

# Resurse lingvistice românești și tehnologii informatice aplicate limbii române

Dan Cristea<sup>♦</sup>, Dan Tufiș<sup>\*</sup>

<sup>♦</sup> Conferențiar doctor

Universitatea "Al. I. Cuza" Iași, Facultatea de Informatică

<sup>\*</sup> Cercetător principal 1, membru corespondent al Academiei Române  
Academia Română – Institutul de Inteligență Artificială

## 1. Un pericol al societății moderne

Atașați, așa cum suntem, ideii de globalizare și deschidere spre Europa și lume, nu putem însă neglija și un pericol, datorat imixtiunii abuzive prin mediile electronice (televiziune, internet) a limbilor intens mediatizate asupra celorlalte. Amenințate în acest mod sunt dialectele locale, deluvionate de tăvălugul lingvistic al limbilor oficiale, dar și limbi ale unor națiuni relativ mari dar, temporar, defavorizate economic, ce pot fi impurificate prin importuri abuzive în anumite registre tehnice sau segmente ale jargonului.

Într-un recent schimb de mesaje între membrii Comitetului Program al conferinței filialei europene a Asociației de Lingvistică Computațională<sup>1</sup>, ce urmează să fie găzduită anul viitor de Budapesta, Claire Gardent – președinta Comitetului Program propunea organizarea unui atelier de lucru cu titlul *Resurse și instrumente lingvistice pentru limbile est-europene*. În încercarea de a face acest atelier de lucru de interes și pentru participanți ce vin din alte zone decât estul Europei, sintagma "limbi mai puțin studiate"<sup>2</sup> a putut doar cu greu fi evitată prin iterații succesive asupra titlului, în final ajungându-se la formularea "limbi de densitate mai mică"<sup>3</sup>. E lesne de înțeles de ce formulări de acest tip, aplicate unei arii lingvistice din care face parte și româna sună atât de neplăcut pentru orice lingvist român, care cunoaște îndelungata strădanie a Academiei Române pentru construirea de dicționare, pentru definirea regulilor de gramatică românească și pentru crearea tezaurilor lingvistice ale limbii noastre. Cu toate acestea, veninoasele sintagme sunt, din păcate, justificate, pentru că tipul de tratament pe care îl insinuează ele ca fiind insuficient este cel de natură informatică.

Societatea contemporană are acum mijloace pentru a împiedica dispariția limbilor mici, pentru a încetini "dizolvarea" lor în limbile mari, ori pentru a stopa deteriorarea lor prin importuri abuzive. Limbile au în ele tot atâta creație divină ca și speciile. O limbă este, într-un anumit fel, un "jurnal de bord" al unui popor, pentru că limba unui popor păstrează urme ale trecerii lui prin timp așa cum o formă de relief memorează uneori, în straturile geologice, istoria unei specii. A pierde o limbă e la fel de dureros cu a pierde o specie. Dar, așa cum selecția naturală, cu sacrificiile ei, a întărit speciile, poate că procesul de "nivelare" a limbilor, pe care noi îl resimțim ca fiind o tragedie, trebuie de fapt privit cu detașare. Pentru că, ar putea spune cineva, el nu reprezintă altceva decât un act de selecție naturală cu care ne-a obișnuit evoluția, actul numit "cel mai tare rezistă", exercițiul eschimos de a izola în noaptea polară un bătrân neputincios și nefolositor pentru ca cei tineri și utili să aibă o îmbucătură în plus, să ducă mai departe făclia speciei. Aceasta a fost evoluția firească a naturii de la facerea ei și e posibil ca fără această "ladă de gunoi" a speciilor extinse ori a limbilor dispărute nici noi să nu fi fost atât de inteligenți și de bine făcuți și nici o limbă atât de frumoasă precum italiană, ca să luăm un exemplu, să nu fi existat.

În acest calcul liniștitor, nu putem însă să lăsăm la o parte angoasa ca nu cumva societatea modernă, prin noile canale de comunicație create, capabile de diseminare planetară și instantanee, să fi denaturat într-un fel compoziția supei primordiale în care se coace dintotdeauna evoluția limbilor. Pentru ca o entitate colectivă (specie sau limbă) ce se supune unor legi compoziționale (de încrucișare sau de comunicare) să evolueze natural spre o identitate "remarcabilă" (singura formă care ne interesează), e nevoie să fie respectate

<sup>1</sup> *European Chapter of the Association for Computational Linguistics (EACL)*

<sup>2</sup> *less-studied languages*

<sup>3</sup> *low density languages* (desigur densitatea aici referindu-se la

simultan: o dimensiune minimă, care să împiedice o proliferare a raporturilor incestuoase, dacă vorbim de specie, sau o plafonare a formei și lexicului datorată placidității imaginative a unei colectivități redusă numeric, dacă vorbim de limbă, dar, credem noi, și o dimensiune maximă, care e totuna cu o anumită izolare geografică sau comunicațională, pentru a împiedica o diluare a trăsăturilor genetice, dacă vorbim de specie, sau o impurificare lingvistică, dacă vorbim de o limbă naturală.

Forțarea dimensiunii maximale poate duce la pierderea identității unei limbi "năcăjite" prin agresiunea arogantă a limbilor tari, a celor ce fac ocolul globului în mediile electronice. În ziua întâi a globalizării, această perspectivă sumbră este cunoscută și guvernele statelor lumii devin din ce în ce mai conștiente de necesitatea apărării identității lingvistice autohtone. Globalizarea nu trebuie să conducă la o amestecătură lingvistică amorfă, un fel de *new-speak* planetar, ci dimpotrivă, la cultivarea limbilor proprii prin ușurarea folosirii lor în comunicarea în mediile electronice. Ori este evident că acest lucru nu poate fi realizat fără aportul tehnologiilor informatice aplicate limbajului natural.

## 2. Comisia și Consorțiul de Informatizare pentru Limba Română

În octombrie 2001, a fost creată în Academia Română, în cadrul Secției de Știința și Tehnologia Informației, Comisia de Informatizare pentru Limba Română (CILR). Scopul acestei comisii îl constituie apărarea identității limbii române prin promovarea studiilor dedicate ei dintr-o perspectivă informațională. Ulterior, la Adunarea Generală a CILR din martie 2002 s-a hotărât împărțirea activității Comisiei în patru sectoare: lingvistică teoretică și formală, prelucrarea limbajului natural, prelucrarea vorbirii și terminologie.

Rezultat al unor întâlniri și discuții desfășurate între lingviști și informaticieni ce au avut loc succesiv în București, Iași, și Chișinău, încă înainte de apariția Comisiei, se prefigura necesitatea creării unui Consorțiu care să asigure un cadru organizat de comunicare între cercetătorii lingviști și informaticieni și de depozitare de informații în domeniile de lucru promovate de Comisie. Momentul decisiv în crearea Consorțiului de Informatizare pentru Limba Română (ConsILR) poate fi considerat eliberarea unei finanțări, în primăvara anului 2002, din partea proiectului național Societatea Informațională – Societatea Cunoașterii, finanțare care a permis construirea unui sit electronic dedicat<sup>4</sup>.

Explozia de informații pe Internet, necesitatea tot mai mare de a accede la ele prin mijloace inteligente utilizând limba naturală proprie, indiferent de limba de depozitare a informațiilor, face astăzi ca domeniul Prelucrării Limbajului Natural să fie considerat prioritar în Comunitatea Europeană. Tot mai multă lume recunoaște că Tehnologia Limbajului va deveni domeniul primordial al Tehnologiei Informației în secolul nostru. Tot mai multe limbi sunt astăzi extrem de bine studiate și reprezentate în mediile electronice. Rămâne în urmă în acest domeniu înseamnă un handicap tehnologic cu urmări care nu se limitează numai asupra limbii proprii, ce poate fi grav afectată, dar poate avea și consecințe economice greu de recuperat. Se manifestă, ca urmare, o necesitate stringentă de a ridica limba română, în privința resurselor și a tipurilor de prelucrări automate, la nivelul altor limbi europene. Acesta este scopul principal al creării Consorțiului.

În al doilea rând, limba română dispune de fonduri de date lingvistice impresionante (dicționare explicative, tezaur, terminologice, bilingve etc.) care, din păcate, din cauza unei politici de protecție deficitară ori a neînțelegerii posibilităților de exploatare oferite de informatică, sunt puțin ori deloc cunoscute publicului larg ori cercetătorilor din lingvistică computațională. Alte tipuri de resurse în limba română, în special corpusuri lingvistice, în stare neprelucrată ori adnotate, există într-un număr insuficient încă în comparație cu alte limbi europene. De crearea lor depinde accesarea la tehnologia lingvistică modernă și este neîndoios că acest proces trebuie accelerat.

În al treilea rând, doar puțini cercetători lingviști sunt astăzi conștienți de avantajele enorme pe care le-ar putea avea în munca lor apelând la modalități de prelucrare automată a

---

<sup>4</sup> Adresa acestui sit este <http://consilr.info.uaic.ro>.

datelor lingvistice sau slujindu-se de calculator ca de un instrument de lucru. Conștientizarea în rândul lingviștilor a posibilităților de utilizare a tehnologiilor informatice în cercetarea lingvistică este încă un motiv în favoarea creării Consorțiului.

În al patrulea rând, puțini cercetători care lucrează în lingvistică computațională ori în prelucrarea limbajului natural ori a vorbirii, în România și Republica Moldova, și care au un fundament informatic, dispun de cunoștințe lingvistice suficiente care să le permită o integrare la un nivel înalt a cunoștințelor din cele două domenii. Atragerea lingviștilor de formație clasică în colective mixte de cercetare alături de informaticieni, cât și promovarea studiilor post-universitare, de masterat ori doctorat în domeniul lingvisticii computaționale și al prelucrării limbajului natural este un alt argument în favoarea creării Consorțiului.

Activitatea ConsILR este caracterizată, în principal, de următoarele trăsături:

- partenerii la Consorțiu se pot asocia pentru alcătuirea de proiecte adresate unor instituții finanțatoare naționale ori internaționale, în vederea solicitării de granturi de cercetare și mobilitate;
- Consorțiul se dorește să devină un forum al ideilor de cercetare teoretică și aplicativă asupra limbii române, inclusiv pentru dezvoltarea de resurse dedicate limbii române;
- Consorțiul se preocupă de includerea pe situl său, într-un format electronic care să permită exploatarea facilă, a acelor componente ce se consideră că au o valoare covârșitoare pentru limba română (dicționare, enciclopedii, corpusuri textuale sau fonetice, gramatici, terminologii etc.). Orientarea generală trebuie să fie una de transparență, militându-se pentru obținerea acordului de acces public asupra a cât mai multe componente considerate de interes național și internațional;
- rezultatele eforturilor de cercetare, prin consensul autorilor, pot fi depozitate (și) pe situl Internet al Consorțiului, autorii contribuțiilor urmând să-și păstreze neștirbite drepturile de autor. Membrii Consorțiului, în nume individual sau colectiv, pot utiliza componente aflate pe sit pentru a le exploata în scopuri didactice ori științifice, sau în scopuri ce se pot finaliza cu profit, după obținerea acordului autorilor și, dacă e cazul, reglementarea cu aceștia a manierei de participare la profit;
- Consorțiul își propune să identifice surse de finanțare (obținute din contracte de cercetare/dezvoltare, exploatarea componentelor de pe sit, cotizații etc.), pe care să le poată folosi pentru cercetare și dezvoltare, inclusiv pentru creșterea resurselor de calcul și depozitare electronică a datelor, sprijinirea participării membrilor săi la conferințe naționale și internaționale, organizarea de întruniri, școli de scurtă durată, cursuri de utilizare a calculatoarelor pentru prelucrări lingvistice, cursuri de masterat și doctorat în domeniul tehnologiei limbajului uman și al lingvisticii computaționale, reclama componentelor sitului în vederea valorificării lor, activități de editare etc.;
- Consorțiul va trebui să ajungă la confortul financiar care să facă posibilă editarea cel puțin a unui periodic științific propriu (inclusiv în format electronic), eventual înființarea unei edituri care să publice cercetările și realizările membrilor săi ca serii de cărți și CD-uri.

### **3. Resurse, instrumente, arhitecturi și standarde lingvistice și fonologice computaționale pentru limba română**

#### **3.1 Metodologie**

În versiunea finală, situl ConsILR va trebui să reflecte patru teme mari: **Resurse – R, Standarde – S, Instrumente – I și Arhitecturi – A**. Fiecare dintre aceste teme va fi divizată în subteme corespunzând componentei **Textuale – T și Fonologice – F**, rezultând subtemele:

**RT - Resurse Textuale,**

**RF - Resurse Fonologice,**

**ST - Standarde Textuale,**

**SF - Standarde Fonologice,**

**IT - Instrumente Textuale,**

**IF - Instrumente Fonologice,**

**AT - Arhitecturi Textuale,**

**AF - Arhitecturi Fonologice**

**AM - Arhitecturi Mixte** (aplicații care integrează atât componente de prelucrări textuale cât și a vorbirii).

Tema de **Resurse** urmărește colectarea unui eșantion semnificativ sub aspect lingvistic (textual și fonologic) al limbii române. Ea va ținti colectarea și/sau crearea de colecții de date scrise și vorbite care să se constituie într-un depozit informațional care să susțină aplicațiile de prelucrare a limbii scrise și vorbite, sau din care să se infereze colecții de reguli și modele statistice pentru limba română. Ca metodologie generală se recomandă un echilibru între utilizarea corpusurilor pentru inferarea colecțiilor de reguli, structuri sintactice ori modele ale limbii și introducerea lor de mână de către experți lingviști. Atragem atenția asupra faptului, în general, recunoscut că o inferare automată din corpusuri a regulilor este rapidă, inducând descrieri consistente, dar prea generale (ceea ce se poate traduce în acceptarea unor construcții agramaticale), în timp ce compunerea lor manuală este consumatoare de timp, poate provoca inconsistențe, poate duce la respingerea unor construcții uzitate dar încă neincluse în normă, pentru că sunt greu de găsit mecanisme adecvate de control, dar poate genera reguli mult mai compacte. Alinierea prin mijloace tradiționale a normei la uz este un proces lent și imperfect, care poate fi îmbunătățit numai prin mijloace automate.

Tema de **Standarde** urmărește codificarea formatelor intermediare în lanțul de prelucrări. Prin crearea acestui set de standarde se urmărește compatibilizarea în intrare/ieșire a modulelor de prelucrări textuale și fonologice. Formatul de codificare promovat este XML (<http://www.w3.org/XML/>), în cadrul căruia se pot adopta standarde internaționale (parțial) acceptate precum TEI (Barnard&Ide, 1997), XCES (Ide&Romary, 2000) Parole (Ruimy *et al.*, 1998) etc. Pentru formatele pentru care nu există standarde sau practici general acceptate, se are în vedere propunerea de noi scheme de adnotare și documentarea lor (sintaxă și semantică). În mod normal, detalierea tipurilor de codificări ar trebui să precedă elaborarea instrumentelor și a arhitecturilor, deși un proces de rafinare a standardelor ulterior elaborării modulelor componente ale aplicațiilor ar putea, de asemenea, fi acceptat în anumite situații. Ideea generală este ca, prin standardizare, diverse module ale aplicațiilor să se poată cupla în lanțuri alternative de prelucrări. Pe de altă parte, ori de câte ori o nouă schemă de codificare se propune, ea trebuie să fie însoțită de o descriere formală (sintaxă și semantică) care să permită conversia (semi-)automată și deterministă ulterioară din/spre scheme de adnotare viitoare sau compatibilizarea cu scheme existente.

Tema de **Instrumente** urmărește realizarea unor capacități de prelucrare textuală și fonologică a limbii române grupate într-un set de module software. Acestea pot fi realizate în totalitate în cadrul Consorțiului dar și preluate din surse externe (eventual adaptate). În principiu, instrumentele, componente software, trebuie să fie independente de limbă, utilizarea lor în tandem cu resurse românești atribuindu-le funcționalitatea specifică prelucrării limbii române.

Tema de **Arhitecturi** urmărește realizarea unui sistem inteligent, capabil să asiste construirea de aplicații complexe de prelucrare a limbii române în formă scrisă, vorbită sau mixtă. Strategia de dezvoltare trebuie să sprijine utilizatorul (informat) în procesul de construire a unei aplicații pentru limba română prin combinații de module din sfera prelucrării textuale/fonologice (independente de limbă), oferită de tema de Instrumente. Alinierea la standarde atât a realizării instrumentelor cât și a resurselor este o premisă esențială pentru realizarea efectivă a obiectivelor Consorțiului.

Vom detalia în cele ce urmează câteva componente din domeniul textual.

### 3.2 Resurse Textuale (subtema RT)

Se urmărește crearea următoarelor resurse textuale pentru limba română:

- **RT-CorpusText**: o colecție structurată de texte românești în format electronic, independent de platforma hardware sau software de prelucrare, acoperind o arie largă de registre lingvistice (ziaristic, beletristic, științific, e-mail etc.), arii geografice (Ardeal, Moldova, Banat, Moldova românească, Basarabia etc.), origini în timp (secolul trecut, limba dintre cele două războaie, limba actuală) ori origini sociale (livresc, limbajul cult, limbajul suburban ori de argou etc.). Textele incluse în această colecție trebuie să conțină obligatoriu informație de identificare a criteriilor de clasificare (sursă, autor, registru, arie geografică, perioadă etc.). Pe baza acestor criterii, explicit codificate, pot fi extrase oricând corpusuri speciale (tematice) sau balansate (de referință), corpusuri diacronice (monitor) etc. Dimensiunea unei astfel de colecții, pentru a fi reprezentativă, trebuie să fie cât mai mare, conținând cel puțin 500 milioane de cuvinte ocurență. O tipologie a corpusurilor general acceptată de comunitatea științifică din domeniul prelucrării limbajului natural este prezentată în (Sinclair, Ball, 1995) și comentată în (Teubert, 1997).
- **RT-CorpusAnnotXXX**<sup>5</sup>: o colecție de corpusuri românești, adnotate în conformitate cu diverse scheme. Adnotarea trebuie realizată în format XML (elementar și din ce în ce mai complex, în diverse formate – v. și TEI/CES). Adnotarea se va realiza fie prin proceduri semi-automate (asistată de calculator), fie complet automată. Tipul abordării depinde de natura și granularitatea informației adăugată textelor (care în fond definesc dificultatea procesului de adnotare). În acest proces vor putea fi folosite instrumente din categoria **IT-Annot** (v. subtema **IT**). Diversele componente ale acestor corpusuri vor constitui sursele principale de antrenare a componentelor temei **IL**. Este recomandabil ca diferitele componente ale resursei **RT-CorpusAnnotXXX** să utilizeze într-o măsură cât mai mare aceleași texte de origine, pentru a facilita în acest fel coeziunea unor lanțuri de antrenare a componentelor din **IT**. Din punct de vedere tehnic, este de asemenea de dorit ca adnotările la diferite niveluri, aplicate aceluiași text, să poată fi separate oricând prin filtre, lucru ce poate fi realizat ușor prin transpunerea adnotărilor în baze de date. Codificările folosite în **RT-CorpusAnnot** trebuie să fie în concordanță cu standardele descrise în secțiunea **ST** a proiectului.
  - **RT-CorpusAnnotMSD**: corpusuri adnotate la nivel morfo-lexical;
  - **RT-CorpusAnnotNP**: corpusuri adnotate la nivelul grupurilor nominale;
  - **RT-CorpusAdnotSyn**: corpusuri adnotate la nivel sintactic, având specificate atât structurile sintactice cât și dependențele dintre ele; astfel de corpusuri adnotate se mai numesc și *bănci de arbori (tree-banks)*;
  - **RT-CorpusAnnotArg**: corpusuri adnotate la nivel semantic având specificate structurile de tip predicat-argument, cu identificarea statutului (opțional sau obligatoriu) și a rolului fiecărui argument (agent, beneficiar, loc, mod etc); astfel de corpusuri se mai numesc și *bănci propoziționale (prop-banks)*, propoziția fiind considerată în accepțiunea ei din logica predicatelor: o expresie logică de predicate *n*-are, fiecare predicat fiind descris de o schemă semantică precizând atât numărul și tipul argumentelor sale cât și numele lor;
  - **RT-CorpusAnnotSeg**: corpusuri adnotate la segmente (propoziții în cadrul frazelor, unități de discurs, adesea confundate);
  - **RT-CorpusAnnotAna**: corpusuri adnotate la relații anaforice: coreferința (echivalența referențială), referințe meronimice explicite (parte-întreg, obiect-

<sup>5</sup> Sufixul **XXX** adăugat numelui unei componente, aici cât și mai departe, semnifică posibilitatea de împărțire a unei componente în mai multe subcomponente, din aceeași clasă.

proprietate, posesor-obiect posedat etc.) sau implicite (*bridge anaphora*) (ce pot fi de aceleași tipuri ca și cele explicite);

- **RT-CorpusAnnotWSD**: corpusuri adnotate la sensuri ale cuvintelor, ce pot fi folosite pentru inferarea regulilor de identificare a sensurilor cuvintelor (*word sense disambiguation*);
- **RT-CorpusAnnotDS**: corpusuri adnotate la structură de discurs, de exemplu *structura relațiilor retorice* (*Rhetorical Structure Theory - RST*).

Adesea corpusurile adnotate cuprind niveluri diferite de adnotare. De exemplu, un corpus din categoria **RT-CorpusAnnotDS** trebuie să cuprindă, cel puțin, următoarele niveluri de adnotare: morfo-lexicală, grupuri nominale, unități de discurs și relații retorice. Clasificarea noastră trebuie înțeleasă ca identificând o listă a elementelor specifice de adnotare, încât un corpus adnotat anumit să poată aparține simultan mai multor tipuri.

- **RT-CorpusParallel**: un corpus de texte paralele (română – o limbă străină) aliniat. Alinierea poate fi realizată la nivel de paragraf, propoziție, sau chiar cuvânt. Un astfel de corpus este extrem de util pentru antrenarea componentelor de traducere automată (v. **IT-TradXXX**).
- **RT-RegSegm**: o colecție de reguli/șabloane pentru segmentarea textelor românești. Segmentarea poate fi făcută la diverse niveluri, incluzând un nivel morfologic, un nivel sintactic și un nivel al discursului. O componentă a acestei resurse va putea fi integrată instrumentului **IT-Tokenizer** (v. subtema **IT**) ce urmărește recunoașterea articolului lexical de prelucrare primară; noțiunea de articol lexical depinde de natura aplicațiilor ce folosesc prelucrarea limbajului natural și, în acest sens, specificarea regulilor de segmentare trebuie să fie parametrizabilă în raport cu finețea și profunzimea prelucrărilor. Regulile de segmentare (cel mai adesea exprimate prin *expresii regulate*) permit recunoașterea abrevierilor, datelor calendaristice, numerelor, valorilor financiare, formulărilor, tabelor etc. De asemenea, pe baza unor resurse lexicale specifice, pot fi descrise structuri multi-cuvânt ce urmează să fie tratate ca un singur articol lexical (*de+jur+împrejur => de\_jur\_împrejur*, *de-a+dreptul => de-a\_dreptul*, *verde+de+Paris => verde\_de\_Paris*) sau dimpotrivă, forme cliticizate ce trebuie interpretate ca două sau mai multe articole lexicale (*dându-mi-le => dând+mi+le*). Un segmentator mai special este cel ce recunoaște numele proprii și în plus le atribuie o etichetă semantică dintr-o mulțime specificată apriori (persoană, instituție, loc etc.)<sup>6</sup>. Alte colecții de reguli/expresii regulate de segmentare urmăresc identificarea grupurilor nominale nerecursive, a propozițiilor componente ale frazelor sau a unităților de discurs<sup>7</sup>. Drept surse pentru această resursă se pot folosi: dicționare de abrevieri, dicționare de cuvinte compuse, lista cliticelor limbii române, liste de nume proprii (*gazetteers*) etc.
- **RT-DictMorphLex**: dicționar morfo-lexical conținând colecția cuvintelor flexionate românești cu descrieri morfo-lexicale atașate incluzând obligatoriu lema, ca structuri de caracteristici nerecursive sau codificate pe câmpuri (*Morpho-Sintactic Description - MSD*). Identificarea unui model de codificare standardizată, din punct de vedere morfo-lexical, aplicabil cât mai multor limbi a fost unul din obiectivele mai multor proiecte europene (EAGLES, MULTEXT, MULTEXT-EAST) pentru limba română, o instanțiere a acestui model fiind descrisă în (Tufiş et al., 1997).
- **RT-MorphPara**: o descriere paradigmatică a morfologiei românești, care să cuprindă colecția completă a paradigmelor flexionare românești (terminații, rădăcini de cuvinte și asocieri valide rădăcini-terminații). Un instrument de analiză care ar lucra cu această

---

<sup>6</sup> Engl. *name entity recognizer*.

<sup>7</sup> Secvențe lexicale ce se constituie în unități atomice ale structurii de discurs, cel mai adesea propoziții de sine stătătoare sau propoziții constitutive ale frazelor. Uneori însă o unitate de discurs poate fi mai largă decât strict o propoziție din constituția unei fraze.

componentă ar trebuie să fie capabil să „spargă” un cuvânt flexionat în rădăcină și terminație prin încadrarea lui într-o paradigmă cunoscută. În cazul cuvintelor necunoscute, componenta ar trebui să „ghicească” categoria morfo-lexicală cea mai plauzibilă. Alegerea celei mai plauzibile interpretări este o problemă contextuală, de multe ori rezolvabilă printr-o analiză strict distribuțională (v. și **RT-Guesser**).

- **RT-Guesser**: o colecție de reguli de inferență (statistică) a celei mai probabile interpretări morfo-lexicale a unui cuvânt necunoscut.
- **RT-LangModel**: un model al limbii române scrise, adică o rețea de probabilități a grupărilor (bigrame, trigrame) de categorii morfo-sintactice, obținută prin antrenare pe un corpus. **RT-LangModel** poate avea realizări specifice, în funcție de registrul lingvistic, autor etc.
- **RT-WordNet**: un *wordnet* românesc. Într-un *wordnet* (Miller *et al.*, 1990), seriile sinonimice (*sinset-uri*), caracterizate de un înțeles comun, sunt noduri ale unei rețele semantice. Relațiile ce corelează nodurile rețelei sunt de natură lexicală sau semantică. Nodurile mai multor rețele semantice monolingve pot fi puse (în mod independent) în relații de echivalență conceptuală cu reprezentări independente de limbă. Folosind proprietățile relațiilor de echivalență, se pot identifica, pornind de la reprezentările conceptuale, lexicalizările acestora în diverse limbi. Realizarea *wordnet*-ului pentru limba română, aliniat la nucleul comun european, constituie un obiectiv însemnat pentru abordarea traducerii automate având limba română la un pol, cât și pentru generarea de dicționare electronice bilingve spre/dinspre limbi pentru care astfel de resurse lipsesc (Tufiş, 2002).
- **RT-DictColocat**: un dicționar de colocații pentru limba română. El trebuie să cuprindă, pentru fiecare sens al unui verb sau substantiv derivat din verb, exemple de colocatori frecvenți sau minimali.
- **RT-DictSubcat**: un dicționar de subcategorizare al limbii române (pentru limba engleză v. Beth Levin). Acest dicționar trebuie să cuprindă (în varianta sa completă) pentru fiecare sens al fiecărui verb (inclusiv derivatele sale nominale) din limba română, structura sa de subcategorizare (argumentele obligatorii, cu restricțiile sintactice și semantice). Construcția dicționarului se poate automatiza într-o mare măsură în condițiile existenței unei corpus reprezentativ adnotat la nivel sintactic și semantic (*tree-bank*, *prop-bank*).
- **RT-DictParalelXXX**: un număr de dicționare paralele pot fi realizate (semi-)automat din momentul în care va fi finalizată componenta **RT-WordNet**, aliniată cu alte limbi.

### 3.3 Standarde Textuale (subtema ST)

Principalele standarde aplicate resurselor textuale pe care le avem în vedere sunt:

- **ST-Txt**: este notația noastră pentru text neadnotat.
- **ST-CorpusAnnotXXX**: o colecție de standarde de adnotare a corpusurilor ca și a textelor prelucrate automat. Includerea în aceeași categorie a acestor două clase de documente este voită și exprimă tendința de uniformizare a codificărilor atașate corpusurilor (de obicei, interactiv, de către experți, prin utilizarea instrumentelor de adnotare) cu cele generate automat de mașină, în urma proceselor de analiză. Trebuie avută în vedere posibilitatea ca mai multe adnotări, cu destinații diferite, să fie combinate pe același corpus sau text, ceea ce impune necesitatea ca un standard să reflecte doar o soluție elementară de adnotare, iar nu o combinație de codificări. O astfel de concepție ar permite compunerea liberă a adnotărilor elementare, dacă restricțiile limbajului XML sunt respectate (de exemplu neintersectarea etichetelor) - v. (Cristea *et al.*, 1998a), (Cristea *et al.*, 1998b). Din acest motiv, unei componente din **RT-CorpusAdnotXXX** îi pot

corespunde una sau mai multe componente din **ST-CorpusAnnotXXX**. Exemple de componente din această categorie sunt:

**ST-CorpusAnnotTok**: formatul de după **IT-Tokenizer** (v. secțiunea **IT**). El marchează cu etichete distincte: cuvinte, clitice, semne de punctuație, date calendaristice, formule, tabele etc.

**ST-CorpusAnnotMSDXXX**: formate pentru descrieri morfo-sintactice, în care fiecărui element lexical (*token*) i se atașează una sau mai multe interpretări morfo-sintactice. Aceste formate pot cuprinde și marcaje pentru cuvinte necunoscute (necuprinse în dicționar). Variante, pot fi:

**ST-CorpusAnnotMSDk**: una sau mai multe interpretări pentru fiecare *token*, posibil și *token*-uri nemarcate;

**ST-CorpusAnnotMSDk+**: una sau mai multe interpretări pentru fiecare *token*, absolut toate *token*-urile marcate;

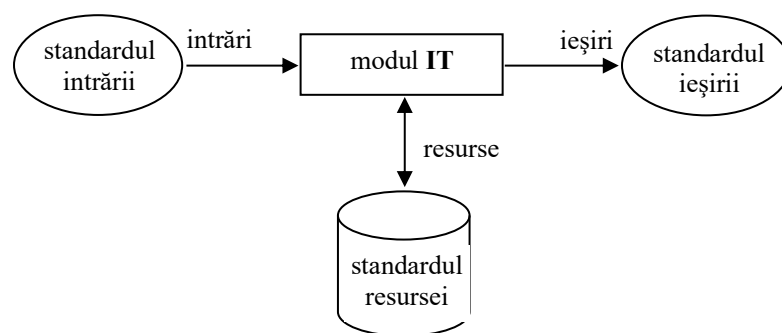
**ST-CorpusAnnotMSD1+**: o unică interpretare pentru fiecare *token*, toate *token*-urile marcate.

**ST-CorpusAnnotNP**: standard de codificate a grupurilor nominale de recursivitate limitată (în sensul că nu conțin subgrupuri verbale; de exemplu în *Omul pe care l-am văzut conducând mașina este directorul firmei* se notează ca grup nominal *omul* iar nu *omul pe care l-am văzut conducând mașina*).

- **ST-RegSegmXXX**: format de codificare a regulilor de segmentare la diferite niveluri.
- **ST-DictSubcat**: formatul de codificare a dicționarelor de subcategorizare.
- **ST-WordnetXXX**: formatul de codificare a diferitelor componente ale unei baze de date de tip *wordnet*.

### 3.4 Instrumente pentru prelucrări textuale (subtema IT)

Modulele temei **I** sînt componente software a căror caracteristică esențială este independența de limbă. Aceasta înseamnă că aplicarea lor la limba română trebuie să se facă prin utilizarea de resurse specifice, componente ale temei **RT**. Fiecare modul realizează o funcție dată, fiind cuplat într-un sistem prin intermediul a trei tipuri de porți: intrări, ieșiri și resurse (v. Figura 1).



**Figura 1:** Interfața unui modul cu sistemul

Pentru simplificarea notării interfețelor modulelor temei **I**, vom utiliza în cele ce urmează o convenție de desemnare a acestora (semnături) de forma:

<nume modul>(I:<listă standarde intrări>, R:<listă resurse>, O:<listă standarde ieșiri>).

Intrările și ieșirile vor fi compatibilizate astfel prin componentele temei de standarde – S.



Dăm în continuare a listă a câtorva dintre cele mai semnificative tipuri de instrumente de prelucrări lingvistice:

- **IT-Tokenizer**: modul de depistare a granițelor dintre cuvinte, de recunoaștere a formulelor, tabelor, cliticilor, abrevierilor, datelor calendaristice, numerelor, valorilor financiare, compușilor și a numelor proprii. Acest modul primește în intrare un text sursă (format **ST-Txt**), accesează **RT-RegSegm** și produce în ieșire un format **ST-Tok**, sau, în convenția noastră:

IT-Tokenizer(I:ST-Txt, R:RT-RegSegm, O:ST-Tok).

- **IT-LookUpMSD**: modul de depistare a interpretărilor morfo-lexicale ale cuvintelor prin căutare în dicționar. Semnătura propusă este:

IT-LookUpMSD(I:ST-Tok, R:RT-DictMorphLex, O:ST-CorpusAnnotMSDk).

- **IT-LookUpPar**: modul de depistare a interpretărilor morfo-lexicale ale cuvintelor prin analiză paradigmatică. Rezultatul trebuie să fie analog lui **IT-LookUpMSD**:

IT-LookUpPar(I:ST-CorpusAnnotTok, R:RT-MorphPara,  
O:ST-CorpusAnnotMSDk).

- **IT-Guesser**: modul de tratare a cuvintelor necunoscute. Trebuie să furnizeze o clasă de interpretări posibile pentru un cuvânt necunoscut, important aici fiind să nu se piardă interpretări. Este de presupus că, la ieșirea din **IT-Guesser**, formatul **ST-CorpusAnnotMSD** nu mai conține *token*-uri neetichetate. În subtema S acest format este notat **ST-CorpusAnnotMSDk+**. Semnătura propusă:

IT-Guesser(I:ST-CorpusAnnotMSDk, R: RT-Guesser, O:ST-CorpusAnnotMSDk+).

- **IT-TaggerXXX**: un set de module de dezambiguizare morfo-sintactică. O componentă de acest tip presupune existența unui model al limbii **RT-LangModel** și a unei funcții de optimizare (un criteriu de decizie). Se pot adopta mai multe soluții: soluția „lacomă”(greedy) **IT-TaggerGre** – în care se atribuie o unică interpretare pentru fiecare *token*; soluția celor mai bune *k* etichete (*k-best tagging*) **IT-TaggerkBT** – în care se acceptă cel mult *k* interpretări atunci când probabilitățile etichetelor posibil de atribuit unui token nu sunt suficient departajate, maniera etichetării pe niveluri (*tiered-tagging*) (Tufiș, 1999), algoritmul Welch Baum (Baum, 1972), algoritmul Brill (Brill, 1994) etc. Câteva semnături posibile sunt:

IL-TaggerGre(I:ST-CorpusAnnotMSDk+, R:RT-LangModel,  
O:ST-CorpusAnnotMSD1+)

IL-TaggerkBT(I:ST-CorpusAnnotMSDk+, R:RT-LangModel,  
O:ST-CorpusAnnotMSDk+).

- **IT-ChunkerXXX**: module de depistare a grupurilor nominale, prepoziționale etc. utilizând reguli de parsare de suprafață (*shallow-parsing*) (Ait-Mokhtar&Chanod, 1997). O semnătură posibilă este:

IT-NPChunkerNP(I:ST-CorpusAnnotMSD1+, R: RT-RegSegm,  
O:ST-CorpusAnnotNP).

- **IT-ChuRulesInducerXXX**: module de inferențiere a regulilor de depistare a grupurilor simple (nerecursive) din corpus. Aceste module prelucrează componente ale corpusurilor adnotat la grupurile respective pentru a genera resurse de folosit de către instrumente de parsare. Un exemplu de modul din această clasă îl reprezintă un inferențiator de reguli de depistare a grupurilor nominale:

IT-ChuRulesInducerNP(I:ST-CorpusAnnotNP, R:∅, O:ST-RegSegmNP).

- **IT-ConcordXXX**: instrumente de găsimă a concordanțelor lexicale și gramaticale într-un corpus. Lexicograful (pe post de utilizator) poate dispune de facilități de alegere a criteriilor de selecție atât pentru elementul central căutat cât și a celor folosite ca filtre

adiționale (contexte ale elementului central). Un concordanțier lexical (**IT-ConcordLex**) lucrează pe un corpus textual (**RT-CorpusText**) pe când unul gramatical (**IT-ConcordGra**) – pe un corpus adnotat la componente sintactice (o componentă a **RT-CorpusAnnotXXX**). Cu instrumente din această clasă se pot forma dicționare de colocații (**RT-DictColocat**) și de subcategorizare (**RT-DictSubcat**). Un concordanțier gramatical ar trebui, în principiu, să permită lexicografului consultarea interactivă a caracteristicilor diverselor elemente de context. Semnături propuse:

ITConcordLex(I:ST-Txt, R:Ø, O:{Ø, ST-DictColocat})

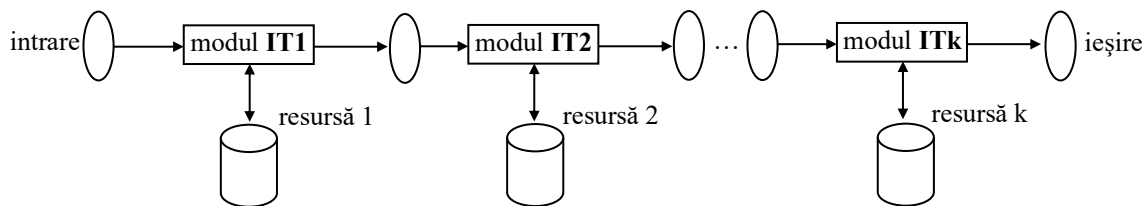
ITConcordGra(I:ST-CorpusAnnotXXX, R:Ø, O:{Ø, ST-DictSubcat}).

Ieșirile vide ce apar în semnăturile de mai sus se referă la utilizarea programelor concordanțiere pentru obținerea de liste.

- **IT-StatistXXX**: pachet de instrumente pentru generarea de estimări statistice pe corpusuri. Exemple de astfel de estimări sunt informațiile mutuale (astfel, dacă frecvența de co-apariții a unui verb împreună cu un grup nominal depășește un anumit prag se poate genera concluzia unei utilizări tranzitive).
- **IT-TradXXX**: pachet de instrumente pentru ajutorarea traducerii în/din limba română.
- **IT-SpellXXX**: pachet de instrumente pentru corectarea ortografiei românești.
- **IT-DiscourseXXX**: pachet de instrumente pentru prelucrări ale discursului ș.a.m.d.

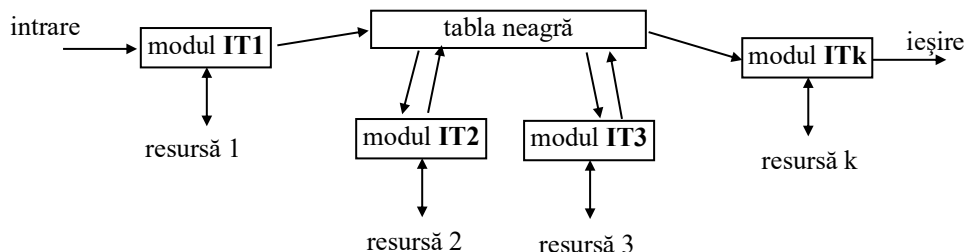
### 3.5 Arhitecturi pentru prelucrări textuale (subtema AT)

Maniera de conectare a unui modul prin cele trei tipuri de porți (v. tema I) permite realizarea de interconexiuni de prelucrare complexe. Unul dintre acestea este cel în cascadă (*pipe-line*), v. Figura 2.



**Figura 2:** Arhitectură în cascadă (ovalurile reprezintă standarde)

Aceași schemă de proiectare individuală a modulelor poate asambla, de asemenea, arhitecturi tablă neagră (*blackboard*), prin utilizarea unor ieșiri specializate să scrie pe o resursă temporară comună, care este tabla (v. Figura 3).

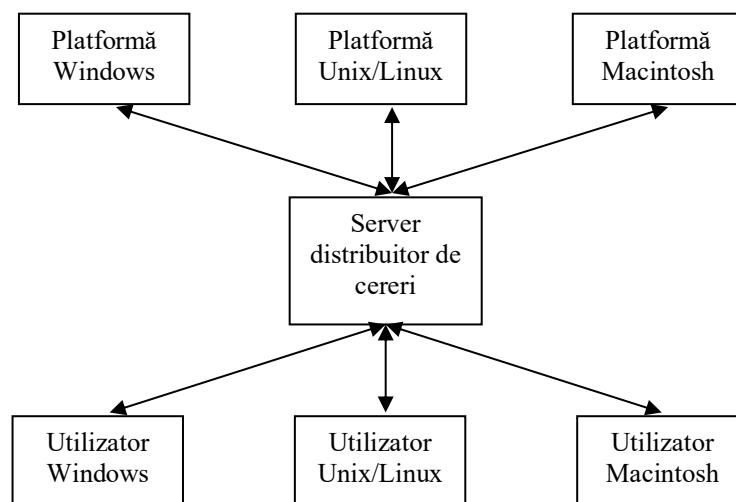


**Figura 3:** O arhitectură tablă neagră

Pentru realizarea de configurații mobile, este recomandabil ca modulele să poată lucra în două maniere: intrare completă (sincron) sau intrare incrementală (asincron). În varianta cu intrare completă, tot textul este consumat de primul modul, înainte de ca ieșirea acestuia să

constituie intrare pentru al doilea modul ș.a.m.d. (v., de exemplu, arhitectura GATE – (Tablan *et al.*, 2002)). Într-o variantă incrementală, intrarea ar putea fi fragmentată, un modul trebuind să rezolve doar o parte din intrare. Acest tip de prelucrare poate fi utilă în conceperea arhitecturilor *tablă-neagră*.

Problema asigurării portabilității aplicațiilor este deosebit de importantă. Neglijarea ei poate transforma o aplicație de succes într-un eșec comercial prin imposibilitatea de a se adresa unei mase de utilizatori ce nu au acces la platforma de calcul pentru care a fost ea construită. Pentru asigurarea portabilității aplicațiilor, recomandăm asigurarea accesului la distanță, pentru execuție de pe orice platformă utilizator, prin intermediul unor interfețe de compatibilizare (v. Figura 4).



**Figura 4:** Arhitectură de comutare, transparentă, pentru asigurarea portabilității aplicațiilor

Zona aplicativă a sitului va trebui să permită și interogări *on-line*, inclusiv configurarea dinamică într-un timp scurt a unor aplicații de prelucrare a textelor sau vorbirii pentru limba română.

## Referințe bibliografice

- Ait-Mokhtar, S., Chanod, J.-P. (1997). Incremental Finite-State Parsing, *Proceedings of the Fifth Conference on Applied Natural Language Processing*.
- Barnard, D. and Ide, N. (1997). The Text Encoding Initiative: Flexible and Extensible Document Encoding. *Journal of the American Society for Information Science*.
- Baum L.E. (1972). An Inequality and Associated Maximization Technique in Statistical Estimation of Probabilistic Functions of Markov Processes, *Inequalities*, 3, 1-8.
- Brill, E. (1994). Some Advances In Rule-Based Part of Speech Tagging. *Proceedings of the AAAI*.
- Cristea, D.; Ide, N.; Romary, L. (1998a): Marking-up multiple views of a Text: Discourse and Reference. *Proceedings of the First International Conference on Language Resources and Evaluation*, Granada, Spain.
- Cristea, D.; Crăciun, O.; Ursu, C. (1998b): GLOSS-A Visual Interactive Tools for Discourse Annotation. *Proceedings of the Workshop on Annotation Tools, ESSLLI'98*, Saarbruecken, Germany.
- Ide, N., Bonhomme, P., Romary, L. (2000). XCES: An XML-based Standard for Linguistic Corpora. *Proceedings of the Second Language Resources and Evaluation Conference (LREC)*, Athens, Greece.

- Miller G.A., Beckwith R., Fellbaum C., Gross D., Miller K.J. (1990) Introduction to WordNet: An On-Line Lexical Database, *International Journal of Lexicography*, Vol. 3, No. 4 (winter 1990).
- Ruimy, N.; Corazzari, O.; Gola, E.; Spanu, A.; Calzolari, N.; Zampolli, A. (1998). LE-PAROLE project: The Italian Syntactic Lexicon, *Proceedings of EURALEX'98*, Université de Liège, Proceedings Volume I.
- Sinclair, J., Ball, J. (1995). Text typology (*External Criteria*). Eagles Report
- Tablan, V.; Ursu, C.; Bontcheva, K.; Cunningham, H.; Maynard, D.; Hamza, O.; McEnery, T.; Baker, P.; Leisher, M. (2002). A Unicode-based Environment for Creation and use of Language Resources. *Proceedings of the Third International Conference on Language Resources and Evaluation*, Las Palmas, Spain.
- Teubert W. (1997) Language Resources for Language Technology. În Tufiş D., Andersen P. *Recent Advances in Romanian Language Technology*, Romanian Academy Publishing House.
- Tufiş D., Barbu A.M., Pătraşcu V., Rotariu G., Popescu C. (1997). Corpora and Corpus-Based Morpho-Lexical Processing. În Tufiş D., P. Andersen (eds.) *Recent Advances in Romanian Language Technology*, Editura Academiei, Bucureşti.
- Tufiş, D. (1999). Tiered Tagging and Combined Classifiers. În F. Jelinek, E. Nöth (eds) *Text, Speech and Dialogue*, Lecture Notes in Artificial Intelligence 1692, Springer.
- Tufiş D., (2002). BALKANET - Tezaur lingvistic multilingv pentru limbile din Balcani, în acest volum.