



Desenvolvimento de Base de Dados para Treinamento de Redes  
Neurais de Reconhecimento de Voz Através da Geração de Áudios  
com Resposta ao Impulso Simuladas por Técnicas de Data  
Augmentation

Bruno Machado Afonso

Projeto de Graduação apresentado ao Curso  
de Engenharia Eletrônica e de Computação  
da Escola Politécnica, Universidade Federal  
do Rio de Janeiro, como parte dos requisitos  
necessários à obtenção do título de Enge-  
nheiro.

Orientador: Mariane Rembold Petraglia

Rio de Janeiro

Julho de 2021

Desenvolvimento de Base de Dados para Treinamento de Redes  
Neurais de Reconhecimento de Voz Através da Geração de Áudios  
com Resposta ao Impulso Simuladas por Técnicas de Data  
Augmentation

Bruno Machado Afonso

PROJETO DE GRADUAÇÃO SUBMETIDO AO CORPO DOCENTE DO CURSO  
DE ENGENHARIA ELETRÔNICA E DE COMPUTAÇÃO DA ESCOLA PO-  
LITÉCNICA DA UNIVERSIDADE FEDERAL DO RIO DE JANEIRO COMO  
PARTE DOS REQUISITOS NECESSÁRIOS PARA A OBTENÇÃO DO GRAU  
DE ENGENHEIRO ELETRÔNICO E DE COMPUTAÇÃO

Autor:

---

Bruno Machado Afonso

Orientador:

---

Prof<sup>a</sup>. Mariane Rembold Petraglia, Ph.D.

Examinador:

---

Prof. José Gabriel Rodríguez Carneiro Gomes, D.Sc.

Examinador:

---

Prof. Julio Cesar Boscher Torres, D.Sc.

Rio de Janeiro

Julho de 2021

## Declaração de Autoria e de Direitos

Eu, *Bruno Machado Afonso* CPF 136.151.347-02, autor da monografia *Desenvolvimento de Base de Dados para Treinamento de Redes Neurais de Reconhecimento de Voz Através da Geração de Áudios com Resposta ao Impulso Simuladas por Técnicas de Data Augmentation*, subscrevo para os devidos fins, as seguintes informações:

1. O autor declara que o trabalho apresentado na disciplina de Projeto de Graduação da Escola Politécnica da UFRJ é de sua autoria, sendo original em forma e conteúdo.
2. Excetua-se do item 1. eventuais transcrições de texto, figuras, tabelas, conceitos e ideias, que identifiquem claramente a fonte original, explicitando as autorizações obtidas dos respectivos proprietários, quando necessárias.
3. O autor permite que a UFRJ, por um prazo indeterminado, efetue em qualquer mídia de divulgação, a publicação do trabalho acadêmico em sua totalidade, ou em parte. Essa autorização não envolve ônus de qualquer natureza à UFRJ, ou aos seus representantes.
4. O autor pode, excepcionalmente, encaminhar à Comissão de Projeto de Graduação, a não divulgação do material, por um prazo máximo de 01 (um) ano, improrrogável, a contar da data de defesa, desde que o pedido seja justificado, e solicitado antecipadamente, por escrito, à Congregação da Escola Politécnica.
5. O autor declara, ainda, ter a capacidade jurídica para a prática do presente ato, assim como ter conhecimento do teor da presente Declaração, estando ciente das sanções e punições legais, no que tange a cópia parcial, ou total, de obra intelectual, o que se configura como violação do direito autoral previsto no Código Penal Brasileiro no art.184 e art.299, bem como na Lei 9.610.
6. O autor é o único responsável pelo conteúdo apresentado nos trabalhos acadêmicos publicados, não cabendo à UFRJ, aos seus representantes, ou ao(s) orientador(es), qualquer responsabilização/ indenização nesse sentido.
7. Por ser verdade, firmo a presente declaração.

---

Bruno Machado Afonso

UNIVERSIDADE FEDERAL DO RIO DE JANEIRO

Escola Politécnica - Departamento de Eletrônica e de Computação

Centro de Tecnologia, bloco H, sala H-217, Cidade Universitária

Rio de Janeiro - RJ CEP 21949-900

Este exemplar é de propriedade da Universidade Federal do Rio de Janeiro, que poderá incluí-lo em base de dados, armazenar em computador, microfilmear ou adotar qualquer forma de arquivamento.

É permitida a menção, reprodução parcial ou integral e a transmissão entre bibliotecas deste trabalho, sem modificação de seu texto, em qualquer meio que esteja ou venha a ser fixado, para pesquisa acadêmica, comentários e citações, desde que sem finalidade comercial e que seja feita a referência bibliográfica completa.

Os conceitos expressos neste trabalho são de responsabilidade do(s) autor(es).

## AGRADECIMENTO

Agradeço, em primeiro lugar, à Deus que guiou todos os meus passos nessa jornada acadêmica e continua a guiar os caminhos da minha vida. Sem Ele, não seria capaz de terminar este projeto e não chegaria nos patamares onde estou. A Ele a glória.

Gostaria de agradecer também aos meus pais, em especial à minha mãe, que graças ao seu infinito empenho, dedicação, amor e carinho que me impulsionavam a cada dia, estou hoje concluindo esta etapa de minha vida.

Não menos importante, agradeço também à minha orientadora, Mariane Rembold Petraglia, que dedicou seu tempo para que eu pudesse me tornar um profissional ainda melhor. Agradeço pela paciência, atenção, boa vontade e o voto de confiança que foram depositados em mim.

Agradeço também aos meus colegas de faculdade, pois sem eles eu não estaria nem no meu último período da graduação para finalizar o meu projeto. Em especial, agradeço aos meus amigos Felipe Claudio e Diogo Nocera, ambos por serem os mais presentes tanto na jornada acadêmica quanto para o resto da minha vida.

## RESUMO

O tema de reconhecimento de voz se torna cada vez mais relevante graças ao seu amplo uso tecnológico na sociedade, desde assistentes pessoais em *smartphones* e automação residencial, até autenticação por voz para aplicações de segurança.

Uma das características mais importantes neste tema é a detecção da Resposta ao Impulso de Salas (RIR), que representa o modelo acústico do ambiente. A RIR é usada no processamento de áudio para identificação e reconhecimento de fontes sonoras em campo distante, que é formada, no caso do tema de reconhecimento de voz, por uma amostra de voz anecóica convoluída com a RIR, acrescida de um ruído, que pode ser pontual ou disperso.

Um dos desafios no reconhecimento de voz é a estimação da RIR em um sinal de voz em campo distante. Além das técnicas tradicionais de processamento de sinais, diversas soluções de *deep learning* foram propostas para a estimação da RIR, contudo estas acabam sendo limitadas devido à falta de variedade e quantidade de bases de RIRs medidas disponíveis para treinamento de redes neurais.

Neste contexto, o objetivo deste projeto é de desenvolver um algoritmo, usando técnicas de *data augmentation*, que gera amostras de voz em campo distante (AVCD), construindo assim uma base de dados para uso em treinamentos de soluções de *deep learning*. O algoritmo é composto por dois segmentos de *data augmentation*: o primeiro modifica as características de razão direto-reverberante (DRR) e tempo de reverberação (T60) do sinal acústico, partindo de RIRs reais, gerando RIRs simuladas (RIRSM); o segundo gera AVCDs, convoluindo amostras de voz anecóicas com as RIRSMs e adicionando ruídos à voz reverberada. Ao final do trabalho, são exibidos exemplos de AVCDs geradas pelo algoritmo proposto, analisando se os dados gerados são válidos para uso em treinamento de redes neurais.

Palavras-Chave: Resposta ao Impulso de sala, *data augmentation*, *deep learning*, reconhecimento de voz.

## ABSTRACT

Speech recognition is a very relevant topic in the present days due to its vast technological usage on modern society from personal assistants on smartphones and residential automated systems, to voice authentication for security applications.

One of the most important characteristics on this topic is the Room Impulse Response (RIR) detection, which represents the acoustic model of the room. The RIR is used on signal processing to identify and recognize far-field audio sources, which for the speech recognition topic, is the anechoic voice sample convolved with the RIR plus noise signal.

One of the challenges when it comes to speech recognition is to estimate the RIR in a far-field voice sample. Beyond the traditional signal processing algorithms, many deep learning solutions are proposed for the RIR estimation, however they end up with limited results due to the lack of variety e quantity of RIR databases available for training.

In this context, the main objective of this project is to develop an algorithm using data augmentation techniques that will generate far-field voice samples, therefore building a database for deep learning training. The algorithm is composed of two segments: the first modifies the real RIRs characteristics of direct-to-reverberant ratio (DRR) and the reverberation time ( $T_{60}$ ) of the acoustic signal generating simulated RIRs (RIRSM); the second generates far-field voice samples using the previously created RIRSMs, anechoic voice samples and noise signals. At the end of this work, examples of the generated far-field voice samples by the algorithm are shown and they are analysed to see if they are valid to be used in neural network training.

Key-words: Room Impulse Response, data augmentation, deep learning, voice recognition.

## SIGLAS

AIR - Aachen Impulse Response

AVA - Amostra de voz anecóica

AVCD - Amostra de voz em campo-distante

AVR - Amostra de voz reverberada

DA - *Data Augmentation*

DL - *Deep Learning*

DRR - Razão Direto-Reverberante

DTMF - Dual-Tone Multi-Frequency

RIR - Resposta ao Impulso de Sala

RIRDA - *Data Augmentation* da Resposta ao Impulso de Sala

RIRO - Resposta ao Impulso de Sala Original

RIRSM - Resposta ao Impulso de Sala Simulada

SNR - Razão Sinal-Ruído

SRF - Sinal de ruído de fundo

SRP - Sinal de ruído pontual

T20 - Tempo de Reverberação (queda de 20 dB)

T30 - Tempo de Reverberação (queda de 30 dB)



T60 - Tempo de Reverberação (queda de 60 dB)

VA - Voz anecóica

VR - Voz reverberada

# Sumário

<b>1</b>	<b>Introdução</b>	<b>1</b>
1.1	Tema . . . . .	1
1.2	Delimitação . . . . .	1
1.3	Justificativa . . . . .	1
1.4	Objetivos . . . . .	2
1.5	Metodologia . . . . .	3
1.6	Descrição . . . . .	4
<b>2</b>	<b>Análise de Fontes Sonoras e seus Desafios</b>	<b>5</b>
2.1	Resposta ao Impulso de Ambiente Acústico e suas Aplicações . . . . .	5
2.2	Desafios correlacionados à RIR . . . . .	7
2.3	<i>Data Augmentation</i> . . . . .	9
<b>3</b>	<b><i>Data Augmentation</i> aplicado à Resposta ao Impulso de Ambientes</b>	<b>10</b>
3.1	Razão Direto-Reverberante (DRR) . . . . .	11
3.2	Tempo de Reverberação (T60) . . . . .	13
<b>4</b>	<b>Sinais de Voz com RIRs Simuladas e com Ruídos</b>	<b>16</b>
4.1	Simulação de fala em campo distante . . . . .	17
<b>5</b>	<b>Bases de Dados</b>	<b>23</b>
5.1	Base de amostras de voz anecóicas . . . . .	23
5.2	Base de RIRs - Aachen Impulse Response database (AIR) . . . . .	24
5.3	Base de ruídos - MUSAN . . . . .	25
<b>6</b>	<b>Resultados Experimentais</b>	<b>27</b>
6.1	Configuração dos parâmetros . . . . .	27

6.2	Resultados . . . . .	28
6.2.1	<i>Data Augmentation</i> do DRR . . . . .	28
6.2.2	<i>Data Augmentation</i> do T60 . . . . .	30
6.2.3	<i>Data Augmentation</i> de fala em campo distante . . . . .	33
<b>7</b>	<b>Conclusões</b>	<b>37</b>
	<b>Bibliografia</b>	<b>39</b>
<b>A</b>	<b>Código Fonte</b>	<b>42</b>

# Lista de Figuras

2.1	Representação de uma sala anecóica e reverberante. . . . .	6
2.2	Gráficos de quantidade de artigos publicados por ano relacionados com <i>Deep Learning</i> . . . . .	8
2.3	Fluxo geral de procedimentos para gerar sinais de voz reverberantes. .	9
3.1	Fluxo de procedimentos para gerar a RIRSM. . . . .	10
3.2	Um exemplo de $h(t)$ com $h_e(t)$ , onde é feita a modificação do DRR de -4,5 para 4, marcado em vermelho. . . . .	12
3.3	Um exemplo de $h_e(t)$ antes e depois da aplicação do algoritmo, origi- nal em azul ( $DRR = -4,5$ dB) e o modificado em vermelho ( $DRR =$ 4 dB). . . . .	13
3.4	Um exemplo de $h(t)$ com $h_l(t)$ , onde é feita a modificação do T60 de 1,38 para 2,6 segundos, marcado em vermelho. . . . .	15
3.5	Um exemplo de um trecho ao final de $h_l(t)$ antes e depois da aplicação do algoritmo, original em azul ( $T60 = 1,38$ s) e o modificado em vermelho ( $T60 = 2,6$ s). . . . .	15
4.1	Fluxo de procedimentos para gerar a AVCD. . . . .	16
4.2	Exemplo de amostra de voz anecóica. . . . .	19
4.3	Exemplo de amostra de voz reverberante, convoluída com uma RIRO, onde $DRR = -4$ e $T60 = 1,38$ s. . . . .	20
4.4	Exemplo de amostra de voz reverberante, convoluída com uma RIRSM, onde $DRR = 10$ e $T60 = 0,50$ s. . . . .	20
4.5	Exemplo de amostra de voz em campo distante com $SNR = 14$ , representada pela voz reverberada mais os ruídos adicionados pelo segundo método de DA. . . . .	21

4.6	Exemplo de amostra de voz em campo distante com $SNR = 4$ , representada pela voz reverberada mais os ruídos adicionados pelo segundo método de DA. . . . .	22
6.1	RIR original do exemplo D1. . . . .	29
6.2	RIR simulada do exemplo D1. . . . .	29
6.3	RIR original do exemplo T2. . . . .	32
6.4	RIR simulada do exemplo T2. . . . .	32
6.5	Amostra de voz original no exemplo N2. . . . .	34
6.6	Amostra de voz reverberada com RIRSM no exemplo N2. . . . .	35
6.7	Amostra de voz em campo distante no exemplo N2. . . . .	35

# Lista de Tabelas

5.1	Descrição dos textos pronunciados por locutor. . . . .	24
5.2	Configurações de RIRs disponíveis na AIR. . . . .	25
5.3	Descrição dos tipos de ruídos pontuais usados da base MUSAN. . . .	26
5.4	Descrição dos tipos de ruídos de fundo usados da base MUSAN. . . .	26
6.1	Faixas dos parâmetros. . . . .	27
6.2	Exemplos de DA de DRR gerados. . . . .	28
6.3	Análise subjetiva de distância. . . . .	30
6.4	Exemplos de DA de T60 gerados. . . . .	31
6.5	Análise subjetiva de eco. . . . .	33
6.6	Exemplos de DA de AVCD gerados. . . . .	34
6.7	Análise subjetiva de nível de ruído. . . . .	36

# Capítulo 1

## Introdução

Neste capítulo, serão introduzidos os principais tópicos do projeto, além de mostrar sua relevância para o escopo da engenharia moderna e as metodologias que são usadas para alcançar seus objetivos. Ao final é descrita a estrutura organizacional do texto.

### 1.1 Tema

O tema do trabalho é sobre o estudo de uma forma de simular Respostas ao Impulso de Salas (RIR) com parametrizações diferentes a partir de amostras de RIR gravadas em ambientes reais, e ainda usar a RIR para gerar amostras de áudio em locais simulados a partir de gravações de voz reais.

### 1.2 Delimitação

O estudo é focado em inferir uma técnica de reforço de dados, tanto em amostras reais de RIR, quanto nas gravações de voz. Este trabalho está delimitado em apenas modificar as RIRs medidas, sem a utilização de programas de simulação acústica para gerar RIRs sintéticas.

### 1.3 Justificativa

Com o avanço das tecnologias de automação residencial, assistentes pessoais nos *smartphones* e comunicação *online*, o estudo de técnicas de processamento de

áudio (no caso específico deste trabalho, relacionados à voz), tornou-se mais relevante para a sociedade. Uma das características mais importantes a ser detectada no processamento de áudio é a Resposta ao Impulso de Salas, que representa o modelo acústico do ambiente, pois através desta é possível extrair informações pertinentes do local em que o áudio foi gravado e também detectar a posição de fontes sonoras e as isolar para reconhecimento. No âmbito da área de reconhecimento de voz, a fala reverberante, ou seja, o sinal de fala combinado com o modelo acústico do ambiente é um dos desafios encontrados para a detecção da voz, tornando a identificação da RIR de vital importância para o reconhecimento de fala [1].

Junto a isso, houve avanços no âmbito do aprendizado de máquina, fornecendo alternativas para os métodos tradicionais de processamento de áudio [2]. Modelos de arquitetura de redes neurais necessitam de um grande volume de dados para que sejam treinados e aprimorados. Um dos mais recentes desafios nessa área, é o fato das bases de RIR não serem extensas, conforme esclarecidas no artigo [3]. Realizar uma grande quantidade de gravações de áudio é uma tarefa de alto custo, tanto financeiro quanto de tempo, necessitando de equipamento especializado e diversos locais com características de modelo sonoro diferentes e pessoas diversas para gerar grande diversidade de amostras de voz.

## 1.4 Objetivos

O objetivo deste trabalho é implementar um algoritmo capaz de gerar amostras de RIR simuladas para diferentes ambientes a partir de uma RIR real e gerar um banco de dados de amostras de voz convoluídas com as RIR simuladas e com ruídos para uso em treinamento de redes neurais. Dessa forma, têm-se como objetivos específicos:

1. Propor um algoritmo que altere as características da RIR para simular diferentes ambientes com RIR diferentes;
2. Elaborar um algoritmo que faça o acréscimo de ruídos pontuais ou ruídos de fundo em uma amostra de voz;
3. Desenvolver um sistema computacional que aplique ambos os algoritmos an-



teriores em sequência para gerar amostras de voz em ambientes ruidosos.

## 1.5 Metodologia

Um sinal de voz gravado em um ambiente pode ser interpretado como a junção de três partes: uma amostra de voz pura, sem nenhum fator externo ou reverberação envolvido, convoluída com a RIR onde ocorre a gravação, somada a um sinal de ruído, podendo este ser pontual ou um ruído de ambiente. A RIR representa um modelo acústico do ambiente, que define como um receptor acústico irá receber caso o áudio seja gerado e percebido de dentro deste ambiente. Uma definição de RIR é a de uma função que registra a pressão sonora temporalmente em um ambiente fechado após ser excitada por um impulso.

Neste trabalho é proposta uma forma de gerar RIRs simuladas partindo de uma RIR real, ou seja, obtida de medições da resposta ao impulso em um ambiente fechado real, alterando suas propriedades. Reproduz-se o que foi proposto no artigo de *data augmentation* para respostas ao impulso usadas na estimação do modelo acústico [4]. São geradas RIRs simuladas, modificando-se as propriedades de Tempo de Decaimento (T60) e de razão entre áudio direto e reverberante (DRR). Através da alteração dessas duas propriedades, obtém-se, artificialmente, uma quantidade representativa de RIRs possíveis de serem utilizadas.

Para gerar as amostras de vozes reverberadas que compõe a base de dados, acompanha-se o que é proposto no artigo de estudo de *data augmentation* em vozes reverberadas [5], onde são convoluídos sinais de voz anecóicos com as RIRs simuladas geradas anteriormente. Além disso, são acrescentados a esse sinal de voz reverberado ruídos diversos, que são caracterizados de duas formas: ruídos pontuais e de ambiente. Os ruídos pontuais são amostras de áudio curtas que podem ser introduzidas em qualquer momento da fala. Os ruídos de ambiente são sons constantemente presentes ao fundo da gravação para simular um ambiente externo ruidoso. Ambos os tipos de ruídos foram extraídos da biblioteca MUSAN [6].

Através desses dois passos, são gerados vários sinais de voz reverberados artificialmente. A simulação da RIR tem por objetivo colocar a amostra de voz em vários ambientes fechados, a inclusão de ruídos ajudam drasticamente no treinamento de

redes neurais, impedindo que as redes fiquem super-treinadas em características muito específicas da fala durante o treinamento, uma vez que tendem a simular os fatores externos que podem estar envolvidos em uma gravação real.

## 1.6 Descrição

O Capítulo 2 apresenta uma breve análise sobre as principais aplicações do tema e os desafios que este trabalho auxilia na solução.

No Capítulo 3 será descrita a metodologia usada para fazer a *data augmentation* de uma RIR já existente.

No Capítulo 4 explica-se a metodologia usada para gerar sinais de voz acústicos a partir de RIRs simuladas anteriormente e da adição de ruídos pontuais ou de fundo.

O Capítulo 5 apresenta as bases de dados que serão usadas para gerar os resultados experimentais.

O Capítulo 6 é focado em exibir os resultados obtidos através dos métodos anteriores e demonstrar sua eficácia.

Por fim, o Capítulo 7 trata das conclusões que podem ser obtidas sobre este projeto, além de propor trabalhos futuros.

Este trabalho possui um apêndice. O Apêndice A contém o link do código fonte da aplicação do algoritmo proposto.

## Capítulo 2

# Análise de Fontes Sonoras e seus Desafios

Este capítulo é dedicado à introdução do leitor ao principal tópico de estudo do projeto e assim mostrar algumas aplicações onde este é usado, além de apresentar os desafios relacionados a estas aplicações.

### 2.1 Resposta ao Impulso de Ambiente Acústico e suas Aplicações

Dentre os diversos tópicos na grande área de estudo de sinais de áudio, destacam-se a detecção e o reconhecimento de fontes acústicas no espaço físico. Um caso frequentemente encontrado em aplicações referentes a este tópico é o de sinais de voz gravados em ambientes fechados, onde um ou mais microfones são posicionados na sala afastados da fonte sonora. Estes sinais são corrompidos pela reverberação do ambiente, que surge a partir da sobreposição da onda sonora anecóica que chega ao microfone com as ondas sonoras atenuadas e refletidas nas paredes do ambiente fechado.

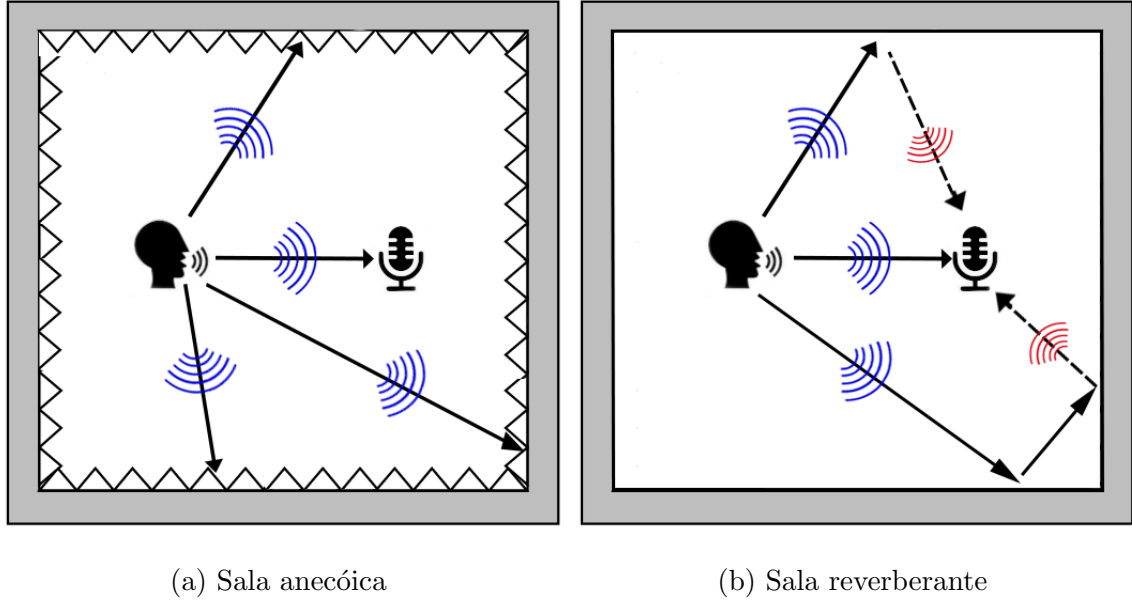


Figura 2.1: Representação de uma sala anecóica e reverberante.

Observa-se na Figura 2.1a uma representação de uma sala anecóica, onde o único áudio capturado pelo microfone é a onda sonora direta enviada pela fonte, sem nenhuma reflexão do ambiente. Na sala reverberante, representada pela Figura 2.1b, nota-se que o áudio capturado será uma combinação da onda sonora direta com as refletidas nas paredes. Este sinal reverberado pode ser modelado da seguinte forma:

$$y(t) = s(t) * h(t) + n(t) \quad (2.1)$$

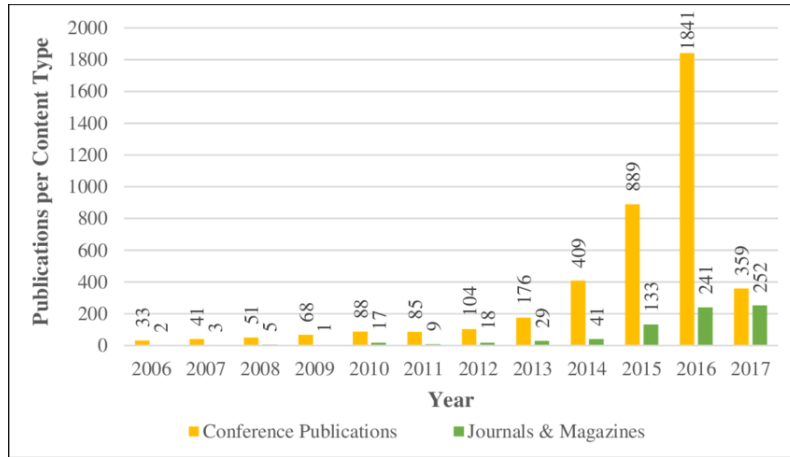
onde  $y(t)$  representa o sinal de voz em campo distante,  $s(t)$  o sinal de voz anecóico,  $h(t)$  a RIR e  $n(t)$  o sinal de ruído que pode estar presente no ambiente. Dessa forma, é possível inferir que a RIR representa o modelo acústico de uma sala, para uma determinada combinação de fatores do ambiente, incluindo: temperatura e umidade relativa do ar, pressão atmosférica, material das paredes e posicionamento de móveis. Reverberação causa degradação do sinal de voz, levando à perda de clareza na comunicação [7] e à redução de desempenho de sistemas de reconhecimento de voz [8]. Este problema demonstra a necessidade de identificar dinamicamente o modelo acústico do ambiente para que possam ser mitigadas as distorções nas amostras de voz gravadas e assim facilitar os algoritmos que usam esses sinais.

Este projeto é focado no estudo de uma forma de gerar RIRs simuladas

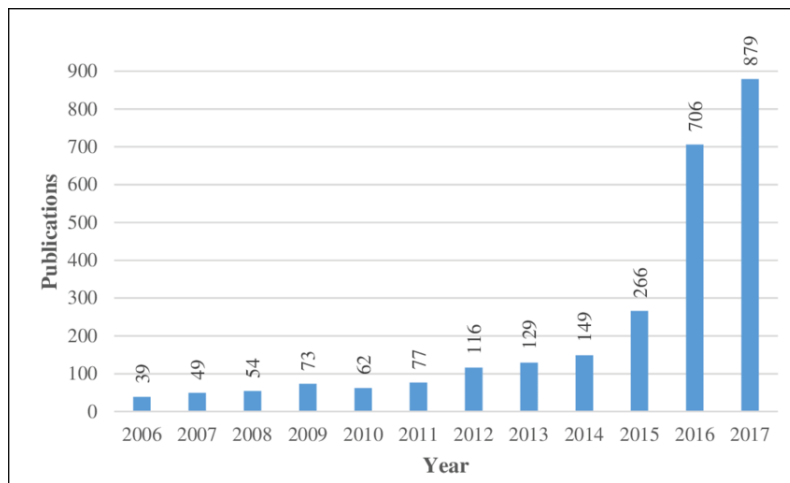
a partir de RIRs reais devido à sua importância para diversas aplicações usadas atualmente na indústria. Uma de suas aplicações é na análise e desenvolvimento de algoritmos de reconhecimento de voz robusta [8], onde é necessário inferir a RIR para que possa ser feita a comparação entre a supressão da reverberação ideal com a cega. Outra aplicação da RIR é para desenvolvimento de algoritmos para localização e separação de fontes sonoras [9], onde as RIRs são usadas no auxílio do mapeamento acústico de ambientes reverberante através de algoritmos de separação de fonte às cegas.

## 2.2 Desafios correlacionados à RIR

É possível notar um recente aumento de pesquisas relacionadas à área de aprendizado de máquina no meio científico, especialmente no tópico de *Deep Learning*) [10]. Observando a Figura 2.2 [11], que apresenta quantidade de artigos publicados por ano relacionados com Deep Learning, nota-se que após 2015, houve um aumento considerável de publicações em conferências do IEEE e em artigos e livros publicados pela editora Springer®. Muitas dessas publicações são dedicadas para áreas de pesquisa relacionadas com áudio [10, 12, 13]; de acordo com o artigo [14], aproximadamente 20% das publicações são voltadas para o tópico de reconhecimento de voz usando técnicas de *Deep Learning* em suas metodologias.



(a) Publicações por ano - IEEE



(b) Publicações por ano - Springer®

Figura 2.2: Gráficos de quantidade de artigos publicados por ano relacionados com *Deep Learning*.

Um dos maiores desafios enfrentados ao utilizar técnicas de *Deep Learning* é de obter uma grande quantidade de dados para treinamento. Podem ser observados exemplos em [3, 15], onde os autores necessitaram agrupar dados de mais de 5 bases contendo RIRs para que fosse possível treinar e avaliar suas redes. No caso de bases de dados que envolvem RIRs, o motivo de não existir uma alta variedade de dados é devido às dificuldades de realizar gravações dos áudios [16]. A RIR pode ser obtida com diversos tipos de sinal de excitação: varreduras de seno, sequencias pseudo aleatórias. De forma aproximada, pode-se usar estouro de balões [17].

Não menos importante, para aumentar a quantidade de amostras na base, deve-se gravar o áudio não só em diferentes posições no ambiente e distâncias fonte-

microfone, como também realizar estes mesmos procedimentos em ambientes diferentes, o que requer o transporte de diversos equipamentos especializados entre localizações físicas.

## 2.3 *Data Augmentation*

*Data Augmentation* (DA) representa um conjunto de técnicas que são usadas em dados já existentes com o intuito de gerar versões modificadas que aumentam a representatividade dos dados disponíveis para uma determinada aplicação. No contexto de *Deep Learning*, essas técnicas tornam-se vitais para incrementar artificialmente bases de dados para treinamento que não possuem uma alta variedade de amostras e isso inclui dados de áudio [18, 19].

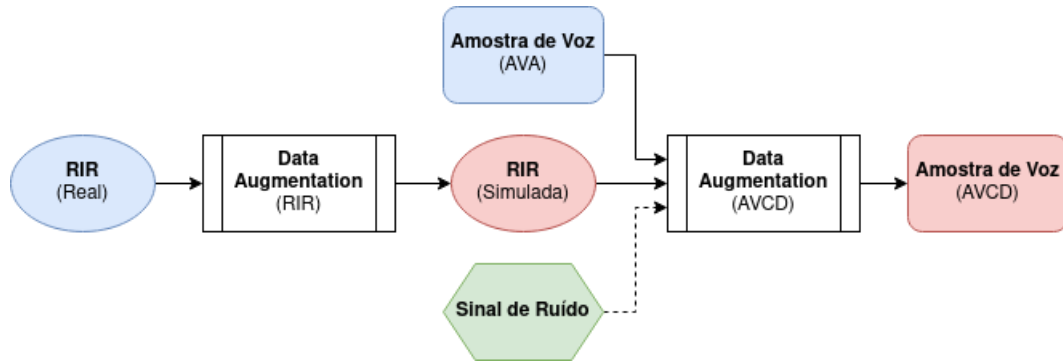


Figura 2.3: Fluxo geral de procedimentos para gerar sinais de voz reverberantes.

No escopo deste trabalho, de acordo com a Figura 2.3 usaremos duas técnicas de *Data Augmentation* para gerar amostras de voz reverberantes. Uma das técnicas é voltada para simulação de RIRs, que altera suas propriedades para que possam ser simuladas diferentes condições e posições em um determinado ambiente.

Outra técnica é desenvolvida para gerar amostras de voz em campo distante, usando as RIRs simuladas geradas pela primeira técnica e sinais de ruído. Na segunda técnica, serão utilizados dois tipos de ruídos: pontual e de fundo. O ruído pontual representa uma fonte sonora, além da amostra de voz anecóica, que está localizada no mesmo ambiente fechado em que ocorre a geração da AVCD, mas que não faz parte do sinal que se deseja detectar. O ruído de fundo representa um fator sonoro externo ao ambiente fechado acrescentado à AVCD que prejudica a detecção da amostra de voz original em uma AVCD.

## Capítulo 3

# *Data Augmentation* aplicado à Resposta ao Impulso de Ambientes

Este capítulo é dedicado à implementação do primeiro algoritmo de *Data Augmentation*, onde são geradas RIRs simuladas (RIRSM) a partir de RIRs originais (RIRO), que foram gravadas em ambientes diversos e compõem o banco de dados AIR [20]. São observados os parâmetros de razão Direto-Reverberante (DRR) e tempo de reverberação (T60), que são inferidos com base em uma RIRO e que serão manipulados pelo algoritmo para gerar RIRs que vão representar modelos acústicos diferentes.

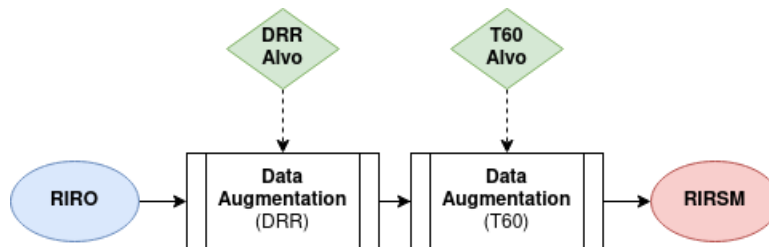


Figura 3.1: Fluxo de procedimentos para gerar a RIRSM.

Este trabalho é uma implementação dos passos demonstrados no artigo [4]. A Figura 3.1 especifica o fluxo de procedimentos implantados por este algoritmo, onde a DRR e T60 alvos são os valores escolhidos pelo usuário ou sorteados aleatoriamente, que determinam as características da RIRSM.

Antes de explicar os métodos usados, é necessário definir duas funções [4]:



$$h_e(t) = \begin{cases} h(t), & t_d - t_0 \leq t \leq t_d + t_0 \\ 0, & \text{caso contrário,} \end{cases} \quad (3.1)$$

$$h_l(t) = \begin{cases} h(t), & t < t_d - t_0 \\ h(t), & t > t_d + t_0 \\ 0, & \text{caso contrário,} \end{cases} \quad (3.2)$$

onde  $t$  representa o tempo discreto,  $t_d$  o tempo que as ondas sonoras diretas, ou seja, sem reflexão, levam da fonte até o destino de gravação,  $t_0$  a janela de tolerância, neste caso definida com o valor 2,5 ms [4],  $h(t)$  uma RIR,  $h_e(t)$  as primeiras reflexões e  $h_l(t)$  as reflexões tardias. Neste algoritmo,  $t_d$  é determinado da forma:

$$\begin{cases} t_d = t_{max}, \\ t_{max}, \text{ onde } h(t_{max}) = \max(|h(t)|) \end{cases} . \quad (3.3)$$

### 3.1 Razão Direto-Reverberante (DRR)

A DRR representa a razão entre a energia sonora da resposta ao impulso que atinge o alvo diretamente e a energia reverberante, ou seja, que é refletida pelas paredes do ambiente fechado. Este parâmetro é calculado pela seguinte expressão:

$$DRR_{dB} = 10 \log_{10} \left( \frac{\sum_t h_e^2(t)}{\sum_t h_l^2(t)} \right). \quad (3.4)$$

Para obter a  $DRR_{dB}$  alvo desejada, aplica-se um fator de ganho escalar  $\alpha$  na função das primeiras reflexões  $h_e(t)$ . De acordo com [4], para evitar descontinuidades durante o cálculo do fator  $\alpha$ , reescreve-se  $h_e(t)$  em duas parcelas, uma que representa a janela direta no pico de intensidade de  $h(t)$  e outra que representa uma janela de resíduo de  $h_e(t)$  formando, assim,

$$h'_e(t) = \alpha w_d(t) h_e(t) + [1 - w_d(t)] h_e(t), \quad (3.5)$$

onde  $w_d(t)$  representa uma janela de Hann de duração de 5 ms, considerando-se uma janela de tolerância  $t_0 = 2,5$  ms. Substituindo na Equação (3.4)  $h_e(t)$  por  $h'_e(t)$  e combinando com a Equação (3.5), obtém-se a seguinte equação quadrática;

$$\begin{aligned}
& \alpha^2 \sum_t w_d^2(t) h_e^2(t) + 2\alpha \sum_t [1 - w_d(t)] w_d(t) h_e^2(t) + \\
& \sum_t [1 - w_d(t)]^2 h_e^2(t) - 10^{DRR_{dB}/10} \sum_t h_l^2(t) = 0,
\end{aligned} \tag{3.6}$$

O parâmetro  $\alpha$  desejado será a raiz de maior valor. Uma ressalva deste procedimento é que se deve atentar para não escolher uma  $DRR_{dB}$  que não seja muito menor que a original, pois após a transformação de  $h_e(t)$  para  $h'_e(t)$ , dependendo do valor de  $\alpha$ , é possível incidir em um caso onde  $\max(h'_e(t)) < \max(h_l(t))$ , tornando a RIRSM impraticável.

A Figura 3.2 exibe um exemplo de sinal  $h(t)$  com o seu  $h_e(t)$  após ser feita a modificação do DRR. Comparando  $h_e(t)$  com  $h'_e(t)$ , de acordo com a Figura 3.3, nota-se que  $h'_e(t)$  está com uma intensidade maior, o que condiz a modificação do DRR de -4,5 dB para 4 dB neste exemplo.

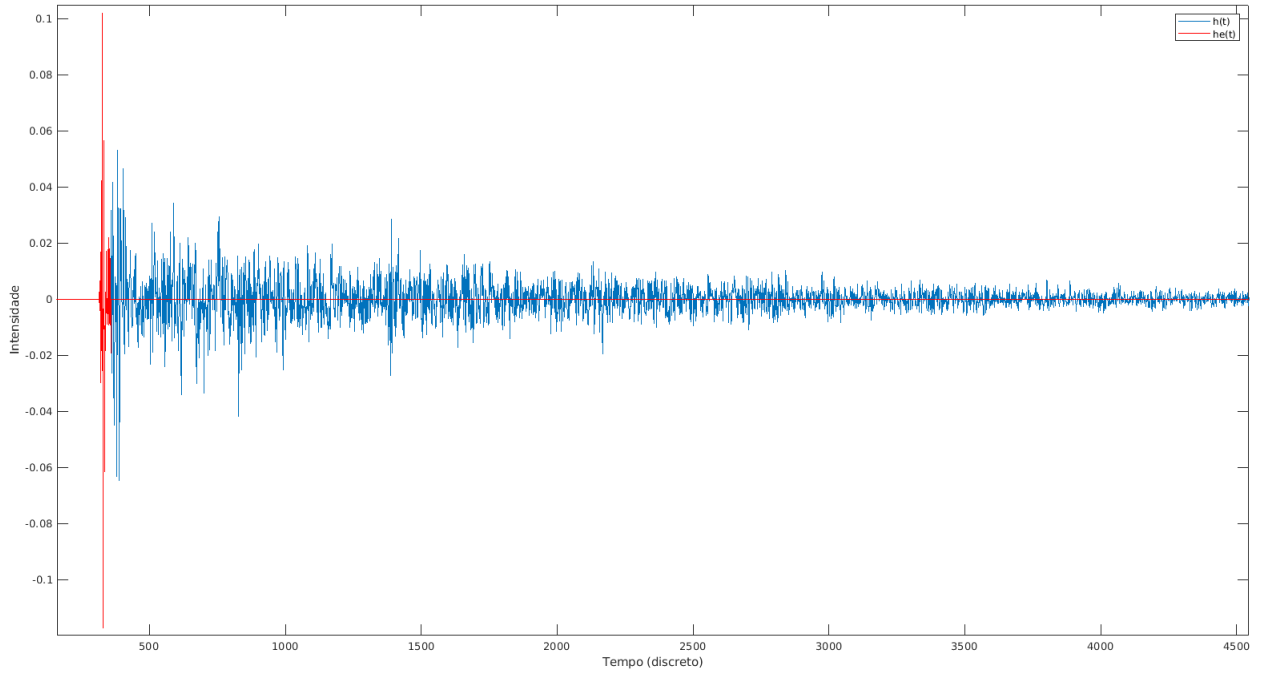


Figura 3.2: Um exemplo de  $h(t)$  com  $h_e(t)$ , onde é feita a modificação do DRR de -4,5 para 4, marcado em vermelho.

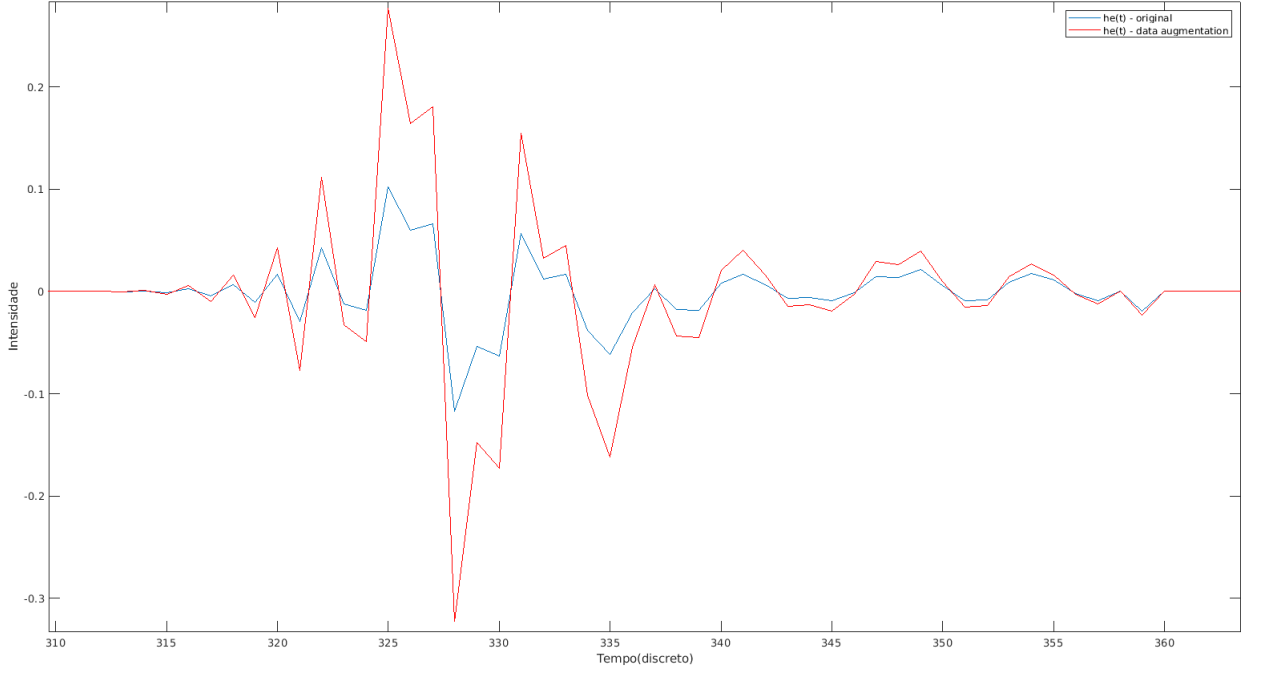


Figura 3.3: Um exemplo de  $h_e(t)$  antes e depois da aplicação do algoritmo, original em azul ( $DRR = -4,5$  dB) e o modificado em vermelho ( $DRR = 4$  dB).

## 3.2 Tempo de Reverberação (T60)

O T60, definido na equação 3.7, representa a duração de tempo que leva para a energia sonora da RIR no alvo decair 60 dB comparado à sua intensidade máxima. Geralmente, devido à dificuldade de se medir uma queda de 60 dB, o parâmetro medido é o T20 ou o T30 e depois multiplica-se os seus valores por 3 e 2, respectivamente, para obter o T60, assumindo-se um decaimento exponencial da envoltória da RIR.

$$\begin{cases} T60 = t_f - t_i, \text{ onde} \\ t_i \rightarrow h(t_i) = \max(h(t)) \\ t_f \rightarrow 10 \log_{10}(h^2(t_i) - h^2(t_f)) = 60\text{dB} \end{cases} \quad (3.7)$$

Para realizar modificações na RIR, é necessário modelar a função das reflexões tardias. De acordo com [4], um modelo normalmente usado é de um ruído gaussiano exponencialmente decadente, acrescido de um ruído residual, ou seja,

$$h_m(t) = Ae^{-(t-t_o)/\tau}n(t)u(t-t_o) + \sigma n(t), \quad (3.8)$$

onde  $A$  representa o ganho da resposta ao impulso,  $\tau$  a taxa de decaimento,  $\sigma_m$  o desvio padrão do ruído residual,  $n(t)$  um ruído gaussiano padrão (média nula e desvio padrão unitário),  $t_o$  o valor temporal onde  $h_l(t)$  tem o seu primeiro valor não nulo e  $u(t)$  um degrau unitário. Neste trabalho, diferente da implementação do algoritmo em [4], é considerada apenas a taxa de decaimento do espectro de frequência por completo da RIR, ao invés de dividi-la em subbandas e analisar a taxa de decaimento para cada faixa de frequência.

Os parâmetros  $A$ ,  $\tau$  e  $\sigma$  são estimados de acordo com o padrão definido na ISO 3382-1 [21]. Seja  $T60_d$  o valor de T60 alvo para DA, é possível inferir a taxa de decaimento através da equação

$$T60_d = \ln(1000)\tau_d T_s, \quad (3.9)$$

onde  $\tau_d$  representa a taxa de decaimento alvo e  $T_s$  o intervalo de amostragem. A DA do tempo de reverberação é feita multiplicando-se  $h_l(t)$  pela exponencial

$$h'_l(t) = h_l(t)e^{-(t-t_o)\frac{\tau-\tau_d}{\tau\tau_d}}. \quad (3.10)$$

Por fim, a RIRSM completa,  $h'(t)$ , pode ser representada pela equação

$$h'(t) = h'_e(t) + h'_l(t). \quad (3.11)$$

A Figura 3.4 exibe um exemplo de sinal  $h(t)$  com o seu  $h_l(t)$  após ser feita a modificação do T60. Comparando  $h_l(t)$  com  $h'_l(t)$ , de acordo com a Figura 3.5, nota-se que  $h'_l(t)$  está com uma intensidade maior, o que condiz a modificação do T60 de 1,38 para 2,6 segundos neste exemplo.

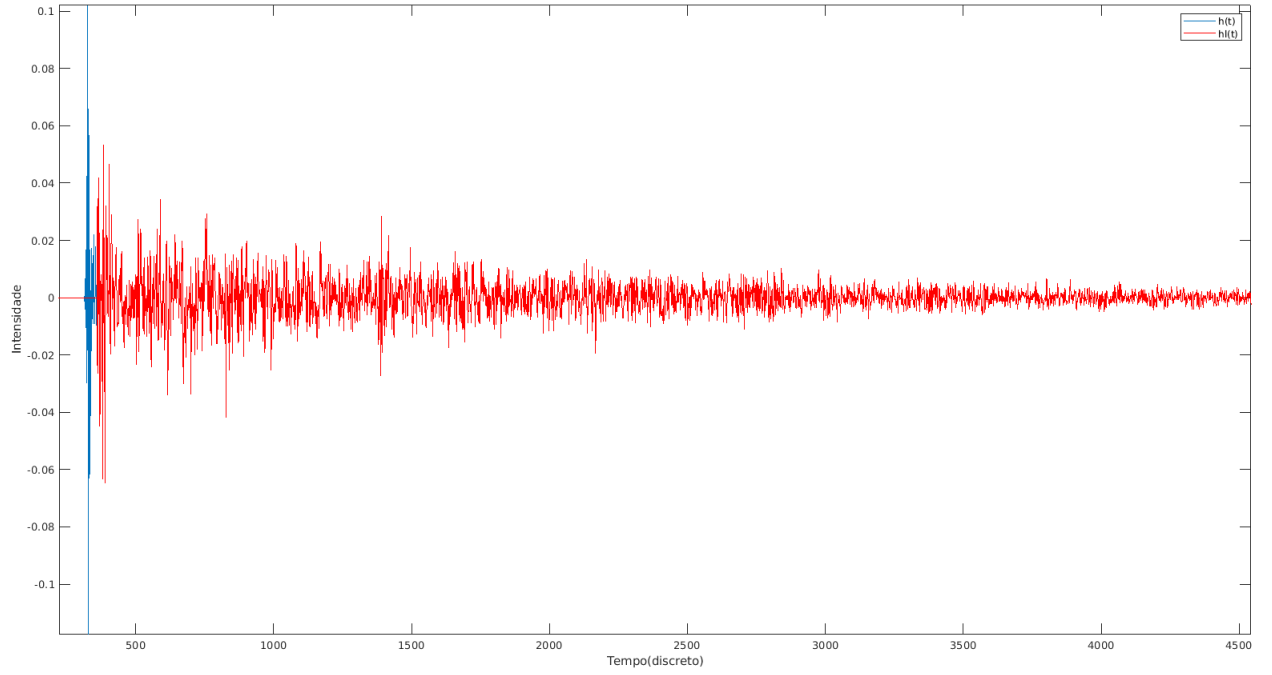


Figura 3.4: Um exemplo de  $h(t)$  com  $h_l(t)$ , onde é feita a modificação do T60 de 1,38 para 2,6 segundos, marcado em vermelho.

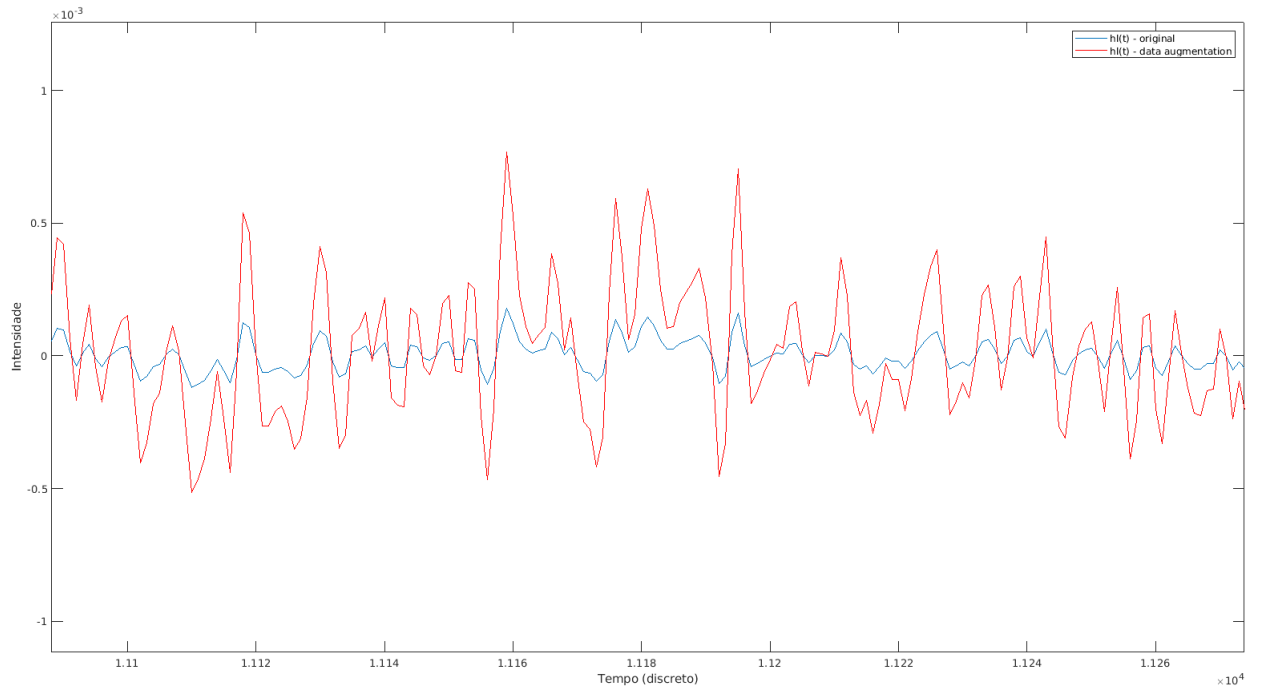


Figura 3.5: Um exemplo de um trecho ao final de  $h_l(t)$  antes e depois da aplicação do algoritmo, original em azul ( $T_{60} = 1,38$  s) e o modificado em vermelho ( $T_{60} = 2,6$  s).

## Capítulo 4

# Sinais de Voz com RIRs Simuladas e com Ruídos

Este capítulo é dedicado ao desenvolvimento do segundo algoritmo de *Data Augmentation*, onde são geradas as amostras de voz reverberadas em campo-distante (AVCDs) usando: amostras de voz anecóicas (AVAs), RIRSM, sinais de ruído pontuais (SRPs) e de fundo (SRFs).

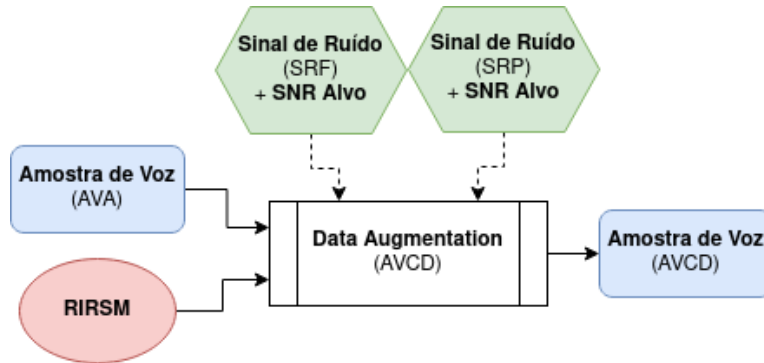


Figura 4.1: Fluxo de procedimentos para gerar a AVCD.

Este trabalho é uma implementação do método de DA proposto no artigo [5]. A Figura 4.1 especifica o fluxo de procedimentos implantados por este algoritmo, onde o SRP, SRF e SNR alvos são aleatoriamente escolhidos, dentro de uma base de dados de ruídos e uma faixa de valores definidas pelo usuário, que determinam as características da AVCD.

## 4.1 Simulação de fala em campo distante

Sinais de voz em campo-distante são tipicamente compostos por uma combinação de VR, SRP (assumindo que a fonte do ruído pontual encontra-se no mesmo ambiente da VR) e SRF (assumindo que o ruído de fundo não é afetado pelo modelo acústico do ambiente). É possível modelar uma AVCD conforme a equação

$$s_{cd}[t] = s_a[t] * h[t] + \sum_i n_{pi}[t] * h[t] + n_f[t], \quad (4.1)$$

onde  $s_{cd}[t]$  representa a AVCD,  $s_a[t]$  a AVA,  $h[t]$  a RIR,  $n_{pi}[t]$  o  $i$ -ésimo SRP e  $n_f[t]$  o SRF. Neste trabalho, diferente da implementação do algoritmo em [5], é considerado apenas uma única RIR para gerar a AVCD, ou seja, os ruídos pontuais são convoluídos com a mesma RIR que é usada para a fonte de voz.

No Algoritmo 1 descreve-se o procedimento que é usado para gerar sinais de voz em campo-distante simulados.

**Input:**  $fl_p$  : Flag de inclusão de ruído pontual

**Input:**  $fl_g$  : Flag de inclusão de ruído de fundo

**Input:**  $m$  : Quantidade de ruídos pontuais

**Input:**  $SNR_{up}$  : Limite superior de SNR

**Input:**  $SNR_{dw}$  : Limite inferior de SNR

$s_r[t] \leftarrow s_a[t] * h[t]$  : Convolução da RIR com AVA

**if**  $fl_p = true$  **then**

**for**  $i = 1$  até  $m$  **do**

Escolha aleatória de um ruído pontual  $n_{pi}[t]$  da biblioteca de ruído.

Escolha aleatória de uma SNR Alvo  $SNR_t$  compreendida dentro do intervalo  $[SNR_{dw}, SNR_{up}]$ .

Dedução do fator  $\alpha$  a partir da  $SNR_t$  para corrigir a intensidade de  $n_{pi}[t]$ .

Escolha aleatória de offset  $o_t$  compreendida dentro do intervalo  $(0, duração(t))$ .

$s_r[t] \leftarrow s_r[t] + \alpha \text{offset}(n_{pi}[t] * h[t], o_t)$  : Adição de SRP na AVR.

**end**

**end**

**if**  $fl_g = true$  **then**

Escolha aleatória de um ruído de fundo  $n_f[t]$  da biblioteca de ruído.

Escolha aleatória de uma SNR Alvo  $SNR_t$  compreendida dentro do intervalo  $[SNR_{dw}, SNR_{up}]$ .

Dedução do fator  $\alpha$  a partir da  $SNR_t$  para corrigir a intensidade de  $n_f[t]$ .

Estender ou encurtar  $n_f[t]$  até que duração  $(n_f[t]) = duração(s_r[t])$

$s_r[t] \leftarrow s_r[t] + \alpha n_f[t]$  : Adição de SRF na AVR.

**end**

**Algoritmo 1:** Procedimentos para gerar AVCD



Neste trabalho, o algoritmo de geração de AVCD usa as RIRSMs geradas através do primeiro algoritmo, diferente do que foi implantado em [5], onde foram geradas RIRs de forma completamente digital [22], ou seja, sem usar RIRs reais como base para *Data Augmentation*. Nota-se também que o algoritmo permite habilitar ou não o uso de cada tipo de ruído para que possa aumentar a variedade de dados gerados, além de acomodar mais propósitos de treinamentos de *Deep Learning*.

A Figura 4.2 exibe uma amostra de voz anecóica que foi usada para gerar os próximos exemplos de aplicação do algoritmo. É feita uma comparação entre a convolução da AVA com a RIRO e com a RIRSM, representadas nas Figuras 4.3 e 4.4, respectivamente. Pode-se notar que a convolução com a RIRSM causa menos modificações no formato do sinal de voz original comparado à RIRO, uma vez que a primeira DA foi configurada para simular um ambiente menos reverberante, partindo dos parâmetros  $DRR = -4$  e  $T60 = 1,38$  s na RIRO para  $DRR = 10$  e  $T60 = 0,50$  s na RIRSM.

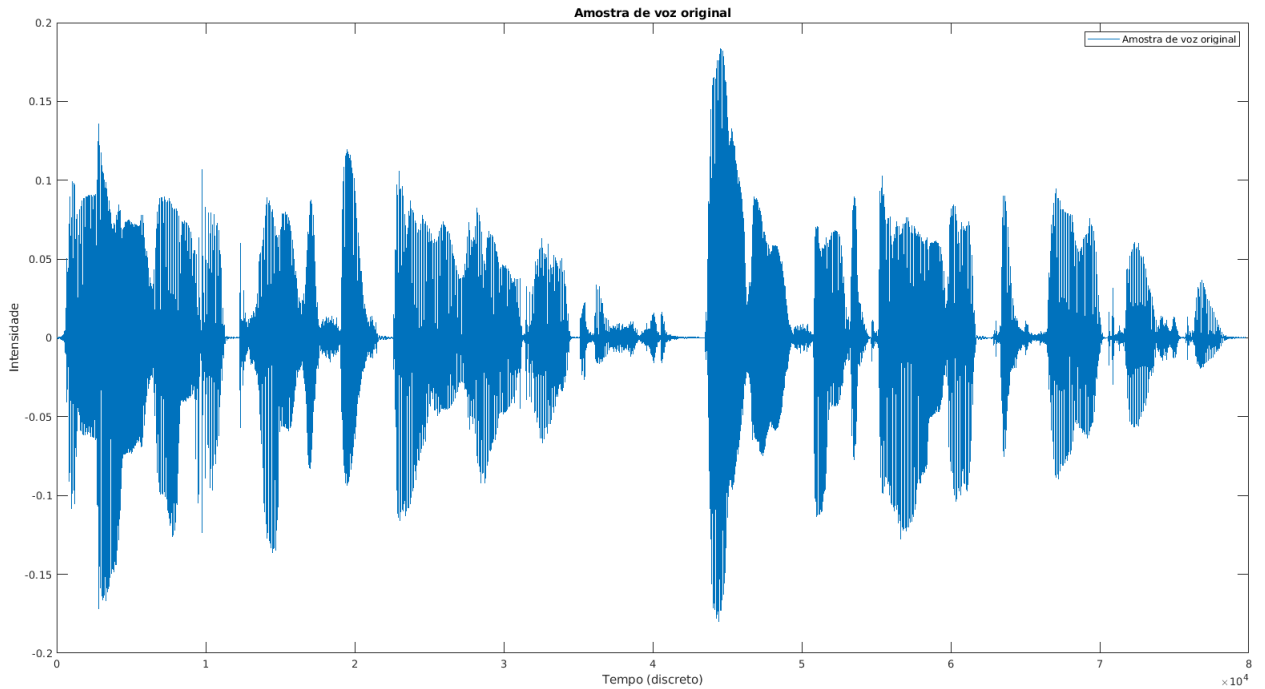


Figura 4.2: Exemplo de amostra de voz anecóica.

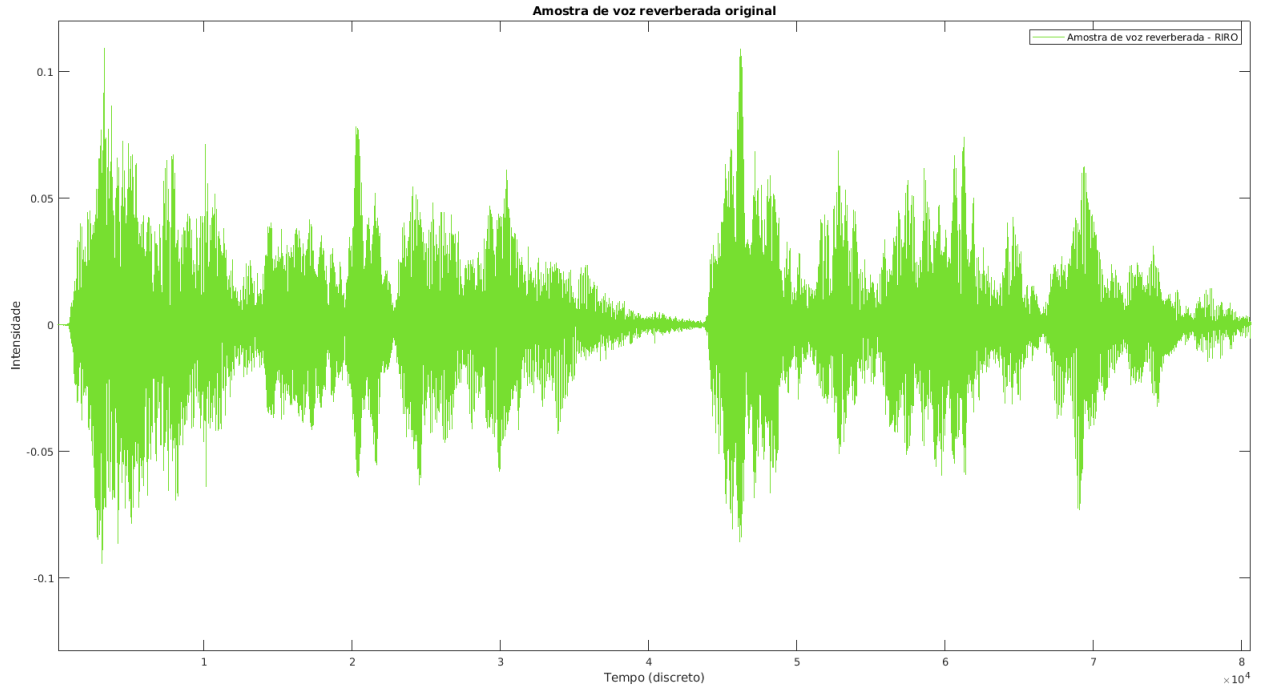


Figura 4.3: Exemplo de amostra de voz reverberante, convoluída com uma RIRO, onde  $DRR = -4$  e  $T60 = 1,38$  s.

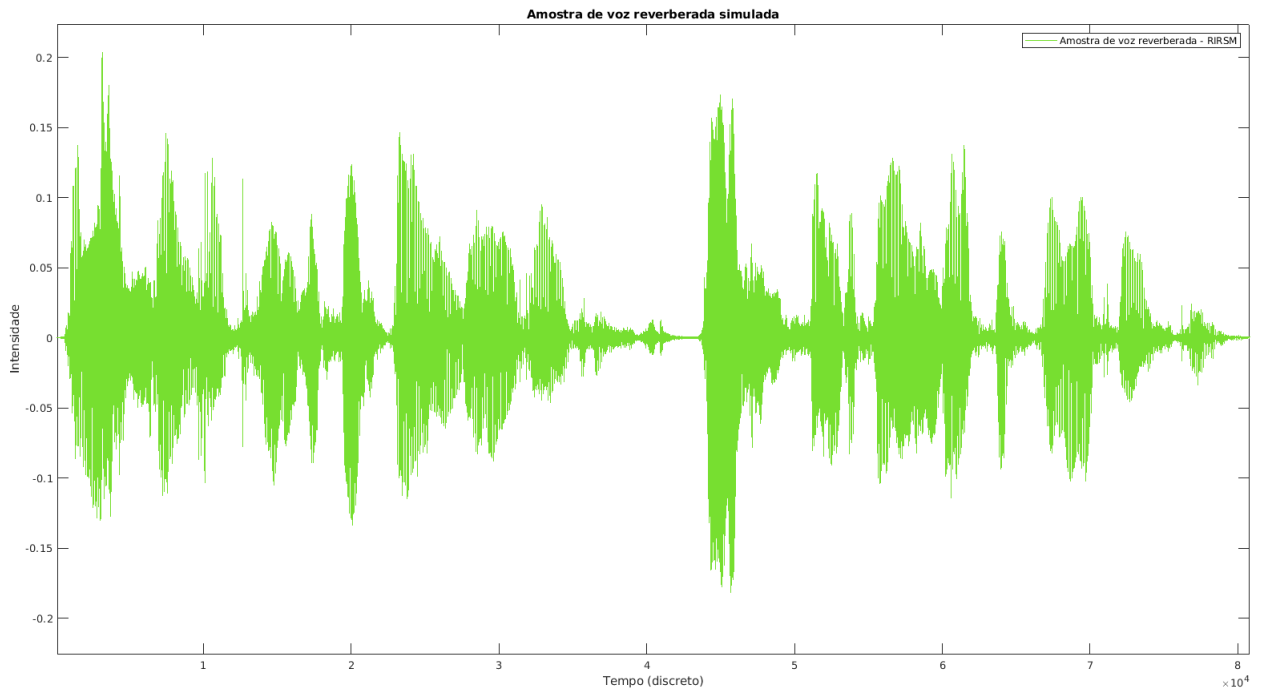


Figura 4.4: Exemplo de amostra de voz reverberante, convoluída com uma RIRSM, onde  $DRR = 10$  e  $T60 = 0,50$  s.

A partir da RIRSM usada nesta aplicação, foram gerados dois exemplos de

AVCDs, representadas nas Figuras 4.5 e 4.6, esta com  $SNR = 4$  e aquela com  $SNR = 14$ . Conforme o esperado, observa-se que a AVCD com o menor SNR possui ruídos claramente mais acentuados comparada à AVCD com o maior SNR.

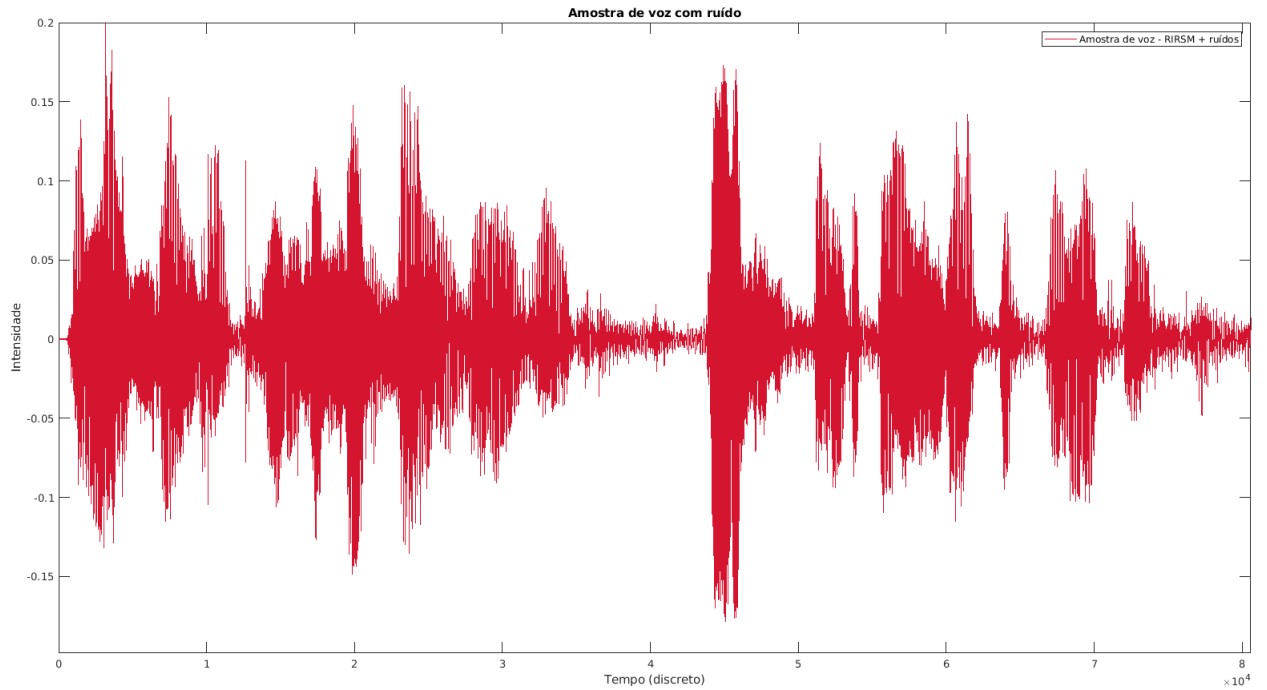


Figura 4.5: Exemplo de amostra de voz em campo distante com  $SNR = 14$ , representada pela voz reverberada mais os ruídos adicionados pelo segundo método de DA.

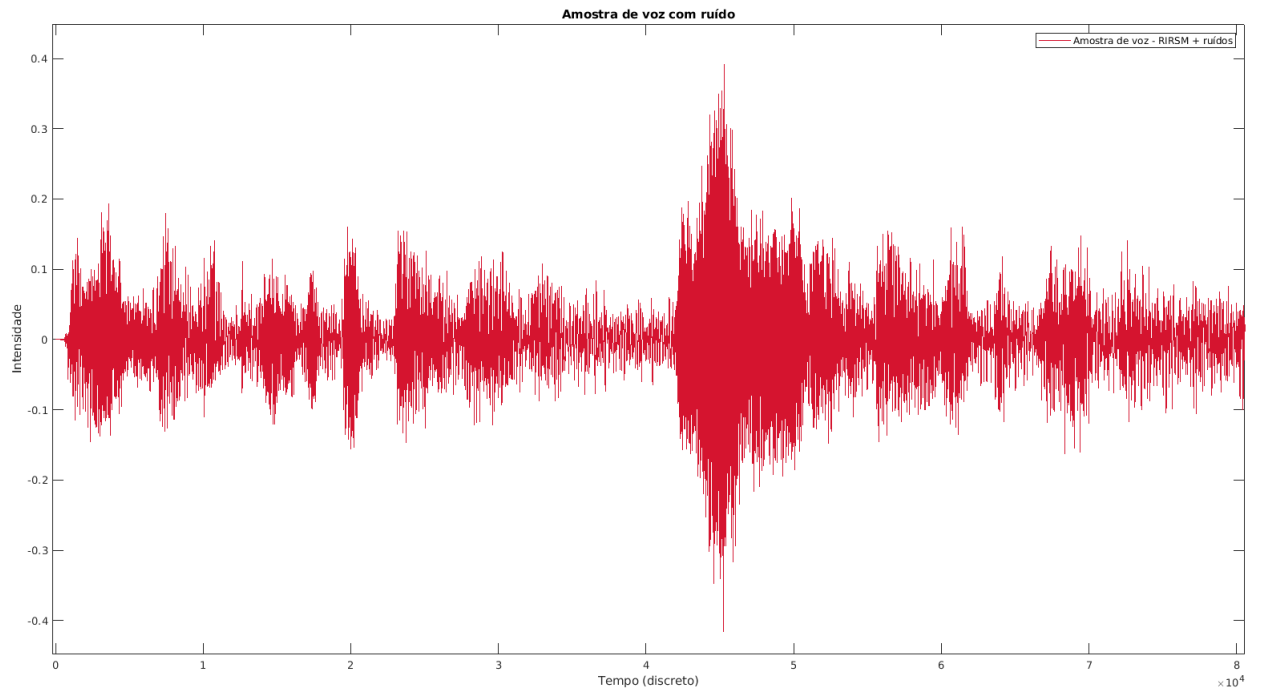


Figura 4.6: Exemplo de amostra de voz em campo distante com  $SNR = 4$ , representada pela voz reverberada mais os ruídos adicionados pelo segundo método de DA.

# Capítulo 5

## Bases de Dados

Este capítulo apresenta as bases de dados que serão usadas para gerar os resultados experimentais. Para a aplicação da metodologia deste trabalho, é necessário três fontes de dados:

- base de dados com amostras de voz anecóicas para convolução com as RIRSMs;
- base de dados com RIRs reais para o primeiro nível de *Data Augmentation*, gerando as RIRSMs;
- base de dados com SRPs e SRFs para o segundo nível de *Data Augmentation*, gerando as AVCDs.

### 5.1 Base de amostras de voz anecóicas

A base de AVAs usada consiste na leitura de textos em inglês por 4 pessoas diferentes (duas vozes masculinas e duas femininas) gravadas em uma câmara anecóica para avaliação de algoritmos de cancelamento de eco acústico. O eco foi adicionado artificialmente, de modo a efetuar experimentos controlados e com diferentes tempos de reverberação. Os arquivos de áudio são disponibilizados em formato WAV, com frequência de amostragem de 16 KHz e cada gravação tem duração em torno de 5 a 6 segundos. No caso deste trabalho, foram concatenadas duas frases por pessoa na mesma gravação devido ao tempo de duração dos ruídos pontuais, que serão adicionados para a geração de AVCDs.

A Tabela 5.1 descreve as gravações usadas neste projeto.

Tabela 5.1: Descrição dos textos pronunciados por locutor.

Nome	Código	Texto
Homen 1 - Texto 1	H1-T1	<i>This food is too spicy he complained. Young man can be very arrogant and rude.</i>
Homen 1 - Texto 2	H1-T2	<i>So Marcus owned a big shipping company. Their eyes met across the table.</i>
Homen 2 - Texto 1	H2-T1	<i>Time is running out for the scientists. If you knew Julie like I know Julie.</i>
Homen 2 - Texto 2	H2-T2	<i>Your new dress is breathtaking darling. Her first book was published last year.</i>
Mulher 1 - Texto 1	M1-T1	<i>Among them are canvases from a young artist. Building from the ground up is very costly.</i>
Mulher 1 - Texto 2	M1-T2	<i>Next year we will see several more exhibitions. The number of works on view will increase.</i>
Mulher 2 - Texto 1	M2-T1	<i>An enourmous quake rocked the island. Eventually he hopes to solve all the problems.</i>
Mulher 2 - Texto 2	M2-T2	<i>Eventually he hopes to solve all the problems. Faulty installation can be blamed for this.</i>

## 5.2 Base de RIRs - Aachen Impulse Response database (AIR)

A base de AIR [20] é um conjunto de respostas ao impulso sonoras que foram medidas em diversas salas por pesquisadores do Instituto de Acústica (ITA) da Universidade RWTH, Aachen, Alemanha. O objetivo dessa base é fornecer dados para estudos de algoritmos de processamento de sinais para ambientes reverberantes.

Ela é composta, primariamente, por RIRs binaurais medidas com uma cabeça de manequim em locais com diferentes propriedades acústicas. É importante frisar também que a base possui gravações com diferentes distâncias entre a fonte sonora e os microfones para a mesma sala, gerando funções com diferentes valores de DRR. A base também possui gravações em diferentes direções com o objetivo de auxiliar

algoritmos de detecção de direção da fonte sonora. Para o escopo deste projeto, tais RIRs não serão usadas.

Tabela 5.2: Configurações de RIRs disponíveis na AIR.

Sala	Descrição	Canais	Cabeça	Distâncias (m)
Booth	cabine de estúdio	E/D	S/N	0,5/1/1,5
Office	escritório comercial	E/D	S/N	1/2/3
Meeting	sala de reuniões	E/D	S/N	1,45/1,7/1,9/2,25/2,8
Lecture	sala de aula	E/D	S/N	2,25/4/5,56/7,1/8,68/10,2
Stairway	escadaria aberta	E/D	S/N	1/2/3
Aula Carolina	igreja de área 570m <sup>2</sup>	E/D	S/N	1/2/3/5/15/20

Os ambientes em que foram feitas as gravações de RIRs e suas respectivas configurações são definidos na Tabela 5.2. Todos os ambientes usados possuem gravações com ambos os canais esquerdo e direito (E/D), com configuração com ou sem a cabeça falsa (S/N) e para diferentes distâncias (em metros) entre a fonte que gera o impulso sonoro e o microfone. As RIRs foram salvas como vetores binários de precisão dupla de ponto flutuante (formato MAT, que pode ser importado via MATLAB®).

### 5.3 Base de ruídos - MUSAN

A base de MUSAN (*A Music, Speech, and Noise Corpus*) [6] consiste em um conjunto de músicas de diversos gêneros, amostras de voz de doze línguas e uma variedade de ruídos técnicos, por exemplo tons DTMF e sons de uma máquina de fax, e ruídos de ambiente, por exemplo sons de animais e chuva. Ela foi criada primariamente para auxiliar no treinamento de modelos voltados para detecção de atividade de voz. Contudo, ela é usada também para teste de algoritmos processamento de sinais na área de áudio, por exemplo, de reconhecimento de voz e orador. Uma das vantagens dessa base é o fato dela ser uma compilação de áudios com fontes em domínios públicos, facilitando a distribuição dos áudios para uso da comunidade científica.

No escopo deste projeto, será usada somente a seção de ruídos da base, que

contém seis horas de áudio no total. A seção de ruídos é composta por sons de curta duração, que são usados como SRPs no segundo processo de *Data Augmentation*, e por sons de ambiente, usados como SRFs no mesmo processo.

Tabela 5.3: Descrição dos tipos de ruídos pontuais usados da base MUSAN.

Código	Descrição
RP-1	miado de gato
RP-2	madeira sendo lixada
RP-3	buzina de automóvel
RP-4	porta abrindo
RP-5	grampeador
RP-6	teclado de forno de microondas
RP-7	<i>zipper</i> sendo fechado
RP-8	latido de cão
RP-9	batendo em uma porta
RP-10	espirro
RP-11	campainha
RP-12	vibrador de celular

Tabela 5.4: Descrição dos tipos de ruídos de fundo usados da base MUSAN.

Código	Descrição
RF-1	avião decolando em aeroporto
RF-2	sala de máquinas
RF-3	estática
RF-4	sons de floresta

As Tabelas 5.3 e 5.4 indicam os ruídos separados da base para gerar AVCDs. Os arquivos de áudio são disponibilizados em formato WAV, com frequência de amostragem de 16 KHz.



# Capítulo 6

## Resultados Experimentais

Este capítulo descreve os procedimentos e as configurações feitas para a implementação empírica da metodologia deste trabalho. Cada seção exibe um conjunto de gráficos relativos ao procedimento implementado para um único exemplo. O código fonte da implementação dos algoritmos propostos se encontra no Apêndice A.

### 6.1 Configuração dos parâmetros

As RIRSMs e AVCDs geradas nestes resultados usam as bases de dados de RIR, amostras de voz anecóicas e ruídos descritas no Capítulo 5. As faixas de DRR, T60 e SNR usados são exibidos na Tabela 6.1.

Tabela 6.1: Faixas dos parâmetros.

Parâmetro	Faixa
$DRR_{alvo}$ (dB)	$-6 \leq DRR_{alvo} \leq 18$
$T60_{alvo}$ (s)	$T60_{org} - 1 \leq T60_{alvo} \leq T60_{org} + 1$ , onde o limite inferior de $T60_{alvo} = 0.2$
$SNR_{alvo}$	$3 \leq SNR_{alvo} \leq 20$

O  $DRR_{alvo}$  segue os valores propostos no artigo [4], já os valores de  $T60_{alvo}$ , devido à maior faixa de valores de T60 das RIRs base de AIR, são limitados a  $\pm 1$  segundo comparado ao valor de  $T60_{org}$  da RIR original usada como base para o algoritmo de DA implementado.

## 6.2 Resultados

### 6.2.1 *Data Augmentation* do DRR

Esta seção apresenta os resultados da *Data Augmentation* do DRR. Para isso, foram gerados três exemplos de RIRSMs. Suas configurações são exibidas na Tabela 6.2, onde  $DRR_{org}$  representa o DRR da RIR original,  $DRR_{alvo}$  o valor de DRR desejado pelo usuário,  $DRR_{res}$  o valor de DRR resultante após DA e  $\rho_{DRR}$  é o erro relativo definido da forma  $\rho_{DRR} = \text{abs}(DRR_{res} - DRR_{alvo})/DRR_{alvo}$ .

Tabela 6.2: Exemplos de DA de DRR gerados.

Exemplo	Sala RIR	Distância (m)	Amostra de Voz
D1	lecture	7.1	H2-T2
D2	booth	1	H2-T1
D3	office	2	M2-T2

Exemplo	$DRR_{org}$ (dB)	$DRR_{alvo}$ (dB)	$DRR_{res}$ (dB)	$\rho_{DRR}$ (%)
D1	-4,5	10	10	0
D2	4,7	-2	-2	0
D3	0,5	18	18	0

Nos exemplos gerados, na faixa determinada para o  $DRR_{alvo}$ , não houve diferença entre o  $DRR_{alvo}$  e o  $DRR_{res}$ . É possível observar na Figura 6.2 que ocorreu um aumento na seção equivalente a  $h_e(t)$  comparado à da Figura 6.1 para o exemplo D1.

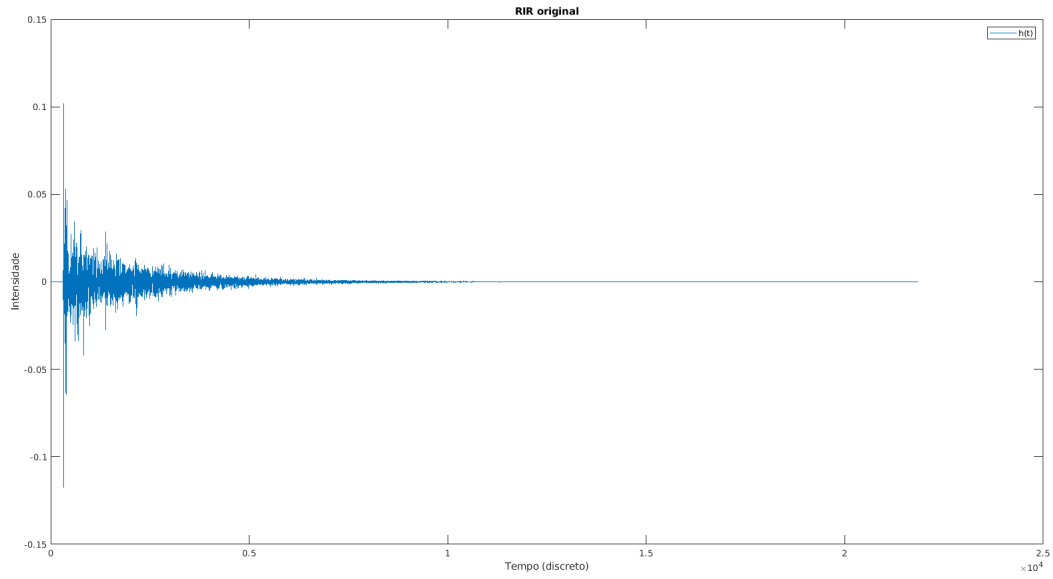


Figura 6.1: RIR original do exemplo D1.

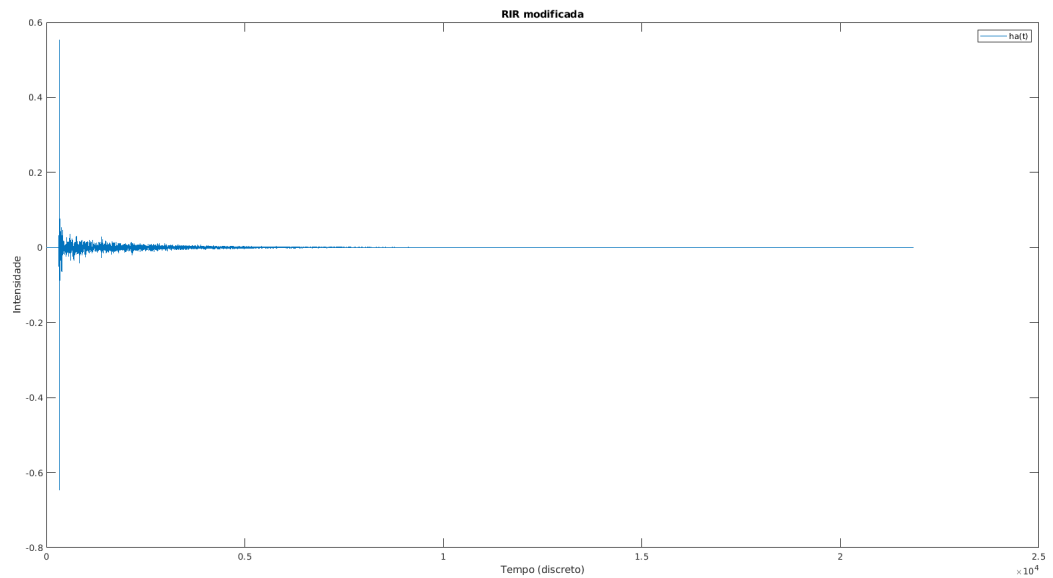


Figura 6.2: RIR simulada do exemplo D1.

Foi realizado também um experimento subjetivo com uma pessoa sem problemas auditivos, cuja idade é de 25 anos na data do experimento, que não possui prévio conhecimento dos resultados gerados. A pessoa foi colocada em um quarto silencioso e foi usado um fone de ouvido para ouvir os áudios gerados para este experimento. O objetivo é avaliar a percepção de “distância” do som, ou seja, a sensação subjetiva do locutor estar falando próximo ou distante do microfone. Foram geradas amostras de voz reverberadas (AVR) usando a RIRO e RIRSM e assim foram feitos dois testes subjetivos exibidos na Tabela 6.3: o primeiro, representado

pela coluna “Comparação”, identifica qual das AVR<sub>s</sub> (a convoluída com RIRO ou a convoluída RIRSM) está mais distante; já o segundo, representado pela coluna “Ordem”, identifica a ordem de mais para menos distante entre as AVR<sub>s</sub> geradas com RIRSMs.

Tabela 6.3: Análise subjetiva de distância.

<b>Exemplo</b>	$DRR_{org}$ (dB)	$DRR_{res}$ (dB)	<b>Comparação</b>	<b>Ordem</b>
D2	4,7	-2	simulado	1 <sup>o</sup>
D1	-4,5	10	original	2 <sup>o</sup>
D3	0,5	18	original	3 <sup>o</sup>

De acordo com a Tabela 6.3, foi possível identificar precisamente a diferença de variações do DRR para as RIRSMs. Na coluna ”Comparação” foi identificado que as RIRs com o menor DRR foram avaliadas subjetivamente como mais distantes. Na coluna ”Ordem” as RIRSMs foram ordenadas conforme o esperado, ou seja, do menor para o maior  $DRR_{res}$ , pois um DRR menor reflete uma sensação de maior distância.

### 6.2.2 *Data Augmentation* do T60

Esta seção apresenta os resultados da *Data Augmentation* do T60. Para isso, foram gerados três exemplos de RIRSMs. Suas configurações são exibidas na Tabela 6.4, onde  $T60_{org}$  representa o T60 da RIR original,  $T60_{alvo}$  o valor de T60 desejado pelo usuário,  $T60_{res}$  o valor de T60 resultante após DA e  $\rho_{T60}$  é o erro definido da forma  $\rho_{T60} = abs(T60_{res} - T60_{alvo})/T60_{alvo}$ .

Tabela 6.4: Exemplos de DA de T60 gerados.

Exemplo	Sala RIR	Distância (m)	Amostra de Voz
T1	lecture	7.1	M2-T1
T2	booth	1	H1-T2
T3	office	2	H2-T2

Exemplo	$T60_{org}$ (s)	$T60_{alvo}$ (s)	$T60_{res}$ (s)	$\rho_{T60}$ (%)
T1	1,38	1,15	1,01	12.1
T2	1,01	1,88	1,89	0,5
T3	0,75	0,61	0,60	1,6

Nos exemplos gerados, na faixa determinada para o  $T60_{alvo}$ , foi observada uma mínima diferença entre o  $T60_{alvo}$  e o  $T60_{res}$  no exemplo T2 e T3. Para o exemplo T1, notamos um erro considerável ao tentar reduzir o T60, o algoritmo proposto tem melhor acurácia para pequenas variações de T60 no caso de redução, contudo o mesmo não é observado para o aumento do T60, mesmo com grandes variações entre  $T60_{alvo}$  e  $T60_{res}$ . As Figuras 6.3 e 6.4 exibem os gráficos das RIRs com intensidade em dB, ou seja,  $10\log_{10}(h(t))$ , para melhor visualização da diferença de decaimento. É possível observar na Figura 6.4 que ocorreu um alongamento da queda exponencial, ou seja, a queda de energia da RIR acontece de forma mais lenta, na seção equivalente a  $h_l(t)$  comparada à Figura 6.3 para o exemplo T2.

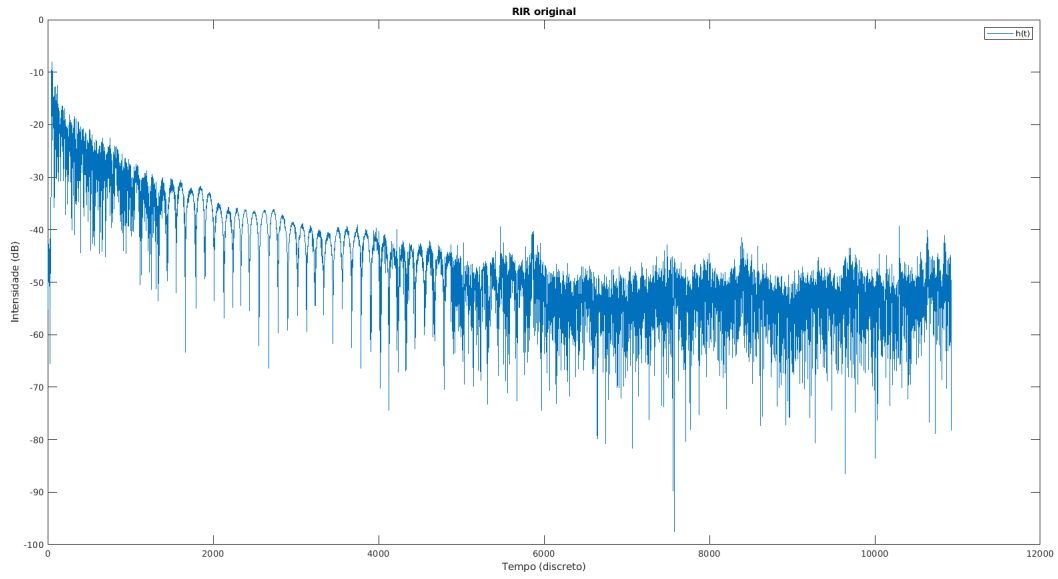


Figura 6.3: RIR original do exemplo T2.

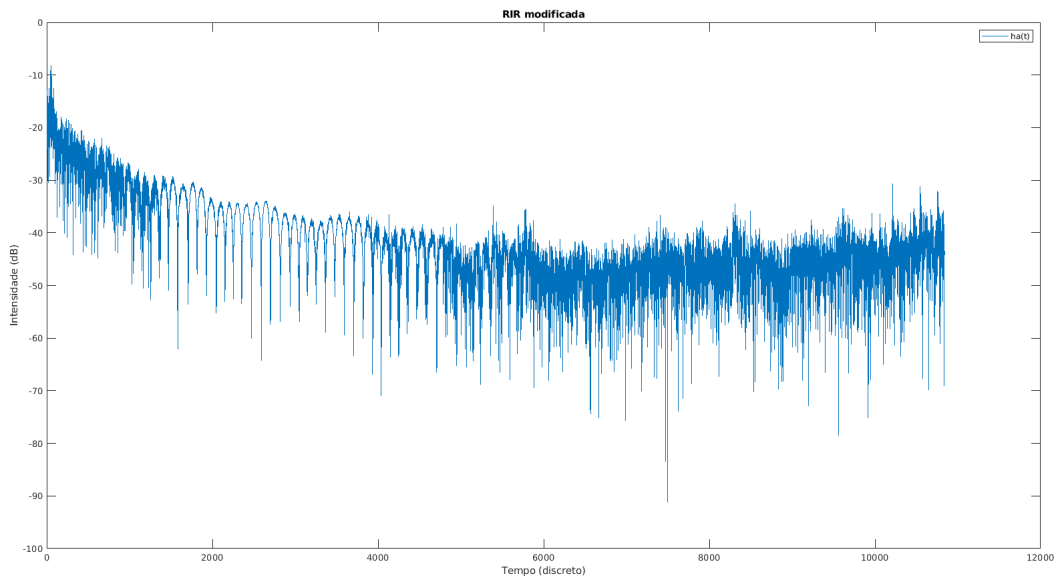


Figura 6.4: RIR simulada do exemplo T2.

De forma análoga aos resultados de DRR, foi feito um experimento subjetivo com o objetivo de avaliar a percepção de “reverberação” do som, ou seja, a sensação subjetiva do locutor estar falando em um espaço fechado mais amplo. Foram geradas amostras de voz reverberadas (AVR) usando a RIRO e RIRSM, e assim foram feitos dois testes subjetivos exibidos na Tabela 6.5: o primeiro, representado pela coluna “Comparação”, identifica qual das AVRs (a convoluída com RIRO ou a convoluída RIRSM) está com mais reverberação. E o segundo, representado pela

coluna “Ordem”, identifica a ordem ,entre as AVRs geradas com RIRSMs, de mais para menos reverberante.

Tabela 6.5: Análise subjetiva de eco.

Exemplo	$T60_{org}$ (s)	$T60_{res}$ (s)	Comparação	Ordem
T2	1,01	1,89	simulado	1 <sup>o</sup>
T1	1,38	1,01	original	2 <sup>o</sup>
T3	0,75	0,60	original	3 <sup>o</sup>

De acordo com a Tabela 6.5, também foi possível identificar precisamente a diferença de variações do T60 para as RIRSMs. Na coluna ”Comparação” foi identificado que as RIRs com o maior T60 foram avaliadas subjetivamente como mais reverberantes. Na coluna ”Ordem” as RIRSMs foram ordenadas conforme o esperado, ou seja, do maior para o menor  $T60_{res}$ , pois um T60 maior reflete a sensação de maior reverberação.

Este resultado demonstra que a DA está ocorrendo, mesmo obtendo o  $\rho_{T60}$  elevado para T1. A discrepância entre  $T60_{alvo}$  e  $T60_{res}$  pode ser inferida pelas diferenças de implementação entre a técnica de DA descrita no artigo [4] e a técnica apresentada neste projeto.

### 6.2.3 *Data Augmentation* de fala em campo distante

Por último, esta seção apresenta os resultados da *Data Augmentation* de AVCDs. Para isso, foram gerados cinco exemplos de AVCDs. Suas configurações são exibidas na Tabela 6.6, onde além dos parâmetros descritos nas seções anteriores, é exibido o  $SNR_{alvo}$  que representa a razão SNR desejada entre a AVCD gerada e os SRP e SRF inseridos.

Tabela 6.6: Exemplos de DA de AVCD gerados.

Exemplo	Sala RIR	Distância (m)	AVA	SRP	SRF
N1	lecture	7.1	M2-T1	RP-6	RF-1
N2	booth	1	H2-T1	RP-12	RF-4
N3	office	2	H1-T1	RP-4	RF-4
N4	meeting	1.7	M1-T2	RP-11	RF-2
N5	stairway	1	H2-T1	RP-7	RF-4

Exemplo	$DRR_{org}$ (dB)	$DRR_{res}$ (dB)	$T60_{org}$ (s)	$T60_{res}$ (s)	$SNR_{alvo}$
N1	-4,5	17	1,38	0,56	5
N2	4,7	17	1,01	1,39	10
N3	0,5	14	0,75	0,60	14
N4	6,0	16	0,81	1,16	19
N5	5,0	18	2,70	3,68	3

Abaixo temos os gráficos relativos ao exemplo N2, contendo a amostra de voz original na Figura 6.5, a amostra de voz reverberada com a RIRSM na Figura 6.6 e a amostra de voz em campo distante na Figura 6.7.

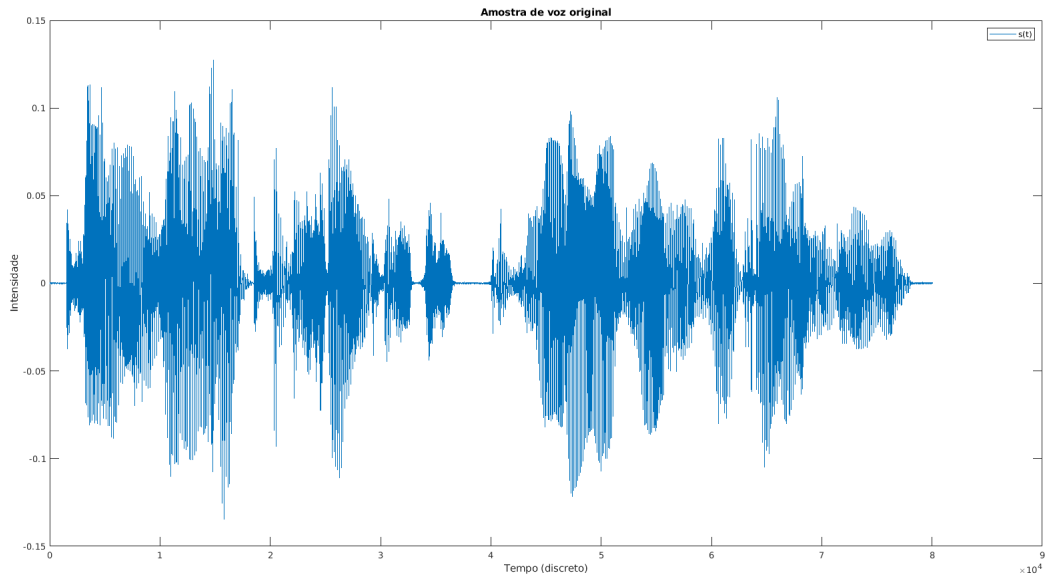


Figura 6.5: Amostra de voz original no exemplo N2.



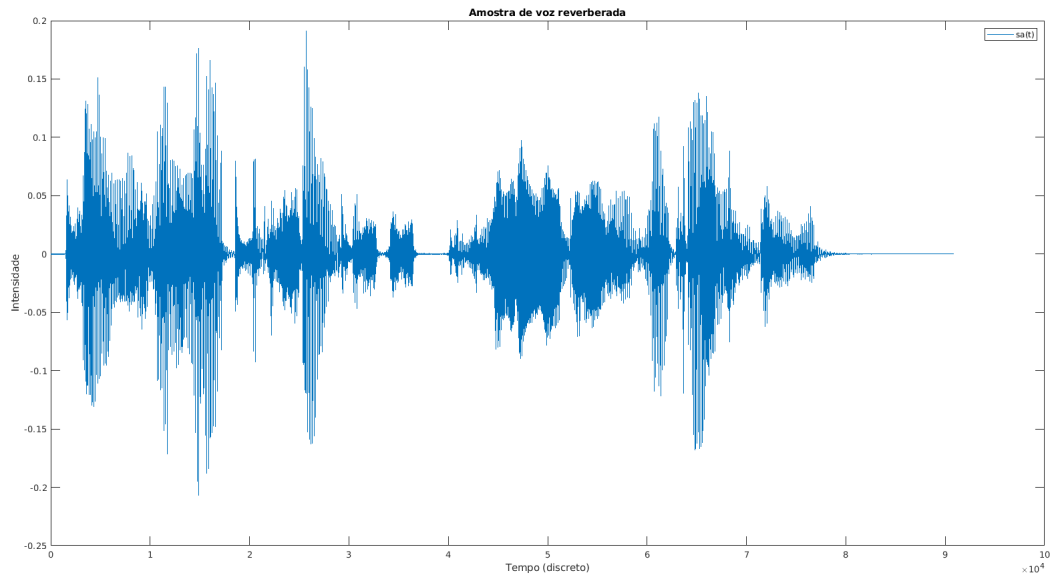


Figura 6.6: Amostra de voz reverberada com RIRSM no exemplo N2.

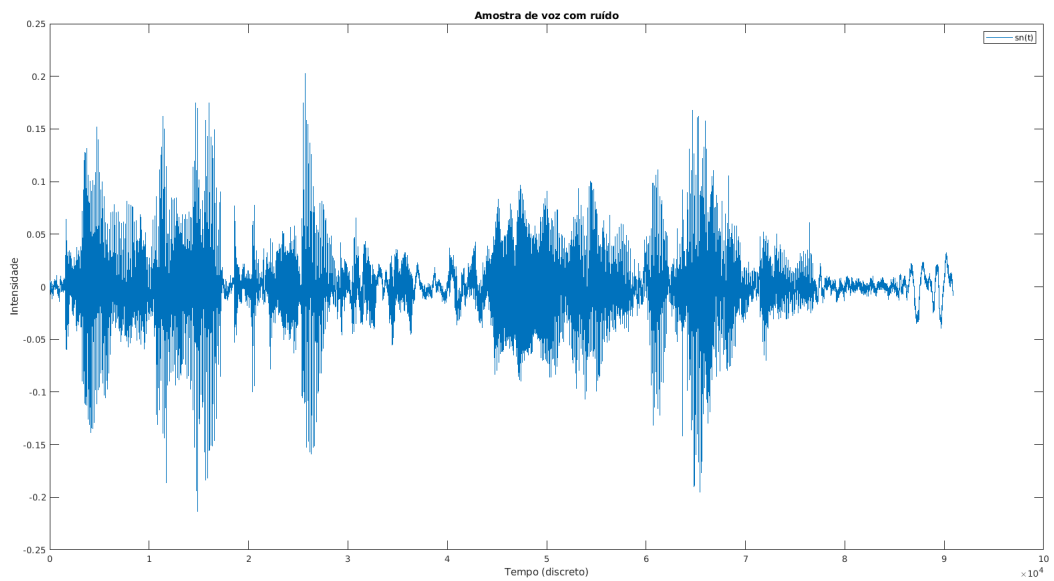


Figura 6.7: Amostra de voz em campo distante no exemplo N2.

Conforme o esperado, observa-se que a Figura 6.7 possui um ruído residual claramente visível comparado à Figura 6.6. Nota-se também que na faixa de tempo entre 2000 e 4000 há picos sonoros que não são observados na amostra reverberada, os quais correspondem ao ruído pontual introduzido no sinal.

De forma análoga, foi feito outro experimento subjetivo com o objetivo de avaliar a percepção de ruído do som, ou seja, a sensação subjetiva de dificuldade de entender a fala do locutor devido aos outros sons misturados na amostra de

voz. Foram usadas as cinco AVCDs apresentadas nesta seção e o teste subjetivo é exibido na Tabela 6.7, onde a coluna “Ordem” identifica a ordem dos sons de mais para menos ruidoso entre as AVCDs geradas.

Tabela 6.7: Análise subjetiva de nível de ruído.

<b>Exemplo</b>	$SNR_{alvo}$ (s)	<b>Ordem</b>
N3	14	1 <sup>o</sup>
N5	3	2 <sup>o</sup>
N1	5	3 <sup>o</sup>
N2	10	4 <sup>o</sup>
N4	19	5 <sup>o</sup>

Observa-se que, os ruídos foram ordenados corretamente de nível decrescente de ruído, onde a única exceção foi o exemplo N3 votado como o mais ruidoso. De acordo com a pessoa que realizou os testes, o exemplo N3 possui um ruído pontual de longa duração (neste caso representado pelo ruído “porta abrindo”), e isso atrapalhou no reconhecimento da fala do locutor.

# Capítulo 7

## Conclusões

Neste trabalho foram propostos dois algoritmos de *data augmentation* com o objetivo de gerar uma base de dados de amostras de voz em campo distante e RIRSMs para treinamento de redes de *deep learning*. Para isso, foi necessário avaliar as principais características e modelos usados nas RIRs para deduzir formas de realizar a modificação das mesmas. Este projeto foi baseado nas técnicas propostas em [4] para DA de RIRs e em [5] para DA de AVCDs.

Ao final do trabalho, foram obtidas diversas RIRSMs e AVCDs geradas através dos algoritmos propostos. Em grande parte, os resultados alcançados estão condizentes com os valores esperados, ou seja, os valores escolhidos durante a geração dos dados. Foi observado uma discrepância considerável entre os valores de T60 modificados para valores abaixo do T60 da RIR original, podendo esta variação ser explicada devido às diferenças de implementação do algoritmo de DA de T60 usados em [4] e o proposto neste projeto.

Quanto às conclusões que podem ser inferidas através dos resultados, nota-se que é possível realizar uma eficaz *data augmentation* de RIRs e AVCDs, mesmo considerando as discrepâncias com os resultados de T60 obtidos, pois foi constatado empiricamente que as variações de “distância” e “eco” são perceptíveis e condizentes com as modificações esperadas.

Para trabalhos futuros, destaca-se a implementação de uma metodologia de *data augmentation* da característica do T60 da RIR que mais se aproxima ao que foi usado em [4]. Neste tópico, seria interessante usar outro modelo de estimativa do T60 e assim observar se há redução nessa discrepância mencionada.

Outra abordagem de trabalho futuro seria comparar os resultados obtidos com as RIRs geradas com o método de DA implementado e RIRs geradas com programas de simulação acústicas, como o RAIOS [23].

Também seria interessante propor um modelo de rede de *deep learning* para estimação de T60 e DRR em AVCDs, realizando dois treinos: um com as RIRSMs e AVCDs geradas pelo algoritmo deste trabalho e outro somente com RIRs e AVCDs reais e assim observar a eficácia da base de dados gerada artificialmente para treinamentos de redes.

# Referências Bibliográficas

- [1] HAEB-UMBACH, R., HEYMANN, J., DRUDE, L., *et al.*, “Far-Field Automatic Speech Recognition”, *Proceedings of the IEEE*, v. 109, n. 2, pp. 124–148, 2021.
- [2] MOKGONYANE, T. B., SEFARA, T. J., MODIPA, T. I., *et al.*, “Automatic Speaker Recognition System based on Machine Learning Algorithms”. In: *2019 Southern African Universities Power Engineering Conference/Robotics and Mechatronics/Pattern Recognition Association of South Africa (SAU-PEC/RobMech/PRASA)*, pp. 141–146, 2019.
- [3] XIONG, F., GOETZE, S., MEYER, B., “Joint Estimation of Reverberation Time and Direct-To-Reverberation Ratio from Speech Using Auditory-Inspired Features”. In: *ACE Challenge Workshop, satellite event of IEEE-WASPAA*, 2015.
- [4] Bryan, N. J., “Impulse Response Data Augmentation and Deep Neural Networks for Blind Room Acoustic Parameter Estimation”. In: *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5, 2020.
- [5] Ko, T., Peddinti, V., Povey, D., *et al.*, “A study on data augmentation of reverberant speech for robust speech recognition”. In: *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5220–5224, 2017.
- [6] SNYDER, D., CHEN, G., POVEY, D., “MUSAN: A Music, Speech, and Noise Corpus”, 2015, <http://www.openslr.org/17/>, visitado última vez em 07/06/2021, arXiv:1510.08484v1.

- [7] BEUTELMANN, R., BRAND, T., “Prediction of speech intelligibility in spatial noise and reverberation for normal-hearing and hearing-impaired listeners”, *The Journal of the Acoustical Society of America*, v. 120, n. 1, pp. 331–342, 2006.
- [8] MAAS, R., HABETS, E. A., SEHR, A., *et al.*, “On the application of reverberation suppression to robust speech recognition”. In: *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 297–300, 2012.
- [9] HELI, H., ABUTALEBI, H. R., “Localization of multiple simultaneous sound sources in reverberant conditions using blind source separation methods”. In: *2011 International Symposium on Artificial Intelligence and Signal Processing (AISP)*, pp. 1–5, 2011.
- [10] NASSIF, A., SHAHIN, I., ATTILI, I., *et al.*, “Speech Recognition Using Deep Neural Networks: A Systematic Review”, *IEEE Access*, v. PP, pp. 1–1, 02 2019.
- [11] VARGAS, R., RUIZ, L., “DEEP LEARNING: PREVIOUS AND PRESENT APPLICATIONS”, *Journal of Awareness*, v. 2, pp. 11–20, 11 2017.
- [12] HAEB-UMBACH, R., WATANABE, S., NAKATANI, T., *et al.*, “Speech Processing for Digital Home Assistants: Combining Signal Processing With Deep-Learning Techniques”, *IEEE Signal Processing Magazine*, v. 36, n. 6, pp. 111–124, 2019.
- [13] TZINIS, E., VENKATARAMANI, S., WANG, Z., *et al.*, “Two-Step Sound Source Separation: Training On Learned Latent Targets”. In: *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 31–35, 2020.
- [14] ALAM, M., SAMAD, M., VIDYARATNE, L., *et al.*, “Survey on Deep Neural Networks in Speech and Vision Systems”, *Neurocomputing*, v. 417, pp. 302–321, 2020.
- [15] PARADA, P. P., SHARMA, D., WATERSCHOOT, T., *et al.*, “Evaluating the Non-Intrusive Room Acoustics Algorithm with the ACE Challenge”, *ArXiv*, v. abs/1510.04616, 2015.

- [16] VIROSTEK, P., “The Quick & Easy Way to Create Impulse Responses”, 2014, <https://www.creativefieldrecording.com/2014/03/19/the-quick-easy-way-to-create-impulse-responses/>, visualizado pela última vez em 16/06/2021.
- [17] NAIR, V., “Recording Impulse Responses”, 2012, <https://designingsound.org/2012/12/29/recording-impulse-responses/>, visualizado pela última vez em 16/06/2021.
- [18] SALAMON, J., BELLO, J. P., “Deep Convolutional Neural Networks and Data Augmentation for Environmental Sound Classification”, *IEEE Signal Processing Letters*, v. 24, n. 3, pp. 279–283, 2017.
- [19] LU, R., DUAN, Z., ZHANG, C., “Metric learning based data augmentation for environmental sound classification”. In: *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pp. 1–5, 2017.
- [20] JEUB, M., SCHäFER, M., VARY, P., “A Binaural Room Impulse Response Database for the Evaluation of Dereverberation Algorithms”. In: *Proceedings of International Conference on Digital Signal Processing (DSP)*, pp. 1–4, IEEE, IET, EURASIP, Santorini, Greece, 2009.
- [21] ISO:, “ISO 3382-1:2009 - Acoustics - Measurement of room acoustic parameters - Part 1: Performance spaces”, <https://www.iso.org/standard/40979.html>.
- [22] ALLEN, J. B., BERKLEY, D. A., “Image method for efficiently simulating small-room acoustics”, *The Journal of the Acoustical Society of America*, v. 65, n. 4, pp. 943–950, 1979.
- [23] TENENBAUM, R., CAMILO, T., TORRES, J. C., *et al.*, “Hybrid method for numerical simulation of room acoustics: Part 2-validation of the computational code RAIOS 3”, *Journal of the Brazilian Society of Mechanical Sciences and Engineering*, v. 29, 04 2007.

# Apêndice A

## Código Fonte

O código fonte da implementação da metodologia deste trabalho está disponível no *GitHub*. Segue o link abaixo.

<https://github.com/afonsobm/RIR-Augmentation-MATLAB>