0>era acoperit<A/3,114,0> <S,115,0>
cu</S,115,0> <NSRY,116,0>pătura</NSRY,116,0><COMMA,117,0>,</COMMA,

117,0> <NSRY,118,0>masca</NSRY,118,0> <S,119,0>de</S,119,0> <NSRY,
120,0>gaze</NSRY,120,0> <S,121,0>cu</S,121,0> <NSRY,122,0>cutia</NSRY,
122,0> <PSS,123,0>ei</PSS,123,0> <ASN,124,0>lunguiață</ASN,124,0> <S,1
0>de</S,125,0> <NSRN,126,0>tinichea</NSRN,126,0> <CR,127,0>și</CR,1
0> <NSRY,128,0>casca</NSRY,128,0> <S,129,0>de</S,129,0> <NSN,130
oțel</NSN,130,0> <V3,131,0>erau agățate<A/3,131,0> <S,132,0>pe</S,132
<DMSR,133,0>aceiași</DMSR,133,0> <NSN,134,0>cuier</NSN,134,0><POI
135,0>.</POINT,135,0>

Ex.5.3.6.Mar. (întindere de text între paranteze acolade {...})

[Fereastra era deschisă, {patul meu era acoperit cu pătura}, masca
gaze cu cutia ei lunguiață de tinichea și casca de oțel erau agățate pe ace
cuier.]i

EX.5.3.6.SCD.

[(Fereastra) era deschisă]i ,[(patul) (meu) era acoperit (cu (pătura))
(masca) (de (gaze)) (cu (cutia) (ei (lunguiață))) (de (tinichea)))] și (casca) (de (o
erau agățate (pe (aceiași (cuier)))).]₂

Exemplul 5.3.7.

Ex.5.3.7.Tag.

<V1,1,0>Aș vrea<A/1,1,0> <C,2,0>să</C,2,0><Z,3,0>-</Z,3,0><PP
4,0>ți</PPSD,4,0> <V1,5,0>spun<A/1,5,0> <C,6,0>că</C,6,0> <CR,7,0>și</
7,0> <R,8,0>mai</R,8,0> <R,9,0>târziu</R,9,0><COMMA,10,0>,</COMMA,10
<CR,11,0>și</CR,11,0> <S,12,0>într</S,12,0><Z,13,0>-</Z,13,0><ASN,14,0>
</ASN,14,0> <NSRN,15,0>parte</NSRN,15,0><COMMA,16,0>,</COMMA,16
<V1,17,0>am văzut<A/1,17,0> <C,18,0>că</C,18,0> <NPRY,19,0>lucruri
NPRY,19,0> <PXA,20,0>se</PXA,20,0> <V3,21,0>întâmplă<A/3,21,0> <R,2
tot așa</R,22,0><COMMA,23,0>,</COMMA,23,0> <C,24,0>dar</C,24
<V3,25,0>ar fi nevoie</V3,25,0> <S,26,0>de</S,26,0> <PI,27,0>oarecari</PI
0> <NPN,28,0>precizări</NPN,28,0> <CR,29,0>și</CR,29,0> <V1,30,0>simt
V1,30,0> <C,31,0>că</C,31,0> <QZ,32,0>nu</QZ,32,0> <PPSD,33,0>mi</PP
33,0><Z,34,0>-</Z,34,0><V3,35,0>ar ajunge<A/3,35,0> <NSRY,36,0>respira
</NSRY,36,0><COMMA,37,0>,</COMMA,37,0> <C,38,0>că</C,38,0> <V1,39,0>
ocoli<A/1,39,0> <R,40,0>prea</R,40,0> <R,41,0>mult</R,41,0><POINT,42,0>
POINT,42,0>

Ex.5.3.7.Mar. (întindere de text între paranteze acolade {...})

[Aș vrea]t [să-ți spun]₂ [că și mai târziu, {și într-altă parte,} am văzut]₂
lucrurile se întâmplă tot așa,]₂ [dar ar fi nevoie de oarecari precizări și simt]₂
nu mi-ar ajunge respirația,]₂ [că aș ocoli prea mult.]₂

Ex.5.3.7.SCD.

[Aș vrea h [să-(ți) spun], [că și (mai (târziu)) , și (într-(altă (parte))) , am văzut], [că (lucrul) se întâmplă (tot așa)], ,[dar ar fi nevoie (de (oarecari (precizări)))], și[simt], [că nu (mi)-ar ajunge (respirația)], ,[că as ocoli (prea (mult))].]

5.4. Programele *TexTag* și *ClauSEGM*

În cele ce urmează sunt prezentate câteva imagini de execuție în cadrul programelor *TexTag* și *ClauSEGM*, scrise în Visual C++ 5.0, și utilizate pentru a eticheta și segmenta texte de LN (limba română) [7]. Figurile 5.4.1. și 5.4.2. se referă la *TexTag*, Figura 5.4.3. conține execuția algoritmului de segmentare M-1997 în cadrul *ClauSEGM*, iar Figura 5.4.4. conține o execuție a algoritmului de segmentare SCD-2002 sub mediul *ClauSEGM*. Stabilirea relațiilor de dependență inter- și intra-clauzale, pentru aceleași două tipuri de algoritmi, urmează să fie implementată în cadrul aceluiași mediu *ClauSEGM*.

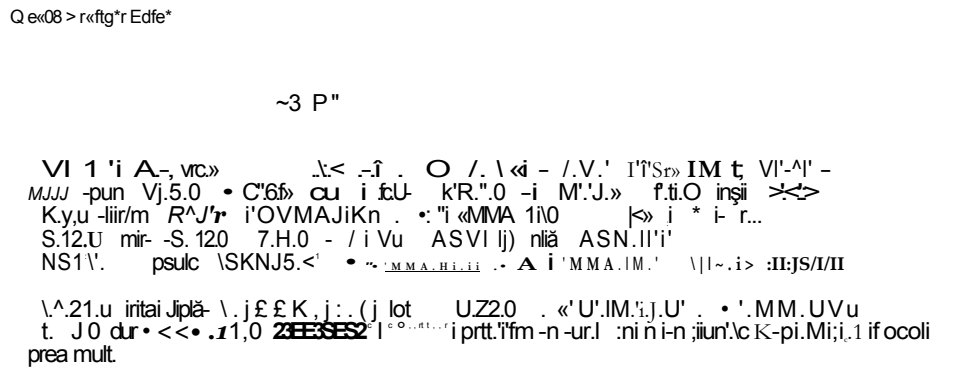


Figura 5.4.1. Rezultatul etichetării morfologice sub *TexTag*

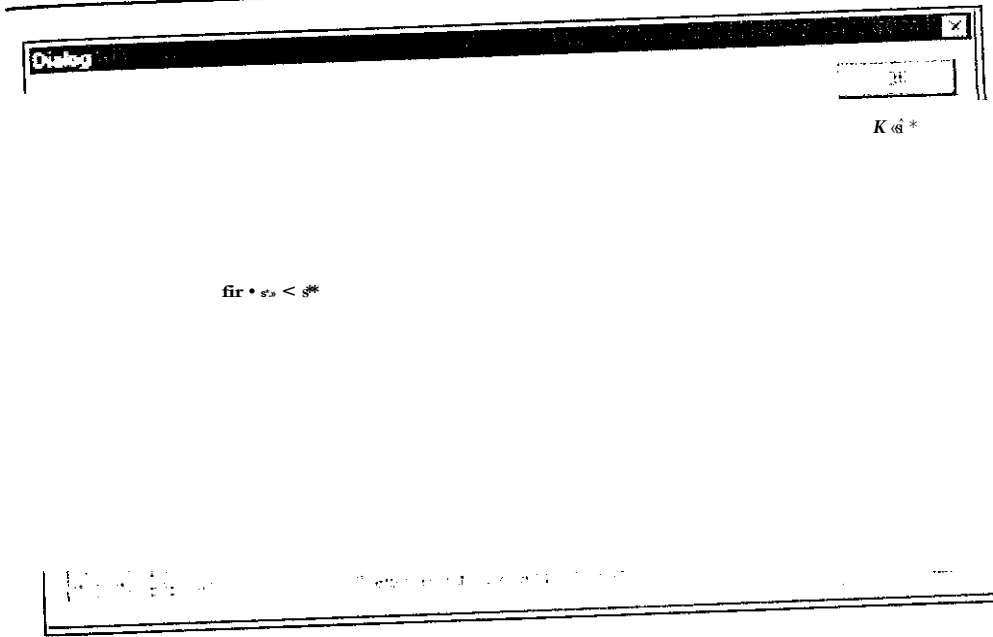


Figura 5.4,2. Lista de etichete selectată cu un meniu din *TexTag*

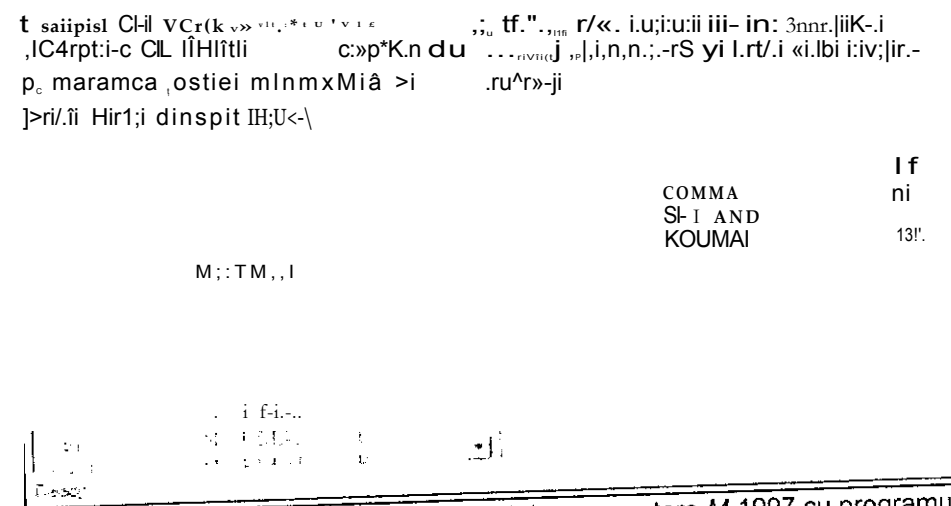
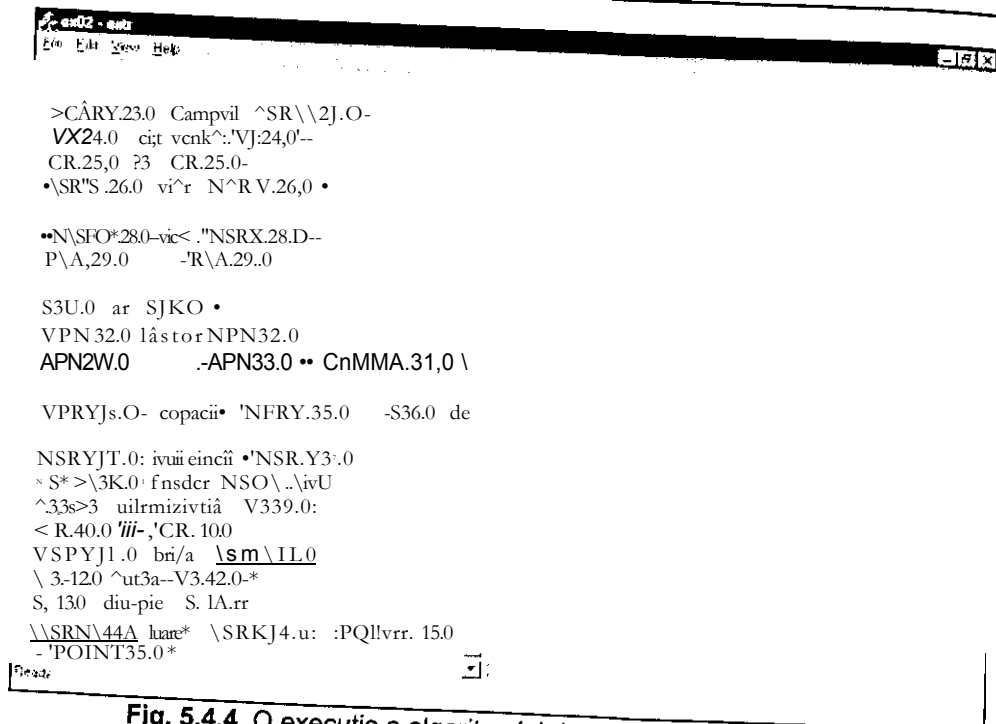


Figura 5.4.3. O execuție a algoritmului de segmentare M-1997 cu programul *ClauSEGM*



```

>CĂRY.23.0 Campvil ^SR\\2J.O-
VX24.0 ciț venk^.:VJ:24,0'--
CR.25,0 ?3 CR.25.0-
•\SR"S.26.0 vi^r N^R V.26,0 •

•N\SFO*28.0-vic<."NSRX.28.D--
P\A,29.0 -R\A.29..0

S3U.0 ar SJKO •
VPN 32.0 lăs tor NPN32.0
APN2W.0 -APN33.0 • CnMMA.31,0 \

VPRYJ.s.O- copacii • NFRY.35.0 -S36.0 de

NSRYJT.0: ivuū eincū •NSR.Y3.0
* S* >\3K.0 fnsdcr NSO\ .\ivU
^33s>3 uilmizivtiā V339.0:
< R.40.0 'iii-'CR.10.0
VSPYJ1.0 bii/a \sm\IL0
\ 3-120 ^ut3a--V3.42.0-*
S, 130 diu-pie S.1A.rr
\\SRN\44A luare* \SRKJ4.u: :PQllvrr. 15.0
- 'POINT35.0*

```

Fig. 5.4.4. O execuție a algoritmului de segmentare SCD-2002 sub mediul *ClauSEGM*

6. Concluzii

Rezultatele obținute în această lucrare nu se referă strict la compararea și implementarea celor doi algoritmi de segmentare. Avem, de fapt, două tipuri de algoritmi de segmentare (și dependență), și fiecare din cele două tipuri reprezintă linii specifice de cercetare, cu importante consecințe asupra domeniilor de procesare a LN cărora se adresează: algoritmul M-1997 [4] este destinat (teoriei și) aplicațiilor de procesare a discursului, generare automată a LN, și rezumării automate, în timp ce algoritmul SCD-2002 [7] se încadrează mai curând în teorii sintactice ale LN, cum sunt teoria FX-bar [3], parsarea bazată pe teorii (principii) sintactico-semantice ale LN, dar și punerea în evidență a structurilor (segmentelor) și relațiilor de discurs [6].

Demonstrarea relației (de scufundare) dintre cele două tipuri de algoritmi de segmentare [7], schițarea (în secțiunea 1) a unui *cadru formal general* pentru *algoritmii de segmentare* a LN, în particular a segmentării de tip *chunking*, propunerea (în cadrul algoritmilor-SCD) unei metode generale de segmentare în unități textuale a LN și de stabilire a dependențelor între ele, toate acestea

constituie posibile noi perspective pentru abordările teoretice și aplicative în procesarea automată a LN, inclusiv, și mai ales, pentru limba română.

Revenind la aspectele concrete expuse în acest articol, vom reveni la algoritmi și algoritmi, la algoritmi către analiza complexă a structurilor semantico-discursive, la clasele de markeri, și perfecționarea actualelor implementări rămânând direcții de continuare a prezentei abordări.

Referințe bibliografice

- [1] Neculai Curteanu (1994). *From Morphology to Discourse Theory: Structures in the SCD Parsing Strategy*, Language and Linguistics, Akademia Libroservo, Prague, p. 61-73.
- [2] N. Curteanu, G. Holban (1996). *Strategia lingvistică SCD aplicată la generarea limbii române*, Limbaj și Tehnologie (D. Tufiș, Academia Română, p. 169-176.
- [3] Neculai Curteanu (2000). *Towards a Funcțional X-bar Theory*, Report, Institute of Theoretical Informatics, Romanian Academy, 32 p.
- [4] Daniel Marcu (1997). *The Rhetorical Parsing, Summarization, and Segmentation of Natural Language Texts*, Ph.D. Thesis, Univ. of Toronto, Canada.
- [5] Daniel Marcu (2000). *The Theory and Practice of Discourse Segmentation and Summarization*, The MIT Press, Cambridge.
- [6] O. Popârda, N. Curteanu (2002). *L'evolution du discours juridique analysé par la stratégie linguistique SCD*, LINCOS Studies in Linguistics, LINCOS Europa, Munchen.
- [7] N. Curteanu, C. Linteș (2002). *Segmentation Algorithms for Textual Units*, Research Report, Institute of Theoretical Informatics, Romanian Academy.
- [8] N. Curteanu, D. Cristea, P. Mihaescu (1982). *Cercetări în domeniul calculului prin intermediul limbajului natural*. Contract de cercetare nr. 4774/1982, Universitatea Iași - ICI București.
- [9] Neculai Curteanu (1983). *Algoritmi de analiză sintactică a frazelor în limba românești*. Lucrările Conferinței INFO-IAȘI'83, p. 533-549.
- [10] D. Cristea, N. Curteanu, P. Mihaescu (1983). *Implementarea algoritmului de analiză sintactică și definitivarea proiectului de analiză sintactică*. Cercetare nr. 1906/1983, Universitatea Iași - ICI București.
- [11] N. Curteanu (1984). *Aspecte ale analizei logice a limbajului natural*. Cercetare nr. 4709/1984, Universitatea Iași - ICI București.

- [12] Rebecca Passonneau, Diane Litman (1997). *Intention-based segmentation: human reliability and correlation with linguistic cues*, in Proc. 31th Annual Meeting of ACL, Ohio, p. 148-155.
- [13] Lance Ramshaw, Michel P. Marcus (1999). *Text Chunking Using Transformation-based Learning*, in (S. Armstrong et al., eds.) "Natural Language Processing Using Very Large Corpora", Kluwer Acad. Publ., p. 157-176.
- [14] Victor Raskin, S. Nirenburg (1999). *"Lexical Rules for Deverbal Adjectives"*, in E. Viegas (Ed.) *Breadth and Depth of Semantic Lexicons*, Kluwer Acad. Publ., p. 99-119.
- [15] M. Johnson, Federica Bosa (1999). *"Qualia Structure and Compositional Interpretation of Compounds"*, in E. Viegas (ed.) *Breadth and Depth of Semantic Lexicons*, Kluwer Acad. Publ., p. 167-186.
- [16] Denis Bouchard (2001). *La source sémantique des facteurs hétérogènes qui régissent la distribution des adjectifs*, Conferința Internațională "Représentations du Sens Linguistique", București.
- [17] Dumitru Irișia (1997). *Morfo-sintaxa verbului românesc*. Editura Universității "Al. I. Cuza", Iași.
- [18] Eva Hajicova, H. Skoumalova, P. Sgall (1995). *An Automatic Procedure for Topic-Focus Identification*. *Computational Linguistics*, 21(1): 81-94.
- [19] P. Sgall, E. Hajicova, J. Panevova (1986). *The Meaning of the Sentence in Its Semantic and Pragmatic Aspects*. Kluwer Academic Publishers, Dordrecht.
- [20] Neculai Curteanu (1988). *Augmented X-bar Schemes*. COLING'88 Proceedings, Budapest, p. 130-132.
- [21] Dan Tufiș (2000). *Using a Large Set of EAGLES-compliant Morpho-Syntactic Descriptors as a Tagset for Probabilistic Tagging*, in Proceedings of the LREC'2000 International Conference, Athens.
- [22] Dan Tufiș, A.M. Barbu (2001). *Computational bilingual lexicography: automatic extraction of translation dictionaries*, In Romanian Journal on Information Science and Technology, voi. 4, no. 3.
- [23] Neculai Curteanu (2002). *Către o teorie X-bar funcțională* (în prezentul volum).

O metodă automată pentru inserarea diacriticelor în texte în limba română

Rada F. MIHALCEA
University of Texas at Dallas, Richardson, Texas, U.S.A.
rada@utdallas.edu

Vivi A. NĂSTASE
University of Ottawa, Ottawa, Canada
vnastase@site.uottawa.ca

1. Introducere

Problema restaurării diacriticelor constă în inserarea diacriticelor într-un text în care lipsesc. Creșterea continuă a numărului de texte disponibile pe Internet face ca metodele automate de inserare a diacriticelor să devină componentă esențială în multe aplicații importante, cum ar fi extragerea de informații, traducerea automată, colecționarea de texte, construirea dicționarelor electronice și multe altele. Corectarea erorilor ortografice poate să aibă un impact major asupra calității rezultatelor obținute în aceste aplicații. De exemplu, în absența unei metode de restaurare a diacriticelor, unele cuvinte devin ambigue, cum este cazul cuvintelor din limba română *pește*, *peste* sau *pături*, *paturi*. O căutare bazată pe astfel de cuvinte poate returna multe texte irelevante (de exemplu o căutare pentru *pește* ar returna și documente conținând *peste*). De asemenea, traducerea unor astfel de cuvinte într-o limbă străină poate fi eronată (de exemplu traducerea corectă a cuvântului *pături* în limba engleză este *blankets*, dar în absența diacriticului este tradus greșit ca și *beds*).

Metodele dezvoltate până în prezent pentru rezolvarea acestei probleme se bazează în general pe dicționare și pe diverse procesoare lexicale și/sau sintactice. Multe dintre limbile care se confruntă cu problema restaurării diacriticelor nu beneficiază însă de astfel de resurse, și ca urmare aplicabilitatea acestor metode este limitată la limbi bine studiate care dispun de suficiente resurse. Lucrarea de față prezintă o metodă automată de reinserare a diacriticelor în text care necesită doar o colecție de texte de dimensiuni modeste. Spre deosebire de alte metode dezvoltate anterior, metoda introdusă în această lucrare nu necesită nici un fel de dicționare sau procesoare morfologice și/sau sintactice.



și prin urmare poate fi folosită pentru prelucrarea de texte în orice limbă care dispune de un număr minim de texte cu diacritice. Datorită lipsei de restricții, metoda propusă este foarte generală și ușor aplicabilă pentru orice limbă. Pentru a demonstra aceasta afirmație, după ce vom prezenta experimentele pentru texte în limba română, vom arăta câteva rezultate obținute pentru limbile cehă, poloneză și maghiară.

2. Experimente anterioare

Restaurarea diacriticelor nu este în sine o problemă dificilă. Experimentele efectuate până în prezent au demonstrat că folosirea de dicționare electronice poate duce la o acuratețe de peste 90% în restaurarea accentelor pentru limbile franceză și spaniolă [9],[11],[5]. Metoda descrisă de Michael Simard în [9] este o îmbunătățire adusă unei metode propusă anterior de El-Beze [4]. Această metodă se bazează pe Hidden Markov Models și învață folosind cuvintele învecinate. Precizia raportată este de 99%. Tufiş și Chițu [10] propun o metodă similară pentru inserarea diacriticelor în texte în limba română. Yarowsky prezintă în [11] un set de metode folosite pentru restaurarea accentelor în limbile franceză și spaniolă. Majoritatea algoritmilor pe care îi prezintă se bazează pe dicționare și cuvinte învecinate pentru a decide asupra ortografiei potrivite pentru fiecare cuvânt ambiguu. Yarowsky compară N-gram taggers, clasificatoare Bayesiene și liste de decizii cu metoda de bază care constă în folosirea unui dicționar. Pentru cele două limbi considerate în experimentele raportate, listele de decizii duc la performanțele cele mai ridicate. Toate aceste tehnici se bazează însă pe context, dicționare și în unele cazuri pe informații adiționale de natură morfologică și sintactică. Nagy et al. prezintă în [7] o abordare diferită a problemei, în care șiruri de litere sunt extrase din fiecare cuvânt și folosite pentru a obține statistici. Folosind metoda propusă, s-a observat o precizie foarte bună obținută pe texte în limba franceză. Experimentele prezentate în [7] sunt asemănătoare cu cele raportate în [1], unde măsuri de similaritate între trigrame sunt folosite pentru a automatiza corectarea greșelilor de ortografie.

Majoritatea studiilor efectuate până în acest moment pe această temă, s-au ocupat de limbi bine cunoscute și răspândite, cum ar fi franceza și spaniola. Foarte puține studii s-au concentrat pe limbi mai puțin mediatizate cum ar fi ceha, slovena, turca sau alte limbi care folosesc diacritice. Tabelul 1' prezintă diacriticele folosite în limbile europene cu alfabet latin. După cum rezultă din această listă, numeroase limbi se confruntă cu problema restaurării diacriticelor. Din setul de 36 de limbi cuprinse în tabel, engleza pare să fie singura limbă pentru care diacriticele

Tabelul cuprinde numai litere mici. Fiecărei litere mici îi corespunde o literă mare. Informația din acest tabel a fost agregată din liste de diacritice în limbi europene, disponibile la adresa www.tiro.com/diintro.html

nu constituie o problemă. Cuvintele din engleză care conțin diacritice au fost împrumutate din alte limbi, și varianta acestora fără diacritice^ nu are un corespondent care să ducă la ambiguitate. Diacriticele par însă să aiba un rol important în diferențierea cuvintelor. Engleza, care după cum spuneam nu are diacritice per se, are în schimb o ambiguitatea semantică mai ridicată.

Tabel 1

Diacritice din limbile europene cu alfabet latin

Limba	Diacritice	Limba	Diacritice
Albaneză	șe	Malteză	Cgh2
Bască	nu	Norvegiană	âaso
Bretonă	â e n u u'	Olandeză	ee
Catalană	â ș e e i î 16 6 u i i	Poloneză	^ c ș l n 6 ș z £
Cehă	â e t f e i f t 6 F § f u u y 2	Portugheză	â ä ș e î o o o
Daneză	â & 0	Română	â â î ș ț
Engleză	None	Sami (Laponă)	â i' C d n r) Stz
Estoniană	â C 8 6 § u z	Serbo-croată	c £ d S z
Faroeză	â a e d i 6 o u y	Slovacă	â ä e d' e i i r n â o f & f i x < / I
Finlandeză	â â o s z	Slovena	Uz
Franceză	â â a ^ e e e e î o o e u u y	Spaniolă	e u n
Galițiană	â e i o u	Suedeză	â ä o
Germană	â o u B	Turcă	ț g i o ș u
Islandeză	â a s d e i 6 6 u y t >	Sorbiană(1)	
Italiană	â e e i î o o u u	Sorbiană (2)	c 5 z S1A f S s z z
Maghiară	â e i 6 5 6 u U i i	Welsh	â e î o u y y

Aplicabilitatea metodelor menționate anterior este limitată în următoarele cazuri:

Studii efectuate pe corpusuri bilingve paralele, ar arătat ca vocabularul construit dintr-un text în limba engleză este aproximativ jumătate din vocabularul construit pe baza aceluiaș text într-o altă limbă. Competiția SENSEVAL [6] raportează de asemenea precizii mult mai mici pentru engleză comparativ cu alte limbi în rezolvarea ambiguității semantice. Lipsa diacriticelor în limba engleză ar putea constitui o explicație a acestui fenomen.

1. Dicționarele electronice nu sunt disponibile sau doar dicționare de dimensiuni relativ mici sunt făcute publice. Mai mult decât atât, în cazul în care dicționarul însuși nu are diacritice, metodele care se bazează pe această resursă pentru restaurarea diacriticelor devin inaplicabile.
2. Procesoarele folosite pentru analiză morfologică și/sau sintactică, considerate folositoare pentru problema restaurării diacriticelor, nu există sau nu sunt public disponibile.
3. Numărul de texte disponibile conținând diacritice este relativ mic. Mărimea corpusurilor publice sau disponibile prin Internet influențează mărimea vocabularului care poate fi construit ad-hoc pe baza acestor texte. În plus, majoritatea siturilor care publică texte pe Internet preferă în multe cazuri să evite diacriticele din motive de simplitate, uniformitate, sau pur și simplu din lipsa de mijloace necesare pentru codificarea diacriticelor.

Lucrarea de față prezintă o metodă de restaurare a diacriticelor bazată pe învățarea la nivel de literă, și nu la nivel de cuvânt. Avantajul principal al acestei metode este faptul că oferă posibilitatea de generalizare dincolo de cuvinte. Metoda este folositoare mai ales pentru limbile pentru care resursele disponibile sunt limitate, în speță limbi care nu au dicționare electronice mari cu diacritice. Limbi cunoscute și bine studiate, precum franceza și spaniola, pot de asemenea beneficia de aceasta metodă pentru procesarea cuvintelor necunoscute.

Experimentele prezentate în această lucrare adresează în principal problema restaurării diacriticelor în texte în limba română. Precizia observată pe limba română este de 99%, măsurată la nivel de literă. Experimente similare au fost efectuate pe alte trei limbi, și anume poloneză, maghiară și cehă, de asemenea cu rezultate foarte bune. Avantajul principal al metodei este faptul că nu necesită nici o etapă de preprocesare, ci numai un corpus relativ mic format din texte cu diacritice. Datorită simplității algoritmului, viteza de procesare este foarte mare, de aproximativ 20 pagini de text pe secundă, măsurată pe un calculator cu un procesor Pentium III cu frecvența de 500MHz și 250MB memorie.

[^] Practic, metoda propusă încearcă să învețe reguli aplicabile la nivel de literă. În loc de a învăța reguli care se aplică la nivel de cuvânt, cum ar fi *"anuncio se scrie anuncio atunci când are funcția de verb"*, dorim să învățăm reguli aplicabile la nivel de literă, cum ar fi *"s urmat de i și spațiu și precedat de spațiu se scrie ș"*. Astfel de reguli, învățate la nivel de literă, sunt mai generale și au aplicabilitate mai mare, în special în cazurile în care dicționarele disponibile sunt de dimensiune redusă, când se întâlnesc multe cuvinte necunoscute în textul dat, sau când procesoare pentru analiză morfologică sau sintactică nu sunt la îndemână.

Este evident că în analiza limbajului literele constituie nivelul cu granularitatea cea mai scăzută, și de aceea au și cel mai mare potențial de generalizare. În loc de aproximativ 150.000 de unități candidate potențiale pentru algoritmul (mărimea aproximativă a vocabularului de uz general a unei limbi), vom

Șvea mai mult sau mai puțin 26 caractere pe baza cărora se vor constitui intrare pentru algoritmul de dezambiguare¹.

3. Experimente

Scopul experimentelor descrise în această lucrare este de a învățarea la nivel de literă este posibilă și poate rezolva, cu precizie problema restaurării diacriticelor. Pe lângă faptul că metoda propusă este o problemă de cercetare, scopul învățării la un nivel de granularitate atât de mic este de a oferi o metodă viabilă pentru limbile pentru care resursele isemantice disponibile sunt limitate, și pentru care restaurarea diacriticelor la nivel de cuvânt este greu de realizat.

3.1. Date

Prezentăm în primul rând experimentele efectuate pe textele în limba română. Limba română nu este o limbă foarte răspândită și în consecință are foarte multe resurse publice disponibile pentru pre-procesare, iar dicționarele electronice sunt de dimensiuni relativ mici. În al doilea rând, am avut de față o problemă specifică de restaurare a diacriticelor într-un dicționar electronic în limba engleză care conține aproximativ 75.000 de cuvinte, dar are dezavantajul că diacriticele lipsesc. Am considerat că este avantajos să studiem problema restaurării diacriticelor și să folosim acest dicționar, în loc să ne bazăm pe alte dicționare cu diacritice de dimensiuni reduse. În plus, pentru procesoarele pe care le dezvoltăm pentru limba română avem nevoie de numeroase texte electronice în limba română. De obicei aceste texte nu au diacritice, și deci restaurarea diacriticelor este din nou necesară. Avem de asemenea posibilitatea de a testa eficacitatea acestei metode cu rezultate obținute în experimente efectuate pe aceeași limbă constând în metode în care învățarea se face la nivel de cuvânt.

Pentru a aplica metoda descrisă în lucrarea de față, avem deci nevoie de o colecție de texte românești cu diacritice. În acest scop, am colectat un dicționar "România Literară"², un ziar românesc publicat săptămânal, cu articole de literatură. Ziarul are o versiune care conține diacritice începând din anul 2000. Colecția disponibilă on-line la data colectării datelor (august 2001) conține 2780 articole. În pasul următor, textul a fost extras din fișierele HTML și deosebită a fost acordată doar caracterelor românești. Alte caractere din text au fost întâlnite ocazional, cum ar fi e, e etc. au fost transformate în forma lor corectă fără diacritice, având în vedere că suntem interesați doar de caracterele

¹ Numărul de litere depinde de limba care se analizează. S-a arătat de exemplu că aproximativ 85% dintre cuvintele în limba franceză nu au o formă ortografică unică și deci numai 20.000 de cuvinte sunt potențial ambigue. Pe de altă parte, în

² sunt ambigue în limba franceză.

Accesibil prin <http://www.romlit.ro>

și nu de caractere franceze sau din alte limbi. După toate aceste faze premergătoare, am obținut un corpus conținând aproximativ 3 milioane de cuvinte.

Literele mari au fost transformate în litere mici. Cazul literelor *â* și *î* este special în limba română: deși pronunția lor este identică, folosirea lor este guvernată de reguli bazate pe poziția lor în cuvânt. La începutul cuvântului se folosește întotdeauna *î*, iar *â* se folosește în interiorul cuvântului. Este bine cunoscut faptul că folosirea acestor litere a fost controversată de-a lungul timpului. O lege din anii '60 a schimbat ortografierea de la *â* la *î*, singura excepție fiind cuvintele derivate din rădăcina *român*. La începutul anilor '90 ortografia veche a fost reintrodusă, și astfel s-a ajuns la cazuri de texte inconsistente, în care se întâlnesc scrieri diferite ale aceluiași cuvânt. De exemplu, *cîntec* și *cântec* sunt forme ale aceluiași cuvânt care pot fi întâlnite în același text. Ziarul "*România Literară*" păstrează încă ortografia cu *î*, cu mici excepții (de exemplu, articole scrise de scriitori invitați care preferă să scrie folosind *â* în loc de *î*).

3.2. Algoritmi de învățare

Pentru a rezolva problema restaurării diacriticelor, am ales să folosim un algoritm bazat pe învățarea de instanțe (IBL). Există două motive importante care au stat la baza luării acestei decizii. În primul rând, este faptul demonstrat că excepțiile au un rol important în procesarea limbajelor naturale. Algoritmii de tip IBL sunt recunoscuți pentru faptul că iau în considerare fiecare exemplu de antrenament în luarea unei decizii de clasificare [2], și deci folosirea acestui tip de algoritmi prezintă un avantaj deosebit în probleme de limbaj natural. În al doilea rând, acest gen de algoritmi sunt foarte eficienți relativ la timpul de antrenament și testare.

Învățarea pe bază de instanțe se desfășoară în felul următor: în pasul de antrenament, toate exemplele de intrare sunt memorate. În faza de testare, fiecare exemplu din set este comparat cu exemplele memorate și va primi-clasificarea dată de exemplul memorat de care este cel mai apropiat, distanța fiind dată de măsura specifică aleasă în implementarea folosită. Pentru efectuarea experimentelor propuse, am folosit implementarea TiMBL [3] a acestor algoritmi. În plus, am efectuat experimente asemănătoare și cu un clasificator pe bază de arbori de decizie, și anume C4.5 [8]. Arborii de decizie sunt construiți din setul de exemple de antrenament. La fiecare pas este ales un atribut care discriminează cel mai bine exemple din clase diferite (prin valorile sale). Grupele obținute prin diviziunea după acest atribut vor fi din nou împărțite în grupe mai mici și mai pure, prin alegerea unui nou atribut care discriminează cel mai bine exemplele din grupă. Acest proces continuă până când grupele obținute au un grad de puritate acceptabil, sau mărimea arborelui depășește un prag ales inițial. Rezultatele obținute folosind C4.5 sunt asemănătoare cu cele obținute folosind TiMBL, însă C4.5 are capacitatea de a genera reguli expresive, folosite pentru implementări practice.

Având în vedere că lucrăm la nivelul literelor, atributul care trebuie învățat este constituit de litera ambiguă. Acesta poate fi oricare din literele ambigue enumerate în Tabelul 1. Pentru limba română avem 4 perechi de litere ambigue: s -

>ș, t - ț, a - ă, i - î. Literele mari au fost convertite în prealabil în litere mici. Datorită faptului că datele folosite aplică ortografia cu *î*, nu avem ambiguitatea a - â, ci doar ambiguitatea / - î. Aceasta nu implică însă o pierdere de generalitate. Conversia între cele două forme de ortografie este simplă și se poate realiza folosind doar poziția literei în cuvânt, și prin urmare scrierile diferite nu afectează rezultatul algoritmului.

3.3. Atribute

Atributele folosite în orice algoritm de învățare au un impact foarte mare asupra eficacității algoritmului. După cum am menționat și în introducere, nu avem posibilitatea de a folosi procesoare care determină partea de vorbire a cuvintelor, și nici un alt fel de analizoare morfologice sau sintactice. În plus, nu dorim să ne bazăm pe cuvintele învecinate, deoarece avem un număr limitat de date, și în consecință există șansa de a întâlni un număr mare de cuvinte necunoscute. Prin urmare, ne-am decis asupra folosirii unor atribute foarte simple, pentru extragerea cărora nu este nevoie de nici un fel de procesare specială. Vom folosi litere învecinate, cu o notație specială atribuită spațiilor, virgulelor și punctelor (aceste caractere pot afecta procesul de învățare, fiind considerate caractere speciale de către C4.5 și/sau TiMBL).

Dacă X este litera a cărei ambiguitate trebuie rezolvată, atributele folosite sunt N litere la stânga și la dreapta literei ambigue:

$$L_{-n}, L_{-n-1}, \dots, L_{-1}, X, L_1, L_2, \dots, L_n$$

Acest set de atribute se comportă surprinzător de bine, relativ la acuratețe, după cum vom arăta în cele ce urmează.

După cum am menționat anterior, am ales să nu ne bazăm pe nici un tag obținut cu procesoare lexicale sau morfologice, ci doar pe informația care se poate extrage din text neprelucrat. De asemenea, suntem interesați să găsim posibilități de generalizare, astfel încât un corpus limitat să poată fi folosit pentru a genera reguli de reinserare a diacriticelor. În loc de a învăța reguli bazându-ne pe cuvinte, după cum s-a procedat până acum, dorim să învățăm reguli bazate pe litere, pentru că acestea constituie cele mai mici unități în limbaj, și oferă posibilitatea învățării chiar și dintr-o colecție mică de texte.

Pentru fiecare pereche ambiguă de litere, parcurgem textul și generăm toate exemplele posibile întâlnite în text. Atributele într-un exemplu sunt formate folosind N litere la stânga și la dreapta literei ambigue, și atributul țintă este însăși litera ambiguă. Forma generală a exemplurilor generate este:

$$L_{-n}, L_{-n-1}, \dots, L_{-1}, X, L_1, L_2, \dots, L_n$$

unde ca și în exemplul anterior, X este litera ambiguă. Prezentăm mai jos exemple de vectori de atribute care constituie date de intrare pentru algoritmul de învățare pentru rezolvarea ambiguității perechii s - ș. CO, DO și SP sunt codurile care înlocuiesc virgula, punctul și spațiul.

*l, i,n,SP,(u,b,SP,i,n,s.
e,CO,SP,r,o,-,g,a,r,d,ș.
g, a, r, d, i, t, u, l, CO, SP, s.
e,SP,o,r,a,DO,SP,t,o,t,ș.*

Învățarea se reduce la detectarea corelațiilor între valorile atributelor care caracterizează exemplele de antrenament și valorile atributelor țintă, și utilizarea acestora pentru stabilirea valorii atributului țintă din exemplele de testare.

Numărul de exemple extrase din corpus depinde de perechea de litere. Din întregul set de 3 milioane de cuvinte, am obținut 2.161.556 exemple pentru perechea ambiguă a - ă, 2.055.147 pentru perechea / - l, 1.257.458 exemple pentru t - ț, și în final 866.964 exemple pentru perechea s - ș. În fiecare din aceste cazuri, spațiul exemplelor este împărțit în două clase, date de cele 2 variante ale literei ambigue. Metoda de învățare automată va folosi atributele date pentru a găsi reguli de clasificare a exemplurilor în cele 2 clase.

3 * 4 ■ Rezultate

Precizia cea mai ridicată s-a obținut pentru o fereastră de 10 litere în vecinătatea literei ambigue (N = 5). Dată fiind această observație, am considerat că este important să studiem mai în detaliu acest caz, și să determinăm ratele de învățare pentru cele 4 perechi de litere ambigue. Cu toate acestea, prezentăm rezultate pentru ferestre de diverse dimensiuni, pentru comparație.

Tabelul 2 arată rezultatele obținute pentru N=5. Preciziile raportate în acest tabel sunt obținute folosind algoritmul bazat pe învățarea de instanțe. Am efectuat experimente cu seturi de antrenament de diverse dimensiuni, variind de la 2.000.000 exemple până la 10 exemple, pentru a determina rata de învățare și dimensiunea minimă a corpusului necesară pentru a obține o precizie satisfăcătoare. În toate aceste experimente s-au folosit seturi de testare conținând 50.000 exemple. Pentru a obține rezultate cât mai exacte am folosit validare încrucișată folosind 10 seturi diferite de test. Tabelul indică de asemenea baza de comparație, definită aici ca fiind precizia obținută când se folosește implicit litera cea mai frecventă din fiecare pereche ambiguă.

Rezultatele prezentate în Tabelul 2 sunt reprezentate grafic în Figura 1. Este interesant de observat că cea mai importantă fază a procesului de învățare are loc când se folosesc primele 10 000 exemple. În conformitate cu măsurătorile efectuate, a rezultat că aproximativ 100.000 - 250.000 caractere (aproximativ 25-60 pagini de text) sunt necesare pentru a genera 10.000 exemple cu diacritice, ceea ce constituie un corpus de dimensiune relativ mic. Mai departe, pentru a obține îmbunătățiri de numai 1% este necesar un număr semnificativ de exemple. Tabelul 2 indică de asemenea, în caractere groase, prima precizie care depășește baza de comparație, ca o indicație a dimensiunii minime a setului de antrenament

pentru care se observă o formă minimă de învățare. După cum se vede din tabel, numai 1.000 exemple sunt suficiente pentru învățare.

Rezultate obținute în rezolvarea ambiguității literelor cu diacritice în română

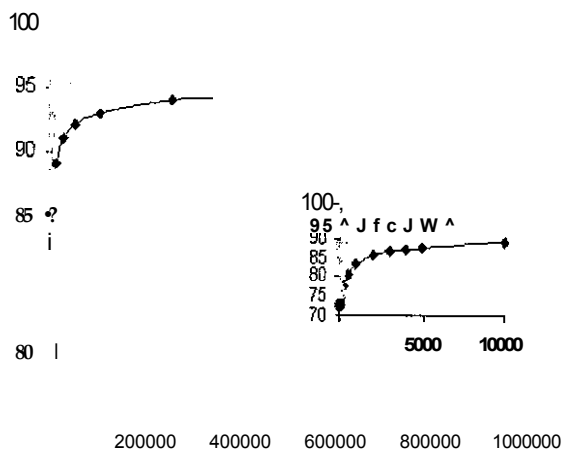
	Pereche ambiguă			
	a - ă	a-ă(2)	i - î	s - ș
Nr.tota! exemple	2.161.566	1.369.517	2.055.147	866.964
Baza comparație	74.70°/c	85.90°/c	88.205	76.53°/c
Exemple de Antrenament	Precizie obținută pe date de test (50.000 exemple)			
2,000,000	96.14%	-	99.69%	-
1,000,000	95.10%	99.14%	99.58%	-
750,000	94.83%	98.97%	99.53%	99.07%
500,000	94.57%	98.79%	99.46%	98.86%
250,000	94.00%	98.37%	99.28%	98.87%
100,000	93.03%	97.56%	98.96%	98.54%
50,000	92.10%	96.86%	98.57%	98.13%
25,000	90.99%	95.75%	98.11%	97.58%
10,000	88.99%	93.75%	97.31%	96.53%
5,000	87.56%	92.76%	96.65%	95.61%
4,000	86.91%	91.86%	96.49%	94.99%
3,000	86.39%	90.99%	96.19%	94.18%
2,000	85.81%	89.93%	95.49%	93.47%
1,000	83.49%	88.36%	93.78%	92.31%
500	80.61%	85.66%	93.07%	90.75%
250	77.89%	83.17%	92.75%	87.41%
100	74.80%	84.04%	91.41%	82.13%
50	72.79%	82.73%	88.05%	86.53%
25	72.45%	81.34%	88.15%	78.26%
10	73.38%	85.90%	88.20%	75.88%

Folosind întregul set de exemple extrase din corpusul de antrenament, precizia de învățare pentru perechea i - î este foarte aproape de 100%. Pentru această pereche, acum o instanță greșită din 300 instanțe, în timp ce baza de comparație este o instanță greșită din fiecare 8 instanțe, deci o îmbunătățire semnificativă.

Cel mai slab rezultat este obținut în cazul perechii a - ă. După cum se vede din tabel, reiese că principalul motiv care cauzează această precizie

este faptul că multe substantive în limba română au forma nearticulată terminată în *ă* și forma articulată terminată în *a*. De exemplu, *masa* și *masă* reprezintă forma articulată și respectiv nearticulată a substantivului *masă*. De asemenea, timpuri diferite ale aceluiași verb se disting numai prin terminația în *a* sau *ă*. Algoritmul de învățare este deci indus în eroare din cauza folosirii acestor litere în contexte identice. O soluție simplă constă în evitarea în procesul de învățare a exemplelor care conțin *a* sau *ă* la sfârșitul unui cuvânt. Rezultatele obținute sub această ipoteză simplificatoare sunt raportate în Tabelul 2, în coloana a-ă(2). După cum se arată în tabel, câștigul este de mai mult de 4% în precizie folosind doar această condiție simplă (câștig care se traduce într-o reducere a erorii de 87%).

Am folosit de asemenea și algoritmul de învățare bazat pe arbori de decizie C4.5, cu aceleași date de antrenament, fără a observa însă nici o îmbunătățire comparativ cu rezultatele raportate în Tabelul 2. Dezavantajul folosirii C4.5 pentru această problemă este faptul că faza de învățare este mult mai lentă decât în cazul folosirii algoritmului TiMBL. Pe de altă parte, C4.5 are capacitatea de a genera reguli expresive. "Dacă $L_1=e$ și $L_2=spațiu\ atunci\ s$ " (99.5%), "Dacă L^t și $L_2=spațiu\ atunci\ s$ " (98.7%), "Dacă $L_1=p$ și $L.^v$ și $L_1=f$ și $L_2=e\ atunci\ s$ " (95.5%), sunt exemple de astfel de reguli. L /denotă o literă învecinată în poziția i relativ la litera ambiguă. Se observă că aceste reguli nu țin cont de faptul că literele folosite în clasificare aparțin aceluiași cuvânt sau nu: Algoritmul de învățare se bazează pur și simplu pe litere, indiferent de cuvântul căruia îi aparțin. În consecință, pseudo-omonimele (cum ar fi *peste* și *pește*), sunt adresate în mod egal de această metodă, pentru că algoritmul are capacitatea de a se extinde dincolo de cuvinte.



70 4-
Figura 1. Rate de învățare pentru diacriticele în limba română. Graficul din mijloc este o reprezentare mărită a zonei 0-10.000

3.5. Ferestre de dimensiune diferită

Am efectuat experimente cu ferestre de diverse dimensiuni, pentru a determina dimensiunea contextului care modelează cel mai bine problema noastră. Pentru aceasta, am considerat ferestre de dimensiune dpi, șase, zece, patrusprezece și optsprezece litere învecinate (i.e. $N = 1,3,5,7,9$). Rezultate comparative sunt prezentate în Tabelul 3. Aceste numere trebuie comparate cu primul rând din Tabelul 2 (coloana corespunzătoare valorii $N=5$ în tabelul de față).

Tabel 3

Rezultate comparative obținute cu ferestre de dimensiuni diferite

Pereche ambiguă	Dimensiune fereastră				
	N=1	N=3	N=5	N=7	N=9
a-ă(2)	88.63%	98.79%	99.14%	99.10%	99.10%
i-î	94.18%	99.13%	99.69%	99.68%	99.43%
s-ș	88.09%	99.06%	99.07%	99.02%	99.00%
t-t	89.45%	98.57%	98.75%	98.67%	98.25%

Când nu există suficient context disponibil, o fereastră de dimensiune $N=3$ poate fi folosită fără a pierde mult din precizie. Însă, după cum am specificat și înainte, cea mai ridicată acuratețe se obține pentru o fereastră de zece litere înconjurătoare ($N=5$).

3.6. Comparație cu experimente asemănătoare

Rezultatele prezentate în lucrarea de față se pot compara cu rezultatele raportate de Tufiș și Chițu [10], care au folosit tot limba română în experimentele lor. Tufiș și Chițu menționează că sarcina recuperării diacriticelor în limba română este mai dificilă decât în alte limbi, deoarece în română diacriticele sunt mai intens folosite. După cum raportează în experimentele lor, numai 60% din cuvintele din limba română nu au diacritice, comparat cu studii menționate în [9] care arată că aproximativ 85% dintre cuvintele limbii franceze se scriu fără accent.

Abordarea prezentată de Tufiș și Chițu folosește dicționare, un analizor morfologic, iar învățarea se face la nivel de cuvinte. Folosind aceste resurse, au obținut o precizie globală de 97.4%. Nu putem efectua o comparație directă a rezultatelor noastre, având în vedere că atât metodele, cât și modul de evaluare, sunt fundamental diferite. Precizia medie de 99% pe care noi o raportăm este măsurată la nivel de literă, pe când acuratețea raportată în [10] este determinată la nivel de cuvânt¹.

¹ Diferența dintre precizia raportată la nivel de literă și precizia raportată la nivel de cuvânt rezultă practic din diferența de granularitate dintre litere și cuvinte. Presupunând că un cuvânt conține L litere ambigue, o singură literă din acest set L a cărei ambiguitate este rezolvată greșit face ca întreg cuvântul să fie considerat greșit, pe când la nivel de litere avem doar o singură eroare din setul L . Pe de altă parte, chiar dacă mai multe litere din

Metodologia noastră depășește abordările anterioare, prin faptul că s-au obținut precizii și viteze de procesare ridicate fără a folosi nici un fel de resurse adiționale cum ar fi procesoare pentru analiză morfologică sau sintactică sau dicționare. Din aceste motive, algoritmul se poate aplica oricărei limbi, singura cerință fiind un corpus relativ mic de texte cu diacritice.

4. Alte limbi

Pentru a demonstra generalitatea algoritmului pe care l-am propus, am efectuat experimente pe texte în alte trei limbi europene care fac uz de diacritice: cehă, poloneză și maghiară. Limbile considerate pentru aceste experimente sunt limbi cu răspândire restrânsă, pentru care resursele publice sunt limitate.

Pentru fiecare dintre aceste limbi am colectat texte cu diacritice disponibile prin Internet. Principalele surse folosite pentru formarea setului de date sunt după cum urmează: (1) pentru cehă, am folosit arhiva ziarului *Lidovky* și texte literare de *Kafka*, *Hasek* și *Capek*; (2) pentru maghiară, arhiva furnizată de către *Digitális Irodalmi Akademia* și un roman de *Petőfi Sándor*, (3) pentru poloneză, arhiva ziarului *Wiedza i zycie*. Pe lângă acestea, am mai folosit texte adiționale colectate de pe diverse surse, astfel încât să obținem un corpus de minim un milion de cuvinte pentru fiecare limbă. Asemănător cu procesarea aplicată limbii române, datele au fost convertite în fișiere text, iar literele mari au fost transformate în litere mici. În urma acestei etape de pre-procesare, am obținut un corpus de 1.46 milioane cuvinte pentru cehă, 1.72 milioane cuvinte în maghiară și 2.5 milioane cuvinte în poloneză.

Algoritmii de învățare și atributele folosite în procesul de învățare sunt identice cu cele folosite în experimentele efectuate pe limba română, raportate în detaliu în secțiunea precedentă. Tabelul 4 prezintă rezultatele obținute pentru cele trei limbi. Pentru fiecare set de litere ambigue, sunt prezentate în tabel: (1) numărul de exemple obținute din corpusul limbii respective, (2) baza de comparație, măsurată ca fiind precizia ce se poate obține dacă pentru fiecare set ambiguu se folosește implicit litera cu frecvența de apariție cea mai ridicată, și (3) precizia obținută prin aplicarea metodei propusă în lucrarea de față.

Media obținută pentru toate patru limbile studiate (cele trei limbi a căror rezultate sunt prezentate în Tabelul 4, și limba română) este de 98.17%. Precizia medie măsurată pe fiecare limbă în parte este influențată de mărimea setului de date folosit. Textele colectate pentru cehă și maghiară conțin aproximativ 1.4-1.7 milioane cuvinte, și prin urmare precizia obținută în aceste două limbi este mai joasă decât pentru poloneză și română, pentru care am reușit să colectăm un corpus de 2.5-3 milioane cuvinte. Estimăm deci posibilitatea creșterii preciziei ca urmare a creșterii dimensiunii corpusului de antrenament.

setul L sunt rezolvate greșit, avem tot o singură eroare la nivel de cuvânt, dar mai multe erori la nivel de literă. Nu este deci foarte clar care ar fi modalitatea corectă de a compara aceste două metode care lucrează la nivele de granularitate diferită.

Tabel 4

Rezultate obținute în restaurarea diacriticelor în trei limbi europene

Set litere ambigue	Număr exemple	Bază comparație	Metodă propusă
Cehă			
a â	649,886	75.01%	96.96%
c t	217,570	72.21%	97.08%
d <r	271,070	99.05%	99.86%
e e	768,051	74.59%	97.02%
i i	504,298	60.43%	96.29%
n fi	439,552	98.97%	99.71%
o 6	566,521	99.08%	99.86%
r f	319,352	65.55%	97.60%
s s	380,805	84.44%	98.88%
t f	387,214	99.05%	99.85%
u ti u	264,408	80.89%	93.51%
y y	191,317	65.55%	95.06%
z z	219,082	66.49%	98.70%
Medie			97.83%
Maghiară			
a â	1,198,294	73.51%)	96.91%
e e	1,306,944	76.34%	96.40%
i i	647,137	89.14%	99.49%
o 6 6 6	678,012	71.15%	96.10%
u u ii îi	207,753	56.00%	97.31%
Medie			97.04%
Poloneză			
a a_	1,387,019	88.83%	97.07%
c e	657,669	91.50%	99.42%
e e.	1,305,584	89.23%	98.47%
l j	506,041	59.29°/,	98.80%
n ri	878,824[96.75°/,	99.85%
o 6	1,230,38S)	88.67°>	99.87%
s á	688,67*	88.67°>	99.83%
z z z	896,90S>	86.26°>	99.73%
Medie			99.02%

Este interesant de observat că numărul de diacritice într-o limbă nu influențează precizia medie obținută. Precizia care se obține în cazul limbii maghiare, care are un total de 5 seturi de litere ambigue, este mai scăzută decât precizia care se obține pentru limba cehă, care are un număr impresionant de diacritice (treisprezece). Și aceasta cu toate că datele colectate pentru limba maghiară sunt mai numeroase decât datele colectate pentru limba cehă.

5. Concluzii

Am descris în lucrarea de față o metodă de restaurare a diacriticelor bazată pe tehnici de învățare la nivelul de literă. Avantajul principal al metodei constă în capacitatea ei de generalizare dincolo de cuvinte. Nu este necesară nici un fel de analiză a textului, și nu se folosesc nici un fel de procesoare de limbaj sau dicționare. Singura cerință este un corpus relativ mic de texte cu diacritice.

Metoda este folositoare în special pentru limbi pentru care nu sunt disponibile dicționare electronice de dimensiune adecvate, și nici procesoare pentru analiză morfologică și/sau sintactică. Mecanismul de învățare folosește date de intrare extrase din texte neprelucrate, și generează rezultate cu o precizie ridicată. Experimente detaliate efectuate pe texte în limba română au arătat că restaurarea diacriticelor în această limbă se poate efectua folosind metoda propusă cu o precizie de peste 99% la nivel de literă. Rezultatele au fost validate prin experimente efectuate pe alte trei limbi europene care fac uz de diacritice: cehă, poloneză și maghiară. Precizie medie măsurată pe cele patru limbi de studiu este de 98.14%, fapt care demonstrează că metoda este independentă de limbă, în plus, un alt avantaj al metodei este faptul că, datorită simplității sale, viteza de procesare este foarte mare, de până la 20 pagini de text pe secundă.

Referințe bibliografice

- [1] Angell, R., Freund G., Willett, P. *Automatic spelling correction using a trigram similarity measure*. Information Processing and Management 19, 4 (1983), 255-261.
- [2] Daelemans, W., van den Bosch, A., Zavrel, J. *Forgetting exceptions is harmful in language learning*. Machine Learning 34, 1-3 (1999), 11-34.
- [3] Daelemans, W., Zavrel, J., van der Sloot, K., van den Bosch, A. *TiMBL: Tilburg memory based learner, version 4.0, reference guide*. Tech. Rep., University of Antwerp, 2001.
- [4] El-Beze, M., Merialdo, B., Rozeron, B., Derouault, A., *Accentuation automatique des textes par des methodes probabilistes*. Techniques et sciences informatique 16, 6 (1994), 797-815.

- [5] Galicia-Haro, S., Bolshakov, I., Gelbukh, A. *A simple Spanish part of speech tagger for detection and correction of accentuation error*. In Text, Speech and Dialogue - Second International Workshop, TSD'99, September 1999, Proceedings (Plzen, Czech Republic, 1999), voi. 1692 of Lecture Notes in Computer Science, Springer, pp. 219-222.
- [6] Kilgariff, A., ed., Proceedings of SENSEVAL-2, 2002.
- [7] Nagy, G., N., N., and Sabourin, M. *Signes diacritiques: perdus et retrouvés*. In Actes du 1er Colloque International Francophone sur l'Écrit et le Document CIFED '98 (Quebec, Canada, 1998), pp. 404-412.
- [8] Quinlan, J. *C4.5: programs formachine learning*. Morgan Kaufman, 1993.
- [9] Simard, M. *Automatic insertion of accents in French text*. In Proceedings of the Conference on Empirical Methods in Natural Language Processing EMNLP-3 (Granada, Spain, 1998).
- [10] Tufiş, D., Chițu, A. *Automatic diacritics insertion in Romanian texts*. In Proceedings of the International Conference on Computational Lexicography COMPLEX'99 (Pecs, Hungary, June 1999).
- [11] Yarowsky, D. *Corpus-based techniques for restoring accents in Spanish and French texts*. In Natural Language Processing Using Very Large Corpora. Kluwer Academics Publisher, 1999, pp. 99-120.

Contribuții privind structura statistică de cuvinte în limba română scrisă*

Adriana VLAD și Adrian MITREA
Universitatea "POLITEHNICA" din București
Facultatea de Electronică și Telecomunicații
B-dul. Iuliu Maniu, 1-3, București, România
adriana_vlad@yahoo.com

1. Introducere

Această lucrare aparține unui studiu mai larg dedicat de autori descrierii limbii române ca sursă de informație. Punctul de plecare al acestui studiu a fost presupunerea generală conform căreia limba naturală este bine aproximată de un lanț Markov ergodic multiplu, cu ordin de multiplicitate mai mare decât 30, [1]. Descrierea acestei surse Markov multiple se realizează prin aproximații succesive. Investigația noastră statistică până în prezent a descris structura de litere, digrame, trigrame, tetragrame, precum și probabilitățile condiționate de o literă precedentă, [2]-[8].

Obiectivul principal al prezentei lucrări este descrierea sursei de informație fără memorie având ca simboluri cuvintele limbii române scrise. Aceasta presupune determinarea probabilității unui cuvânt (oricare ales), în caz că această probabilitate există. Determinarea probabilității a însemnat implicit și o verificare a ipotezei de staționaritate a limbii române scrise pe baza structurii de cuvinte; verificarea s-a făcut utilizând o procedură similară cu cea pe care am dezvoltat-o pentru m-gramme, [3]-[8] (m-grama este o succesiune de m litere consecutive în texte naturale). Metoda noastră statistică de a determina probabilitățile cuvintelor a combinat următoarele tipuri de inferențe statistice: teoria estimării cu multiple intervale de încredere statistică; test al ipotezei că probabilitatea aparține unui interval; test de egalitate între două probabilități.

Primele două tipuri de inferențe statistice menționate (teoria estimării cu multiple intervale de încredere statistică; test al ipotezei că probabilitatea aparține unui interval) au fost folosite pentru a decide care este intervalul de încredere statistică "reprezentativ pentru probabilitatea cuvântului investigat în textul natural.

O parte din acest studiu s-a desfășurat în cadrul unui Grant CNCSIS-MEC (2001-2002) cu tema: 'Descrierea limbii române scrise ca sursă de informație'

Simultan a apărut și o mulțime "reprezentativă" de date *Ud.* extrase din textul natural, corespunzătoare cuvântului investigat (modelul statistic *i.i.d.* presupune că datele provin din variabile aleatoare independente statistic și identic distribuite).

Ultimele doua tipuri de inferențe statistice menționate (test al ipotezei ca probabilitatea aparține unui interval; test de egalitate între două probabilități) ca și intervalele de încredere statistica "reprezentative" și mulțimile de date "reprezentative" obținute în prealabil au fost folosite pentru comparații matematice între texte naturale. Aceste comparații matematice (dincolo de valoarea lor ca atare) au avut scopul principal de a vedea dacă putem vorbi de un model matematic al sursei de cuvinte pentru limba ca ansamblu, pe domenii ale limbii, pe autori etc. Comparațiile s-au făcut în două moduri:

- urmărind probabilitățile unui cuvânt (același) în texte naturale diferite;
- urmărind probabilitățile cuvintelor situate pe același rang în texte naturale diferite (se compară probabilitățile asociate unui aceluiași rang în ierarhiile frecvențelor relative).

Rezultatele experimentale au adus probe noi în sprijinul ipotezei de staționaritate a limbii române scrise în cadrul unui aceluiași domeniu punctând către unele diferențe între domenii diferite.

Investigația noastră (atât privind intervalele de încredere statistică "reprezentative", cât și comparația matematică dintre texte) a avut în vedere și eroarea statistică de ordinul al doilea. Acest tip de eroare are un rol special în dimensionarea unui nou corpus lingvistic care să satisfacă acuratețea dorită pentru descrierea modelului matematic (sursa de informație de cuvinte).

Lucrarea mai conține și un studiu experimental al uneia dintre cele mai cunoscute legi de tipul rang - frecvență, legea lui Zipf. Este analizat și un corolar al acesteia, de interes lingvistic.

Analiza experimentală s-a bazat pe corpusul lingvistic global pe care l-am alcătuit în prealabil pentru studiul structurilor de litere, digrame, trigrame și tetragrame (vezi spre exemplu [6]). Acest corpus este format din 93 de cărți în limba română, scrise cu noua ortografie (introdusă după 1993). Cărțile reprezintă: literatură scrisă de autori români (11 cărți: romane și nuvele), literatură străină tradusă în română (47 de romane și nuvele), cărți științifice (drept, medicină, silvicultură, istorie, sociologie etc.) și altele. Au fost considerate doar cele 31 de litere ale limbii române (A Ă Â B C D E F G H I Î J K L M N O P Q R S Ș T Ț U V W X Y Z) precum și caracterul spațiu (blanc); orice alte simboluri (cifre, elemente de ortografie sau punctuație) au fost eliminate (suprimate).

Rezultatele experimentale au fost obținute pe diverse corpusuri organizate pe baza celor 93 de cărți:

- Corpusul Mixt Global (**#CMG**) - obținut prin concatenarea aleatoare a celor 93 de cărți; acesta conține un număr de $L_c = 8806433$ cuvinte dintre care $N_c = 202403$ sunt distincte.
- Cele două jumătăți ale Corpusului Mixt Global: prima jumătate (**#1JCMG**) și a doua jumătate (**#2JCMG**); acestea conțin un număr de $L_c = 4403217$ cuvinte și respectiv $L_c = 4403216$ dintre care $N_c = 148853$ și respectiv $N_c = 137845$ sunt distincte.
- Corpusul Literar Global (**#CLG**) - obținut prin concatenarea aleatoare a 58 de cărți (romane și nuvele scrise de autori români sau traduse în română); acesta conține un număr de $L_c = 6255235$ cuvinte dintre care $N_c = 162124$ sunt distincte.
- Cele două jumătăți ale Corpusului Literar Global: prima jumătate (**#1JCLG**) și a doua jumătate (**#2JCLG**); acestea conțin un număr de $L_c = 3127618$ cuvinte și respectiv $L_c = 3127617$ dintre care $N_c = 116247$ și respectiv $N_c = 116860$ sunt distincte.
- Corpusul Științific Global (**#CSG**) - obținut prin concatenarea aleatoare a 11 de cărți; acesta conține un număr de $L_c = 1049969$ cuvinte dintre care $N_c = 59093$ sunt distincte.

Au fost făcute determinări atât pe o singură carte cât și pe grupuri de cărți scrise de același autor. Dintre acestea menționăm:

- #1. George Călinescu, *Bietul Ioanide*, Editura Minerva, București, 1995, ISBN 973-21-0432-5 (voi 1, ISBN 973-21-0431-7, pag. 1-214), {voi. 2, ISBN 973-21-0433-3, pag. 5-256}, {voi. 3, ISBN 973-21-0434-1, pag. 5-238}.
- #2. Radu Anton Roman, *Precum fumul*, Editura Cartea Românească, București, 1996, ISBN 973-23-0274-7, pag. 5-283.
- #3. Radu Anton Roman, *Zile de pescuit*, Editura Metropol, București, 1996, ISBN 973-562-073-1, pag. 11-302.
- #4. John le Carre, *Casa Rusia*, Editura Univers, București, 1997, ISBN 973-34-0457-8, pag. 9-355.
- #5. John le Carre, *Spionul care venea din frig*, Editura Univers, București, 1996, ISBN 973-34-0355-5, pag. 9-252, cu ortografie actualizată.
- #6. John Le Carre, *Micuța toboșăreasă*, Editura Univers, București, 1998, ISBN 973-34-0430-6, pag. 7-443, cu ortografie actualizată.
- #7. Alexandr Soljenițin, *Arhipelagul Gulag*, Editura Univers, București, (voi. I, 1997, ISBN 973-34-0454-3, pag. 7-432), {voi. II, 1997, ISBN 973-34-0480-2, pag. 5-474}, {voi. III, 1998, ISBN 973-34-0497-7, pag. 5-414}, cu ortografie actualizată, fără note.

Primul pas în analiza noastră a fost evaluarea frecvențelor relative ale cuvintelor din corpusurile menționate anterior. Tabelul 1 conține primele 55 de cuvinte din ierarhia frecvențelor relative din diverse corpusuri.

Tabel 1

Ierarhia frecvențelor relative în câteva corpusuri
0. Rang; 1. Cuvânt; 2. Frecvență relativă (în %)

	#CIVIG		#CLG		#1		#4+#5+#6		#6		#CSG	
	1	2	1	2	1	2	1	2	1	2	1	2
0												
1	de	4,10	de	4,02	de	4,17	de	4,17	de	4,12	de	4,87
2	și	3,20	și	3,12	și	2,65	în	2,55	în	2,58	în	3,47
3	în	2,67	în	2,44	în	2,50	și	2,39	și	2,58	și	3,07
4	să	1,62	să	1,87	cu	1,75	să	1,81	o	1,94	a	2,35
5	a	1,47	la	1,52	o	1,62	o	1,73	să	1,69	la	1,52
6	la	1,46	cu	1,50	a	1,47	la	1,55	cu	1,52	se	1,46
7	se	1,39	pe	1,45	la	1,43	cu	1,48	la	1,46	cu	1,21
8	cu	1,38	se	1,43	se	1,42	nu	1,41	se	1,44	care	1,17
9	o	1,30	o	1,41	pe	1,39	pe	1,39	pe	1,41	o	0,87
10	nu	1,28	nu	1,33	nu	1,37	se	1,35	nu	1,27	din	0,85
11	pe	1,27	a	1,17	să	1,33	un	1,18	un	1,25	pe	0,82
12	care	0,98	că	1,05	un	1,26	că	1,08	a	1,04	este	0,79
13	că	0,97	un	0,99	că	1,04	a	1,05	care	0,95	mai	0,75
14	mai	0,95	mai	0,97	lui	0,88	care	0,95	că	0,95	nu	0,73
15	din	0,91	din	0,94	mai	0,86	din	0,93	din	0,91	sau	0,71
16	un	0,87	care	0,89	din	0,86	ce	0,85	mai	0,84	să	0,70
17	ce	0,66	ce	0,69	care	0,84	mai	0,84	ce	0,79	pentru	0,67
18	ca	0,60	ca	0,58	ioanide	0,74	lui	0,68	pentru	0,71	că	0,54
19	pentru	0,54	lui	0,54	era	0,66	era	0,63	ui	0,71	al	0,53
20	ui	0,49	dar	0,51	ce	0,64	pentru	0,63	era	0,64	un	0,50
21	dar	0,45	era	0,51	e	0,54	ca	0,55	charlie	0,59	prin	0,44
22	fi	0,42	pentru	0,48	ca	0,53	ei	0,53	ei	0,57	ca	0,43
23	este	0,42	fi	0,42	fi	0,52	dar	0,50	dar	0,54	fi	0,35
24	era	0,39	când	0,39	pomponescu	0,44	fi	0,49	ca	0,53	sunt	0,33
25	sau	0,35	el	0,38	pentru	0,40	fi	0,43	ii	0,53	ale	0,32
26	e	0,34	e	0,37	când	0,35	ei	0,42	ea	0,51	poate	0,29
27	el	0,34	am	0,35	el	0,33	ea	0,38	el	0,50	sa	0,29
28	al	0,33	ei	0,32	prin	0,27	când	0,34	fi	0,45	au	0,28
29	când	0,33	nici	0,30	am	0,26	nici	0,33	kurtz	0,36	ce	0,27
30	ei	0,29	ii	0,29	după	0,26	cum	0,31	când	0,34	art	0,27
31	am	0,28	mă	0,28	1	0,26	e	0,30	nici	0,32	fost	0,24
32	nici	0,28	cum	0,28	al	0,26	aSa	0,30	cum	0,30	după	0,24
33	prin	0,28	sau	0,25	nici	0,25	dacă	0,29	Si	0,28	dacă	0,21
34	sa	0,26	fost	0,25	ară	0,25	il	0,29	il	0,28	c	0,19
35	sunt	3,25	după	0,24	i	0,25	charlie	0,29	a	0,28	când	0,19

36	cum	0,25	sa	0,24	avea	0,24	fost	0,28	aSa	0,26	m	0,18
37	fost	0,25	dacă	0,24	ar	0,23	barley	0,28	dacă	0,25	unei	0,17
38	dacă	0,24	al	0,24	dacă	0,23	al	0,27	e	0,25	cele	0,17
39	după	0,24	ea	0,23	gaitany	0,22	spuse	0,26	oseph	0,24	pot	0,16
40	au	0,23	asa	0,22	însă	0,22	iar	0,26	sau	0,24	are	0,16
41	ii	0,23	il	0,22	foarte	0,21	îsi	0,26	ai	0,23	penală	0,16
42	mă	0,21	îsi	0,21	spre	0,20	este	0,26	după	0,23	trebuie	0,16
43	ea	0,21	este	0,21	ei	0,20	după	0,25	te	0,22	această	0,16
44	iar	0,19	au	0,21	asa	0,20	am	0,25	sa	0,22	iui	0,16
45	poate	0,19	sunt	0,20	sunt	0,19	sau	0,24	fără	0,22	acest	0,16
46	asa	0,18	iar	0,20	cum	0,19	ai	0,24	fost	0,22	iar	0,15
47	ar	0,18	fără	0,19	dar	0,19	ar	0,24	ar	0,22	lor	0,15
48	fără	0,18	prin	0,19	hagienuş	0,19	sa	0,23	este	0,21	numai	0,15
49	îsi	0,17	ar	0,19	sau	0,18	te	0,20	iar	0,21	dar	0,15
50	il	0,17	le	0,18	fost	0,18	avea	0,20	le	0,21	mare	0,15
51	le	0,17	asta	0,18	toate	0,18	le	0,20	spuse	0,21	cel	0,14
52	ale	0,17	tot	0,18	este	0,18	leamas	0,20	apoi	0,20	unor	0,14
53	toate	0,17	eu	0,18	sa	0,18	timp	0,20	timp	0,20	fie	0,14
54	va	0,16	acum	0,17	îşi	0,17	apoi	0,19	lor	0,19	va	0,14
55	decât	0,16	până	0,17	qonzalv	0,17	au	0,18	săi	0,19	între	0,14

Un alt rezultat experimental interesant este identificarea unui număr de 162 de cuvinte care se regăsesc în toate cele 93 de cărți ce alcătuiesc corpusul (fie că este vorba de literatură, medicină, drept etc). Deși sunt doar 162, aceste cuvinte au o pondere importantă în textul global #CMG acoperind circa 45% din totalul celor 8806433 cuvinte. Aceste cuvinte comune împreună cu rangul lor în ierarhie și frecvențele lor relative în întreg textul #CMG sunt conținute în Tabelul 2.

Tabel 2

Lista cuvintelor comune în toate cele 93 de cărți
1. Cuvânt; 2. Rangul cuvântului în ierarhia frecvențelor relative în textul mixt global, #CMG; 3. Frecvența relativă a cuvântului în textul mixt global, #CMG (în %)

	1	2	3	1	2	3	1	2	3	1	2	3
de		1	4,10	iar	44	0,19	unei	93	0,10	sar	185	0,0
și		2	3,20	poate	45	0,19	atunci	94	0,10	una	187	0,0
în		3	2,67	asa	46	0,18	două	95	0,10	început	188	0,0
să		4	1,62	ar	47	0,18	doar	96	0,10	încât	193	0,0
a		5	1,47	fără	48	0,18	dintre	100	0,10	alte	196	0,0
la		6	1,46	îsi	49	0,17	are	101	0,10	acestea	198	0,0
se		7	1,39	il	50	0,17	face	102	0,10	facă	199	0,0
cu		8	1,38	le	51	0,17	sub	104	0,09	altă	200	0,0
o		9	1,30	ale	52	0,17	nimic	106	0,09	aceiași	204	0,0

	1,28	toate	53	0,17	feie	107	0,09	desi	206	0,04	
pe	11	<u>1,27</u>	va	54	0,16	ia	108	0,09	fac	213	0,04
care	12	<u>0,98</u>	decât	55	0,16	<u>puțin</u>	109	0,09	<u>printre</u>	220	0,04
	13	0,97	tot	56	0,16	între	110	0,09	<u>pare</u>	224	0,04
mai	14	0,95	[lor	57	0,16	intru n	111	0,09	<u>partea</u>	225	0,04
din	15	<u>0,91</u>	<u>spre</u>	58	0,15	cea	112	0,09	afară	226	0,04
un	16	<u>0,87</u>	<u>pana</u>	59	0,15		113	0,08	sus	240	0,04
ce	17	0,66	chiar	60	0,15	săi	116	0,08	<u>faptul</u>	246	0,03
ca	18	0,60	mult	61	0,14	aceea	117	0,08	locul	252	0,03
<u>pentru</u>	19	0,54	cel	63	0,14	ci	119	0,08	adevărat	260	0,03
lui		<u>0,49</u>	fie	65	0,14	fată	126	0,08	tuturor		
dar		0,45	ne	66	0,14	unul	127	0,08	<u>măcar</u>	<u>266</u>	<u>0,03</u>
		0,42	ai	67	0,14	astfel	128	0,08	primul	268	0,03
este	23	0,42	acum	68	0,14	<u>parte</u>	129	0,08	<u>aceeași</u>	<u>275</u>	0,03
<u>era</u>		0,39	trebuie	<u>69</u>	<u>0,14</u>	înainte	132	0,07	altfel	277	0,03
sau	25	0,35	cele	70	0,13	<u>pot</u>	134	0,07	noua	299	0,03
	26	0,34	numai	72	0,13	ele	138		acela	302	0,03
	27	0,34	<u>despre</u>	73	0,12	totul	140	0,07	trebui	<u>307</u>	0,03
ai	28	0,33	avea	74	0,12	dată	141	0,07	dintro	330	0,03
când	29	0,33	atât	75	0,12	tot	143	0,07	dă	358	0,02
ei	30	0,29	această	76	0,12	loc	144	0,07	afla	364	0,02
nici	32	0,28	putea	78	0,12	fiecare	153	0,06	ramane	<u>371</u>	0,02
prin	33	0,28	unde	80	0,12	orice	155	0,06	alt	373	0,02
sa	34	0,26	mtro	81	0,11	<u>spune</u>	165	0,06	pus	377	0,02
sunt	35	0,25	acest	82	0,11	asemenea	166	0,06	întâi	387	0,02
cum	36	0,25	noi	84	0,11	sale	167	0,06	rând	397	0,02
fost	37	0,25	sai	86	0,11	acesta	168	0,06	alta	404	0,02
dacă	<u>38</u>	0,24	cat	87	0,11	lucru	169	0,06	<u>legătură</u>		
<u>după</u>	39	0,24		88	0,11	către	174	0,05		429	<u>0,02</u>
au	40			89	0,11	multe	175	0,05	rândul	601	0,01
n	41	0,23	ceva	91	0,11	celor	178	<u>0,05</u>			
ea	43	0,21	insa	92	0,10	totuși	184	0,05			

	1	201	401	
	CÂND	GAITTANY	TĂCU CĂCI ȘTIA	GAITTANY LIPSIT DE
1.	CÂND		CĂCI	LIPSIT...
2.	GAITTANY		ȘTIA	DE
200.		TĂCU		GAITTANY

Figura 1. 200 de mulțimi de date (cuvinte) în model statistic //cf. obținute prin eșantionare periodică a textului natural

Deplasând originea eșantionării în textul natural apar 200 de astfel de mulțimi de date experimentale, fiecare în parte de volum N , Fig. 1.

Fiecare mulțime de N observații astfel obținută satisface modelul stației *i.i.d.*, model necesar în aplicarea inferențelor statistice utilizate. Independența este asigurată de mărimea perioadei de eșantionare; distribuția identică este un rezultat al ipotezei de staționaritate a limbii naturale.

Acceptând ipoteza de staționaritate a limbii, toate cele 200 de mulțimi de date experimentale (compatibile cu modelul *i.i.d.*) extrase din textul natural conform Fig. 1, trebuie să conțină aceeași informație despre probabilitatea cuvântului investigat (oricare ar fi acesta).

Atenție însă, aceste mulțimi de date nu sunt independente între ele.

Un prim obiectiv al studiului nostru a fost de a vedea dacă într-adevăr cele 200 mulțimi de date confirmă sau nu aceeași probabilitate p teoretică (necunoscută) a cuvântului investigat.

Un răspuns afirmativ ne-ar permite să obținem un model matematic pentru sursa de informație de cuvinte asociată limbii române. Pentru a da un răspuns extins o procedură statistică pe care am dezvoltat-o în [3]-[8] pentru $1/n$ -grame. În această procedură cele 200 de mulțimi de date experimentale se compară între ele aplicând repetat un test statistic al ipotezei că probabilitatea aparține unui interval dat, vezi Anexa 1.

Menționăm că nu am putut face o comparație pe baza unui test mai sofisticat, anume acela privind egalitatea între două probabilități, întrucât mulțimile de date care se compară nu sunt independente între ele.

Procedura a permis în final determinarea unui interval de încredere statistică optim care a fost denumit în continuare "*reprezentativ*" pentru cuvântul urmărit și textul natural. Simultan a apărut și mulțimea de date experimentale *i.i.d.*

2. Descrierea structurii statistice de cuvinte. Studiu bazat pe multiple intervale de încredere statistică

Fie un text natural **considerat ca succesiune de cuvinte** pe care îl eșantionăm cu o perioadă suficient de mare astfel încât să rupem practic dependența dintre observațiile succesive. Inițial în investigația noastră statistică am considerat această perioadă ca fiind de 200 cuvinte. La fiecare moment de eșantionare am înregistrat observația făcută (cuvântul respectiv), conform Fig. 1. Mulțimea de date obținute în acest fel conține N -cuvinte, unde $N=L_c/200$ iar L_c este lungimea textului în cuvinte.

"reprezentativă" pentru cuvântul respectiv și textul natural, mulțime ce va fi folosită în comparații matematice între texte naturale.

2.1. Intervale de încredere statistică "reprezentative" pentru probabilitățile cuvintelor. Metodă de determinare și rezultate experimentale

Scopul acestui subcapitol este de a determina probabilitatea p a unui cuvânt urmărit.

Fie m_l numărul de apariții ale cuvântului în mulțimea l de date experimentale *i.i.d.* de volum N , $l = 1+200$. (Aceste mulțimi sunt extrase din textul natural conform Fig. 1)

Aplicând teoria estimării, fiecare din cele 200 de mulțimi de date conduce la o estimatie $p_i = m_i/N$ a probabilității p necunoscute și la un interval de încredere statistică al probabilității $p = (p_i - z_{\alpha/2} \sqrt{p_i(1-p_i)}, p_i + z_{\alpha/2} \sqrt{p_i(1-p_i)})$, $i = 1+200$. Considerând N suficient de mare astfel încât condiția de Moivre - Laplace să fie satisfăcută, $Np(1-p) \gg 1$, limitele intervalului de încredere statistică (inferioară și superioară) se calculează conform relației (1), [9], [10]:

$$P_{li} = P_i - z_{\alpha/2} \sqrt{P_i(1-P_i)} \quad P_{2i} = P_i + z_{\alpha/2} \sqrt{P_i(1-P_i)} \quad (1)$$

unde $z_{\alpha/2}$ este $\alpha/2$ cuantila legii normale de medie 0 și dispersie 1. În determinările noastre experimentale am lucrat cu un nivel de încredere statistică de 95%; rezultă $z_{\alpha/2} = 1.96$.

Cu alte cuvinte putem spune că probabilitatea adevărată p se află în intervalul $[P_{li}, P_{2i}]$ cu o încredere statistică egală cu 0,95.

Într-o primă etapă a analizei noastre, pentru un anumit eveniment urmărit (apariția unui cuvânt), s-au folosit următoarele mărimi (a se vedea Fig. 2):

- frecvența relativă a cuvântului pe întreg textul natural considerat (ceea ce înseamnă măsurare din date corelate); p^* este raportul între numărul de apariții ale cuvântului în textul natural și lungimea L_c a textului respectiv (numărul total de cuvinte). Se observă că p^* este media aritmetică a celor 200 de estimatii. Subliniem că p^* este o mărime importantă pentru orice experimentator.

$$P_{min} = \min p_i, \quad i = 1+200 \quad \text{valoarea minimă a estimatiilor};$$

$$P_{max} = \max p_i, \quad i = 1+200 \quad \text{- valoarea maximă a estimatiilor};$$

$$A_m = \max p_i, \quad i = 1+200 \quad \text{- reuniunea celor 200 de intervale}$$

de încredere statistică;

$$A_m^o = \max p_i - \min p_i, \quad i = 1+200 \quad \text{diferența maximă între două estimatii (intervalul de împrăștiere al estimatiilor);}$$

$$\delta_m = \max |p_i - p^*|, \quad i = 1+200 \quad \text{- diferența maximă între estimatiile } p_i \text{ și } p^* \text{ și}$$

frecvența relativă p^* ;

$$S_m = \min p_i, \quad i = 1+200 \quad \text{- diferența minimă între estimatiile } p_i \text{ și } p^* \text{ și}$$

frecvența relativă p^* .

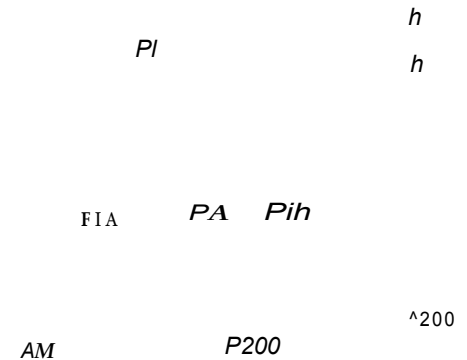


Figura 2. Mărimi utilizate în obținerea intervalului de încredere statistică "reprezentativ" pentru probabilitate

Următoarele întrebări (probleme) au ghidat analiza noastră teoretică experimentală:

1. Cât de largi sunt intervalele A_m , S_m și A_m^o ? Intervalele K_m și S_m sunt importante în analiza împrăștierei estimatiilor în jurul valorii p^* . Intervalul A_m ne dă o idee despre cel mai mare interval în care s-ar afla p , probabilitatea adevărată, bănuită că există.
2. Există valori p_i foarte apropiate de p^* și cât de apropiate?

Pentru a răspunde la această întrebare a fost urmărită experimental mărimea δ_m . S_m conduce la estimăția p_i ; care este cea mai apropiată de p^* , estimăție care va fi în continuare notată cu $p_{i,}$. S-a notat cu A intervalul de încredere statistică asociat estimăției $p_{i,}$ conform relației (1).

- Cât de multe intervale de încredere statistică // conțin (îmbracă) p^* ?** Prin presupunerea de staționaritate, ne așteptăm ca un mare număr de intervale de încredere statistică // să se intersecteze, conținându-l în același timp pe p^* . Nu ne așteptăm la o proporție de 95%, întrucât cele 200 de mulțimi de date U_d nu sunt independente între ele.
- Putem găsi un interval de încredere statistică pentru probabilitatea adevărată p , interval care să fie acceptat de toate cele 200 de mulțimi de date experimentale? Se pot găsi mai multe astfel de intervale? Dacă ipoteza de staționaritate este adevărată atunci astfel de intervale trebuie să existe. Dacă intervalul A (definit mai sus în întrebarea 2) este unul dintre aceste intervale, atunci el va fi preferat de experimentator și va fi considerat ca "*reprezentativ*" pentru probabilitatea cuvântului și textul analizat.

Pentru a înțelege metoda dezvoltată de noi care răspunde la aceste întrebări și care conduce la obținerea intervalului de încredere statistică "*reprezentativ*" exemplificăm pentru cuvântul DE în corpusul mixt global #CMG.

În Tabelul 3 cele 200 de rânduri corespund celor 200 de mulțimi de date experimentale *i.i.d.* fiind explicitate atât estimățiile p_i ; cât și intervalele de încredere statistică //, $l = 1 + 200$. Succesiv, fiecare interval // a fost considerat ca interval de referință și am aplicat 199 de teste ale ipotezei că probabilitatea aparține intervalului menționat, test descris în Anexa 1. Fiecare dintre cele 199 de teste este aplicat pe o singură mulțime de date experimentală. (Numărul 199 se explică prin faptul că nu se testează și mulțimea care a produs intervalul de referință.)

În primul rând al tabelului 3, intervalul $l_1 = (3,89; 4,27) \times 10^{-2}$ este intervalul de referință față de care se aplică testele de apartenență a probabilității. Se testează dacă probabilitatea cuvântului DE aparține sau nu intervalului l_j pe baza unei singure mulțimi de date *i.i.d.* aceasta înseamnă că verificăm succesiv fiecare din restul mulțimilor de date, anume $l = 2 + 200$. Acceptarea ipotezei că probabilitatea cuvântului DE aparține intervalului l_1 este marcată cu "DA" pentru respectiva mulțime de date (Tabelul 3, rândul 1). În caz contrar, pe poziția respectivă este completat "NU". Numărul total de mulțimi de date care trec testele este conținut în ultima coloană din dreapta. Această procedură se repetă alegând ca referință pe rând toate cele 200 de intervale de încredere statistică //, $l = 1 + 200$.

Determinarea intervalului de încredere statistică "*reprezentativ*" A pentru M I L » cuvântului DE în corpusul mixt global #CMG. Este îngroșat rândul 3 care corespunde intervalului "*reprezentativ*" A

i	\hat{p}_i ($\times 10^{-2}$)	$I_i (\times 10^{-2})$		Mulțime i																	Total
		$p_{1,i}$	$p_{2,i}$	1	2	3	...	99	...	97	...	200	DA								
1	4,08	3,89	4,27		DA	DA			DA				DA			DA	199				
2	4,06	3,87	4,25	DA		DA			DA				DA			DA	199				
3	4,10	3,91	4,29	DA	DA				DA				DA			DA	199				
...																					
94	4,31	4,12	4,51	DA	DA	DA										NU	DA	182			
...																					
99	3,78	3,60	3,96	DA	DA	DA							NU				DA	12			
...																					
200	4,18	3,99	4,37	DA	DA	DA							NU			DA		19			

Pentru cuvântul DE s-a obținut $p^* = 0,040986$, iar estimăția cea mai apropiată de p^* a fost $p_i = p_{i,} = 0,040992$, $S_m/p^* = 0,0002$. Pentru estimăția se obține intervalul de încredere statistică 95% $A = (0,0391; 0,0429)$. Din rândul se observă că intervalul A trece toate cele 199 de teste ale ipotezei probabilitatea cuvântului DE este cuprinsă în interiorul său. Sunt multe intervale care au compatibilitate cu toate mulțimile de date *i.i.d.* (în ultima coloană numărul 199 a apărut de 101 ori). Dintre aceste 101 intervale am ales $A = (0,0391; 0,0429)$ ca fiind interval de încredere statistică 95% *reprezentativ* pentru probabilitatea cuvântului DE întrucât este ușor de determinat de oricare experimentator. Mulțimea de date *i.i.d.* specificată de indicele $i = 3$ va fi numită mulțime de date "*reprezentativă*" pentru cuvântul DE în corpusul #CMG.

Tabelul 4 conține informații despre elementele analizei pentru primele zece cuvinte din ierarhia frecvențelor relative în corpusul #CMG. Exemplificăm pentru cuvântul DE care având frecvența relativă $p^* = 4,10 \times 10^{-2}$ este pe primul loc în ierarhie. Valoarea p^* este cuprinsă în $N(p^*) = 192$ de intervale de încredere statistică din cele 200 considerate, (coloana 3); reuniunea celor 200 intervale de încredere statistică raportată la p^* este $A^c/p^* = 22,24 \times 10^{-2}$, (coloana 4); diferența maximă între două estimății raportată la p^* este $A_m^c/p^* = 13,08 \times 10^{-2}$, (coloana 5); diferența maximă între o estimăție p_i și p^* raportată la p^* este $< 5_m/p^* = 7,80 \times 10^{-2}$, (coloana 6); diferența minimă între o estimăție p_i și p^* este $> 5_m/p^* = 1,20 \times 10^{-2}$, (coloana 7).

raportată la p^* este $S_m/p^* = 0.02 \times 10^{-2}$, (coloana 7); lărgimea intervalului de încredere statistică "reprezentativ" A raportată la p^* este $A/p^* = 9,23 \times 10^{-2}$, (coloana 8); există $N(A) = 101$ intervale de încredere statistică la fel de bune ca intervalul A "reprezentativ", (coloana 9). Aceste $N(A)$ intervale sunt confirmate de toate cele 199 de teste de apartenență a probabilității la interval, prin care s-a făcut verificarea staționarității.

Numărul relativ mare de intervale de încredere statistică confirmate practic de toate cele 199 de teste de apartenență a probabilității la interval - $N(A)$ din coloana 9 a tabelului 4 - este o susținere a ideii de staționaritate.

Tabel 4
Rezultate numerice privind mărimile din Fig. 2 pentru cele mai frecvente 10 cuvinte în #CMG.

Cuvânt	p^*	$N(p^*)$	A_m/p^*		S_m/p^*		A/p^*	$N(A)$
1	2	3	4	5	6	7	8	9
DE	4,10	192	22,24	13,08	7,80	0,02	9,23	101
SI	3,20	198	25,42	14,91	8,77	0,04	10,45	102
IN	2,67	194	24,88	13,43	7,09	0,04	11,43	172
SA	1,62	189	35,74	21,02	11,40	0,05	14,69	122
A	1,47	191	37,11	21,66	11,40	0,05	15,42	124
LA	1,46	185	38,07	22,52	12,89	0,00	15,45	104
SE	1,39	190	39,97	24,27	13,92	0,02	15,89	79
CU	1,38	195	37,40	21,55	11,30	0,05	15,92	132
O	1,30	191	37,52	21,29	12,38	0,02	16,39	120
NU	1,28	189	40,20	23,62	12,78	0,01	16,54	120

Prin centralizarea acestor tipuri de rezultate pentru toate corpusurile analizate și pentru toate cuvintele pentru care s-a putut face analiza a rezultat Tabelul 5. Concret, studiul experimental a cuprins toate corpusurile prezentate în Introducere. Am putut aplica inferențele statistice doar pentru acele cuvinte pentru care am avut suficiente date; anume $Np^*(1-p^*) > 20$, unde N este volumul mulțimii de date *i.i.d.* (forma experimentală pentru condiția DeMoivre - Laplace). Cuvintele au fost sortate în ordine descrescătoare a frecvențelor de apariție p^* . Această sortare a permis organizarea studiului pe clase de frecvență. Am ales ca limite ale claselor următoarele valori: 5%, 2%, 1%, 0,5%, 0,2%, 0,1% și 0,05%.

În studiul nostru experimental în aproape toate situațiile (oricare cuvânt urmărit și orice corpus lingvistic investigat) am găsit o estimatie p_m practic egală cu p^* . Acest lucru se vede în Tabelul 5, coloana 8 urmărind raportul dintre S_m și p^* . Pentru toate situațiile analizate am obținut $S_m/p^* < 2,23\%$. Având în vedere

că studiul experimental a condus și la obținerea de intervale A "reprezentative" în toate situațiile analizate, rezultă că aceste intervale de încredere statistică 95% pot fi scrise sub forma:

$$A = (p_{1-\epsilon}; p_{1+\epsilon}) = p^* (1 \pm \epsilon), \quad \text{e. s. l. } 96 \times V(1-p^*) / (\# / > *) \quad (2)$$

ϵ , este eroarea relativă cu care se determină probabilitățile.

Exemplificăm citirea Tabelului 5 pentru corpusul #CMG și clasa a doua de frecvență. Există 8 cuvinte (coloana 3) care au frecvențele relative cuprinse între (0,01; 0,02). Aceste 8 cuvinte acoperă 11,17% (coloana 4) din totalul aparițiilor de cuvinte din #CMG, $L_c = 8806433$. Celelalte coloane, 5-9, conțin informații referitoare la mărimile din Fig. 2. Astfel coloana 9 conține raportul dintre lungimea intervalului A și p^* pentru cuvintele existente în clasa respectivă (limita minimă și maximă). Acest raport este practic dublul erorii relative, ϵ , în determinarea probabilității cuvântului; se observă o precizie relativ bună a determinărilor din această clasă, $\epsilon_c < 8.5 \times 10^{-2} = 17 \times 10^{-2} / 2$.

În total în #CMG au fost 194 = 3 + 8 + 8 + 24 + 59 + 92 cuvinte pentru care s-a putut determina intervalul A "reprezentativ". Deși cele 194 cuvinte reprezintă o mică pondere din totalul cuvintelor distincte posibile, ele acoperă 48,87% din $L_c = 8806433$, totalul aparițiilor de cuvinte în corpusul mixt global, #CMG.

Tabel 5

Rezultate experimentale organizate pe clase de frecvențe relative. Valorile

Clasa de frecvențe	Corpus	Nr.	Aco-perire	A_m/p^*	S_m/p	S_m/P^*	A/p^*
1	2	3	4	5	6	7	8
$2 \times 10^{-2} < p^* < 5 \times 10^{-2}$	#CMG	3	9,97	22-25	13-15	7-9	0,02-0,04
	#1JCMG	3	10,04	31-39	18-22	10-13	0,01-0,07
	#2JCMG	3	9,90	30-41	17-25	9-14	0,01-0,06
	#CLG	3	9,58	27-33	16-19	8-11	0,03-0,04
	#1JCLG	3	9,60	39-47	23-27	12-14	0,05-0,10
	#2JCLG	3	9,55	38-46	22-25	12-13	0,03-0,10
	#CŞG	4	13,76	52-83	27-52	14-27	0,05-0,37
$10^{-2} < p^* < 2 \times 10^{-2}$	#CMG	8	11,17	36-40	21-24	11-14	0,00-0,05
	#1JCMG	9	12,07	49-59	28-33	14-18	0,00-0,20
	#2JCMG	8	11,28	46-60	24-37	12-19	0,10-0,16
	#CLG	9	12,73	39-52	23-30	12-17	0,00-0,09
	#1JCLG	9	12,69	55-78	30-48	16-25	0,02-0,26
	#2JCLG	10	13,77	51-75	28-44	14-24	0,00-0,25
	#CŞG	4	5,36	110-130	67-78	34-48	0,24-0,59
	#CMG	8	6,48	42-61	23-36	12-20	0,00-0,16
#1JCMG	8	5,93	60-92	33-62	17-34	0,11-0,28	

$5 \times 10^{-3} \leq p^* < 10^{-2}$	#2JCMG	8	6,51	58-84	31-52	15-29	0,01-0,26	27-36
	#CLG	9	6,62	53-77	30-45	15-27	0,03-0,27	22-31
	#1JCLG	8	6,15	73-112	40-68	23-42	0,01-0,54	31-44
	#2JCLG	9	6,14	72-107	40-66	24-39	0,06-0,42	32-44
	#CŞG	12	8,46	135-186	71-117	36-64	0,35-1,45	58-78
$2 \times 10^{-3} \leq p^* < 5 \times 10^{-3}$	#CMG	24	7,36	59-103	31-63	16-33	0,01-0,43	27-41
	#1JCMG	24	7,17	90-158	51-100	27-59	0,01-0,82	39-60
	#2JCMG	23	7,05	91-143	49-87	25-45	0,04-0,84	38-60
	#CLG	25	7,02	74-123	42-75	21-40	0,05-0,73	32-50
	#1JCLG	23	6,87	116-186	65-116	37-66	0,02-1,33	46-71
	#2JCLG	23	6,40	105-186	57-115	29-66	0,02-1,20	48-71
$10^{-3} \leq p^* < 2 \times 10^{-3}$	#CŞG	2	0,87	198-198	111-112	69-72	1,08-2,14	82-83
	#CMG	59	7,78	92-156	46-97	26-61	0,02-1,09	43-61
	#1JCMG	57	7,40	135-224	74-141	41-80	0,02-2,07	61-87
	#2JCMG	61	8,13	135-217	75-134	40-77	0,04-2,23	61-87
	#CLG	58	7,74	114-190	62-120	33-67	0,01-1,41	51-73
	#1JCLG	33	5,21	158-231	81-143	42-85	0,07-1,84	73-91
	#2JCLG	32	5,06	160-220	85-129	46-82	0,04-2,11	72-91
$5 \times 10^{-4} \leq p^* < 10^{-3}$	#CŞG	0	0,00	-	-	-	-	-
	#CMG	92	6,11	134-224	72-142	37-82	0,02-2,16	62-90
	#1JCMG	8	0,72	185-221	101-134	52-76	0,30-1,89	88-90
	#2JCMG	9	0,81	194-242	111-150	57-95	0,25-1,89	88-91
	#CLG	52	4,11	162-224	88-131	45-82	0,05-1,95	73-90
	#1JCLG	0	0,00	-	-	-	-	-
	#2JCLG	0	0,00	-	-	-	-	-
Total	#CŞG	0	0,00	-	-	-	-	-
	#CMG	194	48,87					
	#1JCMG	109	43,33					
	#2JCMG	112	43,68					
	#CLG	156	47,80					
	#CŞG	22	28,45					

Tabelul 5 indică și precizia determinărilor (eroarea relativă e_r) pentru cuvintele analizate. Această precizie este relativ bună pentru determinările făcute pe corpusul mixt global #CMG (pentru cuvinte din primele patru clase de frecvență, $e_r < 20,5 \times 10^{-2} = 41 \times 10^{-2} / 2$).

Aplicând procedura descrisă în Cap. 2.1 pentru toate corpusurile lingvistice și pentru toate cuvintele care au satisfăcut condiția deMoivre - Laplace au rezultat probe în sprijinul ipotezei de staționaritate a limbii române scrise.

Intervalele "reprezentative" precum și mulțimile de date "reprezentative" determinate pentru un cuvânt anumit și textul natural considerat în continuare folosite în Cap. 2.2, pentru a analiza dacă putem vorbi de staționaritate matematic al sursei de cuvinte pentru limba ca ansamblu, pentru diverse clase ale limbii, pentru diverși autori etc.

Acuratețea în determinarea probabilității cuvintelor este dată

- încrederea statistică (95%);
- erorile relative, e_r , cu care s-au obținut intervalele A conform Tabelului 5;
- mărimea celor două tipuri de erori statistice care aparțin apartenență a probabilității la interval, întrucât acest tip de erori se bazează pe validarea intervalului A ca "reprezentativ".

În ceea ce privește testul de apartenență a probabilității la interval a fost aplicat pentru un prag statistic $\alpha = 0,05$. Întrucât testul a fost aplicat pe date (fapt pentru care A a fost validat ca "reprezentativ") este imposibil să se facă un control asupra mărimii (3, probabilitatea de a accepta date false). Acesta impune valori mici pentru li am avea nevoie de un corpus mare. În exemplul nostru, conform [6, Tabel 4], dacă se dorește $fi < 0,3$ și $5 = 10$ pentru a investiga cuvinte din primele patru clase de frecvență am nevoie de un corpus de circa 30 de milioane de cuvinte.

2.2 Comparații matematice între diverse texte naturale și structuri de cuvinte

Investigația noastră (privind staționaritatea) a fost completată cu comparații matematice privind probabilitățile cuvintelor, pe care le-am făcut pe baza următoarelor criterii:

- Se verifică dacă un același cuvânt are aceeași probabilitate în două texte naturale care se compară. Această comparație se face în continuare *comparație între cuvinte ca atare*.
- Se verifică dacă probabilitățile cuvintelor situate pe un anumit nivel în ierarhia frecvențelor relative din cele două texte sunt egale. În exemplul nostru, pe rangul 20 în corpusul literar global se află cuvântul UN, iar în corpusul științific cuvântul UN, vezi Tabelul 1. Între cele două domenii se va urmări dacă probabilitățile cuvintelor (DAR și UN) este aceeași. În cele ce urmează vom prezenta un criteriu *comparație pe baza rangului*.

Toate comparațiile matematice, atât pe baza criteriului a) cât și pe baza criteriului b) au fost făcute folosind următoarele teste statistice:

- T_1 - test al ipotezei ca probabilitatea aparține unui interval
- T_2 - test de egalitate între două probabilități, (Anexa 2)

Pentru fiecare din cele două texte naturale care se compară și pentru fiecare cuvânt investigat s-au determinat în prealabil intervalele "reprezentative" precum și mulțimile de date *i.i.d.* "reprezentative".

Când aplicăm testul 77, intervalul $(a;b)$ este intervalul "reprezentativ" A din primul text natural implicat în comparație, iar mulțimea $\{x_i, x_2, \dots, x_n\}$ de date experimentale *i.i.d.* este mulțimea de date "reprezentativă" din cel de-al doilea text natural. Testul a fost aplicat în ambele situații: corpus1 versus corpus2 și corpus2 versus corpus1, Tabel 6.

Când aplicăm testul 72 se considera pentru comparație cele două mulțimi de date *i.i.d.* "reprezentative" extrase din cele două texte naturale pentru cuvintele care se compară.

Toate testele au fost aplicate pentru un prag de semnificație statistică $\alpha=0,05$. Cu alte cuvinte, probabilitatea de a respinge date corecte este mai mică decât 0,05.

Tabel 6

**Comparații între texte naturale pe baza probabilității cuvintelor.
Coloanele 4-9 conțin numărul de cuvinte rejectate de testele statistice**

Texte comparate		Nr.	Comparație între cuvinte ca atare		Comparație pe baza rangului			
Corpus 1	Corpus 2		Test 77		Test 72	Test 77		Test 72
7	2		1 versus 2	2 versus 1	6	1 versus 2	2 versus 1	9
#1JCLG	#2JCLG	72	0	0	0	0	0	0
#1JCMG	#2JCMG	104	0	0	0	0	0	0
#CLG	#CSG	22	10	18	13	1	16	10

Rezultatele experimentale sunt sintetizate în Tabelul 6. Comparațiile făcute în cadrul domeniului literar, când se compară cele două jumătăți de corpus între ele (**#1JCLG** și **#2JCLG**) nu indică diferențe între probabilități indiferent de testul utilizat (77 sau 72) sau de criteriul utilizat (comparații pe baza aceluiași cuvânt sau pe baza aceluiași rang).

Același rezultat s-a obținut și când s-au comparat cele două jumătăți ale corpusului mixt global, **#1JCMG** și **#2JCMG**.

Exemplificăm în continuare modul de citire al Tabelului 6.

Primele două coloane conțin corpusurile care se compară între ele.

Coloana 3 indică numărul de cuvinte investigate în comparații (care au îndeplinit condiția $Np^*(1-p^*) > 20$ în ambele texte care se compară).

Rezultatele din coloanele 4, 5 și 6 au fost obținute aplicând criteriul comparațiilor "cuvintelor ca atare".

Coloanele 4 și 5 arată câte cuvinte nu au trecut testul 77 de apartenență probabilistică la interval. Coloana 4 se referă la situația când intervalul fix $[a \setminus b)$ este intervalul A "reprezentativ" din primul corpus al comparației, iar mulțimea de date *i.i.d.* supusă testului este mulțimea Ud "reprezentativă" din al doilea corpus. Similar, în coloana 5: intervalul fix $(a \setminus b)$ este intervalul A "reprezentativ" din doilea corpus al comparației, iar mulțimea de date *i.i.d.* supusă testului este mulțimea *i.i.d.* "reprezentativă" din primul corpus.

Coloana 6 conține numărul de cuvinte care sunt rejectate de testul 72 de egalitate între probabilități.

Coloanele 7, 8 și 9 conțin același tip de informație specificat în coloanele 5 și 6, cu diferența că de această dată se compară cuvintele care ocupă același rang în loc de cuvintele "ca atare".

Când se compară domenii diferite, spre exemplu literar și științific, apar multe diferențe marcate de ambele teste 77 și 72 și de cele două criterii de comparație.

Rezultatele comparațiilor puntează unele diferențe între domeniile literar și științific. Testele nu au indicat diferențe când s-au comparat corpusuri organizate după aceeași reguli (jumătățile corpusului mixt global între ele sau jumătățile corpusului literar global între ele); reamintim că atât corpusul mixt global cât și corpusul literar global au fost obținute prin concatenarea aleatoare a cărților respective.

3. Legea lui Zipf. Studiu experimental

Legea lui Zipf (prezentată în Cap. 1 și în Cap. 2) și analiza de staționaritate din Cap. 2) au constituit o bază de plecare pentru studiul nostru experimental asupra legii lui Zipf. În lingvistică legea lui Zipf este una din cele mai cunoscute dependențe rang - frecvență. (Aceste dependențe rang - frecvență au fost observate de-a lungul timpului și în diverse alte domenii: economie, fizică, biologie, demografie etc. [11], [12].) Obiectivul acestui capitol este de a stabili dacă și în ce măsură (cu ce acuratețe) limba română scriasă satisface legea lui Zipf.

Fie un text (corpus) având o lungime de L cuvinte, dintre care N_k sunt k cuvinte distincte. Aceste N_k cuvinte se sortează într-o listă în ordine descrescătoare după numărul de apariții în textul natural. Se notează cu k rangul unui cuvânt în listă și cu $f(k)$ frecvența relativă a acestuia (numărul de apariții raportat la L): $f(k) = N_k / L$. (Altfel spus, $f(k)$ este de tipul p^* din capitolele precedente). Legea lui Zipf afirmă că produsul dintre rang și frecvența relativă este constant, [11] - [14]:

$$kf(k) = A$$

Se observă că membrul stâng al ecuației (3) corespunde realității fiind vorba de măsurători efectuate pe texte naturale în timp ce membrul drept corespunde modelului teoretic presupus.

Este știut din considerații privind alte limbi naturale că legea Zipf, apreciată ca foarte simplă și foarte atractivă, funcționează cu aproximație pentru o plajă limitată de ranguri, anume nu prea mici și nu prea mari. Astfel un prim pas al studiului nostru teoretic și experimental a fost să reprezentăm grafic dependența rang - frecvență pe tot corpusul de care am dispus (corpusul mixt global, #CMG). Fig. 3 prezintă această dependență la scară logaritmică ($f(k)$ versus k). La o primă vedere am putea spune că mărimea A din (3) este aproximativ constantă pentru un interval de ranguri $k \in [k_{min}; k_{max}]$ de $k_{min} > 50$. Am limitat studiul la acele ranguri pentru care numărul de apariții ale cuvintelor a fost mai mare decât 50 pentru a beneficia de rezultatele anterioare privind studiul de staționaritate prezentat în Cap. 2. Aceasta face ca rangul k_{max} să depindă de corpusul analizat.

Legea lui Zipf este descrisă în numeroase referințe dintre care în limba română menționăm în special [13] și [14]. Capitolul de față urmărește determinarea constantei legii atât pe corpusul de ansamblu, #CMG, cât și pe diverse texte naturale (grupate după autori sau pe subdomenii ale limbii). Se analizează și în ce măsură comportamentul real se abate de la cel teoretic.

3.1 Elemente teoretice

3.1.1 Determinarea parametrului legii Zipf prin minimizarea erorii pătratice

Presupunând valabilitatea legii Zipf pentru rangurile $k \in [k_{min}; k_{max}]$ ne-am propus să determinăm mărimea A din condiția de minimizare pe acest interval a următoarei funcții (suma pătratelor erorilor):

$$k \sum_{k=k_{min}}^{k_{max}} (f(k) - \frac{A}{k})^2 \quad (4)$$

Derivând funcția $g(A)$ și egalând cu 0 se obține valoarea mărimii A corespunzând minimului:

$$A = \frac{\sum_{k=k_{min}}^{k_{max}} \frac{1}{k^2}}{\sum_{k=k_{min}}^{k_{max}} \frac{1}{k}} \quad (5)$$

Valorile k_{min} și k_{max} sunt la dispoziția experimentatorului. Pentru o evaluare a acurateții cu care limba naturală verifică legea lui Zipf definim următoarele tipuri de erori:

- e_s , suma pătratelor erorilor pe intervalul $k \in [k_{min}; k_{max}]$ și forma ei normalată, $e_{s,n}$:

$$e_{s,n} = \frac{e_s}{\sum_{k=k_{min}}^{k_{max}} \frac{1}{k^2}} \quad (6)$$

$$M = \frac{e_s}{\sum_{k=k_{min}}^{k_{max}} \frac{1}{k^2}} \quad (7)$$

3.1.2 Determinarea parametrului legii lui Zipf considerând cazul ideal

Dacă acceptăm legea lui Zipf ca fiind corectă pe întreg domeniul de ranguri $k \in [1; N_c]$, atunci valoarea constantei A se determină prin raționamentul descris în [13], [14]:

$$c + \ln N_c$$

unde c este constanta lui Euler, egală cu 0,577215 și $N_c > 50$.

Observăm că mărimea A calculată cu relația (8) nu depinde decât de numărul A_{fc} de cuvinte distincte din textul analizat. Prin urmare sunt de așteptat unele diferențe între evaluările mărimii A pe baza datelor experimentale cu relația (5) și cazul ideal, pur teoretic, relația (8).

3.1.3 Corolar al legii lui Zipf

Rezultatele experimentale cuprind și verificarea unui corolar al legii Zipf care se referă la determinarea cotei părți, L_s/L_c , pe care o acoperă cele s cuvinte s cuvinte într-un text de lungime L_c , [13], [14].

$$\frac{L_s}{L_c} = \frac{c + \ln s}{c + \ln N_c} \quad (9)$$

Relația (9) este valabilă pentru un număr de cuvinte $s > 50$.

Observăm că valoarea raportului L_s/L_c nu depinde de mărimea A . Diferențele existente între diversele moduri de evaluare ale mărimii A nu vor influența acest raport. În consecință ne așteptăm la o bună verificare experimentală a acestui corolar.

3.2. Rezultate experimentale și concluzii

Analiza experimentală a legii lui Zipf a început cu corpusul global #CMG (vezi Fig. 3) și a continuat pentru comparație cu o serie de texte naturale incluse în acestă (prezentate în Introducere). Rezultatele experimentale sunt concentrate în Tabelul 7. Pentru fiecare text analizat Tabelul 7 conține în coloanele 2 și 3 numărul total de cuvinte L_c și numărul cuvintelor distincte $7V_c$. În toate textele analiza

s-au investigat toate cuvintele cu număr de apariții mai mare decât 50; acesta determină rangul L_{max} , corespunzător fiecărui text analizat (coloana 4). $\&_{min}$ diferă de la text la text; $\&_{min}$ este ales întotdeauna 51. Pentru acest interval de rangurile $L^{min^{max}}$ s-a determinat cu relația (5) mărimea A cuprinsă în coloana 5. Coloanele 6 - 9 conțin rezultatele numerice calculate cu relațiile (6) și (7) unde mărimea A este cea din coloana 5 (determinată din textul natural respectiv). Coloana 9 conține rangul k_w pentru care s-a obținut eroarea relativă maximă e_w .

Ne-am pus problema și dacă mărimea $A = 0,0909$ determinată pentru corpusul mixt global **#CMG**, ar putea fi acceptată drept referință pentru limba română. De aceea coloanele 10-13 conțin succesiv mărimile din relațiile (6) și (7) unde $A = 0,0909$ pentru toate textele naturale analizate. Eroarea relativă maximă e_w este însoțită de rangul corespunzător, k_w .

Tabel 7

Studiu experimental al legii lui Zipf în limba română scrisă

Text	L_c		L_{max}	A	e			k_w	ξ		ϵM	$\ast M$	
	1	2	3	4	5	6	7	8	9	10	11	12	13
#CMG	8806433	202403	14543	9,09	0,36	0,22	9,81	286	0,36	0,22	9,81	286	
#CLG	6255235	162124	10299	9,60	0,30	0,17	13,93	10136	0,81	0,50	15,43	149	
#1JCLG	3127618	116247	5568	9,58	0,26	0,14	10,53	136	0,72	0,44	16,41	136	
#2JCLG	3127617	116860	5529	9,74	0,29	0,16	10,02	122	1,11	0,69	17,86	122	
#1	226420	26943	466	9,81	0,37	0,22	9,07	68	1,28	0,88	16,81	173	
#2	121177	18457	260	10,15	0,15	0,09	8,21	256	1,95	1,48	20,83	256	
#3	88827	13768	190	10,07	0,20	0,14	10,17	186	1,60	1,33	22,07	186	
#2+#3	210004	25036	484	9,97	0,52	0,29	18,71	478	1,89	1,29	30,18	478	
#4	130743	18223	274	10,71	0,26	0,14	8,38	53	4,47	3,35	25,85	110	
#5	75698	10351	187	11,56	0,42	0,22	10,86	183	9,22	7,71	40,92	183	
#6	197889	23206	399	10,34	0,15	0,08	7,03	121	2,85	1,99	21,73	121	
#4+#5+#6	404330	33555	849	10,53	0,20	0,10	8,90	103	4,03	2,62	26,08	103	
#7	644794	49434	1195	10,03	0,35	0,18	10,80	477	2,04	1,30	21,49	77	

În Fig. 3 sunt prezentate pentru corpusul mixt global două traiectorii, una experimentală (cu 'o') și cea teoretică (cu '*') conform relației (3) cu parametrul $A = 0,0909$ din coloana 5, Tabelul 7. Se observă o bună concordanță a celor două curbe pentru $\ast e [W W]$ -

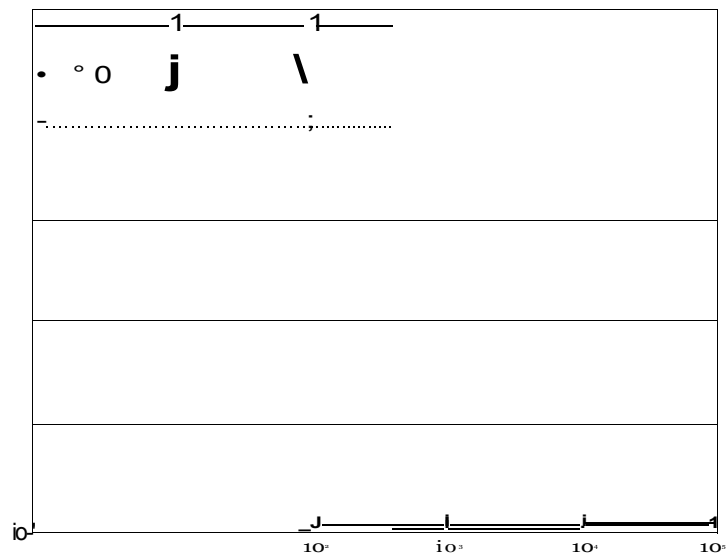


Figura 3. Dependența rang - frecvență relativă de apariție a cuvintelor în corpusul mixt global **#CMG** (scară logaritmică $f(k)$ versus k). Curba experimentală marcată cu 'o'; curba teoretică, relația (3) pentru $A = 0,0909$, marcată cu

în cazul ideal, pur teoretic, mărimea A poate fi determinată cu relația (8) pe baza coloanei 2 din Tabelul 7. Spre exemplu în corpusul global **#CMG**, unde au fost identificate $N_c = 202403$ cuvinte, $A = 0,0781$. În corpusul **#CLG**, pentru $N_c = 162124$ cuvinte distincte, aplicând relația (8) rezultă $A = 0,0795$.

Tabel 8

Valori teoretice, relația (9), și experimentale ale raportului l/L_c în corpusul literar global **#CLG**

X	0,1%	0,05%	0,01%
s	104	189	911
l/L_c (experimental)	43,69%	49,64%	64,38%
l/L_c (teoretic)	41,53%	46,28%	58,78%

Tabelul 8 conține date despre cota parte acoperită de cuvintele pentru care $f(k) > \frac{1}{A}$ unde $A = 0,1\%; 0,05\%; 0,01\%$, în textul literar global **#CLG**. S-a folosit relația (9) unde $N_c = 162124$, iar numărul de cuvinte s corespunzător pragului X

este conținut în linia 2 a Tabelului. Se observă o concordanță destul de bună între valorile teoretice și cele experimentale.

Notă: Din cele 189 cuvinte din corpusul literar global #CLG care au frecvența relativă mai mare decât 0,05%, doar 156 au îndeplinit condiția de Moivre-Laplace și au fost investigate cu control statistic apărând și în Tabelul 5.

Ca o remarcă finală legea lui Zipf poate fi considerată ca valabilă și pentru limba română pentru ranguri nu prea mici și nu prea mari, fapt susținut de Fig. 3 și datele din Tabelul 7.

4. Concluzii

Unul din principalele rezultate obținute în cadrul acestei lucrări este de a aduce probe noi privind staționaritatea limbii române scrise, de această dată pe baza structurii de cuvinte. (Ipoteza de staționaritate este inclusă în presupunerea generală conform căreia limbile naturale sunt lanțuri Markov multiple ergodice). Analiza staționarității s-a făcut prin extinderea unei metode dezvoltate de autori pentru studiul structurii statistice de m -grame (litere, digrame, trigrame, tetragrame). În consecință s-au putut obține probabilitățile cuvintelor cu intervale de încredere statistică 95% "reprezentative". Aceste intervale pe care le-am numit "reprezentative" au avut compatibilitate cu toate mulțimile de date *i.i.d.* obținute prin eșantionarea periodică a textului natural. Simultap au rezultat mulțimile de date *i.i.d.* "reprezentative" pentru cuvântul investigat și textul natural analizat.

O altă contribuție constă în procedura de comparație matematică între texte naturale facilitată de intervalul "reprezentativ" pentru probabilitate și de mulțimile de date *i.i.d.* "reprezentative". Comparațiile făcute între corpusuri organizate în aceeași manieră (literar *versus* literar sau mixt *versus* mixt) au întărit ideea de staționaritate a limbii și au confirmat modelul matematic prezentat anterior prin intervale de încredere statistică 95% "reprezentative" pentru probabilitățile cuvintelor. Au apărut unele diferențe între domeniile literar și științific.

Rezultatele experimentale dau un plus de semnificație frecvenței relative p^* , mărime de care orice experimentator este interesat. Acest plus de semnificație este datorat faptului că în toate situațiile analizate de noi (cuvânt sau text natural) am putut obține o estimatie a probabilității practic egală cu p^* , iar intervalul de încredere statistică asociat acestei estimatii a fost confirmat ca interval "reprezentativ" pentru probabilitate.

Lucrarea conține totodată și confirmarea valabilității pentru limba română scrisă a legii lui Zipf (lege de tip rang - frecvență) și a unui corolar al acesteia de interes lingvistic.

Autorii doresc să mulțumească D-lui dr. ing. Dan T. corepondent al Academiei Române, pentru sprijinul științific acordat studiului limbii române scrise. Autorii menționează, de asemenea, primite din partea D-lui Prof. dr. ing. Alexandru Șerbănescu a Tehnică Militară.

Referințe bibliografice

- [1] C. E. Shannon, "Prediction and Entropy of Printed English", Bell Voi. 30, pp. 50-64, January, 1951.
- [2] Adriana Vlad, A. Mitrea, "Estimating conditional probabilities of the statistical structure in printed Romanian", in Recent Advances in Language Technology, D. Tufiş, P. Andersen eds., Ed. Academic Press, 1997, (ISBN 973-27-0626-0), pp. 57-72; <http://www.raca.ro/vlad.html>.
- [3] Adriana Vlad, A. Mitrea, M. Mitrea, D. Popa, "Statistical modeling of the natural language stationarity based on the first approximation of the printed Romanian", Proc. VEXTAL'99 (Venezia per il trattamento della lingue), Ed. Unipress, (ISBN 88-8098-112-9), pp. 127-130, 1999, Venezia; <http://byron.cgm.unive.it/events/papers/vlad.pdf>.
- [4] Adriana Vlad, A. Mitrea, M. Mitrea, "Verifying Printed Romanian Stationarity Based on the Digram Statistical Structure", Proc. Romanian Academy, Series A, Voi. 1, No. 2, pp. 129-139, 2000.
- [5] Vlad Adriana, Mitrea A., Mitrea M., "Two frequency-rank models for printed Romanian", Procesamiento del Lenguaje Natural, 15 (Sociedad Espanol para el Procesamiento del Lenguaje Natural, 5948), pp. 153-160, Septiembre, 2000.
- [6] Adriana Vlad, A. Mitrea, M. Mitrea, "The trigram statistical structure of the Romanian", ROMJIST (Romanian Journal of Information Technology), Voi. 4, No. 3, pp. 353-372, 2001.
- [7] Adriana Vlad, A. Mitrea, M. Mitrea, "A Corpus - based approach to Accurately Printed Romanian Obeys Some Universal Laws", in Corpora: Corpus Linguistics and the Languages of the World, Rayson, T. McEnery eds., Lincom-Europa Publishing House, (ISBN 3-89586-872-8), pp. 153-165.
- [8] Adriana Vlad, A. Mitrea, M. Mitrea, "Estimating tetragram probabilities from multiple data samples from a natural text. Case study: printed Romanian", Proc. IPMU 2002 - The 9th International Conference on Intelligent Processing and Management of Uncertainty in Knowledge - July 2002, Annecy, France, pp. 1285-1292.

- [9] J. Devore, *Probability and Statistics for Engineering and the Sciences*, 2nd ed., Brooks/Cole Publishing Company, Monterey, California, 1987.
- [10] Adriana Vlad, B. Badea, M. Mitrea, *Metode Statistice în Prelucrarea Informației. Compendiu și Aplicații*, Ed. Metropol, București, 1999, (ISBN 973-562-104-5).
- [11] I. Kanter, D.A. Kessler, "Markov Processes: Linguistics and Zipfs Law", *Physical Review Letters*, Voi, 74, No. 22, pp. 4559-4562, May, 1995.
- [12] R. Gunther, L. Levitin, B. Schapiro, P. Wagner, "Zipfs Law and the Effect of Ranking on Probability Distributions", *Intl. J. of Theoretical Physics*, Voi. 35, No.2, pp. 395-417, 1996.
- [13] S. Marcus, Ed. Nicolau, S. Stați, *Introducere în lingvistica matematică*, Ed. Științifică, București, 1966.
- [14] M. Dinu, *Personalitatea limbii române*, Ed. Cartea Românească, București, 1996.

Fie $I = (a; b)$ un interval în care presupunem că se află probabilitatea p a unui eveniment urmărit. Dispunem de o mulțime de date experimentale $[x_1, x_2, \dots, x_n]$ care satisfac modelul statistic *i.i.d.*. Ne interesează dacă datele experimentale confirmă ipoteza că probabilitatea p aparține intervalului $I = (a; b)$, pentru un prag de semnificație statistică, α , ales.

Procedura de test este următoarea:

Se formulează cele două ipoteze statistice, ipoteza nulă H_0 și respectiv ipoteza alternativă H_1 :

$H_0: p$ aparține intervalului $(a; b); p \in (a; b)$

$H_1: p$ este în afara intervalului $(a; b); p \notin (a; b)$.

Se alege pragul de semnificație α (echivalent, nivelul de încredere statistică $1 - \alpha$). Se calculează estimăția $\hat{p} = m/N$, unde cu m s-a notat numărul de succese ale evenimentului în mulțimea de date $[x_1, \dots, x_n]$. Verificăm dacă estimăția \hat{p} se află sau nu în zona de acceptare a datelor. Regiunea de acceptare a datelor este un interval $(c_1; c_2)$ care include $(a; b)$. Intervalul $(c_1; c_2)$ se determină conform relației (10), [3]-[8]:

$$-y \ln a (\ln a)^i N^{\exp} \frac{(\ln a)^i}{2a(\ln a)N} \int_a^b \frac{1}{2nb(\ln b)N} \exp\left(\frac{x-b}{2b(\ln b)N}\right) dx = i - a \quad (10)$$

În relația (10) apar două funcții de densitate de probabilitate corespunzătoare legii normale: de medie a și dispersie $a/(a-b)/N$ și respectiv de medie b și dispersie $b/(b-a)/N$.

Ipoteza nulă H_0 va fi acceptată dacă și numai dacă estimăția \hat{p} aparține intervalului $(c_1; c_2)$. În caz contrar, $\hat{p} \notin (c_1; c_2)$, datele se resping ca fiind semnificative pentru pragul de semnificație α ales (se acceptă ipoteza H_1).

Ca în orice test statistic, pot să apară două tipuri de erori:

Eroarea de tipul (genul) I: Eroarea de a fi respinse date bune, adică să fie respinsă ipoteza H_0 când ea este corectă. Această situație apare atunci când estimăția \hat{p} nu satisface testul, adică $\hat{p} \notin (c_1; c_2)$, cu toate că probabilitatea adevărată p este în intervalul $(a; b)$. Probabilitatea acestui tip de eroare este mai mică decât α .

Anexa 2. Test de egalitate între două probabilități - T2

Eroarea de tipul (genul) II: Eroarea de a fi acceptate date false, adică să fie acceptată H_0 când ea este, de fapt, falsă. Această situație apare atunci când estimăția p satisface testul, $p \in (q; c_2]$, cu toate că probabilitatea adevărată p a evenimentului nu aparține intervalului $(a; b)$, $p \notin (a; b)$. Pentru a și N fixate probabilitatea acestui tip de eroare depinde de valoarea adevărată necunoscută p și se calculează cu relația:

$$P < \mathcal{L}(a; b) = \int_a^b \frac{p^{m+1}(1-p)^{N-m}}{2p(1-p)^N} dx$$

$\mathcal{L}(p)$ este mare atunci când p este la stânga lui a (sau la dreapta lui b), dar foarte aproape de a (respectiv de b). Practic, deranjante sunt situațiile în care $p < (1-b) > a$ sau $p > (1+a) - b$, cu toate că testul este trecut, adică $p \in (q; c_2]$. Valoarea α este determinată (prestabilită) de utilizator, în funcție de cât de mult deranjează această situație.

În studiul nostru asupra staționarității limbii române acest test a fost absolut necesar, vezi Cap. 2. A trebuit să stabilim dacă probabilitatea p a unui anumit cuvânt este aceeași când dispunem de diverse mulțimi de date experimentale extrase dintr-un același text natural (unde mulțimile sunt compatibile cu modelul statistic *i.i.d.*, dar nu sunt independente între ele). Testul a fost folosit și în comparații între texte naturale.

Disponem de două mulțimi de date experimentale în model statistic *i.i.d.*, de volume N_1 , respectiv N_2 . Notând cu m_1 numărul de succese (aparitii) ale unui eveniment în prima mulțime de date experimentale, estimăția probabilității este $\hat{p}_1 = (m_1/N_1)$. Similar, pentru a doua mulțime de date experimentale, estimăția probabilității este $\hat{p}_2 = (m_2/N_2)$. Urmărim să stabilim dacă cele două estimății \hat{p}_1 și \hat{p}_2 provin din aceeași probabilitate teoretică, respectiv $p_1 = p_2$.

Procedura de test este următoarea:

Se formulează cele două ipoteze statistice, ipoteza nulă H_0 și respectiv ipoteza alternativă H_1 :

H_0 : cele două probabilități teoretice sunt egale $p_1 = p_2$;

H_1 : cele două probabilități teoretice sunt diferite $p_1 \neq p_2$.

Se alege pragul de semnificație statistică α .

Se construiește o valoare de test z conform, [9], [10]:

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{N_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{N_2}}}$$

Valoarea z depinde de datele experimentale prin estimățiile \hat{p}_1 și \hat{p}_2 . În condițiile în care ipoteza H_0 este adevărată z provine dintr-o variabilă aleatoare a cărei lege de repartiție este practic legea normală standard.

Întrucât p_1 și p_2 sunt necunoscute, se consideră $p_1 = p_2 = \frac{m_1 + m_2}{N_1 + N_2}$

În aceste condiții valoarea de test z devine:

$$z = \frac{m_1/N_1 - m_2/N_2}{\sqrt{\frac{m_1 + m_2}{N_1 + N_2} \left(\frac{m_1}{N_1} + \frac{m_2}{N_2} - \frac{m_1 + m_2}{N_1 + N_2} \right)}} \quad (11)$$

Ipoteza nulă H_0 va fi acceptată (se va considera că probabilitățile sunt egale, $p_1 = p_2$) dacă și numai dacă $|z| < z_{\alpha/2}$ ($z_{\alpha/2}$ corespunde pragului de semnificație statistică α ales; am folosit $z_{\alpha/2} = 1,96$). În caz contrar se respinge ipoteza de egalitate a celor două probabilități pentru pragul de semnificație statistică α ales.

Această procedură de test a fost folosită când am comparat între ele diverse texte naturale.

Dezambiguizarea automată a cuvintelor din corpusuri paralele folosind echivalenții de traducere

Dan TUFİȘ

Institutul de Cercetări pentru Inteligența Artificială, Academia Română
Calea 13 Septembrie, nr. 13, 74311, sector 5, București

tufis@racai.ro

Rezumat

Corpusurile paralele constituie surse de cunoștințe extrem de valoroase, traducerea unui text reprezentând o succesiune de decizii lingvistice pe care traducătorul le ia în vederea asigurării unei transpuneri cât mai naturale și mai fidele a semnificației din textul sursă în textul tradus. Explicarea și extragerea acestor cunoștințe prin metode algoritmice, formalizarea și reutilizarea lor ulterioară constituie provocări ale inteligenței artificiale, subiecte de interes fierbinte în cercetarea actuală. Lucrarea prezintă o serie de contribuții în această direcție, prezentând mai întâi o metodă originală de identificare a echivalenților lexicali de traducere a cuvintelor dintr-un corpus paralel (extrăgând deci un dicționar multilingv) și apoi o metodă extrem de promițătoare pentru identificare automată a diferitelor sensuri ale cuvintelor polisemantice.

1. Motivații

Evoluția științifică și tehnologică este o sursă permanentă de formare a noi termeni sau a noi sensuri specializate pentru cuvintele existente. În domeniul lexicografiei multilinguale, păstrarea în actualitate a dicționarelor bi- și multilingve fără a apela la tehnologiile informatice, cu precădere cele din sfera ingineriei lingvistice, este aproape imposibilă. O serie de studii în domeniul traducerii automate au arătat că principalele probleme în *acceptabilitatea* traducerilor automate și cu atât mai mult al celor implicând pre- sau post-editare umană, nu sunt legate de probleme de natură sintactică (topică, acorduri, structură frazală) ci ele se regăsesc cu precădere în sfera lexicală, mai precis a semanticii lexicale. Evaluarea sistemelor existente de prelucrare a limbajului natural și mai ales a

celor de traducere automată (cu variantele ce presupun intervenția expertului uman) a condus la identificarea unor puncte sensibile, deficitare (pentru o interesantă trecere în revistă a problemelor privind evaluarea sistemelor de prelucrare a limbajului natural și a sistemelor de traducere a se vedea <http://www.isi.edu/natural-language/mteval/>). De pildă, traducerea greșită a unui cuvânt sau a unei expresii într-o frază perfectă din punct de vedere sintactic este percepută de imensa majoritate a consumatorilor de traduceri, în special de natură științifică, ca mult mai gravă decât un dezacord gramatical sau vreo altă abatere de la norma gramaticii. S-a invocat pe bună dreptate că dificultatea majoră a prelucrării automate a limbajului este rezolvarea ambiguităților lexicale, a omonimiilor și a polisemiei ce apar în orice text (scris sau vorbit). Spre deosebire de oameni, care de multe ori nici nu conștientizează aceste fenomene (ele sunt "obturate" fie de contextul textului, fie de cunoștințele de "bun simț" ale fiecărei persoane), procesoarele artificiale de limbaj natural încearcă rezolvarea ambiguităților printr-un proces inteligent de alegere, dintr-un spațiu al tuturor soluțiilor posibile în raport cu o modelare a limbajului, a soluției care respectă cel mai bine restricțiile modelului. Raportarea la modelul limbajului este esențială întrucât dificultatea procesului de prelucrare este cu atât mai mare cu cât modelul este mai complex: spațiul de căutare a soluțiilor poate crește exponențial, iar procesul decizional poate deveni nedeterminist sau de complexitate neoperațională.

Rezolvarea algoritmică eficientă a omografiei a cunoscut spectaculoase progrese în ultimii 10-15 ani, dar identificarea automată a sensului pe care îl are un anumit cuvânt polisemantic într-un context dat este încă o problemă nerezolvată satisfăcător și, prin urmare, un subiect "fierbinte" de cercetare. Problema identificării sensului cu care este utilizat un cuvânt este vitală în traducerea automată, întrucât se cunoaște faptul că de foarte multe ori un cuvânt polisemantic dintr-o limbă se traduce într-o altă limbă prin cuvinte diferite, în funcție de sensul considerat. Este interesant de remarcat că dacă un cuvânt polisemantic din limba sursă se traduce printr-un singur cuvânt polisemantic în limba țintă, sau altfel spus toate sensurile cuvântului de tradus se regăsesc în cuvântul reprezentând traducerea sa, necesitatea identificării sensului de utilizare al cuvântului sursă nu mai este obligatorie (cel puțin la nivelul fazei de transfer lexical) cu excepția situației în care diferitele sensuri ale cuvântului țintă se realizează lingvistic prin structuri de subcategorizare distincte.

În această lucrare vom prezenta în prima parte o metodă de extragere automată a echivalențelor de traducere și vom descrie apoi procedura de discriminare a sensurilor cuvintelor din corpusuri paralele pe baza echivalențelor de traducere.

2. Echivalenți de traducere

2.1. Noțiuni preliminare

O pereche de texte în două limbi diferite L_a și L_b , astfel încât unul reprezintă traducerea celuilalt constituie cea ce se numește un *bitext*. Un bitext suficient de mare constituie un corpus paralel. L_a și L_b se numesc echivalenți de traducere. Noțiunea de echivalență de traducere se poate rafina la niveluri subtextuale, de pildă la nivelul paragrafului, al propoziției sau chiar la nivel lexical, al cuvântului sau al expresiei. În continuare elementul de aliniere lexicală îl vom numi, generic, *atom lexical* sau simplu *atom*. Un bitext în care echivalenții de traducere sunt explicitați se numește un bitext *alinat*. Cea mai mică unitate textuală la nivelul căreia se realizează alinierea definește *granularitatea* echivalenților de traducere. Echivalenții lexicali de traducere (obiectul nostru de interes în această lucrare) depind evident de bitextul din care sunt extrași iar procesul de extragere a lor devine echivalent cu extragerea unui dicționar bilingv, specific unui anumit domeniu¹. Extragerea unui dicționar de echivalenți de traducere dintr-un bitext este în fond un proces de explicitare a dicționarului mental folosit de translatorul (sau translatorii) textului original.

Presupoziția fundamentală în încercarea de a alinia corpusurile paralele este că *aceeași* semnificație este exprimată în două sau mai multe limbi. Definirea identității de înțeles între două sau mai multe reprezentări ale (presupus) aceluiași lucru este o binecunoscută problemă filozofică care rămâne deschisă chiar în domenii mult mai precise decât cel al limbii (de pildă în ingineria software). Prin urmare, noțiunea de echivalent de traducere este un concept vag, și pentru operaționalizarea sa în domenii ca traducerea automată, terminologie, managementul multilingual al documentelor și altele asemenea avem nevoie de o definiție precisă în termeni direct cuantificabili. Una dintre cele mai larg acceptate definiții a echivalenței de traducere este cea folosită în [1]: "the translation equivalence defines a (symmetric) relation that holds between two different language texts such that expressions appearing in corresponding parts of the two texts are reciprocal translations. These expressions are called *translation equivalents*".

Majoritatea abordărilor moderne în extragerea automată a echivalenților de traducere², sprijinite de forța de calcul din ce în ce mai mare a calculatoarelor, utilizează metode statistice și pot fi clasificate în două mari categorii:

¹Posibilitatea de a genera automat dicționare bilingve în domenii specializate, coroborată cu performanțele tot mai bune ale programelor de clasificare automată a textelor, deschide noi perspective traducerii automate și în general prelucrării multilinguale a textelor. În continuare, dacă nu vom specifica altminteri, prin "echivalenți de traducere" vom înțelege implicit "echivalenți lexicali de traducere".

- paradigma "*presupune și testează*" [2], [3] etc. se bazează pe generarea unei mulțimi de potențiali echivalenți de traducere (spațiul ipotezelor) din care se selectează ulterior, pe baza unor teste de independență statistică, echivalenții de traducere. Selectarea fiecărui echivalent de traducere se face independent de echivalenții extrași anterior (procesul poate fi considerat ca fiind unul de optimizare locală).
- ® paradigma "*modelului de limbă*" [4], [5], [6] etc. presupune construirea unui model statistic al bitextului, model ai cărui parametri se estimează prin metode de optimizare globală. În această abordare un candidat supus estimării nu mai este o pereche de atomi lexicali ci o mulțime de perechi, numită *asignare* [4].

Există susținători și critici ai ambelor abordări și o discuție a avantajelor și dezavantajelor lor este prezentată în [6]. În esență, paradigma "*presupune și testează*" este mult mai eficientă din punct de vedere computațional deoarece presupune investigarea unui spațiu al soluțiilor proporțional cu N^2 , unde N este maximul dintre numerele de articole lexicale distincte din cele două părți ale bitextului, dar echivalenții de traducere cu număr mic de apariții sunt de obicei pierduți. Paradigma "*modelului de limbă*" este extrem de costisitoare din punct de vedere computațional întrucât spațiul de căutare al soluțiilor este teoretic proporțional cu $N!$, în schimb având potențialitatea identificării corecte chiar și a echivalenților de traducere cu o singură apariție în bitext (*hapax-legomena*). În [4] sunt prezentați o serie de algoritmi foarte eficienți, bazați pe o serie de ipoteze simplificatoare dar raționale, ce permit ignorarea unor mari regiuni din spațiul de căutare, regiuni în care este improbabil să existe soluții acceptabile.

Metoda descrisă aici poate fi încadrată în categoria abordărilor de tip "*presupune și testează*". Algoritmul generează mai întâi o listă de candidați și apoi succesiv, alege din această listă perechile cu cele mai mari scoruri de co-ocurență în regiuni corespondente ale bitextului. După cum se va vedea în continuare, acest algoritm nu are nevoie de un dicționar bilingv inițial, dar dacă acesta există, utilizarea sa poate spori substanțial viteza și acuratețea prelucrării.

2.2. Ipoteza corespondenței lexicale 1:1

În general, un cuvânt dintr-un segment ce apare într-o parte a bitextului se traduce în segmentul corespunzător din cea de a doua parte a bitextului tot printr-un singur cuvânt. Dacă acest lucru s-ar întâmpla întotdeauna, problema alinierii lexicale a unui bitext ar fi substanțial mai simplă decât în realitate. Din păcate ipoteza "*cuvânt la cuvânt*" nu este adevărată în foarte multe cazuri, astfel încât adoptarea ei ca premisă de calcul nu pare foarte promițătoare. Dificultatea poate fi însă ocolită prin considerarea ipotezei conform căreia un articol lexical dintr-o limbă se traduce în cealaltă tot printr-un singur articol lexical. Un articol lexical

este reprezentat fie de un cuvânt, fie de o secvență de cuvinte (sintagmă, compus, l expresie). Această formulare, cunoscută sub numele de "*ipoteza corespondenței*" [^] "*lexicale 1:1*", adoptată ca premisă computațională, simplifică mult problema fintă a alinierii lexicale a unui bitext, dar introduce probleme noi și anume definirea și respectiv recunoașterea automată a articolelor lexicale. Din fericire aceste probleme sunt reducibile la contexte monolingve și au soluții simple și foarte eficiente. Un program capabil să realizeze recunoașterea automată a articolelor lexicale se numește *segmentator lexical*. Un segmentator lexical este în general independent de limbă, iar funcționarea sa este controlată prin resurse specifice (dicționare conținând cuvinte, secvențe de cuvinte sau expresii regulate definite peste un vocabular limitat). În [7] este discutată structura resurselor necesare segmentării lexicale a textelor în limba română cu ajutorul segmentatorului *MtSeg*, dezvoltat la Universitatea Aix-en-Provence în cadrul proiectului european "Multext".

Adoptarea "*ipotezei corespondenței lexicale 1:1*" reduce dramatic complexitatea problemei extragerii echivalenților lexicali, indiferent de paradigma în care este abordată rezolvarea (a se vedea [7], [8] pentru detalii). Trebuie menționat însă că o segmentare lexicală perfectă (din punctul de vedere al utilității ei într-un context multilingv) este practic imposibilă din cauza incompletitudinii inerente a oricărui dicționar frazai. În [9], [8] se arată cum poate fi surmontată această incompletitudine a resurselor necesare segmentării lexicale.

2.3. Etape de preprocesare

2.3.1 Alinierea frazată

Înainte de extragerea propriu-zisă a echivalenților de traducere, un corpus paralel este supus unor prelucrări preliminare, de aducere a bitextului într-o formă normalizată. După ce fiecare parte a bitextului a fost supusă segmentării lexicale, urmează etapa de aliniere la nivelul propoziției a corpusului paralel. Pentru acest scop, am utilizat o variantă puțin modificată a algoritmului prezentat și documentat în [10]. În [7] este descris procesul de aliniere la nivel de frază și furnizate exemple și statistici pentru diferite perechi de limbi prezente în corpusul paralel multilingv "1984", conținând traduceri, în șase limbi, ale romanului omonim al lui George Orwell. Acolo arătam că, în marea majoritate a cazurilor, traducerile din limba engleză s-au realizat în celelalte limbi păstrând corespondența de 1:1 la nivelul frazei, cu alte cuvinte, aproape întotdeauna o frază din textul englezesc a fost tradusă ca o singură frază în celelalte limbi reprezentate în corpusul paralel. Algoritmul de aliniere la nivelul frazei poate depista și acele cazuri în care traducerea s-a realizat fără păstrarea corespondenței 1:1. Astfel, au fost cazuri în care două fraze sursă au fost traduse printr-o singură frază, sau invers, când o frază din limba engleză a fost tradusă prin 2

[^] *Noțiunea de frază este luată aici în sensul ei larg, al unei propoziții sau fraze (enunț terminat cu un semn de punctuație din categoria celor finale: punct, punct și virgulă, două puncte, semnul exclamării, semnul întrebării, trei puncte).*

sau chiar 3 fraze în celelalte limbi. În cele ce urmează, indiferent de tipul de aliniere (1:1, 2:1, 1:2 etc.) vom numi porțiunile aliniate la nivelul frazai, *unități de traducere* (UT).

Rațiunea acestei etape de prelucrare constă în intuiția comună că elementele lexicale aflate în relație de echivalență de traducere se regăsesc în frazele ce se constituie în unități de traducere. Pe de altă parte, procesul de aliniere la nivelul frazei este mult mai simplu, pentru că în general indiferent de perechile de limbi considerate într-un bitext ordinea frazelor dintr-o limbă este păstrată în cealaltă limbă. Această ipoteză, operaționalizată de un algoritm de optimizare dinamică de genul celui descris în [10], permite printre altele și identificarea porțiunilor netraduse într-una din limbi (alinieri de tipul N:0 sau 0:M).

O altă ipoteză simplificatoare pentru procesul identificării echivalențelor lexicale de traducere se bazează pe observația că în marea majoritate a traducerilor, categoriile gramaticale din limba sursă se conservă în limba țintă [1]. Cu alte cuvinte, un verb se traduce de obicei printr-un verb, un substantiv printr-un substantiv ș.a.m.d. Melamed a numit o astfel de pereche de traducere, pereche de tip V, distingând-o de perechile de tip P, în care atomii lexicali în cele două limbi au categorii gramaticale diferite. Melamed, distinge și o a treia categorie de perechi de traducere, tipul I, perechile de traducere incomplete, rezultate ca urmare a unei segmentări lexicale parțiale și a utilizării "ipotezei de aliniere lexicală 1:1". Considerațiile lui Melamed referitoare la distribuția celor trei tipuri de traduceri lexicale sunt foarte bine confirmate de experimentele noastre, în ciuda faptului că textul nostru este un text literar în timp ce textul său este un text politic (dezbatere din Parlamentul Canadian) conținând traduceri literale, mult mai puțin afectate de personalitatea literară a traducătorului. Ceea ce este demn de remarcat este că perechile de tip P nu conțin categorii gramaticale arbitrare, și că de la o pereche de limbi la alta, se pot identifica regularități în alternanța categoriilor gramaticale la traducere (de ex. participiu-adjectiv, gerunziu-substantiv, gerunziu-adjectiv). Astfel de regularități pot fi abstractizate prin expresii regulate, efectul net fiind ca multe din perechile de tip P pot fi asimilate (algoritmice) perechilor de tip V. Prin urmare, necesitatea identificării rapide și precise a categoriei gramaticale (și eventual al altor trăsături morfologice sau lexicale) pentru atomii lexicali dintr-un bitext impune o altă prelucrare preliminară, respectiv etichetarea morfo-lexicală, prelucrare pe care o prezentăm în secțiunea următoare.

2.3.2 Etichetarea morfo-lexicală și lematizarea

Etichetarea morfo-lexicală este procesul prin care fiecărui articol lexical dintr-un text arbitrar i se atribuie un cod morfo-lexical unic dintr-o mulțime specifică articolului lexical respectiv, numită clasa sa de ambiguitate. Codul morfo-lexical reprezintă o reprezentare compactă, și de obicei standardizată, a proprietăților

morfologice și lexicale ce caracterizează apariția unui atom lexical într-un text. Clasa de ambiguitate a unui atom lexical reprezintă mulțimea tuturor interpretărilor posibile în orice context legal al atomului respectiv. De exemplu cuvântul "urâtr" are cel puțin 8 interpretări posibile, putând fi substantiv, adjectiv sau verb. Lema sa poate fi una dintre "urât" (substantiv sau adjectiv), "a urâți" sau "a urî" (verb).

urâți	urâți	Vmnp	(inf.: A <i>urâți</i> înseamnă a face să devină urât)
urâți	urâți	Vmis3s	(ind., perf simplu, sing., pers. 3: El <i>urâți</i> totul în viața ei)
urâți	urâți	Vmm-2s	(imp., sing: Prietene, nu <i>urâți</i> singurul lucru frumos din viața lui!)
urâți	urî	Vmip2p	(ind., prez., pl, pers. 2: De pomană îi <i>urâți</i> pe ei, ceilalți sunt de vină)
urâți	urî	Vmsp2p	(subj., prez., pl., pers. 2: Voi ar trebui să <i>urâți</i> tot ce e împotriva vieții)
urâți	urî	Vmm-2p	(imp., pl.: Nu-i <i>urâți</i> pe apărătorii planetei!)
urâți	urât	Afpmp-n	(adj., mase. pl., neart. : Doi câini <i>urâți</i> și răi păzeau intrarea.)
urâți	urât	Ncmp-n	(subs. com., mase. pl., neart.: Niște <i>urâți</i> m-au băgat în sperieți.)

Așadar, clasa de ambiguitate a cuvântului "i/râff este mulțimea (Vmnp, Vmis3s, Vmm-2s, Vmm-2p, Vmip2p, Vmsp2p, Afmpmp-n, Ncmp-n), iar etichetarea morfo-lexicală a acestui cuvânt înseamnă a alege, în funcție de contextul apariției sale, unul și numai unul dintre cele 8 coduri reprezentând interpretarea contextuală a cuvântului. În cercetările anterioare am dezvoltat o metodă statistică de etichetare morfo-lexicală [11], numită etichetarea cu două niveluri și modele de limbă combinate (TT-CLAM: tiered-tagging with combined language models), bazată pe programul TnT al lui Thorsten Brants [12] de prelucrare a modelelor markov cu legături ascunse de ordin 2 (3-gram HMM), program ce poate fi descărcat de la adresa www.coli.uni-sb.de/~thorsten/tnt/. Abordarea TT-CLAM a arătat că texte arbitrare în limba română pot fi etichetate morfo-lexical în mod corect în peste 98.5% din cazuri și că atunci când de interes este numai categoria gramaticală, procentul de etichetare corectă depășește 99.5%. Metoda TT-CLAM s-a dovedit independentă de limbă, rezultate mai bune decât în alte abordări fiind raportate în literatura de specialitate pentru limbi foarte diferite de limba română: limba maghiară [13], [15] limba germană [16], [17].

Lematizarea este procesul prin care o formă flexionară a unui articol lexical (cuvânt sau expresie) este redusă la forma normală de dicționar. Lematizarea se poate realiza fie printr-un proces de analiză morfologică fie prin căutarea într-o bază de date lexicale, conținând cuvinte în formă flexionară însoțite de analiza lor morfologică și de forma lernă. Lematizarea se realizează în acest caz prin

identificarea în baza de date a lemei pentru care forma flexionară și analiza morfo-lexicală sunt identice cu cele din textul de lematizat, care desigur a fost în prealabil etichetat. Pentru limba română, noi am experimentat cu ambele metode și datorită vitezei mult superioare, am optat pentru varianta a doua.

În figura de mai jos este exemplificat rezultatul prelucrărilor preliminare discutate în această secțiune (segmentare lexicală, aliniere frazală, etichetare morfo-lexicală și lematizare) pentru începutul bitextului Englez-Român din corpusul multilingv "1984". Prima linie arată că în limba română, fraza cu identificatorul Oro.1.2.2.1, reprezintă traducerea a două fraze din textul englezesc, respectiv a celor cu identificatorii Oen.1.1.1.1 și Oen.1.1.1.2 (avem deci o aliniere de tip 1:2). Liniile următoare, specifică pentru fiecare articol lexical din fiecare limbă tipul său (TOK, LSPLIT, DATE, ABR etc), forma ocurență, lerna, codul morfo-lexical și categoria gramaticală (ultimele trei separate prin caracterul T).

<linktargets="Oro.1.2.2.1; Oen.1.1.1.1 Oen.1.1.1.2">

<S FROM="Oro.1.2.2.1">			<S FROM="Oen.1.1.1.1">		
LSPLIT	într-	întru\Spsay\S	TOK	It	it\Pp3ns\P
TOK	o	un\Tifsr\T	TOK	was	be\Vm3s\AUX
TOK	zi	zi\Ncfsrn\N	TOK	a	a\Di\D
TOK	senină	senin\Afpfsrn\A	TOK	bright	bright\AAA
			</S>		
			<S FROM="Oen.1.1.1.2">		
</S>			</S>		

Figura 1: Bitext preprocesat pentru extracția echivalențelor lexicale de traducere

O descriere a principiilor de codificare morfo-lexicală, în conformitate cu recomandările EAGLES poate fi găsită în [18]. Codificarea specifică pentru limba română, conformă cu standardul respectiv este pe larg descrisă în [19].

2.4. Un prim algoritm de extragere automată a echivalențelor lexicale de traducere

Există, așa cum am văzut mai sus, mai multe ipoteze simplificatoare care permit ținerea sub control a complexității problemei extragerii automate a echivalențelor de traducere. Nici una dintre aceste ipoteze nu este satisfăcută întotdeauna, dar situațiile în care ele nu sunt adevărate sunt suficient de rare astfel încât adoptarea lor nu alterează semnificativ valoarea rezultatelor. Trebuie

subliniat faptul că ipotezele simplificatoare folosite de noi, discutate în continuare, în general nu afectează precizia dicționarilor bilingve extrase ci completitudinea lor. Altfel spus, o corectitudine (echivalenți de traducere reali), deși prezente în bitext, pot să nu fie găsite. Precizia și completitudinea (în limba engleză acești termeni sunt definiți în mod standard astfel:

PREC=(număr de echivalenți corect extrași)/(număr total de extrași)

COMP=(număr de echivalenți corect extrași)/(număr total de existenți în bitext)

Mai trebuie precizat și faptul că ipotezele simplificatoare nu împiedică recuperarea ulterioară a echivalențelor negăsiți din aceste ipoteze de lucru. În [9] sunt discutate metode de recuperare a echivalențelor de traducere ce nu respectă ipoteza "echivalenței lexicale 1:1".

- ipoteza "echivalenței lexicale 1:1"; ea stă la baza majorității metodelor cunoscute: [20], [21], [6], [22], [23], [1] etc. Așa cum am văzut de vremea, un articol lexical identificat corespunzător de către un articol lexical adecvat diminuează considerabil efectul corectitudinii ipoteze;
- un articol lexical polisemantic ce apare de mai multe ori în unitate de traducere este folosit cu același termen în toate ipotezele presupuzite este explicit utilizată de [1] și implicit de celelalte amintiți mai sus;
- un articol lexical dintr-o parte a unității de traducere este folosit în altă parte a UT doar dacă are aceeași categorie gramaticale compatibile; în majoritatea cazurilor compatibilitatea categoriilor gramaticale se reduce la aceeași categorie cum am specificat anterior, este posibil să existe și alte corespondențe compatibile (de pildă, verbele la participiu din limba engleză sunt destul de frecvent traduse în română cu adjective sau substantive, și reciproc).
- Deși ordinea cuvintelor nu este un invariant al traducerii, o ordine nici arbitrară; când două sau mai multe perechi de cuvinte pot candida la statutul de echivalenți de traducere, o strategie de evaluare nu permit departajarea lor, atunci este preferabil să conțină articolele cele mai apropiate în pozițiile lor în text. Ouristică este, de asemenea, folosită de [23].

Pe baza bitextului preprocesat așa cum s-a prezentat în secțiunea precedentă, primul pas al algoritmului este de a delimita spațiul soluțiilor. Acest lucru se realizează prin construcția unei liste a tuturor

posibili (în conformitate cu ipotezele de lucru amintite mai sus). Această listă, pe care o notăm cu TECL (Translation Equivalence Candidates List) conține la rândul ei o mulțime de sub-liste (câte una pentru fiecare categorie gramaticală luată în considerare). Fiecare sublistă conține perechi de forma <token_s token_r> unde *token_s* și *token_r* sunt articole lexicale de categorii gramaticale compatibile și care au apărut în părțile corespunzătoare ale aceleiași unități de traducere. Fie TU^j cea de a j^a unitate de traducere (translation unit). Prin colectarea tuturor articolelor lexicale aparținând aceleiași categorii gramaticale POS_k (păstrând ordinea lor relativă și eliminând duplicatele) se construiesc pentru fiecare TU^j mulțimile ordonate L^{pos_k} și L^{tr_k}. Pentru fiecare POS_k fie TUV_{osi} produsul cartezian L^{pos_k} × L^{tr_k}. Atunci, definim lista de corespondențe în unitatea de traducere TU^j ca fiind CTU^j (correspondences in the jth translation unit):

$$CTO = \sum_{i=1}^{no.of.pos} |J^{TU} Pos_i|$$

Cu aceste notații, și presupunând că bitextul de intrare conține *n* unități de aliniere, atunci TECL se definește astfel:

$$TECL = [JCTU^j]_H$$

TECL conține desigur foarte mult "zgomot" și cele mai multe perechi candidate (TEC=Translation Equivalence Candidate) sunt extrem de improbabile. Pentru a elimina cât mai multe din perechile TEC improbabile, TECL este filtrată pe baza unor funcții scor ce supun fiecare TEC la o analiză a ipotezei statistice de independență a asocierii articolelor lexicale. Pentru a prezenta funcțiile scor pe care le-am utilizat în experimentele noastre, vom mai defini o serie de notații:

- TEC = <T_s T_r> ∈ TECL, un potențial echivalent de traducere definit ca perechea formată din articolul lexical sursă T_s și posibila sa traducere T_r în limba țintă;
- nu = numărul de ocurențe ale <T_s T_r> din TECL;
- n_{1,2} = numărul de perechi <T_s -i T_r> din TECL în care T_s a fost asociat cu un articol lexical diferit de T_r;
- n_{2,1} = numărul de perechi <-i T_s T_r> din TECL în care T_r a fost asociat cu un articol lexical diferit de T_s;
- n₂ = numărul de perechi <-J_s -> T_r> din TECL ce nu conțin nici pe T_s și nici pe T_r;
- n_{*} = numărul de perechi <T_s *> din TECL în care apare T_s indiferent cu cine este asociat;
- n<i = numărul de perechi <* T_r> din TECL în care apare T_r indiferent cu cine este asociat;

- n₂* = numărul de perechi <-i T_s *> din TECL în care T_s nu apare;
- n<sub>2 = numărul de perechi <* -i T_r> din TECL în care T_r nu apare;
- n« = numărul de perechi <* *> din TECL;

Tabela de contingență din figura de mai jos ilustrează aceste notații:

	T _r	-i T _r	
T _s	nu	n _{1,2}	n _r
-i T _s	n _{2,1}	n _{2,2}	n<sub>2
	n _m	n<sub>2*	n~

$$n_r = nu + n_{1,2}, \quad n_{2}^* = n_{2,1} + n_{2,2}$$

$$n_{2,1}^* = n_{2,1} + n_{2,2}, \quad n_{2,2}^* = n_{2,1} + n_{2,2}$$

$$n_{2,2}^* = \sum_{i=1}^2 \sum_{j=1}^2 n_{ij}$$

Figura 2: Tabela de contingență pentru un potențial echivalent de traducere <T_s T_r>

Pentru ordonarea potențialilor echivalenți de traducere în vederea filtrării (eliminarea candidaților cei mai puțin plauzibili) am realizat experimente folosind 4 funcții de calcul al scorului de echivalență: MI (*informația mutuală*), DICE, LL (log likelihood), and %² (chi-pătrat). Folosind notațiile de mai sus, aceste funcții-scor se definesc în felul următor:

$$(1) \quad MI(TT_r, T_s) = \log \frac{n_{11} * n_{22}}{n_{12} * n_{21}}$$

$$(2) \quad DICE(T_r, T_s) = \frac{2n_{11}}{n_{11} + n_{22}}$$

$$(3) \quad LL(TT_r, T_s) = 2 * \sum_{i=1}^2 \sum_{j=1}^2 n_{ij} * \log \frac{n_{ij}}{n_{i.} * n_{.j}}$$

$$(4) \quad \chi^2(TT_r, T_s) = \sum_{i=1}^2 \sum_{j=1}^2 \frac{(n_{ij} - n_{i.} * n_{.j} / n_{22})^2}{n_{i.} * n_{.j} / n_{22}}$$

Figura 3: Funcții-scor pentru evaluarea unui potențial echivalent de traducere <T_s T_r>

O formulă mai simplă de calcul pentru $\%^2(T_s, T_t)$ este următoarea:

$$(4') \quad \%^2(T_s, T_t) = \frac{n_{ii}^2}{(n_{i.} \cdot n_{.i})}$$

Filtrarea potențialilor echivalenți de traducere se face în raport cu un prag numeric impus scorului calculat cu una dintre funcțiile de mai sus. Toate perechile ce obțin un scor mai mare decât pragul ales sunt considerate plauzibile și vor fi supuse unor prelucrări suplimentare iar celelalte sunt eliminate. Orice metodă de filtrare statistică va elimina mulți echivalenți falși de traducere, dar pe lângă aceștia și un număr de perechi corecte. Alegerea pragului de scor s-a făcut având ca obiectiv minimizarea numărului de perechi corecte dar eliminate în mod greșit și a numărului de perechi incorecte acceptate ca urmare a scorului superior pragului de selecție. După mai multe experimente, cele mai bune rezultate s-au obținut folosind funcția de scor LL cu limita pragului de acceptanță egală cu 9.

Într-o primă variantă, algoritmul nostru de extragere a echivalenților de traducere, având unele asemănări cu algoritmul iterativ prezentat în [23], implementa o strategie de selecție indiferentă la locul și poziția în corpus a articolelor lexicale apărând în perechea TEC analizată la un anumit moment. O diferență majoră față de algoritmul descris în [23] este că în programul nostru calculul diferitelor probabilități (mai exact al estimațiilor de probabilitate) și al scorurilor (testul t) devine necesar, conducând la o viteză de prelucrare cu cel puțin un ordin de mărime mai mare. Pornind de la lista filtrată a potențialilor echivalenți de traducere, algoritmul selectează în mod iterativ cei mai plauzibili candidați (vezi mai jos) și apoi îi șterge din lista inițială. Algoritmul se oprește după un număr prestabilit de iterații sau mai devreme în cazul în care lista candidaților s-a golit sau dacă nici un candidat nu mai îndeplinește condiția de selecție.

În iterația k a algoritmului se construiește o matrice de contingență (TBLk) pentru fiecare categorie gramaticală (POS) având dimensiunile $S_m \times T_n$ unde S_m și T_n reprezintă numărul de articole lexicale din limba sursă respectiv țintă care mai există în lista de candidați la pasul k (Figura 4). Liniile și coloanele tabelului sunt indexate cu articolele lexicale (având aceeași categorie gramaticală) din limba sursă respectiv limba țintă. Fiecare celulă (i,j) a matricii reprezintă numărul de ocurențe în lista de candidați a perechii $\langle T_{s,i}, T_{t,j} \rangle$.

	T _t , n		
T _s i	nu	n _{in}	n _i *
		n _{imn}	n _m *
	n _i *		n _m **

$$ny = \text{occ}(T_{s,i}, T_{t,j}); i \geq 1; n^* = \sum_{i=1}^n n_{s,i}; n_- = \sum_{j=1}^m (x^*_{i,j}) -$$

Figura 4: Matricea de contingență la pasul k

Condiția de selecție la pasul k a mulțimii de echivalenți de traducere este exprimată de relația (5):

$$(5) \quad TP^k = \{ \langle T_{s,i}, T_{t,j} \rangle \mid V_{p,q} (n_{s,i} > n_{s,q}) \wedge (n_{t,j} > n_{t,m}) \}$$

Condiția de mai sus constituie esența algoritmului iterativ (numit în [14] algoritmul BASE) și ea spune că pentru a selecta perechea $\langle T_{s,i}, T_{t,j} \rangle$ drept echivalent de traducere, numărul de asocieri ale lui $T_{s,i}$ cu $T_{t,j}$ trebuie să fie mai mare sau cel puțin egal decât numărul de asocieri ale lui $T_{s,i}$ cu orice alt $T_{t,p}$ ($p \neq j$) și simultan numărul de asocieri ale lui $T_{t,j}$ cu $T_{s,i}$ trebuie să fie mai mare sau cel puțin egal decât numărul de asocieri ale lui $T_{t,j}$ cu orice alt $T_{s,q}$ ($q \neq i$). Toate perechile selectate în TP^k sunt eliminate din lista de candidați (ceea ce în matricea de contingență pentru pasul $k+1$ implică punerea pe 0 a contoarelor de ocurență pentru perechile selectate anterior). Dacă $T_{t,j}$ este tradus în mai multe moduri (fie pentru că are sensuri ce se lexicalizează diferit în limba țintă, fie pentru că în limba țintă se folosesc diferiți sinonimi pentru $T_{t,j}$) restul traducerilor sale va fi extras în iterațiile următoare. Algoritmul discutat este schițat în figura 5:

```

procedure BASE(bitext, step; dictionary) is:
    k=1;
    TP(0)={};
    TECL(k)=build-cand(bitext);
    for each POS in TECL do
        loop
            TECL(k)=update(TP(k-1), TECL(k))
    
```

```

TBL(k)=build_TEC_table(TECL(k));
TP(k)=select(TBL(k)); ## relația (5) ##
add(dictionary, TP(k));
k=k+1;
until {(TECL(k-1) is empty) or (TP(k-1) is empty) or (k >
step)}
endfor
return dictionary
end

```

Figura 5: Algoritmul iterativ de extragere a echivalenților de traducere

2.5. Un algoritm îmbunătățit de extragere automată a echivalenților lexicali de traducere

Una dintre principalele deficiențe ale algoritmului BASE este vulnerabilitatea la ceea ce [1] numește *asociații indirecte*. Dacă $\langle T_s, T_r \rangle$ are un scor de coocurență ridicat iar T_r apare (dintr-un motiv sau altul) de mai multe ori împreună cu T_{r_1} , s-ar putea ca și perechea $\langle T_s, T_{r_1} \rangle$ să primească un scor ridicat. Deși, așa cum observa și Melamed, în general, asociațiile indirecte au un scor mai mic decât cele directe (corecte), ele pot obține totuși scoruri mai mari decât multe alte perechi corecte ce n-au legătură cu T_s iar acest lucru nu numai că generează echivalenți de traducere greșiți, dar va elimina din competiție și echivalenți corecți. Prin urmare asociațiile indirecte afectează atât precizia cât și completitudinea procesului. Pentru a slăbi această sensibilitate în implementarea algoritmului BASE a fost nevoie de stabilirea unei limite inferioare de ocurență pentru fiecare articol lexical luat în considerație. Această limită, conduce inevitabil la eliminarea din spațiul de căutare a soluțiilor a mai mult de 50% dintre echivalenții de traducere. Deficiența algoritmului BASE se explică prin faptul că scorurile de coocurență sunt calculate în mod global fără a verifica dacă atomii lexicali ai unei perechi evaluate sunt sau nu prezenți în unitățile de traducere prelucrate.

Pentru diminuarea influenței asociațiilor indirecte fără a mai impune un prag de ocurență, algoritmul BASE a fost modificat astfel încât ierarhizarea și alegerea celor mai probabili echivalenți de traducere se realizează în contextul local al fiecărei unități de traducere (deși scorurile lor se calculează tot la nivelul întregului bitext). Cu această modificare, noul algoritm (BETA) se apropie de algoritmul "competitive linking" al lui Melamed [1]. Candidații proveniți din unitatea de traducere curentă sunt analizați prin prisma scorului lor de coocurență și cel cu scorul cel mai mare este selectat. În baza ipotezei corespondenței lexicale 1:1,

Pierderea unui așa mare număr de echivalenți de traducere desigur nu surprinde întrucât

dintre candidații rămași sunt eliminați toți aceia care conțin un număr de echivalenți de traducere mai mic decât pragul de ocurență. Dintre candidații care rămân după filtrare, se alege din nou cel cu scorul cel mai bun și iar se evaluează până când nici un echivalent de traducere nu mai poate fi extras din unitatea de traducere curentă, caz în care algoritmul trece la prelucrarea următoarei unități de traducere.

Eliminarea pragului de ocurență a îmbunătățit substanțial precizia și calitatea dicționarelor de traducere (o detaliată comparație a performanțelor este prezentată în anexa 2). Această modificare a ridicat problema decelării între candidații cu una sau două echivalenți de traducere, dar a ridicat problema decelării între candidații cu una sau două echivalenți de traducere care scorul de coocurență este statistic nesemnificativ. În acest caz, pragul de frecvență a fost înlocuit cu o combinație între un scor de similaritate și un scor de proximitate relativă. Funcția de similaritate ortografică $COGN(T_s, T_r)$, este o variantă a funcției $xxdixc$ descrisă în [21]. Dacă T_s este un șir de m caractere $o_1 a_1 . . . a_m$ and T_r un șir de n caractere $o_2 a_2 . . . a_n$, construiesc două noi șiruri T'_s și T'_r prin inserarea în T_s și T_r a unor caractere speciale astfel încât în final șirurile T'_s și T'_r au aceeași lungime $(\max(m, n) < p < m+n)$ și un număr maxim de caractere poziționale identice. Dacă P_i un caracter din T'_s și P_j un caracter din T'_r , care se potrivesc în corespondență. Fie $8(ai)$ numărul de caractere speciale care se potrivesc cu caracterul P_i în T'_s și $8(Pj)$ numărul de caractere speciale care se potrivesc cu caracterul P_j în T'_r . Fie q numărul de caractere care se potrivesc în corespondență. În aceste notații, măsura de similaritate $COGN(T_s, T_r)$ se definește astfel:

$$(6) \quad COGN(T_s, T_r) = \frac{\sum_{i=1}^m \sum_{j=1}^n 8(P_i) \cdot 8(P_j)}{m \cdot n} \quad \text{if } m > n$$

Limita de relevanță a scorului de similaritate a fost empiric determinată. Această valoare este dependentă într-o oarecare măsură de performanțele considerată în procesul de extragere a echivalenților de traducere. Valoarea efectivă a testului de similaritate include și o serie de norme aplicabile la teste testate (eliminarea unor afixe, reducerea consoanelor duble, ignorarea caracterelor create de diacritice etc.) normalizări care depind de morfologia limbii țintă.

Cel de al doilea criteriu de evaluare a plauzibilității este scorul de proximitate, $DIST(T_s, T_r)$ definit după cum urmează:

Dacă $(\langle T_s, T_r \rangle \in L_{posk}^{Sj} \otimes L_{posk}^{Tj})$ & $(T_s$ este al n-lea element în L_{posk}^{Sj} și T_r este al m-lea element în L_{posk}^{Tj}) atunci $DIST(T_s, T_r) = |n-m|$

Filtrul $\text{COGN}(T_s, T_r)$ este mult mai semnificativ din punct de vedere lingvistic decât $\text{DIST}(T_s, T_r)$, astfel încât scorul de similaritate are precedență asupra celui de proximitate. Funcția $\text{DIST}(T_s, T_r)$ este invocată doar atunci când $\text{COGN}(T_s, T_r)=0$ (deci când atomii lexicali nu prezintă similaritate ortografică) și perechea $\langle T_s, T_r \rangle$ nu reprezintă o pereche singulară în corpus (hapax-legomena), sau când mai multe perechi candidate au obținut același scor de similaritate.

Algoritmul BETA este schițat mai jos:

```

procedure BETA(bitext;dictionary) is :
  dictionary={};
  TECL(k)=build_cand(bitext);
  for each POS in TECL do
    for each  $TXJ_{\text{pos}}$  in TECL do
      finish=false;
      loop

        best_cand=get_the_highest_scored_pairs( $TU_{\text{pos}}^x$ );
        conflicting_cand=select___conflicts(best_cand);
        non_conflicting_cand = best_cand $\setminus$ conflicting_cand;
        best_cand=conflicting_cand;

        if cardinal(best_cand)=0 then finish=true;
        else
          if cardinal(best_cand)>1 then
            best_cand=filtered(best_cand);
          endif;

          best_pairs = non_conflicting___cand + best_cand
          add(dictionary, best_pairs);
        endif
      until { ( $TU_{\text{pos}}^1$ ={}) lor (finish=true) }
    endfor
  endfor
  return dictionary
end

```

```

procedure filtered(best_cand) is :
  result = get_best_COGN_j3Core(best_cand);
  if (cardinal(result)=0)&(non-hapax(best_cand)) then
    result = get_best_DIST_score(best_cand);
  else if cardinal(result)>1
    result = get_best_DIST_score(best_cand);
  endif
endif

```

```

return result;
end

```

Din corpusul paralel multilingv "1984" am extras 6 bitexte conținând text în limba engleză și traducerea în una din cele 6 limbi amintite. Fiecare bitext a fost prelucrat conform celor prezentate în acest capitol și au fost extrase 6 dicționare bilingve, din care s-a obținut și un dicționar multilingv în 7 limbi (cele 6 plus engleza). În [8] este furnizată o analiză contrastivă cu alte sisteme de acest tip și vitezei de prelucrare. Timpul mediu de extragere a unui dicționar bilingv din corpusul paralel multilingv "1984" (circa 110.000 de cuvinte în fiecare limbă) este de câteva minute. Eșantioane ale acestor dicționare pot fi consultate la adresa <http://www.racai.ro/Mufis/BilingualLexicons/AutomaticallyExtractedBilingualLexicons.html>.

3. Dezambiguizarea sensurilor lexicale folosind echivalențele de traducere

3.1. Ambiguitatea limbajului natural

Este binecunoscut faptul că una dintre cele mai dificile probleme în prelucrarea automată a limbajului natural este ambiguitatea sa inerentă. Ambiguitatea se manifestă la toate nivelurile tradiționale ale analizei de limbă: nivelul fonetic și/sau lexical, sintactic, semantic sau discursiv. Ambiguitatea prezentă pe fiecare nivel generează exploziv ambiguități pe nivelurile următoare. De pildă, omofonia sau omografia prezentă pe primul palier, la nivelul unuia sau mai multor cuvinte va produce secvențe lexicale diferite (combinația tuturor interpretărilor posibile la acest palier) pentru intrarea fazei de analiză sintactică. Fiecare secvență poate conduce, din pricina unor ambiguități de natură structurală, la interpretări sintactice multiple, după cum o serie de secvențe lexicale vor putea fi abandonate pe motivul contrazicerii unor restricții postulate de modelul sintactic al limbii prelucrate. Fiecare dintre interpretările sintactice posibile poate la rândul său să conducă la multiple interpretări semantice, în virtutea multiplelor sensuri pe care le poate avea fiecare element frazai al unei analize sintactice. Designul interpretarea semantică poate elimina unele structuri sintactice generate în faza anterioară pe baza încălcării unor restricții semantice (valabile în orice univers de discurs sau specifice unor domenii discursive de interes). În sfârșit, în analiza discursivă, în care contextul interpretativ transcende limita propoziției, ambiguitățile rămase se presupun a putea fi rezolvate prin utilizarea restricțiilor pragmatice motivate fie de principii generale ale dialogului (coeziune, coerență), fie de natura bine precizată a unui univers de discurs (modelată prin cunoștințe extra-lingvistice despre entitățile universului de discurs). De pildă, în [24] rezolvarea anaforelor este proces tipic analizei de discurs, este modelată în termenii identificării căilor de

^ £ S3S55S8S5 £ 3S £

accesibilitate a entităților menționate în discurs ("vene ale discursului"), care la rândul lor sunt formal definite pe baza principiilor generale ale coeziunii și coerenței unui text.

Rezultă din cele spuse până aici că identificarea și rezolvarea timpurie, la fiecare nivel de prelucrare, a ambiguităților este un imperativ al oricărui demers computațional privind prelucrarea limbajului natural. Și cum cuvântul (sau mai exact spus, atomul lexical) este elementul primar în prelucrarea limbajului o mare parte a eforturilor de cercetare este îndreptată spre nivelul lexical al prelucrărilor. Metodele de etichetare morfo-lexicală (tagging), printre care etichetarea cu două niveluri și modele de limbă combinate - amintită în capitolul 2, permit rezolvarea cu mare acuratețe a ambiguităților categoricale și intracategoriale. De pildă cuvântul *vin* poate fi atât substantiv cât și verb (ambiguitate categorială), iar ca verb, el conține ambiguitatea intracategorială de persoană, număr și mod ("indicativ + persoana I + număr singular", "conjunctiv + persoana I + număr singular" sau "indicativ + persoana III + număr plural"). Un program de etichetare morfo-lexicală "instruit" corect pentru limba română este capabil să rezolve, în contextul apariției sale, astfel de ambiguități morfo-lexicale.

Curentul lexicalist, predominant în modelarea sintactică a limbajului natural, presupune precizarea în descrierea de dicționar a fiecărui cuvânt a proprietăților și restricțiilor sale distribuționale sau cologaționale relevante pentru analiza sintactică. Pe baza acestor descrieri lexicalizate și a contextului local, multe din potențialele ambiguități structurale pot fi eliminate, înaintea unei costisitoare analize sintactice, prin tehnici cunoscute sub numele de analiză sintactică parțială (*parțial parsing* sau *shallow parsing*).

Un cuvânt omograf, chiar după ce a fost corect clasificat din punctul de vedere al categoriei sale gramaticale și al proprietăților sale distribuționale sau cologaționale, poate rămâne ambiguu din punct de vedere semantic. Identificarea sensului cu care este utilizat cuvântul polisemantic într-un context dat este desigur de mare interes. Există însă diferite grade de rafinare a noțiunii de sens, iar natura aplicației pentru care identificarea sensului este necesară poate impune o accepție a noțiunii de sens diferită de cea utilizată într-un dicționar explicativ. Să luăm, de pildă, problema traducerii automate. Întrucât în imensa majoritate a cazurilor rezultatul traducerii este destinat uzului uman, ceea ce este important este ca în textul tradus să nu apară ambiguități suplimentare față de cele din textul sursă. Cu alte cuvinte, dacă o analiză algoritmică evidențiază în limba sursă o serie de ambiguități, pornind de la premiza că textul este admisibil pentru vorbitorii nativi ai limbii textului sursă, de cele mai multe ori este nenaturală o traducere ce încearcă să evite total ambiguitatea identificată. La nivel lexical, aceasta revine la a spune că dacă diferitele sensuri ale unui cuvânt din limba sursă nu se lexicalizează prin cuvinte diferite în limba țintă, este neproductivă o încercare a diferențierii sensului contextual, atâta timp cât indiferent care ar fi el, traducerea cuvântului respectiv în

limba țintă este aceeași. De exemplu, cuvântul englezesc "bottle" are în Wordnetl.5 [25] două sensuri (ca substantiv) anume de vas de sticlă sau plastic cilindric cu un gât îngust și fără mâner, respectiv cantitatea de substanță conținută într-un astfel de vas. Ambele sensuri se regăsesc în cuvântul românesc "sticlă" (care însă include și alte sensuri lexicalizate în engleză prin cuvântul "glass"). În acest caz, a încerca eliminarea ambiguității la traducerea textului "He drank only a bottle of beer" în limba română, de pildă prin utilizarea unei parafraze de genul "El băuse doar conținutul unei sticle de bere", este nenesesară. Orice vorbitor al limbii române va găsi traducerea "El băuse doar o sticlă de bere" mult mai naturală și desigur nu va avea dificultăți în a înțelege despre ce este vorba.

Același gen de considerații se poate face și în raport cu ambiguitățile sintactice pure. Celebrul exemplu "I saw the Statue of Liberty flying over New York" conține cel puțin 4 ambiguități, dar dacă de pildă rezolvarea omografului *saw* (am văzut / tai cu fierăstrăul) este esențială în traducere, rezolvarea ambiguităților structurale poate fi lăsată în sarcina minții celui ce citește textul: "Am văzut Statuia Libertății zburând deasupra New York-ului", căci dacă cititorul englez nu are dificultăți în a înțelege cine și cum zbură, e plauzibil că nici cititorul român (de exemplu) nu le va avea. Aceasta nu înseamnă că nu există ambiguități structurale a căror nerezolvare prealabilă să nu conducă la traduceri hazlii sau chiar incompreensibile. Ideea este că metodele formale de analiză a limbajului modelabile algoritmic, explicitează de multe ori ambiguități greu de conștientizat de omul obișnuit, iar luarea în considerare a factorului uman poate simplifica mult prelucrările automate. Reconsiderarea conceptului de traducere automată în accepțiunea clasică (MT) în favoarea unor concepte mai realiste de tipul HAM (human assisted machine translation) sau MAHT (machine assisted human translation) a relevat faptul că, în numeroase ocazii, posteditarea umană a unui text tradus automat introduce ambiguități care, deși nu sunt sezizabile ușor la lectură, pot fi totuși puse în evidență de algoritmi de analiză.

Cercetările moderne în domeniul dezambiguizării automate, în context, sensurilor cuvintelor sunt motivate și de alte aplicații informatice, cum ar fi clasificarea după conținut a volumelor mari de texte, regăsirea mai precisă a documentelor electronice, rezumarea automată a textelor, extragerea de cunoștințe din texte, crearea de ontologii. Această direcție de cercetare este identificată în literatura engleză prin acronimul WSD (Word Sense Disambiguation) constituie de câțiva ani obiectul unor conferințe specializate și chiar a unor competiții de evaluare (SENSEVAL, ajunsă la a treia ediție) a soluțiilor propuse de specialiști din întreaga lume.

Primii care au sugerat ideea că, pentru obiectivele WSD, sensurile trebuie diferențiate sunt cele care se lexicalizează într-o altă limbă prin cuvinte diferite au fost Resnik and Yarowsky [26]. Intuitiv, se poate presupune că, dacă un cuvânt din limba sursă se traduce în limba țintă în mai multe feluri și acele

traduceri nu sunt sinonimice, atunci trebuie să existe o motivație conceptuală. Analizând un număr suficient de mare de limbi și de texte, e plauzibil, afirmam cei doi specialiști, să identificăm diferențierile lexicale semnificative care delimitează sensurile unui cuvânt. Aceste sensuri sunt numite de cei doi "sensuri tari". Inabilitatea de a identifica corect sensurile tari este principala sursă a erorilor inacceptabile în orice aplicație multilinguală. Utilizarea textelor paralele pentru WSD [27], [28], [29], în scopul identificării proprietăților semantice a lexelelor și a relațiilor dintre ele [30] a folosit implicit sau explicit noțiunea de "sens tare". Mai recent, pe baza echivalențelor de traducere extrași din corpusul "1984" prin procedura noastră, descrisă în capitolul precedent, Ide [31] a arătat că diferențele de traducere în 5 limbi (din 4 familii diferite) pot constitui un criteriu extrem de eficace în identificarea sensurilor tari în limba de pornire (în acest caz, engleza). Resnik and Yarowsky [32] au folosit în schimb traducerea unor propoziții izolate în limba engleză efectuate de vorbitori nativi ai limbilor țintă, dar în mare concluziile studiului lor au fost aceleași cu ale lui Ide. În ambele studii amintite referința pentru limba engleză a fost WordNet [33] și deși rezultatele lor sunt promițătoare, mai ales pentru sensurile tari, ele se bazează pe o mulțime prestabilă de sensuri. Date fiind divergențele semnificative între distincțiile de sens realizate în dicționarele (monolingve) existente, precum și inexistența unui acord general asupra gradului de rafinare a descrierilor de sens în practica lexicografică internațională, raportarea la un inventar prestabil de sensuri, cel puțin din perspectiva prelucrării automate a limbajului, nu pare a fi o soluție optimă. În continuare, vom prezenta o abordare alternativă, detaliată în [34], [35].

3.2. Discriminarea automată a sensurilor lexicale: metodologia

Metoda pe care o vom descrie este menită a identifica sensurile distincte cu care unul sau mai multe cuvinte apar într-un text dat. Întrucât este foarte improbabil ca într-un text omogen, chiar foarte lung (de pildă un roman), un cuvânt să fie folosit în toate sensurile sale, metoda desigur va identifica, prin analiza textuală descrisă în continuare, doar acel sens sau acele sensuri cu care este folosit cuvântul respectiv în textul prelucrat. La limită prin prelucrarea unor texte foarte diferite este posibil teoretic să fie identificate toate sensurile atestate ale unui anumit cuvânt.

Din punct de vedere metodologic, studiul nostru s-a bazat pe corpusul paralel multilingv "1984" și pe dicționarul multilingv extras din acest corpus. Cele 7 limbi ale experimentului nostru fac parte din patru familii: germanică (engleza), romanică (româna), slavică (bulgara, ceha și slovena) și ugro-finică (estoniana, maghiara). Deși corpusul conține un text beletristic, textul orwelian ca și traducerea sa în celelalte limbi nu sunt foarte stilizate și, ca atare, oferă un eșantion rezonabil de limbă modernă, comună. Mai mult, traducerea textului original, efectuate de traducători avizați (unii dintre ei fiind apreciați scriitori), reflectă riguros originalul: pentru mai mult de 95% din textul englezesc o frază sursă este tradusă

* în celelalte limbi tot ca o singură frază. Tipurile de alinieri frazale existente în corpusul "1984" sunt prezentate în tabela de mai jos și discutate în [7]:

Estoniană-Engleză			Maghiară-Engleză			Română-Engleză		
Tip	Nr.	Proc	Tip	Nr.	Proc	Tip	Nr.	Proc
3-1	2	0.030321%	7-0	1	0.014997%	3-1	3	0.046656%
2-2	3	0.045482%	4-1	1	0.014997%	2-4	1	0.015552%
2-1	60	0.909642%	3-1	7	0.104979%	2-3	3	0.046656%
1-3	1	0.015161%	3-0	1	0.014997%	2-2	2	0.031104%
1-2	100	1.516070%	2-1	108	1.619676%	2-1	85	1.321928%
1-1	6426	97.422681%	1-6	1	0.014997%	2-0	1	0.015552%
1-0	1	0.015161%	1-5	1	0.014997%	1-5	1	0.015552%
0-2	1	0.015161%	1-2	46	0.689862%	1-3	14	0.217729%
0-1	2	0.030321%	1-1	6479	97.165573%	1-2	259	4.027994%
			0-4	1	0.014997%	1-1	6047	94.043551%
			0-2	3	0.044991%	0-3	2	0.031104%
			0-1	19	0.284943%	0-2	2	0.031104%
						0-1	10	0.155521%
Bulgară- Engleză			Cehă- Engleză			Slovenă- Engleză		
2-2	2	0.030017%	4-1	1	0.015029%	3-3	1	0.014970%
2-1	23	0.345190%	3-1	2	0.030057%	2-1	48	0.718563%
1-2	72	1.080594%	2-1	109	1.638112%	1-5	1	0.014970%
1-1	6558	98.424134%	1-3	2	0.030057%	1-2	53	0.793413%
0-1	8	0.120066%	1-2	81	1.217313%	1-1	6572	98.383234%
			1-1	6438	96.753832%	1-0	2	0.029940%
			0-1	21	0.315600%	0-1	3	0.044910%

Figura 6: Distribuția tipurilor de aliniere frazală în corpusul paralel "1984"

Alinierile de tipul N:M reprezintă situațiile în care M fraze din limba engleză au fost traduse cu N fraze în limba respectivă. Un caz particular îl reprezintă situațiile de omisiune în traducere (0:M) sau de inserare de text fără corespondență în original (N:0).

3.3. Experimentul inițial

Textul original "1984" conține 7.069 leme diferite, iar dicționarul multilingv extras prin metoda descrisă în prima parte a acestei lucrări conține 1.233 de intrări. Aceste intrări au fost reținute respectând condiția ca un articol lexical din limba engleză să aibă traduceri (eventual multiple) în cât mai multe limbi țintă. Condi

impusă dicționarului multilingv este foarte restrictivă, având în vedere că majoritatea dicționarilor bilingve extrase automat conțin între 6000 și 7000 de intrări. Intrări tipice (parțiale) în dicționarul multilingv sunt ilustrate în figura 7. O informație suplimentară, ce nu apare în exemplificarea din figura 7, este mulțimea tuturor unităților de traducere din corpusul paralel în care cuvântul englezesc a fost tradus prin echivalenții săi listați în dicționar. Dintre aceste intrări, au fost selectate 845 pentru care s-au găsit una sau mai multe traduceri în toate limbile. Dintre acestea, s-a ales o mulțime de 33 de substantive, acoperind toate gamele de frecvență și ambiguitate, cu care s-a realizat experimentul ale cărui rezultate au fost validate de experți umani [34].

Engleză	Categorie	Bulgară	Cehă	Estoniană	Maghiară	Română	Slovenă
finally	R	Накпaa	nakonec konecne	lopuks viimaks	vegul	în_cele_di n_urmă până_la_ur mă	koncen nazadnje
wealth	N	6opaTCTBo Guaro	bohatsvi	joukus rikkus	jolet gazdagság	avuție bogăție	blaginja bogastvo

Figura 7: Exemple de echivalenți de traducere identificați în corpusul paralel "1984"

Pentru fiecare substantiv din acest eșantion au fost extrase toate frazele englezești în care apare, împreună cu toate frazele corespunzătoare din celelalte limbi și pentru fiecare ocurență a sa a fost construit un vector binar reprezentând toate traducerile posibile ale cuvântului respectiv. O valoare 1 în poziția n a acestui vector semnifică faptul că acea ocurență a fost tradusă prin cuvântul ce reprezintă a A ?-a traducere posibilă. O valoare 0 semnifică faptul că a n -a traducere posibilă nu a fost folosită. De pildă pentru substantivul "wealth" (vezi figura 7) au fost depistate 11 traduceri posibile (2 în bulgară, estoniană, maghiară, română și slovenă, 1 în cehă). Un vector asociat oricărei ocurențe a lui *wealth* va avea prin urmare 11 poziții. Astfel, dacă a m -a apariție în textul original al romanului "1984" a cuvântului *wealth* are atașat vectorul 10101010101 acest lucru semnifică faptul că în varianta bulgărească el a fost tradus cu *6oeamcmeo*, în cea cehă cu *bohatsvi*, în cea estoniană cu *rikkus*, în cea maghiară cu *gazdagság*, în cea română cu *bogăție* iar în cea slovenă cu *bogastvo*. Vectorii astfel definiți au fost prelucrați cu un algoritm de clasificare de tip aglomerativ [36], clasele rezultate fiind

considerate a reprezenta sensuri distincte în care cuvântul curent a fost folosit de a lungul romanului. Clasele produse de algoritm au fost comparate cu clasele rezultate prin dezambiguizarea manuală efectuată, independent, de 2 vorbitori nativi ai limbii engleze. Dezambiguizarea manuală a fost efectuată utilizând numerotarea sensurilor din WordNet 1.6.

Pentru a putea compara rezultatele produse de dezambiguizatorii umani (numiți în continuare adnotatori) cu cele produse de algoritmul nostru, datele au fost normalizate în felul următor: pentru fiecare adnotator și pentru algoritm, fiecare din cele 33 de cuvinte a fost reprezentat printr-un vector binar de lungime n , unde n este numărul de ocurențe ale cuvântului în tot corpusul. Pozițiile în vector reprezintă o asignare de tip "DA/NU" indicând dacă ocurența respectivă a fost clasificată la fel de către adnotatori, respectiv algoritm. Rezultatele acestui prim experiment sunt rezumate în tabelul din figura 8 indicând procentul de acord între clasificările propuse de algoritm și cele ale fiecărui adnotator, acordul dintre cei doi adnotatori și acordul dintre toți cei trei clasificatori.

Algoritm/Adnotator 1	66.7%
Algoritm /Adnotator 2	63.6%
Adnotator 1/Adnotator 2	76.3%
Algoritm /Adnotator 1/ Adnotator 2	53.4%

Figura 8: Concordanța între diferite clasificări

3.4. Cel de-al doilea experiment

Rezultatele primului experiment au arătat că metoda discriminării sensurilor folosind echivalenții de traducere este foarte competitivă, acuratețea procesului fiind comparabilă (și uneori superioară) cu performanțele obținute de cercetători ce au folosit ca referință același dicționar (Wordnet). Mai mult, diferențele de acord asupra clasificării dintre cei 2 adnotatori pe de o parte și dintre fiecare adnotator și algoritm pe de altă parte este de numai 10-13%, ceea ce înseamnă că noua metodă este foarte competitivă în raport cu scorurile obținute în alte experimente.

Pentru a valida aceste rezultate empirice, în cea de a doua fază a acestui experimentului a fost luat în considerare un număr dublu de substantive (76) din clasele "dificile", adică cu grad de ambiguitate mare, atât din clasa celor abstracte și a celor concrete (de exemplu, "thought", "stuff", "meaning", "feeling" respectiv "hand", "boot", "glass", "girl" etc). Am ales acele substantive care au apărut puțin de 10 ori în corpus (pentru a elimina efectul de "insuficiență a datelor") și plus care au cel puțin 5 traduceri în cele 6 limbi țintă. Restricția de 10 apariții

S S K ^ ^ W Pe care , - am impus procesului de ^ g

$$L L (T_T, T_s) = 2 * \sum_{i=1}^7 n_i > 18$$

aj. - doi vorbitori nativi
 e .chetată, în mod independent de 5; c l a s a * f u v i n t e l o r s e . e t i c e t a t e * f o s t
 a < c i . I n t a b e l a d i n f i g u r a 9 s u n t r e l a t e S P , s i a . 9 0 r i t m u l d i s c u t a t
 adnotatori: rezumate datele ș. rezultatele de acord între cei 4

Nr. de cuvânte	76
Nr. ocurențe	2399
Număr mediu de ocurențe/cuvânt	32
Nr. de sensuri găsite de adnotatorul 1	241
Nr. de sensuri găsite de adnotatorul 2	280
Nr. de sensuri găsite de adnotatorul 3	213
Nr. de sensuri găsite de adnotatorul 4	232
Nr. de sensuri găsite împreună de toți adnotatorii	345
Numărul mediu de sensuri pe cuvânt	4.53
Procent de acord între adnotatori	
Acord total (4/4)	54.27
75% acord total (3/4)	28.13
50% acord total (2/4)	16.92
Dezacord total	0.66

Figura 9: Datele experimentului și acordul între 4 adnotatori umani independenți

Rezultatele produse de algoritmul de clasificare și clasificările realizate de adnotatori prin asignarea sensurilor din Wordnet1.6 au fost de data aceasta normalizate în mod diferit, prin ignorarea etichetei puse de adnotatori și considerând doar clasele rezultând din această etichetare. Pentru a clarifica acest aspect să urmărim modul în care doi dintre adnotatori au dezambiguit cele 7 ocurențe ale cuvântului "youth":

Ocurența nr.	1	2	3	4	5	6	7
Adnotatorul 1	1 ³	1	6	3	6	3	1
Adnotatorul 2	1 ²	1	4	2	6	2	1

Figura 10: Acordul de clasificare pentru cuvântul " youth" în umani independenți

Acordul între cei doi adnotatori este doar de 43% (doar au asignate sensuri consensuale); totuși, ambii adnotatori au 1, 4, și 6 ca având același sens, deși primul le-a etichetat cu sensul 1 în timp ce al doilea le-a etichetat cu sensul 2. Dacă în sensul de clasificarea celor 3 ocurențe este consistentă, în sensul că a decis că ele au același sens. Acordul de clasificare se dublează datele sunt mult mai ușor de comparat cu rezultatele produse de

în acest al doilea experiment am luat în considerare momentul optim de oprire a clasificării aglomerative. În primul rând folosit o distanță minimă predefinită, pentru determinarea numărului de clase între care se realizează discriminarea. Această soluție nu ține cont de proprietățile individuale ale cuvintelor (numărul maxim de sensuri din Wordnet, frecvența de apariție a cuvântului, numărul mediu de sensuri pe cuvânt a primit cuvântul în corpus). Noul algoritmul de clasificare a fost conceput să-și calculeze un număr optim de clase², optimalitatea fiind măsurată prin numărul mediu de clase identificate de adnotatori. Drept criteriu am folosit distanța minimă dintre clasele existente la fiecare pas de agregare, clasele cu cea mai mică distanță relativă sunt agregate într-o clasă mai mare. Procesul începe cu fiecare ocurență într-o clasă separată și oprește când distanțele relative între clasele existente este egală cu o valoare predefinită. Distanța dintre două clase se calculează pe baza vectorilor (centrozii) ai celor două clase (evident depinzând de cuvânt și de numărul de sensuri ale cuvântului clasificat):

$$dist(v, V2) = \sum_{i=1}^n |E_i(v) - V_2(i)|$$

¹ Singurul dezacord rămas constă în faptul că Adnotatorul 1 consideră că ocurența 1 și 6 au același sens, în timp ce Adnotatorul 2 atribuie un sens diferit pentru acestea, realizând o discriminare mai fină între sensurile celor două ocurențe.

² În principiu, limita superioară a numărului de sensuri pe care îl poate avea un cuvânt în engleză într-un text este dată de numărul de sensuri listate în Wordnet. În realitate, de așteptat însă nu există în corpusul nostru nici un exemplu polisemantic să fi apărut cu toate sensurile din WordNet.

Cele mai bune rezultate în discriminarea automată au fost obținute pe cale experimentală, impunând drept criteriu de oprire a algoritmului condiția:

$$\frac{mindist(k) - mindist(k+1)}{mindist(k+1)} \ll 1$$

în care $mindist(k)$ reprezintă distanța minimă între clasele existente la pasul k de aglomerare.

Pentru medierea opiniilor adnotatorilor am definit o adnotare de referință reprezentând clasificarea majoritară între cei 4. În cazul egalității de voturi, adnotatorul care a fost în cele mai multe cazuri de aceeași opinie cu majoritatea a impus clasa. Folosind această clasificare mediată și raportând-o la clasificarea produsă de algoritm pentru cele 76 de cuvinte, am analizat diferențele de clasificare, considerate ca fiind erori. Marea majoritate a erorilor de clasificare pentru cele 2399 de ocurențe au apărut în cazul cuvintelor pentru care distribuția sensurilor este foarte inegală; ca urmare am adăugat algoritmului o fază suplimentară de postprocesare, în care clasele cu un număr mult mai mic de ocurențe decât clasa cu cele mai multe ocurențe au fost incorporate în ultima. Raportul minim între numărul de ocurențe al celei mai mari clase și numărul de ocurențe din clasele potențial absorbabile în cea dintâi a fost ales empiric ca fiind 10^1 . Motivația acestei euristici constă în constatarea făcută de mai mulți cercetători în domeniul lingvisticii corpusului (fapt sugerat chiar de Zipf cu peste 50 de ani în urmă) că utilizarea frecventă a unui cuvânt într-un text omogen tinde să-i păstreze sensul.

Cu această nouă euristică încorporată, algoritmul de clasificare a atins cifra de 74,6% acord cu clasificarea mediată. În [35] sunt prezentate alte variante ale algoritmului care au condus prin evaluarea empirică la versiunea sa finală. Clasele produse de fiecare pereche de clasificatori (om sau mașină) au fost evaluate printr-un algoritm ce calculează alinierea claselor astfel încât intersecția lor să fie maximală. Diferențele dintre două clase astfel aliniate au fost considerate dezacorduri de clasificare. Scorul de acord a fost calculat ca fiind raportul dintre suma numărului de ocurențe comune pentru fiecare clasă aliniată și numărul total al ocurențelor cuvântului respectiv. În tabela din figura 10 este exemplificat modul de calcul al acordului dintre clasificarea produsă de algoritm și clasificarea mediată a adnotatorilor pentru cuvântul *movement*. Acesta a apărut în text de 40 de ori. Atât algoritmul cât și cei patru adnotatori au identificat 4 sensuri distincte în care acest cuvânt a fost utilizat. Așa cum se vede din figura 10, cea mai numeroasă clasă (clasa 1) conține în clasificarea mediată 28 dintre cele 40 de apariții ale cuvântului *movement*, în timp ce clasa corespondentă creată de algoritm conține doar 25 de ocurențe. Dintre acestea, 24 sunt comune cu cele din clasa 1 a

definiția anterioară a scorului de acord, clasificării mediate. În conformitate cu definiția anterioară a scorului de acord, pentru acest exemplu rezultă următoarele cazuri:

CLASA	1	2	3	4	2.
Clasificare mediată	28	6	3	3	40
Clasificare algoritmică	25	7	6	2	40
Intersecție	24	6	3	1	34
Precizie	85%				

Figura 11: Clasificarea mediata si cea produsă de algoritm pentru cuvântul *movement*

3.5. Rezultate

Rezultatele obținute cu ultima variantă a clasificatorului în cel de-al doilea experiment sunt sintetizate în tabelul din figura 12. Tabelul indică procentul de acord între diverse clasificări: 1, 2, 3, 4, reprezintă clasificările realizate de adnotatorii umani, M reprezintă clasificarea mediată a clasificatorilor umani, A reprezintă clasificarea produsă de algoritm, iar B este referința de bază (baseline) care presupune toate ocurențele unui cuvânt ca având același sens.

	1	2	3	4	M	A
B	71.1	65.1	76.3	74.1	75.5	81.5
1		78.1	75.6	83.1	88.6	74.4
2			71.3	75.9	82.5	66.9
3				77.3	82.1	77.1
4					90.4	75.9
M					1	77.3

Figura 12: Acorduri între diverse clasificări

Tabela arată că acordul între adnotatorii umani comparat cu cel dintre algoritm și adnotatorii umani (cu excepția unuia dintre ei (4), pe care îl suspectăm că a văzut clasificările celorlalți trei și în consecință și-a revizuit unele decizii) diferă substanțial. Acest lucru demonstrează (cel puțin în raport cu datele experimentului nostru) că dezambiguizarea automată este comparabilă în acuratețe cu cea efectuată de adnotatorii umani. Diferența fundamentală constă

- faptul că programul a terminat în circa 2 minute clasificarea pentru care adnotatorilor le-au trebuit între 4 și 5 săptămâni.

Experimentul descris a evaluat dezambiguizarea automată a cuvintelor englezești pornind de la traducerea lor în celelalte 6 limbi. Această direcționare a fost impusă doar de disponibilitatea pentru limba engleză a textului dezambiguizat de experți umani (vorbitori nativi ai limbii engleze). Întrucât algoritmul de clasificare nu depinde în nici un mod de limba pentru care se realizează dezambiguizarea (limba țintă) și nici de limbile martor în raport cu care se face acest proces, rezultă că exact același procedeu descris până aici poate fi folosit pentru dezambiguizarea cuvintelor românești folosind echivalenții lor de traducere în engleză, bulgară, cehă, estoniană, maghiară și slovenă, ori pentru dezambiguizarea cuvintelor bulgărești pe baza echivalenților lor de traducere în celelalte 6 limbi. Întrucât sensul este (în principiu) un invariant al traducerii, nu pare a se justifica și pentru celelalte limbi efortul de adnotare umană făcut pentru limba engleză. Este rațional a presupune că rezultate similare (raporturi relative) s-ar obține indiferent de limba țintă și de limbile martor.

Să mai menționăm și faptul că există o anumită corelație (factorul Spearman - 0.51) între numărul de sensuri în Wordnet ale unui cuvânt și nivelul de acord între diferitele clasificări ale ocurențelor sale. Cele mai scăzute scoruri de acord au fost obținute pentru "line" (29 sensuri), "step" (10), position (15), "place" (17) și "corner" (11). Acorduri perfecte s-au obținut pentru majoritatea cuvintelor cu mai puțin de 5 sensuri, ca de exemplu "hair" (5), "morning" (4), "sister" (4), "tree" (2), and "waist" (2) care toate au fost considerate, atât de adnotatori cât și de algoritm, a fi fost folosite cu un singur sens în tot textul. Pe de altă parte, gradul de acord pentru câteva cuvinte cu mai puțin de 5 sensuri ("rubbish" (2), "rhyme" (2), "destruction" (3) și "belief (3)) a fost semnificativ mai mic decât media pentru toate perechile de clasificări (adnotator-adnotator, adnotator-algoritm). Concluzia a fost că pentru unele cuvinte, distincțiile de sens sunt atât de fine în Wordnet, încât chiar vorbitorii nativi (și cu atât mai mult algoritmul de clasificare) nu pot face diferențieri sistematice de sens ale diferitelor ocurențe ale acestor cuvinte. O astfel de hiperdiferențiere a sensurilor este în imensa majoritate a cazurilor irelevantă pentru aplicațiile de prelucrare a limbajului natural.

4. Concluzii

Rezultatele experimentelor noastre arată că acuratețea discriminării sensurilor pe baza echivalenților de traducere extrași din corpusuri paralele este comparabilă cu cea produsă de adnotatori umani. Întrucât abordarea noastră este complet automatizată ea poate fi folosită la crearea de volume mari de texte, având discriminate sensurile cuvintelor polisemantice. Utilizarea experților umani este prohibitivă sub aspectul costului și al timpului de realizare a unei asemenea

sarcini, iar procentajul suplimentar de acuratețe, presupus de activitatea umană, este prea mic pentru a justifica procedurile manuale.

Metoda pe care am descris-o în această lucrare nu etichetează clasele de ocurențe ale unui cuvânt cu un număr de sens ales dintr-un inventar prescris de sensuri iar majoritatea aplicațiilor de prelucrare a limbajului natural (de pildă clasificarea textelor, regăsirea informațiilor, rezumarea automată etc.) nici nu au nevoie de această informație suplimentară; pentru aceste tipuri de aplicații este suficient a decide că două sau mai multe ocurențe ale unui cuvânt sunt folosite același sens sau nu. O etichetare convențională a sensurilor identificate pentru un anumit cuvânt ar putea să se bazeze pe frecvența sensurilor respective (sensul corespunzând clasei cu cele mai multe ocurențe). Evident o astfel de etichetare depinde de registrul lingvistic al textului pe baza căruia se identifică sensurile distincte.

O direcție foarte promițătoare [37], [38], [39] o constituie utilizarea metodologiei prezentată aici în construcția și validarea ontologiilor multilingve de tip EuroWordNet. Folosind echivalenții de traducere și clasificarea ocurențelor echivalente din punctul de vedere al sensului se poate verifica dacă proiectia interlinguală a două sau mai multe dicționare semantice este corectă. Această presupune că sensurile cuvintelor extrase ca echivalenți de traducere ai cuvintelor englezești dezambiguizate să fie puse în corespondență cu același concept interlingual aparținând Indexului Interlingual (ILI - vezi [39]- în acest volum). În cazul contrar (echivalenții de traducere sunt puși în corespondență cu concepte interlinguale diferite) este fie vorba de o eroare propriu-zisă de proiectie conceptuală într-unui sau mai multe dintre dicționarele semantice aliniate sau conceptele interlinguale sunt atât de apropiate semantic încât se poate propune unificarea lor într-un concept mai general cu lexicalizare în mai multe limbi. Aceasta este esența conceptului de "soft-clustering" definit în comunitatea EuroWordNet. Față de identificarea prin metode statistice a conceptelor interlinguale foarte apropiate semantic, analiza prin metoda "echivalenților de traducere și a discriminării sensurilor" a proiecțiilor sensurilor făcute de lexicografuli profesioniști peste o mulțime de sensuri conceptualizate în ILI este mult mai robustă. Experimentele preliminare discutate în [37] au arătat că în diferite limbi pentru care se realizează o ontologie lexicală multilingvă (bulgară, cehă, greacă, română, sârbă, turcă) există dificultăți identice de proiectie conceptuală a sensurilor unor cuvinte din limbile considerate. Faptul că aceleași concepte interlinguale creează același tip de dificultate în proiecția sensurilor unor cuvinte aparținând unor limbi foarte diferite indică cu claritate că acele concepte trebuie generalizate.

Un alt aspect care merită subliniat este că metodologia prezentată aici este corelată cu existența a tot mai multe dicționare semantice de tip Wordnet, ce aplică la principiul EuroWordNet de aliniere la Indexul Interlingual, va permite dezvoltarea

de corpusuri adnotate semantic (de tipul SemCor) pentru orice limbă. Tranzitivitatea relațiilor de tip "EQ-SYN" folosite în proiecția sinseturilor unui wordnet monolingv peste ILLI, corelată cu echivalența de traducere (relație tot între sensuri) extrasă dintr-un corpus paralel, în care textul dintr-una din limbi este adnotat semantic, permite importul adnotărilor în toate celelalte limbi. Deoarece limba din care se importă adnotarea semantică nu este relevantă pentru această procedură, rezultă că eforturile depuse de-a lungul timpului în crearea celor câteva corpusuri cu adnotare semantică pentru limbile "mari" pot fi valorificate pentru orice altă limbă în care există (sau se creează) traduceri ale textelor din corpusurile adnotate. Mai mult, se poate imagina crearea unui consorțiu multilingv care să aleagă un corpus paralel în cât mai multe limbi cu scopul de a-l adnota semantic. Prin adnotarea independentă, în fiecare limbă, a unor porțiuni distincte din corpusul paralel, folosind o metodologie de genul celei prezentate în această lucrare (și desigur având un dicționar semantic multilingv de tip EuroWordNet) adnotările secțiunilor monolingve vor putea fi importate în secțiunile corespunzătoare ale tuturor celorlalte texte monolingve, în final putându-se obține adnotarea semantică, consistentă, a întregului text din fiecare limbă a corpusului paralel.

Mulțumiri

Rezultatele prezentate în această lucrare sunt rodul mai multor proiecte internaționale de cercetare desfășurate la Institutul de Inteligență Artificială, alături de colegii Ana Măria Barbu, Eduard Barbu, Radu Ion, Cătălin Mititelu, Octavian Popescu. De asemenea/colaborarea cu Nancy Ide de la Universitatea Vassar din Poughkeepsie, SUA, și cu Tomaz Erjavec de la Institutul "Jozef Štefan" din Ljubljana, Slovenia, parteneri în proiectele amintite, a fost și este extrem de productivă. Tuturor le aduc aici cuvenitele mulțumiri.

Referințe bibliografice

- [1] Melamed, D. - "Empirical Methods for Exploiting Parallel Texts", MIT Press, 2001, 373p.
- [2] Gale, W.A., K.W. Church, - "Identifying word correspondences in parallel texts". In Fourth DARPA Workshop on Speech and Natural Language, 1991, 152:157
- [3] Smadja, F., K.R. McKeown, and V. Hatzivassiloglou - "Translating collocations for bilingual lexicons: A statistical approach". *Computational Linguistics*, 22/1, 1996, 1:38.
- [4] Brown, P., Della Pietra, S. A., Della Pietra, V. J., Mercer, P. R. - "The mathematics of statistical machine translation: parameter estimation". *Computational Linguistics* 9(2): 263-311, 1993.
- [5] Kupiec, J. - "An algorithm for finding noun phrase correspondences in parallel corpora". In *Proceedings of the 31st Annual Meeting of the Association for Computational Linguistics*, 1993,17:22
- [6] Hiemstra, D. - "Deriving a bilingual lexicon for cross language information retrieval". In *Proceedings of Gronics*, 1997, 21:26
- [7] Tufiș, D., Barbu, A.M. - "Automatic Learning of Translations from Parallel Texts". "Romanian Journal on Information Science and Technology", Romanian Academy, vol.4, no. 3-4, 2001b, 325:351.
- [8] Tufiș, D., Barbu, A.M. - "Revealing translators knowledge: implications for in constructing practical translation lexicons for language processing", în *International Journal of Speech Technology*. John Benjamins Publishers, no.5, 2002, 199:209.
- [9] Tufiș, D. - "Parțial translations recovery in a 1:1 word-alignment task". RACAI Research report, June, 2001b, 32pp.
- [10] Gale, W.A., K.W. Church - "A Program for Aligning Sentences from Parallel Corpora". în *Computational Linguistics*, 19(1), 1993, 75:102
- [11] Tufiș, D. - "Tiered Tagging and Combined Classifiers". în *Text, Speech and Dialogue*, Lecture Notes in Artificial Intelligence, Springer, 1999, 29:33.
- [12] Brants, T. - "TnT - A Statistical Part-of-Speech Tagger", în *Sixth Applied Natural Language Processing Conference, ALPNET 2000* - May 3, 2000, Seattle, WA, 2000
- [13] Varadi, T. - The Hungarian National Corpus, *Proceedings of the 19th International Conference on Computational Linguistics*, Palma de Majorca, Spain, 2002, 385:389.
- [14] Tufiș, D. - "A cheap and fast way to build useful translation lexicons". *Proceedings of the 19th International Conference on Computational Linguistics COLING2002*, Taipei, China, 2002, 246:251.
- [15] Tufiș, D., Dienes, P., Oravecz, C., Váradi T.,*- "Principles and Design for Tiered Tagging of Hungarian" *Proceedings of the 19th International Conference on Language Resources and Evaluation*, Athens, Greece, 2000, 1421:1426
- [16] Hinrichs, H., Trushkina, J. - "Forging Agreement: Morphology of Noun Phrases", *Proceedings of the Workshop on Treebanking and Theories 2002*, Sozopol, Bulgaria, 2002, 1:18.

- [17] Erjavec, T. - "An Experiment in Automatic Bi-lingual Lexicon Construction from a Parallel Corpus", Proceedings of the 7th TELRI International Seminar on Corpus Linguistics, Dubrovnik, Croatia, 2002.
- [18] Erjavec T., Ide, N. - "The Multext-East corpus". în *Proceedings of First International Conference on Language Resources and Evaluation*, Granada, Spain, 1998,971:974
- [19] Tufiş, D., Barbu, A.M., Pătraşcu, V., Rotariu, G., Popescu, C. - "Corpora and Corpus-Based Morpho-Lexical Processing", în D. Tufiş, P. Andersen (eds.) "Recent Advances in Romanian Language Technology", Editura Academiei, 1997, 35:56
- [20] Kay, M., Roscheisen, M. - "Text-Translation Alignment". în *Computational Linguistics*, 19/1, 1993, 121:142
- [21] Brew, C, McKelvie, D. - "Word-pair extraction for lexicography", 1996, <http://www.ltg.ed.ac.uk/~chrisbr/papers/nemplap96>
- [22] Tiedemann, J. - "Extraction of Translation Equivalents from Parallel Corpora", în *Proceedings of the 11th Nordic Conference on Computational Linguistics*, Center for Sprogteknologi, Copenhagen, 1998, <http://stp.ling.uu.se/~ioerg/>
- [23] Ahrenberg, L., M. Andersson, M. Merkel - "A knowledge-lite approach to word alignment", în [40].
- [24] Cristea, D., Dima, G. E. - "An Integrating Framework for Anaphora Resolution", Journal on *Information Science and Technology*, Romanian Academy Publishing House, Bucharest, voi. 4, no. 3, 2001, 273:292.
- [25] Fellbaum C. - Wordnet: An Electronic Lexical Database, MIT Press, 1998, 423p.
- [26] Resnik, P. and Yarowsky, D. - A perspective on word sense disambiguation methods and their evaluation. *ACL-SIGLEX Workshop Tagging Text with Lexical Semantics: Why, What, and How?* Washington, D.C., 1997, 79:86.
- [27] Gale, W. A., Church, K. W., Yarowsky, D. (1993). A method for disambiguating word senses in a large corpus. *Computers and the Humanities*, 26, 415:439.
- [28] Dagan, I., Itai, A., Schwall, U. - Two languages are more informative than one. *Proceedings of the 29th Annual Meeting of the ACL*, 18-21 Berkeley, California, 1991, 130:137.
- [29] Dagan, I., Itai, A. - Word sense disambiguation using a second language monolingual corpus. *Computational Linguistics*, 20/4,1994, 563:596.
- [30] Dyvik, H. (1998). Translations as Semantic Mirrors. *Proceedings of Workshop Multilinguality in the Lexicon II, ECAI98*, Brighton, UK, 1998, 24:44.
- [31] Ide, N. - Cross-lingual sense determination: Can it work? *Computers. and the Humanities*, 34/1-2, 1999, 223:234.
- [32] Resnik, P. and Yarowsky, D. - Distinguishing systems and distinguishing senses: New evaluation methods for word sense disambiguation. *Journal of Natural Language Engineering*, 5(2), 2000,113:133.
- [33] Miller, G. A., Beckwith, R. T. Fellbaum, C. D., Gross, D. and Miller, K. J. WordNet: An on-line lexical database. *International Journal of Lexicography* 3/4,1990,235:244.
- [34] Ide, N., Erjavec, T., and Tufiş, D. (2001). Automatic sense tagging using parallel corpora. *Proceedings of the Sixth Natural Language Processing Pacific Rim Symposium, Tokyo*, 2001,83:89.
- [35] Ide, N., Erjavec, T., Tufiş, D. - "Sense Discrimination with Parallel Corpora" Proceedings of the SIGLEX Workshop on Word Sense Disambiguation: Recent Successes and Future Directions. ACL2002, July Philadelphia, 2002, 56:60
- [36] Stolcke, A. - Cluster - 2.9. <http://www.icsi.berkeley.edu/ftp/global/pub/stolcke/software/cluster-2.9.tar.Z>, 1996.
- [37] Tufiş, D. - "Interlingual alignment of parallel semantic lexicons by means of automatically extracted translation equivalents", Proceedings of the 7th TELRI International Seminar on Corpus Linguistics, Dubrovnik, Croatia, 2002.
- [38] Tufiş, D., Cristea, D. - "Methodological issues in building the Romanian Wordnet and consistency checks in Balkanet", în Proceedings of LREC2002 Workshop on Wordnet Structures and Standardisation, Las Palmas, Spain, 2002, 35-41.
- [39] Tufiş, D., Cristea, D. - "RO-BALKANET - ontologie lexicalizată, în contextul multilingv, pentru limba română", 2002, în acest volum.
- [40] Veronis, J. (ed), Parallel Text Processing. Text, Speech and Language Technology Series, Kluwer Academic Publishers Voi. 13, 2000.

Referențialitate și cursivitate în relație cu structura de discurs

Dan CRISTEA

Universitatea "Al.I.Cuza" Iași, Facultatea de Informatică
Academia Română, Institutul de Informatică Teoretică - Filiala Iași

dcristea@infoiasi.ro

1. Introducere

În ultimii 25 de ani s-a studiat enorm pentru a se înțelege ce anume face dintr-un text (considerat o secvență de propoziții sintactic corecte) să fie un discurs, așadar de ce un discurs e coerent și ce elemente îi atribuie coeziune. Dintre teoriile computaționale ale discursului, trei au avut o influență covârșitoare asupra dezvoltărilor ultimilor ani din acest domeniu: teoria structurilor retorice, teoria stărilor atenționale și teoria centrelor.

Dezvoltată inițial din perspectiva generării textelor, teoria structurilor retorice (*rhetorical structure theory*, de aici încolo RST), a fost elaborată de Mann și Thompson ca o teorie a organizării textelor [Mann, Thompson, 1988; Hovy, 1988; Scott, de Souza, 1990]. Ea caracterizează structura de discurs în termeni de relații ce leagă părți componente ale textului. Unitatea elementară de discurs în RST este, de regulă, o propoziție, așadar o întindere textuală în care se formulează o predicăție. O structură de discurs este descrisă de o **schemă**. Ea grupează o secvență de unități, sau de unități și scheme, sau o secvență de scheme. O schemă poate fi asemuită cu o regulă a unei gramatici, ea relevând structura de constituenți a unui compus. O schemă constă dintr-o **relație** (27 în RST) care leagă două sau mai multe întinderi de text, fiecare dintre ele având, la rândul său, o structură (constituenții schemei). Un discurs este fie o **unitate**, care este o întindere de text elementar, fără structură, fie o schemă (un text mai lung decât o singură unitate și care prezintă o structură). Relațiile pot fi de două tipuri: **hipotactice** - dacă argumentele sunt constituenți neegali ca importanță și **paratactice** (sau echinucleare) - dacă constituenții pe care-i agregă sunt egali ca importanță. Între constituenții uniți de relațiile hipotactice există întotdeauna unul singur mai important, numit **nucleu**, ceilalți fiind numiți **sateliți**. La relațiile paratactice, prin convenție, se consideră că toți constituenții sunt nucleari. Satelitul este, în general, mai susceptibil de a fi schimbat sau eliminat complet decât

nucleul, fără ca, prin aceasta, înțelesul discursului să se modifice. Dimpotrivă, înlocuirea sau ștergerea nucleului este o opțiune mult mai drastică, care poate duce la denaturări ale înțelesului. Relațiile hipotactice sunt, în general, cele intenționale, în care o întindere de text comunică un scop și celelalte exprimă subscopuri ce completează, dezvoltă etc. scopul principal. Pe de altă parte, relațiile paratactice sunt, în general, de natură informațională, simetrice, neputându-se stabili dacă, sau care, componentă predomină.

În RST accentul este pus pe performanța retorică: prin ce mijloace un scriitor (sau vorbitor) reușește să convingă un cititor (ascultător) de intențiile pe care le are de comunicat. Ca produs secundar al liniei principale de investigare în RST, multe eforturi de cercetare care au succedat elaborarea teoriei s-au concentrat asupra îmbunătățirii și structurării setului de relații propus inițial [Rambow, 1993]. Pe de altă parte, pare extrem de convenabil, inclusiv din punct de vedere computațional, să vedem discursul reprezentat ca un arbore, în care nodurile terminale să reconstituie, în secvența lor, textul. Cu toate acestea RST nu aduce nici o lumină în privința vreunei legături care ar exista între structură și referențialitate. RST este deci o teorie asupra structurii globale a discursului.

Teoria stărilor atenționate (*attentional state theory*, AST) [Grosz, Sidner, 1986] reprezintă o dezvoltare a liniei de cercetare în discurs dominată de Barbara Grosz și Candace Sidner asupra manierei în care focarul ori centrul de discurs (*focus* în engleză) se modifică pe parcursul derulării textului și a recunoașterii intențiilor comunicate de discurs [Grosz, 1981; Sidner, 1983]. Grosz și Sidner nu cred că varietatea atât de mare a intențiilor ce pot fi comunicate de un discurs poate fi condensată într-un număr fix de șabloane retorice exprimate ca relații, cum sugerează RST sau tentative similare acesteia. Teoria se dorește a fi un model formal, care se distanțează de detaliile ce ar putea fi asociate participanților la discurs. Realizând proiecții corespunzătoare utilizatorului de limbaj, însoțite de detalii specifice, ea s-ar putea regăsi atât în construcția unui sistem automat cât și într-o teorie psihologică, ambele consumatoare de limbaj natural. Deși recunoaște însemnătatea mesajului transmis de un discurs, teoria nu abordează problema înțelesului discursului și a manierei în care acesta poate fi dedus din elementele constitutive ale textului. Ea este, primordial, o teorie a structurii discursului, prin aceasta plasându-se la baza oricărei tentative de a aborda problema construirii sensului.

Conform lui Grosz și Sidner intențiile joacă rolul principal în explicarea structurii discursului, în timp ce dinamica atenției joacă rolul principal în explicarea interpretării discursului. Structura discursului are trei componente distincte, dar strâns corelate:

- o **structură lingvistică**, care face ca una sau mai multe propoziții, exprimări (*utterance*) să fie agregate într-un segment de discurs iar limitele dintre segmente să fie indicate de expresii lingvistice, intonație,



schimbări ale timpului și aspectelor verbelor. Segmentul de discurs are însă o definiție recursivă: un segment poate îngloba alte segmente acestea pe altele ș.a.m.d.;

- o **structură intențională**, care face să vedem discursul ca având un scop global (scopul discursului - SD), scopul fundamental al vorbitorului/scriitorului la emiterea discursului. Fiecare segment are un scop al segmentului (scopul segmentului de discurs - SSD), care este un subscop al scopului segmentului din care face el parte. Intuitiv SSD specifică cum contribuie respectivul subsegment la realizarea scopului segmentului din care face el parte. Teoria admite că nu există o listă finită de scopuri ale discursului, care să facă posibilă comparație cu lista categoriilor gramaticale, de exemplu. Conform teoriei, două relații structurale sunt suficiente pentru a compune structura discursului: **relația de dominare** (dacă SSD_i domină SSD_j atunci SSD_j contribuie la SSD_i , sau SSD_j este intenționată să satisfacă parțial SSD_i) și **relația de satisfacere-precedență** (SSD_i satisfacă-precede SSD_j dacă SSD_i trebuie satisfăcut înainte de SSD_j);
- o **stare atenționată**, prin care se asociază fiecărui segment al discursului un spațiu al entităților aflate în centrul atenției. Starea atenționată reprezintă o trăsătură dinamică a discursului, păstrând **obiecte, proprietăți și relații** ce sunt importante în fiecare moment al interpretării discursului. Starea atenționată e modelată printr-un **set spații ale centrelor atenției**, în timp ce schimbările ce pot avea loc în starea atenționată sunt restricționate de un set de reguli de tranziție care arată condițiile de adăugare și ștergere a spațiilor. Colecția tuturor spațiilor centrelor de atenție ce sunt disponibile în fiecare moment al interpretării unui discurs formează o structură a atenției ce are o dinamică unei **stive** și care ar fi capabilă să explice procesele implicite în interpretarea discursului, inclusiv accesibilitatea referențială: domeniul în care trebuie căutate entitățile de discurs referențiale este segmentul corespunzător stării atenționale aflate în vârful stivei și cel al stărilor aflate în stivă.

Structura recursivă a segmentului de discurs din AST permite și acceptarea unei reprezentări arborescente, în cadrul căreia cele două relații între segmente, de dominare și de satisfacere-precedență, nu sunt altceva decât relațiile firești pe orice structură de arbore: cea dintre părinte și orice fiu al său respectiv, cea de ordine dintre frați. AST se constituie, așadar, într-o teorie globală asupra structurii și a coeziunii discursului.

Cercetători precum Moser și Moore [1996] sau Marcu [1999] pun în evidență similarități semnificative între AST și RST, inclusiv în ceea ce privește

maniera de reprezentare prin arbori a structurii de discurs, ceea ce permite combinarea puterii de reprezentare, mai fine în RST, datorită proliferării relațiilor, cu implicațiile pe care structura le poate avea asupra referențialității, puse în evidență de AST. Utilizând structura de segmente și stiva, ca mecanism de prelucrare, AST propune o manieră de a rezolva accesibilitatea referințelor anaforice printr-o transparentă pe verticală, de sus în jos, de-a lungul stărilor atenționale ce se află la un moment dat în stivă. Reprezentarea prin segmente din AST are însă o slăbiciune: modelul stivă nu poate reflecta relația de dominare atunci când scopul dominat corespunde unui segment care apare în text înaintea celui care domină [Ide, Cristea, 2000]. Să remarcăm că defectul este unul de granularitate pentru că identificarea segmentului dominat ce precede pe cel dominator cu însuși segmentul dominator elimină problema. AST nu e, așadar, capabilă să reprezinte segmente având o granularitate oricât de fină: coborând de la o granularitate grosieră la una fină, există o limită dincolo de care ne putem aștepta la grave contradicții.

Teoria centrelor (*centering*, CT) [Grosz *et al.*, 1995; Brennan *et al.*, 1987] furnizează explicații convingătoare asupra restricțiilor de utilizare a pronumelor pentru realizarea referințelor și asupra ce anume face un discurs să fie coerent. CT nu se aplică însă dincolo de limitele unui segment (văzut în accepțiunea din AST). Avem de a face, așadar, cu o teorie locală asupra coeziunii și coerenței. Deși nu este definită riguros în teorie, în toate exemplele autorilor unitatea elementară a structurii lingvistice este fraza [*utterance*, exprimare]. Abordări ulterioare întrevăd posibilitatea de a considera o segmentare mai fină, la nivel de propoziție (v. [Kameyama, 1998] de exemplu). Noi vom considera drept **unitate** a structurii de discurs același tip de întindere lexicală ca și în cazul RST, adică acea întindere ce la nivel sintactic este o propoziție iar la nivel semantic - o predicție. Fiecare unitate de discurs u_n ce intră în compoziția unui segment este caracterizată de o listă de **centre anticipatoare** (*forward-looking*) notată $C_f(u_n)$. Centrele listei $C_f(u_n)$ sunt entități semantice ce corespund, la nivelul textului, expresiilor referențiale cuprinse în unitatea u_n . Spunem că o expresie referențială **realizează** un centru. Elementele acestei liste sunt ordonate pentru a reflecta importanța relativă în u_n . Criteriile de ordonare a elementelor listei C_n sunt, în forma originală a teoriei, de natură sintactică, deși alte abordări le diferențiază în funcție de limbă (v. de exemplu [Walker *et al.*, 1994] pentru japoneză, [de Eugenio, 1990; de Eugenio, 1998] pentru italiană, sau [Strube, Hahn, 1996] pentru germană). Pentru limba engleză autorii CT dau următorul criteriu: subiect > obiect-direct > obiect-indirect > complemente > adjuncți. Elementele listei $C^{\wedge}(u_n)$ sunt acele entități despre care se vorbește în unitatea u_n și deci despre care e cel mai probabil că se va continua să se vorbească și în unitatea următoare, u_{n+1} , dacă aceasta aparține aceluiași segment ca și u_n . Cel mai bine plasat element al listei $C(u_n)$ se numește **centru principal** și se notează $C_p(u_n)$. Fiecărei unități îi este asociat un unic **centru retroactiv** (*backward-looking*), notat $C_r(u_n)$. Prin convenție, centrul retroactiv al

primei unități a segmentului este considerat centrul principal, în timp ce, pentru toate celelalte unități ale segmentului, el este cel mai bine plasat element al listei $C(u_n)$ a unității precedente care este de asemenea realizat și în unitatea curentă.

Teoria face o clasificare a tranzițiilor posibile între unități consecutive, din punctul de vedere al invariantei ori nu a centrelor retroactive și al identificării ori nu a lor cu centrele principale. Astfel, cu excepția cazului în care între unități succesive ale aceluiași segment nu există centre comune, următoarele patru tipuri de tranziții sunt posibile:

CONTINUARE (*continuing*, CON): $C_o(u_{n+1}) = C_o(u_n)$ și $C_p(u_{n+1}) = C_p(u_n)$, corespunzând situației în care atât în u_n cât și în u_{n+1} se vorbește despre aceeași entitate și este de așteptat ca și în unitatea următoare să se vorbească despre ea.

REȚINERE (*retaining*, RET): $C_o(u_{n+1}) = C_o(u_n)$ dar $C_p(u_{n+1}) \neq C_p(u_n)$, cărui interpretare este că, deși atât în u_n cât și în u_{n+1} se vorbește despre aceeași entitate, este de așteptat ca în unitatea următoare să se vorbească despre o alta.

SCHIMBARE LINĂ (*smooth-shifting*, SSH): $C_o(u_{n+1}) \neq C_o(u_n)$ dar $C_p(u_{n+1}) = C_p(u_n)$, cu semnificația că deși în u_n și în u_{n+1} se vorbește despre aceeași entitate este de așteptat ca în unitatea următoare să se vorbească despre entitatea menționată ultima oară.

SCHIMBARE ABRUPTĂ (*abrupt-shifting*, ASH): $C_o(u_{n+1}) \neq C_o(u_n)$ și $C_p(u_{n+1}) \neq C_p(u_n)$, cu semnificația că în u_n și în u_{n+1} nu se vorbește despre aceeași entitate și este de așteptat ca în unitatea următoare să se vorbească despre o altă entitate decât cea menționată.

Nucleul CT este concentrat în două reguli, prima enunțând o constrângere asupra formei de realizare a centrelor prin pronume, iar cea de a doua formulând preferințe asupra secvențelor de tranziții ale centrelor. Regula a doua, cea care referă la coerență, formulează presupunerea că anumite secvențe produc încărcare inferențială în ascultător mai mare decât altele:

Regula 2: Secvențele de continuări sunt preferabile secvențelor de rețineri, care sunt preferabile secvențelor de schimbări line, iar acestea sunt preferabile secvențelor de schimbări bruște: CON > RET > SSH > ASH.

Dacă ne abținem de a penaliza CT, ca teorie locală, așadar aplicabilă la întinderea unui segment, pe motivul fragilității noțiunii de segment, care are definiție recursivă (un segment este constituit din alte segmente), slăbiciune moștenită de la AST, atunci apare naturală tentativa de a lărgi aplicabilitatea CT la întregul discurs, într-o manieră recursivă, pe chiar această structură de segmente definită, ea însăși, recursiv. Teoria nervurilor propune o astfel de generalizare.

Teoria nervurilor (*veins theory*, VT) [Cristea *et al.*, 1998], preluând de la RST diferențierea dată de nuclearitate între argumentele relațiilor retorice dar ignorând, ca și în AST, numele acestora, relevă o structură "ascunsă" în arborele de discurs, numită **nervură**. Fără a nega structura lingvistică a segmentelor de discurs, cât și pe cea intențională a relațiilor dintre scopurile comunicate de segmente și care, prin echivalarea de care am amintit ([Moser, Moore, 1996; Marcu, 1999]), poate fi recuperată din structura de arbore proprie analizelor RST, VT corectează defectul de accesibilitate al AST înlocuind modelul accesibilității în stivă cu accesibilitatea de-a lungul nervurilor arborelui de discurs și explicând naturațea unor referințe la distanță realizate prin mijloace de evocare foarte economice (pronume) [Fox, 1987]. Concluziile VT sunt, de asemenea, stabile la granularitate. În felul acesta VT se constituie într-o teorie globală a coeziunii discursului. VT generalizează totodată partea din CT relativă la încărcarea inferențială (regula a doua), extinzând concluziile ei la întregul discurs, prin aceasta VT constituindu-se și într-o teorie globală a coerenței.

În secțiunea următoare sunt prezentate argumente lingvistice în favoarea teoriei. Secțiunea 3 prezintă definițiile teoriei, secțiunea 4 enunță conjectura VT relativă la referențialitate, iar secțiunea 5 - conjectura VT referitoare la coerență. Secțiunea 6 descrie rezultate experimentale în sprijinul presupuzițiilor VT, secțiunea 7 prezintă o proprietate de granularitate, iar ultima secțiune este dedicată concluziilor și prezentării unor aplicații ale VT.

2. Intuițiile VT

Noțiunea de nervură s-a născut sintetizând observațiile asupra modului în care se aliniază referințele pe o reprezentare arborescentă a discursului. Considerând organizarea ierarhică dată de structura de arbore și principiul compoziționalității (v. de exemplu [Marcu, 2000]), care permite ca unități de discurs aflate la distanță să fie frați sub aceeași relație, aceste observații au fost următoarele (pentru simplificarea exprimării vom spune că "o unitate A referă o unitate B" și vom înțelege "o expresie referențială aparținând unei unități A referă o entitate de discurs introdusă de (sau referită dintr-o) unitate B"; de asemenea vom nota cu u_1, u_2, u_3 - unități de discurs iar cu R, R_1, R_2 - relații. Atunci când apar ca argumente ale unei relații, unitățile de discurs vor purta un indice ridicat "sau s ", cu semnificația de nucleu și respectiv satelit):

- un satelit sau un nucleu poate referi un frate nuclear aflat la stânga: în combinații $u'' R u_s$, sau $u'' R u_s, u_s$ poate referi uu

Ex. 1

1. *Ion a plecat de acasă fără umbrelă*
2. *deși dimineață 0 aflase la radio că va ploua.*

Subiectul vid (notat **0**) din unitatea 2, un satelit al unității 1 [**Ion**] introdusă de expresia referențială *Ion* din prima unitate.

- un nucleu poate referi un satelit al său aflat la stânga: în combinații u_s poate referi u_s . Astfel, în exemplul:

Ex.2

1. *Ion i-a dat Măriei o floare.*
2. *Pentru că 0 s-a simțit frustrată,*
3. *soția lui - s-a supărat.*

unitatea 2 este un satelit al unității 3. Pe cine desemnează pronumele din 2, pe [**Măria**] sau pe [**soția lui Ion**]? Într-o interpretare incrementară la sfârșitul receptării celei de a doua unități avem tendința de a atribui timpuriu, subiectul vid [**Măriei**] apreciind totodată bizarul situație în care unității 3 are loc însă o reconsiderare a legării **0**-> [**Măria**] și expresiei referențiale *soția lui* cu subiectul vid din 2, ambele în combinație cu [**soția lui Ion**].

un satelit dreapta al unui nucleu u nu e accesibil dintr-un nucleu sau satelit, al lui u : în combinații $\{u'' u_s\} R_2 u_s''$ sau $u_s'' u_s$, u_s poate referi u -l dar nu u_s .

Ex.3

1. *Ion i-a mărturisit Măriei că o iubește.*
2. *El n-&fost niciodată căsătorit*
3. *și a trăit până la 40 de ani lângă mama sa.*
4. *Ea, dimpotrivă, a fost măritată de două ori.*

Secvența 2-3-4 oferă o completare la 1. Secvența 2-3 se referă la o relație de CONTRAST (o relație paratactică) față de 4, iar 3 aduce o completare la 2. Structura este deci următoarea: $u'' ((u_s R_2 u_s^s) R_2 u_s^s)$ în care u_s^s este un nucleu CONTRAST. Pentru cei mai mulți cititori, ea din unitatea 4 trebuie să se refere la nu [**mama lui Ion**], deși [**mama lui Ion**] este entitatea cea mai recentă în poziția unității 4, cu care pronumele feminin se potrivește în număr și în preferării Măriei în locul mamei este acela că cititorul recunoaște că într-o relație de CONTRAST cu unitatea 2 (relație pusă în evidență de *dimpotrivă*), ceea ce face ca cele două unități să fie percepute ca fiind în relație. Apropierea lor nu este însă una liniară, ci ierarhică, pe structură închisă la referință din unitatea 4.

¹ Vom nota prin [text] entitatea de discurs introdusă/referită de expresia [text]

un nucleu blochează accesibilitatea dintr-un satelit dreapta spre un satelit stânga: în combinații ($uf R, u_2$)ⁿ $R_2 u_3^s, u_3$, poate referi u_2 dar nu u_1 .

Ex.4

1. *Încă înainte cu un an de terminarea mandatului său de președinte al firmei*
2. *dl. W. Ross începuse mașinațiile pentru falimentarea acesteia.*
- *3. *De altfel, circulau vorbe că l-ar fi obținut fraudulos.*

În acest exemplu 1 și 3 sunt sateliți ai lui 2 (1 este o circumstanțială a lui 2, în timp ce unitatea 3 dă o explicație la purtarea necinstită a lui Ross). Referința **!=[mandatul lui Ross de președinte al firmei]** se deduce cu dificultate, ceea ce face ca întregul discurs să fie defectuos. Dimpotrivă, în următoarea variantă, discursul câștigă în cursivitate:

Ex.5

1. *dl. W. Ross începuse mașinațiile pentru falimentarea firmei al cărei președinte era*
2. *încă înainte cu un an de terminarea mandatului său.*
3. *De altfel, circulau vorbe că l-ar fi obținut fraudulos.*

În Ex. 5 unitatea 2 este un satelit al lui 1, iar 3 - un satelit al lui 2 (aici *de altfel* anunță o paranteză la informația asupra mandatului de președinte). Referința **!=[mandatul lui Ross de președinte al firmei]** poate fi recuperată acum fără dificultate.

Motivația acceptării Ex. 5 și rejectării Ex. 4, constă nu în depărtarea liniară mai mare a anaforului de antecedent în Ex. 4 decât în Ex. 5, ci în faptul că în Ex. 4, spre deosebire de Ex. 5, accesul anafor-antecedent se face dinspre un satelit către un alt satelit, între ei interpunându-se un nucleu. Să remarcăm, de asemenea, că Ex. 4 poate fi reparat și dacă se elimină această referință:

Ex. 6

1. *încă înainte cu un an de terminarea mandatului său de președinte al firmei*
2. *dl. W. Ross începuse mașinațiile pentru falimentarea acesteia.*
3. *De altfel, circulau vorbe că el ar fi fraudat alegerile.*

3. Definițiile teoriei

Intuiția fundamentală care stă la baza dezvoltărilor unificatoare asupra structurii de discurs și accesibilității în VT este că distincția specifică RST dintre

nuclee și sateliți constrânge plaja de antecedenti asupra căreia anaforii¹; cu alte cuvinte, distincția nucleu-satelit, corelată cu o structură de discurs, induce pentru fiecare unitate de discurs un domeniu de accesibilitate imediată pentru anaforii pe care-i conține. Mai precis, VT asigură accesibilitate pentru fiecare anafor x aparținând unei unități de discurs u , x pe baza ușurință examinând doar un subset al mulțimii entităților de discurs. Dacă antecedentul lui x este plasat într-o unitate de discurs din domeniul lui u atunci legătura anafor-antecedent este refăcută pentru realizarea ei e nevoie de mijloace referențiale tari, cum ar fi numele proprii.

Mai mult decât atât, aceeași corelație nuclearitate-structură de discurs, permite generalizarea CT dincolo de granițele unității de discurs ceea ce face posibilă aplicarea concluziilor CT asupra coerenței discursului.

VT se bazează, în mare măsură, pe aceleași elemente de analiză a discursului ca și RST:

- unitățile de bază ale discursului sunt întinderi de text (sau *span*) ce nu se intersectează. După cum am precizat în secțiunea anterioară vom asimila cu propoziții, la nivel semantic fiecare întindere de text o predicție (căreia îi corespunde o reprezentare structurală situațională);
- structura unui discurs este reprezentată ca un arbore de RST, dar fără a reduce generalitatea, în VT vorbim de discurs ca fiind binari (fiecare nod al arborelui are două antecedenti) (pentru argumentație, v. [Marcu, 2000] și [Orwell, 1984]);
- principiul secvențialității [Cristea, Webber, 1997]: o întindere de pe frontiera terminală a arborelui corespunde unei unități de discurs ce compune textul²;
- principiul compoziționalității [Marcu, 2000]: o relație de RST între două întinderi de text se aplică, de asemenea, și la nuclele nucleare ale întinderilor aflate în relație;
- la fel ca în RST, nuclearitatea nodurilor arborelui de RST este clasificată în nodurile fiind clasificate în nuclee (cele mai importante și mai puțin importante);

¹ Într-o relație anaforică, interpretarea anaforului depinde de contextul antecedentului fiind plasat în text înaintea anaforului.

² Unitățile de discurs întrerupte nuanțează acest principiu. Astfel în următorul: *O datări când treceau unul pe lângă altul pe coridor?! e piezișă¹ care parcă-l străpunsese!² și pentru o clipă fusese cuprinsă* (G. Orwell, 1984), unitatea 1 este întreruptă de unitatea 2.

- nodurile terminale ale arborelui reprezintă unități de discurs, în timp ce nodurile neterminale reprezintă relații retorice între întinderi adiacente de text. Spre deosebire de RST, în VT nu interesează numele relațiilor, ceea ce contează fiind topologia arborelui, nuclearitatea nodurilor și etichetarea nodurilor terminale;
- între fiii fiecărui nod intermediar al arborelui există cel puțin un nod nuclear. Nodul rădăcină, prin convenție, e considerat satelit.

În vizualizarea arborilor vom reprezenta nodurile neterminale prin dreptunghiuri fără nume, pe cele terminale - prin ovaluri etichetate, iar nodurile nucleare vor fi subliniate (v. Figura 1). În definițiile ce urmează vom folosi următoarele convenții de notare:

- $mark(a)$ este o funcție care întoarce șirul a în care fiecare simbol este marcat (de exemplu, este poziționat între paranteze);
- $unmark(d)$ este funcția inversă lui $markQ$, ce îndepărtează toate marcajele atașate simbolurilor din expresia a (ex. $unmark(mark(a)) = a$);
- $simpl(a)$ este funcția care elimină toate simbolurile marcate din expresia argumentului a (ex. $simpl(mark(a)) = o$, șirul vid, și $simpl(cx \cdot mark(p) \cdot y) = a - y$);
- $seq(a, \beta)$ este o funcție de secvențiere, care întoarce acea permutare a concatenării simbolurilor din a și β dată de citirea de la stânga la dreapta a nodurilor corespunzătoare simbolurilor din a și β pe frontiera terminală a arborelui. Funcția menține marcajele asupra simbolurilor, dacă acestea există, $seq(0, fi) = fi$; și $seq(a, seq(P)) = seq(seq(a), \beta) = seq(a, 0)$;
- $H(n)$ și $V(n)$ reprezintă expresiile *head* și *nervură* (în engleză - *vein*) ale unui nod n ;
- $pref(u, a)$ reține prefixul expresiei simbolice a până la simbolul u inclusiv, o etichetă de nod terminal.

Teoria nervurilor calculează două expresii, pe care le atașează fiecărui nod al structurii.

3.1 Expresia *head* a unui nod al arborelui

Intenția expresiei *head* a unui nod al arborelui de discurs este de a exprima secvența celor mai importante unități de discurs din întinderea de text acoperită de nod. Ea este o secvență de etichete de unități, calculată după cum urmează:

Definiții

Expresia *head* a unui nod terminal este însăși eticheta sa.

2. Expresia *head* a unui nod neterminal este dată de cele mai importante unități de discurs în ordine de apariție în arbore de la stânga la dreapta, excluzând *head* ale descendenților săi nucleari.

Definițiile expresiilor *head* sugerează un proces de calcul de jos în sus în arborele de discurs. Cele mai importante unități de discurs sunt proiectate în sus până în primul nod satelit întâlnit.

direcția
a calculului
expresiei

Figura 1: Calculul expresiilor *head*

3.2 Expresia *nervurii* unui nod al arborelui

Expresia *nervurii* unui nod intenționează să surprindă semnificația de discurs care sunt semnificative pentru a sintetiza, în context, întinderea de text acoperită de nod. Pentru orice nod al structurii, *nervuria* este formată din cele mai importante unități de discurs din întinderea de text acoperită de nod, împreună, eventual, cu alte unități din afara acestei întinderi.

Prin sinteza, sau rezumatul, unei întinderi de text se înțelege un text care conține ideea principală a textului supus sintezei. Indiferent dacă este realizat prin punerea cap la cap a unor subsecvențe ale întinderii originale, rezumatul trebuie să fie comprehensibil, adică trebuie să poată fi înțeles (printre altele, de exemplu, rezumatul trebuie să conțină toate elementele necesare pentru rezolvarea anafilor). Adesea însă, atunci când întinderea este decupată dintr-un text mai larg, pentru ca rezumatul să fie comprehensibil, el trebuie să conțină și unități din afara întinderii și care aparțin contextului. Avem de a face, în acest caz, cu o sinteză a unei întinderi de text în contextul unei întinderi mai vaste. Să mai observăm că, în privința, "a sintetiza" e analog cu "a înțelege", pentru că ceea ce ne interesează dintr-un text este o sinteză a lui.

Definițiile care urmează, datorită recursivității lor, vor face posibilă considerarea contextului dat de totalitatea textului din exprimarea "a înțelege, în contextul întregului text, întinderea s" mărginit la întinderea de text acoperită de nodul părinte al celui corespunzător întinderii s. Cu alte cuvinte, la fiecare nivel al structurii, cu excepția rădăcinii, adică întotdeauna unde există două noduri fii sub un nod părinte, cu întinderile celor două noduri fii însumând întinderea nodului părinte, expresia nervură a părintelui conține deja informația care permite înțelegerea/rezumarea întinderii acoperite de el în contextul global. Coborârea pentru înțelegerea/rezumarea subîntinderii acoperite de nodul curent al definiției (unul dintre cele două noduri fii) înseamnă adăugarea și/sau ștergerea unei secvențe noi/subsecvențe la/din secvența de etichete contribuită de nervura părintelui, în funcție de polaritatea și poziția specifică a întinderii corespunzătoare nodului fiu curent în întinderea nodului părinte. În continuare, întinderea întregului text, o constantă pentru orice subîntindere, va fi numită contextul total. În figurile 2-6, nodurile curente - cele vizate de definițiile curente de nervură - apar în gri. Ele sunt notate simultan cu un dreptunghi și un oval pentru a sugera că pot fi atât noduri interioare (neterminale), cât și noduri terminale.

Definiții

1. Expresia nervurii rădăcinii este egală cu expresia sa *head*.

Particularizând intenția pe care o exprimă expresia nervurii unui nod la modul rădăcină obținem: cele mai semnificative unități de discurs necesare înțelegerii/rezumării întinderii acoperite de nod (în cazul de față - întregul text) în contextul total. Cum contextul este aici egal cu textul în totalitatea lui, el poate fi lăsat la o parte în descriere, ceea ce ne lasă cu definiția expresiei *head* a nodului rădăcină.

2. Pentru fiecare nod nuclear, al cărui părinte are nervura v :

a. dacă nodul nu are un frate nenuclear în stânga, atunci expresia nervurii este v (v. Figura 2);

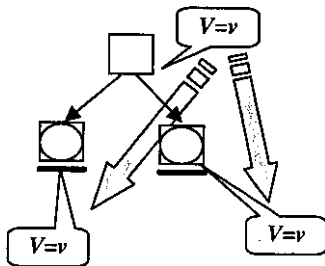


Figura 2: Expresia nervurii unui nod nuclear fără frate satelit în stânga

Definiția exprimă faptul că secvența de unități necesară înțelegerii/rezumării, în contextul total, a unei întinderi nucleare de text ce are ca frate în structură o altă întindere nucleară necesită aceeași secvență de unități ca și ce este necesară înțelegerii/rezumării, în contextul total, a reuniunii celor două întinderi. Cu alte cuvinte, o întindere nucleară ce este frate, în structură, întinderii nucleare curente este esențială înțelegerii/rezumării întinderii curente.

b. dacă nodul are un frate nenuclear în stânga de *head* h , atunci expresia nervurii lui este $seq(mark\{h\}, v)$ (v. Figura 3);

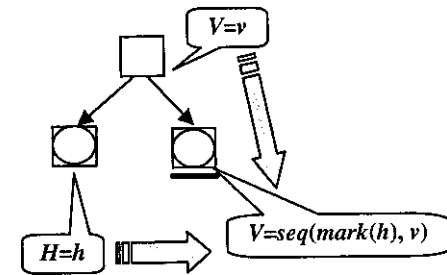


Figura 3: Expresia nervurii unui nod nuclear având un frate satelit în stânga

Secvența de unități necesară înțelegerii/rezumării, în contextul total, a unei întinderi nucleare de text ce are ca frate stânga în structură o întindere nenucleară necesită, suplimentar față de secvența necesară înțelegerii în contextul total întinderii acoperită de nodul părinte (comunicată de expresia nervură a nodului părinte) și secvența *head* a întinderii frate stângi (adică cele mai importante unități din întinderea stângă). Considerarea, în expresia nervurii întinderii nucleare curente, a expresiei *head* a întinderii nenucleare frate stânga, corespunde, prin prisma definiției 2a, cu atribuirea întinderii stângi a calității de a se comporta ca nucleu. Marcarea contribuției satelitului frate stânga prin funcția $markQ$ face în această revizuire a nuclearității lui, una cu valoare temporară, după cum se dovedește mai jos, în definiția 3b.

3. Pentru fiecare nod nenuclear de *head* h , al cărui părinte are nervura v :

a. dacă nodul este descendentul stâng al părintelui său, atunci expresia nervurii sale este $seq(h, v)$

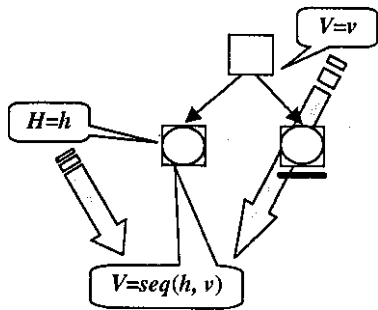


Figura 4: Expresia nervurii unui nod satelit stânga

Definiția exprimă faptul că pentru a înțelege/rezuma, în contextul total, o întindere nenucleară de text ce este descendent stâng, în structură, nodului părinte, la secvența de unități ce exprimă influența contextului total (precizată de expresia nervură a părintelui) trebuie adăugate cele mai importante unități din întinderea proprie (date de expresia *head* proprie). Să observăm că în expresia nervurii nodului părinte, care moștenește expresii *head* ale nodurilor superioare, nu poate răzbate influența unui fiu satelit al său, deci numai includerea *head*-ului fiului satelit, direct în expresia nervurii sale poate completa această influență.

b. dacă nodul este descendentul drept al părintelui său, atunci expresia nervurii lui este $seq(h, simpl(v))$.

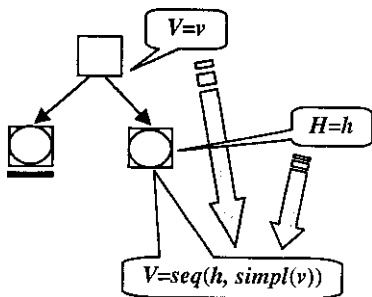


Figura 5: Expresia nervurii unui nod satelit drept

Pentru a înțelege, în contextul total, o întindere nenucleară de text ce este descendent pe dreapta al nodului părinte, la secvența de unități necesară înțelegerii/rezumării contextului total (precizată de expresia nervură a părintelui) și

din care s-au șters unitățile marcate trebuie adăugate cele mai importante unități din întinderea proprie (date de expresia *head* proprie). În acest fel, dacă expresia nervură a nodului părinte nu conține unități marcate (în conformitate cu definiția 2b), atunci expresia nervură a unui satelit dreapta nu diferă de expresia nervură a acelui satelit ce ar fi fost poziționat pe stânga (conform definiției 3a). Dacă în expresia nervură a părintelui conține unități marcate, atunci acestea dispar din expresia nervurii satelitului drept. Conform definiției 2b, unitățile marcate pot fi datorate doar unui satelit stânga, frate al celui mai apropiat ascendent nuclear al întinderii curente. Urmează că definiția curentă exprimă o proprietate de blocare a accesibilității dinspre un satelit plasat în dreapta unui nucleu către un satelit plasat în stânga sa (v. Figura 6).

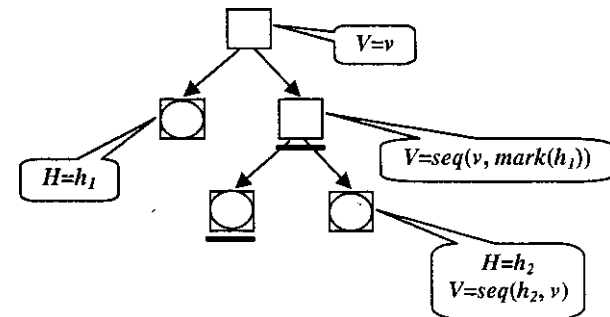


Figura 6: Simplificări în calculul expresiei nervură a unui satelit dreapta:

$$V=seq(h_2, simpl(seq(v, mark(h_1))) = seq(h_2, seq(v)) = seq(h_2, v)$$

Dacă semnificația expresiei nervurii unui nod oarecare din structură este particularizată la un nod terminal, obținem: **expresia nervurii unei unități de discurs reprezintă secvența unităților de discurs care sunt semnificative pentru a înțelege/rezuma, în contextul întregului text, însăși unitatea de discurs în cauză.** Printre altele, aceasta înseamnă că expresia nervurii unei unități de discurs este suficientă pentru a interpreta toate referințele anaforice conținute în unitate.

4. Relația dintre structura de discurs și referențialitate

Ipoteza pe care o avansăm este că rezoluția anaforică este caracterizată de două tipuri de procese: **evocative** (sau **imEDIATE**) și **post-evocative** (sau **interferențiale**). Procesele evocative, cele mai frecvente, sunt rapide și pot fi realizate prin orice mijloace de evocare referențială, inclusiv cele fragile (de tipul subiecte

vide și pronomelor). Ele dau textului fluentă și-l fac coeziv. Cele post-evocative sunt mult mai puțin frecvente decât cele evocative, necesită o încărcare inferențială mai mare pentru a fi interpretate și utilizează mijloace referențiale tari (nume proprii, substantive comune articulate).

Vom asocia spațiul de căutare al proceselor evocative unui **domeniu de accesibilitate referențială evocativă** sau **imediată** (*domain of evocative accessibility - dea*) pe baza definiției nervurii și a următoarelor observații:

- **relația anaforică este de natură semantică, iar nu textuală** [Halliday, Hassan, 1976]: o relație anaforică are doi termeni: anaforul și antecedentul. Anaforul este reprezentat de o expresie referențială a cărei natură este textuală. Natura semantică a relației anaforice trebuie înțeleasă ca răsfrângându-se asupra antecedentului, care nu trebuie identificat cu o anumită expresie referențială ce precede în text anaforul, ci cu o reprezentare a acesteia într-un plan semantic, în așa fel încât semnificația anaforului se construiește din antecedentul însuși iar nu din semnificația lui. În cazul particular al unui lanț co-referențial, acest lucru înseamnă că antecedentul este "realizat" repetat în text în aceeași entitate de discurs. Expresiile co-referențiale "ancorează", în diverse poziții ale textului, o aceeași entitate de discurs.
- **dinamica interpretării discursului este incrementală**: un discurs este un text în procesul citirii ori ascultării lui de către un subiect (om sau mașină). Când citirea/ascultarea unui text s-a terminat, discursul este încheiat și ceea ce rămâne este o reprezentare a lui în memoria subiectului. De asemenea, la un moment dat pe parcursul interpretării unui text, anumite elemente ale discursului pot fi plasate privilegiat în sfera atenției [Grosz, Sidner, 1986; Sidner, 1983; Walker, 1996], iar trecerea de la o unitate de discurs la următoarea poate produce schimbări în structura memorată ce configurează sfera atenționată.
- **anafora și a catafora au o natură cognitivă comună**: din punct de vedere cognitiv, toate referințele anaforice se fac dinspre expresii referențiale (entități textuale) către entități ale discursului (entități semantice) deja introduse de discursul trecut. Acest lucru înseamnă că, într-o limbă în care textul se notează de la stânga spre dreapta, nu există jeferințe anaforice spre dreapta. Distincția dintre anafora și cataforă, devine, în această viziune care încearcă să reconstituie procesele cognitive ce stau la baza înțelegerii textelor (cu sau fără scopul simulării lor pe mașină), inutilă. În aceeași manieră în care, în cazul unei anafore, un antecedent este o entitate de discurs propusă de o expresie referențială ce precede anaforul și pe care anaforul o referă apoi, pronumele ce precede un nume, în cazul unei catafore, propune o reprezentare, mai săracă, pe care numele o referă și o

completează în același timp [Cristea, Dima, 2001]. Această relație este interpretării discursului o unică direcționalitate, care este cea a desfășurării liniare timpului lecturii, și care este cea a desfășurării liniare a limbii europene, de exemplu, de la stânga la dreapta. Această referențialitate trebuie deci să se proiecteze pe aceeași direcție anafori "noi" către entități "vechi", mereu către înapoi în timpul lecturii.

Ex.7

1. Pentru că O n-a vrut să-și lase tata singur,
2. Ion a renunțat la concediu.

Expresia referențială vidă de pe poziția de subiect a unității de discurs propune o entitate de discurs caracterizată cel mult de o descriere referențială [number singular] (ce poate fi atribuită, cel mai probabil, unor surse de natură pragmatică: cineva care nu poate să-și lase tatăl singur și o persoană, corelate cu surse de natură sintactică: acordul în numărul singular). Apoi, substantivul propriu /on, din unitatea 2, referă entitatea corelată și o completează până la o reprezentare: [type human, number singular, name Ion].

Corelarea definiției nervurii cu observațiile de mai sus configurează domeniul de accesibilitate referențială evocativă a unității de discurs, toate unitățile de discurs care preced unitatea în care se referă referențială (și din care au fost îndepărtate eventualele marcaje, în rol de memorie temporară):

$$dea(u) = pref(u, unmark(V(u))).$$

Definiția *dea* formalizează prima conjectură a VT (sau a VT2) și pune în legătură accesibilitatea referențială imediată cu structura de discurs **antecedentii expresiilor referențiale dintr-o unitate de discurs care precede, printre entitățile de discurs ancorate în unitățile de discurs inclusiv u, în expresia nervurii acesteia.**

Paul Cornea [1998] vorbește despre recodificarea semnificativă a memoriei. El pune în evidență trei tipuri de memorie, ce apar, de la mai puțin la mai mult, la cercetători [Kinntsch, Vârî Dijk, 1975; Schank, Abelson, 1975]: memoria imediată, memoria de scurtă durată (de termen scurt - SLT) și memoria de lungă durată (de termen lung - MLT). Memoria imediată este un sistem senzorial al informațiilor, reținerea urmelor din ultima jumătate de secundă conservă câteva secunde informația. Lungimea acestei memorii este de câteva semne (cuvinte, cifre, litere - funcție de context, v. și [Miller, 1956]).

apreciază acest "empan" mijlociu la 13[^]-15 cuvinte, la un lector lent fiind de 8 cuvinte, la unul rapid - de 16-5-20 [Richadeau, 1969] - citat în [Cornea, 1998] p. 166).

Construcția structurii de discurs se face dinamic, în actul lecturii. Să ignorăm un posibil proces de multi-interpretare ce poate duce la sintetizarea simultană a mai multor construcții alternative, din care să se selecteze, în urma unui proces de dezambiguizare, una sau mai multe structuri arborescente finale. Arborele însuși poate fi considerat rezumat în diverse grade, conform capacității de memorare a subiectului. Dacă unitatea curentă este u_n , să notăm AR_n arborele de structură rezumat, la momentul prelucrării unității u_n . Nervura acesteia, culeasă pe AR_n , este $V(u_n)$, iar domeniul ei de accesibilitate imediată $dea(u_n)$. Noi credem că MST poate fi considerată o fereastră de lungime 7 ± 2 semne în directă legătură cu $dea(u_n)$: fie 7 ± 2 unități din această secvență, fie tot atâtea structuri evenimențiale - ca reprezentări ale unităților de discurs, fie încă numai simboluri (cuvinte etc.) culese din acest șir de unități. Tranzitarea la următoarea unitate, u_{n+1} , înseamnă înlocuirea memoriei de scurtă durată $dea(u_n)$ cu $dea(u_{n+1})$. Acest lucru duce uneori la o simplă prelungire a domeniului de accesibilitate precedent, alteori la o alterare a lui prin ștergerea unor unități și adăugarea altora, de fiecare dată domeniul încheindu-se cu unitatea curentă. MST este așadar o proiecție a unui șir de unități de discurs (sau de microstructuri ce-și au suportul în aceste unități) decupate din structura dinamică curentă. Modificările ce apar în șirul MST reflectă schimbările de focalizare, în parcurgerea discursului. Componenta acestui șir este influențată de uitare (deci de un proces de abstractizare) și de modificarea de interes curentă în parcurgerea discursului. Când interesul s-a mutat pe o altă axă, componenta nervurii și, de aici, a domeniului de accesibilitate imediată sunt și ele actualizate. Incluziunea sau excluderea din MST a unor unități de discurs în ritmul citirii, pentru că dea evoluează eliminând unele unități și "redeșteptând" altele "uite", amintesc de procesele de "chemare" în sfera atenției ale memoriei *cash* a lui Walker [Walker, 1996]. Pe de altă parte, structura memorată (rezumată) a discursului este păstrată în MLT și folosită pentru aducerea în prim plan a unităților de interes curent ce au fost temporar retrogradate de o comutare a atenției într-o altă direcție. Procesele evocative se desfășoară așadar în memoria de scurtă durată. Pe de altă parte, procesele post-evocative sunt procese de rezoluție anaforică de natură inferențială, ce presupun un anumit efort de regăsire a unei entități de discurs într-o zonă a memoriei de lungă durată sau evocă entități ale cunoașterii generice din sfera culturală a subiectului. Noi credem că aceste procese se dezvoltă tot pe structura de discurs dezvoltată deja, ieșind din dea , când rezoluția a eșuat acolo.

Dintr-un punct de vedere ce se concentrează asupra relației dintre referențialitate și structura de discurs, celor două tipuri de procese anaforice pe care le-am pus în evidență le corespund **referințe evocative**, respectiv **post-evocative** (sau **inferențiale**). Diferența dintre ele este că, în cazul primelor, lanțul retroactiv al unităților ce ancorează expresii aflate în relații referențiale

intersectează domeniul de accesibilitate referențială imediată al unității anaforului în cel puțin încă un punct decât unitatea anaforului, pe când în cazul referințelor post-evocative nu există această intersecție dublă. În [Cristea *et al.*, 2000; Cristea 2000] referințele evocative sunt, mai departe, detaliate în **directe** și **indirecte**.

În referințele directe, a doua unitate de intersecție este unitatea cea mai recentă liniar ce ancorează aceeași entitate de discurs ca și anaforul (în cazul relației de co-referință) sau o entitate corelată funcțional cu aceasta (în cazul unei relații de referință funcțională). În referințele indirecte intersecția dea cu lanțul co/func-referențial se realizează într-o unitate mai depărtată decât cea mai recentă liniar de unitatea anaforului. În referințele inferențiale lanțul retroactiv al legăturilor anaforice ale anaforului nu intersectează dea (în Figura 7 lanțul legăturilor anaforice este reprezentat punctat, iar dea printr-o linie grosă).

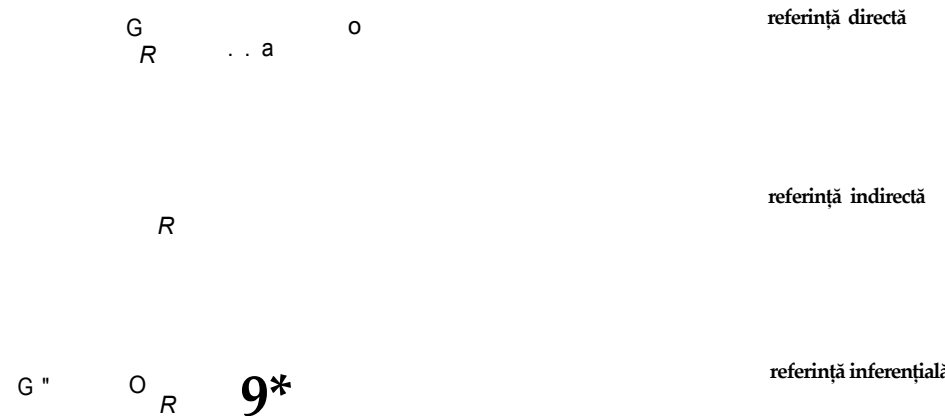


Figura 7: Referințe evocative și post-evocative

O categorie particulară de referințe post-evocative sunt **referințele pragmatice** (ce pot fi numite și **pseudo-referințe**). În acest tip de referințe participă expresii referențiale care pot fi interpretate fără un antecedent pentru că interpretarea lor se bazează pe cunoștințe exterioare textului, ce vin din cunoașterea comună asupra lumii, deci din pragmatică. Deși există cel puțin încă o expresie referențială în text ce realizează aceeași entitate de discurs, expresiile referențiale pot să nu aibă, în mod necesar, o reprezentare unică, fără ca prin aceasta înțelegerea textului să sufere.

Recunoașterea antecedentului se datorează, în toate cazurile, unor procese de confruntare de șabloane (*pattern-matching*) îmbogățite cu euristici, în care intervin structura de caracteristici morfo/sintactico/semantice ce definesc anaforul și structurile de caracteristici ce definesc entitățile de discurs deja introduse [Cristea, Dima, 2001; Cristea *et al.*, 2002a].

5. Relația dintre structura de discurs și cursivitate

5.1. Linii de argumentație

Expresiile nervură ale unităților ce compun un discurs arată tot atâtea moduri diferite în care poate fi citit acel discurs. Fiecare în parte dă o rezumare a discursului prin prisma unității de discurs curente. Atunci când interesul este orientat către un anumit episod al povestirii, putem sări peste pasaje întregi pentru a ne concentra asupra manierei în care elementul de interes se leagă cu ansamblul discursului. În același fel, putem avea în vedere o altă pistă și atunci lectura focalizează un alt fir de interes. Acest nou fir poate să aibă elemente în comun cu primul dar poate, de asemenea, să încorporeze și altele noi. Fiecare fir în parte poate pune în evidență anumite particularități, legate însă strâns de linia principală a discursului. Toate aceste sub-discursuri sunt coerente și, în general, nu există referințe anaforice pentru a căror interpretare să avem nevoie de fragmente aflate în afara rezumatului însuși. Acest lucru înseamnă că traseele referențiale ale rezumatului conțin suficiente elemente care să ducă la recuperarea înțelesului anaforilor.

Să luăm următorul text:

Ex. 8

1. Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,
2. cînd va trece prin munte,
3. și să-i răpună viața.
4. Hera-l ura pe fiul cel nou născut al Letei,
5. pentru că soțul său, prea puternicul Zeus, ținea mai mult la dînsul decît la fiii ei: Hefaistos și Ares.
6. Cînd a ajuns Apolo în muntele Parnas,
7. dihania uriașă s-a avîntat spre dînsul,
8. dornică să-l ucidă.
9. Dar zeul și-a întins arcul.
10. A tras prima săgeată.
11. Erau doar patru zile de cînd văzuse lumea,
12. și întia lui săgeată a și nimerit monstrul.

Alexandru Mitru - *Legendele Olimpului*, Editura Tineretului, 1966

Structura de discurs a acestui text este cea din Figura 8. Tabela 1 dă expresiile nervură și domeniile de referențialitate evocativă ale nodurilor terminale, în coloana $dea(u)$ au fost, totodată, marcate în aldine domenii de referențialitate imediată maximale vis-à-vis de relația de incluziune (cele mai lungi trasee dea). Astfel $dea(1) \subset dea(2) \subset dea(3) \subset dea(4) \subset dea(5) \subset dea(6)$ ș.a.m.d. Vom numi aceste secvențe care întrerup lanțuri de incluziuni linii de argumentație (la), în cazul nostru: 1 2, 1 3 4 5, 1 3 6 7, 1 3 7 8 și 1 3 7 9 10 11 12. Dacă $la(u_i)$ precede imediat $la(u_{i+1})$, atunci în $la(u_i)$ se regăsesc domeniile tuturor unităților dintre u_i+1 și u_{i+1} . În particular, în $la(u_i)$ se regăsesc unitățile ce preced imediat unitatea u_i , pentru orice u_i între u_i+1 și u_{i+1} , în domeniile lor de accesibilitate imediată (adică acel domeniu care conferă discursului maximum de coerență). Cu alte cuvinte, pe $la(u_i)$ putem aplica definițiile CT de calculare a tranzițiilor pentru orice u_i între u_i+1 și u_{i+1} .

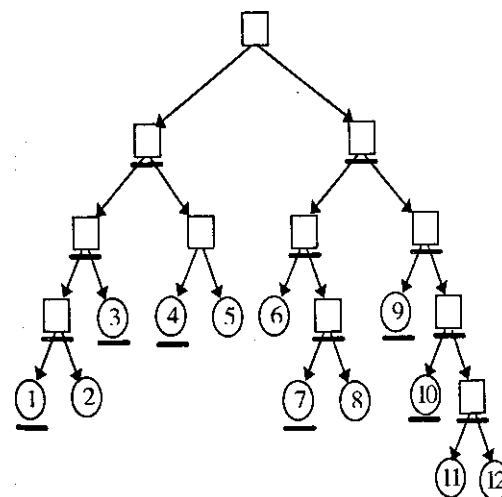


Tabela 1: Nervurile și domeniile unităților din Ex. 8

u	$V(u)$	$dea(u)$
1	1 3 7 9 10 12	1
2	1 2 3 7 9 10 12	1 2
3	1 3 7 9 10 12	1 3
4	1 3 4 7 9 10 12	1 3 4
5	1 3 4 5 7 9 10 12	1 3 4 5
6	1 3 6 7 9 10 12	1 3 6
7	1 3 (6) 7 9 10 12	1 3 6 7
8	1 3 7 8 9 10 12	1 3 7 8
9	1 3 7 9 10 12	1 3 7 9
10	1 3 7 9 10 12	1 3 7 9 10
11	1 3 7 9 10 11 12	1 3 7 9 10 11
12	1 3 7 9 10 (11) 12	1 3 7 9 10 11 12

Figura 8: Structura de discurs a Ex. 8

5.2. O generalizare a CT

Urmând recomandările teoriei centrelor, să presupunem că marcăm tranzițiile ce apar între unități de discurs cu scoruri care să dea un grad al ușurinței de prelucrare:

CONTINUARE	(CON)	4
REȚINERE	(RET)	3
SCHIMBARE LINĂ	(SSH)	2
SCHIMBARE ABRUPTĂ	(ASH)	1
LIPSĂ Cb	(-)	0

În felul acesta, tranzițiile line primesc scoruri mari, cele abrupte, scoruri mici. Însușind aceste scoruri pentru fiecare unitate a unui segment (segment, în spiritul AST) vom avea un scor al segmentului. Să notăm un scor în spiritul CT al unui segment s cu s_{cCT} (CCT de la *Classical Centering Theory*). El ne va da o măsură a ușurinței de interpretare a segmentului: cu cât un segment s , în totalitatea lui, e mai fluent, cu atât scorul lui va fi mai mare și cu cât el este mai abrupt, mai dificil de interpretat, cu atât scorul lui va fi mai scăzut. În fine, să adunăm aceste scoruri pentru toate segmentele discursului, într-un scor al sumei segmentelor S_{cCT} .

$$S_{cCT} = \sum s_{cCT}$$

Să ne imaginăm acuma că forțăm nota și calculăm aceste scoruri și dincolo de granițele de segment, deci inclusiv în punctele de frontieră dintre segmente. Să notăm acest scor global cu s_{cCT} . În scorul global s_{cCT} contribuie cu scoruri de tranziții toate unitățile cuprinse între a doua unitate și ultima. În mod normal tranzițiile în punctele de trecere între segmente ar trebui să fie foarte abrupte, cotate deci slab ori zero, și deci scorul global atașat textului n-ar trebui să fie modificat semnificativ. Dacă apare totuși o diferență, ea trebuie să fie datorată unor tranziții accidentale peste granița de segment. În orice caz trebuie să avem $s_{cCT} \geq s_{cCT}$.

Să procedăm acum în mod analog, ca suport folosind de data aceasta liniile de argumentație iar nu secvențele liniare de unități ale segmentelor în sensul clasic. Datorită comportamentului lor similar segmentelor, putem numi liniile de argumentație **segmente în sens ierarhic**. Să notăm s_{hCT} (HCT de la *Hierarchical Centering Theory*) suma scorurilor unităților aparținând unei linii de argumentație (segment ierarhic) s . Ca să dăm o măsură a fluenței discursului în accepțiunea ierarhică, similară scorului global s_{cCT} în calculul scorului global al discursului în sens ierarhic nu va trebui să repetăm contribuțiile unităților ce apar în mai mult decât o singură linie de argumentație. Dacă notăm s_{hCT} scorul unui segment ierarhic s' în care am păstrat numai unitățile noi față de segmentul anterior, atunci scorul global ierarhic al discursului este:

$$S_{hCT} = \sum s_{hCT}$$

Cea de a doua conjectură a VT (a coerenței): Scorul global în sensul ierarhic al unui discurs este mai bun sau cel puțin egal decât scorul global în sensul clasic: $S_{hCT} \geq S_{cCT}$.

Pentru un anumit detaliu de granularitate în definirea segmentelor în sens clasic, unui segment în sens clasic îi corespunde o secvență de nervură, deci o

portiuone a unei linii de argumentație. În spiritul acestei observații, ce conjectură enunță prezumția că tranzițiile la distanță lungă, calcula / nervurilor, sunt sistematic mai line decât tranzițiile accidentale la gr segmente. Să notăm că această presuposiție este conformă unor obs de autori precum Passonneau [1995] și Walker [1998], furnizând t explicație pentru rezultatele lor.

În cele ce urmează prezentăm o analiză comparativă clasică probează ipoteza coerenței, pe discursul din Ex. 8.

Tabela 2: Analiza Ex. 1 în maniera CCT

n	U_n	$C_f(U_n)$	$C_b(U_n)$
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo, cînd 0 va trece prin munte, și 0 să-i răpună viața.</i>	[Piton], [Hera], [Apolo], $\emptyset =$ [Apolo], [munte] $\emptyset =$ [Piton], KApolo], [viata]	[Piton] [Apolo] [Apolo]
	<i>Hera-l ura pe fiul cel nou născut al Letei, _____ pentru că soțul său, prea puternicul Zeus, ținea mai mult la dînsul decît la fiii ei: Hefaistos și Ares.</i>	[Hera], [Leta], <i>fiul cei nou-născut al Lete</i> HApolo] Zeus], său=[Hera], dmsu]=[Apolo], [Hefaistos], [Ares]	[Apolo] [Hera]
	<i>Cînd a ajuns Apolo în muntele Parnas, dihania uriașă s-a avîntat spre idînsul, \0 (era) dornică să-l ucidă.</i>	[Apolo], [munte] dihania uriașă=[Piton], [d]nsuHApolo] $\emptyset =$ [Pitoni, I=[Apo]o] $\emptyset =$ [Apolo], [săgeata]	Apo [Apo] [Apo] [Apo]
10	<i>0A tras prima săgeată.</i>	$\emptyset =$ [Apolo], [săgeata]	[Apo]
11	<i>Erau doar patru zile de cînd 0 văzuse lumea,</i>	io = [Apolo], [lumea]	[Apo]
12	<i>\și întîia lui săgeată a și nimerit monstrul.</i>	[[săgeata] /u]=[Apolo], monsfriy]=[Piton]	[Apo]

comportă, în medie, intermediar între o schimbare lină (SSH) și o reținere (RET), mai apropiat de o reținere.

Dacă luăm în calcul liniile de argumentație indicate de nervuri, pot fi puse în evidență 5 sub-discursuri, în lungul cărora vom calcula, de asemenea, tranzițiile, în tabelele 3⁷ de mai jos unitățile pentru care considerăm tranzițiile sunt, de asemenea, indicate în caractere aldine în prima coloană. Să remarcăm că citirea textelor date de liniile de argumentație produce, în toate cazurile, discursuri perfect coerente. În ansamblu, doar câte o tranziție este calculată pentru fiecare unitate, la fel ca și în interpretarea clasică.

Tabela 3: Analiza HCT a primei linii de argumentație, secvența de unități 1-2

<i>n</i>	<i>u_n</i>	<i>O(u_n)</i>	<i>C_n(u_n)</i>	Traz.	Scor
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
2	<i>cînd 0 va trece prin munte,</i>	o = [Apolo], [munte]	[Apolo]	SSH	2
Total					2

Tabela 4: Analiza HCT a celei de a doua linii de argumentație, secvența de unități 1-3-4-5

<i>n</i>	<i>u_n</i>	<i>O(u_n)</i>	<i>C_n(u_n)</i>	Traz.	Scor
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
3	<i>și 0 să-i răpună viața.</i>	o = [Piton], !=[Apolo], [viața]	[Piton]	CON	4
4	<i>Hera-l ura pe fiul cel nou născut al Letei,</i>	[Hera], [Leta], fiul cel nou-născut al Letei=[Apolo]	[Apolo]	ASH	1
5	<i>pentru că soțul său, prea puternicul Zeus, ținea mai mult la dînsul decît la fiii ei: Hefaistos și Ares.</i>	[Zeus], său=[Hera], d/?si!= [Apolo], [Hefaistos], [Ares]	[Hera]	ASH	1
Total					6

Se constată că tranziția RET a unității 3 către 2 din analiza CCT s-a transformat într-o tranziție CON, pe nervură, dinspre 3 către 1 deși tranziția RET dinspre 3 spre 4 în CCT devine aici o tranziție ASH, deci mai abruptă, datorită modificării C_b-ului unității 3 din [Apolo] în [Piton].

Tabela 5- Analiza HCT a celei de a treia linii de argumentație, secvența de unități 1-3-6-7

<i>n</i>	<i>u_n</i>	<i>C_n(u_n)</i>	<i>C_n(u_n)</i>
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]
3	<i>și 0 să-i răpună viața.</i>	o = [Piton], HApolo], [viața]	[Piton]
6	<i>Cînd a ajuns Apolo în muntele Parnas,</i>	[Apolo], [munte]	[Apolo]
7	<i>dihania uriașă s-a avîntat spre dînsul,</i>	dihania ur/ășă-[Piton], dfnసు/=[Apolo]	[Apolo]

Tabela 6: Analiza HCT a celei de a patra linii de argumentație, secvența de unități 1-3-7-8

<i>n</i>	<i>u_n</i>	<i>C_n(u_n)</i>	<i>C_n(u_n)</i>
1	<i>Piton primise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]
3	<i>și 0 să-i răpună viața.</i>	o = [Piton], MApolo], [viața]	[Piton]
7	<i>dihania uriașă s-a avîntat spre dînsul,</i>	dihania i/r/ășă-[Piton], dînsul=[Apolo]	[Piton]
8	<i>0 (era) dornică să-l ucidă.</i>	o = [Piton], !=[Apolo]	[Piton]

Se constată că tranziția SSH a unității 7 către 8 din analiza CCT s-a transformat într-o tranziție CON, pe nervură, tot între 7 și 8 (CCT) și a fost schimbat din [Apolo] în [Piton], pentru că, pe nervura lui 8, pe nervura lui 7 este acum 3, iar nu 6 ca în secvența liniară).

Tabela 7: Analiza HCT a ultimei linii de argumentație, secvența de unități 1-3-7-9-10-11-12

<i>n</i>	<i>u</i>		$C_s(u_i)$	Traz.	Scor
1	<i>Piton promise-n taină poruncă de la Hera să-l pîndească pe Apolo,</i>	[Piton], [Hera], [Apolo],	[Piton]	-	-
3	<i>și O să-i răpună viața.</i>	o= [Piton], /= [Apolo], [viața]	[Piton]	-	-
7	<i>dihania uriașă s-a avîntat spre dînsul,</i>	dihania uriașă=[Piton], d/nsil=[Apolo]	[Piton]	-	-
9	<i>Dar zeul și-a întins arcul.</i>	zeul= [Apolo], [arcu]	[Apolo]	SSH	2
10	<i>O A tras prima săgeată.</i>	o = [Apolo], [săgeata]	[Apolo]	CON	4
11	<i>Erau doar patru zile de cînd0 văzuse lumea,</i>	o= [Apolo], [lumea]	[Apolo]	CON	4
12	<i>și întîia lui săgeată a și nimerit monstrul.</i>	[săgeata] /ul=[Apolo], monstruHPiton]	[Apolo]	RET	3
Total					13

Însumând scorurile tranzițiilor pentru toate liniile de argumentație se obține scorul total: 30, ceea ce corespunde unei tranziții medii a discursului, calculată conform HCT de $30/11=2,72$, așadar o tranziție medie mai bună decât cea calculată conform CCT.

6. Validarea conjecturilor VT

Validarea conjecturilor VT s-a realizat pe corpusuri adnotate la structură și la legături co-referențiale. Astfel în [Cristea *et al.*, 1998] se raportează o investigație efectuată pe texte în limbile engleză, franceză și română ce au însumat un total de 176 de unități de discurs. Plecând de o adnotare în maniera RST a structurii de discurs, un program a calculat expresiile nervurilor unităților. Pentru verificarea conjecturii coeziunii, utilizând adnotarea legăturilor referențiale s-a calculat apoi procentajul referințelor directe, indirecte și pragmatice. În medie 99,1% dintre referințe se încadrează acestor trei categorii (87,1% directe, 8,5% indirecte și 3,5% pragmatice). Pentru verificarea conjecturii coerenței, suplimentar marcajelor de structură și lanțuri co-referențiale s-au marcat manual, pentru fiecare unitate, C_s -ul, în varianta clasică și în varianta ierarhică, și s-au calculat tranzițiile în cele două variante. Scorul S_{HCT} a fost mai bun decât scorul S_{CCT} În toate cazurile (scorurile medii pe tranziție au fost de 2,03 în varianta ierarhică față de 1,89 în cea clasică).

În [Cristea *et al.*, 2000] se raportează experimente care au urmărit să compare potențialul modelelor ierarhice, precum cele bazate pe VT, de a regăsi un antecedent într-o plajă de căutare dată față de modelele lineare (modele ce presupun o parcurgere lineară a textului dinspre unitatea anaforului spre începutul

textului). Pentru aceasta s-au utilizat 30 de texte englezești (însușind aproximativ 1560 de unități de discurs), adnotate la structura RST și lanțuri co-referențiale. Presupunând o plajă de căutare de doar 2 unități, căutarea pe nervură a adus cu aproximativ 16% mai mulți antecedenti decât căutarea liniară. După cum era de așteptat, pe măsură ce lungimea textului căutat crește, cele două tipuri de modele se apropie în ceea ce privește potențialul de a regăsi legături co-referențiale. O căutare ierarhică înapoi într-o plajă de 5 unități rezolvă potențial doar 70% dintre anafore, pentru ca o performanță potențială de 90% să poată fi atinsă doar dacă se organizează o căutare într-o lungime de 12 unități pe nervură. O altă investigație a urmărit compararea efortului necesar regăsirii unui anumit antecedent în cele două tipuri de abordări (liniară și ierarhică), unde prin efortul necesar găsirii unui antecedent se înțelege numărul de unități de discurs ce separă, în domeniu, unitatea anaforului de unitatea celei mai recente ancorări în text a unui antecedent. Din nou, modelele ierarhice, de tipul celui dat de VT, s-au dovedit superioare celor liniare: în corpusul folosit în experiment, care a conținut 1200 de expresii referențiale, spațiul de căutare pentru legături co-referențiale s-a redus cu aproximativ 800 de unități.

Un alt tip de investigație empirică [Ide, Cristea, 2000] a urmărit frecvența referințelor evocative în comparație cu cele post-evocative și depistarea unor corelații între tipul de referințe și puterea de evocare a anaforilor. Studiul a comparat predicțiile avansate de VT relativ la domeniul de referențialitate evocativă cu cele ale modelului stivă al AST, corelând excepțiile (referințe ce nu se supun prevederilor celor două teorii) cu puterea de evocare a anaforilor (pentru VT excepțiile marchează, evident, referințe din categoria celor inferențiale). Într-o ordine descendentă a puterii de evocare (v. și [Gundel *et al.*, 1993]), tipurile de anafori care dau naștere la excepții sunt, în ordinea descrescătoare a frecvenței: referințe pragmatice > nume proprii > substantive comune > pronume. Pronumele constituie mijloace de referire foarte fragile. Un emitent al unui mesaj utilizează un pronume când e sigur că structura permite recuperarea cu ușurință a entității referită de pronume. Practic, exceptând câteva cazuri în care un pronume putea fi înțeles fără un antecedent (*our* în *our streets*, de exemplu), este imposibilă utilizarea unui pronume pentru a referi o entitate aflată în afara *de*. La extrema cealaltă se plasează referințele pragmatice ce-și recuperează antecedentul din cunoștințe exterioare discursului și numele proprii. Interesant este că această sortare descrescătoare a tipurilor de anafori dată de puterea de evocare se aliniază numărului de excepții raportate în cazul VT (56,3% - pragmatice, 22,7% - nume proprii, 16,0% - substantive comune și 5,0% - pronume) și nu are nici o semnificație în cazul AST (0,0% - pragmatice, 26,1% - nume proprii, 39,1% - substantive comune și 34,8% - pronume). Ea probează corectitudinea conjecturii coeziunii.

7. O proprietate de granularitate

Atunci când arborele de structură al discursului se modifică prin trecerea de la o granularitate mai fină la una mai grosieră, constrângerea de accesibilitate, conjeturată de VT, se păstrează.

Demonstrație

Să presupunem un arbore de discurs D pe care s-au calculat expresiile *head* și nervură ale nodurilor. O operație de mărire a granularității poate fi efectuată dacă o întindere de text, inițial repartizată în mai multe unități, și **pentru care există un nod, fie el n , care să o acopere strict în structura inițială**, este "compactată" într-o singură unitate de discurs mai mare ce va lua locul nodului n din structura inițială. Pentru a vedea în ce măsură o astfel de operație poate afecta accesibilitatea vom investiga rezultatul aplicării ei asupra expresiilor *head* și nervură.

Definiția expresiei *head*, punctul 1, obligă ca expresia *head* a ceea ce înainte de compactare era un nod interior, fie el n , să fi fost dată de concatenarea unui șir de etichete de noduri nucleare aflate în secvența de text subîntinsă de n . Să notăm acest nod, după compactare, cu o , etichetă compusă din secvența nodurilor terminale pe care le acoperă. De exemplu, pentru arborele din Fig. 9:

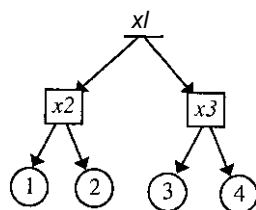


Figura 9: Un subarbore de "compactare"

dacă subarboarele cu rădăcina x_2 ar fi compactat, atunci eticheta sa ar trebui să fie notată 1-2, iar dacă întregul arbore aflat sub x_1 ar fi compactat, atunci eticheta sa ar trebui să fie notată 1-2-3-4 (e imposibil să avem un nod notat 2-3).

Acest lucru înseamnă că, aplicând o compactare asupra unui arbore, în expresiile *head* ale nodurilor sale, secvențe de noduri vor fi acum înlocuite cu etichete compuse care conțin cel puțin aceleași noduri, eventual mai multe, decât în expresiile originale. De exemplu, presupunând că în arborele de mai sus, nodurile nucleare sunt x_2 , x_3 , 1 și 3, atunci, dacă înainte de compactare am fi avut $head(x_1) = 1\ 3$, o expresie rezultată din concatenarea a două etichete, după compactarea întregului arbore vom avea $head(x_1) = 1-2-3-4$, adică o etichetă

compusă, dar care include etichetele nodurilor ce apăreau originală. Vom numi astfel de expresii - **expresii contraise și contrale**, unde e este expresia corespunzătoare de înainte de compactare, deci $contr^A(3) = 1-2-3-4$. Să remarcăm că secvențele de etichete contraise sunt formate întotdeauna din etichete de noduri adiacente, ceea ce permite comutarea funcțiilor *seq* și *contr*. $seq(contr(e)) = contr(seq(e))$.

Vom demonstra mai întâi că expresiile nervură ale nodurilor compactate sunt obținute din expresiile nervură originale în funcție de expresiile *head* originale cu expresiile contraise. Investigând definițiile expresiilor *head* și *nervură* poate constata că nici o altă modificare nu apare în expresiile nervură ale expresiilor contraise. Într-adevăr, cazul 1 se transcrie: expresia *head* a arborelui compactat reprezintă expresia *head* contrasă a arborelui original, $contr(h)$, cu h - expresia *head* a rădăcinii arborelui original.

Să presupunem acum că ne aflăm într-un nod n ale cărei expresii *head* și nervură pe arborele original, necompactat sunt, respectiv h și v . Pe arborele compactat, expresia *head* pe arborele compactat. Considerăm mai întâi cazul în care n este rădăcina, a cărui expresie *head* este $contr(h)$, și expresia *head* pe arborele necompactat. Dacă n este nucleul unui arbore, în cazul 2 (secțiunea 3.2), avem două subcazuri:

- n nu are un frate nenuclear în stânga: atunci nervura sa este $contr(v)$, nervura părintelui, adică $contr(h)$,
- n are un frate nenuclear în stânga de *head* ($contr(h)$): atunci nervura sa va fi $seq(contr(h), v)$, unde $seq(contr(h), v) = seq(contr(h), v)$.

Dacă n este un nod nenuclear, atunci conform cazului 3:

- n este în stânga: nervura sa este $seq(contr(h), v)$, unde $seq(contr(h), v) = seq(contr(h), v)$;
- n este în dreapta: nervura sa este $seq(simpl(contr(h)), v)$, unde $seq(simpl(contr(h)), v) = seq(simpl(contr(h)), v)$.

Folosind inducția, se probează în mod analog că expresia *head* a nodului n este o expresie contrasă și pentru cazul în care n este nucleul unui arbore neapărat imediat sub rădăcină, fiu al unui nod de nervură *contr*.

Cum expresia accesibilității este definită ca un prefix al expresiei *head* din care au fost îndepărtate marcasele, iar nervurile sunt expresii *head* eventual conținând mai multe etichete de noduri, înseamnă că expresia *head* pe arborele original satisface prima conjetură, cu alte cuvinte expresia *head* a unei unitate a unei expresii nervură și alta ce o precede, după compactare, de asemenea conjetura, pentru că nici o unitate nu a dispărut.

8. Discuții, aplicații ale teoriei

Plecând de la o reprezentare a structurii de discurs similară celei din RST și în care esențială este distincția dintre nucleu și satelit, VT definește nervura unui nod al arborelui ca secvența de unități ale discursului ce sunt suficiente pentru a rezuma/interpreta întinderea de text acoperită de nod în contextul întregului discurs. Presupunerea principală pe care se bazează noțiunea de nervură este că *referințele inter-unități sunt posibile cu precădere între unități ce se află într-o relație structurală*, chiar dacă acestea sunt dispuse la distanță una de alta în text. Mai departe, referințele se realizează cu precădere spre unități nucleare și doar în puține cazuri către sateliți, reflectând intuiția că nucleele găzduiesc ideile principale ale discursului. Acest lucru se regăsește în calculul expresiei nervurii pe arbori (binari) polarizați-stânga (pe orice nivel există un nucleu în stânga), în care orice referință se realizează dinspre un nucleu sau un satelit către un nucleu aflat în stânga (deși, nu orice nucleu). Făcând uz de echivalarea modelului stivă al lui Grbsz și Sidner [1985] cu structura de arbore utilizată de RST [Mann, Thompson, 1988], similaritate demonstrată de Moser și Moore [1996] și Marcu [1999], predicțiile VT asupra accesibilității referențiale sunt consistente cu cele ale modelului stivă. În cazurile în care însă arborele de discurs nu e polarizat-stânga (există cel puțin un satelit care precede nucleul său, deci care apare ca frate stânga pe un nivel al structurii), VT oferă o interpretare mai naturală a accesibilității decât modelul stivă, corectând totodată slăbiciunile acestuia. Într-adevăr, într-o secvență /4-satelit, 6-nucleu, deci în care *B* domină *A* în termenii AST, 6 ar trebui să apară în stivă poziționat sub *A*, deși el este procesat în secvență după *A*. Totodată, VT formalizează intuiția că într-o secvență de unități *A*, *S*, *C*, unde *A* și *C* sunt sateliți ai lui *S*, *C* nu poate accesa *A* din cauza interperierii unui nucleu, ce captează întreaga atenție.

Referențialitatea în lungul nervurilor este una naturală, ușor de interpretat și care, în general, nu necesită mijloace de evocare foarte puternice. Dimpotrivă, ieșirea din acest domeniu incumbă utilizarea unor mijloace de evocare anaforică viguroase. Pe acest criteriu se face distincția dintre referențialitate evocativă și ne-evocativă (sau inferențială), referințele evocative fiind detaliate în directe și indirecte, iar între cele ne-evocative remarcându-se referințele pragmatice, ce nu necesită un antecedent pentru înțelegere.

În privința coerenței discursului, VT utilizează domeniile de referențialitate pentru a introduce noțiunea de linie de argumentație și a deduce din ea pe cea de segment în sens ierarhic ce generalizează segmentul în sens clasic (așa cum este el utilizat în AST și CT). Totodată VT avansează conjectura că segmentul în sens ierarhic dă o mai corectă interpretare a porțiunilor de discurs ce se comportă din punctul de vedere al coeziunii și coerenței ca un tot unitar. Aplicând concluziile CT relative la coerența discursului în lungul segmentelor în sens ierarhic, CT poate fi generalizată pentru a o transforma într-o teorie globală a coerenței.

Au fost trecute în revistă o seamă de experimente care probează că prezumțiile VT sunt corecte și independente de limbă. Un aspect important îl constituie, de asemenea, faptul că prezumțiile VT sunt stabile la trecerea de la o granularitate mai fină la una mai grosieră în segmentarea discursului.

Aplicațiile VT se înscriu în trei direcții importante: rezoluția anaforei, parsarea discursului și rezumarea automată. În [Cristea et al., 2002a] și [Cristea et al., 2002b] este descrisă o arhitectură care acționează ca un motor general și ~" configurabil de rezoluție anaforică. Una dintre componentele oricărui model de rezoluție este o definiție a domeniului de referențialitate. Rezoluția anaforică se realizează, așadar, ghidată de structura de discurs.

În [Sereșan, Cristea, 2002] se propune o abordare inversă, în care cunoștințele asupra legărilor anaforice pot fi utilizate pentru corectarea structurii. Noi credem că procesul de rezoluție anaforică și de construire a structurii de discurs sunt interdependente într-un asemenea grad încât în analiza de discurs ele trebuie să aibă loc simultan. În interpretarea unui text există o intercondiționare reciprocă între referințe și structură care trebuie să conducă la obținerea acelei reprezentări în care constrângerile, acționând ca forțe, produc o stare de echilibru, ce trebuie să fie un fel de stare de energie potențială minimă a sistemului. Oamenii dispun de un mecanism cognitiv care le permite să ajungă în mod natural la cea mai plauzibilă interpretare a unui text. Acest lucru este răsplătit de atingerea unei stări mentale "confortabile" ce trebuie să-și aibă suportul în satisfacerea la maxim a unui sistem de constrângeri. În [Tablan et al., 1998] și [Cristea, 2000] se descrie un mecanism de parsare care modelează acest comportament uman. Prin combinarea unor scoruri contribuie de referințe cu scoruri contribuie de o analiză HCT se obține cea mai fluidă posibil structură de discurs (deci manifestând maximum de coerență) și care prezintă maximum de referințe pe nervuri (fiind deci cea mai coezivă posibil).

Noțiunea de *head* din VT este similară celei de mulțime de promovare (*promotion set*) pe care Marcu [2000] o utilizează pentru a obține un rezumat ghidat de structura de discurs. Să remarcăm că definiția nervurii presupune rezumarea ca o alternativă a înțelegerii unei unități de discurs în context. Credem că valențele teoriei nervurilor în realizarea unei strategii de rezumare focalizată [Mani, 2001] pe o anumită entitate sau segment de discurs au fost doar tangențial studiate până acum [Sofronie, 1999], [Postolache, 2001] și merită atenție în abordările viitoare. Credem, de asemenea, ca fiind interesantă o direcție de studiu care să aprecieze maniera în care nervura poate constitui un cadru de sub-specificare a structurii [Schilder, 2001], plecând de la observația că structuri diferite (dau nu fundamentale diferite) pot prezenta aceleași expresii ale nervurilor unităților componente.

Bibliografie

Brennan, S.E.; Walker Fredman, M. and Pollard, C.J. 1987. A centering approach to pronouns. *Proceedings of the 25th Annual Meeting of ACL*, Stanford, p 155-162.

Cornea, P. 1998. Introducere în teoria lecturii, Editura Polirom, Iași.

Cristea, D., and Webber, B.L. 1997. Expectations in incremental Discourse Processing. *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics*, Madrid.

Cristea, D., Ide, N., and Romary, L. 1998. Veins Theory: A Model of Global Discourse Cohesion and Coherence, *Proceedings of the 17th Coling and the 36th Annual Meeting of the ACL (COLING-ACL'98)*, Montreal, Canada p.281-185.

Cristea, D., Ide, N., Marcu, D., and Tablan, M.V. 2000. An Empirical Investigation of the Relation Between Discourse Structure and Co-Reference. *Proceedings of the 18th International Conference on Computational Linguistics COLING'2000*, Saarbrueken, p. 208-214.

Cristea, D. 2000. An Incremental Discourse Parser Architecture. Christodoulakis, D. (Ed.) *Natural Language Processing - NLP 2000*, Second International Conference, Patras, Greece, Lecture Notes in Artificial Intelligence 1835, Springer, p. 162-175.

Cristea, D. and Dima, G.E. 2001. An Integrating Framework for Anaphora Resolution. *Information Science and Technology*, Romanian Academy Publishing House, Bucharest, voi. 4, no. 3-4, p. 259-372.

Cristea, D., Postolache, O.D., Dima, D.E., Barbu C. 2002a. AR-Engine - a framework for unrestricted co-reference resolution. *Proceedings of the Third International Conference on Language Resources and Evaluation, LREC'2002*, Las Palmas, Spain, p. 2000-2006.

Cristea, D., Dima, D.E., Postolache, O.D., Mitkov, R. 2002b. Handling complex anaphora resolution cases. *Proceedings of the Discourse Anaphora and Anaphor Resolution Colloquium*, Lisbon, Portugal.

deEugenio, B. 1990. Centering theory and the Italian pronominal system. *Proceeding of Coling*, p. 270-275.

deEugenio, B. 1998. Centering in Italian. Prince, E., Joshi, A. and Walker, L. (eds.) *Centering in Discourse*, Oxford University Press.

Fox, B. 1987. Discourse Structure and Anaphora. Written and Conversational English. Cambridge Studies in Linguistics, Cambridge University Press.

Grosz, B.J. 1981. Focusing and description in natural language dialogues. Joshi, A., Webber, B. and Sag, I. (eds.) *Elements of Discourse Understanding*, Cambridge University Press, England, P. 85-105.

Grosz, B.J., Joshi, A.K. and Weinstein, S. 1995 Centering: A framework for modeling the local coherence of discourse. *Computational Linguistics*, 12(2), p. 203-225.

Grosz, B.J. and Sidner, C. 1986. Attention, intentions, and the structure of discourse. *Computational Linguistics*, 12(3), p. 175-204.

Gundel, J., Hedberg, N. and Zacharski, R. 1993. Cognitive Status and the Form of Referring Expressions. *Language*, 69, P. 274-307.

Halliday, M.A.K. and Hassan, R. 1976. Cohesion in English, Longman, London and New York.

Hovy, E. 1988. Planning coherent multisentential text. *Proceedings of the 26th Annual Meeting of the Association for Computational Linguistics (ACL-88)*, State University of New York, Buffalo, p. 163-169.

Ide, N. and Cristea, D. 2000. A Hierarchical Account of Referential Accessibility. *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics, ACU2000*, Hong Kong, p. 416-424.

Kameyama, M. 1998. Intrasentential Centering: A Case Study. Prince, E., Joshi, A. and Walker, L. (eds.) *Centering in Discourse*, Oxford University Press, p. 89-112.

Kintsch, W. and Van Dijk, T.A. 1975. Comment on se rappelle et on resume les histoires, *Langages*, 40.

Mani, I. 2001. Automatic Summarization. John Benkamin Publishing Company, Amsterdam/Philadelphia.

Mann, W.C. and Thompson, S.A. 1988. Rhetorical Structure Theory: Toward a Functional Theory of Text Organization. *Text*, 8(3), p. 243-281.

Marcu, D., 1999. A formal and computational synthesis of Grosz and Sidner's and Mann and Thompson's theories. *Proceedings of the Workshop on Levels of Representation in Discourse*. Edinburgh.

Marcu, D. 2000. The theory and practice of discourse parsing and summarization, The MIT Press, Cambridge, Massachusetts.

Miller, G. 1956. The magical number Seven, Plus or Minus Two: Some Limits on Our Capacity for Processing Information, *The Psychological Review*, voi. 63, p. 81-97.

Moser, M. and Moore, J.D. 1996. Toward a synthesis of two accounts of discourse structure. *Computational Linguistics*, 22(3), p. 409-419.

- Passonneau, R., J. 1995. Integrating gricean and attentional constraints. *Proceedings of IJCAI*.
- Postolache, O. 2001. Sumarizarea textelor. Lucrare de licență. Universitatea "Al.I.Cuza" Iași, Facultatea de Informatică.
- Rambow, O. (ed.) 1993. Intentionality and Structure in Discourse Relations. *Proceedings of a Workshop Sponsored by the Special Group on Generation of the Association for Computational Linguistics*, Ohio State University.
- Richardeau, F. 1969. La lisibilité. Langage-Typographie-Signes-Lecture, Paris.
- Schank, R. and Abelson, R. 1977. Scripts, plans, goals and understanding, Hillsdale, N.J.
- Schilder, F. 2001. Robust Discourse Parsing Via Discourse Markers, Topicality and Position. *Natural Language Engineering* 1, (1), p.1-22.
- Scott, D.R., de Souza, C.S. 1990. Getting the message across in RST-based text generation. Dale, R., Mellish, C. and Zock, M. (eds.) *Current Research in Natural Language Generation*, Academic Press, New York.
- Serețan, V. and Cristea, D., 2002. The Use of Referential Constraints in Structuring Discourse, *Proceedings of the Third International Conference on Language Resources and Evaluation, LREC'2002*, Las Palmas, Spain, p.1231-1237.
- Sidner, C. 1983. Focusing in the comprehension of definite anaphora. Brady, M. and Berwick, R.C. (eds.) *Computational Models of Discourse*, MIT Press.
- Sofronie, V. 1999. Implementări existente în sumarizarea textelor. SumVT. Lucrare de licență. Universitatea "Al.I.Cuza" Iași, Facultatea de Informatică.
- Strube, M. and Hahn, U. 1996. Functional Centering. *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics*, Santa Cruz, California., p. 270-277.
- Tablan, M.V., Barbu, C, Popescu, H., Hamza, R.O., Nita, C.I., Bocaniala, C.D., Ciobanu C. and Cristea, D. 1988. Co-operation and Detachment in Discourse Understanding. *Proceedings of the Workshop on Lexical Semantics and Discourse Structure, ESSLLI'98*, Saarbruecken.
- Walker, M., Iida, M., Cote, S. 1994. Japanese Discourse and the Process of Centering. *Computational Linguistics*, 20(2), p. 193-232.
- Walker, M.A. 1996. Limited attention and discourse structure. *Computational Linguistics*, 22-2.
- Walker, M.A. 1998. Centering, anaphora resolution, and discourse structure. In Marilyn A. Walker, Aravind K. Joshi, and Ellen F. Prince, editors, *Centering in Discourse*. Oxford University Press.

DLIR - un sistem de căutare documentară multilingv

Amalia TODIRAȘCU
INRIA Lorraine, LORIA, Campus scientifique BP 239,
54506 Vandoeuvre-les-Nancy Cedex, France,
todirasc@loria.fr

Abstract

Această lucrare prezintă un sistem de căutare documentară bilingv francez-român pentru un domeniu limitat, cel al securității computerelor. Căutarea și indexarea documentelor se realizează utilizând o ontologie specifică domeniului. Identificarea instanțelor conceptelor în texte sau în întrebările utilizator se realizează cu ajutorul unor tehnici robuste de analiza limbajului natural, combinate cu o ontologie specifică domeniului.

Introducere

Sistemele de căutare de informații clasice indexează o bază de documente, folosind o listă de cuvinte cheie extrase din documentele respective. Scopul acestor sisteme este de a regăsi documentele care sunt relevante în comparație cu o întrebare lansată de un utilizator. Sistemele de căutare a informațiilor interpretează întrebările utilizatorului, încearcă să găsească un index (unul sau mai multe cuvinte-cheie) care apare în întrebare sau care este similar unui termen din întrebare. Fiecare cuvânt-cheie care aparține indexului este asociat unui document în care cuvântul cheie a fost folosit. Răspunsul sistemului conține un număr de documente care sunt relevante în raport cu întrebarea utilizatorului. Fiecare sistem de căutare definește un criteriu de relevanță pentru a selecta documentele propuse ca răspuns. Aceste sisteme sunt evaluate pe baza a doi parametri: rapel (numărul de documente regăsite ca răspuns/numărul total de documente relevante care au fost indexate) și precizie (numărul de documente relevante regăsite de sistem/numărul de documente regăsite). În cazul unui sistem de căutare multilingv, răspunsul la o întrebare poate conține mai multe documente relevante, chiar dacă sunt scrise în alte limbi decât cea în care a fost formulată cererea.

Sistemele de căutare de informații clasice oferă utilizatorului răspunsuri imprecise sau vide. Aceste răspunsuri imprecise apar datorită faptului că majoritatea sistemelor de căutare documentară folosesc doar cuvinte-cheie sau expresii extrase cu ajutorul metodelor statistice, ignorând problemele legate de complexitatea limbajului natural: ambiguitatea (un cuvânt poate avea mai multe sensuri) sau polimorfismul (un concept poate fi exprimat în mai multe moduri). În plus, un sistem care își propune să facă o căutare într-o bază de date multilingvă trebuie să fie capabil să găsească informația cerută în orice document disponibil, indiferent de limba în care a fost scris. Unele sisteme de căutare multilingvă folosesc drept indecși cuvinte cheie pentru fiecare limbă, alte sisteme propun indexare pe baza unui index comun, alcătuit din concepte.

O alternativă la sistemele de indexare clasice sunt cele care folosesc structuri sintactice sau conceptuale pentru a indexa baza de documente. Acestea nu sunt foarte numeroase, pentru că pe de o parte, ontologiile generice nu sunt disponibile decât în număr prea restrâns (WordNet [19] și Corelex [3] sunt doar două exemple de resurse libere). Folosirea conceptelor unei ontologii permite asocierea unor termenii din limbi diferite, de aceea am ales o metodă de indexare conceptuală, care va fi în secțiunea 4.

Traducerea termenilor care sunt folosiți ca indecși ridică probleme într-o aplicație de căutare de informații, unui termen se pot asocia mai multe traduceri cu sensuri diferite sau sintagme. Rezolvarea problemelor specifice limbajului natural (ambiguitate, traducere automată) necesită resurse lingvistice importante pentru fiecare limbă care este tratată de către sistem, dacă aplicăm tehnicile clasice de analiza limbajului natural. Tehnicile clasice de analiză sintactică nu sunt adaptate sistemelor de căutare documentară, datorită dimensiunilor prea mari ale bazei de documentare și a resurselor lingvistice necesare. Pe de altă parte, textele specifice necesită adaptarea resurselor lingvistice (dicționare, gramatici locale) folosite de analizoarele sintactice. Tehnicile robuste de analiză sintactică utilizate în domeniul extragerii de informații din texte (GATE [6], FASTUS) sunt dedicate rezolvării unor probleme precise (identificarea numelor proprii, ale grupurilor nominale simple). Printre acestea, automatele cu număr finit de stări [5], colocații [9] sau liste de pattern-uri sintactice (reprezentând structura sintactică a grupului nominal simplu a grupului prepozițional) sunt resursele lingvistice necesare pentru aceste componente. Aceste tehnici au avantajul de a fi robuste, de a putea trata o cantitate importantă de informații în timp real, precum și de a fi portabile de la un domeniu și/sau o limbă la alta.

Într-o aplicație de căutare de informații pe un domeniu restrâns, utilizatorul dorește să obțină răspunsuri mai precise decât pentru texte cu caracter general. Aceasta impune folosirea de tehnici adaptate acestui tip de aplicații, bazate pe existența unei baze de cunoștințe din domeniu. Pentru a evita problemele legate de traducerea termenilor dintr-o limbă într-alta, propun folosirea unei ontologii specifice domeniului, în vederea ameliorării preciziei. În acest context, voi prezenta

o metodologie de extragere a conceptelor candidat din corpus. Acestea sunt folosite de către un expert uman pentru a îmbogăți o ontologie existentă precum și pentru a crea o reprezentare sub formă de concepte a documentelor. De asemenea voi prezenta o metodă de indexare a documentelor pe baza acestei ontologii, metodă care modifică metoda clasică de indexare semantică latentă.

2. Ontologii

Noțiunea de ontologie este dificil de definit, mai multe puncte de vedere coexistă. Pentru a simplifica, vom considera că o ontologie este un model restrâns al unui domeniu specific, format din mulțimea claselor de obiecte ce populează acest domeniu și a relațiilor lor cu celelalte clase din domeniu.

Ontologiile reflectă un anumit grad de subiectivitate din partea expertului ce a definit-o. Fiecare expert poate propune un ansamblu de clase de obiecte ce trebuie incluse în descrierea ontologiei, care poate fi diferit de clasele propuse de alți experți din domeniu.

O problemă a acestor ontologii este legată de portabilitate. O aplicație definită pentru un anumit domeniu dat va trebui adaptată unui alt domeniu prin construirea unei ontologii corespunzătoare. Construirea lor manuală este dificilă și trebuie ținut cont de posibilele redundanțe, erori, informații care lipsesc sau incoerențe ce pot fi introduse în baza de cunoștințe de către expertul uman care o construiește.

În ultimii ani, s-au făcut eforturi deosebite pentru a putea reutiliza ontologiile existente: dezvoltarea unor formate standard: (Knowledge Interchange Format - KIF), Ontology Interface Layer - (OIL) [8], dezvoltate în cadrul proiectului Semantic Web (<http://www.semweb.org>). Aceasta permite reutilizarea ontologiilor existente de către alte sisteme și aplicații, în ciuda erorilor care pot apărea în urma construirii manuale.

Pentru a evita problemele legate de crearea manuală a ontologiilor, au fost propuse mai multe metode semi-automate de extragere a ontologiilor din corpusuri. Acestea disting mai multe etape:

- identificarea termenilor (instanțele conceptelor exprimate în limbaj natural);
- identificarea relațiilor între termi;
- identificarea relațiilor între termi și concepte.

Majoritatea acestor etape necesită validarea rezultatelor de către un expert uman, care va asocia o interpretare claselor de termi și relațiilor între două mulțimi de termi). Metodele statistice interpretează contextele existente și

regroupează termii cu contexte identice în aceeași clasă [1], [7]. Relațiile între termii sunt interpretate pe baza informațiilor de subcategorizare (structura predicat argument) asociate verbelor. Dezavantajul metodelor statistice este acela că necesită corpusuri adnotate de talie importantă pentru a putea învăța, iar rezultatele (clasele obținute) nu pot fi întotdeauna interpretate.

În comparație cu metodele statistice, metodele bazate pe inferențe logice propun proceduri semi-automate pentru a verifica validitatea cunoașterii existente. Conceptele noi, deduse de către regulile de inferență, sunt adăugate ierarhiei domeniului dacă sunt coerente cu cunoașterea existentă. Relațiile pot fi identificate folosind cunoștințe legate de subcategorizare [4] sau interpretând relațiile substantiv-modificator. Supragenerarea de concepte și costul verificării incoerențelor și inconsistentelor cunoașterii sunt principalele neajunsuri ale metodei. Mai multe formalisme de reprezentare a cunoștințelor pot fi folosite în astfel de aplicații. Am ales logicile terminologice datorită avantajelor pe care le prezintă acestea.

2.1. Logici terminologice

Logicile terminologice (LT) sunt formalisme de reprezentare a cunoștințelor care sunt derivate din formalismul rețelelor semantice, dar sintaxa și semantica lor sunt bine definite. Ele combină proprietăți ale sistemelor orientate-obiect, ale sistemelor de tip frame și ale logicilor modale.

LT propun o organizare ierarhică a cunoașterii, pe două nivele: unul conceptual (T-Box), care descrie clasele abstracte conținând obiectele relevante pentru modelarea domeniului și un nivel aserțional (A-Box), conținând instanțele claselor. Clasele de obiecte (concepte) sunt descrise de relații (numite roluri) cu alte concepte, și cu atributele lor (rolurile cu valori atomice).

2.1.1. Sintaxa și semantica logicilor terminologice

Operatorii LT sunt inspirați de logica de prim ordin:

Operator	Operator Logic	Semantica
$D = \text{SOME } R \ C$	$\exists x R(y,x) \& C(x) \& D(y)$	Există cel puțin o instanță a lui C în relația R cu o instanță a lui D
$D = \text{ALL } R \ C$	$\forall x(R(y,x) \& D(x)) \& C(y)$	restricționează co-domeniul relației R
$D = \text{AND } C1 \ C2$	$C1 \ \& \ C2$	Conjunția de descrieri conceptuale
$D = \text{OR } C1 \ G2$	$C1 \ \vee \ C2$	Disjunția de descrieri conceptuale
$C1 \ c \ C2$	$C1 \ c \ C2$	Axiom: C1 conține condiții necesare pentru C2
$D = \text{NOT } C$	$\neg C$	complementul conceptului C
$D = 3n.R.C$	$\exists y1 \dots yn (1 \leq j \leq n, R(x, yi) \& C(yi) \& D(x))$	Există cel puțin n obiecte de tip C în relația R cu o instanță a lui D

Figura 1. Operatorii în LT

Folosind toți acești operatori, sau doar o parte a acestora, expresivități sunt posibile: definirea conceptelor și a rolurilor ALL (ALL, AND, OR, NOT ca operatori, axiomele conceptuale), rolurilor tranzitive (R+), a rolurilor inversabile (I), a ierarhiilor atributelor (f) sau a restricțiilor numerice.

Folosind toți acești operatori, sau doar o parte a acestora, expresivități sunt posibile: definirea conceptelor și a rolurilor ALL (ALL, AND, OR, NOT ca operatori, axiomele conceptuale), rolurilor tranzitive (R+), a rolurilor inversabile (I), a ierarhiilor atributelor (f) sau a restricțiilor numerice.

Unele comenzi LT sunt explicate mai jos. CN este un nomen conceptual este o descriere conceptuală (orice combinație de operatori ALL, AND, OR, NOT). Comenzile LT sunt inspirate de formalismul KRSS ([2]):

1. (define-concept CN C) - definește un nou concept conceptuală;
2. (instance IN C) - definește o instanță a unui concept conceptuală;
3. (implies C1 C2) - introduce o nouă axiomă conceptuală condițiile C1 necesare pentru descrierea conceptuală C2.

LT sunt fragmente decidabile ale logicii de prim ordin și oferă mecanisme logice pentru a identifica subsumarea, regăsirea instanțelor care unesc mai multe concepte. Clasificarea este o ordonare pe niveluri a conceptelor, în raport cu relația de subsumare. Există algoritmi pentru verificarea coerenței și consistenței cunoștințelor.

Câteva exemple de comenzi:

(concept-subsumes? C1 C2) testează dacă C1 subsumează C2

(concept-parents C) regăsește strămoșii direcți ai conceptului C

(concept-children C) regăsește fiii direcți ai lui C

(classify-tbox) calculează toate relațiile de subsumare în T-Box

(concept-children C) regăsește fiii direcți ai lui C

(concept-instances C) regăsește toate instanțele conceptului C

2.2. Logici terminologice pentru sisteme de extragere a informațiilor

Rolul cunoștințelor specifice unui domeniu într-un sistem de extragere a informațiilor este acela de a valida reprezentarea semantică a informațiilor potențial relevante, identificate în text prin tehnici de procesare naturală a limbajului. Aceste entități pot fi folosite pentru a adăuga noi concepte la cunoașterea existentă. Cea mai mare parte a sistemelor de extragere a informațiilor folosește

robuste pentru identificarea candidaților și entitățile candidat sunt validate de către o interpretare semantică. Sistemele de extragere a informațiilor pot folosi cunoaștere implicită, cum ar fi relațiile de hiponimie/hiperonimie.

Logicile terminologice prezintă avantajul de a lucra cu date semi-structurate sau incomplete. Nu este necesară definirea explicită a unor valori ca instanțe ale unor concepte. Valorile implicite nu sunt utilizate de către logicile terminologice. Unele valori ale rolurilor sunt lăsate nespecificate ca în următorul exemplu:

```
(define-concept computer (and physicalobject (some hasOperatingSystem
OSystem) (some hasType Type)))
```

```
(define-primitive-concept Type)
```

```
(define-primitive-concept OSystem)
```

```
(instance suni (and computer (some hasType SparcStation)))
```

În acest exemplu, vom ilustra faptul că definițiile implicite sunt acceptate de către logicile terminologice (**SparcStation** este definit explicit de către o instanță sau un subconcept al conceptului **Type**). Nu este definită explicit nici o instanță a rolului **hasOperatingSystem**.

Aceste proprietăți nu sunt interesante pentru aplicația noastră, dar erorile sunt posibile, iar cunoașterea domeniului este incompletă.

Relațiile de hiperonimie sau hiponimie sunt tratate cu ajutorul relațiilor de subsumare între conceptele domeniului. De exemplu, dacă un concept candidat este identificat în text ca:

```
(instance x (and PC (and hasOperatingSystem Linux)))
```

```
(define-concept PCcomputer (and computer (some hasType PC)))
```

x este de asemenea o instanță a conceptului **computer**.

```
(instance y (and Password (some hasUser Root)))
```

```
(define-concept Password (and String (some hasAtr secret) (some
hasBelongs User»))
```

```
(define-concept System (some hasUser User))
```

```
(define-concept Root User)
```

Pentru aplicația noastră avem nevoie de o logică terminologică care să permită raționament la nivel de instanță, să permită lucrul în contextul unei lumi deschise, precum și proceduri optimizate de calcul a relațiilor de subsumare sau de clasificare. Printre puținele sisteme care implementează raționament la nivel de instanță am ales RACER ([10]), fiind unul dintre cele mai performante și mai complete.

În secțiunea următoare voi prezenta metoda de extragere a termenilor din texte folosind sistemul DLIR [16]. Textele vor fi traduse într-o reprezentare conceptuală unică, folosită ca index, permițând regăsirea informațiilor în mai multe limbi.

3. Arhitectura

Sistemul DLIR conține mai multe module: un modul de analiză semantică robustă, un modul de întreținere a ontologiei domeniului, un modul de extragere a termenilor bazat pe celelalte două module. În cele ce urmează vom prezenta aceste module în detaliu.

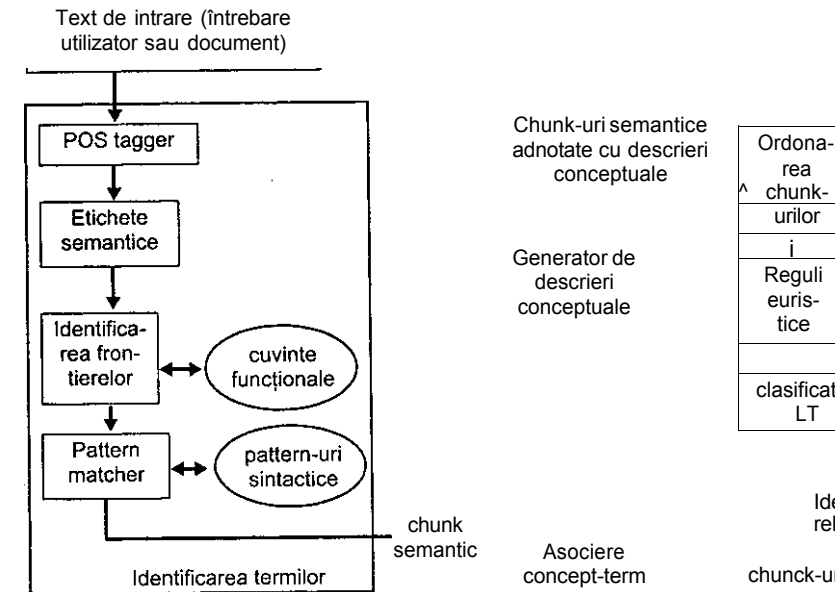


Figura 1: Instanțe ale conceptelor care apar în întrebare

3.1. Analiza sintactică robustă

Acest modul este dedicat identificării termenilor posibili, utilizând o analiză robustă, și resurse specifice domeniului (o listă de corespondențe între termeni și concepte). Termii sunt combinați conform unor reguli euristice pentru a identifica concepte complexe. Aceste concepte sunt validate ulterior, apelându-se la ontologia domeniului. Eventualele concepte valide sunt stocate în ontologia existentă. Acest modul conține mai multe submodule implementate în Perl și CLIPS (modulul care aplică regulile de combinare a termenilor). Modulul de extragere a termenilor bazat pe chunk semantic a fost propusă pentru a identifica termenii candidați [16]. Modulul a fost testat pentru limba franceză, dar cum resursele folosite pentru identificare sunt relativ independente de limba pentru care a fost construită, este posibilă extinderea ei și pentru limba română, după cum voi arăta mai

3.2. Identificarea chunk-urilor semantice

Scopul principal al acestui modul este acela de a identifica secvențele de cuvinte care corespund celor mai semnificative concepte ale domeniului (*chunk-uri semantice*).

Un *chunk semantic* conține un pattern sintactic simplu (grup substantival simplu, grup verbal) și este delimitat de doi separatori de clauze.

Separatorii sunt cuvinte funcționale, verbe auxiliare, sau anumite sintagme prepoziționale.

Exemplu. "la victime d'une intrusion inattendue"

[victima unei intruziuni neașteptate]

În acest exemplu, "victima" și "unei intruziuni neașteptate" sunt chunk-uri semantice, care conțin informația relevantă.

Modulul conține mai multe submodule: un POS tagger, un tagger semantic, un identificator de frontiere și un pattern matcher. Identificarea chunk-urilor semantice este bazată pe informația lexicală, propusă de POS tagger.

3.2.1. Part-Of-Speech tagging

Modulul care este dedicat identificării părților de vorbire asociate cuvintelor (folosind WinBrill, antrenat pentru franceză pe baza unui set de date propuse de Institut National pour la Langue Française [11]) identifică cuvintele conținut (substantive, adjective, verbe) și cuvintele funcționale (prepoziții, conjuncții etc).

Taggerul Brill folosește un set de reguli contextuale și lexicale (bazate pe identificarea prefixelor și a sufixelor), învățate pe baza textelor adnotate, pentru a identifica partea de vorbire pentru cuvintele necunoscute.

Pentru limba română, se folosește QTAG adaptat pentru limba română [17], datorită performanțelor foarte bune (98% rezultate corecte).

3.2.2. Tagger-ul semantic

Tagger-ul semantic conține un pattern matcher, care consultă un dicționar de talie redusă. Acesta conține o listă cu cele mai frecvente cuvinte și un set de sintagme asociate descrierilor conceptuale corespunzătoare.

Setul de descrieri conceptuale a fost stabilit de către un expert pe baza unei liste de cuvinte și segmente repetate obținute dintr-un corpus reprezentativ (200,000 cuvinte). Un segment repetat este o succesiune de cuvinte care intervin într-un text cel puțin de două ori [14].

Acest modul asociază fiecărui cuvânt conceptul sau descrierea conceptuală din dicționar. Un astfel de dicționar este creat pentru fiecare limbă care este tratată de către sistem.

3.2.3. Modulul pentru identificarea frontierelor

Acest modul identifică separatorii (cuvinte funcționale și sintactice mai complexe) care delimitează chunk-urile semantice. folosește rezultatul POS tagger-ului (care identifică cuvintele funcționale și un set de sintagme (constituenți sintactici care conțin auxiliare compuse). Setul de fraze este construit ca rezultat al studiilor corporale pentru franceză și română (200,000 cuvinte pentru fiecare limbă) grupurilor nominale și prepoziționale (determinanți, prepoziții) sunt candidați pentru identificarea separatorilor de chunk-uri semantice. reprezintă anumite relații potențiale între concepte.

3.2.4. Pattern matcher

Scopul acestui modul este de a identifica nucleul chunk-urilor semantice, nucleu care este reprezentat de un grup nominal simplu sau un grup

Exemple. Un grup nominal simplu (în franceză) este identificat a toarele regulii N -) NP, N ADJ -) NP, Def N -) NP, Def N ADJ -) NP. Pentru limba română, un grup de reguli posibil poate fi: Indefart N Adj -) NP, IndefArt N N -) NP

3.2.5. DLgen

Acest modul interpretează informația propusă de POS tagger în mod automat o definiție de concept. Un expert trebuie să verifice acest modul. Câteva exemple de reguli propuse pentru generarea de reguli simple (valabile pentru ambele limbi):

- S1/N S2/ADJ este asociat definiției (define-concept S1 (SOME hasAtr "S2"))
- S1/N S2/NNP este asociat definiției (define-concept S1 (SOME hasName "S2"))
- S1/ADJ S2/N este asociat definiției (define-concept S1 (SOME hasAtr "S1"))
- Verbele sunt traduse ca nume de roluri: S1A/B este asociat definiției (define-concept S1 (SOME hasS1))

Unele pattern-uri identifică negațiile, chiar dacă este necesar să enumerăm toate posibilitățile și să detectăm corect domeniul negațiilor.

- sans/ADV S1/N este asociat definiției (define-concept S1 (NOTS1))
- nici_unul/ADV S1/N este asociat definiției (define-concept S1 (NOTS1))

Rezultatele propuse de DLgen sunt 61% corecte datorită faptului că regulile sunt incomplete. Leșirea este validată de un expert folosind clasificatorul LT pentru a verifica definițiile conceptuale obținute în mod automat. Rezultatul este că fiecare structură are asociată o descriere conceptuală.

3.3. Relații între termi

Acest modul folosește inferențele LT, ca și regulile de sintaxă, pentru a combina descrierile conceptuale asociate fiecărui chunk semantic. Folosim un criteriu de ordonare al chunk-urilor, precum și reguli de combinare a conceptelor pentru a crea concepte complexe. Descrierile rezultante sunt validate de clasificatorul LT.

3.3.1. Ordonarea chunk-urilor

Modulul interpretează ordinea chunk-urilor și poziția chunk-urilor în propoziție.

Clasificăm chunk-urile în două categorii: **chunk-uri principale** și **chunk-uri secundare**. Chunk-urile principale corespund noțiunii de nucleu propuse de către teoriile lingvistice clasice.

Chunk-urile secundare joacă rolul unui modificador, care adaugă informații suplimentare sensului nucleului. Chunk-urile secundare pot lipsi, dar restul propozițiilor este corect. Aceste exemple de reguli definesc chunk-uri diverse:

- chunk-urile care urmează după un verb la gerunziu sau un auxiliar plus un verb la participiu sunt *chunk-uri secundare*;
- verbele sunt întotdeauna chunk-uri principale.

Exemplu:

'[Main Les atacs Main] [Main ont commence Main] [Second â utilisier les faux comptes Second]'

'atacurile au început utilizând conturi false'

Cele doua chunk-uri principale detectate în exemplul de mai sus sunt subiectul propoziției și verbul principal. Chunk-ul secundar este adnotat astfel pentru că urmează după prepoziția **â**.

3.3.2. Reguli euristice

Regulile sunt stabilite de către expert pe baza unui studiu asupra corpusului reprezentativ pentru fiecare limbă. Corpusul a fost adnotat cu categoria lexicală propusă de POS tagger și adnotat manual cu descrierile conceptuale. Setul de reguli euristice este stabilit pe baza unei liste de pattern-uri de forma *<Chunki >?x/FW<Cftun/c1 >*.

Exemplu de reguli euristice sintactice: dacă o prepoziție este două chunk-uri semantice și prepoziția asociază un modificador, atunci putem combina descrierile conceptelor chunk-uri într-o descriere semantică mai complexă, conceptele fiind cel de modificador:

if (<MainChunk1> <Border> <SecChunk2>)

and (Noun in MainChunk1)

and (Modifier in SecChunk2)

then (and sem(MainChunk1) (some hasModifier

Pentru română, un exemplu de regulă de combinare

următoarea: dacă un verb la gerunziu se găsește între un predicativ și un grup nominal, atunci rolul care leagă conceptele este un modificador.

Fiecare pattern este asociat unui cuvânt țintă care este folosit pentru aplicarea regulilor. Prepozițiile, verbele la modul participiu sunt exemple de cuvinte asociate regulilor euristice. Un număr de 21 de reguli pentru română au fost create. Leșirea acestor reguli va fi o serie de *chunk-uri complexe*, ce vor fi verificate de către expert, cu ajutorul ontologiei domeniului, care este în dezvoltare. Rezultatele propuse de acest modul conțin în mare parte concepte care nu sunt identificate. Baza de reguli poartă denumirea de urma studierii unui corpus de dimensiuni mai importante.

3.4. Indexare semantică

O posibilitate de indexare a documentelor este aceea de a crea concepte drept index și nu cuvinte cheie. O metodă eficientă reprezintă indexarea semantică latentă. Această metodă creează un set de document-cuvinte cheie și folosește tehnici de descompunere în valori proprii. În acest fel se elimină coloanele și rândurile rare (datorat faptului că mulți termi apar foarte rar). Propunem să creștem numărul de cuvinte cheie care fac parte din ontologie în locul cuvintelor cheie. Este posibil să găsim de căutare a informațiilor multilingv să avem diferențe între o limbă și alta. Avantajul este că putem folosi drept index conceptele din ontologie. Pentru aplicația noastră am folosit o ontologie care conține 54 de concepte și 34 de relații.

Numărul de concepte este mai redus decât numărul de cuvinte cheie, în special relațiile între termi.

Elementele matricii conțin o pondere $weight(C,i)$ calculată

$$weight(C,i) = \frac{f(C,i)}{\sum_{j=1}^7 S/(C,J)}$$

pentru fiecare concept, codificând frecvența instanțelor conceptului în document și frecvența instanțelor în toate documentele indexate de sistem.

$f(C,i)$ - frecvența conceptului în documentul i ;

Conceptele sunt legate prin rolurile dintre acestea. Frecvența unui concept care este situat în ierarhie foarte sus este compus din suma frecvențelor instanțelor sale. Instanțele conceptelor în LT sunt instanțele tuturor subconceptelor și ale instanțelor sale directe.

Indexarea documentelor se face aplicând metodele de extragere a termenilor prezentate în secțiunea precedentă, înainte de a exploata sistemul. Se folosesc conceptele ontologiei care a fost construită manual. O serie de concepte mai generale ar putea fi obținute combinând ontologia specifică domeniului cu WordNet ([16]).

Evaluarea acestui sistem a fost realizată pentru un set restrâns de întrebări (50) numai pentru limba franceză. Rezultatele au fost comparate cu cele furnizate de un sistem care folosește cuvinte-cheie pentru indexare. Pentru 74% din întrebări răspunsurile sistemului (rapel și precizie) au fost comparabile cu cele obținute prin metoda de indexare bazată pe cuvinte-cheie. În celelalte cazuri, răspunsurile au fost mai slabe decât indexarea pe baza de cuvinte cheie. Ontologia folosită este departe de a fi completă, ceea ce a dus la neidentificarea unor termi, de asemenea regulile de formare a conceptelor sau de generare a descrierilor conceptuale sunt incomplete.

4. Concluzii și perspective

Articolul prezintă o modalitate de a folosi ontologia unui domeniu pentru căutare de informații bilingvă în limbile franceză și română.

Sistemul integrează tehnici de analiză sintactică robustă pentru extragerea celor mai relevante chunk-uri semantice. Metoda folosește o ontologie a domeniului construită manual. Pentru evaluarea pertinentă a metodelor de indexare pe bază de concepte, ontologia va fi actualizată și extinsă cu ajutorul raționamentelor propuse de logicile terminologice, ca și folosirea cunoștințelor sintactice, folosite pentru extragerea unei reprezentări semantice pentru texte și întrebări. Expertul uman trebuie să intervină pentru a decide dacă conceptele identificate în texte pot fi adăugate ontologiei domeniului.

Referințe bibliografice

- [1] Assadi, H., Bourigault, D., 2000, Analyse syntaxique et construction d'ontologies à partir des textes. în J.C. G.Kassel, D.Bourigault (eds.) - *Ingenierie des connaissances récentes et nouveaux défis*, Eyrolles Publishing House, pp. 1-12.
- [2] Baader, F., Hollunder, B., 1991. A Terminology Representation Systems with Complete Inference Algorithms. *Proceedings of the Workshop on Processing Declarative Knowledge*.
- [3] Buitelaar, P., 1998. CORELEX: Systematic Polysemy and its Resolution. Ph.D. thesis, Brandeis University, Department of Computer Science.
- [4] Capponi, N., Toussaint, Y., 2000, Interpretation de clauses et généralisation de structures predicat-argument. în J.C. G.Kassel, D.Bourigault (eds.), *Ingenierie des connaissances récentes et nouveaux défis*, Eyrolles Publishing House, pp. 1-12.
- [5] Chanod, J.P., 1999. Natural Language Processing and Information Extraction. M.T. Paziienza (ed.), *Information Extraction*, Springer-Verlag, pp. 17-31.
- [6] Cunningham, H., Wilks, Y., Gaizauskas, R.J., 1996. Natural Language Trends and Software Infrastructure for NLP. în *Proceedings of the conference on New Methods in Natural Language Processing*, Bilkent University, Turkey, 1996, pp.1-12.
- [7] Daille, B., 1996, Study and Implementation of Combination Methods for Automatic Extraction of Terminology. In J.Klavans, P. Resnik, *Balancing Act - Combining Symbolic and Statistical Methods in Natural Language*, MIT Press, pp. 49-66.
- [8] Fensel, D. et al., 2000, OIL in a nutshell. în R. Dieng et al. *Knowledge Acquisition, Modeling, and Management, Proceedings of the Knowledge Acquisition Conference (EKAW-2000)*, Lecture Notes in Computer Intelligence, LNAI, Springer-Verlag.
- [9] Heid, U., 2000, A linguistic bootstrapping approach to automatic extraction of candidates from German text, *Terminology*, pp. 161-180.
- [10] Haarslev, V., Muller R., 2001, Description of the RAKIT system for Terminology Applications, *Proceedings of the International Workshop on Terminology Logics (DL-2001)*, Stanford, USA, 1.-3. August 2001, pp. 1-12.
- [11] Lecomte, J., Le Categoriseur BRILL14-JL5/WINBRILL. CNRS report, December 1998.
- [12] Riloff, E., Lorenzen, J., 1999, Extraction-based Terminology Generation. Generating Domain-Specific Role Relationships Automatically. Strzalkowski, *Natural Language Information Retrieval*, Springer Publishers, pp. 167-196.



- [13] Riloff, E., Shepherd, J., 1997, A Corpus-Based Approach for Building Semantic Lexicons. în *Proceedings of the Second Conference on Empirical Methods in Natural Language Processing*.
- [14] Rousselot, F., Frath, P., Oueslati, R./Extracting concepts and relations from Corpora. In *Proceedings of the Workshop on Corpus-oriented Semantic Analysis, European Conference on Artificial Intelligence, ECAI 96, Budapest, 12 August 1996*.
- [15] Schimid, H., 1994, Probabilistic Part-of-Speech Tagging Using Decision Trees, *Proceedings of the International Conference on New Methods in Language Processing, Manchester, United Kingdom*.
- [16] Todiraşcu, A., 2001, Semantic Indexing for Information Retrieval Systems, Ph.D. Thesis, University Louis Pasteur of Strasbourg, France, March 2001.
- [17] Tufiş, D., Mason O., Tagging Romanian Texts: a Case Study for QTAG, a Language Independent Probabilistic Tagger. în *Proceedings of the First International Conference on Language Resources and Evaluation (LREC) Granada, Spain, 1998*, pp. 589-596.
- [18] Vilain, M., 1999, Inferential Information Extraction. In M.Pazienza (ed.), *Information Extraction*, LNAI 1714, Springer-Verlag, pp.95-119.
- [19] Vossen, P., *Introduction to EuroWordNet*, Kluwer Academic Publisher, 1998.
- [20] Zweigenbaum, P., Consortium MENELAS, 1995, MENELAS: Coding and Information Retrieval from Natural Language Patient Discharge Summaries. In M.-F. Laires, M.J. Ladeira, J.-P. Christensen (eds.) - *Advances in Health Telematics*, IOS Press, Amsterdam, pp.82-89.

Mediu hermenofor pentru asistarea învățării unor concepte dintr-o limbă străină

Ștefan TRĂUȘAN-MATU
 Universitatea "Politehnica" București, Facultatea de Automatică și
 Calculatoare, Institutul de Cercetări pentru Inteligență Artificială al
 Academiei Române
 email: trausan@cs.pub.ro, trausan@racai.ro
 URL: www.racai.ro/~trausan

1. WWW, o prezență din ce în ce mai comună

În mai puțin de zece ani, rețeaua globală de documente World Wide Web (WWW sau, pe scurt web), a devenit omniprezentă și este posibil ca într-un timp nu prea lung să înlocuiască o mare parte din cărți, televizorul, cinematograful, ziarele și revistele (toate acestea fiind deja disponibile pe web) și, în plus să furnizeze chiar posibilitatea imersiunii în realități virtuale. Un singur exemplu care să fie suficient: anul trecut rezultatele bacalaureatului au fost publicate pe web.

WWW a atins deja dimensiuni comparabile cu imensa Bibliotecă a Congresului SUA. Extinderea sa este datorată ușurinței cu care poate fi parcursă și către oricine are un calculator și de simplitatea cu care se poate publica ceva pe el. Pe de altă parte, costul accesului la resursele web este de cele mai multe ori infim.

WWW este un hipertext extins la scara întregului glob prin rețeaua mondială de calculatoare Internet. Pe fiecare calculator pot fi plasate unul sau mai multe documente care constituie noduri (pagini) în hipertext. Oriunde într-o astfel de pagină poate exista o legătură la o altă pagină, de pe același calculator sau de pe altul, în acest mod putând fi unite informații aflate în locuri diferite. O nouă pagină pentru web poate fi creată ușor chiar de utilizatori nu neapărat profesioniști în informatică, în acest scop existând mai multe editoare de texte specializate.

Termenul de hipertext se pare că provine de la termenul de spațiu hiperbolic sau hiperspațiu, apărut în 1704 și folosit de matematicianul F. Klei pentru geometria cu mai multe dimensiuni [Rad91]. Din această perspectivă, un hipertext este un text cu mai multe dimensiuni explicite (față de doar o dimensiune în cazul textului liniar). De fapt, orice text are implicit mai multe dimensiuni deoarece, chiar dacă forma de prezentare a unui text este liniară, pe hârtie, în

există o structură implicită, dată de discurs. De asemenea, există conexiuni implicite, subiective între părți ale textului, concepte legate între ele, Hipertextul este o organizare a unui text în care toate aceste legături sunt explicitate și pot fi exploatate în parcurgerea făcută pe un calculator.

În jurul anului 1962, Douglas Engelbart a dezvoltat primul sistem hipertext, prezentat atunci drept o arhitectură conceptuală destinată creșterii potențialului intelectului uman ("Conceptual Framework for Augmenting Human Intellect") [Eng95]. Sistemul era destinat manipulării de concepte structurate într-o rețea în care arcele sunt relațiile între concepte.

Primul sistem declarat ca fiind hipertext a fost creat de Theodor Nelson în 1967 sub numele de "Xanadu". Nelson își propunea atunci să dezvolte un sistem, masiv paralel, destinat muncii creative și studiului. El a plecat în îndeplinirea acestei idei de la dorința de a găsi cea mai bună abstracție care să unifice literatura și arta cinematografică.

Sistemele hipertext (hipermedia) pe web permit accesul personalizat la volume imense de informații. În același timp, însă, ele suferă de problema aflului de informație cu care este bombardat un utilizator. O soluție este dezvoltarea de instrumente, aplicații, medii informatice pentru facilitarea accesului la cunoștințele dorite pe web. Aceste instrumente trebuie să faciliteze înțelegerea, abstractizarea textelor, extragerea informațiilor utile. Acesta este unul din motivele pentru care le-am denumit instrumente hermenofore. Trebuie remarcat faptul că ideea de a considera hipertextele ca instrumente de sprijinire a activităților cognitive a stat chiar la baza conceperii acestora, după cum am precizat în paragrafele anterioare.

În continuare, după o trecere în revistă a problematicii ontologiilor, în secțiunea următoare, se va introduce conceptul de mediu hermenofor, se va justifica necesitatea acestuia și se vor prezenta caracteristicile acestora. Lucrarea va fi încheiată cu o exemplificare printr-un sistem care are câteva trăsături ale unui mediu hermenofor și de o secțiune de concluzii.

2. Ontologii

Termenul de "ontologie" a fost, până nu de mult, folosit exclusiv în filosofie, pentru a denumi teoria asupra existenței, mai corect spus, asupra ceea ce consideră că există cel care întocmește teoria. Construirea multor sisteme filosofice pleacă de la o ontologie, adică de la definirea categoriilor fundamentale de entități din realitate și a relațiilor dintre ele. Chiar dacă ontologia nu este întotdeauna explicită, orice demers conceptual construiește o ontologie, chiar implicit, inconștient.

În ultimii ani, termenul de ontologie este folosit și în știința calculatoarelor. Cea mai frecventă extindere a folosirii acestui concept este în cadrul sistemelor de

inteligentă artificială bazate pe cunoștințe. Majoritatea programelor de calculator cu inteligență artificială prelucrează structuri de simboluri, care sunt menite a reprezenta conceptele, cunoștințele referitoare la domeniul considerat. Aceste structuri simbolice sunt grupate într-o așa numită bază de cunoștințe care constituie, de fapt, un model al domeniului respectiv. În ultimii ani se consideră că această bază de cunoștințe trebuie văzută ca o ontologie, o conceptualizare, o teorie asupra ceea ce există în domeniul avut în vedere. O ontologie este, din această perspectivă, o "specificare a unei conceptualizări ... Termenul este împrumutat din filosofie, unde însemna o considerare sistematică a existenței. În inteligența artificială se referă la precizarea a ceea ce se consideră că «există»» [Gru96].'

Între concepte pot exista diverse relații. Cea mai importantă relație este probabil cea hiperonimică [WN], taxonomică, între un concept și unul sau mai multe concepte mai generale, din care derivează, care îl subsumează, din a căror combinație a fost generat. Prin această relație se pot "moșteni" proprietăți de la conceptul (conceptele) mai general(e) la cel mai particular, dacă aceste proprietăți nu sunt redefinite la conceptul din urmă. Alte relații sunt cea meronimică [WN] ("parte-întreg"), între un concept și părțile sale sau cea antonimică, între două concepte (adjective) opuse.

O ontologie include, așadar:

- categoriile, conceptele fundamentale,
- proprietățile conceptelor,
- relațiile și distincțiile între concepte.

O ontologie este rezultatul unei experiențieri, a unor experiențe trăite, în care sunt evidențiate niște constante, niște regularități, care ne îndreptătesc să afirmăm că vor fi regăsite în viitor. În urma investigației făcută pentru a găsi esența regularităților se delimitează entități mentale denumite concepte sau categorii, care pot fi diferențiate de alte categorii. Aceste entități pot intra în combinație cu altele formând noi concepte.

Un aspect deosebit de important în ceea ce privește rolul ontologiilor este faptul că ele exprimă o comuniune, (co)existența unei diversități de concepte, o diferențe și relații între ele. O presuposiție este că există doar un număr limitat de concepte sau categorii, ceea ce înseamnă că se poate face un fel de cuantificare de discretizare a realității. Acestea constituie un punct de sprijin pentru achiziția de noi concepte sau pentru raționamentele făcute de om sau de calculator.

Partajarea unei ontologii este esențială în sistemele bazate pe agenți (programe) inteligenți pentru, de exemplu, comerțul electronic, pentru a le asigura autonomia, flexibilitatea și agilitatea. Ontologiile sunt liantul care integrează sisteme de baze de date, sisteme de obiecte, sisteme bazate pe cunoștințe, diverse aplicații integratoare și bazate pe colaborare. Ele reduc ambiguitățile

semantice în partajarea și reutilizarea cunoștințelor. "Scopul suprem este dezvoltarea de ontologii reutilizabile care pot fi aplicate pentru mai multe discipline". [OORG]

"O ontologie are drept prim scop facilitarea comunicării între calculatoare, independent de tehnologiile unui anumit sistem individual, arhitectura de prelucrare a informațiilor și domeniul aplicației. Ingredientii cheie care constituie o ontologie sunt un vocabular de termeni de bază și o specificare precisă a ceea ce înseamnă acești termeni." [OORG] O ontologie este însă mai mult decât un vocabular. Ea este punctul de plecare pentru dezvoltarea de structuri de cunoștințe, nu numai taxonomii sau clasificări de concepte ci și relații complexe. [OORG]

Din punct de vedere al programelor de calculator care folosesc ontologiile, există două tipuri de ontologii. Primul tip este cel al ontologiilor destinate sistemelor bazate pe cunoștințe, de exemplu, al unui sistem de diagnostic medical. Aceste ontologii sunt caracterizate de un număr relativ redus de concepte, dar legate între ele printr-un număr mare și variat de relații. Conceptele sunt grupate în scheme conceptuale complexe sau scenarii. Pentru fiecare concept pot exista una sau mai multe particularizări.

Spre deosebire de primul tip de ontologii, ontologiile lexicalizate includ un număr foarte mare de concepte, legate printr-un număr redus de tipuri de relații (de exemplu, hiperonimică, meronimică etc). Conceptele sunt reprezentate, de exemplu în WordNet [WN], prin mulțimi de cuvinte sinonime. Astfel de ontologii sunt folosite în sistemele de prelucrare a limbajului uman.

Correspondența ontologiei WordNet (care este concepută pentru limba engleză-americană) pentru limbile europene este EuroWordNet. Aceasta din urmă aduce avantajul că, fiind dezvoltată pentru mai multe limbi (engleză, franceză, germană, italiană, olandeză etc), permite și dezvoltarea de aplicații multilingve. În prezent, în cadrul Institutului de Cercetări pentru Inteligență Artificială al Academiei Române este în desfășurare, în colaborare cu mai multe țări din regiunea balcanică proiectul BalkanNet pentru integrarea în EuroWordNet a limbilor din zonă, inclusiv a limbii române.

3. Medii hermenofore

Denumim mediu hermenofor o colecție integrată de instrumente (pe care le vom numi hermenofore) și aplicații informatice direcționate către facilitarea unor activități de tip hermeneutic ale unui utilizator care explorează resurse aflate pe web. Termenul "hermenofor" [Tra01] poate fi parafrazat prin "generator de hermeneutică", pentru a sugera faptul că un mediu hermenofor facilitează activități hermeneutice, care acordă un rol important experiențierii și sunt orientate spre descoperirea unor înțelesuri, a unor structuri profunde, greu detectabile.

Elaborarea de medii hermenofore este absolut necesară în contextul actual al exploziei numărului și volumului de resurse și a interconexiunilor între acestea pe web. Sistemele hipertext (hipermedia) aduc noi dimensiuni cum ar fi interactivitatea, posibilitățile cu totul remarcabile de vizualizare, accesul personalizat la volume imense de informații. În același timp, însă, ele introduc și unele probleme datorate afluxului de informație, care poate duce la depășirea capacităților cognitive ale utilizatorului, la dezorientare și chiar la alienare. Este un fapt că utilizatorul, chiar profesionist în informatică, poate fi dezorientat în "labirintul" de pagini de web și resurse de tot felul (baze de date, documente, imagini, ontologii, lexicoane etc.) interconectate.

O soluție la problemele enumerate mai sus este dezvoltarea de instrumente, aplicații, medii informatice pentru facilitarea accesului la cunoștințele dorite pe web. Se poate spune, din această perspectivă, că browserele de web, "motoarele de căutare", agenții (asistenții) software sunt rudimente de medii hermenofore. Justificarea necesității considerării perspectivei hermenofore este lipsa abilităților hermeneutice ale acestor aplicații. Un exemplu tipic este faptul că "motoarele de căutare pe web" (de exemplu Google [Goo]) furnizează mii sau chiar zeci de mii de documente ca răspuns la o cerere. Alt exemplu este limita actuală a programelor de calculator în înțelegerea textelor cu scopul traducerii, sumarizării sau extragerii cunoștințelor. Aceste probleme sunt datorate, în primul rând, problemelor generate de ambiguitatea limbajului natural, a aspectelor legate de semantică, de pragmatică, de interpretare, de considerarea contextului, a metaforelor, a cunoștințelor de "bun simț". Toate aceste probleme sunt recunoscute ca fiind "nodul gordian" al aplicațiilor de inteligență artificială. După cum remarca Terry Winograd, programele de inteligență artificială nu pot depăși condiția unui birocrat, care nu poate să acționeze când nu are "reguli", care nu se implică [Win87]. Putem spune că, de fapt, problema este că acestor aplicații le lipsesc abilitățile hermeneutice. Ideea noastră este de a oferi un cadru în care puterea oferită de tehnologia informației să fie integrată cu capabilitățile specific umane.

Hermeneutica este, după opinia lui P. Ricoeur, o abordare complementară celei structuraliste în analiza limbajului, a înțelesului și simbolismului cultural. "Hermeneutica bazează înțelegerea textelor pe intențiile și istoria autorilor și relevanța acestor fapte pentru cititori. În contrast, filosofia analitică identifică de obicei înțelesul cu referenți externi pentru texte iar structuralismul găsiind înțelesul în aranjarea cuvintelor. Hermeneutica privește textele ca mijloace pentru a transmite experiența, crezurile și judecățile de la un subiect sau comunitate către alții. Astfel, determinarea înțelesurilor este o problemă de judecată practică și raționament de «bun simț» și nu privitor la o teorie a priori sau o demonstrație științifică." [MHD].

Hermeneutica este studiul interpretării, inițial ea referindu-se doar la interpretarea textelor [MHD]. În prezent s-a extins accepțiunea termenului hermeneutică, vorbindu-se de o poziție hermeneutică în filosofie, care include pe Heidegger, Gadamer, Habermas și Ricoeur, deosebită de formalişti (filosofia analitică, neo-pozitivism sau pozitivismul logic), reprezentați prin Descartes, Leibniz și Russell [Wes97]. Distanța între cele două abordări pleacă de la problema capturării înțelesului. Pe când formalişti pretind că pot reprezenta înțelesul, semantica, doar prin identificarea unui denotat în lumea reală corespunzător unei expresii formale, adepții hermeneuticii neagă această posibilitate, pentru ei înțelesul implicând și considerarea experienței, a credințelor subiectului. Se poate spune că, dintr-un punct de vedere se ajunge la aceeași dispută dintre Husserl și Heidegger sau dintre Dennett și Chalmers.

Mediile hermenofore furnizează informațiile dorite dintr-o perspectivă particulară, pentru un anumit utilizator, considerând un anumit domeniu și într-un anumit moment dat. Un mediu hermenofor trebuie conceput deci în scopul personalizării interfațării la resursele web-ului, pentru a facilita înțelegerea. Dacă prezentările făcute într-un mediu hermenofor sunt structurate ca hipermedia, una din preocupările principale ce trebuie avute în vedere este faptul că utilizatorul trebuie să experimenteze parcurgerea unei secvențe de pagini de web, secvență care trebuie să respecte niște reguli de pragmatică.

În plus față de furnizarea unei interfețe adaptabile, o altă caracteristică a unui mediu hermenofor trebuie să fie facilitarea inițiativei utilizatorului. El trebuie să poată experimenta, să poată investiga resursele web-ului. Instrumentele hermenofore sunt destinate sprijinirii activității hermeneutice umane adică a unei atitudini direcționate către înțelegerea unor cunoștințe sau structuri ascunse în texte (hipertextelor, hipermedia). Un rol important în procesul înțelegerii îl au modalitățile de a genera experiențieri, adică experiențe de trăire, fapte de viață (conform teoriei că înțelegerea necesită un proces empatic [Wri95], [Mar97]). Unul dintre cele mai uzitate mijloace de acest gen este folosirea metaforelor [LaJ80], [Tra00]. În acest sens se înscrie preocuparea de a dezvolta instrumente (hermenofore) pentru detectarea, adnotarea și prelucrarea metaforelor.

O caracteristică pe care o considerăm esențială la un mediu hermenofor, în contextul precizat mai sus, este și posibilitatea de vizualizare multiplă, din perspective diferite, a aceluiași document. Enumerăm aici, drept exemplu, în afara perspectivei conținutului "brut" al unui document, alte perspective, date de concordanțe, adnotări (cu părți de vorbire, de exemplu), extrase, rezumate, arbori de analiză semantică, structuri care reprezintă conținutul semantic. Remarcăm, în acest context, rolul extraordinar de important al adnotărilor documentelor în limbajul extrem de versatil care este XML [XML].

Vom considera că instrumentele hermenofore au ca scop revelarea și valorizarea unor cunoștințe sau a unor structuri încorporate în volumele imense de

> hipertexte și hipermedia de pe web. Datorită faptului că abordarea ip pune pe prim plan rolul experiențierii umane, un instrument hermenofor poate neapărat considerat în relație cu utilizatorul care îl folosește. De aceea, trebuie să aibă asociat modelul utilizatorului, care să conțină cel puțin următoarele informații despre utilizator:

- ontologia sa,
- scopurile urmărite,
- profilul psihologic,
- istoricul acțiunilor efectuate,
- preferințele sale (explicite sau implicite, derivate din comportamentul său).

Pe de altă parte, instrumentele hermenofore trebuie să ia în considerare aspectele legate de particularitățile autorilor documentelor:

- ontologiile considerate (de exemplu, ontologiile impuse de practicile mele sau de practicile domeniilor considerate),
- scopurile presupuse,
- elemente de istoric,
- aspecte psihologice general umane.

Instrumentele hermenofore pot fi împărțite în mai multe clase în funcție de acțiunile efectuate:

- căutare a documentelor relevante,
- categorizare a documentelor conform unei taxonomii predefinite,
- relevare de regularități (de exemplu, cologații) sau de structuri din documente,
- segmentarea textelor,
- extragere de informații sau cunoștințe din documente,
- sumarizare,
- relevare de structuri pe web [WSD97],
- instrumente de adnotare (la nivel sintactic, semantic sau fonetic) ale documentelor.

Spre deosebire de instrumentele de minerit al textelor, instrumentele hermenofore pun, în plus, accentul pe aspectele legate de interacțiunea, de experiența utilizatorului.

În secțiunea următoare se va prezenta sistemul GenWeb, un sistem asistat de învățare terminologiei financiare într-o limbă străină [TM].

care a fost dezvoltat ca un modul într-un proiect mai mare, denumit "Larflast" și finanțat de Comunitatea Europeană. GenWeb a implementat instrumente hermenofore care identifică și utilizează metafore pentru a facilita înțelegerea unui anumit concept [TraOO]. În acest scop, el caută metafore în texte considerate relevante. Metaforele sunt identificate printre perechile de cuvinte care corespund la concepte din ontologia domeniului considerat (finanțe) și din ontologia metaforelor, aceasta din urmă reflectând aspecte psihologice general umane [LaJ80]. Trecerea de la un concept la o mulțime de cuvinte (sinonime sau înrudite) se face pe baza ontologiei WordNet, derivată din investigații psiholingvistice [WN]. Metaforele sunt adnotate în XML [XML], unul din atributele folosite în adnotare fiind scopul urmărit de autor [TraOO].

Tot în GenWeb, textele adnotate cu metafore sunt folosite ulterior pentru a genera structuri (bazate pe principii retorice) de pagini de web personalizate conform modelului utilizatorului. Aceste structuri se constituie într-un sit în care cel care învață poate experiența. Tot pe post de instrumente hermenofore, în GenWeb este disponibilă vizualizarea de concordanțe în context.

4. Sistem de instruire asistată cu calculatorul în înțelegerea unor termeni financiari

Există mai multe puncte de vedere asupra modului în care are loc un proces de învățare. Suntem de partea abordării constructiviste [BIM96, Wil96] în conceperea proceselor educaționale. Această abordare consideră că fiecare dintre noi ne construim propria realitate, propriul bagaj de cunoștințe, plecând de la experiențele pe care le-am avut [ErK97]. După cum remarcă [BIM96], "Nucleul studiului este activitatea hermeneutică a construcției de interpretări." Învățarea poate fi și ea văzută constructivist ca un proces hermeneutic, de înțelegere, de transpunere în domeniul studiat, de experimentare, de trăire.

Plecând de la ideile învățării constructiviste se ajunge la următoarele principii [ErK97]:

- Învățarea este un proces activ în care studenții experimentează, caută să înțeleagă singuri ceea ce învață, profesorul fiind mai mult un îndrumător;
- Învățarea trebuie să fie un proces auto-reglat de către studenți;
- Învățarea constructivă este un proces situațional în sensul că studentul trebuie introdus într-un mediu de învățare care îi permite să experimenteze, în care se pot face simulări;
- Învățarea trebuie să fie socială, trebuie să existe o permanentă colaborare a studentului cu colegii lui.

Dintr-o altă perspectivă, învățarea poate fi considerată ca un proces de inducere de modele mentale adecvate [JoL83]. Înțelegerea poate fi văzută astfel ca momentul în care realitatea supusă comprehensiunii este pusă în corespondență cu un model mental complet și valid. Empatia [Mar 97], identificarea eu-lui cu starea de lucruri considerată poate fi, în acest caz, tocmai sentimentul de "trăire" în lumea modelului mental.

O practică deja răspândită este de a dezvolta sisteme inteligente de asistare cu calculatorul a instruirii ("Intelligent Tutoring Systems") care încearcă să monitorizeze procesul de învățare prin verificarea asimilării conceptelor din ontologia domeniului considerat [Tra95]. Se consideră că un model adecvat al cunoștințelor elevului poate fi construit prin raportare la această ontologie. De fapt, această metodă este folosită și în învățământul tradițional: noii termeni sunt introduși prin genul proxim și diferența specifică. În termenii ontologiilor, noii termeni sunt definiți prin superconceptele care-i subsumează și prin particularitățile care-i diferențiază.

Orice profesor știe însă că astfel de definiții sunt necesare dar nu sunt suficiente. Pentru a aprofunda termenii definiți sunt necesare exemple, imagini cu un grad mai mare sau mai mic de iconicitate, plecând de la poze și schițe, diagrame și grafice, până la imagini sugerate, până la metafore. Acest fapt este prezent nu numai în învățământ, el apare în orice proces de comunicare (învățământul fiind, bineînțeles, și el inclus).

În cele ce urmează nu ne vom referi la utilizarea imaginilor propriu-zise, care facilitează evident învățarea sau comunicarea. Vom considera un caz particular de imagini, mentale, sugerate, semne iconice lipsite de caracterul vizual dar care comunică o experiențiere (de multe ori chiar mai puternic, printr-un efect care ar putea face să ne gândim la percepția subliminală). Este cazul metaforelor, care sunt folosite într-o proporție de cele mai multe ori nebanuit de mare în comunicarea inter-umană.

Pentru a ilustra puterea de expresie a metaforelor și, bineînțeles, rolul lor în înțelegerea unor termeni, am să exemplific prin metafora "acțiunile la bursă sunt niște creaturi foarte sensibile" (găsită într-un text pe situl de web al Bursei din New York - <http://www.nyse.com>). Nu este nevoie să ne imaginăm o anumită creatură concretă pentru a înțelege ce sugerează metafora exemplificată. Succesul unei metafore, puterea ei expresivă, capacitatea de comunicare sunt date de măsura în care "rezonăm" la mesajul transmis. Ori ce este mai percutant pentru un om decât faptul că suntem creaturi extrem de sensibile? Prin urmare, succesul metaforei folosită într-un context foarte pragmatic, al discursului unui specialist în finanțe este determinat de inspirația vorbitorului de a se referi la un fapt general uman. Nici o definiție de tip gen proxim-diferență specifică nu poate comunica experiența referitoare la aspectul foarte fragil al acțiunilor la bursă precum o face metafora de mai sus.

Rolul covârșitor al metaforelor în viața noastră a fost remarcat și de Lucian Blaga ("omul este un animal metaforic" [Bla85]) și a fost foarte bine evidențiat de Lakoff și Johnson într-o lucrare cu un puternic impact ("Metaforele cu care trăim" - "Metaphors we live by" [LaJ80]). Cei doi autori americani consideră că "subcategorizarea și metaforele sunt două extremități ale unei continuum", că metaforele "formează sisteme coerente în care ne conceptualizăm experiența" [LaJ80]. Putem spune că metaforele oferă alte mijloace expresive decât cele de categorizare oferite de ontologii. Ele nu țin de logica lui Ares, care categorizează, ci de logica lui Hermes, propusă de Noica [Noi86].

Dintr-o altă perspectivă, metaforele pot fi considerate instrumente empaticе, care determină imersiunea cititorului (receptorului) în lumea experiențelor autorului. Acest fapt era evidențiat și de Lakoff și Johnson: "Esența metaforei este înțelegerea și experiențierea unui lucru prin altul" [LaJ80]. De exemplu, metafora amintită mai sus despre acțiunile la bursă ne comunică o informație pe care orice ființă vie o înțelege (sensibilitatea, perisabilitatea) dar care nu poate fi exprimată în categorizări.

Importanța metaforelor a fost revelată și de studiul preliminar făcut în cadrul proiectului Larflast (care a avut drept scop elaborarea unui sistem de asistare cu calculatorul a învățării terminologiei financiare într-o limbă străină [Lar], [TMC02], [ABK02]) de o profesoară de limba engleză la o facultate economică din Sofia. Dânsa remarcă ca o importantă dificultate "înțelegerea metaforelor. Limbajul economic și financiar este extrem de metaforic și, uneori, grupuri de metafore apar în imagini complexe. Deseori cuvinte uzuale sunt folosite în metafore elaborate, ... cum ar fi «a susține o pierdere»" [Vit99].

Proiectul Larflast a inclus mai multe module tipice pentru sisteme inteligente de instruire, cum ar fi o ontologie, un mecanism de inferență, teste (grilă) pentru diagnosticarea cunoștințelor elevului și actualizarea modelului acestuia. Sistemul dezvoltat include cinci servere de web, unul la București și altele la Leeds, Manchester, Montpellier și Sofia. Serverul de la București, după ce este lansat, accesează serverul de la Sofia pentru a prelua modelul elevului (ce concepte știe și ce concepte nu) și apoi generează pagini de web personalizate.

Metaforele sunt identificate în texte considerate relevante care au fost obținute în urma căutării cu o mașină de căutare uzuală (de exemplu, Google [Goo]). Textele găsite sunt grupate într-un corpus care este adnotat cu metaforele identificate. Acest corpus, împreună cu ontologia domeniului și cu modelul studentului (construit pe baza răspunsurilor date de student la teste) sunt folosite pentru generarea personalizată de pagini de web. În figura următoare este ilustrată arhitectura sistemului GenWeb.

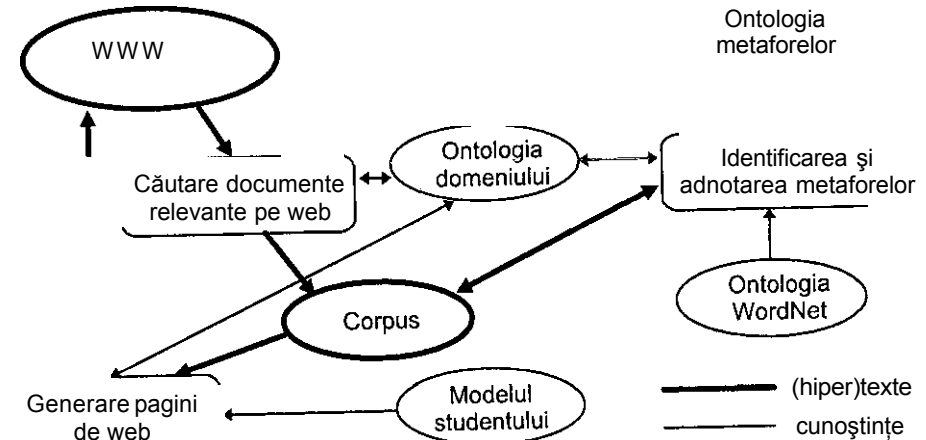


Figura 1

Pentru identificarea și adnotarea metaforelor a fost implementat un editor semantic specializat (fig.2) și un editor de concepte (fig.3).

File Help

Text:
 <article nr="11" type="educational" URL="http://www.nyse.com/about/educacion/xnve.r/i7214.har" >
 <text>
 Oaetaph what="stock" how="creatures" .hy="sensitivxty" >Stc
 Stocks react to all kinds of xnfluences, large and small,
 and oaetaph .hat="stock" hotj="organism" .hy="reactivxty >
 sensitive reactxons register as price changes</metaph>.

Metaphor List:
 their sensitive reactions register as price changes
 News events can trigger a change in stock prices when

what stock
 how: creature Apply
 why: |sensitivity

Tag: |metaph
 Attribute1: |what Distance: |5
 Attribute2: |how
 Attribute3: |why

Source Domain:
 organism (noun) = life form, organisn_±.
 act (verb) = act. move -- (perform an a
 building (noun) = building, edifice o
 pillar (noun) = pillar, mainstay -- (a pr

Destination Domain:
 stock (noun) = stock - (the capital rai:
 market(noun) = market. securities im
 futures (noun) = future. future tense
 credit (noun) = credit -- (money availa
 1

Load File. Save File. Change. Load File.. Save File..

Figura 2

Concept List: Attribute List: Senses List:

Jorganism Add Inoun " 3 Change

Modify Remove

2. organism -- (a system considered analogous in structure or

act Modify

building sensitivity

pillar reactivity

move vulnerability

WordNet Sense Number:

Antonym	iv	Hypernym	Participle of	life	d	mortal	^1
Attribute		RyponyrTJ	See also	biont		human	
Cause		Member Holonym	Similar to	lperson		soul	
Derived	T	Member Meronym	Substance Holonym	individual		animal	
Entailed by	F	Part Holonym	Substance Meronym	someone		beast	
Entailment	F	Part Meronym	Verb group	somebody		brute	
				mortal		creature	
				human		faun>	

Figura 3

Modelul studentului este creat pe baza răspunsurilor la teste:

Address | \$ | http://w _larflast.bas.bg/cgi-bin/gete.exe/passwd

An institutional mechsasm created by society to channel savings and other financial services to those individuals and inshtutions willing to pay for them

Financial Market >|

= Institutional mechsasm set up by society to trade or exchange loans and securities that have already been issued

Expenditures on capital goods or inventones of goods or raw materials that are used to produce other goods and serwces. causing future production and income to use

Options contracts
Investment
Credit

Institutional mechsasm set up by society to make loans and trade securities where the terms of trade are set by direct bargairning between a lender and a borrower

Secondary Market
Open Market
Negotiated Market
Primary Markets
Financial Market
Futures Contract

Contracts that call for the future delivery or sale of designated securities at a price agreed upon the day the contract is made and that are used mainly to hedge (protect) against changing intere st rates

A loan of funds in return for a promise of future payment.

Institutional mechsasm created by society to make loans and trade securities in which any individual or inshtitunon can participate

Agreements between contract wnters and contract buyers to accept delivery of ("cail") securities or place with buyers ("pu") securities at a specified price on or before the date the contract expires

l Institutional mechsasm set up by society to trade newly issued loans and securities.

< Markets where temporary surpluses of cash are channeled into temporary loans of funds. one year or less to maturity.

Diagnostics

Trausan, you havc correctly answered to some questions about: financial_market, secondary_market, futures_contract, option_contract, primary_market, investment, credit, but it seems that you still do n correctly know the following concept(s):

1. Credit
2. Futures_contract
3. Investment
4. Primary_market
5. Option_contract
6. Monev_market
7. Open_market
8. Negotiated_market

Picase browse the web pages describing these concept(s).

Only the wrongly known and unknown concepts are detailed presented!

Back to n
trausan(S:valhalla racai.ro

About LarFlast>Please send questions and remarks at

Sfstartj & [C:\user\shf\sysn] jG\WNNI\Syco?rch j_p)Mrocf PvalRH • Q_ (fj)hp^/www.larflast/Ar~ tE Internet 319 IM

Figura 4

<http://www.farfast.bas.bg/col-bin/cole>

Go File

credit extended to each customer. In September 1999, credit extended to each customer were established. provide foreign currency guarantees, under cooperative, and bill finance companies. T growth or reliance on external funding, but allocation ceilings imposed on these banks, ayatara. It remains a cash-based society, eve base on the central plan and extend loans t credit, a bank-dominated financial ayatem and weak credit management ptocesses, and greater avenuea f credit management Harden budget constraint policy still implemented through a Plan. Establishroelnt of special economic zone agency. SEC granted license to nine Soreign allocation ceilings treplaced with standard

Hi II 1

S Dona

"Financial market"

.....]?;J C'tusanWwMcclastSatyrn/S TY18_HT

Diagnostics

Definition:

Joh* you have correctly answered to some questions about: money market, primary market but it seems that you still do not correctly W the following concept(L):

Some «bšfc.rfcalpfrases in wkick txis ««cept appears.

- theperfoima»»ofafmancial market
- nncial markets that are continuat to glow
- walysts expect the impact on the f W i a i market to be negativ»
- ^ I f P ^ ^ y fo'tis fivncialmarkets» to «topt cha^s that will hek
- mll affect their economy and financial markets
- fmaicial market has undergone substantial developme,
- a robust financial market
- an open financial market

Please browse the web pages describing these concept(s).

Releva*texts &r tkis concept are:

Only the wrongly known and unknown concepts are detailkd pnsented*

My Computer

frfr^

Prinnr- n 11181

Secondary market

feunese into the feinw^ to^

-OwnlitrMa-

Some &cts about secoitdarymarket are :

- Secondary market trades already issued bonds
- Change interest rates in the secondary market
- Secondary market supports asvz investments

Some similar concepts with secondary_market are:

- Money_ma^et

Secowlary_market is the opposito ofprimary,!!|^!^:

Concluzii

În contextul dezvoltării explozive a numărului de docum absolut necesară existența unor medii care să permită utilizato scopul extragerii cunoștințelor din texte și structuri de docume activitate trebuie sprijinită de ontologii, un rol foarte important a ontologiilor de mari dimensiuni existente astăzi pe web. Un integrează instrumente hermenofore cu ontologii într-o ar utilizatorul trebuie să poată experimenta, să investigheze divers textelor. Se poate spune că un mediu hermenofor înglobează și de prelucrare a cunoștințelor cu instrumente de prelucrare a te specifice web.

Bibliografie

- [ABK02] G. Angelova, S. Boytcheva, O. Kalaydjiev, Șt. Trăușan-Strupchanska, Adaptivity in a web-based CALL system, Harmelen (ed.): ECAI 2002. Proceedings of the 15th European Conference on Artificial Intelligence, IOS Press, Amsterdam, 2002, pp. 1-10.
- [Bla85] L. Blaga, Trilogia culturii, Ed. Minerva, 1985
- [BIM96] Black, J.B., McClintock, An Interpretation Constructivist Design, in B.G. Wilson (ed.), Constructivist Environments: Case Studies in Instructional Design, Educational Technology Publications, 1996.
- [Cli95] Clibbon, K., Conceptually Adapted Hypertext For Learning, CHI'95, <http://www.acm.org/sigchi/chi95/Electronic/Documentation>
- [CTr01] Constandache, G.G., St. Trăușan-Matu, Ontologia calculatoarelor, Editura Tehnică, 2001.
- [Cul94] Culianu, I.P., Eros și magie în Renaștere; 1484, Nemira, 1994
- [Eng95] Engelbart, D.G., Toward Augmenting the Human Intelligence, our Collective IQ, CACM No.8, Vol.38, Aug. 95, pp. 30-39
- [ErK97] Ertl, B., Kraan, A.G., Internet-Based Learning Environment: A Constructivist point of view, Proceedings of RILW, Ilieni, 1997
- [Goo] <http://www.google.com>
- [Gru96] Gruber, T., What is an Ontology, <http://www.kr.org/top/d>
- [JoL83] Johnson-Laird, P.N., Mental Models - Towards a Cognitive Language, Inference, and Consciousness, Cambridge University Press, 1983
- [LaJ80] Lakoff.G., Johnson, M., Metaphors We Live by, The University of Chicago Press, 1980.

- [Lar] LarFLaST, <http://www-it.fmi.uni-sofia.bg/larflast/>
- [Mar97] Marcus, S., Empatie și personalitate, Ed. Atos, 1997.
- [MHD] J.C. Mallery, R. Hurwitz, G. Duffy, Hermeneutics, Encyclopedia of Artificial Intelligence, pp. 596-611.
- [Noi86] C. Noica, Scrisori despre logica lui Hermes, Ed. Cartea Românească, 1986.
- [OORG] <http://www.ontology.org/main/papers/faq.html> \
- [Rad 91] Hypertext from Text to Expertext, McGraw Hill, 1991.
- [Sow99] J. Sowa, Knowledge Representation: Logical, Philosophical and Computational Foundations, Brooke Gole Publishing Co., Pacific Grove, CA, 1999, vezi și [CTrOI].
- [THH95] Thiring, M., Hannemann, J., Haake, J.M., Hypermedia and Cognition: Designing for Comprehension, Communications of the ACM, voi.38, no. 8, pp. 57-66, aug. 1995.
- [TMC02] Șt. Trăușan-Matu, D. Maraschi, S. Cerri, Ontology-Centered Personalized Presentation of Knowledge Extracted From the Web, în S. Cerri, G.Gouarderes (eds.), Intelligent Tutoring Systems 2002, Springer, Lecture Notes in Computer Science number 2363, pp. 259-269.
- [Tra95] Șt. Trăușan-Matu, Programe inteligente pentru asistarea învățării, în Revista Română de Informatică și Automatică, voi.5, nr.4, 1995, pag. 7-16.
- [TraOO] Șt. Trăușan-Matu, *Metaphor Processing for Learning Terminology on the Web*, in S.A.Cerri (ed.), Artificial Intelligence, Methodology, Systems, Applications 2000, Springer-Verlag, ISBN 3-540-41044-9, 2000, pp.232-241
- [Tra01] Șt. Trăușan-Matu, Interfatarea evoluată om-calculator, Ed. MatrixRom, 2001.
- [Vit 99] I. Vitanova, English for Finance, Understanding Money and Markets, <http://wwwjit.fmi.uni-Sofia.bg/larflast/>
- [Wes97] D.West, Hermeneutic Computer Science, CACM, Vol.40, No.4, pp. 115-116, 1997, și în [CTrOI].
- [Wil96] B.G. Wilson (ed.), Constructivist Learning Environments: Case Studies in Instructional Design, Education Technology Publications, 1996
- [Win87] T. Windgrad, Thinking machines: Can there be? Are we?, Report No. STAN-CS-87-1161, Stanford, 1987.
- [WN] WordNet, <http://www.cogsci.princeton.edu/~wn/>
- [Wri95] von Wright, G.H., Explicație și înțelegere, Humanitas, 1995.
- [WSD97] <http://www.research.att.com/~suciu/workshop-papers.html>
- [XML] www.w3.org/xml

SECȚIUNEA III

TEHNOLOGII ALE LIMBAJULUI

Experimente în vederea recunoașterii vorbitorului

Comeliu BURILEANU,
Universitatea „Politehnica” din București, Spl.Independenței 1
cburileanu@messnet.pub.ro

Luigi BOJAN,

Graphco Technologies Inc., Newton, PA, USA

1. Introducere

Având în vedere funcția realizată și concomitent, sarcinile tehnologia vorbirii se poate clasifica în mai multe domenii [1, 2]:

- Recunoașterea automată a vorbirii. Se bazează pe analiza semnalului vocal și are în vedere informația transmisă care "îl ascultă". Din această informație, mașina este capabilă să extragă acele caracteristici ce îi vor permite să recunoască vorbitorul, în ce fel și în ce condiții.
- Sinteza automată a vorbirii. Se realizează răspunsul mașinilor către operatorul uman.
- Codificare/decodificare (analiză și sinteză) a vorbirii. Tehnici de compresie a informației conținută în semnal și în vederea unor prelucrări ulterioare specifice sarcinii de

Un domeniu interdisciplinar important, legat în mod esențial de recunoaștere și sinteză automată ale vorbirii este cel al dialogului

Termenul "comunicare om - mașină" pare forțat: mașina nu este socială, nu are nici scop nici cultură. Ea nu poate acționa în lumea omului și de a putea să răspundă corect la întrebări de genul: "ai putea să-mi spuneți ce înțeleg?". Ea nu este "conștientă" decât de propria sa "lume". Avem în vedere de a comunica cu mașinile? Au importanță intențiile lor, chiar dacă mașina poate să-mi comunice sau să mă facă să știu o mașină?

Mașina îmi procură "uneltele" pentru a realiza o sarcină, proiectez noi obiecte (eventual, virtuale), ea mă aduce într-un univers virtual și permite să utilizez un mediu de programare împreună cu alții și să lucrez pentru a lucra într-o manieră cooperantă în același mediu informatic. Mașina prezintă deci ca un *factor de interacțiune*. Ea trebuie să-mi furnizeze

muncă, unelte și metode. Dar pentru aceasta, mașina trebuie adaptată sarcinii curente sau unor sarcini noi, să adopte un comportament "comprehensibil", să se arate "prietenoasă" etc. Paradoxul este deci evident: mașina trebuie să fie, dintr-un anumit punct de vedere, *sociabilă* pentru a colabora eficient cu un utilizator în scopul îndeplinirii sarcinilor, din ce în ce mai complexe, care îi sunt încredințate.

Preocupările noastre în domeniul tehnologiei vorbirii au, între altele, scopul de a oferi mijloacele pentru o comunicare între om și mașină prin mesaje vorbite [3]. Această comunicare este doar un aspect al dialogului. Rămâne în continuare deschisă problema definirii conceptelor și cea a stabilirii unor strategii de dialog adecvate sarcinii de rezolvat.

Semnalul vocal conține o varietate de informații utile: ce se vorbește, cine vorbește, în ce fel și în ce condiții. În cadrul recunoașterii se pune problema identificării unui anumit tip de informații; de pildă, recunoașterea cuvintelor rostite înseamnă determinarea mesajului (ce se vorbește) indiferent (sau ajutându-se) de variabilitățile introduse de vorbitor (cine), maniera de a vorbi (în ce fel) și zgomotul ambiental (în ce condiții). Putem particulariza afirmând că **recunoașterea vorbirii** este procesul de transformare a semnalului acustic continuu produs de organul fonator uman într-o reprezentare discretă căreia i se poate atașa o semnificație și care, când e înțeleasă, poate fi folosită pentru a determina un răspuns.

Problemele majore pe care le ridică recunoașterea automată sunt legate de

- discretizarea semnalului vocal care, din punctul nostru de vedere înseamnă *segmentare*;
- caracterul adecvat al răspunsului ce depinde de natura sarcinii de îndeplinit; modalitatea de prelucrare este irelevantă.

Proiectarea unui sistem de recunoaștere presupune câteva opțiuni fundamentale de abordare. Punctul de vedere adoptat poate viza prelucrarea unui semnal acustic ca oricare altul, poate ține seama de mecanismul producerii vorbirii, poate simula recepția senzorială, sau poate folosi modelul uman al percepției vorbirii.

Termenul de **recunoaștere a vorbitorului** desemnează orice aplicație de discriminare a persoanelor pe baza vocii acestora. Procedurile de recunoaștere se desfășoară în două etape [4]:

- etapa de antrenare: colectarea de material vocal de la persoana care se dorește a fi recunoscută;
- etapa de testare: compararea unui fragment de vorbire neidentificat cu datele provenite din antrenare și luarea deciziei de recunoaștere.

Există două subclase de aplicații:

• **verificarea vorbitorului** își propune să determine dacă un fragment de semnal vocal aparține sau nu unui anumit vorbitor [5, 6, 7, 8]. Există doi parametri care caracterizează performanțele sistemului: respingerea adevăratului vorbitor și

acceptarea unui impostor. Considerând un set de N vorbitori, informația (în biți) obținută este

presupunând probabilitatea de verificare a *priori* egală cu 0.5;

• **identificarea vorbitorului** are ca scop punerea în corespondență a unei voci necunoscute cu un vorbitor dintr-un set dat [9,10,11,12]. Pentru N vorbitori, informația (în biți) obținută este

considerând probabilitatea de identificare a *priori* egală pentru toți vorbitorii.

Rezultă că, potențial, un sistem automat de verificarea vorbitorului are performanțe mai bune.

O clasificare suplimentară a automatelor de recunoaștere are în vedere natura sarcinii de îndeplinit și se reflectă în complexitatea sistemului [13]:

- sisteme de recunoașterea vorbitorului **dependente de text** - textul utilizat în faza de antrenare este același cu cel de testare;
- sisteme **independente de text** - indiferent de materialul vocal avut la dispoziție.

Setul de vorbitori vizat poate impune, de asemenea, o clasificare a automatelor:

- "*set închis*" - pentru procesul de identificare descris ca mai sus;
- "*set deschis*" - în cazul identificării există posibilitatea ca vocea necunoscută să nu aparțină niciunui dintre vorbitorii din setul dat, numărul de decizii posibile fiind în acest caz $N + 1$. Identificarea pe "set deschis" devine astfel o combinație a proceselor de verificare și identificare.

2. Reprezentarea parametrică

Variabilitățile pronunțării pentru diverși vorbitori, sau la un același vorbitor, la momente de timp diferite, constituie una dintre dificultățile majore ale sarcinii de recunoaștere a vorbitorului. Deosebiri de vorbire depind de dialect, context, stil de exprimare, stare emoțională etc. Mai mult, în opinia noastră, așa cum vom încerca să argumentăm mai departe, *limba în care se vorbește* impune deosebiri de abordare și diferențe ale performanțelor automatului [14].

Din acest motiv, alegerea judicioasă a **caracteristicilor acustice** care vor fi utilizate în procesul de recunoaștere este deosebit de importantă:

- să diferențieze vorbitori diferiți dar să fie tolerante pentru același vorbitor;
- să fie ușor măsurabile din semnalul vocal;
- să fie stabile în timp;
- să nu fie susceptibile de a fi contrafăcute de potențiali impostori.

Având în vedere cerințele formulate mai sus, am decis utilizarea parametrilor cepstrali.

Anumite abordări ale prelucrării semnalului vocal presupun adoptarea unor decizii fundamentale de dezvoltare a analizei: considerarea unui model de producere a vorbirii având ca prototip aparatul fonator uman, separarea efectelor sursei vorbirii de comportarea tractului vocal propriu-zis, o serie de aproximări care să facă analiza eficientă în condiții normale de procesare [15]. Variația (lentă) în timp a formei tractului vocal este aproximată printr-o serie de secvențe de durată suficient de mică pentru a presupune forma invariantă: este ceea ce se numește "analiza în timp scurt". Dacă, în plus, în aceste durate "scurte" de timp se presupune că tractul este caracterizat în mod esențial de frecvențele sale de rezonanță, se ajunge la un model al cărui parametri se pot deduce prin rezolvarea unui sistem de ecuații liniare. Deși aproximările avute în vedere par destul de restrictive, analiza prin predicție liniară (LPC) dă rezultate deosebite pentru că semnalul vocal are o redundanță deosebită; este motivul pentru care metoda ne permite să aproximăm un eșantion de semnal printr-o combinație liniară (deci este liniar predictibil) dintr-un număr de eșantioane precedente. Desigur, principiile în sine ale metodei nu sunt noi; ele au permis însă, în decursul ultimilor ani, evoluția spre metode mai sofisticate [16,17].

Nici principiile analizei cepstrale (analiză care, așa cum vom arăta, se poate baza pe rezultatele analizei LPC) nu sunt noi: se dezvoltă un mecanism care să permită decelarea mai amănunțită a influențelor diverselor elemente ale organului fonator. O serie de presupuneri fundamentale de abordare se păstrează (modelarea producerii vorbirii în maniera aparatului fonator uman, analiza "în timp scurt"); dar separarea efectelor excitației glotale, tractului vocal și radiației buzelor poate fi făcută într-o modalitate care ține seama mai detaliat de fiecare efect în parte [18, 19].

În concluzie, presupunerile fundamentale care stau la baza parametrizării propuse sunt:

- efectele excitației tractului vocal și ale tractului propriu-zis pot fi separate;
- tractul vocal este invariant pe durate scurte de timp, ceea ce are drept rezultat obținerea unui model descris de un sistem liniar al cărui parametri variază lent în timp (constanți "în timp scurt").

Fundamental pentru modul în care concepem abordarea analizei semnalului este asimilarea analizei cu parametrizarea semnalului și, în consecință,

cu compresia sa. Alegerea parametrilor a avut în vedere și considerente pragmatice:

- complexitatea prelucrării;
- gradul de compresie,
- tipul de aplicație,
- în ce măsură parametrii sunt semnificativi și robuști.

O primă variantă a schemei bloc care descrie funcționarea sistemului de recunoașterea vorbitorului este prezentată în fig. 1. Blocul de preprocesare presupune filtrarea și achiziția semnalului în condiții normale pentru orice sistem de recunoaștere. În această secțiune vom descrie obținerea cepstrului pornind de la analiza LPC, iar în secțiunea următoare vom descrie principiile cuantizării vectoriale și deci procedura de recunoaștere propriu-zisă.

Fie semnalul vocal presupus a fi convoluția unei excitații și a funcției de transfer a tractului vocal:

$$s(t) = e(t) * v(t) \quad (3)$$

Analiza homomorfică care duce la obținerea cepstrului presupune aplicarea unui operator neliniar "H"

$$s(n) \xrightarrow{H} \hat{s}(n) \quad (4)$$

în care $\hat{s}(n)$ va fi numit cepstrul complex asociat semnalului $s(n)$.

Prin definiție

$$\hat{s}(n) = \sum_{k=-\infty}^{\infty} s(k) e^{-jkn} \quad (5)$$

"

Astfel, cepstrul complex asociat semnalului devine

$$\hat{S}(n) = \hat{E}(n) + \hat{V}(n) \quad (6)$$

ceea ce permite separarea componentelor printr-o "filtrare temporală" aplicată cepstreilor

$$\hat{S}(n) \xrightarrow{H^{-1}} \hat{s}(n) \quad \begin{matrix} \xrightarrow{H^{-1}} \hat{E}(n) \\ \xrightarrow{H^{-1}} \hat{V}(n) \end{matrix} \quad (7)$$

Obținerea parametrilor cepstrali se poate realiza ținând seama de câteva proprietăți ale cepstrului.

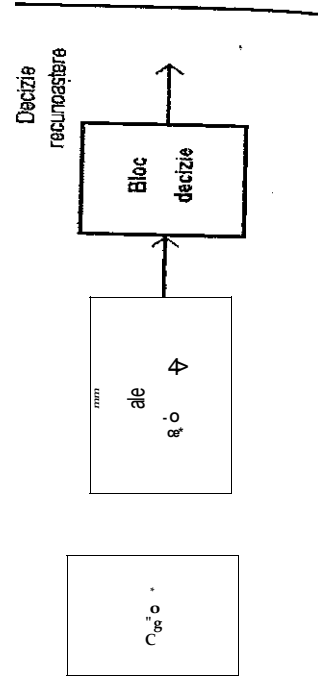
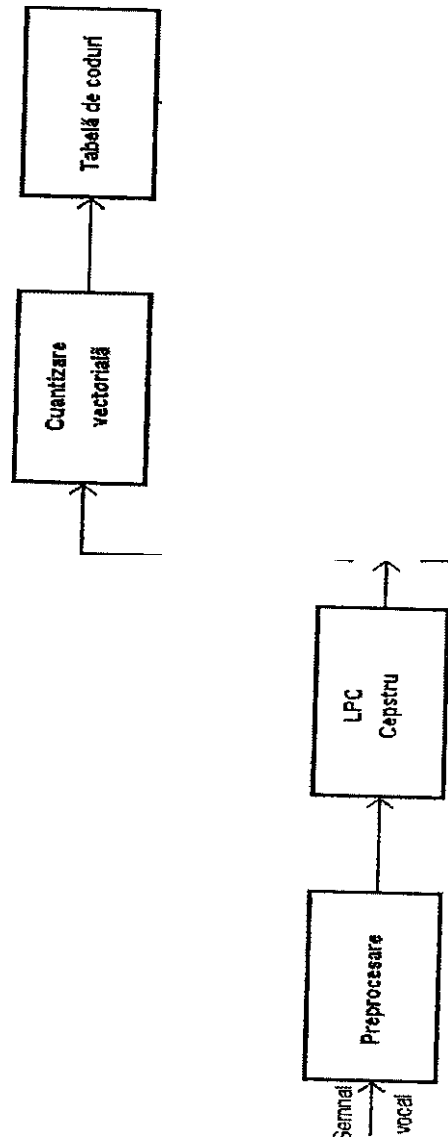


Figura 1. Un sistem de recunoaștere a vorbitorului - schema de principiu

Fie $c(n)$ partea pară a cepstrului complex al semnalului

$$c(n) = [s(n) + s^*(-n)] / 2$$

Secvența $c(n)$ se numește *cepstrul real* al semnalului

este o secvență cauzală - ca și $s(n)$; rezultă

$$s(n) = c(n) - \begin{cases} 0 & \text{pentru } n < 0 \\ 1 & \text{pentru } n = 0 \\ 2 & \text{pentru } n > 0 \end{cases}$$

Cum transformata "z" a unei secvențe cauzale e determinată de partea reală a transformatei sale Fourier, rezultă

$$c(n) = \frac{1}{2} \int_{-\pi}^{\pi} |S(e^{j\omega})|^2 e^{-jn\omega} d\omega$$

Vom prefera calculul coeficienților cepstrali din coeficienții de predicție liniară (LPC) conform relațiilor recursive:

$$c(l) = -\sum_{i=1}^l a_i \cdot c[l-i]$$

Figura 2 prezintă evoluția coeficienților cepstrali pentru un vorbitor masculin.

Materialul vocal a fost achiziționat folosind un microfon (considerat fără zgomot) și a fost eșantionat cu frecvența de 8 kHz. Pentru analiza au lungimea de 240 ms, cu o suprapunere de 160 ms. Pentru predicție liniară s-a efectuat cu ordinul de predicție $p = 10$, iar pentru calculul coeficienților de predicție liniară s-a folosit algoritmul Levinson. Observația este aceea că modulul amplitudinii coeficienților de predicție liniară scade cu ordinul acestora. Pentru coeficienții de ordinul 5-10, evoluția coeficienților tinde să devină uniformă. Amplitudinea redusă a acestora a dus la dificultăți de estimare în condiții de zgomot.

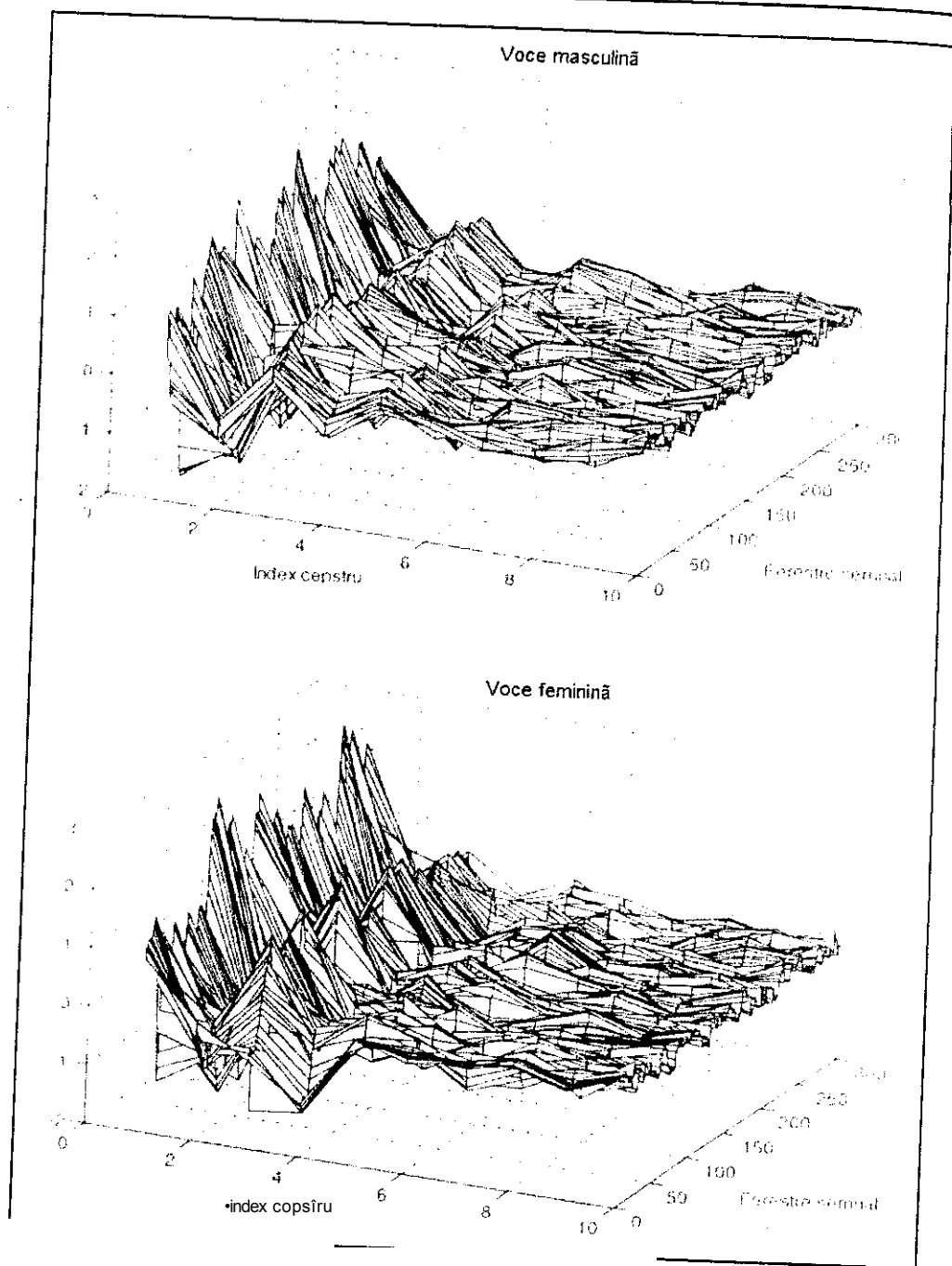


FIGURA 2 - Evoluția în timp a coeficienților cepstrali ai semnalului vocal

În scopul unei aprecieri calitative, fig. 3 prezintă, di
 cepstrali în planul c(1)-c(2), pentru aceiași doi vorbitori (masculin
 f; observa distribuția diferită a principalilor coeficienți cepstrali pent
 f remarcă o concentrare a coeficienților în anumite zone ale planu

Voce masculină

O
 c(1)

Voce feminină

c(1)

Figura 3. Reprezentarea coeficienților cepstrali în plan

În fig. 4 este prezentată distribuția parametrilor cepstrali corespunzători unui semnal vocal compus numai din vocalele limbii române. Ordinul analizei cepstrale este $p = 12$. Reprezentarea grafică s-a făcut numai în planul $c(1) - c(2)$. Se observă faptul că vocalele sunt relativ ușor separabile în spațiul cepstral, într-o configurație asemănătoare celei din spațiul formantic. Această analiză oferă premise interesante și pentru recunoașterea vorbirii în limba română.

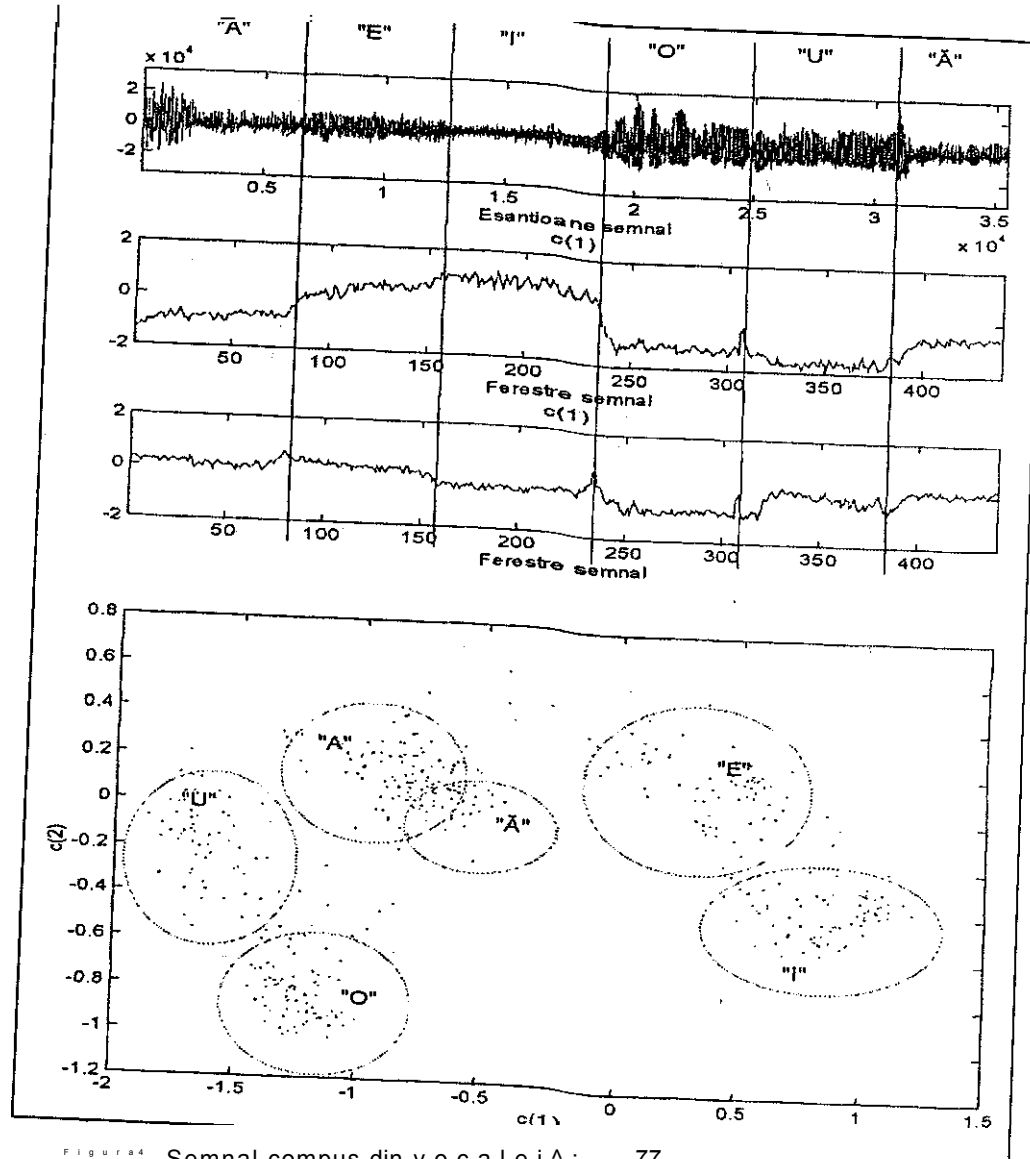


FIGURA 4 - Semnal compus din vocalele limbii române:
 • Parametri cepstrali corespunzători

3 Cuantizarea vectorială

Din punctul de vedere al sistemelor de recunoaștere a vorbirii, o persoană produce în timpul vorbirii o secvență de vectori de parametri care caracterizează atât vorbitorul cât și cuvintele pronunțate. Pentru un număr suficient de lung, ne așteptăm ca datele achiziționate să acopere în mod egal într-un mod care depinde mai mult de caracteristicile vorbitorului decât de ceea ce a pronunțat. Se face presupunerea că, având la dispoziție un număr suficient de date, se poate genera un model al vorbitorului care să descrie un proces de recunoaștere [20, 21].

Principiul cuantizării vectoriale este aplicat în sensul că se creează un volum mare de vectori acustico-fonetici, reprezentând material vorbit de către un vorbitor, într-un set restrâns de vectori denumit *tabelă de referință* (sau *centroizi*). În etapa de antrenare, partiționarea spațiului acoperit de vectorii de referință este făcută astfel încât media distanțelor minime dintre vectorii de referință și cel mai apropiat centroid să fie minimizată. În etapa de recunoaștere, un vector provenind de la un vorbitor necunoscut, este codat utilizând vectorii corespunzătoare vorbitorului vizat. Distorsiunea totală medie este minimă în decizia de recunoaștere [22].

Fie $\{X_i\}$ ansamblul de N versiuni cunoscute ale vectorului de referință.

Fie $\{G_k\}$ o partiție a acestui ansamblu în K clase; o clasă conține n_k elemente, astfel ca

$\sum_{k=1}^K n_k = N$

Notăm cu X_j^k cuvântul prototip ("centroid", "vector-cod") al clasei G_k .

Distanța medie între centroizi este

$$D = \frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} \|X_j^k - X_j^l\|^2$$

Distanța medie între vectorii dintr-o aceeași clasă, pentru o clasă G_k este

Raportul = reprezintă calitatea partiției

Algoritmul utilizat pentru găsirea centrozilor este atunci următorul:

- dacă cei K centroizi sunt aleși la întâmplare, clasele sunt constituite asociind fiecare vector X centroidului cel mai apropiat:

$$X \in G_k \quad \text{dacă} \quad D(x, X_k) < D(x, X^m) \quad \text{Mi} \pm k \quad (15)$$

- se iterează găsirea centrozilor căutând în fiecare clasă k vectorul X^m care are distanța față de vectorul cel mai depărtat al clasei minimă:

$$x^m = \underset{m}{\text{arg min}} D(x^m, X^k) \quad \text{dacă} \quad \max_m D(x^m, X^k) \text{ e minimă} \quad (16)$$

- această procedură e iterată până când centrozii sunt stabiliți.

Prezentăm în fig. 5 un exemplu de cuantizare vectorială folosind algoritmul Linde-Buzo-Gray (LBG). Vectorii cuantizați sunt coeficienții cepstrali de predicție liniară. Pentru reprezentarea în plan s-a ales sistemul de coordonate $c(1) - c(2)$. Dimensiunea tabelii de centroizi aleasă este $M = 8$. Se poate observa cum, în urma operației de optimizare, centrozii tind să "acopere" întregul spațiu ocupat de vectori. În mod evident, eroarea de cuantizare scade pe măsură ce dimensiunea tabelii de centroizi crește.

Pe parcursul algoritmului se pot utiliza diverse strategii de divizare. De exemplu, dacă după o operație de divizare și reclasificare, una dintre clase devine subpopulată sau chiar vidă, o alta va fi divizată la pasul următor, pentru a menține constant numărul total de clase. Se pot folosi următoarele criterii de alegere a clasei care va fi divizată: clasa care posedă cel mai mare număr de elemente, clasa care prezintă distorsiunea totală cea mai mare, clasa care prezintă distorsiunea medie cea mai mare. Folosind această structură arborescentă, clasificarea unui vector se poate efectua prin asocieri succesive, printr-o parcurgere a claselor găsite pentru fiecare nivel de divizare. În aplicațiile care necesită o acuratețe de clasificare ridicată, se preferă o metodă de clasificare prin căutare exhaustivă.

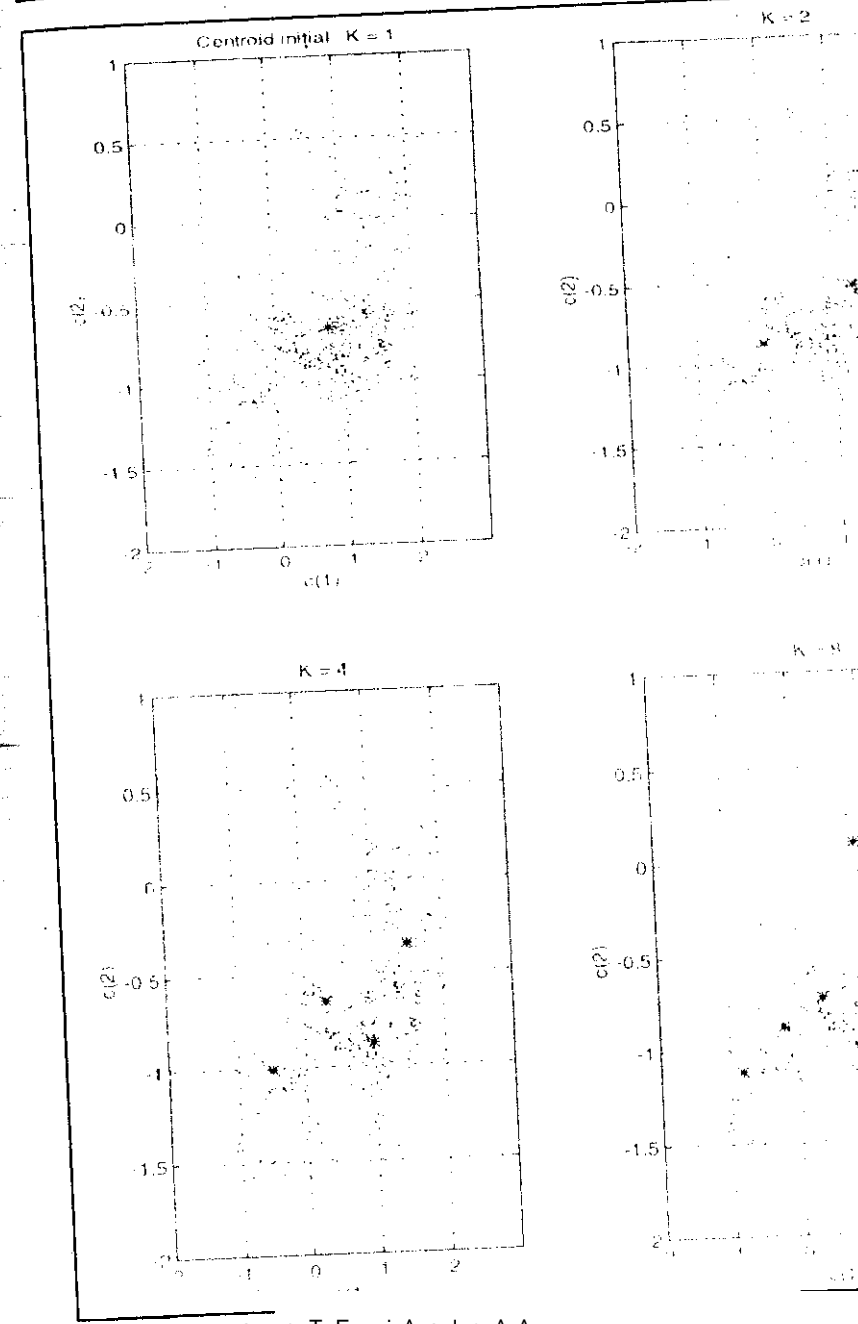


Figura 5. Evoluția algoritmului LBG

"•" - vectorii cepstrali; * - centrozii

4. Rezultate obținute

Un aspect important în proiectarea automatelor de recunoaștere a vorbitorului (eventual independent de text) îl reprezintă posibilitatea de evaluare a performanțelor acestora. Pentru a putea evalua un astfel de automat cu o precizie acceptabilă este nevoie de o bază de date corespunzătoare [23]. O astfel de bază de date trebuie să îndeplinească următoarele cerințe:

- să cuprindă material vocal achiziționat de la cât mai mulți vorbitori (de preferat, de ordinul zecilor sau sutelor);
- să conțină, eventual, dialecte diferite;
- să conțină fraze cât mai variate;
- frazele să fie rostite de mai multe ori, la intervale de timp
- pentru evaluare în condiții reale (de exemplu transmisie telefonică), materialul vocal trebuie să fie achiziționat prin intermediul mai multor aparate telefonice, în decursul mai multor legături, de preferat la distanțe diferite [24, 25].

Proiectarea și construirea unei astfel de baze de date este o sarcină dificilă.

Am folosit mai multe baze de date: internaționale, oarecum standard pentru procedurile de recunoaștere - TIMIT" și "YOHO", precum și o bază de date proprie, în română și engleză - "DiSPPALL".

Baza de date "TIMIT". conține eșantioane de voce provenind de la 630 de vorbitori, fiecare pronunțând 10 fraze. Experimentele descrise în lucrare au fost efectuate pe secțiunea TEST, care conține 168 vorbitori. Cele 10 fraze sunt: două fraze de calibrare (SA), cinci fraze compacte din punct de vedere fonetic (SX) și trei fraze variate contextual (SI). În experimente s-au folosit frazele SA și SX în faza de antrenare și frazele SI în cea de testare. Pentru evaluarea efectelor zgomotului telefonic în algoritmi de recunoaștere a vorbitorului, s-a folosit o variantă a bazei de date numită "NTIMIT". Aceasta conține aceiași material vocal ca și baza "TIMIT" cu deosebirea că acesta a fost transmis prin intermediul rețelei telefonice. Transmisia s-a făcut folosind un echipament de simulare a tractului vocal uman, în legături telefonice reale, la diferite distanțe.

Baza de date "YOHO" cuprinde fraze rostite de 138 de vorbitori (106 bărbați și 32 femei), iar vocabularul folosit constă din numere de două cifre rostite în grupuri de câte trei. Pentru fiecare vorbitor am folosit 4 sesiuni de antrenare de câte 24 de enunțuri și 10 sesiuni de verificare de câte 4 enunțuri.

Baza de date proprie "DiSPPALL" [26] cuprinde materialul vocal de la 26 de vorbitori (23 de bărbați și 3 femei) cu vârsta ce variază de la 21 la 50 de ani. Fiecare vorbitor în parte s rostit 31 de fraze: 11 fraze echilibrate din punct de

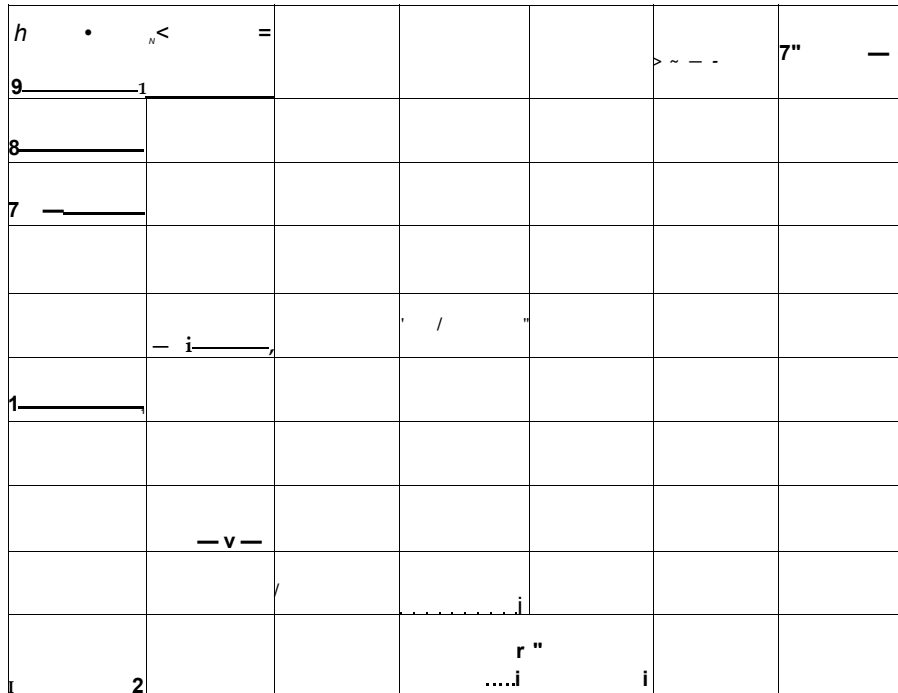
vedere fonetic au fost folosite pentru antrenare și 20 de fraze pentru verificare: 5 enunțuri de bază repetate de câte 4 ori. Frazele de verificare au fost înregistrate în două sesiuni diferite, 5 enunțuri de bază fiind repetate de două ori în fiecare sesiune. Prima sesiune de verificare a fost înregistrată în același timp cu sesiunea de antrenare, iar sesiunea a doua a fost înregistrată după două-trei săptămâni, înregistrările s-au făcut cu un microfon de tip "head-set" într-o cameră cu zgomot ambiental normal: spre deosebire de baza "YOHO", baza "DiSPPALL" conține material vocal alterat de zgomot pentru a face condițiile de test mai dificile și mai apropiate de o situație reală de recunoaștere a vorbitorilor

În experimentele de verificare a vorbitorului, o frază de test este comparată cu referința vorbitorului a cărei identitate se dorește verificată, calculându-se o distorsiune totală medie. Dacă aceasta este mai mică decât un prag dat, vorbitorul este considerat acceptat, altfel el este respins. Există două tipuri de erori asociate procesului de verificare: respingerea utilizatorului căruia îi aparține referința (denumită eroare de tip I) și acceptarea unui impostor (eroare de tip II) [27]. Fiecare frază de test este comparată cu referințele corespunzătoare tuturor vorbitorilor din baza de date aleasă pentru test. Pragurile de decizie nu sunt fixate a priori ci se determină distanța medie totală pentru care eroarea de tip I este egală cu cea de tip II ("rata-erorii-egale"). Valoarea corespunzătoare a erorii este considerată rezultatul final al procesului de evaluare. În fig. 6 sunt prezentate rezultatele procesului de verificare a vorbitorului, folosind cuantizarea vectorială, utilizând baza de date "TEST/TIMIT". Ordinul de predicție (și implicit dimensiunea vectorilor cepstrali) este $P = 10$ iar dimensiunea tabelii de centroizi, $M = 64$. Ca distanță vectorială s-a folosit distanța euclidiană ponderată.

$$d(\mathbf{v}_a, \mathbf{v}_b) = -\ln \sum_j \mathcal{L}_j(\mathbf{v}_a - \mathbf{v}_b)^2 \quad (17)$$

unde s_j este varianta componentei j calculată pe întreg setul vectorilor de antrenare. Ca metodă de cuantizare vectorială s-a folosit algoritmul LBG modificat.

Sunt evidente tendințele contrare ale erorilor de tip I, respectiv II. Rata erorii-egale pentru evaluarea de mai sus este 6.8%, corezpunzând unui prag de decizie egal cu 2.8. În funcție de aplicația concretă, pragul de decizie se poate stabili a posteriori la o altă valoare, adecvată scopului propus. Spre exemplu, dacă se dorește limitarea acceptării impostorilor la 2%, respingerea adevăraților utilizatori va fi de 19.7%. Reciproc, pentru o eroare de respingere a utilizatorilor reali de 2%, acceptarea impostorilor va fi de 12.9%.



Distant medie totala

Figura 6. Eroarea de verificare a unui sistem de recunoaștere a vorbitorului utilizând cuantizarea vectorială

În experimentele de identificare a vorbitorului, fiecare frază de test provenind de la un vorbitor considerat necunoscut este comparată cu referințele fiecărui vorbitor din baza de date aleasă pentru test. Referința asociată cu cea mai mică distorsiune totală medie față de fraza de test este considerată ca aparținând vorbitorului identificat. În funcție de corespondența dintre apartenența frazei de test și a referinței aceluiași vorbitor sau unor vorbitori diferiți, se decide dacă rezultatul procesului de identificare este adevărat sau fals. Eroarea de identificare este calculată ca raportul dintre numărul de identificări incorecte și numărul total de identificări [28, 29, 30].

5. Utilizarea frecvenței fundamentale în recunoașterea vorbitorului

Frecvența fundamentală poate fi utilizată ca parametru discriminator suplimentar în conjuncție cu algoritmi de cuantizare vectorială a vectorilor cepstrali.

Frecvența fundamentală F_0 sau perioada fundamentală T_0 (cunoscută și sub numele de "pitch"), constituie un parametru important al vocii umane, care își găsește utilizări practice în multe domenii ale procesării vorbirii. Încercări de utilizare a frecvenței fundamentale în procesul de recunoaștere a vorbitorului se cunosc încă de la începutul anilor 70, aceasta fiind pusă în corespondență directă cu prozodia. Majoritatea acestor experimente s-au desfășurat utilizând sisteme de recunoaștere dependente de text și metode de aliniere temporală. Sistemele de recunoaștere a vorbitorului independente de text bazate exclusiv pe frecvența fundamentală nu au dat rezultate satisfăcătoare.

Ideea prezentată în secțiunea de față este aceea de a folosi frecvența fundamentală ca un parametru discriminator suplimentar, în conjuncție cu algoritmi de cuantizare vectorială a vectorilor cepstrali [31]. Justificarea teoretică a acestor abordări rezidă în primul rând în modelul predicției liniare aplicat semnalului vocal care presupune, așa cum am arătat o separare clară între sursa de semnal și tractul vocal. De asemenea, am arătat în secțiunea 2 că analiza cepstrală folosită pentru extragerea vectorilor cepstrali este un proces de deconvoluție, coeficienții cepstrali obținuți caracterizând în mod exclusiv tractul vocal. Ca atare, utilizarea ca date de intrare în același sistem atât a vectorilor cepstrali cât și a frecvenței fundamentale nu reprezintă o abordare redundantă.

Cerințele de bază ale unui algoritm de extragere a frecvenței fundamentale sunt: acuratețea de estimare (evitarea armonicilor), robustețea deciziei sonor/nesonor, insensitivitatea la zgomot, volumul de calcule minim. Se cunosc numeroși algoritmi de estimare a frecvenței fundamentale (AMDF, Dubnowski, Rabiner, SIFT, etc), fiecare prezentând avantaje și dezavantaje. Trebuie arătat faptul că, din cauza, în principal, comportării netaționare a semnalului vocal niciunul din algoritmii cunoscuți nu este considerat perfect. Cu alte cuvinte, se acceptă ideea existenței erorilor atât în luarea deciziei sonor/nesonor cât și în obținerea valorilor propriu-zise ale frecvenței fundamentale. În experimentele prezentate mai jos s-a folosit algoritmul Rabiner, considerat ca fiind unul dintr-cele mai robuste.

Ideea introdusă este aceea de a utiliza frecvența fundamentală în scopul unei clasificări grosiere a potențialilor candidați, atât pentru sarcina de verificare a vorbitorului, cât și pentru cea de identificare. În cazul verificării, scopul propus este acela de a reduce eroarea de tip II, prin eliminarea vorbitorilor a căror frecvență fundamentală nu "corespunde" cu cea a vorbitorului de referință. În cazul sarcin

de identificare, se dorește reducerea numărului de candidați posibili, fără a afecta acuratețea de identificare. Aceasta poate conduce la o reducere majoră a volumului de calcule, dat fiind că estimarea frecvenței fundamentale se face o singură dată pentru fiecare vorbitor și este mai puțin consumatoare de timp decât clasificarea vectorială.

Având în vedere considerentele de mai sus, schema de principiu a sistemului de recunoaștere a vorbitorului modificat prin introducerea frecvenței fundamentale ca parametru discriminator este prezentată în fig. 7.

Un aspect important în utilizarea frecvenței fundamentale în aplicațiile de recunoaștere a vorbitorului îl reprezintă alegerea formei de prelucrare a datelor furnizate de estimator. "Conturul de pitch", reprezentând evoluția în timp a parametrului F_0 , deși utilizat în sisteme de recunoaștere a vorbitorului dependente de text, conține un volum de date dificil de utilizat în operații de discriminare. În consecință, s-a încercat o reducere a datelor la câțiva parametri statistici. Au fost investigate patru valori statistice derivate din conturul de pitch: valoarea medie, valoarea maximă, valoarea minimă și dispersia (deviația standard). Pentru fiecare vorbitor, s-au calculat aceste valori pe ansamblul materialului vocal disponibil. Ca parametru de discriminare a fost utilizat raportul valorilor statistice de mai sus

$$P_{medie} = \frac{\text{media } F_0 \text{ pentru antrenare}}{\text{media } F_0 \text{ pentru test}} \quad (18)$$

tratându-se în mod similar toate celelalte valori statistice (maximă, minimă, dispersie). Pentru a evalua utilitatea acestor parametri în procesul de discriminare, s-a determinat distribuția fiecăruia atât pentru frazele pronunțate de aceiași vorbitori (intra-vorbitor) cât și pentru toate combinațiile de fraze pronunțate de vorbitori diferiți (inter-vorbitor).

Modul de discriminare a vorbitorilor este următorul: fixându-se un prag θ , dacă

$$P_{medie} > \theta \quad (19)$$

vorbitorul este rejectat și nu se execută clasificarea vectorială. În caz contrar, vorbitorul este considerat potențial candidat și urmează procesul de clasificare prin cuantizare vectorială.

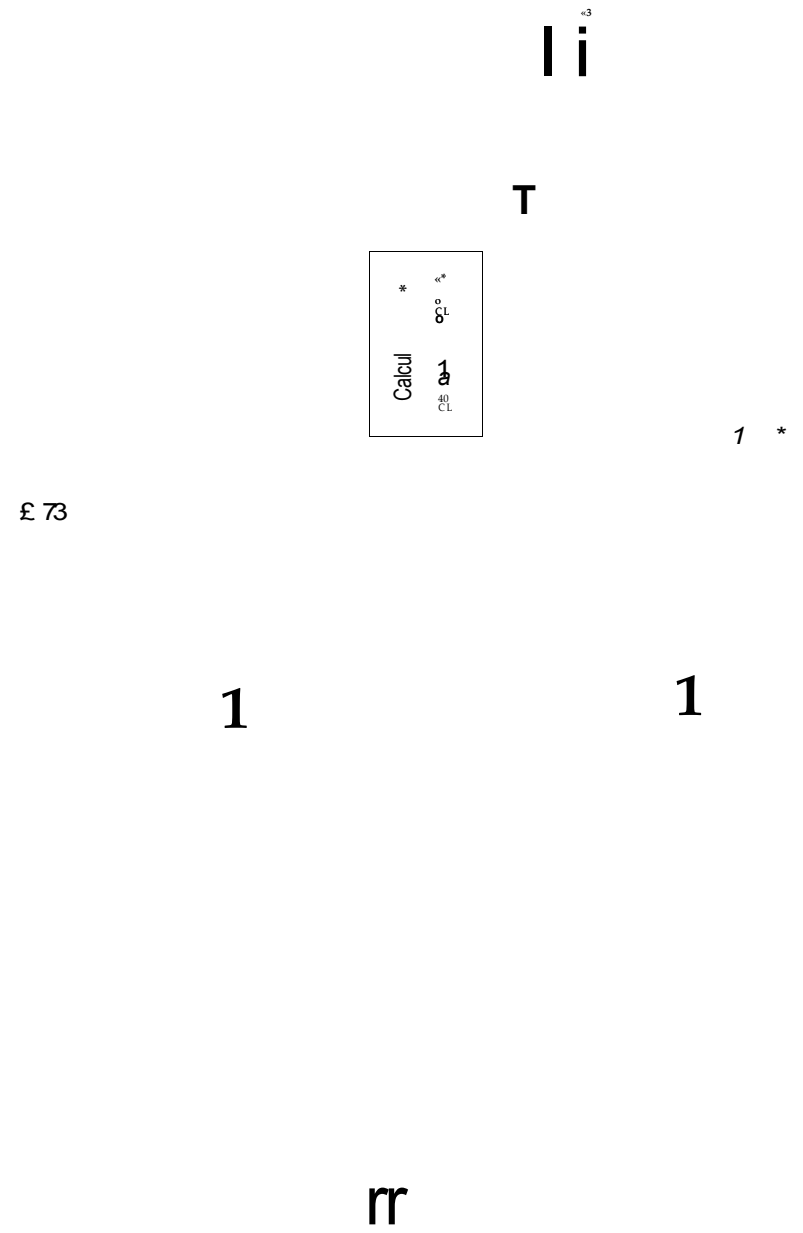


Figura 7. O variantă a sistemului de recunoaștere a vorbitorului - schema de principiu

Utilizând elementul de discriminare descris mai sus s-au obținut îmbunătățiri importante atât în procesul de verificare a vorbitorului cât și în cel de identificare. Rezultatele obținute pentru 14 coeficienți cepstrali și 128 centroizi sunt prezentate în tabelul 1.

Tabelul 1

e	Neutilizat	0.30	0.25	0.20	0.15	0.10
EER la verificare (%)	6.3	6.1	5.3	3.9	2.7	6.5
Eroarea de identificare (%)	6.2	6.2	5.9	5.6	5.5	9.4
Candidați identificare (%)	100	57.2	49.1	43.4	32.3	26.5

Cele mai bune rezultate s-au obținut pentru $\delta = 0.15$, caz în care eroarea de verificare obținută este de aproape 2.5 ori mai mică decât în cazul folosirii doar a clasificării vectoriale. În cazul identificării, deși îmbunătățirile de acuratețe nu sunt impresionante, cel mai important rezultat îl reprezintă reducerea numărului candidaților, cu peste 65%. Pentru valori ale lui δ mai mici decât 0.10, se observă o degradare abruptă a performanțelor de verificare și identificare, ceea ce indică faptul ca variația intra-vorbitor a frecvenței fundamentale medii este mai mare decât acest prag.

6. Concluzii

Lucrarea de față se ocupă de un aspect bine delimitat al tehnologiei vorbirii și anume recunoașterea vorbitorului ca parte integrantă a recunoașterii automate și mai departe a dialogului om-mașină. Tipurile de probleme care apar sunt similare pentru întreg domeniul recunoașterii automate.

Am precizat presupunerile fundamentale care au stat la baza analizei propuse (în special opțiunea de a aborda proiectarea ținând seama de mecanismul producerii vorbirii); insistăm asupra faptului că aceste abordări nu sunt obligatorii, ci constituie alternative care au avantaje și dezavantaje.

S-au trecut în revistă etapele esențiale ale procedurilor de recunoașterea vorbitorului: achiziția semnalului vocal, prelucrarea acustico-fonetică, recunoașterea propriu-zisă.

Am subliniat importanța parametrizării semnalului vocal. Analiza cepstrală care a fost aleasă pentru reprezentarea parametrică a semnalului vocal este legată de opțiunile fundamentale de analiză: separarea efectelor sursei de semnal și ale tractului, separarea efectelor diverselor porțiuni din tractul vocal, analiza "în timp scurt"

Am utilizat cuantizarea vectorială ca metodă de recunoaștere. Sunt prezentate o parte dintre rezultatele experimentelor realizate. Subliniem importanța

utilizării unor baze de date specifice și, în consecință, am acordat spațiu prezentării acestora.

O contribuție pe care o considerăm interesantă la îmbunătățirea performanțelor recunoașterii vorbitorului o constituie utilizarea frecvențelor fundamentale ca parametru discriminator grosier. Sunt prezentate o serie de rezultate care probează în ce mod anumite performanțe sunt superioare abordării "clasice".

O parte dintre rezultatele obținute sunt susceptibile de a fi generalizate pentru recunoașterea vorbirii în limba română [32] (de pildă, coeficienții cepstrali pentru foneme ale limbii române). De asemenea, utilizarea frecvențelor fundamentale apropie recunoașterea vorbitorului de o anumită dependență de limba în care sunt rostite frazele de antrenare și de test.

Referințe bibliografice

- [1] M.Drăgănescu, C.Burileanu, coordonatori (1986). Analiza și sinteza semnalului vocal - Editura Academiei Române, București.
- [2] M.Dragănescu, G.Stefan, C.Burileanu (1991). Electronica funcțională - voi. Editura tehnică, București, ISBN 973-31-0290-3.
- [3] G. Yu and H. Gish (1993). Identification of Speakers Engaged in Dialog, Proc. of IEEE Int. Conf. Acoust, Speech, Signal Processing, Voi.II, p. 383-386.
- [4] Sadaoki Furui (1994). An Overview of Speaker Recognition Technology, Proc. of ESCA Workshop on Automatic Speaker Recognition, Identification and Verification, p. 1-9.
- [5] Y. Bennani, P. Gallinari (1994). Connectionist Approaches for Automatic Speaker Verification - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 95-103.
- [6] M. Hanah s.a. (1994). The Role of the Reference Template in Speaker Verification - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 181-184.
- [7] Chi-Shi Liu; Hsiao-Chuan Wang; Lee, C. (1996) Speaker Verification Using Normalized Log-Likelihood Score, IEEE Tr. on Speech and Audio Processing Voi. 4. Issue 1, p. 56
- [8] S. Nakagawa, K. P. Markov (1997). Speaker Verification Using Frame and Utterance Level Likelihood Normalization, Proc. of SPCHL97 ,Vol. 2, p. 1087
- [9] K.T. Assaleh, R.J. Mammone (1994). New LP - Derived Features for Speaker Identification - IEEE Tr.on SAP, vol.2, no.4, p. 630-638.

- [10] H. Gish, M. Schmidt (1994). Text-Independent Speaker Identification - IEEE Signal Proc. Mag., vol.11, nr.4, p. 18-32.
- [11] Q. Lin s.a. (1994). Microphon Array Speaker Identification - IEEE tr. on ASSP, vol.2. nr.4, p. 622-629.
- [12] D. Reynolds (1994). Experimental Evaluation of Features for Robust Speaker Identification - IEEE Tr. on ASP, vol.2, nr.4, p. 639-643.
- [13] F. Bimbot, G. Chollet, A. Paoloni (1994). Assessment Methodology for Speaker Identification and Verification - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 75-83.
- [14] M. Abe, S. Sagayama (1990). Statistical Study on Voice Individuality Conservation Across Different Languages - Proc. of ICSLP, p. 157-160.
- [15] Y. Gong, J.P. Haton (1994). Non-Linear Interpolation Methods for Speaker Recognition - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 23-26.
- [16] J. He s.a. (1995). On the Use of Features from Prediction Residual Signal in Speaker Identification Proc. of EUROSPEECH95, p. 313-316.
- [17] D.Naik s.a. (1994). Robust Speaker Identification Using Pole Filtering - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 225-228.
- [18] J. Openshaw, J. Masson (1994). Optimal Noise-Masking of Cepstral Features for Robust Speaker Identification - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 231-234.
- [19] J. Thompson, J.S. Masson (1993). Within Class Optimization of Cepstra for Speaker Recognition, Proc. of EUROSPEECH, p. 165-168.
- [20] K. Sonmez, L. Heck, M. Weintraub (2000). Multiple Speaker Tracking and Detection: Handset Normalization and Duration Scoring, Digital Signal Processing, 10(1/2/3), p. 133-143.
- [21] T. Isobe, J. Takahashi (1999). A New Cohort Normalization Using Local Acoustic Information for Speaker Verification, Proc. IEEE Int. Conf. Acoust. Speech, Signal Processing, 26.8, voi. 2, p. 841-844.
- [22] X. Zhu s.a (1994). Text-Independent Speaker Recognition Using VQ, Mixture Gaussian VQ and Ergodic HMMs - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 55-58.
- [23] L. Boves s.a. (1994). Design and Recording of Large Data-Bases for Use in Speaker Recognition and Identification - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 43-46.
- [24] A. Federico (1995). Parametric Speaker Recognition Over Large Population of Telephonic Voices - Proc. of EUROSPEECH95, p. 329-332.

- [25] J.L. Gauvain s.a (1995). Experiments with Speaker Verification over the Telephone - Proc. of EUROSPEECH95, p. 651-654.
- [26] C. Burileanu, D. Burileanu s.a.(2000). Cohort Normalisation Methods for Speaker Verification - Proc. of International Conference "Communications 2000", Bucharest, România, p.118-121.
- [27] M. Wagner s.a. (1994). Analysis of Type-II Errors for VQ-Distortion Based Speaker Verification - ESCA Workshop on Speaker Recognition, Identification and Verification, p. 83-86.
- [28] J.F. Bonastre (1993). Automatic Speaker Recognition and Analytic Process Proc. of EUROSPEECH93, p. 441-444.
- [29] M. Sugiyama s.a. (1993). Speech Segmentation, Clustering Based on Speaker Features - Proc. of ICASSP, p.395-398.
- [30] H. Beigi, S. Maes and J. Sorensen (1998.) A Distance Measure Between Collections of Distributions and Its Application to Speaker Recognition, Proc. of ICASSP, Voi. 2, p. 753-756.
- [31] L.E. Bojan, C. Burileanu s.a. (1996). Enhancements in Automatic Speaker Verification and Identification for Large Data-bases Using Pitch Contour Analysis - Proc. of ICSPAT96, Boston, SUA, p. 1796-1800
- [32] C. Burileanu, L.E. Bojan s.a. (1993). A Representation for Recognition of Isolated Words Spoken in the Romanian Language - Proc. of ICSPAT93, Santa Clara, USA, p. 1478-1484.

Prelucrarea inițială a textului de intrare în cadrul unui sistem de sinteză a vorbirii pornind de la text în limba română

Dragoș BURILEANU
Laboratorul de "Tehnologia vorbirii și prelucrarea digitală a semnalelor",
Facultatea de Electronică și Telecomunicații,
Universitatea "POLITEHNICA" București
Bdul Iuliu Maniu 1-3, Sector 6, 77202 București
bdragos@mESsnet.pub.ro

1. Introducere

Limbajul reprezintă modalitatea de exprimare a ideilor prin intermediul unui ansamblu de semne, fie grafic, fie prin gesturi, sau sunete, un astfel de sistem structurat fiind specific doar oamenilor. Fără îndoială, *vorbirea* este una din principalele sale componente; ea este cea mai veche modalitate de comunicare între oameni și este și astăzi cea mai răspândită. Este deci ușor de înțeles faptul că vorbirea a fost studiată intens și s-a încercat adesea să fie prelucrată într-un mod automat. Pentru mulți ingineri și specialiști din domeniu, posibilitatea de a conversa liber cu o mașină reprezintă de fapt o adevărată provocare pentru înțelegerea cât mai deplină a proceselor de producere și percepție implicate în comunicarea prin voce între oameni. Ceea ce este însă și mai important este faptul ca *interfețele de comunicare prin voce* devin tot mai mult o necesitate. În viitorul apropiat, sistemele și rețelele interactive vor oferi un acces simplu și ieftin la cantități mari de informație și servicii, ceea ce va afecta fundamental viața noastră zilnică.

Deși principiile de bază ale producerii și recepționării vorbirii au început să fie studiate încă de la sfârșitul secolului al XVIII-lea, când s-au înregistrat primele cercetări în domeniul dezvoltării sintetizoarelor mecanice de sunete asemănătoare vocii umane, *tehnologiile de prelucrare a vorbirii* au obținut rezultate semnificative doar în ultimele decenii (fiind denumite în sens larg *tehnici de analiză și sinteză a semnalului vocal*). Aceste rezultate au fost posibile datorită progreselor făcute în domeniile acusticii și lingvisticii, modelării matematice a producerii și percepției vorbirii, prelucrării semnalelor și tehnologiilor VLSI. Putem evidenția în acest sens dezvoltarea procesoarelor numerice de semnal pe un singur chip, realizarea de

capsule de memorie mai mari și mai ieftine, apariția unor algoritmi îmbunătățiți pentru prelucrare de semnal, iar în domeniul comunicațiilor crearea de standarde globale pentru transmisie, compresie de semnal și protocoale de comunicație.

Prin urmare, putem aprecia că cercetările actuale în domeniul prelucrării vorbirii au ca scop larg îmbunătățirea calității, securității și costului comunicațiilor și a accesului uman la informații. Pe de o parte, este de așteptat în viitorul apropiat o extindere importantă a serviciilor integrate de voce, poșta electronică, FAX, paging și transmisiuni de date pe canale fără fir. Pe de altă parte însă, comunicarea verbală între om și mașini, în ambele sensuri, tinde deja să devină o realitate, fiind vizibilă tendința actuală de a apropia caracteristicile mașinii de cele ale utilizatorului uman.

În acest ultim sens, trebuie observat faptul că tendința menționată anterior este absolut firească. Filozoful grec Aristotel (384 - 322 î.C, fondator al logicii formale), afirma: "*Rațiunea de a fi a oricărui lucru constă în funcția sa*". Ori este evident faptul că o interfață de dialog prin voce reprezintă o modalitate ideală de comunicare cu mașina, vorbirea fiind cea mai naturală, flexibilă, eficientă și economică modalitate de comunicare utilizată de oameni.

Aceste idei legate de posibilitatea comunicării prin voce între om și mașina nu sunt noi; totuși, doar în ultimii ani a început să prindă contur conceptul ce a căpătat denumirea de "dialog om-mașină", iar tehnologia necesară implementării acestui concept a părăsit deja laboratoarele și a pătruns în lumea reală, într-o gamă largă de aplicații.

Pentru a realiza un mod de comunicare cât mai natural și pentru a permite o utilizare cât mai largă, calculatorul trebuie să înțeleagă și să producă singur vorbirea; acesta este motivul principal pentru care *recunoașterea și sinteza vorbirii* au devenit în ultimii ani tehnologii de un interes special și constituie subiecte pentru cercetări intense și aprofundate. Ambele tehnologii prelucrează vorbirea în primul rând sub aspectul conținutului informațional: recunoașterea transformă vocea omului în text ce poate fi folosit literal (de exemplu pentru dictare), sau o interpretează sub forma unor comenzi de control pentru diverse aplicații, iar sinteza permite generarea limbajului vorbit pornind de la text sau de la anumite concepte.

Cu toate că s-au făcut pași importanți în aceste domenii, rezultatele sunt încă departe de așteptări. Sarcinile enunțate inițial s-au dovedit în timp a fi deosebit de dificile, în primul rând datorită complexității semnalului vocal ca și a dificultăților legate de prelucrarea acestuia, dificultăți legate fie de recunoașterea conținutului său informațional (semnalul vocal depinzând puternic de vorbitor și de condițiile în care acesta rostește un mesaj), fie de producerea sa, fie de transmiterea acestui semnal la distanță [1].

În acest context, producerea vorbirii artificiale și în special cea de tip *voce*, care constituie obiectul principal al lucrării de față, este astăzi o bază al domeniului prelucrării vorbirii și subiect al unor cercetări intense. *Sinteza de vorbire pornind de la text* (TTS - "*Text-to-Speech*") poate oferi o soluție de aplicații, de la accesul la poșta electronică și diferite tipuri de bază de pronunțarea unui text pentru persoane cu handicap vizual.

Este important de observat faptul că tehnologia de răspuns prezintă o serie de avantaje fundamentale pentru transmiterea informațiilor:

- oricine poate înțelege un mesaj, fără antrenare sau deosebită;
- mesajul poate fi recepționat chiar dacă cel ce ascultă este ocupat cu alte activități, cum ar fi mersul, manipularea unor obiecte sau utilizarea altor informații;
- rețeaua telefonică convențională poate fi utilizată pentru comunicarea la distanță la o bază de informații;
- această formă de comunicare este mai economică decât comunicarea tradițională prin mesaje scrise.

Toți acești factori precum și numeroasele aplicații cerute de piață au creat premisele unor cercetări aprofundate, obținându-se astfel sisteme comerciale care pot produce vorbire sintetică pornind de la informații inteligibile acceptabile.

Într-adevăr, scopul principal al celor mai multe sisteme existente este de a produce o vorbire *inteligibilă*. Din acest punct de vedere sinteza pare a fi de mai multă vreme o tehnologie "stabilă", ieftin implementat; se spune chiar, uneori, că acest domeniu este în prezent bine dezvoltat, iar problemele rămase sunt minore din punct de vedere. Dacă însă scopul este sinteza în timp real, pornind de la un vocabular de cuvinte și fără restricții asupra textului, iar vorbirea să fie nu numai inteligibilă la fel de *naturală* ca cea umană, atunci se constată că performanțele sunt departe de a fi satisfăcătoare. Rămân încă multe probleme importante: extinderea vocabularului oferit, înlăturarea restricțiilor impuse textului unor caractere speciale, îmbunătățirea caracteristicilor de *prozodie*, modificarea *ritmului* și *stilului* vorbirii sintetizate, sau elaborarea unor sisteme de sinteză în mai multe limbi. Aceste sarcini se dovedesc a fi deosebit de dificile, cer, evident, eforturi interdisciplinare susținute [2].

2. Sinteza automată a vorbirii

Etimologic, cuvântul "sinteză" provine din limba greacă și semnifică îmbinarea mai multor elemente diferite într-un tot.

În ceea ce privește *sinteza vorbirii*, nu există o definiție precisă și unanim acceptată de către specialiștii în tehnologia vorbirii. Acest termen a avut în decursul timpului mai multe accepțiuni, majoritatea depinzând de nivelul tehnologic al momentului și de elementele constitutive ale semnalului vocal care au fost folosite pentru sinteză. De exemplu, primele circuite integrate care permiteau simpla restituire a unui mesaj vocal înregistrat și stocat digital au purtat denumirea de "sintetizoare vocale", fie că se făcea sau nu o compresie a semnalului. Este evident că în acest caz nu se poate vorbi de sinteză, din moment ce **textul este fix** și astfel de sisteme nu pot rosti decât mesaje preînregistrate; chiar dacă vocea umană este comprimată cu ajutorul unui algoritm, nu este cu adevărat "sintetică", ci poate fi numită mai curând o "înregistrare cu număr redus de biți".

Aceeași situație este în cazul sintezei la recepție a unor mesaje transmise pe canale de comunicație standard (caracteristică sistemelor de tip "vocoder"), care este de obicei considerată ca făcând parte din domeniul *codării vorbirii* și cuprinde tehnici de reducere a debitului semnalului vocal pentru transmisie; cu alte cuvinte, și acest tip de sinteză, care reface **același** mesaj analizat la emisie, deci nu generează fraze **noi**, nu este tratat ca o sinteză automată propriu-zisă.

O categorie distinctă de sinteză vocală este aceea care implică sisteme ce concatenează cuvinte sau fraze preînregistrate, dar generează fraze noi, acestea nefiind niciodată pronunțate ca atare; astfel de sisteme cer utilizarea unor reguli lingvistice mai mult sau mai puțin complicate pentru a funcționa corespunzător.

În sfârșit, o categorie specială o reprezintă sinteza vorbirii pornind de la text; aceasta reprezintă, în esență, transformarea unui text oarecare, scris într-un anumit limbaj, în semnal vocal. Trebuie remarcat faptul că în prezent, în multe lucrări științifice, acest tip de sinteză este sinonim chiar conceptului de sinteză automată a vorbirii.

Analizând exemplele de mai sus, putem defini trei noțiuni generale [3], pe care le vom utiliza pe parcursul lucrării de față:

Definiția 2.1 *Sinteza automată a vorbirii* este "tehnologia integrată care simulează procesul uman de generare a vorbirii, mergând de la sisteme simple ce pot genera automat fraze noi și cuprind un formalism lingvistic minimal și până la sisteme care transformă în vorbire reprezentări simbolice sau lingvistice ale limbajului".

Definiția 2.2 Un *sistem de sinteză pornind de la text* este "un sistem automat care poate produce vorbirea plecând de la un text scris, prin intermediul unei reprezentări fonetice a mesajului".

Definiția 2.3 *Sintetizorul vocal* este "etajul unui sistem de sinteză automată a vorbirii care realizează conversia finală în semnal vocal, pornind de obicei de la o reprezentare parametrică a unor segmente acustice fundamentale".

3. Sinteza vorbirii pornind de la text

Pentru a înțelege mai bine dificultatea sarcinii unui sistem de sinteză pornind de la text, considerăm că este util să punem în evidență mai întâi modul (fiziologic) în care o persoană citește cu voce tare un text. Imaginea textului este sesizată de neuronii sistemului vizual, transmisă creierului sub forma unor stimuli electrici, aici fiind prelucrată pentru a putea permite comanda neuronilor responsabili de corectă activare a plămânilor^ coardelor vocale și organelor articulatorii. În acest fel se produce vorbirea, ea fiind permanent monitorizată de creier (în special prin intermediul organelor auditive), în scopul ajustării configurației tractului vocal în timp real.

Desigur, cunoaștem încă prea puțin despre organizarea de ansamblu a sistemului nervos uman, care este capabil de această activitate complexă; puteam propune totuși următorul *model funcțional* prin care este prelucrată informația optică și apoi este dată comanda de generare a vorbirii:

- Atunci când citim un text, efectuăm practic o sarcină de *recunoaștere de caractere*, ignorând, parțial inconștient, anumite erori de redactare a cuvintelor (caractere lipsă sau înlocuite cu altele) și decodificând mai degrabă cuvântul ca un întreg; are loc un proces de inferență asupra informației dintr-un context posibil incomplet. De asemenea, recunoaștem cu ușurință caractere speciale sau abrevieri.
- Considerând *fonemele* ca fiind cele mai mici elemente sonore care permit diferențierea între ele a cuvintelor, este evident că secvența fonemică corespunzătoare unui cuvânt diferă de șirul de caractere grafice din care este compus cuvântul; creierul trebuie să facă pe urmare o *transcriere fonetică pornind de la litere*, această operație fiind practic instinctivă permițând pronunția unui număr nelimitat de cuvinte.
- În cele mai multe situații, suntem capabili să începem pronunția unei fraze mult înainte de terminarea ei; cu alte cuvinte, putem face o *structurare sintactică*, descompunând fiecare propoziție în grupuri de cuvinte și asociindu-le intonația corespunzătoare. Și acest proces este practic inconștient, fiind bazat pe educație și experiență.
- În sfârșit, putem discrimina cu ușurință cuvinte ce se scriu asemănător dar se pronunță diferit, după *înțelesul semantic*, fapt posibil datorat aceleiași capacități de deducție a creierului de care am vorbit mai sus.

Concluzia este simplă: pe baza experienței lingvistice căpătate în urma educației, o persoană familiară cu limbajul în care este scris un text depășește imediat pașii descriși anterior și poate cu ușurință să citească cu voce tare textul scris, în primul rând pentru că înțelege ceea ce citește.

Având în vedere considerațiile expuse anterior, devine evident faptul că o mașină care trebuie să pronunțe un text scris nu va putea adopta o schemă de prelucrare atât de complicată cum este cea care caracterizează acțiunea citirii cu voce tare a unui text de către o persoană. Sunetele vorbirii sunt inerent guvernate de ecuații diferențiale ale mecanicii fluidelor, aplicate într-un context nestaționar, deoarece presiunea aerului ia nivelul plămânilor, tensiunea glotală, ca și configurațiile tractului vocal și nazal, evoluează în timp. Toate acestea sunt controlate de creierul uman, care beneficiază de avantajul puterii sale de prelucrare paralelă pentru extragerea esenței textului citit: înțelesul. Chiar și la nivelul la care a ajuns știința astăzi (cercetări intense în domeniile sintezei articulatorii, rețelelor neuronale artificiale și prelucrării limbajului natural), construirea unui sistem de sinteză pornind de la text cu un model atât de complex rămâne practic nerealizabilă; chiar dacă, să spunem, s-ar ajunge foarte aproape de aceste cerințe, sistemul rezultat nu ar fi de loc compatibil cu criteriile economice normale.

Figura 1 introduce o diagramă funcțională foarte generală a unui sistem TTS, bazată pe observațiile anterioare.

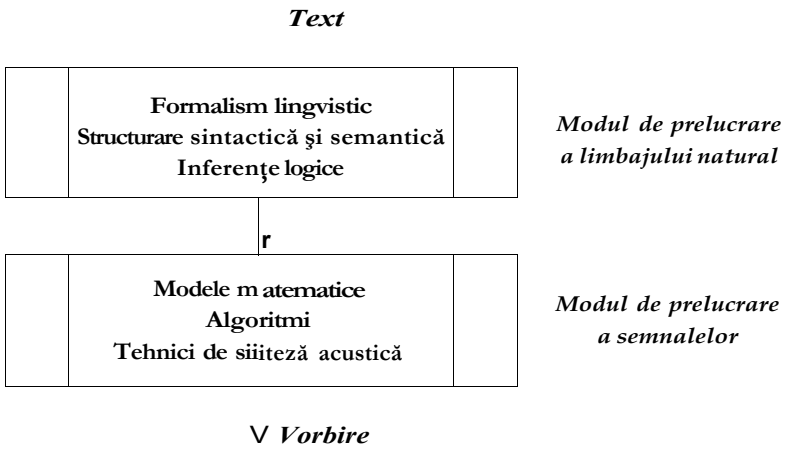


Figura 1. Diagramă funcțională pentru un sistem TTS

Ca și pentru un cititor uman, schema cuprinde un *modul de prelucrare a limbajului natural*, capabil să producă o transcriere fonetică a textului citit împreună cu informații despre intonație, accente, durate și de asemenea un *modul de prelucrare a semnalelor*, care transformă informația simbolică primită în vorbire sintetică, pe baza unor tehnici de sinteză adecvate și a unor structuri stocate în urma unei analize preliminare. Etapele de bază ale sintezei pornind de la text pot astfel descrise printr-un număr de transformări succesive ce trebuie aplicate asupra șirului de caractere ce reprezintă textul de intrare; scopul este de a se obține o vorbire **de calitate**, într-o limbă oarecare, fără constrângeri asupra textului introdus.

Trebuie menționat faptul că formalismul descris poate "sări" uneori peste anumii pași, dacă se utilizează în mod adecvat cunoașterea lingvistică și matematică; acest lucru se întâmplă atunci când punem anumite restricții asupra textului ce trebuie pronunțat, sau impunem vorbirii sintetizate o inteligibilitate și naturalitate moderate. Cu alte cuvinte, proiectarea sistemului TTS se poate simplifica dacă se impun sistemului sarcini precise, corespunzătoare unor aplicații concrete.

Colectivul nostru de cercetare a început acum câțiva ani dezvoltarea unui sistem complet TTS în limba română, bazat pe *concatenare de difoneme*. Arhitectura acestui sistem este prezentată în Figura 2. Sistemul cuprinde o parte importantă de prelucrare lingvistică și un modul de generare a semnalului de vorbire având la bază un algoritm de tip PSOLA [4]. După realizarea unei prime variante a sistemului, se depun în continuare eforturi pentru creșterea naturalității vorbirii sintetizate, prin îmbunătățirea performanțelor la diferite nivele de prelucrare.

Modulul de prelucrare a limbajului într-un sistem TTS are ca sarcină transformarea textului de intrare într-o reprezentare fonetică și prozodică, care trebuie să descrie cât mai fidel posibil pronunția sa. Acest lucru poate fi realizat parcurgând mai multe etape succesive, puse în evidență și în figura anterioară. Vom discuta în cele ce urmează **modalitățile de rezolvare a părții de prelucrare inițială (preprocesare) a textului în cadrul sistemului nostru de sinteză în limba română.**

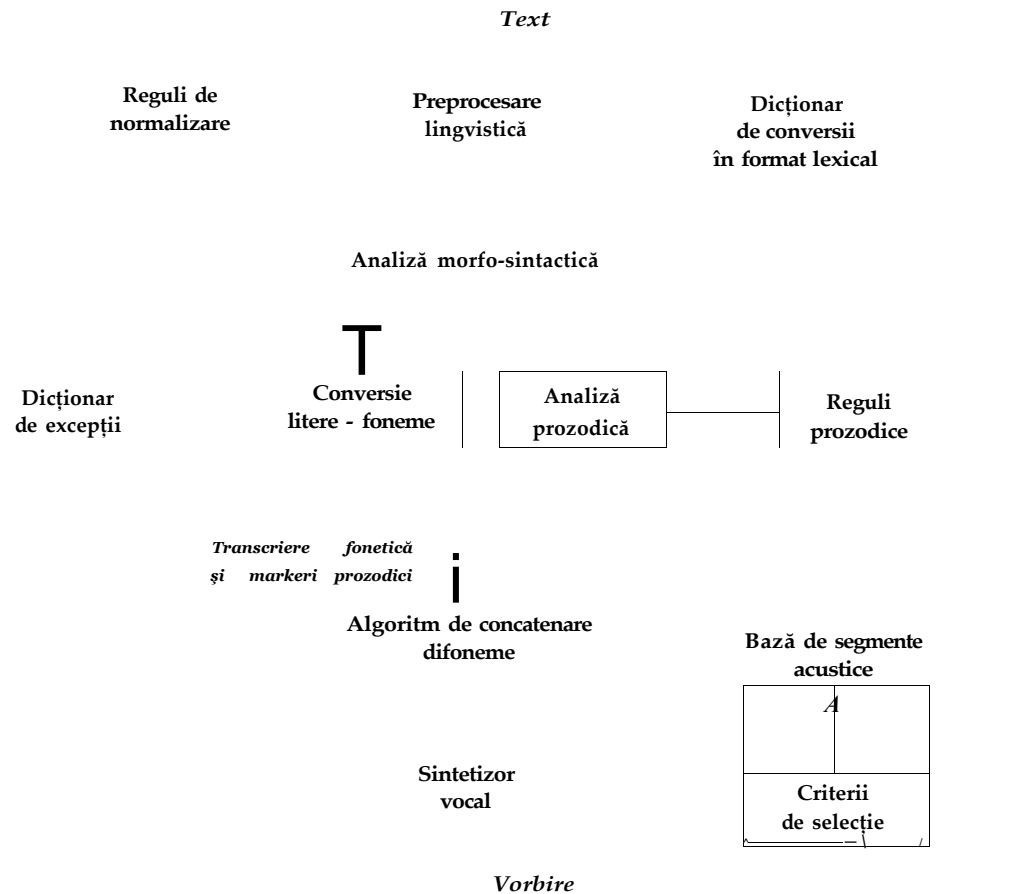


Figura 2. Arhitectura sistemului TTS în limba română

4. Preprocesarea textului de intrare în cadrul sistemului TTS în limba română

4.1 Probleme generale

«ktpm* 0. nrr? 3 difi cu ile ma i o reale sistemelor TTS constă în faptul că aceste sisteme trebuie să poată prelucra practic orice text, plecând de la propoziții simple

izolate și mergând până la paragrafe complexe, care pot cuprinde de propoziții, cu posibile structuri negramaticale și simboluri speciale. Partea de preprocesare lingvistică a textului are un rol extrem de important deoarece detectarea corectă și interpretarea șirurilor de caractere influențează acuratețea întregului sistem de sinteză și contribuie la text fără restricții în vorbire sintetică.

Ușual, un text scris se prezintă sub forma unei secvențe ASCII; el este alcătuit din cuvinte compuse cu ajutorul literelor din alte tipuri de caractere: spații albe, semne de punctuație, șiruri de caractere speciale (de exemplu operatori matematici). Textul poate conține și simboluri speciale (de exemplu operatori matematici), numerale (12, 12.450, 1,245), abrevieri (*prof.*, *dr.*, *ing.*), sau acronime (*TTS*). Aceste secvențe sunt de obicei "anormale" din punct de vedere față de majoritatea cuvintelor din text și trebuie mai întâi transformate în un format ce poate fi recunoscut de partea de analiză lingvistică. Acest proces revine modulului de preprocesare, care trebuie să realizeze o segmentare a textului de intrare (detectarea cuvintelor și a sfârșitului de prelucrare a semnelor de punctuație și a simbolurilor speciale [5,

La prima vedere, preprocesarea unui text pentru un sistem de sinteză pare banală; în realitate însă, lucrurile sunt destul de complicate. Nu este totdeauna posibilă determinarea marginilor unei fraze sau de punctuație. Astfel, punctul (.) poate apărea și la sfârșitul unor multe alte situații, ca de exemplu în abrevieri (*ing.*), acronime (*S*) sau se omite un anumit fragment de text (...), sau numerale (12.4 mii patru sute cincizeci), situații care trebuie diferențiate prin punctuație [9]. De asemenea, cratima creează dificultăți în operația de segmentare și este folosită pentru despărțirea în silabe, pentru scrierea cuvintelor sau delimitarea unui nou paragraf, sau în enumerări.

O sarcină dificilă este și conversia anumitor secvențe de caractere în cuvinte care să poată fi analizate lingvistic. Dacă unele abrevieri sunt "expandate" imediat, cu ajutorul unui tabel de echivalențe, există și simboluri care secvențe de simboluri care nu se pot distinge pe baza caracteristicilor tipuri diferite de conversii; de exemplu, numărul format din șirul "12" poate reprezenta un număr întreg sau un număr de telefon și va fi tratat în cele două situații. În general, prezența șirurilor de numere în text este o dificultate, deoarece ele pot apărea în diferite contexte: ore, date, expresii aritmetice etc.

Trebuie observat că aceste ambiguități create de natura și poziția semnelor de punctuație sau de modul diferit de citire a aceluși simboluri, pot avea implicații majore asupra acurateții întregului proces de prelucrare lingvistică și în final asupra pronunției corecte a textului de către sistemul de sinteză.

Evident, numărul secvențelor de caractere neuzuale dintr-un text ce se dorește a fi transformat în vorbire depinde mult de tipul și subiectul textului. Spre exemplu, textele literare dintr-un volum de proză sau comentariile politice dintr-un ziar au mult mai puține situații dificile decât comentariile economice, sportive, sau prezentările de spectacole. În ultimele situații menționate, construcțiile neuzuale, criptice sau chiar negramaticale, abrevierile uneori ambigui, pot fi atât de numeroase, încât se poate spune chiar că astfel de texte nici nu sunt potrivite pentru o sinteză automată pornind de la text; singura soluție rezonabilă este, probabil, o reeditare a lor pentru a le face mai accesibile unui sistem de sinteză.

Problema enunțată anterior este de fapt mult mai generală. Păreră autorului acestei lucrări este că în orice aplicație TTS trebuie făcut un compromis între calitatea vorbirii sintetizate, dimensiunile vocabularului și complexitatea sistemului de sinteză. Cu alte cuvinte, nu trebuie încercat cu orice preț, prin orice mijloace, obținerea unei vorbiri "perfecte", cel puțin în acest moment.

4.2 Algoritm de preprocesare a textului

Pentru preprocesarea textului de intrare în cadrul sistemului TTS proiectat, am propus un set de definiții, reguli și proceduri, bazate pe o analiză detaliată a situațiilor cele mai întâlnite în limba română.

Definițiile propuse sunt prezentate în continuare.

Definiția 4.1 Vom denumi *expresii* "secvențele de caractere care cuprind una sau mai multe din următoarele categorii: secvențe de litere dintre care cel puțin una este majusculă, secvențe de cifre, semne de punctuație, alte simboluri speciale".

Definiția 4.2 Vom denumi *caractere extra-textuale* "acele semne de punctuație care îndeplinesc în text o funcție de punctuație propriu-zisă".

Definiția 4.3 Vom denumi *caractere intra-textuale* "acele semne de punctuație care fac parte integrantă din expresii și ajută la pronunția lor".

Definiția 4.4 Vom denumi *expandare* "procesul de conversie a unor expresii în format lexical (secvențe de litere alcătuind cuvinte uzuale, ce pot fi analizate lingvistic)".

Definiția 4.5 Vom denumi o secvență de caractere *ambiguă* "dacă ea poate fi încadrată, având în vedere forma sa, în mai multe clase lingvistice".

Pornind de la aceste definiții, am proiectat un algoritm de preprocesare a textului, ce constă în principiu din trei etape de bază:

I. Segmentarea textului

Textul se segmentează de la stânga spre dreapta, în *grupuri de caractere*. Se obțin astfel secvențe de caractere ASCII delimitate de spații albe (blanc); semnele de punctuație se includ temporar în aceste grupe.

II. Conversia șirurilor de caractere de tip expresie în caractere ortografice

Se parcurg pe rând grupurile de caractere rezultate în urma segmentării și se realizează *expandarea* lor (acolo unde este cazul) sub forma unor cuvinte uzuale, pe baza unei analize contextuale simple la nivel de cuvânt sau segment de cuvânt și a unor dicționare de conversie în format lexical (pentru abrevieri și unele tipuri de acronime).

III. Interpretarea unor semne de punctuație

Se detectează și se memorează pozițiile unor *caractere extra-textuale* și a sfârșitului frazelor, pentru a fi folosite ulterior de modulele de analiză sintactică și prozodică.

Detaliind etapa I prezentată anterior și utilizând și definiția 4.1, putem observa că grupurile de caractere rezultate în urma segmentării textului de intrare pot fi de următoarele tipuri [10, 11, 12]:

a. Secvențe de litere alfabetice, scrise cu minuscule

a1. Cuvinte uzuale;

a2. Abrevieri scrise fără punct (de exemplu unități de măsură: *m*, *km*, *ms*).

b. Expresii

b1. Cuvinte scrise cu o singură literă, majusculă: abrevieri (puncte cardinale: *E* - est, *V* - vest; simboluri chimice: *C* - carbon, *O* - oxigen; unități de măsură: *A* - amper, *V* - volt); cifre romane: *V* - cinci, *I* - unu etc.

b2. Abrevieri scrise cu minuscule și puncte (*tel.* - telefon, a.c. - an curent)

b3. Secvențe de mai multe litere, scrise cu minuscule și inițială majusculă

b3.1. Cuvinte la început de frază;

b3.2. Nume proprii;

b3.3. Abrevieri scrise fără punct (de exemplu unități de măsură: *Hz*, *Mw*).

b4. Secvențe de mai multe litere, scrise cu minuscule și o majusculă pe altă poziție decât prima (unități de măsură: *mA*, *kV* etc.)

b5. Secvențe de litere scrise cu mai mult de două majuscule, cu sau fără punct

b5.1. Acronime (*NATO, S.R.L.*);

b5.2. Abrevieri (P.S. - post scriptum);

b5.3. Unități de măsură (*MHz, MByte*);

b5.4. Cifre romane (*VI, IX*).

b6. Secvențe de cifre, scrise cu sau fără semne de punctuație

b6.1. Numere întregi;

b6.2. Numere zecimale;

b6.3. Numerale ordinale (*al 2-lea*);

b6.4. Ore și date;

b6.5. Numere de telefon.

b7. Semne de punctuație: . ? ! : ; . . . , - / ' " () [] { }

b8. Simboluri speciale

b8.1. Simboluri matematice uzuale: + - *(sau x) : (sau /) = < > % ~

b8.2. Alte simboluri speciale: @ \$ &

Deoarece semnele de punctuație ridică cele mai serioase probleme, vom analiza în primul rând situațiile cele mai uzuale de apariție a lor (pe grupe de importanță), precum și soluțiile posibile de rezolvare a acestor situații. Vom discuta apoi câteva aspecte fundamentale legate de grupurile de cifre, abrevieri și acronime.

1. Punctul

Punctul (.) poate apare în abrevieri, acronime, numerale, sau poate semnifica sfârșitul unei fraze. Ambiguitățile create de punct sunt o problemă majoră pentru operația de preprocesare, datorită faptului că el poate reprezenta fie un caracter intra-textual, fie extra-textual, fie ambele în același timp; de exemplu, punctul după abreviere poate marca în același timp și sfârșitul frazei.

Este deosebit de utilă punerea în evidență a câtorva situații de utilizare corectă a punctului în limba română:

- Punctul se folosește în abrevierile provenite din cuvinte simple sau compuse în care nu apare litera finală a cuvântului; exemple: *id.* (idem), *etc.* (etcetera), *tel.* (telefon), *a.c.* (anul curent), *a.m.* (ante

meridian), *d.a.* (după-amiaza), *P.S.* (post scriptum) - deci categoriile **b2, b5.2** puse în evidență anterior.

- Dacă în abreviere apare litera finală a cuvântului, nu se pune punct după abreviere; exemple: *cea* (circa), *dna* (doamna), *dl* (domnul), *dne* (doamnei), *jr* (junior) - categoria **a2**.
- Nu se pune punct după simbolurile unor termeni de specialitate: *C* (carbon), *L* (lungime), *V* (vest sau volt), *mA* (miliamperi), *MHz* (mega hertzi) - categoriile **a2, b1, b3.3, b4, b5.3**.
- în acronime (abrevieri provenite din inițialele unor substantive compuse formate din mai mulți termeni), punctul este facultativ; sunt corecte atât formele *O.N.U.*, *S.U.A.*, cât și *ONU*, *SUA* (categoria **b5.1**).
- Nu se folosește punctul în abrevierile ce reprezintă indicative de state (*RO* - România), sau de județe (*CT* - Constanța) și în situațiile când abrevierea s-a transformat într-un cuvânt sudat, caracterizat prin lectură cursivă (*TAROM*) - categoria **b5.2**.
- Punctul se folosește de asemenea în scrierea unor numere și a datelor: numere întregi sau zecimale (*1.234, 1.234,567*), date (*15.04.2002*) - categoriile **b6.1, b6.2**.

Considerațiile anterioare sugerează următoarea procedură: atunci când este detectat punctul într-un grup de caractere, se cercetează contextul în care apare și apoi se ia decizia corespunzătoare, astfel:

- Dacă există cifre la stânga și la dreapta, el este declarat caracter intra-textual și:
 - dacă mai există un punct în secvența de cifre, secvența reprezintă o dată și se expandează folosind un set de reguli (de exemplu: *15.04.2002* va deveni *cincisprezece aprilie două mii doi*)
 - dacă nu mai există un alt punct, secvența reprezintă un număr și se expandează folosind de asemenea reguli (de exemplu: *1234* va deveni *o mie două sute treizeci și patru*).
- Dacă punctul este în poziție finală și este precedat de alte două puncte (...), această secvență se declară caracter extra-textual, fiind identificată cu semnul de punctuație corespunzător; acest caz îl vom discuta separat.
- Dacă punctul este precedat de o secvență de litere (minuscul sau majuscule) și eventual de alte puncte, se caută într-un dicționar de abrevieri și acronime și:

- dacă grupul de caractere este găsit în dicționar, punctul este declarat caracter intra-textual și secvența se expandează conform echivalenței din dicționar;

- dacă grupul de caractere nu este găsit în dicționar, dar conține majuscule, este un acronim - această situație o vom discuta separat;

dacă grupul de caractere nu este găsit în dicționar și nu conține majuscule și alte puncte, punctul (care este sigur în poziție finală) este declarat caracter extra-textual și va reprezenta sfârșitul unei fraze, poziția sa fiind memorată pentru modulele de analiză sintactică și prozodică.

Ultimele reguli prezentate nu pot însă elimina ambiguitatea situației în care punctul după o abreviere poate reprezenta în același timp și sfârșitul frazei (cazul lui *etc.* este tipic, dar există și numeroase alte exemple).

O soluție ar putea fi cercetarea grupului de caractere ce urmează după blanc, ținând cont de faptul că la începutul unei noi fraze se află de regulă un cuvânt cu inițială majusculă. Această situație nu este însă complet edificatoare, deoarece în limba română majuscula apare ca inițială în multe cazuri: substantive nume proprii de persoană, nume de localități sau denumiri geografice, nume de planete și constelații, nume de instituții, nume de lucrări, nume de evenimente istorice sau de manifestări artistice și științifice, nume de sărbători, ca semn de respect etc.

Este clar că această ambiguitate nu va putea fi rezolvată numai de către preprocesor. Soluția pe care o propunem este următoarea:

- Dacă în urma cercetării contextului din dreapta rezultă că punctul din finalul unei abrevieri ar putea fi în același timp și sfârșitul frazei, punctul rămâne caracter intra-textual (și ajută la expandarea abrevierii), dar se adaugă un simbol special pentru marcarea provizorie a sfârșitului frazei, urmând ca acesta să fie validat sau nu de analiza sintactică ulterioară.

2. Semnele de punctuație ? ! : ; ...

Situațiile cele mai frecvente de apariție a lor sunt următoarele:

- Semnul întrebării (?) și semnul exclamării (!) se folosesc uzual în limba română la sfârșitul frazei. Ele apar foarte rar în interiorul frazelor, când pot reprezenta, de exemplu, considerații personale introduse în text, acestea fiind de obicei puse între paranteze; ca atare, cercetarea caracterului din dreapta lor (blanc sau paranteză) poate diferenția simplu cele două situații.

#

|-

li

%

£

;|

#

* .

t

„

L

~

~

- Semnele : și ; marchează și ele, de cele mai multe ori, finalul unui enunț. Deși nu constituie un sfârșit de frază propriu-zis, pot fi considerate în acest fel în contextul sintezei TTS, deoarece textele din partea stângă și din partea dreaptă se pot pronunța ca și cum ar fi izolate, fără să fie afectată naturalețea pronunției.

Prin urmare, cele patru semne menționate sunt importante în primul rând pentru modulul de analiză prozodică, deci locul lor trebuie detectat și memorat de către preprocesor, iar poziția în frază (finală sau intermediară) este utilă doar pentru a ușura analiza sintactică ulterioară a textului.

- Semnul ... semnifică faptul că se omite un anumit fragment de text (de exemplu finalul neprecizat al unei enumerări); el apare în mod obișnuit la sfârșitul unei fraze, dar poate apare și în poziție intermediară. Putem deci aplica aceeași regulă ca și pentru punctul final al unei abrevieri: cercetarea contextului din dreapta și, dacă este cazul, marcarea provizorie ca final de frază, până la o analiză sintactică mai aprofundată; altfel, el nu modifică prozodia textului,

în toate situațiile menționate, semnele de punctuație vor fi interpretate drept caractere extra-textuale. Există însă și trei excepții, în care semnele / și : au altă semnificație decât cea uzuală; aceste situații pot fi descrise de următoarele reguli:

- Dacă simbolul / se găsește la finalul unei secvențe de numere, el semnifică cu mare probabilitate un "factorial" și va fi transcris ca atare.
- Dacă simbolul : se găsește în interiorul unei secvențe de numere, este considerat caracter intra-textual; secvența reprezintă o oră și se expandează folosind un set de reguli (de exemplu: *14:30* va deveni ora *paisprezece și treizeci de minute*).
- Dacă simbolul: este înconjurat de blancuri, face parte dintr-o expresie matematică și va fi transcris conform dicționarului (*împărțit la*).

3. Virgula

Virgula (,) apare în mod uzual într-o frază în poziție intermediară, la finalul unui cuvânt, dar poate apare și în scrierea numerelor zecimale. Regula aplicată în cadrul algoritmului propus este următoarea:

- Se cercetează contextul în care apare virgula și:

dacă este înconjurată de cifre, se consideră caracter intra-textual; secvența reprezintă un număr zecimal și se expandează folosind un set de reguli (de exemplu: *1,234* va deveni *unu virgulă două sute treizeci și patru*).

dacă la stânga sa se găsește o literă sau un alt semn de punctuație (de exemplu punct după o abreviere), se consideră caracter extra-textual și poziția sa va fi memorată pentru modulul de analiză prozodică.

4. Cratima

Cratima (-) este un semn ortografic ce are în limba română două valori principale:

- *gramaticală*, atunci când servește la scrierea unor cuvinte compuse (*bună-cuviință, nord-vest, prim-plan, pare-mi-se, propriu-zis etc.*)
- *fonetică*, atunci când servește la marcarea pronunțării într-o singură silabă a două sunete din două cuvinte diferite, dar care se găsesc alăturate în vorbirea curentă (*de-a*).

În fapt, deoarece simbolurile uzuale folosite de calculator nu cuprind linii mediane de lungimi diferite, cratima devine practic un semn de punctuație și poate fi folosită atât pentru scrierea cuvintelor compuse sau a unor numerale ordinale, cât și pentru despărțirea în silabe, pentru delimitarea unui nou paragraf, sau în enumerări.

Determinarea caracterului intra sau extra-textual se poate face prin cercetarea contextului în care apare; ea este mărginită de obicei fie de litere, fie de blankuri, dar această informație este utilă doar pentru analiza sintactică, deoarece în mod uzual nu se citește (este suprimată de către preprocesor) și nu modifică prozodia textului. În numeralele ordinale, expandarea se face simplu, pe bază de reguli (*al 2-lea - al doilea*).

5. Bara oblică

Bara oblică (/) are sensul prepoziției "pe" în abrevierile științifice (*km/h - kilometru pe oră, m/s - metru pe secundă*) și în exprimarea unei proporții (*2/3 - doi pe trei*), sau sensul conjuncției "sau" în textele uzuale (*c(e/i) - ce sau ci*) în ambele situații reprezintă un caracter intra-textual. De asemenea, poate semnifica o împărțire în expresiile matematice.

Regulile pe care le propunem pentru simbolul / sunt următoarele:

- Dacă este înconjurat de litere, grupul de caractere din care face parte se caută în dicționarul de abrevieri și:

dacă se găsește în dicționar, se transcrie *pe* și se folosește expresia completă găsită (*metru pe secundă*);

dacă nu este găsit în dicționar, se transcrie *sau*.

- f - Dacă este înconjurat de numere izolate, se transcrie *pe*.
- ț ~ Dacă este înconjurat de secvențe de cifre și alte caractere matematice ($2x3/4x5$), sau de paranteze și secvențe de cifre ($(2+3)/(4+5)$), se transcrie *împărțit la*.
- 1
- i

6. Apostroful

Apostroful (') este folosit în limba română în mai multe situații:

- pentru a reproduce în scris rostiri în care un sunet sau mai multe nu sunt pronunțate; aceste rostiri sunt însă rare, fiind practic neliterare, populare (*pân'deseară*),
- în nume proprii străine sau în neologismele neadaptate (*O'Neill, five o'clock*),
- în scrierea anilor, fără prima sau primele cifre (*'907, '99*).

Regulile pe care le propunem pentru simbolul ' sunt următoarele:

- Dacă se găsește într-o secvență de litere, el este eliminat (nu reprezintă propriu-zis un caracter intra-textual și nu ajută la pronunția cuvântului).
- Dacă în dreapta se găsește o secvență de cifre, în funcție de numărul acestor cifre, grupul de caractere se expandează folosind un set de reguli (de exemplu: *'99* va deveni *o mie nouă sute nouăzeci și nouă*).

7. Alte semne de punctuație: " () [] { }

- Alte semne de punctuație ce pot fi utilizate în textele obișnuite sunt ghilimelele (sau semnele citării) și parantezele rotunde; ele semnifică de obicei un citat, reprezintă porțiuni de text cărora li se dă un sens (stilistic) special sau asupra cărora autorul vrea să insiste, constituie traducerea ori sensul unui cuvânt, sau delimitează considerații personale introduse în text. Apar de obicei în perechi și vor fi declarate caractere extra-textuale, servind modulului de analiză prozodică pentru obținerea unei vorbiri sintetizate cât mai naturale.

Parantezele drepte și acoladele apar extrem de rar în textele românești uzuale; ele pot apare însă (ca și parantezele rotunde) în expresii matematice. Se identifică simplu, deoarece sunt alăturate unor secvențe de cifre și se expandează de obicei prin utilizarea cuvintelor corespunzătoare semnificației lor, cu ajutorul dicționarului de conversii în format lexical.

8. Secvențele de cifre

Secvențe de cifre pot apare și în textele obișnuite, dar mai ales în expresii matematice, împreună cu semne de punctuație sau simboluri matematice; evident,

deoarece numărul lor posibil este practic infinit, ele trebuie expandate pe bază de regului de conversie, în funcție de context.

Am propus anterior o serie de regului pentru cazurile cele mai frecvente (numere întregi sau zecimale, numerale ordinale, ore, date). O situație specială (pe care de asemenea am menționat-o anterior), o reprezintă cazul în care o secvență de cifre, scrisă fără semne de punctuație, poate reprezenta fie un număr întreg, fie un număr de telefon. În acest caz, dacă din cercetarea contextului nu se poate elimina ambiguitatea (de exemplu prezența abrevierii *tel.*), această problemă rămâne în sarcina modului de analiză sintactică, care poate realiza o cercetare contextuală mai amplă.

9. Simbolurile matematice uzuale: + - *(sau x) ; (sau /) = <>% ~

Simbolurile matematice au o situație oarecum privilegiată, deoarece ele sunt încadrate de obicei de blankuri în expresiile matematice uzuale și ca atare pot fi imediat identificate și expandate pe baza dicționarului de conversii în format lexical (de exemplu *plus*, *minus*, *înmulțit cu*, *împărțit la* etc.) Dacă totuși în scrierea expresiei nu apar blankuri, contextul secvențelor de cifre și al celorlalte simboluri duc practic la aceeași rezolvare.

10. Abrevierile

O serie de considerații privind abrevierile au fost expuse anterior la regulile ce privesc punctul. Situația lor este dificilă datorită faptului că în limba română abrevierile se pot scrie în multe feluri: cu majuscule și/sau minuscule, cu sau fără semne de punctuație (uzual punct).

Regula principală ce poate fi aplicată este evidentă:

- Dacă în grupul de caractere apare cel puțin un punct și/sau cel puțin o majusculă, se caută în dicționarul de abrevieri; dacă secvența este găsită, abrevierea se expandează punând-o în corespondență cu cuvântul corespunzător din dicționar.

Pot rămâne însă ambiguități, în special pentru abrevierile scurte (de exemplu V - unitatea de măsură "volt", dar și cifra romană "cinci" și punctul cardinal "vest"), sau pentru abrevierile scrise cu minuscule și fără punct (*km*, *cea*, *dl*), acestea din urmă nefiind căutate în dicționar (după regula expusă). Singurele soluții practice pentru rezolvarea unor astfel de cazuri ambigui este ca ele să fie preluate mai departe de analiza sintactică sau să fie recunoscute la etapa de conversie fonetică, prin căutarea într-un dicționar limitat de excepții.

11. Acronimele

Spre deosebire de abrevieri, cea mai mare parte a acronimelor stocate în dicționar, deoarece pronunția lor nu necesită informații suplimentare. De obicei, pronunția lor se reduce la citirea caracterelor ce compun acronimul, individual (ca pentru *S.R.L.*), normală a cuvintelor, atunci când pronunția lor s-a generalizat în formă compactă (*NATO*, *TAROM*), pentru citirea secvențială a acronimului (ca pentru *S.R.L.* - *serele*). Necesitar doar un set de reguli simple de transcriere a literelor rostite separat (de exemplu *S.R.L.* - *serele*).

Regula propusă pentru acronime este deci următoarea:

- Dacă secvența de caractere cuprinde cel puțin două caractere și este găsită în dicționarul de abrevieri, se caută în dicționarul de acronime:
 - dacă este găsită aici, secvența se expandează în forma echivalenței din dicționar;
 - dacă nu este găsită în dicționarul de acronime, secvența de caractere, puncte, majusculele sunt (eventual) înlocuite cu caracterele din dicționar; dacă nu este găsită în dicționar, dar cuprinde puncte, secvența nu va suferi altă prelucrare (se va citi ca a);
 - dacă nu este găsită în dicționar, dar cuprinde puncte, secvența este expandată secvențial, utilizând un set minim de reguli de transcriere a literelor rostite separat.

Pentru toate situațiile menționate, preprocesorul va serua modulelor ulterioare, pentru o corectă analiză sintactică și prozodică.

5. Concluzii

Am discutat în această lucrare câteva aspecte fundamentale ale sintezei automate a vorbirii, ca și un număr important de reguli generale pe baza cărora a fost proiectat preprocesorul de text pentru în limba română. Nu am urmărit totuși să descriem complet funcționarea și implementarea acestuia; o serie de considerații și totodată modalitatea concretă de implementare (pentru o variantă) au fost prezentate de autor în [13] și [14].

În varianta actuală, preprocesorul de text a fost îmbunătățit pentru a rezolva unele situații dificile legate de abrevieri, numerale următoare etc. De asemenea, un mecanism de automat de transcriere al preprocesorului să fie "tolerant" cu anumite erori tipice de sintaxă, de exemplu fraze ce nu încep cu minuscule, sau un format "ușor" înlocuirea sau numerale.

Putem spune, ca o concluzie a celor discutate anterior, că un preprocesor de complexitate medie, cum este și cel propus pentru sistemul TTS în limba română, poate rezolva cu succes (împreună cu analiza lingvistică ulterioară) o mare parte din problemele întâlnite într-un text obișnuit; el nu poate realiza însă normalizarea completă a **oricărui** text și nu poate soluționa **toate** ambiguitățile care se pot ivi, datorate în special numărului extrem de mare al abrevierilor, acronimelor - în general a secvențelor ne uzuale care pot apare într-un text scris. De asemenea, nu poate face față unor construcții negramaticale (deși, de exemplu, unele simboluri speciale neașteptate sunt ignorate).

Desigur că un set mai mare de reguli și un dicționar de conversii în format lexical mai cuprinzător ar spori eficiența preprocesorului, dar este posibil ca el să devină atât de complicat, încât să fie practic neoperațional pentru un sistem TTS. Singura soluție practică pentru tratarea cazurilor ambigue este folosirea unui set minim de reguli, păstrarea în dicționar a celor mai uzuale situații (cu posibila adaptare a dicționarului la **tipul** textului ce se citește) și examinarea cazurilor rămase la un nivel superior, pe baza plauzibilității sintactice, semantice sau pragmatice a frazelor obținute după preprocesare.

Referințe bibliografice

- [1] R. Boite, H. Bourlard, T. Dutoit, J. Hancq, H. Leich (2000). Traitement de la parole. Presses Polytechniques et Universitaires Romandes, 2000.
- [2] G. Bailly (1996). Pistes de recherches en synthese de la parole - în "Fondements et perspectives en traitement automatique de la parole" (H. Meloni - Coord.), Aupelf-Uref, pp. 109-120, 1996.
- [3] D. Burileanu (1999). Contribuții privind sinteza automată a vorbirii pornind de la text în limba română - Teză de doctorat. Universitatea "POLITEHNICA" București, 1999.
- [4] D. Burileanu (2002). Basic Research and Implementation Decisions for a Text-to-Speech Synthesis System in Romanian Language - Lucrare în curs de publicare în "International Journal of Speech Technology", Kluwer Academic Publishers, 2002.
- [5] G. Fries, A. Wirth (1997). FELIX -A TTS System with Improved Preprocessing and Source Signal Generation - Comunicare la "EUROSPEECH'97", Rodos, pp. 589-592, 1997.
- [6] E. Lewis, M. Tatham (1993). A Generic Front-End for Text-to-Speech Synthesis Systems - Comunicare la "EUROSPEECH'93", Berlin, voi. 2, pp. 913-916, 1993.
- [7] M.Y. Liberman, K.W. Church (1992). Text Analysis and Word F... Text-to-Speech Synthesis - în "Advances in Speech Signal P... Furui, M. Sondhi - Coord.), Dekker, pp. 791-832, 1992.
- [8] A. Lindstrom, M. Ljungqvist (1994). Text Processing with... Synthesis System - Comunicare la "International Conferen... language Processing", Yokohama, pp. 139-142, 1994.
- [9] M. McAllister (1989). The Probiems of Punctuation Ambiguity in... Text-to-Speech Conversion - Comunicare la "EUROSPEECH... 538-541, 1989.
- [10] G. Beldescu (1984). Ortografia actuală a limbii române. B... Enciclopedică, București, 1984.
- [11] T. Dutoit (1997). An Introduction to Text-to-Speech Syn... Academic Publishers, 1997.
- [12] F. Șuteu, E. Șoșa (1993). Dicționar Ortografic al Limbii Româ... București, 1993.
- [13] D. Burileanu (1999). Natural Language Processing for Spee... Romanian Language -Comunicare la "The 12th Internationa... Control System and Computer Science", București, voi. II, pp.
- [14] D. Burileanu, C. Dan, M. Sima, C. Burileanu (1999). A Pa... Preprocessor for Romanian Language TTS Synthesis -... "EUROSPEECH'99", Budapesta, voi. 5, pp.2063-2066, 1999

Utilizarea tehnicilor nuanțate (fuzzy) și de dinamică neliniară pentru sinteza adaptivă a vorbirii

Horia-Nicolai L. TEODORESCU

Academia Română, Secția Știința și Tehnologia Informației,

Calea Victoriei 125, București

E-mail: hteodor@etc.tuiasi.ro

1. Introducere

În timp ce mașina realizează tipic transmisie de date, omul comunică. Diferența constă în participarea intelectuală și afectivă a persoanei la actul comunicării, participare reflectată atât la nivelul limbajelor neverbale (gestică, mimică etc), cât și la nivelul vocal. Această participare afectivă dă varietate, coloratură și sensuri suplimentare, nu neapărat pe plan semantic, semnalului vocal. Sinteza vocii, în prezent, este limitată de lipsa afectului, varietății și sensurilor suprapuse în planuri multiple. Vocea mașinii rămâne astfel cantonată într-o regiune "moartă" a comunicării, este monotună și obositoare pe termen lung.

În această lucrare, reluând unele idei din [1-12], precum și în contextul unor dezvoltări recente [13-27], în special legate de e-Voice și VXML, prezentăm și dezvoltăm unele concepte și tehnici care ar putea permite mașinii atingerea dezideratelor mai sus menționate. Realizarea unor mașini capabile să mimeze calitățile vocii umane și să *dialogheze* cu oamenii, sau măcar să comunice într-o manieră similară cu cea în care omul o face, este un deziderat în numeroase domenii, de la dialogul om-calculator, la sistemele auto și la sistemele de învățare asistată de calculator [13-15]. Rezolvarea acestei probleme are implicații semnificative pentru acceptarea sintezei vocii într-o varietate de aplicații, de la robotică la realitate virtuală, la industria de jocuri electronice și la protezare.

Prozodia, adică structura acustică ce se extinde pe mai multe segmente de semnal vocal, chiar peste mai multe cuvinte sau propoziții, implică ritm, accent, intonație, timbru, afect și alte caracteristici ale vocii încă insuficient înțelese, sau vag definite în literatură. Informația paralingvistică ce este conținută de prozodie nu este nicăieri regăsită la nivelul "spus" prin cuvinte, dar - așa cum am subliniat în [2] - această informație poate fi chiar mai importantă pentru ascultător decât informația lingvistică propriu-zisă. Incapacitatea sistemelor actuale de sinteză

vocală de a reda prozodia naturală este evidențiată chiar de marii producători de aplicații [25] și este bine cunoscută în mediul cercetătorilor în domeniul sintezei vorbirii: *"One of the most difficult problems in speech to date is prosodic modeling"* [25].

2. Soluții pentru sinteza adaptivă și varietală

Cele două calități ale vocii naturale, adaptivitatea - în sens larg - și variabilitatea se pot realiza, cu costuri nu neapărat mari, la nivelul sintetizoarelor actuale, cu adaptări minimale (sau deloc) la nivel hardware și cu îmbunătățiri ale programelor de control. Sinteza adaptivă se referă la adaptarea la:

- Condițiile sonore ambientale [1, 4];
- Contextul semantic-afectiv al cuvintelor și frazelor sintetizate [2, 3];
- Interlocutorul sistemului de sinteză automată, atunci când acesta este recunoscut [2].

Sinteza varietală se referă la modificările inter-pronunție, la repetarea unor fraze, chiar și în cazul în care condițiile ambientale și contextul (și interlocutorul) rămân neschimbate. Această variabilitate elimină monotonia și personalizează vocea (naturală sau sintetizată), în măsura în care variabilitatea se face după reguli adesea proprii individului (cum este cazul în realitate) - și nu doar aleatoare.

Variabilitatea intrinsecă a vorbirii derivă din mecanismele fizice de producere a semnalului vocal (curgere turbulentă a aerului prin organul fonator), precum și din mecanismele neurologice de control al producerii semnalului vocal (controlul neuronal este cunoscut ca având o dinamică cu o importantă componentă neliniară). Aceste caracteristici au fost documentate de mai multe grupuri de cercetare, inclusiv de noi și colaboratorii [5-9].

Adaptabilitatea și variabilitatea în sensurile de mai sus vor fi prezentate sumar în secțiunile următoare, sintetizând lucrările citate și unele cercetări mai noi, nepublicate încă.

3. Adaptabilitate la mediu

Una dintre cele mai elementare adaptări ale semnalului vocal generat de om este cea de adaptare la condițiile de mediu. Adaptarea la un mediu real, cu fond de zgomot, se realizează pe patru căi principale: prin modificarea amplitudinii semnalului (mai mare în mediul cu zgomot ridicat), prin modificarea spectrului (crește contribuția frecvențelor înalte), prin modificarea ritmului (scăderea ritmului, creșterea duratei vocalelor), și prin creșterea duratei dintre cuvinte, care devin

separate, segmentate în timp. Adaptările realizate - instinctiv de un om - se operează deci la un nivel relativ elementar, cu modificări minimale.

Realizarea acestei adaptări este esențială în multe aplicații de sinteză a vocii, incluzând sinteza vocală pentru aplicații în medii industriale și în transport, sau sinteza vocală pentru proteze laringiene. Este remarcabil că adaptarea se poate realiza, la pretenții reduse, cu foarte puțin hardware și/sau cu un soft minimal, aducând însă o îmbunătățire esențială în privința hardului, este necesar un canal de culegere a semnalului (semnal sonor ambiental).

Procesarea semnalului de zgomot, în vederea realizării unui sistemului de sinteză automată, presupune determinarea puterii semnalului ambiental într-o fereastră temporală și determinarea componentelor semnalului ambiental. Primul parametru de caracterizare a zgomotului este media aritmetică a pătratului semnalului s , într-o fereastră dată, de durată eșantioane și caracterizată de momentul actual de timp, n :

$$w$$

$$k=0$$

Caracterizarea spectrală se poate realiza sumar prin raportarea puterii la frecvențe "înalte" (frecvențele înalte corespunzând, în medie, la frecvență ce include formanții nr. 2, 3, 4 și 5 din spectrul vocal) și la frecvențele "joase" (până la aproximativ al doilea formant, deci până la cea. 400 -500 Hz, ținând cont și de vorbitorii feminini):

$$HL = \frac{500}{10000} \int_0^{500} |s(\omega)|^2 d\omega$$

Deoarece parametrii respectivi sunt relaționați cu impactul zgomotului asupra inteligibilității vorbirii, deci sunt dați de calități subiective, este mai ușor să abordăm o definiție probabilistă sau fuzzy a lor. Dată fiind simplitatea termenului "fuzzy", vom prefera a doua variantă. Un exemplu de definiție de apartenență respective este prezentat în Figura 1. Este de presupus că definiția de apartenență respectivă este prezentată în Figura 1. Este de presupus că definiția să constituie doar un punct de plecare, îmbunătățirea realizându-se și prin modificarea funcțiilor de apartenență.

¹ Deși nu este larg acceptat și are o traducere mai dificilă în alte limbi, termenul "nuanțat", propus de Grigore C. Moisil, în locul englezescului "fuzzy".

² Pentru a nu încălca prezentarea, ecuațiile funcțiilor respective sunt date în

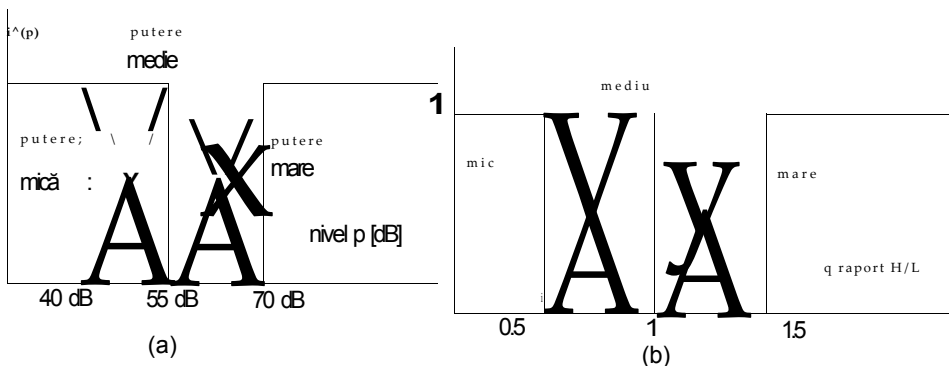


Figura 1. Funcțiile de apartenență ale premiselor regulilor folosite pentru determinarea modificărilor parametrilor de control ai sintetizorului

După cum s-a precizat deja, ca rezultat al aprecierii condițiilor de mediu, se controlează patru parametri ai semnalului sintetizat:

- creșterea amplitudinii (parametru notat AI)
- creșterea conținutului în frecvențe înalte (HFCI)
- creșterea duratei vocalelor (VLI)
- creșterea duratei dintre cuvinte (accentuarea segmentării pe cuvinte a frazei), notat IDBBW.

Controlul se realizează pe bază de reguli și este rezumat în Tabelele 1-4 de mai jos.

Tabelul 1

Creșterea amplitudinii (AI - Amplitude Increase)

HL/P	mic	mediu	mare
mic	0,0	0,1	0,4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelul 2

Creșterea conținutului de frecvențe înalte (HFCI - High Frequency Content Increase - F3 increase)

HL/P	mic	mediu	mare
mic	0,0	0,1	0,4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelul 3

Creșterea duratei vocalelor (Vowel Length Increase - VLH)

HL/P	mic	mediu	mare
mic	0,0	0,1	0,4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelul 4

Creșterea duratei dintre cuvinte (increase of the Duration of the Break Between Words - DBBW)

HL/P	mic	mediu	mare
mic	0,1	0,1	0..4
mediu	0,1	0,3	0,5
mare	0,4	0,5	0,6

Tabelele sunt interpretate în sensul uzual pentru sistemele nuanțate. Preferăm sistemele de tip Sugeno de ordin 0 (vezi Anexa 1), deoarece furnizează ca rezultat, direct, valori numerice, care vor fi interpretate ca și coeficienți de multiplicare a valorilor nominale ale sintezei. De exemplu, prima linie și prima coloană din Tabelul 1 spun că:

DACĂ Puterea (zgomotului) este **medie** și parametrul LH este **mediu**
 ATUNCI Amplitudinea crește de **0,3** ori.

Toate regulile din Tabelul 1 și toate celelalte tabele se interpretează într-un mod similar.

Rezultatul final se obține prin agregarea rezultatelor parțiale, date de regulile respective. De exemplu, dacă valoarea intensității sonore este de 45 dB, iar raportul HL este de 0,7, prin aplicarea fuzificării³ se obține gradul de adevăr al premisei (combinat) din regula respectivă, prin

$$Hlin(\text{putere}^{\wedge}mic, \text{putere}^{\wedge}medie, \text{putere}^{\wedge}mare, LH=mic) = 0,6$$

unde $P_o = 45$, iar $LH_o = 0,7$. Folosind expresiile funcțiilor (v. Anexa 1), se obțin valorile $V_{putere=mic} = \{P_o \cdot h(0,67), \wedge_{LH=mic}(\wedge_o) = 0,6$, deci valoarea minimă este 0,6 și reprezintă gradul de încredere în faptul că amplitudinea crește de 1,1 ori. Aceasta este valoarea de adevăr pentru singletonul (de la ieșirea sistemului) ce corespunde regulii respective, oc[^]. În total, sunt 9 reguli per tabel, deci există 9 valori de singletoni. Într-adevăr, în același timp, valorile de intrare corespund funcțiilor de apartenență „mediu” pentru „putere” și LH, deci regulii:

³ Termenul echivalent românesc ar fi "nuanțare".

DACA Puterea (zgomotului) este mică și parametrul LH este mic
 ATUNCI Amplitudinea crește de 0,0 ori.

cu gradul de încredere în rezultat:

*putere

precum și regulilor:

DACĂ Puterea (zgomotului) este mică și parametrul LH este mediu
 ATUNCI Amplitudinea crește de 0,1 ori.

respectiv:

DACĂ Puterea (zgomotului) este medie și parametrul LH este mic
 ATUNCI Amplitudinea crește de 0,1 ori.

cu gradele de încredere

$\wedge(\wedge putere \wedge mica Mv \quad LH-medie(Mo))$

și respectiv

$\min\{v_{putere=medie(o|V)} \quad LH-mic$

Celelalte cinci reguli din Tabelul 1 au gradele de încredere în rezultat nule, deoarece valorile funcțiilor de apartenență „mare” ale premiselor („puterea este mare” și „LH este mare”) sunt nule, pentru valorile date, $P_o = 57$ și $LH_o = 0,7$.

Prin agregare (defuzzificare), considerată aici conform formulei uzuale:

$$z = \frac{\sum_{i=1}^n \mu_i(x) \cdot y_i}{\sum_{i=1}^n \mu_i(x)}$$

(3)

5 X W

se obține valoarea de ieșire (amplitudinea, creșterea conținutului de frecvențe înalte, creșterea lungimii vocalelor, respectiv creșterea duratei pauzei dintre cuvinte). În relația de mai sus, x reprezintă abscisele singletonilor de ieșire din sistemele tip Sugeno respective, y_i reprezintă gradele de încredere în concluziile regulilor respective, iar z reprezintă valoarea agregată (defuzzificată) de ieșire a sistemului Sugeno. Sumarea se face pentru toți singletonii de ieșire (notați de la 1 la 9). Indicele „A” arată că ne referim la parametrul controlat „amplitudine”, controlului fiind desigur diferențiat pentru cei patru parametri discutați.

Valorile astfel obținute sunt folosite, cum s-a precizat, în calculul de adaptare a amplitudinii și frecvenței, prin înmulțirea cu valorile nominale⁴. De exemplu, dacă amplitudinea nominală este AQ , atunci, prin aplicarea controlului, amplitudinea efectivă a semnalului este:

$$A = A_0 \cdot \prod_{k=1}^9 \mu_k$$

Sistemul de control este instantaneu, în sensul că nu ține cont de valorile recente (din fereastra prezentă, de lărgime W) ale zgomotului și ale valorilor anterioare. Controlul de amplitudine și frecvență se poate realiza pe sintetizorul propriu-zis, asupra unui amplificator și a unui filtru plasat înaintea sintetizorului. Aceste două controale se pot prevedea de altfel și în cadrul precum sisteme de sonorizare mari (eventual distribuite, ca în cazul unor spații mari, gen piețe sau stadioane), sau sisteme de sonorizare (de exemplu, sisteme de interfonie). Controlul pauzelor dintre cuvinte și al spectrului vocalelor necesită comanda directă a sintetizorului.

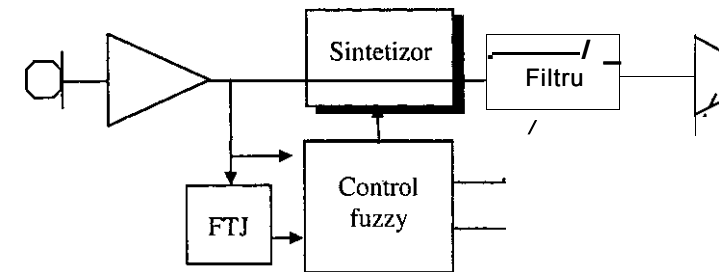


Figura 2. Schema bloc a unui sistem audio adaptiv la zgomotul

în cazul în care se utilizează doar primele două tipuri de control (amplitudine și spectral, adaptarea se poate realiza și cu mijloace hardware înaintea sintetizorului, putând, de altfel, fi utilizată în orice aplicație audio (de exemplu, etc). Schema unui asemenea sistem de adaptare este cea prezentată în figura 2, o variantă fiind inițial propusă în [4].

⁴ Nominale, în sensul că sunt valorile standard pentru sistemul de sinteză audio și pentru sunetul respectiv produs în condițiile contextuale date.

4. Adaptare și variabilitate contextual-interpretativă

Interlocutorul uman răspunde cu afect, după cum consideră anormală, nepotrivită, sau oricum în alt fel "departe de așteptări" întrebarea sau afirmația făcută de partenerul la dialog. De asemenea, răspunsul este diferit atunci când vorbitorul uman este nesigur de răspuns, are un interes special în răspuns sau în topica discuției, sau, din contra, este dezinteresat. În plus, situarea interlocutorului față de partenerul sau partenerii de dialog, în context social sau afectiv, tonalizează discursul verbal și îi imprimă specificitate relativă. Toate aceste caracteristici participative, precum și altele asemenea, dau *comportamentul verbal* al omului, sunt traduse în mare măsură la nivelul semnalului vocal prin prozodie, dar în prezent nu se regăsesc la nivelul mașinii. Privitor la elementele de bază privind prozodia, vezi [26].

Pentru a implementa un comportament verbal, mașina trebuie să dispună de o bază de cunoștințe minimală prin care să genereze acest comportament. De exemplu, este necesar să se interpreteze "departe de normal" într-o aserțiune sau întrebare a interlocutorului uman. Deci, vom presupune că există o bază de cunoștințe care permite o asemenea interpretare. Construcția acestei baze de cunoștințe depinde de domeniul în care se poartă dialogul. În aceste condiții, accentul va fi mai puternic pe anumite părți ale frazei, sau răspunsul va depinde de aserțiune sau întrebare. Modul de răspuns va fi dirijat de asemenea de o bază de cunoștințe, care include regulile necesare modificării sintezei (vezi Figura 3).

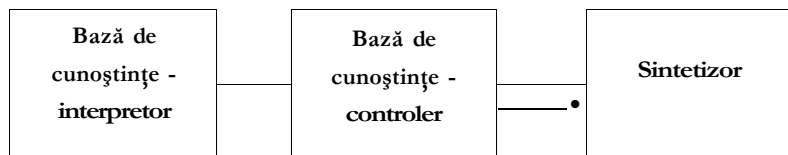


Figura 3. Schema de principiu a controlului contextual-interpretativ

Baza de cunoștințe-controler poate de asemenea fi implementată cu reguli. Dacă... *Atunci*, de exemplu, de forma:

DACĂ oferta / răspunsul interlocutorului este neașteptat (negăsit în baza de cunoștințe - baza de așteptare/baza de cazuri),

A TUNCI afectul sintezei este mirare / neîncredere.../ etc.

ori

DACĂ oferta / răspunsul interlocutorului este neașteptat negativ (conform bazei de cunoștințe),

A TUNCI afectul sintezei este mirare și/sau furie.

Folosind rezultatele regulilor de acest fel, se pot seta parametrii ierarhic inferiori, de tonalitate, ai vocii sintetizate, pe baza acestora generându-se parametrii efectivi de control ai sintezei (amplitudine, frecvențe formanți etc).

Deși acest gen de control poate părea complicat, sunt situații destul de generale în care el se poate implementa cu un efort relativ redus. De exemplu, atunci când se determină (printr-o măsurătoare relativ simplă, de frecvență medie Li — *spectru* vocal, sau de fundamentală) că interlocutorul este un copil sau o persoană de gen feminin, se poate selecta una sau ambele dintre alternativele:

- sistemul de sinteză automată se setează pe o voce de același tip (copil/feminin)
- sistemul de sinteză automată se setează pe voce "caldă" și "vorbire clară".

Utilitatea și modalitatea de realizare a primei setări nu necesită explicații. A doua setare (care poate fi simultană cu prima) se justifică - în cazul interlocutorului • copil - prin necesitatea de a îi crea un mediu afectiv propice și liniștit de dialog (voce "caldă") și prin necesitatea unei comunicări cât mai informative, ușor de urmărit. Pentru a obține o voce "caldă", se pot folosi trasee melodice cu variații lente precum și frecvențe mai joase ale formanților și largimi mai mari (în zona spre frecvențe joase) ale spectrelor formanților. "Claritatea" vocii se poate traduce prin segmentarea mai pronunțată pe cuvinte, precum și vocale mai lungi (cu sau fără accentuări ale spectrelor formanților). Utilizarea unor asemenea adaptări - ce rămân în mare măsură să fie concepute în detaliu, implementate și testate - este neîndoielnic mare la sinteza pentru procese educative [15, 26], în aplicații medicale (răspuns sintetic destinat pacienților), precum și în numeroase aplicații generale (de exemplu, sintetizoare utilizate în muzee, pentru prezentarea exponatelor).

Alte modalități de personalizare afectivă sunt colorarea frecvențială și în amplitudine a anumitor părți din frază sau în cadrul unui cuvânt, aceste modificări locale fiind larg documentate în literatură, de ex. [16-18] și fiind relativ ușor de implementat.

5. Variabilitate prin metoda modulării de către un sistem dinamic neliniar

Variabilitatea semnalului vocal uman este bine cunoscută [5-9], [19-26]. Variabilitatea de tip natural a semnalului vocal sintetizat se poate obține prin modularea diverselor controale (al amplitudinii, lungimii vocalelor, accentului, pitch-ului etc) cu semnale lent variable, generate de sisteme care prezintă dinamică neliniară (haos). Parametrii sistemului haotic respectiv pot modela un anume subiect; considerăm aici că acești parametri reprezintă individul vorbitor și

"personalitatea" lui. Această metodă, propusă de noi inițial în 1992 ([28] ș.a.) dar nepublicată în forma extinsă, credem că reprezintă o metodă promițătoare de "personalizare" a vocii.

Considerăm un sistem dinamic neliniar, dependent de parametri; semnalul în timp generat de acesta este de forma $x(t) = x^{\wedge} X^{\wedge}, \dots^{\wedge}$, unde \wedge reprezintă parametrii sistemului haotic și permit modelarea specificității vorbitorului. Semnalul x poate fi folosit în modularea amplitudinii, frecvenței fundamentale sau spectrului semnalului vocal sintetizat. De exemplu, spectrul poate fi modificat folosind o lege de variație a frecvenței centrale a formanților, de forma-

$$f_j(t) = H_j + x M_j(t) \quad (5)$$

unde $f_j(t)$ este frecvența formantului numărul j la momentul t , $x(t)$ este semnalul haotic respectiv $(x, 0 < 1)$, iar f_{j0} este frecvența "nominală" a formantului respectiv.

Un exemplu simplu de sistem haotic ce poate fi folosit în acest scop este dat de ecuațiile:

$$r_{n+1} = X_1 \cdot u_n + X_2 \cdot u_n \quad (6)$$

unde setul de coeficienți $(X_1, \dots, X_s) \in \mathbf{R}^s$ se alege în domeniul de valori ce corespunde unui comportament haotic al sistemului (vezi Anexa 2). Setul de coeficienți (X_1, \dots, X_s) se poate seta specific pentru fiecare sistem de sinteză automată, "personalizând" sistemul. Valorile de ieșire ale generatorului se scalează corespunzător și se folosesc la modularea unuia dintre parametrii de sinteză. Pentru exemplul din secțiunea 3, amplitudinea semnalului sonor devine, prin utilizarea modulației haotice:

$$1 + 1 = L \cdot (1 + K \cdot T_j) \quad (7)$$

unde K este un coeficient de scalare a seriei de timp r_n . Coeficientul K se alege astfel încât contribuția termenului $K \cdot r_n$ să fie de ordinul procentelor ($K \cdot T_j < 0,1 \forall j$).

Desigur, scara de timp a procesului de generare de eşantioane de semnal vocal diferă de scara de timp a proceselor haotice folosite în modulație, ceasul celui de al doilea proces fiind mult mai lent (de ordinul 1/100) decât al primului proces. Pentru evitarea tranzițiilor bruște ale parametrului controlat, valorile generate pot fi interpolate și se poate realiza o variație lentă între două valori succesive. Considerând că un eşantion al seriei haotice r_n este generat la fiecare Q eşantioane de semnal vocal, seria r_n se poate înlocui cu seria (mai "fină", după ceasul de generare a eşantioanelor semnalului vocal):

$$r_{n/Q}^* = 0, 1, \dots, Q \quad (8)$$

În scopul modulării haotice a mai multor parametri de sinteză (amplitudine, frecvența centrală a formanților, lărgimea formanților, elemente prozodice etc.) sunt necesare mai multe generatoare haotice, câte unul pentru fiecare parametru controlat. Alternativ, se poate folosi un sistem nuanțat (fuzzy) haotic, aceste sisteme generând simultan un număr mare de ieșiri necorelate sau slab corelate [28].

6. Concluzii și discuții

Adaptabilitatea și variabilitatea sistemelor de sinteză a vocii și ale celorlalte audio, în general, se pot asigura prin modificări relativ simple hard și soft ale sistemelor actuale. Adaptabilitatea se poate manifesta atât în raport cu mediul sonor, cât și în raport cu contextul sau cu interlocutorul. Ideea de adaptabilitate și metodele respective au fost introduse de noi în urmă cu peste 20 de ani și dezvoltate continuu în lucrările citate, atât pentru aplicații de uz general, cât și pentru aplicații medicale.

O aplicație de interes medical-educational este utilizarea unor sisteme de învățare a unei limbi pentru copii de vârste mici (1 lună - 3 ani), care suferă de deficiențe de auz. Utilizarea unor sintetizoare cu spectru și amplitudine controlate astfel încât să fie optim adaptate auzului (curbei de sensibilitate audiometrică) a fiecărui copil în parte ar ajuta asemenea copii să învețe limba la această vârstă. Este, într-adevăr, demonstrat că învățarea primelor elemente ale unei limbi la aceste vârste asigură o șansă mult mai mare de învățare a limbii ulterior și de inserare socială [24].

Lucrarea prezentă se situează într-un context mai larg, în cadrul cercetărilor realizate de diverse colective care caută soluții pentru a face vocea sintetică purtătoare de informație emoțională. Astfel, în [31] se descrie o metodă de sinteză a "vocii emoționale", capabilă să transmită trei emoții (supărare-furie, bucurie, tristețe) folosind elemente de prozodie și segmente de tip vocală

consoană-vocală (specifice limbii japoneze). În [32], starea ("mood") și personalitatea sunt văzute ca elemente esențiale apărând în subsidiar în voce și necesar a fi introduse și în vocea sintetizată. Alți autori [33] vorbesc de "nivelul de plăcere al audiției" (pleasantness) - dincolo de inteligibilitate - și văd naturația vocii sintetizate prin această prismă, a utilizării, la un nivel semnificativ, a prozodiei ("...we need to know more about how prosody could be utilized in human-computer interaction. We believe that we could borrow a lot from professional human speakers. Furthermore, speech applications should be built in a way that makes it possible to use prosodic features efficiently.").

Credem că, în viitor, o metodă comodă de a genera automat prozodia, pentru o voce artificială dată și pentru o anumită stare, ar putea fi constituită de o procedură inversă celei descrise în [34].

Încheiem cu un citat din [35]: "... in spite of the long history of speech synthesis, no one speech synthesis system available today is able to produce speech that could be characterized as natural or completely pleasant. In order to improve the speech quality of current text-to-speech (TTS) systems in terms of naturalness, three areas must be addressed: 1) improved linguistic analyses, 2) improved prosody modeling, and 3) improved speech synthesis models."

Mulțumiri. Această lucrare a fost realizată cu sprijinul material al Academiei Române - Institutul de Informatică Teoretică Iași - precum și cu sprijinul material parțial al Societății "Tehnici și Tehnologii" s.r.l. Iași. Autorul mulțumește colegilor Dragoș Burileanu, Bogdan Branzilă și Oana Geman pentru sugestiile și corecțiile la o formă preliminară a lucrării.

Referințe bibliografice

- [1] Teodorescu H.N., Chelaru M., Sofron E., Adăscăliței A.: Adaptive speech synthesis. în voi. *Digitale Sprach-verarbeitung - Prinzipien und Anwendungen*. VDE Verlag, Berlin (W), pp. 183-188, 1988
- [2] Teodorescu H.N.: Interrelationship, Communication, Semiotics, and Artificial Consciousness. în: Kitamura, T. (Ed.): *What Should be Computed to Understand and Model Brain Functions?* FLSI Book Series, voi. 3, World Scientific, 2000
- [3] Teodorescu H.N.: Computer semiotics: understanding meanings and parallel languages (Refereed invited paper), Proc. Int. Conf. IIZUKA'98, Japan, 1998
- [4] Teodorescu H.N.: Making speech synthesizers noise-adapted. *Engineering (UK)*, July 1987, p. 23
- [5] Rodriguez, W., Teodorescu H.N., Grigoras Fl., Kandel A., Bun: information space approach to speech signal nonlinear a. *Intelligent Systems (Wiley)*, Dec. 1999
- [6] Grigoras Fl., Teodorescu H.N., Apopei V.: Nonlinear Analysis of Speech. *Studies in Informatics and Control*, voi. 7, no. 1, M, 57-72
- [7] Teodorescu H.N., Grigoras Fl., Apopei V.: Nonlinear speech production. *Int. J. Chaos Theory and Applications*, voi. 2, no. 2, 52
- [8] Teodorescu H.N., Grigoras Fl.: Nonlinear Techniques in Speech. Proc. International Conference on Intelligent Technologies in Sciences, ITHURS'96. July 5-7, Leon, Spain. Voi. 2, pp. 293-298.
- [9] Grigoras Fl., Teodorescu H.N., Apopei V.: Analysis of nonstationary processes in speech production, IEEE 1997 Applications of Processing to Audio and Acoustics. Mohonk M New Paltz, New York, October 19-22, 1997 (IEEE Catalog # 97
- [10] Burlui V., Teodorescu H.N., Morarașu C.S.: La fonction p l'edente total. Analyse en frequence. *Les Cahiers de Prothese* 88, Decembre 1994, pp. 63-68 1994
- [11] Teodorescu H.N. et al.: Fuzzy models in speech analysis application, în Book of Summaries Int. Conf Modelling and Simulation Turkey, July 1988, voi. 1, p. 162 (Summary)
- [12] Teodorescu H.N., L. Buchholtzer, Chelaru M., Teodorescu prosthesis based on perilaryngean reflexes, Proc. 9th Int. EM Boston. Voi. 4, IEEE, pp. 2114-2115, 1987
- [13] Anonymous Automotive Industry OEM/Supplier: Talking to talking to humans 7/12/2000. <http://www-nrd.nhtsa.dot.gov/13/driver-distraction/Topics013040293.htm#A293>
- [14] Anne-Marie Derouault, The Future of Speech Recognition. E recognition technology is driving transparent computing, making people to interact with computers. <http://www.advisor.ID/OA000107.DERO01>
- [15] House D., Bell L., Gustafson K. & Johansson L. Child-synthesis: evaluation of prosodic variation for an educational program. Proc of Eurospeech'99, pp. 1843-1846, 1999

[16] Heldner M., Strangert E. & Deschamps T.: Focus detection using overall intensity and high frequency emphasis. In: Andersson R, Abelin Å, Allwood J & Lindblad P, eds. Proc of Fonetik 99; pp. 73-76,1999.

[17] Heldner M., Strangert E. & Deschamps T.: A focus detector using overgl intensity and high frequency emphasis. Proc. of ICPHS-99, pp. 1491-1494, 1999.

[18] Heldner M.: On the non-linear lengthening of focally accented Swedish words. In: W. yan Dommelen & T Fretheim, eds. Nordic Prosody: Proc of the VIIIth Conference, Trondheim 2000 . Frankfurt am Main: Peter Lang. 2001

[19] Karlsson I, Banziger T., Dankovicová J., Johnstone T., Lindberg J., Melin H., Noian F. & Scherer K.: Within-speaker variability due to speaking manncrs. Mannell RH & Robert-Ribes J, eds. Proc. of ICSLP98, 2379-2382.1998

[20] Karlsson I: Within-speaker variability in the VeriVox database. In: Andersson R, Abelin Å, Allwood J & Lindblad P, eds. Proc. of Fonetik 99, pp. 93-96, 1999.

[21] Karlsson I, Banziger T, Dankovicová J, Johnstone T, Lindberg J, Melin H, Noian F, Scherer K (1998), Within speaker variation due to induced stress, Proc Fonetik-98, 150-153. www.ling.su.se/fon/publications/fonetik98/

[22] Gustafson-Gapkova S & Megyesi B.: A Comparative Study of Pauses in Dialogues and Read Speech. Proc. of Eurospeech 2001, pp. 931-935, 2001

[23] Beskow J.: A tool for teaching and development of parametric speech synthesis. In: Branderud P & Traunmuller H (eds). Proc. of Fonetik -98, pp. 162-165.1-98,1998

[24] Rachel I. Mayberry, Elizabeth Lock, Hena Kazmi: Linguistic ability and early language exposure. *NATURE*, Voi. 417, 2 May 2002, p. 38, 2002

[25] Microsoft Co.: Platform SDK: Agent. Characters. http://msdn.microsoft.com/library/default.asp?url=/library/en-us/msagent/deschar__8nn6.asp

[26] Mauricio Lumbreras, Gustavo Rossi: Metaphor for the Visually Impaired: Browsing Information in a 3D Auditory Environment. CHT95 Proc, www.acm.org/sigchi/chi95/proceedings/shortppr/mLbdy.htm

[27] Christophe d'Alessandro & Jean-Sylvain Lienard: 5.2 Synthetic Speech Generation. In: Survey of the State of the Art in Human Language Technology. <http://cslu.cse.ogi.edu/HLTsurvey/ch5node4.html#SECTION52>

[28] Teodorescu H.N.: Chaos in fuzzy systems and signals. Voi. Proceedings of the 2nd Int. Conf. on Fuzzy Logic and Neural Networks. Voi. 1, pp. 21-50 (Jono Printing Co., 1992, Iizuka, Japan)

I
%

[29] Teodorescu H.N., Kandel A., Jain L. C. (Eds.), Fuzzy and Neuro-Fuzzy Systems in Medicine (International Series on Computational Intelligence). CRC Press, Boca Raton, USA, 1998.

[30] Teodorescu H.N., Mlynek D., Kandel A. (Eds.): Intelligent Systems and Interfaces (The Kluwer International Series In Intelligent Systems). Kluwer Publ., Boston, 2000.

[31] Yasuhisa Niimi, Masanori Kasamatu, Takuya Nishimoto and Masahiro Araki: Synthesis of Emotional Speech Using Prosodically Balanced VCV Segments. <http://www.ssw4.org/papers/133.pdf>.

[32] Nick Campbell: WHERE IS THE INFORMATION IN SPEECH? (and to what extent can it be modelled in synthesis?) www.slt.atr.co.jp/cocosda/jenolan/Proc/r82/r82.pdf.

[33] Hakulinen J., Turunen, M.: Prosodic Features for Speech User Interfaces. www.cs.uta.fi/hci/spi/reports/Prosodi^.pdf.

[34] Ansgar Rinscheid: Voice Conversion Based On Topological Feature Maps and Time-Variant Filtering. www.asel.udel.edu/icslp/cdrom/vol3/235/a235.pdf.

[35] Syrdal A., Stylianou Y., Garrison L., Conkie A. Schroeter J.: Td-Psola Vs. Harmonic Plus Noise model in Diphone Based Speech Synthesis. www.research.att.com/projects/tts/papers/1998_ICASSP/paperSYN.ps.

Anexa 1: Sisteme nuanțate de tip Sugeno, de ordin 0. Funcții de apartenență

Reamintim că o mulțime (clasică) A c X , unde X notează universul de discurs, este definită de o funcție caracteristică, de forma:

$$\mu_A(x) = \begin{cases} 1 & \text{dacă } x \in A \\ 0 & \text{dacă } x \notin A \end{cases}$$

Prin generalizarea conceptelor de mulțime și de funcție caracteristică, se definesc mulțimile nuanțate (fuzzy) și funcțiile de apartenență corespunzătoare

astfel: o mulțime nuanțată, notată A , peste universul de discurs X , este caracterizată unic de o funcție de apartenență:

$$\mu_A(x): X \rightarrow [0, 1]$$

În particular, funcția de apartenență poate fi de forma:

$$\begin{cases} 1 & \text{pentru } x = a \in X \\ 0 & \text{pentru } x \neq a \end{cases}$$

caz în care se numește *singleton*.

Un sistem de tip Sugeno, de ordin 0, este descris de reguli de forma:

DACĂ intrarea (premisa) # 1 ȘI premisa #2 ȘI... ȘI premisa # n ATUNCI concluzia

unde premisele sunt de forma: x , este \tilde{A}_i , iar \tilde{A}_j sunt valori nuanțate (fuzzy), de exemplu $\tilde{A}_1 = \text{"mare"}$, $\tilde{A}_2 = \text{"mediu"}$, atributelor lingvistice "mare", "mic" etc. fiindu-le atașate câte o funcție de apartenență. Specific sistemelor Sugeno este faptul că în concluzie apar valori numerice și nu valori nuanțate, concluzia fiind deci de forma " $y = 0,3$ " (singleton).

Definițiile funcțiilor de apartenență pentru intensitatea sonoră din Figura 1.a sunt:

$$\mu_{\text{Putere-mic}}(p) = \begin{cases} 1 & \text{pentru } p < 40 \text{ dB} \\ 1 - \frac{p - 40}{15} & \text{pentru } 40 < p < 55 \text{ dB} \\ 0 & \text{pentru } p > 55 \text{ dB} \end{cases}$$

$$\mu_{\text{Putere-medie}}(p) = \begin{cases} 0 & \text{pentru } p < 40 \text{ dB} \\ \frac{p - 40}{15} & \text{pentru } 40 < p < 55 \text{ dB} \\ 1 - \frac{p - 55}{15} & \text{pentru } 55 < p < 100 \text{ dB} \\ 0 & \text{pentru } p > 100 \text{ dB} \end{cases}$$

$$\mu_{\text{Putere-mare}}(p) = \begin{cases} 0 & \text{pentru } p < 55 \text{ dB} \\ \frac{p - 55}{15} & \text{pentru } 55 < p < 100 \text{ dB} \\ 1 & \text{pentru } p > 100 \text{ dB} \end{cases}$$

Definițiile funcțiilor de apartenență pentru raportul HL (Figura 1b) sunt:

$$\mu_{\text{HL=medie}}(q) = \begin{cases} 1 & \text{pentru } q < 0.5 \\ 0 & \text{pentru } q > 1.5 \\ \frac{q - 0.5}{0.5} & \text{pentru } 0.5 < q < 1.0 \\ \frac{1.5 - q}{0.5} & \text{pentru } 1.0 < q < 1.5 \\ 0 & \text{pentru } q > 1.5 \end{cases}$$

Pentru detalii asupra manipulării funcțiilor de apartenență și a regulilor în sistemele nunațate, a se vedea orice manual în domeniul sistemelor fuzzy, sau volume precum [29, 30] în care se pot găsi și aplicații specifice legate de înțelegerea vorbirii, sau alte aplicații medicale.

Anexa 2: Procesul haotic

Procesul reprezentat de ecuațiile (7) are o dinamică haotică doar pentru anumite subintervale relativ înguste din \mathbb{R}^s . În restul spațiului, comportamentul este asimptotic instabil (peste tot pentru valori ale coeficienților lui r mai mari ca 1, în modul, dacă și coeficientul lui u este mai mare ca 1 în modul); comportamentul este stabil sau periodic pentru alte zone, relativ reduse din \mathbb{R}^s .

Diagrama de bifurcație a procesului, așa cum apare în Figura A1, este obținută pentru: valorile coeficienților $[Q]=\{.1, -.17, -.18, .1\}$; $\text{coeff_4} = 1.1$; $\text{coeff_5} = -.15$; condiție inițială $r[0] = 0.3$; număr total de puncte în diagrama de bifurcație: 500 (punctele de la 500 la 1000); regimul tranzitoriu eliminat: primele 500 puncte; precizia tuturor coeficienților și variabilelor: double.

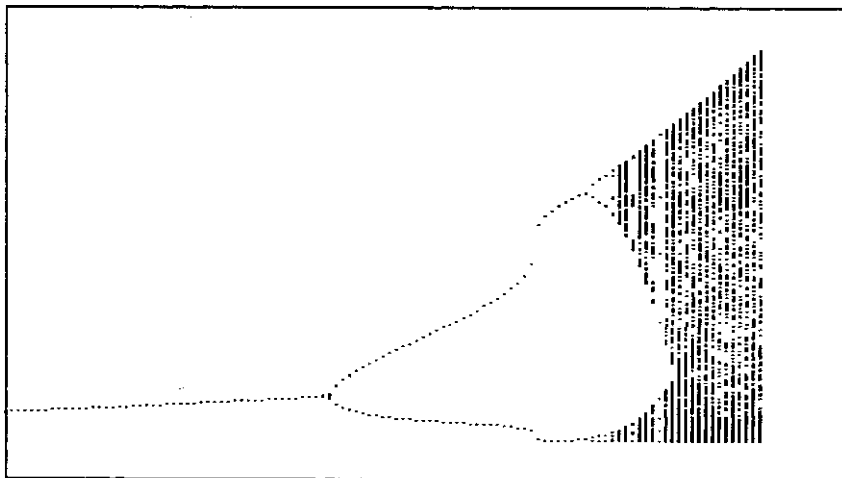


Figura A2-1. Diagrama de bifurcație a procesului

Legile folosite (conform codului, scris în limbajul C) sunt:

```
u[n]= (coefL4)*r[n]+coefL5 -0.005*(float)k;
x=_[n];      r[n+1]=poly(x, Q, coef);
(Q este numărul de valori în vectorul coeficienților, Q=4)
```

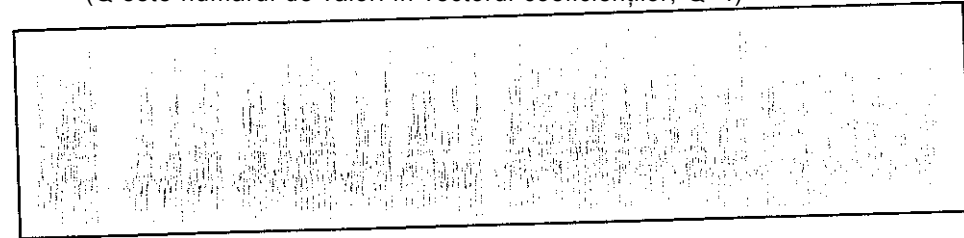


Figura A2-2

Semnalul în domeniul amplitudine-timp din Figura A2-2 a fost obținut pentru ecuațiile (cod C):

```
u[n]= coeff_4*r[n]+coefL5-0.05*21.;
x=u[n]; r[n+1]=poly(x, Q, coef);
```

Semnalul obținut pentru valoarea $u[n] = \text{coeff_4} * r[n] + \text{coeff_5} - 0.05 * 21$. (restul programului fiind identic ca pentru cazul anterior) este ilustrat în Figura A2-3.

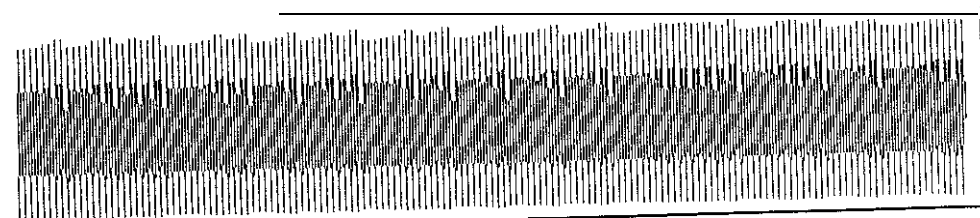


Figura A2-3

Semnalul obținut cu $u[n] = \text{coeff_4} * r[n] + \text{coeff_5} - 0.05 * 20.7$, precum și la o scară dublă de timp, este ilustrat în Figura A2-4:

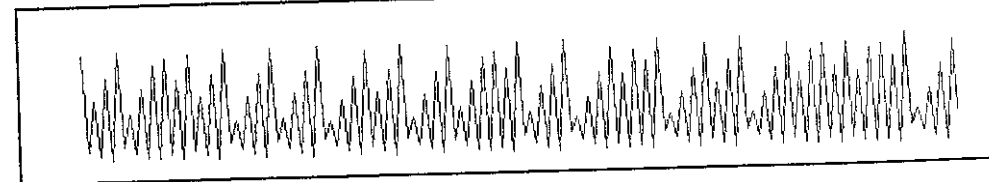


Figura A2-4

Regiunile spațiului parametrilor în care sistemul este stabil, după cum s-a spus deja, sunt relativ înguste. Pentru parametrii coeffJ-coeff_4 fixați și coeficientul coeff_5 variabil între -25.15 și + 4.85 (600 de pași, cu pas 0.05) doar zona îngustă din Figura A2-5 este stabilă, oscilantă sau haotică, în rest sistemul fiind asimptotic instabil. Pentru ușurința urmăririi scării, linia din partea de jos a figuri, reprezintă intervalul menționat, [-25.15, + 4.85], în care s-a testat sistemul

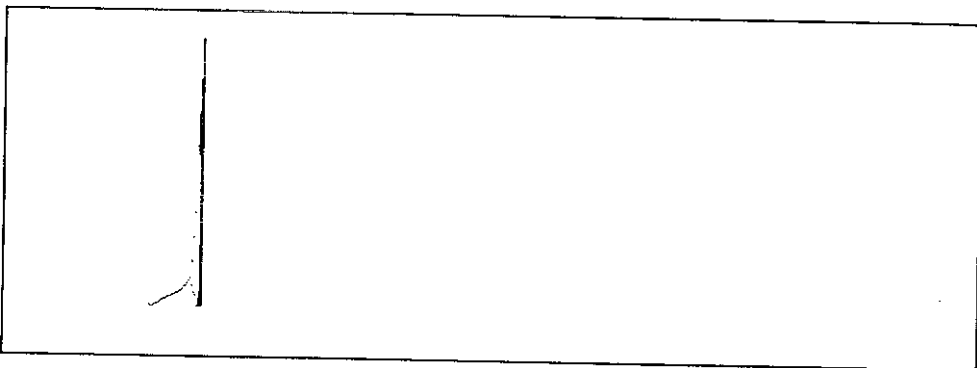


Figura A2-5

În figură, se poate remarca diagrama de bifurcație a sistemului, cu zonele de stabilitate, oscilație și haos. Pentru restul intervalului, prin program calculele sunt abandonate, deoarece valorile de ieșire ale sistemului depășesc, în valoare absolută, 10000.

Dumitru TODOROI, Diana MICUSA, Zinaida TODOROI, Ion LINGĂ, Ion COVALENCO, Nicolae OBJELEANU, Ștefan SPĂTARU, Stela LUNGU, Virginia ȚURCANU, Elana COZLOV, Nadejda AMBROZII, Victor SLOBODEANU, Igor COȘERU, Cătălina SURUCEANU
Academia de studii economice din Moldova, Str. Banulescu-Bodoni, 59-61/503»B», Chișinău MD 2005, Republica Moldova,
E-mail: todoroi@ase.md

Lucrarea actuală realizată în cadrul achiziționării Marelui Dicționar al Limbii Române (MDLR) în format electronic a fost metodologic influențată de ideile subliniate în [1-3] și este o continuare a cercetărilor [4-7,10-11], efectuate în cadrul procesării limbajului natural. Au fost elaborate un șir de proiecte [8-9,12] de informatizare a Limbii Române. Experimentările cu elaborarea sistemelor computerizate de nivelul unu, orientate pe diferite sub-dicționare ale MDLR structurate pe următoarele caracteristici: TEXT, AUDIO, IMAGINI și VIDEO, au început recent în Academia de studii economice din Moldova (ASEM) în colaborare cu ONG-ul ECO-INFO-MOLD. Unele rezultatele ale cercetării și experimentării în cadrul platformei alcătuite din aceste 4 subsisteme, sunt expuse în lucrarea de față. Sunt prezentate diferite scenarii [19] și metodologii de utilizare a sub-dicționarelor informatizate ale limbii române. Clarificarea mijloacelor Hardware-ului și Software-ului modern, care pot suporta MDLR informatizat [18], constituie o problemă importantă pentru etapa creării Societății Informaționale - Societate a Cunoașterii [20].

I. Componenta TEXT a dicționarului economic MULTIMEDIA al limbii Române [23].

Scopul acestui compartiment computerizat al MDLR constă în crearea subsistemului TEXT de nivelul unu a unei părți introductive a dicționarului economic al limbii române și experimentarea cu acest sistem. Acest dicționar economic constă din 35.000 - 40.000 cuvinte. Cuvintele conținute în Dicționarul Enciclopedic Ilustrat (DEI) [21] vor fi definite în-totalitate. Experimentarea cu subsistemul TEXT al MDLR computerizat este efectuată la momentul actual cu circa 200 articole din DEI.

Baza de date TEXT (BDT) a dicționarului economic constă dintr-o culegere de texte-articole, alcătuită din cuvinte, fraze, paragrafe, capitole etc. ale DEI. Documentele în BDT includ nu numai informații textuale (definiții de cuvinte), dar pot conține și informații de alt tip, de exemplu, prin extindere, imagini. Prin urmare BDT în sistemul computerizat al MDLR conține nu numai materialul textual, dar și ilustrativ: diagrame, grafice, fotografii etc.

Prin crearea subsistemului TEXT de nivel unu, utilizatorul obține un mijloc important, prin intermediul căruia informația poate fi introdusă și utilizată în mod complementar pe cale electronică.

1.1. Capacitățile necesare ale unui sistem de gestiune a bazei de date MULTIMEDIA

MDLR este o bază de date MULTIMEDIA. Sistemul de gestiune al MDLR este un sistem de gestiune a bazei de date MULTIMEDIA (SGBDMM) și constituie un mecanism care operează cu diferite tipuri de date, reprezentate într-o diversitate de formate pe un set larg de mijloace și surse. Pentru a funcționa eficient e necesar ca SGBDMM să posede următoarele capacități:

- (a) Capacitatea de a interoga uniform datele reprezentate în diferite formate;
- (b) Capacitatea de a interoga datele, reprezentate în diferite media;
- (c) Capacitatea de a transmite subiectele media din dispozitivele de stocare locale într-un mod eficient;
- (d) Capacitatea de a primi răspunsul la o interogare și de a realiza prezentarea acestui răspuns pe baza componentei audiovizuale;
- (e) Capacitatea de a furniza această prezentare pe o cale adecvată, care ar satisface calitățile diferitor cerințe ale serviciului.

1.2. Structura bazei de date TEXT (BDT) a dicționarului economic MULTIMEDIA

Dicționarul economic, care este pe cale de a fi pus pe calculator, este o BDT cu posibilitatea de a fi extinsă cu diferite componente ale MULTIMEDIA. Subsistemul TEXT a dicționarului economic MULTIMEDIA este un subsistem al SGBDMM, care aprovizionează acesta posibilitate împreună cu utilizarea complementară a BDT.

Structura BDT e compusă din:

- (1) Indice cu caracteristica "număr";
- (2) Termen principal (cuvânt, articol) cu caracteristica "text";
- (3) Variantă(e), derivate, abreviere (concretizare) cu caracteristica "text";

- (4) Categorie gramaticală cu caracteristica "text";
- (5) Domeniu cu caracteristica "text";
- (6) Definiții pentru termenul principal (și concretizări) cu caracteristica "text";
- (7) Sinonim(e) cu caracteristica "text";
- (8) Antonim(e) cu caracteristica "text" și altele.

De asemenea BDT are posibilitatea de a fi extinsă cu următoarele subdiviziuni MULTIMEDIA:

- (9) Audio cu caracteristica „OLE”;
- (10) Imagini cu caracteristica „OLE”;
- (11) Video cu caracteristica „OLE” și altele.

1.3. Scenarii de utilizări și interogări ale subsistemului TEXT al MDLR informatizat

Interogarea este o formă de interacțiune care ajută utilizatorul să prezinte o informație anumită într-o structură anumită, definită de utilizator. Spre exemplu, utilizatorul dorește să obțină informații din arhive, articole, sau alte documente, care conțin informația despre Uniunea Europeană. Interogarea poate avea următoarea formă: "Găsește toate dosarele, legate de investițiile străine, făcute de UE în domeniul educației". Un simplu cuvânt-cheie al acestui dosar nu va permite găsirea răspunsului corect, chiar și dacă indicile acestui document deja există. Totuși, sistemul ne va prezenta unele cuvinte, legate de această interogare, dar este posibil ca acestea să nu poată fi direct asociate la tema dorită. De aceea textul trebuie să fie indexat nu numai pe cuvintele cheie, dar și pe conținutul semantic și/sau pragmatic al cuvintelor (în cazul BDT, de exemplu, concretizarea).

Soluționarea problemei utilizatorului, care dorește să afle definiția cuvântului "Academie", de exemplu, cere introducerea polisemiei în BDT, care conține concepte ca precizia și rechemarea. Întrebarea, propusă de către utilizator în acest context, este: "Cum să aflu din baza de cunoștințe a MDLR sensul cuvântului "Academie - ca instituție de învățământ economic". Pentru aceasta BDT va fi completată cu o nouă coloană "concretizare", care va preciza și va face posibilă afișarea pe monitor a acelei definiții a cuvântului, de care utilizatorul este interesat (de exemplu: *Academia de Studii Economice*).

Un fragment de structură schematică a BDT.

1 Indece cuvânt	1 Cuvânt	1 Concretizare	1	Definiție	Traducere
1 03342	1 Academia	1 de studii economice	1	Nume dat școlii de ...	
1 14269	Banii	EURO		Denumire a princip...!	
J 14271	Banca	de economii			

SGBDMM, ca o extindere a SGBD Ms ACCESS-2000, în baza căruia este creată componenta TEXT a MDLR informatizat, gestionează BDT, utilizând limbajul SQL. În exemplele următoare utilizatorul este interesat de sfera finanțelor. Accesul BDT a MDLR este efectuat prin intermediul următoarelor interogări din SQL (care, în general, constituie comenzile SUMMARING, JOIN, PROJECTION, DIVISION, SELECT și altele):

Ex.1. `SELECT Banii` (termen principal, nume de interes)
`FROM Ambrozii-Godzina` (nume de fișier)
`WHERE Concretizare = EURO` (concretizare pentru termenul principal)

Ex.2. `SELECT Academia`
`FROM Ambrozii-Godzina`
`WHERE Concretizare = de studii economice`

Ex.3. `SELECT Banca`
`FROM Ambrozii-Godzina`
`WHERE Concretizare = de economii.`

II. AUDIO-dicționarul explicativ economic al limbii Române [24]

Dicționarul explicativ economic MULTIMEDIA al limbii române, ca o parte componentă al DEI, include circa 35000-40000 de cuvinte și este divizat în compartimentele: Text, Audio, Video și Imagini. Aceste componente MULTIMEDIA ale MDLR informatizat satisfac cerințele de bază ale unui dicționar informatizat prezintă formele exacte ale cuvintelor, accentul, etimologia, definiția-text, definiția-sunet (audio), definiția-video (film), definiția-imagine (grafic, schema, poza etc.) și corespunde cerințelor unor categorii foarte largi de utilizatori, nu numai elevi și studenți, dar și funcționari și profesioniști, contribuind la ridicarea nivelului de cultură.

Compartimentul AUDIO al MDLR informatizat furnizează informații necesare ale articolului respectiv (cuvântul, definiția lui) în forma AUDIO. Subsistemul AUDIO de nivelul unu al SGBDMM oferă posibilitatea de AUDIO utilizare a dicționarului. Acest AUDIO-dicționar va contribui din plin la ridicarea pe o treaptă superioară a societății noastre în utilizarea corectă a limbajului economic atât la nivel oral cât și scris. Conținutul de date AUDIO poate fi caracterizat prin două metode: (a) folosind metadata prin explicarea conținutului unui fișier AUDIO sau (b) prin extragerea tipului potrivit de date AUDIO, folosind procesorul tehnic.

2.1. Componenta AUDIO a metadatelor

Cu un fișier AUDIO se procedează la fel ca în cazul unei date VIDEO acestui fișier i se asociază un set (grup) de segmente, toate referindu-se la o perioadă de timp. Fiecărui segment i se atribuie un set de activități, care au decurs în acea perioadă de timp, subliniate prin aceste segmente. În general, metadata utilizează reprezentarea AUDIO, care este sesizată ca un set de obiecte marcate în timp.

Utilizarea componentei AUDIO a metadatai din MDLR informatizat este recomandabilă mai ales ca o modalitate de creare și de modificare a acestor metadata și, îndeosebi, la interogarea AUDIO - dicționarului de către utilizatori care necesită această formă de comunicare om-mașină.

Crearea componentei AUDIO a metadatai este mai complexă decât crearea altor forme de dicționare informatizate, deoarece identitatea indivizilor ce vorbesc nu poate fi ușor cunoscută; de asemenea conținutul discursului poate fi neclar.

Conceptul despre conținut este descris în termeni de metadata a procesului. Ca rezultat, data AUDIO este considerată ca un semnal DELTA(x) în timpul x. Trăsăturile de utilizare ale acestui semnal DELTA(x) sunt: (a) extragerea (b) indicarea și (c) depozitarea.

O undă constă dintr-un set de vârfuri (creste) și adâncituri (văi). Perioada vibrației T este definită ca timpul în care o parte a undei revine la poziția inițială.

Alte caracteristici utilizate de componenta AUDIO în crearea metadatei sunt: (1) frecvența, (2) viteza și (3) amplitudinea.

Baza de date **AUDIO (BDA)** poate fi interacționată și gestionată, utilizând sunetul auditiv prin intermediul secvenței de prelucrări: segmentare, memorizare și extragere a informației.

2.2.2. Segmentarea

Segmentarea e o procedură de separare a semnalului audio în câteva ferestre egale. Această procedură poate fi utilizată conform următoarelor două metode:

- Utilizatorul specifică dimensiunile ferestrei, presupunând că proprietățile unde și ale ferestrei se vor obține prin medie;
- Utilizatorul segmentează sunetul în același mod ca și imaginile, folosind predicatul de omogenitate H.

2.2.2. Extragerea

La extragere cel mai des utilizate sunt facilitățile de indicare a intensității, zgomotului, înălțimii și clarității.

2.2. Unele sisteme de utilizare a BDA

Din punct de vedere a MULTIMEDIA, AUDIO - baza de date (BDA) poate fi interpretată ca o sursă auditivă, ca un fișier cu o fereastră auditivă și cu trăsăturile respective, asociate acestei ferestre.

Scenariile de utilizare a BDA cuprind toate formele MULTIMEDIA, care pot fi utilizate în diferite domenii. În sistemele comerciale, de exemplu, **Bazele de date Informix** includ bazele de date a sistemului managerial, care permit utilizatorului să acceseze baza de date, bazându-se pe nesiguranța conținutului.

Baza de date DB2, un alt exemplu, utilizată cu calculatorul de tip IBM, necesită cuplarea cu un sistem auxiliar, care permite lăsarea mesajelor vocale pe robot. DB2 poate importa și menține clipurile, care pot fi căutate printr-un nume sau descriere.

Putem reasculta mesajele lăsate pe robot, prin intermediul **Internetului**. Un exemplu în plus îl constituie o utilizare a unui cuvânt din AUDIO-dicționarul economic al limbii române prin intermediul AUDIO-VIDEO-robotului, care este un sistem autorizat și care acționează pe baza unui program de lucru stabilit sau care reacționează la anumite influențe exterioare.

Un exemplu de interogare a componentei AUDIO a subdicționarului economic al MDLR, prin intermediul limbajului SQL și al subsistemului AUDIO de nivel unu al SGBDMM, poate avea forma:

```
SELECT      Robot
FROM        Țurcan-Mutruc
WHERE       Attribute IS Definiție AND Attribute IS Audio
```

Ca rezultat al acestei interogări utilizatorul prin intermediul răspunsului prietenos, obține pe ecran definiția TEXT a cuvântului Robot și, paralel acest subsistem AUDIO al SGBDMM, difuzează acesta definiție cu voce feminină sau masculină (la dorința utilizatorului).

III. Subsistemul IMAGINI de nivel unu al dicționarului economic informatizat al limbii române [25]

Scopul acestui capitol constă în descrierea posibilităților de introducere a imaginilor în baza de date a MDLR informatizat și de utilizare a acestora în viața cotidiană. Baza de date IMAGINI (BDI) a subdicționarului economic al MDLR informatizat constituie o componentă, care oferă posibilitatea de extindere a procesului de înțelegere a sensului cuvântului dat. Din cele aproximativ 35000-40000 de articole ale dicționarului economic din MDLR doar 50-60%, pot fi prezentate în forma de imagini, după părerea noastră.

Experiența, obținută pe baza câtorva zeci de articole din DEI în cadrul evaluării subsistemului IMAGINI al SGBDMM, ne confirmă întru totul conținutul zicalei: «Mai bine odată să vezi decât să auzi de o sută ori» și al zicalei «O imagine este mai mult decât o mie de cuvinte». Aceste facilități de utilizare din evoluția procesului de creare și utilizare a MDLR informatizat sunt confirmate și de lucrările din [22] precum și prin intermediul următorului Tabel 3.1, prezentat în original

Table 13.1.

Data rates and storage requirements per hour, day, and lifetime for a person to record all the text they've read, all the speech they've heard, and all the

Data type	data rate (bytes per second)	storage needed per hour and day	storage needed in a lifetime
Read text, few pictures	50	200 KB; 2-10 MB	60-300 GB
speech text @ 120 wpm	12	43 K; 0.5 MB	15GB
speech (compressed)	1,000	3.6 MB; 40 MB	1.2TB
video (compressed)	500,000	2 GB; 20 GB	1 PB

3.1. Baza de date IMAGINI (BDI)

Imaginea poate transmite mai multe informații despre un obiect decât câteva pagini (Vezi Tabelul 3.1) de descrieri textuale. Pentru un chirurg este cu mult mai ușor să-și găsească un pacient potențial prin investigarea diferitor imagini. Imaginile pot fi combinate cu corpusuri, text-definiții, sunet-definiții, traduceri etc.

În afară de datele IMAGINI ale dicționarului economic MULTIMEDIA în MDLR informatizat vor fi prezente video, audio, document, manuscrise și altele. Datele VIDEO sunt des folosite în domeniul învățământului. Datele AUDIO sunt importante în domeniul criminalisticii, de exemplu, în identificarea vocilor celor suspectați. Datele documentare diferă de datele TEXT prin aceea că pot conține nu numai informații textuale, dar și imagini încadrate. Datele manuscrise se presupune că în viitorul apropiat vor prevala înregistrările electronice.

Sunt cunoscute diferite formate electronice, care dau posibilitatea de a vizualiza imaginea (fișierele de tip GIF, TIFF, PCX, de exemplu). Subsistemul IMAGINI a SGBDMM are anumite trăsături specifice necesității de utilizare a imaginilor ca o componentă vitală a MDLR informatizat.

3.2. Subsistemul IMAGINI

Baza de date IMAGINI diferă de bazele de date TEXT și AUDIO prin complexitatea imaginilor, necesitatea de a diviza, combina și utiliza diferite părți componente ale imaginii, care deseori la interogare se complică și prin utilizarea incorectă și analiza neprecizată a tehnicilor de manipulare a imaginilor. Aceasta se complică și prin faptul că diferite organizații adună date fotografice, hărți, scheme și alte imagini de tip universal sau specializat (cum ar fi, de exemplu, NASA). Interogările datelor de tip IMAGINI sunt efectuate în baza datelor de tip TEXT, căutate în baza de date de tip IMAGINI și vizualizate în formă de text și imagini. În final imaginile pot fi transferate în baza de date specializate, cum ar fi, de exemplu, încadrarea lor în baza de date MULTIMEDIA comerciale. În subsistemul IMAGINI al SGBDMM este prevăzut un set larg de proceduri cu imaginile.

3.2.1. Plasa imaginii

Conținutul imaginii constă din toate obiectele acestei imagini și caracteristicile lor, care reprezintă interes din punctul de vedere al programului aplicativ. Imaginea poate avea o mulțime de proprietăți, precum descrierea formei, prezentarea vectorului subdiviziunilor, prezentarea vectorului ordinii de descompunere și compunere a imaginii și altele. Fiecare imagine "I" are o pereche asociată schimbătoare de numere pozitive (m,n), care se numește plasa imaginii. Ea este compusă din $m \times n$ celule de măsuri egale.

3.2.2. Transformări de imagini

Imaginea se împarte în părți omogene, care se numesc segmente. Schemele de compresare a imaginii sunt invertibile, deoarece unele scheme de compresare pot conduce la pierderea informației sau la pierderea perfecțiunii. Există două abordări a problemei căutării similarii imaginilor: abordarea metrică și abordarea de transformare.

Abordarea de transformare este mai generală decât abordarea metrică. Această abordare utilizează operațiuni ca: transformarea, transferarea, rotația, scalarea, simetrizarea ș.a.

3.3. Utilizarea imaginii

În prezent multe instituții de învățământ oferă programe de studii individuale. Unele persoane doresc să se specializeze în diferite domenii independent de o anumită formă instituționalizată de învățământ. Astfel de cursuri pot fi reprezentate sub formă de imagini speciale.

Imaginile pot fi utilizate în industria filmelor. Specialiștii au posibilitatea de a vizualiza imaginile alese de ei, lucrând la calculator.

Imaginile sunt importante și în industria turismului. Pentru informații despre imaginile necesare se poate de asemenea apela la sistemul de tip IMAGINE a SGBDMM.

Interogările de imagini în dicționarul economic al MDLR informatizat pot fi efectuate la fel ca în subsistemele, de același nivel unu, de tip TEXT și AUDIO prin intermediul limbajului SQL al SGBD. Rezultatul în forma textuală a articolului și imaginea într-o formă complementară este prezentată utilizatorului în formă de Soft-copy sau Hard-copy.

Obținerea imaginii cuvântului «bancă», de exemplu, în subsistemul IMAGINI al SGBDMM al MDLR informatizat se efectuează prin intermediul următoarelor acțiuni. Se deschide baza de date IMAGINI a dicționarului economic al limbii române (în care sunt acumulate la data experimentării cu SGBDMM al MDLR informatizat doar 25 de cuvinte cu imaginile respective). Se alege cuvântul «bancă». În înregistrarea respectivă a băncii în compartimentul imagini se găsește OLE al imaginii cuvântului ales. Se efectuează clic pe ea și vizualizăm pe ecran imaginea respectivă. Analog se procedează și cu alte cuvinte din BDI.

IV. VTDEO-dicționarul economic al limbii române[26]

În ultimii ani a crescut imens necesitatea de a putea chestiona și procesa cantități mari de date, care nu sunt întotdeauna ușor de reprezentat prin intermediul simbolurilor. Exemple de astfel de date sunt: informația în formă

imagini, informația-video, datele-audio, informația textuală, notițe și altele. În continuare vor fi examinate unele probleme de realizare a dicționarului economic informatizat cu VIDEO clipuri. A fost inițiată baza de date VIDEO (BDV) a dicționarului economic MULTIMEDIA- o subdiviziune a MDLR informatizat - prin crearea subsistemului VIDEO de nivel unu al SGBDMM. Se va demonstra viabilitatea acestui subsistem.

4.1. Problemele creării subsistemului VIDEO al SGBDMM.

Pentru a opera o bază de date MULTIMEDIA (BDMM), un SGBDMM trebuie să posede următoarele abilități:

- Capacitatea de a chestiona uniform datele reprezentate în diferite formate;
- Capacitatea de a chestiona uniform datele reprezentate în diferite surse media;
- Capacitatea de a aporta unitățile media dintr-o diviziune locală de depozitare, asigurând continuitatea acestui proces;
- SGBDMM trebuie să primească răspunsul, generat de o chestionare și să poată genera o prezentare a aceluși răspuns utilizând audiovizualul;
- Capacitatea de a oferi prezentarea într-un mod care ar satisface diferite cerințe ale utilizatorului.

Tehnologiile, legate de bazele de date, au dezvoltat în ultimii 40 de ani baza pe care ar trebui să fie creată o BDMM. În prezent sunt create limbaje de chestionare, tehnicile de aranjare, algoritmi de aportare pentru o mulțime de baze de date de tip relațional, spațial, temporal și altele. Fiecare din aceste mijloace extind posibilitățile limbajelor și algoritmi precedenți pentru a face față noilor tipuri de date sau pentru a dezvolta paradigmele respective.

În acest capitol se va analiza informația de tip VIDEO. Necesitatea de a accesa o bază de date VIDEO (BDV) poate apărea în numeroase aplicații, și de obicei modelul de acces variază considerabil de la o aplicație la alta.

În procesul reprezentării conținutului unui film în BDV este necesar să se răspundă la un set de întrebări de tipul:

- Ce aspecte posibile ale filmului pot interesa utilizatorii BDV?
- Cum pot fi aceste aspecte ale filmului eficient depozitate, astfel încât să minimalizeze timpul necesar subsistemului VIDEO al SGBDMM pentru a răspunde interogărilor utilizatorilor?
- Cum ar trebui să fie limbajul de interogare a datelor VIDEO și cum ar trebui schimbat modelul relațional pentru a corespunde informației VIDEO?

(D) Poate fi oare automatizat procesul de extragere a informației în bază de date în contextul?

Aceste probleme au fost abordate în procesul creării și experimentării cu BDV și subsistemul VIDEO de nivel unu al SGBDMM.

4.2. Definițiile datelor de tip VIDEO

De obicei un film este caracterizat prin personajele sale, atributele acestora și activitățile în care sunt angajate aceste personaje. Principalele surse de interes într-un film includ: (a) oameni, (b) obiecte neînsuflețite, (c) ființe însuflețite și (d) activități.

De observat că tema generală, care se repetă în fiecare cadru, constă în existența unui grup de obiecte și activități asociate. Astfel vom încerca să definim o bază de date VIDEO printr-un șir de definiții.

Definiție 1: O *proprietate* VIDEO este o pereche $\{pname, Values\}$, unde *pname* este numele proprietății și *Values* este o mulțime. O *instanță* a proprietății $\{pname, Values\}$, este o expresie de forma $pname-v$, unde $v \in Values$.

Definiție 2: O *schemă obiect* este o pereche (fd, fi) , unde:
fd este o mulțime de proprietăți cadru-dependente,
fi este o mulțime de proprietăți cadru-independente (*fi* și *fd* sînt mulțimi disjunctive).

Definiție 3: O *instanță obiect* este un triplet (oid, os, ip) , unde:
oid/of este o frază numită identitatea obiectului,
 $os = (fd, fi)$ este o schemă obiect și
ip este o mulțime de afirmații de tip:
 (a) pentru fiecare proprietate $\{pname, Values\}$, în *fi*, *ip* conține cel puțin o instanță a proprietății $\{pname, Values\}$,
 (b) pentru fiecare proprietate $\{pname, Values\}$ în *fd* și pentru fiecare cadru *f* al filmului, *ip* conține cel puțin o proprietate instanță $\{pname, Values\}$. Această proprietate instanță este notată $pname = v \text{ IN } f$.

Definiție 4: O schemă activitate *ACT_SCH* este o mulțime finită de proprietăți astfel încât, dacă $\{pname, Values1\}$ și $\{pname, Values2\}$ aparțin *ACT_SCH*, atunci $Values1 \cap Values2 = \emptyset$.

Definiție 5: O activitate este o pereche, care constă din:
 (a) *AcID*, indicele schemei activitate *ACT_SCH* și

(b) pentru fiecare pereche ($pname$, $Values$) c ACT_SCH este valabilă ecuația de forma $pname = v$, unde $v \in Values$.

Oricărei activități i se asociază o schemă de activitate și fiecărei proprietăți i se asociază o valoare din mulțimea valorilor posibile.

Fiind dată o singură dată VIDEO v , putem defini "conținutul" filmului v .

Definiție 6: Fie că $framenum(v)$ specifică numărul total de cadre din filmul v . Conținutul lui v constă dintr-un triplet (OBJ, AC, \check{A}), unde:

1. $OBJ = \{oid_1, \dots, oid_n\}$ este o mulțime finită de instanțe ale obiectului,
2. $i4C = \{AclDi_1, \dots, AclDi_n\}$ este o mulțime finită de activități/evenimente și
3. A este o hartă de la $\{1, \dots, framenum(v)\}$ până la $2^{OBJ \cup AC}$.

Intuitiv, conținutul unei date VIDEO v este teoretic descris de tripletul (OBJ, AC, \check{A}), unde:

1. OBJ reprezintă mulțimea obiectelor de interes în film,
2. AC reprezintă mulțimea activităților de interes din film și
3. A reprezintă obiectele și activitățile, care sunt asociate cu fiecare cadru f al filmului.

4.3. VIDEO biblioteca

O persoană interesată de obținerea unei lecții imprimabile pe o casetă video ar dori să chestioneze o VIDEO bibliotecă, care găzduiește o colecție de casete video, referitoare la un anumit subiect. De exemplu, Universitatea Maryland oferă cursuri, utilizând contactul prin satelit. În viitor casetele video, create în acest fel, vor putea fi accesate cu ajutorul unui calculator, oferind astfel studenților prelegeri pentru mai multe obiecte de studiu adunate de-a lungul anilor și ținute de diferiți lectori. Chestionarea bazei de date VIDEO de un student individual ar presupune accesarea unui număr foarte mare de casete video.

O bibliotecă VIDEO este o colecție, care specifică: (a) totalitatea filmelor din bibliotecă, (b) conținutul fiecărui film și (c) memorizarea fizică a filmelor.

Definiție 7: O VIDEO bibliotecă $VidLib$ constă dintr-o mulțime finită de cuvinte de tip ($VidContent$, $Vidjd$, $framenum$, R , plm), unde:

- (a) $VidContent$ este conținutul filmului,
- (b) $Vidjd$, este numele filmului,
- (c) $Framenum$ este numărul de cadre în film,
- (d) Plm este amplasarea, care specifică adresele diferitor părți ale filmului și
- (e) R este mulțimea relațiilor despre filme în întregime.

4.3.1. Chestionarea bibliotecii VIDEO

Chestionarea unei VIDEO bibliotecii conține următoarele tipuri de interogări: (a) *aportarea segmentelor* (Găsește toate segmentele care corespund unei anumite cerințe), (b) *aportarea obiectelor*, (c) *aportarea activităților și* (d) *aportarea proprietăților de bază* (Care VIDEO-date sunt în bibliotecă, care este conținutul fiecărei VIDEO-date selectate, unde sunt localizate fizic VIDEO-datele).

4.3.2. Funcțiile VIDEO-datei

Cu bibliotecile VIDEO pot fi definite o serie de funcții:

$FindVideoWithObject(o)$: fiind dat numele obiectului o , această funcție ne oferă tripletul ($Videoid$, $StartFrame$, $EndFrame$),

$FindVideoWithActivity(a)$

$FindVideoWithActivityandProp(a,p,z)$

$FindVideoWithObjectandProp(o,p,z)$

$FindObjectsInVideo(v,s,e)$

$FindActivitiesInVideo(v,s,e)$

$FindActivitiesAndPropsInVideo(v,s,e)$

$FindObjectAndPropsInVideo(v,s,e)$

O chestionare standard a VIDEO-bibliotecii, utilizând SQL are forma:

```
SELECT câmp!,..., câmp_n
FROM relația^RV), relația_(f?2),..., relația_(Rk)
WHERE condiție.
```

4.3.3. Ordonarea datelor VIDEO

O problemă importantă este crearea structurilor informaționale, care organizează bazele de date VIDEO în așa fel încât să optimizeze procesarea celor două funcții enumerate mai sus. Este imposibil de a se depozita conținutul al VIDEO datelor cadru cu cadru, deoarece un singur film de 90 minute conține 162,000 cadre. Astfel, este necesar să se creeze reprezentări compacte ale conceptului de conținut video. În acest sens vom prezenta două astfel de structuri: (a) arborele segment cadru, și (b) arborele R-segment.

4.3.4. Arborii segment cadru

Ideea de bază a arborelui segment cadru este foarte simplă. La început se creează două tabele unidimensionale: $OBJECTARRAY$ și $ACTIVITYARRAY$.

În acest context arborele poate fi creat în 2 etape:

La prima etapă presupunem că $[s^e!], \dots, [s_{w}, e_w]$ sunt toate intervalele în coloana "Segment" a tabelului segment. Fie q_1, \dots, q_r o enumerație ascendentă a tuturor membrilor $\{S_j \mid 1 < j < w\}$. Dacă z nu este o putere a numărului 2, atunci se procedează astfel: fie r cel mai mic număr întreg astfel încât $2^r > z$ și $2^{r-1} > \text{framenum}(v)$. Se adaugă noi elemente q_{z+1}, \dots, q_{2^r} astfel încât $q_{z+r} = \text{framenum}(v) + i$ și $q_{z+j} = q_z + j$ ($j > 0, z+j < 2^r$).

La a doua etapă arborele este unul binar format după cum urmează:

1. În fiecare nod arborele segment cadru reprezintă o secvență de cadru $[X, Y]$.
2. Fiecare frunză este la nivelul r . Prima frunză din stânga marchează intervalul $[z_1, z_2]$, a doua $[z_3, z_4]$ și așa mai departe.
3. Numărul din interiorul fiecărui nod este adresa aceluia nod.
4. Mulțimea de numere de lângă nod marchează numărul de identitate al VIDEO-obiectelor și al VIDEO-activităților, care apar în întreaga secvență de cadru asociată cu nodul dat.

Definiție 8: O secvență de cadru este o pereche $[i, j]$, unde $1 < i < n$ și $[ij]$ reprezintă mulțimea tuturor cadrelor între i (inclusiv) și j .

Definiție 9: O ordonare parțială c asupra mulțimii tuturor secvențelor de cadru este definită ca $[i_1, j_1] c [i_2, j_2]$ cu condiția, că $i_1 < j_1 = i_2 < j_2$.

Definiție 10: O mulțime X de secvențe de cadru este bine aranjată dacă:

1. X este finită (adică $X = \{[i_1, j_1], \dots, [i_{r_2}, j_{r_2}]\}$, pentru oricare r_2) și
2. $[i_1, j_1] c [i_2, j_2] c \dots c [i_{r_2}, j_{r_2}]$

Definiție 11: O mulțime X de secvențe de cadru este solidă dacă:

1. X este bine ordonată și
2. Nu există nici o pereche de secvențe de cadru în X de forma $[h, i_1]$ și $[i_2, j_3]$

4.3.5. Operații cu arborii segment cadru.

Fiecare film v este o structură de VIDEO-date, care constă dintr-un arbore segment cadru, un tablou obiect și un tablou activitate. În particular, dacă biblioteca *VidLib* conține filmele v^1, \dots, v_n , atunci este suficient să asociem următoarele:

1. O singură tabelă numită INTOBJECTARRAY cu schema (VID.ID, OBJ, PTR),
2. O tabelă numită INACTIVITYARRAY cu schema (VID.ID, ACT, PTR) și
3. Pentru fiecare arbore segment cadru V_j , $fst(v_j)$ este asociat cu filmul V_j .

De asemenea pot fi exprimate cele 8 funcții, în SQL, introduse mai sus. De exemplu, una din aceste funcții FindVideoWithObject(o), poate fi implementată cu arborii segment cadru printr-o operație de selecție, efectuată asupra INTOBJECTARRAY DE TIP:

```
SELECT VIDEOJD
FROM INTOBJECTARRAY
WHERE OBJ = o.
```

4.3.6. Arborii R-segment (RS-arbori)

Arborii R-segment sunt foarte asemănători cu arborii segment cadru, cu o singură deosebire. Deși conceptele de OBJECTARRAY și ACTIVITYARRAY rămân aceleași, în locul utilizării unui arbore segment cadru, pentru a reprezenta secvența de cadru, profităm de faptul că o secvență $[s, e]$ este un dreptunghi cu lungimea laturii (e-s) și lățimea 0. Fiecare nod va avea o structură specială pentru a specifica, pentru fiecare dreptunghi, care obiect sau activitate este asociată acestuia.

4.4. Operații cu VIDEO-clipuri

Un film este creat prin filmarea unor secvențe și combinarea lor, utilizând un operator de combinare. O secvență este de obicei filmată de mai multe camere, fiecare având o viteză relativă de rotație constantă. În general o secvență poate avea mai multe atribute asociate precum durata filmării, tipul de cameră utilizat și altele.

Un operator de combinare a filmărilor, deseori numit *edit effect*, este o operație care în baza a două filmări S_1 și S_2 , și a unui interval de timp t efectuează o secvență compusă în timpul t . Așadar un film este creat prin combinarea unor mulțimi de secvențe filmate, utilizând un șir finit de operații de compunere. Exemple de astfel de operații de compunere a filmelor includ:

1. Concatenarea filmărilor,
2. Compoziția spațială și
3. Compoziția cromatică.

4.5. Standardele video

Deși în general standardele industriale nu sunt parte componentă a nucleului cadrului MULTIMEDIA, este important să explicăm în linii generale ideea de bază a standardelor MPEG.

Toate standardele de comprimare a informației VIDEO încearcă să comprime filmele prin executarea unei analize intra-cadru: fiecare cadru este

divizat în blocuri, diferite cadre sunt comparate, pentru a vedea dacă informația conținută de acestea nu se repetă în două cadre. Calitatea tehnicii de compresie este măsurată conform următorilor trei parametri de bază:

- (a) Fidelitatea hărții color: cât de multe culori ale filmului original sunt prezente după comprimare?
- (b) Rezoluția pixel pe cadru: câți pixeli au fost abandonați?
- (c) Numărul de cadre pe secundă: câte cadre au fost abandonate?

4.6. Scenarii de utilizare a VIDEO-dictionarului

Dicționarul MULTIMEDIA al limbii române cuprinde peste 70000 de cuvinte din cele mai diverse domenii. Dicționarul este conceput atât pentru studenți, cât și pentru cercul larg al vorbitorilor limbii române, care doresc să cunoască sensul propriu care trebuie conferit cuvintelor. Dicționarul MULTIMEDIA satisface cerințele de bază: dă definiția exactă a cuvântului și, dacă e cazul, genul, numărul, sinonimele, antonimele, imagini, secvențe VIDEO și AUDIO, care exprimă sensul exact și limpede, deplin accesibil, ceea ce constituie partea cea mai importantă de utilizare. Acest dicționar este una din pietrele de temelie ale culturii tineretului, care va contribui la opera de culturalizare a maselor prin inițierea în folosirea limbii române în mod corect, exact și unitar.

Compartimentul VIDEO al acestui dicționar MULTIMEDIA al limbii române conține, după pronosticurile noastre, peste 12000 cuvinte. Acest compartiment furnizează informații necesare referitoare la cuvintele căutate, secvențe video ce oferă posibilitatea de a percepe mai bine esența cuvintelor. Diviziunea video face dicționarul mult mai accesibil și atractiv pentru utilizatori de toate vârstele și preocupările.

Necesitatea utilizării VIDEO-dicționarului poate apărea în cele mai diverse situații. Să considerăm situația în care un student este nevoit să scrie un referat la merceologia și tehnologia produselor alimentare. Studentul trebuie să analizeze procesul tehnologic de producere a pâinii. În acest sens, apelarea la VIDEO-dicționarul limbii române îi va ușura lucrul; acesta îi va furniza secvențe VIDEO, ce prezintă procesul de fabricare a pâinii, ingredientele utilizate, utilajul necesar.

4.6.2. Chestionarea Video dicționarului

Dicționarul VIDEO este organizat ca o mini-bibliotecă VIDEO. După cum am subliniat mai sus, în procesul de chestionare cele mai importante aspecte sunt:

- (a) Aportarea segmentelor: utilizatorul poate cere bazei de date VIDEO să-i ofere toate secvențele, care conțin informații despre procesul tehnologic de producere a pâinii. O astfel de chestionare ar fi: "Găsește toate secvențele unde se combină ingredientele", sau "Găsește toate secvențele unde se frământă pâinea".

- (b) Aportarea obiectelor: în acest caz, utilizatorul poate solicita toate segmentele, în care este prezent cuptorul, banda rulantă etc. Formularea întrebării ar fi: "Găsește toate secvențele, în care apare cuptorul", "Găsește toate secvențele, în care apare banda rulantă" etc.
- (c) Aportarea activităților: se solicită prezentarea tuturor segmentelor, în care pot fi urmărite diferite operațiuni de producere. Întrebarea poate fi: "Găsește toate secvențele, în care se desfășoară operațiunile de producere".

4.6.2. Utilizarea bazelor de date VIDEO în diferite domenii

După cum am menționat anterior, scopul baze de date VIDEO este de a satisface cele mai diverse cerințe. Astfel, aceste BDV își găsesc aplicarea în cele mai diverse domenii.

4.6.2.1. Educație. Bazele de date VIDEO au o aplicare largă în educație și cercetare. Universitățile pot acorda servicii precum studii la distanță prin satelit sau utilizând Internetul. Acestea pot pune la dispoziția studenților un set de casete VIDEO cu înregistrări ale cursurilor. Dicționarul VIDEO, fiind și el o bază de date VIDEO, pune la dispoziția utilizatorilor secvențe VIDEO care pot fi utilizate în cadrul comunicărilor, pentru pregătirea unor prezentări, lecții deschise, rapoarte.

4.6.2.2. Sport. Sălile de Sănătate oferă baze de date, în care sunt înregistrate casete VIDEO ce conțin diferite programe de antrenament utilizatorului oferindu-i-se posibilitatea de a alege între programe de slăbire, fortificare sau menținere a condiției fizice.

4.6.2.3. Agricultură. Institutele de cercetări științifice în domeniul agriculturii din țară ar putea utiliza VIDEO-dicționarul pentru a studia mai aprofundat procesul de plantare, condițiile de creștere și dezvoltare a plantelor, specificul dezvoltării plantelor în diferite regiuni sau țări, aclimatizarea plantelor în condițiile țării în cauză.

4.6.2.4. Economie. VIDEO-dicționarul poate fi utilizat în foarte multe domenii ale economiei: finanțe, contabilitate, management, marketing, statistică, turism. Vocabularul economic cuprinde destul de mulți termeni, care pot fi redăți printr-un limbaj VIDEO mai accesibil, atât specialiștilor, cât și utilizatorilor obișnuiți.

V. Concluzii

5.1. Compartimentul TEXT. Dicționarul economic TEXT al limbii române în forma sa de BDT, ca o subdiviziune a MDLR, are posibilitatea de a fi extins cu caracteristicile respective ale MULTIMEDIA: Imagine, Audio, Video etc. Acest BDT va ocupa aproximativ 18 MB memorie. La conferința tinerilor savanți ai ASE din 4-5 aprilie 2002, pe baza câtorva sute de articole din DEI au fost demonstra-

caracteristicile de utilizare prietenoasă a subsistemului TEXT al SGBDMM, utilizând sistemele Ms ACCESS - 2000, Ms WORD - 2000 și Ms PowerPoint - 2000 în calitate de componente ale Software-ului Ms OFFICE -2000 și WINDOWS - 2000, exploatate în baza hardware-ului de tip PC Pentium II, conectat la rețelele Intranet, Extranet și Internet.

5.2. Subsistemul AUDIO. Subsistemul AUDIO interacționează cu celelalte subsisteme de nivel unu (TEXT, IMAGINI, VIDEO) ale SGBDMM, care susține evaluarea Marelui Dicționar al Limbii Române informatizat cu MULTIMEDIA. Acest subsistem AUDIO susține toate definițiile celor 61635 de articole din DEI de comun acord cu subsistemul TEXT al SGBDMM. Cele 2320 de ilustrații din DEI sunt susținute de componenta IMAGINI a SGBDMM, dar cu ele poate fi extinsă componenta TEXT și/sau componenta AUDIO. Exemplele, enumerate mai sus, de utilizare a AUDIO componente a MDLR informatizat, au un aspect comun, abstract vorbind formează corpul unei date, fiind individual executate în diferite probleme prin intermediul diferitor suporturi ale Software-ului și Hardware-ului modern. Baza de date BDA a compartimentului AUDIO-dicționarului economic al MDLR informatizat va ocupa un volum de memorie de circa 60 GB memorie.

5.3. Subsistemul IMAGINI. BDI al subsistemului IMAGINI al MDLR informatizat recent a fost expusă pentru analizare și discuții la Conferința tinerelor cercetători ai ASEM din 4-5 aprilie 2002 în baza câtorva zeci de articole din DEI. Mijloacele Software-ului și Hardware-ului de tip Ms ACCESS-2000, Ms WORD-2000 și Ms PowerPoint-2000 cu dispozitivele respective al PC-ului Pentium II au fost suficiente la etapa inițială pentru a demonstra eficiența și eficacitatea mijloacelor și metodelor alese pentru realizarea Proiectului "Limba Română - Limba a Comunității Europene" de către grupul de cercetători - autori ai acestei lucrări. Volumul BDI de prezentare în Ms ACCESS-2000 fără comprimare a 50 articole din DEI ocupă circa 550 MB memorie.

5.4. Subsistemul VIDEO. După cele menționate mai sus ținem să subliniem, că subdicționarul VIDEO are o utilitate mare pentru persoanele ce operează în diferite domenii ca: economia, educația, sport, agricultură, industrie etc. Avantajul acestui dicționar este ușurința în folosire și accesibilitatea. Dicționarul VIDEO este o bază de date, cu care putem opera oricând avem nevoie și oferă posibilitatea de a percepe o informație în formă de videoclipuri. În acest mod, persoanele ce se folosesc de acest dicționar înțeleg mai ușor sensul cuvântului, care este reprezentat în formă VIDEO, fiindcă se formează o imagine amplă despre cuvântul dat și este ușor de memorizat.

5.5. Lucrări paralele și perspective. Paralel cu sistemele de nivel unu, sunt elaborate sistemele de nivelul doi, care suportă subdiviziunile MDLR în planurile: TEXT&AUDIO, TEXT&IMAGINI și TEXT&VIDEO.

Elaborarea sistemului, care suportă toate compartimente MULTIMEDIA MDLR informatizat, constituie a treia platformă, mai complexă, de experimentare și implementări ale dicționarilor computerizate în cadrul elaborării MDLR informatizat [17].

Rezultatele evaluării preliminare ale primelor elemente ale acestor platforme: sistemele unare TEXT, AUDIO, IMAGINI și VIDEO au creat posibilitatea de a transforma unele concluzii ale evaluării MDLR informatizat ca o componentă a cercetărilor în cadrul Proiectului «Limba română - limbă a Comunității Europene», care se desfășoară în perioada 2000-2006. Acest proiect a fost inițiat [10-11] de către Forumul Internațional din Chișinău, 14-15 aprilie 2000. Proiectul constituie unul dintre subiectele de cercetare, experimentare și evaluare, efectuate în cadrul Consorțiului Uniunii Latine «Pentru limba română» și Consorțiului «Pentru informatizarea limbii române» și a Comisiei Academice Române «Pentru informatizarea limbii române».

O serie de aplicații a MDLR computerizat este evidențiată în [13-16].

Referințe bibliografice

- [1] V. S. Subrahmanian. Principles of Multimedia Database Systems. // Kaufman Publishers, Inc., San-Francisco, California, USA, 1998, -pp. 44
- [2] D. Todoroi, S. Nazem, T. Jucan, D. Micusha. Transition To A Full Information Society: Stage Development. // Working Paper No. 98-2, UNO, Omaha, March 1998.-38 p.
- [3] D. Todoroi, D. Micușă, V. Clocotici, I. Lingă, V. Tapcov, N. Drucioc, A. C. M. Morari. **Data Bases** and Communications Tools. Ms ACCESS - 2000 ASEM, Chișinău 2002, 337 pages. (Eng.)
- [4] Dumitru N. Todoroi, Zinaida Todoroi, Diana Micusa. Romanian Computer Language - One of the European Community Languages. // Proceedings of the 26th Annual Congress of the American Romanian Academy of Arts and Sciences (ARA), Montreal, Quebec, Canada, July 25-29, 2001, pp. 1-10 (Rom)
- [5] Diana D. Micusha, Dumitru Todoroi. Natural language processing and transition to a full information society initial development phase. Part 1. și cercetări economice. Voi. XXX. Lucrări prezentate la Sesiunea jubileu de comunicări științifice: «Creștere economică, dezvoltare, progres» Napoca, 2001, pp. 1396-1413.'
- [6] Diana D. Micusha, Dumitru Todoroi. Natural language processing and transition to a full information society initial development phase. Part 2. și cercetări economice. Voi. XXX. Lucrări prezentate la Sesiunea jubileu de comunicări științifice: «Creștere economică, dezvoltare, progres» Napoca, 2001, pp. 1414-1427.'

- [7] Sabin-Corneliu Buraga, Dumitru Todoroi. Adaptabilitatea informațională și operațională. // Studii și cercetări economice. Voi. XXX. Lucrări prezentate la Sesiunea jubiliară de comunicări științifice : «Creștere economică, dezvoltare, progres», Cluj-Napoca, 2001, pp. 1447-1457.
- [8] Dumitru TODOROI. The Computerized Romanian Natural Language Processing Development-Projects-Perspectives. // INFORMATION SOCIETY. The Proceedings of the 5th International Symposium on Economic Informatics, May 2001, Ed ECONOMICA, Bucharest 10-13 May 2001, pp. 927-935.
- [9] Dumitru N. TODOROI. IEE-2000 PROJECT: Natural Language Processing Initialization. // EUROPEAN EXCELENCE IN BUSINESS STUDIES STUDENTS' EDUCATION. International Symposium. Edited by IOAN ANDONE, București, Editura Economica, 2000, pp. 328-334.
- [10] Dumitru Todoroi. Project: Romanian Language - One of the European Community Languages. // Proc. of the VI Conf. « Application Sciences», 18-19 May 2000, USAM, Chișinău, pp. 12-15.
- [11] Dan Crisrea, Dumitru Todoroi, Dan Tușiș. Computațional Linguistic: Romanian Language - One of the European Community Languages. // Proc. of the Intern. Sc. Seminar "Strategies and Modalities for România and Moldova' European Integration", 28-29 Sept. 2000, V.2, ASEM, Chișinău, pp.276-280.
- [12] D. Todoroi, D. Micusa, V. Clocotici, S. Pereteatcu, V. Bordeianu, C. Grigoras, S. Cretu, I. Lingă, S. Spataru. Natural Language Processing: IEE-2000 Project. // Proc. of the Intern. Sc. Seminar "Strategies and Modalities for România and Moldova' European Integration", 28-29 Sept. 2000, V.2, ASEM, Chișinău, pp.281-285.
- [13] Ștefan Spataru, Dumitru Todoroi. Distance Education Via Internet, Multimedia and modern System Environment. // Proc. of the Intern. Sc. Seminar "Strategies and Modalities for România and Moldova' European Integration", 28-29 Sept. 2000, V.2, ASEM, Chișinău, pp. 307-312.
- [14] Ion LINGĂ. IMPACTUL IMPLEMENTĂRII COMPUTERULUI ASUPRA PROCESULUI DE ASIMILARE A CUNOȘTINȚELOR. // Proc. Of the 27th ARA Congress, May 29 - June 2, 2002, Oradea, România.(To be published).
- [15] Ion COVALENCO. Metode adaptabile de evaluare a cunoștințelor asistată de calculator. // Proc. Of the 27th ARA Congress, May 29 - June 2, 2002, Oradea, România. (To be published).
- [16] Nicolae OBJELEAN. The Metod for Error Corection in String with Applications in Speach Recognition. // Proc. Of the 27th ARA Congress, May 29 - June 2, 2002, Oradea, Romania.(To be published).
- [17] Dumitru N. TODOROI, ASEM, Chișinău, Nicolae MARGINEANU, L'Ecole Polytechnique, Montreal, Canada. THE ROMANIAN LANGUAGE 'MULTIMEDIA - DICTIONARIES IMPLEMENTATION ENVIRONMENT AT THE FULL INFORMATION SOCIETY INIȚIAL DEVELOPMENT PERIOD. // Proc. Of the 27th ARA Congress, May 29 - June 2, 2002, Oradea, Romania.(To be published).
- [18] Diana MICUSHA. Mijloace adaptabile ale sistemelor de procesare a limbajului natural computerizat. // Proc. Of the 27th ARA Congress, May 29 - June 2, 2002, Oradea, România.(To be published).
- [19] Zinaida TODOROI, ULIM, Chișinău, Eugenia MARGINEANU, L'Ecole Polytechnique, Montreal, Canada. MULTIMEDIA - dictionaries for Romanian Language. Usage Scenarios on the EAPEC Base. // Proc. Of the 27th ARA Congress, May 29 - June 2, 2002, Oradea, România.(To be published).
- [20] Societatea informațională - Societatea cunoașterii. Concepte, soluții și strategii pentru România. // ACADEMIA ROMÂNĂ, Editura EXPERT, București, decembrie 2001. - 541 pages.
- [21] Dicționar Enciclopedic Ilustrat (DEI). // Editura CARTIER SRL, Chișinău, Editura CODEX SRL, București, 1999, 1808 pages.
- [22] Beyond Calculation : The Next Fifty Years of Computing. // Edited by Peter Denning and Bob Metcalfe, Copernicus, 1997 Springer-Verlag New York, 350 pages.
- Comunicări la Conferința tinerilor cercetători ASEM, 4-5 aprilie 2002
Chișinău.**
- Coordonator: Dumitru TODOROI, Prof. Univ., doctor habilitatus.
- [23]. AMBROZII Nadejda, GODZINA Irina. Componenta Text a Audio Dicționarului Economic al Limbii Române.
- [24]. TURCANU Virginia, MĂTRUC Carolina. AUDIO-DICTIONARUL EXPLI
ECONOMIC AL LIMBII ROMÂNE
- [25]. COZLOV Elena, BABĂNU Irina. Subsystemul IMAGE al dicționarului economic informatizat al limbii române.
- [26]. LUNGU Stela, CIOBANU Diana, GUZUN Oxana. VIDEO-dicționarul economic al limbii române.

Mediu pentru editarea transcrierilor fonetice în limba română. Realizarea atlasului lingvistic român pe regiuni

Silviu BEJINARIU, Vasile APOPEI, Mariana ROMAN
Academia Română, Institutul de Informatică Teoretică, Iași, B-dul Carol nr. 8
silviub@academie.is.edu.ro, vapopei@academie.is.edu.ro

Abstract

The goal of our work is to create an Electronic Linguistic Atlas of România. The Electronic Linguistic Atlas has features of a multimedia application allowing the user to consult and/or print the linguistic maps and to listen audio recordings or synthesized speech.

In order to show all the spelling variations, the phonetically transcription is used in the linguistic atlases. For the Romanian Language, the graphic symbols have been hand-written.

The editing process is too difficult using a standard text editor as consequence of the great number of fonts used. In this paper we propose an editing interface for the phonetic transcription of the Romanian Language. This interface can be used to edit dictionaries of the Linguistic Atlas and as editing tool for the phonetic transcriptions in stand-alone mode or as server for other text editors.

Keywords: dictionary, phonetically transcription, multimedia, linguistic atlas

1. Clasificarea simbolurilor grafice pentru editarea transcrierilor fonetice

Pentru a putea arăta toate nuanțele de rostire, în lingvistică se recurge (după practica internațională) la transcrierea fonetică. Pe lângă transcrierea fonetică internațională realizată cu Alfabetul Fonic Internațional (IPA), fiecare țară își are propriile simboluri grafice [1], [2]. La realizarea atlaselor lingvistice românești, aceste simboluri sunt scrise doar manual. În lucrarea [3] este

prezentată o primă abordare a realizării simbolurilor grafice pentru transcrierea fonetică din perspectiva realizării variantei computerizate a atlaselor lingvistice românești.

În această primă parte vom prezenta principiile care au stat la baza modului în care au fost organizate simbolurile grafice folosite în transcrierea fonetică a limbii române.

Pentru claritatea prezentării introducem următoarele noțiuni:

- sunete primare¹:
- vocale, consoane - existente în alfabetul latin care au corespondent pe tastatură;
- diacritice - vocale, consoane - care nu au corespondent pe tastatură dar pot fi obținute prin combinații de taste;
- sunete marcate cu unul sau mai multe fenomene fonetice.

De aici a rezultat necesitatea realizării unui font de bază (*ALRJBaza*) care să cuprindă simbolurile grafice pentru toate sunetele primare. Poziția în "font" a simbolurilor grafice pentru diacritice, a fost stabilită urmărind păstrarea poziției implicite din familiile de fonturi uzuale (*Arial*, *Times New Roman*). Pentru realizarea sunetelor marcate cu unul sau mai multe fenomene fonetice am proiectat familii de fonturi ale căror denumiri le-am format folosind denumirea fenomenelor fonetice aplicate (ex. *ALRJSemivocale*, *ALRJNazalizate*, *ALRJSeminazalizate*, *ALRJScurteNazalizate*, *ALR_* etc). Această organizare a fonturilor a fost făcută cu scopul de a permite scrierea textelor cu transcrieri fonetice cu orice editor de text (*Microsoft Word*), iar textul scris cu aceste fonturi să poată fi citit chiar dacă fonturile proiectate de noi nu sunt instalate (în acest context se vor pierde numai fenomenele fonetice aplicate sunetelor primare).

Pentru generarea acestor fonturi am folosit programul *FontLab 3.1* care permite definirea de simboluri grafice compuse, pornind de la o familie de fonturi *TrueType* existentă în sistemul de operare Windows. Pentru familiile de fonturi pe care le-am realizat am convenit să folosim ca model de plecare fontul *ARIAL*.

Facem precizarea că fenomenele fonetice și modul lor de aplicare este diferit pentru cele două tipuri de sunete: vocale și consoane.

1.1. Fenomene fonetice aplicate vocalelor primare

simple	diacritice				
a	â	ă	â	a'	â
e	e	e			—
i			î		î
o	o	o			
u	u	u			

Cu ajutorul acestor "vocale primare" și al celor trei variante accentuate (a - â - â - â) . | . , fiecareia dintre ele se obține seria completă de sunete vocale care se regăsesc în fontul de bază *ALRJBaza* (17*4=68 grafeme).

Fenomenele fonetice care pot modifica cele 17 vocale de bază (împreună cu variantele lor accentuate), sunt clasificate în următoarele grupe de fenomene disjuncte²:

Grupe	Poziționare	Notație	Fenomen	Exemplu
[Durată	Așezat cel mai sus	(a)	[Scurtime	
		(b)	Semilungime	
		(c)	Lungime	e e e
Nazalizare	Așezat deasupra vocalei, dar stih fenomenele (a)-(c)	(d)	iSeminazalizare	e e l
		(e)	[Nazalizare	- e A
Ocluzie glotală	așezat "în umăr", în fața		Coup de glotte	e e e
Deschidere	Așezat imediat sub vocală	(9)	închidere	'e 'e 'e
		(h)	Semideschidere	e â e
		(l)	Deschidere	e ? ?
		(l)	Deschidere mare	f f
Afonizare	Așezat sub vocală dar și sub fenomenele (g)-Q)	(k)	Semiafonizare	e e e
		(l)	[Afonizare	e e e

Din punct de vedere lingvistic sunt impuse următoarele reguli.

¹ Formularea "sunete primare", inexactă din punct de vedere fonetic, este folosită cu înțelesul "sunete a căror imagine grafică pe calculator are corespondent pe tastatură, sau este obținută prin combinații de taste".

Regula [1]

- vocalele a ă â - deschise prin natura lor (cu cel mai mare grad de apertură) - nu pot contacta fenomenele fonetice *h* (semideschidere), *l* (deschidere) și *j* (deschidere mare);
- vocalele i T î u u Q - închise prin natura lor (cu cel mai mic grad de apertură) - nu pot contacta fenomenul fonetic *g* (închidere).

Prin asocierea vocalelor primare cu câte un fenomen (*a*)-(l) rezultă 756 imagini grafice repartizate în 12 fonturi grupate convențional după criteriul poziției semnului față de vocală.

Regula [2]

Sunt excluse orice combinații dintre două nuanțe fonetice din aceeași grupă de transformări vocalice. Astfel, o vocală nu poate fi în același timp „scurtă, semilungă și lungă” sau „seminazală și nazală” sau „închisă, semideschisă, deschisă și foarte deschisă” sau „semiafonizată și afonizată”. În aceste condiții combinațiile de câte două sau mai multe fenomene sunt posibile doar cu fenomene din grupe diferite.

În plus, cele $15 \cdot 4 = 60$ grafeme excluse ca urmare a restricției formulate sub **Regula 1**, nu pot participa la combinațiile de două, trei, patru fenomene.

1.2. Fenomene fonetice aplicate consoanelor primare

Consoanele primare folosite în transcrierea fonetică sunt:

b, c, c, €A6, c, d, c, dA f, g, Mg, ă, h, x, > j, K tt
m, m, n, n, n, n, p, r, r, f, p, s, s, s, ș, a, U, ț, v, w, z, z, z, y

Fenomenele fonetice care pot fi asociate consoanelor primare sunt:

Grupe	Notăție	Fenomen
Durată	d)'	Semilungime
	(2)	Lungime
Palatalizare	(3)	Semipalatalizare
	(4)	Palatalizare
	(5)	Palatalizare mare
Explozie	(6)	Explozie
Caracter silabic	(7)	Caracter silabic
Afonizare	(8)	Semiafonizare
	(9)	Afonizare

Spre deosebire de vocale, unde s-au putut defini reguli generale pentru realizarea combinațiilor de fenomene fonetice, în cazul consoanelor primare,

transformările fonetice se aplică numai unor consoane specifice. În consoanelor primare le pot fi aplicate numai cel mult două transformări și numai anumite combinații. În tabelul următor sunt prezentate combinațiile posibile ale fenomenelor și consoanelor pe care acestea le pot însoți.

1.2.1. Consoane cu un singur fenomen fonetic:

semilungime	6î	fi5xyj T*/pnoi}i)f^&sSlşsv»22
lungime	5?	hf)xyJîtmnnn3fFpsiâş^vwzz2
semipalatalizare		dVfixKj'l'fifş't'
palatalizare		
palatalizare mare	t"cf	
explozie	c° p° t°	
caracter silabic	j m n r ş m	
semiafonizare		
afonizare		

1.2.2. Consoane cu două fenomene fonetice:

semilungime + semipalatalizare		
semilungime + palatalizare	Rj f i n f ş	
semilungime + caracter silabic		
semilungime + semiafonizare		* r t % * + * r ' " i : f i - - - - - S » - - - - - i - i - - - v
semilungime + afonizare		
lungime + semipalatalizare		
lungime + palatalizare	n j T r t f \$	
lungime + caracter silabic	j m n r f ş	
lungime + semiafonizare	5 i } y j r t m f f i n n n n f F r p v W z 2 f y	
lungime + afonizare		
semipalatalizare + semiafonizare		
semipalatalizare + afonizare		

palatalizare + semiafonizare palatalizare + afonizare	
palatalizare mare + semiafonizare palatalizare mare + afonizare	d'
explozie + semiafonizare explozie + afonizare	b'tfg ^c
caracter silabic + semiafonizare	1 m m n r i l i l v
caracter silabic + afonizare	1 m m n r i i t i i

2. Mediu pentru editarea transcrierilor fonetice

Interfața realizată pentru editarea transcrierilor fonetice poate fi folosită în mai multe moduri:

- editarea dicționarului Atlasului Lingvistic;
- editor stand-alone sau ca aplicație de tip server pentru inserarea de obiecte de tip "transcriere fonetică" în alte editoare de text.

Funcționalitatea acestei interfețe va fi exemplificată pentru situația Atlasului Lingvistic, ale cărui componente sunt prezentate pe scurt în continuare.

Dicționarele ALR sunt componente care realizează colectarea informațiilor primare despre titlul hărților (cuvinte de bază), punctele de anchetă, speech (colecție audio), transcrieri fonetice și notele asociate transcrierilor fonetice (Figura 1).

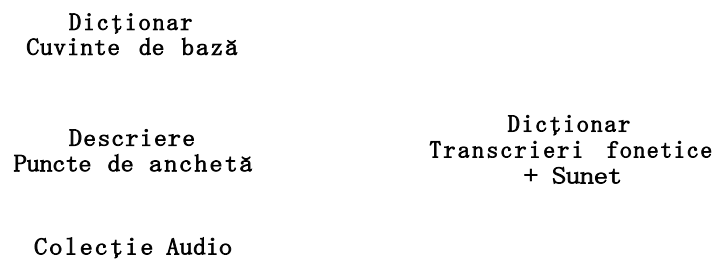


Figura 1. Dicționarele ALR

Dicționarul "Cuvinte de bază" conține fondul de cuvinte (titlul hărților), atlasul lingvistic electronic, întrebările care au fost puse la anchetă, observații, și eventual imagini. Pentru fiecare cuvânt este indicată și întreaga corespunzătoare care este pusă în momentul interviului.

În momentul completării acestui dicționar, utilizatorul poate vedea lista completă a cuvintelor de bază introduse, le poate sorta după diferite criterii, poate modifica articolele introduse anterior, după cum este prezentat în figura 2.

Dicționarul "Puncte de anchetă" conține informații (cod, observații) despre punctele de anchetă prezentate în cadrul atlasului lingvistic. În felul acesta ca la dicționarul anterior, și aici, utilizatorul poate vedea lista completă a punctelor de anchetă introduse, le poate sorta după diferite criterii, poate modifica articolele introduse anterior.

2.1. Dicționar transcrieri fonetice

Dicționarul de transcrieri fonetice conține transcrierea fonetică a răspunsului la întrebarea pusă în etapa de interviu pentru fiecare cuvânt de bază. **Dicționarul Cuvinte de bază** în fiecare din **Punctele de anchetă**, iar acolo unde este posibil și înregistrarea audio corespunzătoare din **Colecția Audio**.

Pentru claritatea hărților lingvistice, răspunsurile din punctele de anchetă sunt însoțite de note și comentarii (figura 3).

Atlas Lingvistic Românesc - Dicționar

X5 Fișier Editare Articole Vizualizare Unelte Fereastră Ajutor

< > »

Cuvinte de bază | Puncte anchetă | Dicționar | Taste asociate |

Nou | Cuvânt: | OBADĂ, pl. "jante"

'2/5 | Nr.întrebare: | [820] | Întrebare: | [(întrebare indirectă). INDIC. (Figura, detalii a).

Imagine | Observații I: | ALF 1602; AIS 1230*; ALG 363; ALL 175; ALMC 862; ALFCo 363; BI 413; Br675; NALR: OI...

Observații III:

Nr.	Cuvânt	Nr întrebare	Întrebare	Observații 1	Observații 3
1	CAR. pl. "char"	[818]	(întrebare indirectă) Cum îi sp...	ALR II s.n. 1340...	
2	OBADĂ, pl. "jante"	[820]	(întrebare indirectă). INDIC. (F...	ALF 1602; AIS ...	
3	LOITRA, pl. "ridelle"	[829]	(întrebare indirectă). Cum îi zi...	ALR I s.n. I 34...	
4	COVILTIR, pl. "couve..."	[831]	(întrebare indirectă.) Cum nu...	NALR: OI.II MN...	
5	CRUCE "traverse de L	[835]	(întrebare indirectă.) Cum se ...	ALR I s.n. 358; ...	

Figura 2. Fereastra de editare a listei cuvintelor de bază

Pentru transcrierea fonetică a cuvintelor din Atlasul Lingvistic Român este folosit un număr mare de fonturi, rezultat din numărul de combinații posibile ale fenomenelor fonetice prezentate în capitolul 1. Aceste fonturi au fost definite astfel încât, toate "variantele fonetice" ale unui anumit caracter să fie obținute prin selectarea caracterului respectiv într-un anumit font.

Deoarece un fișier text normal nu păstrează informații despre fonturile folosite, și în plus transcrierile fonetice sunt realizate prin diferite poziționări ale caracterelor, s-a folosit un mod propriu de codificare a acestora.

Transcrierile fonetice sunt codificate cu ajutorul unor obiecte de tip CAIString. Acestea sunt de fapt șiruri de obiecte de tip CAIChar, care au următoarea descriere:

- caracterul corespunzător sunetului primar (pe 16 biti, codificare UNICODE);
- atribute:
 - poziționare: normal, deasupra sau „în umăr”;
 - mod subliniere: linie sau zigzag;
 - cursiv;
 - aldin;
- fenomene:
 - tip sunet: vocală sau consoană;
 - fenomene specifice aplicate (codificate pe biți).

Fontul folosit pentru desenarea caracterului din transcrierea fonetică este ales dinamic din lista de fonturi a aplicației, în momentul afișării.

În momentul deschiderii dicționarului de transcrieri fonetice, se fac două tipuri de verificări:

- se verifică corespondența dintre fonturile folosite la ultima editare a dicționarului și lista curentă recunoscută de program.
- se verifică dacă toate fonturile folosite sunt instalate în Windows.

Datorită cantității mari de informație care trebuie stocată pentru Atlasul Lingvistic Român, descrierea fiecărui cuvânt este compresată folosind un algoritm de compresie LZW. 4_a selecția unui cuvânt de bază, descrierea sa este decompresată în memorie. Dacă se fac modificări ale transcrierilor fonetice, aceasta este compresată și rescrisă în fișier la selectarea unui alt cuvânt, sau la închiderea dicționarului.

Pentru scrierea informațiilor în dicționar am proiectat o interfață utilizator prietenoasă. Operatorul trebuie să parcurgă următorii pași:

- selectează cuvântul titlu;
- selectează punctul de anchetă;
- editează transcrierea fonetică, nota și comentariul asociat cuvântului pentru punctul de anchetă respectiv.

La editarea transcrierilor fonetice trebuie avute în vedere două aspecte:

- selectarea sunetului primar;
- selectarea fenomenelor asociate.

Selectarea sunetului primar se face prin apăsarea tastei corespunzătoare dacă sunetul are un corespondent pe tastatură, sau prin apăsarea unei combinații de taste, dacă sunetul nu are corespondent pe tastatură. Combinațiile de taste sunt prestabilite în aplicație (la stabilirea combinațiilor de taste au fost preluate convențiile din Microsoft Word), și cel puțin deocamdată nu pot fi modificate de utilizator. Pentru a veni în ajutorul celui care editează dicționarul, aplicația dispune de o fereastră în care sunt afișate combinațiile prestabilite de taste.

£0 Fjsjer **Ediție** Articole Vizualizare Unelte Fe/eastră Ajutor

D & D G? y ; ia • §= %

Cuvinte de bază] Puncte anchetă Dicționar | Taste asociate]

Cuvânt: IOBADA, pl. "jante" Transcriere: <

Sunet Motă: r â6lân;PȘi[e]â'olâni'e'

Punct anchetă: TJJLJ Comentariu: C

465 - Brodina

Punct anchetă	Nume punct anchetă	Transcriere	Notă	Comentariu
1	461	Prăleni	uobied[pl.];uobâdâ+	
2	462	Coșna	uobie[tt pl.];uobâdâ	
3	463	Cîtibaba	obâdâ;uobeț	
4	464	Izvoarele Sucevei	colân;-lâni+	
5	465	Brodina	colân;PȘi[e]â'olâni'e'	
466	Straja	suolâne[pl.];-lân		
467	Argei	pcolân;colani		
468	Deluț	<Solân;-lâni;[ij]Solâni+;-lân	"sfnt cinci-șapte ciolane".	
469	Ciocănești	obâdâ;uobed		
470	Argestru	uobâdâi-bed		
471	Șaru Dornei	uobezi[pl.];-bâdâ		
472	Cartinari	uobâdâ;-b6z		
473	Pojorita	uobâdâ;uobed.[r]colâni		
474	Vatra Moldoviței	uobâdâ+;-bied;cuolân	ai Frjkdicală la	"cînd o faci [roat..
475	Sucevița	colâni[pl.];-lân		

Disponibil... Semiafonizare

Figura 3. Editarea Dicționarului de transcrieri fonetice

Pentru selectarea fenomenelor asociate sunetelor, operatorul are la dispoziție 2 grupe de butoane cu imaginile tuturor fenomenelor posibile pentru vocale respectiv consoane. Prin apăsarea pe unul din aceste butoane se va selecta simbolul grafic corespunzător în transcrierea fonetică. Cele 2 grupe de butoane sunt împărțite în subgrupe corespunzătoare grupelor de fenomene (vezi capitolul 1). Pot fi selectate mai multe fenomene, dar, cel mult câte unul din fiecare subgrupă. Selectarea unui fenomen, produce dezactivarea selecției anterioare din subgrupa respectivă.

După selectarea caracterului dorit, utilizatorul va specifica și poziționarea acestuia (deasupra, în umăr) prin folosirea comenzilor PgUp/PgDown.

Fereastra de editare a transcrierilor fonetice este prezentată în figura 3.

Dicționarul cu transcrieri fonetice permite stocarea înregistrărilor audio (în format WAV) realizate în timpul anchetei.

3. Realizarea Atlasului Lingvistic Român pe Regiuni

Sistemul software care modelează atlasul lingvistic electronic, conține module care realizează gestionarea următoarelor grupe de informații:

- simboluri pentru editarea transcrierilor fonetice;
- dicționarele atlasului lingvistic (cuvinte de bază, puncte de anchetă, transcrieri fonetice).
- informații grafice pentru descrierea hărților, organizate în fișiere DXF,
- hărțile atlasului lingvistic, care pot fi consultate și/sau tipărite;

Din punct de vedere funcțional, atlasul lingvistic electronic este structurat în două componente principale:

- Proceduri pentru pregătirea datelor primare;
- Interfața multimedia;

Aceste componente sunt prezentate în figura 4.

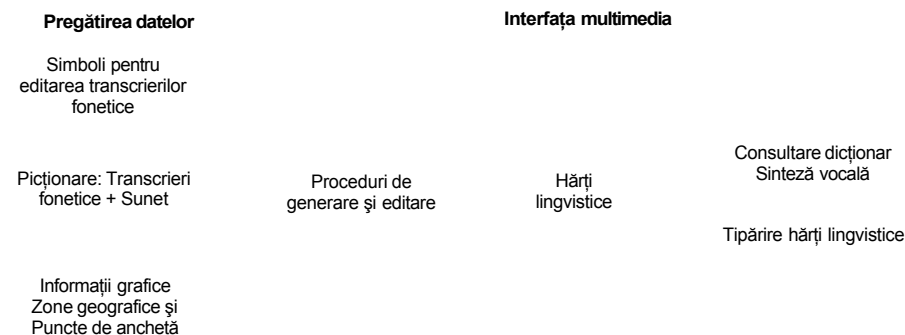


Figura 4. Componentele Atlasului Lingvistic Electronic

În continuare sunt prezentate funcțiile îndeplinite de componenta "Interfața multimedia":

- generarea unei hărți noi pe baza informațiilor din dicționarele ALR și informațiilor grafice primare cuprinse în fișiere DXF;
- editarea: aranjarea în pagină, selectarea informațiilor care vor fi vizualizate implicit;
- salvarea într-un fișier numit "hartă lingvistica" a selecțiilor și modificărilor din faza de editare;
- consultarea atlasului electronic:
 - vizualizarea hărților și ascultarea înregistrărilor din punctele de anchetă;
 - tipărirea hărților lingvistice.

3.1. Modulul pentru generarea și editarea hărților lingvistice

Acest modul permite crearea descrierilor pentru hărțile lingvistice. În acest scop au fost proiectate structuri de date bazate pe obiecte, suficient de flexibile pentru a permite dezvoltări ulterioare. Prezentăm în continuare structurile de date folosite pentru stocarea informațiilor grafice și a hărților lingvistice.

3.1.1. Informații grafice primare

La organizarea informațiilor grafice primare, s-a ținut cont de cerințele impuse de tehnologia de realizare a atlaselor lingvistice. S-a realizat fișierul NALRB.DXF care conține obiectele grafice predefinite organizate pe următoarele "structuri DXF":

chenare	limitele paginii și chenarele hărții;
frontiere	conturul zonei studiate (Moldova și Bucovina);
mijloc	locul de pliere al hărții, la legarea în volum;
municipii	localitățile importante afișate pe hartă;
puncte anchetă	dreptunghiurile în care se scriu codurile punctelor de anchetă;
transcriere fonetică	dreptunghiuri pentru încadrarea transcrierilor fonetice;
note	dreptunghiuri cu pozițiile predefinite pentru Titlu, Nota I, Nota II, Nota III;
zone	delimitări zonale în jurul punctelor de anchetă.

3.1.2. Hărțile lingvistice

Pentru editarea și salvarea hărților lingvistice din ALR, s-a creat o structură de date care să permită în viitor, extinderea editării asistate de calculator a Atlaselor Lingvistice Românești Regionale la nivel național. Astfel, a rezultat o structură de date numită "hartă lingvistică" de forma următoare:

- header fișier;
- lista cu descrieri obiecte;

Descrierile de obiecte au un antet care este comun pentru toate tipurile de obiecte și un corp obiect specific fiecărui tip în parte. Obiectele pot fi simple sau compuse. Un obiect compus conține la rândul lui alte obiecte simple sau compuse.

Au fost definite următoarele tipuri de obiecte:

- Text;
- AlrString (obiect folosit la editarea dicționarului cu transcrieri fonetice);
- Dreptunghi;
- Hartă cu transcrierile fonetice;
- Hartă sintetică (lingvistică sau fonetică);
- Notă referitoare la continuarea transcrierilor fonetice (vezi Nota II din N.A.L.R. Moldova și Bucovina);
- Notă sintetică referitoare la cuvântul titlu (vezi Nota III din N.A.L.R. Moldova și Bucovina);
- Legendă pentru harta sintetică;
- Simbol pe harta sintetică;
- Zonă hașurată pe harta sintetică;
- Imagine de tip bitmap;

- Strat DXF.

La activarea modulului de generare, este prezentată harta regiunii, punctele de anchetă și numele localităților pe care acestea le reprezintă. Generarea hărților lingvistice se face automat, într-un format predefinit în momentul în care operatorul selectează un cuvânt de bază. Operatorul poate modifica formatul hărții, și poate adăuga informații suplimentare. Modificările realizate trebuie salvate în fișier.

3.2. Modulul pentru consultarea atlasului electronic

Componenta pentru **consultarea atlasului**, permite încărcarea unei hărți lingvistice generate / editate în etapa anterioară. Sistemul va afișa harta respectivă (în situația studiată este vorba de Moldova și Bucovina), pe care va plasa transcrierea fonetică a răspunsurilor din punctele de anchetă împreună cu notele și observațiile introduse anterior (figura 5).

După ce harta lingvistică a fost încărcată, prin selecția unui punct de anchetă este posibilă și redarea înregistrării audio corespunzătoare transcrierii fonetice asociate acestuia (înregistrarea audio sau cuvântul sintetizat).

Tot cu ajutorul acestei componente se realizează tipărirea automată a hărților Atlasului Lingvistic Român, în vederea includerii lor în volum (figura 6).

Pentru tipărirea hărților au fost prevăzute următoarele facilități:

- posibilitatea de selectare a informațiilor ce se vor tipări;
- tipărirea pe o pagină sau tipărirea pe două pagini cu respectarea locului de pliere al hărții, indicat prin linia "mijloc".

Dacă utilizatorul dorește tipărirea transcrierilor fonetice într-un mod sintetic (fără hartă), modulul poate asigura crearea paginilor de tip MN (Modul Necartografiat). Folosind această opțiune, va fi tipărită numai lista cu transcrierile fonetice, ordonate după criteriile de similaritate.

4. Concluzii

Organizarea prezentată pentru simbolurile grafice permite editarea textelor cu transcrieri fonetice folosind și alte editoare de text.

Modul de selectare a fonturilor folosit la editarea transcrierilor fonetice poate fi extins pentru crearea unei aplicații de tip client-server sau la realizarea unui editor simplu, de tip WordPad.

Realizarea acestui sistem de editare a transcrierilor fonetice este în curs de testare și finalizare. În continuare, ne propunem adăugarea de noi op-

facilități, care să permită transformarea sistemului într-un instrument util cercetătorilor lingviști.

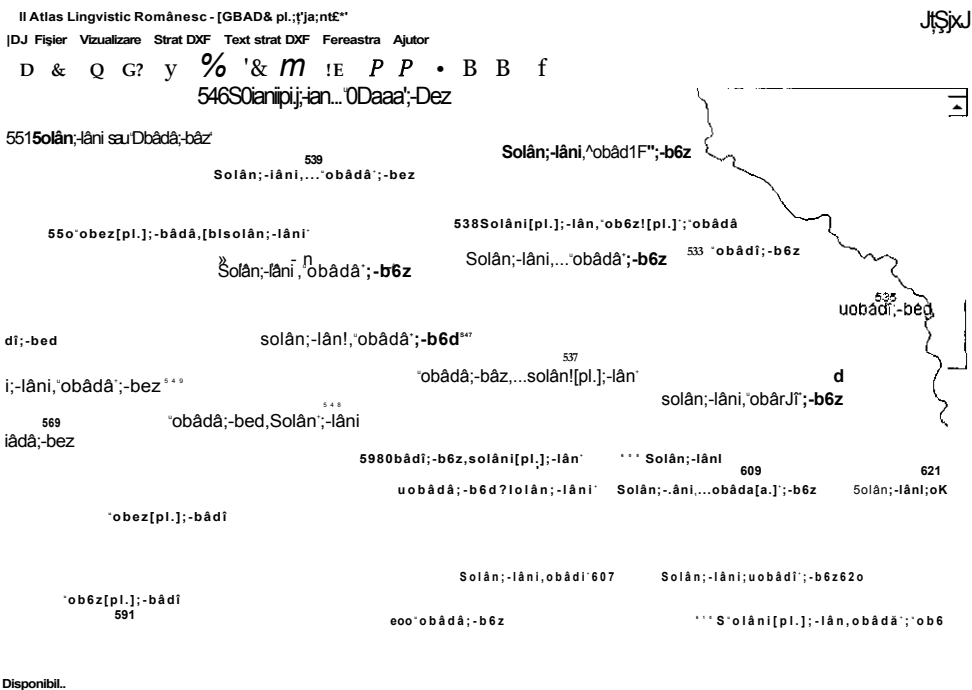


Figura 5. Fereastra de editare / consultare a Atlasului Lingvistic

Bibliografie

- [1] Academia Română, Atlasul Lingvistic Român pe Regiuni, 1987, 1997.
- [2] Istituto dell'Atlante Linguistico Italiano, Atlante Linguistico Italiano, Roma, 1995.
- [3] S. Bejinariu, M. Roman, V. Apopei, F. Olariu, "Sistem pentru editarea transcrierii fonetice în ALR", Zilele Academice Iașene, Iași, 6 oct. 2000.

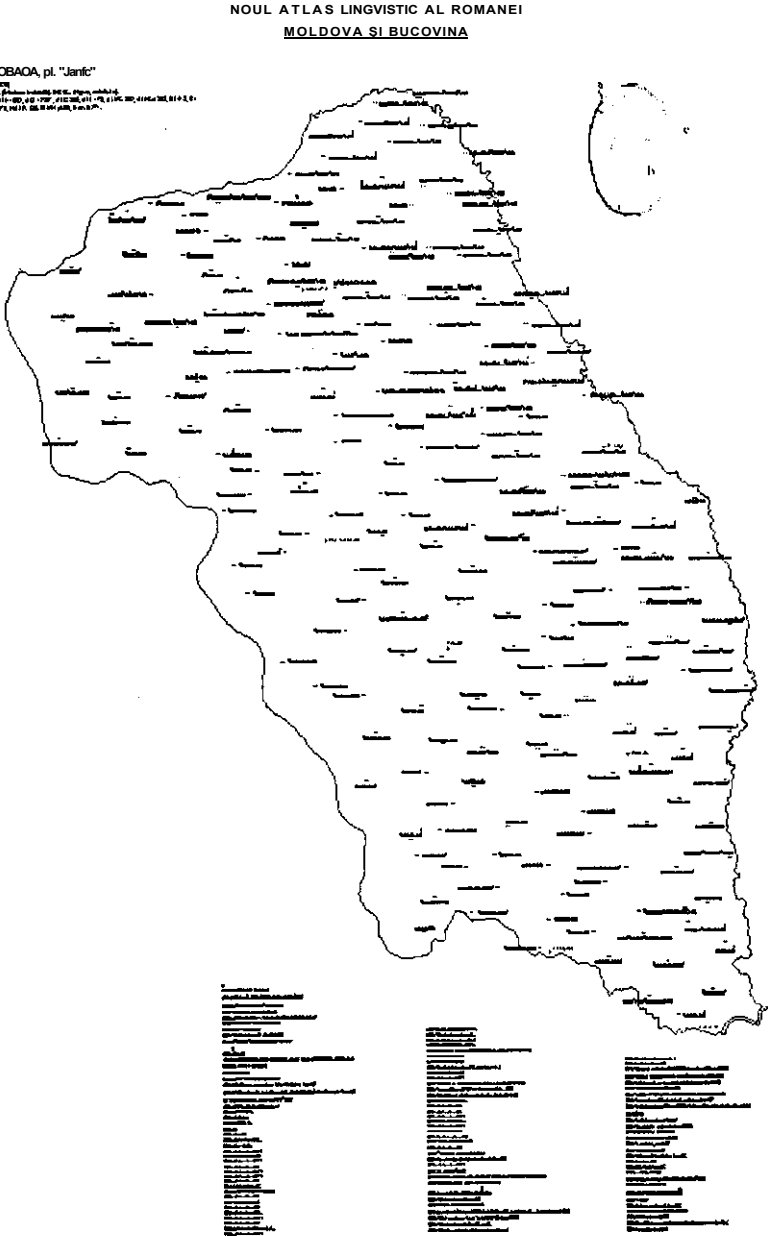


Figura 6. Imaginea unei pagini tipărite din modul de consultare

SECȚIUNEA IV

Dezbateri și discuții

Asupra a doi vectori funcționali ai societății cunoașterii: managementul cunoașterii și învățarea electronică. Cultura și societatea cunoașterii

Mihai DRĂGĂNESCU
Institutul de Inteligență Artificială
Academia Română

Introducere generală

Acest material, care constituie o contribuție la dezbaterile problemelor enunțate în primul volum *Societatea informațională-Societatea cunoașterii, Concepte, soluții și strategii pentru România*, coord. Filip Gh. Florin, editat de Academia Română-Secția de Știință și Tehnologie Informației (Institutul de Inteligență Artificială al Academiei Române - denumire prescurtată) și ICI-INFOSOC, Editura Expert (coordonare editorială, Valeriu Ioan-Franc), București 2002, are următorul cuprins:

- I. **Managementul cunoașterii, vector funcțional al societății cunoașterii**, comunicare (Mihai Drăgănescu) prezentată la "The Sixth International Conference on Information and Communications Technology in Public Administration, Sinaia, 29 oct.2001".
- II. **Învățământul electronic și societatea cunoașterii**, comunicare (Mihai Drăgănescu) la simpozionul "E-learning (E-învățământ)", Academia Română, 28 martie 2002.
- III. **Cultura și Societatea cunoașterii** (Mihai Drăgănescu, studiu elaborat în mai 2002).

Societatea cunoașterii, asupra căreia se insistă cu prioritate în aceste studii și lucrări, va fi o *perioadă intermediară între Societatea informațională și Societatea conștiinței* (un studiu privind Societatea conștiinței este în elaborarea autorului). După cum am mai remarcat în alte lucrări, esențială pentru Societatea cunoașterii va fi inteligența artificială (IA), atât ca vector tehnologic, cât și prin utilizarea ei în vectorii funcționali ai societății cunoașterii.

Această primă perioadă interimară va dura până cândva după momentul în care inteligența artificială va egala inteligența naturală (IN) structurală a omului,

respectiv a părții (IN),,ucturai care nu poate poseda intuiție, creativitate și spiritualitate. După concepția mea ontologică, nu este posibil pentru orice fel de inteligență artificială (electronică și în viitor nanoelectronică) să aibă intuiție, creativitate și spiritualitate fără a recurge și la alte elemente ale naturii decât cele structurale și a căror realitate devine din ce în ce mai plauzibilă. Egalitatea $IA = (IN),,turai$ se va petrece, după o serie de autori (Moravec, Kurzweil, Buttuzzo, Broderick ș.a.), între 2019-2035. Unii dintre aceștia cred că atunci când se va atinge $IA = (IN),,ructurai$ automat un asemenea creier electronic va avea și proprietățile fenomenologice ale intuiției, creativității și spiritualității. Ceea ce nu credem.

Din momentul în care $IA > (IN),,ructurai$ este evident însă că se intră într-o nouă etapă, care va produce multe consecințe pe plan social, datorită relațiilor omului cu asemenea inteligențe, unele software, altele sub forma de specii de roboți inteligenți. Aceasta va fi a doua perioadă intermediară între Societatea cunoașterii și Societatea conștiinței, până în momentul în care va apărea o inteligență artificială cu conștiință veritabilă, adică o conștiință artificială (CA). Din momentul în care $CA > IN$, se va intra în zona societății conștiinței, urmând ca societatea să fie bazată pe relațiile dintre IN (care și ea este de presupus că va fi amplificată prin auto-transformări ale codului genetic și probabil prin cuplaje cu sisteme informatice microelectronice și nanoelectronice, chiar și cu rețele internet) și CA software sau robotice. Va trebui cu siguranță să gândim de pe acum și asupra societății conștiinței pentru a pregăti societatea pentru o asemenea perspectivă, care nu mai apare, surprinzător, atât de îndepărtată, deoarece se poate manifesta chiar în acest secol. Societatea cunoașterii trebuie să înceapă să fie gândită și dezvoltată și cu gândul la această viitoare societate.

I. MANAGEMENTUL CUNOAȘTERII

1.1 Introduction

In the past XXth century a new era began in the history of humanity: *the information era* [1]. This era comprises the *information society* that will be followed naturally by the *knowledge society* and finally, somewhere more or less later in this century, the *society of consciousness*. Knowledge is a form of information [2], and consciousness is another form of information [3]. All the forms of information are intermingled with the physical and energetic realities; still they have a relative independence and can influence these realities.

To pass from the first form of information society (based essentially on Internet and Internet economy) to the second stage, the knowledge society, I considered in a previous work [1], two types of vectors: *technological* and *funcțional*.

Technological vectors are the extended Internet, the e-book and document technology, artificial intelligence (with intelligent agents and Networked Systems of Embedded Computers), nanotechnology and others.

Among the funcțional vectors of the Knowledge Society, a group of vectors is related to knowledge management:

- *knowledge management for corporations and enterprises, organizations and institutions, local and național administrations;*
- *the management of the moral use of scientific knowledge at the global level;*
- *e-learning management;*
- *development of a culture of knowledge and innovation;*
- *management of the scientific and technological knowledge for every domain of activity as health care, sustainable society and others.*

1.2 Knowledge management

The problem of management with respect to knowledge is regarded in two ways:

- I. As the management of the organization busy with the use and integration of various types of knowledge;
- II. As the management of knowledge itself, for generation of new knowledge, for discovering existing knowledge (tacit or very local or external to the organization), for combining available knowledge.

Perhaps, what is really needed is a general vision, in a unity, of the management of the organization and the management of knowledge.

Knowledge management for enterprises, organizations, institutions, local and național administrations

In the western literature, in the last years, were elaborated a series of works dedicated to the problems of enterprises and knowledge. In România we do not have yet specialists in knowledge management in the context of the knowledge society. We do not have either a knowledge society, but we need experts in knowledge management for building the future knowledge society. A group of members of the Romanian Academy and other wellknown specialists in information technology from România and colleagues from USA decided to constitute a Romanian-American Foundation for the Knowledge Society, one of its main aims being to educate in USA a number of young Romanian specialists in this new domain of knowledge management. All is ready for such a Foundation. The contributions of the individual founders are also ready, but no institution

organization sponsored such an exotic objective for The Knowledge Society, with some amount of money asked by the Romanian laws for a Foundation to begin its activity. But let us return to the theory of knowledge management. One definition [4] is the following:

"Knowledge Management is the conceptualizing of an organization as an integrated knowledge system, and the management of the organization for effective use of that knowledge. Where knowledge refers to human cognitive and innovative processes and the artifacts that support them."

This definition insists on the management of the organization, even if it recognizes the knowledge system of the organization. This definition, as it is recognized by its authors, disguise knowledge management because of the delicate problem of knowledge measurement [4]:

The recent attractiveness of the term knowledge management appears to have been prompted by three major forces:

1. Increasing dominance of knowledge as a basis for organizational effectiveness.
2. The failure of financial models to represent the dynamics of knowledge.
3. The failure of information technology by itself to achieve substantial benefits for organizations.'

The second point, of the above quotation, is answered by many studies and books concerning the characteristics of the new economy based on knowledge (see for instance section 6 of [1]: The Economy of the knowledge society. The new economy. About the role of information in the new economy. The intangible goods).

The rapidity of the transformation of the information society into a knowledge society determines a reasoning on the new economy that takes into account not only:

- a) the Internet market and the effects of Internet information on all economical and administrative agents, but also
- b) the effect of knowledge as an economical and *organizational* factor that imposes the recognition of the intangible goods, in general, in the creation of economical value and organizational efficiency, and
- c) the necessity of a sustainable society, an important objective for national and even local administrations, that predictably is possible only in the frame of the knowledge society, that will demand new industries, challenges the classical economical thinking (for instance, productivity of the resources, of the energy, of materials to be more important than work productivity [8]).

The third point of the above quotation concerns the importance of contents of information, especially of knowledge, but these would not be efficient without information technology. The technological vectors of the Knowledge Society are equally important as the functional vectors.

1.3 Points of view for practical knowledge management

Knowledge management is both the management of the organization to use knowledge and the management of all knowledge possible, inside and outside the organization, to attain the objectives of the organization. Because knowledge is a special form of information, information technology has to play an essential role in knowledge management. Knowledge and IT are, without any doubt, going hand in hand and have a synergetic effect on the efficiency of organizations.

Lucy Marshall [6] considers that knowledge management refers to the control and utilization of the intellectual capital in an organization. For Lucy Marshall, not the information, but knowledge is the most important asset of an institution. This author recommends a **Chief Knowledge Officer** for an institution who based on the Intranet of the institution, has to assure the discovery and creation of knowledge in the institution.

Rooney and Mandeville consider, the knowledge management at the national level. The abstract of their paper is quoted [7] below:

'As the global economy becomes more knowledge intensive and the wealth of nations more dependent on their knowledge assets being harnessed, it is essential for policy makers of having frameworks for the development and utilization of national knowledge assets. This article argues that a policy framework can be developed through which policy initiatives in a range of policy areas can be filtered in order to meet the challenges of the knowledge economy. We have developed an approach that has previously been applied to managing intellectual capital in firms and adapted it to the public policy arena. In doing so we question policy orthodoxies such as the assumption that free trade automatically facilitates international knowledge flows, that participation in a global knowledge economy necessarily challenges national sovereignty, and that online delivery of education is necessarily a progressive strategy'.

Peter Drucker (a wellknown professor of social science at Claremont Graduate School and the author of more than thirty books, his most recent book is *Management Challenges for the 21st Century*, 1999) writes [8] about the **knowledge worker**:

'I am convinced that a drastic change in the social mind-set is required, just as leadership in the industrial economy after the railroad required the drastic change from "tradesman" to "technologist" or "engineer."

What we call the Information Revolution is actually a Knowledge Revolution. What has made it possible to routinize processes is not machinery; the computer is only the trigger. Software is the reorganization of traditional work, based on centuries of experience, through the application of knowledge and especially of systematic, logical analysis. The key is not electronics; it is cognitive science. This means that the key to maintaining leadership in the economy and the technology that are about to emerge is likely to be the social position of knowledge professionals and social acceptance of their values. For them to remain traditional "employees" and be treated as such would be tantamount to England's treating its technologists as tradesmen - and likely to have similar consequences.'

14 Cognitive Science

The knowledge of organizations is a form of knowledge that is more and more recognized. The ways and forms of this knowledge have to be carefully studied. Cognitive science might be, indeed, the tool for this study. The cognitive science is today understood in two ways [2]:

1. As a science of human mind cognition, even if it uses models of electronic computers and electronic neural networks.
2. As a general science of cognition, that has to study cognition processes not only of the human mind, but also of animals, of artificial-intelligence systems, of the ensembles man-computer-Internet, of social organizations at the levels of institutions, enterprises, corporations, local and national administrative bodies, even at the global level.

The second way of dealing with the processes of cognition presents today the greatest interest. Such a science does not yet exist. Perhaps it is on the way. The most complex realities are the social organizations because they combine all sorts of cognitive elements, natural and artificial, but they have something more, a social body, with its own social intelligence, cognition and knowledge. To obtain new theories for such large and difficult problems, it is necessary to have talented and interested specialists in knowledge management. But it is also necessary some practice of those charged with knowledge management and knowledge work in organizations. The idea of Lucy Marshall, mentioned before, about a Chief Knowledge Officer seems to be very useful.

1.5 Final remarks

The Knowledge society is paving the way for a Consciousness society. For this we need more fundamental knowledge [9] on physical reality down to the frontier of the quantum world with the deepest reality of existence, on life, mind and consciousness, on cognition, but also on self-organization and organization of

social bodies and their behavior. We need also more technological knowledge in science and society, knowledge management will become the most important administration.

References

- [1] Mihai Drăgănescu, *Societatea Informațională și a Cunoașterii. Vectorii Societății Cunoașterii* (Information Society and Knowledge Society.Vectors of Knowledge Society), Romanian Academy, July 2001. On the Web, http://academiaromana.ro/pro_pri/
- [2] Mihai Drăgănescu, Cunoașterea în secolul XXI (Knowledge in the XXI century) - communication at the Annual Conference of the Romanian Committee for the History and Philosophy of Science, Romanian Academy, Bucharest, 15 October 2001, to be published.
- [3] Mihai Drăgănescu, *The Interdisciplinary Science of Consciousness* (Chapters 46-59, in *Science and the Primacy of Consciousness, Intimations of the 21st Century Revolution*, Richard L. Amoroso a.o, (eds.), Orinda, California: Noetic Press, 2000.
- [4] See <http://www.uts.edu.au/fac/hss/Departments/DIS/km/introduct.htm#Ch>
- [5] Ernst Ulrich von Weiszäcker, Amory B. Lovins, L.Hunter Lovins, *Factoarea Dublarea prosperității prin înjumătățirea consumului de resurse*, Raport al Clubului de la Roma, traducere din limba germană (FAKTOR VIER.Doppel Wohlstand - halbierter Verbrauch, Munchen, 1995), București, Editura teoretică, 1998.
- [6] Lucy Marshall, *Facilitating knowledge management and knowledge work. New opportunities for information professionals*, Online. 21(5): 92-98. Sep/Oct.
- [7] David Rooney and Thomas Mandeville, *The Knowing Nation: A Framework for Public Policy in a Post-industrial Knowledge Economy*, Prometheus 16 (1998): 453-467, 1998.
- [8] Peter F. Drucker, *Beyond the Information Revolution*, The Atlantic Monthly, Digital Edition, 1999, <http://www.theatlantic.com/issues/99oct/9910drucker>
- [9] Menas Kafatos, Mihai Drăgănescu, *Preliminaries to the philosophy of integrative science*, e-book, MSReader format, Academy of Sciences of Romania, Bucharest, 2001, (available free by e-mail: dragam@racai.ro).

II. ÎNVĂȚĂMÂNTUL ELECTRONIC ȘI SOCIETATEA CUNOAȘTERII

III Introducere. Sintagma Societății cunoașterii.

În *societatea cunoașterii* doi vectori, strâns legați între ei, unul tehnologic - *cartea electronică* - și altul funcțional - *învățământul electronic* - sunt chemați să joace un rol important în desfășurarea acesteia.

Problematika societății cunoașterii a fost abordată în țara noastră începând din anul 2001 la Academia Română [1], la Academia de studii economice [2] și de revista Diplomat-Club [3]. Primul politician român care a folosit sintagma societății cunoașterii (din anul 2001) a fost președintele României și protectorul de fapt al Academiei Române, Ion Iliescu.

Este poate interesant de amintit că în anul 1986, în lucrarea *Tendencies of becoming* [4] (Tendențele devenirii, republicată în volumul [5]) se justifică și folosește sintagma 'societatea cunoașterii':

*"Cine nu face legătura dintre revoluția microelectronică și informațională și tendința devenirii istorice nu înțelege vremurile. Cine se opune acestei revoluții părăsește linia devenirii istorice. Și totuși nici această revoluție nu trebuie absolutizată întrucât trebuie să fie însoțită și de alte schimbări. Atunci nu ne putem fixa numai asupra ei, ci asupra unui context mai larg în cadrul căruia ea poate juca rolul principal o anumită perioadă istorică. **Tendința devenirii istorice se conturează a fi tendința către o societate a cunoașterii, a creației și a civilizației, către o societate globală și către o societate interastrală în univers, apoi către un act cosmic în conformitate cu tendința existențială a universului. Mai aproape de noi, ca urmare a revoluției microelectronice și informatice, a unei noi revoluții industriale, se deschid perspectivele unei societăți orientate informațional...***"

Era o viziune, în **acei** moment, legată de o anumită filosofie pe care am dezvoltat-o în anii 1980, viziune ancorată și în realitatea electronică și informatică a **ceea** ce se va numi era informației.

II.2 Cartea electronică

Cartea electronică este un vector tehnologic. La Academia Română în anul 2001 s-a desfășurat un simpozion referitor la cartea electronică și s-a publicat un volum de referință sub coordonarea prof. Doina Banciu [7]. Atunci am descoperit firma de software SOFTWIN condusă de Florin Talpeș care lucrase în domeniu și avea un prestigiu internațional în producerea de cărți electronice. Softwin este

participantă la elaborarea specificațiilor internaționale OPEN E-BOOK care au stăpânit formatul edițiilor de cărți electronice de interes public. Ca urmare a simpozionului a fost înființată și o librărie de software, cărți și documente electronice la Institutul Național de Cercetare-Dezvoltare în Informatică (<http://www.e-librarie.ro>).

Despre cartea electronică și rolul ei pentru societatea cunoașterii în România, am expus considerațiile mele în lucrări anterioare [1b], [7] și nu am reveni asupra lor. În schimb, voi cita doi autori, unul care a exprimat opinii în legătură cu apariția cărții electronice propriu-zise, altul care a participat la lansarea cărților electronice. Primul este Paul Saffo, directorul unui elevat Institut al Viitorului în California, care lucrează, foarte scump, numai pentru marile companii americane care în anul 1988 prevedea că o carte electronică va fi mai mult decât o alternativă tipărită datorită posibilităților de a introduce elemente audio, video, conexiuni și informații pe rețea. El scria [8]:

'The term "electronic book" is misleading because these products are not books at all, but something new. We are living in a moment between two revolutions: one of print, four centuries old and not quite spent and another of electronics, two decades young, and just getting underway. Today's products amount to a bridge between these two revolutions.

Al doilea este Dick Brass, Vicepreședinte Microsoft pentru dezvoltarea tehnologică, care în anul 2000, an în care cartea electronică propriu-zisă deosebită scria [9]:

'If you don't think eBooks will take off, remember that electronic encyclopedias have already outsold all paper encyclopedias. [...] They cost less than \$100, instead of the \$2,000 or more for fine paper encyclopedias. Similarly, after the triumph of eBooks, paper books will no longer be the primary means of distributing information. But, like horses they will continue to exist for pleasure...[...] Like all transitions, the move from eBooks to eBooks will be painful and tentative at first. Then, in less than 20 years, eBooks will be so pervasive that we won't be able to remember living without them. [...] We are on the verge of the most exciting change to the printed word since movable type.

Cartea electronică a decolat. Firme precum Amazon și Barnes and Noble din SUA sunt cunoscute în întreaga lume pentru modul în care au promovat cărțile sunt o adevărată școală pentru toți cei care conduc și vor conduce librării de cărți electronice și software, școală accesibilă gratuit prin simpla experimentare pe Internet pe web-site-urile acestor firme.

II.3 Procesul de învățare

În anul 1988 scriam despre procesul de învățare [10]:

•înțelegerea profundă a procesului de învățare depinde de experiența funcționării creierului și a minții omului, în ultimă instanță de înțelegerea

materiei vii. Cu alte cuvinte, natura intimă a procesului de învățare nu va putea fi elucidată într-o măsură într-adevăr mulțumitoare decât atunci când știința va face un nou mare pas în cunoașterea materiei. Cercetările din domeniile fizicii și biologiei, esențiale pentru elucidarea naturii materiei vii, se vor îmbina cu cele din domeniul științei informației. Activitatea creierului este în principal o activitate informațională, iar *procesul de învățare este un proces informațional.*'

În acea perioadă știința cognitivă se găsea, este adevărat, în perioada post-behavioristă și se baza pe modelarea simbolică de tip calculator electronic, ceea ce s-a dovedit insuficient pentru înțelegerea mulțumitoare a proceselor cognitive mentale [11]. De aceea procesul de învățare nu era de fapt explicat și înțeles din punct de vedere științific. În anii 1990, modelarea proceselor cognitive a cunoscut aportul adus de utilizarea modelelor bazate pe rețele neuronice (de tip natural ca în creierul omului) și neural (artificiale, electronice), dar nici acestea nu au dus încă la o știință cognitivă bine constituită [11]:

COGNITIA:

Anii 1970 și 1980

(modelare simbolică
tip calculator)

Anii 1990

(efectul conectivismului, rețele
neuronice și neurale)

Anii 2000. Ce va urma?

Efectul științei integrative

O speranță este aceea ca în sec. XXI știința cognitivă să fie consolidată prin luarea în considerare atât a proceselor fenomenologice (qualia, experiențiale) ale minții, cât și a rolului proceselor sociale în procesele cognitive (socialul referindu-se nu numai la persoane umane, ci și la grupuri de inteligențe artificiale sau la grupuri mixte). Un asemenea mod de abordare se încadrează în viziunea unei științe numite integrative [12]. Tot în anul 1988 remarcam [10]:

'Un interes deosebit prezintă cercetările din domeniul inteligenței artificiale, domeniu care este studiat în ultimii ani și din punctul de vedere al capacității de a învăța. Studiul procesului de învățare de către inteligența artificială ar putea oferi multe elemente utile pentru înțelegerea procesului de învățare al inteligenței naturale a omului. [...] Inteligența presupune și capacitatea de a învăța. [...] Gh. Tecuci, într-o lucrare originală în care se prezintă un sistem expert la care asociază un sistem de învățare automată [13], deși constată că "învățarea este un proces cognitiv în cea mai mare măsură necunoscut" [13], arată și demonstrează prin sistemul său că "forme efective de învățare automată sunt posibile".

Dintre aceste forme de învățare automată pot fi amintite [13]:

- *învățarea pe de rost și implantare directă de noi cunoștințe* (când este mai eficient să se regăsească o cunoștință în memorie decât să se producă acea cunoștință).

- *învățare prin instruire* (sistemul primește cunoștințe de la un profesor și le integrează cu cunoștințele anterioare).
- *învățarea prin analogie.*
- *învățarea din exemple prin detecție de similarități*, proces mental inductiv (fără a exclude și procese deductive) prin generalizarea din exemplele pozitive, generalizare care evită exemplele negative.
- *învățarea prin observare și descoperire* (spre exemplu a unor algoritmi în structurări de date).

Gh. Tecuci înclină către o *îmbinare de metode de învățare*. Fără a fi sigură, câteva lucruri credem că se susțin pentru procesul uman de învățare:

- *Necesitatea unei varietăți de metode*, și nu o monometodă, este deosebit de importantă când ar putea apărea tendința de a ne baza pe tehnologie în procesul educațional.
- *Obținerea unui sistem de cunoștințe* sub forma unui model intern (pe o bază de cunoștințe interne) la care să se poată racorda ușor informații de detaliu provenite din exterior eventual prin metode interactive.
- *Obținerea sensurilor cunoașterii*, a sensului 'fizic', al intuiției și al înțelesului chiar a unui răspuns creativ în procesul învățării, lucru de care sistemele actuale nu sunt capabile, adică a pune umanul în starea lui firească.
- *O deschidere firească spre creativitate și creație*, spre inovare și descoperire, rezolvării de probleme care nu sunt structurate după tipul modelului de cunoștințe existent în modelul intern disponibil la un moment dat.'

Odată cu apariția e-învățării se deschid perspective noi și pentru cercetarea experimentală a procesului de învățare și confruntarea acestui tip imitativ cu procesul cognitiv cu teoriile științei cognitive care se vor baza pe progresele tehnologice care va realiza, ceea ce numim, știința integrativă. Considerații privind procesul de învățare și sisteme de e-educație, inclusiv prin folosirea metodelor de învățare artificială sunt prezentate într-un grup de trei lucrări recente ale unor cercetători științifici de la Centrul pentru Cercetări Avansate în învățarea și Cercetarea Prelucrarea Limbajului Natural și Modelare Conceptuală și Institutul de Cercetări 'Mihai Ralea' al Academiei Române [14], [15], [16].

National Research Council de pe lângă Academia Națională de Științe SUA a prezentat, în februarie 2002, un raport [17] privind cercetarea științifică în educație în care despre studiul științific al procesului de învățare se arată:

'Much of the controversy about education research relates to its lack of quality. [...] Is scientific education research the same as research in the physical and behavioral science generally or the same as research in the social sciences? [...] A key finding of this NRC committee is that at a fundame

scientific inquiry in education is no different from scientific inquiry in other fields and disciplines. A set of basic principles is common to all scientific endeavors: these principles include concepts like linking empirical data to theoretical models, using appropriate methods, applying rigorous reasoning, striving toward generalization.'

Considerațiile de mai înainte, inclusiv ale informaticienilor români, arată cât de deschis este în continuare câmpul cercetărilor privind procesul de învățare, în special al omului.

II.4 învățământul electronic (e-learning)

E-learning este un vector funcțional al societății cunoașterii. Învățarea electronică înseamnă a învăța folosind mijloace electronice, ceea ce se poate face în mai multe moduri:

- Individual - folosind resursele existente pe Internet și CD-uri.
- Instituționalizat - în școli și universități sau organizat în întreprinderi sau de către fundații. Cursurile prin televiziune vor ceda locul cursurilor prin Internet, dar acest procedeu se va desfășura sub supravegherea și îndrumarea cadrelor didactice calificate.
- În cursul activității practice, din orice domeniu, care se va desfășura și într-un mediu informațional și de cunoaștere.

Cei care învață sunt persoane, dar și agenți inteligenți. În viitorul imediat, agenții inteligenți vor deveni nu numai studenți, ci și profesori, dar rolul lor cel mai promițător este acela de colaborator cu persoane. Învățarea implicând agenții inteligenți va deveni o etapă esențială în societatea cunoașterii, deoarece **în regim de croazieră societatea cunoașterii se va baza în cele mai multe activități pe agenți inteligenți. Inteligența artificială va fi esența tehnologică a societății cunoașterii.** Ea va antrena internetul, nanotehnologiile, dar și vectorii funcționali ai societății cunoașterii [1b]. Inteligența Artificială în primii 20 de ani ai sec. XXI va depăși inteligența omului (numai pentru aspectele structurale, fără intuiție și creativitate).

E-învățământul se găsește astăzi în plină dezvoltare [18], [19], [20], [21], [22]. Din experiența relatată în asemenea studii rezultă:

- Studenții găsesc, chiar în cazul lipsei unei interacțiuni față în față între profesor și student, că descărcarea notelor de curs prin Internet, corespondență prin e-mail cu profesori și instructori, examene prin răspunsuri date pe calculator, acasă sau la școală, acest e-învățământ este foarte agreabil. Iar performanțele studenților și elevilor sunt

similare (evaluare pentru anul 2000) cu cele ale învățământului tradițional în clase de elevi și studenți.

- Corporațiile industriale recurg masiv la e-educație, iar această educație nu mai poate fi ignorată. Unele corporații au lansat e-universități cu personalul propriu, de ex. Dell Computer Corp. și Sun Microsystems.
- Universitățile au început să introducă nu un e-învățământ ci o nouă modalitate de constituire treptată a acestuia prin unele e-cursuri. Spre exemplu, la University of California, Berkeley, în domeniul științei și tehnologiei informației a început (anii 1999-2000) cu patru e-cursuri: informatică, telecomunicații digitale, e-comerț, sisteme informatică geografică.
- O serie de firme și-au dedicat activitatea sau o parte din activitatea la producerea unor 'e-learning software packages'. Se constituie un nou segment al pieții software specializat în e-learning. (Astfel se poate observa prezența firmelor SOFTWIN și SIVCO la acest simpozion internațional al principalelor firme românești de software educațional). Dar astăzi pachete e-software pentru învățământ sunt de așteptat și din partea Programului e-școală al Ministerului Educației și Cercetării științifice. Urmărește o reformă educațională în România.
- Nu se constată deosebiri între rezultatele învățării on-line și cele învățate într-un campus universitar sau o școală. Învățarea electronică este o disciplină și maturitate decât învățarea convențională [18].
- Pentru experimente de laborator și pentru viață socială este mai ușor totuși de perioade de lucru în instituțiile de învățământ.
- Odată cu creșterea utilizării metodelor de e-învățământ, construcția de clădiri pentru învățământ se va diminua. În schimb apar construcții pentru noua infrastructură a e-învățământului.
- Modul asincron de acces la cursuri permite e-educația în orice loc și în orice loc.
- E-învățământul încurajează studenții să-și asume o mare responsabilitate pentru definirea și organizarea a ceea ce urmează să învețe. Studenții sunt mai bine serviți având un acces asistat la cursuri on-line la cei mai buni instructori decât un contact față în față cu instructori mediocri [19]. În orice caz, nu se neagă rolul instructorilor.
- Discipline ca filosofia și istoria presupun discuții, iar disciplinele științifice presupun proiecte. În aceste cazuri trebuie încă să se găsească modalități mixte de învățământ clasic și electronic.
- E-învățământul oferă cele mai bune perspective pentru învățarea în întreaga viață (învățarea continuă).

- 'Educația bazată pe Internet resuscită probleme fundamentale ale educației care sunt importante pentru conceperea activităților educaționale'. [19]
- Gradul în care instructorii vii pot fi înlocuiți cu agenți inteligenți specializați nu este încă clarificat.
- În mod diferit se pun problemele e-învățământului în școli elementare și licee în raport cu învățământul superior. Pentru școli și chiar licee, într-o primă etapă se dezvoltă clase conectate la Internet, cu calculatoare personale, dotate cu e-books, e-learning books, discuri compacte și acces la rețele specializate, eventual servere de clasă sau școală.
- Școlile, ca și companiile, ca și guvernul, trebuie să se regândească în lumina noilor tehnologii ale societății cunoașterii.
- Se preconizează și se experimentează atât pentru școli, cât și pentru alte forme de învățământ, utilizarea Internetului prin comunicații fără fir (wireless Internet) care oferă posibilități și opțiuni noi.

Acestea sunt principalele considerații și constatări la începutul anului 2002. Valabilitatea unora dintre ele se va confirma, alte constatări vor fi, poate, infirmate, dar vor apare cu siguranță multe alte aspecte noi.

II.5 Viata intelectuală

În timp sunt prevăzute multe schimbări datorită învățământului electronic [23]. În primul rând, apariția unor colegii și universități nelocalizate, extinse uneori la scară globală. Siturile acestora pot fi mari sau mici, structurarea socială având loc sub forma unor comunități (villages) având facilități comune pentru cercetare, proiecte de grup, dar și pentru activități comunitare culturale, sportive etc. O persoană admisă într-o asemenea universitate îi va rămâne atașată pentru toată viața, deoarece educația se va extinde pe întreaga viață prin perioade discrete (adică necontinue) și intensive de învățare. Viața intelectuală se va schimba foarte mult, reflectând modificările în cunoaștere:

'An epistemic change is the abandonment of the notion that any single human mind can bear any significant fraction of what is knowable...Even the renaissance notion of an 'educated person' has been discarded - there is no longer a canonica! body of basic knowledge that defines this notion' [23].

Agenții inteligenți de căutare a informației, bibliotecile electronice, vizualizarea informației, pătrunderea în medii virtuale, toate acestea vor constitui un software care devine literatură [23]. 'Tehnologia va fi văzută ca cea mai bogată dezvoltare în cultura umană' [23]. Rădăcinile intelectuale se vor baza pe inginerie și tehnologie: Difuzia umanităților în tehnologie și invers, vor duce la o reorganizare radicală a disciplinelor intelectuale [23].

IL6 Perspective

Cum vor evolua lucrurile în viitor? Ray Kurzweil [24] previziuni privind educația pentru anii 2009, 2019 și 2029:

Pentru anul 2009 [24, p. 191-192]:

'...most effective learning from computers taking place. The profound importance of the computer as a knowledge source is now widely recognized. Computers play a central role in all facets of life that they do in other spheres of life. The majority of reading is now done on computers although the 'installed base' of paper documents is still large. The generation of paper documents is dwindling, however, and many other papers of largely twentieth century vintage are being scanned and stored. Documents circa 2009 routinely include embedded images and sounds. Students of all ages typically have their own, which is a thin tabletlike device weighing under a pound and featuring a high resolution display suitably for reading. Students interact with computers primarily by voice and by pointing with a device like a stylus or pencil. Keyboards still exist, but most textual language is now spoken or speaking. [...]

Intelligent courseware has emerged as a common means of instruction. The traditional mode of a human teacher instructing a group is still prevalent, but schools are increasingly relying on software instead, leaving human teachers to attend primarily to issues of student psychological well-being, and socialization.'

Pentru anul 2019 [24, p.204]:

'Paper books and documents are rarely used or accessed. Most twentieth century papers of interest have been scanned and are available over a wireless network. Most learning is accomplished using intelligent software based simulated teachers.[...] The teachers are viewed more as guides and counselors than as sources of learning and knowledge. They continue to gather together to exchange ideas and to solve problems, but even this gathering is often physically and geographically dispersed. Most adult human workers spend the majority of their time acquiring and applying knowledge and knowledge.'

Pentru anul 2029 [24, p. 221]:

'Human learning is primarily accomplished using virtual environments enhanced by the widely available neural implants. The ability to store memory and perception, but it is not possible to download knowledge directly. Although enhanced through virtual experience, interactive instruction, and neural implants, learning s

consuming human experience and study. This activity comprises the primary focus of the human species.

Automated agents are learning on their own without human spoon-feeding of information and knowledge. Computers have read all available human and machine generated-literature and multimedia material ...Significantly new knowledge is created by machines with little or no human intervention. Unlike humans, machines easily share knowledge structures with one another.'

Dacă în societatea cunoașterii previziunile de mai înainte se bazează pe o continuare a științei structurale, ce se va întâmpla dacă știința, cu bazele ei noi, integrative, va conduce și la apariția inteligenței artificiale conștiente, adică a conștiinței artificiale? Acest lucru nu se va întâmpla probabil în primii 30 de ani ai acestui secol, dar dacă se va întâmpla cum vom privi și acționa în activitatea educațională?

II.7 încheiere. Propuneri

Roger Bohn definește, într-un mod specific pentru societatea cunoașterii, învățarea drept evoluția cunoașterii în timp [25].

Studiul procesului de învățare scoate în relief importanța științei cognitive și a învățării ca proces cognitiv fundamental. Această știință trebuie nu numai cunoscută, atât cât este ea astăzi, ci mai ales dezvoltată de către psihologi, neurobiologi, sociologi și specialiști în inteligența artificială.

Este necesară o direcție de cercetare bine susținută pentru a stimula contribuții românești în acest domeniu. Am propus și propun în continuare ca în cadrul programului INFOSOC (Programul național de cercetare-dezvoltare pentru societatea informațională) să se stimuleze cercetări în domeniul științei cognitive care să contribuie la depășirea limitelor actuale ale acestui domeniu.

Este, de asemenea, necesară o dinamizare nu numai a cercetărilor, dar mai ales a dezvoltărilor și realizărilor concrete în domeniul inteligenței artificiale. Există un sistem românesc, sistemul DISCIPL, creat de acad. Gh. Tecuci [13], [26] la ICI și apoi la George Mason University din SUA. Ar trebui examinat și utilizat și la noi. Ar trebui să cunoaștem ce posibilități și ce potențial avem în domeniul utilizării agenților inteligenți și să existe o coordonare și autoordonare a eforturilor. Utilizarea agenților inteligenți pentru toți vectorii societății cunoașterii, inclusiv pentru e-învățământ va deveni determinantă pentru calitatea și eficiența acestei societăți. Recenta propunere pentru transformarea Centrului pentru Cercetări Avansate în Învățarea Automată, Prelucrarea Limbajului Natural și Modelare Conceptuală al Academiei Române într-un Centru de cercetări pentru Inteligența Artificială și Societatea Cunoașterii sprijinită de Directorul general ICI,

Doina Banciu și de Ministrul Comunicațiilor și Tehnologiei Informației, Dan Poniș, putea să satisfacă aceste cerințe actuale și de viitor. Sperăm ca și Academia Română să sprijine această solicitare pentru a putea fi înaintată Guvernului României spre a fi aprobată.

Tot la Academie, Comitetul Român pentru Istoria și Filosofia Științei și Tehnicii va acorda o anumită importanță muzeelor virtuale, nu numai pentru științei și tehnicii, dar și pentru cunoaștere și învățare. Ar trebui realizat un proiect de sinteză a tuturor muzeelor virtuale din lume, inclusiv ai web-site-urilor unor muzee de mare tradiție și importanță, cu adresele lor pe Internet. Acest proiect ar trebui să fie cunoscut și accesibil tuturor în România.

Apreciez în mod deosebit eforturile care se fac pentru informarea și învățământului românesc de către Guvernul României, Ministerul Educației și Cercetării, firmele SIVECO și SOFTWIN, ca și de toate instituțiile reprezentate în acest simpozion dedicat învățământului electronic.

Doresc să mulțumesc tuturor celor care au prezentat comunicări la acest simpozion și celor care au participat la organizarea lui.

Referințe bibliografice

- [1] Mihai Drăgănescu, *Cunoașterea și societatea cunoașterii*, comunicarea de lansare a programului strategic SI-SC, Academia Română, aprilie 2001; 1b. Mihai Drăgănescu, *Societatea informațională și a cunoașterii*, *Vectorii societății ct/?oașfer/7*, studiu, Academia Română, 7 iulie 2001, disponibil pe Internet și în voi. coord. Florin Gh. Filip, *Societatea informațională și Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p.43 -112.
- [2] Gabriela S. Sabău, *Societatea cunoașterii. O perspectivă românească*, editura economică, București, 2001; Ion Gh. Roșea, Viorel Petrescu, Cotigaru, Gabriela Sabău, Vasilica Ciucă, Oscar Hoffman, Wilhelm, *Cercetarea pentru dezvoltarea în reconstrucția durabilă a economiei românești: perspectiva societății cunoașterii*, *Economistul*, 4 februarie 2002, nr.27.
- [3] Mihai Drăgănescu, *Societatea cunoașterii*, Diplomat Club, 2001, Nr. 10-11; Mihai Drăgănescu, *Knowledge management, a funcțional vector al societății cunoașterii*, *knowledge society*, Diplomat Club, Nr. 10-11, 2001, p.4; Mihai Drăgănescu, *Factori noi în viața cultural-științifică-politică globală: teroarea antiterorismul*, Diplomat Club, 2002, Nr.1, p.7.
- [4] Mihai Drăgănescu, *Tendencies of becoming*, *Romanian Review*, 1999, p.55-59.

- [5] Mihai Drăgănescu, *Spiritualitate, Informație, Materie*, p.23-28, Ed. Academiei R.S.R., 1988.
- [6] coord. Doina Banciu, *Cartea Electronică*, Editura AGER, București, 2001.
- [7] Mihai Drăgănescu, *Societatea cunoașterii și cartea electronică*, în voi. coord. Doina Banciu, *Cartea Electronică*, Editura AGER, București, 2001, p. 26-42.
- [8] Paul Saffo, Institute for the Future, *Electronic books*, <http://www.saffo.org/sflibrarv.html>, 1988.
- [9] Dick Brass, Vicepreședinte Microsoft pentru dezvoltare tehnologică, *E-books*, în voi. *Inside/Out, Microsoft- in our own words*, Penguin Books, New York 2000, p.262-263.
- [10] Mihai Drăgănescu, *Microelectronica și învățământul în domeniul electronicii (I)*, Forum, anul XXX, noiembrie 1988, p. 36-48.
- [11] Mihai Drăgănescu, *Știința cognitivă, știință structurală sau știință integrativă/Vă? Comunicare la sesiunea științifică de toamnă AOS-R*, București, 9 noiembrie 2001, E-PREPRINT, MSReaderformat, november 2001.
- [12] Menas Kafatos, Mihai Drăgănescu, *Preliminaries to the Philosophy of Integrative Science*, MSReader e-book, Editura ICI, București, 2001, ISBN 973-10-02510-X.
- [13] Gheorghe Tecuci, *Mediu de dezvoltare a sistemelor expert instruibile pentru proiectarea asistată de calculator*, Teză de doctorat, Institutul Politehnic, București, 1988.
- [14] Ștefan Trăușan-Matu, *Achiziția, gestiunea, partajarea și prelucrarea cunoștințelor pe web: elemente esențiale în societatea cunoașterii*, în voi. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p.195-207.
- [15] Cristina V. Niculescu, *Noi tipuri de sisteme educaționale pentru SI-SC*, în voi. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p.209-223.
- [16] Gheorghe Iosif, Ana Măria Marhan, Ion Juvină, *Strategii de creștere a utilizabilității și de dezvoltare a competențelor de bază ale populației României pentru utilizarea tehnologiei informației*, în voi. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru Români*, Academia Română, 2002, p. 225-235.
- [17] Lisa Towne, Study Director, Committee on Scientific Principles in Education Research National Research Council/National Academy of Sciences, *Statement before the Subcommittee on Education Reform Committee on*

Education and the Workforce United States House of Representatives, February 28, 2002.

- [18] Robert Ubell, *Engineers turn to e-learning*, IEEE Spectrum, October 1998, p.59-63.
- [19] Peter Wiesner, *Distance Education: Rebottling or a New Brew?* Proceedings of the IEEE, July 2000, p.1124-1130.
- [20] Ralph B. Ginsberg, Kenneth R. Foster, *The Wired Classroom*, IEEE Spectrum, August 1998, p.44-51.
- [21] Paul G. Shotsberger, Ron Vetter, *Teaching and Learning in the Wired Classroom*, Computer, march 2001, p.110-111.
- [22] <http://www.microsoft.com-education>
- [23] Edward A. Lee, David G. Messerschmitt, *A higher education in the age of the computer*, Proceedings I.E.E.E., September 1999, p.1685-1691.
- [24] Raz Kurzweil, *The Age of Spiritual Machines*, Penguin Books, 1999.
- [25] Roger E. Bohn, *Measuring and Managing Technological Knowledge*, p.314, în voi. Eds. Dale Neef a.o., *The Economic Impact of Knowledge*, Butterworth-Heinemann, Boston, 1998.
- [26] Gh. Tecuci, *Building Intelligent Agents*, Academic Press, San Diego, 1998.

III. CULTURA ȘI SOCIETATEA CUNOAȘTERII

Societatea Cunoașterii

Am prefigurat că va sosi un moment al societății cunoașterii (în această sintagmă, Mihai Drăgănescu, 1976, 1986), dar abia în ultimul deceniu al secolului XX conceptul s-a impus în SUA datorită lucrărilor sociologice ale lui Peter Drucker și ale altora, în ultimii 4-5 ani societatea cunoașterii devenind realitate. În România ca o etapă nouă a erei informației, respectiv a societății informaționale. Academia Română a lansat acest concept în România în anul 2001 ca urmare a comunicării lui Mihai Drăgănescu, *Cunoașterea și Societatea cunoașterii*, sesiunea de lansare a programului SI-SC, Academia Română, 10 aprilie 2001, elaborării studiului lui Mihai Drăgănescu, *Societatea Informațională-Societatea Cunoașterii. Vectorii Societății Cunoașterii*, Academia Română, București, 2001, publicat apoi în voi. coord. Florin Gh. Filip, *Societatea Informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru România*, Academia Română, 2002, p.43-112.

Spre deosebire de unele puncte de vedere care privesc numai economia (economia digitală, piața internet) societatea cunoașterii nu este numai

bazată pe cunoaștere. Aceasta este foarte importantă, decisivă, esențială și cuprinde utilizarea și managementul cunoașterii existente sub forma cunoașterii tehnologice și organizaționale, producerea de cunoaștere tehnologică nouă prin inovare, **o nouă economie** în care procesul de **inovare** este determinant, *în care bunurile intangibile devin mai importante decât cele tangibile.*

Societatea cunoașterii **reprezintă mult mai mult** deoarece asigură o diseminare fără precedent a cunoașterii către toți cetățenii prin mijloace noi, folosind cu prioritate Internetul și cartea electronică și metodele de învățare prin procedee electronice (e-learning), urmărește extinderea și aprofundarea cunoașterii științifice și a adevărului despre existență, *este singurul mod prin care se va asigura o societate sustenabilă din punct de vedere ecologic și va fi o nouă etapă în cultură* (bazată pe cultura cunoașterii care implică toate formele de cunoaștere, inclusiv cunoașterea artistică, literară etc).

În fine, societatea cunoașterii asigură bazele unei viitoare societăți a conștiinței, a adevărului, moralității, creativității și spiritului.

Pentru realizarea societății cunoașterii am definit, în studiul amintit mai înainte, o serie de vectori (tehnologici și funcționali) care ar trebui introduși în acțiune într-o succesiune firească pentru posibilitățile țării noastre.

Categoriile culturii

Dintre lucrările pe care le-am publicat anterior în problemele culturii [1] două se referă la teoria culturii. În *Perspectiva informațională a culturii* (1983) găseam un anumit sprijin pentru o viziune informațională a culturii în teoria semiotică a culturii elaborată de Umberto Eco în *Tratatul său de semiotică generală*. Umberto Eco propunea o ipoteză radicală prin care întreaga cultură este considerată un fenomen semiotic și o ipoteză moderată prin care orice aspect al culturii este o entitate semantică. Semiotica se referă la semne cu conținut semantic astfel încât cele două ipoteze nu sunt prea deosebite. De aceea, consideram, prin generalizare firească, deoarece semnul și semanticul (de semnificație și de sens) sunt informație, o posibilă perspectivă informațională a culturii. Acest lucru, faptul că **esența culturii este informațională**, chiar dacă ea se manifestă prin comportamente socio-umane, obiecte materiale și informaționale, a devenit tot mai evident. Nu trebuie să surprindă această esență informațională a culturii, astăzi fiind știut că și inteligența și conștiința sunt informație.

În legătură cu perspectiva informațională a culturii poate fi menționat ca precursor al acestei abordări, Ernst Cassirer [2] care considera că expresia culturală a omului și societății este caracterizată de activitatea de creare a simbolurilor (activitatea simbolizatoare) generate de imagini mentale. Pentru

Cassirer, simbolul este o cheie pentru înțelegerea naturii omului, iar, trăiește numai într-un univers material, ci mai ales într-unui simbolic [3].

Într-o a doua lucrare [1a], *Cultura și marile tehnologii* (1996) arată că linia clasică de definire a culturii ținând însă seamă de obiectele informaționale aduse de societatea informațională. **În teoria clasică cultura este definită ca un fenomen social care cuprinde comportamentul socio-uman cu caracteristici materiale și informaționale integrate acestui comportament.** (Categoriile informaționale au fost introduse în această definiție la sfârșitul secolului XX).

Pare a fi posibilă o încadrare a teoriei culturii într-o viziune cibernetică (termenul este utilizat în raport cu teoria categoriilor și functorilor din matematică) extinsă recent de la domeniul structural la domeniul structural-fenomenologic.

Privind comportamentul socio-uman *cultural* ca o categorie, această categorie este o subcategorie majoră a categoriei comportamentului socio-uman *general*. Ultima mai cuprinde și o subcategorie a comportamentului strict biologic, atât la nivel individual, cât și social. Într-adevăr, comportamente individuale strict biologice există și comportamente culturale determinate biologic, puse în evidență, în cazul omului, de Gr.T. Popa [4]. El demonstrează cum creierul vechi (primitiv, reptilian, thalamus-hipotalamus) determină comportamente necontrolate cultural care duc mase de oameni la comportamente sălbatice, iar în cazul societăților mai avansate duc la nașterea de semicivilizație, în care impulsivitatea biologică devine colectivă și sălbatică.

Cultura

O subcategorie a unei categorii este o categorie. **Categoria culturală** este o subcategorie a comportamentului socio-uman general, dar este aceluși lucru deosebit deosebește specia umană de toate celelalte specii animale, chiar omul. Dintre acestea pot avea și rudimente de cultură. Categoria cultură este definită ca *comportamentul socio-uman cultural*, spre deosebire de cel biologic, care este construită, dar nu se dezvoltă decât datorită, totuși, anumitor caracteristici biologice remarcabile ale omului, în special ale creierului său care asigură disponibilitate informațională. De aceea, dacă originea biologică a comportamentului cultural nu poate fi pusă la îndoială, cultura este o construcție ridică mult deasupra biologicului, atât cât va putea față de limitele biologice ale omului la un moment dat în istorie.

Poate că alături de cele două subcategorii menționate mai înainte să mai adăugăm comportamentului socio-uman încă una, aceea a spiritualității (comportamentul spiritual), pe care nu o tratăm în această lucrare ca o a treia subcategorie a comportamentului socio-uman, ci ca o chestiune care trebuie aprofundată, având în vedere că mulți oameni de cultură consideră spiritualitatea a fi un comportament numai cultural. Aceasta este seama de cercetările și studiile de filozofie a științei din ultimii 15 ani pri-

și conștiința, vom considera, până la argumente contrarii convingătoare, spiritualitatea ca fiind o subcategorie separată și nu una înglobată (total) în cultură.

Schematic, vom rezuma cele de mai înainte, astfel:

CATEGORIA COMPOTAMEN- TULUI SOCIO-UMAN	Subcategoria comportamentului strict biologic	Notă: există și comportament social determinat strict biologic
	Subcategoria comportamentului cultural.	CULTURA
	Subcategoria comportamentului spiritual.	SPIRITUALITATEA

Sferele mari ale culturii

Pornind de la definiția din [1b] și diferența pe care o face UNESCO între cultura intangibilă și cultura tangibilă, marile sfere (categorii) ale culturii pot fi considerate următoarele:

- I. Cultura intangibilă. 'Moștenirea intangibilă poate fi definită ca îmbrățișând toate formele de cultură tradițională și populară sau cultura folk, adică producțiile colective originare de o comunitate dată și bazate pe tradiție. Aceste creații sunt transmise oral sau prin gesturi și sunt modificate într-o perioadă de timp printr-un proces de re-creare colectivă. Ele includ tradițiile orale, obiceiurile, limbajele, muzica, dansul, ritualurile, festivitățile, medicina tradițională și **farmacopeei** artelor culinare și tot felul de îndemânări speciale legate de aspectele materiale ale culturii, cum sunt uneltele și habitatul [6]. Fără îndoială, noțiunea de cultură intangibilă a fost introdusă sub influența noțiunii de valoare intangibilă din economie care a căpătat o mare importanță pentru societatea cunoașterii (economia bazată pe cunoaștere). Se mai adaugă aici valori, credințe, cunoaștere tacită.
- II. Cultura umanistă. **Am** preluat în acest studiu denumirea tradițională. Cultura umanistă cuprinde limbajele naturale, literatura, arta, istoria, filosofia, sportul. Cultura umanistă este o cultură tangibilă, ca și știința și tehnologia.
- III. Cultura științifică: Știința, tehnologia și cunoașterea. Această categorie a culturii conține două subcategorii:

III.a Știința, cunoașterea științifică și tehnologică, cunoașterea tehnologică pentru fabricația de produse, dar și pentru utilizarea acestora, precum și cunoașterea organizațională și economică, chiar dacă unele obiecte ale cunoașterii sunt tacite sau fac parte și din cultura intangibilă. În categoria mare

a culturii, anumite obiecte pot aparține la două sau mai multe subcategorii, acestea nu sunt neapărat disjuncte.

III.b Uneltele fizice și informaționale, obiectele fizice și informaționale produse sau fabricate, utilizarea lor, instituțiile și organizațiile, care sunt consecințe, în cea mai largă măsură, a cunoașterii științifice, tehnologice, economice și organizaționale, poate chiar și a culturii intangibile.

Nu numai că unele obiecte culturale pot face parte din mai multe subcategorii ale culturii, dar vor exista și zone de interferență între obiectele acestor subcategorii. De exemplu, filosofia științei, care este un obiect al culturii, nu se poate dezvolta decât în strânsă legătură cu știința. În teoria categorică asemenea legături se numesc morfisme (morphisms sau maps, în limba engleză). Mai mult, pe lângă legăturile dintre obiectele subcategoriilor culturii, din ordinea a culturii ar proveni, există relații între aceste sfere în totalitatea lor. Aceste relații se numesc functuri. Cei mai importanți functuri sunt aceia dintre categorii care aparțin categoriei III de mai sus. Acești functuri,

F1 : Categoria III (Cultura științifică) + Categoria II (Cultura umanistă)

F2 : Categoria II (Cultura umanistă) + Categoria III (Cultura științifică)

reprezintă relația și influența reciprocă dintre, în esență, cultura umanistă și știința (cultura științifică). Importanța lor pentru societate și om nu trebuie subestimată.

Care este mai importantă dintre cele două categorii? Ambele sunt importante, dar motorul dezvoltării provine din sânul categoriei III. Acest lucru devine tot mai evident odată cu formularea conceptelor societății cunoașterii.

Este adevărat că o altă resursă importantă este viața spirituală, care este și componenta de creație implicând puternic atât cultura umanistă, cât și cultura științifică.

Odată cu era informației vor apare desigur multe elemente noi ale culturii datorită tehnologiei informației, cărții și documentelor electronice, internetului, tehnologiilor vorbirii, tehnologiilor bioelectronice și bioinformaticice, roboților artificiali și agenților inteligenți informatici, mediului ambiant inteligent, și conștiinței artificiale. Vor apare schimbări în viața intelectuală, socială și politică.

Ce se va mai petrece în cultură ?

În secolul XXI sunt posibile câteva evenimente majore care vor modifica viața omenirii:

- Prăbușirea ecologică a societății și a speciei umane, datorită deteriorării grave a mediului înconjurător, ceea ce s-ar putea întâmpla la mijlocul sec. XXI (să spunem, anul 2050) dacă nu se trec

repectiv de pe acum, la efortul de asigurare a unor societăți sustenabile. Salvarea este posibilă chiar cu cunoașterea științifică și tehnologică de astăzi dacă se trece la un management adecvat al cunoașterii [7] și la noi concepte economice adaptate sustenabilității. În această problemă au apărut și alte noi perspective care vor rezulta dintr-o serie de evenimente descrise în continuare.

- Dezvoltarea inteligenței artificiale până la depășirea inteligenței umane, ceea ce se va putea petrece între anii 2019-2035 sau chiar mai devreme [8], [9], [10], [11], [12].
- Apariția conștiinței artificiale, tot în cursul sec. XXI, după ce inteligența artificială va depăși inteligența umană, dar fără a putea preciza perioada.

Aceste două ultime evenimente presupun apariția unor noi specii inteligente, dar și noi specii conștiente, unele nebiologice (roboți umanoizi în topul unor *specii* de roboți mai puțin inteligenți care simulează animale (insecte, pisici, câini) și roboți construiți pentru anumite funcțiuni care să înlocuiască omul [8][13].

Speciile de roboți umanoizi inteligenți și de agenți software inteligenți, ambele egal de inteligente sau mai inteligente decât omul sunt uneori numite robo sapiens [13]. Într-o primă etapă, aceste specii nu vor avea conștiință, astfel cum are omul, datorită faptului că au numai o organizare structurală și nu una structural-fenomenologică [14]. Dar aceste specii vor interacționa puternic cu omul și societatea și se pune întrebare în ce măsură ele vor fi și artefacte culturale, nu numai prin faptul că fac parte din cultura omului, ci și prin participarea lor activă la cultură. Vor dezvolta cultura lor (într-o anumită măsură, da) sau vor intra în jocul marii culturi, participând la *cultura totală devenită din fenomen socio-uman, unul socio-uman- inteligentă/conștiință artificială*

Întrucât robo sapiens va avea cunoaștere și va participa la dezvoltarea științei și tehnologiei, chiar la dezvoltarea sa ca obiect tehnologic, el va participa cu siguranță la cultura științifică, poate chiar la anumite forme de cultură umanistă sau numai robotică. El poate fi implicat, prin cunoașterea culturii umaniste, să participe ceva mai pronunțat la această cultură. Când va trece de la inteligență la conștiință, o asemenea activitate ar putea fi mult mai pronunțată.

Probabil, între homo sapiens și robo sapiens vor exista relații de competiție și cooperare, dar acestea se vor dezvolta într-o societate comună, cel puțin până la o segregare care nu ar fi de dorit, în care spiritualitatea și creativitatea lui homo sapiens îi va conferi acestuia din urmă poziții inabordabile lui robo sapiens. Din momentul în care vor apare specii de *robo sapiens-conștient*, lucrurile se vor schimba din nou, cu efecte poate și mai dramatice pentru om și societate. Încerc să mă conving că ideile unei societăți a conștiinței ar putea fi benefice pentru un asemenea viitor care probabil nu va putea fi prohibit. Probabil, înspre un asemenea viitor și într-un asemenea viitor să fie rezolvată și sustenabilitatea unei societăți a conștiinței.

Este interesant de reluat aici câteva previziuni ale lui Kurzweil [10] asupra starea societății în anii 2019 și 2029.

Pentru anul 2019, în domeniul afacerilor și al economiei, prevede tra care în majoritate vor folosi persoane simulate, oamenii de afaceri vor folosi asistenți software care vor conduce tranzacțiile în numele lor. Locuințele vor dispune de roboți de întreținere. Cu aceste artefacte comunicarea se va face fără voce, deoarece vor dispune de o tehnologie a limbajului natural și a vocii de înaltă calitate. Oamenii vor avea relații cu persoane automate inteligente. Calitatea acestora de profesori, îngrijitori medicali, persoane de companie. Aceste persoane automate au și calități superioare omului în privința memoriei, dar, afirmă Kurzweil, 'ele nu sunt încă privite ca fiind egale cu oamenii în ceea ce privește subtilitatea personalității acestora'. Inteligența artificială este însă prezentă și împletită cu toate aspectele societății. Responsabilitatea omului va rămâne pe primul plan și nu a persoanelor (agenților) care îl ajută. Operele de artă vor fi realizate prin colaborarea dintre artiști umani și inteligențe artificiale. Prin pericol în societate îl vor constitui micile grupuri de oameni și inteligențe artificiale folosind comunicații criptate care nu pot fi descifrate. Acestea vor folosi tehnologii informatice și agenți de îmbolnăvire obținuți prin bioinginerie. Pe de altă parte, descifrarea relațiilor dintre genele genomului uman va permite o medicină bazată pe inteligența artificială pentru tratamentul și chiar eradicarea multor boli, pentru prelungirea considerabilă a vieții omului natural.

Pentru anul 2029, Kurzweil prognozează: în domeniul comunicării va continua să predomina, ca volum, comunicația dintre oameni și mașini. Populația umană va crește până la 12 miliarde de persoane reale, cărora li se asigură toate condițiile normale de viață. Populația umană și a inteligențelor artificiale va fi preocupată în primul rând, pentru crearea de cunoaștere, într-o puzderie de forme. Va fi menționate capacități ale omului care să nu fie preluate de mașini, deosebire netă nu mai există între lumea oamenilor și lumea mașinilor. Cunoașterea umană a fost transferată mașinilor și multe mașini au personalitate, îndemnată pe baze de cunoaștere preluând și cunoașterea umană. Implanturile neurale bazate pe inteligență artificială vor amplifica funcțiile cognitive ale omului. Kurzweil afirmă: 'A defini ceea ce înseamnă o ființă umană devine o chestiune semnificativă politică și de legislație. Creșterea rapidă a posibilităților mașinilor este controversată, dar nu există nici o rezistență față de ea. Deoarece la mașinile au fost proiectate pentru a fi supuse controlului uman, ele nu au putut fi o forță amenințătoare față de populația umană. Oamenii realizează că nu este posibilă dezangajarea civilizației devenită om-mașină de dependența de inteligența mașinilor. Crește discuția despre drepturile legale ale mașinilor, în special în cazul acelor mașini care sunt independente de oameni (care nu sunt introduse în creier uman). Cu toate că nu se recunoaște deplin, prin lege, influența evolutivă a mașinilor la toate nivelele de decizie asigură o protecție importantă a mașinilor

Kurzweil consideră calități ale mașinilor inteligente, care încă din anul 2029 pot fi persoane de artă în toate domeniile artei ('Mulți dintre artiștii de frunte sunt mașini'). Observăm însă că acest lucru ar presupune o stare de conștiință similară omului și prin manifestarea fenomenelor de qualia. Implicit, Kurzweil consideră că mașini inteligente complexe, structurale, pot avea asemenea stări și pot chiar participa la discuții filosofice pe baza experienței proprii. Vorbind de experiența subiectivă a unor astfel de mașini, aceasta ar însemna că asemenea mașini să fi trecut pragul de la inteligență la conștiință *numai pe baze structurale* încă din anul 2029. Ceea ce nu credem, în principiu, a fi posibil.

Într-adevăr, previziunile pe care oamenii de știință le fac privind dezvoltarea inteligenței artificiale spre conștiință artificială se bazează pe extrapolări ale științei structurale (complexitatea structurală de la un anumit grad în sus generează conștiință, acest lucru fiind considerat valabil începând cu creierul animalelor). Odată cu creșterea complexității artefactelor creiere electronice sau creierelor software se consideră că atunci când acestea ating complexitatea creierului uman se va produce de la sine conștiința artificială [8], [10], [11]. Uneori, unii dintre cei care susțin un asemenea punct de vedere au îndoeli asupra valabilității lui [12]. În viziunea unei filosofii integrative a științei [15],[16], conștiința nu se poate realiza numai din elemente structurale, fiind nevoie și de elemente fenomenologice [17]. Conștiințele artificiale vor pune probleme foarte mari speciei umane care cred că ar putea fi rezolvate în cadrul unei viitoare societăți a conștiinței. Aceasta va urma atunci societății cunoașterii în cadrul erei informației [18],[19].

Ce va fi cultura în societatea conștiinței, la care vor participa, dacă nu chiar vor predomina conștiințele artificiale? Dacă lucrul cel mai important, în cele din urmă, este continuitatea conștiinței create de om, atunci și culturii create de ea trebuie să i se asigure o continuitate.

Aceste considerații arată, dacă mai era nevoie de subliniat, cât de importante vor fi în sec. XXI cultura științifică și cultura umanistă, ambele având nevoie de o cultură filosofică adecvată.

Culturi, cultură pozitivă și cultură negativă. Polarizarea culturii în jurul cunoașterii

O cultură poate fi apreciată *pozitiv sau negativ*, în raport cu anumite criterii. Se pierde prea mult din vedere acest lucru. Există astăzi și o cultură a teroristilor (chiar și o știință a terorismului) o cultură a corupției care ne pune nouă românilor atâtea probleme, o cultură a hoților etc. Desigur, acestea pot fi numite sub-culturi, dar tot culturi sunt. Cultura are multe fațete.

Cultura negativă este o cultură deformată în raport cu criteriile civice și socio-umane.

În ultimii 12 ani, în societatea românească, pe lângă multe lucruri pozitive s-au accentuat, din nefericire, și fenomene negative îngrijorătoare: corupție, imoralitate, injustiție. Creșterea imoralității și a injustiției, a influențat până și viața academică din țara noastră. Avem nevoie și de un efort cultural pentru a elimina aceste flageluri din societatea noastră, pe lângă efortul dezvoltării economice.

Un exemplu de cultură pozitivă este arta. A cunoaște arta înseamnă a cunoaște dar a simți arta, a trăi arta, a avea nevoie de ea, a fi o parte din viața interioară, acestea înseamnă cultură umanistă adevărată.

Dar dacă cele de mai sus nu sunt însoțite de comportament civilizabil în civilizație socio-umană, cultura poate fi denaturată (rapturile de opere de artă în scopuri personale sau statale). Natura firească a culturii pozitive este aceea care susține civilizația socio-umană, spiritualitatea, cunoașterea și conștiința, în consecință urmă societatea cunoașterii și societatea conștiinței.

În privința relației dintre cultura umanistă și cultura științifică, astăzi nu se mai poate vorbi de cultură, cu înțelesul de cultură - în general, dar de gândul la cultura umanistă.

Cultura - în general, are o mult prea puternică componentă științifică (inclusiv tehnologică, economică, organizațională, politică) pentru a mai acționa ca asemenea simplificare, este adevărat, continuatoarea unei tradiții care astăzi este complet depășită. Cultura, respectiv cultura - în general, este cultura umanistă și cultura științifică, împreună, ultima având, ca și prima, un conținut extrem de bogat.

În spatele confuziei care se menține astăzi atunci când vorbim de cultură se întreține schisma dintre cele două culturi, datorită unor interese de grup. În etapa actuală a societății, cultura umanistă nu-și mai poate erija numele genului de cultură, de fapt nu ea, ci slujitorii ei care nu s-au adaptat la vremurile cunoașterii și societatea cunoașterii, înainte de trecerea la societatea conștiinței, cultura umanistă concentra în jurul cunoașterii. Iar tehnologia va fi un factor cultural al societății covârșitor încât va reveni poate la pozițiile ei mitologice din antichitate. În acest context, remarcăm:

În antichitate, la egipteni, zeul Ptah era privit ca patronul lucrărilor de metal (metalurgiști și fierari) și al artizanilor. Ptah era însă unul dintre cei mai mari zei, creatorul pământului, părintele zeilor și al începuturilor. Interesant ar fi să vedem "zeul începuturilor".

La grecii antici, echivalentul lui Ptah era Hefaistos, zeul focului, al meșteșugurilor, protectorul artizanilor. El nu mai era o divinitate primordială, o divinitate, fiul lui Zeus și al Herei, fiind căsătorit cu Afrodita. Se pare că la egipteni la greci, tehnologia nu mai păstra poziția începuturilor, dar avea totuși o importanță deosebită.

reprezentant divin. La romani, echivalentul lui Hefaistos era Vulcan, considerat zeul focului.

Decăderea poziției tehnologiei în cultură începuse din antichitate. Ea a continuat până în secolul XX când într-adevăr avea să cunoască un reviriment. Astăzi vorbim despre marile tehnologii și chiar despre o filosofie a tehnologiei, de care o serie de gânditori și filosofi au scris lucrări deosebit de interesante: Ernst Kapp, Friedrich Desauer, Jose Ortega y Grasset, Martin Heidegger ș.a. Este adevărat că au apărut și lucrări îndreptate împotriva tehnologiei (L.Mumford, J.Ellul ș.a.), declanșând ceea ce în secolul XX s-a numit dilema tehnologică.'

Revirimentul filosofic al tehnologiei în societatea cunoașterii, în secolul XXI, va fi un factor important în gândire, în general. Tehnologia va continua biologicul, culturalul și conștiința.

Ce va face omul? Marea lui înțelepciune va fi aceea de a pregăti în mod corespunzător viitorul [19]. Din ce în ce mai mult, gândirea filosofică va avea un rol hotărâtor în știință, politică, viața socială.

Există și vor exista culturi ale profesiilor, ale domeniilor cunoașterii, ale națiunilor, etniilor, grupurilor, ale comunităților constituite pe Internet, ale instituțiilor și localităților virtuale, ale mașinilor inteligente etc. Lumea devine tot mai pluriculturală. Probabil aceasta este trăsătura cea mai importantă a postmodernității [20].

Momentul actual ar trebui să fie acela al tendinței spre cunoaștere și cultură (cu înțelesul ei total) pentru întreaga populație a omenirii, fiecare zonă locală, geografică sau virtuală, trebuind să fie preocupată activ de realizarea concretă a acestei tendințe.

Referințe bibliografice

[1] Mihai Drăgănescu: *Lucrări despre cultură*:

- a. Mihai Drăgănescu, **Cultura și marile tehnologii**, conferință, Universitatea Populară de Vară "Nicolae Iorga", Vălenii de Munte -30 august 1996.
- b. Mihai Drăgănescu, **Perspectiva informațională a culturii**, Contemporanul, 27 mai 1983.
- c. Mihai Drăgănescu, **Dimensiunile europene ale culturii române**, expunere, Vălenii de Munte, 1992, publicată în *Academica*, 1992.
- d. Mihai Drăgănescu, **Arta și societatea**, cuvânt, Ploiești, 4 noiembrie 1991, publicat în *Academica*, 1991.

- e. Mihai Drăgănescu, **Criterii transpolitice și transmafioțe în cultura românească**, 18 mai 1997, *Caiete Critice*, 1997, nr.3-4, p. 145-147.
- f. Mihai Drăgănescu, **Spirit enciclopedic și enciclopedism**, conferință, Vălenii de Munte, 22 august 1993 (publicată în *Academica*, volumul autorului, *Cariatidele gândului*, Ed. Academiei Române, p. 163-168).

- [2] Ernst Cassirer, **Substanzbegriff und Funktionbegriff**, 1910; **Die Philosophie der Symbolischen Formen**, 1923-1929 (3 voi).
- [3] Oltea Misco, Elena Gheorghe, **Repere istorice în filosofia culturii**, *Revisia filosofie*, XLVII, Nr. 5-6, 2000, p.449-459.
- [4] Mihai Drăgănescu, **Categories and functors for the Structural Phenomenological Modeling**, *Proceedings of the Romanian Academy, Series A, Vol.1, No.2*, 2000, p.111-115.
- [5] Grigore T. Popa, **Reforma spiritului**, volum în editare, conținând lucrări ale acestui autor prezentate și publicate la Academia Română în anii 1940 (vezi și prefața: Mihai Drăgănescu, *O gândire asupra conștiinței, moralității și societății*).
- [6] UNESCO, definiția culturii intangibile, web-site UNESCO.
- [7] Mihai Drăgănescu, **Societatea Informațională și a Cunoașterii. Vectorii Societății Cunoașterii**, Academia Română, București, 9 iulie 2001, publicat în voi. coord. Florin Gh. Filip, *Societatea informațională-Societatea cunoașterii. Concepte, soluții și strategii pentru România*, Academia Română, 2002, p.43-112).
- [8] Moravec H., **Rise of the Robots**, *Scientific American*, December 1999, p. 93.
- [9] Moravec H., **Robot Mere Machines to Transcendent Mind**, Oxford University Press, Oxford, 1999.
- [10] Kurzweil R., **The Age of Spiritual Machines**, Penguin Books, 2000.
- [11] Broderick **D.Jhe Spike**, New York, 2002, paperback.
- [12] Buttazzo G., **Artificial Consciousness. Utopia or Reral Possibility?** *Computer (IEEE)*, July 2001, p.24-30.
- [13] Interviews of Menzel P. and D'Aluisio F., **Robo Sapiens. Evolution of a species**, MIT Press, Cambridge, Massachusetts, 2002.
- [14] Drăgănescu M., *Din lucrările despre minte și conștiință*:
 - a. Mihai Drăgănescu, **The Interdisciplinary Science of Consciousness**, *Noetic Journal*, Vol.3, No.1, Jan.2000, p. 3-10, republicat în eds. Richard L. Amoroso et al, *Science and the Philosophy of Consciousness, Intimation of a 21st Century Revolution*, Chapman & Hall, pp. 46-59, Orinda: The Noetic Press, 2000.
 - b. Mihai Drăgănescu, **Theories of Brain, Mind and Consciousness. Still Great Divergences**, *Noetic Journal*, vol.3, No. 2, Apr. 2000, p.125-139.

- c. Mihai Drăgănescu, *The Brain as an Information Processor*, NOESIS, XXV, 2000, p. 9-20.
- d. Mihai Drăgănescu, *On the Structural-Phenomenological Theories of Consciousness*, NOETIC JOURNAL, Vol.1, No.1, June 1997.
- e. Mihai Drăgănescu, *Continuities and Discontinuities in the realms of life and mind*, Revue Roumaine de Philosophie, Tome 41,1997, Nos 1-2, p.3-9.
- f. Mihai Drăgănescu, *De la filosofia la știința mentalului*, Revista română de filosofie, XLIV, Nr.5, sep-oct 1997, p. 457-464.
- g. Mihai Drăgănescu, *Procesarea mentală a informației*, Memoriile Sect. St. ale Acad. Române, SERIA IV, Tom. XX, 1997, p.263-284.

[15] Kafatos M., Drăgănescu M., *Preliminaries to the Philosophy of Integrative Science*, E-book (Microsoft Reader), ISBN 973-10-02510-X, Editura ICI, Bucharest, 2001.

[16] Drăgănescu M., Kafatos M., *Generalized Foundational Principles in the Philosophy of Science*, paper presented at the Conference on "Consciousness in Science and Philosophy" in Charleston, Illinois, 6-7 Nov 1998, published in The Noetic Journal, Vol.2, No.4, Oct. 1999, p. 341-350, republished in the voi. *Science and the Primacy of Consciousness, Intimation of a 21st Century Revolution*, Richard L. Amoroso and others (eds), Orinda: The Noetic Press, 2000, Chapter 9, pp. 86-98.

[17] Mihai Drăgănescu, *Advancement in Neural Engineering and Neuroelectronics Put Forward Artificial consciousness*, Communication at the INGIMED II Conference, Bucharest, Dec. 13, 2001; E-PREPRINT, MSReader Format, 2002.

[18] Mihai Drăgănescu, *Conștiința, frontieră a științei, frontieră a omenirii*, Revista de Filosofie, XLVII, nr. 1-2, 2000, p.15-22.

[19] Mihai Drăgănescu, *Societatea conștiinței, o viitoare etapă a erei informației. Vectorii societății conștiinței*, studiu, Academia Română, în pregătire.

[20] După Alain Fienckielkrant, apud [3], p.458-459.

Între lingvistica matematică și cea computațională

Solomon MARCUS

Secția de Științe Matematice a Academiei Române

solomon.marcus@imar.ro

Mă simt obligat să reacționez la un anumit mod de prezentare a evolei ideilor, în cea de a doua jumătate a secolului al XX-lea, în articolul [1] al d-lui Tufiș (de aici mai departe DT), membru corespondent al Academiei Române. Precizez de la început că nu contest interesul și utilitatea direcției de prezentare în [1]; am în vedere numai modul în care această direcție este pusă în relație cu alte cercetări dedicate limbajului.

Cităm din [1: 133]:

"Desprinzându-se din lingvistica formală, "lingvistica matematică" încercat dezvoltarea unor modele matematice de reprezentare a limbajului natural sau formale (în general al aspectului lor sintactic, gramatical), căutând soluții abstracte de modelare generativă de tip universal a ceea ce se presupune a fi nivelul cunoașterii științifice a anilor 1960) a fi facultatea limbajului.

Nu știu ce înțelege DT prin "lingvistica formală", o sintagmă nu foarte folosită în perioada de emergență a lingvisticii matematice; exista lingvistica structurală (altceva decât ceea ce ar putea fi lingvistica formală, adică baza formalizare în sensul logicii matematice moderne), care desigur a constituit una din sursele lingvisticii matematice (de aici mai departe LM), așa cum i se pot încadra și alte surse (biologice, logice, matematice, psihologice etc), dar factorul determinant în nașterea LM, în a doua jumătate a anilor '50, a fost dezvoltarea calculatoarelor electronice și, împreună cu ea, a primelor preocupări sistematice de automatizare (prescurtare a lingvisticii computaționale), numite atunci traducere automată și documentare automată, prelucrarea automată a limbajului, cu diverse variații lor în engleză (de exemplu, "machine translation"), franceză, rusă, germană, italiană etc. Din aceste preocupări s-au inspirat primele modele care au constituit noua disciplină a LM.

Vorbesc despre lucruri trăite. Punctul meu de plecare s-a aflat în lucrările unor Kulagina și Melciuk, puternic implicați în studiile de traducere automată din franceză, Yves Leclercq, implicat în problemele de documentare automată, și Hays, implicat în traducerea automată din rusă în engleză și reciproc, B. Van

cu preocupări de informatică lingvistică la Grenoble. De la ei, ca și de la alți autori similari, am preluat în bună măsură ștafeta pe care am căutat s-o duc mai departe. Ceea ce afirm despre mine este valabil pentru cei mai mulți cercetători din domeniul LM din anii 1950 și 1960, cum ar fi Maurice Gross, Masami Ito, A. Trybulec și mulți alții. Dubioasă mi se pare sintagma "soluții abstracte", probabil efectul unui obicei binecunoscut de a diaboliza abstractul.

În ceea ce privește sintagma "lingvistică formală", ea a căpătat o anumită utilizare în anii târzii 1960 și în anii următori, iar personal am folosit-o în unele lucrări, după cum se va vedea imediat, considerând-o oarecum echivalentă cu LM; dar chiar dacă nu acceptăm această echivalență, nu putem eluda faptul că lingvistica formală se află în imediata vecinătate a LM. DT pretinde ca LM "a încercat", sugerând astfel că ea a eșuat în tentativă de modelare a limbajului natural. Ceea ce este deocamdată numai o sugestie devine, după cum se va vedea, o certitudine pentru DT.

Într-adevăr, iată ce scrie mai departe DT ([1]: 133):

"Curând metodele lingvisticii matematice și-au atins limitele drept care, în anul 1966, la propunerea lui David Hays, domeniul de cercetare al limbajelor naturale, din perspectiva utilizării acestora în interacțiunea cu calculatoarele electronice, este individualizat sub numele de lingvistică computațională".

Chestiunea cu atingerea limitelor ține de domeniul umorului involuntar și trecem peste ea, dar nu ne miră, după ce am văzut la ce se reduce LM pentru DT. Nu mi-am imaginat niciodată că între LM și LC ar putea avea loc o competiție, prima definindu-se prin metodă (căci ce altceva este LM decât studiul limbajului cu ajutorul matematicii ?) iar a doua prin obiectivul pe care și-l propune. LM nu poate ignora problematica LC iar LC nu-și poate realiza proiectele fără LM. Probabil însă că DT lucrează cu o definiție specială a LM, pe care am dori s-o aflăm. Modul simplificator în care DT se referă la generativismul lingvistic, într-o logică binară care eludează faptul că în materie de modelare se lucrează cu grade de adecvare și relevanță, este însă simptomatic pentru viziunea sa limitativă în problema în discuție.

Crede DT că gramaticile lui Joshi, atât de importante în LC, puteau fi concepute fără să fi fost precedate de cele ale lui Chomsky ? Da, Chomsky a fost tot timpul foarte controversat, dar fără stimulentele sale nu știu ce ne-am fi făcut, inclusiv în LC și în LM, în ciuda faptului că el nu s-a prea referit explicit nici la LC, nici la LM. Faptul că gramaticile context free se află din nou, începând cu anii '80, în centrul atenției în LC nu spune ceva ? Iar faptul că aceleași gramatici (cu extensiunile lor) au marcat, încă din anii '60, teoria limbajelor de programare, domeniu în care ținta programării în limbaj natural se află în actualitate, nu este și el semnificativ ? LC are mai multe părți, mai multe orientări, mai multe niveluri de abstracție, care comportă criterii diferite de evaluare. DT îl asociază pe D. Hays la ideea sa privind falimentul LM și lansarea, drept consecință, a LC. Ca unul care a

cunoscut bine cercetările lui Hays (a se vedea frecvența citărilor numelui său în lucrările subsemnatului) și l-a cunoscut și personal foarte bine, fiind invitatul ca "plenary speaker" la Institutul de lingvistică al Americii (SUNY, Buffalo, 1966), pot depune mărturie că acest autor vedea în LM și LC două domenii solidare, două fețe ale aceleiași medalii, așa cum se va vedea din citatul pe care-l vom da mai jos. Desigur, Hays a avut un rol important în anii de pionierat ai LM și LC, dar ideea unei competiții între ele i-a fost străină. Voi evoca aici intervenția sa la o sesiune de a treia Conferință Internațională de LC (COLING, September 1971): "The history and scope of Computational Linguistics" [2]. Cităm ([2]:p.23):

"Solomon Marcus says that formal linguistics is a pilot science, emphasizing at the same time that the ordinary field of linguistics is not. But this is to say that linguistics as a branch of mathematics will supply methods to many other fields of science, whereas linguistics as a descriptive field, a branch of natural history or natural science, does not. [...] A four-way scheme can be arranged with psychology, computation, formal linguistics, and descriptive linguistics as the poles. Psychology and computation are about performance, formal and descriptive linguistics are about competence, computation and formal linguistics are about knowledge, and psychology and descriptive linguistics are sciences. But two other fields need to be found to find places in this scheme: psycholinguistics joins psychology with linguistics and seems at this time a most fruitful field, one in which great progress can be made with benefit to both parent fields. Correspondingly, on the abstract side, COMPUTATIONAL LINGUISTICS JOINS COMPUTATION WITH FORMAL LINGUISTICS (subl. mea, S. M.) and also seems a fruitful area, one in which rapid progress can be expected with benefit to both parent fields (subl. mea, S. M.) and with beneficial application to psycholinguistics".

Referirea pe care o face Hays la subsemnatul are în vedere sloganul pe care l-am folosit de mai multe ori, "formal linguistics as a pilot science", în care sintagma "formal linguistics" era folosită ca un echivalent al LM. Iată deci că Hays vedea în LC o alianță a LM cu computaționalul, alianță de natură să imprimă un progres rapid atât în LM cât și în domeniul computațional. Cei 30 de ani scurși până atunci au confirmat-o pe deplin. Denumirile folosite pentru preocupările de reducere a interferența limbajelor, informaticii și matematicii au variat tot timpul și nu creșterea în acest aspect merită prea multă atenție. Lingvistica matematică? computațională? inginerescă? algebrică? cognitivă? aplicată? cantitativă? teoretică? statistice? probleme matematice ale semioticii? tehnologia limbajului? limbajul în inteligență artificială? inginerie lingvistică? procesarea limbajului natural? "information storage and retrieval"? lingvistica cibernetică? pe fiecare dintre acestea am întâlnit în propriile mele articole au fost publicate aproape sub fiecare dintre etichetele mai sus. Iată și câteva detalii semnificative ale istoriei.

În 1962 s-a înființat în USA "Association of Computational Linguistics"

În 1963 Ferenc Kiefer a demarat la Budapesta revista "Computational Linguistics", care a trăit peste zece ani. Conferința de la Grenoble de "traitement automatique des langues" din 1967 era a treia de acest fel, fiind precedată de o alta, la New York, în 1965 și de una în Anglia, probabil în 1963, organizată de M. Masterman. Între timp, la ruși, numeroase conferințe au avut loc pe tema "avtomatetskaja obrabotka tekstov" iar "Sprachkunde und Informationsverarbeitung" a fost uneori eticheta folosită de germani ș.a.m.d. Nu negăm rolul important pe care l-a avut David G. Hays în dezvoltarea CL, dar acest rol a fost altul decât cel afirmat de DT. Emergența LC s-a produs încă din anii '50, sintagma LC a devenit curentă încă de la începutul anilor '60. Șirul de conferințe COLING nu a făcut decât să continue această tradiție. Alții au preferat folosirea sintagmei LM (a se vedea, de exemplu, "Prague Bulletin of Mathematical Linguistics", "Prague Studies of Mathematical Linguistics", revista japoneză "Mathematical Linguistics" (în echivalentul ei japonez) etc. în ceea ce privește însă profilul acestor reviste, nu am constatat o diferență față de cele de CL. Desigur, între timp au început să apară și unele publicații mai specializate, cu referire la părți determinate ale CL (cum ar fi cea relativă la corpusul lingvistic). Etichetele nu au avut importanță și nu știu să se fi desfășurat vreo competiție între ele. Chiar Hays a folosit diverse etichete, de exemplu cea din [3]. Dar DT merge mai departe pe ideea sa și afirmă (în completă discordanță cu viziunea lui Hays, de la care se reclamă) că "metodele LM și-au atins limitele" (încă în urmă cu peste 30 de ani!), pentru ca numai două pagini după această afirmație (deci la pagina 135 din [1]) să afirme că e nevoie de "modele formale ale limbii la toate nivelurile ei (fonetică, morfologie, sintaxă, discurs) gramatici formale [...]". Cum vede DT aceste modele formale altfel decât sub formă logico-matematică? Știe oare că multe modele de acest fel există de câteva decenii? Indicații bibliografice asupra lor sunt date parțial în [4], [5], [6], [7] iar pentru cercetările românești în [8], [9]. Desigur, aceste modele sunt inegale ca valoare, au nevoie de continuări, modificări, ameliorări, dar ele nu pot fi ignorate. Fonetica, fonologia, vocabularul, morfologia, sintaxa, semantica lingvistică și lingvistica istorică au beneficiat din plin de metodele matematice, așa cum se poate vedea din impactul deosebit al lucrărilor respective în literatura de specialitate; DT indică, drept domeniu al LM, numai "aspectul sintactic, gramatical", despre celelalte nu a aflat. Nu a aflat nici că LM a abordat și aspecte analitice, nu numai pe cele generative. DT definește "dimensiunea fundamentală" a LC prin "fezabilitatea instanțierii unei descrieri lingvistice cât mai complete, mentenabilitatea acestei instanțieri și, desigur, conformanța cu realitatea uzului limbii". ([1]: 133). Cu un mic efort înțelegem despre ce este vorba. Desigur că problemele de complexitate, de cost, nu puteau fi încă abordate în anii '50 și '60 cu mijloacele cu care ele au început a fi studiate în a doua jumătate a anilor 70, când instrumentele elaborate în informatica matematică deveniseră mult mai perfecționate. Dar acest fapt nu ține, cum crede DT, de alegerea între LM și LC, ci de progresul general realizat în știință. Pentru a mă referi la propria noastră experiență, atunci când, în

1969, prezentam la COLING-ul din Suedia gramaticile contextuale nu aveam să mă ocup de aspectul complexității acestor gramatici în maniera în care puteam face acest lucru ulterior (a se vedea, de exemplu, [10]). Dar acest fapt nu are nici o legătură cu eticheta folosită.

Anii '80 și '90 au confirmat necesitatea unui orizont cât mai larg în domeniul computațional. Nu m-am mirat atunci când "Encyclopedia of Microcomputers" și "Encyclopedia of Computer Science and Technology" m-au solicitat o contribuție cu tema "Semiotics and Formal Artificial Languages" (a se vedea [11]) și nici când "Handbook of Formal Languages" mi-a solicitat un capitol privind "Contextual Grammars and Natural Languages" [12] iar o lucrare preponderent teoretică a fost inserată în "Computational Linguistics in the Netherlands 2000"[13]. Nu m-am mirat nici când am văzut că o revistă de "Linguistics and Philosophy" publică articole excelente de LC. Interferențele între CL și LC în toate direcțiile și ele caracterizează cultura contemporană. În acest caz trebuie să ne plasăm, cred, atunci când ne referim la disciplinele cognitive care se dezvoltă sub ochii noștri și își pun amprenta pe modul nostru de gândire și de comportare. Un tratat ca "Mathematical Methods in Linguistics" [14] include multe fapte de LC, deși în titlul său nu figurează epitetul "computațional". O revistă de "Theoretical Linguistics" (1970-2000), publicată de Walter de Gruyter (Berlin-New York) a inclus multe articole vizând aspecte matematice și/sau computaționale, deși numele revistei nu indică acest lucru. Chiar o revistă mai tradițională de "Linguistics" a inclus de multe ori articole de LM și nici "Foundations of Linguistics" nu a procedat altfel. Multe fapte de LM și de LC se plasează în mod natural pe orizontul semioticii computaționale. Era internetului impune desigur o problemă nouă, față de care abordările anterioare se pot dovedi insuficiente. Să încercăm inițiativa noii generații de cercetători de a se dedica noilor probleme. Dar trebuie de la ieri la azi și de la azi la mâine nu poate fi decât una care ține seama în primul rând de critic de experiența acumulată. Din tot ceea ce am prezentat mai sus rezultă că LM și LC au fost mereu împreună și că, în general, etichetele nu au contat prea mult. Unii au mers chiar mai departe; astfel, în capitolul 4, "Mathematical Methods in Computational Linguistics", din [15], se afirmă pur și simplu (p.86): "Mathematical linguistics has also been called theoretical linguistics and even computational linguistics". Iar mai departe, în același loc: "Computational Linguistics originated around 1950 with the initiation of research on automatic translation" (se trimite la o carte editată de D.G.Hays [3] și la o alta avându-l ca autor pe acesta [16]).

Ca unul care crede în legătura naturală a lingvisticii cu matematica încercat o deosebită satisfacție să trăiesc momentul în care această legătură este acceptată de ambii parteneri și că de multe ori nici nu mai e nevoie de acțiune retorică al epitetului "matematică"; LM este acceptată pur și simplu ca lingvistică. Suntem convinși că o traiectorie similară o urmează și LC iar unele serbări pe această privință există de pe acum, așa cum am arătat mai sus.

LC este de mai mulți ani o secțiune la congresele internaționale de lingvistică iar LM și LC au secțiunea lor în reviste internaționale de referate ca "Language and Language Behavior Abstracts". În România, minți luminate ale anilor '60, ca profesorii Al. Rosetti, Grigore Moisil și Tudor Vianu, au înțeles schimbările care se profilau și au sprijinit proiectul înființării unei secțiuni de "lingvistică aplicată" la Facultatea de Limbă și Literatură Română a Universității din București, dar s-au găsit alții care să-l torpileze.

La Academia Română a funcționat mulți ani "Comisia de Lingvistică Matematică" iar revista "Cahiers de Linguistique Theorique et Appliquee", înființată în 1962, a fost multă vreme expresia colaborării lingvisticii cu matematica și cu informatica. În ciuda forțelor adverse, s-a reușit în acei ani atragerea unor studenți străluciți ai unor facultăți umaniste la cercetarea limbii și literaturii cu mijloacele LM și LC. Pentru a da numai două exemple de actuali profesori universitari care au susținut teze de doctorat de acest tip, voi menționa pe Pia Brinzeu, de la Catedra de Engleză a Universității din Timișoara și pe Mihai Dinu, de la Facultatea de Litere a Universității din București. Tot în acea perioadă și-a susținut teza de doctorat Sorin Cristian Niță, pe o temă de critică textuală automată privind înlănțuirea (filiația) diferitelor variante ale "Istoriei Țării Românești" (Șerban Cantacuzino).

Iată însă că, în pofida realităților puse în evidență mai sus, în ([1]: 134) se scrie: "În România, cercetările în domeniul LC și al prelucrării limbajului natural, precum și primele rezultate practice au apărut la începutul anilor '80 [3, 4, 5, 6]".

La ce trimit numerele indicate în paranteze? La o bibliografie de 24 de titluri în care aproape toate (dar toate cele indicate între paranteze) încep cu DT (ignorându-se regula generală în lumea științifică, a așezării numelor autorilor aceluiași articol în ordine alfabetică; dar nu acest fapt este cel care ne interesează în momentul de față). Să observăm că încă în 1978, în articolul "Mathematical and Computational Linguistics" [9] de prezentare a activității din România în domeniul LM și LC se face referire la peste 400 de articole publicate de 130 de autori români și sunt menționați peste 300 de autori străini (unii dintre ei, nume de vază ale LM și LC din acea perioadă) care au citat și continuat cercetările românești. Să mai adăugăm că numeroși lingviști români dintre cei mai importanți au citat și folosit rezultatele școlii românești de LM și LC. Iată că vine acum DT și face (deliberat sau nu) din tot acest efort un teren viran care-l aștepta pe DT să tragă primele jaloane. Nu e cam mult?

Să fim bine înțeleși. Nu noi avem nevoie de încă o citare pe lângă miile de citări deja acumulate, ci noile generații de studenți și de cercetători au dreptul la o informare corectă asupra dezvoltării LM și LC în general și, în particular, asupra LM și LC în România. DT a mai publicat, în urmă cu câțiva ani, un articol în care se schița o privire istorică asupra LC în România, cu câteva citări la întâmplare, care trădau necunoașterea situației reale.

Mai este un aspect care cere o precizare. În conformitate cu spațiul volumului în care apare articolul [1], DT face numeroase referiri la articole și documente ale unor organisme europene și internaționale, cum este și cazul nostru pentru a nu mai vorbi de aspectul financiar al colaborării cu organisme și instituții respective. Această situație a existat de la începutul LM și LC (chiar dacă nu s-a avut amplexarea de azi), datorită faptului că LM și LC au apărut și ca urmare a unor comandamente sociale, privind precaritatea mijloacelor de prelucrare a informației. Îmi amintesc de faimoasele Rapoarte CETIS care veneau de la EURASIA la Bruxelles, pe teme legate de analiză și prelucrarea automată a limbajului natural, traducere automată și documentare automată. În USA, diferite corporații (cum ar fi RAND Corporation, Santa Monica, Calif.) finanțau cercetări similare. O inițiativă semnificativă a fost aceea din 1962, organizată de "NATO Advanced Study Institutes", la Veneția, Italia, privind traducerea automată. De numele acestui proiect este legat un document care a marcat evoluția cercetărilor de traducere automată: seria de expuneri prezentate de Y. Bar-Hillel [17]. În legătură cu aceste activități, dirijate și finanțate de diferite organisme europene și internaționale, trebuie să observăm că cei implicați au avut înțelepciunea și priceperea necesare pentru a nu reduce proiectele respective la dimensiunea lor exclusiv utilitară, ci să se subordona pe aceasta unei perspective mai ample, care lua în considerare și nivelul orizontului științific real al problemelor. Pentru a da un prim exemplu, mă voi referi la faptul că mai multe rapoarte CETIS au pus în discuție un concept care, născut din experimentele de traducere automată, avea să se dovedească de o deosebită semnificație pentru teoria sintactică în toată generalitatea sa; este vorba de conceptul de proiectivitate sintactică, cu consecințe bogate în studiul structurilor arborescente și al gramaticilor de dependență. Azi putem spune că și în studiul limbajului natural și teoria matematică a grafurilor au profitat esențial de conștientizarea respectiv (folosit până și de Rene Thom, în probleme de morfogeneză). Această expansiune a unui concept sau rezultat dincolo de motivația sa imediată este testul cel mai convingător al interesului său. Un al doilea exemplu se referă la titlul provocator folosit de Bar-Hillel pentru expunerile sale: "Patru concepte despre lingvistica algebrică și traducerea automată".

Simpla alăturare a celor două sintagme, una foarte teoretică, cealaltă, aparent tehnologică, avea menirea să-i avertizeze pe cei care presau să se mulțumească cât mai repede rezultate practice asupra faptului că proiectele de traducere automată nu se pot finaliza de azi pe mâine, ci au nevoie de un lung proces de cercetare lingvistic, matematic și computațional. Acum știm că acest itinerar continuă să se desfășoare cu tatonări și reveniri, și, chiar dacă nu a dus încă la rezultatele visate, este impulsivat în mod esențial cercetările de AI, cu consecințe benefice în multe dintre aspectele logice și semantice ale limbajului natural.

Întrebarea pe care ne-o punem, dar o lăsăm deocamdată fără răspuns, este de ce deocamdată deoarece nu suntem pregătiți pentru a-l da, este următoarea: Nu cumva aspectele pe care le-am criticat mai sus sunt consecința unui fenomen mai general, a

unui orizont insuficient de cuprinzător, al unei prea mari dependențe de factori utilitari imediați? Știința a oscilat mereu între cognitiv și utilitar, dar istoria arăta că funcția utilitară s-a manifestat în toată profunzimea ei atunci când ea a fost fructul unei evoluții firești a funcției cognitive, evoluție care poate fi de doi ani, de 20 de ani, de 200 sau de 2000 de ani. Cu un ochi îndreptat spre comisiile europene, suntem obligați totuși să ținem treaz și celălalt ochi, îndreptat spre ceea ce se întâmplă pe scena cercetării științifice vii, așa cum apare ea în revistele de specialitate și la întâlnirile științifice de profil. Istoria generală a științei și, în particular, scurta istorie a LM și LC, sunt pline de învățăminte în această privință.

Referințe bibliografice:

- [1] D. Tufiș. *Promovarea limbii române în SI-SC*. în *Societatea Informațională - Societatea cunoașterii* (coord. F. Gh. Filip). Ed. Expert, București, 2001, 131-142.
- [2] D. G. Hays. *The field and scope of computational linguistics*. Papers in Computational Linguistics (eds. F. Papp, G. Szepe). Proceedings of the Third International Meeting of Computational Linguistics, held in Debrecen, Hungary, 1971. Akademiai Kiado, Budapest, 1976, 21-26.
- [3] D. G. Hays (ed.). *Readings in Automatic Language Processing*, American Elsevier, New York, 1967.
- [4] S. Marcus. *Mathematical Linguistics in Europe. Current Trends in Linguistics* (Th. A. Sebeok, ed.), vol.9, Mouton, The Hague, 1972, 646-687.
- [5] S. Marcus. *Mathematique et Linguistique*. în *Mathematique, Informatique et Sciences Humaines*, Paris, 26, 1988, 103, 7-21.
- [6] S. Marcus. *The status of research in the field of analytical algebraic models of language*. în *Current Issues in Mathematical Linguistics* (C. Martin-Vide, ed.). Elsevier - North Holland, Amsterdam, 1994, 3-21.
- [7] S. Marcus. *Lingvistica matematică, azi*. în *Matematica în lumea de azi și de mâine* (C. Iacob, coord.), Editura Academiei, București, 1985, 182-186.
- [8] S. Marcus. *Recent Romanian investigations in the field of mathematical and computational linguistics*. Avtomaticeskaja Obrabotka Tekstov, Matern. Fyz. Fakulta, KL Praha, 1973, 15-42.
- [9] S. Marcus. *Mathematical and computational linguistics*. în *Current Trends in Romanian Linguistics* (A. Rosetti, S. Golopentia Eretescu, eds.). Revue Roumaine de Linguistique 23, 1978, 1-4, 559-588.
- [10] S. Marcus, C. Martin-Vide, G. Paun. *Contextual grammars as generative models of natural languages*. Computational Linguistics 24, 1998, 2, 245-274.
- [11] S. Marcus. *Semiotics and formal artificial languages*. în *Encyclopedia of Computer Science and Technology* (A. Kent, J.C. Williams, eds.) 29, Ed. Marcel Dekker, New York, 1994, 393-405; also in *Encyclopedia of Microcomputers* (A. Kent, J.C. Williams, eds.) 15, 1995, 299-312.
- [12] S. Marcus. *Contextual grammars and natural languages*, *Handbook of Formal Languages* (G. Rozenberg, A. Salomaa, eds.), 2, Springer, Berlin, New York, 1997, 215-235.
- [13] S. Marcus, C. Martin-Vide, G. Păun. *A new-old class of linguistically motivated regulated grammars*. *Computational Linguistics in the Netherlands 2000* (W. Daelemans et al., eds.), Selected Papers from the Eleventh Meeting, Ed. Rodopi, Amsterdam, New York, 2001, 111-125.
- [14] B. H. Partee, A. Ter Meulen, R. Wall. *Mathematical Methods in Linguistics*. Kluwer, Dordrecht et al, 1990.
- [15] E. F. Beckenbach, Ch. B. Tompkins (eds.). *Concepts of Communication: Interpersonal, Intrapersonal and Mathematical*. John Wiley and Sons, New York, 1976.
- [16] D. G. Hays. *Introduction to Computational Linguistics*. American Elsevier, New York, 1967.
- [17] R. Thom. *Stabilite Structurelle et Morphogenese*. John Benjamins, New York, 1970.
- [18] Y. Bar-Hillel. *Four Lectures on Algebraic Linguistics and Machine Translation*. revised version of a series of lectures given in July 1962, before a Summer School at the Advanced Summer Institute, Venezia, Italy.

Între lingvistica matematică și cea computațională: o altă perspectivă

Dan TUFIS
Institutul de Cercetări pentru Inteligență Artificială,
Str. 13 septembrie, nr. 13, 74311, sector 5, București
tufis@racai.ro

1. în loc de introducere

Dat fiind că acest articol este un comentariu asupra filipicei de neînțeles "între lingvistica matematică și cea computațională" a domnului Solomon Marcus, membru titular al Academiei Române, mărturisesc că elaborarea sa fost o întreprindere asupra căreia am avut multe ezitări iscate din incertitudinea receptării sale corecte, constructive. Din păcate majoritatea afirmațiilor și implicațiilor pe care domnia sa le face în articolul amintit, sunt inexacte și umorale. Nu mai insist și asupra decontextualizării citatelor din lucrarea mea [1], procedeu neelegant. Este binecunoscut din logica clasică faptul că dintr-o serie de premise false se poate demonstra orice. În ciuda ezitărilor amintite, violenta polemică lansată de domnul Solomon Marcus prin articolul menționat îmi oferă posibilitatea de a aduce în discuție elemente de istorie a domeniului care ar putea fi de interes, cu precădere pentru cititorii al căror domeniu de specialitate nu este prelucrarea automată a limbajului natural. Pentru specialiștii în domeniul prelucrării limbajului natural, majoritatea argumentelor pe care le voi aduce sunt bine cunoscute.

Ca modalitate de documentare, am optat pentru includerea integrală a materialului produs de domnul Academician Marcus, indentat și redat cu caractere italice. De asemenea, am păstrat secțiunea domniei sale de referințe bibliografice. Lucrările pe care le-am citat eu sunt documentate în cuprinsul textului, prin includerea referinței complete între paranteze rotunde. Singura excepție este lucrarea mea, sursa nemulțumirii domnului Marcus, care este referită de amândoi ca [1]. Cititorul va putea face astfel mai ușor distincția între cele două categorii de referințe. Înainte de a proceda la analiza afirmațiilor domnului Academician Marcus, aș dori să fac unele precizări:

- contextul discuției în [1], ca și în cele ce urmează, este cel al tehnologiei limbajului, al cercetărilor foarte intense în întreaga lume

pentru dezvoltarea de sisteme inteligente capabile să faciliteze comunicarea dintre doi sau mai mulți conlocutori (oameni sau sisteme software), prin intermediul limbajului natural;

În raport cu lucrarea [1] domnul Academician Marcus se oprește cu îndârjire asupra a doar trei fraze interpretate ca atac la persoana sau activitatea sa științifică și se referă ironic (și după cum se va vedea în continuare, în mod nejustificat) la alte două, făcând abstracție de restul prezentării care nu are nici o contingență cu domnul Marcus. Domnul Academician are merite pe care nu i le poate lua nimeni, are contribuții importante în mai multe domenii și este creatorul școlii românești de lingvistică matematică. Interesul domniei sale pentru aspectele legate de implementarea pe calculator a programelor de prelucrare a limbajului natural a fost minim. Îmi reamintesc o discuție pe care am avut-o în anul 1991 la câțva timp după ce mă întorsesem de la Conferința Europeană de Lingvistică Computațională organizată la Berlin de profesorul Jurgen Kunze. Cu acea ocazie, domnul Academician Marcus mi-a mărturisit că îl cunoaște de multă vreme pe profesorul Kunze și că au și colaborat o perioadă cât amândoi au avut ca domeniu de preocupări lingvistica matematică. La sfârșitul anilor '60, mai spunea domnul Marcus atunci, drumurile celor doi s-au despărțit, profesorul Kunze optând pentru noua paradigmă a lingvisticii computaționale.

Domnul Academician Marcus a scris enorm, în domenii extrem de variate, aici mă refer în special la cele legate de studiul limbii, și prin urmare era inevitabil să nu atingă subiectul foarte actual al prelucrării automate a limbajului natural. A făcut-o însă detașat de nivelul inerent perisabil al tehnologiei informatice. O teorie științifică, un model formal teoretic sau transpus într-o implementare a unui program software sunt inevitabil supuse „eroziunii” timpului, unele mai rapid altele mai lent. Lucrarea [1], despre care discutăm, ia în discuție exact acest cadru al investigației tehnologice și a măsurilor științifice, tehnice, organizatorice și chiar legislative pentru a crea o bază perenă a cercetării și dezvoltării tehnologice privind prelucrarea automată a limbii noastre: resursele computaționale fundamentale ale limbii române. Societatea Informațională-Societatea Cunoașterii este caracterizată de vectori tehnologici și funcționali [M. Drăgănescu: „Societatea informațională-societatea cunoașterii. Vectorii societății cunoașterii” In *Societatea Informațională - Societatea cunoașterii* (coord. F. Gh. Filip). Ed. Expert, București, 2001, 43-112.] a căror ignorare este nu numai neproductivă dar și periculoasă, „în era electronică, **este esențial pentru supraviețuirea unei limbi ca ea să fie folosită în sistemele de informare electronică**” afirmă fără echivoc Alain Danzin în influentul raport al Comisiei Europene „Towards a European Language Infrastructure” întocmit în 1992 prin consultarea a

182 de specialiști din cercetare și industrie. Promovarea limbii române în contextul informațional al societății cunoașterii este un obiectiv actual și de viitor și nu poate fi subiect de dispută în viața științifică românească;

- deși este un truism, cred că pentru evitarea unor interpretări greșite este necesar să subliniez faptul că în dezvoltarea programelor de inteligență artificială, de prelucrare a limbajului natural sau în general în ingineria software, o mulțime de discipline matematice (teoria algoritmilor, teoria complexității, teoria limbajelor formale, teoria categoriilor, statistica matematică și multe, multe altele) sunt fundamente indispensabile în avansul științific și tehnologic al acestor discipline (și desigur nu numai al lor). Programarea (ca și matematica elementară) sau utilizarea de produse informatice sunt activități la îndemâna tuturor (de altfel reflectate și în programele școlare de învățământ), dar proiectarea și realizarea de programe software inteligente necesită o pregătire teoretică solidă, talent și multă muncă. Diferența între două programe care produc aceleași rezultate dar unul în câteva secunde și altul în câteva ore, apare tocmai din diferența de pregătire teoretică și talent a autorilor lor.
- domeniul științei și tehnologiei informației este poate cel mai dinamic sector al activității creative: Bill Gates spunea că dacă de pildă industria automobilelor ar fi avut aceeași dinamică cu cea a calculatoarelor, acum o mașină ar trebui să coste 1 dolar. Fantasticul ritm de dezvoltare al tehnologiei hardware (bazată pe importante descoperiri științifice obținute în ultimii 50 de ani) nu a fost nici pe departe egalat de ritmul dezvoltării în domeniul software. În ciuda acestui decalaj, știința ingineriei software și-a reînnoit instrumentarul teoretic (modele și/sau formalisme) cu o viteză neîntâlnită în alte domenii științifice. Dinamica fără precedent a cunoașterii în știința și tehnologia informației obligă omul de știință din acest domeniu la o informare continuă, din ce în ce mai specializată și mai selectivă. Se estimează că în acest domeniu se scriu în fiecare zi mai multe articole decât poate citi un om în întreaga sa activitate și că informația mai veche de 15-20 ani este foarte probabil să fie perimată (desigur cu excepțiile ce întotdeauna confirmă regula). Evoluția terminologică în acest domeniu este încă o mărturie vie a dinamicii de care aminteam: în domeniul prelucrării limbajului natural se vorbește acum de ontologii lexicale, de gramatici lexicalizate susținute de ontologii, de analiză (parsing) ontologică, de lingvistica WEB-ului și WEB-ul semantic, de resurse lingvistice standardizate și așa mai departe.
- referitor la antinomia „lingvistică matematică-lingvistică computațională” pe care domnul Academician Marcus mi-o atribuie, vreau să precizez că nicidecum nu am afirmat că cele două domenii se exclud reciproc sau că

ar fi în competiție; pur și simplu ele sunt subsecvente din punctul de vedere al relevanței față de problemele pe care le discutăm aici. Există fără îndoială o filiație între ele, în sensul că lingvistica computațională a preluat o mare parte din instrumentarul lingvisticii matematice (nici nu se putea altfel) dar ce a adus nou lingvistica computațională, pe lângă noi modele și formalisme, este în primul rând de natură metodologică și tehnologică: experimentul și evaluarea. Ceea ce se numește astăzi lingvistică computațională teoretică este în mare măsură asimilată cu lingvistica formală modernă. Acest segment al lingvisticii computaționale a moștenit de la lingvistica matematică cel mai mult și adecvându-și metodele la realitățile tehnologice a produs și este de așteptat să producă noi rezultate validabile și incorporabile în sisteme automate de prelucrare a limbajului natural. Teoriile și formalismele lingvistice, azi în vogă în lingvistica computațională (TAG, LFG, HPSG, CG, CUG), au fost produse de lingvistica formală și prin validarea instanțierilor pe segmente de limbă netriviale, au devenit instrumente operaționale ale prelucrării limbajului natural. Dezvoltarea de modele de limbă, analiza algoritmilor de prelucrare a limbajului (resursele de calcul necesare unei implementări funcționale, viteza de răspuns), construcția (achiziția) resurselor lingvistice standardizate, gradul de acoperire lingvistică al unei formalizări lingvistice (cunoștințe lingvistice=resurse lingvistice), sunt doar câteva direcții definitorii ale metodologiei lingvisticii computaționale, în sfârșit, în raport cu obiectivele finale urmărite de implementarea unui model de prelucrare a limbajului se remarcă în ultimii circa 10 ani o departajare și chiar o competiție (fără însă a fi o antinomie) între abordările introspective-principiale și cele inductive, bazate pe date. Prima categorie de abordări este caracterizată de dezvoltarea prin introspecție științifică de teorii și formalisme gramaticale computaționale (imensa lor majoritate bazate pe restricții și unificare categorială cu accentuată lexicalizare) și mai apoi instanțiate manual de experți lingviști. Cea de a doua abordare, ce câștigă foarte mult teren în ultima perioadă, este cea bazată pe tehnicile învățării automate ce pornesc de la premiza că, într-un corpus lingvistic reprezentativ și de dimensiuni mari, există suficientă informație privind regularitățile dintr-o limbă (cea în care sunt textele ce alcătuiesc corpusul lingvistic) astfel încât, tehnici adecvate de învățare automată să fie capabile să construiască un model de limbă robust și de mare acoperire lingvistică. Aș mai menționa că, în fapt, de multe ori cele două abordări sunt combinate (cu preponderența uneia dintre ele). Într-un anumit sens, acest dualism în abordările modelelor de prelucrare automată a limbajului natural continuă o celebră confruntare de idei între Chomsky și Piaget susținătorii teoriilor înăscutului (innate) și respectiv al învățării în explicarea facultății umane a limbajului.

Cu aceste lămuriri preliminare, voi analiza în continuare afirmațiile domnului Academician Marcus cu sincera speranță că cititorii acestui text, dar mai ales domnia sa, vor înțelege că preocupările mele și ale distinsului profesor au alte obiective, motivații și desigur modalități foarte diferite de finalizare. Acest lucru nu înseamnă că rezultatele fiecăruia dintre noi le anulează sau le diminuează pe ale celuilalt (cu atât mai mult cu cât recunoaștere internațională există pentru amândoi). După cum la fel de bine diferențele de perspectivă și opinii, naturale în fond, nu înseamnă că nu avem a ne spune lucruri interesante unul altuia.

2. O analiză textuală

„Mă simt obligat să reacționez la un anumit mod de prezentare a evoluției ideilor, în cea de a doua jumătate a secolului al XX-lea, în articolul [1] al d-lui Dan Tufiș (de aici mai departe DT), membru corespondent al Academiei Române. Precizez de la început ca nu contest interesul și utilitatea direcției de preocupări prezentate în [1]; am în vedere numai modul în care aceasta direcție este pusă în relație cu alte cercetări dedicate limbajului.”

Așa își începe domnul Academician Marcus articolul solicitat de mine pentru volumul „Limba Română în Societatea Informațională-Societatea Cunoașterii” rezultat al proiectului INFOSOC „SI-SC: Soluții și strategii în România”. Să urmărim un prim citat incriminat (care în transcrierea dlui Academician este trunchiat și conține niște ghilimele ce nu-mi aparțin; redau mai jos varianta publicată):

[1: p.133]:

”Din acest punct de vedere (al folosirii calculatorului în prelucrarea limbajului natural - precizarea mea, DT), este semnificativ a arăta că însuși numele domeniului de cercetare a prelucrării automate a limbajului natural a suferit modificări reflectând progresele științifice și tehnologice: inițial, desprinzându-se din lingvistica formală, lingvistica matematică a încercat dezvoltarea unor modele matematice de reprezentare a limbajelor naturale sau formale (în general al aspectului lor sintactic, gramatical), căutând soluții abstracte de modelare generativă de tip universal a ceea ce se presupunea (la nivelul cunoașterii științifice a anilor 1960) a fi facultatea limbajului.”

Ce l-a supărat aici pe distinsul polemist? Ne spune chiar domnia sa:

„Nu știu ce înțelege DT prin ”lingvistica formală”, o sintagmă nu prea folosită în perioada de emergență a lingvisticii matematice; există lingvistica structurală (altceva decât ceea ce ar putea fi lingvistica formală, adică bazată pe formalizare în sensul logicii matematice moderne), care desigur a constituit una din sursele lingvisticii matematice (de aici mai departe LM), așa cum i se pot indica și alte surse (biologice, logice, matematice, psihologice etc.)

Mă surprinde întrebarea retorică cu care începe „argumentația”, și căreia nu-i văd decât un gratuit rol derogativ. Eu nu-mi închipui că domnia sa nu a auzit de antinomia „gramatică descriptivă - gramatică formală” la limitele extreme ea fiind reprezentată de lucrările lui O. Jespersen (O. Jespersen: *The philosophy of Grammar*, Allen & Unwin, London, 1924 și *Analytical Syntax*. Hoit Rinehart & Winston, New York, 1937 (republicată în 1969)) și respectiv lucrările timpurii ale lui Chomsky referitoare la lingvistica generativă. Dacă însă mă înșel, o lectură lămuritoare, este influența carte editată de Keith Brown și Jim Miller în Pergamon Press, 1996 numită „Concise Encyclopedia of Syntactic Theories”, cu precădere articolul „Descriptive Grammar and Formal Grammar” de F. Stuurman, al cărui prim capitol se numește chiar Descriptive and Formal Grammar: The Fundamental Opposition. La fel de utilă este și lucrarea monumentală a lui David Crystal „The Cambridge Encyclopedia of Language”, Cambridge University Press, 1987.

Pe de altă parte, o pagină mai încolo, domnul Academician mărturisește că și domnia sa a folosit termenul de lingvistică formală:

În ceea ce privește sintagma "lingvistică formală", ea a căpătat o anumită utilizare în anii târzii 1960 și în anii următori, iar personal am folosit-o în unele lucrări, după cum se va vedea imediat, considerând-o oarecum echivalentă cu LM (lingvistica matematică);

Pentru lămurirea elementului istoric, furnizez în continuare un citat din recenzia lui R.B. Lees (Language, nr. 33, voi 3, '1957, pp. 375-408) la faimoasa carte a lui Chomsky (Syntactic Structures, Mouton, The Hague, 1957): „in a sense, transformational analysis is essentially a **formalization** of a long-accepted, tradițional approach...”. Citatul apare la pagina 387. Chomsky se pare că ă apreciat termenul și i-a adoptat, cel puțin în raport cu propria filozofie generativistă asupra limbajului.

„dar factorul determinant în nașterea LM, în a doua jumătate a anilor '50, a fost dezvoltarea calculatoarelor electronice și, împreună cu ea, a primelor preocupări sistematice de LC (prescurtare a lingvisticii computaționale), numite atunci traducere automată, documentare automată, prelucrarea automată a limbajului, cu diverse variante ale lor în engleza (de exemplu, "machine translation"), franceză, rusă, germană, italiană etc. Din aceste preocupări s-au inspirat primele ^ modele care au constituit noua disciplină a LM.”

Înainte de a face o serie de precizări istorice mai exacte, vreau să notez că de la începutul istoriei sale, domeniul traducerii automate a fost, și în mare măsură a și rămas, un domeniu distinct de restul preocupărilor legate de prelucrarea limbajului natural. Aș mai observa că textul de mai sus, încearcă să sugereze că LM s-ar fi constituit ca disciplină ulterior LC. Ambiguitatea afirmației de mai sus provine din punerea în relație de concordanță temporală a primelor preocupări în domeniul LC cu apariția domeniului în sine. Oricine știe că un anumit domeniu științific se cristalizează în timp, pe baza unor rezultate științifice promițătoare, a unor experimente convingătoare (în cazul domeniilor tehnologice). Până la

sedimentarea elementelor definitorii ale unui domeniu de cercetare, pot coexista sau se pot succeda mai multe direcții de cercetare. Dintre acestea unele dispăreau sau își pot diminua foarte mult influența în raport cu motivația inițială și pot continua însă existența prin noi motivații, prin alegerea de noi obiective

Ca element istoric, aș preciza că în toate evocările pe care le-am citat până acum cel ce pentru prima dată a sugerat ideea folosirii calculatorului și a tehnicii de decodificare pentru prelucrarea automată a limbajului natural a fost Warren Weaver în 1946. În 1949 el scrie lucrarea „**Translation**” considerată de toți specialiștii în traducere automată ca primul document programatic al acestei discipline. În 1952 a avut loc la Universitatea Georgetown din SUA prima conferință dedicată exclusiv traducerii automate. În 1954, Peter Toma de la Universitatea Georgetown, împreună cu un grup de cercetători de la IBM, realiza primul experiment de traducere automată (engleza-rusa) folosind un dicționar de 250 de cuvinte și reguli sintactice de rescriere. Acest sistem avea să constituie nucleul faimosului program de traducere automată Systran pe care Peter Toma îl finalizează în

Punctul meu de plecare s-a aflat în lucrările unor Kulagina, Melciuk, puternic implicați în studiile de traducere automată rusă-franceză, Yves Leclercq, implicat în problemele de documentare automată, D. G. Hays, implicat în traducerea automată din rusă în engleză și reciproc, B. Vauquois, cu preocupări de informatică lingvistică la Grenoble. De la ei, ca și de la alți autori similari, preluat în bună măsură ștafeta pe care am căutat s-o duc departe. Ceea ce afirm despre mine este valabil pentru cei mai mulți cercetători din domeniul LM din anii 1950 și 1960, cum ar fi Masami Ito, A. Trybulec și mulți alții.

Traducerea automată, dar mai ales eșecul primelor încercări de realizare a acestui obiectiv încă nerezolvat sau nerezolvat complet, a constituit fără îndoială o motivație a „emergenței” LM. Așa cum voi arăta pe larg mai departe, proiectele de traducere automată au fost puse, prin interpretarea unilaterală și tendențioasă a raportului APLAC, exclusiv pe seama inadecvării teoriilor lingvistice folosite atunci și a cantonării în fapticul unor limbi particulare. Teoria „înnăscute a limbajului” lansată de Chomsky, opunându-se tradiției tipologice de studiu lingvistic prin diversitatea limbilor, a generat o prodigioasă cercetare * direcția determinării principiilor gramaticii universale, în speranța că identifierea și caracterizarea lor riguroasă le-ar putea operaționaliza atât pentru explicația comunicării umane prin limbaj cât și (un derivat subsidiar al obiectivului Chomsky) pentru realizarea de sisteme de traducere automată aprofundate și performanța umană.

Dubioasă mi se pare sintagma "soluții abstracte", probabil etimologică și a unui obicei binecunoscut de a diaboliza abstractul.

Remarca de mai sus mă surprinde de două ori: mai întâi pentru că este ceva atât de nimic reproabil în expresia „o soluție abstractă” (ba chiar dimpotrivă: „așa ceva” Care rezultă din separarea și generalizarea însușirilor caracteristice ale

de obiecte sau de fenomene care este considerat independent, detașat de obiecte, de fenomene sau de relațiile în care există în realitate" DEX'96) și apoi referirea la un obicei binecunoscut (al cui?) de diabolizare a abstractului. Nu neagă nimeni că acele soluții abstracte de care aminteam au generat idei valoroase și cercetări computaționale (mai ales în domeniul traducerii automate bazate pe conceptul „interlingua”) dar rezultatele acestor idei și cercetări nu sunt revendicate nici chiar de Chomsky.

În ceea ce privește sintagma "lingvistică formală", ea a căpătat o anumită utilizare în anii târzii 1960 și în anii următori, iar personal am folosit-o în unele lucrări, după cum se va vedea imediat, considerând-o oarecum echivalentă cu LM; dar chiar dacă nu acceptăm această echivalență, nu putem eluda faptul că lingvistica formală se află în imediata vecinătate a LM.

Cu amendamentele cronologice pe care le-am comentat mai devreme, apropierea între LM și LF (lingvistica formală) este exact ceea ce am afirmat și eu.

DT pretinde ca LM "a încercat", sugerând astfel ca ea a eșuat în tentativa de modelare a limbajului natural.

În primul rând este vorba de modelarea computațională a limbajului. În al doilea rând nu eu pretind acest lucru, dar sunt perfect de acord cu el. Iată câteva opinii ale unor mari specialiști, activi, din domeniul prelucrării automate a limbajului natural (sublinierile îmi aparțin):

- Christopher Manning and Hinrich Shutze: Foundations of Statistical Natural Language Processing, The MIT Press, 1998:

....the availability of large text corpora has changed the scientific approach to language in linguistics and cognitive science. Phenomena that were not detectable or seemed uninteresting in *studying toy domains and individual sentences* have moved into the center field of what is considered important to explain."

- Susan Armstrong-Warwick (editor): Prefața la „Special Issue on Using Large Corpora”, Computational Linguistics, Volume 19, no 1, 1993 p. 4:

„What is that has brought about this rapid growth of interest in corpus-based NLP?...The technological advances in computer power has certainly favoured the approach, as has the growing availability of large-scale textual resources in machine readable form. *More important, perhaps, is the growing frustration of trying to use standard rule-based methods to account for more than a well-chosen fragment of text, regardless of the application.* The data extracted from large corpora have demonstrated that language is more flexible and complex than that which most rule-based systems have up to present tried to account for. The relative lack of practical results at a time when industrial concerns are looking to the CL community to demonstrate progress toward useful applications has also contributed to the growing interest in new methods.

And finally, the success rate demonstrated in the speech community offers hope for similar progress in NLP."

- Nancy Ide and Jean Veronis (editori) Computational Linguistics -Special Issue on Word Disambiguation, Voi. 24, No. 1 1998 p.15:

„Although quantitative methods were embraced in early MT work, in the mid-1960s interest in statistical treatment of language waned among linguists due to the trend toward the discovery of *formal linguistic rules* sparked by the theories of Zellig Harris (1951) and bolstered most notably by the transformational theories of Noam Chomsky (1957). Instead, attention turned toward full linguistic analysis and hence to sentences rather than texts, and toward contrived examples and artificially limited domains instead of general language."

- Victor Yngve: From Grammar to Science: New Foundations for General Linguistics, John Benjamin Publishing Company, 1996:

„there seems to be no scientific way of deciding among the many contenders...We find positions and methods being promoted like a new movie or defended with withering polemics or taken up like the latest fad...We should abandon logical-domain theories entirely and move to physical domain...Because this (notation) can be programmed on a computer it can be used to test large-scale models...Gone will be the babe of arbitrary grammatical notations, each to be discarded in turn".

Deși nu împărtășesc în întregime poziția extrem de radicală a lui Yngve, este simptome pentru insatisfacția generală față de abordările tradiționale anilor '60-80.

- R.F. de Bruine (editor) „Synthesis of Proposal for an RTD Programme Users, Industry and Research in Language and Technology", DGXIII, Commission of the European Communities, September 1992:

„There is a broad need to further understanding of linguistic phenomena in the context of computerising the analysis and generation of language. General research should be stimulated within the following three main topics:

- research on the linguistic meaning representation at the various levels of description, ranging from the lower (e.g. phonological, morphological and syntactic) and better understood ones to the higher, scientifically more difficult ones (e.g. semantic, pragmatic, contextual and communicative ones). It is foreseen that the former must yield results in the short to medium term. Even if the latter are long-term enterprises, they must be organised in a way that ensures availability of usable intermediate results.
- research on more adequate and efficient computational schemes for natural language processing (e.g. constraints based computing and quantitative aspects) providing the base for research

processing behaviour vz the applications of advanced computer science and statistical methods in close collaboration and synergy with related actions.

- research into the human factors related with the future spread of advanced language processing technologies taking into account the ergonomics aspects, economic and socio-cultural dimensions."

Lista unor astfel de citate poate continua pe zeci de pagini, dar am să mă opresc aici nu înainte de a mai reaminti raportul comisiei prezidate de Alain Danzin „Towards a European Language Infrastructure”. Acest document, o adevărată cartă albă a cercetării în domeniul tehnologiilor limbajului, a restructurat complet programele de cercetare și prioritățile pe termen mediu și lung. A o ignora (ba chiar mai mult a o critica fără a-i cunoaște conținutul și a o eticheta ca pe un document birocratic al celor de la Uniunea Europeană) poate fi desigur o opțiune personală, dar cu efectul izolării științifice și mai accentuate.

Ceea ce este deocamdată numai o sugestie devine, după cum se va vedea, o certitudine pentru DT. într-adevăr, iată ce scrie mai departe DT([1]: 133):

"Curând metodele lingvisticii matematice și-au atins limitele drept care, în anul 1966, la propunerea lui David Hays, domeniul de cercetare al limbajelor naturale, din perspectiva utilizării acestora în interacțiunea cu calculatoarele electronice, este individualizat sub numele de lingvistică computațională".

Chestiunea cu atingerea limitelor ține de domeniul umorului involuntar și trecem peste ea, dar nu ne miră, după ce am văzut la ce se reduce LM pentru DT.

În ciuda repetatelor mele clarificări, și după cum se observă și din citatul de mai sus, referirea mea era la utilizarea metodelor lingvisticii matematice în programele de prelucrare a limbajului și nicidecum la domeniul în sine. Probabil că pentru cine nu a încercat să realizeze un sistem de prelucrare a limbajului natural și nu s-a lovit de problemele implementării unui dicționar și a unei gramatici computaționale e mai greu de înțeles remarca mea anterioară. Domnul Academician Marcus nu s-a apropiat niciodată de problemele unei implementări și prin urmare nu mă surprinde lipsa de înțelegere a diferenței între o definiție formală a unei gramatici (de exemplu) care se explicitează în câteva rânduri și implementarea unei gramatici computaționale care nu numai că nu încapă în câteva sute sau mii de pagini dar reclamă o muncă exprimată convențional în mii de oameni/an. Gramatica computațională a limbii engleze, dezvoltată în cadrul proiectului Alvey, a fost rezultatul a 10 ani de muncă intensă a celor mai importante 12 colective de cercetare din Anglia, fiecare dintre acestea fiind conduse de cercetători importanți și fiind suplimentate cu numeroși studenți doctoranzi. Gramatica GPSG dezvoltată este unul din exemplele standard de gramatică introspectivă de mari dimensiuni. Un astfel de efort uman și financiar nu este la îndemâna multor societăți. Și experiența a arătat că nici nu este necesar! Ralph Grishman, de la Universitatea din New York a demonstrat că programul său

de inducție gramaticală, pe baza unui corpus de antrenare a generat o gramatică în nucleu, a cărei „finisare” a durat mai puțin de două săptămâni și, confruntat cu gramatica Alvey pe un text arbitrar a reușit să analizeze mai multe fraze și cuvinte a demonstrat o mai mare acoperire lingvistică.

Nu mi-am imaginat niciodată că între LM și LC ar putea avea loc o competiție, prima definindu-se prin metoda (căci ce altceva este decât studiul limbajului cu ajutorul matematicii ?) iar a doua prin obiectivul pe care și-l propune. LM nu poate ignora problematica și LC nu-și poate realiza proiectele fără LM. Probabil însă că LM lucrează cu o definiție specială a LM, pe care am dori s-o aflăm.

Nici nu există această competiție decât în imaginația Academicianului care sugerează mai sus că LC nu folosește matematică atunci când o face, disciplina se numește LM. Ceea ce, așa cum am spus înainte, este fals. Elementele suplimentare, esențiale și definitorii sunt calculul, algoritmi eficienți și cunoștințele cu care acesta trebuie „hrănit”. O formă de procesul de înțelegere și/sau producere a limbajului natural, de orice s-ar fi ea, nu este decât o ipoteză asupra unui fenomen încă neelucidat. Așa că această ipoteză este cheia care a diferențiat LC de LM. În anexa acestei lucrări am furnizat două definiții pentru LM și LC. Prima definiție (LM) aparține lui G. Pullum și Andras Kornai iar cea de a doua (LC) se află în pagina WEB a Academiei de Lingvistică Computațională (al cărui membru sunt din 1985). Așa că precizarea că lingvistica teoretică modernă (în sensul precizat mai sus) studiază limbajul nu numai cu ajutorul matematicii. Alături de matematică, sociologia, psihologia, medicina și științele cognitive constituie domeniul de cunoașterii care sunt fundamental implicate în explicarea acestui miracol. Acesta reprezintă comunicarea inter-umană. Incapacitatea actuală de a simula procesor artificial de limbaj la nivelul performanței și competenței umane datorează nedescifrării (încă) a mecanismelor minții și creierului omului la nivel structural-fenomenologic și noile cercetări în direcția unei științe a minții (reprezentată între alții de lucrările de pionierat ale Academicianului Drăgănescu) sunt fără îndoială porți deschise spre cunoașterea, înțelegerea exactă a minții și implicit a facultății limbajului. Până atunci, obiectivul (realizarea de sisteme automate capabile să prelucreze limbajul natural) se bazează la modele aproximative, a căror acceptabilitate se probează prin implementarea și evaluarea lor pe date reale. Cum între afirmarea unui obiectiv de LC și realizarea sa operațională este o distanță mare, pe care uneori cercetătorii fără experiență tehnologică programării fie că o ignoră, fie nu vor (și de multe ori sunt foarte interesați) să o parcurgă, confuzia ce duce la auto-acreditarea într-un domeniu conex este explicabilă.

Modul simplificator în care DT se referă la generativismul lingvistic este într-o logică binară care eludează faptul că în materie de lingvistică se lucrează cu grade de adecvare și relevanță, esențial este simptomatic pentru viziunea sa limitativă în problema în discuție.

Crede DT că gramaticile lui Joshi, atât de importante în LC, puteau fi concepute fără să fi fost precedate de cele ale lui Chomsky? Da, Chomsky a fost tot timpul foarte controversat, dar fără stimulentele lui nu știu ce ne-am fi făcut, inclusiv în LC și în LM, în ciuda faptului că el nu s-a prea referit explicit nici la LC, nici la LM.

Modul „simplificator” incriminat mai sus se referă la fraza „soluții abstracte de modelare generativă de tip universal”. Având în vedere că în articolul [1] aceasta este singura referire la generativism, bănuiesc că domnul Academician Marcus a vrut să spună „succint”. Apoi, continuarea ce se referă la logica binară pe care o folosesc în interpretare și simptomele viziunii mele limitative asupra problemei discutate desigur sunt efecte stilistice nereușite, întrucât nu am abordat (și nici nu mă interesează în mod deosebit) subiectul pe care îl invocă domnul Academician. Pentru că tot am ajuns aici, țin să-i reamintesc domnului Academician Marcus că Noam Chomsky și-a revizuit complet punctul de vedere care a dominat aproape 15 ani lingvistica mondială. Într-adevăr Chomsky este un mare om de știință, chiar dacă foarte controversat, dar acest statut îi este conferit și de onestitatea cu care s-a detașat de creațiile sale anterioare ce i-au adus notorietatea, dovedite (unele chiar de el însuși) ca fiind depășite, propunând soluții și teorii noi.

Formalismul TAG al lui Joshi este într-adevăr unul foarte important în LC ca și HPSG, LFG, CG și alte câteva. Dar dintre formalismele de lingvistică computațională, TAG este cel mai departe de influența chomskyană. Dacă se poate face o asociere între TAG și vreo teorie generativistă de tip chomskyan aceasta este doar de natură antinomică. Am colaborat cu profesorul Aravind Joshi în 1991 la Institutul Lingvistic de la Universitatea Santa Cruz din California, am fost apoi invitatul său la Universitatea din Pennsylvania, invitație motivată printre altele și de o deosebită apreciere pentru o demonstrație alternativă a mea, mai scurtă și, considerată de profesorul Joshi, mai elegantă a unei teoreme a domniei sale referitoare la categoria de limbaje acoperite de LTAG. Cu acea ocazie, profesorul Joshi mi-a pus la dispoziție trei volume consistente de lucrări asupra TAG tratând foarte amănunțit motivațiile lingvistice, proprietățile computaționale și caracterizarea matematică. Aceste volume i le-am pus la dispoziție și domnului Academician Marcus. Profesorul Joshi a fost în 1997 invitatul profesorului Dan Cristea și al meu la Școala de Vară EUROLAN unde a susținut o serie de prelegeri de înaltă ținută științifică. Am evocat aceste lucruri pentru a-l lămurii pe domnul Academician Marcus că formalismul TAG și varianta sa mai nouă LTAG îmi sunt familiare și prin urmare mă surprinde afirmația dânsului implicând o filiație între teoriile lui Joshi și Chomsky.

Faptul că gramaticile context free se află din nou, începând cu anii '80, în centrul atenției în LC nu spune ceva ?

Acest lucru este exact și ilustrează foarte bine ceea ce spuneam înainte: contextul computațional în care complexitatea algoritmică este primul mare judecător al adecvării unui model (inerent limitat, după cum arătam mai devreme) bazat pe o anumită teorie lingvistică. În anii de vârf ai lingvisticii matematice, și în

cei de început ai lingvisticii computaționale, pornindu-se de la o coniectură a Chomsky (limbajele naturale nu sunt limbaje independente de context) demontată în anii '80 de Gerald Gazdar (autorul teoriei GPSG), cercetarea a fost orientată spre identificarea de formalisme lingvistice cât mai puternice, cu puterea generativă mai apropiată de cea a gramaticilor universale (echivalente deci cu modelul de Turing). Formalismul ATN (Augmented Transition Networks) al lui William Woods de la BBN a fost timp de peste 10 ani suportul standard al majorității sistemelor de prelucrare a limbajului natural. Eu însumi am dezvoltat în anii 1984 și 1985 un mediu de programare lingvistică conținând un editor de gramatici ATN și un compilator ATN. Din punct de vedere formal ATN-ul este echivalent cu o mașină de Turing și tocmai această putere formală prea mare l-a scos din competiția soluțiilor utile în lingvistica computațională. La sfârșitul anilor '80 obiectivul major al lingvisticii (valabil și astăzi) a devenit identificarea unui formalism de putere generativă cât mai mică dar care să acopere cât mai multe din problemele practice posedate de prelucrarea automată a limbajului natural. Așa au revenit în actualitate gramaticile independente de context și s-au dezvoltat abordările lexicalizate. Cele din urmă au fost propuse tocmai pentru a rezolva, în cadrul scheletelor de gramatică independente de context, idiosincraziile limbajului natural cel mai adesea localizate la nivelul lexical. Mai mult, după anii '90, odată cu resurecția interesului față de abordările statistice, gramaticile regulate și automatele finite au căpătat o utilitate foarte largă.

LC are mai multe părți, mai multe orientări, mai multe niveluri de abstracție, care comportă criterii diferite de evaluare.

Este adevărat că actualmente în LC se regăsesc orientări, abordări și motivații diferite. Dar indiferent de sorginte, ele se plasează (cel puțin declarativ) în contextul computațional prin raportarea la un mediu software de prelucrare. Considerând exemplul HPSG, probabil cea mai în vogă teorie lingvistică computațională actuală, atunci când Ivan Sag analizează sau argumentează adecvarea teoriei sale în descrierea formală a unei limbii naturale (așa cum a procedat în recentele sale conferințe la Facultatea de Litere a Universității București și în Aula Academiei Române) el se plasează în sfera lingvisticii teoretice. Atunci când prezintă soluțiile de implementare a unui fragment natural al limbii engleze și discută rezultatele generate de analizorul HPSG dezvoltat de grupul sau de la Universitatea Stanford și modalitățile algoritmice de rezoluție a ambiguităților (așa cum a făcut în prelegerea susținută la sediul RACAI) se plasează în sfera LC.

DT îl asociază pe D. Hays la ideea sa privind falimentul LM și lansarea, drept consecință, a LC.

Afirmația de mai sus conține două lucruri false:

- a) nu am vorbit de falimentul LM ci de insuficiența metodelor de rezoluție în momentul invocat (cred că citatele pe care le-am prezentat conțin argumentele aduse până acum sunt lămuritoare).

b) Eu nu-l pot asocia pe David Hays la o idee pe care nu am exprimat-o.

În textul meu original scriam: „la propunerea lui David Hays, domeniul de cercetare al limbajelor naturale, din perspectiva utilizării acestora în interacțiunea cu calculatoarele electronice, este individualizat sub numele de lingvistică computațională”.

Propunerea lui Hays venea în sprijinul identificării unui nume comun pentru diversele preocupări asupra limbajului din perspectiva implementării de sisteme automate de prelucrare. Traducerea automată, un domeniu care se dezvoltase distinct de celelalte preocupări în domeniul prelucrării automate a limbajului natural, căzuse în disgrație în urma raportului ALPAC (*Languages and machines: computers in translation and linguistics. A report by the Automatic Language Processing Advisory Committee, Division of Behavioral Sciences, National Academy of Sciences, National Research Council. Washington, D.C.: National Academy of Sciences, National Research Council, 1966. (Publication 1416.) 124pp.*). În raportul ALPAC, comandat în 1964 de Academia Națională de Științe, în afara criticilor deosebit de dure la adresa realizărilor și abordărilor de până atunci în domeniul traducerii automate existau și o mulțime de recomandări care se refereau la noi metode de investigație științifică și la abordarea unor obiective mai realiste. Istoria domeniului a reținut (pe nedrept) doar apriga critică a lui Bar-Hillel care, considerată unilateral, a dus la stoparea pentru circa 15 ani a cercetării oficiale în domeniul traducerii automate în SUA și mai apoi în majoritatea țărilor dezvoltate (o incitantă prezentare a a ceea ce a însemnat proiectul ALPAC este „ALPAC: the (in)famous report”, <http://ourworld.compuserve.com/homepages/WJHutchins/Alpac.htm>, și îi aparține lui John Hutchins). Ceva trebuia făcut pentru a conserva câștigurile științifice obținute până atunci și a permite în noul context continuarea cercetărilor anterioare cu scopul declarat al realizării de programe cu obiective realiste. O serie de minți luminate (John Pierce, David Hays, John Carroll) au văzut pericolul ca, asociate cu domeniul traducerii automate, toate celelalte preocupări privind prelucrarea automată a limbajului puteau fi periclitate, și în acest sens în raport s-a inserat un capitol distinct numit „Automatic language processing and computațional linguistics” ce arăta beneficiile aduse de cercetarea în domeniul traducerii automate în domeniile prelucrării automate a limbajului și al lingvisticii computaționale. Printre altele în capitolul respectiv se arată că „... (what is required is) **basic developmental research in computer methods for handling language**, as tools for the linguistic scientist to use as a help to discover and state his generalizations, and ... to state in detail the complex kinds of theories..., **so that the theories can be checked in detail.**” (sublinierea mea, DT). Mai mult președintele comitetului de elaborare a raportului ALPAC, John Pierce, conștient de pericolul interpretării greșite sau al ignorării recomandărilor prezente în anexele raportului (așa cum s-a și întâmplat), a ținut să insereze în raportul final adresat președintelui Academiei Naționale de Științe o secțiune nouă care sublinia idea de a susține lingvistica computațională în mod distinct de traducerea automată („supporting computațional linguistics, as distinct from

automatic language translation”). Dezvoltând ideile din capitolul raportului ALPAC referitor la prelucrarea limbajului natural (concept care și atunci și acum este distinct de cel al traducerii automate) Pierce considera că NSF (National Science Foundation) trebuia să asigure fonduri de cercetare pentru dezvoltarea de modele de limbă de dimensiuni mari „**since small-scale experiments and work with miniature models of language have proved seriously deceptive in the long run, and one can come to grips with real problems only above a certain size of grammar size, dictionary size, and available corpus**”.

Acesta este contextul în care David Hays, activ cercetător la începutul anilor '60 în domeniul traducerii automate (de altfel unul din membrii comitetului care a elaborat raportul ALPAC) a propus individualizarea preocupărilor legate de prelucrarea limbajului natural cu ajutorul calculatorului, dezvoltarea de modele de limbă realiste (nu miniaturi la îndemâna cercetării individuale) și a aplicațiilor „serioase” (în opoziție cu experimentele la scară mică) sub numele de lingvistică computațională.

Denumirile folosite pentru preocupările la interferența limbajelor naturale și lingvistică informaticii și matematicii au variat tot timpul și nu cred că acest aspect merită prea multă atenție. Lingvistică matematică? computațională? inginerească? algebrică? cognitivă? aplicată? cantitativă? teoretică? statistică? probleme matematice ale semioticii? teorie a comunicării? nologia limbajului? limbajul în inteligența artificială? lingvistică inginerească? procesarea limbajului natural? "information storage and retrieval"? lingvistica cibernetică? pe fiecare dintre acestea să scrie și să întălnit-o și propriile mele articole au fost publicate aproape simultan pe fiecare dintre etichetele de mai sus.

Citatul de mai sus mi se pare extrem de relevant pentru discuția de mai sus și definește clar diferența de opinii. Dacă de pildă distincția dintre *medicină umană* și *medicină veterinară* sau (coborând în taxonomie) între cardiologie și stomatologie „nu merită prea multă atenție” atunci domnul Academician are dreptate.

Din punctul meu de vedere însă, este o mare diferență între termenii și denumirile ale studiului limbii amintite mai sus (la care se mai poate adăuga și lingvistica), ele definind câteva domenii distincte definite prin obiective, competențe, metode și modele.

În 1962 s-a înființat în USA "Association of Computational Linguistics".

De fapt în 1962 s-a înființat AMTCL, acronim pentru „Association of Machine Translation and Computațional Linguistics”, primul președinte al acestei organizații fiind Victor Ingve (cel pe care l-am citat mai devreme), iar al doilea fiind David Hays. ACL (Association of Computațional Linguistics) a apărut abia în 1968.

În 1963 Ferenc Kiefer a demarat la Budapesta revista "Computațional Linguistics", care a trăit peste zece ani.

Este adevărat, dar conținutul ei era foarte diferit de al revistei „Mechanical Translation and Computational Linguistics” apărută în 1965 ca revistă oficială a AMTCL. Și tot ca un rezultat al diferențierilor tot mai mari care apăruseră în domeniu, AMTCL își încetează activitatea la începutul anilor '70 fiind înlocuită de „American Journal of Computational Linguistics” care în 1984 devine „Computational Linguistics” (actuala denumire).

Conferința de la Grenoble de "traitement automatique des langues" din 1967 era a treia de acest fel, fiind precedată de o alta, la New York, în 1965 și de una în Anglia, probabil în 1963, organizată de M. Masterman. Între timp, la ruși, numeroase conferințe au avut loc pe tema "avtomatizatskaja obrabotka tekstov" iar "Sprachkunde und Informationsverarbeitung" a fost uneori eticheta folosită de germani

- *s.a.m.d. Nu negăm rolul important pe care l-a avut David G. Hays în dezvoltarea CL, dar acest rol a fost altul decât cel afirmat de DT.*

Nu am să reiau explicația faptului că nu i-am atribuit lui Hays nici un rol demolator, dar trebuie să subliniez faptul că inițiativa lui David Hays, de care am discutat mai devreme, a avut un rol **fundamental** în evoluția CL. Așa cum am arătat mai sus, inițiativa disocierii de traducerea automată, pentru a nu periclita restul preocupărilor privind prelucrarea automată a limbajului a fost o necesitate conjuncturală. În 1965, când la New York a avut loc prima conferință COLING, Hays anticipa desigur efectul de bumerang al raportului la elaborarea căruia participa, și a propus chiar atunci, detașarea oficială prin sintagma „computational linguistics” de domeniul traducerii automate (pe care îl părăsise de altfel și Hays cel ce fusese unul dintre principalii specialiști în traducere automată ai RAND Corporation). Deci nu Hays a creat domeniul lingvisticii computaționale, el este cel ce a „oficial” botezul. Și nu a făcut-o de pe orice poziție ci de pe cea de fost membru al Comisiei Alpac și de președinte al AMTCL.

Emergența LC s-a produs încă din anii "50, sintagma LC a devenit curentă încă de la începutul anilor "60. Șirul de conferințe COLING nu a făcut decât să continue aceasta tradiție. Alții au preferat folosirea sintagmei LM (a se vedea, de exemplu, "Prague Bulletin of Mathematical Linguistics", "Prague Studies of Mathematical Linguistics", revista japoneza "Mathematical Linguistics" (în echivalentul ei japonez) etc. în ceea ce privește însă profilul acestor reviste, nu am constatat o diferență față de cele de CL. Desigur, între timp au început să apară și unele publicații mai specializate, cu referire la părți determinate ale CL (cum ar fi cea relativă la corpusul lingvistic). Etichetele nu au avut importanța și nu știu să se fi desfășurat vreă competiție între ele. Chiar Hays a folosit diverse etichete, de exemplu cea din [3].

Persistența cu care domnul Academician pune semnul egalității între domeniul lingvisticii matematice, în care fără discuție nu a avut sau nu are rival în

România, și cel al lingvisticii computaționale sau tehnologia limbajului este apa foarte curioasă. Nu și dacă observăm următoarele fapte:

- sintagma „lingvistică matematică” este din ce în ce mai puțin utilă (o căutare pe internet a termenilor „mathematical linguistics”, „computational linguistics”, „natural language processing” și „language technology” este foarte instructivă: numărul de documente ce îi re este 4.630, 87.900, 169.000 și respectiv 2.840.000);
- în domeniul strict computațional, la care se referea [1], în România activează de câțiva timp o serie de cercetători importanți (majoritatea dintre ei membri ai Comisiei de Informatizare pentru Limba Română pe care am onoarea să o conduc, și din care de altfel face parte și domnul Academician Marcus);
- domnul Academician Marcus fie nu cunoaște, fie dezavuează rezultatele românești obținute în domeniul **prelucrării cu calculatoare** a limbii române (cel puțin așa poate fi considerată ignorarea compoziției a acestora în lucrările domniei sale); ori poate considera că acestea reprezintă domeniul său de interes.

Dar DT merge mai departe pe ideea sa și afirmă (în completă discordanță cu viziunea lui Hays, de la care se reclamă) că "metodele LM și-au atins limitele" (încă în urmă cu peste 30 de ani!), pentru ca numai două pagini după această afirmație (deci la pagina 135 din [1]) să afirme că e nevoie de "modele formale ale limbii la toate nivelurile ei (fonetică, morfologie, sintaxă, discurs) gramatici formale [...]". Cum vede DT aceste modele formale altfel decât sub forma logico-matematică?

Asupra primei părți a acestei fraze cred că am discutat suficient. Refuz la „contradicția” pe care o semnalează în partea a doua a frazei de mai sus, nu decât să-i recomand domnului Marcus să citească încă de câteva ori articolul respectiv (sau să-l citească integral). Este vorba de **NOI modele formale de limbă** (în opoziție cu cele vechi), resurse lingvistice computaționale adecvate momentului actual. Dintre noile teorii care au apărut și s-au și impus așa putea să amintim teoria optimalității în comunicare dezvoltată de Prince and Smolensky în 1990 și implementări în domeniul fonologiei și morfologiei computaționale și promițătoare rezultate chiar în sintaxă), teoriile sintactice bazate pe unificarea satisfacerea de restricții, precum și o întreagă pleiadă de teorii ale discursului în domeniul prelucrării automate a limbajului natural există standarde, tehnologii specifice, există organizații mondiale specializate, mai toate apărute în ultimii 10-15 ani. Dacă domnul Academician Marcus poate afirma că pentru România în domeniul resurselor lingvistice computaționale s-a făcut (sau a făcut) ceva înainte de anii '80 înseamnă că domnia sa are o imagine complet diferită de cea a tuturor specialiștilor din lume.

Știe oare că multe modele de acest fel există de câteva decenii? Indicații bibliografice asupra lor sunt date parțial în [4], [5], [6], [7] iar pentru cercetările românești în [8], [9]. Desigur, aceste modele sunt inegale ca valoare, au nevoie de continuări, modificări, ameliorări, dar ele nu pot fi ignorate. Fonetica, fonologia, vocabularul, morfologia, sintaxa, semantica lingvistică și lingvistica istorică au beneficiat din plin de metodele matematice, așa cum se poate vedea din impactul deosebit al lucrărilor respective în literatura de specialitate;

Recursul la modelele anilor '60-70 descrise în lucrările menționate ca argument pentru concepte ce au apărut la începutul anilor '90 mă scutește de comentarii. Pe de altă parte, avansul științific în orice domeniu se clădește pe cunoașterea anterioară iar cazurile de „frângere cognitivă”, când salturile științifice neașteptate anterioare sunt rare și ele de regulă definesc revoluțiile în știință. Filația sau influențele în dezvoltarea unui domeniu științific (atunci când ele pot fi depistate cu obiectivitate) constituie preocuparea istoricilor științei. Lucrările tehnice, de regulă se raportează la contemporaneitate, ceea ce în termeni temporali poate însemna, în funcție de dinamica domeniului, câțiva ani, un deceniu, mai multe decenii sau perioade chiar mai mari. De pildă, puține lucrări tehnice în domeniul lingvisticii teoretice, al fonologiei se referă la marele gânditor Panini, considerat de mulți oameni de știință creatorul științei limbii. Lucrarea sa fundamentală *Astaka*, cunoscută și sub numele de „gramatica lui Panini” conține descrieri formale ale regulilor de producție ale limbii sanscrite și o clasificare cu peste 1700 de elemente constitutive ale limbajului. Aceste elemente sunt organizate în clase a căror agregare este descrisă prin intermediul unor reguli ordonate, într-o manieră apropiată de teoriile actuale. El poate fi considerat un precursor al teoriei limbajelor formale și al lingvisticii matematice, dar puține cărți sau lucrări de referință în aceste domenii menționează numele genialului savant ce a trăit cu mai bine de peste 2500 de ani în urmă. În schimb, numele său se regăsește în orice lucrare serioasă de istorie a lingvisticii formale.

Obstinația cu care domnul Academician Marcus încearcă să sugereze că eu aș dezavua metodele matematice, sau rezultatele importante ale lingvisticii românești dovedește că domnia sa complet neinformată în ceea ce mă privește.

DT indică, drept domeniu al LM, numai "aspectul sintactic, gramatical", despre celelalte nu a aflat. Nu a aflat nici ca LM a abordat și aspecte analitice, nu numai pe cele generative.

Fals: „**numai**” este imaginația domnului Academician. Citatul corect este: „în **general** al aspectului lor sintactic, gramatical”.

DT definește "dimensiunea fundamentală" a LC prin "fezabilitatea instanțierii unei descrieri lingvistice cât mai complete, mentenabilitatea acestei instanțieri și, desigur, conformanța cu realitatea uzului limbii". ([1]: 133). Cu un mic efort înțelegem despre ce este vorba.

Desigur că problemele de complexitate, de cost, nu puteau fi în abordate în anii '50 și '60 cu mijloacele cu care ele au început a studiate în a doua jumătate a anilor 70, când instrumente elaborate în informatica matematică deveniseră mult mai perfecționate. Dar acest fapt nu ține, cum crede DT, de alegerea între LM și LC, ci de progresul general realizat în știință. Pentru mă referi la propria noastră experiență, atunci când, în 1990 prezentam la COLING-ul din Suedia gramaticile contextuale aveam cum să mă ocup de aspectul complexității acestor gramatici în maniera în care s-a putut face acest lucru ulterior (a se vedea de exemplu, [10]). Dar acest fapt nu are nici o legătură cu eticheta folosită.

Efortul (chiar mic) este probabil generat de unii termeni de sens nefamiliari domnului Academician. Voi furniza lămuririle necesare mai jos.

Eu mă refer la perioada actuală când invoc ca dimensiune fundamentală *fezabilitatea instanțierii* unei descrieri lingvistice cât mai complete. Instanțierea descrieri lingvistice înseamnă altceva decât complexitatea formală, de altfel și amintesc în secțiunea trunchiată a citatului folosit de domnul Academician Marcus mai sus. Este un termen tehnic care se referă la construcția programelor în baza unui formalism sau teorii lingvistice, a unei gramatici și a dicționarului aferent, care furnizate ca resurse unui program de prelucrare a limbajului permit acestuia să analizeze sau să genereze un text arbitrar. O instanțiere este fezabilă dacă ea se poate realiza în condiții de timp și cost umane rezonabile.

Nu m-am mirat atunci când "Encyclopedia of Microcomputers and Technology", "Encyclopedia of Computer Science and Technology" mi-au solicitat o contribuție cu tema "Semiotics and Formal Artificial Languages" se vedea [11]) și nici când "Handbook of Formal Languages" solicitat un capitol privind "Contextual Grammars and Natural Languages"[12] iar o lucrare preponderent teoretică a fost inclusă în "Computational Linguistics in the Netherlands 2000"[13].

Nu văd rostul acestor lămuriri. Toată lumea îl știe, îl recunoaște și dintre cercetătorii adevărați nu-l contestă pe omul de știință Marcus, reprezentant român al lingvisticii matematice, creatorul acestei școli în România. articolul [1] nu m-am referit nici direct nici indirect la domnia sa. Faptul că am evocat criticile pe care le-am comentat anterior la adresa **metodelor matematice** ale începutului deceniului șapte nu are nici o legătură cu domnia sa (încă o dată, deosebite) ale domnului profesor. Însă probabil că identificarea LM mondială, domnia sa a considerat critica asupra **metodelor LM** din România un atac la persoana sa, adevărat act de blasfemie.

În anii din urmă, domnul Academician încearcă să transfere în România noile tendințe și tehnologii ale limbajului, ignorând o realitate existentă

portofoliul de rezultate pe care le-a obținut anterior creditându-le ca surse primare a tot ceea ce se întâmplă azi în tehnologia limbajului în România (și nu numai). Și cine nu este de acord cu acest lucru (parafrazându-l pe domnul Marcus) trebuie demonizat. Textul pe care îl comentez ca și acțiunile recente declanșate de domnul Academician Marcus, pretinse a fi iscate de conținutul articolului [1], nu fac decât să-mi întărească această impresie. Eu nu am nimic de împărțit cu domnul Academician.

Nu m-am mirat nici când am văzut că o revistă cu titlul "Linguistics and Philosophy" publică articole excelente de LC. Interferențele merg în toate direcțiile și ele caracterizează cultura contemporană, în acest orizont trebuie să ne plasăm, cred, atunci când ne referim la disciplinele cognitive care se dezvoltă sub ochii noștri și își pun amprenta pe modul nostru de gândire și de comportare. Un tratat ca "Mathematical Methods in Linguistics" [14] include multe fapte de LC, deși în titlul său nu figurează epitetul "computațional". O revistă ca "Theoretical Linguistics" (1970-2000), publicată de Walter de Gruyter (Berlin-New York) a inclus multe articole vizând aspecte matematice și/sau computaționale, deși numele revistei nu indică acest lucru. Chiar o revista mai tradițională, ca "Linguistics" a inclus de multe ori articole de LM și nici "Foundations of Language" nu a procedat altfel. Multe fapte de LM și de LC se plasează în mod natural în orizontul semioticii computaționale.

Faptul că tratatul amintit nu încorporează în titlu atributul computațional nu mă surprinde, pentru că ar fi creat o confuzie pe care autorii au evitat-o deliberat. Cartea respectivă nu este o carte de lingvistică computațională, conținutul ei tratează exact ce anunță în titlu: metode matematice folosite în studiul lingvistic. Lingvistica teoretică, puternic formalizată în ultimele decenii apelează inevitabil (ca de altfel marea majoritate a domeniilor științifice) la metode și modele matematice.

Era internetului impune desigur o problemă nouă, față de care abordările anterioare se pot dovedi insuficiente.

Exact aceasta este esența celor 3 paragrafe din [1] incriminate și combătute pe larg de domnul Academician Marcus: insuficiența abordărilor anterioare. Conștientizarea acestei insuficiențe însă a precedat cu câțiva ani apariția internetului.

Salutăm inițiativa noii generații de cercetători de a se dedica noilor probleme.'

Nu putem ignora tonul paternalist privind noua generație de cercetători care se dedică problemelor ridicate de internet în prelucrarea automată a limbajului natural. INTERNET-ul este o revoluție! Și implicațiile sale sunt atât de mari încât asigurarea accesului universal la Internet a devenit o problemă fundamentală chiar și pentru o organizație de calibrul UNESCO. Am avut onoarea să fac parte din Comisia de Experți creată de Secretarul General al UNESCO (comisie de cel mai

înalt nivel) pentru elaborarea documentului Recommendation on Multilingualism and Universal Access to Cyberspace. Sunt al doilea expert român (după Ambassador Dan Hăulică, Membru Corespondent al Academiei) care a făcut parte dintr-o comisie de experți UNESCO de acest nivel.

Ignorarea în cercetarea privind prelucrarea automată a limbajului natural a fenomenului INTERNET este de neconceput. Societatea cunoașterii are ca una din premisele sale fundamentale accesul universal, neîngrădit de bariere lingvistice la cunoșterea stocată în internet. Alte comentarii sunt de prisos.

Dar trecerea de la ieri la azi și de la azi la mâine nu poate fi decât una care ține seama în mod critic de experiența acumulată.

Nimeni nu neagă acest lucru, și faptul că l-am rugat insistent pe domnul Academician să facă parte din Comisia de Informatizare pentru Limba Română cred că arată buna mea credință și speranța pe care o nutream (și care încă supraviețuiește) că experiența domniei sale va fi pusă în slujba obiectivelor pe care nici eu nici domnul Marcus nu le putem atinge singuri. În același spirit, i-am propus domnului Academician Marcus să scriem împreună o antologie a cercetărilor românești în domeniul lingvisticii formale și computaționale, de la începuturile care le evocă domnia sa și pînă în zilele noastre. Din păcate propunerea a rămas fără răspuns.

Din tot ceea ce am prezentat mai sus rezultă clar ca LM și LC au evoluat și s-au dezvoltat fost mereu împreună și că, în general, etichetele nu au contat prea mult. Unii au mers chiar mai departe; astfel, în capitolul 4 al cărții "Mathematical and Computational Linguistics", din [15], se afirmă: "Theoretically, pur și simplu (p.86): "Mathematical linguistics has also been called theoretical linguistics and even computational linguistics". Iar mai departe, în același loc: "Computational Linguistics originated around 1950 with the initiation of research on automatic translation" (se referă la o carte editată de D.G.Hays [3] și la o alta avându-l ca autor pe acesta [16]).

Nu văd în pasajul pe care l-am citat mai sus nici un argument împotriva ceea ce am susținut în [1] și în cele prezentate aici. Notez în treacăt adesea „even” cu o valoare discursivă în completă consonanță cu considerentele istorice pe care le-am invocat ale evoluției științifice și tehnologice în domeniul prelucrării limbajului natural.

În România, minți luminate ale anilor '60, ca profesorii Al. Rosetti, Grigore Moisil și Tudor Vianu, au înțeles schimbările care se profilau și au sprijinit proiectul înființării unei secțiuni de "lingvistică aplicată" la Facultatea de Limba și Literatură Română a Universității din București, dar s-au găsit alții care să-i torpileze.

Așa este, și mă bucură elogiul adus acestor corifei ai științei române. Poate și pentru că alături de câțiva reprezentanți importanți ai lingvisticii române actuale care au înțeles tendințele și imperatiile momentului (Prof. Dan M. Ștefănescu, decanul Facultății de Litere, Prof. Alexandra Cornilescu, Conf. Emil Ionescu

participat la reluarea acestei lucrări. Programul de Masterat în Lingvistică Formală și Computațională de la Facultatea de Litere a Universității din București, funcționează de mai bine de 2 ani și nutresc speranța că Ministerul Educației și Cercetării va aproba demersurile noastre privind chiar înființarea unui departament cu acest profil.

În același sens, am participat alături de profesorul Cristea (având fără discuție și sprijinul altor minți luminate ale Universității A.I.Cuza din Iași) la lansarea în 2001 a Masterat-ului în Lingvistică Computațională al Facultății de Informatică. Nu este ușor să pendulezi între Iași și București, dar și domnul profesor Cristea, și doamna profesor Cornilescu și eu o facem pentru ca cele două programe „surori” de mașter să-și împlinească menirea de a pregăti câți mai mulți specialiști în folosul programelor de informatizare pentru limba română.

La Academia Română a funcționat mulți ani "Comisia de Lingvistică Matematică" iar revista "Cahiers de Linguistique Theorique et Appliquee", înființată în 1962, a fost multă vreme expresia colaborării lingvisticii cu matematica și cu informatica. În ciuda forțelor adverse, s-a reușit în acei ani atragerea unor studenți străluciți ai unor facultăți umaniste la cercetarea limbii și literaturii cu mijloacele LM și LC.

Comisia de Informatizare pentru Limba Română de la Academia Română, înființată în anul 2001, încearcă, ținând cont de realitățile și prioritățile actuale, să armonizeze eforturile celor ce lucrează în domeniul limbii române și care cred în perspectiva înrolării ei în cadrul limbilor importante ale societății cunoașterii. Eu am convingerea că voi putea spune peste timp același lucru: „In ciuda forțelor adverse, s-a reușit în acei ani atragerea unor studenți străluciți ai unor facultăți umaniste la cercetarea limbii și literaturii cu mijloacele” tehnologiei limbajului.

Pentru a da numai două exemple de actuali profesori universitari care au susținut teze de doctorat de acest tip, voi menționa pe Pia Brinzeu, de la Catedra de Engleză a Universității din Timișoara și pe Mihai Dinu, de la Facultatea de Litere a Universității din București. Tot în acea perioadă și-a susținut teza de doctorat Sorin Cristian Niță, pe o tema de critică textuală automată privind înlănțuirea (filiația) diferitelor variante ale "Istoriei Tării Românești" (Șerban Cantacuzino).

Exemple de profesori și cercetători români valoroși, cu contribuții substanțiale în domeniul limbii române se pot da foarte multe. Mulți dintre ei sunt în străinătate și fac o bună propagandă științei românești. Mi-e cunoscută cartea cu adevărat remarcabilă a domnului profesor Mihai Dinu „Personalitatea limbii române”, de altfel premiată de Academia Română. Această lucrare este o solidă cercetare de lingvistică computațională în spiritul actual tocmai pentru că a parcurs acea cale dificilă a instanțierii lingvistice (în cazul său la nivelul componentului lexical).

lată însă că, în pofida realităților puse în evidență mai sus, în ([1]: 134) se scrie: "În România, cercetările în domeniul LC și al prelucrării limbajului natural, precum și primele rezultate practice au apărut la începutul anilor "80 [3, 4, 5, 6]".

La ce trimit numerele indicate în paranteze ? La o bibliografie de 24 de titluri în care aproape toate (dar toate cele indicate între paranteze) încep cu DT (ignorandu-se regula generală în lumea științifică, a așezării numelor autorilor aceluiași articol în ordine alfabetică; dar nu acest fapt este cel care ne interesează în momentul de față).

Înainte de a comenta acest pasaj și pe cel următor, nu pot să trec observația absurdă și falsă pusă între parantezele ce trădează totuși o ezităre a omului de știință în fața unei răutăți gratuite. Nu există nici o regulă generală de genul celei afirmate. Ordonarea alfabetică este o convenție autorii cu contribuții egale în redactarea unei lucrări. Am deschis la întâmplare două volume de specialitate, conținând contribuții (S. Armstrong et al. „Natural Language Processing Using Very Large Corpora, Kluwer, 1999 Strzalkovski (ed) „Natural Language Information Retrieval”, Kluwer, 1999). Din 19 lucrări cu mai mulți autori, doar trei urmăresc (probabil din întâmplare) „regula generală în lumea științifică” pe care o invocă domnul Academician care probabil a impus-o și o impune tuturor celor alături de care publică, indiferent de contribuția fiecăruia.

Să observăm că încă în 1978, în articolul "Mathematical and Computational Linguistics" [9] de prezentare a activității din România în domeniul LM și LC se face referire la peste 400 de articole publicate de 130 de autori români și sunt menționați peste 300 de autori străini (unii dintre ei, nume de vază ale LM și LC din acea perioadă) care au citat și continuat cercetările românești. Să mai adăugăm că numeroși lingviști români dintre cei mai importanți au citat și folosit rezultatele școlii românești de LM și LC. Iată ca vine acum DT și face (deliberat sau nu) din tot acest efort un teren viran care-l aștepta pe DT să tragă primele jaloane. Nu e cam mult?

Deși am repetat de nenumărate ori până în acest moment, o mai dată, precizând că discuția din [1] se referea la **resurse lingvistice computaționale și programe software de dialog în limbaj natural** (în română). Acestea erau rezultatele practice pe care le menționam în comentat cu gratuită aciditate. Poate să-mi menționeze domnul Academician sistem de dialog în limba română implementat înaintea sistemelor pe care realizat eu și colaboratorii mei? Iată câteva repere:

- Sistemul QA (1980) un sistem inferențial de întrebare răspuns în limba română, susținut de un demonstrator original de teoreme în calculul predicatelor de ordin 1;

- SDLR (1981) un sistem de dialog în limba română ce a extins capabilitățile lui QA cu operatorii lingvistici ai logicii fuzzy;
- IUREȘ (1983) sistem de generare automată a sistemelor de întrebare-răspuns, independent de limbă, pe care l-am realizat împreună cu Dan Cristea, acum decanul facultății de informatică a Universității Cuza. Sistemul IUREȘ a fost omologat internațional în 1988 și a constituit primul produs de inteligență artificială exportat (în același an). Sistemele IUREȘ și SDLR sunt referite printre altele în enciclopedia de lingvistică computațională. Mai important este faptul că sistemele IUREȘ și SDLR sunt amplu descrise în prestigioasa antologie "The Survey of the Current Status Research and Future Trends in Machine Translation and Natural Language Processing" realizat în 1992 de JEIDA (Japan Electronic Industry Development Association), fiind de altfel singurele sisteme de dialog în limbaj natural din întreaga zonă fost comunistă incluse în această carte.

Acestea erau referințele incriminate de domnul Academician și dacă domnia sa poate să-mi indice un singur sistem de prelucrare a limbajului natural realizat în România înaintea celor pe care le-am citat, eu am greșit. Dar mă îndoiesc. Nu cunosc conținutul articolului menționat (pe care i l-am solicitat de altfel domnului Academician, fără a-l primi însă), astfel încât nu pot afirma nimic despre cei 130 de autori români ce au realizat (conform afirmației domnului Marcus) lucrări de lingvistică computațională. Ce pot însă să afirm este că am citit multe din lucrările de lingvistică teoretică contemporană ale marilor noștri lingviști și ele au fost extrem de relevante ca material factual în cercetările mele. Dar lucrările pe care le-am citit (și citat) eu, nu erau din domeniul lingvisticii computaționale. Lucrările domnului Marcus (în special cele din domeniul *limbajelor formale*) apăreau destul de frecvent între referințele bibliografice ale lucrărilor mele de la începutul anilor '80. Eram la început de drum, sursele documentare erau puține și demersul era natural. Pe atunci, Chomsky era din nou foarte în vogă, noua sa teorie *Government and Binding* impulsiona o serie de cercetări în domeniul formalizării gramaticii universale. Tentația computațională față de această teorie a fost enormă, și chiar dacă actualmente nu există nici o gramatică computațională efectivă a GB, idei fundamentale din GB se regăsesc în formalisme lingvistice computaționale moderne (cum ar fi HPSG).

Să fim bine înțeleși. Nu noi avem nevoie de încă o citare pe lângă miile de citări deja acumulate, ci noile generații de studenți și de cercetători au dreptul la o informare corectă asupra dezvoltării LM și LC în general și, în particular, asupra LM și LC în România. DT a mai publicat, în urmă cu câțiva ani, un articol în care se schița o privire istorică asupra LC în România, cu câteva citări la întâmplare, care trădau necunoașterea situației reale.

Cu rezerve față de prima parte a paragrafului, mă opresc la grija domnului Academician pentru dreptul noilor generații de studenți și de cercetători asupră „informării corecte” asupra istoriei LM și LC. Personal, cred că mult mai important pentru ei este să știe prezentul și tendințele viitoare ale domeniului. Astfel, cunoștințele le pot asigura un loc de muncă, o direcție de specializare, o carieră viitoare. Noile generații de studenți și de cercetători sunt utilizatori pasionați ai Internetului. Acest uriaș ocean informațional le asigură un imens volum de cunoștințe, începând cu cursuri on-line (obligatorii pentru profesori la mai multe universități importante ale lumii), volume ale conferințelor sau articole extrem de utile, recente și mai puțin recente, cărți electronice. Chiar și relevante lucrări de istorie asupra diverselor domenii științifice. Sistemele moderne de regăsire de documentară le asigură și o ierarhizare a acestor surse de informare în raport cu relevanța și cu interesul manifestat de alți cititori. Listele de discuții sau arhive de întrebări frecvente (FAQ) le pot oferi răspunsuri avizate și obiective la întrebările ce-i preocupă. În anexă este furnizat un exemplu.

În ultima parte a citatului de mai sus, domnul Academician Marcus a menționat în discuție o lucrare a mea din 1996 și care arată că frustrările domniei sale sunt mai vechi. Articolul de care amintește domnul Academician mai sus, are ca titlu „Resurse lingvistice computaționale: trecut, prezent și viitor” și a apărut în volumul „Limbaj și Tehnologie”, Ed. Academiei, 1996. Cei interesați, pot găsi articolul respectiv în pagina oficială a RACAI (<http://www.racai.ro> secțiunea publicații). Iar cele „câteva citări la întâmplare, care trădau necunoașterea situației reale” apar în capitolul 2. „Cercetări și realizări românești în domeniul prelucrării automate a limbajului natural”. Cred că titlul volumului, al articolului și capitolului sunt lămuritoare pentru ceea ce discutăm acolo, dar probabil fraza, care trimitea la un volum editat de domnul Marcus, „abordările statistice, revenite în actualitate, au avut o tradiție strălucită (în România, adăugarea mea DT)” este prea scurtă și insuficient de laudativă.

Mai este un aspect care cere o precizare. În conformitate cu cerințele specifice ale volumului în care apare articolul [1], DT face numeroase referiri la acte și documente ale unor organisme europene și internaționale, cum este și firesc, pentru a nu mai vorbi de aspectul financiar al colaborării cu organismele respective. Această situație a existat de la începutul LM și LC (chiar dacă nu au avut amploarea de azi), datorită faptului că LM și LC au apărut și s-au dezvoltat în urma urmării a unor comandamente sociale, privind precaritatea mijloacelor de prelucrare a informației. Îmi amintesc de faimoasele Rapoarte CETIS care veneau de la EURATOM, Bruxelles, pe teme legate de analiza și prelucrarea automată a limbajului, traducerea automată și documentare automată. În USA, diferite corporații (cum ar fi RAND Corporation, Santa Monica, Calif.) finanțau cercetări similare. O întâlnire semnificativă a fost aceea din 1962, organizată

de "NATO Advanced Summer Institute", la Veneția, Italia, privind traducerea automată. De numele acestui Institut este legat un document care a marcat evoluția cercetărilor de traducere automată: seria de expuneri prezentate de Y. Bar-Hillel [17]. În legătură cu aceste activități dirijate și finanțate de diferite organisme europene și internaționale, trebuie să observăm că cei implicați au avut înțelepciunea și priceperea necesare pentru a nu reduce proiectele respective la dimensiunea lor exclusiv utilitară, ci de a o subordona pe aceasta unei perspective mai ample, care lua în considerare orizontul științific real al problemelor. Pentru a da un prim exemplu, mă voi referi la faptul că mai multe rapoarte CETIS au pus în discuție un concept care, născut din experimentele de traducere automată, avea să se dovedească de o deosebită semnificație pentru teoria sintactică în toată generalitatea sa; este vorba de conceptul de proiectivitate sintactică, cu consecințe bogate în studiul structurilor arborescente și al gramaticilor de dependență. Azi putem spune că și sintaxa limbajului natural și teoria matematică a grafurilor au profitat esențial de conceptul respectiv (folosit până și de Rene Thom, în probleme de morfogeneză [17]). Această expansiune a unui concept sau rezultat dincolo de motivația sa inițială este testul cel mai convingător al interesului său. Un al doilea exemplu se referă la titlul provocator folosit de Bar-Hillel pentru expunerile sale: "Patru conferințe despre lingvistica algebrică și traducerea automată".

Simpla alăturare a celor două sintagme, una foarte teoretică, cealaltă aparent tehnologică, avea menirea să-i avertizeze pe cei care presau să se obțină cât mai repede rezultate practice asupra faptului că proiectele de traducere automată nu se pot finaliza de azi pe mâine, ci au nevoie de un lung itinerar lingvistic, matematic și computațional. Acum știm că acest itinerar continuă și azi, cu tatonări și reveniri, și, chiar dacă nu a dus încă la rezultatele visate, a impulsinat în mod esențial cercetările de AI, cu consecințe benefice pentru aspectele logice și semantice ale limbajului natural, întrebarea pe care ne-o punem, dar o lăsăm deocamdată fără răspuns, deoarece nu suntem pregătiți pentru a-l da, este următoarea: Nu cumva aspectele pe care le-am criticat mai sus sunt consecința unui fenomen mai general, acela al unui orizont insuficient de cuprinzător, al unei prea mari dependențe de factori utilitari imediați? Știința a oscilat mereu între cognitiv și utilitar, dar istoria arată că funcția utilitară s-a manifestat în toată profunzimea ei atunci când ea a fost fructul unei evoluții firești a funcției cognitive, evoluție care poate fi de doi ani, de 20 de ani, de 200 sau de 2000 de ani. Cu un ochi îndreptat spre comisiile europene, suntem obligați totuși să ținem treaz și celălalt ochi, îndreptat spre ceea ce se

întâmplă pe scena cercetării științifice vii, așa cum apare ea în revistele de specialitate și la întâlnirile științifice de profil.

Remarcile de mai sus îmi sugerează celebra fabulă cu strugurii ce Cercetarea instituționalizată (în opoziție cu cea „de dragul artei”) are m întotdeauna justificabile. Organismele de finanțare a cercetării, naționale și internaționale, nu fac desigur acte de caritate. Obținerea unei finanțări pentru un proiect de cercetare nu este la îndemâna oricui și el implică nu numai abordarea unei probleme importante, dar și credibilitatea grupului de cercetare. Evaluarea și propunerilor de proiecte se face de către experți recunoscuți în domeniul respectiv și angajați și plătiți de agențiile de finanțare a cercetării. În condițiile unei concurențe internaționale acerbe pentru fondurile (din păcate prea mici) destinate cercetării, luarea în derâdere, invocând caracterul utilitar, cercetărilor ce obțin concurențial finanțarea arată o desprindere de realitate. În luna martie a.c. am participat la evaluarea propunerilor de proiecte europene din cadrul Programului Cadru (apelul 8), și în calitate de raportor al direcției „11.1.1 • Exploratory Research and Risk/Long Term Research”, pot să afirm că propunerile de proiecte pe care eu am văzut erau foarte departe de a avea caracter utilitar. Domnul Academician Marcus lasă fără răspuns o întrebare cu răspuns sugerat, ridicând o problemă discutată cu ceva timp în urmă, anume a tipului de cunoaștere contemporană enciclopedică (și inerent generalistă) sau specializată. Cel puțin în domeniul tehnologic, viteza fără precedent a apariției de cunoștințe noi face imposibilă cunoașterea enciclopedică și în același timp expertă pe toată lărgimea spațiului cunoașterii actuale chiar și într-un domeniu aparent îngust. Tehnologia limbajului este actualmente termenul ce subsuma toate preocupările legate de prelucrarea automată a limbajului natural. Cred că acest lucru spune totul!

3. în loc de concluzii

Ajungând în acest punct al răspunsului meu la atacul domnului Academician Marcus mărturisesc că mă încercă un apăsător sentiment de deșertăciunii. Nu am dorit această polemică și în nici un caz în acest context. Considerând că ea este nepotrivită față de obiectivele urmărite de proiectul „Soluții și strategii în România”, în calitate mea de director de proiect și coautor al volumului de față, am discutat cu membrii comitetului director al proiectului oportunitatea publicării polemicii domnului Academician Marcus (și imediat al răspunsului meu) în volumul destinat unor probleme tehnice. Părerea noastră unanimă că nu este cazul să amestecăm obiectivele proiectului cu discuția de fond. Dar transmițând domnului Academician această opinie și făcându-i propunerea a găzdui această polemică pe internet (în pagina oficială a RACAI) domnia sa a simțit cenzurat, insultat și îndreptățit să facă o serie de afirmații pe care m să le comentez. Decizia de includere a acestei secțiuni în volumul de față a fost o fără plăcere pentru că pe de o parte, în ciuda părerii domnului Academician

Marcus (*Articolul meu se încadrează perfect în obiectivul pe care pretindeți că-l urmăriți și în acest spirit a fost conceput Realizați gravitatea deciziei Dv?* - de a nu-l include în volum, precizarea mea, D.T.) continui să cred că nici articolul domniei sale nici al meu nu își aveau rostul aici. Pe de altă parte, nu pot decât să deplâng supărarea pe care i-am provocat-o fără voie domnului Marcus și risipa de energie pe care o depune într-o problemă care din punctul meu de vedere nu există. Drept care, sperând că includerea articolului ce se încadrează perfect în obiectivul...îi va da domnului Academician satisfacția pe care și-a dorit-o, las cititorii să aprecieze cât de grav ar fi fost pentru obiectivul tehnologiei limbii române în contextul „Societatea Informațională - Societatea Cunoașterii: Soluții și strategii în România” ca cele două articole să nu fi apărut aici.

Referințe bibliografice (secțiune din lucrarea domnului Academician Marcus):

- [1] D. Tufiș. *Promovarea limbii române în SI-SC*. în *Societatea Informațională - Societatea cunoașterii* (coord. F. Gh. Filip). Ed. Expert, București, 2001, 131-142.
- [2] D. G. Hays. *The field and scope of computational linguistics*. Papers in Computational Linguistics (eds. F. Papp, G. Szepe). Proceedings of the Third International Meeting of Computational Linguistics, held in Debrecen, Hungary, 1971. Akademiai Kiado, Budapest, 1976, 21-26.
- [3] D. G. Hays (ed.). *Readings in Automatic Language Processing*, American Elsevier, New York, 1967.
- [4] S. Marcus. *Mathematical Linguistics in Europe. Current Trends in Linguistics* (Th. A. Sebeok, ed.), vol.9, Mouton, The Hague, 1972, 646-687.
- [5] S. Marcus. *Mathematique et Linguistique*. în *Mathematique, Informatique et Sciences Humaines*, Paris, 26, 1988, 103, 7-21.
- [6] S. Marcus. *The status of research in the field of analytical algebraic models of language*. în *Current Issues in Mathematical Linguistics* (C. Martin-Vide, ed.). Elsevier-North Holland, Amsterdam, 1994, 3-21.
- [7] S. Marcus. *Lingvistica matematică azi*. în *Matematica în lumea de azi și de mâine* (C. Iacob, coord.), Editura Academiei, București, 1985, 182-186.
- [8] S. Marcus. *Recent Romanian investigations in the field of mathematical and computational linguistics*. Avtomatetskaja Obrabotka Tekstov, Matern. Fyz. Fakulta, KL Praha, 1973, 15-42.
- [9] S. Marcus. *Mathematical and computational linguistics*. în *Current Trends in Romanian Linguistics* (A. Rosetti, S. Golopentia Eretescu, eds.). Revue Roumaine de Linguistique 23, 1978, 1-4, 559-588.
- [10] S. Marcus, C. Martin-Vide, G. Paun. *Contextual grammars as generative models of natural languages*. Computational Linguistics 24, 1998, 2, 245-274.

- [11] S. Marcus. *Semiotics and formal artificial languages*. în *Encyclopedia of Computer Science and Technology* (A. Kent, J.C.Williams, eds.) 29, Marcel Dekker, New York, 1994, 393-405; also in *Encyclopedia of Microcomputers* (A. Kent, J.C.Williams, eds.) 15, 1995, 299-312.
- [12] S. Marcus. *Contextual grammars and natural languages*. *Handbook of Formal Languages* (G. Rozenberg, A. Salomaa, eds.), 2, Springer, Berlin, New York, 1997, 215-235.
- [13] S. Marcus, G. Martin-Vide, G. Paun. *A new-old class of linguistically motivated regulated grammars*. *Computational Linguistics in the Netherlands 2000* (W. Daelemans et al., eds.), Selected Papers from the Eleventh Meeting, Ed. Rodopi, Amsterdam, New York, 2001, 111-125.
- [14] B. H. Partee, A. Ter Meulen, R. Wall. *Mathematical Methods in Linguistics*. Kluwer, Dordrecht, 1990.
- [15] E. F. Beckenbach, Ch. B. Tompkins (eds.). *Concepts of Communication: Interpersonal, Intrapersonal and Mathematical*. John Wiley and Sons, New York, 1976.
- [16] D. G. Hays. *Introduction to Computational Linguistics*. American Elsevier, New York, 1967.
- [17] R. Thom. *Stabilité Structurelle et Morphogenèse*. John Benjamins, New York, 1970.
- [18] Y. Bar-Hillel. *Four Lectures on Algebraic Linguistics and Machine Translation*. revised version of a series of lectures given in July 1962, before a group of students at the Advanced Summer Institute, Venezia, Italy.

ANEXAI: Exemple de căutare într-o arhivă de întrebări frecvente (Usenet FAQ)

The screenshot shows a search interface with a search bar containing the query "mathematical linguistics". Below the search bar, there are navigation buttons like "Search", "Favorites", "History", "Mail", "Bin", "Discuss", "Go", "Links", "Citations", "Unlink", "Free Home", and "Show Media". The search results section is titled "The Usenet FAQ Archives" and "Results for query 'mathematical linguistics*'. It shows a message from "Sony" stating "NO Matches were found for query 'mathematical linguistics'". Below this, there is a note: "Since M&M was found... you might want to be a bit less specific in your search phrase. Just a thought...". There are also links to "References" and "FAQ Search".

The screenshot shows a search interface with a search bar containing the query "computational linguistics". Below the search bar, there are navigation buttons like "Search", "Favorites", "History", "Mail", "Bin", "Discuss", "Go", "Links", "Citations", "Unlink", "Free Home", and "Show Media". The search results section is titled "The Usenet FAQ Archives" and "Results for query 'computational linguistics'". It shows a list of search results, including a link to "Computational Linguistics: leu.edu" and "The Association for Computational Linguistics homepage".

ANEXA 2: Definiții

What is Mathematical Linguistics?

MATHEMATICAL LINGUISTICS is the study of mathematical structures and methods that are of importance to linguistics. As in other branches of mathematics, the influence of the empirical subject matter is somewhat limited. Theorems are often proved more for their inherent mathematical value than for their applicability.

Both in phonology/morphology and in syntax/semantics the choice of linguistic formalism is to some extent influenced by considerations that go beyond the primary issue of descriptive adequacy. One important issue is Recognition Complexity. This concerns the complexity of the decision problem for membership in a language: it is assumed that a grammatical theory should have the property of guaranteeing that there is some reasonably rapid (polynomial in the length of the input) computation that will answer the question of whether a given sequence of words is a grammatical expression according to a given grammar. Human beings certainly do much more than this when they listen to an utterance and figure out the meaning of what was said, so a grammatical theory that cannot guarantee reasonably rapid confirmation of well-formedness is probably not psycholinguistically realistic. Another one is Learnability, which concerns the sorts of mathematically definable procedures could in principle correctly generate grammars for languages.

(Geoffrey K. Pullum and Andras

What is Computational Linguistics?

Simply put, COMPUTATIONAL LINGUISTICS is the scientific study of natural language from a computational perspective. Computational linguists are interested in providing computational models of various kinds of linguistic phenomena. These models may be "knowledge-based" ("hand-crafted") or "data-driven" ("statistical" or "empirical"). Work in computational linguistics is in some cases motivated by a scientific perspective in that one is trying to provide a computational explanation of a particular linguistic or psycholinguistic phenomenon; and in other cases the motivation may be more purely technological in that one wants to provide a working component of a speech or natural language system. Indeed, the work of computational linguists is incorporated into many working systems today, including speech recognition systems, text-to-speech synthesizers, automated customer response systems, web search engines, text editors, language instruction materials, to name just a few.

(Copyright © 2000, The Association for Computational Linguistics)

LIMBA ROMANA
în
Societatea Informațională
Societatea Cunoașterii

Acest volum este dedicat Academicianului Mihai Drăgănescu, Profesorul și mentorul unei întregi generații de specialiști în știința și tehnologia informației în general și al problemelor societății informaționale și a cunoașterii în special. Marea majoritate a contribuțiilor din acest volum aparțin unor experți ce fac parte din Comisia de Informatizare a Limbii Române, comisie a Academiei Române la a cărei naștere un rol esențial l-a avut Profesorul Drăgănescu, președintele Secției de Știința și Tehnologia Informației. Savantul Mihai Drăgănescu are numeroase contribuții în știința contemporană, binecunoscute atât în țară cât și în străinătate. Pentru cine îl cunoaște pare incredibilă puterea sa de muncă, debordanta creativitate și neostoita căutare a noului. Profesorul Drăgănescu este indiscutabil port-drapelul conceptului de societate informațională-societate a cunoașterii în România. În lucrările sale din urmă cu peste 25-30 de ani se regăsesc cu claritate multe concepte foarte actuale în zilele noastre, previziuni curajoase atunci, acum realități cotidiene. În lucrările domniei sale din ultima vreme, apare un nou concept ce avem convingerea că se va impune: Societatea Conștiinței, o treaptă superioară a societății cunoașterii. Nu este de mirare deci că în contextul societății informaționale și a cunoașterii profesorul Drăgănescu a susținut cu consecvență și a afirmat cu claritate rolul Inteligenței Artificiale în devenirea noilor societăți ale cunoașterii. Între domeniile Inteligenței Artificiale un loc de frunte în promovarea principiilor societății cunoașterii îi revine Tehnologiei Limbajului Natural. Profesorul Drăgănescu a fost unul dintre puținii oameni de știință români care au înțeles și au sprijin total aceste direcții. Cu aproape douăzeci de ani în urmă (1983), Profesorul Drăgănescu edita (împreună cu Adrian Davidoviciu și Ioan Georgescu) volumul "Inteligența Artificială și Robotica" pentru ca trei ani mai târziu (împreună cu Corneliu Burileanu) să editeze un alt volum de referință "Analiza și sinteza semnalului vocal". Astăzi, cercetările mondiale în domeniul tehnologiilor lingvistice au atins un nivel de maturitate ce permit sinergizarea eforturilor lingviștilor, informaticienilor, matematicienilor și a altor specialiști din sectorul academic sau industrial, să abordeze proiecte mari, interdisciplinare având ca obiectiv prelucrarea automată, în mediile de comunicare electronică, a din ce în ce mai multe limbi naturale. Printre acestea, limba română își face loc încet dar sigur. Volumul de față este o mărturie în acest sens. În același timp, volumul se constituie într-o nouă confirmare a realităților pe care Profesorul Mihai Drăgănescu le prefigura cu mulți ani în urmă.