

# STATISTICI PARȚIALE LA ÎNCHEIEREA PROIECTULUI eDLR - DICȚIONARUL TEZAUR AL LIMBII ROMÂNE ÎN FORMAT ELECTRONIC

DAN CRISTEA

Universitatea „Alexandru Ioan Cuza” din Iași

GABRIELA HAJA

Academia Română, Filiala din Iași

ALEX MORUZ

Academia Română, Filiala din Iași

MARIUS RĂSCHIP

Universitatea „Alexandru Ioan Cuza” din Iași

MĂDĂLIN PĂTRAȘCU

Universitatea „Alexandru Ioan Cuza” din Iași

## 1. Dicționarul Academiei

Pe parcursul a mai bine de un secol, a fost elaborată sub egida Academiei Române<sup>1</sup> opera lexicografică cea mai importantă pentru cultura noastră, *Dicționarul limbii române*, comparabil, prin dimensiuni și perspectiva abordării lexicografice, cu lucrări similare din lexicografia europeană. O prezentare statistică a unor date generale puse paralel cu cele referitoare la alte două mari dicționare este relevantă: *Dicționarul limbii române* a fost redactat și editat în două etape (seria cunoscută sub sigla DA, în perioada 1906–1944; seria nouă, DLR, în perioada 1965–2010), în 37 de volume și cuprinde cca 175.000 de cuvinte și variante, cu peste 1.300.000 de citate. Elaborarea variantei electronice s-a efectuat în intervalul 2007–2010. *Deutsches Wörterbuch „der Grimm”* (DWB), realizat și publicat în perioada 1838–1961, în 32 de volume, cuprinde 350.000 cuvinte și variante. Elaborarea variantei electronice s-a derulat între 1997 și 2004. *Trésor de la Langue Française* (TLFi), cuprinzând lexicul scris al secolelor XIX și XX, elaborat în perioada 1971–1994 (prima ediție tipărită), în 16 volume, cuprinde 100.000 cuvinte, 270.000 definiții, 430.000 exemple, conceput *ab initio* în variantă informatizată.

### 1.1 Dicționarul limbii române (DA)

*Dicționarul limbii române*, coordonat de Sextil Pușcariu între 1906–1944, a rămas cunoscut sub numele de *Dicționarul Academiei* (DA). Elaborarea lucrării a fost încredințată, după două încercări nefinalizate, unui colectiv mai larg, alcătuiindu-se două echipe, alături de autorul principal, Sextil Pușcariu. În timp, însă, a rămas câte un lexicograf în fiecare echipă: C. Lacea și, respectiv, Th. Capidan.

---

<sup>1</sup> La 1 aprilie 1866, se organizează „Societatea Literară Română”, care, peste un an, devine „Societatea Academică Română” (Pascu 1991: 54–55). *Ortografia, Gramatica și Dicționarul* limbii române au constituit de la bun început repere ale existenței Societății Academice Române, ceea ce a făcut ca Secția filologică-literară să ocupe un loc privilegiat, atât prin numărul de membri, cât și prin problemele dezbătute și rezolvate (Pascu 1991: 105). Societatea Academică Română a fost declarată „Institut Național”, sub numele de *Academia Română*, la 29 martie 1879.

Folosind fișierul preluat de la A. Philippide, precum și mai multe zeci de mii de fișe extrase de Pușcariu și echipa lui, s-a trecut la redactarea propriu-zisă, urmând ca îmbogățirea materialului lexical să se facă paralel cu operația de redactare. Vreme de 38 de ani, din 1906 până în 1944, s-a realizat aproape jumătate din lucrare (3071 de pagini tipărite însumând literele *A-C*, *F-K* și o bună parte din litera *L*, *l* – *lojniță*, iar începutul literei *D* (*d* – *de*) a rămas în fază de tipar). Pentru întâia oară în istoria lexicografiei românești redactarea propriu-zisă este dublată de o preocupare lexicologică permanentă; de aceea, *Introducerea* semnată de Sextil Pușcariu la primul tom din DA constituie una dintre cele mai largi expuneri de motive din istoria lexicografiei române (Seche 1969: 37). Sextil Pușcariu prezintă principiile de alcătuire a listei de cuvinte, pornind de la dicționarele anterioare precum și de la extrasele din literatură, aparținând unor texte cât mai variate, din epoci diverse. În privința criteriului de introducere a cuvintelor în lista dicționarului, DA, ca dicționar istoric și general, include atât cuvinte vechi, populare și regionale, cât și împrumuturi recente care „exprimă o idee sau nuanța unei idei pentru care limba noastră nu are un termen neechivoc” (DA 1913: XV). Numărul cuvintelor-titlu este mare: dintr-o statistică întreprinsă de S. Pușcariu rezultă că fasciculele cuprinzând patru litere, medii ca întindere (*A*, *B*, *F* și *G*), însumează 15.444 de intrări. În DA elementele componente ale fiecărui articol au un loc stabil: după cuvântul-titlu și indicarea categoriei lui morfologice urmează, în ordine, precizarea domeniului căruia îi aparține cuvântul, traducerea sensurilor termenului românesc în limba franceză, indicații asupra situației sale istorice sau asupra răspândirii diatopice ori indicații stilistice, definiția, izvoarele ei, unitățile frazeologice, formele gramaticale, variantele lexicale, etimologia. De la litera *B* se schimbă modul general de organizare a materialului, adoptându-se metoda „cuiburilor” lexicale, potrivit căreia derivatele sunt subsumate cuvântului-titlu, figurând în același articol cu acesta, ceea ce determină simplificarea excesivă a analizei semantice a termenilor subsumați. Sensurile au fost despărțite în numeroase ramificații, cu cele mai diverse sisteme de izolare (cifre romane și arabe, litere mari și mici, bare duble și simple etc.), împrumutate din lexicografia franceză. Este o premieră, în istoria lexicografiei românești, organizarea semantică a materialului, în mod sistematic, pe două linii dominante: prin *coordonare* și prin *subordonare*. Unul dintre dezideratele coordonatorului, acela de a realiza, prin modul de explicare a termenilor, o operă lexicografică lingvistică și nu una enciclopedică, a fost în bună măsură atins. Lista bibliografică a dicționarului a fost mult îmbogățită, ajungându-se de la cele 199 de izvoare folosite de A. Philippide, la circa 553 de titluri în 1926. DA se deosebește de versiunile anterioare și prin faptul că numărul citatelor crește în cadrul unui articol, potrivit unor necesități diferențiate, până la câteva zeci (față de cele trei folosite de Philippide), astfel încât fiecare segment dintr-o epocă istorică este reprezentat prin texte (Seche 1969: 36-39). Autorul a acordat, însă, o mai mică atenție neologismelor în ansamblul lucrării, potrivit concepției sale privitoare la neologizarea exagerată a limbii. În dicționarul condus de Pușcariu este continuată linia lui Philippide cu privire la unitățile frazeologice, în sensul integrării lor în filiația semantică a articolului, subordonându-le direct, pe fiecare, sensurilor din care se dezvoltă.

### **1.2 Dicționarul limbii române (DLR). Serie nouă**

Lucrările la *Dicționarul limbii române*, întrerupte brusc de război și, mai cu seamă, de schimbarea regimului politic, au fost reluate, într-o nouă formulă, începând din 1949, după reorganizarea Academiei Române. S-a considerat util ca noua versiune a dicționarului academic, fără a fi propriu-zis o continuare a celui anterior, să înceapă cu litera *M*, urmând ca, ulterior, după terminarea porțiunii până la *Z*, să fie elaborată, pe aceleași baze, întreaga porțiune de la *A* la *L*; obiectivul a fost atins prin volumele ce cuprind literele *D*, *E*, *K*, *L*, *Q*. Noua serie a dicționarului se deosebește mult de cea imediat anterioară prin condițiile în care este elaborat: trei echipe largi, la București, la Cluj și la Iași, care au lucrat paralel la redactarea Dicționarului. La București s-a efectuat, sub conducerea lui Iorgu Iordan, Alexandru Graur și Ion Coteanu, revizia finală a materialului. Ulterior, revizia finală a fost făcută în fiecare dintre cele trei centre de cercetare, coordonarea lucrării fiind asumată, din 1997, de Marius Sala și Gh. Mihăilă. O altă deosebire față de DA o constituie materialul lexical, astfel că, numai la litera *M*, care începe noua

versiune, se folosește o bibliografie cu 1627 de titluri. DLR înregistrează toate cuvintele atestate în limba literară generală și în limbajul literaturii artistice, în vorbirea populară și regională, în textele vechi; își propune (ca și DA) să fie mai puțin deschis la cuvintele aparținând terminologiei tehnico-științifice. În noua versiune, se aplică un criteriu riguros de selectare a neologismelor: au dreptul să figureze în dicționar acei termeni ai limbajelor tehnico-științifice care au pătruns sau care manifestă tendința evidentă de a pătrunde în limba literară generală, în limbajul literar artistic și în cel popular (DLR 1965: VI); hotărâtoare este așadar atestarea circulației unui cuvânt în cel puțin două dintre stilurile funcționale ale limbii române. În DLR, sistemul de concentrare a termenilor în „cuiburi” lexicale a fost abandonat, fiecare lexem fiind înregistrat ca articol independent. Cele două opere lexicografice se deosebesc între ele și prin modul de folosire a izvoarelor. Diferența cea mai evidentă rezultă din volumul diferit de extrase de care a dispus fiecare dintre lucrări; în plus, dimensiunea bibliografiei de referință a Dicționarului a sporit de la un tom la altul, încât, dacă în 1965 erau citate 1627 de titluri, la finele primei ediții numărul acestora ajungea la 2587. Fișarea s-a făcut continuu, pentru fiecare literă în lucru, în fiecare dintre cele trei centre. În DLR, ordinea tuturor izvoarelor ilustrative este strict cronologică, încadrându-se astfel în caracterul de bază al lucrării, acela de dicționar istoric. În DLR alineatul final al articolului, cuprinzând etimologia, este semnificativ mai redus decât în versiunea anterioară, pentru că numeroase explicații de amănunt au fost lăsate pe seama *Dicționarului etimologic al limbii române*, în curs de elaborare. Însă chiar DLR suportă adăugiri și completări, de la un volum la altul, datorită faptului că multe neologisme, pătrunse în limbă de la începutul secolului al XX-lea, regionalisme incluse în *Atlasele lingvistice* sau în culegeri de texte recente, precum și termeni vechi extrași din documente aparținând secolelor al XVI-lea – al XVIII-lea editate în ultimele decenii (Sala, Mihăilă 2000: VI) nu au fost cuprinse în primele volume ale noii serii și lipsesc, parțial, și din celelalte tomuri.

Redactarea intrărilor s-a făcut până în anii '90 în manieră tradițională, textul scris de mână era dactilografiat și corectat înainte de a fi trimis la editură, unde se culegea din nou, pentru a fi corectat și în șpalturi. Abia din anul 1995 tehnoredactarea intrărilor a început să se facă pe calculator, utilizând programe de editare nespecializate (Microsoft Word). Textul acestor litere (*D, E, K, L, Q, V, W, X, Y, Z*) a fost trimis la editură pe suport electronic, acolo având loc (potrivit tradiției) ultimele corecturi, formatul final electronic fiind prelucrat printr-un program de editare (PageMaker).

Ca punct de plecare în achiziționarea formatului electronic al marelui Dicționar, s-a dispus așadar de două tipuri de surse primare: un număr de volume, publicate înainte de 1997 (Haja et al. 2005: 23), care nu existau decât în forma tipărită și restul de volume, pentru care se dispunea de formatul electronic ante-tipar, adică cel păstrat la sediile celor trei institute de lingvistică, în momentul predării materialului spre editare. Formatul electronic al variantelor finale, tipărite de Editura Academiei, nu a putut fi recuperat de la păstrătoarea acestuia.

## **2. Prezentarea proiectului eDTLR – Dicționarul tezaur al limbii române în format electronic**

În primăvara anului 2007, răspunzând apelului Ministerul Educației, Cercetării și Tineretului de formare de Parteneriate în Domenii Prioritare, se constituia un consorțiu din 7 parteneri aparținând Universității „Alexandru Ioan Cuza” din Iași și Academiei Române. Partenerii din cadrul Universității ieșene au fost Facultatea de Informatică și Facultatea de Litere, iar din cadrul Academiei Române, Institutul de Lingvistică „Iorgu Iordan – Al. Rosetti”, București, Institutul de Filologie Română „A. Philippide”, Iași, Institutul de Lingvistică și Istorie Literară „Sextil Pușcariu”, Cluj-Napoca, Institutul de Cercetări pentru Inteligență Artificială, București și Institutul de Informatică Teoretică, Iași. Aceștia au câștigat finanțarea unui proiect care s-a numit **eDTLR – Dicționarul tezaur al limbii române în format electronic**. La începerea proiectului, în septembrie 2007, se mai lucra încă la finalizarea volumelor din marea serie a Dicționarului (literele *D, E, K, L, Q*). În proiect s-a pornit însă de la premisa că activitatea de elaborare a Dicționarului avea să se încheie înainte de terminarea proiectului eDTLR, ceea ce s-a și împlinit (ultimul volum al seriei a fost tipărit în aprilie 2010).

Pe parcursul a 3 ani, proiectul eDTLR ar fi trebuit să materializeze următoarele obiective principale: **1.** realizarea unei variante electronice a marelui *Dicționar Tezaur*; **2.** realizarea unei baze de date care să cuprindă izvoarele textuale ale Dicționarului; **3.** realizarea legăturii dintre intrările Dicționarului și paginile de unde au fost excerptate exemplele; **4.** construirea unei interfețe care să permită consultarea interactivă a Dicționarului și regăsirea exemplelor în contextele lor originare.

În ciuda mai multor dificultăți, dintre care cele mai mari au fost generate de reducerea finanțării și de procesul de adaptare a lingviștilor lexicografi la nou createle tehnologii și metode de lucru, precum și de găsirea unui limbaj comun de comunicare între specialiștii celor două domenii – informatică și lingvistică, la terminarea lui, proiectul raporta următoarele realizări:

– un mediu complex de prelucrare a formatelor electronice ale Dicționarului, care include: interfețe de editare online a formei obținute în urma OCR<sup>2</sup>-izării; programe de curățare și de conversie a formatului HTML, obținut după editarea online, într-o variantă simplificată XML<sup>3</sup>; două variante (pentru formatul DA și formatul DLR) ale unui program de recunoaștere a câmpurilor din intrările de dicționar (parser) și etichetarea lor corespunzătoare în XML utilizând convențiile TEI-P5<sup>4</sup>; un lanț de post-procesare care realizează o curățare a codurilor XML rezultate în urma parsării; programe de actualizare a bazei de date a intrărilor de dicționar și de acces la baza de date; o interfață care dă acces utilizatorilor la eDTLR, permițând căutări condiționate și generarea de statistici; programe de verificare a integrității surselor scanate; programe de segmentare a imaginilor paginilor scanate pentru recuperarea coordonatelor spațiale ale fiecărui cuvânt în pagină; programe de regăsire a citatelor în formatele txt imperfecte ale surselor (*approximate string matching*) și de indexare a citatelor în imaginile paginilor scanate;

– o bază de date a Dicționarului, conținând Dicționarul, integral, în următoarele formate: **a.** pentru partea veche (secțiunile pentru care nu s-a dispus decât de suportul pe hârtie): imaginile paginilor scanate, formatul rtf<sup>5</sup> rezultat în urma OCR-izării, formatul HTML<sup>6</sup> rezultat în urma corectării, formatul XML TEI-P5 rezultat în urma parsării; **b.** pentru partea nouă (secțiunile pentru care s-a dispus de formatul electronic): formatul rtf original al pre-tipăriturilor, formatul XML TEI-P5 rezultat în urma parsării; o bază de date a surselor bibliografice, conținând lista tuturor surselor cu numele acestora, autorii, siglele și cronologia lor;

– o bază de date a surselor bibliografice ale Dicționarului, conținând, pentru fiecare volum: forma imagine a volumului scanat și forma rtf obținută în urma OCR-izării lui (așadar texte cu imperfecțiuni).

eDTLR a fost realizat în mai multe etape, rezumate mai jos în această secțiune.

## **2.1 Procesări preliminare aplicate volumelor existente doar pe hârtie**

---

<sup>2</sup> OCR = *Optical Character Recognition*, operațiunea de recunoaștere a codului caracterelor scrise de mână ori tipărite, dintr-o pagină scanată. În funcție de calitatea hârtiei, de alfabet și de limbă, recunoașterea poate genera mai multe ori mai puține erori.

<sup>3</sup> XML = *eXtensible Markup Language*, un limbaj de adnotare, cu aplicare la formate text, imagine, sunet și multimedia.

<sup>4</sup> TEI – *Text Encoding Initiative*, un standard de adnotate a documentelor, cuprinzând un capitol special, relativ la dicționare online. P5 este versiunea ultimă, cuprinzând și o descriere a convențiilor de marcare a intrărilor de dicționar.

<sup>5</sup> Un format care cuprinde adnotări de formatare, utilizat curent de editoarele de texte, precum *Microsoft Word*, de exemplu.

<sup>6</sup> HTML – *HyperText Markup Language*, un format de adnotare utilizat cu precădere în realizarea paginilor Web.

Proiectul a început prin scanarea volumelor din seria veche a Dicționarului, care nu existau decât în formă tipărită. Fiecare imagine astfel obținută includea două pagini A4, care au fost separate, segmentate pe coloane, curățate de margini negre, de adnotări manuscrise marginale, de pete etc., reduse apoi de la 400 la 300 dpi, înainte de a fi OCR-izate. Dificultăți suplimentare la transpunerea în text a dicționarului au fost adăugate de calitatea inegală a tipăriturilor și de existența mai multor alfabet utilizate în redarea etimoanelor: latin, grecesc, chirilic, turcesc, polonez, rusesc etc.

Simultan, după o scurtă perioadă în care s-a actualizat fișierul surselor bibliografice, s-a început lucrul la scanarea surselor bibliografice, acest proces derulându-se apoi pe toată durata proiectului.

## **2.2 Faze de corectare**

Întrucât operația de OCR-izare introduce erori, volumele Dicționarului care au suferit acest tip de prelucrare au necesitat corecții. În proiect s-a hotărât ca aceste volume să fie corectate în două etape, mai întâi de către novici, printr-un proces colaborativ (Cristea et al. 2008), iar ulterior de către experți lexicografi – parteneri în proiect. A fost creat un portal dedicat corecturii online<sup>7</sup> și, pe diverse căi, s-a adresat chemarea de a contribui la această acțiune<sup>8</sup>.

Întrucât Academia Română nu putea permite, fapt lesne de înțeles, ca versiuni ale Dicționarului care ar fi putut fi extrase de pe site-ul proiectului să circule neautorizat, s-a adoptat soluția de a împărți fiecare pagină a Dicționarului în 12 segmente, doar un singur segment fiind accesibil pentru corectură la un moment dat, iar extragerea lor în secvență făcându-se aleatoriu. Prin această strategie, orice tentativă de a recompune chiar și o singură pagină a fost descurajată.

S-a dovedit că prima corectură a fost destul de inegală, o parte dintre experți raportând la corectura a doua un număr exagerat de erori găsite pe anumite porțiuni, când așteptările ar fi fost ca doar puține erori să mai fi rămas în urma corectorilor benevoli. Cu toate acestea, considerăm că antamarea acestui efort colaborativ a fost o alegere corectă a proiectului, pentru că, fără el, corectura a doua ar fi fost mult mai greoaie. La asamblarea finală a lucrării, s-a constatat totuși, în câteva locuri, lipsa unor secțiuni. Acestea nu puteau fi datorate decât unor omisiuni în scanările originale, care nu fuseseră verificate pagină cu pagină, așa cum s-a procedat ulterior cu sursele bibliografice. În consecință, în ultimele două luni ale derulării proiectului aproximativ 200 de pagini au fost introduse în baza de date a Dicționarului, direct în format HTML, de un număr mic de cercetători care s-au oferit să refacă aceste lipsuri.

Au participat la corectarea celor 11410 pagini, prelucrate după cum am descris mai sus, peste 900 voluntari, dintre care: 392 de curioși și 508, cu următoarea distribuție, după numărul de secvențe corectate: 127 au corectat între 10 și 50 de secvențe, 90 au corectat 50 – 100 de secvențe, 102 au corectat 100 – 200 de secvențe, 102 au corectat 200 – 500 de secvențe, 40 au corectat 500 – 1000 de secvențe și 47 au corectat peste 1000 de secvențe.

## **2.3 Parsarea intrărilor**

După aceste două faze de corectură, textul electronic al volumelor vechi ar fi trebuit să fie identic cu cel tipărit. Ultima fază în procesare, parsarea, urma să explicitizeze semnificația câmpurilor intrărilor, operație extrem de dificilă pentru că presupune interpretarea marcajelor de formatare și de demarcare a câmpurilor în context. Respectând practica internațională, structura intrărilor a fost evidențiată în notații ale limbajului de marcare XML, cu aplicarea standardului utilizat în descrierea dicționarilor, TEI-P5. Parsarea a fost pregătită prin desemnarea categoriilor de marcaje utilizate în Dicționar pentru a pune în evidență diferite tipuri de câmpuri ale intrărilor. Aceste marcaje, clasificate pe niveluri, ajută parserul să descopere mai întâi un arbore general al sensurilor, pentru ca mai apoi fiecare nod al arborelui să fie detaliat suplimentar (Curteanu et al. 2008), rezultând o recunoaștere fidelă a

---

<sup>7</sup> <https://consilr.info.uaic.ro/edtlr/>

<sup>8</sup> Aproximativ 900 de profesori, doctoranzi, studenți de la universitățile din Iași, Suceava, Bacău, Baia Mare, Galați și altele, precum și din Republica Moldova au răspuns apelului nostru.

schemei semantice a fiecărui articol din Dicționar, fiecare sens și subsens fiind însoțit de șirul citatelor ilustrative cu precizarea izvoarelor.

Această tehnică are instanțieri diferite pentru volumele din seria DA și cele din seria DRL. Mai mult decât atât, pentru că volumele seriei DLR provin, la rândul lor, din două surse diferite (cele obținute prin OCR-izare din scanări și corectate, și cele recuperate din formate electronice originare), cele două serii de volume au necesitat prelucrări specifice de pregătire a parsării. Prelucrarea volumelor editate cu ajutorul calculatorului (seria nouă DLR, tipărită după 1997) a fost terminată înaintea volumelor editate cu mijloace vechi (identificând aproximativ seria veche DA). Diferențele de structură dintre DA și DLR ne-au obligat să elaborăm două versiuni diferite de parsere, unele dintre informațiile codificate în varianta DA trebuind să rămână nespecificate în descrierea TEI finală.

Parsarea reprezintă cel mai complicat proces din întregul lanț de prelucrare al Dicționarului. În mod normal, programele de parsare sunt realizate cu ajutorul unor gramatici speciale, dependente de context, extrem de dificil de elaborat. De aceea, am optat pentru o soluție care se bazează pe recunoașterea și ierarhizarea marcajelor delimitatoare ale câmpurilor din intrare, care a permis reducerea semnificativă a dimensiunii gramaticii și elaborarea parserului într-un timp mult mai scurt decât cel necesar aplicării metodelor clasice. Cu toate acestea, acum, în momentul încheierii oficiale a proiectului, parserul nu e considerat încă finisat suficient. Erorile de structură găsite de noi au, în esență, trei surse: greșeli de dactilografie a textului, încălcarea normelor de redactare la realizarea unui articol și situații rare, încă netratate de către parser. Parsarea scoate la iveală imediat încălcarea normelor. Acestea produc fie structuri paradoxale, fie erori de parsare. Parserul a fost proiectat să aibă un comportament extrem de robust, o eroare de parsare afectând, în general, doar local structura. În astfel de situații, parserul își revine din starea de eroare, ca după o sincopă, primul indice semnificativ găsit din acel loc înainte fiind folosit pentru a repune parserul pe nivelul corect al arborelui de structură, procesul putând fi astfel continuat.

Inițial, proiectul prevedea și o fază de corectare a structurii (care nu a mai fost finanțată), în care erorile de structură ar fi trebuit observate și remediate de către experții lexicografi, capabili să interpreteze formatul XML și să opereze modificări asupra lui. Soluția pe care o gândim acum are în vedere o altă strategie, care ar presupune, din nou, o activitate benevolă a unei mase mari de utilizatori cunoscători, care se va limita, de astă dată, strict la *semnalarea* erorilor. Interfața de consultare a Dicționarului încorporează deja o funcționalitate de semnalare a erorilor: utilizatorul selectează o porțiune a textului pe care o consideră eronată și, dacă dorește, completează și un comentariu. Un mesaj e-mail pleacă apoi spre server, care depozitează astfel de semnale ce urmează a fi tratate ulterior de experți informaticieni și lexicografi. Operațiile de corectare nu vizează o corectare de suprafață a structurii generate greșit, ci, în funcție de caz, fie o intervenție asupra textului pentru corectarea unei erori evidente de dactilografie, fie una asupra parserului însuși. Reluarea parsării pentru intrarea respectivă va produce structura corectă.

#### **2.4 Utilizarea surselor bibliografice ale Dicționarului**

Două au fost motivele care ne-au determinat să scanăm sursele bibliografice ale marelui Dicționar. În primul rând, pentru că am dorit să legăm citatele din Dicționar de imaginile paginilor de unde acestea au fost extrase. Utilizatorul interesat de ocurența cuvântului este acum în măsură să consulte contextul fiecărui citat ilustrativ, care, în funcție de situația drepturilor de autor a sursei respective, poate fi egal ca întindere cu citatul însuși ori nelimitat. În al doilea rând, pentru că sursele se constituie astfel într-o bibliotecă online care ar putea fi deschisă spre consultare interactivă lexicografilor, pentru căutarea de contexte ale cuvintelor. Biblioteca încorporează în prezent fișiere cu imaginile paginilor volumelor și cu textele rezultate în urma OCR-izării. Variantele de text nu sunt corectate, pentru că nu ne-am propus acest lucru. Nici nu era posibil această operație, deoarece, baza textuală conține 2.600 de volume (cărți, reviste, dicționare etc.) din cele aproximativ 4.000, câte are *Bibliografia* DLR. În prezent, sunt prelucrate complet și accesibile doar 842 de volume, dar se lucrează în continuare la postarea în Biblioteca eDTLR a celor rămase. Numărul total de pagini scanate și OCR-izate este de 972.000.

Pentru a contracara lipsa de acuratețe a textelor izvoarelor, obținute după OCR-izare, s-au dezvoltat algoritmi de căutare aproximativă, care fac posibilă regăsirea citatelor, inclusiv în cazurile în care cuvântul căutat a fost OCR-izat imperfect. Utilizatorul are sub ochi doar pagina scanată, deși interfața ascunde sub ea formatul text. Toate operațiile de selecție și de copiere realizate de utilizatorul lexicograf se realizează asupra textului (așa cum precizăm, imperfect), însă confruntarea cu imaginea originalului aflată pe ecran lasă posibilitatea corectării eventualelor erori.

Dintre cele 2.600 de volume scanate, 2.148 fac parte din Biblioteca DLR a Institutului de Filologie Română „A. Philippide”, 308 volume au fost scanate în cadrul Bibliotecii Centrale Universitare din Iași, unde imaginile și textele rezultate sunt depozitate și accesibile online, iar 144 au fost scanate în cadrul Institutului de Lingvistică „Iorgu Iordan – Al. Rosetti” din București, care le-a pus la dispoziția proiectului.

### **3. Analize statistice ale materialului din DLR**

Consultarea Dicționarului și obținerea statisticilor de mai jos este posibilă grație interfeței create în cadrul proiectului și care oferă o serie de opțiuni pe care le prezentăm aici pe scurt. Așa cum rezultă și din nume, interfața este un instrument electronic ce permite legătura dintre utilizator și imensa bază de date și programe care este eDTLR. În prezent, sunt accesibile doar volume corespunzătoare seriei noi a Dicționarului. Celelalte tomuri, inclusiv cele din seria DA există în baza eDTLR, asupra lor derulându-se procesul de prelucrare (parsare) a articolelor, acestea urmând a fi puse, treptat, la dispoziția celor ce vor dori să le consulte.

Pe lângă vizualizarea informației din *Dicționar*, interfața permite o serie de statistici, operate pe întregul material lexicografic. În chip firesc, pe măsură ce se adaugă datele din volumele în curs de prelucrare, rezultatele analizelor statistice se modifică. Este permisă, în mod firesc, operația de consultare la nivel de literă ori de cuvânt.

Criteriile de consultare au fost definite în funcție de normele de structurare lexicografică a Dicționarului, așa cum au fost ele integrate în parser. Pentru consultarea după criterii cronologice a fost necesară o prelucrare suplimentară a bazei de date eDLTR, în sensul adnotării fiecărui citat cu indicația anului, așa cum a putut fi el extras automat, după mai multe elemente auxiliare, conținute în *Siglarul* DLR, completat și transpus în format XML. Potrivit *Siglarului* DLR, fiecare autor sau lucrare (dicționar, periodic etc.) din Bibliografia DLR are atribuit un indice numeric ce precizează rangul său cronologic. În unele situații, nu este precizat acest indice cronologic, deoarece lucrarea respectivă este o colecție de documente dintr-o perioadă istorică extinsă ori o publicație periodică; pentru aceste cazuri, s-a luat în considerare anul precizat între paranteze rotunde la finele citatului, înainte de siglă (ex. FURNICĂ, I. C.) sau anul din siglă (V. ROM. iunie 1970). În alte situații, indicelui numeric atribuit unui autor sau unei lucrări nu îi corespunde un an (ex. NEGRUZZI, A., NEGRUZZI, A. P.) și atunci a trebuit să se găsească o rezolvare de compromis, prin echivalarea indicelui (124, în cazul nostru) cu un interval de ani. O a treia situație problematică, este aceea a citatelor populare din colecții mixte, care cuprind texte literare și populare, dintr-o perioadă istorică veche (ex. ȘIO). Citatelor populare extrase din aceste culegeri nu li s-a putut atribui niciun an, rămânând ca soluție realizarea unui algoritm de deducere a cronologiei în funcție de cronologia vecinătăților citatului.

Toate criteriile de consultare a Dicționarului pot fi utilizate singular ori în combinații. Aceste formule de căutare pot fi salvate de utilizator, pentru căutări ulterioare.

Pentru a sugera posibilitățile de documentare pe care formatul electronic al Dicționarului le oferă, am operat o serie de interogări după diverse criterii asupra a trei litere: E (5.885 de intrări redactate la Iași), L (6.039 de intrări redactate la Iași și Cluj) și S (14.347 de intrări redactate la București).

Statisticile care urmează sunt realizate la nivel de intrare, urmărindu-se situația cuvintelor în contexte, lucru care poate antrena concluzii de ordin semantic și funcțional.

Criterii		E		L		S		
Criterii morfologice	substantive	masculine	616	10,48%	401	6,64%	2140	14,92%
		feminine	2269	38,56%	1278	21,16%	5209	36,31%
		neutre	1093	18,57%	592	9,80%	2395	16,69%
	verbe		544	9,24%	482	7,98%	1731	12,07%
	adjective		1151	19,56%	621	10,28%	2495	17,39%
	adverbe		82	1,39%	80	1,32%	117	0,82%
Criteriu etimologic	franceză		2439	41,44%	1008	16,69%	1388	9,67%
	latină		612	10,40%	344	5,70%	453	3,16%
	albaneză		1	0,02%	5	0,08%	9	0,06%
	slava veche		6	0,10%	15	0,25%	71	0,49%
	maghiară		24	0,41%	100	1,66%	63	0,44%
	turcă		49	0,83%	44	0,73%	160	1,12%
Criteriul circulației	învechit		518	8,80%	802	13,28%	1870	13,03%
frecvență	rar		279	4,74%	304	5,03%	1046	7,29%
diacronie și frecvență	învechit+rar		33	0,56%	74	1,23%	249	1,74%
spațiu	regional		1	0,02%	655	10,85%	1367	9,53%
	și Moldova		1	0,02%	107	1,77%	236	1,64%
	numai Moldova		0	0	44	0,73%	105	0,73%
	numai Transilvania		9	0,15%	86	1,42%	202	1,41%
Îmbinări stabile	intrări sub care apar expresii		40	0,68%	177	2,93%	336	2,34%
Compuse	intrări sub care apar cuvinte compuse		0	0	42	0,70%	62	0,43%
După autori	CORESI		18	0,31%	120	1,99%	267	1,86%
	EMINESCU		421	7,15%	356	5,90%	427	2,98%
	PREDA		194	3,30%	201	3,33%	498	3,47%
După cronologie	înainte de 1600		2	0,05%	7	0,12%	13	0,09%
	între 1700 și 1800		4	0,07%	8	0,13%	22	0,15%
	între 1800 și 1900		70	1,19%	139	2,30%	352	2,45%
	după 1900		451	7,66%	421	6,97%	979	6,82%

Procentele au fost calculate raportând numărul de intrări rezultat în urma căutării la numărul total de intrări pe fiecare literă. Criteriul morfologic poate fi exploatat mai fin, dar am optat aici pentru căutarea după doar patru părți de vorbire (substantiv, verb, adjectiv, adverb) distingând, pentru substantive, categoria gramaticală a genului. Suma rezultată și lista cuvintelor cuprind acele intrări sub care se află, cel puțin o dată, una dintre aceste folosiri (încât, de exemplu, un substantiv care are forme de masculin și de feminin va apărea în două liste). Prima observație evidentă, în urma comparării ponderii părților de vorbire analizate în cele trei litere, este aceea că în lista de cuvinte domină numeric și proporțional substantivele, iar, dintre acestea, substantivele feminine sunt cele mai numeroase. De asemenea, este interesant de analizat faptul că o literă de dimensiuni reduse spre medii, precum litera *E* depășește, în privința ponderii substantivelor, o literă de mari dimensiuni cum este *S*.

La criteriul etimologie valorile raportate cuprind cuvinte cu etimon cert, unic, sau cu etimologie multiplă; sunt luate în considerare și etimonele complexe, conform cărora un cuvânt se formează pe teren românesc după un model străin. Interfața nu face, deocamdată, distincție între cuvintele de origine



latină împrumutate și cele moștenite. Rezultatele obținute în urma căutărilor după criteriul etimologic demonstrează importanța influenței limbii franceze asupra lexicului românesc, începând cu secolul al XIX-lea, când intră în circulație majoritatea neologismelor cu această origine. Firește, și în cazul căutărilor după criteriul etimonului, rezultatele obținute trebuie înțelese astfel: litera *E* cuprinde 2439 de cuvinte care provin numai sau și din limba franceză. Semnificativ este faptul că aproape jumătate dintre cuvintele literei *E* sunt (și) de origine franceză (41,44%), în vreme ce pentru celelalte litere se schimbă raportul. Diferențele au explicații lingvistice, dar și conjuncturale (anul / intervalul redactării și al publicării fiecărei litere). Aceste rezultate pot fi coroborate cu cele obținute în urma căutării după vechimea primei atestări, cu cele obținute în urma căutării regionalismelor ori a sensurilor regionale și cu cele privitoare la numărul de cuvinte care au dezvoltat expresii. Remarcăm că îmbogățirea lexicului limbii române a crescut semnificativ după 1900, dar această concluzie este determinată și de faptul că numărul izvoarelor textuale editate după 1900, incluse în Bibliografia DLR, este mare, iar fișierul extras pentru această perioadă este foarte bogat. Din analiza datelor numerice și a proporțiilor rezultate se impune, cu evidență, caracterul neologic al literei *E*, prin comparație cu celelalte două, mai echilibrate din punctul de vedere al distribuției lexicului, după vechime, spațiu, uz etc.

Fie și parțiale, aceste rezultate și analiza lor sumară oferă specialiștilor interesați de studiul lexicului limbii române, așa cum este el analizat și prezentat în Dicționar, câteva repere pentru modalitatea de consultare și de valorificare a bazei eDTLR.

#### **4. Concluzii**

Așa cum era de așteptat, semnificația ieșită din comun pentru cultura românească a marelui Dicționar tezaur al limbii române și apariția formatului lui electronic, face acest produs extrem de dorit publicului vorbitor de limba română (incluzând aici toate categoriile de utilizatori, de la oameni obișnuiți care doresc să se informeze asupra sensului unui cuvânt, până la cercetători ai limbii). Membrii consorțiului eDTLR au avut semnale repetate care arată nerăbdarea cu care este așteptată apariția nerestricționată a eDTLR în internet.

Este de necontestat că eDTLR trebuie să devină un bun public din mai multe motive, cel mai important fiind acela că atât realizarea conținutului cât și realizarea formatului digital au fost plătite din bani publici. Pe de altă parte, finanțarea lui parțială, sub nivelul promis la momentul încheierii contractului, a împiedicat ducerea la îndeplinire a tuturor etapelor de corectare prevăzute, iar ieșirea în spațiul public a unui bun imperfect, sub girul Academiei Române, nu este o opțiune plauzibilă. Drept urmare, intenționăm ca pentru o perioadă, atâta vreme cât se vor mai constata erori de conținut ori structură, eDTLR să fie deschis spre consultare numai specialiștilor, într-o etapă ulterioară *Dicționarul* urmând a fi deschis spre consultare publicului larg.

Mai sunt, de asemenea, de rezolvat câteva probleme de suport tehnologic. Accesul la baza de date în căutare angrenează procese de căutare sofisticate, mari consumatoare de timp. Nu știm cum se va comporta serverul atunci când un număr mare de utilizatori vor face interogări în același timp. Vor fi de considerat soluții de distribuire a bazei de date și de paralelizare a calculelor. Amintim totodată și imensa responsabilitate pe care gestionarea acestei importante baze de date o incumbă, ea trebuind să fie intangibilă la atacuri maligne ori la tentative de scurgeri de informații.

#### **5. Mulțumiri**

Proiectul eDTLR a fost finanțat prin contractul CNMP nr. 91\_013/18.09.2007. Mulțumim membrilor consorțiului eDTLR, fără a căror dedicație proiectul nu s-ar fi realizat. Adresăm mulțumiri, de asemenea, sutelor de contribuabili (de cele mai multe ori anonimi) care au ajutat la corectarea surselor, utilizând interfața online. Cu prisosință, merită toată recunoștința noastră firma PIM SRL Iași și directorul acesteia d-l Marius Petrariu, care, înțelegând importanța activității noastre și resursele limitate cu care suntem obligați să o realizăm, nu numai că ne-a oferit un preț extrem de convenabil pentru scanări, dar a acceptat să mute, pentru o perioadă de câteva luni, un scanner și să delege o persoană care să se ocupe

exclusiv de scanarea surselor aflate în sediul Bibliotecii Universitare Iași. Mulțumim conducerii acestei biblioteci, pentru că a acceptat să ne ofere, pentru a fi scanate la sediul lor, un număr de 308 volume care nu se află sub incidența Legii Patrimoniului Național și nici a Legii Drepturilor de autor<sup>9</sup>. Mulțumim doamnei Cătălina Marănduc pentru cele 144 de volume care au fost scanate de dânsa la sediul Institutului de Lingvistică din București și oferite consorțiului. Mulțumim următorilor cercetători care au acceptat să introducă în baza de date porțiunile care lipseau: Cristina Florescu, Elena Dănilă, Laura Manea, Mioara Dragomir și Monica Corodeanu, toți de la Institutul Philippide Iași, precum și Daniela Gifu și Iulia Scutariu de la Facultatea de Informatică.

### **Bibliografie**

1. Cristea, Dan, Corina Forăscu, Marius Răschip, Michael Zock, 2008, „How to Evaluate and Raise the Quality in a Collaborative Lexicographic Approach”, în *Proceedings of LREC-2008*, Marrakech.
2. Curteanu, Nicolae, Alex Moruz, Diana Trandabăț, 2008, „Extracting Sense Trees from the Romanian Thesaurus by Sense Segmentation & Dependency Parsing”, în *Proceedings of CogAlex Cognitive Aspects of the Lexicon: Enhancing the Structure, Indexes and Entry Points of Electronic Dictionaries*, COLING 2008, pp. 55–63, ISBN 978-1-905593-56-9
3. Haja, Gabriela, Elena Dănilă, Corina Forăscu, Bogdan-Mihai Aldea, 2005, *Dicționarul limbii române (DLR) în format electronic. Studii privind achiziționarea*, Iași, Editura Alfa.
4. Iordan, Iorgu, Al. Graur, Ion Coteanu (coord.), 1965, *Dicționarul limbii române (DLR)*. Serie nouă. Tomul VI, *Litera M*, București, Editura Academiei R.P.R.
5. Pascu, Ștefan, 1991, *Istoricul Academiei Române. 125 de ani de la înființare*, Editura Academiei Române, București.
6. Pușcariu, Sextil, *Dicționarul limbii române (DA)*, 1913, Întocmit și publicat după îndemnul și cu cheltuiala Maiestății Sale regelui Carol I, Tomul I, Partea 1, A–B, București, Librăria Socec & Comp. și C. Sfetea.
7. Sala, Marius, Gheorghe Mihăilă, 2000, *Cuvânt înainte*, în *Dicționarul limbii române (DLR)*. Serie nouă. Tomul XIV, *Litera Z*, București, Editura Academiei Române.
8. Seche, Mircea, 1969, *Schiță de istorie a lexicografiei românești*, vol. II, București, Editura Științifică.

### *PARTIAL STATISTICS AT THE END OF THE PROJECT "THESAURUS DICTIONARY OF THE ROMANIAN LANGUAGE IN ELECTRONIC FORM"*

#### *Abstract*

*The paper describes the project aimed at creating the electronic form of the big Dictionary of Romanian Language, an outstanding work elaborated by the Romanian Academy, all along a century. A short comparison with other similar dictionaries is first presented. Then, the two series of the Dictionary are described, the old series, elaborated before the second war and known as the Dictionary of the Academy, and the new series, elaborated between 1965 and 2010, known as the Dictionary of the Romanian Language. The methodology of building the electronic form of the Dictionary is then presented: scanning the volumes of the Dictionary which were found only on paper support, OCR-ing them and correcting the resulting editable form on two passes – one performed by novices, in a collaborating activity, the second operated by experts. The corrected manuscript as well as the newer volumes, originally typed in before printing, were passed to a parsing process which produced the final XML form encoding the structure of the entries. This form is deposited in a large database and can be accessed by the user through an online interface. To show the searching and interpretation possibilities opened by this technology, a number of statistics are presented, as an exercise, on three letters: E, L and S. Finally, come conclusions are drawn.*

---

<sup>9</sup> Conform convenției încheiate între Universitatea din Iași și Biblioteca Universitară din Iași, aceste surse se găsesc depozitate pe un server aflat la sediul BCU-Iași, de unde sunt accesate prin internet.