

Testarea versiunii beta a unei componente de recunoaștere dinamică a limbii – rezultate preliminare

Costin Pribeanu¹, Paul Fogarassy-Neszly²

¹ Institutul Național de Cercetare-Dezvoltare în Informatică – ICI București
Bd. Mareșal Averescu, Nr. 8-10, 011455, București
E-mail: pribeanu@ici.ro

² BAUM Engineering
Str. Traian Moșoiu nr. 8, 310175 Arad
E-mail: pf@baum.ro

Rezumat. Recunoașterea automată a limbii în care este scris un text face parte din categoria generală a sarcinilor (recunoașterea este o problemă, o sarcină, o aplicație, care se rezolvă pe baza unui algoritm), de clasificare a textelor și are numeroase aplicații. Unele implementări recente au în vedere sinteza vocală diferențiată lingvistic. În acest caz, este necesară o recunoașterea dinamică a limbii, care presupune implicarea unor algoritmi diferiți, capabili să identifice limba pe baza unui flux continuu de date. Această lucrare prezintă rezultatele preliminare obținute în testarea a două versiuni succesive (alpha și beta) a unei componente de recunoaștere dinamică a limbii, cu accent pe a doua versiune. În testare au fost manipulați doi parametri: volumul de text analizat și inerția la schimbarea limbii. Evaluarea a fost făcută din perspectiva utilizabilității aplicațiilor software, cu accent pe eficacitatea și eficiența recunoașterii limbii. Rezultatele testării versiunii beta sunt utile pentru determinarea valorii optime a parametrilor și sugerează noi direcții de optimizare a metodei și a algoritmilor.

Cuvinte cheie: recunoaștere automată a limbii, recunoaștere dinamică a limbii, testare beta, sinteză vocală, utilizabilitate, tehnologii asistive.

1. Introducere

Recunoașterea automată a limbii, un domeniu de interes atât pentru lingvistica computațională cât și pentru accesibilitatea sistemelor informatice, are numeroase aplicații, între care menționăm sinteza vocală diferențiată lingvistic și traducerea automată (Kłosowski & Dustor, 2013; Devadoss, 2010).

În acest articol se prezintă rezultatele preliminare obținute la testarea unei componente care are ca scop recunoașterea dinamică a limbii în care

este scris un text, în vederea sintezei vocale. Necesitatea recunoașterii dinamice este impusă de modul de generare a fluxului audio, aproape simultan cu primirea fluxului de text (Dunning, 1994; Zissman, 1996). Deși există aplicații software capabile să recunoască limba, precum și aplicații capabile să redea prin voce sintetică un anumit text (TTS – Text-To-Speech), o aplicație care să îndeplinească ambele sarcini nu există.

Sinteza vocală este folosită mai ales de către persoane cu deficiențe de vedere, sau dificultăți de citire (persoane cu deficiențe cognitive, dislexici sau analfabeți), în scopul accesibilizării documentelor. Atunci când există posibilitatea ca texte scrise în limbi diferite să alterneze, este necesară discriminarea limbilor 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. Cu alte cuvinte, este necesară selectarea sintezei vocale corespunzătoare limbii, pentru ca rezultatul să fie comprehensibil.

Între sistemele software care au nevoie de această componentă menționăm aplicațiile informatice pentru nevăzători (de exemplu, cititoarele de ecran), dispozitivele de citire automată pentru deficienți de vedere, aplicațiile care utilizează sisteme GPS și aplicațiile cu redare vocală (En: self-voicing applications).

2. Metode și algoritmi de recunoaștere a limbii

2.1 Metode de recunoaștere a limbii

Există numeroase metode de recunoaștere a limbii (Athias, 2007; Baldwin & Lui, 2010) precum și abordări care utilizează metode combinate (Gînscă et al., 2011). Întrucât într-o lucrare anterioară (Fogarassy-Neszly & Gherhes, 2013) a fost făcută o prezentare mai detaliată, în cele ce urmează le vom aminti succint împreună cu avantajele și dezavantajele fiecăreia.

Identificarea limbii în care este scris un text poate fi făcută uneori cu ușurință prin simpla identificare a codificării caracterelor (character encoding detection) pentru textul analizat (Baldwin & Lui, 2010). Metoda nu poate fi utilizată (de exemplu) pentru discriminarea limbilor engleză și franceză sau a limbilor rusă și bulgară, deoarece aceste perechi de limbi folosesc caractere din același interval Unicode.

O altă metodă constă în căutarea cuvintelor (inclusiv a formelor flexionare) în dicționare specifice limbilor respective. Deși are avantajul de a fi precisă chiar și pentru texte scurte, metoda nu este practică datorită volumului foarte mare al dicționarelor și faptului că implică algoritmi lenți.

O altă metodă intuitivă și simplă, se bazează pe identificarea unor caractere speciale tipice unei limbi. Acestea sunt de obicei caracterele cu diacritice și ligaturile. Folosirea caracterelor speciale pentru discriminarea limbilor are aplicabilitate limitată, deoarece în cazul unor texte scurte acestea pot să nu apară. De asemenea, unele dintre caracterele speciale sunt comune mai multor limbi.

O altă abordare a plecat de la identificarea unor secvențe de caractere specifice exclusiv unei singure limbi (Churcher, 1994; Churcher et al., 1994, Dunning, 1994). Deși tentantă prin simplitate, asemenea secvențe de caractere izolate nu oferă suficientă încredere în identificarea limbii pentru care se presupune că sunt specifice, datorită unui număr mare de excepții posibile, precum și din cauza cuvintelor preluate în alte limbi și scrise conform regulilor din limba de origine.

Metoda secvenței de caractere specifice a condus la cea mai utilizată metodă și anume metoda *n-gramelor* (Dunning, 1994; Cavnar & Trenckle, 1994). 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. 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.

2.2 Aspecte specifice privind implementarea algoritmilor

Discriminarea limbilor pe baza *n-gramelor* (în general) pleacă de la observația că pentru fiecare limbă anumite *n* grame apar mai frecvent decât altele. Identificarea limbii se face cel mai simplu prin compararea frecvenței de apariție a trigramelor 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). Rezultatele unor studii (Dunning, 1994; Ljubesi et al., 2007) arată că un corpus de circa 50.000 de cuvinte oferă o precizie foarte bună care nu mai crește semnificativ prin mărirea volumului.

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.

Spre deosebire de implementarea algoritmilor generici de identificare a limbii, algoritmi folosiți la sinteza vocală diferențiată sunt utilizați în condițiile limitării la minimum a limbilor posibile; în mod normal, aplicațiile care funcționează în medii multiculturale trebuie să distingă între două, mai rar trei și foarte rar patru limbi. Prima condiție specifică de implementare constă în configurarea aplicației doar pentru limbile care trebuie discriminate, conform specificațiilor beneficiarului. Aceasta permite o viteză de răspuns mult mai mare decât a algoritmilor cu caracter general.

A doua condiție specifică constă în selecția metodei de discriminare în funcție de intervalul Unicode caracteristic limbilor specificate. Dacă cele două limbi folosesc intervale Unicode diferite atunci acesta poate fi singurul criteriu de discriminare. Numai în cazul în care cel puțin două dintre limbile specificate folosesc același interval de caractere, se va folosi un algoritm bazat pe spectrul n-gramelor specific acestor limbi.

A treia condiție specifică de implementare constă în definirea unei limbi implicite. Acest lucru este util pentru situațiile în care rezultatul analizei este incert, mai ales datorită fragmentelor de text foarte scurte, sau datorită amestecului de cuvinte din limbi diferite în interiorul aceluiași fragment de text; tot rezultat incert poate să apară și în cazul în care textul este scris într-o altă limbă decât cele pentru care algoritmul a fost configurat. În conjuncție cu inerția de schimbare a limbii, acest lucru permite obținerea unui rezultat bun, chiar și în cele mai dificile condiții.

Ultima particularitate la implementarea algoritmilor este posibilitatea de realizare a unor liste de n-gramă specifice pentru fiecare limbă analizată. În aceste liste vor fi păstrate doar n-gramele caracteristice care apar cu o pondere semnificativă în limba respectivă, de exemplu prin eliminarea n-gramelor care apar cu o frecvență relativă mai mică decât un anumit prag față de media frecvenței tuturor n-gramelor identificate în corpusul limbii respective; valoarea acestui procent rezultă în urma testelor.

2.3 Rezultatele testării versiunii alpha

Într-o lucrare recentă (Fogarassy-Neszly & Gherheș, 2014) au fost prezentați succint algoritmi și interfața de testare a versiunii alpha (proof-of-concept). Versiunea are un rol intermediar, scopul fiind recunoașterea limbii în care este scris un text. Textul este scris într-o singură limbă iar recunoașterea se face la nivelul întregului text. Testarea a fost făcută pe corpusuri de mici dimensiuni.

Procedura de testare a prevăzut mai multe sesiuni, corespunzător unor corpusuri de dimensiuni diferite, pentru fiecare din cele șase limbi utilizate: română, franceză, engleză, germană, spaniolă și italiană. Mărirea succesivă a corpusului nu a îmbunătățit sensibil rezultatele decât pentru primele două iterații (până la 6.500 de caractere).

În vederea recunoașterii au fost pregătite trei seturi a câte cinci paragrafe de dimensiuni diferite. Primele două seturi sunt în limba română (cu diacritice și fără diacritice) iar al treilea în limba engleză. Primele șase propoziții sunt scurte (trei cuvinte), următoarele două de dimensiuni medii (6-11 cuvinte) și ultima propoziție scurtă (patru cuvinte).

Ca indicator al eficacității a fost utilizată rata de succes, calculată ca număr de fragmente de text pentru care limba este corect identificată raportat la numărul total de fragmente de text testate. Rezultatele au arătat că la un text de peste 6.500 de caractere, rata de succes a fost între 67.67% și 73.33%. Rata de succes a fost mai mare pentru textul în limba română scris cu diacritice. Majoritatea eșecurilor au fost identificate în cazul propozițiilor scurte (de trei cuvinte). În toate cazurile, limba a fost recunoscută pentru propozițiile de 10-11 cuvinte.

3 Testarea versiunii beta

3.1 Componenta de testare

Prima versiune funcțională (beta 01) a componentei pentru recunoașterea dinamică a limbii permite testarea algoritmilor implementați, precum și analiza influenței diversilor parametri și factori asupra rezultatului obținut pentru diferite tipuri de texte. Componenta de recunoaștere a limbii necesită

o fază de antrenare înainte de utilizare. Metoda recomandată pentru introducerea unui corpus lingvistic semnificativ este prin selectarea unui fișier de tip document. Astfel, se pot introduce corpusuri lingvistice largi care pot caracteriza suficient de exact o limbă.

Figura 1 prezintă interfața componentei de antrenare. Corpusul lingvistic poate fi introdus direct în fereastra „Analyze text”, sau poate fi introdus sub forma unui fișier document de tip text (txt) sau Word (doc sau rtf). Pentru a se evita erorile triviale, datorită unor corpusuri lingvistice mult prea mici pentru a fi caracteristice, aplicația cere introducerea unui text de minim 1.000 de caractere. Datorită limitării interfeței, se pot introduce direct (prin tastare sau prin copiere) maxim 100.000 de caractere.

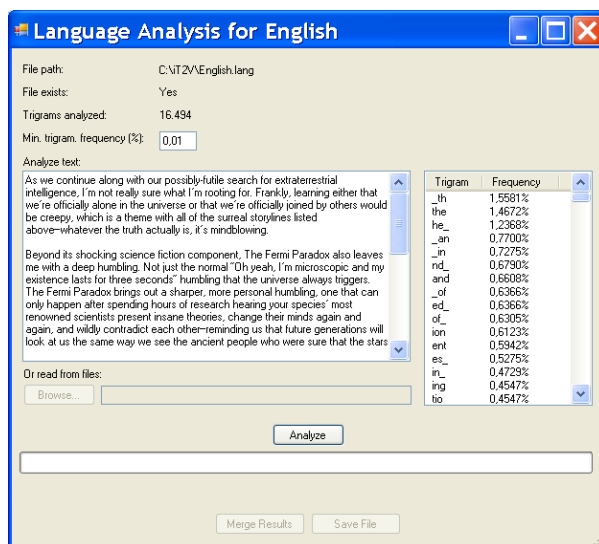


Figura 1. Interfața componentei de antrenare

Pe lângă calitatea corpusului lingvistic (volum și reprezentativitate lingvistică) și nivelul de ignorare a trigramelor nesemnificative, această componentă permite testarea influenței volumului de text analizat la un moment dat în timpul procesului de recunoaștere LA (*Look-ahead*), precum și inerția la schimbarea limbii I (*Inertia*). Un alt factor (mai puțin evident, dar extrem de important) care influențează calitatea rezultatului este numărul de limbi activate.

3.2 Strategie de testare și măsuri colectate

Fragmentul de text pentru care se face testarea cuprinde mai multe propoziții, în două limbi: română și engleză. Recunoașterea se face la nivel de cuvânt. Testarea a fost făcută din perspectiva utilizabilității aplicațiilor software. Din punctul de vedere al utilizabilității, în această etapă de dezvoltare interesează eficacitatea și eficiența componentei. Eficacitatea este definită de acuratețea și completitudinea cu care este recunoscută limba. Eficiența este definită de resursele consumate în acest scop.

Ca măsuri ale eficacității au fost colectate, pentru fiecare limbă, numărul de cuvinte (NC) și numărul de propoziții (NP) pentru care limba este corect identificată și numărul de cazuri când schimbarea limbii este corect făcută (NS). Ca indicator al eficacității a fost utilizată rata de succes, calculată la nivel de propoziție (EFP) și la nivel de cuvânt (EFC). Rata de succes la nivel de propoziție este importantă întrucât propoziția este o unitate semantică având o semnificație pentru utilizator.

Ca măsuri ale eficienței au fost utilizați cei doi parametri care pot fi modificați în testare, respectiv numărul de cuvinte analizate (LA) și inerția (I), precum și numărul de limbi candidat.

Testarea a cuprins 3 etape, având următoarele obiective:

- Analiza influenței mărimii corpusului de antrenare, analiza influenței numelor proprii și analiza variației rezultatelor în funcție de valoarea parametrilor, pentru un număr de șase limbi candidat;
- Testarea unui text din care s-au eliminat numele proprii, analiza influenței pe care o are ordinea propozițiilor și influența pe care o are variația parametrilor;
- Identificarea valorilor optime a parametrilor și confirmarea rezultatelor prin testarea cu un text preluat din lumea reală (articole din presă), având fragmente în patru limbi.

În primele două etape au fost folosite texte cu fragmente în două limbi: română și engleză. În acest scop, pentru prima etapă au fost pregătite două seturi, a câte opt propoziții fiecare. Primul set de propoziții este în limba română, iar al doilea în limba engleză. Dintre cele 16 propoziții, șase sunt scurte (patru din trei cuvinte și două din patru cuvinte), șase sunt de 5-7 cuvinte și numai patru propoziții sunt mai lungi. Primele patru propoziții

sunt în limba română și următoarele opt în engleză, astfel încât sunt două schimbări de limbă. Pentru limba română, numărul total de cuvinte este 44. Dintre acestea, șase sunt nume proprii, iar unul este neologism. Pentru limba engleză, numărul total de cuvinte este 45.

Pentru etapa a doua, textul de recunoscut a fost adaptat corespunzător, păstrând propozițiile și eliminând majoritatea numelor proprii. Ordinea a două propoziții a fost inversată.

4. Rezultate

4.1 Rezultate etapa I

Prima etapă de testare a avut un rol preliminar. Obiectivele au fost: analiza influenței mărimii corpusului de antrenare, analiza influenței numelor proprii și analiza variației rezultatelor în funcție de valoarea parametrilor, pentru un număr de șase limbi candidat.

Procedura a prevăzut zece sesiuni de testare. Primele patru, corespunzător corpusurilor de diferite dimensiuni (5, 10, 20, 30 mii de cuvinte), utilizează patru limbi (română, engleză, franceză și spaniolă). La fiecare iterație, trigramele din textul de antrenare au fost integrate cu cele existente în baza de date. Următoarele două sesiuni utilizează ultimul corpus (de dimensiune maximă), dar testează numai cu trei limbi candidat (română, engleză și spaniolă), respectiv două limbi (română și engleză). Sesiunea 7 utilizează limba italiană în locul limbii franceze, iar sesiunea 8 utilizează limba maghiară în locul limbii franceze. Următoarele două sesiuni utilizează numai trei limbi (română, engleză și spaniolă, respectiv română, engleză și maghiară). Toate cele patru sesiuni utilizează corpusul de dimensiune maximă disponibil (30 de mii de cuvinte).

În fiecare sesiune, a fost testată influența celor doi parametri, făcând câte patru teste, cu parametrii la valorile extreme. Rezultatele au fost mai bune pentru valorile mari ale parametrilor și sugerează o testare cu variații mai fine. De asemenea, rezultatele au fost mai bune la testarea cu mai puține limbi candidat.

În toate sesiunile, diferența între recunoașterea limbii române și cea a limbii engleze a fost mare. În Tabelul 1 sunt prezentate rezultatele obținute la trei sesiuni de testare: valoarea parametrilor (LA/I), eficacitatea la nivel de propoziție și cuvânt (EFP, EFC) pentru fiecare limbă în care este scris

textul (ro și en) și eficacitatea globală, ca medie a eficacității pentru cele două limbi din textul testat (EFP, EFC)

Se observă diferența dintre rezultate în funcție de numărul de limbi candidat. De asemenea, se observă influența mai mare a inerției (I) față de volumul de text analizat (LA).

Tablul 1 Rezultate testare (corpus de 30Kcuvinte)

LA/I	EFPPro %	EFPen %	EFCro %	EFCen %	EFP %	EFC %
Sesiune 4: patru limbi candidat						
1/0	62.50	0.00	75.00	42.22	31.25	37.50
15/0	62.50	0.00	75.00	46.67	31.25	37.50
1/5	87.50	12.50	93.18	64.44	50.00	52.84
15/5	100.00	12.50	100.00	28.89	56.25	56.25
Sesiune 5: trei limbi candidat						
1/0	37.50	12.50	81.82	33.33	25.00	47.16
15/0	62.50	12.50	81.82	51.11	37.50	47.16
1/5	87.50	50.00	93.18	77.78	68.75	71.59
15/5	100.00	50.00	100.00	77.78	75.00	75.00
Sesiune 6: două limbi candidat						
1/0	75.00	25.00	93.18	57.78	50.00	59.09
15/0	62.50	25.00	86.36	55.56	43.75	55.68
1/5	100.00	50.00	100.00	80.00	75.00	75.00
15/5	100.00	62.50	100.00	77.78	81.25	81.25

4.2 Rezultate etapa II

Obiectivele au fost analiza influenței numelor proprii, ordinea propozițiilor și variația parametrilor. Având în vedere că numele proprii ridică unele probleme, majoritatea au fost eliminate pentru a analiza schimbarea limbii. În ceea ce privește ordinea propozițiilor, s-a eliminat influența pe care o au unele cuvinte care sunt comune în mai multe limbi. În ceea ce privește parametrii, interesează valoarea cea mai scăzută pentru care se păstrează precizia.

În vederea testării au fost pregătite fișiere cu text de analizat pentru patru limbi: română, engleză, franceză și germană. Corpusurile utilizate au fost de dimensiuni mai mari, între 45 și 120 de mii de cuvinte. Textul de recunoscut a fost adaptat corespunzător, păstrând propozițiile și eliminând majoritatea numelor proprii. Ordinea a două propoziții a fost inversată.

Procedura de testare a prevăzut trei sesiuni a câte 24 iterații fiecare.

Prima sesiune utilizează patru limbi (ro, en, fr și de). Următoarele două sesiuni testează numai cu trei limbi candidat (ro, en și fr), respectiv două limbi (ro și en). În fiecare sesiune, a fost testată influența celor doi parametri. Au fost făcute opt teste, mărind de fiecare dată valoarea parametrilor cu o unitate. În acest fel poate fi identificat intervalul cu valoarea cea mai scăzută a parametrului pentru care s-a păstrat precizia.

Rezultatele sunt mult mai bune decât în etapa precedentă. În Tabelul 2 sunt integrate rezultatele din primele două sesiuni, după cazurile în care parametrii au valori extreme.

Tabelul 2 Rezultate testare etapa II

LA/I	EFPr %	EFPe %	EFCr %	EFCe %	EFP %	EFC %
Sesiune 1: patru limbi candidat						
5/1	37.50	37.50	76.09	81.25	37.50	56.79
10/1	37.50	50.00	78.26	83.33	43.75	64.13
5/2	75.00	62.50	91.30	89.58	68.75	76.90
10/2	75.00	75.00	93.48	93.75	75.00	84.24
5/3	75.00	62.50	89.13	85.42	68.75	75.82
10/3	75.00	62.50	91.30	91.67	68.75	76.90
5/4	75.00	100.00	89.13	100.00	87.50	94.57
10/4	75.00	100.00	89.13	100.00	87.50	94.57
Sesiune 2: trei limbi candidat						
5/1	75.00	37.50	89.13	81.25	56.25	63.32
10/1	75.00	50.00	89.13	83.33	62.50	69.57
5/2	75.00	75.00	89.13	91.67	75.00	82.07
10/2	75.00	87.50	89.13	95.83	81.25	88.32
5/3	75.00	75.00	89.13	87.50	75.00	82.07
10/3	75.00	75.00	89.13	93.75	75.00	82.07
5/4	75.00	100.00	89.13	100.00	87.50	94.57
10/4	75.00	100.00	89.13	100.00	87.50	94.57

Rezultatele sesiunii 3 nu au fost păstrate fiind mai puțin concludente, deși s-a observat același patern în rezultate. Pentru ambele sesiuni, cele mai bune rezultate se obțin pentru valori ale LA între 5 și 10 și valori ale inerției în jurul valorii de 2. La valori ale inerției de 3 se observă, în mod surprinzător, o scădere a eficacității. Este posibil să fie ca urmare a particularității cuvintelor din textul testat.

4.3 Rezultate etapa III

Obiectivele au fost identificarea valorilor cele mai scăzute ale parametrilor pentru care se păstrează precizia și confirmarea rezultatelor pe un alt text.

Pentru testare au fost utilizate aceleași corpusuri.

Pentru identificarea valorii optime a parametrilor, s-a testat cu patru limbi candidat (ro, en, fr, de) în două sesiuni. În prima sesiune au fost efectuate șapte iterații, variind parametrul LA cu pas 1, între 5 și 11 (parametrul I constant la valoarea 2). În cea de a doua sesiune au fost efectuate zece iterații variind parametrul I cu pas 0.1, între 1.5 și 2.4 (parametrul LA constant la valoarea determinată în prima iterație).

Rezultatele testării sunt prezentate în Tabelul 3.

Tabelul. Rezultate testare etapa III

LA/I	EFPr %	EFPe %	EFCr %	EFCe %	EFP %	EFC %
Lb=4						
10/1.9	75.00	75.00	89.13	93.75	75.00	82.07
10/2.0	75.00	75.00	93.48	93.75	75.00	84.24
11/1.9	87.50	75.00	95.65	93.75	81.25	85.33
11/2.0	87.50	75.00	95.65	93.75	81.25	85.33
Lb=3						
10/1.9	75.00	87.50	89.13	95.83	81.25	88.32
10/2.0	75.00	87.50	89.13	95.83	81.25	88.32
11/1.9	75.00	87.50	89.13	93.75	81.25	88.32
11/2.0	75.00	87.50	89.13	93.75	81.25	88.32
Lb=3						
10/1.9	75.00	75.00	93.48	93.75	75.00	84.24
10/2.0	75.00	87.50	93.48	97.92	81.25	90.49
11/1.9	87.50	75.00	95.65	93.75	81.25	85.33
11/2.0	87.50	75.00	95.65	97.92	81.25	85.33
Lb=2						
10/1.9	75.00	87.50	89.13	95.83	81.25	88.32
10/2.0	75.00	100.00	89.13	100.00	87.50	94.57
11/1.9	75.00	87.50	89.13	95.83	81.25	88.32
11/2.0	75.00	100.00	89.13	100.00	87.50	94.57

Cele mai bune rezultate au fost obținute pentru o valoare a parametrului LA de 11. Întrucât și pentru valoarea lui LA de 10 rezultatele sunt apropiate, următoarea sesiune a utilizat această valoare, testând apoi și pentru valoarea 11. În cea de a doua sesiune a fost variată inerția. Cele mai bune rezultate s-au obținut pentru valorile 1.9 și 2.0. La următoarele iterații s-a observat o ușoară scădere a eficacității la recunoașterea limbii române.

Se observă că la utilizarea a trei limbi candidat (primul caz este cu română, engleză, franceză și al doilea cu română, engleză, germană) rezultatele cele mai bune se obțin la valori mai mici ale parametrului LA (10). De asemenea, la utilizarea a numai două limbi candidat rezultatele cele mai bune se obțin pentru LA=10 și I=2.0

Se cuvine a fi menționat faptul că eficacitatea la nivel de propoziție este un indicator care are în mod inerent valori scăzute, deoarece presupune ca toate cuvintele dintr-o propoziție să fie corect identificate. Având în vedere existența multor cuvinte comune pentru mai multe limbi (de exemplu „cinema”) și existența numelor proprii, este greu să se obțină valori ridicate pe texte din lumea reală.

Din aceleași motive, sunt greu de identificat în mod corect toate schimbările de limbă (la majoritatea iterațiilor s-a identificat una singură).

Pentru confirmare pe un text din lumea reală, testarea s-a făcut cu un text format din patru paragrafe care conțin propoziții preluate din presă (Adevărul, Times, Le Monde și Deutche Welle). Fiecare paragraf se referă la cel puțin două domenii diferite (propozițiile au fost preluate din rubrici diferite). Paragraful în limba română conține opt propoziții și 173 cuvinte. Paragraful în limba engleză are zece propoziții și 203 cuvinte. Paragraful în limba franceză are nouă propoziții și 217 cuvinte. Paragraful în limba germană are zece propoziții și 161 cuvinte. Numărul de schimbări de limbă care ar trebui detectate este trei.

Tabelul 4 Rezultate test de confirmare

Limba	LA=10, I=2.0				LA=11, I=1.9			
	NP	NC	EFPr %	EFPe %	NP	NC	EFP %	EFC %
română	5	153	62.50	88.44	5	151	62.50	87.28
engleză	8	203	80.00	98.07	7	202	70.00	97.58
franceză	9	217	100.00	100.00	9	216	100.00	99.54
germană	2	126	20.00	78.26	3	137	30.00	85.09
			65.63	91.19			65.63	92.38

Testarea s-a făcut pentru valori ale parametrului LA de 10 și 11 și ale inerției de 1.9 și 2.0. Rezultatele testării sunt prezentate mai jos. Așa cum se observă din tabel, rezultatele au valori apropiate în cele două iterații, cu puțin mai mari pentru LA=11. În ceea ce privește recunoașterea fiecărei limbi, pentru rezultatele au fost mai bune pentru limbile franceză și engleză și mai slabe pentru limbile română și germană. Schimbarea limbii este corect detectată în două cazuri (en/de și de/fr).

O comparație cu rezultatele anterioare arată că eficacitatea la nivel de cuvânt este mai mare (92.38% față de 85.33%). Eficacitatea la nivel de propoziție este ceva mai mică (65.63% față de 81.25%). Pe ansamblu, rezultatele testării pe un text din lumea reală confirmă rezultatele testării pe un textul utilizat anterior.

Analiza detaliată arată că erorile provin în primul rând din identificarea unor cuvinte din română, engleză și germană ca fiind franceze. Este posibil ca aceasta să se datoreze dimensiunilor diferite ale corpusurilor (cel în germană este mai mic) și / sau sursei acestora.

În toate cazurile însă, numărul de cuvinte dintr-o propoziție care sunt corect identificate este mai mare decât cel al cuvintelor care nu sunt corect identificate. Aceasta sugerează investigarea posibilității de a introduce în algoritmul de recunoaștere o analiză la nivel de propoziție.

5. Concluzii și direcții de continuare

În cadrul acestui articol au fost prezentate rezultatele testării primei versiuni funcționale a unei componente software de recunoaștere dinamică a limbii cu aplicație în sinteza vocală diferențiată lingvistic.

O contribuție a articolului este evaluarea performanțelor componente din perspectiva utilizabilității. În acest scop, a fost propusă eficacitatea recunoașterii limbii la nivel de propoziție, eficacitatea la nivel de cuvânt și eficacitatea globală. Ca parametri au fost utilizați: numărul de limbi candidat, numărul de cuvinte analizate și inerția la schimbarea limbii.

Rezultatele testării arată o influență scăzută a mărimii corpusurilor (posibil și datorită utilizării unor corpusuri de dimensiuni relativ mici). Rezultatele sunt utile deoarece au permis identificarea unor valori optime a celor doi parametri care determină eficiența recunoașterii limbii (volum de text analizat și inerție la schimbarea limbii).

Cercetările vor continua prin investigarea posibilității de a introduce în algoritmul de recunoaștere o analiză la nivel de propoziție. De asemenea, se intenționează antrenarea componente cu corpusuri de dimensiuni mai mari.

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

- Athias, L. (2007). *Statistical Machine Translation and Automatic Speech Recognition under Uncertainty*. PhD thesis, Johns Hopkins University.
- Baldwin, T., Lui, M. (2010) Language Identification: The Long and the Short of the Matter, *Human Language Technologies: 2010 Annual Conference of the North American Chapter of the ACL*, 229–237.
- Cavnar, W., and Trenkle, J. (1994). N-gram-based text categorization. *Proc. 3rd Symp. on Document Analysis and Information Retrieval (SDAIR-94)*.
- Churcher, G. (1994) *Distinctive character sequences*. Personal communication.
- Devadoss, J.M. (2010) Advanced Natural Language Translation System. *Global Journal of Computer Science and Technology* 9(5), 114-122.
- Dunning, T. (1994) *Statistical Identification of Language*. Technical Report MCCS 94-273, New Mexico State University.
- Fogarassy-Neszly, P., Gherhes, V. (2013) Recunoașterea automată a limbii cu aplicație în sinteza vocală diferențiată lingvistic. *Revista Română de Interacțiune Om-Calculator* 6(2), 155-168.
- Fogarassy-Neszly, P., Gherhes, V. (2014) Aplicații pentru recunoașterea dinamică a limbii. Popovici, D.M., Iordache D.D. (Eds.). *Proc. Conferința Națională de Interacțiune Om-Calculator – RoCHI 2014*, 51-54.
- Gînscă, A. L., Boros, E., Iftene, A. (2011). Adapting Statistical Language Identification Methods for Short Queries. *Notebook Paper CLEF 2011 LABs Workshop*, 19-22 September, Amsterdam, Netherlands.
- Kłosowski, P., Dustor, A. (2013) Automatic Speech Segmentation for Automatic Speech Translation. *Computer Networks Communications in Computer and Information Science*, 370, 466-475.
- Kłosowski, P., & Dustor, A. (2013). Automatic Speech Segmentation for Automatic Speech Translation. In *Computer Networks*, 466-475.
- Ljubešić, N., Mikelić, N., Boras, D. (2007) Language identification: How to distinguish similar languages. In Lužar-Stifter, V. and Hljuz Dobrić, V. (Eds), *Proceedings of the Intl Conference on Information Technology Interfaces*, 541–546.
- Zissman, M. A. (1996) Comparison of four approaches to automatic language identification of telephone speech. *IEEE Transactions on Speech and Audio Processing* 4(1), 31-44.