

**eLEXICOGRAFIA – INTERDISCIPLINARITATEA CA PREMISĂ
PENTRU CERCETAREA LIMBII ROMÂNE¹**
*eLexicography – Interdisciplinarity as a Premise for the Research of
Romanian Language*

**Scient. Researchers Dr. Elena TAMBA,
Scient. Researchers Dr. Marius CLIM,
Scient. Researchers Dr. Ana CATANĂ-SPENCHIU,
Scient. Researchers, PhD Mădălin PĂTRAȘCU
“A. Philippide” Institute of Philology of the Romanian Academy,
Iași Branch**

Abstract

During the last years, the Romanian linguistic research in general and the lexicographic one in particular have known an intensive growth in the natural language processing area. In this direction, interdisciplinary steps have been taken to provide Romanian electronic dictionaries and text corpora.

For researchers, the Romanian Thesaurus Dictionary in electronic format (eDTLR) and the Essential Romanian Lexicographic Corpus (ERLC) will provide important and awaited tools for the Romanian language study.

Keywords: *interdisciplinarity, digitalized_lexicography, lexicographic_corpus, linguistic_resources, e-lexicography*

În ultimii ani, atât în domeniul lexicografiei, cât și în cel al informaticii aplicate s-a simțit nevoia unei evoluții interdisciplinare, ca urmare a cerințelor de dezvoltare pe cât mai multe direcții. Astfel, într-o primă fază la propunerea cercetătorilor informaticieni cu preocupări în domeniul prelucrării limbajului natural, au fost inițiate unele proiecte comune, interdisciplinare, care au dus la realizarea unor corpusuri de limbă română sau a unor programe care să ajute la prelucrarea textelor în limba română, cu rezultate eficiente în domeniul cercetării acesteia.

Istoricul proiectelor de tip academic, care au presupus o interdisciplinaritate sporită, începe cu demersurile de informatizare a *Dicționarului Tezaur al limbii române*, care reprezintă cea mai importantă lucrare lexicografică apărută sub egida Academiei Române – editarea sa a început acum 105 ani; a apărut în două serii: DA (1907-1944), DLR (1965-2010), fiind publicat în 14 tomuri / 37 volume, totalizând cca 18.000 pagini și mai mult de 175.000 intrări, cu tot cu variante.

Astfel, începând cu anul 2005, la Iași, s-au desfășurat unele proiecte interdisciplinare, în care au fost implicați lexicografii de la Institutul de Filologie Română „A. Philippide”, Iași, lingviști de la Facultatea de Litere de la Universitatea „Alexandru Ioan Cuza” din Iași și informaticienii de la Facultatea de Informatică, Universitatea „Alexandru Ioan Cuza” din Iași și de la Institutul de Informatică Teoretică, Academia Română, Iași:

¹ *Acknowledgement:* Acest articol a fost realizat în cadrul proiectului *CLRE. Corpus lexicografic românesc esențial. 100 de dicționare de bază din Bibliografia DLR aliniate la nivel de intrare și la nivel de sens*, CNCS-UEFISCDI, cod TE_246/2010, 2010 – 2013TE_246 (55/2010).

a) *Dicționarul limbii române (DLR) în format electronic. Studii privind achiziționarea* (cod CNCISIS 1815), grant finanțat de Ministerul Educației și Cercetării (MEC) prin Consiliul Național al Cercetării Științifice din Învățământul Superior (CNCISIS), desfășurat în perioada 2003–2005. Prin acest proiect s-a verificat și demonstrat posibilitatea transformării *Dicționarului limbii române* din text tipărit în text electronic adnotat², prelucrat cu ajutorul unui program specific, DLReX – un instrument de achiziționare, prelucrare și consultare a DLR, bazat pe o euristică prin care sunt recunoscute diferitele câmpuri formale ale textului unui articol, putându-se identifica automat textul definițiilor, al citatelor sau al siglelor³.

b) *Resurse lingvistice în format electronic: Monumenta linguae Dacoromanorum. Biblia 1688. Regum I, Regum II – Ediție critică și corpus adnotat. (MLD. Biblia 1688)* (cod CNCISIS 1454), desfășurat în perioada 2006–2007. Prin acest proiect a fost găsită o posibilă metodă de achiziționare în format electronic a unor cărți vechi din Bibliografia DLR, cu aplicație asupra a două cărți din *Biblia* tipărită la București în anul 1688, *A împărățiilor cea dentâiu*, *A împărățiilor a doua*, precum și crearea unor instrumente de indexare și adnotare automată, la nivel de cuvânt, a textelor românești vechi.

c) *DLRI. Bază lexicală informatizată. Derivate*. (cod CNCISIS nr. 1609), desfășurat în perioada 2007–2008. Prin acest proiect s-a prelucrat un eșantion lexicografic format din derivatele pe terenul limbii române cu sufixul *-ime* – de origine latină, și cele cu *-iște* – de origine veche slavă, din seria veche (DA) și din seria nouă a dicționarului-tezaur (DLR), și s-a demonstrat posibilitatea unificării tehnico-lexicografice a articolelor DA – DLR, prin mijloace informatice.

Aceste trei proiecte interdisciplinare au anunțat și pregătit realizarea unui proiect cu implicare națională, care a vizat realizarea variantei electronice a *Dicționarului Tezaur al limbii române* – eDTLR. Echipa proiectului a fost coordonată de Facultatea de Informatică, Universitatea „Alexandru Ioan Cuza” din Iași, prin prof. dr. Dan Cristea și a inclus ca parteneri: Institutul de Lingvistică „Iorgu Iordan – Al. Rosetti”, Academia Română, București; responsabil de proiect acad. Marius Sala (prin dr. Monica Busuioc); Institutul de Filologie Română „A. Philippide”, Academia Română, Iași; responsabil de proiect dr. Gabriela Haja; Institutul de Lingvistică și Istorie Literară „Sextil Pușcariu”, Academia Română, Cluj-Napoca; responsabil de proiect dr. Rodica Marian; Institutul de Cercetări pentru Inteligență Artificială, Academia Română, București; responsabil de proiect acad. Dan Tufiș; Institutul de Informatică Teoretică, Academia Română, Iași; responsabil de proiect acad. Horia Neculai Teodorescu; Facultatea de Litere, Universitatea „Alexandru Ioan Cuza” din Iași; responsabil de proiect dr. Eugen Munteanu.

Acest proiect⁴ a beneficiat inițial de o finanțare națională din partea CNMP, iar ulterior, după decembrie 2010, lucrându-se la finalizarea site-ului eDTLR doar cu ajutorul unor voluntari – cercetători informaticieni și lexicografi. Proiectul menționat, după încheierea sa, ar trebui, pe de o parte, să pună la dispoziția tuturor cunoscătorilor sau celor interesați de

² Textul electronic adnotat este un text analizat și marcat din punct de vedere formal astfel încât să poată fi consultat, corectat, modificat etc. de către specialiștii lexicografi, cu ajutorul calculatorului. Există posibilitatea extragerii din formatul complet a unei forme destinate numai consultării, care să se adreseze unui public mai larg decât cel al specialiștilor propriu-ziși. Pentru detalii, vezi și Haja, Dănilă *et alii*, 2005.

³ Pentru o prezentare detaliată a rezultatelor acestui proiect, vezi Haja, Dănilă *et alii* 2005.

⁴ Pentru o prezentare detaliată a proiectului, vezi Dănilă 2010: 37–46.

limba română formatul electronic al *Dicționarului Academiei*, pe suport electronic – și, poate ulterior, în funcție de schimbarea politicii lingvistice din România, și on-line, în acces liber sau condiționat – și, pe de altă parte, să pună la dispoziția deocamdată doar a cercetătorilor o arhivă electronică care să cuprindă toate textele din Bibliografia DLR. Prin rezultatele sale, eDTLR ar pune lexicografia fundamentală pentru limba română într-o situație de egalitate cu cea a limbilor care au deja dezvoltate astfel de resurse: *Le Trésor de la Langue Française Informatisé* (TLFi – <http://atilf.atilf.fr/>); *Diccionario de la lengua española de la Real Academia Española* (DRAE – <http://buscon.rae.es/draeI/>); *Tesoro della lingua italiana delle origini* (TLIO – <http://tlio.oiv.cnr.it/TLIO/index2.html>); *Deutsches Wörterbuch der Grimm* (DWB – <http://germazope.uni-trier.de/Projects/DWB>); *Oxford English Dictionary* (OED – <http://www.oed.com/>) ș.a.

Pornind de la eDTLR a apărut necesitatea realizării unui alt proiect interdisciplinar: crearea unui corpus lexicografic românesc – *CLRE. Corpus lexicografic românesc esențial. 100 de dicționare din bibliografia DLR alinate la nivel de intrare și la nivel de sens*, care este finanțat de CNCS-UEFISCDI pentru perioada 2010 – 2013 și care se desfășoară tot la Institutul de Filologie Română „A. Philippide” din Iași.

Cu ajutorul unei echipe formate din doar 4 cercetători (3 lexicografi – Elena Tamba, Marius-Radu Clim și Ana-Veronica Catană-Spenchiu – și un informatician – Marius Răschip) și cu aportul unui alt informatician voluntar – Mădălin Pătrașcu, proiectul își propune să realizeze o bază de date care să cuprindă dicționarele esențiale din Bibliografia DLR, alinate la nivel de intrare și parțial la nivel de sens; să construiască un mediu de programe care să permită consultarea interactivă a acestui corpus și care să se constituie într-un cadru modern de lucru și cercetare lexicografică, ușor adaptabil la o diversitate de obiective; să realizeze o listă de cuvinte cvasi-exhaustivă, pentru limba română, pornind de la corpusul aliniat. Astfel, prin acest demers se vizează: realizarea unui corpus scanat, format din dicționarele de referință ale DLR (cu respectarea legislației în vigoare în ceea ce privește drepturile de proprietate intelectuală); scanarea și prelucrarea (OCR-izarea⁵; parsarea⁶ textului la nivel de intrare și, parțial, la nivel de sens) a acestor dicționare; realizarea unei interfețe on-line pentru validarea/corectarea parsării, precum și validarea alinierii între textul *Dicționarului Tezaur al limbii române* (în format electronic, rezultat al proiectului eDTLR) și dicționarele de referință din Bibliografia DLR. Dicționarele alese au fost incluse în trei categorii:

1. *dicționare generale*, de tipul:

DA = *Dicționarul limbii române*, tom I-II, Tipografia ziarului „Universul”, București, Imprimeria Națională, 1907-1944;

DLR = *Dicționarul limbii române*, Serie nouă, tom VI-XIV, București, Editura Academiei, 1965-2010;

DEX = *Dicționarul explicativ al limbii române*. București, Editura Academiei, 1975;

2. *dicționare auxiliare* (care sunt strâns legate de redactarea *Dicționarului Tezaur*), de tipul:

⁵ Convertirea din format imagine în format text.

⁶ Identificarea automată a intrărilor din dicționarele scanate și OCR-izate anterior.

Alexandru Ciorănescu, *Dicționarul etimologic al limbii române*. Ediție îngrijită și traducere din limba spaniolă de Tudora Sandru-Mehedinți și Magdalena Popescu Marin. București, Editura Saeculum I. O., 2002.

*** *Dicționarul ortografic, ortoepic și morfologic al limbii române*. Ediția a II-a revăzută și adăugită, București, Univers Enciclopedic, 2005.

Florin Marcu, *Noul dicționar de neologisme*. București, Editura Academiei Române, 1997.

3. *dicționare speciale* (enciclopedice ori dicționare speciale, alese după criteriul importanței lor pentru perspectiva diacronică asupra limbii), de tipul:

Dicționar enciclopedic. [Vol.] I: A–C (1993), [vol.] II: D–G (1996), [vol.] III: H–K (2000), [vol.] IV: L–N (2001), [vol. V]: O–Q (2004). [vol.] VI: R–S (2006). [vol.] VII: T–Z (2009). București, Editura Enciclopedică;

I.-Aurel Candrea – Gh. Adamescu, *Dicționarul enciclopedic ilustrat. Partea I: Dicționarul limbii române din trecut și de astăzi* de I.-Aurel Candrea. *Partea II: Dicționarul istoric și geografic universal* de Gh. Adamescu. București, Editura Cartea Românească, [1926–1931];

Lexiconul tehnic român. I ș. u. Elaborare nouă. București, Editura Tehnică, 1957 ș. u.

În proiect se utilizează, astfel, atât metode lingvistice clasice / tradiționale (de exemplu, transliterarea intrărilor în alfabet chirilic sau de tranziție ori studiul comparativ, la nivel semantic, al dicționarelor), cât și metode noi, de lexicografie computațională.

Pe plan internațional, există deja astfel de corpusuri, care facilitează foarte mult cercetarea atât în domeniul studiului limbii / limbilor, cât și în cel al informaticii aplicate, al prelucrării limbajului natural:

– *Le rayon des dictionnaires* (<http://www.atilf.fr/>) – colecție de dicționare informatizate franceze, din secolul al XVI-lea până în secolul al XX-lea;

– *Nuevo tesoro lexicográfico de la lengua española* (<http://buscon.rae.es/ntlle/SrvltGUILoginNtlle>) – bază de date cuprinzând versiunile facsimilate ale tuturor dicționarelor editate și publicate de Real Academia Española;

– *Das Wörterbuchnetz* (<http://germazope.uni-trier.de/Projects/WBB/>) – rețea de dicționare de limba germană, creată la universitatea Trier din Germania.

Pentru realizarea corpusului de dicționare (CLRE) este nevoie de parcurgerea mai multor etape care necesită o adaptare în funcție de specificul lucrărilor lexicografice avute în vedere. În primul rând este indispensabilă transpunerea în format electronic a celor 100 de dicționare. Acest lucru presupune scanarea și prelucrarea textelor. Astfel se impune utilizarea atât a unor echipamente performante care să faciliteze achiziționarea electronică a dicționarelor, cât și a unor programe de prelucrare a scanărilor și de recunoaștere a caracterelor care să permită realizarea în bune condiții a bazei de date. Pentru aceasta s-a achiziționat un scanner special pentru cărți, Atiz BookDrive DIY⁷. Varianta aleasă s-a dovedit a fi cea mai potrivită soluție de digitalizare a cărților, din perspectiva costurilor și a eficienței, aceasta din urmă fiind dată de performanțele camerelor digitale SLR și de suportul de carte în forma literei V. Scannerul Atiz este prevăzut cu două camere Canon EOS 450D cu lentile de 35 mm. Lentilele EF 35 mm permit o focalizare mai mare și sunt folosite în special pentru cărțile

⁷ Mai multe informații despre acest produs sunt disponibile pe site-ul <http://diy.atiz.com/>.

în format A3 sau A2. Scannerul poate fi folosit pentru toate tipurile de cărți, indiferent de dimensiune, grosime sau tipul de legătură. În plus, este ideal pentru scanarea și prelucrarea cărților vechi care necesită o atenție aparte.

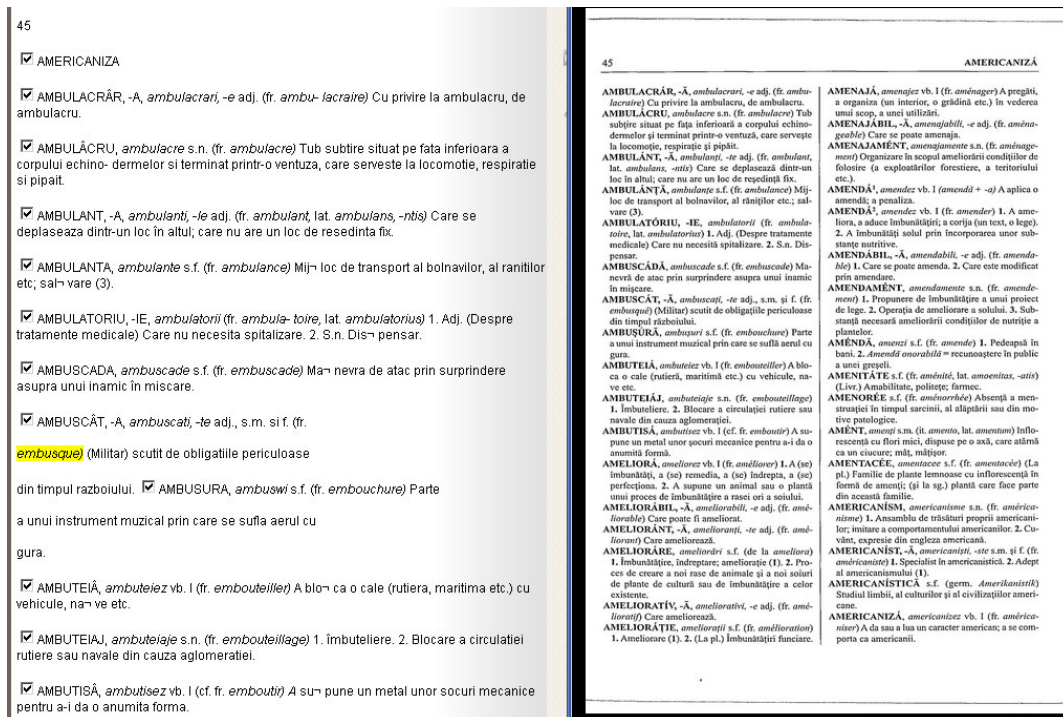
Cea de-a doua etapă presupune transformarea imaginilor scanate în text editabil. Pentru recunoașterea optică a caracterelor din imaginile scanate și pentru procesarea de text a fost achiziționat ABBYY FineReader Engine, care înglobează atât sistemul de recunoaștere optică a caracterelor (OCR), cât și recunoașterea inteligentă a caracterelor (ICR), recunoașterea mărcii optice (OMR), recunoașterea codurilor de bare (OBR), procesare de imagini și conversia în format .pdf. Acest program permite convertirea rapidă și cu acuratețe a documentelor scanate, a fișierelor în format .pdf sau a documentelor în format imagine într-o varietate de formate Office în care se pot realiza căutări și care sunt ușor de editat.

Toate aceste echipamente și programe facilitează o prelucrare informatică de mare acuratețe a materialului lexicografic avut în vedere.

Cel de-al treilea pas a presupus realizarea unei interfețe de validare care să permită verificarea corectitudinii parsării și validarea intrărilor. Așadar, dicționarele scanate sunt introduse într-o bază de date și, după parsare, sunt validate de către lexicografi. Interfața de validare permite vizualizarea fiecărei pagini de dicționar în parte. Pentru fiecare dicționar este afișată pagina-titlu, iar lexicograful poate insera toate informațiile despre dicționar (autor, editura etc.). În baza de date dicționarele sunt identificate după sigla din bibliografia DLR.

Lexicografii validează fiecare pagină, iar în cazurile în care parsarea s-a realizat cu erori, pagina respectivă este validată doar parțial, urmând a fi reprelucrată și trimisă din nou spre validare.

Paginile validate corect au toate intrările recunoscute, așa cum se poate observa și în imaginea următoare.

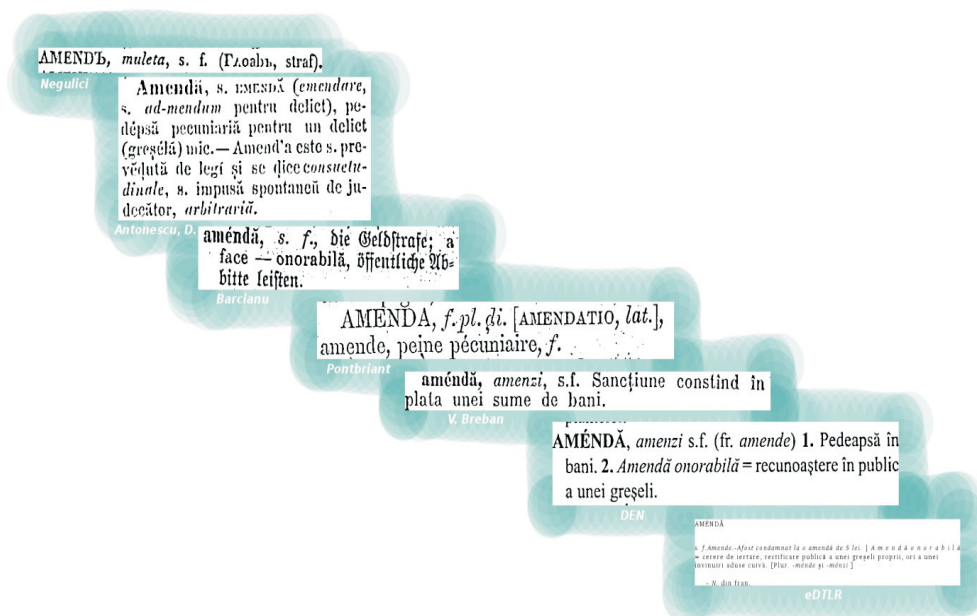


În această etapă fiecare dicționar trebuie tratat separat întrucât există caracteristici specifice în ceea ce privește formatul cuvântului-titlu, felul în care autorul dicționarului a delimitat fiecare intrare în parte, care este ordinea informațiilor în cadrul intrării, cum sunt prezentate definițiile etc. Dificultățile întâmpinate în această fază sunt determinate și de cantitatea mare de material care trebuie să fie tratat din punct de vedere informatic – parsat și aliniat (cca 100 de dicționare / aproximativ 150.000 de pagini de dicționar), dar mai ales de situația dicționarelor scrise în alfabet chirilic sau în alfabet de tranziție.

În cazul dicționarelor în alfabet chirilic și a celor în alfabet de tranziție s-a optat pentru introducerea manuală de către lexicografi, a cuvintelor-titlu, în interfața de validare. Aceasta s-a datorat erorilor de recunoaștere automată a caracterelor, întrucât nu există încă un program de OCR-izare pentru texte vechi românești. Dat fiind faptul că grafia chirilică românească este destul de diferită pentru texte din epoci și surse variate, este practic imposibil, pentru moment, să se recunoască automat un text vechi românesc. De aceea s-a renunțat la recunoașterea automată a textului din dicționarele de acest tip, optându-se pentru o validare manuală făcută de lexicografi, care, practic, au atașat o „etichetă” cu transcrierea în alfabet latin a cuvintelor-titlu. În plus, dicționarele vechi vor fi aliniate la nivel de imagine și nu de text, fiind afișată în format imagine doar porțiunea aferentă cuvântului căutat.

Ultima etapă din acest proiect vizează alinierea la nivel de intrare și, parțial, la nivel de sens a informației din dicționarele vizate.

La sfârșitul proiectului, ne dorim să oferim utilizatorului posibilitatea de a vizualiza, pentru un anumit cuvânt, toate intrările corespunzătoare din cele 100 de dicționare, într-o structură de tipul celei prezentate mai jos.



De asemenea, pentru alinierea parțială la nivel de sens, corelațiile vor fi făcute în funcție de definițiile din DLR, folosindu-se astfel rezultatele din proiectul eDTLR.

Prin acest proiect se are în vedere, aşadar, obţinerea unor rezultate care vor permite ulterior dezvoltarea de aplicaţii de anvergură privind dezambiguizarea semantică a cuvintelor, selecţii de tipuri de intrări în vederea elaborării de noi dicţionare specializate (tematice, etimologice etc.), corelarea cu alte resurse lingvistice ori multimedia, ceea ce ar aduce lexicografia românească la un nivel comparabil cu lexicografia europeană sau chiar mondială. CLRE reprezintă, astfel, şi un punct de plecare pentru cercetări viitoare, cu precădere în domeniile interdisciplinare.

În lexicografia românească actuală se pot observa câteva tendinţe care, în general, implică activităţi interdisciplinare şi care pot fi grupate astfel:

- se continuă realizarea (cu termen neprecizat) a unor ediţii „clasice” (v. continuarea DLR, reeditarea DEX etc.);
- se lucrează (deocamdată nesistematic) la realizarea unor corpusuri de texte pentru limba română;
- se realizează corpusuri lexicografice (vezi proiectul CLRE);
- a început utilizarea sistemelor de scriere de dicţionare.

Concluzii

Rezultate ale unor demersuri interdisciplinare, versiunea informatizată a *Dicţionarului Academiei* (eDTLR) şi corpusul lexicografic românesc esenţial (CLRE) vor facilita accesul specialiştilor la nişte instrumente de lucru indispensabile, mult timp aşteptate, foarte utile pentru studiul limbii române.

Bibliografie

- OED = *Oxford English Dictionary* – <http://www.oed.com/>.
DWB = *Deutsches Wörterbuch “der Grimm”* – <http://germazope.uni-trier.de/Projects/DWB>.
TLFI = *Le Trésor de la Langue Française Informatisé* – <http://atilf.atilf.fr/>
TLIO = *Tesoro della lingua italiana delle origini* – <http://tlio.oiv.cnr.it/TLIO/index2.html>
- Atkins, B.T.S.; Rundell, M. 2008: *The Oxford Guide to Practical Lexicography*. Oxford, Oxford University Press.
- Clim, Dănilă *et alii* 2008: Marius Clim, Elena Dănilă, Gabriela Haja, *Premise ale informatizării cercetării lexicografice academice româneşti* în volumul *Limba română. Dinamica limbii, dinamica interpretării*, Editura Universităţii din Bucureşti, p. 585 – 591.
- Cristea, Răschip *et alii* 2007: Dan Cristea, Marius Răschip, Corina Forăscu, Gabriela Haja, Cristina Florescu, Bogdan Aldea, Elena Dănilă, *The Digital Form of the Thesaurus Dictionary of the Romanian Language*, în vol. *Advances in Spoken Language Technology* (editors Corneliu Burileanu, Horia-Nicolai Teodorescu), Bucureşti, Editura Academiei Române, p. 195-206.
- Dănilă 2010a: Elena Dănilă, *eDTLR – base de données et instrument pour la recherche lexicographique roumaine*, în „*Philologica Jassyensia*”, An VI, Nr. 1 (11), 2010, p. 37–46.
- Dănilă 2010b: Elena Dănilă, *Despre necesitatea realizării unui corpus lexicografic românesc esenţial*, în „*Philologica Jassyensia*”, An VI, Nr. 2 (12), 2010, p. 41-50.

Haja, Dănilă *et alii* 2005: Gabriela Haja, Elena Dănilă, Corina Forăscu, Bogdan-Mihai Aldea, *Dicționarul limbii române (DLR) în format electronic. Studii privind achiziționarea*, Iași, Editura Alfa, publicat și electronic pe www.consilr.info.uaic.ro.

Tamba Dănilă, Clim *et alii* 2012: Elena Tamba Dănilă, Marius-Radu Clim, Mădălin Pătrașcu, Ana Catană-Spenchiu, *The Evolution of the Romanian Digitalized Lexicography. The Essential Romanian Lexicographic Corpus*, in *Proceedings of the 15 th EURALEX International Congress*, 7-11 august 2012, Oslo, eds. Ruth Vatvedt Fjeld, Julie Matilde Torjusen, Press Representrales, UiO, p. 225, *in extenso* (pe suport electronic) p. 1014-1017.