

Lucrările atelierului

*Resurse lingvistice și instrumente pentru
prelucrarea limbii române*

Iași, 3 noiembrie 2006

Editura Universității Alexandru Ioan Cuza Iași

Volum apărut cu sprijinul Ministerului Educației și Cercetării,
prin Autoritatea Națională pentru Cercetare Științifică

Descrierea CIP a Bibliotecii Naționale a României
Lucrările atelierului Resurse lingvistice și instrumente
pentru prelucrarea limbii române / Corina Forăscu, Dan
Tufiș, Dan Cristea (editori) –
Iași, Editura Universității „Alexandru Ioan Cuza”, 2006

ISBN: 978-973-703-208-9

I. Forăscu, Corina

II. Tufiș, Dan

III. Cristea, Dan

Lucrările atelierului
*Resurse lingvistice și instrumente pentru
prelucrarea limbii române*

Iași, 3 noiembrie 2006

Editori:

Corina Forăscu

Dan Tufiș

Dan Cristea

Organizatori:

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

Institutul de Cercetări pentru Inteligență Artificială
Academia Română – București

Institutul de Informatică Teoretică
Academia Română, Filiala Iași

Comitetul de Program

Corneliu Burileanu, Facultatea de Electronică, Universitatea Politehnica București și Institutul de Cercetări pentru Inteligență Artificială, A.R., București, România
Constantin Ciubotaru, Institutul de Matematică și Informatică, Academia de Științe a Moldovei, Chișinău, R. Moldova
Svetlana Cojocar, Institutul de Matematică și Informatică, Academia de Științe a Moldovei, Chișinău, R. Moldova
Dan Cristea, Facultatea de Informatică, Universitatea “Al. I. Cuza” și Institutul de Informatică Teoretică, A.R., Iași, România
Nicolae Curteanu, Institutul de Informatică Teoretică, A.R., Iași, România
Cristina Florescu, Institutul de Filologie Română "Al. Philippide", A.R., Iași, România
Corina Forăscu, Facultatea de Informatică, Universitatea “Al. I. Cuza”, Iași și Institutul de Cercetări pentru Inteligență Artificială, A.R., București, România
Gabriela Haja, Institutul de Filologie Română "Al. Philippide", A.R., Iași, România
Cătălina Hallett, Open University, UK
Radu Ion, Institutul de Cercetări pentru Inteligență Artificială, A.R., București, România
Rada Mihalcea, Universitatea North Texas, SUA
Constantin Orăsan, Universitatea Wolverhampton, Anglia
Oana Postolache, ISI - Universitatea California, SUA
Irina Prodanoff, ILC-Pisa și Universitatea Pavia, Italia
Georgiana Pușcașu, Universitatea Wolverhampton, Anglia
Valentin Tablan, Universitatea Sheffield, Anglia
Amalia Todirașcu, Universitatea Marc Bloch, Strasbourg, Franța
Dumitru Todoroi, Academia de Studii Economice, Chișinău, R. Moldova
Doina Tătar, Universitatea „Babeș-Bolyai”, Cluj-Napoca, Romania
Horia-Nicolai Teodorescu, Institutul de Informatică Teoretică, A.R. și Universitatea Tehnică, Iași, România
Dan Tufiș, Institutul de Cercetări pentru Inteligență Artificială, A.R., București și Universitatea “Al. I. Cuza”, Iași, România
Adriana Vlad, Facultatea de Electronică, Universitatea Politehnica București și Institutul de Cercetări pentru Inteligență Artificială, A.R., București, România

Comitetul de Organizare

Alexandru Ceaușu, ICIA-AR (alceasu@racai.ro)
Dan Cristea, FII-UAIC și IIT-AR (dcristea@info.uaic.ro)
Corina Forăscu, FII-UAIC și ICIA-AR (corinfor@info.uaic.ro)
Adrian Iftene, FII-UAIC (adiftene@info.uaic.ro)
Elena Irimia, ICIA-AR (elena@racai.ro)
Ionuț Pistol, FII-UAIC (ipistol@info.uaic.ro)
Dan Ștefănescu ICIA-AR (danstef@racai.ro)
Horia-Nicolai Teodorescu, IIT-AR și UT Iași (hteodor@etc.tuiasi.ro)
Diana Trandabăț, FII-UAIC și IIT-AR (dtrandabat@info.uaic.ro)
Dan Tufiș, ICIA-AR și FII-UAIC (tufis@racai.ro)

Introducere

Credem că limba română este răsplătită de eforturile de analiză, documentare, păstrare și publicare ale institutelor de lingvistică și universităților în mai bine de 100 de ani de cercetare (pentru a remarca numai perioada inaugurată de Hașdeu prin activitatea la dicționarul tezur). În acești ani s-au elaborat și tipărit dicționare, s-au emis și dezbătut teorii, s-au constituit puncte de vedere oficiale și personale și a fost suficient timp chiar și pentru contestarea unora dintre ele și perpetuarea unor dispute.

Între timp, limba română nu a stat nici ea pe loc, iar mijloacele de a studia limba s-au schimbat de asemenea. Dacă, pentru a-i studia evoluția sau pentru a găsi filioane lingvistice încă nedescoperite, atenția cercetătorului rămâne captată în continuare de aspecte de fonologie, sintaxă, semantică, lexicologie, terminologie etc., randamentul și precizia observațiilor lui crește dacă face apel la metode de investigare informatică a limbii. De câțiva timp accesul la o carte se poate face și altfel decât ținând-o în mână și deschizând-o. Și nu mai e nevoie ca ea să existe în biblioteca de lângă noi ca s-o putem citi. Dintr-o dată a devenit posibil să ne uităm la o carte și altfel decât parcurgând-o în secvența ei liniară. Rafturile cu fișe de ocurențe ale lexicografilor, care luau ani pentru a fi completate, sunt acum generate automat prin metode de indexare de către programe și regăsirea unui context se face cât ai clipi...

Dar domeniile lingvisticii computaționale și ale tehnologiilor limbajului uman au repercusiuni și de altă natură decât ca metode de cercetare asupra unei limbi. Aplicații de prelucrare a limbajului natural care să deschidă un nou tip de acces la informații pot fi acum concepute. Textul, chiar și în format electronic, începe să fie privit și altfel decât ca un șir de caractere sau de cuvinte. Au început să apară metode de a pătrunde în structura lui sintactică și semantică încât structura și înțelesul textului să poată fi relevate mașinii și ea să poată opera cu ele așa cum operează cu numere, de când a fost ea inventată. Începem să știm cum să facem mașinile noastre să execute un alt tip de „calcul”, mai apropiat de modul nostru de gândire, și care-și găsește originea în text...

Limba română trebuie să ajungă la nivelul de tehnologizare de care se pot mândri astăzi alte limbi intens studiate. Rostul acestei cărți, pe care o dorim prima dintr-o serie, trebuie atașat acestei ambiții. Ea este scrisă de lingviști și informaticieni români care, spre marea noastră bucurie, încep să se înțeleagă din ce în ce mai bine. Este exact ceea ce a urmărit acel grup de constituire a Comisiei de Informatizare pentru Limba Română, când, în martie 2001, s-a reunit pentru prima dată în sediul de pe Calea Victoriei al Academiei Române. Ulterior, această întâlnire a devenit o tradiție prin organizarea anual în București, Iași și Chișinău a unor sesiuni de lucru ale unui grup largit, care, din acest motiv s-a numit Consorțiu. De doi ani am dorit să invităm la aceste întâlniri și cercetători aflați la mai mare distanță de noi. Ca urmare, ultimele două întâlniri au căpătat caracterul de ateliere de lucru și au fost organizate în regim de teleconferință. Am putut asculta astfel glasuri de români care lucrează în universități din America, Germania, Italia, Franța și Anglia, după cum și ei ne-au putut urmări pe noi.

Întâlnirea din 3 noiembrie 2006 a Atelierului, a fost găzduită de Biblioteca Facultății de Informatică a Universității „Al.I.Cuza” din Iași și a beneficiat de implicarea MEC în finanțare. Această generoasă contribuție bănească ne-a permis să-i îmbunătățim organizarea, dar mai ales, să tipărim această carte. Îi suntem recunoscători pentru acest ajutor, cu precădere d-nei Veronica Bubulete. Mulțumim totodată participanților la atelier, aflați în sală sau conectați prin Internet, cât și colectivului de recenzori care ne-au ajutat să îmbunătățim calitatea lucrărilor.

Cuprins

Introducere

Capitolul 1. Resurse lingvistice pentru prelucrarea vorbirii1

Situl ‘Limba Română Vorbită’ Horia-Nicolai Teodorescu, Monica Feraru, Diana. Trandabăț.....	3
Schemă XML de adnotare a intonației în cadrul corpusurilor de text Vasile Apopei, Doina Jitcă	9

Capitolul 2. Dicționare și corpusuri adnotate pentru prelucrarea textelor.....15

Noi dezvoltări ale wordnet-ului românesc Dan Tufiș, Verginica Barbu Mititelu, Alexandru Ceaușu, Luigi Bozianu, Cătălin Mihăilă, Margareta Manu Magda	17
Framenet român: tentativă de elaborare Victoria Bobicev, Victoria Maxim, Tatiana Zidrașco, Alina Iaciurinschi.....	23
DEI Multimedia: evoluții, perspective Dumitru Todoroi, Adrian Chiorescu	29
Maparea cuvintelor dintr-un lexicon pe ontologie Natalia Burciu, Antonina Bîrlădeanu	35
Crearea resurselor lingvistice cu ajutorul unui limbaj specializat Ștefan Diaconescu	39
Resurse lingvistice românești în format electronic. <i>Biblia 1688</i> Bogdan-Mihai Aldea, Gabriela Haja	45
Resurse românești în cadrul proiectului LT4eL Diana Trandabăț, Adrian Iftene, Ionuț Pistol, Corina Forăscu, Dan Cristea.....	51
Tehnici de validare și corecție focalizată a adnotării morfo-sintactice în corpusuri de mari dimensiuni Dan Tufiș, Elena Irimia	57
RoGER – un corpus paralel aliniat Monica Gavrilă, Natalia Elița.....	63
TimeBank 1.2: O versiune adnotată în limba română Corina Forăscu, Radu Ion.....	69
Resurse lingvistice reutilizabile Constantin Ciubotaru, Svetlana Cojocar, Elena Boian, Alexandru Colesnicov, Ludmila Malahova, Valentina Demidov, Oleg Burlaca.....	75
Capitolul 3. Aplicații ale tehnologiilor lingvistice81	
Sisteme de Întrebare Răspuns pentru limba română Adrian Iftene, Ionuț Pistol, Diana Trandabăț, Georgiana Pușcașu, Corina Forăscu, Dan Cristea.....	83
Identificarea și extragerea automată a cologațiilor din texte Dan Ștefănescu, Dan Tufiș, Elena Irimia	89

Spre o extragere automată a colocațiilor: cazul verbului “a face” Amalia Todirașcu	95
Rezoluția anaforei pentru limba română Gabriela Pavel, Oana Postolache, Ionuț Pistol, Dan Cristea	101
Instrumente pentru consultarea Atlasului Lingvistic și editarea textelor dialectale Silviu Bejinariu, Vasile Apopei, Ramona Luca, Luminița Botoșineanu, Florin Olariu ...	107
Generare de concordanțe pentru dicționarul limbajului poetic eminescian Mihaela Brut, Dumitru Irimia, Oana Panait	113
Crearea unui generator morfologic pentru verbele din limba română Antonina Bîrlădeanu, Natalia Burciu	119
Parsarea predicatului (verbal / nominal) și a clauzei (finite / nefinite) în limba română. Aplicare la parsarea FDG Alex Moruz, Neculai Curteanu, Diana Trandabăț, Iustin Dornescu, Cecilia Bolea	123
Prelucrarea resurselor românești în cadrul proiectului LT4eL Ionuț Pistol, Adrian Iftene, Diana Trandabăț, Dan Cristea, Corina Forăscu	129
Sistem de instruire asistată de calculator pentru morfologia limbii române Elena Boian, Constantin Ciubotaru, Svetlana Cojocar, Galina Magariu, Tatiana Verlan, Iuri Rogojin	135
Capitolul 4. Modelare lingvistică	141
Structura grupului verbal, predicția lexicală și reprezentarea logică a predicatului în limba română Neculai Curteanu, Diana Trandabăț, Mihai Moruz	143
Perspective semantice din nou: cum și sub ce formă avansăm lexicologic spre DLRI Cristina Florescu	149
Modelarea relațiilor semantice într-un dicționar de simboluri Cristina Ciocârlău, Mihaela Brut	155
Dreptul de publicare pe web Noemi Bomher	161
Modelare cu ontologii și adnotări Radu Cibotaru	165
Cadre pentru o implementare PC-PATR a verbelor tranzitive din limba română Nadia Luiza Huțuliac	171
Index de autori	
177	

Capitolul 1

Resurse lingvistice pentru prelucrarea vorbirii

SITUL 'LIMBA ROMÂNĂ VORBITĂ'

HORIA-NICOLAI TEODORESCU^{1,2}, MONICA FERARU², DIANA TRANDABĂȚ^{1,3}

¹*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

²*Universitatea Tehnică „Gh. Asachi”, Iași*

³*Facultatea de Informatică, Universitatea “Al.I.Cuza”, Iași*

{hteodor, mferaru}@etc.tuiasi.ro, dtrandabat@info.uaic.ro

Rezumat

Inițiativa construirii unei arhive publice a sunetelor a fost determinată de lipsa unei asemenea resurse pentru limba română, lipsă resimțită atât în cercetare, cât și în învățământ. Situl include peste 600 de înregistrări, în diverse formate de precizie și codare, ale sunetelor limbii române vorbite.

1. Introducere

Situl (http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.htm) a fost creat prin colaborarea dintre Institutul de Informatică Teoretică al Academiei Române - Grupul de Prelucrarea Vorbirii, Universitatea "Al. I. Cuza" Iași - Facultatea de Informatică și Universitatea Tehnică "Gh. Asachi" Iași - Centrul de Excelență în cercetare "CERFS" (coordonat de primul autor) în ideea realizării unui suport pentru un "dicționar al sunetelor" limbii române. Situl cuprinde, pe lângă sunetele propriu-zise (vocale, consoane, diftongi, scurte fraze), informații despre fonetica limbii române, protocoale de documentare și de înregistrare, instrumente de analiză, trimiteri la lucrări referitoare la prelucrarea limbii vorbite etc.

Scopul acestei inițiative a fost realizarea unei arhive pentru sunetele limbii române cu următoarele caracteristici:

- Bază de date cu voci atât profesionale (pronunții "perfecte"), cât și ne-profesionale ("vocea omului de pe stradă"), din zona Iași, iar apoi, pe cât posibil, cu pronunții (sunete) specifice diverselor regiuni.
- Pe baza acestor date, realizarea unui studiu statistic amplu al sunetelor limbii române, care să includă, de exemplu "triunghiul formanților limbii române", "caracteristici statistice ale pronunțiilor regionale" etc.
- Corectarea unor sisteme de sinteză ne-concatenativă, pe baza datelor din arhivă.
- Îmbunătățirea unor sisteme de recunoaștere acustică, pe baza datelor din arhivă. (Arhiva ar putea deveni, sperăm, un banc de probă pentru asemenea sisteme).

- Realizarea unei baze de date (sunete) a limbii vorbite pentru persoane cu diverse patologii (neurologice, laringiene, nazale, bucale, respiratorii) (Teodorescu et al., 2006a). Această bază de sunete produse de voci patologice va fi parțial utilă și în cercetările pentru un grant CEEEX.
- Realizarea unei baze de date de tip silabe și cuvinte (sursă pentru sintetizoare concatenative și banc de probă pentru sisteme de recunoaștere de cuvinte).
- Un dicționar electronic al pronunțiilor din limba română. Corelarea cu Atlasele limbii române.
- Pagina de referințe, care să prezinte toate titlurile de volume și lucrări ce au ca obiect sunetele limbii române (fonetică, sinteză, recunoaștere etc.).

S-au efectuat peste 600 de înregistrări, în diverse formate de precizie și codare (Teodorescu et al., 2006b). Fiecare vorbitor a rostit de trei ori fiecare frază, propoziție, cuvânt, vocală, consoană, diftong, triftong, hiat precum și grupuri de sunete specifice limbii române (ex. *ce, ci, che, chi* etc). Înregistrările au fost efectuate folosind programul Goldwave 5.0 la o frecvență de eșantionare de 22050Hz, codate pe 16 și 24 biți, mono.

Fișierele sunt grupate în clase, după cum urmează:

A. Sunete de bază: i) fișiere de vocale; ii) fișiere de consoane, înregistrate conform standardului IPA, în forma VCV, unde V este vocala *a*; iii) fișiere de diftongi, triftongi și hiatusuri; iv) fișiere de sunete specifice, care în scrierea în limba română corespund grupurilor *ce, ci, che* (ke), *chi* (ki), *ge, gi, ghe, ghi*. Subiecții sunt atât bărbați, cât și femei, persoane cu vârsta cuprinsă între 26-31 ani, proveniți (născuți și educați) din zona Moldovei de mijloc (județele Iași, Vaslui, Bacău), cu educație superioară și fără patologii manifestate. Vocalele sunt înregistrate atât în varianta scurtă (pronunție uzuală), cât și în varianta „susținută”.

B. Scurte propoziții sau segmente de fraze, cu încărcătură emoțională diferită

Pe lângă sunete simple au fost înregistrate și fraze scurte: *Cine a făcut asta, Vine mama, Aseară*. Fiecărui subiect i s-a cerut să rostească fiecare frază simulând următoarele emoții: fericire, tristețe, bucurie, ură, optimism, pesimism, ton exclamativ, ton interogativ, ton plat și starea de supărare. Ulterior, aceste stări au fost reduse la patru: fericire, supărare, furie și ton neutru.

2. Metodologia de înregistrare

Analiza zgomotului de fond

Zgomotul de fond este un semnal aleator, arareori staționar, mixat cu unele perturbații determinate de tipul perturbațiilor de frecvență ale rețelei. Pentru a asigura calitatea înregistrării, amplitudinea zgomotului trebuie să fie mult mai mică decât amplitudinea semnalului. Înregistrările au fost efectuate într-un laborator cu zgomot redus, dar nu s-a dispus de un spațiu total izolat fonic. O înregistrare de bună calitate are un zgomot în al cărui spectru nu se depășește valoarea de -80dB pentru nici o componentă spectrală, în timp ce formantul F1 are un nivel cu cel puțin 30dB mai mare, iar formanții superiori au

tipic amplitudini cu peste 15-20dB peste nivelul zgomotului (Teodorescu, 2006a). Fișierele sunt în curs de verificare, urmând a fi eliminate cele care nu satisfac nivelul de calitate pe care ni l-am impus.

Alegerea microfonului

Microfonul este primul element în convertirea sunetelor și are un rol esențial în calitatea înregistrărilor. Un microfon de bandă limitată, de sensibilitate redusă sau cu zgomot mare poate compromite înregistrarea. Un microfon de calitate are zgomot redus și raportul semnal / zgomot bun. Caracteristica omnidirecțională conduce la o sensibilitate mare la zgomotele ambientale și nu este de dorit. Sunetele au fost înregistrate folosind căști cu microfon SONIC Stereo Dynamic Headphones HP-259 cu caracteristicile: frecvența de răspuns: 20-20.000 Hz; impedanța microfon: $U=3V$, $R=1,5K \Omega$; impedanța căști: 32Ω ; sensibilitate microfon: $-58dB \pm 2$; sensibilitate căști: $100dB/mw$; putere: $100mW$.

Poziționarea microfonului

O atenție specială trebuie acordată poziției microfonului, deoarece apar zgomote sau distorsiuni introduse prin poziționarea deficitară. Ținerea microfonului prea aproape de gură poate duce la efectul de saturare a amplificatorului, cu rezultat de puternică distorsionare a semnalului. Se recomandă menținerea microfonului mai jos de gură, aproximativ în dreptul bărbiei, la câțiva centimetri de aceasta. Distanța de la bărbie trebuie să fie aproximativ egală cu distanța până la buze.

Placa de sunet și driverele corespunzătoare

Majoritatea calculatoarelor actuale conțin pe placa de bază circuite de preluare a semnalelor de la microfon și de generare de semnale audio la căști sau difuzoare (calitatea acestor circuite diferă substanțial de la o placă la alta). Placa de bază a calculatorului pe care au fost efectuate înregistrările este MB FOXCONN 760 GXK8MC-S, având încorporată o placă de sunet Sound MAX Digital Audio produsă de Analog Devices cu caracteristicile: procesor de semnal SiS964, standard AC '97.

3. Adnotări

Sunetele și frazele disponibile pe sit au fost adnotate la diferite niveluri cu scopul de a avea un corpus pentru analiza statistică a datelor. Adnotarea s-a realizat folosind utilitarul Praat™ (www.praat.org), ales datorită eficienței, recunoașterii internaționale și ușurinței în folosire. Primul pas a fost segmentarea la nivel de fonem. Ulterior, au fost grupate fonemele pentru a se realiza o segmentare în silabe, cuvinte și propoziții. Principala problemă a fost faptul că este dificil de stabilit întotdeauna cu exactitate unde se află granița dintre foneme. În figura 1 este prezentat un exemplu de adnotare pentru propoziția *Vine mama!*. În viitor, adnotările vor fi validate prin analiza efectuată de mai mulți adnotatori.

Informația de la nivel fonologic va fi completată cu informații prozodice (ton, intensitate, durată etc.), deoarece unul dintre obiectivele noastre este de a detecta parametri prozodici care fac diferența dintre vorbirea umană și cea sintetică.

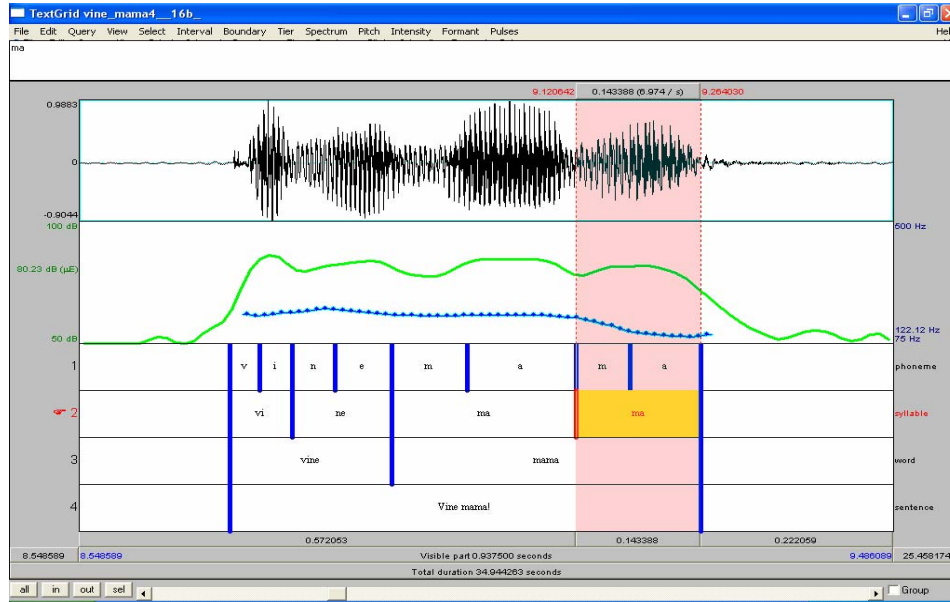


Figura 1: Exemplu de adnotare a propoziției *Vine mama!*

4. Alte elemente ale resursei: documentații, instrumente etc.

Subiecții au fost informați anterior înregistrărilor despre obiectivele proiectului, fiind asigurați de confidențialitatea datelor personale. Subiecții au semnat un consimțământ informat în conformitate cu „Protocolul de Protecție a Subiecților Umani” al *U.S. Food and Drug Administration* (http://www.fda.gov/cdrh/devadvice/ide/informed_consent.shtml) și cu „Principiile etice ale Asociației Acustice Americane privind cercetările care implică ființa umană” (<http://asa.aip.org/ethical.html>). Vorbitorilor li s-au explicat în prealabil condițiile de înregistrare: poziția microfonului, susținerea vocalelor pe o durată cât mai lungă, dar fără a se realiza vreun efort etc. De asemenea, fiecare subiect a completat o fișă personală, care include date despre vârsta, sexul, limba, educația vorbitorului, patologii, precum și evaluarea calității subiective a vocii. Pe lângă sunetele propriu-zise și fișele vorbitorilor, situl mai conține și instrumente de analiză a semnalului vocal.

5. Metode suplimentare de analiză: măsurări accelerometrice

Printre elementele specifice introduse de grupul nostru pentru analiza procesului vorbirii, se afla metoda accelerometrică. Metoda constă în determinarea accelerațiilor mandibulei în timpul vorbirii și corelarea mișcărilor cu tipul sunetelor pronunțate, cu duratele sunetelor și cu energia sonoră. Deși gradul de corelare constatat până în prezent este relativ redus, sperăm că metoda poate ajuta la segmentarea automată a semnalului vocal și la evidențierea unor corelații între caracteristicile pronunției cu mișcărilor

fonatorii. Fără să fie absolut nouă (există un număr mic de lucrări cu abordări oarecum similare, dar nu în scopul segmentării), abordarea sperăm să aducă elemente suplimentare în explicarea proceselor vorbirii. Rezultatele preliminare obținute (Teodorescu 2006a, b) evidențiază caracteristici ale tranzițiilor sunet nazal - vocală, *o-a* în diftongul *oa*, precum și consoană plozivă-vocală.

6. Concluzii și direcții viitoare

Considerăm că resursa este utilă în prezent ca mijloc educațional, iar în viitorul apropiat ca suport în cercetarea lingvistică și în realizarea de aplicații informatice (sinteză și recunoaștere). Credem că resursa impune și noi standarde de calitate în realizarea unor resurse similare.

Arhiva sunetelor limbii române va fi dezvoltată prin adăugarea de noi înregistrări și de prelucrări statistice ale sunetelor. Ulterior, se va urmări adăugarea unor înregistrări cu ușoare patologii, cum ar fi tremurul vocii (de natura emoțională sau patologică), adnotate și prelucrate (Teodorescu et al., 2006c).

Referințe bibliografice

- Teodorescu H.N. (2006a). Gnatofonia și Gnatosonia. Analiza semnalelor vocale, Capitolul 2, *Ed. Performatica*, Iași, România, pag. 29-40.
- Teodorescu H.N. (2006b). Gnatophonetics – A New Discipline Analyzing Relations between Speech and the Stomato-Gnathic System. *Zilele Academice Ieșene, Simp Inventica. Simpozionul național "Bazele performanței și inventică"* organizat în cadrul "Zilelor Academice Ieșene" ISBN 973-730-244-3, 978-973-730-244-1, 9 September 2006.
- Teodorescu H.N., Zbancioc M., Mihăilescu E. (2006a), Speech Technology and Bio-Medical Engineering Teaching Based on the Web – A New Tool and Case Study. *Conference ICL 2006*, Villach, 27 -29 September 2006, Proceedings CD 2005 Ambient and Mobile Learning, Kasset University Press, Editors Michael Auer, Ursula Auer and R. Mittermeir, ISBN 3-89958-136-9.
- Teodorescu H.N., Tandabăț D., Feraru M., Zbancioc M., Luca R.(2006b). A corpus of the sounds in the Romanian spoken language for language-related education. *International Conference on Human and Material Resources in Foreign Language Learning – RFL 2006*, Murcia, Spania, 12-14 iulie 2006.
- Teodorescu H.N., Feraru M., Tandabat D. (2006c), Nonlinear Assessment of Professional Voice 'Pleasantness', Conference *BIOSIGNAL 2006*, ISBN 80-214-3152-0, Brno, 28-30 June 2006, pag. 63-66.
- Voiced Sounds of Romanian Language Project. [http://iit.iit.tuiasi.ro/romanain_spoken_language/index.htm].

SCHEMĂ XML DE ADNOTARE A INTONAȚIEI ÎN CADRUL CORPUSURILOR DE TEXT

VASILE APOPEI, DOINA JITCĂ

Institutul de Informatică Teoretică, Academia Română, Filiala Iași

{vapopei, jdoina}@iit.tuiasi.ro

Rezumat

În lucrarea se propune o schemă de adnotare a unui corpus de text în format XML, cu informație prozodică rezultată din analiza rostirilor respectivului text. În secțiunea 2 se prezintă ierarhia unităților intonaționale pe care s-a bazat structurarea textului și setul de etichete folosit în adnotarea evenimentelor tonale din conturul F0. În secțiunea 3 se prezintă schema XML de adnotare a intonației în cadrul unui corpus de text, prin prezentarea tagurilor și a atributelor acestora. În secțiunea 4 se prezintă un exemplu ce ilustrează corespondența între evenimentele de intonație marcate pe conturul F0 și structurarea XML a aceleași informații.

1. Introducere

Lucrarea prezintă o modalitate de introducere a nivelului intonațional în adnotarea XML a corpusurilor de text, structurate deja la nivel morfologic și sintactic. O structurare multinivel a unui corpus de text este realizată în cadrul proiectului MULI la care adnotarea s-a efectuat pe 3 nivele: sintactic, prozodic, discurs, în vederea intercorelării trăsăturilor corespunzătoare acestora (Baumann et al. 2004). În adnotarea intonației autorii au avut în vedere trei unități de structurare: unitate intonațională, unitate intonațională intermediară și cuvântul, care se regăsesc și în adnotarea noastră. Evenimentele intonaționale din cadrul acestora, la fel ca și în aplicația noastră, s-au adnotat folosind o variantă a sistemului de etichete ToBI (GToBI).

Analiza intonației asupra corpusului de voce s-a efectuat din perspectiva modelului autosegmental-metric concretizat într-o structură ierarhică a unităților intonaționale (prezentată în secțiunea 2) stabilită pe baza structurilor folosite de autori precum: Selkirk (1984), și Di Cristo (2004).

Pentru marcarea evenimentelor intonaționale, de pe conturul frecvenței fundamentale F0, din cadrul unităților intonaționale, am folosit etichete ale sistemului ToBI de adnotare a intonației: (Beckman, Ayers, 1997).

În afară evenimentelor luate în considerație de sistemul ToBI s-au marcat și alte tonuri semnificative din conturul F0 folosind etichetele H+ și L+. Aceste tonuri se pot afla fie pe silaba anterioară unei silabe accentuate, fie pe silaba următoare.

Proiectarea schemei XML de structurare a unui text din punct de vedere al intonației s-a bazat pe analize anterioare referitoare la structurarea diferitelor rostiri în unități

intonaționale, în diverse contexte de accentuare a cuvintelor, atât în propoziții afirmative cât și interogative (Apopei et al., 2006).

Textul folosit în realizarea corpusului a fost extras din romanul “1984” al autorului George Orwell și rostit de doi vorbitori. Adnotarea rostirilor s-a realizat cu un program dezvoltat la Institutul de Informatică Teoretică Iași iar informația rezultată a fost apoi convertită în format XML, folosind schema de adnotare prezentată în secțiunea 2.

2. Prezentarea ierarhiei unităților intonaționale

Adnotarea prozodică pe care o propunem se bazează pe ierarhia unităților intonaționale prezentată în figura 1. În cadrul acestei ierarhii, cea mai mică unitate căreia i se poate asocia un eveniment din conturul frecvenței F0 este silaba. Silabele constituie părți componente ale cuvintelor. Cuvintele sunt purtătoare ale accentelor sintactice sau gramaticale. În general, unitățile de accentuare cuprind un cuvânt cu accent și unul sau mai multe cuvinte clitice. Există situații în care unitățile de accentuare pot include pe lângă cuvântul accentuat, un alt cuvânt neclitic, dar care și-a pierdut complet accentul în vecinătatea acestuia.

O unitate de accentuare purtătoare de accent puternic se grupează în cadrul acestei ierarhii cu alte unități care includ cuvinte purtătoare de accente mai slabe, formând unități ritmice (Di Cristo, 2004). Grupările cuvintelor realizate de unitățile ritmice corespund la nivel semantic, sintagmelor.

Una sau mai multe unități ritmice compun o frază intonațională (intonational phrase, în limba engleză, și notată, IP) sau o frază intonațională intermediară (intermediate phrase, în limba engleză și notată, ip).

Diferența dintre unitățile IP și cele notate ip o constituie durata tonurilor finale și a pauzei de după acestea. În delimitarea unităților IP/ip se ține cont de modificările tendințelor în evoluția tonurilor țintă asociate evenimentelor din cadrul conturului frecvența F0. Tendințele pot fi descrescătoare (downstepping) sau crescătoare (upstepping). De exemplu, când frecvența F0 părăsește tendința de coborâre începută, înainte de a atinge un final al frazei intonaționale se consideră că s-a încheiat o frază intermediară .

Evenimentele intonaționale din conturul frecvenței F0 avute în vedere în adnotare sunt următoarele: accentele de pitch, produse pe durata silabelor accentuate (în engleză „Pitch Accent”); tonurile de sfârșit ale frazelor intonaționale intermediare; tonurile graniță ale frazelor intonaționale; alte tonuri semnificative din conturul F0 (în engleză „Target Ton”), care se pot afla fie pe silaba anterioară silabei accentuate, fie pe silaba următoare. Pentru marcarea primelor trei tipuri de evenimente s-au folosit etichetele sistemului de adnotare ToBI, iar pentru ultima categorie s-au adăugat două etichete, H+ și L+, care au fost folosite și în alte aplicații de adnotare prozodică (Baumann et al., 2004).

SCHEMĂ XML DE ADNOTARE A INTONAȚIEI ÎN CADRUL CORPUSURILOR DE TEXT

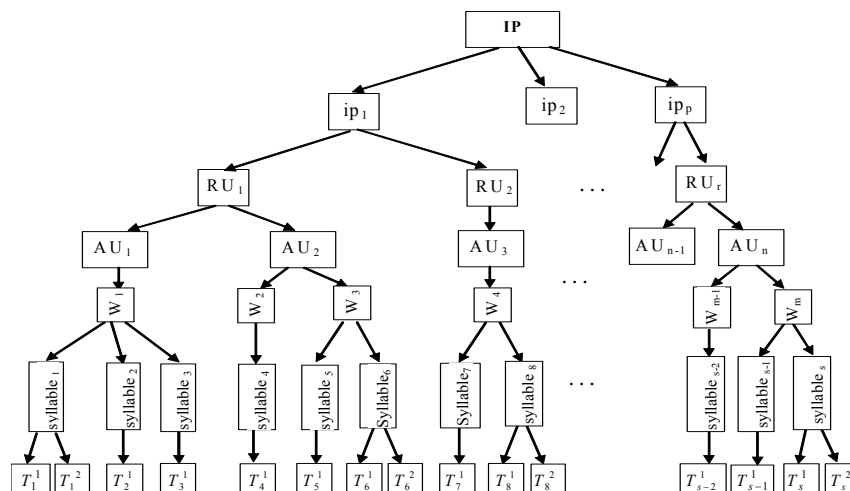


Figura 1: Ierarhia structurii intonaționale

3. Schema XML de adnotare a intonației

În stabilirea tag-urilor pentru adnotarea intonației în format XML s-au avut în vedere toate unitățile din ierarhia prezentată în figura 1, creând câte un tag pentru marcarea unităților de pe fiecare nivel.

Atributele tag-urilor conțin pe lângă secvențe de etichete referitoare la tipul de evenimente și informații cantitative legate de nivelul tonurilor țintă asociate acestora. În acest scop s-a împărțit gama de variație a frecvenței F0, în cadrul rostirii de adnotat, în semitonuri și s-a realizat o scală de măsurare a tonurilor cu baza la nivelul tonului celui mai înalt. Astfel atributele de tip *ToneValues* se exprimă prin numere întregi cuprinse între 0 și 20 corespunzătoare celui mai înalt nivel de ton și respectiv, celui mai scăzut din cadrul rostirii.

Tag-urile și atributele corespunzătoare sunt prezentate în tabelul 1. Unele atribute reprezintă etichete de ton, cum ar fi: cele de început și respectiv, sfârșit ale unităților IP (*BeginToneLabel*, *BoundaryToneLabel*), cele de sfârșit ale unităților intermediare (*PhraseToneLabel*) sau cele legate de accentele cuvintelor și a tonurilor adiacente (*TonalGroupLabel*- când sunt marcate la nivelul cuvintelor și *Accent*, *TargetTone* -când sunt marcate la nivelul silabelor).

Secvența de numere întregi, separate prin virgulă, ce constituie valori pentru atributele de tip *ToneValues* corespunde tonurilor marcate prin etichete în cadrul atributelor aceluiași unități. Sincronizarea secvenței de etichete stabilită prin atributul *TonalGroupLabel* cu secvența de silabe a cuvântului, se realizează prin asocierea etichetei de accent de pitch cu silaba accentuată. În cazul cuvintelor compuse cu mai multe accente se impune marcarea separată a fiecărei silabe din cadrul acestora, cu tag-ul *<Syllable>*. În cadrul unui astfel de cuvânt vor exista silabe marcate cu atributele *Accent* și *ToneValues*, silabe marcate cu atributele *TargetTone* și *ToneValues*, și silabe fără nici o indicație de ton. Oricare silabă poate fi caracterizată prin nivelul de energie, sau printr-o măsură a duratei acesteia folosind atributele *Energy* și respectiv, *Length*.

Tag-ul <RU> nu are atribute pentru că unităților ritmice nu conțin tonuri semnificative diferite de cele ale subunităților componente (unități de accentuare, cuvinte, silabe).

Tabel 1: Atributele și valorile tag-urilor de adnotare prozodică

Tag	Atribut	Valoare	Unitate intonațională
<IP>	<i>BeginToneLabel</i>	%L, %H	frază intonațională
	<i>BoundaryToneLabel</i>	L%, H%	
	<i>ToneValues</i>	Sir de numere întregi	
<ip>	<i>PhraseToneLabel</i>	L-, H-	frază intonațională intermediară
	<i>PhraseToneValue</i>	Număr întreg	
<RU>			unitate ritmică
<AU>	<i>Break</i>	No, Short, Large	unitate de accentuare
	<i>PunctSign</i>	/, !: / ; !: / ! / ? /	
<W>	<i>TonalGroupLabel</i>	Secvențe de etichete	Cuvânt
	<i>ToneValues</i>	Sir de numere întregi	
	<i>ID W</i>	Numeric	
<Syllable>	<i>Accent</i>	H*, L*, L+H*, H+!H*, H+L*	Silabă
	<i>TargetTone</i>	H+, L+	
	<i>ToneValues</i>	Sir de numere întregi.	
	<i>Length</i>	Small, Medium, Large	
	<i>Energy</i>	Small, Medium, Large	

Atributul “*Break*” indică absența sau prezența pauzei după unitatea de accentuare. Prezența pauzei este marcată în termenii *Short* sau *Large* după cum aceasta este de durată mai scurtă sau mai lungă. Valoarea implicită este *No*, adică absența pauzei.

Atributul “*PunctSign*” indică prezența unui semn de punctuație în text după ultimul cuvânt din unitatea de accentuare.

Atributul *ID_W* al tag-ului <W> poate fi folosit pentru a face legătura cu alte tipuri de structuri ale aceluiași text cum ar fi cele care marchează categoriile morfologice, sintactice sau de discurs. Folosind Tag-urile prezentate în acest paragraf, se poate marca un text cu informație relativă la intonația unei rostiri a acestuia.

4. Exemplu de adnotare a intonației în format XML

Exemplificarea adnotării intonației se va face prin corelarea conturului F0 corespunzător rostirii textului „*Avem de discutat lucruri serioase, zece minute nu-i nevoie să mai faci pe valetul*” (figura 2) cu marcarea textului folosind categoriile XML prezentate în secțiunea 3 (figura 3). Intonația este formată din două fraze intonaționale. În cadrul primei unități *IP* se formează două unități ritmice iar în cadrul celei de a doua unității *IP* se formează două fraze intermediare *ip* ce cuprind cele trei unități ritmice.

SCHEMĂ XML DE ADNOTARE A INTONAȚIEI ÎN CADRUL CORPUSURILOR DE TEXT

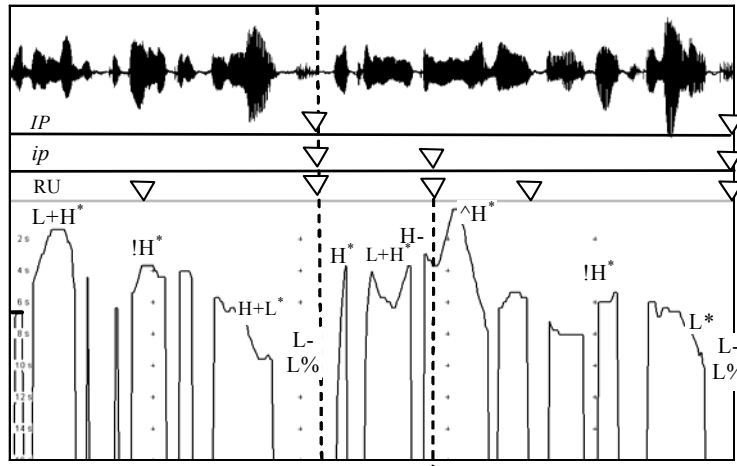


Figura 2. Unda vocală și conturul F0 al rostirii

“Avem de discutat lucruri serioase, zece minute nu-i nevoie să mai faci pe valetul”

Delimitările realizate pe corpusul de voce, corespunzătoare unităților intonaționale ale rostirilor, au fost aplicate textului într-un fișier în format XML generând structurarea acestuia din punct de vedere intonațional.

```
<IP BeginTonLabel="%L" BoundaryTonLabel="L%">
  <RU>
    <AU>
      <W TonalGroupLabel="L+H*" ToneValues="4,1">Avem</W>
    </AU>
    <AU>
      <W>de </W>
      <W TonalGroupLabel="!H*" ToneValues="3">discutat</W>
    </AU>
  </RU>
  <RU>
    <AU>
      <W>lucruri</W>
    </AU>
    <AU Break="Short" PunctSign=",">
      <W TonalGroupLabel="H+L*" ToneValues="8,10">serioase</W>
    </AU>
  </RU>
</IP>
<IP BeginTonLabel="%L" BoundaryTonLabel="L%">
  <ip PhraseTonLabel="H-">
    <RU>
      <AU>
        <W TonalGroupLabel="H*" ToneValues="3">zece</W>
      </AU>
      <AU>
        <W TonalGroupLabel="L+H*" ToneValues="6,3">minute</W>
      </AU>
    </RU>
  </ip>
  <ip PhraseTonLabel="L-">
    <RU>
      <AU>
        <W TonalGroupLabel="^H*" ToneValues="0">nu-i</W>
      </AU>
      <AU>
        <W>nevoie</W>
      </AU>
    </RU>
    <RU>
      <AU>
        <W>să</W>
        <W>mai</W>
        <W TonalGroupLabel="!H*" ToneValues="6">faci</W>
      </AU>
    </RU>
  </ip>
</IP>
```

```
</AU>  
<AU Break="Large" PunctSign=".">  
<W>pe</W>  
<W TonalGroupLabel="L*" ToneValues="9">valetul</W>  
</AU>  
</RU>  
</ip>  
</IP>
```

Figura 3. Adnotarea intonației rostirii textului

„Avem de discutat lucruri serioase, zece minute nu-i nevoie să mai faci pe valetul”

5. Concluzii

Schema de adnotare prezentată în această lucrare, a fost dezvoltată în scopul de a realiza un corpus de text în limba română cu informație relativă la intonația rezultată din rostirea acestuia. Alegerea unor fragmente din romanul “1984” al autorului George Orwell este justificată de existența unor adnotări la nivel morfologic și sintactic pentru acest text. Folosind identificatorul de cuvânt <ID> aceste resurse pot fi alinate și pe baza lor se pot îmbunătăți sau dezvolta noi aplicații lingvistice pentru limba română.

Corpusurile de text adnotate la nivel intonațional sunt utile în deducerea de reguli ce vizează implementarea intonației în sinteza vocală în limba română.

Referințe bibliografice

- Albert Di Cristo (2004). La prosodie au carrefour de la phonétique, de la Phonologie et de l’articulation formes-fonctions, *Travaux Interdisciplinaires du Laboratoire Parole et Langage*, vol. 23, p. 67-211
- Beckman M., Ayers G. (1997). *Guidelines for ToBI Labelling* (version 3), [http://www.ling.ohio-state.edu/~tobi/ame_tobi/labelling_guide_v3.pdf].
- Baumann S, Brinckmann C., et al (2004). Multi-dimensional annotation of linguistic corpora for investigating information structure, *In: Proceedings Frontiers in Corpus Annotation Workshop at HLT/NAACL*, Boston, USA, p. 39-46.
- Apopei V, Jitcă D, Turculeț A. (2006). Intonational structures in Romanian Yes-No Questions, *Jurnal of Computer Science of Moldavia*, pp. 113-137, vol 14, no 1(40), Chișinău
- Selkirk, E.O. (1984). *Phonology and Syntax: The Relation between Sound and Structure*. Cambridge, MA: MIT Press.

Capitolul 2

Dicționare și corpusuri adnotate pentru prelucrarea textelor

NOI DEZVOLTĂRI ALE WORDNET-ULUI ROMÂNESC

DAN TUFIȘ, VERGINICA BARBU MITITELU, ALEXANDRU CEAUȘU, LUIGI BOZIANU, CĂTĂLIN MIHĂILĂ, MARGARETA MANU MAGDA

Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București

{tufis, vergi, alceausu, bozi, cata}@racai.ro; mmagda@csb.ro

Rezumat

Ontologiile lexicale de tip wordnet sunt dintre cele mai importante resurse lexicale folosite în aplicațiile de prelucrare a limbajului natural. O astfel de ontologie este disponibilă și pentru limba română. Lucrarea descrie o parte din activitățile de îmbunătățire cantitativă și calitativă a wordnet-ului românesc.

1. Dezvoltarea wordnet-ului românesc

Una dintre cele mai importante resurse lingvistice computaționale pentru limba română este, fără îndoială, ontologia lexicală de tip WordNet (Fellbaum 1998), a cărei dezvoltare a început, în colaborare cu Facultatea de Informatică a Universității Al. I. Cuza din Iași, în anul 2001, în cadrul proiectului BalkaNet¹. Dezvoltarea wordnet-ului românesc a continuat la Institutul de Cercetări pentru Inteligență Artificială și după anul 2004 (Tufiș et al., 2006), când proiectul BalkaNet s-a încheiat.

Experiența acumulată în timpul proiectului BalkaNet, precum și o serie de noi instrumente de achiziție lexicală au făcut ca productivitatea echipei de lingviști să crească substanțial, astfel încât, actualmente, wordnet-ul românesc are 33422 de sinseturi (i.e. serii sinonimice), conținând un număr de 31246 de literalii unici.

Tabel 1: Date statistice despre wordnet-ul românesc

33421 serii sinonimice (1289 nelexicalizate)	Relații semantice:
<ul style="list-style-type: none">• 24640 substantive• 7096 verbe• 851 adjective• 834 adverbe	<ul style="list-style-type: none">• hypernym 32041• holo_part 2096• holo_member 1029• holo_portion 199• category_domain 1861• also_see 508• similar_to 899
53160 literalii (31246 literalii unici)	Relații lexicale:
163 domenii	<ul style="list-style-type: none">• near_antonym 1976• be_in_state 566• verb_group 1196• causes 148• subevent 264
1773 categorii SUMO/MILO	

Principala strategie de dezvoltare a wordnet-ului românesc constă în implementarea în limba română a seriilor sinonimice din wordnet-ul englezesc. Pentru selectarea

¹ <http://www.ceid.upatras.gr/Balkanet/>

sinseturilor ce urmau a fi implementate au fost urmărite criteriile ce țin de acoperirea cu măcar un sens pe literal a cuvintelor din corpusul Acquis-ului comunitar. Ca modalități de stabilire a relevanței unui sinset au fost considerate: (i) domeniul și caracterizarea SUMO a seriei sinonimice (fiind avantajate domenii ca politic, legislativ etc., specifice Acquis-ului); (ii) numărul de apariții ale literalilor în corpus; iar în cazul literalilor monosemantici, (iii) rangul de ocurență întors de motorul de căutare Google.

Folosind metoda de dezambiguizare automată descrisă în (Ion, Tufiș, 2004), dar, de data aceasta, pe un bitext diferit (SemCor), s-au depistat circa 7.000 de literalii absenți din sinseturile deja implementate. Actualizarea acestor sinseturi incomplete este în curs de realizare (până la data scrierii acestui articol au fost completate 600 de sinseturi), în paralel cu adăugarea de noi sinseturi.

2. Alinierea wordnet-ului românesc la versiunea PWN 2.1

În timpul scurs de la încheierea proiectului BalkaNet, cercetătorii de la Princeton au dezvoltat versiunea 2.1 a wordnet-ului american, PWN, care a adus o serie de modificări majore față de versiunea precedentă. Deși numărul de sinseturi nu a crescut substanțial, a apărut distincția între instanță și clasă, iar mai multe relații au fost redenumite sau rafinate, ceea ce a condus la necesitatea regenerării formatului XML pentru PWN folosit de VisDic (Horak, Smrz, 2004) în cadrul proiectului BalkaNet.

După cum se știe, wordnet-ul românesc este aliniat cu Princeton WordNet (PWN), cu ontologia SUMO/MILO (Niles, Peace, 2003) precum și cu taxonomia DOMAINS (Bentivogli et al., 2004). Au fost corectate mai multe erori de mapare între PWN2.0 și SUMO/MILO și DOMAINS.

Ontologia WN DOMAINS 2.0 are, față de versiunea 1.0, noi etichete de domeniu, iar unele dintre domenii au suferit o schimbare de granularitate. Aceste modificări au fost operate pentru a asigura o distribuție uniformă a etichetelor de domeniu pentru seriile sinonimice ale wordnet-ului englezesc.

2.1. Generarea în format XML a Princeton WordNet versiunea 2.1

Față de versiunea XML a PWN2.0, versiunea 2.1 include și numărul de ocurențe ale sensurilor adnotate în corpusul SemCor. În plus, la nivelul relațiilor lexicale au fost specificați literalii între care acestea se stabilesc. Deși specificarea acestor literalii s-a făcut încă de la versiunea 1.5 a Princeton WordNet, ea lipsea din versiunea XML a PWN 2.0.

Interfața de vizualizare/editare VisDic (Horak, Smrz, 2004), interfață dezvoltată în cadrul proiectului EuroWordNet², folosește un format XML propriu de reprezentare a wordnet-ului. Versiunea XML a PWN 2.1 are în plus tagurile <FREQ>, <SRCL> și <TRGL> pentru reprezentarea numărului de ocurențe din SemCor și pentru specificarea literalilor între care se stabilesc relațiile lexicale:

² <http://www.illc.uva.nl/EuroWordNet/>

NOI DEZVOLTĂRI ALE WORDNET-ULUI ROMÂNESC

```

<SYNSET>
  <ID>ENG21-06551177-n</ID>
  <POS>n</POS>
  <SYNONYM>
    <LITERAL>wordnet<SENSE>1</SENSE><FREQ>1</FREQ></LITERAL>
  </SYNONYM>
  <DEF>any of the machine-readable lexical databases modeled
  after the Princeton WordNet</DEF>
  <ILR>ENG21-06550617-n<TYPE>hypernym</TYPE></ILR>
  <DOMAIN>computer_science</DOMAIN>
  <SUMO>ContentBearingObject<TYPE>+</TYPE></SUMO>
</SYNSET>

```

Tabel 2: Date statistice despre Princeton WordNet 2.1

117597 serii sinonimice <ul style="list-style-type: none"> • 81426 substantive • 13650 verbe • 18877 adjective • 3644 adverbe 	155327 literali <ul style="list-style-type: none"> • 117097 substantive • 11488 verbe • 22141 adjective • 4601 adverbe
Relații semantice: <ul style="list-style-type: none"> • hypernym 88258 • instance 8515 • holo_part 8874 • region_domain 1327 • usage_domain 1258 • category_domain 6534 • holo_portion 793 • holo_member 12262 • also_see 3272 • similar_to 22622 	Relații lexicale: <ul style="list-style-type: none"> • derived 8065 • eng_derivative 71914 • near_antonym 8029 • verb_group 1748 • particle 124 • be_in_state 1286 • subevent 409 • causes 219

2.2. Alinierea wordnet-ului românesc la versiunea PWN 2.1

După generarea versiunii XML a PWN2.1, alinierea wordnetului românesc la noua versiune s-a realizat folosind ca resurse WN-Map (Daudé et al., 2000) și maparea pentru substantive și verbe disponibilă pe site-ul Princeton WordNet.

WN-Map folosește un algoritm iterativ pentru optimizarea unei funcții bazate pe un set de criterii ce descriu un context local. Criteriile pot fi eticheta morfologică a sinsetului, definiția sinsetului, locul pe care îl ocupă acesta în ierarhie etc. De asemenea, aplicația WN-Map mai folosește și reguli care pot decide asupra compatibilității sau incompatibilității candidaților la aliniere, reguli bazate pe criteriile enunțate anterior.

În cazul sinseturilor pentru care ambiguitatea de mapare nu a putut fi rezolvată automat (destul de puține, de altfel), dezambiguizarea s-a făcut manual, de către experții lingviști implicați în proiect.

Problemele de aliniere s-au datorat, în principal, modificărilor operate asupra sinseturilor adjectivale și adverbiale. O altă problemă de aliniere a provenit din transferul conceptelor asignate indexului interlingual specific țărilor balcanice (Balkanet Interlingual Index – BILI). Aceste concepte nu au echivalent în wordnet-ul englezesc, dar sunt integrate în ontologia acestuia.

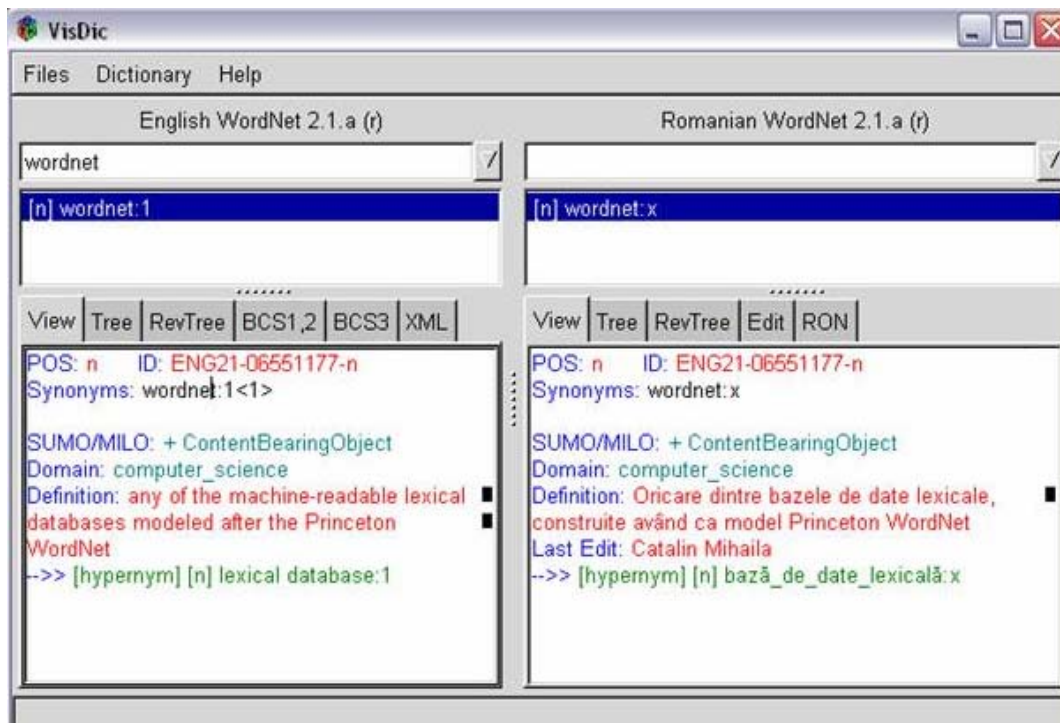


Figura 1: Exemplu de serie sinonimică din PWN 2.1 și din wordnet-ul românesc, vizualizat în VisDic

3. Perspective

Corpusul JRC-Acquis ce cuprinde setul de legi, dispoziții, tratate comune tuturor statelor membre ale Uniunii Europene este unul dintre cele mai mari corpusuri paralele disponibile la momentul actual. Mărimea corpusului și numărul mare de limbi componente îl fac instrumentul perfect de validare a unei ontologii lexicale multilinguale. Astfel, strategia de implementare a seriilor sinonimice din limba engleză în limba română urmărește, în principal, acoperirea cu cel puțin un sens per literal a tuturor cuvintelor din JRC-Acquis.

Îmbunătățirea la nivel cantitativ a wordnet-ului – mai multe serii sinonimice, mai mulți literalii echivalați – este de maximă importanță pentru aplicațiile de traducere automată din și în limba română. Nu a fost, însă, ignorată îmbunătățirea calitativă a wordnet-ului românesc (au fost rezolvate conflictele de asignare a sensurilor, au fost adăugați literalii absenți din unele serii sinonimice, identificați cu ajutorul unui algoritm aplicat asupra unui corpus paralel englez-român, dezamabiguitat la nivel semantic).

De asemenea, folosind corpusuri paralele, sunt în curs de implementare sisteme pentru importul automat de sinseturi conținând instanțe și, respectiv, substantive monosemantice.

Mulțumiri. Autorii sunt recunoscători finanțatorilor proiectului BalkaNet (Comisia Europeană și Ministerul Educației și Cercetării), în cadrul căruia a debutat dezvoltarea wordnet-ului românesc, ai proiectului ROTEL (Ministerul Educației și Cercetării), în cadrul căruia continuă dezvoltarea cantitativă și calitativă a wordnet-ului românesc, precum și Academiei Române, pentru finanțarea temelor de plan în cadrul cărora s-a lucrat și la sporirea calității acestei resurse lingvistice pentru limba română.

Referințe bibliografice

- Bentivogli, L., Forner, P., Magnini, B., Pianta, E. (2004). Revising WordNet Domains Hierarchy: Semantics, Coverage, and Balancing. *Proceedings of COLING 2004 Workshop on "Multilingual Linguistic Resources"*, Geneva, Switzerland, 101-108.
- Daudé, J., Padró, L., Rigau, G. (2000). Mapping WordNets Using Structural Information. *38th Annual Meeting of the Association for Computational Linguistics (ACL'2000)*, Hong Kong, 504-511.
- Fellbaum, Ch. (ed.) (1998). *WordNet: AN Electronic Lexical Database*, MIT Press.
- Horak, A., Smrz, P. (2004). New Features of Wordnet Editor VisDic. *Romanian Journal of Information Science and Technology*, volume 7, Numbers 1-2, 1-13.
- Niles, I., Pease, A. (2003). Linking Lexicons and Ontologies: Mapping WordNet to the Suggested Upper Merged Ontology. *Proceedings of the 2003 International Conference on Information and Knowledge Engineering (IKE '03)*, Las Vegas, Nevada, June 23-26.
- Ion, R., Tufiș, D. (2004). Multilingual Word Sense Disambiguation Using Aligned WordNets. *Romanian Journal of Information Science and Technology*, volume 7, Numbers 1-2, 183-200.
- Tufiș, D., Barbu Mititelu, V., Bozianu, L., Mihăilă, C. (2006). Romanian WordNet: New Developments and Applications. *Proceedings of the 3rd Conference of the Global WordNet Association*, Seogwipo, Jeju, Republic of Korea, January 22-26, 337-344

FRAMENET ROMÂNESC: TENTATIVĂ DE ELABORARE

VICTORIA BOBICEV, VICTORIA MAXIM, TATIANA ZIDRAȘCO, ALINA IACIURINSCHI

Universitatea Tehnică din Moldova

vika@rol.md, maxivica@yahoo.com, tzidrashco@mail.md, ialinca@yahoo.com

Rezumat

În lucrarea de față sunt prezentate rezultatele lucrului efectuat la crearea variantei românești a resurselor multilingve 'Romance FrameNet'. Au fost evidențiate unele întrebări și probleme nerezolvate. În concluzie sunt menționate unele sugestii privind metodologia creării FrameNet-ului multilingv.

1. Introducere

FrameNet (Johnson et al., 2003) este un proiect de cercetare lexicografică inițiat în cadrul International Computer Science Institute Berkeley, California¹. Elementele constructive ale bazei de date FrameNet sunt propozițiile marcate atât semantic, cât și sintactic, fapt ce permite extragerea automată a cadrelor sintactico-semantice. Aceste cadre vizualizează legătura dintre sens și structura sintactică prin care acesta este redat. Ca bază pentru marcarea semantică sunt utilizate cadrele semantice (*frames*) – structuri conceptuale ce reprezintă evenimente, obiecte, proprietăți. Fiecare cadru este dotat cu un set de elemente semantice (*frame elements*) ce caracterizează cadrul dat. Cu fiecare cadru este asociat un număr de cuvinte – unități ale lexiconului (*lexicon units*) care evocă sensul reprezentat de cadrul dat.

În legătură cu creșterea interesului în reprezentarea sensului folosind semantica cadrelor (Gildea & Jurafsky, 2002), FrameNet a fost folosit în exercițiul SensEval-3 (Litkowski, 2004). Sarcina propusă participanților a fost marcarea automată a rolurilor semantice într-un set de propoziții de testare având un subcorpus de propoziții marcate ale Framenet-ului pentru antrenare. Rezultatele celor 20 sisteme care au participat în exercițiu sunt destul de bune (precizia medie 80%), ceea ce caracterizează FrameNet-ul ca o resursă semantică excelentă. Astfel, eforturile depuse în scopul creării acestei resurse au un rezultat de mare valoare.

2. FrameNet în alte limbi

Utilizând principii similare cu Framenet-ul englezesc, s-au creat resurse lingvistice pentru alte limbi, și anume:

- **German FrameNet** este realizat de trei echipe în colaborare:

¹ <http://framenet.icsi.berkeley.edu/>

- Scopul proiectului **SALSA** (Erk et al., 2003) este crearea lexiconului limbii germane cu informația semantică și sintactică bazată pe teoria cadrelor semantice și analiza posibilităților de a folosi lexiconul dat pentru procesarea limbii germane. Etapa precedentă **SALSA I** avea ca scop marcarea semantică a unui corpus în limba germană și cercetarea metodelor de utilizare a lui în procesarea textului.
- Altă echipă (Boas, 2004) folosește baza de date FrameNet, înlocuind în întregime părțile englezești dependente de limbă cu părțile germane.
- În Stuttgart un grup cercetează analiza și extragerea datelor din corpusul adnotat, în special cologații și nominalizări (*nominalizations*).
- **Spanish FrameNet** (Subirats & Petruck, 2003) este un proiect național la care participă un număr de cercetători din diferite universități guvernate de universitatea independentă din Barcelona. FrameNet-ul spaniol se creează în colaborare cu echipa framenet-ului englezesc și se bazează pe semantica cadrelor. Scopul proiectului este adnotarea semantică a propozițiilor spaniole dintr-un corpus de limbă spaniolă.
- **Japanese FrameNet** – JFN (Ohara et al., 2004) are ca scop crearea unui lexicon care înregistrează descrierea valențelor cuvintelor japoneze, bazată pe semantica cadrelor. Scopul final al JFN, creat în colaborare cu echipa FrameNet-ului englezesc, este crearea unei baze de date de tip FrameNet pentru cuvintele japoneze.
- **French FrameNet²** are ca scop crearea lexiconului semantic francez utilizând metodologia originală FrameNet. La fel este creat un corpus de propoziții adnotate conform metodologiei FrameNet. În proiectul dat este studiată intens posibilitatea de utilizare repetată a resurselor semantice lexicale franceze. Numai în acest proiect din cele patru descrise se cercetează posibilitatea de traducere a propozițiilor FrameNet-ului englezesc utilizând resurse valabile de traducere.

După cum se vede din descrierile date, numai în proiectul francez se cercetează posibilitatea de a crea o sursă într-adevăr bilingvă, care să conțină un set de propoziții paralele adnotate folosind metodologia FrameNet. Astfel o sursă se creează în cadrul proiectului chinez (Fung & Chen, 2004), însă rezultatele încă nu au fost publicate.

3. *Romance FrameNet*

Cu scopul extinderii ariei de acoperire a proiectului FrameNet s-a propus inițierea proiectului „Romance FRAMENET”³, care reprezintă o resursă multilingvă pentru limbile romanice (franceză, spaniolă, italiană, română, portugheză, catalană) bazată pe cadrele semantice FrameNet (Lowe et al., 1997). Această inițiativă are ca scop crearea unei resurse care ar reprezenta un corpus paralel, aliniat și adnotat semantic.

Pentru crearea acestei resurse s-a propus o metodologie, asemănătoare metodologiei folosite în crearea lui MultiSemCor (Lupu et al., 2005) care implică câteva etape. Inițial se traduc în mod manual propozițiile din limba engleză în limbile sus-numite. Traducerea este efectuată de echipele din țările corespunzătoare. Apoi în mod automat propozițiile se aliniază la nivel de cuvânt și adnotarea semantică se transferă din propozițiile englezești la propozițiile din altă limbă folosind alinierea efectuată.

² <http://libresource.inria.fr/projects/framenet/>

³ <http://www.icsi.berkeley.edu/~vincenzo/rfn/index.html>

Rezultatul va fi o resursă multilingvă aliniată la nivel de cuvânt și marcată semantic în baza teoriei cadrelor FrameNet-ului englezesc. O astfel de resursă va fi apreciată de cei care studiază limbile implicate, de lingviștii care cercetează diferențele structurale dintre limbile acestea. În afară de aceasta, resursa dată poate fi folosită în traducerea automată și în interpretarea semantică multilingvă. Însă crearea resursei date folosind metodologia propusă necesită un volum mare de adnotari manuale, fapt ce duce la încetinirea obținerii rezultatelor.

4. *FrameNet românesc*

Echipa noastră participă în crearea părții românești a resursei date. Pentru aprecierea metodologiei propuse inițial au fost alese 100 de propoziții din FrameNet-ul englezesc care s-au tradus în limba română și s-au marcat cu rolurile semantice. Pe parcursul lucrului au fost observate următoarele probleme:

- pot exista mai multe variante de traducere;
- poate să nu fie găsită nici o variantă de traducere;
- calitatea transferului automat al marcajelor semantice nu este suficient de bună;
- diferențele lexicale și sintactice dintre limbi cauzează diferite probleme în transferarea marcării semantice.

La următoarea etapă studiul s-a efectuat asupra a 1000 de propoziții, care au fost traduse de traducători profesioniști. Prin intermediul aliniatorului lexical dezvoltat de RACAI (Tufiș et al., 2006) s-a realizat alinierea automată la nivel de cuvânt. Și în final marcarea semantică de la propozițiile englezești s-a transferat automat la cele românești. Pentru transferarea marcării semantice, două persoane au creat în mod independent două programe de transferare (unul în Perl, altul în C#). Procesul de transferare a fost complicat din cauza semnelor diferite în propozițiile date, și anume apostrofuri, semne care reprezintă bani, semne de punctuație care nu coincideau și altele. În multe cazuri au apărut probleme din cauza diferenței între modul de marcare și modul de aliniere. Marcarea fragmentelor de propoziții este executată la nivel de caractere, alinierea este executată la nivel de cuvinte. În afară de aceasta, în cadrul alinierii sunt formate unități de traducere care conțin câteva cuvinte unite cu cu semnul ‘_’. În unele cazuri unirea cuvintelor în propoziția marcată și cea aliniată a fost diferită ce complica procesul de transferare a rolurilor. Rezultatele transferării au fost validate manual. Au fost considerate corecte numai propozițiile în care nu a fost găsită nici o greșală în transferarea rolurilor semantice. Primul experiment, efectuat asupra unui număr de propoziții mai simple, a arătat un rezultat de 59% corectitudine (total propoziții transferate - 424, corecte - 252). Al doilea experiment, care a implicat un număr mai mare de propoziții, a dat rezultatul de 36% (total propoziții transferate - 600, corecte - 219). O parte din propozițiile acestea au fost analizate mai detaliat. În primul rând, în 7% de propoziții transferarea nu a fost efectuată corect. În restul propozițiilor a fost controlată alinierea. În 9% din propozițiile acestea nu era marcat corect cuvântul de bază (*target*). În 31% a fost marcat greșit numai un cuvânt. În majoritatea cazurilor greșit sunt marcate propoziții, pronume sau articole care este problematic de aliniat corect. În multe cazuri un rol în propoziția englezească conținea o parte din propoziție din trei-cinci și mai multe cuvinte. În astfel de cazuri transferarea rolului pe cuvinte deseori îl întreprupea

și o parte din cuvintele date nu au fost marcate. Uneori în fragmentul dat toate cuvintele erau marcate în afară de un articol sau un pronume. În 27% de propoziții au fost observate mai multe erori. Majoritatea din erorile acestea au fost cauzate absența alinierii pentru cuvintele englezești marcate. Astfel, există necesitatea verificării rezultatului procesării automate.

În scopul creării variantei românești a resursei multilingve sus-numite la catedra noastră a fost efectuat un considerabil volum de lucru cu ajutorul studenților. Au fost selectate câteva cadre și repartizate astfel, ca propozițiile dintr-un cadru să fie traduse și marcate de studenții unui grup și în final toate propozițiile cadrului să fie traduse. În cele ce urmează sunt descrise datele referitoare la cadrele prelucrate.

Cadrul: **Removing**; unitățile lexicale cu numărul propozițiilor adnotate: pluck – 38, prised – 3, evacuated – 6, evacuation – 25, purge – 28, remove – 30, extract – 8, oust – 13, expunge – 6, expulsion – 24, evict – 3, eviction – 25, excise – 2, elimination – 13, empty – 16, ejection – 18, eliminate – 21, drain – 2, eject – 19, clear – 15, confiscate – 15; în total – 330 propoziții.

Cadrul: **Sensation**; unitățile lexicale cu numărul propozițiilor adnotate: fragrance – 15, sight – 24, taste – 16, bouquet – 8, incense – 6, reek – 14, savour – 8, whiff – 12, scent – 24, sensation – 27, noise – 34, sense – 16, aroma – 6, odour – 16, perception – 18, stink – 13, feeling – 21, flavour – 30, perfume – 12, smell – 39, sound – 50; în total – 409.

Cadrul: **Commerce**; unitățile lexicale cu numărul propozițiilor adnotate: buyer – 60, purchaser – 31, seller – 37, vendor – 35, retailer – 25; în total -188 propoziții.

Cadrul: **Change-of-phase**; unități lexicale cu numărul propozițiilor adnotate: condense – 13, thaw – 15, evaporate – 20, defrost – 5, solidify – 12, freeze – 21, vaporize – 5, melt – 27, sublime – 1; în total 119 propoziții.

În total – 1047 propoziții.

Apoi propozițiile traduse s-au marcat în mod manual. Cu scopul facilitării lucrului s-a creat un produs soft FRAME STUDIO 1.06 care permite transferarea, marcarea și corectarea rolurilor semantice. Un avantaj al acestei aplicații este o interfață foarte simplă în utilizare și obținerea rezultatului în format XML. Un exemplu al acestui rezultat este prezentat mai jos:

```
<?xml version="1.0" encoding="utf-16" ?>
  <frames>
    <sentence> <frame>Hear</frame>
      <lexunit startpos="22" endpos="25">hear.v</lexunit>
      <markups>
        <markup felement="Hearer" startpos="27" endpos="28" />
        <markup felement="Message" startpos="30" endpos="36" />
      </markups>
      <text>Pat a spus că până să audă ea aceasta, ea nu realizase
      faptul cât de anti-feminist este fratele ei.</text>
    </sentence> . . .
```

În urma lucrului efectuat am constatat că cea mai complicată parte a lucrului este traducerea, care nu poate fi executată automat. Traducerea se complica din cauza faptului că propozițiile sunt rupte din context și în unele cazuri nu este clar sensul lor.

Exemplul 1: nu este clar dacă este vorba despre o persoană sau despre luna aprilie.

We love the April approach : seven good-value , gentle products with a luxurious apricot FRAGRANCE.

În unele cazuri nu este clar din ce tip de text propozițiile sunt extrase. Însă este important de știut domeniul și topica textului, uneori tonalitatea traducerii depinde de tipul textului. În exemplul 2 propoziția poate fi tradusă în diferite modalități:

REMOVE the dumplings with a slotted spoon and serve them the same way as the fried dumplings.

Scoate gogoășele ...Scoateți gogoășele ...Gogoășele se scot ...

După cum am menționat deja, metodologia care include traducerea a fost adaptată în baza creării resursei traduse MultiSemCor. Însă diferența constă în faptul că MultiSemCor conținea texte întregi ce permitea înțelegerea lor mai amplă. Sensul fragmentului alcătuit din propoziții legate este mai clar decât sensul unei propoziții înafară contextului. Totuși, propozițiile din FrameNet au fost scoase din text și inițial au fost scrise fără o intenție să fie percepute înafară contextului lor. Chiar și traducerea pronumelui „it” creează probleme, fiindcă se referă la ceva în afară propoziției date. Astfel, traducerea propozițiilor devine mai complicată decât traducerea textelor și este nevoie să fie executată de traducători cu experiență și cunoștințe profunde în limba engleză. După cum am menționat, nici un proiect de creare a FrameNet-ului descris mai sus nu include traducerea propozițiilor englezești. În cazuri când posibilitatea aceasta este menționată (BiFrameNet și FrameNet franțuzesc) rezultatele nu au fost anunțate.

O problemă nerezolvată rămâne întrebarea dacă traducătorii trebuie să fie informați despre scopul traducerii sau nu. O condiție de bază pentru propozițiile traduse este că propozițiile trebuie să fie corect alcătuite în limba română. Însă în multe cazuri traducătorii interpretează propozițiile destul de liber, reformulând sensul redat de propoziția englezească, și atunci deseori în traducere nu rămân elementele necesare pentru adnotare. Totuși, este posibil de tradus în română făcând cât mai mult posibilă similaritatea cu propoziția englezească.

La catedra noastră traducerile au fost efectuate de către studenți și verificate de profesori. Și studenții și profesorii cunoșteau scopul traducerii. Pe parcursul lucrului am observat că există un număr de propoziții prea complicate. Una din ele este prezentată în exemplul 1. Sensul propozițiilor acestea poate fi interpretat greșit din cauza lipsei contextului. După părerea noastră, astfel de propoziții nu trebuiesc traduse. În schimb, pentru fiecare unitate lexicală a fost ales un număr (10-20) de propoziții românești, ce conțin unitatea lexicală respectivă și au fost marcate cu elementele cadrului respectiv. Nedispunând de un corpus reprezentativ, noi am folosit Internetul pentru obținerea acestor propoziții.

5. Concluzii și discuții

În urma lucrului efectuat am ajuns la unele concluzii și întrebări nerezolvate. După părerea noastră, în metodologia adoptată, cea mai complicată parte a lucrului este traducerea. Pentru ameliorarea calității rezultatelor obținute ar fi bine ca traducerea să fie executată de cel puțin trei traducători în mod separat. Problema selectării variantei de traducere mai convenabile pentru marcarea rămâne deschisă. E posibil ca această alegere să fie efectuată de către adnotator. O problemă nerezolvată rămâne întrebarea dacă

traducătorii trebuie să fie informați despre scopul traducerii sau nu. Propozițiile pentru care traducerile au fost prea diferite sau neclare trebuie excluse din varianta românească. Propozițiile problematice nu trebuie traduse, pot fi traduceri incorecte. Dacă dorim să creăm o sursă calitativă, trebuie să tratăm atent fiecare propoziție. În schimb, pentru fiecare unitate lexicală ar trebui de ales un număr de propoziții românești, ce conțin unitatea lexicală respectivă, cu scopul marcării conform metodei FrameNet.

Referințe bibliografice

- Boas, H. C. (2004) Semantic Frames as an Interlanguage for Multilingual Lexical Databases, *First Global FrameNet Workshop*, ISI, Berkeley, California.
- Erk, K., Kowalski, A., Pado Sebastian and Pinkal Manfred. (2003) Towards a Resource for Lexical Semantics: A Large German Corpus with Extensive Semantic Annotation. *Proceedings of ACL 2003*, Sapporo.
- Fung P., Benfeng C. (2004) "BiFrameNet: bilingual Frame Semantics Resource Construction by Cross-lingual Induction", in *Proceedings of the 20th International Conference on Computational Linguistics (COLING 2004)*, August 2004.
- Gildea, D., and Daniel Jurafsky. (2002) Automatic Labeling of Semantic Roles. *Computational Linguistics*, 28 (3), 245-288.
- Johnson, C., Miriam Petruck, Collin Baker, Michael Ellsworth, Josef Ruppenhofer, and Charles Fillmore, (2003). *FrameNet: Theory and Practice*. Berkeley, California.
- Litkowski, K. C. (2004) Senseval-3 Task: Automatic Labeling of Semantic Roles, *Proceedings of Senseval-3: The Third International Workshop on the Evaluation of Systems for the Semantic Analysis of Text, ACL 2004*, Barcelona, pp. 9-12.
- Lowe, J.B., Baker, C.F. and Fillmore, C.J. (1997): A frame-semantic approach to semantic annotation, In *Proceedings of the SIGLEX Workshop on Tagging Text with Lexical Semantics: Why, What, and How?*, Washington, USA.
- Lupu M., Trandabăț .D., Husarciuc M. (2005) A Romanian SemCor Aligned to the English and Italian MultiSemCor, *1st International ROMANCE FrameNet Workshop*, în cadrul școlii de vară EUROLAN 2005, România, pp. 20-27.
- Ohara, K. H., Fujii S., Ohori T., Suzuki R., Saito H., Ishizaki S. (2004). "The Japanese FrameNet Project: An introduction." *LREC 2004. Proceedings of the Workshop "Building Lexical Resources from Semantically Annotated Corpora"*, p.9-11.
- Subirats, Carlos; Petruck, Miriam. (2003) Surprise: Spanish FrameNet. *International Congress of Linguists. Workshop on Frame Semantics*, Prague, Czech Republic.
- Tufiș Dan, Ion Radu, Alexandru Ceaușu, Dan Ștefănescu. (2006): Improved Lexical Alignment by Combining Multiple Reified Alignments. In *Proceedings of the 11th Conference of the EACL*, Trento, Italy, 3-7 April, 2006, pp. 153-160.

DEI MULTIMEDIA: EVOLUȚII, PERSPECTIVE

DUMITRU TODOROI¹, ADRIAN CHIORESCU²

¹ *Academia de Studii Economice, Chișinău*

² *Universitatea "Al.I.Cuza", Facultatea de Informatică, Iași*

todoroi@ase.md , chiorescu@yahoo.com

Rezumat

Lucrarea prezintă evoluțiile proiectului **DEI Multimedia** de la lansarea ideii în anul 2000 până în prezent: crearea bazei de date, scenarii de utilizare, posibilități de conversie, exemplu de interfață (web), concluzii și perspective de viitor.

1. Introducere

Dicționarul Explicativ Ilustrat Multimedia este rezultatul unui proiect al Academiei de Studii Economice din Moldova denumit „*Romanian Language for the European Community*”.

În cadrul unui lung șir de proiecte de informatizare a limbii române, acest proiect a realizat o *bază de date multimedia* pentru Dicționarul Enciclopedic Ilustrat al Limbii Române (DEI) intitulată **DEI Multimedia**.

Concret, este vorba de o bază de date în SGBD Microsoft Acces 2000. S-a ales acest format datorită posibilităților de stocare a obiectelor (OLE) multimedia cum ar fi imagini, sunete, video.

Structura unui tabel are 7 câmpuri astfel:

cuvânt, **cuvânt_accent**, **definiție** (câmpuri de tip text)

cuvânt_audio, **definiție_audio**, **image**, **video** (câmpuri OLE)

Lucrările asupra creării Sistemului **DEI Multimedia** au început după Conferința din 14 - 15 aprilie 2000 de la Chișinău cu inițiativa, planificarea, organizarea și coordonarea lucrărilor de către Prof. univ. dr. hab. Dumitru TODOROI. Lucrările au fost planificate și organizate în activități pe grupe, compuse din profesori, doctoranzi, masteranzi și studenți, în total activând până în prezent 53 de persoane.

Grupul TEXT: Zinaida Todoroi, Claudia Vasilache, Ion Linga și un colectiv compus din 18 doctoranzi, masteranzi și studenți de la ASEM. Acest grup a realizat **componenta TEXT** (preluarea din DEI), câmpurile sistemului DEI MULTIMEDIA (în aplicația ACCESS), cu etichetele CUVÂNT, CUVÂNT-ACCENT și CUVÂNT-DEFINIȚIE și verificarea acestor câmpuri.

Grupul AUDIO: Nicoleta Todoroi, Silvia Donici, Ștefan Spătaru și un colectiv compus din 12 doctoranzi, masteranzi și studenți de la ASEM - a realizat **componenta AUDIO**

(înregistrări audio), câmpurile sistemului DEI MULTIMEDIA (în ACCESS), cu etichetele CUVÂNT-ACCENT - AUDIO și CUVÂNT-DEFINIȚIE - AUDIO și verificarea acestor câmpuri. S-au cercetat posibilitățile de inițiere, introducere, menținere, utilizare și distribuire a componentelor AUDIO în sistemul DEI MULTIMEDIA.

Grupul IMAGINI: Diana Micusa, Dumitru Micusa, Nicoleta Todoroi, Igor Coseru și un colectiv compus din 15 doctoranzi, masteranzi și studenți de la ASEM. Acest grup a realizat **componenta IMAGINI**, compusă din câmpul sistemului DEI MULTIMEDIA (în ACCESS) cu eticheta CUVÂNT-IMAGINE și verificarea acestui câmp. În prima etapă au fost examinate posibilitățile de stocare a imaginilor în ACCESS, optimizări necesare etc. În etapa următoare au fost introduse 2630 de imagini color din DEI și mai târziu au fost adăugate din diferite surse câteva sute de imagini suplimentare.

Grupul VIDEO: Victor Andronatiev, Zinaida Todoroi, Dumitru Micusa, Nicoleta Todoroi și un colectiv compus din 8 doctoranzi, masteranzi și studenți de la ASEM a cercetat posibilitatea de a compune câmpul VIDEO din sistemul DEI MULTIMEDIA cu suport ACCESS. În continuare s-au înregistrat și extras secvențe VIDEO de 10-30 secunde, care ar explica acele cuvintele din DEI, diverse surse video, inclusiv Internet. Astfel de circa 100 de secvențe s-au acumulat în DEI MULTIMEDIA până în prezent.

La ora actuală, **DEI Multimedia** este doar această bază de date, care, pentru un utilizator neavizat, este cam dificil de folosit. Tocmai această problemă se dorește a fi principalul subiect tratat de lucrarea de față, dl. Adrian Chiorescu realizând chiar o interfață web ce este descrisă pe larg în capitolul 4.

2. Generalități privind dicționarele informatizate

Dicționare se realizează la noi în țară de foarte multă vreme. Cele mai mari dicționare, cele lingvistice în special, au necesitat ani și chiar zeci de ani de muncă, unele dintre ele nefiind nici acum terminate. Alta ar fi fost situația poate dacă atunci ar fi existat calculatoarele și programele de azi. În ultimele decenii s-au făcut progrese enorme în domeniul procesării limbajului natural, lingvisticii computaționale și astfel s-au reluat și finalizat multe astfel de proiecte.

Crearea de baze de date lexicale (LDB) este o preocupare a mai multor țări balcanice. În acest scop a fost inițiat proiectul CONCEDE ce își propune crearea de metode universale de realizare a acestor baze de date. Proiectul se desfășoară pentru 6 limbi central europene: bulgară, cehă, estonă, ungară, română și slovenă.

În ultimii ani numeroase proiecte de informatizare a limbii române s-au desfășurat și în Republica Moldova, printre cele mai importante fiind cel de informatizare a Marelui Dicționar al Limbii Române (MDLR). Cel mai recent este cel al Academiei de Studii Economice din Moldova și anume realizarea unei *baze de date multimedia* pornind de la Dicționarul Explicativ Ilustrat. Astfel rezultatul a fost baza de date în MS Acces 2000 intitulată **DEI Multimedia** care face și obiectul de studiu al acestei lucrări în următoarele capitole.

La baza realizării **DEI Multimedia** au stat ideile și tehnologiile folosite pentru MDLR, adică reunirea a diferite sub-dicționare pe anumite direcții cum ar fi: TEXT, AUDIO, IMAGINI și VIDEO.

3. Scenarii de utilizare – posibile interfețe utilizator

Utilizarea de către publicul larg a bazei de date multimedia în SGBD MS Acces 2000 (**DEI Multimedia**) este la ora actuală destul de anevoioasă pentru că trebuie să se folosească facilitățile MS Acces-ului pentru navigare, căutare, interogare etc. Ori nu toată lumea cunoaște limbajul SQL de exemplu pentru o căutare cât mai precisă.

Din acest motiv ar trebui create aplicații care să interacționeze cu această bază de date și să dispună de interfețe cu utilizatorul cât mai prietenoase și ușor de folosit. Interogările cu baza de date trebuie să fie transparente pentru utilizator.

Este imposibil de realizat o „aplicație perfectă” care să placă și să folosească eficient absolut tuturor utilizatorilor și de aceea ar trebui create mai multe astfel de aplicații, fiecare având un „grup țintă” de utilizatori. Am identificat 5 astfel de grupuri țintă care să acopere cât mai mult din publicul larg vorbitor de limbă română și nu numai:

3.1. Publicul preșcolar și școlar de clase mici (I, II)

DEI Multimedia ar fi ideal în acest scop dacă peste el s-ar construi o aplicație cu o interfață cât mai atrăgătoare pentru copii, cât mai veselă, viu colorată, animată. Aici modul de prezentare este cel care contează foarte mult și interfața trebuie astfel organizată încât să acorde un spațiu mai larg afișării imaginilor și a clipurilor video.

3.2. Elevi de ciclu primar (III, IV) și ciclu gimnazial

O interfață pentru acest grup de utilizatori ar trebui să fie, la fel ca și în cazul anterior, veselă, viu colorată, atrăgătoare și sugestivă, ușor de folosit, pentru ca elevii să nu aibă nevoie de ajutorul profesorului și să o poată utiliza singuri. În acest caz, organizarea interfeței trebuie să ofere în egală măsură spațiu de afișare atât textelor cât și imaginilor. De asemenea nu trebuie să lipsească opțiunea de căutare.

3.3. Elevi de liceu, studenți, utilizatori cu studii de nivel mediu sau superior

La acest grup deja interfața ar trebui să fie „serioasă”. Aplicația nu mai trebuie să fie orientată *mod de prezentare* ci orientată *funcționalitate*. Bineînțeles că nu trebuie pierdut din vedere aspectul estetic, interfața trebuie să fie atrăgătoare, dar în același timp simplă și ergonomică, să asigure rapiditatea funcționării.

3.4. Cei ce învață limba română, turiști străini, studenți străini

Deși pare un grup restrâns, nu trebuie deloc pierdut din vedere, mai ales că în acest caz este necesară o aplicație cu totul specială. Pe lângă faptul că interfața ar trebui prezentată în mai multe limbi de circulație internațională, aceasta trebuie să scoată în evidență

foarte mult controalele audio pentru a se putea studia cu mare atenție pronunția cuvintelor.

3.5. Mediul academic lingvistic, studenți sau profesori de litere

Deși aceștia sunt „creatorii” dicționarelor, chiar ei au nevoie de multe ori de un dicționar cu un acces foarte rapid la informație atât pentru uz personal cât și pentru uz didactic.

O aplicație destinată acestui grup de utilizatori va trebui să ofere tot felul de posibilități de generare și listare de rapoarte, analize, statistici, de salvare, exportare în diverse formate a datelor pentru ca aceștia să le folosească ulterior mai ales în scop didactic.

4. Aplicația web: DEI Multimedia online

Studiind necesitățile fiecărui grup de utilizatori și făcând o „medie” a acestora am creat aplicația **DEI Multimedia online** ce se dorește a fi universală, adică să poată fi utilizată cu succes de utilizatori din toate grupurile.

Este o aplicație web (online) și tocmai de aceea prezintă unele avantaje suplimentare: nu trebuie actualizată de utilizatori (este actualizată permanent de echipa DEI Multimedia), asigură foarte ușor feedback-ul, este accesibilă oricând, de oriunde și de către oricine are o conexiune Internet (independentă de platformă), în schimb, fiind online, nu excelează cu modul grafic de prezentare fiind orientată mai mult spre funcționalitate. Este o aplicație cu un design extrem de simplu dar plăcut (Figura 1).

Am hotărât ca acest exemplu să fie un site web, pentru că nu există unul asemănător pentru limba română. Deși există numeroase dicționare online complete ale limbii române (www.dexonline.ro), ele prezintă numai definiții în mod text, nici unul nu este multimedia.

DEI Multimedia online este o aplicație web distribuită, construită pe arhitectura Client/Server în triplă legătură. Această structură este cea mai întâlnită la aplicațiile distribuite și se constituie din 3 părți: *interfața cu utilizatorul*, *funcționalitatea* și *baza de date*, fiecare dintre părți putând fi stocată pe calculatoare diferite (de aici distribuită). Astfel, utilizatorul, operează cu interfața prin intermediul browserului web de pe calculatorul propriu. Acesta emite cereri către *serverul web* pe care se găsește propriu-zis codul sursă (*funcționalitatea*) și primește informațiile dorite. Browserului web îi rămâne sarcina de a formata informațiile pentru afișare. Dacă sunt necesare informații din baza de date, atunci serverul web este cel care le cere de la *serverul de date* (interogări SQL) și le prelucrează, acest lucru fiind transparent pentru calculatorul client (utilizator).

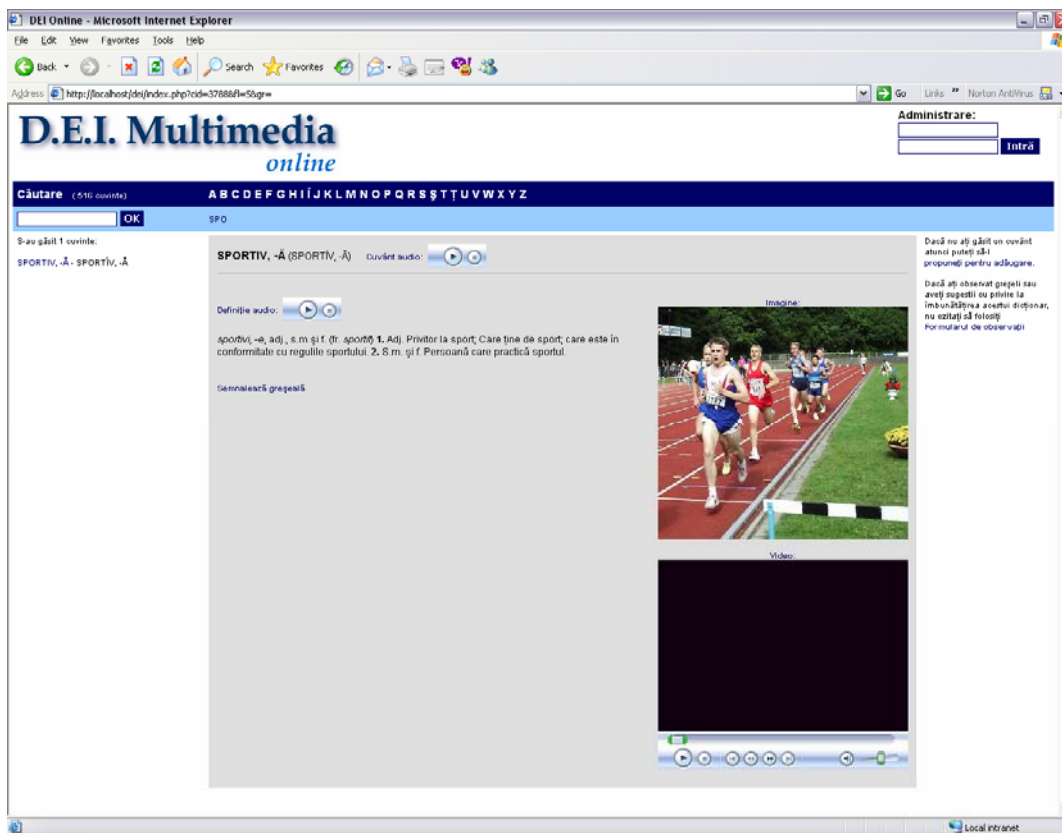


Figura 1 – Interfața cu utilizatorul

Aplicația are de fapt două interfețe: cea prezentată până acum, dedicată utilizatorilor obișnuiți ce caută cuvinte (modulul front-end) și o interfață dedicată administratorilor site-ului, unde aceștia, după ce se autentifică, pot opera modificări, adăugări, ștergeri etc. din baza de date. Acesta se numește *Modul de administrare* (modulul back-end).

Tot în acest modul, administratorii analizează propunerile de cuvinte făcute de utilizatori și le acceptă sau nu în baza de date, de asemenea, tot aici pot vedea și greșelile sau observațiile trimise de utilizatori.

5. Concluzii și perspective

Așa cum și-a propus, această lucrare tratează modalitățile prin care baza de date MS Acces 2000, **DEI Multimedia**, poate fi făcută accesibilă publicului larg vorbitor de limbă română și nu numai.

Acest lucru se poate realiza direct construind aplicații (în orice limbaj de programare vizual sub Windows) ce se conectează la baza de date și facilitează interacțiunea utilizatorului cu aceasta folosind interfețe utilizator cât mai accesibile.

DEI Multimedia acum este introdus și în Consorțiul pentru informatizarea limbii române.

Ca perspective de viitor, **DEI Multimedia** dorește să participe într-un proiect INTAS pentru 2007-2009 sub titulatura “*MULTIMEDIA Multilingual Dictionaries System for Republic of Moldova, Georgia, Armenia, and Azerbaijan in the Process of its Adhesion to the European Community Structures*”.

Referințe bibliografice

- Todoroi D., Todoroi, Z., Micusa, D. (2004). Procesarea limbajului natural în baza limbii computerizate române. // *România și Republica Moldova: Problemele competitivității economiilor naționale. Posibilități de valorificare pe piața internă, europeană și internațională*, București, INCE, 2004, p. 369-375.
- Todoroi, N., Todoroi, Z., Todoroi, D. (2004). Complexity Degrees of Illustrated Encyclopaedic Dictionary MULTIMEDIA. *Proceedings of the Int. Symp. “Inovative Applications of Information Technologies in Business and Management”*, October 22-23, 2004, Iași, Romania, pp. 23-27.
- Todoroi, D., Cristea, D., Tufiș, D., Todoroi, Z. (2003). Limba Română – Limba comunității europene. (LR – LCE – 2000). *Economica*, Nr. 1(44), p. 99-105.
- Todoroi D., Micusa, D., Todoroi, Z., I. Linga, I. Covalenco, N. Objelean, S.Spataru, S.Lungu, V. Turcanu, E. Cozlov, N. Ambrozii, V. Slobodeanu, I. Coseru, C. Suruceanu. (2002). Dictionarele multimedia ale limbii române. Secvențe de implementări și experimentări. *Limba Româna în Societatea Informațională – Societatea Cunoașterii*, Ed. Expert, Academia Română, București, p. 401-421.
- Todoroi, D., Todoroi, Z., Micusa, D. (2001). Romanian Computerized Language – One of the European Community Languages. *Proceedings of the 26th Annual Congress of the American Romanian Academy of Arts and Sciences (ARA)*, Montreal, Quebec, Canada, July 25-29, 2001, pp. 133-137. (Rom)
- Todoroi, D. (2001). The Computerized Romanian Natural Language Processing Development-Projects-Perspectives. // *INFORMATION SOCIETY. The Proceedings of the 5th International Symposium on Economic Informatics*, May 2001, Ed Economica, Bucharest 10-13 May 2001, pp. 927-935.
- Micusa, D., Jucan, T., Todoroi, D. (2002). The E-T-M Formalism for NLP Adaptable Processors’ Interactions. *Proc. of the Intern. Conf. “Globalisation And University’ Economic Education”*, Vol. II, October 24-27, 2002, Iasi, Romania, pp. 200-218.

MAPAREA CUVINTELOR DINTR-UN LEXICON PE ONTOLOGIE

NATALIA BURCIU, ANTONINA BÎRLĂDEANU

*Universitatea Tehnică a Moldovei, Facultatea Calculatoare Informatică și
Microelectronică, Chișinău*

natusicb@gmail.com, toni_birlad@yahoo.com

Rezumat

Mapare înseamnă corespondență, adică fiecărui element dintr-o mulțime îi corespunde un alt element din altă mulțime. În termenii acestui proiect – Maparea cuvintelor dintr-un lexicon pe ontologie – maparea este un proces complex care constă în crearea unei ontologii a termenilor juridici, crearea unui lexicon pe baza unui corpus de texte din domeniul juridic, crearea unui adnotator care adnotează fiecare termen juridic din text cu conceptul corespunzător ontologiei.

1. *Introducere*

Aplicațiile Software pentru Procesarea Limbajului Natural sunt în continuă dezvoltare, în special în domeniile Ontologiilor, Rezumării Informației, Extragerii Informației, Webului Semantic, Traducerii Automate etc. (van Harmelen, Fensel, 1999). Aplicațiile soft pentru Traduceri Automate, Extragerea Informației se dezvoltă, în special, în domeniul guvernamental și juridic pentru a obține traduceri (documente, texte) mai utile și mai perfecte decât la nivel de conversație sau texte mai puțin standardizate.

Sunt multe realizări în aceste domenii dar foarte puține pentru limba română. Proiectul „Maparea cuvintelor dintr-un lexicon pe ontologie” are ca obiectiv îmbogățirea realizărilor în domeniu și a resurselor sistemelor informaționale pentru limba română.

2. *Descrierea proiectului*

Maparea cuvintelor dintr-un lexicon pe ontologie constă în crearea unei ontologii a termenilor juridici, crearea unui lexicon, crearea unei aplicații – adnotator semiautomat – ce adnotează termenii juridici din texte cu conceptele corespunzătoare din ontologie. Semantici formale sunt de obicei încorporate în ontologii. O „adnotare semantică” prezintă o descriere mult mai precisă a cunoștințelor conținute în texte și a semanticii acestora în domeniul juridic. O adnotare semantică trebuie să fie bine definită, ușor de înțeles de către experții din domeniu și să nu fie ambiguă. Pentru a respecta aceste cerințe, o adnotare semantică trebuie să bazeze pe un model formal al domeniului, de exemplu ontologia. Formalizarea schemei de adnotare utilizând ierarhia ontologică permite adnotatorului semiautomat să aleagă nivelul corect al detaliului de adnotare, să diminueze ambiguitatea și să reducă erorile în timpul procesului de adnotare (van Harmelen, Fensel, 1999)

2.1. Ontologia termenilor juridici. Fragment

Ontologia este definită ca specificație a conceptualizării (Horrocks, 2000). Ontologia termenilor juridici constă din 44 clase amplasate ierarhic, cuprinde peste 140 sloturi și există în format RDFS. Un fragment din ontologie este în figura 1.

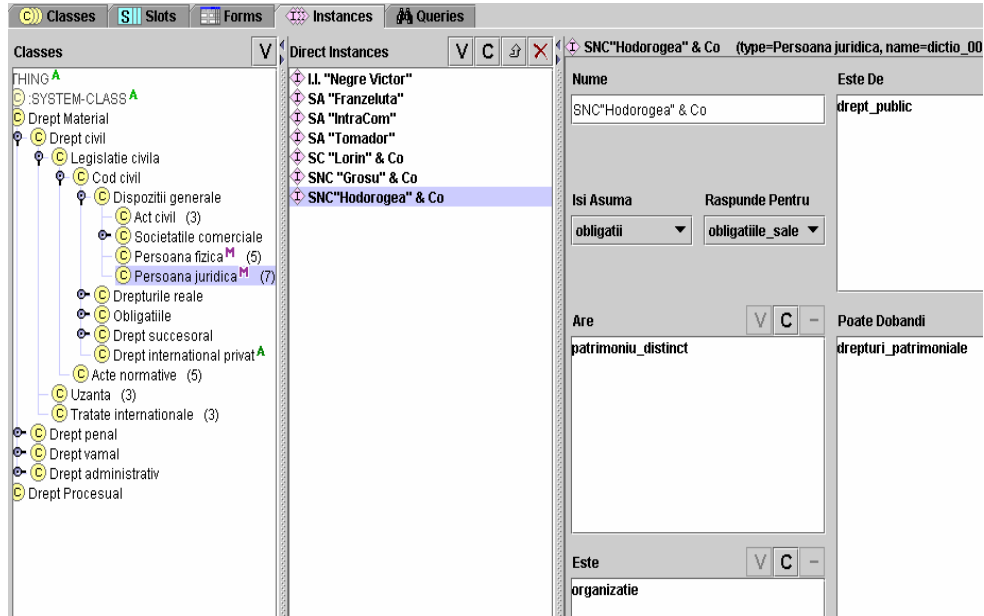


Figura 1: Fragment din ontologia termenilor juridici

2.2. Lexiconul

Lexiconul a fost creat manual dintr-un corpus de texte din domeniul juridic. Fiecare termen juridic este adnotat cu conceptul corespunzător din ontologie via taguri XML ca în tabelul de mai jos:

Decizia	<sem>de ciză </sem>
instanțe	<sem> instanța judecătorească </sem>
recunoașterea	<sem>recunoașterea dreptului</sem>
recunoașterea	<sem>recunoașterea hotărârii</sem>
dreptului	<sem> drepturile reale </sem>
spațiu	<sem>patrimoniului </sem>
locativ	
contractului	<sem> contract de locațiune </sem>

Figura 2: Fragment din lexicon

2.3 Adnotatorul semiautomat

Adnotarea Semantică este o tehnologie de bază pentru conținutul inteligent și este foarte utilă pentru o mulțime de aplicații inteligente orientate pe conținut (Vintar, 2003). Aplicația soft – adnotator semiautomat – a fost creat cu ajutorul limbajului de programare C++. El verifică fiecare cuvânt din text dacă este în lexicon și respectiv dacă este găsit adnotat copie tagul XML cu conceptul și îl alipește termenului dat în text. Deasemenea, dacă adnotatorul găsește în lexicon că termenul aparține mai multor concepte, el permite utilizatorului să aleagă conceptul corect, corespunzător contextului în care se află termenul juridic. Aceasta poate fi văzută în figura următoare:

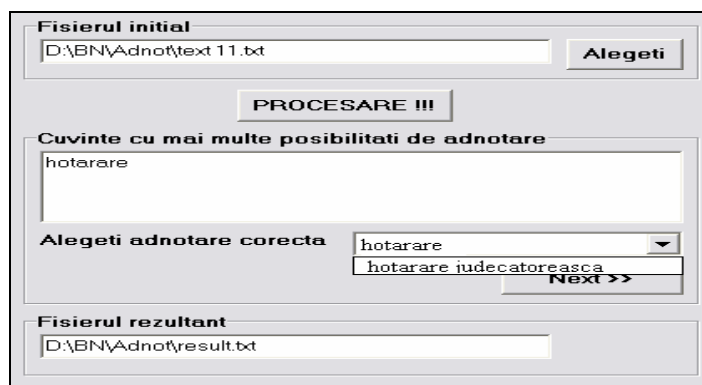


Figura 3: Interfața adnotatorului semiautomat

Fragment de text adnotat:

Contractul<sem> contract de vanzare-cumparare</sem> de vinzare-cumpararea autoturismului<sem> patrimoniu</sem> a fost recunoscut valabil din motivul ca alin. 2 art<sem> acte normative</sem>. 49 C.C. nu stabileste un termen concret in timpul caruia tranzactia executata trebuie sa fie intocmita. C.D.<sem> persoana fizica</sem> a indicat ca in baza tranzactiei orale din 12.08.1997 cu P.V.<sem> persoana fizica</sem> a procurat automobilul<sem> patrimoniu</sem> "Alfa-Romeo" 164 TD numarul de stat CDP-408 cu pretul de 9.400 lei, inasa nu a dovedit sa intocmeasca contractul<sem> contract</sem> la notariat<sem> notar</sem> deoarece masina<sem> patrimoniu</sem> a fost ridicata de la el<sem> persoana fizica</sem> de catre colaboratorii<sem> persoana fizica</sem> Considerind ca in aceasta tranzactie el<sem> persoana fizica</sem> este cumparator<sem> persoana fizica</sem> de buna credinta, reclamantul<sem> reclamant</sem> solicita sa fie recunoscuta tranzactia de vinzare-cumparare a automobilului<sem> patrimoniu</sem> "Alfa-Romeo" dintre el<sem> persoana fizica</sem> si P.V.<sem> persoana fizica</sem> valabila, iar Departamentul sa fie obligat sa-i<sem> persoana fizica</sem> intoarca automobilul<sem> patrimoniu</sem> care a fost ridicat ilegal. Prin hotarirea<sem> hotarare</sem> Judecatoriei<sem> instanta judecatoreasca</sem> sect.Buiucani mun.Chisinau din 07.09.1998 actiunea<sem> actiunea</sem> a fost admisa.

3. Concluzii

Textele cu termenii adnotați cu conceptele corespunzătoare din ontologia termenilor juridici, care descrie acest domeniu, va îmbunătăți procesul de extragere a informației din texte și documnte cu conținut juridic pentru limba română. Textele vor contribui deasemenea la obținerea unor rezultate mult mai calitative în Traduceri Automate prin dezambiguizarea termenilor juridici. Mai mult ca atât, datorită faptului că adnotarea se bazează pe ontologie ne face să utilizăm formalisme standardizate, așa ca RDF și OWL care permit reutilizarea acestor adnotări de către alte instrumente de adnotare sau instrumente de căutare.

Referințe bibliografice

- van Harmelen, F., Fensel, D. (1999). Practical Knowledge Representation for the Web. *In Proceedings of the IJCAI Workshop on Intelligent Information Integration.*
- Horrocks, I. (2000). The ontology interchange language oil: The grease between ontologies. *Technical report*, Dep. of Computer Science, Univ. of Manchester, UK/ Vrije Universiteit Amsterdam, NL/ Administrator, Nederland B.V./ AIFB, Univ. of Karlsruhe, DE.
- Vintar, S. (2003). Using parallel corpora for translation-oriented term extraction. Internet [<http://www2.arnes.si/~svinta/babel.rtf>].

CREAREA RESURSELOR LINGVISTICE CU AJUTORUL UNUI LIMBAJ SPECIALIZAT

ȘTEFAN DIACONESCU

SOFTWIN, București
sdiaconescu@softwin.ro

Rezumat

Lucrarea de față prezintă o metodă ce permite tratarea relativ unitară a mai multor capitole din lingvistică prin intermediul unui limbaj de reprezentare a cunoștințelor lingvistice numit GRAALAN (Grammar Abstract Language). Acest limbaj oferă unui lingvist posibilitatea descrierii eficiente a cunoștințelor lingvistice privind o limbă naturală precum și corespondența între două limbi naturale.

1. Introducere

Există numeroase și fructuoase încercări de uniformizare a reprezentării cunoștințelor lingvistice. O asemenea uniformizare ar oferi un avantaj foarte mare în dezvoltarea unor studii, statistici și, în cele din urmă, aplicații lingvistice care să poată trata într-un mod asemănător diverse limbi naturale sau să poată compara (stabili corespondențe) într-un mod unitar între diverse limbi naturale. Din păcate diversele capitole lingvistice au suferit abordări întrucâtva independente, cum ar fi subcategorizarea (EAGLES, 1996b) adnotarea (EAGLES, 1996b), lexiconul (EAGLES, 1993), etc. astfel încât este uneori dificil de aplicat tratamente unitare.

Comunicarea de față prezintă un limbaj de reprezentare a cunoștințelor lingvistice numit GRAALAN (Grammar Abstract Language). Acest limbaj permite unui lingvist descrierea eficientă a cunoștințelor lingvistice privind o limbă naturală precum și corespondențele între două limbi naturale.

2. Caracteristicile generale ale GRAALAN

Din punct de vedere teoretic, GRAALAN se bazează în special pe următoarele noțiuni: gramatici generative de dependențe (GDG - Generative Dependency Grammar) (Diaconescu, 2002), arbori de dependențe (DT - Dependency Tree) (Diaconescu, 2002) și arbori atribut - valoare (AVT - Attribute Value Tree) (Diaconescu, 2005).

Pornind de la aceste noțiuni, GRAALAN poate descrie diverse capitole lingvistice conforme cu gramaticile convenționale ale limbilor naturale: alfabetul, despărțirea în silabe, morfologia, sintaxa, regulile de flexiune, formele de flexiune, lexiconul, corespondențe lexicale între două limbi (inclusiv între expresii multicuvânt MWE - Multiword Expression), corespondențe morfologice, corespondențe sintactice.

GRAALAN este în esență un limbaj descriptiv care permite însă eventual și legătura cu anumite subprograme de tip procedural scrise în alte limbaje de programare.

În principiu, descrierile GRAALAN vor putea fi convertite printr-un compilator adecvat în formatul XML care este mai adecvat exploataării ulterioare prin diverse programe.

3. Descrierea alfabetului

În GRAALAN se pot preciza pentru o anumită limbă: alfabetul fonetic utilizat în descrierea limbii (care poate fi un subset al IPA (International Phonetic Alphabet) (IPA, 2005), alfabetul normal și caracterele speciale.

În afară de acestea se mai pot defini: i) grupe de caractere (diftongi, triftongi, etc.), transcrise cu caractere normale (eventual speciale) dar și fonetice; ii) clase alfabetice (de exemplu clasa vocalelor, clasa consoanelor, etc.)

Caracterele folosite în GRAALAN se consideră codificate în UNICODE (ISO, 1992).

4. Descrierea despărțirii în silabe

În GRAALAN sunt considerate trei tipuri de despărțire în silabe: i) Despărțirea eufonică a cuvintelor scrise cu alfabetul normal și respectând modul de pronunție; ii) Despărțirea fonetică a cuvintelor scrise cu alfabetul fonetic și respectând de asemenea modul de pronunție; iii) Despărțire morfologică - analogă cu despărțirea eufonică însă respectând și restricții ce țin cont de structura morfematică a cuvântului.

Primele două tipuri au reguli specifice. Ultimele tipuri nu are reguli speciale deoarece ea acționează ca o despărțire eufonică cu restricțiile suplimentare privind morfemele obținute din consultarea lexiconului.

5. Descrierea morfologiei

În GRAALAN, morfologia unei limbi (mai exact ansamblul categoriilor lexicale și al valorilor lor), se reprezintă sub forma unui arbore atribut valoare (AVT) (Diaconescu, 2005) în care nodurile de tip atribut corespund categoriilor lexicale iar nodurile de tip valoare corespund valorilor categoriilor lexicale. În plus, cele două tipuri de noduri mai au atașate diverse alte tipuri de informații: numele, abrevieri, (eventual) atașamente procedurale, etc.

În secțiunea corespunzătoare morfologiei se poate indica de asemenea dacă anumitor situații de flexiune distincte le corespund forme flexionate identice.

6. Descrierea lexiconului

Lexiconul GRAALAN este un ansamblu de intrări de diverse tipuri: i) Morfeme (rădăcini, prefixe, sufixe, prefixoide, sufixoide etc.); ii) Cuvinte care la rândul lor pot fi: intrări principale de tip lemă, intrări suplimentare (care însoțesc o lemă), intrări principale care nu sunt însă leme; iii) MWE-urile cărora li se indică și structura sub forma unui arbore de dependențe (DT); iv) Structuri morfologice analitice sau analitico-sintetice (forme flexionate formate din mai multe cuvinte) analoge MWE-urilor; v) Structuri sintactice de asemenea analoge MWE-urilor.

În funcție de tipul lor, intrările în lexicon mai pot avea asociate și alte tipuri de informații: semantice, etimologice, morfologice, etc.

Lexiconul în general nu este scris direct în GRAALAN ci se creează cu ajutorul unui instrument specializat.

7. Descrierea regulilor de flexiune

Intrarea din lexicon care se poate flexiona (lema de exemplu) identifică o regulă compusă de flexiune aflată în secțiunea GRAALAN a regulilor de flexiune. Regula compusă este o listă de reguli de bază. O regulă de bază este de fapt un arbore atribut valoare care indică mai multe situații de flexiune, câte una pentru fiecare frunză a sa. Fiecare situație de flexiune (deci frunză) are asociată una sau mai multe reguli de flexiune elementare. O regulă de flexiune elementară conține: i) O condiție de aplicare a regulii; ii) O secvență de transformări care trebuie făcute asupra lemei (sau asupra altei forme de flexiune) pentru a obține forma de flexiune curentă exprimată în alfabetul normal; iii) Analog cu (ii) pentru alfabetul fonetic; iv) În cazul formelor analitico-sintetice - o caracterizare sub forma unui AVT pentru fiecare cuvânt component și relațiile care se află între diversele cuvinte componente.

Pe baza regulilor de flexiune aplicate intrărilor din lexicon se pot obține formele din secțiunea GRAALAN a formelor de flexiune.

8. Descrierea formelor de flexiune

Secțiunea GRAALAN corespunzătoare formelor de flexiune conține câte o intrare pentru fiecare formă de flexiune. O intrare conține: i) Forma de flexiune în alfabet normal și fonetic; ii) Identificarea în lexicon a intrării căreia îi corespunde forma respectivă de flexiune; iii) Caracterizarea formei de flexiune sub forma unui ansamblu de categorii lexicale cu valorile lor (AVT); iv) Despărțirea în silabe.

Formele de flexiune nu sunt scrise în general direct în GRAALAN ci se creează cu ajutorul unui instrument specializat.

9. Descrierea sintaxei

Sintaxa se descrie în GRAALAN sub forma unei liste de reguli sintactice etichetate (care respectă principiile gramaticilor de dependențe generative (Diaconescu, 2002)).

O regulă are un membru stâng care conține un neterminal însoțit de un AVT format din categorii lexicale și/sau sintactice) și un membrul drept care conține unul sau mai mulți alternanți. Un alternant este format din trei subsecțiuni:

a) Subsecțiunea sintactică care conține o secvență de NTPA: Neteminali, Terminali, Pseudo terminali, Acțiuni (subprograme procedurale). Neterminalii și terminalii au accepțiunea obișnuită. Pseudoterminalii sunt neterminali care, dacă ar avea reguli care să îi descrie, acestea ar conține direct terminali din lexicon. Acțiunile sunt subprograme procedurale care ar putea fi utilizate în anumite tratamente specifice dacă este cazul.

Fiecare NTPA conține un nume, un AVT format din categorii lexicale și/sau sintactice, modul de legare (relaționare) cu alți NTPA.

b) Subsecțiunea de dependențe unde se descriu relațiile de dependență între NTPA-uri ale alternantului. Relațiile de dependență pot fi de tip de regență / subordonare sau de tip coordonare.

c) Subsecțiunea de acord care descrie acordul între NTPA-urile alternantului sub forma unor condiții complexe.

Descrierea sintaxei în GRAALAN este reversibilă adică poate fi folosită și în procesul de analiză sintactică prin care se generează din textul de suprafață un arbore de dependențe ca formă de adâncime, și în procesul de generare din arborele de dependențe a textului de suprafață.

10. Descrierea corespondențelor bilingve

Secțiunea GRAALAN privitoare la corespondențele bilingve descrie corespondențe între următoarele tipuri de elemente aparținând la două limbi diferite:

a) Corespondențe între MWE-uri care sunt reprezentate în lexicon sub forma unor arbori de dependențe se exprimă prin echivalarea între expresia sursă și expresia țintă corespunzătoare dar și prin regulile de transformare care indică modul în care extensiile expresiei sursă din instanțe reale sunt preluate de expresia țintă.

b) Corespondențe între cuvinte. Este un caz particular al corespondenței între MWE-uri în care expresiile echivalate au câte un singur cuvânt.

c) Corespondențe între structuri sintactice. Este un caz particular al corespondenței între MWE-uri în care cele două expresii pot avea drept caracterizări de noduri nu numai categorii lexicale (cu valorile lor) ci și categorii sintactice (cu valorile lor).

d) Corespondențe între structuri morfologice. Este un caz particular al corespondenței între MWE-uri în care cel puțin expresia sursă corespunde unei forme flexionate analitico-sintetice.

e) Corespondențe între subarbori morfologici. Este o corespondență între diverse seturi de categorii lexicale (cu valorile lor) organizate sub forma unor AVT-uri.

Informațiile din secțiunea de corespondențe bilingve GRAALAN se pot folosi în aplicații de generare a unor dicționare sau în aplicații de traducere automată.

11. Concluzii

Descrierile de cunoștințe lingvistice pot fi formulate direct în GRAALAN sau, în anumite cazuri (cum ar fi de exemplu pentru formele flexionate sau pentru lexicon) pot fi create cu ajutorul unor instrumente (programe) speciale care generează text GRAALAN. Textul GRAALAN obținut pe o cale sau pe alta se compilează cu un compilator adecvat care traduce textul GRAALAN în XML, creindu-se astfel o Baza de cunoștințe lingvistice XML creată prin intermediul GRAALAN va putea fi exploatată într-un mod unitar pentru diverse studii sau pentru elaborarea de aplicații informatice.

Deoarece textul GRAALAN se realizează în mai multe tranșe, o componentă specială GRAALAN Link va determina legăturile între aceste tranșe și compatibilitatea lor.

Un compilator GRAALAN este în curs de implementare și unele cunoștințe lingvistice privind limba română au fost deja scrise în GRAALAN.

Referințe bibliografice

- Diaconescu, S. (2002). Natural Language Understanding Using Generative Dependency Grammar, în Max Bramer, Alun Preece and Frans Coenen (Eds), *Proceedings of ES2002*, Cambridge UK, Springer, pp.439-452.
- Diaconescu, S. (2003). Natural Language Agreement Description for Reversible Grammars, în Tamás D. Gedeon, Lance Chun Che Fung (Eds.), *Proceedings of AI 2003*, Perth, Australia, pp. 161-172.
- Diaconescu, S. (2004) Multiword Expression Translation Using Generative Dependency Grammar, în *Proceedings of ESTAL 2004 - ESPAÑA for NATURAL LANGUAGE PROCESSING*, Alicante, Spain.
- Diaconescu, S. (2005). Some Properties of the Attribute Value Trees Used for Linguistic Knowledge Representation, *Proceedings of IICAI-05, INDIA*.
- EAGLES (1996a). Recommendations for the Morphosyntactic Annotation of Corpora.
- EAGLES (1996b). Preliminary Recommendations on Subcategorisation.
- EAGLES (1993). Lexicon architecture Draft Report, EAG--LSG/IR--T1.1.
- IPA (2005) International Phonetic Association (2005): Handbook of IPA.
- ISO/IEC 10646 (1992). Information technology -- Universal Multiple-Octet Coded Character Set (UCS).

RESURSE LINGVISTICE ROMÂNEȘTI ÎN FORMAT ELECTRONIC. *BIBLIA 1688*

BOGDAN-MIHAI ALDEA¹, GABRIELA HAJA²

¹*Facultatea de Informatică, Universitatea "Al.I.Cuza", Iași*

²*Institutul de Filologie Română „A. Philippide”, Academia Română, Iași*

{bogdan.aldea, gabihaja}@gmail.com

Rezumat

Crearea resurselor textuale în format electronic prelucrat reprezintă o prioritate pentru procesul de informatizare a cercetării lingvistice românești. Un domeniu puțin cercetat la noi este cel al achiziționării în format electronic și al prelucrării textelor românești vechi. Lucrarea de față descrie rezultatele la care s-a ajuns în prelucrarea a două cărți din *Biblia de la 1688*, ms. 45 și ms. 4389 (sec. XVII) în vederea generării automate a indicelui de cuvinte.

1. Introducere

Grație colaborării științifice dintre cercetători ai Institutului de Filologie Română „A. Philippide”, cercetători ai Institutului de Informatică Teoretică – Academia Română, Filiala Iași, cercetători de la Facultatea de Informatică și de la Facultatea de Litere a Universității „Alexandru Ioan Cuza” din Iași, la inițiativa specialiștilor informaticieni, s-a demarat la Institutul „Philippide” din Iași un amplu proces de informatizare a cercetării filologice românești. Meritul specialiștilor de la acest institut este acela că toate eforturile lor în această direcție s-au concentrat asupra creării unor instrumente și resurse specifice proiectelor Academiei Române, dintre care le amintim pe cele prioritare: *Dicționarul limbii române (DLR)*, *Dicționarul general al literaturii române (DGLR)*, *Noul Atlas lingvistic român, pe regiuni. Moldova și Bucovina (NALR – MB)*.

Cu referire la cercetările lexicografice din domeniul limbii române, a fost finalizat, în 2005, grantul *Dicționarul limbii române (DLR) în format electronic. Studii privind achiziționarea*, rod al colaborării între lingviștii lexicografi de la Iași și cercetători de la Facultatea de Informatică din Iași. Rezultatele acestui proiect au fost publicate (Haja et al., 2005), și trebuie subliniat că necesitatea definitivării unei forme electronice **integrale** a *Dicționarului limbii române* se impune cu o tot mai mare stringență.

Finalitatea cercetărilor din cadrul proiectului încheiat în 2007 ca și ale celui în desfășurare la Institutul „Philippide” este realizarea unui *Dicționar al limbii române informatizat (DLRI)*, creat ca instrument și resursă lexicografică, punct de plecare în constituirea unui dicționar al limbii române care să poată fi permanent actualizat și îmbogățit prin reeditări consecutive, comparabil cu lucrări fundamentale pentru culturile lumii – tezaure ori dicționare ale limbii electronice – precum cele realizate de lexicografia franceză, italiană, spaniolă, ca să amintim doar trei dintre limbile romanice europene, ori de lexicografia anglosaxonă europeană și americană.

După ce a fost stabilită soluția achiziționării în format electronic a formei tipărite a *DLR* și s-a creat instrumentul de achiziție, prelucrare și consultare a acestuia, DLReX, un alt element care este indispensabil realizării *DLRI*, anume **informatizarea colecției de texte** din care sunt excerptate atestările după care se redactează *DLR*, a intrat în atenția noastră. Pentru început, am optat pentru crearea formatului electronic al ediției unei lucrări care face parte din *Bibliografia DLR*: prima traducere integrală în română a *Vechiului și Noului Testament*, citată de literatura de specialitate sub numele de *Biblia de la București (BB)* sau *Biblia de la 1688*, un monument al limbii și culturii românești (Andriescu, 1997; Miron, 1988, 2004).

2. Tradiție și actualitate

2.1. Monumenta linguae dacoromanorum. Biblia de la 1688

Reeditarea critică a *BB*, într-o ediție cu format enciclopedic, în care sunt cuprinse și două variante de traducere realizate în același secol XVII, dar rămase în formă manuscrisă, *Manuscrisul 45* (ms. 45) și *Manuscrisul 4389* (ms. 4389), a fost inițiată de Paul Miron, profesor la Albert-Ludwigs-Universität din Freiburg, Germania, și realizată la Iași, prin concursul specialiștilor români și germani, lingviști, istorici literari, informaticieni și istorici, oameni de cultură implicați într-un proiect amplu de recuperare și valorificare a textelor fundamentale ale culturii române.

Din această ediție, proiectată în 20 de părți, au fost tipărite, în seria *Monumenta linguae dacoromanorum*, la Editura Universității „Alexandru Ioan Cuza” din Iași, șapte volume (în ordinea cărților din biblice: *Pentateuhul*, 1988–1997, *Iosue*, *Judicum*, *Ruth* 2005, *Liber Psalmorum* 2003). În prezent se lucrează, în cadrul grantului *Resurse lingvistice în format electronic. Monumenta linguae dacoromanorum. Biblia 1688. Pars VII. Regum I, Regum II – ediție critică și corpus adnotat*, la cel de-al optulea volum.

Obiectivele acestui proiect sunt, pe lângă continuarea monumentalei ediții în forma sa tipărită, definirea unui format electronic al acestui volum, adnotarea semiautomată, la nivel de cuvânt, a textului românesc vechi, crearea unui program de indexare a textului. La finele proiectului se vor fi creat premisele constituirii unui nucleu de corpus de limbă română veche, necesar cercetărilor lingvistice în genere, a celor lexicografice în special, și va fi definitivat instrumentul de realizare a formatului electronic al întregii ediții.

În proiect sunt implicați cercetători din domeniul lingvisticii și al informaticii care participă la activități diferite (cercetarea filologică a textului și crearea instrumentelor informatice de prelucrare a textului) și de comune (prelucrarea textului), ca în Figura 1.

<p>Lingviști:</p> <ul style="list-style-type: none"> - transcrierea interpretativă a textelor; - stabilirea textului, colaționarea, corectura, revizia; - pregătirea volumului pentru tipar. 	<p>Informaticieni:</p> <ul style="list-style-type: none"> - realizarea unui instrument de achiziționare în format electronic a textelor din sec. XVII; - achiziționarea unor eșantioane de text din <i>Monumenta linguae dacoromanorum. Biblia 1688. Pars VI și Pars VII</i>; - generarea indicelui de cuvinte.
<p>Lingviști + informaticieni:</p> <ul style="list-style-type: none"> - adnotarea la nivel de cuvânt a textelor; - definitivarea formatului electronic al volumului. 	

Figura 1: Distribuția activităților.

Potrivit tradiției create prin editarea volumelor de până acum, pe lângă reproducerea textului tipărit la 1688 însoțit de transcrierea interpretativă a acestuia, alături de cele două manuscrise, se propune o variantă apropiată de limba română contemporană a textului vechi, stabilită de partenerii de proiect din Freiburg. Aparatul critic al ediției constă în realizarea notelor de transcriere, a comentariilor filologice și, acolo unde este cazul, istorice, a indicelui de cuvinte pentru textul tipărit. O inovație pe care și-o propune proiectul în realizarea formei electronice a volumului în curs de finalizare o constituie lărgirea indexării asupra manuscriselor și asupra variantei actualizate a textului. Dintre activitățile descrise mai sus, au fost finalizate următoarele: transcrierea interpretativă a textelor (*BB*, ms. 45, ms. 4389); colaționarea și corectura textelor; achiziționarea unor eșantioane de text din *Pars VI*, achiziționarea unor eșantioane din *Pars VII*; realizarea unui instrument de achiziționare, prelucrare, indexare a textelor numit convențional AdBB (Ad < adnotare; BB < *BB*).

2.2. AdBB și funcționalitățile sale

AdBB reprezintă un instrument de achiziționare, prelucrare și consultare a *BB*, a ms. 45 și a ms. 4389 în format electronic. Principalele funcționalități ale acestei aplicații create, sunt: permite trecerea textului *Bibliei* (sec. XVII) din format .rtf (Word) în format XML; permite vizualizarea și corectarea fișierelor XML; generează grupuri de ocurențe ale unei lemme; permite editarea grupurilor; generează și vizualizează indicele de cuvinte.

2.2.1. Trecerea din format .rtf în XML

Se alege din meniu opțiunea RTF → XML pentru conversia unui text din format Word în format XML:

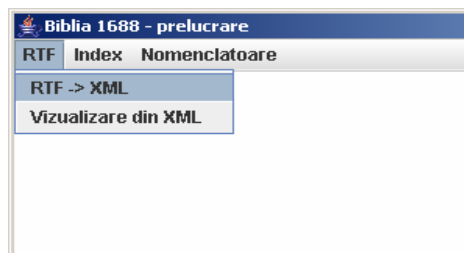


Figura 2: AdBB – meniu. Captură de ecran.

În urma parsării fișierului .rtf, rezultă un fișier XML cu următorul format:

```
<?xml version="1.0" encoding="UTF-8"?>
<biblia>
  <carte name="NUME CARTE">
    <capitol id="NR CAPITOL">
      <verset id="NR_VERSET"> Conținutul primului verset </verset>
      <verset id="NR_VERSET"> Conținutul versetului al doilea </verset>
      <verset id="NR_VERSET"> Conținutul versetului al treilea </verset>
      .....
    </capitol>
    <capitol id="NR_CAPITOL">
      </capitol>
    </carte>
    <carte name="NUME CARTE">
      <capitol id="NR_CAPITOL">
        .....
      </capitol>
    </carte>
  </biblia>
```

2.2.2. Vizualizarea fișierelor XML

Pentru vizualizarea fișierului XML se poate alege cea de-a doua opțiune din meniu, în acest fel putându-se revizui corectitudinea parsării.

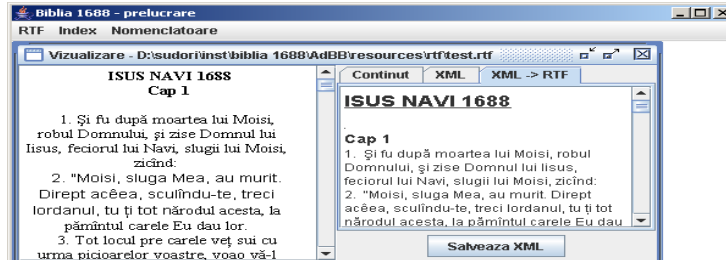


Figura 3: AdBB. Vizualizarea textului .rtf și XML, înainte de corectură. Captură de ecran.

2.2.3. Generarea grupurilor de ocurențe

Pentru a genera grupurile de forme flexionare se alege din meniul „Index” opțiunea „Generare Grup”:

- 1) Pentru gruparea cuvintelor dintr-un fișier, se generează vocabularul fișierului selectat și se elimină o serie de cuvinte cum ar fi: pronumele, articolele ș.a.m.d (vezi nomenclatoare)¹, rezultând un fișier XML, cu grupuri.

2.2.4. Editarea grupurilor forme flexionare

Pentru editarea grupurilor de forme flexionare, se alege cea de-a doua opțiune din meniul „Index”: „Editează Gr. Cuvinte”, și se alege pentru deschidere fișierul XML generat la pasul anterior.

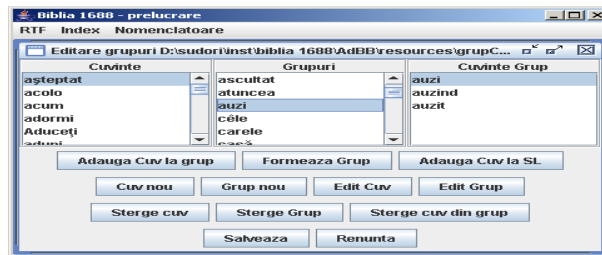


Figura 4: AdBB. Editarea grupurilor de ocurențe generate. Captură de ecran.

Astfel rezultă 3 coloane: Cuvinte, Grupuri și Cuvinte Grup. Prima coloană cuprinde cuvintele cu o singură apariție în text sau care nu prezintă forme flexionare, a doua cuprinde lista tuturor grupurilor generate automat, iar ultima cuprinde formele paradigmatică, existente în text, ale fiecărei intrări din coloana a doua, cu posibilitatea modificării acestora de către specialiștii lingviști.

2.2.5. Generarea și vizualizarea indexului

¹ Această serie de cuvinte este conținută într-un fișier intern ce poate fi editat și anume „romanian.stoplist”. StopList este o lista de cuvinte pentru care nu se dorește generarea indexului și care poate fi modificată selectând din meniul „Nomenclatoare” opțiunea „Editeaza StopList”

Pentru orice fișier XML rezultat în urma parsării și pentru orice fișier de grupuri de ocurențe, se poate genera un index. Se selectează a treia opțiune din meniul „Index”, „Generează index”, rezultând o interfață din care se selectează fișierul XML dorit, precum și fișierul de grupuri aferent acestuia și apoi se generează indicele de cuvinte. În urma generării, rezultă un fișier XML cu următorul format:

```
<?xml version="1.0" encoding="UTF-8"?>
<indecsi>
<data>
<grupFilePath>FISIERUL_XML_CU_GRUPURI</grupFilePath>
<bibleFilePath>FISIERUL_XML_CU_BIBLIA</bibleFilePath>
</data>
<cuvinte>
  <grup name="NUME GRUP #1">
    <cuvant value="CUVANT #1">
      <index>
        <carte>NUME_CARTE</carte>
        <capitol>NR_CAPITOL</capitol>
        <verset>NR_VERSET</verset>
      </index>
    </cuvant>
  </grup>
</cuvinte>
</indecsi>
```

Pentru vizualizarea indicelui, se alege ultima opțiune din meniul „Index” și se selectează fișierul XML generat la pasul anterior. Fiecărui cuvânt fiindu-i precizat locul în text, în formatul următor: [NUME_CARTE ; NR_CAPITOL ; NR_VERSET]

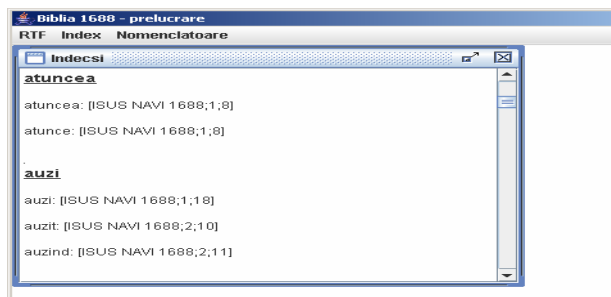


Figura 4: AdBB. Indexul, cu trimiteri, editat. Captură de ecran.

2.3. Necesitatea rafinării AdBB

Pentru o mai bună utilitate a AdBB, este necesară rezolvarea următoarelor chestiuni puse, deocamdată, de specialiștii lingviști consultați cu privire la rezultatele aplicației: **1)** posibilitatea extragerii de contexte (de diferite dimensiuni) în care apar ocurențele. Ideal al fi ca, din fișierul „Index”, să se poată accesa, prin link-uri, textul în format XML sau .rtf, din care să se selecteze citatele, cu posibilitatea editării acestora în DLRI. Ar fi un câștig, în acest sens, corelarea AdBB cu DLReX, astfel încât, pornind de la o formă paradigmatică ce apare în DLRI, aceasta să poată fi căutată automat în corpusul *BB*; **2)** cuvintele cu o singură apariție în text trebuie, de asemenea, indexate; există și în prezent această posibilitate, dar trebuie simplificați pașii necesari; **3)** cuvintelor cu frecvență foarte ridicată (cuvintele de relație, pronumele etc., inventariate în „StopList”) ar trebui să li se precizeze măcar frecvența; până în prezent, au fost prelucrate, prin AdBB, eșantioane mici de text; **4)** randamentul aplicației va putea fi verificat după ce vor fi supuse prelucrării fișiere .rtf care să cuprindă peste 100000 de ocurențe.

3. *În loc de concluzii*

Acest tip de cercetare interdisciplinară nu poate fi decât constructiv. Ceea ce s-a realizat, într-o perioadă relativ scurtă, deschide o cale nouă cercetării filologice românești, în care instrumentele electronice sunt adaptate diverselor tipuri de scriitură, demonstrându-se faptul că vechimea textului nu pune probleme majore mijloacelor actuale de prelucrare a limbii în forma sa scrisă.

Odată create, aceste instrumente facilitează cercetarea filologică, al cărei specific este minuțiozitatea și acribia, studiul comparativ și contrastiv al textelor, analiza grafiei și „arheologia” lingvistică, reducând spectaculos perioada de documentare și aceea de fișare a textelor, oferindu-i specialistului posibilitatea concentrării asupra fazei de analiză și interpretare, permițându-i concluzii mai ferme, pe baza unei evidențe cvsie exhaustive a faptelor de limbă. Acesta este doar un exemplu de ameliorare a cercetării, determinată de actualizarea instrumentelor și de informatizarea resurselor, într-un domeniu atât de vast precum acela descris.

Mulțumiri. Cercetarea descrisă aici s-a desfășurat în cadrul grantului *Resurse lingvistice în format electronic. Monumenta linguae dacoromanorum. Biblia 1688. Pars VII. Regum I, Regum II – ediție critică și corpus adnotat* (2006–2007). Autorii mulțumesc Ministerului Educației și Cercetării, CNCSIS, pentru susținerea financiară a proiectului.

Referințe bibliografice

- Andriescu, Al. (1997). *Locul Bibliei de la București în istoria culturii, literaturii și limbii române literare în Studii de filologie și istorie literară*, Iași, Editura Universității „Alexandru Ioan Cuza”, 90-208.
- Haja, G., Forăscu, C., Dănilă, E., Aldea, B. M. (2005). *Dicționarul limbii române (DLR) în format electronic. Studii privind achiziționarea*, Iași, Editura Universității „Alexandru Ioan Cuza”.
- Miron, Paul (1988). *O nouă ediție a Bibliei lui Șerban în Monumenta linguae dacoromanorum. Biblia 1688. Pars I. Genesis*, Iași, 3-6.
- Miron, Paul (2004). *Prefață la ediția Freiburg und München a cărții Ruth, în Monumenta linguae dacoromanorum. Biblia 1688. Pars VI. Iosue, Iudicum, Ruth*, Iași, Editura Universității „Alexandru Ioan Cuza”, 5-6.

RESURSE ROMÂNEȘTI ÎN CADRUL PROIECTULUI LT4EL

DIANA TRANDABĂȚ^{1,2}, ADRIAN IFTENE¹, IONUȚ PISTOL¹,
CORINA FORĂSCU^{1,3}, DAN CRISTEA^{1,2}

¹*Facultatea de Informatică, Universitatea “Al. I. Cuza”, Iași*

²*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

³*Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București*

{dtrandabat, adiftene, ipistol, corinfor, dcristea}@info.uaic.ro

Rezumat

Proiectul LT4eL are ca scop realizarea unei tehnologii multilingve care să ajute la realizarea și exploatarea obiectelor de învățare utilizate în cadrul unui sistem de eLearning. La baza acestora stă un corpus semnificativ de documente, colectat inițial, apoi adnotat manual și automat pe diferite niveluri. Tehnologiile dezvoltate în proiect vor facilita operațiile de creare a obiectelor de învățare de către profesori cât și de regăsire a lor de către studenți, inclusiv prin criterii de natură semantică. Lucrarea prezintă etapele de colectare și prelucrare a acestui corpus pentru limba română.

1. Introducere

Proiectul LT4eL¹ (Tehnologii Lingvistice pentru eLearning) își propune utilizarea unor tehnologii multilingve, unelte lingvistice și tehnologii ale web-ului semantic pentru a perfecționa regăsirea și accesul la materiale de învățare în sistemele de management al învățării, prin generarea semi-automată a unor metadate descriptive. Astfel, va fi dezvoltat un extractor de cuvinte cheie și un detector de definiții și termeni definiți, adaptate tuturor limbilor implicate în proiect (bulgară, cehă, engleză, germană, malteză, olandeză, poloneză, portugheză și română).

Tehnologia ce va fi dezvoltată în cadrul proiectului va facilita accesul personalizat la cunoștințele din sistemele de management al învățării și va favoriza descentralizarea și cooperarea în managementul conținutului didactic (Monachesi et al., 2006).

După o trecere în revistă în secțiunea 2 a cerințelor specifice proiectului, cu accent pe domeniile din care s-au extras resursele, în secțiunea 3 vom prezenta succint etapele de prelucrare a resurselor, de la forma inițială în care au fost colectate din diferite surse, până la forma în care vor fi folosite drept corpus de antrenare/test în proiect.

2. Colecția de obiecte de învățare

Pentru a îmbunătăți managementul, distribuția și regăsirea materialului de învățare prin atașarea semi-automată de metadate este necesară, într-o primă etapă, observarea modului în care aceste metadate sunt marcate manual. Astfel, prima cerință a proiectului

¹ <http://www.lt4el.eu>

a fost colectarea și normalizarea unor obiecte de învățare, obiectiv realizat prin intermediul unui portal special dezvoltat² de partenerii din Universitatea „Al.I.Cuza” Iași (Pistol et al, 2006).

Cuvântul de ordine al proiectului LT4eL este multilingvismul. Cu nouă limbi implicate, resursele care vor constitui baza de plecare a extractoarelor automate de cuvinte cheie și definiții trebuie să fie comparabile pentru toate limbile (din aceleași domenii). Necesitatea ca domeniile reprezentate în proiect să fie relativ uniform acoperite în toate cele nouă limbi, să aibă o mare deschidere spre sisteme de eLearning, iar documentele să îndeplinească simultan și criteriul de accesibilitate cu restricții minime în privința drepturilor de autor au dus la alegerea domeniului informatic, cu precădere a celui dedicat predării de noțiuni informatice către începători, și a domeniului eLearning. Aceste două domenii mari au fost rafinate în mai multe subdomenii, printre care: *Writing a diploma paper, Making an interview, Using MS Word/Excel/PowerPoint/Latex/XML, Creating Web pages, Accessing the Internet, eLearning, eMarketing, Impact of use of computers in society, Impact of eLearning on education, Calimera Documents* etc.

Corpusul colectat pentru limba română conține 56 de documente din aproape toate domeniile și subdomeniile avute în vedere între partenerii proiectului; la ora actuală acest corpus însumează 683.357 cuvinte. Descrierea fiecărui obiect de învățare se face printr-un nume, un set restrâns de cuvinte cheie³ și datele privind drepturile de autor⁴. Din motive statistice, pentru fiecare resursă se calculează numărul de cuvinte.

3. Formatul obiectelor de învățare

Resursele lingvistice au avut, la momentul colectării lor, diferite formate (nivelul 1 din Figura 1.): *.doc, .pdf, .html, .txt* etc. Pentru o prelucrare unitară, s-a hotărât definirea unui format comun la care să fie aduse toate resursele partenerilor. Acest format a fost unul de tip XML și el urma să păstreze doar puține informații relative la formatarea documentelor (precum fontul subliniat, înclinat sau îngroșat), adică atâtea câte s-ar putea dovedi utile în extragerea automată a cuvintelor cheie sau a definițiilor (de exemplu este foarte mare probabilitatea ca un cuvânt subliniat sau îngroșat să fie un termen cheie).

Deoarece colecția de documente din proiect proveneau, așa cum s-a menționat, din nouă limbi diferite, fiecare cu convențiile proprii asupra setului de diacritice, s-a convenit asupra utilizării formatului UTF-8, care pare a fi cel mai potrivit păstrării unitare a unor colecții de documente multilingve în vederea unor prelucrări similare. Aducerea la formatul XML UTF-8 (notat Base-XML în Figura 1) a reprezentat, așadar, primul nivel de prelucrare a documentelor primare. Până la transformarea lor în obiecte de învățare, acestea au fost suferit ulterior prelucrări lingvistice (nivelul 2 din Figura 1.), adnotări

² http://consilr.info.uaic.ro/uploads_lt4e/

³ În etapele preliminare ale proiectului aceste cuvinte cheie au fost incluse în informațiile adăugate pentru fiecare document, pentru a ajuta la selecția domeniilor, deci implicit a documentelor, cu care se lucrează în etapele ulterioare. În prezent aceste cuvinte cheie selectate inițial manual nu mai sunt incluse în informația atașată fiecărui document.

⁴ Documentele vor putea fi făcute publice la sfârșitul proiectului, dar până atunci majoritatea lor au fost oferite de către autori doar pentru cercetare.

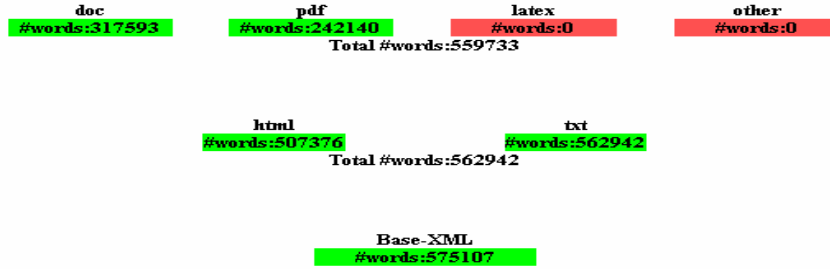
manuale și automate la cuvinte cheie și definiții (nivelul 3) și au fost plasate într-o ierarhie a schemelor de adnotare (Cristea et al., 2006).

Romanian resources

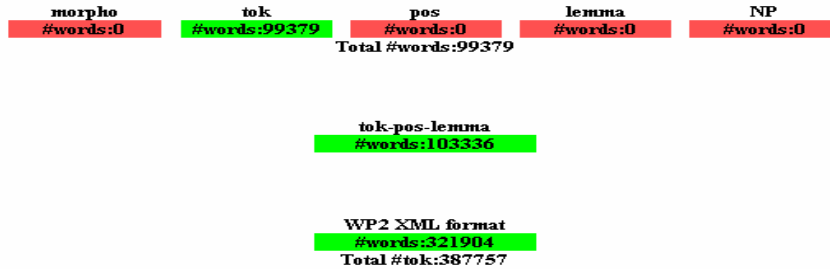
Total words:683357

Total words is computed as the sum of the number of words for all documents that are uploaded on the server in at least one format. In the hierarchy below, **Total #words** is computed as the sum of words for all documents that are uploaded on the server in at least one of the formats on the above level only.

Level 1: Initial document



Level 2: Linguistic annotation



Level 3: Keywords and definitions



Figura 1: Reprezentarea diferitelor nivele de adnotare a resurselor românești.

3.1. Nivelul lingvistic

Procesarea lingvistică a documentelor este importantă pentru a permite extragerea unor informații morfologice ce pot ajuta în detectarea automată a cuvintelor cheie și a definițiilor. Pe acest nivel, resursele au fost adnotate, folosind serviciul Web pus la dispoziție de ICIA⁵, astfel: împărțirea în unități lexicale, marcarea informației morfo-sintactice și marcarea lemelor formelor flexionate (Ion, 2006).

Ținând cont de contextul multilingv și de faptul că fiecare partener a venit cu instrumente de prelucrare diferite, ce manifestă, în general, formate de ieșire diferite, s-a convenit și pentru această etapă asupra unui format comun la care să se aducă adnotarea

⁵ Institutul pentru Cercetări în Inteligență Artificială, Academia Română, <http://www.racai.ro/>

lingvistică. Astfel, există etichete pentru segmentarea în paragrafe, în propoziții și la nivel de cuvânt. Fiecare cuvânt, de exemplu, trebuie marcat folosind etichete `<tok>`:

```
<tok rend="" base="Uniunii_Europene" ctag="Ed" id="t961">Uniunii_Europene</tok>
```

Atributele elementului *tok* sunt: *id*, un identificator unic; *rend*, care conține informația de formatare din formatul XML al primului nivel, dacă ea există; *base*, care conține forma lematizată a intrării lexicale și *ctag*, unde este trecută informația morfo-sintactică.

3.2. Nivelul post-lingvistic

Adnotările care au urmat nivelului lingvistic s-au făcut în două etape: o *adnotare manuală* și o *adnotare automată*. Motivul acestei duble adnotări a fost acela de a permite apoi compararea lor în scopuri de evaluare a adnotării automate. În fiecare din aceste două etape s-a avut în vedere adnotarea la cuvinte cheie și adnotarea la definiții.

3.2.1. Adnotarea cuvintelor cheie

Cuvintele cheie (unități lexicale formate fie dintr-un singur cuvânt, fie din expresii multi-cuvânt) sunt marcate (semi)automat de creatorii obiectelor de învățare; utilizatorii sistemului de învățare folosesc cuvintele cheie pentru a găsi documentele care conțin referi la anumite noțiuni. Din acest motiv cuvintele cheie trebuie să fie reprezentative pentru obiectul de învățare din care provine, să rezume subiectul textului sau să fie un obiectiv central al documentului. În adnotarea cuvintelor cheie s-a avut în vedere și posibilitatea ca aceleași noțiuni să fie uneori referite prin sinonime în același text.

Un exemplu de adnotare a unui cuvânt cheie este:

```
<markedTerm id="k36" comment="" dt="n" kw="y" status="">
<tok rend="" base="Uniunii_Europene" ctag="Ed" id="t961">Uniunii_Europene</tok>
</markedTerm>
```

Cuvintele cheie sunt marcate cu `<markedTerm>`. Pentru că aceasta este o etichetă comună cu cea folosită pentru marcarea termenilor definiți dintr-o definiție, diferența dintre cele două adnotări este dată de atributele *dt* și respectiv *kw*. Acestea pot lua valorile *y* (*yes*) și *n* (*no*). Pentru exemplul de mai sus, valoarea atributului *dt* este *n*, ceea ce înseamnă ca sintagma nu este un termen definit în acest context, iar valoarea atributului *kw* este *y*, ceea ce înseamnă ca sintagma este aici un cuvânt cheie.

Celelalte atribute ale elementului `<markedTerm>` sunt un *id*, a cărui valoare trebuie să fie unică în document, *status*, un atribut de confirmare, care poate lua valoarea *?* sau *??* dacă adnotatorul nu este sigur, respectiv este foarte nesigur, că ceea ce a marcat este corect, și *comment*, care poate conține comentarii.

În ceea ce privește adnotarea automată, în proiect s-au implementat trei metode pentru extragere a cuvintelor cheie: TF/IDF, Residual IDF (RIDF) și o versiune ajustată a RIDF (RIDF este înmulțit cu rădăcina pătrată a frecvenței termenilor). Programul generează un model de limbă, folosind fișierele adnotate manual, și aplică acest model pe restul documentelor (Lemnitzer, Degórski, 2006). Momentan suntem în stadiul de validare a rezultatelor obținute de extractor.

3.2.2. Adnotarea definițiilor

Prin definiție se înțelege o explicație concisă a înțelesului unui cuvânt sau a unei sintagme, o descriere a înțelesului unui concept sau a tipului său. O definiție are două părți: elementul definit și explicația propriu-zisă. Un exemplu de definiție extrasă din corpus este:

[*Cetățenia Uniunii Europene*]_{DEF PART1}, prevăzută în tratatul de la Roma și mai apoi în cel de la Maastricht [*este caracterizată de drepturi, de obligații și de implicarea în viața politică*]_{DEF PART2}.

unde elementul definit este *Cetățenia Uniunii Europene*, iar definiția propriu-zisă este marcată între paranteze []. Se observă că atributiva care determină termenul definit nu a fost considerată ca făcând parte din definiție. Notarea definițiilor care au astfel de întreruperi în secvența textuală (formate din mai multe părți) este exemplificată mai jos:

```
<definingText comment="" id="def37" status="" continue="y" def="dt35" part="1">
  <markedTerm id="dt35" comment="" dt="y" kw="n" status="">
    <tok rend=" /b, /p, p" base="cet&#259;&#355;enie" ctag="Ncfsry" id="t960">
      Cet&#259;&#355;enia </tok>
    <markedTerm id="k36" comment="" dt="n" kw="y" status="">
      <tok rend="" base="Uniunii_Europene" ctag="Ed" id="t961">Uniunii_Europene</tok>
    </markedTerm>
  </markedTerm>
</definingText>
<tok rend="" base="," ctag="COMMA" id="t962">,</tok>
<tok rend="" base="prevedea" ctag="Vmp--sf" id="t963">prev&#259;zut&#259;</tok>
<tok rend="" base="&#238;n" ctag="Spsa" id="t964">&#238;n</tok>
<tok rend="" base="tratat" ctag="Ncmsry" id="t965">Tratatul</tok>
<tok rend="" base="de_la" ctag="Spca" id="t966">de_la</tok>
<tok rend="" base="Roma" ctag="Np" id="t967">Roma</tok>
<tok rend="" base="(0.67)&#351;" ctag="Vmis1s" id="t968">&#351;i</tok>
<tok rend="" base="mai" ctag="Rp" id="t969">mai</tok>
<tok rend="" base="apoi" ctag="Rgp" id="t970">apoi</tok>
<tok rend="" base="&#238;n" ctag="Spsa" id="t971">&#238;n</tok>
<tok rend="" base="acela" ctag="Pd3msr" id="t972">cel</tok>
<tok rend="" base="de_la" ctag="Spca" id="t973">de_la</tok>
<tok rend="" base="Maastricht" ctag="Np" id="t974">Maastricht</tok>
<definingText comment="" id="def38" status="" continue="y" def="dt35" part="2">
  <tok rend="" base="fi" ctag="Vaip3s" id="t975">este </tok>
  <tok rend="" base="caracteriza" ctag="Vmp--sf" id="t976">caracterizat&#259;</tok>
  <tok rend="" base="de" ctag="Spsa" id="t977">de</tok>
  <tok rend="" base="drept" ctag="Ncfn" id="t978">drepturi</tok>
  <tok rend="" base="," ctag="COMMA" id="t979">,</tok>
  <tok rend="" base="de" ctag="Spsa" id="t980">de</tok>
  <tok rend="" base="obliga&#355;ie" ctag="Ncfn" id="t981">obliga&#355;i</tok>
  <tok rend="" base="(0.62)&#351;" ctag="Ncmpry" id="t982">&#351;i</tok>
  <tok rend="" base="de" ctag="Spsa" id="t983">de</tok>
  <tok rend="" base="implicare" ctag="Ncfsrn" id="t984">implicare</tok>
  <tok rend="" base="&#238;n" ctag="Spsa" id="t985">&#238;n</tok>
  <tok rend="" base="via&#355;&#259;" ctag="Ncfsry" id="t986">via&#355;a</tok>
  <tok rend="" base="politic" ctag="Afpfsrn" id="t987">politic&#259;</tok>
</definingText>
```

Pentru elementul definit se folosește eticheta <markedTerm> cu atributul *dt="y"*, după cum s-a arătat în secțiunea anterioară. Pentru marcarea definiției este folosită eticheta <definingText> cu atribute asemănătoare etichetei <markedTerm>, dar având în plus atributul *part* pentru a marca definițiile întrerupte (valoarea acestui atribut este 0 dacă definiția este continuă, sau un număr începând cu 1 și continuând să crească pentru fiecare parte, dacă definiția are alte elemente intercalate) și atributul *def* pentru a indica *id*-ul elementului definit, care ajută și la unificarea definițiilor formate din mai multe părți. În exemplul de mai sus, definiția se referă la termenul definit cu *id*-ul *dt35*.

Adnotarea automată a definițiilor se face utilizând o gramatică realizată de fiecare partener pentru limba respectivă. În afară de problemele care apar datorită dificultății de a surprinde toate modurile de exprimare a unei definiții (mai ales dacă se evită lexicalizarea în exces), de a trata definițiile întrerupte, de a stabili unde se termină o definiție etc., mai apar și probleme de ordin tehnic, care țin de modul de redactare a documentelor. Astfel, în exemplul dat anterior, în textul de intrare virgula nu apare la sfârșitul atributivei, așa cum cer normele gramaticale, lucru care poate crea dificultăți unei reguli de detecție a propozițiilor intercalate.

4. Concluzii

În cadrul proiectului LT4eL s-au colectat 56 de documente care însumează peste 600.000 de cuvinte. Acestea au fost aduse la un format unitar XML și adnotate la nivel lingvistic (segmentare în cuvinte, adnotare morfo-sintactică și lematizare). Ulterior, o parte din ele a fost adnotată manual la cuvinte cheie și definiții. Momentan se lucrează la îmbunătățirea rezultatelor obținute cu extractorul automat de cuvinte cheie și de definiții și la validarea acestora. Una din direcțiile de lucru viitoare implică adnotarea semantică a obiectelor de învățare conform ontologiei dezvoltate în cadrul proiectului, proces deja început pentru limba engleză.

Referințe bibliografice

- Cristea, D., Forăscu, C., Pistol, I. (2006). Requirements-Driven Automatic Configuration of Natural Language Applications. In Bernadette Sharp (Ed.): Natural Language Understanding and Cognitive Science, *Proceedings of the 3rd International Workshop on Natural Language Understanding and Cognitive Science - NLUCS 2006*, in conjunction with ICEIS 2006, Cyprus, Paphos, May 2006. INSTICC Press, Portugal. ISBN: 972-8865-50-3.
- Ion, R. (2006). *Metode de dezambiguizare semantică automată. Aplicații pentru limbile engleză și română*. Teză de doctorat în curs de susținere la Academia Română.
- Lemnitzer, L., Degórski, L. (2006): Language Technology for eLearning -- Implementing a Keyword Extractor. The fourth EDEN Research Workshop "Research into online distance education and eLearning. Making the Difference", 25-28 October, 2006 in Castelldefels, Spain
- Monachesi, P., Cristea, D., Evans, D., Killing, A., Lemnitzer, L., Simov, K., Vertan, C. (2006). Integrating Language Technology and Semantic Web techniques in eLearning. *Proceedings of ICL 2006*.
- Pistol, I., Trandabăț, D., Iftene, A., Cristea, D., Forăscu, C. (2006). Prelucrarea resurselor românești în cadrul proiectului LT4eL. În acest volum.

TEHNICI DE VALIDARE ȘI CORECȚIE FOCALIZATĂ A ADNOTĂRII MORFO-SINTACTICE ÎN CORPUSURI DE MARI DIMENSIUNI

DAN TUFIS, ELENA IRIMIA

Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București

{tufis, elena}@racai.ro

Rezumat

Articolul descrie procesul de realizare și corectare a RoCo-News, corpus jurnalistic pentru limba română, de dimensiune medie, abundent în nume proprii, numerale și entități denumite. Textul, inițial ne-procesat, a fost segmentat cu MtSeg, apoi adnotat morfo-sintactic cu tagger-ul TNT. Ulterior, RoCo-News a fost re-adnotat și lematizat cu tagger-ul TTL realizat la RACAI, iar în final validat și corectat. Datorită resurselor umane limitate, a constrângerilor temporale și a dimensiunilor corpusului, validarea de mână a fiecărei unități lexicale a fost exclusă. Etapa de validare a necesitat o metodologie coerentă pentru identificarea automată a cât mai multe erori de adnotare morfo-sintactică și lematizare. Procesul de validare manuală s-a concentrat apoi doar pe aceste posibile erori descoperite în mod automat.

1. *Introducere. Descrierea corpusului RO-CO News*

Pentru cercetătorii în domeniul lingvisticii computaționale, importanța dezvoltării resurselor, în special a corpusurilor, este evidentă. Reprezentând o “colecție de fragmente de text într-o anumită limbă, selectate și ordonate potrivit unor criterii lingvistice, în scopul de a fi utilizate ca mostre ale limbii respective” (Sinclair, 1991), un corpus ajută la formarea unei imagini comprehensive asupra limbii sincronice.

Un corpus jurnalistic este o bună sursă de informare asupra cuvintelor nou intrate în vocabularul unei anumite limbi, despre entitățile denumite, abrevierile comune și multe alte aspecte ale stilului funcțional. În cele ce urmează, vom descrie corpusul RoCo-News, care, în forma sa actuală, reprezintă rezultatul unei munci de un an de validare și corectare a procesării și adnotării automate.

RoCo-News este un corpus jurnalistic pentru limba română de dimensiuni medii. Conține aproximativ 7 milioane de unități lexicale, numărul de unități distincte depășind 231.000. Diferitele articole din corpus, disponibile inițial în diverse formate (doc, rtf și pdf) au fost convertite în formatul ASCII, cu diacriticele codificate ca entități SGML.

O analiză preliminară a textului a relevat abundența numelor proprii și a expresiilor numerice, date deloc surprinzătoare pentru un registru jurnalistic. Specific acestui tip de text, titlurile și numele de autori apar ca paragrafe distincte, iar datorită structurii gramaticale parțiale pot conține mai multe erori de adnotare decât paragrafele reale.

Textul neprelucrat a fost segmentat cu segmentatorul MtSeg, dezvoltat în contextul proiectului MULTTEXT (<http://aune.lpl.univ-aix.fr/projects/Multext/>), bazându-ne pe

resursele de segmentare dezvoltate de proiectul MULTTEXT-East (<http://nl.ijs.si/ME/>) și adnotat morfo-sintactic folosind tagger-ul TnT (Brants, 1988). Tagger-ul a fost antrenat pe un corpus validat manual care include romanul “O mie nouă sute optzeci și patru” al lui Orwell (aproximativ 110.000 unități lexicale), “Republica” lui Platon (aproximativ 140.000 de unități lexicale) și mai multe articole din câteva jurnale cu acoperire națională (aproximativ 140.000 de unități lexicale). Setul de etichete (tagset) al modelului de limbă este derivat dintr-un tagset mai mare, complet compatibil cu specificațiile morfo-sintactice MULTTEXT-East. Tagset-ul redus utilizat în corpusul RoCo-News este tagset-ul ascuns din metodologia de adnotare stratificată (“tiered tagging”, vezi (Tufiș, 1999) și (Tufiș și Dragomirescu, 2004), pentru mai multe detalii). Tagset-ul redus conține 93 de etichete pentru cuvinte și 10 etichete pentru punctuație.

În continuare, corpusul adnotat a fost lematizat. Procesul de lematizare a fost, în esență, o procedură de căutare într-un lexicon de dimensiuni mari, conținând peste 600.000 de intrări de forma: <formă-cuvânt> <lemă> <etichetă>. Pentru unitățile lexicale care nu se regăsesc în această resursă (și care nu sunt marcate ca nume proprii), lema este furnizată de lematizorul cuprins în modulul statistic de adnotare și lematizare TTL (Ion, 2006), realizat la RACAI. Aplicația folosește un set de reguli (specific fiecărei categorii gramaticale flexionare) induse automat din lexicon, care generează leme candidat pentru cuvântul necunoscut, și apoi modele Markov (antrenate pe leme din lexicon) pentru a ordona candidații. Candidatul cu cea mai mare probabilitate câștigă. Procedura funcționează foarte bine, cele mai multe dintre erori sunt în cazul cuvintelor necunoscute care aparțin paradigmelor flexionare neregulate (Tufiș, 1989) sau atunci când adnotarea morfo-sintactică a formei ocurență a fost greșită. În ansamblu, ținând cont de vasta acoperire a lexiconului și de rata mică de eroare a lematizorului statistic, probabilitatea unei erori de lematizare este neglijabilă.

2. Trei tehnici utilizate pentru identificarea erorilor posibile în RoCo-News

2.1. Lematizare și re-segmentare

Lematizarea este un procedeu mai simplu decât adnotarea și, de aceea, se poate face automat cu mai multă acuratețe. În mare parte dintre cazuri, perechea formată de ocurența unui cuvânt și una dintre etichetele sale legale ar trebui să identifice în mod unic lema acelei unități lexicale, dacă ea este înregistrată în lexicon. Cum majoritatea cuvintelor dintr-un text nou sunt cuvinte de uz general, deci teoretic prezente în lexicon, putem presupune că lematizarea unui text folosind această procedură se face cu mare acuratețe. Lexiconul este actualizat în mod constant, pe măsură ce întâlnim cuvinte noi în textele la care lucrăm. Totuși, lematizorul statistic apelat pentru cuvintele necunoscute poate produce leme greșite, în special dacă unitățile lexicale sunt adnotate în mod eronat. Pentru a nu introduce erori în lexicon, tripletele < **formă-cuvânt lema etichetă** > sunt subiectul validării unui expert, înainte de a fi incluși în lexicon.

Procedura de lematizare, re-segmentare și corectare a erorilor identificate în decursul acestui proces este descrisă pe scurt în cele ce urmează:

a) Dacă unitatea lexicală curentă nu este marcată printr-un asterisc și a fost adnotată cu o etichetă de semn de punctuație sau de categorie gramaticală fără flexiune, lema este identică formei ocurență.

b) Dacă unitatea lexicală curentă este marcată de tagger ca necunoscută, verificăm dacă eticheta morfo-sintactică este NP (nume propriu), caz în care lema este din nou considerată ca identică formei ocurență a lemei. Raționamentul este că în limba română, numele proprii (cele străine și cele masculine) sunt rareori flexionate. Pe de altă parte, numele proprii feminine pot avea flexiune, dar cele mai frecvente dintre ele se află deja în lexicon. Unitățile lexicale adnotate în mod consecutiv prin eticheta NP sunt concatenate și considerate drept o singură unitate lexicală, iar lema sa este concatenarea lemelor unităților lexicale care o compun. Tripletele necunoscute <formă-cuvânt NP lemă> obținute prin concatenare au fost adăugate în fișierul NumeProprii. Toate numele proprii din acest fișier au fost validate, iar erorile corectate. Corecțiile au fost operate și în corpus. Câteva dintre erorile tipice au fost reprezentate de unitățile lexicale ale căror caractere erau în întregime scrise cu majuscule (făceau parte din titluri de articole) sau erau nume proprii independente, a căror concatenare nu era necesară.

c) Dacă o unitate lexicală nerecunoscută de TnT nu este etichetată ca și NP, atunci este căutată în lexicon (mult mai mare decât lexiconul tagger-ului) împreună cu eticheta sa. În cazul în care este găsită, lema este copiată din respectiva intrare în lexicon. Altfel, lematizorul probabilistic este apelat pentru unitatea lexicală curentă iar tripletul <formă-cuvânt etichetă lemă> este salvat în fișierul denumit NuSuntÎnLexicon, pentru inspecție și validare ulterioară. Conținutul acestui fișier a fost clasificat și analizat în ordine descrescătoare a frecvenței tripletelor sale.

Au fost identificate mai mult de 20.000 de erori, majoritatea datorate conversiei eronate a anumitor diacritice în entități SGML. O astfel de eroare sistematică, o dată observată, este relativ ușor de corectat. Un caz special de unități lexicale nerecunoscute de tagger este reprezentat de numere. Ele sunt adnotate în mod sistematic ca și numerale, dar există numeroase cazuri în care segmentarea acestora a fost eronată, datorită utilizării ca separator între grupurile de trei cifre a caracterului spațiu, în loc de virgulă sau punct. Astfel, segmentatorul a considerat că are de a face cu unități lexicale diferite. Pentru astfel de cazuri am procedat la concatenarea grupurilor de cifre.

Printre unitățile lexicale necunoscute am găsit de asemenea și adrese web sau de e-mail. Aceste grupuri de cuvinte speciale au fost adnotate sistematic de către TnT ca și NN. Importanța de necontestat pe care aceste tipuri de unități textuale și-au câștigat-o motivează decizia de a introduce în tagset două etichete noi: NNWEB și NNMAIL. Toate ocurențele adreselor de web și e-mail au fost re-adnotate corespunzător acestui tagset extins.

2.2. Utilizarea analizei clasei închise pentru identificarea erorilor

Divizarea categoriilor lexicale în două tipuri diferite de clase este tradițională în lingvistică: clasele închise sunt acelea enumerabile (ex: clase precum determinatori, prepoziții, verbe modale sau auxiliare), în timp ce clasele deschise sunt categoriile mari și productive precum verbele, substantivele și adjectivele.

(Dickinson și Meurers, 2002) au exploatat ideea că, pentru detectarea erorilor, se poate utiliza în mod practic conceptul de clasă închisă. Poate fi ușor de observat, susțin ei, că aproximativ jumătate din etichetele oricărui tagset corespund claselor lexicale închise. O categorie lexicală clasă închisă conține un număr redus de cuvinte, enumerabil. În mod frecvent, aceste cuvinte pot face parte din mai multe categorii clase închise (ex.: în limba română există cuvinte care pot fi în același timp: prepoziții sau conjuncții; prepoziții și auxiliare etc.). În funcție de granularitatea tagset-ului, o categorie clasă închisă poate acoperi mai multe etichete (ex.: tipuri de conjuncții, prepoziții, pronume). Ținând cont de frecvența mare a cuvintelor clase închise, pentru un corpus mare putem presupune că aceste tipuri de cuvinte apar în cele mai multe (dacă nu în toate) dintre contextele posibile și este, deci, de așteptat ca toate etichetele acestora să fie regăsite în corpus. Bazându-ne pe aceste considerații, am decis să facem o serie de teste cu privire la cuvintele clasă închisă din RoCo-News. Astfel, am extras din lexicon o listă L1, de etichete clasă închisă, fiecare dintre ele indexând mulțimea cuvintelor care ar putea să primească acea etichetă. Din această listă, am calculat o alta, L2, conținând cuvinte din L1 ce indexează două sau mai multe etichete de clasă închisă. Apoi am extras din RoCo-News toate perechile <cuvânt, etichetă> a.î. eticheta să fie de clasă închisă. Dacă *cuvânt* nu a fost în mulțimea din L1 indexată de *etichetă*, am verificat ocurența respectivului cuvânt în context. În marea majoritate a cazurilor, am găsit o eroare de adnotare, dar ocazional am descoperit erori și în lexicon (o posibilă etichetă de clasă închisă nu a fost înregistrată pentru anumite cuvinte). Bazându-ne pe L2, am extras toate cuvintele care au fost văzute în corpus doar cu un subset al etichetelor de clasă închise posibile. Au fost descoperite din nou anumite erori în lexicon (cuvinte care erau în mod eronat asociate cu anumite etichete de clasă închisă).

2.3. Utilizarea evaluării auto-referențiale pentru o mai bună identificare a erorilor

Cea de-a treia tehnică utilizată în corectarea corpusului RoCo-News se bazează pe ipoteza evaluării auto-referențiale (Tufiș, 1999), care spune că un corpus consistent și corect adnotat, re-adnotat cu modelul de limbaj învățat din el însuși (evaluare auto-referențială), ar trebui să aibă majoritatea unităților lexicale adnotate în mod identic. Procentul de etichete identice depinde de dimensiunea corpusului, dar de obicei este mai mare de 97.5%-98%.

După ce am efectuat corecturile descrise în secțiunile precedente, am luat această versiune ca referință pentru procedura de evaluare auto-referențială descrisă în cele ce urmează. Am antrenat tagger-ul TnT pe corpusul referință, construind un nou model de limbă. Am re-adnotat RoCo-News cu acest nou model de limbă și am comparat noua adnotare cu cea de referință. Am descoperit 96.8% unități lexicale adnotate identic și am extras diferențele. Sortând diferențele în ordine inversă a frecvenței lor în corpus, am examinat pe rând, în context, primele 100 tipuri de diferențe (reprezentând aproximativ 8-10.000 de ocurențe diferite) iar expertul însărcinat cu validarea lor a stabilit care dintre cele două etichete era cea corectă (dacă vreuna dintre ele a fost corectă). Unele dintre diferențe au fost explicate prin inconsistența sau incompletitudinea corecturilor din etapele precedente. Alte diferențe au apărut deoarece corecturile au modificat contextele pentru unitățile lexicale vecine și astfel, conform modelului de limbă auto-referențial, multe dintre unități au apărut în contexte diferite și au primit etichete diferite. Corectarea

tuturor erorilor descoperite în analiza primelor 100 tipuri de diferențe încheie procedura. Având în vedere dimensiunea corpusului, realizarea ei necesită foarte mult timp. Procedura a fost repetată de mai multe ori, cu o scădere continuă a numărului de diferențe; la final, numărul de diferențe s-a stabilizat la 0.4% din dimensiunea corpusului (rămânând 25.500 diferențe dintre care 6353 distincte).

3. Concluzie

Am descris o procedură semi-automată prin ale cărei mijloace am construit un corpus jurnalistic pentru limba română cu înalt nivel de acuratețe a adnotării și lematizării. Deși rezultate analizelor și natura erorilor sunt dependente de limbă, de tagger și de tagset, scheletul acestei abordări poate fi adaptat și aplicat cu ușurință într-un alt context. Tipul de analiză pe care am descris-o poate oferi indicații importante despre cuvintele/tipurile de cuvinte care ar putea fi adnotate nesatisfăcător într-un alt corpus din același registru.

Metoda nu asigură eliminarea tuturor erorilor existente, dar câștigul în acuratețe este substanțial iar faptul că nu este nevoie de examinarea cuvânt cu cuvânt constituie, din punct de vedere al economiei de timp și de resurse umane, un mare avantaj, expertul care validează putându-se concentra pe acele erori pe care procedurile le evidențiază ca fiind frecvente și/sau importante.

Referințe bibliografice

- Brants, T. (1998). TnT - A Statistical Part-of-Speech Tagger. Instalation and User Guide, University of Saarland, *Computational Linguistics*, March 1998.
- Dickinson, M., Meurers, W. Detmar (2003). Detecting Errors in Part of Speech Annotation. In *Proceedings of the 11th conference of the EACL-03*, Budapest, Hungary.
- Ion, R. (2006). *Metode de dezambiguizare semantică automată. Aplicații pentru limbile engleză și română*, Teză de doctorat în curs de susținere la Academia Română, București, România, 145 p.
- Sinclair, J. (1991). *Corpus, Concordance, Collocation*, Oxford University Press.
- Tufiș, D. (1989). It Would Be Much Easier If *WENT* Were *GOED*., Harry Somers, Mary McGee Wood (eds.), *Proceedings of the 4th European Conference of the Association for Computational Linguistics*, Manchester.
- Tufiș, D. (1999). Tiered Tagging and Combined Classifiers. F. Jelinek, E. Nöth (eds) *Text, Speech and Dialogue, Lecture Notes in Artificial Intelligence 1692*, Springer, 1999, pp. 28-34.
- Tufiș, D., Dragomirescu, L. (2004). Tiered Tagging Revisited. *Proceedings of the 4th LREC Conference*. (pp. 39—42). Lisbon, Portugal.

ROGER – UN CORPUS PARALEL ALINIAT

MONICA GAVRILĂ, NATALIA ELIȚA

*Departamentul de Informatică, Facultatea de Matematică, Informatică și Științele
Naturii, Universitatea din Hamburg*

{gavrila, elita}@informatik.uni-hamburg.de

Rezumat

Colecțiile de texte în format electronic (corpusurile) sunt folosite foarte des în comunitatea procesării limbajului natural. Deoarece este o resursă primară, alegerea unui corpus potrivit influențează direct rezultatul pe care utilizatorul dorește să îl obțină și, implicit, prin folosirea acestui rezultat, cercetările și rezultatele ulterioare.

Această lucrare prezintă un corpus multilingv (**român-german-englez-rus**), aliniat și paralel, de aproximativ 2300 propoziții.

Lucrarea este organizată în patru secțiuni. Prima secțiune descrie importanța unui corpus în prelucrarea limbajului natural (PLN), cu accente asupra traducerii automate bazate pe exemple. Cea de-a doua secțiune prezintă resurse existente și problemele întâmpinate în găsirea unui corpus adecvat, iar cea de-a treia informații asupra corpusului creat: RoGER. Concluziile și lista de referințe bibliografice încheie lucrarea de față.

1. Importanța corpusului în PLN și în traducerea automată bazată pe exemple

Colecțiile de texte în format electronic (corpusurile) sunt foarte des folosite în comunitatea procesării limbajului natural (PLN). Deoarece este o resursă primară, alegerea unui corpus potrivit influențează direct rezultatul pe care utilizatorul dorește să îl obțină și, implicit, prin folosirea acestui rezultat, cercetările și rezultatele ulterioare. De aceea, câteva aspecte sunt foarte importante:

- dimensiunea corpusului;
- tipul textului;
- relevanța textului;
- tipul corpusului (monolingv / multilingv, aliniat / nealiniat, etc.)

Corpusul este folosit atât pentru antrenarea unui sistem, cât și pentru testarea lui.

Corpusurile sunt utilizate pentru rezolvarea diverselor probleme din PLN: traducerea automată, analiza textelor cu metode statistice, dezambiguizarea sensurilor cuvintelor, construirea modelelor limbilor, etc.

În ceea ce privește traducerea automată bazată pe exemple, aceasta are la bază un corpus ce trebuie să îndeplinească anumite condiții legate de dimensiune, tipul textului, etc.

Pentru a elabora un astfel de sistem, în primul rând este nevoie de un corpus paralel aliniat¹. Din acest corpus sunt extrase exemple de traducere, care vor fi folosite ulterior la traducerea unor texte noi. Aceste exemple pot fi salvate în formate diferite (de exemplu: arbori, șabloane – *engl.*: „templates”, etc), iar în funcție de modul de salvare câteodată este nevoie de adnotarea corpusului.

Pentru extragerea automată de exemple este necesar a avea corpusul, și informația aferentă lui, într-un format ușor de accesat și procesat de calculator.

În traducerea automată bazată pe exemple, o mare importanță o are balanța corectă între lungimea și similaritatea exemplurilor din baza de exemple: cu cât exemplele sunt mai lungi, cu atât este mai greu să găsești un exemplu potrivit pentru textul ce urmează a fi tradus, și cu cât exemplele sunt mai scurte, cu atât crește posibilitatea de ambiguitate.

2. Resurse existente

Aplicația pe care dorim să o realizăm, și pentru care avem nevoie de un corpus, este un sistem de traducere automată bazată pe exemple, în care este specificată informația semantică, prin adnotarea semantică a corpusului. Adnotarea semantică se va face având la bază o ontologie.

În cercetarea noastră, inițial, am încercat să găsim un corpus care să îndeplinească cerințele traducerii automate bazate pe exemple: paralel și aliniat. În plus, doream ca acesta să fie în patru limbi: român-englez-german-rus, și domeniul descris în text să fie restricționat. Am dorit să avem cu corpus al cărui domeniu este restricționat pentru a ușura munca necesară realizării ontologiei.

Din resursele existente, analizate de noi, fac parte și cele din tabelul de mai jos:

Tabel 1. O parte a corpusurilor analizate

	Denumire	Conținut	Observații
1	Corpusul paralel Român-Englez (Rada Mihalcea www.cs.unt.edu/~rada/downloads.html)	Arhive de ziare	Traduceri incomplete, nu sunt toate cele patru limbi
2	Corpus Român-Englez-Rus (www.azi.md)	Colecție de știri	Traduceri incomplete, nu sunt toate cele patru limbi
3	Corpus paralel German-Englez (www.iccs.informatics.ed.ac.uk/~pkoehn/publications/de-news)	Colecție de știri	Traduceri incomplete, nu sunt toate cele patru limbi
4	JRC-Acuis (wt.jrc.it/lt/Acquis), detalii în (Steinberger et al., 2006)	Colecție de texte din legislația UE (1950-2005)	Lipsește rusa

¹ Un corpus paralel este un corpus în două sau mai multe limbi, în care textele dintr-o limbă sunt traduse în celelalte limbi. Un corpus paralel aliniat este un corpus paralel în care este realizată corespondența dintre traduceri.

	Denumire	Conținut	Observații
5	OPUS (logos.uio.no/opus/kdedoc.html)	Documentații, manuale	Traduceri parțiale

Din analiza acestor resurse, unele dintre problemele întâlnite au fost:

- informația eronată asupra conținutului textului;
- traduceri incomplete;
- traduceri incorecte;
- domeniul corpusului este prea extins,
- corpusul nu conține toate cele patru limbi dorite, etc.

3. *RoGER*

Motivația noastră de a crea RoGER constă în faptul că resursele existente și descrise în secțiunea anterioară nu corespund în totalitate cerințelor noastre asupra corpusului, din diverse motive: limbile considerate, traduceri inexacte, domeniul corpusului, etc.

RoGER este un corpus:

- paralel,
- aliniat la nivel de propoziție,
- specializat (domeniu tehnic) - textele sunt preluate dintr-un manual de utilizare a unui aparat electronic,
- multiligv: român – german - englez – rus,
- realizat în proporție de peste 80 % manual,
- neadnotat (la nivel semantic, sintactic sau morfologic),
- în care diacriticele sunt neglijate.

Textul inițial a fost procesat, în sensul că unele noțiuni au fost înlocuite cu "meta-noțiuni", astfel că: numerele au fost înlocuite cu *NUM*, denumirile de pagini web cu *WWW SITE*, imaginile cu *PICT*. De asemenea, pentru a ușura procesul de traducere automată bazată pe exemple (exemple salvate ca șabloane), unele abrevieri au fost extinse.

Pentru alinierea la nivel de propoziție și corectarea traducerilor a fost efectuată o verificare manuală a corpusului. Verificarea și corectarea traducerilor s-au realizat în momentul creării alinierii corpusului.

În tabelul de mai jos se pot găsi câteva date statistice referitoare la corpus:

Tabel 2. RoGER - statistici

	Engleză	Germană	Română	Rusă
Dimensiune corpus (propoziții)	2333	2333	2333	2333
Dimensiune corpus (cuvinte) ²	26096	25850	27142	22383
Dimensiune vocabular	2012	3104	3031	3883
Vocabular (cu nr. de apariții mai mare ca doi)	1231	1575	1698	1904
Lungimea medie a propoziției	11	11	11	9

Referitor la dimensiunea și scopul pentru care a fost creat (traducerea automată bazată pe exemple), RoGER poate fi caracterizat, conform datelor menționate în (Somers, 1999), ca un corpus de dimensiune medie.

El se situează la jumătatea listei (compusă din 30 exemple) menționate în (Somers, 1999), deasupra corpusului folosit în sistemul Gaijin (Veale and Way, 1997) - 1836 exemple.

Corpusul este salvat în format XML. Mai jos se găsește un exemplu:

```
<?xml version="1.0" encoding="UTF-8"?>
<sentences>
.....
  <sentence id="1010">
    <en>Press Options and some of the following options may be
available .</en>
    <de>Druecken Sie Optionen . und einige der folgenden Optionen
sind ggf. verfuegbar .</de>
    <ro>Apasati Optiuni dupa care unele din urmatoarele optiuni
pot fi disponibile .</ro>
    <ru>Нажмите Вар-нты и выберите одну из перечисленных ниже
функций .</ru>
  </sentence>
.....
</sentences>
```

În tabelul 3 sunt incluse câteva exemple de cuvinte foarte frecvente în corpus, și analizându-le, ne putem da seama că în mare parte ele sunt aceleași în toate limbile considerate. Multe dintre ele sunt prepoziții, conjuncții (și (ro), und (ge), and (en), и (ru)), dar sunt și cuvinte purtătoare de sens - *engl.*: "content words" - (selectați (ro), wahlen (ge), select (en), выберите (ru)).

² Cuvintele se numără la nivel de șir de caractere („selectat” și „selectate” se numără ca și două cuvinte diferite).

Tabel 3. Cuvintele cele mai frecvente în RoGER

Romană	Germană	Engleză	Rusă
de, și, pentru, în, la, selectați, apăsați, este, un, dacă, să, pe, nu, o, care, Dvs, pagina, meniu, din, setari	Sie, und, die, der, wahlen, druecken, das, oder, um, zu, den, wenn, auf, fuer, auf, von, in	the, to, and, select, a, in, press, or, you, for, of, is, on, your, service, settings, menu	и, в, для, нажмите, на, выберите, или, Меню, Если

Cuvintele cele mai frecvente sunt:

- română: "de" (1459 ori)
- engleză: "the"(2075 ori)
- germană:"Sie" (1677 ori)
- rusă: "и" (799 ori)

Pentru a fi folosit ulterior în sistemul de traducere automată bazată pe exemple, asupra căruia lucrăm, intenționăm să extindem corpusul până la (minim) 2600 propoziții, introducând artificial ambiguități³, și să îl adnotăm semantic.

4. Concluzii

RoGER, corpusul realizat de noi - multilingv, paralel, aliniat - reprezintă o resursă utilă nu numai în antrenarea și testarea unui sistem de traducere automată bazată pe exemple, ci și pentru alte aplicații de PLN.

Referințe bibliografice

- Steinberger R., Pouliquen B., Widiger A., Ignat C., Erjavec T., Tufiş D., Varga D. (2006). The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. *LREC'2006*, pag. 2142-2147. Genoa, Italia, 24-26 mai.
- Herold S. (1999). Review Article: Example-based Machine Translation. *Machine Translation* 14: 113 – 157.
- Tony V., Way W. (1997). Gaijin: A Template-Driven Bootstrapping Approach to Example-Based Machine Translation. *NeMNL97. New Methods in Natural Language Processing*, Sofia, Bulgaria, septembrie.

³ Pentru a putea demonstra utilitatea semanticii în traducerea automată bazată pe exemple.

TIMEBANK 1.2: O VERSIUNE ADNOTATĂ ÎN LIMBA ROMÂNĂ

CORINA FORĂSCU^{1,2}, RADU ION²

¹*Facultatea de Informatică, Universitatea “Al.I.Cuza”, Iași*

²*Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București*

corinfor@info.uaic.ro, radu@racai.ro

Rezumat

Informația temporală s-a dovedit foarte relevantă mai ales în Prelucrarea Limbajului Natural. De proveniență lingvistică, teoriile temporale au fost studiate și formalizate cu predilecție pentru limba engleză. Lucrarea prezintă cercetările efectuate pentru obținerea corpusului paralel TimeBank, englez-român, care să fie folosit, printr-un import automat, la validarea acestor teorii pentru limba română. Corpusul va constitui și baza de lucru pentru dezvoltarea unor instrumente de prelucrare temporală a limbajului natural.

1. Introducere

Teoria logicii temporale s-a dovedit extrem de relevantă încă din anii '70, mai ales în Inteligența Artificială. Posibilitatea de a identifica și analiza informația temporală este de mare importanță pentru multe dintre aplicațiile Prelucrării Limbajului Natural precum rezumarea multi-document, sistemele de întrebare răspuns, structura temporală a discursului, regăsirea și extragerea informației, traducere automată, etc.

Dacă pe plan internațional în domeniul marcării informației temporale, s-au dezvoltat atât standarde de adnotare precum Timex2 (Ferro et al., 2005a) și TimeML (Sauri et al., 2006), cât și corpusuri adnotate conform cu acestea, precum ACE-TERN (Ferro et al., 2005b) sau TimeBank (Pustejovsky et al., 2006), predominant pentru limba engleză, pentru limba română cercetări anterioare (Forăscu, Solomon, 2004) au arătat că o adnotare manuală a unui corpus în limba română este foarte costisitoare, atât ca timp cât și ca resurse umane implicate și, mai mult, adnotările sunt deseori incomplete.

Lucrarea prezintă activitățile de creare ale unui corpus paralel englez-român, folosind ca sursă corpusul englez de știri TimeBank, pentru ca adnotarea temporală din acesta să fie apoi transferată în varianta română a corpusului, obținându-se astfel atât un corpus paralel cât și o sursă exemplificată de fundamentare a teoriilor temporale pentru română.

Secțiunea a doua a lucrării motivează necesitatea unui corpus paralel adnotat temporal, incluzând și fundamente ale informației temporale și ale principalului standard de adnotare temporală. În continuare sunt prezentate corpusul englezesc, modul de realizare a variantei românești a acestuia (secțiunea 4), prelucrările realizate asupra corpusului paralel (secțiunea 5). În încheiere se dezvăluie câteva obiective de viitor ale cercetării.

2. Necesitatea unui corpus paralel, adnotat temporal

Cînd un corpus este privit ca o colecție de documente selectate și ordonate conform unor criterii lingvistice stabilite, el permite punerea în evidență, informarea și fundamentarea unor teorii lingvistice specializate. Cum la ora actuală engleza este limba cu cea mai densă realizare de documente adnotate, ea este deseori utilizată ca sursă din care să se transfere adnotări specifice asupra altor limbi. Pentru limba română o serie de corpusuri paralele au fost deja create (Cristea, Forăscu, 2006), însă niciunul care să permită evidențierea informației temporale în limbajul natural.

2.1. Informație și adnotare temporală

Cf. (Mani et al., 2005), informația temporală este reprezentată în limbajul natural prin:

- expresii temporale exprimate prin grupuri nominale, prepoziționale sau adverbiale – ore (timp al zilei), date, durate: *acum trei ore, mai 1984, anii 90, 5 februarie 2007*, etc.; aceste expresii temporale referă timpul ca:
 - o punct (moment): Am deschis ușa la ora doisprezece,
 - o interval: *Am fost plecat ieri*.
- expresii ce denotă evenimente exprimate în principal, pe lângă adjective, clauze predicative sau grupuri frazale prepoziționale, prin:
 - o propoziții, mai exact prin centrul (eng. *head*) sintactic, anume verbul principal: *Ion a plecat la munte*.
 - o grupuri nominale: Greva va continua și în zilele următoare.

Expresiile ce denotă timpul pot avea:

- referințe explicite (specificate), care referă la o intrare într-un sistem calendaristic / orar: *amiază, 11.10.2006 (midday, 11.10.2006)* ;
- referințe implicite (sub-specificate) - pot fi evaluate doar prin intermediul unui timp indexat: *anul viitor, săptămâna trecută, acum două ore*;
- referințe vagi (nespecificate, neancorate), care nu pot fi corelate cu un punct sau interval exact de timp: *după-amiază, în câteva săptămâni, acum câteva zile*.

Evenimentele exprimate prin verbe pot fi temporal ancorate:

- indirect, prin categoria morfologică a timpului și
- direct, prin modificatori adverbiali (adverbe de timp și frecvență, grupuri nominale și prepoziționale și clauze subordonate).

Pentru a codifica toate tipurile de expresii temporale, evenimente și relații între acestea, a fost creat standardul TimeML (Pustejovsky et al., 2006), ale cărui fundamente s-au pus încă din 2002. Standardul reunește multe dintre eforturile anterioare de adnotare temporală, diferind de acestea prin separarea reprezentării evenimentelor și a expresiilor temporale de legăturile de ancorare, ordonare sau dependență ce apar în texte.

Standardul TimeML definește 7 etichete: EVENT, MAKEINSTANCE (pentru evenimente și instanțe ale acestora – doar instanțele vor participa în legături temporale), TIMEX3 (pentru expresii temporale de tip DATE, TIME, DURATION, SET, complet specificate, non- și sub-specificate), SIGNAL (pentru elemente lexicale de legătură) și TLINK, ALINK, SLINK (pentru legături temporale, aspectuale și respectiv de subordonare între expresii și evenimente).

3. Corpusul TimeBank – versiunea engleză

Realizarea corpusului TimeBank a început în 2002 în cadrul proiectului TERQAS¹. Corpusul conține în versiunea actuală 183 de fișiere de rapoarte de știri în limba engleză, adnotate conform cu TimeML v.1.2. (Pustejovsky et al., 2006). Documentele provin din evaluarea rezumatelor DUC 2001 și din corpusurile ACE incluse în cataloagele LDC2003T11 și LDC99T42. Documentele conțin și alte marcaje XML: formatul documentelor, informație structurală, nume de entități (ENAMEX, NUMEX din MUC7), marcaje de propoziție.

Adnotarea temporală inițială a corpusului este considerată „preliminară” întrucât s-a arătat (Boguraev, Ando, 2006) că apar greșeli sistematice datorate dimensiunii relativ reduse a corpusului și datorate inconsistențelor în adnotare: legături temporale sau de subordonare inconsistente sau incomplete, clasificarea evenimentelor – în perfectare, adnotare incompletă a timpului și aspectului unor evenimente.

TimeBank 1.2. este versiunea actuală – din 2006² - a corpusului, conformă cu specificațiile TimeML 1.2.1. Structura și adnotarea corpusului sunt, în esență, aceleași cu cele din prima versiune a corpusului. TimeBank 1.2 este distribuit prin LDC (Pustejovsky et al., 2006). Statisticile pe TimeBank 1.2 sunt ilustrate în Tabelul 1.

Adnotarea documentelor a început cu o fază de preprocesare, când unele articole lexicale de tip evenimente (EVENT) și semnale (SIGNAL) au fost marcate cu unele clase, timpuri sau aspecte ale acestora. După această etapă 5 adnotatori umani au verificat preprocesările și corectitudinea adnotărilor conforme cu specificația TimeML 1.2.1.

Tabel 1: Statistici asupra corpusului TimeBank 1.2

TimeML tags	#	General	#
events	7935	propoziții	4715
instances	7940	unități lexicale	61042
timexes	1414	unități lexicale unice	10586
signals	688		
alinks	265		
slinks	2932		
tlinks	6418		
TOTAL	27592		

Corpusul TimeBank este în revizie continuă, pentru următoarele distribuții avându-se în vedere: evenimentele compuse, legăturile dintre argumente, evenimentele generice, relațiile temporale între data creării documentului și evenimentele de tip REPORTING, o distincție mai clară între data creării și data publicării unui articol de știri.

4. Crearea corpusului pentru limba română

Textul englezesc a fost repartizat inițial în vederea traducerii la două masterande în Lingvistică Computațională, Facultatea de Informatică Iași, cu un set minimal de recomandări, pentru a obține traduceri unitare și alinieri satisfăcătoare cu originalele.

Ori de câte ori a fost posibil, traducerile au fost unu la unu. Alinierea la fraze/propoziții s-a obținut astfel direct prin notările care au marcat traducerile. S-a recomandat folosirea

¹ *Temporal and Event Recognition for Question Answering Systems*, <http://www.timeml.org/site/terqas/index.html>

² <http://www.timeml.org/site/timebank/timebank.html>

echivalențelor de traducere cu aceeași parte de vorbire, cuvintele românești trebuind să fie cât mai “apropiate” de corespondentele lor englezești: atunci când cuvântului englezesc îi poate fi asociat în românește un cognat, acesta va fi preferat unei expresii (*sporadic* -> *sporadic* și nu *mai rar*). S-au tradus toate cuvintele și nu s-au introdus în traducere, din motive stilistice, cuvinte sau expresii fără corespondent în engleză. S-a folosit scrierea cu diacritice, conformă cu normele lingvistice în vigoare. Timpurile verbale s-au păstrat pe cât posibil, modificările fiind acceptate doar pe temeuri lingvistice, nu stilistice. S-a păstrat formatul din engleză pentru date, momente ale zilei și numere.

Varianta actuală pentru limba română a fost verificată manual, urmărindu-se evitarea unor inconsistențe și lipsuri în traducere, care nu ar fi permis o aliniere a unor elemente temporale esențiale. În cele 4.715 propoziții sunt 65.375 unități lexicale (inclusiv semne de punctuație), din care 12.640 sunt unice.

5. Prelucrări ale corpusului paralel

5.1. Adnotări ale corpusurilor *TimeBank* englez și română

În vederea alinierii lexicale a celor două jumătăți ale corpusului, s-a utilizat o preadnotare unitară a textelor care să poată fi folosită de aliniatorul lexical YAWA (Tufiş et al., 2006). Această procesare preliminară se referă la segmentarea la nivel de cuvânt, adnotarea cu etichete morfosintactice și lematizarea textelor în engleză și română. Modulul TTL (Ion, 2006) oferă aceste adnotări și în plus, asigură o reprezentare uniformă a textelor adnotate în termenii codificării corpusului paralel într-un format XML similar cu formatul XCES (Ide et al., 2000).

Segmentarea la nivel de cuvânt trebuie să ia în calcul faptul că spațiul nu este singurul delimitator de cuvinte și nici nu este întotdeauna delimitator de cuvinte. Atât în engleză cât și în română există expresii idiomatice care trebuie considerate ca unități lexicale în procesul de aliniere (*bun simț*, *take a look*). Adnotarea morfosintactică se face cu ajutorul unui adnotator probabilistic care implementează adnotatorul TnT (Brants, 2000) bazat pe Modele Markov Ascunse. Setul de etichete morfosintactice este compatibil cu specificațiile MULTEXT-East³ fapt care permite reprezentarea uniformă a informației morfosintactice în engleză și în română.

Lematizorul implementat în TTL este de asemenea unul probabilistic. O leamnă candidată se generează pe baza unei mulțimi de reguli extrase automat dintr-un lexicon care conține pentru fiecare formă ocurență a unui cuvânt, lema și eticheta morfosintactică a acesteia. Lema unei noi forme ocurențe a unui cuvânt de o etichetă morfosintactică dată este lema cea mai probabilă dintre toate lemele candidate după Modelul Markov al tuturor lemelor de aceeași etichetă din lexicon (Ion, 2006).

Tot ca o cerință a alinierii lexicale, s-au recunoscut, folosind expresii regulate peste secvențe de etichete morfosintactice, grupuri nominale și prepoziționale nerecursive, compuși verbali (*s-a dus*), adjectivali (*cea mai frumoasă*) și adverbiali (*tare de tot*).

³ <http://nl.ijs.si/ME>

5.2. Alinierea lexicală a corpusurilor

Alinierea lexicală a corpusului paralel a fost realizată cu YAWA (Tufiş et al., 2006) pe ieşirea modului TTL. Corpusul paralel TimeBank 1.2 a fost aliniat la nivel de unitate lexicală din română în engleză urmându-se patru faze specifice acestui aliniator:

1. alinierea cuvintelor conţinut (substantive, verbe, adjective şi adverbe) folosind un dicţionar de echivalenţi de traducere extras automat (Tufiş, 2002);
2. pe scheletul de aliniere de la pasul anterior se aliniază cuvintele aflate în acelaşi grup sintactic cu cuvintele aliniate utilizându-se reguli de aliniere. De exemplu, dacă avem un substantiv românesc aliniat la unul englezesc care este precedat de un determinant, aliniază determinantul englezesc la substantivul românesc;
3. pe scheletul de aliniere de la pasul 2, aliniază toate blocurile de indecşi consecutivi care au rămas nealiniaţi (Tufiş et al., 2006);
4. corectează alinierea de la 3.

Fazele 2 şi 4 sunt evident dependente de perechea de limbi aliniate dar regulile de aliniere şi cele de corecţie nu sunt integrate în corpul aliniatorului astfel încât să poată fi schimbate atunci când se doreşte alinierea altei perechi de limbi.

6. Obiective viitoare

În vederea obţinerii unui transfer optim al adnotărilor temporale din limba engleză, corpusul paralel aliniat este în prezent validat manual în proporţie de 60%. O aliniere perfectă va fi folosită în continuare atât pentru îmbunătăţirea performanţelor aliniatorului, cât şi la importul adnotărilor TimeML în varianta română a corpusului. După validarea manuală a acestui import, liste de activatori⁴ (*eng. trigger*) lexicali vor fi extrase pentru a fi folosite, eventual în combinaţie cu metode specifice de învăţare automată, pentru crearea şi antrenarea unui adnotator temporal pentru limba română. Pentru evaluarea adnotatorului se vor avea în vedere şi alte domenii pe lângă cel de ştiri, precum beletristică, legislaţie etc. Ca planuri de lungă durată se pot menţiona folosirea adnotărilor temporale combinate cu cele de discurs pentru determinarea structurii temporale a discursului, rezumarea multi-document şi folosirea ontologiilor temporale pentru a obţine inferenţe despre evenimente în timp.

Mulţumiri. Autorii sunt recunoscători Ministerului Educaţiei şi Cercetării, de a cărui finanţare au beneficiat în cadrul proiectului CEEEX 29 ROTEL şi CEEEX 132 InterOb. Pentru sfaturile şi sprijinul primit, autorii mulţumesc coordonatorului comun de doctorat, prof. dr. Dan Tufiş, precum şi prof. dr. Dan Cristea.

Referinţe bibliografice

Armstrong, A. (1996). *Multext: Multilingual Text Tools and Corpora. Lexikon und Text*, pp. 107–119.

⁴ Cuvinte care semnaleză un anumit fenomen lingvistic; în acest caz, de exemplu: expresii temporale (*azi, septembrie*), semnale (*să, şi, că, când*).

- Boguraev, B., Ando, R. (2006). Analysis of TimeBank as a Resource for TimeML Parsing. In *Proceedings of LREC 2006*, Genoa, Italy, pp. 71-76.
- Brants, T. (2000). TnT – a statistical part-of-speech tagger. In *Proceedings of the 6th Applied NLP Conference, ANLP-2000*, Seattle, WA, pp. 224-231.
- Cristea, D., Forăscu, C. (2006). Linguistic Resources and Technologies for Romanian Language. In *Journal of Computer Science of Moldova*, Academy of Science of Moldova, vol. 14, nr. 1(40), pp. 34-73, ISSN 1561-4042.
- Ferro, L., Gerber, L., Mani, I., Sundheim, B., Wilson, G. (2005a). *TIDES 2005 Standard for the Annotation of Temporal Expressions*, April 2005.
- Ferro, L., Gerber, L., Hitzeman, J., Lima, E., Sundheim, B. (2005b). *ACE Time Normalization (TERN) 2004 English Training Data v 1.0*, Linguistic Data Consortium, Philadelphia, ISBN 1-58563-331-3.
- Forăscu, C., Solomon, D. (2004). Towards a Time Tagger for Romanian. In *Proceedings of the ESSLLI Student Session*, August 2004, Nancy, France.
- Ion, R. (2006). *Metode de dezambiguizare semantică automată. Aplicații pentru limbile engleză și română*. Teză de doctorat în curs de susținere la Academia Română.
- Ide, N., Bonhomme, P., Romary, L. (2000). XCES: An XML-based Encoding Standard for Linguistic Corpora. In *Proceedings of the Second International Language Resources and Evaluation Conference.*, pp. 825-830.
- Mani, I., Pustejovsky, J., Gaizauskas, R. (eds.) (2005). *The Language of Time: A Reader*. Oxford University Press, ISBN-13: 978-0-19-926853-5, May 2005.
- Pustejovsky, J., Verhagen, M., Sauri, R., Littman, J., Gaizauskas, R., Katz, G., Mani, I., Knippen, B., Setzer, A. (2006). *TimeBank 1.2*. Linguistic Data Consortium, Philadelphia, ISBN: 1-58563-386-0.
- Sauri, R., Littman, J., Knippen, B., Gaizauskas, R., Setzer, A., Pustejovsky, J. (2006). *TimeML Annotation Guidelines, Version 1.2.1*, January 2006.
- Tufiș, D., Ion, R., Ceaușu, A., Ștefănescu, D. (2006). Improved Lexical Alignment by Combining Multiple Reified Alignments. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics (EACL2006)*, Trento, Italy, 3-7 April, 2006, pp. 153-160.
- Tufiș, D., Barbu, A.M. (2002). Revealing translators knowledge: statistical methods in constructing practical translation lexicons for language and speech processing. In *International Journal of Speech Technology*. Kluwer Academic Publishers, no.5, pp.199-209, 2002, ISSN 1381-2416.

RESURSE LINGVISTICE REUTILIZABILE

CONSTANTIN CIUBOTARU, SVETLANA COJOCARU, ELENA BOIAN,
ALEXANDRU COLESNICOV, LUDMILA MALAHOVA, VALENTINA
DEMIDOV, OLEG BURLACA

Institutul de Matematică și Informatică, Academia de Științe a Republicii Moldova

{chebotar, sveta, lena, kae, mal, demidova, burlaca}@math.md

Rezumat

Lucrarea este executată în cadrul proiectului RoLTech¹ și are drept scop achiziționarea resurselor electronice pentru limba română. Este prezentată structura bazei lexicale (componenta de bază a resurselor), sunt descrise particularitățile gestionării, verificării integrității, corectitudinii și completitudinii ei. Se expun unele tehnici de verificare formală a bazei.

1. Introducere

Informatizarea continuă a societății se manifestă prin aplicarea activă a tehnologiilor informaționale. În acest context procesarea limbajului natural (PLN) devine o problemă actuală. Putem menționa trei direcții importante în PLN:

- elaborarea instrumentarului pentru PLN;
- crearea resurselor lingvistice reutilizabile;
- elaborarea aplicațiilor în baza acestor resurse.

Un efort important în implementarea produselor program pentru limba română îl prezintă pachetul de programe “Produse program pentru aplicații lingvistice” (Boian et al., 2000). Acest pachet a fost utilizat cu succes la implementarea Corectorului de texte pentru limba română RomSP (Boian et al., 2000; Cojocaru, 1997; Colesnicov, 1995). Dezvoltarea lui a condus la implementarea Resurselor Reutilizabile pentru Tehnologia Limbajului Natural (RRTLN), care conține o bază de date cu informație lingvistică la nivel de cuvânt și un set de programe de gestionare (Boian et al., 2005; Boian et al., 2005; Boian et al., 2003).

O trecere în revistă a produselor program create pentru procesarea limbajului natural este prezentată în (Cristea, Forăscu, 2006). RRTLN pot fi adăugate la clasificarea efectuată constituind o colecție de intrări lexicale completată cu informație morfologică, forme flexionate, traduceri în limba engleză și rusă, sinonime.

Corectorul RomSP poate fi considerat o aplicație lingvistică bazată pe RRTLN. Utilizând RRTLN, apare posibilitatea de elaborare a unui set de noi aplicații. De exemplu, dicționare electronice; sisteme educaționale pentru studierea morfologiei limbii române, scanere Web și motoare de căutare capabile să utilizeze formele

¹ Romanian Language Technology, proiect INTAS Ref. Nr. 05-104-7633

flexionate ale cuvintelor limbii române. RRTLN pot fi utilizate și la elaborarea aplicațiilor “clasice” pentru PLN (de exemplu, parsere, tokenizatoare, lematizatoare, etc.).

Extinderea, modernizarea și menținerea resurselor lingvistice existente, precum și a produselor program sunt efectuate în paralel cu elaborarea noilor aplicații. Începând cu anul 2006 se depun eforturi majore în direcția dezvoltării platformei RoLTech pentru tehnologia limbii române. Platforma reprezintă o colaborare dintre cercetători din Republica Moldova, România și Marea Britanie.

2. Proiectul RoLTech

Proiectul RoLTech propune următoarele obiective tehnice de bază:

1. Construirea portalului Web cu resurse lingvistice reutilizabile care vor fi folosite în tehnologia limbajului, produse program pentru tehnologia limbii române (atât surse deschise cât și cu cod autorizat) și referințe la informații utile despre limba română.
2. Elaborarea aplicațiilor bazate pe RRTLN:
 - un sistem de instruire adaptabil pentru morfologia limbii române cu elemente multimedia dedicat nevorbitorilor de limba română;
 - o aplicație dedicată vorbitorilor de limba română care are ca scop îmbogățirea vocabularului în rezultatul căutării în colecțiile de documente în limba română cu ajutorul unui motor de căutare avansat;
 - o aplicație Web ce oferă un serviciu interactiv de corectare a textelor în limba română;
 - servicii pentru utilizatori experți în limba română (de exemplu, un sistem suport pentru elaborarea dicționarilor specializate).

Resursele și aplicațiile create în cadrul proiectului vor fi plasate pe Web-portal. Inițial se vor crea versiuni prototipice, care mai apoi vor fi extinse, finalizate și menținute în continuare.

3. Flexionarea cuvintelor în limba română

Programele de flexionare (Cojocaru, 1997) au contribuit substanțial la acumularea resurselor lingvistice. Flexionarea cuvintelor este efectuată prin două metode: statică și dinamică.

Metoda statică de flexionare se bazează pe clasificarea descrisă în lucrarea (Lombard, 1981). Algoritmul utilizează o gramatică de flexionare care formalizează procesul de realizare a alternanțelor și concatenare a seturilor de terminații. Pentru limba română această gramatică include 866 reguli și 320 seturi de terminații. Această metodă a fost aplicată pentru flexionarea a circa 30000 de cuvinte-leme.

Metoda dinamică nu utilizează liste de cuvinte, dar încearcă să calculeze paradigma de flexionare utilizând clasificările asemănătoare cu cele descrise în (Lombard 1981).

Algoritmul a fost verificat pe câteva mii de cuvinte în limba română, care nu au fost incluse în acele tabele. De asemenea au fost depistate unele iregularități (3% din mulțimea de cuvinte flexionate).

4. Structura bazei de date

Resursele Reutilizabile pentru Limba Română (RRLR) conțin o bază de date (BD) cu informație lingvistică pentru limba română la nivel de cuvânt și un set de programe de gestionare a acestei baze de date. Ca volum RRLR conțin circa un milion de elemente.

În continuare vom descrie unele tehnici formale de verificare a integrității și corectitudinii RRLR ce conțin cuvinte în limba română, derivate morfologice, sinonime, traduceri în limba engleză și rusă. Descrierea mai detaliată a BD și a algoritmilor este expusă în (Cojocaru, 2006).

BD a RRLR are șase tabele de bază și 16 tabele auxiliare. Tabelele de bază sunt: *words*, *words_engl*, *words_rus*, *word_flexies*, *word_synonyms*, *word_translations*. În primele trei tabele sunt cuvinte în limbile română, engleză și rusă, cărora li se pun în corespondență niște coduri numerice. Aceste coduri numerice sunt utilizate în celelalte trei tabele. De exemplu, tabelul *word_synonyms* conține perechea de sinonime, care constă din două numere, ce corespund cuvintelor în limba română situate în tabelul *words_table*.

Tabelele auxiliare conțin diferite coduri utilizate în tabelele de bază: caracteristici morfologice, codurile limbilor, părților de vorbire, etc.

5. Popularea bazei de date

Pentru completarea BD cu informație morfologică s-a utilizat setul de fișiere produse în cadrul proiectelor precedente. Informația pentru traduceri și sinonime a fost luată din diferite surse lexicografice (Boian et al., 2003).

Fișierele existente au fost transformate într-un format unic elaborat special pentru intrările BD.

Programul de populare a BD produce adițional un fișier, care avertizează dacă cuvântul a fost deja inserat în BD, arată codul cuvântului și rezultatul fiecărei operații. Erorile sunt marcate și pot fi ușor depistate. Un alt mijloc de populare a BD cu informație morfologică este un program semi-automat care generează toate formele flexionate în baza cuvântului-lemă indicat.

Menționăm trei surse de erori care apar la popularea BD: erori preluate din surse lexicografice, erori în programele utilizate la procesarea informației și erori produse de operator (factorul uman). O parte din aceste erori pot fi depistate doar cu implicarea experților filologi. O serie de alte erori pot fi depistate în mod automat cu ajutorul unor programe special elaborate.

6. Verificarea BD

Pentru început au fost utilizate metode formale de verificare a validității structurii BD. Aceste metode au fost formulate folosind semantica și interdependențele câmpurilor BD și a tabelelor. De exemplu, câmpul *part_code* din tabelul *words* conține numere – coduri ale părților de vorbire din tabelul *parts_of_speech* – și de aceea ele pot avea numai valori întregi de la 1 (codul pentru verb) până la 10 (codul pentru conjuncție). Următoarea metodă formală a fost aplicată la verificarea cuvintelor. Pentru cuvintele în limba română a fost utilizat corectorul RomSP, care operează cu o listă de cuvinte deja testată de elaboratorii și utilizatorii acestui produs. Cuvintele limbilor română, rusă și engleză au fost testate utilizând corectoarele de texte MS Office pentru limbile corespunzătoare.

O metodă efektivă de verificare a fost utilizarea n-gramelor (părți de cuvinte ce conțin exact n litere). Cuvintele care conțin n-gramă mai puțin frecvente se consideră a fi cele mai suspicioase.

Tabelele atributelor se pot verifica vizual deoarece ele sunt scurte. O altă metodă de verificare constă în căutarea codurilor atributelor care se folosesc rar sau nici nu se folosesc în tabelele de bază.

Tabelele de bază în BD conțin referințe mutuale. În caz ideal, oricărui cuvânt în limba română ar trebui să i se atașeze forme flexionate, sinonime și traduceri. Utilizând codificările din BD se pot căuta, de exemplu, cuvintele care nu au forme flexionate, traduceri, sinonime. Putem obține o listă de cuvinte pentru care lipsește informația corespunzătoare, care ulterior ar putea fi adăugată în BD, sau o listă de cuvinte cu erori, care pot fi corectate.

Dublarea datelor s-a evitat la etapa de completare a BD. Apariția datelor dublate indica prezența unor erori în programele de completare a BD.

Verificarea statistică a procesului de flexionare ne-a permis să depistăm un șir de erori pentru cazurile când numărul formelor flexionate depășea numărul admisibil pentru o anumită parte de vorbire, de exemplu, 35–40 pentru verb. Devierea acestor numere ne indică posibile erori.

A fost efectuată verificarea cuvintelor utilizând dicționare paralele. Vom menționa, că utilizarea resurselor paralele s-a dovedit a fi o metodă utilă în PLN (Tufiș, Barbu, 2002). În cazul nostru acestea au fost traducerile în limba rusă. Limba rusă, ca și limba română, are un grad înalt de flexivitate. Au fost analizate verbe, adjective și adverbe. Aceste părți de vorbire în limba rusă au terminații tipice. Cazurile, când partea de vorbire nu corespundea celei așteptate, au fost clasificate drept suspecte și examinate suplimentar.

7. Concluzii

Proiectul RoLTech, prin natura sa interdisciplinară (combinând informatica cu lingvistica) și crearea portalului Web dedicat resurselor lingvistice, instrumentarului de procesare și referințelor la cele mai importante evenimente și descoperiri, relativ la limba vorbită în România și Republica Moldova, va ajuta la atingerea mult râvnitei coordonări a activităților ambelor categorii de cercetători: a informaticienilor și lingviștilor.

Referințe bibliografice

- Boian, E., Cojocaru, S., Malahova, L. (2000). Instruments pour applications linguistiques. *La terminologie en Roumanie et en Republique de Moldova*, Hors serie, No. 4.
- Boian, E., Ciubotaru, C., Cojocaru, S., Colesnicov, A., Demidova, V., Malahova, L. (2005). Lexical resources for Romanian. *Scientific Memoirs of the Romanian Academy*, ser.IV, vol. XXVI, București, România, pp. 267–278.
- Boian, E., Cojocaru, S., Ciubotaru, C., Colesnicov, A., Demidova, V., Malahova, L. (2005). Technologization of Romanian: linguistic resources, applications, tools. *Proceedings of the 4rd International Conference on Microelectronics and Computer Science*. Vol.II, pp.519–522.
- Boian, E., Ciubotaru, C., Cojocaru, S., Colesnicov, A., Demidova, V., Malahova, L. (2003). Lexical Resources for Romanian – a project overview. In: *Proceedings of Symposium on Intelligent Systems and Application*, September 19-20, Iasi, Romania, 12 pp. ISBN 973–97737–29.
- Cojocaru, S. (1997). Romanian Lexicon: Tools, Implementation, Usage. In: Dan Tufiș, Poul Andersen (eds.). *Recent Advances in Romanian Language Technology*. ISBN 973–27–0626–0, Editura Academiei, I, pp. 107–114.
- Cojocaru, S., Colesnicov, A., Malahova, L. (2006). Integrity and correctness checking of a lexical database. *Computer Science Journal of Moldova*, v. 14, Nr. 1(40), pp. 138–151.
- Colesnicov, A. (1995) The Romanian spelling checker ROMSP: the project overview. *Computer Science Journal of Moldova*, v. 3, Nr. 1(7), pp. 40–54.
- Cristea, D., Forăscu, C. (2006). Linguistic Resources and Technologies for Romanian Language. *Computer Science Journal of Moldova*, v. 14, Nr. 1(40), pp. 34–73.
- Lombard, A., Gadei, C. (1981). *Dictionnaire morphologique de la langue roumaine*, București (în franceză).
- Tufiș, D., Barbu, A.M. (2002). Revealing Translator's Knowledge: Statistical Methods in Constructing Practical Translation Lexicons for Language and Speech Processing. *International Journal of Speech Technology* 5, pp. 199–209.

Capitolul 3

Aplicații ale tehnologiilor lingvistice

SISTEME DE ÎNTREBARE RĂSPUNS PENTRU LIMBA ROMÂNĂ

ADRIAN IFTENE¹, IONUȚ PISTOL¹, DIANA TRANDABĂȚ^{1,2}, GEORGIANA
PUȘCAȘU³, CORINA FORĂSCU^{1,4}, DAN CRISTEA^{1,2}

¹*Facultatea de Informatică, Universitatea "Al.I.Cuza", Iași*

²*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

³*Universitatea Wolverhampton*

⁴*Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București*

{adiftene, ipistol, dtrandabat, corinfor, dcristea}@info.uaic.ro, georgie@wlv.ac.uk

Rezumat

În acest articol vom prezenta pe scurt modul în care au fost abordate problemele apărute în dezvoltarea sistemului de întrebare-răspuns pentru competiția CLEF 2006¹, secțiunea română-engleză. Pe lângă etapele clasice de adnotare sintactică și semantică a corpusurilor, au fost probleme specifice datorate necesității unei traduceri sigure din română în engleză și necesității evaluării rezultatelor folosind cât mai mult posibil metodologia CLEF. Acest articol descrie gradat pașii necesari implementării acestui sistem, pentru a putea evalua mai bine rezultatele rulărilor noastre în formatul CLEF.

1. Introducere

Primul sistem de Întrebare-Răspuns² românesc a fost dezvoltat în anii '80 (Tufiș și Cristea, 1985) și era reprezentat de o interfață ce facilita comunicarea cu o rețea semantică (care codifica cunoașterea). Astăzi sistemele de ÎR folosesc documente text ca bază de cunoaștere și integrează tehnici de prelucrare a limbajului natural (PLN) pentru a găsi (într-o colecție dată de documente sau prin căutare pe web) răspunsul la o întrebare pusă în limbaj natural.

România a participat pentru prima dată la o competiție CLEF în 2006, în cadrul secțiunii QA@CLEF³. Organizatorii au decis că limba sursă (a întrebărilor) să fie româna în timp ce limba țintă (a colecției de documente în care este căutat răspunsul) să fie engleza, datorită inexistenței unui corpus ziaristic din perioada anilor 1994-1995, care ar fi permis și folosirea românei ca limbă țintă. Astfel, la întrebările puse în română s-au căutat răspunsuri sub forma unor fragmente de text în colecția de documente în engleză.

Ca și în celelalte interacțiuni multilingve din competiție, sistemul nostru a fost evaluat pe un set de 200 de întrebări în limba română. Aflați pentru prima dată într-o astfel de competiție, intenția a fost obținerea în primul rând a unui sistem funcțional, calitatea rezultatelor fiind lăsată, în acest an, pe locul al doilea.

¹ Cross-Language Evaluation Forum: <http://www.clef-campaign.org/2006.html>

² Question Answering (QA) – rom.: Întrebare-Răspuns (ÎR)

³ Multilingual Question Answering at CLEF: <http://clef-qa.itc.it/>

2. Descrierea Sistemului

2.1. Prezentare Generală

De regulă, sistemele de ÎR folosesc o arhitectură generală de tip pipe-line, în care prelucrarea parcurge trei etape principale: analiza întrebării, căutarea documentară și extragerea răspunsului (Harabagiu, Moldovan, 2003). Sistemul creat este o variantă a arhitecturii generale, cu particularizări specifice legate de reprezentare și procesare pentru fiecare din componentele amintite mai sus. Un modul aparte inclus în sistemul de ÎR este modulul care traduce cuvintele din română în engleză, pentru a face transferul interlingv. Arhitectura și funcționalitatea sistemului sunt ilustrate în Figura 1.

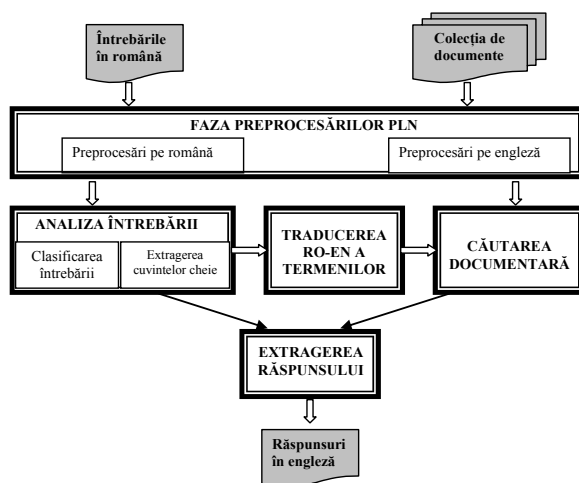


Figura 1: Arhitectura și funcționalitatea sistemului

2.2. Preprocesări asupra corpusului

Inițial, întrebările sunt procesate morfo-sintactic folosind POS tagger-ul românesc dezvoltat de ICIA⁴ (Tufiș, 1999; Ion, 2006). Ulterior, se realizează clasificarea numelor proprii, printr-o procedură de recunoaștere bazată pe șabloane, în următoarele clase: PERSOANĂ, LOCAȚIE, MĂSURĂ și altele. Aceleași operații de preprocesare sunt efectuate și pe colecția de documente englezești, folosind o segmentare la nivel de propoziție. Pentru aceasta s-a folosit același parser de la ICIA, dar cu un alt model de limbă. Corpusul CLEF englezesc constând din articole din ziarele Los Angeles Times anul 1994 și Glasgow Herald anul 1995. Acest corpus a fost segmentat la nivel de propoziție și cuvânt, iar apoi a fost etichetat la partea de vorbire.

2.3. Analiza întrebării

Această etapă are în vedere, în primul rând, identificarea tipului semantic al entității ce ar trebui să sugereze tipul răspunsului așteptat. În plus se identifică focusul întrebării, tipul întrebării și mulțimea cuvintelor cheie relevante pentru întrebare. Pentru a putea atinge aceste scopuri, analiza efectuează următorii pași:

⁴ Institutul de Cercetări în Inteligență Artificială al Academiei Române: <http://www.racai.ro/>

a) Depistarea grupurilor nominale (GN), extragerea numelor de entități (NE), identificarea expresiilor temporale (ET)

Identificatorul de nume de entități construit de ICIA determină numele de entități din întrebarea în românește. Expresiile temporale sunt de asemenea identificate folosind un identificator și un normalizator de ET pentru limba română adaptat după varianta pe limba engleză descrisă în (Pușcașu, 2004).

b) Identificarea focusului întrebării

Focusul întrebării este cuvântul sau secvența de cuvinte care arată ce anume se caută. Se consideră ca focus al întrebării substantivul ulterior pronumelui interogativ din întrebare (ca în *Ce țară*) sau primul grup nominal (GN) al întrebării dacă el apare înainte de verbul principal al întrebării sau de cel ce urmează verbului *a fi*.

c) Găsirea tipului răspunsului așteptat

Sistemul de analiză a întrebării poate face distincție între următoarele clase de răspuns: PERSOANĂ, LOCAȚIE, ORGANIZAȚIE, TEMPORAL, NUMERIC, DEFINIȚIE și GENERIC. Atribuirea unei clase unei întrebări analizate este realizată folosind focusul întrebării și tipul acestuia. De exemplu, în cazul întrebării *În ce oraș a fost omorât Vladislav Listyev?*, focusul întrebării este *oraș*, substantiv ce apare în lista LOCAȚIILOR, și astfel se determină tipul răspunsului ca fiind LOCAȚIE. Tipul focusului întrebării este determinat folosind WordNet (Fellbaum, 1998).

d) Deducerea tipului întrebării

În acest an, competiția QA@CLEF a făcut distincție între patru tipuri de întrebare: *factoid*, *definiție*, *listă* și întrebări cu *restricții temporale*. Deoarece restricțiile temporale se pot atașa oricărui tip de întrebare, alegerea unuia dintre tipurile *factoid*, *definiție* sau *listă* se face înainte de testarea existenței restricțiilor temporale.

e) Generarea mulțimii cuvintelor cheie

Mulțimea cuvintelor cheie este generată automat din lista termenilor importanți ai întrebării în ordinea inversă a relevanței acestora. Prin urmare, mulțimea cuvintelor cheie cuprinde: focusul întrebării, NE și ET identificate, substantivele rămase, și toate verbele diferite de cele auxiliare prezente în întrebare. Această mulțime este apoi trimisă modulului de traducere a termenilor, cu scopul de a obține cuvintele cheie englezești necesare căutării documentare.

2.4. Traducerea termenilor

Pentru a realiza traducerea termenilor s-a folosit ca resursă WordNet-ul, disponibil atât în limba română (Tufiș et al., 2006) cât și în limba engleză. Mulțimea cuvintelor cheie extrase în faza de analiză a întrebării este folosită ca intrare pentru faza de traducere a termenilor, și prin urmare s-au tradus atât substantivele cât și verbele. Cuvintele componente ale GN sunt traduse unul câte unul: synseturile românești care conțin cuvântul în cauză sunt puse în corespondență, prin indexul inter-lingual ILI, cu cele echivalente din engleză, obținându-se astfel mulțimea tuturor traducerilor posibile. Dacă cuvântul ce trebuie tradus nu apare în WordNet, caz destul de frecvent, acesta se caută în alte dicționare disponibile și, dacă este găsit, se păstrează primele trei traduceri. În cazul

verbelor, se extrage pentru fiecare verb traducerea echivalentă din WordNet la fel ca la substantive. Pentru cazurile de verb-substantival am folosit situațiile din (Pekar et al, 2004) în selectarea traducerii.

2.5. Crearea indexului și căutarea documentară

Corpusul englezesc a fost preprocesat inițial folosind instrumente pentru împărțirea în cuvinte, instrumente pentru găsirea lemei și a părții de vorbire, și unelte de recunoaștere a NE. În rulările noastre indexarea și căutarea s-a realizat cu motorul Lucene⁵.

Colecția de documente a fost indexată atât la nivel de document cât și la nivel de paragraf folosind lema cuvintelor conținute și a clasele NE (MĂSURĂ, PERSOANĂ, LOCAȚIE, etc). Când nu este găsit nici un paragraf pentru o anumită interogare, se folosesc două strategii: fie se măresc segmentele de la paragrafe la documente, fie se reformulează interogarea folosind pentru anumite cuvinte alte variante de traducere.

2.6. Extragerea răspunsului

Două module de extragere a răspunsului au fost dezvoltate, unul de către UAIC⁶ și altul de către ICIA. Ambele module au ca intrare tipul răspunsului așteptat, focusul întrebării, mulțimea de cuvinte cheie, părțile de text obținute în urma căutării pe partea de vorbire, lema și informații de tip NE și indicatorul de relevanță al paragrafelor determinat de Lucene. Procesul de extragere depinde de tipul așteptat al răspunsului: când răspunsul are ca tip un NE, modulul de extragere a răspunsului identifică în fiecare propoziție întoarsă de Lucene entitățile de tip NE care au tipul dorit de răspuns. Când tipul răspunsului nu este un NE, procesul de extragere se bazează în principal pe recunoașterea focusului, în acest caz șabloanele sintactice de găsire a răspunsului bazate pe focus fiind cruciale.

3. Descrierea rulărilor înscrise în competiție

Au fost înscrise în competiție trei rulări diferite, cu următoarele detalii:

- **UAIC** - Această rulare a fost obținută prin parsarea și analizarea întrebărilor, traducerea cuvintelor cheie, căutarea pasajelor relevante și căutarea răspunsurilor finale folosind extractorul de răspunsuri realizat de UAIC.
- **RACAI**⁷ - Această rulare a fost obținută de asemenea prin parsarea și analizarea întrebărilor, traducerea cuvintelor cheie, căutarea pasajelor relevante, dar pentru căutarea răspunsurilor finale s-a folosit extractorul de răspunsuri ICIA.
- **DIOGENE** - Cea de a treia rulare a fost obținută prin conversia rezultatelor modulelor de analiză a întrebării și traducere a termenilor în formatul cerut de sistemul de ÎR DIOGENE (Kouylekov et al, 2003), și apoi trimiterea lor ca intrare la modulele DIOGENE de căutare documentară și de extragere de răspunsuri.

⁵ <http://lucene.apache.org/>

⁶ Universitatea "Al.I.Cuza": <http://www.uaic.ro>

⁷ Sigla englezească pentru ICIA.

Datorită numărului mare de rulări înscrise având ca țintă limba engleză doar rulările UAIC și RACAI au fost evaluate. În continuare, sistemul RACAI va fi referit ca Sistemul 1, și sistemul UAIC va fi referit ca Sistemul 2.

4. Analiza Rezultatelor

Rezultatele evaluării oficiale pentru Sistemele 1 și 2 sunt prezentate în Tabelul 1. Fiecare răspuns a fost evaluat ca fiind NECUNOSCUȚ (răspunsurile neevaluate), CORECT (răspunsurile corecte), NEJUSTIFICAT (răspunsuri care nu puteau fi găsite în bucățile de text justificatoare), INCORECT (răspunsurile greșite) sau INEXACT (răspunsuri incomplete).

Tabel 1: Evaluarea rezultatelor pentru cele două sisteme

Evaluarea rezultatelor pentru Sistemul 1			Evaluarea rezultatelor pentru Sistemul 2		
Z	NECUNOSCUȚ	400	Z	NECUNOSCUȚ	543
R	CORECT	35	R	CORECT	22
U	NEJUSTIFICAT	13	U	NEJUSTIFICAT	4
W	INCORECT	184	W	INCORECT	191
X	INEXACT	7	X	INEXACT	1
	TOTAL	639		TOTAL	761

Numărul mare de răspunsuri evaluate ca fiind NECUNOSCUȚ se datorează faptului că s-au determinat 10 răspunsuri pentru aproape toate cele 200 de întrebări, unde 10 a fost numărul maxim de răspunsuri posibile. Cum evaluarea finală a ținut cont doar de primul răspuns pentru majoritatea întrebărilor (doar în cazul întrebărilor de tip *listă* au fost evaluate primele trei răspunsuri), răspunsurile de pe pozițiile de la 2 la 10 au fost etichetate ca NECUNOSCUȚ, indicând faptul că nu s-a făcut nici o încercare pentru a le verifica corectitudinea. Folosind un evaluator dezvoltat special pentru aceasta răspunsul corect a fost găsit în primele zece răspunsuri generate de sistemele noastre pentru 35-40% din întrebări. Acesta este un rezultat promițător, dovedind că extractorul de răspunsuri funcționează, dar enecesară îmbunătățirea ordonării răspunsurilor obținute.

5. Concluzii

Respectând arhitectura clasică a sistemelor de tip ÎR, sistemul dezvoltat implementează cele trei niveluri esențiale ale unui astfel de sistem, ca și modulul specific sistemelor interlinguale care traduce termenii relevanți ai întrebării din română în engleză.

Rezultatele și evaluările, deși nu tocmai satisfăcătoare, vor fi folosite pentru îmbunătățirea sistemului la ediția viitoare de QA@CLEF sau la alte competiții similare precum TREC⁸.

O analiză detaliată a relevat un număr important de direcții pentru îmbunătățirea substanțială a sistemului. Modulul de traducere a termenilor, cheia performanțelor pentru orice sistem interlingv, este principalul obiectiv. Modulul de extragere a

⁸ Text REtrieval Conference: <http://trec.nist.gov/>

răspunsului va fi modificat astfel încât precizia acestuia să crească. O metodă mai bună de ordonare a răspunsurilor candidate este a treia direcție prioritară pentru viitor.

Autorii sistemului

- ◇ **Membrii echipei UAIC:** Dan Cristea, Iustin Dornescu, Corina Forăscu, Maria Husarciuc, Adrian Iftene, Ana Masalagiu, Alex Moruz, Ionuț Pistol, Diana Trandabăț;
- ◇ **Membrii echipei ICIA:** Alin Ceaușu, Radu Ion, Dan Ștefănescu, Dan Tufiș;
- ◇ **Universitatea Wolverhampton:** Georgiana Pușcașu, Constantin Orăsan.

Mulțumim lui Milen Kouylekov și Bernardo Magnini pentru disponibilitatea lor de a procesa ieșirea procesorului nostru de întrebări și a modulului de traducere cu sistemul DIOGENE, dezvoltat de IRST Trento.

Parțial, acest proiect a fost finanțat de Ministerul Educației și Cercetării în cadrul proiectului CEEEX 29 ROTEL și de INTAS în cadrul proiectului RolTech (INTAS ref. 05-104-7633).

Referințe bibliografice

- Fellbaum, C. (1998) (ed.) WordNet: An Electronic Lexical Database. The MIT Press.
- Harabagiu, S., Moldovan, D. (2003). Question answering. In R. Mitkov (ed.), *Oxford Handbook of Computational Linguistics*, pp. 560 - 582. Oxford University Press.
- Kouylekov, M., Magnini, B., Negri, M., Tanev, H. (2003). ITC-irst at TREC-2003: the DIOGENE QA system. *Proceedings of the TREC-12 Conference*.
- Ion, R. (2006). Metode de dezambiguizare semantică automată. Aplicații pentru limbile engleză și română. Teză de doctorat în curs de susținere la Academia Română.
- Pekar, V., Krkoska, M., Staab, S. (2004). Feature weighting for cooccurrence-based classification of words. *Proceedings of the 20th International Conference on Computational Linguistics (COLING-04)*.
- Pușcașu, G. (2004). A Framework for Temporal Resolution. *Proceedings of the 4th Conference on Language Resources and Evaluation (LREC2004)*.
- Tufiș, D. (1999). Tagging with Combined Language Models and Large Tagsets. *Proceedings of the TELRI International Seminar on "Text Corpora and Multilingual Lexicography"*.
- Tufiș, D., Barbu Mititelu, V., Ceaușu, A., Bozianu, L., Mihăilă, C., Manu Magda, M., (2006). Noi dezvoltări ale wordnet-ului românesc. În acest volum.
- Tufiș, D., Cristea, D., Stamou, S. (2004). BalkaNet: Aims, Methods, Results and Perspectives. A General Overview. In D. Tufiș (ed.), *Romanian Journal on Information Science and Technology*. Special Issue on BalkaNet. Romanian Academy.

IDENTIFICAREA ȘI EXTRAGEREA AUTOMATĂ A COLOCAȚIILOR DIN TEXTE

DAN ȘTEFĂNESCU, DAN TUFIS, ELENA IRIMIA

Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București

{danstef, tufis, elena}@racai.ro

Rezumat

Identificarea și extragerea automată a cologațiilor din texte este necesară în rezolvarea multor probleme dificile de prelucrare a limbajului natural, cum ar fi generarea limbajului natural, rezumarea sau traducerea automată. Cologațiile sunt expresii care de obicei nu pot fi traduse cuvânt cu cuvânt (folosind doar un simplu dicționar și un model de limbă). Acest lucru se întâmplă deoarece sunt caracterizate de compoziționalitate limitată – înțelesul expresiei nu se obține întotdeauna însumând înțelesurile cuvintelor ce o compun.

1. Introducere

Diverse definiții au fost propuse pentru noțiunea de *cologație*, mai mult sau mai puțin stricte. Iată câteva dintre cele care sunt folosite în lingvistica computațională:

- expresie formată din două sau mai multe cuvinte ce corespunde unui mod convențional de a afirma, de a exprima, anumite lucruri;
- „Două sau mai multe cuvinte ce apar împreună semnificativ de des în interiorul unei ferestre pre-definite într-un corpus dat” (Quasthoff & Wolff, 2002);
- “O secvență de două sau mai multe cuvinte consecutive, ce are caracteristicile unei unități sintactice și semantice, și a cărei înțeles exact și neambiguu nu poate fi obținut direct din înțelesurile sau conotațiile cuvintelor ce o compun” (Choueka, 1988).

Cologațiile pot fi grupuri nominale (*televizor alb-negru, arme de distrugere în masă, vin roșu, drept de suită*), locuțiuni verbale (*a aduce atingere, a intra în vigoare, a face obiectul, a lua în considerare*) și nu numai (*sărac dar cinstit, tânăr și neliniștit, de jur împrejur*).

Cologațiile se pot caracteriza prin (Manning & Schütze, 1999):

- **Non-compoziționalitate** – atunci când înțelesul întregului este diferit de suma înțelesurilor părților;
- **Non-substituționalitate** – atunci când componentele cologației nu pot fi substituite cu sinonime;
- **Non-modifiabilitate** – atunci când cologațiile nu pot fi modificate prin adăugarea de material lexical adițional sau prin transformări gramaticale.

O cologație pentru care avem îndeplinite toate cele trei condiții de mai sus se apropie foarte mult de noțiunea de idiom. Cologațiile pot fi clasificate după mai multe criterii, ele pot fi de natură lexicală, sintactică, sau de natura semantică, pot fi generale sau specifice unui anumit domeniu, pot avea structură fixă sau structură variabilă.

Literatura de specialitate propune diferite metode pentru găsirea cologațiilor. Justeson și Katz (1995) au folosit doar frecvența de ocurență a cuvintelor în perechi și un filtru pe părțile de vorbire; Smadja (1990) a folosit o metodă bazată pe media și dispersia distanțelor dintre (două) cuvinte în corpus, în timp ce alții (Church et al., 1991) au utilizat Testul *t*, *chi* pătrat, *log-likelihood* sau informația mutuală pentru a găsi cuvinte ce apar împreună, în text, mai des decât ne-am putea aștepta să apară întâmplător.

2. Modelarea cologațiilor

În modelarea noastră, cologațiile sunt succesiuni de cuvinte (nu neapărat adiacente) care respectă două criterii statistice:

- distanța dintre cuvinte este relativ constantă;
- apar în aceleași contexte de un număr de ori semnificativ din punct de vedere statistic.

Primul criteriu este evaluat folosind abordarea lui Smadja (1990) iar cel de-al doilea se bazează pe calculul raportului Log-Likelihood (LL).

Pentru identificarea cologațiilor Smadja propune utilizarea mediei și dispersiei distanțelor dintre cuvintele din corpus (2 câte 2). Dispersia măsoară deviația distanțelor de la medie:

$$\sigma^2 = \frac{\sum_{i=1}^n (d_i - \mu)^2}{n - 1}$$

Dacă două cuvinte se găsesc în corpus mereu la aceeași distanță, dispersia este 0. Dacă distanțele au o distribuție aleatoare (cum este cazul atunci când cuvintele apar împreună întâmplător), dispersia are valori ridicate. Putem afirma astfel că media și deviația standard (rădăcina pătrată a dispersiei) sunt mărimi care caracterizează distribuția distanțelor dintre două cuvinte într-un corpus. Smadja demonstrează că se pot descoperi cologații căutând perechi de cuvinte pentru care avem deviații standard mici.

Scorul Log-Likelihood calculează raportul probabilităților a două ipoteze statistice care pot fi emise în descrierea datelor observate într-un text. Ipotezele pe care le luăm în considerare sunt (i) H_0 : cele două cuvinte nu au nici o legătură între ele și apar întâmplător împreună și (ii) H_1 : cele două cuvinte sunt cumva corelate și apariția lor împreună nu este întâmplătoare.

- $H_0 : P(w_2|w_1) = p = P(w_2|\neg w_1)$
(presupunere de independență)
- $H_1 : P(w_2|w_1) = p_1 \neq p_2 = P(w_2|\neg w_1)$
(w_1 și w_2 nu sunt independente)

IDENTIFICAREA ȘI EXTRAGEREA AUTOMATĂ A COLOCAȚIILOR DIN TEXTE

Pentru calculele efective se folosește un tabel de contingență ca mai jos, în care fiecare celulă conține numărul de apariții ale diferitelor combinații de cuvinte pentru care se evaluează scorul LL. Astfel, n_{11} reprezintă numărul de apariții împreună ale cuvintelor w_1 și w_2 , n_{12} reprezintă numărul de apariții ale cuvântului w_1 în contextele în care cuvântul w_2 lipsește, etc.

	w_2	$\neg w_2$
w_1	n_{11}	n_{12}
$\neg w_1$	n_{21}	n_{22}

Notând cu: $n_{1*} = n_{12} + n_{11}$, $n_{*1} = n_{21} + n_{11}$ și cu $n_{**} = n_{22} + n_{11}$, formula de calcul este:

$$LL = 2 * \sum_{j=1}^2 \sum_{i=1}^2 n_{ij} * \log \frac{n_{ij} * n_{**}}{n_{i*} * n_{*j}}$$

Dacă scorul obținut este mai mare decât un anumit prag ipoteza de nul (H_0) este respinsă cu un anumit grad de certitudine.

Rezultatele obținute folosind o combinație a celor două metode descrise mai sus indică un lucru interesant: utilizarea scorului LL calculat pentru perechi de cuvinte care îndeplinesc anumite criterii ce țin de partea de vorbire, cât și de media distanței dintre cuvinte, constituie o abordare eficientă. În ceea ce ne privește, suntem interesați de extragerea de colocații ce sunt de tip substantiv-substantiv (S-S), substantiv-adjectiv (S-A) / adjectiv-substantiv (A-S) și substantiv-verb (S-V)/ verb-substantiv (V-S). Trebuie să remarcăm că, în timp ce colocațiile din prima și a doua categorie se încadrează în general în categoria compușilor terminologici (termenii cheie se găsesc aici în marea lor majoritate), cele din a treia categorie caracterizează structurile de sub-categorizare verbale.

Următoarele rânduri descriu metoda folosită de noi pentru extragerea colocațiilor de tip substantiv-verb (verb-substantiv). Inițial, textul este lematizat și adnotat la părți de vorbire. Apoi, o fereastră de 11 cuvinte (acesta este contextul în care se consideră co-ocurențele) parcurge fiecare propoziție din text în așa fel încât fiecare cuvânt devine la un moment dat centrul ferestrei¹. Cuvintele ce se introduc în fereastră sunt substantive sau verbe; celelalte părți de vorbire sunt ignorate. Lungimea a fost aleasă astfel încât fereastra să poată cuprinde orice pereche de cuvinte interesantă care ar exista. Am considerat că o distanță de 5 (stânga/dreapta) pentru o astfel de fereastră, în care se găsesc doar cuvinte ce sunt verbe sau substantive, (pentru alte tipuri de colocații considerăm doar substantivele și adjectivele, sau doar substantivele) este suficientă pentru a găsi perechile interesante. Deși ar putea exista exemple în care distanța dintre cuvinte este mai mare de 5 (numărând doar cuvintele din categoriile gramaticale care ne interesează pe noi), aceste cazuri sunt rare și se datorează probabil intercalării unor expresii între cuvintele ce formează perechea interesantă². Toate perechile de cuvinte (sub formă de lemă) ce se formează între centrul ferestrei și celelalte cuvinte din fereastră, împreună cu distanța³ dintre cuvintele ce formează aceste perechi, sunt

¹ În acest pas aplicăm metoda lui Smađja. Aceasta ne permite să identificăm perechi interesante de cuvinte ce nu sunt neapărat adiacente.

² Este evident că această tehnică funcționează pentru limbi ca româna, engleza, franceza și multe altele; pentru limba germană, unde verbul stă uneori la sfârșitul propoziției, căutarea perechilor interesante verb-substantiv în funcție de distanța dintre cuvinte trebuie regândită, iar contextul de căutare trebuie extins la nivelul întregii propoziții.

³ Distanța este negativă dacă perechea e formată de cuvântul-centru împreună cu un cuvânt ce stă în fața sa.

introduse într-o bază de date. După ce a fost parcurs întreg textul, pentru fiecare pereche din baza de date, calculăm media și dispersia luând în calcul ocurențele la diferite distanțe. Dispersia reprezintă variația distanței dintre cele două cuvinte în jurul mediei. În cazul în care avem o dispersie mică, media ne indică distanța uzuală la care se află cele două cuvinte în text. Dispersia este pătratul deviației standard. În cazul nostru, am păstrat ca perechi interesante pe cele pentru care deviația standard este sub un prag de 1,5. Valoarea de 1,5 este îndeajuns de mare încât să prindem toate perechile interesante⁴. Pentru perechile interesante extrase, calculăm scorul LL. În acest calcul intră doar perechile de cuvinte de tip S-V / V-S care se află la o distanță egală cu media perechii pentru care se calculează acest scor. Dacă scorul LL depășește pragul de 9, spunem că perechea este o cologație. Pentru un scor de 9, probabilitatea de eroare este mai mică de 0,004⁵. Dacă dorim îndeplinirea condițiilor, putem micșora pragul pentru dispersie și / sau ridica pragul pentru scorul LL.

Trebuie să remarcăm că, dacă din lista cologațiilor obținute, luăm primele x în ordinea tăriei cu care sunt legați constituenții din cologații⁶, tărie care este dată de scorul LL, obținem o listă de termeni cheie relevanți pentru documentul din care i-am extras.

3. *Evaluarea metodei*

Am testat această metodă în contextul ambelor proiecte de care am amintit chiar la începutul acestui articol. Am folosit un corpus românesc, lematizat și adnotat la părți de vorbire, cu o mărime de aproximativ 350Mb ce conține articole ce fac parte din Acquis-ul Comunitar.

În cazul extragerii de termeni cheie, trebuie să facem precizarea că ne așteptăm ca ei să fie o submulțime a cologațiilor de tip S-S și S-A. Pentru o evaluare adecvată am fi avut nevoie de documente în care termenii cheie să fi fost etichetați. Cum nu am avut posibilitatea să folosim astfel de documente, ne-am orientat spre utilizarea tezaurului Eurovoc, un tezaur poli-tematic, multilingv, folosit pentru indexarea Acquis-ului Comunitar (AC) (legislația și tratatele interne ale comunității europene). Tezaurul conține 6645 de termeni organizați în structuri arborescente. Dintre aceștia, 519 termeni sunt foarte generali, constituind rădăcinile arborilor. Din punctul nostru de vedere, termenii Eurovoc-ului sunt doar o submulțime a întregii mulțimi a termenilor cheie ce pot caracteriza documentele AC. Datorită faptului că ierarhia tezaurului nu este una adâncă, există un anumit grad sau nivel de generalizare la care se opresc termenii din Eurovoc. Cu alte cuvinte, se pot găsi termeni cheie pentru anumite documente din AC, foarte specifici, care însă nu se regăsesc în Eurovoc. Un exemplu în acest sens este *tratat de instituire* care este un anumit tip de tratat, dar care nu face parte din tezaur.

tratat instituire 2(distanța) 71286.61852(scorul LL) 6175(ocurențe în text)

⁴ A se vedea exemplele din Manning & Schütze – într-o pereche cu o deviație standard de peste 2, cuvintele nu au nici o legătură; ajung împreună din întâmplare.

⁵ Pragul de 0,004 înseamnă că aproximativ odată la 250 de cazuri avem o situație în care deși două cuvinte apar împreună întâmplător, scorul LL este 9. Pentru un prag de 0,001 scorul LL trebuie ales 10,83 (a se consulta tabelele cu pragurile date pentru distribuțiile *chi* pătrat cu un grad de libertate).

⁶ În funcție de cât de caracteristici vrem să fie termenii cheie, x poate fi mai mare (chiar câteva sute) sau mai mic (până la câteva zeci).

IDENTIFICAREA ȘI EXTRAGEREA AUTOMATĂ A COLOCAȚIILOR DIN TEXTE

Perechea **tratat – instituire** așa cum apare în text:

tratatul/tratat/nsry de/de/s instituire/instituire/nsm 6052
 tratatului/tratat/nsoy de/de/s instituire/instituire/nsm 70
 tratatele/tratat/npy de/de/s instituire/instituire/nsm 29
 tratatelor/tratat/npoy de/de/s instituire/instituire/nsm 13
 tratatul/tratat/nsry prevede/prevedea/v3 instituirea/instituire/nsry 3
 tratat/tratat/nsn prevede/prevedea/v3 instituirea/instituire/nsry 2
 tratate/tratat/npn de/de/s instituire/instituire/nsm 2
 tratatului/tratat/nsoy este/fi/v3 instituirea/instituire/nsry 1
 tratat/tratat/nsn figurează/figura/v3 instituirea/instituire/nsry 1
 tratat/tratat/nsn menționează/menționa/v3 instituirea/instituire/nsry 1
 tratatul/tratat/nsry privind/privi/vg instituire/instituire/nsm 1

Se observă că *tratat de instituire* apare foarte des în documentele AC și, în plus, scorul LL este foarte ridicat. Aceste lucruri susțin ideea că *tratat de instituire* poate fi considerat cuvânt cheie chiar dacă el nu apare în tezaurul oficial.

Iată alte colocații, de tip S-S, considerate de noi termeni cheie pentru documentele AC. Acești termeni se regăsesc și în Eurovoc:

stat membru	1	100653.529	23854	program lucru	2	2893.112575	439
directivă consiliu	1	42393.08842	6556	tratat aderare	2	690.3592391	153

Iată exemple de colocații de tip S-V extrase folosind aceeași metodă:

aduce atingere	1	51567.34864	4959	ține seamă	1	22825.70709	2357
înlocui text	3	43992.3067	4114	adopta bruxelles	2	21792.22915	2610
intra vigoare	2	42527.03736	4473	adopta regulament	2	20847.73793	2951
avea tratat	3	32050.11219	5816	lua măsuri	1	19207.12849	2491
face obiect	1	30729.47663	3898	pune aplicare	2	13186.20796	1564
modifica regulament	4	29141.39454	3098	face referire	1	11854.14299	1486
modifica dată	2	27658.4116	3213	informa privire	4	10586.88849	1175
lua considerare	2	27062.0349	2621	îndeplini condiție	1	8382.436218	1334
ține cont	1	26635.12649	2868	intra incidență	2	8119.841768	985
adresa membru	2	25844.0428	2362	îndeplini cerință	1	7473.851703	1223

O parte din acestea, cum ar fi *a intra în vigoare*, *a lua în considerare*, *a fi adoptat la bruxelles* pot fi considerați termeni specifici limbajului juridic reflectat în corpusul AC.

4. Concluzii și continuarea cercetărilor

Au fost prezentate o parte din rezultatele obținute, până în prezent, în cadrul unui proiect internațional la care participă: ICIA (RO), Universitatea March Bloch (FR) și IMS Stuttgart (GE), care are ca obiectiv construirea unui dicționar de colocații pentru cele 3 limbi. Totodată, această tehnică de extragere a colocațiilor a fost folosită cu succes și în cadrul proiectului CEEX-ROTEL, pentru extragerea automată a termenilor cheie multi-cuvânt din documente arbitrare.

Testele efectuate până acum arată că metoda noastră este una eficientă. În plus, ea nu depinde de limba naturală în care sunt redactate documentele prelucrate. În continuare, sunt prevăzute o serie de experimente și evaluări noi:

- vor fi extrase și colocații de tip S-A evaluarea extragerii termenilor cheie se va face în raport cu tezaurul Eurovoc;
- se va testa o metodă de extragere a colocațiilor similară cu cea prezentată, dar în care cele două faze principale sunt interschimbate;
- se vor extrage colocațiile din textele de limbă franceză și germană;
- se vor extrage diverse statistici din colocațiile deja extrase;
- se vor compara cross-lingual colocațiile extrase din documentele AC, folosind tehnologia alinierii la nivel de cuvânt (Tufiș et al., 2006) a corpusului paralel.

Referințe bibliografice

- Choueka, Y. (1988), Looking for needles in a haystack, *Proceedings of RIAO '88*, 609 – 623.
- Church, K., Gale, W., Hanks, P., Hindle, D. (1991). Parsing, word associations and typical predicate-argument relations, *Current Issues in Parsing Technology*. Kluwer Academic, Dordrecht, Olanda.
- Justeson, J. S., Katz, S. M. (1995). Technical terminology: Some linguistic properties and an algorithm for identification in text, *Natural Language Engineering*, 1:9-27.
- Manning, C., Schütze, H. (1999), *Foundations of Statistical Natural Language Processing*, MIT Press, Cambridge.
- Quasthoff, U., Wolff, C. (2002). The Poisson collocation measure and its application, *Workshop on Computational Approaches to Collocations*, Viena, Austria
- Salton, G., McGill, M. J. (1983). *Introduction to modern information retrieval*, McGraw-Hill.
- Salton, G., Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management* 24(5), 513–523.
- Smadja, F. A., McKeown, K. R. (1990). Automatically extracting and representing collocations for language generation. *Proceedings of the 28th annual meeting on Association for Computational Linguistics*, 252-259, Pittsburgh, Pennsylvania.
- Steinberger, R., Pouliquen, B., Widiger, A., Ignat, C. Erjavec, T., Tufiș, D., Varga, D. (2006). The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. *Proceedings of the 5th LREC Conference*, pp.2142-2147.
- Tufiș, D., Ion, R., Ceaușu, A., Ștefănescu, D. (2006). Improved Lexical Alignment by Combining Multiple Reified Alignments. *Proceedings of the 11th EACL Conference*, pp. 153-160.

SPRE O EXTRAGERE AUTOMATĂ A COLOCAȚIILOR: CAZUL VERBULUI 'A FACE'

AMALIA TODIRAȘCU

¹LILPA, Université Marc Bloch, Strasbourg

todiras@umb.u-strasbg.fr

Rezumat

Articolul prezintă rezultatele unui studiu al proprietăților morfologice, sintactice și semantice ale locuțiunilor verbale, concentrându-se asupra celor care folosesc verbul 'a face'. Studiul este realizat pe baza unor corpusuri specializate (texte juridice) dar și al unor corpusuri generale (ziare, texte literare, manuale de utilizare).

1. Introducere

Colocațiile sunt expresii des utilizate, care au un sens diferit de cuvintele care o compun, reprezentând un element dificil în învățarea și folosirea unei limbi străine, cât și în cadrul unor sisteme de procesare automată a limbajului natural. Articolul de față își propune să prezinte rezultatul studiului proprietăților morfo-sintactice pentru o clasă particulară de cologații, cele care implică folosirea verbelor speciale, care intră în componența multor locuțiuni verbale (a face, a lua etc.), în vederea creării unor resurse lingvistice electronice complete pentru procesare automată. Acest studiu este realizat în cadrul unui proiect de cercetare internațional, implicând ca parteneri Institutul de Cercetări în Inteligența Artificială al Academiei Române (România), IMS Stuttgart (Germania), INSA Strasbourg și Universitatea Marc Bloch Strasbourg (Franța). Proiectul "Collocations en contexte: étude et analyse contrastive" (finanțat de către Agenția Universitară pentru Francofonie) are ca obiectiv realizarea unui dicționar de cologații multilingv (română, germană și franceză), precum și al unui sistem de extragere automată a cologațiilor, pe baza informațiilor contextuale (morfologice, sintactice). Proiectul își propune să identifice principalele proprietăți morfo-sintactice ale celor mai frecvente cologații, precum și ai constituenților sintactici care se combină cu acestea, în vederea creării unui dicționar electronic multilingv. Metodologia aleasă este deja aplicată pentru limba germană (Heid & Ritz, 2005), fiind bazată pe interpretarea informațiilor contextuale asociate cologațiilor și va fi aplicată pentru limba română și franceză.

În cadrul acestui articol, ne concentrăm asupra studiului unei clase specifice de locuțiuni verbale, care implică verbul a face (*a face obiectul, a face apel la, a face referire la* etc.).

2. Metodologia propusă

O astfel de resursă lingvistică care descrie comportamentul cologațiilor este absolut necesară pentru un sistem de traducere automată sau pentru o platformă e-learning pentru învățarea limbilor străine. Folosirea corectă a cologațiilor pune probleme persoanelor care învață și utilizează o limbă străină, datorită faptului că acestea au un sens care nu poate fi dedus întotdeauna pe baza sensului părților componente. De asemenea, acest tip de expresii pune

probleme deosebite unui sistem de traducere automată, o traducere cuvînt cu cuvînt nu este întotdeauna adaptată. Din aceste motive, mai multe studii s-au orientat spre o identificare automată a acestora. Astfel, metodele statistice fac ipoteza de lucru următoare, cologațiile sunt expresii care se repetă în mod frecvent, deci criteriul frecvenței este important pentru a detecta în mod automat cologațiile într-un text (Smadja, 1993). O altă categorie de metode pleacă de la principiul ca o cologație este caracterizată de proprietăți morfo-sintactice specifice (Hausmann, 2004). Astfel, între elementul de bază al cologației și constituentul asociat se stabilesc relații de dependență sintactică (substantiv modificat de un adjectiv, verb și complementul său, etc.), chiar dacă expresiile sunt discontinue. O serie de metode de extragere a cologațiilor se concentrează asupra relațiilor sintactice care au loc între bază și constituentul asociat (Seretan et al., 2004), dar nici una din metodele acestea nu reușește identificarea corectă a cologațiilor.

Ipoteza de lucru pe care am adoptat-o în acest proiect este aceea că vom combina metodele statistice cu cele bazate pe informație sintactică pentru a obține o precizie mai mare în cazul extragerii automate a cologațiilor. În afară de relațiile sintactice care există între diversele elemente ale unei cologații, facem ipoteza că putem stabili o serie de proprietăți morfologice și sintactice care permit identificarea cologației cu precizie. Astfel, anumite cologații preferă folosirea articolului definit (*face obiectul* dar nu **face un obiect*) sau al pluralului, se pot combina doar cu anumite clase de prepoziții, sau acceptă unele adverbe între verb și complementul său direct. Pentru a identifica proprietățile cele mai interesante, am studiat un corpus paralel multilingv, pentru identificarea proprietăților pertinente care permit identificarea unor clase de cologații specifice. Proprietățile identificate pe baza acestui studiu vor fi verificate pe baza unor corpusuri disponibile în fiecare din limbile studiate. Aceste proprietăți vor fi selecționate pentru a fi reprezentate în dicționarul de cologații care va fi construit. Dicționarul va fi integrat în cadrul unui sistem de extragere automată a cologațiilor. În continuare ne concentrăm asupra tipului de analiză lingvistică pe care dorim să îl efectuăm pentru identificarea proprietăților pertinente.

3. *Locuțiuni verbale, predicate complexe*

3.1. *Locuțiuni verbale - o analiză generativă*

Între cologațiile care prezintă un interes deosebit din punctul de vedere al analizei pe care o aplicăm, ne-am oprit la o categorie specială de cologații, locuțiunile verbale, deoarece acestea reprezintă o clasă foarte numeroasă de expresii în toate limbile europene:

avoir besoin (a avea nevoie), porter bonheur (a purta noroc), to make a decision (a lua o decizie), a-și aduce aminte, a face obiectul

După Gledhill (2006), aceste expresii sunt formate dintr-un verb care descrie un proces și un complement care precizează sensul expresiei.

În cadrul unei analize lingvistice generative, verbele care intră în componența locuțiunilor verbale sunt considerate ca fiind golite de sens, verbele capătă un rol de verb auxiliar, sensul locuțiunii fiind propus de către complementul acesteia. Întîlnim astfel noțiunea de verb suport (Storrer, 2006), « *light verb* » (Kearns, 1989) sau „constructions converses » (Gross, 1989), care consideră că verbele din această categorie trebuie analizate doar din punct de vedere sintactic, obiectul direct fiind cel care dă sensul locuțiunii. Aceste analize ignoră constituenții

sintactici care se combină cu acestea și care pot oferi de asemenea informații importante despre gradul de libertate pe care îl avem în folosirea expresiei respective. Dacă putem folosi o locuțiune combinată cu o anumită clasă de prepoziții, atunci probabil avem de a face cu o adevărată locuțiune verbală. Astfel, am constatat că analiza clasică de tip generativ nu este suficient de completă, și vom considera în continuare că locuțiunile verbale pot fi tratate ca predicate complexe, cu proprietăți sintactice și semantice de sine stătătoare.

3.2. *Predicate complexe*

(Gledhill, 2006) consideră aceste construcții verb-substantiv (notate VS) având o serie de proprietăți sintactice și semantice care sunt o rezultată a proprietăților verbului și al substantivului considerate separat. Astfel, construcțiile de tip VS au proprietăți similare unui verb simplu (morfologie, diateza, complemente). În unele situații putem deriva un verb plecând de la substantiv (*a face apel – a apela, a lua o decizie – a decide*), dar acest lucru nu este întotdeauna posibil (*a purta ghinion – *a ghinion*), sau (*a face obiectul – a obiecta?*). În ceea ce privește folosirea diatezei active sau pasive, formele pasive nu sunt întotdeauna posibile (*a lua o decizie - decizia a fost luată, a face obiectul – *obiectul a fost făcut*).

În același timp, construcțiile VS au proprietăți specifice unui substantiv : poate fi definit sau nu, poate fi modificat de către o propoziție relativă sau poate fi transformat în substantiv. Astfel, putem observa folosirea sistematică a articolului definit sau nedefinit (*a face obiectul - - a face un obiect?; a face apel – a face un apel?*). Substantivul nu poate fi întotdeauna modificat de către o clauză relativă (*a luat decizia care se impunea, dar *a făcut referirea care trebuia*)

Substantivul, care este complementul verbului, joacă un rol semantic important, precizând sensul verbului. Cum sensul este acela de proces (stare sau eveniment), complementul nu este doar obiectul verbului, ci reprezintă tipul de proces care are loc (mental sau material).

Terenul a făcut obiectul unui litigiu, care s-a rezolvat la tribunal.

Comisia de disciplină a luat o decizie rapidă privind suspendarea jucătorului.

Astfel, procesele exprimate de verbele *a face* și *a lua* sunt procese abstracte, exprimate de data aceasta de către complementele directe (obiectul, o decizie).

Verbele *a face, a lua* sunt foarte productive, făcând parte din componența multor locuțiuni verbale. De aceea am studiat proprietățile morfologice și sintactice ale celor mai des utilizate locuțiuni în corpusul de lucru.

4. *Cîteva rezultate*

4.1. *Un corpus specializat*

Deoarece proiectul este orientat spre un studiu comparativ al colocațiilor în franceză, română, germană, avem nevoie de corpusuri paralele în cele 3 limbi, alinate la nivel de cuvânt și de propoziție. Un corpus care îndeplinește condițiile este corpusul AcquisCommunaire (ACC) (<http://langtech.jrc.it>). Corpusul conține 17 milioane cuvinte în limba română, 16 milioane de cuvinte în limba franceză, 15 milioane cuvinte în germană. Documentele conțin articole de lege și directive legate de legislația europeană. Limbajul este specific textelor juridice, stilul

este impersonal, iar expresiile fixe sunt foarte numeroase (*se face trimitere, se face apel la...*).

Avem la dispoziție o versiune a corpusului neetichetată precum și o versiune etichetată, care ne-a fost pusă la dispoziție de către Institutul de Cercetări pentru Inteligență Artificială al Academiei Române. Pentru a putea realiza o analiză statistică corectă, am eliminat din corpus unele elemente de structură a documentelor (grafice, tabele etc.), deoarece modifică rezultatele analizei efectuate. Versiunea etichetată a fost realizată aplicând TreeTagger (Schmid, 1994) pentru limbile franceză și germană, iar pentru limba română TTL și MeTT (Tufiș & Dragomirescu, 2004). Pentru a putea compara proprietățile locuțiunilor interesante identificate în corpusul AcquisCommunaire, a trebuit să creem un corpus general (CG) care să permită verificarea datelor extrase din corpusul specializat (alcătuit din ziare, romane (1984), manuale de utilizare (Php), care însumează 2 milioane de cuvinte.

4.2. Cîteva observații asupra verbului ‘a face’

Am realizat studiul cu ajutorul programului WordSmith care permite identificarea concordanțelor (contextelor unui cuvînt), sortarea acestora în funcție de contextele stîng sau drept. Contextele unui cuvînt sunt reprezentate de o fereastră de n cuvinte (am ales $n=5$). Am realizat o căutare folosind formele *face/fac/făcut/făceam/făceau/face*, și am analizat rezultatele care se găsesc în dreapta verbului, imediat după verb sau la un cuvînt distanță sortate în ordinea descrescătoare a frecvenței. Printre 20 cele mai frecvente cuvinte care apar imediat după verb, regăsim multe articole (*unui, unora*), prepoziții (*dintre, din, pentru*), conjuncții (*sau*). Cum pe noi ne interesează în special construcțiile VS, am selecționat doar substantivele care apar imediat după verb și care au sens în limba română:

Tabel 1: Cele mai frecvente construcții VS pentru verbul ‘a face’, extrase din corpusul ACC

Expresie	Număr de apariții în ACC imediat după verb	Număr de apariții în CG
<i>A face obiectul</i>	2869	6
<i>A face referire</i>	1336	6
<i>A face parte</i>	1038	20
<i>A face trimitere</i>	476	7
<i>A face dovada</i>	209	2
<i>A face față</i>	86	9
<i>A face notificarea</i>	71	0

Pentru corpusul general, care este mult mai limitat decât ACC, frecvențele obținute sunt diferite de cele obținute pentru corpusul ACC. Astfel, marea majoritate ale expresiilor celor mai frecvente în corpusul CG sunt cele de forma *V+Prepoziție*. Doar expresia *a face parte* o regăsim între primele 20 de expresii frecvente.

În ambele corpusuri, am urmărit identificarea unor proprietăți specifice fiecărei construcții în cele două corpusuri pentru limba română. Astfel, am urmărit următoarele aspecte :

- 1) dacă substantivul este articulat sau nu, dacă acceptă articol definit, sau nedefinit ;

- 2) proprietățile complementului indirect;
- 3) folosirea unor prepoziții speciale;
- 4) folosirea adverbelor între verb și substantiv ;

Expresia *a face obiectul* este folosită întotdeauna sub această formă (obiectul este articulat, iar articolul este definit). Ea este urmată de un substantiv în cazul genitiv/dativ în majoritatea cazurilor. Modificatorul substantivului este un substantiv reprezentând un termen juridic sau un proces abstract (*modificării, deciziei, litigiului* etc.). Între verb și substantiv pot apare diverse adverbe caracteristice: *deja, de asemenea, imediat* etc. Un comportament asemănător a fost constatat și în limba franceză. Pentru expresia *a face parte*, substantivul nu este articulat, iar expresia este folosită în mod sistematic împreună cu prepoziția *din*. Pentru expresiile *a face referire*, de asemenea putem constata că substantivul este folosit mereu fără articol (definit sau nu), iar prepozițiile care urmează imediat după această expresie sunt *la* și *în*, urmate de un substantiv indicând locul în document (articol, paragraf, alineat etc.). De asemenea, pentru expresiile *a face față* sau *a face apel*, substantivele nu sunt articulate și se folosesc exclusiv în această formă. Expresia *a face față* este urmată de un substantiv în cazul genitiv/dativ, modificatorul substantivului poate fi considerat ca fiind complement indirect al predicatului complex « *a face față* ». Aceste preferințe pentru una din proprietățile morfologice sau sintactice sunt identice în cele două corpusuri (chiar dacă cel general trebuie încă îmbogățit) arată că putem încerca o caracterizare a claselor de colocații cu ajutorul unui ansamblu de proprietăți identificate pe baza unei analize lingvistice.

5. Perspective

Articolul de față se concentrează doar asupra unei clase specifice de locuțiuni verbale, cele generate de verbul 'a face'. Am identificat unele proprietăți interesante ale locuțiunilor, cum ar fi preferința pentru un anumit tip de articol (definit, nedefinit) sau pentru o anumită prepoziție. Aceste proprietăți vor fi specificate în dicționarul de colocații care este în curs de realizare.

Mulțumiri. Autoarea este recunoscătoare organizației AUF (Agence Universitaire pour la Francophonie), care finanțează acest proiect în cadrul rețelei « Lexicologie, Terminologie, Traduction » pe durata mai 2006-martie 2007. De asemenea, autoarea mulțumește doamnei Rada Mihalcea pentru corpusul românesc pus la dispoziție de către aceasta.

Referințe bibliografice

- Gross, G. (1989). *Les constructions converses du français*, Genève-Paris, Droz.
- Gledhill, C. (2006). Vers une analyse systémique des locutions verbales, constructions verbo-nominales et autres prédicats complexes, *La Linguistique systémique fonctionnelle et la langue française* (D.Banks ed.), ERLA, Brest, Université de Bretagne Occidentale
- Hausmann, F.J. (2004). Was sind eigentlich Kollokationen?, *Wortverbindungen – mehr oder weniger fest* (K.Steyer ed.), pp. 309-334
- Heid, U., Ritz, J. (2005). Extracting collocations and their contexts from corpora, *Proceedings of COMPLEX-2005, Conference on Computational Lexicography and Text Research*,

Budapest, juin 2005.

- Kearns, K.(1989). Predicate Nominals in Complex Predicates, *MIT Working Papers in Linguistics*, 10, 123-134.
- Smadja, F. (1993). Retrieving collocations from text: Xtract.*Computational Linguistics*, 19(1): 143-177.
- Seretan, V., Nerima, L., Wehrli, E.(2004). A tool for multi-word collocation extraction and visualization in multilingual corpora. *Actes du congrès EURALEX'2004*, Lorient, France, Vol. 2, pp.755-766
- Storrer, A. (2006). Corpus-based investigations on German support verb constructions. *Collocations and Idioms: Linguistic, lexicographic, and computational aspects*, London: Continuum Press (Fellbaum, Christiane ed.).
- Tufiș, D., Dragomirescu, L. (2004). Tiered Tagging Revisited. In *Proceedings of the 4th LREC Conference*, Lisabona, pp. 39-42.

REZOLUȚIA ANAFOREI PENTRU LIMBA ROMÂNĂ

GABRIELA PAVEL¹, OANA POSTOLACHE², IONUȚ PISTOL¹, DAN CRISTEA^{1,3}

¹*Facultatea de Informatică, Universitatea "Al.I.Cuza", Iași*

²*Institute of Information Sciences, University of Southern California*

³*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

{pavelg, ipistol, dcristea}@info.uaic.ro, oana@isi.edu

Rezumat

În această lucrare se prezintă un model de rezoluție a anaforei pentru limba română, implementat în motorul general de rezoluție RARE, și pașii făcuți până în momentul de față în direcția rezolvării referințelor anaforice pentru limba română. Au fost adaptate aplicații existente pentru marcarea automată pe un text de intrare a informațiilor necesare rezoluției. Printre acestea, marcatorul de grupuri nominale a fost dezvoltat plecând de la reguli învățate dintr-un corpus adnotat manual. Se prezintă primele rezultate și dezvoltări preconizate.

1. Introducere

Una dintre problemele domeniului lingvisticii computationale, care rezistă încă asalturilor cercetătorilor și după aproape 30 de ani de eforturi continue, o reprezintă rezoluția referințelor anaforice. O referință anaforică este o secvență lexicală (numită și *anafor*), în general realizată printr-un grup nominal, care poate avea diferite interpretări în funcție de contextul în care apare. Secvența lexicală care determină interpretarea anaforului, în general precedându-l pe acesta în text, se numește *antecedent*. Relația dintre anafor și antecedent poartă numele de *relație anaforică*. În general, pentru găsirea acestei relații trebuie luate în considerare o gamă diversă de proprietăți morfologice, sintactice și semantice ale cuvintelor, în contextul lor de utilizare. Rezoluția anaforei are însemnate aplicații în regăsirea documentară inteligentă, în sistemele de întrebare-răspuns, în sistemele de inferențe textuale etc. În grupul de Tehnologii ale Limbajului Natural de la UAIC¹ s-a încercat soluționarea acestei probleme prin crearea unui motor de rezoluție simbolic, numit RARE (Cristea, Postolache, 2002a), conceput a fi suficient de general pentru a oferi soluții de rezoluție practic în orice context. Într-un scenariu de utilizare al motorului, acesta ar urma să primească în intrare texte românești și să scoată în ieșire lanțuri coreferențiale într-o adnotare XML. Cu el au fost efectuate experimente și s-au implementat deja modele de rezoluție pentru limba engleză, care sunt la nivelul altor realizări cunoscute în lume (o rată de succes de 61%, un recall de 73%).

În lucrarea de față vom arăta cum poate fi folosit motorul RARE pentru soluționarea problemelor de rezoluție anaforică în limba română. În continuare se prezintă motorul RARE, maniera lui de funcționare pentru limba engleză, și propunerile pentru dezvoltarea motorului pentru limba română. În final se dau rezultatele preliminare ale procesului de evaluare și se amintesc planurile de dezvoltare ulterioară.

¹ <http://consilr.info.uaic.ro/research/>

2. Motorul de rezoluție RARE

RARE (*Robust Anaphora Resolution Engine*) este un program care funcționează ca un cadru general de rezoluție a referințelor anaforice. El poate fi considerat un motor de rezoluție pentru că, la fel ca un motor de sistem expert, pentru a funcționa asupra unui text de intrare, trebuie să fie alimentat cu un “program” ce descrie comportarea lui în operațiunile de identificare a antecedentilor anaforilor. RARE are la bază o concepție asupra relației anaforice (Cristea et al., 2002a) conform căreia *nivelul textului*, populat cu expresii referențiale (notate RE în Figura 1), trebuie separat de nivelul semantic, sau *al cogniției*, unde rezidă reprezentări ale entităților de discurs (notate DE în Figura 1). Intermediar acestor două niveluri se află plasat un *nivel al proiecțiilor* informațiilor din text. Elementele acestui nivel sunt structuri de atribute (PS în Figura 1). Programul care pune în mișcare motorul RARE, numit **model**, are în componență patru elemente: un *set de atribute* care caracterizează descrierile obiectuale ale elementelor celor trei niveluri; un *set de surse de cunoaștere*, ca proceduri elementare capabile să găsească valorile corespunzătoare setului de atribute ale nivelului proiecțiilor; un *set de reguli sau euristici* capabil să răspundă la întrebarea dacă expresia referențială curentă este menționată pentru prima dată sau a mai fost menționată în textul precedent, caz în care să decidă cărui obiect de tip DE de pe nivelul cogniției îi corespunde obiectul PS curent de pe nivelul proiecției; un *domeniu de accesibilitate referențială*, care descrie un set de reguli de limitare a căutării unui antecedent (distanță și căutare liniară versus ierarhică).

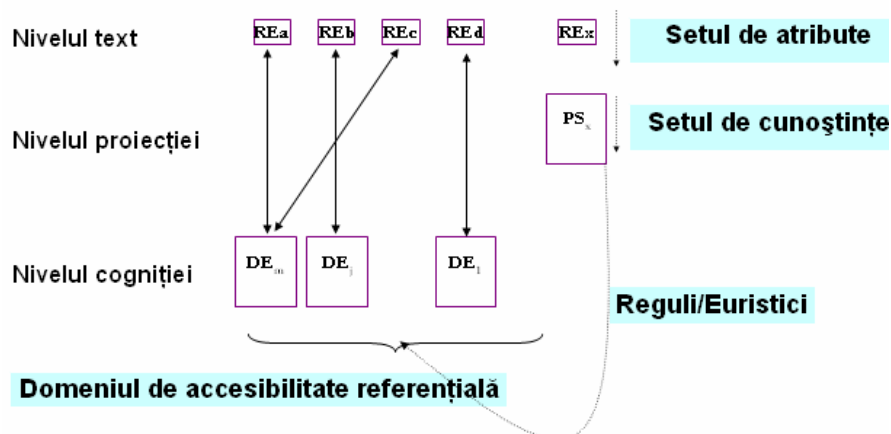


Figura 1: Motorul RARE

3. Implementări RARE ale rezoluției pentru limba engleză

Implementările de sisteme de rezoluție anaforică efectuate până în prezent utilizând RARE (Cristea, Postolache, 2005) au urmărit recunoașterea antecedentilor nepredicaționali pentru anafori pronominali cât și pentru anafori generali de tip grup nominal (nume comune sau proprii) în limba engleză. Problemele de rezoluție au inclus cazuri elementare de rezoluție bazată pe acorduri în gen și număr, dar și cazuri speciale, cum sunt cele în care apar dezacorduri în gen și număr între anafor și antecedent, diferență în leamă, recuperarea antecedentilor distribuiți, ori rezoluții amânate. O parte din aceste cazuri necesită implicarea unor surse de cunoaștere deosebit de sofisticate, mergând până la capacitatea de a recunoaște și opera cu restricții semantice și pragmatice în context ori de a manipula cunoaștere generală

despre lume. În principiu, dacă astfel de surse pot fi găsite, atunci ele ar putea fi incluse în modelul motorului și rezoluții de acest gen s-ar putea realiza. Din păcate însă, resursele posibil a fi angrenate actualmente într-un proces de rezoluție nu sunt capabile de mai mult decât de simulări de mică anvergură, insuficiente pentru a descrie complexitatea situațiilor reale cu care se pot confrunta anumite procese de rezoluție.

Majoritatea sistemelor de rezoluție actuale accesează în intrare un text pe care s-au plasat deja o seamă de notații (Mitkov, 2001), în principal legate de identificarea expresiilor referențiale și de proprietățile morfo-sintactice ale acestora. Astfel de preprocesări influențează semnificativ performanțele sistemelor și vor fi utilizate și în implementarea descrisă în această lucrare, care este o adaptare la limba română a precedentelor create în colectivul nostru, pentru limba engleză.

Detaliem în cele ce urmează câteva elemente ale unui model RARE construit pentru limba engleză (Postolache, Cristea, 2004)². Setul de attribute conține: lema cuvântului, numărul lexical, partea de vorbire, rolul sintactic al RE în propoziție sau o legătură de dependență funcțională, o indicație asupra întinderii de text acoperită de RE, setul de RE incluse (unde e cazul), indicația dacă un nume propriu este nume de familie și dacă numele proprii mici reprezintă nume masculine ori feminine etc. Toate aceste attribute pot primi valori prin accesarea unor proceduri (surse de cunoaștere) ce fac apel la adnotări anterior plasate în formatul de intrare al textului de către preprocesoare, sau le determină ad-hoc în faze incipiente ale procesului de rezoluție. Componenta a treia (setul de reguli ori euristici de rezoluție) implementează trei tipuri de reguli: demolatoare – responsabile de invalidarea unui anumit candidat (de exemplu, RE incluse nu pot niciodată fi coreferențiale); certificatoare – care, dimpotrivă, stabilesc cu precizie un anumit candidat ca antecedent (de exemplu, pe baza identității numelor proprii); reguli cu scor – prin a căror aplicare se mărește sau se micșorează un scor global asociat unei perechi formate din PS-ul curent și un DE candidat a fi considerat antecedent (o astfel de regulă, de exemplu, calculează probabilitatea ca RE-urile corespunzătoare DE-ului candidat să poată fi referite printr-unul din pronumele *he, she, it* sau *they*; o alta mărește scorul dacă anaforul și antecedentul se acordă în număr; legături de sinonimie și hipernimie, determinate prin accesul la WordNet se constituie în alte reguli cu scor, ș.a.m.d.). În sfârșit, domeniul de accesibilitate implementează un tip de căutare liniară, înapoi dinspre anafor spre începutul textului, precum și o limitare la un număr de propoziții (parametru).

4. Implementarea modelului românesc

Intrarea a fost analizată în prealabil cu serviciile web ale ICIA (Institutul de Cercetare în Inteligență Artificială al Academiei Române), care se știe că identifică corect caracteristicile morfo-sintactice în mai bine de 98% din cazuri. Notațiile ANA obținute la ieșire conțin informații morfologice referitoare la parte de vorbire, gen, număr, caz, articol, persoană, mod și timp verbal. Pentru ușurința de tratare ulterioară, etichetele ANA complexe, care înglobează condensat toate aceste informații morfo-sintactice, au fost decodificate la liste de perechi atribut-valoare într-o notație XML alternativă.

Exemplul de mai jos ilustrează secvența de text *o zi frumoasă* cu notațiile ICIA, ca rezultat al preprocesării utilizate de RARE:

² Pentru amănunte suplimentare a se vedea (Postolache, 2004).

```

<root>
  <W id="W0" LEMA="un" ANA="Tdfsr">o</W>
  <W id="W1" LEMA="zi" ANA="Ncfsrn">zi</W>
  <W id="W2" LEMA="frumos" ANA="Afpfsry">frumoasa</W>
</root>

```

Recodarea explicatorie a perechii atribut-valoare ANA corespunzătoare lexemului *zi*, produce:

```
<W ID="W1" LEMA="zi" POS="N" NUM="SG" NGEN="F" />
```

Identificarea grupurilor nominale s-a realizat cu un extractor antrenat pe corpusul *1984* (George Orwell) adnotat inițial la parte de vorbire cu *pos-tagger*-ul ICIA. Grupurile, marcate pozițional manual în corpus, au fost selectate la unică apariție și sortate, după care asupra lor s-au aplicat o seamă de reguli de generalizare. Rezultatul a fost o listă de șabloane, care notează pozițional marcase ANA. De exemplu, șablonul de mai jos:

```
{0={ana=ts}, 1={ana=nsrn}, 2={ana=a}, flagPos=[0, 1, 2]}
```

este capabil să recunoască o secvență lexicală formată dintr-un articol nehotărât, un substantiv și un adjectiv, ca în secvența *o zi frumoasă*. Aplicatorul de șabloane va încadra apoi grupul între etichete <NP></NP>. Astfel, secvenței menționate i se asociază următoarea adnotare ca ieșire a detectorului de grupuri nominale:

```

<NP HEADID="W1" ID="NP1">
  <W ID="W0" LEMA="un" POS="DET" NUM="SG" NGEN="F" />
  <W ID="W1" LEMA="zi" POS="N" NUM="SG" NGEN="F" />
  <W id="W2" LEMA="frumos" POS="A" NUM="SG" NGEN="F" />
</NP>

```

Doar o parte dintre atributele modelului englezesc au fost reținute în modelul RARE de rezoluție anaforică pentru limba română: partea de vorbire, numărul, genul și lema. Datorită inexistenței la momentul actual al unui parser sintactic pentru limba română³ atributul care indică rolul sintactic al RE-ului a fost eliminat.

În modelul românesc au fost menținute toate regulile demolatoare și certificatoare, ele aplicându-se identic în română ca și în engleză. O parte a regulilor englezești cu scor au fost menținute în implementarea românească, ca de exemplu cele de testare a numărului și a lemei. În acest set au trebuit însă operate și modificări care să reflecte diferențele care există între cele două limbi. De exemplu, în română a putut fi adăugată o regulă care mărește scorul în cazul unei potriviri în gen între anafor și antecedent, atribut inexistent pentru substantivele limbii engleze. Testarea potrivirii în gen poate fi luată în considerare în cazul în care pronume referă substantive (grupuri nominale) de același gen (LuperFoy, Rich, 1988). De exemplu, în secvența:

Maria scrie poezii... Versurile ei sunt frumoase.

aplicarea acestei reguli poate duce la concluzia corectă că pronumele personal în dativ *ei* referă *Maria*, dacă există disponibilă o sursă de cunoaștere specializată capabilă să recunoască drept feminin genul substantivului propriu *Maria*. Această sursă ar trebui să genereze următoarea notație pentru aceste două secvențe lexicale:

```

<NP ID="N1">
  <W ID="W7" NUM="SG" NGEN="F" POS="N" LEMA="Maria">Maria</W>
</NP>

```

³ În curs de elaborare, v. (Moruz et al., 2006).

REZOLUȚIA ANAFOREI PENTRU LIMBA ROMÂNĂ

```
<NP ID="N9">  
  <W ID="W14" NUM="SG" NGEN="F" ROLE="" POS="N" LEMA="ea">ei</W>  
</NP>
```

Aplicarea regulii de coreferențialitate pe criterii de acord în gen produce o ieșire RARE de genul:

```
<DE ID="2" reList="N1,N9" />
```

care indică că grupurile nominale N1 și N9 sunt coreferențiale.

De asemenea, regula care verifică posibilitatea ca un antecedent nume comun să fie referit printr-unul din pronumele *he, she, it, they* a fost modificată pentru a putea lucra cu pronumele românești, *el, ea, ei, ele*, cât și cu toate variantele lungi ori prescurtate ale acestora.

S-au menținut regulile de verificare a sinonimiei și hipernimiei prin accesul la WordNet-ul românesc⁴. De exemplu, pentru aceeași secvență ca mai sus, în WordNet se găsește că lemele *poezie* (la singular) și *versuri* (la plural) sunt sinonime (fac parte dintr-un același synset):

```
<W LEMA="poezie" SYN_ID="ENG20-05981555-n" />  
<W LEMA="versuri" SYN_ID="ENG20-05981555-n" />
```

Secvențele corespunzătoare șirurilor *poezii* și *versurile* au următoarea reprezentare în intrarea motorului:

```
<NP ID="N2">  
  <W ID="W9" NUM="PL" NGEN="F" ROLE="" POS="N"  
    LEMA="poezie">poezii</W>  
</NP>  
  
<NP ID="N8">  
  <W ID="W13" NUM="PL" NGEN="F" ROLE="" POS="N" LEMA  
    ="versuri">Versurile</W></NP>
```

5. Evaluare și dezvoltări ulterioare

Textul folosit în acest exercițiu a fost construit de autori în ideea de a fi scurt dar foarte bogat în expresii referențiale. El conține doar 33 de grupuri nominale, relațiile anaforice fiind adnotate manual ca lanțuri coreferențiale (Mitkov, 2001). În această etapă am fost interesați numai de realizarea unui prototip, îmbunătățirea lui urmând a se realiza în continuare.

Ca întotdeauna, probarea performanțelor unui sistem de rezoluție trebuie făcută prin compararea ieșirii motorului, care notează lanțurile coreferențiale determinate automat, cu cele ce s-au notat manual pe corpusul considerat standard. Așa cum s-a exemplificat mai sus, în urma rulării motorului se obține o listă de entități de discurs (DE-uri), fiecare dintre acestea având asociată o listă de expresii referențiale (RE-uri) găsite a fi în relația de coreferință. În felul acesta fiecare DE are semnificația unui lanț de coreferențialitate (trivial, acestea putând a avea și lungimea 1, reprezentând entități cu unică menționare). Corpusul de probă a fost marcat în maniera identică ieșirii motorului RARE: o listă de DE-uri, fiecare listând RE-urile lanțului.

Pentru evaluare s-a folosit evaluatorul asociat motorului. Rezultatele evaluării au următoarele valori:

```
SUCCESS_RATE: 0.696969696969697  
MUC_PRECISION = 0.25
```

⁴ <http://multiwordnet.itc.it/online/multiwordnet>

MUC_RECALL = 0.6
MUC_F-Measure = 0.35294117647058826

Pe viitor se urmărește în primul rând adăugarea de noi reguli în model. O astfel de regulă avută în vedere este *WhRule*, care, în varianta pentru limba engleză, se referă la depistarea antecedentilor pronomelor relative.

Lista de erori obținută din rulările pe acest corpus inițial va fi folosită pentru corectarea regulilor și recalcularea ponderilor acelor celor cu scor. Etapa următoare va consta în utilizarea unui corpus de mari dimensiuni, probabil o parte a romanului „1984”, deja adnotat parțial la coreferințe.

Avem în vedere, totodată, realizarea unui mecanism de reglare automată a ponderilor asociate regulilor, mecanism care va folosi tehnici de învățare automată. Sperăm ca prin acest mod să realizăm un sistem hibrid simbolic-statistic, care ar trebui să aibă performanțe superioare atât unuia pur simbolic cât și unuia pur statistic.

Nu în ultimul rând se dorește integrarea motorului în alte proiecte care folosesc rezoluția anaferei, unul dintre acestea fiind un proiect de întrebare-răspuns pe limba română.

Referințe bibliografice

- Cristea, D., Postolache, O., Dima, G.E., Barbu, C. (2002). AR-Engine – a framework for unrestricted coreference resolution. Appeared in *Proceedings of Language Resources and Evaluation Conference - LREC 2002*, Las Palmas de Gran Canaria, Spain, 29-31 May 2002, vol. VI, p.2000-2007.
- Cristea D., Postolache O.D. (2005). How to deal with wicked anaphora, in António Branco, Tony McEnery and Ruslan Mitkov (editori): *Anaphora Processing: Linguistic, Cognitive and Computational Modelling*, Benjamin Publishing Books.
- Moruz, A., Curteanu, N., Trandabăț, D., Dornescu, I., Bolea, C. (2006). Parsarea predicatului (verbal / nominal) și a clauzei (finite / nefinite) în limba română. Aplicare la parsarea FDG. În acet volum.
- Postolache, O. (2004) *RARE – Robust Anaphora Resolution Engine*. Teză de disertație în Lingvistică Computațională, Facultatea de Informatică, Universitatea „Al.I.Cuza” Iași.
- Postolache, O., Cristea, D. (2004): Designing Test-beds for General Anaphora Resolution, in *Proceedings of the Discourse Anaphora and Anaphor Resolution Colloquium – DAARC*, St. Miguel, Portugal.
- Mitkov, R. (2001). *Outstanding issues in anaphora resolution*. Benjamin Publishers.
- LuperFoy, S., Rich, E. (1988). An Architecture for Anaphora Resolution, *ANLP 1988*: 18-24.

INSTRUMENTE PENTRU CONSULTAREA ATLASULUI LINGVISTIC ȘI EDITAREA TEXTELOR DIALECTALE

SILVIU BEJINARIU¹, VASILE APOPEI¹, RAMONA LUCA¹, LUMINIȚA
BOTOȘINEANU², FLORIN OLARIU²

¹*Institutul de Informatică Teoretică, Academia Română, Filiala Iași,*

²*Institutul de Filologie Română „A. Philippide”, Academia Română, Filiala Iași*

*{silviub, vapopei, ramonad}@iit.tuiasi.ro, lumi_florin@personal.ro,
olariuft@yahoo.com*

Rezumat

Această lucrare prezintă ultimele rezultate ale cercetărilor desfășurate în cadrul proiectului „Proiectarea și implementarea unui sistem integrat de aplicații software pentru editarea textelor dialectale și realizarea Atlasului Lingvistic Român, pe regiuni”, proiect interdisciplinar al Academiei Române. Sunt prezentate noile funcții implementate: generarea automată a indexului de cuvinte și forme pentru dicționarele atlasului, respectiv funcția care permite comunicarea editorului de texte dialectale cu alte editoare de text.

1. Introducere

Proiectul intitulat *Noul Atlas lingvistic român, pe regiuni (NALR/ALRR)* a fost inițiat în scopul radiografierii situației actuale a graiurilor vorbite pe teritoriul României, anchetele făcându-se într-un număr de aproximativ 1000 de localități din România. Aceste atlase concretizează rezultatul cercetărilor dialectologi din București, Cluj-Napoca, Iași și Timișoara, ele fiind importante atât pentru lingviști, cât și pentru istorici, geografi sau etnologi (Arvinte et al., 1987, 1997).

Scopul principal al cercetărilor actuale, în care lingvistica descoperă relevanța și flexibilitatea metodelor de lucru furnizate de informatică, a fost de a crea un instrument de tip „computer aided design” care să faciliteze publicarea noilor volume ale atlasului lingvistic regional românesc. Varianta proiectării asistate de calculator a planșelor atlaselor lingvistice, orientare de mare actualitate în geolingvistica internațională, are o serie de avantaje nete, care o recomandă spre a înlocui maniera clasică de editare: pe lângă faptul că înlătură copierea manuală, fiind mai economică sub aspectul costurilor și al timpului de execuție, cel mai important aspect este că ea poate constitui o sursă de informații stocate în format electronic pentru alte lucrări interdisciplinare.

Sistemul prezentat a stat la baza realizării planșelor pentru prospectul celui de-al 3-lea volum al *Noului Atlas lingvistic român pe regiuni. Moldova și Bucovina*, care a fost publicat în anul 2005 sub formă de volum și CD multimedia (Arvinte et al., 2005), (http://iit.tuiasi.ro/editare_td/atlas/atlas.html).

2. Componentele atlasului lingvistic electronic

Proiectarea acestui sistem (aplicația ALR) a fost începută în urmă cu 6 ani și principalele funcții, prezentate deja cu alte prilejuri sunt următoarele:

- crearea și întreținerea dicționarelor atlasului (Bejinariu et al., 2002),
- generarea automată, editarea și tipărirea în diferite formate a planșelor de tip „hartă lingvistică”, respectiv „material necartografiat” (Florea et al., 2002),
- editarea, formatarea și tipărirea de texte dialectale (aplicația EditTD) (Apopei et al., 2003),
- generarea automată de pagini HTML pentru conținutul dicționarelor,
- consultarea sincronizată a edițiilor mai vechi ale atlaselor.

Transcrierea fonetică specifică limbii române

Pentru a putea reda cât mai fidel toate nuanțele rostirii, transcrierea fonetică s-a dovedit a fi un instrument fiabil și, de aici, varietatea alfabetelor fonetice utilizate de specialiști. Pe lângă transcrierea fonetică internațională realizată cu Alfabetul Fonetic Internațional (IPA), specialiștii au dezvoltat sisteme proprii anumitor domenii lingvistice.

În cazul limbii române, transcrierea fonetică presupune folosirea de simboluri asociate sunetelor primare: 68 de variante vocalice (17 vocale simple, fiecare dintre ele având și câte 3 variante accentuate) și 50 de variante consonantice (fig. 1). Simbolurilor primare le sunt asociate semne diacritice ilustrând fenomene fonetice specifice. Acestea sunt în număr de 12, organizate în 5 grupe, în cazul vocalelor, respectiv în număr de 9, organizate tot în 5 grupe, în cazul consoanelor (fig. 2).

simple	diacritice				
a	ă	ă̂	ă̄	ă̆	ă̈
e	ë	ë̂			
i			î		î̂
o	ö	ö̂			
u	ü		û		

b, c, ĉ, c̄, c̆, c̈, d, d̂, d̄, d̆, d̈, f, g, ĝ, ḡ, ğ, g̈, h, ĥ, h̄, h̆, ḧ, j, k, l, l̂, l̄, l̆, l̈, m, m̂, m̄, m̆, m̈, n, n̂, n̄, n̆, n̈, p, r, r̂, r̄, r̆, r̈, s, ŝ, s̄, s̆, s̈, t, t̂, t̄, t̆, ẗ, v, w, z, ẑ, z̄, z̆, z̈, y

Figura 1. Simbolurile primare folosite în transcrierea fonetică a limbii române
Sunetul vocalic marcat cu (*) nu permite aplicarea de fenomene fonetice

Datorită numărului foarte mare de caractere primare, s-a decis folosirea modului Unicode pentru codificarea caracterelor. În cazul vocalelor, un calcul simplu arată că există un număr de 359 combinații de fenomene fonetice. De aici rezultă necesitatea de a proiecta un număr de 359 de fonturi, respectiv $359 \cdot 17 \cdot 4 = 24412$ simboluri grafice, doar pentru vocale.

Fenomenele fonetice pot fi plasate deasupra, sub, dar și lateral față de simbolul de bază. În plus, unele simboluri pot fi plasate deasupra, sau dreapta-sus în raport cu simbolul precedent din text.

Fenomene asociate vocalelor		Fenomene asociate consoanelor	
Grupă	Fenomen	Grupă	Fenomen
Durată	Scurtime	Durată	Semilungime
	Semilungime		Lungime
	Lungime	Palatalizare	Semipalatalizare
Nazalizare	Seminazalizare		Palatalizare
	Nazalizare		Palatalizare mare
Ocluzie glotală	Coup de glotte	Explozie	Explozie
Deschidere	Închidere	Caracter silabic	Caracter silabic
	Semideschidere	Afonizare	Semiafonizare
	Deschidere		Afonizare
	Deschidere mare		
Afonizare	Semiafonizare		
	Afonizare		

Figura 2. Fenomenele fonetice folosite în transcrierea fonetică

Din acest motiv, am decis realizarea unui sistem de generare on-line a imaginii simbolurilor cărora le sunt aplicate fenomene fonetice, prin sinteza imaginilor componente. Pentru editarea textelor folosind transcrierea fonetică specifică limbii române, utilizatorul trebuie să introducă din tastatură simbolul de bază, aplicarea fenomenelor realizându-se prin selectarea acestora direct de pe bara de instrumente. Acest mod de desenare a simbolurilor a permis reducerea numărului de fonturi folosite la numai 2, ambele derivate din fontul „Arial”.

3. *Generarea automată a indexului de cuvinte și forme*

Unul dintre instrumentele care s-au dovedit a fi necesare cercetătorilor lingviști este generatorul automat al indexului de cuvinte și forme. Generarea unui astfel de index presupune identificarea aparițiilor unui sunet sau grup de sunete în dicționarul de transcrieri fonetice.

Prima etapă a acestui proces constă în stabilirea parametrilor de căutare:

- **filtrul „cuvânt”** – stabilește dacă căutarea se realizează între transcrierile fonetice ale unui anumit cuvânt de bază, sau în întreg dicționarul;
- **filtrul „punct de anchetă”** – permite restrângerea căutării la un singur punct de anchetă;
- **căutare în...** – permite stabilirea câmpurilor din dicționar în care se face căutarea;
- **mod căutare** – este folosit pentru a specifica modul în care se realizează căutarea, ca „text” sau ca „transcriere fonetică”;
- **forma de căutat** – permite utilizatorului să editeze textul ale cărui apariții dorește să le identifice.

Caracterele conținute în textele transcrise fonetic sunt însoțite de două grupe de atribute: pe de o parte, fenomenele fonetice, pe de altă parte, atribute care specifică modul de poziționare și desenare. În funcție de parametrii specificați înainte de generarea indexului, aceste atribute sunt luate sau nu în considerare în cursul procesului de căutare.

Procesul de căutare este finalizat prin sintetizarea informației privitoare la ocurențele formei căutate sub forma unei liste (fig. 3), care poate fi tipărită sau poate fi folosită pentru identificarea poziționării acestora în dicționar.

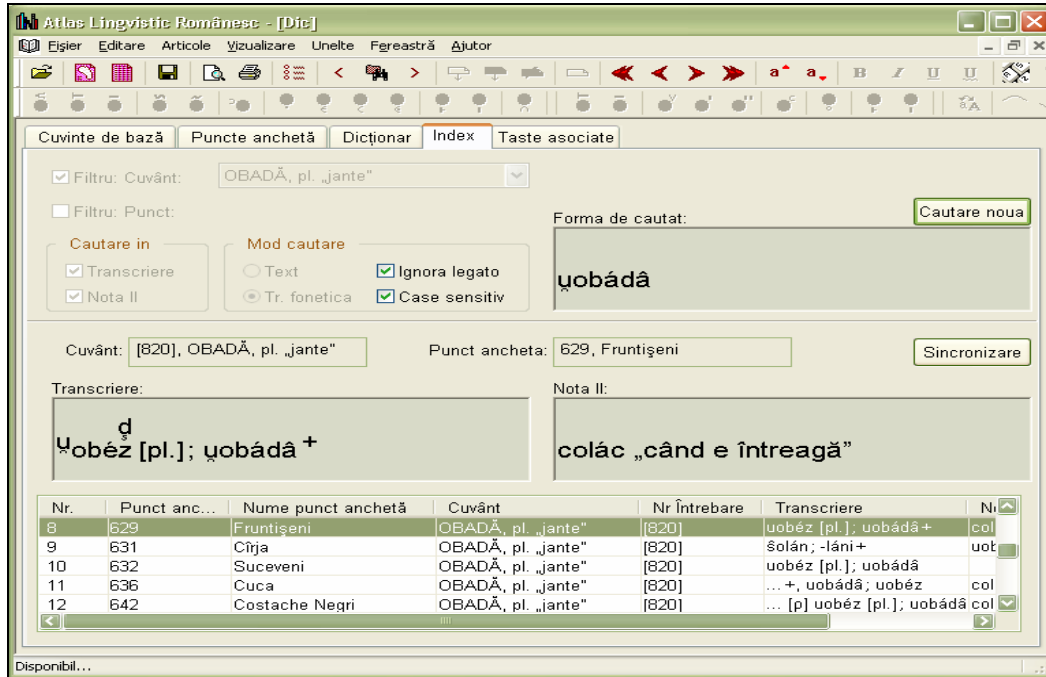


Figura 3. Indexul de forme, generat de aplicația ALR

4. Comunicarea editorului de texte dialectale cu alte editoare de text

Una dintre probleme apărute în cursul proiectării sistemului pentru editarea textelor dialectale a fost aceea de a oferi posibilitatea transferului de fișiere între aplicația EditTD și alte editoare de text, având în vedere faptul că există deja în alte proiecte, texte editate folosind alte modalități de operare, care folosesc fonturi proiectate special în acest scop. În general, au fost folosite două tipuri de fonturi:

- fonturi conținând imaginile caracterelor împreună cu un fenomen fonetic specific fontului respectiv;
- fonturi prezentând doar fenomenele fonetice sub formă de caractere separate, introduse după caracterul căruia îi sunt aplicate.

Pentru realizarea conversiei între textele editate cu alte editoare și aplicația EditTD a fost ales ca intermediar formatul RTF (Rich Text Format), recunoscut pe mai multe platforme și acceptat de editoarele din familia Microsoft.

Realizarea efectivă a importului fișierelor de tip RTF în formatul specific aplicațiilor ALR și EditTD presupune parcurgerea următoarelor etape:

- selecția fișierului de intrare,
- defnirea tabelii de conversie (selecția acesteia în cazul în ea există deja),
- stabilirea regulilor de conversie globale, la nivel de font,
- stabilirea regulilor de conversie la nivel de detaliu (caracter),
- conversia propriu-zisă,

- analiza rezultatului conversiei și eventuale corecții.

Tabela de conversie RTF este o colecție de reguli prin care fiecărei perechi de forma *{caracter_inițial, font_inițial}* i se asociază o structură de forma

{ caracter, font, fenomen_fonetic, atribut_poziție }

Sistemul permite specificarea regulilor de conversie:

- la nivel global, în cazul când caracterelor din fontul inițial le sunt asociate aceleași caractere din fontul de bază, cărora le sunt aplicate unul sau mai multe fenomene fonetice,
- la nivel de detaliu, pentru indicarea excepțiilor existente între regulile definite la nivel global sau pentru indicarea unor noi reguli.

Definirea unei reguli de traducere presupune selectarea fontului inițial și a caracterului inițial folosit în regula de traducere, selectarea fontului și a caracterului în care se face traducerea, selectarea fenomenelor fonetice specifice ce urmează a fi aplicate în momentul traducerii, precum și poziția caracterului în cazul în care aceasta este modificată în momentul traducerii.

După definirea tabelii de traducere se poate trece la conversia efectivă a fișierului RTF (fig. 4).

Menționăm că procesul de traducere din formatul RTF în formatul intern al aplicației pentru editarea textelor dialectale este realizat cu păstrarea atributelor de formatare a caracterelor și paragrafelor. În plus, sunt tratate special caracterele diacritice specifice limbii române, deoarece aplicația noastră folosește codificarea Unicode a caracterelor, în timp ce în formatul RTF codificarea caracterelor se face pe 8 biți, fiind specificat însă setul de caractere regional care trebuie folosit.

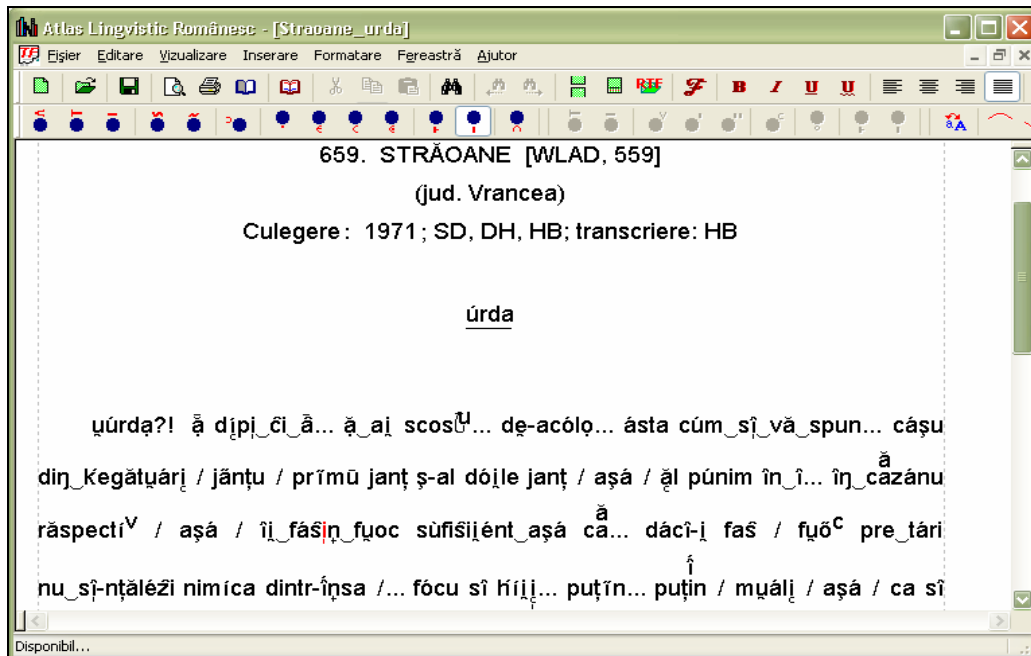


Figura 4. Rezultatul aplicării procedurii de conversie

5. Concluzii

Sistemul proiectat pentru modelarea *Atlasului Lingvistic Român pe Regiuni* este funcțional și a fost folosit pentru publicarea prospectului celui de al 3-lea volum al *Atlasului Lingvistic Român pe Regiuni, Moldova și Bucovina*. Acest prospect a fost publicat sub formă de volum tipărit și CD multimedia, fiind publicat în același timp și pe internet, la adresa http://iit.tuiasi.ro/editare_td/atlas/atlas.html.

Noile funcții implementate, generarea automată a indexului de cuvinte și forme, respectiv funcția care permite comunicarea editorului de texte dialectale cu alte editoare text, se dovedesc foarte utile în cercetare pentru lingviști și dialectologi.

În concluzie, putem spune că sistemul realizat reprezintă la ora actuală o încununare a unor eforturi care au început acum mai bine de 130 de ani (primul sistem de transcriere fonetică aplicat limbii române a fost pus la punct de Émile Picot, în anul 1873).

Aplicarea tehnologiei informatice la transcrierea fonetică a etno- și socio-textelor ușurează extrem de mult procesul editării permițând în același timp realizarea unor colecții de informații lingvistice în format electronic, ce pot fi folosite în cadrul altor cercetări interdisciplinare.

Referințe bibliografice

- Arvinte, V., Dumistrăcel, S., Florea, I.A., Nuță, I., Turculeț, A. (1997). *Noul Atlas lingvistic român, pe regiuni. Moldova și Bucovina*, București, Editura Academiei, vol. I, 1987; vol. II.
- Arvinte, V., Dumistrăcel, S., Florea, I.A., Nuță, I., Turculeț, A. (2005). *Noul Atlas lingvistic român, pe regiuni. Moldova și Bucovina, III. Prospect*, Iași, Editura Universității „Alexandru Ioan Cuza”.
- Florea, I.A., Apopei, V., Olariu, F.T., Bejinariu, S. (2002). „Editarea asistată de calculator a atlaselor lingvistice și a textelor dialectale”, în *Identitatea limbii și literaturii române în perspectiva globalizării*, Editura Trinitas, Iași, 2002, p. 211-232
- Bejinariu, S., Apopei, V., Roman, M. (2002). *Mediu pentru editarea transcrierilor fonetice în Limba Română. Realizarea Atlasului Lingvistic Român pe Regiuni, Limba Română în Societatea Informațională, Societatea Cunoașterii*, Editura Expert, București, p. 423-440.
- Apopei, V., Rotaru, F., Bejinariu, S., Olariu, F. (2003). *Electronic Linguistic Atlases, Proceedings of the International Conference on Information and Knowledge Engineering IKE'03*, June 23, Las Vegas, Nevada, USA, Volume 2, ISBN 1-932415-08-4, p. 628-633
- [http://iit.tuiasi.ro/editare_td/atlas/atlas.html] *Noul Atlas Lingvistic Român pe Regiuni, Moldova și Bucovina III*.

GENERARE DE CONCORDANȚE PENTRU DICȚIONARUL LIMBAJULUI POETIC EMINESCIAN

MIHAELA BRUT¹, DUMITRU IRIMIA², OANA PANAIT²

¹*Facultatea de Informatică, Universitatea "Al.I.Cuza" Iași*

²*Facultatea de Litere, Universitatea "Al.I.Cuza" Iași*

mihaela @info.uaic.ro, tirimia@uaic.ro, oana_pan@yahoo.fr

Rezumat

În urma unui efort susținut de câțiva ani buni de zile, dicționarul limbajului poetic eminescian a fost finalizat, fiind deja disponibil specialiștilor și tuturor iubitorilor de poezie. La dezvoltarea lui a fost utilizată aplicația „Concordanțe eminesciene”, dezvoltată de un colectiv din Cluj la începutul anilor 90. Lipsită de scalabilitate din cauza tehnologiei depășite pe care o utilizează, această aplicație a fost totuși singurul instrument disponibil pentru procesarea limbii române de care s-a putut folosi colectivul care a lucrat la dezvoltarea dicționarului.

Articolul de față își propune să prezinte mecanismul intern de funcționare a acestei aplicații, pașii de procesare a textului eminescian, încercând să ofere o propunere de rescriere a ei apelându-se la tehnologiile Web-ului semantic.

1. Introducere

Dicționarul limbajului poetic eminescian este un proiect inițiat în anii '90 de Facultatea de Litere a Universității „Al. I. Cuza” din Iași, vizând două componente majore:

- *Concordanțele poeziilor eminesciene* - urmându-se modelul din lexicografia poetică europeană;
- *Semne și sensuri poetice* - avându-se drept model *Dictionnaire des Symboles* (Chevalier & Gheerbrant, 1994)

În privința primei componente a proiectului, până acum au fost finalizate și tipărite *Concordanțele poeziilor antume* (Irimia, 2004) și *Concordanțele poeziilor postume* (Irimia, 2006), de un real folos celor care doresc să se apropie de poezia eminesciană, să înțeleagă și urmărească modul în care cuvintele limbii române au fost încărcate de poeticitate, de noi sensuri și semnificații. O pătrundere de profunzime a forței semantice a limbajului poetic eminescian nu poate fi atinsă, însă, fără a urmări întreaga creație a poetului, de aceea dezvoltarea pe viitor a concordanțelor din teatru, proză literară, publicistică, critică literară/teatrală /muzicală, corespondență, însemnări manuscrise ar fi de bun augur pentru critica literară și pentru iubitorii de poezie. În alte literaturi sunt demult disponibile concordanțele complete ale operei poetilor reprezentativi: G.Leopardi, E. Montale, G. Pascoli, D'Annunzio, G. Ungaretti (Italia); William Blake, Lord George Gordon Byron, John Butler Yeats (Anglia); Federico Garcia-Lorca (Spania); Emily Dickinson (SUA).

Vom prezenta în continuare instrumentele informatice utilizate în generarea primelor două seturi de concordanțe eminesciene, accentuând deficiențele acestor instrumente și prezentând și o propunere de rescriere și îmbunătățire a lor.

2. *CONCORD și SILEX*

Pentru obținerea concordanțelor eminesciene au fost utilizate două instrumente informatice dezvoltate de colective aparținând Universității “Babeș-Bolyai” din Cluj:

- *CONCORD* – *Sistem de lematizare automată și generare a concordanțelor*, coordonat de prof. dr. Sanda Cherata (Cherata, 1996);

- *SILEX* – *Sistem lexical informatizat*, sub coordonarea cerc. Teodor VUȘCAN; acest sistem a fost utilizat pentru realizarea analizei morfologice automatizate a cuvintelor din poezia eminesciană, fiind de fapt integrat în sistemul CONCORD (Vușcan, 1996).

Ambele sisteme sunt proiectate utilizându-se sistemul de gestiune a bazelor de date Foxpro 2.6. A fost luată ca reper *Poezia eminesciană* - edițiile Perpessicius (1952) și D.Murărașu (1970-1972) -, care a fost transformată într-o bază de date FoxPro numită ME.POE, conținând câte o înregistrare pentru fiecare poezie. Structura acestei baze de date include codul autorului, codul volumului, codul ciclului de poezii, codul subciclului, codul poeziei, titlul poeziei, subtitlul acesteia, dedicația scrisă de Eminescu pe marginea poeziei, *motto*-ul acesteia, notele de final ale poeziei, precum și textul integral al poeziei (inclus într-un câmp de tip memo):

Fiecare poezie este lematizată separat, lemele sunt analizate sintactic utilizându-se SILEX, iar rezultatul acestor operații este inclus în câte o bază de date asociată fiecărei poezii. Într-o astfel de bază de date, este alocată câte o înregistrare fiecărui cuvânt din poezie, fiind incluse: codul volumului, al ciclului, al subciclului, al poeziei, numărul versului în cadrul poeziei, numărul liniei pe care este afișat acest vers, numărul cuvântului în cadrul versului, codul contextului (numărul versului asociat cuvântului - diferă în cazul cuvintelor afișate pe rândul următor al versului din care fac parte, în cadrul versurilor mai lungi), lema de care aparține cuvântul curent, categoria gramaticală, atributul eventual al acestei categorii (de exemplu, în cazul categoriei *adjectiv* poate exista atributul *posesiv* etc.), variantele eventuale ale lemei (de exemplu, pentru seară - sară), numărul caracterului din cadrul versului de la care începe cuvântul și numărul caracterului la care acesta se sfârșește.

Se poate observa cantitatea imensă de informație redundantă, repetată la fiecare cuvânt în parte. În cazul în care ar fi fost utilizată o structură XML pentru stocarea poeziilor eminesciene, fiecare dintre informații ar fi trebuit furnizată o singură dată, deoarece informațiile ce țin, de exemplu, de un volum ar fi fost furnizate ca sub-elemente ale elementului <Cod_vol> etc.

Utilizându-se lematizarea și analiza sintactică automatizată a poeziilor în maniera descrisă mai sus, se poate efectua o primă generare a concordanțelor, utilizându-se opțiunile *Gen_Conc* → *Integrare Volum*, apoi *Gen_Conc* → *Integrare Operă* puse la dispoziție de CONCORD.

Din cauza faptului că unele leme nu sunt recunoscute corect de analizor, lingviștii trebuie să efectueze o primă corectură pe hârtie a concordanțelor generate, în special în privința analizei gramaticale și a recunoașterii corecte a unei leme în formele flexionare în care aceasta apare în poezii. Aceste corecturi sunt apoi operate electronic utilizând opțiunea *Lem_Poem* → *Corect_Lem* din cadrul CONCORD. Pentru operarea corecturii asupra unei leme, trebuie furnizat codul volumului, codul poeziei, numărul versului în care apare lema respectivă (aceste informații se găsesc pe pagina listată unde s-a efectuat corectura), dar și numărul cuvântului ce urmează a fi corectat în cadrul versului (în acest caz, numărarea făcându-se manual și fiind dificilă mai ales în cazul cuvintelor aflate după poziția a zecea în cadrul versului).

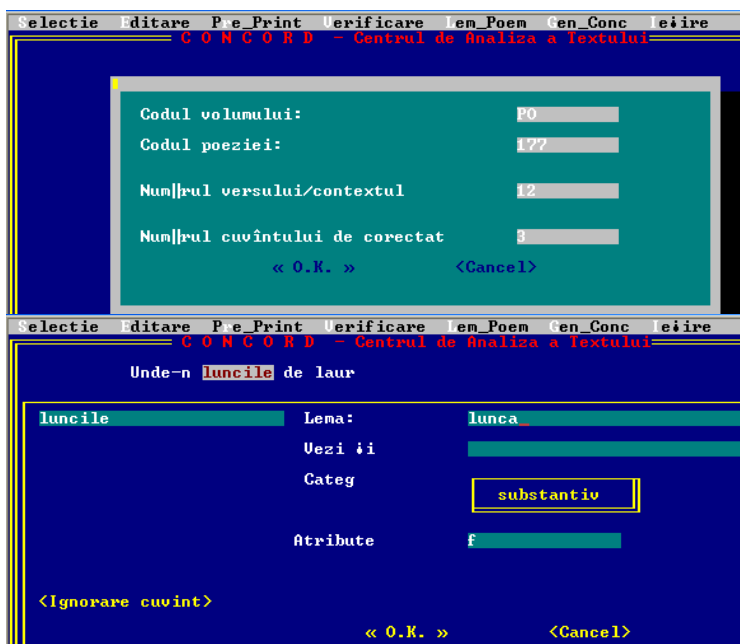


Figura 1: Operarea corecturii asupra unei leme în cadrul CONCORD

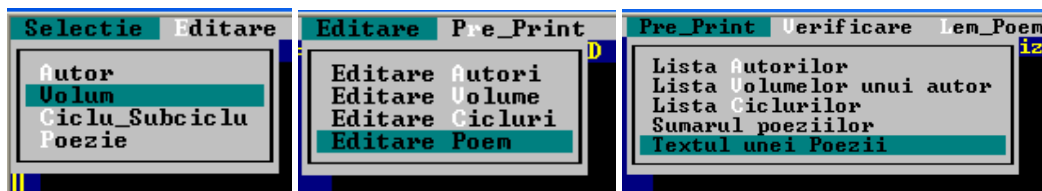


Figura 2: Facilități oferite de CONCORD

CONCORD pune la dispoziție câteva facilități suplimentare ce pot fi utile pentru a ușura munca centrală de operare a corecturilor. De exemplu, opțiunea *Selecție* → *Volum* permite selectarea unui volum, al cărui cod va apărea automat în caseta *Codul volumului* din fereastra de corectare a unei leme, nemaitrebuind să fie introdus manual la fiecare leamnă în parte. Pentru a efectua corecturi asupra textului unei poezii (în cazul în care au fost observate erori), în locul modificării manuale a bazei de date corespunzătoare poate fi utilizată opțiunea *Editare* → *Editare Poem*. Pentru a vizualiza într-o manieră formatată textul unei poezii există facilitatea *Editare* → *Editare Poem*.

După operarea corecturii asupra tuturor lemelor, se realizează o nouă generare a concordanțelor, utilizându-se aceleași opțiuni *Gen_Conc* → *Integrare Volum*, apoi *Gen_Conc* → *Integrare Operă* oferite de CONCORD. Rezultatul acestor operații este alcătuit din următoarele fișiere:

- Me_C.tlm → conținând concordanțele întregii poezii eminesciene;
- Me_ct_C.txt → lista lemelor, ordonate pe clase lexico-gramaticale;
- Me_fr_C.txt → lemele ordonate descrescător după frecvența de utilizare;
- Me_Lm_C.txt → lemele ordonate alfabetic.

Deoarece aceste fișiere sunt de tip text, au fost folosite anumite codificări, anumite convenții de reprezentare a caracterelor românești, a caracterelor speciale sau a indicațiilor de formatare. În plus, lemele a căror categorie gramaticală nu a fost recunoscută de către program au fost marcate cu mențiunea „,??” pentru a fi analizate cu atenție la corectura manuală.

Codificările utilizate în fișierele prezentate mai sus sunt interpretate în cadrul aplicației *MS Word*, utilizând două macrocomenzi:

- *ConvConc* → aplicată fișierului Me_C.tlm și având drept rezultat documentul Word formatat incluzând concordanțele întregului volum.
- *Conv_Car_Conc* → aplicată fiecăruia dintre fișierele Me_ct_C.txt, Me_fr_C.txt, Me_Lm_C.txt pentru a se obține listele cu diverse ordonări ale lemelor sub formă tot de document Word formatat.

109	Domnul	sb pr	25	0.00026
PO	27	55	Tron de-argint și strălucire,/Maica Domnului zîmbea;/Iar pe schițele de scară	
PO	28	25	«Mărire ție Doamne! O, Iehova, mărire!	
PO	29	249	Stinge, puternic Doamne , cuvîntul nimicirii	
PO	29	288	Mai bine stinge, Doamne , viața ginții mele,	
PO	43	180	Genunchind să-i ierte Domnul osînditul lui păcat.	
PO	43	870	Cerul lumea o cuprinde cu sinistru-i mîndru Domn .	
PO	44	318	De Domnul eu trimisu-s, căci te iubește mult,	
PO	44	328	A pus în tine Domnul nemargini de gîndire.	
PO	44	433	Gîndite de Domnu -ntr-a sorilor nimb.	
PO	52	9	Era vremi acelea, Doamne , cînd gravura grosolană	
PO	104	28	În albă mantie de Domn .	
PO	128	11	«Pietrea, Doamne Sfînte, căzu în orice colț,	
PO	159	43	Și azi cînd am puterea ce-o are numai Domnul ,	
PO	179	39	Decis-ai oare, Doamne , ca în etern să fie	
PO	179	58	Mai bine stinge, Doamne , viața națiunii mele,	
PO	206	10	A dreptății, dragă Doamne ,	
PO	240	60	Astei stranie povești.../ - Doamne, Doamne , mătucico,/Nu ți-i greu să mai vorbești?	
PO	240	193	Doamne! limpede mai știe	
PO	240	276	Doamne , cum nu ai un dascal, să te tragă de urechi.	
PO	240	311	MUȚI: Cum aș începe... Doamne? Doamne -ah, mîmucută,	
PO	240	385	Doamne , iartă-mă! Dar bine, serios vorbești tu, fată?	
PO	240	412	- Doamne, Doamne , mătucică, hai să zicem că-l iubesc.	
110	Don	ls	3	0.00003
PO	183	4	spaniolește de Signorul Don Lopez de Poeticales	
PO	183	42	Știi: Don Manuel, perfidul,	
PO	183	48	Căci Don Manuel ș-acuma	

Figura 3: Extrase din documentele de concordanțe

Documentele formatate obținute sunt din nou parcurse atent de către lingviști. Dacă mai apar corecturi, acestea trebuie din nou operate în cadrul CONCORD, fiind necesară o nouă generare a concordanțelor conform pașilor expuși mai sus. În cazul în care lingviștii își dau acordul în privința formei finale a concordanțelor, cele patru documente Word formatate trebuie integrate într-unul singur, ce urmează a fi aranjat în pagină conform cerințelor specificate de editura la care va fi tipărit volumul de concordanțe.

3. Concluzii

Programele SILEX și CONCORD au oferit suportul necesar pentru obținerea concordanțelor poeziei eminesciene, însă ar fi necesară o reabordare a acestor sisteme din perspectiva tehnologiilor moderne. Principalele deficiențe ale celor două aplicații, constatate pe parcursul utilizării efective a lor în cadrul *Dicționarul limbajului poetic eminescian* sunt:

- Lipsa de documentare a codului FoxPro 2.6, astfel încât efectuarea unor modificări pe codul sursă al aplicațiilor este foarte dificilă;
- Stocarea redundantă a informațiilor, ilustrată de structura bazelor de date ce conțin poezia eminesciană;

Rang	Lema	Cl	Atrib.	Nr. ap	Frecv.	Rang	Lema	Cl	Atrib.	Nr. ap	Frecv.
90	Roma	sb	pr	44	0.00023	1	în	pp		6558	0.03360
	barbă	sb	f	44	0.00023	2	i	cj		6432	0.03295
	casă	sb	f	44	0.00023	3	fi	vb	p	5864	0.03004
	copac	sb	m	44	0.00023	4	de	pp		7322	0.03751
	laur	sb	m	44	0.00023	5	cu	pp		3558	0.01823
	masă	sb	f	44	0.00023	6	eu	pn	pers	3546	0.01817
	mei	sb	m	44	0.00023	7	al	ar	pos	3280	0.01680
	patimă	sb	f	44	0.00023	8	tu	pn	pers	3200	0.01639
	pânză	sb	f	44	0.00023	9	pe	pp		2806	0.01437
	preot	sb	m	44	0.00023	10	avea	vb	a	2752	0.01410
	răs	sb	n	44	0.00023	11	să	cj		2628	0.01346
	suspîn	sb	n	44	0.00023	12	ea	pn	pers	2566	0.01315
	veac	sb	n	44	0.00023	13	sine	pn	refl	2532	0.01297
	vin	sb	n	44	0.00023	14	nu	av		2406	0.01233
91	amar	sb	n	42	0.00022	15	ce	pn	relat	2400	0.01229
	arbore	sb	m	42	0.00022	16	un	ar	m	1856	0.00951
	bucurie	sb	f	42	0.00022	17	un	ar	f	1840	0.00943
	fluviu	sb	n	42	0.00022	18	din	pp		1742	0.00892
	iarnă	sb	f	42	0.00022	19	el	pn	pers	1722	0.00882
	joc	sb	n	42	0.00022	20	la	pp		1374	0.00704
	miere	sb	f	42	0.00022	21	ca	av		1336	0.00684
	mijloc	sb	n	42	0.00022	22	lume	sb	f	1236	0.00633
	mister	sb	n	42	0.00022	23	ei	pn	pers	1116	0.00572
	plan	sb	n	42	0.00022	24	cç	cj		1076	0.00551
	rând	sb	n	42	0.00022	25	tău	aj	pos	1042	0.00534
	sărutare	sb	f	42	0.00022	26	când	av		934	0.00478
	sicriu	sb	n	42	0.00022	27	mai	av		848	0.00434
	sunet	sb	n	42	0.00022	28	ochi	sb	m	828	0.00424
	uitare	sb	f	42	0.00022	29	vedea	vb		790	0.00405
	vers	sb	n	42	0.00022	30	viață	sb	f	786	0.00403
	zid	sb	n	42	0.00022	31	prin	pp		758	0.00388
92	Ana	sb	pr	40	0.00020	32	vrea	vb		754	0.00386
						33	cum	av		740	0.00379
						34	lui	ar	m	712	0.00365
						35	meu	aj	pos	666	0.00341

Figura 4: Generarea concordanțelor în format text utilizând CONCORD

- Fluxul de operații necesare obținerii concordanțelor este unul inflexibil și mare consumator de timp. Astfel, pentru a vizualiza o corectură efectuată în cadrul CONCORD trebuie repetați toți pașii de generare a concordanțelor, vechile instanțe de documente conținând concordanțele fiind neutilizate.

Precum am evidențiat și mai sus, structurarea poeziei eminesciene în format XML ar elimina deficiența menționării redundante a informațiilor. Utilizarea unui sistem de lematizare cu marcatori XML ar face disponibilă analiza gramaticală a lemelor și altor aplicații de procesare a textelor românești. Sistemul de generare a concordanțelor ar putea utiliza, în acest caz, o procesare DOM, beneficiind de bibliotecile de funcții de procesare existente deja în mai multe limbaje de programare.

Actualmente există mai multe instrumente de generare a concordanțelor¹. Problema majoră este aceea că fiecare astfel de instrument impune anumite restricții legate de formatul fișierelor ce vor fi procesate. O posibilă soluție de reproiectare a CONCORD ar presupune selectarea unui astfel de sistem (sau proiectarea de la zero), integrarea unui analizor gramatical, precum și reproiectarea sistemului de reprezentare a datelor, în cazul în care cel original nu corespunde formatului XML în care avem la dispoziție opera poetică pe care dorim să o procesăm.

Mulțumiri. Autorii sunt recunoscători colectivului de cercetători de la Universitatea “Babeș-Bolyai” din Cluj care le-a pus la dispoziție sistemul CONCORD pentru procesarea concordanțelor eminesciene.

Referințe bibliografice

- Cherata, S., (1996). CONCORD: Sistem de realizare a concordanțelor textelor poetice românești. *Limba și Tehnologie*, Dan Tufiș editor, Ed. Academiei Române, București, 1996, pp. 215-220.
- Chevalier, J., Gheerbrant, A. (1994). *Dictionnaire des symboles : Mythes, rêves, coutumes, gestes, formes, figures, couleurs, nombres*, Robert Laffont Edition.
- Irimia, D. coord (2004). *Dicționarul limbajului poetic eminescian. Concordanțele poeziilor antume*, vol. I-II, Editura Hyperion, Botoșani.
- Irimia, D. coord (2006). *Dicționarul limbajului poetic eminescian. Concordanțele poeziilor postume*, vol. I-IV, Editura Univ. „Alexandru Ioan Cuza” Iași, 2006.
- Irimia, D. coord (2005). *Dicționarul limbajului poetic eminescian. Semne și sensuri poetice. I. Câmpul semantica ARTE*, Editura Univ. „Alexandru Ioan Cuza” Iași.
- Vușcan, T., (1996). SILEX - sistem lexico-morfologic computerizat pentru limba română. *Limba și Tehnologie*, Dan Tufiș editor, Ed. Academiei Române, București, 1996, pp. 209-214.

¹ Exemple de instrumente de generare a concordanțelor: *AntConc*, Waseda University, Japonia: http://www.antlab.sci.waseda.ac.jp/antconc_index.html; *Concordance*, R.J.C. Watt of Dundee University: <http://www.dundee.ac.uk/english/wics/wics.htm>; *Monoconc*, Athelstan: <http://www.athel.com/mono.html#monopro>; *Wordsmith*, Mike Scott, Oxford University: <http://www.oup.com/elt/catalogue/isbn/6890?cc=gb>

CREAREA UNUI GENERATOR MORFOLOGIC PENTRU VERBELE DIN LIMBA ROMÂNĂ

ANTONINA BÎRLĂDEANU, NATALIA BURCIU

Facultatea de Calculatoare, Informatică și Microelectronică, Universitatea Tehnică a Moldovei, Chișinău

{toni_birlad, natusicb}@yahoo.com

Rezumat

Resursele computerizate pentru limba română reprezintă un suport de bază pentru dezvoltarea instrumentelor automate și a aplicațiilor lingvistice dedicate procesării informației lingvistice specifice pentru gramatica, fonetica, și lexicul limbii române. În acest articol prezentăm rezultatul unui studiu efectuat asupra morfologiei verbelor din limba română, precum și etapele elaborării unui generator morfologic în baza rezultatelor obținute în urma acestui studiu. Rezultatele obținute de acest generator pot fi folosite pentru corectarea greșelilor de flexiune, sau chiar pentru prevenirea acestora, în cadrul traducerii automate în limba română, sau în alte aplicații lingvistice.

1. Introducere

Una din cercetările efectuate asupra morfologiei limbii române a fost crearea unui model de formalizare a morfologiei limbii române. În cadrul acestui proiect, cercetătorii și-au propus ca scop realizarea unui analizor (corector și generator) morfologic, și realizarea unei baze de date ce conține numai atributele specifice morfologiei. E de menționat că atributele introduse au permis crearea unei prime versiuni de corector morfologic, cu posibilități de realizare a unui analizor morfologic complet, și pe de altă parte a unui generator de paradigme concretizat într-un program de învățare automată a conjugării verbelor. În cadrul acestui proiect s-au definit clasele flexionare, s-au determinat clasele flexionare la verb, și clasele flexionare pentru nume, și s-a determinat codificarea rădăcinii cuvintelor și structurile de date (Peev et al., 1997).

O altă aplicație ce studiază structura morfologică a cuvintelor este aplicația Anmor. Această aplicație reprezintă un mediu de dezvoltare/actualizare pentru modelul morfologic paradigmatic al limbii române, iar preocupările sale esențiale sunt asigurarea corectitudinii și completitudinii datelor. Componentele principale sunt: un verificator de erori sintactice și de inconsistențe, un editor al dicționarului, și un asistent în procesul de îmbogățire cu noi cuvinte a bazei de date. (Cosman, 2002). Studiarea acestor aplicații lingvistice a și stat la baza creării generatorului nostru morfologic.

2. Generatorul morfologic

Scopul acestui proiect a fost dezvoltarea unui generator morfologic bazat pe reguli pentru verbele din limba română. Am ales limba română fiindcă este una din limbile flexionare, care face ca obiectivul pe care ni l-am pus să fie mai dificil. Scopul principal a fost crearea regulilor pentru partea de vorbire flexionară verbul, prin atribuirea acestuia a caracteristicilor sale formale sub formă de atribut-valoare, pentru ca mai apoi aceste caracteristici să poată fi folosite pentru generatorul nostru morfologic. Inițial am început prin crearea regulilor pentru verbe, iar aceasta s-a dovedit a fi o sarcină foarte complexă, având în vedere că nu au fost de găsit reguli deja existente pentru această parte de vorbire.

Ideea de bază a acestui generator morfologic a fost crearea unor reguli de formare a verbelor, depinzând de atributele și valorile acestora: conjugarea (conjugarea I, II, III, IV), modul (Indicativ, Conjunctiv, Condițional-Optativ, Infinitiv, Gerunziu, Participiu, Supin), timpul (prezent, imperfect, perfectul compus, perfectul simplu, mai mult ca perfectul, viitorul (viitorul simplu, viitorul anterior)), persoana (I, II, III), precum și numărul (singular, plural), iar rezultatul (output-ul) acestui generator să fie salvat într-un fișier *.txt, pentru ca mai apoi acesta să poată fi folosit pentru alte aplicații lingvistice.

În Tabelul 1 sunt prezentate un set de reguli create folosind formalismul atribut-valoare, pentru atributul *Indicativ* și valoarea *Mai-mult-ca-perfect*. Au fost create manual circa 1700 reguli, dintre care 280 de reguli pentru verbele de bază și mai mult de 1400 pentru verbele neregulate (verbele auxiliare, verbele modale, etc.). Făcând o paralelă la numărul de reguli pentru alte părți de vorbire, e de menționat că acestea constituie aproximativ 200 de reguli pentru părțile de vorbire substantiv și adjectiv.

Tabel 1. Crearea regulilor pentru atributul *Indicativ* și valoarea *Mai-mult-ca-perfect*

<i>Verbul la Infinitiv</i>	<i>Conjugarea</i>	<i>Terminația conjugării</i>	<i>Persoana</i>	<i>Numărul</i>	<i>Verbul</i>	<i>Regulile create</i>
a învăța	I	V-a	1	Sg	învătasem	V+sem
			2	Sg	învătaseși	V+seși
			3	Sg	învătase	V+se
			1	Pl	învătaseram	V+serăm
			2	Pl	învătaserați	V+serăți
			3	Pl	învătasera	V+seră
a dormi	IV	V-i, V-î	1	Sg	dormisem	V+sem
			2	Sg	dormiseși	V+seși
			3	Sg	dormise	V+se
			1	Pl	dormiserăm	V+serăm
			2	Pl	dormiserăți	V+serăți
			3	Pl	Dormiseră	V+seră

3. Descrierea algoritmului de funcționare a generatorului morfologic

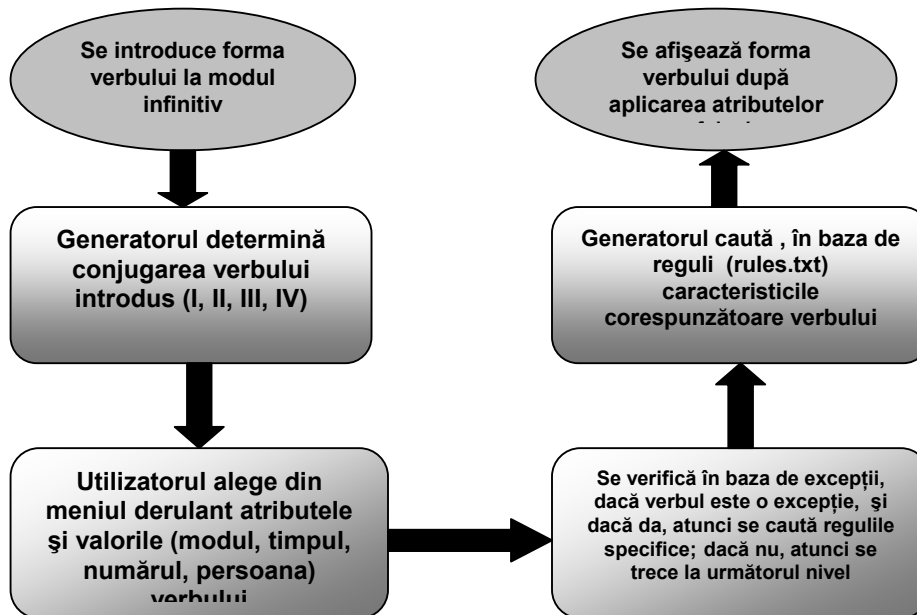


Figura 1. Algoritmul de funcționare a generatorului morfologic

În figura 1 este prezentat schematic algoritmul de funcționare a generatorului morfologic. Sistemul începe prin introducerea verbului la modul infinitiv. După care urmează determinarea conjugării verbului, pentru a putea trece la următoarea etapă.

În cazul programului nostru dispunem de mai multe fișiere care conțin informații morfologice, de tipul conjugările verbelor, regulile de formare a acestora, excepțiile (de exemplu, verbele auxiliare, etc.). După determinarea conjugării se cere alegerea caracteristicilor morfologice ale verbului, din meniul derulant al programului.

Apoi generatorul verifică în fișierul de excepții dacă verbul introdus este o excepție, și dacă da, atunci se extrag și se adaugă regulile proprii acestui verb-excepție, iar dacă nu, atunci se adaugă regulile ce depind de caracteristicile (atribute – valori) cerute de către utilizator.

În final se afișează rezultatul generării morfologice și acest rezultat este salvat într-un fișier *.txt, pentru ca aceste rezultate să poată fi folosite în alte aplicații lingvistice.

Codul

```

II Indicativ Perfectul Simplu 2 Plural V - ea + urați |
II Indicativ Perfectul Simplu 3 Plural V - ea + ură |
ex01III Indicativ Perfectul Simplu 1 Singular V - ea + ui |
ex01III Indicativ Perfectul Simplu 2 Singular V - ea + uși |
    
```

În figura 2 este prezentată o formă în care sunt incluse atributele și valorile verbelor, precum și rezultatele generării acestora.

Figure 2. Afişări

4. Concluzii și cercetări ulterioare

În acest proiect am creat prin mijloace lingvistice și computaționale un generator morfologic pentru o categorie flexionară a limbii române, și anume, verbul. În cadrul acesteia au fost create reguli de generare a verbelor, bazându-ne pe regulile gramaticale de formare a verbului în limba română. Au fost create aproximativ 1700 de reguli de generare. Odată cu integrarea în Uniunea Europeană este necesară crearea mult mai multor aplicații lingvistice, aceasta reprezentând una din motivațiile proiectului nostru. În viitor planificăm să adăugăm noi părți de vorbire pentru acest generator morfologic, și să rezolvăm problema diacriticelor.

Referințe bibliografice

- Dicționarul ortografic, ortoepic și morfologic al limbii române* (2000). Editura Academiei Române, București.
- Popescu, Ș. (1997). *Gramatica practică a Limbii Române*, Editura Lider, București.
- Boatcă, M., Crihană, M. (1996). *Manual preparator de Gramatică a Limbii Române*, Editura Mondan, Moldova.
- Peev, L., Bibolar, L., Jodal, E. (1997). Un Model De Formalizare A Morfologiei Limbii Române. In: Dan Tufiș, Poul Andersen (eds.). *Recent Advances in Romanian Language Technology*. ISBN 973-27-0626-0, Editura Academiei Române.
- Cosman C. M. (2002). *Morfologia paradigmatică a limbii române. Mediu de dezvoltare / actualizare*. Lucrare de licență. Facultatea de Informatică, Universitatea Al.I. Cuza Iași.

PARSAREA PREDICATULUI (VERBAL / NOMINAL) ȘI A CLAUZEI (FINITE / NEFINITE) ÎN LIMBA ROMÂNĂ. APLICARE LA PARSAREA FDG

ALEX MORUZ^{1,2}, NECULAI CURTEANU¹, DIANA TRANDABĂȚ^{1,2}, IUSTIN DORNESCU^{1,2}, CECILIA BOLEA¹

¹*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

²*Facultatea de Informatică, Universitatea "Al.I.Cuza" Iași*

{[curteanu](mailto:curteanu@iit.tuiasi.ro), [mmoruz](mailto:mmoruz@iit.tuiasi.ro), [dtrandabat](mailto:dtrandabat@iit.tuiasi.ro)}@iit.tuiasi.ro

Rezumat

Lucrarea prezintă abordări și experimente de parsare FDG pentru limba română. Sunt puse în evidență metode de parsare a structurilor globale (inter-clauzale) și locale (intra-clauzale), cu accent pe parsarea grupului verbal.

1. Introducere

Articolul de față își propune să schițeze folosirea de strategii și programe de procesare a limbajului natural (LN), deja existente în cadrul colectivelor de cercetare din Iași și București, pentru proiectarea unui parser FDG (*Functional Dependency Grammar*) la nivel de frază (*sentence*), pentru limba română. Se pleacă de la strategia SCD (*Segmentare-Coeziune-Dependență*) de segmentare și parsare a clauzei (și sintagmelor subclauzale) (Curteanu *et al.*, 2005), intrarea în program fiind reprezentată de text (multiplu) adnotat la categorii morfologice și marcheri SCD. Conform algoritmului de segmentare-parsare SCD, o fază esențială în stabilirea corectă a clauzelor finite este *determinarea predicatelor finite*, verbale sau nominale.

Să precizăm de la început că în abordarea noastră folosim termenul de „*finit*” pentru toate *formele verbale* ce corespund unui *mod personal*, în timp ce termenul „*nefinit*” este atribuit formelor verbale ce corespund „modurilor” nepersonale, cunoscute și ca *forme absolute* ale verbului (*infinitiv*, *gerunziu*, *participiu*, *supin*). În acest sens, clauza (sau propoziția) al cărei predicat are ca nucleu semantic un verb *finit* (formă predicativ/verbală), o vom numi *clauză finită*. Astfel, *clauza finită* dintr-o frază este definită ca fiind întinderea de text aflată între doi marcheri SCD de nivel clauzal (sau de discurs), care conține (exact) un predicat finit (cu nucleu semantic predicativ/personal). *Clauza nefinită* corespunde formelor verbale *nefinite*, predicativ/nepersonale ale categoriilor lexicale majore V (verb), N (substantiv), și A (Adjectiv-Adverb), care posedă trăsătura de *predicaționalitate* (*deverbalitate*). Parsarea predicatului finit (verbal sau nominal) revine la determinarea *grupului verbal finit* (Verbal Group, VG). VG mai este cunoscut în literatura de specialitate și sub denumirea de *Complex Verbal* (Monachesi, 2005), (Barbu, 1999). Facem observația că determinarea VG (finit) de natură copulativă, (notat TASG – *Tense Auxiliary SubGroup* în (Curteanu & Trandabăț, 2006)), este o etapă esențială în parsarea predicatului nominal.

Lanțul operațiilor de parsare este următorul: la intrare, textul este adnotat morfologic și apoi la markeri SCD, rezultatul acestui proces fiind un text multiplu (*heavy*) adnotat. Pe adnotarea morfologică a acestui text se realizează parsarea VG (Curteanu *et al.*, 2006). Pentru determinarea clauzelor din componența unei fraze, cât și a relațiilor inter-clauzale dintre acestea, este folosit un program de segmentare-parsare la clauză bazat pe algoritmul de parsare SCD (Curteanu *et al.*, 2005). În urma execuției acestui program se obține un arbore de dependență a clauzelor (*arbore clauzal*) din cadrul unei fraze date. Pentru realizarea arborilor de dependență intra-clauzali propunem două soluții complementare: o *abordare deterministă*, bazată pe reguli, rezultată din algoritmul de parsare SCD intra-clauzală (Moruz, 2006), și o *abordare statistică*, bazată pe algoritmi de învățare automată. Scopul este de a obține un parser mai performant prin *combinarea* unor algoritmi de parsare diferiți (atât statistici cât și bazați pe reguli). În acest moment se află în desfășurare un proiect care are ca scop crearea unui corpus de fraze adnotate la dependențe funcționale pentru limba română.

2. Etape și Instrumente

2.1. Parsarea VG

Primul pas în adnotarea FDG a unei fraze date este determinarea predicatelor din textul de intrare. Această operație este necesară nu numai pentru segmentarea frazei în clauze, ci și pentru găsirea proprietăților *grupurilor verbale* (VGs) în vederea determinării corecte a dependențelor intra-clauzale (de exemplu, în diateza pasivă subiectul gramatical devine *obiect semantic*, iar complementul direct devine *agent*). În urma parsării predicatelor, datele obținute sunt următoarele:

(a) *Nucleul semantic* al fiecărui VG. Această componentă a predicatului este importantă deoarece face diferența dintre *predicatele nominale* și *predicatele verbale*. În cazul în care nucleul semantic este de tip *predicațional (deverbal)*, VG reprezintă un *predicat verbal*; în cazul în care avem un nucleu verbal de tip *copulativ*, VG reprezintă nucleul sintactic al unui *predicat nominal*. În acest al doilea caz, *nucleul semantic* al predicatului nu mai este verbul copulativ din VG ci argumentul acestuia, *numele predicativ*. În “*Ion pleacă acasă.*”, nucleul semantic este verbul “*pleacă*”, în timp ce predicatul propoziției “*Ion a fost student.*” este nominal, având ca nucleu semantic substantivul nepredicațional “*student*”, pe când nucleul VG este copulativul “*fost*”.

(b) *Diateza formal-sintactică (de suprafață)* a predicatului. Aceasta este de fapt diateza gramaticală clasică. În multe cazuri, diateza sintactică și cea semantică nu coincid, astfel argumentele *directe* (ce corespund valenței) ale verbului predicațional nu sunt aranjate corect în lista SUBCAT a argumentelor. De exemplu, diateza sintactică a predicatului în propoziția “*Mașina se spală.*” este cea *reflexivă*, dar *diateza semantică* este cea *pasivă*. Procesul de determinare a VGs împreună cu proprietățile lor este descris în (Curteanu *et al.*, 2006). În urma găsirii nucleului semantic și a diatezei sintactice, dacă nucleul este predicațional (echivalând cu situația că nu este verb copulativ), se coboară în lexicon prin proiecția *FX-bar inversă* în vederea determinării *diatezei semantice* a VG, după care se revine de la nivelul lexiconului la nivelul textului de suprafață cu noua diateză și cu restricțiile de *linking* asociate acesteia prin

mecanismul *proiecției FX-bar directe*. Procesul de proiecție FX-bar inversă a VG către lexicon, prin care se determină nucleul semantic al VG, și proiecția FX-bar directă a nucleului semantic al VG către VG, și al VG către clauza finită, este descris detaliat în (Curteanu, Trandabăț, 2006). Câteva exemple de transformare a *diatezei formal-sintactice* în *diateză lexical-semantică* (de observat că *diateza sintactică* pentru toate aceste exemple este cea reflexivă):

- (e1) *Ion și Maria se știu de mici copii.* (diateză semantică = *reciprocă*);
 (e2) *Se știe vinovat de moartea mamei sale.* (diateză semantică = *reflexivă*);
 (e3) *Se știe că pisicile fugăresc șoriciei.* (diateză semantică = *impersonală*).

2.2. Parsarea Clauzelor

Parsarea clauzelor (unei fraze) este realizată prin intermediul unui program de separare a unităților clauzale și subclauzale bazat pe strategia de parsare SCD (*Segmentare-Coeziune-Dependență*) (Curteanu *et al.*, 2005). Separarea unităților clauzale este realizată pe baza claselor de markeri SCD, a ierarhiei de tip graf a acestor clase, și a proprietății de *predicaționalitate (deverbalitate)* a categoriilor lexicale majore N (Noun), V (Verb), și A (Adjectiv-Adverb). De aici rezultă necesitatea determinării markerilor SCD și a predicatelor din frază ca o preprocesare în vederea aplicării algoritmului SCD. Pentru obținerea *arborelui de dependență clauzal* este necesară, pe lângă parsarea clauzelor, și determinarea relațiilor de dependență inter-clauzale. Aceasta se face cu ajutorul markerilor de tip SCD, reprezentați într-o bază de date în care este specificat și tipul de relație (inter-clauzală) pe care o introduc.

2.3. Relații de Dependență pentru Limba Română

În Tabelul 1 sunt exemplificate o parte din relațiile de dependență FDG intra-clauzală ce au fost determinate pentru limba română. Cuvintele scrise cu litere îngroșate reprezintă *fi* în relațiile descrise, iar categoria „X” - orice categorie morfologică.

Tabel 1: Exemple de relații de dependență funcțională FDG pentru limba română

NUCLEU	FIU	Cuv. Urm.	RELAȚIE	ABREV.	EXEMPLU
Substantiv	Prepoziție	Nominal	Atribut substantival	a.subst	<i>Praf de pușcă</i>
Substantiv	Articol	X	Determinant	det	Un om
Substantiv	Verb	X	Atribut verbal	a.verb	Om care merge
Substantiv	Adjectival	X	Atribut adjectival	a.adj	Om tânăr
Adjectiv	Adverb	X	Comparativ	comp.	Mai mare
Verb	Nominal	X	Subiect	Sbj	Ion merge
Verb	Nominal	X	Complement direct	c.d.	Îl văd
Verb	Verb aux	X	Auxiliar	aux	Am mers
Verb	Nominal, adj.	X	Nume predicativ	n.pred	Sunt tânăr
Verb	Negație	X	Negație	neg	Nu stau
Adverb	Adverb	X	Comparativ	comp	Mai repede
Prepoziție	X	X	Rel. prepozițională	prep.	De vorbă
Coordonator	X	X	Rel. de coordonare	coord	Mare și tare

2.4. Parsarea FDG Bazată pe Reguli

Pentru rezolvarea problemei dependențelor funcționale FDG au fost folosite două *abordări complementare*, una bazată *pe reguli* și una *statistică*, bazată *pe învățare automată*. Programul de parsare FDG (*Functional Dependency Grammar*) pornește de la o formalizare a gramaticilor de dependență descrisă de (Järvinen, Tapanainen, 1997). Algoritmului folosește reguli (de fapt, *expresii regulate*) în vederea stabilirii de dependențe; pentru rezolvarea problemei dependențelor la distanță s-a folosit paradigma *Island Parsing*. Island Parsing este o strategie de parsare multidirecțională, utilizată atât în cadrul prelucrării limbajului natural cât și în alte domenii în care robustețea este importantă sau resursele de procesare sunt limitate, atașată gramaticilor independente de context. Este o strategie bidirecțională, în sensul că elemente de parsare incomplete, care corespund părții drepte ale unei reguli de producție independentă de context, pot fi extinse în ambele direcții.

Regulile de parsare au fost determinate experimental, în urma studierii unei serii ample de exemple pentru limba română (Moruz, 2006); din acest motiv, expresiile regulate astfel determinate nu pot fi utilizate pentru parsarea la dependențe funcționale a textelor în alte limbi. În *parsarea bazată pe reguli*, pașii urmași pentru a parsă la dependențe funcționale FDG un text sunt următorii: **(1)** *segmentarea textului* primit la intrare în clauze și unități lexicale (acest pas este realizat pe baza strategiei de segmentare-parsare SCD a textului); **(2)** *delimitarea constituenților* (extragerea de grupuri verbale, nominale, adjectivale, etc.); **(3)** *determinarea elementelor relaționale* de nivel inter-clauzal și de discurs, și atașarea lor la elementele subordonate. Delimitarea constituenților și determinarea elementelor relaționale împreună cu legăturile lor funcționale se realizează cu ajutorul expresiilor regulate. Expresiile regulate în sine nu sunt suficient de particulare pentru a putea determina o structură arborescentă neambiguă. În scopul scăderii ambiguității au fost create reguli de procesare contextuală (cum ar fi *acordul*), ce cresc eficiența determinării structurilor (*ruleNounAdj(Tree t1, Tree t2)*) – dacă rădăcina arborelui *t1* este *substantiv* și rădăcina arborelui *t2* este *adjectivală*, iar cele două rădăcini sunt în acord morfo-sintactic, atunci *t2* devine subarbore pentru *t1*).

2.5. Parsarea Statistică

Rațiunea principală în *abordarea statistică* este că putem folosi o serie de algoritmi de parsare deja dezvoltati pentru alte limbi. Proiectarea unei gramatici formale pentru limba română presupune un efort foarte mare în crearea unui corpus adnotat pe baza căruia să se facă antrenarea modelului, ținând cont de particularitățile fiecărui formalism. Întrucât pentru limba română nu avem la dispoziție un corpus de dependențe funcționale FDG, scopul este de a folosi parsere statistice existente pentru dezvoltarea unui astfel de corpus. Propunem o dezvoltare *iterativ-incrementală* a unui corpus, după cum urmează: la început adnotăm manual un număr relativ mic de fraze; pe baza acestui micro-corpus, antrenăm trei parsere statistice și rulăm cele trei modele obținute pentru adnotarea automată a unui set nou de fraze. Acestea vor fi *doar corectate* de adnotatori umani iar corpusul astfel obținut va fi folosit pentru reantrenarea parserelor; procesul e reluat până când întreg corpusul este adnotat. Clasa parserelor *deplasare-reducere*

folosește un algoritm simplu: parcurge textul cu o fereastră conținând un număr fix de cuvinte. Alte metode de adnotare FDG statistică sunt *arborii parțiali de cost maxim* și *meta-parserele*. *Arborele parțial de cost maxim* pornește prin a construi un graf orientat complet care are ca noduri cuvintele din frază, iar ca arce numărul de apariții în corpus al unei muchii similare. Algoritmul caută un arbore parțial de cost maxim, iar acest arbore reprezintă cea mai probabilă parsare FDG a frazei. *Meta-parserele* sunt modele de combinare a mai multor *parsări*, și nu a mai multor parsere propriu-zise. Deși, intuitiv, o metodă ierarhică de a combina parsările în mod diferențiat, în funcție de diverse criterii, ar trebui să ducă la o creștere a preciziei, astfel de metode au dezavantajul costului computațional destul de ridicat.

3. Performanțe

Datorită absenței unui corpus adnotat la arborii de dependență funcțională pentru limba română, programul nu a putut fi încă testat suficient. Din acest motiv nu se pot oferi estimări realiste în legătură cu precizia atașată procesului de adnotare la dependențe funcționale. În urma verificării manuale a unei părți din arborii obținuți, rezultatele conduc către o precizie destul de mare pentru frazele formate dintr-o singură propoziție.

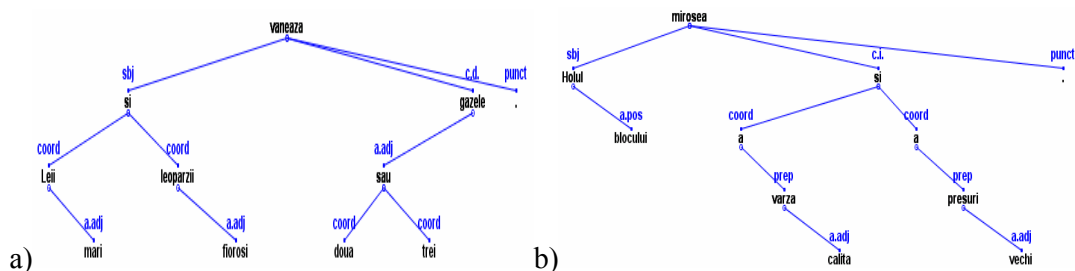


Figura 1: Exemple de parsare FDG: a) „Leii mari și leoparzii fioroși vânează două sau trei gazele.”; b) ”Holul blocului mirosea a varză călită și a preșuri vechi.”;

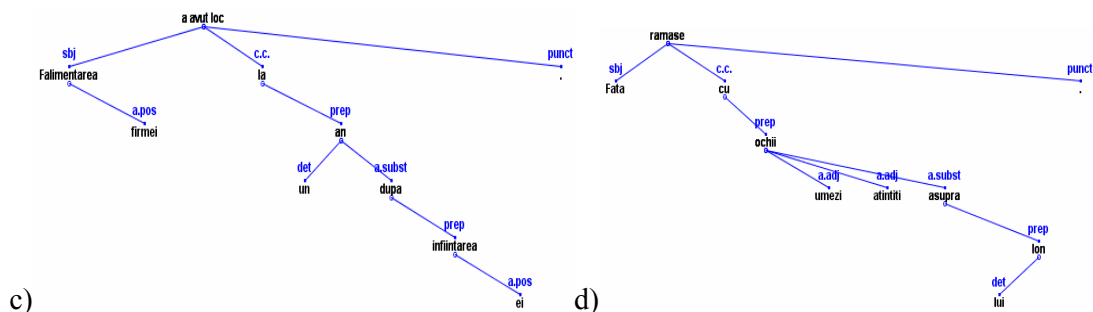


Figura 2: c) ”Falimentarea firmei a avut loc la un an după înființarea ei.”; d) ”Fata rămasese cu ochii umezi ațintiți asupra lui Ion.”;

PRELUCRAREA RESURSELOR ROMÂNEȘTI ÎN CADRUL PROIECTULUI LT4EL

IONUȚ PISTOL¹, ADRIAN IFTENE¹, DIANA TRANDABĂȚ^{1,2}, DAN CRISTEA^{1,2},
CORINA FORĂSCU^{1,3}

¹ *Facultatea de Informatică, Universitatea “Al. I. Cuza”, Iași*

² *Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

³ *Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București*

{ipistol, adiftene, dtrandabat, dcristea, corinfor}@info.uaic.ro

Rezumat

Proiectul LT4eL are ca scop realizarea unei tehnologii multilingve, utilizate în cadrul unui sistem de eLearning, care să faciliteze operațiile de creare a obiectelor de învățare de către profesori și de regăsire a lor de către studenți, inclusiv prin criterii de natură semantică. Până în momentul de față, în cadrul proiectului a fost creat un corpus semnificativ (peste 600.000 cuvinte) de texte românești în diferite formate de adnotare, plecând de la textul original și mergând până la un format XML ce pune în evidență informații morfo-sintactice, cuvinte cheie și definiții. În lucrare sunt prezentate etapele de prelucrare ale documentelor corpusului și modulele de prelucrare utilizate.

1. Introducere

Colectarea și prelucrarea unui corpus românesc semnificativ ca dimensiune și nivel de adnotare este una din cele mai importante etape în aducerea tehnologiilor și resurselor lingvistice pentru limba română la un nivel comparativ cu limbile vest-europene. În cadrul proiectelor inițiate în această direcție poate fi inclus și LT4eL¹ (Monachesi et. al., 2006), proiect susținut de Comunitatea Europeană prin departamentul *Information Society and Media Directorate, Learning and Cultural Heritage*. În acest proiect sunt implicate universități reprezentând 9 limbi europene (engleza, germana, olandeza, portugheza, poloneza, ceha, bulgara, malteza și româna). Limba română este reprezentată în proiect prin autori, membri ai grupului de cercetare în Tehnologiile Limbajului Uman² din cadrul Facultății de Informatică). Scopul principal al proiectului este dezvoltarea de instrumente și resurse lingvistice concepute pentru folosirea în învățarea asistată de calculator (eLearning).

O primă etapă a LT4eL, pentru care responsabilă a fost echipa română, a constat în conceperea și realizarea unui mediu de adăugare, acces și actualizare a tuturor corpusurilor, resurselor lingvistice și sistemelor de prelucrare ce urmau a fi dezvoltate în cadrul proiectului. În acest scop s-a creat un portal³ prin intermediul căruia au fost colectate de la începutul proiectului (decembrie 2005) resurse, totalizând aproape 9

¹ Eng: Language Technologies for E-Learning, pagina principală la <http://www.lt4el.eu>.

² <http://consilr.info.uaic.ro/research/>

³ http://consilr.info.uaic.ro/uploads_lt4el/

milioane de cuvinte, și 30 aplicații dezvoltate integral în cadrul proiectului sau adaptate necesităților proiectului.

În secțiunea a doua a acestei lucrări va fi prezentat procesul general de prelucrare a corpusului românesc, pentru ca în secțiunile 3 și 4 să fie descrise mai pe larg două dintre etape, semnificative prin noutatea modulelor de procesare dezvoltate. Secțiunea a 5-a conține o descriere a etapelor următoare ale proiectului, precum și a resurselor și instrumentelor ce urmează să fie create pentru limba română.

2. Etapele de prelucrare a corpusului românesc

Etapa inițială de formare a corpusului românesc a constat în colectarea de documente din 14 domenii convenite în cadrul proiectului, domenii ce țin de predarea informaticii, pedagogie și organizarea studiului universitar. Un prim nivel de adnotare propus a fost XML⁴ cu marcate de formatare a documentului, format definit de un DTD⁵ general. Numai marcasele de formatare care pot facilita prelucrări ulterioare, cum ar fi marcarea cuvintelor cheie, au fost păstrate. După transformarea întregului corpus în format XML, etapa a doua a fost cea de adnotare lingvistică în care s-a realizat: segmentarea în unități lexicale, propoziții, fraze și paragrafe, marcarea informațiilor morfo-sintactice, a formelor de bază (leme) ale cuvintelor flexionate, marcarea grupurilor nominale.

A treia etapă a constat în marcarea manuală a cuvintelor cheie și a definițiilor din corpus, relevante pentru domeniul general al proiectului. Drept cuvinte cheie au fost marcate cuvintele și expresiile considerate de adnotator ca fiind relevante în contextul conținutului și al scopului lucrării respective. Toate aceste cuvinte cheie au fost traduse din celelalte 8 limbi în engleză, centralizate și sortate, obținându-se astfel colecția lexicală, ca prim pas în construirea unei ontologii generale pentru domeniile de interes ale LT4eL. Formatul convenit pentru realizarea ontologiei a fost DOLCE⁶. Într-o etapă ulterioară, forma finală a acestei ontologii și lexiconul aferent vor fi mapate integral la Princeton WordNet⁷, lucru ce va permite integrarea ei cu alte ontologii de domenii, cum ar fi SUMO, dar și utilizarea ei în cadrul altor proiecte.

În prezent, corpusul proiectului conține peste 600.000 cuvinte, în variantele de adnotare descrise mai sus, până la marcasele de cuvinte cheie și definiții.

3. Conversia corpusului la formatul XML

Documentele primare ale proiectului au trebuit aduse de la formatele inițiale la un format unitar XML în codificare UTF-8 (ce permite notarea diacriticelor în toate limbile proiectului), înainte de abordarea etapelor adnotării lingvistice și a celei semantice. Transformarea întregului corpus LT4eL la un format XML standard, definit de o specificație DTD, s-a făcut, în parte, utilizând o serie de convertoare disponibile, în parte, implementând trei noi convertoare. Formatul html s-a ales ca nivel intermediar în

⁴ Extensible Markup Language (<http://www.w3.org/XML/>)

⁵ Document Type Definition

⁶ <http://www.loa-cnr.it/DOLCE.html>

⁷ <http://wordnet.princeton.edu/>

transformarea între formatele inițiale doc și pdf și formatul final XML. Motivul a fost că html păstrează în bună măsură indicațiile de formatare ale documentelor inițiale.

Prezentăm mai jos, pe scurt, doar câteva din problemele întâlnite la conversia corpusului din formatul html la formatul XML ce respectă restricțiile convenite în cadrul LT4eL:

- Unele cuvinte consecutive apar alipite: cauza trebuie pusă fie pe seama conversiei la html (text pe două coloane, probleme date de justify, cuvinte despărțite de newline), fie pe seama convertorului la xml. Cuvintele ce nu sunt despărțite textual, ci doar prin interpretarea etichetelor html, sunt principalele cauza ale apariției acestei erori. Dacă textul are numeroase elemente de formatare (fonturi multiple, culori, siluri), aceste probleme sunt frecvente.
- Unele cuvinte apar cu spații în interior: în general, aceasta este o consecință a conversiei de la pdf la html și apare în textele pe două coloane, cu *justify* sau consecință a despărțirii în silabe. Rezolvarea acestei probleme necesită resurse lingvistice pentru limba respectivă, pentru evitarea creării de erori opuse, adică apariția unor cuvinte alipite.
- Unele caractere au coduri greșite în xml: fie din cauza html-ului de origine unde caracterul este notat printr-un alt cod decât UTF-8, fie pentru că un singur caracter are uneori mai multe codificări (de obicei UTF-8 și UTF-16), sau pentru că același caracter are coduri diferite în seturi diferite de caractere UTF-8 (cum ar fi Latin-1 și Arabic). O soluție optimă, cum ar fi crearea unui instrument automat de conversie din orice codare la UTF-8, pare a fi imposibil de realizat (de exemplu UTF-16 *unpaired*⁸ nu poate fi convertit la UTF-8).
- Uneori apar spații multiple în loc de unul singur. Corectarea prin înlocuirea spațiilor multiple cu unul singur poate uneori strica formatarea documentului.
- Uneori convertorul păstrează în XML atributele unor etichete eliminate: problema apare în special la etichetele ce apar doar deschise și la atributele ale căror valoare apare fără ghilimele. Soluția are în vedere, în principal, rezolvarea acestor două situații.

Pe lângă problemele menționate anterior mai apar câteva, semnificative, ce sunt datorate procesului de conversie. În primul rând, conversia corectă a textelor puternic formate (tabele, formule, imagini grupate) este o problemă uneori extrem de dificil de rezolvat, lucru afirmat și în documentațiile convertoarelor de firmă (*MS Word, Adobe, pdf2html*). O soluție viabilă pentru acest tip de documente, a căror formatare poate fi uneori relevantă pentru adnotarea lingvistică, ar impune o prelucrare adițională, în vederea transformării obiectelor problematice într-un format mai ușor de prelucrat. Acest lucru poate implica și un efort manual din partea adnotatorului, posibil facilitat de dezvoltarea unui mediu vizual de editare. Există de asemenea documente doc și pdf care provin din conversia unor documente scanate. Acestea vor rămâne practic imposibil de convertit la XML cu păstrarea integrală a conținutului, chiar și după utilizarea unui software de tip OCR (Optical Character Recognizer) performant. O altă

⁸ <http://unicode.org/unicode/faq/>

problemă ce implică o prelucrarea adițională este cea a documentelor ce conțin, în original, secvențe de adnotare html/XML a textului (ca cele din manuale, de exemplu) și care se pot confunda cu metadatele XML ale formatului final.

4. Adnotarea lingvistică

Adnotarea lingvistică (îmbogățirea corpusului cu informații sintactice și morfologice) urmează etapei de transformare a corpusului în format XML. Într-o primă etapă s-a făcut o evaluare a instrumentelor de prelucrare disponibile în vederea preluării în proiect a unora dintre ele, apoi s-a luat decizia implementării unora noi. În final, pentru adnotarea lingvistică au fost folosite:

- *tokenizator* (marcator de unități lexicale de bază), dezvoltat de echipa UAIC;
- *POS-tagger* (adnotator morfo-sintactic), adaptat după o implementare ICIA⁹;
- *lemmatizer* (marcator de rădăcini morfologice neflexionate), realizat la ICIA;
- *NP-chunker* (marcator de grupuri nominale), dezvoltat de echipa UAIC, utilizând un corpus adnotat manual pentru a genera un set de reguli, ce au fost apoi revizuite în parte înainte de a fi utilizate de marcatorul de grupuri nominale.

Cele patru module de procesare de mai sus obțin rezultate foarte bune pentru limba română. *Tokenizator*-ul și *POS-tagger*-ul obțin scoruri *F-measure*¹⁰ de aproximativ 98% (Tufiș și Dragomirescu, 2004), *Lemmatizer*-ul obține un scor *F-measure* de aproximativ 95%, iar *NP-chunker*-ul aproximativ 75%. Evaluările pentru instrumentele ICIA au fost preluate din documentația aplicațiilor, iar pentru cele dezvoltate la UAIC au fost calculate automat utilizând un corpus standard adnotat manual. Acest corpus (ca și o primă variantă a modulelor respective) a fost dezvoltat în cadrul unui proiect anterior, ce a presupus dezvoltarea unui sistem de întrebare-răspuns (Iftene et al, 2006) în context bilingv română-engleză. Corpusul utilizat ca standard în evaluare este romanul '1984' de G.Orwell în varianta tradusă în limba română și cuprinde peste 100.000 unități lexicale.

Pentru combinarea și simplificarea adnotărilor și aducerea la formatul proiectului a fost folosită o aplicație dezvoltată de echipa UAIC, ce va sta la baza unui viitor sistem de adnotare automată ce va combina module de prelucrare în scopul obținerii automate a unor formate complexe.

5. Etape viitoare

Etapele viitoare ale proiectului LT4eL constau, în primul rând, în implementarea și antrenarea unui sistem de recunoaștere automată a cuvintelor cheie și a definițiilor. Cuvintele cheie și definițiile din texte vor fi principalele elemente de legătură între obiectele de învățare (eng: *Learning Objects*) și vor permite integrarea și combinarea acestor resurse în scopul creării unui mediu dinamic de generare de materiale pentru învățarea asistată de calculator. Pentru recunoașterea cuvintelor cheie va fi folosită o

⁹ Institutul de Cercetări în Inteligență Artificială din cadrul Academiei Române – București.

¹⁰ Calculat ca $2 * P * R / (P + R)$, unde P (precizia) = numărul de obiecte corect identificate de program raportat la numărul de obiecte identificate de program și R (recall) = numărul de obiecte corect identificate de program raportat la numărul de obiecte existente.

adaptare a algoritmului *tf-idf* (Salton și Buckley, 1988), cu antrenare pe corpusul adnotat manual la cuvinte cheie existent (aproape 25 mii de cuvinte cheie marcate manual, din care 7772 în corpusul românesc). Pentru recunoașterea automată a definițiilor s-a elaborat o primă variantă a unei gramatici ce utilizează informațiile morfo-sintactice și de formatare existente în corpus pentru a identifica și marca granițele definițiilor din text.

Îmbogățirea corpusului românesc LT4eL este de asemenea una din direcțiile de lucru, scopul echipei românești fiind acela de a utiliza tehnologiile și formalismele dezvoltate în cadrul LT4eL pentru a colecta și prelucra un corpus mult mai mare decât cel necesar pentru proiect, ce ar putea fi folosit mai târziu în cadrul altor proiecte de cercetare în care este implicată limba română, cum ar fi dezvoltarea unui sistem de traducere automată, a unui rezumator, a unui parser de discurs etc.

Transformarea ierarhiei de formate într-un mediu de prelucrare automată a documentelor este un alt scop avut în vedere de echipa românească, în afara cerințelor stricte ale proiectului. Proiectarea unui astfel de mediu, cu posibilitatea îmbogățirii sale ulterioare cu noi module de prelucrare și resurse lingvistice poate facilita dezvoltarea unor sisteme complexe de procesare lingvistică, cu aplicații ce-și pot găsi utilitatea și în prelucrările dedicate altor limbi decât limbii române. Formalismele teoretice pentru dezvoltarea acestui mediu sunt în curs de îmbunătățire, după ce a fost elaborată o primă propunere (Cristea et. al, 2006). Se speră ca până la terminarea proiectului LT4eL să fie finalizată o variantă a acestui mediu care să acopere cel puțin procesările necesare LT4eL.

Referințe bibliografice

- Iftene A., Pistol I., Trandabăț D., Pușcașu G., Forăscu C., Cristea D. (2006) Sisteme de Întrebare Răspuns pentru Limba Română. În acest volum.
- Monachesi P., Lemnitzer L., Simov K. (2006). Language Technology for eLearning to appear in *Proceedings of EC-TEL 2006*, Springer LNCS (<http://www.ectel06.org/index.html>).
- Tufiș D., Dragomirescu L. (2004). Tiered Tagging Revisited. In *Proceedings of the 4th LREC Conference*, Lisabona, 2004, pp. 39-42.
- Salton G., Buckley C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management* 24(5): 513–523.
- Cristea D., Forăscu C., Pistol I. (2006). Requirements-Driven Automatic Configuration of Natural Language Applications, *Proceedings of the 3rd International Workshop on Natural Language Understanding and Cognitive Science NLUCS-2006*, ICEIS 2005, 23 - 27 May 2006, Paphos, Cyprus.

SISTEM DE INSTRUIRE ASISTATĂ DE CALCULATOR PENTRU MORFOLOGIA LIMBII ROMÂNE

ELENA BOIAN, CONSTANTIN CIUBOTARU, SVETLANA COJOCARU,
GALINA MAGARIU, TATIANA VERLAN, IURI ROGOJIN

Institutul de Matematică și Informatică, Academia de Științe a Moldovei

{lena, chebotar, sveta, gmagariu, tverlan, rogozhin}@math.md

Rezumat

În lucrare¹ este descrisă concepția de elaborare a sistemului de instruire asistată de calculator pentru limba română. Pentru implementarea acestei concepții s-au utilizat mijloacele sistemului Claroline, o platformă Open Source. Concepția propusă presupune studierea morfologiei limbii române pentru trei categorii de utilizatori cu diferite niveluri de pregătire.

1. Introducere

Dezvoltarea continuă a tehnologiilor informaționale în sfera educației se manifestă printr-o vastă orientare spre utilizarea produselor program deschise, tendință manifestată și în domeniul sistemelor de instruire asistată de calculator (SIAC). În lucrare sunt expuse principiile și particularitățile de proiectare și realizare a unui SIAC pentru morfologia limbii române, sistem bazat pe o arhitectură ierarhică orientată pe categorii de utilizatori și integrată cu resursele lingvistice reutilizabile plasate pe Internet.

Sistemul de instruire presupune nu numai cursul materialului pentru studiere (inclusiv și exerciții), ci și mijloace de prezentare și expunere a materialului de către profesor, modalitatea de organizare a comunicării între student și profesor, tehnicile de însușire a materialului de către student.

Elaboratorii SIAC, fiind în postură de profesor, selectează materialul teoretic și exercițiile necesare, ținând cont de metodică și principiile de prezentare a obiectului și de nivelul de pregătire a studentului; în funcție de student, elaborează o modalitate optimă pentru însușirea materialului propus. Se propune o secvență de activități – alternarea materialului teoretic și a exercițiilor de autotestare, a testelor controlate prin intermediul sistemului de instruire, de asemenea, a testelor gestionate de profesor în procesul de comunicare cu studentul (în special, pentru exercițiile creative).

Așadar, un sistem de instruire nu este numai o simplă colecție și nu numai o colecție de material teoretic ce ține de disciplina pentru instruire, cum s-ar părea la prima vedere. Acesta este o noțiune mai vastă și mai profundă și, în cazul nostru, reprezintă un sistem consistent, multilateral și valoros.

SIAC propus presupune elaborarea:

¹ Lucrarea este efectuată în cadrul proiectului RoLTech INTAS Ref. Nr. 05-104-7633

- materialului teoretic propriu zis,
- exercițiilor pentru formarea deprinderilor de utilizare a materialului teoretic,
- metodicii de predare a materialului de către profesor,
- modalității de comunicare între student și profesor,
- metodelor de asimilare a materialului de către student.

2. Concepția de construire a cursurilor pentru studierea limbii române

În baza discuțiilor cu persoanele nevorbitoare de limba română, care au frecventat cursurile de studiere a limbii române, a cercetării literaturii de specialitate existente în domeniul gramaticii (Acsan et al., 2004), (Bărbuță, 1998), (Bărbuță et al., 2000, 2003), (Berteza, 1996), (Cazacu & Vrabie, 2006), (Cruceru & Teodorescu, 2005), (Dumeniuc & Matcaș, 1989), (Gramatica limbii române, 1963), (Irimia, 1997), (Irimia, 2004), (Limba română contemporană, 1999), (Molan et al., 1995, 1996), (Nastasenco, 1996), (Pop, 2000), (Zaiuncikovski & Repina, 1989) și a metodicii de predare a limbii române (Consiliul Național pentru Curriculum, 2002), (Ionescu, 2001), (Ivănuș, 1997), (Nuță, 2000), a fost elaborată o concepție de construire a cursurilor asistate de studiere a limbii române.

Pentru implementarea acestei concepții s-au utilizat mijloacele sistemului Claroline².

Concepția propusă presupune studierea morfologiei limbii române pentru trei categorii de utilizatori:

- **prima categorie** – persoanele care nu cunosc gramatica și au un vocabular sărac;
- **categoria a doua** – persoanele care înțeleg limba vorbită, au cunoștințe gramaticale nesistematizate;
- **categoria a treia** – persoanele care cunosc limba română, dar vor să-și aprofundeze cunoștințele referitor la gramatică, particularități lingvistice, derivare a cuvintelor etc.

2.1. Prima categorie – nivelul 1

Scopurile lecțiilor de nivelul I:

- dezvoltarea limbii vorbite, începând cu fraze elementare și expresii uzuale;
- îmbogățirea treptată a vocabularului și complicarea frazelor;
- studierea paralelă a categoriilor gramaticale la nivel elementar. Aceste categorii sunt predate treptat în funcție de necesitatea utilizării lor în vorbire.

Așadar, studierea gramaticii nu este prioritară și nu reprezintă scopul principal al lecțiilor de nivelul I.

În concordanță cu scopurile lecțiilor a fost elaborată structura lecțiilor de nivelul I. O lecție conține următoarele compartimente:

² <http://www.claroline.net>

- *cuprinsul*, care ne permite prin referințele interactive să accesăm compartimentul selectat în cadrul lecției;
- *vocabularul*, care conține cuvintele noi din cadrul lecției. Fiecare cuvânt este însoțit de sugestii cu privire la: informația despre partea de vorbire, (de exemplu, pentru substantiv se indică genul, formele derivate pentru numărul plural al substantivelor, pentru verb – conjugarea, sufixul verbului etc.); traducerea în limba maternă a studentului; exemple de utilizare a cuvântului indicat în fraze uzuale. De asemenea, se specifică două referințe interactive: la un fișier sonor cu pronunția corectă a cuvântului și la resursele din Internet cu o informație mai amplă despre acest cuvânt;
- *gramatica*, ce conține descrierea regulilor gramaticale la un nivel elementar pentru lecția prezentată. Această informație trebuie să fie în limba română și, la cererea studentului, să fie accesibilă în limba vorbită de student (de exemplu, în limba rusă, engleză etc.);
- *textul*, care este alcătuit în baza cuvintelor și a gramaticii din lecțiile precedente cu utilizarea cuvintelor și gramaticii din lecția curentă. Fiecare propoziție din text este însoțită cu două referințe interactive: la fișierul sonor cu pronunția corectă a frazei și la fișierul cu traducerea în limba studentului;
- *modele comunicative, dicționar de contexte minime, proverbe, fragmente din folclor* etc.;
- *exerciții* elaborate în baza textului. Ele au ca scop întărirea deprinderilor de utilizare a noului vocabular și a gramaticii în cadrul lecției. Pentru fiecare lecție se propune elaborarea a trei tipuri de exerciții:
 - pentru autotestare, care se propun imediat după însușirea materialului teoretic (de exemplu, un fișier aparte cu sarcinile exercițiilor și răspunsurile corecte).
 - elaborate de profesor cu ajutorul mijloacelor sistemului Claroline;
 - care utilizează „corespondența” profesorului cu studentul.

În acest mod, în cadrul acestor lecții se pune accent atât pe pronunțarea cuvintelor și a frazelor uzuale, cât și pe îmbogățirea vocabularului. În acest scop se utilizează procedee interactive, asistate de mijloace multimedia.

2.2. Categoria a doua – nivelul 2

Scopul de bază al cursului de nivelul II este studierea sistematizată a gramaticii, în special, a morfologiei limbii române. Se presupune că informația despre categoriile gramaticale va fi consistentă, adică, dacă se studiază substantivul, pe parcursul a câtorva lecții consecutive se descrie informația despre toate categoriile gramaticale ce țin de substantiv: definiția substantivului, gen, număr, caz, declinarea substantivului și alte particularități ale lui. Ca și în cadrul cursului de nivelul I, se va propune material teoretic, exemple de utilizare și exerciții practice. Titlurile lecțiilor se vor specifica în conformitate cu materialul gramatical expus. Materialul teoretic se prezintă cu utilizarea categoriilor și a terminologiei uzuale. Se presupune că studentul cunoaște această terminologie în virtutea nivelului său de pregătire. Pentru termenii de bază utilizați la lecția curentă se vor folosi sugestii pentru reamintirea definițiilor acestor termeni.

Pentru o sistematizare și o structurizare efectivă a materialului și, de asemenea, pentru o înțelegere mai bună, se propun scheme și diagrame. Prezentarea regulilor gramaticale se

va explica prin utilizarea tabelelor care vor ajuta la asimilarea completă și eficientă a materialul expus pentru instruire.

2.3. *Categoria a treia – nivelul 3*

Nivelul III presupune că, în principiu, studentul știe bine limba română, cunoaște regulile normative și utilizarea uzuală a cuvintelor. La lecțiile de nivelul III o atenție sporită se va acorda devierilor de la reguli, particularităților și excepțiilor de la situațiile morfologice normative (regulate). De asemenea pot fi propuse anumite reguli (situații) nenormative, momente speciale de utilizare a unor cuvinte și a regulilor gramaticale. Se va explica sensul semantic și utilizarea neologismelor. Limbajul de prezentare și expunere a materialului se va deosebi de cel utilizat pentru descrierea lecțiilor de nivelul I și II. Pentru cursul de nivelul III, acest limbaj este mai complicat, mai rafinat, la un nivel științific avansat. În calitate de exemple se vor utiliza citate din operele scriitorilor clasici români și folclor. Ca și în cazul lecțiilor de nivelul II, materialul se va expune cu utilizarea categoriilor și terminologiei indicând sugestii pentru definirea termenilor.

3. *Concluzii*

În comparație cu alte cursuri publicate în cărți, cursul de lecții propus va putea fi accesat pe Internet de orice doritor de a învăța limba română și va conține multe posibilități specifice aplicațiilor Web. Iar în comparație cu materialele „pasive” în formă electronică, lecțiile se vor deosebi prin utilizarea interacțiunii studentului și a profesorului (comunicării cu profesorul și cu colegii din grup, schimbului de informație între profesor și student) și a mijloacelor de evaluare a cunoștințelor.

Acest ciclu de lecții nu-l va înlocui pe profesor, ci îl va ajuta în activitatea sa didactică cu material adițional în predarea lecțiilor ce țin de morfologia limbii române. Pentru studenții ce nu cunosc limba, cursul de lecții va deveni o sursă adițională în procesul de studiere a limbii române, iar pentru alte categorii de utilizatori va oferi posibilități de aprofundare a cunoștințelor în morfologia limbii române.

Referințe bibliografice

- Acsan, A., Cojocaru-Zavadschi, A., Cucu, L. (2004). *Limba care ne unește. Caiet de exerciții. Nivelul 1*. Departamentul Relații Interetnice, Programul Națiunilor Unite pentru Dezvoltare (PNUD). Chișinău.
- Bărbuță, I. (1998). *Gramatica limbii române. Scurt îndrumar*. Academia de Științe a Republicii Moldova, Inst. de Lingvistică. Chișinău, Litera, 152p. (în l. rusă)
- Bărbuță, I., Callo, T., Cojocaru-Zavadschi, A., Constantinovici, E., Cucu, L. (2003). *Limba care ne unește. Manual. Nivelul 1*. Departamentul Relații Interetnice, PNUD. Chișinău, (<http://cnt.dnt.md/undp/>).
- Bărbuță, I., Cicală, A., Constantinovici, E., Cotelnic, T., Dîrul A. (2000). *Gramatica uzuală a limbii române*. AȘM, Inst. de Lingvistică. Litera, Chișinău. 326 p.

- Berteau, M. (1996). *Gramatica explicativă a limbii române plus Vocabular Ortografie Ortoepie*. Ed. a IV, revăz. și adăug., Ed. Venus. București, Î.E.P. Știința. Chișinău.
- Cazacu, T., Vrabie, D. (2006) *Româna: eficient și atractiv. Gramatica limbii române în scheme și tabele*. Integritas, Chișinău. 52 p.
- Consiliul Național pentru Curriculum. (2002). *Ghid metodologic pentru aplicarea programelor de limba și literatura română – Învățământ primar și gimnazial*. C.N.C., București, România.
- Cruțeru, C., Teodorescu, V. (2005). *Gramatica limbii române*. Editura 100+1 GRAMAR, București, 151 p.
- Dumeniuc, I.Z. Matcaș, N.G. (1989). *Limba Moldovenească. Manual pentru autodidacți*. Chișinău, Lumina, 362 p. (în l. rusă)
- Gramatica limbii române*. (1963). Academia Republicii Române, v1, v2, București.
- Ionescu, M. (2001). *Didactica modernă*, Radu I. coord., Dacia, Cluj, România.
- Irimia, D. (1997). *Morfo-sintaxa verbului românesc*. Universitatea „A.I.Cuza”, Iași, România, 410 p.
- Irimia, D. (2004). *Gramatica limbii române*. Polirom, ed. II, Iași, România.
- Ivănuș, D. (1997). *Metodica predării limbii și literaturii române în gimnaziu și liceu*. Avrămeanca, Craiova, România.
- Limba română contemporană*. (1999). Îndrumar pentru persoane nevorbitoare de limba română. Chișinău, Litera, 336p. (în l. rusă)
- Molan, V., Părvulescu, L., Teodorescu, I. (1995). *Limba noastră-i o comoară. Exerciții de limbă română pentru ciclul primar*. Ediția a II-a. Ed. Petrion. București.
- Molan, V., Teodorescu, I., Dobrin, E. (1996). *Gramatică, ortografie și punctuație pentru toți copiii*. (Clasele II-IV). Editura Petrion. București. Chișinău.
- Nastasenco, O. (1996). *Gramatica limbii române în tabele*. Chișinău. VIRT. (în l. rusă)
- Nuță, S. (2000). *Metodica predării limbii române la clasele primare*. Aramis, București, România.
- Pop, L. (2000). *Româna cu sau fără profesor*. Ed. a IV, revăz. și adăug. Echinoc, Cluj.
- Zaiuncikovski, I.P., Repina, T.A. (1989). *Limba Română. Curs superior. Manual pentru anii de studii II-III ai facultăților filologice ale universităților*. Moscova. Visșaiia școla, (în l. rusă)

Capitolul 4

Modelare lingvistică

STRUCTURA GRUPULUI VERBAL, PREDICAȚIA LEXICALĂ ȘI REPREZENTAREA LOGICĂ A PREDICATULUI ÎN LIMBA ROMÂNĂ

NECULAI CURTEANU, DIANA TRANDABĂȚ^{1,2}, ALEX MORUZ^{1,2}

¹*Institutul de Informatică Teoretică, Academia Română, Filiala Iași*

²*Facultatea de Informatică, Universitatea "Al.I.Cuza" Iași*

{curteanu, dtrandabat, mmoruz}@iit.tuiasi.ro

Rezumat

Articolul include o prezentare elementară a *teoriei proiecțiilor funcționale* FX-bar, introduce *predicația lexicală* în locul predicației clasice, pune în evidență principalele *substructuri (subgrupuri) ale grupului verbal*, și propune o *reprezentare unitară* a predicatului verbal și nominal în logică intensională / extensională.

1. Introducere

Lucrarea își propune să prezinte succint câteva rezultate legate de substructurile sintactice / semantice ale *grupului verbal* românesc (*verbal group*, VG) (Monachesi, 2005), (Barbu, 1999), (Dobrovie-Sorin, 1994), pornind de la instrumente și argumente cunoscute în literatură, cărora li se adaugă mecanismele *teoriei FX-bar a proiecțiilor funcționale* (Curteanu, 2003-2004, 2005). Sunt schițate probleme și soluții în cadrul teoriilor sintactice ale VG, cu accent pe semantica lexicală, interesul major fiind orientat către *parsarea VG*, *subgrupurile verbale* ale VG, *proiecțiile FX-bar* (directe și inverse) ale VG, definirea *predicației lexicale* (în locul celei clasice), și *reprezentarea unitară*, în logică intensională / extensională, a predicatului verbal și nominal.

2. Proiecții FX-bar directe și inverse, de nivel local și global

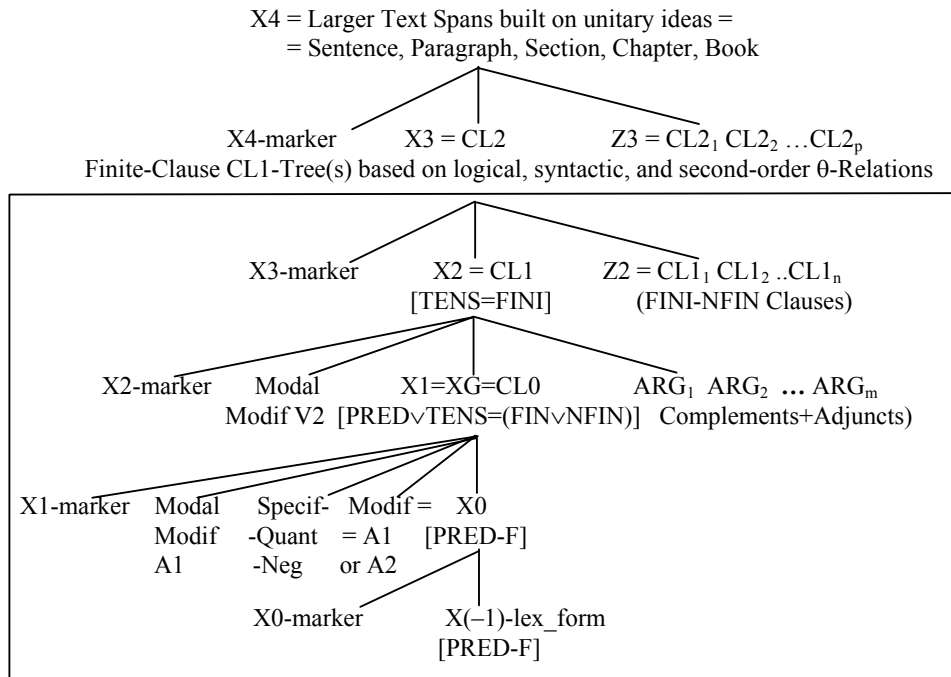
În (Curteanu, 2005) a fost (re)introdusă schema generală FX-bar din Fig. 1, ce folosește clasele de marcheri SCD și ierarhia de tip graf a acestor clase, un instrument esențial pentru a reprezenta structurile sintactico-semantice de nivel propozițional (clauzal) și pentru a stabili dependențele dintre ele. Partea inferioară a schemei, încadrată în chenar, reprezintă nivelul local, intra-clauzal, iar partea superioară indică nivelul global, inter-clauzal și de discurs, bazat pe relații retorice RST (Mann & Thompson, 1988).

Să precizăm că în abordarea noastră folosim termenul de „*finit*” pentru toate *formele verbale* ce corespund unui *mod personal*, în timp ce termenul „*nefinit*” este atribuit formelor verbale ce corespund „*modurilor*” nepersonale, cunoscute și ca *forme absolute* ale verbului (*infinitiv*, *gerunziu*, *participiu*, *supin*). În acest sens, clauza (sau propoziția) al cărei predicat are ca nucleu semantic un verb în formă predicativ/verbală (*finită*), o vom numi *clauză finită*. *Clauza nefinită* corespunde formelor verbale predicativ/nepersonale (*nefinite*), ale categoriilor lexicale majore V (verb), N

(substantiv), și A (Adjectiv-Adverb), care posedă trăsătura de *predicaționalitate* (*deverbalitate*). Parsarea predicatului verbal sau nominal revine la determinarea *grupului verbal* (Verbal Group, VG). Facem observația că determinarea VG (finit) de natură copulativă (nepredicativ/personală), notat TASG – *Tense Auxiliary SubGroup* în (Curteanu & Trandabăț, 2006), este o etapă esențială în parsarea predicatului nominal.

Interpretarea schemei FX-bar se face *bottom-up* pentru trecerea la o sintagmă mai complexă, dar și *top-down* pentru stabilirea dependențelor și pentru adăugarea de elemente lexicale și/sau sintagme de același nivel. Nivelul X0 este cel lexical-textual, din care se coboară la nivelul intrărilor de lexicon X(-1), unde categoriilor *lexicale majore* N (Noun), V (Verb), A (Adjectiv-Adverb) li se poate atribui trăsătura de *predicaționalitate* (sau *deverbalitate*), notată PRED-F. Elementele de lexicon sunt (FX-bar) proiectate pe nivelul X0 după ce li se aplică X0-marcherul de *flexionare morfologică*. Prin adăugarea *on-line* de elemente lexicale noi se constituie (incremental) structurile sintagmatice de nivel X1, având ca nucleu sintactic categoriile lexicale N, V, A: *grupul nominal, verbal, adjectival-adverbial*, notate respectiv NG, VG, AG. FX-bar proiecția lui X0 pe nivelul X1 (sau XG) poate fi însoțită de modificator (verb modal), specificator (negație sau cuantificator), sau auxiliare. V1 (sau VG) poate fi finit (formă predicativ/personală) sau nu, proiectându-se mai departe în schema FX-bar pe nivelul X2, alături de posibili modificatori modali și argumente directe și indirecte (adjuncți), formând propoziția (clauza finită) pe nivelul X3. Împreună cu alte clauze, se trece la proiecția inter-clauzală. Direcția ascendentă ilustrează mecanismul *proiecției* FX-bar *directe*, incluzând FX-bar proiecția VG în clauză.

Figura 1: FX-bar schema pentru structuri de nivel local (clauzal) și global (discurs)



Direcția de *proiecție* FX-bar *descendentă*, *top-down* în schema din Fig. 1, reprezintă *proiecția inversă*, care include *proiecția* FX-bar *inversă* a VG. Mecanismul este folosit pentru a stabili în mod corect proiecția FX-bar directă a unui VG în clauză, prin

specificarea de restricții fonologice, morfologice, sintactice și semantice asupra nucleului semantic (eventual predicțional) și asupra *theta*-argumentelor *directe* ale VG.

2.1. Predicația clasică și predicația lexicală

Pentru orice intrare N, V, A, sunt reprezentate la nivelul X(-1) de lexicon: (i) trăsătura de *predicționalitate*; (ii) *valența* (*aritatea* și *tipul* – ordinul logic al – argumentelor); (iii) *diateză semantică*; (iv) *restricții* sintactico-semantice ale argumentelor (directe) ale nucleului semantic al VG. *Predicația clasică*, cunoscută ca fiind perechea (Subiect-Gram, Predicat), poate fi considerată doar una din fațetele VG al cărui nucleu semantic poartă *trăsătura de predicționalitate* PRED-F (Curteanu, 2003-2004). Această pereche corespunde fie rolului tematic (*theta*-rol) de „Actor” / „Actant”, fie de „Obiect” / „Pacient”, în funcție de *diateza activă*, respectiv *pasivă* a VG. Un al treilea argument de natură *theta*-semantică poate fi „Beneficiar” / „Adresant” al predicației din VG. Predicația clasică (Subiect-Gram, Predicat) este rescrisă ca fiind perechea (SUBJ_{Obliqueness} = 0, PREDF_verb) corespunzând *theta*-rolului de „Actor” / „Actant” sau „Obiect” / „Pacient” în lista SUBCAT a argumentelor directe cerute de valența nucleului semantic al VG (Pollard&Sag; 1994). Trebuie specificat că în mod normal există cel puțin două liste SUBCAT: SUBCAT_{oblic_order}, ce conține argumentele sintactice ale verbului predicțional PREDF_verb în ordinea crescătoare a oblicității, și SUBCAT_{theta_order} care conține argumentele în ordinea sistemică a argumentelor *theta*-semantice. De obicei, cele două liste SUBCAT sunt identice numai pentru *diateza activă*. La predicația clasică (Subiect-Gram, Predicat) se adaugă predicațiile din Fig. 2, justificate de comportamentul lor perfect similar față de nucleul semantic al argumentelor directe:

$$\begin{array}{l}
 (\text{SUBJ}_{\text{Obliqueness} = 0}, \text{PREDF_Verb} \left[\begin{array}{l} [\text{VG Tense_Aspect}] \\ \text{Semantic_Diathesis}(\text{SUBJ}, \text{OBJD}, \text{OBJI}) \\ = (\theta(\text{SUBJ}), \theta(\text{OBJD}), \theta(\text{OBJI})) \\ \text{Agreement}(\text{SUBJ}, \text{Inflection_VG}) \end{array} \right] \\
 (\text{OBJD}_{\text{Obliqueness} = 1}, \text{PREDF_Verb} \left[\begin{array}{l} [\text{VG Tense_Aspect}] \\ \text{Semantic_Diathesis}(\text{SUBJ}, \text{OBJD}, \text{OBJI}) \\ = (\theta(\text{SUBJ}), \theta(\text{OBJD}), \theta(\text{OBJI})) \\ \text{Agreement}(\text{OBJD}, \text{CliticOBJD_VG}) \end{array} \right] \\
 (\text{OBJI}_{\text{Obliqueness} = 2}, \text{PREDF_Verb} \left[\begin{array}{l} [\text{VG Tense_Aspect}] \\ \text{Semantic_Diathesis}(\text{SUBJ}, \text{OBJD}, \text{OBJI}) \\ = (\theta(\text{SUBJ}), \theta(\text{OBJD}), \theta(\text{OBJI})) \\ \text{Agreement}(\text{OBJI}, \text{CliticOBJI_VG}) \end{array} \right]
 \end{array}$$

Figura 2: Predicația lexicală extinde predicația clasică pentru toate argumentele directe

Aceste noi predicații „tradiționale” se bazează pe trăsătura de *predicționalitate* (deverbalitate) PRED-F, ce este atribuită la nivel lexical în nucleul VG. Problema rolului ’special’ al subiectului în lista SUBCAT a argumentelor este rezolvată în structura propusă prin demonstrarea similarității relației *Subiect – PREDF_verb* cu relațiile *Complement (Direct / Indirect) – PREDF_verb* și prin folosirea funcției de *diateză semantică* (Irimia; 1997). *Acordul* dintre VG cu nucleu predicțional și

argumentele sale directe, fie că se face prin *flexionare morfologică* (pentru subiectul gramatical SUBJ) sau prin *pronume clitice* (pentru OBJD și OBJI), este un alt argument major ce susține *echivalența* elementelor constitutive ale *predicației lexicale*.

2.2. *Diatezele formal-sintactică și lexical-semantică*

Diatezele *formal-sintactică* și *lexical-semantică* sunt interpretări semantice diferite ale aceluiași șir de elemente lexicale (de suprafață) din VG. Se recunoaște una din cele 3 diateze sintactice, se coboară la descrierile lexical-semantică din lexicon, iar nucleul semantic selectat se proiectează FX-bar într-una din cele 6 diateze semantice ale VG. Funcția *Semantic_Diathesis(Actor, Patient, Addressee)* primește ca intrare diateza sintactică a VG, reprezentată de lista de (cel mult 3) argumente directe, în ordinea oblicității $SUBCAT_{oblic_order}$ și returnează diateza semantică împreună cu argumentele în ordinea $SUBCAT_{theta_order}$. Această soluție obligă *actorul*-subiect și subiectul gramatical (cel mai puțin oblic element) să ocupe fiecare poziția corectă. Transformarea argumentelor din diateza sintactică în cea semantică este stabilită de un tabel de corespondență prezentat în (Curteanu & Tradabăț, 2006).

3. *Substructurile grupului verbal VG*

3.1. *Natura gramaticală a componentelor grupului verbal*

Nucleul semantic verbal antrenează în jurul său intensificatori, modificatori, particule clitice etc. formând *grupul verbal VG*. Componentele VG sunt: (a) verbul lexical; (b) auxiliar (de timp, de pasiv); (c) verb semi-auxiliar; (d) verbe de restructurare (modale, de aspect, de mișcare); (e) adverbe speciale (*mai, cam prea, și, tot*), fiecare având poziții sintactice și semantice bine definite în VG; (f) clitice pronominale (ce pot apare cu sau fără *dublarea* prin argumentele directe); (g) negații. Exemple de VG-uri: *am plecat, nu că nu mi-l va mai și plăti greu; nu i-ar fi trecut; n-ar trebui să putem cheltui; ar fi trebuit să nu mai poată trece.*

Dacă în literatura de specialitate este acceptat în general că *cliticele pronominale* și *intensificatorii* au comportament de *afix*, în timp de *complementizatorii* și *negația* au proprietăți de *cuvânt*, natura *auxiliarelor* este încă controversată. Ele sunt cuvinte morfo-sintactice după (Monachesi, 2005), afixe după (Barbu, 1999), sau clitice simple după (Zwicky, 1985).

3.2. *Subgrupul auxiliarului de timp și subgrupul verbului modal*

Cel mai frecvent și mai natural *subgrup verbal* (Verbal SubGroup, VSG) este *Subgrupul Auxiliarului de Timp* (Tense Auxiliary SubGroup, TASG): *voi fi, aș fi, am fi, sunt*. TASG poate conține, pe lângă auxiliar, *adverbe speciale* sau *negație*. TASG este *nesaturat* (i.e. cere un argument lexical) și *nepredicațional* deoarece nucleul său semantic este un verb copulativ. Din acest motiv, asociem TASG cu o *funcție de atribuire* $x := y$, unde x și y sunt variabile sau constante *intensionale* sau *extensionale*.

TASG este un subgrup al grupului *Auxiliarului de Pasiv* (Passive Tense Auxiliary Subgroup, PTASG). PTASG are același sens de atribuire copulativă, nefiind saturat, la

fel ca TASG. Exemple de subgrupuri PTASG sunt: *i-a fost dată; nu-i va fi recunoscută (diploma), trebuia să fie arestat.*

O altă substructură tipică a VG se referă la modalitate. *Subgrupul modal* (ModVSG) derivă din TASG, cu diferența că nucleul semantic pe care îl are acest subgrup este un verb modal (*a trebui, a putea*), cu posibile inserții de adverbe speciale sau negații. Exemple de ModVSG sunt: *Am fi putut alerga; Îi trebuie apă [ca] să crească; Nu-l mai puteam reține peste noapte; Ar fi putut-o vedea tot satul; Nu i l-ar fi putut da cu împrumut.* ModVSG este nesaturat și cere ca argument o altă predicatie, de ordin unu sau doi.

3.3. Modelare intensională / extensională a predicatului verbal și nominal

Să considerăm următoarea serie de predicate pentru analiza intensională / extensională.

(a) *a fost predată.* Acest VG este un predicat verbal în diateza pasivă, la timpul trecut. „*Predată*”, nucleul căruia *i* se aplică TASG, este un predicat intensional de valență 3 (nefinit, *i.e.* formă verbală absolută). Astfel, reprezentarea pentru o propoziție de genul „*Lucrarea a fost predată.*” este $lucrarea(Y) :=_{past} predată_{passive}(x, Y, z)$, unde *Y* este o variabilă extensională, iar *x* și *z* sunt predicate extensionale. Aici am considerat înțelesul extensional al substantivului „*lucrarea*”, dar acesta poate avea și un sens intensional-predicațional: *Lucrarea cu miză a pereților exteriori de către meșterii populari.*

(b) *a fost plecată.* Acest VG este în gramatica clasică un predicat nominal. Totuși poate fi interpretat și ca predicat verbal al cărui nucleu semantic este o predicatie reprezentată de un verb intransitiv (valență 1). O asemenea categorie are reprezentarea $plecată(x(X))$, unde *x* este un predicat extensional iar *X* este o variabilă extensională.

(c) *a fost frumoasă.* Acesta este un predicat nominal clar, al cărui nucleu semantic „*frumoasă*” nu mai este o categorie predicațională. Totuși, deoarece orice modifier adjectival, predicațional sau nu, necesită (cel puțin) un argument nominal, exprimat ca un predicat extensional $x(X)$, reprezentarea corectă este $frumoasă(x(X))$, cu *x* un predicat extensional și *X* variabilă extensională.

(d) *a fost elevă.* Acesta este un predicat nominal clasic, nucleul semantic fiind reprezentat de predicatul extensional $elevă(X)$, *X* variabilă (sau constantă) extensională.

(e) *va fi trădarea.* Și acest exemplu este un predicat nominal, constând dintr-un subgrup TASG al cărui nucleu semantic „*trădarea*” este un substantiv predicațional. Reprezentarea intensională a acestei structuri este $P :=_{future; act\ diat} trădarea(x, y)$, unde *P* este o variabilă intensională corespunzătoare predicatului intensional „*trădarea*” iar *x* și *y* sunt predicate extensionale corespunzând categoriei nominale predicaționale „*trădarea*”, de valență 2. De exemplu, *P* ar putea reprezenta pronumele demonstrativ (și implicit anaforă intensională) „*aceasta*” în exemplul *Aceasta a fost trădarea.*

Recapitulând această serie de exemple, se poate observa că predicatul are un subgrup TASG care se aplică unei sintagme verbale sau nominale a cărui nucleu nefinit (predicativ/nepersonal) variază astfel:

- *predată* = verb predicațional (intensional), valență 3, nesaturat, exemplul **(a)**;
- *plecată* = verb predicațional (intensional), valență 1 (intransitiv), nesaturat, ex. **(b)**;

- *frumoasă* = adjectiv nepredicațional (extensional), nesaturat, care necesită un nucleu nominal (predicat extensional): exemplul (c);
- *elevă* = substantiv nepredicațional (extensional), saturat, exemplul (d);
- *trădarea* = substantiv predicațional (intensional), valență 2, nesaturat, ex. (e).

Punctul de tranziție de la predicatul verbal către cel nominal este localizat în exemplele (b) și (c). Aceste două predicate (intuite corect de gramatică ca *nominale*) sunt de fapt un *predicat verbal* și unul *nominal*, oferind, din motive diferite, aceeași reprezentare la nivelul logicii intensionale / extensionale: *plecată* este verb predicațional-intensional, cu un singur *argument* (fiind intransitiv), în timp ce *frumoasă* este o adjectiv nepredicațional, deci extensional; fiind însă o categorie *nesaturată*, necesită ca *nucleu nominal* tot un predicat extensional.

Referințe Bibliografice

- Barbu, A.M. (1999). The Verbal Complex. *Studii și Cercetări Lingvistice*, L, no.1, București, p. 39-84.
- Curteanu, N. (2003-2004). Contrastive Meanings of the Terms “Predicative” and “Predicational” in Various Linguistic Theories (I, II). *Computer Science Journal of Moldova*, Vol. 11, No. 4, 2003 (I); Vol. 12, No. 1, 2004 (II).
- Curteanu, N. (2005). Functional FX-bar Theory Extended to Discourse (Rhetorical) Structures. In *‘Intelligent Systems’ Conference Volume*, H.-N. Teodorescu et al. (Editors), Performantica Press, Iași, pp. 169-182.
- Curteanu, N., Trandabăț, D. (2006). Functional (F)X-bar Projections for Local and Global Text Structures. The Anatomy of Predication. *Revue Roumaine de Linguistique*, București.
- Dobrovie-Sorin, C. (1994). *The syntax of Romanian. Comparative Studies*. Berlin: Mouton de Gruyter.
- Irimia, D. (1997). *Morfosintaxa verbului în limba română*. Ed. Univ. “Al. I. Cuza” Iași.
- Mann, W., Thompson, S. (1987). *Rhetorical Structure Theory: A Theory of Text Organization*. Research Report RS-87-190, Information Sciences Institute, University of Southern California, Marina del Rey, California, 80 pp.
- Monachesi, P. (2005). *The Verbal Complex in Romance. A Case Study in Grammatical Interfaces*. Oxford University Press, Oxford Studies in Theoretical Linguistics.
- Pollard, C., Sag, I. (1994). *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago & London.
- Zwicky, A. (1985). Clitics and Particles. *Language*, Vol. 62, No. 2, pp. 283-305.

PERSPECTIVE SEMANTICE DIN NOU: CUM ȘI SUB CE FORMĂ AVANSĂM LEXICOLOGIC SPRE *DLRI*

CRISTINA FLORESCU

Institutul de Filologie Română „Al. Philippide”, Academia Română, Filiala Iași

cristinafl24@yahoo.fr

Rezumat

Din perspectivă lingvistică, plecându-se de la rezultatele analizelor anterioare privind tratamentul informatic al DLR, sunt propuse două modalități de organizare și valorificare a unor eşantioane semnificative lexical din DLR (serie nouă) + DA (seria veche) în ideea verificării randamentului tratării în format electronic a acestui corpus lexicografic al limbii române, unic și din punct de vedere semantic.

Faptul că limba română își caută intens vadul academic informatizat reprezintă unul dintre elementele evidente din arealul lingvisticii contemporane. Această afirmație este un loc comun. Unghiul de abordare a respectivei probleme diferă însă, deschizând un evantai de posibilități aflate adesea în raport de opoziție.

Ne vom opri asupra a trei chestiuni.

I. Prima chestiune pune în corelație considerațiile noastre de astăzi cu cele avansate anul trecut, despre „schița unui viitor lexicografic imediat”, respectiv suita de propuneri și sugestii privind realizarea *Dicționarului limbii române informatizat și unificat* (Florescu, 2005)¹. În respectivul material se propune o modalitate de eşalonare pe parcursul a zece ani a complicatelor operații lexicografice și informatice care sunt avute în vedere în procesul de actualizare lingvistică a vechiului DA; se consideră că întreprinderea respectivă trebuie să aibă ca punct de plecare introducerea în format electronic a DA și DLR.

1. Cea mai mare parte a factorilor punctuali, a elementelor concrete și a reperelor temporale din seria de soluții cuprinse în lucrarea prezentată anul trecut, s-a modificat, unii factori – în mod previzibil, alții – imprevizibil. Un exemplu din ultima categorie: s-a micșorat drastic numărului lingviștilor lexicografi, autori ai DLR, dispuși să reia conform aceluiași parametri (fie și în format electronic) laborioasa și tehnica redactare lexicografică, extrem de tensionată profesional în ultimii ani.

Se adaugă o problemă deloc de neglijat: transformarea DLR în text adnotat presupune în mod obligatoriu desfășurarea unui demers filologic deosebit de minuțios, cronofag, care uzează specialistul în mod neprofitabil și neperformant: corectarea lingvistico-lexicografică a unui număr de câteva zeci de mii de pagini tip (2000 semne/pagină).

¹ Cf. <http://consilr.info.uaic.ro/ro/resources/pre/DLRI/DLRI.ppt>

Din cadrul „propunerilor și sugestiilor privind *Dicționarul limbii române informatizat și unificat*” a rămas cu certitudine punerea în circulație a inițierii unei planificări din perspectivă lingvistică și informatică, factor care a aplecat încă o dată benefic spiritele asupra problemei.

La ora actuală conducerea Academiei Române lucrează la un acord de colaborare dintre instituțiile filologice și informatice, acord care va reprezenta primul pas în realizarea DLRI.

2. În limba română, corpusurile de texte existente, fie ele ediții ale unor texte vechi, fie corpusuri ale limbii vorbite (limbaje regionale, familiare, televizuale, radiofonice, socio-profesionale, argotice, eșalonate stilistic sau pragmatic, psiholingvistic sau vericondițional, neologice sau conservative ș. a. m. d.), toate acestea nu sunt supuse *totdeauna* unor normative / norme lingvistice prealabile. Conceptul de normă nu este utilizat aici în sens scolastic, ci în sens larg, desemnând reguli care corespund sistemului limbii, uzuri cuprinse obiectiv în anumite categorii, tipologii, paradigme, valabile și incorporate caracteristicilor limbii române.

Se înțelege că un instrument informatizat în sensul discuției noastre este cu atât mai elastic (și cu atât mai performant prin urmare) cu cât înmagazinează (adecvându-se prin prelucrare informatică) un număr cât mai mare de texte. Este scopul oricărui demers informatic aplicat unei limbi /unor limbaje.

Această aplecare informatică asupra limbii pune totdeauna problema unei prealabile analize lingvistice care să direcționeze principal și metodologic faptele lingvistice.

Limba română în speță este inegal prelucrată lingvistic. Are gramatici performante, ortografie elaborat actualizată. Fără discuție, caracterul specific al limbii române în areal romanic este marcat și analizat. La nivel romanic, istoria limbii române, gramatica și ortografia limbii române, prin urmare partea mai puțin elastică (dacă se acceptă acest calificativ) a limbii române este vizualizată sistemic în cea mai mare parte a ei. Este suficient să ne referim la lucrările de tipul MDA, ^{1,2}DOOM sau ^{1,2}*Gramatica limbii române*. Partea de maximă elasticitate, efervescență a limbii, respectiv semantismul limbii române are cele mai multe lacune analitice. Celelalte limbi romanice dispun (dincolo de o suită, cu vechime apreciabilă, de lexicoane tezaur de diferite tipuri) de tratate de semantică și/sau de corpusuri bogate de texte eșalonate pe secole, curente culturale, scriitori, limbaje de specialitate etc.

În mod evident, DLR reprezintă la ora actuală, pentru limba română privită în totalitatea ei, cea mai bogată sursă semantică structurată cu fermitate și normată lingvistic adecvat.

Faptul pare de la sine înțeles și cunoscut. Lucrurile nu stau însă totdeauna astfel. Caracterul de „fără precedent” al extensiei semantice a DLR este de multe ori foarte greu acceptat.

Exemplele sunt nenumărate. Ne vom rezuma la un singur caz relativ recent.

În cadrul unei îndelungate cercetări lingvistice, cuprinse într-un volum (Florescu, 1999), s-a demonstrat faptul că în limba română există pentru verbul *a lăsa* o situație

specială, o idee semantică specifică, cea „a coborârii, a deplasării pe plan înclinat”², în cadrul tuturor celorlalte limbi romanice în care există bine dezvoltăți urmașii lat. *laxare*. Demonstrația punctuală³ este lungă, minuțioasă și tehnică. Au fost purtate câteva discuții pe această temă și cu o serie de specialiști romaniști neromâni. Faptele au fost acceptate, demonstrația a convins, prin urmare.

Punctul de minimă credibilitate a fost numărul de ocurențe⁴, suma de fapte lingvistico-semantică cu ajutorul cărora s-a realizat, în cadrul DLR, redactarea articolului lexicografic corespunzător analizei semantico-lingvistice a verbului românesc *a lăsa*. Pentru un lexicograf autor al DLR, a selecta, ierarhiza, încheia pentru prima oară semic (la acest nivel de minuție) în jur de câteva mii de ocurențe este ceva obișnuit. Pentru un lingvist de aiurea numărul se dovedește incredibil.

Ceea ce demonstrează încă o dată, și din exteriorul sistemului, că structura semantică a DLR este excepțională.

II. Al doilea aspect se referă la modalitatea concretă de unificare a celor două structuri lexicografice (DA și DLR).

Realizarea unui grup lexical „de probă” (prin care să se actualizeze articolele DA) este obligatorie și semnificativă.

Au existat câteva încercări meritorii în acest sens, nefăcute public. Ele însă nu au dispus de instrumentul informatic. Prin urmare, timpii de apreciere sunt în afara sferei noastre de interes.

Un prim eșantion care se dezvoltă în direcția introducerii în format electronic a *Dicționarului limbii române* a realizat, în cadrul unui grant⁵, un instrument de lucru adecvat informatic și lingvistic chestiunii numit DLReX (ale cărui caracteristici și performanțe sunt prezentate pe larg în (Haja et al., 2005).

Am propus, ca o a doua întreprindere realizată în sensul discuției noastre, un demers de același calibru ca cel de mai sus, demers care se oprește, pentru început, asupra unui număr de 142 de cuvinte, eșalonate în DLR și în DA, derivate pe teren românesc cu sufixul de origine slavă (slavă veche, slavonă, bulgară și sârbă) în *-iște*. În funcție de sursa de finanțare, cercetarea se va putea extinde și la alte categorii de derivate.

Lingvistic problema este complicată. O categorie lingvistică este o mulțime „continuă” pe când un grup / eșantion lexicografic este discret. Trecerea de la un grup ale cărui limite sunt în mod inerent și obiectiv vagi, la un număr bine limitat de intrări

² Cităm câteva structuri prototipale: *păsările se lasă pe câmp, X se lasă în fântână / la vale, s-a lăsat ceața, i s-au lăsat măruntaiele*.

³ Cf. (Florescu, 1999: 46 – 57, 128 - 137)

⁴ Compară, în acest sens, două studii paralele: în (Soares, 1999:77) pentru port. *deixar* (și forma veche *leixar*) cifra citată este de aproape 4.000 ocurențe, în (Florescu, 1999: 48) pentru rom. *a lăsa* cifra citată se ridică la un număr de peste 15.000 de ocurențe; reamintim faptul că ambele cuvinte sunt echivalente lingvistice (cu etimologie comună) în areal romanice.

⁵ Cf. grantul CNCSIS nr. 1415 pe 2003-2005, prezentat și valorificat în (Haja et al., 2005).

lexicografice reprezintă o cercetare lingvistică în sine, complexă, cu multe probleme lexicologice de rezolvat.

Plecând de la faptul că acest aspect lexical (stabilirea limitelor categoriei lexicale a derivatelor în discuție) este deja rezolvat în cadrul unei cercetări lingvistice (Florescu, 2006), proiectul propune introducerea - prin scanare, OCR-izare și transformare în limbaj XML - în memoria computerului a unei sume de cca 200 pp. DLR și DA - pagini echivalente articolelor care cuprind derivatele în studiu; va urma prelucrarea paginilor respective cu ajutorul DLReX, rafinarea acestuia din urmă și realizarea unei serii de articole lexicografice actualizate din seria veche DA, conform principiilor lexicografice ale seriei noi DLR.

Existența cercetării informatice anterioare oferă actualului demers multe șanse de reușită.

În final vom dispune de un mini-corpus lexicografic al limbii române tip DLR, adnotat în format XML, unificat lingvistic prin mijloace informatice.

III. Al treilea aspect se referă la modalitatea de valorificare a masei lexicale fără precedent, a mulțimii de citate / contexte / texte existente în seria nouă DLR.

Dacă admitem ca premiză de lucru faptul că prima parte a travaliului legat de introducerea seriei noi DLR în format electronic va fi realizată, ne găsim în fața a zece tomuri (34 de volume) DLR scanate, adnotate și prelucrate în XML cu ajutorul DLReX-ului.

Cum poate fi valorificată această masă lexicală, acest corpus de micro-texte, pentru actualizarea seriei vechi DA?

Chestiunea este lexicologic și lingvistic extrem de delicată. Reprezintă una dintre frânele care blochează eventualele aproximări și planificări. Nu se poate avansa un proiect academic fără o relativ corectă evaluare a acestui fapt.

Părerile lingviștilor sunt extrem de împărțite: unii consideră că, în urma efortului preluării citatelor (micro-contexte) din DLR, grupul lexical aparținând seriei vechi (cuvintele din porțiunile A-C, F-J) va avea lacune lingvistice (sensuri, etimologie etc.) prea mari pentru a merita efortul respectivei preluări. Alți lingviști consideră că, în urma valorificării citatelor (micro-contexte) din DLR, „zonele albe” corespunzătoare articolelor din DA vor fi relativ restrânse și că trebuie făcut efortul preluării acestor citate. De exemplu, articolul lexicografic corespunzător cuvântului *casă* (redactat, prin urmare, în seria veche, în DA), prin preluarea tuturor citatelor existente în DLR în care apare cuvântul, va putea utiliza, să zicem, și contextele în care apare sintagma *casa mare*, contexte care pot fi recuperate din DLR tom VI, Litera M, articolul MARE.

Există tentația de a spune că adevărul este undeva la mijloc. Nimic mai fals. Oricare dintre ipostazele probabile sunt posibile. Singura soluție este realizarea unui proiect parțial care să verifice în mic realitatea lingvistico-lexicografică.

Un asemenea proiect se poate opri asupra tomului ieșean XIII al literelor: V, W, X și Y, tom care însumează 1340 pagini DLR și 6327 de intrări (cca 5000 de cuvinte). Dintre acestea numai 325 pagini și 1747 intrări (cca 1500 cuvinte) nu se găsesc în memoria computerului.

Scanarea și adnotarea întregului tom XIII va fi urmată de „vărsarea” ocurențelor / citatelor / microcontextelor - în funcție de cuvintele relevate și delimitate corect lingvistic și lexicografic - în ordine alfabetică. Se înțelege că se va avea în vedere confruntarea, prin operații lexicologico-lexicografice specifice, cu o sumă de liste alfabetice lexicale anterioare, inclusiv acelea din DA și din MDA.

Stabilirea listei de cuvinte în funcție de variantele lexicale ale unui cuvânt și de structurile omonimice va presupune eșalonarea a cel puțin trei operații de identificare.

1. Prima identificare va fi făcută la nivelul intrărilor.
2. Următoarea: la nivel diacronic – modalitatea de acoperire a timpilor filologi originari, vechi, moderni și contemporani în cadrul fiecărui articol / cuvânt.
3. Ultima identificare comparativă și corelativă va privi extensivitatea semică (sensuri, subsensuri, construcții fixe, variații dialectale, gramaticale și stilistice).

Se va realiza un instrument de lucru adecvat informatic și lingvistic, cu ajutorul căruia informaticianul și lexicograful vor putea face o apreciere corectă privind gradul de performanță lingvistico-informatică a demersului respectiv.

Rezultatele acestor două întreprinderi - a) eșalonul DLRI al derivatelor în *-iște* și b) posibilitățile de acoperire semantică a nevoilor de actualizare a DA - vor limpezi o suită de probleme care, la ora actuală, în lexicografia academică românească, par insurmontabile nu atât prin imposibilitatea realizării, cât prin numărul prea mare de variante posibile.

Referințe bibliografice

Haja, G., Dănilă, E., Forăscu, C., Aldea, B. (2005). *Dicționarul limbii române (DLR) în format electronic. Studii privind achiziționarea*. Editura Alfa, Iași.

Cristea, D. (2005). *Resurse lingvistice și tehnologii ale limbajului natural. Cazul limbii române*. Prelegerile Academiei Române. Filiala Iași.

^{1,2}*Gramatica limbii române*. Academia Română. Institutul de Lingvistică “Iorgu Iordan-Al.Rosetti”, București. Editura Academiei Române. Ediția I, 1963. Ediția II, 2005.

DA = *Dicționarul limbii române*. Academia Română. București. Librăriile Socec, Universul. Tom I-II, 1913-1937.

DLR = *Dicționarul limbii române. Serie nouă*. București. Editura Academiei Române. Tomul VI ș. u., 1965 etc.

- ^{1,2}DOOM = *Dicționarul ortografic, ortoepic și morfologic al limbii române*. Ediția I, Editura Academiei Române, București, 1982. Ediția a II-a, Univers Enciclopedic, București, 2005.
- Florescu, C. (1999). *Gîndire specifică și gîndire europeană în semantismul românescului a lăsa*, Iași, Document.
- Florescu, C. (2005). Propuneri și sugestii privind Dicționarul limbii române unificat și informatizat (DLRI). *Atelierul de lucru ConsILR. Resurse lingvistice și instrumente pentru prelucrarea limbii române*, Iași, 3 noiembrie 2005.
- Florescu, C. (2006). *Liniește și derivatele pe teren românesc în -iște*. În Volumul Omagial Mioara Avram, București, Editura Academiei, 2006 (sub tipar).
- MDA = *Micul Dicționar Academic*, Academia Română, Institutul de Lingvistică "Iorgu Iordan", București, Univers Enciclopedic, vol. I-IV, 2001-2002.
- Soares da Silva, A. (1999). *A Semântica de 'Deixar'. Uma Contribuição para a Abordagem Cognitiva em Semântica Lexical*. Braga, Fundação Calouste Gulbenkian.

MODELAREA RELAȚIILOR SEMANTICE ÎNTR-UN DICȚIONAR DE SIMBOLURI

CRISTINA CIOCÂRLĂU, MIHAELA BRUT

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

{cciocarlau, mihaela}@info.uaic.ro

Rezumat

Dicționarele de simboluri sunt instrumente de lucru foarte importante pentru cei care lucrează în domenii precum critica literară, lingvistica, literatura, dar și relații publice, marketing sau arte vizuale. Considerând poezia eminesciană un punct de plecare pentru descoperirea bogăției de semnificații ale cuvintelor din limba română, aplicația prezentată în articolul de față oferă posibilitatea întreținerii și utilizării unui dicționar de simboluri, organizat pe arii semantice. Adăugarea unui nou cuvânt în dicționar de către administrator are ca efect evidențierea lui în toate poeziile eminesciene, utilizatorul obișnuit având posibilitatea de a-i vizualiza diversele accepțiuni și de a-și dezvolta propriile comentarii, definiții, adnotări. În plus, alături de definițiile simbolului selectat, utilizatorului îi sunt semnalate și simbolurile din aceeași arie semantică. Astfel, având acces la diversele valențe semantice ale unui cuvânt, utilizatorului îi va fi facilitată utilizarea mai nuanțată și mai percutantă a acestuia.

Aplicația este implementată în XHTML, CSS, XML, DOM PHP, asocierea de metadata XML simbolurilor poetice constituind un prim pas spre dezvoltarea unei aplicații de Web semantic pentru simbolurile înmagazinate de limba română.

1. Introducere

Operele artistice de valoare propun o uimitoare dialectică a semnificațiilor și contextelor, incitând la interpretări multiple, oferind astfel delectări intelectuale și spirituale. Simbolurile poetice consacrate apar în operele marilor creatori îmbogățite semantic sau chiar resemantizate, iar miturile sunt reinterpretate sau le este extinsă semnificația. Deseori, artiștii își dezvoltă propriile simboluri poetice, care contribuie într-un mod insolit și fascinant la stabilirea corespondențelor simbolice subterane ale operei lor. Privită prin prisma unei creații artistice, *constelația simbolurilor*¹ are o perspectivă personalizată. Metamorfoza semnificațiilor se dezvoltă studiind maniera în care zestrea culturală a umanității a fost asimilată și trasfigurată de o experiență artistică particulară.

Dicționarele de simboluri sunt punctul de plecare în acest demers. Dicționare precum *Dicționarul de simboluri* al lui Jean Chevalier și Alain Gheerbrant, sau *Elsevier's Dictionary of Symbols and Imagery* scris de Ad de Vries și revizuit de Arthur de Vries,

¹ Termen consacrat de Gilbert Durand, *Structurile antropologice ale imaginarului*, Editura Univers Enciclopedic, București, 1998.

ori *The Complete Dictionary of Symbols* de Jack Tresidder sunt instrumente de lucru uzuale ale criticilor literari și artistici.

Disponibilitatea on-line a acestui tip de dicționare ar putea facilita accesul la informație. Avantajul dicționarului de simboluri propus de noi este acela că permite criticilor literari - sau altor tipuri de utilizatori interesați - să își construiască un sistem de adnotări referitoare la semnificațiile particulare ale simbolurilor în opera eminesciană, aplicația fiind ușor extensibilă la operele diverșilor artiști sau la varii contexte. Astfel, paralel cu definițiile furnizate de dicționar pentru un anumit simbol, utilizatorul își gestionează un set de comentarii legate de aceste semnificații. În plus, pentru fiecare simbol, sunt oferite și referințele la simbolurile din aceeași arie semantică.

Deoarece simbolurile sunt foarte importante în dezvoltarea unor diverse tipuri de discursuri, considerăm util dicționarul propus de noi și celor care lucrează în publicitate, marketing, relații publice, arte vizuale.

2. Funcționalitățile oferite de dicționarul de simboluri

Aplicația *Dicționar de simboluri on-line* oferă suport pentru două tipuri de utilizatori: utilizatori obișnuiți, care doresc să acceseze suportul oferit de dicționar, și administrator, care are posibilitatea gestionării globale a aplicației.

<p>A C D E F G I Î J K L M</p> <p>A</p> <p>Adio Afară-i toamnă Amorul unei marmure Atât de fragedă..</p> <p>C</p> <p>Când amintirile Călin (file din poveste) Ce e amorul? Ce te legeni... Ce-ți doresc eu ție, dulce Românie Când însuși glasul (din Sonete) Copii eram noi amândoi Crăiasa din povești Criticilor mei Cugetările sârmanului Dionis Cu mâne zilele-ți adaogi</p>	<p>Ce te legeni...</p> <p>Ce te legeni, codrule, Fara ploaie, fara vant, Cu crengile la pamant? De ce nu m-as legana, Daca trece vremea mea! Zna scade, noaptea creste Si frunzisul mi-l rareste. Bate vantul frunza-n dunga - Cantareti mi-i alunga; Bate vantul dintr-o parte - Iarna-i ici, vara-i departe. Si de ce sa nu ma plec, Daca pasarile trec! Peste varf de ramurile Trec in stoluri randurele, Ducand gandurile mele Si norocul meu cu ele. Si se duc pe rand, pe rand, Zarea lumii-ntunecand, Si se duc ca clipele, Scuturand aripele, Si ma lasa pustiit, Vestejit si amortit Si cu doru-mi singurel, De ma-ngan numai cu el!</p> <p>Simbolul noapte are definitiile:</p> <ol style="list-style-type: none"> Nocturnul este cadrul de predilectie al romanticilor; noaptea dispar granitele dintre lumea reala si cea a misterelor, dintre taramul vietii si al mortii. Fiinta traieste experiente capitale, se manifesta demoniac, sau capata forta sa se inalte la ceruri. <p>Adaugă</p> <ol style="list-style-type: none"> Valentele simbolice ale noptii sunt extrem de variate in creatia eminesciana; noaptea terestra, poleita de lumina de luna, intunericul cosmic, noaptea-moarte sau timpul halucinant al "somniai" se constituie in cateva ipostaze ale nocturnului. Feeria lunara alterneaza cu intunericul de nepatruns, creind un univers magic al existentei lucrurilor. <p>Adaugă</p> <ol style="list-style-type: none"> Timp oracular, miezul noptii reprezinta simbolic si o poarta spre eternitate. Tot acum se produc metamorfoze stranii, iar mortii urca in luna. (Fat-Frumos din lacrima) <p>Adaugă</p> <p>Simboluri conexe: luna; stea; vis; soare;</p>
--	--

Figura 1: La afișarea unei poezii sunt marcate simbolurile existente în dicționar

Utilizatorul obișnuit are acces, după autentificare, la lista poeziilor eminesciene. La afișarea unei poezii, simbolurile care se regăsesc în dicționar vor apărea marcate în text. Selectarea unui simbol din poezie va determina afișarea definițiilor simbolului, așa cum

sunt ele memorate în dicționar, precum și enumerarea simbolurilor din aceeași arie semantică. În partea de jos a ecranului, utilizatorul are la dispoziție o zonă în care poate adăuga diferitele definiții ale simbolului ales, și în același timp poate aduce modificări și completări acestor definiții, creând un mic comentariu aplicat la contextul curent. Operația de redactare a zonei respective poate continua și după selectarea altui simbol din poezia curentă sau dintr-o altă poezie; informația este persistentă. În plus, alături de afișarea definițiilor pentru simbolul curent selectat, utilizatorului îi sunt puse la dispoziție referințe către simbolurile din aceeași arie semantică, astfel încât, dacă se dorește un studiu asupra unei problematice particulare a poeziei eminesciene, să se ajungă la toate fragmentele și simbolurile elocvente. Această zonă de editare rămâne activă pe parcursul navigării în alte poezii, iar conținutul editat poate fi salvat, fiind accesibil ulterior doar utilizatorului curent. Pe viitor ar putea fi adăugată o opțiune de partajare între utilizatori a comentariilor făcute unor diverse simboluri, figuri de stil, poezii.

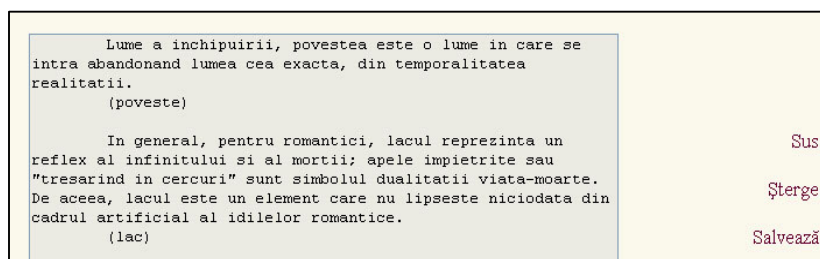


Figura 2: Zona de editare în care se pot adăuga atât definiții de simboluri, cât și comentarii proprii

Funcția principală a administratorului aplicației este gestionarea intrărilor în dicționar. Adăugarea unui nou simbol în dicționar presupune, alături de precizarea definițiilor acestuia, indicarea simbolurilor conexe, făcând parte din aceeași arie semantică. Desigur, pot fi indicate drept simboluri conexe doar simbolurile existente la momentul curent în dicționar, includerea ulterioară a unui simbol ce poate avea această însușire trebuind făcută cu indicarea simbolului curent drept conex. Ar putea fi dezvoltat pe viitor un mecanism automatizat de moștenire sau de tranzitivitate a apartenenței la aceeași arie semantică.

Modificarea definiției unui simbol existent este o altă funcționalitate accesibilă administratorului. Simbolul este selectat dintr-o listă derulantă, iar definiția nou completată o va substitui pe cea veche. Un procedeu similar se aplică în cazul în care administratorul dorește îmbogățirea colecției de definiții corespunzătoare unui anumit simbol cu o nouă definiție.

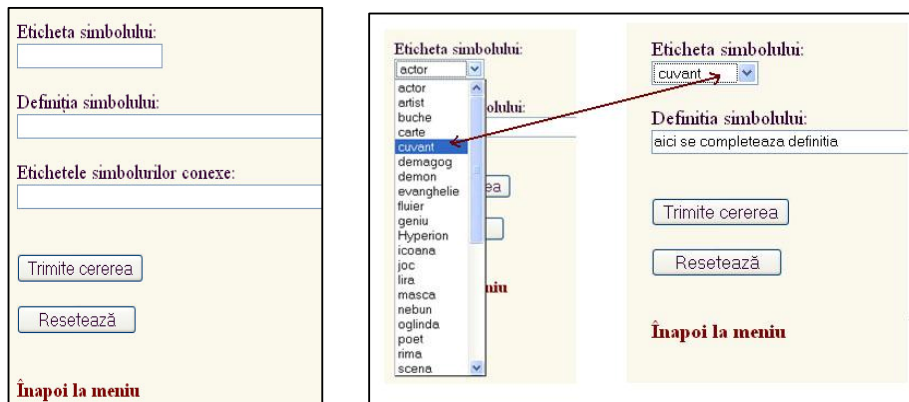


Figura 3: Adăugarea unui nou simbol în dicționar; completarea definițiilor pentru un simbol existent

După inserarea unui nou simbol în dicționar, administratorul trebuie să efectueze operația de parsare a poeziilor, pentru ca noul simbol să apară marcat și în text. În cadrul acestei operații trebuie precizate formele flexionare ale noului simbol, pentru ca și diversele apariții ale acestora în text să fie recunoscute și marcate corespunzător. Unele studii ne arată că această operație poate fi efectuată de un modul de generare automată a formelor flexionare ale unui cuvânt. Administratorul mai poate face listarea simbolurilor existente, fiecare însoțit de identificatorul unic asociat în mod automat la includerea în dicționar.

3. Tehnologiile utilizate în implementare

Interfața utilizator a fost organizată și formatată utilizând XHTML și CSS2. Flexibilitatea soluției noastre este conferită de utilizarea XML pentru reprezentarea și stocarea datelor folosite, precum și de utilizarea implementării PHP a DOM (*Document Object Model*) pentru prelucrarea informațiilor din documentele XML.

Poeziile sunt stocate în fișierul *poezii.xml*, a cărui structură este descrisă în imaginea alăturată. Acest format ne permite evitarea informației redundante, proprietate mai greu de obținut, de exemplu, în cazul memorării poeziilor într-o bază de date, în care s-ar alocă o întregă înregistrare pentru un singur vers sau chiar un singur cuvânt, trebuind să fie precizat - la fiecare înregistrare - codul poeziei, al volumului, al autorului etc.

```

- <poezii>
- <poezie nume="adico">
  <titlu>Adico</titlu>
  + <strofa></strofa>
- <strofa>
  <vers>De astazi dar tu fa ce vrei,</vers>
  <vers>De astazi nu-mi mai pasa</vers>
- <vers>
  Ca cea mai dulce-ntre
  <simb id="femeie">femeie</simb>
  </vers>
  <vers>Ma lasa.</vers>
</strofa>
+ <strofa></strofa>

```

În cadrul fiecărei poezii sunt adnotate cuvintele care constituie intrări în dicționar. Dicționarul în sine este reprezentat tot în format XML, fiecare simbol având un identificator unic generat automat la crearea acestuia. De asemenea, fiecare simbol are indicate mai multe simboluri conexe, prin intermediul identificatorilor acestora.

MODELAREA RELAȚIILOR SEMANTICE ÎNTR-UN DICȚIONAR DE SIMBOLURI

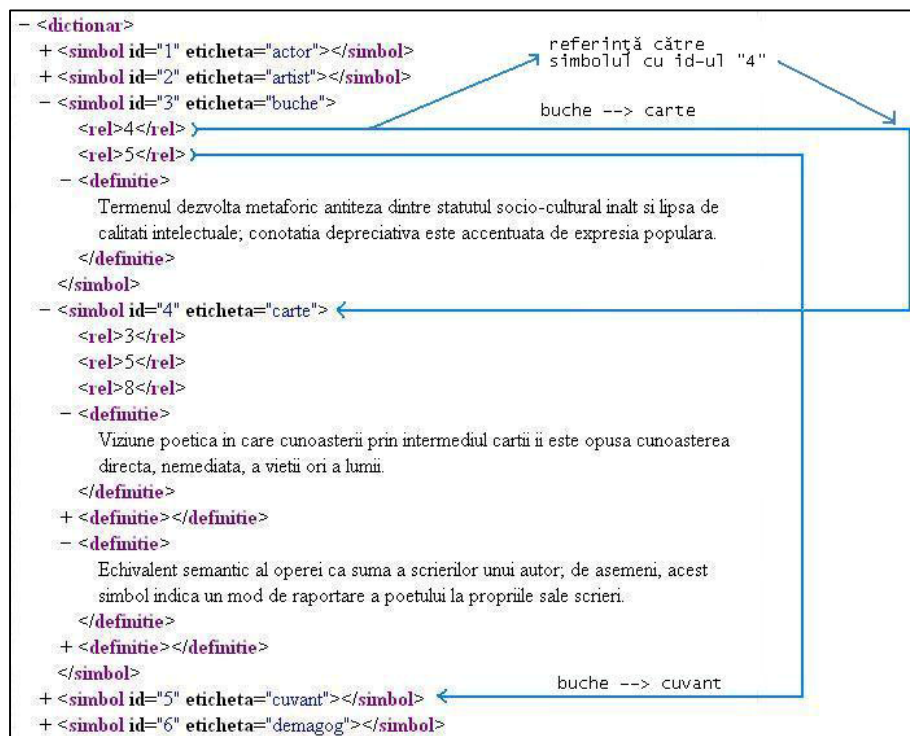


Figura 4: informațiile în format XML din dicționarul de simboluri

4. Concluzii și direcții viitoare

Prima extindere a aplicației ar trebui să aibă în vedere integrarea simbolurilor din dicționar într-o ontologie de simboluri, oarecum complementară ontologiilor de concepte generale de tip *WordNet*, *SUO* sau *SUMO*. Acest lucru este viabil deoarece fiecare simbol are o serie de accepțiuni, înregistrate de dicționarele de simboluri, evident diferite de accepțiunea comună a cuvântului care desemnează respectivul simbol. Pe acest schelet ontologic pot fi grefate accepțiunile particulare dobândite de fiecare simbol în contextul unei opere sau a unui poem, și pot fi stabilite de asemenea diverse relații (de exemplu, de apartenență la un câmp semantic comun) între două sau mai multe simboluri.

Deci, urmând ideile din (Niles & Pease, 2003) și (Huang, 2004), ar putea fi făcută trecerea de la un simplu text la un lexicon (de exemplu, în genul unui dicționar de simboluri obișnuit) și apoi la o ontologie dacă s-ar ține cont de delimitarea între:

- *Ontologie generală*: ontologie de nivel înalt partajată de toate domeniile (precum SUMO și Wordnet)
- *Ontologie specifică*: pentru un domeniu, perioadă istorică, autor, eventual câmp semantic etc.

Un dicționar general de simboluri ar putea constitui o ontologie generală, iar adaptările și modificările de semnificații ale simbolurilor în contextul unei opere particulare, sau a

unui scriitor particular ori a unui curent literar particular, ar putea fi structurate într-o ontologie specifică.

Aplicația *Dicționar de simboluri on-line* se dorește a fi extinsă și prin adăugarea spre parsare a altor opere literare, aparținând unor autori diferiți, sau chiar a fragmentelor diverse, aflate în afara creației artistice. Dicționarul poate fi dezvoltat într-o altă direcție prin asocierea și integrarea de simboluri grafice alături de simbolurile textuale actuale. Pentru utilizatori ar putea fi utilă și gestionarea unui istoric al comentariilor, în care să fie marcate simbolurile și fragmentele de plecare pentru fiecare comentariu, efectuându-se astfel cu ușurință corelarea simbolurilor/fragmentelor cu adnotările.

Deși reprezintă doar un început, aplicația *Dicționar de simboluri on-line* oferă în acest moment un cadru general de lucru pentru exegeții operei eminesciene. În perspectiva oferită de diversele direcții de dezvoltare, putem vorbi despre această aplicație ca despre o unealtă care vine în sprijinul specialiștilor din domeniul literaturii și chiar al lingvisticii.

Referințe bibliografice

- Abdoullaev, A.Sh. (2006). *Ontology, Semantic Technology, and Knowledge Society: World Wide Intelligent Web*, <http://www.eis.com.cy/>
- Chevalier, A., Gheerbrant, A. (1998). *Dicționar de simboluri*, Ed. Artemis.
- Huang, Chu Ren (2004). Text-based Construction and Comparison of Domain Ontology: A Study Based on Classical Poetry, *Proceedings of the 18th Pacific ASIA Conference on Language, Information and Computation*.
- Niles, I & Pease A. (2001). Towards A Standard Upper Ontology. In *Proceedings of FOIS 2001*, October 17-19, Ogunquit, Maine, USA.
- Niles, I & Pease A. (2003). Linking Lexicons and Ontologies: Mapping WordNet to the Suggested Upper Merged Ontology, *Proceedings of the IEEE International Conference on Information and Knowledge Engineering*, 412-416.
- *** SUMO (*Suggested Upper Merged Ontology*) <http://ontology.teknowledge.com/>
- *** SUO - *Standard Upper Ontology*, The IEEE Standard Ontology: <http://suo.ieee.org>
- *** *Wordnet. A lexical database for the English language*:
<http://wordnet.princeton.edu/>

DREPTUL DE PUBLICARE PE WEB

NOEMI BOMHER

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

noemibom@yahoo.com

Rezumat

În contextul societății informaționale se observă, chiar la o sumară cercetare a site-urilor românești de pe Internet, prezența extrem de redusă a resurselor on-line de literatură română, privity din perspectiva istoriei literare, a teoriei literaturii, a teoriilor predării-învățării de înalt nivel științific etc. Situația se impune a fi semnalată în vederea găsirii soluțiilor de ameliorare, în limitele legii dreptului de autor în vigoare, printr-o strategie stabilită, la nivel instituțional, de către cei în măsură să o facă.

1. Introducere

Problema dreptului de difuzare a informației prin mijloace electronice nu este nouă. Ea a fost dezbătută, la nivel european, iar liniile generale ale unei legislații specifice au fost stabilite (Huét, Mails, 1989).

Chestiunea pe care dorim să o punem, în cele ce urmează, nu se referă neapărat la restricțiile legislative impuse de legea dreptului de autor și implicațiile acestora în difuzarea *on-line* a textelor literare românești, ci la faptul că, în limitele acestei legi, prezența autorilor studiați de istoria literaturii române, la nivel universitar, este nepermis de săracă.

În condițiile de învățământ actuale, în care numărul studenților de la facultățile umaniste, în speță de la facultățile de litere, este în creștere, posibilitatea de documentare oferită de bibliotecile tradiționale este insuficientă. Chiar la nivel preuniversitar, accesarea resurselor electronice de către tinerii învățăcei este tot mai frecventă și este cunoscut faptul că, urmare a cererii mari, s-au creat site-uri de tipul „www.referate-online.ro” unde sunt postate lucrări, cu autori cvasianonimi, redactate după criterii diferite, în domenii diferite, fără însă a avea girul unui specialist.

Un program educativ, dirijat la nivel național de un for academic, ar trebui să aibă în vedere realizarea unor lucrări-model, de referință atât pentru profesori cât și pentru elevi, cu prevederea explicită a interdicției reproducerii / însușirii materialelor în nume personal.

2. Educația în societatea informațională

2.1. Societatea informațională și schimbările impuse de aceasta

Ideea coordonării instituționalizate a promovării literaturii române prin mijloace electronice se impune cu atât mai mult cu cât există în derulare, la nivel guvernamental, un amplu program de informatizare a instituțiilor și a domeniilor de desfășurare a activităților economice și științifice. Or, aceste direcții de implementare a informatizării instrumentelor și resurselor ar trebui gândite în ansamblu, în funcție de implicațiile sociale ale fenomenului.

Societatea informațională nu este o societate nouă, ruptă de societatea existentă. Ea reprezintă un stadiu evoluat al acesteia din urmă. Ceea ce nu trebuie să se uite este faptul că fundamentele stabilității sociale rămân aceleași: democrație, civilitate, solidaritate, egalitate, muncă. Doar forma de expresie, modul de exercitare a acestor drepturi și intensitatea diferitelor aspecte ale vieții democratice pot fi altfel decât cele de până în prezent, dar nu contradictorii. Este adevărat că internetul schimbă tradițiile și obiceiurile de a gândi și de a acționa, cere repunerea în discuție a multor probleme, impune renașterea cunoștințelor și a practicilor. Există instituții pe care noua conjunctură comunicațională le îngrijorează: anumite administrații statale se tem de tulburări; poliția și justiția se simt depășite de proceduri legale care îi împiedică să urmărească eficientele criminalii care se joacă cu frontierele și cu distanțele doar printr-un *click de mouse...*, instituțiile fiscale doresc să impună tranzacții care să restrângă dezvoltarea comerțului electronic etc., etc.

În mod firesc, dreptul de comunicare la distanță interactivă ar trebui să ofere răspunsuri satisfăcătoare unor interese principial opuse: libertatea comerțului și dreptul consumatorului; dreptul de autor și acela de utilizator; libertatea expresiei și libertatea de a limita această libertate; obligația de a asigura securitatea persoanelor, a bunurilor și dreptul la viața privată și deci la protejarea / criptarea mesajelor.

Dinamismul societății contemporane crește ca rezultat al dezvoltării tehnologiilor informatice și de telecomunicație. Legăturile devin mai strânse, schimbările sunt rapide și inovațiile apar succesiv, cu o mare frecvență. Această dinamică provoacă o stare de stupefacție sistemului legislativ, pentru că, tradițional, legea nu este un instrument care să se poată adapta suficient de repede, „din mers”, rigorilor proceselor dinamice prin care trecem și schimbărilor tehnologice rapide. Dezvoltarea care afectează cel mai puternic fundațiile legii o reprezintă evanescența sau dispariția frontierelor teritoriale, procesele de dematerializare și declinul intervențiilor umane nemijlocite.

Principiile fundamentale care asigură o protecție legală cetățenilor rămân importante și un obiectiv capital constă în asigurarea că aceste principii sunt respectate în noua societate. Trebuie să subliniem faptul că noile tehnologii nu sunt un scop social în sine și că legislatorii nu trebuie să fie biruiți de această nouă tehnologie.

O societate în care distanțele și granițele dispar, presupune modificarea viziunii asupra multor domenii: lucrul „la distanță” schimbă relațiile între lucrători, între lucrător și angajator; informațiile publicate pe net de diversele ministere modifică relațiile între

guvern și cetățeni; comerțul *on-line* creează alte relații între producători și consumatori; serviciile bancare prin internet facilitează relația client – bancă ș.a.m.d.

2.2. Educația – domeniu fundamental al oricărei societăți – în contextul informatizării

Unul dintre cele mai însemnate domenii la care societatea informațională actuală trebuie să mediteze atent este acela educativ.

Învățământul la distanță (IDD) propus de instituțiile de tip universitar este rezultatul modificărilor impuse de noul sistem de comunicare.

Pornind de la materialele documentare realizate pentru această formă de învățământ, de la situația de criză a fondului de carte pentru bibliotecile tradiționale, de la nevoia tot mai mare de facilitare a accesului la informație, la sursele textuale propriu-zise și bazându-ne pe o experiență de peste treizeci de ani în domeniul predării istoriei și a teoriei literare în cadrul Facultății de Litere a Universității „Alexandru Ioan Cuza”, credem că trebuie susținută necesitatea creării unui departament guvernamental care să coordoneze, prin intermediul Academiei și al instituțiilor de învățământ superior, sistemul informațional de pe internet, adecvând la noul sistem legile privitoare la învățământ, legile privitoare la publicarea operelor literare pe web, la modul de promovare a acestora, cu o concentrare specială asupra sistemelor de legi referitoare la traducerea / publicarea textelor literare pe internet, precum și la răspândirea acestora.

În domeniul literaturii, se impune crearea unor programe de nivel național care să finanțeze proiecte ce ar trebui să aibă ca finalitate publicarea operelor fundamentale ale literaturii române (de la primele scrieri literare românești până la literatura interbelică) în ediții critice definitive (în parte existente, în bună parte urmând a fi realizate de specialiști filologi), instrumente indispensabile cercetării literare fundamentale.

Desigur, acest demers nu trebuie să impiețeze legislația în vigoare și este necesară menținerea drepturilor privitoare la taxa de timbru ș.c.l., astfel încât instituții tradiționale, precum Uniunea Scriitorilor din România, să nu aibă de suferit. În orice caz, sistemul informatic permite găsirea unor soluții de acces condiționat la textele publicate pe net.

Pași în această direcție au fost făcuți, ca urmare a unor inițiative izolate ale unor entuziaști cercetători ori ale unor edituri. Ne referim aici la editarea integrală, în format electronic, a *Operele* lui Mihai Eminescu și ale celor semnate de I. L. Caragiale. Dar aceste demersuri sunt izolate și accesibile în mod restrictiv doar celor care își permit achiziționarea respectivelor CD-uri. Or, se știe, specialiștii în literaturi clasice sunt favorizați din acest punct de vedere. Ediții adnotate la cuvânt, corpusuri paralele de texte latine și grecești sunt la îndemâna utilizatorilor de net, fără să fie impuse nici un fel de restricționări. Iar acesta este doar un exemplu. Marile biblioteci europene au început crearea unor imense baze de texte în vederea conservării fondului de carte veche și rară, dar și a unor ediții epuizate ale operelor marilor autori ai literaturii lumii. Fiecare dintre aceste instituții au stabilit modalități proprii de facilitare a relației cititor – carte, în

funcție de legislațiile naționale și de cele europene, în funcție, mai ales, de politicile specifice de promovare a propriilor culturi și limbi.

Publicarea textelor literare pe net ar facilita, în bună măsură, crearea unor proiecte de traducere a acestora în limbile de circulație internațională mai ușor realizabile și profitabile din multe puncte de vedere. O bibliotecă electronică a literaturii române – în original și tradusă – ușor de consultat, eliminând distanțele și rigorile impuse de contractul cititor-biblioteca obișnuită, ar face mai vie relația carte-receptor, indiferent de granițele fizice și cele subiective.

O cultură „minoră”, ca să cităm epitetul folosit de Lucian Blaga, precum este a noastră nu poate avea decât de câștigat de pe urmă utilizării constructive a noilor mijloace de conservare și de promovare a limbii și literaturii.

Pe lângă publicarea textelor literare propriu-zise – a celor mai vechi de 25 de ani și a celor contemporane, cu îngăduința autorilor interesați de această formă de promovare a propriei creații – este necesară și publicarea instrumentelor necesare predării / învățării. Ne referim aici, pe de o parte, la dicționarele de specialitate, și, pe de altă parte, la cursuri și la volume de istorie, critică literară, teorie și poetică, la manualele (atât de multe!) alternative, compendii și antologii comentate ș.a.

3. Închinare seniorilor

Un câștig îl constituie deocamdată publicarea, și în format electronic, a unor periodice de cultură consacrate (d. ex. „România literară”, „Dilema veche”, „22”, de la București, „Timpul” de la Iași) într-un proces ce nu dezechilibrează economic aceste reviste, susținute și de Ministerul Culturii și Cultelor. Dar, în orice țară civilizată, cultura este finanțată și promovată de instituții guvernamentale naționale. Sunt de semnalat, de asemenea, publicațiile cu orientare literară, ori cenacluri literare românești (d. ex. www.clublitar.com; www.agero-stuttgart.de; www.onlinegallery.ro) care există numai în spațiul virtual al calculatorului, indiferent de notația finală de după punct (.com, .de, .ro), fără ca acest statut „imaterial” să determine vreo scădere a frecvenței și, deci, a cunoașterii lor.

Ne înclinăm dinaintea seniorilor ce pot decide soarta literaturii române în societatea informațională rugându-i să ia aminte, *sine ira et studio*, la mersul vremurilor și la destinul propriei noastre culturi și să-și amintească faptul că portretul nostru în ansamblul colajului global singuri ni-l facem, singuri îl putem expune, dacă vrem ca el să ne reprezinte și să nu arate asemenea unei caricaturi glumețe, în cel mai blând dintre cazuri.

Referințe bibliografice

Huét, J., Mails, H. (1989). *Dreptul informaticii și telecomunicațiilor. Situația întrebărilor. Texte și jurisprudență, studii și comentarii*, Paris, Litec, 1011 p.

MODELARE CU ONTOLOGII ȘI ADNOTĂRI

RADU CIBOTARU

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

deceneu@info.uaic.ro

Rezumat

Lucrarea prezintă un studiu de caz privind crearea, popularea și utilizarea unei ontologii. Studiul se rezumă la un domeniu restrâns și anume, la camere digitale, descrie structura ierahică de clase și relațiile dintre proprietăți ale ontologiei. Prezintă o modalitate de populare automată a ontologiei, bazată pe șabloane. Face un studiu comparativ a metodei descrise cu o altă metodă bazată pe adnotări. Lucrarea este concentrată asupra modului de organizare și structurare a informațiilor.

1. Introducere

Diverse studii au demonstrat necesitatea utilizărilor sistemelor de adnotare în modelarea relațiilor semantice. Sistemele de adnotare s-au dovedit utile deoarece un document adnotat a fost stocat și utilizat ulterior de același sistem sau un altul, fapt ce a permis reutilizarea și extinderea sistemului. Un astfel de sistem este Gate (Cunningham et al., 2006). Problema principală ce apare la comunicarea între două aplicații diferite prin intermediul documentelor adnotate este compatibilitatea adnotărilor. În general un document adnotat are o structură XML, auto-descriptivă. Datorită acestei flexibilități fiecare sistem își organizează datele într-un format propriu. De aici apare și problema incompatibilității adnotărilor. Pentru ca două aplicații, ce utilizează sisteme de adnotare diferite, să poată comunica este necesară translatarea unui document dintr-un sistem de adnotare în celălalt sau în altul intermediar. Această problemă se accentuează când adnotările marchează informații semantice prin intermediul referințelor.

Odată ce sistemul crește în dimensiune și trebuie să răspundă la mai multe cerințe este necesară o organizare riguroasă. Cu ajutorul unei structuri ierarhice se pot defini concepte, entități și relații între entități (Perez et al., 2005). Această soluție rezolvă unele probleme, dar introduce altele noi:

- se confruntă în continuare cu problema translatării adnotărilor;
- necesită studiul și înțelegerea sistemului creat;
- devine tot mai complicată și dificil de utilizat, pe măsură ce se introduc concepte noi.

Comunicarea cu documente adnotate seamănă cu serviciile web, ce se confruntă cu aceeași problemă: a incompatibilității celor doi participanți la comunicare. În cazul serviciilor web se folosește o soluție orientată obiect. S-au creat librării ce convertesc documentele XML în obiecte specifice limbajului de programare¹. Astfel aplicațiile au

¹ <http://java.sun.com/webservices>

posibilitatea să lucreze obiectual făcându-se abstracție de formatul mesajului din comunicare. Această tehnică are mai multe avantaje:

- introduce un standard pentru formatul mesajelor utilizate în comunicare, rezolvând problema translatării;
- se lucrează la nivel de limbaj de programare fără a introduce concepte noi;
- nivelul de dificultate se limitează la cel al limbajului utilizat;

Pentru a face față noilor cerințe impuse de aplicații pentru fișierele de configurare în format XML², s-a impus standard-ul OWL³. Fără a face o introducere în OWL, este de remarcat faptul că se apropie tot mai mult de conceptul „orientat obiect”, prin toate aspectele sale: clasa, moștenire, polimorfism. Prin faptul că se introduce conceptul de clasă și instanță, devine posibilă separarea informației de structură. Iar prin apropierea de conceptele orientate-obiect, s-a facilitat crearea de librării pentru procesarea documentelor OWL în diferite limbaje de programare. O implementare reușită este oferită gratis de HP, numele proiectului este – Jena. Utilizarea și exemplificarea conceptului orientat-obiect, reprezintă scopul de bază al lucrării. Pe lângă avantajele enunțate la serviciile web, la OWL remarcăm:

- spațiile de nume, ce permit re-utilizarea în alte contexte;
- sistem de inferențe ce permite interogarea în mod standard a ontologiei;

2. Arhitectura și funcționalitatea sistemului

Aplicația „Domain Ontology” a fost creată pentru a exemplifica practic utilizarea ontologiilor în locul adnotărilor. Sursa datelor pentru popularea ontologiei e un număr restrâns de situri ale producătorilor de camere digitale (<http://www.kodak.com>, <http://www.canon.co.uk>). Datele sunt restrânse numai la camere digitale, pentru a ușura analiza și a permite compararea rezultatelor, proiectul a fost divizat în două părți:

- crearea structurii ontologiei (s-a utilizat Protege (Horridge et al., 2004));
- popularea automată a ontologiei cu informațiile extrase de pe situri web;

Fiecare parte conține anumite etape de dezvoltare.

Partea I:

- analiza datelor tehnice legate de camere digitale;
- organizarea și gruparea datelor tehnice în clase și subclase;
- crearea ierarhiei de clase și a proprietăților claselor;
- detalierea claselor de bază cu subclase particulare datelor tehnice;
- crearea de relații între clase și ierarhii de clase;
- popularea manuală a unui șablon aferent unei pagini web cu detalii tehnice la o singură cameră, pentru a permite analiza rezultatelor;
- repetarea procesului de populare a câte unui șablon pentru fiecare producător;

Partea II:

- crearea unui procesor general pentru o pagină web;
- crearea extensiilor la procesorul general, pentru fiecare site în parte;
- crearea unui crawler general pentru un singur site;

² <http://www.w3.org/XML>

³ <http://www.w3.org/2004/OWL>

- crearea extensiilor la crawler-ul general pentru fiecare site în parte;
- crearea punții de comunicare între crawler și procesor;
- crearea punții de comunicare între procesor și modulul de populare a ontologiei;
- generarea modelului obiectual pentru ontologia creată;
- implementarea funcției abstracte „createIndividual” din clasa abstractă de bază a modelului generat, pentru a conecta acest model generat cu cel de populare.

3. Descrierea sistemului

Din analiza specificațiilor proiectului se distinge ușor modularizarea și elementele de comunicare între module.

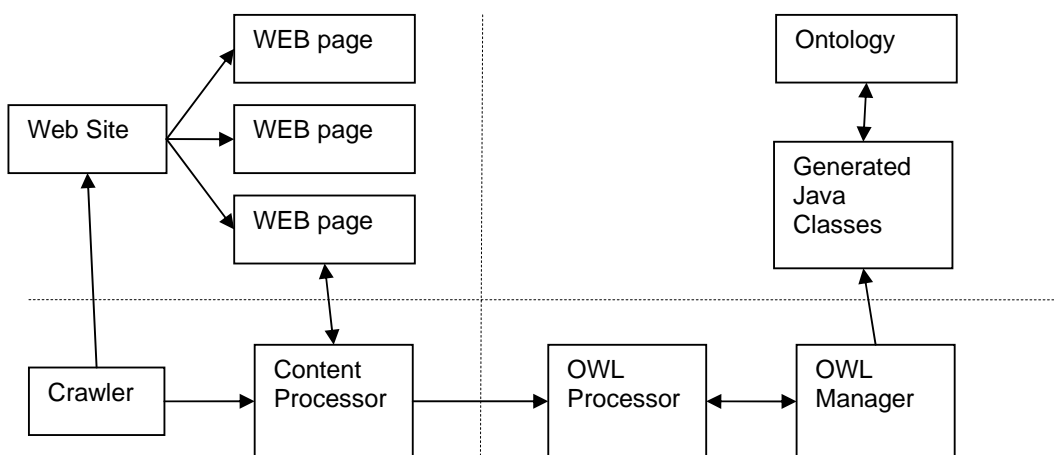


Figura 1: Arhitectura aplicației

Scopul principal impus acestui proiect a fost re-utilizarea și posibilitatea extinderii ulterioare. Arhitectura generală a aplicației este prezentată în figura 1. S-a luat în calcul organizarea tip plug-in, încât să fie posibilă adăugarea unui nou procesor pentru un alt site fără a modifica funcționalitatea celorlalte module. Iar scoaterea unui modul specific

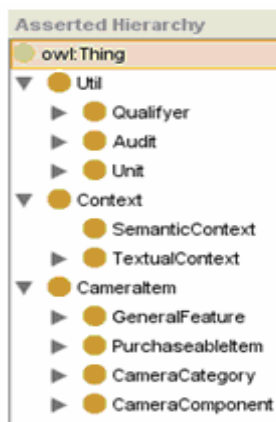


Figura 2: Clasele de bază ale ontologiei

unui site să nu afecteze funcționalitatea aplicației. Clasele aferente funcționalității generale sunt scrise astfel încât conțin operații aplicabile oricărei surse de date (site producător). Astfel, se reduce numărul operațiilor efectuate de către clasele particulare unui singur producător și totodată se micșorează efortul dezvoltării unui nou modul (pentru alt producător). Clasele particulare unui producător, efectuează operații specifice aceluși producător. În figura 1 se observă fiecare nivel; situl producătorului și modulele de cautare și extragere a conținutului, ontologia și modulele de populare automată a ontologiei. Avantajul major adus de aplicație e faptul că modulul de populare rămâne același indiferent de modulul utilizat pentru procesarea conținutului web. Acest lucru a fost posibil prin introducerea unui nivel abstract pentru procesarea conținutului web ce nu intră în detalii aferente sitului web.

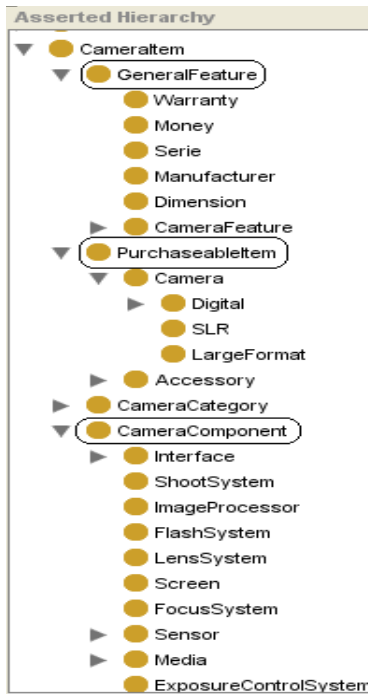


Figura 3 : Detalierea claselor de bază cu clase particulare

Ontologia a fost creată pentru a acoperi toate aspectele tehnice ale unei camere digitale. Datele tehnice au fost grupate în categorii generale (figura 2) și detaliate cu altele particulare (figura 3). Toate aspectele tehnice au o singură clasă de bază: „*CameraItem*” (figura 2). Această clasă a fost creată pentru a permite relaționarea diferitor aspecte tehnice între ele, precum și compunerea unui aspect tehnic din mai multe detalii (figura 3).

Clasa „*TextualContext*” (figura 2), subclasă a clasei „*Context*” este folosită pentru a stoca șabloane pentru fiecare producător în parte. Mecanismul e simplu: pentru fiecare producător se alege un produs și pentru acel produs se crează manual un șablon. Șablonul conține instanțe ale claselor din ontologie ce sunt folosite în procesul de populare automată a ontologiei. Avantajul utilizării acestui sistem este faptul că se crează un singur șablon pentru un singur produs și pe baza acestui șablon se procesează întreg situl producătorului.

Fiecare instanță a clase „*TextualContext*” stochează câte un șablon pentru un singur site al unui producător de camere digitale. Iar proprietatea „*cameraItemContext*” relaționează șablonul cu instanțele ce conțin informațiile

tehnice din ontologie. Fiecare instanță, a unei informații tehnice, este marcată cu o proprietate „*label*” utilizată în momentul populării ontologiei pentru a identifica tipul informației tehnice conținută în pagina web. Pe baza acestei proprietăți se face legătura între conținutul paginii web și șablon, pentru a identifica tipul informației. Iar șablonul este utilizat în continuare pentru a identifica tipul clasei pentru care se va crea un individual, ce va stoca informația din pagina web. Aici se poate observa un aspect important al ontologiilor, și anume, de a modela atât datele cât și structura lor, precum și relațiile dintre ele.

4. Rularea sistemului

Procesul de rulare se desfășoară în mai multe etape:

- configurarea directoarelor de resurse și a ontologiei cu șabloane;
- crearea și popularea manuală a unui șablon la o pagină web pentru un singur produs de la un producător de camere digitale și repetarea operației pentru fiecare alt producător;
- lansarea în execuție a aplicației;

În procesul de completare a șabloanelor se crează câte o instanță a clasei „*TextualContext*” pentru fiecare producător. Pentru fiecare site se selectează câte un produs pentru popularea șablonului. Pentru fiecare informație tehnică se crează câte o instanță a unei clase din ontologie. Tipul instanței este ales de utilizator pe baza categoriei din care face parte informația tehnică, această categorie este folosită pentru a popula proprietatea „*label*” a instanței create, iar informația propriu zisă este folosită

pentru a completa proprietățile instanței. Instanța nou creată și populată cu informații este atașată proprietății „*hasTextualContext*” pentru a face legătura cu șablonul.

Aplicația încarcă întâi ontologia cu clase, iar la pasul doi încarcă ontologia cu șabloane. Lista de șabloane este inspectată pentru a instanția crawlerul (figura 1) specific sitului producătorului. Acest crawler are singura funcție de a naviga pe site și a extrage paginile cu toate produsele oferite de producător, a le corecta și a le transforma într-un document XML valid. Paginile sunt transmise unui procesor (figura 1) particular producătorului ce are doar funcția de a extrage informațiile din pagină și a le transmite într-un format unic, procesorului pentru ontologie (figura 1). Acest ultim procesor are rolul de a popula ontologia cu instanțe ale claselor din ontologie. În procesul de populare se folosește șablonul pentru a identifica tipul instanței ce trebuie creată, iar crearea relațiilor între individuali se face apelând metodele prin reflexie. În cazul în care relaționarea trebuie făcută folosind o proprietate a unei clase de bază, se folosește inspectarea claselor de baza pentru a identifica nivelul la care se poate crea legătura.

5. *Discuții și evaluări*

Primă remarcă privitoare la rezultatele populării ontologiei este faptul că s-a urmărit crearea corectă de individuali aferenți datelor tehnice conținute în paginile web, nu și a caracteristicilor unei date tehnice. Corectitudinea populării ierarhiei „*CameraItem*” este dată de fidelitatea corespondenței template - site web producător. Problemele ce au apărut au fost legate de șabloane. Șablonul a fost creat pentru un singur produs, iar acel produs nu conține neapărat toate caracteristicile tehnice ce le poate avea un produs de la acel producător. Astfel pentru un anumit produs pot apărea caracteristici tehnice ce nu sunt cuprinse în șablon. Soluția ar putea fi crearea unei aplicații ce ar examina conținutul sitului și l-ar compara cu cel din șablon. În procesul de evaluare a rezultatelor s-a utilizat ontologia populată automat cu instanțe. Din această ontologie s-au analizat instanțele create automat pentru produsul ce a fost folosit în șablon. S-a constatat o precizie la populare de 100% pentru acel produs, însă această precizie nu s-a menținut pentru toate produsele. În funcție de producător s-a constatat o precizie de peste 90% pentru produse din aceeași categorie și o precizie sub 50% pentru produse din altă clasă. Acest procentaj se poate îmbunătăți considerabil prin popularea șablonului cu informații ale produselor din categorii diferite, modificând doar datele de intrare. O altă problemă deosebit de gravă este faptul că siturile producătorilor se modifică în timp. Pentru rezolvarea acestei probleme, s-a recurs la utilizarea unui fișier de configurare în care sunt păstrate informațiile sensibile aplicației. Pentru extragerea informației dintr-o pagină s-a folosit Xpath. Un aspect interesant este compararea ontologiei create cu o altă ontologie⁴ accesibilă pe web. Primul lucru ce poate fi observat este faptul că această ontologie a fost o sursă de inspirație. Însă ontologia nu făcea diferență între categorii de obiecte, ci doar le relaționa. Ontologia nou creată organizează mai bine datele tehnice în categorii, detaliază fiecare aspect tehnic și introduce ierarhii de proprietăți pentru a specifica relațiile între instanțe ale ontologiei. Ultimul aspect important este compararea ontologiei cu sistemul de adnotări. Pentru a pune în lumină toate diferențele vom recurge la un exemplu foarte simplu extras din ontologia populată automat și se va compara cu o

⁴ <http://protege.cim3.net/file/pub/ontologies/camera/camera.owl>

soluție posibilă bazată pe adnotărilor. Cea mai mare problemă cu care se confruntă o soluție bazată pe adnotări este necesitatea unui editor vizual care să genereze o astfel de structură. Mai important este că, în momentul în care apare necesitatea modificării structurii documentului xml, este necesară modificarea aplicației vizuale. Deși, la ora actuală există adnotatoare foarte bune, ele sunt capabile să realizeze marcări elementare. Iar problema cel mai des întâlnită apare la crearea relațiilor între etichetele ce marchează textul. În urma analizei se constată că toate problemele discutate sunt generate de structura documentelor xml, cu care se lucrează, precum și de incompatibilitățile ce apar între ele. Înlocuirea formatului xml cu OWL aduce avantajul separării conținutului de structură, iar utilizarea unui editor vizual de documente OWL ar face posibilă editarea oricărui document indiferent de structura sa internă.

6. Concluzii

Lucrarea de față prezintă un caz practic de utilizare a ontologiilor și face o comparație cu un sistem de adnotare aplicat aceleiași situații. S-au luat în discuție probleme actuale cu care se confruntă aplicațiile ce folosesc documente adnotate, utilizate în procesarea limbajului natural. Pentru fiecare problemă s-au propus soluții utilizate în aplicații ce nu țin de domeniul lingvisticii computaționale și anume, domeniul comercial, care e cel mai sensibil la probleme minore. Soluțiile propuse au fost prezentate în ordinea apariției lor, sau în ordinea în care s-au impus pe piață, cu avantajele sau dificultățile aduse. S-a descris arhitectura aplicației, fără a da detalii legate de implementare. S-a dat o atenție deosebită structurii ontologiei, precum și a modului de organizare a conceptelor, fapt foarte important în dezvoltarea unei ontologii. S-a constatat că utilizarea ontologiilor duce la eliminarea multor probleme legate de modul de structurare a unui document XML, dar introduce o problemă importantă: nu vor mai putea fi folosite vechile instrumente dezvoltate. Această problemă poate fi ușor rezolvată prin crearea unei punți pentru conversia din xml în owl și invers. Acest proiect ar putea fi punctul de plecare pentru noile cercetări în domeniul lingvisticii computaționale.

Referințe bibliografice

- Cunningham, H., Maynard, D., Bontcheva, K., Tablan, V., Ursu, C., Dimitrov, M., Aswani, M.N., Roberts, I. (2006) University of Sheffield, *Developing Language Processing Components with GATE Version 4 (a User Guide)*
- Perez, D., Postolache, O., Alfonseca, E., Cristea, D. and Rodriguez, P. (2005): *Hierarchical XML Layers Representation for Heavily Annotated Corpora* In Proceedings of the RANLP-2005 Conference, Borovets, Bulgaria, 21-23 September 2005, pp. 380-386.
- Horridge, M., Knublauch, H., Rector, A., Stevens, R., Wroe, C. (2004) A Practical Guide To Building OWL Ontologies Using The Protege-OWL Plugin and CO-ODE Tools.

CADRE PENTRU O IMPLEMENTARE PC-PATR A VERBELOR TRANZITIVE DIN LIMBA ROMÂNĂ

NADIA LUIZA HUȚULIAC

Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București

hnadia_luiza@hotmail.com

Rezumat

Ne propunem să prezentăm din perspectiva formalismului lingvistic **PATR** diferite situații gramaticale întâlnite în analiza verbelor tranzitive: structuri verbale bivalente și trivalente; acordul trăsăturilor morfologice între complementul direct și forma neaccentuată a pronumelui personal, care îl anticipează; introducerea informației semantice în structurile de trăsături asociate grupului verbal. Realizarea unor reguli complete de implementare folosind PC-PATR este o sarcină uriașă, deoarece există foarte multe posibilități în limbajul real de care cercetătorul nu poate ține cont în a analiza fiecare propoziție. Pentru un grup de propoziții însă, această implementare se dovedește utilă deoarece arborii sintactici rezultați pot fi transformați ulterior în exemplele de traducere ale unei baze de traducere.

1. Introducere

Scopul acestei cercetări este de a furniza o implementare **PC-PATR**¹ pentru verbele tranzitive din limba română, ca punct de plecare pentru un proiect ulterior de extragere a cunoștințelor de traducere dintr-un corpus paralel, aliniat propozițional. Tehnologiile lingvistice sunt puse la dispoziție de către **SIL** International și se caracterizează prin modularitate, deoarece instrumentul **CARLASTUDIO**, abreviere pentru *Computer-Assisted Related Language Adaptation*, folosește o gramatică generată de sistemul expert **PAWS** (*Parser and Writer for Syntax*) și un lexicon creat cu ajutorul programului **ToolBox**.

2. Tehnologii lingvistice

CarlaStudio² (**CS**) este un program cu dublă funcționalitate: pe de o parte, permite modelarea unui limbaj anume de către lingvistul cercetător, iar pe de altă parte, implică modelul lingvistic creat în analiza textelor sau în adaptarea lor pentru un alt limbaj. Versiunea **2.9.0.4 Unicode** utilizată include trei programe importante pentru scopul cercetării noastre, fiecare cu o funcție specifică de procesare: analizorul morfologic **AMPLE**, analizorul sintactic **PC-PATR**, responsabil de dezambiguizarea

¹ <http://www.sil.org/pcpatr/>

² <http://www.sil.org/>

și crearea arborilor sintactici, **JOINCOMP**, instrumentul ce recunoaște cuvintele compuse sau locuțiunile morfologice.

Autorii sistemului **CS** (John Hatton, Andy Black, Bob Eaton, Marius Doornenbal) au gândit modelarea limbajului ca o gestionare sinergică a mai multor tipuri de informație lingvistică. Diacriticele din sistemul românesc sunt tratate la primul nivel de procesare, unde există posibilitatea de a descrie asociații de caractere dependente de limbă. Categoriile gramaticale aparțin nivelului de analiză și sunt adaptate după recomandările sistemului expert **PAWS**, astfel încât co-există categorii tradiționale, de tipul : nume, verb, pronume, adverb, adjectiv, și categorii adaptate- cuantificator, demonstrativ, auxiliar.

CS are la bază paradigma Analiză-Transfer-Sinteză. Analiza e focalizată pe morfologie și fonologie. Din cauza ambiguității analizei cuvintelor independent de contextul lor, s-a urmărit modelarea construcțiilor sintactice prin folosirea ordinii permise a lexemelor în sintagme și construcții gramaticale speciale, prin marcarea trăsăturilor de acord și a celor de entități numite.

Una din principalele motivații legate de crearea sistemului expert **PAWS** a fost de a adăuga un instrument de dezambiguizare bazată pe sintaxă pentru instrumentele de analiză morfologică existente în **CS**. În scenariul propus de Andy Black și Cheryl A. Black, utilizatorii folosesc un lexicon de morfeme pe care programul morfologic **AMPLE** le utilizează pentru a analiza cuvintele dintr-un text în constituenți. Rezultatul poate fi avansat apoi către **PC-PATR**, inclus în **CS**, împreună cu un fișier de gramatică.

Ultimul program utilizat în dezvoltarea implementării **PATR** a verbelor tranzitive este lexiconul **Toolbox**, responsabil de crearea unui fișier cu extensia **.lex**, apelat de **CarlaStudio** în analiza cuvintelor. Primul câmp (introdus prin `\w`) furnizează forma grafică a cuvântului, al doilea (`\c`) introduce categoria morfo-sintactică a intrării, al treilea (`\g`) este glosa, al patrulea (`\f`) conține trăsăturile de subcategorizare ale intrării lexicale. Ultimul câmp permite o evidență a actualizărilor datelor introduse de lingvist.

3. *Morfologia verbului*

Modurile *indicativ* și *declarativ* sunt sinonime în contextul **PAWS**, dar opțiunea pentru lexemul *indicativ* s-a realizat prin scrierea tipului corespunzător:

```
(1) Let indicative be <head infl mood indicative> = +
```

Modurile imperativ și conjunctiv există în descrierea morfologică **PAWS**, iar condiționalul a fost introdus pentru limba română:

```
(2) Let conditional be <head infl mood conditional> = +
```

Pentru modul indicativ am adăugat timpurile imperfect și mai-mult-ca-perfect, iar timpurile compuse au fost definite prin reguli și restricții particulare:

```
(3) Let imperfect be <head infl tense imperfect> = +
```

```
(4) Let pluperfect be <head infl tense pluperfect> = +
```

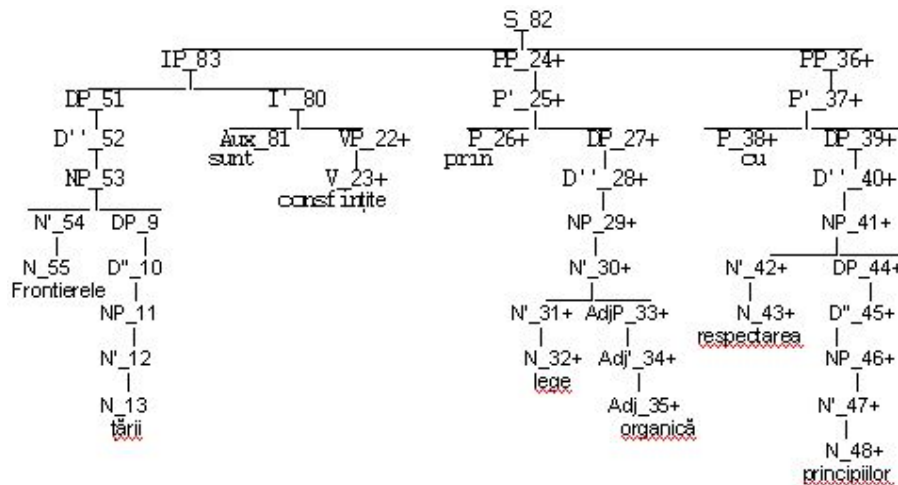
4. Structuri verbale de subcategorizare

4.1. Reguli sintagmatice de expansiune a grupului verbal

Regulile sintagmatice de expansiune a grupului verbal înlocuiesc simbolul unic al sintagmei verbale printr-unul sau mai multe simboluri, rezultând două reguli generale, diferențiate după valența verbului:

- (A) VP = V DP³
- <VP head> = <V head>
 - <V head object> = <DP>
 - <V head type transitive> = +
 - <V head type copular> = -
 - <V head type passive> = -
- (B) VP = V DP₁ DP₂
- <VP head> = <V head>
 - <V head object> = <DP₁>
 - <V head type ditransitive> = +
 - <DP₁ head case> = accusative
 - <DP₂ head case> = dative

În figura următoare, pot fi urmărite relațiile de dependență sintactică stabilite între nodurile fiică și mamă din perspectiva regulilor de expansiune a grupului verbal, cu nucleu pasiv și determinări prepoziționale:



Cazul special al anticipării obiectului direct prin forme pronominale neaccentuate este rezolvat prin introducerea a două trăsături gramaticale- **Case**, respectiv, **ObjAgr**- și prin

³ Ca terminologie, am păstrat structura mai extinsă **Determiner Phrase** (DP), motivul constituindu-l existența numeroaselor poziții pentru tipuri diferite de modificatori ai numelui.”

modificarea regulii de expansiune a grupului nominal. Deși marcată prepozițional, sintagma cu funcția de obiect nu e considerată grup prepozițional, ci nominal, unificându-se trăsăturile de acord și de caz:

```

Let Case be <cat> = Case | for pe
    <head case> = accusative
    <head type proper> = +

Let ObjAgr be <cat> = ObjAgr
    <head type proper> = + |for whole set
    of object agreement pronouns that go
    before the aux and/or verb

rule {IBar option3- ObjectAgr initial}
I' = ObjAgr VP
    <I' head> = <VP head>
    <I' head type auxiliary> = +
    <I' head object head agr> = <ObjAgr head agr>
    <I' head object head type proper> = +
    <I' head type transitive> = +
    <I' head type prefix> = <ObjAgr head type prefix>

rule {DP option Case- prepositional marker for DO}
DP = Case DP_1
    <DP head> = <DP_1 head>
    <DP head type proper> = +
    <Case head case> = accusative
    <DP head case> = accusative
    <DP option> = Case
    
```

4.2. Roluri semantice

Rolurile semantice nu interesează în mod strict gramatica generată de PAWS, însă ele pot fi adăugate șabloanelor de descriere, cu observația că există o diferență între a indica rolul complinit de un verb printr-un obiect și a predetermina rolurile unui substantiv într-o structură de subcategorizare.

Codificarea informației semantice în structura de trăsături se realizează îndeosebi în lexicon, deoarece fiecare verb va avea un cadru de subcategorizare, iar argumentele trebuie să aibă asignate fiecare un rol semantic potrivit. Aceste trăsături introduse pentru intrările lexicale vor avea specificate tipuri care să le modifice în structuri de trăsături ce se vor unifica apoi corect prin regulile sintagmatiche.

Inițial, gramatica nu includea și roluri semantice, iar descrierea tipului de tranzitivitate nu prezenta o structură specifică de subcategorizare. Pentru nominale, soluția cea mai

CADRE PENTRU O IMPLEMENTARE PC-PATR A VERBELOR TRANZITIVE DIN LIMBA ROMÂNĂ

satisfăcătoare este să se marcheze ce roluri nu pot fi îndeplinite, de exemplu- rolul de AGENT sau de EXPERIENCER pentru obiecte inanimate:

```
|templates for semantic roles on nominals
Let -AGENT be <head role AGENT> = -
                <head type animate> = -
Let -EXP be <head role EXPERIENCER> = -
                <head type animate> = -
```

Pentru o structură verbală bitranzitivă, perechea de roluri tematice AGENT- THEME este evidențiată prin următoarele declarații de descriere :

```
Let AGENT_THEME be <head subject head role AGENT> = +
                <head object head role THEME> = +
                <head indirectobject> = none
```

Cadrele de subcategorizare sunt foarte complexe. În general, schimbările la nivelul regulilor nu sunt recomandate pentru fișierul de gramatică, pentru că trăsăturile nou introduse vor unifica apoi cu alte trăsături, măbind numărul de analize și de arbori de analiză și, implicit, timpul de procesare sintactică.

5. Concluzie

Rezultatele recente din NLP datorează mult metodelor statistice sau celor bazate pe corpus. Traducerea automată nu face excepție de la această afirmație (Carl & Way, 2003), deoarece paradigmele traducerii bazate pe exemple și cea statistică folosesc lexicoane și reguli de traducere achiziționate din corpusul paralel și compilate apoi în motorul de traducere. Conceput inițial ca o analiză PC-PATR a verbelor tranzitive, proiectul nostru va fi dezvoltat în continuare cu o dublă funcționalitate. Pe de o parte, intenționăm realizarea unor modele lingvistice pentru română și engleză prin procesarea unui corpus paralel, aliniat propozițional. Pe de altă parte, vom utiliza modelele lingvistice PC-PATR pentru a obține cunoștințe de traducere, exprimate în unități lexicale sau sintagmatice echivalente, prin care să fie validată traducerea automată.

Referințe bibliografice

- Black, C., (1997). A PC-PATR Implementation of GB Syntax. *SIL Electronic Working Papers 1997-0006*.
- Carl, M., Way, A. (eds.) (2003). *Recent advances in Example-Based Machine Translation*. Kluwer Academic Publishers, Netherlands.

Index de autori

- Aldea, Bogdan-Mihai: 45
Apopei, Vasile: 9, 107
Barbu Mititelu, Verginica: 17
Bejinariu, Silviu: 107
Bîrlădeanu, Antonina: 35, 119
Bobicev, Victoria: 23
Boian, Elena: 75, 135
Bolea, Cecilia: 123
Bomher, Noemi: 161
Botoșineanu, Luminița: 107
Bozianu, Luigi: 17
Brut, Mihaela: 113, 155
Burciu, Natalia: 35, 119
Burlaca, Oleg: 75
Ceaușu, Alexandru: 17
Chiorescu, Adrian: 29
Cibotaru, Radu: 165
Ciocârlău, Cristina: 155
Ciubotaru, Constantin: 75, 135
Cojocar, Svetlana: 75, 135
Colesnicov, Alexandru: 75
Cristea, Dan: 51, 83, 101, 129
Curteanu, Neculai: 123, 143
Demidov, Valentina: 75
Diaconescu, Ștefan: 39
Dornescu, Iustin: 123
Elița, Natalia: 63
Feraru, Monica: 3
Florescu, Cristina: 149
Forăscu, Corina: 51, 69, 83, 129
Gavrilă, Monica: 63
Haja, Gabriela: 45
Huțuliac, Nadia Luiza: 171
Iaciurinschi, Alina: 23
Iftene, Adrian: 51, 83, 129
Ion, Radu: 69
Irimia, Dumitru: 113
Irimia, Elena: 57, 89
Jitcă, Doina: 9
Luca, Ramona: 107
Magariu, Galina: 135
Malahova, Ludmila: 75
Manu Magda, Margareta: 17
Maxim, Victoria: 23
Mihăilă, Cătălin: 17
Moruz, Alex: 123, 143
Olariu, Florin: 107
Panait, Oana: 113
Pavel, Gabriela: 101
Pistol, Ionuț: 51, 83, 101, 129
Postolache, Oana: 101
Pușcașu, Georgiana: 83
Rogojin, Iuri: 135
Ștefănescu, Dan: 89
Teodorescu, Horia-Nicolai: 3
Todirașcu, Amalia: 95
Todoroi, Dumitru: 29
Trandabăț, Diana: 3, 51, 83, 123, 129, 143
Tufiș, Dan: 17, 57, 89
Verlan, Tatiana: 135
Zidrașco, Tatiana: 23

