



(12) CERERE DE BREVET DE INVENȚIE

(21) Nr. cerere: a 2010 01247

(22) Data de depozit: 29.11.2010

(41) Data publicării cererii:  
29.06.2012 BOPI nr. 6/2012

(71) Solicitant:  
• UNIVERSITATEA "POLITEHNICA" DIN  
BUCUREȘTI, SPLAIUL INDEPENDENȚEI  
NR.313, SECTOR 6, BUCUREȘTI, B, RO

(72) Inventatori:  
• BURILEANU DRAGOȘ, STR. JOHANNES  
KEPLER NR. 2, BL. 2, SC. 2, AP. 61,  
SECTOR 2, BUCUREȘTI, B, RO;  
• UNGUREAN CĂTĂLIN, STR. MOLDOVIȚA  
NR. 2, BL. M2D9/4, AP. 44, SECTOR 4,  
BUCUREȘTI, B, RO

(54) **METODĂ DE INSERARE AUTOMATĂ A SEMNELOR  
DIACRITICE PENTRU SINTEZA VORBIRII PORNIND DE LA  
TEXT ÎN LIMBA ROMÂNĂ ÎN APLICAȚII CE AU CA SUPORT  
REȚELE DE DATE**

(57) Rezumat:

Invenția se referă la o metodă de inserare automată a semnelor diacritice în texte în limba română scrise fără diacritice, destinată utilizării în sinteza vorbirii pornind de la text, pentru aplicații ce trebuie să prelucreze documente în format electronic și au ca suport rețele de date. Metoda conform invenției are la bază o strategie statistică axată pe  $n$ -grame, și constă în esență într-un proces de filtrare secvențială pe mai multe niveluri, bazat pe probabilitatea de apariție a cuvintelor scrise cu diacritice și pe contextul în care acestea apar. Fiecare

criteriu de filtrare aduce o îmbunătățire suplimentară, fără a mări caracterul de ambiguitate. Astfel, după ce tratează inițial cuvintele fără diacritice sau pentru care nu există ambiguități, metoda utilizează incremental trei etaje de filtrare: la nivel de unigrame, la nivel de bigrame și la nivel de trigrame, informația contextuală fiind adăugată gradual, după fiecare etaj, simultan cu eliminarea variantelor incorecte.

Revendicări: 1



39

OFICIUL DE STAT PENTRU INVENȚII ȘI MĂRCI
Cerere de brevet de invenție
Nr. .... a 2010 01247
Data depozit ... 25.11.2010

## **METODĂ DE INSERARE AUTOMATĂ A SEMNELOR DIACRITICE PENTRU SINTEZA VORBIRII PORNIND DE LA TEXT ÎN LIMBA ROMÂNĂ ÎN APLICAȚII CE AU CA SUPORT REȚELE DE DATE**

Invenția se referă la o metodă de inserare automată a semnelor diacritice în texte în limba română scrise fără diacritice, destinată utilizării în sinteza vorbirii pornind de la text pentru aplicații ce trebuie să prelucreze documente în format electronic și au ca suport rețele de date, cum ar fi spre exemplu ascultarea mesajelor de tip e-mail/SMS folosind rețelele telefonice publice sau citirea unor informații de pe pagini web.

Sunt cunoscute un număr de studii pe plan mondial în problema refacerii caracterelor scrise cu semne diacritice, realizate în special pentru limbile franceză, spaniolă, ungară și doar parțial pentru limba română. Au fost astfel propuse câteva metode, dintre care unele se bazează pe algoritmi pur-statistici, fără alte cunoștințe lingvistice, iar altele folosesc modele de limbaj și diferite niveluri de procesare lingvistică.

– De exemplu, din categoria metodelor statistice se poate menționa un algoritm de restaurare a diacriticelor care folosește tehnici de învățare automată și care operează la nivel de caracter, propus în [R. Mihalcea, V. Nastase, “Letter Level Learning for Language Independent Diacritics Restoration”, in *Proc. CoNLL 2002*, Taipei, Taiwan, 2002, pp. 105-111]. Algoritmul a fost implementat pentru limbile cehă, ungară și poloneză iar pentru limba română performanța raportată (la nivel de caracter) a fost de aproximativ 98,3%.

Dezavantajele acestei metode sunt legate de următoarele aspecte: în primul rând, considerăm că abordarea bazată doar pe inspectarea literelor ce înconjoară caracterul analizat, nu este adecvată limbii române, unde multe ambiguități pot fi rezolvate doar la nivel de cuvânt (un exemplu simplu este cel legat de substantivele feminine articulate/nearticulate); în al doilea rând, rezultatele raportate sunt puțin relevante, deoarece baza de date folosită pentru testarea metodei este foarte mică (aproximativ 50.000 de cuvinte).

– Din cea de-a doua categorie de metode menționate, în [D. Tufiș, A. Chițu, “Automatic Diacritics Insertion in Romanian Texts”, in *Proc. of COMPLEX'99*, Pecs, Hungary, June 16-19, 1999, pp. 185-194] este propus un algoritm mai complex de inserare a diacriticelor pentru textele scrise în limba română, care funcționează pe baza informației de parte de vorbire și conduce la o acuratețe raportată de 97,4%.

Dezavantajul major al acestei metode este complexitatea mare a algoritmului, bazat pe resurse externe (dicționare) de dimensiuni mari și prelucrări lingvistice adiționale, fapt ce nu este compatibil cu utilizarea sa eventuală într-un sistem de sinteză automată a vorbirii.

Problema tehnică pe care și-a propus să o rezolve invenția poate fi formulată pornind de la constatarea că *sinteza automată a vorbirii* reprezintă o tehnologie importantă pentru aplicațiile care au ca suport rețelele de date. Astfel de aplicații, cum ar fi cele de citire a mesajelor de tip text (e-mail, SMS, pagini HTML etc.) trebuie să rezolve numeroase probleme, una dintre acestea fiind necesitatea refacerii semnelor diacritice în textul original. De obicei, majoritatea utilizatorilor nu folosesc diacritice (semne distinctive plasate deasupra, în interiorul sau dedesubtul caracterelor), chiar dacă terminalul (sau sistemul de operare) permite acest lucru; pe de altă parte, în cazul proceselor de lucru care utilizează o codare limitată sau neunitară a caracterelor (de exemplu cea pe 7 biți), diacriticele sunt înlocuite fie cu semne indescifrabile fie cu caracterele de bază (în cel mai bun caz). Această situație este valabilă atât pentru limba română, cât și pentru alte limbi vorbite care folosesc alfabetul latin (franceza, spaniola, germana, daneza, ungara, poloneza, suedeza etc.) și care trebuie să-și suplimenteze setul alfabetic de bază prin folosirea literelor scrise cu semne diacritice (sau *markeri* de accent) pentru a indica sunete (sau accente) diferite. Pe de altă parte, sinteza automată a unui text scris fără diacritice duce la o inteligibilitate scăzută; uneori, ambiguitățile sintactice și semantice fac ca o întreagă frază rostită să fie aproape de neînțeles. Repoziționarea automată a diacriticelor reprezintă însă o problemă extrem de dificilă, deoarece practic nu există reguli lingvistice care să permită rezolvarea acestei sarcini.

Metoda care face obiectul prezentei invenții rezolvă problema tehnică menționată anterior prin faptul că la intrarea în sistemul de sinteză a vorbirii, se verifică mai întâi dacă textul de intrare (scris la tastatura calculatorului sau provenind dintr-un fișier stocat anterior) conține caractere scrise cu diacritice, iar dacă acest lucru nu se întâmplă, metoda se aplică prin intermediul unui algoritm eficient de inserare automată a semnelor diacritice.

Metoda este originală, în contextul preocupărilor similare pe plan național și internațional. Principalul aspect original ce poate fi menționat constă într-un proces de *filtrare*

*secvențială*, pe mai multe niveluri, bazat pe probabilitatea de apariție a cuvintelor scrise cu diacritice și pe contextul în care acestea apar. Fiecare criteriu de filtrare aduce o îmbunătățire suplimentară, fără însă a mări caracterul de ambiguitate (deci formele posibile teoretic de scriere a unui cuvânt cu sau fără diacritice). Astfel, metoda utilizează practic trei etaje de filtrare: la nivel de *unigrame* (pe cuvânt), la nivel de *bigrame* (de asemenea pe cuvânt) și la nivel de *trigrame* pentru terminații de cuvânt. Un al doilea aspect important este reprezentat de faptul că se realizează în fapt o filtrare minimală la nivelul bazei de date de antrenare, îndepărtând, într-o primă fază, numai formele asupra cărora se decide că sunt “sigur greșite”.

Prezenta invenție înlătură dezavantajele metodelor cunoscute, în primul rând prin faptul că se aplică ca un modul de pre-procesare înaintea unui sistem de sinteză a vorbirii pornind de la text, permițând astfel în final sinteza corectă a unui text scris fără diacritice. În acest sens, trebuie subliniat faptul că metodele abordate pe plan mondial (pentru diferite limbi unde sunt utilizate în scriere semne diacritice), ca și cele pentru limba română prezentate anterior, sunt destinate în marea lor majoritate doar aplicațiilor și studiilor strict lingvistice. În al doilea rând, metoda care face obiectul invenției este extrem de performantă, atât din punct de vedere al acurateții obținute, cât și a timpului de calcul și a resurselor de memorie necesare (ambele extrem de importante în situația în care întreg sistemul de sinteză a vorbirii trebuie să funcționeze în timp real). Astfel, metoda utilizează o abordare statistică axată pe *n*-grame, se bazează pe cunoștințe lingvistice limitate și necesită un corpus de antrenare de dimensiuni medii.

Se prezintă în continuare în detaliu obiectul invenției, un exemplu de aplicare a sa constând din utilizarea metodei de către autorii invenției ca modul de pre-procesare în cadrul unui sistem complet de sinteză a vorbirii pornind de la text în limba română (dezvoltat în cadrul colectivului din care fac parte autorii), sistem folosit la rândul său într-o aplicație de tip “cititor de e-mail/SMS” pentru telefoane mobile.

Limba română folosește trei semne diacritice, și anume *accentul circumflex*, *sedila* și *accentul grav*, rezultând cinci litere scrise cu diacritice: **ă**, **â** / **î** (folosite pentru același sunet dar în circumstanțe diferite), **ș** și **ț**. Unele dintre diacritice pot face diferența între două forme distincte ale aceluiași cuvânt, de exemplu **casa** – **casă** (substantiv articulat – nearticulat), **înalta** – **înaltă** (adjectiv articulat – nearticulat), sau între două cuvinte complet diferite ca înțeles, de exemplu **fată** – **față**. Trebuie menționat faptul că în textele scrise în limba română,

cuvintele scrise cu diacritice au un aport semnificativ, mai mare decât în alte limbi. Astfel, se estimează că frecvența lor de apariție este între 25% și 40%, în timp ce în limba franceză aproximativ 15% dintre cuvinte pot avea diacritice. În alte limbi, ca olandeza, ceha și slovacă, aportul cuvintelor scrise cu diacritice este semnificativ mai mare decât în limba română.

După o analiză făcută pe un volum mare de texte, extrase din diferite domenii, am decis că putem clasifica cuvintele limbii române, din punctul de vedere al interacțiunii lor cu semnele diacritice, în următoarele cinci categorii (le vom folosi ulterior cu denumirea de “tipul 1”, “tipul 2” etc.):

- 1) Cuvinte scrise fără diacritice; exemplu: **lemn**.
- 2) Cuvinte scrise întotdeauna cu diacritice: **câteva, științific**.
- 3) Cuvinte ambigue, cu două forme: **masa, masă**.
- 4) Cuvinte ambigue în care un anumit număr de diacritice sunt întotdeauna prezente: **cămașă, cămașa**.
- 5) Cuvinte ambigue în care oricare semn diacritic poate să apară sau nu: **pana, pană, până**.

Această împărțire pe clase distincte a cuvintelor scrise cu diacritice stă la baza realizării metodei de inserare automată a semnelor diacritice, care constă în două etape mari: *antrenare și testare*.

**Etapa de antrenare** presupune realizarea următorilor pași:

1. Construirea manuală a unui dicționar  $D1$  care conține cele mai frecvent folosite cuvinte din limba română și care conține de asemenea cât mai multe forme flexionate ale acestora.
2. Pe baza lui  $D1$  se construiește o structură dicționar  $D2$  (“*hash table*”) care pune în corespondență fiecare cuvânt din  $D1$ , cu diacriticele înlăturate, cu toate formele sale posibile de scriere cu diacritice (extrase din  $D1$ ).
3. Pe baza unui corpus de antrenare  $A$ , care conține forme corecte scrise cu diacritice, se construiesc la nivelul cuvintelor un set de trei dicționare, astfel: o structură  $U$  care conține *unigramele* extrase din dicționarul de antrenare împreună cu frecvențele lor de apariție, o structură  $B$  care conține toate *bigramele* posibile extrase din corpusul  $A$ , și o structură  $T$ , care conține *trigramele* la nivelul sufixelor cuvintelor din  $A$ .

Trebuie menționat faptul că în cazul metodei propuse, termenul de *sufix* al cuvintelor, cu care realizăm dicționarul  $T$  nu se referă strict la conceptul morfologic cunoscut, ci se referă

de fapt la ultimele caractere ale fiecărui cuvânt. Acest număr de caractere finale a fost testat în intervalul 1 la 4 pe baza rezultatelor oferite în etapa de testare, în așa fel încât erorile de inserare a diacriticelor să fie minimizate. S-a observat că inițial are loc o scădere a erorilor de refacere a diacriticelor odată cu creșterea numărului de litere finale cu care se formează trigramele, proces urmat apoi de o creștere a numărului acestor erori simultan cu obținerea unui dicționar  $T$  de dimensiuni foarte mari.

**Etapă de testare** (cu alte cuvinte **de funcționare propriu-zisă**) este compusă în principal din trei procese de filtrare în cascadă, în care datele de la ieșirea fiecărui etaj superior constituie date de intrare în etajul imediat următor. Intrarea generală constă în secvențe de text în care semnele diacritice au fost înlăturate. Etapa de testare funcționează astfel:

1. Pe baza dicționarelor  $D1$  și  $D2$  sunt inserate formele care nu conțin diacritice și formele care nu introduc ambiguități, altfel spus, cuvintele de **tipul 1** respectiv de **tipul 2**. Etapa este una de tip determinist și presupune o înlocuire biunivocă a formei fără diacritice cu cea corectă (cu semnele diacritice inserate).
2. Pentru fiecare pereche de cuvinte consecutive din textul test, care conține cel puțin o formă ambiguă, se formează setul de variante posibile cu toate aceste cuvinte; dacă cel puțin o pereche din acest set este găsită în corpusul de bigrame  $B$ , sunt eliminate toate perechile de variante care nu apar în  $B$ , altfel sunt lăsate toate variantele nemodificate.
3. Din textul de intrare care conține variantele de ieșire de la pasul 2, se formează toate tripletele de câte trei cuvinte consecutive și din acestea se extrag sufixele; se obțin astfel perechi alternative de triplete de tip cuvinte-sufixe; dacă există cel puțin o tripletă de sufixe în dicționarul  $T$ , atunci renunță la toți membrii din setul de cuvinte corespunzătoare tripletelor de sufixe care nu sunt în  $T$  și furnizează spre ieșire doar variantele pentru care s-au găsit membri în  $T$ ; dacă nu s-a găsit niciuna dintre variante în  $T$ , lasă tot setul de cuvinte neschimbat.
4. Folosind setul de unigrame  $U$ , păstrează pentru fiecare cuvânt doar varianta cea mai probabilă, pe baza frecvențelor de apariție calculate din corpusul de antrenare, dintre variantele furnizate la pasul 3.

În acest punct al descrierii metodei ce face obiectul invenției, trebuie făcute câteva precizări referitoare la modul de obținere a dicționarelor de prelucrare și a corpusurilor de  $n$ -grame.

### a. Obținerea dicționarilor de prelucrare

– Se pornește de la un dicționar care conține cuvinte scrise cu diacritice, cu cât mai multe forme flexionate. Concret, dicționarul astfel creat are un număr de peste 330.000 de intrări.

– Se elimină semnele diacritice păstrându-se câte o singură formă din cuvintele astfel obținute și se compară cu dicționarul de start. Sunt separate în acest fel cuvintele cu probleme la refacerea diacriticelor de cele fără ambiguități. Dicționarul cuvintelor “cu probleme” va fi folosit și în etapa de reducere a dimensiunilor fișierelor cu *n*-grame.

După această primă etapă, se obțin două dicționare distincte, la care se mai adaugă un dicționar de nume proprii și unul de abrevieri, create *offline*, formându-se un așa-zis *bloc al dicționarilor de prelucrare*.

### b. Obținerea fișierelor reduse de *n*-grame

Se formează *offline* un corpus, prin juxtapunerea de texte scrise cu diacritice. Pe baza acestui corpus de start de aproximativ  $10^7$  intrări, se extrag *unigramele*, *bigramele* și *trigramele*.

Studiile realizate de autorii invenției au arătat că în limba română cele mai mari probleme pentru refacerea diacriticelor sunt cele create de ambiguitatea **a – ă** la finalul cuvintelor. Acest fenomen este produs cu o frecvență mare de apariție de unele construcții lingvistice scurte precum **ca – că, sa – să, sau – său** etc. și de formele articulate/nearticulate ale substantivelor feminine, de exemplu: **casa – casă, piața – piață** etc.

Din perspectiva ideii expuse anterior, am introdus un etaj de filtrare care folosește trigramele obținute din terminațiile cuvintelor. Câștigul adus de acest etaj este de aproximativ 0,8%, mult mai bun decât cel al trigramelor obținute din cuvintele întregi, tocmai prin abordarea unei probleme care apare cu o incidență foarte mare.

Pentru reducerea dimensiunilor fișierelor cu *n*-grame (implicit pentru scăderea resurselor necesare în prelucrare), se poate folosi dicționarul cuvintelor cu probleme din blocul dicționarilor de prelucrare. Cu alte cuvinte, se păstrează numai *n*-gramele cuvintelor care introduc ambiguități la refacerea diacriticelor, celelalte forme fiind înlăturate.

Pentru exemplul concret discutat de aplicare a invenției vom prezenta în continuare câteva aspecte importante referitoare la metoda expusă.

În primul rând trebuie specificat faptul că în etajele implementate cu bigrame și trigrame nu au fost folosite probabilitățile de apariție pentru aceste *n*-grame. Deși o

implementare realizată la nivel de cuvânt, folosind probabilitățile trigramelor obținute pe baza unui corpus de antrenare, ar fi fost binevenită, dimensiunile corpusului de antrenare ar trebui să fie uriașe, de ordinul a câteva sute de MB, pentru a se obține un corpus de trigrame eficient, care să conțină cât mai multe dintre trigramele posibile (dată fiind și dimensiunea mare a dicționarului de start  $D1$ ). În aceste condiții fenomenele nedorite produse de insuficiența datelor de antrenare ("*data sparseness*") ar fi fost clare, lucru pus în evidență și în cazul metodei propuse, în care la nivelul trigramelor s-a lucrat doar cu terminații de cuvinte pentru care singura condiție impusă a fost cea de existență a trigramelor.

Astfel s-a putut observa că, odată cu creșterea lungimii trigramelor, după ordinul 4, etajul trigramelor introduce erori mari din cauza fenomenului amintit. Acest lucru este ușor de intuit deoarece, odată cu creșterea lungimii trigramelor, lungimea sufixelor devine comparabilă cu dimensiunea cuvintelor, iar dimensiunea vocabularului de sufixe tinde să se apropie de dimensiunea lui  $D1$ . Etajul trigramelor este eficient atâta vreme cât se realizează un bun echilibru între dimensiunea vocabularului de sufixe și dimensiunea corpusului de antrenare. Chiar și în aceste condiții minimale impuse bigramelor și trigramelor – care presupun numai existența acestora, s-a putut observa încă efectul nedorit al insuficienței datelor de antrenare, tradus prin lipsa unor  $n$ -grame mai puțin întâlnite și furnizarea către etajul unigramelor a unor variante incorecte. Acest etaj final va decide varianta cea mai probabilă, urmărind probabilitățile de apariție ale fiecărei variante calculate pe baza corpusului de antrenare  $A$ . Două aspecte sunt de remarcat referitor la etajul final: în primul rând, trebuie spus că dacă în  $U$  nu există nici o variantă, cuvântul de intrare este lăsat nemodificat. În al doilea rând, deoarece acesta este un etaj final, o singură variantă – cea mai probabilă, va fi furnizată la ieșire.

De asemenea, trebuie menționat faptul că metoda de refacere a diacriticelor ce face obiectul invenției a fost dezvoltată într-o manieră incrementală, tratând inițial cuvintele fără diacritice sau pentru care nu există ambiguități, după care se adaugă gradual informație contextuală furnizată prin nivelul bigramelor de cuvinte și/sau cel al trigramelor – simultan cu eliminarea variantelor incorecte, pentru ca în final să se filtreze variantele rămase prin folosirea probabilităților unigramelor calculate pe baza frecvențelor de apariție obținute din corpusul de antrenare. Fiecare nivel de procesare aduce un plus în reducerea ratei de eroare la nivel de cuvânt. Deși s-au încercat diferite variante de succesiune ale etajelor de filtrare, am constatat că cea mai bună configurație este cea în care etajul unigramelor este ultimul. Acest lucru poate fi demonstrat și teoretic deoarece fiecare etaj de filtrare, cu excepția celui de unigrame ia în calcul numai apariția perechilor de cuvinte sau tripletelor de sufixe, fără a se



ține cont de frecvențele de apariție ale acestora. În acest fel, pentru un text de intrare dat, nivelurile de filtrare realizate cu bigrame respectiv trigrame pot întoarce una, două sau chiar toate variantele ambigue, printre care se află și varianta corectă. Fiecare set de filtrare reduce, teoretic, acest set de variante oferite către ieșire, dar acest lucru nu este garantat.

Ceea ce se urmărește cu primele două niveluri este ca, din punct de vedere statistic, numărul variantelor incorecte furnizate către ultimul etaj să fie sensibil diminuat. Prin urmare, etajul bigramelor și cel al trigramelor reduce numărul de variante incorecte furnizat către etajul final, însă numai etajul unigramelor va întoarce câte o singură variantă pentru fiecare set de variante posibile ale fiecărui cuvânt, adică pe cel mai frecvent, eliminându-le pe celelalte. Toate rezultatele experimentale pe care le vom raporta în continuare vor avea ca ultim nivel de prelucrare etajul unigramelor.

În cadrul metodei ce face obiectul invenției, antrenarea a fost făcută pe un corpus extras din texte literare scrise în limba română, care conține peste  $10^7$  cuvinte. Aceste texte au contribuit de asemenea la extinderea dicționarului *D1*, construit inițial manual, pe baza *Dicționarului Explicativ al Limbii Române*. Din păcate acest dicționar conține numai formele de bază ale cuvintelor, adică sub 70.000, iar pentru rularea algoritmului am avut nevoie de un dicționar care să conțină cât mai multe forme morfologice flexionate. La ora actuală dicționarul conține aproximativ 330.000 de forme și include aproximativ 6.000 de nume proprii. Dicționarele (vom folosi mai departe și denumirea de *corpusuri*) de unigrame, bigrame și trigrame au fost extrase din textele de antrenare și salvate ca fișiere text separate.

Corpusul de testare a fost corectat manual și conține propoziții extrase din articole, texte literare și teze de doctorat din mai multe domenii de activitate, disponibile în format electronic în limba română. Acest corpus are următoarele caracteristici: (i) numărul de cuvinte de intrare: **1.200.000**, (ii) procentul de cuvinte scrise fără diacritice: 57,41%; procentul de cuvinte scrise întotdeauna cu diacritice: 16,33%; procentul de cuvinte ambigue la refacerea diacriticelor: 26,26%.

Pe baza celor descrise anterior, au fost calculate următoarele mărimi de performanță standard, utilizate în testele statistice clasice:

1. *Precizia*: definită ca raportul dintre numărul total de diacritice inserate corect și numărul total de diacritice inserate.
2. *Recall*: definit ca raportul dintre numărul total de diacritice inserate corect și numărul de diacritice din baza de test corectată manual.
3. *F-measure*: definită ca media armonică dintre *precizie* și *recall*.

În Tabelul 1 sunt prezentate principalele rezultate obținute, detaliate pentru fiecare dintre cele patru clase de ambiguitate enunțate anterior. Rezultatele au fost calculate pentru trei versiuni de rulare ale metodei propuse, astfel:

- (i) Ca bază de comparație s-a luat metoda de restaurare a diacriticelor pe baza dicționarului *D2*, pentru cuvintele care nu prezintă ambiguități, fără a fi necesare date de antrenare.
- (ii) A doua metodă folosește bigramele extrase din corpusul de antrenare și probabilitățile unigramelor calculate pe baza aceluiași corpus.
- (iii) Metoda completă, care constă din toate cele trei etaje de filtrare așezate succesiv: etajul bigramelor, cel al trigramelor cu sufixe, și cel ce utilizează unigramele.

Clasa de ambiguitate	Precizie (%)	Recall (%)	F-measure (%)	F-measure medie (%)
Baza de comparație				
a / ă / â	78.03 / 92.69 / 99.50	99.65 / 35.27 / 61.86	87.53 / 51.10 / 76.29	77.11
i / î	98.63 / 99.67	99.96 / 89.54	99.30 / 94.34	98.71
s / ș	94.83 / 97.89	99.48 / 81.48	97.10 / 88.93	95.25
t / ț	95.30 / 91.25	98.82 / 71.63	97.03 / 80.25	94.57
Bigrame și unigrame				
a / ă / â	97.82 / 91.58 / 99.24	96.49 / 94.66 / 99.57	97.15 / 93.10 / 99.40	96.19
i / î	99.97 / 99.57	99.94 / 99.82	99.95 / 99.69	99.93
s / ș	99.93 / 99.49	99.85 / 99.76	99.89 / 99.63	99.83
t / ț	99.86 / 98.07	99.66 / 99.22	99.76 / 98.64	99.60
Bigrame, trigrame și unigrame				
a / ă / â	98.22 / 93.30 / 99.26	97.22 / 95.62 / 99.55	97.72 / 94.45 / 99.41	<b>96.93</b>
i / î	99.98 / 99.56	99.94 / 99.82	99.96 / 99.7	<b>99.93</b>
s / ș	99.93 / 99.51	99.85 / 99.77	99.89 / 99.64	<b>99.84</b>
t / ț	99.87 / 98.2	98.69 / 99.27	99.78 / 98.74	<b>99.63</b>

**Tabelul 1.** Măsurarea performanțelor metodei de inserare automată a diacriticelor

În Tabelul 1 sunt prezentate de asemenea valorile medii statistice ponderate ale *F-measure* pentru fiecare clasă de ambiguitate în care ponderile sunt calculate la nivel de caracter. Se observă că la nivel de caracter se obține o valoare medie ponderată pentru *F-measure* de **99,34%**.

Din analiza datelor obținute putem trage concluzia că etajul de filtrare cu bigrame aduce o îmbunătățire substanțială de performanță la nivel de caracter, pentru toate clasele de ambiguitate. În plus, etajul de filtrare cu trigrame cu sufixe aduce o îmbunătățire evidentă numai pentru clasa de ambiguitate **a / ă / â**, ceea ce este un lucru pozitiv dacă ținem cont de faptul că această clasă este sursa cea mai importantă de erori în procesul de refacere a diacriticelor. Pentru clasa **i / î** de ambiguitate rezoluția metodei este cea mai bună. Explicația vine din faptul că, în limba română, **î** are o poziționare quasi-regulată, putând să apară fie la început de cuvânt, fie în interiorul unor cuvinte compuse. Se poate merge chiar mai departe să se spună că pentru această clasă se poate oferi o metodă de refacere a diacriticelor de tip determinist.

După cum se observă, etajul trigramelor nu alterează rezultatele obținute bune obținute prin filtrarea cu bigrame pentru toate cele trei clase de ambiguitate.

Au fost efectuate și alte experimente folosind trigramele obținute la nivelul prefixelor, mai exact a caracterelor de început ale cuvintelor, dar această metodă nu a adus nici un plus evident de performanță. Explicația vine de acolo că, prin natura inflexiunilor din limba română, prefixele apar mai degrabă independent de context.

Ca o primă concluzie se poate observa că cel mai problematic caracter, care aduce cel mai slab nivel de performanță la refacerea diacriticelor, este reprezentat de clasa **a / ă / â**. În ciuda faptului că o mare parte dintre erori au fost corectate, rămân unele situații pe care algoritmul nu le tratează corect. Cele mai semnificative astfel de erori sunt prezentate în Tabelul 2. De aici se observă că cel mai mare număr de erori este adus de perechea ambiguă **ca – că** pe care o putem regăsi cu o frecvență mare de apariție în corpusul de test. Această situație de ambiguitate este oarecum explicabilă dacă se ține cont de faptul că există numeroase cazuri în care **ca** și **că** apar în aproximativ aceleași contexte.

Totodată trebuie arătat faptul că există și alte limitări pe care metoda nu le poate soluționa, printre care pot fi enumerate:

1. Unele nume proprii, care conțin diacritice, sau unele cuvinte rare, care nu apar în dicționarul *D1*. Soluția poate fi cea de creștere a dicționarului *D1* prin adăugarea de noi cuvinte. Trebuie menționat faptul că adăugarea de noi forme în *D1* nu

poate duce la creșterea dicționarelor de bigrame respectiv de trigrame. Aceste structuri au fost obținute din corpusul de antrenare  $A$  iar pentru mărirea lor este necesar să îl extindem pe acesta și cu alte contexte care să conțină noile forme introduse în  $D1$ .

2. Forme puternic dependente de context, la care nu există doar la nivelul frazei informații suficiente pentru a dezambiguiza varianta corectă. Ca exemplu se poate da situația ambiguă *Am văzut o față frumoasă la geam – Am văzut o față frumoasă la geam*. Este evident că pentru a clarifica o astfel de situație confuză sunt necesare informații suplimentare, la nivel de discurs.

Cuvânt	Bază de comparație	Bigrame și unigrame	Bigrame, trigrame și unigrame
ca / că	19,718	4,890	2,478
fată / față	2,199	912	495
sau / său	2,548	441	350
sa / să	30,064	395	341
lua / luă	547	330	272
ușa / ușă	1009	362	267
banca / bancă	507	248	243

**Tabelul 2.** Numărul de erori frecvente corectate gradual de metodă

## REVENDICĂRI

**Metodă de inserare automată a semnelor diacritice pentru sinteza vorbirii pornind de la text în limba română în aplicații ce au ca suport rețele de date,** caracterizată prin aceea că la intrarea unui sistem de sinteză a vorbirii destinat unor aplicații cum ar fi cele de citire a mesajelor de tip text (e-mail, SMS etc.) pe telefoane mobile, se verifică mai întâi dacă textul de intrare conține caractere scrise cu diacritice, iar dacă acest lucru nu se întâmplă, metoda se aplică prin intermediul unui algoritm eficient de inserare automată a semnelor diacritice.

Metoda are la bază o strategie statistică axată pe  $n$ -grame și constă în esență într-un proces de *filtrare secvențială*, pe mai multe niveluri, bazat pe probabilitatea de apariție a cuvintelor scrise cu diacritice și pe contextul în care acestea apar. Fiecare criteriu de filtrare aduce o îmbunătățire suplimentară, fără însă a mări caracterul de ambiguitate (deci formele posibile teoretic de scriere a unui cuvânt cu sau fără diacritice). Astfel, după ce tratează inițial cuvintele fără diacritice sau pentru care nu există ambiguități, metoda utilizează incremental trei etaje de filtrare: la nivel de *unigrame* (pe cuvânt), la nivel de *bigrame* (de asemenea pe cuvânt) și la nivel de *trigrame* pentru terminații de cuvânt, informația contextuală fiind adăugată gradual, după fiecare etaj, simultan cu eliminarea variantelor incorecte. Acest lucru se realizează practic prin folosirea probabilităților  $n$ -gramelor, calculate pe baza frecvențelor lor de apariție obținute din corpusul de antrenare.