



UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE
CENTRO DE CIÊNCIAS EXATAS E DA TERRA
DEPARTAMENTO DE INFORMÁTICA E MATEMÁTICA APLICADA
BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO



Geração de prosódia para o português brasileiro em sistemas *text-to-speech*

Felipe Cortez de Sá

Natal-RN
Junho de 2018

Felipe Cortez de Sá

Geração de prosódia para o português brasileiro em
sistemas *text-to-speech*

Monografia de Graduação apresentada ao
Departamento de Informática e Matemática
Aplicada do Centro de Ciências Exatas e da
Terra da Universidade Federal do Rio Grande
do Norte como requisito parcial para a ob-
tenção do grau de bacharel em Ciência da
Computação.

Orientador

Carlos Augusto Prolo, doutor

UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE – UFRN
DEPARTAMENTO DE INFORMÁTICA E MATEMÁTICA APLICADA – DIMAP

Natal-RN

Junho de 2018

Monografia de Graduação sob o título *Geração de prosódia para o português brasileiro em sistemas text-to-speech* apresentada por Felipe Cortez de Sá e aceita pelo Departamento de Informática e Matemática Aplicada do Centro de Ciências Exatas e da Terra da Universidade Federal do Rio Grande do Norte, sendo aprovada por todos os membros da banca examinadora abaixo especificada:

Professor Doutor Carlos Augusto Prolo
Orientador
Departamento de Informática e Matemática Aplicada
Universidade Federal do Rio Grande do Norte

Professor Doutor Antônio Carlos Gay Thomé
Departamento de Informática e Matemática Aplicada
Universidade Federal do Rio Grande do Norte

Professora Doutora Erica Reviglio Iliovitz
Departamento de Letras
Universidade Federal do Rio Grande do Norte

Natal-RN, 20 de junho de 2018.

Dedicado à minha família e a Antônia

Agradecimentos

Agradeço à minha família, a Antônia, a Marc, a Prolo e a todos os meus amigos e professores.

Some few people are born without any sense of time. As consequence, their sense of place becomes heightened to an excruciating degree. They lie in tall grass and are questioned by poets and painters from all over the world. These time-deaf are beseeched to describe the precise placement of trees in the spring, the shape of snow on the Alps, the angle of sun on a church, the position of rivers, the location of moss, the pattern of birds in a flock. Yet the time-deaf are unable to speak what they know. For speech needs a sequence of words, spoken in time.

Alan Lightman, *Einstein's Dreams*

Geração de prosódia para o português brasileiro em sistemas *text-to-speech*

Autor: Felipe Cortez de Sá

Orientador(a): Professor Doutor Carlos Augusto Prolo

RESUMO

Com a cada vez mais forte presença de *smartphones* e *home assistants* no cotidiano, grandes empresas de tecnologia vêm desenvolvendo sistemas de conversação baseados em fala, denominadas *voice user interfaces*. Apesar dos avanços, é perceptível que os sistemas de síntese de voz, especialmente para o português brasileiro, deixam a desejar quanto à naturalidade da fala gerada. Um dos fatores principais que contribuem para isso é a prosódia, isto é, entoação, ritmo e acento da fala. Este trabalho investiga sistemas *text-to-speech* existentes através do estudo de seus algoritmos para síntese de voz e geração de prosódia para diversas línguas, com foco no português brasileiro. São explicitados os desafios encontrados, é feito um levantamento de modelos de análise prosódica na fonologia e propõem-se possíveis soluções para tornar a geração de voz mais próxima à humana.

Palavras-chave: text-to-speech, prosódia, voice user interfaces

Prosody generation for Brazilian Portuguese in text-to-speech systems

Author: Felipe Cortez de Sá

Advisor: Carlos Augusto Prolo, Ph.D.

ABSTRACT

With the evergrowing presence of smartphones and home assistants in our daily lives, technology companies have been developing two-way conversation systems, that is, voice user interfaces. Despite its recent improvements, text-to-speech programs still sound artificial, especially for their Brazilian Portuguese voices. A big contributing factor for that is the lack of accurate prosody, that is, pitch, length and emphasis. This thesis explores existing text-to-speech systems, especially those for which there are Brazilian Portuguese voices, focusing on their prosody generation modules. We highlight challenges of prosody generation, review prosodic analysis in the intonational phonology field and propose possible solutions for improving text-to-speech quality.

Keywords: text-to-speech, prosody, voice user interfaces

Lista de figuras

1	Arquitetura geral de sistemas TTS	p. 20
2	Representações temporal e espectral de uma gravação de voz	p. 23
3	Exemplo de arquivo de entrada para MBROLA	p. 26
4	Oração analisada utilizando o modelo ToBI	p. 28
5	Oração analisada utilizando o modelo DaTo	p. 30
6	Exemplo de texto anotado com SSML	p. 31
7	Exemplo de texto anotado com SSML e EmotionML (adaptado de Schröder e Burkhardt (2014))	p. 31
8	Anotações no modelo ToBI para o sistema TTS Festival	p. 32
9	Utilização do programa espeak e saída correspondente	p. 33
10	Extrato de linhas da tabela de conversão	p. 34
11	Arquitetura para o programa da interface gráfica	p. 36
12	Utilização por linha de comandos	p. 37
13	Editor gráfico desenvolvido	p. 37

Lista de tabelas

1	Regras para INTSINT	p. 29
2	Regras para DaTo	p. 30

Lista de códigos

A.1	Servidor	p. 43
A.2	Conversor eSpeakNG-MBROLA	p. 44
A.3	Script para editor gráfico	p. 48
A.4	Exemplo de resposta para <i>endpoint</i> do eSpeakNG	p. 50
A.5	Exemplo de resposta para <i>endpoint</i> do MBROLA	p. 51
A.6	Gerador de f_0 a partir do modelo INTSINT	p. 51

Lista de abreviaturas e siglas

TTS – *text-to-speech*

INTSINT – *International Transcription System for Intonation*

HMM – *Hidden Markov Model*

SSML – *Speech Synthesis Markup Language*

XML – *eXtensible Markup Language*

W3C – *World Wide Web Consortium*

MBRPSOLA – *Multi-Band Resynthesis Pitch Synchronous Overlap and Add*

API – *Application Programming Interface*

JSON – *JavaScript Object Notation*

REST – *Representational State Transfer*

AJAX – *Asynchronous JavaScript And XML*

Lista de símbolos

f_0 (frequência fundamental da fala)

ms (milissegundos)

Hz (Hertz)

Lista de algoritmos

- 1 Pseudocódigo para geração de frequências utilizando modelo INTSINT p. 35

Sumário

1	Introdução	p. 17
1.1	Objetivos	p. 17
1.2	Organização do trabalho	p. 18
2	Fundamentação teórica	p. 19
2.1	Prosódia	p. 19
2.1.1	Função prosódica	p. 19
	Afetiva	p. 19
	Suprasegmental	p. 20
	Aumentativa	p. 20
2.2	Sistemas <i>text-to-speech</i>	p. 20
2.2.1	Estrutura	p. 20
2.2.2	<i>Front end</i>	p. 21
2.2.2.1	Normalização de texto	p. 21
2.2.2.2	Conversão grafema-fone	p. 21
2.2.2.3	Geração de prosódia	p. 22
2.2.3	<i>Back end</i>	p. 22
2.2.3.1	Síntese articulatória	p. 22
2.2.3.2	Síntese por formantes	p. 22
2.2.3.3	Síntese concatenativa	p. 23
2.2.4	Síntese por <i>Hidden Markov Models</i> e <i>Deep Neural Networks</i> . .	p. 24

3	Revisão da literatura	p. 25
3.1	Sistemas TTS para o português brasileiro	p. 25
	Aiuruetê (BARBOSA et al., 1999)	p. 25
	eSpeakNG (DUNN, 2006)	p. 25
	(COUTO et al., 2010)	p. 25
	LianeTTS (SERPRO, 2011)	p. 26
3.1.1	MBROLA	p. 26
	3.1.1.1 Formato	p. 26
3.2	Modelos de análise entoacional	p. 27
3.2.1	Teoria métrica-autossegmental	p. 27
3.2.2	INTSINT	p. 28
3.2.3	DaTo (<i>Dynamic Tones</i>)	p. 29
3.3	Prosódia afetiva em sistemas TTS	p. 29
3.3.1	SSML	p. 30
3.3.2	EmotionML	p. 31
3.3.3	Anotação manual	p. 32
4	Módulo de prosódia	p. 33
4.1	Implementação	p. 33
4.1.1	espeak-ng	p. 33
4.1.2	MBROLA	p. 34
4.1.3	INTSINT	p. 34
4.1.4	Editor gráfico	p. 35
	4.1.4.1 Arquitetura	p. 36
	4.1.4.2 <i>Endpoints</i>	p. 36
	[POST] /api/espeak	p. 36
	[POST] /api/mbrola	p. 36

4.1.5	Utilização	p. 36
5	Resultados	p. 38
6	Considerações finais	p. 39
6.1	Trabalhos futuros	p. 39
	Referências	p. 40
	Apêndice A – Primeiro apêndice	p. 43

1 Introdução

Interfaces humano-computador que utilizam a voz, denominadas *voice user interfaces*, antigamente vistas apenas na ficção científica, hoje são uma realidade e estão disponíveis em *smartphones* e ambientes *desktop*. De acordo com Dutoit (1997), Jurafsky e Martin (2009), há uma grande área de aplicação para essas interfaces, como a acessibilidade, permitindo que deficientes visuais possam ouvir texto sem a necessidade de gravação prévia de seu conteúdo, ensino de linguagens e auxílio à pesquisa na linguística, por exemplo. Além disso, com o aumento da presença de sistemas embarcados no cotidiano – computadores de bordo, eletrodomésticos inteligentes, *smartphones*, entre outros –, é importante investigar novas formas de interação humano-máquina. Dessa forma, a síntese de fala, juntamente com o seu reconhecimento computacional, permite comunicação de duas vias com esses sistemas.

Os serviços mais populares e com resultados considerados mais naturais que temos atualmente são implementações proprietárias de grandes empresas, como Siri (Apple Inc., 2011), Cortana (Microsoft Corp., 2014) e Alexa (Amazon.com, Inc., 2014). Apesar das vantagens destacadas providas por *voice user interfaces*, os serviços disponíveis sintetizam voz com resultados perceptivelmente artificiais, principalmente para a língua portuguesa. Múltiplos trabalhos (HIRSCHBERG, 2006; RAJESWARI; UMA, 2012; TAYLOR, 2009) apontam como uma das maiores causas da artificialidade a prosódia empregada, isto é, aspectos rítmicos e melódicos da fala gerada.

1.1 Objetivos

Assim sendo, o presente trabalho tem como objetivo geral propor melhorias para o módulo de prosódia de sistemas *text-to-speech* (doravante TTS) com suporte ao português brasileiro, com foco na função afetiva e aumentativa da prosódia. A metodologia utilizada para atingir esse objetivo será de uma revisão da área da fonologia e estado da arte de algoritmos para sistemas TTS, com atenção especial à modelagem de prosódia, servindo

de base para o desenvolvimento de um módulo de geração de prosódia a ser integrado a um sistema TTS *open-source* existente.

1.2 Organização do trabalho

O presente trabalho está dividido em cinco capítulos, sendo o primeiro esta introdução.

No capítulo 2, são explicados em detalhes a arquitetura e funcionamento de sistemas TTS, destacando-se seus componentes principais e os algoritmos empregados em cada subsistema. Também é feita uma apresentação à área da fonologia entoacional, explicando-se conceitos fundamentais.

O capítulo 3, por sua vez, consiste de uma revisão da literatura, sendo realizado um levantamento dos sistemas *text-to-speech* existentes para o português brasileiro, a fim de demonstrar a forma como a prosódia é modelada em cada um deles. Além disso, explicita-se como é abordada a síntese de fala em trabalhos recentes e a análise de prosódia em um contexto não necessariamente computacional.

Já no capítulo 4, é descrito o sistema desenvolvido, justificando-se a abordagem com base na revisão da literatura. Explica-se a implementação do *software*, descrevendo-se sua arquitetura bem como algoritmos, ferramentas e linguagens de programação utilizados e os modos de operação do programa.

Finalmente, no capítulo 5, apresentam-se as conclusões feitas a partir da revisão bibliográfica e dos resultados obtidos com a implementação do módulo de prosódia para o português brasileiro.

2 Fundamentação teórica

2.1 Prosódia

Dutoit (1997), Ladd (2008) descrevem como parâmetros principais da prosódia o *pitch* – termo intercambiável com frequência fundamental ou f_0 –, intensidade e duração. Outros nomes para os mesmos fenômenos também são utilizados: Moraes (1998), por exemplo, separa as características principais da prosódia em *stress*, *accent* e *rhythm*. O que as definições têm em comum, porém, é que a prosódia compreende melodia e ritmo da fala, permitindo expressão comunicativa e realização de fenômenos “sintáticos”, como a acentuação de sílabas tônicas. Como os elementos prosódicos estão alinhados, isto é, ocorrem paralelamente a segmentos como sílabas e fones, a prosódia é dita um fenômeno suprasegmental (LADD, 2008).

Além desses componentes principais, Taylor (2009) destaca como aspecto relevante para análise entoacional o *downdrift* (declinação), isto é, a queda gradual do valor de f_0 durante uma enunciação.

2.1.1 Função prosódica

Taylor (2009) argumenta que uma das dificuldades do desenvolvimento de um bom modelo de prosódia se deve à falta de consideração da função comunicativa da fala, isto é, comumente análises são feitas ignorando o contexto da intenção do locutor. Apresenta, então, as três principais funções comunicativas para prosódia, descritas a seguir.

Afetiva A prosódia afetiva é utilizada para demonstrar emoções, atitude e intenção do locutor, e está praticamente ausente da forma escrita. Mesmo para diferentes línguas, os contornos melódicos associados a diferentes emoções são semelhantes, de forma que conseguimos deduzir relativamente bem uma intenção ou emoção de um locutor que fala uma língua estrangeira mesmo sem compreender o significado de cada palavra individualmente.

Suprasegmental Quando uma mensagem é dita de maneira declarativa, sem conteúdo afetivo significativo – descrita como *discourse neutral* –, ainda é possível observar variação de *pitch*, intensidade e duração. Na abordagem de Taylor (2009), essa parte da prosódia, dita suprasegmental, não é considerada conteúdo prosódico verdadeiro, mas sim um aspecto da fonética verbal. É possível ainda pensar em prosódia “real”, ou seja, afetiva e aumentativa, como desvios dos parâmetros suprasegmentais.

Aumentativa Além da prosódia afetiva, é possível desviar da prosódia padrão para assegurar a comunicação efetiva de uma mensagem sem adicionar informação ao conteúdo que está sendo dito. É usada, por exemplo, para enfatizar palavras, desambiguando uma mensagem que poderia ser interpretada de diferentes formas.

2.2 Sistemas *text-to-speech*

Burnett e Shuang (2010) definem *text-to-speech* como “o processo de geração automática de fala a partir de texto ou texto anotado” (tradução nossa). São utilizados em leitores digitais, assistentes pessoais para *smartphones*, aprendizagem de linguagens, entre outros.

Sistemas TTS são compostos por múltiplos subsistemas. Como veremos na seção 3.1, algumas implementações de TTS são modulares, permitindo desenvolvimento paralelo de cada componente individual. Isso possibilita que pesquisas possam focar em melhorias de uma parte específica do sistema sem necessidade de entendimento total. Neste trabalho, por exemplo, estudamos especificamente o módulo de prosódia.

2.2.1 Estrutura

Na literatura, cada sistema TTS normalmente emprega algoritmos diferentes, resultando em múltiplas arquiteturas possíveis para conversão de texto em fala. Apesar disso, Dutoit (1997) propõe uma arquitetura geral, dividindo os sistemas em duas partes principais descritas na Figura 1.

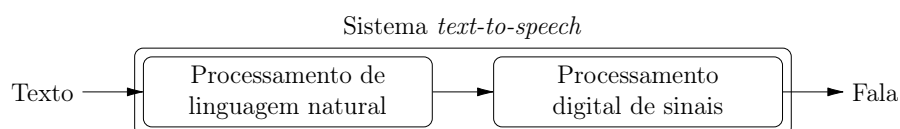


Figura 1: Arquitetura geral de sistemas TTS. Fonte: Adaptado de Dutoit (1997)

É comum encontrar em outros trabalhos o termo *front end* para o bloco de processamento de linguagem natural e *back end* para o bloco de processamento digital de sinais. Doravante utilizaremos essa nomenclatura.

2.2.2 *Front end*

O *front end* de um sistema TTS é responsável pela conversão do texto em sua representação em fones juntamente com parâmetros prosódicos. Em outras palavras, o bloco é responsável por determinar a pronúncia de cada palavra do texto, incluindo o contorno melódico e ritmo da fala. As definições seguintes são adaptadas de Jurafsky e Martin (2009).

2.2.2.1 Normalização de texto

O processo de conversão de texto para fala começa com o processamento do texto a fim de gerar uma representação grafêmica. A primeira etapa da normalização é a separação de sentenças, isto é, devem ser encontrados os limites de cada frase do texto. Frases que terminam com siglas ou abreviações podem dificultar o processo, como pode ser observado no exemplo seguinte.

Posteriormente, devem-se transformar símbolos, abreviações, siglas e outras *non-standard words* em suas representações pronunciáveis. Como exemplo, “V. Exa. me deve R\$ 50” deve ser lido como “vossa excelência me deve cinquenta reais”.

Por último, deve ser realizada a desambiguação de homônimos heterófonos: em “Gosto de pão”, “gosto” pode ser pronunciada como “gôsto” ou “gósto”, por exemplo.

2.2.2.2 Conversão grafema-fone

Com o texto normalizado, é preciso converter as letras em uma representação pronunciável, ou seja, fonemas ou fones. Para isso, normalmente é composto um conjunto de regras *letter-to-sound* ou letra-som, contendo as pronúncias comuns para sequências de letras, juntamente a um dicionário de pronúncia com palavras que não se adequam às regras. Para cada palavra, é feita uma busca no dicionário, e caso não seja encontrada, utilizam-se as regras.

2.2.2.3 Geração de prosódia

A partir do texto e fones gerados nas etapas anteriores, deve-se estimar *pitch*, intensidade e duração. Taylor (2009) explica que o desafio para implementação deste componente é que o texto praticamente não possui informação prosódica. Dutoit (1997) divide as abordagens para geração de prosódia automática em três técnicas: heurísticas derivadas a mão, sistemas baseados análise gramatical e métodos baseados em corpus.

Destacamos na seção 3.3 uma quarta abordagem possível, que consiste na anotação automática ou manual de prosódia, ênfase ou emoção, fornecendo “dicas” a sistemas TTS para melhor geração de prosódia.

2.2.3 *Back end*

Com o texto de entrada transformado em fones e informação prosódica, um *back end* é responsável por gerar uma forma de onda, ou seja, o áudio a ser reproduzido pelos alto-falantes. Jurafsky e Martin (2009), Taylor (2009) separam algoritmos de síntese em três classes, listadas a seguir.

2.2.3.1 Síntese articulatória

Sintetizadores articulatórios sintetizam fala através de modelos matemáticos, aproximando o aparelho fonador por uma série de tubos abertos. Pequenas alterações nos parâmetros de entrada podem gerar uma grande variação de sons. Em contrapartida, é difícil encontrar uma correspondência entre o texto de entrada e os parâmetros necessários para o sintetizador.

2.2.3.2 Síntese por formantes

Na Figura 2, podem-se ver duas representações para uma gravação de fala. Acima tem-se a amplitude em função do tempo e abaixo, frequências em função do tempo, isto é, a decomposição espectral. O conteúdo frequencial é majoritariamente composto por múltiplas senoides – também ditas harmônicas – de frequências f_0 , $f_1 = f_0 * 2$, $f_2 = f_0 * 3$ e $f_3 = f_0 * 4$. A frequência mais baixa é dita frequência fundamental ou apenas f_0 , emitida pela glote. As três senoides imediatamente acima, f_1 , f_2 e f_3 , são chamadas formantes, geradas pela ressonância do filtro resultante da posição da língua, queixo e lábio, alterando a intensidade de cada componente frequencial. A variação de intensidade de cada formante

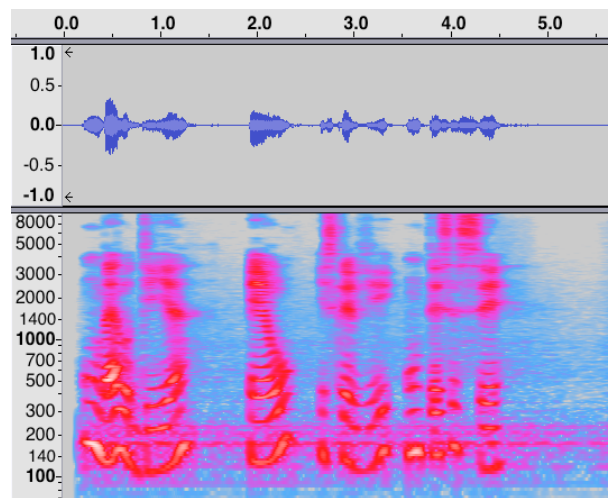


Figura 2: Representações temporal (acima) e espectral (abaixo) de uma gravação de voz. Fonte: própria

determina a vogal percebida. Sintetizadores por formantes são tentativas de modelagem da fala humana pela geração computacional dessas ondas. Apesar da simplicidade do modelo, além das senoides principais, pode ser visto na Figura 2 conteúdo harmônico significativo no sinal de uma fala humana. Ao ignorar-se essa parte, dita residual, a fala sintetizada é perceptivelmente robótica.

2.2.3.3 Síntese concatenativa

Devido aos problemas destacados para as outras duas abordagens, a maior parte dos sistemas TTS, até alguns anos atrás, utilizava síntese concatenativa, que consiste na gravação de múltiplas frases posteriormente divididas em fragmentos de granularidade variável, como ilustrado a seguir.

Para vocabulários com poucas palavras, como um sistema de anúncios de um metrô, é suficiente gravar frases completas e palavras que podem variar dentro da frase. Ao enunciar “Este metrô para nas estações Catete, Glória, Cinelândia [...]”, por exemplo, o fragmento “Este metrô para nas estações” pode advir de uma única gravação, enquanto os nomes individuais das estações são provenientes de outras gravações, concatenadas. Para obter sistemas TTS mais flexíveis, porém, as gravações devem ser separadas em unidades menores, como fones ou dífonos – pares de fones – que podem ser então concatenados para formar até mesmo palavras que não foram gravadas.

Mais recentemente, foi popularizada uma nova classe de algoritmos para síntese concatenativa denominada *unit selection synthesis*, abordagem probabilística que funciona através da minimização de duas função de custo, uma para fones individuais e outra para

a junção entre dois fones. Black (2002) argumenta que, com advento desse tipo de síntese, a fala gerada passou a ser natural o suficiente para ser possivelmente confundida com a de um humano.

2.2.4 Síntese por *Hidden Markov Models* e *Deep Neural Networks*

Zen, Tokuda e Black (2009) descrevem síntese por *Hidden Markov Models* como uma maneira paramétrica de produzir fala. Nela, um banco de dados com gravações é utilizado para treinar um modelo a partir de múltiplos parâmetros de entrada divididos em duas categorias: *excitation parameters* e *spectral parameters*. Com este método de síntese de fala, há unificação dos módulos de *front end* e *back end*. Mais recentemente, esse modelo, denominado *Statistical Parametric Speech Synthesis* foi adaptado para utilizar redes neurais profundas no lugar de cadeias de Markov (ZE; SENIOR; SCHUSTER, 2013). Os resultados mais naturais para síntese de fala até o momento foram obtidos por este paradigma.

3 Revisão da literatura

3.1 Sistemas TTS para o português brasileiro

São destacados a seguir sistemas TTS *open-source* com suporte ao português brasileiro, ordenados cronologicamente a fim de evidenciar a evolução das tecnologias utilizadas e as tendências para trabalhos futuros.

Aiuruetê (BARBOSA et al., 1999) Desenvolvido com o objetivo de obter um sistema TTS com fala natural sem custo computacional elevado, os autores descartam síntese articulatória e formantes justificando que o custo de desenvolvimento seria muito alto, optando por uma solução concatenativa. Para a prosódia, a parte de *pitch* utiliza curvas f_0 associadas a frases declarativas com declinação e as durações segmentais são obtidas por redes neurais treinadas a partir de parâmetros como número de vogais e clíticos.

eSpeakNG (DUNN, 2006) Projeto *open-source* com suporte a múltiplas linguagens. É modular, permitindo *back end* com síntese por formantes ou concatenativa. Também possibilita que seja utilizado apenas o *front end*, gerando representação no formato X-SAMPA (*Extended Speech Assessment Methods Phonetic Alphabet*), isto é, um alfabeto fonético simplificado. No seu módulo de prosódia, determina contorno f_0 a partir de pontuação e tabela com regras para diferentes posições silábicas: *pre-head*, *head*, *nucleus* e *tail*. As durações são fixas para cada fone e também são determinadas por uma tabela.

(COUTO et al., 2010) Foi desenvolvido com base no *framework* MaryTTS um sistema completo para português brasileiro baseado em HMM. O projeto iniciou com uma implementação de um módulo de conversão grafema-fone LaPS-G2P (BARBOSA et al., 2003), tornando-se então um sistema completo (COUTO et al., 2010), e foi subsequentemente estendido por Costa et al. (2012) para funcionar de maneira *stand-alone*, isto é, ser utilizado sem necessidade de instalação do *framework* MaryTTS. Como anteriormente explicado

na seção 2.2.4, sistemas construídos a partir HMM unificam front end e back end, deste modo a prosódia é gerada pelos parâmetros das cadeias de Markov, ou seja, determinada probabilisticamente a partir do corpus de treinamento.

LianeTTS (SERPRO, 2011) Projeto da SERPRO, o LianeTTS utiliza síntese concatenativa através do programa MBROLA, descrito na seção 3.1.1. Assim como o eSpeakNG, os subsistemas são baseados em regras e tabelas. O cálculo de curva f_0 é feito com base em *parts-of-speech tagging*, ou seja, atribuição de classes gramaticais a cada palavra do texto de entrada.

3.1.1 MBROLA

MBROLA é uma ferramenta para geração de voz baseada em síntese concatenativa por dífonos (pares de fones) desenvolvida com o objetivo de fomentar pesquisas acadêmicas em geração de prosódia (DUTOIT et al., 1996). É utilizada como *back end* para diversos sistemas TTS, como MaryTTS, Festival (BLACK; TAYLOR, 1997) e eSpeakNG, e possui três vozes disponíveis para o português brasileiro. A seguir, descrevemos como a ferramenta é utilizada.

Figura 3: Exemplo de arquivo de entrada para MBROLA

```

_ 150 50 150
t 70 50 125
e 125 50 75
c 70 50 125
e 125 50 75
c 70 50 125
e 116 20 232 80 300
_ 150 50 150

```

3.1.1.1 Formato

Como pode ser observado no Código 3, em cada linha tem-se um fone ou um silêncio (representado pelo *underscore*) seguido por uma duração em milissegundos e, por último, um ou mais pares de porcentagem e frequência em Hertz determinando alvos para a curva f_0 . Como exemplo, na penúltima linha temos o fone “e” com duração de 116 ms e dois alvos para altura; 232 Hz em 20% e 300 Hz em 80%. Para sintetizar a fala, a tabela é lida sequencialmente, concatenando pares de fones pré-gravados, e é feita uma interpolação linear entre cada alvo de *pitch*. Os dífonos pré-gravados têm frequência alterada utilizando

o algoritmo MBRPSOLA. Cada voz gravada fornece uma tabela com os difonos que podem ser utilizados.

3.2 Modelos de análise entoacional

Botinis, Granström e Möbius (2001) definem entoação como “a combinação de características tonais em unidades estruturais maiores associadas com o parâmetro acústico da frequência fundamental da fala e suas variações” (tradução nossa). Em outras palavras, a entoação diz respeito à variação da frequência base emitida pela glote quando falamos, como visto na seção 2.2.3.2. Múltiplas teorias tentam descrever o contorno melódico a partir de rótulos associados às variações de frequência. Destacamos a seguir três delas.

3.2.1 Teoria métrica-autossegmental

O modelo ToBI – *Tone Breaks and Indices* (SILVERMAN et al., 1992) faz parte da teoria métrica-autossegmental e descreve prosódia através de duas camadas. Para a entoação, tem-se o *Tonal tier*, que consiste na identificação de *pitch accents*: variações de f_0 que podem ser categorizadas a partir de combinações de dois rótulos “superior” (H, de *high*) e “inferior” (L, de *low*). As combinações possíveis do modelo ToBI são H^* , L^* , $H^* + L$, $H + L^*$, $L + H^*$ ou $L^* + H$, onde a estrela sinaliza o tom associado à sílaba acentuada. Além disso, *boundary tones* (ou tons de fronteira) podem ser representados com o símbolo % e demarcam fins e começos de sintagmas entoacionais.

A outra camada, *Break Index tier*, refere-se ao componente rítmico da prosódia, com rótulos de 0 a 6, onde 0 indica que duas palavras são pronunciadas juntas, e 6 representa uma separação máxima.

Alguns sistemas TTS, como MaryTTS e Festival, têm suporte a notação ToBI, como se pode ver no Código 3.3.3.

O modelo ToBI é utilizado por Moraes (2008) para analisar uma mesma frase com múltiplas intenções, mudando apenas o contorno melódico. Um exemplo de oração analisada utilizando o modelo pode ser visto na Figura 4.

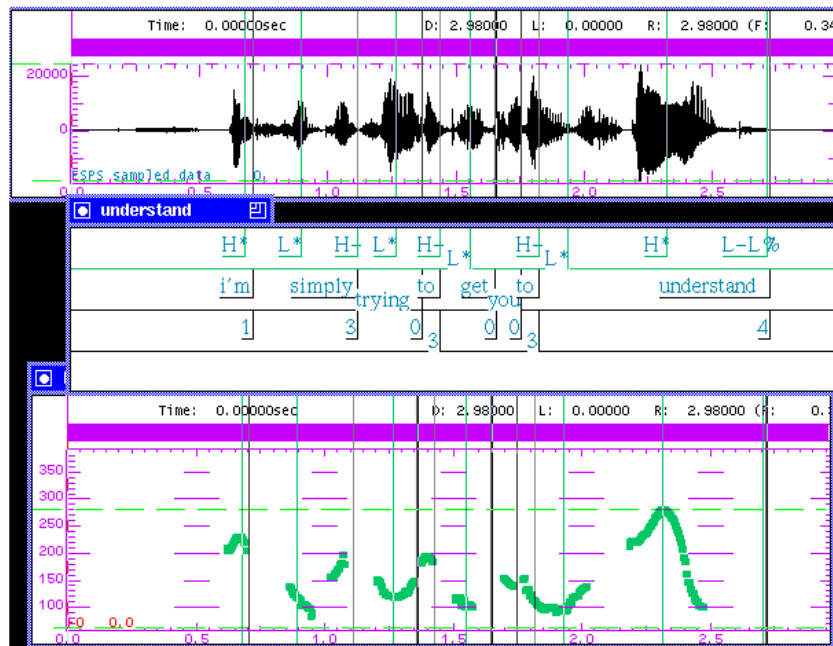


Figura 4: Oração analisada utilizando o modelo ToBI. Fonte: (CMU, 1994)

3.2.2 INTSINT

Proposto por Hirst (1987) com objetivo de criar uma representação equivalente ao alfabeto fonético IPA (*International Phonetic Alphabet*) para prosódia, o modelo INTSINT é um sistema de transcrição para entoação capaz de ser utilizado para analisar múltiplas línguas.

A entoação é anotada através de símbolos que indicam frequências-alvo. Os símbolos absolutos, Top (T), Mid (M) e Bottom (B), referem-se respectivamente ao limite superior, a base e o limite inferior para a frequência fundamental. Há adicionalmente três símbolos que modificam a marcação anterior, Higher (H), Same (S) e Lower (L), e, finalmente, dois símbolos denominados iterativos: Upstepped (U) e Downstepped (D), indicando subidas e descidas suaves.

Para calcular as frequências-alvo, utilizam-se dois parâmetros fixos *key*, definido em Hertz, que indica a f_0 média do locutor, e *range*, definido em oitavas, determinando o alcance máximo da entoação, isto é, quão intensas são as variações de *pitch*.

Apesar de ser inicialmente proposto como um modelo de análise, trabalhos subsequentes como Hirst (2007), Véronis et al. (1998) utilizaram a notação para sintetizar contornos f_0 . A Tabela 1 detalha como cada marcação é convertida para uma frequência em Hertz para determinação de uma sequência de valores f_0 .

O modelo INTSINT foi utilizado por Celeste e Reis (2012) e Moraes (1998) para analisar entoação para o português brasileiro.

Tabela 1: Regras para INTSINT

Regra	Cálculo
Top	$\text{key} \times \sqrt{2^{range}}$
Middle	key
Bottom	$\text{key} / \sqrt{2^{range}}$
Higher	$\sqrt{P_{i-1} \times T}$
Same	P_{i-1}
Lower	$\sqrt{P_{i-1} \times B}$
Upstepped	$\sqrt{P_{i-1} \times \sqrt{P_{i-1} \times T}}$
Downstepped	$\sqrt{P_{i-1} \times \sqrt{P_{i-1} \times B}}$

3.2.3 DaTo (*Dynamic Tones*)

Apesar de o modelo ToBI ter sido utilizado para analisar o português brasileiro em diversos trabalhos, Lucente (2014) argumenta que há características fonéticas que a notação da escola métrica-autossegmental não consegue expressar, propondo um novo modelo, DaTo, com base na entoação do português brasileiro.

No DaTo, em vez do direcionamento para a identificação de *pitch accents*, tem-se o conceito de “tarefa linguística por meio dos contornos entoacionais”, isto é, as transições entre diferentes pontos de ênfase tornam-se o meio de análise, em vez de simples preenchimento entre eventos estáticos, justificando o nome *Dynamic Tones*. As anotações possíveis para análise entoacionais são mostradas na Tabela 2 e um exemplo de oração analisada utilizando o modelo pode ser visto na Figura 5.

3.3 Prosódia afetiva em sistemas TTS

Destacamos a seguir ferramentas que podem auxiliar na geração de prosódia – principalmente com função afetiva – em sistemas TTS.

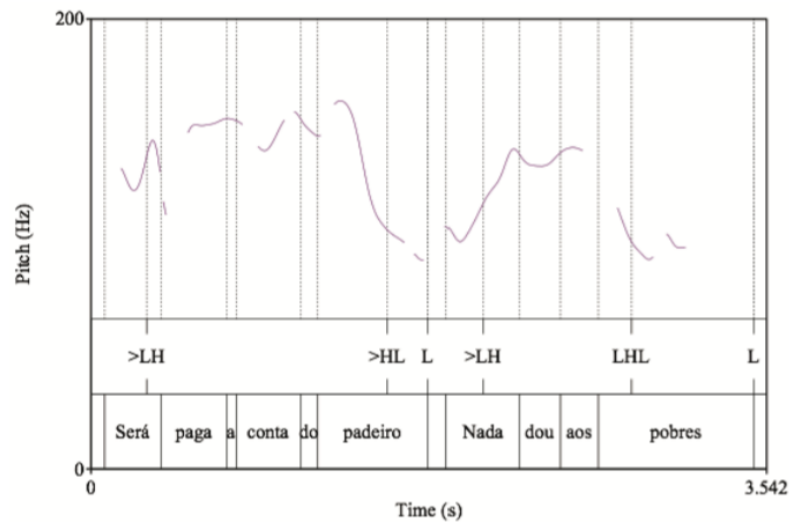


Figura 5: Oração analisada utilizando o modelo DaTo. Fonte: (LUCENTE, 2014)

Tabela 2: Regras para DaTo

Contorno	Notação
Rising	LH
Late rising	>LH
Compressed rising	vLH
Rising falling	LHL
Late falling	>HL
Compressed falling	vHL
Falling rising	HLH

3.3.1 SSML

A linguagem de marcação *Speech Synthesis Markup Language* foi criada motivada pela dificuldade de predição computacional de pronúncia (TAYLOR; ISARD, 1997). Quando originalmente proposta, alguns sistemas TTS permitiam anotações extra-textuais para auxiliar a estimação de parâmetros, mas cada sistema tinha sua própria notação. Consequentemente, usuários tinham que aprender um sistema de anotação para cada programa diferente. A proposta da SSML é padronizar notações, permitindo que sistemas TTS recebam o texto com marcações adicionais feitas utilizando uma linguagem unificada. A linguagem foi adotada por soluções *open-source* como MaryTTS e eSpeakNG, além dos sistemas proprietários encontrados em Alexa, Google Assistant e Cortana. A especificação

3.3.3 Anotação manual

Uma solução mais simples para geração de contornos f_0 para prosódia efetiva é permitir que o usuário especifique na entrada do programa TTS marcações prosódicas utilizando um dos modelos de análise entoacional vistos anteriormente na seção 3.2. Apesar de requerer conhecimento de fonologia, é uma opção viável enquanto não são desenvolvidos algoritmos para determinação de prosódia a partir de marcação emocional. O sistema Festival (BLACK; TAYLOR, 1997) disponibiliza uma maneira de especificar entoação seguindo o modelo ToBI, como pode ser visto no Código 3.3.3.

```
(Utterance Words
(The
(boy ((accent L*))
saw
the
(girl ((accent H*) (tone L-)))
with
the
(telescope ((accent H*) (tone H-H%))))))
```

Figura 8: Anotações no modelo ToBI para o sistema TTS Festival

4 Módulo de prosódia

Baseado na revisão da literatura, decidiu-se criar um módulo de prosódia que permite anotações prosódicas manuais, tornando possível variar a prosódia afetiva e aumentativa para um sistema TTS.

4.1 Implementação

4.1.1 `espeak-ng`

Foi utilizado o programa *open-source* eSpeakNG (DUNN, 2006) para realizar as primeiras etapas do *front end*: normalização de texto e conversão grafema-fone, ou seja, obter a partir do texto de entrada uma representação em fones. A saída gerada pelo programa é subsequentemente passada para o programa desenvolvido neste trabalho através da biblioteca `subprocess` da linguagem Python. O comando utilizado para obter os fones pode ser visto no Código 4.1.1.

- A *flag* `-v` seleciona uma voz. Neste caso, `pt-br`.
- `-x` e `-q` escrevem fones na saída em vez da fala sintetizada.
- `-sep=/` separa fones utilizando o caractere `/`.
- O caractere `'` indica *primary stress*, ou seja, o fone tônico da palavra.

```
$ espeak-ng -v pt-br 'Bom dia' -x -q --sep=/
$ b/'o/N dZ/'i/&
```

Figura 9: Utilização do programa `espeak` e saída correspondente

Apesar da existência de outras ferramentas para *front end* para o português brasileiro, optamos pelo eSpeakNG pela disponibilidade e facilidade de instalação.

4.1.2 MBROLA

O programa `sampa_mbrola.py` foi desenvolvido para converter a saída gerada pelo eSpeakNG mostrada anteriormente em fones possíveis de serem sintetizados pela voz utilizada no *back end*. O equivalente ao fone `/&/` gerado pelo eSpeak é `/a/` na voz do MBROLA, por exemplo.

A ferramenta de conversão lê a tabela do arquivo `sampa_mbrola.tbl` e substitui cada fone em notação XSAMPA pelo equivalente. A tabela possui três colunas: o fone que o eSpeak produz, o fone equivalente para o MBROLA e um exemplo numa palavra do português. No Código 4.1.2 é possível ver algumas linhas do arquivo. Utilizamos neste trabalho a voz `br3` desenvolvida por Denis R. Costa disponível no site oficial do projeto MBROLA.

b		b		_b_arco
k		k		_c_om
d		d		_d_ose
&		a		v_a_le
6		@		tam_a_nho
n		n		_n_unca

Figura 10: Extrato de linhas da tabela de conversão

Além da conversão XSAMPA-MBROLA, são determinadas as durações padrão para cada fone a partir da Tabela `durations.tbl`, adaptada do código utilizado pelo LianeTTS (SERPRO, 2011). Cada linha contém o fone e uma duração em milissegundos.

4.1.3 INTSINT

Como visto na seção 3.2.2, o INTSINT serve não apenas para análise, mas também para síntese de contornos f_0 . Escolhemos utilizá-lo na implementação pela simplicidade da notação e flexibilidade.

Em `intsint.py` (Código A.6) é feita a determinação de valores f_0 a partir do texto de entrada anotado e das regras descritas na Tabela 1. O algoritmo 1 ilustra o funcionamento do processo de conversão.

Algoritmo 1 Pseudocódigo para geração de frequências utilizando modelo INTSINT

Entrada: lista de rótulos P , $\forall i : P_i \in \{T, M, B, H, S, L, U\}$
Saída: lista de frequências f (em Hertz)

 $\text{freq_t} \leftarrow \text{key} \times \sqrt{2^{\text{range}}}$
 $\text{freq_m} \leftarrow \text{key}$
 $\text{freq_b} \leftarrow \text{key} / \sqrt{2^{\text{range}}}$
se rótulo $P_0 \notin \{T, M, B\}$ **então**

retorne erro % Primeiro rótulo deve ser absoluto

fim se
para todo rótulo P_i **faça**
se $P_i = T$ **então**
 $f_i \leftarrow \text{freq_t}$
senão se $P_i = M$ **então**
 $f_i \leftarrow \text{freq_m}$
senão se $P_i = B$ **então**
 $f_i \leftarrow \text{freq_b}$
senão se $P_i = H$ **então**
 $f_i \leftarrow \sqrt{f_{i-1} \times \text{freq_t}}$
senão se $P_i = S$ **então**
 $f_i \leftarrow f_{i-1}$
senão se $P_i = L$ **então**
 $f_i \leftarrow \sqrt{f_{i-1} \times \text{freq_b}}$
senão se $P_i = U$ **então**
 $f_i \leftarrow \sqrt{f_{i-1} \times \sqrt{f_{i-1} \times \text{freq_t}}}$
senão se $P_i = D$ **então**
 $f_i \leftarrow \sqrt{f_{i-1} \times \sqrt{f_{i-1} \times \text{freq_b}}}$
fim se
fim para

4.1.4 Editor gráfico

Para permitir refinamento de prosódia após cálculo inicial de parâmetros, foi desenvolvido um editor gráfico para *web* utilizando HTML, CSS e JavaScript. A duração e altura de cada fone podem ser especificadas manualmente arrastando barras de controle, gerando um novo arquivo de áudio com a fala sintetizada a partir dos novos parâmetros.

4.1.4.1 Arquitetura

O editor se comunica com o `espeak-ng` e MBROLA através de um servidor programado em Python utilizando o *framework* Flask para prover *endpoints* de uma API *RESTful*. Foi utilizada uma arquitetura cliente-servidor (FIELDING, 2000) comumente observada em aplicações *web* atuais. A maneira como os módulos se comunicam pode ser vista na Figura 11. Os módulos implementados neste trabalho estão com bordas espessas.

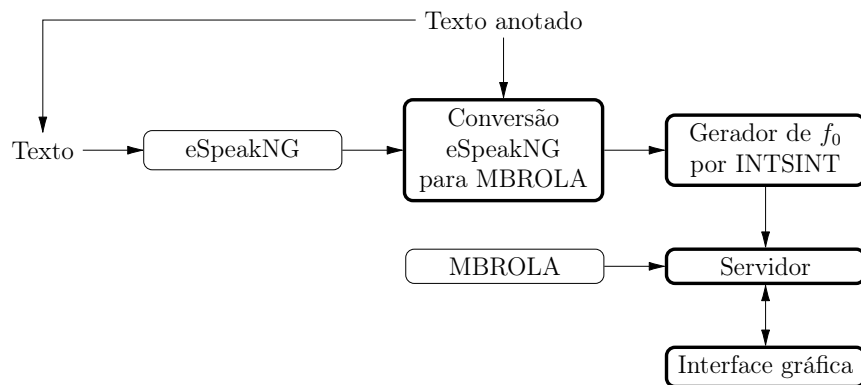


Figura 11: Arquitetura para o programa da interface gráfica

4.1.4.2 Endpoints

Foram criados dois *endpoints* para a API, possibilitando a comunicação entre interface gráfica e servidor.

[POST] /api/espeak Recebe um texto como entrada e gera uma resposta no formato JSON com lista de fones. Cada fone, por sua vez, possui campos `duration`, `phone_mbrola`, `phone_sampa` e `pitch_changes`. Um exemplo de resposta pode ser visto no Código A.4.

[POST] /api/mbrola Recebe uma lista de fones no formato MBROLA descrito na seção 3.1.1, passa a lista para o programa para criar o arquivo de áudio final e gera uma resposta com o nome do arquivo de áudio contendo a fala sintetizada.

4.1.5 Utilização

O programa pode ser executado através de uma interface de linha de comandos, como ilustrado no Código 4.1.5.

```
$ echo "Testando" | python3 sampa_mbrola.py > out.pho  
$ mbrola ../br3/br3 out.pho out.wav; afplay out.wav
```

Figura 12: Utilização por linha de comandos

Alternativamente, pode-se utilizar um navegador e ter acesso à interface gráfica descrita anteriormente e vista na Figura 13, permitindo edição de parâmetros para fones individuais.

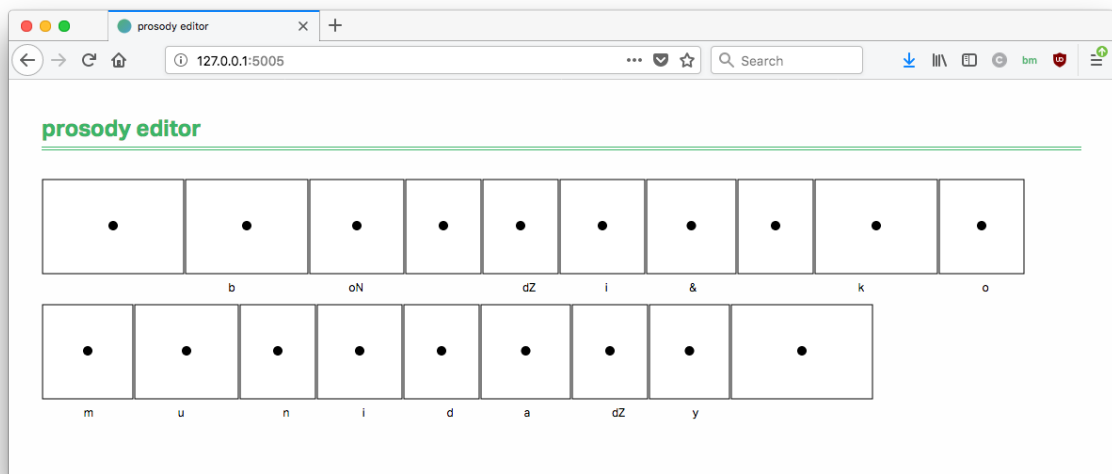


Figura 13: Editor gráfico desenvolvido

5 Resultados

Com o levantamento dos sistemas TTS desenvolvidos para o português brasileiro, é possível observar que há carência de suporte à síntese expressiva. Seguindo o modelo de funções prosódicas de Taylor (2009), os sistemas estudados concentram-se na geração de prosódia suprasegmental, mas não há suporte algum à geração de prosódia afetiva e aumentativa atualmente.

Linguagens de marcação como SSML e EmotionML já começaram a ser integradas a *frameworks open-source* e sistemas comerciais de TTS para múltiplas línguas, mas não foram encontrados trabalhos para o português brasileiro.

Estudando os modelos de anotação entoacional, percebe-se que ainda não há uma solução considerada mais apropriada para analisar o português brasileiro. Há, por outro lado, uma grande quantidade de trabalhos de análise de contornos melódicos que podem ser futuramente adaptados e convertidos em parâmetros para sistemas TTS expressivos, servindo como base para a geração manual ou automática de prosódia afetiva.

Identificamos lacunas que precisam ser preenchidas, como a criação de corpus anotados prosodicamente e algoritmos de compreensão textual para geração de marcação EmotionML.

Com o desenvolvimento dos programas para linha de comandos e para *web* neste trabalho, dá-se um passo inicial para a melhoria de prosódia afetiva e aumentativa em sistemas TTS para o português brasileiro.

6 Considerações finais

Neste trabalho foi apresentada a definição de um sistema TTS, bem como a descrição de seus componentes principais e possíveis implementações para cada módulo. Também foi realizada uma busca dos sistemas TTS *open-source* e comerciais existentes. Com o objetivo de melhorar a geração de prosódia, foi feita uma revisão do estado da arte, descrevendo os principais sistemas de anotação usados na fonologia entoacional, principalmente aplicados ao português brasileiro. Investigamos como a prosódia funciona nos sistemas encontrados e identificamos que, enquanto a geração de prosódia com função suprasegmental em trabalhos recentes já produz resultados satisfatórios, ainda há desafios quanto à síntese de fala expressiva, isto é, capaz de emitir contornos melódicos relacionados a emoções.

6.1 Trabalhos futuros

Considerando o que foi estudado e concluído, os próximos passos para melhorar a naturalidade da fala em sistemas TTS a partir do módulo de prosódia incluem:

- Expansão do sistema desenvolvido neste trabalho para adicionar suporte a mais modelos de notação entoacional, como ToBI e DaTo.
- Avaliação estatística da qualidade de falas geradas com a aplicação desenvolvida comparada a outros sistemas *open-source* e comerciais.
- Adicionar suporte a SSML e EmotionML ao sistema desenvolvido, gerando curvas f_0 a partir de marcações.
- Uso técnicas de técnicas de *Natural Language Understanding* para gerar notação SSML e EmotionML automaticamente.
- Desenvolvimento de corpus anotados com prosódia para o português brasileiro.

Referências

- Amazon.com, Inc. *Alexa*. 2014. Disponível em: <<https://developer.amazon.com/alexa>>. Acesso em 25 de setembro de 2017.
- Apple Inc. *Siri*. 2011. Disponível em: <<https://www.apple.com/ios/siri/>>. Acesso em 25 de setembro de 2017.
- BARBOSA, F. et al. Grapheme-phone transcription algorithm for a Brazilian Portuguese TTS. In: SPRINGER. *International Workshop on Computational Processing of the Portuguese Language*. [S.l.], 2003. p. 23–30.
- BARBOSA, P. A. et al. Aiuruete: A high-quality concatenative text-to-speech system for brazilian portuguese with demisyllabic analysis-based units and a hierarchical model of rhythm production. In: *Sixth European Conference on Speech Communication and Technology*. [S.l.: s.n.], 1999.
- BLACK, A. W. Perfect synthesis for all of the people all of the time. In: IEEE. *Speech Synthesis, 2002. Proceedings of 2002 IEEE Workshop on*. [S.l.], 2002. p. 167–170.
- BLACK, A. W.; TAYLOR, P. A. *The Festival Speech Synthesis System: System Documentation*. 1.1. ed. Scotland, UK, 1997. Available at <http://www.cstr.ed.ac.uk/projects/festival.html>.
- BOTINIS, A.; GRANSTRÖM, B.; MÖBIUS, B. Developments and paradigms in intonation research. *Speech communication*, Elsevier, v. 33, n. 4, p. 263–296, 2001.
- BURNETT, D.; SHUANG, Z. W. *Speech Synthesis Markup Language (SSML) Version 1.1*. [S.l.], 2010. Disponível em: <<http://www.w3.org/TR/2010/REC-speech-synthesis11-20100907/>>. Acesso em 5 de junho de 2018.
- CELESTE, L.; REIS, C. Análise entonativa formal: INTSINT aplicado ao português. *Journal of Speech*, v. 2, n. 2, p. 3–21, 2012.
- CHARFUELAN, M.; STEINER, I. Expressive speech synthesis in MARY TTS using audiobook data and emotionML. In: *INTERSPEECH*. [S.l.: s.n.], 2013. p. 1564–1568.
- CMU. *Guidelines for ToBI Labelling*. 1994.
- COSTA, E. S. et al. Um sintetizador de voz baseado em hmms livre: Dando novas vozes para aplicações livres no português do Brasil. In: *Workshop de Software Livre*. [S.l.: s.n.], 2012.
- COUTO, I. et al. An open source HMM-based text-to-speech system for Brazilian Portuguese. In: *7th international telecommunications symposium*. [S.l.: s.n.], 2010.

DUNN, R. H. *espeak-ng*. 2006. Disponível em: <<https://github.com/espeak-ng/espeak-ng>>. Acesso em 29 de outubro de 2017.

DUTOIT, T. *An introduction to text-to-speech synthesis*. [S.l.: s.n.], 1997. (Text, Speech and Language Technology 3).

DUTOIT, T. et al. The MBROLA project: Towards a set of high quality speech synthesizers free of use for non commercial purposes. In: IEEE. *Spoken Language, 1996. ICSLP 96. Proceedings., Fourth International Conference on*. [S.l.], 1996. v. 3, p. 1393–1396.

FIELDING, R. T. REST: architectural styles and the design of network-based software architectures. *Doctoral dissertation, University of California*, 2000.

HIRSCHBERG, J. Speech synthesis, prosody. *Encyclopedia of language & linguistics*, v. 7, p. 49–55, 2006.

HIRST, D. J. *La représentation linguistique des systèmes prosodiques: une approche cognitive*. Tese (Doutorado) — Aix-Marseille 1, 1987.

HIRST, D. J. A praat plugin for momel and intsint with improved algorithms for modelling and coding intonation. In: *Proceedings of the XVIth International Conference of Phonetic Sciences*. [S.l.: s.n.], 2007. v. 12331236.

JURAFSKY, D.; MARTIN, J. H. *Speech and Language Processing (2nd Edition)*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 2009. ISBN 0131873210.

LADD, D. R. *Intonational phonology*. [S.l.]: Cambridge University Press, 2008.

LUCENTE, L. Uma abordagem fonética na fonologia entoacional. *Fórum Linguístico*, v. 11, n. 1, p. 79–95, 2014.

Microsoft Corp. *Cortana*. 2014. Disponível em: <<https://www.microsoft.com/en-us/windows/cortana>>. Acesso em 25 de setembro de 2017.

MORAES, J. A. de. Brazilian portuguese. *Intonation systems: A survey of twenty languages*, Cambridge University Press, p. 179, 1998.

MORAES, J. A. de. The pitch accents in Brazilian Portuguese: analysis by synthesis. In: *Proc. Speech Prosody*. [S.l.: s.n.], 2008. p. 389–397.

RAJESWARI, K.; UMA, M. Prosody modeling techniques for text-to-speech synthesis systems—a survey. *International Journal of Computer Applications*, v. 39, n. 16, p. 8–11, 2012.

SCHRÖDER, M.; BURKHARDT, F. *Emotion Markup Language (EmotionML) 1.0*. [S.l.], maio 2014. Disponível em: <<http://www.w3.org/TR/2014/REC-emotionml-20140522/>>. Acesso em 8 de junho de 2018.

SERPRO. *LianeTTS*. 2011. Disponível em: <<http://intervox.nce.ufrj.br>>. Acesso em 7 de junho de 2018.

SILVERMAN, K. et al. Tobi: A standard for labeling english prosody. In: *Second international conference on spoken language processing*. [S.l.: s.n.], 1992.

TAYLOR, P. *Text-to-speech synthesis*. [S.l.]: Cambridge University Press, 2009.

TAYLOR, P.; ISARD, A. SSML: A speech synthesis markup language. *Speech communication*, Elsevier, v. 21, n. 1-2, p. 123–133, 1997.

VÉRONIS, J. et al. A stochastic model of intonation for text-to-speech synthesis1. *Speech Communication*, Elsevier, v. 26, n. 4, p. 233–244, 1998.

ZE, H.; SENIOR, A.; SCHUSTER, M. Statistical parametric speech synthesis using deep neural networks. In: IEEE. *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. [S.l.], 2013. p. 7962–7966.

ZEN, H.; TOKUDA, K.; BLACK, A. W. Statistical parametric speech synthesis. *Speech Communication*, Elsevier, v. 51, n. 11, p. 1039–1064, 2009.

APÊNDICE A – Primeiro apêndice

```

import sqlite3
import subprocess
import sampa_mbrola
import secrets
from flask import Flask, g, jsonify, render_template, abort, request, send_file
from flask_cors import CORS

app = Flask(__name__)
CORS(app)

@app.route('/api/espeak', methods=['POST'])
def frontend():
    text = request.form['text']
    converter = sampa_mbrola.Converter()
    sentence = converter.convert_sentence(text)
    return jsonify(sentence.dictify())

@app.route('/api/gen_audio', methods=['POST'])
def gen_audio():
    text = request.form["text"]
    converter = sampa_mbrola.Converter()
    sentence = converter.convert_sentence(text)

    token = secrets.token_hex(nbytes=4)
    with open(f"out/{token}.pho", 'w') as f:
        f.write(str(sentence))

    mbrola_str = f"mbrola br3/br3 out/{token}.pho out/{token}.wav"
    p = subprocess.Popen(mbrola_str, stdout=subprocess.PIPE, shell=True)
    (output, err) = p.communicate()
    p_status = p.wait()

    return jsonify({"token": token, "sentence": sentence.dictify()})

```

```

@app.route('/api/get_audio/<token>', methods=['GET'])
def get_audio(token):
    return send_file(
        f"out/{token}.wav",
        mimetype="audio/wav",
        as_attachment=True,
        attachment_filename=f"out/{token}.wav")

if __name__ == '__main__':
    app.run(debug=True)

```

Código A.1: Servidor

```

import sys
import re
import subprocess
import random
import json
import uuid
from math import sqrt
from collections import OrderedDict
from typing import List

def flatten(lst: list):
    return [item for sublist in lst for item in sublist]

class Phone():
    def __init__(
        self,
        phone_sampa: str,
        phone_mbrola: str,
        duration: int,
        pitch_changes: list,
        stress: bool = False
    ):
        self.phone_sampa = phone_sampa
        self.phone_mbrola = phone_mbrola
        self.duration = duration
        self.pitch_changes = pitch_changes
        self.stress = stress

    def as_line(self):
        return "{} {} {}".format(

```

```

        self.phone_mbrola,
        self.duration,
        " ".join([str(item) for item in flatten(self.pitch_changes)])
    )

class Sentence():
    def __init__(self, phones: List[Phone]=None):
        if phones is None:
            self.phones = []
        else:
            self.phones = phones

    def mbrola_phones(self):
        return [phone.phone_mbrola for phone in self.phones]

    def dictify(self):
        return [vars(phone) for phone in self.phones]

    def __repr__(self):
        return "\n".join([phone.as_line() for phone in self.phones])

class Converter():
    def __init__(self):
        self.load_sampa_mbrola()
        self.load_durations()

    def load_sampa_mbrola(self):
        equivs = {}
        with open("sampa_mbrola.tbl") as f:
            for line in f:
                k, v, _ = line.split()
                equivs[k] = v

        self.equivs = OrderedDict(
            sorted(equivs.items(), key=lambda t: -len(t[0])))

    def load_durations(self):
        durations = {}
        with open("durations.tbl") as f:
            for line in f:
                k, v = line.split()
                durations[k] = v

```

```

self.durations = durations

def convert_phoneme(self, sentence: str) -> tuple:
    """Returns first phone from the sentence"""

    if sentence[0] == " ":
        return ("_", 1)
    elif self.equivs:
        # s is a special case, needs to peek next
        if sentence[0] == "s":
            try:
                if sentence[1] in "aeiouyAEIOUY&":
                    return ("s", 1)
            except IndexError:
                pass

        for equiv in self.equivs.items():
            if re.match(re.escape(equiv[0]), sentence):
                # print("match:", equiv[0], phoneme, "=", equiv[1])
                return (equiv[1], len(equiv[0]))

    return (equiv[0], 1)

def get_duration(self, phoneme: str, factor: float) -> float:
    if self.durations and phoneme in self.durations:
        return int(float(self.durations[phoneme]) * factor)
    else:
        return 100 * factor

def convert_sentence(self, input_str: str) -> Sentence:
    key = 125
    octaves = 1.2
    decl = 0

    factor = 1 / 1.0
    freq = key

    freq_t = key * sqrt(2 ** octaves)
    freq_m = key
    freq_b = key / sqrt(2 ** octaves)

```

```

sentence = Sentence()
sentence.phones.append(Phone(" ", "_", 150 * factor, [[50, key]]))

ignored = ["@", "\n", ",", "'", "^", ";"]

for word in input_str.split():
    print(word)
    labeled = re.match(r"\[([TMBHSLUD])]", word)
    if labeled:
        word = re.sub(r"\[([TMBHSLUD])]", "", word)
        label = labeled.group(1)
        if label == "T":
            freq = freq_t
        elif label == "M":
            freq = freq_m
        elif label == "B":
            freq = freq_b
        elif label == "H":
            print(freq)
            freq = sqrt(freq * freq_t)
            print(freq)
        elif label == "S":
            freq = freq
        elif label == "L":
            freq = sqrt(freq * freq_b)
        elif label == "U":
            freq = sqrt(freq * sqrt(freq * freq_t))
        elif label == "D":
            freq = sqrt(freq * sqrt(freq * freq_b))

    sampa = self.text_to_sampa(word)
    sampa = sampa.replace("'", "")
    print(";; ", sampa)

    while sampa:
        if sampa[0] not in ignored:
            converted = self.convert_phoneme(sampa)
            phone_sampa = sampa[:converted[1]]
            sampa = sampa[converted[1]:]

            duration = self.get_duration(converted[0], factor)
            phone = Phone(
                phone_sampa = phone_sampa,

```



```

        phone_mbrola = converted[0],
        duration = duration,
        pitch_changes = [[50, freq + decl]],
        stress = False# percentage, Hz
    )

    sentence.phones.append(phone)
else:
    sampa = sampa[1:]
    decl -= 1

sentence.phones.append(Phone(" ", "_", 150 * factor, [[50, freq]]))
return sentence

def text_to_sampa(self, sentence: str) -> str:
    espeak_str = "espeak-ng -v pt-br '{}' -x -q".format(sentence)

    p = subprocess.Popen(espeak_str, stdout=subprocess.PIPE, shell=True)
    (output, err) = p.communicate()
    p_status = p.wait()
    output = output.decode("utf-8")
    output = output.replace("\n", "").strip()

    return output

if __name__ == "__main__":
    converter = Converter()

    for line in sys.stdin:
        print(";;", line)
        print(converter.convert_sentence(line))

```

Código A.2: Conversor eSpeakNG-MBROLA

```

function getRandomInt(max) {
    return Math.floor(Math.random() * Math.floor(max));
}

var self = this;
const requestURL = "http://127.0.0.1:5000/api/espeak";
let XHR = new XMLHttpRequest();
let FD = new FormData();

let textBox = document.querySelector("#textBox");

```

```

let audio = document.querySelector("#audio");
let audioSource = document.querySelector("#audioSource");

let postDiv = document.querySelector("#post");
postDiv.style.display = "none";

let genBtn = document.querySelector("#generateBtn");

let drawArea = document.querySelector("#drawArea");

genBtn.addEventListener("click", function( event ) {
  while (drawArea.firstChild) {
    drawArea.removeChild(drawArea.firstChild);
  }

  console.log(textBox.value);
  const requestURL = "http://127.0.0.1:5000/api/gen_audio";

  let XHR = new XMLHttpRequest();
  let FD = new FormData();
  FD.append("text", textBox.value);

  XHR.open("POST", requestURL);
  XHR.send(FD);

  XHR.onreadystatechange = function() {
    if (XHR.readyState === XMLHttpRequest.DONE) {
      if (XHR.status === 200) {
        const results = JSON.parse(XHR.responseText);
        console.log(results);
        console.log(results.token);

        let curLen = 0;
        let height = 250;

        for (let phone of results.sentence) {
          console.log(phone);
          let txtX = curLen + (phone.duration / 2);

          var svgns = "http://www.w3.org/2000/svg";
          var rect = document.createElementNS(svgns, 'rect');
          rect.setAttributeNS(null, 'x', curLen + phone.duration);

```

```

    rect.setAttributeNS(null, 'y', 0);
    rect.setAttributeNS(null, 'height', height);
    rect.setAttributeNS(null, 'width', '2');
    rect.setAttributeNS(null, 'fill', '#CCCCC');
    drawArea.appendChild(rect);

    var circ = document.createElementNS(svgns, 'circle');
    var cx = curLen + ((phone.pitch_changes[0][0] / 100.0) * phone.duration);
    circ.setAttributeNS(null, 'cx', cx);
    circ.setAttributeNS(null, 'cy', height - phone.pitch_changes[0][1]);
    circ.setAttributeNS(null, 'r', 5);
    circ.setAttributeNS(null, 'fill', '#777777');
    circ.setAttributeNS(null, 'class', "dot");
    drawArea.appendChild(circ);

    var text = document.createElementNS(svgns, 'text');
    var tx = curLen + (phone.duration / 2);
    text.setAttributeNS(null, 'x', cx);
    text.setAttributeNS(null, 'y', 20);
    text.setAttributeNS(null, 'text-anchor', "middle");
    text.setAttributeNS(null, 'style', "font: 14px sans-serif");
    text.setAttributeNS(null, 'fill', '#777777');
    text.textContent = phone.phone_mbrola;
    drawArea.appendChild(text);

    curLen += phone.duration;
  }
  console.log(curLen);
  drawArea.setAttribute("width", `${curLen}px`);

  audioSource.src = 'http://127.0.0.1:5000/api/get_audio/${results.token}';
  audio.load();
  postDiv.style.display = "block";
}
}
};

}, false);

```

Código A.3: Script para editor gráfico

```

[
  {
    "duration": 150,

```

```

        "phone_mbrola": "_",
        "phone_sampa": " ",
        "pitch_changes": [[50, 150]
    ]
},
{
    "duration": 80,
    "phone_mbrola": "n",
    "phone_sampa": "n",
    "pitch_changes": [[50, 150]]
},
{
    "duration": 110,
    "phone_mbrola": "u",
    "phone_sampa": "U",
    "pitch_changes": [[50, 150]]
},
...
{
    "duration": 150,
    "phone_mbrola": "_",
    "phone_sampa": " ",
    "pitch_changes": [[50, 150]]
}
]

```

Código A.4: Exemplo de resposta para *endpoint* do eSpeakNG

```

{
    "mp3_file": "07acc5c4ba924294.mp3"
}

```

Código A.5: Exemplo de resposta para *endpoint* do MBROLA

```

from math import sqrt

def intsint(labels: list, key: float = 150.0, range: float = 1.0):
    if labels[0] not in ["T", "M", "B"]:
        raise ValueError("First item in list can't be relative")

    freq_t = key * sqrt(2 ** range)
    freq_m = key
    freq_b = key / sqrt(2 ** range)

    freqs = []

```

```

for idx, l in enumerate(labels):
    if l == "T":
        freqs.append(freq_t)
    elif l == "M":
        freqs.append(freq_m)
    elif l == "B":
        freqs.append(freq_b)
    elif l == "H":
        freqs.append(sqrt(freqs[idx - 1] * freq_t))
    elif l == "S":
        freqs.append(freqs[idx - 1])
    elif l == "L":
        freqs.append(sqrt(freqs[idx - 1] * freq_b))
    elif l == "U":
        freqs.append(sqrt(freqs[idx - 1] * sqrt(freqs[idx - 1] * freq_t)))
    elif l == "D":
        freqs.append(sqrt(freqs[idx - 1] * sqrt(freqs[idx - 1] * freq_b)))

return freqs

```

Código A.6: Gerador de f_0 a partir do modelo INTSINT