

**METANET4U** 

**D2.3.ro**  
**Language Report for**  
**Romanian**  
**(Romanian version)**

Version 1.2

2011-07-31



# METANET4U

[www.metanet4u.eu](http://www.metanet4u.eu)

The central objective of the Metanet4u project is to contribute to the establishment of a pan-European digital platform that makes available language resources and services, encompassing both datasets and software tools, for speech and language processing, and supports a new generation of exchange facilities for them.

This central objective is articulated in terms of the following main goals:

**Assessment:** to collect, organize and disseminate information that permits an updated insight into the current status and the potential of language related activities, for each of the national and/or language communities represented in the project. This includes organizing and providing a description of: language usage and its economic dimensions; language technologies and resources, products and services; main actors in different areas, including research, industry, government and society in general; public policies and programs; prevailing standards and practices; current level of development, main drivers and roadblocks; etc.

**Collection:** to assemble and prepare language resources for distribution. This includes collecting languages resources; documenting these language resources; upgrading them to agreed standards and guidelines; linking and cross-lingual aligning them where appropriate.

**Distribution:** to distribute the assembled language resources through exchange facilities that can be used by language researchers, developers and professionals. This includes collaboration with other projects and, where useful, with other relevant multi-national forums or activities. It also includes helping to build and operate broad inter-connected repositories and exchange facilities.

**Dissemination:** to mobilize national and regional actors, public bodies and funding agencies by raising awareness with respect to the activities and results of the project, in particular, and of the whole area of language resources and technology, in general.

METANET4U is a project in the META-NET Network of Excellence, a cluster of projects aiming at fostering the mission of META. META is the Multilingual Europe Technology Alliance, dedicated to building the technological foundations of a multilingual European information society.



## Deliverable D2.3.ro: Language Report for Romanian (Romanian version)

METANET4U is co-funded by the participating institutions and the ICT Policy Support Programme of the European Commission



and by the participating institutions:



Faculty of Sciences, University of Lisbon



Instituto Superior Técnico



University of Manchester



University *Alexandru Ioan Cuza*



Research Institute for Artificial Intelligence,  
Romanian Academy



University of Malta



Technical University of Catalonia



Universitat Pompeu Fabra

Revision History

Version	Date	Author	Organisation	Description
1.0	03-07-2011	Diana Trandabăț, Elena Irimia, Verginica Mititelu- Barbu, Dan Tufiș, Dan Cristea	UAIC, RACAI	Draft version
1.1	15-07-2011	Diana Trandabăț, Elena Irimia, Verginica Mititelu- Barbu, Dan Tufiș, Dan Cristea	UAIC, RACAI	Pre-final version
1.2	31-07-2011	Diana Trandabăț, Elena Irimia, Verginica Mititelu- Barbu, Dan Tufiș, Dan Cristea	UAIC, RACAI	Final version

Statement of originality:

This deliverable contains original unpublished work except where clearly indicated otherwise. Acknowledgement of previously published material and of the work of others has been made through appropriate citation, quotation or both.



**METANET4U**

**D2.3.ro**  
**Language Report for**  
**Romanian**  
**(Romanian version)**

Document METANET4U-2011-D2.3.ro  
EC CIP project #270893

Deliverable  
Number: D2.3.ro  
Completion: Final  
Status: Submitted  
Dissemination level: Public

Responsible: Asunción Moreno (WP2 coordinator)

Contributing Partners: UAIC, RACAI

Authors: Diana Trandabăț, Elena Irimia, Verginica Mititelu-Barbu,  
Dan Tufiș, Dan Cristea

Reviewer: Amália Mendes, Adrian Iftene, Ionuț Pistol

© all rights reserved by FCUL on behalf of METANTE4U



## Cuprins

<b>Rezumat</b> .....	<b>7</b>
<b>Un risc pentru limbile noastre și o provocare pentru tehnologia limbajului</b> .....	<b>9</b>
Frontierele lingvistice frânează crearea unei societăți informaționale europene.....	9
Limbile noastre sunt în pericol.....	10
Tehnologia limbajului este cheia activării tehnologiei .....	11
Oportunități ale tehnologiei limbajului.....	11
Provocările tehnologiei limbajului .....	12
Achiziția limbii .....	12
<b>Limba română în societatea informațională europeană</b> .....	<b>14</b>
Fapte generale .....	14
Particularitățile limbii române.....	14
Dezvoltări Recente .....	16
Cultivarea limbii în România.....	17
Limba în educație .....	18
Aspecte Internaționale.....	18
Limba română pe Internet .....	19
Bibliografie selectivă .....	19
<b>Suport tehnologic pentru limba română</b> .....	<b>20</b>
Tehnologiile limbajului .....	20
Arhitecturile aplicațiilor din tehnologia limbajului .....	20
Principalele domenii de aplicații .....	21
<i>Corector de limbă</i> .....	21
<i>Căutarea pe Web</i> .....	23
<i>Interacțiunea vocală</i> .....	24
<i>Traducerea automată</i> .....	26
Tehnologiile limbajului .....	28
Tehnologiile limbajului în educație .....	31
Industria TL și programe.....	31
Cercetarea în domeniul TL și educația.....	32
Situația instrumentelor și resurselor pentru limba română .....	33
Tabel al instrumentelor și resurselor pentru limba română .....	34
Concluzii .....	36
<b>Despre META-NET</b> .....	<b>38</b>
Linii de acțiune.....	38
Organizații membre .....	40
<b>References</b> .....	<b>43</b>





## Rezumat

Multe limbi europene riscă să devină victimele erei digitale, fiind sub-reprezentate și dispunând de insuficiente resurse online. Oportunități uriașe de piață regională rămân neutilizate astăzi datorită barierelor de limbă. Dacă nu acționăm acum, vorbirea propriei limbi native va deveni un dezavantaj economic și social pentru mulți dintre cetățenii europeni.

Inovativă, Tehnologia Limbajului (TL) este un intermediar care va permite tuturor cetățenilor europeni să participe la o societate a cunoașterii și informației egalitară, inclusivă și de succes economic. Tehnologia Limbajului va fi poarta către comunicare și interacțiune dincolo de granițele limbii, de o manieră instantanee, ieftină și fără efort.

Astăzi, serviciile lingvistice sunt oferite în general de către furnizori comerciali din Statele Unite. Google Translate, un serviciu gratuit de traducere automată, este doar un exemplu. Succesul recent al lui Watson, un sistem IBM care a câștigat un episod din spectacolul-joc TV Jeopardy împotriva unor candidați umani, ilustrează imensul potențial al tehnologiei limbajului. Ca europeni, trebuie să ne punem câteva întrebări urgente:

- ❑ Infrastructura noastră de comunicare și cunoaștere trebuie să fie dependentă de companii monopoliste?
- ❑ Ne putem cu adevărat baza pe servicii lingvistice care pot fi foarte ușor întrerupte de către alții?
- ❑ Concurăm activ pe piața globală a cercetării și dezvoltării în domeniul tehnologiei limbajului?
- ❑ Sunt părțile terțe de pe alte continente dispuse să adreseze problemele noastre de traducere precum și alte subiecte legate de multilingvismul European?
- ❑ Este moștenirea culturală europeană capabilă de a modela societatea cunoașterii oferindu-i tehnologie de înaltă calitate, mai bună, mai sigură, mai precisă, mai inovativă și mai robustă?

Acest raport asupra limbii române demonstrează că cercetarea în universități și în mediul academic a reușit să proiecteze software specific de înaltă calitate precum și modele și teorii aplicabile pe scară largă. Totuși, pentru ca dezvoltarea TL să continue, o implicare mai intensă a Guvernului prin finanțare adecvată ar trebui obținută, iar colaborările atractive cu industriile care folosesc sau furnizează TL ar trebui promovate.

META-NET contribuie la construirea unui spațiu digital informațional European puternic și multilingv. Prin realizarea acestui deziderat, o uniune multiculturală a națiunilor poate prospera și deveni un model pentru cooperarea internațională pașnică și egalitaristă. Dacă acest țel nu poate fi atins, Europa va trebui să aleagă între a-și sacrifica identitățile culturale sau a suferi o înfrângere economică.



## Un risc pentru limbile noastre și o provocare pentru tehnologia limbajului

Suntem martorii unei revoluții digitale care are un impact dramatic asupra comunicării și societății. Dezvoltările recente din tehnologia comunicațiilor digitizate în rețea sunt uneori comparate cu invenția de către Gutenberg a tiparului. Ce ne poate spune această analogie despre viitorul societății informaționale europene și, în special, despre viitorul limbilor noastre?

Ulterior invenției lui Gutenberg au avut loc progrese reale în comunicare și schimbul de informație, datorită unor eforturi precum traducerea de către Luther a Bibliei din Latină în limba enoriașilor. În secolele următoare, au fost dezvoltate tehnici culturale pentru a îmbunătăți prelucrarea limbajului și schimbul de cunoștințe:

- Standardizarea ortografică și gramaticală a limbilor importante a permis diseminarea rapidă a noilor idei intelectuale și științifice;
- Dezvoltarea limbilor oficiale a făcut posibilă comunicarea cetățenilor în interiorul anumitor granițe (adeseori politice);
- Predarea și traducerea au facilitat schimbul între limbi;
- Crearea de principii jurnalistice și bibliografice a asigurat calitatea și disponibilitatea materialelor imprimate;
- Crearea diferitelor tipuri de media, precum ziarele, radioul, televiziunea, cărțile etc. a satisfăcut nevoile de comunicare.

În ultimii 20 de ani, tehnologia informației a contribuit la automatizarea și facilitarea multora dintre procese:

- Software-ul pentru tehnoredactare computerizată înlocuiește dactilografierea și culegerea;
- Microsoft PowerPoint înlocuiește retroproiectorul;
- E-mail-ul trimite și primește documente mai repede decât un fax;
- Skype permite convorbirile prin Internet și găzduiește întâlnirile virtuale;
- Formatele de codificare audio și video ușurează schimbul de conținut multimedia;
- Motoarele de căutare oferă acces la pagini web bazat pe cuvinte cheie;
- Serviciile online precum Google Translate produc traduceri rapide și aproximative;
- Platformele de media socială facilitează colaborarea și partajarea de informații.

Deși astfel de instrumente și aplicații sunt utile, acestea nu sunt suficiente pentru a implementa o societate informațională europeană multilingvă și sustenabilă, o societate inclusivă și modernă unde informația și bunurile pot circula liber.

## Frontierele lingvistice frânează crearea unei societăți informaționale europene

Nu putem ști cu precizie cum va arăta viitoarea societate informațională. Când vine vorba despre discutarea unei strategii energetice sau a unei politici externe comune, am dori să-i ascultăm pe miniș-

*Suntem martorii unei revoluții digitale comparabilă cu inventarea tiparului de către Gutenberg.*

trii de externe europeni vorbind în limbile materne. Preferăm o platformă în care oamenii, vorbitori de limbi diferite și cu diverse competențe lingvistice, pot discuta un subiect specific în timp ce tehnologia le colectează automat opiniile și generează scurte rezumate. Am dori de asemenea să putem vorbi cu un birou de asistență pentru asigurări de sănătate localizat într-o altă țară.

Este limpede că este nevoie de o calitate a comunicării diferită celei de acum câțiva ani. Într-un spațiu de informare și economic global, suntem confrunțați cu mai multe limbi, mai mulți vorbitori și mai mult conținut și suntem nevoiți să interacționăm rapid cu noi tipuri de media. Popularitatea actuală a mediilor sociale (Wikipedia, Facebook, Twitter și YouTube) reprezintă doar vârful icebergului.

Astăzi putem recepționa gigaoceteți de text din orice colț al planetei în câteva secunde, doar pentru a afla că textul este într-o limbă pe care nu o înțelegem. Potrivit unui raport recent solicitat de Comisia Europeană, 57% dintre utilizatorii de Internet din Europa achiziționează bunuri și servicii în limbi diferite de cea maternă. (Engleza este cea mai comună limbă străină, urmată de Franceză, Germană și Spaniolă.) 55 % dintre utilizatori citesc conținut într-o limbă străină în timp ce doar 35% utilizează o altă limbă pentru a scrie e-mail-uri sau a publica comentarii pe web. Cu câțiva ani în urmă, engleza era privită ca *lingua franca* (limba de lucru) a Internetului – o vastă majoritate a conținutului era scrisă în această limbă – dar situația s-a schimbat drastic acum. Cantitatea de conținut online în alte limbi, în special asiatice și arabe, a explodat.

În mod surprinzător, ubicua diviziune digitală provocată de frontierele lingvistice nu a câștigat încă prea multă atenție în discursul public; totuși, ea ridică o întrebare foarte presantă: „Care dintre limbile europene va reuși să ajungă și apoi să persiste în societatea virtuală a informației și cunoașterii?”

## Limbile noastre sunt în pericol

Tiparul a contribuit la un inestimabil schimb de informații în Europa, dar a condus de asemenea la extincția multora dintre limbile europene. Limbile regionale și minoritare au fost tipărite arareori. În consecință, multe limbi precum dalmata sau limba din Cornwall erau adeseori limitate la forme orale de transmitere, care le-au restricționat adoptarea, răspândirea și utilizarea.

Cele aproximativ 60 de limbi vorbite astăzi în Europa reprezintă unul dintre cele mai bogate și importante bunuri culturale ale sale. Această multitudine de limbi este în același timp o componentă importantă a succesului său social<sup>ii</sup>. În timp ce limbi populare precum Engleza sau Spaniola vor rămâne cu siguranță prezente în societatea și pe piața digitală emergentă, multe limbi europene ar putea fi deconectate de la comunicarea digitală și ar putea deveni irelevante pentru societatea Internetului.

O astfel de evoluție ar fi cu siguranță de nedorit. În primul rând, ea ar însemna că o oportunitate strategică rămâne nefolosită, slăbind poziția Europei pe piața globală. În al doilea rând, o astfel de evoluție ar fi în contradicție cu obiectivul esențial al participării egale a fiecărui cetățean european, indiferent de limba lui. Potrivit unui raport recent al UNESCO privind multilingvismul, limbile străine reprezintă un mediu esențial pentru exercitarea drepturilor fundamentale precum exprimarea politică, educația și participarea în societate<sup>iii</sup>.

*Un spațiu economic și de comunicare global ne confruntă cu mai multe limbi, mai mulți vorbitori, mai mult conținut.*

*Care dintre limbile europene va reuși să ajungă și apoi să persiste în societatea virtuală a informației și cunoașterii?*

*Larga varietate de limbi a Europei este unul dintre cele mai importante bunuri culturale ale sale și o componentă esențială a succesului său social.*

## Tehnologia limbajului este cheia activării tehnologiei

În trecut, eforturile de investiții financiare au fost concentrate asupra educației lingvistice și a traducerii. De exemplu, potrivit anumitor estimări, piața europeană de traducere, interpretare, software de localizare și globalizare a paginilor de Internet a valorat 8.4 miliarde € în 2008 și este de așteptat să crească cu 10% pe an<sup>iv</sup>. Totuși, această capacitate nu este suficientă pentru a satisface nevoile curente și viitoare.

Tehnologia limbajului este o tehnologie cheie care poate proteja și promova limbile europene. Ea ajută oamenii să colaboreze, să conducă afaceri, să împărtășească cunoștințe și să participe în dezbateri politice și sociale, independent de barierele lingvistice sau de competențele de lucru cu calculatorul. Tehnologia limbajului asistă deja sarcinile de zi cu zi, precum scrierea e-mail-urilor, căutările online sau rezervarea unor bilete de avion. Beneficiem de pe urma tehnologiei limbajului atunci când:

- ❑ Descoperim informații cu un motor de căutare pe Internet;
- ❑ Verificăm corectitudinea ortografică și gramaticală într-un editor de text;
- ❑ Vizualizăm recomandări de produse într-un magazin online;
- ❑ Ascultăm instrucțiunile verbale ale unui sistem de navigare;
- ❑ Traducem pagini web cu un serviciu online.

Tehnologiile limbajului detaliate în continuare reprezintă o componentă esențială a aplicațiilor inovative viitoare. Tehnologia limbajului este de obicei o tehnologie auxiliară în interiorul unui cadru de aplicare mai larg, precum un sistem de navigare sau un motor de căutare. Cărțile albe realizate de METANET se concentrează pe disponibilitatea unor tehnologii lingvistice de bază pentru fiecare dintre limbi.

În viitorul apropiat, avem nevoie de tehnologie a limbajului care să fie disponibilă pentru toate limbile europene, accesibilă și bine integrată în medii software mai largi. O experiență de utilizator interactivă, multimedia și multilingvă nu este posibilă fără tehnologia limbajului.

## Oportunități ale tehnologiei limbajului

Tehnologia limbajului face posibilă, pentru toate limbile europene, aplicații precum traducerea automată, generarea de conținut, procesarea informației și managementul cunoștințelor. De asemenea, ea poate continua dezvoltarea de interfețe intuitive bazate pe limbă pentru dispozitive electrocasnice, utilaje, vehicule, computere și roboți. Deși există deja multe prototipuri, aplicațiile comerciale și industriale sunt încă în stadii incipiente de dezvoltare. Realizări recente în cercetare și dezvoltare au creat o adevărată fereastră de oportunitate pentru TL. De exemplu, traducerea automată (TA) oferă o acuratețe rezonabilă pentru domenii specifice iar aplicații experimentale asigură managementul informației și cunoștințelor precum și producerea de conținut în multe din limbile europene.

Aplicații precum interfețe bazate pe limbă și pe voce sau sisteme de dialog sunt utilizate în mod tradițional în domenii specializate și prezintă adeseori performanțe limitate. Un câmp activ de cercetare este reprezentat de tehnologia dedicată operațiilor de salvare în zonele sinistrate. În astfel de medii cu risc înalt, acuratețea traducerii poate fi o problemă de viață și de moarte. La fel în ceea ce

*Tehnologia Limbajului ajută oamenii să colaboreze, să conducă afaceri, să împărtășească cunoștințe și să participe în dezbateri politice și sociale în diverse limbi.*

privește utilizarea tehnologiei limbajului în sectorul de îngrijire a sănătății. Roboți inteligenți cu capabilități trans-linguale au potențialul de a salva vieți.

Există oportunități uriașe de piață în sectorul educației și al divertismentului pentru integrarea tehnologiei limbajului în jocuri, oferte “edutainment” (educație prin divertisment), medii de simulare sau programe de formare. Serviciile mobile de informații, software-ul pentru învățarea de limbi străine asistată de calculator, mediile e-learning, instrumentele de auto-evaluare și cele de detectare a plagiatului sunt doar câteva exemple de zone de aplicare în care tehnologia limbajului poate juca un rol important. Popularitatea aplicațiilor de media socială precum Twitter și Facebook sugerează încă o ocazie în care tehnologii sofisticate ale limbajului sunt necesare pentru monitorizarea publicărilor, rezumarea discuțiilor, identificarea unor curente de opinie, detectarea răspunsurilor emoționale, descoperirea încălcărilor drepturilor de autor sau a situațiilor de abuz.

Tehnologia limbajului reprezintă o oportunitate uriașă pentru Uniunea Europeană, atât din punct de vedere economic cât și din perspectivă culturală. Multilingvismul a devenit regulă în Europa. Companiile, organizațiile și școlile europene sunt de asemenea multinaționale și diverse. Cetățenii doresc să comunice dincolo de frontierele de limbă care persistă pe Piața Comună Europeană. Tehnologia limbajului poate ajuta la depășirea acestor bariere, sprijinind în același timp utilizarea liberă și deschisă a limbilor. Mai mult, o tehnologie a limbajului europeană inovativă și multilingvă ne poate ajuta să comunicăm cu partenerii noștri globali și comunitățile lor multilingve. Tehnologiile limbajului sprijină o varietate de oportunități economice internaționale.

*Multilingvismul este regula, nu excepția.*

## Provocările tehnologiei limbajului

Deși tehnologia limbajului s-a dezvoltat considerabil în ultimii ani, ritmul actual al progresului tehnologic și al inovării este prea lent. Nu putem aștepta zece sau douăzeci de ani pentru a fi martorii unor îmbunătățiri semnificative în direcția unei mai bune comunicări și productivități în mediul nostru multilingv. Tehnologiile elementare care sunt utilizate pe scară largă, precum opțiunile de corectare gramaticală și ortografică din editoarele de text, sunt de obicei monolingve și sunt disponibile doar pentru câteva limbi.

*Ritmul actual al progresului tehnologic și al inovării este prea lent pentru a obține produse software substanțiale în următorii douăzeci de ani.*

Aplicațiile pentru comunicare multilingvă necesită un anumit nivel de sofisticare. Serviciile online de traducere automată precum Google Translate sau Bing Translator sunt excelente în a produce o bună aproximare a conținutului documentelor. Totuși, asemenea servicii, ca și aplicațiile profesionale de TA, întâmpină multe dificultăți atunci când este nevoie de traduceri precise și complete. Există multe exemple bine-cunoscute de traduceri greșite amuzante (ca de exemplu traducerile literale ale unor nume precum „Bush” sau „Kohl”) care ilustrează provocările cărora tehnologia limbajului încă trebuie să le facă față.

## Achiziția limbii

Ca să ilustrăm modul în care computerele interacționează cu limbajul și să explicăm de ce achiziția limbajului este o sarcină foarte dificilă, vom arunca o scurtă privire asupra modului în care oamenii achiziționează prima și a doua limbă, iar ulterior vom schița funcționarea sistemelor de traducere automată — există un motiv pentru care tehnologia limbajului este strâns legată de domeniul inteligenței artificiale.

Oamenii achiziționează competențele lingvistice în două moduri distincte. În primul rând, un copil învață o limbă ascultând interacțiuni între vorbitorii acelei limbi. Expunerea la exemple lingvistice concrete produse de utilizatorii limbii precum părinții, frații sau alți membri ai familiei ajută copiii să producă primele lor cuvinte sau fraze scurte la vârsta aproximativă de doi ani. Acest lucru este posibil pentru că oamenii au o predispoziție genetică specială pentru învățarea limbilor. Învățarea unei a doua limbi presupune un efort mult mai mare atunci când copilul nu este introdus într-o comunitate lingvistică de vorbitori nativi. La vârsta școlară, limbile străine sunt achiziționate de obicei prin învățarea structurii lor gramaticale, a vocabularului și a ortografiei din cărți și materiale educaționale care descriu cunoașterea lingvistică prin reguli abstracte, tabele sau texte exemplu. Învățarea unei limbi străine presupune mult timp și efort și devine din ce în ce mai dificilă cu înaintarea în vârstă.

*Oamenii achiziționează competențe lingvistice în două modalități diferite: învățând exemple și învățând regulile care stau la baza limbii.*

Cele două tipuri principale de sisteme de TL achiziționează capacități lingvistice într-o manieră similară oamenilor. Abordările statistice obțin cunoștințe lingvistice dintr-o colecție vastă de exemple concrete într-o anumită limbă sau din așa numitele texte paralele disponibile în una sau mai multe limbi. Algoritmii de învățare automată modelează acea facultate a limbajului care poate deriva șabloane de utilizare corectă a cuvintelor, frazelor scurte sau a propozițiilor complete, într-o anumită limbă sau în traducerea dintr-o limbă în alta. Numărul de propoziții necesare abordărilor statistice este uriaș iar calitatea performanței crește odată cu creșterea numărului de texte analizate. Nu este rar întâlnită antrenarea unor astfel de sisteme pe texte care conțin milioane de propoziții. Acesta este unul dintre motivele pentru care furnizorii de motoare de căutare sunt dornici de a colecta cât mai mult material scris. Corectarea erorilor de ortografie în editoarele de text, disponibilitatea informației online precum și servicii ca Google Search și Google Translate se bazează pe abordări statistice (orientate către date).

*Cele două tipuri principale de sisteme de TL achiziționează informația lingvistică într-o manieră similară oamenilor.*

Sistemele bazate pe reguli reprezintă a doua direcție majoră în domeniul TL. Experți din Lingvistică, Lingvistică Computațională sau Știința Calculatoarelor codifică gramatici și compilează liste de tip vocabular (lexicoane). Realizarea unui sistem bazat pe reguli este o activitate care necesită mult timp și efort intens, dar și experți cu specializare înaltă. O parte dintre cele mai performante sisteme de traducere automată bazate pe reguli se află în dezvoltare constantă de mai mult de douăzeci de ani. Avantajul acestor sisteme este că experții pot avea un control mai detaliat asupra procesării limbajului. Asta face posibilă corectarea sistematică a greșelilor din software și furnizarea de răspunsuri detaliate către utilizator, în special când sistemele bazate pe reguli sunt folosite pentru învățarea unei limbi. Datorită constrângerilor financiare, tehnologia limbajului bazată pe reguli este fezabilă doar pentru limbile majore.



## Limba română în societatea informațională europeană

### Fapte generale

Vorbită de aproximativ **29.000.000 milioane de vorbitori**, limba română este limba maternă a 25.000.000 de vorbitori: în jur de 21.500.000 de vorbitori în România<sup>v</sup> plus aprox. 3.500.000 de vorbitori în Republica Moldova<sup>vi</sup> (unde limba este denumită în mod oficial moldovenească). Țările vecine României (Albania, Bulgaria, Croația, Grecia, Ungaria, Fosta Republică Iugoslavă a Macedoniei, Serbia, Ucraina) și comunități de imigranți din Australia, Canada, Israel, America Latină, Turcia, S.U.A. și alte țări Europene și Asia-tice însumează în jur de 4.000.000 de vorbitori nativi de română<sup>vii</sup>.

Româna este de asemenea o limbă oficială în Provincia Autonomă Voivodina din Serbia, în Muntele Athos autonom din Grecia, în Uniunea Europeană și în Uniunea Latină; este recunoscută ca limbă minoritară în Ucraina.

Limba română are **4 dialecte**<sup>viii</sup>: Daco-Româna/ Româna, Aromâna (vorbită de aproximativ 600.000 de vorbitori în Albania, Bulgaria, Grecia și Macedonia), Istro-Româna (15.000 de vorbitori în 2 zone mici din Peninsula Istria, Croația) și Megleno-Româna (în jur de 5.000 de vorbitori în Grecia și Macedonia). Datorită numărului mic de vorbitori, aceste dialecte sunt incluse în Cartea Roșie a Limbilor pe Cale de Dispariție *UNESCO*.

În România există 18 minorități etnice recunoscute oficial; la ultimul recensământ (2002), cei mai numeroși erau ungurii (1.431.807) și rroma (535.140), urmași de germani, ucrainieni, ruși lipoveni, turci, sârbi, croați, sloveni, tătari, slovaci, bulgari, evrei, cehi, polonezi, greci, armeni, etc. Pentru toate aceste minorități, politicile lingvistice oficiale în România garantează drepturile acestora de a fi protejate în calitate de comunități lingvistice și de a utiliza limba maternă în medii private și publice, culturale și sociale, economice și de comunicare. Totuși, articolul 13 al Constituției prevede că "în România, limba oficială este româna". Mai mult, Legea nr. 500 din 12 noiembrie 2004 stipulează obligația ca orice text (fie el oral sau scris) de interes public să fie tradus sau adaptat în limba română<sup>ix</sup>.

### Particularitățile limbii române

Limba română este o limbă romanică orientală, care s-a format la distanță de surorile sale occidentale. Elemente ale latinei populare, din care a evoluat, sunt mai bine păstrate în această limbă izolată geografic: s-au moștenit structura morfo-sintactică latinească, particularități pe care alte limbi romanice le-au pierdut (precum declinările), au fost întărite elemente morfologice (reflexivul) sau au fost preluate elemente non-romanice (vocativul în -o).

Cea mai mare parte a vocabularului limbii române are origine latină, fie moștenit din latina vulgară, fie împrumutat pe cale savantă, în epoca modernă. 60% din vocabularul fundamental (i.e. cuvintele cunoscute și folosite curent de toți vorbitorii) este moștenit din latină.

În timpul colonizării Daciei de către romani (106-271 d.Hr.), coloniștii au impus limba latină ca limbă oficială. Cu toate acestea, studii comparative între vocabularul românesc și cel albanez dovedesc existența unui număr de aproximativ 100 de cuvinte păstrate din substratul traco-dac. Aceste cuvinte denumesc concepte fun-



damentale, precum părți ale corpului, elemente naturale, hrană. Ele sunt folosite și astăzi, sunt foarte frecvente, au dezvoltat o poliemie și familii lexicale bogate.

În timpul migrației triburilor slave pe teritoriul României de astăzi, limba română a suferit un proces de transformare în toate compartimentele: fonetică, vocabular, morfologie și sintaxă. Cu toate acestea, morfologia, care dă esența unei limbi, a rămas latinească în cele mai multe aspecte ale sale. Alfabetul chirilic a fost adoptat în această perioadă, mai ales datorită influenței bisericești. Slavona a fost limba în care s-a oficiat serviciul religios în biserica ortodoxă până în secolul al XVIII-lea, când româna a început un proces de latinizare, modernizare și occidentalizare. Atunci, multe cuvinte de alte origini au fost înlocuite de cuvinte latinești, împrumutate direct sau indirect, prin intermediul altor limbi romanice (franceză și italiană). Franceza ca limbă de cultură în ultimele două secole și Franța ca țara în care aristocrația română își trimitea copiii la învățătură justifică existența unui număr extrem de mare de cuvinte de această origine în limba română. În ultimul timp, limba engleză a luat locul francezei, iar româna are multe anglicisme, adaptate total, parțial sau deloc la sistemul său fonetic și morfologic.

Aspecte politice, economice și sociale din istoria poporului român explică existența cuvintelor de diverse origini: turcă, greacă, germană, maghiară, bulgară, rusă etc. În română au fost create cuvinte noi mai ales prin sufixare, deși studiile recente reflectă creșterea importanței pe care a căpătat-o în ultima vreme prefixarea.

Limba română are cinci litere cu diacritice: *ă, î, â, ș, ț*. Pentru ultimele două au circulat două variante: una cu virgulă sub literă, alta cu sedilă. Însă numai prima variantă este recomandată astăzi de Asociația de Standardizare din România (ASRO). Multe texte electronice nu sunt scrise cu diacritice, însă au fost create programe pentru a introduce diacriticele în mod automat în astfel de texte.

*Limba română are cinci litere cu diacritice: ă, î, â, ș, ț. Pentru ultimele două au circulat două variante: una cu virgulă sub literă, alta cu sedilă. Însă numai prima variantă este recomandată.*

Sistemul flexionar al limbii române este destul de bogat: pentru substantive, pronume și adjective există cinci cazuri și două numere, pentru verbe sunt două numere, fiecare cu câte trei persoane, cinci timpuri sintetice plus infinitivul, gerunziul și participiul. În medie, un substantiv poate avea cinci forme, un pronume personal șase, un adjectiv șase, iar un verb peste treizeci. În afară de sufixele morfologice și de desinențe, flexiunea cuvintelor mai prezintă și alternanțe fonetice în interiorul rădăcinii.

Româna este o limbă care permite nelexicalizarea subiectului pronominal, ca cele mai multe limbi romanice, de altfel:

*Știu.*

Explicația rezidă în sistemul flexionar bogat al verbelor, care au desinențe diferite pentru persoane și numere diferite.

Cu toate acestea, și dublarea subiectului este posibilă în română, atunci când un pronume personal dublează un grup nominal lexical:

*Vine el tata imediat!*

Structura este caracteristică limbajului familiar, marcând o anumită atitudine ilocutionară a vorbitorului: amenințare, promisiune, asigurare verbală.

Româna are în comun cu anumite dialecte spaniole și cu câteva limbi balcanice o structură cunoscută sub numele de „dublare cliti-

că”. Dublarea clitică pronominală în română se face cu pronume neaccentuate de dativ, de acuzativ sau ambele. De exemplu, în propoziția

*“I<sub>i</sub> l<sub>i</sub>-am dat mamei<sub>i</sub> pe Ion<sub>i</sub> la telefon.”*

substantivul *mamei* și cliticul de dativ *i* se referă la aceeași persoană, iar cliticul de acuzativ *l-* și substantivul în acuzativ *pe Ion* sunt tot coreferențiale. Prezența cliticelor în asemenea construcții este obligatorie, deși ele nu complinesc valențe verbale. Însă atunci când substantivele nu sunt prezente, pronumelor le revine sarcina de a satura valențele verbale:

*“I l-am dat la telefon.”*

Este obligatorie dublarea numelor proprii și a substantivelor articulate hotărât, cu funcție sintactică de complement direct sau indirect.

Limba română prezintă atât fenomenul concordanței negative, cât și dubla negație. Prezența marcatorului negativ *nu* în grupul verbal imprimă caracter negativ întregii propoziții și autorizează cuvintele negative în respectiva propoziție (concordanță negativă):

*„Nu am văzut pe nimeni niciodată aici.”*

Totuși, anumite configurații în care apar marcatorii și cuvintele negative trebuie interpretate ca având dublă negație (adică, în ciuda formei negative a verbului predicativ, enunțul are un conținut afirmativ). De exemplu, o propoziție principală negativă urmată de o subordonată cu verbul la forma negativă a modului conjunctiv este o astfel de configurație cu sens afirmativ:

*„Maria nu a vrut să nu spună nimic.”=*

*„Maria a vrut să spună ceva.”*

Cazul este sintetic în limba română: substantivul își schimbă forma pentru exprimarea cazului. Cu toate acestea, există și trei prepoziții care marchează cazul: *pe* pentru acuzativ (condiționată de trăsăturile animat, hotărât și specific ale grupului nominal), *la* pentru dativ și *a* pentru genitiv (ambele condiționate de prezența unui numeral în grupul nominal):

*L-am văzut pe colegul meu.*

*Am dat cărțile la trei dintre copii.*

*Cărțile a trei copii erau noi.*

## Dezvoltări Recente

Asemănător cu procesul de relatinizare din secolul al nouăsprezecelea, de după eliberarea de sub dominația greacă și otomană, limba română a parcurs în secolul al XX-lea un proces de trecere de la limbajul totalitar („limba de lemn”, discursul unidirecțional etc.) la utilizarea deschisă, în care noi modele lingvistice trebuie să se adapteze la tranziția socială și culturală. Astfel, asemănător multor altor limbi, româna traversează un proces continuu de internaționalizare, sub influența vocabularului anglo-saxon.

În domenii esențiale, precum științele politice, administrative și economice, în presă, publicitate, știința calculatoarelor etc., au fost împrumutate numeroase cuvinte sau cuvinte existente au căpătat sensuri noi, după model englezesc; terminologiile domeniilor noi

*Limba română este o limbă cu un sistem bogat de flexionare, cu diferite particularități lingvistice: permite elipsa subiectului, dublarea cliticelor, permite concordanța negativă și negație dublă.*

se bazează pe împrumuturi din engleză, vocabularul activ al oamenilor instruiți conține din ce în ce mai multe anglicisme, se pot observa noi modele intonaționale (mai ales în presa vorbită).

În anumite domenii, anglicismele au început să înlocuiască vocabularul limbii române. Un exemplu este folosirea titlurilor englezești pentru anunțuri de locuri de muncă, în special pentru poziții de conducere, de ex. ‘Human Resource Manager’ în loc de *Director de Resurse Umane*. O tendință puternică de a exagera folosirea anglicismelor poate fi observată în reclame. Bănci din România folosesc slogane promoționale de genul: *Cu cine faci banking?* Sau *Prima modalitate de plată contactless*, deși *banking* sau *contactless* snt anglicisme care nu au intrat în vocabularul comun și cu care majoritatea românilor nu sunt obișnuiți.

Exemplul de mai sus demonstrează importanța ridicării unui semnal de alarmă asupra unei dezvoltări care riscă să excludă din societatea informațională o mare parte a populației, care nu este familiară cu limba engleză.

## Cultivarea limbii în România

Academia Română, cel mai înalt forum cultural al țării, are printre obiectivele sale principale cultivarea limbii naționale. Scopul major al institutelor sale lingvistice a fost crearea și publicarea *Dicționarului Tezaur al Limbii Române*, proces care a durat aproape un secol. Seria mai veche, cunoscută sub numele de *Dicționarul Academiei* (DA), include 5 volume cu 3146 de pagini și 44890 de intrări lexicale și a fost realizată între anii 1913 și 1947. După o întrerupere, lucrul a fost reluat la mijlocul deceniului al șaptelea al secolului trecut cu o serie nouă, cunoscută sub numele de *Dicționarul Limbii Române* (DLR). Ultimul volum a fost publicat la Editura Academiei la începutul lui 2009. În total, DA și DLR au 33 de volume, peste 15000 de pagini și în jur de 175000 de intrări. Dicționarul a fost creat în stilul tradițional, „cu creionul pe hârtie”, cu citate adunate din peste 2500 de volume de literatură română scrisă.

Institutul de Lingvistică „Iordun Iordan – Al. Rosetti” are un program de cercetare ce urmărește cultivarea limbii, elaborează dicționare normative (*Dicționar ortografic, ortoepic și morfologic al limbii române, Dicționarul împrumuturilor neadaptate, Dicționarul termenilor oficiali*) și gramatici (*Gramatica limbii române, Dinamica limbii române actuale*).

Legea 500 din 12 noiembrie 2004 prevede ca toate textele scrise sau orale în limba română, care servesc interesul public, să respecte normele academice.

Institutul Limbii Române a fost creat cu scopul de a promova învățarea limbii române peste hotare, de a-i sprijini pe cei care învață limba română și de a le atesta cunoștințele de română\*. Există în străinătate peste 70 de centre în care se predă limba română ca limbă străină de către cadre didactice din învățământul universitar românesc.

*Există în străinătate peste 70 de centre în care se predă limba română ca limbă străină de către cadre didactice din învățământul universitar românesc.*

Și în România se manifestă un interes crescând pentru studiul limbii române printre străini, nu doar la nivel diplomatic (de către reprezentanții misiunilor diplomatice ale diverselor țări), ci și în mediul de afaceri. În afară de universități, care oferă cursuri de română ca limbă străină de obicei pentru studenții străini din România, există și numeroase firme particulare cu oferte mai ales pentru străini implicați în sectorul economic. Sunt organizate cur-

suri de vară de limba română pentru toate nivelurile, anual, în diverse locuri din țară, de Fundația Culturală Română, precum și de câteva instituții de învățământ superior (precum Universitatea „Al. I. Cuza” din Iași).

Cultivarea limbii în contextul înnoirii accelerate este o prioritate și pentru presă. Canalele naționale de radio și televiziune au emisiuni în care sunt discutate împreună cu specialiști și explicate publicului aspecte mai complicate ale limbii.

## Limba în educație

Conform Noului curriculum național (2000), româna se predă 4-5 ore obligatorii pe săptămână în școala gimnazială și 3-4 ore în liceu. Aspectele prescriptive ale conservării limbii se combină cu comunicarea, comportament axat pe competențe, accentuându-se relația limbă-cultură. Limba și literatura română reprezintă o materie obligatorie la examenele naționale (la absolvirea ciclului gimnazial și liceal; bacalaureatul cuprinde două probe de limba română: una orală și alta scrisă).

Limba și literatura română se studiază ca specializări principale sau secundare în peste 30 de universități de stat și particulare din România.

## Aspecte Internaționale

România este recunoscută pe plan internațional pentru literatura sa, lucrările principale ale lui Eminescu (marele poet național al României) fiind traduse în peste 60 de limbi. Alte nume cunoscute din literatura română sunt: Mircea Eliade, primul istoric care a scris o istorie a religiilor; Eugen Ionesco, unul dintre promotorii Teatrului Absurdului, sau Emil Cioran, cunoscut pentru filosofia lui.

În prezent, o mare parte din publicațiile științifice din domeniul TL sunt scrise în limba engleză, deși Consorțiul pentru Informatizarea Limbii Române - ConsILR – organizează anual un atelier științific dedicat cercetărilor în TL, cu lucrările atelierului publicate în limba română. O situație similară se regăsește și în celelalte domenii, mai puțin proeminentă pentru discipline precum filozofie, lingvistică, teologie sau pentru domeniul juridic.

Aceeași situație se întâlnește și în lumea afacerilor. În multe companii mari internaționale, engleza a devenit *lingua franca*, atât în comunicarea scrisă (e-mail și documente) cât și în cea orală, în special în companii multinaționale cu directori străini.

Tehnologiile limbajului pot rezolva această provocare din altă perspectivă prin oferirea unor servicii precum traducerea automată sau regăsirea de informații multilingvă în texte în limbi străine, ajutând astfel la diminuarea dezavantajelor personale și economice cu care se confruntă vorbitorii care nu au cunoștințe avansate de limbă engleză.

Minorități române trăiesc în țările vecine și în diaspora peste tot în lume. România promovează politici pentru păstrarea identității lingvistice și culturale de către comunitățile românești. Centrul Euxodius Hurmuzachi oferă sute de burse anual în România pentru minoritățile române din țările vecine. Sunt multe schimburi școlare și academice, mai ales cu Republica Moldova. Primele extinderi în sistem franciză ale școlilor și universităților din România au apărut în Republica Moldova în 2000.

În diferite comunități din diaspora, există inițiative diverse, prin care cei interesați pot studia limba și cultura română. De exemplu, Școala de limba română din Kitchener, Canada, oferă ore de limbă și cultură română pentru copii și adolescenți.

Institutele Culturale Române există în 19 orașe din lume (inclusiv București, New York, Paris, Londra, Roma, Istanbul etc.) și toate au drept preocupare importantă promovarea limbii române prin cursuri și evenimente culturale de diverse tipuri.

### Limba română pe Internet

Piața Internetului în România este în continuă creștere. În 2010, 44,2% dintre români aveau acces la un calculator acasă, iar 35,5% (i.e. 7.786.700 de români) erau utilizatori de Internet<sup>xi</sup> (aproximativ 60% dintre ei fiind utilizatori zilnici), ceea ce plasează România pe locul 8 într-un top 10 al utilizatorilor de Internet din Europa<sup>xii</sup>. Peste 500.000 de website-uri sunt înregistrate cu domeniul .ro.

Comparând aceste date cu cele din 2000, când numai 3,6% din populație (adică 800.000 de români) foloseau Internetul, observăm o creștere de aproape zece ori.

Un studiu al Uniunii Latine din 2007 arată că, similar cu tendința celorlalte limbi neolatine, prezența limbii române pe Internet crescut din 1998 până în 2007. Împărțind procentul de pagini web pentru fiecare limbă la procentul de prezență relativă a vorbitorilor limbii din lumea reală, s-a calculat vigoarea fiecărei limbi (sau prezența limbilor studiate în spațiul virtual). Deși acest coeficient este considerat unul redus pentru limba română (0,6 în 2007, în comparație cu 4,44 pentru engleză, 2,24 pentru franceză și 2,93 pentru italiană), româna este singura limbă care a cunoscut o creștere în vigoare în perioada 2005-2007 (înaintea integrării în Uniunea Europeană).

### Bibliografie selectivă

Grigore Brâncuș, *Vocabularul autohton al limbii române*, Bucharest, Scientific and Encyclopaedic Ed., 1983.

Alf Lombard, *La langue roumaine: Une presentation*, Paris, Klincksieck, 1974.

Ioana Vintilă-Rădulescu, *Limba română din perspective integrării în Uniunea Europeană*, [http://www.unibuc.ro/ro/limba\\_romn\\_din\\_perspectiva\\_integrri\\_europene](http://www.unibuc.ro/ro/limba_romn_din_perspectiva_integrri_europene)

Uniunea Latină în colaborare cu Funredes, *Limbile și culturile pe internet 2007* [http://dtil.unilat.org/LI/2007/index\\_ro.htm](http://dtil.unilat.org/LI/2007/index_ro.htm)

## Suport tehnologic pentru limba română

### Tehnologiile limbajului

Tehnologiile limbajului sunt tehnologii informatice specializate pentru lucrul cu limbile naturale. De aceea ele sunt adesea subsumate termenului Tehnologia limbajului uman. Limbajul uman se manifestă în scris și oral. În timp ce vorbirea este modul cel mai vechi și mai natural al comunicării umane, informațiile complexe și cea mai mare parte a cunoștințelor omenești sunt păstrate și transmise prin texte scrise. Tehnologia vorbirii și a textelor scrise prelucrează și produce limbaj în aceste două modalități de realizare. Dar limba are și aspecte comune vorbirii și scrierii, precum lexicul, cea mai mare parte a gramaticii și semantica. De aceea, o mare parte a tehnologiilor limbajului nu poate fi subsumată nici tehnologiei vorbirii, nici tehnologiei textelor scrise. Printre acestea se află tehnologiile care leagă limba de cunoaștere. Imaginea din dreapta ilustrează peisajul tehnologiilor limbajului. În comunicare, oamenii combină limbajul cu alte moduri de comunicare și cu alte mijloace de informare. Îmbinăm vorbirea cu gesturile și expresiile faciale. Textele electronice se combină cu imagini și sunete. Filmele pot conține limbaj în formă scrisă și vorbită. De aceea, tehnologia vorbirii și a textelor scrise se suprapune și interacționează cu multe alte tehnologii care facilitează prelucrarea comunicării multimodale și a documentelor multimodale.

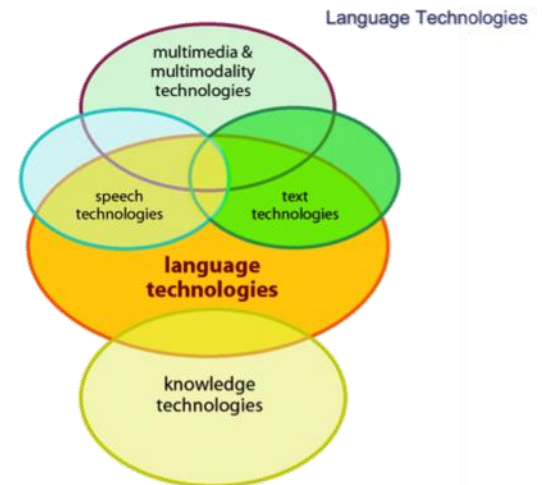


Figura 1: Suprapunerea tehnologiilor limbajului

### Arhitecturile aplicațiilor din tehnologia limbajului

Aplicațiile software tipice pentru prelucrarea limbii constau din câteva componente care reflectă diferite aspecte ale limbii și ale sarcinii pe care o implementează. Figura din dreapta prezintă arhitectura foarte simplificată a unui sistem de prelucrare a textelor. Primele trei module abordează structura și sensul textului introdus:

- Preprocesarea: curățarea datelor, eliminarea formatărilor, recunoașterea limbii introduse, înlocuirea diacriticelor greșite cu cele recomandate (de exemplu, înlocuirea lui ș cu sedilă cu ș cu virgulă).
- Analiza gramaticală: găsirea verbelor și a argumentelor sale, a modificatorilor etc.; recunoașterea structurii propoziționale.
- Analiză semantică: dezambiguizare (cu ce sens sunt folosite cuvintele în context?), rezoluția anaferei și a expresiilor referențiale precum „ea”, „mașina” etc.; reprezentarea sensului unei propoziții într-un mod accesibil calculatorului.

Modulele specifice sarcinii efectuează apoi mai multe operații diferite precum rezumare automată a unui text introdus, căutări în baze de date și multe altele. Mai jos vom ilustra principalele domenii de aplicații și vom evidenția anumite module ale diferitelor arhitecturi în fiecare secțiune. Arhitecturile sunt foarte simplificate și idealizate, servind pentru ilustrarea complexității aplicațiilor tehnologiei limbajului într-un mod inteligibil, la modul general.

După introducerea principalelor domenii de aplicații, vom face o scurtă prezentare a situației din cercetarea și predarea din domeniul tehnologiei limbajului, încheind cu o enumerare a programelor de finanțare (din trecut). La finalul acestei secțiuni vom prezenta evaluarea de către experți a instrumentelor și resurselor din techno-

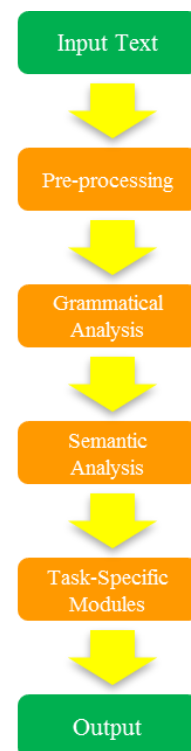


Figure 2: A Typical Text Processing Application Architecture



logia limbajului, pe baza unor criterii precum disponibilitate, maturitate sau calitate. Acest tabel reprezintă un tablou corect al situației tehnologiei limbajului pentru limba română.

Cele mai importante instrumente și resurse din domeniu sunt subliniate în text și pot fi găsite de asemenea în tabelul de la sfârșitul capitolului. Secțiunile care descriu principalele aplicații conțin de asemenea o trecere în revistă a companiilor active în domeniul respectiv în România.

## Principalele domenii de aplicații

### Corector de limbă

Oricine folosește un instrument de prelucrare a textului precum Microsoft Word a întâlnit o componentă care verifică ortografia, identifică greșelile de scriere și face sugestii de corectură. După 40 de ani de la apariția primului corector ortografic al lui Ralph Gorin, astfel de programe nu doar compară lista cuvintelor extrase cu cele dintr-un dicționar cu cuvinte scrise corect, ci au devenit extrem de sofisticate în zilele noastre: utilizează algoritmi dependenți de limbă pentru tratarea morfologiei (de exemplu, formarea pluralului), sunt acum capabile să recunoască greșeli de sintaxă, precum lipsa unui verb sau dezacordul în număr și persoană dintre verb și subiect, de exemplu „Ei scrie o scrisoare”. Cu toate acestea, cele mai multe corectoare ortografice disponibile (inclusiv cel din Microsoft Word) nu vor găsi nicio greșeală în următoarea strofă a unei poezii de Jerrold H. Zar (1992), poem bazat pe omofonie (pronunția identică a unor cuvinte cu grafie diferită) și lipsit de sens:

*Eye have a spelling chequer,*

*It came with my Pea Sea.*

*It plane lee marks four my revue*

*Miss Steaks I can knot sea.*

Pentru a da seama de astfel de greșeli este necesară analiza contextului în multe cazuri, de exemplu, pentru a decide dacă un cuvânt trebuie scris cu sau fără cratimă în română, precum în

*Plouă întruna de ieri.*

*Într-una din zile am să merg la Paris.*

Aceasta presupune fie formularea unor reguli gramaticale specifice limbii, adică un nivel înalt de expertiză și muncă manuală, fie utilizarea așa-numitelor modele lingvistice statistice. Acestea pot calcula probabilitatea ca un cuvânt să apară într-un anumit context (i.e. cuvintele dinainte și de după). De exemplu, *într-una din zile* este o secvență de cuvinte mult mai probabilă decât *întruna din zile*, iar *plouă întruna* este mai frecventă decât *plouă într-una*, deci în al doilea caz se recomandă scrierea fără cratimă. Un model de limbă statistic poate fi creat automat pe baza unei cantități mari de date (corecte) de limbă (adică un corpus). Și totuși, sunt cazuri când nici măcar acesta nu este util::

*Plouă întruna din primele zile ale lui martie.*

*Plouă într-una din primele zile ale lui martie.*

Singurul element discriminatoriu aici este verbul. În prima propoziție acesta este la prezent, având un sens durativ. În a doua este la trecut. Cele două forme verbale sunt omografe (se scriu la fel), dar

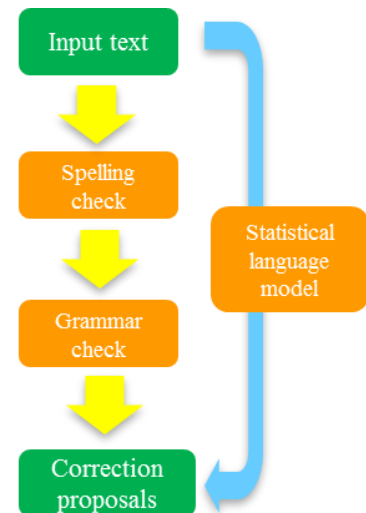


Figure 3: Language Checking (left: rule-based; right: statistical)

nu și omofone în română. Numai adnotarea morfo-sintactică are valoare discriminatorie în asemenea cazuri.

Până acum, aceste abordări au fost dezvoltate și aplicate mai ales pe date de limbă engleză. Ele nu pot fi transferate direct în limba română, care are o morfologie mai bogată și construcții specifice.

Utilitatea verificării limbii nu se limitează la instrumentele de prelucrare a textelor, ci se regăsește și în sistemele suport pentru autori. Urmare a creșterii numărului de produse tehnice, cantitatea de documentație tehnică a crescut vertiginos în ultimele decenii. Pentru a evita reclamațiile clienților în legătură cu utilizarea incorectă și pretențiile pentru pagube rezultate din instrucțiunile greșite sau din înțelegerea greșită a acestora, companiile au început să se concentreze tot mai mult pe calitatea documentației tehnice, ținând în același timp piețele internaționale. Evoluția prelucrării limbajului natural a dus la dezvoltarea unui software pentru sprijinul autorilor, care îl asistă pe cel care scrie documentații tehnice să folosească vocabularul și structurile sintactice conforme cu anumite reguli și restricții terminologice (ale corporațiilor).

Astăzi nu există companii românești sau furnizori de servicii lingvistice care să ofere astfel de produse, deși cercetătorii din diverse grupuri de prelucrare a limbajului natural au dezvoltat modele de limbă ajustate la particularitățile României. La Institutul de Cercetări pentru Inteligență Artificială al Academiei Române (RACAI) au fost create modele de limbă pentru română pe baza unor corpusuri de mari dimensiuni. Întrucât majoritatea textelor de pe Web sunt scrise fără diacritice, RACAI a mai dezvoltat și o aplicație de inserare a diacriticelor<sup>xiii</sup>, care are scopul de a indica diacriticele corecte ale unui cuvânt scris inițial fără diacritice; această aplicație folosește un lexicon românesc de mari dimensiuni dezvoltat în cadrul Institutului și un model de ferestre de 5 caractere pentru a găsi cea mai probabilă interpretare în termeni de diacritice a unui cuvânt necunoscut. Metoda de lucru ia în considerare contextul unui cuvânt în faza de preprocesare prin adnotare morfo-sintactică, esențială pentru alegerea cuvântului corect din lexicon. De exemplu, cuvântul „peste” este transformat în „pește” în

*Am cumparat peste.*

dar este păstrat ca „peste” în

*Era un pod peste rau.*

Această decizie se bazează pe o etapă anterioară de adnotare morfo-sintactică, în care „peste” din prima propoziție este adnotat cu o etichetă substantivală, iar același cuvânt din a doua propoziție este adnotat cu o etichetă prepozițională.

În română, cel puțin 30% dintre cuvintele dintr-o propoziție folosesc semne diacritice, cu o medie de 1.16 semne diacritice per cuvânt. Doar aproximativ 12% dintre aceste cuvinte pot fi transformate imediat în versiunea lor cu diacritice (întrucât forma fără diacritice nu este un cuvânt valid în dicționarul limbii române). Pentru restul cuvintelor, este util programul de descoperire a diacriticelor.

În afară de corectoarele de limbă și sistemele suport pentru autori, verificarea limbii este importantă și în domeniul învățării limbilor cu ajutorul calculatorului și se folosește la corectarea automată a întrebărilor introduse în motoarele de căutare pe Web: vezi sugestiile „Ați vrut să scrieți...” din Google.



## Căutarea pe Web

Căutarea pe Web, în Intranet sau în biblioteci digitale este, probabil, cea mai folosită și totuși cea mai subdezvoltată tehnologie a limbajului astăzi. Motorul de căutare Google, care a apărut în 1998, este folosit pentru 80% dintre căutările la nivel mondial<sup>xiv</sup>. Nici interfața de căutare, nici prezentarea rezultatelor căutării nu s-au schimbat semnificativ de la prima versiune. În actuala versiune, Google oferă corectarea grafică a cuvintelor scrise greșit și, din 2009, a încorporat abilități de căutare semantică elementară în pachetul de algoritmi<sup>xv</sup>, ceea ce poate îmbunătăți acuratețea căutării prin analiza sensului termenilor din fraza de interogare în context. Povestea de succes a Google dovedește că, dispunând de o cantitate uriașă de date și de tehnici eficiente de indexare a acestora, o abordare în principal statistică poate conduce la rezultate satisfăcătoare.

Cu toate acestea, pentru o căutare mai sofisticată de informații, este esențială integrarea cunoștințelor lingvistice mai profunde. În laboratoarele de cercetare, experimentele care folosesc tezaure în format electronic și resurse lingvistice de tip ontologie precum WordNet (sau echivalentul său românesc Romanian WordNet<sup>xvi</sup>) au demonstrat îmbunătățiri ale rezultatelor prin apelul la sinonime ale termenilor de căutare, de exemplu *energie atomică* ori *energie nucleară* sau chiar termeni mai îndepărtați semantic.

Generația următoare a motoarelor de căutare va trebui să includă tehnologii ale limbajului mult mai sofisticate. Dacă o frază de interogare constă dintr-o întrebare sau alt tip de propoziție, nu dintr-o listă de cuvinte-cheie, găsirea răspunsurilor relevante la această frază necesită analiza propoziției la nivel sintactic și semantic, precum și existența unui index care să permită găsirea rapidă a documentelor relevante. De exemplu, imaginați-vă că un utilizator introduce fraza de interogare:

*Dă-mi o listă cu toate companiile care au fost preluate de alte companii în ultimii cinci ani.*

Pentru a găsi un răspuns satisfăcător, trebuie efectuată analiza sintactică pentru a identifica structura propoziției și a stabili faptul că utilizatorul caută companii care au preluat alte companii. De asemenea, expresia *în ultimii cinci ani* trebuie prelucrată pentru a stabili la ce ani se referă.

În sfârșit, trebuie încercată potrivirea dintre fraza prelucrată și o cantitate uriașă de date nestructurate pentru a găsi informația căutată de utilizator. Acest proces este cunoscut sub numele de regăsirea informației și presupune căutarea și ordonarea documentelor relevante. În plus, generând o listă de companii trebuie să extragem și informația că un anumit șir de cuvinte dintr-un document se referă la numele unei companii. Acest tip de informație ne este furnizat de programele de recunoaștere a entităților numite.

Și mai solicitantă este încercarea de potrivire a frazei de interogare cu documente scrise în altă limbă. Pentru regăsirea informației la nivel interlingual trebuie să traducem automat fraza de interogare în toate limbile sursă posibile și să transferăm informația regăsită în limba țintă. Procentul în creștere de date disponibile în format non-text necesită servicii care să permită regăsirea informației multimedia, i.e. căutarea informației în date de tip imagine, audio și video. Pentru fișierele audio și video, aceasta presupune un modul de recunoaștere a vorbirii care convertește conținutul de vorbi-

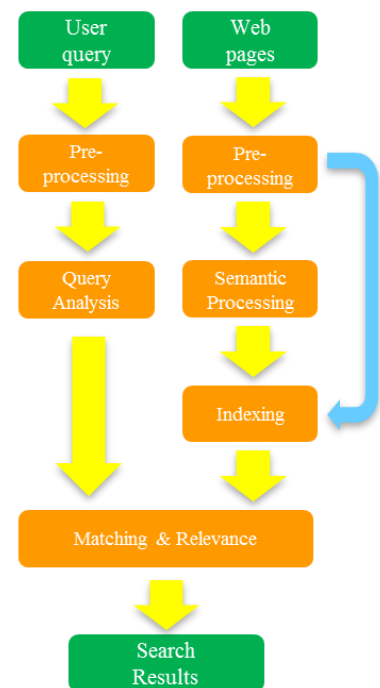


Figure 4: Web Search Architecture

re în text sau într-o reprezentare fonetică, la care se încearcă potrivirea frazei de interogare.

În România, tehnologiile de căutare bazate pe limba naturală nu sunt încă vizate de aplicațiile industriale. În schimb, tehnologiile de tip *open source* precum Lucene sunt adesea folosite de companiile care fac căutări pentru a furniza infrastructura elementară de căutare. Totuși, grupurile de cercetare de la Universitatea Al. I. Cuza (UAIC) și de la RACAI au dezvoltat diverse module care constituie partea centrală a unui instrument de căutare semantică, precum analizoare morfo-sintactice, analizoare sintactice, analizoare semantice, programe de recunoaștere a entităților numite, instrumente de indexare, programe de regăsire a informației multimedia etc. Acoperirea și eficiența lor sunt, totuși, destul de limitate.

Astfel, la RACAI, un analizor morfo-sintactic capabil să identifice forma de dicționar și partea de vorbire a cuvintelor din text este disponibil ca serviciu web<sup>xvii</sup>. De exemplu, dacă fraza de interogare a unui utilizator pentru o căutare pe web conține cuvântul „evenimente”, poate fi utilizată rădăcina („eveniment”) pentru a efectua căutarea<sup>xviii</sup>.

Alt modul dezvoltat de cercetătorii de la UAIC și de la RACAI este un program de recunoaștere a entităților numite, care, dându-se un text conținând nume de persoane, de companii, de organizații, de evenimente etc. (toate acestea cunoscute ca entități numite), poate să le recunoască în text. Pentru propoziția

*Maria și-a luat bilet la concertul trupei din vară de la Paris.*

acest sistem recunoaște „Maria” ca nume de persoană, „din vară” ca o referință temporală, iar „Paris” ca nume de loc.

Un analizor semantic dezvoltat la UAIC<sup>xix</sup>, disponibil pentru limba română, poate identifica într-o propoziție rolurile semantice diferite ale entităților. De exemplu, pentru propoziția de mai sus, sistemul identifică „Maria” ca persoana care face acțiunea, „bilet la concertul trupei” ca obiectul care a fost cumpărat. Asemănător, în exemplul

*Maria și-a luat fără ezitare bilet pentru a-și vedea trupa preferată*

sistemul recunoaște „fără ezitare” ca *modalitatea* în care Maria și-a cumpărat biletul, iar „pentru a-și vedea trupa preferată” reprezintă *scopul* pentru care biletul a fost achiziționat.

Recent, un grup de cercetători de la UAIC au început o cercetare pentru detectarea și adnotarea automată a imaginilor, în vederea dezvoltării unui instrument de căutare a imaginilor<sup>xx</sup>. Sistemul este încă într-o fază incipientă.

### **Interacțiunea vocală**

Tehnologia interacțiunii vocale reprezintă baza pentru crearea de interfețe care să permită utilizatorului să interacționeze cu mașinile utilizând limba vorbită mai degrabă decât, de exemplu, o interfață grafică, o tastatură și un mouse. Astăzi, interfețele vocale cu utilizatorul (VUI – Vocal User Interface) sunt de obicei utilizate pentru servicii complet sau parțial automatizate furnizate de companii, prin telefon, clienților, angajaților, sau partenerilor. Domenii de afaceri care se bazează foarte mult pe VUI sunt băncile, logistica, transportul public și telecomunicațiile. Alte utilizări ale tehnologiei interacțiunii vocale sunt interfețele pentru anumite dispozitive, ca

de exemplu sistemele de navigare ale autovehiculelor, și utilizarea limbii vorbite ca alternativă la modalitățile de input/output ale interfețelor grafice, ca de exemplu pe smartphone-uri.

Interacțiunea vocală cuprinde următoarele patru tehnologii diferite:

- ❑ Recunoașterea automată a vorbirii (RAV) este responsabilă pentru identificarea cuvintelor care au fost rostite într-o secvență de sunete rostite de utilizator.
- ❑ Analiza sintactică și interpretarea semantică presupune analiza structurii sintactice a enunțului utilizatorului și interpretarea acestuia conform scopurilor sistemului în care este integrată tehnologia.
- ❑ Managementul dialogului este necesar pentru determinarea acțiunii care va fi efectuată dat fiind inputul utilizatorului și funcționalitatea sistemului.
- ❑ Sinteza vorbirii (Text-to-Speech – TTS) este utilizată pentru transformarea cuvintelor unui enunț în sunetele care vor constitui output-ul pentru utilizator.

Una dintre provocările majore este realizarea unui sistem RAV care să recunoască cuvintele enunțate de utilizator cât mai precis cu putință. Acest lucru necesită fie o restrângere a domeniului enunțurilor posibile la un set limitat de cuvinte cheie, fie crearea manuală a unor modele de limbă care să acopere un interval larg de enunțuri în limbaj natural. În timp ce prima soluție rezultă într-o utilizare a unei VUI mai degrabă rigide și inflexibile și produce o acceptare slabă din partea utilizatorilor, opțiunea creării, reglării și menținerii unor modele de limbă poate crește semnificativ costurile. Totuși, VUI care utilizează modele de limbă și permit inițial utilizatorului să-și exprime intenția în mod flexibil prezintă atât o rată mai mare de automatizare cât și o acceptare mai mare din partea utilizatorului și pot fi, de aceea, considerate ca avantajoase în raport cu o abordare mai puțin flexibilă, precum cea a dialogului dirijat.

Pentru componenta de output a unei VUI, companiile tind să utilizeze enunțuri pre-înregistrate ale unor vorbitori profesioniști. Pentru enunțuri statice, în care rostirea nu depinde de contextul particular de utilizare sau de datele personale ale unui anumit utilizator, această soluție va conduce la o experiență utilă pentru utilizator. Totuși, cu cât conținutul unui enunț este mai dinamic, cu atât experiența utilizatorului are de suferit datorită unei prozodii sărace, rezultate din concatenarea diferitelor fișiere audio conținând diverse silabe-cuvinte. Prin contrast, sistemele TTS de azi se dovedesc a fi superioare (chiar dacă optimizabile) în ceea ce privește naturalitatea prozodică a enunțurilor dinamice.

În ceea ce privește piața tehnologiei interacțiunii vocale, ultimul deceniu a adus o puternică standardizare a interfețelor dintre diferitele componente tehnologice. A avut loc o puternică consolidare a pieței în ultimii zece ani, cu precădere în domeniile RAV și TTS. Aici, piețele naționale din țările blocului G20 – altfel spus țări puternice din punct de vedere economic și cu o populație considerabilă – sunt dominate de mai puțin de 5 actori mondiali, Nuance și Loquendo fiind cei mai proeminenți din Europa.

Domeniul recunoașterii și analizei vorbirii este unul dintre cele mai puțin reprezentate în România. Pe piața românească de TTS, există soluții comercializate de companii internaționale (precum MBROLA sau IVONA), dar rezultatele prezintă o acuratețe și o

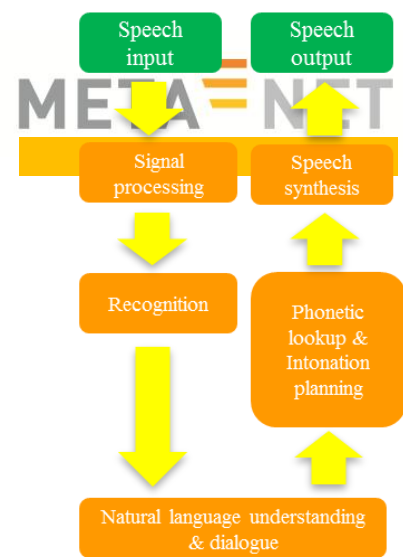


Figure 5: Simple Speech-based Dialogue Architecture

fluență redusă. Companiile de echipamente auto sau de telecomunicații, precum Continental și Orange, au început recent să aloce resurse pentru departamente specializate în procesarea vorbirii, adaptând soluții deja existente nevoilor lor specifice. Pe de altă parte, cercetări în această direcție au loc la Universitatea București și la Institutul de Informatică Teoretică al Academiei Române, Filiala Iași. Majoritatea cercetătorilor se concentrează pe sinteza vorbirii, în timp ce interpretarea vorbirii nu este încă dezvoltată.

Privind dincolo de starea actuală a tehnologiei, preconizăm schimbări semnificative datorită răspândirii smartphone-urilor ca o nouă platformă pentru administrarea relațiilor cu clienții, alături de canalele mai vechi precum telefon, internet și email. Această tendință va afecta și utilizarea tehnologiei pentru Interacțiune prin Voce. Pe de o parte, cererea pentru VUI bazate pe telefonie va scădea pe termen lung. Pe de altă parte, utilizarea limbii vorbite ca o modalitate facilă de input pentru smartphone-uri va căpăta o importanță semnificativă. Această tendință este sprijinită de progresul evident al acurateții recunoașterii vorbirii independente de vorbitor în cadrul serviciilor de dictare a vorbirii, care sunt deja oferite ca servicii centralizate utilizatorilor de smartphone-uri. Având în vedere “externalizarea” sarcinii de recunoaștere a vorbirii către infrastructura aplicațiilor, utilizarea tehnologiilor lingvistice într-o manieră specifică unei anumite aplicații va câștiga în importanță în raport cu situația actuală.

### **Traducerea automată**

Ideea de a folosi calculatoarele pentru traducere a apărut în 1946 la A.D. Booth și a fost urmată de finanțare substanțială pentru cercetări în acest domeniu între anii 1950 - 1980. Cu toate acestea, traducerea automată (TA) încă nu s-a ridicat la nivelul așteptărilor ridicate stabilite în primii ani de la apariția domeniului.

În cea mai simplă formă, TA înlocuiește pur și simplu cuvintele dintr-o limbă cu cuvintele din altă limbă. Acest lucru poate fi util în domenii cu limbaj foarte restrâns, formalizat, cum sunt de exemplu rapoartele meteo. Însă, pentru o traducere bună a unor texte mai puțin standardizate, trebuie potrivite elemente mai lungi din text (expresii, propoziții, sau chiar pasaje întregi) cu fragmentele lor echivalente din limba țintă. Dificultatea majoră aici constă în faptul că limbajul uman este ambiguu, ceea ce ridică provocări pe mai multe niveluri, de exemplu dezambiguizarea sensurilor cuvintelor la nivel lexical (Jaguar poate însemna fie o mașină fie un animal) sau atașarea corectă a grupurilor prepoziționale la nivel sintactic, ca în:

*Polițistul a văzut omul cu telescopul.*

*Polițistul a văzut omul cu arma.*

Una din modalitățile de abordare a traducerii automate se bazează pe reguli lingvistice. Pentru traduceri între limbi strâns legate, o traducere directă poate fi fezabilă în cazuri precum cele din exemplul de mai sus. Dar, de cele mai multe ori, sistemele bazate pe reguli (sau bazate pe cunoștințe) analizează textul de intrare și creează o reprezentare intermediară, simbolică, pe baza căreia este generat textul pentru limba țintă. Succesul acestor metode depinde în mare măsură de disponibilitatea unor lexicoane extinse cu informații morfologice, sintactice și semantice, precum și de existența unor seturi mari de reguli gramaticale atent proiectate de lingviști calificați.

Începând cu sfârșitul anilor 1980, pe măsură ce puterea de calcul a crescut și a devenit mai puțin costisitoare, au început să atragă interes modelele statistice pentru TA. Parametrii acestor modele statistice sunt derivați din analiza corpusurilor de texte bilingve, cum este corpusul paralel Europarl, care conține lucrările Parlamentului European în 11 limbi europene. Având date suficiente, TA statistică funcționează suficient de bine pentru a obține un înțeles aproximativ al unui text într-o limbă străină. Cu toate acestea, spre deosebire de sistemele bazate pe cunoștințe, sistemele de TA statistică (sau bazate pe date) generează de multe ori texte incorecte gramatical. Pe de altă parte, pe lângă avantajul că este necesar mai puțin efort uman pentru scris reguli gramaticale, sistemele de TA bazate pe date pot de asemenea acoperi particularități ale limbii care lipsesc în sistemele cunoștințe bazate pe cunoștințe, de exemplu expresii idiomatice.

Deoarece avantajele și dezavantajele sistemele de TA bazate pe date și a celor bazate pe cunoștințe sunt complementare, cercetătorii folosesc în prezent aproape în unanimitate abordări hibride, care combină cele două metodologii. Acest lucru poate fi realizat în mai multe moduri. Unul presupune folosirea atât a sistemelor de TA bazate pe cunoaștere cât și a celor bazate pe date, iar apoi un modul de selecție decide care dintre cele două traduceri este mai bună pentru fiecare propoziție. Cu toate acestea, pentru propoziții mai lungi, nici una dintre traduceri nu va fi perfectă. O soluție mai bună este de a combina, pentru fiecare propoziție, secvențele traduse corect de sisteme diferite. Sarcină destul de complexă, deoarece nu este totdeauna evident care sunt părțile traduse corect de fiecare sistem, și este necesară o aliniere.

Folosirea traducerii automate poate îmbunătăți semnificativ productivitatea cu condiția unei bune adaptări la terminologia specifică utilizatorului. Sisteme speciale pentru ajutor la traducerea interactivă au fost dezvoltate, în vreme ce portaluri lingvistice oferă acces la dicționare și terminologii specific companiilor, memorii de traducere și suport pentru TA.

Tabelul 1, prezentat în cadrul proiectului European Euromatrix+, arată performanțele obținute de sisteme la traducerile automate încrucișate în cele 22 de limbi oficiale ale Uniunii Europene (galeza irlandeză lipsește), raportate la scorul BLEU<sup>xxi</sup>.

Cele mai bune rezultate (în verde și albastru) le au limbile care beneficiază de eforturi de cercetare considerabile în domeniul TA în cadrul unor programe coordonate, și de existența unor corpusuri substanțiale paralele (ex. engleză, franceză, olandeză, spaniolă, germană). Rezultatele cele mai slabe (în roșu) sunt obținute de limbi care nu beneficiază de eforturi similar, sau care sunt foarte diferite din punct de vedere al comportamentului lingvistic față de alte limbi (ex. Ungară, malteză, finlandeză).

Calitatea sistemelor de TA este considerată ca având potențial de îmbunătățire. Modificările includ adaptabilitatea resurselor lingvistice la diferite domenii sau utilizatori, precum și integrarea în platforme existente cu memorii de traducere sau forme de bază a termenilor.

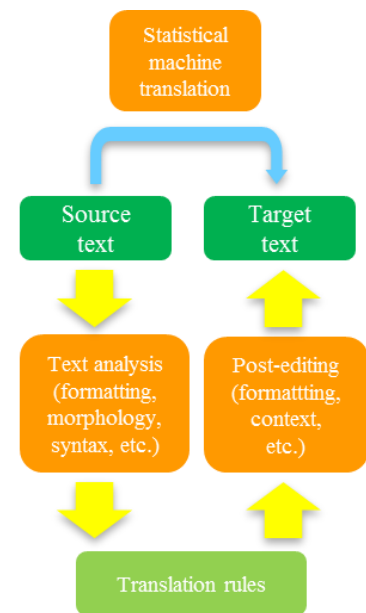


Figure 6: Machine translation (top: statistical; bottom: rule-based)



		Target Language																					
		en	bg	de	cs	da	el	es	et	fi	fr	hu	it	lt	lv	mt	nl	pl	pt	ro	sk	sl	sv
en	–	40.5	46.8	52.6	50.0	41.0	55.2	34.8	38.6	50.1	37.2	50.4	39.6	43.4	39.8	52.3	49.2	55.0	49.0	44.7	50.7	52.0	
bg	61.3	–	38.7	39.4	39.6	34.5	46.9	25.5	26.7	42.4	22.0	43.5	29.3	29.1	25.9	44.9	35.1	45.9	36.8	34.1	34.1	39.9	
de	53.6	26.3	–	35.4	43.1	32.8	47.1	26.7	29.5	39.4	27.6	42.7	27.6	30.3	19.8	50.2	30.2	44.1	30.7	29.4	31.4	41.2	
cs	58.4	32.0	42.6	–	43.6	34.6	48.9	30.7	30.5	41.6	27.4	44.3	34.5	35.8	26.3	46.5	39.2	45.7	36.5	43.6	41.3	42.9	
da	57.6	28.7	44.1	35.7	–	34.3	47.5	27.8	31.6	41.3	24.2	43.8	29.7	32.9	21.1	48.5	34.3	45.4	33.9	33.0	36.2	47.2	
el	59.5	32.4	43.1	37.7	44.5	–	54.0	26.5	29.0	48.3	23.7	49.6	29.0	32.6	23.8	48.9	34.2	52.5	37.2	33.1	36.3	43.3	
es	60.0	31.1	42.7	37.5	44.4	39.4	–	25.4	28.5	51.3	24.0	51.7	26.8	30.5	24.6	48.8	33.9	57.3	38.1	31.7	33.9	43.7	
et	52.0	24.6	37.3	35.2	37.8	28.2	40.4	–	37.7	33.4	30.9	37.0	35.0	36.9	20.5	41.3	32.0	37.8	28.0	30.6	32.9	37.3	
fi	49.3	23.2	36.0	32.0	37.9	27.2	39.7	34.9	–	29.5	27.2	36.6	30.5	32.5	19.4	40.6	28.8	37.5	26.5	27.3	28.2	37.6	
fr	64.0	34.5	45.1	39.5	47.4	42.8	60.9	26.7	30.0	–	25.5	56.1	28.3	31.9	25.3	51.6	35.7	61.0	43.8	33.1	35.6	45.8	
hu	48.0	24.7	34.3	30.0	33.0	25.5	34.1	29.6	29.4	30.7	–	33.5	29.6	31.9	18.1	36.1	29.8	34.2	25.7	25.6	28.2	30.5	
it	61.0	32.1	44.3	38.9	45.8	40.6	26.9	25.0	29.7	52.7	24.2	–	29.4	32.6	24.6	50.5	35.2	56.5	39.3	32.5	34.7	44.3	
lt	51.8	27.6	33.9	37.0	36.8	26.5	21.1	34.2	32.0	34.4	28.5	36.8	–	40.1	22.2	38.1	31.6	31.6	29.3	31.8	35.3	35.3	
lv	54.0	29.1	35.0	37.8	38.5	29.7	8.0	34.2	32.4	35.6	29.3	38.9	38.4	–	23.3	41.5	34.4	39.6	31.0	33.3	37.1	38.0	
mt	72.1	32.2	37.2	37.9	38.9	33.7	48.7	26.9	25.8	42.4	22.4	43.7	30.2	33.2	–	44.0	37.1	45.9	38.9	35.8	40.0	41.6	
nl	56.9	29.3	46.9	37.0	45.4	35.3	49.7	27.5	29.8	43.4	25.3	44.5	28.6	31.7	22.0	–	32.0	47.7	33.0	30.1	34.6	43.6	
pl	60.8	31.5	40.2	44.2	42.1	34.2	46.2	29.2	29.0	40.0	24.5	43.2	33.2	35.6	27.9	44.8	–	44.1	38.2	38.2	39.8	42.1	
pt	60.7	31.4	42.9	38.4	42.8	40.2	60.7	26.4	29.2	53.2	23.8	52.8	28.0	31.5	24.8	49.3	34.5	–	39.4	32.1	34.4	43.9	
ro	60.8	33.1	38.5	37.8	40.3	35.6	50.4	24.6	26.2	46.5	25.0	44.8	28.4	29.9	28.7	43.0	35.8	48.5	–	31.5	35.1	39.4	
sk	60.8	32.6	39.4	48.1	41.0	33.3	46.2	29.8	28.4	39.4	27.4	41.8	33.8	36.7	28.5	44.4	39.0	43.3	35.3	–	42.6	41.8	
sl	61.0	33.1	37.9	43.5	42.6	34.0	47.0	31.1	28.8	38.2	25.7	42.3	34.6	37.3	30.0	45.9	38.2	44.1	35.8	38.9	–	42.7	
sv	58.5	26.9	41.0	35.6	46.6	33.3	46.6	27.4	30.9	38.9	22.7	42.0	28.2	31.0	23.7	45.6	32.2	44.2	32.7	31.3	33.5	–	

Tabel 1: Performanțe obținute la traducerile automate în cele 22 de limbi oficiale ale Uniunii Europene (sursa: Euromatrix+)

Domeniul traducerii automate este, în ochii investitorilor, cel mai atractiv domeniu dintre tehnologiile limbajului. Astfel, companii precum Language Weaver lucrează în domeniul traducerilor din/spre română folosind diferite tehnologii lingvistice. Sistemul major online de traducere automată cuprinde limba română atât ca limbă sursă, cât și ca limbă țintă. De asemenea, există în mediul online o multitudine de dicționare pentru limba română.

Eforturi importante de cercetare au fost și continuă să fie dedicate de cercetători din diferite centre domeniului traducerii automate cu româna ca limbă sursă sau limbă țintă. Rezultate bune sunt raportate pentru un experiment de traducere bazată de date pentru perechea de limbi română-engleză, raportate la rezultatele sistemului Google Translate pentru aceeași pereche de limbi.<sup>xxii</sup>

La RACAI, de mai bine de 5 ani se experimentează cu diferite abordări de genul traducere automată bazată pe exemple, traducere automată statistică, extragerea de traduceri din corpusuri paralele etc. Două teze de doctorat, însoțite de mai multe articole științifice și susținute de diferite proiecte naționale sau internaționale, precum STAR și ACCURAT, sunt dedicate acestui domeniu<sup>xxiii, xxiv</sup>.

## Tehnologiile limbajului

Construirea de aplicații bazate pe tehnologiile limbajului implică o varietate de sub-probleme care nu apar întotdeauna la nivelul interacțiunii cu utilizatorul, dar oferă funcționalități semnificative “în cilise” sistemului. Din acest motiv, ele constituie domenii importante de cercetare care au devenit discipline de sine statătoare ale lingvisticii computaționale.

Sistemele de Întrebare-Răspuns (ÎR) reprezintă o zonă importantă a cercetării, pentru care au fost construite corpusuri adnotate și au fost inițiate competiții științifice. Ideea este trecerea de la căutarea bazată pe cuvinte cheie (la care sistemul răspunde print-o colecție de documente cu posibilă relevanță) la scenariul în care utilizatorul pune o întrebare concretă și sistemul oferă un singur răspuns : ‘La ce vârstă a pășit Neil Armstrong pe lună?’ - ‘38’. În vreme ce acest domeniu este în mod evident legat de domeniul Căutării pe Internet, sistemele ÎR sunt acum un termen general pentru cercetări de

genul ce *tipuri* de întrebări există și cum trebuie tratate, cum poate o colecție de documente cu un posibil răspuns să fie analizată și comparată (în cazul răspunsurilor conflictuale) și cum poate fi extras dintr-un document răspunsul (informație specifică) fără a ignora contextul. ÎR pot fi folosite cu succes pentru a identifica răspunsuri de tipul locație, persoană, organizație, dată, măsură, număr.

Acest domeniu este strâns legat de cel al extragerii informației (EI), o zonă care a fost extrem de populară și influențială în perioada statistică a lingvisticii computaționale, la începutul anului 1990. EI tinde să identifice bucăți de informație în diferite clase de documente; de exemplu, detectarea persoanelor cheie în preluările de companii, după cum sunt raportate în ziare. Alt scenariu care a fost luat în considerare e reprezentat de rapoartele asupra incidentelor teroriste, unde problema se reduce la potrivirea pe text a unui șablon care specifică atentatorul, ținta, locația și momentul incidentului, precum și rezultatul acestuia. Caracteristica principală a EI este completarea unor șabloane specifice fiecărui domeniu, din acest motiv fiind un exemplu de tehnologie din culise care constituie o arie de cercetare bine delimitată, dar pentru a avea aplicabilitate practică trebuie inclusă într-un mediu specific.

Două zone de limită, care uneori joacă rolul de aplicații independente, alte ori de componentă din culise, sunt rezumarea și generarea de text. Rezumarea se referă evident la scurtarea unui text lung, și este oferită ca funcționalitate de exemplu în MS Word. Una din abordările rezumării automate are baze statistice, identificând cuvinte “importante” din text (de exemplu cuvinte care au frecvență mare în text și care sunt mai puțin frecvente în utilizarea comună a limbajului) și apoi determinând acele propoziții care conțin aceste cuvinte importante. Propozițiile sunt apoi marcate în document, sau extrase din el, pentru a constitui rezumatul. În acest scenariu, rezumatul este o extragere de propoziții iar textul este redus la un subset din propozițiile sale. Majoritatea sistemelor de rezumare automată folosesc această metodă.

Un dezavantaj al acestei abordări este faptul că ignoră expresiile referențiale care pot apărea în textul inițial și să fie păstrate în rezumat. Dacă datorită eliminării de propoziții, antecedentul acestor referințe nu mai este prezent, rezumatul rezultat poate deveni de neînțeles. De exemplu, pentru textul:

*Hercule, dintre toți copiii nelegitimi ai lui Zeus, părea să fie centru mâniei Herei. Pe când el era doar un copil, ea a trimis un șarpe cu două capete să-l atace.*

rezumatul acestui fragment ar putea fi, folosind metoda de eliminare a propozițiilor:

*Ea a trimis un șarpe cu două capete să-l atace.*

ceea ce este destul de greu de înțeles dacă nu există nici o explicație despre cine este “ea” sau “el” (din cliticul –l se înțelege că există o persoană atacată care este de genul masculin).

O modalitate de a spori coerența acestor rezumate este de a deriva inițial structura de discurs a textului și de a ghida selecția propozițiilor care urmează a fi introduse în rezumat folosind un scor care să țină cont și de relevanța propoziției în discurs, dar și de coerența textului, rezultată din rezoluția anaforelor<sup>xxv</sup>. Pentru rezumatul dat ca exemplu mai sus, rezoluția anaforelor presupune identificarea

relației dintre “ea” și “Hera” și dintre “-l” și “Hercule”. Astfel, rezumatul devine inteligibil:

*Hera a trimis un șarpe cu două capete să-l atace pe Hercule.*

Sistemul de rezumare automată dezvoltat de UAIC a adoptat această metodă, producând rezumate foarte bune pentru texte de dimensiuni reduse<sup>xxvi</sup>.

O metodă alternativă căreia îi sunt dedicate multe cercetări este sintetizarea de *noi* propoziții, adică construirea unui rezumat din propoziții care nu sunt neapărat în forma din rezumat și în textul inițial. Această metodă necesită o înțelegere mai profundă a textului (ceea ce este mai costisitor din punct de vedere al resurselor computaționale și mai greu de realizat), dar poate fi aplicată cu succes pentru texte mai lungi. De exemplu, pentru romane nu este relevant calculul celor mai frecvente cuvinte (pentru că acestea vor fi cuvintele funcționale gen *și, iar, dar, al, etc.*), și nici structura de discurs (aceasta fiind mult prea stufoasă). În aceste cazuri, alte metode, exemplu expandarea unui set de șabloane flexibile predefinite (bazate de exemplu pe genul romanului, sau pe anumite informații despre personajele principale, timpul sau locația intrigii) pot fi aplicate.

Un generator de text nu este în majoritatea cazurilor o aplicație de sine stătătoare, ci inclusă într-o platformă software mai largă, așa cum într-un sistem de management medical sunt colectate, stocate și procesate informații despre pacient, iar generarea rapoartelor este doar o funcționalitate.

Limba română, ca limbă țintă pentru cercetările din toate aceste domenii, este mai puțin atractivă decât limba engleză, unde sistemele de întrebare-răspuns, de extragere de informații sau de rezumare automată au fost încă din anii 1990 subiectul a numeroase competiții, precum cele organizate de DARPA/NIST în Statele Unite sau campaniile CLEF în Europa. Totuși, echipe de cercetători români de la UAIC și RACAI au participat începând cu anul 2006 la competiții de întrebare-răspuns cu rezultate destul de bune<sup>xxvii</sup>. Principalul dezavantaj este dimensiune redusă a corpusurilor adnotate sau alte resurse necesare dezvoltării acestor domenii. Sistemele de rezumare automată, dacă folosesc doar metode statistice, sunt în mare măsură independente de limbă, astfel că există prototipuri care pot fi aplicate și pentru limba română. La UAIC, un instrument de rezumare bazat pe structura de discurs și rezoluția anaferei este disponibil pentru texte în limba română.

Domenii adiacente în care cercetători români au fost implicați cuprind lexicologie computațională, e-learning și analiza sentimentelor și a opiniilor.

Un consorțiu de cinci institute de cercetare din România și o universitate (UAIC) a fost implicat recent în transformarea în format electronic a Dicționarului Tezaur al Limbii Române, care însumează 35 de volume, redactate din 1913 până în prezent. Obiectivul principal a fost transformarea celor aprox. 15.000 de pagini ale dicționarului într-un format electronic structurat, care să permită căutări complexe, dar și o editare și activitate de actualizare mai ușoară<sup>xxviii</sup>.

Alt acces la materialul lexicografic al limbii este facilitat de rețelele semantice sub formă de wordnets (rețele de cuvinte). WordNetul românesc a fost în lucru timp de 8 ani, și totalizează mai mult de 52.000 de serii sinonimice (synset-uri) în care apar aprox. 60.000



de cuvinte, distribuite în patru categorii de părți de vorbire: substantive, verbe, adjective și adverbe. Fiecare synset conține un set de cuvinte (cu un număr de sensuri asociate) care sunt sinonime. Synset-urile sunt noduri ale rețelei, în timp ce arcele sunt relațiile semantice dintre synset-uri: hiponimie (relația *is-a*, care specifică că X este un fel de Y), meronimie, înlănțuire, cauză, și altele. WordNetul românesc este aliniat cu Princeton WordNet<sup>xxxix</sup> (varianța pentru limba engleză), primul și cel mai mare wordnet dintre cele existente pentru diferite limbi. Synset-urile au etichete DOMENIU: fiecare synset este etichetat cu numele domeniului în care este folosit. Mai mult, WordNetul românesc este aliniat cu cea mai mare ontologie disponibilă gratuit, SUMO&MILO<sup>xxx</sup>, și este folosit în diverse aplicații dezvoltate pentru limba română: sisteme de întrebare-răspuns, dezambiguizare a sensurilor cuvintelor, traducere automată.

Un alt domeniu în care cercetătorii din UAIC au fost implicați este e-learning, prin încorporarea instrumentelor multilingve de tehnologie a limbajului și tehnici de semantică web pentru îmbunătățirea regăsirii de materiale de învățare. Tehnologia dezvoltată facilitează accesul personalizat la cunoaștere în cadrul sistemelor de gestionare a învățării și ajută la operarea colectivă a datelor în gestionarea conținutului.

Cel mai nou domeniu de interes pentru tehnologiile limbajului este analiza sentimentelor și a opiniilor. Astfel, fiind dat un text, un program identifică dacă textul are o încărcătură emoțională pozitivă sau negativă. Cercetări în acest domeniu au început la RACAI cu dezvoltarea SentiWordNet, o adnotare la sentimente a WordNet-ului românesc<sup>xxxix</sup>. La UAIC, cercetări în această direcție au implicat colaborarea cu o companie privată, Sentimatrix, pentru dezvoltarea unui sistem capabil să monitorizeze web-ul și să extragă opinia utilizatorilor (din forumuri, bloguri, rețele sociale, etc.) referitoare la diferite produse<sup>xxxii</sup>.

## Tehnologiile limbajului în educație

Tehnologiile limbajului sunt un domeniu interdisciplinar, care implică expertiza lingviștilor, informaticienilor, statisticienilor, psiholingviștilor. Până acum nu și-a stabilit un loc fix în sistemul de învățământ din România. Multe universități din România și din Republica Moldova au introdus recent cursuri de prelucrare a limbajului natural și lingvistică computațională la nivelul studiilor universitare, de masterat și doctorat. Din 2001, un masterat în lingvistică computațională a fost introdus în curricula Facultății de Informatică a Universității Alexandru Ioan Cuza din Iași. Totuși trebuie conceput un sistem consolidat de educație superioară în procesarea limbajului natural și lingvistică computațională.

## Industria TL și programe

Industriile care folosesc și furnizează TL în România sunt cu siguranță importante (BitDefender, Continental, Nokia, etc.), dar este necesară o mai bună colaborare între ele. O problemă importantă care trebuie rezolvată este “caracterul secret” al TL, care ar putea fi rezolvată printr-o strategie bună de marketing. Industria limbajului nu este un angajator important în România, puține companii din domeniul Tehnologiilor Informației și Comunicării (TIC) având deja departamente de TL.

Programe naționale anterioare au avut un impuls inițial, dar lipsa ajutorului financiar consecvent sau destul de atractiv a dus la pierderea interesului marilor companii de TIC și a tinerilor cercetători,

formați de universități și de institutele de cercetare. Unul dintre programele de colaborare dintre industrie și educație care a avut un impact pozitiv și rezultate bune în România în domeniul TL este Alianța Academică MSDN, care oferă acces gratuit studenților la diferite tehnologii Microsoft.

Principalele laboratoare de cercetare cu activitate în domeniul TL în România sunt: RACAI la Academia Română din București; Departamentul de cercetare al Facultății de Informatică al Universității Alexandru Ioan Cuza University din Iași; Institutul de Informatică Teoretică al Academiei Române din Iași, care găzduiește arhiva Sunetele Limbii Române – un repozitoriu online de sunete ale limbii române înregistrate. În ceea ce privește programele de cercetare, UAIC și RACAI au fost implicate în mai multe proiecte de cercetare naționale și internaționale care-și propun să dezvolte tehnologii ale limbajului existente sau noi. Printre acestea pot fi menționate proiectele europene: ACCURAT-RO (Analiza și evaluarea corpusurilor comparabile pentru domenii cu puține resurse pentru traducere automată), STAR (Sistem pentru traducere automată pentru limba română), proiectul PC7 CLARIN (Infrastructură interoperabilă de resurse lingvistice pentru limba română), BALKANET (Construirea unei rețele de wordnet-uri pentru limbile balcanice), proiectul PC6 LT4eL (Tehnologii ale limbajului pentru e-learning), proiectul INTAS RoLTech (platformă pentru tehnologiile limbajului pentru limba română: resurse, instrumente și interfețe) etc. Au existat de asemenea proiecte cu finanțare națională precum: SIR-RESDEC (Sistem de întrebare răspuns pentru domeniu deschis pentru limbile română și engleză), ROTEL (Sisteme inteligente pentru web semantic, bazate pe logica ontologiilor și pe TL), eDTLR (Dicționarul Tezaur al Limbii Române în format electronic), printre altele.

Piața pentru tehnologiile limbajului poate fi doar estimată, și mai mult ca sigur va primi un impuls prin platformele mobile, de tipul Apple iPad și alte produse similare, jocuri (educaționale) etc.

### Cercetarea în domeniul TL și educația

Cele mai reprezentative centre în lingvistica computațională a limbii române sunt în România la București, Iași, Cluj-Napoca, Timișoara și Craiova, iar în Republica Moldova la Chișinău.

Punctele comune de întâlnire ale celor mai mulți cercetători din domeniul TL sunt, pe lângă conferințele internaționale din străinătate, o serie de evenimente internaționale și naționale care adună tinerii și cercetători cu experiență, lingviști și informaticieni ținute periodic în România: Evenimentele consorțiului pentru Informatizarea Limbii Române – ConsILR<sup>xxxiii</sup>, seria de școli de vară internaționale EUROLAN, conferințele SPED – Tehnologiile vorbirii și dialog om – calculator, conferințele KEPT – Ingineria cunoașterii: principii și tehnici, conferințele ECIT – conferința europeană pe domeniul sistemelor și tehnologiilor inteligente etc.

Lingvistica computațională este un domeniu exotic și este localizat fie la facultăți de informatică fie la facultăți de științe umaniste. Acest lucru este un dezavantaj pentru domeniul LT, deoarece studiul lingvisticii computaționale este astfel orientată fie pe aspecte lingvistice, fie pe cele de inginerie, iar suprapunerile sunt doar parțiale. Alt dezavantaj al acestui peisaj este implicarea minoră a companiilor din domeniul TIC în cercetarea în TL (deși recent au început să fie prezente în viața educațională).

## Situația instrumentelor și resurselor pentru limba română

Tabelul următor oferă o privire de ansamblu asupra situației prezente a tehnologiilor limbajului pentru limba română. Evaluarea tehnologiilor și resurselor existente este bazată pe estimarea mai multor experți din domeniu, folosind următoarele criterii (fiecare de la 0 la 6).

- 1 **Cantitate:** Există pentru limba română instrumentul sau resursa respectivă? Cu cât există mai multe instrumente /resurse cu atât este mai mare scorul.
  - 0: nici un instrument/resursă;
  - 6: multe instrumente/resurse, varietate mare.
- 2 **Disponibilitate:** Sunt accesibile instrumente/resurse, sunt Open Source, gratuite, pot fi folosite pe orice platformă sau sunt disponibile doar pentru un preț ridicat sau în condiții foarte restrictive?
  - 0: toate instrumentele/resursele sunt disponibile pentru un preț foarte mare;
  - 6: o mare parte dintre instrumente/resurse sunt disponibile gratuit cu licență Open Source sau Creative Commons care permit re folosirea lor.
- 3 **Calitate:** cât de bine sunt respectate criteriile de performanță ale instrumentelor și indicatorii de calitate ale resurselor de către cele mai bune instrumente, aplicații sau resurse? Sunt actualizate și menținute aceste resurse/instrumente?
  - 0: resursă/instrument de tip proof-of-concept;
  - 6: instrument de calitate înaltă, resursă cu adnotare manuală verificată.
- 4 **Acoperire:** În ce măsură cele mai bune instrumente îndeplinesc criteriile respective pentru acoperire (stil, gen, tipuri de texte, fenomene lingvistice, tipuri de intrări/ieșiri, număr de limbi pentru sistemele de TA etc.)? În ce măsură sunt reprezentative pentru limba țintă resursele?
  - 0: resursă sau instrument pentru scop specific, caz particular, acoperire foarte mică, utilizabil doar pentru cazuri foarte specifice;
  - 6: acoperire foarte largă, instrument foarte robust, aplicare pe scară largă, multe limbi suportate.
- 5 **Maturitate:** poate fi considerat matur instrumentul / resursa, stabil, pregătit pentru piață? cele mai bune instrumente / resurse pot fi folosite așa cum sunt sau trebuie adaptate? Performanțele unei astfel de tehnologii sunt adecvate și pregătite pentru folosirea în producție sau este doar un prototip? Un indicator poate fi dacă instrumentul/ resursa sunt acceptate de comunitate și folosite cu succes în sisteme de TL.
  - 0: prototip preliminar, exemplu de resursă;
  - 6: componentă integrabilă/aplicabilă imediat.
- 6 **Sustenabilitate:** Cât de bine poate fi integrat în sistemele IT existente? Resursa/instrumentul îndeplinește un anumit nivel de susținută în ceea ce privește documentația/manuale, explicarea diferitelor cazuri de utilizare, interfețelor etc.? Folosește medii de programare standard (pre-

cum Java EE)? Există standarde de cercetare /industriale și dacă da instrumentul / resursa le respectă (format al datelor etc)?

□ 0: format al datelor, API specific, ad hoc;

□ 6: respectă standardele, documentat.

7 **Adaptabilitate:** Cât de ușor poate fi adaptat / extins la domeniul nou / problemă nouă / tip diferit de text instrumentul respectiv?

□ 0: practic imposibil de adaptat instrumentul/resursa la altă problemă, chiar dacă sunt implicate multe resurse umane sau financiare;

□ 6: nivel înalt de adaptabilitate, adaptare foarte ușoară și eficientă.

## Tabel al instrumentelor și resurselor pentru limba română

	Cantitate	Disponibilitate	Calitate	Acoperire	Maturitate	Sustenabilitate	Adaptabilitate
<b>Tehnologia Limbajului (instrumente, tehnologii, aplicații)</b>							
Tokenizare, Morfologie (tokenizare, parsarea părților de vorbire, analiza/generarea morfologică)	5	4	5	5	5	4	4
Parsare (sintactică de adâncime)	3	3	4	4	4	3	4
Semantica Propoziții (WSD, structura argumentelor, roluri semantice)	4	3	4	4	3	4	4
Semantica Textului (coreferență, rezoluția anaferei, context, pragmatică, inferență)	3	3	5	4	5	5	5
Procesare avansată de discurs (structura textuală, coerență, statura retorică /RST, argumentare, șabloane de text, tipuri de text etc.)	3	3	4	3	3	3	3
Regăsirea de informații (indexarea de text, sisteme ÎR multimedia, sisteme ÎR multilingve)	3	4	5	5	5	5	5
Extragerea de informații (recunoașterea numelor de entități, extragerea relației dintre evenimente, recunoașterea sentimentelor)	3	4	4	5	4	4	5
Generarea de Limbaj (generarea de propoziții, generarea de rapoarte, generarea de text)	0	0	0	0	0	0	0
Rezumare, Sisteme de Întrebare-Răspuns, Tehnologii avasate de acces informațional.	5	4	5	4	5	4	4
Traducere automată	3	4	4	3	4	4	4
Recunoașterea vorbirii	2	1	3	2	2	2	2

	Cantitate	Disponibilitate	Calitate	Acoperire	Maturitate	Sustenabilitate	Adaptabilitate
Sinteza vorbirii	1	1	2	2	2	2	1
Gestionarea dialogului	0	0	0	0	0	0	0
Resurse Lingvistice (Resurse, date, baze de date, baze de cunoștințe)							
Corpus de referință	1	1	1	2	2	1	2
Corpus sintactic (treebanks, arbori de dependență)	3	3	5	4	4	4	4
Corpusuri semantice	2	3	3	2	3	3	3
Corpusuri de discurs	2	2	3	2	2	3	2
Corpusuri paralele, memorii de traducere	5	6	5	4	6	6	5
Corpusuri pentru vorbire	3	2	4	2	3	3	3
Date multimedia și multimodale (date text combinate cu date audio)	0	0	0	0	0	0	0
Modele de limbă	4	1	4	4	4	3	3
Lexicoane, terminologii	4	3	5	4	5	4	4
Gramatici	2	2	3	2	2	3	3
Thezaure, WordNets	4	3	4	4	5	5	4
Resurse ontologice pentru cunoașterea asupra lumii	1	2	2	2	2	2	2

Principalul scop al tabelului nu este de a oferi o privire exhaustivă a domeniului, ci de a sprijini mesaje de nivel înalt, care sunt explicate în această secțiune:

- Chiar dacă în general toate domeniile TL sunt acoperite, există domenii care nu sunt încă avute în vedere pentru limba română de cercetători: generarea de limbaj, sisteme de gestionare a dialogului și construirea de corpusuri multimodale.
- Deși sunt disponibile diferite tehnologii de parsare pentru limba română, un corpus de referință care să fie refolosit pentru evaluarea automată a parsărilor nu există încă
- Procesarea vorbirii este momentan mult mai puțin dezvoltată decât alte domenii ale TL, în ceea ce privește corpusurile și instrumentele.
- Dacă poate fi observată o atenție semnificativă pentru domenii precum tokenizarea, semantica propozițiilor sau sisteme de întrebare-răspuns, nu același lucru este valabil și pentru domenii care se ocupă de domenii mai complexe precum analiza sintactică de adâncime sau procesarea avansată a discursului.

- Resursele pentru limba română sunt mai puțin reprezentative decât instrumentele, deși sunt esențiale pentru testarea instrumentelor create.
- Cu câteva excepții cum ar fi serviciile web pentru procesări de bază ale limbajului, analiză morfologică, instrumente de întrebare-răspuns și sisteme de traducere automată, sistemele existente pentru limba română nu sunt disponibile.
- Instrumentele pentru limba română au o acoperire largă pentru domenii privind semantica propoziției și regăsirea de informații, dar sunt restrânse pentru celelalte probleme.
- Printre instrumentele existente de TL pentru limba română, cele mature sunt disponibile gratuit.
- Dacă instrumentele nu sunt în mod necesar menținute activ, resursele pentru limba română au o calitate bună și sunt în general sustenabile.
- Deoarece majoritatea instrumentelor sunt bazate pe modele de limbă sau folosesc tehnici de învățare automată, adaptarea lor este în general posibilă, ceea ce nu se întâmplă în cazul resurselor.
- Scorurile pe care diferiți experți le-au dat aceluiași domeniu din TL au fost în general asemănătoare, în special în ceea ce privește disponibilitatea, ceea ce indică faptul că instrumentele și resursele existente pentru limba română sunt diseminate pe scară largă. Uneori totuși, privind sustenabilitatea și acoperirea, experții au dat scoruri care diferă cu mai mult de jumătate din scorul total. Principalele zone de dezacord au fost: corpusul de referință, corpusuri semantice, gramatici și resurse ontologice.
- Rândul care conține informații despre modele de limbă poate fi interpretat diferit, deoarece unii experți au dat scoruri ținând cont de modele pentru limbajul scris, în timp ce alții au dat scoruri mai mici gândindu-se la modele pentru limbajul vorbit.

## Concluzii

Acest document descrie stadiul actual al tehnologiei limbajului în general cu aplicații pentru limba română în particular, și evidențiază sprijinul pe care TL îl pot aduce limbii.

Cercetările din universități și institute de cercetare au dus la dezvoltarea de sisteme de înaltă calitate, precum și modele și teorii aplicabile pe scară largă. Totuși este aproape imposibil să introducă soluții sustenabile și standardizate în contextul dezvoltării reduse a resurselor lingvistice. Există o necesitate majoră de resurse, de la texte în limba română până la corpusuri adnotate, în care fenomene lingvistice particulare sunt evidențiate de experți. Cum cea mai bună sursă de texte sunt copiile electronice ale publicațiilor, o campanie de conștientizare adresată editurilor, cu scopul de a le convinge să doneze o parte din textele lor pentru cercetare este mai mult decât necesară<sup>xxxiv</sup>.

Multe resurse nu sunt standardizate; programe și inițiative de standardizare a datelor și a formatelor interschimbabile sunt necesare.

Generarea de limbaj și sistemele de gestionare a dialogului sunt domenii ale TL la început de drum pentru limba română, pentru care se pot dezvolta încă multe tehnologii, aplicații și resurse. Tehnologiile și corpusurile pentru vorbire necesită o atenție deosebită în vederea alinierii limbii române la standardele celorlalte limbi europene.

Pentru dezvoltarea domeniilor TL în România, trebuie obținută o implicare mai accentuată a guvernului prin mecanisme politice și financiare adecvate. Acest document nu conține o evaluare a fondurilor alocate de Guvernul României domeniilor TL, dar sprijinul redus al statului pentru acest domeniu a fost observat în mod repetat de cercetătorii activi din domeniu, și se reflectă în numărul redus de proiecte finanțate din fonduri naționale, în implicarea timidă în crearea de infrastructuri europene de genul CLARIN-ERIC, și în lipsa interesului pentru cofinanțarea proiectelor ICT-PSP din domeniu (precum proiectul METANET4U, în cadrul căruia a fost luată inițiativă creării acestui raport).

## Despre META-NET

META-NET este o rețea de excelență finanțată de către Comisia Europeană. Rețeaua conține în prezent 47 de membri din 31 de țări europene. META-NET promovează Alianța Tehnologică pentru o Europă Multilingvă (Multilingual Europe Technology Alliance - META), o comunitate de profesioniști și organizații din domeniul tehnologiei limbajului din Europa în continuă creștere.

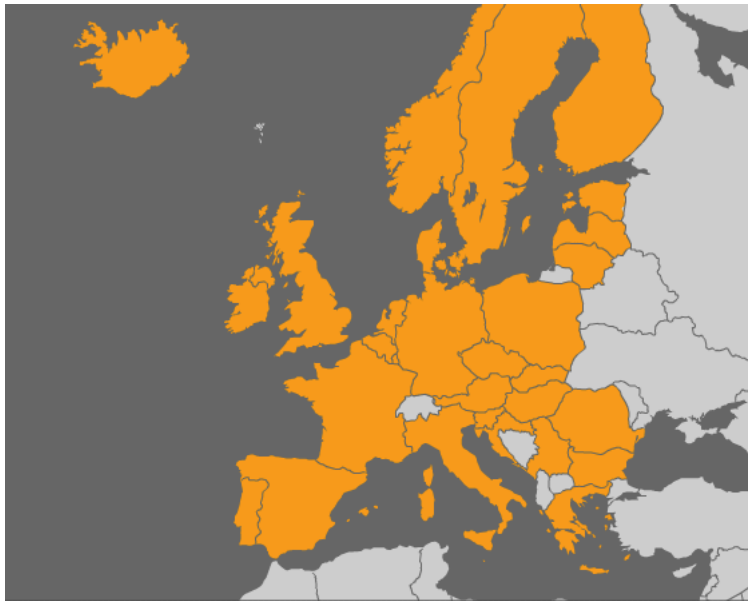


Figura 7: Țările reprezentate în META-NET

META-NET cooperează cu alte inițiative precum Infrastructura Comună pentru Resurse și Tehnologii ale Limbajului (Common Language Resources and Technology Infrastructure - CLARIN), care ajută la stabilirea cercetării în domeniul științelor umaniste digitalizate în Europa. META-NET promovează fundamentele tehnologice pentru stabilirea și menținerea unei societăți informaționale europene cu adevărat multilingve, care:

- ❑ vor facilita comunicarea și cooperarea între limbi diferite;
- ❑ vor asigura acces egal la informații și cunoaștere în orice limbă;
- ❑ vor oferi funcționalități ale tehnologiei informației cetățenilor europeni.

META-NET stimulează și promovează tehnologiile multilingve pentru toate limbile europene. Tehnologiile limbajului activează traducerea automată, producerea de conținut, procesarea informațiilor și gestionarea cunoștințelor pentru o gamă largă de aplicații și domenii. Rețeaua își propune să îmbunătățească abordările curente, pentru a facilita o comunicare și cooperare mai bună între limbi. Cetățenii Europei au drept egal la informație și cunoaștere, indiferent de limba vorbită.

### Linii de acțiune

META-NET a fost lansată pe data de 1 Februarie 2010 cu scopul de a avansa cercetarea din domeniul tehnologiilor limbajului. Rețeaua sprijină o Europă unită într-o singură piață digitală și spațiu informațional. META-NET a desfășurat mai multe activități care au



*The Multilingual Europe Technology Alliance (META)*



avansat scopurile sale. META-VISION, META-SHARE și META-RESEARCH sunt cele trei linii de acțiune ale rețelei.



Figura 8: Cele trei linii de acțiune ale META-NET

**META-VISION** promovează o comunitate dinamică și influențiază, unită în jurul unei viziuni comune și a unei agende strategice comune de cercetare. Principalul scop al acestei activități este constituirea unei comunități de TL coerente și coezive în Europa prin persoane cheie din diferite grupuri reprezentative. În primul an al META-NET, prezentările de la FLaReNet Forum (Spania), Language Technology Days (Luxemburg), JIAMCATT 2010 (Luxemburg), LREC 2010 (Malta), EAMT 2010 (Franța) și ICT 2010 (Belgia) s-au axat pe ridicarea numărului de adepți din public. În conformitate cu estimările inițiale, META-NET a contactat deja mai multe de 2.500 profesioniști din domeniul TL pentru a-și dezvolta viziunea împreună cu ei. La META-FORUM 2010 în Bruxelles, META-NET a comunicat rezultatele inițiale ale procesului de construire a viziunii unui public de peste 250 de participanți. În cadrul unor serii de sesiuni interactive, participanții au avut ocazia să discute viziunea prezentată de rețea.

**META-SHARE** creează o infrastructură publică distribuită pentru schimbul și partajarea de resurse. Rețeaua de depozite va conține date lingvistice, instrumente și servicii web documentate cu metadate de nivel înalt, organizate în categorii standardizate. Resursele pot fi accesate și permit căutări uniformizate. Resursele disponibile includ materiale gratuite, cu acces Open Source sau restricționat, precum și resurse disponibile contra unei sume. META-SHARE țintește datele lingvistice, instrumentele și sistemele existente, dar și produsele noi necesare pentru construcția sau evaluare de noi tehnologii sau servicii. Refolosirea, combinarea sau remodelarea datelor și instrumentelor lingvistice joacă un rol crucial. META-SHARE va deveni o componentă critică a pieței de tehnologii ale limbajului pentru dezvoltatori, experți în localizare, cercetători, traducători și profesioniști ai limbajului, de la firmele mici și mijlocii până la cele mari. META-SHARE vizează întregul proces de dezvoltare a TL – de la cercetare la produse și servicii inovatoare. Un aspect important al acestei activități este stabilirea META-SHARE ca componentă de valoare a unei infrastructuri globale europene pentru comunitate TL.

**META-RESEARCH** construiește punți de legătură între domenii tehnologice învecinate. Această activitate încearcă să aplice descoperirile recente din alte domenii pentru TL. În particular, această activitate este dedicată incorporării de mai multă semantică în Traducerea Automată (TA), optimizării diviziunii muncii în TA hibridă, exploatării contextului în vederea traducerii și pregătii unei fundații empirice pentru TA. META-RESEARCH lucrează cu alte discipline, precum învățarea automată și comunitatea web-ului semantic. META-RESEARCH se concentrează pe colectarea date-

lor, pregătirea seturilor de date și organizarea resurselor lingvistice pentru evaluare; crearea unor inventare de instrumente și metode; și organizarea de ateliere și evenimente de formare pentru membrii comunității. Această activitate a identificat deja aspecte ale TA unde informația semantică poate avea un efect substanțial în sistemele existente. În plus, activitatea a dezvoltat recomandări privind integrarea semanticii în TA. META-RESEARCH finalizează o nouă resursă lingvistică pentru TA, corpusul hibrid adnotat de mostre de TA, care oferă date pentru perechile de limbi engleză-germană, engleză-spaniolă și engleză-cehă. META-RESEARCH a dezvoltat de asemenea un instrument de colectare a corpusurilor multilingve ascunse pe web.

## Organizații membre

Tabelul următor prezintă lista organizațiilor și a reprezentanților lor care participă în META-NET.

Țara	Organizația	Participant/Participanți
Austria	Universitatea din Viena	Gerhard Budin
Belgia	Universitatea din Antwerp	Walter Daelemans
	Universitatea din Leuven	Dirk van Compernelle
Bulgaria	Academia de Științe din Bulgaria	Svetla Koeva
Croatia	Universitatea din Zagreb	Marko Tadić
Cipru	Universitatea din Cyprus	Jack Burston
Republica Cehă	Universitatea Charles din Praga	Jan Hajic
Danemarca	Universitatea din Copenhaga	Bolette Sandford Pedersen and Bente Maegaard
Estonia	Universitatea din Tartu	Tiit Roosmaa
Finlanda	University Aalto	Timo Honkela
	Universitatea din Helsinki	Kimmo Koskenniemi and Krister Linden
Franța	CNRS/LIMSI	Joseph Mariani
	Agencia de Evaluare și Distribuire a Resurselor Lingvistice (ELRA)	Khalid Choukri
Germania	DFKI	Hans Uszkoreit and Georg Rehm
	RWTH Aachen University	Hermann Ney
	University din Saarland	Manfred Pinkal
Grecia	Institutul pentru Prelucrarea Limbajului și a Vorbirii, "Athena" R.C.	Stelios Piperidis
Ungaria	Academia de Științe din Ungaria	Tamás Váradi

Țara	Organizația	Participant/Participanți
	Budapest Universitatea de tehnologii și Științe din Budapesta	Géza Németh and Gábor Olaszy
Islanda	Universitatea din Islanda	Eiríkur Rögnvaldsson
Irlanda	University din Dublin City	Josef van Genabith
Italia	Consiglio Nazionale Ricerche, Istituto di Linguistica Computazionale "Antonio Zampolli"	Nicoletta Calzolari
	Fondazione Bruno Kessler	Bernardo Magnini
Letonia	Tilde	Andrejs Vasiljevs
	Institutul de Matematică și Informatică, Universitatea din Letonia	Inguna Skadina
Lituania	Institutul pentru Limba Lituaniană	Jolanta Zabarskaitė
Luxemborg	Arax Ltd.	Vartkes Goetcherian
Malta	Universitatea din Malta	Mike Rosner
Olanda	University din Utrecht	Jan Odijk
	Universitatea din Groningen	Gertjan van Noord
Norvegia	Universitatea din Bergen	Koenraad De Smedt
Polonia	Academia de Științe din Polonia	Adam Przepiórkowski and Maciej Ogrodniczuk
	Universitatea din Lodz	Barbara Lewandowska-Tomaszczyk and Piotr Pezik
Portugalia	Universitatea din Lisabona	Antonio Branco
	Institutul pentru Ingineria Sistemelor și Calculatoare	Isabel Trancoso
România	Academia Română	Dan Tufis
	Universitatea Alexandru Ioan Cuza din Iași	Dan Cristea
Serbia	Universitatea din Belgrad	Dusko Vitas, Cvetana Krstev and Ivan Obradovic
	Institutul Mihailo Pupin	Sanja Vranes
Slovacia	Academia de Științe din Slovacia	Radovan Garabik
Slovenia	Institutul Jozef Stefan	Marko Grobelnik
Spania	Barcelona Media	Toni Badia
	Universitatea Tehnică Catalonia	Asunción Moreno
	Universitatea Pompeu Fabra	Núria Bel

Țara	Organizația	Participant/Participanți
Suedia	Universitatea din Gothenburg	Lars Borin
Regatul Unit	Universitatea din Manchester	Sophia Ananiadou
	Universitatea din Edinburgh	Steve Renals

## References

---

- <sup>i</sup> Diresctoratul general pentru Societatea Infromațională și Media, Comisia Europeană, *User language preferences online*, Flash Eurobarometer #313, 2011 ([http://ec.europa.eu/public\\_opinion/flash/fl\\_313\\_en.pdf](http://ec.europa.eu/public_opinion/flash/fl_313_en.pdf)).
- <sup>ii</sup> Conisia Europeană, *Multilingualism: an asset for Europe and a shared commitment*, Bruxelles, 2008 ([http://ec.europa.eu/education/languages/pdf/com/2008\\_0566\\_en.pdf](http://ec.europa.eu/education/languages/pdf/com/2008_0566_en.pdf)).
- <sup>iii</sup> Director-General UNESCO, *Intersectoral mid-term strategy on languages and multilingualism*, Paris, 2007 (<http://unesdoc.unesco.org/images/0015/001503/150335e.pdf>).
- <sup>iv</sup> Directoratul General pentru traduceri al Comisiei Europene, *Size of the language industry in the EU*, Kingston Upon Thames, 2009 (<http://ec.europa.eu/dgs/translation/publications/studies>).
- <sup>v</sup> Anuar statistic 2009 al Institutului Național de Statistică [http://www.insse.ro/cms/files/Anuar\\_statistic/02/02\\_Populatie\\_en.pdf](http://www.insse.ro/cms/files/Anuar_statistic/02/02_Populatie_en.pdf)
- <sup>vi</sup> Bază de date statistic a Biroului Național de Statistici al Republicii Moldova <http://statbank.statistica.md>
- <sup>vii</sup> [http://en.wikipedia.org/wiki/Romanian\\_diaspora](http://en.wikipedia.org/wiki/Romanian_diaspora)
- <sup>viii</sup> Marius Sala (ed.), *Encyclopaedia of the Romanian Language* (in Romanian), 2nd Edition, Bucharest, Univers Enciclopedic Publishing House, 2006.
- <sup>ix</sup> <http://www.efnil.org/documents/language-legislation-version-2007/romania>
- <sup>x</sup> <http://www.ilr.ro/plr.php?lmb=1>
- <sup>xi</sup> <http://www.internetworldstats.com/eu/ro.htm>
- <sup>xii</sup> <http://www.internetworldstats.com/stats9.htm>
- <sup>xiii</sup> Tufiș Dan și Ceaușu Alexandru (2008). *DIAC+: A Professional Diacritics Recovering System*, In Proceedings of Language Resources and Evaluation Conference, LREC 2008, Marakkech, Morocco. ELRA - European Language Resources Association. ISBN: 2-9517408-4-0
- <sup>xiv</sup> <http://www.spiegel.de/netzwelt/web/0,1518,619398,00.html>
- <sup>xv</sup> [http://www.pcworld.com/businesscenter/article/161869/google\\_rolls\\_out\\_semantic\\_search\\_capabilities.html](http://www.pcworld.com/businesscenter/article/161869/google_rolls_out_semantic_search_capabilities.html)
- <sup>xvi</sup> Tufiș, D., Radu I., Bozianu, L., Ceaușu, 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, ISBN 978-963-482-854-9.
- <sup>xvii</sup> [www.racai.ro/WebServices](http://www.racai.ro/WebServices).
- <sup>xviii</sup> Tufiș, D., Radu I., Ceaușu, A., Ștefănescu, D. (2008). *RACAI's Linguistic Web Services*, In Proceedings of Language Resources and Evaluation Conference, LREC 2008, Marakkech, Morocco. ELRA - European Language Resources Association. ISBN: 2-9517408-4-0
- <sup>xix</sup> Trandabăț D. (2011) *Towards automatic cross-lingual transfer of semantic annotation*, in 6e Rencontres Jeunes Chercheurs en Recherche d'Information RJCRI-CORIA 2011, 16-18 March, Avignon, France.
- <sup>xx</sup> Iftene, A., Vamanu, L., Croitoru, C. (2010). *UAIC at ImageCLEF 2009 Photo Annotation Task*. In C. Peters et al. (Eds.): CLEF 2009, LNCS 6242, Part II (Multilingual Information Access Evaluation Vol. II Multimedia Experiments). Pp. 283-286. ISBN: 978-3-642-15750-9. Springer, Heidelberg.

- xxi Cu cât scorul este mai mare, cu atât traducerea este mai bunăș un traucător uman ar obține în jur de 80. K. Papineni, S. Roukos, T. Ward, W.-J. Zhu. BLEU: A Method for Automatic Evaluation of Machine Translation. In Proceedings of the 40th Annual Meeting of ACL, Philadelphia, PA.
- xxii Munteanu D.S. și D. Marcu. (2005). *Statistical Machine Translation: An English-Romanian Experiment*. Tutorial invitat la EUROLAN 2005 Summer School on NLP & HLT: The Multilingual Web: Resources, Technologies, and Prospects, Cluj Napoca, Romania.
- xxiii Tufiș, Dan, Ceaușu Alexandru (2009). *Factored Phrase-Based Statistical Machine Translation*, In Corneliu Burileanu, Horia Nicolai Teodorescu (eds.) Proceedings of the 5th Conference "Speech Technology and Human-Computer Dialogue" SpeD 2009, IEEE Catalogue number:CFP095H-CDR
- xxiv Irimia. Elena (2009). *EBMT experiments for the English-Romanian Language Pair*. International Joint Conference Intelligent Information Systems (IIS 2009). Kraków, Poland, June 15-18, 2009.
- xxv O analiză detaliată a coerenței diferitelor texte este prezentată în Cristea, D., Iftene, A. (2011) *If you want your talk be fluent, think lazy! Grounding coherence properties of discourse*. Lucrare invitată la the University of Sussex, March.
- xxvi Cristea, D., Postolache, O., Pistol, I. (2005): *Summarisation through Discourse Structure*. In Alexander Gelbukh (Ed.): Computational Linguistics and Intelligent Text Processing, Proceedings of CICLing 2005, LNCS, vol. 3406, ISBN 3-540-24523-5, pp. 632-644.
- xxvii Vezi de exemplu sistemul prezentat în Iftene, A., Trandabăț, D., Moruz, A., Pistol, I., Husarciuc, M., Cristea, D. (2010). *Question Answering on English and Romanian Languages*. In C. Peters et al. (Eds.): CLEF 2009, LNCS 6241, Part I. Pp. 229-236. ISBN 978-3-642-15753-0. Springer, Heidelberg.
- xxviii Vezi mai multe detalii despre informatizarea Dicționarului tezaur al Limbii Române în Cristea, D. (2009): *Steps towards an electronic version of the Tesauros Dictionary of the Romanian language* (in Romanian), ASTRA 2009, Iași.
- xxix <http://wordnet.princeton.edu/>
- xxx Conținutul acestei resurse lexical poate fi găsit la <http://www.racai.ro/wnbrowser/>
- xxxi Baccianella Stefano, Andrea Esuli, Fabrizio Sebastian (2008). *SENTIWORDNET 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining*, In Proceedings of LREC 2008, pp. 2200-2204.
- xxxii Gînscă, A. L., Boroș, E., Iftene, A., Trandabăț, D., Toader, M., Corîci, M., Perez, C. A., Cristea, D. (2011). *Sentimatrix - Multilingual Sentiment Analysis Service*. In Proceedings of the 2<sup>nd</sup> Workshop on Computational Approaches to Subjectivity and Sentiment Analysis (ACL-WASSA2011) Portland, Oregon, USA, June 19-24, 2011.
- xxxiii Vezi colecția de volume cu lucrările publicate la Atelierul de lucru Instrumente și resurse lingvistice pentru limba română din 2006 până în 2010, Editura Universității "A.I. Cuza" Iasi, ISBN 978-973-703-208-9.
- xxxiv Vezi detalii ale soluției propuse în Cristea, Dan (2010). *Resurse lingvistice în flux continuu*. In Iftene, A. et al. (ed). Linguistic resources and instruments for Romanian language processing, Bucharest, University "Al. I. Cuza" Iași Ed. ISSN 1843-911X.