

Ângela Daniela Carneiro Teixeira

# DETEÇÃO E CORREÇÃO DE DISFLUÊNCIAS EM CRIANÇAS

Setembro de 2012



UNIVERSIDADE DE COIMBRA





FACULDADE DE CIÊNCIAS E TECNOLOGIA  
UNIVERSIDADE DE COIMBRA

MESTRADO INTEGRADO EM ENGENHARIA ELETROTÉCNICA E DE  
COMPUTADORES

# DETEÇÃO E CORREÇÃO DE DISFLUÊNCIAS EM CRIANÇAS

---

**Ângela Daniela Carneiro Teixeira**

## **Júri**

Presidente: Professora Doutora Teresa Martinez dos Santos Gomes

Orientador: Professor Doutor Fernando Santos Perdigão

Vogal: Professor Doutor Marco Alexandre Cravo Gomes

**Setembro de 2012**



“Begin at the beginning,’ the King said, very gravely, ‘and go on till you come to the end: then stop.”

— Lewis Carrol, ‘Alice’s Adventures in Wonderland’



## Agradecimentos

“I almost wish I hadn't gone down the rabbit-hole--and yet--and yet--...”\*

Em primeiro lugar, gostaria de deixar o meu mais profundo agradecimento ao Professor Doutor Fernando Santos Perdigão por todo o apoio ao longo destes meses. Todos os conselhos e conhecimento que me transmitiu foram decisivos para a execução deste trabalho. Obrigada por nunca me ter deixado desamparada!

Gostaria também de agradecer aos colegas de laboratório, Dirce, Jorge, Arlindo e Sara, pelo companheirismo e apoio. Agradeço, em particular, à Carla Lopes pela dedicação com que me aconselhou e co-orientou.

Não posso também deixar de agradecer aos meus colegas e amigos de curso. Obrigada por todos os momentos que comigo partilharam, dentro e fora destas paredes. Num futuro próximo, comeremos mimos doces todos juntos mais uma vez!

Também aos meus amigos gostaria de deixar umas palavras de apreço. Vocês são o meu Norte e, como uma vez alguém me escreveu: ‘Posso viver sem vocês, mas quando estão longe torna-se difícil respirar...’. Obrigada por me terem deixado entrar como eu vos deixei entrar também!

Resta-me agradecer à minha família: ao meu pai por ser uma pedra de força e estabilidade; à minha mãe por me ralhar e exigir de mim tudo aquilo que sabe de que eu sou capaz e por ser o meu coração; ao meu irmão por ser tão simplesmente ele, brincalhão, chato e sempre presente para mim; também ao mais recente membro da família, o meu Sleepy, a bola de pelo mais parva e carente, nos momentos mais inconvenientes.

“I can't go back to yesterday because I was a different person then.”\*

Obrigada a todos por me terem feito ‘diferente’.

---

\* Lewis Carrol , ‘Alice’s Adventures in Wonderland’





## Resumo

A fala é o meio de comunicação por excelência, sendo o seu processo de aprendizagem um dos mais complexos que se conhece. No entanto, as crianças são extremamente eficientes neste processo. Ao longo do seu desenvolvimento, a criança desenvolve uma série de mecanismos mentais (ao nível da percepção e ao nível do controlo motor) que lhe permitem a aprendizagem e aperfeiçoamento da fala. No início do processo de aprendizagem, existem sempre alguns sons cuja replicação é mais difícil, sendo expectável que a criança apenas adquira o domínio de alguns fones mais complexos numa fase mais avançada do seu desenvolvimento.

Com a evolução das tecnologias desenvolvidas na área do reconhecimento automático de fala houve um investimento da sua aplicação a áreas como as da aprendizagem da língua ou da terapia da fala.

Um trabalho realizado anteriormente, que propunha o desenvolvimento de um sistema de auxílio à terapia da fala para detecção de disfluências em crianças em idade pré-escolar, serviu de ponto de partida para o presente estudo. Foi objectivo deste estudo implementar o módulo desse sistema correspondente à identificação e correcção de disfluências para sua posterior inserção no sistema completo.

Em termos de processamento do sinal de fala, o sistema que se baseia nos Modelos de Markov Não Observáveis (HMM – *Hidden Markov Models*) é um dos mais utilizados, tendo sido a escolha para o desenvolvimento de modelos de reconhecimento de fala para crianças a aplicar neste estudo.

No âmbito deste estudo procedeu-se ao tratamento da base de dados recolhida no estudo referido, para criar os modelos de reconhecimento de fala necessários para solucionar o problema principal proposto, o da identificação de dificuldades de pronúncia. Os modelos foram criados com recurso às ferramentas facultadas pelo HTK (*Hidden Markov Model Toolkit*), produzindo resultados de reconhecimento satisfatórios. Procedeu-se então ao desenvolvimento de um método que permitisse a detecção de disfluências em palavras proferidas por crianças, usando para isso os valores de verosimilhança calculados no processo de reconhecimento da locução a analisar.

**Palavras-chave: Fala, reconhecimento, crianças, modelos, verosimilhança, HMM.**



## Abstract

Speech is the prime mean of communication, with its learning process being one of the most complex and yet, one where humans excel at. Throughout its growth, the child develops a series of mental mechanisms (either at a perception level or at a motor control level) which allow the learning and improvement of speech. At the beginning of the learning process, there are always some sounds whose replication is more difficult, it is expected that the child will only acquire the domain of some of the more complex phones at a later stage of its development.

As the technologies in the area of automatic speech recognition evolved, there was an increase in the tendency to apply it to areas such as language learning or speech therapy.

A previous work, which proposed a Computer Aided Speech Therapy system directed toward the detection of disfluencies for children in pre-school age, served as the starting point for this study. The development of this system's module that corresponded to the identification and correction of disfluencies, for subsequent insertion into the complete system, was the main goal of this thesis.

In terms of speech signal processing, the system based on Hidden Markov Models (HMM - Hidden Markov Models) is one of the most widely used and was chosen in this study for the development of models of speech recognition for children.

For the purpose of this study, the processing of the database collected in the above mentioned work was carried out. This necessary for creating the speech recognition models needed to implement this study's main goal of identifying pronunciation problems. The models were created using the tools provided by HTK (Hidden Markov Model Toolkit), having produced satisfactory recognition results. A method that allowed the detection of disfluencies on words spoken by children was also developed, using for this the likelihood values calculated in the recognition of the utterance to be tested.

**Keywords: Speech, recognition, children, likelihood, HMM.**



# Índice

1. Introdução .....	1
1.1. Motivação .....	1
1.2. Objectivos .....	2
1.3. Organização dos Conteúdos.....	2
2. O Ensino Assistido da Língua.....	3
2.1. Tecnologia Actual .....	3
2.2. Considerações Sobre a Fala de Crianças .....	6
2.3. Proposta do Sistema a Desenvolver .....	9
3. Técnicas de Análise e Reconhecimento da Fala .....	11
3.1. Coeficientes Cepstrais em Escala Mel .....	11
3.2. Hidden Markov Models .....	12
4. Reconhecimento de Fala de Crianças.....	17
4.1. Recolha de uma Base de Dados de Fala de Crianças .....	17
4.2. Sistema de Reconhecimento .....	23
4.2.1 Treino dos Modelos.....	23
4.2.2 Teste dos Modelos.....	26
5. Verificação.....	29
6. Conclusão e Trabalho Futuro.....	39
7. Anexos .....	41
7.1. Anexo A – Ficheiro de Configuração do HCopy.....	41
7.2. Anexo B - Lista de Fones e Dicionário .....	41
7.3. Anexo C - Scripts em DOS/BATCH de Treino dos Modelos .....	48
7.4. Anexo D – Scripts e Funções do Matlab .....	50
8. Referências Bibliográficas .....	61

## Lista de Acrónimos e Abreviações

ASR	Automatic Speech Recognition
CALL	Computer-Assisted Language Learning
DET	Detection Error Tradeoff
FSM	Finite State Machine
GMM	Gaussian Mixture Model
HMM	Hidden Markov Model
HTK	Hidden Markov Model Toolkit
LPC	Linear Predictive Coding
MAPLR	Maximum A Posteriori Linear Regression
MFCC	Mel-Frequency Cepstrum Coefficients
MLF	Master Label File
MMF	Master Macro File
ROC	Receiver Operating Characteristic
SAMPA	Speech Assessment Methods Phonetic Alphabet
SLF	Standard Lattice Format
VTLN	Vocal Tract Length Normalisation
WER	Word Error Rate

## Lista de Figuras

FIGURA 2. 1 – GRÁFICO LPC DE UMA TRAMA DA VOGAL ‘I’ .....	7
FIGURA 2. 2 – POSICIONAMENTO DAS VOGAIS EM FUNÇÃO DE F1-F2 PARA CRIANÇAS (A) E ADULTOS (B) [17].....	7
FIGURA 2. 3 – VARIÇÃO DA FREQUÊNCIA DE F1 E F2 COM A IDADE [15]’ .....	8
FIGURA 2. 4– DIAGRAMA DE BLOCOS DO SISTEMA PROPOSTO [26].....	9
FIGURA 3. 1 – EXEMPLO DE UM MODELO HMM SIMPLES COM TOPOLOGIA ESQUERDA-DIREITA [35].....	13
FIGURA 3. 2 – EXEMPLO DE UM MODELO HMM COM TOPOLOGIA ESQUERDA-DIREITA E COM FUNÇÕES DENSIDADE DE PROBABILIDADE COMPOSTAS POR GMMS [26].....	14
FIGURA 3. 3 – EXEMPLO DE UM DIAGRAMA <i>TRELLIS</i> (A) E DO HMM QUE O ORIGINOU (B) [38].....	15
FIGURA 4. 1 - OCORRÊNCIA DE CADA FONE NA LISTA DE 200 PALAVRAS.....	17
FIGURA 4. 2 – VIEWLABS EM EXECUÇÃO COM UM FICHEIRO DE ÁUDIO E AS RESPECTIVAS ETIQUETAS E MARCAÇÕES TEMPORAIS .....	18
FIGURA 4. 3 – EXEMPLO DE UM FICHEIRO MLF (EXCERTO). .....	19
FIGURA 4. 4 – EXCERTO DO DICIONÁRIO COM TRADUÇÃO SAMPA E COM EXEMPLO DE MULTIPRONUNCIÇÃO. ...	20
FIGURA 4. 5 – REPRESENTAÇÃO DOS 14 FILTROS NA FREQUÊNCIA.....	21
FIGURA 4. 6 - REPRESENTAÇÃO DOS 14 FILTROS DA FIGURA ANTERIOR EM ESCALA MEL.....	21
FIGURA 4. 7 – EXEMPLO DO USO DA FERRAMENTA HLED PARA A PALAVRA ‘COPO’ QUE TEM A TRANSCRIÇÃO SAMPA ‘KOPU’ .....	24
FIGURA 4. 8 – DIAGRAMA DO PROCESSO EXECUTADO NO HTK PARA TREINO DOS MODELOS DE FONES. AS ELIPSES INDICAM FERRAMENTAS OU PROCESSOS, SENDO QUE A MAIORIA SÃO FERRAMENTAS DO HTK ESTANDO TAMBÉM PRESENTE UM PROCESSO MANUAL. OS FICHEIROS MLF ENCONTRAM-SE REPRESENTADOS A COR- DE-ROSA, ENQUANTO OS FICHEIROS UTILIZADOS PELO HTK SE ENCONTRAM REPRESENTADOS A LILÁS. O DICIONÁRIO COM A LISTA DE PALAVRAS CORRECTAS E SUA CORRESPONDENTE TRANSCRIÇÃO EM SAMPA ESTÁ REPRESENTADO A VERDE. ESTE DIAGRAMA REPRESENTA DE UMA FORMA CONCISA O FUNCIONAMENTO DO SCRIPT ‘TREINO.BAT’ (SECÇÃO DE ANEXOS 7.3) .....	25
FIGURA 4. 9 – DIAGRAMA DA SEQUÊNCIA AUTORIZADA NO RECONHECIMENTO. ....	26
FIGURA 4. 10 – EXEMPLO DO <i>OUTPUT</i> DA ANÁLISE ESTATÍSTICA PRODUZIDA PELA FERRAMENTA. ....	27
FIGURA 4. 11 – EXEMPLO DE ERROS DE RECONHECIMENTO (SUBSTITUIÇÕES) PRESENTES NO FICHEIRO DE <i>OUTPUT</i> DO <i>HRESULTS</i> . ....	28
FIGURA 5. 1 – DIAGRAMA DA SEQUÊNCIA LIVRE DE FONES AUTORIZADA PELA GRAMÁTICA LIVRE. ....	29
FIGURA 5. 2 – EXCERTO DO FICHEIRO GERADO QUE COMPARA AS DIFERENTES VEROSIMILHANÇAS E AS <i>STRINGS</i> DE FONES PARA FICHEIROS REF E SAID IGUAIS.....	30
FIGURA 5. 3 – RESULTADOS PRODUZIDOS PELO <i>HVITE</i> PARA O FICHEIRO ‘S0054L152’ COM RECURSO À GRAMÁTICA ‘FIXA’ (A) E À GRAMÁTICA ‘LIVRE’ (B), PARA FICHEIROS REF E SAID IGUAIS.....	32

FIGURA 5. 4 – RESULTADOS PRODUZIDOS PELO HVITE PARA O FICHEIRO ‘S0097L000’ COM RECURSO À GRAMÁTICA ‘FIXA’ (A) E À GRAMÁTICA ‘LIVRE’ (B), PARA FICHEIROS REF E SAID IGUAIS. ....	32
FIGURA 5. 5 – EXCERTO DO FICHEIRO GERADO QUE COMPARA AS DIFERENTES VEROSIMILHANÇAS E AS <i>STRINGS</i> DE FONES PARA FICHEIROS REF E SAID DIFERENTES. ....	33
FIGURA 5. 6 – RESULTADOS PRODUZIDOS PELO HVITE PARA O FICHEIRO ‘S0097L000’ COM RECURSO À GRAMÁTICA ‘FIXA’ (A) E À GRAMÁTICA ‘LIVRE’ (B), PARA ETIQUETAS REF E SAID DIFERENTES. ....	33
FIGURA 5. 7 – EXEMPLO DA MATRIZ DE CONFUSÃO 2X2 QUE REPRESENTA O CORRENTE PROBLEMA DE CLASSIFICAÇÃO. [44]. ....	34
FIGURA 5. 8 – GRÁFICO DAS CURVAS NO ESPAÇO ROC QUE COMPARA A DISTÂNCIA L1-L2 PARA AS TRANSCRIÇÕES DITAS CORRECTAMENTE E AS ERRADAS. ....	35
FIGURA 5. 9 – CURVA ROC. ....	36
FIGURA 5. 10 – CURVA DET. ....	37
FIGURA 5. 11 – CURVA ROC PARA O <i>SCORE</i> DE ALINHAMENTO. ....	38
FIGURA 5. 12 – CURVA DET PARA O <i>SCORE</i> DE ALINHAMENTO. ....	38



# 1. Introdução

## 1.1. Motivação

A fala e o seu processo de aprendizagem são um dos sistemas mais complexos que se conhece. No entanto, o ser humano fá-lo de modo inato, com extrema rapidez e eficiência. A aprendizagem de uma língua é um processo iterativo baseado em imitação, como é norma no decurso de desenvolvimento da criança. Quando esta começa a falar, surgem dificuldades em pronunciar alguns sons e palavras. Com o tempo, a criança desenvolve capacidade de controlo dos meios intervenientes na produção da fala, desenvolvendo aptidões de pronúncia e articulação. Porém, nem sempre estas dificuldades são ultrapassadas em tempos considerados normais. É nestes casos que surge a necessidade de intervenção terapêutica, de modo a impedir ou minimizar problemas na comunicação que daí advenham. Este problema torna-se particularmente pertinente quando se observam as elevadas taxas de insucesso escolar que podem ser reflexo de problemas como o acima descrito.

Desde o aparecimento dos computadores e sistemas informatizados surgiu uma tendência de aplicar as suas mais-valias aos aspectos mais vulgares da vida. A fala, sendo o meio de comunicação por excelência, não foi excepção, tendo sido motivo do desenvolvimento de inúmeros estudos e aplicações. Os sistemas de reconhecimento automático de fala (ASR - *Automatic Speech Recognition*) são talvez o expoente máximo do investimento feito neste campo. A aplicação de todo este conhecimento e tecnologia ao processo de ensino da língua surgiu de forma natural. O uso de computadores aplicado ao ensino/aprendizagem de línguas data da década de 60 sob o acrónimo CALL (*Computer Assisted Language Learning*) e faz uso de inúmeras aplicações ICT (*Information and Communication Technologies*).

Em 2002 surgiu um projecto de terapia da fala financiado pela União Europeia [1] que motivou o desenvolvimento de diversas ferramentas para auxiliar os terapeutas, para avaliação da produção de fala, entre outras.

Os sistemas de auxílio à terapia da fala podem ser segmentados em grupos, dependendo do público-alvo. Releva-se de entre eles o tratamento de problemas de dislalia<sup>1</sup>. Muitos destes sistemas usam tecnologias ASR na sua implementação.

## 1.2. Objectivos

O principal objectivo do presente trabalho foi o da criação de um sistema que possibilitasse diagnóstico de disfluências de pronunciação, dificuldade articulatória, ou mesmo dislexia em crianças em idade pré-escolar, para posterior correcção ou terapia da fala, com recurso ao reconhecimento automático de fala.

Foi ainda objectivo deste estudo criar modelos de reconhecimento de fones adaptados à faixa etária em questão, necessários para o desenvolvimento do supracitado sistema.

## 1.3. Organização dos Conteúdos

No capítulo 2 deste estudo apresentam-se algumas considerações sobre o contexto actual do uso de computadores adaptado à terapia de fala e dos sistemas de reconhecimento automático de fala. São ainda apresentadas algumas particularidades da fala de crianças e expõe-se uma visão geral sobre o sistema a desenvolver.

O capítulo 3 contém uma breve exposição sobre uma das técnicas mais comumente utilizadas em reconhecimento automático de fala, os Modelos de Markov Não-Observáveis (HMM – *Hidden Markov Models*), e sobre uma técnica de parametrização espectral.

No capítulo 4 apresenta-se o desenvolvimento do processo de criação dos modelos de fones, incluindo o tratamento da base de dados e a parametrização dos ficheiros de áudio.

O capítulo 5 inclui a descrição do sistema de verificação e identificação de disfluências desenvolvido, bem como a proposta de uma solução para o problema de classificação binária inerente ao problema de identificação de disfluências.

No último capítulo são apresentadas as conclusões, a análise dos resultados obtidos e propostas sobre a direcção trabalhos futuros de continuidade.

---

<sup>1</sup> Dificuldades na fala de origem não neurológica associadas aos órgãos externos de produção de fala [26].

## 2. O Ensino Assistido da Língua

### 2.1. Tecnologia Actual

O uso de computadores aplicado ao ensino/aprendizagem de línguas é conhecido sob o acrónimo CALL. A fase em que a tecnologia das aplicações CALL se encontra varia de acordo com a perspectiva de quem faz esta análise. Por exemplo, enquanto Warschauer defende que, de 2000 em diante, se entrou na fase ‘Integrativa’<sup>2</sup> [2], já Bax [3] considera que a análise anterior possui falhas e algumas contradições e propõe uma nova visão onde considera que a presente fase é a ‘Open CALL’<sup>3</sup>. Este defende ainda que só deveria ser considerada a fase ‘Integrativa’ quando as aplicações CALL forem de uso tão simples e comum como os mais tradicionais e correntes métodos de ensino. A tabela seguinte resume a opinião de Warschauer sobre as fases das aplicações CALL. A Tabela 2 apresenta a visão de Bax sobre o mesmo assunto.

Tabela 1 – As fases das aplicações CALL segundo Warschauer [2].

<i>Stage</i>	1970s-1980s: Structural CALL	1980s-1990s: Communicative CALL	21st Century: Integrative CALL
<i>Technology</i>	Mainframe	PCs	Multimedia and Internet
<i>English-Teaching Paradigm</i>	Grammar-Translation & Audio-Lingual	Communicate Language Teaching	Content-Based, ESP/EAP
<i>View of Language</i>	Structural (a formal structural system)	Cognitive (a mentally-constructed system)	Socio-cognitive (developed in social interaction)
<i>Principal Use of Computers</i>	Drill and Practice	Communicative Exercises	Authentic Discourse
<i>Principal Objective</i>	Accuracy	Fluency	Agency

(Based on Kern & Warschauer, 2000; Warschauer, 1996; Warschauer, 2000a)

<sup>2</sup> Fase ‘Integrative’ – baseada na alargada dispersão de meios devido aos grandes desenvolvimentos dos computadores e da Internet.

<sup>3</sup> ‘Open CALL’ – aberto em todas as dimensões desde o *feedback*, tipos de software, papel do professor, particularmente quando comparado com a fase ‘Restricted CALL’.

Tabela 2 As fases CALL de acordo com Bax [3].

Table 2  
Restricted, Open and Integrated CALL: an outline

Content	Type of task	Type of student activity	Type of feedback	Teacher roles	Teacher attitudes	Position in curriculum	Position in lesson	Physical position of computer
<i>Restricted CALL</i>								
Language system	Closed drills Quizzes	Text reconstruction Answering closed questions Minimal interaction with other students	Correct/incorrect	Monitor	Exaggerated fear and/or awe	Not integrated into syllabus—optional extra  Technology precedes syllabus and learner needs	Whole CALL lesson	Separate computer lab
<i>Open CALL</i>								
System and skills	Simulations Games CMC	Interacting with the computer Occasional interaction with other students	Focus of linguistic skills development Open, flexible	Monitor/ facilitator	Exaggerated fear and/or awe	Toy Not integrated into syllabus—optional extra Technology precedes syllabus and learner needs	Whole CALL lesson	Separate lab—perhaps devoted to languages
<i>Integrated CALL</i>								
Integrated language skills work Mixed skills and system	CMC WP e-mail  Any, as appropriate to the immediate needs	Frequent interaction with other students Some interaction with computer through the lesson	Interpreting, evaluating, commenting, stimulating thought	Facilitator Manager	Normal part of teaching—normalised	Tool for learning Normalised integrated into syllabus, adapted to learners' needs <i>Analysis of needs and context precedes decisions about technology</i>	Smaller part of every lesson	In every classroom, on every desk, in every bag

O desenvolvimento de aplicações CALL é também assegurado por associações profissionais, a maior parte das quais sem fins lucrativos, que se dedicam a promulgar e discutir assuntos relacionados com a investigação e desenvolvimento das tecnologias de ensino de línguas. De entre estas podem referir-se algumas como a APACALL (The Asia-Pacific Association for CALL), a CALICO (The Computer Assisted Language Instruction Consortium), a EUROCALL (European Association for Computer Assisted Language Learning), ou a JATCALL (The Japan Association for Language Teaching - CALL).

De entre as aplicações CALL, têm interesse para o presente estudo aquelas que usam o ensino da língua falada. Neste contexto as aplicações CALL consideram-se divididas em dois grupos: aquelas que se assentam maioritariamente na tecnologia do reconhecimento de fala e as que usam apenas a síntese de fala.

No presente estudo, as tecnologias ASR (Automatic Speech Recognition) são de extrema importância uma vez que se pretende perceber se as crianças têm ou não problemas de dicção. A confiabilidade e a eficiência com que a aplicação a desenvolver é capaz de reconhecer correctamente a fala, dita o sucesso da aplicação desenvolvida.

Actualmente, existem vários pacotes de *software* comerciais destinados ao ensino de línguas,

No que diz respeito aos produtos CALL presentes no mercado, o Rosetta Stone [4] é provavelmente um dos mais conhecidos e bem-sucedidos, contendo aplicações que fazem uso de ASR (proprietário). Abrange uma ampla gama de línguas com uma vasta selecção de exercícios e tem tido uma boa aceitação por parte dos utilizadores.

A empresa LANGMaster distribui cursos de línguas *online* gratuitamente em diversos idiomas, sendo que alguns dos exercícios utilizam a tecnologia de reconhecimento de fala ViaVoice [5] da IBM<sup>5</sup>. Oferece centenas de horas de aprendizagem gratuitas em 75 idiomas, incluindo Português, embora apenas possua aplicações que fazem recurso a sistemas ASR para inglês. Babbel [6] é outro *software* de ensino de línguas que pode ser usado *online* e faz uso de ASR [7], necessitando o utilizador apenas do *browser* de Internet e da última versão do Adobe Flash Player. É multiplataforma e suporta vários idiomas, incluindo Português do Brasil (a norma PT-BR).

---

<sup>5</sup> - O sistema de reconhecimento desenvolvido pela IBM foi absorvido pela Nuance, já não existindo de forma independente.

Embora os programas CALL proliferem no mercado, os algoritmos de reconhecimento de fala atingem o seu auge em termos de eficiência e confiabilidade em aplicações para dicção de texto e aplicações de controlo do computador por fala.

O produto ‘Dragon NaturallySpeaking’, agora da Nuance [8] encontra-se, de longe, no topo [9] da lista no que diz respeito a eficiência e desempenho em reconhecimento da fala. A Nuance adquiriu o sistema ViaVoice da IBM, tornando-se líder de mercado. Também a companhia LumenVox [10] fornece soluções comerciais para aplicações de reconhecimento de fala.

A Microsoft também fez apostas nesta área, possuindo um sistema de reconhecimento de fala incorporado nos seus sistemas operativos desde o Windows Vista, tendo visto o seu desempenho e recepção melhorar com as optimizações incluídas no Windows 7 [11] [12].

Os produtos comerciais mencionados e a actual tecnologia ASR em geral, têm respostas bastante boas para os objectivos a que se propõem, embora os seus resultados com utilizadores não nativos (diferenças na pronúncia) ou com crianças, sofram uma considerável quebra.

Desde que surgiram estudos apontando o facto do fraco desempenho dos sistemas ASR com crianças, houve um incremento na atenção e interesse nessa área, tendo aparecido posteriormente muitos estudos sobre este assunto.

## **2.2. Considerações Sobre a Fala de Crianças**

Uma vez que as tecnologias de reconhecimento de fala apresentam limitações quando aplicadas à fala de crianças, tem interesse referir alguns aspectos fundamentais sobre as características da fala das crianças. De salientar dois desses aspectos: o tom e os formantes. O tom da fala é a frequência de vibração das cordas vocais que é especialmente alta nas crianças. Os formantes são os picos espectrais associados às ressonâncias do tracto vocal [13].

A análise LPC (Linear Predictive Coding) é usada frequentemente para representar um sinal na sua envolvente espectral, e é uma das ferramentas mais eficientes na análise de voz. A sua aplicação a sinais de fala permite uma fácil observação dos formantes, como é de resto perceptível na Figura 2. 1, representativa de uma vogal ‘i’. O conhecimento das posições espectrais dos formantes F1 e F2 (os formantes com as frequências mais baixas), possibilita a distinção das vogais, determinando as suas características como a altura (posição vertical da língua) – vogal aberta ou fechada - e a posteridade (posição horizontal da língua) – anterior ou posterior [14].

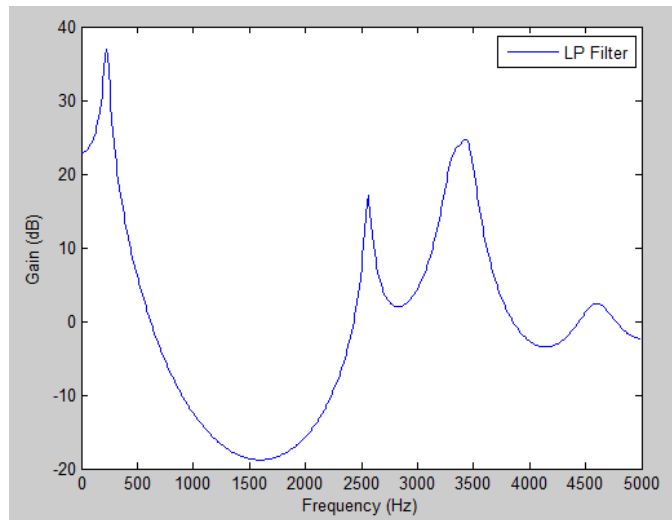


Figura 2. 1 – Gráfico LPC de uma trama da vogal ‘i’

No caso de crianças, os formantes possuem valores bastante mais elevados que no caso de adultos, aparecendo os picos na frequência que representam os formantes com frequências maiores. Este facto é devido ao reduzido tamanho do seu tracto vocal. Além disso a frequência de tom é também mais elevada que no caso de adultos, o que causa um aumento no espaçamento entre harmónicos [15] [16] em sons vozeados. Enquanto num adulto os três primeiros formantes se encontram geralmente entre os 300 Hz e os 3.2 kHz, para uma criança apenas dois podem ser encontrados nesta gama [17]. Estas particularidades, entre outras, provocam uma deterioração nos resultados dos sistemas de reconhecimento automático de fala.

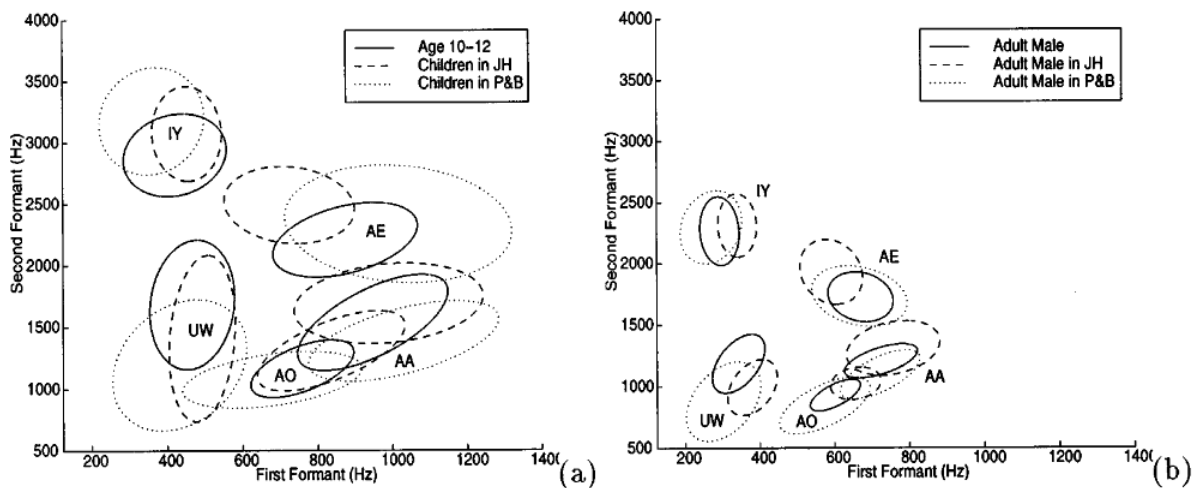


Figura 2. 2 – Posicionamento das vogais em função de F1-F2 para crianças (a) e adultos (b) [18]

Como é visível na figura anterior, há um considerável aumento da variação da frequência dos formantes com o decréscimo da idade. Também a Figura 2. 3 apresenta a variação espectral dos formantes com a variação da idade. É possível observar um decréscimo nos valores de F1 e de F2.

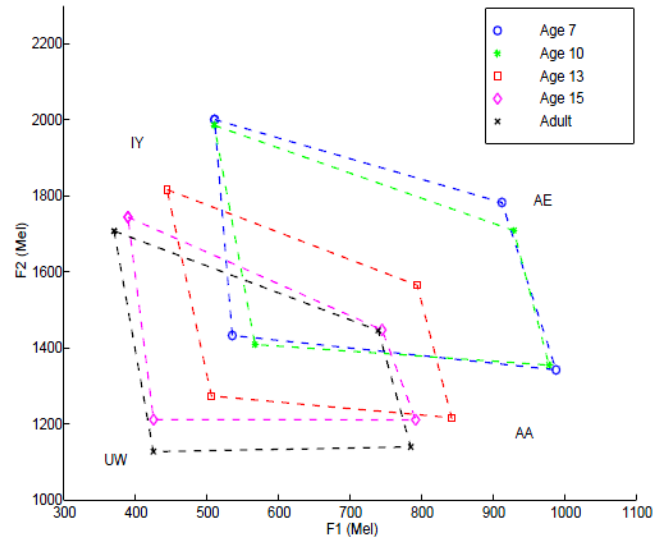


Figura 2. 3 – Variação da frequência de F1 e F2 com a idade [15].

Adicionalmente, crianças com idades inferiores a 10 anos apresentam uma maior variação da duração de vogais comparativamente aos valores normais em adultos além de uma maior variação espectral [19].

Além destes aspectos, a fala das crianças apresenta fortes variações provocadas pelas alterações de crescimento e desenvolvimento, que são atribuídas a diferenças anatômicas e morfológicas na geometria no tracto vocal. Adicionalmente, as crianças têm, de um modo geral, menor controlo dos articuladores e variações de prosódia muito acentuadas. Estudos foram realizados relacionando as frequências fundamentais e dos formantes com a idade [20], impulsionando a aplicação de algumas técnicas para o desenvolvimento de aplicações mais robustas.

Estas variações apresentam-se como dificuldades na obtenção de bons resultados quando se pretende fazer o reconhecimento automático da fala das crianças. Algumas técnicas podem ser aplicadas para minimizar estes efeitos como é o caso do VTLN (*Vocal Tract Length Normalization*) ou ‘frequency warping’ em que se faz uma distorção em frequência dos parâmetros espectrais de forma a apresentarem valores de referência. A utilização de modelos para normalização da frequência dos formantes em vogais de crianças para os níveis de adultos foi considerada em alguns estudos [21], [22], mostrando resultados aquém dos obtidos com fala



de adultos, particularmente para crianças mais novas. Mostram ainda diferenças nos valores dos formantes no que diz respeito ao género, durante o período de desenvolvimento.

Outra destas técnicas de normalização é conhecida como MAPLR (*Maximum A Posteriori Linear Regression*), os resultados produzidos, apesar de revelarem melhorias notórias no desempenho, situam-se aquém dos desejáveis [19] [23] [24] [17]. O estudo [25] representa um bom resumo dos avanços conseguidos nesta área.

### 2.3. Proposta do Sistema a Desenvolver

A proposta de criação de um ‘Sistema para Terapia da Fala Auxiliada por Computador’ [26] serviu como ponto de partida da presente dissertação. O sistema pretende auxiliar (não substituindo) o terapeuta de fala, na identificação de dificuldades de articulação de fonemas, fazendo para isso uso da tecnologia do reconhecimento automático de fonemas num ambiente lúdico. O esquema do sistema proposto pode ser observado na figura a seguir apresentada.

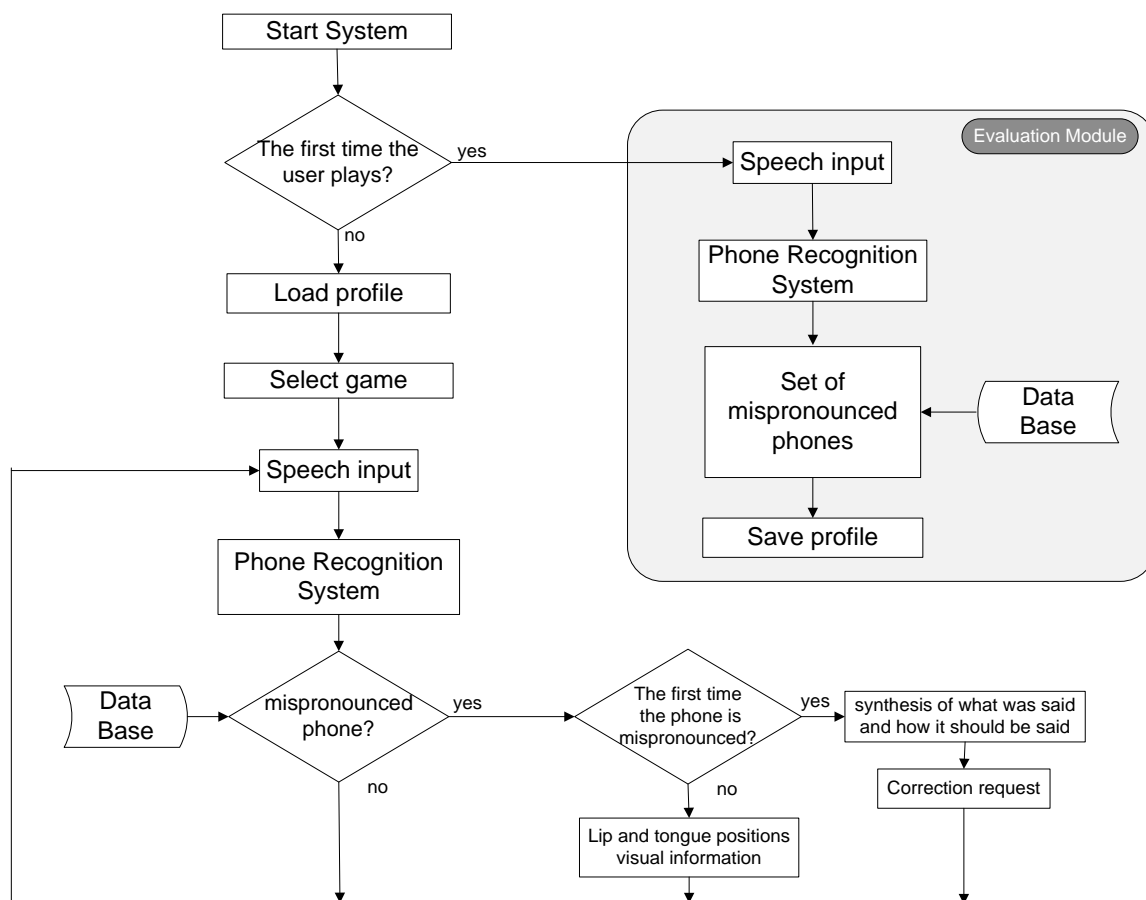


Figura 2. 4– Diagrama de blocos do sistema proposto [26].

Este sistema propõe vários módulos que se articulam de modo a poder analisar a fala de uma criança, identificando os seus problemas de dicção, e fornecer auxílio descritivo (informação visual) para a correcção das disfluências detectadas.

No presente estudo foi implementado o módulo de reconhecimento de fones para fala de crianças (através de modelos gerados e treinados para o efeito) e um módulo de verificação da correcta pronúncia de palavras.

A aplicação de um sistema de reconhecimento de fones para fala de crianças revela-se assim de inegável importância na elaboração e desenvolvimento do sistema. Para a correcta identificação de dificuldades de articulação das crianças em idade pré-escolar, é essencial que o reconhecimento do que a criança disse seja o mais preciso possível e que o sistema de verificação detecte correctamente os acertos e erros de articulação.

Para a criação do referido sistema foi necessário criar modelos para o reconhecimento de fones, especificamente para a faixa etária em questão. Estes modelos foram gerados a partir de uma base de dados de fala de crianças recolhida no âmbito de um projecto anterior [26]. A descrição desta base de dados será detalhada em pontos posteriores deste estudo (secções 4.1 e 4.1.1).

Torna-se, neste contexto, importante definir o conceito de fone e fonema. No contexto da fonética, um fone é a unidade mínima acústica do sistema fonético, que corresponde à realização física ou real de um fonema. São os sons que são distinguíveis pelo seu contraste dentro de uma palavra. Um fonema representa então uma abstracção do fone, sendo considerado a unidade básica de uma língua. [27] [28].

### 3. Técnicas de Análise e Reconhecimento da Fala

O problema do reconhecimento automático de fala não se encontra ainda completamente resolvido, particularmente no que respeita à fala contínua (por oposição ao reconhecimento de palavras isoladas). Porém, como foi já referido no capítulo anterior, apresenta bons resultados (nomeadamente para fala de adultos). Neste capítulo faz-se uma descrição sumária sobre uma das técnicas mais utilizadas em reconhecimento automático de fala, a modelação com os modelos de Markov não observáveis, bem como a técnica de parametrização espectral dos sinais de fala.

#### 3.1. Coeficientes Cepstrais em Escala Mel

O *cepstrum* em escala mel é um esquema de representação de sinais largamente utilizado na análise de sinais de fala. Tem-se vindo a tornar uma escolha bastante popular para extracção de características devido à sua capacidade de caracterizar eficientemente o sinal de fala [29].

A definição original de ‘*Cepstrum*’ foi proposta em 1963 [30] e consiste na transformada inversa de Fourier do logaritmo do espectro estimado de um sinal, sendo os seus coeficientes calculados através da seguinte fórmula:

$$c_x[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log |X(e^{j\omega})| e^{j\omega n} d\omega \quad (1)$$

Os coeficientes cepstrais em escala mel<sup>6</sup>, conhecidos pela sigla MFCC (*Mel-Frequency Cepstrum Coefficients*), diferem dos normais coeficientes de *cepstrum* em dois aspectos essenciais: o uso de uma descrição espectral baseada num banco de filtros e no facto de as bandas de frequência destes filtros se encontrarem igualmente espaçadas na escala mel. A expressão (2), proposta por Douglas O’Shaughnessy [31], estabelece a conversão de frequência de uma escala linear para a escala mel. Este mapeamento aproxima-se da resposta do sistema auditivo humano que tem uma relação não linear entre a percepção do tom e a frequência, sendo esta relação aproximadamente linear até cerca de 700Hz e logarítmica para frequências superiores [29].

---

<sup>6</sup> Escala perceptual de tons onde os ouvintes definem os tons espaçados igualmente.

$$m = 2595 \log_{10} \left( 1 + \frac{f}{700} \right) \quad (2)$$

Os filtros apresentam uma forma triangular cobrindo toda a gama de frequências de interesse. De notar ainda o facto de as larguras de banda dos filtros serem crescentes com o aumento da frequência segundo a escala mel, adaptando-se, assim, à diminuição de resolução de frequência do ouvido humano.

Os coeficientes MFCC são, desde há longa data, os parâmetros escolhidos em sistemas de reconhecimento automático de fala, uma vez que se traduzem numa descrição compacta do sinal de fala com boas características de redução de ruído. Devido a estas características, foram então os parâmetros escolhidos para a análise das locuções de fala das crianças.

## 3.2. Hidden Markov Models

Os modelos HMM (Hidden Markov Models) são modelos estatísticos apresentados inicialmente por Baum [32] e seus colaboradores em 1966, cuja ampla e reconhecida aplicação em reconhecimento automático da fala foi em muito impulsionada por Rabiner nos seus estudos [33] [34].

O sinal de fala é um sinal não estacionário que apresenta variações na sua composição espectral. Os modelos de Markov podem ser usados para descrever esta variação temporal de propriedades do sinal, assumindo um conjunto finito de configurações espectrais ou estados. Um HMM é um processo que descreve um sinal como uma cadeia de  $N$  estados ( $S_1, S_2, \dots, S_N$ ), alterando o seu estado a intervalos regulares de tempo, podendo mudar de estado ou permanecer nele próprio. A probabilidade de transição de um estado para cada um dos outros possíveis estados encontra-se definida numa matriz de distribuição de probabilidade – Matriz de Transição (também conhecida como Matriz de Markov -  $M$ ) – com elementos  $a_{ij}$  (a probabilidade de transitar do estado  $i$  para o estado  $j$ ). A transição para um novo estado fica assim, apenas dependente do estado anterior. A figura seguinte representa um modelo HMM simples com as possíveis transições e probabilidades associadas, sequências observadas ( $o_1, o_2, \dots, o_6$ ) e suas probabilidades de observação ( $b_j(o_t)$ ).

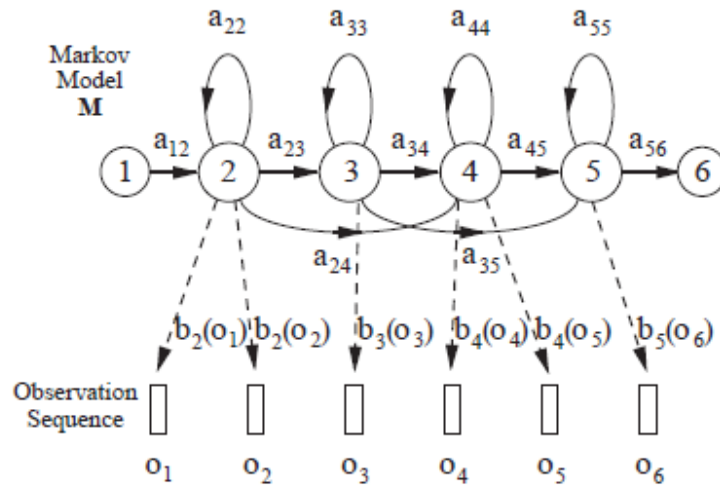


Figura 3.1 – Exemplo de um modelo HMM simples com topologia esquerda-direita [35].

Assim, cada estado pode gerar uma sequência de observações, de entre um conjunto, de acordo com uma determinada distribuição que é associada a cada estado do modelo. É portanto possível obter a mesma sequência de observações, com probabilidades diferentes, através de diferentes sequências de estados. No caso dos Modelos Não Observáveis (*Hidden Markov Models*), a sequência de estados do sistema não é conhecida, apenas a sequência de observações gerada. Define-se então  $b_j(x_t)$  como essa distribuição de probabilidade associada a cada estado: a distribuição condicional de observar o vector de características do sinal (observação  $x_t$ ) no instante  $t$  para o estado  $j$ . Para sistemas modelados por distribuições contínuas, a função densidade de probabilidade contínua é composta por uma mistura pesada de componentes Gaussianas (GMM – *Gaussian Mixture Model*).

Um HMM é, portanto, um duplo processo estocástico com uma cadeia de Markov e funções densidade de probabilidade por estado, que pode ser representada pelas equações (3) e (4), onde:

- $Q$  representa a dimensão do vector de observação  $\mathbf{x}$ ;
- $M$  - é o número de componentes Gaussianas;
- $\boldsymbol{\mu}_{jm}$  - é o vector de médias da componente  $m$  do GMM do estado  $j$ ;
- $\boldsymbol{\Sigma}_{jm}$  - é a matriz de covariância da componente  $m$  do GMM do estado  $j$ ;
- $c_{jm}$  - é o peso de cada componente da mistura;
- $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{jm}; \boldsymbol{\Sigma}_{jm})$  - é a função de densidade normal (Gaussiana).

$$b_j(\mathbf{x}) = \sum_{m=1}^M c_{jm} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{jm}, \boldsymbol{\Sigma}_{jm}), \quad \sum_{m=1}^M c_{jm} = 1 \quad (3)$$

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{jm}, \boldsymbol{\Sigma}_{jm}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}_{jm}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{jm})^T \boldsymbol{\Sigma}_{jm}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{jm})\right) \quad (4)$$

Para efeitos de aplicação em sistemas ASR, a topologia usada é a esquerda-direita, onde apenas são autorizadas transições para o mesmo estado ou para o estado subsequente, o que é facilmente perceptível quando se analisa o processo de produção da fala como uma sequência de fonemas.

Considera-se ainda que a produção da sequência de observações  $\mathbf{x}$  resulta de um HMM, onde apenas é observável a sequência emitida por cada estado e não a sequência de estados inerente ao modelo. Torna-se então possível estimar iterativamente os parâmetros do modelo usando um número considerável de sequências de observação e, por conseguinte, estimar a mais provável sequência de estados e a probabilidade de esta ter sido gerada por um dado modelo.

A figura seguinte apresenta um modelo HMM de topologia esquerda-direita para representar um fonema, com GMMs e observações associadas aos estados.

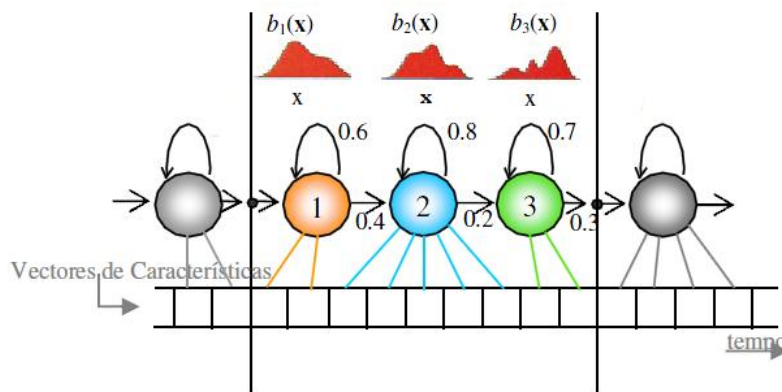


Figura 3. 2 – Exemplo de um modelo HMM com topologia esquerda-direita e com funções densidade de probabilidade compostas por GMMs [26].

Neste contexto surge o algoritmo de Viterbi [36] [37], de relevante importância neste estudo pois é usado na ferramenta HVite, facultada pelo HTK (*Hidden Markov Model Toolkit*)<sup>7</sup>, para efeitos de alinhamento e reconhecimento. Este algoritmo de programação dinâmica calcula a sequência de estados mais provável para uma dada sequência de observações (sendo essa

<sup>7</sup> Toolkit criado na Universidade de Cambridge para a geração e manipulação de HMM.

sequência conhecida como o ‘caminho de Viterbi’), podendo então ser considerado um algoritmo de decodificação da máxima verossimilhança. O algoritmo acha o caminho mais provável através de uma *trellis*, um gráfico de estados e tramas, ou seja, calcula o caminho mais curto num conjunto de observações, como exemplifica a figura seguinte.

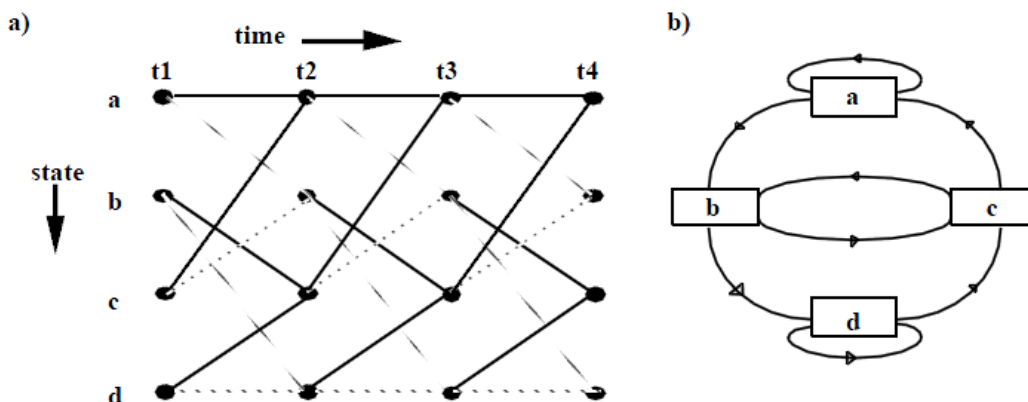


Figura 3. 3 – Exemplo de um diagrama *trellis* (a) e do HMM que o originou (b) [38].

Cada ponto na imagem anterior representa a probabilidade de um estado ser observada naquele instante e os arcos representam a probabilidade de cada transição. A probabilidade de cada caminho é então calculada através da soma das probabilidades de transições e as probabilidades de *output* ao longo do caminho.

De relevante importância é também o algoritmo Baum-Welch [39], que calcula os parâmetros de um HMM através do algoritmo *forward-backward*. Este algoritmo usa o critério de maximização da verossimilhança para estimar os parâmetros. Em geral, para um conjunto fixo de dados e o modelo estatístico subjacente, o método de máxima verossimilhança selecciona valores dos parâmetros do modelo, que produzem uma distribuição que fornece os dados observados com maior probabilidade (ou seja, os parâmetros que maximizam a função de probabilidade). O valor do logaritmo da verossimilhança (*log-likelihood*) é então calculado usando a função de densidade de probabilidade, como se pode observar na equação seguinte.

$$\sum_{i=1}^n \ln(f(x_i | \theta)) \quad (5)$$





## 4. Reconhecimento de Fala de Crianças

Como dito anteriormente, um aspecto essencial deste trabalho é a criação de modelos acústicos de fonemas, aplicados a fala de crianças, para o sistema de verificação de disfluências. Para criar estes modelos é necessária uma base de dados de fala de crianças. A secção seguinte descreve a base de dados usada e as outras secções descrevem o processo de criação de modelos.

### 4.1. Recolha de uma Base de Dados de Fala de Crianças

A inexistência de corpora para fala de crianças em português Europeu propulsionou a criação de uma base de dados que pudesse servir o propósito da definição de modelos para sistema como o supra-referido. Esta base de dados, usada no presente trabalho, foi recolhida em 2010, no âmbito de um estudo anterior, e está descrita com detalhe em [26]. Os dados de áudio foram recolhidos em infantários da zona centro de Portugal com um total de 111 crianças, 55 do sexo masculino e 56 feminino, com idades compreendidas entre os 5 e os 6 anos. A cada criança foi pedido que identificasse 40 palavras (representadas por imagens) escolhidas aleatoriamente de um conjunto de 200 palavras, o que resultou num total de 3726 ficheiros válidos. A frequência de amostragem usada foi de 8 kHz. A figura seguinte representa o histograma da distribuição dos diferentes fones presentes na base de dados.

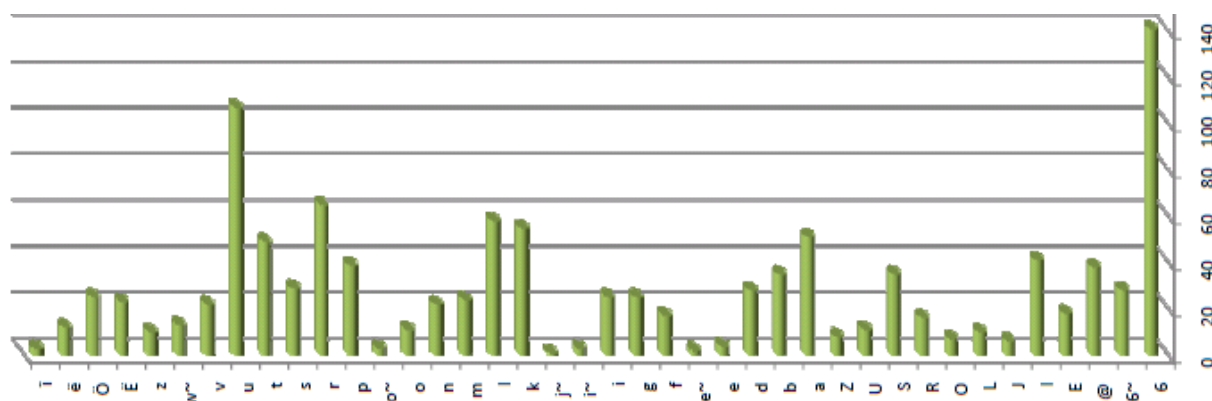


Figura 4.1 - Ocorrência de cada fone na lista de 200 palavras.

As condições da recolha das locuções e o facto de se lidar com crianças comprometeram, de alguma forma, a qualidade e as condições acústicas das gravações, o que implicou uma anotação da base de dados com bastante detalhe. No âmbito deste projecto as etiquetas ortográficas desta base de dados foram reanalisadas com o intuito de as verificar e eliminar inconsistências.

### 4.1.1 Tratamento da Base de Dados

Inicialmente, todos os ficheiros de áudio ('wav') foram filtrados com um filtro passa-alto com frequência de corte de 100 Hz, num script de Matlab, com o intuito de lhe retirar parte dos ruídos de baixa frequência. Também, por existirem inconsistências em termos temporais e de etiquetas atribuídas, e até alguns erros ortográficos, houve necessidade de rever todos os ficheiros da base de dados. Para tal foi usado o um *software* de anotação, em ambiente Matlab ('Viewlabs', ver Figura 4. 2), que permitiu a audição, visualização e edição das etiquetas e dos seus tempos.

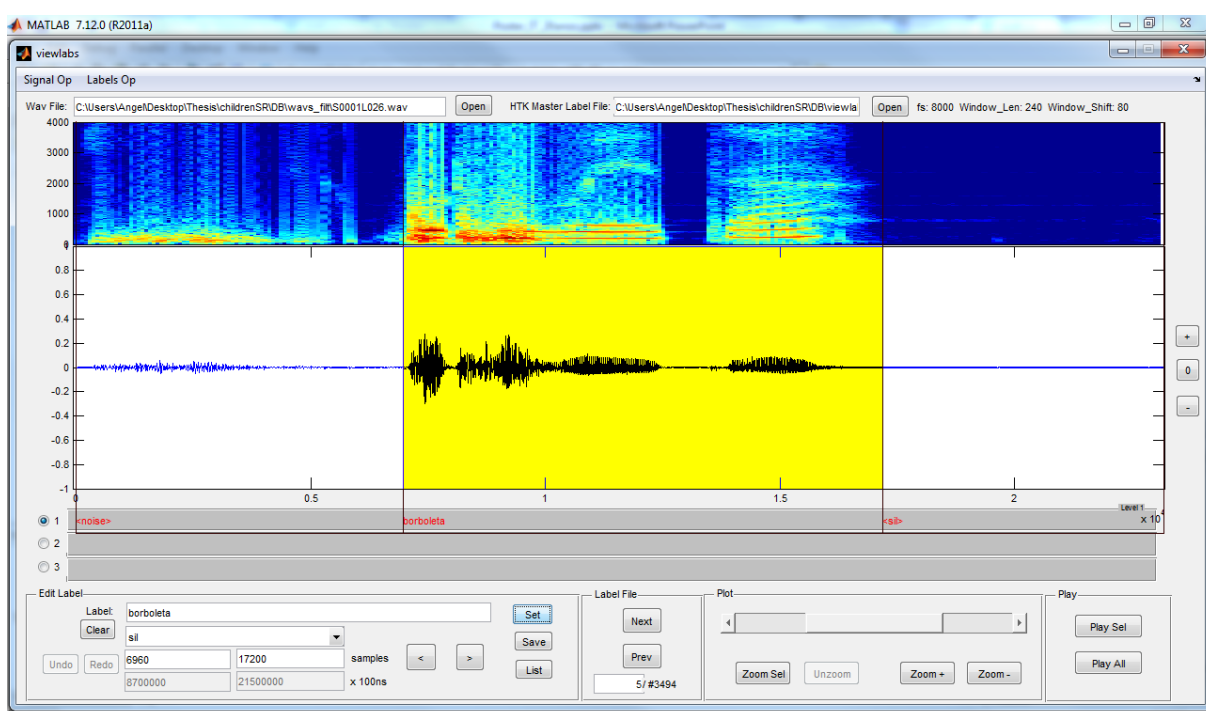


Figura 4. 2 – Viewlabs em execução com um ficheiro de áudio e as respectivas etiquetas e marcações temporais

Os ficheiros de anotação seguem um formato definido pelo *software* de criação de modelos acústicos, o HTK. Trata-se de ficheiros de texto com indicação de tempos iniciais e finais para cada etiqueta. A cada ficheiro de áudio corresponde um ficheiro de etiquetas com o mesmo nome e uma extensão diferente ('.lab'). É possível compilar as transcrições de muitos ficheiros de áudio num único ficheiro, chamado ficheiro "mestre" (MLF – *Master Label File*). Este ficheiro possui, como se pode observar na Figura 4. 3, uma lista de transcrições '.lab' (a indicação de início de ficheiro é representada pelo nome entre aspas). Cada uma das etiquetas presente nos ficheiros é representada por um triplete (tempo inicial, tempo final e a própria etiqueta) sendo o fim de cada transcrição indicado por um ponto final.

```

#!MLF!#
"/s0001L047.lab"
0          7600000  <sil>
7600000    17000000  cenoura
17000000   19200000  <breath>
19200000   25998750  <sil>
.
"/s0001L060.lab"
0          9800000  <sil>
9800000    20500000  copo
20500000   26998750  <sil>
.
"/s0001L068.lab"
0          9500000  <sil>
9500000    20300000  espada
20300000   25000000  <noise>
25000000   28998750  <sil>
.

```

Figura 4.3 – Exemplo de um ficheiro MLF (excerto).

No âmbito do presente trabalho todos os ficheiros áudio foram atentamente ouvidos, sendo as respectivas etiquetas ortográficas verificadas. Nesta fase foram feitas correcções, nomeadamente algumas etiquetas foram redefinidas e os seus limites temporais reajustados.

A notação usada nas etiquetas foi normalizada para facilitar o seu processamento. As etiquetas são definidas de forma ortográfica com a transcrição das palavras ditas pelas crianças. Além disso, foram usadas outras etiquetas para representar eventos extralinguísticos, resultando na seguinte lista de etiquetas com a respectiva correspondência em termos de modelos de eventos.

- <breath> resp
- <garbage> noise
- <laugh> noise
- <noise> noise
- <resp> resp
- <resp\_exp> resp
- <resp\_insp> resp
- <sil> sil

Adicionalmente, procedeu-se à catalogação das locuções quanto à sua qualidade (condições de ruído e/ou inteligibilidade) através de um processo manual onde se adicionaram caracteres específicos no fim das transcrições de modo a poder construir as marcações de ruído. Deste

processo resultaram diferentes versões do ficheiro de etiquetas (de acordo com o nível de ruído e de inteligibilidade).

O último e mais limpo ficheiro MLF (em termos de qualidade e percepção do áudio, com o nome 'ChildCAST\_clean.mlf') tem 2650 transcrições, tendo sido removidas todas as locuções/transcrições que possuíam demasiado ruído e cuja inteligibilidade era baixa <sup>8</sup>.

A partir deste ficheiro MLF foi produzida uma lista com as ocorrências de todas as etiquetas, incluindo as palavras correctas (palavras bem pronunciadas), erros de pronúncia/articulação, bem como as palavras de ligação (artigos, conjunções e verbos, como por exemplo em "é uma bruxa") e sons/palavras de hesitação (ver *script* 'formatmlf.m', presente na secção de anexos 7.4).

Um dicionário com todas as palavras da lista (num total de 522 ocorrências, das quais 346 foram consideradas palavras correctas) e sua tradução no alfabeto fonético SAMPA<sup>9</sup>, foi então gerado. O alfabeto tem 33 fones que estão listados no anexo 7.2. A transcrição em SAMPA foi inicialmente feita com o auxílio do conversor de grafemas para fonemas (G2PConverter.bat) entretanto desenvolvido no laboratório [40]. Posteriormente todas as transcrições foram manualmente verificadas e corrigidas quando necessário. A este dicionário foram adicionadas as etiquetas extralinguísticas e variações que permitem multipronúnciação em algumas palavras, como demonstra o exemplo da figura seguinte.

caramelo	k <sup>a</sup> r <sup>a</sup> m e l u
caranguejo	k <sup>a</sup> r <sup>a</sup> ~ g <sup>a</sup> Z u
caranguejo	k <sup>a</sup> r <sup>a</sup> ~ g <sup>a</sup> i Z u
carrinha	k <sup>a</sup> R i J <sup>a</sup>
carro	k a R u
carrocel	k <sup>a</sup> R O s E l

Figura 4. 4 – Excerto do dicionário com tradução SAMPA e com exemplo de multipronúnciação.

#### 4.1.2 Parametrização da Base de Dados

<sup>8</sup> Este resultado foi obtido com o script de Matlab 'cleannoisemlf.m' (ver anexo), que remove as etiquetas com as marcações de ruído, seguido pelo script 'emptynoisemlf.m', que remove do MLF todas as etiquetas que contêm apenas marcas acústicas (sem locuções).

<sup>9</sup> SAMPA (Speech Assessment Methods Phonetic Alphabet) – sistema de escrita fonético que tem como base o Alfabeto Fonético Internacional (IPA).

Procedeu-se então à fase de parametrização dos ficheiros de áudio com MFCC. Devido às características do sinal de fala típico de uma criança, cujos formantes e tom são muito superiores aos de um adulto, optou-se por um espaçamento de filtros de 143 mel, o que se traduziu na aplicação de 14 filtros, na escala mel desde 100 Hz até metade da frequência de amostragem - 4 kHz (ver Figura 4. 5 e Figura 4. 6). Este valor do espaçamento foi calculado dividindo a gama de frequências de interesse pelos 14 filtros. Abaixo de 100 Hz considera-se que as características espectrais são irrelevantes.

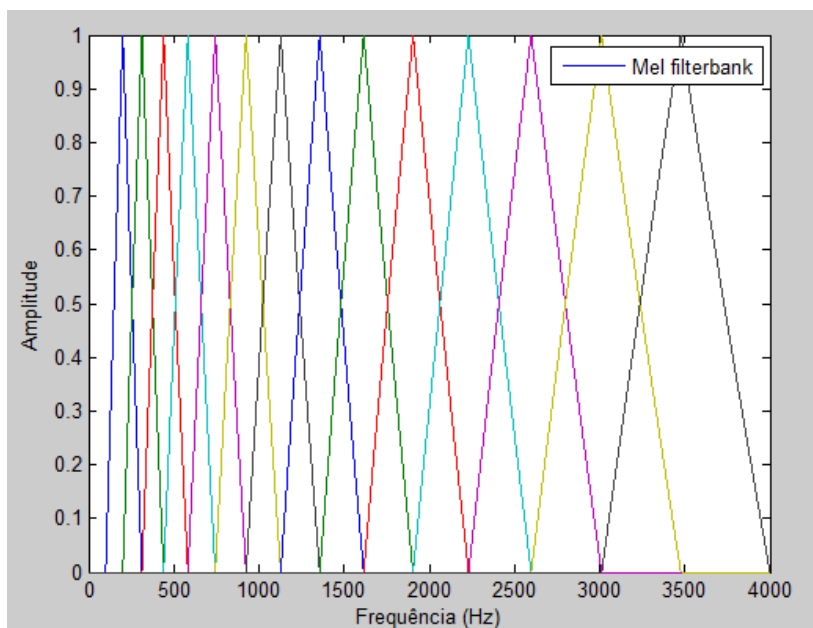


Figura 4. 5 – Representação dos 14 filtros na frequência.

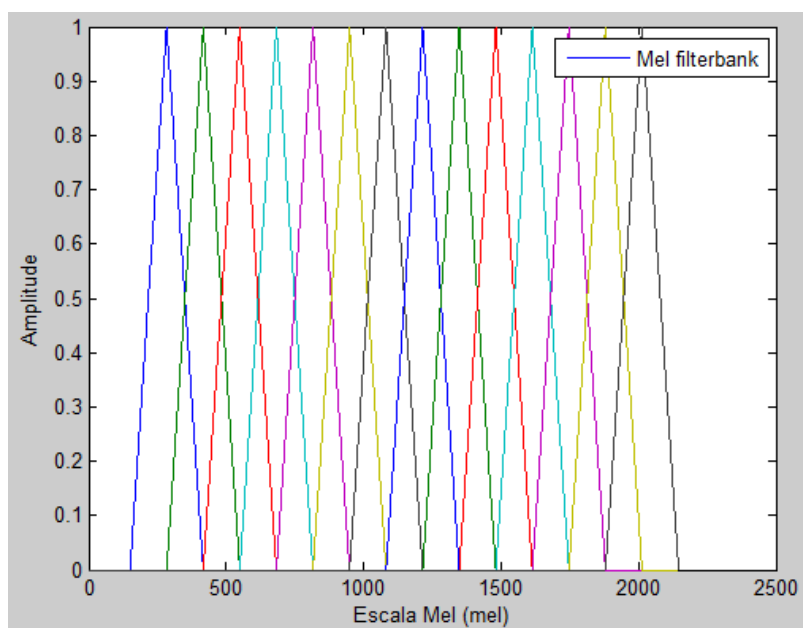


Figura 4. 6 - Representação dos 14 filtros da figura anterior em escala MEL.

Numa primeira abordagem, utilizou-se o Praat<sup>10</sup> [41] para calcular os MFCCs. Foram criados *scripts* para correr todos os ficheiros áudio e, para cada, recolher a informação relativa à harmonicidade<sup>11</sup>, tom e os próprios MFCCs.

Pelo facto de a ferramenta de reconhecimento escolhida, onde se iria proceder ao grosso do trabalho, ter sido o HTK, considerou-se mais proveitoso proceder ao cálculo dos coeficientes MFCC com as ferramentas providenciadas pelo HTK, evitando-se assim possíveis problemas de concordância de formatos ou outros.

Para o cálculo dos coeficientes foi escolhida uma janela de Hamming<sup>12</sup> de 15 ms de duração, 14 filtros (com forma triangular) e passos (avanço da janela) de 10 ms para uma obter o *frame-rate* tradicionalmente usado de 100 tramas por segundo. O número de coeficientes cepstrais usado foi 12, o que também é usual em processamento da fala. O uso de janelas mais curtas que o usual deve-se a proporcionar uma melhor adaptação à fala de crianças que têm tom mais alto.

Além dos 12 parâmetros MFCC e do coeficiente de ordem zero,  $c_0$ , os vectores de características associados a cada trama de sinal de fala têm também parâmetros ditos "dinâmicos". Estes correspondem a uma derivada temporal discreta (inclinação de regressão linear numa janela de 5 tramas) de cada componente do vector. São designados usualmente por parâmetros "DELTA". A variação dos parâmetros DELTA (parâmetros de aceleração) é também incluída. Desta forma os vectores de parâmetros têm dimensão 39 (13 coeficientes MFCC,  $12+c_0$ , mais 13 coeficientes DELTA e mais 13 coeficientes de aceleração).

Em anexo apresenta-se o ficheiro '`config_Hcopy.txt`', com os parâmetros de configuração da ferramenta HCopy do HTK. Esta ferramenta HCopy gera os ficheiros com os parâmetros MFCC a partir dos ficheiros áudio e das etiquetas correspondentes no ficheiro MLF<sup>13</sup>.

O conjunto de ficheiros MFCC correspondentes às locuções de fala foi dividido em dois conjuntos: um conjunto de 'treino' e de 'teste' dos modelos de fones. A divisão foi de quatro quintos para o conjunto de 'treino' e o restante um quinto para o conjunto de 'teste'. O script '`mfclist.m`' gerou as duas listagens de ficheiros de treino e de teste.

---

<sup>10</sup> Software (freeware) concebido por Paul Boersma e David Weenink para análise de fala em termos fonéticos.

<sup>11</sup> Representa o grau de periodicidade acústica, também chamado de Harmonics-to-Noise Ratio (HNR).

<sup>12</sup> O 'cosseno-elevado' desta janela serve para minimizar as laterais do lobo, enfatizando a zona central.

<sup>13</sup> Parametriação indicada no ficheiro de configuração (Anexo A - 7.1) como:

TARGETKIND = MFCC\_0\_D\_A

O ficheiro MLF foi seguidamente separado em dois novos ficheiros MLF, um com as transcrições e outro contendo apenas as etiquetas extralinguísticas, através do *script* 'cutmlf.m'. Paralelamente à criação dos dois MLF, foram também gerados através do mesmo *script* novos ficheiros de MFCC de acordo com os ficheiros separados MLF e as suas respectivas marcações temporais das etiquetas.

Optou-se por esta abordagem de separar as locuções da restante componente não-falada dos ficheiros de parâmetros em diferentes ficheiros MLF e, por conseguinte, em diferentes grupos de ficheiros MFCC, devido à elevada complexidade na sequência da grande maioria dos ficheiros de áudio. A generalidade das sequências possuíam padrões consideravelmente mais complexos do que o normal '<etiqueta> + palavra + <etiqueta>'. Durante as gravações de áudio, as crianças por diversas vezes hesitavam ou alteravam o que estavam a dizer a meio da locução, riam audivelmente ou respiravam pesadamente e, somente em raras ocasiões, diziam unicamente a palavra requerida, resultando assim em anotações com sequências de etiquetas das mais variadas e estranhas formas.

Uma vez cortados os ficheiros, procedeu-se à criação de grupos de 'treino' e de 'teste' para os modelos, a partir dos recentemente gerados ficheiros MLF. Assim, dois novos ficheiros foram criados contendo uma lista dos ficheiros MFCC que correspondem às etiquetas existentes no ficheiro MLF. O script 'mfclist.m' gerou os dois referidos ficheiros, dividindo o número total de etiquetas no MFL e atribuindo quatro quintos dos correspondentes MFCCs à lista de 'treino' e o restante quinto como uma lista de 'teste'.

## **4.2. Sistema de Reconhecimento**

### **4.2.1 Treino dos Modelos**

O processo de criação e treino dos modelos de fones seguidamente descrito, foi baseado e adaptado do *tutorial* descrito no 'HTK Book 3.4' [35] e as instruções presentes na página da VoxForge [42].

A informação contida nas etiquetas dos ficheiros MLF são palavras. Considerando o objectivo do presente estudo de avaliar a correcta articulação de palavras na sua sequência de fones, tornou-se necessário transcrever as palavras nos seus fones constituintes. Para isso bastou usar o

dicionário e manipular o ficheiro de etiquetas. Assim, dando o ficheiro MLF, foi criada uma nova versão do mesmo mas, desta feita, com um fone por linha e dividindo os tempos (do início ao fim da palavra) em partes iguais por cada fone. Foi ainda gerada uma outra versão do mesmo ficheiro, também com um fone por linha, mas sem indicações temporais. Isto foi feito usando a ferramenta HLEd do HTK. Na figura seguinte mostra-se um exemplo do resultado da aplicação desta ferramenta. Neste caso, a ferramenta tomou a locução do ficheiro como ‘copo’ e procurou a sua tradução em SAMPA no ficheiro dicionário, gerando então o novo MLF com a tradução e dividindo equitativamente os tempos atribuídos a cada fone.

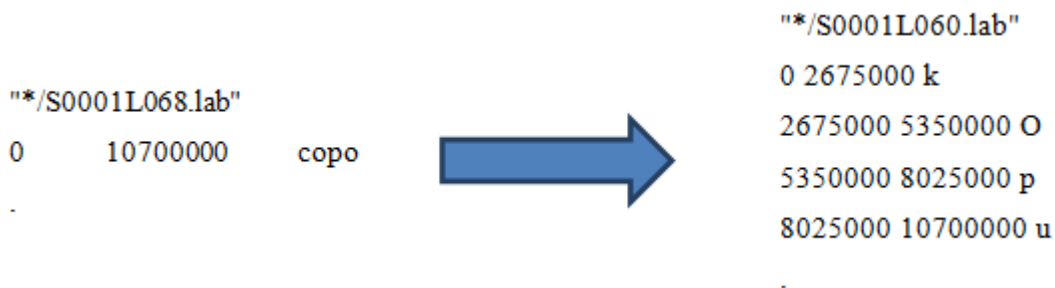


Figura 4. 7 – Exemplo do uso da ferramenta HLEd para a palavra ‘copo’ que tem a transcrição SAMPA ‘kOpu’.

O passo anterior é essencial para iniciar o processo de treino, uma vez que os modelos de reconhecimento escolhidos para este estudo são modelos de fones. Os modelos de fones iniciais foram gerados usando uma estratégia de *flat-start*, onde todos os modelos são iguais na fase inicial. Significa que todos os modelos de fones são iniciados com vectores de médias e matrizes de covariância iguais à média e variância global de todo o material acústico de treino. Dá-se o nome de protótipo (`'proto'`) a este modelo inicial, que também tem definida a topologia esquerda-direita dos modelos com 3 estados. Este *'flat-start'* é feito usando a ferramenta HCompV.

A partir destes novos modelos de fones, foi gerado um ficheiro MMF (*Master Macro File*) que é simplesmente um ficheiro com o conjunto de todos os modelos criados.

Uma vez criado este ficheiro MMF procedeu-se à re-estimação dos modelos em sete passos. Para tal foi usada a ferramenta HERest que executa o algoritmo Baum-Welch, re-estimando todo o conjunto de HMMs.

O passo seguinte consiste na escolha da pronúncia mais adequada sempre que existem variações de pronúncia autorizadas no dicionário de pronúncias. Isso é feito forçando um alinhamento entre as observações acústicas e as sequências alternativas de fones



correspondentes. A ferramenta HVite em modo de alinhamento é usada para este efeito. O HVite faz uso do algoritmo de Viterbi (descrito na secção 3.2) para escolher qual o modelo (quando existe alternativas) que mais provavelmente gerou aquela sequência de observações. Como foi já anteriormente mencionado, este algoritmo tenta calcular o melhor caminho através de uma grelha onde uma dimensão corresponde aos possíveis estados, e outra às tramas (tempo).

Após o alinhamento, foram executadas mais duas re-estimações aos modelos.

O processo acima descrito encontra-se esquematizado na figura seguinte.

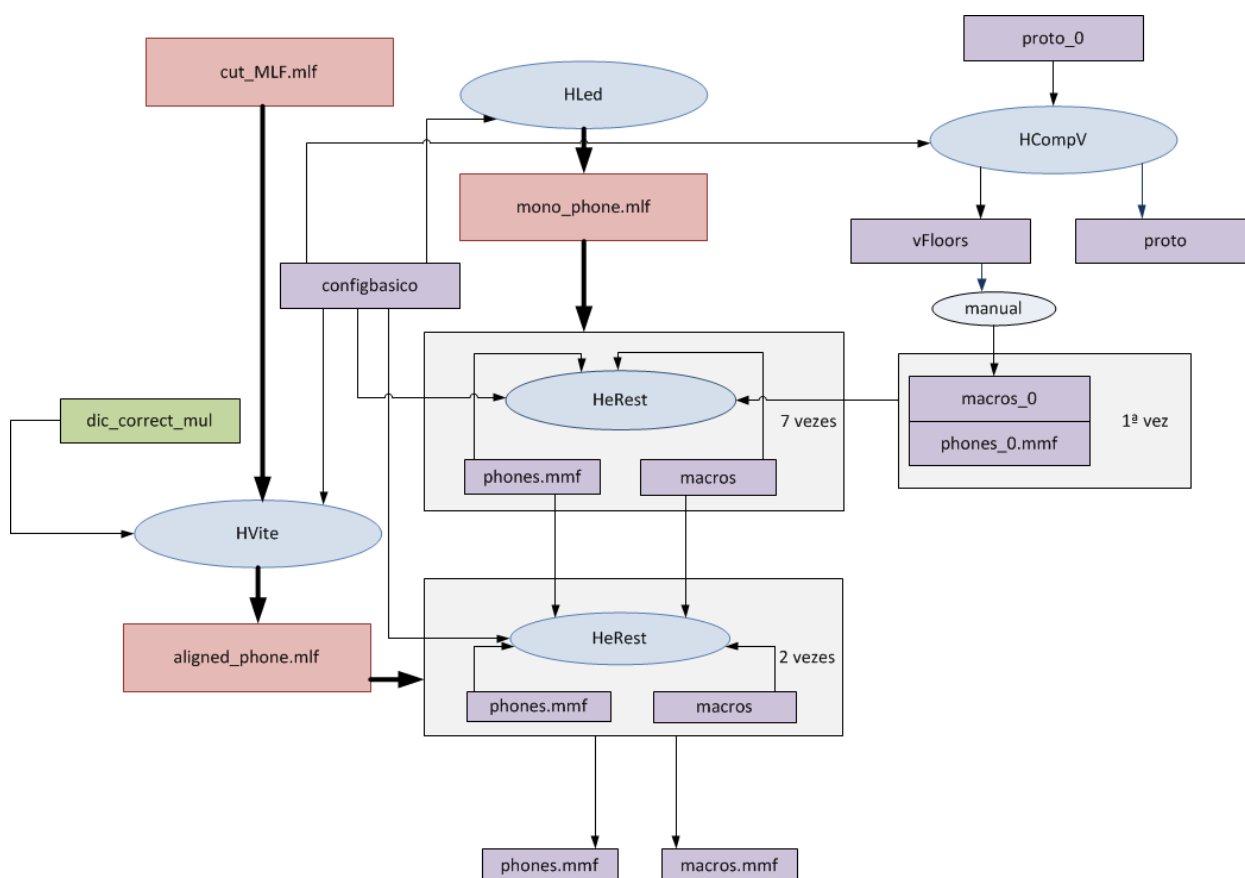


Figura 4. 8 – Diagrama do processo executado no HTK para treino dos modelos de fones. As elipses indicam ferramentas ou processos, sendo que a maioria são ferramentas do HTK estando também presente um processo manual. Os ficheiros MLF encontram-se representados a cor-de-rosa, enquanto os ficheiros utilizados pelo HTK se encontram representados a lilás. O dicionário com a lista de palavras correctas e sua correspondente transcrição em SAMPA está representado a verde. Este diagrama representa de uma forma concisa o funcionamento do script 'treino.bat' (secção de anexos 7.3).

O passo seguinte consiste no incremento das componentes Gaussianas das funções densidade de probabilidade de cada estado (GMMs). O número de Gaussianas foi incrementado fazendo uso da ferramenta HHed que executa um processo chamado *mixture splitting* quando recebe como

argumento "MU". Este processo foi executado quinze vezes para todos os modelos, fazendo três re-estimações entre cada incremento, de forma a obter GMMs com 16 componentes (secção de anexos 7.3).

#### 4.2.2 Teste dos Modelos

Com os modelos devidamente treinados, torna-se necessária uma avaliação preliminar da taxa de reconhecimento do sistema. Para efectuar reconhecimento, o sistema necessita conhecer *a priori*, as palavras que pode reconhecer, ou seja, uma gramática da qual faz parte uma lista de todas as palavras permitidas e as sequências autorizadas. Para tal, faz uso de uma "gramática" e da "rede" correspondente.

Para gerar a gramática, uma lista de palavras foi criada manualmente a partir da lista de todas as locuções recolhidas das etiquetas do MLF, seleccionando apenas as palavras correctas ou bem pronunciadas (lista igual à constante do ficheiro de dicionário `'dic_correct_mul'` presente no anexo 7.2). Esta lista foi atribuída à variável `'$word'` na figura seguinte. Também foi criada uma lista com as palavras/expressões de ligação mais frequentes, atribuída à variável `'$garbage'` na mesma figura. As sequências autorizadas são aquelas indicadas pelo grafo (ou rede) da figura, por exemplo, "Ahm um copo", onde "ahm" e "um" estão incluídos nos dois nodos opcionais do grafo. De notar que uma palavra do vocabulário tem de ocorrer uma e uma só vez, e sempre na última posição da sequência gerada.

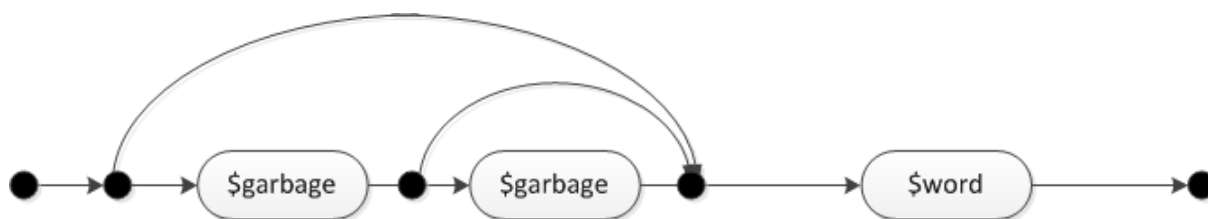


Figura 4. 9 – Diagrama da sequência autorizada no reconhecimento.

Com esta gramática, a ferramenta `HParse` pode então gerar o ficheiro de rede no formato SLF (Standard Lattice Format). Este formato serve para fornecer ao reconhecedor as múltiplas hipóteses de ocorrência de etiquetas e possui um conjunto de nós que representam as palavras e um conjunto de arcos que representam as possíveis transições entre palavras.

Fazendo uso dos modelos treinados, da lista de fones e da rede criada no passo anterior, o `HVite` pode então ser usado para efectuar o reconhecimento dos ficheiros de 'teste'. Assim, o

HVite pode gerar um ficheiro MLF de reconhecimento ('REC') com as sequências reconhecidas de acordo com a gramática. O ficheiro gerado no reconhecimento poderá ser comparado ao MLF com as transcrições reais para efeitos de avaliação de desempenho.

O HResults é uma ferramenta de análise de desempenho. Compara ficheiros de etiquetas, por norma o ficheiro de *output* produzido pelo HVite, com um ficheiro de referência das transcrições, calculando as estatísticas da comparação como a correcção e a precisão. Para o efeito usa um procedimento de alinhamento de *strings*. Um exemplo do seu *output* pode ser observado na figura seguinte.

```
===== HTK Results Analysis =====  
Date: Sat Sep 08 19:57:27 2012  
Ref : MLF_cut2_ltx.mlf  
Rec : ..\test\recout_FINAL.mlf  
----- Overall Results -----  
SENT: %Correct=79.66 [H=423, S=108, N=531]  
WORD: %Corr=78.58, Acc=78.58 [H=433, D=20, S=98, I=0, N=551]  
=====
```

Figura 4. 10 – Exemplo do *output* da análise estatística produzida pela ferramenta.

Nesta avaliação do desempenho dos modelos foi obtida uma taxa de correcção de 78.58% com 433 palavras reconhecidas correctamente (H - hits), 20 apagamentos (D - *deletions*), 0 inserções (I - *insertions*) e 98 substituições (S - *substitutions*), num total de 551 palavras reconhecidas. Em grande parte dos casos onde ocorreram erros de reconhecimento, a palavra dada como output é semelhante à palavra do MLF de referência (com o que foi realmente dito), como se pode constatar em alguns dos exemplos presentes na figura seguinte.

```
Aligned transcription: S0098L082.lab vs S0098L082.rec
LAB: garrafa vinho
REC:          passarinho
Aligned transcription: S0053L103.lab vs S0053L103.rec
LAB: livro
REC: livros
Aligned transcription: S0089L178.lab vs S0089L178.rec
LAB: talhere
REC: talheres
Aligned transcription: S0078L079.lab vs S0078L079.rec
LAB: fralda
REC: fada
Aligned transcription: S0061L120.lab vs S0061L120.rec
LAB: mola
REC: mala
```

Figura 4. 11 – Exemplo de erros de reconhecimento (substituições) presentes no ficheiro de *output* do *HResults*.

O MLF gerado segundo esta gramática possui apenas uma etiqueta (que será obrigatoriamente uma única palavra da lista de palavras consideradas correctas), por cada transcrição. O *HVite* ignora todos os resultados da lista  $\$garbage$ , não os incluindo no ficheiro MLF gerado (REC).

De relevante importância é ainda o facto de se ter adoptado a sequência da gramática apresentada na Figura 4. 9. A escolha desta gramática permite um reduzido peso computacional no reconhecimento porque diminui o número de possíveis sequências (arcos). No entanto, o custo associado é a redução da taxa de palavras correctamente reconhecidas. Como foi mencionado anteriormente, muitas das transcrições apresentavam sequências de etiquetas invulgares, existindo mesmo algumas com mais do que uma etiqueta com palavras consideradas correctas. O resultado global do reconhecimento sofre assim, em certo modo, visto estas transcrições não poderem ser correctamente reconhecidas com a gramática escolhida.

## 5. Verificação

Uma vez criado o sistema de reconhecimento de fones, pode então ser gerado o procedimento para avaliação e verificação das palavras proferidas de acordo com o que era esperado. Um estudo com objectivos semelhantes foi desenvolvido para a linguagem sueca. Foi proposta a criação de um sistema para segmentação e alinhamento automático de fala obtendo resultados que, apesar de promissores, revelam ainda bastantes falhas em termos da correcção do alinhamento produzido. [43]

Trata-se de um problema de teste de hipóteses e consiste em saber se o locutor disse ou não aquilo que lhe era pedido para dizer. Posteriormente, e em caso negativo, é também importante saber quais os fones que foram mal pronunciados.

O HVite, ao gerar o ficheiro de reconhecimento, pode também fornecer a informação das verosimilhanças associadas a cada palavra ou fone reconhecido. Este valor do logaritmo da verosimilhança de palavra calculado pelo HVite será referido deste ponto em diante como  $L_1$ . De notar o facto de este valor de  $L_1$  ser negativo, portanto, quanto mais próximo de zero for esse valor, maior será a verosimilhança calculada.

Como medida de desempenho para o teste de hipóteses é usada uma diferença de verosimilhanças obtidas com o algoritmo de Viterbi. A primeira dessas verosimilhanças,  $L_1$ , é aquela obtida com o reconhecedor segundo a gramática da Figura 4. 9. A outra verosimilhança,  $L_2$ , é a obtida através da gramática representada pela figura seguinte.

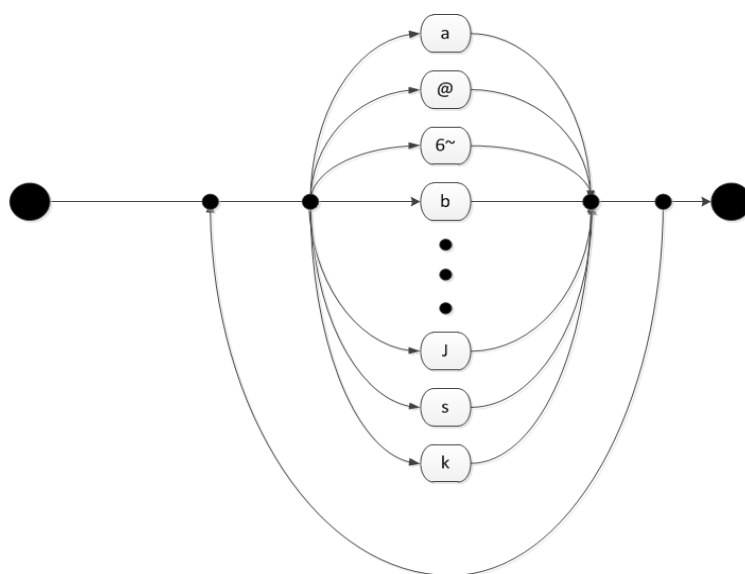


Figura 5. 1 – Diagrama da sequência livre de fones autorizada pela gramática livre.

Para a identificação e análise de disfluências e incorrecções nas locuções considerou-se que a estratégia mais vantajosa a adoptar seria a da utilização de uma gramática ‘livre’ (por oposição à gramática descrita na secção 4.2.2, na Figura 4. 9, a qual passará a ser referida como gramática ‘fixa’), onde o reconhecedor pudesse reconhecer uma e uma só palavra, com uma qualquer sequência de fones. A figura anterior descreve esta gramática ‘livre’. Assim, a sequência de fones gerada desta forma produz uma verosimilhança,  $L_2$ , sempre maior ou igual a  $L_1$  uma vez que a sequência de fones é óptima segundo os modelos.

Se a diferença de verosimilhanças  $L_2-L_1$  for baixa significa que a palavra foi bem pronunciada. A sequência de fones não tem de ser exactamente igual; podem até existir repetições, apagamentos ou inserções de alguns fones. No caso de uma diferença elevada, pode significar que o locutor não disse a palavra que era suposto dizer.

Para proceder à avaliação da veracidade da premissa criou-se um script de Matlab que analisasse todos os ficheiros MFCC contemplados no conjunto de teste, comparando o resultado obtido pelo reconhecedor com a o que foi inicialmente solicitado ao locutor (criança), imprimindo as informações relevantes num ficheiro. Temos assim 3 etiquetas potencialmente diferentes:

- REF: a palavra solicitada à criança;
- SAID: a palavra que a criança disse (possivelmente uma palavra mal pronunciada);
- REC: a palavra reconhecida automaticamente pelo sistema.

Assim, numa primeira instância, considerou-se que a criança disse exactamente o que lhe foi pedido (MLF de referência igual ao MLF das palavras a dizer), fazendo o reconhecedor passar pelas duas gramáticas anteriormente descritas (gramática ‘fixa’ e gramática ‘livre’) e comparando as verosimilhanças e as sequências de fones obtidas pelos dois processos.

REF	SAID	REC	phones	Freephones	Nframes	L1	L1-L2	align_score	refph_algn	freeph_algnt
coelho	coelho	coelho	ku <sup>a</sup> Lu	u <sup>v</sup> fboa <sup>a</sup> ~Lu	94	-5.245491e+001	-1.3038e+000	6	ku6 <sup>·</sup> ·Lu	ufboaãLu
caixa	caixa	caixa	kais <sup>a</sup>	kaS <sup>a</sup>	59	-5.522873e+001	-350.1587e-003	1	kaiS6	ka·S6
anel	anel	anel	<sup>a</sup> nEl	<sup>a</sup> nE <sup>a</sup> ~gr	49	-6.122777e+001	-30.1912e-003	3	6nEl <sup>·</sup> ·	6nEãgr
ovelha	ovelha	pilha	Ov <sup>a</sup> L <sup>a</sup>	O <sup>a</sup> fiL <sup>a</sup>	77	-5.395169e+001	-1.1648e+000	3	O <sup>·</sup> v6 <sup>·</sup> ·L6	O6fiL6
lâmpada	lâmpada	bola	l <sup>a</sup> ~p <sup>a</sup> d <sup>a</sup>	poObld <sup>a</sup> ~ <sup>a</sup>	113	-5.809059e+001	-1.1063e+000	6	lãp6 <sup>·</sup> d <sup>·</sup> 6	poObldã6
prato	prato	prato	pratu	pratu	44	-6.170283e+001	0.0000e+000	0	pratu	pratu
tomate	tomate	tomate	tumat@	RmaGt@	68	-5.971878e+001	-312.9497e-003	3	tuma <sup>·</sup> t@	·RmaGt@

Figura 5. 2 – Excerto do ficheiro gerado que compara as diferentes verosimilhanças e as strings de fones para ficheiros REF e SAID iguais.

O ficheiro gerado, cujo excerto pode ser observado na figura anterior, contém informação relevante a cada etiqueta:

- As transcrições nos diferentes ficheiros MLF: REF, SAID e REC (o que foi pedido à criança, o que ela de facto disse e o que o `HVite` reconheceu, respectivamente);
- A sequência de fones correspondente à palavra pedida segundo o dicionário (phones) e a sequência de fones dada pelo reconhecedor com recurso à gramática ‘livre’ (Freephones);
- O número de *frames* correspondente à palavra (Nframes);
- O valor de verosimilhança para a gramática com a tradução SAMPA da palavra –  $L_1$ ;
- O valor da diferença entre a verosimilhança anterior e o obtido com uma gramática livre:  $L_1-L_2$ ;
- Uma medida (`align_score`) de alinhamento entre as duas *strings* de fones: a sequência de fones da palavra e a sequência de fones ‘livre’. O alinhamento usa o algoritmo ‘edit distance’<sup>14</sup>.

Os valores de  $L_1$  e  $L_2$  apresentados no ficheiro representam a soma dos valores de verosimilhança calculados pelo `HVite` para cada fone da palavra em questão, normalizada pelo número de *frames* (Nframes).

No caso do uso da gramática ‘fixa’, o reconhecedor é forçado a calcular o caminho de Viterbi por uma sequência de fones específica, de modo a que corresponda exactamente à sequência ditada pela tradução SAMPA da palavra presente no dicionário. Se for proferido um fone não esperado na sequência, a verosimilhança será obrigatoriamente menor, isto é, um valor de  $L_1$  menor.

Já no caso da gramática ‘livre’, o reconhecedor gerará uma sequência de fones com os modelos mais parecidos com os parâmetros em análise e portanto, um caminho óptimo segundo os modelos.

---

<sup>14</sup> *Script* de Matlab desenvolvido no Laboratório de Processamento de Sinal que compara duas *strings* atribuindo pontos por cada inserção, supressão ou troca de carácter. O resultado obtido como *output* será tanto maior quanto maiores forem as diferenças entre as *strings*.

"/temp.rec"				"/temp.rec"			
0	300000	p	-247.054306	0	300000	p	-247.054306
300000	800000	r	-382.939423	300000	800000	r	-382.939423
800000	2000000	a	-607.079712	800000	2000000	a	-607.079712
2000000	3900000	t	-1149.680542	2000000	3900000	t	-1149.680542
3900000	4400000	u	-328.170441	3900000	4400000	u	-328.170441
.			(a)	.			(b)

Figura 5. 3 – Resultados produzidos pelo HVite para o ficheiro ‘S0054L152’ com recurso à gramática ‘fixa’ (a) e à gramática ‘livre’ (b), para ficheiros REF e SAID iguais.

No caso demonstrado na figura anterior, os parâmetros obtidos da locução da criança encontram-se perfeitamente enquadrados com a sequência de fones esperada para aquela palavra (‘prato’). Esse facto, conjuntamente com um valor total (do somatório) de verosimilhanças relativamente baixo, permite inferir que a articulação da palavra foi correcta.

No exemplo da figura seguinte é possível verificar que a dicção da palavra (‘coroa’) não foi perfeita, indicando o reconhecedor com a gramática ‘livre’ uma supressão (apagamento) do fone ‘u’ e uma inserção do fone ‘a’.

"/temp.rec"				"/temp.rec"			
0	800000	k	-576.471863	0	900000	k	-660.738403
800000	1100000	u	-251.689270	900000	1600000	r	-542.880615
1100000	1600000	r	-390.381073	1600000	3800000	o	-1222.837402
1600000	4000000	o	-1345.375000	3800000	5300000	a	-783.361755
4000000	6500000	a	-1556.693115	5300000	6500000	a	-816.888062
.			(a)	.			(b)

Figura 5. 4 – Resultados produzidos pelo HVite para o ficheiro ‘S0097L000’ com recurso à gramática ‘fixa’ (a) e à gramática ‘livre’ (b), para ficheiros REF e SAID iguais.

Por análise dos resultados compilados no ficheiro, é possível observar que, como expectável, os valores de  $L_2$  são sempre iguais ou superiores aos de  $L_1$ . Também se verifica uma tendência para que os menores valores no *score* do alinhamento (sequências com menos necessidade de alterações no alinhamento) correspondam a valores maiores de  $L_2$  (maior verosimilhança). É também notória uma diminuição nos valores de  $L_1$  e  $L_2$  (menor verosimilhança) com o aumento do número de fones e complexidade da palavra.



Posteriormente procedeu-se à mesma análise mas simulando que a criança não disse o que lhe foi pedido (os dados fornecidos para reconhecimento não correspondem ao que se diz ao reconhecedor que a criança disse). Para tal, gerou-se um ficheiro MLF a partir do REF, onde as transcrições foram alteradas de modo aleatório por outras palavras constantes do dicionário. Para o efeito foi criado o *script* 'changeMLF lab.m'.

REF	SAID	REC	phones	Freephones	Nframes	L1	L1-L2	align	score	refph	alqn	freeph	alqnt
cereais	coelho	coelho	s@riaiS	u~fboa^~Lu	94	-5.979505e+001	-8.6440e+000	7		s@riaiS	'	ufboaãLu	
gomas	caixa	caixa	gom^S	kaS^a	59	-6.165912e+001	-6.7806e+000	4		gom6S	'	kaS6	'
cubo	anel	anel	kubu	^nE^~gr	49	-7.215646e+001	-10.9589e+000	6		kubu	'	6nEãgr	
cenoura	ovelha	pilha	s@nor^a	O^fiL^a	77	-6.251457e+001	-9.7277e+000	5		s@nor6	'	O6fiL6	
chavina	lâmpada	bola	S^vin^a	poOblD^~^a	113	-6.204074e+001	-5.0565e+000	7		S6vin	'	poOblDã6	
maçã	prato	prato	m^s^a~	pratu	44	-6.562010e+001	-3.9173e+000	5		m6sã	'	pratu	
retrato	tomate	tomate	R@tratu	RmaGt@	68	-6.345029e+001	-4.0445e+000	5		R@tratu	'	RmaGt@	

Figura 5.5 – Excerto do ficheiro gerado que compara as diferentes verosimilhanças e as *strings* de fones para ficheiros REF e SAID diferentes.

As informações presentes no ficheiro gerado permitiram verificar uma diminuição generalizada (mais negativos) nos valores de  $L_1$  (e por consequência na diferença  $L_1-L_2$ ) e nos *scores* do alinhamento, o que corrobora a tese de que poderá ser possível identificar disfluências com base nos valores de verosimilhança.

No exemplo da figura seguinte é possível observar os resultados do reconhecedor, segundo as duas gramáticas, para uma etiqueta onde a palavra pedida (REF toma o valor 'carta') foi alterada de modo a não coincidir com o que a criança realmente disse (SAID) que, neste exemplo, foi a palavra 'praia'.

"/temp.rec"				"/temp.rec"			
0	1300000	k	-1041.198120	0	400000	p	-286.109436
1300000	2600000	a	-731.584534	400000	1400000	r	-700.334595
2600000	4000000	r	-945.395264	1400000	2600000	a	-667.383606
4000000	4300000	t	-288.709442	2600000	4100000	i	-927.261780
4300000	4900000	a	-428.117767	4100000	4900000	a	-550.425964
.			(a)	.			(b)

Figura 5.6 – Resultados produzidos pelo HVite para o ficheiro 'S0097L000' com recurso à gramática 'fixa' (a) e à gramática 'livre' (b), para etiquetas REF e SAID diferentes.

Procedeu-se então à análise ROC (*Receiver operating characteristic*) dos valores de  $L_1-L_2$  para se poder escolher um limiar de aceitação mais adequado ao problema de classificação em análise – ‘Aceitar ou não a locução como correcta.’. Em detecção de sinais, as curvas ROC são gráficos de descrição de um sistema de classificação binária e do seu desempenho [44].

Considerou-se então que a hipótese  $H_0$  é: ‘A criança disse correctamente o que lhe foi pedido.’, havendo quatro possíveis soluções, em função da resposta do reconhecedor, como demonstra o quadro da figura seguinte:

- Correctamente aceite (ou *True Positive* – TP) - se a hipótese  $H_0$  foi aceite e estava certa;
- Correctamente rejeitado (ou *True Negative* – TN) - se a hipótese  $H_0$  não foi aceite e estava errada;
- Falsamente aceite (ou *False Positive* - FP) - se a hipótese  $H_0$  foi aceite e estava errada;
- Falsamente rejeitado (ou *False Negative* - FN) - se a hipótese  $H_0$  não foi aceite e estava certa.

		actual value		total
		$p$	$n$	
prediction outcome	$p'$	True Positive	False Positive	$P'$
	$n'$	False Negative	True Negative	$N'$
total		$P$	$N$	

Figura 5. 7 – Exemplo da matriz de confusão 2x2 que representa o corrente problema de classificação. [44]

Idealmente, o reconhecedor deveria produzir o maior número de correctamente aceites e correctamente rejeitados e um número mínimo de falsamente aceites e falsamente rejeitados, minimizando assim os erros de reconhecimento. Em termos práticos, a escolha do limiar de decisão não é mais que um compromisso que torna o reconhecedor mais ou menos flexível na correcção (benefícios), em detrimento do número de erros (custos).

Assim, são usadas diversas medidas, cujos valores podem ser inferidos directamente da matriz de confusão, das quais se destacam:

- *Sensitivity* ou *True Positive Rate* -  $TPR = TP / P = TP / (TP + FN)$ ;
- *False Positive Rate* -  $FPR = FP / N = FP / (FP + TN)$ ;
- *Accuracy* -  $ACC = (TP + TN) / (P + N)$ .

A escolha do limiar de decisão/aceitação deve ser tomada dependendo de qual a medida a atribuir maior relevância. A prioridade a atribuir à minimização dos resultados falsamente aceites (falsos positivos) ou à redução dos falsamente rejeitados (falsos negativos) dita a localização do limiar, visto serem medidas de crescimento inversamente proporcional.

Trata-se portanto, de uma decisão de estratégia de reforço psicológico positivo ou negativo, pesando se será mais profícuo detectar mais erros (correctamente rejeitados), implicando isso um maior número de casos rejeitados em que a criança de facto disse a palavra correcta (falsamente rejeitados), ou, adoptando a estratégia inversa, maximizar o número de correctas aceitaçãoes, acarretando um aumento no número de falsas aceitaçãoes.

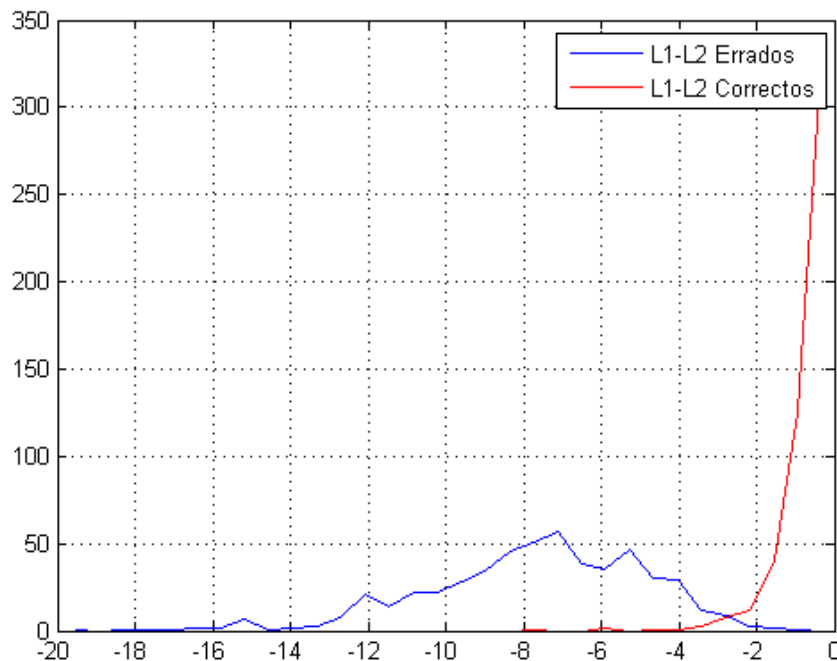
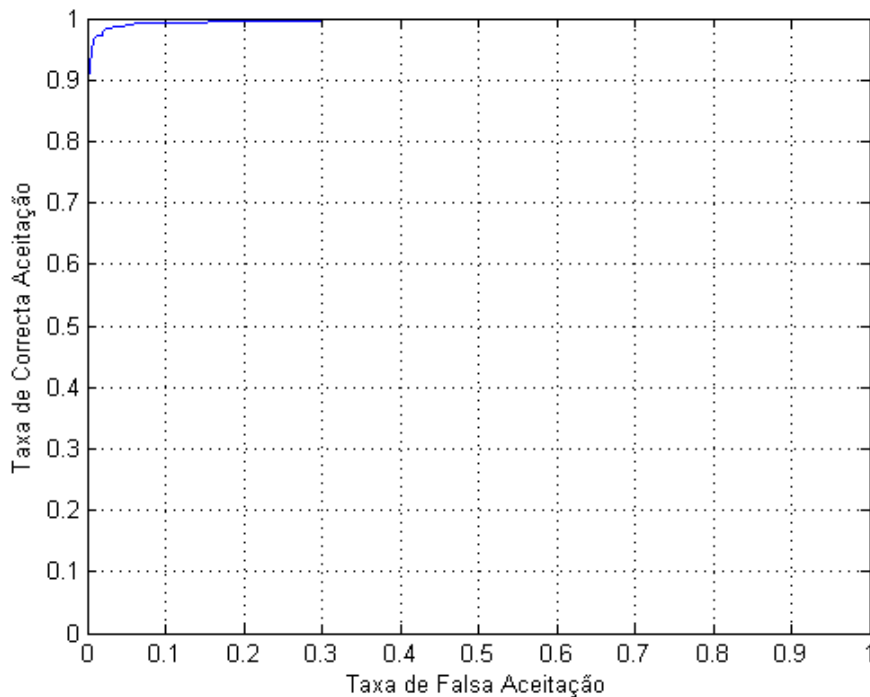


Figura 5.8 – Gráfico das curvas que compara a distância L1-L2 para as transcrições correctas e as erradas.

O traçado das duas curvas da diferença  $L_1-L_2$  para as palavras correctamente pronunciadas e para as palavras erradas permite observar as verosimilhanças para os dois casos, dando uma boa perspectiva da sua distribuição. O ponto de intersecção das duas curvas ROC é um bom indicador de onde se deverá localizar um bom limiar.



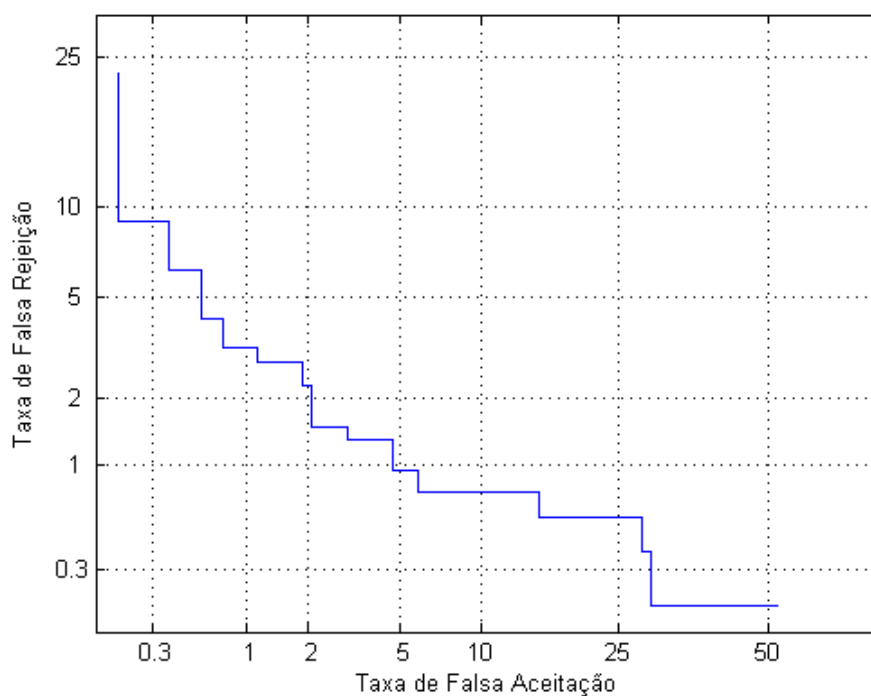
**Figura 5.9 – Curva ROC**

A curva ROC relaciona a Taxa de correctas aceitações com a taxa de falsas aceitações, sendo que o seu ponto ideal de operação se encontra nas coordenadas (0,1), onde todas as correctas aceitações seriam bem detectadas e não produzindo nenhum falso alarme (falsa aceitação).

No entanto, para melhor proceder à avaliação do problema, uma comparação do compromisso entre os dois tipos de erro, revela-se mais eficaz.

Para tal, recorreu-se também a gráficos DET (*Detection Error Tradeoff*) [45] cuja curva representa a relação da taxa de falsas rejeições em função da taxa de falsas aceitações (compromisso entre os diferentes tipos de erro). Este gráfico produz curvas mais lineares que as curvas ROC, conseguindo um melhor aproveitamento da área da imagem para a região crítica de operação.

Considerando ambas as taxas de erro com o mesmo peso, a solução ideal para a optimização do problema seria o ponto no gráfico DET com a menos distância à origem do referencial.



**Figura 5. 10 – Curva DET.**

Para o efeito deste estudo, e tendo em conta que é direccionado a crianças, será mais interessante minimizar o número de falsas rejeições (a criança disse correctamente a palavra mas o sistema reconheceu-a como erro), permitindo comprometer em certo ponto a taxa de falsas aceitações (erros da criança aceites pelo sistema como correctos). Assim, um limiar correspondente a 5% de falsa aceitação contra 0.9% de falsa rejeição ou 10% de falsa aceitação contra 0.8% de falsa rejeição poderão ser uma boa escolha.

Como foi já referido, o sistema não é excepcionalmente preciso ao nível do reconhecimento de fones individuais, podendo mesmo apresentar inserções ou substituições (em relação à sequência esperada da palavra) que podem ser ignoradas sempre que os seus valores de verosimilhança sejam relativamente baixos. Para tornar o sistema mais robusto é possível usar os valores ‘align\_score’ calculados como uma segunda medida de classificação.

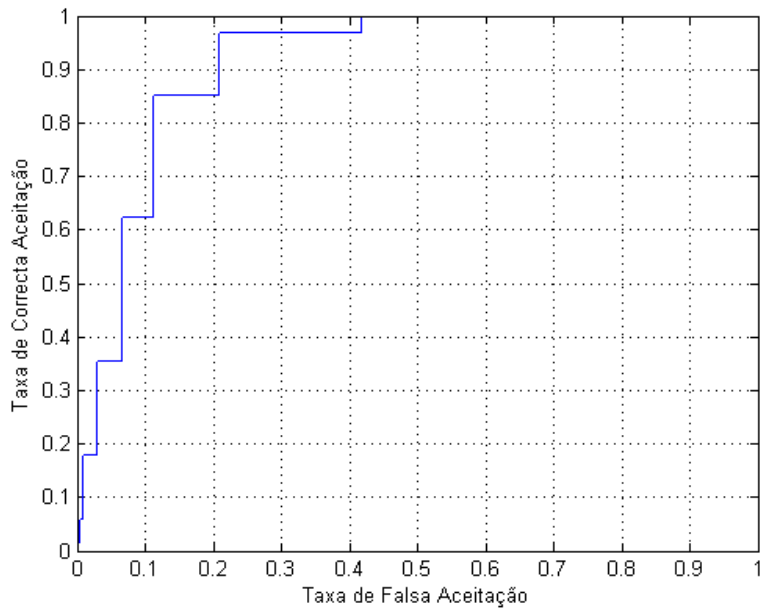


Figura 5. 11 – Curva ROC para o *score* de alinhamento.

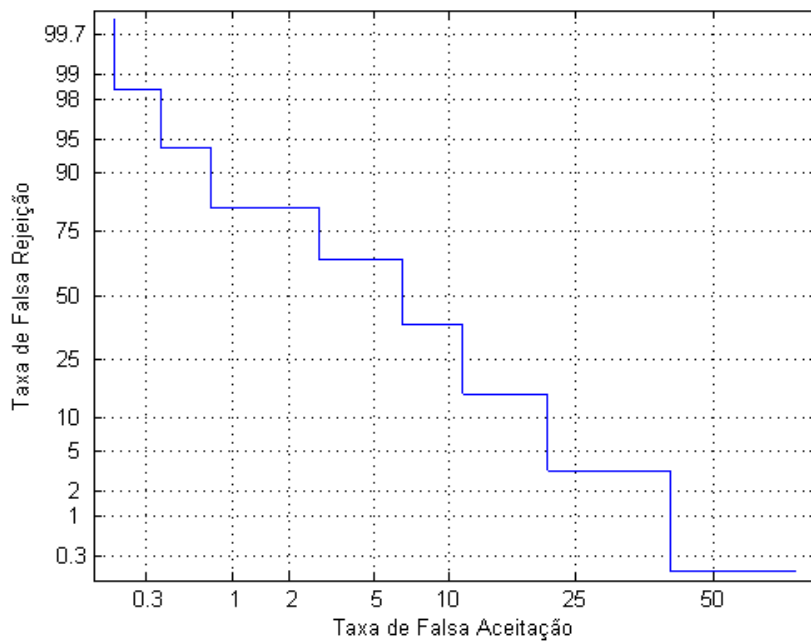


Figura 5. 12 – Curva DET para o *score* de alinhamento.

A Figura 5. 11 e a Figura 5. 12 representam a curva ROC e a curva DET, respectivamente, para os valores de *score* de alinhamento, à imagem do que foi feito com os valores de verosimilhança.

## 6. Conclusão e Trabalho Futuro

Este estudo tinha como objectivo principal a criação de um sistema de identificação de disfluências. Este objectivo foi implementado com sucesso. Foi também desenvolvido um sistema de avaliação da correcção da dicção das palavras proferidas em função dos valores de verosimilhança calculados com os modelos de fones. A decisão baseada nesta medida mostra um bom compromisso entre a taxa de falsas aceitações e falsas rejeições. A decisão pode ainda ser melhorada com base nos valores de alinhamento de *strings* de fones calculados com base no reconhecimento de fones sem restrições e com a transcrição fonética da palavra esperada em questão. Foi ainda aplicado um método de classificação binária e sugeridos valores para os limiares de decisão.

Com base no presente trabalho é agora possível desenvolver o sistema global de auxílio à terapia da fala em crianças, tal como foi descrito na secção 2.3. Com efeito, esta tese constitui o núcleo tecnológico base que torna possível e útil o referido sistema. Para tal dever-se-á aplicar os modelos e o método de verificação desenvolvidos, embebidos no sistema que usará jogos ou outras actividades lúdicas na interacção com a criança.

Seria ainda interessante melhorar os modelos de reconhecimento através da recolha de mais locuções para compor a base de dados. Uma maior amostra de locuções tornará o sistema ainda mais robusto. O próprio método de identificação das disfluências pode também ser melhorado, por exemplo, aplicando alguns conhecimentos de linguística sobre a evolução das capacidades de fala em crianças. O estudo [46] pode servir de fundamento para as referidas implementações.





## 7. Anexos

### 7.1. Anexo A – Ficheiro de Configuração do HCOPY

#### **config\_Hcopy.txt**

```
1 NATURALREADORDER = TRUE
2 NONUMESCAPES = TRUE
3 NATURALWRITEORDER = TRUE
4
5 # Coding parameters
6 SAVECOMPRESSED = F
7 SAVEWITHCRC = F
8 SOURCEFORMAT=WAV
9 TARGETKIND = MFCC_0_D_A
10 TARGETRATE = 100000.0
11 WINDOWSIZE = 150000.0
12 USEHAMMING = T
13 PREEMCOEF = 0.97
14 NUMCHANS = 14
15 NUMCEPS = 12
16 ENORMALISE = F
```

### 7.2. Anexo B - Lista de Fones e Dicionário

#### **Lista de fones com exemplo de uso**

<u>Fone</u>	<u>Exemplo de palavra</u>
a	<b>cama</b>
a~	mel <b>ancia</b>
@	microf <b>one</b>
a	p <b>ato</b>
b	<b>bola</b>
d	r <b>ádio</b>
E	car <b>ame</b> lo
e	m <b>esa</b>
e~	p <b>ente</b>
f	<b>faca</b>
g	<b>golo</b>
i	<b>igreja</b>
i~	<b>í</b> ndio
J	gol <b>finho</b>
k	<b>ca</b> sa
l	<b>l</b> ápis
L	ab <b>elha</b>
m	<b>m</b> ola
n	ca <b>ne</b> ca

O mota  
 o mosca  
 o~ bombeiro  
 p pato  
 r barco  
 R carro  
 s cenoura  
 S espada  
 t tomate  
 u boneco  
 u~ um  
 v velho  
 Z casa  
 z queijo

### Dicionário de transcrição SAMPA - diccorrect\_mul.txt

abacaxi	<sup>a</sup> b a k a k s i
abelha	<sup>a</sup> b <sup>a</sup> L <sup>a</sup>
abelha	<sup>a</sup> b e L <sup>a</sup>
abóbora	<sup>a</sup> b O b u r <sup>a</sup>
ananas	<sup>a</sup> n <sup>a</sup> n a S
anel	<sup>a</sup> n E l
aranha	<sup>a</sup> r <sup>a</sup> J <sup>a</sup>
arco-íris	a r k u i r i S
astronauta	<sup>a</sup> S t r O n a u t <sup>a</sup>
autocarro	a u t O k a R u
ave	a v @
avental	<sup>a</sup> v e~ t a l
avioneta	<sup>a</sup> v i u n E t <sup>a</sup>
avião	<sup>a</sup> v i <sup>a</sup> ~ u~
avô	<sup>a</sup> v o
babete	b a b E t @
bailarina	b a i l <sup>a</sup> r i n <sup>a</sup>
balde	b a l d @
baleia	b <sup>a</sup> l <sup>a</sup> i <sup>a</sup>
baliza	b <sup>a</sup> l i z <sup>a</sup>
baloiço	b <sup>a</sup> l o i s u
balão	b <sup>a</sup> l <sup>a</sup> ~ u~
banana	b <sup>a</sup> n <sup>a</sup> n <sup>a</sup>
bananas	b <sup>a</sup> n <sup>a</sup> n <sup>a</sup> S
bandeira	b <sup>a</sup> ~ d <sup>a</sup> i r <sup>a</sup>
banho	b <sup>a</sup> J u
baralho	b <sup>a</sup> r a L u
barco	b a r k u
barriga	b <sup>a</sup> R i g <sup>a</sup>
baton	b a t O n
bebé	b E b E
berbequim	b @ r b @ k i~
bicho	b i S u
bicicleta	b i s i k l E t <sup>a</sup>
bicicletas	b i s i k l E t <sup>a</sup> S
boca	b o k <sup>a</sup>
bola	b O l <sup>a</sup>

bolo	b o l u
bombeiro	b o~ b <sup>a</sup> i r u
bombom	b o~ b o~
bombons	b o~ b o~ S
boneca	b u n E k <sup>a</sup>
boneco	b u n E k u
borboleta	b u r b u l e t <sup>a</sup>
bota	b O t <sup>a</sup>
branco	b r <sup>a</sup> ~ k u
braço	b r a s u
brinquedo	b r i~ k e d u
bruxa	b r u S <sup>a</sup>
burro	b u R u
cabelo	k <sup>a</sup> b e l u
cabide	k <sup>a</sup> b i d @
cachorro	k <sup>a</sup> S o R u
cadeira	k <sup>a</sup> d <sup>a</sup> i r <sup>a</sup>
caixa	k a i S <sup>a</sup>
caixote	k a i S O t @
calça	k a l s <sup>a</sup>
calças	k a l s <sup>a</sup> S
calções	k a l s o~ i~ S
cama	k <sup>a</sup> m <sup>a</sup>
camelo	k <sup>a</sup> m e l u
camioneta	k a m i u n E t <sup>a</sup>
camisola	k <sup>a</sup> m i z O l <sup>a</sup>
camião	k a m i <sup>a</sup> ~ u~
candeeiro	k <sup>a</sup> ~ d i <sup>a</sup> i r u
caneca	k <sup>a</sup> n E k <sup>a</sup>
caneta	k <sup>a</sup> n e t <sup>a</sup>
canguru	k <sup>a</sup> ~ g u r u
cangurus	k <sup>a</sup> ~ g u r u S
capacete	k <sup>a</sup> p <sup>a</sup> s e t @
caracol	k <sup>a</sup> r <sup>a</sup> k O l
caramelo	k <sup>a</sup> r <sup>a</sup> m e l u
caranguejo	k <sup>a</sup> r <sup>a</sup> ~ g <sup>a</sup> Z u
caranguejo	k <sup>a</sup> r <sup>a</sup> ~ g <sup>a</sup> i Z u
carrinha	k <sup>a</sup> R i J <sup>a</sup>
carro	k a R u
carrocel	k <sup>a</sup> R O s E l
carta	k a r t <sup>a</sup>
cartas	k a r t <sup>a</sup> S
casa	k a z <sup>a</sup>
cavalo	k <sup>a</sup> v a l u
caçador	k <sup>a</sup> s <sup>a</sup> d o r
cebola	s @ b o l <sup>a</sup>
cenoura	s @ n o r <sup>a</sup>
cereais	s @ r i a i S
cereja	s @ r <sup>a</sup> Z <sup>a</sup>
cereja	s @ r <sup>a</sup> i Z <sup>a</sup>
cerejas	s @ r <sup>a</sup> Z <sup>a</sup> S
cerejas	s @ r <sup>a</sup> i Z <sup>a</sup> S
cerveja	s @ r v <sup>a</sup> Z <sup>a</sup>
cerveja	s @ r v <sup>a</sup> i Z <sup>a</sup>
cesta	s e S t <sup>a</sup>
cesto	s e S t u
chapéu	S <sup>a</sup> p E u

chave	S a v @
chavina	S ^ a v i n ^ a
chinelo	S i n E l u
chinelos	S i n E l u S
chocolate	S u k u l a t @
chupeta	S u p E t ^ a
chá	S a
cisne	s i Z n @
cobra	k O b r ^ a
coelho	k u ^ a L u
coelho	k u e L u
cofre	k O f r @
cogumelo	k u g u m E l u
colheres	k u L E r @ S
comboio	k o ~ b O i u
computador	k o ~ p u t ^ a d o r
copo	k O p u
coração	k u r ^ a s ^ a ~ u ~
corneta	k O r n e t ^ a
coroa	k u r o ^ a
cortina	k u r t i n ^ a
coruja	k O r u Z ^ a
cruzeta	k r u z e t ^ a
cubo	k u b u
cueca	k u E k ^ a
cuecas	k u E k ^ a S
cão	k ^ a ~ u ~
dado	d a d u
dedo	d e d u
dinossauro	d i n O s a u r u
direito	d i r ^ a i t u
doces	d o s @ S
doninha	d O n i J ^ a
dormir	d u r m i r
doutor	d o t o r
duas	d u ^ a S
elefante	e l @ f ^ a ~ t @
escada	@ S k a d ^ a
escova	@ S k o v ^ a
espada	@ S p a d ^ a
esquerda	@ S k e r d ^ a
esquilo	@ S k i l u
estrela	@ S t r e l ^ a
estrelinha	@ S t r @ l i J ^ a
faca	f a k ^ a
fada	f a d ^ a
fato	f a t u
ferro	f E R u
figo	f i g u
flauta	f l a u t ^ a
flor	f l o r
foca	f O k ^ a
foguete	f u g @ t ^ a ~ u ~
fogão	f u g ^ a ~ u ~
folha	f o L ^ a
formiga	f u r m i g ^ a
forno	f o r n u

fralda	f r a l d a
fruta	f r u t a
frutos	f r u t u S
galinha	g a l i J a
garfo	g a r f u
garfos	g a r f u S
garrafa	g a R a f a
gato	g a t u
gelado	Z @ l a d u
gelo	Z e l u
girafa	Z i r a f a
golfinho	g o l f i J u
golfinhos	g o l f i J u S
gommas	g o m a S
gravata	g r a v a t a
guitarra	g i t a R a
helicóptero	e l i k O p t @ r u
hipopótamo	i p O p O t a m u
homem	O m a ~ i ~
igreja	i g r a Z a
igreja	i g r a i Z a
iogurte	i O g u r t @
janela	Z a n E l a
jarra	Z a R a
jarro	Z a R u
joaninha	Z u a n i J a
lanterna	l a ~ t E r n a
laranja	l a r a ~ Z a
lata	l a t a
lego	l e g u
legos	l e g u S
legumes	l @ g u m @ S
leão	l i a ~ u ~
limão	l i m a ~ u ~
linha	l i J a
livro	l i v r u
livros	l i v r u S
lixo	l i S u
lobo	l o b u
lua	l u a
luva	l u v a
luz	l u S
lápiz	l a p i S
lâmpada	l a ~ p a d a
macaco	m a k a k u
mala	m a l a
mar	m a r
martelo	m a r t E l u
maçã	m a s a ~
melancia	m @ l a ~ s i a
melão	m @ l a ~ u ~
menino	m @ n i n u
mesa	m e z a
microfone	m i k r O f O n @
milho	m i L u
mochila	m u S i l a
mocho	m O S u

moinho	m u i J u
mola	m O l <sup>a</sup>
morango	m u r <sup>a~</sup> g u
mosca	m o S k <sup>a</sup>
mosquito	m u S k i t u
mota	m O t <sup>a</sup>
mão	m <sup>a~</sup> u~
médico	m E d i k u
música	m u z i k <sup>a</sup>
ninho	n i J u
nuvem	n u v <sup>a~</sup> i~
nuvens	n u v <sup>a~</sup> i~ S
olho	o L u
ovelha	O v <sup>a</sup> L <sup>a</sup>
ovelha	O v e L <sup>a</sup>
ovo	o v u
ovos	O v u S
palhaço	p <sup>a</sup> L a s u
panela	p <sup>a</sup> n E l <sup>a</sup>
papagaio	p <sup>a</sup> p <sup>a</sup> g a i u
pássaro	p <sup>a</sup> s <sup>a</sup> r u
passarinho	p <sup>a</sup> s <sup>a</sup> r i J u
patim	p <sup>a</sup> t i~
patins	p <sup>a</sup> t i~ S
pato	p a t u
pedra	p E d r <sup>a</sup>
peixe	p <sup>a</sup> i S @
pendurar	p e~ d u r a r
pente	p e~ t @
penteado	p e~ t i a d u
pequena	p @ k e n <sup>a</sup>
piano	p i <sup>a</sup> n u
piscina	p i S s i n <sup>a</sup>
pilha	p i L <sup>a</sup>
pipocas	p i p O k <sup>a</sup> S
pirata	p i r a t <sup>a</sup>
piratas	p i r a t <sup>a</sup> S
pizza	p i z <sup>a</sup>
pião	p i <sup>a~</sup> u~
polícia	p u l i s i <sup>a</sup>
pomba	p o~ b <sup>a</sup>
porco	p o r k u
porta	p O r t <sup>a</sup>
portugal	p u r t u g a l
praia	p r a i <sup>a</sup>
pranchas	p r <sup>a~</sup> S <sup>a</sup> S
prato	p r a t u
princesa	p r i~ s e z <sup>a</sup>
pulseira	p u l s <sup>a</sup> i r <sup>a</sup>
páscoa	p a S k u <sup>a</sup>
pássaro	p a s <sup>a</sup> r u
pão	p <sup>a~</sup> u~
pé	p E
pés	p E S
pêra	p e r <sup>a</sup>
pêssego	p e s @ g u
pêssegos	p e s @ g u S

pónei	p O n <sup>a</sup> i
quadrado	k u <sup>a</sup> d r a d u
quadro	k u a d r u
queijo	k <sup>a</sup> i Z u
quente	k e ~ t @
quivi	k i v i
rainha	R <sup>a</sup> i J <sup>a</sup>
raposa	R <sup>a</sup> p o z <sup>a</sup>
raquete	R <sup>a</sup> k E t @
rato	R a t u
recta	R E k t <sup>a</sup>
rectângulo	R E k t <sup>a</sup> ~ g u l u
rede	R e d @
regador	R @ g <sup>a</sup> d o r
rei	R <sup>a</sup> i
relógio	R @ l O Z i u
rena	R e n <sup>a</sup>
renas	R e n <sup>a</sup> S
retrato	R @ t r a t u
rinoceronte	R i n O s @ r o ~ t @
robô	R o b o
rocha	R O S <sup>a</sup>
rolha	R o L <sup>a</sup>
roupa	R o p <sup>a</sup>
rádio	R a d i u
rã	R <sup>a</sup> ~
régua	R E g u <sup>a</sup>
sanita	s <sup>a</sup> n i t <sup>a</sup>
sapato	s <sup>a</sup> p a t u
sapo	s a p u
secador	s @ k <sup>a</sup> d o r
senhor	s @ J o r
sentada	s e ~ t a d <sup>a</sup>
sereia	s @ r <sup>a</sup> i <sup>a</sup>
seta	s E t <sup>a</sup>
símbolo	s i ~ b o l u
sinal	s i n a l
sino	s i n u
sofá	s u f a
sol	s O l
sumo	s u m u
tacho	t a S u
talher	t <sup>a</sup> L E r
talheres	t <sup>a</sup> L E r @ S
tambor	t <sup>a</sup> ~ b o r
tampa	t <sup>a</sup> ~ p <sup>a</sup>
tangerina	t <sup>a</sup> ~ Z @ r i n <sup>a</sup>
tartaruga	t <sup>a</sup> r t <sup>a</sup> r u g <sup>a</sup>
telefone	t @ l @ f O n @
telemóvel	t E l E m O v E l
televisão	t @ l @ v i z <sup>a</sup> ~ u ~
tesoura	t @ z o r <sup>a</sup>
tigre	t i g r @
tocar	t u k a r
tomate	t u m a t @
triângulo	t r i <sup>a</sup> ~ g u l u
unha	u J <sup>a</sup>

urso	u r s u
uva	u v a
uvas	u v a S
vaca	v a k a
varanda	v a r a~ d a
vaso	v a z u
vassoura	v a s o r a
veado	v i a d u
vela	v E l a
velhinho	v E L i J u
velho	v E L u
velhote	v E L O t @
vestido	v @ S t i d u
vinho	v i J u
viola	v i O l a
visão	v i z a~ u~
xilofone	S i l O f O n @
zebra	z E b r a
água	a g u a
áriel	a r i E l
árvore	a r v u r @
índio	i~ d i u
íris	i r i S
óculos	O k u l u S

### 7.3. Anexo C - Scripts em DOS/BATCH de Treino dos Modelos

#### treino.bat

```

1  REM cd para directoria de treino
2
3  HLEd -l * -C configbasico.txt -d dic_newword.txt -i mlf\cut\mono_phones.mlf mkphones1.led
4  mlf\cut\MLF_cut2.mlf
5
6  REM set cdir=modelos\hmm0
7  REM - criar proto manualmente
8  REM list_train.txt criada com mfclist.m
9  REM - HCompV cria vFloors e actualiza proto
10 HCompV -C configbasico.txt -f 0.01 -S mlf/cut/MLF_cut2_scp_train.scp -m -M modelos/hmm0 proto_
11 REM - cuidado: o nome em ~h (neste caso "proto_") tem ser o mesmo do ficheiro (HCompV cracha)
12
13
14 REM criacao do MLF de treino com phonemlf.m
15 REM NOTA: os ficheiros MFC foram cortados para não terem <sil>,<noise>, etc. Criado MLF em
16 conformidade mantendo a ref do tempo inicial de cada label.
17 REM ver ficheiro: cutmlf.m
18
19
20 REM - fazer modelos\hmm0\macros com opcoes e vFloors
21 echo macros - vFloors %i - !date! !time!
22 echo.
23 echo ~o >modelos\hmm0\macros
24 echo ^<STREAMINFO^> 1 39 >>modelos\hmm0\macros
25 echo ^<VECSIZE^> 39^<NULLD^>^<MFCC_0_D_A^>^<DIAGC^> >>modelos\hmm0\macros
26 type modelos\hmm0\vFloors >>modelos\hmm0\macros
27
28
29 REM - fazer hmmdefs = phones_GV.mmf -> clonagem dos modelos: ver protoscpcclone.m
30 REM palavras

```



```

31 HERest -C configbasico.txt -T 7 -I mlf\cut\monophone_notime.mlf -t 250.0 150.0 1000.0 -S
32 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm0\macros -H modelos\hmm0\phones_GV.mmf -M
33 modelos\hmm1_phone_list_GV.txt
34 REM set cdir=modelos\hmm1
35 HERest -C configbasico.txt -T 7 -I mlf\cut\monophone_notime.mlf -t 250.0 150.0 1000.0 -S
36 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm1\macros -H modelos\hmm1\phones_GV.mmf -M
37 modelos\hmm2_phone_list_GV.txt
38 REM set cdir=modelos\hmm2
39 HERest -C configbasico.txt -T 7 -I mlf\cut\monophone_notime.mlf -t 250.0 150.0 1000.0 -S
40 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm2\macros -H modelos\hmm2\phones_GV.mmf -M
41 modelos\hmm3_phone_list_GV.txt
42 REM set cdir=modelos\hmm3
43 HERest -C configbasico.txt -T 7 -I mlf\cut\monophone_notime.mlf -t 250.0 150.0 1000.0 -S
44 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm3\macros -H modelos\hmm3\phones_GV.mmf -M
45 modelos\hmm4_phone_list_GV.txt
46 REM set cdir=modelos\hmm4
47 HERest -C configbasico.txt -T 7 -I mlf\cut\monophone_notime.mlf -t 250.0 150.0 1000.0 -S
48 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm4\macros -H modelos\hmm4\phones_GV.mmf -M
49 modelos\hmm5_phone_list_GV.txt
50 REM set cdir=modelos\hmm5
51 HERest -C configbasico.txt -T 7 -I mlf\cut\monophone_notime.mlf -t 250.0 150.0 1000.0 -S
52 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm2\macros -H modelos\hmm5\phones_GV.mmf -M
53 modelos\hmm6_phone_list_GV.txt
54 REM set cdir=modelos\hmm6
55 HERest -C configbasico.txt -T 7 -I mlf\cut\monophone_notime.mlf -t 250.0 150.0 1000.0 -S
56 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm2\macros -H modelos\hmm6\phones_GV.mmf -M
57 modelos\hmm7_phone_list_GV.txt
58
59 REM - Realinhar phones
60 HVite -l * -C configbasico.txt -o SW -H modelos\hmm7\macros -H modelos\hmm7\phones_GV.mmf -i
61 mlf\cut\aligned_phone.mlf -m -t 250.0 150.0 1000.0 -y lab -a -I mlf\cut\MLF_cut2.mlf -S
62 mlf\cut\MLF_cut2_scp_train.scp dic_newword.txt phone_list_GV.txt > logs\HVite_align_phone.log
63
64 REM - mais 2x HERest com alinhado
65 HERest -C configbasico.txt -T 7 -I mlf\cut\aligned_phone.mlf -t 250.0 150.0 1000.0 -S
66 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm7\macros -H modelos\hmm7\phones_GV.mmf -M
67 modelos\hmm8_phone_list_GV.txt
68 HERest -C configbasico.txt -T 7 -I mlf\cut\aligned_phone.mlf -t 250.0 150.0 1000.0 -S
69 mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm8\macros -H modelos\hmm8\phones_GV.mmf -M
70 modelos\hmm9_phone_list_GV.txt

```

---

## gmixture.bat

```

1 @ECHO OFF
2
3 set init=hmm9
4 for /L %%i in (2,1,16) do (
5
6     md modelos\hmm_mix%%i_0
7     md modelos\hmm_mix%%i_1
8     md modelos\hmm_mix%%i_2
9     md modelos\hmm_mix%%i_3
10
11     REM apenas uma questão de um nome mais explicito:
12     REM copy phone_list3.txt phone_list_notag.txt
13
14     REM splitting de uma para 2 componentes Gaussianas:
15     HHed -H modelos\%init%\phones_GV.mmf -H modelos\%init%\macros -M modelos\hmm_mix%%i_0
16     hed\mix%%i.hed phone_list_GV > logs\logHHed_mix%%i.txt
17
18     REM 3 passos de reestimação embebida:
19     HERest -C configbasico.txt -T 7 -I mlf\cut\aligned_phone16.mlf -t 250.0 150.0 1000.0 -S
20     mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm_mix%%i_0\phones_GV.mmf -H
21     modelos\hmm_mix%%i_0\macros -M modelos\hmm_mix%%i_1 phone_list_GV.txt
22
23     HERest -C configbasico.txt -T 7 -I mlf\cut\aligned_phone16.mlf -t 250.0 150.0 1000.0 -S
24     mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm_mix%%i_1\phones_GV.mmf -H
25     modelos\hmm_mix%%i_1\macros -M modelos\hmm_mix%%i_2 phone_list_GV.txt
26
27     HERest -C configbasico.txt -T 7 -I mlf\cut\aligned_phone16.mlf -t 250.0 150.0 1000.0 -S
28     mlf\cut\MLF_cut2_scp_train.scp -H modelos\hmm_mix%%i_2\phones_GV.mmf -H
29     modelos\hmm_mix%%i_2\macros -M modelos\hmm_mix%%i_3 phone_list_GV.txt
30

```

```

31         set init = hmm_mix%i_3
32     )
33
34     @ECHO ON

```

---

## teste.bat

```

1  REM cd para directoria de teste
2
3  REM creates de grammar network from the file taskgram generated with writehtkGrammar.m
4  HParse -C configbasico.txt ..\gram\gram ..\gram\net
5
6  HVite -T 01 -C configbasico.txt -H ..\test\modelos\hmm_mix16_3\macros -
7  H ..\test\modelos\hmm_mix16_3\phones_GV.mmf -S mlf\cut\MLF_cut2_scp_test.scp -l * -i recout.mlf -
8  w ..\gram\net.txt dic_lix.txt phone_list_GV.txt > log\HVite_rec.log
9
10 HResults -T 07 -e ??? lix -t -I MLF_cut2_lix.mlf phone_list_GV.txt recout.mlf > log\HRes.log

```

## 7.4. Anexo D – Scripts e Funções do Matlab

### filtra\_db.m

```

1  function filtra_bd(listfile,dir_orig,dir_dest)
2  %
3  %Filtra BD passa-alto com freq. de corte 100Hz
4  %
5  %filtabd('list.txt','wavs','wav_hp')
6
7  if dir_orig(end)~='\', dir_orig=[dir_orig,'\']; end
8  if dir_dest(end)~='\', dir_dest=[dir_dest,'\']; end
9
10 list = textread(listfile,'%s');
11 N=length(list);
12 fs=8000;
13 f0=100;
14 wcn = (2*pi*f0/fs)/pi;
15 b = fir1(256,wcn,'high');
16 [H,w]=freqz(b,1);
17 plot(w/pi*fs/2,db(H)), grid
18
19 for n=1:N
20     x=wavread([dir_orig,list{n}]);
21
22     y = conv(x,b);
23     maxy=max(abs(y));
24     if maxy>0.99
25         y = y./(maxy+0.0001);
26     end
27     y(1:128)=[];
28     y(end-127:end)=[];
29
30     wavwrite(y,fs,16,[dir_dest,list{n}]);
31 end

```

---

## cleannoisemlf.m

```
1 function cleannoisemlf(in_mlf,out_mlf)
2 %function cleannoisemlf(in_mlf,out_mlf)
3 %
4 %Takes an mlf file and substitutes the labels that end
5 % with $ and # by <noise>.
6
7 [pats,labs] = mlfRead2(in_mlf,'%f%f%s',char(9)); % del = char(9)
8 pats2 = pats;
9 labs2 = labs;
10 N = length(pats);
11
12 for n=1:N,
13     for m=1:length(labs(n).lab)
14         s = labs(n).lab{m};
15         if s(end) == '#' || s(end) == '$'
16             labs2(n).lab{m} = '<noise>';
17         end
18     end
19 end
20
21 mlfWrite2(pats2,labs2,out_mlf);
22
23 end
```

---

## emptynoisemlf.m

```
1 function emptynoisemlf(mlf_file, out_mlf,delimiter)
2 %function emptynoisemlf(mlf_file, out_mlf,delimiter)
3 %Verifies if any of the pats of the mlf_file no spoken labs
4 % and only tags: <noise>, <sil>, <resp>, ...
5
6 fe = fopen(out_mlf,'wb');
7 empty = '';
8 if nargin<3, delimiter=char(9); end
9 [pats,labs] = mlfRead2(mlf_file,'%f%f%s',delimiter);
10
11 N = length(pats);
12 numpats = 0;
13
14 f = fopen(out_mlf, 'wb');
15 fprintf(f,'#!MLF!\n');
16
17 for n=1 : N %corre os ficheiros
18     cont = 0;
19     M = length(labs(n).lab);
20     for m=1:M
21         s = labs(n).lab{m};
22         if s(1) ~= '<' % spoken
23             cont = cont + 1;
24         end
25     end
26     if cont ~= 0
27         fprintf(f, "%s\n",pats{n});
28         for m=1:M
29             fprintf(f, '%d\t%d\t%s\n',labs(n).ti(m),labs(n).tf(m),labs(n).lab{m});
30         end
31         fprintf(f, '\n');
32         numpats = numpats + 1;
33     end
34
35 end
36
37 numpats
38 fclose(f);
```

**formatmlf.m**

```

1 function formatmlf(mlf_file,dir_dest,delimiter)
2 %function formatmlf(listfile,mlf_file,dir_dest)
3 % limpa as anotações de wav sem nada (lixo)
4 %cria lista de palavras e conta ocorrências
5 %cria lista de tags
6
7 MLFfile = mlf_file;
8 %list = textread(listfile,'%s');
9 if nargin<2, dir_dest=pwd; end
10 if nargin<3, delimiter=char(9); end
11 [pats,labs] = mlfRead2(MLFfile,'%f%f%s',delimiter);
12
13 N = length(pats)
14 words = ''; %locações
15 tags = ''; %tags
16 flixo = ''; %pats de lixo
17 ilixo = []; %índices de lixo
18 nwords = 0; %índice locuções
19 ntags = 0; %índice tags
20 tagc = []; %num tags
21 wordc = []; %num locuções
22 fpats = ''; %pats
23 npats = ''; %pats noise
24
25 pl = strcat(dir_dest,'list_clean.txt');
26 pw = strcat(dir_dest,'new_word.txt');
27 pt = strcat(dir_dest,'new_unique_tags.txt');
28 pg = strcat(dir_dest,'new_lixo.txt');
29 pn = strcat(dir_dest,'new_noise.txt');
30 fl = fopen(pl,'wb');
31 fw = fopen(pw,'wb');
32 ft = fopen(pt,'wb');
33 fg = fopen(pg,'wb');
34 fn = fopen(pn,'wb');
35
36 cont = 1;
37 for n=1 : N %corre os ficheiros
38     M = length(labs(n).lab);
39     nome = pats{n};
40     name = horzcat(nome(3:(end-4)), '.wav');
41     fpats = strvcats(fpats, name);
42     for m=1 : M %corre as anotações para cada ficheiro
43         s = labs(n).lab{m}; %guarda o que está escrito
44
45         if (M==1) && (strcmp(s,'lixo')) % não tem nada
46             flixo = strvcats(flixo, nome); % guarda as ocorrências só com 'lixo'
47             ilixo(cont) = n;
48             cont = cont + 1;
49         else %guarda esta locução!
50
51
52             if s(1) == '<' % é uma tag
53                 i = strmatch(s, tags,'exact');%verifica se já existe a tag
54                 if isempty(i) %se não existir : isempty=true
55                     tags = strvcats(tags, s); %acrescenta à lista de tags
56                     ntags = ntags + 1; %total tags
57                     tagc(ntags) = 1;
58                 else %já existe a tag : isempty=false
59                     tagc(i) = tagc(i)+1; %incrementa ocorrências
60                 end
61             else % é uma palavra
62                 if s(end) == '$' % é para desprezar
63                     labs(n).lab{m} = '<noise>';
64                     fprintf(fn, '%s\t%s\n', s, nome);
65                 else
66                     i = strmatch(s,words,'exact');%verifica se já existe a palavra
67                     if isempty(i) %se não existir : isempty=true
68                         fprintf(fw, '%s\n', s);
69                         words = strvcats(words, s);

```

```

70         nwords = nwords +1; % total words
71         wordc(nwords) = 1;
72     else
73         wordc(i) = wordc(i)+1;
74     end
75 end
76 end
77 end
78 end
79 end
80 pats(ilixo) = []; %elimina locuções só com lixo
81 labs(ilixo) = [];
82 new_mlf = strcat(dir_dest,'clean_mlfd.mlf');
83 mlfWrite2(pats, labs, new_mlf);
84 tags = sortrows(tags);
85 flixo = sortrows(flixo);
86
87 fclose(fw);
88 pw2 = strrep(pw, '_word', 'word');
89 f = fopen(pw2, 'wb');
90 wlist = textread(pw, '%s', 'delimiter', char(9));
91 wlist = sortrows(lower(wlist));
92 wlist = unique(wlist);
93 for n = 1 : length(wlist)
94     fprintf(f, '%s\n', wlist{n});
95 end
96 fclose(f);
97
98 for row = 1:size(pats,1)
99     fprintf(fl, '%s\n', fpats(row,1:end));
100 end
101 for row = 1:size(tags,1)
102     fprintf(ft, '%s\n', tags(row,1:end));
103 end
104 for row = 1:size(flixo,1)
105     fprintf(fg, '%s\n', flixo(row,1:end));
106 end
107
108 fclose(fn);
109 fclose(fl);
110 fclose(ft);
111 fclose(fg);

```

---

## cutmlf.m

```

1 function cutmlf (mlf_file,mfc_dir,dest_dir)
2 % Cuts mlf files to get only the relevant info (cuts noise/breath/...)
3 % Separates mlf file into 2 files
4 % - mlf with only the words
5 % - mlf with remainder tags
6 % Uses two separate directories for the utterances and the remainder:
7 % - .\cut;
8 % - .\rest;
9 %
10 % the new /mlf files has an additional fourth parameter with
11 % time information of the labels that refers to their original initial time
12 % in the original mlf file
13 % EXAMPLE:
14 % ti tf lab #original_ti
15 % 12 17 <sil> #40
16
17 %% cutmlf('clean_mlfd_v2.mlf','..\mfc')
18
19 MLFfile = mlf_file; %'clean_mlfd_v2.mlf'
20 [pats,labs] = mlfRead2(MLFfile,'%f%f%s',char(9)); % del = char(9)
21 N = length(pats);
22
23 mkdir(dest_dir,'cut');
24 mkdir(dest_dir,'rest');
25 MLF_cut = strcat(dest_dir,'\cut\MLF_cut.mlf');
26 MLF_cut2 = strcat(dest_dir,'\cut\MLF_cut2.mlf');
27 MLF_rest = strcat(dest_dir,'\rest\MLF_rest.mlf');
28 MLF_rest2 = strcat(dest_dir,'\rest\MLF_rest2.mlf');

```

```

29
30 fec = fopen('MLF_cut_empty.txt', 'wb');
31 fer = fopen('MLF_rest_empty.txt', 'wb');
32 fc = fopen(MLF_cut, 'wb');
33 fc2 = fopen(MLF_cut2, 'wb');
34 fr = fopen(MLF_rest, 'wb');
35 fr2 = fopen(MLF_rest2, 'wb');
36
37 fprintf(fc, '#!MLF!\n');
38 fprintf(fc2, '#!MLF!\n');
39 fprintf(fr, '#!MLF!\n');
40 fprintf(fr2, '#!MLF!\n');
41
42 for n=1 : N
43     %disp(n);
44     M = length(labs(n).lab);
45     %nome = strrep(pats{n}, '*', mfc_dir);    % ex: mfc_dir = '..\mfc'
46     nome = strrep(pats{n}, '/', '\');
47     nome(end-3:end) = '.mfc';
48     nome_o = strrep(nome, '*', mfc_dir);
49
50     nome_c = strrep(nome, '*', [mfc_dir, '\cut_clean']);
51     nome_r = strrep(nome, '*', [mfc_dir, '\rest_clean']);
52
53     [X, nSamples, sampPeriod, sampSize, parmKind]=readHTK(nome_o, 0);
54     X_r=[];
55     X_c=[];
56     %f_mfc_c = fopen(nome_c, 'wb');
57     %f_mfc_r = fopen(nome_r, 'wb');
58
59     fprintf(fc, "%s\n", pats{n}); %ESTÁ AQUI '/' A MAIS!!!
60     fprintf(fc2, "%s\n", pats{n});
61     fprintf(fr, "%s\n", pats{n}); %ESTÁ AQUI '/' A MAIS!!!
62     fprintf(fr2, "%s\n", pats{n});
63
64     ti_r = 0;
65     ti_c = 0;
66
67     nSamples_c = 0;
68     nSamples_r = 0;
69     mt = 0;
70     mc = 0;
71     for m=1 : M    %corre as anotações para cada ficheiro
72         lab = labs(n).lab(m);    %gets the comment
73         ti = labs(n).ti(m);
74         tf = labs(n).tf(m);
75         d = tf - ti;    %duration
76
77         %get correspondent sample number
78         fi = floor(ti / sampPeriod) + 1;
79         ff = floor(tf / sampPeriod);
80         if ff>size(X, 2)
81             fprintf('Tempo final errado em %s\n', pats{n});
82             ff=size(X, 2);
83         end
84
85         if lab(1) == '<' % it's a tag -> MLF_rest
86             mt = mt + 1;
87             fprintf(fr, '%.0f\t%.0f\t%s\t%.0f\n', ti_r, ti_r + d, lab, ti);
88             fprintf(fr2, '%.0f\t%.0f\t%s\n', ti_r, ti_r + d, lab);
89             X_r = [X_r, X(:, fi:ff)];
90             ti_r = ti_r + d;
91             nSamples_r = nSamples_r + ff-fi + 1;
92         else % it's an utterance -> MLF_cut
93             mc = mc + 1;
94             fprintf(fc, '%.0f\t%.0f\t%s\t%.0f\n', ti_c, ti_c + d, lab, ti);
95             fprintf(fc2, '%.0f\t%.0f\t%s\n', ti_c, ti_c + d, lab);
96             X_c = [X_c, X(:, fi:ff)];
97             ti_c = ti_c + d;
98             nSamples_c = nSamples_c + ff-fi + 1;
99         end
100     end
101     if mt == 0
102         fprintf(fer, '%s\n', nome);
103     end
104     if mc == 0
105         fprintf(fec, '%s\n', nome);
106     end
107

```

```

108     writeHTK(nome_c,X_c,nSamples_c,sampPeriod,sampSize,parmKind, 0);
109     writeHTK(nome_r,X_r,nSamples_r,sampPeriod,sampSize,parmKind, 0);
110
111     fprintf(fc, '\n');
112     fprintf(fr, '\n');
113     fprintf(fc2, '\n');
114     fprintf(fr2, '\n');
115 end
116
117 fclose(fc);
118 fclose(fr);
119 fclose(fc2);
120 fclose(fr2);
121 fclose(fec);
122 fclose(fer);

```

---

### mfclist.m

```

1  function mfclist(mlf_file)
2  %function mfclist(mlf_file)
3  %
4  %Given the directory of the .mlf file, generates 2 .scp files
5  % from the list of all pats in a mlf file.
6  %
7  %It generates a training list with 4/5 of the files and a
8  % test list with the remainder 1/5
9
10
11 nametrain = strrep(mlf_file,'.mlf','_scp_train.scp');
12 nametest = strrep(mlf_file,'.mlf','_scp_test.scp');
13
14 [pats,labs]=mlfRead2(mlf_file,'%f%f%s');
15 %files = dir (mfc_dir);
16 nfiles = '';
17
18 ftrain = fopen(nametrain,'wb');
19 fttest = fopen(namettest,'wb');
20 %N = length (files);
21 N = length (pats);
22 Ntrain = 4*round (N/5);
23 Ntest = N - Ntrain;
24
25 %list_train = strrep(fntransc, '.lab', '.wav');
26
27 % for n=3:N % start at n=3 to ignore .\ and ..\
28 % scp_file = strrep(files(n).name,'S','..\mfc\cut\S');
29 % fprintf(fmfc,'%s\n',scp_file);
30 % end
31
32
33 rand('seed',0);
34 i=randperm(N);
35 j=i(1:N);
36 for k=1:N
37     scp_file = pats{j(k)};
38     scp_file = strrep(scp_file,'*/S','../mfc/cut_clean/S');
39     scp_file = strrep(scp_file, '.lab', '.mfc');
40     if k < Ntrain
41         fprintf(ftrain, '%s\n', scp_file);
42     else
43         fprintf(fttest, '%s\n', scp_file);
44     end
45 end;
46
47 fclose(ftrain);
48 fclose(fttest);

```

---

## protoscpcclone.m

```
1 function protoscpcclone(in_proto,out_mmffile,dic)
2 %Clonagem
3 %ler proto em modelos\hmm0
4 %Escrever MMF com todos os modelos
5 [lab mod] = textread(dic,'%s%s','delimiter',char(9));
6
7 plist='';
8 N=length(mod);
9 for n=1:N
10     c=textscan(mod{n},' %s');
11     for m=1:length(c{1})
12         cs=c{1}(m);
13         i = strmatch(cs{1},plist,'exact');
14         if isempty(i)
15             plist=[plist,cs];
16         end
17     end
18 end
19
20 M=length(plist);
21 f=fopen('phonelist.txt','wb');
22 if f<3,error('mal'); end
23 for m=1:M
24     fprintf(f,'%s\n',plist(m));
25 end
26 fclose(f);
27
28 f=fopen(in_proto,'rt');
29 if f<3, error('mal2'); end
30 t=char(fread(f));
31 fclose(f);
32 i=strfind(t,"proto");
33 t1 =t(1:i);
34 t2 = t(i+7:end);
35
36 f=fopen(out_mmffile,'wb');
37 if f<3,error('mal4'); end
38 fprintf(f,'%s',t1(1:end-4));
39
40 for m=1:M
41     name = plist(m);
42     fprintf(f,'%h "%s"%s',name,t2);
43 end
44 fclose(f);
```

---

## compareL1L2.m

```
1 function compareL1L2(ref_mlf,said_mlf,rec_mlf,dic,out_csv,mfc_dir)
2 %function compareL1L2(ref_mlf,said_mlf,rec_mlf,dic,out_csv,mfc_dir)
3 %
4 % ref_mlf -> the reference mlf with what was asked of the child
5 % said_mlf -> mlf with what the child actually said
6 % rec_mlf -> mlf output of the HVite recognition
7 %
8 % dic -> dictionary
9 % mfc_dir -> path to the mfc files : mfc_dir = '..\mfc\cut_clean\';
10 %
11 %Generates a CVS (TAB separated)file with a list of:
12 % REF SAID REC phones_Ref phones_free Nframes L1 L1-L2 ali_score al_sref al_sfree
13 %
14 %Ver. igual a compareL1L mas onde os tempos de início e fim de comando são
15 %tomados com base nos indicados em "MLF said_mlf"
16 % usage ex:
17 %compareL1L2('MLF_swlix.mlf','MLF_cut2_lix.mlf','recout.mlf','dic_newword.txt','swsaid.csv','..\m
18 fc\cut_clean\')
```



```

19 %compareL1L2('MLF_cut2_lix.mlf','MLF_cut2_lix.mlf','recout.mlf','dic_newword.txt','eqsaid.csv','.
20  \mfc\cut_clean\')
21 %NOTA: o rec_mlf só serve para ir buscar os 'pats' de interesse (test)
22 %NOTA: assumé que existe a mesma ordem dos pats e labs em SAID e REF.
23 %
24 %Serve para gerar out_csv para fazer histogramas de hipóteses H0/H1 com
25 %L1-L2, etc.
26
27 %cd C:\Users\Angel\Desktop\Thesis\childrenSR\DB;
28 [patref,labsref]=mlfRead2(ref_mlf,'%f%f%s');
29 [patsaid,labsaid]=mlfRead2(said_mlf,'%f%f%s');
30 [patrec,labsrec]=mlfRead2(rec_mlf,'%f%f%s%f');
31 Nrec = length(patrec);
32 %Ler dic (só deve ter as palavras-comando (sem [lix])
33 [words,transcs]=textread(dic,'%s%s','delimiter',char(9));
34
35 f = fopen(out_csv,'wb');
36 fprintf(f,'REF\tSAID\tREC\tphones\tFreephones\tNframes\tL1\tL1-
37 L2\talign_score\trefph_algn\tfreeph_algn\n');
38 %toma todos os ficheiros MFC (áudio) em rec_mlf
39 for nrec=1 : Nrec
40     Mrec = length(labsrec(nrec).ti);
41     %no rec_mlf o comando reconhecido é sempre o último.
42     labrec = labsrec(nrec).lab{Mrec};
43     pat = strrep(patrec(nrec),'.rec','.lab');
44
45     %finds same pat in REF
46     iref = strmatch(pat,patref,'exact');
47     m = strcmp('lix',labsref(iref).lab);
48     m = find(m==0);
49     m=m(end);
50     labref=labsref(iref).lab{m};
51
52     % sref = labsref(iref).lab{m};
53     % %finds same pat in SAID
54     % isaid = strmatch(wrec,patsaid,'exact');
55     % m = 1;
56     % while ~isempty(strmatch(labssaid(isaid).lab{m},'lix','exact'))
57     %     m = m + 1;
58     % end
59
60     %NOTA: assume que existe a mesma ordem em SAID e REF dos pats e labs
61     isaid = strmatch(pat,patsaid,'exact');
62     labsaid = labssaid(isaid).lab{m};
63
64     ti = labssaid(isaid).ti(m);
65     tf = labssaid(isaid).tf(m);
66
67     f1 = fix(ti*1e-5)+1; %frame rate de 100 frames por segundo
68     f2 = fix(tf*1e-5)-1;
69     Nframes = (f2-f1)+1;
70
71     % mfc_dir = '..\mfc\cut_clean\';
72     mfc = strrep(patrec{nrec},'*/',mfc_dir);
73     mfc = strrep(mfc,'.rec','.mfc');
74
75     [X,nSamples,sampPeriod,sampSize,paramKind]=readHTK(mfc,0);
76     writeHTK('temp.mfc',X(:,f1:f2),Nframes,sampPeriod,sampSize,paramKind,0);
77
78     %Corre HVite com temp.mfc para conhecer a verosimilhança L2 (com gramática fones livre
79     (sequência qualquer de qualquer fone + GV)
80
81     !HVite -C configbasico.txt -H hmm_mix16_3\phones_GV.mmf -l * -i tempfree.rec -w
82     phoneNet_GV.txt dic_phone_GV.txt phone_list_GV.txt temp.mfc
83
84     %Corre HVite com uma gramática que é a sequência de fones de REF.
85     %Para isso, vai buscar labref ao dicionário de transcrição dos comandos)
86     %e gera uma REDE com HParse.
87
88     %finds the labs.lab in the dictionary
89     i = strmatch(labref,words,'exact');
90     if isempty(i)
91         error('não há');
92     end
93     i=i(1);
94     % writes SAMPA of the word to make the grammar
95     fg=fopen('gtemp.txt','wb');
96     fprintf(fg,'%s\n',transcs{i});
97     fclose(fg);

```

```

98     phoneref = transcs{i};
99     phoneref(phoneref==' ')=[];
100
101     % makes wordnet from the grammar
102     !HParse -C configbasico.txt gtemp.txt ntemp.txt
103     !HVite -T 1 -C configbasico.txt -H hmm_mix16_3\phones_GV.mmf -l * -i tempcomm.rec -w
104     ntemp.txt dic_phone_GV.txt phone_list_GV.txt temp.mfc > log\HVite_2.log
105     flog = fopen('log\HVite_2.log','rt');
106     A='';
107     while 1
108         tline = fgetl(flog);
109         if ~ischar(tline), break, end
110         A=tline;
111     end
112     fclose(flog);
113     if strcmp(A,'No tokens survived to final node of network')
114         %Nframes too small to generate .rec file: put L1 =-infinity
115         L1 = -1e10;
116     else
117         [pat1, labs1]=mlfRead2('tempcomm.rec','%f%f%s%f');
118         %Tem a sequência de fones indicada em "ref_mlf" e em gtemp.txt/ntemp.txt
119         L1 = sum(labs1.score)/Nframes; %soma dos likes dos fones normalizada pelo Nframes
120     end
121
122     [pat2, labs2]=mlfRead2('tempfree.rec','%f%f%s%f');
123     L2 = sum(labs2.score)/Nframes;
124     phonefree = '';
125     for l=1 : length(labs2.lab)
126         phonefree = [phonefree, labs2.lab{l}];
127     end
128
129     phonefree_uc = phonefree;
130     phonefree_uc(find(diff(double(phonefree_uc))==0)+1)=[];
131     phonefree_uc = SAMPA2SAMPAuc(phonefree_uc);
132
133     phoneref_uc = phoneref;
134     phoneref_uc(find(diff(double(phoneref_uc))==0)+1)=[];
135     phoneref_uc = SAMPA2SAMPAuc(phoneref_uc);
136
137     if length(phoneref) < 2 || length(phonefree) < 2
138         score_ed = NaN;
139         s1 = phoneref_uc; s2=phonefree_uc(1:length(phoneref_uc));
140     else
141         [score_ed, s1, s2] = editdistance(phoneref_uc, phonefree_uc);
142     end
143
144     fprintf(f, '%s\t%s\t%s\t%s\t%s\t%d\t%d\t%d\t%d\t%s\t%s\n', ...
145         labref, labsaid, labrec, phoneref, phonefree, Nframes, L1, L1-L2, score_ed, s1, s2);
146 end
147
148 fclose(f);

```

---

## script para análise ROC e DET

```

1 % ROC and DET analysis
2
3
4 % errados
5 [wref1, ws1, recl1, ph1, phr1, nf1, L11, d1, s1, sa, sb]=textread('sw.csv', '%s%s%s%s%d%f%f%d%s%s', 'delimiter', '\t', 'headerlines', 1);
6
7 % correctos
8 [wref2, ws2, rec2, ph2, phr2, nf2, L12, d2, s2, sc, sd]=textread('eq.csv', '%s%s%s%s%d%f%f%d%s%s', 'delimiter', '\t', 'headerlines', 1);
9
10
11 % TP = sum(strcmp(ws1, recl1)); % true positive = Correctly accepted
12 % FP = sum(strcmp(wref2, rec2)); % false positive = Falsely accepted
13 % TN = sum(~strcmp(wref2, rec2)); % true negative = Correctly rejected
14 % FN = sum(~strcmp(ws1, recl1)); % false negative = Falsely rejected
15 %
16 % P = TP + FN;
17 % N = FP + TN;
18 % TPR = TP / P; % sensitivity
19 % FPR = FP / N;
20 % ACC = (TP + TN) / (P + N);
21

```

```

22 m=find(d1>100);
23 dl(m)=[];
24 Nbins=32;
25
26 %% L1-L2 correctos VS L1-L2 errados
27
28 % ROC
29 binmax=ceil(max([d1;d2]));
30 binmin=min([d1;d2]);
31 deltabin = (binmax-binmin)/Nbins;
32 bin = binmin + deltabin/2 +deltabin*(0:Nbins-1);
33 [count1] = hist(dl,bin); %errados
34 [count2] = hist(d2,bin); %correctos
35 plot(bin,count1,bin,count2,'r');
36 legend('L1-L2 Errados','L1-L2 Correctos');
37 grid on;
38
39 % DET
40 Ntrue = sum(count2);
41 Nfalse = sum(count1);
42 Ntotal = Ntrue + Nfalse;
43 fom = [d2;d1];
44 class0=false(size(fom));
45 class0(1:Ntrue)=true;
46 [NCA,NCR,NFA,NFR] = DET_DetectionErrorTradeoff(fom,class0);
47
48 %% align_score correctos VS align_score errados
49
50 % ROC
51 binmax2=ceil(max([s1;s2]));
52 binmin2=min([s1;s2]);
53 deltabin2 = (binmax2-binmin2)/Nbins;
54 bin2 = binmin2 + deltabin2/2 +deltabin2*(0:Nbins-1);
55 [count1] = hist(s1,bin2); %errados
56 [count2] = hist(s2,bin2); %correctos
57 plot(bin2,count1,bin2,count2,'r');
58 legend('Score de alinhamento - Errados','Score de alinhamento - Correctos');
59 grid on;
60
61 % DET
62 Ntrue = sum(count2);
63 Nfalse = sum(count1);
64 Ntotal = Ntrue + Nfalse;
65 fom = [s2;s1];
66 class0=false(size(fom));
67 class0(1:Ntrue)=true;
68 [NCA,NCR,NFA,NFR] = DET_DetectionErrorTradeoff(fom,class0);
69
70 %% Alguns valores de análise
71 avgL1 = mean(L11); % média L1 errados
72 avgd1 = mean(d1); % média d errados
73 avgL2 = mean(L12); % média L1 correctos
74 avgd2 = mean(d2); % média d correctos
75 medL1 = median(L11); % mediana L1 errados
76 medd1 = median(d1); % mediana d errados
77 medL2 = median(L12); % mediana L1 correctos
78 medd2 = median(d2); % mediana d correctos
79
80 [minL1,iminL1] = min(L11); % mínimo L1 errados
81 [mind1,imind1] = min(d1); % mínimo d errados
82 [minL2,iminL2] = min(L12); % mínimo L1 correctos
83 [mind2,imind2] = min(d2); % mínimo d correctos
84 [maxL1,imaxL1] = max(L11); %máximo L1 errados
85 [maxd1,imaxd1] = max(d1); %máximo d errados
86 [maxL2,imaxL2] = max(L12); %máximo L1 correctos
87 [maxd2,imaxd2] = max(d2); %máximo d correctos

```



## 8. Referências Bibliográficas

- [1] A. Öster, D. House, A. Protopapas and A. Hatzis, “Presentation of a new EU project for speech therapy: OLP (Ortho-Logo-Paedia),” in *Proceedings Fonetik 2002*, 2002.
- [2] M. Warschauer, “Technological Change and the Future of CALL,” em *New Perspectives on CALL for Second and Foreign Language Classrooms*, Mahwah, NJ, Lawrence Erlbaum Associates, 2004, pp. 15-25.
- [3] S. Bax, “CALL-Past, Present, and Future,” *System*, vol. 31, n.º 1, pp. 13-28, 2003.
- [4] A. Stoltzfus, “RosettaStone,” Rosetta Stone (UK) Limited, [Online]. Available: <http://www.rosettastone.co.uk/>. [Acedido em 2012].
- [5] “IBM Desktop ViaVoice,” IBM, [Online]. Available: <http://www-01.ibm.com/software/pervasive/viavoice.html>. [Acedido em 2011].
- [6] L. Heine, M. Witte, T. Diepstraten e T. Holl, “Babbel,” Lesson Nine GmbH, 2008. [Online]. Available: <http://pt.babbel.com/>. [Acedido em 2011].
- [7] “Babbel blog,” Lesson Nine GmbH, 10 June 2010. [Online]. Available: <http://blog.babbel.com/tech-background-babbel-speech-recognition/>. [Acedido em 2012].
- [8] “NUANCE,” Nuance Communications, [Online]. Available: [www.nuance.com](http://www.nuance.com). [Acedido em 2011].
- [9] I. Lessner, “The Street,” TheStreet, 02 10 2007. [Online]. Available: [http://www.thestreet.com/s/nuance-buys-new-york-software-firm/newsanalysis/techsoftware/10382384.html?puc=\\_googlen?cm\\_ven=GOOGLLEN&cm\\_cat=FREE&cm\\_ite=NA](http://www.thestreet.com/s/nuance-buys-new-york-software-firm/newsanalysis/techsoftware/10382384.html?puc=_googlen?cm_ven=GOOGLLEN&cm_cat=FREE&cm_ite=NA). [Accessed 2012].
- [10] “LumenVox Speech Understood,” LumenVox, [Online]. Available: <http://www.lumenvox.com>. [Acedido em 2012].
- [11] “Accessibility - Windows 7 features - Microsoft Windows,” Microsoft, [Online]. Available: <http://windows.microsoft.com/en-US/windows7/products/features/accessibility>.

[Acedido em 2012].

- [12] N. Anderson, “Win 7's built-in speech recognition: a review | Ars Technica,” *ars technica*, 1 June 2010. [Online]. Available: <http://arstechnica.com/information-technology/2010/05/win-7s-built-in-speech-recognition-a-review/>. [Acedido em 2012].
- [13] G. Fant, *Acoustic Theory of Speech Production*, The Hague: Mouton & Co, 1960.
- [14] “Vowel - Wikipedia, the free encyclopedia,” Wikipedia, [Online]. Available: <http://en.wikipedia.org/wiki/Vowel>. [Accessed 2012].
- [15] S. Ghai e R. Sinha, “Exploring the Effect of Differences in the Acoustic Correlates of Adults’ and Children’s Speech in the Context of Automatic Speech Recognition,” *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2010, 2010.
- [16] M. Gerosa, D. Giuliani e F. Brugnara, “Acoustic variability and automatic recognition of children’s speech,” *Speech Communication*, vol. 49, n.º 10-11, pp. 847-860, 2007.
- [17] S. Ghai e R. Sinha, “An Investigation into the Effect of Pitch Transformation on Children Speech Recognition,” em *TENCON*, Hyderabad, 2008.
- [18] A. Potamianos and S. Narayanan, “A Review of the Acoustic and Linguistic Properties,” in *Proc. Intern. Workshop on Multimedia Signal Processing*, Chania, Crete, Greece, 2007.
- [19] A. Hagen, B. Pellom and R. Cole, “Children’S Speech Recognition With Application to Interactive Books and Tutors,” in *Automatic Speech Recognition and Understanding*, 2003.
- [20] R. D. Kent, “Anatomical and neuromuscular maturation of the speech,” *Journal of Speech and*, vol. 9, p. 421–447, 1976.
- [21] P. Martland, S. P. Whiteside, S. W. Beet and L. Baghai-Ravary, “Estimating Child and Adolescent Formant Frequency Values From Adult Data,” in *Internat. Conf. Speech Language Processing*, Philadelphia, PA, 1996.
- [22] A. Potamianos and S. Narayanan, “Robust Recognition of Children’s Speech,” *IEEE Transactions on Speech and Audio Processing*, vol. 11, no. 6, pp. 603-616, 2003.

- [23] S. Narayanan e A. Potamianos, “Creating conversational interfaces for children,” *IEEE Transactions on Speech and Audio Processing*, vol. 10, n.º 2, pp. 65-78, 2002.
- [24] D. Elenius and M. Blomberg, “Comparing speech recognition for adults and children,” in *FONETIK*, Stockholm, 2004.
- [25] M. Gerosa, D. Giuliani, S. Narayanan e A. Potamianos, “A Review of ASR Technologies for Children’s Speech,” em *Proceedings of the Workshop on Child, Computer and Interaction*, Cambridge, MA, 2009.
- [26] C. Lopes, A. Veiga and F. Perdigão, “A European Portuguese Children Speech Database for Computer Aided Speech Therapy,” in *PROPOR 2012*, Coimbra, 2012.
- [27] “fone - Infopédia,” Porto Editora, [Online]. Available: URL: [http://www.infopedia.pt/\\$fone](http://www.infopedia.pt/$fone). [Acedido em 2012].
- [28] C. A. C. Lopes, “CLASSES FONÉTICAS ALARGADAS NO RECONHECIMENTO AUTOMÁTICO DE FONES,” Universidade de Coimbra, Coimbra, 2011.
- [29] H. P. Combrinck e E. C. Botha, “On The Mel-Scaled Cepstrum,” em *Seventh Annual South African Workshop on Patter Recognition*, Pretoria, 1996.
- [30] B. Bogert, M. Healy e J. Tukey, “The quefrency alanalysis of time series for echoes: Cepstrum, Pseudo-Autocovariance, Cross-Cepstrum and Saphe Cracking,” em *Proceedings of the Symposium on Time Series Analysis*, John Wiley & Sons, 1963, pp. 209-243.
- [31] D. O’Shaughnessy, “Speech Communications: Human and Machine,” in *Speech Communications: Human and Machine*, Addison-Wesley, 1987, p. 150.
- [32] L. E. Baum e T. Petrie, “Statistical Inference for Probabilistic Functions of Finite State Markov Chains,” *The Annals of Mathematical Statistics*, vol. 37, n.º 6, pp. 1554-1563, 1966.
- [33] S. E. Levinson, L. R. Rabiner e M. M. Sondhi, “Speaker Independent Isolated Digit Recognition Using Hidden Markov Models,” em *International Conference on Acoustics, Speech, and Signal Processing*, Boston, 1983.

- [34] L. R. Rabiner, “A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition,” em *Proceedings of the IEEE*, Orlando, 1989.
- [35] S. J. Young, G. Evermann, M. J. F. Gales, T. Hain, D. Kershaw, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev e P. C. Woodland, *The HTK Book*. Revised for HTK version 3.4, Cambridge: Cambridge University Engineering Department, 2006.
- [36] A. J. Viterbi, “Error bounds for convolutional codes and an asymptotically optimum decoding algorithm,” *Information Theory, IEEE Transactions on*, vol. 13, n.º 2, pp. 260-269, 1967.
- [37] G. D. Forney, “The Viterbi Algorithm,” *Proceedings of the IEEE*, vol. 61, n.º 3, pp. 268-278, 1973.
- [38] M. S. Ryan e G. R. Nudd, “The Viterbi Algorithm,” University of Warwick Coventry, Coventry, 1993.
- [39] L. R. Welch, “Hidden Markov Models and the Baum–Welch Algorithm,” *IEEE Information Theory Society Newsletter*, vol. 53, n.º 4, pp. 1,10-13, 2003.
- [40] A. Veiga, S. Candeias and F. Perdigão, “Generating a Pronunciation Dictionary for European Portuguese Using a Joint-Sequence Model with Embedded Stress Assignment,” in *STIL*, Cuiabá, Mato Grosso, 2011.
- [41] P. Boersma and D. Weenink, “Praat: doing phonetics by computer,” University of Amsterdam, Janeiro-Março 2010. [Online]. Available: <http://www.fon.hum.uva.nl/praat/>. [Accessed 2012].
- [42] K. MacLean, “VoxForge,” Ken MacLean. [Online]. Available: <http://www.voxforge.org/home>. [Acedido em 2012].
- [43] K. Sjölander, “An HMM-based system for automatic segmentation and alignment of speech,” em *Fonetik*, 2003.
- [44] “Receiver operating characteristic - Wikipedia, the free encyclopedia,” Wikipedia, [Online]. Available: [http://en.wikipedia.org/wiki/Receiver\\_operating\\_characteristic](http://en.wikipedia.org/wiki/Receiver_operating_characteristic). [Acedido em 2012].



- [45] A. Martin, G. Doddington, T. Kamm, M. Ordowski and M. Przybocki, "THE DET CURVE IN ASSESSMENT OF DETECTION TASK PERFORMANCE," in *Eurospeech*, Rhodes, Greece, 1997.
- [46] L. d. C. Moutinho e R. M. Lima, "Desempenho Fonético em Crianças dos 3 a 7anos de idade no P.E.," em *III Congresso Internacional de Fonética e Fonologia, Universidade Federal Minas Gerais*, Minas Gerais, Brasil, 2007.
- [47] S. Lee, A. Potamianos and S. Narayanan, "Acoustics of children's speech: Developmental changes of temporal and spectral parameters," *Journal of the Acoustical Society of America*, vol. 105, no. 3, 1999.
- [48] K. Shobaki, J.-P. Hosom e R. A. Cole, "The OGI KIDS' Speech Corpus and Recognizers," em *International Conference on Spoken Language Processing (ICSLP)*, Beijin, 2000.
- [49] A. Batliner, M. Blomberg, S. D'Arcy, D. Elenius, D. Giuliani, M. Gerosa, C. Hacker, M. Russell, S. Steidl and M. Wong, "The PF\_STAR Children's Speech Corpus," in *Interspeech*, Lisboa, 2005.
- [50] S. Young, "HTK," Cambridge University Engineering Department, [Online]. Available: <http://htk.eng.cam.ac.uk/>. [Acedido em 2011-2012].
- [51] A. F. Martin and e. al., "The DET Curve in Assessment of Detection Task Performance," in *Eurospeech*, Rhodes, Greece, 1997.