

# Identificarea automată a afixelor românești. Studiu de caz: identificarea sufixelor

Verginica Barbu Mititelu

Institutul de Cercetări pentru Inteligență Artificială, Academia Română – ICIA București  
Calea 13 Septembrie, nr. 13, 050711, București  
E-mail: [vergi@racai.ro](mailto:vergi@racai.ro)

**Rezumat.** Din perspectiva identificării automate a cuvintelor derivate și a bazelor lor de derivare în wordnetul românesc, în vederea îmbogățirii acestuia cu relații derivaționale și etichete semantice asociate acestora, prezentăm în acest articol rezultatele unui studiu de caz al cărui scop a fost identificarea automată a sufixelor cu care se formează cuvinte în limba română. O prezentare succintă a fenomenului derivării așa cum se manifestă el în română deschide acest articol, anticipând și provocările pe care le întâmpină studiul nostru. Pentru identificarea automată a sufixelor în cuvinte am folosit arborii de sufixe generalizați sub forma cărora am reprezentat lemele dintr-un lexicon electronic românesc. Am impus o serie de filtre, stabilite prin metoda observației, cu care am îmbunătățit rezultatele. Pentru evaluarea sufixelor identificate automat am folosit liste de referință de sufixe, grupate în funcție de partea de vorbire a cuvântului obținut prin derivare. Evaluarea rezultatelor am făcut-o în trei etape, prin comparare doar cu lista de sufixe, apoi cu lista de sufixe îmbogățită cu sufixoide și, în final, cu lista unificată a sufixelor și sufixoidelor pentru părțile de vorbire care prezintă omonimie în română. Am prezentat felul în care precizia și acuratețea algoritmului variază în funcție de limita impusă asupra productivității sufixelor. Studiul de față pune în evidență ansamblul de cunoștințe necesare pentru recunoașterea structurii morfematice a cuvintelor, dificultățile întâmpinate în acest proces, dar și importanța acestui studiu pentru cercetările lingvistice, pentru domeniul Inteligenței Artificiale, în sarcini precum regăsirea informației, rezumarea textelor, răspunsul automat la întrebări și toate celelalte care se bazează pe prelucrarea limbajului natural.

**Cuvinte cheie:** sufix, derivare, identificarea radicalului cuvintelor, analiză morfematică, arbore de sufixe generalizat.

## 1. Introducere

Vocabularul (i.e. totalitatea cuvintelor) unei limbi se modifică de-a lungul timpului, iar explicația rezidă în schimbările care se petrec în viața oamenilor, a căror comunicare ține pasul cu aceste modificări.

Compartimentul limbii cel mai susceptibil de a suferi variații în timp este vocabularul. Cuvintele existente se pot perima, își pot schimba sau îmbogăți sensurile, iar alte cuvinte pot apărea. În lingvistică se vorbește despre două modalități de îmbogățire a vocabularului: o modalitate internă (reprezentată de derivare, compunere și schimbarea categoriei gramaticale) și alta externă (împrumutul); alături de ele mai există calcul lingvistic, un procedeu mixt. Acestea sunt principalele mijloace de îmbogățire a vocabularului, alături de ele existând și altele, considerate secundare și pe care le lăsăm deoparte aici, întrucât obiectul preocupărilor noastre îl formează derivarea.

Aceasta se definește ca procedeu de formare a unui cuvânt nou având ca bază un alt cuvânt existent în limbă (Marouzeau, 1933: 63 ș.u.). Are două subtipuri: derivarea progresivă sau propriu-zisă și derivarea regresivă. Cea dintâi constă din adăugarea unor afixe (sufixe sau prefixe) la un cuvânt bază. În funcție de poziția din cuvântul bază în care se atașează afixul, vorbim despre prefixare (afixul, i.e. prefixul, se atașează înaintea bazei, de ex. *străbun*<prefixul *stră-*+*bun*) sau sufixare (afixul, i.e. sufixul, se atașează la sfârșitul bazei, de ex. *scriitor*<*scrie*+sufixul *-tor*). Derivarea regresivă constă în suprimarea unor afixe de la un cuvânt bază (de ex., *deranj*<*deranja*).

Derivarea este studiată de lingviștii interesați de evoluția vocabularului, de identificarea afixelor unei limbi și de caracterizarea acestora sub diverse aspecte (etimologie, sens, variații fonetice, productivitate etc.).

Derivarea prezintă interes și pentru inginerii din domeniul prelucrării limbajului natural. Sarcini precum regăsirea informației (engl. *information retrieval*), răspunsul automat la întrebări (engl. *question answering*) și altele folosesc informații despre cuvintele care aparțin aceleiași familii lexicale (cuvinte derivate de la aceeași bază cu afixe diferite). Asemenea informații se obțin prin reducerea cuvintelor derivate la baza lor, prin înlăturarea tuturor afixelor conținute; dacă în urma acestei înlăturări se obține același șir de litere, atunci cuvintele de la care s-a pornit au mari șanse de a aparține aceleiași familii lexicale<sup>1</sup> și sunt tratate împreună. În unele limbi, printre ele și româna, în acest șir de litere sunt permise alternanțe vocalice și/sau

---

<sup>1</sup> Cele două cuvinte nu aparțin aceleiași familii lexicale dacă bazele lor nu reprezintă același cuvânt, ci sunt omonime.

consonantice. De exemplu, cuvântul *băiețel* este derivat de la baza *băiat* cu sufixul *-el*; rădăcina prezintă atât o alternanță vocalică (ia:ie), cât și una consonantică (t:ț).

În acest articol ne-am propus să investigăm în ce măsură este posibilă identificarea automată a afixelor limbii române pornind de la un inventar de cuvinte la formele lor de bază și supunându-le unui proces de analiză și filtrare a rezultatelor.

Teoretic, putem folosi același algoritm (cu adaptări impuse doar de poziția din cuvânt asupra căreia ne concentrăm, i.e. început sau sfârșit) pentru identificarea prefixelor și sufixelor. Practic, însă, am constatat că rata de eroare este mult mai mare în cazul în care încercăm să identificăm prefixele. Acest lucru ne-a obligat să ne rezumăm în acest articol la un studiu de caz, și anume posibilitatea identificării automate a sufixelor în cuvinte.

Pentru început, se impun câteva precizări terminologice. O distincție importantă este aceea între rădăcină (engl. *root*) și radical (engl. *stem*). Rădăcina unui cuvânt este segmentul comun tuturor cuvintelor care formează o familie lexicală. Este un morfem indivizibil, minimal, purtător de sens lexical. De exemplu, în cazul familiei de cuvinte *pădure*, *pădurice*, *pădurar*, *păduros*, *împăduri*, *împădurire*, *împădurit*, *reîmpăduri*, *reîmpădurire*, *reîmpădurit* rădăcina este *pădur*, comună tuturor membrilor familiei. Radicalul se definește ca acea parte a cuvântului în structura căreia se regăsește rădăcina plus/minus toate afixele derivative. Astfel, în cazul membrilor familiei lexicale de mai sus, avem următorii radicali: *pădur* (identic cu rădăcina, deci) în cazul lui *pădure*, *păduric* (rădăcină plus sufix) în cazul lui *pădurice*, *pădurar* (rădăcină plus sufix) în cazul lui *pădurar*, *păduros* (rădăcină plus sufix) în cazul lui *păduros*, *împădur* (prefix plus rădăcină) în cazul lui *împăduri*, *împădurir* (prefix plus rădăcină plus sufix plus sufix) în cazul lui *împădurire* și așa mai departe. În concluzie, în cazul cuvintelor nederivate, rădăcina și radicalul sunt identice, iar în cazul cuvintelor derivate, radicalul conține rădăcina și toate afixele derivative.

În Prelucrarea Limbajului Natural (PLN) se face distincția între lematizare (engl. *lemmatization*) (i.e. identificarea formei de bază, de dicționar a unui cuvânt flexionat) și identificarea rădăcinii cuvintelor (engl. *stemming*). De exemplu, de la șirul de caractere *sfînțeniei* obținem prin lematizare *sfînțenie*, iar prin stemming rezultatul trebuie să fie *sfânt*, deoarece acestuia i s-a atașat sufixul *-enie* pentru a forma *sfînțenie*. În cazul

cuvintelor nederivate, precum *sfântului*, ambele procese dau același rezultat: atât prin lematizare, cât și prin stemming obținem *sfânt*.

## 2. Sistemul derivativ românesc

În limba română principalul mijloc intern de formare a cuvintelor noi este derivarea, iar dintre subtipurile sale cel mai frecvent întâlnit este sufixarea, după cum atestă lucrări ale specialiștilor din momente diverse din evoluția limbii (Pușcariu 1940, Hristea 1978, Avram 1989, Stoichițoiu-Ichim 2007). Prefixarea nu a fost foarte productivă în trecut, însă în cercetările care se concentrează asupra ultimelor decenii autorii constată o creștere a productivității acestui fenomen lingvistic, explicabil, printre altele (după părerea noastră), prin dorința sau nevoia vorbitorilor de concizie în comunicare.

Prefixele se atașează la începutul unui cuvânt direct sau prin modificarea corpului lor fonetic, prin fenomene de asimilare (de exemplu, *-în* devine *-îm* în contextul unei consoane labiale) sau disimilare (de exemplu, *des-* devine *de-* în *desăra*). Ele nu modifică structura fonetică a bazei.

Sufixe se atașează la sfârșitul cuvântului bază direct sau în urma unor fenomene fonetice petrecute în corpul bazei: eliminări de sunete (de ex. *bucura* > *bucuros*), alternanțe vocalice sau/și consonantice (de ex. *sfânt* > *sfîntenie*). În general, rădăcina unui derivat este o formă de dicționar (sau leamnă, după engl. *lemma*). Există însă și cazuri când aceasta este o formă flexionară (de ex. forma de plural a unor substantive: *noduri* > *noduros*, participiul verbelor: *plâns* > *plânset*).

Sufixarea poate schimba clasa morfologică a bazei. Astfel, cu ajutorul sufixelor putem forma substantive de la alte substantive (*prunc* > *pruncie*), de la verbe (*vorbă* > *vorbire*), de la adjective (*tineri* > *tineret*) sau de la adverbe (*împrejur* > *împrejurime*); verbe putem obține prin sufixare de la alte verbe (*linge* > *linguși*), de la substantive (*mușama* > *mușamaliza*), de la adjective (*orb* > *orbeca*) sau de la adverbe (*în luntre* > *în luntroșă*); adjective obținem de la alte adjective (*nițel* > *nițeluș*), de la substantive (*copil* > *copilăresc*), de la verbe (*întuneca* > *întunecos*), de la adverbe (*anevoie* > *anevoios*) sau de la numerale (*patru* > *pătrar*); adverbe obținem de la alte adverbe (*încet* > *încetinel*), de la substantive (*cruce* > *cruciș*), de la verbe (*fura* > *furiș*) sau de la adjective (*chior* > *chiorăș*).

Prefixarea singură nu schimbă, de obicei, clasa morfologică a bazei.

În studiul nostru am avut în vedere numai clasele deschise de cuvinte, adică substantivele, adjectivele, verbele și adverbele, conștienți de limitarea analizei. Dată fiind însă productivitatea slabă a derivării implicând restul părților de vorbire, considerăm că rezultatele nu sunt decât în mică măsură afectate.

Alături de prefixe și sufixe, în lingvistică se discută și despre prefixoide și sufixoide: elemente formative care dau impresia unor prefixe, respectiv sufixe. Sunt foarte slab productive în limba română, dar se pot recunoaște în structura unor cuvinte create în alte limbi sau după modelul cuvintelor din alte limbi. Sfera lor de întrebuintare este limitată, în general, la domeniul științific.

Din punct de vedere funcțional, ca și afixe, afixoidele se pot atașa la începutul sau la sfârșitul unui cuvânt. Din punct de vedere lexico-semantic, rezultă un nou cuvânt, la al cărui sens contribuie sensul tuturor elementelor participante la formarea sa.

Spre deosebire de afixe care se atașează numai la cuvinte, afixoidele se pot combina între ele, rezultând astfel cuvinte formate exclusiv din prefixoide și sufixoide: *kinetoscop*. În plus, sufixoidele pot deveni prefixoide și invers, schimbându-și locul în structura cuvântului rezultat: *rusofil* versus *filorus*.

Din punctul de vedere al productivității, afixoidele sunt mult mai slab productive decât afixele, care participă la formarea a numeroase cuvinte în română.

În general, afixele nu au utilizare restrânsă la anumite domenii de activitate sau stiluri funcționale. Excepție fac trei sufixe din domeniul medical: *-ită*, *-om*, *-oză*; creațiile spontane de felul lui *chiulangită* demonstrează statutul lor de sufixe recunoscute de vorbitori, care le transferă în limba comună, chiar în aspectul ei neîngrijit. În schimb, afixoidele sunt utilizate preponderent în limbajele științifice.

Date fiind aceste deosebiri între afixe și afixoide, lingviștii români tratează cuvintele formate cu afixe ca fiind derivate, iar pe cele cu afixoide ca fiind la limita între compunere și derivare, mulți înclinând spre includerea lor în cel dintâi fenomen.

### 3. Resurse și instrumente folosite

Studiul efectuat se bazează pe puține resurse lingvistice pe care le-am avut la dispoziție, în comparație cu cantitatea mare și diversă a informațiilor pe care le folosește lingvistul pentru a analiza morfematic un cuvânt. Considerăm că cea mai importantă cunoștință este de tip semantic: cuvintele derivate și baza lor au un conținut semantic parțial comun, conținut cu care participă baza de derivare. Chiar și un elev din școala primară va analiza cuvântul *bălos* ca derivat de la *bale* cu sufixul *-os*, nu de la *bal* cu sufixul *-os*, chiar dacă sufixul *-os* formează frecvent adjective de la baze substantivale, iar alternanța fonetică a:ă este frecventă. Studiul nostru nu a avut la dispoziție informații despre semantica cuvintelor implicate, fapt ce și-a lăsat amprenta în mod vizibil asupra rezultatelor.

#### 3.1 Resurse lingvistice

Baza studiului nostru o reprezintă un lexicon ce conține formele flexionare ale cuvintelor din limba română, însoțite de forma lor de dicționar (lema) și de descrierea morfosintactică. Lexiconul a fost creat în cadrul ICIA în proiectul MULTEXT-EAST și conține 1.222.832 de intrări (la data scrierii acestui articol), reflectând toate părțile de vorbire. O intrare constă într-o formă ocurență (flexionată în cazul cuvintelor flexionare), lema acesteia (cuvântul-titlu) și descrierea morfosintactică (Tufiș et al., 1997). Iată un fragment:

Tabelul 1. Fragment din lexiconul cu cuvinte-titlu, forme flexionare și etichete morfosintactice

Formă flexionată	Cuvânt-titlu	Etichetă morfosintactică
circuitului	circuit	Nemsoy
circul	circ	Nemsry
circul	circula	Vmip1s
circulă	circula	Vmip3
circulă	circula	Vmis3s

Prima coloană conține forme ocurență. În cea de-a doua se află lemele sau cuvintele-titlu corespunzătoare formelor din prima coloană, iar apoi urmează descrierea morfosintactică a acestora. Astfel, *circuitului* are lema *circuit* și din descrierea morfosintactică aflăm că este substantiv (N) comun

(c) de genul masculin<sup>2</sup> (m), la singular (s), are formă de caz oblic (o) și este articulat hotărât (y).

Nu am folosit acest lexicon în întregime, ci ne-am rezumat doar la formele de bază ale cuvintelor aparținând claselor deschise (substantive, adjective, verbe, adverbe). Am ignorat, deci, formele flexionare, cu riscul asumat de a nu regăsi fapte de tipul *colțuri*>*colțuros*, adică derivate de la forme flexionare; considerăm însă că sufixul *-os*, cu același sens, se combină și cu alte cuvinte, la forma lor de bază, deci va putea fi regăsit (și chiar a fost regăsit) în acelea: de exemplu: *băț*>*bățos*. De asemenea, nu vom putea ști în câte derivate apare fiecare sufix. Dar, după cum se va vedea mai jos, în prezentarea rezultatelor, nici acest lucru nu afectează studiul de caz, el interferând cu alte aspecte.

Pe lexicon am rulat un algoritm de identificare a sufixelor. Rezultatele au fost comparate, în vederea evaluării, cu o listă de sufixe creată manual, pe baza unor lucrări de referință despre derivare în limba română (*Formarea cuvintelor în limba română* 1970, 1978, 1989, Pascu 1916, Coteanu 2007, Philippide 2011, Tudose 1978, Iordan 1939, *Studii și materiale despre formarea cuvintelor în limba română* 1959, 1967, 1969). Listele acestea create manual le numim liste de referință (engl. *gold standard*) și pornim de la asumptiunea că ele au caracter exhaustiv. În tabelul de mai jos ilustrăm numărul de sufixe conținute de fiecare listă.

Tabelul 2. Numărul de sufixe conținute de listele standard.

Tipul de sufixe	Numărul de sufixe din lista de referință
substantivale	260
adjectivale	104
verbale	104
adverbiale	14
TOTAL	482

<sup>2</sup> Substantivul *circuit* este de genul neutru. Conform specificațiilor de adnotare, substantivele neutre la singular sunt marcate ca având genul masculin, iar la plural ca având genul feminin. Distincția între substantivele masculine, feminine și neutre se face astfel: substantivele marcate cu m și i la singular și la plural sunt masculine, cele adnotate cu f și i la singular și la plural sunt feminine, iar cele adnotate cu m la singular și cu f la plural sunt neutre.

### 3.2. Descrierea algoritmului

Pentru identificarea automată a sufixelor în cuvintele din lexicon am decis să lucrăm cu arbori de sufixe (engl. *suffix trees*) (Ukkonen, 1995). Aceștia prezintă două avantaje majore: rapiditate și utilizare eficientă a memoriei.

Un arbore de sufixe este o reprezentare arborescentă a unui string (în cazul nostru, un cuvânt) și a tuturor sufixelor sale. Sufixele sunt considerate a fi ceea ce rămâne dintr-un cuvânt după eliminarea succesivă a câte unei litere de la începutul acestuia. Astfel, pentru cuvântul *banană* avem reprezentarea arborescentă din Figura 1. Sufixele acestui cuvânt sunt: *banană*, *anană*, *nană*, *ană*, *nă*, *ă*.

Menționăm că figurile 1 și 2 au fost create cu ajutorul Generalized Suffix Tree Java Applet disponibil la adresa <http://illya-keeplearning.blogspot.com/2009/06/generalized-suffix-trees-java-applet.html> (consultată în august 2011).

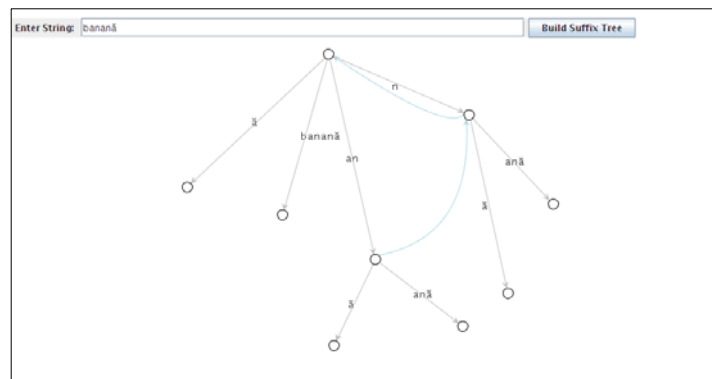


Figura 1. Arborele de sufixe pentru cuvântul *banană*.

Ipoteza studiului nostru este că sufixele sunt productive și că există, deci, în lexicon o mulțime de cuvinte formate cu același sufix. De aceea, mult mai potrivit pentru studiul nostru este un tip special de arbori de sufixe, și anume arborii de sufixe generalizați (engl. *generalized suffix trees*) (Gusfield, 1999). Aceștia se construiesc pentru o mulțime de cuvinte și reprezintă toate sufixele acestora. Iată cum arată arborele de sufixe generalizat pentru cuvintele *cădere* și *citire*. De obicei, pentru a construi un asemenea arbore, se adaugă un simbol ca marcator de sfârșit de cuvânt. Cuvintele din mulțimea pentru care se creează arborele sunt împărțite în sufixe, dar pentru



toate aceste cuvinte se construiește un singur arbore. Sufixele comune (aici *-re*) sunt recunoscute ca fiind acele sufixe care diferă doar prin simbolul terminal (aici \$ și #).

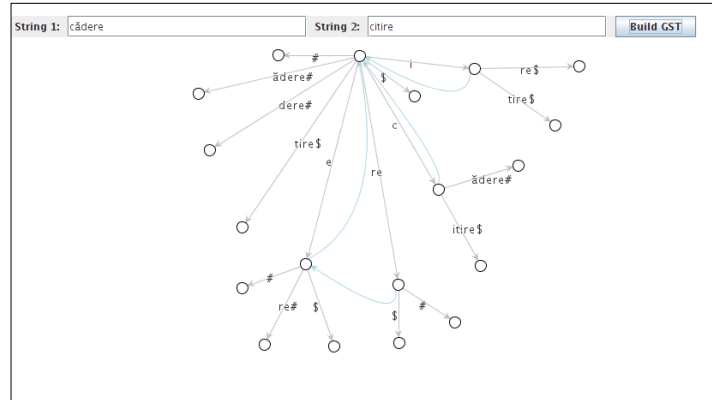


Figura 2. Arborele de sufixe generalizat pentru cuvintele cădere și citire.

Primul pas în identificarea sufixelor îl constituie crearea unui index Lucene pentru lexiconul cu care lucrăm. Acest index conține documente. Fiecare document are trei câmpuri: forma de bază sau lema cuvântului, rădăcina și eticheta morfosintactică. Am impus trei restricții asupra cuvintelor din documentele din index: o restricție asupra părților de vorbire (am lucrat numai cu clasele de cuvinte deschise, adică cu substantive, verbe, adjective și adverbe), o restricție morfologică (am selectat numai lemele acestor cuvinte, adică: forma de nominativ singular nearticulat a substantivelor, forma de infinitiv a verbelor, de masculin singular nominativ a adjectivelor; adverbele au o singură formă în română, nu fac parte dintre părțile de vorbire flexibile) și o restricție de formă (adică am exclus: cuvintele scrise cu cratimă, de ex. *floarea-soarelui*, cu liniuță de subliniere – engl. *underscore* – de ex. *praf de pușcă*, termeni considerați compuși prin alăturare, cu puncte, de ex. *pp.* cu lema *pagină*, care a fost adnotat ca substantiv, în contradicție cu toate celelalte prescurtări care sunt marcate cu o etichetă specială pentru prescurtări, și anume Y, sau cuvintele cu apostrof, de ex. *dom'* pentru *domn*).

Am creat acest index pentru a face mai rapidă căutarea și pentru a putea modifica tipul căutării, atunci când este nevoie.

În Tabelul 2 prezentăm numărul de leme cu care am lucrat, grupate pe părți de vorbire. Chiar dacă nu este un fenomen extrem de frecvent, omonimia (i.e. relația între cuvinte cu aceeași formă, dar sens diferit) există în română și se stabilește între diverse părți de vorbire: substantive și adjective (de exemplu *nebun*), adverbe și adjective (de pildă *liniștit*), adjectiv, adverb și substantiv (precum *frumos*) etc. În consecință, aceeași leme poate avea mai multe adnotări morfosintactice. Pe ultima linie a acestui tabel am notat numărul de leme unice cu care am lucrat.

Tabelul 3. Numărul de leme utilizate, distribuite pe părți de vorbire.

Partea de vorbire	Numărul de leme
substantiv	67255
verb	8473
adjectiv	35285
adverb	1255
TOTAL	112.268
LEME UNICE	103.856

Indexul Lucene este folosit pentru a alege cu ușurință categoria gramaticală cu care să lucrăm și pentru a valida radicalii obținuți după îndepărtarea afixelor identificate.

Radicalul cuvintelor derivate cu sufixe se obține după următoarea procedură (inspirată din algoritmul Snowball (<http://snowball.tartarus.org/algorithms/romanian/stemmer.html>) de identificare a rădăcinilor din cuvinte românești): (i) atunci când *i* și *u* apar în poziție intervocalică sunt protejate prin marcarea lor cu literă majusculă, pentru a nu fi eliminate; (ii) cel mai lung dintre șirurile *a*, *e*, *i*, *ie*, *ă* este înlăturat; (iii) *i* și *u* sunt de-protejate prin transformarea lor înapoi în litere mici.

Apoi extragem din indexul Lucene lista tuturor cuvintelor aparținând unei părți de vorbire. Eliminăm primele două litere din acestea (deci nu vom trata cuvintele alcătuite doar din două litere; ele nu sunt susceptibile de a fi derivate), iar formele astfel scurtate le introducem într-un arbore de sufixe generalizat. Menționăm că am folosit implementarea *open source* a algoritmului de arbori de sufixe generalizați de la carrot<sup>2</sup> (disponibilă pe [situl project.carrot2.org](http://project.carrot2.org)).

Apoi am căutat în arbore ramuri comune pentru mai multe cuvinte. Acestea sunt sufixele potențiale, pe care le-am introdus într-o tabelă (engl. *map*), împreună cu cuvintele în care apar, i.e. cuvintele probabil derivate cu

aceste sufixe. Am inspectat manual această tabelă și am constatat că este absolut necesară aplicarea unor filtre în procesul de identificare a sufixelor.

În primul rând, am impus o restricție asupra lungimii sufixelor, bazată pe observațiile noastre asupra listei standard de sufixe verbale: ele nu trebuie să fie mai lungi de cinci caractere pentru verbe, adjective și adverbe și șase caractere pentru substantive.

În al doilea rând, validăm radicalii cuvintelor: verificăm dacă șirul de caractere rămas după îndepărtarea potențialului sufix este un radical sau o lemă în indexul Lucene. Păstrăm doar acele sufixe care apar cu cel puțin un radical valid sau cu o lemă validă.

În al treilea rând, atunci când sunt posibile mai multe sufixe pentru același cuvânt (de ex. *spălătoreasă* are două sufixe: *-tor* și *-easă*) păstrăm ca valid sufixul cel mai scurt. Procedăm astfel considerând că în cazul unui cuvânt cu două sufixe, derivarea s-a produs succesiv, în două etape. Suntem conștienți că astfel nu mai identificăm eventualele sufixe compuse din limba română.

În al patrulea rând, păstrăm doar sufixele care apar în cel puțin zece cuvinte, nu neapărat valide toate. Motivația este că scopul nostru este să găsim sufixe productive pe teren românesc. În finalul acestui articol prezentăm și variații ale preciziei și acurateței algoritmului în funcție de modificarea limitei de productivitate a sufixelor: am rulat programul impunând, pe rând, productivitate 20, 30 și 40 asupra potențialelor sufixe.

În al cincilea rând, am observat că sufixele de patru sau cinci caractere care apar în mai puțin de șapte radicali valizi nu sunt cu adevărat sufixe și le-am eliminat.

Pentru fiecare sufix, am calculat numărul total de cuvinte în care apare și numărul de cuvinte (cu radical sau lemă) valide în care apare. Apoi am calculat un scor pentru fiecare sufix, împărțind numărul de ocurențe valide la numărul total de ocurențe. Acest scor reprezintă probabilitatea ca respectivul șir de litere aflat la finalul unui cuvânt să reprezinte un sufix.

#### 4. Rezultate și interpretări

Prezentăm mai jos câteva exemple din fișierul de sufixe substantive:

oar ( $19 / 35 = 0.54$ ): [balans+oar, ben+oar, bud+oar, cont+oar, cul+oar, eboș+oar, ferm+oar, fum+oar, hematoz+oar, hidroz+oar, lav+oar, muz+oar, patin+oar, pis+oar, polis+oar, poș+oar, raz+oar, remont+oar, urin+oar]

uleț (79 / 147 = 0.54): [bob+uleț, ciob+uleț, cioc+uleț, colț+uleț, corb+uleț, corn+uleț, cort+uleț, coș+uleț, crâng+uleț, cub+uleț, cuib+uleț, drug+uleț, dâmb+uleț, foc+uleț, frig+uleț, ghiob+uleț, golf+uleț, grup+uleț, hud+uleț, ied+uleț, joc+uleț, lemn+uleț, moș+uleț, muc+uleț, murg+uleț, mânz+uleț, nor+uleț, om+uleț, orz+uleț, partid+uleț, piept+uleț, plug+uleț, pod+uleț, pom+uleț, porc+uleț, prost+uleț, prunc+uleț, prânz+uleț, puf+uleț, punct+uleț, pâlc+uleț, roib+uleț, rug+uleț, runc+uleț, răc+uleț, schit+uleț, scoc+uleț, sfert+uleț, sfânt+uleț, smoc+uleț, snop+uleț, solz+uleț, somn+uleț, soț+uleț, spin+uleț, stog+uleț, strop+uleț, struț+uleț, strămb+uleț, sturz+uleț, stâlp+uleț, stârc+uleț, sur+uleț, sânt+uleț, sân+uleț, tiv+uleț, toc+uleț, tren+uleț, trunchi+uleț, tub+uleț, turc+uleț, turn+uleț, urs+uleț, vin+uleț, șip+uleț, șir+uleț, șoim+uleț, șorț+uleț, ștuc+uleț]

asă (136 / 173 = 0.79): [arhondăr+asă, arăm+asă, aurăr+asă, baron+asă, blănăr+asă, boier+asă, bor+asă, brutăr+asă, brânzăr+asă, bucătăr+asă, bumbăcăr+asă, băbăr+asă, băcăn+asă, cenușăr+asă, ciorăpăr+asă, ciubotăr+asă, cizmăr+asă, clucer+asă, cofetăr+asă, cojocăr+asă, colonel+asă, covrigăr+asă, croitor+asă, ctitor+asă, cucul+asă, culbec+asă, cusător+asă, cârciumăr+asă, cârnățăr+asă, călcător+asă, cămățăr+asă, căpităn+asă, cărtură+asă, doctor+asă, doftor+asă, dubălăr+asă, fermecător+asă, florăr+asă, franzelăr+asă, făină+asă, gal+asă, giuvaierger+asă, gornic+asă, grădinăr+asă, găină+asă, hătmăn+asă, irod+asă, isprăvnic+asă, jimblăr+asă, jitar+asă, jud+asă, judec+asă, judecător+asă, județ+asă, jupân+asă, legător+asă, lenjer+asă, leurd+asă, lăptăr+asă, maior+asă, mezelăr+asă, meșter+asă, mir+asă, morăr+asă, moșier+asă, măcelăr+asă, neder+asă, negustor+asă, neguțător+asă, notăr+asă, năsturăr+asă, ofițer+asă, olăr+asă, ospățăr+asă, pair+asă, pantofăr+asă, papucăr+asă, pescăr+asă, pețitor+asă, pităr+asă, plugăr+asă, plăpumăr+asă, poftitor+asă, polcovnic+asă, pop+asă, popicăr+asă, portăr+asă, porumbăr+asă, postelnic+asă, postăvăr+asă, potcovăr+asă, potropop+asă, preot+asă, primăr+asă, profesor+asă, protopop+asă, puric+asă, pâină+asă, pârcălăb+asă, păcurăr+asă, păhărnice+asă, pălărier+asă, rotăr+asă, rudăr+asă, serdăr+asă, servitor+asă, sfârănăr+asă, sluger+asă, spițer+asă, spoitor+asă, spălător+asă, spătăr+asă, stolnic+asă, strungăr+asă, stăpân+asă, sânger+asă, săpunăr+asă, sărăr+asă, tinichiger+asă, tutunger+asă, tâmplăr+asă, tăbăcăr+asă, vier+asă, vod+asă, voievod+asă,

voivod+easă, vornic+easă, vrăjitor+easă, văcăr+easă, zidăr+easă, zlătăr+easă, învățător+easă, șoric+easă, șuncăr+easă, șurubăr+easă, șătrăr+easă]

Primul este un exemplu de terminație care nu este sufix în limba română: *-oar*. El este un fals sufix, identificat de algoritm întrucât șirurile de caractere rămase prin îndepărtarea lui *-oar* din cuvintele enumerate sunt leme sau radicali valizi în limba română. Cu toate acestea, nu există nicio legătură derivativă între acești radicali și cuvintele de la care s-a pornit, deși una semantică poate fi găsită în anumite cazuri: între *balans* și *balansoar*, între *fum* și *fumoar* etc. Vorbitorii de limbă română nu au decupat din aceste împrumuturi terminația *-oar* pe care să o folosească apoi la crearea altor cuvinte, adică să o trateze ca pe un sufix.

În schimb, *-uleț* și *-easă* sunt sufixe în limba română, chiar dacă nu în toate cuvintele validate de algoritm și enumerate mai sus.

După aplicarea filtrelor descrise mai sus am obținut următoarele rezultate:

Tabelul 4. Rezultatele după rularea algoritmului și aplicarea filtrelor.

Partea de vorbire	Total sufixe identificate	Sufixe corecte	Precizie (% corecte din total)	Acuratețe (% corect identificate din listele standard)
substantive	888	161	18,13	61,92
verbe	209	81	38,76	77,88
adjective	320	54	16,88	51,92
adverbe	46	3	6,52	21,43

Printre rezultate găsim și false pozitive. Unele terminații ale cuvintelor sunt identificate ca sufixe greșite. De exemplu, *metresă* este identificat ca derivat de la *metru* cu sufixul *-esă*. Explicația rezidă în precaritatea diversității resurselor lingvistice folosite: singura resursă pe care se bazează este lexiconul de forme. Nu folosește nicio resursă fonetică sau semantică. Un verb ca *noroi* este analizat ca fiind format din *nor* cu sufixul *-oi*. Cuvintele derivate și bazele lor au în comun un anumit conținut semantic. Or, cum algoritmul nostru nu folosește nicio informație de tip semantic, astfel de analize sunt inevitabile.

Parcurgând manual listele de sufixe identificate, am constatat următoarele:

- Efectuând și o validare manuală a rezultatelor, am identificat și sufixe inexistente în listele de referință. De exemplu, *-ăt* în *strigăt*

(<*striga*+*-ăt*) sau *-atec*, variantă a sufixului *-atic* : *singuratec*. Mult mai interesant este cazul sufixului *-iadă* (nediscutat în lucrările de specialitate), care a devenit productiv în română mai târziu, creând cuvinte precum *universiadă*, *mineriadă*. Algoritmul nostru identifică aceste sufixe, dar nu le regăsește în lista de referință, ca atare nu le validează.

- Lipsa unei resurse fonetice care să îmbunătățească rezultatele duce la analize greșite: *barochist* este analizat în *baroc* + *-hist*. În limba română grupurile de litere *ch* și *gh* nu reprezintă niciodată două sunete. Algoritmul poate fi îmbunătățit adăugând acest filtru: *ch* sau *gh* nu se separă astfel încât *c* sau *g* să aparțină rădăcinii, iar *h* sufixului.
- Unele cuvinte prezintă omonimie: aceeași formă are valori aparținând mai multor clase morfologice. Neincluderea unui sufix decât în lista standard a unei părți de vorbire, ignorând astfel conversiunea, scade precizia și acuratețea. Soluția noastră a fost să efectuăm o evaluare și prin unificarea listelor de referință ale părților de vorbire ce prezintă omonimii între ele. Rezultatele se găsesc mai jos, în Tabelul 5.
- Printre sufixele identificate, mai ales pentru substantive și adjective, se află și multe sufixoide, pe care nu le incluseserăm în lista standard. Ca atare, ele nu au fost validate de algoritm. De exemplu: *-fag* din *fotofag*, *-cid* din *zoocid*, *-urie* din *alcoolurie* și altele. Pentru unele, analiza este greșită: în *aglutinogen* este identificat sufixul *-ogen* (corect ar fi *-gen*), în *anglofob* sufixul *-ofob* (corect este *-fob*); lingviștii trec vocala de legătură (în exemplele de mai sus *o*) la partea inițială a cuvântului, nu la sufixoid. În fond, atunci când această vocală nu face parte din radicalii cuvintelor, ea poate fi decupată alături de elementul din față sau de cel următor.

Inițial, lista noastră de referință a inclus numai sufixe. Întrucât unei mașini îi este greu să distingă între sufix și sufixoid (mai ales că unele au statut controversat, unii lingviști considerându-le sufixe, alții sufixoide), am decis să includem și sufixoidele alături de sufixe, pentru a verifica în ce măsură se modifică rezultatele.

În consecință, listele standard ale sufixelor se modifică după cum urmează: cea a substantivelor a fost îmbogățită cu 151 de sufixoide,

ajungând astfel la 411 elemente, iar cea a adjectivelor cu 73 de sufixoide, numărând în final 177 de elemente. Celelalte liste (ale sufixelor verbale și adverbiale) au rămas nemodificate. Comparând rezultatele algoritmului cu aceste noi liste standard constatăm următoarele rezultate în cazul substantivelor și adjectivelor (pentru celelalte părți de vorbire nu înregistrăm modificări, în mod evident):

Tabelul 5. Evaluarea prin luarea în considerare și a sufixoidelor.

Partea de vorbire	Total sufixe identificate	Sufixe corecte	Precizie (% corecte din total)	Acuratețe (% corect identificate din listele standard)
substantive	888	188	21,17	45,74
adjective	320	77	24,06	43,50

Comparând datele din Tabelul 4 cu cele din Tabelul 5 reiese că precizia a crescut ușor (peste 3% pentru substantive și cu aproape 8% pentru adjective), în detrimentul acurateței (care a scăzut cu aproape 17% pentru substantive și cu peste 8% pentru adjective).

Pentru a acoperi și cazurile de omonimie între părți de vorbire diferite, am comparat sufixele identificate de algoritm cu listele de referință unificate astfel: substantive cu adjective și adjective cu adverbe. Adică, lista sufixelor substantivale identificate de algoritm am comparat-o cu lista unificată a sufixelor (și sufixoidelor) substantivale și adjectivale, lista celor adjectivale am comparat-o tot cu lista unificată a sufixelor (și sufixoidelor) substantivale și adjectivale, lista sufixelor adverbiale am comparat-o cu lista unificată a sufixelor (și sufixoidelor) adjectivale și adverbiale. Am obținut datele din tabelul următor:

Tabelul 6. Evaluarea prin luarea în considerare și a omonimiilor

Parte de vorbire	Total sufixe (și sufixoide) unice în lista standard unificată	Total sufixe găsite	Bune găsite	Precizie (% Bune găsite din Total găsite)	Acuratețe (% Bune găsite din lista standard)
substantive	513 (subst.+adj.)	888	214	24,10	41,72
adjective	513 (subst.+adj.)	320	89	27,81	17,35
adverbe	183 (adj.+adv.)	46	16	34,78	8,74

și în felul acesta observăm, comparând datele din Tabelul 6 cu cele din Tabelele 4 și 5, că precizia crește, pentru substantive, cu aproape 3% față de datele din Tabelul 5, deci față de cazul cu luarea în considerare și a sufixoidelor, respectiv cu 7% față de cazul inițial; pentru adjective, precizia crește cu aproape 4% față de cazul din Tabelul 5, deci cu aproape 11% față de cazul inițial; pentru adverbe precizia crește de peste 5 ori față de cazul inițial; pentru toate părțile de vorbire acuratețea scade: cu 4%, respectiv cu 20% pentru substantive față de cazul din Tabelul 5, respectiv din Tabelul 4; pentru adjective ea scade dramatic, de la 45% în Tabelul 5, respectiv de la 51% în Tabelul 4, la 17% în ultimul caz; și la adverbe scade de 3 ori față de cazul inițial.

În graficele din figurile 3, 4 și 5 prezentăm variația preciziei și acurateții în funcție de limita impusă asupra productivității sufixelor: toate datele din tabelele de mai sus sunt redată pentru productivitate 10. Adică permitem un număr minim de 10 cuvinte în care potențialul sufix să apară, pentru a fi luat în considerare în procesul de validare. În graficele de mai jos am analizat variațiile celor două mărimi atunci când am impus o productivitate diferită asupra prefixelor, adică atunci când ele apar în 10, 20, 30, respectiv 40 de cuvinte. După cum se observă, cu cât creștem limita impusă productivității sufixelor, cu atât precizia este mai mare, iar acuratețea mai mică. Modul în care ele variază diferă de la o parte de vorbire la alta, după cum se poate observa după formele curbelor în cele trei figuri. În cazul adverbilor, oricum numărul de sufixe găsite este foarte mic; dacă modificăm limita impusă productivității sufixelor, găsim tot 3 sufixe corecte pentru productivitate 20, iar pentru productivitate 30 găsim doar 2, în vreme ce pentru productivitate 40 găsim doar un sufix. Pentru adverbe nu am mai redat un grafic.



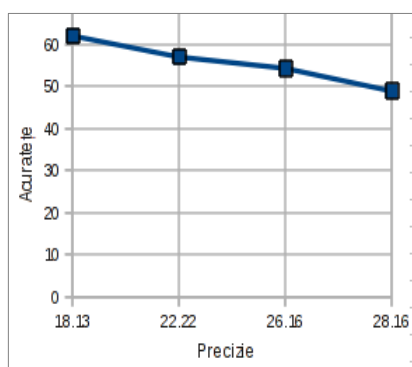


Figura 3. Evoluția preciziei și acurateții în cazul substantivelor.

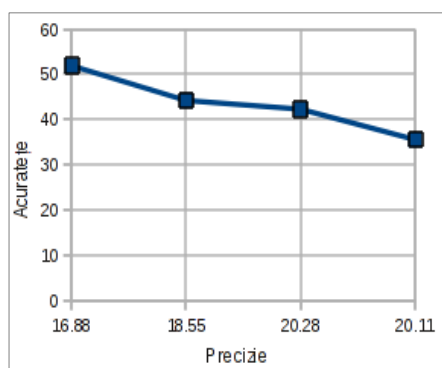


Figura 4. Evoluția preciziei și acurateții în cazul adjectivelor.

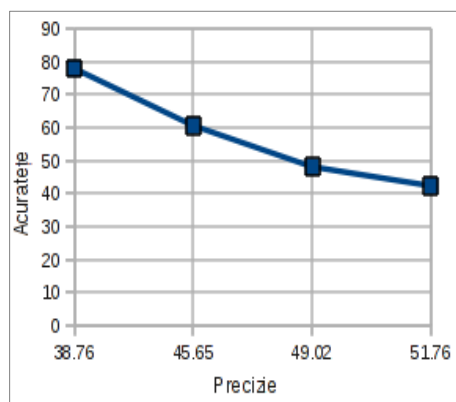


Figura 5. Evoluția preciziei și acurateții în cazul verbelor.

## 5. Constatări finale

Analiza morfematică a cuvintelor prezintă interes atât pentru lingviști (datori, printre multe altele, să dea explicații despre modul de apariție sau formare a cuvintelor), precum și pentru inginerii care, propunându-și să rezolve variate sarcini ce implică prelucrarea textelor, conștientizează valoarea legăturii lexico-semantice dintre cuvintele din aceeași familie lexicală și se angajează la îmbunătățirea rezultatelor activității lor prin valorificarea informațiilor oferite de aceasta.

Limba română este, în continuare, deficitară în ceea ce privește resursele lingvistice și instrumentele de prelucrare a ei: corpusurile existente nu au dimensiuni foarte mari, nu acoperă domenii și stiluri funcționale variate, nu au adnotări pe mai multe niveluri (morfologic, sintactic, semantic, pragmatic etc.). În diverse aplicații se simte nevoia unui analizor sintactic, fie el și de suprafață.

Având ca scop final al unui proiect mai amplu dezvoltarea unei rețele semantico-derivaționale pentru limba română, ce va putea fi folosită și ca dicționar de cuvinte derivate, ne-am propus să întocmim, pentru început, o listă cu afixele românești. O provocare în realizarea acestui obiectiv a constituit-o încercarea de identificare automată a acestora într-un dicționar de leme. În acest articol am prezentat modul în care am efectuat acest experiment și rezultatele obținute.

Bazându-ne exclusiv pe informații lexicale (am avut la dispoziție un lexicon de cuvinte românești din care am ales să lucrăm numai cu formele de bază ale substantivelor, adjectivelor, verbelor și adverbilor), ignorând orice fel de informație de tip fonetic (alternanțe), grafic (convenții de notare a unor sunete) și semantice (relații semantice între cuvinte), și folosind metoda arborilor de sufixe generalizați am extras terminații ale cuvintelor, adică potențiale sufixe. Pentru evaluarea metodei, prin calcularea preciziei și acurateții, am comparat aceste terminații cu listele de referință ale sufixelor românești, pe care le-am extras din lucrările dedicate derivării, mai exact sufixării, în limba română. Aceste liste de referință sunt disponibile la adresa [www.racai.ro/~vergi](http://www.racai.ro/~vergi) în categoria Research.

Analizând rezultatele obținute constatăm câteva lucruri esențiale: în primul rând, se observă imediat care sufixe sunt productive și care nu, precum și cât de productive sunt.

În al doilea rând, scorurile calculate de noi (prin împărțirea numărului total de cuvinte în care apare un sufix la numărul de cuvinte valide)

reprezintă probabilitatea ca terminația respectivă să fie sufix într-un cuvânt oarecare. Această informație și cea de dinainte sunt extrem de importante în diversele sarcini din Inteligența Artificială, în care identificarea derivatelor și reducerea lor la rădăcină au un rol important (de exemplu, rezumarea automată, răspunsul automat la întrebări etc.)

În al treilea rând, algoritmul a identificat și sufixe despre care lingviștii încă nu au scris (vezi *-iadă*). Astfel, am îmbogățit listele de referință ale sufixelor românești, este adevărat, cu un număr mic de elemente (*-ăt*, *-atec*, *-iadă*), însă ele sunt valoroase.

În al patrulea rând, se observă consecințele de neevitat ale nefolosirii informațiilor semantice: este identificat sufixul *-ier* din mai multe cuvinte (*afișier*, *textier*), iar la validare se acceptă, în mod eronat, *vistier* ca derivat de la *vist* (joc de cărți).

În al cincilea rând, se observă tratarea identică a derivării progresive, a celei regresive, a substituției de sufixe și a împrumuturilor: ceea ce contează este posibilitatea decupării sufixului dintr-un șir de caractere, nu și proveniența sa. Astfel, *problematiza*, considerat de lingviști derivat prin substituția sufixului *-ic* din *problematic* cu *-iza*, este analizat, în absența unui radical valid *problemat*, ca derivat de la *problemă* cu sufixul *-atiza*.

Așa cum am menționat pe parcursul articolului, metoda de lucru mai poate fi îmbunătățită. În plus, o nouă provocare poate fi identificarea prefixelor cuvintelor, după un algoritm asemănător, aducând ajustările necesare.

Din cunoștințele noastre, aceasta este prima raportare a unui experiment de identificare a sufixelor românești în cuvinte derivate. Alte experimente au vizat recunoașterea cuvintelor românești derivate, cunoscut fiind un inventar de afixe (Petic, 2011).

□i pentru alte limbi, s-a manifestat interesul pentru procesul de reducere a derivatelor la baza lor, util pentru îmbunătățirea rezultatelor în domeniul regăsirii informației (Porter, 1980). Totuși, lipsa resurselor lingvistice pentru diverse limbi a impulsivat cercetările legate de găsirea unei metode de identificare a morfemelor în cuvinte (Goldsmith, 2001, Majgaonker și Siddiqui, 2010, Freitag, 2005 ș.a.). Metodele variază de la învățarea nesupervizată (Goldsmith, 2001), la folosirea grupurilor de cuvinte (Freitag, 2005), a regulilor de segmentare (Porter, 1980), a împărțirii cuvintelor în n-gramme dintre care să se recunoască sufixele (Majgaonker și Siddiqui, 2010).

În contextul cercetărilor despre wordnet, îmbogățirea acestuia cu relații derivaționale a urmărit fie identificarea perechilor bază-derivat dat fiind un inventar de afixe (în general sufixe), pentru care au fost generate apoi etichetele semantice (Fellbaum et al., 2007, Bilgin et al., 2004), fie generarea derivatelor de la baze existente în wordnet cu ajutorul unor sufixe productive (Pala și Hlavackova, 2007, Kahusk et al., 2010).

Interesul pentru identificarea sufixelor unui cuvânt se explică, din perspectivă lingvistică, prin nevoia de a crea, într-un mod rapid și ieftin, resurse lingvistice pentru un număr cât mai mare și cât mai divers de limbi naturale. Din perspectiva interacțiunii om-calculator, această etapă de prelucrare a limbii are utilizări în mai multe subdomenii: generarea limbii naturale, răspunsul automat la întrebări, sumarizarea, regăsirea informației și altele.

## Mulțumiri

Această lucrare a fost realizată în cadrul proiectului "Valorificarea identităților culturale în procesele globale", cofinanțat de Uniunea Europeană și Guvernul României din Fondul Social European prin Programul Operațional Sectorial Dezvoltarea Resurselor Umane 2007-2013, contractul de finanțare nr. POSDRU/89/1.5/S/59758. Mulțumim pentru observațiile formulate pe marginea manuscrisului dlui academician Grigore Brâncuș și recenzorilor anonimi ai RRIOC, precum și lui Emanuele Pianta pentru discuțiile fructuoase purtate pe marginea acestui experiment.

## Referințe

- Avram, M., *Introducere în studiul sufixelor*, în Graur, Al., Avram, M. (coord.), *Formarea cuvintelor în limba română*, București, Editura Academiei, vol. III 1989.
- Bilgin, O., Cetinoglu, O., Oflazer, K., *Morphosemantic Relations in and across Wordnets: A Study Based on Turkish*, în P. Sojka, K. Pala, P. Smrz, C. Fellbaum, P. Vossen (eds.), *Proceedings of GWC*, 2004.
- Coteanu, I., *Formarea cuvintelor în limba română: derivarea, compunerea, conversiunea*, ed. Narcisa Forăscu, Angela Bidu-Vrânceanu, București, Editura Universității din București, 2007.
- Fellbaum, C., Osherson, A., Clark, P. E., *Putting Semantics into Wordnet's "Morphosemantic" Links*, în *Proceedings of the 3<sup>rd</sup> Language and Technology Conference*, Poznan, 2007.

- Freitag, D., *Morphology induction from term clusters*, în Proceedings of the ninth conference on computational natural language learning, p. 128-135, 2005.
- Goldsmith, J., *Unsupervised Learning of the Morphology of a Natural Language*, în Computational Linguistics, vol. 27, nr. 2, p. 153-198, 2001.
- Graur, Al. și Avram, M. (coord.), *Formarea cuvintelor în limba română*, București, Editura Academiei, vol. I 1970, vol. II 1978, vol. III 1989.
- Gusfield, D., *Algorithms on Strings, Trees and Sequences: Computer Science and Computational Biology*, Cambridge University Press, SUA, 1999.
- Hristea, Th., *Romanian Vocabulary and Etymology in Current Trends in Romanian Linguistics*, RRL, no. 1-4, 1978.
- Iordan, I., *Sufixe românești de origine recentă (neologisme)*, în Buletinul Institutului de filologie română „Alexandru Philippide”, vol VI, Iași, p. 1-59, 1939.
- Kahusk, N., Kerner, K., Vider, K., *Enriching Estonian WordNet with Derivations and Semantic Relations*, în Proceedings of the 2010 Conference on Human Language Technologies – The Baltic Perspective, 2010.
- Majgaonker, M. M., and Siddiqui, T. J., *Discovering Suffixes: A Case Study for Marathi Language*, în International Journal on Computer Science and Engineering, vol. 2, nr. 8, p. 2716-2720, 2010.
- Marouzeau, J., *Lexique de la terminologie linguistique*. Librairie orientaliste Paul Geuthner, Paris, 1933.
- Pala, K., Hlavackova, D., *Derivational Relations in Czech Wordnet*, în Proceedings of the Workshop on Balto-Slavonic, 2007.
- Pascu, G., *Sufixele românești*, București, Ed. Academiei Romane, 1916.
- Petic, M., *Automatizarea procesului de creare a resurselor lingvistice computaționale*, teză de doctorat, Chișinău, Academia de Științe a Moldovei, Institutul de Matematică și Informatică, 2011.
- Philippide, Al., *Istoria limbii române*, Iași, Polirom, 2011.
- Porter, M., *An Algorithm for Suffix Stripping*, în Program, 14, nr. 3, p. 130-137, 1980.
- Pușcariu, S., *Limba română*. Vol. I Privire generală, Fundația pentru Literatură și Artă "Regela Carol II", București, 1940, vol. II. Rostirea, Editura Academiei RSR, București, 1959.
- Stoichițoiu-Ichim, A., *Vocabularul limbii române actual. Dinamică, influențe, creativitate*, București, Editura BIC ALL, 2007.
- Studii și materiale privitoare la formarea cuvintelor în limba română*, vol I, Editura Academiei, vol I 1959, vol IV 1967, vol V 1969.
- Tudose, C., *Derivarea cu sufixe în româna populară*, București, Editura Universității din București, 1978
- Tușiș, D., Barbu, A.-M., Pătrașcu, V., Rotariu, G., Popescu, C., *Corpora and Corpus-Based Morpho-Lexical Processing*, în D. Tușiș, P. Andersen (eds.), Recent Advances in Romanian Language Technology, București, Editura Academiei Române, p. 35-56,

1997.

Ukkonen, E., *On-line construction of suffix trees*, în *Algorithmica*, 14(3), p. 249-260, 1995.