

Construcția unui sistem simplu de recunoaștere a vorbirii continue

- recunoașterea secvențelor audio ce conțin cifre -

I. Introducere

Procesul recunoașterii vorbirii continue (RVC) vizează transformarea unui semnal audio ce conține vorbire într-o succesiune de cuvinte. Un sistem de recunoaștere a secvențelor audio de cifre este un sistem de recunoaștere a vorbirii continue cu vocabular redus, în sensul că singurele cuvinte ce pot fi recunoscute de către sistem sunt cele zece cifre ale sistemului zecimal: zero, unu, ..., nouă.

Arhitectura generală a unui sistem de recunoaștere a vorbirii continue

Arhitectura generală a unui sistem de recunoaștere automată a vorbirii (RAV) este prezentată în Figura 1. Din figură reies două lucruri esențiale în ceea ce privește procesul de recunoaștere a vorbirii: a) recunoașterea se face utilizând o serie de parametri vocali extrași din semnalul vocal (nu direct, utilizând mesajul vorbit) și b) recunoașterea se face pe baza unor modele (acustic, fonetic și lingvistic) dezvoltate în prealabil.

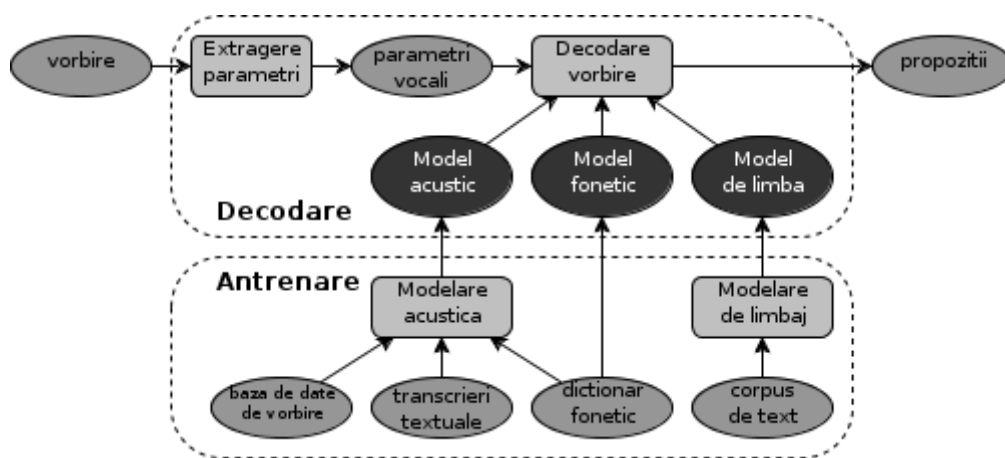


Figura 1. Arhitectura generală a unui sistem de recunoaștere a vorbirii

Modelul acustic are rolul de a estima probabilitatea unui mesaj vorbit, dată fiind o succesiune de cuvinte. În sistemele de RVC state-of-the-art modelul acustic nu folosește cuvinte ca unități acustice de bază pentru că: a) fiecare sarcină de RAV are un vocabular de cuvinte diferit pentru care nu există modele deja antrenate și nici date de antrenare disponibile și b) numărul de cuvinte diferite dintr-o limbă este prea mare. În locul cuvintelor, se utilizează unități acustice de bază sub-lexicale, de exemplu foneme, sau, mai recent, chiar unități sub-fonemice numite senone. Prin urmare, modelul acustic este format dintr-un set de modele pentru foneme (sau senone) care se conectează în timpul procesului de decodare pentru a forma modele pentru cuvinte și apoi modele pentru succesiuni de cuvinte. Acestea sunt utilizate în final pentru a estima probabilitatea ca mesajul rostit

de vorbitor să fie format dintr-o succesiune sau alta de cuvinte. Experiența a arătat că această abordare generativă pentru modelele acustice poate fi foarte bine implementată utilizând modele Markov ascunse (hidden Markov models - HMM). Pentru mai multe informații despre modul cum funcționează modelele acustice consultați bibliografia îndrumarului de proiect.

Modelul de limbă este utilizat în timpul decodării pentru a estima probabilitățile tuturor secvențelor de cuvinte din spațiul de căutare. În general, rolul unui model de limbă este de a estima probabilitatea ca o secvență de cuvinte $W = w_1, w_2, \dots, w_n$, să fie o propoziție validă a limbii. Probabilitatea acestor secvențe de cuvinte ajută foarte mult modelul acustic în procesul de decizie. De exemplu, aceste două fraze: *casa ta e mare?* și *casat a Ema re?* sunt similare din punct de vedere acustic, însă cea de-a doua nu are niciun sens. Rolul modelului de limbă este acela de a asigura o probabilitate mult mai mare primei secvențe de cuvinte și, prin urmare, de a ajuta sistemul de RAV să decidă în favoarea acesteia.

Modelul fonetic are rolul de a conecta modelul acustic (care estimează probabilitățile acustice ale *fonemelor*) cu modelul de limbă (care estimează probabilitățile secvențelor de *cuvinte*). De cele mai multe ori modelul fonetic este un dicționar de pronunție care asociază fiecărui cuvânt din vocabular una sau mai multe secvențe de foneme adecvate, reprezentând modul în care se poate pronunța respectivul cuvânt.

Figura 1 ilustrează, de asemenea, procesele implicate în dezvoltarea unui sistem de RAV, dar și resursele necesare creării modelelor acustice, lingvistice și fonetice. Modelul acustic se construiește strict pe baza unui set de clipuri audio înregistrate asociate cu transcrierea textuală a mesajelor vorbite și a unui dicționar fonetic ce cuprinde toate cuvintele din respectiva transcriere textuală. În cazul sistemelor de RVC cu vocabular extins se folosesc modele de limbă statistice, care se construiesc utilizând corpusuri de text de dimensiuni cât mai mari și cât mai adaptate domeniului din care fac parte mesajele vorbite ce trebuie decodate. În sistemele de RVC cu vocabular redus (cum este și sistemul pe care îl veți dezvolta în continuare) se folosesc preponderent modele de limbă de tip gramatică cu stări finite, iar pentru construcția acestora nu este nevoie de resurse textuale.

II. Înregistrarea unei baze de date de clipuri audio

Așa cum am precizat mai sus, pentru a construi modelul acustic al unui sistem de recunoaștere a vorbirii este nevoie în primul rând de o bază de date de clipuri audio înregistrate. Pentru proiectul de recunoaștere de secvențe audio ce conțin cifre veți înregistra un set de 100 de clipuri audio a câte 12 cifre fiecare. Pentru înregistrarea clipurilor audio veți folosi o aplicație de înregistrări accesibilă on-line la adresa <http://speed.pub.ro/speech-recorder>. Pentru a accesa această aplicație aveți nevoie de numele de utilizator și de parola care v-au fost comunicate la prima ședință de proiect.

Conectați un microfon la calculator și verificați ca volumul acestuia să fie la maxim.

Porniți browserul și încărcați pagina <http://speed.pub.ro/speech-recorder>. În cazul în care aplicația nu se încarcă va trebui să instalați plug-in-ul Java pentru browser.

După ce aplicația s-a încărcat, urmați instrucțiunile de pe primul ecran pentru a calibra microfonul și sistemul audio. Este foarte important ca înregistrările să se facă într-o cameră în care nu există surse de zgomot constant (de exemplu un ventilator sau un frigider), dar nici surse de zgomot intermitent (de exemplu un televizor, alte persoane care vorbesc, etc.). Etapa de calibrare are ca scop configurarea corectă a setărilor microfonului în condiții de liniște. În cazul în care aveți probleme la etapa de calibrare și nu reușiți să treceți de această etapă modificând volumul microfonului, vorbind mai tare/încet, utilizând posibilele setări speciale +10/20/30dB ale driverului audio, contactați îndrumătorul de proiect trimițându-i un email.

După ce etapa de calibrare a fost depășită cu succes, completați în câmpurile corespunzătoare numele de utilizator și parola ce v-au fost comunicate de îndrumătorul de proiect la prima ședință. Aplicația va afișa un al doilea ecran în care trebuie să selectați din câmpul *Select speaker* numele vorbitorului curent. În cazul în care rulați aplicația pentru prima dată va trebui să selectați *new speaker* pentru a crea un nou vorbitor (pentru persoana în cauză). Completați apoi câmpurile cerute și apăsați butonul *Save* pentru a salva datele noului vorbitor. Înainte de a putea continua va trebui să acceptați termenii și condițiile de utilizare ale aplicației apăsând butonul *I agree*.

Odată ce noul vorbitor a fost creat aplicația va reveni la ecranul de bază, iar noul vorbitor va fi activ. Selectați în acest moment, din câmpul *Phrase group*, grupul de fraze corespunzător proiectului vostru: *rodigits*. Aplicația va afișa fraza 1 a acestui grup de fraze (7779 6569 0341) și statusul ei (*not recorded yet*). În acest moment puteți înregistra această primă frază apăsând butonul *Record*, pronunțând cele 12 cifre **cât mai natural** una câte una (șapte, șapte, șapte, nouă, ..., unu) și apăsând în final butonul *Stop*. Dacă verificarea semnalului audio se încheie cu succes fraza va fi uploadată pe server, iar aplicația va afișa automat informațiile pentru fraza a doua. Continuați cu înregistrările în acest mod până când terminați de înregistrat toate cele 100 de fraze.

Ori de câte ori aveți impresia că înregistrarea curentă nu a fost făcută corect (a apărut un zgomot neașteptat, v-ați bâlbâit, ați pronunțat un cuvânt în plus sau în minus, ați ezitat, etc.) verificați corectitudinea ei apăsând butonul *Play last* și ascultând înregistrarea. În cazul în care înregistrarea nu a fost făcută corect apăsați butonul <, reveniți astfel la fraza anterioară și reînregistrați-o. După ce ați terminat de înregistrat toate frazele puteți utiliza butoanele <, >, <<, >> pentru a parcurge și reasculta toate clipurile audio, reînregistrându-le pe cele care vi se par incorecte.

După ce ați terminat de înregistrat și de verificat tot setul de fraze „rodigits” îl veți anunța pe îndrumătorul de proiect trimițându-i un email în care să specificați numele vorbitorului (așa cum l-ați creat în aplicația de înregistrări on-line) pentru care doriți extragerea fișierelor audio din baza de date. În cel mai scurt timp veți fi contactat pentru a vi se comunica locația (de pe serverul de dezvoltare) unde se află clipurile audio. Până atunci, puteți continua cu crearea dicționarului fonetic și a listei de foneme (de clipurile audio veți avea nevoie când veți ajunge la capitolul IV. Antrenarea modelului acustic).

III. Crearea dicționarului fonetic și a listei de foneme

Un dicționar fonetic este un instrument lingvistic care specifică modul în care se pronunță cuvintele unei limbi. Altfel spus, dicționarul fonetic face corespondența între forma scrisă și forma fonetică a

cuvintelor unei limbi. Forma fonetică a unui cuvânt este o succesiune de foneme (fonemul fiind unitatea de sunet fundamentală a unei limbi). Fonemele limbii române, codificarea internă utilizată în continuare și câteva exemple de cuvinte transcrise fonetic sunt prezentate în Tabelul 1.

Într-un sistem de recunoaștere a vorbirii continue dicționarul fonetic are rolul de a face legătura între modelul acustic (care modelează modul de producere a *sunetelor specifice limbii*) și modelul de limbă (care modelează modul în care se succed *cuvintele limbii*). În consecință dicționarul fonetic trebuie să conțină toate cuvintele posibile pentru o anumită sarcină de recunoaștere a vorbirii alături, bineînțeles, de transcrierile fonetice ale acestor cuvinte. Pentru mai multe informații despre rolul dicționarului fonetic în cadrul unui sistem de recunoaștere a vorbirii continue consultați bibliografia îndrumarului de proiect.

Pentru sarcina curentă (recunoașterea secvențelor audio ce conțin cifre) dicționarul fonetic trebuie să conțină numai cele zece cifre ale sistemului zecimal: zero, unu, doi, trei, patru, cinci, șase, șapte, opt și nouă. Sistemul de recunoaștere a vorbirii continue CMU Sphinx utilizează dicționare fonetice în format text, ce conțin câte un cuvânt pe fiecare linie a fișierului. Transcrierea fonetică a fiecărui cuvânt trebuie să apară pe aceeași linie. Fonemele trebuie separate de câte un spațiu. Crearea dicționarului fonetic în formatul agreeat este descrisă mai jos.

[Logați-vă la serverul de dezvoltare](#), apoi creați un director de lucru temporar:

```
mkdir phonetics
cd phonetics
```

Creați și editați un nou fișier numit [rodigits.dic](#) utilizând editorul vim:

```
vi rodigits.dic
```

Utilizând Tabelul 1 editați acest fișier astfel încât să aveți pe fiecare linie câte o cifră alături de transcrierea ei fonetică. Câteva instrucțiuni despre modul de editare al unui fișier text utilizând editorul vim sunt prezentate [aici](#). În final fișierul ar trebui să arate în felul următor:

```
zero z e r o
unu u n u
...
nouă n o 1 w a 1
```

Notă: în cazul în care nu puteți tasta caractere cu diacritice utilizând sistemul de operare și tastatura calculatorului pe care lucrați, puteți copia denumirile cifrelor ce conțin diacritice din lista următoare: șase, șapte, nouă. În consolă se poate da paste utilizând combinația de taste ctrl+alt+V.

Fonemul			Exemple de cuvinte	
Nr.	Simbol IPA	Simbol intern	Forma scrisă	Forma fonetică
1	a	a	sat	s a t
2	ə	a1	gură	g u r a1
3	e	e	mare	m a r e
4	i	i	lift	l i f t
5	j	i1	tari	t a r i1
6	ɨ	i2	între	i2 n t r e
7	o	o	loc	l o c
8	u	u	șut	s1 u t
9	y	y	ecru	e c r y
10	ø	o2	bleu	b l o2
11	ɛ̃	e1	deal	d e1 a l
12	j	i3	fiară	f i3 a r a1
13	ɔ̃	o1	oase	o1 a s e
14	w	w	sau	s a w
15	c	k2	chem	k2 e m
16	b	b	bar	b a r
17	p	p	par	p a r
18	k	k	acum	a k u m
19	tʃ	k1	cenușă	k1 e n u s1 a1
20	g	g	galben	g a l b e n
21	ɕ	g1	girafă	g1 i r a f a1
22	ʒ	g2	unghi	u n g2
23	d	d	dar	d a r
24	t	t	tot	t o t
25	f	f	fața	f a t1 a
26	v	v	vapor	v a p o r
27	h	h	harta	h a r t a
28	ʒ	j	ajutor	a j u t o r
29	ʃ	s1	coș	k o s1
30	l	l	lac	l a c
31	m	m	măr	m a1 r
32	n	n	nas	n a s
33	s	s	sare	s a r e
34	z	z	zar	z a r
35	r	r	risc	r i s k
36	ts	t1	țaran	t1 a1 r a n

Tabelul 1. Fonemele limbii române

În cazul sistemelor de recunoaștere a vorbirii cu vocabular redus (cum este și cel pe care tocmai îl construiți) manualul CMU Sphinx recomandă utilizarea de foneme specifice fiecărui cuvânt. Altfel spus ni se recomandă să modelăm diferit un același fonem X în funcție de cuvântul (și de poziția în cuvânt) în care apare acesta. În consecință veți modifica în continuare fișierul `rodigits.dic` adăugând în numele fiecărui fonem numele cuvântului din care face parte, cât și poziția fonemului în cuvânt (numai acolo unde este cazul). După această operație dicționarul fonetic ar trebui să arate astfel:

```
zero z_zero e_zero r_zero o_zero
unu u_unu1 n_unu u_unu2
...
nouă n_nouă o1_nouă w_nouă a1_nouă
```

Notă: în cazul în care nu puteți tasta caractere cu diacritice utilizând sistemul de operare și tastatura calculatorului pe care lucrați, puteți copia denumirile cifrelor ce conțin diacritice din lista următoare: șase, șapte, nouă. În consolă se poate da paste utilizând combinația de taste ctrl+alt+V.

În continuare trebuie să creăm lista tuturor fonemelor ce vor fi modelate de modelul acustic. Este necesar să scriem un fișier text care să conțină câte un fonem pe fiecare linie, iar liniile trebuie să fie sortate în ordine alfabetică. Crearea acestui fișier se poate face manual (parcurgând dicționarul fonetic creat anterior) sau automat (urmând instrucțiunile de mai jos).

Creați un nou fișier (`rodigits.phones.temp`) care să conțină toate cuvintele și toate fonemele din dicționarul fonetic (câte un cuvânt/fonem pe linie):

```
sed 's/ /\n/g' rodigits.dic > rodigits.phones.temp
```

Din această listă selectați numai fonemele (profitând de faptul că ele conțin caracterul “_”) și listați-le într-un nou fișier numit `rodigits.phones.temp.onlyPhones`:

```
grep "_" rodigits.phones.temp > rodigits.phones.temp.onlyPhones
```

Adăugați la sfârșitul acestui fișier “fonemul” SIL (acest “fonem” este utilizat pentru a modela zonele de liniște):

```
echo "SIL" >> rodigits.phones.temp.onlyPhones
```

Sortați fonemele în ordine alfabetică, redenumiți fișierul rezultat și ștergeți toate celelalte fișiere temporare:

```
sort -u rodigits.phones.temp.onlyPhones > rodigits.phones.temp.onlyPhones.sorted
mv rodigits.phones.temp.onlyPhones.sorted rodigits.phones
rm rodigits.phones.temp*
```

IV. Antrenarea modelului acustic

Antrenarea modelului acustic se va face în cadrul unui proiect CMU Sphinx și, așa cum am precizat în prima secțiune, presupune existența următoarelor resurse:

- un set de clipuri audio ce conțin vorbire (set pe care l-ati înregistrat în prealabil),

- transcrierea textuală corespunzătoare cuvintelor pronunțate în clipurile audio (deja existentă),
- un dicționar fonetic cuprinzând toate cuvintele (dicționar pe care l-ați creat în prealabil),
- un dicționar cu elemente acustice care nu sunt foneme (liniște, tuse, râs, muzică, etc.);

Crearea unui proiect CMU Sphinx numit rodigits

Logați-vă pe serverul de dezvoltare și creați un director numit `projects` în care veți crea ulterior toate proiectele:

```
cd ~
mkdir projects
```

Intrați în directorul `projects` și creați un nou director numit `rodigits` pentru proiectul curent. Apoi intrați în directorul `rodigits`:

```
cd ~/projects
mkdir rodigits
cd ~/projects/rodigits
```

Configurați directorul `rodigits` pentru a putea fi utilizat ca proiect CMU Sphinx:

```
sphinxtrain_biosinf -t rodigits setup
mkdir logdir
```

Adăugarea resurselor necesare în directorul proiectului

Logați-vă pe serverul de dezvoltare și intrați în directorul proiectului:

```
cd ~/projects/rodigits
```

Copiați clipurile audio înregistrate în prealabil în directorul `wav` (important: modificați comanda următoare pentru a include propriul `SPEAKER_ID`, așa cum vi l-a specificat îndrumătorul de proiect):

```
mkdir wav
cp -r /home/biosinfShare/resources/wav/SPEAKER_ID/ ~/projects/rodigits/wav/
```

Copiați transcrierea textuală a cuvintelor pronunțate în clipurile audio în directorul `etc`:

```
cp /home/biosinfShare/resources/transcription/rodigits.data ~/projects/rodigits/etc/
```

Copiați dicționarul fonetic și lista de foneme create în prealabil în directorul `etc`:

```
cp ~/phonetics/rodigits.dic ~/projects/rodigits/etc/
cp ~/phonetics/rodigits.phones ~/projects/rodigits/etc/
```

Copiați dicționarul cu elemente acustice care nu sunt foneme în directorul `etc`:

```
cp /home/biosinfShare/resources/phonetic/rodigits.filler ~/projects/rodigits/etc/
```

Copiați scripturile `createFileIds.sh` și `createTranscriptions.sh` în directorul `etc`:

```
cp -r /home/biosinfShare/scripts/createFileIds.sh ~/projects/rodigits/etc/  
cp -r /home/biosinfShare/scripts/createTranscriptions.sh ~/projects/rodigits/etc/
```

Transformarea fișierelor necesare procesului de antrenare în formatul corespunzător

Înainte de a putea începe antrenarea modelului acustic, trebuie creat un fișier cu lista clipurilor audio ce vor fi folosite la antrenare și un alt fișier cu transcrierea textuală a respectivelor clipuri audio. Din motive de eficiență vom crea tot acum și lista clipurilor audio ce vor fi folosite la evaluare, cât și fișierul cu transcrierea textuală a acestora.

Logați-vă pe serverul de dezvoltare și intrați în directorul proiectului, în subdirectorul `etc`:

```
cd ~/projects/rodigits/etc/
```

Vizualizați scriptul `createFileIds.sh` utilizând editorul vim:

```
vi createFileIds.sh
```

Intrați în modul de editare (tastând “i”) și schimbați valoarea variabilei `SPEAKER_ID` pentru a utiliza propriul `SPEAKER_ID`, așa cum vi l-a specificat îndrumătorul de proiect. Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și părăsiți editorul (tastând “:q” și apoi Enter).

Creați lista tuturor clipurilor audio înregistrate:

```
./createFileIds.sh > rodigits.fileids.all
```

Creați lista clipurilor audio ce vor fi utilizate pentru antrenare selectând primele 50 și ultimele 30 de clipuri audio din lista tuturor clipurilor (`rodigits.fileids.all`):

```
head -50 rodigits.fileids.all > rodigits.fileids.train  
tail -30 rodigits.fileids.all >> rodigits.fileids.train
```

Celelalte clipuri audio vor fi utilizate pentru evaluare. Creați un fișier cu lista acestora:

```
diff rodigits.fileids.train rodigits.fileids.all | grep '>' | sed 's/> //' > rodigits.fileids.test
```

Vizualizați scriptul `createTranscriptions.sh` utilizând editorul vim:

```
vi createTranscriptions.sh
```

Intrați în modul de editare (tastând “i”) și schimbați valoarea variabilei `SPEAKER_ID` pentru a utiliza propriul `SPEAKER_ID`, așa cum vi l-a specificat îndrumătorul de proiect. Ieșiți din modul de editare

(tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și părăsiți editorul (tastând “:q” și apoi Enter).

Creați transcrierea textuală (în formatul corespunzător cerut de CMU Sphinx) pentru toate clipurile audio:

```
./createTranscriptions.sh > rodigits.transcription.all
```

Creați un fișier care să conțină numai transcrierea textuală pentru clipurile audio ce vor utilizate pentru antrenare selectând primele 50 și ultimele 30 de linii din transcrierea textuală:

```
head -50 rodigits.transcription.all > rodigits.transcription.train  
tail -30 rodigits.transcription.all >> rodigits.transcription.train
```

Celelalte linii din transcrierea textuală vor fi utilizate pentru evaluare. Creați un fișier cu lista acestora:

```
diff rodigits.transcription.train rodigits.transcription.all | grep '>' | sed 's/> //' >  
rodigits.transcription.test
```

Configurarea procesului de antrenare

Orice proiect de recunoaștere de vorbire CMU Sphinx are un fișier de configurare se specifică parametrii modelului acustic ce va fi antrenat și locația resurselor necesare în procesul de antrenare. Pentru primul proiect pe care îl veți face nu veți pastra valorile implicite ale parametrilor modelului acustic, specificând numai locația resurselor.

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului, în subdirectorul `etc` și vizualizați fișierul de configurare al proiectului cu ajutorul editorului vim:

```
cd ~/projects/rodigits/etc/  
vi sphinx_train.cfg
```

Identificați în acest fișier liniile care specifică locația resurselor utilizate în procesul de antrenare (dicționarul fonetic, lista de foneme, dicționarul de sunete care nu sunt foneme, lista de clipuri audio de antrenare și transcrierea textuală a acestor clipuri audio):

```
$CFG_DICTIONARY = "$CFG_LIST_DIR/$CFG_DB_NAME.dic";  
$CFG_RAWPHONEFILE = "$CFG_LIST_DIR/$CFG_DB_NAME.phone";  
$CFG_FILLERDICT = "$CFG_LIST_DIR/$CFG_DB_NAME.filler";  
$CFG_LISTOFFILES = "$CFG_LIST_DIR/${CFG_DB_NAME}_train.fileids";  
$CFG_TRANSCRIPTFILE = "$CFG_LIST_DIR/${CFG_DB_NAME}_train.transcription";
```

Intrați în modul de editare (tastând “i”) și modificați aceste linii astfel:

```
$CFG_DICTIONARY = "$CFG_LIST_DIR/rodigits.dic";
$CFG_RAWPHONEFILE = "$CFG_LIST_DIR/rodigits.phones";
$CFG_FILLERDICT = "$CFG_LIST_DIR/rodigits.filler";
$CFG_LISTOFFILES = "$CFG_LIST_DIR/rodigits.fileids.train";
$CFG_TRANSCRIPTFILE = "$CFG_LIST_DIR/rodigits.transcription.train";
```

Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și apoi închideți editorul (tastând “:q” și apoi Enter).

Crearea fișierelor cu coeficienți MFCC corespunzătoare clipurilor audio de antrenare

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului și creați fișierele cu coeficienți cepstrali executând următoarea comandă:

```
cd ~/projects/rodigits/
/usr/local/sphinx/lib/sphinxtrain/scripts/000.comp_feat/slave_feat.pl
```

Notă: este normal ca în urma execuției comenzii de mai sus să apară eroarea [Failed to open control ... : No such file or directory at .../make_feats.pl line 89](#). Ea ne atenționează asupra faptului că lista de fișiere audio de evaluare nu a fost găsită (acest lucru este normal, pentru că această listă nu a fost încă generată).

Execuția acestei comenzi are ca efect crearea a 80 de fișiere cu coeficienți cepstrali corepunzătoare celor 80 de clipuri audio de antrenare. Verificați faptul că fișierele cu coeficienți cepstrali au fost create corect (important: modificați comenzile următoare pentru a include propriul [SPEAKER_ID](#), așa cum vi l-a specificat îndrumătorul de proiect):

```
ls feat/SPEAKER_ID/*
ls feat/SPEAKER_ID/* | wc -l
ls wav/SPEAKER_ID/* | wc -l
```

Prima comandă din grupul celor de mai sus va afișa toate fișierele din subdirectorul [feat/SPEAKER_ID](#). A doua comandă va afișa numărul de fișiere din respectivul subdirector (rezultatul ar trebui să fie 80 pentru că doar 80 de clipuri audio se folosesc pentru antrenare). A treia comandă va afișa numărul de fișiere din subdirectorul [wav/SPEAKER_ID](#) (rezultatul ar trebui să fie 100 pentru că în total ar trebui să aveți 100 de clipuri audio înregistrate).

Antrenarea modelului acustic

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului și porniți antrenarea modelului acustic executând următoarea comandă:

```
cd ~/projects/rodigits/
sphinxtrain_biosinf run
```

Așteptați terminarea procesului de antrenare (ar trebui să dureze aproximativ 2 minute). În cazul în care apar erori verificați faptul că fișierele cu resurse sunt în formatul corespunzător și că nu ați

modificat greșit fișierul de configurare. În cazul în care nu reușiți să identificați și să rezolvați problema discutați cu îndrumătorul de proiect.

V. Crearea modelului de limbă (gramaticii)

Sistemele de recunoaștere de vorbire cu vocabular mare de ultimă generație utilizează modele de limbă statistice de tip n-gram. Aceste modele de limbă se construiesc pe baza unor corpusuri mari de text specific sarcinii de recunoaștere, estimând probabilitățile de apariție ale cuvintelor și secvențelor de cuvinte specifice respectivei sarcini. Modelele de limbă de tip n-gram se utilizează apoi în procesul de decodare (recunoaștere a vorbirii) pentru a calcula probabilitățile secvențelor de cuvinte propuse de modelul acustic. Pentru mai multe informații despre modele de limbă statistice și despre rolul modelului de limbă în cadrul unui sistem de recunoaștere a vorbirii continue consultați bibliografia îndrumarului de proiect.

Sarcina de recunoaștere a secvențelor audio ce conțin cifre este o sarcină de recunoaștere cu vocabular foarte redus pentru care nu se pretează modelele de limbă statistice. Mai mult, cifrele și succesiunile de cifre din cadrul clipurilor audio înregistrate apar cu probabilități aproximativ egale (nu se poate afirma că o cifră este pronunțată sistematic mai des decât alta). În aceste condiții, modelele de limbă de tip gramatică cu stări finite (FSG – finite state grammar) sunt mult mai potrivite. O gramatică cu stări finite este un model de tip graf în care nodurile reprezintă cuvinte ale limbii, iar tranzițiile între cuvinte sunt reprezentate de arcele grafului. Un astfel de model de limbă specifică în mod explicit toate secvențele de cuvinte permise de gramatica sarcinii de recunoaștere. Mai mult, fiecărui arc în parte îi poate fi asignat un cost specificând astfel probabilitatea ca un cuvânt să fie precedat de un altul (sau altfel spus probabilitatea respectivei secvențe de două cuvinte).

În Figura 2 este reprezentată grafic gramatica cu stări finite a sarcinii noastre de recunoaștere. Se poate observa că modelul este format din 14 noduri, dintre care numai zece reprezintă cuvinte ale limbii (cifrele de la zero la nouă), iar celelalte patru sunt utilizate pentru a realiza „intrarea”, respectiv „ieșirea” din graf și pentru tranziția înapoi. Tranzițiile arată modul în care se poate parcurge această gramatică și implicit secvențele de cuvinte permise: în fiecare clip audio se pot pronunța una sau mai multe cifre.

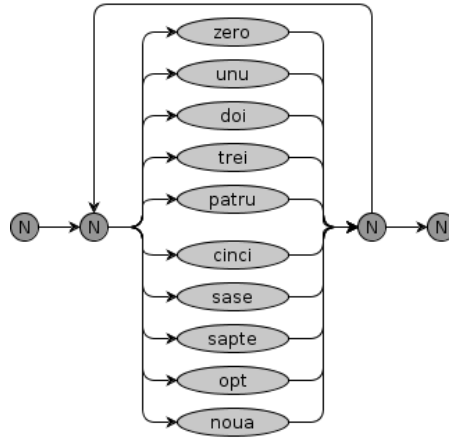


Figura 2. Gramatica cu stări finite a sarcinii de recunoaștere rodigits

Implementarea unei astfel de gramatici FSG poate fi făcută foarte simplu utilizând formatul Java Speech Grammar (JSGF) așa cum vom vedea în continuare.

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului, în subdirectorul [etc](#) și creați fișierul [rodigits.jsgf](#) utilizând editorul vim:

```
cd ~/projects/rodigits/etc/
vi rodigits.jsgf
```

Intrați în modul de editare (tastând “i”) și scrieți următoarele trei linii:

```
#JSGF V1.0;
grammar rodigits;
public <numbers> = (zero | unu | doi | trei | patru | cinci | șase | șapte | opt | nouă) * ;
```

Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și apoi închideți editorul (tastând “:q” și apoi Enter). Prima linie pe care ați scris-o în fișier specifică numele gramaticii, iar cea de-a doua linie specifică cuvintele și succesiunile de cuvinte permise astfel: [numbers](#) este *zero* sau *unu* sau ... *nouă* de oricâte ori. Sau-ul este specificat prin caracterul |, iar faptul că cifrele se pot repeta de oricâte ori este specificat de caracterul *. Pentru mai multe detalii despre formatul Java Speech Grammar și modul în care puteți crea și alte tipuri de gramatici cu stări finite consultați site-ul JSGF precizat în bibliografia îndrumarului de proiect.

Transformați (cu ajutorul utilitarului [sphinx_jsgf2fsg](#)) gramatica sarcinii de recunoaștere din format JSGF (Java Speech Grammar Format) în formatul intern FSG (Finite State Grammar) utilizat de CMU Sphinx:

```
sphinx_jsgf2fsg -jsgf rodigits.jsgf -fsg rodigits.fsg
```

Vizualizați fișierul rezultat ([rodigits.fsg](#)) utilizând editorul vim:

```
vi rodigits.fsg
```

Observați modul în care se specifică nodurile, tranzițiile și probabilitățile de tranziție ale gramaticii cu stări finite în formatul intern utilizat CMU Sphinx. Ca exercițiu, desenați pe o bucată de hârtie graful definit în acest fișier și comparați-l cu graful prezentat în Figura 2.

VI. Evaluarea sistemului de recunoaștere și interpretarea rezultatelor

Cele trei componente fundamentale ale unui sistem de recunoaștere a vorbirii (modelul acustic, modelul de limbă și modelul fonetic) sunt în acest moment disponibile. În consecință se poate trece la decodarea clipurilor audio de evaluare și apoi la compararea transcrierii textuale rezultate cu transcrierea textuală de referință în vederea evaluării sistemului de recunoaștere de secvențe audio ce conțin cifre.

Configurarea procesului de decodare (și evaluare)

Pentru a configura procesul de decodare vom edita același fișier de configurare care a fost modificat și pentru antrenare ([sphinx_train.cfg](#)). De această dată nu ne vom concentra asupra parametrilor de antrenare, ci asupra parametrilor de decodare.

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului, în subdirectorul [etc](#) și vizualizați fișierul de configurare al proiectului cu ajutorul editorului vim:

```
cd ~/projects/rodigits/etc/  
vi sphinx_train.cfg
```

Identificați în acest fișier linia care specifică numele scriptului utilizat pentru decodare:

```
$DEC_CFG_SCRIPT = 'psdecode.pl';
```

Intrați în modul de editare (tastând “i”) și modificați această linie astfel:

```
$DEC_CFG_SCRIPT = 'psdecode_fsg.pl';
```

Identificați apoi liniile care specifică numele și locația resurselor utilizate în procesul de decodare (dicționarul fonetic, dicționarul de sunete care nu sunt foneme, lista de clipuri audio de evaluare și transcrierea textuală a acestor clipuri audio) și numele subdirectorului în care se vor salva rezultatele:

```
$DEC_CFG_DICTIONARY = "$DEC_CFG_BASE_DIR/etc/$DEC_CFG_DB_NAME.dic";  
$DEC_CFG_FILLERDICT = "$DEC_CFG_BASE_DIR/etc/$DEC_CFG_DB_NAME.filler";  
$DEC_CFG_LISTOFFILES = "$DEC_CFG_BASE_DIR/etc/${DEC_CFG_DB_NAME}_test.fileids";  
$DEC_CFG_TRANSCRIPTFILE =  
"$DEC_CFG_BASE_DIR/etc/${DEC_CFG_DB_NAME}_test.transcription";  
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result";
```

Modificați aceste linii astfel:

```
$DEC_CFG_DICTIONARY = "$DEC_CFG_BASE_DIR/etc/rodigits.dic";
$DEC_CFG_FILLERDICT = "$DEC_CFG_BASE_DIR/etc/rodigits.filler";
$DEC_CFG_LISTOFFILES = "$DEC_CFG_BASE_DIR/etc/rodigits.fileids.test";
$DEC_CFG_TRANSCRIPTFILE = "$DEC_CFG_BASE_DIR/etc/rodigits.transcription.test";
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result.cd_cont_200_8";
```

Identificați apoi linia care specifică numele modelului de limbă ce va fi utilizat la decodare:

```
$DEC_CFG_LANGUAGEMODEL =
"$DEC_CFG_LANGUAGEMODEL_DIR/${DEC_CFG_DB_NAME}.lm.DMP";
```

Și modificați-o astfel:

```
$DEC_CFG_LANGUAGEMODEL = "$DEC_CFG_LANGUAGEMODEL_DIR/rodigits.fsg";
```

Identificați apoi linia care specifică aplicația cu care se va face alinierea textului rezultat în urma recunoașterii cu textul de referință (în vederea evaluării ratei de eroare la nivel de cuvânt):

```
$DEC_CFG_ALIGN = "builtin";
```

Și modificați-o astfel:

```
$DEC_CFG_ALIGN = "sclite";
```

Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și apoi închideți editorul (tastând “:q” și apoi Enter).

Crearea fișierelor cu coeficienți MFCC corespunzătoare clipurilor audio de evaluare

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului și creați fișierele cu coeficienți cepstrali executând următoarea comandă:

```
cd ~/projects/rodigits/
/usr/local/sphinx/lib/sphinxtrain/scripts/000.comp_feat/slave_feat.pl
```

Execuția acestui script are ca efect crearea a 100 de fișiere cu coeficienți cepstrali corepunzătoare tuturor clipurilor audio înregistrate (atât cele utilizate pentru antrenare, cât și cele utilizate pentru evaluare). Verificați faptul că fișierele cu coeficienți cepstrali au fost create corect (important: modificați comenzile următoare pentru a include propriul `SPEAKER_ID`, așa cum vi l-a specificat îndrumătorul de proiect):

```
ls feat/SPEAKER_ID/*
ls feat/SPEAKER_ID/* | wc -l
ls wav/SPEAKER_ID/* | wc -l
```

Prima comandă din grupul celor de mai sus va afișa toate fișierele din subdirectorul `feat/SPEAKER_ID`. A doua comandă va afișa numărul de fișiere din respectivul subdirector (rezultatul

ar trebui să fie 100 pentru că în total sunt 100 de clipuri audio înregistrate). A treia comandă va afișa numărul de fișiere din subdirectorul `wav/SPEAKER_ID` (rezultatul ar trebui să fie tot 100).

Decodarea clipurilor audio de evaluare

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului și porniți procesul de evaluare executând următoarea comandă:

```
cd ~/projects/rodigits/  
/usr/local/sphinx/lib/sphinxtrain/scripts/decode/slave.pl
```

Așteptați terminarea procesului de decodare (ar trebui să dureze aproximativ 1 minut). În cazul în care apar erori încercați să identificați problema vizualizând fișierul de log al procesului de decodare:

```
vi logdir/decode/rodigits-1-1.log
```

În cazul în care nu reușiți să identificați și să rezolvați problema discutați cu îndrumătorul de proiect. În cazul în care procesul de decodare se încheie cu succes ar trebui ca pe ecran să vă apară un mesaj similar cu acesta:

```
MODULE: DECODE Decoding using models previously trained (decode script: psdecode_fsg.pl)  
Decoding 20 segments starting at 0 (part 1 of 1)  
0%  
Aligning results to find error rate  
SENTENCE ERROR: 50.0% (10/20) WORD ERROR RATE: 12.9% (30/240)
```

Vizualizarea și interpretarea rezultatelor evaluării

Evaluarea unui sistem de recunoaștere a vorbirii se face comparând în mod automat cele două transcrieri textuale ale clipurilor audio de evaluare: transcrierea textuală de referință (cea despre care se știe a priori că este corectă) și transcrierea textuală ipotetică (cea rezultată în urma procesului de decodare). Analiza se face frază cu frază, în două etape: a) fraza ipotetică și fraza de referință sunt aliniate printr-un algoritm care urmărește să minimizeze numărul de erori de transcriere, după care b) se numără toate erorile de recunoaștere la nivel de cuvânt (cuvinte inserate, cuvinte substituite, respectiv cuvinte șterse). Se disting astfel două criterii de performanță standard folosite în evaluarea sistemelor de recunoaștere a vorbirii continue: rata de eroare la nivel de propoziție (SER - sentence error rate) și rata de eroare la nivel de cuvânt (WER - word error rate). Rata de eroare la nivel de propoziție se calculează ca raport între numărul de propoziții corecte (propoziții ce nu conțin nicio greșeală) și numărul total de propoziții. Rata de eroare la nivel de cuvânt se calculează cu formula:

$$WER[\%] = \frac{\#Erori\ de\ inserare + \#Erori\ de\ substituție + \#Erori\ de\ ștergere}{\#Cuvinte\ în\ transcrierea\ de\ referință} \times 100$$

Așa cum se observă din mesajul afișat în urma execuției procesului de decodare (dat ca și exemplu mai sus), în urma evaluării transcrierii ipotetice s-a obținut o rată de eroare la nivel de propoziție

de 50% și o rată de eroare la nivel de cuvânt de 12.9%. Mai multe informații despre rezultatele evaluării sistemului de recunoaștere puteți afla vizualizând întreg raportul de evaluare:

```
cd ~/projects/rodigits/  
vi result.cd_cont_200_8/rodigits.align
```

Primul tabel din raportul de evaluare afișează pe câte o linie sumarul rezultatelor pentru diverșii vorbitori din baza de date de evaluare. În cazul de față aveți un singur vorbitor în baza de date, deci o singură linie în tabel. Pe coloane aveți mai multe informații: id-ul vorbitorului, numărul de propoziții de evaluare, numărul de cuvinte de evaluare, procentul de cuvinte recunoscute corecte, procentul de cuvinte substituie/șterse/inserate, rata de eroare la nivel de cuvânt și rata de eroare la nivel de propoziție.

La secțiunea *Sentence Recognition Performance* se oferă mai multe detalii cu privire la erorile la nivel de propoziție, iar la secțiunea *Word Recognition Performance* se oferă mai multe detalii cu privire la erorile la nivel de cuvânt.

Urmează alte cinci secțiuni foarte interesante care prezintă statistici despre perechile de cuvinte confundate cel mai frecvent (*Confusion Pairs*), cuvintele inserate cel mai frecvent (*Insertions*), cuvintele șterse cel mai frecvent (*Deletions*), cuvintele substituie cel mai frecvent (*Substitutions*) și cuvintele recunoscute fals cel mai frecvent (*Falsely Recognized*).

În final, raportul de evaluare prezintă toate frazele ipotetice aliniate la frazele de referință corespunzătoare specificând, de asemenea, și tipurile de erori apărute la fiecare frază în parte.

Toate aceste informații și rezultate din raportul de evaluare al sistemului de recunoaștere vor fi comparate cu rezultatele ce vor fi obținute în secțiunea următoare și vor fi în final prezentate și discutate în raportul final al proiectului de cercetare-dezvoltare.

VII. Optimizarea modelului acustic

Înainte de a putea face câteva optimizări simple ale modelului acustic trebuie să studiați și să înțelegeți câteva aspecte importante privind modul de construcție al unui astfel de model:

- arhitectura generală a modelului acustic (platforma HMM-GMM). Informații despre acest subiect găsiți în bibliografia îndrumarului de proiect: [3] pag. 31, [4] pag. 23, [5] pag. 8, [6] pag. 306.
- alegerea unităților de vorbire (foneme, trifoneme, senone) și dependența/independența de context. Informații despre aceste subiecte găsiți în bibliografia îndrumarului de proiect: [3] pag. 33, [5] pag. 52, [6] pag. 310.

Utilizând deja noțiunile amintite mai sus (și explicate pe larg în bibliografie), vom sumariza în continuare modul particular în care toolkitul de recunoaștere a vorbirii pe care îl folosim (CMU Sphinx) realizează antrenarea modelului acustic. În mod implicit modelul acustic creat de CMU Sphinx este construit cu unități de vorbire dependente de context de tip trifonem (fonem căruia i se precizează vecinii stânga-dreapta). Tot în mod implicit aceste trifoneme sunt implementate ca HMM-uri cu trei stări emise, fiecare având o distribuție de probabilitate de ieșire implementată cu

un GMM. Fiecare astfel de stare-model se numește *senone* și ea poate fi comună mai multor trifoneme (în funcție de similaritățile dintre ele: de exemplu trifonemele p-a-r și t-a-v pot avea *senone*-ul inițial comun). În mod implicit numărul de *senone* este configurat la 200, iar numărul de densități de probabilitate pentru fiecare GMM este configurat la 8.

În consecință modelul acustic antrenat în prealabil modelează cele câteva foneme specificate în dicționarul fonetic utilizând modele fonetice dependente de context având în total 200 de *senone*, fiecare *senone* folosind câte 8 densități Gaussiene de probabilitate pentru a modela parametri acustici de ieșire. Procesul de antrenare implementat de CMU Sphinx conduce la obținerea acestui model acustic astfel:

1. inițializarea modelelor fonetice independente de context (câte un model cu câte trei *senone* folosind câte o singură densitate Gaussiană pentru fiecare *fonem* în parte)
2. antrenarea modelelor fonetice (independente de context) rezultate la pasul anterior
3. transformarea modelelor fonetice independente de context în modele fonetice dependente de context (se obține câte un model cu câte trei *senone* folosind câte o singură densitate Gaussiană pentru fiecare *trifonem* în parte; în total modelul acustic poate avea *oricâte senone*!)
4. antrenarea modelelor fonetice (dependente de context) rezultate la pasul anterior
5. legarea stărilor (*senonelor*) modelelor fonetice dependente de context pe baza similarităților dintre ele (se obține câte un model cu câte trei *senone* folosind câte o singură densitate Gaussiană pentru fiecare trifonem în parte; în total modelul acustic va avea *200 senone*!)
6. antrenarea modelelor fonetice (dependente de context) rezultate la pasul anterior
7. dublarea numărului de densități Gaussiene în cadrul fiecărei *senone* (se obține câte un model cu câte trei *senone* folosind câte *două densități Gaussiane* pentru fiecare trifonem în parte; în total modelul acustic va avea 200 *senone*!)
8. antrenarea modelelor fonetice (dependente de context) rezultate la pasul anterior
9. dublarea numărului de densități Gaussiene în cadrul fiecărei *senone* (se obține câte un model cu câte trei *senone* folosind câte *patru densități Gaussiane* pentru fiecare trifonem în parte; în total modelul acustic va avea 200 *senone*!)
10. antrenarea modelelor fonetice (dependente de context) rezultate la pasul anterior
11. dublarea numărului de densități Gaussiene în cadrul fiecărei *senone* (se obține câte un model cu câte trei *senone* folosind câte *opt densități Gaussiane* pentru fiecare trifonem în parte; în total modelul acustic va avea 200 *senone*!)

În consecință, pentru a obține modelul acustic cu modele fonetice dependente de context având 200 de *senone* și 8 densități Gaussiene per *senone* s-a trecut și prin etapele în care modelul acustic era format din:

- modele fonetice independente de context (după pasul 2)
- modele fonetice dependente de context cu o singură densitate Gaussiană per *senone* (după pasul 6)
- modele fonetice dependente de context cu 2 densități Gaussiane per *senone* (după pasul 8)
- modele fonetice dependente de context cu 4 densități Gaussiane per *senone* (după pasul 10)

Toate aceste modele se află în directorul proiectului, în subdirectorul `model_parameters`. Puteți verifica acest lucru astfel:

```
cd ~/projects/rodigits
ls -all model_parameters
```

Lista subdirectoarelor afișate ar trebui să fie următoarea:

```
rodigits.cd_cont_200
rodigits.cd_cont_200_1
rodigits.cd_cont_200_2
rodigits.cd_cont_200_4
rodigits.cd_cont_initial
rodigits.cd_cont_untied
rodigits.ci_cont
rodigits.ci_cont_flatinitial
```

În mod implicit decodarea realizată în capitolul VI a utilizat modelele din subdirectorul `rodigits.cd_cont_200` (modelele fonetice dependente de context cu 8 densități Gaussiană per senone). În continuare ne propunem să configurăm procesul de decodare pentru a utiliza pe rând modelele din subdirectoarele:

- `rodigits.ci_cont` (modelele fonetice independente de context)
- `rodigits.cd_cont_200_1` (modelele fonetice dependente de context cu o densitate Gaussiană per senone)
- `rodigits.cd_cont_200_2` (modelele fonetice dependente de context cu 2 densități Gaussiană per senone)
- `rodigits.cd_cont_200_4` (modelele fonetice dependente de context cu 4 densități Gaussiană per senone)

Decodarea clipurilor audio de evaluare folosind modele fonetice independente de context

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului, în subdirectorul `etc` și vizualizați fișierul de configurare al proiectului cu ajutorul editorului vim:

```
cd ~/projects/rodigits/etc/
vi sphinx_train.cfg
```

Identificați în acest fișier linia care specifică modelele acustice utilizate pentru decodare:

```
# Models to use.
$DEC_CFG_MODEL_NAME = "$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}";
```

Intrați în modul de editare (tastând “i”), copiați linia respectivă și comentați una dintre copii (pentru a păstra un back-up) și modificați cealaltă copie astfel încât să specificați locația modelelor fonetice independente de context. Această zonă a fișierului ar trebui să arate în final ca mai jos:

```
# Models to use.
#$DEC_CFG_MODEL_NAME =
"$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}";
$DEC_CFG_MODEL_NAME = "$CFG_EXPTNAME.ci_${CFG_DIRLABEL}";
```

Identificați apoi linia care specifică numele subdirectorului în care se vor salva rezultatele:

```
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result.cd_cont_200_8";
```

Și modificați-o astfel încât numele subdirectorului în care se vor salva rezultatele să fie sugestiv:

```
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result.ci_cont";
```

Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și apoi închideți editorul (tastând “:q” și apoi Enter).

Intrați apoi în directorul proiectului și porniți procesul de evaluare executând următoarea comandă:

```
cd ~/projects/rodigits
/usr/local/sphinx/lib/sphinxtrain/scripts/decode/slave.pl
```

Așteptați terminarea procesului de decodare (ar trebui să dureze aproximativ 1 minut). În cazul în care apar erori încercați să identificați problema vizualizând fișierul de log al procesului de decodare:

```
vi logdir/decode/rodigits-1-1.log
```

În cazul în care nu reușiți să identificați și să rezolvați problema discutați cu îndrumătorul de proiect. În cazul în care procesul de decodare se încheie cu succes ar trebui ca pe ecran să vă apară un mesaj similar cu acesta:

```
MODULE: DECODE Decoding using models previously trained (decode script: psdecode_fsg.pl)
  Decoding 20 segments starting at 0 (part 1 of 1)
  0%
  Aligning results to find error rate
  SENTENCE ERROR: 55.0% (11/20)  WORD ERROR RATE: 6.3% (15/240)
```

Aceste rezultate, precum și informațiile și rezultatele detaliate din raportul de evaluare generat la acest pas ([result.ci_cont/rodigits.align](#)) vor fi comparate cu rezultatele obținute în prealabil și cu rezultatele ce vor fi obținute în continuare și vor fi în final prezentate și discutate în raportul final al proiectului de cercetare-dezvoltare.

Decodarea clipurilor audio de evaluare folosind modele fonetice dependente de context cu mai puține densități Gaussiene per senone

Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului, în subdirectorul [etc](#) și vizualizați fișierul de configurare al proiectului cu ajutorul editorului vim:

```
cd ~/projects/rodigits/etc/  
vi sphinx_train.cfg
```

Identificați în acest fișier linia care specifică modelele acustice utilizate pentru decodare:

```
# Models to use.  
#$DEC_CFG_MODEL_NAME =  
"$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}";  
$DEC_CFG_MODEL_NAME = "$CFG_EXPTNAME.ci_${CFG_DIRLABEL}";
```

Intrați în modul de editare (tastând "i"), ștergeți linia decommentată (cea care specifica locația modelelor fonetice independente de context), apoi copiați linia comentată (pentru a păstra un back-up) și decommentați una dintre copii. Modificați copia decommentată astfel încât să specificați locația modelelor fonetice dependente de context ce folosesc o singură densitate Gaussiană per senone. Această zonă a fișierului ar trebui să arate în final ca mai jos:

```
# Models to use.  
#$DEC_CFG_MODEL_NAME =  
"$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}";  
$DEC_CFG_MODEL_NAME =  
"$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}_1";
```

Identificați apoi linia care specifică numele subdirectorului în care se vor salva rezultatele:

```
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result.cd_cont_200_8";
```

Și modificați-o astfel încât numele subdirectorului în care se vor salva rezultatele să fie sugestiv:

```
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result.cd_cont_200_1";
```

Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând ":w" și apoi Enter) și apoi închideți editorul (tastând ":q" și apoi Enter).

Copiați definiția modelului acustic din subdirectorul ce conține modelele fonetice dependente de context ce folosesc 8 densități Gaussiane per senone ([rodigits.cd_cont_200](#)) în subdirectorul ce conține modelele fonetice dependente de context ce folosesc o singură densitate Gaussiană per senone ([rodigits.cd_cont_200_1](#)):

```
cd ~/projects/rodigits  
cp model_parameters/rodigits.cd_cont_200/mdef model_parameters/rodigits.cd_cont_200_1/
```

Porniți apoi procesul de evaluare executând următoarea comandă:

```
/usr/local/sphinx/lib/sphinxtrain/scripts/decode/slave.pl
```

Așteptați terminarea procesului de decodare (ar trebui să dureze aproximativ 1 minut). În cazul în care apar erori încercați să identificați problema vizualizând fișierul de log al procesului de decodare:

```
vi logdir/decode/rodigits-1-1.log
```

În cazul în care nu reușiți să identificați și să rezolvați problema discutați cu îndrumătorul de proiect. În cazul în care procesul de decodare se încheie cu succes ar trebui ca pe ecran să vă apară un mesaj similar cu acesta:

```
MODULE: DECODE Decoding using models previously trained (decode script: psdecode_fsg.pl)
  Decoding 20 segments starting at 0 (part 1 of 1)
  0%
  Aligning results to find error rate
  SENTENCE ERROR: 40.0% (8/20) WORD ERROR RATE: 9.6% (23/240)
```

Aceste rezultate, precum și informațiile și rezultatele detaliate din raportul de evaluare generat la acest pas ([result.cd_cont_200_1/rodigits.align](#)) vor fi comparate cu rezultatele obținute în prealabil și cu rezultatele ce vor fi obținute în continuare și vor fi în final prezentate și discutate în raportul final al proiectului de cercetare-dezvoltare.

Repetăți pașii de mai sus realizând decodarea clipurilor audio de evaluare cu modele fonetice dependente de context cu 2 și respectiv 4 densități Gaussiene per senone.

Antrenarea unui model acustic având în total 100 de senone

În continuare ne propunem să verificăm dacă numărul de 200 de senone (configurat implicit într-un proiect CMU Sphinx) este optim sau nu pentru sarcina de recunoaștere a vorbirii propusă în acest îndrumar (recunoașterea secvențelor audio ce conțin cifre).

Vom începe investigația antrenând un nou model acustic având în total 100 de senone disponibile pentru modelarea fonemelor ce formează cuvintele specifice sarcinii de recunoaștere. Pentru a realiza acest lucru ar trebui urmărit în principiu pașii descriși în capitolul IV. Însă, pentru această sarcină nu vom crea un nou proiect CMU Sphinx, în consecință nu va trebui nici să adăugăm alte fișiere de resurse și nici să transformăm fișierele cu resurse în formatul cerut de CMU Sphinx. Vom modifica numai fișierul de configurare al proiectului [rodigits](#) astfel încât numărul de senone să fie 100 (nu 200, așa cum era în mod implicit):

```
cd ~/projects/rodigits/etc/
vi sphinx_train.cfg
```

Identificați în acest fișier linia care specifică numărul de senone:

```
# Number of tied states (senones) to create in decision-tree clustering
$CFG_N_TIED_STATES = 200;
```

Intrați în modul de editare (tastând “i”) și modificați-o astfel:

```
# Number of tied states (senones) to create in decision-tree clustering
$CFG_N_TIED_STATES = 100;
```

Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și apoi închideți editorul (tastând “:q” și apoi Enter).

Intrați în directorul proiectului și porniți antrenarea modelului acustic executând următoarea comandă:

```
cd ~/projects/rodigits  
sphinxtrain_biosinf run
```

Așteptați terminarea procesului de antrenare (ar trebui să dureze aproximativ 2 minute). În cazul în care apar erori verificați faptul că fișierele cu resurse sunt în formatul corespunzător și că nu ați modificat greșit fișierul de configurare. În cazul în care nu reușiți să identificați și să rezolvați problema discutați cu îndrumătorul de proiect.

După încheierea cu succes a procesului de antrenare verificați faptul că au fost generate modelele acustice corespunzătoare listând subdirectoarele din directorul proiectului, subdirectorul `model_parameters`:

```
cd ~/projects/rodigits  
ls -all model_parameters
```

Lista subdirectoarelor afișate ar trebui să fie următoarea:

```
rodigits.cd_cont_100  
rodigits.cd_cont_100_1  
rodigits.cd_cont_100_2  
rodigits.cd_cont_100_4  
rodigits.cd_cont_200  
rodigits.cd_cont_200_1  
rodigits.cd_cont_200_2  
rodigits.cd_cont_200_4  
rodigits.cd_cont_initial  
rodigits.cd_cont_untied  
rodigits.ci_cont  
rodigits.ci_cont_flatinitial
```

Observați că, pe lângă subdirectoarele existenete anterior (cele care conțin modele acustice cu 200 de senone și număr variabil de densități Gaussiene), în această listă apar acum și subdirectoarele `rodigits.cd_cont_100*` ce conțin modele acustice cu 100 de senone și număr variabil de densități Gaussiene.

Decodarea clipurilor audio de evaluare folosind modele acustice cu 100 de senone

Urmează ca acum să evaluăm aceste 4 noi modele acustice. Vom începe cu modelul acustic cu 100 de senone și 8 densități Gaussiene per senone. Logați-vă pe serverul de dezvoltare, intrați în directorul proiectului, în subdirectorul `etc` și vizualizați fișierul de configurare al proiectului cu ajutorul editorului vim:

```
cd ~/projects/rodigits/etc/  
vi sphinx_train.cfg
```

Identificați în acest fișier linia care specifică modelele acustice utilizate pentru decodare:

```
# Models to use.  
#$DEC_CFG_MODEL_NAME =  
"$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}";  
$DEC_CFG_MODEL_NAME =  
"$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}_4";
```

Intrați în modul de editare (tastând “i”), ștergeți linia decommentată (cea care specifică locația modelelor fonetice dependente de context cu 4 densități Gaussiene), apoi decommentați linia comentată (ce specifică locația modelelor fonetice implicite). Această zonă a fișierului ar trebui să arate în final ca mai jos:

```
# Models to use.  
$DEC_CFG_MODEL_NAME = "$CFG_EXPTNAME.cd_${CFG_DIRLABEL}_${CFG_N_TIED_STATES}";
```

Identificați apoi linia care specifică numele subdirectorului în care se vor salva rezultatele:

```
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result.cd_cont_200_4";
```

Și modificați-o astfel încât numele subdirectorului în care se vor salva rezultatele să fie sugestiv:

```
$DEC_CFG_RESULT_DIR = "$DEC_CFG_BASE_DIR/result.cd_cont_100_8";
```

Ieșiți din modul de editare (tastând Esc), salvați modificările (tastând “:w” și apoi Enter) și apoi închideți editorul (tastând “:q” și apoi Enter).

Intrați apoi în directorul proiectului și porniți procesul de evaluare executând următoarea comandă:

```
cd ~/projects/rodigits  
/usr/local/sphinx/lib/sphinxtrain/scripts/decode/slave.pl
```

Așteptați terminarea procesului de decodare (ar trebui să dureze aproximativ 1 minut). În cazul în care apar erori încercați să identificați problema vizualizând fișierul de log al procesului de decodare:

```
vi logdir/decode/rodigits-1-1.log
```

În cazul în care nu reușiți să identificați și să rezolvați problema discutați cu îndrumătorul de proiect. În cazul în care procesul de decodare se încheie cu succes ar trebui ca pe ecran să vă apară un mesaj similar cu acesta:

```
MODULE: DECODE Decoding using models previously trained (decode script: psdecode_fsg.pl)
Decoding 20 segments starting at 0 (part 1 of 1)
0%
Aligning results to find error rate
SENTENCE ERROR: 5.0% (1/20) WORD ERROR RATE: 0.4% (0/240)
```

Aceste rezultate, precum și informațiile și rezultatele detaliate din raportul de evaluare generat la acest pas ([result.cd_cont_100_8/rodigits.align](#)) vor fi comparate cu rezultatele obținute în prealabil și cu rezultatele ce vor fi obținute în continuare și vor fi în final prezentate și discutate în raportul final al proiectului de cercetare-dezvoltare.

Repetăți pașii de mai sus (ținând cont și de capitolul care prezintă decodarea folosind modele cu mai puține densități Gaussiene per senone) realizând decodarea clipurilor audio de evaluare cu modele fonetice dependente de context cu 100 de senone și 1, 2 și respectiv 4 densități Gaussiene per senone.

VIII. Construcția unui sistem de recunoaștere independent de vorbitor

Acest capitol vă propune realizarea unui al doilea proiect de recunoaștere a secvențelor audio ce conțin cifre, de această dată *independent de vorbitor* (îl vom numi [rodigitsIndep](#)). Pentru a justifica necesitatea unui sistem independent de vorbitor ar trebui să evaluați sistemul [rodigits](#) creat anterior pe fișierele de evaluare ale celorlalți colegi de grupă (fișiere ce se află pe serverul de dezvoltare în directorul [/home/biosinfShare/resources/wav/](#)). În urma evaluării veți constata că rezultatele de recunoaștere a vorbirii sunt mai slabe pentru ceilalți vorbitori decât pentru voi înșivă. Explicația este simplă: modelul acustic din proiectul [rodigits](#) a fost antrenat doar cu vocea voastră și, în consecință, vă recunoaște bine numai pe voi!

Crearea unui proiect independent de vorbitor presupune reluarea tuturor etapelor prezentate în acest îndrumar, începând cu capitolul IV, cu o singură modificare: clipurile audio de antrenare și evaluare nu vor mai fi doar clipurile înregistrate de voi, ci toate clipurile înregistrate de toți colegii de grupă. Bineînțeles că această modificare poate duce și la alte modificări (de exemplu adaptarea scripturilor [createTranscriptions.sh](#) și [createFilelds.sh](#)), însă pe toate acestea le veți aborda și eventual discuta cu îndrumătorul de proiect la momentul respectiv.

IX. Redactarea raportului final al proiectului de cercetare-dezvoltare

După efectuarea experimentelor prezentate în acest îndrumar veți sumariza toate rezultatele obținute și veți face o discuție asupra lor într-un raport final aferent proiectului de cercetare-dezvoltare.

Raportul de cercetare-dezvoltare trebuie să înceapă cu o scurtă secțiune care să cuprindă noțiuni introductive cu privire la arhitectura sistemelor de recunoaștere a vorbirii, sub-sistemele implicate, etc. Această primă secțiune are rolul de a evidenția nivelul de înțelegere (asupra sistemelor de recunoaștere a vorbirii) pe care l-ați dobândit pe parcursul realizării proiectului.

O a doua secțiune a raportului de cercetare-dezvoltare va prezenta pe scurt cerința (tema) proiectului și va enumera pașii ce au fost urmați pentru a dezvolta sistemul de recunoaștere a secvențelor audio ce conțin cifre.

În final, partea cea mai consistentă a raportului trebuie să prezinte **toate** rezultatele obținute și să le discute în mod comparativ. În funcție de stadiul în care ați ajuns cu experimentele veți avea rezultate de recunoaștere pentru:

- modelul acustic dependent de vorbitor:
 - cu 200 de senone:
 - independente de context
 - dependente de context, cu 1/2/4/8 densități Gaussiene per senone
 - cu 100 de senone:
 - independente de context
 - dependente de context, cu 1/2/4/8 densități Gaussiene per senone
- modelul acustic independent de vorbitor (cu aceleași variații ale parametrilor)

Prezentarea individuală a rezultatelor (pentru fiecare experiment în parte) va trebui să se facă punând accentul pe rata de eroare la nivel de cuvânt (Percent Total Error = Word Error Rate), dar se pot discuta și alte metrice adiacente cum sunt rata de cuvinte inserate (Percent Insertions), rata de cuvinte substituite (Percent Substitutions), rata de cuvinte șterse (Percent Deletions), rata de eroare la nivel de propoziție (Sentences with errors), etc. De asemenea, în funcție de rezultatele obținute, pot fi prezentate perechile de cuvinte confundate foarte frecvent, cuvintele inserate/substituite/șterse cel mai frecvent și eventual 1-2 exemple de fraze (evidențiind varianta corectă a frazei și varianta obținută la ieșirea sistemului).

Compararea rezultatelor obținute în diversele experimente se va face a) numai în funcție de Word Error Rate (WER), dacă ați efectuat multe experimente sau b) în funcție de mai multe metrice, dacă aveți puține rezultate. În primul caz rezultatele vor fi prezentate folosind tabele și grafice de variație a ratei de eroare la nivel de cuvânt (WER). Rezultatele se vor discuta pe baza acestor tabele și grafice și, în final, se vor trage concluzii cu privire la dependența/independența de context, numărul optim de senone, numărul optim de densități Gaussiene, etc. pentru modelul acustic dependent de vorbitor, respectiv independent de vorbitor.

Notă: redactarea raport de cercetare-dezvoltare este obligatorie, iar transmiterea lui (prin e-mail) către îndrumătorul de proiect trebuie să se facă până pe data de **9 mai 2012**.

Notă: raportul final trebuie să fie redactat și transmis îndrumătorului de proiect indiferent de stadiul în care ați ajuns cu experimentele.

Notă: toate rezultatele prezentate în raportul de cercetare-dezvoltare vor fi însoțite de referințe către rapoartele de evaluare de pe serverul de dezvoltare (calea către fișierele corespunzătoare).

X. Anexa 1. Desfășurarea activității pe server-ul de dezvoltare

Proiectul de recunoaștere a secvențelor audio ce conțin cifre va fi realizat în întregime pe un server de dezvoltare aflat în cadrul laboratorului de cercetare SPEED. Studentul va trebui să dispună de un calculator conectat la Internet care va fi folosit pentru conectarea de la distanță (prin intermediul Internetului) la serverul de dezvoltare.

Serverul de dezvoltare utilizează un sistem de operare de tip unix (Ubuntu 11.10). În consecință studentul trebuie să fie familiarizat sau să se familiarizeze cu modul de operare pe un astfel de sistem. În continuare sunt date câteva informații despre modul de conectare la server și despre utilizarea editorului vim.

Conectarea la serverul de dezvoltare

Conectarea la serverul de dezvoltare se face extrem de simplu în cazul în care calculatorul de care dispuneți utilizează un sistem de operare de tip unix/linux. În acest caz conectarea la server se face utilizând un terminal în care executați următoarea comandă (important: modificați comanda următoare pentru a utiliza propriul nume de utilizator `USER_NAME`, așa cum vi l-a comunicat îndrumătorul de proiect):

```
ssh -p 10022 USER_NAME@dev.speed.pub.ro
```

Serverul va raspunde cerându-vă să tastați parola contului vostru. Pentru a continua va trebui să tastați parola contului vostru, așa cum v-a fost specificată îndrumătorul de proiect. În cazul conectării cu succes pe ecran vă va apărea următorul prompt (unde veți executa comenzile precizate în acest îndrumar):

```
USER_NAME@alpha:~$
```

În cazul în care calculatorul de care dispuneți utilizează un sistemul de operare Windows atunci conectarea la serverul de dezvoltare se va face cu ajutorul unei aplicații intermediare numită *putty*. Această aplicație poate fi descărcată gratuit de pe siteul: <http://putty.en.softonic.com/>. După descărcare, porniți aplicația și completați în primul ecran *Host Name:* `dev.speed.pub.ro` și *Port:* `10022`. Pentru a salva aceste informații și a nu fi nevoiți să le completați de fiecare dată completați *Saved Sessions:* `biosinfServer` și apăsați butonul *Save*. Pentru ca această aplicație să poată afișa corect caracterele diacritice, intrați în meniul *Window*, submeniul *Translation* și selectați *Remote Character Set:* `UTF-8`. Apoi reveniți la ecranul inițial, selectând meniul *Session* și apăsați butonul *Open* pentru a stabili conexiunea. Serverul vă va cere numele de utilizator și apoi parola (informații pe care vi le-a specificat îndrumătorul de proiect). În cazul conectării cu succes pe ecran vă va apărea următorul prompt (unde veți executa comenzile precizate în acest îndrumar):

```
USER_NAME@alpha:~$
```

Recomandare: primul lucru pe care ar trebui să îl faceți după ce vă conectați pe server pentru prima dată este să vă schimbați parola. Pentru aceasta executați următoarea comandă și introduceți, pe rând, parola actuală, noua parola și încă o dată noua parolă:

`passwd`

Notă: indiferent de modul de conectare la serverul de dezvoltare (terminal linux sau putty) aveți posibilitatea și de multe ori este chiar foarte util să deschideți mai multe conexiuni simultane către server.

Utilizarea editorului în mod text vim

În cadrul acestui proiect veți lucra extensiv (creare, vizualizare, editare, modificare, etc.) cu fișiere text. Cel mai simplu mod de a manipula fișierele text aflate pe un server linux este cu ajutorul editorului vim.

Pentru a vizualiza un fișier deja existent sau a crea un fișier nou (cu numele filename) tastați comanda:

`vi filename`

În acest moment vizualizați prima parte a fișierului, iar cu ajutorul tastelor cu săgeți și a tastelor PgUp/PgDn puteți vizualiza și alte zone ale fișierului. Pentru a căuta un cuvânt (sau orice alt șir de caractere) tastați `/`, apoi cuvântul respectiv și în final Enter. Pentru a căuta același cuvânt din nou tastați `/` și apoi Enter. Pentru a intra în modul de editare (în vederea modificării fișierului) tastați `i`. Acum puteți șterge și scrie în acest fișier ca în oricare alt editor de text. Pentru a ieși din modul de editare tastați Esc. Pentru a salva modificările făcute tastați `:w` și apoi Enter. Pentru a părăsi editorul tastați `:q` și apoi Enter. Pentru mai multe informații de utilizare și alte comenzi de editare, copiere, etc. consultați site-ul vim precizat în bibliografia îndrumarului de proiect [8].

XI. Bibliografie

- [1] TODO Recunoașterea vorbirii continue, <http://speed.pub.ro/biosinf/RecunoastereaVorbiriiContinue.pdf>.
- [2] Toolkitul CMU Sphinx și tutorialele asociate, <http://cmusphinx.sourceforge.net/wiki/>
- [3] H. Cucu, "Towards a speaker-independent, large-vocabulary continuous speech recognition system for Romanian," PhD Thesis, Bucharest, Romania, 2011.
- [4] A. Buzo, "Automatic Speech Recognition over Mobile Communication Networks," PhD Thesis, Bucharest, Romania, 2011.
- [5] D. Jurafsky, J. Martin, "Automatic Speech Recognition," Chapter 9 in *Speech and Language Processing: An introduction to natural language processing, computational linguistics, and speech recognition (2nd Ed.)*, Pearson Education, 2009.
- [6] S. Renals, T. Hain, "Speech Recognition," in *Handbook of Computational Linguistics and Natural Language Processing*, 2010.
- [7] Java Speech Grammar Format, <http://java.sun.com/products/java-media/speech/forDevelopers/JSGF/index.html>.
- [8] Manualul editorului de text vi, <http://www.cs.fsu.edu/general/vi/manual.html>.
- Notă:** documentele [3], [4], [5] și [6] se găsesc la <http://speed.pub.ro/pcdtv/>.