

Lucrările atelierului  
*Resurse lingvistice și instrumente pentru*  
*Prelucrarea limbii române*  
Iași, 19-21 noiembrie 2008

Volum apărut cu sprijinul Ministerului Educației și Cercetării,  
prin Consiliul Național al Cercetării Științifice din Învățământul  
Superior (CNCSIS)

ISSN 1843-911X

Lucrările atelierului  
*Resurse lingvistice și instrumente pentru  
Prelucrarea limbii române*  
Iași, 19-21 noiembrie 2008

Editori:  
Diana Maria Trandabăț  
Dan Cristea  
Dan Tufiș

Organizatori:  
Facultatea de Informatică,  
Universitatea „Alexandru Ioan Cuza” Iași

Institutul de Cercetări pentru Inteligență Artificială  
Academia Română, București

Institutul de Informatică Teoretică  
Academia Română, Filiala Iași

Editura Universității “Alexandru Ioan Cuza” Iași

**COMITETUL DE PROGRAM:**

**Corneliu Burileanu**, Facultatea de Electronică, Universitatea Politehnica București și Institutul de Cercetări în Inteligență Artificială, Academia Română, București

**Constantin Ciubotaru**, Institutul de Matematică și Informatică, Academia de Științe a Moldovei, Chișinău, R. Moldova

**Svetlana Cojocaru**, Institutul de Matematică și Informatică, Academia de Științe a Moldovei, Chișinău, R. Moldova

**Dan Cristea**, Facultatea de Informatică, Universitatea "Al. I. Cuza" și Institutul de Informatică Teoretică, Academia Română, Iași

**Nicolae Curteanu**, Institutul de Informatică Teoretică, Academia Română, Iași

**Cristina Florescu**, Institutul de Filologie Română "Al. Philippide", Academia Română, Iași

**Corina Forăscu**, Facultatea de Informatică, Universitatea "Al. I. Cuza", Iași și Institutul de Cercetări în Inteligență Artificială, Academia Română, București

**Gabriela Haja**, Institutul de Filologie Română "Al. Philippide", Academia Română, Iași

**Radu Ion**, Institutul de Cercetări în Inteligență Artificială, Academia Română, București

**Rada Mihalcea**, Universitatea North Texas, SUA

**Vivi Năstase**, EML Research, Germania

**Constantin Orăsan**, Universitatea Wolverhampton, Anglia

**Oana Postolache**, ISI - Universitatea California, SUA

**Irina Prodanoff**, ILC-Pisa și Universitatea Pavia, Italia

**Georgiana Pușcașu**, Universitatea Wolverhampton, Anglia

**Violeta Serețan**, Departamentul de lingvistică, Universitatea Geneva, Elveția

**Valentin Tablan**, Universitatea Sheffield, Anglia

**Amalia Todirașcu**, Universitatea Marc Bloch, Strasbourg, Franța

**Doina Tătar**, Universitatea "Babeș-Bolyai", Cluj-Napoca

**Horia-Nicolai Teodorescu**, Institutul de Informatică Teoretică, Academia Română și Universitatea Tehnică „Gh. Asachi”, Iași

**Dan Tufiș**, Institutul de Cercetări în Inteligență Artificială, Academia Română, București și Universitatea "Al. I. Cuza", Iași

**Adriana Vlad**, Facultatea de Electronică, Universitatea Politehnica București și Institutul de Cercetări în Inteligență Artificială, Academia Română, București

**COMITETUL DE ORGANIZARE:**

**Dan Cristea**, FII-UAIC și IIT-AR (dcristea@info.uaic.ro)

**Corina Dima**, FII-UAIC (cdima@info.uaic.ro)

**Maria Husarciuc**, LITERE-UAIC și FII-UAIC (mhusarciuc@gmail.com)

**Adrian Iftene**, FII-UAIC (adiftene@info.uaic.ro)

**Mihai-Alex Moruz**, FII-UAIC (mmoruz@info.uaic.ro)

**Ionuț Pistol**, FII-UAIC (ipistol@info.uaic.ro)

**Diana Trandabăț**, FII-UAIC și IIT-AR (dtrandabat@info.uaic.ro)

**Dan Tufiș**, ICIA-AR și FII-UAIC (tufis@racai.ro)

# Cuprins

<b>Cuvânt înainte</b>	<b>7</b>
-----------------------	----------

## **Capitolul 1: Resurse lingvistice și instrumente pentru prelucrarea vorbirii** **9**

<i>Adrian Turculeț, Vasile Apopei, Doina Jitcă</i> Studiul variației intonaționale în limba română literară folosind o măsură a distanței prozodice .....	11
<i>Horia-Nicolai Teodorescu, Monica Silvia Feraru</i> De ce nu place vocea sintetizată? – câteva elemente de comparație cu vocea umană .....	21
<i>Diana Hanes, Cristina Petrea, Andi Buzo, Vladimir Popescu, Corneliu Burileanu</i> Baza de date în limba română pentru recunoașterea vorbirii spontane .....	31
<i>Marius-Dan Zbancioc, Horia-Nicolai Teodorescu</i> Metodă ierarhică de detecție a fundamentalei .....	41

## **Capitolul 2: Platforme, dicționare și corpuri adnotate pentru prelucrarea textelor** **53**

<i>Dan Cristea, Ionuț Cristian Pistol</i> Limba română în perspectiva Clarin .....	55
<i>Neculai Curteanu, Alex Moruz, Diana Trandabăț, Cecilia Bolea, Mădălina Spătaru, Maria Husarciuc</i> Parsarea arborilor de sensuri și segmentarea la definiții în dicționarul tezaur eDTLR .....	65
<i>Radu Ion</i> Segmentarea în unități textuale atomice a intrărilor din dicționarul limbii române în vederea analizei structurale .....	75
<i>Doina Spiță, Claudia Bîzdîgă</i> Platformă plurilingvă de formare și autoformare în domeniul limbilor romanice .....	83
<i>Nadia Luiza Dincă</i> Considerații teoretice asupra aplicabilității unei baze de date cu exemple de traducere .....	93

## **Capitolul 3: Aplicații ale tehnologiilor lingvistice** **103**

<i>Adrian Iftene, Ancuța Rotaru, Dana-Alina Marcu</i> Evaluarea răspunsurilor oferite de un sistem de tip întrebare-răspuns pentru limba română .....	105
<i>Maria Husarciuc</i> Echivalarea în limba română a unităților frazeologice infinitivale din limba franceză .....	115
<i>Alexandru Ceaușu</i> Colectarea și procesarea documentelor românești ale corpului Jrc-Acquis .....	125
<i>Irimia Elena</i> Experimente de traducere automată bazată pe exemple pentru limbile engleză/română .....	131
<i>Dan Ștefănescu, Dan Tufiș</i> CONAN – detecția posibilelor conotații ale unui text .....	141
<i>Petic Mircea</i> Completarea automată a resurselor lingvistice românești .....	151

## **Index de autori** **161**



## CUVÂNT ÎNAINTE

Acest volum include lucrările celei de a șasea ediții a Atelierului, în seria de manifestări ale Consorțiului de Informatizare pentru Limba Română, devenită oglindă a rezultatelor cercetărilor ce sunt dedicate în fiecare an domeniului Tehnologiei Limbajului (TL) din perspectiva limbii române. Volumul identifică principalele direcții și progresele realizate în acest domeniu pe parcursul anului 2008.

De data aceasta am grupat cele 15 lucrări acceptate (din 20 primite) în 3 capitole: Resurse lingvistice și instrumente pentru prelucrarea vorbirii, Platforme, dicționare și corpusuri adnotate pentru prelucrarea textelor și Aplicații ale tehnologiilor lingvistice, modificând așadar structura din anii trecuți, pentru că am considerat că această nouă împărțire reflectă mai adecvat orientările actuale ale lucrărilor ce au primit acceptul Comitetului de Program.

Lucrările incluse în volum descriu cercetări desfășurate în diverse proiecte naționale sau internaționale precum și rezultatele obținute de doctoranzi în elaborarea tezelor lor. Îmbucurător este faptul că un număr din ce în ce mai mare de tineri talentați își aleg TL ca domeniu de cercetare. În cursul anului 2009, cel puțin șase dintre tinerii cercetători care au lucrări incluse în acest volum își vor finaliza tezele de doctorat, având contribuții majore în prelucrarea limbii române în contextul multilingv al societății informaționale.

La aproximativ o lună după Atelier, în 14 și 15 ianuarie 2009, a avut loc la Luxemburg manifestarea *Language Technology Days*, organizată de *DG Information Society and Media Unit INFISO/E1 – Language Technologies, Machine Translation* și prezidată de Roberto Cencioni și Kimmo Rossi. Această importantă întâlnire a urmărit informarea în rândul potențialilor propunători europeni de proiecte în domeniul TL asupra noilor oportunități de finanțare deschise de apelul 4, publicat în noiembrie 2008, al FP7-ICT (*Framework Program 7 of the Information and Communication Technologies*) și de apelul 3 al ICT-PSP (*Information and Communication Technologies Policy Support Programme*), ce va fi publicat la sfârșitul lunii ianuarie 2009. Cele două direcții principale menționate în aceste apeluri se referă la traducerea automată și la exploatarea multilingvă a web-ului (modele, arhitecturi și instrumente pentru sisteme de traducere text și voce auto-adaptive, standarde *de jure* și *de facto* în gestiunea multilingvă a web-ului). În prima zi a evenimentului de la Luxemburg au fost invitați să prezinte situația domeniului TL în țările lor trei experți din țări nou intrate în UE: România, Polonia și Ungaria. România a fost apreciată de mulți participanți ca având un nivel competitiv, raportat chiar față de țări cu tradiție în domeniul TL, atât din punctul de vedere al învățământului care pregătește specialiști, cât și din cel al nivelului instrumentelor și resurselor specifice dezvoltate deja. Credem că la această situație favorabilă au contribuit în mare măsură întâlnirile organizate de către ConsILR, schimbarea pozitivă de atitudine a Ministerului Educației, Cercetării și Tineretului față de problemele tehnologiei limbii române prin susținerea unor proiecte semnificative ca ambiție științifică și nivel de finanțare, precum eDLR (Dicționarul Tezaur al Limbii Române în format electronic – <https://consilr.info.uaic.ro/edtlr/wiki>) ori SIR-RESDEC (Sistem de Întrebare-Răspuns în limbile Română și Engleză cu Spații Deschise de Căutare – <https://sir-resdec.racai.ro:450/>), precum și participarea specialiștilor români în proiecte europene importante, cum ar fi CLARIN (*Common Language Resources and Technology Infrastructure* – <http://www.clarin.eu/>), ALEAR (*Artificial Language Evolution on Autonomous Robots*) ori FlareNet (*Fostering Language Resources Network* – <http://www.ilc.cnr.it/flarenet/>).

Schimbările de experiență facilitate de seria de întâlniri ale Atelierului contribuie la apropierea nivelului cercetărilor în tehnologia limbii române de cel mondial, la statornicirea unei terminologii coerente în limba română pentru acest domeniu și, nu în ultimul rând, la cunoașterea eforturilor comune desfășurate de specialiști din diferite centre ale României ori din afara ei. În acest sens, ne propunem ca la edițiile viitoare ale atelierului să includem o secțiune nouă de prezentări care să fie dedicată descrierilor de proiecte naționale sau internaționale, furnizând informații la zi (participanți, obiective, realizări etc.); o altă secțiune nouă ar putea fi dedicată discuțiilor de tip „brainstorming” pe marginea unor propuneri preliminare de proiecte comune ale membrilor Consorțiului.

În deschiderea manifestării din decembrie 2008, ale cărei lucrări le oferim publicului în această carte, Consorțiul a hotărât ca pe viitor manifestarea să se transforme din atelier de lucru în conferință cu participare internațională, pe considerentul maturității științifice dobândite deja și al notorietății ei în rândul cercetătorilor de pretutindeni care se preocupă de limba română prin prisma metodelor computaționale.

Ianuarie 2009

Editorii



## **CAPITOLUL 1**

### **RESURSE LINGVISTICE ȘI INSTRUMENTE PENTRU PRELUCRAREA VORBIRII**



# STUDIUL VARIAȚIEI INTONAȚIONALE ÎN LIMBA ROMÂNĂ LITERARĂ FOLOSIND O MĂSURĂ A DISTANȚEI PROZODICE

ADRIAN TURCULEȚ<sup>1</sup>, VASILE APOPEI<sup>2</sup>, DOINA JITCĂ<sup>2</sup>

<sup>1</sup>*Facultatea de Litere, Universitatea „Al. I. Cuza” Iași*

<sup>2</sup>*Institutul de Informatică Teoretică, Academia Română - Filiala Iași*

[aturcu@uaic.ro](mailto:aturcu@uaic.ro), [vapopei@iit.tuiasi.ro](mailto:vapopei@iit.tuiasi.ro)

## Rezumat

În lucrare se prezintă posibilitatea de a evalua în mod obiectiv distanța prozodică dintre contururile melodice observate la vorbitorii de limbă română proveniți din aceeași zonă sau din zone diferite. Calculul distanței prozodice se bazează pe utilizarea coeficientului de intercorelație dintre curbele frecvenței F0, avându-se în vedere de asemenea, și cele ale duratelor și intensităților (energia) segmentelor vocalice. Aplicat la un corpus de enunțuri, calculul coeficienților de corelație confirmă, în general, rezultatele analizei auditive și vizuale. Cel mai apropiat de perceperea auditivă și vizuală este coeficientul de corelație aplicat contururilor frecvenței F0.

## 1. Introducere

În cercetările lingvistice actuale a prozodiei s-a simțit nevoia de a evalua obiectiv (a confirma sau a infirma) concluziile percepției auditive asupra diferențelor melodice ale unor rostiri, cu ajutorul foneticii instrumentale acustice. Fără a deveni superfluă sau a putea fi înlocuită, analiza auditivă a primit, în ultimele decenii, un sprijin solid și obiectivat în sensul bazării pe cercetarea cantitativă obținută prin analiza și sinteza computerizată a semnalului vocal.

Într-o primă etapă a studierii prozodiei, lingviștii au beneficiat de rezultatele metodelor de analiză a semnalului vocal referitoare la extragerea componentelor armonice ale semnalului vocal și în special a traseului frecvenței F0. Atât percepția auditivă, cât și cea vizuală, ambele categoriale, subiective, par însă a avea nevoie de cuantificarea matematică. Astfel s-a născut ideea de măsurare a „distanței prozodice” dintre curbele melodice ale vorbirii pe baza variației tonului fundamental, a duratei și intensității segmentelor vocalice din cadrul unităților intonaționale.

În lipsa unei definiții clare a conceptului de *distanță prozodică*, acesta se poate raporta la conceptul mai larg de *distanță lingvistică*, utilizat de cei care se ocupă cu tipologia lingvistică în încercările de a cuantifica similitudinea sau diferențele dintre limbi diferite sau dintre varietățile aceleiași limbi, de exemplu, între dialectele ei. De exemplu, cercetările auditive și acustice arată că intonația moldovenilor este mai apropiată de cea a muntenilor, în timp ce ardelenii au o intonație mai deosebită, având deci o „distanță prozodică” mai mare.

Compararea curbelor frecvenței F0 cu ajutorul statisticii matematice<sup>1</sup> oferă posibilitatea

---

<sup>1</sup> De aceea, utilizarea restrânsă a termenului *distanță prozodică* se referă la o măsură a acesteia calculată cu ajutorul unor formule matematice și nu la evaluarea auditivă sau vizuală a acesteia.

de a obiectiva „distanța prozodică” vizuală. Un sprijin în acest sens a primit cercetarea prozodiei de la domeniul aplicat al învățării limbilor străine sau al deprinderii intonației de către persoane cu hipoacuzie. Având în vedere situațiile de antrenament în care se urmărește corectitudinea reproducerii unui contur tonal, cercetătorul olandez Dik J. Hermes (1998a, b) propune mai multe tipuri de măsurători pentru evaluarea automată a gradului de corectitudine, printre care și indicele de corelație. El vorbește mai ales de „similaritatea/disimilaritatea auditivă și vizuală a conturilor tonale”, dar, uneori, și de „distanța” dintre contururi.

În cadrul proiectului AMPER<sup>2</sup> a fost reactivat conceptul de distanță prozodică în cercetarea variației diatopice a prozodiei, pentru evaluarea matematică a similarității/disimilarității a două contururi prozodice aparținând unor varietăți diferite. A. Romano (1999\*2001: 226-235) utilizează la compararea modelelor intonaționale ale unor dialecte apropiate coeficientul de intercorelație a curbelor după formula lui Pearson; în timp ce indicii de similaritate a curbelor F0, precum și indicii de corelație a energiei au o valoare ridicată (peste 0.90), indicii de corelație a duratei ating și valori ceva mai scăzute. Ulterior, A. Romano și R. Miotti (2008) utilizează indicii de intercorelație a curbelor F0 pentru evaluarea *distanței prozodice* dintre o varietate venețiană (levantina) și varietatea iberică din Malaga. Pentru limba română s-a realizat un studiu al variației diatopice a intonației folosind o metodologie similară (Turculeț et al. 2008).

În lucrarea de față, în secțiunea 2, se va prezenta metodologia de realizare a corpusului de analiza prozodică iar în secțiunea 3, metodologia de prelucrare a datelor. Secțiunile 4 și 5 conțin câteva rezultate ale analizei prozodice perceptuale-vizuale, și respectiv, cele ale analizei realizată pe baza coeficienților de intercorelație.

## **2. Prezentarea metodologiei de realizare a corpusului de analiză prozodică**

Înregistrările corpusului de analiză prozodică au fost realizate prin anchete pe teren în opt centre culturale din țară: Baia Mare, Brașov, București, Cluj, Craiova, Iași, Oradea, Sibiu, Timișoara și unul din Republica Moldova: Chișinău. Subiecții selectați, câte unul de sex masculin și feminin, au studii universitare și utilizează în mod curent varietatea cotidiană a limbii române standard. Din anchete s-au selectat înregistrările cu rostiri „neutre” (adică rostiri fără focalizarea expresă a unui constituent și fără conotații expresiv-afective evidente) ale unor enunțuri asertive și interogative, ambele în variante afirmative și negative.

Înregistrările au fost stocate în fișiere de tip „wav” al căror nume conține informații despre punctele de anchetă, vorbitori (masculin sau feminin), structura morfologico-sintactică și silabico-accentuală a enunțului precum și, de modalitățile de rostire: afirmativă (*a*), negativă (*n*), interogativă (*i*) și interogativă negativă (*m*). Fiecare enunț a fost repetat de fiecare subiect de cel puțin trei ori, rezultând, la un corpus de 45 de fraze, un număr de 540 de fișiere de sunet. Aceste fișiere au fost adnotate la nivelul segmentelor vocalice cu ajutorul programului Praat, realizându-se un corpus de 540 de fișiere cu semnal vocal de tip *wav* și același număr de fișiere cu etichete de tip *TextGrid*. Pentru calcularea distanțelor prozodice (vezi tabelele de mai jos) s-au selectat trei

<sup>2</sup> Acest proiect are ca obiectiv realizarea primului atlas prozodic romanic: *Atlas Multimédia Prosodique de l'Espace Roman*.

STUDIUL VARIAȚIEI INTONAȚIONALE ÎN LIMBA ROMÂNĂ LITERARĂ  
FOLOSIND O MĂSURĂ A DISTANȚEI PROZODICE

enunțuri cu structuri morfologico-sintactice și silabico-accentuale diferite, rostite cu cele patru modalități diferite (codificate mai anterior cu  $a, n, i, m$ ):

- *Nevasta vede-un căpitan* - având codificarea structurii prin grupul de litere *twk*;
- *Nevasta frumoasă vede-un căpitan* - codificată prin *swk*;
- *Nevasta vede-un căpitan elegant* - codificată prin *twg*.

### 3. Prezentarea metodologiei de prelucrare a corpusului

Prelucrările asupra corpusului de voce au fost realizate în mediul de programare Matlab, conform metodologiei stabilite în cadrul proiectului AMPER, și vizează extragerea următorilor cinci parametrii prozodici pentru fiecare segment vocalic (relația 1): durata, valoarea maximă a energiei și frecvența F0 măsurată în câte trei puncte (inițial, mijloc, final).

$$\{D_1, E_1, (F_1^i, F_1^m, F_1^f); D_2, E_2, (F_2^i, F_{21}^m, F_2^f); \dots; D_n, E_n, (F_n^i, F_n^m, F_n^f)\} \quad (1)$$

unde  $n$  este numărul de segmente vocalice din cadrul rostirii.

Evoluția valorilor acestor parametrii pe durata unei rostiri generează următoarele trei contururi pe care au fost folosite în calculul coeficienților de corelație: curba frecvenței F0 (un contur stilizat), curba duratelor și a energiilor segmentelor vocalice. Duratale segmentelor vocalice se calculează prin diferența reperelor de timp asociate fiecărei etichete. Din valorile energiei calculate pe ferestre de 160 de eșantioane, cu factor de suprapunere  $\frac{1}{2}$ , se selectează valoarea maximă pe fiecare segment vocalic.

Din aceste valori se extrag cele trei curbe folosite pentru calculul coeficientului de intercorelație:

- curba frecvenței F0 :  $\{(F_1^i, F_1^m, F_1^f); (F_2^i, F_{21}^m, F_2^f); \dots; (F_n^i, F_n^m, F_n^f)\}$
- curba duratelor:  $\{D_1; D_2; \dots; D_n\}$
- curba de energie  $\{E_1; E_2; \dots; E_n\}$ .

Pentru fiecare rostire (fișier tip *wav*) rezultă un fișier cu un număr de parametri egal cu numărul de vocale multiplicat cu 5. Valorile rezultate pentru cele trei repetiții realizate de fiecare tip de rostire (cu aceeași structură, aceeași modalitate și același tip de subiect) Pentru cele cinci valori de parametri prozodici pentru fiecare vocală s-a calculat media și dispersia. Valorile medii rezultate sunt salvate în fișiere de tip text, pentru fiecare segment vocalic, împreună cu reperele de timp ale acestora.

Informațiile din aceste fișiere constituie datele pe baza cărora s-au calculat coeficienții de similaritate pentru analiza comparativă a prozodiei pe baza textelor cu valorile medii a câte trei repetiții. Coeficienții de similaritate s-au calculat cu formula coeficientului de intercorelație a lui Pearson,  $C_{xy}$ , aplicată valorilor de pe curbele de frecvență, durată și energie din două înregistrări selectate pentru comparație (relația 2).

$$C_{xy} = \frac{\sum_{i=1}^N (x_i - \bar{x}) * (y_i - \bar{y})}{N * \sigma_x * \sigma_y} \quad (2)$$

unde:  $N$  = numărul de puncte ale curbelor pentru care calculează coeficientul de intercorelație;

$x_i, y_i$  = valorile cu indicele  $i$  de pe curbele supuse analizei;

$\bar{x}, \bar{y}$  = valorile medii corespunzătoare curbelor supuse analizei;

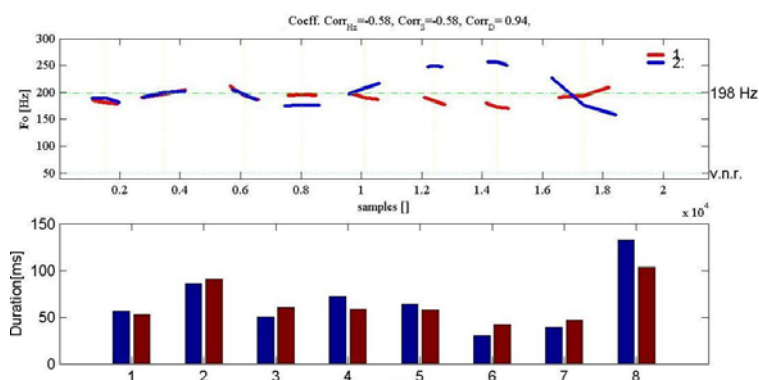
$\sigma_x, \sigma_y$  = valorile dispersiilor corespunzătoare curbelor supuse analizei

Analiza variației diatopice a intonației la nivelul limbii române standard a fost realizată prin valorile coeficienților de corelație calculați între contururile melodice ale aceluiași enunț (cu aceeași structură sintactico-lexicală și fonologică segmentală) rostit de către fiecare cuplu de subiecți din celelalte nouă localități, raportate la rostirea subiecților bucureșteni considerată ca reprezentând intonația standard. Unele concluzii ale raportării modelelor intonaționale stabilite prin analiza auditivă și acustică la distanțele prozodice reprezentate de coeficienții de corelație sunt prezentate în secțiunea 5.

#### 4. Câteva rezultate ale analizei prozodice perceptuale

Analiza acustică atestă prezența unor modele intonaționale regionale în vorbirea literară. Aceste particularități prozodice sunt chiar mai persistente decât particularitățile fonetice segmentale sau lexicale, permițând identificarea zonei de proveniență a vorbitorului.

Intonația subiecților din București, a constituit punctul de reper pentru comparația cu modelele intonaționale folosite de ceilalți subiecți. Contururile intonaționale stilizate prin eliminarea unor efecte de micro-prozodie<sup>3</sup> au fost reprezentate prin grafice care prezintă valorile medii rezultate din prelucrarea a trei rostiri ale fiecărui enunț (cf. Fig. 1-3). Subiecții bucureșteni utilizează modele intonaționale standard (v. Dascălu-Jinga (2001, 2005), distingându-se de vorbitorii din alte zone și printr-o intonație mai „economică”, cu o extensiune tonală (pitch range) mică, utilizând doar două accente tonale, cel de pe prima silabă accentuată și cel ce determină modalitatea de realizare a propoziției (asertivă/ interogativă). Alte accente lexicale sunt dezaccentuate în modelele neutre. O ușoară focalizare se poate observa, uneori, mai ales pe cuvântul final al enunțului.



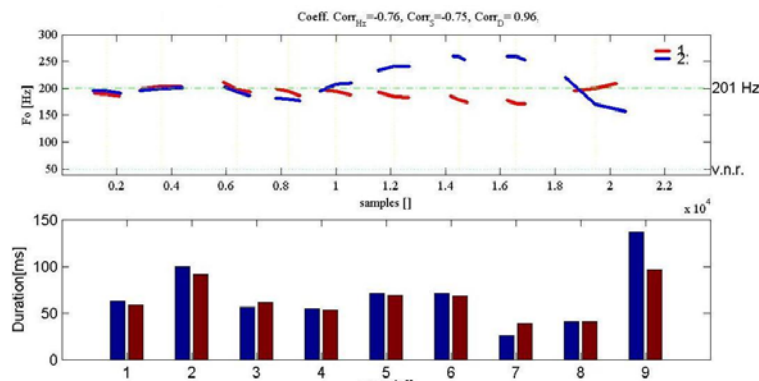
**Figura 1:** Contururile frecvenței F0 și ale duratei în rostirea subiectului feminin din București (culoare deschisă), respectiv a subiectului feminin din Timișoara (culoare închisă) a enunțului *Nevasta vede-un căpitan?*

<sup>3</sup> Analiza acustică a fost realizată după programe elaborate în cadrul proiectului AMPER de către A. Romano și A. Rilliard.

STUDIUL VARIAȚIEI INTONAȚIONALE ÎN LIMBA ROMÂNĂ LITERARĂ  
FOLOSIND O MĂSURĂ A DISTANȚEI PROZODICE

La polul opus se află vorbitorii ardeleni, care prezintă aproape regulat o extensiune tonală largă (indiferent de registrul tonal al vorbitorului) în cadrul enunțului și schimbarea pronunțată (de obicei urcarea) tonului fundamental pe silaba accentuată a fiecărui grup accentual. Aceasta dă impresia de „cântat” și de „emfază” pe care o percep interlocutorii neardeleni ai vorbitorilor ardeleni<sup>4</sup>.

Vorbitorii originari din Transilvania în sens larg (cu provinciile istorice adiacente Banat, Crișana, Maramureș) păstrează, atunci când folosesc stilul colocvial (adesea, chiar cel formal) al limbii standard, contururi intonative specifice, în special la interogativele totale afirmative și negative. Conturul interogativ ardelenesc, în special în sintagma verbală, are o formă generală concavă, opusă formei convexe observate la vorbitorii din Muntenia și Moldova, iar conturul melodic terminal este descendent, opus celui ascendent al modelului muntean-moldovean. Pe silaba accentuată a verbului începe o urcare treptată amplă a tonului fundamental, care se extinde asupra întregii sintagme verbale și se termină printr-o coborâre abruptă pe ultima silabă accentuată a enunțului.



**Figura 2.** Contururile frecvenței F0 și ale duratei în rostirea subiectului feminin din București (culoare deschisă), respectiv a subiectului feminin din Timișoara (culoare închisă) a enunțului „*Nevasta nu vede-un căpitan?*”.

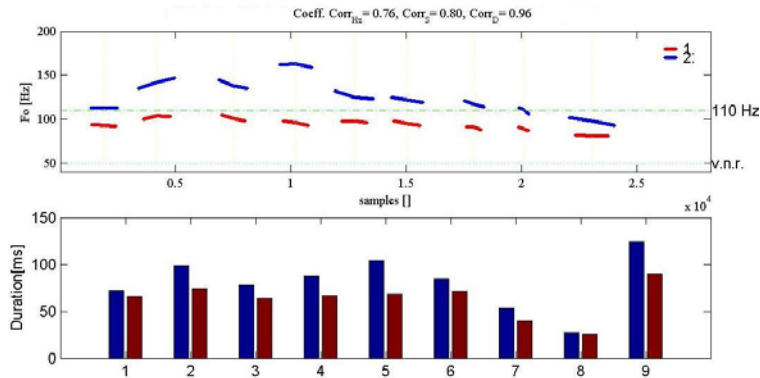
În enunțurile cu structura SVO, tonul silabei accentuate a verbului, care poate fi și plat sau ușor descendent, realizează o cezură între sintagma nominală și cea verbală; cu un aliniament întârziat, F0 urcă pe sintagma verbală sub formă de cupolă sau de platou până la ultima silabă accentuată. În cazul interogativelor totale negative, adverbul „*nu*” poartă accentul sintagmei verbale, realizând cezura ritmico-sintactică urmată de conturul tonal cunoscut (figura 2). În partea de sud a Transilvaniei (Brașov, Sibiu), dar și la subiecții din Oradea, Baia Mare sau Chișinău, modelul acesta „ardelenesc”<sup>5</sup> este concurat de modelul muntenesc-moldovenesc standard.

Moldovenii, mai ales cei din partea de nord (subiecții din Iași), se disting în special la enunțativele negative printr-o emfază puternică a adverbului „*nu*”, care poartă accentul tonal principal (nuclear) al enunțului (v. figura. 3, culoarea albastră). Această trăsătură,

<sup>4</sup> Unele teste de percepție realizate cu ajutorul studenților originari din diferite zone ale țării vor fi completate și expuse cu altă ocazie.

<sup>5</sup> Am numit acest model intonațional „ardelenesc” (între ghilimele), deoarece apare și la subiecții noștri originari din Chișinău.

emfatică la origine, s-a „gramaticalizat”, fiind prezentă în enunțurile rostite obișnuit, neutral.



**Figura 3:** Contururile frecvenței F0 și ale duratei în rostirea subiectului masculin din București (culoare deschisă), respectiv a subiectului masculin din Iași (culoare închisă) a enunțului *Nevasta nu vede-un căpitan*.

### 5. Rezultatele analizei prozodice pe baza măsurii distanței prozodice

Înainte de a aplica măsurarea distanței prozodice la compararea contururilor melodice realizate de vorbitori din zone diferite, pentru a determina mai obiectiv variația diatopică a distanței prozodice, am probat aplicarea distanței prozodice (măsurată prin coeficienții la corelație ai frecvenței F0) la: rostirile repetate ale aceluiași enunț de către același vorbitor; același enunț rostit de către fiecare dintre cei doi subiecți (feminin și masculin) din aceeași localitate; enunțuri cu modalități diferite (asertivă și interogativă, afirmativă și negativă) rostite de către același subiect.

Cum era de așteptat, în primul caz, indicele de corelație a F0 este ridicat: peste 0.80, în cele mai multe cazuri peste 0.90; în cel de al doilea caz, există diferențe (în afara registrului melodic diferit) care fac să coboare indicele de corelație spre 0,70 și chiar spre 0,50. În al treilea caz, indicii sunt mai mari dacă se compară modalități cu contururi asemănătoare (de exemplu, *swkn: Nevasta nu vede-un căpitan* și *swkm: Nevasta nu vede-un căpitan?* la subiectul feminin din București: 0.61), dar sunt mici dacă se compară contururi tonale diferite (la același subiect, *swka: Nevasta vede-un căpitan* și *swki: Nevasta vede-un căpitan?* : 0.22). Cu ajutorul indicelui de corelație a frecvenței F0 se poate deci preciza dacă două enunțuri rostite de același vorbitor sau de vorbitori diferiți au același *pattern* intonativ sau au *pattern*-uri diferite.

Compararea coeficienților de intercorelație corespunde, în cea mai mare parte, cu evaluarea distanței/apropierii prozodice dintre localități realizată auditiv (la ascultarea enunțurilor) și vizual (prin compararea traseelor F0 rezultate din analiza acustică). Rezultatele acestei comparații confirmă și fundamentează concluziile lucrării privind existența unor *pattern*-uri intonaționale diferite la vorbitori ai limbii literare provenind din zone diferite. Valorile absolute ale acestor coeficienți se referă la contururile intonaționale stilizate și pot fi apreciate ca având o valoare relativă care depinde de reprezentativitatea subiecților selectați și de relevanța enunțurilor „neutrale” obținute de anchetator.



STUDIUL VARIAȚIEI INTONAȚIONALE ÎN LIMBA ROMÂNĂ LITERARĂ  
FOLOSIND O MĂSURĂ A DISTANȚEI PROZODICE

Cele trei tabele de mai jos, conțin coeficienții de intercorelație pentru contururile frecvenței F0, duratelor și energiilor ale celor trei propoziții selectate (*twka*, *swka*, *twga*), rostite în cele patru modalități de către doi subiecți (prima valoare aparține subiectului feminin, iar cea de a doua subiectului masculin) din cele nouă localități menționate, raportate la rostirile corespunzătoare ale subiecților bucureșteni. Aceste valori reprezintă *distanțele prozodice* (de frecvență F0, durată și energie) dintre contururile intonaționale ale aceluiași enunțuri rostite de vorbitori bucureșteni și de vorbitori din celelalte nouă centre culturale românești.

Tabel 1: Coeficienții de corelație între contururile frecvenței F0.

	Baia Mare	Brașov	Cluj	Oradea	Sibiu	Timișoara	Iași	Chișinău	Craiova
<b>twka</b>	0.82/ 0.85	0.88/ 0.94	0.98/ 0.91	0.88/ 0.94	0.98/ 0.91	0.97/ 0.89	0.74/ 0.80	0.75/ 0.93	0.95/ 0.55
<b>twki</b>	0.11/ 0.63	0.72/ -0.36	0.00/ -0.06	0.29/ 0.51	0.66/ 0.00	-0.58/ -0.42	0.70/ 0.58	0.61/ 0.01	0.81/ 0.73
<b>twkm</b>	0.43/ -0.37	0.66/ -0.58	-0.80/ -0.62	0.29/ -0.56	0.69/ 0.04	-0.76/ -0.58	0.85/ 0.73	-0.66/ -0.30	0.73/ 0.59
<b>twkn</b>	0.78/ 0.84	0.81/ 0.89	0.84/ 0.88	0.87/ 0.73	0.85/ 0.80	0.94/ 0.83	0.81/ 0.76	0.85/ 0.76	0.92/ 0.30
<b>swka</b>	0.74/ 0.69	0.94/ 0.84	0.87/ 0.87	0.90/ 0.62	0.92/ 0.83	0.93/ 0.87	0.84/ 0.52	0.91/ 0.65	0.89/ 0.64
<b>swki</b>	0.35/ 0.32	0.67/ -0.11	0.11/ 0.20	0.31/ 0.05	0.72/ 0.22	-0.49/ -0.10	0.69/ 0.66	0.01/ -0.01	0.84/ 0.40
<b>swkm</b>	0.07/ -0.52	0.76/ -0.54	-0.62/ -0.53	0.01/ -0.42	0.45/ 0.01	-0.61/ -0.43	0.72/ 0.72	-0.62/ 0.12	0.68/ 0.31
<b>swkn</b>	0.68/ 0.83	0.93/ 0.88	0.79/ 0.92	0.78/ 0.67	0.90/ 0.77	0.66/ 0.86	0.63/ 0.76	0.84/ 0.76	0.94/ 0.79
<b>twga</b>	0.76/ 0.75	0.90/ 0.80	0.58/ 0.83	0.71/ 0.84	0.93/ 0.85	0.94/ 0.61	0.59/ 0.47	0.87/ 0.85	0.63/ 0.71
<b>twgi</b>	0.80/ 0.53	0.62/ -0.35	0.35/ 0.27	0.29/ -0.39	0.60/ 0.33	-0.56/ -0.23	0.62/ 0.48	-0.56/ -0.13	0.45/ 0.60
<b>twgm</b>	0.37/ -0.32	0.73/ -0.46	-0.77/ -0.63	-0.70/ -0.63	0.53/ 0.60	-0.73/ -0.54	0.84/ 0.70	-0.63/ -0.13	0.64/ -0.65
<b>twgn</b>	0.76/ 0.81	0.89/ 0.83	0.78/ 0.85	0.77/ 0.89	0.83/ 0.78	0.95/ 0.88	0.86/ 0.76	0.72/ 0.77	0.90/ 0.65

Tabel 2: Coeficienții de corelație între curbele duratelor vocalelor

	Baia Mare	Brașov	Cluj	Oradea	Sibiu	Timișoara	Iași	Chișinău	Craiova
<b>twka</b>	0.86/ 0.61	0.91/ 0.75	0.88/ 0.57	0.92/ 0.83	0.86/ 0.51	0.79/ 0.63	0.89/ 0.88	0.84/ 0.68	0.89/ 0.72
<b>twki</b>	0.91/ 0.52	0.98/ 0.77	0.85/ 0.43	0.98/ 0.61	0.90/ 0.66	0.94/ 0.55	0.90/ 0.77	0.86/ 0.71	0.91/ 0.61
<b>twkm</b>	0.90/ 0.83	0.86/ 0.84	0.84/ 0.81	0.90/ 0.86	0.88/ 0.78	0.96/ 0.88	0.90/ 0.93	0.80/ 0.89	0.96/ 0.85
<b>twkn</b>	0.84/ 0.88	0.83/ 0.92	0.86/ 0.74	0.91/ 0.83	0.87/ 0.76	0.78/ 0.89	0.80/ 0.96	0.65/ 0.74	0.93/ 0.92
<b>swka</b>	0.87/ 0.86	0.97/ 0.89	0.89/ 0.79	0.85/ 0.82	0.92/ 0.80	0.89/ 0.88	0.91/ 0.95	0.84/ 0.88	0.86/ 0.63
<b>swki</b>	0.93/ 0.82	0.89/ 0.92	0.83/ 0.66	0.88/ 0.78	0.86/ 0.88	0.94/ 0.80	0.95/ 0.92	0.76/ 0.78	0.91/ 0.79
<b>swkm</b>	0.85/ 0.43	0.89/ 0.88	0.74/ 0.84	0.83/ 0.81	0.81/ 0.78	0.90/ 0.92	0.98/ 0.93	0.81/ 0.88	0.90/ 0.85
<b>swkn</b>	0.89/ 0.78	0.93/ 0.80	0.81/ 0.78	0.77/ 0.72	0.83/ 0.83	0.77/ 0.80	0.92/ 0.88	0.83/ 0.84	0.90/ 0.81
<b>twga</b>	0.90/ 0.79	0.93/ 0.70	0.82/ 0.71	0.89/ 0.68	0.88/ 0.78	0.92/ 0.58	0.89/ 0.72	0.89/ 0.67	0.95/ 0.65

<b>twgi</b>	0.88/ 0.87	0.92/ 0.94	0.85/ 0.75	0.94/ 0.86	0.91/ 0.69	0.94/ 0.74	0.84/ 0.83	0.85/ 0.84	0.90/ 0.75
<b>twgm</b>	0.86/ 0.88	0.86/ 0.92	0.85/ 0.78	0.94/ 0.96	0.91/ 0.76	0.89/ 0.84	0.86/ 0.92	0.86/ 0.91	0.91/ 0.74
<b>twgn</b>	0.87/ 0.83	0.82/ 0.87	0.86/ 0.84	0.84/ 0.90	0.82/ 0.79	0.87/ 0.79	0.79/ 0.84	0.79/ 0.89	0.88/ 0.71

Tabel 3: Coeficienții de corelație între curbele energiei vocalelor

	Baia Mare	Brașov	Cluj	Oradea	Sibiu	Timișoara	Iași	Chișinău	Craiova
<b>twka</b>	0.65/ 0.81	0.78/ 0.77	0.82/ 0.90	0.91/ 0.87	0.93/ 0.90	0.54/ 0.74	0.69/ 0.75	0.56/ 0.81	0.67/ 0.78
<b>twki</b>	0.71/ 0.51	0.64/ 0.60	0.37/ 0.35	0.75/ 0.77	0.79/ 0.42	0.35/ 0.47	0.43/ 0.80	0.31/ 0.85	0.71/ 0.84
<b>twkm</b>	0.52/ 0.29	0.76/ 0.42	0.61/ -0.03	0.73/ 0.83	0.88/ 0.17	0.39/ 0.62	0.52/ 0.62	0.35/ 0.83	0.87/ 0.85
<b>twkn</b>	0.86/ 0.79	0.69/ 0.70	0.88/ 0.87	0.95/ 0.91	0.97/ 0.93	0.62/ 0.76	0.71/ 0.69	0.77/ 0.87	0.93/ 0.80
<b>swka</b>	0.65/ 0.54	0.52/ 0.59	0.60/ 0.74	0.96/ 0.71	0.92/ 0.85	0.72/ 0.45	0.53/ 0.48	0.56/ 0.72	0.46/ 0.84
<b>swki</b>	0.31/ 0.17	0.46/ 0.42	0.41/ 0.31	0.53/ 0.60	0.88/ 0.15	0.23/ 0.31	0.19/ 0.30	0.45/ 0.34	0.39/ 0.36
<b>swkm</b>	0.09/ 0.32	0.63/ 0.76	0.30/ 0.47	0.48/ 0.75	0.86/ 0.56	0.01/ 0.76	0.44/ 0.85	0.31/ 0.75	0.71/ 0.64
<b>swkn</b>	0.80/ 0.67	0.85/ 0.67	0.90/ 0.52	0.86/ 0.79	0.93/ 0.71	0.18/ 0.56	0.49/ 0.54	0.57/ 0.67	0.80/ 0.67
<b>twga</b>	0.94/ 0.68	0.83/ 0.62	0.78/ 0.56	0.91/ .51	0.95/ 0.81	0.55/ 0.59	0.57/ 0.59	0.65/ 0.75	0.58/ 0.57
<b>twgi</b>	0.57/ 0.63	0.22/ 0.38	0.58/ 0.64	0.82/ 0.75	0.81/ 0.54	0.20/ 0.73	0.13/ 0.59	0.20/ 0.46	0.62/ 0.59
<b>twgm</b>	0.66/ 0.32	0.66/ 0.57	0.58/ 0.74	0.48/ 0.68	0.90/ 0.69	0.29/ 0.47	0.45/ 0.64	0.04/ 0.66	0.69/ 0.31
<b>twgn</b>	0.81/ 0.66	0.90/ 0.67	0.90/ 0.67	0.56/ 0.86	0.91/ 0.90	0.54/ 0.76	0.64/ 0.80	0.81/ 0.73	0.66/ 0.69

Dintre cele trei tipuri de corelații prezentate în tabelele 1-3: pe baza variației frecvenței F0, a duratei și a energiei, corelația frecvenței F0 corespunde, în gradul cel mai înalt, evaluării perceptivă a asemănărilor / deosebirilor dintre modelele intonative comparate. Pe baza coeficientului de corelație general (luând în considerație toate cele patru modalități) dintre subiecții bucureșteni, pe de o parte, și subiecții din celelalte nouă orașe cercetate, pe de altă parte), se poate vorbi de o distanță prozodică relativ mică la subiecții din Iași: 0.71, Sibiu: 0.63, Craiova: 0.62; medie la Brașov: 0.51 și Baia Mare: 0.41; mare la Chișinău: 0.31, și foarte mare la Cluj: 0.29, Oradea: 0.23 și mai ales la Timișoara: 0.18.

O comparație între intonațiile vorbitorilor din București și Ardeal conform tabelului 1 indică grade de similaritate diferite pentru asertive și interogative, ultimele având o variație mult mai mare. Asertivele afirmative și negative (codificate cu *twka*, *twkn*, *swka*, *swkn*, *twga*, *twgn*) au un indice de corelație mediu ridicat: 0.81; în schimb interogativele totale afirmative (*twki*, *swki*, *twgi*) au un coeficient de corelație mediu mic: 0.24, iar cele negative (*twkm*, *swkm*, *twgm*) au chiar coeficient negativ: -0.12. Ultimele distanțe prozodice, considerabile, se datorează modelului specific ardelenesc al intonației enunțurilor interogative. Distanța prozodică mai mică a subiecților moldoveni față de cei bucureșteni se explică prin menținerea unui echilibru între toate cele patru modalități de rostire avute în studiu.

## STUDIUL VARIAȚIEI INTONAȚIONALE ÎN LIMBA ROMÂNĂ LITERARĂ FOLOSIND O MĂSURĂ A DISTANȚEI PROZODICE

Un alt procedeu de a calcula coeficientul de corelație a curbelor, plecând de la exprimarea frecvenței  $F_0$  în semitonuri, a dat valori identice sau foarte apropiate (cu diferențe de 1-2 sutimi de procent).

Coeficientului de corelație între contururile intonaționale ca măsură a distanței prozodice i s-a reproșat faptul de a reflecta *în mod global* similaritatea curbelor  $F_0$ , fără a ține seama de punctele cele mai importante din punct de vedere funcțional ale conturului. De exemplu, o deosebire funcțională cum este plasarea tonului nuclear pe *nu* în asertivele negative ale ieșenilor nu se reflectă suficient în coeficientul mediu de corelație ridicat: 0.80; în fig.2, coeficientul este ceva mai scăzut: 0.76 (cu frecvența exprimată în Hz) și 0,80 (cu frecvența exprimată în semitonuri).

Coeficienții corelației ai conturilor duratelor vocalelor sunt mult mai mari decât cei ai corelației curbelor de frecvență  $F_0$ , depășind, în toate cazurile procentul de 0,80. Menționăm, ca o curiozitate, faptul că subiecții din Iași au și în privința duratei cea mai bună corelație (0.88) cu subiecții bucureșteni.

Datele din tabelul 3 pun în evidență o corelație scăzută între curbele de energie la nivelul global al rostirilor. Analiza perceptuală pune în evidență o posibilă corelație a curbelor de energie doar pe anumite porțiuni ale enunțului; de exemplu, la începutul asertivelor, unde se observă o concentrare a energiei sonore care scade treptat odată cu declinația, în timp ce la interogativele totale, energia crește odată cu ridicarea finală a tonului. De asemenea, intensitatea vocalelor poate crește considerabil pe constituenții care poartă focusul contrastiv.

### **6. Câteva concluzii și perspective**

Folosirea coeficientului de intercorelație ca măsură a similarității conturilor frecvenței  $F_0$  confirmă evaluările perceptive (auditive și vizuale) cu privire la variabilitatea conturilor melodice utilizate de vorbitorii limbii române literare originari din zone diferite. Cele mai mari valori pentru gradul de similaritate se obțin în cazul conturilor medii care prezintă aceleași secvențe de creșteri-descreșteri pe toată durata lor, indiferent de gama de variație în care se desfășoară aceste secvențe. Valoarea coeficientului de intercorelație este mai scăzută în cazul conturilor intonaționale care manifestă tendințe de evoluție opuse pe perioade mai mari. La valori intermediare ale coeficientului de intercorelație avem de a face cu cazuri de similaritate pe grupuri sintactice.

Vom continua cercetările în vederea perfecționării procedeelelor de calculare a diferențelor prozodice dintre contururile prozodice utilizate de vorbitori în enunțuri care să reflecte atât variația diatopică a rostirii, cât și alte tipuri de variație: diastratică, diafazică, utilizând pentru aceasta corpusurile aflate la Seminarul de dialectologie și sociolingvistică al Facultății de Litere și corpusul SRoL (Teodorescu H.N, ș.a). Interesul actual pentru utilizarea distanței prozodice și ritmice între contururile intonaționale poate prilejui apariția unei metode noi în cercetarea prozodiei: *intonometria*, cu două direcții principale de dezvoltare: după modelul dialectometriei, care își propune cuantificarea asemănărilor și deosebirilor dintre dialecte, găsirea unor procedee cât mai adecvate pentru calcularea *distanțelor prozodice* (și *ritmice*) dintre varietăți lingvistice în scopul

clasificării tipologice a acestora; după modelul fonometriei, care își propune, stabilirea cu ajutorul statisticii, a *normelor de rostire* ale unei comunități lingvistice. Intonometria poate avea ca obiectiv precizarea *modelelor intonaționale neutre* ale unei comunități, precum și a abaterilor de la acestea în scopuri funcționale, pragmatice, expresiv-afective.

### Referințe bibliografice

- Dascălu-Jinga, Laurenția (2001), *Melodia limbii române vorbite*, Univers Enciclopedic, București.
- Dascălu-Jinga, Laurenția, Organizarea prozodică a enunțului, în *Gramatica limbii române, II, Enunțul*, Editura Academiei Române, București, 2005, p. 902-946.
- Lai, Jean-Pierre & Albert Rilliard (2008): *Distance prosodiques entre les variétés occitanes et sardes* (sub tipar).
- Hermes D.J. (1998a), Auditory and Visual Similarity of Pitch Contours, *Journal for Speech, Language, and Hearing Research*, vol. 41, 63-72, p.63-72.
- Hermes D.J. (1998b), Measuring the Perceptual Similarity of Pitch Contours, *Journal for Speech, Language, and Hearing Research*, vol. 41, 73-82, p.73-81.
- Mairano, Paolo & Antonio Romano (2008), *Distances rythmiques entre variétés romanes* (sub tipar).
- Miotti, Renzo & Antonio Romano (2008), *Distanze prosodiche tra varietà friulane, romene e ispaniche* (sub tipar).
- Rilliard, Albert & Jean-Pierre Lai (2007), La base de données AMPER et ses interfaces : structure et formats de données, exemple d'utilisation pour une analyse comparative de la prosodie de différents parlers romans, *I Jornadas Científicas AMPER-POR. Actas, Universidade de Aveiro*, 2007, p. 127-139.
- Romano, Antonio (1999\*2001), *Analyse des structures prosodiques des dialectes et d'italien régional parlés dans le Salento (Italie): approche linguistique et instrumentale*, Presses Universitaires du Septentrion, CEDEX France, 2001.
- Romano, Antonio & Renzo, Miotti (2008), *Un contributo per il confronto tra l'intonazione veneta e quella andalusa* (sub tipar).
- Teodorescu Horia-Nicolai, ș.a., *SRoL - Proiectul Sunetele Limbii Române*, [www.etc.tuiasi.ro/sibm/romanian\\_spoken\\_language/index.htm](http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.htm)
- Turculeț Adrian, Botoșineanu Luminița., Minuț Ana-Maria, Mladin Ioan-Constantin. (2008), Aspects de la variation diatopique de l'intonation au niveau de la langue roumaine littéraire, Simpozionul internațional *La Variation diatopique de l'intonation dans le domaine roumain et roman*, Iași , 20-21 octombrie 2008 (sub tipar).

# DE CE NU PLACE VOCEA SINTETIZATĂ? – CÂTEVA ELEMENTE DE COMPARAȚIE CU VOCEA UMANĂ

HORIA-NICOLAI TEODORESCU<sup>1,2</sup>, MONICA SILVIA FERARU<sup>1,2</sup>

<sup>1</sup> *Institutul de Informatică Teoretică al Academiei Române,  
Filiala Iași a Academiei Române*

<sup>2</sup> *Universitatea Tehnică Gheorghe Asachi din Iași*

*{hteodor ,mferaru}@etc.tuiasi.ro*

## Rezumat

Prezentăm o scurtă analiză a diferențelor care apar, la nivel de formați, între vocea sintetizată și vocea naturală. După expunerea scopului analizei și a metodologiei, prezentăm comparativ date privind valorile medii ale formațiilor pentru mai multe cuvinte și fraze scurte. Pentru comparație s-a utilizat o voce umană din corpusul adnotat SRoL, pentru care frecvența fundamentală este foarte apropiată de cea a vocii sintetice. Comparația privește singurul sintetizor comercial autohton pentru limba română.

## 1. Introducere

Interesul nostru este, dincolo de cel strict practic, unul legat de aspectele cognitive, anume legat de întrebarea “cât de mult influențează pattern-urile formantice învățate de sistemul nervos calitatea perceptivă a vocii”? În măsura în care inteligibilitatea vocii (a mesajului transmis) este bună sau foarte bună, de ce sunt încă respinse în aplicații sintetizoarele vocalice? Care calități sunt lipsă și cum sunt metrizabile aceste calități? Răspunsurile la aceste întrebări pot avea impact semnificativ de ordin teoretic (îmbunătățirea modelelor cognitive ale audiției), ca și de ordin aplicativ, privind creșterea calității sintetizoarelor comerciale.

## 2. Metodologie

Metoda de comparație privește două aspecte la nivel strict formantic static: valorile absolute ale frecvențelor formațiilor și valorile relative ale raporturilor frecvențelor formațiilor raportate la frecvența fundamentalei pentru acea vocală. Analiza este statică în sensul că suntem interesați pe moment doar de valorile medii, dar nu și de traseele formantice. Pentru comparație, s-a folosit o singură voce umană, aleasă pe criteriul similarității valorii frecvenței fundamentale. Anume, s-a folosit înregistrarea cu codul 55555f din cadrul SRoL ([www.etc.tuiasi.ro/sibm/romanian\\_spoken\\_language/index.htm](http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.htm)). Vocea respectivă este feminină; fișa vorbitorului indică un vorbitor cu educație, fără patologie, voce regională din zona Moldovei, fără accent dialectal evident etc.

Rezultatele prezentate în secțiunea a treia se referă punctual la o voce umană și la una sintetică. O comparație mai judicioasă ar trebui făcută pe mai multe niveluri:

- punctual, vocea sintetică cu câteva voci umane cu frecvențe F0 apropiate, dar cu tonalități diferite (cu conținut diferit de formați superiori);

- statistic, vocea artificială cu media și intervalul de dispersie al valorilor formanților superiori;

- aceleași, pentru cazul dinamic, al traseelor formantice.

Cuvântul analizat este “Aseară”, iar propozițiile comparate sunt “Vine mama” și “Cine a făcut asta?”.

În finalul lucrării sunt menționate sumar și câteva observații privind dinamica formanților, așa cu rezultă dintr-o analiză vizuală preliminară (Praat™ și Wasp™). Aceste rezultate sunt doar calitative.

Pentru analiză, s-a folosit utilitarul Praat™; cuvintele au fost segmentate manual pe foneme și, folosind zona centrală a vocalelor, am determinat cu utilitarul Praat™ valorile medii pe vocale ale frecvențelor fonemelor. Aceste valori au constituit baza analizei raportată aici.

Sintetizorul discutat este unul de tip concatenativ, destinat utilizării de către persoane fără vedere, sau cu vedere redusă. Sintetizorul este descris sumar în cadrul sitului comercial (<http://www.baum.ro/index.php>). Firma producătoare, BAUM Engineering, se prezintă ca având obiectul de activitate “Dezvoltare de Produse pentru Nevăzători și Ambliopi”; firma a realizat demonstrativ sintetizorul TTS Online “Voce sintetică românească profesională Ancutza” v3.6.16., disponibil la adresa <http://www.baum.ro/index.php>. Pentru sintetizor sunt folosite setările default, anume Viteza: 60, Intonația: 60, Format: .WAV. Nu sunt date de către autorii parametrii sau modul de realizare ai sintetizorului respectiv și prin urmare nu putem face aprecieri asupra posibilităților de îmbunătățire a sintezei, la nivel tehnic.

### 3. Rezultatele analizei și discuție

Prezentăm preponderent sub formă grafică, rezultatele comparației la nivel punctual și static. Nu tratăm dinamica formanților.

Rezultatele pentru cuvântul “Aseară” (Teodorescu & Feraru, 2007), pronunțat independent – de exemplu, ca în cazul scurtului dialog “- Când ai fi dorit să mergi la film? – Aseară.”, sunt rezumate în Tabelul 1.

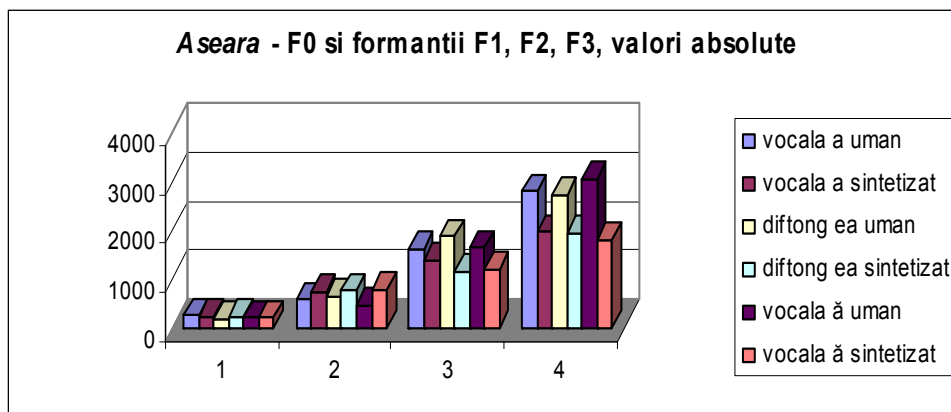
Tabelul 1. Valorile absolute și relative, raportate la F0, ale formanților, în cuvântul “Aseară”. Voce de tip feminin sintetizată cu TSS comparativ cu vocea umană feminină cu indicativul 55555f din SRoL

	Aseară					
	a	a	diftong ea	diftong ea	ă	ă
	uman	sintetizat	uman	sintetizat	uman	sintetizat
<b>F0</b>	254	226	200	228	207	219
<b>F1</b>	593	752	652	769	456	785
<b>F2</b>	1634	1403	1892	1176	1658	1212
<b>F3</b>	2806	1978	2736	1963	3060	1827
<b>F1/F0</b>	2.33	3.33	3.26	3.37	2.20	3.58
<b>F2/F0</b>	6.43	6.21	9.46	5.16	8.01	5.53
<b>F3/F0</b>	11.05	8.75	13.68	8.61	14.78	8.34

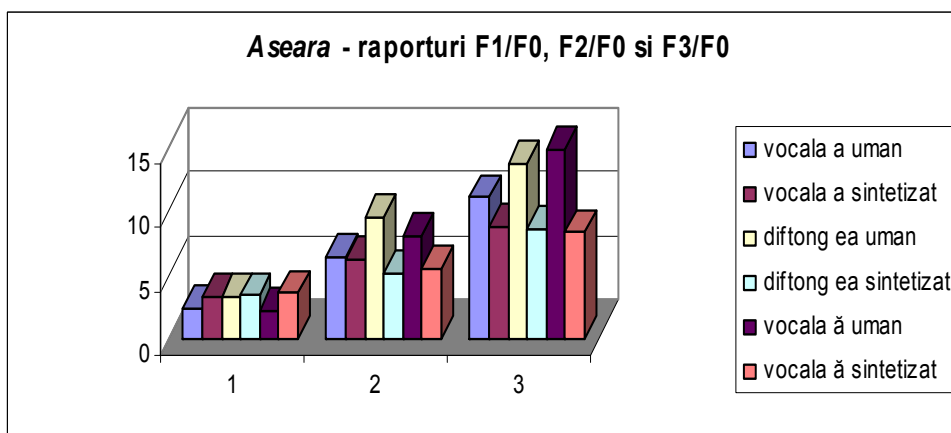
## DE CE NU PLACE VOCEA SINTETIZATĂ? – CÂTEVA ELEMENTE DE COMPARAȚIE CU VOCEA UMANĂ

Din tabelul 1, se observă că valoarea formantului F3 pentru vocea umană (3060Hz, pentru vocala *ă*, din cuvântul *aseară*) este semnificativ mai mare decât valoarea formantului F3 pentru vocea sintetică (1827Hz, pentru vocala *ă*, din cuvântul *aseară*).

În cazul raportului F1/F0, pentru diftongul *ea*, nu există diferențe semnificative (3.26 față de 3.37) în comparație cu raportul F2/F0, care este ca valoare aproape dublu (9.46 față de 5.16) în cazul vocii umane comparativ cu vocea sintetizată.



**Figura 1** Evoluția valorilor absolute ale frecvenței fundamentale și ale formanților pe vocala “a”, diftongul “ea” și vocala “ă” în cuvântul “Aseară”



**Figura 2** Evoluția valorilor relative ale raporturilor F1/F0, F2/F0 și F3/F0 pe vocala “a”, diftongul “ea” și vocala “ă” în cuvântul “Aseară”

În cazul valorilor relative (Teodorescu et al., 2007), se constată mari diferențe la toate vocalele, la nivelul formantului F3, dar și pentru diftongul *ea* și vocala *ă*, între vocea sintetică și cea umană (v. Fig. 2). Se observă că pentru formantului F1 nu există diferențe semnificative, iar pentru formantul F2 diferențe mai semnificative sunt pentru diftongul *ea* și vocala *ă*, între vocea umană și cea sintetizată.

Tabelul 2. Valorile absolute și relative, raportate la F0, ale formațiilor, în propoziția simplă “Vine mama”. Voce de tip feminin sintetizată cu TSS comparativ cu vocea umană feminină cu indicativul 55555f din SRoL

	i	i	e	e	a1	a1	a2	a2
	uman	sintetizat	uman	sintetizat	uman	sintetizat	uman	sintetizat
<b>F0</b>	222	234	242	226	196	225	200	218
<b>F1</b>	348	463	518	661	873	816	921	807
<b>F2</b>	2690	864	2242	1054	1255	1348	1427	1255
<b>F3</b>	3792	2450	3116	2083	2718	1928	3070	1642
<b>F1/F0</b>	1.57	1.98	2.14	2.92	4.45	3.62	4.60	3.70
<b>F2/F0</b>	12.11	3.69	9.26	4.66	6.40	5.99	7.13	5.75
<b>F3/F0</b>	17.08	10.47	12.87	9.21	13.87	8.56	15.35	7.53

În cazul vocalei *i* din cuvântul *vine*, valoarea raportului F2/F0 pentru vocea umană este de patru ori mai mare decât pentru vocea sintetică; în cazul vocalei *e* din cadrul aceluiași cuvânt este de două ori mai mare, iar în cazul ultimei vocale din cuvântul *mama*, valoarea raportului F3/F0 este dublă.

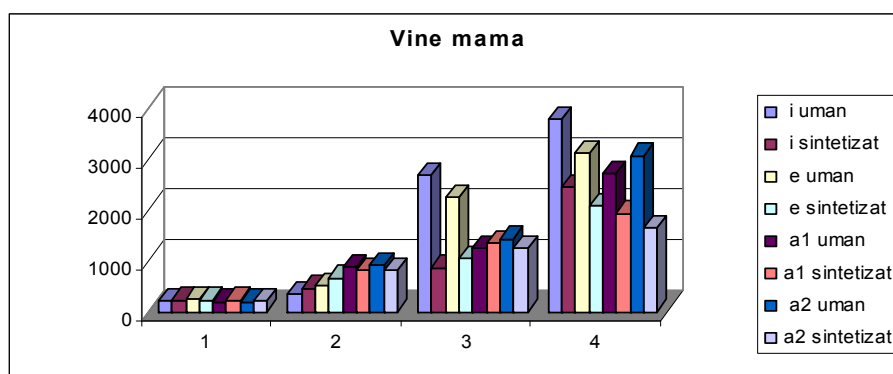


Figura 3 Evoluția valorilor relative ale raporturilor F1/F0, F2/F0 și F3/F0 pe vocala “i”, “e”, “a1” și vocala “a2” în propoziția “Vine mama”

Se observă diferențe semnificative mai mari în cazul formațiilor F2 și F3, pentru vocea umană comparativ cu cea sintetizată pentru următoarele vocale din cadrul propoziției “Vine mama”: *i*, *e* și ultimul *a* din cadrul cuvântului *mama*.

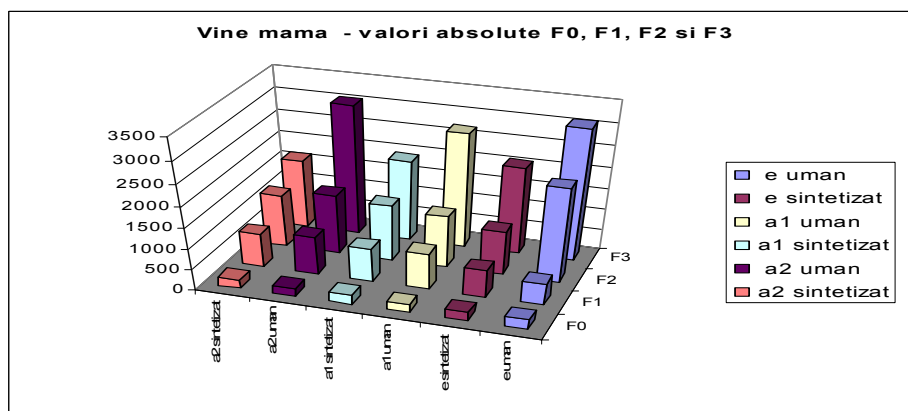
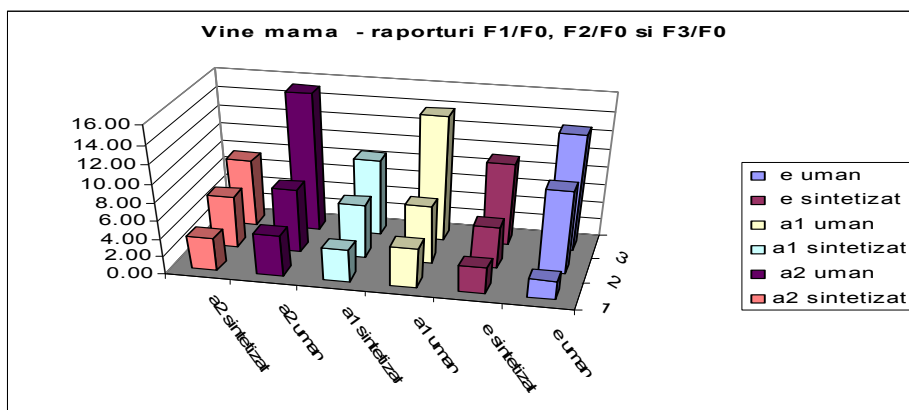


Figura 4 Evoluția valorilor absolute ale frecvenței fundamentale și ale formațiilor pe vocala “i”, “e”, “a1” și vocala “a2” în propoziția “Vine mama”



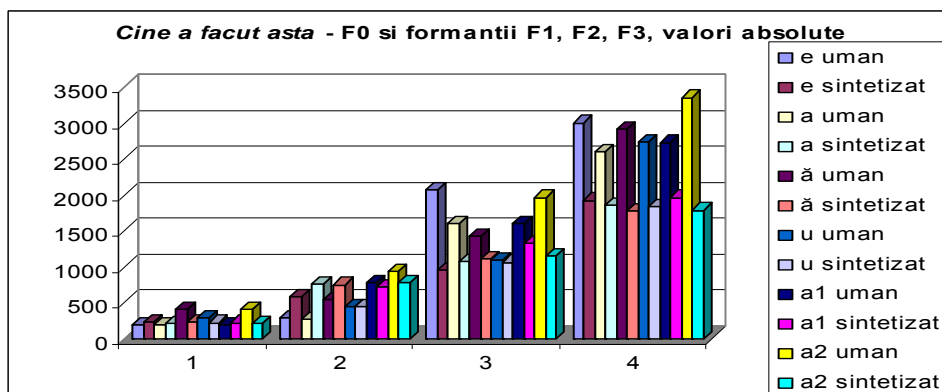
DE CE NU PLACE VOCEA SINTETIZATĂ? – CÂTEVA ELEMENTE DE COMPARAȚIE  
CU VOCEA UMANĂ

În figura 4, se observă pentru vocala *a* (primul *a* și ultimul *a* din cuvântul *mama*) la vocea umană, că valoarea formantului F3 este mult mai mare comparativ cu vocea sintetizată. Referitor la F0, formantul F1 și F2, diferențele sunt ne semnificative pentru toate vocalele din propoziția “Vine mama”.



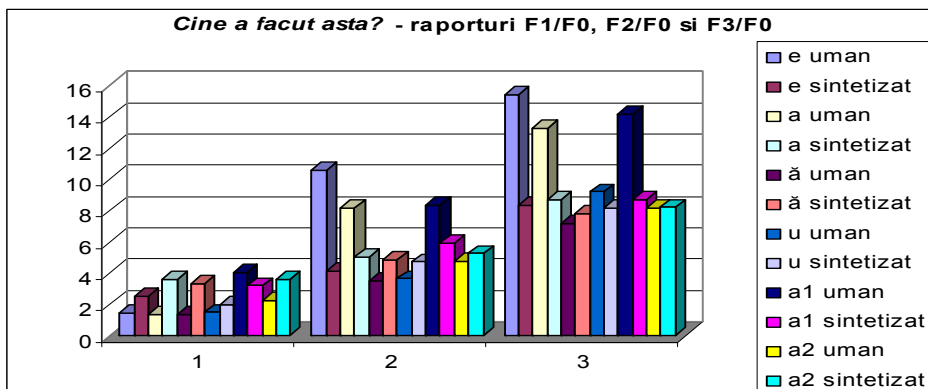
**Figura 5** Evoluția valorilor relative ale raporturilor F1/F0, F2/F0 și F3/F0 pe vocala “e”, “a1” și vocala “a2” în propoziția “Vine mama”

Similar, pentru propoziția interogativă “Cine a făcut asta?” – pe care de altfel sintetizorul discutat nu o poate produce la forma interogativă, rezultatele comparative sunt prezentate în Fig. 6 și Fig. 7.



**Figura 6** Evoluția valorilor absolute ale frecvenței fundamentale și ale formanților pe vocala “e”, “a”, “ă”, “u”, “a1” și vocala “a2” în propoziția “Cine a făcut asta”

Graficele din Fig. 6, 7 corespund valorilor din Tabelul 3. Valorile F0 sunt semnificativ mari (duble) pentru vocala *ă* (din cuvântul *făcută*) și vocala *a* (ultimul *a* din cuvântul *asta*) din propoziția “Cine a făcut asta” în comparație cu vocea sintetizată. În cazul formantului F1, valorile sunt aproape duble în cazul vocii sintetizate comparativ cu vocea umană pentru vocalele *e* (din cuvântul *cine*) și *a* din cadrul propoziției “Cine a făcut asta”. Valorile F3 sunt ca valoare mari în cazul vocii umane comparativ cu vocea sintetică, iar raporturile F2/F0 și F3/F0, în cazul vocii umane, pentru vocala *e* (din cuvântul *cine*) sunt aproape duble ca valoare comparativ cu vocea sintetizată.

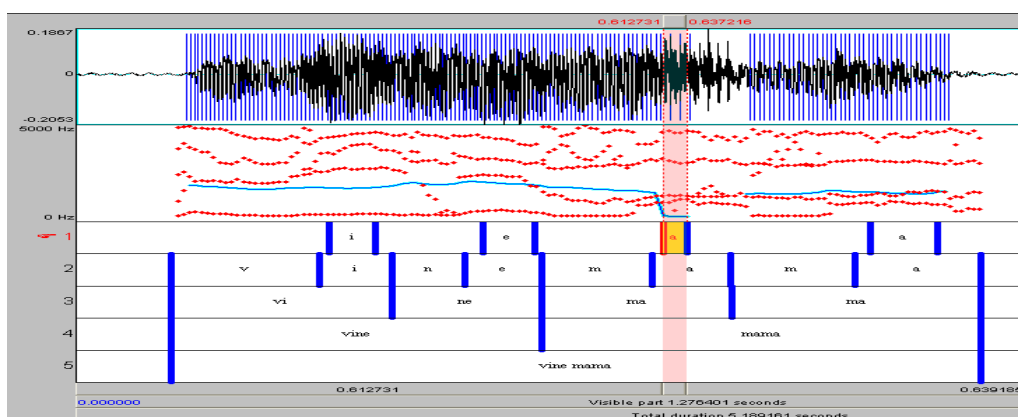


**Figura 7** Evoluția valorilor relative ale raporturilor F1/F0, F2/F0 și F3/F0 pe vocala “e”, “a”, “ă”, “u”, “a1” și vocala “a2” în propoziția “Cine a făcut asta”

Tabelul 3. Valorile absolute și relative, raportate la F0, ale formanților, în propoziția simplă “Cine a făcut asta”. Voce de tip feminin sintetizată cu TSS comparativ cu vocea umană feminină cu indicativul 55555f din SRoL

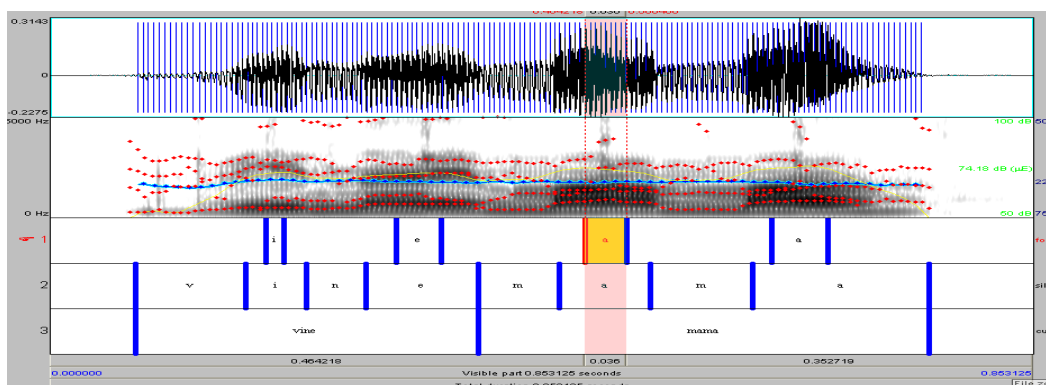
	e		a		ă		u		a1		a2	
	Um	Sint	Um	Sint	Um	Sint	Um	Sint	Um	Sint	Um	Sint
<b>F0</b>	195	231	196	213	406	227	298	225	192	225	411	218
<b>F1</b>	290	591	268	765	546	753	448	452	776	723	931	783
<b>F2</b>	2068	951	1609	1076	1424	1110	1090	1064	1604	1332	1966	1148
<b>F3</b>	3006	1925	2605	1864	2916	1781	2744	1840	2729	1968	3353	1791
<b>F1/F0</b>	1.49	2.56	1.37	3.59	1.34	3.32	1.50	2.01	4.04	3.21	2.27	3.59
<b>F2/F0</b>	10.61	4.12	8.21	5.05	3.51	4.89	3.66	4.73	8.35	5.92	4.78	5.27
<b>F3/F0</b>	15.42	8.33	13.29	8.75	7.18	7.85	9.21	8.18	14.21	8.75	8.16	8.22

Privitor la dinamica formanților (traseele formantice), sunt prezentate în Fig. 8, 9, 10 și 11 imagini obținute cu utilitarul Praat™ pentru propozițiile “Vine mama” și “Cine a făcut asta” (Feraru & Teodorescu, 2008), pentru vocea umană și pentru cea sintetizată, iar în Fig. 12 și 13 imagini similare obținute cu utilitarul Wasp™, pentru propoziția “Vine mama.”



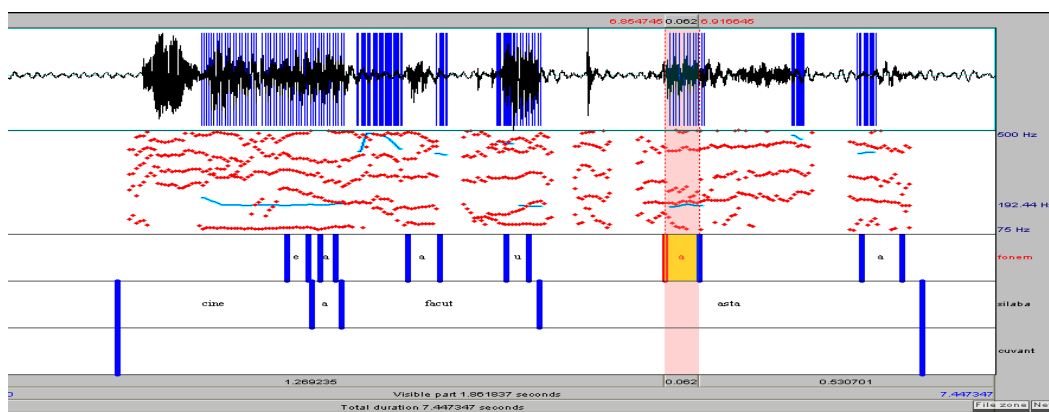
**Figura 8** Exemplu de adnotare manuală folosind utilitarului Praat™; voce umană – propoziția “Vine mama”

DE CE NU PLACE VOCEA SINTETIZATĂ? – CÂTEVA ELEMENTE DE COMPARAȚIE  
CU VOCEA UMANĂ

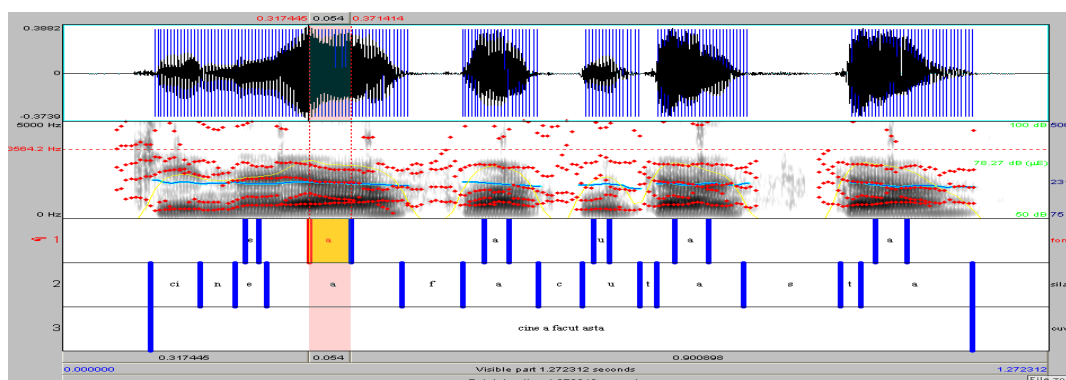


**Figura 9** Exemplu de adnotare manuală folosind utilitarului Praat™; voce sintetizată – propoziția “Vine mama”

În figura 8, în exemplul de adnotare pentru propoziția “Vine mama”, voce umană se observă o bună demarcare pentru formanții F3, F4 și o variabilitate bogată a traseelor, indicând o modulare prozodică bună, ne-monotonă.



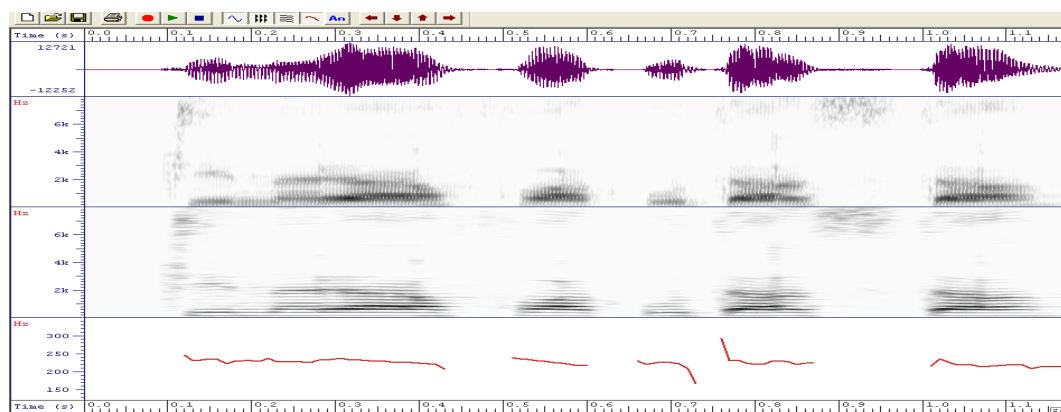
**Figura 10** Exemplu de adnotare manuală folosind utilitarului Praat™; voce umană – propoziția “Cine a făcut asta”



**Figura 11** Exemplu de adnotare manuală folosind utilitarului Praat™; voce sintetizată – propoziția “Cine a făcut asta”

Din compararea traseelor formanților superiori în Fig. 8-11, se constată că evoluția acestor formanți este discontinuă în cazul vocii sintetice, în timp ce în cazul vocii naturale traseele urmează curbe cu variație relativ lină. Chiar și atunci când traseele lui

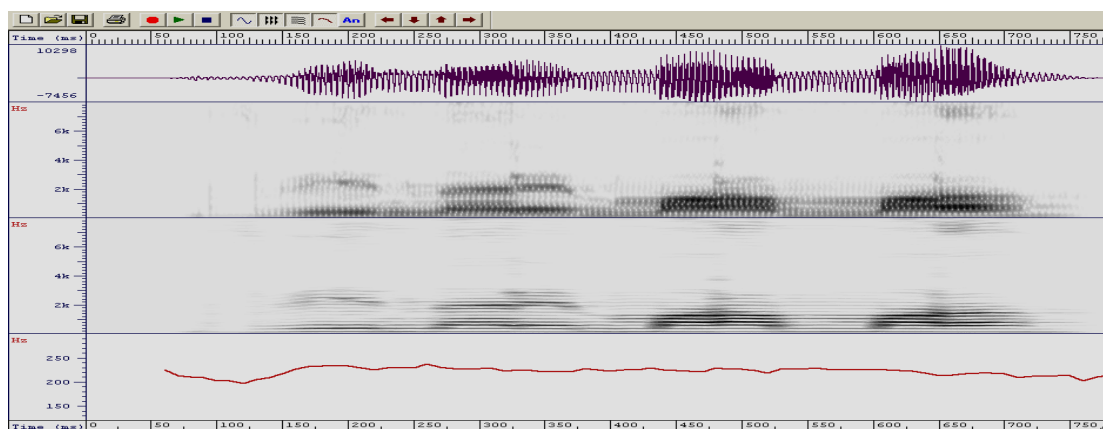
F0 și F1 sunt asemănătoare pentru vocea sintetică și cea naturală, traseele formanților superiori pentru cele două voci sunt total diferite. Considerăm că aceste diferențe explică cel puțin parțial de ce vocea sintetică este “neplăcută”, induce senzația de “nenatural”.



**Figura 12** Traseul frecvenței fundamentale (linia roșie) folosind utilitarul Wasp<sup>TM</sup>; voce sintetizată - propoziție “Cine a făcut asta”

În figura 12 se observă pentru vocea sintetică, discontinuități pentru F0, F1, F2 precum și variații de tip „ruptură”, indicate de prima săgeată. În figura de mai jos (Fig.13), pentru vocea sintetizată se observă că formanții superiori (F3, F4) sunt slab demarcați, aproape inexistenți. Benzile de frecvență sunt înguste, foarte bine demarcate.

Se observă în Fig. 12 comparativ cu Fig. 13 că traseul F0 este, în cazul vocii sintetice, mai monoton, iar structura de formanți este mult mai grosieră și mai săracăcioasă. Traseul F0 în cazul vocii umane prezintă variații mai mari (probabil datorită accentuării corecte pe cuvinte) comparativ cu traseul F0 în cazul vocii sintetizate, ultimul fiind mai plat, mai uniform (ceea ce poate indica absența unui bloc de stabilire a modului corect în care sunt plasate accentele pe cuvinte în cadrul propozițiilor).



**Figura 13** Traseul frecvenței fundamentale (linia roșie) folosind utilitarul Wasp<sup>TM</sup>; voce sintetizată - propoziție “Vine mama”

Spectru vocii umane din Fig. 12 se observă că are un domeniu de variație al frecvenței mult mai larg și mai difuz; este un spectru bogat în informație comparativ cu spectru

vocii sintetizate, care este într-un domeniu de variație al frecvenței mai scăzut, și sărac în informație.

Suplimentar, se constată salturi bruște (“rupturi”) și “lărgiri” bruște ale valorilor respectiv benzilor formanților (vezi locurile indicate de săgeți), salturi care fac vocea neplăcută.

Rezultatele menționate se confirmă și pentru alte voci umane. Precizări suplimentare vor fi date într-o lucrare viitoare.

#### **4. Concluzii și direcții viitoare**

Analiza sumară realizată indică cel puțin o serie de diferențe majore între formanții superiori la vocea sintetică față de cea umană. Deși – se știe – inteligibilitatea este dată doar de formanții inferiori (F1, F2), calitatea pronunției, nuanțele și bogăția vocii sunt datorate în mare măsură de formanții superiori, care sunt sensibil mai puțin “corect” produși de sintetizor.

Sub nici o formă nu sugerăm că metoda concatenativă nu poate da rezultate excelente în simularea vocii naturale – scopul nostru a fost să arătăm o metodă de verificare cantitativă a “naturaletii” și să sugerăm un “benchmark” în verificări. Creșterea naturaleții necesită, foarte probabil, un număr mai mare de fragmente (elemente) în memorie, astfel încât selecția să se facă și ținând cont de traseele pentru formanții superiori.

În această lucrare nu ne referim la sintetizatoarele de cercetare, ci numai la cele din domeniul comercial; ca urmare ne-am restrâns doar la singurul produs comercial autohton, realizat pentru limba română. În acest context, amintim că numeroase colective de cercetare s-au preocupat de vocea sintetică și de îmbunătățirea calității sintezei concatenative în limba română precum, colectivele de la Universitatea Tehnică din București, prof. dr. D. Burileanu și colaboratorii săi, de la Universitatea Tehnică Cluj-Napoca, prof. dr. G. Todorean, de la Academia Tehnică Militară, etc.

Această lucrare reprezintă o raportare preliminară a unor analize începute recent. Ne propunem în viitor să facem un număr mai mare de comparații pentru limba română – folosind diverse setări pentru F0 pentru sintetizorul respectiv și folosind mai multe voci naturale corespunzător alese. De asemenea, este necesar să realizăm analize la nivel dinamic, pentru trasee ale formanților (prozodie completă).

**Mulțumiri.** Primul autor mulțumește colegului V. Apopei pentru menționarea sitului <http://www.baum.ro/index.php>. Mulțumim de asemenea recenzorilor anonimi pentru observațiile pertinente făcute.

Analiza prezentată rezumativ în acest scurt raport a fost sprijinită de către Academia Română, Secția de Știință și Tehnologia Informației, în cadrul temei interne „Procese de cogniție, limbaj și calcul”, subtema 2.3 “SRoL: optimizare modele statistice disponibile on-line (folosind date din urma analizei formanților). Realizarea, în cadrul sitului, a facilității de recunoaștere (specificare) on-line a stărilor emoționale”.

### Referințe bibliografice

- H.N. Teodorescu, M. Feraru, D. Trandabăț, M. Zbancioc, R. Luca, A. Verbuță, M. Hnatiuc, R. Ganea, O. Voroneanu, L. Pistol, “*Proiectul Sunetele Limbii Române*”, [www.etc.tuiasi.ro/sibm/romanian\\_spoken\\_language/index.htm](http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.htm), [http://www.etc.tuiasi.ro/sibm/romanian\\_spoken\\_language/ro/fisa\\_2.htm](http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/ro/fisa_2.htm)
- BAUM Engineering, *TTS Online, Voce sintetică românească profesională* Ancutza v3.6.16., [http://www.baum.ro/index.php?language=ro&pagina=despre\\_noi](http://www.baum.ro/index.php?language=ro&pagina=despre_noi), <http://www.baum.ro/index.php?language=ro&pagina=ttsonline>
- P. Boersma, D. Weenink, Institute of Phonetic Science, University of Amsterdam, *Praat: doing phonetics by computer*, [www.praat.org](http://www.praat.org)
- WasP – the Wind Atlas Analysis and Application Program*, [www.wasp.dk](http://www.wasp.dk)
- H.N. Teodorescu, M. Feraru, (2007) A study on Speech with Manifest Emotions, *10th International Conference on Text, Speech and Dialogue*, TSD 2007, Pilsen, Czech Republic, Lecture Notes in Computer Science, Springer Verlag, vol. 4629/2007, ISBN 978-3-540-74627-0, p. 254-262
- H.N. Teodorescu, M. Feraru, D. Trandabăț, (2007) Studies on the Prosody of the Romanian Language: The Emotional Prosody and the Prosody of Double-Subject Sentences. In Corneliu Burileanu and H.N. Teodorescu (Eds.), *Advances in Spoken Language Technology*, The Publishing House of the Romanian Academy, ISBN 978-973-27-1516-1, p.171-182
- M. Feraru, H.N. Teodorescu, (2008) *Speech Corpus for the Romanian Language: the Emotional Speech Section*, Inventica 2008, Ed. Performantica, Iași, România ISBN 978-973-730-491-9, p. 261-273

# BAZA DE DATE ÎN LIMBA ROMÂNĂ PENTRU RECUNOAȘTEREA VORBIRII SPONTANE

DIANA HANES<sup>2</sup>, CRISTINA PETREA<sup>2</sup>, ANDI BUZO<sup>2</sup>, VLADIMIR POPESCU<sup>1,2</sup>,  
CORNELIU BURILEANU<sup>2</sup>

<sup>1</sup> *Laboratoire d'Informatique de Grenoble, Grenoble INP - France*

<sup>2</sup> *Universitatea Politehnică, Facultatea de Electronică, Telecomunicații și Tehnologia  
Informației, București – România;*

*[cburileanu@mesnet.pub.ro](mailto:cburileanu@mesnet.pub.ro) [vladimir.popescu@imag.fr](mailto:vladimir.popescu@imag.fr)*

## Rezumat

Recunoașterea vorbirii spontane reprezintă un domeniu mai puțin cercetat, comparativ cu recunoașterea vorbirii continue în general. În această lucrare sunt prezentate o serie de rezultate privind proiectarea, achiziția și adnotarea unei baze de date de vorbire spontană în limba română. Este propusă o metodologie pentru achiziția datelor, evidențiind etapele importante: definirea unui lexic și a unui dicționar de unități fonetice, achiziția propriu-zisă a semnalului vocal, precum și adnotarea materialului vocal în unități lingvistice.

## 1. Introducere

Interacțiunea personalizată între subiectul uman și calculator constituie o provocare de primă importanță la ora actuală, în contextul în care serviciile și aplicațiile informatice devin din ce în ce mai mult centrate pe utilizator.

În limbile de mare circulație (engleză, franceză) există sisteme complete de dialog persoană-calculator. În alte limbi (limba română) realizarea sistemelor de dialog reprezintă un demers de lungă durată.

Recunoașterea vorbirii poate fi privită drept un proces de recunoaștere a formelor, iar acest lucru se poate realiza fie pe bază de reguli, fie prin metode statistice (Russell & Norvig, 2003).

Un obstacol major în calea realizării unei recunoașteri fiabile constă în variabilitatea semnalului vocal care provine din: variabilitate lingvistică, variabilitatea vorbitorilor, variabilitatea canalului. Un sistem de recunoaștere a vorbirii spontane trebuie să ia în considerare: independența de vorbitor; dimensiunea vocabularului; caracterul continuu al vorbirii; spontaneitatea vorbirii (enunțurile rostite de către utilizator sunt de regulă spontane, neplanificate, caracterizate de disfluențe, ezitări, interjecții).

Recunoașterea vorbirii presupune găsirea unei secvențe de cuvinte, folosind un ansamblu de modele determinate, achiziționate într-o fază anterioară de antrenare, și potrivirea acestor modele cu semnalul de vorbire incident. Sisteme bazate pe abordări statistice sunt disponibile atât în comunitățile academice (sistemul SPHINX de la Universitatea Carnegie Mellon, ansamblul de utilitare HTK – „Hidden Markov Modelling Toolkit” - Universitatea Cambridge, sistemul RAPHAEL - Laboratoire d'Informatique de Grenoble) cât și în domeniul comercial (sistemele produse de Nuance, Dragon, Microsoft în Statele Unite – (Huang & Acero, 2001)).

## **2. Arhitectura sistemului de recunoaștere în limba română**

Un sistem tipic de recunoaștere a vorbirii funcționează în două moduri de lucru: antrenare (prin crearea cunoștințelor necesare funcționării sale, a modelelor acustice și lingvistice utilizate) și recunoaștere (presupune utilizarea resurselor create la antrenare pentru conversia enunțului provenit de la subiectul uman într-o secvență de cuvinte).

Figurile 1 și 2 prezintă propunerea pentru arhitectura fazelor de antrenare și de testare care stau la baza sistemului de recunoaștere de vorbire spontană. Întrucât scopul acestei lucrări este reprezentat de descrierea procedurii de construire a bazei de date, nu se va insista aici asupra arhitecturii concepute pentru recunoașterea vorbirii.

La finalul prelucrărilor, rezultatul recunoașterii este reprezentat de un număr de șiruri de cuvinte alternative pentru un enunț rostit. Alegerea alternativei celei mai pertinente în raport cu contextul cade în sarcina altor componente ale sistemului de dialog. Sistemul de recunoaștere propus în această lucrare se bazează pe modele Markov ascunse (MMA-uri) antrenate pentru fiecare trifonem.

Sistemul de recunoaștere a cărui dezvoltare se dorește este bazat pe pachetul de utilitare HTK versiunea 3.0. Mare parte din funcționalitatea HTK este încapsulată într-un set de biblioteci statice, realizate în limbajul C, care asigură faptul că fiecare aplicație se interfațează cu celelalte într-o manieră puternic controlată și reproductibilă. De regulă, fiecare aplicație HTK implementează un nivel de analiză a vorbirii, însă uneori este necesar ca mai multe aplicații să contribuie la realizarea unui nivel de analiză.

Utilitarele HTK, versiunea 3.0, considerate relevante pentru realizarea sistemului de recunoaștere sunt: HLEd, HInit, HCompV, HERest, HCopy, HVite, HResults (Evermann, 2005).

### **2.1. Etapa de antrenare**

În faza de antrenare s-au folosit ca intrare fișiere audio înregistrate. Pentru a construi dicționarul fonetic s-au folosit convențiile SAMPA (Munteanu, 2006), (Burileanu, 2002) în scopul realizării transcrierii fonetice. Antrenarea s-a realizat la nivel de trifonem.

Un pas important în faza de antrenare a constat în definirea prototipurilor pentru modelele Markov. Fiecărui trifonem îi corespunde un model Markov cu 6 stări. Ca observații pentru modelele Markov s-au folosit două mixturi gaussiene (cu parametrii mediile și varianțele). În prima etapă a fost ales un prototip pentru toate trifonemele.

Modelele Markov ascunse au fost inițializate folosind o matrice implicită de observații și tranziții. Parametrii de inițializare au fost calculați folosind toate fișierele de semnal. Acești pași au fost realizați prin intermediul utilităților HInit și HCompV.

Etapa de antrenare s-a făcut folosind metoda „embedded” (înglobată) (Evermann, 2005). Această metodă constă în faptul că etichetarea nu se face la nivel de trifonem ci la nivel de fișier. Utilitarul HERest construiește un MMA care conține toate MMA-urile din interiorul unui fișier, în ordinea specificată în fișierul de etichetare. Antrenarea MMA-urilor pentru fiecare trifonem s-a realizat folosind parametrii MFCC și trifonemele etichetate. S-a realizat o aliniere iterativă Viterbi pentru modelele Markov



În scopul obținerii probabilității maxime pentru ca un anumit MMA să reprezinte trifonemul corespunzător.

### Parametrizarea semnalului vocal

În cadrul fazei de antrenare o primă subetapă parcursă a fost aceea de parametrizare a semnalului vocal. S-a utilizat parametrizarea cu 12 coeficienți MFCC, energia și derivatele corespunzătoare. În total au fost utilizați 26 parametri. Lungimea perioadei a fost considerată 10ms. Transformata FFT a folosit fereastră Hamming de 20ms iar semnalului i s-a aplicat un filtru de preaccentuare de ordinul întâi cu coeficientul 0.97. Bancul de filtre are 26 canale și 12 coeficienți MFCC la ieșire.

### Generarea modelelor Markov prototip

S-au considerat modele Markov ascunde prototip caracterizate printr-o topologie Bakis cu șase stări, cu tranziții de la stânga la dreapta, dintre care starea inițială și starea finală sunt neemisive. Într-o prima fază s-a optat pentru modelarea probabilităților de ieșire cu două mixturi gaussiene per stare emisivă.

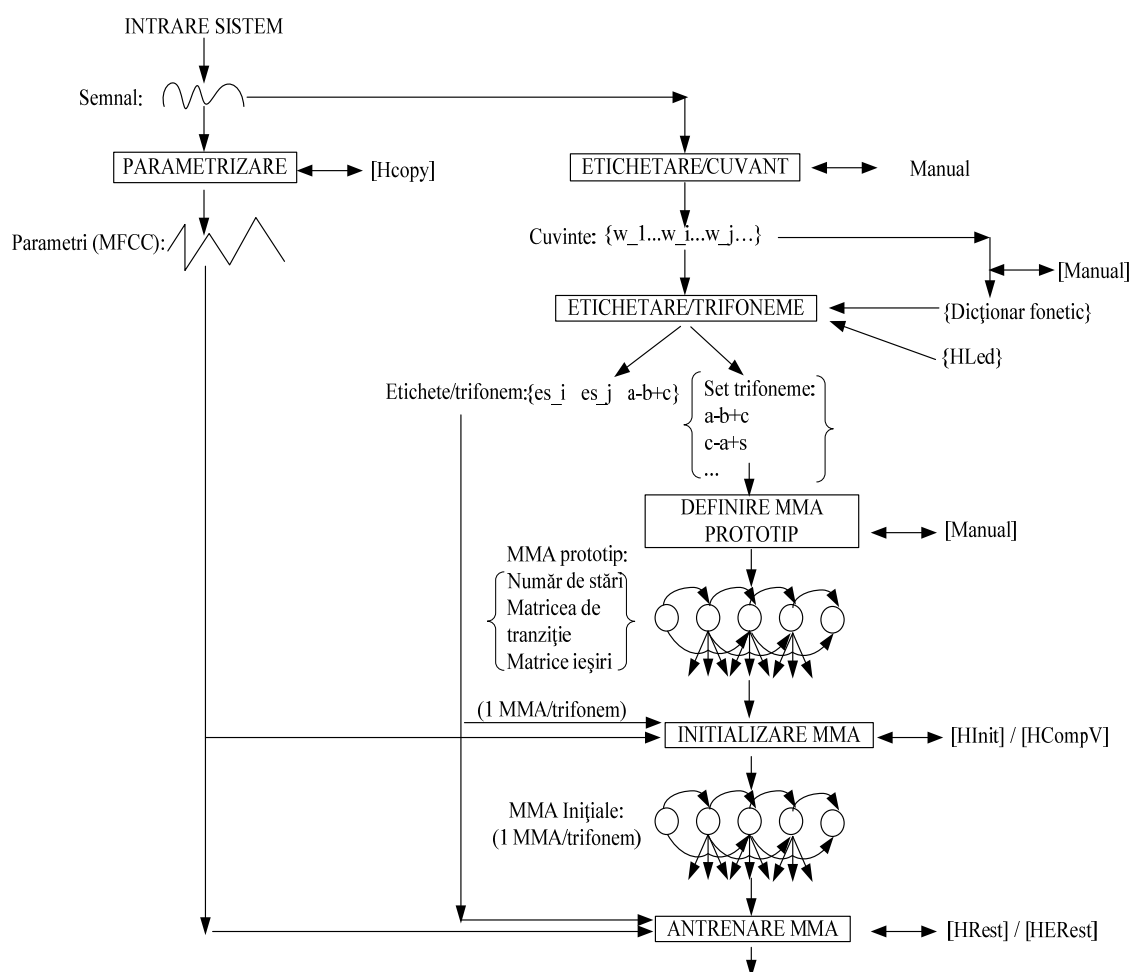


Figura 1: Arhitectura sistemului pentru antrenarea modelelor Markov ascunde

### **Inițializarea globală a modelelor Markov ascunse prototip**

La inițializarea globală a modelelor Markov ascunse, mediile mixturilor devin egale cu media globală a vectorilor acustici și varianțele mixturilor devin egale cu varianța globală a vectorilor acustici.

Pentru inițializarea globală s-a folosit utilitarul HCompV (Evermann, 2005). S-au folosit ca intrări: setul de vectori cepstrali filtrați Mel obținut cu HCopy; fișierele cu etichetele și setul de MMA-uri prototip. Rezultatul rulării îl reprezintă un set de MMA-uri inițializate grupate într-un singur fișier „mmf”.

S-a constatat că rularea utilitarului HCompV a durat pentru un singur prototip MMA mai mult de 10s. Numărul de MMA-uri prototip este egal cu numărul de trifoneme obținut pentru baza de date creată și anume 5095 trifoneme. Timpul de rulare a inițializării globale pentru toate trifonemele a fost de aproximativ 50950s, adică 850 minute, ceea ce înseamnă estimativ 14 ore.

### **Antrenarea propriu-zisă**

Pentru antrenarea Baum-Welch „embedded” s-a folosit utilitarul HERest. S-au utilizat ca intrări: fișierele „mfc” care conțin parametrii MFCC, fișierele cu etichete la nivel de trifonem grupate într-un singur fișier „mlf”, fișierele MMA rezultate în urma rulării lui HCompV grupate într-un singur fișier „mmf” și o listă cu toate trifonemele. La ieșire a rezultat un fișier cu extensia „mmf” care conține toate MMA-urile antrenate.

### **2.2. Etapa de testare**

În faza de testare se folosesc rezultatele obținute în cadrul etapei de antrenare (modele acustice antrenate), resurse folosite la antrenare (dicționarul fonetic) precum și alte resurse (gramatica). Faza de testare împreună cu evaluarea rezultatelor reprezintă faza terminală a procesului de recunoaștere.

Fișierele audio folosite la intrarea sistemului au fost parametrizate și au fost extrași parametrii MFCC folosind utilitarul HCopy (Evermann, 2005). Parametrii MFCC constituie intrare pentru decodarea trifonemelor.

Decodarea parametrilor acustici se realizează folosind MMA-urile antrenate la nivel de trifonem. Secvența de trifoneme obținută va fi transformată într-o secvența de cuvinte folosind gramatica cu un număr finit de stări.

În faza de testare se va parcurge procesul de „definire gramatici”. Astfel se va defini o gramatică-bucă în care se plasează toate cuvintele, care ar putea urma unul după altul cu egale șanse de apariție. Se folosește utilitarul HParse (Evermann, 2005). Decodarea utilizează HVite, HParse și apoi HResults pentru evaluare.

Utilitarul HVite are ca intrări un set de .mfc-uri care constituie semnalul de test, setul de MMA-uri, gramatica constituită din cuvintele din dicționarul fonetic, dicționarul fonetic și lista de modele Markov ascunse (ca în cazul HERest).

Pașii pentru rularea testării sunt: pregătirea datelor (plasarea într-un același director a dicționarului fonetic, a fișierului cu MMA-urile antrenate) și apoi construirea manuală a gramaticii "grammar-orig" folosind lista de cuvinte din dicționar. Între fiecare două

cuvinte se introduce separatorul "|". Urmează rularea propriu-zisă și evaluarea rezultatelor folosind HResults.

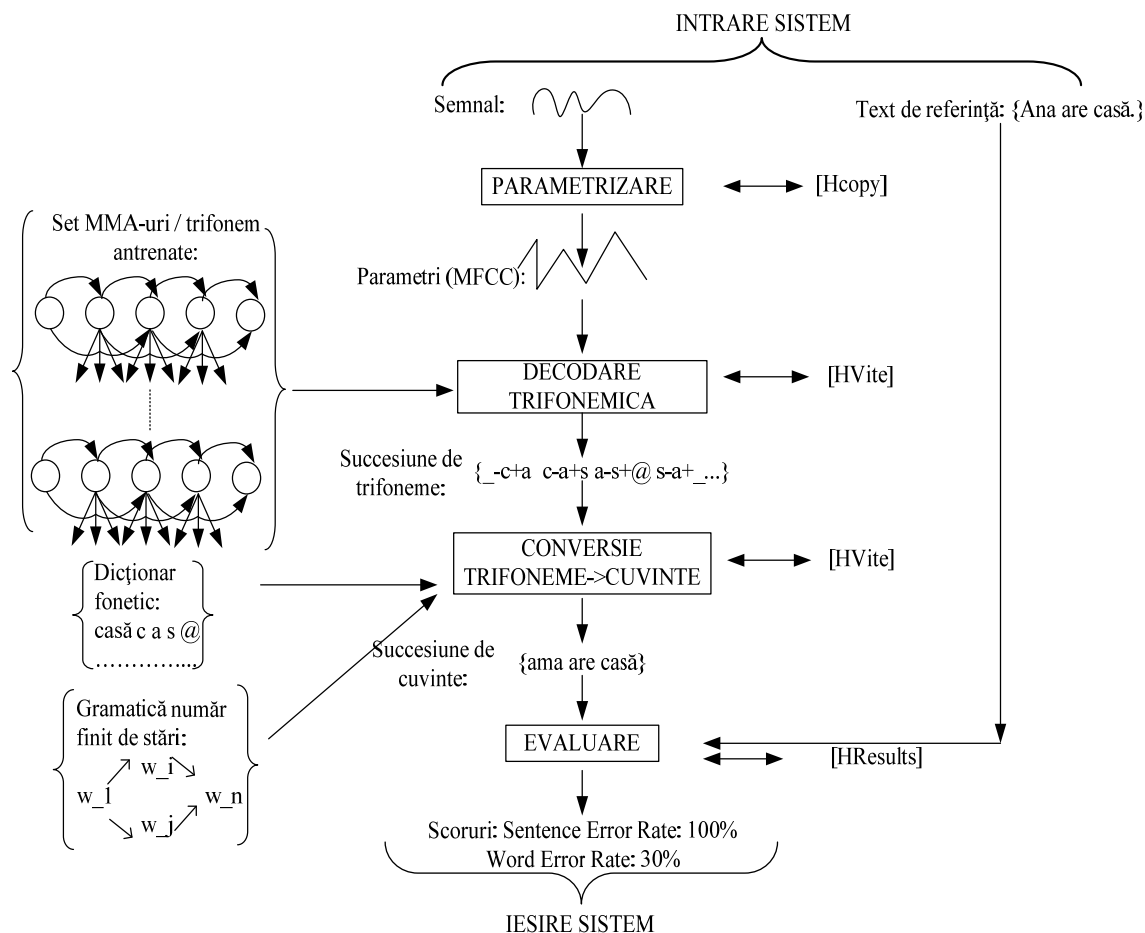


Figura 2: Arhitectura sistemului pentru decodarea vorbirii utilizând modele Markov ascunse

Semnalul de referință este comparat cu rezultatele obținute pentru a determina performanțele utilitarului de recunoaștere construit. Analiza performanțelor acestui utilitar se va face folosind: Sentence Error Rate și Word Error Rate.

### 3. Construirea bazei de date în limba română pentru recunoașterea vorbirii spontane

#### 3.1. Caracteristici

Vocabularul folosit în recunoașterea de vorbire spontană trebuie să fie pe cât posibil de cuprinzător. De la început s-a urmărit construirea vocabularului într-o manieră scalabilă făcând posibilă adăugarea de cuvinte noi ulterior acestei faze. În faza de creare a bazei de date au fost realizate convenții care să fie refolosite în etapele ulterioare de lărgire a bazei de date.

### 3.2. Probleme specifice realizării bazelor de date utilizate în recunoașterea de vorbire spontană

Construirea bazelor de date pentru recunoașterea vorbirii spontane se caracterizează prin elemente specifice. Recunoașterea de vorbire poate fi considerată un proces de recunoaștere a formei iar acest lucru se obține pe baza unor reguli sau metode statistice (Russell & Norvig, 2003). Cea din urmă variantă este cea preferată în acest moment datorită rezultatelor bune obținute cu costuri de producție acceptabile. Metodele statistice presupun folosirea unor date de intrare în procesul de antrenare astfel încât sistemul generează informație pe care o folosește în etapele ulterioare.

Următoarele caracteristici sunt relevante pentru recunoașterea vorbirii: tipul sistemului (dependent de vorbitor sau independent de vorbitor), dimensiunea vocabularului (Peinado & Segura, 2006) (vocabular mic: 10-100 cuvinte, vocabular mediu: 100-1000 cuvinte, vocabular mare: 10.000-100.000 cuvinte).

Opțiunile disponibile pentru realizarea bazei de date sunt: înregistrările directe, fișiere audio din cadrul programelor de televiziune sau radio care au fost difuzate pe Internet, utilizarea de fișiere audio înregistrate direct de la radio sau de la televizor.

Parametrii necesari în construirea unei baze de semnal vocal pentru recunoașterea de vorbire spontană sunt: fișiere audio cu semnal vocal; extragerea caracteristicilor semnalului vocal din cadrul înregistrărilor audio; etichetarea fișierelor audio; parametrii acustici ai fișierelor audio (coeficienții de predicție liniară (LPC), coeficienții cepstrali).

Tabelul 1 indică principalele caracteristici pentru baza de date în limba română.

Tabel 1: Caracteristicile bazei de date realizate.

Proprietate	Valoare
Procedura colectare fișiere audio	Înregistrări preluate de pe Internet ale unor emisiuni românești
Limba folosit	Limba română, vorbire orală/spontană
Durată înregistrări	~ 4 ore
Vorbitori	12 (4 bărbați, 8 femei)
Sesiuni per vorbitor	3-20
Număr total de cuvinte	37604
Număr de cuvinte unicat	8068
Frecvența de eșantionare semnal vocal	8kHz

Pe baza elementelor menționate în tabel au fost evidențiate următoarele aspecte: segmentarea semnalului vocal – este de preferat ca lungimea fișierelor vocale să fie de 60s; etichetarea semnalului – etichetarea poate fi realizată la nivel de cuvânt (proces realizat manual și consumator de timp) sau la nivel de fonem/trifonem (proces semi-automat care se bazează pe etichetarea inițială manuală, este lipsit de stabilitate); parametrizarea semnalului vocal – pentru aceasta anumite criterii trebuie îndeplinite: maximizarea dispersiei inter-fonem și minimizarea dispersiei intra-fonem.

### 3.3. Culegerea și structurarea datelor

Pentru a construi baza de date de vorbire spontană în limba română, au fost utilizate înregistrări ce conțin știri, povești, show-uri de televiziune, discuții medicale, financiare, previziuni meteo și alte tipuri de informații toate fiind transmise pe Internet sau la radio.

Baza de înregistrări audio conține fișiere audio care provin de la 12 vorbitori cu pregătiri în diverse domenii; sunt persoane care au diferite stiluri de viață, experiență, obiceiuri. Vorbirea este fluentă și înregistrările audio conțin diferite tipuri de informații. Pentru fiecare vorbitor există în baza de înregistrări audio între 5 și 38 de fișiere audio.

### 3.4. Construirea bazei de date

#### Transcrierea din fișiere audio în fișiere text

Fișierele audio au fost prelucrate înainte de a fi utilizate efectiv la crearea bazei de date. Înregistrările audio au fost divizate pe vorbitori; au fost eliminate zonele în care vocile s-au suprapus; fiecare fișier audio a fost divizat în fișiere audio de durată 60s.

S-a ascultat fiecare fișier audio, s-a extras mulțimea de cuvinte rostite ținând cont de diacriticele specifice limbii române: ă -> @; â -> i\_; î -> i\_; ș -> S; ț -> ts. Cratimele au fost suprimate, cuvintele sub formă de acronim au fost fonetizate astfel: "bcr" devine "becere"; "bvb" devine "bevebe" etc. În cazul onomatopeelor prin convenție s-a dublat vocala mai lungă astfel "breee..." a devenit "bree".

Întrucât se dorește construirea unei baze de cuvinte pentru vorbire spontană, au fost incluse în fișierele .txt cuvintele care sunt pronunțate incomplet, incorect și eventualele bâlbe. Acestea sunt entități recurente ce trebuie considerate ca atare.

Din mulțimea de cuvinte rezultată au fost generate două fișiere: un fișier care include totalitatea cuvintelor din înregistrările audio prelucrate (37604 apariții cuvinte) și respectiv al doilea fișier care conține doar o reprezentare a cuvintelor din lista generată inițial (8068 cuvinte).

S-a constatat că în cadrul vocabularului construit 12147 de cuvinte au peste 100 de apariții. În figura 3 sunt reprezentate cuvintele care au mai mult de 250 de apariții în cadrul vocabularului realizat până în acest moment.

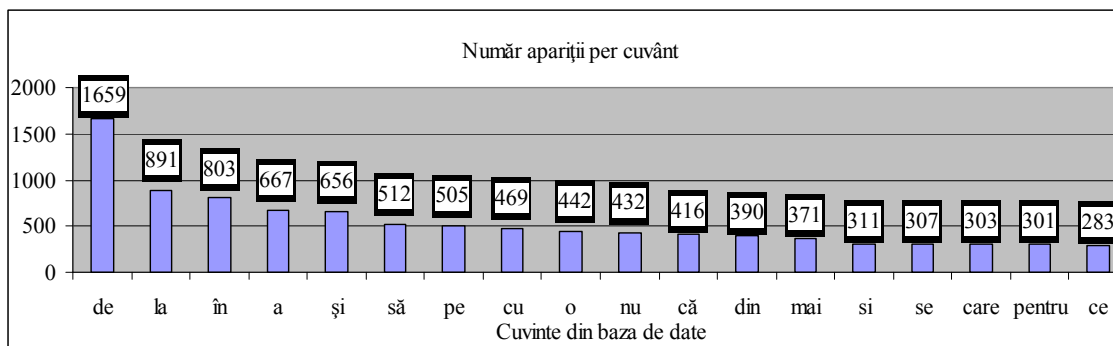


Figura 3: Cuvintele cu cele mai multe apariții în baza de date.

### Transcriere fonetică

Etapa următoare în crearea bazei de date a constat în transcrierea fonetică a fiecărui cuvânt. Transcrierea fonetică s-a realizat pe baza regulilor SAMPA (Burileanu, 2002). Astfel a rezultat dicționarul fonetic.

### Transcriere trifonetică

Conversia transcrierilor fonetice în transcrieri trifonetice s-a realizat cu utilitarul HLEd. Pentru fiecare trifonem s-a construit un model Markov în formatul HTK.

### Etichetare la nivel de fișier audio

Etichetarea s-a realizat la nivel de fișier audio folosind utilitarul Wavesurfer disponibil gratuit la: <http://www.speech.kth.se/wavesurfer/download.html>. În cadrul etichetării la nivel de fișier audio, pauzele scurte între cuvinte au fost marcate explicit prin „sp” și pauzele mai lungi au fost marcate prin „sil”. S-a obținut o etichetare mai puțin fiabilă decât cea la nivel de cuvânt însă cu prețul îmbunătățirii timpului de etichetare.

## 4. Statistici pentru limba română

Baza de date descrisă în secțiunile anterioare este destinată utilizării în cadrul unei aplicații de recunoaștere a vorbirii spontane. Rezultatele sunt prezentate prin intermediul statisticilor care fac referire la numărul de apariții al cuvintelor și al trifonemelor. În lucrarea de față în locul fonemelor au fost utilizate trifoneme. Având în vedere faptul că între cuvinte a fost folosit ca separator „sp”, care nu reprezintă un fonem, atunci se poate considera că începutul cuvintelor și sfârșitul acestora sunt realizate ca difoneme.

Exemplu: cuvântul „c a s @” are următoarea reprezentare fonetică: „c-a”, „c-a+s”, „a-s+@”, „s+@”. Motivația pentru care se folosesc trifonemele este că acestea permit analizarea contextului, întrucât sunt entități care păstrează informația despre ce găsim înaintea și după un anumit fonem.

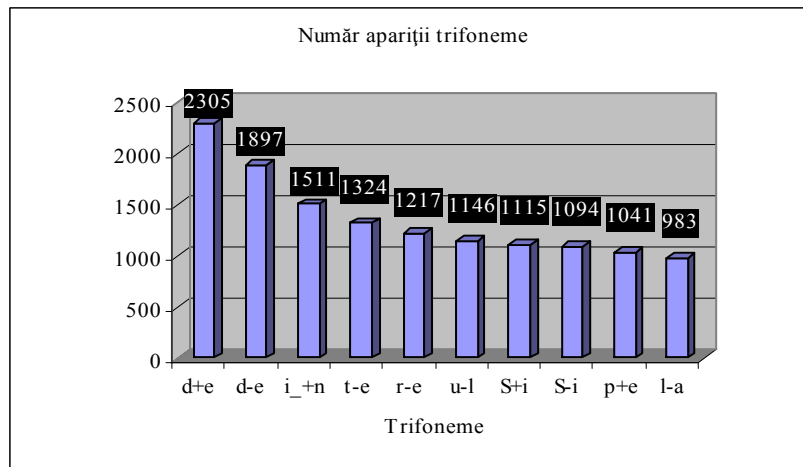
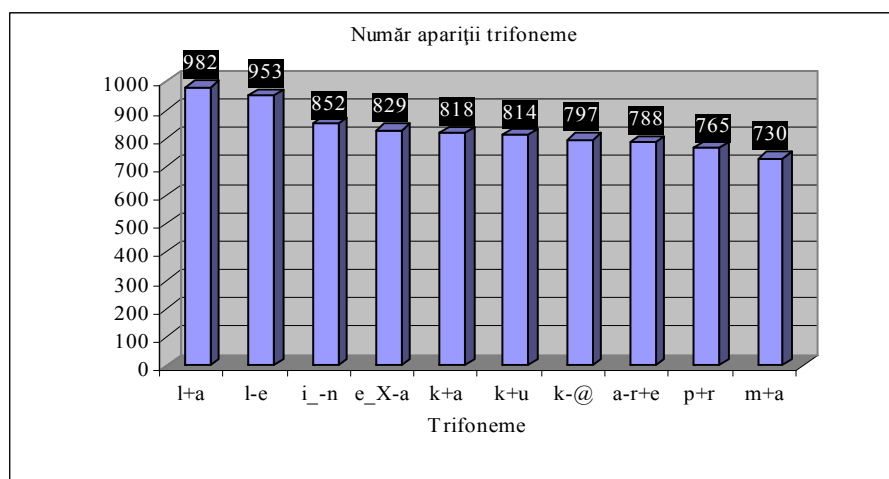


Figura 4: Trifoneme cu cele mai multe apariții.

În continuare, se vor prezenta statistici referitoare la trifonemele din cadrul bazei de date construite în limba română. Astfel, în figura 4 sunt ilustrate trifonemele care au cele mai multe apariții în cadrul bazei de date având între 983 și 2305 de apariții. Predominante sunt trifonemele „d+e” și „d-e”.

Figura 5 ilustrează următoarele 10 trifoneme din cadrul bazei de date care au între 982 de apariții și 730 apariții în baza de date. Predominante sunt trifonemele „l+a” și „l-e”.



**Figura 5:** Trifoneme cu apariții multiple.

Din cele 5095 de trifoneme care caracterizează vocabularul construit se pot extrage următoarele caracteristici: 1% din numărul total de trifoneme au între 500 și 1000 de apariții; 7% din trifoneme au între 100 și 500 de apariții; 8% din trifoneme au între 50 și 100 de apariții; alte 8% din trifoneme au între 30 și 40 de apariții; 19% din trifoneme au între 10 și 30 de apariții; 16% din trifoneme au între 5 și 10 apariții iar 41% din trifoneme au între 1 și 5 apariții.

## 5. Concluzii

Baza de date construită pentru limba română în scopul recunoașterii vorbirii spontane are o dimensiune relativ medie. Pentru munca de cercetare viitoare și pentru a putea realiza o bună analiză a caracteristicilor vorbirii spontane, obiectivul principal în viitorul apropiat îl constituie creșterea numărului total de cuvinte și creșterea numărului de apariții al cuvintelor.

Caracteristicile principale ale bazei de date în limba română, realizate în scopul utilizării în cadrul proiectului sunt următoarele: durata înregistrărilor este de 4 ore, înregistrările sunt preluate din diverse medii de lucru, vorbitorii folosesc vorbirea spontană. Au fost folosiți 12 vorbitori diferiți (8 voci feminine, 4 voci masculine), fiecare vorbitor având mai multe sesiuni de înregistrare, în medie 20 de sesiuni per vorbitor. Baza de cuvinte numără în total 37.604 cuvinte, 8068 de cuvinte au apariție singulară iar numărul total de trifoneme este 5095.

Într-o primă fază, testele de recunoaștere efectuate au condus la obținerea unor rezultate cu o performanță scăzută. Cauzele principale au fost identificate și localizate după cum urmează: numărul foarte mic de apariții pentru fiecare trifonem; calitatea redusă a

anumitor fișiere audio; erori ale intervenției umane în etapa de prelucrare a fișierelor audio.

Din analiza cauzelor menționate anterior a rezultat necesitatea parcurgerii următoarelor etape:

- prelucrarea manuală a materialului audio pentru a corecta eventualele erori comise la pregătirea datelor;
- verificarea fișierelor cu etichete;
- antrenarea la nivel de fonem a fonemelor care au mai multe apariții și apoi refolosirea anumitor parametri calculați la antrenarea fonemelor și a trifonemelor;
- segmentarea fișierelor etichetate la nivel de fonem în fișiere cu durata sub un minut;
- introducerea unor constrângeri gramaticale care să amelioreze semnificativ calitatea recunoașterii.

**Mulțumiri.** Cercetările prezentate în această lucrare au fost finanțate de Guvernul României, prin grantul de cercetare IDEI, nr. 930/2007.

### Referințe bibliografice

- Peinado, A., Segura, J. (2006). *Speech Recognition Over Digital Channels: Robustness and Standards*, Chapter 2, Pages 7-30
- Russell, S., Norvig, P. (2003). *Artificial Intelligence: A Modern Approach*, Prentice Hall (Second Edition)
- Burileanu, D. (September, 2002). Basic Research and Implementation Decisions for a Text-to-Speech Synthesis System in Romanian, *International Journal of Speech Technology*, Volume 5, Number 3, Pages 211-225
- Evermann, G., et al. (2005). *The HTK Book*, Version 3.0, Cambridge University Engineering Department;
- Munteanu, D. (2006). *Contribuții la realizarea sistemelor de recunoaștere a vorbirii continue pentru limba română*, Teză de doctorat, Academia Tehnică Militară din București
- Huang, X., Acero, A., Hon, H.-W. (2001). *Spoken Language Processing - A Guide to Theory, Algorithm and System Development*, Prentice Hall



# METODĂ IERARHICĂ DE DETECȚIE A FUNDAMENTALEI

MARIUS-DAN ZBANCIOC<sup>1,2</sup>, HORIA-NICOLAI TEODORESCU<sup>1,2</sup>

<sup>1</sup> *Institutul de Informatică Teoretică al Academiei Române - Filiala Iași*

<sup>2</sup> *Universitatea Tehnică "Gheorghe Asachi", Iași – România*

{zmarius ,hteodor}@etc.tuiasi.ro

## 1. Introducere

Problematica analizei prozodice nu este încă rezolvată, datorită naturii nestaționare a semnalului vocal și inexistenței unui suport matematic pentru definirea conceptului de frecvență fundamentală,  $F_0$ , respectiv a conceptului de formant. Colective din institutul nostru de cercetare au implementat mai multe instrumente automate de extragere a lui  $F_0$ , pe baza unor metode clasice din literatură (Rowden, 1991; Rabiner & Juang 1993; Rabiner & Schafer 1978; Calliope, 1989; O'Shaughnessy, 1987; Cristea & Valsan, 1999): AMDF, autocorelație, HPS și metoda cepstrală, respectiv a unei metode hibride propuse în (Teodorescu, 2006). Recent s-a realizat și un instrument de extragere a formanților superiori  $F_1, \dots, F_4$ , validarea rezultatelor de ieșire fiind realizată într-un bloc decizional neuro-fuzzy. Sistemul ierarhic hibrid, cu un bloc neuro-fuzzy care înglobează mai multe metode de detecție a  $F_0$ , ponderează diferit fiecare extractor funcție de performanțele lui și controlează astfel influența fiecărei metode asupra valorilor finale.

Îmbunătățirile recent aduse de noi instrumentelor de analiză au vizat blocul de pre-procesare, care realizează operații de filtrare, respectiv segmentare a zonei de interes (vocalice) de zona de background (consonantică sau de pauze între rostiri). Valorile de prag (threshold) folosite anterior erau determinate empiric, noile valori fiind determinate pe baza regulilor furnizate de un arbore de decizie. Instrumentele de analiză (codurile sursă și executabilele) sunt disponibile pe site-ul SRoL - „Proiectul Sunetele Limbii Române” ([http://iit.iit.tuiasi.ro/romanain\\_spoken\\_language/index.htm](http://iit.iit.tuiasi.ro/romanain_spoken_language/index.htm)). Aceste instrumente de analiză a informației prozodice bazate pe valorile frecvenței fundamentale și ale formanților sunt utile în aplicații de recunoaștere și sinteză a semnalului vocal (prin modulele de studiu al caracteristicilor fonemelor), în aplicații de studiu al intonației și al altor informații paralingvistice, sau în aplicații de identificare de limbă și de particularități ale acesteia (dialecte) etc.

## 2. Descrierea modulelor componente ale aplicației

Pentru implementarea unui instrument pentru detecție de  $F_0$  și pentru determinarea valorilor formanților, s-au conceput modular mai multe aplicații (programe), fiecare dintre acestea fiind apelabilă independent de restul aplicațiilor (principiul modularității). Rezumăm aceste aplicații, fiecare fiind asociată uneia dintre etapele de analiză și principalele lor funcții:

**Modulul de extragere de trăsături** din fișiere de sunet prin corelare cu informația extrasă din fișierele adnotate „\*.TextGrid”. Se determină un vector de timp folosit pentru delimitarea fonemelor și se extrag pattern-uri de trăsături (energie benzi spectrale, valoare medie energie în domeniul timp, deviația standard, rata trecerilor

prin zero etc.), care să fie folosite ulterior de o metodă de clasificare automată (arbori de decizie, rețele neuronale, algoritmi genetici etc.).

### **Modulul de preprocesare**

- Filtrarea semnalului (folosind un filtru median sau un filtru de mediere, respectiv un filtru trece bandă [70, 5000 Hz]) pentru a izola mai bine banda de frecvențe de interes în care se caută fundamentală, respectiv formanții.
- Segmentarea V/C zonei vocalice de zona consonantică (pe baza energiei din fereastra de analiză în domeniul timp, respectiv a energiei spectrale dintr-o bandă de frecvențe raportată la valoarea întregii energii spectrale).

### **Modulul de prelucrare statistică**

- Poate fi utilizat pe fișierele de ieșire, pentru determinarea unor valori ce caracterizează un anumit grup de foneme (de exemplu consoanele), sau pentru extragerea unui vector de trăsături specifice unui anumit fonem.
- Este apelat pentru determinarea automată a pragurilor utilizate de modulul de segmentare C/V, atunci când se dorește o minimizare a erorii de clasificare a zonei de interes, de zona de background (fundal sonor).

### **Modulul de detecție de valori formantice F0, F1,..., F4 folosește**

- metoda autocorelației (analiză în domeniul timp);
- metoda diferențelor AMDF (analiză în domeniul timp);
- metoda produsului spectrelor armonice HPS (analiză în domeniul frecvențelor);
- metoda cepstrală (analiză în domeniul que-frecvențelor) – utilizată și pentru căutarea formanților superiori.

**Modulul decizional** pentru ponderarea ieșirilor furnizate de fiecare metodă de detecție de formanți. Sunt eliminate valorile eronate prin comparare cu un număr de N ieșiri anterioare, respectiv ieșiri ulterioare (fără a mai respecta condiția de cauzalitate). Printr-un algoritm multicriterial, funcție de performanțele fiecărei metode se asociază acestora ponderi, determinate astfel încât să se apropie cât mai mult statistic valoarea finală de valorile „reale” ale formanților.

## **3. Modulul de preprocesare**

Etapa de preprocesare este una esențială a instrumentului de analiză a informației prozodice. În această etapă se realizează filtrarea semnalului pentru eliminarea zgomotelor suprapuse peste semnalul util, cum ar fi zgomotul indus de rețea (50 Hz) și zgomotul echipamentului de înregistrare. Banda pentru căutarea formanților se consideră a fi [70-5000Hz]. Extragerea acestei benzi se realizează printr-un filtru digital trece bandă (FTB), iar eliminarea zgomotului uniform prin aplicarea unui filtru de mediere (FTJ). Pentru o filtrare bună a semnalului, fără a afecta mult semnalul original, am preferat alegerea unui ordin mai mic pentru filtrul de mediere și aplicarea înseriată (repetată) a filtrului.

În această etapă se realizează și segmentarea semnalului, încercându-se o cât mai bună separare a zonelor vocalice de zonele consonantice și de zonele care sunt pauze între rostiri. Căutarea valorilor formantice se face doar pe zonele vocalice extrase. Realizarea unui instrument automat de analiză este dificilă din mai multe puncte de vedere. Înregistrările diferă prin energia semnalului (vorbitorul poate vorbi mai încet sau mai tare), prin raportul semnal / zgomot, prin spectrul de frecvențe specific fiecărui vorbitor (se știe că, adesea, altfel arată traseele formantice în cazul unui vorbitor masculin și altfel pentru un vorbitor de sex feminin, în special pentru primii doi formați  $F_0$  și  $F_1$ ). O segmentare cât mai precisă este esențială pentru obținerea de rezultate bune în final.

Separarea zonei vocalice de zona consonantică se realizează cu două metode:

- Criteriu global: Se compară energia din fereastra curentă de analiză  $E_w = \sum_{n=1}^N x_n^2$  cu energia maximă  $E_{\max}$  calculată după ce se parcurge tot semnalul. Dacă  $E_w > procent_1 \cdot E_{\max}$ , atunci se consideră că semnalul din fereastra curentă este vocalic și că se poate determina  $F_0$ . S-a folosit  $procent_1 = 20\%$ .
- Criteriu local: Dacă energia spectrală în banda  $[70, 1000\text{Hz}] > procent_2$  din toată energia spectrală, se consideră că avem vocală. Premisa acestei segmentări este că energia spectrală este mare în zona formanților. Pragul folosit este de  $50\%$  ( $procent_2$ ).

Primul criteriu este *global*, căci după ce se caută în tot semnalul fereastra cu energie maximă se consideră că doar secvențele de semnal care au energia ferestrelor de analiză mai mare de  $20\%$  sunt de interes în analiză. Valorile de prag utilizate, de  $0.2$  și  $0.5$ , au fost estimate empiric, după încercări succesive de ajustare a acestora.

Segmentarea folosind doar cele două criterii nu este satisfăcătoare, motiv pentru care s-au extras vectori cu mai multe trăsături (descriși mai jos) și au fost introduși într-un *arbore de decizie* (<http://www.rulequest.com/see5-win.html>), instrument capabil să furnizeze un set de reguli, care să minimizeze eroarea de clasificare.

Folosind *modulul de căutare fișiere* de adnotare s-au obținut vectori de timp  $[t_0, t_1, \dots, t_n]$  care delimitează fiecare fonem din secvența rostită. Dintre nivelurile de structurare a informației folosite la adnotare (foneme, silabe, cuvinte, propoziții etc.) s-a utilizat doar segmentarea la nivel de fonem (nivel 1 în figura 1).

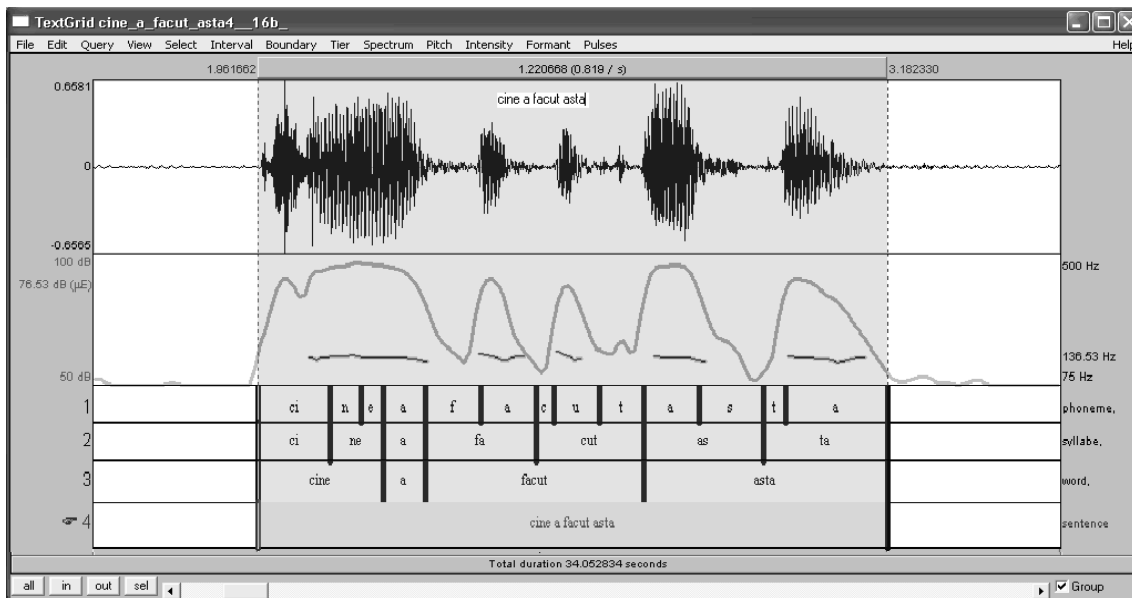
După citirea unui fișier adnotat *nume.TextGrid* se identifică fișierul de sunet *nume.wav* corespunzător, verificând dacă dimensiunea acestuia coincide cu limitele de timp  $[t_0, t_{end}]$ .

$$\frac{\text{dimensiune\_fișier} - \text{dimensiune\_header}}{\text{nr\_biti\_per\_esantion} \cdot \text{Frecv\_esantionare}} = t_{end} - t_0$$

Parametrii de intrare în sistemul de clasificare automată sunt determinați pentru fiecare fonem din secvență rostită analizată și sunt următorii:

- *nume fonem*;

- *id\_class* identificator clasa (1-vocală, 2-consoană, 3-pauză rostiri), *identificator fișier*;
- *zcr* (rata trecerilor prin zero, valoare normalizată pentru durata de o secundă);
- *avg\_e* energia medie a ferestrelor de analiză, din secvența de semnal ce delimitează fonemul curent;
- *std\_e* deviația standard a energiei pe durata fonemului, de la valoarea medie *avg\_e*;
- *B1* energia în banda [70, 500]Hz exprimată procentual față de toată energia spectrală;
- *B2, B3, B4* energii spectrale în benzile [500, 1000] Hz, [1000, 2000] Hz și [2000, 5000] Hz



**Figura 1.** Adnotarea unui fișier de sunet pe mai multe niveluri folosind Praat™

Ferestrele de analiză utilizate au dimensiunea de  $W=1024$  eșantioane, ceea ce corespunde, pentru o frecvență de eșantionare  $F_s$  de 22050 Hz, la o durată de 46,44 ms, respectiv pentru  $F_s$  de 16000 Hz la 64 ms. Pentru fonemele ale căror durate  $[t_i, t_{i+1}]$  au fost mai mici decât  $W$  nu s-au extras vectori de trăsături. Într-o rafinare ulterioară a algoritmului de segmentare se vor asocia funcții de apartenență (f.a.) fuzzy trapezoidale fiecărui fonem și se vor pondera rezultatele ferestrelor de analiză cu gradul de apartenență descris de f.a.

Rezultatele prezentate mai jos sunt realizate pe un set de 12 fraze adnotate ("A trecut așa un răstimp", "O ști el careva cum să rezolve asta", „Mama vine și ea mai târziu”, „Mama știe ea ce face”, „Chiar știe el ce face?”, „Vine ea mama!”).

See5 [Release 1.15]

Rule-based classifiers-50% data for training

```
Rule 1: (270/120, lift 1.4)
  E_MED > 0.000516
  B4 <= 0.068311
  -> class 1 [0.555] | vocala
```

## METODĂ IERARHICĂ DE DETECȚIE A FUNDAMENTALE

```

Rule 2: (55/2, lift 1.8)
  E_MED > 0.008076
  E_MED <= 0.095878
  B2 <= 0.053437
  B4 > 0.00501
  -> class 2 [0.947] | consoana
Rule 9: (120/16, lift 1.6)
  E_STD <= 0.085847
  B2 <= 0.014152
  -> class 2 [0.861] | consoana
Rule 10: (16, lift 11.1)
  E_MED <= 0.008076
  B2 > 0.014152
  B2 <= 0.07656
  B4 <= 0.068311
  -> class 3 [0.944] | pauza rostiri
Rule 11: (36/5, lift 9.9)
  E_MED <= 0.000516
  -> class 3 [0.842] | pauza rostiri
    
```

Tabel 1. Evaluarea matricei de confuzie pe un set de date de antrenare de 400 cazuri: (eroare 7.8%)

clasa 1	clasa 2	clasa 3	clasificate ca:
133	21	0	(a): clasa 1 vocale
7	203	2	(b): clasa 2 consoane
0	1	33	(c): clasa 3 pauză rostiri

Tabel 2. Evaluarea matricei de confuzie pe un set de date de test de 400 cazuri: (eroare 23.8%)

clasa 1	clasa 2	clasa 3	clasificate ca:
130	54	0	(a): clasa 1 vocale
31	157	7	(b): clasa 2 consoane
0	3	18	(c): clasa 3 pauză rostiri

Nu s-au luat în calcul diftongii, triftongii și nici sunetele 'ghe', 'ghi', 'che', 'chi', 'ș', 'ț'. Din cele 11 reguli furnizate de sistem, una a fost folosită pentru clasificarea vocalelor, 8 pentru clasificarea consoanelor și 2 pentru pauze. Se justifică separarea consoanelor într-un studiu viitor în mai multe clase, funcție de particularitățile acestora, de exemplu grupul consoanelor plozive, care se confundă ușor în segmentare cu zona de pauză, precum și al consoanelor semivocalice ('l', 'm', 'n', 'r'), pentru care are sens să se extragă informație formantică și care se confundă des în segmentare cu vocalele (vezi Tabel 1, 2).

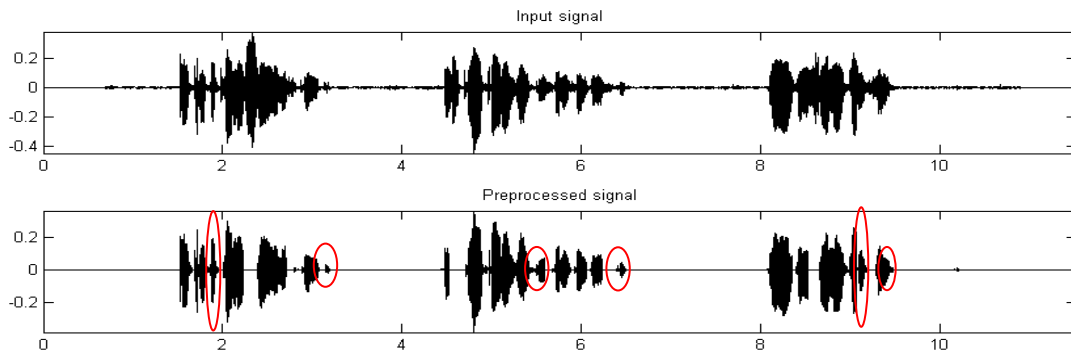
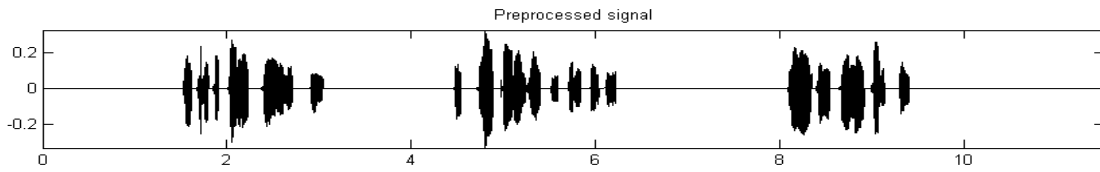


Figura 2. Selectarea zonelor vocalice de interes în urma aplicării noilor algoritmi de segmentare



**Figura 2.** Segmentarea V/C folosind algoritmi de segmentare anteriori

În urma aplicării algoritmilor de segmentare folosind valorile prag determinate automat de regulile generate de arborele de decizie, zona de interes este mai precis selectată, apar zone vocalice, marcate pe figura 2, pe care algoritmii anteriori nu reușeau să le identifice. S-au realizat extrageri de reguli, folosind instrumentul **See5** pe mai multe seturi de fișiere adnotate, o parte din regulile obținute și erorile de clasificare fiind date mai sus. În figura 2 s-a reprezentat segmentarea obținută aplicând negata regulii R1, pentru eliminarea zonelor nevocalice:

```
IF (E_MED < 0.000516) OR (B4 > 0.068311) => nu este vocala
```

Pentru determinarea zonei de pauză între rostiri s-a folosit o clasificare, utilizând doar rata trecerilor prin zero -  $zcr$  și energia medie în domeniul timp -  $e_{med}$ . S-au folosit de această dată 276 de fișiere adnotate de tipul „b\_ba\_aba” (consoană, consoană urmată de vocala 'a', consoană încadrată între două vocale 'a'). Au rezultat cca. 2000 de vectori de trăsături. Regulile generate pentru clasificarea pauzei dintre rostiri au permis identificarea a 119 secvențe de pauză din 124. Eroarea globală, de 9,8%, este afectată de confuziile între regiunile vocalice și cele consonantice. Ulterior, vom încerca minimizarea acestei erori, prin introducerea mai multor clase (categorii) de consoane și introducerea în vectorii de trăsături a unor noi parametri de intrare.

#### 4. Modulul de detecție a valorilor formantice $F_0, F_1, \dots, F_4$

Metodele de extragere a informației prozodice sunt clasificate, funcție de domeniul de analiză a datelor, în două categorii:

*Metode de analiză în domeniul timp:* a) autocorelația; b) metoda diferențelor AMDF;

*Metode de analiză în domeniul frecvențelor* (informație spectrală): a) HPS - produsul spectrelor armonice; b) analiza cepstrală.

Toate metodele de analiză folosesc la rulare aceiași parametri (dimensiunea ferestrei de analiză,  $W$  și valoarea pasului de deplasare a ferestrei) pentru ca vectorii de ieșire să aibă aceeași lungime și să poată fi comparați în final în modulul decizional. Se recomandă ca durata ferestrei de analiză să fie de minim 4-5 ori perioada fundamentalei maxime pentru o determinare bună a  $F_0$ . Pentru  $F_s = 22050$  Hz, considerând limita inferioară a domeniului de interes în detecția frecvențelor formantice de 80Hz, se obține o fereastră de minim 1100 eșantioane, iar pentru  $F_s = 16000$ Hz,  $W$  trebuie să fie de 800 de eșantioane.

Metodele de analiză în domeniul frecvențelor necesită ferestre de analiză de durate mai mici, dar chiar și în cazul acestora, pentru metoda cepstrală durata minimă este de 615 eșantioane. Deci dimensiunea ferestrei de analiză va fi de minim 1024 de eșantioane (durată de 46 – 64 ms). Dacă  $W$  este prea mare, e posibil ca ea să includă mai multe foneme, sau tranziții rapide în fluctuațiile valorilor formantice și deci există riscul ca

prin medierea rezultatelor detecția sa aibă o imprecizie mare, inclusiv să redea informații asociate unei întregi grupări de foneme.

**Metode de analiză în timp (autocorelația, AMDF)**

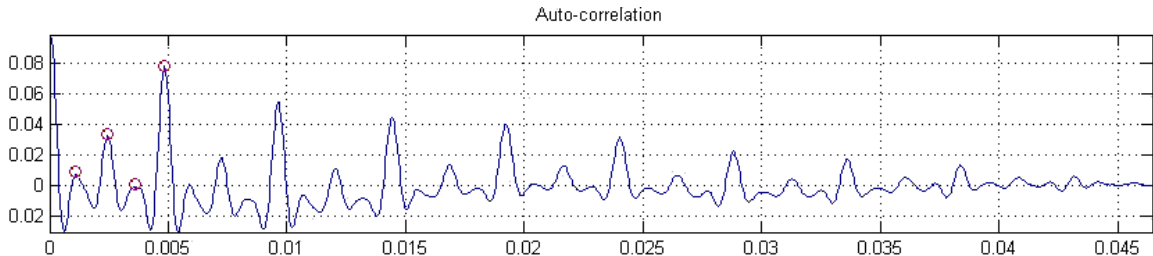
Aceste metode sunt des întâlnite în literatura de specialitate și sunt considerate metode care conduc la o bună detecție a frecvenței fundamentale. Pentru extractorii de frecvență fundamentală pe care i-am implementat, metodele de analiză în domeniul timp au dat mai puține erori în detecția  $F_0$  decât cele în domeniul frecvențelor.

Funcția de corelație aplicată pe două semnale  $x$  și  $y$  oferă informații legate de similitudinile dintre acestea. Această metodă de comparație este utilizată pentru detecția unor regularități (legate de periodicitatea semnalelor, în cazul în care acestea pot fi considerate periodice sau cvasi-periodice).

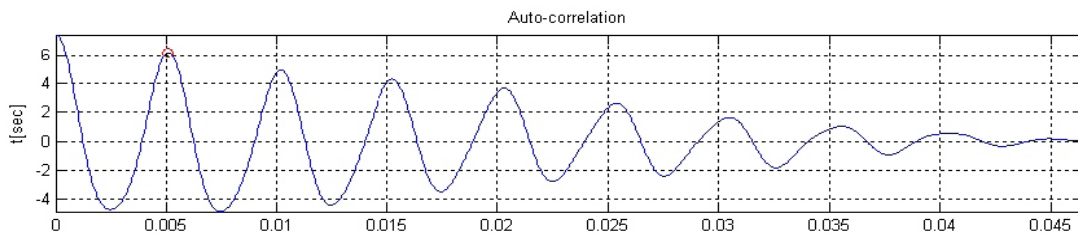
$$R_{XY}[n] = \sum_{k=0}^{W-n} x[k] \cdot y[k+n]; \quad R_{YX}[n] = \sum_{k=0}^{W-n} y[k] \cdot x[k+n], \quad n = \overline{0, W} \quad (1)$$

$$C[n] = R_{XX}[n] = \sum_k x_k \cdot x_{n+k} \quad (\text{autocorelația}) \quad (2)$$

Valorile de maxim local sunt date de periodicitatea semnalului. Valoarea maximă găsită în  $R_{XY}[0]$  nu trebuie luată în considerație; ea reprezintă energia semnalului într-o fereastră de analiză de dimensiune fixată  $W$ . Următorul maxim local este asociat cu perioada fundamentală și căutarea lui poate fi limitată într-un interval de valori dat de banda de căutare a frecvenței fundamentale [70, 500] Hz. Pentru o frecvență de eșantionare de  $F_s=16000$  Hz se va căuta în semnalul furnizat de funcția de autocorelație între valorile  $R_{XX} [F_s/F_{0max}, F_s/F_{0min}]$ , adică  $R_{XX} [16000/500, 16000/70]=R_{XX} [32, 229]$ .



**Figura 3.** Detecție valori formanți prin metoda autocorelației  
a) extragere maxime corespunzătoare valori formantice  $F_0, F_1$



**Figura 3.** b) caz defavorabil - extragere valoare frecvență fundamentală  $F_0$

Valorile frecvenței fundamentale  $F_0$  sunt ușor de extras folosind această metodă, dar există situații când informațiile corelate cu ceilalți formanți  $F_1$  și  $F_2$  nu sunt prezente în semnalul funcției de corelație. În fig. 3, pentru primul semnal există mai multe maxime locale cu periodicități diferite corespunzătoare primilor formanți, pentru al doilea

semnal se observă doar distribuția periodică a unui singur maxim local (corespunzător lui  $F_0$ ).

Metoda AMDF – magnitudinea medie a funcției diferență (Average Magnitude Difference Function) se aseamănă ca algoritm cu funcția de autocorelație; diferența între cele două constă în faptul că funcția diferență nu necesită operații de înmulțire:

$$D_n = \sum_k (x_n - x_{n+k}), \quad n, k \in \overline{1, W} \quad (3)$$

La metoda diferențelor AMDF, minimele locale sunt cele care servesc la calcularea valorii perioadei fundamentale  $T_0$ , respectiv a primilor formanți  $T_1, T_2$ , spre deosebire de metoda autocorelației, unde maximele locale erau folosite în determinarea lui  $T_0$ .

Metoda permite o mai bună detecție a formanților, în special a primului formant  $F_1$ , față de funcția de autocorelație, pentru care detecția este mai dificilă datorită estompării valorilor corespunzătoare formanților, prin funcția de multiplicare.

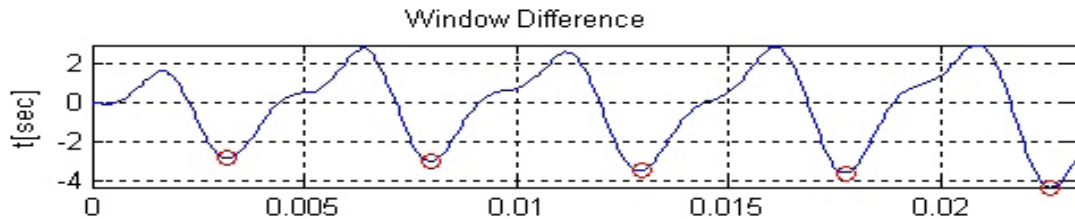


Figura 4. Detecție valori formantice prin metoda funcției diferență AMDF

Aceleași dificultăți întâlnite la metoda autocorelației sunt valabile și pentru metoda diferențelor, vectorul pe care se caută valorile formantice fiind construit de aceasta dată din distanțele dintre minimele locale ale semnalului diferență.

### Metode de analiză spectrale (HPS, metoda cepstrală)

Metoda cepstrală se bazează pe separarea componentelor spectrale care țin de modul în care este generat sunetul  $H_g$  (depind de frecvența de rezonanță a corzilor vocale, dimensiunea tubului generator și pot oferi informații despre frecvența fundamentală), de cele care depind de modul de filtrare a semnalului vocal  $H_f$  (și care descriu modelul rezonator al cavităților în care se formează sunetul vocal). În formula de calcul al cepstrului (spectrul spectrului logaritmat), operația de înmulțire dintre spectrul semnalului excitator și spectrul funcției de transfer este transformată prin logaritmare într-o operație de adunare. Cele două componente sunt separabile; căutarea maximului corespunzător frecvenței fundamentale se face în banda [70, 500] Hz.

$$H(\omega) = H_g(\omega) \cdot H_f(\omega), \quad \text{cepstrum} = \text{IFFT}(\log|\text{FFT}(s)|) \quad (4)$$

$$\text{cepstrum} == \mathfrak{T}^{-1}(\log|H_g(\omega) \cdot H_f(\omega)|) = \mathfrak{T}^{-1}(\log|H_g(\omega)|) + \mathfrak{T}^{-1}(\log|H_f(\omega)|)$$

Deoarece calculul cepstrului implică trecerea în domeniul quefrențelor (spectrul spectrului), algoritmul de calcul are cea mai mare complexitate dintre cele 4 metode de detecție de informație prozodică. Se recomandă  $W$  de minim 1024 eșantioane, deoarece pentru o frecvență de eșantionare de  $F_s=22050$  pentru banda de frecvențe [70, 500]Hz avem nevoie ca vectorul cepstral (corespunzător frecvențelor pozitive) să conțină minim  $F_s/F_{0\min}=22050/70=315$  eșantioane. Deoarece jumătate din spectrul unui semnal este



METODĂ IERARHICĂ DE DETECȚIE A FUNDAMENTALE

asociat frecvențelor pozitive, iar cealaltă jumătate frecvențelor negative,  $W$  minimă este de  $2 \cdot F_s / F_{0min} = 650$  eșantioane.

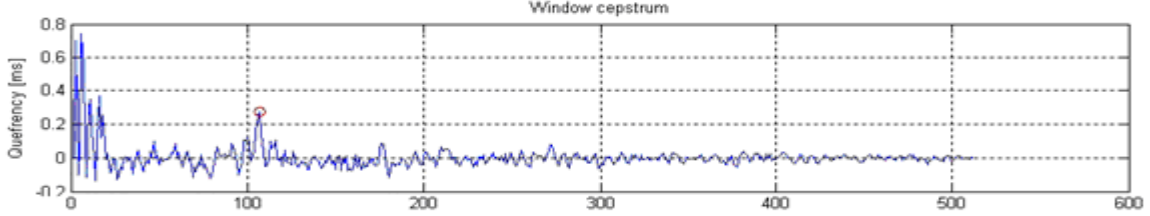


Figura 5. Detecție valoare frecvența fundamentală F0 prin metoda cepstrală

Pentru determinarea formațiilor prin metoda cepstrală se calculează „spectrul netezit” (o anvelopă a spectrului) prin aplicarea unui FTJ în cepstru. Se păstrează din cepstru doar primele  $L$  valori (restul fiind anulate) și se aplică transformata inversă, obținând un semnal spectral cu tranziții lente („componenta frecvențelor înalte” fiind eliminată).

$$cepstrum^* = \begin{cases} cepstrum[k] & , k \leq L \text{ sau } c > N_w - L \\ 0 & , \text{altfel} \end{cases}$$

$$S^*(\omega) = \exp(FFT(cepstrum^*)) \quad (5)$$

Funcție de dimensiunea ferestrei de liftare  $L$  avem mai multe sau mai puține tranziții în spectru (frecvența de tăiere a FTJ este mai mică sau mai mare). În anumite situații pentru valori mici ale lui  $L$ , unii formați pot fi greu de găsit sau nu apar în spectru, cum este și cazul lui  $F_1$  când  $L=30$  (Figura 6.a). Pentru valori mai mari ale lui  $L$  sunt mai mulți candidați, fiind necesară alegerea unei valori reprezentative.

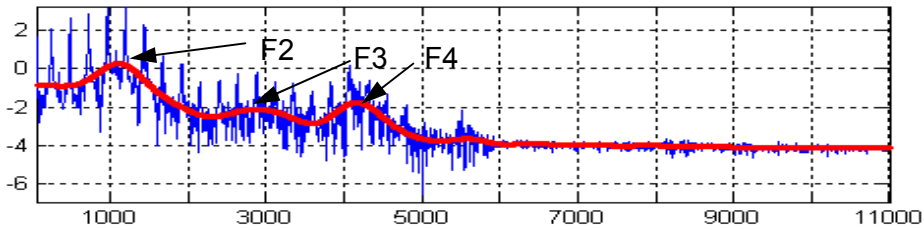


Figura.6a. Spectru + spectru „netezit” vocala 'a', pentru L=30

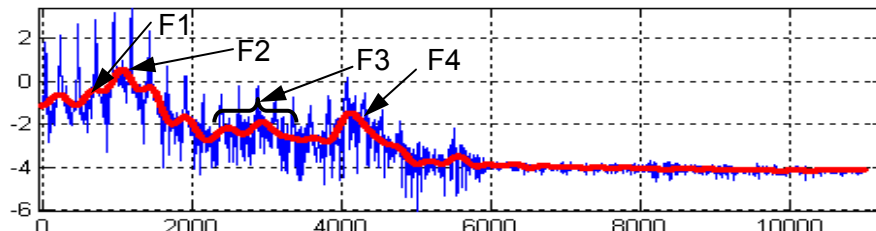


Figura.6b. Spectru + spectru „netezit” vocala 'a', pentru L=60

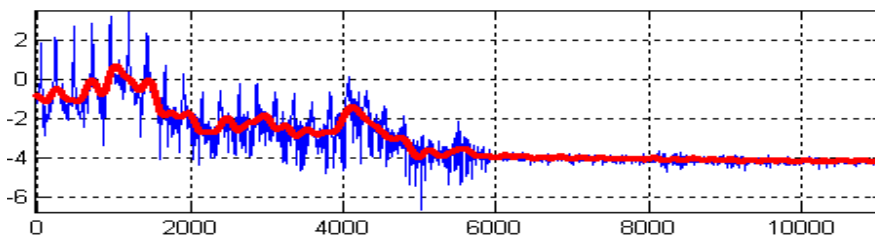


Figura 6c. Spectru + spectru „netezit” vocala 'a', pentru L=90

Ca direcție viitoare de cercetare se va determina valoarea optimală a lui  $L$  (dimensiunea ferestrei de liftare), pentru o detecție mai bună a formanților, printr-un algoritm care să permită asocierea de mulțimi fuzzy în care să fie căutați formanții. Alegerea mulțimilor fuzzy va fi realizată funcție de valoarea lui  $F_0$ , determinată anterior de extractorul de frecvență fundamentală. Al doilea autor a propus o metodă ce utilizează o funcția diferență modificată combinată cu cea a spectrului produs de armonici pentru a realiza o detecție mai bună a valorii frecvenței fundamentale. Metoda este descrisă pe larg în (Teodorescu H.N., 2006).

Metoda HPS constă în determinarea spectrului semnalului, decimarea acestuia cu factori de decimare  $1/2, 1/3, 1/4, \dots$ , și realizarea produsului între semnalele realizate. Decimarea se realizează printr-o parcurgere a semnalului și selectarea eşantioanelor cu un pas (2,3,4,...). Metoda HPS – Harmonic Product Spectrum se bazează pe proprietatea că în spectrul unui semnal periodic cu frecvența fundamentală  $F_0$ , apar maxime la multiplii acestei frecvențe  $2 \cdot F_0, 3 \cdot F_0, 4 \cdot F_0, \dots$  (armonicele fundamentale). Dacă semnalul este rescalat cu factori  $1/2, 1/3, 1/4, \dots$ , după operația de decimare, prin produsul semnalelor rezultate care au toate un maxim spectral în jurul frecvenței fundamentale  $F_0$ , celelalte maxime vor dispărea sau vor fi puternic atenuate. Decimarea cu un factor  $1/k$  a valorilor spectrale se poate face fie selectând o valoare (de obicei prima) dintr-un set de  $k$  valori consecutive, fie realizând media celor  $k$  valori,

$$H_n^k = H_{k \cdot n}^0 \text{ (decimare) sau } H_n^k = \frac{1}{k} \sum_{i=0}^{k-1} H_{k \cdot n+i}^0 \quad (6)$$

unde  $H^0$  este spectrul semnalului și  $H^k$  semnalul rezultat după scalarea cu un factor  $1/k$ .

Metoda HPS ridică probleme atunci când avem subarmonici de amplitudine mare ale lui  $F_0$ . În special prima subarmonică este cea care conduce la detecții eronate de  $F_0$ . Aplicarea metodei HPS pentru determinarea formanților este anevoioasă, chiar și în lui  $F_1$  și  $F_2$  deoarece la fiecare înjumătățire a spectrului, banda de frecvențe rămasă este tot mai îngustă. De exemplu, la un semnal achiziționat la o frecvență de eşantionare de 16 kHz, spectrul util (al frecvențelor pozitive) este  $[0-8000]$  Hz. După aplicarea algoritmului de decimare HPS de 3 ori, banda de frecvențe rămasă este de  $[0-1000]$  Hz. La a patra aplicare, banda de frecvențe rămasă nu mai include  $F_1$ ! Metoda HPS este utilă doar atunci când se dorește accentuarea valorilor formantice de frecvență joasă, fiind acceptabilă în cazul frecvenței fundamentale.

### 5. Metoda hibridă comparativă

Pentru a putea realiza compararea rezultatelor, toate metodele de extragere  $F_0$  sunt rulate cu aceiași parametri la intrare (dimensiunea ferestrei de analiză, pas de deplasare etc.). Algoritmii de selecție se aplică pentru situațiile în care diferența dintre valoarea detectată printr-o metoda  $vF0\_m1$  este cu un procent de 20% mai mică sau mai mare decât valoarea furnizată de metoda altă metodă  $vF0\_m2$ . În caz contrar, se consideră valoarea lui  $F_0$  ca fiind media celor două valori.

```
IF (vF0_m1 > vF0_m2*(1+percent)) OR (vF0_m1 < vF0_m2*(1-percent))
THEN compară cu N vecini la stânga și [optional cu N vecini la dreapta]
```

## METODĂ IERARHICĂ DE DETECȚIE A FUNDAMENTALE

```
calculează media  $\overline{vF0}$  și abaterea standard  $\sigma_{F0}$   
IF  $|vF0\_m1 - \overline{vF0}| < |vF0\_m2 - \overline{vF0}|$  THEN  $vF0=vF0\_m1$   
ELSE  $vF0=vF0\_m2$   
ELSE  $vF0 = (vF0\_m1+vF0\_m2)/2$ 
```

Pentru a decide valoarea frecvenței fundamentale  $F_0$ , se selectează dintre valorile comparate cea valoare care este mai apropiată de valorile determinate anterior. Numărul de vecini  $N$  folosiți pentru această comparație depinde de dimensiunea pasului de deplasare a ferestrei de analiză, și se alege astfel încât să nu comparăm valori ale  $F_0$  aflate la o distanță mai mare de 3-5 ms. Sunt comparate valori ale lui  $F_0$  pe durate mici de timp, ca să nu apară fluctuații mari ale lui  $F_0$ . Stabilirea ponderilor ce sunt asociate fiecărui extractor este realizată statistic pe baza estimării raportului dintre detecțiile eronate de fundamentală și detecțiilor corecte. Astfel, o metodă cu mai puține erori are o influență mai mare asupra rezultatului final.

### 6. Concluzii

Validarea rezultatelor furnizate în final este realizată prin mai multe metode de detecție de  $F_0$ , respectiv pe baza comparațiilor cu valorile anterioare ale lui  $F_0$  (pe intervale mici de timp nu pot avea loc variații bruște). S-au comparat vizual valorile fundamentale date de sistemul nostru hibrid cu rezultatele furnizate de alte instrumente de detecție puse la dispoziția utilizatorilor pe Internet: WASP (<http://www.phon.ucl.ac.uk/resource/sfs/wasp.htm>), Praat (<http://www.praat.org>) și s-a constatat că rezultatele noastre sunt mai bune. Studiul s-a realizat pentru un număr de 12 pronunții de fraze aflate pe situl SROL la secțiunea „Fraze->Particularități lingvistice” [[www.etc.tuiasi.ro/sibm/romanian\\_spoken\\_language/ro/fraze\\_sd\\_arhiva.htm](http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/ro/fraze_sd_arhiva.htm)]. Urmează ca pe viitor să se realizeze o statistică automată pe un număr mai mare de fișiere.

Din simulările efectuate s-a constatat că, din punct de vedere al erorilor de detecție  $F_0$ , cele mai robuste sunt metodele de analiză în domeniul timp. Dintre acestea, metoda autocorelației este cea care oferă cele mai bune rezultate în detecție. Dintre metodele de analiză în domeniul spectral, atât metoda HPS, cât și metoda cepstrală ridică probleme mai ales la nivelul primei armonici, respectiv a primei subarmonici. O detecție mai bună a frecvenței fundamentale permite o extragere mai exactă a valorilor formantice (definirea intervalelor de căutare a formaților se face în funcție de  $F_0$ ). Dintre metodele de extragere  $F_1$ ,  $F_2$ ,  $F_3$  metoda cepstrală s-a dovedit cea mai sigură. O problemă doar parțial rezolvată, de care urmează să ne ocupăm, o constituie segmentarea mai precisă a zonelor vocalice.

**Mulțumiri.** Autorii mulțumesc recenzorilor pentru observațiile pertinente.

### Referințe bibliografice

- Rowden C. (1991), *Speech Processing*, McGraw - Hill Book Company, Chapter 2, pp.35-74.
- Rabiner L.R., Juang B.H. (1993), *Fundamentals of Speech Recognition* Englewood Cliffs, N.J.

- Rabiner L.R. Schafer R. W. (1978), *Digital Processing of Speech Signal*, Prentice-Hall, Inc. Englewood Clifford, pp. 11-65
- Calliope (1989), *La parole et son traitement automatique*, ISBN 2-225-81516-X, Masson, France
- O'Shaughnessy, D.O. (1987), *Speech Communication Human and Machine*, INRS-Telecom.
- Cristea P., Valsan, Z. (1999) *New Cepstrum Frequency Scale for Neural Network Speaker Verification*, Proc. of the VI<sup>th</sup> International Conference on Electronics, Circuits and Systems, ICECS, 5-8 sept. Cyprus.
- Teodorescu H.N., (2006), *Aplicații ale analizei și sintezei semnalului vocal*, Iași, Capitolul 2.
- Teodorescu H.N., Trandabăț D., Feraru M., Zbancioc M., Luca R., (2006a) “*A Corpus of the Sounds in the Romanian Spoken Language for Language-Related Education*”, I<sup>st</sup> International Conference on Human and Material Resources in Foreign Language Learning - HMRFL, Murcia, Spain
- Teodorescu H.N., Feraru M., Trandabat D., Zbancioc M. (2006b), “*Limba română vorbită*”, Atelierul Resurse lingvistice și instrumente pentru prelucrarea limbii române, ConsILR-06, 3-4, Iași, România, Editura Universitatii “A.I. Cuza” Iasi
- Teodorescu H.N., Zbancioc M., Mihailescu E. (2006c), “*Speech Technology and Bio-Medical Engineering Teaching Based on the Web – A New Tool and Case Study*”, International Conference on Interactive Computed Aided Learning, ICL, September 27 - 29, Villach, Austria
- Zbancioc M. (2006), *Tools for the Archive of the Romanian Language Sounds Project*, 4<sup>th</sup> European Conference on Intelligent Systems and Technologies, ECIT'2006, sept.20-23, Iași, Romania
- Proiectul Sunetele Limbii Române, [http://iit.iit.tuiasi.ro/romanain\\_spoken\\_language/index.htm](http://iit.iit.tuiasi.ro/romanain_spoken_language/index.htm)
- Praat, Boersma P., Weenink D., Institute of Phonetic Sciences, Amsterdam: <http://www.praat.org>
- WASP web page, <http://www.phon.ucl.ac.uk/resource/sfs/wasp.htm>
- Arbore de decizie See5 <http://www.rulequest.com/see5-win.html>

## **CAPITOLUL 2**

### **PLATFORME, DICȚIONARE ȘI CORPUSURI ADNOTATE PENTRU PRELUCRAREA TEXTELOR**



# LIMBA ROMÂNĂ ÎN PERSPECTIVA CLARIN

DAN CRISTEA<sup>1,2</sup>, IONUȚ CRISTIAN PISTOL<sup>1</sup>

<sup>1</sup>*Facultatea de Informatică, Universitatea “Al. I. Cuza” Iași,*

<sup>2</sup>*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

*{dcristea, ipistol}@info.uaic.ro*

## Rezumat

CLARIN<sup>1</sup> este un proiect PC7 care își propune dezvoltarea unei infrastructuri de resurse și tehnologii lingvistice, reunind ultimele progrese în domeniul prelucrării limbajului natural, într-o formă accesibilă celor din afara domeniului, cum ar fi specialiștii în științele umane, științele sociale și chiar publicului larg. Lucrarea propune o abordare teoretică și aplicativă de integrare a instrumentelor de procesare lingvistică dedicate limbii române, care se încadrează spiritului CLARIN. La baza acestei abordări stă ALPE<sup>2</sup>, un meta-sistem de configurare de soluții în problemele de tratamente aplicate limbilor naturale. Se propune o ierarhie care reunește instrumentele cunoscute de procesare lingvistică disponibile pentru limba română. Beneficiile acestei abordări, odată demonstrate, pot fi extinse la nivelul comunității globale a “consumatorilor” de procesări lingvistice, deziderat central al proiectului CLARIN.

## 1. Introducere

Ultimii ani au fost martorii unui interes crescând pentru prelucrarea lingvistică a limbilor europene vorbite în țările “noului val” al Uniunii Europene, interes ghidat în special prin proiecte europene de cercetare în domeniul prelucrării limbajului natural în diferite scopuri, de la sisteme de traducere automată la sisteme de e-learning. Printre alte manifestări de importanță națională ori internațională, întâlnirile Consorțiului de Informatizare pentru Limba Română (ConsILR<sup>3</sup>), care de câțiva ani s-au transformat în Ateliere de lucru “Resurse lingvistice și instrumente pentru prelucrarea limbii române”, au relevat rezultate tot mai interesante în ceea ce privește procesarea limbii române cât și în crearea de resurse lingvistice românești.

De mai mult timp ConsILR militează pentru colectarea informațiilor de natură lingvistică și a instrumentelor informatice capabile să proceseze texte în limba română. Aceste eforturi sunt în perfect acord cu preocupări similare manifestate în afara țării care vizează crearea de infrastructuri la nivel internațional pentru stocarea resurselor lingvistice și procesarea limbajului natural. Un exemplu în această direcție este recent lansatul proiect FP7 CLARIN, în care România este reprezentată prin două instituții cu statut de partener și alte două instituții ca membri. De mare actualitate, în acest context, sunt și preocupările de dezvoltare de meta-sisteme de procesare lingvistică. Exemple

---

<sup>1</sup> Common Language Resources and Technology Infrastructure Network

<sup>2</sup> Automated Linguistic Processing Architecture

<sup>3</sup> <http://consilr.info.uaic.ro>

sunt GATE (Cunningham et al., 2002), UIMA (Ferrucci și Lally, 2004) și ALPE (Cristea și Pistol, 2008), care permit conceperea unor structuri de procesare complexe, prin integrarea de module existente, în vederea construirii de aplicații ce presupun procesări lingvistice. Dintre aceste meta-sisteme, ALPE promite inclusiv facilitarea interacțiunii utilizatorilor nespecialiști cu modulele de procesare.

Capitolul doi al acestei lucrări descrie pe scurt proiectul CLARIN, prezentând obiectivele și avantajele includerii limbii române în această inițiativă europeană. Capitolul trei prezintă principalele funcționalități oferite de ALPE, iar capitolul patru propune o încercare de sistematizare într-o ierarhie a resurselor de procesare ce se cunosc pentru limba română. Capitolul cinci conține comentarii referitoare la aspectele practice ale interacțiunii cu o ierarhie ALPE. Capitolul șase descrie planul de lucru și obiectivele pe termen scurt și lung, precum și posibile direcții noi de dezvoltare din perspectiva limbii române în CLARIN.

## **2. CLARIN**

CLARIN (Váradi et al., 2008) este un proiect-program finanțat de Comisia Europeană, structurat în trei etape, care se desfășoară pe parcursul anilor 2008-2018. Scopul acestuia este de a crea și pune la dispoziția celor interesați, cu precădere cercetătorilor din domeniul umanist și al științelor sociale, resurse lingvistice și tehnologii de prelucrare a limbajului, în toate formele lui de manifestare (textuală, vorbire, semne etc.). Ca purtător al conținutului cultural și al cunoașterii civilizațiilor, ca instrument de comunicare și componentă a identității naționale, precum și ca obiect de studiu, limbajul invită acum, din ce în ce mai imperios, la o abordare care să beneficieze de suportul tehnologiilor informaționale. CLARIN își propune crearea unei infrastructuri de cercetare care să facă posibilă partajarea și reutilizarea resurselor precum și prelucrarea lor prin instrumente specializate, la o scară care să justifice standardizarea. Totodată, CLARIN urmărește să transforme tehnologia actuală, extrem de fragmentată, precum și resursele existente, ori ce vor fi create în viitor, în servicii interoperabile și stabile, pe care utilizatorii să le poată accesa sau adapta după nevoi. Ambiția proiectului este de a crea o arhitectură orientată spre servicii care să faciliteze comunității de cercetători umaniști sau din domeniul științelor sociale obținerea de extensii și adaptări în orice manieră imaginabilă, pentru accesul la resurse, pentru prelucrări asupra lor, ori pentru consultații. Serviciile vor avea la bază o rețea de centre de mărime și tipuri diferite, distribuite în Europa.

CLARIN speră să reunească într-o comunitate specializată toate instituțiile din Europa care, într-un fel ori altul, dispun de resurse lingvistice în format scris, vorbit ori multimodal, ori de tehnologii lingvistice. De asemenea, inițiativa CLARIN nu poate fi considerată de succes dacă nu va reuși să atragă masa mare de utilizatori care fac uz de astfel de resurse ori tehnologii în cercetarea lor. Comunitatea CLARIN include actualmente ca membri mai mult de 130 de instituții din 32 de țări europene. Dintre acestea, 24 de țări au exprimat deja acordul de a cofinanța proiectul.

Beneficiile principale aduse de includerea limbii române între limbile proiectului CLARIN constau, printre altele, în reconsiderarea eforturilor de tehnologizare a limbii române prin prisma standardelor ce vor fi adoptate la nivelul întregii Europe, accesul la tehnologii moderne, posibil de adaptat și pentru prelucrarea limbii române, oferirea de



servicii de informare relative la tehnologii și colecții de resurse, mărirea gradului de utilizare a acestor resurse prin facilitarea accesului la ele a specialiștilor în științe umane și sociale și, nu în ultimul rând, mărirea vizibilității eforturilor de creare de tehnologii de procesare și de resurse specifice limbii române, o dată cu includerea lor în colecțiile proiectului, ce se așteaptă să fie larg accesate de cercetători. Se speră, totodată, ca feedback-ul oferit de utilizarea mai frecventă și în situații noi a tehnologiilor și resurselor să ducă și la îmbunătățirea calității acestora.

### 3. *ALPE*

ALPE este un meta-sistem care permite unui utilizator dispunând doar de minime abilități informatice să exploateze configurații de procesare a documentelor adnotate XML, deja create anterior, sau chiar să creeze altele noi. Generarea unei arhitecturi de procesare (workflow) se realizează ca un proces de navigare într-o ierarhie de scheme de adnotare XML (Cristea et al., 2006, Cristea și Pistol, 2008). Ierarhia este un graf direcționat în care nodurile reprezintă scheme de adnotare iar arcele sunt relații de subsumare. În acest context spunem că nodul A subsumă nodul B dacă:

- schema B conține toate elementele schemei A;
- schema B include cel puțin un element (tag sau atribut) ce nu este cuprins în schema A.

Direcția arcelor din graf este dată de relația de subsumare: de la nodul care subsumă către nodul subsumat. Un nod poate subsuma mai multe noduri și poate fi subsumat de mai multe. Considerând elementul rădăcină ca fiind adnotarea XML vidă (cuprinzând numai identificatorul de format XML), el subsumă toate celelalte noduri, ceea ce înseamnă că nu există noduri izolate.

Arcelor grafului li se pot atașa module de procesare lingvistică capabile să transforme un fișier ce corespunde formatului de intrare (nodului origine al arcului) într-un fișier corespunzând formatului de ieșire (nodului destinație al arcului). De notat că nu întotdeauna un arc, căruia îi corespunde așadar o relație de subsumare de scheme, are atașat un modul de procesare. Arcele cărora le sunt atașate minimum un modul de procesare se numesc arce de procesare (processing edges), iar cele care nu au nici un modul atașat se numesc arce purtătoare (carrier edges).

Pe un graf de acest tip, ALPE definește o serie de operații ce permit calculul automat al unor lanțuri de procesare. Aceste lanțuri de procesare (processing flows) sunt capabile să transforme automat un document dintr-un format în altul, dacă formatul de intrare și cel de ieșire corespund la două noduri ale ierarhiei și dacă modulele corespunzătoare arcelor de procesare sunt disponibile. Cât despre arcele purtătoare, există două moduri în care ele pot interveni într-un lanț: dacă, incorporate unui lanț de procesare, ele se combină cu alte arce în nodul de ieșire, atunci ele mixează informația din intrare cu cea provenită din celelalte arce, altfel ele blochează lanțul.

ALPE oferă două funcționalități de bază:

- Un utilizator poate îmbogăți o ierarhie ALPE deja existentă prin:

- oferirea unui nou format de adnotare, ceea ce duce la includerea unui nod nou în ierarhie și legarea lui automată de nodurile existente, astfel încât relația de subsumare să fie respectată;
- oferirea unui nou modul de procesare și a informațiilor privind formatele sale de intrare și ieșire (modalitatea de apel, resursele adiționale necesare, condiții de acces, etc.). Acest modul de procesare va fi automat integrat în ierarhia existentă ca arc de procesare.
- Un utilizator poate procesa un document, transmițându-l ierarhiei ca document de intrare și indicând un format de ieșire. ALPE calculează toate lanțurile de procesare posibile. Din acest evantai de soluții, utilizatorul poate alege pentru rulare efectivă pe acelea care corespund propriilor lui criterii de cost/eficiență.

O primă variantă a sistemului este deja implementată, fiind folosită pentru configurarea unui sistem de tip Întrebare-Răspuns și pentru procesarea lingvistică a documentelor într-un sistem de e-Learning. ALPE va fi disponibil ca serviciu web, prin intermediul lui utilizatorii urmând a fi capabili să creeze și utilizeze propriile ierarhii, sau să contribuie la dezvoltarea unei ierarhii globale. Se are în vedere și promovarea ALPE pentru a fi adoptat ca help-desk interactiv în CLARIN. O cale de dezvoltare atractivă o constituie adaptarea ALPE pentru a lucra în rețele de tip GRID, lucru care ar aduce importante îmbunătățiri relativ la viteza de procesare și operabilitate în condițiile măririi numărului de utilizatori și a dimensiunii ierarhiilor.

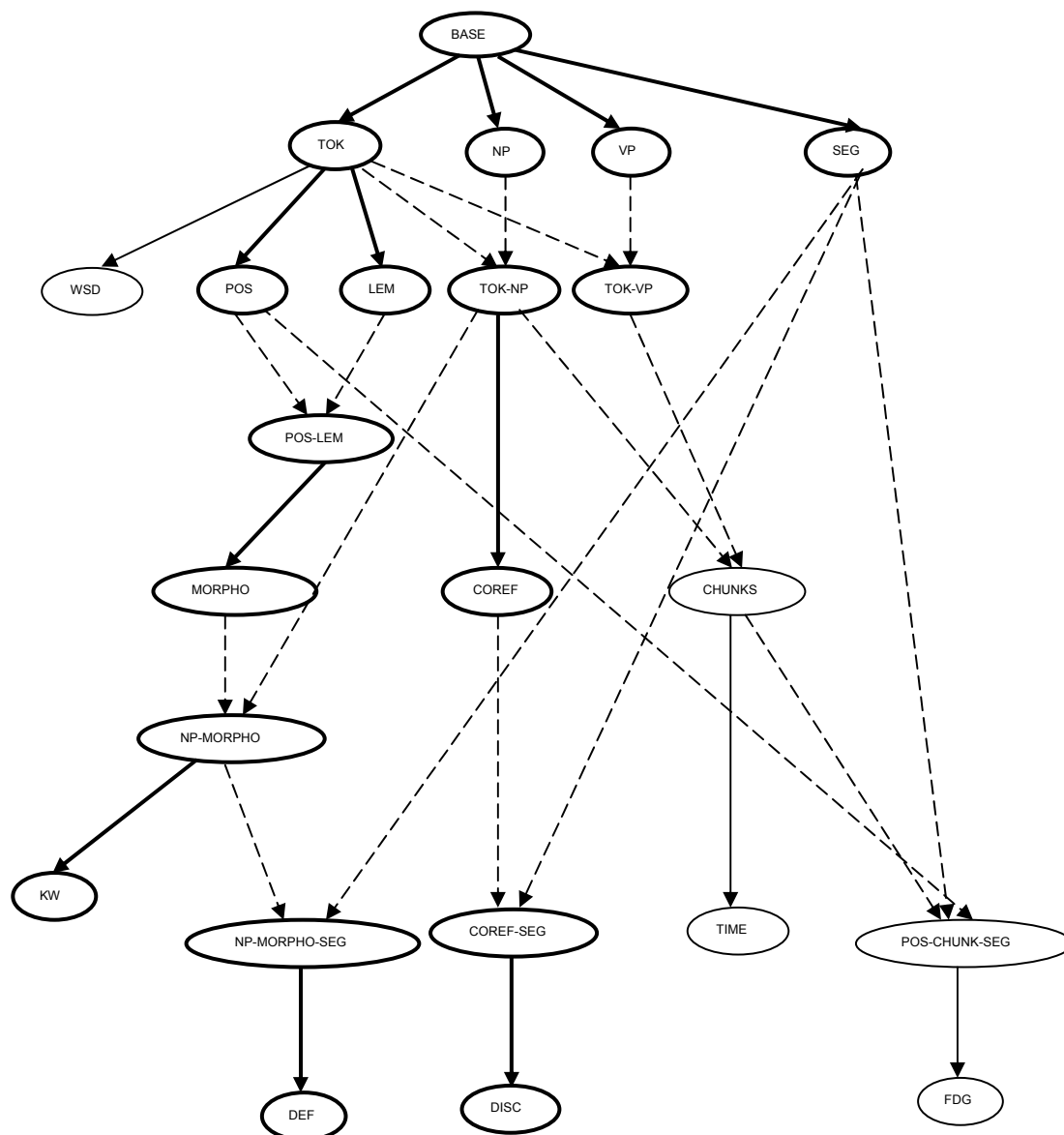
#### ***4. Ierarhia de resurse de procesare pentru limba română***

Întâlnirile ConsILR precedente (Forăscu et al., 2006; Pistol et al., 2007), participările cu succes la competiții dedicate limbii române (Iftene et al., 2008; Orasan et al., 2008) și includerea limbii române în proiecte de cercetare europene (Tufiș et al., 2004; Lemnitzer et al., 2007) indică atât mărirea interesului arătat limbii române de un grup tot mai numeros de cercetători, cât și existența deja a unui set important de instrumente de prelucrare dedicate limbii române. Colectarea acestor instrumente, facilitarea accesului la ele și a utilizării lor în diverse aplicații noi se integrează atât obiectivelor ConsILR cât și ale CLARIN. În această secțiune încercăm să dăm o caracterizare ca o ierarhie ALPE a unui set de formate XML dedicate aplicațiilor de prelucrare a limbajului natural, pe care le-am putut identifica în literatura dedicată limbii române.

Nodurile din Figura 1, cu excepția nodului rădăcină, trebuie înțelese ca reprezentând clase de formate de adnotare cu același conținut semantic. De exemplu, diferitele marcaje utilizate pentru elementele lexicale identifică o clasă de formate de adnotare denumită TOK. Definițiile formatelor din noduri sunt, pe scurt, următoarele:

- BASE: format cu marcaje minimale de început și sfârșit de document XML;
- TOK: clasă de formate ce marchează elementele lexicale de bază (cuvinte, unități de punctuație);
- NP: clasă de formate ce marchează grupurile nominale;
- VP: clasă de formate ce marchează grupurile verbale;

## LIMBA ROMÂNĂ ÎN PERSPECTIVA CLARIN



**Figura 1:** Ierarhia ALPE pentru limba română

- SEG: clasă de formate ce marchează fraze sau unități elementare de discurs (propoziții);
- WSD: clasă de formate ce marchează sensuri ale unităților lexicale;
- POS: clasă de formate ce marchează părțile de vorbire ale unităților lexicale;
- LEM: clasă de formate ce marchează formele de bază ale unităților lexicale (leme);
- TOK-NP: clasă de formate ce reunește TOK și NP;
- TOK-VP: clasă de formate ce reunește TOK și VP;
- POS-LEM: clasă de formate ce reunește POS și LEM;

- MORPHO: clasă de formate ce reunește informațiile morfo-sintactice;
- COREF: clasă de formate ce marchează lanțuri coreferențiale;
- CHUNKS: clasă de formate ce reunește NP și VP;
- NP-MORPHO: clasă de formate ce reunește TOK, NP și MORPHO;
- KW: clasă de formate ce marchează termeni (cuvinte cheie);
- NP-MORPHO-SEG: clasă de formate ce adaugă la NP-MORPHO și informații de segmentare;
- COREF-SEG: clasă de formate ce reunește COREF și SEG;
- TIME: clasă de formate ce adaugă marcaje pentru adnotarea temporală;
- POS-CHUNK-SEG: clasă de formate ce reunește POS, CHUNK și SEG;
- DEF: clasă de formate ce adaugă marcaje pentru definiții;
- DISC: clasă de formate ce adaugă marcaje pentru structura de discurs;
- FDG: clasă de formate ce adaugă marcaje pentru dependențele funcționale.

Acele îngroșate indică existența unuia sau a mai multor module de procesare corespunzătoare. O parte din modulele considerate în această ierarhie provin din:

- Serviciile web ale ICIA (Tufiș et al., 2007);
- Adnotatorul de expresii temporale (Forăscu și Solomon, 2004);
- Rezolvitoarele de anaforă AR-Engine și RARE (Cristea et al., 2002; Pavel et al., 2007);
- Parserul de discurs (Cristea et al., 2005);
- POS taggere (Tufiș și Dragomirescu, 2004);
- Dezambiguitorul semantic (Ion și Tufiș, 2004).

Alte detalii privind resursele de procesare disponibile pentru limba română pot fi găsite în (Cristea și Tufiș, 2002; Forăscu et al. 2007; Pistol et al., 2008).

Acele subțiri indică posibila existență a unor module de procesare corespunzătoare, dar indisponibile integrării în ierarhie în perspectiva imediată. Arcele marcate cu linii întrerupte sunt arce *carrier*. Nodurile îngroșate sunt noduri ce pot fi atinse de lanțuri de procesare din orice alt nod al ierarhiei. Nodurile marcate cu linie subțire nu pot fi atinse de lanțuri de procesare decât plecând din noduri subsumate lor.

În forma din Figura 1, ierarhia ar permite numeroase prelucrări semnificative, cum ar fi: adnotarea automată a structurii de discurs, marcarea definițiilor, marcarea unui text cu informație temporală, informație sintactică, precum și operații de combinare și simplificare a adnotărilor.

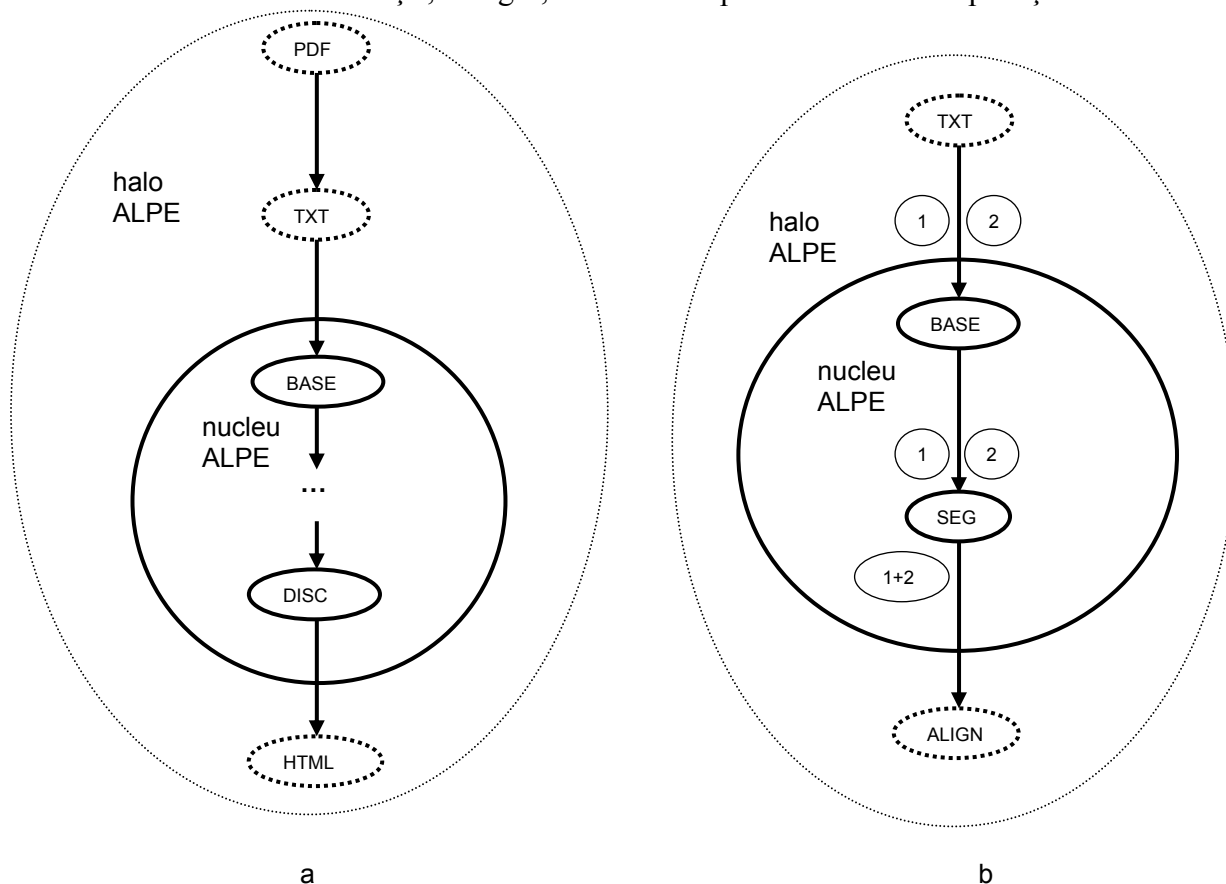
## **5. Interacțiunea cu o ierarhie ALPE**

Premizele care stau la baza interacționării cu o ierarhie ALPE sunt:

a. pe parcursul procesărilor, documentul de bază (*hub*) este neschimbat și doar adnotările XML aplicate lui suferă transformări;

b. o schemă XML (un nod al ierarhiei) trebuie înțeleasă ca purtătoarea unui mesaj de adnotare, cu semnificație cunoscută, care completează un document.

Premiza *a* vede spațiul nodurilor ierarhiei ca adnotări diferite aplicate aceluiași document. Această restricție, desigur, nu este compatibilă cu multe aplicații din NLP.



**Figura 2:** Conectarea nucleului ALPE cu haloul

Exemplele în care documentul de bază se modifică pe parcursul procesării sunt numeroase: transformarea textelor din formate diferite de XML în formate XML (de exemplu, conversoare pdf-txt, html-xml), programe care „impurifică” textul inițial (de exemplu, prin adăugarea de blancuri ori newline-uri), POS-taggere ori lematizoare care primesc un text în intrare și întorc ieșiri tabulare (câte un element lexical pe linie, de exemplu), programe care modifică substanțial textul din intrare (rezumatoare, traducătoare automate etc.), ori care convertesc informația dintr-un mediu în altul (conversoare *text-to-speech*, *speech-to-text* etc.).

Acomodarea diversității imense de aplicații în universul modelului nostru presupune percepția ierarhiei altfel decât unicul spațiu al prelucrărilor. În jurul acestuia trebuie înțeles că gravitează, ca un halo, o serie întreagă de standarde care nu reprezintă notații XML aplicate unui document, ori dacă sunt adnotări XML, atunci ele nu se aplică asupra aceluiași document. Dacă nodurile din afara nucleului ALPE sunt legate de

noduri ale nucleului, înseamnă că există procese capabile de astfel de transformări între nodurile corespunzătoare.

Spre exemplu, Figura 2a arată situația unei aplicații de rezumare în care textul de intrare are formatul pdf, iar rezumatul este postat ca document html pe Web.

În cazul unui aliniator de texte (Figura 2b), situația este diferită: două documente separate sunt prelucrate (în paralel ori serial) de la un format TXT (exterior nucleului ALPE) până la un format XML în care sunt marcate frazele (segment), după care ele sunt date unui aliniator care, ieșind din nou în afara nucleului ALPE, produce alinierea. Ambele fișiere de intrare sunt prelucrate de același lanț de procesare ALPE, rulat pe fiecare fișier în parte. Rezultatul celor două lanțuri de procesare este combinat în afara ierarhiei ALPE de către aliniator. De notat că aceasta poate fi tot un fișier XML.

Alte exemple de aplicații ce presupune o migrare între nucleul și haloul ALPE sunt date de generatorul text-voce și analizorul voce-text. Înregistrările sonore ce au rol de ieșire, respectiv intrare în aceste aplicații se regăsesc în noduri din afara nucleului ALPE, pentru că se referă la formate diferite de cele prelucrate în nucleu (textuale). Ele diferă însă atât prin format (ce nu mai este XML) cât și ca tip de document, trecând de la document scris la înregistrare audio. ALPE poate astfel funcționa ca suport pentru aplicații multimodale ce includ etape de procesare pe text.

Premiza *b* prevede că unei ierarhii ALPE îi corespunde o colecție de standarde de notații XML. Diversitatea extravagantă dată de un spațiu al numelor neconstrâns de un standard la care să adere majoritatea utilizatorilor de XML în aplicațiile de prelucrări lingvistice trebuie însă să găsească o expresie în model. Soluția constă în a vedea un nod ALPE ca pe un nor de notații care poartă aceeași semnificație semantică. De exemplu, un cuvânt (*token*) poate fi notat ca un element XML TOK, sau WRD, sau W. O marcă morfologică se poate exprima printr-o diversitate de notații, de la includerea tuturor informațiilor într-un singur atribut (ca în cazul tag-setului MULTTEXT) până la separarea atributelor morfologice într-un set de etichete (ca în cazul tagsetului XCES-EAGLES). Este clar că o procesare similară, indiferent de formatul de intrare a documentelor, duce la economisirea de resurse de calcul și permite uniformizarea lanțurilor de prelucrări. Este, de aceea, de dorit ca nucleul ALPE să integreze notații universale acceptate, dacă se poate, consolidate ca standarde, fără însă ca celelalte notații să fie interzise uzului utilizatorilor. Programe de conversie (*wrappere*) vor asigura compatibilitatea intrărilor și ieșirilor din sistem cu cerințele utilizatorilor, prelucrarea în nucleul sistemului realizându-se în conformitate cu standardele acceptate.

## 6. Concluzii

Integrarea resurselor dedicate procesării limbii române într-o ierarhie ALPE promite atât sporirea vizibilității și utilizării acestor resurse cât și mărirea eforturilor dedicate prelucrării electronice a limbii române. Posibilitatea de a utiliza și compara resursele existente poate duce la crearea de sisteme tot mai complexe de prelucrare și la ideea îmbunătățirii modulelor existente ce au performanțe scăzute.

Dezvoltarea ierarhiei ALPE urmează să fie făcută, prin colaborare, în colectivele implicate în procesarea limbii române, în marea lor majoritate agreând întâlnirile ConsILR. Succesul acestui efort va contribui atât la materializarea unuia din scopurile

originare ale ConsILR-ului, dar și la integrarea cercetărilor românești dedicate domeniului ingineriei lingvistice în consorțiul exigent și elevat al CLARIN. Cum independența de limbă este un deziderat care se apropie tot mai mult de realitate în realizarea de aplicații în zilele noastre, devine posibil ca propunerea unei ierarhii ALPE în contextul unei anumite limbi, cum e româna, să servească drept iterația zero într-un ciclu de dezvoltare a unei ierarhii generale de instrumente de procesare, care să servească toate limbile implicate în proiect.

### Referințe bibliografice

- Cristea D., Forăscu C., Pistol I. (2006). Requirements-Driven Automatic Configuration of Natural Language Applications. In Bernadette Sharp (Ed.): *Proceedings of the 3rd International Workshop on Natural Language Understanding and Cognitive Science - NLUCS 2006*, Paphos, INSTICC Press, Portugal. ISBN: 972-8865-50-3.
- Cristea, D., Pistol, I. (2008). Managing Language Resources and Tools Using a Hierarchy of Annotation Schemas. In *Proceedings of the Workshop on Sustainability of Language Resources*, LREC-2008, Marakesh.
- Cristea, D., Postolache, O., Dima, G.E., Barbu, C. (2002). AR-Engine – a framework for unrestricted coreference resolution. *Proceedings of Language Resources and Evaluation Conference - LREC 2002*, Las Palmas, vol. VI, 2000-2007.
- Cristea D., Postolache O., Pistol I. (2005). Summarisation through Discourse Structure, In Alexander Gelbukh (Ed.): *Computational Linguistics and Intelligent Text Processing*, 6th International Conference CICLing 2005, Mexico City, February 2005, Springer LNCS, vol. 3406, ISBN 3-540-24523-5, pp. 632-644.
- Cristea, D.; Tufiș, D. (2002): Resurse lingvistice românești și tehnologii informatice aplicate limbii române. In Ichim, O și Olariu F.-T. (eds.): *Identitatea limbii și literaturii române în perspectiva globalizării*, Academia Română, Institutul de Filologie Română „A. Philippide”, Editura Trinitas, Iași, pp. 211-234.
- Cunningham H., D. Maynard, K. Bontcheva, V. Tablan. (2002): GATE: A framework and graphical development environment for robust NLP tools and applications. In *Proceedings of the 40th Anniversary Meeting of the ACL (ACL'02)*. Philadelphia, US.
- Ferrucci D. și Lally A. (2004): UIMA: an architectural approach to unstructured information processing in the corporate research environment, *Natural Language Engineering* 10, No. 3-4, 327-348.
- Forăscu C., Cristea D., Tufiș D. (eds) (2007) *Lucrările atelierului Resurse lingvistice și instrumente pentru prelucrarea limbii române* Iași, Editura Universitatii “Al.I. Cuza” Iasi, România, ISBN 978-973-703-208-9.
- Forăscu C., Solomon D. (2004). Towards a Time Tagger for Romanian. In *Proceedings of the ESSLLI Student Session*, Nancy, France.
- Iftene A., Pistol I., Trandabăț D. (2008). UAIC Participation at QA@CLEF2008. In *Proceedings of the CLEF 2008 Workshop*. 17-19 September. Aarhus, Denmark.

- Ion R., Tufiș D. (2004). Multilingual Word Sense Disambiguation Using Aligned Wordnets. In *Romanian Journal on Information Science and Technology*, Dan Tufiș (ed.) Special Issue on BalkaNet, Romanian Academy, vol. 7, no. 2-3, pp. 198-214, ISSN 1453-8245.
- Lemnitzer, L., Vertan, C., Killing, A., Ivanov Simov, K., Evans, D., Cristea, D., Monachesi, P. (2007): Improving the Search for Learning Objects with Keywords and Ontologies. In *Creating New Learning Experiences on a Global Scale*, EC-TEL 2007, Lecture Notes in Computer Science, vol. 4753/2007, pp. 202-216, ISBN 978-3-540-75194-6.
- Orasan, C., Cristea, D., Mitkov, R., Branco, A. (2008). Anaphora Resolution Exercise - An Overview. In *Proceedings of LREC-2008*. Marakesh.
- Pavel, G., Postolache, O., Pistol, I., Cristea, D. (2007): Rezolutia anaforei pentru limba română. In C. Forăscu, D. Tufiș, D. Cristea (eds.): *Resurse lingvistice și instrumente pentru prelucrarea limbii române*, Editura Universitatii "Al.I. Cuza" Iasi, România, ISBN 978-973-703-208-9.
- Pistol I.C., Cristea D., Tufiș D. (eds) (2008) *Lucrările atelierului Resurse lingvistice și instrumente pentru prelucrarea limbii române* Iași, 14-15 decembrie 2007, Editura Universitatii "Al.I. Cuza" Iasi, România, ISBN 978-973-703-208-9.
- Tufis, D., Cristea, D., Stamou, S. (2004): BalkaNet: Aims, Methods, Results and Perspectives. A General Overview. In *Romanian Journal of Information Science and Technology*, Romanian Academy, Bucharest, Romania, special issue on Balkanet, July, pp. 9–43, ISSN 1453-8245.
- Tufis D., Dragomirescu L. (2004). Tiered Tagging Revisited. In *Proceedings of the 4th LREC Conference*, Lisabona.
- Tufiș D., Ion R., Ceaușu A., Ștefănescu D. (2007). Servicii web lingvistice ale ICIA. *Lucrările atelierului "Resurse lingvistice și instrumente pentru prelucrarea limbii române"* Iași, 14-15 decembrie 2007, Editura Universitatii "Al.I. Cuza" Iasi, România, ISBN 978-973-703-208-9.
- Váradi T., Krauwer S., Wittenburg P., Wynne M. and Koskenniemi K. (2008): CLARIN: Common Language Resources and Technology Infrastructure. In *Proceedings of Language Resources and Evaluation Conference - LREC 2008*. Marakesh.



## PARSAREA ARBORILOR DE SENSURI ȘI SEGMENTAREA LA DEFINIȚII ÎN DICȚIONARUL TEZAUR eDTLR

NECULAI CURTEANU<sup>1</sup>, ALEX MORUZ<sup>1,2</sup>, DIANA TRANDABĂȚ<sup>1,2</sup>,  
CECILIA BOLEA<sup>1</sup>, MĂDĂLINA SPĂTARU<sup>1,2</sup>, MARIA HUSARCIUC<sup>1,2</sup>

<sup>1</sup> *Institutul de Informatică Teoretică Iași, Academia Română*

<sup>2</sup> *Facultatea de Informatică, Universitatea "Al. I. Cuza", Iași;*

[curteanu@iit.tuiasi.ro](mailto:curteanu@iit.tuiasi.ro), [mmoruz@iit.tuiasi.ro](mailto:mmoruz@iit.tuiasi.ro), [dtrandabat@info.uaic.ro](mailto:dtrandabat@info.uaic.ro)

### Rezumat

Lucrarea de față prezintă o nouă metodă de parsare a textului de dicționar, bazată pe *configurații* de tip SCD (Segmentare-Coeziune-Dependență), utilizate pentru transformarea eficientă a unui tezaur într-un lexicon structurat. Astfel, strategia dezvoltată reunește două configurații diferite de parsare: una care identifică și extrage, pentru fiecare intrare din dicționar, arborele specific de sensuri (ierarhia sensurilor principale și secundare), și o altă configurație care parsează fiecare nod din arborele de sensuri cu scopul de a clasifica *definițiile* aceluși sens din dicționar. Spre deosebire de metodele standard de parsare a textului de dicționar, în care *toate* câmpurile unei intrări de dicționar sunt analizate *secvențial*, noua metodă reușește detașarea procesului de construire a arborelui de sensuri (*prima configurație* SCD) de procesul parsării la definițiile sensurilor (*a doua configurație* SCD). Separarea celor două procese se face în principal prin selectarea *breadth-first* a tuturor *marcherilor* la sensuri, urmată de analiza lor *depth-first*, în fiecare intrare de dicționar. Pentru clasificarea tipurilor de definiții la sensurile din Dicționarul Tezaur al Limbii Române și parsarea lor a fost realizată o modelare lexical-semantică a definițiilor, iar strategia de parsare propusă se aplică în cei doi pași mai sus menționați. Sunt discutate analiza erorilor și rezultatele parsării la sensuri și definiții, precum și posibilitatea aplicării parserului pe un dicționar tezaur dintr-o altă limbă.

### 1. Introducere

Parsarea unui dicționar presupune transformarea intrărilor ce conțin text sub formă de glosă, într-un format indexabil. Astfel, fiecare intrare de dicționar este transpusă într-o structură complexă, care conține atât sensurile definite, precum și descrieri detaliate ale formei intrării, cu referire la ortografie, morfologie, fonetică, etimologie, uz etc. Scopul acestei lucrări este introducerea unei metode noi de parsare a unui tezaur, bazată pe configurații de tip Segmentare-Coeziune-Dependență (SCD), (Curteanu, 2006). Spre deosebire de orientările standard în parsarea intrărilor de dicționar (Neff and Boguraev, 1989), de exemplu sistemul *LexParse* (Hauser and Storrer, 1993; Kammerer, 2000; Lemnitzer and Kunze, 2005) sau gramaticile lexicografice (Curteanu and Amihăesei, 2004; Tufiș *et al.*, 1999), metoda folosită de noi detașează complet procesul de *construire a arborilor de sensuri* de procesul *parsării definițiilor* la sensuri.

O *configurație* SCD are următoarele componente:

- Un set de clase de marcheri: un *marcher* reprezintă o graniță pentru o categorie lingvistică specifică;
- O ierarhie de tip arbore, care stabilește dependențele dintre clasele de marcheri;
- Un algoritm de parsare, care execută următorii pași: recunoașterea marcherilor, identificarea structurilor dintre doi marcheri și clasificarea acestor structuri ținând cont de ierarhia claselor de marcheri. Algoritmul poate fi aplicat pe diferite clase sau ierarhii de marcheri, depinzând strict de *semantica* textului ce urmează a fi parsat.

Prima expunere a ideilor de bază ale parsării *Dicționarului Tezaur al Limbii Române* (referit în continuare prin DTLR) cu o metodă derivată din strategia SCD a fost făcută în (Curteanu *et al.*, 2007), unde este schițat algoritmul de parsare DSSD (Dictionary Sense Segmentation and Dependency) cu extragerea *ab initio* a marcherilor de sensuri din intrare, în prima etapă, și parsarea definițiilor la sensuri într-o a doua etapă. Este inclusă o primă formă a grafului de dependențe între clasele de marcheri la sensuri. În (Curteanu *et al.*; 2008) sunt prezentate rezultatele teoretice și de implementare ale parsării DTLR. Lucrarea de față aduce următoarele noutăți: **(1)** Folosim o *primă configurație* SCD (SCD-*config1*) pentru a extrage arborele de sensuri din DTLR. SCD-*config1* corespunde algoritmului de parsare DSSD (Curteanu *et al.*; 2008), care obține arborii de sensuri din DTLR cu o precizie de 91.18%. **(2)** Pentru a rafina conținutul lexical-semantic al sensurilor din DTLR, este necesară coborârea analizei în sensurile secundare, la nivelul *definițiilor* DTLR, care constituie întinderea de text situată între două noduri consecutive din arborele de sensuri al unei intrări. Lucrarea de față prezintă și modelarea definițiilor din DTLR. Această a doua etapă a parsării DTLR reprezintă o nouă configurație SCD, notată SCD-*config2*, ce constă dintr-un set specific de clase de marcheri pentru *segmentarea la definiții*. Rezultatul final al parsării va fi aplicarea în cascadă a SCD-*config1* și SCD-*config2*, în această ordine.

Configurațiile SCD propuse de noi sunt prezentate în următoarele secțiuni, și aplicate pentru parsarea DTLR. Astfel, Secțiunea 2 prezintă parsarea la arborii de sensuri (Sense Tree Parsing), în timp ce Secțiunea 3 expune modelarea lexical-semantică pentru diferitele tipuri de definiții care pot fi găsite într-un nod din arborele de sensuri. Secțiunea 4 analizează rezultatele parsării utilizând cele două configurații SCD, urmând apoi prezentarea concluziilor și dezvoltările pe care le avem în vedere.

## **2. Parsarea arborilor de sensuri din intrările de dicționar**

Parsarea arborelui de sensuri folosind configurații SCD a fost inspirată de comparația între clasele de marcheri de sens din DTLR și clasele de marcheri SCD pentru parsarea textului general (Curteanu, 2006).

*Clasele de marcheri* folosite în procesul de *parsare a arborilor de sensuri* din intrările DTLR sunt descrise mai jos:

Clasa de marcheri ce conține *majuscule* (**A.**, **B.**, etc.) reprezintă *nivelul cel mai de sus* al ierarhiei de sensuri a marcherilor DTLR (Fig. 1) pentru orice intrare de dicționar dată. Când apare, acest marcher desemnează înțelesul cel mai general și desemnează un *sens principal* al cuvântului definit. Dacă acest nivel are doar un element de acest tip, atunci marcherul corespunzător lipsește, el fiind înlocuit de un marcher de nivel inferior.

PARSAREA ARBORILOR DE SENSURI ȘI SEGMENTAREA  
LA DEFINIȚII ÎN DICȚIONARUL TEZAUR eDTLR

Clasa de marcheri ce conține *cifre romane* (I., II., etc.) reprezintă *al doilea nivel* de sens pentru o intrare DTLR. Acest nivel este subsumat de un marker de tip *majusculă*, dacă acesta există; dacă majuscula nu există (sau nu este reprezentată în mod explicit), markerul de tip *cifră romană* apare pe nivelul cel mai de sus al arborelui de sensuri. Dacă intrarea lexicală are doar un sens pentru acest nivel al analizei, markerul nu este reprezentat în mod explicit, el fiind înlocuit de un marker de nivel inferior.

Clasa de marcheri ce conține *cifre arabe* (1., 2., etc.) reprezintă *al treilea nivel* de sens pentru o intrare DTLR. Acest nivel este subsumat de un marker de tip *cifră romană*, dacă există; dacă acesta nu este reprezentat în mod explicit, este subsumat de primul marker explicit de nivel superior. Dacă intrarea are doar un sens pentru acest nivel al analizei, markerul nu este reprezentat în mod explicit.

Aceste prime *trei niveluri* codifică *sensurile principale* ale unei intrări DTLR.

*Rombul plin* reprezintă *al patrulea nivel* de sens și este folosit pentru a enumera *sensurile secundare* ale unei intrări din DTLR. Acest nivel este subsumat de orice marker de sens de nivel superior (oricare dintre markerii unui sens principal).

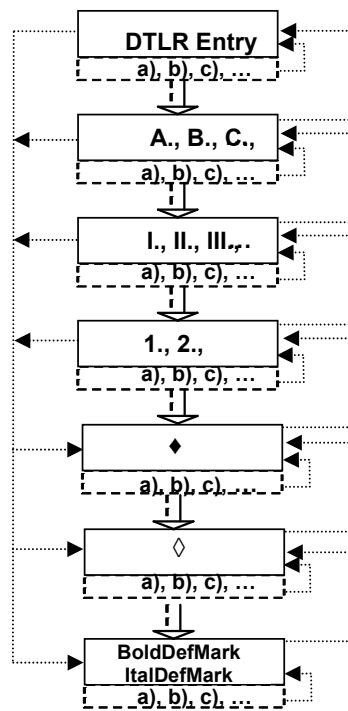


Figura 1: Ierarhia marcherilor DTLR

Markerul *romb gol* reprezintă *al cincilea nivel* al analizei sensurilor și este folosită pentru a enumera expresii pentru un *sub-sens secundar* dat. Acest nivel este subsumat de un marker de tip *romb plin* sau de orice marker de sens principal.

Markerii de tip *BoldDefMark* și *ItalDefMark* sunt delimitatorii definițiilor de tip *BoldDef*, respectiv *ItalDef* ( v. Secțiunea 3).

Markerii de tip *litere mici* (a), b), etc.) nu reprezintă, în realitate, o clasă distinctă de marcheri de sens, ci mai curând o *procedură* folosită pentru rafinarea, prin *enumerare*

*literală*, a unui sens sau sub-sens. Un marker de acest tip nu reprezintă un nivel specific în ierarhia claselor de markeri, deoarece aparține nivelului de sens al părintelui. Regulile de bază ale procedurii de *enumerare literală* în DTLR sunt următoarele: **(i)** se asociază cu nivelul ierarhiei claselor de markeri în cadrul căreia apare și **(ii)** poate îngloba sub-sensuri de nivel mai mic (decât nivelul nodului părinte).

Fig. 1 prezintă ierarhia claselor de markeri din DTLR. *Săgețile* cu linie întreruptă indică faptul că orice nivel de sens este opțional; *săgețile continue* indică ierarhia claselor de markeri. Datorită caracteristicilor sale specifice, *enumerarea literală* este ilustrată pe un nivel *atașat* nivelului căruia îi este asociat în ierarhie.

Exemplul de mai jos prezintă rezultatul parsării arborelui de sensuri al unei intrări DTLR. Se observă că exemplul de intrare prezentat (*VENIT*<sup>2</sup>) reprezintă numai secvențele markerilor de sens din DTLR (în dicționar această intrare întinzându-se pe mai mult de două pagini):

```
<entry>
  <hw><VENIT2, -Ă </hw>
  <pos>adj. </pos>
  <senses>
    <definition>...</definition>
    <marker>1.
    <definition>...</definition>
    </marker>
    <marker>2.
    <definition>...</definition>
    <marker>∅
      <marker> a)
      <definition>...</definition>
      </marker>
      <marker> b)
      <definition>...</definition>
      </marker>
      <marker> c)
      <definition>...</definition>
      </marker>
    </marker>
    <marker>∅
      <marker> a)
      <definition>...</definition>
      </marker>
      <marker> b)
      <definition>...</definition>
      </marker>
    </marker>
  </senses>
</entry>
```

### 3. Modelarea sensurilor din dicționar cu tipuri specifice de definiții

Analiza conținutului semantic al sensurilor în DTLR este realizată cu o *a doua configurație* SCD, bazată pe un set de *markeri de definiții*. Markerii sunt folosiți pentru a crea un număr finit de șabloane care delimitează subsensuri prin informații lingvistice distinct detaliate, prin trăsături morfologice, sintactice, semantice sau

pragmatice și prin aspecte lexicale, etimologice, diacronice și prozodice. Analizând intrările DTLR, au fost găsite următoarele tipuri de definiții:

1. *MorfDefs* – definiții morfologice;
2. *RegDefs* – definiții scrise cu font *regular*;
3. *BoldDefs* – definiții scrise cu *bold*;
4. *ItalDefs* – definiții scrise cu *italic*;
5. *SpecDefs* – definiții ce conțin specificații;
6. *SpSpecDefs* - definiții scrise cu litere spațiate, ce conțin *anumite* specificații;
7. *DefExems* – exemple la definiții, cu rolul de a întregi înțelesurile unei definiții.

Tipurile de definiții propuse aici primesc roluri funcționale specifice în descrierea sensurilor principale, secundare, sau de granularitate semantică mai fină (Curteanu *et al.*; 2008). Se pot distinge două taxonomii ale definițiilor din DTLR. Prima conține următoarele clase:

**(obli)** *definiții obligatorii*, care conține, de exemplu, *MorfDef*-uri și, pentru fiecare sens din DTLR, *una* din următoarele trei definiții: *RegDef*, *BoldDef* sau *ItalDef*. Nu există nici o intrare de dicționar care să nu conțină *MorfDef* și (cel puțin) *una* dintre definițiile ce aparțin mulțimii {*RegDef*, *BoldDef*, *ItalDef*}.

**(opti)** *definiții optionale*, de exemplu *SpecDefs*, *SpSpecDefs* și *DefExems*, care pot să apară ca *modificatori* sau *specificatori* în fața unei definiții obligatorii.

Cealaltă taxonomie împarte definițiile în:

**(auto)** *definiții autonome*, care sunt *RegDef*, *BoldDef* și *ItalDef*, aceste definiții având un rol de sine stătător în introducerea sensurilor din DTLR;

**(cont)** *definiții contingente*, de exemplu *MorfDefs*, *SpecDefs*, *SpSpecDefs* și *DefExems*, care nu pot fi folosite independent, având înțeles doar în *contextul* altor definiții DTLR.

*MorfDef* apare obligatoriu în rădăcina oricărei intrări din DTLR, fiind moștenită la nivelurile inferioare ale arborelui de sensuri. *SpecDefs*, *SpSpecDefs* și *DefExems* sunt definiții *contingente*, deoarece ele nu pot defini un (sub)sens în mod autonom, ci numai ca instrumente auxiliare de modificare a altor definiții autonome sau contingente.

### 3.1. Definițiile morfologice – *MorfDefs*

*Definițiile morfologice* (*MorfDefs*) sunt formate din una sau mai multe etichete care descriu categorii morfologice la diferite niveluri ale arborelui de sensuri. Primul element într-o intrare de dicționar, după cuvântul titlu, este un *MorfDef* complex, o listă de *MorfDef*-uri care detaliază toate categoriile morfologice posibile pentru cuvântul de intrare. Cu cât sensurile devin mai rafinate, cu atât *MorfDef*-urile ulterioare devin mai specifice (sub-liste ale *MorfDef*-ului complex), până când ajung să desemneze o singură categorie morfologică. Dacă definiția unui sens nou nu conține și un *MorfDef*, atunci definiția morfologică se moștenește de la primul sens regent care are un *MorfDef*. În cele ce urmează, exemplele de un anumit tip de definiții DTLR sunt evidențiate cu gri:

**VERZIȘÓR, -OĂRĂ** adj., subst. **I. Adj.** Diminutiv al lui *v e r d e* (**I 1**)... **II. Subst. 1. S. m.** (La pl.) Corp de trupă al cavaleriei... **2. S. m. și f.** (lht.; prin Munt.) Boiștean... **3. S. n.** (Prin Mold.; în forma *verdișor*) Rachiu cu mentă... **4. S. f.** (Regional) Varietate de struguri... **5. S. n.** (Familiar) Bancnotă de culoare verde...

Expresia regulată care descrie *MorfDef*-urile este:

$(x)^+, x \in \{\text{"subst."}, \text{"adj."}, \text{etc.}\}^1$

### 3.2. Definiții de tip regular – *RegDefs*

Definițiile scrise cu font *regular* (*RegDefs*) reprezintă cel mai frecvent instrument lingvistic folosit în DTLR pentru a descrie sensuri. *RegDef* corespunde glosei cuvântului de intrare sau unor sintagme care îl conțin, reprezentând descrierea standard a sensurilor în majoritatea dicționarilor. Expresia regulată de mai jos descrie cea mai generală formă a unei definiții *RegDef*.

$(([A-Z] | [a-z]) + (\ . + \ ))^* +$

O definiție de tip *RegDef* poate să apară în rădăcina intrării, în sensurile principale sau secundare, poate fi *moștenită* în aceste sensuri (cum ar fi sensul **I.** din **VENIRE**) sau poate face parte, ca *explicație*, din corpul altor două tipuri de definiții autonome din DTLR: *BoldDef* și *ItalDef*, ca în exemplul de *BoldDef* “**Bun venit**” sau de *ItalDef* “*Venit național*” de mai jos.

**VENIRE** s. f. Acțiunea de a v e n i și rezultatul ei. **I. 1.** Deplasare către cineva sau către ceva; parcurgere a unui traseu pentru a ajunge la un anumit loc,...

◇ *Ex p r.* **Bun venit** = formulă de salut prin care se exprimă mulțumirea în legătură cu sosirea, cu prezența cuiva.

◇ *Venit național* = parte a produsului economiei naționale dintr-o perioadă de timp, care rămîne după...

### 3.3. Definiții de tip bold – *BoldDefs*

O definiție de tip *BoldDef* este folosită cu scopul de a explica sensul unei *sintagme* sau al unei *exprimări* specifice; expresia este scrisă cu litere îngroșate, urmată de un *separator BoldDef* (în general “=”<sup>2</sup>) și de un *RegDef*. De obicei, *BoldDef*-urile rafinează subsensuri specifice, cum ar fi sensurile secundare introduse în DTLR prin ♦ și ◇. Expresia regulată care descrie forma generală a unui *BoldDef* este dată mai jos:

$(\text{bold}(\ . + )^* + (\text{separator}) (\text{RegDef}))$

◇ **A semăna în verde** = a semăna imediat după arat, cînd arătura este încă proaspătă. ... **A ara în verde** = a ara un pămînt care este încă jilav. ...

Există situații în care *BoldDef*-urile pot să apară în sensurile principale, inclusiv pe *nivelul-rădăcină* al unei intrări lexicale din DTLR. Un *BoldDef* poate fi foarte complex, conținând numeroase variante ale expresiei marcate cu caractere bold.

**3. A se duce** (sau **a merge**, **a se lăți**, învechit și regional, **a ieși**) **vestea** (cuiva, a ceva, de ceva etc.) sau **a i se duce** (ori **a-i merge**, **a i se lăți**, învechit, rar, **a i se ridica**, regional, **a-i ieși** cuiva) **vestea**, **a-i merge** (sau **a i se duce** cuiva) **vestea** și **povestea**, (învechit și regional) **a ieși veste** (de cineva sau de ceva) = a deveni foarte bine cunoscut, a i se duce faima;...,

### 3.4. Definițiile de tip italic – *ItalDefs*

*ItalDef*-urile sunt similare din punct de vedere sintactic *BoldDef*-urilor, dar sunt diferite din punct de vedere semantic, deoarece ele descriu în general locații, spre deosebire

<sup>1</sup> Lista posibilelor elemente morfologice este evident finită.

<sup>2</sup> Uneori, separatorul “=” este înlocuit cu expresii echivalente, cum ar fi “vezi”, “v.”, “se spune”, etc., și introduce o relație de echivalență semantică între expresia din stînga și secvența din dreapta.

de *BoldDef*-uri care descriu expresii. Partea de definiție ce conține colocația este codificată cu caractere cursive (italice).

(*italic*(.+)\*)+ (separator) (RegDef)

5. *Verde antic* = matostat. ...

**VERZÉR** subst. (Regional; în sintagma) *Verzerul tilegii* = schimbătoare la roțile plugului...

### 3.5. Definiții de specificare – *SpecDefs*

*SpecDef*-urile sunt definiții *contingente* scrise cu font *regular* și cuprinse, în general, între paranteze. Multe dintre ele sunt abrevieri, cuvinte sau expresii rezervate care denotă diferite contexte de utilizare ale intrării DTLR, cum ar fi: “(Regional)”, “(Argou)”, “(Fam.[iliar])” etc. Uneori *SpecDef*-urile nu apar între paranteze, dar acest lucru se întâmplă numai în cazurile în care acestea reprezintă cuvinte rezervate sau abrevieri. Expresia regulată care recunoaște acest tip de definiție este:

\( ([a-z] | [A-Z])+ \) | x; x ∈ {abrevieri}

*SpecDef*-urile sunt folosite la orice nivel în arborele de sensuri și au ca scop specificarea, ’modificarea’ definițiilor, ca în exemplele ce urmează.

(1) În rădăcina intrării de dicționar, imediat după *MorfDef*:

**VENIÁL**, -Ă adj. (Livresc; despre păcate<sup>2</sup>, greșeli etc.) Care poate fi iertat (de Biserică); ușor, fără importanță...

(2) În rădăcina unui sens principal:

2. (Învechit și regional; despre lichide, substanțe etc.) Veninos (2). ...

(3) În rădăcina unui sens secundar:

♦ F i g. (Despre oameni) Rău (A I 1); dușmănos; (despre manifestări, stări, acțiuni etc. ale oamenilor) care trădează răutate (I 1),...

### 3.6. Definiții de specificare scrise spațiat – *SpSpecDefs*

Un alt tip de definiție contingentă este *SpSpecDef*, care precizează mai multe trăsături standard. *SpSpecDef* se scrie cu litere spațiate și conține elemente dintr-o listă prestabilită de abrevieri, având următoarea formă:

(( [A-Z] | [a-z] ) ) +

*SpSpecDef* poate să apară la toate nivelurile de sensuri din DTLR, uneori împreună cu alte definiții contingente, ca în exemplul de mai jos:

2. **T r a n z. și r e f l. F i g.** A (se) amărî, a (se) supăra, a (se) necăji, a (se) mînia. *A sa prea iubită inimă ș-a veninat.* PANN, E. II, 94/18.

Unele *referințe externe* (către sensuri din alte intrări DTLR) sunt scrise tot cu font spațiat, din acest motiv trebuie verificat întotdeauna dacă un cuvânt scris spațiat face parte din lista prestabilită de abrevieri care formează *SpSpecDef*-uri sau nu.

Plantă erbacee din familia scrofulariacee, cu florile albe sau trandafirii, care crește în locuri umede sau mlăștinoase și care este folosită în medicină pentru proprietățile ei iritante și purgative; avrămeasă, (regional) milostivă (v. m i l o s t i v III 2), potroacă1 (4), mila-Domnului (v. m i l ă I 6) ( *Gratiola officinalis*). Cf. hem 2182, conv. lit. xxiii, 1060, brandza, fl. 349, damé, t. 188, barcianu, jahresber. viii, 101,...

### 3.7. Exemplificări de definiții – DefExems

Definițiile *autonome* pot primi unul sau mai multe *exemple* din surse bibliografice referite prin *sigle* sau create de către autorii dicționarului. *DefExem*-urile au rolul de a rafina sensul definițiilor autonome și a tuturor sensurilor mai generale decât ele (sensuri secundare și principale). O secvență de *DefExem*-uri, fiecare fiind urmat de o *siglă*, este următoarea:

**A intra în viață = a)** (despre oameni; și în forma **a pași în viață**) a începe să se confrunte cu realitatea. *Cum a intrat el în viață? Cât amor de drept și bine, Cât sinceră frăție adusese el cu sine?* EMINESCU, O. I, 53. ...; **b)** (rar) a începe să activeze, să funcționeze. *Guvernul cel nou... va intra în luna lui martiu în viață.* VASICI, ap. BARIȚIU, C. II, 47.

## 4. Modelarea sensurilor din dicționar cu tipuri specifice de definiții

### 4.1. Analiza arborelui de parsare

SCD-*config1* a fost testată pe mai mult de 500 intrări din dicționar, de dimensiuni medii și mari. Rata de succes a fost de 91.18%, fiind calculată prin compararea fișierului de ieșire din program cu fișierul adnotat manual la arbori de sens. Cauzele erorilor găsite în urma parsării intrărilor de dicționar pot fi grupate în două clase de bază:

#### I. Inconsecvențe în scrierea articolului în DTLR

O primă sursă de erori de parsare este lipsa monotoniei valorilor marcherilor la același nivel din ierarhia marcherilor de sens:

Ex.1. **A.** [**B.** lipsește] ... **C.** etc.;

Ex.2. **2.** [în loc de **1.**]... **2.** etc.;

Ex.3. **a)**... **b)** ... **c)** ... **b)** [în loc de **d)**]etc.

O soluție este verificarea *monotoniei stricte* a valorilor marcherilor. Astfel, înainte de definitivarea arborelui de sensuri, este necesară verificarea validității succesiunilor de marcheri de pe fiecare nivel de sens.

#### II. Ambiguități în stabilirea regentului și a subordonatului unui sens

În următoarea secvență de marcheri de sens apare o ambiguitate inerentă:

Ex.4. **1. a) b) c) ◇ [◇]**

Problema apare atunci când nu se poate stabili dacă romburile “◇” trebuie considerate ca depinzând de **c)** sau de marcherul de un nivel superior (**1.**). Rezolvarea acestei ambiguități depinde de contextul semantic al perechilor de marcheri implicate.

### 4.2. Analiza tipurilor de definiții din DTLR

Până în prezent a fost realizată segmentarea elementelor dintre doi marcheri de sens succesivi, ținând cont de marcherii de definiții DTLR și de expresiile regulate cu care pot fi recunoscute. Evaluarea segmentării la definiții a fost abordată folosind două metrici: *potrivire exactă* și *suprapunere*. Potrivirea exactă (*exact-match* metric) reprezintă numărul de segmente corect extrase (folosind *precizia*, *recall*-ul și *F-measure*); suprapunerea (*overlap* metric) reprezintă procentul de cuvinte clasificate



PARSAREA ARBORILOR DE SENSURI ȘI SEGMENTAREA  
LA DEFINIȚII ÎN DICȚIONARUL TEZAUR eDTLR

corect (folosind, de asemenea, *precizia*, *recall*-ul și *F-measure*). Deoarece există situații în care segmente de același tip sunt consecutive, primul și ultimul cuvânt al fiecărui segment sunt marcate, pentru a putea penaliza clasificarea cuvintelor într-un segment mare în locul unei succesiuni de segmente mai mici.

Pentru evaluare au fost utilizate 52 de intrări din dicționar, de dimensiuni diferite, ca *standard-gold*, însumând un număr de aproximativ 2000 de segmente și 22.000 de cuvinte. Rezultatele sunt prezentate în Tabelul 1.

Tabel 1: Evaluarea pentru segmentarea la nivel de definiții DTLR

Tipul evaluării	Precizie	Recall	F-measure
Potrivire exactă	93.24%	85.41%	89.15%
Suprapunere	97.86%	97.80%	97.83%

După analiza rezultatelor evaluării, am observat că cele mai frecvente erori se datorează segmentării greșite a siglelor. Tabelul 2 prezintă cele mai frecvente zece tipuri de erori.

Tabel 2: Cele mai frecvente zece erori în segmentarea definițiilor din DTLR

Rezultatul parsării	Parsare gold	% erori introduse
Sigle	Început de Siglă	29.45%
Sigle	Sfârșit de Siglă	28.08%
RegDef	SpecDef	6.39%
RegDef	Sfârșit de RegDef	3.88%
RegDef	Început de RegDef	2.96%
DefExem	ItalMarker	2.73%
RegDef	Început de SpecDef	2.51%
Sigle	RegDef	2.28%
RegDef	Sigle	2.05%
RegDef	SpSpecDef	2.05%

Corectarea segmentării siglelor duce la o *F-measure* de 94.43% pentru metrica de potrivire-exactă și de 98.01% pentru metrica de suprapunere.

## 5. Concluzii

Această lucrare a prezentat o metodă nouă de parsare a intrărilor de dicționar, în mod concret a tezaurului DTLR, bazată pe configurații SCD. Prima configurație a exploatat setul de marcheri de sensuri DTLR pentru construirea arborelui de sensuri, obținând o acuratețe de 91.18%. A doua configurație SCD este folosită pentru a parsă definițiile din DTLR cuprinse în fiecare nod din arborele de sensuri. Acuratețea pentru clasificarea definițiilor depășește 93%.

Parserul bazat pe configurații SCD are avantajul că, odată stabilite în mod adecvat clasele de marcheri și ierarhia lor pentru un anumit dicționar tezaur, oricât ar fi acesta de complex (cum este cazul DTLR), programul poate parsă foarte eficient acel tezaur.

**Mulțumiri.** Rezultatele din această lucrare au fost obținute în cadrul cercetărilor la grantul eDTLR – PNCDI 2, No. 91\_013/18.09.2007. Mulțumiri speciale sunt datorate cercetătoarelor Gabriela Haja și Elena Dănilă, Institutul de Filologie Română "Al.

Philippide” Iași, pentru discuțiile consistente privind funcționarea, dependențele și moștenirea definițiilor la sensuri în DTLR (Secțiunea 3).

### Referințe bibliografice

- Curteanu, N., E. Amihăesei (2004): Grammar-based Java Parsers for DEX and DTLR Romanian Dictionaries. *Proc.of ECIT-2004 Conference*, Iasi, Romania.
- Curteanu, N. (2006): Local and Global Parsing with Functional (F)X bar Theory and SCD Linguistic Strategy. (I.+II.), *Computer Science Journal of Moldova*, Academy of Science of Moldova, Vol. 14 no. 1 (40):74-102 and no. 2 (41):155-182.
- Curteanu, N., G. Pavel, C. Vereștiuc, D. Trandabăț (2007). Parsarea eDTLR cu gramatici în mediul JavaCC. Stadiul actual, probleme și soluții de dezvoltare. (Ed. I. Pistol, D. Cristea, D. Tufiș) *Resurse lingvistice și instrumente pentru prelucrarea limbii române*, ConsILR-2007, Ed. Univ. ”Al. I Cuza” Iași, p. 87-96.
- Curteanu, N., Moruz, A., Trandabăț, D. (2008). Extracting Sense Trees from the Romanian Thesaurus by Sense Segmentation & Dependency Parsing, *Proceedings of CogAlex Workshop, COLING 2008*, pp. 55-63, ISBN 978-1-905593-56-9.
- Hauser, R., Storrer, A. (1993). Dictionary Entry Parsing Using the LexParse System. *Lexikographica* 9 (1993), 174-219
- Kammerer, M. (2000): *Wörterbuchparsing Grundsätzliche Überlegungen und ein Kurzbericht über praktische Erfahrungen*, <http://www.matthias-kammerer.de/content/WBParsing.pdf>
- Lemnitzer, L., Kunze, C. (2005): *Dictionary Entry Parsing*, ESSLLI 2005 Tutorial.
- Neff, M., Boguraev, B. (1989) Dictionaries, Dictionary Grammars and Dictionary Entry Parsing, *Proc. of the 27th ACL Vancouver*, British Columbia, Canada Pages: 91 - 101
- Tufiș, D., Rotaru, G., Barbu, A.M. (1999). Data Sampling, Lemma Selection and a Core Explanatory Dictionary of Romanian. *Proc. of the 5th International Workshop on Computational Lexicography COMPLEX*, Pecs, Hungary, pp. 219-228, 1999

# SEGMENTAREA ÎN UNITĂȚI TEXTUALE ATOMICE A INTRĂRILOR DIN DICȚIONARUL LIMBII ROMÂNE ÎN VEDEREA ANALIZEI STRUCTURALE

RADU ION

*Institutul de Cercetări pentru Inteligență Artificială, Academia Română,  
București – România*

[radu@racai.ro](mailto:radu@racai.ro)

## Rezumat

Lucrarea de față prezintă un algoritm de adnotare a unităților textuale atomice care compun definițiile intrărilor de dicționar din Dicționarul Limbii Române (DLR) al Academiei Române. Algoritmul care va fi prezentat se bazează pe colecții de expresii regulate care sunt aplicate succesiv (atât colecțiile cât și expresiile din fiecare colecție) pe intrarea de dicționar. Ca rezultat, fiecare expresie regulată va „recunoaște” secvențe continue de text care au anumite semnificații în cadrul definiției. Această fază de procesare poate fi folosită ulterior de un parser al cărei gramatici va utiliza adnotările pe post de simboluri terminale. În acest fel, se va simplifica scrierea gramaticii care acceptă o intrare de dicționar.

## 1. Introducere

Dicționarul Limbii Române (DLR) este continuarea Dicționarului Academiei<sup>1</sup> (DA) a cărui construcție a început în 1913. El reia enumerarea minuțioasă a fondului lexical de la intrarea *Lojniță* cu scopul declarat de a inventaria tezaurul lexical al limbii române. Cele două lucrări sunt astfel colectiv cunoscute sub denumirea de Dicționarul Tezaur al Limbii Române (DTLR) care este „cea mai amplă lucrare lexicografică românească, considerat nu o dată o operă de importanță națională” (Sala, 1996).

Proiectul eDTLR<sup>2</sup>, început în anul 2007, are drept scop transpunerea DTLR în format electronic cu urmări benefice evidente pentru comunitatea lexicografică românească implicată în dezvoltarea lui dar și pentru comunitatea lingvisticii computaționale românești (Cristea et al., 2007). În ce privește lucrul la DTLR, formatul electronic permite operații ca interogarea (pe diverse criterii) și vizualizarea intrărilor cu o ușurință de neimaginat pentru lexicografii secolului trecut. Lingvistica computațională românească are însă, probabil, cel mai mult de câștigat de pe urma unei astfel de resurse lexicografice monumentale. De la analizele morfologice până la diversitatea enormă de sensuri inventariate, eDTLR este util pentru o multitudine de probleme precum dezambiguizarea semantică automată (engl. *Word Sense Disambiguation*), analiză și generare morfologică, adnotarea morfosintactică (engl. *Part Of Speech Tagging*). De asemenea, eDTLR este o sursă neprețuită de validare semantică și extindere a ontologiei

---

<sup>1</sup> Române.

<sup>2</sup> [https://consilr.info.uaic.ro/edtlr/wiki/index.php?title=Despre\\_proiect](https://consilr.info.uaic.ro/edtlr/wiki/index.php?title=Despre_proiect)

lexicale pentru limba română, RoWN (Tufiş et al., 2008) prin diversitatea sensurilor inventariate pentru fiecare cuvânt.

Formatul electronic al unui dicționar este valoros în măsura în care evidențiază prin adnotări structura intrărilor. Astfel, de cele mai multe ori, o intrare de dicționar este divizată logic într-o mulțime de sensuri. Fiecare sens conține o definiție, exemple de utilizare a cuvântului în sensul respectiv, sensuri secundare etc. Structura intrărilor de dicționar poate fi utilizată de aplicații de Prelucrare Automată a Limbajului Natural (PLN) sau poate fi utilă diferitelor tipuri de interogări care se pot imagina (de exemplu, definiția primului sens al celui de al doilea omonim al cuvântului „mină”). În general, formatul electronic se obține printr-o analiză structurală (engl. *Parsing*) a unei intrări de dicționar furnizată sub formă de text electronic<sup>3</sup>.

În cele ce urmează vom descrie pe scurt câteva metode de generare a formatelor electronice ale dicționarelor și apoi vom prezenta colecția noastră de expresii regulate care segmentează o intrare din DLR în unități textuale atomice în vederea analizei structurale.

## 2. Analiza automată a structurii unei intrări de dicționar

Transformarea intrărilor de dicționar din format text (care este un format electronic nestructurat) în format electronic care evidențiază structura este o problemă care interesează comunitatea lexicografiei computaționale în măsura în care se pot crea automat resurse lexicografice computaționale din diversele formate text ale dicționarelor.

„Dictionary Parsing Project (DPP)<sup>4</sup>” derulat de grupul de PLN din cadrul USC Information Sciences Institute își propune să extragă relații semantice (hipernimie, holonimie, relații sintagmatice de tipul șofer–vehicul etc.) din „Noah Webster's 1913 Dictionary of the English Language<sup>5</sup>”. Pentru aceasta, structura unei intrări în formă text a dicționarului este inițial procesată pentru a se obține forma descrisă cu expresii regulate din Figura 1. Din această formă, intrarea de dicționar este mai departe prelucrată în direcția depistării unor unități de text denumite „frazе” care au semnificații bine-stabilite pentru compoziția intrării de dicționar. De exemplu, fraza care delimitează începutul intrării de dicționar și care conține cuvântul-titlu, partea de vorbire și numărul sensului, se definește cu următoarea expresie regulată:

```
HEADWORDLINE := <hw> ( { (WORD{WORD}*) }+ ) </hw> <pos>POS</pos>
<sn>NUMBER</sn>
```

în care simbolurile „{,},\*,+” fac parte din limbajul de specificare a expresiilor regulate (vezi Figura 1) iar WORD, POS și NUMBER sunt definițiile altor expresii regulate care descriu un cuvânt, partea sa de vorbire și, respectiv, un număr de sens.

Cea mai mare parte a literaturii care se referă la achiziția și prelucrarea de MRD (engl. *Machine Readable Dictionaries*) descrie metodologiile de a transforma MRD (formă electronică nestructurată sau text în accepțiunea noastră) în LDB (engl. *Lexical*

<sup>3</sup> Acesta obținându-se la rândul său prin transformarea textului tipărit în imagine electronică cu recunoașterea automată a caracterelor și generarea textului electronic corespunzător celui tipărit.

<sup>4</sup> <http://www.isi.edu/natural-language/dpp/>

<sup>5</sup> [http://humanities.uchicago.edu/orgs/ARTFL/forms\\_unrest/webster.form.html](http://humanities.uchicago.edu/orgs/ARTFL/forms_unrest/webster.form.html)

SEGMENTAREA ÎN UNITĂȚI TEXTUALE ATOMICE A INTRĂRILOR DIN DICȚIONARUL LIMBII ROMÂNE ÎN VEDEREA ANALIZEI STRUCTURALE

*DataBases* – formatul electronic în care structura unei intrări de dicționar este evidențiată prin adnotări specifice).

	Legend
<pre> &lt;entry&gt; {&lt;hw&gt;words&lt;/hw&gt; &lt;stress&gt;words&lt;/stress&gt;?}+ &lt;plural&gt;words&lt;/plural&gt;? &lt;colref&gt;word&lt;/colref&gt;? {&lt;asp&gt;word&lt;asp&gt;* &lt;sense&gt; &lt;pos&gt;word&lt;/pos&gt; &lt;uniqhw&gt;word&lt;/uniqhw&gt;? &lt;uniqsn&gt;word&lt;/uniqsn&gt;? &lt;sn&gt;word&lt;/sn&gt; &lt;subsense&gt;? &lt;subsn&gt;word&lt;/subsn&gt;? &lt;uniqhw&gt;word&lt;/uniqhw&gt;? &lt;fld&gt;words&lt;/fld&gt;? &lt;uniqdef&gt;words&lt;/uniqdef&gt;? &lt;def&gt;words&lt;/def&gt;+ &lt;as&gt;words&lt;/as&gt;* &lt;note&gt;words&lt;/note&gt;* &lt;quote&gt;words&lt;/quote&gt;* &lt;comment&gt;words&lt;/comment&gt;* &lt;/subsense&gt;? &lt;/sense&gt;}+ &lt;/entry&gt; </pre>	<pre> ? zero or one * zero or more + one or more {} grouping  as example usage asp alternate spelling comment remark def definition fld sense field hw headword note note plural pluralization pos part of speech quote example quotation stress headword with stress markings sn sense number subsn subsense number uniqhw unique headword ID uniqdef unique definition number uniqsn unique sense ID </pre>

Figura 1: Structura unei intrări din dicționarul Webster 1913<sup>6</sup>.

Neff și Boguraev (1989) disting două tipuri de sisteme de analiză structurală a intrărilor de dicționar:

1. sistemele monolit în care regulile de analiză sunt conținute în aplicație și care, din acest motiv, nu pot fi adaptate să funcționeze pe alte dicționare;
2. sistemele bazate pe gramatici independente de context (GIC) în care sistemul este format dintr-un parser și o gramatică. Gramatica are producții cu care recunoaște intrări de dicționar, „marele avantaj” fiind acela că sistemul se poate adapta când avem de-a face cu un alt dicționar, prin scrierea altei gramatici.

Bineînțeles că „scrierea unei alte gramatici” echivalează practic cu scrierea unui nou sistem monolit de analiză întrucât intrări din dicționare diferite, diferă substanțial în termenii convențiilor de alcătuire a unei intrări.

În ce privește tratamentul aplicat dicționarelor în limba română, putem exemplifica transformarea Dicționarului Explicativ al Limbii Române (DEX), (Tufiș et al., 1999) din format text în format bază de date XML prin utilizarea unei GIC special dezvoltată pentru DEX. Experiența autorilor demonstrează faptul că adnotarea conformă cu TEI<sup>7</sup> nu a fost posibilă în cazul DEX fără a sacrifica din informația lexicală (TEI nu are

<sup>6</sup> Captură de imagine de la <http://www.isi.edu/natural-language/dpp/>

<sup>7</sup> Text Encoding Initiative, <http://www.tei-c.org/index.xml>

elemente care să descrie informația lexicală din DEX) și/sau cea editorială (ordinea în care elementele unei intrări sunt date și diversele notații care se pierd – devin redundante – prin adnotare). Acest lucru ne îndreptățește să credem că nici DTLR nu va putea fi 100% reprezentat în această codificare.

### 3. Segmentarea unei intrări din DLR

Dacă ar fi să adoptăm sistemul de analiză structurală care folosește GIC, putem afirma că o metodă evidentă de a obține gramatici mai simple (și astfel, probabil, mai ușor de generalizat) ar fi să *simplificăm limbajul pe care gramatica trebuie să-l accepte*. Altfel spus, fiecare parser consumă un șir de simboluri care îi este furnizat la intrare pentru ca la final să raporteze dacă a acceptat șirul sau nu (dacă l-a acceptat, poate de asemenea să prezinte structura arborescentă a șirului de simboluri). Dacă simplificăm limbajul (reducem numărul de simboluri posibile) vom simplifica implicit gramatica care îl acceptă, rezultând firesc o structură mai simplă. *Ideea principală* este următoarea: structura simplificată ar trebuie să fie cât mai apropiată de structura logică pe care o putem vedea într-un dicționar: un cuvânt are mai multe sensuri, fiecare sens are o definiție și exemple de utilizare, etc.

Pentru o intrare DLR în format text codificat UTF-8 (deci fără marcajele bold, italic, superscript, etc.), am imaginat această simplificare printr-o operație de segmentare: identificarea pasajelor de text continue care reprezintă unități atomice (pe care le vom numi în continuare „*tokeni*”) în alcătuirea unei intrări (de exemplu sursa unui citat, anul atestării documentare a unui citat, un identificator de sens, o parte de vorbire, cuvântul titlu, etc.). Pentru segmentarea unei intrări DLR, am utilizat o listă ordonată de colecții de expresii regulate. Fiecare expresie regulată dintr-o colecție este menită a identifica un tip de token care se realizează sub o anumită formă în text. De exemplu, pentru a identifica anul atestării unui exemplu de utilizare (care este citatul) a sensului cuvântului titlu, textul ne poate oferi tokeni cum ar fi: „(cca. 1550)”, „(a. 1742)”, „(cca 1569–1575)”, etc. Când sunt întâlniți, acești tokeni sunt recunoscuți de diversele expresii regulate din colecția dedicată acestui tip de token: anul atestării documentare.

Prima problemă cu un astfel de tip de abordare este că diversele expresii regulate dintr-o colecție pot recunoaște *tokeni care se suprapun*. Pentru a elimina acest inconvenient, fiecare expresie regulată din colecție are asociată o prioritate (un număr natural) iar expresiile sunt „încercate” în ordinea crescătoare a acestor priorități. Vom prefera bineînțeles expresiile regulate care recunosc *tokeni cât mai lungi*. Alături de mecanismul priorităților, o a doua metodă, cea a verificării argumentelor, este folosită pentru a ne asigura că o anumită expresie regulată recunoaște un token întreg și nu unul parțial. Fiecare expresie regulată conține o serie de capturi (secvențe incluse între paranteze rotunde „(” și „)”) în cadrul expresiei regulate) pe care le numim „*atributele*” tokenului. De exemplu, tokenul „(cca. 1550)” are ca atribut anul în care s-a făcut atestarea, anume 1550. Aceste atribute sunt verificate automat cu ajutorul unor liste de atribute posibile în momentul în care expresia regulată a recunoscut un token. Pentru exemplificare, fie două expresii regulate în Perl care recunosc tokeni de tipul surse de citat:

```
#Recunoaste o citare de tipul BELEA, P. A. 148
"author_2" => qr/((\${RXAUTH}),\s*(\${RXWORK})\s+(\${RXNOINT}))/,
```

## SEGMENTAREA ÎN UNITĂȚI TEXTUALE ATOMICE A INTRĂRILOR DIN DICȚIONARUL LIMBII ROMÂNE ÎN VEDEREA ANALIZEI STRUCTURALE

```
"author_2_args" => { "_text" => 1, "author" => 2, "source" =>
3, "pages" => 4 },
"author_2_check" => { "author" => \%DLRAUTHORS, "source" =>
\%DLRSOURCES },
"author_2_rank" => 400
#Recunoaste o citare de tipul MARCOVICI, D. 154/13
"author_3" =>
qr/((\${RXAUTH}),\s*(\${RXWORK})\s+(\${RXNOFRAC}))/,
"author_3_args" => { "_text" => 1, "author" => 2, "source" =>
3, "pages" => 4 },
"author_3_check" => { "author" => \%DLRAUTHORS, "source" =>
\%DLRSOURCES },
"author_3_rank" => 200
```

Expresia regulată „author\_2” recunoaște tokenul „BELEA, P. A. 148” în care „BELEA” este autorul, „P.A.” este abrevierea lucrării iar „148” este numărul de pagină în lucrarea respectivă. Aceste 3 atribute sunt delimitate cu „()” în expresia regulată<sup>8</sup>. În momentul în care motorul de aplicare a expresiilor regulate a recunoscut tokenul „BELEA, P. A. 148” cu expresia „author\_2”, se extrag atributele tokenului

```
author="BELEA", source="P.A.", pages="148", _text="BELEA, P.
A. 148"
```

prin inspecția listei de atribute corespunzătoare (author\_2\_args) iar valorile atributelor author și source se verifică căutându-se în listele DLRAUTHORS și DLRSOURCES<sup>9</sup> (author\_2\_check). Numai în cazul în care valorile au fost validate, se acceptă tokenul și se continuă procesul de recunoaștere. Recunoașterea unui token înseamnă marcarea lui în textul intrării de dicționar cu o notație de tip XML care specifică atât tokenul cât și atributele sale. Pentru exemplul nostru, tokenul va fi adnotat ca

```
<AUTHCITE source="P.A." author="BELEA" pages="148">BELEA, P.
A. 148</AUTHCITE>
```

Deocamdată, segmentatorul nostru recunoaște 6 tipuri de tokeni:

- cuvinte-titlu împreună cu terminații și părți de vorbire (DLREENTRY, colecție cu 3 expresii regulate);
- atestare documentară (ATTESTED, 6 expresii regulate);
- trimiteri la sensurile altor cuvinte (ALSOSEE, 7 expresii regulate);
- citare cu autor (AUTHCITE, 14 expresii regulate);
- marcaje de sens (SENSE, 3 expresii regulate);

<sup>8</sup> \$RXAUTH, \$RXWORK și \$RXNOINT sunt variabile ale căror valori sunt alte expresii regulate care recunosc un nume de autor, un titlu de lucrare și respectiv un număr de pagină. De exemplu, valoarea variabilei \$RXNOINT este qr/(?:[0-9]+)”.  
<sup>9</sup> La momentul scrierii acestor rânduri lista autorilor are aprox. 500 de nume iar cea de lucrări, cca. 400 de intrări.

- citare fără autor (SRCCITE, 19 expresii regulate).

În general în fiecare colecție există câte o expresie regulată pentru fiecare formă a tipului de token întâlnită în practică. Evident, aceste colecții nu sunt complete și vor trebui îmbogățite cu expresii regulate pentru fiecare formă necunoscută de token. Plătim astfel prețul unei gramatici de analiză a unei intrări de dicționar mai simplă care altfel ar fi trebuit să conțină reguli de producție pentru astfel de tokeni.

Pentru a exemplifica ieșirea segmentatorului cu adnotarea tokenilor de tipurile descrise mai sus, dăm începutul intrării REVĂRSĂT

```
<DLRENTRY suffix="-Ă" pos="adj." hword="REVĂRSĂT"
note="2">REVĂRSĂT2, -Ă adj.</DLRENTRY> <SENSE ind="Despre ape
curgătoare" subsense="1">1. (Despre ape curgătoare)</SENSE>
```

Care s-a vărsat peste margini, care a ieșit din albie; care a inundat.

Cf. <ALSOSEE word="revărsa" subsense="1">revărsa (1)</ALSOSEE>.

Agiunsă de a trece apele, revărsate.

```
<AUTHCITE source="S. L." volume="II" author="ASACHI"
pages="19">ASACHI, S. L. II, 19</AUTHCITE>, ...
```

#### 4. Concluzii

Despre evaluarea corectitudinii segmentării putem spune doar că exemplele de segmentare verificate de noi au fost corecte (când am întâlnit erori, am ajustat prioritățile de aplicare și/sau am modificat expresiile astfel încât să obținem rezultatele scontate). În momentul în care parserul va putea folosi această segmentare, vom putea da un procent de intrări analizate corect. Trebuie să spunem că acest algoritm de segmentare este destul de lent întrucât fiecare expresie regulată este încercată pe fiecare intrare de dicționar. Timpul mediu de segmentare a unei intrări din DLR folosind cele 52 de expresii regulate existente acum este de aproximativ 1.1 secunde dar va crește cu creșterea numărului de expresii regulate.

Acum dispunem de o listă de aproximativ 3800 de sigle bibliografice<sup>10</sup> și de o listă cu toate abrevierile folosite în DLR, resurse care vor mări considerabil recall-ul segmentatorului. Următorul pas în dezvoltarea unui parser DLR este să ne concentrăm pe segmentarea completă a 100 de intrări extrase aleatoriu din DLR urmată de scrierea unei gramatici-nucleu care să accepte aceste intrări. Apoi, vom analiza DLR intrare cu intrare și vom modifica segmentatorul/gramatica astfel încât noile intrări să fie acceptate. Metoda de analiză structurală alternativă dezvoltată de colegii noștri (Curteanu et al., 2008) va produce analize care vor fi de referință pentru parserul nostru. Analize identice pentru o aceeași intrare vor putea fi considerate corecte întrucât au fost generate independent de două metode diferite.

<sup>10</sup> Nu am știut de existența acestei liste când am lucrat la segmentator. A fost creată odată cu DLR de autorii dicționarului. Mulțumim domnului Victor Celac pentru o copie a acestei liste.



### Referințe bibliografice

- Cristea, D., Răschip, M., Forăscu, C., Haja, G., Florescu, C., Aldea, B., Dănilă, E. (2007). The Digital Form of the Thesaurus Dictionary of the Romanian Language. In C. Burileanu, H.N. Teodorescu (eds.), *Proceedings of the 4th International IEEE Conference SpeD 2007, "Advances in Spoken Language Technology"*, Iași, 10-12 Mai 2007. Editura Academiei Române, ISBN 978-973-27-1516-1.
- Curteanu, N., Moruz, A.-M., Trandabăț, D. (2008). Extracting Sense Trees from the Romanian Thesaurus by Sense Segmentation & Dependency Parsing. In *Proceedings of the COLING 2008 Workshop on Cognitive Aspects of the Lexicon (COGALEX 2008)*, pp. 55–63, Manchester, UK, 24 August, 2008. ISBN: 978-1-905593-56-9.
- Neff, M.S., Boguraev, B.K. (1989). Dictionaries, Dictionary Grammars and Dictionary Entry Parsing. In *Proceedings of the 27rd Annual Conference of the Association for Computational Linguistics*, pp. 91-101.
- Sala, M. (1996). Dicționarul limbii române (DLR). In Dan Tufiș (ed.), *Limba și Tehnologie*. Editura Academiei Române, București, 1996. ISBN 973-27-0542-6. 270 p.
- Tufiș, D., Rotariu, G., Barbu, A.-M. (1999). TEI-Encoding of a Core Explanatory Dictionary of Romanian. In Ferenc Kiefer, Gábor Kiss, and Júlia Pajzs (eds.), *Proceedings of the 5th International Workshop on Computational Lexicography (COMPLEX 1999)*, pp. 219-228, Pecs, Hungary, May 1999. Linguistics Institute, Hungarian Academy of Sciences.
- Tufiș, D., Ion, R., Bozianu, L., Ceașu, A., Ștefănescu, D. (2008). Romanian Wordnet: Current State, New Applications and Prospects. In Attila Tanacs, Dora Csendes, Veronika Vincze, Christiane Fellbaum, and Piek Vossen (eds.), *Proceedings of 4th Global WordNet Conference, GWC-2008*, pp. 441-452, Szeged, Hungary, January 2008. University of Szeged, Hungary. ISBN 978-963-482-854-9.



**PLATFORMĂ PLURILINGVĂ DE FORMARE ȘI AUTOFORMARE  
ÎN DOMENIUL LIMBILOR ROMÂNICE**

DOINA SPIȚĂ ȘI CLAUDIA BÎZDÎGĂ

*Universitatea "Al.I.Cuza", Facultatea de Litere, Iași - România*

[doinaspita@yahoo.com](mailto:doinaspita@yahoo.com), [claudia\\_bizdiga@yahoo.es](mailto:claudia_bizdiga@yahoo.es)

**Rezumat**

"GALAPRO"<sup>1</sup> este un program transversal *Langues* "Formation de formateurs à l'intercompréhension en Langues Romanes", componentă a Programului pentru educație și formare pe tot parcursul vieții, N° 135470-LLP-1-2007-1-PT-KA2-KA2MP, 2008-2010. El reunește universități din Portugalia – Universitatea din Aveiro fiind coordonatorul de proiect, Belgia, Franța, Italia, Spania, precum și Universitatea "Al.I.Cuza" din Iași. Proiectul își propune să dezvolte o rețea de formare specializată în domeniul cunoscut de mulți ani în cercetarea lingvistică aplicată sub denumirea de *intercompréhension des langues romanes*, limbi între care, până în 2008, româna era cvasi ignorată ca referință. Obiectivul este acela de a crea și de a experimenta, cu ajutorul unei platforme digitale, scenarii colaborative de formare centrate pe sarcini capabile să răspundă nevoilor și așteptărilor (în termeni de competențe profesionale și profil lingvistic și comunicativ) diverselor categorii de public țintă. Pe termen lung, proiectul vizează resurse conceptuale și practice (glosare plurilingve, publicații, bază de date) transferabile către alte familii de limbi și către alte discipline.

**1. O soluție de comunicare lingvistică pentru comunitățile plurilingve**

Demersul propus de "Galapro" se înscrie în perspectiva plurilingvistului și a didacticii acționale, așa cum sunt ele definite în documentele programatice ale Consiliului European în domeniul politicilor lingvistice. Libertatea de circulație și mobilitatea pe piața muncii, armonizarea sistemelor de învățământ și intensificarea schimburilor economice, culturale și științifice în Uniunea Europeană impun cu evidență cunoașterea mai multor limbi ca pe o soluție strategică prioritară pentru prezervarea diversității lingvistice și culturale.

Patru orientări posibile sunt recomandate de către instituțiile europene rețelelor de cercetare care servesc această prioritate: difuzarea bunelor practici cu privire la învățarea limbilor de către adulți; inventarierea nevoilor actuale și identificarea celor ale viitorului, legate de cooperarea europeană în domeniul învățării limbilor de către adulți, prin metode formale, nonformale și informale; elaborarea de strategii care să acopere lipsurile existente în acest domeniu în care oferta nu mai corespunde exigențelor actuale și mai ales ale celor de perspectivă; în sfârșit, difuzarea programelor și instrumentelor

---

<sup>1</sup> **Membrii echipei românești de cercetare:** Doina Spiță – coordonator; Claudia Tărnauceanu, Mihaela Lupu, Dana Nica, Maria Husarciuc, Paula Onofrei, Claudia Bîzdîgă.

permițând formarea persoanelor implicate în învățarea limbilor de către adulți. Așa cum vom vedea, programul "Galapro" răspunde tuturor acestor obiective specifice.

## 2. *Istoricul cercetării*

Primele proiecte vizând elaborarea unor metode de dezvoltare a competențelor de înțelegere în scris în domeniul limbilor romanice au apărut la sfârșitul anilor '80, imediat după intrarea Spaniei și Portugaliei în Uniunea Europeană. Reținem trei dintre acestea: *EuRom4*, coordonat de Claire Blanche Benveniste, la Aix-en-Provence; *EuroComRom*, coordonat de Horst G. Klein, la Frankfurt; în fine *Galateea*, coordonat de Louise Dabène, la Grenoble, ale cărei cercetări au condus la *Galanet*, proiect Socrates-Lingua coordonat de Christian Degache, de la Universitatea Stendhal Grenoble 3, în anii 2001-2004. Nici unul dintre aceste proiecte nu a vizat în mod explicit limba română.

Platforma de comunicare plurilingvă a fost inițiată în cadrul proiectului *Galanet*. Obiectivul era acela de a pune la dispoziția vorbitorilor de portugheză, spaniolă, italiană și franceză un instrument de formare la distanță pe Internet, care să le ofere posibilitatea de a comunica. Originalitatea platformei consta în oportunitatea oferită vorbitorilor de diferite limbi romanice de a practica *intercomprehenșiunea*, înțelesă ca o formă de comunicare plurilingvă în care fiecare înțelege limbile vorbite de ceilalți, dar se exprimă în limba / limbile romanice pe care el însuși le cunoaște, dezvoltând în acest fel competențe de nivel diferit de cunoaștere a diverselor limbi. Pentru stimularea comunicării, s-a recurs la formularea unor sarcini comune de lucru, participanții fiind puși în situația de a interacționa pentru elaborarea unui proiect colectiv. Intervenția directă a "tutorilor", ca și numeroasele "resurse" puse la dispoziție pe platformă aveau rolul de a facilita comunicarea și de a-i permite fluidizarea. Publicul țintă vizat era constituit din studenți în învățământul superior sau în centre de limbi, liceeni și adulți cunoscători a cel puțin o limbă romanică de referință ca limbă maternă sau străină, fără însă a fi în mod necesar cunoscători, fie și la nivel de debutant, al celorlalte trei limbi.

De la începutul anului 2008, partenerii s-au angajat într-un nou proiect (2008-2010, LLP KA2), numit *Galapro*, vizând de această dată formarea de formatori în domeniul intercomprehenșiunii. Coordonatorul de proiect este Maria Helena de ARAUJO e SA, de la Universitatea din Aveiro<sup>2</sup>. Pentru a servi noii finalități, platforma este în curs de revizuire, cu atât mai mult cu cât două alte limbi romanice au fost invitate în echipă: româna și catalana.

### 2.1. *Ce este deci intercomprehenșiunea?*

Așa cum puncta Jean-Pierre Chavagne de la Universitatea Lumière Lyon 2, parteneră în proiect, *intercomprehenșiunea* înseamnă înțelegere încrucișată, înțelegere reciprocă, faptul că, în situație de dialog, fiecare se poate exprima în limba sa înțelegând-o, în același timp, pe a celorlalți, ceea ce este mult mai avantajos, cel puțin în termeni de

<sup>2</sup> În perioada cuprinsă între 2 – 4 octombrie 2008, a avut loc la Iași a doua întâlnire a echipei internaționale. Au participat reprezentanții Universității din Aveiro – coordonatorul proiectului, Stendhal Grenoble 3, Lumière Lyon 2, ai Universității Autonome din Barcelona, Universității Complutense din Madrid, Universității din Cassino și Universității din Mons-Hainaut - Departamentul de Tehnologie a Educației.

## PLATFORMĂ PLURILINGVĂ DE FORMARE ȘI AUTOFORMARE ÎN DOMENIUL LIMBILOR ROMANICE

randament lingvistic, decât a încerca să te exprimi într-o limbă care nu este a ta, cu riscul de a nu te face înțeles decât rudimentar. În sprijinul utilizării acestei strategii plurilingve de comunicare ce exploatează proximitatea lingvistică pot fi evocate și alte argumente. Mai întâi, acela că este mai ușor și mai rapid să înveți să înțelegi o limbă decât să o vorbești. Apoi, conversația este mai echilibrată și mai eficientă: cele două persoane sunt în poziție de egalitate, fiecare se poate exprima cu un plus de claritate și finețe, căci se exprimă în limba pe care o cunoaște cel mai bine. Situația creată prezintă, în același timp, o importantă dimensiune de convivialitate, interlocutorii apreciindu-și reciproc efortul investit în a-l înțelege pe celălalt.

### 2.2. Platforma colaborativă – un concept spațial

Vă invităm să o vizitați la adresa [www.galanet.eu](http://www.galanet.eu) (Figura 1), iar noi vă vom fi ghizi.



Figura 1: Pagina principală a platformei

Disponibilă în toate limbile proiectului, platforma este un *concept în același timp spațial și temporal* – cum afirmă autoarele *Manualului de instrucțiuni*, pentru că ea presupune, pe de o parte, un scenariu cronologic – *sesiunile de formare* și, pe de altă parte, un *spațiu de învățare virtual*, cu săli și instrumente de exersare. Astfel:

*Zona A* este rezervată scenariului cronologic al sesiunilor. Este un scenariu conceput în patru faze (vezi cele patru butoane), corespunzătoare etapelor de derulare a unei sesiuni, respectiv unui anumit interval de timp și unui forum. Concepută după acest model, formarea urmărește un proces gradual, care îi conduce pe participanți spre sarcini de lucru din ce în ce mai complexe.

*Zona B*, numită și „Ochiul”, îndeplinește o dublă funcțiune: permite să știi cine mai este conectat și, dacă dorești, să angajezi o comunicare tip *chat*.

*Zona C* este zona barelor de opțiuni: cea din dreapta sus permite alegerea limbii de lucru (catalana, franceza, italiana, portugheza, româna sau spaniola), accesul la mesagerie și

la anunțurile de pe panoul de afișaj; cea din centru jos permite compunerea echipelor, schimbarea statutului de participare și modificarea fazei sesiunii în curs de desfășurare.

*Zona D* este spațiul propriu-zis de lucru și permite accesul la o serie de „instrumente” având funcții bine precizate: unele servesc comunicării între participanți (chatul, mesageria, forumul); altele ajută la arhivare (a chaturilor sau a ultimelor conexiuni); altele permit organizarea sesiunilor („Cine este cine?”, „Profilul meu”, „Profilul echipei mele”, „Preferințele mele”); în fine, altele au funcția de „facilitator” al auto-formării (modulele și resursele).

Platforma dispune de mai multe săli polivalente:

*Forumul* (Sala a) este spațiul central al interacțiunilor, locul unde se desfășoară scenariul pedagogic al sesiunii; aici aflăm date despre participanți, despre profilul echipei, de aici se accesează mesageriile personale.

*Biroul meu* (Sala b) este locul de unde îmi pot trimite mesajele, unde pot fi cunoscut după „Profilul” și „Preferințele mele”: limbi de referință, parolă, documentele din forum pe care doresc să le primesc.

*Biroul echipei mele* (Sala c) este spațiul în care sunt propuse și votate temele proiectelor comune („Alegerea temei”) și din care poți avea acces la „Profilul” și la „Chatul echipei”.

*Sala de redactare* (Sala d) este locul unde se concepe și se editează „Dosarul de presă”, cu acces la chat pentru echipa redacțională.

*Sala de reuniune* (Sala e) este un *chat* rezervat întâlnirilor coordonatorilor diferitelor echipe.

*Biblioteca* (Sala f) este locul în care pot fi consultate fișierele.

În *Sala tehnică* (Sala g) descoperim „Profilul echipei tehnice”.

*Sălile pentru chat* (Sala h) sunt trei saloane diferit colorate (albastru, galben și roșu) în care participanții care intervin în diverse echipe își pot da întâlnire pe chat. Aceste întâlniri sunt automat arhivate (cu excepția celor private) în spațiul denumit *Arhiva chaturilor* (Sala i). Pentru conversații care nu se doresc a fi arhivate, platforma ne invită în spațiul numit *Bar* (Sala j). În fine, pentru a avea acces la compoziția și profilul diferitelor echipe poți merge în Sala k („Cine este cine?”), pentru a te documenta poți consulta spațiul de „Resurse”, unde vei găsi, de exemplu, compendii de gramatică și de fonetică în limbile proiectului sau abordări comparative (Sala n), iar ca să fii mereu informat asupra actualităților din proiect, poți consulta *Panoul de afișaj* din holul central.

### **2.3. Platforma colaborativă – un concept temporal**

Așa cum am menționat mai sus, platforma „Galanet” prezintă oportunitatea de a-ți permite accesul la ceea ce se numesc *sesiuni de formare*, construite pe bază de scenarii cronologice de activități care se succed în patru faze, controlate cu ajutorul celor patru butoane din Zona A. Astfel:

*Faza 1* permite participanților să se cunoască, pentru ca în final să poată propune și apoi alege, prin vot, o temă de lucru comună. Ea presupune desfășurarea mai multor

## PLATFORMĂ PLURILINGVĂ DE FORMARE ȘI AUTOFORMARE ÎN DOMENIUL LIMBILOR ROMÂNICE

activități, după cum urmează: înscrierea participanților – îi vom numi „stagiați” -, de către coordonatori, în diverse echipe (vezi bara de opțiuni de jos); elaborarea profilului echipei și a fiecăruia dintre stagiați (are loc în spațiul denumit *Biroul meu*, care cuprinde „profilul” și „preferințele mele”); cunoașterea profilului celorlalți („Cine este cine?”); pregătirea temelor de discuție (chaturi, bar, forumuri); votarea temei (se desfășoară în *Biroul echipei mele*); afișarea temei alese de către responsabilul de sesiune (în *Biroul echipei mele* și în forum); desemnarea, de către responsabilul de sesiune, a echipei redacționale (în cadrul forumului și prin curier electronic).

*Faza 2* permite schimburile de opinii între stagiați pe tema aleasă, via chaturi și forumuri. Ea presupune următoarele activități succesive, monitorizate de către animatori: discuții privind alegerea subtemelor și fixarea acestora de către echipa redacțională, care crează și forumurile de discuții; redactarea sintezelor diferitelor discuții desfășurate pe subteme; redactarea liniei editoriale și repartizarea sarcinilor între membrii echipei redacționale (se realizează în sala de reuniune, bar, forumuri și chaturi).

În *faza 3* au loc schimburile de opinii între stagiați și colectarea documentelor, în funcție de rubricile definite de către comitetul de redacție. Ea se desfășoară în trei timpi: adunarea și redactarea documentelor de către stagiați, în funcție de linia editorială hotărâtă (în spațiul numit forum); cunoașterea și discutarea documentelor propuse de către celelalte echipe (chat, forum, bar) și redactarea sintezelor de discuții.

*Faza a 4-a* este rezervată pregătirii și publicării „Dosarului de presă”. Acesta se constituie din sintezele obținute în finalul fazei precedente și poate cuprinde extrase din schimburi plurilingve, precum și documente ilustrative. Activitățile succesive vizează: sintetizarea forumurilor pe echipe (forumuri, bar și chaturi); redactarea sintezei propriilor discuții de către fiecare echipă (sala de redactare); discutarea sintezelor de către toți participanții (forumuri, bar și chaturi); bilanțul final.

### **2.4. Dimensiunea colaborativă**

Așa cum sperăm că a reieșit din cele prezentate anterior, organizarea și desfășurarea unei sesiuni de formare este rezultatul unui efort de concertare între echipele diferitelor limbi prezente pe platformă și a căror funcționare, pentru o etapă de timp dată, se bazează pe acordul negociat și consimțit de către participanți. Aceștia pot avea statute diverse:

*Responsabilul de sesiune* este cel care deschide o nouă sesiune. El este cel care decide asupra diverselor responsabilități de coordonare, precum aceea de înscriere a stagiarii și de repartizare a lor în echipe, asupra datelor de deschidere și de încheiere a diverselor faze, precum și asupra termenului până la care se pot accepta cererile spontane de înscriere la o sesiune.

*Coordonatorul local* constituie echipele (una sau mai multe), înscrie stagiarii, poate invita animatori și poate primi cererile de înscriere depuse pe pagina principală.

*Animatorul* este responsabil de dinamica grupului. El facilitează identificarea, de către stagiați, a celor mai bune strategii de autoformare și răspunde la întrebări.

*Stagiarul* participă la toate activitățile propuse în sesiune, încearcă să citească în toate limbile în care participanții intervin, îi interpelează pe ceilalți stagiaari și pe animatori în legătură cu diverse probleme de formă sau de fond. El este dator să manifeste aceeași considerație pentru toate limbile de comunicare folosite pe platformă.

În fine, *Vizitatorul* poate să intre în forumuri și în spațiile de autoformare, poate să consulte diferitele arhive, dar nu poate interveni.

### 3. De la Galanet la Galapro

Obiectivul principal al proiectului *Galapro* este acela de a difuza și valorifica achizițiile *Galanetului*, exprimate în informații, strategii și instrumente de formare deja experimentate, prin formarea unor agenți educativi specializați în tehnicile de sensibilizare lingvistică prin *intercomprehensiune*.

Cercetările recente din domeniul didacticii limbilor străine abordează conceptul de *intercomprehensiune* fie în raport cu politicile de constituire a unei Europe unite și coerente, fie în raport cu nevoile de comunicare ale diverselor comunități sociale sau profesionale. În această optică, modulele de formare pentru *intercomprehensiune* propuse de *Galapro* vor viza cu prioritate sensibilizarea mediilor educative față de nevoia pregnantă de dezvoltare a unor competențe plurilingve și pluriculturale. Principiile conducătoare vor fi acela de *diversificare* (reflectată în formularea sarcinilor de lucru și activităților, metodologiei și instrumentelor didactice) și de *flexibilitate* (prin crearea unui cadru funcțional de dezvoltare autonomă, a unei gestionări curriculare adaptate, a conceperii programelor și sistemului de evaluare în manieră suplă etc.). În acest spirit, se preconizează realizarea unei ample anchete vizând identificarea nevoilor și așteptărilor diferitelor tipuri de public țintă (a se consulta pagina principală a sitului, pe care se regăsesc Chestionarele 1 și 2, redactate în cele șase limbi ale proiectului).

Cui i se adresează *Galapro* ? Mai întâi, profesorilor de limbi, în formare inițială, debutanți sau experimentați. Apoi profesorilor de alte discipline, cum ar fi istorie și geografie, arte sau turism, interesați în descoperirea și dezvoltarea conceptului de Învățământ Disciplinar Integrat printr-o Logică a Intercomprehensiunii (*Enseignement de Matières Intégrée par une Logique d'Intercompréhension*). Apoi tutorilor și animatorilor de formări la distanță, în exercițiu sau potențiali, orientați către *intercomprehensiune*. În sfârșit, studenților, specialiști sau nu în domeniul lingvistic, precum masteranzii în științele limbajului, limbi și culturi străine, psiho-pedagogie, științele educației, științele comunicării.

Acestui public, *Galapro* îi propune participarea la o serie de sesiuni prototipice de formare de formatori în limbi romanice (catalana, franceza, italiana, portugheza, româna și spaniola), concepute pe baza a două principii integrate:

- formarea în didactica *intercomprehensiunii* prin practicarea *intercomprehensiunii*;
- difuzarea *intercomprehensiunii* prin formarea de agenți sau de viitori agenți educativi, pregătiți să acționeze în contexte diverse.

Dincolo de competențe didactice, *Galapro* ambiționează diseminarea principiilor ce fundamentează însuși conceptul de *intercomprehensiune*, adică *plurilingvism*,



*diversificare și flexibilitate*, prin inducerea unui comportament lingvistic și social adecvat nevoilor unui public specific, acela care are apetitul experiențelor de mobilitate geografică și / sau virtuală.

Fiecare itinerar – sau sesiune – de formare (prevăzut a priori a se desfășura pe o durată cuprinsă între 4 și 15 săptămâni, în funcție de particularitățile contextuale) va viza îndeplinirea uneia sau a mai multor sarcini finale, ale căror produse colaborative vor fi publicate pe site. Diversele sesiuni succesive vor contribui astfel la constituirea progresivă, prin capitalizare, a unui glosar plurilingv al principalelor concepte ce fundamentează *Galapro*. Această bază de resurse va fi deschisă „vizitatorilor”, atât pentru rațiuni de perfectibilitate – prin înregistrarea reacțiilor, cât și pentru difuzare.

Perspectiva partenerială va permite echipei internaționale de proiect construirea, până la finele anului 2010, a unui scenariu de formare modular și flexibil, adaptabil nevoilor specifice ale publicului propriu fiecărei sesiuni, un demers în același timp colectiv, colaborativ și coerent.

#### **4. *Limba română în dispozitivul propus de Galapro***

Devenind, în 2008, membră a unei echipe de proiect ce funcționează de mai bine de zece ani, echipa românească este preocupată ca, într-o primă etapă, să recupereze, în plan personal, informația și competențele dezvoltate de Galanet, iar la nivel de echipă să contribuie la integrarea limbii române în diversele spații și etape de lucru propuse de platformă. Vor fi create, după modelul activităților și traseelor deja experimentate pentru celelalte limbi, module de auto-formare în limba română. Echipa românească va contribui apoi la efortul colaborativ de reconstruire și readaptare a platformei față de exigențe derivate din formularea noului obiectiv, formarea de formatori.

Principalele sarcini de lucru, în curs de realizare în perioada actuală, vizează traducerea întregii platforme în limba română, construirea modulelor de învățare specifice și constituirea bazei de documentare asupra acestei limbi neacoperită de câmpul de cercetare Galanet.

Un prim exercițiu de inițiere s-a desfășurat în cadrul reuniunii de proiect pe care am găzduit-o la Iași, la începutul lunii octombrie. S-a vorbit atunci despre urmele modelului latin la nivel de lexic, morfologie, sintaxă, fonetică și fonologie, s-a vorbit despre evoluția istorică specifică limbii române, despre procesul de *reromanizare* sau *relatinizare* sau *occidentalizare* din secolele al XVIII-lea și al XIX-lea, despre „miracolul” existenței acestei „insule de latinitate în mijlocul unei mări slave”. Evocarea rădăcinilor adânci ale latinității noastre, pe bază de exemple edificatoare, a provocat discuții și a determinat întărirea sentimentului de coeziune a grupului, prin conștientizarea, o dată în plus, a apartenenței la o matcă comună.

Au urmat apoi două secvențe interactive, concepute în spiritul pedagogiei intercomprehenșivității: prima secvență a fost propusă de echipa ieșeană, cealaltă de echipa universității din Barcelona. Aceasta din urmă a propus un model de didacticizare a unei înregistrări audio-video de română vorbită. Secvența poate fi vizionată accesând platforma conform traseului: pasul 1 – „Session en préparation”; pasul 2 – click pe una dintre sesiuni; pasul 3 – se merge la Salonul cu 16 scaune (spațiul de autoformare);

pasul 4 – se ajunge la meniu (dreapta sus), de unde se optează pentru FR>RO, apoi se alege modulul pentru limba română.

Secvența interactivă a echipei românești a ales ca suport un citat din Lucian Blaga:

“*După ce am descoperit că viața nu are nici un sens, nu ne rămâne altceva de făcut decât să-i dăm un sens*”.

Participanții, vorbitori de limbi romanice altele decât limba română, au fost invitați să reconstituie originea latină a termenilor românești și să găsească echivalentele în limbile lor respective. S-au obținut rezultatele de mai jos:

<p><i>după</i> &lt; lat. <i>de post</i></p> <p><i>după ce</i> + indicativ ~ lat. <i>postquam</i> + indicativ</p> <p><i>ce</i> &lt; lat. <i>quid</i> (pronume interogativ)</p>	<p><b>RO:</b> după ce <b>ES:</b> después de <b>GA:</b> despois de <b>CAT:</b> després d' <b>PT:</b> depois de <b>FR:</b> après <b>IT:</b> dopo</p>
<p><i>a descoperi</i> (IV) &lt; lat. <i>disco(o)perio, -ire, -operui, -opertum</i> (IV)</p> <p><i>am descoperit</i> (indicativ, perfect compus, pers. I pl.) &lt; lat. <i>habemus</i> + participiul trecut, la pasiv</p>	<p><b>RO:</b> am descoperit <b>ES:</b> haber descubierto <b>GA:</b> descubrir (<i>limba galiciană nu prezintă timpuri compuse</i>) <b>CAT:</b> haver descobert <b>PT:</b> haver descoberto <b>FR:</b> avoir découvert <b>IT:</b> avvere scoperto</p>
<p><i>că</i> &lt; lat. <i>quod</i></p>	<p><b>RO:</b> că <b>ES:</b> que <b>GA:</b> que <b>CAT:</b> que <b>PT:</b> que <b>FR:</b> que <b>IT:</b> che</p>
<p><i>viața</i> &lt; *<i>vivitia</i> &lt; <i>vivus, -a, -um</i> (&lt; <i>viu</i> + <i>-eață</i>)</p> <p><i>-a</i> (articol hotărât enclitic, feminin, singular, nominativ) &lt; <i>illa</i> (pronume demonstrativ), feminin, singular</p>	<p><b>RO:</b> <u>viața</u> (<i>articol enclitic</i>) <b>ES:</b> <u>la</u> vida (<i>articol hotărât proclitic în cazul celorlalte limbi</i>) <b>GA:</b> <u>a</u> vida <b>CAT:</b> <u>la</u> vida <b>PT:</b> <u>a</u> vida <b>FR:</b> <u>la</u> vie <b>IT:</b> <u>la</u> vita</p>
<p><i>nu</i> &lt; <i>non</i></p> <p><i>a avea</i> (II) &lt; <i>habere</i> (II)</p> <p><i>are</i> – indicativ prezent, pers. a III-a, sg. &lt; <i>haberet</i> (conjunctiv imperfect) sau <i>habuerit</i> (conjunctiv perfect)</p>	<p><b>RO:</b> nu are <b>ES:</b> no tiene <b>GA:</b> non ten <b>CAT:</b> no té <b>PT:</b> não tem <b>FR:</b> n'a pas <b>IT:</b> non ha</p>
<p><i>nici</i> &lt; lat. <i>neque</i></p> <p><i>un</i> &lt; lat. <i>unus</i> (pronume nehotărât) → latina vulgară: valoare apropiată de cea a articolului din limbile romanice <i>nu ... nici</i> (dubla negație) vs. lat. <i>duplex negatio est affirmatio</i></p>	<p><b>RO:</b> nici un <b>ES:</b> ningún <b>GA:</b> ningún <b>CAT:</b> cap <b>PT:</b> nenhum <b>FR:</b> n'a pas de = aucun <b>IT:</b> nessun</p>

PLATFORMĂ PLURILINGVĂ DE FORMARE ȘI AUTOFORMARE  
ÎN DOMENIUL LIMBILOR ROMÂNICE

<p><i>sens</i> &lt; fr. <i>sens</i> &lt; lat. <i>sensus</i></p>	<p><b>RO:</b> sens <b>ES:</b> sentido <b>GA:</b> sentido <b>CAT:</b> sentit <b>PT:</b> sentido <b>FR:</b> sens <b>IT:</b> senso</p>
<p><i>ne</i> (pronume personal, pers. I, pl., dativ) &lt; <i>nă</i> &lt; <i>nobis</i> (pronume personal, pers. I, pl., dativ)</p>	<p><b>RO:</b> ne <b>ES:</b> nos <b>GA:</b> nos <b>CAT:</b> ens <b>PT:</b> nos <b>FR:</b> nous <b>IT:</b> ci</p>
<p><i>a rămâne</i> (III) &lt; <i>remaneo</i>, <i>-ēre</i> (II)</p>	<p><b>RO:</b> rămâne <b>ES:</b> queda <b>GA:</b> queda <b>CAT:</b> queda <b>PT:</b> fica <b>FR:</b> reste <b>IT:</b> rimane</p>
<p><i>altceva</i> &lt; <i>alt</i> (&lt;<i>alter</i>) + <i>ceva</i> (<i>ce</i> &lt;<i>quid</i> + <i>va</i> &lt; <i>vare</i> &lt;<i>vare</i> &lt;*<i>voare</i> &lt;<i>volet</i>) (variante pentru <i>va</i> &lt; <i>vra</i> &lt; <i>vrea</i>)</p>	<p><b>RO:</b> altceva <b>ES:</b> otra <u>cosa</u> <b>GA:</b> outra <u>cousa</u> <b>CAT:</b> una altra <u>cosa</u> <b>PT:</b> outra <u>coisa</u> <b>FR:</b> autre <u>chose</u> <b>IT:</b> altro</p>
<p><i>a face</i> (III) &lt; <i>facere</i> (III)  <i>de făcut</i> (supin) &lt; prep. <i>de</i> + participiu trecut <i>făcut</i> &lt; *<i>facutus</i>, participiu trecut, forma pasivă, de la <i>facere</i></p>	<p><b>RO:</b> de făcut <b>ES:</b> por hacer <b>GA:</b> por facer <b>CAT:</b> per fer <b>PT:</b> por fazer <b>FR:</b> à faire <b>IT:</b> da fare</p>
<p><i>decât</i> &lt; <i>de</i> + <i>quantum</i></p>	<p><b>RO:</b> decât <b>ES:</b> que <b>GA:</b> que <b>CAT:</b> que <b>PT:</b> que <b>FR:</b> que <b>IT:</b> che</p>
<p><i>să</i> &lt; lat. <i>si</i>, *<i>se</i> (it. <i>se</i>); devine morfem al conjunctivului  <i>a da</i> (I) &lt; <i>dare</i> (I) <i>dăm</i> (indicativ prezent, pers. I pl.) &lt; <i>damus</i></p>	<p><b>RO:</b> să(-i) dăm <b>ES:</b> dar(le) <b>GA:</b> dar(lle) <b>CAT:</b> donar(-li) <b>PT:</b> dar(-lhe) <b>FR:</b> (lui) donner <b>IT:</b> dar(le)</p>
<p><i>-i</i> &lt; <i>îi</i> (pronume personal, pers. a III-a, sg., dativ) &lt; lat. <i>illi</i> (pronume demonstrativ, dativ)</p>	<p><b>RO:</b> (să)-i (dăm) <b>ES:</b> (dar)le <b>GA:</b> (dar)lle <b>CAT:</b> (donar)-li <b>PT:</b> (dar)-lhe <b>FR:</b> lui (donner) <b>IT:</b> (dar)le</p>

<p><i>un</i> &lt; lat. <i>unus</i> (pronume nehotarât) → lat. vulgară: valoare apropiată de cea a articolului din limbile romanice</p> <p><i>sens</i> &lt; fr. <i>sens</i> &lt; lat. <i>sensus</i></p>	<p><b>RO:</b> un sens <b>ES:</b> un sentido <b>GA:</b> un sentido <b>CAT:</b> un sentit <b>PT:</b> um sentido <b>FR:</b> un sens <b>IT:</b> un senso</p>
--	--

**RO:** După ce am descoperit că viața nu are nici un sens, nu ne rămâne altceva de făcut decât să-i dăm un sens.  
**ES:** Después de haber descubierto que la vida no tiene ningún sentido, no nos queda otra cosa por hacer que darle un sentido.  
**GA:** Depois de descubrir que a vida non ten ningún sentido, non nos queda outra cousa por facer que darlle un sentido.  
**CAT:** Després d'haver descobert que la vida no té cap sentit, no ens queda una altra cosa per fer que donar-li un sentit.  
**PT:** Depois de haver descoberto que a vida não tem nenhum sentido, não nos fica outra coisa por fazer que dar-lhe um sentido.  
**FR:** Après avoir découvert que la vie n'a aucun sens, il ne nous reste autre chose à faire que lui donner un sens.  
**IT:** Dopo avere scoperto che la vita non ha nessun senso, non ci rimane altro da fare che darle un senso.

**Listă abrevieri:** RO: română, ES: spaniolă, GA: galiciană, CAT: catalană,  
PT: portugheză, FR: franceză, IT: italiană

### Referințe bibliografice

- Andrade, Ana Isabel, Maria Helena de Araujo e Sa, Covadonga Lopez Alonso, Silvia Melo, Arlette Séré. (2005). *Manuel d'Instructions*, Projecto Socrates Lingua 2.
- Chavagne, Jean-Pierre. (2008). *L'intercompréhension en langue romanes – la plateforme Galanet* (ppt).
- Pagina oficială a platformei Galapro: [www.galanet.eu](http://www.galanet.eu)

# CONSIDERAȚII TEORETICE ASUPRA APLICABILITĂȚII UNEI BAZE DE DATE CU EXEMPLE DE TRADUCERE

NADIA LUIZA DINCĂ

*Institutul de Cercetare pentru Inteligența Artificială*

[hnadia\\_luiza@hotmail.com](mailto:hnadia_luiza@hotmail.com)

## Rezumat

Una dintre regulile interne ale traducerii bazate pe exemple este dependența calității traducerii de lungimea și modul de reprezentare a exemplelor de traducere. La rândul lor, acestea sunt gestionate de către o bază de exemple, în proiectarea căreia lingvistul este obligat să răspundă la două întrebări cheie:

- ce mod de reprezentare va alege pentru exemplul de traducere?
- care sunt posibilitățile de generalizare a exemplului de traducere stocat în baza de date?

În acest articol propun două posibile răspunsuri ale acestor întrebări, orientându-mă, pentru limbile română și engleză, spre reprezentarea exemplelor de traducere ca arbori de dependență și, respectiv, spre generalizarea lor prin informația semantică introdusă de clasele verbale descrise de Levin.

## 1. Introducere

În momentul în care a fost introdusă, (Nagao, 1984), traducerea bazată pe exemple era definită ca o traducere prin analogie, care utilizează o bază neadnotată de exemple, colectată, de regulă, dintr-un dicționar bilingv. Echivalențele erau exprimate sub forma perechilor de cuvinte, exceptând echivalențele verbale formalizate prin cadre cazuale.

Ulterior, în categoria sistemelor structurale de traducere bazată pe exemple, se introduce reprezentarea exemplelor de traducere ca arbori de dependență cu legături explicite între subarbori (incluzând **nodurile frunză** care corespund unităților lexicale). Aceste legături permit folosirea **fragmentelor de exemplu** sau **subarborilor** pentru recunoașterea corespondențelor exacte cu segmente sau structuri ale intrării în limba sursă, și pentru identificarea și combinarea unităților de traducere echivalente în limba țintă.

Sistemul de traducere automată MBT2, dezvoltat de S. Sato și M. Nagao în 1990, utilizează arborii de dependență pentru reprezentarea exemplelor de traducere și consideră trei operații de bază aplicabile subarborilor de dependență existenți în baza de date:

- operația de ștergere a unui subarbore;
- operația de înlocuire a unui subarbore cu o expresie corespondentă intrării de traducere;

c. operația de adăugare a unei expresii corespondente intrării de traducere ca fiică pentru nodul rădăcină al unui subarbore.

Existența mai multor unități de traducere candidate și generate prin backtracking solicită o selecție, al cărei principal criteriu este mărimea unității de traducere: se preferă, astfel, o unitate de traducere mai mare, accepția conceptului de mărime fiind numărul de noduri din unitatea de traducere. În principiu, sistemul MBT2 rezolvă compromisul dintre lungimea și similaritatea corespondențelor, considerat ca fiind baza creării de traduceri corecte. Problemele specifice acestei metode de lucru se referă la necesitatea unor calcule laborioase pentru determinarea scorului unității de traducere optimale și la obligativitatea existenței unui tezaur care să determine corect valorile de similitudine dintre cuvinte.

Sistemul de traducere automată bazată pe exemple propus de Kaji în 1992 se particularizează prin două subsisteme: învățarea de modele de traducere și, respectiv, traducerea bazată pe confruntarea modelelor cu datele.

Un model de traducere este o pereche de propoziții bilingve, în care unitățile echivalente (cuvinte și sintagme) sunt înlocuite prin variabile cărora le sunt asociate restricții sintactice și semantice. Se caută mai întâi acel model care unifică partea de limbă sursă cu o propoziție de intrare, apoi se înlocuiesc cuvintele și sintagmele prin variabilele modelului de traducere. Se face transferul pe limba țintă, iar cuvintele și sintagmele legate de variabile sunt traduse folosindu-se o metodă convențională.

Procedura de învățare a modelelor de traducere cuprinde etapa de antrenare a corpusului (pentru fiecare pereche de propoziții se construiesc modelele de traducere potrivite) și etapa de generare (se rafinează mulțimea de modele de traducere pentru a rezolva conflictele apărute în situații în care un model în limba sursă ar avea mai mulți candidați de traducere).

Kaji exemplifică rafinarea prin două șabloane generate pentru exemplele de traducere „play baseball” și „play the piano”. Cele două structuri generalizate sunt marcate prin categoriile semantice „sport” și „instrument”, astfel: (1) play X [NP / sport]; (2) play X [NP / instrument].

Combinarea unei metode de lucru bazate pe exemple cu o analiză a textului sursă formalizată prin reguli gramaticale reprezintă alternativa propusă de O. Furuse și H. Iida în 1992, respectiv 1994. Cunoștințele de transfer se aplică asupra șirului de intrare, executându-se simultan parsarea structurală și confruntarea șabloanelor cu datele.

Din perspectiva celor doi cercetători, un șablon elementar este o secvență formată din variabile și simboluri pentru fixarea granițelor constituenților lingvistici. Nu există simboluri gramaticale de tipul grupului nominal sau grupului verbal, ci un set de părți de vorbire specificate potrivit rolului gramatical (substantivul comun, substantivul propriu, verbul-a fi, verbul auxiliar, etc.). Adicional cuvintelor funcționale, un șablon de traducere folosește și bigramul parte de vorbire, exemplificat în expresia „I sing” prin construcția echivalentă „pron-verb” : I sing -> șablon::= **I pron-verb sing**.

Algoritmul propus de Furuse și Iida pentru identificarea șabloanelor de traducere cuprind, în esență, următorii pași:

a. asignarea de informații morfologice fiecărui cuvânt din propoziție;

- b. introducerea marcatorelor de graniță pentru constituenții parsați morfologic;
- c. derivarea structurilor posibile prin confruntarea șabloanelor cu datele, de la nivelul cel mai înalt de descriere lingvistică (propoziția introductivă), până la nivelurile inferioare: propoziția compusă, propoziția simplă, sintagma verbală, sintagma nominală, cuvântul compus.

Dincolo de diferențele dintre cele trei metode de reprezentare a exemplelor de traducere, se constată, ca invariantă, descompunerea propozițiilor în constituenți pentru a realiza o corespondență parțială în care părțile ce diferă se substituie prin variabile. Un avantaj al constituenților astfel obținuți îl constituie posibilitatea de prelucrarea a acestora în manieră independentă și, în consecință, flexibilitatea traducerii.

În acest articol, proiectarea unei baze de exemple de traducere, pentru limbile română și engleză, se realizează în maniera următoare:

- Exemplul de traducere este reprezentat prin arbori de dependență între care se stabilesc legături de corespondență. Sunt identificate, totodată, tipurile de relații sintactice de dependență dintre unitățile constituente ale unui grup verbal.
- În scopul generalizării, verbul primește o clasă semantică după tipologia creată de Levin. În situația în care există o breșă în găsirea unei corespondențe între șirul de intrare și subșirurile din baza de exemple, aceasta este rezolvată prin apelarea clasei semantice a verbului și, implicit, a listei de verbe cu care verbul căutat contractează o relație sinonimică.

## 2. *Reprezentarea exemplurilor de traducere*

### 2.1. *Descrierea exemplului de traducere*

Exemplul de traducere este un grup de cuvinte, uneori cu un înțeles diferit decât cel rezultat din însumarea înțelesurilor fiecărui cuvânt, căruia i se atribuie în limba țintă o traducere și un înțeles exact, favorizând o echivalență cu un nivel calitativ al traducerii ridicat.

Ca formă de reprezentare, un exemplu de traducere este compus din trei părți:

- un arbore de dependență în limba sursă (în acest articol, limba sursă este limba română);
- un arbore de dependență în limba țintă (limba engleză, în articol);
- legături de corespondență.

Cele trei părți sunt evidențiate în grupul verbal următor, extras din romanul lui G. Orwell, „1984”, subiect al unui amplu proiect lingvistic, **Multext-East**:

*își imaginase orice ↔ had imagined everything*  
*ro\_e ([ro1, [imagina, v],*  
*[ro1.1, [își, pron]],*  
*[ro2, [orice, pron]]])*

*en\_e* ([*en1*, [*have*, *aux*],  
           [*en1.2*, [*imagine*, *v*],  
           [*en2*, [*everything*, *pron*]]]])  
*clinks* ([[*ro1*, *en1*], [*ro2*, *en2*]])

Arborii de dependență *ro\_e* și *en\_e* afișează, pe fiecare linie, după prefixul de limbă, un număr (ro 1-2), (en 1-2), acesta etichetând un **nod** din subarbore ce conține forma bază a cuvântului și categoria sintactică asociată. Drumurile, în fiecare din cei doi subarbori construiesc, pe fragmentele traducibile posibile, unități de traducere (ro1- ro1.1, ro1-ro2, en1-en1.1, en1-en1.1-en2).

Fluxul procesului de traducere pleacă de la un arbore de dependență în limba sursă, pe care îl descompune, stabilește corespondențele sursă, realizează transferul către corespondențele țintă, după care le combină pentru a obține arborele de dependență echivalent în limba țintă.

## 2.2. Tipurile de relații de dependență

Toate unitățile constituente dintr-un enunț sunt aranjate de către vorbitor în construcții bine formate, pe baza dependențelor create între acestea: un cuvânt depinde de un altul prin poziția sa lineară și prin forma gramaticală.

Structura sintactică de suprafață, cea care interesează în lucrare, este un arbore ale cărui noduri sunt etichetate cu lexemele din propoziție, iar **arcele**, denumite și ramuri, primesc numele unei relații sintactice specifice, exemplificate mai jos.

Cele trei clase mari de dependențe sintactice, și anume: complementaritate, modificare și coordonare, organizează, la rândul lor, un număr mare de relații sintactice, dispuse la nivelul grupului verbal astfel:

### I. Relația de subordonare

#### a. obiect direct:

(*cumpărase* – **ob-dir** → *cartea*) ↔ (*bought* – **ob-dir** → [*the*] **book**)

(*luă* – **ob-dir** → [*o*] **țigară**) ↔ (*took* – **ob-dir** → [*a*] **cigarette**)

(*o* ← **ob-dir** – *ura*) ↔ (*hated* – **ob-dir** → **her**)

#### b. obiect indirect în Dativ

([*să*] *spună* – **ob-indir** → **i**) ↔ ([*should*] *tell* – **ob-indir** → **him**)

#### c. obiect prepozițional în Acuzativ

([*se simțea*] *atras* – **ob-prep** → **de** [*el*]) ↔ ([*felt*] *drawn* – **ob-prep** → **to** [*him*])

(*vorbea* – **ob-prep** → **despre** [*ea*]) ↔ (*referred* – **ob-prep** → **to** [*it*])

#### d. obiect infinitival

([*le*] *putea* – **ob-inf** → **vedea**) ↔ (*could* – **ob-inf** → **see**)

### II. Relația de coordonare

*scoase* – **ob-dir** → [*un*] **toc** – **coord** → [*o*] **sticlă** [*de cerneală*] – **coord** → **și** [*un volum*]  
 ↔ *took down* – **ob-dir** → [*a*] **penholder** – **coord** → [*a*] **bottle** [*of ink*] – **coord** → **and**  
 [*a book*]



### 2.3. Arborii de dependență și legăturile de corespondență

În crearea arborilor de dependență pentru grupurile verbale în limbile română și engleză am urmărit trei criterii de existență a unei relații sintactice de dependență între două cuvinte dintr-o propoziție:

- criteriul de conectivitate dintre două forme lexicale;
- criteriul de dominanță între două cuvinte;
- criteriul tipului specific de dependență sintactică dintre două lexeme.

Criteriile sunt dependente de limbă, ceea ce face uneori dificilă stabilirea unei corespondențe. Este cazul pronomelor reflexive din limba română, de exemplu, nerealizate, în planul enunțului, în limba engleză:

*își turnă o ceașcă de ceai ↔ poured out a teacupful*

*ro\_e* ([ro1, [turna, v],  
[ro2, ob-indir, [își, pron]],  
[ro3, ob-dir,  
[ro3.1, [ceașcă, n],  
[ro3.2, [o, art]],  
[ro3.3, [de, prep],  
[ro3.4, [ceai, n]]]]]])

*en\_e* ([en1,  
[en1.1., [pour, v],  
[en1.2, jonctiv, [out, prep]],  
[en2, ob-dir,  
[en2.1, [teacupful, n],  
[en2.2, [a, art]]]]]])

*clinks* ([[ro1,en1], [ro3, en2]])

Tipul de dependență specifică obiectului indirect este preluat, în limba engleză, de către subiect, ca agent al acțiunii descrise de verb. În același timp, criteriul de conectivitate pentru limba română admite o abatere de la regula generală a dispunerii lineare a cuvintelor, de aceea verbul, ca nucleu al sintagmei sintactice, va impune pronumelui relația de obiect indirect. Nodul 3 din limba română, extins în nodurile: 3.1.- „ceașcă”, 3.2.- „o”, 3.3.- „de”, 3.4.- „ceai”, are legături de corespondență cu nodul 2 din limba engleză, dezvoltat într-un nod părinte și un altul fiică.

În exemplul de traducere următor, pronumele reflexiv din română este realizat în engleză prin forma pronumelui personal, o explicație găsindu-se în structura *dativ+verb+substantiv*, unde formele neaccentuate de dativ, ale pronumelui reflexiv sau personal, exprimă ideea de posesie:

*își întinsese brațele către ecran ↔ extended her arms towards the screen*

*ro\_e* ([ro1, [întinde, v],  
           [ro2, posesie, [își, pron]],  
           [ro3, ob-dir, [brațe, n]],  
           [ro4, direcție,  
               [ro4.1, [cătred, prep],  
               [ro4.2, [ecran, n]]]])  
*en\_e* ([en1, [extend, v],  
           [en2, ob-dir,  
               [en2.1., posesie, [her, pron]],  
               [en2.2., [arms, n]],  
           [en3, direcție,  
               [en3.1, [towards, prep],  
               [en3.2, [screen, n],  
               [en3.3, [the, art]]]])  
*clinks* ([[ro1, en1], [ro3, en2.2], [ro4, en3]])

Nodul 2 din arborele în limba engleză cumulează două relații de dependență de rang diferit. Cea dominantă este de obiect direct, aplicată numelui comun, care impune o relație de posesie pronumelui personal în genitiv, întreaga structură construind, împreună cu verbul, unitatea traductibilă “extended her arms”. Subarborele în engleză respectă criteriul de aranjare lineară a formelor lexicale în scopul stabilirii legăturii de dependență sintactică. În schimb, unitatea traductibilă corespondentă în română “își întinsese brațele” consideră criteriul de dominanță sintactică pentru a identifica orientarea relației de dependență.

Generalizând structurile în care dispunerea lineară a cuvintelor cedează în fața dominanței sintactice, se observă două principale relații de dependență pe care verbul tranzitiv le impune formelor lexicale predecesoare: obiectul direct, respectiv obiectul indirect.

O altă situație gramaticală cu un regim aparte de echivalență între limbile română și engleză o constituie anticiparea sau reluarea obiectului direct sau indirect prin forme pronomiale personale. În exemplul de traducere următor, anticiparea se realizează prin pronumele personal neaccentuat “îl”, fără un corespondent lexical în limba engleză:

*îl văzuse pe O'Brien* ↔ *had seen O'Brien*  
*ro\_e* ([ro1, [vedea, v],  
           [ro2, ob-dir,  
               [ro2.1., binar, [îl, pron]],  
               [ro2.2., [pe, prep],  
               [ro2.3., [O'Brien, n]]]])  
*en\_e* ([en1, [had, aux]],  
           [en2, [see, v],

*[en3, [O'Brien, n]]])*  
*clinks ([[ro1,en2],[ro2, en3])*

Introducerea relației de binaritate pentru anticiparea obiectului direct este de natură să rezolve problema satisfacerii valenței verbale, deoarece verbul tranzitiv „a vedea” nu permite două componente directe. Pronumele „îl” și constructul „pe O'Brien”, cu dominanță pe lexemul prepozițional, au același rang sintactic și contribuie la complinirea verbului. În limba engleză, dependența binară nu mai există, verbul dominând un obiect direct realizat nominal.

### 3. Generalizarea exemplilor de traducere

#### 3.1. Preliminarii

Una dintre principalele probleme pe care trebuie să le înfrunte traducerea bazată pe exemple este necesitatea de a folosi un exemplu de traducere pentru mai mult de o situație de intrare. În mod obișnuit, primul pas îl constituie identificarea unei corespondențe între șirul lexical de intrare sau subșiruri ale acestuia și exemplele de traducere stocate în baza de date. Aceasta poate cauza uneori frustrări în ce privește calitatea traducerii, deoarece baza de exemple, oricât de complexă ar ajunge la un moment dat, nu reușește să acopere flexibilitatea lingvistică.

O posibilă soluție pentru acest inconvenient rezidă în combinarea relațiilor semantice și structurilor sintactice, astfel încât un lexem din exemplul de traducere să deschidă posibile instanțieri pentru lexemele din seriile semantice, urmând generarea relației de sinonimie.

Fie următoarea structură sintactică de tradus:

*Ceruse libertatea cuvântului.*

După deflexionare și dezambiguizare, etape ce nu formează subiectul propriu-zis al acestei lucrări, algoritmul trebuie să treacă la căutarea corespondențelor. În baza de exemple însă, verbul „a cere” nu intră în nicio combinație cu grupul nominal „libertatea cuvântului”. În schimb, structura nominală este identificată ca fiind în relație de dependență față de un alt verb, „a solicita”, în construcția sintactică: „solicitate libertatea cuvântului”, cu echivalent de traducere: „was advocating freedom of speech”. De aceea, la pasul următor se verifică existența unei relații semantice între verbele „a cere” și „a solicita”. Se identifică astfel clasa semantică a *verbelor de transfer al unui mesaj*<sup>1</sup>, clasă ce instanțiază, pentru limba română, lexemul verbal solicita:5, iar pentru limba engleză preach:2, advocate:2. Se oprește căutarea și se validează corespondența dintre „ceruse libertatea cuvântului”, respectiv „solicitate libertatea cuvântului”.

Trebuie precizat faptul că relațiile semantice pentru limba română sunt identificate cu ajutorul dicționarului de sinonime, iar pentru limba engleză prin intermediul ontologiei lexicale WordNet, interogată cu ajutorul editorului Visdic, versiunea 1.3.50.

<sup>1</sup> În taxonomia descrisă de Beth Levin, este vorba de clasa 37.1- *Verbs of Transfer of a Message*.

### 3.2. Rolul introducerii relațiilor semantice la nivelul sintagmei verbale

Există o relație profundă între proprietățile semantice ale unui verb și cele sintactice, astfel explicându-se de ce, dezambiguizând înțelesul unui lexem, vorbitorii completează modelul comunicațional intuindu-i tiparul sintactic.

Aceasta este, în esență, și motivația introducerii relațiilor semantice la nivelul sintagmei verbale. Verbul ce guvernează o relație de dependență nu este izolat în mulțimea tuturor verbelor, ci este actant al unei relații de sinonimie cu alte lexeme verbale. Nu toate sensurile verbelor participă la crearea sinsetului, ci doar acelea care sunt ordonate în jurul unui înțeles comun.

Precizarea tipurilor de relații de dependență sintactică și crearea sinseturilor grupate în clase verbale sunt de natură să stimuleze calitatea traducerii prin mai buna adaptare la flexibilitatea limbii și la condițiile de bună formare a unei propoziții.

În exemplul de traducere următor sunt puse în evidență proprietățile sintactice și semantice ale centrului verbal:

*dădea o muzică stridentă, militărească.* ↔ *had played a strident military music*

*ro\_e* (*[ro1, [Verbe de Reprezentație-> da:9, transmite:13], [da, v],*

*[ro2, ob-dir,*

*[ro2.1, [muzică, n],*

*[ro2.2, [o, art]],*

*[ro2.3, [stridentă, adj]],*

*[ro2.4, [militărească, adj]]]])*

*en\_e* (*[en1, [have, aux],*

*[en1.1, [Performance Verbs -> play:7, perform:3], [play, v],*

*[en2, ob-dir,*

*[en2.1, [music, n],*

*[en2.2, [a, art]],*

*[en2.3, [strident, adj]],*

*[en2.4, [military, adj]]]])*

*clinks*(*[[ro1, en1], [ro2, en2]])*

Crearea arborilor de dependență se realizează, după cum se observă mai sus, prin descrierea, pentru verbul principal a clasei semantice, a sinsetului asociat și a tipurilor de dependență guvernate de verb. Se generalizează astfel posibilitățile de corespondență între șirul de intrare și exemplele din baza de date, dar se impune, totodată, un filtru de validare a acestora din perspectiva respectării tipurilor de relații de dependență. Dintr-o mulțime de candidați potențiali la stabilirea corespondenței, sunt selectați doar cei care domină aceleași tipuri de dependență sintactică. Împreună, cele două descrieri - sintactică și respectiv, semantică- au rol în dezambiguizarea traducerii.

#### 4. Concluzii

Acest articol prezintă câteva considerații teoretice despre aplicabilitatea unei baze de exemple de traducere. În dezvoltarea lui, am plecat de la premisa relațiilor sintactico-semantică dintre cuvinte urmărind două idei esențiale, și anume utilitate și generalizare. Avem, astfel, pe de o parte, relațiile de dependență sintactică dintre verb și celelalte valori morfo-sintactice dominate de el, iar pe de altă parte, relația de sinonimie dintre verb și lexemele din același sinset, respectiv clasă verbală.

La prima vedere, o bază de exemple de traducere încărcată cu toate aceste informații poate fi dificil de manipulat, prin mărimea numărului de căutări. În vederea diminuării acestui inconvenient propun trei criterii de selecție a exemplurilor pentru a construi o bază de date. În primul rând trebuie selectate expresiile și structurile considerate a fi cele mai frecvente în limbile sursă și țintă. Odată constituit acest nucleu, urmează, cu prioritate de rang doi, secvențele propoziționale al căror înțeles este diferit de compunerea sensurilor unităților constituente. Cu prioritate de rang 3 pentru completarea bazei de date sunt grupurile verbale extrase din corpusul „1984”.

În condițiile asigurării unei bune acoperiri lexicale și a unor structuri ordonate semantic, respectiv sintactic, se poate considera că și numărul de căutări în baza de exemple este mai mic. În același timp, crește posibilitatea identificării printre lexemele verbale cele mai frecvente, a acelor care sunt în relație de sinonimie cu verbul de intrare.

Un alt avantaj al proiectării bazei de exemple prin relații de sinonimie, filtrate de relații de dependență sintactică, se găsește în dezambiguizare, semantică și sintactică. Există verbe care au, fără îndoială, mai multe sensuri, unele dintre ele particularizând relații de dependență diferite. În momentul în care sinonimia dintre sensurile a două verbe este evaluată ca având aceleași tipuri de dependență sintactică, se reține un anumit sens din mulțimea de candidați ai verbului de intrare. Operația este validă și pentru argumentele selectate de verb: dacă verbele aparțin unui sinset, iar unul are o linie sintagmatică deja cunoscută, celălalt îi imită comportamentul sintactic.

#### Referințe bibliografice

- Levin, B. (1993). *English Verb Classes and Alternations- A Preliminary Investigation*, The University of Chicago Press.
- Mel'čuk, I. (2003). Levels of Dependency in Linguistic Description: Concepts and Problems. In V. Agel, L. Eichinger, H.-W. Eroms, P. Hellwig, H. J. Herringer, H. Lobin (eds): *Dependency and Valency. An International Handbook of Contemporary Research*, vol. 1, Berlin - New York, W. de Gruyter, 188-229.
- Multext-East Home Page: <http://nl.ijs.si/ME/>
- Nagao, M. (1984), A framework of a mechanical translation between Japanese and English by analogy principle, *Proceedings of the international NATO symposium on Artificial and human intelligence*, Lyon, France, 173 – 180.

- Furuse, O, Iida, H. (1992), Cooperation between transfer and analysis in example-based framework. *Proceedings of the 14th International Conference on Computational Linguistics*, Nantes, France. ACL, 645-651.
- Furuse, O., Iida, H. (1994), Constituent boundary parsing for example-based machine translation, *Proceedings of the 15th conference on Computational linguistics*, vol. 1, Kyoto, Japan, 105-111.
- Kaji H., Kida, Y., Morimoto, Y. (1992), Learning translation templates from bilingual text, *Proceedings of the Fourteenth International Conference on Computational Linguistics*, Nantes, 672-678,.
- Sato, S., Nagao, M. (1990), Towards memory-based translation. *Proceedings of COLING-90*, Helsinki, Finland, vol. 3, 247-252.
- Seche, M., Seche, L. (2002), *Dictionar de sinonime*, Editura Litera Internațional

Visdic: <http://nlp.fi.muni.cz/projekty/visdic/>

## **CAPITOLUL 3**

### **APLICAȚII ALE TEHNOLOGIILOR LINGVISTICE**





# EVALUAREA RĂSPUNSURILOR OFERITE DE UN SISTEM DE TIP ÎNTREBARE RĂSPUNS PENTRU LIMBA ROMÂNĂ

ADRIAN IFTENE, ANCUȚA ROTARU, DANA-ALINA MARCU

*Universitatea "Al.I.Cuza", Facultatea de Informatică, Iași – România*

*{adiftene, ancuta.rotaru, dana.marcu}@info.uaic.ro*

## Rezumat

În cadrul competiției QA@CLEF2008<sup>1</sup> una din principalele provocări a fost exercițiul de validare a răspunsurilor AVE (Answer Validation Exercise). Lansat cu 3 ani în urmă acesta are ca scop evaluarea răspunsurilor oferite de către un sistem de tip întrebare răspuns, dorind astfel să mărească și calitatea acestora. Anul acesta, ca și în ediția din 2007, printre cele 5 limbi participante a fost prezentă și limba română, iar noi am participat pentru prima dată cu un sistem dedicat acesteia.

Articolul de față prezintă structura exercițiilor, principalele componente ale sistemului construit de noi pentru competiția din acest an, precum și rezultatele acestei ediții. Important de remarcat este faptul că pe 3 limbi participante (engleză, germană și română) sistemele AVE au obținut rezultate mai bune decât sistemele de tip Întrebare-Răspuns, în ordonarea răspunsurilor oferite de un sistem de tip Întrebare-Răspuns.

## 1. Introducere

AVE<sup>2</sup> a avut loc pentru prima oară în 2006 (Peñas et al., 2007) din nevoia de a promova dezvoltarea și evaluarea sub-sistemelor care aveau ca scop validarea corectitudinii răspunsurilor oferite de sistemele de tip Întrebare-Răspuns (ÎR). Din start s-a dorit ca AVE să îmbunătățească și calitatea sistemelor de tip ÎR, dar în primul rând să verifice dacă răspunsurile alese corespund sau nu fragmentelor de texte ajutătoare existente.

De la an la an metodologia de evaluare s-a modificat în încercarea de a surprinde cât mai bine factorii care ar duce efectiv la îmbunătățirile sistemelor de tip ÎR. Astfel, în 2007 sistemele trebuiau să selecteze doar un singur răspuns valid pentru fiecare întrebare dintr-o mulțime de răspunsuri posibile, spre deosebire de ediția din 2006 când era posibil să se aleagă mai multe răspunsuri valide. În 2008, s-a observat că această metodologie are o problemă: nu se știa cum se vor comporta sistemele în cazul în care toate răspunsurile posibile ar fi fost incorecte. Se dorea ca ele să poată cere alte răspunsuri de la sistemele de tip ÎR, în speranța că vor putea obține măcar un răspuns corect. Ediția de anul acesta a avut ca obiectiv eliminarea acestor neajunsuri descoperite în edițiile precedente.

În continuare vom prezenta caracteristicile competiției de anul acesta și modul în care am construit sistemul folosit de noi pe limba română. În partea de final vom prezenta rezultatele și concluziile.

---

<sup>1</sup> CLEF: <http://www.clef-campaign.org/2008.html>

<sup>2</sup> AVE: <http://nlp.uned.es/clef-qa/ave/>

## 2. Descrierea Exercițiului

### 2.1. Formatul datelor de intrare

Urmărind tiparul propus în ediția din 2007, în ediția din 2008 (Rodrigo et al., 2008) sistemele trebuiau să ia în considerare triplete de forma (*Întrebare, Răspuns, Fragment de Text*) și să hotărască dacă răspunsul la întrebare este corect și poate fi dedus din fragmentul de text atașat. Astfel, pentru fiecare tripletă de această formă, participanții trebuie să stabilească o valoare care are ca semnificație faptul că tripletul este validat sau respins. Pentru limba română, fișierul de intrare conține 119 întrebări, iar pentru fiecare întrebare sunt între 1 și 9 răspunsuri posibile, în total fiind 497 de triplete.

Tabel 1: Formatul datelor de intrare

```

<q id="1" lang="RO">
  <q_str>Câte zile avea aprilie înainte de 700 î.Hr.?</q_str>
  <a id="0001_1" value="">
    <a_str>30</a_str>
    <t_str doc="Aprilie.html">
      Înainte de anul 700 î.Hr., luna aprilie era a doua lună a anului în calendarul roman și
      avea 29 de zile. După ce Iuliu Cezar a introdus calendarul iulian în 45 î.Hr., luna aprilie
      avea 30 de zile și devenea a patra lună a anului.
    </t_str>
  </a>
  <a id="0001_6" value="">
    <a_str>29 de zile</a_str>
    <t_str doc="1">
      Hr., luna aprilie era a doua lună a anului în calendarul roman și avea 29 de zile. După ce
      Iuliu Cezar a introdus calendarul iulian în 45 î .
    </t_str>
  </a>
  <a id="0001_7" value="">
    <a_str>de anul 700</a_str>
    <t_str doc="1">
      Numele lunii aprilie (latină: Aprilis ) vine de la cuvântul latinesc aperio, ire = a
      deschide, deoarece în aprilie se deschid mugurii plantelor. Înainte de anul 700 î .
    </t_str>
  </a>
  <a id="0001_8" value="">
    <a_str>cu aceeași zi a săptămânii în toți anii</a_str>
    <t_str doc="1">
      Aprilie începe cu aceeași zi a săptămânii ca și Iulie în toți anii și ca Ianuarie în anii
      bisecți.
    </t_str>
  </a></q>

```

În tabelul 1 putem vedea formatul datelor de intrare (unde tag-ul “**q\_str**” conține întrebarea, tag-urile “**a**” corespund fiecărui răspuns posibil, acesta fiind propriu-zis conținut în tag-ul “**a\_str**”, iar fragmentele de text apar în tag-ul “**t\_str**”).

## 2.2. *Formatul datelor de ieșire*

Răspunsurile oferite de participanți trebuie să fie în următorul format:

*q\_id a\_id [VALIDAT|SELECTAT|RESPINS] scor\_de\_încredere*

unde semnificația răspunsurilor este următoarea:

- **VALIDAT**: indică faptul că răspunsul este corect și este suportat de paragraful de text asociat. Nu există nici o restricție asupra numărului de răspunsuri validate (pot fi toate validate sau nici unul).
- **SELECTAT**: indică faptul că răspunsul este VALIDAT și reprezintă cel mai probabil răspuns al unui posibil sistem de tip ÎR. Fiecare întrebare va avea doar un singur răspuns selectat. Cel puțin unul dintre răspunsurile valide trebuie să fie selectat.
- **RESPINS**: indică faptul că răspunsul este incorect (sau că nu există suficiente dovezi care să-i demonstreze corectitudinea). Nu există nici o restricție asupra numărului de răspunsuri respinse (pot fi toate sau nici unul).
- **scor\_de\_încredere**: Opțional, pentru fiecare tripletă se poate acorda un scor de încredere (care poate lua valori din intervalul [0, 1]): unde 0 – reprezintă faptul că suntem nesiguri de răspunsul dat, iar 1 – reprezintă faptul că suntem siguri de răspunsul oferit.

## 2.3. *Proveniența datelor de intrare*

Ca și în edițiile precedente ale competiției AVE, datele folosite în antrenarea și testarea sistemelor de apreciere a răspunsurilor provin din fișierele cu evaluarea sistemelor de tip ÎR folosite în competiția QA@CLEF, cu unele mici completări și modificări.

Transformarea evaluării răspunsurilor sistemelor de tip ÎR în date de test pentru competiția AVE2008 s-a făcut în modul următor (Rodrigo et al., 2008):

- un răspuns care a fost evaluat CORECT în competiția sistemelor de tip ÎR va fi evaluat ca fiind VALID în datele de test AVE;
- un răspuns evaluat ca fiind GREȘIT sau NESUPORTAT în QA@CLEF va fi considerat RESPINS în AVE;
- un răspuns evaluat ca fiind INEXACT sau NEEVALUAT în QA@CLEF va avea valoarea NECUNOSCUȚ în AVE (și nu va fi considerat în evaluarea sistemelor).

Deoarece colecția datelor de test pentru competiția AVE s-a construit pe baza tuturor fișierelor trimise de participanții pe o anumită limbă, a fost nevoie să se stabilească niște reguli pentru îmbunătățirea calității acestora:

- eliminarea răspunsurilor redundante;

- dacă între răspunsurile posibile pentru o întrebare există răspunsuri care se conțin unele pe altele se va recurge la următoarea abordare: se vor elimina răspunsurile care au lungimea cea mai scurtă;
- eliminarea întrebărilor fără răspuns (care au răspuns nul).

Spre deosebire de competiția sistemelor de tip întrebare-răspuns, unde întrebările au fost grupate pe domenii, întrebările nu au fost grupate în același mod în competiția AVE.

Testarea colecțiilor folosite s-a făcut în 9 runde de test, pentru fiecare limbă fiind generat un test individual. Pentru limba română sistemul creat a avut următoarele rezultate pentru 119 întrebări și 497 de răspunsuri:

- 48,58 % de răspunsuri au fost evaluate ca fiind VALIDE din mulțimea de răspunsuri posibile;
- 52 de răspunsuri SELECTATE, 406 de răspunsuri RESPINSE și 39 de răspunsuri NECUNOSCUTE.

### 3. Sistemul pentru limba română construit pentru competiția AVE

Structura sistemului construit pentru limba română este asemănătoare cu structura sistemului construit pe limba engleză (Iftene, Balahur-Dobrescu, 2008). Diferențele față de acesta sunt legate de modulele și resursele folosite. Sistemul primește la intrare triplete de forma (*întrebare, răspuns, fragment de text*) și oferă la ieșire evaluarea fiecărui răspuns în parte (Vezi figura 1 de mai jos).

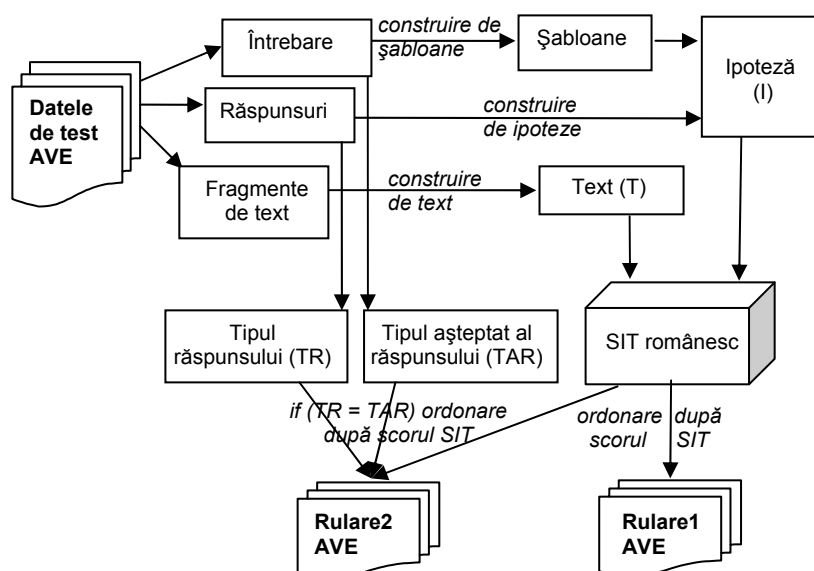


Figura 1: Sistemul AVE românesc

Principalele componente ale acestui sistem realizează următorii pași:

- Construiesc *ipoteze* necesare sistemului de *inferențe textuale* (SIT) folosind șabloanele construite din întrebări și răspunsurile din datele de test AVE.
- Consideră fragmentele de text ca fiind *textul* necesar unui SIT.

EVALUAREA RĂSPUNSURILOR OFERITE DE UN SISTEM DE TIP ÎNTREBARE RĂSPUNS  
PENTRU LIMBA ROMÂNĂ

- Calculează folosind SIT-ul românesc un scor de potrivire pentru fiecare pereche (*Text, Ipoteză*).
- Prelucreează întrebarea cu tehnici specifice sistemelor de tip întrebare-răspuns și identifică *tipul așteptat al răspunsului* (TAR).
- Aplică șabloane specifice și resurse de entități de tip nume pentru a identifica *tipul răspunsului* (TR).
- În final, pe baza scorului de potrivire, a tipului răspunsului și al tipului așteptat al răspunsului am trimis două rulări.

Vom vedea în continuare modul de funcționare a fiecărei componente din sistemul prezentat mai sus.

### 3.1. Construirea șabloanelor

Pentru a putea folosi sistemul de inferențe textuale am construit din întrebările inițiale o mulțime de șabloane folosind tehnici asemănătoare celor prezentate în (Bar-Haim et al., 2006). Astfel, pentru întrebarea 1 din datele de test:

**Întrebarea 1:** *Câte zile avea aprilie înainte de 700 î.Hr.?*

șablonul construit are forma:

**Șablon 1:** *Aprilie înainte de 700 î.Hr. avea NUMĂR zile.*

unde NUMĂR reprezintă o variabilă ce va fi înlocuită cu toate răspunsurile posibile pentru această întrebare. Pentru limba română am considerat șabloane specifice pentru următoarele tipuri de răspunsuri: DATA CALENDARISTICĂ (DATA, AN), DEFINIȚIE, MĂSURĂ, LOCAȚIE (ȚARĂ, ORAȘ), NUMĂR, PERSOANĂ, ORGANIZAȚIE, ALTCEVA. Se observă că atunci când a fost posibil am identificat tipuri cât mai specifice pentru tipul răspunsului. Tabelul de mai jos prezintă exemple de astfel de șabloane pentru fiecare tip în parte:

Tabel 2: Șabloane asociate întrebărilor

Tipul așteptat al răspunsului	Exemplu de întrebare	Șablon
NUMĂR	Câți jucători participă la jocul de bridge?	NUMĂR jucători participă la jocul de bridge.
MĂSURĂ	Ce lungime are Biserica Neagră din Brașov?	Biserica Neagră din Brașov are MĂSURĂ.
LOCAȚIE	Unde s-a născut Emil Constantinescu?	Emil Constantinescu s-a născut în LOCAȚIE.
ORAȘ	În ce oraș s-a născut Charlie Chaplin?	Charlie Chaplin s-a născut în ORAȘ.
PERSOANĂ	Ce zeiță, soră a lui Ares, este fiica lui Metis?	PERSOANĂ, soră a lui Ares, este fiica lui Metis.
ORGANIZAȚIE	Din ce organizație teroristă face parte Osama bin Laden?	Din ORGANIZAȚIE teroristă face parte Osama bin Laden.
DATA	Când s-a vândut primul produs	La DATA s-a vândut primul

	Apple?	produs Apple.
AN	În ce an a fost produs filmul românesc Furia?	În AN a fost produs filmul românesc Furia.
ALTCEVA	Din ce se produce cașcavalul?	Cașcavalul se produce din ALTCEVA.

O situație deosebită a fost pentru cazul întrebărilor de tip definiție. În acest caz am considerat doar răspunsul ca fiind ipoteza ce va fi trimisă sistemului de inferențe textuale, fără a mai lua ceva din întrebarea inițială.

### 3.2. Construirea ipotezelor și a textelor

În șabloanelor construite ca mai sus am înlocuit variabilele folosind răspunsurile din datele de intrare și am construit ipotezele. Astfel, pentru întrebarea 1 în șablonul “*Aprilie înainte de 700 î.Hr. avea NUMĂR zile.*” am înlocuit variabila NUMĂR cu toate cele 4 valori posibile ale răspunsurilor corespunzătoare din tabelul 1. În urma înlocuirii am obținut cele 4 ipoteze de mai jos:

*I<sub>1\_1</sub>: Aprilie înainte de 700 î.Hr. avea 30 zile.*

*I<sub>1\_6</sub>: Aprilie înainte de 700 î.Hr. avea 29 de zile zile.*

*I<sub>1\_7</sub>: Aprilie înainte de 700 î.Hr. avea de anul 700 zile.*

*I<sub>1\_8</sub>: Aprilie înainte de 700 î.Hr. avea cu aceeași zi a săptămânii în toți anii zile.*

Pentru aceste ipoteze, cele 4 texte le obținem din fragmentele de text corespunzătoare din tabelul 1:

*T<sub>1\_1</sub>: Înainte de anul 700 î.Hr., luna aprilie era a doua lună a anului în calendarul roman și avea 29 de zile. După ce Iuliu Cezar a introdus calendarul iulian în 45 î.Hr., luna aprilie avea 30 de zile și devenea a patra lună a anului.*

*T<sub>1\_6</sub>: Hr., luna aprilie era a doua lună a anului în calendarul roman și avea 29 de zile. După ce Iuliu Cezar a introdus calendarul iulian în 45 î.*

*T<sub>1\_7</sub>: Numele lunii aprilie (latină: Aprilis) vine de la cuvântul latinesc aperio, ire = a deschide, deoarece în aprilie se deschid mugurii plantelor. Înainte de anul 700 î.*

*T<sub>1\_8</sub>: Aprilie începe cu aceeași zi a săptămânii ca și Iulie în toți anii și ca Ianuarie în anii bisecți.*

### 3.3. Folosirea sistemului de inferențe textuale pentru limba română

Sistemul de inferențe textuale folosit pentru limba română (Iftene, Balahur-Dobrescu, 2007) primește la intrare perechi de tip (*ipoteză, text*) și oferă la ieșire un scor de potrivire, iar în plus precizează dacă există probleme cu entitățile de tip nume. Problemele de acest tip apar în cazurile în care în ipoteză avem o entitate de tip nume căreia nu-i găsim corespondent în text. Acest lucru se întâmplă dacă entitatea cu probleme apare în întrebare, dar nu apare în fragmentul de text care ar trebui să justifice alegerea răspunsului curent, sau în cazul în care entitatea nu apare în fragmentul de text justificator. În ambele cazuri considerăm ca nejustificată alegerea răspunsului curent și stabilim răspunsul final ca fiind RESPINS.

În tabelul de mai jos avem scorurile asociate celor 4 perechi (text, ipoteză) de mai sus:

EVALUAREA RĂSPUNSURILOR OFERITE DE UN SISTEM DE TIP ÎNTREBARE RĂSPUNS  
PENTRU LIMBA ROMÂNĂ

Tabel 3: Scorurile asociate perechilor (T, I) corespunzătoare întrebării 1

Perechea	Scor de potrivire	Entitatea de tip nume cu probleme
(T <sub>1 1</sub> , I <sub>1 1</sub> )	0.727	-
(T <sub>1 6</sub> , I <sub>1 6</sub> )	0.889	-
(T <sub>1 7</sub> , I <sub>1 7</sub> )	0.636	î.Hr.
(T <sub>1 8</sub> , I <sub>1 8</sub> )	0.563	î.Hr.

### 3.4. Identificarea tipurilor răspunsurilor și a tipului așteptat al răspunsurilor

Scopul acestui pas este de a elimina din start cazurile în care aceste valori sunt diferite pentru întrebarea curentă și un răspuns curent al acesteia.

Pentru identificarea *tipurilor răspunsurilor* (TR) pentru limba română am folosit din GATE<sup>3</sup> următoarele tipuri de entități de tip nume: Oraș, Companie, Țară, Organizație, Persoană, Regiune. În plus am folosit șabloane specifice pentru identificarea NUMERELOR, DATELOR CALENDARISTICE, ANILOR și a MĂSURILOR.

La identificarea *tipului așteptat al răspunsului* (TAR) am utilizat aceleași valori ca cele folosite la construirea șablonelor din întrebări, prezentate în tabelul 2: DATA CALENDARISTICĂ (DATA, AN), DEFINIȚIE, MĂSURĂ, LOCAȚIE (ȚARĂ, ORAȘ), NUMĂR, PERSOANĂ, ORGANIZAȚIE, ALTCEVA.

Pentru întrebarea 1 avem următoarele valori:

Tabel 4: Întrebarea 1 - Valoarea TAR și valorile TR asociate răspunsurilor

TAR	Răspuns	TR	Scor potrivire între TAR și TR
NUMĂR	30	NUMĂR	1
	29 de zile	MĂSURĂ	0.5
	de anul 700	ALTCEVA	0.25
	cu aceeași zi a săptămânii în toți anii	ALTCEVA	0.25

unde scorul de potrivire dintre TAR și TR s-a calculat similar modului în care am calculat această valoare pentru limba engleză. Tabelul 5 ne prezintă principalele situații întâlnite:

Tabel 5: Calcularea scorului de potrivire dintre TAR și TR

Situație	Scor de potrivire
TAR = TR	1
(TAR = "DEFINIȚIE") și (TR = "ALTCEVA")	1
TAR și TR sunt în aceeași clasă de entități: {ORAȘ, ȚARĂ, REGIUNE, LOCAȚIE} sau {AN, DATA} sau {NUMĂR, MĂSURĂ, AN}	0.5
(TR = "ALTCEVA") sau (TAR = "ALTCEVA")	0.25
În celelalte cazuri	0

<sup>3</sup> GATE: <http://www.gate.ac.uk/>

### 3.5. Caracteristicile rulărilor trimise

Pe limba română am trimis două rulări, diferența dintre ele constând în faptul că am folosit sau nu tabelul 5 pentru a compara valorile TAR cu valorile TR.

**Rularea 1:** nu folosește comparația dintre TAR și TR, ci doar ieșirea oferită de SIT-ul românesc. Astfel, răspunsurile pentru care avem probleme cu entitățile de tip nume sunt considerate ca fiind RESPINSE (cazurile răspunsurilor cu id-urile 7 și 8 de la întrebarea 1). Toate celelalte sunt considerate ca fiind VALIDATE (cazurile răspunsurilor cu id-urile 1 și 6 de la întrebarea 1). Ca SELECTAT este considerat răspunsul VALIDAT care are scorul de potrivire cel mai mare oferit de SIT-ul românesc (răspunsul cu id-ul 6 de la întrebarea 1). Toate valorile pentru întrebarea 1 sunt prinse în acest caz în tabelul 6. Putem observa cum în acest caz din cele 4 răspunsuri am obținut valoarea corectă în 3 cazuri.

Tabel 6: Rularea 1: Valorile obținute pentru răspunsurile de la întrebarea 1

Răspuns	Valoare obținută	Valoare corectă
30	VALIDAT	RESPINS
29 de zile	SELECTAT	VALIDAT
de anul 700	RESPINS	RESPINS
cu aceeași zi a săptămânii în toți anii	RESPINS	RESPINS

**Rularea 2:** ca mai sus, dar folosește în plus comparația dintre TAR și TR. Comparația dintre TAR și TR este folosită astfel: se consideră răspunsuri RESPINSE cele care au probleme cu entitățile de tip nume sau cele pentru care scorul de potrivire dintre TAR și TR este 0. Pentru întrebarea 1, deoarece în tabelul 4 nu avem scoruri de potrivire 0, vom considera ca fiind RESPINSE cazurile răspunsurilor cu id-urile 7 și 8, care au probleme cu entitățile de tip nume. Restul răspunsurilor se consideră ca fiind VALIDE (cazurile răspunsurilor cu id-urile 1 și 6 de la întrebarea 1). Dintre răspunsurile VALIDE pentru a decide care dintre răspunsuri este SELECTAT folosim și comparația dintre TAR și TR. Vom considera ca fiind SELECTAT, răspunsul cu cel mai mare scor de potrivire întors de SIT-ul românesc, dintre cele cu cel mai mare scor de potrivire dintre TAR și TR. Deoarece avem un singur răspuns cu cel mai mare scor de potrivire dintre TAR și TR, răspunsul 1, care are scorul de potrivire 1, este și cel ales ca fiind selectat. Valorile pentru întrebarea 1 sunt prinse în tabelul 7 de mai jos. De observat faptul că numărul răspunsurilor în care am dat răspunsul corect este 3.

Tabel 7: Rularea 2: Valorile obținute pentru răspunsurile de la întrebarea 1

Răspuns	Valoare obținută	Valoare corectă
30	SELECTAT	RESPINS
29 de zile	VALIDAT	VALIDAT
de anul 700	RESPINS	RESPINS
cu aceeași zi a săptămânii în toți anii	RESPINS	RESPINS

## 4. Rezultate

Rezultatele oficiale obținute pentru limba română sunt prezentate mai jos:



EVALUAREA RĂSPUNSURILOR OFERITE DE UN SISTEM DE TIP ÎNTREBARE RĂSPUNS  
PENTRU LIMBA ROMÂNĂ

Tabel 8: Rezultatele obținute în competiția AVE2008 pe limba română

Rezultate	Rularea 1	Rularea 2
<i>F-measure</i>	0.22	0.23
<i>Precizia</i>	0.12	0.13
<i>Recall</i>	0.92	0.92
<i>qa accuracy</i>	0.17	0.24
<i>estimated qa performance</i>	0.17	0.25

unde *precizia*, *recall* și *F-measure* au ca scop evaluarea unui sistem care ordonează și filtrează răspunsurile. Formulele care au fost aplicate pentru acestea sunt următoarele:

$$precizie = \frac{\text{prezise\_corect\_de\_sistem\_ca\_SELECTATE\_sau\_VALIDATE}}{\text{prezise\_de\_sistem\_ca\_SELECTATE\_sau\_VALIDATE}}$$

$$recall = \frac{\text{prezise\_corect\_de\_sistem\_ca\_SELECTATE\_sau\_VALIDATE}}{\text{multimea\_raspunsurilor\_SELECTATE\_sau\_VALIDE}}$$

$$F - measure = \frac{2 \times recall \times precizie}{recall + precizie}$$

Iar *qa accuracy* și *estimated qa performance* au ca scop compararea performanțelor unui sistem de tip ÎR cu un sistem ipotetic de tip ÎR care ar folosi și un sistem AVE, formulele aplicate fiind:

$$qa\_accuracy = \frac{\text{Raspunsuri\_SELECTATE\_Corect}}{\text{Numarul\_Intrebarilor}}$$

$$qa\_rej\_accuracy = \frac{\text{Raspunsuri\_RESPINSE\_Corect}}{\text{Numarul\_Intrebarilor}}$$

$$estimated\_qa\_performance = qa\_accuracy + qa\_accuracy * qa\_rej\_accuracy$$

După cum se observă din tabelul 8 rularea a doua, care folosește și comparația dintre TAR și TR, este mai bună. Analizând rezultatele obținute am putut observa cum din cele 52 de răspunsuri VALIDE aflate în fișierul de test sistemul nostru oferă 48 dintre ele, din care 28 sunt SELECTATE (de aici valoarea foarte mare a *recall*-ului). Precizia mică obținută se datorează faptului că sistemul nostru prin modul în care e construit are condiții foarte stricte pentru a stabili dacă un răspuns este RESPINS și prin urmare oferă foarte multe răspunsuri VALIDE și SELECTATE și foarte puține răspunsuri RESPINSE. Pe de altă parte, trebuie să precizăm faptul că din cele 73 de răspunsuri RESPINSE date de sistemul nostru, 69 au fost corecte.

## 5. Concluzii

Lucrarea prezintă principalele componente ale sistemului folosit în competiția AVE de anul acesta pe limba română. Sistemul construit a avut o comportare foarte bună pentru răspunsurile VALIDE, dar nu a tratat aproape deloc răspunsurile RESPINSE.

De remarcat este faptul că din cele două rulări, rularea a doua are valoarea *estimated qa performance* de 0.25 care este superioară preciziei celui mai bun sistem din competiția sistemelor de tip întrebare-răspuns pe limba română de anul acesta. Acest

lucru ne indică faptul că folosirea acestui sistem de ordonare a răspunsurilor în cadrul unui sistem de tip Întrebare-Răspuns ar duce la creșteri semnificative ale preciziei pentru un astfel de sistem.

Pe viitor avem ca principal obiectiv eliminarea a două neajunsuri majore ale sistemului: Primul este datorat faptului că am considerat o mulțime de șabloane pentru identificarea tipului răspunsului, iar situațiile noi care pot apare ar fi tratate incorect. Cel de-al doilea este datorat faptului că am considerat condiții prea stricte pentru identificarea răspunsurilor RESPINSE, iar numărul acestora este prea mic.

**Mulțumiri.** Mulțumim colegilor de la Facultatea de Informatică care ne-au ajutat la construirea anumitor componente ale sistemului. Lucrul din cadrul acestui proiect este parțial finanțat de proiectul PNCDI II, SIR-RESDEC și de firma Siemens VDO Iași.

### Referințe bibliografice

- Bar-Haim, R., Dagan, I., Dolan, B., Ferro, L., Giampiccolo, D., Magnini B., Szpektor, I. (2006). The Second PASCAL Recognising Textual Entailment Challenge. *In Proceedings of the Second PASCAL Challenges Workshop on Recognizing Textual Entailment*. Venice. Italy.
- Iftene, A., Balahur-Dobrescu, A. (2007). Improving a QA System for Romanian Using Textual Entailment. *In Proceedings of RANLP workshop "A Common Natural Language Processing Paradigm For Balkan Languages"*. Pages 7-14, September 26, Borovets, Bulgaria.
- Iftene, A., Balahur-Dobrescu, A. (2008). Answer Validation on English and Romanian Languages. *In Proceedings of the CLEF 2008 Workshop*. 17-19 September. Aarhus, Denmark.
- Peñas, A., Rodrigo, Á., Sama, V., Verdejo, F. (2007) Overview of the Answer Validation Exercise 2006. CLEF 2006, *Lecture Notes in Computer Science LNCS 4730*. Springer, Berlin.
- Rodrigo, Á., Peñas, A., Verdejo, F. (2008). Overview of the Answer Validation Exercise 2008. *In working notes of CLEF2008*. September, Aarhus, Denmark.

# ECHIVALAREA ÎN LIMBA ROMÂNĂ A UNITĂȚILOR FRAZEOLOGICE INFINITIVEALE DIN LIMBA FRANCEZĂ

MARIA HUSARCIUC<sup>1,2,3</sup>

<sup>1</sup>*Facultatea de Litere, Universitatea "Al. I. Cuza", Iași*

<sup>2</sup>*Facultatea de Informatică, Universitatea "Al. I. Cuza", Iași*

<sup>3</sup>*Institutul de Informatică Teoretică Iași, Academia Română*

[mhusarciuc@gmail.com](mailto:mhusarciuc@gmail.com)

## Rezumat

Lucrarea de față prezintă o metodă de identificare a unităților frazeologice (expresii idiomatice, locuțiuni, structuri proverbiale etc.) într-un corpus paralel francez-român. Etapele urmărite sunt: adnotarea unităților frazeologice în textul în limba franceză și importul acestei adnotări pentru cele două variante de traducere în limba română. Se au în vedere criteriile specifice de identificare a acestor unități și a tipurilor lor specifice, precum și dificultățile care pot să apară în etapa de import, aplicate pe structuri cu verbul *faire* la infinitiv, în limba franceză, cu diferite posibilități de traducere în limba română.

## 1. Introducere

În realizarea resurselor lexical-semantice superioare dicționarelor informatizate, problemele cele mai mari apar în cazul mutațiilor semantice în contexte specifice și în cazul unităților frazeologice. Un imperativ, în acest caz, este descoperirea unui mod, flexibil și riguros în același timp, de a „manipula” îmbinările stabile de cuvinte. Astfel a luat naștere un subdomeniu al lingvisticii computaționale, cunoscut sub denumirea de *frazologie computațională* (engl. *computational phraseology*, v. Heid, 2005), disciplină complementară așa numitei *frazologii tradiționale*. Lucrarea de față prezintă o metodă de identificare a unităților frazeologice într-un corpus bilingv, francez-român, cu particularitatea că, pentru același text în limba franceză, sunt avute în vedere două versiuni de traducere în limba română. Identificarea acestor unități și stabilirea tipurilor în care ele se încadrează se realizează pornind de la criteriile precise.

Lucrarea conține, în a doua secțiune, o prezentare a premiselor teoretice, cu accent pe criteriile de identificare a unităților frazeologice, stabilind totodată o tipologie a acestora. A treia secțiune are în vedere analiza bazată pe corpus, cu accent pe structurile care conțin, în limba franceză, verbul *faire* la infinitiv și pe probleme care apar la importul adnotărilor în textele în limba română. Lucrarea se încheie cu prezentarea concluziilor și a perspectivelor viitoare de lucru.

## **2. Premise teoretice**

### **2.1. Criterii de identificare a unităților frazeologice**

Lăsând la o parte criteriul implicit, conform căruia orice unitate frazeologică trebuie să conțină minim două lexeme, reprezentative sunt frecvența, instituționalizarea, stabilitatea, caracterul idiomatic, variația și caracterul gradual (Corpas Pastor, 1996).

#### **2.1.1. Frecvența**

Principalul criteriu de identificare a unităților frazeologice este frecvența cu care ele apar, văzută atât ca frecvență de co-ocurență a elementelor componente într-un text sau corpus de texte, cât și ca frecvență de utilizare a structurii respective în vorbire. Se pleacă de la ideea conform căreia, cu cât este mai des folosită o structură sintactico-semantică într-o limbă dată, cu atât sunt mai mari șansele ca acea structură să fie sau să devină îmbinare stabilă de cuvinte sau chiar expresie idiomatică.

(Jean David, 1988) demonstrează, cu exemple, faptul că trecerea de la ne-idiomatic la idiomatic în limbă depinde de frecvența de utilizare a anumitor structuri. Orice îmbinare stabilă de cuvinte, care are o semnificație globală cunoscută, poate avea și o lectură „compozițională” și, reciproc, orice îmbinare liberă de cuvinte, a cărei semnificație de bază este semnificația elementelor componente, poate câștiga o semnificație globală (conotație) datorită unor elemente situate în realitatea extra-lingvistică (premise sociale, de exemplu).

#### **2.1.2. Instituționalizarea**

Prin utilizare frecventă, structurile neologice dobândesc un caracter oficial, ajungând să fie reproduse în vorbire fără modificarea formei. Fiind o etapă tranzitorie între intrarea în uz a unei structuri și stabilizarea ei morfo-sintactică și semantică, instituționalizarea constituie un criteriu destul de ambiguu în identificarea unităților frazeologice.

#### **2.1.3. Stabilitatea**

Proces complex, stabilitatea se realizează în două etape: la nivelul formei (fr. *figement*, sp. *fijacion*) și la nivelul conținutului (specializare semantică/ lexicalizare).

Stabilitatea este considerată aproape unanim principala trăsătură definitorie a unităților frazeologice. Totuși, contrar semnificației implicite a termenului, stabilitatea este o trăsătură greu de hotărât cu precizie, fiind mereu relativă la procesul continuu de construire și re-construire a faptelor de limbă (v. Eric Beaumatin, 1988). De aceea, între stabilitate și idiomaticitate limitele sunt destul de fragile.

#### **2.1.4. Idiomaticitatea**

În literatura de specialitate, stabilitatea și idiomaticitatea sunt de regulă prezentate împreună, ca două fațete ale aceluiași proces, definindu-se una prin cealaltă. Prin stabilizarea unei îmbinări libere de cuvinte, aceasta are șanse mari să se “idiomatizeze” și, reciproc, nici o expresie idiomatică nu poate exista dacă nu a fost în prealabil stabilizată morfo-sintactic.

(Harald Burger, 1988), analizând cele două concepte de bază ale cercetării frazeologice, pe care el le numește stabilitate (*Festigkeit*) și metaforicitate/ caracter metaforic (*Metaphorizität*), accentuează faptul că numai primul este un concept din sfera frazeologiei, al doilea fiind mai general. Există însă o inter-dependență între ceea ce exprimă aceste două concepte, fapt foarte bine ilustrat în lucrările lexicografice în care, adesea, unitățile frazeologice sunt marcate cu trăsătura „metaforic” (sau „figurat”), uneori specificându-se că structura respectivă este o expresie sau o locuțiune.

(Bernd Spillner, 1988), făcând distincția între *figement syntaxique* (termen asimilabil stabilității sintactice și care stă la baza formării cologațiilor) și *figement phraséologique* (termen sinonim cu idiomatizarea, dând naștere frazeologismelor), consideră caracterul idiomatic drept definitoriu pentru unitățile frazeologice, cologațiile fiind doar fenomene tranzitorii.

### 2.1.5. Variația

Variația frazeologică este o regulă lingvistică pe baza căreia se poate stabili gradul de regularitate al unui sistem frazeologic dat. Deși majoritatea autorilor vorbesc de existența a două tipuri de variație frazeologică (variantele și modificările creative), (Duneton & Claval, 1990) demonstrează, cu exemple pentru limba franceză, că așa-numitele variante frazeologice sunt în realitate expresii vecine, similare morfo-sintactic, și nu variante ale aceleiași expresii. În lucrarea citată sunt analizate structurile *être au bout de son rouleau* și *être au bout du rouleau*, care, deși foarte asemănătoare structural (singura diferență este prezența sau absența pronumelui reflexiv) au origini și sensuri diferite.

Spre deosebire de variante, modificările creative sunt contaminări conștiente între unități frazeologice din aceeași sferă semantică, de genul *Cine sapă groapa altuia departe ajunge*, dar nu și modificările inconștiente, datorate lipsei de cultură, care apar în structuri ca *a unsprezecea minune a lumii*.

### 2.1.6. Gradația

Gradația este o proprietate a însuși sistemului frazeologic, bazat pe îmbinări de cuvinte cu diferite grade de idiomatizare. Trecerea de la ne-idiomatic la idiomatic este posibilă numai datorită caracterului gradual al acestui sistem.

Acestor criterii le mai putem adăuga și capacitatea unităților frazeologice de a conserva termeni ieșiți din uz: *pe de rost, avoir la berlue*.

## 2.2. Tipuri de unități frazeologice

Unitățile frazeologice sunt enunțuri aparținând discursului repetat, „prefabricate de vorbire”, care reprezintă „tot ceea ce în vorbirea unei comunități se repetă într-o formă mai mult sau mai puțin identică de discurs deja făcut” (Coșeriu, 2000).

Dintre termenii folosiți pentru a desemna unitatea minimală a frazeologiei (izolare, frazeologism etc.), cel mai adecvat din punctul de vedere al consecvenței modului de definire este termenul *unitate frazeologică*. Poate fi unitate frazeologică orice îmbinare stabilă de cuvinte, fie expresie sau locuțiune expresivă, acestea la rândul lor fiind

subcategorizabile. (Dimitrescu, 1958) stabilește o serie de trăsături de diferențiere a locuțiunilor de expresii, recunoscându-le locuțiunilor proprietatea de a se comporta ca o singură parte de vorbire, funcția gramaticală unică și posibila coloratură expresivă, în timp ce expresiile sunt întotdeauna marcate stilistic, sunt variabile și fără o funcție gramaticală precisă. Criterii diferite de identificare a locuțiunilor găsim în (Branca-Rosoff, 1997): imposibilitatea realizării mutațiilor interne sau a substituțiilor, caracterul inseparabil al elementelor, „pierderea sentimentului de analicitate” datorită închegării. În plus, locuțiunile verbale oferă constrângeri asupra determinantilor și asupra posibilităților de pronominalizare.

Alți autori stabilesc subcategoriile de expresii. (Slave, 1966) clasifică expresiile în două tipuri: „consacrate și în uz propriu, având accepție tehnică sau fiind folosite ca îmbinări curente, banale” și „îmbinări folosite numai cu accepție figurată”. Exemplele pe care le alege autoarea (*a aduce la același numitor*, *a se spăla pe mâini* etc. pentru prima clasă și *a avea păr pe limbă*, *a scoate vorba cu cleștele* pentru a doua) corespund celor două clase de expresii denumite de (Dumistrăcel, 2001) „coppii ale realității”, respectiv „imaginare”. În prezent, se folosește foarte mult termenul *colocație* pentru a desemna îmbinări stabile de cuvinte care nu au în mod obligatoriu sens conotativ. (Todirașcu *et al.*, 2007)

În linii mari, termenii vehiculați reprezintă fie

a) tipuri complementare de structuri lexical-sintactice: expresii (având de regulă sens conotativ), locuțiuni (structuri morfo-sintactice complexe cu rol funcțional, uneori putând avea și conotații), colocații (îmbinări stabile de cuvinte care nu sunt nici expresii, nici locuțiuni, de genul formulelor de salut și a clișeeleor lingvistice), realizându-se astfel o clasificare cvasi-exhaustivă a tuturor îmbinărilor stabile de cuvinte,

fie

b) tipuri (cvasi-)concentrice de structuri lexical-semantice: unul dintre termeni reprezintă, în acest caz, o categorie supra-ordona(n)tă căreia i se subordonează celelalte.

### 3. Analiza bazată pe corpus

#### 3.1. Metoda de lucru

##### 3.1.1. Adnotarea unităților frazeologice în textul în limba franceză

Structuri sintactice similare pot reprezenta sau nu unități frazeologice. Având în vedere dificultatea diferențierii îmbinărilor stabile de cuvinte de îmbinările libere, metoda pe care am ales-o este adnotarea manuală a acestor structuri în textul original (în limba franceză) și importarea adnotărilor în traducerea românească aliniată cu originalul. Textul folosit este *Madame Bovary* de Gustave Flaubert, cu două versiuni de traducere în limba română: a lui Ludovic Dauș (Flaubert, 1915) și a lui Demostene Botez (Flaubert, 1968). Adnotarea unităților frazeologice se realizează pe textul francez, tokenizat în prealabil, folosind adnotatorul PALinkA.

## ECHIVALAREA ÎN LIMBA ROMÂNĂ A UNITĂȚILOR FRAZELOGICE INFINITIVALE DIN LIMBA FRANCEZĂ

Următorul exemplu prezintă o secvență din fișierul XML de ieșire, ce conține adnotarea unei unități frazeologice din textul în limba franceză.

```
<MWE
  DEF="v.TLFI: FAIRE-VALOIR, subst. masc. invar. Mode de
gestion, d'exploitation d'un capital immobilier."
  HEADID="faire"
  ID="0"
  OBS="Structură folosită aici cu sensul de bază"
  OTHER_TYPE=""
  TYPE="IDIOM">
<W id="1298">faire</W>
<W id="1299">valoir</W>
</MWE>
```

Dintre elementele avute în vedere în adnotarea manuală, important este tipul (TYPE), care trebuie ales dintr-o listă reprezentând IDIOM (expresie idiomatică), EXPRESSION (expresie, structură cu sens conotativ) COLLOCATION (colocație sau locuțiune), PROVERB (proverb) sau OTHER (în cazul în care o structură nu corespunde niciuneia din clasele de mai sus, caz în care completarea unei valori pentru OTHER\_TYPE este foarte importantă. Definiția (DEF) va fi preluată din *Le Trésor de la Langue Française Informatisé*. ID-ul unității frazeologice adnotate este incrementat automat, iar HEADID-ul este reprezentat de nucleul sintactico-semantic al respectivei structuri. În OBS se pot nota, în timpul adnotării, observații suplimentare ce pot fi utile la o prelucrare ulterioară.

### ***3.1.2. Importarea adnotărilor în traduceri în limba română, aliniate cu textul original***

Etapa următoare adnotării manuale a unităților frazeologice în textul francez este importarea lor în traduceri în limba română. Acest lucru se poate realiza numai după alinierea, la nivel de cuvânt, a textelor. Premisa de la care se pleacă este că importarea unei structuri se poate realiza chiar și în cazul lipsei unei corespondențe literale a expresiilor în cele două limbi, deoarece elementele care preced și succed expresia sunt, de regulă, aliniate corespunzător. Este important să se țină cont de contextul de apariție, o aliniere a elementelor din interiorul unor unități frazeologice echivalente fiind greu de realizat, datorită unor neconcordanțe de structură (a se vedea, pentru detalii, subsecțiunea 3.3).

Realizarea unor șabloane cu ajutorul cărora să se precizeze structurile morfo-sintactice vizate (de exemplu, precizarea faptului că unei structuri infinitivale din franceză îi poate corespunde, în limba română, o structură cu verbul la conjunctiv) va crește acuratețea identificării lor în traduceri românești.

### ***3.2. Tipuri de structuri ce-l conțin pe faire în limba franceză***

Verbul *faire*, relevant pentru polisemantismul său, apare frecvent în *Madame Bovary* la modul infinitiv. Structurile care îl conțin sunt dintre cele mai variate. Exemplele extrase din corpusul adnotat au fost foarte utile în marcarea distincției între diferite tipuri de unități frazeologice, ținând cont de tipologia lor și de criteriile de identificare prezentate în secțiunea a doua.

### 3.2.1. Expresii idiomatice

O expresie idiomatice este caracterizată prin frecvență redusă de apariție în același text, stabilitate morfo-sintactică, fiind intrată în uz, deci recunoscută de vorbitorii nativi; este o structură cu sens figurat, metaforic, dezvoltat pe baza unui sens primar. O expresie idiomatice poate fi specializată pentru un anumit domeniu sau poate să apară numai în contexte specifice. Adesea, structurile din această categorie sunt, în corpusul ales, fie traduse prin parafraze (a), fie absente (b).

*L'officier de santé, chemin faisant, compris aux discours de son guide que M. Rouault devait être un cultivateur des plus aisés. Il s'était cassé la jambe, la veille au soir, en revenant de **faire les Rois**, chez un voisin.*

(a) *Pe drum, din tot ce-i spunea călăuza, ofițerul sanitar își dădu seama că domnul Rouault trebuie să fie un gospodar din cei mai înstăriți. Își rupsesse piciorul, în ajun, seara, întorcându-se de la un vecin unde **serbase Boboteaza**. (Botez)*

(b) *Ofițerul de sănătate, pe drum, înțelese din vorbele călăuzei sale că d. Rouault trebuia să fie un cultivador dintre cei mai bogați. Își frânsese piciorul, seara din ajun, întorcându-se de la un vecin. (Dauș)*

### 3.2.2. Expresii

Structuri intermediare între colocații și expresii idiomatice, expresiile sunt frecvente ca mod de utilizare, stabile din punct de vedere morfo-sintactic și au sens conotativ. În general, echivalenții de traducere pentru aceste structuri sunt fie calcuri lexical-semantice după limba franceză, fie expresii deja existente și frecvente.

*A l'encontre des tendances maternelles, il avait en tête un certain idéal viril de l'enfance, d'après lequel il tâchait de former son fils, voulant qu'on l'élevât durement, à la spartiate, pour **lui faire une bonne constitution**.*

*Potrivnic înclinațiilor materne, avea în cap un anume ideal viril despre copilărie, după care încerca să-și educe feciorul, voind să fie crescut cu asprime, după moda spartană, **ca să ajungă voinic**. (Botez)*

*În contra tendințelor materne, avea în cap un oare-care ideal viril al copilăriei, după care căuta să-și formeze feciorul, voind să-l crească cu asprime după moda spartacă **ca să-l facă tare de constituție**. (Dauș)*

### 3.2.3. Colocații

Structuri morfo-sintactice complexe cu rol funcțional, colocațiile pot prezenta uneori și sensuri conotative. Diferențele între colocații și locuțiuni nu sunt foarte clar formalizabile, de aceea am păstrat o singură categorie. Cele mai frecvente structuri din această categorie care-l conțin în limba franceză pe *faire* la infinitiv sunt sintagmele cu determinări nominale: *faire sa toilette*, *faire du punch*, *faire de la tapisserie*, care sunt traduse în general literal în limba română.



### 3.2.4. *Predicate verbale compuse*

Predicatele verbale compuse din limba franceză, în care verbul *faire* la infinitiv are valoare de auxiliar și este urmat de un alt infinitiv, sunt de regulă traduse în limba română prin predicate verbale simple, fiind diferite de *predicatele complexe*, care reprezintă structuri compuse din verb la infinitiv și grup substantival sau prepozițional, de genul *a aduce atingere, a intra în vigoare* etc. (Todirașcu *et al.*, 2007).

Un exemplu de predicat verbal compus din limba franceză este *faire étudier la médecine*, structură ale cărei traduceri în limba română vor fi analizate în secțiunea 3.3.1.

### 3.3. *Probleme la alinierea traducerii în limba română a structurilor cu faire*

**3.3.1.** Adesea pot să apară probleme la alinierea celor două versiuni de traducere în limba română, mai ales datorită faptului că Ludovic Dauș recurge mult mai des la traduceri literale ale expresiilor decât Demostene Botez. Un exemplu ar fi structura infinitivală „*pour lui faire une bonne constitution*”, care în traducerea lui Dauș apare sub forma „*ca să-l facă tare de constituție*”, în timp ce D. Botez o traduce prin: „*ca să ajungă voinic*”.

Ludovic Dauș păstrează de regulă structura sintactică a frazei din originalul franțuzesc și traduce adesea literal îmbinările de cuvinte (de exemplu, traduce *lui faire étudier la médecine* prin *să-l facă să învețe medicina*). Demostene Botez simplifică, în unele cazuri, structura frazei, recurge la inversiuni sau parafraze (traduce, de exemplu, aceeași structură prin *să-l înscrie la Medicină*).

**3.3.2.** Problematice sunt și locuțiunile cu sens general, particularizat în text datorită existenței unei referințe anaforice concrete. Prin traducere, se poate păstra referința anaforică la o structură exprimată anterior (cum procedează L. Dauș în fraza (2), traducere fidelă a frazei din franceză) sau se poate folosi un echivalent particular, cu explicitarea sensului (ca în (3), unde D. Botez traduce locuțiunea *faire la demande* prin *s-o ceară în căsătorie*, deși în paragraful anterior este prezentată intenția lui Charles de a se căsători cu Emma).

(1) „*Charles se promet de faire la demande quand l'occasion s'en offrira*”.

(2) „*Carol își făgăduie să facă cererea când ocazia se va prezenta*”.

(3) „*Charles se hotărăie s-o ceară în căsătorie de îndată ce va avea ocazie*”.

Aceste diferențe nu impun dificultăți în lectura și înțelegerea textelor, dar în automatizarea procesului de extragere a diferitelor tipuri de unități frazeologice, în cazuri similare celui din (3) nu se va putea găsi echivalentul locuțiunii din (1) decât în situația în care corpusul va fi procesat suplimentar cu un sistem de rezoluție a anaforelor.

**3.3.3.** Contextul este absolut necesar în unele cazuri pentru a discerne între sensul propriu al unei anumite unități frazeologice și un sens particular. Locuțiunea verbală *faire valoir* are în limba franceză atât un sens propriu, general (ca în structura: *faire valoir une excuse*), cât și sensul particular: „*Loc. Faire valoir ses droits. Demander (pour soi) à une administration, à une hiérarchie, l'application d'une norme, d'un règlement*” (TLFI). Acest sens particular se regăsește în contextul „*se retira dans la campagne, où il voulut «faire valoir»*”, tradus de L. Dauș: „*se retrase la țară unde căută să exploateze*” și în mod similar, dar mai explicit, de D. Botez: „*se retrase la țară, unde voi să exploateze singur o moșie*”.

**3.3.4.** Există cazuri în care anumite expresii fără echivalent idiomatice în limba română sunt pur și simplu omise la traducere. În exemplul următor, expresia *faire les Rois*, are, în TLFI, următoarea definiție: „*Faire les Rois (vieilli), tirer les Rois. Se réunir pour une fête qui consiste à partager la galette des Rois contenant la fève qui rend roi ou reine celui ou celle qui la trouve. [M. Rouault] s'était cassé la jambe, la veille au soir, en revenant de faire les Rois, chez un voisin (FLAUB., M<sup>me</sup> Bovary, t. 1, 1857, p. 13). Au jour de l'Épiphanie, M<sup>me</sup> Bavretel conviait les amis d'Armand à venir « tirer les rois » (GIDE, Si le grain, 1924, p. 474).*”

Expresia idiomatice, întâlnită în *Madame Bovary* în fraza citată în definiția din TLFI, nu are echivalent în traducerea lui Dauș („*Își frânsese piciorul, seara din ajun, întorcându-se de la un vecin.*”), în timp ce Demostene Botez îi găsește un echivalent ne-idiomatice, un fel de definiție prescurtată: „*Își rupsesse piciorul, în ajun, seara, întorcându-se de la un vecin unde serbase Boboteaza.*”

**3.3.5.** Expresiilor onomatopice le sunt uneori găsite echivalente de traducere ce nu au în componența lor interjecții, ca în propoziția „*Ce n'est pas la peine de faire tant de fla-fla*”, tradusă prin „*Pentru atâta lucru nu face să-ți iei aere*” (L. Dauș) sau prin „*Nu-i cazul să-ți dai atâtea ifose*” (D. Botez). Importarea automată a adnotării în situații de acest gen este posibilă numai dacă în faza de verificare a alinierii textelor se recurge la un artificiu, aliniindu-se fie interjecția din original cu substantivul din traducere, fie direct întreaga expresie cu echivalentul ei în cele două versiuni românești.

#### 4. Concluzii

Încercând să prezinte o metodă de identificare a unităților frazeologice într-un corpus beletristic bilingv, cu marcarea structurilor echivalente în franceză și română, această lucrare a avut în vedere și o încercare de tipologizare a acestor tipuri de unități, precum și a dificultăților întâmpinate.

În perspectivă, se urmărește definitivarea adnotării și găsirea de soluții optime pentru realizarea unei resurse bilingve a unităților frazeologice, cu informații cât mai complexe.

**Mulțumiri** Cercetarea întreprinsă beneficiază de finanțare CNCSIS (Grant TD, cod 492) și nu s-ar putea realiza fără sprijinul continuu al îndrumătorilor de doctorat, Prof. Dr. Eugen Munteanu și Prof. Dr. Dan Cristea, cărora țin să le mulțumesc pentru

viziunea critică oferită. Mulțumesc de asemenea membrilor grupului de lingvistică computațională din Iași, care m-au ajutat în diferite etape ale cercetării mele de până acum.

### Referințe bibliografice

- Beaumat, Eric (2000). *Langue/ discours/ texte à l'épreuve des faits de figement*. Greciano, Gertrud (Ed.), *Micro- et macroléxemes et leur figement discursif*, Actes du colloque international CNRS URA 1035 Langue-Discours-Cognition, 6-7-8 décembre 1998, Saverne, Editions Peeters, Louvain / Paris, p. 3-12.
- Branca-Rosoff, Sonia (1997). Modèles de locutionarité et effets de figement dans le discours politique de l'an II. „*La locution: entre lexique, syntaxe et pragmatique. Identification en corpus, traitement, apprentissage*”, Textes réunis par Pierre Fiala, Pierre Lafon, Marie-France Piguet, Editions Klincksieck, Paris, 1997, p. 285-286.
- Burger, Harald (1988). „BILDHAFT, ÜBERTRAGEN, METAPHORISCH...“. Zur Konfusion um die semantischen Merkmale von Phraseologismen. „*EUROPHRAS 88. Phraseologie Contrastive*“, p. 17-29.
- Corpas Pastor, Gloria (1996). *Manual de fraseología española*. Editorial Gredos, Madrid.
- Coșeriu, Eugeniu (2000). *Lecții de lingvistică generală*. Editura Arc, Chișinău.
- David, Jean (1988). Tous les predicats ne meurent pas idiomes. Mais nul n'est à l'abri. „*EUROPHRAS 88. Phraseologie Contrastive*”, Actes du Colloque International Klingenthal – Strasbourg, 12-16 mai 1988, Collection Recherches Germanique N.2, Strasbourg, p. 75-82.
- Dimitrescu, Florica (1958). *Locuțiunile verbale în limba română*. Editura Academiei, București.
- Dumistrăcel, Stelian (2001). *Până-n pânzele albe. Expresii românești*. Biografii-motivații. Institutul European, Iași.
- Duneton, Claude; Claval, Sylvie (1990). *Le Bouquet des expressions imagees. Encyclopedie thematique des locutions figurees de la langue francaise*. Editions du Seuil, Paris.
- Flaubert, Gustave (1968). *Doamna Bovary*. În românește de Demostene Botez. Prefață de Aurelian Tănase. Editura pentru Literatură Universală, București.
- Flaubert, Gustave (1915). *Doamna Bovary*. Traducere de Ludovic Dauș. Ediția a II-a, revăzută. Editura Minerva, București.
- Heid, Ulrich (2005). Computational Phraseology. Approaches to the computational analysis and representation of phraseological units and to their extraction from text corpora. „*PHRASEOLOGIE 2005 – La phraséologie dans tous ses états (Colloque interdisciplinaire)*”, Louvain-la-Neuve, p.13-15.
- Slave, Elena (1966). *Structura sintagmatică a expresiilor figurate*, LL, XI, p. 397-413.

Spillner, Bernd (2000). Phraséologie et textologie comparées français – allemand. Greciano, Gertrud (Ed.), *Micro- et macrolexemes et leur figement discursif*, p. 23-32.

Todirașcu, Amalia, Dan Ștefănescu, Christopher Gledhill (2007). Un sistem de extragere a colocațiilor. *Lucrările atelierului Resurse lingvistice și instrumente pentru prelucrarea limbii române*. Ed. Universității Al. I. Cuza, Iași, p. 119-129.

PALinkA, <http://clg.wlv.ac.uk/projects/PALinkA/>

*Le Trésor de la Langue Française Informatisé* <http://atilf.atilf.fr/tlf.htm>

# COLECTAREA ȘI PROCESAREA DOCUMENTELOR ROMÂNEȘTI ALE CORPUSULUI JRC-ACQUIS

ALEXANDRU CEAUȘU

<sup>1</sup>*Institutul de Cercetări pentru Inteligența Artificială, Academia Română*

[aceausu@racai.ro](mailto:aceausu@racai.ro)

## Rezumat

Partea românească a corpusului JRC-Acquis continuă să crească odată cu creșterea numărului de documente ale “Acquis Communautaire” traduse în română. De asemenea, calitatea corpusului se îmbunătățește deoarece multe din documentele deja traduse trec prin câteva faze de corectare. Lucrarea prezintă o nouă versiune românească a corpusului JRC-Acquis-Ro conținând peste 30 de milioane de cuvinte în 19211 documente. Numărul de documente românești prezente în noua versiune este de trei ori mai mare decât cel din versiunea precedentă.

## 1. Introducere

Necesarul de corpuri paralele pentru aplicațiile de procesare a limbajului natural a cunoscut un trend ascendent accentuat pe parcursul ultimilor ani. Corpusurile paralele sunt folosite în aplicațiile de traducere automată sau categorizare multilinguală; pentru a produce resurse lexicale sau semantice multilinguale, cum sunt dicționarele sau ontologiile; pentru a testa consistența procesului de traducere, etc. Cele mai multe corpuri paralele existente conțin limbi de largă circulație și au un număr redus de perechi de limbi. Dintre acestea, cel mai cunoscut este corpusul francez-englez Hansards (German, 2001). Corpusurile paralele ce conțin mai multe perechi de limbi sunt de mici dimensiuni sau pentru texte foarte specializate cum ar fi biblia (Resnik et al. 1999) sau romanul 1984 al lui George Orwell (Erjavec 2004). Unul dintre cele mai importante corpuri multilinguale este EuroParl (Koehn, 2005) disponibil pentru 11 din limbile comunității europene.

Pentru aplicațiile de procesare a limbajului natural din limba română cel mai important corpus este JRC-Acquis (Steinberger et al., 2006). Acesta este în prezent cel mai mare corpus multilingual disponibil, conținând 22 de limbi. Corpusul este disponibil în format XML conform specificațiilor TEI (Text Encoding Initiative). De asemenea, conține și alinierea celor mai mult de 230 de perechi de limbi conținute de JRC-Acquis. Corpusul crește pe măsura traducerii legislației europene și în limbile noilor candidați.

„UE Acquis Communautaire” este termenul prin care se face referire la corpul comun de legi și obligativități care leagă toate statele membre ale Comunității Europene. Acquis-ul conține principii și obiective politice ale diverselor tratate semnate în cadrul Uniunii Europene (UE), legislație UE, declarații și rezoluții, acorduri internaționale și obiective comune. Toate țările acceptate în Uniunea Europeană trebuie să ratifice „Acquis Communautaire”. Pe lângă cele 22 de limbi ale Comunității Europene, Acquis este tradus și în limbile croată și turcă. Datorită efortului depus la ICIA (Institutul de Cercetări pentru Inteligența Artificială, Academia Română) pentru colectarea și

adnotarea documentelor Acquis-ului românesc, limba română a fost prezentă în pachetul de distribuție JRC-Acquis încă de la prima versiune a acestuia.

Numărul de documente comune în perechea de limbi Engleză-Română este de 11 469. Documentele comune constituie un important corpus paralel conținând 59 986 838 de cuvinte.

**Tabelul 1:** JRC-Acquis versiunea 3.0 și noua versiune de corpus românesc

<b>Limbă</b>	<b>Documente</b>	<b>Caractere</b>	<b>Cuvinte</b>
<b>bulgară</b>	11384	104522671	30146967
<b>cehă</b>	21438	148972981	46832312
<b>daneză</b>	23624	213468135	50944626
<b>germană</b>	23541	232748675	50929652
<b>greacă</b>	23184	239583543	55887003
<b>engleză</b>	23545	210692059	55537910
<b>spaniolă</b>	23573	238016756	62132608
<b>estoniană</b>	23541	192700704	40953424
<b>finlandeză</b>	23284	212178964	40107981
<b>franceză</b>	23627	234758290	62100432
<b>ungară</b>	22801	213804614	46188364
<b>italiană</b>	23472	230677013	57217002
<b>lituaniană</b>	23379	199438258	44392842
<b>letonă</b>	22906	196452051	44703607
<b>malteză</b>	10545	128906748	37883562
<b>olandeză</b>	23564	231963539	56771856
<b>poloneză</b>	23478	214464026	49253537
<b>portugheză</b>	23505	227499418	59606203
<b>română</b>	<b>19211</b>	<b>182631277</b>	<b>30832212</b>
<b>slovacă</b>	21943	179920434	46211035
<b>slovenă</b>	20642	178651767	47643215
<b>suedeză</b>	20243	199004401	46974192
<b>Total</b>	<b>476430</b>	<b>4288962348</b>	<b>1053305415</b>

Pentru indexarea lor, documentele din JRC-Acquis au fost clasificate manual cu ajutorul unui sistem de clasificare (EUROVOC) conținând peste 6 000 de descriptori organizați ierarhic. Versiunea 4.2 a EUROVOC este disponibilă în 21 de limbi ale țărilor din UE printre care și limba română. În cadrul Eurovoc-ului, termenii se împart în două categorii: descriptori și non-descriptori. Descriptorii sunt cuvinte sau expresii care denotă concepte din domeniile tezaurului într-un mod ne-ambiguu, pe când non-descriptorii sunt cuvinte sau expresii reprezentate deja în tezaur de un descriptor echivalent. Versiunea în limba engleză conține 6 645 de descriptori iar, în comparație, tezaurul în limba română conține doar 4 625 (aproximativ 70% din cel englezesc). Acești descriptori sunt organizați în 21 de domenii (de la politică și relații internaționale până la mediu, industrie sau geografie) ce conțin la rândul lor micro-tezaure. Există un total de 519 micro-tezaure (în română doar 508), fiecare din aceștia constituind un arbore în nodurile căruia se găsesc descriptori. Domeniile și micro-tezaurele au

identificatori unici, independenți de limbă, asigurându-se astfel o inter-relaționare multilingvă.

## **2. Colectarea și convertirea documentelor JRC-Acquis-Ro**

Fișierele românești și bulgărești nu sunt disponibile în același format (HTML) ca și documentele celorlalte limbi din Acquis, neputând fi procesate de aceleași instrumente de convertire HTML-TEI. Fișierele în format HTML disponibile pentru celelalte limbi ale JRC-Acquis conțin și informații cu privire la structura documentului, cum ar fi secțiunile de anexe și semnături, secțiunile cu textul și titlul documentului etc. Această structură nu se regăsește în formatul Microsoft Word, format în care sunt disponibile documentele românești ale JRC-Acquis.

Pentru a constitui colecția de documente în limba română, fișierele au fost descărcate de pe situl „CCVista Translation Database” folosind drept adresă „<http://ccvista.taix.be/Fulcrum/CCVista/RO/<celex>>” unde <celex> este numărul unic de identificare al documentului. Numărul total de documente în limba română disponibile pe situl CCVista este de 19 286.

Fișierele au fost convertite din formatul Microsoft Word în formatul XML conform specificațiilor TEI. Conversia celor 19 286 fișiere a fost făcută automat fiind folosit pachetul de funcții „Visual Studio Tools for Office”. Aceste funcții permit interacțiunea directă cu aplicația Microsoft Office direct din mediul de programare. Datorită particularităților formatului, conversia documentelor a implicat și o serie de etape intermediare:

- au fost înlăturate comentariile traducătorilor;
- au fost șterse notele de subsol și secțiunile de cap de pagină;
- a fost normalizată folosirea caracterelor diacritice (unele documente foloseau „ș” și „ț” cu cedil iar altele foloseau „ş” și „ţ” cu virgulă).

Dintre cele 19 286 de fișiere în format Microsoft Word au fost convertite 19 211 (restul de documente având erori de format).

În formatul TEI-XML al documentelor românești au fost adăugate și datele de indexare EUROVOC acolo unde acestea erau disponibile.

## **3. Comparație între corpusul JRC-Acquis-Ro și alte corpusuri românești**

Pentru a compara diferența de vocabular între domeniul legislativ și alte domenii am folosit corpusul românesc Agenda (colecție de articole din săptămânalul timișorean *Agenda*). Acest corpus conține 8 408 185 de cuvinte. După cum se observă din tabelul de mai jos, datorită domeniilor diferite abordate în cele două corpusuri, listele primelor cuvinte-conținut, ordonate după rangul de frecvență, diferă într-o proporție considerabilă.

Tabelul 2: Primele 25 de cuvinte, sortate în funcție de ocurență, din corpusurile Jrc-Acquis-Ro și Agenda

Jrc-Acquis Ro		Agenda	
art.	Yn	Timișoara	Np
articolul	Ncmsry	este	Vmip3s
nr.	Yn	ora	Ncfsry
membre	Ncftp-n	Timiș	Np
regulamentul	Ncmsry	ani	Ncmp-n
alin.	Yn	pot	Vmip3p
privind	Vmg	România	Np
trebuie	Vmip3s	sunt	Vmip3p
statele	Ncfpry	zona	Ncfsry
este	Vmip3s	mare	Afpfsm
regulament	Ncms-n	vând	Vmip3p
vedere	Ncfsrn	privind	Vmg
având	Vmg	data	Ncfsry
CEE	Np	fost	Vmp--sm
comisiei	Ncfsoy	perioada	Ncfsry
consiliului	Ncmsoy	an	Ncms-n
prezentul	Afpmsry	apartament	Ncms-n
directiva	Ncfsry	piața	Ncfsry
ce	Np	persoane	Ncftp-n
comisia	Ncfsry	anul	Ncmsry
prezenta	Afpfsm	poate	Vmip3s
în special	Rgp	muncă	Ncfsrn
comunității	Ncfsoy	astfel	Rgp
prevăzute	Vmp--pf	are	Vmip3s
membru	Ncms-n	București	Np

Aplicațiile de procesarea limbajului natural care vor avea la bază corpusul JRC-Acquis-Ro trebuie să ia în considerare și zgomotul pe care un astfel de corpus îl conține. Experimentele noastre ne-au arătat că pentru un lexicon de peste 1 200 000 de cuvinte (incluzând aici și entități denumite), în corpusul JRC-Acquis-Ro încă mai găsim foarte multe cuvinte necunoscute - din 30 832 212, 2 796 473 sunt cuvinte necunoscute.

Tabelul 3: Primele 20 de cuvinte necunoscute din JRC-Acquis-Ro

CEE	126951	NC	9872
ex	30926	and	9522
Amtsgericht	25646	BCE	8510
JO	25632	THE	8223
see	22721	the	8201
year	17553	en	4965
please	13805	comarca	4910
pct.	12978	del	4464
EUR	12870	Euratom	4455
from	10951	CECO	3875

Cuvintele necunoscute sunt, în marea lor majoritate, cuvinte aparținând altor limbi – multe din pasajele din JRC-Acquis-Ro sunt copii din originalul limbii din care au fost



COLECTAREA ȘI PROCESAREA DOCUMENTELOR ROMÂNEȘTI  
ALE CORPUSULUI JRC-ACQUIS

traduse. În tabelul 3 sunt prezentate primele 20 de cuvinte necunoscute ordonate după rangul de frecvență.

Pentru a observa diferențele dintre JRC-Acquis-Ro și alte corpusuri care au la bază „Acquis Communautaire” am testat modelele de limbă construite folosind SRILM (Stolcke, 2002) din JRC-Acquis-Ro, DGT-TM (Directorate-General for Translation – Translation Memory – <http://langtech.jrc.it/DGT-TM.html>) și SEEERANET (Tufiș et. al., 2008). Modelele de limbă sunt de ordinul 3 și au fost antrenate folosind forma de ocurență a cuvintelor. Textul pe care a fost evaluată perplexitatea este unul din fișierele JRC-Acquis-Ro (4 271 de cuvinte). În tabelul 4 se poate observa că deși corpusurile diferă mult considerând modul în care au fost colectate, perplexitatea raportată pe fișierul de evaluare se îmbunătățește în funcție de mărimea acestora. Numărul de cuvinte necunoscute este un indicator al uni-gramelor neîntâlnite în corpusul de antrenament.

Tabelul 4: Evaluarea modelelor de limbă ale corpusurilor JRC-Acquis-Ro, DGT-TM Ro și SEEERANET Ro

	<b>JRC-Acquis Ro</b>	<b>DGT-TM Ro</b>	<b>SEEERANET Ro</b>
Număr de cuvinte în corpus	30 832 212	2 528 584	1 442 915
Cuvinte necunoscute	79	218	362
Perplexitate	141	174	262

#### 4. Concluzii

Corpusul JRC-Acquis-Ro, cu peste 30 de milioane de cuvinte, este cel mai mare corpus monolingual și multilingual disponibil pentru limba română. Încă de la lansarea sa, în 2005, odată cu prima versiune a JRC-Acquis, corpusul JRC-Acquis-Ro este folosit în experimente de traducere automată și pentru dezvoltarea sistemelor de întrebare-răspuns multilinguale, devenind un corpus de referință pentru limba română.

Deși mărimea sa îl recomandă pentru orice aplicație de procesare a limbajului natural pentru limba română, corpusul necesită o etapă de filtrare pentru a fi înlăturate paragrafele din altă limbă decât româna. De asemenea, trebuie avută în vedere limitarea strictă la discursul juridic, limitare care nu permite o generalizare a fenomenelor de limbă observate în corpus.

#### Referințe bibliografice

- Steinberger Ralf, Bruno Pouliquen, Anna Widiger, Camelia Ignat, Tomaz Erjavec, Dan Tufiș, Dániel Varga (2006). *The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages*. Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC'2006). Genoa, Italy, 24-26 May 2006
- Germann Ulrich (ed.) (2001). *Aligned Hansards of the 36th Parliament of Canada - Release 2001-1a*. <http://www.isi.edu/natural-language/download/hansard/>
- Koehn Philipp (2005). *EuroParl: A Parallel Corpus for Statistical Machine Translation*. Machine Translation Summit 2005. Phuket, Thailand. <http://people.csail.mit.edu/koehn/publications/europarl/>

- Resnik Philip, Mari Broman Olsen & Mona Diab (1999). *The Bible as a Parallel Corpus: Annotating the ,book of 2000 Tongues'*. Computers and the Humanities, 33(1-2), pp. 129-153. <http://www.umiacs.umd.edu/users/resnik/>
- Andreas Stolcke. (2002). "SRILM - An Extensible Language Modeling Toolkit", in *Proc. Intl. Conf. Spoken Language Processing*, Denver, Colorado, September 2002
- Erjavec Tomaž (2004). MULTEXT-East Version 3: Multilingual Morphosyntactic Specifications, Lexicons and Corpora. *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'2004)*, pp. 1535-1538, Paris. <http://nl.ijs.si/ME/CD/docs/1984.html>.
- Dan Tufiș, Svetla Koeva, Tomaž Erjavec, Maria Gavrilidou, Cvetana Krstev. (2008). Building language resources and translation models for machine translation focused on South Slavic and Balkan languages, *FASSBL 2008: The Sixth International Conference Formal Approaches to South Slavic and Balkan Languages*, Dubrovnik, Septembrie 25-28, 2008

# EXPERIMENTE DE TRADUCERE AUTOMATĂ BAZATĂ PE EXEMPLE PENTRU LIMBILE ENGLEZĂ/ROMÂNĂ

IRIMIA ELENA

*Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București*

[elena@racai.ro](mailto:elena@racai.ro)

## Rezumat

Acest articol descrie arhitectura unui sistem de traducere automată bazată pe exemple care a fost implementat în procesul de cercetare doctorală a autorului. Aplicația nu este neapărat inovativă în cadrul mai larg al paradigmei EBMT (Example Based Machine Translation), ci reprezintă un experiment care a dorit să combine în mod eficient o parte dintre resursele și instrumentele pentru prelucrarea automată a limbajului natural dezvoltate la ICIA cu tehnici și algoritmi consacrați în domeniul traducerii automate.

### 1. Introducere

Orice aplicație de traducere automată este o întreprindere care necesită investirea a importante cantități de timp și energie. De aceea este util ca o astfel de aplicație să se construiască pe fundația muncii unei echipe de lucru (cum este grupul de PLN de la ICIA), care să poată pune la dispoziție instrumentele și resursele indispensabile. Articolul descrie în detaliu aplicația până la stagiul actual de implementare și prezintă pe scurt obiectivele încă neatinse. Din lipsă de spațiu, numărul exemplurilor este foarte redus iar rezultatele sunt prezentate doar cantitativ și evaluate sumar.

### 2. Resurse lingvistice utilizate și aplicații de preprocesare ale acestora

Ca resursă fundamentală pentru aplicația de traducere bazată pe exemple pe care am implementat-o am ales corpusul paralel multilingv **JRC-Acquis** (Steinberger et al., 2006). Am considerat foarte potrivit faptul că acest corpus este: **omogen** - dedicat unui domeniu specific (cu un conținut de natură juridică); **consistent** - cel puțin în teorie, orice expresie juridică din corpus trebuie să fie tradusă întotdeauna în același fel, într-o manieră validată de comunitatea profesioniștilor în domeniu; **actual**: JRC-Acquis este o colecție dinamică de documente juridice extrase din *Acquis Communautaire* (AC), care reprezintă corpul total de legi ale Uniunii Europene aplicabile în toate țările membre UE. AC, și implicit JRQ-Acquis, se îmbogățește constant cu noi documente, pe măsură ce Uniunea Europeană se extinde și țările membre își aliniază legislația la cea comunitară. JRC-Acquis este disponibil în 22 dintre cele 23 de limbi oficiale ale Uniunii Europene (traducerile irlandeze nu sunt încă disponibile) și reprezintă cel mai mare corpus paralel existent în acest moment, atât ca dimensiune cât și ca număr de limbi implicate. În forma utilizată de aplicația pe care am construit-o, corpusul conține doar perechea de limbi română-engleză și este rezultatul unor acțiuni consecutive de procesare: segmentare și aliniere la nivel de propoziție, segmentare la nivel de cuvânt, analiză morfo-sintactică, lematizare, adnotare sintactică de suprafață (chunking), aliniere lexicală și analiză a dependențelor sintactice între cuvinte. Documentele sunt

codificate XML conform DTD-ului (atributele acestuia captează toate adnotările produse de aplicațiile de pre-procesare, cu excepția alinierilor lexicale care sunt furnizate într-un fișier separat):

```
<!DOCTYPE text [
  <!ELEMENT text (body)>
  <!ATTLIST text id CDATA #REQUIRED>
  <!ELEMENT body (tu+)>
  <!ELEMENT tu (seg+)>
  <!ATTLIST tu id CDATA #REQUIRED>
  <!ELEMENT seg (s)>
  <!ATTLIST seg lang (en | ro) #REQUIRED>
  <!ELEMENT s (w | c)+>
  <!ATTLIST s id ID #REQUIRED>
  <!ELEMENT c (#PCDATA)>
  <!ELEMENT w (#PCDATA)>
  <!ATTLIST w
    ana CDATA #REQUIRED
    lemma CDATA #REQUIRED
    chunk CDATA #IMPLIED
    wns CDATA #IMPLIED
    head CDATA #IMPLIED
  >
]>
```

Strategia de extragere de exemple de traducere din corpusul JRC-Acquis pe care am implementat-o se bazează pe existența unor alinieri la nivel de cuvânt între 2 propoziții pereche, precum și pe adnotarea acestora cu dependențe sintactice; aceste proceduri sunt asigurate de aplicațiile următoare (implementate în cadrul ICIA):

- **YAWA**. Presupunând că avem de-a face cu o propoziție  $p_1$  într-o limbă  $l_1$  și traducerea ei  $p_2$  în limba  $l_2$ , o aliniere lexicală presupune stabilirea de corespondențe între cuvintele din  $p_1$  și cele din  $p_2$  astfel încât acestea să reprezinte traduceri reciproce. Într-o formă ușor prelucrabilă de către alte aplicații, structura de alinieri din Figura 1 este reprezentată de către YAWA prin lista:  $\{(1,1), (2,2), (3,3), (4,4), (5,5), (6,6), (7,7), (8,8), (9,8), (10,9), (11,10), (13,10), (12,11), (14,12), (15,13), (17,13), (16,14), (18,15)\}$ . Elementele acestei liste sunt perechi de poziții din cele două propoziții care sunt asociate prin corespondența de traducere. Dacă un cuvânt dintr-o propoziție nu are corespondent în propoziția echivalentă, atunci absența acestuia se marchează prin cifra 0.

- **LexPar** (Ion, 2007) este o aplicație ce se bazează pe Modelul de Atracție Lexicală a lui Yuret (MAL, (Yuret, 1998)) pentru a analiza legăturile sintactice între cuvinte. În viziunea lui Yuret, atracția lexicală este o măsură a afinității de combinare a două cuvinte într-o propoziție. Dacă două cuvinte sunt „atruse lexical” într-o propoziție, atunci probabilitatea ca ele să se combine și în alte contexte este semnificativă. De aceea, două sau mai multe cuvinte care se atrag lexical, împreună cu traducerea lor într-o altă limbă, se constituie într-un exemplu bun de traducere. Prin intermediul aplicației LexPar, corpusul de lucru este adnotat cu atributul „head”, a cărui valoare (un număr întreg) reprezintă poziția unui cuvânt din propoziție de care se „leagă” printr-o dependență sintactică forma adnotată. De exemplu, pentru Figura 1, prezența în corpus a atributului  $head=„3”$  pentru cuvântul „payments” indică faptul că „payments” se leagă de „all”, aflat în poziția 3 în propoziție. Atributul „head” nu indică centrul

EXPERIMENTE DE TRADUCERE AUTOMATĂ BAZATĂ PE EXEMPLE  
PENTRU LIMBILE ENGLEZĂ/ROMÂNĂ

constituentului sintactic și sensul relației de dependență dintre cele două cuvinte, ci doar faptul că aceste cuvinte sunt legate între ele din punct de vedere sintactic.

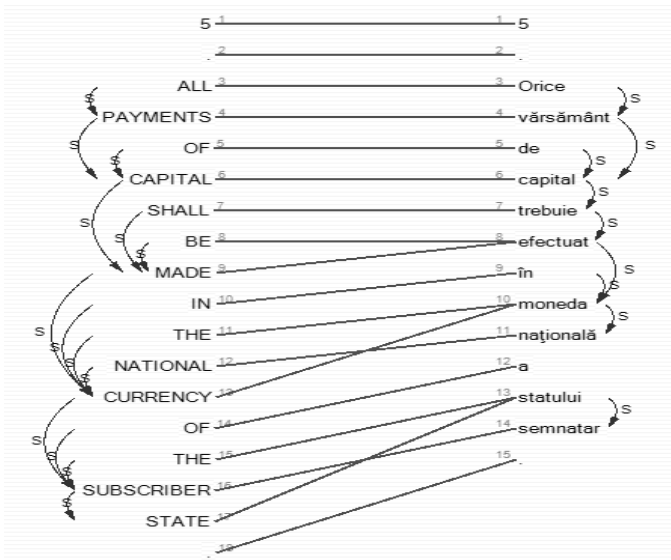


Figura 1. Vizualizarea alinierilor și legăturilor pentru o unitate de traducere din corpus. Numerele reprezintă pozițiile cuvintelor în propoziție. Corespondențele de traducere sunt marcate prin linii. O săgeată marchează existența unei legături de dependență sintactică între cele două cuvinte pe care le unește. Sensul săgeții este irelevant pentru sensul legăturii de dependență.

### 3. Baza de date cu exemple

În construcția bazei de date cu exemple am dorit să pornim de la ceea ce Yuret numea atracție lexicală sau afinitate de combinare între cuvinte. Am menționat deja că această atracție lexicală crește probabilitatea de asociere a două sau mai multe cuvinte și în alte contexte decât cel în care s-au identificat legăturile, ceea ce face ca secvența respectivă de cuvinte să se constituie într-un exemplu de traducere mai bun decât o simplă n-gramă. De asemenea, atunci când descompunem propoziția de tradus în subsecvențe care se suprapun peste baza de date folosim același concept de atracție lexicală pentru a decide granițele subsecvențelor. Se știe că pentru extragerea unui exemplu de traducere nu este îndeajuns să stabilim o strategie de divizare a propozițiilor în subsecvențe de cuvinte, ci este necesară stabilirea de corespondențe între subsecvențele dintr-o propoziție și traducerile lor în propoziția echivalentă. În acest scop, vom utiliza alinierea lexicală produsă de YAWA .

Revenind la exemplul din Figura 1, se observă că legăturile trasate cu LexPar tind să se grupeze prin imbricare și să descompună propoziția prin înlănțuire. Aceste proprietăți sugerează mai multe descompuneri posibile ale propoziției și implicit extragerea unor subsecvențe de lungimi diferite, dar care să fie compuse din cuvinte între care există fenomenul de atracție lexicală de care am vorbit.

Am denumit *superlegătură* un vector de forma  $S = (poziție_1, \dots, poziție_s)$  unde:

- $poziție_i$ , cu  $i \in [1, s]$  reprezintă poziția unui cuvânt într-o propoziție dată P;
- $s$  este lungimea lui S;
- există un vector de legături de forma  $[(poziție_1, poziție_2), (poziție_1, poziție_3), \dots, (poziție_1, poziție_s)]$  sau un vector de legături de forma  $[(poziție_1, poziție_s), (poziție_2, poziție_s), \dots, (poziție_{s-1}, poziție_s)]$  ce caracterizează secvența de cuvinte descrisă de

pozițiile din  $S$ ; deoarece legăturile care se intersectează sunt filtrate de LexPar, o superlegătură nu poate avea decât o formă imbricată precum cea din Figura 1.

Am denumit *lanț* un vector de forma  $L = (\text{poziție}_1, \dots, \text{poziție}_l)$  unde:

- $\text{poziție}_i$ , cu  $i \in [1, l]$  reprezintă poziția unui cuvânt într-o propoziție dată  $P$ ;
  - $l$  este lungimea lui  $L$ ;
  - există un vector de legături de forma  $[(\text{poziție}_1, \text{poziție}_2), (\text{poziție}_2, \text{poziție}_3), \dots, (\text{poziție}_{s-1}, \text{poziție}_s)]$  ce caracterizează secvența de cuvinte descrisă de pozițiile din  $L$ .
- Trebuie remarcat că o pereche de poziții  $(\text{poziție}_i, \text{poziție}_{i+1})$  nu reprezintă în mod obligatoriu poziții consecutive în propoziție.

Extractorul de exemple de traducere pe care l-am construit, *ExTract*, primește ca date de intrare corpusul de lucru precum și un fișier care conține pentru fiecare propoziție din corpusul de lucru, alinierea lexicală asociată de YAWA Vom descrie modul de procesare a unei singure unități de traducere din corpul  $U$ , care se repetă până la prelucrarea întregului document. O unitate de traducere din corpul de lucru este alcătuită dintr-o propoziție în limba engleză *Pen* și traducerea ei *Pro* în limba română, împreună cu toate adnotările ce rezultă în urma prelucrărilor descrise în capitolul 2. Aplicația lucrează în două etape:

**Etapa 1.** Se construiesc superlegăturile și lanțurile posibile atât pentru *Pen* cât și pentru *Pro*. Procedul este independent de limbă, deci putem simplifica descrierea acestuia prin generalizare. Fie  $P$  o propoziție oarecare dintre cele două propoziții ale unei unități de traducere. Informația despre atracția lexicală între cuvintele lui  $P$  este conținută de către atributul „head”. Pentru a face această informație accesibilă unor prelucrări ulterioare, construim vectorul de legături  $L$ , care, similar formalizării aliniierilor din secțiunea 3.3. conține *perechi de poziții* ale unor cuvinte. Elementele unei perechi de poziții din  $L$  sunt însă poziții ale unor cuvinte din *aceeași propoziție (P)* iar relația care le aduce împreună în aceeași pereche este cea de atracție lexicală între cuvintele pe care le indexează. Vectorul de legături se construiește parcurgând propoziția  $P$  cuvânt cu cuvânt iar pentru fiecare cuvânt care deține atributul „head”, se introduce în  $L$  perechea  $(\text{poziție\_cuvânt}, \text{valoare\_atribut\_head})$ . De exemplu, pentru propoziția în limba engleză din Figura 1 vectorul  $L$  conține perechile: (3,4), (4,6), (5,6), (6,9), (7,9), (8,9), (9,13), (10,13), etc. Vectorul rezultat al etapei 1) este similar vectorului  $L$ , dar poate conține nu doar perechi de poziții, ci liste de dimensiune variabilă care păstrează proprietatea de atracție lexicală între cuvintele pe care le indexează în  $P$ . Voi denumi acest vector *Lfinal* iar popularea lui cu elemente are loc în doi pași:

### **Pasul 1. Identificarea superlegăturilor și introducerea lor în *Lfinal*.**

$L_{final} \leftarrow \text{null};$

Pentru fiecare  $p$ , unde  $p$  este o poziție în  $P$

{ Pentru fiecare pereche  $(x, x')$  din  $L$

{ Dacă  $((x=p) \text{ sau } (x' = p))$

{ Lista  $\leftarrow \text{null};$

Dacă  $((x' - x) > 1$  și  $(x' - x) \leq 4)$  Lista  $\leftarrow x, x+1, \dots, x-1, x$  //completează pozițiile care lipsesc;

Altfel Lista  $\leftarrow x, x'$ ;

*Dacă* (Lista  $\neq$  null și Lista  $\notin$  Lfinal) Lfinal  $\leftarrow$  Lista} } }

*Observație:* În cazul în care avem de a face cu o legătură la distanță mare, probabilitatea ca aceasta să reprezinte o legătură simplă și nu un indiciu pentru o superlegătură crește. De exemplu, în propoziția „(1)*Dacă* (2)*la* (3)*încheierea* (4)*exercițiului* (5)*financiar* (6)*se* (7)*constată...*”, o legătură de tipul (1,7) nu atrage după sine existența unor legături care să producă superlegătura (1,2,3,4,5,6,7). Algoritmii de extragere a superlegăturilor a fost gândit pe baza observațiilor legăturilor din corpusul de lucru și reflectă proprietățile sintactice ale limbilor implicate. Pragul peste care o legătură nu poate fi tratată de către algoritm ca posibilă superlegătură este 4.

## **Pasul 2. Identificarea lanțurilor și introducerea lor în Lfinal**

Pentru identificarea unor subsecvențe de tip lanț în care să descompunem propoziția P, vom parcurge Lfinal de la cap la coadă efectuând următoarele operații:

- fie  $l_i$  un element din Lfinal; din lista de elemente  $l_{i+1}, \dots, l_n$  ce  $i$  se succed, alegem prima listă  $l'$  care verifică proprietățile  $l_i \cap l' \neq \emptyset$  și  $l' - l_i \equiv \emptyset$ ;

- lanț  $\leftarrow l_i \cup l'$ ;

- Lfinal  $\leftarrow$  lanț;

- Repetăm operațiile 1), 2) și 3) substituind  $l_i$  cu lanț;

Astfel, vom construi lanțuri formate din legături simple și superlegături, limitând numărul de elemente care se înlănțuiesc la 3. Această limitare, care este motivată de dorința de a nu supraîncărca baza de date în mod inutil, a fost stabilită tot pe baza observațiilor datelor din corpusul de lucru și reflectă faptul că probabilitatea ca un cuvânt să fie atras lexical de un cuvânt aflat la o distanță mai mare de 2 verigi ale unui lanț este foarte mică. De asemenea, proprietățile morfologice precum acordul între substantiv și verb sau între substantiv și adjectiv se transmit arareori la distanțe care să nu fie acoperite de un lanț cu 3 elemente. De exemplu, în Figura 1 acordul între substantivul “vărsământ” și adjectivul “efectuat” este surprins de lanțul cu 3 verigi (vărsământ de capital, trebuie, efectuat).

**Etapa 2.** Presupunând că am calculat (conform Etapei 1) cei doi vectori Lfinal<sub>en</sub> și Lfinal<sub>ro</sub>, în acest moment trebuie să stabilim corespondențe între elementele acestora pentru a construi exemplele de traducere. Fie A vectorul de alinieri asociat unității de traducere U. Am specificat în secțiunea 3.3 că un astfel de vector conține perechi de poziții ale unor cuvinte, (poziție<sub>i</sub>, poziție<sub>j</sub>), unde poziție<sub>i</sub> reprezintă poziția unui cuvânt  $w_i$  din Pen, poziție<sub>j</sub> reprezintă poziția unui cuvânt  $w_j$  din Pro, iar  $w_j$  este traducerea lui  $w_i$ . Structura care va formaliza corespondențele dintre listele lui Lfinal<sub>en</sub> și cele ale lui Lfinal<sub>ro</sub> este o listă ET de perechi de forma : (( $p_1, p_2, \dots, p_k$ ), ( $p_1', p_2', \dots, p_s'$ )), unde ( $p_1, p_2, \dots, p_k$ )  $\in$  Lfinal<sub>en</sub>, ( $p_1', p_2', \dots, p_s'$ )  $\in$  Lfinal<sub>ro</sub>, iar pentru fiecare  $p_i$  există un  $p_j'$  astfel încât ( $p_i, p_j'$ )  $\in$  A. Algoritmii de construire a listei ET este :

**Pasul 1.** Grupează toate alinierea 1:n și n:1 din A. Introduce în ET toate alinierea 1:n și n:1, precum și alinierea din A care nu sunt implicate în alinierea multiple. Acest pas

este necesar deoarece se dorește ca baza de date cu exemple să conțină și echivalențe de traducere la nivel lexical, care vor compensa absența unui dicționar;

## Pasul 2.

Pentru fiecare element listă  $l$  din  $L_{final_{en}}$

{ Lista  $l' \leftarrow null$ ;

Pentru fiecare element  $p$  al listei  $l$

{ Extrage din  $A$  lista  $C$  de perechi  $(x, x')$  pentru care  $p = x$ ;

Pentru fiecare pereche  $(p, x')$  din  $C$ ,  $l' \leftarrow x'$ ; }

Dacă  $l' \in L_{final_{ro}}$ ,  $ET \leftarrow (l, l')$  }

**Etapa 3** presupune o simplă recuperare a secvențelor de cuvinte indexate de listele de poziții din  $ET$  și afișarea lor într-un fișier ale cărui intrări au forma descrisă la pagina 4. Informațiile legate de lema și eticheta MSD a cuvântului sunt extrase din atributele „lemma” și „ana” cu care este adnotat corpul de lucru. În plus, în această etapă, fiecărui exemplu de traducere  $i$  se asociază un scor de încredere după cum urmează: dacă cele două liste de poziții (corespunzătoare membrului în limba engleză, respectiv membrului în limba română al exemplului de traducere) respectă condiția de consecutivitate (nu există elemente  $p_i, p_{i+1}$  astfel încât  $p_{i+1} - p_i > 1$ ) atunci scorul de încredere este 10; dacă cel puțin una dintre cele două liste de poziții nu respectă condiția de consecutivitate (implică legături la distanță), scorul de încredere este 5.

EXEMPLE DE TRADUCERE: LIMBA ENGLEZĂ	EXEMPLE DE TRADUCERE: LIMBA ROMÂNĂ
ALL(all,Di3)	Orice(orce,Di3--r---e)
PAYMENTS(payment,Ncnp)	vărsământ(vărsământ,Ncms-n)
OF(of,Sp)	de(de,Spsa)
CAPITAL(capital,Ncns)	capital(capital,Ncms-n)
SHALL(shall,Vaip)	trebuie(trebui,Vmip3s)
BE(be,Van) MADE(make,Vmps)	efectuat(efectua,Vmp--sm)
THE(the,Dd) NATIONAL(national,Afp) CURRENCY(currency,Ncns)	moneda(monedă,Ncfsry) națională (național,Afpfsm)
ALL(all,Di3) PAYMENTS(payment,Ncnp)	Orice(orce,Di3--r---e) vărsământ (vărsământ,Ncms-n)
OF(of,Sp) CAPITAL(capital,Ncns)	de(de,Spsa) capital(capital,Ncms-n)
PAYMENTS(payment,Ncnp) OF(of,Sp) CAPITAL(capital,Ncns)	vărsământ(vărsământ,Ncms-n) de(de,Spsa) capital(capital,Ncms-n)
SHALL(shall,Vaip) BE(be,Van) MADE(make,Vmps)	trebuie(trebui,Vmip3s) efectuat(efectua,Vmp--sm)
IN(in,Sp) THE(the,Dd) NATIONAL(national,Afp) CURRENCY(currency,Ncns)	în(în,Spsa) moneda(monedă,Ncfsry) națională(național,Afpfsm)
THE(the,Dd) SUBSCRIBER(subscriber,Ncns) STATE(state,Ncns)	statului(stat,Ncmsoy) semnatar(semnatar,Ncms-n)
PAYMENTS(payment,Ncnp) OF(of,Sp) CAPITAL(capital,Ncns) SHALL(shall,Vaip) BE(be,Van) MADE(make,Vmps)	vărsământ(vărsământ,Ncms-n) de(de,Spsa) capital(capital,Ncms-n) trebuie(trebui,Vmip3s) efectuat(efectua,Vmp--sm)
NATIONAL(national,Afp) CURRENCY(currency,Ncns) OF(of,Sp) THE(the,Dd) SUBSCRIBER(subscriber,Ncns)	moneda(monedă,Ncfsry) națională(național,Afpfsm) a(al,Tsfs) statului(stat,Ncmsoy) semnatar(semnatar,Ncms-n)
PAYMENTS(payment,Ncnp) OF(of,Sp) CAPITAL(capital,Ncns) SHALL(shall,Vaip) BE(be,Van) MADE(make,Vmps) IN(in,Sp) THE(the,Dd) NATIONAL(national,Afp)	vărsământ(vărsământ,Ncms-n) de(de,Spsa) capital(capital,Ncms-n) trebuie(trebui,Vmip3s) efectuat(efectua,Vmp--sm) în(în,Spsa) moneda(monedă,Ncfsry) națională(național,Afpfsm)



EXPERIMENTE DE TRADUCERE AUTOMATĂ BAZATĂ PE EXEMPLE  
PENTRU LIMBILE ENGLEZĂ/ROMÂNĂ

CURRENCY(currency,Ncns)
-------------------------

Tabelul 1. O parte dintre rezultatele ExTract pentru unitatea de traducere din Figura 3.2.

Corpusul de lucru a fost împărțit în corpusul de date pentru extracție de exemple (99% din corpusul total) și corpusul de date de test (1% din corpusul total). După rularea *ExTract pe* unitățile de traducere rezervate construirii bazei de date cu exemple, rezultatele au fost numărate și a rezultat un fișier cu 900.000 de exemple de traducere diferite, asociate frecvențelor lor în corpus. *ExTract* este o aplicație de sine stătătoare și nu face parte din fluxul de traducere. Construcția bazei de date cu exemple se face o singură dată și este urmată de o procedură de reorganizare a informației în 5 fișiere diferite, conectate printr-un index comun, pentru ca procedura de matching să nu supraîncarce memoria și să nu dureze foarte mult. Astfel, unei intrări din fișierul de ieșire al lui *ExTract* i se asociază un index de exemplu de traducere (*iet*) și toate informațiile conținute de intrarea respectivă sunt distribuite în cele 5 fișiere după cum urmează: fișierul **en\_forme**: *iet "enf\_1 enf\_2 ...enf\_n" MD5("enf\_1 enf\_2 ...enf\_n")*; fișierul **en\_leme**: *iet "enl\_1 enl\_2 ...enl\_n" MD5("enl\_1 enl\_2 ...enl\_n")*; fișierul **ro\_forme**: *iet "rof\_1 rof\_2 ...rof\_n" MD5("rof\_1 rof\_2 ...rof\_n")*; fișierul **ro\_leme**: *iet "rol\_1 rol\_2 ...rol\_n" MD5("rol\_1 rol\_2 ...rol\_n")*; fișierul **info**: *iet enm\_1 enm\_2 ... enm\_n rom\_1 rom\_2... rom\_3 frecvență, scor încredere*;

unde: **MD5** este o funcție hash utilizată des în criptografie, care asociază unui șir de caractere un număr natural pe 16 octeți, reprezentat în mod uzual ca o secvență de 32 de cifre hexazecimale; **enl** = leamnă cuvânt în limba engleză, **enf** = formă cuvânt în limba engleză, **enm** = MSD cuvânt în limba engleză, **\_1, \_2, \_n** = pozițiile cuvintelor în șirul extras în limba engleză; analog, **rol** = leamnă cuvânt în limba română, **rof** = formă cuvânt în limba română, **rom** = MSD cuvânt în limba română, **\_1, \_2, \_m** = pozițiile cuvintelor în șirul extras în limba română .

#### 4. Suprapunerea propoziției de tradus peste baza de date cu exemple (matching-ul)

Datorită formei în care a fost organizată baza de date, procedura de matching devine una de căutare într-o listă de numere naturale și este mult mai eficientă. De asemenea, separarea informației legată de forma cuvintelor și cea legată de leamnă în fișiere diferite este utilă în etapa de matching, când aplicația poate urma una dintre următoarele două direcții:

**1) suprapunere la nivel de formă ocurență:** *pasul 1)* încarcă fișierul *en\_forme*; *pasul 2)* descompune propoziția de tradus *P* în fragmente utilizând algoritmi de identificare a superlegăturilor și lanțurilor din secțiunea 4.1 (de fapt, etapa de descompunere este corespunzătoare Etapei 1 din secțiunea 4.1 aplicată membrului în limba engleză al unității de traducere și produce un vector similar cu  $L_{final_{en}}$ ); *pasul 3)* pentru fiecare dintre elementele vectorului de legături: recuperează din propoziție secvența de forme pe care o indexează – *șir\_forme* – și calculează  $MD5(șir\_forme)$ ; caută  $MD5(șir\_forme)$  în lista MD5 din *en\_forme* și extrage indexul *iet* pentru toate intrările identificate astfel; un singur identificator MD5 poate avea mai mulți indecși *iet* asociați, deoarece unui șir de forme în limba engleză îi pot corespunde mai multe

traduceri diferite în baza de date cu exemple; pentru fiecare *iet*, recuperează informația din *ro\_forme* și *info*.

**2) suprapunere la nivel de leamnă:** se execută aceiași pași ca în 1), înlocuind *en\_forme* cu *en\_leme*, *ro\_forme* cu *ro\_leme*, *șir\_forme* cu *șir\_leme*.

Pentru eficientizare, informația asociată unui exemplu de traducere candidat pentru traducerea finală se organizează într-un obiect din clasa *trans\_ex*, cu următoarele proprietăți: *iet*, *en\_formă*, *en\_lemă*, *ro\_formă*, *ro\_lemă*, *en\_msd*, *ro\_msd*, *frecvență*, *scor\_încredere*, *md5*, *poziție*: aceste atribute primesc valori în etapa de *matching*; atributul *poziție* are ca valoare lista de poziții din  $L_{final_{en}}$  corespunzătoare listei de cuvinte *en\_formă* și distinge între secvențe de cuvinte identice aflate în aceeași propoziție, în poziții diferite; *scor\_aliniere*, *scor\_traducere*, *lungime\_suprapunere\_en*, *lungime\_suprapunere\_ro*, *cel\_mai\_bun\_scor\_traducere*: aceste atribute primesc valori în procesul de *recompunere a propoziției* (vezi secțiunea 5). La finalul etapei de *matching* se construiește o listă *Frg* de obiecte de tip *trans\_ex* care reprezintă mulțimea tuturor fragmentelor în care P se poate descompune ce au fost găsite în baza de date cu exemple. Din această listă *Frg* se va recompune cea mai bună traducere posibilă a propoziției P în limba țintă.

### 5. Etapa de recombinare și adaptare

Pentru această etapă am ales Metoda Suprapunerii Maximale (Hutchinson et al., 2003), care combină fragmente care se suprapun și ale căror traduceri sunt consistente”. Autorii acestei metode exploatează intuiția că, atunci când două exemple de traducere se suprapun atât la nivelul fragmentelor sursă cât și la cel al fragmentelor țintă, probabilitatea ca o combinație a acestor exemple să producă o traducere corectă este crescută. În accepțiunea algoritmului (Hutchinson et al., 2006), suprapunere a două șiruri înseamnă de fapt suprapunere la stânga, adică: două șiruri  $s = \{w_1, w_2, \dots, w_n\}$  și  $s' = \{w_1', w_2', \dots, w_m'\}$  se suprapun dacă există un întreg  $p < m, n$ , astfel încât  $w_1' = w_{n-p}$ ,  $w_2' = w_{n-p+1}$ , ...,  $w_p' = w_n$ .

*Exemplu: (in this Regulation, în acest regulament), (this Regulation, prezentul Regulament): combinarea acestor exemple nu este indicată deoarece: 1. pentru fragmentul în limba engleză, condiția se suprapunere este îndeplinită doar parțial – există un întreg  $p=2$  care reprezintă lungimea suprapunerii „this Regulation”, dar  $p=m$ ; 2. pentru fragmentul în limba română condiția de suprapunere nu este îndeplinită.*

Combinarea exemplurilor de traducere este ghidată de o funcție de evaluare  $s(E)$  – E este un exemplu de traducere reprezentat ca un obiect din clasa *trans\_ex*.  $s(E)$  este calculată ținând cont doar de un alt exemplu de traducere, considerat a fi predecesorul lui E în soluția finală și depinde de următorii parametri: *overlap\_length\_en* (respectiv *overlap\_length\_ro*) – lungimea suprapunerii între membrul în limba engleză (respectiv română) al lui E și membrul în limba engleză (respectiv română) al predecesorului său; *length\_en*: lungimea fragmentului în limba engleză al lui E (un fragment mai lung este preferat); *gap*: distanța dintre primul cuvânt în fragmentul în limba engleză al lui E și ultimul cuvânt al fragmentului în limba engleză din predecesorul lui E; *alignment*:

scorul de aliniere, calculat ca medie ponderată a frecvenței, scorului de încredere și a unui scor MSD dat de un model de traducere pe MSD-uri (scorul este calculat pentru perechea (*en\_msd*, *ro\_msd*) asociată lui E iar modelul este extras din fișierul *info*).

$$s(E) = g * gap + s'(E),$$

$$s'(E) = 1 / (a * alignment + or * overlapp\_length\_ro + oe * overlapp\_length\_en + l * length\_en + 1).$$

Coefficienții *g*, *a*, *oe*, *or* și *l* sunt optimizați în mod experimental. Funcția de evaluare totală  $s(P)$ , este aditivă pe mulțimea exemplilor de traducere care descompun *P*. Pentru a minimiza numărul de calcule, (Hutchinson et al., 2003) propun o *tehnică best-first cu rază limitată de acțiune*, expandând primul cel mai bun candidat neevaluat la un moment dat și păstrând în memorie numai primii cei mai buni *n* candidați neevaluați. Algoritmul produce un vector de exemple de traducere (l-am numit *Soluție*) a căror combinație ar trebui să formeze cea mai bună traducere posibilă pentru propoziția *P* în condițiile unei anumite baze de date cu exemple. Indicațiile (Hutchinson et al., 2003) au fost implementate destul de exact, cu doar câteva modificări: s-a introdus ca parametru și lungimea suprapunerii pentru fragmentele în limba română; s-a fixat un coeficient  $l = \sqrt{length\_en}$  care poate favoriza (mai bine decât un scalar) un exemplu mai lung, asigurând valori mai bune pentru funcția de evaluare decât suma funcțiilor de evaluare a mai multor exemple scurte; s-a fixat o rază de acțiune  $n = 40$ .

O etapă finală de adaptare este necesară pentru transformarea informației conținute de vectorul *Soluție* într-o propoziție în limba țintă. Atât pentru opțiunea de matching pe leme cât și pentru cea de matching pe forme, adaptarea va implica eliminarea dublurilor produse de suprapunere, concatenarea secvențelor de cuvinte și câteva reguli de reordonare bazate pe secvențele de MSD-uri asociate secvențelor de cuvinte. Aceste proceduri au fost implementate, dar lista de reguli de rescriere poate fi extinsă. Pentru opțiunea de matching pe leme, este necesară și integrarea unui mecanism de generare a formelor ocurență înainte de etapa de adaptare. Această etapă nu a fost încă implementată, dar se va baza pe resursele și instrumentele dedicate generării morfologice dezvoltate la ICIA (vezi (Irimia, 2007) și (Tufiș et al., 2008)).

## 6. Concluzii

Am rulat componentele de matching la nivel de formă, recombinație și adaptare pe datele de test (600 de exemple de traducere) și am calculat scorul BLEU în raport cu o singură traducere referință. Această evaluare este prematură (scorul BLEU este sensibil la numărul de traduceri referință), dar un rezultat de 0,232 reprezintă un bun motiv pentru a continua implementarea componentei de matching pe leme și a concluziona printr-o evaluare și o comparare a celor două direcții. Analiza bazei de date cu exemple a condus la observarea unor erori sistematice produse de aliniatorul lexical (de exemplu, articolul hotărât în limba engleză este, adeseori, nealiniat cu substantivul în limba română sau aliniat incorect cu un alt substantiv decât cel corespunzător), a unor erori generate de LexPar (în special legături care nu pot fi captate) precum și a unui mic număr de erori generate de adnotatorul morfologic. Testarea algoritmului de extragere a exemplilor pe un corpus de mici dimensiuni (200 de unități de traducere) corectat la

nivel de aliniere lexicală a produs exemple de traducere corecte în proporție de 99%. Considerăm că o creștere a performanțelor instrumentelor de aliniere și analiză a legăturilor poate îmbunătăți semnificativ calitatea exemplilor din baza de date și, implicit, a rezultatelor aplicației de traducere. O altă soluție evidentă este creșterea dimensiunilor corpusului din care se extrage baza de date, dar această abordare presupune scăderea vitezei de răspuns a aplicației (spațiu mai mare de căutare) și necesită găsirea unor metode de optimizare a acesteia.

### Referințe bibliografice

- Brants, Thorsten (2000). TnT – a statistical part-of-speech tagger. In *Proceedings of the 6<sup>th</sup> Applied NLP Conference*, ANLP-200, pages 224-231, Seattle, WA.
- Erjavec T., R.Pavlov, L.Dimitrova, L.Sinapova, K.Simov, M.Tadi, V.Petkevi, HJ.Kaalep, N.Ide, G.Priest-Dorman, L.Tihanyi, T.Vradi, C.Oravecz D.Tufiş, AM.Barbu, P Holozan, V Gorjanc, M. Stabej (2001). Specifications and Notation for MULTEXT-East Lexicon Encoding. *MULTEXT-East Report*, Concede Edition D.
- Hutchinson Rebecca, Paul N. Bennett, Jaime Carbonell, Peter Jansen, Ralf Brown (2003). *Maximal Lattice Overlap in Example-Based Machine Translation*. School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213
- Ion, Radu (2007). *Word Sense Disambiguation Methods Applied to English and Romanian*. PhD thesis (in Romanian). Romanian Academy, Bucharest.
- Irimia Elena (2007). ROG- a Parafigmatic Morphological Generator for Romanian. *Proceedings of the 3<sup>rd</sup> Language Conference: Human Languages Technologies as a Challenge for Computer Science and Linguistics*, Poznan, Poland, pages 408-412 ISBN 978-83-7177-407-2.
- Steinberger Ralf, Bruno Pouliquen, Anna Widiger, Camelia Ignat, Tomaz Erjavec, Dan Tufiş, Dániel Varga (2006). The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. *Proceedings of the 5<sup>th</sup> International Conference on Language Resources and Evaluation* (Genoa, Italy).
- Tufiş Dan, Elena Irimia, Radu Ion, Alexandru Ceaușu (2008). Unsupervised Lexical Acquisition for Part of Speech Tagging. In *Proceedings of LREC 2008*, May 26 - June 1, Marakkech, Morocco. ELRA - European Language Resources Association. ISBN: 2-9517408-4-0.
- Yuret Deniz (1998). *Discovery of linguistic relations using lexical attraction*. PhD thesis, Department of Computer Science and Electrical Engineering, MIT.

# CONAN – DETECȚIA POSIBILELOR CONOTAȚII ALE UNUI TEXT

DAN ȘTEFĂNESCU, DAN TUFIȘ

*Institutul de Cercetări pentru Inteligență Artificială  
Academia Română*

*{danstef, tufis}@racai.ro*

## Rezumat

Ambiguitatea limbajului natural trece adesea neobservată, fiind de cele mai multe ori generată în mod involuntar. O propoziție în cazul căreia intenția autorului este evidentă în contextul original, poate fi înțeleasă total diferit atunci când este pusă într-un alt context, mai ales în situația în care ea conține ambiguități neobservate. Uneori acest lucru poate fi amuzant, alteori stânjenitor. Lucrarea prezintă o aplicație pe care am numit-o CONAN, realizată la ICIA, cu ajutorul căreia astfel de ambiguități pot fi detectate și, în funcție de intenția autorului, înlăturate, diminuate sau amplificate.

## 1. Introducere

Există diverse metode de modelare a proceselor de clasificare a opiniilor și diferite grade de granularitate în definirea acestor modele. De exemplu, în cazul clasificării recenziilor se poate evalua opinia/aprecierea generală (pozitivă, negativă sau neutră) a autorului cu privire la un anumit subiect al discuției. Clasificarea opiniei la nivelul documentelor este însă considerată a fi prea puțin granulară în cazul majorității aplicațiilor (de pildă un document poate exprima diferite opinii în raport cu diverse aspecte ale subiectului tratat - un produs poate fi excelent dar prețul său mult prea mare, o piesă de teatru poate avea o distribuție foarte bună dar un scenariu prost, etc.). De aceea mai toate sistemele avansate de analiză/clasificare opiniilor iau în considerare nivelul propozițional. La acest nivel, problemele tipice care apar includ identificarea propozițiilor ce exprimă o opinie, a relevanței acestei propoziții față de subiectul de interes, a agentului care exprimă opinia (care poate fi autorul textului sau o sursă citată de autor), polaritatea (pozitivă, negativă sau neutră) opiniei precum și intensitatea ei (puternică, slabă).

Toate metodele sau algoritmii (abia la început) folosiți pentru analiza subiectivității exploatează cuvinte sau expresii deja procesate, acestea fiind unități lexicale purtătoare de opinie sau sentiment. Aceste unități lexicale (pe care le vom numi senti-cuvinte) sunt în general manual codificate, extrase din corpusuri sau marcate în lexicoane precum General Inquirer sau SentiWordNet (Esuli & Sebastiani, 2006). În SentiWordNet fiecare sens al fiecărui cuvânt are asociat un marcaj <O, P, N> în care O, P și N reprezintă scoruri a priori (independente de context) de obiectivitate (O), subiectivitate pozitivă (P) și respectiv subiectivitate negativă (N) și în plus, pentru oricare triplet <O, P, N> ce marchează înțelesurile din dicționar există relația  $O+P+N=1$ .

În timp ce stabilirea faptului că o propoziție are caracter subiectiv (exprimă o opinie) și a detectării faptului că această propoziție se referă la subiectul de interes al investigației sunt mai puțin controversate, polaritatea ei poate fi în schimb problematică. Dificultatea

este generată de polisemia majorității cuvintelor și de faptul că în multe cazuri polaritatea subiectivității apriori (engl. "*prior subjectivity*") a senti-cuvintelor depinde de sensul lor contextual (uneori local, alteori global). Aparent, noțiunea de sens, așa cum este ea definită în cadrul SentiWordNet, ar rezolva problema. În realitate însă nu este așa. După cum arată (Tufiș, 2008a), este necesară stabilirea unei distincții clare între cuvintele intrinsec purtătoare de subiectivitate cu polaritate specifică, și cuvintele a căror polaritate ar trebui luată în considerare în funcție de context. Cazul al doilea face referire la situații de felul: „timpul de răspuns al monitorului este *lung*” comparativ cu „viața *lungă* a unui motor”. Se poate observa aici că polaritatea cuvântul *lung* depinde de substantivul pe care îl modifică.

Cercetările în această direcție au fost până nu demult de tip monolingv, orientate în marea lor majoritate către limba engleză datorită bogatelor resurse existente, necesare acestor tipuri de analiză<sup>1</sup> (Mihalcea et al., 2007). În ultimii ani însă, pentru tot mai multe limbi (printre care și limba română) au fost dezvoltate astfel de resurse, în general prin exploatarea textelor paralele și a lexicoanelor multilingve. O ipoteză fundamentală în transferul cros-lingual al caracterizării apriorice de subiectivitate a unui sens al unui element lexical dintr-o limbă este că aceasta (caracterizarea) este validă și pentru sensul corespunzător al echivalentului de traducere într-o altă limbă. După cunoștințele noastre nu există nici un experiment care să infirme această ipoteză, cel puțin la nivelul calitativ. Desigur, există posibilitatea ca raportul dintre subiectivitatea apriori pozitivă și cea negativă să varieze între echivalenții de traducere pentru o pereche de limbi, dar acest subiect necesită cercetări și experimente care depășesc cadrul investigației noastre prezente. Deocamdată, adoptând ipoteza corespondenței subiectivității apriori a echivalenților lexicali de traducere, metoda transferului cros-lingual al informației de subiectivitate reprezintă o modalitate comodă de creare a unor resurse lingvistice necesare detectării și evaluării automate a opiniilor exprimate textual. Deși SentiWordNet a fost dezvoltat pentru limba engleză, având în vedere existența de dicționare semantice monolingve<sup>2</sup> cu structură identică sau foarte asemănătoare cu Princeton WordNet (Fellbaum, 1998) pentru mai mult de 40 de limbi și că majoritatea acestor dicționare folosesc ca index interlingual chiar Princeton WordNet, marcajul de subiectivitate poate fi transferat (în virtutea echivalenței interlinguale a înțelesurilor) și exploatat în oricare dintre limbile respective (inclusiv limba română). Vom denumi în continuare orice wordnet astfel îmbogățit cu informația de subiectivitate, cu termenul de sentiwordnet. Apariția resurselor de tip sentiwordnet va conduce la amplificarea cercetărilor de subiectivitate pentru tot mai multe limbi, iar în viitorul imediat, la realizarea de noi aplicații în acest domeniu.

O astfel de aplicație va fi prezentată în cele ce urmează.

## **2. *Detectia posibilelor conotații textuale***

Majoritatea reclamelor publicitare pe care le vedem zilnic exploatează în mod inteligent ambiguitatea limbajului utilizând jocuri de cuvinte, asocieri surprinzătoare, imagini care

<sup>1</sup> citat: "mainly explained by the availability of resources for subjectivity analysis, such as lexicons and manually labeled corpora", (Mihalcea et al., 2007).

<sup>2</sup> <http://www.globalwordnet.org/>

împing către un anumit context de interpretare, pentru a promova diverse produse și/sau servicii. Multe din scurtele propoziții folosite în acest sens, atunci când sunt folosite în texte obișnuite, pot avea o parte din posibilele conotații mascate de context și astfel să fie nesesizate de către cititorul obișnuit. Observația este validă însă și în sens invers: anumite propoziții luate din contextul original și plasate în contexte noi, eventual cu grijă alese, pot fi purtătoare de mesaje complet noi, adesea nedorite. Este relativ ușoară identificarea, mai ales în interiorul textelor argumentative, a unor astfel de propoziții ce pot fi folosite într-un mod malițios în contexte care să poată induce o interpretare complet diferită și chiar opusă față de cea originală.

Metodele de analiză a subiectivității decid în general dacă o propoziție este subiectivă sau nu, iar în caz afirmativ stabilesc polaritatea și scorul de subiectivitate. Acest lucru se poate face ușor atât timp cât fiecare cuvânt din propoziție este dezambiguit la nivel de sens, iar pentru fiecare sens identificat există înregistrat într-un dicționar sentiwordnet un scor de subiectivitate. Pe lângă acest tip de analiză, aplicația descrisă în continuare, numită CONAN (CONnotation ANalyzer), poate fi folosită pentru o prelucrare textuală mai complexă: estimarea *variabilității conotative* a unei propoziții, reprezentând potențialul unui enunț textual de a-și modifica, în funcție de context, intensitatea sau chiar polaritatea opiniei subiective. Aplicația estimează pe o scară de la 0 la 1, măsura în care o propoziție, independent de contextul său curent de ocurență, poate fi interpretată obiectiv (O), subiectiv-pozitiv (P) sau subiectiv-negativ (N), elementele de interes pentru această estimare fiind senti-cuvintele din propoziția prelucrată. În mai toate cazurile, aceste scoruri sunt diferite. Experimentele noastre arată că propozițiile pentru care scorurile de subiectivitate (pozitivă și respectiv negativă) sunt ridicate și comparabile pot fi ușor folosite în „jocuri conotaționale” de care autorul poate fi conștient sau nu.

### 3. CONAN (CONotation ANalyzer)

Sistemul CONAN a fost implementat astfel încât să fie independent de limbă și ca atare el poate fi folosit pentru diferite limbi, atât timp cât textele de analizat sunt preprocesate corespunzător și atât timp cât există sentiwordnet-uri pentru aceste limbi.

Preprocesarea textelor, așa cum este ea necesară sistemului CONAN, include: segmentarea la nivel de unitate lexicală (tokenizare), adnotarea cu descriptori morfo-sintactici (tagging) și grupuri gramaticale (chunking). Acestea constituie operații fundamentale pentru aproape orice aplicație NLP, existând implementări pentru majoritatea limbilor. Instrumente ce realizează astfel de prelucrări pentru limba română au fost implementate în mai multe colective din țară sau din străinătate și există o bibliografie semnificativă. Pentru o prezentare detaliată a unor astfel de instrumente, implementate la Institutul de Cercetări pentru Inteligență Artificială (ICIA), a se vedea (Tufiș, 2008b). Recent, majoritatea programelor de preprocesare realizate la ICIA au fost făcute publice prin intermediul unor servicii web<sup>3</sup> nu numai pentru limba română, dar și pentru limba engleză (Tufiș et al., 2008).

După preprocesarea textelor, în faza a doua sunt identificate toate senti-cuvintele – cuvinte cu cel puțin o interpretare subiectivă (scorul de obiectivitate este mai mic decât

<sup>3</sup> <http://tutankhamon.racai.ro/ttlws.wsdl>

1) – folosind sentiwordnet-ul specific limbii prelucrate (în prezent, limba română și limba engleză). Literatura de specialitate indică faptul că, în problema analizei subiectivității, abordările ce ignoră ordinea și relațiile dintre cuvinte (eng.: *bag-of-words* – BoW) nu sunt potrivite deoarece subiectivitatea atribuită inițial poate fi schimbată de către contextul local al propoziției prin așa-numiții modificatori de valență (eng.: *valence-shifters*): intensificatori, moderatori și negații. Primii doi modificatori cresc și respectiv descresc scorurile de subiectivitate în timp ce negația completează valoarea subiectivității apriori a senti-cuvântului de sub incidența modificatorilor. Cum modificatorii nu acționează în mod necesar doar asupra senti-cuvântului din imediata vecinătate, adnotarea la grupuri gramaticale este importantă pentru a delimita raza de influență a acestora. De exemplu, în propoziția „NU este FOARTE *simpatic*”<sup>+</sup>, cuvântul *simpatic*<sup>+</sup> este un senti-cuvânt pozitiv, în timp ce cuvintele scrise cu majuscule (NU, FOARTE) sunt modificatori: negație, respectiv intensificator. Intensificatorul acționează asupra senti-cuvântului, în timp ce negația acționează asupra rezultatului NU (este FOARTE(*simpatic*)). În consecință, propoziția de mai sus are un scor negativ de subiectivitate. În (Tufiș, 2008a) am arătat că majoritatea scorurilor de subiectivitate care sunt greșit atribuite în SentiWordNet, se datorează abordării de tip BoW în cazul analizei definițiilor sensurilor. Cele mai multe mulțimi de sinonime (eng. *synset*) cu astfel de scoruri greșit calculate au în interiorul definițiilor modificatori care aparent au fost ignorați. Acest lucru ar putea explica de ce cuvintele *honest* (sensul 1) și sinonimul său *honorable* (primul sens) sunt considerate ca având o conotație mult mai negativă (0.5) decât pozitivă (0,25). Glosa atașată acestei serii sinonimice este: NOT DISPOSED to *cheat* or *defraud*; NOT *deceptive* or *fraudulent*.

CONAN acceptă texte de prelucrat atât de la tastatură cât și din fișiere. În cazul în care la intrare avem un fișier, aplicația presupune că fișierul este deja preprocesat și codificat în același mod în care platforma de servicii web a ICIA codifică documentele prelucrate – formatul XCES. În figura 1 se poate observa codificarea unei propoziții aparținând corpusului SEMCOR<sup>4</sup>, preprocesat de platforma TTL (Ion, 2007).

Figura 2 prezintă capturi de ecran ale aplicației având la intrare un fișier (ce conține și propoziția din figura 1). Utilizatorul specifică un mod de interpretare a propozițiilor din textul de intrare (obiectivă/positivă/negativă).

Fereastra din stânga-jos afișează analiza propozițiilor din fișierul de intrare. Propozițiile sunt ordonate în funcție de interpretarea obiectivității sau a polarității selectate (detaliate în continuare). Fereastra din mijloc-jos afișează scorurile de interpretare ale propozițiilor din fereastra din partea stângă-jos. Fereastra din dreapta-jos afișează informații precum indecșii interlinguali și definițiile seriilor de sinonime din care face parte cuvântul selectat de către utilizator în fereastra de analiză (fereastra stânga-jos).

<sup>4</sup> <http://www.cs.unt.edu/~rada/downloads.html>



CONAN – DETECȚIA POSIBILELOR CONOTAȚII ALE UNUI TEXT

```

<s id="br-a01.4.4.ro"><c>"</c>
<w lemma="doar" ana="14+,Rgp" chunk="Ap#1" wns="ili:ENG20-
00004331-b">Doar</w>
<w lemma="un" ana="21+,Timsr" chunk="Np#1">un</w><w
lemma="num&abreve;r" ana="1+,Ncms-n"
chunk="Np#1">num&abreve;r</w>
<w lemma="relativ" ana="14+,Rp" chunk="Np#1,Ap#2">relativ</w>
<w lemma="mic" ana="1+,Afpms-n" chunk="Np#1,Ap#2">mic</w>
<w lemma="de" ana="5+,Spsa" chunk="Pp#1,Ap#3">de</w>
<w lemma="asemenea" ana="1+,Afp"
chunk="Pp#1,Ap#3,Np#2">asemenea</w>
<w lemma="raport" ana="1+,Ncfn-n"
chunk="Pp#1,Np#2">rapoarte</w>
<w lemma="avea" ana="3+,Va--3s" chunk="Vp#1">a</w>
<w lemma="fi" ana="3+,Vap--sm" chunk="Vp#1">fost</w>
<w lemma="primi" ana="1+,Vmp--sm" chunk="Vp#1,Ap#4"
wns="ili:ENG20-00508949-v">primit</w><c>"</c><c>,</c>
<w lemma="avea" ana="3+,Va--3s" chunk="Vp#2">a</w>
<w lemma="spune" ana="1+,Vmp--sm" chunk="Vp#2,Np#3,Ap#5"
wns="ili:ENG20-00976600-v">spus</w>
<w lemma="juriu" ana="1+,Ncmsry" chunk="Np#3" wns="ili:ENG20-
07903245-n">juriul</w><c>,</c><c>"</c>
<w lemma="considera" ana="1+,Vmg"
chunk="Vp#3">consider&acirc;nd</w>
<w lemma="interes" ana="1+,Ncmsry" chunk="Np#4"
wns="ili:ENG20-05354775-n">interesul</w>
<w lemma="r&abreve;sp&acirc;ndi" ana="1+,Vmp--sm"
chunk="Np#4,Ap#6,Vp#4">r&abreve;sp&acirc;ndit</w>
<w lemma="&icirc;n" ana="5+,Spsa" chunk="Pp#2">&icirc;n</w>
<w lemma="alegere" ana="1+,Ncfn-n" chunk="Pp#2,Np#5"
wns="ili:ENG20-00171672-n">alegeri</w><c>,</c>
<w lemma="num&abreve;r" ana="1+,Ncmsry" chunk="Np#6"
wns="ili:ENG20-12816962-n">num&abreve;rul</w>
<w lemma="aleg&abreve;tor" ana="1+,Ncmpoy" chunk="Np#6"
wns="ili:ENG20-10058086-n">aleg&abreve;torilor</w>
<w lemma="&scedil;i" ana="31+,Crssp">&scedil;i</w>
<w lemma="m&abreve;rime" ana="1+,Ncfsry" chunk="Np#7"
wns="ili:ENG20-04819645-n">m&abreve;rimea</w>
<w lemma="acest" ana="2+,Dd3mso---e" chunk="Np#8">acestui</w>
<w lemma="ora&scedil;" ana="1+,Ncms-n" chunk="Np#8"
wns="ili:ENG20-08005407-n">ora&scedil;</w><c>.</c><c>"</c></s>

```

Figura 3: Propoziție codificată conform standardului XCES

The image displays the CONAN software interface, which is used for detailed analysis of text. It is divided into two main sections: 'Detailed Analysis' (top) and 'CONAN' (bottom).

**Detailed Analysis (Top Window):**

- Left Panel:** Shows a tree view of the analyzed text. The root node is 'Doar un număr relativ mic de asemenea rapoarte a fost primit a spus juriul considerând interesul răș...'. Subsequent nodes include 'un număr relativ mic', 'relativ mic', 'relativ', 'mic', 'de asemenea rapoarte', 'de asemenea rapoarte', 'asemenea rapoarte', 'a fost primit', 'a fost primit', and 'a spus juriul'. The text is highlighted in yellow.
- Right Panel:** Displays a list of English codes and their corresponding descriptions. For example:
  - ENG20-1303499-n: Raport între dimensiunile unor obiecte, între dimensiunile părților unui întreg sau între fiecare dintre aceste părți și întreg.
  - ENG20-1303761-6-n: Căut dîntre două matriți de același fel, exprimate în aceeași unități.
  - ENG20-00027929-n: o caracteristică comună pentru două părți sau entități.
  - ENG20-1300361-n: Legătură între două sau mai multe persoane, obiecte, fenomene, noțiuni pe care gândirea omenească o poate constata și stabili.
  - ENG20-0577461-2-n: Comunicare scrisă sau oral făcută de cineva în fața unei adunări, a unei autorități etc., cuprinzând o relație (oficială) asupra unei activități personale sau colective.
  - ENG20-0577154-n: Comunicare scrisă sau oral făcută de cineva în fața unei adunări, a unei autorități etc., cuprinzând o relație (oficială) asupra unei activități.

**CONAN (Bottom Window):**

- Left Panel:** Shows a tree view of the analyzed text, similar to the top window, but with a different set of nodes: 'Doar un număr relativ mic de asemenea rapoarte a fost primit a spus juriul considerând interesul răș...', 'un număr relativ mic', 'număr', 'relativ mic', 'relativ', 'mic', 'de asemenea rapoarte', 'de asemenea rapoarte', 'asemenea rapoarte', 'rapoarte', 'a fost primit', and 'a spus juriul'.
- Right Panel:** Displays a list of numerical scores and an 'Average Score'.
  - 4. 0.2975
  - 2. 0.278845153846154
  - 11. 0.2265625
  - 6. 0.203125
  - 5. 0.2
  - 7. 0.160576923076923
  - 8. 0.15
  - 1. 0.1425
  - 9. 0.106275
  - 3. 0.0833333333333333
  - Average Score = 0.171188082750583

The interface includes a menu bar with 'File', 'Input', 'Analysis', 'Utils', 'Test', 'About', and 'Default LILs'. The bottom status bar indicates 'Aici trebuie introdus textul de analizat.'

Figura 2: CONAN

CONAN exploatează informația furnizată de TTL despre grupurile gramaticale (eventual imbricate) identificate în propoziția curentă și construiește structuri arborescente similare structurilor recursive de constituenți. Reprezentările astfel obținute sunt folosite pentru a calcula scoruri pentru interpretări de subiectivitate și obiectivitate. Primul pas presupune selecționarea înțeleșurilor senti-cuvintelor cu scorurile cele mai ridicate relativ la interpretarea selectată de către utilizator. Algoritmul calculează apoi recursiv scorurile de interpretare pentru fiecare nod al arborelui făcând media aritmetică a scorurilor nodurilor sale copil. Pornind de la frunze (ce conțin senti-cuvintele), scorurile se propagă până când scorul întregii propoziții este calculat. Prin selecția din meniu a opțiunii *Analysis*, utilizatorul are posibilitatea de a-și concentra analiza pe o propoziție aleasă (vezi figura 2, partea superioară). Aplicația permite multiple alegeri de acest fel cu posibilitatea deschiderii mai multor asemenea ferestre concomitent. Astfel, utilizatorul poate compara diversele grade de subiectivitate ale diferitelor propoziții. După cum afirmam, programul oferă opțiuni multiple de interpretare în direcții de polaritate sau obiectivitate diverse: interpretare pozitivă, negativă sau obiectivă (pe acestea le vom numi interpretări principale), iar pe lângă acestea se oferă posibilitatea forțării interpretărilor în direcții dorite: forțează cea mai pozitivă interpretare, forțează cea mai negativă interpretare, forțează cea mai obiectivă interpretare, forțează cea mai non-negativă interpretare, forțează cea mai non-negativă pozitivă, forțează cea mai non-subiectivă interpretare. Cea mai simplă operațiune constă în afișarea polarității tuturor propozițiilor textului, în eventualitatea în care cuvintele au fost în prealabil dezambiguizate în ceea ce privește sensul.

În cazul în care textul este introdus de la tastatură, aplicația detectează dacă acesta (format din una sau mai multe propoziții) este sau nu preprocesat. În cazul în care nu este, textul brut este trimis serviciilor lingvistice web ala ICIA, servicii de care aminteam mai sus. Restul operațiilor se petrec ca în descrierea de la paragraful anterior.

Așadar, utilizatorul poate cere o analiză a tipurilor de interpretare principale. În acest caz sensurile considerate pentru senti-cuvinte sunt cele cu scorurile de polaritate/obiectivitate cele mai ridicate, pentru interpretarea dorită. Cuvintele care nu sunt senti-cuvinte sunt considerate de obiectivitate maximă: 1.

În timp ce scorurile de mai sus pot fi calculate doar pentru întreg textul de la intrare, opțiunea *Analysis* oferă posibilitatea schimbării rapide între diferite interpretări ale unei propoziții. Mai mult, utilizatorul poate cere două tipuri de interpretări forțate: în direcția unei polarități, sau opusă direcției unei polarități. Aceste două tipuri de interpretare sunt mai elaborate și le vom numi *tipuri de analiză complexă*, deoarece, spre deosebire de cazul interpretărilor principale, aplicația nu numai că face aceeași analiză a propozițiilor, dar și sugerează înlocuirea anumitor cuvinte în funcție de direcția interpretării cerute. Așadar, pentru ambele tipuri de analiză complexă, pasul inițial este același cu cel al tipurilor principale: cuvintelor le sunt atribuite înțeleșurile cu scorurile cele mai mari pentru interpretarea dorită și apoi scorurile pentru toate propozițiile sunt calculate. Pentru partea a doua însă, în cazul primului tip de analiză complexă, cuvintele sunt înlocuite de sinonime selectate din seriile sinonimice corespunzătoare înțeleșurilor deja atribuite, sinonime reprezentate de literalii ce pot avea sensuri cu scoruri mai mari (evident în alte serii sinonimice) pentru interpretarea selectată. Pentru a formaliza, să presupunem că avem cuvântul  $w$  cu sensurile  $m_1, m_2, \dots, m_n$  (corespunzând evident la  $n$

serii sinonimice). Pentru o anumită interpretare  $I$ , să presupunem că sensul cu scorul cel mai ridicat este  $m_i$ . Pentru acest sens,  $w$  are următoarele sinonime  $s_1, s_2, \dots, s_k$ . În mod evident, sinonimele corespund unor sensuri ale altor literali care mai pot avea și alte sensuri. Cu alte cuvinte, fiecare  $s_i$  este un sens pentru un literal  $L_i$  care poate avea multiple alte sensuri:  $m_{i1}, m_{i2}, \dots, m_{it}$ . Algoritmul selectează acel literal care are un sens având cel mai mare scor pentru interpretarea curentă (dintre toate celelalte sensuri ale tuturor literalilor considerați). Literalul câștigător este acela care corespunde expresiei:

$$\max_{i \text{ pentru } I} (\text{score}(m_{i1}), \text{score}(m_{i2}), \dots, \text{score}(m_{i|L_i|}))$$

Este clar că, în anumite cazuri, selecția literalilor pe acest criteriu poate conduce schimbarea sensului propozițiilor. Este tocmai ce ar dori să obțină un utilizator ce folosește această opțiune. În acest caz, utilizatorul vrea să-l forțeze pe eventualul cititor să interpreteze textul într-o anumită direcție. Înlocuirea cuvintelor se face la nivelul formei ocurență prin utilizarea lemei literalului (furnizată de WordNet), descriptorul morfo-sintactic (msd) asociat cuvântului original în faza de preprocesare și un tabel conținând forme de ocurență a cuvintelor limbii, împreună cu lemele și msd-urile corespunzătoare (tblwordform).

Așa cum am mai menționat, am numit al doilea tip de analiză complexă *interpretare care forțează în direcția opusă direcției unei polarități*. Diferența dintre această interpretare și cea anterioară constă în modul în care se face selecția literalilor care să înlocuiască cuvintele originale. În acest caz, pentru a înlocui un cuvânt, algoritmul selectează literalul având sensul cu cel mai mic (și nu mai mare) scor, în interpretarea inversă (și nu curentă). Literalul câștigător corespunde expresiei:

$$\min_{i \text{ pentru } I} (\text{score}(m_{i1}), \text{score}(m_{i2}), \dots, \text{score}(m_{i|L_i|}))$$

Motivația acestei opțiuni constă în dorința de a selecta sinonime pentru cuvintele originale astfel încât să se poată evita interpretările cu anumite polarități.

În cadrul aplicației se folosesc nuanțe de culoare cu ajutorul cărora utilizatorul are posibilitatea de a identifica imediat gradul în care diferite noduri din structura arborescentă corespunzătoare unei propoziții contribuie la scorul final al acesteia. Un simplu clic pe un nod afișează scorul corespunzător nodului.

O altă opțiune pe care CONAN o oferă utilizatorilor săi este aceea de a calcula scorul de *interpretativitate* a propozițiilor unui text. Definem scorul de interpretativitate al unei propoziții (ISP) ca fiind o mărime cantitativă a potențialului unei propoziții de a-și schimba conotația. Interpretativitatea unei propoziții se calculează ca sumă normalizată a scorurilor de interpretabilitate subiectivă ale propoziției (Eq2). Scorul de interpretativitate al unei propoziții este strâns legat de scorurile de interpretativitate ale cuvintelor ce o compun. Cu cât scorurile pentru cuvinte sunt mai mari, cu atât cel calculat pentru propoziție va fi mai mare.

*Interpretativitatea* unui cuvânt (ISW) se definește ca în ecuația (Eq1) pe baza scorurilor de maximă pozitivitate și respectiv maximă negativitate asociate sensurilor sale.

$$ISW(sw_k) = \frac{0.5 * (\max P(sw_k) + \max N(sw_k))}{1 + |\max P(sw_k) - \max N(sw_k)|} \quad (\text{Eq1})$$

$$ISP(\text{propoziție}_k) = \frac{0.5 * (\max P(\text{propoziție}_k) + \max N(\text{propoziție}_k))}{1 + |\max P(\text{propoziție}_k) - \max N(\text{propoziție}_k)|} \quad (\text{Eq2})$$

Justificarea intuitivă a acestor formule empirice constă în faptul că, în condițiile în care un senti-cuvânt are cel puțin două sensuri, unul marcat cu un scor mare de interpretare pozitivă iar celălalt cu un scor mare de interpretare negativă, acel cuvânt are un impact major în schimbarea conotației unei propoziții. O propoziție are o variabilitate conotativă cu atât mai mare cu cât scorul ei de interpretativitate este mai mare.

Scorul de interpretativitate al unui senti-cuvânt poate fi maxim 1 când are cel puțin două sensuri, unul de pozitivitate maximă (1), iar altul de negativitate maximă (1). În SentiWordNet-ul românesc, cuvintele cu interpretativitatea cea mai ridicată sunt adjectivele *prost* și *imoral* și substantivul *generozitate* (toate având 0.875).

Trebuie să menționăm că modificatorii de valență se găsesc în trei fișiere externe, ușor editabile, care sunt citite la fiecare lansare a aplicației. În momentul de față toți modificatorii au o influență uniformă asupra senti-cuvintelor: intensificatorii și modificatorii măresc sau scad scorurile argumentelor lor cu 20%, în timp ce negația comută între scorurile P și N. O abordare mai elaborată, aflată în construcție, va specifica mai multe trăsături pe care le poate avea un modificator: categoria gramaticală preferată, argumentul preferat și chiar numărul de sens preferat al argumentului, dacă este cazul. În plus se vor putea defini diferite grade de influență a modificatorului în funcție de argument.

#### 4. Rezultate și Concluzii

Majoritatea experimentelor au fost realizate pe corpusul SEMCOR, un corpus paralel Română-Engleză preprocesat, cu conținut divers. Versiunea de lucru utilizată conține 8146 de propoziții a căror analiză, indiferent de interpretarea dorită, durează doar câteva minute. Au fost efectuate analize de variabilitate a conotației atât pentru limba engleză cât și pentru limba română. Valorile obținute pentru limba română sunt ușor diferite de cele obținute pentru limba engleză, deși ordinea propozițiilor sortate după scorul de interpretabilitate se păstrează (lucru ușor previzibil datorită ipotezei de lucru enunțate în secțiunea de introducere). Diferența între scorurile de interpretativitate se datorează în primul rând numărului mult mai mic de adjective (principalele senti-cuvinte) existente în wordnetul pentru limba română față de wordnetul pentru limba engleză. De asemenea, cuvintele din wordnetul pentru limba română au, în general, mai puține sensuri implementate decât cuvintele din Princeton WordNet.

Deși analiza rezultatelor acestui experiment este abia la început, au fost identificate o serie de propoziții (în ambele limbi) cu scoruri de interpretabilitate contrazicând intuiția comună. Prin inspectarea mai amănunțită a acestor propoziții și a senti-cuvintelor componente au fost detectate multe marcaje <O,P,N> din SentiWordNet cu valori cel puțin discutabile (a se vedea discuția din secțiunea 3).

Cercetările noastre viitoare vor avea în vedere, pe lângă o evaluare cros-linguală a ipotezei de validitate a importului cros-lingual de marcaje de subiectivitate pe baza echivalenței de traducere, o mărire substanțială a inventarului de adjective în wordnetul

românească precum și o îmbunătățire a metodei de calcul al scorurilor de subiectivitate lexicală (apriori) pentru sensurile adjectivelor.

### Referințe bibliografice

- Esuli A. & F. Sebastiani. (2006). SENTIWORDNET: A Publicly Available Lexical Resource for Opinion Mining. *Proceedings of the 5th Conference on Language Resources and Evaluation LREC-06*, Genoa, Italy, pp. 417-422.
- Fellbaum C., editor. (1998). *WordNet: An Electronic Lexical Database*. MIT Press, Cambridge, MA.
- Ion, Radu (2007). *Word Sense Disambiguation Methods Applied to English and Romanian*. PhD thesis (in Romanian). Romanian Academy, Bucharest.
- Mihalcea R.; Banea C.; Wiebe J. (2007). Learning Multilingual Subjective Language via Cross-Lingual Projections. *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, Prague, Czech Republic, June, pp. 976-983.
- Polanyi L. & Zaenen A. (2006). Contextual Valence Shifters. In James G. Shanahan, Yan Qu, and Janyce Wiebe, editors, *Computing Attitude and Affect in Text: Theory and Application*. Springer Verlag.
- Tușiș D.; Ion R.; Ceașu A.; Ștefănescu D. (2008). RACAI's Linguistic Web Services. In *Proceedings of 6th Conference on Language Resources and Evaluation LREC-08*, Marrakech, Morocco.
- Tușiș D. (2008a). *Subjectivity mark-up in WordNet: does it work cross-lingually? A case study on Romanian Wordnet*. Invited talk on the Panel "Wordnet Relations" at the Global WordNet Conference, January 22-25, 2008.
- Tușiș D. (2008b). Algorithms and Data Design Issues for Basic NLP Tools. In Sergei Nirenburg and Oleg Kapanatze (eds). *Advances in Language Engineering for Low- and Middle-Density Languages*, NATO-ASI, 48 pages, IOS Press

# COMPLETAREA AUTOMATĂ A RESURSELOR LINGVISTICE ROMÂNEȘTI

PETIC MIRCEA

*Institutul de Matematică și Informatică, Academia de Științe a Republicii Moldova*

*[mirsha@math.md](mailto:mirsha@math.md)*

## Rezumat

În articol sunt examinate unele aspecte de completare a resurselor lingvistice utilizând proceduri de derivare automată. De asemenea este prezentată o descriere a particularităților afixelor românești, precum și a metodelor de generare automată a derivatelor cu prefixe și sufixe pentru a completa ulterior Resursele Reutilizabile pentru Tehnologia Limbajului Natural.

*Cuvinte cheie:* generarea automată a derivatelor, flexionare automată, validarea cuvintelor

## 1. Introducere

Aplicațiile ce țin de prelucrarea limbajului natural presupun crearea, completarea și folosirea resurselor lingvistice electronice. Completarea automată sau/și semiautomată a acestor resurse cu cuvinte generate în baza celor deja existente reprezintă o sursă importantă de îmbogățire a vocabularului prin mijloace exclusiv interne.

Scopul acestui articol va fi studierea particularităților afixelor românești și a metodelor de generare automată a derivatelor cu sufixe și prefixe pentru a completa ulterior Resursele Reutilizabile pentru Tehnologia Limbajului Natural<sup>1</sup> (RRTLN).

Inițial vom efectua o trecere în revistă a definițiilor și noțiunilor ce țin de derivare împreună cu anumite clasificări, precum și a metodelor existente de derivare automată. Un compartiment aparte este destinat resurselor lingvistice electronice folosite în studiul mecanismelor de recunoaștere și generare a derivatelor noi cu prefixe și sufixe. Ne vom opri asupra trei momente:

- a) Identificarea derivatelor,
- b) Derivarea cu prefixe și sufixe, urmată de flexionarea derivatelor obținute,
- c) Validarea derivatelor generate automat.

În această ordine de idei procesul de identificare a derivatelor este caracterizat de unele momente relevante în proiectarea ulterioară a algoritmilor de generare automată a derivatelor cu prefixe și sufixe.

În completarea resurselor lingvistice prin derivare automată apare tendința firească de a folosi cele mai frecvente afixe. Însă practic cele mai productive afixe se dovedesc a fi problematice datorită comportamentului neregulat. De aceea pentru cercetările noastre am ales acele afixe, care ne-au permis să stabilim niște legități de comportament mai simple, fără a invoca prea multe excepții.

---

<sup>1</sup> Lexiconul se conține pe site-ul <http://imi201.math.md/elrr/>

Din aceste considerente am operat cu prefixele *ne-* și *re-*, precum și cu sufixele *-tor* și *-bil*, ultimele fiind, la rândul său, frecvente în procesul de derivare cu prefixul *ne-*. Am inclus în examinarea noastră sufixul lexical verbal *-iza*, care este de origine neologică și foarte productiv la momentul de față cu o relație foarte strânsă cu sufixele lexicale *-ism* și *-ist*.

## 2. Particularitățile procesului de derivare

Derivarea reprezintă unul din mijloacele de îmbogățire a vocabularului care folosește resursele proprii ale limbii, pornind de la cuvinte existente în limbă, în particular, formarea de cuvinte noi ori cu sens nou prin adăugarea unor afixe la bazele lexicale existente. Prin *afix* se înțelege orice morfem care rămâne în afara rădăcinii, atunci când segmentăm un cuvânt. În denumirea globală de afixe se includ prefixele și sufixele. Cuvântul care este format prin adăugarea unui prefix sau sufix se numește *cuvânt derivat*. Unitățile de la care se formează cuvintele derivate se numesc *baze* sau *primitive*. De multe ori, prefixele și sufixele nu se adaugă direct la rădăcină, ci la așa-numita temă lexicală. Aceasta este comună tuturor formelor flexionare sau gramaticale ale unui cuvânt și este formată, în mod obligatoriu, dintr-o rădăcină și, cel puțin, un sufix sau prefix (Hristea, 1984).

În conformitate cu structura lor morfologică afixele sunt divizate în: *simple* (când nu pot fi divizate în unități mai mici) și *complexe* (când structura lor permite identificarea unor unități mai mici, dar întregul complex funcționează ca un element unic de derivare).

După poziția pe care o ocupă față de rădăcină, elementele adăugate la bază se împart în două mari categorii, și anume: unele care sunt plasate înaintea rădăcinii și se numesc *prefixe*, iar altele care sunt atașate la sfârșitul ei poartă denumirea de *sufixe*.

Sufixele sunt *lexicale* sau *derivative* în cazul în care ele servesc la formarea de noi cuvinte și *flexionare*, *morfologice* sau *gramaticale* când servesc la realizarea unor forme ale aceluiași cuvânt; sufixele lexicale se mențin în toate formele flexionare ale derivatului respectiv, pe când cele flexionare caracterizează anumite forme. În structura cuvintelor derivate sufixele lexicale sunt urmate de cele flexionare, iar acestea de desinențe. De cele mai multe ori, sufixele conferă cuvintelor noi create o anumită valoare semantică și morfologică.

Numărul prefixelor simple, incluse în (Graur&Avram, 1978) în urma identificării în cel puțin un derivat este de 86, iar a sufixelor lexicale de peste 600 (Hristea, 1984).

## 3. Metode existente de derivare automată

Dicționarele moderne se confruntă cu anumite deficiențe, care sunt obiecte de cercetare pentru lexicografi. Cu toate că dicționarele sunt permanent completate cu intrări noi, grație dezvoltării continue a limbii, sarcina elaborării unui vocabular complet rămâne una practic imposibilă. Mai mult ca atât, cunoaștem că în fiecare zi sunt create ad-hoc o mulțime de cuvinte noi, totuși în majoritatea cazurilor, dar nu în totalitate, ele rămân pentru o perioadă îndelungată nevalidate, cei care creează neologismele fiind mai puțin preocupați de „legalitatea” prezenței acestor cuvinte într-o limbă sau alta.



## COMPLETAREA AUTOMATĂ A RESURSELOR LINGVISTICE ROMÂNEȘTI

În cadrul modelului de derivare cu sufixe pentru limba italiană descris în (Carota, 2006) sunt investigate principalele modalități de formare a cuvintelor complexe prin sufixare în baza cuvintelor italiene morfologic simple. Pentru principalele tipuri de derivate italiene cu sufixe au fost identificate unele nuclee semantice. Astfel, pe de o parte se tratează interfața dintre morfo-sintaxă și semantică, iar pe de altă parte o interfață între sufixare și flexionare (Carota, 2006).

Caracteristicile relației dintre morfologia derivațională și sinonimie raportată la un dicționar electronic sunt studiate în lucrarea (Duško&Krstev, 2005) fiind ilustrate în baza derivării cuvintelor în limba sârbă. În acest context au fost generate noi leme cu sensuri previzibile. Acest procedeu a fost numit derivare regulată. Acest tip de derivare este utilizat în prelucrarea textului folosind dicționarul electronic morfologic al limbii sârbe și o colecție de traducere cu constrângeri lexicale.

În (Vilares et al., 2001) se descrie modul de lucru al unei aplicații de prelucrare a limbajului natural pentru extragerea informației. Autorii aplicației propun generarea familiilor morfologice ale unui cuvânt, fapt care va reduce varietatea lingvistică de documente indexate în limba spaniolă. Principalele caracteristici ale acestui sistem sunt: utilizarea minimă a resurselor lingvistice, costul mic „computațional” și independență față de motorul de indexare.

Pentru limba arabă a fost proiectat un sistem MORPHE (Leavitt, 1994) care reprezintă un analizor/generator elaborat ca o componentă a tehnologiei de traducere automată KANT (Nyberg&Mitamura, 1992). Deși MORPHE a fost proiectat ca să fie folosit atât în analiză, cât și în generare, în practică el este folosit doar pentru generare și doar în limbile care conțin prefixe și sufixe.

În (Santana et al., 2004) este descris un instrument capabil să recunoască, să genereze și să manipuleze relațiile morfo-lexicale ale cuvintelor cât și stabilirea cuvintelor primitive de la care au fost formate derivatele. Totodată permite lucrul cu prefixele și sufixele în parte cât și stabilirea relațiilor între ele și cuvinte derivate cu afixele respective.

Una din primele aplicații de derivare automată pentru limba română a fost sistemul FAVR, realizat în mediul Mac-ELU (Tufiș et al., 1996) care a avut drept scop acoperirea completă a morfologiei flexionare. Odată însă cu migrarea descrierii FAVR în acest mediu s-a abordat și descrierea proceselor lexicale. Sub aspect derivativ, sufixele și prefixele lexicale au un potențial productiv mare. În urma analizei atributelor specifice fiecărei părți de vorbire în parte, în descrierea morfologică implementată în Mac-ELU s-au utilizat 20 de categorii gramaticale. Clasificarea s-a efectuat nu numai în baza cerințelor prelucrărilor morfo-lexicale, ci și a granularității necesare analizei și, respectiv, a generării sintactice.

O altă aplicație ce merită atenție este AnMorph (Cristea&Forăscu, 2006). Ea reprezintă un mediu de dezvoltare și actualizare a modelului morfologic paradigmatic al unei limbi neaglutinative (modelul de bază al cuvântului flexionat consideră cuvântul ca fiind compus din rădăcină și o terminație). Programul compară formele introduse de utilizator cu acele care pot fi generate pornind de la o paradigmă deja existentă în baza de date a programului, și dacă asemănarea este confirmată, se generează restul formelor. Când se întâmplă acest lucru, utilizatorul doar verifică și validează partea tabelului generată automat. În afară de interfața pentru dezvoltare-actualizare, mediul oferă un editor

pentru dicționar și o colecție de paradigme, o componentă care permite verificări de consistență a datelor și a lematizorului.

#### **4. Resursele lingvistice electronice ale limbii române**

Aplicațiile ce țin de procesarea unui limbaj natural necesită în mare parte resurse lingvistice, care reprezintă cunoștințele lingvistice împreună cu datele suport (structurate într-o formă prestabilită) și programe asociate (Tufiș&Barbu, 2002).

O resursă importantă este dicționarul morfologic de limba română (DMLR) (Lombard&Gâdei, 1981). Acest dicționar conține 28932 de cuvinte care sunt împărțite în clase de flexionare în dependență de modul de formare a acestora. Pornind de la DMLR, au fost elaborate programe de flexionare pentru limba română (Cojocaru, 1997). Ele au contribuit substanțial la acumularea resurselor lingvistice.

Pachetul de programe “Produce program pentru aplicații lingvistice” a fost utilizat cu succes la implementarea corectorului de texte pentru limba română RomSP (Boian et al., 2000). Dezvoltarea celui din urmă a condus la implementarea RRTLN, care conține o bază de date cu informație lingvistică la nivel de cuvânt și un set de programe de gestionare (Boian et al., 2005a, 2005b). Astfel, lexiconul conține nu doar reprezentarea grafică a cuvântului, dar și informația despre partea de vorbire al lui. RRTLN are aproximativ 100000 de cuvinte de bază și circa un milion de flexiuni. De menționat, că un cuvânt poate avea mai multe intrări pentru diferite părți de vorbire.

#### **5. Identificarea derivatelor**

Drept sursă pentru recunoașterea derivatelor cu prefixe a servit lexiconul RRTLN și o listă de prefixe simple cu formele lor fonologice care sunt înregistrate în (Graur&Avram, 1978). Ținând cont de particularitățile prefixelor precum și a derivatelor lor, a fost elaborat un algoritm de extragere automată a cuvintelor derivate cu prefixe simple.

Cu mici schimbări algoritmul menționat mai sus s-a folosit pentru studierea problemelor de extragere a derivatelor cu sufixe din lexiconul RRTLN, în baza derivatelor cu sufixele *-tor* și *-bil*. Selectarea acestor sufixe e motivată prin existența unui număr mare de cuvinte cu aceste particule în lexicon. În urma verificării s-a stabilit că nu au fost găsite toate cuvintele derivate. Motivul este prezența alternanțelor vocalice în desinențele verbelor, fapt ce nu este luat în calcul în algoritmul (Petic, 2007a). În afară de aceasta lexiconul nu conține toate verbele de la care au fost formate derivatele cu sufixele *-tor* și *-bil*. În plus, s-a constatat că verbele de la care s-au format derivate cu sufixul *-tor* fără alternanțe vocalice/consonantice se termină în *a*, *i*, *ă* și *î*. Acest algoritm a fost implementat într-un program în limbajul de programare C++, în mediul Windows (Petic, 2007a). Totodată, programul este util în extragerea atât a prefixelor compuse cât și a sufixelor lexicale compuse în cuvintele derivate din același lexicon.

Verbele care se termină în *e* formează substantive în *-tor* doar cu ajutorul alternanțelor vocalice/consonantice. Numărul verbelor care se termină în *i* este foarte mic, în primul caz 7, în cel de-al doilea doar 1 – *dogorî*. În ceea ce privește litera *a* se observă că dacă verbul se termină în *ja*, *ua*, *va*, *xa*, atunci există o probabilitate destul de mare că se vor

## COMPLETAREA AUTOMATĂ A RESURSELOR LINGVISTICE ROMÂNEȘTI

forma derivate cu *-tor* fără alternanțe vocalice și consonantice. Numărul unor astfel de terminații pentru litera *i* este mai mare: *îi, și, țî, ai, bi, di, ei, fi, hi, ii, ji, li, mi, ni, si, vi și zi*. Totodată este ambiguă situația cu verbele în alte terminații, precum *a și i*, din cauza că pot fi atestate cuvinte derivate atât cu alternanță cât și fără alternanțe. În plus, numărul unor astfel de derivate este destul de mic.

S-a constatat că verbele de la care s-au format derivate cu sufixul *-bil* fără alternanțe vocalice/consonantice se termină *a și i*. Verbele care formează derivate cu sufixul *-bil* cu alternanțe vocalice/consonantice se termină în *a, e și i*. Este cert că verbele care se termină în *e* formează substantive în *-bil* doar cu ajutorul alternanțelor vocalice/consonantice. Se pune în evidență repetarea literelor *a și i*. În ceea ce ține de litera *a* se poate spune că dacă verbul se termină în *ș a, ț a, b a, g a, j a, l a, m a, u a, v a și x a* atunci există o probabilitate destul de mare că se vor forma derivate cu *-bil* fără alternanțe vocalice și consonantice. Numărul unor astfel de terminații pentru litera *i* este mai mic: *ăi, li, ri, ni, si și ti*. Ca și în cazul sufixului *-tor* rămâne ambiguă situația cu verbele în alte terminații, precum *a și i*, din cauza că pot fi atestate cuvinte derivate atât cu alternanță cât și fără alternanțe. În plus, numărul unor astfel de derivate este destul de mic.

Este interesantă situația cu alternanța literei *e* la sfârșitul verbului pentru formarea derivatelor în sufixele *-tor* și *-bil*. În cazul derivării sufixale de la verbe cu afixele *-tor* și *-bil* se observă unele alternanțe vocalice și consonantice atât la desinențe cât și în rădăcină, unele nu se atestă la flexionare.

### 6. Derivarea automată

#### 6.1. Derivarea automată cu prefixele *ne-* și *re-*

Studiind particularitățile prefixelor *ne-* și *re-* (Iordan, 1970), s-au obținut următoarele legități pentru prefixul *ne-* care permit îmbogățirea resurselor lingvistice:

- de la adjectivele derivate cu sufixele *-tor, -bil, -os* se formează adjectivele derivate cu prefixul *ne-* (de exemplu: *neconductor, nenobil, neinvidios*)
- de la participiile terminate în alomorfele *-at, -it, -ut* se formează adjectivele derivate cu prefixul *ne-* (de exemplu: *nelaureat, neiubit, nenăscut*);
- de la gerunzii se formează adjectivele derivate cu prefixul *ne-* (de exemplu: *nesuferind*).

Respectiv, legitățile pentru prefixul *re-* sunt următoarele:

- de la infinitivul verbelor se formează verbe derivate cu prefixul *re-* (de exemplu, *a regenera*);
- de la infinitivul verbelor se formează substantive derivate atât în sufixul *-re*, cât și în prefixul *re-* (de exemplu, *recitare*).

Legitățile formulate mai sus necesită cunoașterea doar a reprezentării grafice a cuvântului și a părții de vorbire a lui. La stabilirea lor s-a operat cu DMLR. Algoritmul de derivare analizabilă cu prefixele *ne-* și *re-* constă în examinarea cuvintelor din

lexicon și concatenarea cu prefixe din clasa celor, care se încadrează în categoriile stabilite de legitățile de mai sus (Petic, 2007b).

Astfel, programul elaborat în baza algoritmului menționat a îmbogățit lexiconul cu 397 cuvinte derivate cu prefixul *ne-* și 8556 cuvinte derivate cu prefixul *re-* (Petic, 2008a).

## 6.2. *Derivarea automată cu sufixe lexicale*

### 6.2.1. *Cazul sufixului lexical verbal -iza*

Analizând particularitățile de derivare a sufixului lexical verbal *-iza* s-a constatat:

- cuvintelor care se termină în *-an* sau *-ian* le poate fi atașat sufixul *-iza*, fără alternanțe vocalice sau consonantice, așa cum sufixul *-an* este mai scurt, deci afixul *-ian* poate fi inclus în *-an* (de exemplu: *alcaniza*);
- nu este similar cazul cu *-ean*, deoarece aici apare alternanța vocalică *ea->e* (de exemplu: *europeniza*);
- în cazul terminațiilor *-atic* (de exemplu: *dramatiza*), *-etic* (de exemplu: *cosmetiza*), *-otic* (de exemplu: *patriotiza*) și *-ific* (de exemplu: *științifica*) se înlătură ultimele două litere și se adaugă sufixul *-iza*, iar în alte cazuri la *-ic* se alipește, pur și simplu, *-iza*;
- în cazul terminației *-ură* ultima vocală este înlăturată și se alipește *-iza* (de exemplu: *caricaturiza*);
- există o relație strânsă între verbele în *-iza* și substantivele și adjectivele în *-ism* și *-ist*, care se manifestă prin apariția a numeroase serii de derivate de la aceleași teme (Petic, 2008b).

Ținând cont de cele expuse mai sus pentru generarea noilor cuvinte se vor verifica nu doar terminațiile, dar și existența substantivelor și adjectivele respective cu sufixele lexicale neologice *-ist* și *-ism*. În plus, se va verifica existența cuvintelor obținute în lexicon. Drept sursă pentru generarea noilor verbe în *-iza* a servit lexiconul RRTLN. Deoarece în lexicon sunt 1178 de cuvinte în *-ism* și 1285 în *-ist* este posibilă verificarea multiplă a terminațiilor cuvintelor înainte de a fi formate cuvinte noi.

Examinând cuvintele din lexicon și concatenându-le sufixul lexical *-iza* respectiv celor, care se încadrează nu doar în categoriile stabilite terminațiilor în (Petic, 2008b), dar care lipsesc în lexicon și permit formarea de serii de la aceleași teme cu unul din sufixele lexicale neologice *-ism* sau *-ist*, cu unele alternanțe vocalice, s-a construit un algoritm de derivare cu sufixul lexical verbal *-iza*, în baza căruia a fost elaborat un modul în limbajul C în mediul de programare KDevelop (sistemul de operare Linux OpenSuse 10.3), care generează cuvinte noi. În baza algoritmului expus în (Petic, 2008c) în cazul terminațiilor *-atic*, *-etic*, *-otic* și *-ific* s-au generat automat 420 de verbe.

### 6.2.2. *Cazul sufixelor lexicale -bil și -tor*

Din cele stabilite anterior se formează derivate cu sufixele *-tor* și *-bil* de la infinitivul prezent al verbelor (de exemplu: *cititor*, *caracterizabil*), în unele cazuri cu careva alternanțe vocalice sau consonantice (de exemplu: *dogorî*→*dogoritor*,

## COMPLETAREA AUTOMATĂ A RESURSELOR LINGVISTICE ROMÂNEȘTI

*dispune*→*disponibil*). Drept sursă pentru generarea noilor derivate cu sufixe a servit lexiconul RRTLN, în care sunt înregistrate 7796 de verbe distincte. Ținând cont că în lexicon sunt derivate cu sufixul *-tor* și *-bil*, pentru care au fost utilizate verbe, rămâne să fie folosite celelalte verbe pentru a fi generate noi cuvinte pentru lexiconul RRTLN.

Astfel pentru *-tor*, mai întâi, se vor genera derivate de la verbele care se termină în *a* și *i* fără alternanțe vocalice/consonantice, ca după această să se încerce să se genereze cele în *a*, *i* și *e* cu alternanțe vocalice și consonantice. Numărul de cazuri posibile pentru *-tor* de la verbe în *a* și *i* fără alternanțe este de 1140.

Pentru *-bil* se va proceda la fel doar pentru verbele în *a* și *i* fără alternanțe și după aceea cu alternanțe pentru literele *a*, *i* și *e*. Numărul de cazuri posibile pentru *-bil* de la verbe în *a* și *i* fără alternanțe este de 1962.

### 7. Problema validării și flexionării derivatelor generate automat

Derivatele noi, care au fost generate, ar trebui să fie corecte din punct de vedere morfologic și semantic. Unul din procedeele de validare a derivatelor constă în validarea manuală a fiecărui cuvânt generat în corespundere cu cerințele regulilor morfologice și semantice. Garantând calitatea rezultatului (în cazul când procedeul este efectuat de către un specialist în domeniu) ne confruntăm cu dezavantajele specifice unui lucru manual: resurse considerabile de timp, precum și posibilitatea comiterii unor erori.

Un alt mod de validare constă în verificarea prezenței cuvintelor derivate în documentele electronice existente pe Internet. Căutarea în Internet trebuie să se realizeze pentru documentele culese doar pentru limba română. Aici însă ne confruntăm cu o serie de dificultăți. Chiar cu opțiunea cu privire la limba setată este posibil să se găsească cuvinte în alte limbi. Este cazul cuvintelor *maciza* (limba spaniolă), *bariza* (limba arabă), *neautomobil* (limba cehă), *nemonolit* (limba croată) care au fost găsite de către motorul Google la căutare pentru limba română. În plus, apar deficiențe create de o eventuală segmentare a cuvântului căutat. Astfel, de exemplu, la încercarea de a valida în acest mod verbul *fataliza* s-a găsit „...o fată, Liza...”, în loc de *crisianiza* s-a găsit *Cristian Iza*. Există și cuvinte care reprezintă substantive proprii, în particular denumiri de companii, validitatea cărora trezește dubii, de exemplu, *SRL „Daniza”* și *SRL „Cariza”* găsite de către motorul de căutare nu pot confirma validitatea verbelor omonime generate automat.

Dincolo de cele expuse mai sus apare și dificultatea stabilirii condițiilor în care un cuvânt este valid. S-ar părea că numărul de apariții ale cuvântului în listă ar fi un criteriu obiectiv. De exemplu, pentru cuvântul *catiza* s-au găsit mai multe intrări, dar cu greu s-ar găsi argumente pentru a-l valida.

Mai pot apărea cazuri în care se formează un cuvânt derivat cu alt sens, de exemplu *negros* format de la adjectivul *gros*, ar trebui să aibă sensul „subțire”. În DEX există un astfel de cuvânt *negros* dar are alt sens „brun, brunet, negricios”.

O altă problemă este cea a stabilirii părții de vorbire pentru derivatele obținute în mod automat și flexionării lor ulterioare.

Tabel 1: Date statistice despre numărul cuvintelor derivate automat

Afixul	Num. de cuvinte generate automat	Num. de cuv. validate manual	Num. de cuv. validate automat	Num. de flexiuni
<i>ne-</i>	397 (100%)	362 (91%)	187 (47%-52%)	3740
<i>-iza</i>	420 (100%)	317 (75%)	76 (18%-24%)	2920
<b>Total</b>	817 (100%)	679 (82%)	263 (32%-39%)	6660

Cuvintele derivate cu prefixul *ne-* validate în mod manual, cât și folosind mijloacele *google.com* ca parte de vorbire sunt adjective. Aceste cuvinte vor moșteni clasa de flexionare de la baza derivatului. În procesul de flexionare au apărut unele situații ambigue soluționate cu ajutorul mijloacelor motorului de căutare *google.com*.

În cazul flexionării cuvintelor derivate cu sufixul lexical verbal *-iza* a apărut o situație problematică la stabilirea verbelor personale și impersonale. Determinarea acestui lucru s-a realizat manual, prin consultarea cu un specialist filolog.

La flexionarea derivatelor verbale cu prefixul *re-* de asemenea trebuie de stabilit dacă verbul este personal sau impersonal. Cuvintele derivate, conform legităților de derivare cu prefixe, vor moșteni clasa de flexionare de la baza derivatului.

Sufixul *-re*, care transformă verbele la infinitiv în substantive de genul feminin, se flexionează ca și toate celelalte substantive care se termină în *-re*, astfel cuvintele derivate nu vor moșteni clasa de flexionare de la baza derivatului.

Deoarece cuvintele sufixate cu *-tor* și *-bil* pot fi atât substantive cât și adjective, în fiecare caz aparte trebuie să decidem, dacă într-adevăr le putem flexiona ca substantive sau ca adjective.

## 8. Concluzii

Aplicațiile ce țin de derivarea automată s-au dovedit a fi destul de utile, în particular, pentru completarea resurselor lingvistice electronice, stabilirea legăturii între derivatele cu sufixe și semantică, generarea prin derivare a unor leme cu sensuri previzibile, generarea familiilor morfologice ale unui cuvânt pentru a reduce varietatea lingvistică de documente indexate.

Studierea particularităților procesului de derivare a permis stabilirea asemănarilor și deosebirilor care apar la procesul de prefixare și sufixare în limba română. Aceasta a condus la elaborarea algoritmilor necesari în prelucrarea cuvintelor derivate atât cu sufixe cât și cu prefixe. Totuși n-a fost posibil de a recunoaște toate cuvintele derivate, în special, cu sufixe, cauza majoră fiind prezența alternanțelor vocalice în procesul de derivare și lipsa în lexicoane a cuvintelor lemă ale tuturor cuvintelor derivate.

Rezolvarea problemei generării unor derivate noi inexistente în dicționare a fost ilustrată în baza unor afixe concrete. Aceasta a permis evaluarea rezultatelor pentru fiecare afix în parte. Totodată nu toate cuvintele generate pot fi considerate acceptabile, ci doar cele care au trecut printr-un proces de validare.

Cuvinte considerate acceptate au fost flexionate în mod automat cu ajutorul programelor existente, astfel valorificând posibilitățile și completând arsenalul lexicografic al lexiconul RRTLN al limbii române.

## COMPLETAREA AUTOMATĂ A RESURSELOR LINGVISTICE ROMÂNEȘTI

**Mulțumiri.** Sunt recunoscător dnei. dr. hab. Svetlanei Cojocarului și dnei dr. Elena Boian pentru ajutorul acordat la realizarea acestei lucrări.

### Referințe bibliografice

- Boian, E., Cojocarului, S., Malahova, L. (2000). Instruments pour applications linguistiques. *La terminologie en Roumanie et en Republique de Moldova*, Hors serie, No. 4.
- Boian, E., Ciubotaru, C., Cojocarului, S., Colesnicov, A., Demidova, V., Malahova, L. (2005). Lexical resources for Romanian. *Scientific Memoirs of the Romanian Academy*, ser.IV, vol. XXVI, București, România, pp. 267-278.
- Boian, E., Cojocarului, S., Ciubotaru, C., Colesnicov, A., Demidova, V., Malahova, L. (2005). Technologization of Romanian: linguistic resources, applications, tools. *Proceedings of the 4rd International Conference on Microelectronics and Computer Science*. Vol.II, pp. 519-522.
- Carota F. (2006). Derivational Morphology of Italian: Principles of Formalization, *Literary and Linguistic Computing*, Vol. 21, Suppl. Issue, pp. 41-53.
- Cojocarului, S. (1997). Romanian Lexicon: Tools, Implementation, Usage. In: Dan Tufiș, Poul Andersen (eds.). *Recent Advances in Romanian Language Technology*. ISBN 973-27-0626-0, Editura Academiei, I, pp. 107-114.
- Cristea D., Forăscu C. (2006). Linguistic Resources and Technologies for Romanian Language, *Computer Science Journal of Moldova*, Volume 14, Nr. 1 (40), pp. 34-73.
- Duško V., Krstev C. (2005). Derivational Morphology in a E-Dictionary of Serbian, In Zygmunt Vetulani (ed.), *Proceedings of the 2nd Language & Technology Conference*, Poznan, Poland, pp. 139-143.
- Graur Al., Avram M. (1978). *Formarea cuvintelor în limba română*, vol. II Editura Academiei, București.
- Hristea T. (1984). *Sinteze de limba română*, București, pp. 66-99.
- Iordan I. (1970). *Limba română contemporană*, Editura Academiei, București.
- Leavitt, JR. (1994). MORPHE: A Morphological Rule Compiler. Technical Report, CMU-CMT-94-MEMO.
- Lombard A., Gâdei C. (1981). *Dictionnaire morphologique de la langue roumain*, București, Editura Academiei, 232 p.
- Nyberg, E. H., Mitamura, T. (1992). The KANT System: Fast, Accurate, High Quality Translation in Practical Domains. In: *Proceedings of COLING92*.
- Petic M. (2007). Automatic extraction of the analysable formations with simple prefixes. *Proceedings of the Second International Conference of Young Scientists „Computer Science and Engineering-2007”*, Lvov, pp. 215-217.
- Petic M. (2007). Derivarea automată cu prefixele ne- și re- pentru adjective și verbe. *Proceedings of the International Conference BIT+2007*, Chișinău.
- Petic M. (2008). Specific features in automatic processing of the formations with prefixes, *Computer Science Journal of Moldova*, 4 1(7), pp. 209-222.

- Petic M. (2008). Probleme în popularea resurselor lingvistice electronice prin derivarea automată cu sufixul lexical verbal –iza. *Proceedings of the International Conference BIT+2008*, Chișinău.
- Petic M. (2008). Generarea automată a verbelor cu sufixul lexical –iza. *The 2nd International Conference „Telecommunications, Electronics and Informatics. Proceedeings*. Volume I, Chișinău, pp. 441-446.
- Santana O, Perez J., Carreras F., Rodrigues G. (2004). Suffixal and Prefixal Morpholexical Relationships of Spanish. *Lecture Notes in Artificial Intelligence*, Ed. Springer-Verlag, pp. 407-418.
- Tufiș, D., Barbu, A.M. (2002). Revealing Translator's Knowledge: Statistical Methods in Constructing Practical Translation Lexicons for Language and Speech Processing. *International Journal of Speech Technology* 5, pp. 199-209.
- Tufiș D., Diaconu L., Barbu A. M., Diaconu C. (1996). Morfologia limbii române, o resursă lingvistică reversibilă și reutilizabilă, *Limbaj și Tehnologie*, Editura Academiei Române, București, pp. 59-65.
- Vilares J., Cabrero D. M., Alonso A. (2001) Applying Productive Derivational Morphology to Term Indexing of Spanish Texts Source Lecture Notes In Computer Science; Vol. 2004, *Proceedings of the Second International Conference on Computational Linguistics and Intelligent Text Processing*, pp. 336 — 348.



## INDEX DE AUTORI

Apopei Vasile	11
Bîzdîgă Claudia	83
Bolea Cecilia	65
Burileanu Corneliu	31
Buzo Andi	31
Ceașu Alexandru	125
Cristea Dan	55
Curteanu Neculai	65
Dincă Nadia Luiza	93
Feraru Monica Silvia	21
Hanes Diana	31
Husarciuc Maria	65, 115
Iftene Adrian	105
Ion Radu	75
Irimia Elena	131
Jitcă Doina	11
Marcu Dana-Alina	105
Moruz Alex	65
Petic Mircea	151
Petrea Cristina	31
Pistol Ionuț Cristian	55
Popescu Vladimir	31
Rotaru Ancuța	105
Spătaru Mădălina	65
Spiță Doina	83
Ștefănescu Dan	141
Teodorescu Horia-Nicolai	21, 41
Trandabăț Diana	65
Tufiș Dan	141
Turculeț Adrian	11
Zbancioc Marius-Dan	41