



(12)

BREVET DE INVENȚIE

(21) Nr. cerere: **a 2014 00346**

(22) Data de depozit: **07/05/2014**

(45) Data publicării mențiunii acordării brevetului: **30/09/2020** BOPI nr. **9/2020**

(41) Data publicării cererii:
29/01/2016 BOPI nr. **1/2016**

(73) Titular:

- **PETRICĂ LUCIAN, INTRAREA VÂSLEI**
NR. 1, BL. PM63, SC. 2, ET. 5, AP. 77,
SECTOR 3, BUCUREȘTI, B, RO;
- **CUCU HORIA, ALEEA POLITEHNICII**
NR. 6, BL. 3, SC. 4, ET. 2, AP. 42,
SECTOR 6, BUCUREȘTI, B, RO;
- **BUZO ANDI,**
BD. GENERAL VASILE MILEA NR. 8,
BL. B2, SC. 3, ET. 4, AP.57, SECTOR 6,
BUCUREȘTI, B, RO

(72) Inventatori:

- **PETRICĂ LUCIAN, INTRAREA VÂSLEI**
NR. 1, BL. PM63, SC. 2, ET. 5, AP. 77,
SECTOR 3, BUCUREȘTI, B, RO;

- **CUCU HORIA, ALEEA POLITEHNICII**
NR. 6, BL. 3, SC. 4, ET. 2, AP. 42,
SECTOR 6, BUCUREȘTI, B, RO;
- **BUZO ANDI,**
BD. GENERAL VASILE MILEA NR. 8, BL.
B2, SC. 3, ET. 4, AP.57, SECTOR 6,
BUCUREȘTI, B, RO

(74) Mandatar:

CABINET D.NICOLAESCU, STR.TURDA,
NR.102, BL.30A, ET.7, AP.28, BUCUREȘTI

(56) Documente din stadiul tehnicii:

**"ASR domain adaptation methods for
low-resourced languages:Application to
Romanian language" in Proceedings of the
20th IEEE European Signal Processing
Conference (EUSIPCO), H. Cucu, L.
Besacier, A. Buzo, pp. 1648-1652:
US8543382B2: US7966173B2**

(54) **METODĂ PENTRU RESTAURAREA AUTOMATĂ
A SEMNELOR DIACRITICE FOLOSIND TEXTE
ACHIZIȚIONATE ELECTRONIC, UTILIZATĂ ÎN PROCESAREA
LIMBAJULUI NATURAL**



RO 130875 B1

1 Invenția se referă la o metodă de restaurare automată a semnelor diacritice într-un
2 corpus de text format din fișiere multiple, pus la dispoziție de utilizator, destinată a fi utilizată în
3 domeniul sistemelor de procesare a limbajului natural.

4 Invenția aparține domeniului sistemelor de procesare a limbajului natural, sub formă
5 vorbită sau scrisă, care includ sistemele de corectare automată a resurselor text, recunoaștere
6 automată a vorbirii și sinteză automată a vorbirii.

7 Semnele diacritice sunt utilizate în multe limbi (franceză, spaniolă, română etc.) pentru
8 extinderea alfabetului dincolo de caracterele latine. Limba română utilizează trei semne
9 diacritice, care sunt folosite pentru construirea a cinci caractere speciale: ă, â, î, ș, ț. Semnele
10 diacritice modifică pronunția celor patru caractere de bază: a, i, s, t. Literatura științifică a
11 domeniului raportează că în textele din limba română, până la 40% din cuvinte folosesc semne
12 diacritice.

13 Din diverse motive, care includ lipsa unor tastaturi dedicate pentru multe din limbile care
14 utilizează diacritice, din majoritatea documentelor text în format electronic lipsesc semnele
15 diacritice. De cele mai multe ori, caracterele speciale ce utilizează semne diacritice sunt
16 înlocuite de caracterele de bază corespunzătoare. Acest fapt duce la scăderea inteligibilității
17 textului scris în format electronic, dar și la apariția unor cuvinte greșite din punct de vedere
18 gramatical.

19 Corectarea acestor texte poate fi automatizată, folosind un sistem informatic de
20 restaurare a semnelor diacritice. Pentru cuvintele care au o singură formă diacritică, restaurarea
21 este trivială și se poate face pe baza unui dicționar. Un exemplu ar fi restaurarea diacriticelor
22 cuvântului câine, a cărui singură formă corectă este câine. În cazul în care există mai multe
23 forme cu diacritice ale unui cuvânt, restaurarea necesită folosirea contextului pentru alegerea
24 variantei corecte de restaurare. De exemplu, în limba română formele articulate și nearticulate
25 ale multor cuvinte diferă doar prin diacritice: mașină și mașina sunt forme corecte, cu semne
26 diacritice, ale cuvântului de bază mașina, dar alegerea uneia din cele două variante la
27 restaurare necesită cunoașterea măcar a cuvântului anterior celui restaurat, în propoziție.

28 Un astfel de sistem de restaurare utilizează un model statistic de limbă ce specifică
29 probabilitățile de apariție ale fiecărei forme diacritice, dat fiind cuvintele anterioare cuvântului
30 de bază în textul restaurat. Un model statistic de limbă este construit (antrenat) folosind cantități
31 mari de text de înaltă calitate.

32 În acest context, un text de înaltă calitate este acel text care conține un număr mic de
33 erori din punctul de vedere al utilizării semnelor diacritice.

34 Se cunosc sistemele de restaurare a diacriticelor folosind abordări statistice.

35 În [D. Burileanu, C. Ungurean, "Metodă de inserare automată a semnelor diacritice pentru
36 sinteza vorbirii pornind de la text în limba română în aplicații ce au ca suport rețele de date",
37 cerere de brevet de invenție RO 127582 A2] se descrie un sistem de restaurare a semnelor
38 diacritice care necesită: (i) un dicționar al celor mai frecvent folosite cuvinte din limba română
39 și a formelor flexionate ale acestora, și (ii) un corpus de text care conține forme corecte scrise
40 cu diacritice.

41 În literatura de specialitate, au fost propuse diverse metode pentru restaurarea semnelor
42 diacritice, dintre care cele mai notabile sunt [C. Ungurean and D. Burileanu, "An advanced NLP
43 framework for high quality Text-to-Speech synthesis", 6th IEEE Conference on Speech
44 Technology and Human-Computer Dialogue (SpeD), pp. 1-6, 2011] care prezintă concepte simi-
45 lare cu cererea de brevet menționată anterior și raportează o performanță de 1,4% eroare la
46 nivel de cuvânt și 0,4% eroare la nivel de caracter. O altă lucrare științifică relevantă este [H.
47 Cucu, L. Besacier, C. Burileanu, and A. Buzo, 'ASR domain adaptation methods for low-
48 resourced languages: Application to Romanian language", in Proceedings of the 20th IEEE

RO 130875 B1

European Signal Processing Conference (EUSIPCO), pp. 1648-1652, 2012] care folosește pentru restaurarea semnelor diacritice un model statistic de limbă de tip 3-grame, antrenat folosind un corpus de text corectat manual. Autorii raportează o eroare de 1,99% la nivel de cuvânt și 0,48% la nivel de caracter. 1 3

Prezenta invenție se referă la o metodă de restaurare automată a semnelor diacritice într-un corpus de text format din fișiere multiple, pus la dispoziție de utilizator, metoda cuprinzând etapele de împărțire automată a corpus-ului în secțiuni de calitate înaltă, respectiv scăzută, folosind un prag fix de frecvență de apariție a semnelor diacritice ca principal criteriu de filtrare și, utilizare a secțiunilor de calitate înaltă pentru antrenarea unui model probabilistic de limbă care este folosit pentru restaurarea semnelor diacritice în secțiunile de calitate scăzută; metoda mai poate cuprinde o etapă de determinare automată a pragului optim de frecvență de apariție a semnelor diacritice, prin repetarea iterativă a primelor două etape enunțate mai sus și ajustarea pragului de filtrare la fiecare iterație, prin variere între zero și o valoare specificată de utilizator, în vederea minimizării erorii de restaurare a semnelor diacritice, calculată folosind un text de evaluare pus la dispoziție de utilizator, pragul optim fiind considerat acela care corespunde erorii minime de restaurare și putând fi identificat speculativ. 5 7 9 11 13 15

În continuare, se prezintă detaliat, în legătură și cu fig. de la 1 la 4, principiile și realizarea invenției și un exemplu de aplicare a sa constând în utilizarea metodei pentru corectarea unui corpus de știri în limba română, achiziționate de pe un portal online, și generarea unui model statistic de limbă pentru restaurarea semnelor diacritice în alte texte din limba română. Fig. 1... 4 reprezintă: 17 19 21

- fig. 1, metoda de creștere a calității unui corpus de text; 23
- fig. 2, algoritmul de filtrare folosind un prag fix de separare; 23
- fig. 3, algoritmul de identificare a pragului optim de separare; 25
- fig. 4, algoritmul de identificare speculativă a pragului optim de separare. 25

Invenția se referă la o metodă pentru curățarea unui corpus de text (**CT** 101, achiziționat din surse predispușe la erori legate de semnele diacritice (de exemplu, de pe internet), folosind un sistem de restaurare a semnelor diacritice (**SRD** 102, care se bazează pe un model statistic de limbă (**ML**) 103, generat prin procesarea **CT** 101. Invenția are la bază faptul că generarea modelului statistic de limbă 103 se face pornind de la **CT** 101, fără a avea garanția corectitudinii textului din **CT** 101 în ceea ce privește utilizarea semnelor diacritice. Metoda are ca rezultat un corpus augmentat de text (**CTA**) 104, de calitate ridicată, ce poate fi folosit pentru aplicații legate de recunoașterea vorbirii, sinteza vorbirii și alte procesări de limbaj natural. 27 29 31 33

Pentru a asigura calitatea modelului statistic de limbă 103, textul **CT** 101 este procesat de un filtru de separare 105, folosind anumite criterii de calitate, concretizate în valoarea unui prag de separare 106. În metoda propusă, acest prag de separare poate fi calculat automat prin metoda de selecție a pragului de separare 107, folosind un corpus de evaluare (**CE**) 108, sau furnizat direct de către utilizator. Filtrarea împarte textul 101 în două categorii: 35 37

- text ce conține cu probabilitate ridicată semne diacritice utilizate corect, care formează corpusul de calitate înaltă (**CI**) 109; 39

- text care, cu probabilitate ridicată, conține greșeli din punctul de vedere al utilizării semnelor diacritice, care formează corpusul de calitate redusă (**CR**) 110 41

În urma separării corpusului de text **CT** 101 în corpusurile de calitate înaltă **CI** 109, respectiv redusă **CR** 110, **CI** 109 este folosit pentru generarea modelului probabilistic de limbă 103. Acest model este generat folosind un bloc de analiză statistică 111. Apoi, modelul de limbă 103 este utilizat la restaurarea semnelor diacritice în **CR** 110, folosind **SRD** 102, bazat pe metodologia standard enunțată în literatura de specialitate și stadiul tehnicii. Textul obținut din restaurarea semnelor diacritice asupra **CR** 110 este adăugat la **CI** 109, rezultând un corpus final, **CTA** 104, de dimensiune egală cu **CT** 101, dar de calitate superioară. 43 45 47 49

RO 130875 B1

1 Un element esențial al invenției este algoritmul de filtrare folosit de filtrul de separare
2 **105**, care separă corpusul în două secțiuni de calitate înaltă **109**, respectiv scăzută **110**. Inven-
3 ția propune o metodă de separare bazată pe observații statistice asupra limbilor ce utilizează
4 semne diacritice. În aceste limbi, frecvența de apariție a semnelor diacritice în text este relativ
5 constantă, indiferent de natura subiectului descris. Astfel, este de așteptat ca dintr-o analiză
6 asupra unui text achiziționat online, de dimensiuni rezonabile (minim sute de cuvinte), să reiasă
7 o frecvență de utilizare a semnelor diacritice apropiată de media generală pentru limba textului
8 respectiv.

9 În practică, multe din textele achiziționate online nu folosesc semne diacritice sau le
10 folosesc doar parțial, având o frecvență a diacriticelor mult sub media limbii. Filtrul de separare
11 **105** folosește această proprietate pentru identificarea fișierelor text de calitate redusă. Fig. 2
12 prezintă filtrul **105** în detaliu. Atât timp cât mai există fișiere în **CT 101**, blocul **201** de extragere
13 de fișiere extrage un fișier text **202** din **CT 101** și livrează acest fișier către blocul **203** de
14 calculare a frecvenței semnelor diacritice. Dacă frecvența calculată este mai mică decât pragul
15 de separare **106**, fișierul de text se consideră a fi de calitate redusă și este adăugat la corpus-ul
16 **CR 110** de către blocul de adăugare **204**. Dacă frecvența depășește acest prag, fișierul se
17 consideră a fi de calitate ridicată și este adăugat la corpus-ul **CI 109** de către blocul de
18 adăugare **205**. Procesul se încheie atunci când nu mai există fișiere în **CT 101**.

19 Alegerea pragului de separare **106** determină calitatea corpus-ului final **104**. Dacă pragul
20 de separare este ales cu o valoare prea mică, fișiere de calitate redusă vor fi considerate a fi
21 de înaltă calitate, scăzând calitatea totală a **CI 109** și, prin urmare, a modelului statistic de limbă
22 **103**. Alternativ, dacă pragul de separare este ales cu o valoare prea mare, numărul fișierelor
23 considerate a fi de calitate înaltă va fi mic, reducând dimensiunea **CI 109**, ceea ce duce la
24 scăderea calității modelului statistic de limbă **103**.

25 În prezenta invenție, metoda prezentată anterior este completată și cu o etapă de
26 selecție automată a pragului optim de separare. Algoritmul de selecție automată este prezentat
27 în fig. 3. Algoritmul menține și ajustează iterativ un prag temporar **301**, care este inițial setat la
28 o valoare minimă, specificată de utilizator, în blocul **302**. Pragul temporar este folosit pentru
29 filtrarea **CT 101** și generarea modelului de limbă **111**. Sistemului din fig. 1 i se adaugă modulul
30 **301** de evaluare a modelului statistic de limbă folosit pentru restaurarea semnelor diacritice.
31 Acest modul folosește corpus-ul de evaluare **108**, de dimensiuni relativ mici și calitate bună
32 pentru a determina eroarea de restaurare. Pragul de separare temporar **301** este înregistrat ca
33 optim **304** dacă eroarea de evaluare este mai mică decât a tuturor iterațiilor anterioare. Pragul
34 temporar este incrementat în blocul **305** cu un increment specificat de utilizator. Pentru fiecare
35 valoare a pragului de separare se execută pașii de filtrare, generare a modelului statistic de
36 limbă folosind corpus-ul **CI** și evaluare a erorii de restaurare, folosind acest model de limbă. La
37 final, pragul optim identificat **304** este folosit ca prag de separare **106**.

38 Un potențial dezavantaj al metodei prezentate anterior este faptul că filtrarea se execută
39 de un număr mare de ori, pentru fiecare variantă posibilă a pragului de separare, între valorile
40 sale minimă, respectiv maximă. În funcție de dimensiunea corpus-ului **101**, timpul de execuție
41 a acestui proces pe un calculator personal poate deveni foarte mare.

42 Pentru reducerea acestui timp de execuție, metoda prezentată anterior este completată
43 și cu un mecanism de întrerupere speculativă a căutării pragului optim de separare. Decizia
44 întreruperii căutării se bazează pe observația că eroarea de evaluare va scădea, odată cu
45 creșterea valorii pragului de separare, către o valoare minimă, după care va crește pe măsură
46 ce pragul se incrementează și valoarea sa se îndepărtează de pragul optim. Astfel, se poate
47 considera că, odată ce eroarea de evaluare a crescut cu o anumită valoare față de eroarea
48 minimă înregistrată anterior, căutarea poate fi oprită. Fig. 4 prezintă algoritmul cu întrerupere
49 speculativă.

RO 130875 B1

Metoda care face obiectul prezentei invenții prezintă avantajul că permite restaurarea semnelor diacritice în textul de intrare, și generarea unui model de limbă ce poate fi folosit pentru restaurarea diacriticelor în alte texte, fără a necesita resurse de text adiționale, de înaltă calitate (ridicată). Pentru multe din limbile ce utilizează semne diacritice, nu există astfel de resurse de text de înaltă calitate, singurele resurse disponibile fiind cele ce pot fi achiziționate online, a căror calitate nu poate fi garantată.

Metoda conform invenției mai are avantajul că permite dezvoltarea de aplicații de procesare a limbajului natural sub formă textuală și vorbită, pentru limbi cu resurse subdezvoltate, unde majoritatea resurselor text sunt achiziționate online. Metoda nu este specifică limbii române, ci poate fi aplicată oricărei limbi care utilizează semne diacritice.

În cele ce urmează, vom exemplifica utilizarea metodei propuse, împreună cu optimizările destinate selecției pragului de filtrare, pentru curățarea greșelilor de folosire a diacriticelor în limba română, comparând rezultatele cu starea tehnicii. Metoda este aplicată asupra unui corpus de text care constă în articole de știri, achiziționate online de pe www.libertatea.ro. Parametrii folosiți pentru evaluarea metodei sunt următorii:

- TC conține 325.000 de fișiere și peste 70 de milioane de cuvinte în total. Din totalul fișierelor, aproximativ 153.000 nu conțin deloc diacritice, iar 6000 de fișiere conțin diacritice în mod sporadic.

- Pentru evaluarea filtrării cu prag fix furnizat de utilizator, s-a folosit un prag de separare cu valoarea de 20% (frecvența semnelor diacritice în text)

- Pentru evaluarea metodelor de selecție automată a pragului de separare, s-a folosit un corpus de evaluare de 4 milioane de cuvinte, verificate și corectate manual.

- Pentru metodele de selecție automată a pragului de separare, s-a folosit o gamă de valori a pragului de separare între 0 și 25%, în increment de 1%

- Pentru metoda speculativă de oprire a căutării pragului optim, s-a decis oprirea căutării atunci când eroarea crește cu 5% față de minim.

Comparația cu stadiul tehnicii s-a făcut folosind ca referință sistemul de restaurare a diacriticelor din [H. Cucu, L. Besacier, C. Buriianu, and A. Buzo, "ASR domain adaptation methods for low-resourced languages: Application to Romanian language", in Proceedings of the 20th IEEE European Signal Processing Conference (EUSIPCO), pp. 1648-1652, 2012.].

Tabelul 1 prezintă rezultatele evaluării metodelor propuse și comparația cu stadiul tehnicii.

Comparație cu stadiul tehnicii

Tabelul 1

Metoda	Cantitate necesară de text de calitate (10^6 cuvinte)	Eroare de restaurare (%)	Prag de separare identificat (%)	Timp de găsire a pragului de separare (ore)
Stadiul tehnicii	15	1,13	-	-
Prag fix	0	1,11	20	0
Căutare exhaustivă	4	0,92	13	16
Căutare speculativă	4	0,92	13	10

Se observă că metoda propusă duce la restaurarea semnelor diacritice cu eroare mai mică decât stadiul tehnicii, în toate cele 3 cazuri evaluate. Metoda cu prag fix de evaluare duce la o eroare de restaurare de 1,11%, având avantajul timpului de rulare minim, datorită faptului că nu necesită găsirea pragului de separare. De asemenea, această metodă nu necesită text de calitate.

RO 130875 B1

- 1 Ambele metode de căutare automată a pragului optim de separare duc la obținerea unor rezultate mai bune din punctul de vedere al erorii de restaurare. În cazul ambelor metode, se
- 3 obține o scădere de 17% a erorii, identificându-se pragul de separare optim la 13%. Dintre cele două metode de căutare automată a pragului optim de separare, căutarea speculativă are timpul
- 5 de rulare cel mai mic, 10 ore comparativ cu 16 ore în cazul căutării exhaustive.

RO 130875 B1

Revendicări

1. Metodă pentru restaurarea automată a semnelor diacritice folosind texte achiziționate electronic, utilizată în procesarea limbajului natural, implementată ca program software executabil pe un calculator sau sub formă de circuit integrat, **caracterizată prin aceea că** această metodă cuprinde etapele de:
- achiziționarea online a unui corpus de text **(101)** format din fișiere multiple, din surse disponibile, ce pot fi afectate de erori legate de semnele diacritice,
 - împărțirea automată a corpusului de text **(101)** într-o secțiune de calitate înaltă **(109)** și o secțiune de calitate redusă **(110)**, prin procesarea fiecărui fișier în ceea ce privește diacriticele, prin intermediul unui filtru de separare **(105)** folosind drept criteriu de filtrare un prag de separare **(106)** corespunzător frecvenței de apariție a semnelor diacritice, prag ce poate fi calculat automat folosind un corpus de evaluare **(108)** sau care exprimă în procente frecvența de apariție a semnelor diacritice în limba respectivă,
 - utilizarea secțiunii de calitate înaltă **(109)** pentru antrenarea unui model probabilistic de limbă **(103)**, generat prin intermediul unui bloc de analiză statistică **(111)**, și
 - utilizarea modelului probabilistic de limbă **(103)**, astfel antrenat, pentru restaurarea semnelor diacritice a secțiunii de calitate redusă **(110)**, obținându-se în final un corpus de text augmentat **(104)** de calitate ridicată, ce poate fi folosit pentru aplicații legate de procesări de limbaj natural, în special recunoașterea și sinteza vorbirii.
2. Metodă conform revendicării 1 **caracterizată prin aceea că**, pentru determinarea automată a unei valori optime a pragului de separare **(106)**, etapele metodei de la revendicarea 1 sunt repetate iterativ, pornind de la un prag temporar **(301)** setat la o valoare minimă și ajustarea incrementală a acesteia la fiecare iterație, până la o valoare maximă, valoarea optimă a pragului de separare **(106)** corespunzând erorii minime de restaurare a semnelor diacritice, calculată folosind corpusul de evaluare **(108)**.
3. Metodă conform revendicării 2 **caracterizată prin aceea că** pentru reducerea timpului de determinare automată a valorii optime a pragului de separare **(106)**, procesul de iterare este oprit atunci când eroarea de restaurare a iterației curente depășește cu un procent prestabilit valoarea minimă a erorii de restaurare înregistrată pentru iterațiile anterioare, pragul optim fiind considerat cel care corespunde respectivei erori minime de restaurare.
4. Sistem informatic pentru restaurarea automată a semnelor diacritice folosind texte achiziționate electronic, utilizat în procesarea limbajului natural, **caracterizat prin aceea că** acesta cuprinde un circuit integrat configurat pentru realizarea etapelor metodei din revendicarea 1.

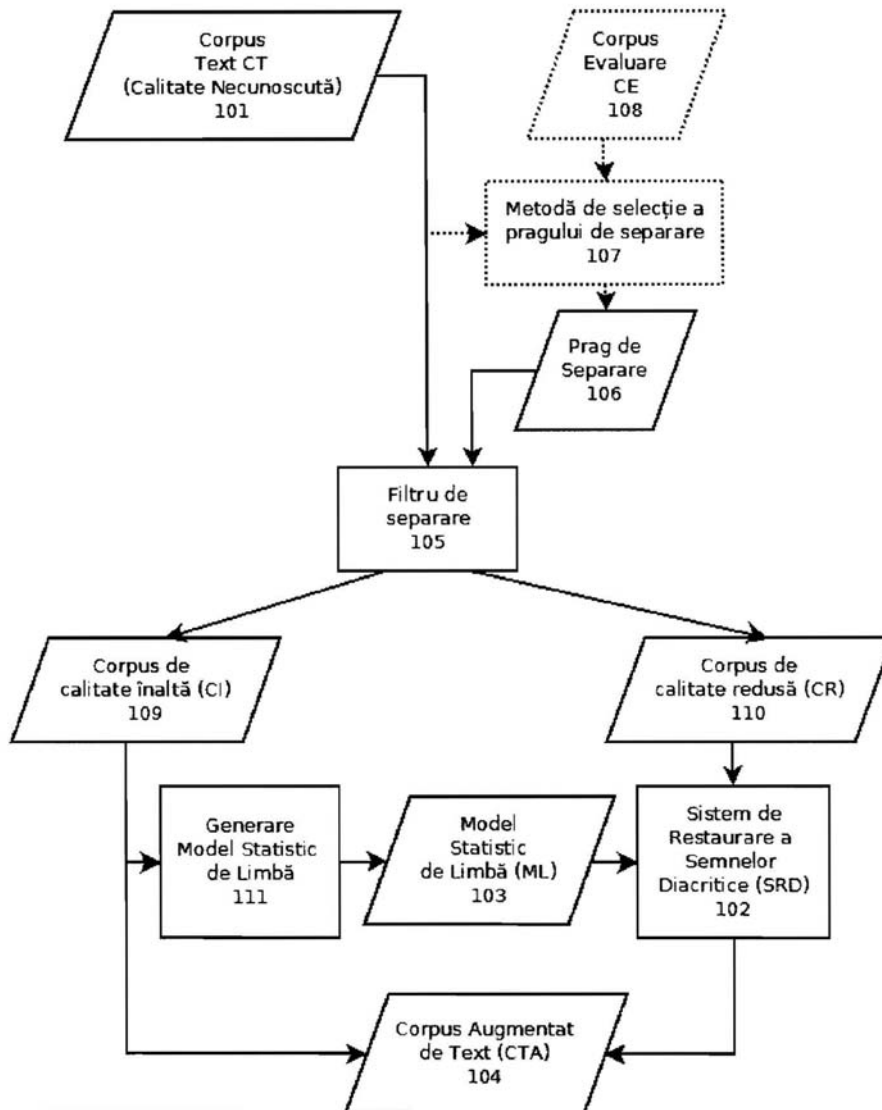


Fig. 1

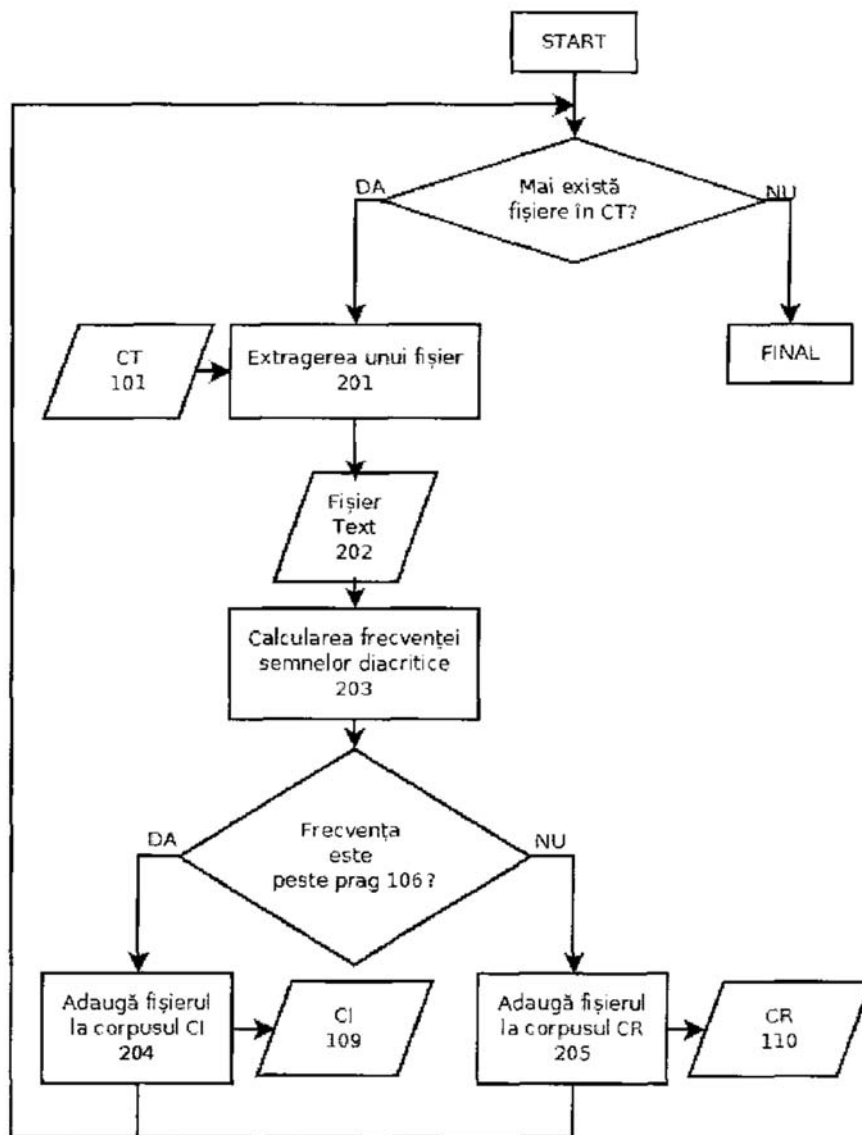


Fig. 2

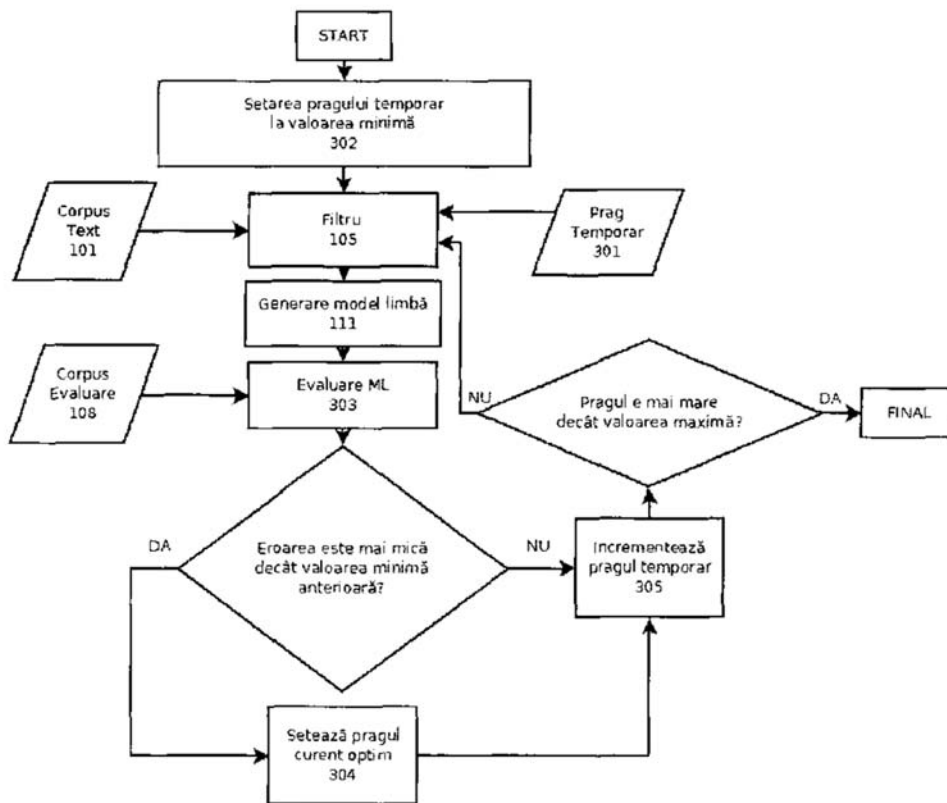


Fig. 3

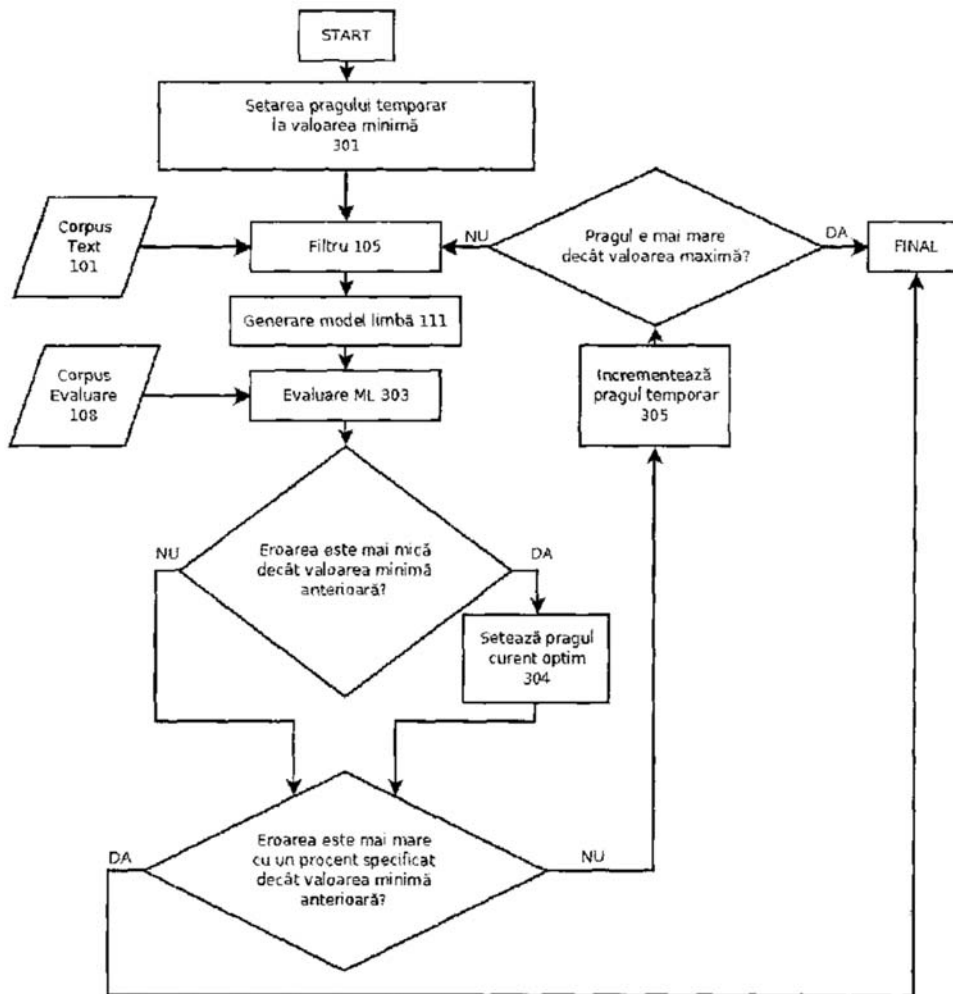


Fig. 4

