

# Lingvistica românească teoretică și computațională

**Dan Tufis**

*Membru corespondent al Academiei Române*

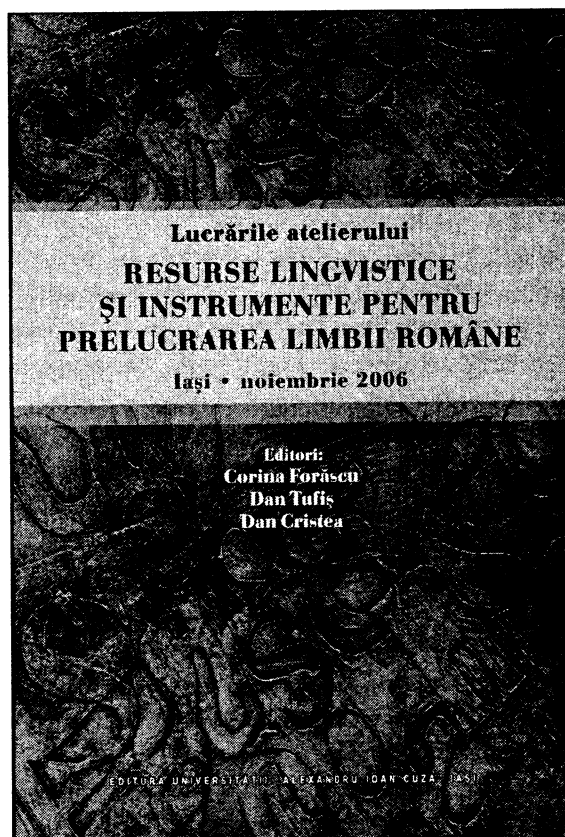
În contextul actual, al comunicării mediate de tehnologia informației și comunicațiilor, prin contrast cu termenul deja încetățenit de „limbă de circulație internațională“, putem vorbi, de „limbă de circulație electronică“ în ciberspațiu. Acest concept, pe lângă semnificația lui directă, are profunde implicații culturale, sociale, economice și nu în ultimul rând etice statuând dreptul fiecărui cetățean de a avea acces în propria limbă la cunoștințele, informațiile și serviciile ciberspațiului.

Voi încerca, în limita timpului disponibil, să schițez câteva dintre elementele semnificative ale evoluției în țara noastră în ultimii 5-6 ani a domeniului.

Această evoluție a urmat în mod firesc prioritățile internaționale născute din provocările și necesitățile celei de a doua generații a comunităților bazate pe web și servicii web, de trecerea de la un web al informațiilor la un web al cunoștințelor, al opiniilor și sentimentelor, dar și al afacerilor, al comerțului, al serviciilor de tot felul. Caracterul multilingual și multicultural al noului web acordă tehnologiilor lingvistice o importanță strategică.

Întâlnirea de azi este un moment important al evaluării domeniului tehnologiilor lingvistice pentru limba română, atât prin prisma onorantei participării la această masă rotundă, a domnului Comisar european Leonard Orban, prin deosebitul interes dovedit de prezența dumneavoastră numeroasă în Aula Academiei cât și prin faptul că cercetarea românească de profil începe să culegă roadele unei îndelungate perioade, pe care aș numi-o de construcție infrastructurală.

Educația de specialitate a fost și este o componentă esențială a acestei construcții. În anul 2001 au început primele programe de masterat în lingvistica teoretică și computațională (la Universitatea București) și prelucrarea limbajului natural (la Universitatea „A. I. Cuza“ din Iași). Tot în anul 2001, Academia Română a aprobat înființarea Comisiei de Informatizare pentru Limba Română, for consultativ al specialiștilor români. Câteva luni mai târziu Comisia a decis înființarea unui organism mai larg, având responsabilități organizatorice și executive: Consiliul de



Informatizare pentru Limba Română (ConsILR). ConsILR este o organizație deschisă, incluzând între membrii săi pe lângă specialiști români din țară și din străinătate, numeroși studenți, masteranzi sau doctoranzi de la principalele instituții de cercetare, învățământ sau industrie. Majoritatea manifestărilor științifice importante din România în domeniul tehnologiilor lingvistice au fost organizate de sau cu participarea nemijlocită a Comisiei și Consiliului pentru Informatizarea Limbii Române: Conferințele anuale ale ConsILR, școlile internaționale de Vară EUROLAN sau Conferințele internaționale SPED. Numeroși absolvenți români ai programelor de masterat amintite, sau a școlilor de Vară EUROLAN sunt actualmente membri ai unor echipe de cercetare în universități sau institute prestigioase, majoritatea lor lucrând la proiecte multilinguale ce includ limba română sau fiind înscriși în programe doctorale, pregătindu-și teze asupra prelucrării limbii române

la importante universități din SUA, Canada, Italia, Spania, Marea Britanie, Franța, Germania, Elveția ș.a.

O serie de tineri specialiști români și-au făcut deja un nume important în viața științifică internațională a domeniului și mulți dintre cei care acum sunt în străinătate sunt membri activi ai Consiliului de Informatizare pentru Limba Română, făcând o foarte bună propagandă cercetării și școlii românești. Acești tineri deosebiți păstrează o strânsă legătură cu formatorii lor din țară încercând și reușind de multe ori să atragă echipe din România în proiectele în care instituțiile lor actuale sunt implicate. În mediile academice internaționale se vorbește cu deosebit respect despre fenomenul românesc în domeniul PLN. Nu există conferință majoră a acestui domeniu, în care printre protagoniști să nu fie prezenți și cercetătorii români, din țară sau din străinătate.

Ca o consecință firească a vizibilității cercetării românești, limba română a intrat de curând în circuitul select al limbilor pentru care se organizează competiții de evaluare a tehnologiilor lingvistice, de regulă în context multilingual: alinierea lexicală – organizată de asociația nord-americană de lingvistică computațională Edmonton – Canada 2003, dezambiguizarea semantică automată a cuvintelor, organizată de asociația mondială de lingvistică computațională, Barcelona, Spania 2004, alinierea lexicală organizată de asociația mondială de lingvistică computațională Ann-Arbor, SUA 2005, întrebare-răspuns cross-lingual română-engleză, organizată de Forumul European de Evaluare Cross-linguală CLEF 2006 (Alicante, Spania), întrebare-răspuns în limba română organizată de Forumul European de Evaluare Cross-linguală CLEF 2007 (Budapesta, Ungaria). Sistemele dezvoltate de specialiștii români nu numai că au câștigat toate competițiile pentru limba română, dar s-au dovedit printre cele mai performante și în competițiile pentru limba engleză organizate în jurul unor probleme dificile ale PLN: rezolvarea automată a anafelor, organizată de DAARC (Discourse Anaphora and Anaphor Resolution Colloquium) în 2007, detectarea implicației lexicale și respectiv dezambiguizarea semantică automată a cuvintelor organizate de asociația mondială de lingvistică computațională, Praga, Republica Cehă 2007.

Prezența cercetătorilor români în proiectele de cercetare europeană din domeniul tehnologiilor lingvistice s-a îmbunătățit (de pildă numai ICIA și FII, instituții a căror activitate o cunosc foarte bine, participând în 18 astfel de proiecte, 4 aflându-se în derulare, iar alte 4 urmând să înceapă în lunile următoare). O mențiune specială aș dori să fac asupra programului CLARIN, unul din cele 28 de programe incluse în planul strategic de dezvoltare a infra-

structurii europene de cercetare. CLARIN (Common Language Research Infrastructure) are un orizont de realizare de 10 ani și un buget de aproape 1/4 miliard de Euro, din care contribuția UE va fi de peste 50%. Programul CLARIN reprezintă cea mai amplă acțiune europeană de până acum pentru armonizarea cercetărilor intra- și intercomunitare în domeniul resurselor și a tehnologiilor lingvistice. Acesta include actualmente 92 instituții din 33 de țări europene (3/4 din România). Ca o recunoaștere a contribuțiilor românești în domeniu, la nivelul decizional al acestui program România este reprezentată printr-un vicepreședinte al consiliului științific și printr-un membru al comitetului executiv.

Afirmam la începutul discursului meu că până nu demult cercetarea românească în domeniul tehnologiilor putea fi caracterizată ca traversând perioada construcției infrastructurale. În această perioadă de acumulare teoretică și experimentale sprijinul Comisiei Europene a fost fundamental.

Un câștig remarcabil al perioadei la care mă refer îl reprezintă introducerea cercetării în domeniul tehnologiilor limbajului, ca una dintre priorități în mai toate programele naționale ale principalilor finanțatori ai cercetării românești.

În contextul programului european CLARIN, multe din decalajele infrastructurale ale tehnologiilor limbajului pentru limba română vor fi recuperate mai repede și mai coerent decât până acum. Cu atât mai mult cu cât începând cu programul CEEEX și continuând cu PNII, finanțarea internă a devenit și sperăm că va rămâne semnificativă.

Ceea ce s-a realizat până în prezent în domeniul resurselor lingvistice standardizate (corpusuri adnotate, lexicoane multilinguale, modele statistice ale limbii române etc.) precum și numeroase programe specializate în prelucrarea unor segmente critice ale lanțului complet de prelucrare lingvistică presupus de realizarea unor aplicații practice, constituie „cărămizi” cu care se poate trece la „asamblarea” sau proiectarea și construirea unor sisteme complexe pentru prelucrarea textelor sau vorbirii în limba română.

Traducerea automată este numită pe drept cuvânt regina tehnologiilor lingvistice multilinguale, subsumând practic toate tipurile de resurse lingvistice și cele mai robuste metode și tehnici de prelucrare lingvistică.

Avansurile științifice din ultimii 10-15 ani, sprijinite de marile agenții de finanțare a cercetării (Comisia Europeană, NSF, Darpa, NICT în Japonia etc.) sau de a marile companii (IBM, Microsoft, Google, Altavista, ș.a.) au permis ca pentru acele limbi pentru care existau resurse lingvistice

compuționale adecvate să se realizeze în intervale de timp relativ scurte, sisteme automate de traducere foarte promițătoare, mult mai robuste, cu o mai largă acoperire lingvistică, mai flexibile, mult mai rapide și evident mai economice decât toate încercările celor 50 de ani de istorie anterioară a domeniului. Deși sistemele dezvoltate sau în curs de realizare în Uniunea Europeană, în SUA sau în Asia, nu reușesc deocamdată să atingă acuratețea și finețea traducătorilor umani, ele produc traduceri inteligibile, uneori chiar de bună calitate. Iar aceste traduceri se obțin în timpi de sute sau chiar mii de ori mai scurți și cu costuri neglijabile în raport cu traducerile umane.

Pentru limba română încercările, foarte puține care au fost până acum, au eșuat sau sunt foarte departe de performanțele sistemelor de traducere amintite. Credem însă că stadiul actual al cunoașterii științifice și tehnologice, ca și realizările obținute în România constituie pentru prima dată premise realiste de lansare a unui proiect de traducere automată pentru limba română. Dintre aceste realizări, în ultima parte a întâlnirii de azi, vor fi prezentate câteva instrumente pentru prelucrarea automată a textelor, componente esențiale în implementarea unui sistem competitiv de traducere automată pentru limba română:

- platforma web de servicii lingvistice bilingve (română și engleză) incluzând identificarea automată a limbii unui text (actualmente sunt recunoscute 22 de limbi), segmentarea lexicală, dezambiguizarea morfosintactică și lematizarea;

- ontologia lexicală Ro-wordnet, aliniată la nivel conceptual cu ontologiile lexicale pentru limbile engleză, bulgară, cehă, greacă, sârbă și turcă. Prin utilizarea tranzitivității relației de echivalență conceptuală cu limba engleză, se pot deriva alinieri conceptuale ale intrărilor din ontologia lexicală Ro-wordnet cu intrările ontologiilor lexicale ale altor zeci de limbi;

- sistemul DIAC<sup>+</sup> un program complex de inserare automată a diacriticelor în textele românești scrise fără sau parțial cu diacritice;

- MT-KIT – un ansamblu de instrumente pentru construcția modelelor de traducere, pentru alinierea lexicală automată, vizualizarea și editarea structurilor de control într-un sistem de traducere automată bazat pe abordări statistice.

Demonstrațiile vor fi făcute prin intermediul unor scurte filme, dar persoanele interesate ne pot contacta ulterior pentru demonstrații „live“ și detalii suplimentare.