



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

Dr. Carlos Augusto Prolo

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:

Dr. Carlos Augusto Prolo

Orientador

Departamento de Informática e Matemática Aplicada

Universidade Federal do Rio Grande do Norte

Dr. Antônio Carlos Gay Thomé

Departamento de Informática e Matemática Aplicada

Universidade Federal do Rio Grande do Norte

Dra. Erica Reviglio Iliovitz

Departamento de Letras

Universidade Federal do Rio Grande do Norte

Natal-RN, 20 de junho de 2018.

Dedicado a várias pessoas

Agradecimentos

Obrigado várias pessoas

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): Dr. 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, entonação, 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 linguística 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 Linguistics 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 (adaptado de Dutoit (1997))	p. 18
2	Arquitetura do sistema desenvolvido	p. 25

Lista de tabelas

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*

Lista de símbolos

ms (milisegundos)

Hz (Hertz)

Sumário

1	Introdução	p. 15
1.1	Objetivos	p. 15
1.2	Organização do trabalho	p. 16
2	Fundamentação teórica	p. 17
2.1	Prosódia	p. 17
2.1.1	Componentes	p. 17
2.1.2	Função prosódica	p. 17
2.1.3	Prosódia como elemento extra-textual	p. 17
2.2	Sistemas TTS	p. 18
2.2.1	Estrutura	p. 18
2.2.2	<i>Front end</i>	p. 18
2.2.2.1	Normalização de texto	p. 18
2.2.2.2	Conversão grafema-fone	p. 18
2.2.2.3	Geração de prosódia	p. 19
2.2.3	Back end	p. 19
3	Revisão da literatura	p. 20
3.1	Sistemas TTS	p. 20
3.1.1	MaryTTS	p. 20
3.1.2	LianeTTS	p. 20
3.1.3	espeak	p. 20

3.1.4	Microsoft	p. 20
3.1.5	Aiuruetê	p. 20
3.1.6	LINSE-DNN	p. 21
3.1.6.1	MBROLA	p. 21
3.1.6.2	Formato	p. 21
3.1.7	Prosódia em sistemas TTS	p. 22
3.1.7.1	SSML	p. 22
3.1.7.2	EmotionML	p. 22
3.1.8	Modelos de prosódia	p. 22
3.1.8.1	Teoria métrica-autossegmental	p. 22
3.1.8.2	IPO	p. 22
3.1.8.3	INTSINT	p. 23
3.1.8.4	DaTo	p. 23
4	Editor de prosódia	p. 24
4.1	Implementação	p. 24
4.1.1	espeak-ng	p. 24
4.1.2	MBROLA	p. 24
4.1.3	Arquitetura	p. 25
4.1.4	Módulo de prosódia	p. 25
4.1.5	Editor gráfico	p. 25
4.1.6	Endpoints	p. 25
	[POST] api/espeak	p. 26
	[POST] api/mbrola	p. 26
5	Resultados	p. 27
6	Considerações finais	p. 28

6.1	Trabalhos futuros	p. 28
	Referências	p. 29

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; 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, ferramenta de ensino de linguagens e auxílio à pesquisa na linguística. Além disso, com o aumento da popularidade de sistemas embarcados, é importante investigar novas formas de interação humano-máquina, e a síntese de fala, juntamente com o reconhecimento, permitem comunicação de duas vias com esses sistemas.

Os serviços mais populares e robustos 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 da praticidade e ganho de acessibilidade providos por essas interfaces, os serviços disponíveis sintetizam voz com resultados perceptivelmente artificiais, principalmente para a língua portuguesa. Uma das causas da artificialidade é a prosódia empregada, isto é, o ritmo, entonação e acento.

1.1 Objetivos

Este trabalho tem como objetivo geral propor melhorias para o módulo de prosódia afetiva para sistemas TTS com suporte ao português brasileiro através de uma revisão da área da fonologia e estudo de estado da arte de sistemas TTS, com atenção especial a como a prosódia é abordada.

1.2 Organização do trabalho

No capítulo 2, é feita uma breve explanação dos conceitos abordados neste trabalho, com foco em TTS e prosódia.

No capítulo 3, é feita uma revisão da literatura, fazendo um levantamento dos sistemas *text-to-speech* existentes tanto para o inglês quanto para o português brasileiro e como a prosódia é modelada em cada um deles. Mostramos como trabalhos recentes abordam síntese de fala. Também são descritos os trabalhos existentes em análise de prosódia em um contexto não necessariamente computacional.

No capítulo 4, é descrito o sistema desenvolvido, justificando a abordagem com base na revisão da literatura. Explica-se a implementação do software, descrevendo sua arquitetura, as linguagens de programações utilizadas e o funcionamento.

No capítulo 5, apresentamos os resultados obtidos com a implementação.

2 Fundamentação teórica

2.1 Prosódia

2.1.1 Componentes

(CAGLIARI et al., 1981, p. 150) cita. Stress (loudness and phonatory force) Syllabic length F0, intensidade, duração Intonational tune Sintagma entoacional Além desses componentes principais, Taylor (2009) descreve downdrift e microprosódia.

2.1.2 Função prosódica

Taylor (2009) divide a prosódia em três tipos distintos:

- Afetiva: prosódia pura, emoção, estado mental e atitude.
- Suprasegmental: quando uma mensagem é dita de maneira declarativa, sem conteúdo afetivo significativo, dá o nome de *discourse neutral*. utiliza um modelo de prosódia em que parte supra-segmental não é conteúdo prosódico verdadeiro, mas sim um outro aspecto do componente verbal.
- Aumentativa: assegura a comunicação efetiva de uma mensagem.

2.1.3 Prosódia como elemento extra-textual

Taylor (2009) argumenta que a geração de prosódia afetiva é árdua e depende de uma compreensão do texto. Ainda fala que, até a data de publicação do livro, nenhuma solução satisfatória foi encontrada e os sistemas TTS atuais pecam nesse aspecto.

2.2 Sistemas TTS

Burnett e Shuang (2010) define *text-to-speech* como “o processo de geração automática de fala a partir de texto ou texto anotado” (tradução nossa). Sistemas TTS são compostos por múltiplos subsistemas, como mostramos a seguir.

2.2.1 Estrutura

Na literatura, encontramos diversas arquiteturas específicas de um sistema TTS. Dutoit (1997) propõe um diagrama geral, dividindo sistemas em duas partes principais:

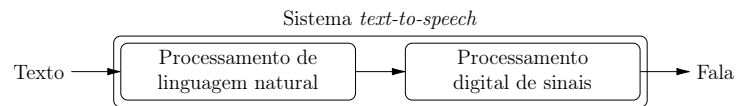


Figura 1: Arquitetura geral de sistemas TTS (adaptado de Dutoit (1997))

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*

2.2.2.1 Normalização de texto

O texto de entrada pode conter números e símbolos. A primeira parte do processo de conversão de texto para fala é transformar em sua representação por extenso. Sentence tokenization, isto é, descobrir as boundaries de cada sintagma. Depois, normalização de palavras não-padrão? abreviações, acrônimos e siglas. Desambiguação de homônimos heterófonos (e.g.: “Gosto de pão”). Análise fonética (normalized word strings para pronúncia).

2.2.2.2 Conversão grafema-fone

Palavras não-padrão são colocadas num dicionário de pronúncia. O resto é calculado de acordo com regras letra-som. Context-dependent e independent. Abordagens por machine learning ou regras. A conversão grafema-fone (ou fonema) consiste em transformar o texto normalizado em fones, ou seja, uma sequência de caracteres em uma sequência de fones.

2.2.2.3 Geração de prosódia

A partir do texto e fones gerados nas etapas anteriores, deve-se estimar uma curva F0, ritmo e ênfase.

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. Jurafsky e Martin (2009) separam algoritmos de síntese em três paradigmas: síntese concatenativa, síntese por formantes e síntese articulatória.

- Síntese articulatória
- Síntese por formantes
- Síntese concatenativa Síntese concatenativa pode trabalhar com diferentes atomicidades: nível de palavra, dífonos e fones individuais.

Abordagem utilizada pelo programa MBROLA. Consiste em gravar fala, separar pedaços de dois em dois.

Algumas vozes para o MaryTTS (SCHRÖDER; TROUVAIN, 2003) utilizam HMMs, isto é, Modelos ocultos de Markov para *unit selection*. Há, inclusive, uma voz brasileira feita a partir de HMMs: (COUTO et al., 2010). Projetos mais recentes como (WU; WATTS; KING, 2016; ZE; SENIOR; SCHUSTER, 2013) utilizam redes neurais para estimação de parâmetros. O trabalho da Google informa? que a estimação da curva F0 é uma possível melhoria.

3 Revisão da literatura

3.1 Sistemas TTS

3.1.1 MaryTTS

(COUTO et al., 2010), MaryTTS.

3.1.2 LianeTTS

Projeto da SERPRO LianeTTS. Utiliza MBROLA como *back end*. A geração de prosódia é simples, baseada em tabela.

3.1.3 espeak

Projeto do Dunn. Síntese por formantes (com consoantes gravadas). É possível obter apenas até a parte de g2p como saída.

3.1.4 Microsoft

Braga et al. (2008) descrevem sistema TTS desenvolvido pela Microsoft baseado em HMM. A stochastic-based LTS (Letter-to-Sound) converter to predict phonetic transcriptions for out-of-vocabulary words and the prosody models, which are data-driven using a prosody tagged corpus of 2000 sentences.

3.1.5 Aiuruetê

Usa PSOLA pra recodar unidades armazenadas (ortofones)

3.1.6 LINSE-DNN

statistical parametric speech synthesis

Popularização de redes neurais. Redes neurais utilizadas no back end. Trabalho pro português: (MAIA; SEARA, 2017)

3.1.6.1 MBROLA

MBROLA é uma ferramenta para geração de voz baseada em síntese por dífonos 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, e possui três vozes disponíveis para o português brasileiro.

```

_ 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

```

Listing 3.1: Exemplo de arquivo de entrada para MBROLA

3.1.6.2 Formato

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 F0. 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%.

Cada voz gravada provê uma tabela com os fones que podem ser utilizados. Utilizamos neste trabalho a voz br3 desenvolvida por Denis R. Costa disponível no site oficial do projeto MBROLA.

3.1.7 Prosódia em sistemas TTS

3.1.7.1 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, diferentes sistemas TTS permitiam anotações extra-textuais para auxiliar a estimação de parâmetros, mas usuários tinham que aprender um sistema de anotação para cada programa diferente. A proposta da SSML é que os sistemas TTS recebam o texto anotado numa linguagem unificada. A linguagem foi adotada pelo MaryTTS, além dos sistemas proprietários Alexa, Google Assistant e Cortana. A especificação é mantida pela W3C (BURNETT; SHUANG, 2010). A especificação cita os elementos *emphasis*, *break* e *prosody* como marcadores que podem auxiliar o processador de linguagem natural a gerar parâmetros prosódicos apropriados.

```
< speak >
  Siga
  < emphasis level="strong">aquele</ emphasis >
  carro.
</ speak >
```

Listing 3.2: Exemplo de texto anotado com SSML

3.1.7.2 EmotionML

EmotionML (SCHRÖDER; BURKHARDT, 2014) foi criada para várias coisas, uma delas é ajudar algoritmos a determinarem prosódia. (CHARFUELAN; STEINER, 2013) descreve um *framework* para implementação de determinação de prosódia a partir de anotações em emotionml.

3.1.8 Modelos de prosódia

3.1.8.1 Teoria métrica-autossegmental

3.1.8.2 IPO

(MIRANDA, 2015) analisa a prosódia para o português brasileiro através do modelo IPO, seguindo crítica de Lucente que o sistema ToBI é inapropriado.

3.1.8.3 INTSINT

INTSINT é um sistema de anotação para prosódia.

3.1.8.4 DaTo

Apesar de o modelo ToBI para anotação entoacional ter sido utilizado para analisar o português brasileiro em diversos trabalhos, ??) argumenta que há características perceptíveis que a notação não consegue expressar, propondo o modelo DaTo com base na entonação do português brasileiro.

4 Editor de prosódia

4.1 Implementação

4.1.1 `espeak-ng`

Foi utilizado o programa *open-source* `espeak-ng` (DUNN, 2006) para realizar a normalização de texto e realizar a conversão grafema-fone, ou seja, obter a partir do texto de entrada uma representação em fones. A saída gerada é passada para o programa desenvolvido neste trabalho através da biblioteca `subprocess` da linguagem Python. O comando utilizado para obter os grafemas pode ser visto no código 4.1.

- 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 /

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

Listing 4.1: Utilização do programa `espeak` e saída correspondente

Apesar da existência de outras ferramentas para *front-end* para o português brasileiro, optamos por esta pela facilidade de instalação, marcação de ênfase disponível, etc.

4.1.2 MBROLA

Cada voz do MBROLA tem suporte a um conjunto específico de dífonos. O programa `sampa_mbrola.py` foi desenvolvido para converter a saída gerada em fones suportados pela voz do MBROLA utilizada. 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 pelo equivalente corrigido. A tabela possui três colunas: o fone que

o `espeak` produz, o fone equivalente para o MBROLA e um exemplo numa palavra do português. Em 4.2 é possível ver algumas linhas do arquivo.

```
& | a | v_a_le
6 | @ | tam_a_nho
n | n | _n_unca
```

Listing 4.2: Linhas da tabela de conversão

4.1.3 Arquitetura

Foi utilizada uma arquitetura cliente-servidor utilizada por aplicações web contemporâneas. A interface gráfica consome uma API REST. A maneira como os módulos se comunicam pode ser vista na figura 2.

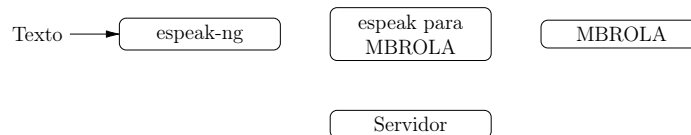


Figura 2: Arquitetura do sistema desenvolvido

4.1.4 Módulo de prosódia

4.1.5 Editor gráfico

O programa foi codificado em Python em sua versão 3.6. Pega resultado do `espeak-ng`, processa com editor gráfico e gera MBROLA.

Para alterar a prosódia manualmente, foi desenvolvido um editor gráfico para *web* utilizando HTML, CSS e JavaScript. A duração e altura de cada fone pode ser especificado arrastando barras de controle. 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 REST.

4.1.6 Endpoints

Foram criados dois endpoints para a API, tornando possível a comunicação entre interface gráfica e servidor.

[POST] api/espeak Recebe um texto como entrada e gera uma resposta JSON com lista de fones

[POST] api/mbrola Recebe uma lista de fones no formato MBROLA descrito na seção 3.1.6.1 e gera uma resposta com o nome do arquivo de áudio contendo a fala sintetizada

5 Resultados

Necessidade de um corpus anotado com prosódia para o português.

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 de prosódia na linguística, 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 suprasegmental em trabalhos recentes já produz resultados satisfatórios, ainda há desafios quando à síntese de fala expressiva relacionada à prosódia afetiva e aumentativa.

6.1 Trabalhos futuros

Foi identificado que ainda há um significativo desafio para a geração de prosódia afetiva automática a partir de um texto não anotado. Possíveis soluções para avanço nesta área incluem:

- Expansão do sistema desenvolvido neste trabalho para adicionar suporte a mais modelos notação entoacional.
- Adicionar suporte a SSML e EmotionML ao sistema desenvolvido, gerando curvas F0 a partir de linguagens de marcação
- Uso técnicas de técnicas de *Natural Language Understanding* para gerar notação 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.
- BRAGA, D. et al. HMM-based Brazilian Portuguese TTS. *Braga et al.(eds), Propor*, p. 47–50, 2008.
- 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.
- CAGLIARI, L. C. et al. Elementos de fonética do português brasileiro. [sn], 1981.
- 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.
- 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.
- JURAFSKY, D.; MARTIN, J. H. *Speech and Language Processing (2nd Edition)*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 2009. ISBN 0131873210.
- MAIA, R.; SEARA, R. Um sistema TTS baseado em redes neurais profundas usando parâmetros síncronos de pitch. p. 388–392, 2017.
- Microsoft Corp. *Cortana*. 2014. Disponível em: <<https://www.microsoft.com/en-us/windows/cortana>>. Acesso em 25 de setembro de 2017.
- MIRANDA, L. *Análise da entoação do português do Brasil segundo o modelo IPO*. Tese (Doutorado) — Dissertação de mestrado em Língua Portuguesa. Rio de Janeiro: UFRJ, 2015.

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.

SCHRÖDER, M.; TROUVAIN, J. The german text-to-speech synthesis system mary: A tool for research, development and teaching. *International Journal of Speech Technology*, Springer, v. 6, n. 4, p. 365–377, 2003.

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.

WU, Z.; WATTS, O.; KING, S. Merlin: An open source neural network speech synthesis system. *Proc. SSW, Sunnyvale, USA*, 2016.

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.