

# Menținerea identității prin limbă: instrumente și resurse lingvistice în format electronic

Gabriela Haja

Cercetările în planul prelucrării limbajului natural și a limbii vorbite, al tehnologiei limbajului în genere, sunt de maximă importanță în contextul exploziei informaționale contemporane, în cadrul căreia Internetul funcționează ca un mijloc cu posibilități încă insuficient exploatate de noi. De aceea, implicarea specialiștilor români din domenii devenite complementare – informatica și lingvistica – este firească.

În anul 2001 s-au creat și circumstanțele acestei colaborări, atunci când a fost înființată, la Academia Română, Comisia de Informatizare pentru Limba Română (CILR), cu scopul declarat al apărării identității, a individualității și a particularităților limbii române, prin susținerea studiilor lingvistice dintr-o perspectivă informațională. Încă înainte de această dată, a fost semnalată, la întâlniri și discuții care au avut loc în București, Iași și Chișinău, necesitatea creării unui Consorțiu care să realizeze un cadru propice comunicării dintre specialiștii lingviști și informaticieni. Grație unei finanțări asigurate printr-un proiect național, intitulat „Societatea Informațională – Societatea Cunoașterii”, s-a realizat un sit electronic dedicat Consorțiului de Informatizare pentru Limba Română (<http://consilr.info.uaic.ro>). Principiile și obiectivele acestui Consorțiu au fost prezentate de Dan Cristea în comunicarea *Resurse lingvistice românești și tehnologii informatice aplicate limbii române*<sup>1</sup>, la una dintre secțiunile Simpozionului Internațional *Identitatea limbii și literaturii române în perspectiva globalizării*, organizat de Institutul de Filologie Română „A. Philippide” și Asociația pentru Literatura Română și Cultura Poporului Român „Astra” – Despărțământul „Mihail Kogălniceanu”, la Iași, în 2002.

Potrivit acestor principii și obiective, cercetători din domeniul informaticii de la Facultatea de Informatică a Universității „Alexandru Ioan Cuza” și de la Institutul de Cercetări pentru Inteligență Artificială al Academiei Române din București, împreună cu lingviști de la Institutul de Filologie Română „A. Philippide” din Iași, s-au implicat în proiecte naționale și internaționale, care contribuie la realizarea de resurse, instrumente și tehnologii de prelucrare a limbii române, în toată complexitatea acesteia.

---

<sup>1</sup> Dan Cristea, Dan Tufiș, *Resurse lingvistice românești și tehnologii informatice aplicate limbii române*, în Ofelia Ichim, Florin-Teodor Olariu (eds.), *Identitatea limbii și literaturii române în perspectiva globalizării*, Iași, Editura „Trinitas”, 2002, p. 193–210.

Unul dintre aceste proiecte, care poate asigura o importantă cale de integrare a limbii române în circuitul global al comunicării informatizate, este BalkaNet<sup>2</sup>. Un alt proiect, în care sunt implicate Academia Română, prin trei dintre institutele sale (Institutul de Informatică Teoretică și Institutul de Filologie Română „A. Philippide” din Iași și Institutul de Cercetări pentru Inteligență Artificială din București) și Universitatea ieșeană, prin Facultatea de Informatică, este cel intitulat *Tehnologii lingvistice pentru Web semantic*, care își propune să realizeze, pe de o parte, corpusuri adnotate și, pe de altă parte, tehnologii pentru realizarea unui Web semantic pentru limba română. În esență, acesta va funcționa ca instrument indispensabil traducerilor automate.

S-au făcut până acum eforturi în acest sens, dar, datorită insuficienței resurselor, rezultatele sunt modeste<sup>3</sup>.

Este evident că participarea lingviștilor la realizarea acestor resurse (baze de date) este esențială și salutară. Proiectul amintit mai sus, *Tehnologii lingvistice pentru Web semantic*, lansat în acest an, sub programul-cadru INFOSOC, are drept principal scop ameliorarea situației actuale. Obiectivele detaliate ale proiectului nostru sunt realizarea de *corpusuri de texte adnotate*, de *standarde textuale* și de *instrumente de prelucrare* a limbii.

Adnotarea textelor (scrise și vorbite) se va face la mai multe niveluri: *fonetic*, *morfo-sintactic*, la nivelul *grupurilor nominale*, la nivel *sintactic* (specificându-se atât structurile sintactice, cât și relațiile dintre ele, așa numitele *tree-banks*), la nivel *semantic* (bănci propoziționale, *prop-banks*, care presupun specificarea structurilor de tip predicat-argument, cu identificarea statutului – opțional sau obligatoriu – și a rolului fiecărui argument: agent, beneficiar, loc, mod etc.), la nivelul „formal”, *al segmentării* (propoziții în cadrul frazelor, unități de discurs), la nivelul *relațiilor anaforice*: coreferința (echivalențe referențiale), referințe de tip parte-întreg, obiect-proprietate, posesor-obiect posedat etc., explicite sau implicite, la nivelul *sensului cuvintelor* (nivel care

<sup>2</sup> Vezi în acest sens: Dan Tufiș, *BalkaNet – tezaur lingvistic multilingv pentru limbile din Balcani*, în *op. cit.*, p. 177–192 și Dan Tufiș, Dan Cristea, *RO-BALKANET – ontologie lexicalizată, în context multilingv, pentru limba română* în volumul *Societatea Informațională – Societatea Cunoașterii*, București, Academia Română, Secția de Știință și Tehnologia Informației, Institutul de Cercetări pentru Inteligența Artificială, 2003, p. 137–164.

<sup>3</sup> Există în Internet mașini de tradus, dintre care amintim câteva: Softchim Translator, Millenium Computers, Edison Translator și cele realizate de InterTran, care face următoarea precizare: „InterTran pentru perechile română – engleză funcționează ceva mai bine decât pentru română – [celelalte limbi]. Traducerile sunt de slabă calitate (chiar și pentru română – engleză), translatorul fiind în curs de dezvoltare pentru limba română și alte limbi din țările în curs de dezvoltare. Actualmente *nu există translatoare on-line profesionale, gen SysTran, care să includă limba română*. InterTran-ul are interfața on-line, dar *translatorul pentru română are baza de date extrem de mică*” (s.n.) (<http://webenterprise.hast.gk/php/rotrans.php>).

implică chestiuni de *dezambiguizare semantică*<sup>4</sup>) și la nivelul *structurilor de discurs* (de exemplu, structura relațiilor retorice).

Aceste corpusuri pot fi adnotate semi-automat: pe baza unui model de analiză/adnotare realizat de specialistul lingvist, se pot genera adnotări automate, a căror corectitudine trebuie, firește, verificată. Adnotarea automată este facilitată de existența în format electronic a unor dicționare precum *Dicționarul explicativ al limbii române* (DEX), realizat de Institutul de Lingvistică „Iorgu Iordan” din București, și *Dicționarul de sinonime al limbii române* (DSR), redactat de Luiza și Mircea Seche. La acest nivel intră în funcțiune standardele textuale, care facilitează adnotarea elementară automată.

Pentru prelucrarea inteligentă/informatizată a limbajului sunt necesare instrumente. Acestea sunt componente software, a căror caracteristică este independentă de limbă, iar aplicarea lor la limba română este dependentă de folosirea resurselor lingvistice specifice, românești. Domeniul instrumentelor aparține informaticienilor, dar buna lor funcționare nu poate fi realizată decât cu ajutorul asistenței specialiștilor lingviști.

Realizarea acestor elemente implică, desigur, eforturi considerabile. Dar odată puse în practică rezultatele cercetărilor de acest tip, problema amenințării identității lingvistice și implicit culturale a unei națiuni nu se mai pune, de vreme ce comunicarea vorbitorului de limbă română, prin Internet, cu orice vorbitor, indiferent de limba lui maternă, devine astfel posibilă.

### **The preservation of identity through language: instruments and linguistic resources in electronic format**

*In the context of new informational communities and of globalization, it appears very clear that it is necessary to create different types of instruments for the Romanian (spoken and written) language. The collaboration between specialists in these two fields – linguistics and computer science – is very important in order to accomplish this task. BalkaNet and Linguistic Technologies are two of the results of this collaboration. The former is an international project and its aim is to create a lexical ontology for the Balcanic languages which could be related to EuroWordNet. The Romanian Academy, by three of its institutes (the Institute of Theoretical Computer Science, the “A. Philippide” Institute of Romanian Philology, Iasi and the Research Institute for Artificial Intelligence, Bucharest), and the “Al.I. Cuza” University of Iasi achieve the latter. Its main goal is to create a large database for computerized translations. If these goals are fulfilled, the problem of the threat of the globalization against linguistic and cultural identity becomes irrelevant.*

---

<sup>4</sup> Rezultate și soluții foarte interesante în domeniul dezambiguizării semantice prezintă Dan Tufiș, în articolul *Dezambiguizarea automată a cuvintelor din corpusuri paralele, folosind echivalenții de traducere*, din volumul *Societatea Informațională – Societatea Cunoașterii* ed. cit., p. 235–268.