

## TENDINȚE ACTUALE ÎN PROCESAREA LIMBAJULUI NATURAL PENTRU LIMBA ROMÂNĂ

Mircea PETIC, dr, lect. sup. univ.,  
Universitatea de Stat „Alecru Russo” din Bălți

**Summary:** *The article presents some trends in natural language processing research topic for Romanian language such as the use of computational linguistic resources. In addition, some useful applications are demonstrated, that use computational linguistic resources in both pre-processed texts as well as unstructured texts.*

**Key-words:** *natural language processing, computational linguistic resources, linguistic web services.*

În zilele noastre foarte multe domenii ale activității umane se informatizează. Nu este o excepție nici limbajul natural. Mai mult, este esențial pentru supraviețuirea unei limbi ca ea să fie folosită în sistemele de informare electronică. În acest sens, tehnologiile limbajului uman devin tot mai importante în era informațională. Limba română începe să se contureze ca una din limbile semnificative în privința resurselor informatice și a tehnologiilor aplicate ei. Totuși, chiar și cei mai fideli dintre susținătorii informatizării limbajului sînt de acord cu privire la imensa cantitate de necunoscut care însoțește acest proces. Astfel tehnologizarea unei limbi necesită consolidarea bazei teoretice din domeniul lingvisticii.

Eforturi pentru dezvoltarea resurselor lingvistice și instrumente de procesare a limbajului natural pentru limba română sînt depuse în puține centre academice din România (în general București, Iași și Cluj-Napoca) și Republica Moldova (Chișinău), cît și în unele centre din afara acestor regiuni. Astfel cele mai importante centre de cercetare în procesarea limbajului natural sînt:

- Institutul de Matematică și Informatică al AȘM, Chișinău;
- Universitatea Tehnică a Republicii Moldova, Chișinău;
- Universitatea „A. I. Cuza” din Iași;
- Institutul de Informatică Teoretică al Academiei Române (Filiala Iași);
- Institutul de Cercetări în Inteligență Artificială, București.

Una din primele aplicații de procesare a limbajului natural a fost cea de despărțire în silabe a cuvîntului. Automatizarea procesului de divizare în silabe nu reprezintă o problemă ordinară, deoarece, la etapa de prelucrare a cuvîntului, informația fonetică este inaccesibilă, deși anume această informație este importantă (Demidov&Verlan, 1996). Algoritmii elaborați și programele corespunzătoare, dezvoltate ulterior, se bazează pe regulile clasice de divizare a cuvintelor în silabe, avînd drept reper semnificația fonetică a literelor. Totuși, caracterul specific al limbii române nu permite formalizarea completă a regulilor de divizare a cuvintelor în silabe, în mare parte din cauza vocalelor. Acestea pot fi simple și complexe (așa-numitele semivocale), accentuate și neaccentuate, și regula de divizare în silabe depinde de categoria vocalei (frea-măt, re-cre-at, ghiont, ghi-o-cel etc). În acest fel, unicul lucru care poate fi stabilit în momentul în care se cunoaște cuvîntul este consecutivitatea vocalelor și a consoanelor. Informația fonetică nu este accesibilă. În această ordine de idei, este imposibilă acoperirea completă a regulilor de împărțire în silabe. Cele mai complicate situații apar în procesarea consecutivității de vocale (au, ea, ia, ie, io, ii, oa, ua), care pot fi (sea-ră, ier-ta, iod, au-gust) sau nu diftongi. Aceeași problemă e și cu

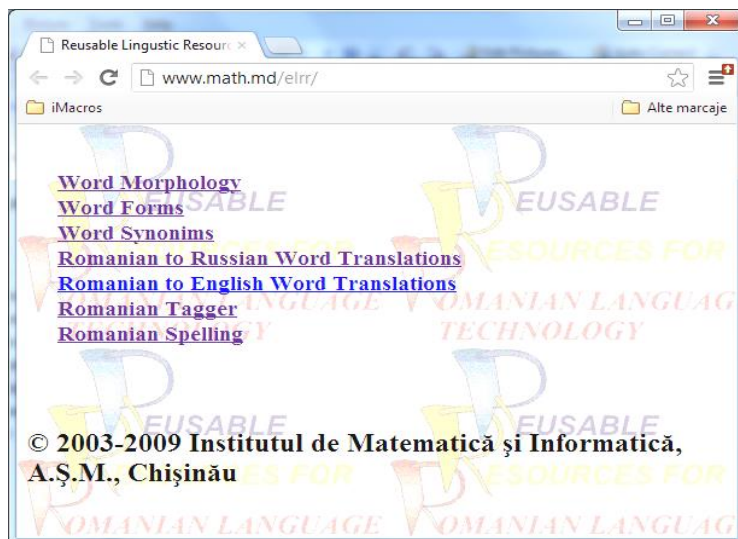
procesarea triftongilor. În acest sens, algoritmul de împărțire a cuvintelor în silabe din limba română a prevăzut procesarea a unui număr destul de mare de situații. Chiar dacă nu se pretinde la completitudinea algoritmului, totuși circa 70% de cuvinte din aproximativ 70000 de cuvinte din textele din domeniul științei și artei sînt împărțite corect în silabe (Demidov&Verlan, 1996).

Un alt tip de aplicații de procesare a limbajului natural sînt cele de flexionare automată. Prima încercare de acoperire a morfologiei flexionare pentru limba română a fost modelarea FAVR în mediul Mac-ELU. 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 (Tufiș, 1996).

O altă aplicație de flexionare automată ce merită atenție este sistemul AnMorph, descris în (Cristea&Forăscu, 2006). În prezent baza de date a aplicației acoperă doar parțial morfologia limbii române. În sistemul morfologic se compară formele introduse de utilizator cu acele care pot fi generate, pornind de la o paradigmă deja existentă în baza de date a sistemului și, dacă asemănarea este confirmată, se generează restul formelor. Dacă se întîmplă acest lucru, utilizatorul doar verifică și validează partea tabelului generată automat.

O altă abordare poate fi găsită în metodele de flexionare propuse de S. Cojocaru, care ține cont de clasificarea descrisă de A. Lombard și C. Gâdei în (Lombard&Gâdei, 1981). Algoritmul utilizează o gramatică funcțională de flexionare, care formalizează procesul de realizare a alternanțelor vocalice și consonantice și de formare a flexiunilor (prin concatenarea temei cuvîntului de bază cu seturile de terminații corespunzătoare). Pentru limba română această gramatică include 866 de reguli și 320 de seturi de terminații. Această metodă, numită metodă statică de flexionare, a contribuit substanțial la acumularea resurselor lingvistice, fiind aplicată pentru flexionarea a circa 30.000 de cuvinte pentru care era cunoscut grupul de flexionare din (Lombard&Gâdei, 1981). Metoda statică diferă de cea dinamică, în cadrul căreia se încearcă să se calculeze modelul de flexionare pentru fiecare cuvînt în parte. Algoritmul a fost testat pe 30.000 de cuvinte, care nu au fost incluse în clasificarea statică. De asemenea, au fost depistate unele iregularități (3% din numărul de cuvinte flexionate). Dincolo de aceasta s-a elaborat un set de programe care în regim semiautomat, generează toate formele flexionate. În anumite cazuri pentru obținerea corectă a flexiunilor dorite, se solicită implicarea utilizatorului pentru selectarea modelului potrivit (Ciubotaru et al, 2007).

Aplicațiile morfologiei flexionare descrise mai sus îmbogățesc resursele lingvistice, dar numai cu formele acestor cuvinte deja existente, fără a obține cuvinte cu sensuri noi, această funcție revenindu-i procesului de derivare pentru care în baza regulilor, constrîngerilor, schimbului de afixe, precum și a metodelor de proiectare a derivatelor a fost elaborat un algoritm de derivare automată care, fiind aplicat asupra a trei resurse lingvistice computaționale preprocesate, a produs 11191 cuvinte derivate noi. Aceste cuvinte, fiind ulterior flexionate, au completat dicționarul RRTLN (*Resurse reutilizabile ale tehnologiei limbajului natural*) cu 123106 intrări noi (Fig. 1) (Petic, 2011).



**Figura 1.** Interfața Web a RRTLN

RRTLN – conține o bază de date cu informație lingvistică la nivel de cuvânt și un set de programe de gestionare (Boian et al, 2005). Astfel, lexiconul conține nu doar reprezentarea grafică a cuvântului, dar și informația despre categoriile morfologice ale lui și posibilele funcții sintactice. RRTLN are aproximativ 100.000 de cuvinte leme și circa 1.000.000 de flexiuni. De menționat, că un cuvânt poate avea mai multe intrări pentru diferite părți de vorbire, astfel avînd o altă semantică, de exemplu, *bun* ca adjectiv în sens de binevoitor, amabil, *bun* ca adverb reprezentînd o aprobare și *bun* ca substantiv în sens de bunic sau de proprietate. Alături de informația lingvistică pentru limba română, complexul RRTLN conține și informații cu privire la sinonime și traduceri (engleză, rusă). Un set de programe de gestionare a bazei de date (Colesnicov, 2010) asigură efectuarea operațiilor necesare. Popularea BD cu informație morfologică se efectuează cu ajutorul unui program, care în baza lemei indicate, permite analiza și generarea automată a circa 90% din numărul total al cuvintelor leme (Boian et al, 2010), celelalte 10% prezentînd iregularități.

Rezultatele de mai sus au permis elaborarea corectorului ortografic pentru limba română care a fost elaborat și integrat în editorul de text Microsoft Word. În plus, în colaborare cu grupul editorial LITERA un CD cu corectorul ortografic cu o baza lingvistică de 1 milion de cuvinte a fost publicată.

O altă tendință actual de cercetare în procesare a limbajului natural constă în procesarea corpusurilor lingvistice computaționale. În lucrarea V. Bobicev a fost descris modul de procesare a mai multor corpusurilor cu ajutorul unor metode statistice. În experimentele cu restabilirea semnelor diacritice în text au fost obținute rezultate foarte bune de restabilire, aproximativ 98%. În experimentele de clasificare a documentelor aproximativ 97% de documente au fost clasificate corect (68 din 70), ceea ce poate fi considerat un rezultat destul de reușit. În experimentul de dezambiguizare morfologică programul a dezambiguizat absolut corect 95,3% de cuvinte. Din cele 4,7% de dezambiguizare greșită, numai în 2,5% a fost comisă o eroare în partea de vorbire. În restul cazurilor, 2,2%, a fost greșit detectat numai o caracteristică morfologică din cod. În concluzie experimentele efectuate demonstrează abilita-

tea metodelor statistice de a rezolva o clasă mare de probleme practice de prelucrare a textului (Bobicev, 2007).

În teza de doctorat a Victoriei Maxim este tratat un alt subiect promițător pentru cercetare ce ține de analiza computațională a corespondenței de afaceri. Drept rezultat au fost create următoarele resurse lingvistice computaționale: corpusul paralel bilingv englez-român cu texte ale scrisorilor de afaceri, adnotarea lui automată și dezambiguizarea manuală și o ontologie sumară a conceptelor de business. În plus s-a reușit formalizarea sintactică a Grupurilor Nominale ce conțin termeni și extragerea din corpus a unităților terminologice, validarea lor și organizarea într-un dicționar englez-român de termeni de afaceri. Ca rezultat a fost elaborat un program de traducere asistată de calculator a textelor comerciale din limba engleză în cea română (Maxim, 2008).

Un alt domeniu care se dezvoltă actualmente în procesarea limbajului natural ține de dezvoltarea serviciilor Web lingvistice. În continuare vor fi descrise succint posibilitățile <http://www.racai.ro/webservices/>. De menționat că toate aplicațiile care vor fi pomenite există ca aplicații de sine stătătoare (Tufiș et al, 2007).

Ceea ce ține de identificarea limbii se poate de spus că acest serviciu asigură identificare a limbii unui text scris într-una dintre cele 22 de limbi ale Uniunii Europene. Textul ar trebui să conțină un număr minim de 10-15 cuvinte.

Adnotarea morfo-lexicală se face cu o mare acuratețe (în jur de 98%) atât pentru limba română cât și pentru limba engleză. Există două implementări diferite pentru sarcina de adnotare morfo-lexicală: una se bazează pe paradigm HMM (Hidden Markov Models – Modele Markov Ascunse) iar cealaltă folosește abordarea Maximum Entropy (Entropie Maximă).

Un alt serviciu garantează codificarea de corpusuri paralele în format XCES pornind de la texte neprelucrate. O altă posibilitate a serviciilor este de a vizualiza cu ajutorul grafurilor hiporbolice la elementele celei mai mari ontologii lexical pentru limba română Ro-WordNet.

DIAC<sup>+</sup> este un serviciu care permite recuperarea automată a diacriticilor în texte în limba română scrise fără – sau scrise doar parțial – cu caractere diacritice. DIAC<sup>+</sup> este de asemenea disponibil și ca o aplicație de sine stătătoare, în forma unui DLL pentru MSOffice.

Unul din cele mai importante centre universitare care se ocupă cu succes de cercetarea în domeniul de procesare a limbajului natural este Universitatea „Alexandru Ioan Cuza” din Iași. Unele din direcțiile prioritare sînt:

- Parsarea articolelor din DEX pentru obținerea arborelui lexico-semantic al unei intrări de dicționar într-un fișier XML care reprezintă codificarea definițiilor din intrarea respectivă;
- Realizarea inferențelor textuale pe limba română presupune ca în baza a două fragmente de text se cere precizarea dacă înțelesul unuia din texte poate fi dedus din celălalt.
- Extragerea automată a definițiilor din texte este o parte a unui sistem complex de întrebare răspuns, pentru regăsirea răspunsurilor la întrebări de tip definiție.

Elaborarea și dezvoltarea unui sistem de întrebare răspuns este una din cele mai tentante sarcini în cercetarea din domeniul procesării limbajului natural. Primul sistem de întrebare răspuns românesc a fost dezvoltat în anii 80 și era reprezentat de o

interfață ce facilitează comunicarea cu o rețea semantică (care codifică cunoașterea). Astăzi astfel de sisteme folosesc documentele text ca baza de cunoaștere și integrează tehnici de prelucrare a limbajului natural pentru a găsi răspunsul la o întrebare pusă în limbajul natural. Arhitectura generală a unui sistem întrebare-răspuns constă din: o colecție de texte din Wikipedia (XML, HTML), componentele de filtrare a documentului (html→txt), pre-procesarea lingvistică (tokenizare, lematizare, etc.), analiza întrebării, crearea indexului și căutarea documentară, extragerea răspunsului și analiza rezultatului (Iftene et al, 2007).

Din cele menționate în articol constatăm că limba română începe să se contureze ca una din limbile semnificative în privința resurselor informatice și a tehnologiilor aplicate ei. Totuși, chiar și cei mai fideli dintre susținătorii informatizării limbajului natural sînt de acord că acest proces este încă supus hazardului. De aceea sînt deosebit de importante cercetările și elaborările, care permit să înaintăm cu anumiți pași concreți în acest domeniu.

### Referințe bibliografice

1. V. Demidov T. Verlan, An approach to the word division into syllables for Romanian language, *Computer Science Journal of Moldova*, v.4, n.1 (10), 1996, pp. 59-68.
2. Tufiș D. ș. a. Morfologia limbii române, o resursă lingvistică reversibilă și reutilizabilă. În: *Limbaj și Tehnologie*. București: Editura Academiei, 1996, p. 59-65.
3. Cristea D., Forăscu C. Linguistic Resources and Technologies for Romanian Language. In: *Computer Science Journal of Moldova*, Vol. 14, nr. 1 (40), 2006, p. 34-73.
4. Lombard A., Gâdei C. (1981) *Dictionnaire morphologique de la langue roumaine*, București, Editura Academiei, 232 p.
5. Ciubotaru C. ș. a. Contribuții la proiectul RoLTech: Platforma pentru tehnologia limbii române: resurse, instrumente, interfețe. În: *Lucrările atelierului „Resurse Lingvistice și Instrumente Pentru Prelucrarea Limbii Române*. Iași: Editura Universității, 2007, p. 171-179.
6. Petic M. Lexical derivation approaches for functional extension of computational linguistic resources. In: *Proceedings of the 8th International Conference “Linguistic Resources and Tools for processing of the Romanian language”* 8-9 december 2011, 26-27 april 2012. Bucharest, Editura Universitatii “Alexandru Ioan Cuza” Iasi, pp. 29-38
7. Boian E. ș. a. Technologization of Romanian: linguistic resources, applications, tools. In: *Proceedings of the 4th International Conference on Microelectronics and Computer Science*. Chisinau, Vol. II, 2005, p. 519-522.
8. Colesnicov A. Selection of processing tools in digital research environments. In: *Proceedings of „Actual problems of mathematics and informatics”*, UST – 80 ani. Chișinău, 2010, p. 88-89.
9. Boian E. ș. a. Application based on Reusable Linguistic Resources. In: *Multilinguality and Interoperability in Language Processing with Emphasis on Romanian*. București: Editura Academiei, 2010, p. 461-475.
10. Bobicev V. Metode și algoritmi statistici de procesare a textelor (în baza textelor în limba română). Teză de dr. în informatică, Chișinău, 2007. 229 p.
11. Maxim V. Analiza computațională a corespondenței de afaceri. Teză de dr. în informatică, Chișinău, 2008. 220 p.
12. Tufiș D. ș. a. Servicii web lingvistice ale ICIA. În: *Lucrările atelierului „Resurse Lingvistice și Instrumente Pentru Prelucrarea Limbii Române*. Iași: Editura Universității, 2007, p. 61-68.
13. Iftene A. ș. a. Construirea unui sistem de Întrebare Răspuns pentru limba română. În: *Lucrările atelierului „Resurse Lingvistice și Instrumente Pentru Prelucrarea Limbii Române*. Iași: Editura Universității, 2007, p. 109-118.