

# Aplicații Pentru Recunoașterea Dinamică a Limbii

**Paul FOGARASSY-NESZLY**

BAUM Engineering

Arađ 310175, str. T.Moșoiu nr. 8  
pf@baum.ro

**Vasile GHERHES**

Universitatea "Politehnica" Timișoara

Timișoara 300006, Piața Victoriei nr.2  
vasile.gherhes@upt.ro

## REZUMAT

În acest articol sunt prezentate două aplicații care implementează două metode diferite de recunoaștere automată a limbii, în scopul sintezei vocale. Sunt descriși algoritmi implementați, condițiile specifice de funcționare a aplicațiilor de acest gen și sunt prezentate interfețele aplicațiilor.

În timp ce identificarea limbii prin intervalul Unicode al textului consumă cel mai puțin resursele sistemului și este foarte rapid, algoritmi de recunoaștere prin analiza statistică a textului consumă mai multe resurse și sunt ceva mai lenți. Găsirea unui compromis optim cade atât în sarcina dezvoltatorului, cât și în sarcina utilizatorului care trebuie să decidă în faza de configurare a aplicației care sunt condițiile concrete de lucru.

În ultima parte a lucrării sunt prezentate o serie de considerente privind inerția necesară la schimbarea limbii atunci când aplicația este utilizată în contextul sintezei vocale diferențiate lingvistic a textului.

## Cuvinte cheie

Recunoașterea limbii, sinteză vocală, accesibilitate, tehnologii asistive.

## Clasificare ACM

H5.2. Information interfaces and presentation (e.g., HCI): Miscellaneous.

## INTRODUCERE

Spre deosebire de recunoașterea automată a limbii (RAL) în care este scris un text, recunoașterea dinamică a limbii (RDL) presupune implicarea unor algoritmi diferiți, capabili să identifice limba pe baza unui flux continuu de date. Această restricție majoră de abordare a problemei este impusă de aplicațiile de tip sinteză vocală (Text To speech - TTS).

O altă aplicație a RDL are în vedere sinteza vocală diferențiată lingvistic la traducerea automată a unui text, bazată pe recunoașterea vocală [1,2]

Sinteza vocală este folosită mai ales de către persoanele cu deficiente de vedere, sau dificultăți de citire (dislexici sau analfabeți) pentru accesibilizarea documentelor în format electronic. Atunci când există posibilitatea ca texte scrise în limbi diferite să alterneze, este necesar un algoritm capabil să discrimineze limba în care este scris fragmentul de text, deoarece sinteza vocală se bazează pe particularitățile fonetice ale limbii în care este scris textul, precum și pe regulile de scriere specifice limbii respective. Așadar, este obligatorie folosirea sintezei vocale

corespunzătoare limbii, pentru ca rezultatul să fie comprehensibil.

Algoritmii utilizați la RDL se bazează pe identificarea intervalului de codificare a caracterelor (character encoding detection) [3] și/sau pe modele statistice [4,5].

## APLICAȚIE PENTRU RECUNOAȘTEREA AUTOMATĂ A LIMBII PRIN METODA INTERVALULUI DE CODIFICARE A CARACTERELOR

Una dintre metodele pentru identificarea limbii în care este scris un text (în mod dinamic sau nu) este prin identificarea codificării caracterelor (character encoding detection) pentru textul analizat [3]. Unicode este un format definit de către Unicode Consortium pentru codarea, stocarea și interpretarea textelor în mediul informatic; acesta este standardul de codificare de facto utilizat la interpretarea datelor binare în format text [6].

Aplicația permite definirea limbii pentru intervalele Unicode stabilite, ceea ce permite schimbarea dinamică a limbii, în timpul citirii textului. Dezavantajul metodei constă în faptul că se comportă haotic în cazul în care texte scurte scrise cu caractere diferite sunt prezente frecvent în interiorul aceluiași text. Utilizatorul experimentează într-o asemenea situație modificări rapide a vocilor cu care este citit textul, ceea ce produce de obicei un rezultat greu de înțeles.

Figura 1 prezintă interfața de configurare a aplicației; după cum se poate vedea, se pot configura vocile cu care sunt citite textele scrise cu caractere latine obișnuite, cu caractere arabice, grecești, ebraice și chirilice. De asemenea, se pot configura parametrii de viteză și volum a fiecărei voci în parte, cu scopul de a armoniza vocile cu care este citit un text.

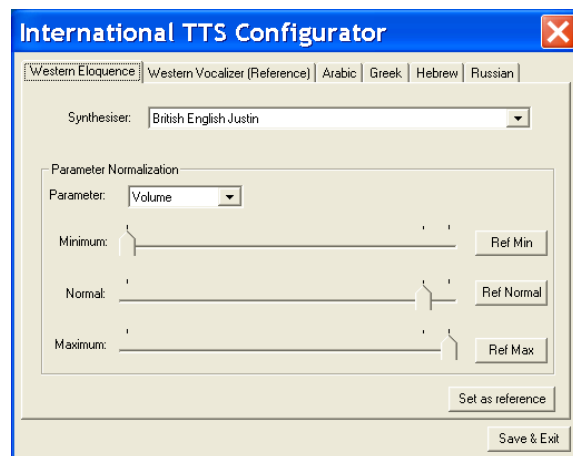


Figura 1 Interfața de configurare a aplicației

Avantajul metodei utilizate este că încarcă foarte puțin procesorul și nu necesită memorie; dezavantajul constă că trebuie prevăzut în mod explicit o oarecare inerție la schimbarea vocii, pentru a evita schimbarea prea frecventă a vocii în cazul textelor insuficient prelucrate, așa cum rezultă acestea de obicei din procesul de recunoaștere optică a caracterelor (OCR - Optical character Recognition).

**APLICAȚIE PENTRU RECUNOAȘTEREA LIMBII PRIN METODE STATISTICE**

Cea mai avansată metodă de identificare a limbii constă în analiza n-gramelor. Aceasta este cea mai utilizată metodă statistică și are mai multe variante, fiecare cu avantajele și dezavantajele sale, aplicabilă pentru RAL și RDL

O n-gramă este o sub-secvență de n elemente dintr-o secvență dată; în general secvența de elemente poate fi orice, de la caractere și până la cuvinte. În analiza lingvistică n-gramele sunt utilizate mai mult pentru cuvinte sau pentru caractere. În această lucrare, prin n-gramă se înțelege o secvență de n-caractere succesive dintr-un text. Atunci când este vorba de două caractere (n = 2) se mai folosește termenul de bigramă (sau digramă), iar când este vorba de succesiuni de trei caractere (n = 3) termenul consacrat este trigramă. În Tabelul 1 este prezentat un exemplu de descompunere în n-grame a cuvântului „analiză”.

Tabelul 1 Bigramele, trigramele și 4-gramele cuvântului „analiză”

	analiză
2-gram	_a , an , na , al , li , iz , ză , ă_
3-gram	_an , ana , nal , ali , liz , iză , ză_
4-gram	_ana , anal , nali , aliz , liză iză_

În Tabelul 3.1, începutul de cuvânt și sfârșitul de cuvânt a fost marcat printr-un caracter special, ne-literal ( \_ ). Este relevantă din punct de vedere lingvistic această marcare, deoarece localizarea unui grup de caractere la începutul sau la sfârșitul cuvântului este semnificativă statistic pentru caracterizarea unei limbi.

Discriminarea limbilor pe baza n-gramelor (în general) pleacă de la observația că pentru fiecare limbă anumite n-gramme apar mai frecvent decât altele. Studii experimentale realizate de [7] au arătat că utilizarea trigramele conduce la cele mai bune rezultate. Identificarea limbii se face prin compararea frecvenței de apariție a trigramele în textul analizat cu frecvența acestora în corpusurile limbilor care sunt avute în vedere.

În faza de „antrenare” a aplicației se construiește spectrul de frecvențe al n-gramelor pentru fiecare limbă în parte. Acesta se bazează pe un corpus relevant pentru limba avută în vedere și domeniul de aplicare (dacă este cazul). Din studii realizate de [8, 9] rezultă că un corpus de circa 50.000 de cuvinte oferă o precizie foarte bună care nu mai crește semnificativ prin mărirea volumului. Frecvențele relative ale n-gramelor reprezintă caracteristica fiecărei limbi, iar frecvența n-gramelor din textul analizat se realizează în timpul rulării, în același mod. În funcție de modul de construire a acestui spectru, metoda poate fi mai rapidă sau mai lentă, mai precisă sau mai puțin precisă.

Corpusul trebuie să fie omogen din punct de vedere al limbii caracterizate de acesta și trebuie să fie corect gramatical și sintactic; calitatea corpusului are o influență hotărâtoare asupra preciziei de identificare a limbii.

Compararea spectrului analizat cu cele de referință se poate face în diverse feluri. Cel mai simplu criteriu este suma abaterilor absolute (ecuația 1) sau suma abaterilor pătrăte (ecuația 2).

$$A_L = \sum_{i=1}^m |f_{ai} - f_{Li}| \tag{1}$$

unde  $A_L$  este abaterea frecvențelor pentru limba  $L$ ,  $m$  este numărul de n-gramme din textul analizat,  $f_{ai}$  este frecvența n-grammei  $i$  din textul analizat, iar  $f_{Li}$  este frecvența n-grammei  $i$  din spectrul de frecvențe al limbii  $L$ .

$$A_L = \sum_{i=1}^m (f_{ai} - f_{Li})^2 \tag{2}$$

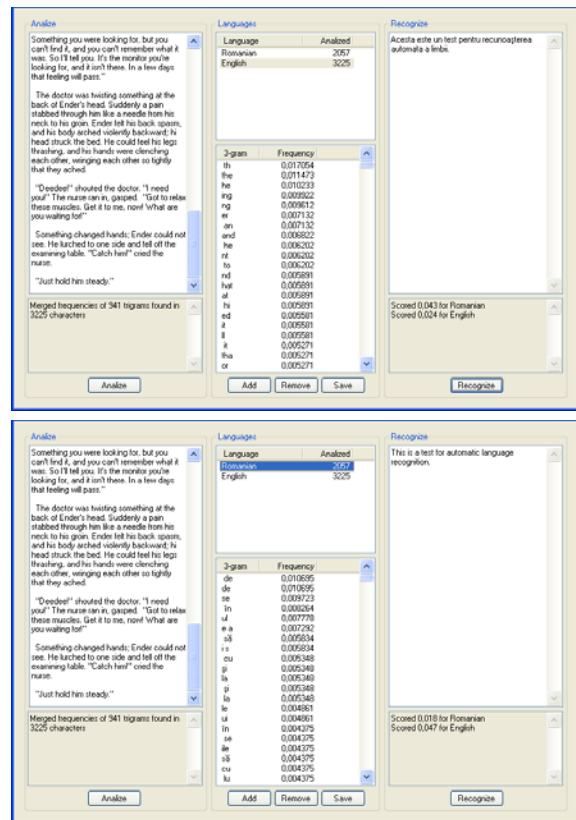


Figura 2. Interfața aplicației de testare a algoritmilor de analiză statistică

Limba identificată cea mai probabil este limba  $L$  pentru care  $A_L$  este minim.

Criterii mai sofisticate de calcul al abaterii  $A_L$  țin cont printr-un coeficient de pondere de probabilitatea mult mai mare sau mult mai mică (chiar zero) a unor n-gramme pentru o anumită limbă. În acest caz, abaterea  $A_L$  se calculează conform relației (3).

$$A_L = \sum_{i=1}^m k_{Li} (f_{ai} - f_{Li})^2 \tag{3}$$

unde  $k_{Li}$  este coeficientul de pondere pentru n-grama  $i$  din limba  $L$ ; spre deosebire de ecuațiile (1) și (2), deoarece  $k_{Li}$  poate fi pozitiv sau negativ, și  $A_L$  poate avea valori pozitive sau negative.

Figurile 3.1 prezintă interfața aplicației de testare a algoritmilor de analiză statistică. În exemplul prezentat este evaluat cel mai simplu algoritm; sunt introduse două mini-corpusuri de numai câteva mii de caractere pentru limba română (2057 caractere) și limba engleză (3225 caractere). În urma analizei acestora, sunt identificate toate trigramele și se calculează frecvența acestora. Textul analizat este la rândul său descompus în trigrame, iar pentru fiecare dintre acestea se construiește o sumă a frecvenței în cele două corpusuri. Limba în care este scris textul este probabil cea care are cea mai mare sumă caracteristică. După cum se poate vedea, diferența dintre cele două sume este semnificativă pentru ambele texte analizate.

Tabelul 2 prezintă rezultatul exemplului prezentat în condițiile din Figura 2.

Tabelul 2. Rezultatele testului preliminar

Text analizat	Limba	Scor
Acesta este un test pentru recunoașterea automată a limbii.	Română	0,043
	Engleză	0,024
This is a test for automatic language recognition.	Română	0,018
	Engleză	0,047

Algoritmii care se bazează pe analiza statistică a textului în vederea recunoașterii limbii sunt mai complecși și pentru a nu încălca procesorul în timpul rulării este necesară optimizarea bazelor de date obținute în faza de antrenare. De asemenea, algoritmii trebuie optimizați pentru un consum de memorie cât mai redus, în special în cazul rulării aplicației pe dispozitive mobile.

Optimizările sunt necesare deoarece în mod normal în paralel cu aplicația de recunoaștere funcționează o aplicație de sinteză vocală, aplicația client care folosește sinteza vocală, foarte probabil mai funcționează un cititor de ecran (cu sau fără funcție de magnificare a ecranului) și în sfârșit, aplicația utilizatorului.

#### CONSIDERENTE PRIVIND COMPORTAMENTUL APLICAȚIILOR DE RECUNOAȘTERE A LIMBII

În mod intuitiv, prin recunoașterea limbilor în care este scris un text se dorește etichetarea tuturor fragmentelor acestuia în mod corespunzător. Acest lucru poate fi util pentru aplicații de clasificare și analiză bibliografică.

Scopul aplicațiilor prezentate în această lucrare este de a permite sinteza vocală diferențiată a unui text scris în două sau mai multe limbi diferite, așa cum se întâlnește relativ frecvent în zonele multiculturale.

Teste preliminare au arătat că o etichetare excesiv de riguroasă a unui text poate duce la modificarea limbii în care este citit textul mult prea frecvent. Fiecare limbă conține împrumuturi lingvistice sau substantive proprii din alte limbi. La realizarea unei cărți audio sau la lectura prin voce sintetică, comună nevăzătorilor, dar și dislexicilor, a citi cu o voce diferită un nume de origine străină poate fi foarte deranjant. La fel de deranjant poate fi și modificarea

vocii pentru redarea unor fragmente scurte, sau chiar a unor cuvinte disparate.

De exemplu, în limba română au pătruns în ultimele decenii (justificat sau nu) multe cuvinte și expresii din limba engleză. Dacă azi a devenit comună prezența în textele tehnice a unor cuvinte cum sunt „mouse” sau „webpage”, în texte comune găsim frecvent cuvinte sau expresii cum sunt „weekend”, „city break”, „team building”, „fitness” etc. De asemenea, în limba română se păstrează (aproape întotdeauna) ortografia originală a toponimelor și a numelor proprii; în română este firesc să scriem „New-York” și nu „Nowy Jork”, cum se scrie în poloneză.

Algoritmii de analiză statistică vor identifica repede faptul că fragmentele care conțin literele „k”, „y” și „w” sunt mult mai probabil în limba engleză, iar literele duble „ee” sau „ss” apar foarte rar în limba română, în timp ce în engleză sunt relativ frecvente. Citirea în engleză a acestor cuvinte presărate într-un text scris aproape integral în română poate fi însă foarte supărătoare.

Pentru a evita asemenea situații există mai multe posibilități: fie este introdusă o anumită inerție la schimbarea limbii, prin care se evită citirea cu voce diferită a unor cuvinte sau fragmente relativ scurte, fie se definește o listă de cuvinte străine frecvente în limba de bază, care sunt ignorate de aplicație și deci vor fi citite cu vocea considerată implicită.

Pe de altă parte, aplicațiile de sinteză vocală lucrează sincron, pe baza unui flux de date (text în acest caz) și generează cu o mică întârziere tot un flux de date (audio în cazul sintezei vocale). Datorită acestui lucru și în funcție de decalajul temporar admis între cele două fluxuri, este aproape imposibil ca aplicația de recunoaștere să reacționeze imediat și să schimbe corespunzător vocea. Inevitabil, la schimbarea limbii va exista cel puțin un fragment care va fi citit cu vocea anterioară. În cazul unor cuvinte izolate sau a unor fragmente scurte, acest lucru poate duce la rezultate greu de înțeles pentru utilizator.

Inerția la modificarea vocii reprezintă cea mai bună metodă de evitare a unor asemenea situații și este perfect admisibilă în cazul sistemelor care funcționează sincron.

Desigur, în cazul aplicațiilor de conversie text-audio, utilizate la realizarea cărților audio sau a cărților Daisy, textul poate fi parcurs în întregime și etichetat corespunzător înainte de conversia audio propriu-zisă. În acest gen de aplicații întârzierile mari între fluxul text de la intrare și fluxul audio rezultat sunt acceptabile.

Soluția cea mai bună pentru managementul erorilor de conversie admisibile datorită inerției se rezolvă în mai multe feluri. Pe lângă inerția, menționată anterior, s-au elaborat algoritmi care păstrează un fragment de text sintetizat anterior ca și referință. În cazul în care limba unui fragment de text este discriminată cu incertitudine (definită statistic în faza de configurare), atunci limba de referință este cea din fragmentul anterior. O altă opțiune de configurare permite definirea unei limbi implicite, atunci când fragmentele scrise într-o limbă diferită sunt puțin frecvente.

Din testele preliminare, a rezultat că algoritmi de recunoaștere a limbii sincroni se comportă mult mai bine în cazul în care sunt definite doar două limbi posibile. În cazul în care un text poate fi scris în trei limbi, acest gen de algoritmi necesită configurări specifice suplimentare pentru ca rezultatul să rămână foarte bun.

#### CONCLUZII ȘI DIRECȚII DE CERCETARE

Această lucrare prezintă două aplicații generice de RAL aplicabile și pentru RDL. Sunt prezentate particularitățile algoritmilor în timpul rulării împreună cu cititoarele de ecran, pe dispozitive mobile, relativ sărace în resurse hardware (memorie și procesor), precum și faptul că textul care trebuie analizat este scurt.

Analiza performanțelor algoritmilor implementați, precum și propunerile de îmbunătățire urmează să fie identificate în timpul activității de testare în diverse condiții, cum sunt de exemplu texte scurte dar care conțin totuși cuvinte din limbi diferite. Un amplu program experimental semi-automatizat urmează să valideze algoritmi care vor fi propuși de analizați. De asemenea, tot experimental, se vor identifica pragurile optime de incertitudine.

#### CONFIRMARE

Această lucrare a fost elaborată în cadrul contractului 29DPST/13.09.2013, „Aplicație pentru Conversia din Text în Voce Sintetică cu Recunoașterea Automată a Limbii”, în Programului Inovare, Dezvoltare Sisteme-Produse-Tehnologii a UEFISCDI.

#### REFERINȚE

1. Mathias, L. (2007). Statistical Machine Translation and Automatic Speech Recognition under Uncertainty. PhD thesis, Johns Hopkins University

2. Moore, R. C. and Quirk, C. (2008). Random restarts in minimum error rate training for statistical machine translation. In Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008), pages 585–592, Manchester, UK. Coling 2008 Organizing Committee
3. Timothy Baldwin and Marco Lui, Language Identification: The Long and the Short of the Matter, Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the ACL, pages 229–237, Los Angeles, California, June 2010
4. Gavin Churcher, „Distinctive character sequences”, 1994, personal communication
5. Dunning, T. "Statistical Identification of Language". Technical Report MCCS 94-273, New Mexico State University, 1994
6. The Unicode Standard, Version 5.0, Fifth Edition, The Unicode Consortium, Addison-Wesley Professional, 27 October 2006
7. Cavnar, W., and Trenkle, J. (1994). N-gram-based text categorization. Proc. 3rd Symp. on Document Analysis and Information Retrieval (SDAIR-94)
8. Dunning, T. (1994) Statistical Identification of Language. Technical Report MCCS 94-273, New Mexico State University
9. Ljubešić, N., Mikelić, N., and Boras, D. (2007) Language identification: How to distinguish similar languages. In Lužar-Stifter, V. and Hljuz Dobrić, V., editors, Proceedings of the 29th International Conference on Information Technology Interfaces, pages 541–546, Zagreb SRCE University Computing Centre