



(12)

## CERERE DE BREVET DE INVENȚIE

(21) Nr. cerere: **a 2010 01246**

(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 PREDICȚIE A ACCENTULUI LEXICAL ÎN  
SINTEZA VORBIRII PORNIND DE LA TEXT ÎN LIMBA  
ROMÂNĂ**

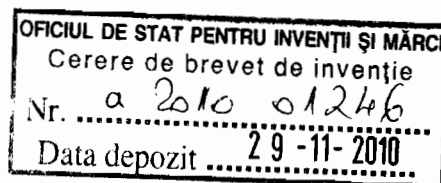
(57) Rezumat:

Invenția se referă la o metodă de predicție a accentului lexical, destinată utilizării în cadrul etapei de analiză prozodică, din cadrul unui etaj de prelucrare a limbajului natural într-un sistem de sinteză a vorbirii pornind de la un text în limba română, cu aplicații în dezvoltarea interfețelor inteligente de comunicare om-mașină. Metoda conform invenției utilizează o strategie statistică, bazată pe  $n$ -gramele construite la nivel de caracter, și folosește

un algoritm de netezire de tip "Katz back-off" modificat. Metoda necesită cunoștințe lingvistice limitate, are la bază un corpus de antrenare de dimensiuni medii și este foarte precisă în predicția poziției accentului lexical în limba română.

Revendicări: 1





## METODĂ DE PREDICȚIE A ACCENTULUI LEXICAL ÎN SINTEZA VORBIRII PORNIND DE LA TEXT ÎN LIMBA ROMÂNĂ

Invenția se referă la o metodă de predicție a accentului lexical, destinată utilizării în cadrul etapei de analiză prozodică din cadrul etajului de prelucrare a limbajului natural într-un sistem de sinteză a vorbirii pornind de la text în limba română, cu aplicații diverse în dezvoltarea interfețelor inteligente de comunicare om-mașină.

Sunt cunoscute un număr de studii pe plan mondial în problema poziționării automate a accentelor lexicale pentru sisteme de sinteză a vorbirii pornind de la text (TTS – “Text-to-Speech”), care se bazează fie pe o serie de reguli (“rule based”), fie sunt de tip statistic (“data driven”).

– De exemplu, rezultate recente foarte bune au fost raportate în [Q. Dou, S. Bergsma, S. Jiampojamarn, G. Kondrak, “A ranking approach to stress prediction for letter-to-phoneme conversion”, *Proceedings of the 47th Annual Meeting of the ACL and the 4th IJCNLP of the AFNLP*, Suntec, Singapore, 2009, pp. 118-126], folosind un algoritm de tip *data driven*, prin reprezentarea fiecărui cuvânt ca subșir de caractere. Algoritmul, cunoscut sub denumirea de SVM – “Support Vector Machine”, este antrenat a-priori pentru a clasifica toate modelele de accent. Autorii au raportat o acuratețe de 98% pentru predicția accentului primar (calculat la nivelul fonemelor) pentru limba engleză, 97,1% pentru germană și 93,1% pentru olandeză.

Dezavantajul principal al acestei metode constă în faptul că ea este destinată altor limbi decât limba română, fiind bazată pe particularități lingvistice complet diferite și de asemenea pe resurse lingvistice în format electronic ce nu sunt disponibile pentru limba română.

– În ceea ce privește limba română, singura încercare de soluționare a problematicii plasării accentelor lexicale din perspectiva folosirii în sistemele TTS, este descrisă în [E. Oancea, A. Badulescu, “Stressed syllable determination for Romanian words within speech synthesis applications”, *International Journal of Speech Technology*, vol. 5, no. 3, Kluwer,

2002, pp. 237-246.]. Autorii propun și evaluează experimental o metodă bazată pe reguli prin care calculează poziția silabelor accentuate în funcție de mulți parametri morfologici, fonetici și lexicali, și anume: numărul de caractere și silabe, tipul silabelor, secvențele de trei caractere de la începutul și sfârșitul cuvintelor. Acești parametri sunt folosiți în vederea antrenării sistemului prin împărțirea cuvintelor de intrare pe clase, iar în vederea recunoașterilor accentelor cuvintele de test sunt clasificate după clasa de apartenență. A fost raportată o eroare maximă ("Word Error Rate") de 6% pentru un corpus de test de 4.500 cuvinte, bază de test care a fost folosită și pentru generarea regulilor.

Dezavantajele metodei, în afară de rezultatele modeste, sunt un set incomplet de reguli și lipsa unor teste mai complexe și adecvate.

Problema tehnică pe care și-a propus să o rezolve invenția poate fi formulată pornind de la constatarea că modelarea prozodică a vorbirii generate de un sistem TTS este o etapă fundamentală în vederea obținerii unei vorbiri sintetizate cât mai naturale plecând de la un text de intrare oarecare. Realizarea unei prozodii corecte este esențială pentru un sistem TTS deoarece oferă informații referitoare la semnificațiile posibile ale mesajului vocal și la structura vorbirii, care nu rezultă implicit dintr-o simplă rostire segmentată. Pot fi date numeroase exemple în care mesajul vocal sintetizat devine ambiguu sau chiar neinteligibil din cauza faptului că informația prozodică fie lipsește, fie este incorect poziționată în fișierul audio generat. Pe de altă parte, trebuie menționat că prozodia este dificil de prezis în sistemele TTS din cauza faptului că textul de intrare conține puțină informație referitoare la structura și semnificația frazelor și că asemenea informații sunt greu de dedus automat. În acest context, este cunoscut faptul că *accentul lexical*, concept care apare la nivelul silabelor cuvintelor, are o influență majoră asupra conturului intonațional generat de un sistem TTS: silabele accentuate sunt pronunțate cu un efort mai mare decât celelalte silabe din același cuvânt și sunt caracterizate de un nivel mai înalt al frecvenței fundamentale și deseori de o durată mai mare a fonemelor.

Metoda care face obiectul prezentei invenții rezolvă problema tehnică menționată anterior prin faptul că utilizează o strategie statistică, bazată în principal pe  $n$ -gramele construite la nivel de caracter, destinată poziționării automate a accentelor lexicale pentru cuvintele scrise în limba română; abordarea statistică a fost aleasă ținând cont de faptul că poziția accentului lexical în limba română este variabilă, fiind practic imposibil de prezis pe bază de reguli.

În vederea rezolvării problemelor generate de insuficiența datelor de antrenare, metoda folosește un algoritm de netezire de tip *Katz back-off* modificat. În plus, în vederea reducerii efectelor nedorite cauzate de poziționarea incorectă a accentului lexical în vorbirea sintetizată, se utilizează un prag de decizie variabil. Cuvintele monosilabice sunt considerate ca fiind fără accent lexical și sunt puse în evidență cu ajutorul unui algoritm de despărțire automată în silabe. Metoda necesită cunoștințe lingvistice limitate, este bazată pe un corpus de antrenare de dimensiuni medii și este foarte precisă în predicția accentului lexical în limba română.

Se prezintă în continuare în detaliu obiectul invenției, un exemplu de aplicare a sa constând din utilizarea metodei de către autorii invenției în cadrul etajului de prelucrare a limbajului natural pentru un sistem complet de sinteză a vorbirii pornind de la text în limba română (dezvoltat în cadrul colectivului din care fac parte autorii).

Se poate spune că realizarea unei prozodii corecte este practic cea mai mare provocare pentru sistemele TTS actuale. În acest sens, un aspect foarte important, puternic legat de prozodie și de naturalitatea percepută a vorbirii sintetizate, este *accentul lexical*, concept care apare la nivelul silabelor cuvintelor și care are o influență majoră asupra conturului intonațional generat de un sistem TTS: silabele accentuate sunt pronunțate cu un efort mai mare decât celelalte silabe din același cuvânt și sunt caracterizate de un nivel mai înalt al frecvenței fundamentale și deseori de o durată mai mare a fonemelor. Astfel, semnalul vocal sintetizat la ieșirea unui sistem TTS este substanțial îmbunătățit prin prezența accentului lexical prin faptul că vorbirea sună mult mai natural.

Există mai multe definiții, teorii și modele de accent lexical. Deși unii autori apreciază că există mai multe niveluri de accent lexical (*primar, secundar* etc.) în interiorul aceluiași cuvânt, diferențierea fiind făcută după nivelul de accentuare, metoda propusă ia în considerare cea mai comună abordare, și anume aceea că **un cuvânt poli-silabic are o singură silabă accentuată (corespunzătoare accentului lexical primar)**. În multe alte limbi vorbite ca engleza, germana, spaniola sau româna, **accentul lexical are o poziție variabilă**, fiind extrem de dificil de prezis.

Înainte de a descrie metoda care face obiectul invenției, vom explica succint câteva particularități ale accentului lexical pentru limba română.

După cum am menționat anterior, accentul în limba română este liber (sau dinamic) și teoretic poate să fie poziționat pe oricare dintre silabele unui cuvânt (Tabelul 1).

Model	Exemplu
--	ca-fea
--	ca-să
---	pi-ja-ma
-- --	pe-re-te
---	re-pe-de
----	pa-ta-la-ma
--- --	pro-fe-soa-ră
-- ---	pe-re-te-le
----	la-cu-ri-le
-----	(a) în-tre-bu-in-ța
--- ---	as-cu-ți-toa-re
-- ---	li-pi-to-ri-le
-- ----	fa-mi-li-i-le
-----	(al) șap-te-spre-ze-ce-lea

**Tabelul 1.** Diferite situații de prezență a accentului lexical în limba română

Pe de altă parte, în limba română există multe cuvinte care sunt neaccentuate lexical. În această categorie sunt introduse o serie de *cuvinte funcționale* ca: articolele, prepozițiile și conjuncțiile, dar și cuvinte din alte categorii ca verbe și substantive. Deoarece există excepții de la aceste reguli, nu poate fi trasă concluzia că toate cuvintele funcționale sunt neaccentuate iar restul sunt accentuate. În urma unui studiu aprofundat, autorii invenției au găsit totuși o regulă care poate fi enunțată astfel: ***cuvintele monosilabice ale limbii române pot fi considerate ca fiind neaccentuate***; această regulă este consistentă cu definiția accentului dată anterior și a fost introdusă în metoda propusă. În plus (ca o consecință directă a regulii anterioare), am considerat că fiecare cuvânt polisilabic, poartă un singur accent lexical. Accentul acționează de fapt la nivelul unei singure vocale dintr-o silabă, iar caracterele vecine creează un context limitat în care efectul acestuia poate fi pus în evidență. După găsirea vocalei accentuate, cu ajutorul unui modul de despărțire automată în silabe (care nu face obiectul prezentei invenției), evidențierea silabei accentuate este imediată.

Metoda de predicție a poziției accentului lexical pentru limba română ce face obiectul invenției are la bază o abordare statistică. Metoda este de tip *data driven*, bazată pe *n*-grame, statisticile de limbaj fiind realizate la nivel de caracter.

Modelele sunt construite pe baza unui corpus de antrenare. Din datele de antrenare au fost extrase inițial toate  $n$ -gramele, după care modelele de  $n$ -grame au fost antrenate prin calcularea distribuțiilor de probabilitate. În implementarea propusă în cadrul invenției, numărul total de caractere este de 34, și anume: 31 de litere ale alfabetului limbii române (cinci litere fiind scrise cu diacritice), semnul (*markerul*) de accent ‘, simbolurile de start și de stop, “<” și respectiv “>”.

În mod evident, numărul total de  $n$ -grame posibile este foarte mare (pentru o lungime de 7 caractere acest număr este  $34^7$ ). Una dintre problemele cele mai importante care apare în astfel de implementări statistice este cea cauzată de insuficiența datelor de antrenare care necesită realizarea unui compromis între istoria  $n$ -gramelor (lungimea acestora) și dimensiunea corpusului de antrenare. Cu o istorie mai lungă, deci o lungime mai mare a  $n$ -gramelor, rezultatele pot fi sensibil mai bune, dar poate să apară o insuficientă antrenare din cauza fenomenului de *data sparseness* tradus prin lipsa unui număr mare de  $n$ -grame. Soluțiile cele mai folosite sunt cele de utilizare a tehnicilor de netezire (*smoothing*) pentru o re-estimare a probabilităților  $n$ -gramelor care lipsesc sau care sunt întâlnite foarte rar în corpusul de antrenare. Este general acceptată ideea că folosind  $n$ -grame de lungime mai mare, împreună cu tehnici de tip netezire, se obțin rezultate sensibil mai bune decât folosind  $n$ -grame de lungime mai mică dar pentru care modelul de limbaj este mai consistent, dacă se ia în calcul o aceeași dimensiune a corpusului de antrenare.

Pentru descrierea metodei să presupunem că fiecare cuvânt  $\mathbf{w}$  este compus din  $N$  litere:

$$\mathbf{w} = l_1 l_2 \dots l_N \quad (1)$$

Prin adăugarea caracterelor de start și stop,  $\mathbf{w}$  devine:

$$\underline{\mathbf{w}} = \langle l_1 l_2 \dots l_N \rangle \quad (2)$$

Vom exemplifica aceste notații cu cuvântul **casă**, pentru care  $N = 4$  (cuvântul are 4 caractere) și  $M = 2$  (cuvântul are două vocale: **a** și **ă**):  $\mathbf{w} = \mathbf{casă}$  devine  $\underline{\mathbf{w}} = \langle \mathbf{casă} \rangle$ .

După cum menționam anterior, accentul lexical acționează la nivelul vocalelor unei silabe, iar un cuvânt are un singur accent lexical. În abordarea propusă dacă  $\mathbf{w}$  are  $M$  vocale, fiecare dintre aceste  $M$  caractere poate avea marca accentului asociată cuvântului de start. Mai mult, o singură vocală poate fi accentuată. Soluția de predicție a poziției accentului lexical constă în găsirea celei mai bune variante dintre cele  $M$  posibile. Aceste variante sunt obținute prin plasarea markerului de accent în fața fiecărei vocale.

Prin urmare, pentru un cuvânt pentru care luăm decizia că poate avea accent lexical (fie acesta  $\underline{w}'$ ), generăm toate variantele posibile  $\mathbf{g}_j$  prin plasarea markerului de accent pe pozițiile  $j$ , cu  $j = 1, \dots, M$ ; pentru exemplul anterior vom avea două variante:

$$\mathbf{g}_1 = \langle \mathbf{c}'\mathbf{as}\mathbf{\check{a}} \rangle \text{ cu } g_{11} = \langle, g_{12} = \mathbf{c}, g_{13} = ', g_{14} = \mathbf{a}, g_{15} = \mathbf{s}, g_{16} = \mathbf{\check{a}}, g_{17} = \rangle, \text{ și}$$

$$\mathbf{g}_2 = \langle \mathbf{cas}'\mathbf{\check{a}} \rangle \text{ cu } g_{21} = \langle, g_{22} = \mathbf{c}, g_{23} = \mathbf{a}, g_{24} = \mathbf{s}, g_{25} = ', g_{26} = \mathbf{\check{a}}, g_{27} = \rangle$$

Se poate observa că numărul total de caractere pentru  $\underline{w}'$  este  $N + 3$  (pentru un cuvânt care are  $N$  caractere), dar pentru a simplifica notațiile, pentru restul explicațiilor vom considera că  $\underline{w}'$  are  $N$  caractere. Soluția găsirii celei mai bune variante dintre cele  $M$  posibile este dată prin folosirea modelelor de  $n$ -game.

### Etapa de antrenare

Pentru antrenarea modelelor de  $n$ -game am folosit un corpus  $T1$ , compus din cuvinte având accentul lexical poziționat manual iar delimitatorii de start și stop poziționați automat. Un corpus de tip *held-out*, denumit  $T2$  a fost folosit în etapa de antrenare în vederea optimizării parametrilor aleși în etapa de netezire.

Modelele au fost construite folosind datele de antrenare. Toate  $n$ -gamele posibile au fost mai întâi extrase din  $T1$  după care modelele de  $n$ -game au fost antrenate prin calculul distribuțiilor de probabilitate.

Pentru fiecare cuvânt din  $T1$  am extras toate variantele posibile de  $n$ -game și am calculat probabilitățile acestora. De exemplu, pentru cuvântul  $\langle \mathbf{c}'\mathbf{as}\mathbf{\check{a}} \rangle$ ,  $n$ -gamele extrase (de lungime maximă 7) sunt figurate în Tabelul 2.

Lungime $n$ -game	$n$ -game
1	$\langle, \mathbf{c}, \mathbf{a}, ', \mathbf{s}, \mathbf{\check{a}}, \rangle$
2	$\langle \mathbf{c}, \mathbf{c}', ' \mathbf{a}, \mathbf{as}, \mathbf{s}\mathbf{\check{a}}, \mathbf{\check{a}} \rangle$
3	$\langle \mathbf{c}', \mathbf{c}'\mathbf{a}, ' \mathbf{as}, \mathbf{as}\mathbf{\check{a}}, \mathbf{s}\mathbf{\check{a}} \rangle$
4	$\langle \mathbf{c}'\mathbf{a}, \mathbf{c}'\mathbf{as}, ' \mathbf{as}\mathbf{\check{a}}, \mathbf{as}\mathbf{\check{a}} \rangle$
5	$\langle \mathbf{c}'\mathbf{as}, \mathbf{c}'\mathbf{as}\mathbf{\check{a}}, ' \mathbf{as}\mathbf{\check{a}} \rangle$
6	$\langle \mathbf{c}'\mathbf{as}\mathbf{\check{a}}, \mathbf{c}'\mathbf{as}\mathbf{\check{a}} \rangle$
7	$\langle \mathbf{c}'\mathbf{as}\mathbf{\check{a}} \rangle$

**Tabelul 2.**  $n$ -gamele extrase pentru cuvântul  $\langle \mathbf{c}'\mathbf{as}\mathbf{\check{a}} \rangle$

În general, pentru un cuvânt  $\mathbf{w} = l_1 l_2 \dots l_N$ , un model de limbaj cu  $n$ -game poate fi extras ca o distribuție de probabilitate  $p(\mathbf{w})$ , obținută pe baza probabilităților condiționate astfel:

$$p(\mathbf{w}) = p(l_1)p(l_2 | l_1)p(l_3 | l_1l_2)\dots p(l_N | l_1l_2\dots l_{N-1}) = \prod_{i=1}^N p(l_i | l_1l_2\dots l_{i-1}) = \prod_{i=1}^N p(l_i | l_i^{i-1}) \quad (3)$$

unde  $l_i^{i-1}$  este o notație pentru  $l_1l_2\dots l_{i-1}$ .

În modelul  $n$ -gramelor facem aproximarea că probabilitatea unui caracter depinde numai de  $n-1$  caractere precedente:

$$p(\mathbf{w}) \approx \prod_{i=1}^N p(l_i | l_{i-n+1}l_{i-n+2}\dots l_{i-1}) = \prod_{i=1}^N p(l_i | l_i^{i-n+1}) \quad (4)$$

unde  $l_i^{i-n+1}$  reprezintă *istoria* caracterului  $l_i$ , iar lungimea  $n$ -gramelor ( $n$ ) mai este cunoscută și ca *ordinul modelului*. De exemplu, pentru  $n = 3$  se obține modelul *trigramelor*:

$$p(\mathbf{w}) \approx p(l_3 | l_1l_2)p(l_4 | l_2l_3)\dots p(l_N | l_{N-2}l_{N-1}) \quad (5)$$

În ecuația (4), probabilitatea condiționată  $p(l_i | l_i^{i-n+1})$  poate fi estimată pe baza principiului probabilității maxime ("*maximum likelihood*"), prin numărarea frecvențelor de apariție ale  $n$ -gramelor în corpusul de antrenare:

$$p(l_i | l_i^{i-n+1}) = \frac{c(l_i^{i-n+1})}{\sum_{l_i} c(l_i^{i-n+1})} \quad (6)$$

unde  $c(l_i^{i-n+1})$  numără de câte ori  $n$ -grama  $l_{i-n+1}l_{i-n+2}\dots l_i$  a apărut în corpusul de antrenare, iar  $\sum_{l_i} c(l_i^{i-n+1})$  este numărul total de apariții ale istoriei  $l_{i-n+1}l_{i-n+2}\dots l_{i-1}$ .

Totuși, o modelare care se ghidează numai după principiul ML nu este o alegere suficient de bună. Problema datelor insuficiente de antrenare obținute din  $T1$  va fi echivalentă cu un număr nul de apariții pentru  $n$ -gramele neobservate în timpul antrenării, iar emisia de probabilitate va fi  $p(\mathbf{w}) = 0$ , ceea ce va conduce la o probabilitate finală nulă, semnificând faptul că întregul cuvânt nu apare în datele de antrenare, deși doar o singură  $n$ -gramă nu a fost găsită. Rezultatul final va fi că la ieșire cuvântul de test va fi fără accent lexical (algoritmul nu va putea decide care este varianta corectă), sau va avea accentul poziționat greșit (ceea ce este și mai rău pentru sarcina TTS de care ne ocupăm) din cauza faptului că s-a obținut o probabilitate ne-nulă pentru o variantă eronată.

Pentru reducerea efectelor nedorite cauzate de insuficiența datelor de antrenare, am implementat în metoda propusă o tehnică de netezire de tip *Katz backoff*, care extinde algoritmul *Good-Turing* de estimare a probabilităților prin combinarea modelelor de  $n$ -game de ordin mai mare cu cele de ordin mai mic. Frecvențele de apariție pentru  $n$ -gramele generice  $l_i^{i-n+1}$  sunt modificate după formula:



$$c_{katz}(l_{i-n+1}^i) = \begin{cases} d_r r & \text{dacă } r > 0 \\ \alpha(l_{i-n+1}^{i-1}) c_{katz}(l_{i-n+2}^i) & \text{dacă } r = 0 \end{cases} \quad (7)$$

unde  $r$  este numărul de  $n$ -grame observate:

$$r = c(l_{i-n+1}^i), \quad (8)$$

$\alpha(l_{i-n+1}^{i-1})$  sunt coeficienți *backoff* de ponderare care asigură că suma probabilităților este egală cu 1, iar  $d_r$  este o rată de reducere dată de formula:

$$d_r = \begin{cases} 1 & \text{dacă } r > k \\ \frac{r^* \frac{(k+1)n_{k+1}}{n_1}}{1 - \frac{(k+1)n_{k+1}}{n_1}} & \text{dacă } r \leq k \end{cases} \quad (9)$$

unde  $k$  este o valoare de prag.

Trebuie specificat faptul că pentru  $r \leq k$  metoda Katz este obținută din algoritmul de estimare Good-Turing, iar pentru  $r > k$  numărul de apariții este cel obținut prin estimarea nemodificată a probabilității maxime. De asemenea, termenul  $r^*$  din ecuația (9) reprezintă numărul de  $n$ -grame, modificat (netezit) prin algoritmul *back-off*, și este dat de formula:

$$r^* = (r+1) \cdot \frac{n_{r+1}}{n_r} \quad (10)$$

unde  $n_r$  este numărul total de  $n$ -grame diferite care apar de  $r$  ori în datele de antrenare.

Noua probabilitate este calculată după formula:

$$p_{katz}^*(l_i | l_{i-n+1}^{i-1}) = \frac{c_{katz}(l_{i-n+1}^i)}{\sum_{l_i} c_{katz}(l_{i-n+1}^i)} \quad (11)$$

Filozofia metodei de tip Katz *backoff smoothing* este aceea că încearcă să mențină eficiența de rezoluție a  $n$ -gramelor, atunci când nu există un corpus de antrenare consistent, prin reducerea lungimii istoriei unei  $n$ -grame de ordin  $n$  pentru caracterul  $l_i$  la o istorie mai scurtă  $l_{i-n+2}^{i-1}$  pentru care există suficiente date de antrenare. Algoritmul coboară recursiv, dacă este nevoie, până la nivelul unigramelor.

În metoda *Katz backoff*, evenimentelor neobservate li se alocă o probabilitate diferită de zero prin extragerea acestora din evenimentele nenule (prin intermediul lui  $d_r$ ). Extragerea probabilității se face numai din  $n$ -gramele cu o frecvență de apariție sub ordinul  $k$  de prag.

Aparițiile extrase sunt distribuite către  $n$ -gramele neobservate de lungime  $n$  prin intermediul distribuțiilor de ordin  $n-1$ , dar cu păstrarea numărului total de apariții (*counts*).

În algoritmul Katz original, coeficienții de ponderare  $\alpha(l_{i-n+1}^{i-1})$  sunt calculați pentru fiecare istorie  $l_{i-n+1}^{i-1}$  astfel încât numărul total de apariții să se conserve, ceea ce se traduce prin

$$\sum_{l_i} c_{katz}(l_{i-n+1}^{i-1}) = \sum_{l_i} c(l_{i-n+1}^{i-1}).$$

Autorii acestei invenții **propun un algoritm simplificat** dar eficient care pleacă de la ipoteza că pentru aceeași lungime a  $n$ -gramelor, coeficienții de ponderare  $\alpha(l_{i-n+1}^{i-1})$  pot fi aceeași. Alegerea lor se face pe baza corpusului de tip *held-out T2* prin alegerea acelor coeficienți  $\alpha_{held-out}(n)$  care minimizează eroarea de predicție a accentului lexical. Se va observa totuși că, deși algoritmul nu garantează o sumă a distribuțiilor de probabilitate egală cu 1, acest lucru nu afectează performanța globală a metodei propuse.

Formula de calcul a emisiilor de probabilitate finale va fi dată de formula:

$$p_{katz}(l_{i-n+1}^{i-1}) = \begin{cases} p_{katz}^*(l_{i-n+1}^{i-1}) & \text{dacă } c(l_{i-n+1}^{i-1}) > 0 \\ \alpha_{held-out}(n) \cdot p_{katz}(l_{i-n+2}^{i-1}) & \text{dacă } c(l_{i-n+1}^{i-1}) = 0 \end{cases} \quad (12)$$

Într-un asemenea model recursiv probabilitățile de ordin  $n$  sunt calculate pe baza probabilităților de ordin  $n-1$ . Algoritmul se oprește la nivelul unigramelor prin alegerea distribuției de tip ML, astfel încât în această etapă probabilitățile bigramelor sunt calculate direct din probabilitățile unigramelor de tip "*Maximum Likelihood Estimation*":

$$p_{katz}(l_{i-1}^{i-1}) = \begin{cases} p_{katz}^*(l_{i-1}^{i-1}) & \text{dacă } c(l_{i-1}^{i-1}) > 0 \\ \alpha_{held-out}(n) \cdot p_{MLE}(l_i) & \text{dacă } c(l_{i-1}^{i-1}) = 0 \end{cases} \quad (13)$$

Pentru a putea realiza o recursie completă trebuie numărate toate aparițiile  $n$ -gramelor, până la nivelul unigramelor. Frecvențele de apariție sunt salvate în fișiere separate și apoi re-estimate prin algoritmul Good-Turing.

Pentru a însuma explicațiile anterioare, **etapa de antrenare** are următorii pași:

1. Din corpusul de antrenare  $T1$  (compus din cuvinte cu accentele poziționate corect și cu delimitatorii de start și stop pentru fiecare cuvânt), sunt extrase toate  $n$ -gramele, la nivel de caracter. Pentru o istorie de lungime  $n$ , sunt extrase de asemenea  $n$ -gramele de ordin  $n-1, n-2, \dots, 1$ .
2. Se numără toate aparițiile  $n$ -gramelor de ordin  $n, n-1, \dots, 1$  și apoi se calculează și se salvează probabilitățile de tip ML (sau numărul de apariții) rezultate.

3. Pe baza corpusului *held-out*  $T_2$  se estimează coeficienții de ponderare  $\alpha_{held-out}(n)$  care minimizează eroarea de predicție a accentului lexical în  $T_2$  (care în această etapă joacă rolul *de facto* al corpusului de test).
4. Folosind coeficienții de ponderare estimați în etapa a 3-a, se recalculează probabilitățile potrivit metodei Katz *backoff*. Se renunță la datele de antrenare anterioare.

### Etapa de testare

În etapa de test fiecare cuvânt din textul de intrare este comparat mai întâi cu toate intrările dintr-un dicționar de forme omografe; dacă un cuvânt este găsit în dicționar, acesta va fi lăsat nemodificat și trimis către un alt modul al sistemului TTS pentru o procesare ulterioară. În cazurile rare în care un cuvânt intră în algoritmul de poziționare a accentelor, rezultatul poate să nu fie o formă corectă semantic.

După cum am menționat anterior, modulul de despărțire în silabe din sistemul TTS este folosit pentru a identifica dacă un cuvânt este de tip monosilabic, situație în care acesta va fi lăsat nemodificat (deoarece astfel de cuvinte nu prezintă accent). Dacă un cuvânt are cel puțin două silabe, atunci acesta este procesat mai departe, după următorul algoritm:

1. Se pleacă de la forma  $\mathbf{w}$ , și se obține forma  $\underline{\mathbf{w}}$  prin plasarea delimitatorilor de start și stop:

$$\underline{\mathbf{w}} = l_1 l_2 \dots l_j \dots l_N, \text{ în care } l_1 = < \text{ și } l_N = > .$$

2. Pentru fiecare formă  $\underline{\mathbf{w}}$ , se generează toate variantele posibile accentuate  $\mathbf{g}_j$ ,  $j = 1, \dots, M$  (în care  $M$  este numărul de vocale din  $\underline{\mathbf{w}}$ ), prin inserarea unui caracter de accent înaintea fiecărei vocale.
3. Din fiecare variantă  $\mathbf{g}_j$  se extrag  $n$ -gramele de ordin  $n, n-1, \dots, 1$  și se calculează emisiile de probabilitate finale după formula generală din algoritmul Katz *backoff*, în care probabilitățile condiționate sunt cele obținute din datele salvate în timpul etapei de antrenare.
4. Varianta  $\mathbf{g}_{j,\max}$  care are probabilitatea maximă este furnizată la ieșirea sistemului.

Datele folosite în **experimentele de evaluare** ale metodei ce face obiectul invenției au fost extrase din texte literare scrise în limba română, totalizând peste 10 milioane de cuvinte, iar corpusul  $T$  de start conține peste 330.000 de cuvinte diferite (cuprinzând multe forme flexionate ale cuvintelor de bază). Acest dicționar a fost procesat mai departe prin plasarea

manuală a semnelor corespunzătoare accentelor lexicale înaintea vocalelor accentuate. După această etapă au fost înlăturate cuvintele omografe prin eliminarea automată a cuvintelor care au două sau mai multe forme accentuate, din aceste cuvinte rezultând un dicționar omografic. În plus, cuvintele monosilabice nu au fost introduse în algoritm, după cum am explicat anterior. În final,  $T$  a fost împărțit în trei părți egale (și diferite):  $T1$ ,  $T2$  și  $T3$ , potrivit Tabelului 3:  $T1$  a fost folosit pentru antrenarea modelelor  $n$ -gramelor,  $T2$  a fost folosit ca dicționar *held-out* pentru estimarea coeficienților de ponderare în algoritmul de tip *Katz-backoff*, iar  $T3$  a fost folosit pentru testare. Evaluarea a fost făcută în mod automat, prin compararea rezultatelor obținute în diferite condiții de testare, pe corpusul  $T3$  cu și fără markerii de accente la intrare.

Corpus	Denumire	Mărime din T(%)	Mărime în cuvinte
Antrenare	$T1$	75	234.700
<i>Held-out</i>	$T2$	10	31.300
Testare	$T3$	15	47.000

**Tabelul 3.** Corpusurile extrase din dicționarul  $T$

Ținând cont de faptul că principiile de funcționare ale etapelor de antrenare respectiv de testare ale metodei de poziționare ale accentelor au fost discutate anterior în detaliu, vom adăuga doar două aspecte particulare.

Astfel, în algoritmul de netezire *Katz backoff*, valoarea inițială propusă de Katz pentru pragul  $k$  a fost egală cu 5. În cazul metodei noastre, am testat rezultatele de ieșire pe corpusul  $T2$ , pentru  $k = 1, \dots, 20$  și am ales în final  $k = 10$  ca valoare optimă. Pe de altă parte vrem să subliniem ideea amintită anterior, și anume că o abordare originală propusă de algoritmul nostru este cea de calculare a coeficienților de ponderare  $\alpha_{held-out}(n)$ : **pentru fiecare  $n$ , cu  $n = 1, \dots, 8$  este evaluată valoarea optimă care minimizează erorile de predicție a accentelor pe corpusul *held-out*.** Această abordare duce la un algoritm mai eficient (adică mult mai rapid), cu o descreștere a acurateței rezultatelor de numai 0,02% față de algoritmul Katz original.

Ca măsură a performanței metodei propuse, am folosit **acuratețea la nivel de cuvânt**, definită ca raportul dintre numărul de cuvinte accentuate corect și numărul total de cuvinte de test (cu alte cuvinte, fiecare cuvânt cu un accent lexical poziționat incorect este considerat ca

fiind eronat). Această măsură de performanță este relevantă și consistentă cu definiția dată în literatura de specialitate, deoarece decizia (la nivel de cuvânt) poate avea doar două variante de ieșire, care se referă la poziționarea corectă sau incorectă a accentului. Metoda decide din mai multe variante, acestea fiind, după cum am menționat deja, egale cu numărul de vocale pentru cuvintele polisilabice.

Metoda a fost testată mai întâi cu diferite lungimi ale  $n$ -gramelor extrase la nivel de caracter, până la lungimea maximă de 8 caractere. Peste această lungime, potrivit datelor experimentale, nu am obținut o creștere de performanță. Principalele rezultate sunt ilustrate în Tabelul 4. Se poate observa că acuratețea crește rapid pentru  $n$ -gramele de ordin mic și se saturează pentru cele de ordin mare.

Lungime $n$ -grame	Acuratețe (%)
2	50,61
3	75,57
4	90,45
5	95,42
6	97,92
7	99,04
<b>8</b>	<b>99,11</b>

**Tabelul 4.** Acuratețea pentru diferite lungimi ale  $n$ -gramelor

Am obținut astfel o acuratețe maximă de **99,1%** pentru sarcina de predicție a accentului lexical, pentru  $n$ -gramele de lungime 8. Este evident că lungimea 8, aleasă ca ordinul maxim de mărime pentru  $n$ -grame, este mai mare decât dimensiunea unor cuvinte și se poate naște întrebarea modului în care sunt abordate aceste situații. Soluția a fost deja amintită și constă în folosirea unui algoritm de tip Katz *backoff*, care va reduce recursiv dimensiunea  $n$ -gramelor în cazul formelor care lipsesc din dicționarul de antrenare. Un alt aspect care merită amintit este acela că lungimea medie, exprimată în caractere, pentru corpusul de antrenare este de aproximativ 10,3 – valoare destul de apropiată de maximumul de 8 pentru care am obținut cele mai bune rezultate.

De asemenea, pentru o mai bună integrare a metodei cu celelalte componente ale sistemului TTS pentru limba română, am efectuat un **al doilea set de teste**.

Deoarece etajul de predicție a accentelor lexicale este o parte componentă a modului de analiză și generare prozodică, trebuie menționat faptul că o poziționare greșită a accentului lexical va duce la generarea unui contur prozodic eronat, ceea ce va conduce în final la producerea unei vorbiri care sună nenatural (sau chiar neinteligibil). Din această cauză, după o analiză a erorilor generate de metoda de predicție a accentelor lexicale, am decis să tratăm cazurile de cuvinte ambigue ca fiind neaccentuate, deoarece sistemul TTS, în ansamblul său, a dovedit că funcționează mai bine cu cuvintele neaccentuate decât cu cele având accentul poziționat eronat.

Am decis să analizăm în profunzime erorile de predicție pentru cuvintele din corpusul T3, analiză care a arătat faptul că pentru aproape 60% dintre cuvintele cu accent incorect poziționat, probabilitățile de emisie sunt comparabile până la un ordin de mărime. Pentru a aborda aceste cazuri, am propus un algoritm de testare ușor modificat, prin introducerea unui *prag de confidență* (nivel de decizie),  $\beta$ . Astfel, algoritmul va lăsa nemodificate cuvintele pentru care nivelul de decizie nu este trecut. Dezavantajul acestui tip de testare cu prag este că eroarea de predicție raportată va fi puțin mai mare deoarece cuvintele lăsate voite neaccentuate vor fi cuantificate ca fiind eronate.

Am efectuat o serie de experimente cu diferite valori pentru nivelul de prag  $\beta$ , constatând, așa cum era de așteptat, că la o creștere a lui  $\beta$  se obține o creștere a numărului de cuvinte furnizate neaccentuat la ieșirea algoritmului, prin netratarea situațiilor ambigue, și rezultând totodată o ușoară descreștere a acurateței. Se poate pune întrebarea: care ar putea fi cel mai bun criteriu de alegere a nivelului de prag? Deși răspunsul corect ar putea fi dat numai pe baza unei analize subiective a vorbirii generate (eventual printr-o estimare de tip *Mean Opinion Score*), o soluție teoretică imediată și rezonabilă (deși pot fi alese și altele) ar fi aceea de oprire a creșterii lui  $\beta$  atunci când numărul de cuvinte corecte, lăsate neaccentuat ca urmare a modificării pragului de decizie, îl egalează pe cel al cuvintelor eronate lăsate neaccentuat deoarece sunt clasificate ca fiind ambigue. În aceste condiții, am obținut o valoare de prag  $\beta = 30$ , pentru o acuratețe de aproximativ 97%. În opinia noastră această valoare nu trebuie tratată ca fiind prea mică deoarece ieșirea sistemului va fi foarte echilibrată, multe dintre cazurile de erori de predicție fiind de fapt clasificate ca fiind ambigue și lăsate neaccentuat. Rezultatele testelor sunt prezentate în Tabelul 5.

$\beta$	Acuratețe(%)
1	99,11
1,5	98,71
2	98,32
3	97,99
5	97,76
10	97,47
15	97,27
20	97,12
<b>30</b>	<b>96,92</b>
40	96,74
50	96,60

**Tabelul 5.** Rezultatele testării după varierea pragului de decizie

În acest punct al experimentelor, prin analiza erorilor obținute din corpusul *T3*, am încercat să realizăm o clasificare a acestora în scopul găsirii unor soluții de reducere a acestor erori.

Am putut identifica patru tipuri importante de erori:

**1. Erori de poziționare a accentului lexical, cauzate de metodă.** Reprezintă aproape 97,5% din numărul total de erori (0,86%).

Exemplu: **căm'ăși** a fost accentuat incorect ca fiind **cămăș'i**.

Astfel de situații sunt cauzate de existența în corpusul de antrenare a unui spectru larg de *n*-grame asemănătoare, provenind de regulă de la diverse categorii gramaticale (în exemplul menționat de la substantive și verbe).

Am efectuat, totodată, o analiză a modului de variație a acurateței metodei cu dimensiunea corpusului de antrenare. Concluzia a fost că acuratețea crește odată cu mărirea corpusului de antrenare, după cum rezultă din Tabelul 6. Rezultatele au fost obținute folosind ca bază de test corpusul *T3*, iar ca bază de antrenare diferite segmente extrase din *T1*, în mod aleator.

Dimensiune T1 (cuvinte)	Acuratețe (%)
10.161	91,61
54.818	94,42
99.526	95,53
144.234	96,62
188.942	97,37
<b>234.700</b>	<b>99,11</b>

**Tabelul 6.** Creșterea acurateții cu variația dimensiunii corpusului de antrenare

**2. O decizie incorectă referitoare la cuvintele monosilabice inserată de modulul de despărțire în silabe.** Aceste erori însumează aproximativ 1,7% din numărul total de erori (0,015% din totalul cuvintelor).

Exemplu: **urci** a fost introdus în algoritm, deși are o singură silabă, fiind incorect despărțit în silabe (adică **ur-ci**), la ieșirea algoritmului fiind furnizată varianta **urc'i**.

**3. Forme omografe incorecte introduse în corpusul de test.** Aceste erori totalizează aproximativ 0,4% din numărul total de erori (ceea ce înseamnă 0,004% din totalul cuvintelor).

Exemplu: **p'oza** nu a fost găsit în dicționarul de forme omografe și a fost accentuat ca fiind **poz'a**. Aceste tipuri de erori sunt provocate de faptul că dicționarul de forme omografe este deocamdată incomplet.

#### **4. Introducerea manuală incorectă a formelor accentuate în corpusul de test**

Exemplu: Cuvântul **biciuie** a fost introdus în corpusul de test ca **bici'ui** (cuvânt lipsit de semnificație). Algoritmul a inserat corect markerul corespunzător accentului adică **b'iciuie** însă datorită comparației cu cuvântul din baza de test, a fost raportat eronat ca fiind o eroare. Aceste tipuri de erori apar de multe ori în corpusurile mari adnotate manual, și duc la o mică scădere a performanțelor (în cazul nostru scăderea a fost de 0,3% (0,0024% din totalul cuvintelor).



## REVENDICĂRI

**Metodă de predicție a accentului lexical în sinteza vorbirii pornind de la text în limba română**, caracterizată prin aceea că ajută la modelarea corectă a conturului intonațional pentru sinteza automată a vorbirii de bună calitate prin poziționarea corectă a accentelor lexicale la nivelul cuvintelor, asigurând astfel o prozodie corespunzătoare vorbirii generate.

Metoda utilizează o strategie statistică, bazată în principal pe  $n$ -gramele construite la nivel de caracter, fiind destinată poziționării automate a accentelor lexicale pentru cuvintele scrise în limba română; abordarea statistică a fost aleasă ținând cont de faptul că poziția accentului lexical în limba română este variabilă, fiind practic imposibil de prezis pe bază de reguli. În vederea rezolvării problemelor generate de insuficiența datelor de antrenare, metoda folosește un algoritm de netezire de tip *Katz back-off* modificat, după următorul principiu: din fiecare variantă posibilă se extrag  $n$ -gramele de ordin  $n, n-1, \dots, 1$  și se calculează emisiile de probabilitate finale, probabilitățile condiționate fiind cele obținute în timpul etapei de antrenare. În plus, în vederea reducerii efectelor nedorite cauzate de poziționarea incorectă a accentului lexical în vorbirea sintetizată, se utilizează un prag de decizie variabil. Cuvintele monosilabice sunt considerate ca fiind fără accent lexical și sunt puse în evidență cu ajutorul unui algoritm de despărțire automată în silabe. Metoda necesită cunoștințe lingvistice limitate, este bazată pe un corpus de antrenare de dimensiuni medii și este foarte precisă în predicția poziției accentului lexical în limba română.