



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



MINISTERUL
EDUCAȚIEI
CERCETĂRII
ȘI SPORTULUI

OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

**Universitatea Tehnică "Gheorghe Asachi" din Iași
Facultatea de Electronică, Telecomunicații și Tehnologia
Informației**

Raport de cercetare II

**CORPUSURILE ADNOTATE DE SEMNALE
VOCALE DE UZ MEDICAL ȘI INSTRUMENTE
AFERENTE**

Doctorand:

Bioing. Alina Untu (Hulea)

Conducător științific:

Prof.dr.ing. Horia Nicolai Teodorescu, m.c. A.R.



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

Cuprins

- 1. Introducere**
- 2. Corpusuri de semnale vocale și instrumente aferente**
 - 2.1. Corpusuri comerciale existente pe site-ul LDC (Linguistic Data Consortium)
 - 2.2. Metodologii utilizate la crearea de corpusuri de semnal vocal și instrumente aferente
 - 2.2.1. Metodologii pentru crearea de corpusuri destinate recunoașterii vocale
 - 2.2.2. Metodologii pentru crearea de corpusuri destinate sintezei de voce
 - 2.2.3. Metodologie pentru crearea unui corpus de dialecte
 - 2.2.4. Metodologie pentru crearea corpusului multilingv „PFSTAR”
 - 2.2.5. Metodologie pentru crearea corpusului „CIAIR inCar”
- 3. Corpusuri de voci patologice și instrumente aferente**
 - 3.1. Baza de date de voci patologice dezvoltată de MEEI (MEEI VDDb)
 - 3.2. Metode de detecție de voci patologice care utilizează ca referință MEEI VDDb
 - 3.3. Alte corpusuri de voci patologice și instrumente aferente
- 4. Microcorpusuri și metode de detecție a vocilor patologice**
 - 4.1. Aplicații în stomatologie
- 5. Cazul SRoL (Sunetele Limbii Române) și extensie în limba franceză**
 - 5.1. Microcorpus de sunete gnatofonice în limba română
 - 5.2. Microcorpus de sunete gnatofonice în limba franceză
- 6. Discuții și concluzii**
- 7. Direcții viitoare**



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

1. Introducere

La ora actuală există o serie de consorții care colaborează în vederea construirii unei baze de cunoștințe și instrumente de limbă care să fie disponibile cercetărilor din domenii precum fonetică, acustică, informatică și medicină. Câteva dintre cele mai importante sunt CLARIN, ELRA, ELDA, LDC și TELRI.

CLARIN este o infrastructură europeană care are ca scopuri principale: stocarea de corpusuri, instrumente și resurse lexicale mono sau bilingve, pe domenii specifice sau generale; investigarea standardelor existente, adaptarea lor și propunerea de sugestii pentru schimbări; crearea unei taxonomii și investigarea formatului de codare a resurselor existente; investigarea caracteristicilor instrumentelor disponibile; stabilirea criteriului de evaluare a calității resurselor de limbă (<http://www.clarin.eu/>). Grupul CLARIN își propune să ofere servicii persistente sigure și să furnizeze acces facil la resursele de procesare a limbajului.

ELRA este Asociația Europeană de Resurse pentru Limbaj incluzând ca parte operativă **ELDA** (Agenția de evaluare și distribuție a resurselor de limbaj), care are rolul de a identifica, clasifica, colecta, valida și produce resurse de limbaj. Se ocupă de asemenea de dezvoltarea ariei științifice din acest domeniu (<http://www.elra.info/>).

LDC (Linguistic Data Consortium) este un consorțiu alcătuit din universități, companii și laboratoare de cercetare guvernamentale care crează, colectează și distribuie baze de date de text și voce, lexicoane și alte resurse comerciale în scop de cercetare și dezvoltare (<http://www ldc.upenn.edu/>).



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

2. Corpusuri de semnale vocale și instrumente aferente

2.1. Corpusuri comerciale existente pe site-ul LDC

Corpusul acustico-fonetic de vorbire continuă „TIMIT” este creat pentru a furniza înregistrări de voce pentru studii acustico-fonetice și pentru dezvoltarea și evaluarea sistemelor de recunoaștere automată. Este alcătuit din înregistrări realizate cu microfon provenite de la 630 vorbitori din opt dialecte majore de engleză americană, fiecare înregistrare conținând câte zece propoziții citite. Au fost colectate informații legate de vârstă, sex, etnie, nivel de educație, înălțime și dialect corespunzătoare fiecărui vorbitor.

Corpusul TIMIT include ortografie aliniată în timp, transcripție la nivel de fonem și cuvânt și câte un fișier .wav pe 16 biți cu frecvența de eșantionare de 16 kHz pentru fiecare pronunție. A fost dezvoltat de către Massachusetts Institute of Technology (MIT), SRI International (SRI) și Texas Instruments, Inc. (TI) și este disponibil pe CD-ROM (J. S. Garofolo et al., 1986).

Corpusul „**Articulation Index Corpus**,” (J. Wright, 2005) a fost dezvoltat în scopul evaluării gradului de percepție corectă a silabelor în mediu zgomotos de către vorbitori. Alte aplicații sunt identificarea, modelarea și procesarea de limbă și modelarea pronunției. Corpusul conține înregistrări de limbă engleză la microfon, eșantionate cu 16 kHz în format pcm. Fiecare vorbitor a pronunțat silabe ale limbii engleze dintre care unele sunt cuvinte altele sunt silabe nonsens. Scopul a fost ca fiecare vorbitor să pronunțe un set de 2.000 de silabe comune tuturor vorbitorilor și 20 de silabe unice fiecărui vorbitor. Corpusul este disponibil pe DVD contracost.

„**CSLU: Kids' Speech Version 1.1**” (K. Shobaki et al., 2007) dezvoltat de LDC este o colecție de înregistrări în limba engleză realizate pe vorbire spontană și citită provenind de la 1100 copii de grădiniță și până în 10 ani. Fiecare copil a citit aproximativ 60 dintr-o listă de 319 cuvinte, propoziții sau litere. Fiecare pronunție de vorbire spontană conține la început o recitare a alfabetului și un monolog de un minut. Corpusul conține 1017 fișiere de 8-10 minute pentru fiecare vorbitor, digitizate pe 16 biți cu o frecvență de eșantionare de 16 kHz utilizând carduri audio Soundblaster 16 PnP. Pentru înregistrări au fost folosite căști cu microfon. Acest corpus a fost dezvoltat pentru a facilita analiza caracteristicilor vocilor copiilor la diverse stadii ale vârstei și la antrenarea și evaluarea sistemelor de recunoaștere utilizate în învățarea limbii și alte sarcini care implică copii



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

incluzînd dezvoltarea limbii în cazul copiilor fără auz. Informații cu privire la vîrstă, sex, limba vorbită și condițiile fizice care influențează vorbirea au fost de asemenea incluse în corpus.

„**Emotional Prosody Speech and Transcripts**” (M. Liberman et al., 2002) este un corpus care conține înregistrări audio în limba engleză cu microfon și transcripții corespunzătoare colectate pe o perioadă de opt luni între anii 2000-2001 pe două canale și cu rată de eșantionare de 22050 Hz. Înregistrările aparțin unor actori profesioniști care au citit o serie de enunțuri neutre din punct de vedere semnatic, acoperind paisprezece categorii de emoții. Corpusul are aplicații în modelare de pronunție, prozodie și recunoaștere vocală.

NTIMIT (Jankowski et al., 1990) este un alt corpus de limbă engleză cu aplicații în recunoaștere vocală care a fost dezvoltat de către Grupul NYNEX Science and Technology Speech Communication Group și este disponibil pe DVD. Acesta a fost realizat pentru a furniza înregistrări telefonice asociate corpusului TIMIT. NTIMIT a fost creat prin transmiterea a 6300 înregistrări originale printr-un telefon pe diverse canale către rețeaua telefonică NYNEX care au fost redigitizate. Înregistrările au fost transmise prin zece arii de transport și acces local dintre care pentru jumătate din acestea era necesară utilizarea purtătorilor de distanță mare. Formele de undă reînregistrate au fost aliniate în timp cu formele de undă originale ale corpusului TIMIT astfel încât transcripțiile TIMIT să poate fi utilizate concomitent cu cele NTIMIT.

Alte corpusuri existente în catalogul LDC care au fost create din înregistrări de conversații telefonice cu aplicații în recunoaștere vocală sunt:

- „**CALLHOME American English Lexicon (PRONLEX)**” (P. Kingsbury et al., 1997), conține 90,988 cuvinte selectate din texte aparținând jurnalului „Wall Street Journal” utilizate recent în corpusul ARPA de recunoaștere a vorbirii continue.
- „**CALLHOME German Lexicon**” (K. Karins et al., 1997), conține 318,807 cuvinte din care 315,503 sunt adaptate din lexiconul german CELEX produs de Centrul pentru informație lexicală, Max Planck Institute for Psycholinguistics in Nijmegen și 3,304 sunt cuvinte adiționale provenite de la 80 de transcripții de antrenare și 20 de dezvoltare a câte zece minute fiecare din cadrul corpusului LDC German CALLHOME. Lexiconul german conține câmpuri separate cu informații despre ortografie, morfologie, fonologie, accent, sursă și frecvență pentru fiecare cuvânt.



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

- **„CALLHOME Spanish Lexicon”** (S. Garrett et al., 1996) conține 45,582 cuvinte în limba spaniolă și câmpuri separate cu informații fonologice, morfologice și despre frecvență pentru fiecare cuvânt.

În mod similar au fost dezvoltate corpusuri pentru limba japoneză (M. Kobayashi et al., 1996), coreană (N. Han et al., 2003), egipteană (Kilany et al., 2002) și mandarină (S. Huang et al., 1996). Acestea sunt disponibile pe web contracost.

O altă categorie de corpusuri sunt cele create din înregistrări ale emisiunilor de radio, TV, care au ca scop detecția și urmărirea de subiect (Topic Detection and Tracking (TDT)). Printre acestea se enumeră:

- **„TDT4 Multilingual Broadcast News Speech Corpus”** conține înregistrări radio cu o frecvență de eșantionare de 16 kHz, pe două canale, în limbile arabă, engleză și mandarină și este disponibil pe DVD contracost. A fost utilizat în anii 2002-2003 pentru detectarea și urmărirea de subiect în știri. Adiacent la acest corpus au fost dezvoltate tehnici automate de găsire a materialelor cu subiect comun în știri de date provenind din știri difuzate la emisiuni radio. Tehnica de evaluare utilizată include segmentarea unor știri sursă în știri separate, urmărirea subiectelor cunoscute, detectarea subiectelor necunoscute, detectarea știrilor inițiale cu subiect necunoscut și detectarea de perechi de știri cu subiect comun (J. Kong și D. Graff, 2005).

- **„TDT3 English Audio”** conține înregistrări audio ale emisiunilor de știri colectate zilnic de la șase surse de știri în engleză americană pe o perioadă de trei luni (1998). Acestea sunt înregistrări pe un singur canal de 30 – 60 minute care au fost eșantionate cu 16 kHz utilizând eșantioane pe 16 biți și compresate cu algoritmul „shorten”. Acest corpus a fost creat ca suport pentru tehnologia TDT3 care include segmentare, detecție și urmărire, detectarea primei știri și a legăturii dintre știri (D. Graff, 2001a).

- **„TDT3 Mandarin Audio”** (D. Graff, 2001b) realizat după aceeași metodologie cu „TDT3 English Audio” conține înregistrări de știri în limba mandarină preluate de pe postul de radio „Voice of America”.

O altă categorie de corpusuri sunt destinate recunoașterii de voce în medii zgomotoase cum ar fi traficul aerian. Un astfel de corpus este **„Air Traffic Control”** (ATC0) disponibil în mai multe variante : „Air Traffic Control DFW” (J. J. Godfrey, 1994a), „Air Traffic Control BOS” (J. J. Godfrey, 1994b), „Air Traffic Control DCA” (J. J. Godfrey, 1994c), „Air Traffic Control Complete” (J. J. Godfrey, 1994d). ATC0 conține opt discuri de înregistrări de voce destinate utilizării în cercetare și activităților de dezvoltare în



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

aria recunoașterii vocale robuste, în domenii similare controlului traficului aerian (mai mulți vorbitori, canale cu zgomot, vocabular mic). Înregistrările se bazează pe comunicarea între diverse turnuri de control și piloți. Fișierele audio sunt eșantionate cu 8 kHz pe 16 biți și reprezintă monitorizarea continuă fără eliminarea zonelor cu liniște, pe o durată de la una la două ore. Sunt de asemenea disponibile fișiere care indică amplitudinea semnalului primit la intervale de 10 milisecunde. Setul complet conține 70 ore de transmisii colectate prin antene și receptoare radio care sunt localizate în vecinătatea aeroporturilor.

Au fost create corpusuri de înregistrări audio în cadrul întâlnirilor în scopul analizei de discurs, extracției de metadata, identificare de vorbitor și recunoaștere vocală. Dintre acestea se enumeră **„2004 Spring NIST Rich Transcription (RT-04S) Evaluation Data”** (J. Fiscus et al., 2007a) care conține material de test (înregistrări în cadrul întâlnirilor și transcripții referință) utilizat în evaluarea RT-04S. Transcripția bogată este definită ca o fuziune între tehnologia speech-to-text și tehnologii de extracție metadata creat pentru a furniza o bază pentru generarea celor mai utilizate transcripții în cadrul întâlnirilor dintre oameni. RT-04S include două tipuri de sarcini și anume transcripție vorbire – text (STT) și „diarization” (cine și când vorbește) (SPKR). Sarcina STT include selectarea condițiilor în care se fac înregistrările și condițiile de procesare. Se pot utiliza unul sau mai multe microfoane la distanță sau căști individuale cu microfon, iar timpul de procesare poate fi nelimitat, mai mic sau egal cu 10, 20 de ori sau o dată din timpul real. În cadrul SPKR se poate utiliza unul sau mai multe microfoane la distanță, timpul de procesare este similar cu cel pentru STT și trebuie luate în calcul condițiile de input și anume se poate utiliza ca input doar voce sau voce plus o transcripție de referință. Aceiași autori au dezvoltat în mod similar și un corpus de dezvoltare **„2004 Spring NIST Rich Transcription (RT-04S) Development Data”** (J. Fiscus et al., 2007b). Ambele sunt disponibile pe DVD contracost.

„2006 NIST Spoken Term Detection Development Set” (NIST Multimodal Information Group, 2011b) și **„2006 NIST Spoken Term Detection Evaluation Set”** (NIST Multimodal Information Group, 2011a) sunt două corpusuri dezvoltate în scopul detecției de termen (expresie) în limba vorbită. Acestea conțin înregistrări ale emisiunilor radio de știri, întâlniri în săli de conferințe și conversații telefonice. Fișierele provenite de la înregistrările emisiunilor de știri sunt pe un singur canal, codate pcm, cu frecvență de eșantionare de 16 kHz, în format SPHERE. Cele provenite de la conversațiile telefonice sunt pe două canale, eșantionate pe 8 kHz, iar cele provenite de la conferințe sunt pe un singur canal. Aceste corpusuri au fost dezvoltate în scopul găsirii tuturor ocurențelor unui



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

anumit termen (o secvență de unul sau mai multe cuvinte) într-un anumit corpus de semnal vocal. Evaluarea s-a realizat în scopul dezvoltării unei tehnologii pentru căutare rapidă în cantități mari de date audio. Sistemele au fost implementate în două faze: indexare și căutare. În faza de indexare sistemul procesează datele audio fără cunoașterea termenilor, iar în faza de căutare sistemul utilizează termenii, indexul și opțional datele audio pentru a detecta ocurențele de termen. Cele două corpusuri sunt disponibile pe DVD contracost.

2.2. Metodologii utilizate la crearea de corpusuri de semnal vocal și instrumente aferente

2.2.1. Metodologii pentru crearea de corpusuri destinate recunoașterii vocale

O. Salor et al. (O. Salor et al., 2007) a dezvoltat un corpus de semnal vocal care conține 193 înregistrări și instrumente de recunoaștere vocală (un aliniator fonetic și un sistem de recunoaștere de fonem) pentru limba turcă bazate pe motorul de recunoaștere vocală SONIC. Metoda utilizată de autori în dezvoltarea corpusului constă în analiza lingvistică în vederea determinării alfabetului, crearea corpusului de text și ulterior crearea corpusului audio. Pentru obținerea unui set echilibrat de propoziții în limba turcă au fost considerate drept unități de bază trifonemele deoarece s-a raportat că sistemele de recunoaștere vocală care utilizează modelul Markov ascuns (HMM) modelează cu o acuratețe mai mare trifonemele în comparație cu fonemele, cuvintele sau silabele.

Metoda utilizată de autori în vederea creării unui corpus audio din trifoneme a constat în determinarea celor mai frecvente ocurențe dintr-un corpus de text suficient de mare pentru a modela limba în mod satisfăcător. Pe baza acestora a fost creat un set de propoziții care a stat la baza dezvoltării corpusului audio. Setul de propoziții a fost obținut prin traducerea în limba turcă a primelor 2000 de propoziții din corpusul TIMIT și convertirea acestora în simboluri METUbet care reprezintă o mapare a simbolurilor IPA într-un nou set de simboluri ASCII. Ocurențele trifonemelor în setul de propoziții au fost obținute și comparate cu cele găsite pentru un corpus alcătuit din 2.5 milioane simboluri de cuvinte. A fost adăugat un extra set de 462 propoziții pentru a acoperi cele mai frecvente 5000 de trifoneme existente în corpusul de text.

Corpusul audio a fost creat după metoda utilizată la dezvoltarea corpusului TIMIT. Pentru fiecare vorbitor au fost selectate în mod aleator 40 din 2462 propoziții care au fost citite o singură dată. Au fost colectate înregistrări provenind de la 193 vorbitori (89 de gen



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOS DRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

feminin și 104 de gen masculin) cu vârste cuprinse între 19 și 50 ani, cu o medie de 23.9 ani care au fost realizate într-un mediu fără zgomot cu un set de căști prevăzute cu microfon, pe un calculator, cu o frecvență de eșantionare de 16 kHz.

Sistemele de aliniere fonetică și de recunoaștere de fonem dezvoltate de autori utilizează toolkit-ul SONIC care conține un set de instrumente de recunoaștere a vorbirii continue utilizat ca bază de testare a instrumentelor dezvoltate de specialiști în domeniu. Ambele sisteme utilizează o listă de foneme, un dicționar fonetic, modele fonetice și modele acustice în limba turcă. Modelele acustice incluse în SONIC sunt modele Markov ascunse (HMM) bazate pe arbori de decizie la nivel de stare conținând funcții asociate densitate de probabilitate gamma pentru a modela duratele stării acestuia. Portul SONIC în limba turcă dezvoltat de autori este primul creat pentru o altă limbă decât limba engleză.

Corpusul dezvoltat de autori a fost acceptat de LDC în octombrie 2005 și este disponibil pentru cercetători în domeniu. Sistemele de aliniere și de recunoaștere fonetică au fost testate și dau rezultate comparabile cu cele obținute pentru limba engleză. Rezultatele obținute pentru recunoașterea fonetică au arătat că erorile de recunoaștere sunt mai mici la cuvintele finale din propoziții care în general sunt verbe. Matricile de confuzie create pentru foneme au arătat că plozivele, fricativele și africaterile (P, T, B, D, YH, CH, S) din limba turcă necesită un algoritm de extragere de trăsături mai bun, iar vocala O este singura vocală care are o probabilitate de confuzie ridicată.

M. Boldea et al. (M. Boldea et al., 1998) a creat un corpus de semnal vocal annotat în limba română destinat în principal dezvoltării de instrumente de recunoaștere vocală independente de vorbitor. Acesta este alcătuit din 100 vorbitori, iar durata totală de înregistrare este de 10 ore. Textele utilizate pentru înregistrări au fost construite prin traducerea versiunii din limba engleză a 40 pasaje din corpusul EUROM-1 care au fost împărțite în 10 clustere printr-o metodă euristică care calculează apriori numărul de ocurențe a fiecărui fonem în fiecare cluster. Pentru a crește variația fonetică și pentru a furniza modele de foneme dependente de context au fost adăugate un număr de 550 propoziții individuale, numere întregi, cuvinte în context CVC și în contexte izolate sau controlate.

Înregistrările s-au realizat într-o cameră fără zgomot, cu softul EUROPEC pe un calculator compatibil și echipat cu o placă de conversie OROS AU-21 A/D-D/A. Microfonul utilizat de tip electret a fost plasat la 25 cm față de gura vorbitorului, iar ieșirea lui a fost legată la placa OROS printr-un preamplificator cu câștig fix. Frecvența de



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

eșantionare setată a fost de 20.000 Hz, cu 16 biți pe eșantion. Fiecare fișier rezultat a fost acompaniat de fișiere de configurare și descriere.

Adnotarea fișierelor audio a fost realizată în mod similar cu cea pentru corpusul TIMIT și anume transcripție manuală, segmentare și aliniere automată utilizând HMM, verificare și corectare manuală. Vorbitorii înregistrați cu vârste cuprinse între 20 și 50 de ani au fost împărțiți pe cinci categorii de vârstă.

2.2.2. Metodologii pentru crearea de corpusuri destinate sintezei de voce

Stan et al. (Stan et al., 2011) au dezvoltat un corpus de limbă română numit RSS pe care l-au utilizat la crearea unui sistem de sinteză voce de înaltă calitate bazat pe HMM-uri, utilizând o frecvență de eșantionare înaltă. Metodologia constituirii corpusului constă în realizarea înregistrărilor într-o cameră fără ecou, cu trei microfoane: un microfon cu diafragmă mare (Neumann u89i), un microfon cu diafragmă mică și bandă largă de frecvențe (Sennheiser MKH 800) și un set de căști cu microfon (DPA 4035). Toate înregistrările au fost realizate cu o frecvență de eșantionare de 96 kHz pe 24 biți și reeșantionate cu 48 kHz. Această metodă numită supra-eșantionare care presupune eșantionare cu o frecvență de patru ori mai mare decât frecvența Nyquist ajută la eliminarea zgomotului, raportul semnal-zgomot îmbunătățindu-se cu un factor egal cu 4. Pentru înregistrări și reeșantionare autorii au utilizat instrumentele hardware și software „Pro Tools HD”. Toate înregistrările aparțin unui vorbitor de gen feminin tânăr, vorbitor nativ de limbă română. S-au realizat opt sesiuni într-o lună, în care s-au înregistrat câte 500 de propoziții.

Corpusul constituit este alcătuit din două seturi de propoziții, unul de antrenare și unul de testare. Durata totală de înregistrare pentru setul de antrenare a fost de 3.5 ore și conține 3500 propoziții din care 1500 sunt alese aleator din propoziții de ziar, 1000 de propoziții sunt alese astfel încât să acopere cele mai frecvente difoneme și 1000 provin din pasaje de basme. Timpul de înregistrare pentru setul de testare a fost de 0.5 ore și conține 200 respectiv 100 de propoziții selectate aleator din ziare și povești, iar 200 sunt propoziții nepredictibile semantic. Rezultatele obținute arată că, corpusul creat este potrivit pentru utilizarea în sisteme de sinteză voce bazate pe HMM-uri și că sistemul de sinteză voce dezvoltat prezintă o bună inteligibilitate. Autorii au constatat că o reeșantionare la 32 kHz nu afectează rezultatele, iar reeșantionarea la 16 kHz degradează similaritatea vocii vorbitorului. Corpusul creat este disponibil gratuit online pentru cercetători în domeniu.



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

J. Matousek și J. Romportl (J. Matousek și J. Romportl, 2007) au creat un corpus de limbă cehă în vederea utilizării în sinteza de semnal vocal. În mod identic autorilor din (Stan et al., 2011), înregistrările s-au realizat pe o singură voce de gen feminin și au inclus pasaje din ziare pe diverse domenii. Autorii au utilizat un algoritm de selecție automat pentru selectarea propozițiilor de interes. Din 524.472 propoziții au fost selectate 5000 astfel încât setul rezultat să fie echilibrat din punct de vedere fonetic și prozodic. Au mai fost adăugate o serie de propoziții afirmative, interogative, exclamative și imperative și au fost eliminate cele care aveau un număr mai mic de trei sau mai mare de 30 cuvinte. Înregistrările s-au realizat într-un studio audio cu un microfon cu diafragmă mare prevăzută cu un filtru pentru a reduce forța aerului expirat în timpul pronunției plozivelor bilabiale sau a altor consoane. Un card de captură cu fidelitate înaltă a fost utilizat pentru conversia AD de 48 kHz. Concomitent a fost măsurat semnalul glotal în vederea detectării momentelor de închidere a glotei utilizate la estimarea cu acuratețe a conturului frecvenței fundamentale (F0), la sinteza de voce sincronizată cu F0 sau la detectarea precisă a zonelor vocalice/nevocalice.

Adnotarea corpusului creat s-a făcut în două faze cu ajutorul soft-ului TranscriberTM. În prima fază un adnotator a realizat adnotarea având ca referință seturile de propoziții selectate cu sistemul automat, iar în a doua fază un alt specialist a verificat adnotările facute de primul adnotator și le-a corectat unde a fost necesar. În procesul de adnotare au fost luate în considerare o serie de reguli referitoare la simbolurile utilizate și la porțiunile din înregistrări care trebuiau adnotate. Au fost etichetate fonemele, regiunile de liniște, lipsa segmentelor de liniște înaintea și la sfârșitul înregistrărilor, zgomotele de tip respirație, cele produse de buze sau alte zgomote datorate sistemului de achiziție. În urma adnotării finale au rezultat 62.332 cuvinte din care 7.6% evenimente fără semnal vocal și 2.62% excepții de tipul abrevierilor. Autorii concluzionează că deși diferențele dintre cele două faze ale adnotării nu sunt foarte mari și anume în prima fază au fost adnotate corect un procent de 99.62 cuvinte, luând în considerare numărul total de cuvinte din corpus, 237 cuvinte au fost adnotate greșit ceea ce ar produce probleme în sistemul de sinteză voce în timpul selectării unităților.

2.2.3. Metodologie pentru crearea unui corpus de dialecte

C. G. Clopper și D. B. Pisoni (C. G. Clopper și D. B. Pisoni, 2006) au dezvoltat în cadrul proiectului „The Nationwide Speech Project” un corpus de limbă engleză în



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

scopul analizei variațiilor la nivel de dialect. Acesta este alcătuit din cuvinte, propoziții, pasaje și interviuri citite de a câte cinci vorbitori de gen feminin și cinci vorbitori de gen masculin din șase regiuni dialectale ale Statelor Unite. Înregistrările au fost realizate cu echipament de înregistrare de înaltă calitate într-un mediu fără zgomot. În constituirea corpusului autorii au ținut cont de următorii factori: informații demografice ale vorbitorilor, condiții de înregistrare și tipul de material înregistrat. Au fost utilizate patru tipuri de material și anume cuvinte izolate, propoziții, pasaje și interviuri, timpul total de înregistrare per vorbitor fiind de aproximativ o oră. Cuvintele izolate au fost împărțite în seturi de cuvinte în context hVd, CVC și cuvinte multisilabice. Cuvintele hVd conțin cinci repetiții a 10 vocale din engleza americană, iar cuvintele CVC sunt alcătuite din 76 cuvinte monosilabice. Fiecare din cele 14 vocale de tip monoftong și diftong au fost incluse de cel puțin patru ori în lista CVC, iar contextul următor vocalei a fost variat astfel încât să includă consoane lichide, nazale, nesonore și sonore obstruative.

Propozițiile înregistrate au fost împărțite în funcție de probabilitatea de ocurență în propoziții cu probabilitate de ocurență ridicată, mică și propoziții greșite din punct de vedere semnatic. Au fost înregistrate două interviuri pentru fiecare vorbitor dintre care unul includea informații legate de rezidență, activități extracuriculare (5 min), iar al doilea a fost structurat astfel încât să conțină cuvinte țintă rostite de vorbitor în mod natural în timpul conversației (7-12 min).

Vorbitorii au fost recrutați din cadrul comunității Universității Indiana, au vârste cuprinse între 18 și 25 ani, o singură limbă nativă, majoritatea nelicențiați și fără probleme de auz sau patologii ale vocii la momentul înregistrării. Materialul înregistrat respectiv cuvinte izolate, propoziții, interviuri au fost separate pe blocuri și afișate pe rând pe un ecran LCD legat la un laptop. Vorbitorul situat în fața ecranului a purtat un microfon atașat la cap, situat la 2.5 cm de colțul stâng al gurii. Ieșirea microfonului a fost legată la un preamplificator, iar câștigul acestuia a fost reglat de către experimentator în timpul pronunției unei propoziții test de către vorbitor. Ieșirea preamplificatorului a fost conectată printr-o o interfață audio USB Roland UA-30 la laptop. De această interfață a fost legat și un set de căști purtate de către experimentator pentru a auzi semnalul de input. Fiecare pronunție a fost înregistrată într-un fișier audio digital de tip AIFF pe 16 biți, la o frecvență de eșantionare de 44.1 kHz.

Corpusul de semnal vocal creat este organizat în funcție de vorbitor și ulterior în funcție de tipul stimulului într-o structură ierarhică de fișiere numite printr-un identificator de trei caractere pentru vorbitor, un caracter pentru tipul stimulului și patru caractere pentru



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

numărul stimulului. Pentru fiecare vorbitor s-a creat un fișier în care este definită ordinea în care i s-au prezentat blocurile de cuvinte și un fișier log care conține ordinea prezentării stimulilor și legături între numărul stimulului, transcripția ortografică a stimulului și numele fișierului audio.

Propozițiile înregistrate care prezintă o probabilitate ridicată de ocurență au fost utilizate de către autori ca materiale de stimul într-o serie de sarcini perceptuale. Rezultatele obținute arată că ascultătorii de limbă nativă pot percepe și reprezenta proprietăți foneto-acustice remarcabile ale dialectelor limbii engleze americane și pot face judecăți explicite despre vorbitori bazate pe dialectul regional.

2.2.4. Metodologie pentru crearea corpusului multilingv „PFSTAR”

Corpusul de semnal vocal „PF Star” multilingv dezvoltat de A. Batliner et al. (A. Batliner et al., 2005) conține 60 de ore de înregistrări provenite de la 611 copii cu vârste cuprinse între 4-14 ani, vorbitori nativi de limbă engleză, germană, italiană și suedeză. Înregistrările provin din citire, imitare, vorbire spontană sau emoțională. Corpusul este creat pentru a servi ca bază în dezvoltarea sistemelor de recunoaștere automată a vocilor de copii cu limbă nativă sau nu, a vorbirii spontane și emoționale. Au fost utilizate materiale și procedee de înregistrare comune pentru toate limbile, iar metodologia AIBO a fost utilizată pentru ambele site-uri de colectare a vorbirii spontane și emoționale. Materialul utilizat pentru înregistrări constă în 220 cuvinte izolate în limba engleză, 50 de propoziții în limba engleză bogate din punct de vedere fonetic și 400 propoziții generice. Toate acestea au fost utilizate la înregistrarea copiilor vorbitori nativi de limbă engleză, italiană și suedeză. Pentru constituirea corpusului de voci spontane și emoționale a fost utilizată metodologia „AIBO” care constă în utilizarea unui robot controlat căruia copii îi dau instrucțiuni crezând că vorbește cu ei. S-au realizat trei experimente și anume cu robot ascultător care îndeplinește sarcinile primite, cu robot neascultător care este utilizat pentru a elicită vorbirea emoțională și cu robot care după instrucțiunile primite trebuie să localizeze obiecte.

Ca metodă generală de înregistrare a fost utilizat un set de căști cu microfon prevăzut cu preamplificator, frecvența de eșantionate a fost de 44.1 kHz pe 16 biți și semnalul a fost transmis prin intermediul placii audio la un calculator

Pentru înregistrările de voce spontană și emoțională s-a utilizat un set de căști cu microfon wireless, un DATrecorder – casetă audio digitală, o frecvență de eșantionare de



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

48 kHz și cuantizare de 16 biți. Fiecare înregistrare a durat 30 de minute, iar timpul total al înregistrărilor eliminând secvențele de liniște datorate răspunsului robotului a fost de 9.2 ore. Toate înregistrările au fost adnotate cu transcripții fonetice marcându-se informațiile referitoare la zgomot sau sunete extralingvistice.

2.2.5. Metodologie pentru crearea corpusului „CIAIR inCar”

Un alt tip de metodologie a fost utilizat la crearea corpusului de semnal vocal „CIAIR In Car” (N. Kawaguchi et al., 2005) destinat analizei vorbirii în timpul conducerii unui autovehicul și realizarea diverselor manevre în trafic. Autorii au dezvoltat un sistem numit „Data Collection Vehicle” (DCV) care suportă înregistrări sincrone de date multicanal provenite de la 12 microfoane plasate în interiorul vehiculului, înregistrări de la trei camere video multicanal, și colectarea datelor referitoare la vehicul cum ar fi viteză, accelerare, frână, numărul de rotații pe minut ale motorului, etc. Corpusul este alcătuit din 812 subiecți și conține 1960 sesiuni totalizând un număr de 187.5 ore. Autorii au studiat un corpus format din 1.06 milioane morfeme pentru a determina informații fundamentale specifice vocii și au utilizat informațiile referitoare la autovehicul și manevrele făcute pentru a studia modurile în care acestea influențează pronunțiile vorbitorului. Au fost utilizate trei tipuri de sisteme dialog: un sistem bazat pe operator uman (HUM), un sistem de tip Wizard of Oz (WOZ) și un sistem de dialog de limbă vorbită (SYS). Toți vorbitorii au completat un chestionar cu informații despre experiența în conducere, în utilizarea sistemelor electronice, de recunoaștere vocală și de navigație. Corpusul a fost transcris fonetic la nivel de morfem în vederea analizei caracteristicilor specifice limbii vorbite în timpul conducerii de autovehicule.

În urma acestui studiu autorii au concluzionat că viteza de pronunție pentru șoferi a fost mai mică decât în cazul dialogului regulat; dialogul cu sisteme de dialog a constatat în pronunții mai scurte comparativ cu dialogul cu operator uman; condițiile de conducere (în timpul conducerii sau în pauză) nu influențează dialogul; operațiile efectuate de tip accelerare sau manevră volan influențează semnificativ pronunțiile.



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

3. Corpusuri de voci patologice și instrumente aferente

În vederea dezvoltării de instrumente de evaluare, educare și diagnostic automat a vocii patologice au fost dezvoltate corpusuri și microcorpusuri de semnal vocal aparținând persoanelor cu diverse afecțiuni care influențează producerea și articularea corectă a sunetelor cum ar fi: afecțiuni ale laringelui, afecțiuni respiratorii, neurologice, stomatologice sau de comunicare.

Pentru constituirea acestor corpusuri este necesar un protocol și o metodologie specială în funcție de aplicațiile urmărite. Suplimentar corpusurilor realizate pentru aplicații de sinteză și recunoaștere voce este necesară colectarea de informații medicale referitoare la diagnostic, boli asociate și simptome specifice care duc la modificarea vocii. În ceea ce privește materialul text de înregistrat trebuie ales în funcție de zona din canalul fonator afectată. Spre exemplu dacă se dorește evaluarea funcției corzilor vocale sau a aparatului respirator se vor face înregistrări de vocale susținute.

În vederea identificării unor patologii de dentiție se va avea în vedere includerea în corpusul text de consoane a căror locuri de articulare sunt dental, labio-dental sau alveolar. Va fi de asemenea necesară semnarea unui protocol suplimentar de protecție a datelor pacientului. În ceea ce privește protocolul de înregistrare acesta nu diferă semnificativ de cele utilizate la crearea corpusurilor cu alte aplicații. Este necesară utilizarea unei camere cu zgomot redus și instrumente de înregistrare performante cu respectarea protocoalelor standard legate de frecvența de eșantionare, cuantizare și achiziție.

La metoda de înregistrare standard cu microfon se pot adăuga instrumente suplimentare cum ar fi electroglotograf, stroboscop sau aparat de măsurare a fluxului de aer expirat în funcție de aplicație. După tipul analizei efectuate este necesară repetarea înregistrărilor înainte și după tratament sau operație pentru a evalua îmbunătățirea sau înrăutățirea calității vocii. În cazul realizării corpusului de către un cercetător din domeniu tehnic-informatic este necesară colaborarea cu un medic specialist în vederea colectării informațiilor cu privire la diagnostic.



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

3.1. Baza de date de voci patologice dezvoltată de MEEI (MEEI VDDb)

O bază de date cu voci patologice în limba engleză comercială, disponibilă la ora actuală, este distribuită de Kay ElemetricsTM și a fost dezvoltată de Massachusetts Eye and Ear Infirmary Voice and Speech Lab (**MEEI VDDb**) (Kay Elemetrics Corp, 1994). Corpusul conține 53 respectiv 657 fișiere audio de înregistrări cu voci normale, respectiv patologice ale vocalei susținute /ah/ și 53 respectiv 661 fișiere de înregistrări cu voci normale respectiv patologice de vorbire continuă. Baza de date conține și fișiere cu informații personale și clinice pentru fiecare pacient și rezultatele obținute în urma analizei acustice obținute cu programul MDVP dezvoltat de Kay ElemetricsTM. Înregistrările au fost efectuate cu ajutorul programului CSL (Computerized Speech Lab.) cu frecvențe de eșantionare de 25 și 50 kHz pe 16 biți. Fiecare înregistrare conține trei repetiții de vocală susținută /ah/ a câte 3 s fiecare. Un specialist pe voce patologică a ales din aceste trei pronunții cea mai bună variantă pentru a fi inclusă în corpus. Această bază de date este larg utilizată de către specialiști în domeniu la dezvoltarea de sisteme și metode de detecție de voce patologică, însă conform observațiilor făcute de (N. S. Lechon et al., 2006) sunt câteva puncte cheie care trebuie luate în considerare în vederea utilizării acesteia în cercetare:

- nu toți pacienții cu voce patologică au înregistrări sau diagnostic asociat și sunt unii pacienți care au mai multe înregistrări din diferite vizite la clinică;
- fișierele au frecvențe de eșantionare diferite. Cele provenite de la subiecții normali și un mic procent din cele provenite de la subiecții patologici sunt eșantionate cu 50 Hz, iar celelalte cu 25 kHz;
- vocile normale și patologice au fost înregistrate în locații diferite, iar subiecții normali nu au fost evaluați clinic;
- fișierele sunt deja editate astfel încât să includă doar partea stabilă din fonație și astfel se pierde informațiile cu privire la raportul semnal-zgomot și la părțile de onset și offset care conform unor studii în domeniu conțin mai multă informație acustică de interes;
- fișierele cu vocale susținute și cele cu vorbire continuă a subiecților normali au o durată de 3 s respectiv 12 s, iar cele aparținând subiecților cu voce patologică au durate de 1 s respectiv 9 s;
- există o singură fonație a unei singure vocale per pacient și per vizită;



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

- există un număr heterogen de patologii în baza de date, în jur de 200 diagnostice diferite și sunt multe fișiere etichetate cu mai multe diagnostice, aparținând la categorii diferite de patologii ale vocii;
- în comparație cu înregistrările de voci patologice există un număr mic de voci normale;
- înregistrările nu sunt evaluate din punct de vedere perceptiv cum ar fi indice GRABAS sau prin alte metode și nu sunt înregistrări video asociate (stroboscopii, endoscopii).

3.2. Metode de detecție de voci patologice care utilizează ca referință MEEI VDDb

Corpusul de voci patologice dezvoltat de Massachusetts Eye and Ear Infirmary Voice and Speech Lab a fost larg utilizat ca referință în dezvoltarea și testarea de metode și sisteme de recunoaștere automată a vocilor disfonice, a dizartriei, a patologiei de laringe, a paraliziei și edemului corzilor vocale, etc.

N. S. Lechon et al. (N. S. Lechon et al., 2006) prezintă o serie de preocupări metodologice care ar trebui luate în calcul la crearea sistemelor automate de detecție a vocii patologice în vederea comparării cu experimente precedente sau viitoare. Autorii subliniază că orice experiment trebuie să aibă o strategie de cros-validare, iar rezultatele trebuie să includă o matrice de confuzie, intervale de confidență pentru toate măsurătorile și grafice de detecție a performanței sistemului cum ar fi curba DET (Detector Error Tradeoff = eroarea detectorului de compromis) și curba ROC (receiver operating characteristic = caracteristica de operare a receptorului). Autorii prezintă un exemplu de metodologie și un experiment bazat pe parametri pe termen scurt și percepșoni multi-nivel. Experimentul este realizat pe un grup de 53 voci normale și 173 voci patologice care au fost împărțite în două seturi: unul de antrenare și unul de testare și validare a rezultatelor, reprezentând 70% respectiv 30% din numărul total de voci, astfel încât sexul și vârsta să fie uniform distribuite între cele două clase. După antrenarea sistemului sau calcularea modelelor, setul de test este utilizat pentru estimarea performanței detectorului. În acest scop sunt definite trei măsuri:

- (tp) rata adevărat-positiv (senzitivitatea) care reprezintă raportul dintre numărul vocilor patologice corect clasificate și numărul total de voci patologice;
- (fn) rata fals-negativ care este raportul dintre numărul vocilor patologice greșit clasificate și numărul total de voci patologice;



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

- (tn) rata adevărat-negativ (specificitatea) care reprezintă raportul dintre numărul vocilor normale clasificate corect și numărul total de voci normale;
- (fp) rata fals-positiv care este raportul dintre numărul total de voci normale clasificate incorect și numărul total de voci normale.

Acuratețea finală a sistemului se obține prin calculul raportului dintre toate vocile clasificate corect de sistem și numărul total de voci analizate.

Pentru evaluarea capacității de generalizare a sistemului trebuie adoptată o schemă de cros-validare. Cea mai simplă constă în a repeta fiecare experiment de N ori, cu un set de test diferit ales random din întreg setul de fișiere sau o altă metodă constă în împărțirea random a setului de date în K subseturi și repetarea experimentului de K ori utilizând de fiecare dată un set diferit de testare a performanței. Când numărul K este egal cu numărul F de fișiere disponibile metoda este cunoscută sub numele crosvalidare „leave-one-out”. Experimentul este repetat de F ori, iar de fiecare dată sistemul este antrenat cu $F-1$ fișiere, lăsând fișierul rămas pentru testare. După crosvalidare rezultatele finale se mediază între repetiții și intervalele de confidență pot fi calculate utilizând deviațiile standard ale măsurătorilor. O schemă a matricii de confuzie este dată în Tabelul 1.

Tabel 1. Ilustrarea unei matrici tipice de confuzie a sistemului de clasificare (N. S. Lechon et al., 2006)

Decizia detectorului	Diagnostic actual	
	Patologic	Normal
Patologic	tp	fp
Normal	fn	tn

Matricea de confuzie se calculează pe baza comparației dintre scorul obținut în urma testării și o valoare prag. Dacă modificăm această valoare prag obținem o serie de puncte posibile de operare ale sistemului care pot fi reprezentate printr-o curbă DET frecvent utilizată în verificarea vorbitorului. Această curbă se obține prin reprezentarea grafică a valorilor fals pozitive în funcție de cele fals negative la diferite valori prag. O altă curbă des utilizată în sistemele medicale de decizie este curba ROC care se obține prin reprezentarea valorilor fals pozitive în funcție de cele adevărat pozitive. Cea mai utilizată măsură este aria de sub curba ROC (AUC) care în cazul în care performanța sistemului este slabă are valoare în jur de 0.5, iar dacă performanța este bună are valoarea 1. Eficiența sistemului este ridicată când curba DET se apropie de colțul stâng de jos al graficului și curba ROC se apropie de colțul stâng de sus al graficului.

Utilizând metodologia descrisă mai sus autorii au dezvoltat un sistem de recunoaștere a vocii patologice bazat pe o rețea neuronală care modelează cele două clase



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

de voci (normală și patologică). Detectorul este un perceptron de tip „feedforward” cu mai multe nivele. Nivelul de intrare este alcătuit din mai multe intrări ale parametrilor MFCC, nivelul ascuns are 12 neuroni, iar nivelul de ieșire are două noduri. Cele două ieșiri sunt construite astfel încât să se obțină o rată a similitudinii sau un scor pentru fiecare pattern de intrare. Este utilizat un algoritm de învățare supervizat bazat pe metoda backpropagation cu moment și regulă delta. Funcțiile de activare ale tuturor nodurilor sunt de tip logic, iar ponderile de conectare sunt inițializate cu valori random dintr-o distribuție gauss cu medie zero și deviație standard inversă numărului de ponderi al fiecărui neuron. Au fost realizate 40 de iterații ale algoritmului de antrenare, iar experimentul a fost repetat de 10 ori utilizând seturi random diferite de antrenare și testare. Acuratețea totală a sistemului este de 89.6% cu o deviație standard între experimente de $\pm 2.49\%$, iar procentajul de voci clasificate corect este de 90.42%.

Un alt sistem de detectare automată a vocii normale și patologice utilizând eșantioane de voce din corpusul MEEI VDDb a fost creat de **A. A. Dibazar și S. Narayanan** (A. A. Dibazar și S. Narayanan, 2002). Sistemul de detecție este un clasificator de tip HMM care modelează coeficienții Mel cepstrali și dinamicile frecvenței fundamentale prin mixuri gaussiene. Metoda a fost evaluată utilizând caracteristici ale vocalei susținute /a/ provenind de la 700 subiecți cu voci normale și patologice și a constat în asocierea fiecărui subiect a unui vector de trăsături care include 42 trăsături extrase din informația spectrală și a F0. Pentru antrenarea HMM-ului a fost utilizat toolkit-ul HTK care a fost modificat pentru adaptarea la frecvența fundamentală. Autorii au comparat rezultatele obținute cu această metodă cu cele obținute prin clasificarea vocilor utilizând caracteristici în domeniul timp frecvență extrase cu programul MDVP (Multi Dimensional Voice Program) și patru clasificatori: LDC (analiză liniară discriminantă), NCM (clasificatorul bazat pe cea mai apropiată medie), HMM (model cu trei mixturi gaussiene) și o rețea neuronală cu 34 neuroni de intrare, un nivel ascuns și doi neuroni de ieșire.

Acuratețea sistemului bazat pe mixturi gaussiene cu utilizarea coeficienților Mel cepstrali și a frecvenței fundamentale este de 99.40%, iar cea obținută în cazul utilizării parametrilor extrași cu MDVP (F0 minim, F0 maxim, F0 mediu, jitter, shimmer, etc.) este de 97.97%. Autorii au realizat și o clasificare la nivel de patologie, împărțind vocile patologice în patru clase și anume patologii minoră, ușoară, moderată și severă. La evaluarea patologiei utilizând parametrii MDVP pentru vocala susținută /a/ și în cazul utilizării parametrilor spectrali și F0 s-au obținut rate de clasificare de 57.99%, respectiv 72.55%. Atât la discriminarea de voce normală-patologică cât și la identificarea de tip de



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

patologie sistemul HMM cu mixturi gaussiene utilizând trăsături spectrale și F0 a arătat o îmbunătățire semnificativă a rezultatelor de clasificare.

R. J. Moran et al. (R. J. Moran et al., 2006) au construit un sistem de detectare a patologiilor corzilor vocale prin înregistrări telefonice. Sistemul are la bază un clasificator liniar, măsurători de procesare a perturbațiilor frecvenței fundamentale, a amplitudinii și a raportului semnal-zgomot (HNR). Pentru dezvoltarea și validarea sistemului au fost utilizate eșantioane de voce din MEEI VDDb .

Plecând de la baza de date originală au fost create cinci corpusuri. Patru din ele alcătuite din 631 înregistrări de voce au fost create prin reeșantionarea la 10 kHz și adăugarea a diverse distorsiuni care pot fi transmise prin telefon. Acestea au fost utilizate la evaluarea acurateții de discriminare între vocile normale și patologice cu diferite seturi de trăsături. Al cincilea corpus a fost construit prin transmiterea tuturor înregistrărilor prin intermediul unui canal telefonic la distanță mare și înregistrarea fiecărui semnal primit la receptor.

Trăsăturile utilizate în construirea clasificatorului au fost perturbațiile de F0 și de amplitudine și HNR. Dintre trăsăturile comune ale F0 și amplitudinii sunt parametrii statistici medie ($F0_{av}, Amp_{av}$), minim ($F0_{lo}, A_{lo}$), maxim ($F0_{hi}, A_{hi}$) și deviație standard ($F0_{sd}, A_{sd}$). Alte trăsături derivate din F0 sunt []:

- PFR (Phonatory frequency range):

$$PFR = \frac{\log\left(\frac{F0_{hi}}{F0_{lo}}\right)}{\log 2} \times 12, \quad (1)$$

- MAJ (Media absolută a jitter-ului):

$$MAJ = \frac{1}{n-1} \sum_{i=1}^{n-1} |F_{i+1} - F_i|, \quad (2)$$

- Jitter (%):

$$JITT = \frac{MAJ}{F0_{av}}, \quad (3)$$

- RAP (Media perturbării relative regulată pe 3 perioade):

$$RAP = \frac{1}{n-2} \sum_{i=2}^{n-1} \left| \frac{F_{i+1} + F_i + F_{i-1}}{3} - F_i \right| \times 100, \quad (4)$$

- PPQ_5 (Coeficientul de perturbație a F0 regulat pe 5 perioade):

$$PPQ_5 = \frac{1}{n-4} \sum_{i=3}^{n-2} \left| \frac{\sum_{k=i-2}^{i+2} F(k)}{5} - F_i \right| \times 100, \quad (5)$$

- PPQ₅₅ (Coeficientul de perturbație a F0 regulat pe 55 perioade):

$$PPQ_55 = \frac{1}{n-54} \sum_{i=28}^{n-27} \left| \frac{\sum_{k=i-27}^{i+27} F(k)}{55} - F_i \right| \times 100, \quad (6)$$

- PPF (factorul de perturbație a F0):

$$PPF = \frac{N_{p \geq \text{threshold}}}{N_{\text{voice}}}, \quad (7)$$

unde N_p este perturbația în magnitudine a perioadei în timp mai mare decât 0.5 msec.

- DPF (factorul direcțional de perturbație):

$$DPF = \frac{N_{\Delta \pm}}{N_{\text{voice}}} \times 100, \quad (8)$$

unde $N_{\Delta \pm}$ este perturbația epocii în timp pentru care există o schimbare de semn.

Trăsăturile derivate din perturbațiile amplitudinii sunt:

- Shimmer (SHIM)(%):

$$SHIM \% = \frac{MAS}{Amp_av}, \quad (9)$$

unde MAS este media absolută a shimmer-ului care se calculează similar cu MAJ .

- Shimmer (SHIM)(db):

$$SHIM = \frac{1}{n-1} \sum_{i=1}^{n-1} 20 \log \frac{A_i}{A_{i+1}}, \quad (10)$$

- ARP₃ (media perturbării relative a amplitudinii regulată pe 3 perioade), APQ₅, APQ₅₅, APF și ADPF care se calculează similar cu RAP, PPQ₅, PPQ₅₅, PPF și DPF.

Histogramele aferente unui număr de 35 trăsături utilizate realizate pe grupul de voci normale și cel de voci patologice au arătat că acestea prezintă o distribuție aproximativ gaussiană ceea ce a înlesnit construirea unui clasificator de discriminare liniară (LDC).

Metoda generală de antrenare a clasificatorului utilizată de către autori a constatat în următorii pași: (a) crearea unui vector coloană x care conține d trăsături care trebuie asignat la una din cele două clase; se consideră că există N_1 vectori de trăsături pentru antrenarea clasificatorului din clasa 1 și N_2 trăsături din clasa 2; al n -lea vector de trăsături pentru



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

antrenare în clasa k este notat cu $x_n^{(k)}$; (b) determinarea vectorilor medie μ_1, μ_2 care condiționează clasele cu formulele []:

$$\mu_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} x_n^{(1)}, \mu_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} x_n^{(2)}. \quad (11)$$

și a matricii comune de covarianță:

$$\Sigma = \frac{1}{N_1 + N_2 - 2} \sum_{k=1}^2 \sum_{n=1}^{N_k} (x_n^{(k)} - \mu_k)(x_n^{(k)} - \mu_k)^T. \quad (12)$$

(c) clasificarea unui vector de trăsături x prin setarea unei probabilități a priori π_1 ($\pi_1 = \pi_2 = \pi_1 = 0.5$) la clasa 1 și calculul valorii discriminante y :

$$y = (\mu_1 - \mu_2)^T \Sigma^{-1} x - \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 + \mu_2) + \log\left(\frac{\pi_1}{1 - \pi_1}\right). \quad (13)$$

(d) calcularea probabilităților posterioare:

$$P(1|x) = \frac{\exp(y)}{\exp(y) + 1}, P(2|x) = 1 - P(1|x). \quad (14)$$

(e) alegerea probabilității posterioare maxime.

Clasificatorul a fost testat prin metoda cross-validării. Rezultatele arată că o fonație susținută înregistrată într-un mediu controlat poate fi clasificată drept normală sau patologică cu o acuratețe de 89.1 %, în timp ce vocea telefonică poate fi clasificată drept normală sau patologică utilizând aceeași metodă cu o acuratețe de 74.2%. Caracteristicile referitoare la perturbațiile de amplitudine s-au dovedit a fi cele mai robuste în clasificare. Înregistrările au fost împărțite pe patru categorii și anume voci normale, cu patologie neuromusculară, fizică și mixtă (neuromusculare și fizice). A fost dezvoltat un clasificador separat pentru discriminarea grupului normal de fiecare din cele trei patologii. Rezultatele au arătat că patologiiile neuromusculare pot fi detectate de la distanță cu o acuratețe de 87%, cele fizice cu 78%, iar cele mixte cu 61%.

3.3. Alte corpusuri de voci patologice și instrumente aferente

Un corpus de voce patologică și normală în limba germană numit **COPAS** (Dutch Corpus of Pathological and Normal Speech) a fost dezvoltat de către G. Van Nuffelen et al. și a fost făcut public prin intermediul Uniunii Limbii Germane (http://www.inl.nl/tst-centrale/images/stories/producten/documentatie/copas_manual.pdf). Acesta a fost dezvoltat pentru a constitui baza unui instrument de evaluare a vocii patologice bazat pe diverse teste de inteligibilitate. Corpusul este alcătuit din opt categorii distincte de diagnostic și anume



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

voce normală, voce patologică, laringectomii, disartrii, glosectomii, defecte de auz, fisuri de palat și buză și defecte de articulare, totalizând un număr de 319 cazuri din care 122 normale.

Informațiile cu privire la vorbitor, de tip cod, vârstă, sex, numărul înregistrării, perioada de la prima înregistrare, tipul microfonului utilizat, subseturile de cuvinte și indicele de inteligibilitate aferent și tipul patologiei sunt stocate într-un fișier ExcelTM. Au fost înregistrate o serie de 11 teste, însă nu toți subiecții au participat la fiecare test și nu toate materialele text utilizate pentru teste au fost adnotate. Primul test a constat în înregistrarea a 50 de cuvinte în context CVC și calcularea inteligibilității la nivel de fonem, al doilea test a constat în citirea a 11 pasaje de text diferite, iar testele teste au constat în citirea unui text care conține un număr echilibrat de foneme în limba germană, evaluarea articulării prin numirea de desene, măsurarea ratei diadochokinetică care constă în repetarea fonemului /p/ timp de opt secunde, evaluarea tranzițiilor formantice prin repetarea alternativă timp de 6 secunde a vocalelor /i/ și /u/, în vorbire spontană și semispontană, și susținerea vocalei /a/ timp de 5 s.

Toate înregistrările au fost realizate într-o cameră clinică liniștită, cu două microfoane, unul de masă și unul cu cască și au fost salvate pe minidiscuri și apoi transferate pe un notebook. Transferul s-a făcut cu un editor audio gratuit Audacity®. Fiecare eșantion audio a fost salvat în format .wav pe 16 biți, cu o frecvență de eșantionare de 16 kHz. Toate eșantioanele audio înregistrate pentru un anumit test sunt salvate în același document. Adnotările au fost realizate cu programul PraatTM și au fost stocate în fișiere .TextGrid. Fiecare fișier audio și de adnotare a aceluiași vorbitor a fost numit cu un caracter de început semnificând patologia, un cod vorbitor și un string referitor la subtestul înregistrat.

J. F. Bonastre et al. (J. F. Bonastre et al., 2007) a dezvoltat tehnici complementare de evaluare a vocii patologice utilizând o selecție de înregistrări din baza de date a Departamentului ENT al Centrului Universitar Spitalicesc Timone din Marseille. Primul studiu s-a bazat pe un număr de 449 subiecți din care 391 pacienți cu voci patologice (308 femei și 20 bărbați) și 58 subiecți de control cu voce normală. Pacienții prezentau o serie de afecțiuni ale vocii întâlnite frecvent în practica clinică și anume 96 pacienți cu noduli vocali, 91 cu polipi, 65 cu disfonie paralică, 55 cu edem Reinke, 27 cu chiști, 24 cu disfonie funcțională, 19 cu displazie și 14 cu „Sulcus Glotidis”.

S-a realizat o evaluare perceptuală de către patru ascultători cu experiență după scala GRBAS, însă doar componenta G a fost utilizată pentru acest studiu. Analiza



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

obiectivă s-a efectuat cu ajutorul stației de lucru EVA[®]. Acest sistem permite realizarea de măsurători simultane ale parametrilor acustici și aerodinamici utilizând o piesă specială situată în dreptul gurii prevăzute cu microfon și pneumotahograf. Presiunea intraorală este măsurată cu ajutorul unui senzor de presiune montat în interiorul piesei.

Metoda constă în pronunția consecutivă a trei vocale susținute (/a/) care sunt analizate ulterior la nivel de F0, intensitate, jitter, shimmer, SNR (raportul semnal – zgomot), fluxul de aer oral și coeficientul Lyapunov. Au fost realizate trei sesiuni de măsurători pentru fiecare parametru și datele corespunzătoare pe fiecare parametru sunt o mediere a celor trei măsurători. S-au măsurat de asemenea presiunea subglotică în timp ce subiectul a pronunțat opt grupuri CV (/pa/) consecutive cu nivel normal de pich și claritate, minimul și maximul F0 și durata maximă a fonației.

Rezultatele au arătat că o combinație neliniară de șapte parametri a permis clasificarea identică cu cea a juriului a 82% eșantioane de voce. Al doilea studiu realizat de autori se referă la adaptarea tehnicilor de recunoaștere automată de vorbitor la evaluarea vocii patologice. Sistemul construit pentru această sarcină are la bază modele de mixturi gaussiene și este testat pe 80 de voci de gen feminin din care 20 sunt normale. Rezultatele obținute în acest studiu arată că gradul cel mai mare de confuzie apare între grade adiacente de patologie.

J. B. Tomblin (J. B. Tomblin, 2010) a dezvoltat corpusul EpiSLI care conține înregistrări de semnal vocal provenite de la copii de grădiniță cu defecte de limbaj. Acesta este creat în scopul utilizării pentru detecția automată de defecte de limbaj și voce patologică. Grupul target a fost un cluster stratificat în funcție de locul de rezidență și școală format din 6000 de copii vorbitori nativi de limbă engleză. Acesta a fost împărțit în trei categorii urban, suburban și rural. Părinții au semnat un consimțământ cu privire la participarea copiilor lor la acest studiu. Un procent de 2% (161) au refuzat să participe.

Copii au fost evaluați mai întâi în ceea ce privește performanța limbajului cu ajutorul unui test din 40 de puncte dezvoltat de autori, care a durat în jur de 10 minute. Copii din grupul de control au fost selectați din aceeași școală cu cei care nu au trecut primul test. Din 3877 copii care au fost selectați pentru faza de diagnostic doar 2084 au primit permisiunea părinților. 70 dintre aceștia sunt vorbitori de o a doua limbă și au fost excluși. Grupul final a fost alcătuit din 2009 copii vorbitori nativi de limbă engleză.

Scopul diagnosticării a constat în identificarea copiilor cu defecte de limbaj și a celor cu limbaj normal pentru grupul de control. Faza de diagnostic a inclus evaluarea auzului, a limbajului, a vorbirii, a înțelegerii și observații motorii și a durat aproximativ



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

două ore la fiecare copil. Baza de date creată este disponibilă la cerere fără cost, făcând parte din inițiativa Institutelor Naționale de împărțire a datelor.

4. Microcorpusuri și metode de detecție a vocilor patologice

G. Schlotthauer și M.E. Torres (G. Schlotthauer și M.E. Torres, 2006) a realizat un sistem de clasificare automat a vocilor disfonice și normale, care discriminează și între două tipuri diferite de disfonie și anume disfonia spasmodică (SD) de natură neurologică și disfonia datorată tensiunii musculare (MTD) de natură fiziologică, care poate fi tratată. Clasificatorul automat utilizează măsurători acustice ale vocalei susținute /a/ și instrumente de recunoaștere de pattern-uri bazate pe rețele neuronale.

Corpusul utilizat este alcătuit din 89 vorbitori împărțiți în două grupuri: 36 cu disfonie (15 cu MTD și 21 cu SD) și 53 cu voci normale care au fost instruiți să susțină vocala /a/ cel puțin 3 s. Înregistrările au fost realizate cu o frecvență de eșantionare de 22 kHz pe 16 biți. În acest studiu autorii au utilizat un singur eșantion de voce per pacient din care au extras opt parametri acustici care formează un pattern, incluzând perturbația pe termen scurt a frecvenței fundamentale și a intensității și zgomotul glotal. Toate pattern-urile au fost clasificate în trei categorii: normale SD și MTD utilizând o rețea neuronală cu perceptron multistrat. Pentru a reduce numărul de intrări ale rețelei neuronale autorii au realizat o analiză a componentelor principale (PCA) în urma căreia a rezultat că un număr de cinci componente a contribuit cu 99.5% la varianța setului de date. Astfel vectorii de intrare ai rețelei neuronale au fost reduși de la opt la șase. Dimensiunea nivelului ascuns a fost variată între 8 și 34 neuroni și s-au rulat 100 experimente pentru fiecare caz în vederea selecției celei mai bune topologii. Nivelul de ieșire are trei neuroni, câte unul pentru fiecare clasă (SD, MTD și normal), iar cel câștigător este cel cu valoare maximă.

Patologiile SD și MTD au fost recunoscute cu o acuratețe de 95.24%, respectiv 86.67%. În ceea ce privește vocile normale, acestea au fost recunoscute în procent de 100%. Rezultatele obținute de către autori utilizând o rețea neuronală cu perceptron multistrat cu nivelul ascuns alcătuit din opt neuroni indică o acuratețe de recunoaștere de 98.94%, depășind alte raportări din literatura de specialitate bazate pe clasificator de tip SVM (Support Vector Machine) (96.5%).

R. F. Pozo et al. (R. F. Pozo et al., 2009) au creat un microcorpus de voci patologice și normale în vederea detectării automate a pacienților cu apnee obstructivă în



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

somn (OSA), utilizând analiza semnalului vocal și tehnici de recunoaștere automată de voce și vorbitor. Acesta este alcătuit din 80 de subiecți de gen masculin dintre care 50% suferă de apnee severă și 50% sunt fie subiecți sănătoși fie cu forme ușoare de apnee. Subiecții din ambele grupuri au caracteristici fizice cum ar fi vârsta și indicele de masă corporală similare. Pentru grupul subiecților cu apnee s-au realizat două sesiuni de înregistrări, una înainte de a fi diagnosticați și una după câteva luni de tratament.

Corpusul conține înregistrări a patru propoziții în limba spaniolă care au fost repetate de trei ori de fiecare vorbitor. Acestea au fost selectate astfel încât să poată fi studiate anomaliile de rezonanță, de fonație și de articulare. S-a urmărit gradul de nazalizare a vocalelor cu și fără context nazal, a sunetelor sonore continue în vederea măsurării pattern-urilor de voce neregulată cauzate de oboseala mușchilor și a sunetelor sonore afectate de câteva foneme precedente care au locul de articulare velar (/g/).

Înregistrările au fost efectuate la o frecvență de eșantionare de 48 kHz într-o încăpere izolată acustic, iar aparatura utilizată a constat dintr-un laptop prevăzut cu placă audio și un microfon de tip cască prevăzut cu conversie A/D și transmiterea datelor digitale prin USB. În mod adițional pentru fiecare subiect au fost colectate și două imagini faciale pe fond alb cu iluminare controlată cu ajutorul unei camere digitale întrucât inspecția vizuală reprezintă un prim pas în diagnosticul pacienților suspecți de OSA.

În ceea ce privește anomaliile de articulare autorii au constatat că în cazul subiecților cu apnee distanța dintre formanții F2 și F3 ai vocalei /i/ este mai mare decât în cazul pacienților sănătoși. Pentru detectarea anomaliilor de fonație au fost calculați parametrii HNR (Raportul armonici – zgomot) și disperioicitatea care se referă la anomaliile semnalului generat de excitația glotală. O voce normală ar trebui să aibă un HNR mare și o disperioicitate mai mică decât o voce patologică. Acești parametri au fost calculați pentru a patra frază din baza de date deoarece conține mai multe sunete sonore. Se cunoaște faptul că pacienții cu apnee prezintă o nazalitate anormală prin prezența în voce unei extracomponente de frecvență joasă. Autorii au găsit o variație a nazalizării mai mică la subiecții cu OSA decât la cei normali. O ipoteză emisă de autorilor în urma acestui rezultat este aceea că pacienții cu apnee au un nivel de hipernazalitate global mai mare datorat disfuncției velofaringeale și astfel diferența între vocalele orale și nazale este mai mică deoarece și cele orale sunt nazalizate.

Luând în calcul ca trăsătură de bază nazalitatea autorii au testat un model de mixturi gaussiene (GMM) în vederea discriminării între pacienții cu OSA și cei sănătoși. Rezultatele au arătat că există diferențe semnificative între cele două grupuri în ceea ce



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

privește nivelele relative de nazalizare între diferite contexte lingvistice. Autorii au testat de asemenea puterea discriminativă a modelului GMM bazat pe tehnici de recunoaștere de vorbitor adaptate la detecția apneei severe obținând o rată de clasificare corectă de 81%

L. Salhi et al. (L. Salhi et al., 2010) au creat o nouă metodă de detectare a vocii patologice bazată pe o rețea neuronală multistrat (MNN). Algoritmul de procesare are la bază o tehnică hibridă care utilizează ca intrare energia coeficienților wavelet. Primul pas a constat în antrenarea supervizată a sistemului în vederea discriminării între voci normale și patologice și al doilea pas a constat în clasificarea patologiilor neurale și vocale de tip Parkinson, laringeale, Alzheimer etc. A fost utilizată o bază de date care conține voci normale și patologice colectate de la spitalul „Rabta-Tunis”. Inițial autorii au testat o metodă de determinare a vocii patologice și normale prin analiză formantică și de F0. Algoritmul a constat în extragerea F0 și a formațiilor prin metoda cepstrumului respectiv LPC și compararea cu valorile normale și împărțirea vocilor în două clase o clasă cu voci normale și o clasă cu voci patologice împărțită la rândul ei în două subclase cu patologie neurală și organică.

Algoritmul utilizat la metoda hibridă de clasificare constă în preaccentuarea și segmentarea semnalului cu o fereastră Hamming și aplicarea transformatei wavelet în vederea extragerii energiei coeficienților care vor constitui vectorul de intrare în rețeaua MNN. Sistemul propus de autori a fost dezvoltat în Matlab și este alcătuit din trei nivele: un nivel de intrare alcătuit dintr-un număr de neuroni egal cu numărul componentelor vectorului de trăsături, un nivel ascuns alcătuit din 15 neuroni și un nivel de ieșire alcătuit dintr-un singur neuron corespunzător deciziei patologic sau normal. Fiecare neuron din stratul ascuns este conectat la neuronii de intrare și nu există nici o conexiune între celulele aceluiasi strat. Funcția de activare utilizată este cea de tip sigmoid. Pentru extragerea coeficienților wavelet a fost utilizat un banc de filtre și ulterior energia fiecărui nivel a fost normalizată în raport cu energia totală a semnalului:

$$E_N(i) = \frac{E_i}{\sum_i E_i}, \quad (15)$$

unde E_i este energia la nivelul i .

Autorii au testat sistemul creat cu trei, cinci și șapte coeficienți de energie wavelet pe un total de 60 cuvinte pronunțate de vorbitori diferiți, din care 30 cu voce normală și restul cu voce patologică. Pentru antrenare au fost utilizate câte 20 de cuvinte din cele două



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOS DRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

seturi. Rata de clasificare corectă cu sistemul dezvoltat a fost între 80% și 100%, cea maximă fiind obținută în cazul utilizării a șapte coeficienți de energie.

J. Lee și M. Hahn (J. Lee și M. Hahn, 2009) propun o nouă metodă de detectare a vocii patologice prin analiză statistică de ordin înalt (HOS) în domeniul timp, bazată pe coeficienți de predicție liniară reziduali (LPC). Ca trăsături de bază au fost utilizați parametrii statistici skewness și kurtosis normalizați. Studiul se bazează pe un corpus format din 83 eșantioane de voce care conțin vocala susținută /a/, furnizate de către Societatea Logopedică și Foniatică Japoneză. Acestea au fost evaluate de către terapeuți de limbă conform scalei GRBAS și împărțite pe categorii în funcție de severitatea disfoniei. Drept clasificator a fost utilizat un arbore de decizie regresiv, iar cele mai bune rezultate cu o acuratețe de clasificare de 92.9% au fost obținute prin combinarea celor doi parametri statistici.

4.1. Aplicații în stomatologie

La ora actuală nu este disponibilă o bază de date de referință cu înregistrări de voci patologice cauzate de diverse patologii ale aparatului stomatognat care să poată fi utilizată în evaluarea calității vocii în stomatologie. S-au realizat studii pe grupuri mici de pacienți, iar bazele de date sunt utilizate doar de grupul de cercetare care le-au dezvoltat. Evaluarea eșantioanelor de voce s-a realizat fie prin metode subiective și anume percepție auditivă de către specialiști (K. Stevens et al., 2011), fie prin utilizarea de software-uri de analiză vocală (P. Jindra et al., 2002), fie prin dezvoltarea de metode automate proprii (Bocklet et al., 2010).

Tehnici de recunoaștere vocală automată au fost utilizate în evaluarea și diagnosticul automat al pacienților edentați, cu dantură insuficient fixată. Astfel **Bocklet et al.** (Bocklet et al., 2010) au testat trei sisteme de clasificare diferite, în vederea determinării dacă dantura unui pacient edentat este fixată corect. Baza de date utilizată în acest studiu este alcătuită din 13 pacienți edentați cu vârste cuprinse între 54 și 74 ani. Aceștia au fost diagnosticați ca având dantură insuficientă și li s-au construit danturi adiționale. Tehnica a constat în înregistrarea pacienților fără dantură, cu dantură insuficientă și ulterior cu dantură suficientă. Materialul de înregistrat este un text standard în limba germană, echilibrat fonetic, care conține 108 cuvinte. Înregistrările au fost realizate cu o frecvență de eșantionare de 16 kHz și cuantizate pe 16 biți. Autorii au creat un vector de trăsături de dimensiune 24 care conține energia pe termen scurt, 11 coeficienți mel cepstrali și



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

derivatele delta de ordin 1 ale acestora, pe care l-au utilizat la ASR și la două din cele trei sisteme de clasificare.

Sistemul de clasificare de bază are un vector de trăsături unidimensional reprezentat de acuratețea de recunoaștere a cuvintelor (WA) cu un sistem de recunoaștere vocală automată (ASR). În urma testării ASR pe cele trei seturi de înregistrări (fără dantură, cu dantură insuficientă și suficientă) s-au obținut rezultate diferite pentru WA (60.06%, 64.35%, respectiv 70.91%). Din acest motiv autorii au utilizat WA ca parametru pentru discriminarea celor trei clase de vorbitori cu un sistem de clasificare de tip SVM (Support Vector Machine).

Al doilea sistem de clasificare utilizează tehnica DTW (Dinamic Time Warping) pentru a extrage vectorul de trăsături care codează distanțele în spațiul MFCC dintre înregistrarea unui vorbitor de test cu dantură și un vorbitor de referință fără patologie. Vectorul de trăsături de dimensiune 2314 este ulterior clasificat cu un SVM. Al treilea sistem de clasificare utilizează ca trăsături media vectorilor unui model de mixturi gaussiene (GMM) de dimensiune 128. Supervectorul GMM de dimensiune 3072 a fost ulterior clasificat cu un SVM.

Rezultatele obținute în urma clasificării celor trei seturi de înregistrări relevă o acuratețe a clasificării de 61.5%, 73.1% respectiv 80.1% cu sistemele de clasificare bazate pe WA, DTW, respectiv GMM.

K. T. Bressmann et al. (K. S. Bressmann et al., 2011) au studiat impactul expanderilor palatali rapizi (RPE) asupra articulării sunetelor. Aceștia sunt prevăzuți cu atașamente cimentate la dinți și au o porțiune care acoperă palatul. Datorită poziției pe care o ocupă și dimensiunii pot afecta vorbirea. Scopul studiului a constat în evaluarea perturbației și adaptării vocii de-a lungul timpului, relaționate la aplicarea expanderilor. Au fost înregistrați 22 de pacienți cu vârste cuprinse între 9 și 19 ani, care urmau să fie tratați cu RPE. Înregistrările s-au realizat în șase etape: înainte de plasarea expanderilor, după plasare, în timpul expansiunii, în timpul retenției, după eliminarea expanderului și după patru săptămâni de la eliminare.

Materialul înregistrat a fost alcătuit din 35 de propoziții din care trei au fost alese pentru analiză. Eșantioanele de voce au fost evaluate perceptiv de către 10 ascultători specialiști, iar formanții vocalei /i/ și spectrul fricativelor /s/ și /ʃ/ au fost măsurate cu software-ul de analiză vocală Wavesurfer™.



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

Rezultatele au arătat că după plasarea expanderului vocea s-a degradat, în timp s-a îmbunătățit și după eliminarea expanderului a revenit la normal. Primul formant al vocalei /i/ a crescut și al doilea formant a scăzut în frecvență indicând centralizarea vocalei. Formanții au revenit la valorile dinaintea aplicării expanderului în timpul tratamentului. Pentru fricativele /s/ și /ʃ/ rapoartele frecvențelor joase și înalte au arătat prezența unei distorsiuni care a dispărut după îndepărtarea expanderului.

Rezultatele au arătat că media spectrală a scăzut iar skewness-ul a devenit pozitiv. Analiza repetată a arătat o varianță semnificativă a măsurătorilor acustice. Vocea a fost alterată când expanderul a fost aplicat prima dată, s-a îmbunătățit pe parcurs și a revenit la normal după îndepărtarea lui.

Gradul de îmbunătățire a articulării sunetelor la pacienți edentați cărora li s-a aplicat dantură parțială sau completă a fost studiat și de P. Jindra et al. (P. Jindra et al., 2002). Au fost investigați parametrii semnalului vocal pentru un set de silabe test pronunțate de câteva ori de către un grup de zece pacienți dintr-o clinică stomatologică dintre care cinci bărbați și cinci femei cu vârste cuprinse între 58 și 81 ani. Silabele analizate au fost : "si, že, ša, ta, da, ře, ra, fa, vi, pa, bo", pronunțate de către pacienți cu și fără proteză parțială sau completă aplicată pe maxilar sau mandibulă. Odată cu înregistrările a fost realizată și o anamneză a pacienților care conținea probleme subiective întâmpinate de pacienți în vorbirea fără proteză. Analiza Fourier a înregistrărilor realizată cu programul "SoundForge 5.0" a arătat că există diferențe în ceea ce privește parametrii acustici ai consoanelor dar și a vocalelor pronunțate de pacienți fără proteze, în special la silabele care conțin consoane alveolare.

Analiza silabei "si" a arătat o scădere în limita frecvențelor superioare și o descreștere a maximului spre frecvențele joase când articularea s-a realizat fără proteză la toți pacienții examinați. În mod similar s-a constatat o schimbare în pronunția consoanei /f/ din silaba "fa" la șase pacienți. La articularea fără proteză au dispărut total frecvențele mai mari de 3.5 kHz, iar la articularea cu proteză a apărut o creștere în zona frecvențelor de peste 13 kHz. La pronunțarea cu proteză formanții sunt mai clar delimitați și frecvențele în jur de 3.5 kHz sunt mai puțin frecvente. În cazul consoanei /r/ pronunțată fără proteză s-a constatat de asemenea schimbări în spectru la toți pacienții. În cazul câtorva pacienți vocalele în special /i/ și /u/ își schimbă frecvența fundamentală și arată apariția unui zgomot manifestat ca o delimitare neclară a formanților.

Autorii concluzionează că factorii cei mai importanți în construirea danturii din punct de vedere fonetic sunt overbite, overjet, înălțimea palatului, grosimea materialului din



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

care este confecționat palatul, poziția incizorilor și modelarea palatului dur. Acest studiu arată că silabele care conțin consoane alveolare frontale, labiodentale și laterale /s/, /f/, /v/, /r/ și vocale înalte /i/, /u/ sunt indicatori în evaluarea articulării sunetelor în stomatologie.

Un alt studiu referitor la evaluarea vocii pacienților tratați cu proteză totală care poate fi scoasă a fost realizat de **F. Stelzle et al.** (F. Stelzle et al., 2010). Autorii au utilizat un ASR pentru discriminarea între un grup de 28 de pacienți și un grup de control alcătuit din 40 de subiecți fără patologii de dentiție prin utilizarea parametrului WA:

$$WA = \frac{C - W}{R} \times 100, \quad (16)$$

unde C este numărul de cuvinte recunoscute corect, W numărul de cuvinte recunoscute greșit și R numărul total de cuvinte analizate. S-au realizat două înregistrări per pacient cu proteză și fără proteză, iar subiecții examinați au fost împărțiți în funcție de unele simptome specifice în 18 pacienți cu proteză completă suficientă și 10 pacienți cu proteză completă insuficientă.

Rezultatele evaluării vorbirii de către ascultători specialiști și cu ASR au arătat o corelație de 0.71. Acuratețea de recunoaștere a cuvintelor a fost semnificativ redusă la pacienții edentați (55.42%) comparativ cu grupul de control (69.79%). În urma acestui studiu s-a constatat că după pierderea completă a danturii calitatea vocii este semnificativ alterată, iar ASR s-a dovedit a fi un instrument util și ușor de aplicat în evaluarea automată a vorbirii într-un mod standardizat.

5. Cazul SRoL (Sunetele Limbii Române) și extensie în limba franceză

Site-ul web intitulat ”**Sunetele limbii române**” este o colecție care include mii de înregistrări de vocale, consoane, diftongi, propoziții cu voci emoționale și dialectale și o arhivă de sunete gnatofonice și gnatosonice dezvoltată de (H.N. Teodorescu et al., 2005-2007). Înregistrările sunt adnotate și documentate conform metodologiei și protocoalelor concepute de autori. Site-ul include de asemenea documentație despre limba română, tehnologia vorbirii și instrumente dezvoltate de echipa SRoL pentru analiză vocală. Resursele și instrumentele sunt dedicate învățării virtuale a limbii române, aplicațiilor fonetice, de tehnologie a vorbirii, medicale și de reabilitare a vorbirii și fac parte din Rețeaua Europeană de resurse pentru limbaj CLARIN. SRoL este singura bază de date disponibilă pe internet, de voci emoționale în limba română care conține peste 1500 de înregistrări adnotate în diverse formate (.wav., .ogg, .txt) eșantionate cu 22 kHz și precizie



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IASI

de 16 sau 24 biți. Conform (S.M. Feraru et al., 2010) site-ul conține șase secțiuni principale:

- Pronunții standard de vocale susținute, consoane, diftongi, cuvinte și propoziții scurte în limba română destinate învățării corecte a pronunției limbii române și cercetării statistice în fonetică;
- Construcții sintactice speciale cum ar fi propoziții cu subiect dublu și apozitie destinate cercetării;
- Voci emoționale (M. Feraru și H.N. Teodorescu, 2008);
- Comparație între vocea normală și cea sintetică (H.N. Teodorescu și M. Feraru, 2008);
- Pronunții dialectale;
- Microcorpus de sunete gnatofonice și gnatosonice (H.N. Teodorescu și M. Feraru, 2007), (H. N. Teodorescu și A. Untu, 2010), (A. Untu și H.N. Teodorescu, 2011).

Instrumentele dezvoltate de echipa SRoL se referă la extragerea de pattern-uri din semnalul vocal și de calcul a traseelor frecvenței fundamentale și a formanților F1, F2 și F3 (M. Zbancioc, 2006). Pe lângă programe executabile autorii pun la dispoziție și o descriere a acestora în vederea înțelegerii funcționării lor.

5.1. Microcorpus de sunete gnatofonice în limba română

Microcorpusul de sunete gnatofonice dezvoltat de echipa SRoL, prezentat în (H.N. Teodorescu și M. Feraru, 2007), și apoi dezvoltat în (H. N. Teodorescu și A. Untu, 2010) și (A. Untu și H.N. Teodorescu, 2011) conține la ora actuală 29 înregistrări provenind de la 24 vorbitori din care 10 de gen feminin și 14 de gen masculin. Pe SRoL sunt disponibile doar 19 din cele 29 înregistrări restul fiind protejate. Din setul total de vorbitori 16 sunt fără patologii de dentiție și 8 cu diverse patologii. În continuare sunt detaliate statistica corpusului la momentul respectiv și metodologiile utilizate și prezentate în (H. N. Teodorescu și A. Untu, 2010).

Metodologia de culegere a datelor

Metodologia, protocolul de înregistrare și protocolul de documentare sunt cele prezentate pe sit-ul "Sunetele limbii române". Subiecții au fost informați anterior înregistrărilor de obiectivele proiectului, fiind asigurați de confidențialitatea datelor personale. Subiecții au semnat un consimțământ informat în conformitate cu protocolul de protecție a subiecților umani și cu principiile etice ale cercetărilor care implică ființa umană



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

existente la nivel național și internațional. Cercetarea pe subiecți umani privind analiza sunetelor gnatofonice și gnatosonice a fost aprobată de către consiliul Facultății de Electronică, Telecomunicații și Tehnologia Informației din cadrul Universității Tehnice "Gheorghe Asachi" din Iași (H.N. Teodorescu et al., 2005-2007).

Metodologia de înregistrare

Sunetele gnatofonice au fost înregistrate cu ajutorul unui microfon prevăzut cu căști, A4 Tech Stereo HS-60, având caracteristicile: frecvență de răspuns 20 Hz-20 kHz, impedanță: $U=3V$, $R=1,5k\Omega$, sensibilitate: $-58dB\pm 2$. Placa de bază a calculatorului pe care au fost efectuate înregistrările este Sony Vaio MBX-189 având încorporată o placă de sunet Intel® High Definition Audio compatible 3D audio (Direct Sound 3D support) cu următoarele caracteristici: procesor de semnal 44-kHz / 16-bit stereo CD quality, mod de ieșire a sunetului 8 canale, 192 kHz / 32 bit, standard Intel HD audio.

Pentru înregistrări s-a utilizat programul GoldWave™, versiunea 5.54, la o frecvență de eșantionare de 22050 Hz cu atributele PCM signed (16-24 bits mono). Culegerea de semnal vocal se realizează în condiții de zgomot redus (amplitudinea zgomotului trebuie să fie mai mică cu cel puțin 20 dB decât amplitudinea frecvenței fundamentale). Conform metodologiei de înregistrare de pe situl SRoL, se recomandă ca poziția microfonului să fie mai jos de gură, aproximativ în dreptul bărbiei (la câțiva centimetri de aceasta), iar distanța de la bărbie să fie aproximativ egală cu distanța până la buze (Teodorescu et al., 2005-2007).

Metodologia de documentare

Fișele subiecților "Profil vorbitor", "Chestionar Patologie Vocală și Factori Obiectivi" și "Fișă dinți subiect" au fost completate de către al doilea autor. Pentru "Fișă dinți subiect" s-au luat în considerare informații ce pot fi cunoscute de către subiectul înregistrat sau vizualizate de către persoana care completează fișa, de tip prezență / absență dinte, plombă, implant, coroană, punte. Am considerat elemente de interes și plombele întrucât acestea dacă nu sunt realizate corect pot modifica structura dintelui având repercursiuni asupra spațiului interdentar, în consecință și asupra vorbirii (spațierile dintre incisivi pot produce în vorbire efectul de "fonfăit"). Pentru a valida aceste idei este necesară o analiză statistică pertinentă (cel puțin 20 de subiecți). Un exemplu de fișă ce conține informații referitoare la patologiile arcadei superioare și inferioare este dat în



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

figura 1. Subiectul prezintă opt plombe simple la nivelul maxilarului și opt plombe simple la nivelul mandibulei, care afectează primii molari de pe fiecare hemiarcadă.

Fișă dinți vorbitor															
Cod: 1234															
Hemiarcada dreaptă superioară								Hemiarcada stângă superioară							
1.8	1.7	1.6	1.5	1.4	1.3	1.2	1.1	2.1	2.2	2.3	2.4	2.5	2.6	2.7	2.8
	Pb	Pb	Pb	Pb							Pb	Pb	Pb	Pb	
4.8	4.7	4.6	4.5	4.4	4.3	4.2	4.1	3.1	3.2	3.3	3.4	3.5	3.6	3.7	3.8
	Pb	Pb	Pb	Pb							Pb	Pb	Pb	Pb	
Hemiarcada dreaptă inferioară								Hemiarcada stângă inferioară							
Legendă: Dinte lipsă - X Implant - I Plombă - Pb Punte - Pu Coroană - C															

Figura 1: Exemplu de fișă dinți subiect.

Listă cuvinte

Pentru realizarea corpusului de sunete gnatofonice s-a utilizat setul de cuvinte stabilit anterior de primul autor (©), cuvinte ce conțin consoanele *s, f, ș, v, z, j*. Analiza comparativă a pronunției consoanelor *s, f, ș*, și a celor semi-vocalice *v, z, j*, furnizează informații despre afectarea dentiției.

Cuvintele utilizate pentru înregistrările gnatofonice sunt: vată / fată, var / far, vuiet (pronunțat vvvvuiet) / vuiet (pronunțat normal, scurt, vuiet) / fffffui / fui, vvvvvalet / valet / fffffaieton / faieton, vecin / fecior, vvvvvânt / vânt / fffffân / fân, vvvvvine / vvvvvine / vvvvine / fffffine, vine / fine, vehement / ferment, vierme / fierbe, vâjâit / gâjâit / sâsâit / fâsâit, bâzâie / zâzâie, bâzzzzzâie / zâzzzzzâie, fâșâit / fâlfâit, vâjâie / fâlfâie / sâsâie / fâșâie, vâjjjjjâie / fâfffffâie / sâsssssâie / fâșșșșșâie, vâjjjjjâit / sâsssssâit / fâsssssâit / fâșșșșșâit / fâfffffâit, vorbit / fortuit / sortit, suit / vuit.

Scopul cercetării noastre este de a identifica și analiza modificările ce apar în dinamica formașilor / pseudo-formașilor consoanelor fricative pronunțate de către două clase de vorbitori: o clasă martor formată din vorbitori cu dentiție normală (neafectată de patologii ale aparatului stomatognat) și o clasă de vorbitori cu defecte de dentiție. Setul de cuvinte întocmit în acest scop include cuvinte ce conțin consoanele *v, f, s, ș, j, z* în context consoană-vocală (CV) sau vocală-consoană-vocală (VCV).

Pronunțiile celor 56 de cuvinte pentru fiecare din cei 10 vorbitori sunt înregistrate într-un singur fișier .wav care este salvat cu un nume mnemonic, în formatul "cod_subiect_sex" (ex. 1234_m). Fișierele au fost filtrate și re-salvate în format .wav și în format .ogg (16 și 24 biți). Înregistrările au fost ascultate pentru o analiză perceptuală a zgomotului, iar pronunțiile cu trunchiere (saturație) în amplitudine au fost eliminate.



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

Discuția corpusului

Corpusul de sunete gnatofonice este disponibil pe sit-ul "Sunetele limbii române", în cadrul arhivei pentru aplicații de gnatofonie și gnatosonie (Figura 2). În paralel am realizat și un corpus de sunete gnatosonice pe care îl vom prezenta într-o lucrare ulterioară.

Arhiva pentru aplicații de gnatofonie și gnatosonie

Metodologia de culegere a semnalelor gnatofonice este identică cu cea de culegere de semnal vocal.

Cuvintele utilizate pentru înregistrările gnatofonice sunt alese astfel încât să se poată analiza comparativ modificările de siflante, fricative și de consoane semi-vocalice.

Cuvintele utilizate pentru înregistrările gnatofonice sunt: vată / fată; var / far; vuiet (pronunțat wwwuiet) / vuiet (pronunțat normal, scurt, vuiet) / fui / vaiet (pronunțat wwwaiet) / vaiet (pronunțat vaiet) / faieton / vecin / fecior / vânt (pronunțat wwwânt) / vânt (pronunțat vânt) / fân / wwwine, wvine, wvine / vine / fine / vehement / ferment / vierme / fierbe / vâjăit / wwwwâjăit / wwwâjăie / ffffâșșșșăie / ffffâșșșșăit / fâșăit / sâșăit / sssssâșșșșăie / găjăit / zăzăie / bââzzzzăăăie / báz&226;ie.

Un număr de 10 înregistrări gnatofonice provenite de la 5 subiecți de gen feminin și 5 subiecți de gen masculin sunt adnotate, iar fișierele TextGrid și wav corespunzătoare se pot vizualiza (accesa) la adresa (http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/ro/sunete_gnatofonice.htm) sau aici.

Fișă vorbitor	Gnatofonie		Gnatofonie - filtrate		Gnatosonie		Gnatosonie - filtrate	
	wav	ogg	wav	ogg	wav	ogg	wav	ogg
Fișă vorbitor #394715	394715_f	394715_f	394715_f	394715_f				
	394715_f_v1	394715_f_v1	394715_f_v1	394715_f_v1				
	394715_f_v2	394715_f_v2	394715_f_v2	394715_f_v2				
Fișă vorbitor #231515	231515_m	231515_m	231515_m	231515_m				
Fișă vorbitor #280269	280269_m	280269_m	280269_m	280269_m				
Fișă vorbitor #20048					20048_f_v0	20048_f_v0	20048_f_v0	20048_f_v0
Fișă vorbitor #55555	55555_f	55555_f	55555_f	55555_f	55555_f	55555_f	55555_f	55555_f
	55555_f_v1	55555_f_v1	55555_f_v1	55555_f_v1				
Fișă vorbitor #1234	1234_m	1234_m	1234_m	1234_m	1234_m	1234_m	1234_m	1234_m
Fișă vorbitor #1357	1357_m	1357_m	1357_m	1357_m	1357_m	1357_m	1357_m	1357_m
Fișă vorbitor #2001	2001_f	2001_f	2001_f	2001_f	2001_f	2001_f	2001_f	2001_f

Figura 2: Arhiva pentru aplicații de gantofonie și gnatosonie.

Statistica înregistrărilor și patologiilor

Am realizat înregistrări de sunete gnatofonice provenind de la cinci subiecți de gen masculin și cinci subiecți de gen feminin, cu vârste cuprinse între 21 și 46 ani, majoritatea cu studii superioare, fără afecțiuni respiratorii, laringeale, neurologice, sau psihologice.

Corpusul realizat are la bază înregistrări provenind de la nouă subiecți fără defecte majore de dentiție și fără deficiențe de vorbire detectabile prin percepție auditivă și un subiect cu edentație majoră (13 molari lipsă). Starea dentiției subiecților înregistrați este ilustrată în Tabelul 2.



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

Tabel 2. Starea dentiției subiecților înregistrați

Cod subiect	Sex	Vârș-tă (ani)	Nr. Pb (plombă) / localizare		Nr. X (lipsă dinte) / localizare		Nr. C (coroană) / localizare		Nr. Pu (punte) / localizare		Tip voce / patologii
2404	M	30	1	primul incisiv HSS	-	-	-	-	-	-	Profesională, fără patologii
2001	F	26	7	primul molar HDI, al 2-lea și al 3-lea molar HDS și HSS, al 2-lea molar HSI, al 3-lea molar HDI	-	-	1	al 2-lea molar HDI	-	-	Neprofesională, fără patologii
4312	F	29	3	primii incisivi HSS și HDS, al 3-lea molar HDS	1	al 3-lea molar HDI	2	al 2-lea molar, HSS, al 3-lea molar HSI	-	-	Profesională, fără patologii
1357	M	30	2	canin HSS, al 2-lea molar HSS	2	ultimul molar HDI și HSI	2	al doilea molar HDI și HSI	1	molarii 1, 3, HDS	Neprofesională, fără patologii
2202	F	26	3	al 2-lea molar HSS, al 3-lea molar HDS, al 4-lea molar HSI	1	al 3-lea molar HSI	-	-	-	-	Neprofesională, fără patologii
1234	M	30	16	primii 4 molari de pe fiecare hemiarcadă	-	-	-	-	-	-	Neprofesională, fără patologii
3371	F	21	4	primul incisiv HSS, al 2-lea molar HSS, al 4-lea molar HDI, al 3-lea molar HSI	2	al 3-lea molar HDI și HSI	-	-	-	-	Neprofesională, fără patologii
3298	M	30	3	2 pe primii incisivi de pe HSS și HDS, al 3-lea molar HDS	-	-	2	al 2-lea molar HDI și HSI	-	-	Profesională, fără patologii
01321	F	30	2	canin și primul molar HDI	-	-	1	primul molar HSS	-	-	Neprofesională, fără patologii
5343	M	46	1	ultimul molar HSI	13	incisivii 1, 2 HSS, primii 3 molari HDS, ultimii 3 molari HSS, molarii 2,3,4 HSI, molarii 3, 4 HDI	-	-	-	-	Neprofesională, fără patologii

Metodologia de analiză

Metoda de prefiltrare

Toate fișierele au fost prefiltrate cu un filtru trece bandă cu frecvențele de tăiere de cca. 70 Hz și 7 kHz și cu atenuare de 100 dB la frecvențele de 50 Hz și 12 kHz (în afara benzii de trecere), setat în utilitarul GoldWaveTM (Figura 3). Filtrul se poate aplica o dată sau, repetat, de două ori (echivalent, două filtre înseriate) astfel ca nivelul final de zgomot la 50 Hz să fie cu cel puțin 20 dB mai mic decât nivelul la frecvența fundamentală (F0). În cazul înregistrărilor curente, filtrul s-a aplicat o singură dată.

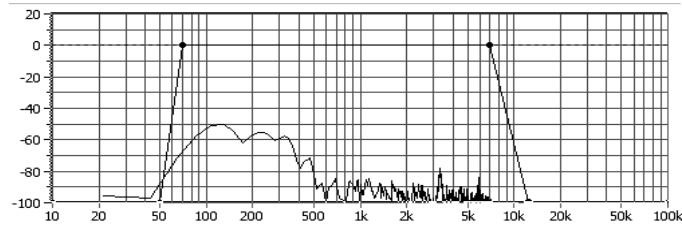


Figura 3: Filtru trece bandă de la 70 Hz la 7 kHz setat în utilitarul GoldWave™.

Metodologia de adnotare

Pentru o analiză ulterioară a diferențelor / similitudinilor ce apar la nivelul consoanelor fricative, am segmentat și adnotat cele 10 fișiere la nivel de fonem / silabă / cuvânt cu ajutorul utilitarului Praat™. Am realizat împărțirea fișierelor în subfișiere ce conțin consoanele *f*, *s*, *ș*, *v*, *z*, *j*, în context CV și VCV. Au rezultat câte nouă sau 11 subfișiere la fiecare subiect, numite astfel: *vf_CV(a)_cod_subiect_sex* – consoanele *v* și *f* în context CV, *V=vocala a*, *jsfshz_VCV(a-)_cod_subiect_sex* – consoanele *j*, *s*, *f*, *ș*, *z* în context VCV, *V=vocala â*, *vfs_CV(o)_cod_subiect_sex* – consoanele *v*, *f*, *s* în context CV, *V=vocala o*, etc. Pentru *ș*, *â* și *ă* am utilizat notațiile *sh*, *a-* și *a+*.

În cadrul procesului de adnotare am luat în considerare și pauzele intravorbire (ce apar în rostirea de silabă, cuvânt) notate cu simbolul \$, pauzele intervorbire (dintre cuvinte) fără notație (blanc) și pauzele care nu se aud, notate cu simbolul %. Adnotarea s-a realizat ținând cont de percepția auditivă, prezența / absența frecvenței fundamentale (ce evidențiază caracterul vocalic, sonor, cu activarea corzilor vocale), forma de undă, energia și modificările ce apar pe spectrogramă de la un fonem la altul.

Metodologia de analiză formantică și temporală

Cu ajutorul utilitarului Praat™ am extras F0, formații / pseudo-formații F1, F2, F3, F4 și timpii corespunzători fiecărui fonem și fiecărui cuvânt în parte.

Analiza a constat din următoarele etape: i) Segmentarea, în vederea analizei modificărilor duratelor fonemelor. Segmentarea este efectuată la nivelurile fonem, silabă și cuvânt, pentru a se determina momentele de timp corespunzătoare granițelor dintre foneme și cuvinte. ii) Generarea fișierelor .txt cu instrumentul Praat, cu valorile formațiilor în evoluția lor; iii) Importarea fișierelor într-un utilitar (de ex., Excel™) și separarea manuală a fonemelor în funcție de limitele temporale corespunzătoare fiecărui fonem; iv) Analiza dinamicii și analiza statistică a formațiilor pe segmentele de interes, comparativ între subiecții cu afecțiuni și tratamente ale aparatului stomatognat și subiecții sănătoși. Am analizat dinamica formațiilor / pseudo-formațiilor consoanelor *v*, *f*, *s* din cuvintele *vată*,



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

fată, sâsâit, fâsâit, sâsssâit, fâsssâit. Acestea s-au aplicat la patru vorbitori (doi de gen feminin și doi de gen masculin).

5.2. Microcorpus de sunete gnatofonice în limba franceză

Complementar microcorpusului de sunete gnatofonice în limba română am dezvoltat un microcorpus de sunete gnatofonice în limba franceză raportat în (Untu (Hulea), A., et al., 2011) care conține 51 de înregistrări de cuvinte care încep cu fricativele /v/ și /f/ în context CV (consoană-vocală) provenind de la 17 vorbitori de gen feminin și 34 vorbitori de gen masculin. Subiecții sunt studenți ai Institutului Universitar de Tehnologie din Angouleme și au vârste cuprinse între 18 și 24 ani, fără patologii de tip respirator, laringeal, neurologic sau psihologic care ar influența vocea. Microcorpusul este alcătuit din studenți cu dantură normală și subiecți cu patologii dentare. Patru respectiv opt subiecți de gen feminin au malocluzie clasa I respectiv II și 15 subiecți de gen masculin au malocluzie clasa I.

Pentru crearea microcorpusului în limba franceză am urmat aceeași metodologie ca și în cazul microcorpusului în limba română. Am creat o fișă de dentiție în limba franceză care include arcadele dentare, litere simbol pentru patologii dentare posibile și ilustrarea tipurilor de malocluzie posibile. Studenții au marcat pe arcade în dreptul fiecărui dinte/molar și pe pozele ilustrative tipul de patologie respectiv tipul de malocluzie prezentă. Câteva exemple de cuvinte din lista creată pentru limba franceză sunt: “varier” (a varia) / “fatiguer” (a obosi), “vol” (zbor) / “folle” (nebun).

Înregistrările le-am realizat cu sistemul de achiziție ”Harmonie” produs de ”01db” utilizând software-ul dBBATI32. Am utilizat un microfon de tip G.R.A.S. 40 AC, cu un răspuns în frecvență de 3.15 Hz - 45 kHz at ± 2 dB și sensibilitate de 12.5 mV/Pa. Înregistrările au fost realizate în timpul unui stagiu de cercetare în cadrul Laboratorului de electroacustică și acustică a mediului al IUT ANGOULEME.

6. Discuții și concluzii

Corpusurile de semnal vocal sunt indispensabile în vederea dezvoltării și testării instrumentelor de recunoaștere vocală, sinteză vocală, a sistemelor de recunoaștere automată de voce patologică, de detecție a patologiiilor vocale induse de dentiție sau de evaluare a calității aplicării unui tratament în stomatologie.



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

Dezvoltarea unui corpus de semnal vocal presupune respectarea unei metodologii bine puse la punct și adaptată în funcție de aplicație. De o mare importanță sunt parametrii utilizați la realizarea înregistrărilor și anume frecvența de eșantionare și precizia. Semnalul vocal se întinde pe o bandă de la câteva zeci de hertzi la 20 kHz. Deși se consideră că limita superioară din semnalul vocal a intervalului în care este conținută cea mai importantă informație este la 8 kHz, o eșantionare cu 16 kHz ar altera consoanele fricative care au un spectru bogat în frecvențe înalte de peste 8 kHz. Conform teoremei Shannon care impune utilizarea unei frecvențe de eșantionare de cel puțin de două ori mai mari decât frecvența maximă a semnalului, pentru înregistrări de calitate este necesară o frecvență de eșantionare de peste 20 kHz, fapt ce nu este respectat la zece din corpusurile prezentate și sintetizate în Anexă cum ar fi (J.S. Gorofolo et al., 1986), (K. Shobaki et al., 2007), (O. Salor et al., 2007) etc. De asemenea înregistrările trebuie filtrate pentru eliminarea zgomotelor introduse de aparatură, alimentare, cum ar fi zgomotul de 50 Hz. În nici unul din materialele referitoare la corpusurile prezentate cu excepția SRoL nu este specificat dacă înregistrările au fost filtrate. Din punct de vedere statistic, corpusurile trebuie să fie compuse dintr-un număr cât mai mare de vorbitori de ambele genuri și un set diversificat de foneme, cuvinte, propoziții. Corpusuri limitate din punct de vedere al numărului de vorbitori sunt (A. Stan et al. 2011) și (J. Matousek și J. Romportl, 2007) care au un singur vorbitor. Deși aplicația de bază este sinteză de voce, sistemele dezvoltate ar trebui testate pe mai multe voci.

Informațiile demografice cu privire la subiecții înregistrați și anamneza în cazul constituirii corpusurilor de voci patologice sunt de asemeni importante. Un îndeajuns a câtorva corpusuri și instrumente dezvoltate din Anexa este acela că în articolele aferente nu sunt specificate întotdeauna parametrii de înregistrare (frecvență de eșantionare, precizie), număr de vorbitori, informații despre vorbitori, adnotare și limba utilizată cum ar fi în (J. F. Bonastre et al., 2007), (J. B. Tomblin, 2010). În cazul dezvoltării unui corpus de voci patologice este necesară dezvoltarea complementară a unui corpus de voci normale de control pentru utilizarea la sisteme de clasificare de voce sau detecție de voce patologică. (F. Bocklet et al., 2010) și (K. T. Bressmann et al., 2011) au evaluat articularea sunetelor la pacienți edentați înainte de tratament și după tratament fără a compara cu persoane fără patologii de dentiție.

La ora actuală nu există un corpus de referință statistic valabil cu voci provenite de la pacienți cu diverse patologii de dentiție care să poată fi utilizat de specialiști în domeniu în dezvoltarea de sisteme de diagnostic automat, de evaluare a vocii înainte, în timpul și după



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

tratament sau de detectare a protezelor incorect fixate sau nepotrivite. Câteva limitări ale SRoL (Teodorescu et al. 2005-2007) sunt un număr redus de vorbitori în special la înregistrările de vocale susținute și consoane, necesitatea utilizării unei frecvențe de eșantionare mai mari în cazul sunetelor gnatofonice pentru a nu tăia componentele spectrale ale consoanelor fricative de interes și un număr insuficient de înregistrări cu voci patologice. Extinderea corpusului de vocale susținute ar lărgi aplicabilitatea site-ului în domeniul testării sistemelor de recunoaștere automată a vocilor disfonice, putând fi utilizat ca material de control.

7. Direcții viitoare

Ca direcții viitoare ne propunem dezvoltarea microcorpusului de sunete gnatofonice în limba română atât cu înregistrări de voci normale cât și cu înregistrări de voci patologice induse de defectelor de dentiție, în vederea constituirii unei baze de cunoștințe care să înglobeze caracteristicile acusto-fonetice ale sunetelor gnatofonice.

Alte scopuri sunt identificarea de trăsături specifice vocilor patologice utile în dezvoltarea de sisteme de clasificare de voce, design-ul teoretic și dezvoltarea unui sistem de detectare automată a patologiilor de dentiție sau de evaluare a vocii persoanelor protezate.



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ĂȘACHI"
DIN IAȘI

Autori	Nr. Vb.	Limba	Material text	Adnotare/Transcripție	Frecv. eșant.	Precizie	Tip înregistrări	Instrum.	Aplicații	Comercial
J.S. Gorofolo et al., 1986-TIMIT	630	6 dialecte lb. engleză	10 propoziții citite/subiect	DA	16 kHz	16 biți	microfon	---	-recunoaștere vocală; studii acustico-fonetice	DA
J. Wright, 2005	---	lb. engleză	Silabe-2000 comune, 20 unice/subiect	---	16 kHz	---	microfon	---	-modelare pronunție; identificarea, modelarea și procesarea de limbă	DA
K. Shobaki et al., 2007	1100	lb. engleză copii	Vorbire spontană și citită	---	16 kHz	16 biți	microfon	---	-recunoaștere vocală, dezvoltarea limbii copiilor fără auz, învățării limbii	DA
M. Liberman et al., 2002	---	lb. engleză, actori	Enunțuri neutre, 14 emoții	DA	22 kHz	---	microfon	---	-recunoaștere vocală, modelare pronunție, prozodie	DA
Jankowski et al., 1990 - NTIMIT	630	Lb. engleză	Text TIMIT	DA	---	---	telefon	---	Recunoaștere vocală	DA
P. Kingsburg	---	Lb.	90988	---	---	---	telefon	---	Recunoaștere vocală	DA



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

et al., 1997		engleză	cuvinte de jurnal							
K. Karins et al., 1997	---	Lb. germană	318807 cuvinte	---	---	---	telefon	---	Recunoaștere vocală	DA
S. Garrett et al., 1996	---	Lb. spaniolă	45582 cuvinte	---	---	---	telefon	---	Recunoaștere vocală	DA
J. Kong și D. Graff, 2005	---	Lb. engleză, arabă, mandarină	știri	---	16 kHz / 2 canale	---	radio	DA	Detectare subiect comun în știri, detectare perechi de știri cu subiect comun etc.	DA
D. Graff, 2001a	---	Lb engleză	știri	---	16 kHz/1 canal	16 biți	radio	DA	Detectare subiect comun în știri, detectare perechi de știri cu subiect comun etc.	DA
D. Graff, 2001b	---	Lb. mandarină	știri	---	16 kHz/1 canal	16 biți	radio	DA	Detectare subiect comun în știri, detectare perechi de știri cu subiect comun	DA



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013MINISTERUL
EDUCAȚIEI
CERCETĂRII
TINERETULUI
ȘI SPORTULUI

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

									etc.	
J.J. Gofrey, 1994a,b,c,d	---		Comunicare control aerian	---	8 kHz	---	Transmisie aeriană	---	Recunoaștere vocală în medii zgomotoase	DA
J. Fiscus et al., 2007a,b	---	---	Discursuri în cadrul întâlnirilor	---	---	---	microfon	---	-analiză discurs, recunoaștere vocală, extragere metadata	DA
NIST Multimodal Information Group,2011a,b	---	---	știri, conversări telefonice, comunicări	---	16 kHz- știri radio; 8 kHz/2 canale- telefon	---	Radio, Telefon, microfon	---	-detectare subiect, recunoaștere vocală	DA
O. Salor et al., 2007	193	Lb. turcă	Traducere 2000 propoziții TIMIT	---	16 kHz	---	microfon	DA	Recunoaștere vocală	DA
M. Boldea et	100	Lb.	Traducere 40	DA	20 kHz	16 biți	microfon	NU	Dezvoltarea de	---



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

al., 1998		română	pasaje lb. engleză corpus EUROM-1, 550 propoziții						instrumente de recunoaștere vocală automată independentă de vorbitor	
A. Stan et al. 2011	1	Lb. română	3800 propoziții din ziare, povești	DA	96 kHz/ 44 kHz	24 biți	microfon	DA	Sinteza voce	---
J. Matousek și J. Romportl, 2007	1	Lb. cehă	5000 pasaje ziare	DA	-	-	microfon	NU	Sinteza voce	---
C.G. Clopper și D.B. Pisoni, 2006	60	6 dialecte lb. engleză	Cuvinte, propoziții, pasaje	DA	44 kHz	16 biți	microfon	DA	Determinare dialecte	---
A. Batliner et al., 2005	611	Lb. engleză, suedeză, germană, italiană	220 cuvinte izolate, 50 propoziții engleză, 400 propoziții	DA	44 kHz	16 biți	microfon	NU	Recunoașterea automată a vocii spontane și emoționale de copii	---



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

		copii	generice, voce spontană, citită							
N. Kawaguchi et al., 2005	812	---	1 milion morfeme	DA	---	---	microfon, video	DA	Analiza vocii în timpul conducerii de autovehicule	---
Kay Elemetrics Corp, 1994	53 N, 657 P	Lb. engleză	Vocala susținută /a/	---	25 kHz/50 kHz	16 biți	microfon	---	Deteție voce patologică	DA
G. Nuffelen et al.,	122 N, 197 P	Lb. germană	50 cuvinte context CVC, vocala susținută /a/ vorbire spontană și semispontană	DA	16 kHz	16 biți	microfon	NU	Instrumente de evaluare a vocii patologice, teste de inteligibilitate	---
J. F. Bonastre et al., 2007	58 N,	---	Vocala susținută /a/	---	---	---	microfon	DA	Evaluare voce patologică	---



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

	391 P									
J. B. Tomblin, 2010	2009	Lb. engleză, copii	---	---	---	---	microfon	DA	Evaluare voce patologică și detecție defecte de limbaj	---
G. Schlottauer și M.E. Torres, 2006	53 N, 36 P	---	Vocala /a/ susținută	---	22 kHz	16 biti	microfon	DA	Clasificarea automată a vocii disfonice	---
F. Pozo et al., 2009	40 N, 40 P	Lb. spaniolă	Patru propoziții	---	48 kHz	---	microfon	DA	Detectarea automată a apneei obstructive	---
L. Salhi et al., 2010	30 N, 30 P	---	cuvinte	---	---	---	microfon	DA	Detecție voce patologică	---
J. Lee și M. Hahn, 2009	83 P	---	Vocala susținută /a/	---	---	---	---	---	Detecție voce patologică	---
F. Bocklet et al., 2010	13 P	Lb. germană	Text standard din 108 cuvinte	---	16 kHz	16 biți	microfon	DA	Determinarea automată a eficienței tratamentului	---
K. T.	22 P	---	35 propoziții	---	---	---	microfon	NU	Determinarea	---



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

Bressmann et al., 2011									modificării vocii după aplicarea de expanderi	
P. Jindra et al., 2002	10 P	---	Silabe care conțin consoane bilabial, alveolare, labiodentale	---	---	---	microfon	DA	Detectarea aplicării corecte de proteză parțială și completă	---
F. Stelzle et al., 2010	40 N, 28 P	---	---	---	---	---	microfon	DA	Recunoaștere automată a vocii patologice	---
H.N. Teodorescu și S.M. Feraru, 2007; H.N. Teodorescu și A. Untu, 2010; A. Untu, H.N. Teodorescu,	16 N, 8P	Lb. română	Cuvinte care conțin consoane fricative	DA	22 kHz	16 biți	microfon	DA	Diagnosticarea automată a vocilor patologice de cauză stomatologică	DA



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



MINISTERUL
EDUCAȚIEI
CERCETĂRII
TINERETULUI
ȘI SPORTULUI
OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

2011										
A. Untu et al., 2011	24 N, 27 P	Lb. franceză	Cuvinte care conțin consoane fricative	DA	22 kHz	16 biți	microfon	DA	Diagnosticarea automată a vocilor patologice de cauză stomatologică	NU

--- = nespecificat



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



MINISTERUL
EDUCAȚIEI
CERCETĂRII
TINERETULUI
ȘI SPORTULUI

OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

Referințe

1. Batliner, A., Blomberg, M., D'Arcy, S., Elenius, D., Giuliani, D., Gerosa, M., Hacker, C., Russell, M., Steidl, S. and Wong, M., *The PF-STAR Children's Speech Corpus*, in Proc. Eurospeech, pp. 2761–2764, 2005.
2. Bocklet T., Hönig, F., Haderlein, T., Stelzle, F., Knipfer C. and Nöth E., *Automatic Detection and Evaluation of Edentulous Speakers with Insufficient Dentures*, Lecture Notes in Computer Science, , vol. 6231, pp. 243-250, 2010.
3. Boldea, M., Munteanu, C. and Doroga, A., *Design, Collection, and Annotation of a Romanian Speech Database*, In Proceedings LREC Workshop on Speech Database Development for Central and Eastern European Languages, Granada, Spain, 1998.
4. Bonastre, J.F., Fredouille, C., Ghio, A., Giovanni, A., Pouchoulin, Revis, G. J., Teston, B. and Yu, P., *Complementary Approaches for Voice Disorder Assessment*, in Proc. Interspeech, pp. 1194–1197, 2007.
5. Bressmann, K.T., Gong, S. and Tompson, B.D., *Impact of a Rapid Palatal Expander on Speech Articulation*, American Journal of Orthodontics and Dentofacial Orthopedics, vol. 140, nr. 2, pp. 67-75, 2011.
6. Clopper, C.G., Pisoni, D.B., *The Nationwide Speech Project: A new corpus of American English dialects*, Speech Communication., vol.48, pp. 633–644, 2006.
7. Dibazar, A.A. and Narayanan, S.S., *A System for Automatic Detection of Pathological Speech*, in 36th Asilomar Conf. Signal, Systems, and Computers, Asilomar, CA, USA, 2002.
8. Feraru S.M., Teodorescu H.N., Zbancioc M.D., *SRoL-Web-based Resources for Languages and Language Technology e-Learning*, International Journal of Computers Communications & Control, vol. 5, nr. 3, pp. 301-313, 2010.
9. Feraru, M., Teodorescu, H.N., *The Emotional Speech Section of the Romanian Spoken Language, Archive*, Conf. On Intelligent Systems and Technologies, Proc. 5th European, Iasi, Romania, 2008.
10. Fernández-Pozo R, Murillo J.L.B., Gómez L.H., Gonzalo E.L., Ramírez J.A., Toledano D.T., *Assessment of Severe Apnoea through Voice Analysis, Automatic Speech, and Speaker Recognition Techniques*, EURASIP Journal on Advances in Signal Processing - Special issue on recent advances in biometric systems: a signal processing perspective, 2009.
11. Fiscus, J. et al., *2004 Spring NIST Rich Transcription (RT-04S) Development Data*, Linguistic Data Consortium, Philadelphia, 2007a.
12. Fiscus, J. et al., *2004 Spring NIST Rich Transcription (RT-04S) Evaluation Data*, Linguistic Data Consortium, Philadelphia, 2007b.
13. Garofolo, J. S., Lamel, L. F., Fisher, W. M., Fiscus, J. G., Pallett, D. S. and Dahlgren N. L., *The DARPA TIMIT Acoustic-Phonetic Continuous Speech Corpus CDROM*, NIST, 1986 [www ldc.upenn.edu/ol/docs/TIMIT.html].
14. Garrett, S., Morton, T. and McLemore, C., *CALLHOME Spanish Lexicon*, Linguistic Data Consortium, Philadelphia, 1996.
15. Godfrey, J. J., *Air Traffic Control BOS*, Linguistic Data Consortium, Philadelphia, 1994a.
16. Godfrey, J. J., *Air Traffic Control Complete*, Linguistic Data Consortium, Philadelphia, 1994b.
17. Godfrey, J. J., *Air Traffic Control DCA*, Linguistic Data Consortium, Philadelphia, 1994c.



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRU



Fondul Social European
POSDRU 2007-2013



Instrumente Structurale
2007-2013



OIPOSDRU



UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

18. Godfrey, J. J., *Air Traffic Control DFW*, Linguistic Data Consortium, Philadelphia, 1994d.
19. Graff, D., *TDT3 English Audio*, Linguistic Data Consortium, Philadelphia, 2001a.
20. Graff, D., *TDT3 Mandarin Audio*, Linguistic Data Consortium, Philadelphia, 2001b.
21. Han, N., Graff, D. and Kim, M., *Korean Telephone Conversations Lexicon*, Linguistic Data Consortium, Philadelphia, 2003.
22. http://www.inl.nl/tst-centrale/images/stories/producten/documentatie/copas_manual.pdf
23. <http://www.clarin.eu/>
24. <http://www.elra.info/>
25. <http://www ldc.upenn.edu/>
26. Huang, S., Bian, X., Wu, G. and McLemore, C., *CALLHOME Mandarin Chinese Lexicon*, Linguistic Data Consortium, Philadelphia, 1996.
27. Jankowski et al., *NTIMIT: A Phonetically Balanced, Continuous Speech, Telephone Bandwidth Speech Database*, Proc. ICASSP-90, April 1990.
28. Karins, K. et al., *CALLHOME German Lexicon*, Linguistic Data Consortium, Philadelphia, 1997.
29. Kawaguchi, N., Matsubara, S., Takeda, K. and Itakura, F., *CIAIR In-Car Speech Corpus -Influence of Driving Status*, in IEICE Transactions on Information and Systems, vol. E88-D, nr. 3, pp. 578–582, 2005.
30. Kay Elemetrics Corp, *Disordered Voice Database*, Version 1.03, 1994.
31. Kilany, et al., *Egyptian Colloquial Arabic Lexicon*, Linguistic Data Consortium, Philadelphia, 2002.
32. Kingsbury, P. et al., *CALLHOME American English Lexicon (PRONLEX)*, Linguistic Data Consortium, Philadelphia, 1997.
33. Kobayashi, M., Crist, S., Kaneko, M. and McLemore, C., *CALLHOME Japanese Lexicon*, Linguistic Data Consortium, Philadelphia, 1996.
34. Kong, J. and Graff, D., *TDT4 Multilingual Broadcast News Speech Corpus*, Linguistic Data Consortium, Philadelphia, 2005.
35. Lee, J., Hahn, M., *Automatic Assessment of Pathological Voice Quality Using Higher-Order Statistics in the LPC Residual Domain*, EURASIP Journal on Advances in Signal Processing, vol. 2009, 2009.
36. Liberman, M. et al., *Emotional Prosody Speech and Transcripts*, Linguistic Data Consortium, Philadelphia, 2002.
37. Matousek, J., Romportl, J., *Recording and Annotation of Speech Corpus for Czech Unit Selection Speech Synthesis*, Lecture Notes in Artificial Intelligence, vol. 4629, Springer, Berlin-Heidelberg, pp. 326–333, 2007.
38. Moran, R. J., Reilly, R. B., de Chazal, P., and Lacy, P. D., *Telephony-based Voice Pathology Assessment using Automated Speech Analysis*, IEEE Trans. Biomed. Eng., vol. 53, nr. 3, pp. 468–477, 2006.
39. NIST Multimodal Information Group, *2006 NIST Spoken Term Detection Evaluation Set*, Linguistic Data Consortium, Philadelphia, 2011.
40. NIST Multimodal Information Group, 2011, *2006 NIST Spoken Term Detection Development Set*, Linguistic Data Consortium, Philadelphia
41. P. Jindra, M. Eber, J. Pešák, *The Spectral Analysis of Syllables in Patients using Dentures*, Biomed. Papers, vol. 146, nr. 2, pp. 91–94, 2002.
42. Saenz-Lechon, N., Godino-Llorente, J.I., Osma-Ruiz, V., Gomez-Vilda, P., *Methodological Issues in the Development of Automatic Systems for Voice Pathology Detection*, Biomed. Signal Process. Control 1, pp.120–128, 2006.



UNIUNEA EUROPEANĂ

GUVERNUL ROMÂNIEI
MINISTERUL MUNCII, FAMILIEI
ȘI PROTECȚIEI SOCIALE
AMPOSDRUFondul Social European
POSDRU 2007-2013Instrumente Structurale
2007-2013

OIPOSDRU

UNIVERSITATEA TEHNICĂ
"GHEORGHE ASACHI"
DIN IAȘI

43. Salhi, L., Mourad, T. and Cherif, A., *Voice Disorders Identification Using Multilayer Neural Network*, The International Arab Journal of Information Technology, vol. 7, nr. 2, 2010.
44. Salor, O., Pellom, B. L., Iloglu, T., Demirekler, C. M., *Turkish Speech Corpora and Recognition Tools Developed by Porting Sonic: Towards Multilingual Speech Recognition*, Comput. Speech. Lang., vol. 21, nr. 4, pp. 580–593, 2007.
45. Schlotthauer, G., Torres, M.E., *Automatic Diagnosis of Pathological Voices*, Proceeding of the 6th WSEAS International Conference on Signal, Speech and Image Processing, Lisbon, Portugal, pp. 150-155, 2006.
46. Shobaki, K., Hosom, J. P. and Cole, R., *CSLU: Kids` Speech Version 1.1*, Linguistic Data Consortium, Philadelphia, 2007.
47. Stan, A., Yamagishi, J., King, S., and Aylett, M., *The Romanian Speech Synthesis (RSS) Corpus: Building a High Quality HMM-based Speech Synthesis System using a High Sampling Rate*, Speech Communication, vol. 53, pp. 442–450, 2011.
48. Stelzle, F., Ugrinovic, B., Knipfer, C., Bocklet, T., Noth, E., Schuster, M., Eitner, S., Seiss, M. și Nkenke, E., *Automatic, Computer-Based Speech Assessment on Edentulous Patients with and without Complete Dentures – Preliminary Results*, Journal of Oral Rehabilitation, vol. 37, pp. 209–216, 2010.
49. Teodorescu, H.N., Feraru, S.M., Trandabăț, D., Zbancioc, M., Luca, R., Verbuță, A., Hnatiuc, M., Ganea, R., Voroneanu, O., Pistol, L., Șcheianu, D., *Situl Web Sunetele Limbii Române*, http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.html, 2005-2007
50. Teodorescu, H.N., Feraru, M., *De ce nu Place Vocea Sintetizată - Câteva Elemente de Comparație cu Vocea Umană*, Trandabăț, Cristea, Tufiș (eds.), ConsILR 2008, Resurse lingvistice și instrumente pentru prelucrarea limbii române, Editura Universității "Al. I. Cuza", 19-21 noiembrie, pp. 21-30, Iași, România, 2008.
51. Teodorescu, H.N., Feraru, M., *Micro-corpus de Sunete Gnatosonice si Gnatofonice*, Pistol, Cristea, Tufis (Eds.), Resurse lingvistice si instrumente pentru prelucrarea limbii romane, Ed. Universitatii "Al.I.Cuza" Iasi, pp.21-30, 2007.
52. Teodorescu, H.N., **Untu, A.**, *Corpus pentru Gnatofonie: Protocol, Metodologie, Adnotare*, ConsILR 2010, București, 6-7 mai, Editori: Adrian Iftene, Horia-Nicolai Teodorescu, Dan Cristea, Dan Tufiș, Editura Univesității "Al. I. Cuza" Iași, pp. 13-22, 2010.
53. Tomblin, J. B., *The EpiSLI Database: A Publicly Available Database on Speech and Language*, Language, Speech, and Hearing Services in Schools, vol. 41, nr. 1, pp. 108-117, 2010.
54. **Untu (Hulea), A.**, Duvanaud, C., Teodorescu, H.N., *A Study of the Relationship between the Statistical Acoustical Features of the /v/ and /f/ French Fricative Consonants and the Dentistry Pathologies*, (acceptată, IEEE International Conference on E-Health and Bioengineering, 24-26 nov., 2011, Iasi, Romania).
55. **Untu, A.**, Teodorescu, H.N., *Pattern Analysis of /v/ Pronunciations in _v/V Contexts*, The Eighth IASTED Int. Conf. on Signal Processing, Pattern Recognition, and Applications SPPRA 2011, Innsbruck, 2011.
56. Wright, J., *Articulation Index*, Linguistic Data Consortium, Philadelphia, 2005.
57. Zbancioc, M., *Tools for the Archive of the Romanian Language Sounds Project*, 4th European Conf. On Intelligent Systems and Technologies, Iasi, Romania, 2006.