



# **Speaker Anonymization**

## **Representation, Evaluation and Formal Guarantees**



**Brij Mohan Lal Srivastava**

**Supervisors:** Dr. Aurélien Bellet  
Dr. Emmanuel Vincent  
Prof. Marc Tommasi

Department of Engineering  
Université de Lille

This dissertation is submitted for the degree of  
*Doctor of Philosophy*

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

I would like to dedicate this thesis to my loving parents ...

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## **Declaration**

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text and Acknowledgements. This dissertation contains fewer than 65,000 words including appendices, bibliography, footnotes, tables and equations and has fewer than 150 figures.

Brij Mohan Lal Srivastava  
September 2021

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## **Acknowledgements**

And I would like to acknowledge ...

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## **Abstract**

Speaker anonymization refers to the task of removing speaker-related information from the speech signal.

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

# Table of contents

<b>List of figures</b>	<b>xv</b>
<b>List of tables</b>	<b>xix</b>
<b>Nomenclature</b>	<b>xxi</b>
<b>1 Introduction (7 pages)</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Scope & Objectives . . . . .	3
1.3 Summary of contributions . . . . .	5
1.4 Publications . . . . .	7
1.5 Thesis structure . . . . .	8
<b>2 Background and Related Work</b>	<b>9</b>
2.1 A brief historical overview of speech processing and privacy . . . . .	9
2.2 Principles and tools of speech processing . . . . .	11
2.2.1 Fundamentals of speech processing . . . . .	11
2.2.2 Artificial neural networks . . . . .	16
2.2.3 Automatic speech recognition . . . . .	22
2.2.4 Speech synthesis . . . . .	26
2.2.5 Automatic speaker recognition . . . . .	30
2.3 Techniques to transform speaker information . . . . .	33
2.3.1 Adversarial Learning for Speech . . . . .	34
2.3.2 Speech transformation . . . . .	35
2.3.3 Voice conversion . . . . .	37
2.4 Machine learning based anonymization methods . . . . .	39
2.5 The speaker anonymization task . . . . .	42
2.6 Summary of techniques . . . . .	44
<b>3 Privacy evaluation using Informed attackers</b>	<b>47</b>
3.1 Attack model and the notion of attackers' knowledge . . . . .	47
3.2 Voice conversion methods . . . . .	49
3.2.1 VoiceMask . . . . .	50
3.2.2 VTLN-based voice conversion . . . . .	50
3.2.3 Disentangled representation based voice conversion . . . . .	50
3.3 Target selection strategies and exploitable parameters . . . . .	51
3.3.1 Target selection strategies . . . . .	51

---

3.3.2	Exploitable parameters . . . . .	51
3.4	Performance metrics . . . . .	52
3.4.1	Privacy measures . . . . .	52
3.4.2	Utility measures . . . . .	53
3.4.3	Comparison of privacy metrics . . . . .	54
3.5	Experimental setup . . . . .	54
3.5.1	Data and evaluation setup . . . . .	54
3.5.2	Voice conversion settings . . . . .	56
3.5.3	Design of attackers using ASV . . . . .	57
3.6	Experimental comparison with different attackers . . . . .	58
3.7	Experimental comparison of privacy metrics . . . . .	60
3.7.1	Exhibiting differences and blindspots through simulation . . . . .	60
3.7.2	Evaluation on real anonymized speech . . . . .	63
3.8	Summary . . . . .	64
<b>4</b>	<b>Adversarial Learning based Anonymization</b>	<b>67</b>
4.1	Proposed model . . . . .	68
4.1.1	Baseline end-to-end ASR model . . . . .	69
4.1.2	Speaker-adversarial model . . . . .	69
4.2	Experimental evaluation . . . . .	70
4.2.1	Datasets . . . . .	70
4.2.2	Evaluation metrics . . . . .	71
4.2.3	Network architecture and training . . . . .	72
4.3	Results and Discussion . . . . .	72
4.4	Summary . . . . .	73
<b>5</b>	<b>X-vector based Anonymization</b>	<b>75</b>
5.1	Fixed-pool voice conversion . . . . .	75
5.2	Flexible-pool voice conversion . . . . .	76
5.3	The first VoicePrivacy challenge . . . . .	76
5.3.1	Anonymization task . . . . .	77
5.3.2	Datasets . . . . .	78
5.3.3	Objective and subjective metrics . . . . .	79
5.3.4	Anonymization baselines . . . . .	81
5.3.5	Results . . . . .	82
5.4	Design choices in x-vector space . . . . .	85
5.4.1	Anonymization framework . . . . .	86
5.4.2	Distance metric . . . . .	87
5.4.3	Proximity . . . . .	88
5.4.4	Gender selection . . . . .	88
5.4.5	Assignment . . . . .	89
5.4.6	Experimental setup . . . . .	89
5.4.7	Results and Discussion . . . . .	92
5.5	Large-scale speaker study . . . . .	103
5.5.1	Remark on utility . . . . .	104
5.5.2	Privacy evaluation metrics . . . . .	104
5.5.3	Gender identification . . . . .	105

## Table of contents

xiii

---

5.5.4	Experimental setup . . . . .	105
5.5.5	Experiments and results . . . . .	106
5.5.6	Worst-case analysis . . . . .	110
5.6	Summary . . . . .	115
<b>6</b>	<b>Towards removing residual speaker information and provable guarantee</b>	<b>117</b>
6.1	Individual impact on privacy and utility . . . . .	117
6.2	Adding Differentially-Private Noise in F0 and BN . . . . .	117
<b>7</b>	<b>Usability of anonymized speech for training ASR (20 pages)</b>	<b>119</b>
7.1	Impact of re-training ASR . . . . .	119
7.2	Data augmentation . . . . .	119
7.3	Model adaptation . . . . .	119
<b>8</b>	<b>Conclusion and Perspectives (4 pages)</b>	<b>121</b>
<b>References</b>		<b>123</b>
<b>Appendix A</b>	<b>Extra top-<math>k</math> results</b>	<b>143</b>
<b>Appendix B</b>	<b>Worst-case analysis of the anonymization scheme</b>	<b>145</b>

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

# List of figures

2.1	Waveform, magnitude spectrogram, MFCC and pitch contour for the word “privacy”, pronounced by a male (left) and a female speaker (right). . . . .	13
2.2	Perceptron model of a neuron with $A = 3$ . . . . .	17
2.3	Fully-connected feed-forward neural network (multilayer perceptron). . . . .	18
2.4	Time delay neural network architecture with dilation in layers 2 and 3. Dotted lines indicates the connections and the nodes which are not included in the computation due to dilation applied to the layers. . . . .	20
2.5	Factorized TDNN (TDNN-F) architecture showing the linear bottleneck inserted between the hidden layer and the feature map. Matrices $\mathbf{P}_1$ and $\mathbf{P}_2$ are constrained to be semi-orthogonal and $\oplus$ represents the concatenation of the linear bottleneck and the feature map layer linked by the skip connection. . . . .	21
2.6	Generative model for ASR. . . . .	23
2.7	Network architecture for ASR acoustic modeling composed of TDNN-F layers followed by the fully-connected bottleneck layer $\mathbf{B}$ which branches into the computation of the two loss functions, LF-MMI and cross-entropy. The skip connections between TDNN-F layers are not shown for the sake of simplicity. . . . .	24
2.8	End-to-end ASR architecture with multi-objective training consisting of CTC and attention-based loss functions combined by a hyperparameter $\beta$ , where $0 < \beta < 1$ . The attention layer is shown as the heatmap in the attention decoder which assigns combination weights to the bottleneck representation $\mathbf{B}$ before processing it by the LSTM layer. . . . .	26
2.9	General schema of a TTS system with the style features as optional input to the frontend. . . . .	26
2.10	Autoregressive network architecture for the TTS acoustic model [179]. The feedforward layers have Tanh activation and the linear layers have identity activation. It has skip connections to reinforce the noisy signal, and the time delay block passes the current output Mel-spectrogram back to the LSTM layer with a random dropout. . . . .	28
2.11	NSF model architecture. . . . .	29
2.12	Schema for two types of automatic speaker recognition. Red arrows indicate the enrollment or training flow, and green arrows indicate the authentication or testing flow. Note that the speaker model trained to classify speakers in the case of ASI can also be used to extract speaker embeddings for ASV, just like x-vectors. . . . .	31
2.13	Score distribution and threshold. . . . .	32
2.14	General architecture for domain adversarial training of neural networks. Black arrows indicate forward propagation; purple, teal and red arrows indicate backpropagation of gradients for the primary task classifier, the encoder and the adversarial branch, respectively. The red GRL block refers to the gradient reversal layer with $\lambda$ as the gradient reversal coefficient. . . . .	34

2.15	Bilinear function warping the frequency $\omega \in [0, \pi]$ using positive and negative values of $\alpha \in \{-0.7, -0.3, 0, 0.3, 0.7\}$ . . . . .	36
2.16	General schema for traditional voice conversion with parallel data. Red arrows indicate training flow, while green arrows show the flow during conversion. . . . .	37
3.1	Threat model . . . . .	48
3.2	Attacker's knowledge continuum . . . . .	49
3.3	Three target selection strategies: const, perm and rand. Trial utterances are publicly released data set, while enroll utterances are found data used by the attacker. Utterances are shown by small green balls, and the arrows indicate the mapping between original to target speakers. . . . .	51
3.4	Utility evaluation . . . . .	56
3.5	Privacy evaluation . . . . .	57
3.6	I-vector score distribution for trials conducted on VTLN (strategy <i>random</i> ) converted data by <i>Ignorant</i> , <i>Semi-Informed</i> , or <i>Informed</i> attackers. The orange distribution indicates impostor scores, while the blue distribution indicates genuine scores. The crossing between the two curves indicates the threshold for EER. More overlap means greater confusion, hence greater privacy protection. . . . .	59
3.7	$C_{llr}^{\min}$ vs. $1 - D_{\leftrightarrow}^{\text{sys}}$ on simulated Gaussian scores. . . . .	62
3.8	Simulated ‘non-mated in-between’ data. Top: x-vectors visualized in 2D. Bottom: resulting score distributions. . . . .	62
3.9	$C_{llr}^{\min}$ vs. <i>EER</i> on real data. The color scale $\mu - \bar{\mu}$ is the difference of the means of mated and non-mated scores. . . . .	63
3.10	$C_{llr}^{\min}$ vs. $1 - D_{\leftrightarrow}^{\text{sys}}$ on real data. The color scale $\mu - \bar{\mu}$ is the difference of the means of mated and non-mated scores. . . . .	64
4.1	Threat model related to speech-to-text provided by cloud-based services. . . . .	67
4.2	Architecture of the proposed model. The speaker-adversarial branch is shown as a red box. The teal arrow going from GRL to encoder indicates <i>gradient reversal</i> . When the model is deployed, the encoder could reside at the client side, while the decoder can be hosted by cloud services. . . . .	70
4.3	Visualization of x-vector representations of 20 utterances of 10 speakers computed by t-SNE (perplexity equals to 30). Males are represented by circles and females by triangles. . . . .	73
5.1	Speech synthesis based VC framework conditioned upon continuous speaker representation that can be replaced by unseen targets. . . . .	76
5.2	Objective evaluation of privacy protection provided by the two baseline systems in the first VoicePrivacy challenge. Higher EER indicates better protection. . . . .	84
5.3	Subjective evaluation of privacy protection provided by the two baseline systems in the first VoicePrivacy challenge. Bar plots of <i>subjective speech naturalness</i> , <i>intelligibility</i> , and <i>speaker similarity</i> obtained from the normalized scores. For naturalness and intelligibility, scores from target and non-target anonymized data are pooled; for similarity, scores for anonymized target and non-target speakers data are separately plotted in 3rd and 4th sub-figures, respectively. Numbers indicate mean values over all the three data sets. Higher values for naturalness and intelligibility correspond to better utility, and lower scores for similarity to target speaker with anonymized data from target speaker — to better privacy	85
5.4	General architecture of the anonymization system. . . . .	86

5.5	Zoomed-in view of the x-vector anonymization step in Fig. 5.4 showing the design choices for the generation of the target x-vector. . . . .	87
5.6	Relationship between EER, linkability and $C_{llr}^{\min}$ in <i>Lazy-Informed</i> setting for various combination of design choices. . . . .	91
5.7	Privacy against <i>Ignorant</i> and <i>Lazy-Informed</i> attackers depending on the distance choice and the gender of the original speaker. Each swarm plot shows the 24 linkability values for each gender on the development set resulting from all combinations of proximity (excluding <i>random</i> ), gender selection, and assignment choices. . . . .	93
5.8	Privacy against <i>Ignorant</i> and <i>Lazy-Informed</i> attackers depending on the proximity choice and the gender of the original speaker. Distance is fixed to PLDA. Each swarm plot shows the 6 linkability values for each gender on the development set resulting from all combinations of gender selection and assignment choices. . . . .	93
5.9	Privacy against <i>Ignorant</i> and <i>Lazy-Informed</i> attackers depending on the gender selection choice and the gender of the original speaker. Distance is fixed to PLDA and proximity to <i>dense</i> or <i>random</i> . Each swarm plot shows the 4 linkability values for each gender on the development set resulting from the assignment choice and the 2 proximity choices. . . . .	94
5.10	Average cosine distance on the development set between the x-vectors of the original utterance and the anonymized utterance in Step 4 plotted against the average cosine distance between the x-vector of the original utterance and the target x-vector in Step 2. Distance is fixed to PLDA. For each choice of proximity (color) and gender selection (marker shape), four values are shown depending on the gender of the original speaker and the assignment choice. . . . .	95
5.11	Privacy against <i>Ignorant</i> and <i>Lazy-Informed</i> attackers depending on the assignment choice and the gender of the original speaker. Distance is fixed to PLDA, proximity to <i>dense</i> or <i>random</i> , and gender selection to <i>random</i> . Each swarm plot shows the 2 linkability values for each gender on the development set resulting from the 2 proximity choices. . . . .	96
5.12	Utility of anonymized speech in terms of WER compared to the original (baseline) speech depending on the different design choices and the gender of the original speaker. Each swarm plot shows the WER values on the development set for each gender and for a given design choice. The remaining design choices are fixed in the same way as in Figs. 5.7, 5.8, 5.9 and 5.11. . . . .	97
5.13	t-SNE visualization of speaker-level x-vectors from the LibriSpeech <i>train-clean-360</i> data set transformed using different anonymization schemes. . . . .	98
5.14	Performance of ASR <sub>eval</sub> <sup>anon</sup> models re-trained using anonymized speech corpus obtained by two methods, namely Random proximity and Dense proximity. . . . .	99
5.15	Performance of informed ASV <sub>eval</sub> <sup>anon</sup> models re-trained using anonymized speech corpus obtained by two methods, namely Random proximity and Dense proximity. BL = Original (baseline), Ign = <i>Ignorant</i> , Lazy-I = <i>Lazy-Informed</i> and Semi-I = <i>Semi-Informed</i> attacker. . . . .	101
5.16	Performance of ASR <sub>eval</sub> <sup>anon</sup> and ASV <sub>eval</sub> <sup>anon</sup> after the two types of target pitch transformation as compared to original pitch. . . . .	102
5.17	t-SNE representation of speaker x-vectors in common voice data set. . . . .	106
5.18	Speaker gender distribution observed in the common voice data set. The exact number of speakers for each gender are indicated in the brackets. Four speakers are discarded due to lack of data. . . . .	107
5.19	Open-set ASV performance of different attackers in terms of EER and Linkability as the speakers in the population double at each step. . . . .	108

5.20	Closed-set ASI performance in terms of un-normalized and normalized rank obtained by different attackers as the number of speakers in the population double at each step. . . . .	109
5.21	Top-20 precision of ASI for different attackers as the speaker population is doubled at each step. The number of speakers needed before anonymization ( $N$ on blue curve) and after anonymization ( $n$ on red curve) to achieve equivalent drop in precision are highlighted. . . . .	110
5.22	The normalized rank distribution for the baseline case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of $U^{\text{worst}}$ (red), $S^{\text{worst}}$ (green), $\bar{U}_S^{\text{worst}}$ (black). . . . .	111
5.23	The normalized rank distribution for the <i>Semi-Informed</i> case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of $U^{\text{worst}}$ (red), $S^{\text{worst}}$ (green), $\bar{U}_S^{\text{worst}}$ (black). . . . .	112
5.24	Line plot in original case showing the normalized rank (y-axis) against the duration of utterances (x-axis). The dashed horizontal lines show the value of $U^{\text{worst}}$ (red), $S^{\text{worst}}$ (green), $\bar{U}_S^{\text{worst}}$ (black). . . . .	113
5.25	Line plot in <i>Semi-Informed</i> case showing the normalized rank (y-axis) against the duration of utterances (x-axis). The dashed horizontal lines show the value of $U^{\text{worst}}$ (red), $S^{\text{worst}}$ (green), $\bar{U}_S^{\text{worst}}$ (black). . . . .	114
A.1	Speaker identification top-1 performance as enrollment speakers increase. . . . .	143
A.2	Speaker identification top-10 performance as enrollment speakers increase. . . . .	144
A.3	Speaker identification top-50 performance as enrollment speakers increase. . . . .	144
B.1	The normalized rank distribution for the <i>Ignorant</i> case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of $U^{\text{worst}}$ (red), $S^{\text{worst}}$ (green), $\bar{U}_S^{\text{worst}}$ (black). . . . .	146
B.2	The normalized rank distribution for the <i>Lazy-Informed</i> case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of $U^{\text{worst}}$ (red), $S^{\text{worst}}$ (green), $\bar{U}_S^{\text{worst}}$ (black). . . . .	147
B.3	Normalized rank for the worst-performing utterances, i.e., $U^{\text{worst}}$ . . . . .	148
B.4	Normalized rank for the worst-performing speaker, i.e., $S^{\text{worst}}$ . Whiskers indicate the standard deviation of the normalized rank for the utterance of the worst-performing speaker.	148
B.5	Normalized rank for the worst-performing utterances of each speaker, i.e., $\bar{U}_S^{\text{worst}}$ . Whiskers indicate the standard deviation of the normalized rank for the worst-performing utterance of each speaker. . . . .	149

# List of tables

2.1	Original network architecture for x-vector speaker classification model as presented in [262, Table 1]. . . . .	32
3.1	The subsets of the LibriSpeech data set along with their total duration in hours, duration per speaker in minutes, and number of male and female speakers. . . . .	54
3.2	Detailed description of the trial set for speaker verification experiments. . . . .	55
3.3	EER (%) achieved using x-vector based ASV <sub>eval</sub> for <i>Ignorant</i> attacker, and ASV <sub>eval</sub> <sup>anon</sup> for <i>Semi-Informed</i> and <i>Informed</i> attackers. Bold face indicates the best privacy protection strategy against <i>Informed</i> attackers. . . . .	58
3.4	EER (%) achieved using i-vector based ASV <sub>eval</sub> for <i>Ignorant</i> attacker, and ASV <sub>eval</sub> <sup>anon</sup> for <i>Semi-Informed</i> and <i>Informed</i> attackers. Bold face indicates the best privacy protection strategy against <i>Informed</i> attackers. . . . .	59
3.5	Additional i-vector results with VTLN <i>Semi-Informed</i> attackers against the protection strategies. Bold face indicates the best performing attacker's strategy against a protection strategy. . . . .	60
3.6	WER (%) achieved using end-to-end ASR <sub>eval</sub> <sup>anon</sup> . . . . .	60
3.7	$C_{llr}^{\min}$ and EER with discrete scores in $\{1, \dots, 8\}$ . $H$ (resp. $\bar{H}$ ) denote mated (resp. non-mated) scores. . . . .	61
4.1	Splits of Librispeech used in our experiments. . . . .	71
4.2	ASR and speaker recognition results with different representations. WER (%) is reported on <i>test-clean</i> set, ACC (%) on <i>test-adv</i> set and EER (%) on <i>test-clean-trial</i> . . . . .	73
5.1	Statistics of the training data sets. . . . .	78
5.2	Statistics of the development data sets. . . . .	78
5.3	Statistics of the evaluation data sets. . . . .	79
5.4	Statistics of the training data set for the ASV <sub>eval</sub> and ASR <sub>eval</sub> evaluation systems. . . . .	79
5.5	Number of speaker verification trials in objective evaluation on speaker verifiability. . . . .	80
5.6	Number of trials evaluated in subjective evaluation on verifiability, intelligibility, and naturalness. Anonymized trials for subjective evaluation are from 9 anonymized systems (baseline and primary participants' systems as described in [289]). The number of speakers is 30 (15 male and 15 female) in each data set. . . . .	80
5.7	Baseline-1 system: model architectures, objective functions, output features that are used in the anonymization pipeline, and training corpora are mentioned. Superscript numbers represent feature dimensions. . . . .	81

5.8	Speaker verifiability achieved by the pretrained ASV <sub>eval</sub> model in original, <i>Ignorant</i> and <i>Lazy-Informed</i> scenarios, and the ASV <sub>eval</sub> <sup>anon</sup> model in <i>Semi-Informed</i> case. Baseline-1 is used for anonymization. . . . .	83
5.9	ASR decoding error achieved by the pretrained ASR <sub>eval</sub> model. Baseline-1 is used for anonymization. . . . .	83
5.10	ASR decoding error achieved by the pretrained ASR <sub>eval</sub> model. Baseline-2 is used for anonymization. . . . .	84

# Nomenclature

## Roman Symbols

$A$	Number of input features
$g$	Activation function
$\mathbf{a}$	Alignment sequence
$M$	Number of possible alignments
$\mathbf{B}$	Bottleneck features
$F$	Frequency in Hz
$F_0$	Fundamental frequency
$\mathcal{F}$	Fourier coefficients
$\vec{s}$	Speech time frame
$F_s$	Sampling frequency
$G$	Grapheme sequence
$\mathcal{G}$	Set of grapheme symbols
$y$	Ground truth
$\mathcal{X}$	Input data set
$L$	Frame length
$\mathcal{L}$	Loss function
$N$	Number of observations in a data set
$\mathbf{o}$	Acoustic features of a time frame
$\mathbf{O}$	Acoustic features of an utterance
$\hat{y}$	Estimated output
$\mathcal{Y}$	Output data set

$s$	Time-domain speech signal
$T$	Number of time frames in an utterance
$w$	Synaptic weight of a neuron
$W$	Word sequence

**Greek Symbols**

$\xi$	Frame energy
$\lambda$	Gradient reversal coefficient
$\eta$	Learning rate
$\theta$	Neural network parameters
$\rho$	Phoneme sequence
$\Psi$	Short-time analysis window

**Acronyms / Abbreviations**

ASI	Automatic Speaker Identification
ASR	Automatic Speech Recognition
ASV	Automatic Speaker Verification
IPA	International Phonetic Alphabet
MFCC	Mel Frequency Cepstral Coefficient
STFT	Short-Time Fourier Transform
TTS	Text-to-Speech

# Chapter 1

1

## Introduction (7 pages)

2

**Deadline: Feb 26, 2021**

Civilization is the progress toward a society of privacy.

---

Ayn Rand

### 1.1 Motivation

3

Speaking and listening are the most convenient, non-tactile and expressive forms of human communication. During oral conversations, we transmit not only the linguistic content, but also paralinguistic and extra-linguistic cues such as our emotional state, age, gender, personality, health state, etc. to our interlocutors [282]. Hence the rich nature of verbal dialog makes it a natural choice as the interface between humans and machines. For over two centuries, researchers have been intrigued by the process of speech generation by humans and machines [136]. The earliest known attempt of producing human-like sounds from a mechanical model was made in the later part of the 18th century by the Russian scientist Christian Kratzenstein [157] and soon after by an Austrian inventor named Wolfgang von Kempelen [302] who is famous for his “speaking machine”.

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

Fast forward to the 1980s, Artificial Neural Networks were successfully introduced to recognize a few words in a speech signal [175]. The field of speech processing has come a long way since then as four decades later and with several groundbreaking advances, Deep Neural Networks (DNNs) have become the state-of-the-art [108, 112] for large vocabulary continuous speech recognition and several associated tasks, such as text-to-speech, speaker recognition, etc. These gigantic networks with several millions of parameters have surpassed human level performance in speech recognition [12], but there are new and extraordinary challenges with the emergence of speech interfaces in the marketplace.

20

21

22

23

24

25

Several smart digital assistants are available in the market today which are powered by the decades of advances in speech recognition and conversational models. The goal of their manufacturers is to make conversations between humans and digital assistants as seamless as possible.<sup>1</sup> They are designed to handle a wide range of commands and questions in a jesterful manner, and are usually provided with a name and gender so that users can personify them. Such realizations can build trust between the digital assistant and the human, which helps to enhance the engagement [181]. Users can now control their home appliances,

---

<sup>1</sup><https://developer.amazon.com/alexaprize/about>

1 play music, request a joke, shop online and of course send messages among several other functions using  
2 the digital assistant. The users of digital assistants are growing at the rate of about 35% per year. Till early  
3 2019, Amazon had already sold 100 million devices worldwide [35] with Alexa<sup>2</sup>, a cloud-based voice-enabled  
4 virtual assistant, and the projected market of digital assistant by the end of 2021 is expected to reach 843  
5 million users worldwide, which amounts to a revenue of \$15.8 billion [231].

6 Unfortunately, the most determining factor in the success of the advanced statistical models running  
7 the digital assistant ecosystem is the enormous size of their training data sets [129], closely followed by the  
8 availability of high computing infrastructure such as Graphical Processing Units (GPUs) [120]. With the  
9 availability of pervasive Internet and smart devices, large quantities of speech data is being collected by  
10 digital assistant manufacturers like Google, Amazon, Apple and Microsoft. This data is stored at centrally  
11 located servers and, depending on the needs, it can be made available to developers, annotators and managers.  
12 Among the consumers who own a digital assistant, 65% claim that they do not know everything the device  
13 can do [225] and, among those who do not own a digital assistant, only 16% cite privacy reasons not to  
14 purchase one. This lack of awareness further opens the doors to a massive privacy breach.

15 Privacy is considered as a fundamental human right in many regions of the world [24]. It is intimately  
16 linked with human dignity and freedom of thought and expression. The Indian Constitution lists the “right  
17 to privacy” under Article 21 which deals with protection of life and personal liberty [43]. Yet there is no  
18 universally accepted definition of privacy [182] which makes it hard to enforce it as a legally protected right.  
19 The extent of technological intrusion into people’s lives as described previously poses a severe threat to  
20 individual privacy but there is little consensus between technological and legal communities to legislate  
21 strong laws for data protection. In 1890 Warren and Brandeis [309] defined privacy for the first time as the  
22 “general right of the individual to be let alone”. Since then the European Union and several countries such as  
23 the United States of America and Canada [169] have included some privacy articles in their constitution.

24 In 2016, the European Union passed the General Data Protection Regulation (GDPR) [52] setting  
25 a historical precedent for data privacy law worldwide. The GDPR is listed under the EU Charter of  
26 Fundamental Rights which stipulates that European citizens have right to protection of their personal  
27 data.<sup>3</sup> The law clearly holds companies accountable for users’ data, users have complete control over the  
28 usage and distribution of their data and can request deletion at any time they want. Moreover, the law is  
29 applicable over the data of citizens of all the member states even if it was processed overseas. The violators  
30 are heavily fined [18] if found guilty, causing companies to block users from Europe [155] to avoid non-  
31 compliance issues. Although Section 2 of the GDPR mentions guidelines to ensure the “security of personal  
32 data” through pseudonymisation and encryption, there is a clear lack of understanding with respect to the  
33 capture, storage and processing of speech data. Recently, Nautsch et al. [207] launched a collaborative effort  
34 to harmonise the terminology between speech researchers and legal experts so that the sensitive attributes in  
35 speech signals could be clearly understood by the legislators.

36 Recently the French Data Protection Authority (CNIL) published a white paper [60] to explore legal,  
37 technical and ethical issues associated with voice assistants. The paper briefly mentions the work done during  
38 this thesis as the potential solution to some of the technical issues. At the time of writing this thesis, the  
39 European Data Protection Board (EDPB) also released guidelines [34] on virtual voice assistants for different  
40 stakeholders involved in their production and use. Although the guidelines focus on the legal bases which  
41 empower digital assistant users to request data erasure and voice sanitisation techniques to remove situational  
42 information and background noise, they are not very encouraging of voice anonymization methods due  
43 to several open challenges that impede the evolution of the technology. Nevertheless, the guidelines were

<sup>2</sup><https://developer.amazon.com/en-US/alexa>

<sup>3</sup>[https://ec.europa.eu/info/law/law-topic/data-protection/data-protection-eu\\_en](https://ec.europa.eu/info/law/law-topic/data-protection/data-protection-eu_en)

---

1.2 Scope & Objectives

3

open for public consultation calling for further technological as well as legal development in this direction to explore exact grounds for processing speech data and clear jurisdiction upon violation.

Speech is a biometric characteristic of human beings [134], which can produce distinguishing and repeatable biometric features. According to Article 4(14) of the GDPR, voice data is inherently biometric personal data which relates to physical, physiological or behavioural characteristics of a natural person. Voiceprints [234], which are used to identify speakers, are also deployed in payment services [28] for authenticating transactions. The wealth of personal information present in speech signals and the availability of efficient techniques to identify that information pose a severe privacy risk for the users of speech interfaces. In particular, recent advances in voice cloning [300, 178, 209] and synthesis [276, 241] techniques that leverage “found speech” call for efficient speaker anonymization schemes.

Several open challenges pertaining to the privacy of speech data arise due to the emergence of large-scale data collection by voice-enabled apps and SPA devices. Privacy breaches by corporates like Samsung<sup>4</sup>, Apple<sup>5</sup> and Google<sup>6</sup> have made headlines in the newspapers. The concern for malicious usage of such sensitive data has widely alarmed individual citizens, researchers and the legal community. Recently, governments have also shown political will to support the efforts towards effective formulation of laws to achieve voice data protection by design and by default, along with the supporting technological advances that can secure the rights and interests of common citizens. With the above mentioned motivation, this thesis is a timely effort to propose speaker anonymization methods which aim to remove speaker identity from speech signals while keeping other linguistic attributes and speech quality intact. For widespread adoption of this technology, it is also important that the transformed speech remains usable for downstream tasks such as training an automatic speech recognition (ASR) model.

## 1.2 Scope & Objectives

Although there have been a few efforts to protect speech signals against external attacks, the topic of privacy-preserving speech processing itself has attracted quite limited interest so far. The methods proposed in the last decade can be broadly classified into four categories: deletion, encryption, distributed learning and anonymization. Deletion refers to the blurring or obfuscation of sensitive segments of speech [50, 104] while retaining the acoustic scene, but has limited scope in terms of diverse speech applications. Encryption based methods [331, 36] aim to secure the transmission channel and perform operations in the encrypted domain, but incur a high computational cost and may require special hardware. Distributed learning methods such as federated learning [168] are machine learning techniques where training is performed by averaging gradients coming from several distributed nodes, thereby ensuring decentralization of training data to avoid central ownership of massive datasets, but may not protect the privacy of the speaker due to information leaking through gradients [98]. Finally, *anonymization* refers to the task of suppressing personally identifiable attributes of speech signal, leaving all other attributes intact. This thesis is a consolidated effort to propose effective methods and rigorous evaluation schemes for speaker anonymization. Hence we briefly review deletion, encryption and distributed learning based approaches here to present our arguments against using these approaches, thereby clearly defining the exact scope of our methods.

The earliest attempt at processing speech data in a private and secure manner mostly assumed a client-server model for speech applications, where the two parties communicate through mutually understood

<sup>4</sup><https://www.bbc.com/news/technology-31296188>

<sup>5</sup><https://www.theguardian.com/technology/2019/jul/26/apple-contractors-regularly-hear-confidential-details-on-siri-recordings>

<sup>6</sup><https://www.bloomberg.com/news/articles/2019-04-10/is-anyone-listening-to-you-on-alexa-a-global-team-reviews-audio>

<sup>7</sup><https://nos.nl/artikel/2292889-google-medewerkers-luisteren-nederlandse-gesprekken-mee.html>

1 encrypted speech tokens and the transmission channel is secured by cryptographic methods, such as secure  
2 multiparty computation [260], secure two-party computation [27], hash functions and homomorphic  
3 encryption to represent and process speech [217], as well as nearest neighbour audio query search [233] or  
4 phonetic search [100] in the encrypted domain. Recently, with the introduction of cryptographic methods  
5 like homomorphic encryption and Paillier cryptosystem in neural networks [213, 13, 58], sensitive attributes  
6 in speech such as emotions have been identified in a secure manner [66]. Finally, Intel's<sup>8</sup> Software Guard  
7 Extensions (SGX) [54] provides a trusted execution environment in private regions of memory, called  
8 enclaves, for carrying out sensitive operations. VoiceGuard [36] presents a proof of concept by executing a  
9 speech recognition engine within the SGX enclave.

10 Although cryptographic methods have advanced manyfold, their strength is hinged upon the future  
11 breakability of the underlying encryption algorithm and the hardware resilience to adversaries. They require  
12 additional computational overhead and special hardware for successful implementation. Similar methods  
13 that protect privacy by reducing the speech signal to hash tokens destroy the paralinguistic and extra-linguistic  
14 characteristics of the signal, thereby losing all the utility for downstream tasks. Recent advances in federated  
15 learning [171] with the goal of decentralized ownership of user's data have enabled researchers to apply it to  
16 speech processing applications such as keyword spotting [168] and emotion recognition [161]. Although  
17 federated learning claims to protect users' privacy by not requiring them to share their data, it is not resistant  
18 to membership inference attacks [206]. Moreover, it has been shown that it is possible to reconstruct user's  
19 data given the knowledge of the received gradients [98].

20 This quick review of speech related privacy-preserving methods reveals that deletion, encryption or  
21 distributed learning based methods do not address the primary concern of this thesis, that is, to obtain an  
22 anonymous yet useful representation of speech with strong privacy guarantees. These methods do not focus  
23 specifically on the biometric speaker information present in the signal, instead they securely obfuscate the  
24 whole speech signal or devise a trusted data sharing mechanism. Hence, these methods do not align with our  
25 objectives, which are as follows: to recognise specific biometric identifiers present in the speech signal which  
26 makes it linkable to the speaker, to learn a *global transform* which could remove these identifiers effectively  
27 from the signal without affecting the linguistic content, and to evaluate the effectiveness of identity removal  
28 through strong attack measures and formal protection guarantees. Certainly, there have been earlier attempts  
29 to study speech transformation methods that claim to have removed speaker's identity up to a certain extent  
30 with varying loss of utility. We present an in-depth review of the research material on such anonymization  
31 techniques that closely align with our objectives in Section 2.4.

32 In essence, we try to answer the following central question in this thesis:

33 *While maintaining the usefulness of the signal, how to effectively remove the biometric identity  
34 of the speaker from any speech utterance?*

35 With the above stated central goal as the "holy grail" of privacy protection in speech, we reiterate that  
36 usability as well as privacy are the most important objectives of speaker anonymization. Here "usefulness" is  
37 a broad term that encompasses the ability to train models for downstream tasks such as ASR, human-level  
38 intelligibility for listening and transcription as well as the presence of the inherent varibilities of the natural  
39 speech signal. The intent to preserve these qualities of the speech signal while carrying out the process of  
40 anonymization emerges from the perspective of the potential users of this transformed speech data. Without  
41 the usability of the transformed speech corpus, the widespread adoption of speaker anonymization as the first  
42 step before speech data collection by digital assistant manufacturers and other service providers will not be  
43 possible. The hesitation to adopt anonymization techniques would directly lead to the non-compliance with

---

<sup>8</sup><https://www.intel.com/content/www/us/en/architecture-and-technology/software-guard-extensions.html>

---

1.3 Summary of contributions

5

the recent guidelines [34] put forth by the EDPB which cites the GDPR [52] and the e-Privacy Directive [51] to achieve privacy by design and by default while implementing and introducing virtual voice assistants in the market.

We focus our efforts towards the development of speech transformation or representation learning based techniques to produce anonymous speech representations. This development aligns with our objective of identifying speaker-related information in the speech signal and leads us to explore different representations either in a client-server setting or independently of a fixed speech processing architecture, such as anonymous waveforms. We impose the constraint of being *global* over these transformations so that speaker information is identified and removed across gender, accents, domains or recording conditions. We evaluate the privacy/utility trade-off of these representations in strong attack conditions, and prove their consistency through formal protection guarantees.

The combined goal of preserving usability along with privacy requires us to answer several fundamental questions. What constitutes the speaker’s identity in a speech signal? How to identify and remove it without affecting the usefulness? To what extent does it make the speech signal linkable to the speaker? Can we disentangle speaker information from other attributes in the signal such as emotional states & traits, communicative acts, syntactic content, intonation, etc.? How to confirm the removal of identity among other attributes with high confidence?

Studies on speaker identification research have attributed speaker-related information to the human speech production mechanism [40]. Most of the speaker-related factors arise due to the physiology and shape of the vocal tract. Although the mechanism is well studied by linguists and the speech processing community, it is hard to formulate a global rule-based approach which modifies the speaker information alone and does not affect other attributes. In this thesis we explore machine learning based approaches which enable us to identify and remove speaker-dependent features in the speech signal either in the spectral domain or in learned features such as neural network representations.

Succinctly, the successful implementation of our central goal will translate not only into increased personal data protection but also increased trust by citizens and service providers. It will enable them to satisfy the requirements set by the law and eventually build society’s trust in future private-by-design voice-based applications.

---

**1.3 Summary of contributions**

29

Our main focus is to obtain an anonymous speech representation which conceals the speaker’s identity while retaining the linguistic content such that it can be used for further processing such as linguistic analysis, decoding the content (i.e. ASR) or training an ASR model. This representation of a speech utterance must be unlinkable to the person who spoke it, hence preventing an adversary to perform membership or linkage attacks. Clearly, there are many stakeholders in the process of anonymization, therefore we start by defining our attack model to concretize the goal of the anonymization process and demarcate the roles of the associated stakeholders in Chapter 3. We introduce the three actors affected by anonymization, namely: the speaker, the user and the attacker. In this chapter, we formally define the privacy and utility metrics that we use throughout the thesis and some preliminary experiments with a diverse set of voice transformation algorithms. We shift the paradigm of speaker anonymization evaluation schemes by formulating the premise of attacker’s knowledge and gradually increase this knowledge to establish the idea of continuity from *Ignorant* to *Informed* attackers. All the subsequent representations are subjected to the rigorous evaluation regime proposed in this chapter and their performance is reported as per the established metrics.

First, we present our experiments with adversarial learning to remove speaker information from the intermediate representation of the ASR network. Following the lead of previous approaches, we assume a

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

1 client-server model, where private information is removed from the signal at the client side and the output  
2 is sent to server for decoding the text. We consider an end-to-end ASR model for this approach and the  
3 intermediate encoder representation is anonymized using speaker-adversarial learning. This approach is  
4 described in Chapter 4. We evaluate the anonymous representations obtained using this approach using  
5 a *Semi-Informed* attacker i.e. an attacker who possesses auxiliary knowledge about the anonymization  
6 mechanism. Such an attacker could easily diminish the strength of a weak anonymization process by learning  
7 the vulnerable discriminative patterns exhibited by the speakers after anonymization. The trend observed in  
8 the representations obtained after adversarial learning was that the anonymity does not generalize to unseen  
9 speakers. Moreover, this representation is tied to a fixed client-server architecture which restricts its usability  
10 to ASR by a fixed decoder server.

11 Given the rigid shortcomings of the client-server model and the intent to generalize the usability of  
12 anonymized representation to any arbitrary downstream task, we explore voice conversion techniques whose  
13 output is a speech waveform, i.e., an intelligible speech signal. More specifically, we experiment with x-vector  
14 based speaker anonymization where the speaker’s identity is assumed to be perfectly disentangled from  
15 other factors of variation, such as the linguistic content and intonation, and concentrated only in the x-  
16 vector component. We extend the original idea and propose a baseline for the first VoicePrivacy challenge  
17 with flexible design choices to select the target identity for the source speaker. This approach is explained  
18 in Chapter 5. We conduct extensive experiments with this approach and establish the superior privacy  
19 protection and utility achieved by it against *Semi-Informed* attackers, even in presence of thousands of  
20 speakers.

21 Further experiments with the x-vector based anonymization framework reveal that the assumption of  
22 perfect disentanglement does not hold true, and in practice it is far from being perfect. The linguistic features  
23 and the prosodic pattern indeed retain some residual speaker information which makes the synthesized  
24 speech linkable to the original speaker even after anonymization. Moving forward, we drop the assumption  
25 of perfect disentanglement in x-vector based anonymization and undertake the task to measure the residual  
26 speaker information that might be present in the features extracted to represent the linguistic content and  
27 the prosodic pattern. We represent the linguistic content of an utterance using the bottleneck features  
28 (intermediate layers) extracted from the ASR network that can efficiently decode the textual content present  
29 in the utterance. We assume that such bottleneck features capture the relevant phonetic information in the  
30 signal. The prosodic pattern is represented by the pitch (or fundamental frequency) contour which also  
31 demarcates the voiced-unvoiced regions of speech. We add Laplace noise to these features to make them  
32 differentially private and show that the anonymization can be improved by further removal of biometric  
33 identity from all the features. This approach is investigated in Chapter 6, where we first measure the individual  
34 contribution of each feature towards speaker’s identity and then exhibit superior privacy protection by  
35 formal noise addition.

36 Finally, we conduct an in-depth analysis and comprehensive experiments to explore the usability of the  
37 anonymized speech data in Chapter 7. We first study the impact of re-training ASR models with anonymized  
38 data to ascertain whether they can demonstrate comparable performance to the baseline model. Then we  
39 apply training approaches such as data augmentation, improved neural network architecture and model  
40 adaptation to exhibit that state-of-the-art acoustic models can be trained for ASR without requiring large  
41 scale un-anonymized (original) data. Such investigation repudiates the claim that original (untransformed)  
42 data is needed for training ASR systems.

43 To summarize, the contributions of this thesis are the following:

- 44 1. We define the attack model associated with privacy threats to speech interfaces. The three actors (the  
45 speaker, the user and the attacker) who are concerned with the anonymization task are mentioned and  
46 the anonymization techniques are evaluated from each of their perspective. These roles are analogous

## 1.4 Publications

7

- to real-world entities and the assumption of the knowledge they possess substantiates their capacity during the evaluation of anonymization schemes. 1
- 2
2. We propose a new regime of evaluation for the speaker anonymization methods using the idea of *Informed* attacker. The attacker may possess auxiliary knowledge of the anonymization algorithm based on which he/she can design effective linkage functions to discover the true identity of the speaker. We establish the feasibility of such attacks by simulating several attackers with varying degree of knowledge about the anonymization scheme and show that previous studies have used an inferior model to evaluate their algorithms, namely the *Ignorant* attacker. 3
- 4
- 5
- 6
- 7
- 8
3. We conduct a comprehensive study of design choices for x-vector based speaker anonymization to select the target *pseudo-speaker* from an external pool of identities. This technique was also proposed as a strong baseline for the first VoicePrivacy challenge. We perform a complete analysis of the privacy/utility trade-off of each design choice in different attack scenarios as well as measure the sustainability of the best combination against re-identification when the attacker possesses the data of thousands of speakers. 9
- 10
- 11
- 12
- 13
- 14
4. We investigate the residual speaker-related information in the pitch contour and the bottleneck (BN) features and show that it can be removed through transformations or differentially-private noise addition. We propose a new percentile based transformation to transform the source pitch to target *pseudo-speaker* pitch. We also propose an adversarial learning based transformation for BN features. Additionally, we put forth new neural network architectures for adding differentially-private noise to pitch as well as BN features and measure its impact of the privacy/utility trade-off. 15
- 16
- 17
- 18
- 19
- 20
5. We study the claim that the anonymized speech corpus is usable and propose techniques to train a viable ASR model which performs equally well for original and anonymized evaluation sets. Concretely, we explore data augmentation, model adaptation and transfer learning based methods to minimize the use of original data and generalize the performance of ASR models using a large amount of anonymized speech data. 21
- 22
- 23
- 24
- 25

## 1.4 Publications

26

## Publications as first author:

27

1. Brij Mohan Lal Srivastava, Aurélien Bellet, Marc Tommasi and Emmanuel Vincent. “Privacy-preserving adversarial representation learning in ASR: reality or illusion?” *In Proc. Interspeech*, pp. 3700–3704, 2019. 28
- 29
- 30
2. Brij Mohan Lal Srivastava, Nathalie Vauquier, Md Sahidullah, Aurélien Bellet, Marc Tommasi and Emmanuel Vincent. “Evaluating voice conversion-based privacy protection against informed attackers” *In Proc. 45th International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pp. 2802–2806, 2020. 31
- 32
- 33
- 34
3. Brij Mohan Lal Srivastava, Natalia Tomashenko, Xin Wang, Emmanuel Vincent, Junichi Yamagishi, Mohamed Maouche, Aurélien Bellet and Marc Tommasi. “Design choices for x-vector based speaker anonymization” *In Proc. Interspeech*, pp. 1713–1717, 2020. 35
- 36
- 37
4. Brij Mohan Lal Srivastava, Mohamed Maouche, Md Sahidullah, Emmanuel Vincent, Aurélien Bellet, Marc Tommasi, Natalia Tomashenko, Xin Wang and Junichi Yamagishi. “Privacy and utility of 38
- 39

1       x-vector based speaker anonymization” *Submitted to IEEE/ACM Transactions on Audio, Speech, and*  
2       *Language Processing*, 2021.

3       **Other publications:**

- 4       1. Mohamed Maouche, Brij Mohan Lal Srivastava, Nathalie Vauquier, Aurélien Bellet, Marc Tommasi  
5       and Emmanuel Vincent. “A comparative study of speech anonymization metrics” *In Proc. Interspeech*,  
6       pp. 1708–1712, 2020.
- 7       2. Natalia Tomashenko, Brij Mohan Lal Srivastava, Xin Wang, Emmanuel Vincent, Andreas Nautsch,  
8       Junichi Yamagishi, Nicholas Evans, Jose Patino, J.-F Bonastre, Paul-Gauthier Noé and Massimiliano  
9       Todisco. “Introducing the VoicePrivacy initiative” *In Proc. Interspeech*, pp. 1693–1697, 2020.
- 10      3. Ali Shahin Shamsabadi, Brij Mohan Lal Srivastava, Aurélien Bellet, Nathalie Vauquier, Emmanuel  
11      Vincent, Mohamed Maouche, Marc Tommasi, Nicolas Papernot. “Differentially private speaker  
12      anonymization” *Accepted at 43rd IEEE Symposium on Security and Privacy (IEEE S&P)*, 2022.

13      **1.5 Thesis structure**

14      This thesis lies at the crossroads of speech processing, privacy and machine learning. In **Chapter 2**, we  
15      review the fundamentals and principles of these domains which are relevant for our work. First, we describe  
16      the established tools of speech processing, such as speech recognition and speaker identification. Then we  
17      review the existing methods of speaker anonymization which are close predecessors of our techniques. We  
18      formally define the tasks for which we propose the potential solutions along with the metrics that we report  
19      in order to evaluate their effectiveness. We further recapitulate some key technologies that we employ while  
20      proposing our solutions, such as voice transformation (VT) and voice conversion (VC).

21      In **Chapter 3**, we define the attack model and the actors involved in the process of anonymization. We  
22      firmly establish the concept of attacker’s knowledge through preliminary experiments with voice conversion  
23      based techniques against *Informed* attackers. We briefly present a comparative analysis of the metrics that  
24      are used for privacy evaluation.

25      **Chapter 4** describes our effort to learn an anonymous representation of speech using adversarial learning  
26      which can be processed locally and then transmitted to a server for decoding. We evaluate the privacy protection  
27      achieved by this representation using closet-set speaker identification as well as open-set speaker verification  
28      against an *Informed* attacker. **Chapter 5** introduces the VoicePrivacy initiative and the x-vector based  
29      speaker anonymization framework. We explore several design choices associated with it to choose a robust  
30      target *pseudo-speaker* x-vector. We conduct an extensive evaluation of the representations obtained using this  
31      technique in strong attack scenarios as well as against re-identification attacks in the presence of thousands  
32      of speakers.

33      We further propose to measure and remove the residual speaker-related information in the inputs of  
34      the x-vector based speaker anonymization, i.e., the BN features and the pitch contour. These analyses are  
35      mentioned in **Chapter 6**. We add differentially-private noise in these features and measure the privacy/utility  
36      trade-off.

37      We conduct a thorough investigation of the usability of the anonymized speech corpus for training a  
38      viable ASR model in **Chapter 7**. We propose data augmentation and transfer learning based approaches to  
39      show that a reasonable ASR model can be trained with a minimum amount of original data.

40      Finally, we summarize our contributions and reflect some of our perspectives towards the future direc-  
41      tions opened up by the research done in the course of this thesis. **Chapter 8** concludes the thesis with these  
42      reflections.

# Chapter 2

## Background and Related Work

If I have seen further it is by standing on the shoulders of Giants.

---

*Isaac Newton*

There is abundant research material present in each of the domains relevant for this thesis. It uses previous knowledge from speech processing, privacy, and machine learning, and their associated sub-domains. In this chapter, we present selected background details of the sub-domains which are relevant to the techniques and terminology proposed in this thesis, and use the existing literature to elucidate the full picture of the problem statement and how the proposed solutions are perceived by the speech and privacy communities. We start by giving a brief account of the key historical advancements and current perspectives in speech processing and digital assistant technology that led to the crisis of privacy. We then describe the principles and tools of speech processing that will help the reader understand the basic terminology we use in the course of this thesis. Next, we present a brief review of the relevant literature which describes the previously proposed machine learning-based anonymization methods that are closely related to our work. Thereafter we formally define the task of privacy-preserving speech processing as presented in previous studies and how it was traditionally evaluated in terms of privacy and utility. Finally, we give a detailed explanation of the core techniques used in our proposed solutions, such as adversarial learning, voice transformation, and voice conversion.

### 2.1 A brief historical overview of speech processing and privacy

Speech processing came a long way since 1881 when the earliest device for recording speech was invented by Alexander Graham Bell. It used a rotating cylinder coated with wax over which up-and-down grooves could be cut by a stylus responding to the acoustic pressure generated by the sound wave. One can only imagine the tremendous challenges posed by this device to record, process, and store speech signals. Thankfully it has been replaced by microphones which capture the acoustic pressure from sound waves and record it as a relative change in voltage. There are several such historical advancements in speech technology that facilitated convenient and large-scale speech processing, eventually leading to the current privacy crisis. Particularly, Homer Dudley's work [72] inspired several generations of researchers to focus on making speech the mainstream medium for human-computer interaction, which propelled the large-scale storage of speech data and overall, the domain of speech signal processing forward.

1 Recall the “speaking machine” invented by von Kempelen introduced in Section 1.1 which could produce  
2 a few human-like sounds. In the mid-1800s, Sir Charles Wheatstone improved upon its design [313] using  
3 adjustable and configurable leather resonators capable of producing many more speech-like sounds. This  
4 model was adopted by Homer Dudley to design an electrical speech synthesizer [70] for Bell Labs. The  
5 synthesizer could be operated as a piano with hand controls to switch between voiced and unvoiced sounds,  
6 keys to control the characteristics of the signal and a foot pedal to control the pitch. It was called the  
7 VODER (Voice Operation Demonstrator) and was first demonstrated at the New York World’s Fair in  
8 1939. This event attracted the focus of researchers worldwide leading to several speech interest groups in  
9 the community. Dudley also pioneered the field of speech coding [263] which aims to represent speech  
10 signals for efficient storage and transmission by exploiting their inherent redundancies. He provided the  
11 analysis-synthesis [71, 73] method for speech coding.

12 The initial usage of speech technology was predominantly envisioned in a controlled setting, such as  
13 offices and research labs, where storage is limited, and through experience and training the people being  
14 recorded gradually became cautious to not divulge private information in the collected data. The recent  
15 advances that have led speech interfaces to enter our homes at the consumer level are quite new, and the  
16 privacy-related implications of this technology are still being explored. As of today, speech interfaces are  
17 present in personal mobile phones as well as digital assistants which have a widespread consumer base.  
18 Exposing an unaware user to such advanced technology will open the doors for potential adversaries to  
19 exploit the sensitive attributes present in the speech signal.

20 Several researchers have studied the security and privacy vulnerabilities of digital assistants [80, 162, 88],  
21 and their third-party applications [167]. The two most concerning privacy issues are the “always listening”  
22 feature and the cloud storage of the audio queries. The device remains in the inert state of buffering and  
23 re-recording until the wake word is spotted [133], it then records the audio and sends it to a cloud-based  
24 service for ASR and natural language understanding (NLU). All the audio files are usually stored in the  
25 user’s account and can be accessed by logging into the account. This data may contain sensitive details about  
26 the user’s life, such as bank details. A compromised account can lead to a user’s private speech data being  
27 leaked to the public. Due to the rich nature of speech signal as we described earlier, not only the linguistic  
28 content but many other attributes of the speaker may become known to a malicious entity.

29 Extensive surveys of digital assistant users have been conducted to understand their mental models,  
30 beliefs, attitudes, and concerns towards their devices. Some studies [4, 162] show that users have an incorrect  
31 understanding of the working of digital assistants and the third-party services with which their sensitive data  
32 is shared. They are also unaware of the existing privacy controls in the digital assistant architecture. Malkin  
33 et al. [188] show that half of the users are not aware of the permanent retention policy of audio queries in  
34 the user’s account. Users are not aware of existing privacy features and they express the need for automatic  
35 deletion of their recordings. Huang et al. [130] studied users’ behaviour and privacy concerns when a digital  
36 assistant is shared among several housemates. Bispham et al. [32] present a taxonomy of attacks on speech  
37 interfaces which motivates future research on voice privacy to focus on exact vulnerabilities present in such  
38 devices. These surveys make some recommendations to users and manufacturers such as turning off the  
39 microphone when not in use, updating the firmware with the latest release, strict data deletion policies, and  
40 screening of sensitive content.

41 The above studies are indicative of the fact that users of speech interfaces are gradually becoming more  
42 aware of the underlying mechanisms and more concerned about their privacy being leaked through the  
43 interface. In this thesis, we aim to propose speaker anonymization techniques that will protect users’ identity  
44 at the source, without requiring them to put in significant effort. These techniques can be built in directly  
45 into the device firmware by the manufacturers.

## 2.2 Principles and tools of speech processing

Now let us introduce some basic principles and tools of speech processing behind the proposed methods. We start with the basics of speech as a signal, and how it is processed to extract relevant features with physiological and phonetic considerations. We give a brief account of artificial neural networks due to their pervasive use as statistical models in speech processing tasks. Then, we describe the technology behind the three most popular speech applications that enable the design and evaluation of our proposed methods: automatic speech recognition, speech synthesis, and automatic speaker recognition.

### 2.2.1 Fundamentals of speech processing

In this section, we briefly discuss the mechanism of human speech production followed by its representation and processing as a discrete-time signal.

**Vocal tract.** The physiological apparatus that generates speech is called the *vocal tract* [115], which starts at the lungs and ends at the lips and the nostrils. The larynx (also called the voice box) separates the vocal tract into two anatomical regions: the lower part is called the sublaryngeal region and the upper part is called the supralaryngeal region. The sublaryngeal region of the vocal tract is composed of the diaphragm, the lungs, and the trachea (also called the windpipe). The air flows outward from the lungs and encounters a pair of flap-like structures in the larynx, called the *vocal folds*. When the vocal folds are held at an intermediate tension so that they are not too close or too far apart, the movement of the air induces ripples along their length. This causes them to vibrate, and the result is *voicing*. Voicing is the cause of periodic segments in the speech signal which are called *voiced* regions. On the contrary, when the vocal folds are held at sufficient distance from each other so that air flows freely through them, they do not vibrate, which results in voicelessness. This can be observed in the speech signal as aperiodic segments which look like random noise and are known as *unvoiced* regions.

The supralaryngeal region, which is composed of the oral cavity and the nasal cavity, plays a major role in determining the exact nature and quality of the sounds that are produced. The different parts of the supralaryngeal region that contribute towards the articulation of different vowels and consonants are referred to as articulators. The major articulators in the oral cavity are the lips, the teeth, the tongue, the alveolar ridge, the hard palate, and the velum. Among these, the tongue and the lower lip are the active articulators, whereas the others are passive and immobile. The complex interaction between active and passive articulators to completely stop the airflow, constrict it through a narrow channel, or allow it to pass through without restriction gives us the vast variety of speech sounds found in all of the world's languages.

**Phonemes.** The vocal tract is a continuous system capable of producing infinitely many sounds. These sounds are called *phones*. The exact physical mechanism of producing phones by the vocal tract, their transmission in acoustic space, and their auditory perception by the human ear are studied under the branch of linguistics called *phonetics* [115], which is independent of language. A given language can have only a small, finite number of sound units that can be used to compose words in that language and have some grammatical significance. These sounds must be perceptually distinct from each other for effortless communication and are called *phonemes*. The organization of phonemes, their combinations to produce words, and their semantic role in language are studied under the branch of linguistics called *phonology* [41]. Phonology categorizes the continuous signal produced by the vocal tract into discrete phoneme classes based on their acoustic, articulatory, and perceptual characteristics. Most languages feature two broad classes of phonemes, namely *vowels*, that are voiced sounds produced with no obstruction by the articulators, and *consonants*, that are produced by obstructing the airflow passing through the vocal tract. Although every language has a different

<sup>1</sup> set of phonemes, the International Phonetic Alphabet (IPA) [16] describes the universal set of phonemes  
<sup>2</sup> based on their articulatory characteristics.

<sup>3</sup> Vowels are described based on the position of the tongue and the roundedness of the lips. The tongue is  
<sup>4</sup> a highly active articulator and is subdivided into the front, central and back parts which can move somewhat  
<sup>5</sup> independently of each other. It can also be placed at different heights to control the width of the constriction  
<sup>6</sup> in the vocal tract. For example, /i/ as in “feed” is made by placing the front part of the tongue close to the  
<sup>7</sup> hard palate, hence it is categorized as a close front vowel without rounding, while /o/ as in “foe” is made  
<sup>8</sup> by raising the back of the tongue up to a certain height and rounding the lips, hence it is a close-mid back  
<sup>9</sup> vowel with rounding. Consonants are categorized based on the presence or absence of voicing, the place  
<sup>10</sup> of articulation that indicates the place of constriction in the vocal tract, and the manner of articulation  
<sup>11</sup> which is the method of air release. For instance, /p/ as in “pan”, is a voiceless consonant made by completely  
<sup>12</sup> blocking the airflow using the lips, hence it is categorized as a voiceless bilabial plosive, whereas /z/ as in  
<sup>13</sup> “zoo”, is a voiced consonant produced by making a narrow constriction by placing the tip of the tongue close  
<sup>14</sup> to the alveolar ridge, therefore it is a voiced alveolar fricative.

<sup>15</sup> Such categorization of phonemes also helps us understand the linguistic behaviour of speakers when a  
<sup>16</sup> sound is missing in their language [151, 159]. Generally, the non-native speakers retain voicing and manner,  
<sup>17</sup> but replace the place of articulation, for example, the sound of consonant /ð/ as in the English word “the”  
<sup>18</sup> is a voiced *dental* fricative that is not available in the French language, hence most native French speakers  
<sup>19</sup> replace it with /z/ which is a voiced *alveolar* fricative [152]. Similarly, some dialects of Hindi do not have  
<sup>20</sup> the phoneme /ʃ/ as in “sheep” which is a voiceless postalveolar fricative, hence they replace it with /s/ as in  
<sup>21</sup> “sun” that is a voiceless alveolar fricative.

<sup>22</sup> **Speech in the time domain.** Sound is a pressure wave traveling through the air as the medium of propagation. It can be recorded by measuring the variation in pressure at a single point in space over time. As  
<sup>23</sup> mentioned before, a microphone is used to record the acoustic wave which measures the relative change in  
<sup>24</sup> pressure as the electrical signal that is proportional to the pressure variation. Figure 2.1(a) shows the output  
<sup>25</sup> of the microphone (for the word “privacy”), also called a waveform or a time-domain signal, pronounced  
<sup>26</sup> by a male or a female speaker. The *duration* of both waveforms is shorter than one second. To represent  
<sup>27</sup> a speech signal digitally, we must select the *bit depth* which is the finite precision needed to encode the  
<sup>28</sup> amplitude values, and the *sampling rate* (denoted as  $F_s$ ) which defines how many times per second the  
<sup>29</sup> actual waveform is sampled to obtain discrete values of the amplitude. The duration, the bit depth, and  
<sup>30</sup> the sampling rate decide the memory requirement to store the audio file. The audio file can be stored in a  
<sup>31</sup> lossless uncompressed format (e.g., “.wav”) or a lossy compressed format where the file size is reduced while  
<sup>32</sup> maintaining good audibility (e.g., “.mp3”).

<sup>34</sup> A discrete-time speech signal can be represented as  $s$  and  $s[n]$  denotes a single sample of instantaneous  
<sup>35</sup> amplitude value, where  $n = 0, \dots, N_s - 1$ . As described further, the speech signal is generally analyzed to  
<sup>36</sup> determine its frequency components in a short duration.

<sup>37</sup> **Short-term analysis.** A spoken sentence is also called an *utterance* which is a sequence of phonemes  
<sup>38</sup> (note the phoneme annotations in Figure 2.1(a)). Depending on the recording conditions, whether the  
<sup>39</sup> speaker is reading a given text or engaged in a spontaneous conversation, the utterance may or may not be  
<sup>40</sup> grammatically correct. In any given utterance, except for global utterance-level information such as duration  
<sup>41</sup> or speaker-related characteristics, other properties of a speech signal, like amplitude, voicing, etc. vary over  
<sup>42</sup> time. We also know that different sounds are produced by different configurations of the articulators, so  
<sup>43</sup> the system producing the signal itself is changing over time. Hence, in order to process the speech signal,  
<sup>44</sup> it is divided into uniform regions called *time frames* that are individually analyzed. The speech signal  $s$  is

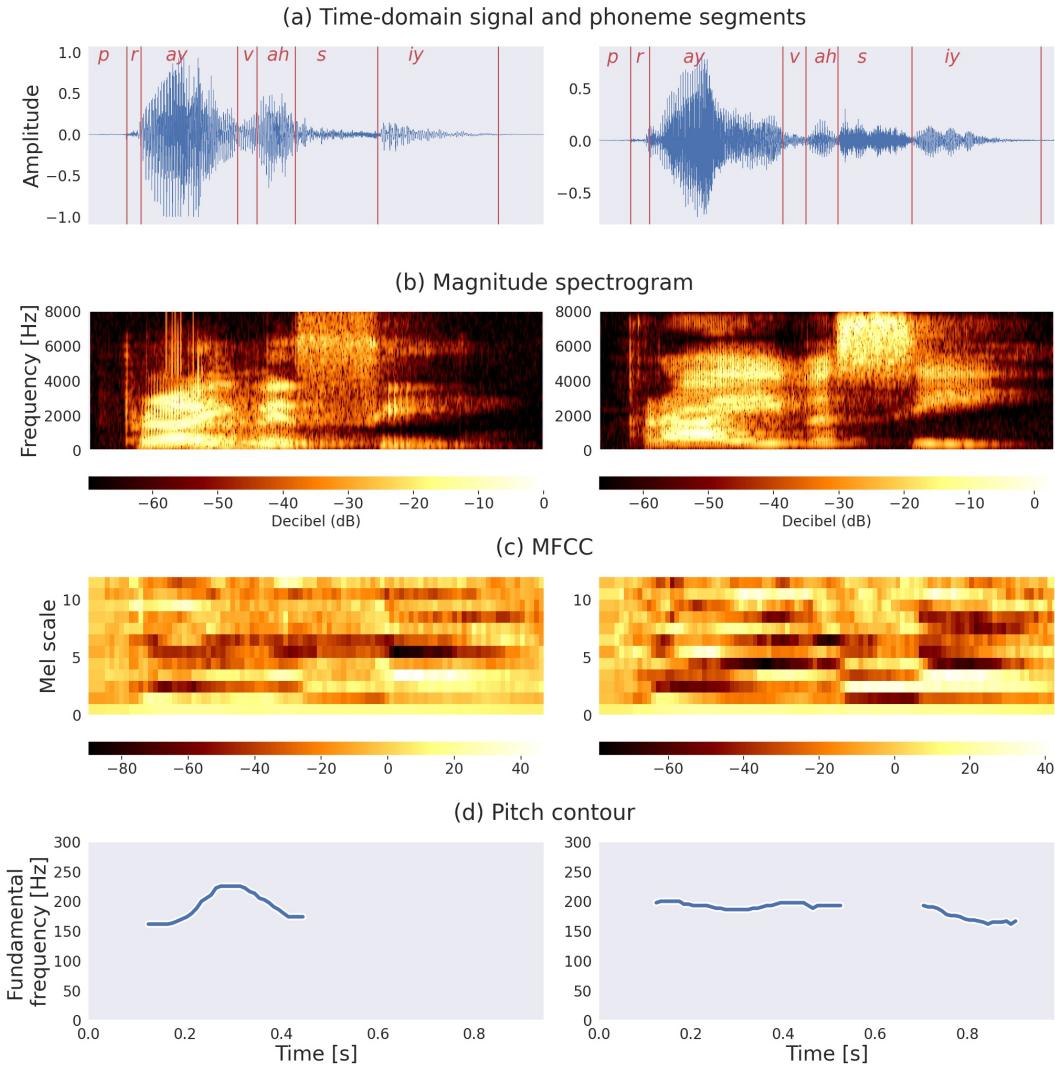


Fig. 2.1 Waveform, magnitude spectrogram, MFCC and pitch contour for the word “privacy”, pronounced by a male (left) and a female speaker (right).

divided into  $T$  frames (subsequences of  $\mathbf{s}$ ) each of length  $L$  samples with an overlap of  $L/2$  samples, thereby obtaining  $[\vec{s}_0, \dots, \vec{s}_{T-1}]$ , where for each  $t \in [0, \dots, T-1]$

$$\vec{s}_t = \left[ \mathbf{s} \left[ t \cdot \frac{L}{2} + l \right] \right]_{l=0}^{L-1}. \quad (2.1)$$

The time frames  $\vec{s}_t$  are multiplied with a short window  $\Psi$ , of the same length as the time frame, i.e.,  $L$ . The value of  $L$  is typically 25 to 30 ms. It is assumed that the system (vocal tract shape) remains stationary over the duration of the window, so that the signal’s spectral properties are constant in this region. The window function [214] is not rectangular but tapering at the beginning and end, such as the Hamming window:

$$\Psi[l] = 0.54 - 0.46 \cos\left(\frac{l\pi}{L}\right), \quad (2.2)$$

1 to avoid introducing artifacts in the original signal. Due to the tapering window, we might lose some information present in the signal during analysis, hence a small overlap ( $L/2$ ) is introduced between consecutive frames. The analysis of the speech signal after splitting it into a sequence of overlapping frames is called short-term analysis.

5 **Speech in the frequency domain.** By analysis, it is implied that we want to determine the frequency  
6 content of a particular frame of the speech signal. This can be done by faithfully reconstructing the signal  
7 as the weighted sum of simple sine waves called the basis functions. The basis functions are orthogonal to  
8 each other, i.e., no energy at the frequency of one sine wave is present in another. This property helps to get  
9 a unique solution for the coefficients of the weighted sum. Each basis function with unit amplitude and  
10 only a single fixed frequency is correlated with the signal to determine the exact magnitude of this frequency  
11 present in the signal. The frequency of the basis functions ranges from the lowest value, where a single cycle  
12 of the sine wave fits the entire analysis frame, to the highest possible frequency which is half of the sampling  
13 rate, also called the *Nyquist frequency* ( $= F_s/2$ ). The process of computing the frequency content of the  
14 original time-domain signal is called the *Fourier transform*. The discrete Fourier transform of the  $t$ -th time  
15 frame  $\vec{s}_t$  is computed to extract the Fourier coefficients  $\mathcal{F}_t$  for the  $k$ -th frequency component:

$$16 \quad \mathcal{F}_t[k] = \sum_{l=0}^{L-1} \vec{s}_t[l] \cdot \Psi[l] \cdot e^{-j \frac{2\pi}{L} kl}, \quad 0 \leq k \leq L - 1. \quad (2.3)$$

17 The value of  $k$  corresponds to the frequency bin center  $F(k) = kF_s/L$  in Hz (with zero at the start).  
18 The frequency bin centers are used to arrange the Fourier coefficients as meaningful frequency-domain  
19 representations, like a spectrum, and are also used to derive the features that warp the frequency axis, as  
20 described in Equation (2.5). Since  $e^{j\omega} = \cos \omega + j \sin \omega$ ,  $\mathcal{F}_t[k]$  is a complex number for each of the  $L$   
21 frequency bands which encode the signal's magnitude and phase. This process can be sped up by the fast  
22 Fourier transform (FFT) [212] algorithm when  $L$  is a power of 2.

23 The vector of Fourier coefficients ( $\mathcal{F}_t$ ) of a given frame is called a *spectrum* which has the size of  $L \times 1$ .  
24 For analysis purposes, it is common practice to record only the magnitude  $|\mathcal{F}_t[k]|$  of the coefficients and  
25 discard their phase. The resulting vector is called the magnitude spectrum. Stacking the magnitude spectra  
26 of all frames results in a 2D representation of the whole utterance called the *magnitude spectrogram* (denoted  
27 as  $|\mathcal{F}|$ ) which has the size of  $L \times T$ . The dynamic range is generally compressed by expressing the magnitude  
28 on a logarithmic scale called decibels (dB). Figure 2.1(b) depicts the magnitude spectrogram of the given  
29 speech signal.

30 **Features.** Frequency-domain representation has proven to be very powerful in order to inspect the properties of speech sounds. For example, the smooth curve that follows the peaks of the spectrum for any given analysis frame is called the *spectral envelope*, and it is governed by the shape of the vocal tract. The dominant peaks corresponding to the resonant frequencies in the spectral envelope are called *formants*. Each phoneme is characterized by specific spectral properties which can be used as a template to recognize it [119]. Using the spectrum or the spectral envelope directly for speech recognition may however not be optimal due to the inherent covariance between frequency bands, and the range of magnitude which does not linearly correspond to loudness.

38 To alleviate these shortcomings, researchers have proposed several transformations of the spectrum such  
39 that the resulting features correspond to a compressed representation which is motivated by the perceptual  
40 mechanism of the human ear [137]. It is well known that humans can perceive sound within a defined  
41 frequency range of 20 Hz to 20 kHz. The human auditory system is more discriminative between tones at

lower frequencies and increasingly less discriminative at higher frequencies [173]. The sensitivity to higher frequencies also reduces as we age. Hence speech signals are generally processed at the sampling rate of 16 kHz or lower, limiting the information to 8 kHz based on the Nyquist frequency. After the FFT, the magnitude spectrum  $|\mathcal{F}_t|$  is obtained for each frame, which is then warped according to a non-linear perceptual scale called the *Mel scale*. The linear frequency  $F$  (in Hz) can be converted to Mel scale using the following formula:

$$\text{Mel}(F) = 1127 \cdot \ln(1 + F/700). \quad (2.4)$$

The Mel scale aims to mimic the sensitivity of the human auditory system by warping the frequency scale using closely-spaced narrow bandpass filters at lower frequencies, and increasingly wider and sparsely-spaced filters at higher frequencies.<sup>1</sup> The Mel scale filters capture the general shape of the spectral envelope needed for speech recognition and smoothes out the harmonics,<sup>2</sup> thereby losing the fundamental frequency information. The warped magnitude spectrum is segmented into frequency bands according to a Mel filter bank which consists of a fixed number of overlapping triangular filters, typically 40 to 80, defined by their center frequencies  $F_c(m)$  on the linear scale. The Mel filter bank is parameterized by the number of filters  $N_M$ , the minimum frequency  $F_{\min}$ , and the maximum frequency  $F_{\max}$ . The center frequencies of the Mel filters are the integer multiples of the fixed frequency resolution  $\delta_{\text{Mel}}$  in the Mel scale which is computed using  $\delta_{\text{Mel}} = (\text{Mel}(F_{\max}) - \text{Mel}(F_{\min}))/N_M$ . Hence, the center frequencies are given by  $\text{Mel}(F_c(m)) = m \cdot \delta_{\text{Mel}}$  for  $m = 1, \dots, N_M$ . The center frequencies of the triangular filters are converted to the linear scale using the inverse mapping:  $F_c(m) = 700 \cdot (e^{\text{Mel}(F_c(m))/1127} - 1)$ . The Mel filter bank  $M(m, k)$  is given by [154]:

$$M(m, k) = \begin{cases} 0 & \text{for } F(k) < F_c(m-1), \\ \frac{F(k)-F_c(m-1)}{F_c(m)-F_c(m-1)} & \text{for } F_c(m-1) \leq F(k) < F_c(m), \\ \frac{F(k)-F_c(m+1)}{F_c(m)-F_c(m+1)} & \text{for } F_c(m) \leq F(k) < F_c(m+1), \\ 0 & \text{for } F(k) \geq F_c(m+1). \end{cases} \quad (2.5)$$

The Mel filter bank  $M(m, k)$  is a matrix of size  $N_M \times L$  which, when multiplied by the power spectrum (i.e., squared magnitude spectrum), yields a set of coefficients called the *Mel-filterbank coefficients*. To further enhance their usability, a logarithm is applied to compress the dynamic range such that it is more directly related to the perceptual loudness. The resulting logmel coefficients are sometimes directly used for speech recognition [105] and synthesis [139]:

$$\text{logmel}_t(m, k) = \ln \left\{ \sum_{k=0}^{L-1} M(m, k) \cdot |\mathcal{F}_t[k]|^2 \right\}. \quad (2.6)$$

Finally, the discrete cosine transform (DCT) can be applied to these coefficients to approximately de-correlate them from each other and obtain *Mel-frequency cepstral coefficients* (MFCCs) as shown in Figure 2.1(c). MFCCs are widely used in speech applications. Note that logmel or MFCC are real-valued vectors obtained per frame: they are often concatenated over time to get the whole Mel spectrogram or MFCC sequence for an utterance.

<sup>1</sup>The Mel scale is the most popular perceptual scale, but there are other such scales like the Bark scale.

<sup>2</sup>A periodic signal with frequency  $F$  is only composed of the frequencies that are integer multiples of  $F$ , i.e.,  $F, 2F, 3F$ , etc. These frequencies are called harmonics.

**1 Pitch.** An important property of speech is the presence of *pitch* in the voiced regions. Strictly speaking,  
 2 pitch is the perceptual property that relates to the rising-falling tonal pattern, or the intonation of speech. It  
 3 highly correlates with a physical property of the speech signal, called the fundamental frequency (denoted as  
 4  $F_0$ ), which is the rate of vibration of the vocal folds. The range of pitch is determined by the physiological  
 5 factors of the vocal folds, such as their mass and length, hence it depends on the speaker and is typically  
 6 lower for male than female [30]. The pitch sequence governs the *intonation* of the spoken utterance, and it  
 7 significantly contributes towards the message that is being conveyed to the listener; for instance, a rising pitch  
 8 at the end of the sentence may convey to the listener that a question is being asked. It is a key component of  
 9 prosody (together with stress and rhythm) which determines the utterance expressiveness, and is crucial  
 10 for speech synthesis. Prosody is a useful tool of communication in language as it indicates the prominence  
 11 of different linguistic units that compose the utterance, and hence contribute towards the naturalness of  
 12 speech. It is important to note that pitch is the rate of vibration of the vocal folds hence it is only defined for  
 13 voiced phonemes such as /a/, /b/, /z/, etc. It is pointless to compute pitch for silence, noise or unvoiced  
 14 regions of an utterance since there is no vibration of the vocal folds, hence by convention the pitch value in  
 15 these regions is set to zero.

**16** Pitch can be estimated from the speech signal, without having physical access to the vocal folds, using  
 17 pitch estimation algorithms [258]. Pitch estimation is a difficult task due to erroneous observation of  
 18 harmonics causing pitch doubling/halving [322]. It is also difficult to estimate the pitch when the quality  
 19 of speech is distorted due to noise or channel effects. A fairly robust and widely used algorithm for pitch  
 20 tracking is called Yet Another Algorithm for Pitch Tracking (YAAAPT) [144]. It is a hybrid pitch tracking  
 21 method as it considers both the time and the frequency domain to estimate the value of  $F_0$ . It comprises  
 22 a nonlinear preprocessing step on the squared speech signal, followed by  $F_0$  estimation using Spectral  
 23 Harmonics Correlation from the spectrogram of the nonlinearly processed signal. A crucial components of  
 24 YAAAPT is the normalized cross-correlation function (NCCF) which is used to extract prominent peaks  
 25 corresponding to  $F_0$  candidates in the time domain. The NCCF for a given time frame and lag-index  $q$  is  
 26 defined as [277]:

$$\text{NCCF}_t(q) = \frac{1}{\sqrt{\xi_0 \xi_q}} \sum_{l=0}^{L-Q} \vec{s}_t[l] \vec{s}_t[l+q]. \quad (2.7)$$

**27** The NCCF $_t$  is computed for  $0 \leq q < Q$ , where the value of maximum lag  $Q$  is generally lesser than  
 28 the frame length  $L$ , and the frame energy is given by  $\xi_q = \sum_{l=q}^{q+L-Q} (\vec{s}_t[l])^2$ . The final pitch contour is  
 29 obtained using dynamic programming and a normalized low frequency energy ratio function is applied to  
 30 make voiced/unvoiced decision. Figure 2.1(d) shows the estimated pitch for the word “privacy” produced  
 31 by a male or a female speaker.

### **33** 2.2.2 Artificial neural networks

**34** At this point, we digress a little bit to explain the core concepts of artificial neural networks (ANNs) and  
 35 deep learning, which form the key components of modern speech processing tasks, as we will see later in this  
 36 chapter. ANNs are inspired by the working and structure of the biological neural network present in the  
 37 brain. The idea of ANNs emerged from the *connectionist* [97] school of cognitive science which hopes to  
 38 simulate human intelligence through a large network of connections between the neurons. Neurons are the  
 39 smallest unit in ANNs and are represented as nodes in this large computational graph. Actions are triggered  
 40 when a specific combination of neurons are fired together. ANNs derive their power from the general  
 41 framework proposed in the seminal work of parallel and distributed processing [242], which describes  
 42 the parallel nature of neural information processing and the distributed nature of neural representation  
 43 in ANNs, which are similar to how the brain processes and stores information. Although we are still not

aware of the complete working of the human brain, storage of memories, production of thoughts, etc., we know that there are special regions for processing different types of sensory inputs, such as the visual and auditory cortex. Thus, ANNs are typically used to learn specific tasks, such as object detection or phoneme recognition, which they can accomplish quickly and sometimes more precisely than human beings, rather than designing a general-purpose intelligent machine, like the brain.

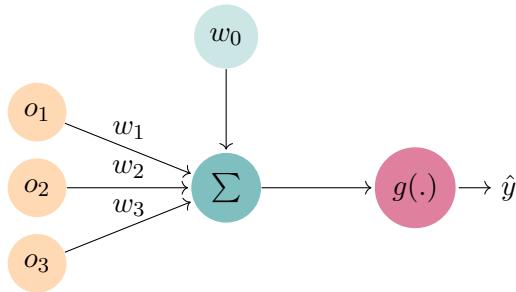


Fig. 2.2 Perceptron model of a neuron with  $A = 3$ .

**Structure** As mentioned before, a neuron is the fundamental unit of an ANN, sometimes also called a *perceptron* model [239]. It is nothing but the weighted sum of its inputs, which is transformed by a non-linear *activation function*, denoted by  $g(\cdot)$  as shown in Figure 2.2. The activation function decides whether the neuron fires or not based on the cumulative influence of the input features and the activation threshold, thereby filtering the information that is passed on to the output. The activation function should be non-linear so that the neuron can learn the complicated non-linear relationship between the input and the output. It should also be differentiable so that the gradients of the error can be computed and *backpropagated* to optimize the model parameters,  $\theta$ . The output activation of a perceptron can be simply written as:

$$\begin{aligned}\hat{y} &= g \left( w_0 + \sum_{i=1}^A w_i o_i \right) \\ &= g(\theta^T \mathbf{o}).\end{aligned}\tag{2.8}$$

Here, we define  $\mathbf{o} = [1, o_1, \dots, o_A]$  as a single sample (observation) from the data set containing  $A$  features and an additional 1, and the parameters  $\theta = \{w_0, w_1, \dots, w_A\}$  include one synaptic weight  $w_i$  for each feature and a bias  $w_0$  to ensure that the decision boundary isn't fixed at the origin. In real-world applications, we encounter complex multi-class problems such as speaker identification, phoneme recognition, etc. that require better expressivity and the ability to learn complex non-linear mappings/representations. Therefore, in practice we use neural networks with a more complicated architecture than a perceptron and several interconnections that exist between millions of neurons represented by the weights and biases. The input is propagated forward sequentially through the consecutive layers, where each *layer* contains multiple neurons. This is referred to as *forward propagation*, which is used to compute the output of the neural network.

A *fully-connected* network involves directed connections from each neuron in the current layer to every neuron in the subsequent layer as shown in Figure 2.3. It is sometimes also referred to as a *multilayer perceptron*. The first layer which receives the features directly is called the *input layer*, and the last layer is called the *output layer*. The remaining layers in between are called *hidden layers*. The neural network shown in Figure 2.3 includes two hidden layers, where  $h_j^{(k)}$  represents the activation value for the  $j$ -th neuron in

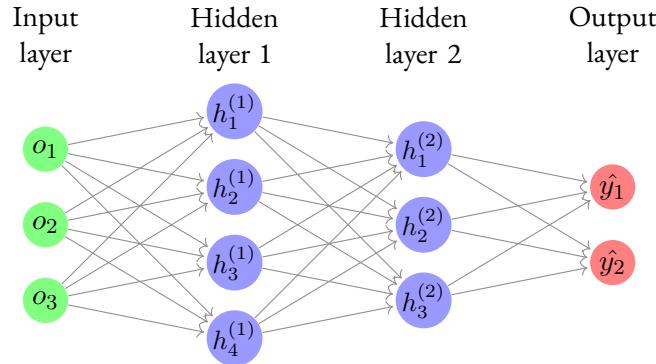


Fig. 2.3 Fully-connected feed-forward neural network (multilayer perceptron).

the  $k$ -th layer. Although it has been proved that a neural network with a single hidden layer and enough neurons can approximate any computable function [57], it is more costly to add neurons in a single hidden layer than to add more hidden layers. Moreover, it has been repeatedly shown that the sequential hidden layers learn representations of data with multiple levels of abstraction [163], which gave rise to the field of *deep learning*.

**Training** Let there be input samples  $\{\mathbf{o}_i\}_{i=1}^N$  in a given data set and the corresponding desired outputs  $\{y_i\}_{i=1}^N$  where  $\mathbf{o}_i \in \mathcal{X}$  and  $y_i \in \mathcal{Y}$ . Then, an ANN, like most machine learning algorithms, can be simply considered as a non-linear function  $f_\theta : \mathcal{X} \mapsto \mathcal{Y}$  with some set of parameters  $\theta$ , which maps an input sample  $\mathbf{o}_i$  to an output  $y_i$ . Although not limited to, among other tasks ANNs are generally used to solve two types of problems in machine learning: regression and classification. When  $y_i$  takes continuous values, such as commodity prices, the size of a tumor, spectral amplitudes, etc., the task is a *regression* problem. On the contrary, when  $y_i$  is a class within a discrete set, such as speaker identities, phonemes, tumor presence or absence, etc., it is a *classification* problem.

ANNs learn the mapping from  $\mathcal{X}$  to  $\mathcal{Y}$  up to a reasonable bound of accuracy by adjusting their parameters  $\theta$  based on the total amount of error between the ground truth  $y_i$  and the estimate  $\hat{y}_i$ . The function that measures this error is called by different names: the *cost*, the *loss* or the *objective* function, denoted by  $\mathcal{L}(\theta)$ . Indeed, the value of  $\theta$  determines the current value of  $\mathcal{L}$ , and the goal of training algorithms is to minimize the value of  $\mathcal{L}$  until convergence. The loss function used for real-world problems is usually a non-convex function of  $\theta$  [333]. If  $\mathcal{L}$  is differentiable, then we can make a step in the steepest direction by simply finding the gradient of  $\mathcal{L}$  with respect to each element in  $\theta$  over the whole data set, subtract the gradient  $\nabla \mathcal{L}$  from corresponding element in  $\theta$  iteratively to nudge it in the direction where  $\mathcal{L}$  is smaller. The gradients are generally scaled by a small value  $\eta$ , called the *learning rate*, which decides the step size to avoid missing the minimum when it is too close. The process of analytically computing the partial derivative of the error with respect to each parameter, and updating those parameters to minimize the loss function is the workhorse of machine learning and is referred to as *backpropagation using gradient descent* [240].

A commonly used loss function is the mean squared error (MSE) which measures the average squared distance between the desired and the actual output of the neural network:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \|y_i - \hat{y}_i\|^2. \quad (2.9)$$

A step in gradient descent to update each parameter  $\theta_j$  is given as follows:

$$\theta_j \leftarrow \theta_j - \eta \frac{\partial \mathcal{L}}{\partial \theta_j}. \quad (2.10)$$

In practice, the huge size of the training data set makes it intractable to compute the gradient over the whole data set at once as shown in Equation (2.9). Alternatively, one may compute gradients over individual training examples and update the parameters each time, this is called *stochastic gradient descent* (SGD). Each step of SGD can be computed much faster than regular gradient descent, but the variance of updates is much higher, leading to fluctuations in the loss function. As a trade-off between the two methods, the gradients can be computed over disjoint subsets of training data called *mini-batches*. Computing the gradient over a mini-batch is computationally efficient and leads to more stable convergence. This modified version of the algorithm is called the *mini-batch gradient descent*, but it is interchangeably referred to as SGD in the literature so we will call it SGD from here on. When SGD has seen the whole training set, i.e., it has computed and backpropagated the gradients over all the mini-batches, this is referred to as the completion of an *epoch*. For better results, systems are trained for several epochs until the loss does not change significantly any further.

**Activation function** There are several choices for the activation function, such as the sigmoid function, hyperbolic tangent ( $\tanh$ ), rectified linear unit (ReLU), softmax, etc. The sigmoid function was traditionally used because it squashes inputs to the  $[0, 1]$  range, but its derivative is upper bounded by 0.25 which implies that the magnitude of the gradient values reduces by at least 75% at each layer. This leads to the *vanishing gradient* problem [278] in deep networks, which also arises with the  $\tanh$  activation function. Hence, in recent years ReLU [324] has become the preferred choice of activation:

$$g(x) = \max(x, 0). \quad (2.11)$$

The derivative of the ReLU function is as follows:

$$g'(x) = \begin{cases} 0 & \text{for } x < 0, \\ 1 & \text{for } x \geq 0. \end{cases} \quad (2.12)$$

The softmax function is another popular activation function that ensures that the outputs are positive and that they sum up to 1. It is commonly used as the activation for the output layer in classification problems. The outputs can then be interpreted as a probability distribution over the categorical classes.

**Relevant deep neural network models** Deep learning has enabled researchers to explore several complex network architectures which may be suitable for specific tasks. Here we briefly discuss some of the models that are relevant for processing speech data. As described in Section 2.2.1, speech data is processed by decomposing an utterance into fixed-length overlapping segments called frames. In applications such as speech recognition or speaker identification, MFCC or logmel features are computed for each frame, and the whole sequence is fed as input to a neural network that accounts for the temporal dynamics of speech.

The simplest of all deep neural network (DNN) architectures is the *feed-forward* network as shown in Figure 2.3. Simple feed-forward architectures are great function approximators for data that can be represented by independent factors, such as predicting loan application outcomes using a person's financial attributes, but they fail to efficiently capture the local spatial and temporal relationships that exists in image

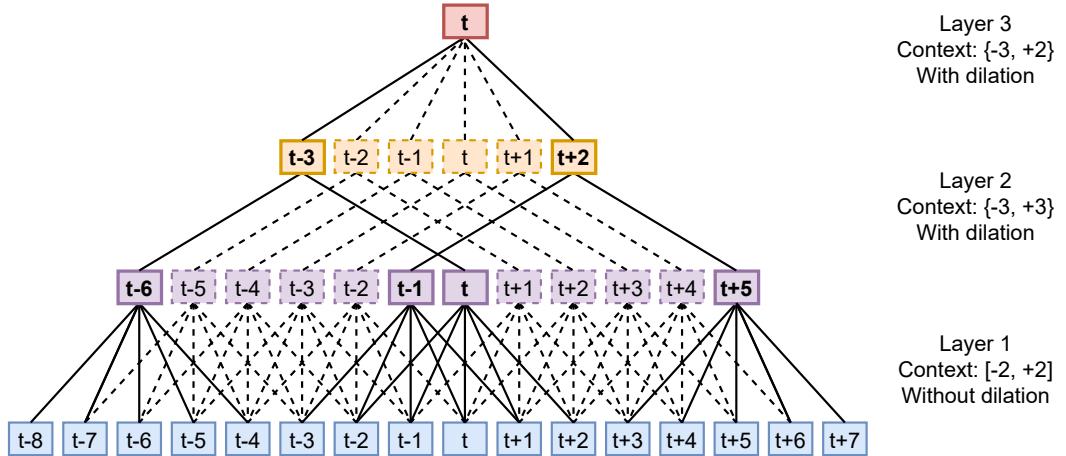


Fig. 2.4 Time delay neural network architecture with dilation in layers 2 and 3. Dotted lines indicates the connections and the nodes which are not included in the computation due to dilation applied to the layers.

1 and speech data. Hence other architectures, such as convolutional neural networks (CNNs) [164], have  
2 been proposed to model these relationships effectively.

3 CNNs are a special type of feed-forward network, which is best suited for image data since it is inspired by  
4 the working of the visual cortex. The fundamental neuron in a CNN acts like a *kernel* having a pre-specified  
5 2D receptive field, that moves over the input image and performs a convolution operation with the area  
6 covered by it. Several kernels in a layer convolve with the same input image to learn different spatial properties  
7 of the image. They are followed by a non-linear activation and a pooling operation to transform the kernel  
8 output into a *feature map*. At each successive layer, the kernels learn to discriminate between hierarchical  
9 features such as edges, geometrical shapes, objects and eventually lead to a fully-connected layer that predicts  
10 the output classes.

11 Another feed-forward architecture was proposed to model the temporal dependencies present in speech  
12 data, called the time delay neural network (TDNN) [303]. It can also be seen as a CNN with 1D kernels,  
13 where each layer operates at a different temporal resolution as shown in Figure 2.4. The bottom layers of a  
14 TDNN learn an affine transform for a narrow context window at each time step, and the context becomes  
15 wider in upper layers. Due to shared kernel weights across time steps, TDNNs are capable of learning  
16 translation invariant feature transforms. It has been observed that TDNNs can be made computationally  
17 efficient by sub-sampling the activations that are passed on to the next layer due to a large overlap between  
18 neighbouring contexts [219]. This process of sampling non-contiguous frames for building the context is  
19 called *dilation*.

20 Factorized TDNN (TDNN-F) models as depicted in Figure 2.5 have been proposed by Povey et al. [222]  
21 to reduce the number of parameters and the computational cost. Each unit of a TDNN hidden layer acts  
22 as a 1D kernel which produces a feature map by processing several time frames together depending on the  
23 temporal resolution of the layer. Let  $\mathbf{W}$  be the weight matrix between the hidden layer and the feature map,  
24 then TDNN-F factorizes  $\mathbf{W} = \mathbf{P}\mathbf{Q}$  into two factors using the singular value decomposition, and imposes a  
25 constraint such that  $\mathbf{P}$  is semi-orthogonal, i.e.,  $\mathbf{P}\mathbf{P}^T = \mathbf{I}$  or  $\mathbf{P}^T\mathbf{P} = \mathbf{I}$ . The interior dimension between  $\mathbf{P}$   
26 and  $\mathbf{Q}$  is much smaller than the number of units in the hidden layer or the feature map and it is referred to  
27 as the linear bottleneck dimension. It is assumed that even with a reduced number of parameters, no model  
28 strength is lost if one of the factors is constrained to be semi-orthogonal. This constraint is imposed every  
29 four training iterations by updating matrix  $\mathbf{P}$  such that it is closer to being semi-orthogonal by using the

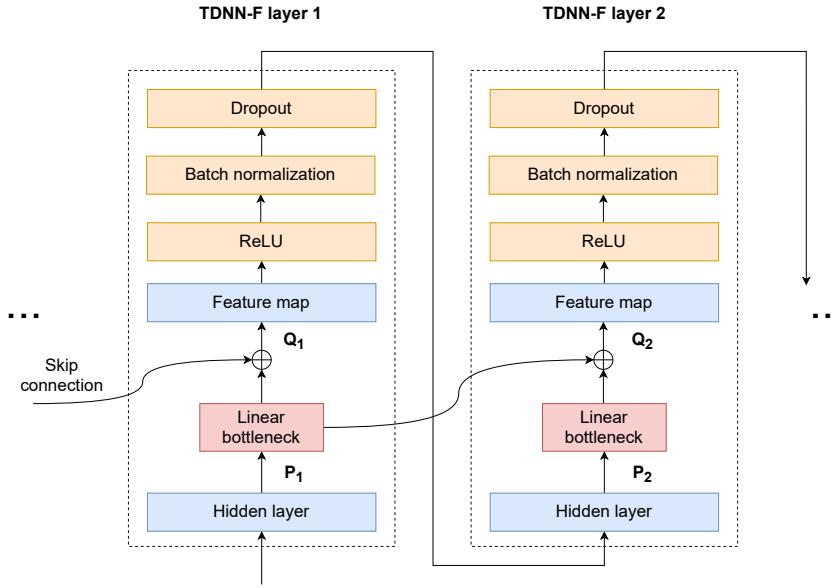


Fig. 2.5 Factorized TDNN (TDNN-F) architecture showing the linear bottleneck inserted between the hidden layer and the feature map. Matrices  $\mathbf{P}_1$  and  $\mathbf{P}_2$  are constrained to be semi-orthogonal and  $\oplus$  represents the concatenation of the linear bottleneck and the feature map layer linked by the skip connection.

following rule:

$$\mathbf{P} \leftarrow \mathbf{P} - \frac{1}{2\eta^2} (\mathbf{P}\mathbf{P}^\top - \eta^2 \mathbf{I})\mathbf{P}. \quad (2.13)$$

Here  $\eta$  modifies the constraint such that a scaled version of  $\mathbf{P}$  is expected, and it is similar to the learning rate hyperparameter of a neural network because it controls how fast the layer parameters are changing in a consistent manner. It is claimed that a network composed of TDNN-F layers does not require pretraining, but the training may be unstable if  $\mathbf{P}$  is too far from being semi-orthogonal. Hence, it is initialized using the Glorot mechanism [101] and the learning rate is carefully chosen as  $\eta = \sqrt{\text{tr}(\mathbf{K}\mathbf{K}^\top)/\text{tr}(\mathbf{K})}$ , where  $\mathbf{K} \equiv \mathbf{P}\mathbf{P}^\top$  and  $\text{tr}(\cdot)$  computes the trace of a matrix. Another feature that helps stabilize the training of TDNN-F networks is *skip connections* which append<sup>3</sup> the linear bottleneck of previous layers to the feature map of the current layer as shown in Figure 2.5. Such a network is much faster to train using parallel computing (GPUs) than other neural networks which model temporal dependencies, such as recurrent neural networks (RNN), due to their feed-forward architecture.

TDNNs model temporal dependencies by merging contextual information with the input at the current time step and estimating the output through a feed-forward mechanism. In contrast, RNNs do not follow a feed-forward mechanism, but incorporate some kind of memory or *hidden state* of a sequence that remembers information in previous time steps and is used for subsequent computations. The hidden state for the current time step is obtained by combining the hidden state in the previous time step and the current input. The parameters of current and previous hidden states are optimized based on the feedback from the current output using a modified version of backpropagation, called “backpropagation through time” [312]. One limitation of vanilla RNNs is that they can be very inefficient at learning relevant information in the sequence due to gradients vanishing across time, hence several RNN variants have been proposed to retain the feedback

<sup>3</sup>Generally, a skip connection *adds* the input of the current layer (or a previous layer) to the output of the current layer, but [222] refers to concatenation as the skip connection.

1 signal, based on long short-term memory (LSTM) [121], bidirectional LSTM [106] (BLSTM) or gated  
 2 recurrent unit (GRU) [48] layers. Another limitation is the sequential nature of computation at each time  
 3 step which cannot leverage the enormous parallelism offered by advanced computing infrastructure like  
 4 GPUs. Due to these limitations, RNN architectures are increasingly being replaced by convolutional or  
 5 Transformer [297] based architectures for sequential data like speech and natural language. We will describe  
 6 some of these architectures in later sections when we apply them to our applications.

7 Although applying DNNs to speech data has undeniably benefitted the community by improving the  
 8 performance of the systems, it has significantly raised the demand for large-scale data speech collection.  
 9 Previous studies have shown that without requiring new data, neural networks can avoid overfitting if the  
 10 speech signals in the training set are artificially augmented with reverberation or noise [199, 153]. This way,  
 11 the existing data set can be multiplied several folds with diverse settings contributing to the enrichment  
 12 and robustness of the model. The experiments performed in this thesis rely on these techniques to achieve  
 13 state-of-the-art results.

### 14 2.2.3 Automatic speech recognition

15 Automatic speech recognition (ASR) aims to convert an utterance into its textual content, also called  
 16 transcription. The output of ASR is used by natural language understanding systems that take speech as  
 17 input, and is widely deployed in commercial applications ranging from cloud servers to mobile devices. We  
 18 mentioned before that speech utterances are of varying duration and the system producing them also varies  
 19 through time, hence they are processed as a sequence of  $T$  overlapping time frames of fixed duration. The  
 20 input to ASR is a sequence represented as a matrix,  $\mathbf{O} = [\mathbf{o}_1, \dots, \mathbf{o}_T]^\top \in \mathbb{R}^{T \times A}$  of length  $T$  time frames,  
 21 where  $\mathbf{o}_t \in \mathbb{R}^A$  are feature vectors derived from the speech signal, e.g., MFCCs or logmel spectra, and the  
 22 output is the estimated word sequence  $\hat{W}$ . This problem can be formulated as [320]:

$$23 \quad \hat{W} = \underset{W}{\operatorname{argmax}} P(W|\mathbf{O}). \quad (2.14)$$

24 Researchers in the domain of ASR have tried to solve the problem defined in Equation (2.14) through  
 25 two main approaches. The first one is called the *conventional HMM-based* pipeline and the more recent one  
 26 is the *end-to-end ASR* approach. We use both of them in different parts of this thesis, therefore we give a  
 27 brief overview of both of them below.

28 **Conventional approach** In this approach, it was noted that it is infeasible to directly model the conditional distribution of the most probable word sequence given the acoustic features. To simplify this problem,  
 29  $P(W|\mathbf{O})$  can be decomposed into a simpler probabilistic model by defining a generative process, and then  
 30 the true word sequence is inferred from it. The generative model is depicted in Figure 2.6 and defined as  
 31 follows: we know that an utterance is a sequence of spoken words that are distributed according to the  
 32 language model. Spoken words are in turn made up of a sequence of fundamental sounds called phonemes  
 33 ( $\rho$ ), but the same phoneme can manifest itself differently in the signal due to natural variation of the vocal  
 34 tract and the context surrounding it, also known as the coarticulation effect. The different manifestations of  
 35 a phoneme in the presence of varying contexts can be represented by triphones (i.e., tied context-dependent  
 36 phonemes) where each triphone is modeled by its own hidden Markov model (HMM) and the speech  
 37 features follow a Gaussian probability density function within each state. This is called the GMM-HMM  
 38 approach. Alternatively, deep neural networks, instead of GMMs, can be used to model the density of  
 39 HMM states [120], which is referred to as the hybrid DNN-HMM approach that is used in this thesis and  
 40 described further. To handle data scarcity issues, the triphone HMM states are clustered using a decision

tree and the same emission probability is shared by all the states in a given cluster  $S$ , which is also called a “tied state”. Hence, the ASR problem is reformulated as [189]:

$$\hat{W} = \operatorname{argmax}_W P(\mathbf{O}|W)P(W) \quad (2.15)$$

$$\approx \operatorname{argmax}_W \sum_{S,\rho} P(\mathbf{O}|S)P(S|\rho)P(\rho|W)P(W). \quad (2.16)$$

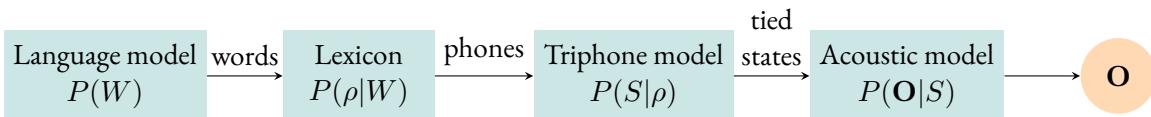


Fig. 2.6 Generative model for ASR.

Here  $P(W)$  is the so-called language model that represents the prior distribution of word sequences,  $P(\rho|W)$  is the lexicon which maps all the words in the vocabulary to their corresponding phoneme sequences,  $P(S|\rho)$  maps a phoneme sequence to the corresponding tied state sequence  $S = [S_1, \dots, S_T]$ , and  $P(\mathbf{O}|S) \propto \prod_{t=1}^T P(S_t|\mathbf{O})/P(S_t)$  where the tied state posterior probabilities  $P(S_t|\mathbf{O})$  are given by the so-called DNN acoustic model and  $P(S_t)$  is the prior probability of each tied state. These models are learned independently and composed together as a graph using finite state transducers (FST). In some use-cases, a sequence of phonetic features called bottleneck (BN) features, denoted as  $\mathbf{B}$ , can be extracted from an intermediate layer of the ASR acoustic model [321] and used, possibly in combination with other features, for other tasks. The above mentioned generative model is also used to “synthesize” speech utterances as explained in the next section.

The DNN acoustic model  $P(S_t|\mathbf{O})$  is trained on acoustic features  $\{\mathbf{O}_i\}_{i=1}^N$  extracted from the utterances  $\{\mathbf{s}_i\}_{i=1}^N$  in some annotated data set  $\mathcal{D}$  and the corresponding transcriptions  $\{W_i\}_{i=1}^N$ . The cost function  $\mathcal{L}_{\text{ASR}}$  which can be carefully crafted to predict the accurate triphone sequence, is minimized to optimize the parameters of the acoustic model. For example, one popular [125, 111, 99] cost function is the following:  $\mathcal{L}_{\text{ASR}} = \mathcal{L}_{\text{MMI}} + 0.1 \cdot \mathcal{L}_{\text{CE}}$ , which is composed of two terms. The dominant term,  $\mathcal{L}_{\text{MMI}}$ , is the lattice-free maximum mutual information (LF-MMI) [224] cost which aims to maximize the posterior probability of the ground truth word sequence  $W_i$ :

$$\mathcal{L}_{\text{MMI}} = - \sum_{i=1}^N \log \frac{P(\mathbf{O}_i|W_i)P(W_i)}{\sum_{W'} P(\mathbf{O}_i|W')P(W')}. \quad (2.17)$$

The numerator is the joint likelihood of the acoustic features  $\mathbf{O}_i$  and the ground truth word sequence  $W_i$ , while the denominator is the likelihood of the acoustic features marginalized over all possible word sequences. The numerator is computed by summing over all tied state sequences corresponding to  $W_i$ :  $P(\mathbf{O}_i|W_i) = \sum_{S_i,N_i} P(\mathbf{O}_i|S_i)P(S_i|N_i)P(N_i|W_i)$  where  $P(\mathbf{O}_i|S_i) \propto \prod_{t=1}^{T_i} P(S_{i,t}|\mathbf{O}_i)/P(S_{i,t})$ , and  $P(S_{i,t})$ ,  $P(S_i|N_i)$ ,  $P(N_i|W_i)$  and  $P(W_i)$  are fixed. The numerator is computed in a similar way, except that the (intractable) sum over all possible word sequences with a word-level language model is approximated by a (tractable) sum over all possible phoneme sequences with a phoneme-level language model. The second term  $\mathcal{L}_{\text{CE}}$  of the cost function is the frame-level cross-entropy loss between the true and estimated tied states, which acts as a regularizer [224]:

$$\mathcal{L}_{\text{CE}} = - \sum_{i=1}^N \sum_{t=1}^T \log P(S_{i,t} | \mathbf{O}_i). \quad (2.18)$$

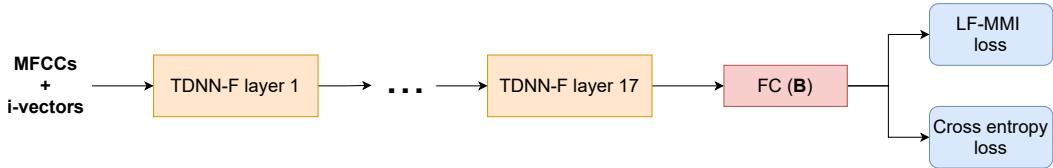


Fig. 2.7 Network architecture for ASR acoustic modeling composed of TDNN-F layers followed by the fully-connected bottleneck layer **B** which branches into the computation of the two loss functions, LF-MMI and cross-entropy. The skip connections between TDNN-F layers are not shown for the sake of simplicity.

The acoustic model, used to design and evaluate the anonymization techniques in and after Chapter 5, is a deep neural network which takes MFCCs appended with i-vectors for speaker adaptation [246]. It is composed of 17 TDNN-F layers<sup>4</sup> having 1536 neurons in the hidden layer and 160 neurons in the linear bottleneck, followed by a 256-dimensional fully connected layer **B** which leads to two branches that compute the LF-MMI loss and the cross-entropy loss over the tied states as shown in Figure 2.7. This network needs alignment between the observations and the HMM state sequence before starting the training, which is obtained using an HMM that is trained using the iterative Baum-Welch algorithm [229] to marginalize over all possible state sequences that could have generated the observations. At test time, the Viterbi algorithm [93] is used to find the most likely sequence of HMM states ( $S$ ) that emit the sequence of acoustic observations  $\mathbf{O}$ , and thereby also give the likelihood of observing  $\mathbf{O}$  given the state sequence  $S$ . The sequence of acoustic observations is passed through the FST which is the composition of the acoustic model, the context dependency, the lexicon, and the language model. There may be several paths that lead to alternative transcriptions for the same input. Picking the most likely node at each time step (i.e., greedy search) may not lead to the best transcription, hence multiple best paths are considered together at each time step and the remaining less likely paths are pruned. This ensures that the most likely path emerges as the winner and is referred to as the beam search algorithm. The number of paths stored at each time step is called the beam width. It is widely used in ASR and machine translation where the most likely output sequence could not be found using the greedy approach (i.e., beam width is equal to 1). A larger beam width ensures better results but requires the storage of more alternate transcriptions thereby increases the computational cost.

The conventional HMM-DNN approach is quite effective and widely used for ASR, but not without its limitations [105]. The major criticism for this approach is the complexity of its pipeline and the requirement for human expertise. A pretrained GMM-HMM model is needed to generate triphone states that are used as the training targets for the cross-entropy branch of the DNN acoustic model. Separately prepared language model, lexicon and acoustic model are glued together which might compound the overall errors of the ASR system. Moreover, a lexicon must be prepared by expert linguistic rules which might not be available for low-resource languages. The end-to-end approach provides a solution to these issues by subsuming the different models into a single neural network. It has been shown that they perform reasonably well as compared to the conventional pipeline [275], and that they are well suited for low-resource settings [314, 253].

**End-to-end approach** As a holistic solution to these limitations, recent years have seen rapid development in the domain of end-to-end speech recognition which aims to directly transcribe graphemes (i.e., lexical

<sup>4</sup>TDNN-F layers are described in Figure 2.5.

characters) from speech instead of phonemes, thereby collapsing all the components of the conventional pipeline into a single neural network which is trained in an end-to-end fashion. Ideally, the end-to-end ASR network optimizes its parameters directly based on the sequence-level transcription accuracy, which is the true measure of ASR performance. In practice, a language model is used to re-score the outputs produced by the ASR network which helps them to achieve competitive performance compared to the conventional pipeline [105]. It is reasonable to use a language model because they are trained on additional text-only data that provides realistic prior distribution over words and corrects the mistakes made by the end-to-end network.

Replacing the composite pipeline of conventional ASR with a single neural network requires innovation, both in terms of the training objective as well as the architecture. Graves et al. [105] proposed to use connectionist temporal classification (CTC) as the training objective for an end-to-end ASR network containing five layers of BLSTM with 500 cells each. It does not require a pre-defined alignment between the acoustic features  $\{\mathbf{O}_i\}_{i=1}^N$  and the corresponding true grapheme label sequence  $\{Y_i\}_{i=1}^N$ , where  $Y_i = [y_1, \dots, y_c, \dots, y_C]$ ,  $y_c \in \mathcal{G}$  with  $\mathcal{G}$  being the set all grapheme symbols including characters, punctuations and space, and  $C$  is the number of characters in the transcription of the  $i$ -th utterance. Instead, it uses the conditional distribution  $P(\mathbf{a}|\mathbf{O})$  which gives the probability of a possible alignment sequence  $\mathbf{a} = [\bar{y}_1, \dots, \bar{y}_T]$  given the acoustic observation sequence  $\mathbf{O}$  of length  $T$  frames, where  $\bar{y}_t \in \mathcal{G}$ . The characters  $\bar{y}_t$  in a given alignment sequence are obtained by repeating  $y_c$  to match the length of  $\mathbf{O}$ . The probability of an alignment can be obtained using the chain rule and it is simplified by a conditional independence assumption:

$$P(\mathbf{a}|\mathbf{O}) = \prod_{t=1}^T P(\bar{y}_t|\bar{y}_1, \dots, \bar{y}_{t-1}, \mathbf{O}) \approx \prod_{t=1}^T P(\bar{y}_t|\mathbf{O}). \quad (2.19)$$

We must marginalize over all possible alignment sequences  $\mathbf{a}$  to get the probability of the grapheme sequence  $G_i$  but in practice, it is not feasible to sample all possible alignments, so Monte-Carlo sampling is used to compute the CTC loss and its gradient. The final loss function  $\mathcal{L}_{\text{ctc}}$  is as follows:

$$\mathcal{L}_{\text{ctc}} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M \log P(\mathbf{a}_j^{(i)} | \mathbf{O}_i). \quad (2.20)$$

Here  $\mathbf{a}_j^{(i)}$  is one of the  $M$  possible alignment sequences for  $i$ -th utterance, sampled using Monte-Carlo.

Attention-based approaches are an alternative to CTC which do not make any conditional independence assumptions. Instead, they consider all the previous outputs and the whole input sequence to estimate the posterior [21, 311]. The attention mechanism does not require an intermediate alignment representation, hence the loss is computed as follows:

$$\mathcal{L}_{\text{att}} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C \log P(y_c | y_1, \dots, y_{c-1}, \mathbf{O}_i). \quad (2.21)$$

Although several different architectures based on recurrent [114, 25, 198] and convolutional [332] neural networks were proposed for end-to-end ASR, Watanabe et al. [311] proposed to combine CTC and attention-based mechanisms through multi-objective training, which is used in Chapters 3 and 4 to design and evaluate speaker anonymization techniques. As shown in Figure 2.8, it follows the so-called encoder-decoder architecture where the encoder, composed of four BLSTM layers with 320 units each, transforms the input sequence into a new bottleneck representation  $\mathbf{B}$ , and the decoder, with a single unidirectional LSTM layer having 320 units, predicts the grapheme sequence from  $\mathbf{B}$ . The attention layer

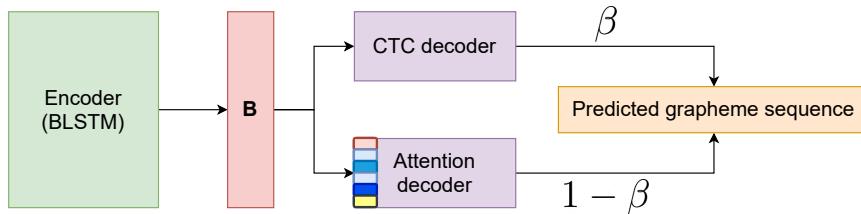


Fig. 2.8 End-to-end ASR architecture with multi-objective training consisting of CTC and attention-based loss functions combined by a hyperparameter  $\beta$ , where  $0 < \beta < 1$ . The attention layer is shown as the heatmap in the attention decoder which assigns combination weights to the bottleneck representation  $\mathbf{B}$  before processing it by the LSTM layer.

- <sup>1</sup> sits in between the encoder and the decoder and tells the decoder at each time step how much weight should
- <sup>2</sup> be assigned to a particular part of  $\mathbf{B}$  to predict the output grapheme at that time step with the least amount
- <sup>3</sup> of error. The weights of the attention layer are learned within the end-to-end framework [143].

#### 4 2.2.4 Speech synthesis

- <sup>5</sup> As opposed to ASR, the goal of speech synthesis also known as text-to-speech (TTS) is to convert a text
- <sup>6</sup> string into a speech waveform. We have already seen how there have been historical efforts to produce speech
- <sup>7</sup> sounds either using a physical model like von Kempelen’s “speaking machine” or electronic models like
- <sup>8</sup> Dudley’s VODER. Modern-day TTS technology has greatly advanced especially with the introduction
- <sup>9</sup> of statistical models, like neural networks. The current state-of-the-art TTS models can produce almost
- <sup>10</sup> natural-sounding speech from any text, in several voices, and multiple languages. The progress in modern
- <sup>11</sup> TTS technology can be divided into three generations of systems [210], namely unit selection, statistical
- <sup>12</sup> parametric speech synthesis (SPSS), and neural speech synthesis. The exact formulation of these three
- <sup>13</sup> systems is out of scope for this thesis, but they are all based on some essential fundamental principles which
- <sup>14</sup> will be briefly covered here. Finally, a small note about the evaluation of TTS systems is mentioned at the
- <sup>15</sup> end of this section.

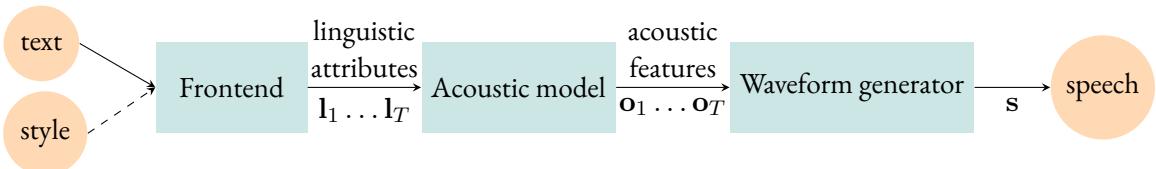


Fig. 2.9 General schema of a TTS system with the style features as optional input to the frontend.

- <sup>16</sup> Figure 2.9 shows a general schematic diagram for TTS systems. It is composed of the frontend, the
- <sup>17</sup> acoustic model, and the waveform generator. Note that unit selection methods directly generate a waveform
- <sup>18</sup> using the output of the frontend. In contrast, SPSS methods employ an acoustic model to first convert the
- <sup>19</sup> frontend’s output into spectral parameters of the target speech and then generate the target speech using a
- <sup>20</sup> waveform model. We briefly describe these three blocks of TTS systems below.

- <sup>21</sup> **The frontend** The first challenge is to predict non-linguistic attributes such as speaking style, emotions,
- <sup>22</sup> and prosodical cues simply from the plain text without any auxiliary information. Speaker traits, which are
- <sup>23</sup> generally derived from the fixed number of speakers in the training data set and are otherwise impossible to
- <sup>24</sup> predict simply from text, are also needed as input to the frontend. The frontend contains a large number of

meticulous rules hand-crafted by linguists for different languages which convert written form of words to their spoken form by normalizing them into tokens found in a dictionary, assigning them a part-of-speech class, getting the exact sequence of phonemes to pronounce, as well as predicting the intonation pattern and phrase breaks. It is costly to build and maintain the frontend for different languages since it requires careful analysis of phonetic inventory and different ways to pronounce certain complicated tokens like abbreviations, numbers, currency symbols, etc. Some end-to-end TTS systems [306, 15] aim to replace this complex pipeline with a neural network that encodes linguistic attributes in its hidden layers and requires only the sequence of letters as input, but they are still in their infancy and make some crucial pronunciation mistakes [252] which prevent their commercial use. There have also been efforts to capture general speaking style including speaking rate, emotions, prosodical patterns using neural network embeddings [307]. Given the input text and the style embedding, expressive speech can be synthesized. Moreover, intonation and pronunciation mistakes can be corrected by using morphological features [281] that are present in the text. Nevertheless, all commercial TTS systems still maintain the traditional frontend because it can be easily understood and corrected when it makes mistakes.

At this point, we know that it is hard to reproduce the exact linguistic and paralinguistic attributes that were present in the original speech simply from the text content. Moreover, the generated utterance may not be as diverse as the natural speech due to limited training data for TTS. Therefore, given these drawbacks, we discard the seemingly simple solution to achieve the privacy objective, i.e. to convert speech to text using ASR and then synthesize speech from this text using TTS, due to the destruction of usable information which is contrary to the goals of this thesis. Having addressed this possibility, we move on to describe the remaining components in the TTS pipeline that aim to produce a waveform using the linguistic and non-linguistic features generated by the frontend.

**Acoustic model** The TTS acoustic model converts linguistic features generated by the frontend into a sequence of acoustic features, such as a magnitude spectrogram, through a regression task that can be easily solved using a neural network. One popular approach is to use an autoregressive<sup>5</sup> neural network-based acoustic model, such as the one proposed by Lorenzo-Trueba et al. [179], to generate the Mel-spectrogram  $\mathbf{O}$  of an utterance given the linguistic features  $[\mathbf{l}_1, \dots, \mathbf{l}_T]$  of  $T$  frames extracted from the text. This approach generates the Mel-spectrogram corresponding to the target speaker.

The autoregressive acoustic model, as shown in Figure 2.10, is a sequence model with feed-forward and recurrent LSTM layers, where the output acoustic feature at time  $t$  is produced depending on the whole input sequence and some of the acoustic features at previous times. Hence, the probability of observing output acoustic features is defined as follows:

$$P(\mathbf{o}_1, \dots, \mathbf{o}_T | \mathbf{l}_1, \dots, \mathbf{l}_T) = \prod_{t=1}^T P(\mathbf{o}_t | \mathbf{o}_{t-T'}, \dots, \mathbf{o}_{t-1}, \mathbf{l}_1, \dots, \mathbf{l}_T). \quad (2.22)$$

The acoustic model generates  $\mathbf{o}_t$  based on the previous  $T'$  outputs and the whole sequence of input linguistic features.

**Waveform generation** The final component of the TTS pipeline is the waveform generator, which processes the acoustic features to produce an intelligible waveform. Unit selection speech synthesis bypasses the acoustic modeling and instead, generates the waveform using a concatenation approach. It assumes that a large data set of natural speech spoken by real humans already exists, and at synthesis time a plausible sequence of speech segments that satisfy the given linguistic criteria are selected and stitched together to

<sup>5</sup> Autoregressive model relies on its past outcomes to predict the current one as in Equation (2.22).

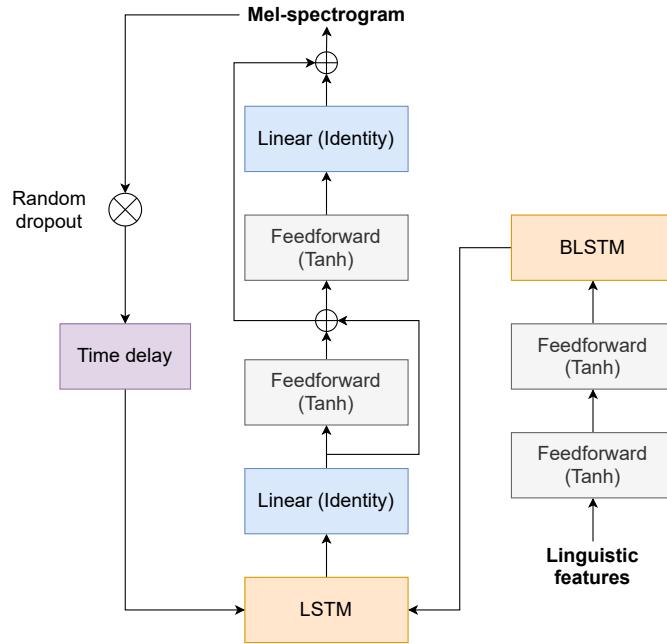


Fig. 2.10 Autoregressive network architecture for the TTS acoustic model [179]. The feedforward layers have Tanh activation and the linear layers have identity activation. It has skip connections to reinforce the noisy signal, and the time delay block passes the current output Mel-spectrogram back to the LSTM layer with a random dropout.

1 generate the waveform. The units that are stitched together are usually diphones, which are nothing but  
 2 the second half of one phone and the first half of another to retain the overlap between the two, thereby  
 3 capturing the coarticulation boundary. Therefore the data set must ideally include multiple instances of all  
 4 possible diphones in the considered language. This approach used to be popular due to the naturalness of  
 5 the generated speech. Yet there are several glaring limitations, for example, the concatenation of waveform  
 6 segments may not be smooth at the joins which may result in perceptual glitches while listening to the  
 7 produced audio. It is also not possible to have all the possible speaking styles in the data set and it is certainly  
 8 quite cumbersome to scale this approach to include new speakers.

9 SPSS methods try to alleviate the abovementioned limitations of unit selection by estimating the acoustic  
 10 parameters of the target speech, instead of using pre-recorded samples. The acoustic parameters can then  
 11 be manipulated or used directly to generate the waveform with desired properties. The task of waveform  
 12 generation just using the acoustic features is still challenging because they are composed of logmel features  
 13 only. These features, derived from the magnitude spectrogram, are not sufficient to reconstruct the original  
 14 signal since we also require the exact phase of each sine wave corresponding to that particular speech sound.  
 15 The waveform generator, also called the vocoder, needs to predict this missing information for producing an  
 16 intelligible speech signal. Traditional vocoders like STRAIGHT [145] or WORLD [201] provide various  
 17 analysis algorithms to be applied on the speech signal to efficiently extract the spectra, the fundamental  
 18 frequency, and the aperiodicity from all frames. They also provide synthesis algorithms that can combine  
 19 this information and produce a good quality speech waveform in real-time, but they ignore phase prediction  
 20 and instead make a minimum phase assumption, which leads to noticeable speech distortion in low F0  
 21 regions [201].

Autoregressive neural networks, such as Wavenet [211] and SampleRNN [194], have shown promising results as waveform generators, but they are highly inefficient and hard to parallelize due their sequential generation process. This limitation is relieved by Neural source-filter (NSF) models that are waveform generators [305] inspired by the classical source-filter paradigm of speech synthesis [117, 193]. The traditional source-filter model aims to mimic the speech production mechanism by assuming a *source* of sound that produces a discrete-time excitation signal  $\bar{e}$ , and a *filter* which modulates the frequency components of the excitation signal to convert it into a phoneme-like speech signal  $s$ . The source can produce either a periodic signal, such as an impulse train or a sine wave, to mimic the voiced excitation that contributes to the inherent harmonic structure in natural speech, or a noise signal for the unvoiced turbulence across the vocal tract. The filter acts like the vocal tract which produces formants in the spectral envelope due to its characteristic shape, and therefore must be individually designed for each type of sound. In its simplest form, it is nothing but the linear transformation of the excitation signal and the previous output, and the output at each time instance is specified by the following equation:

$$s[n] = \bar{e}[n] + \sum_{k=1}^p c_k \cdot s[n-k]. \quad (2.23)$$

Here,  $c_k$  are the filter coefficients that may vary for different speech sounds, and the output is produced at each time instance by considering  $p$  previous outputs, which is also known as the order of the filter.

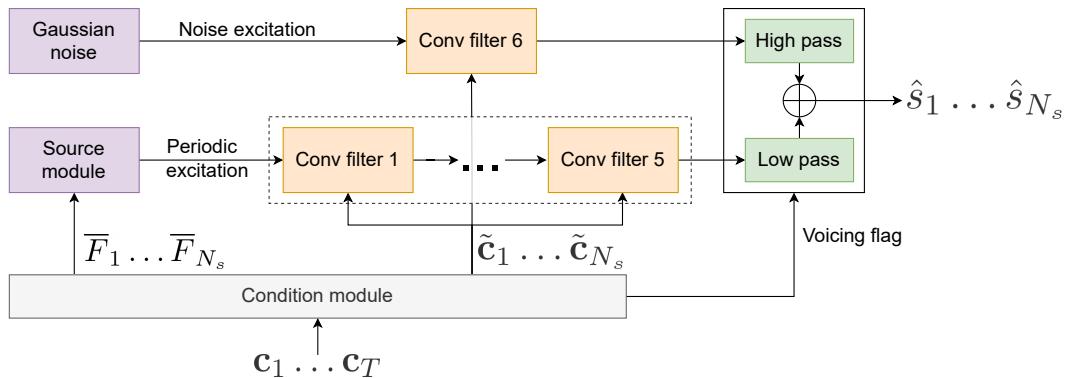


Fig. 2.11 NSF model architecture.

The NSF model<sup>6</sup> is much more advanced than the simple, linear source-filter formulation. First, it contains a condition module which takes  $[c_1, \dots, c_T]$  as input, where  $c_t = [F_t, \mathbf{o}_t]^\top$  is composed of the fundamental frequency  $F_t$  and the acoustic feature  $\mathbf{o}_t$  for the  $t$ -th time frame. It upsamples the fundamental frequency and outputs  $[\bar{F}_1, \dots, \bar{F}_{N_s}]$  to match the length of the target waveform. It also processes the acoustic features  $\mathbf{o}_t$  using BLSTM and convolutional layers, and concatenates the output with the upsampled fundamental frequency to get the condition feature sequence  $[\tilde{c}_1, \dots, \tilde{c}_{N_s}]$ .

Next, the source module which accepts  $[\bar{F}_1, \dots, \bar{F}_{N_s}]$  as input and generates a periodic signal for voiced sounds (i.e., a mixture of sine waves parametrized by their amplitude and phase) as an excitation based on the value of  $\bar{F}_n$ . Another module generates a separate noise excitation for unvoiced sounds. The periodic signal, combined with the condition features  $[\tilde{c}_1, \dots, \tilde{c}_{N_s}]$ , is transformed through a series of *Conv filters* which contain multiple dilated convolutional layers with residual connections, while the noise excitation is transformed using a single Conv filter. The periodic signal is subjected to a lowpass filter to preserve

<sup>6</sup>We describe here the state-of-the-art NSF model which is referred to as Harmonic-plus-Noise NSF model in [305].

dominant regions in higher frequencies, while the noise excitation is passed through a highpass filter to preserve lower frequencies. There are two configuration of the bandpass filters depending on the value of the voicing flag. They are configured such that the higher frequencies are preserved for voiced regions, while the lower frequencies for unvoiced regions. The output of the two bandpass filters is summed up to get the estimated target waveform  $[\hat{s}_1, \dots, \hat{s}_{N_s}]$ .

The Conv filters are learned to minimize the log spectral amplitude distance:

$$\mathcal{L}_{\text{NSF}} = \frac{1}{2TL} \sum_{t=1}^T \sum_{k=1}^L \left[ \log \frac{|\mathcal{F}_t[k]|^2}{|\hat{\mathcal{F}}_t[k]|^2} \right]^2. \quad (2.24)$$

Here,  $\mathcal{F}_t[k]$  and  $\hat{\mathcal{F}}_t[k]$  denote the  $k$ -th short-time Fourier transform coefficient of the  $t$ -th time frame obtained from the original and the predicted waveform, respectively.

**Evaluation** The evaluation of output speech is usually performed using mean opinion scores (MOS) obtained using subjective listening tests by human subjects [235]. Several such subjects with different gender and age profiles rate the speech on a perceptual scale based on its quality, intelligibility, and naturalness. This method of evaluation is costly, labour intensive, and slow, hence in this thesis we objectively evaluate the generated speech samples using ASR systems, which may not be a very effective measure of qualitative attributes of speech but highly correlate with human intelligibility [19, 92, 110].

A flexible high-quality TTS system that generates personalized utterances of a given target speaker can be used to clone the identity of any arbitrary person. But these threats are not investigated in this work and the solutions to such issues are beyond the scope of this thesis.

### 2.2.5 Automatic speaker recognition

Automatic speaker recognition is the task of recognizing the speaker of a given speech utterance. As mentioned in Section 1.1, speaker information in the speech signal is quite sensitive since it describes several attributes related to the speaker's identity and personality. Campbell, in his seminal tutorial on speaker identification [40], lists several factors responsible for the speaker-dependent characteristics present in speech signal due to speech production mechanism. Most of the factors arise due to the physiology and the shape of the vocal tract of the speaker. When the acoustic wave passes through the vocal tract, its frequency is modulated by the dominant resonances (i.e., formants), which can be easily observed in the spectral envelope of the signal. There are other speaker-dependent factors that emerge due to the source of excitation, generated by lungs and are then carried across the trachea over to the vocal folds. The source is responsible for phonetic features such as voicing, frication, whisper, etc. along with the fundamental frequency F0. The mass and length of the vocal folds are the defining properties for the fundamental frequency, hence it is a speaker- as well as gender-dependent characteristic. Some other characteristics that describe the style of speaking like the speaking rate, the dialectal shift in frequencies for speakers of similar language, and general prosodic patterns that emerge from conversations with specific vocabulary such as technical or professional settings, can also contribute to the speaker's identity.

Speaker recognition technologies are largely deployed in forensic studies [10, 283, 11] and telephone banking systems.<sup>7,8,9</sup> The techniques for automatic speaker recognition can be categorized into automatic

<sup>7</sup><https://www.us.hsbc.com/customer-service/voice/>

<sup>8</sup><https://www.chase.com/personal/voice-biometrics>

<sup>9</sup><https://www.lloydsbank.com/contact-us/voice-id.html>

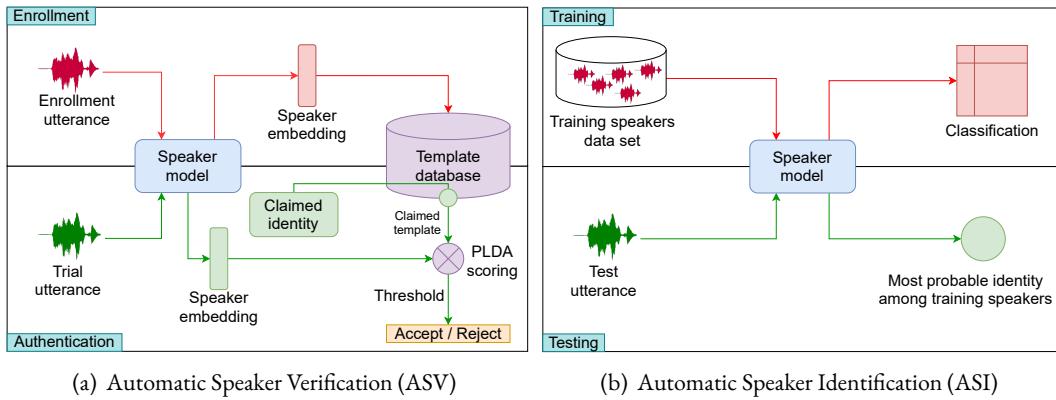


Fig. 2.12 Schema for two types of automatic speaker recognition. Red arrows indicate the enrollment or training flow, and green arrows indicate the authentication or testing flow. Note that the speaker model trained to classify speakers in the case of ASI can also be used to extract speaker embeddings for ASV, just like x-vectors.

speaker verification (ASV), i.e., authenticating the identity claimed by the speaker, and automatic speaker identification (ASI), i.e., determining the identity within a set of known speakers [31]. The schematic diagram of these two methods is depicted in Figure 2.12. ASV (Figure 2.12(a)) comprises two successive phases: enrollment and authentication. In the former, speakers are enrolled using discriminative speaker embeddings which are extracted from enrollment utterances. Speaker embeddings must have some characteristic features [150]. They must have large between-speaker variability and small within-speaker variability to get similar embeddings for different utterances from the same speaker. They must be easy to extract and difficult to impersonate, as well as robust against noise and distortion. The most popular embeddings called x-vectors [262] are obtained from an intermediate layer of a neural network trained to perform speaker classification. In the latter phase, the x-vector extracted from the utterance of an unknown speaker (called trial utterance) is compared with the x-vector of the speaker whose identity is being claimed, and a log-likelihood ratio score is computed by probabilistic linear discriminant analysis (PLDA) [147]. The ASV system then decides whether the trial utterance is from that speaker or not by comparing the obtained score with a threshold. As opposed to the open-set (rejection/acceptance) ASV task, ASI (Figure 2.12(b)) is a closed-set task in which a speaker classifier (e.g., similar to the one used to obtain x-vectors) is trained on training utterances from multiple speakers to later classify the identity of each test utterance as one of the known training identities.

X-vectors [262], which formulate ASI as a sequence classification task, have made the training and evaluation of an efficient ASI model quite straightforward. The architecture of the neural network used for extracting x-vectors is presented in Table 2.1, where the input is a sequence of speech features extracted from the utterance of a speaker, such as MFCCs, and the output is the posterior over the speaker classes. This task is accomplished using five TDNN layers followed by a statistical pooling layer and a fully connected classifier. Since the speaker information is present throughout the utterance, the statistical pooling layer computes the mean and standard deviation of the feature sequence produced by the preceding TDNN layers to retain the global speaker-related characteristics and diminish the local linguistic variations in speech. The output of intermediate layer just after statistical pooling (i.e., segment6), being rich in speaker information, is used as the x-vector. The design and evaluation of the ASV system are not so trivial as it requires the speaker embedding to capture relevant speaker information which generalizes to unseen speakers. The scoring mechanism must also produce values that truly discriminate closer speakers from farther ones. And finally,

Table 2.1 Original network architecture for x-vector speaker classification model as presented in [262, Table 1].

Layer name	Layer type	Layer context	Dilation	Total context	input × output
frame1	TDNN	$[t - 2, t + 2]$	No	5	$120 \times 512$
frame2	TDNN	$\{t - 2, t, t + 2\}$	Yes	9	$1536 \times 512$
frame3	TDNN	$\{t - 3, t, t + 3\}$	Yes	15	$1536 \times 512$
frame4	TDNN	$\{t\}$	No	15	$512 \times 512$
frame5	TDNN	$\{t\}$	No	15	$512 \times 1500$
stats pooling	Linear	$[0, T]$	NA	T	$1500T \times 3000$
segment6	Linear	$\{0\}$	NA	T	$3000 \times 512$
segment7	Linear	$\{0\}$	NA	T	$512 \times 512$
softmax	Output	$\{0\}$	NA	T	$512 \times N$

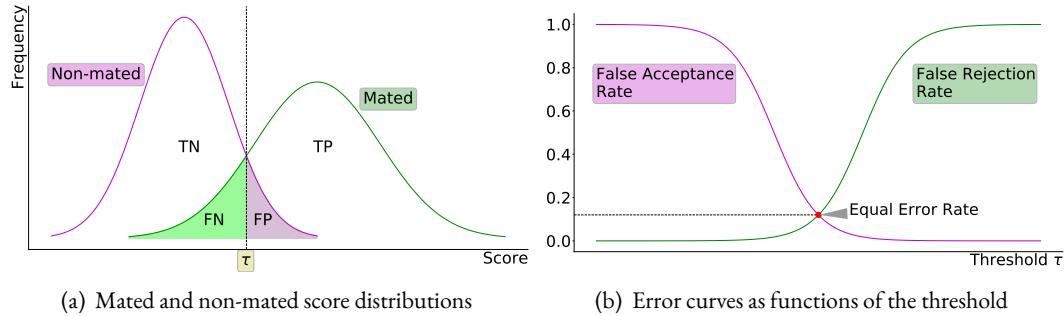


Fig. 2.13 Score distribution and threshold.

- 1 the selection of the decision threshold is critical for the performance of the ASV system. It must be carefully
- 2 calibrated based on the application, for example speaker authentication for bank transactions requires a
- 3 stringent threshold to avoid frauds while mobile phone unlocking requires a liberal threshold to allow quick
- 4 access to users.

The ASV model is conventionally evaluated using performance metrics, such as accuracy, which measure how frequently the model outputs the correct class. On the contrary, ASV systems are evaluated based on the errors they make. Suppose there are  $N$  trials, i.e., pairs of enrollment and trial utterances in a data set and each trial is either *genuine* (mated) or *impostor* (non-mated). Genuine, here referred as positive, trial represents the case when the speaker is really who he/she claims to be, hence the enrollment and the trial utterances belong to the same speaker. Contrary to that, in case of impostor or negative trial, the enrollment and the trial utterances belong to different speakers. If the model outputs ‘positive’ for TP examples which are truly positive, ‘positive’ for FP examples which are negative, ‘negative’ for FN examples which are positive, and ‘negative’ for TN examples which are truly negative, then the accuracy is given by:  $\frac{TP+TN}{TP+TN+FP+FN}$ , where  $TP + TN + FP + FN = N$ . It can be computed by simply counting the instances corresponding to TP, FP, TN, and FN, but the efficacy of ASV authentication is dependent upon two types of errors which are false acceptance rate (FAR), also known as type-I error and is given by  $\frac{FP}{FP+TN}$ , and false rejection rate (FRR), also known as type-II error and is given by  $\frac{FN}{FN+TP}$ . FAR and FRR are in turn dependent on the selected decision threshold,  $\tau$ .

Figure 2.13(a) shows the mated and non-mated score distributions. All biometric authentication systems must choose a threshold  $\tau$  by observing these two distributions such that the values of FN and FP are

minimized. It would be ideal if there were no overlap between mated and non-mated score distributions, but in practice, there is always some overlap due to the diversity of enrolled persons, imperfect scoring mechanisms, oversimplified modeling assumptions and limited training data. Hence, there is always a tradeoff between FAR and FRR as observed in Figure 2.13(b). Authentication systems may choose the threshold based on the sensitivity of their applications, but to evaluate them a fixed point is often chosen where the FAR is equal to the FRR. The error at this point is called the Equal Error Rate (EER). A lower EER indicates a better ASV system.

**Implications for the assessment of privacy** ASV and ASI systems are extensively used to design and simulate privacy attackers for the evaluation of the techniques proposed in this thesis as described in the next chapter. It is imperative to understand how to interpret the obtained results when the level of privacy protection is measured under strong hostile conditions. This thesis endeavors to show that it is possible to fine-tune the degree of hostility or the strength of the malicious entity who is trying to re-identify protected speakers. Hence, the anonymization techniques are evaluated based on their resilience against such strong criteria. Broadly speaking, anonymization does not imply that the features conveying speaker information are *completely* deleted from the speech signal or that it is even possible to do so. Instead, it means that the confusion for the attacker has been significantly increased by transforming the features of a given speaker such that they are indistinguishable from the other speakers. In that case, the degree of privacy is contingent upon the speaker discrimination capacity possessed by the attacker after anonymization and the set of enrollment speakers who are shortlisted as the possible targets. To that end, Section 3.1 gives a detailed guideline for designing the best possible attackers so that the degree of protection is truly measured.

The EER is widely used to evaluate biometric authentication systems, and a higher EER may indicate the inadequacy of the attacker who is trying to infer the true identity of the speaker. But some properties of EER may have limiting implications for it to be used as a general measure of privacy protection and evaluate the strength of a vast range of attackers. First, it assigns the same cost to false alarms (FP) and misses (FN) which may not be an optimal assumption to model the attacker. The attacker may choose a more relaxed setting to shortlist all possible speaker identities by lowering the threshold, which calls for a generalized evaluation where all possible priors over error costs are considered. Second, it only considers a single point of overlap between the mated and non-mated distributions, whereas it is possible that they are not monotonic and overlap at several points. The scores at the points of overlap may be present on either side of the EER and can be leveraged by the attacker to strengthen the attack. Finally, the EER indicates the overall efficacy of the biometric authentication system in the presence of several enrolled speakers while the actual level of protection for a particular speaker may slightly differ from the EER based on the indistinguishability of their individual score distribution from the overall scores.

In this thesis, some of the abovementioned limitations are alleviated by reporting other suitable metrics of privacy that are described and compared in the next chapter (Section 3.4). Later, some analysis is also provided to measure the best and worst-case privacy protection using re-identification metrics and formal methods such as differential privacy.

## 2.3 Techniques to transform speaker information

In this section, we review the fundamentals of the three most relevant techniques, i.e., adversarial learning, speech transformation, and voice conversion, that are widely used by researchers to modify speaker information in speech data. The remaining chapters of this thesis employ these techniques for their potential to hide/remove speaker-related biometric information from speech.

### 2.3.1 Adversarial Learning for Speech

Domain adversarial training helps to adapt neural network classifiers to a new domain without requiring labeled data in the new domain. The original paper [96] shows promising results on a handwritten digit classification task where the features learned using domain adversarial training are distributed identically, whether the image is grayscale or RGB (colored). It has been extensively applied to speech data since its inception. The idea of adversarial training enables neural networks to learn intermediate features in the form of hidden layers that are indiscriminate towards data belonging to similar classes but originating from different domains. For instance, the intermediate features for a particular phoneme class learned using domain adversarial training will be distributed almost identically, whether it is spoken in different ambient noise or by different speakers. Domain adversarial training is implemented as a neural network architecture with a so-called *adversarial branch* that predicts the domain, and a special layer just before this branch called the *gradient reversal layer* which scales the gradients that are being backpropagated through this layer using a negative scalar value. This layer ensures that the parameters of the preceding network are shifted in such a way that they implicitly remove the domain information from the representation, thereby making it invariant to irrelevant factors of variation that do not contribute towards the main task. This property of domain adversarial training is quite appealing to the speech community because speech signals are very expressive and full of factors of variation such as speaker's identity, channel information, emotions, language, accents, etc., but generally, machine learning models are trained to classify or predict only a single attribute from data, hence it is desirable to get rid of irrelevant attributes.

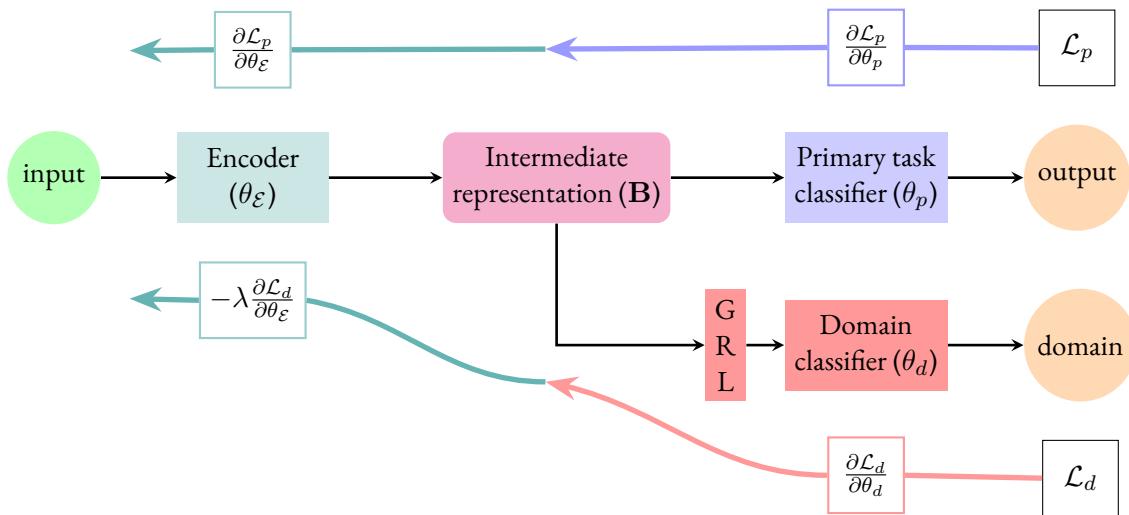


Fig. 2.14 General architecture for domain adversarial training of neural networks. Black arrows indicate forward propagation; purple, teal and red arrows indicate backpropagation of gradients for the primary task classifier, the encoder and the adversarial branch, respectively. The red GRL block refers to the gradient reversal layer with  $\lambda$  as the gradient reversal coefficient.

Speech researchers have employed adversarial training not only for domain adaptation [291, 62], but for enhancing speech quality [196, 172], training noise-robust ASR [254, 264], and learning representations which are invariant towards speaker [8, 195, 293], language [318, 6] and accents [272]. Figure 2.14 shows a general architecture for domain adversarial training which takes input speech features and transforms them into an intermediate representation (**B**) using an encoder neural network with parameters  $\theta_E$ . The intermediate representation is fed to the primary task classifier with parameters  $\theta_p$  which estimates the desired output, such as the transcription, emotional valence, etc., and computes the loss  $\mathcal{L}_p$ . In parallel, **B** is

## 2.3 Techniques to transform speaker information

35

also fed through the gradient reversal layer, which leaves it unchanged during the forward propagation and passes it to the adversarial branch, also called the domain classifier with parameters  $\theta_d$ , which predicts the domain label and computes the adversarial loss  $\mathcal{L}_d$ . The parameters  $\theta_{\mathcal{E}}$ ,  $\theta_p$ , and  $\theta_d$  are jointly estimated by solving the following minimax optimization problem:

$$\min_{\theta_{\mathcal{E}}, \theta_p} \max_{\theta_d} \mathcal{L}_o(\theta_{\mathcal{E}}, \theta_p, \theta_d). \quad (2.25)$$

Here,  $\mathcal{L}_o$  is the overall loss given by:  $\mathcal{L}_o(\theta_{\mathcal{E}}, \theta_p, \theta_d) = \mathcal{L}_p(\theta_{\mathcal{E}}, \theta_p) - \lambda \mathcal{L}_d(\theta_{\mathcal{E}}, \theta_d)$ , and  $\lambda$  is the gradient reversal coefficient which decides the trade-off between the primary and the adversarial objectives. A higher  $\lambda$  increases robustness of  $\mathbf{B}$  towards the domain but may decrease its efficacy towards the primary task.

During the backward pass, the parameters of the primary task classifier and the domain classifier are updated according to their respective losses  $\mathcal{L}_p$  and  $\mathcal{L}_d$ , but the encoder's parameters are updated with respect to both losses as opposing goals. The goal of domain adversarial training is to make  $\mathbf{B}$  invariant towards the domain label, hence the encoder parameters must be shifted in the direction such that the new  $\mathbf{B}$  maximizes  $\mathcal{L}_d$ , while minimizing  $\mathcal{L}_p$  at the same time. The gradient descent update rules for each part of the network are as follows:

$$\theta_p \leftarrow \theta_p - \eta \frac{\partial \mathcal{L}_p}{\partial \theta_p}, \quad (2.26)$$

$$\theta_d \leftarrow \theta_d - \eta \frac{\partial \mathcal{L}_d}{\partial \theta_d}, \quad (2.27)$$

$$\theta_{\mathcal{E}} \leftarrow \theta_{\mathcal{E}} - \eta \left( \frac{\partial \mathcal{L}_p}{\partial \theta_{\mathcal{E}}} - \lambda \frac{\partial \mathcal{L}_d}{\partial \theta_{\mathcal{E}}} \right), \quad (2.28)$$

where  $\eta$  is the learning rate hyperparameter.

In the context of this thesis, we will investigate whether speaker-related information can be removed using domain adversarial training. Meng et al. [195] train an ASR acoustic model with an additional speaker adversarial branch and show that while it improves the performance of ASR, the intermediate features from the same phoneme belonging to different speakers are qualitatively more identically distributed after using adversarial training. Adi et al. [8] observe that the accuracy of speaker classification performed using the intermediate representation of an end-to-end ASR network is significantly reduced after training it with speaker adversarial branch. Tu et al. [293] also show an increase in emotion recognition when the classifier is trained with a speaker adversarial branch, and the effect of the speaker becomes negligible over the performance of the network. These studies demonstrate the potential of domain adversarial training for designing a privacy-preserving mechanism to identify and remove speaker-related information from speech. Chapter 4 presents our investigation in this direction.

## 2.3.2 Speech transformation

Speech transformation, also called voice transformation [270], is a general modification of speech that aims to shift the perceivable physical attributes of an utterance in a certain direction while leaving the linguistic content unchanged. It is often used as a complementary step after speech synthesis to make the output sound more natural through careful rule-based manipulations of the signal. Speech transformation systems are widely used in speech toolkits due to their efficient real-time nature and general applicability. Some of the interesting applications of speech transformation algorithms are emotion simulation in synthetic speech [39] and conversion of speech to sound like songs [64].

Speech transformation algorithms ease the manipulation of prosodic features of speech, such as speaking rate, loudness, pitch, stress pattern, and in effect the overall speaking style, which originate from the source

part of the vocal tract (i.e., lungs and vocal folds). Although the speaking style is an abstract concept, speech transformation algorithms can be used to map the style of one speaker over the utterance of another speaker. They are typically not designed to achieve a predefined target but a relative shift from the source instead, such as a faster time-scale, a lower pitch, etc. Speech transformation algorithms are also designed to modify the filter characteristics (recall the source-filter model described in Section 2.2.4), which describe the frequency response of the vocal tract. The frequency response can be simply warped in a certain direction by applying, for example, a bilinear function, expressed as  $f(\omega, \alpha) = | -j \ln \frac{z-\alpha}{1-\alpha z} |$ , where  $\omega \in [0, \pi]$  is the normalized frequency,  $\alpha \in (-1, 1)$  is the warping factor, and  $z = e^{j\omega}$ . It is applied to the spectrum with a predefined domain over it [227] to quickly produce perceptibly different voices. Figure 2.15 shows the response of the bilinear function for positive and negative values of  $\alpha$ . When  $\alpha < 0$ , the lower frequency region is compressed and the higher region is stretched, while the reverse happens when  $\alpha > 0$ .

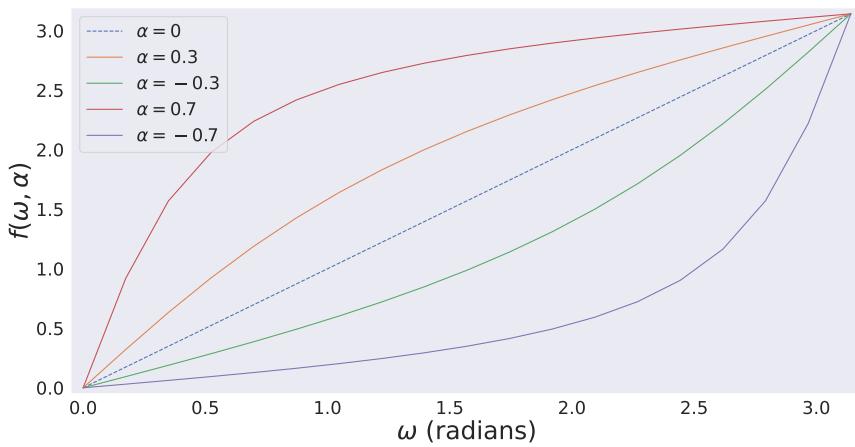


Fig. 2.15 Bilinear function warping the frequency  $\omega \in [0, \pi]$  using positive and negative values of  $\alpha \in \{-0.7, -0.3, 0, 0.3, 0.7\}$ .

Despite their simplicity, speech transformation algorithms are hard to design because they require a clear understanding of the processes related to speech production and perception. Several systematic studies [268, 85] have been conducted to find the acoustic correlates of phonetic processes. Researchers have closely analyzed the source characteristics of speech and defined the parameters that can be modified to perceive certain effects [79, 86, 84]. For example, Stylianou [270] mentions the case that, when someone wants to increase the loudness in a given speech utterance, they must increase the energy of consonant segments rather than vowels because consonants have a short duration yet carry most of the information load in oral communication. This increase in stress also implies an increase in subglottal pressure, thereby an increase in pitch and energy in high-frequency regions. Similarly, increasing the speaking rate also has an effect on pitch. These examples illustrate that the parameters of speech cannot be modified in isolation, and their phonetic and articulatory relationship must be known before implementation, otherwise, the resulting voice loses naturalness.

One of the works that are closely related to this thesis is that of Matrouf et al. [191] who show that speech transformation can be used to transform the voice of impostor speakers such that it closely mimics the mated distribution of the speaker recognition system. The goal of this mechanism is similar to voice spoofing [316] where an impostor is caused to be accepted as a legitimate target speaker by exploiting the vulnerabilities of the speaker recognition system. Such a mechanism can also be used to fool the algorithms used by an attacker to identify the true speaker from a data set by making the mated and non-mated distributions indistinguishable from each other, that is not robustly achieved by the given mechanism in the case when the

attacker is aware of the transformation algorithm. Most of the proposed approaches in this thesis attempt to do this in rigorous settings for achieving robust anonymization.

### 2.3.3 Voice conversion

Voice conversion (VC) [259] aims to convert the voice of a speaker to sound like another, while leaving the linguistic content unchanged. VC has a more specific goal than speech transformation in terms of the precision of the target voice achieved after conversion, thereby it also requires identification and modeling of the exact characteristics that are particular to the source as well as target speaker. We have already described in Section 2.2.5 the different factors which contribute towards the vocal identity of each speaker. They include the overall spectral information, sometimes called *timbre*, and prosodic information, such as pitch, duration, and intensity. VC attempts to learn a mapping from source to target speaker characteristics by modifying these relevant factors. Traditional VC approaches use a vocoder at their core to analyze and synthesize the voice with a feature conversion module to map the source to the target speaker.

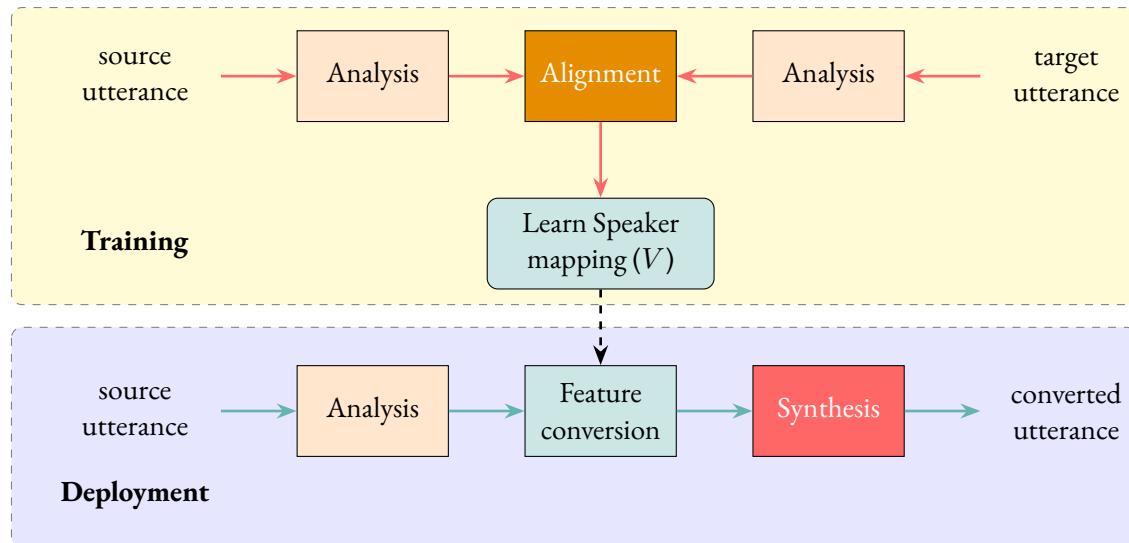


Fig. 2.16 General schema for traditional voice conversion with parallel data. Red arrows indicate training flow, while green arrows show the flow during conversion.

Figure 2.16 shows a general schema for traditional VC where red arrows show the training process of the speaker mapper. The orange blocks are parts of the vocoder (as described in Section 2.2.4) which provides algorithms to extract features, such as pitch, spectrogram, and aperiodicity from speech and re-synthesize the original signal, given these features. During the training phase, the parameters of the frame-wise mapping function  $V$  are learned to convert the source acoustic features into a given target speaker's acoustic features. To learn this mapping effectively, traditional VC approaches require a training set that contains parallel data, i.e., several different speakers uttering the same linguistic content. Parallel data reduces the complexity of VC mapping by keeping the same linguistic content in the source and target utterances, but it raises the issue of alignment because each speaker may have a different speaking rate. As an additional step during training, the source and target features are aligned in time using dynamic time warping (DTW) to deal with varying speaking rates. The learned function  $V$  is eventually used to transform the source features at test time to a particular target speaker's features frame-by-frame. A glaring limitation of this approach is that the parameters of  $V$  have to be learned separately for different source-target speaker pairs. And in order to scale

1 this system, new speakers must be recorded while they utter the same sentences which already exist in the  
 2 training data set.

3 Existing statistical models can be used to design the mapping function  $V$ , such as joint modeling  
 4 of source and target using Gaussian mixture models [286] or exemplar dictionary-based source-target  
 5 mapping [317, 94, 315]. Of course, deep learning has also been employed extensively to build the mapping  
 6 function: feedforward networks [63], RNNs [204], and recently Transformer based models [127] have  
 7 been proposed. The traditionally used DTW has also been replaced by the neural network attention  
 8 mechanism [280] which takes the context into account and gives superior alignment between source and  
 9 target features. However, the breakthrough in VC research has been achieved by *non-parallel* techniques that  
 10 obviate the requirement for parallel data, such as CycleGAN-VC [141], which involves a generator  $G_{S \rightarrow T}$   
 11 which maps source features to target features and a discriminator  $D_T$  that classifies whether the generated  
 12 features belong to the target.  $G_{S \rightarrow T}$  and  $D_T$  are jointly trained to learn an efficient mapping function.  
 13 Since the source and target utterance may not contain the same linguistic content, another generator  $G_{T \rightarrow S}$   
 14 and discriminator  $D_S$  are added to the architecture which learn the inverse mapping from target to source.  
 15 During the training phase, the source features are converted to the target and then back to the source in a  
 16 cycle, and the converted features are verified at each step using their corresponding discriminator. All four  
 17 models are jointly trained in an end-to-end fashion using the cycle-consistency loss function so that the  
 18 linguistic content is preserved.

19 One of the limitations of CycleGAN-VC is that it only supports one pair of source-target speakers per  
 20 model. This limitation is alleviated using a simple modification in the architecture, where the generator  
 21 converts source features to target conditioned upon a speaker embedding provided at training and test time.  
 22 This architecture supports many-to-many voice conversion and is referred to as StarGAN-VC [140]. Another  
 23 class of non-parallel VC is based on the idea of *disentanglement*, which means that during training it learns to  
 24 decompose speech utterances such that the speaker and content information are perfectly separated from each  
 25 other. At test time, source speaker information can be replaced by the target while copying the source content,  
 26 and speech is re-synthesized using the new features to perform VC. Variational autoencoders [122, 124]  
 27 have been used to learn speaker-independent content features, which can generate voice converted speech  
 28 based on a target speaker one-hot vector. Several extensions to this approach have been proposed, where  
 29 speaker embeddings are also learned to have a continuous representation [123], and better disentanglement  
 30 is achieved using speaker adversarial training [128, 46]. Another interesting approach to disentangle speaker  
 31 identity from content is to leverage the fact that the speaker information stays constant throughout the  
 32 utterance, hence it can be removed from the content embeddings by using instance normalization which  
 33 averages out the global statistics [45].

34 Despite the progress made by non-parallel and disentanglement approaches, VC systems suffer from  
 35 a lack of training data and require careful tuning of parameters. Recently researchers have observed the  
 36 similarities between the neural architectures of TTS, ASR, and VC, which have enabled them to propose a  
 37 parameter sharing mechanism either through joint training of these systems or re-using certain components  
 38 which may be well optimized for their task due to the availability of large data sets. Zhang et al. [329] jointly  
 39 train TTS and VC systems in an encoder-decoder architecture by having two encoders, one for source text  
 40 and another one for source speech, which are merged into a single decoder that takes the speaker-independent  
 41 intermediate representation and produces speech in the voice of the target speaker. Zhang et al. [330] propose  
 42 to train a state-of-the-art TTS system and then transfer the decoder of the TTS to the VC system while  
 43 supervising the encoder output of the VC system to be similar to the encoder output of the TTS system  
 44 in order to maintain speaker independence. The similarities between the goals of TTS and VC, and the  
 45 fact that they produce speaker-independent features have allowed several such techniques to be proposed  
 46 recently [327, 126, 183]. Park et al. [216] propose Cotatron VC which leverages speaker-independent

linguistic features coming out of a pre-trained TTS system and targets the speaker's identity to convert the given content into the target speaker's voice. This approach is similar to phonetic posteriorgram-based techniques which we describe below.

Speaker-independent linguistic features can be produced not only using TTS, but also ASR. When extracted from the output layer of an ASR system, they are referred to as *phonetic posteriograms* since they are optimized to classify phonetic information. Phonetic posteriogram features are similar to the BN features (**B**) described in Section 2.2.3, except that phonetic posteriograms are obtained from the output layer of the ASR network. ASR systems have matured quite a lot due to decades of research in this domain. They are also trained on large data sets which are readily available for some popular languages. This makes them desirable for the extraction of rich linguistic features that are also assumed to be speaker-independent. One of the earliest approaches to use phonetic posteriograms for VC was proposed by Sun et al. [271], who used them to generate target acoustic features, but a full VC model is required to be trained for different target speakers. Tian et al. [285] extend [271] to propose a model averaging and adaptation technique which reduces the data requirement for scaling the VC system to new target speakers. Phonetic posteriogram-based VC approaches are gaining popularity [284, 334, 328], and in fact, we also use them for the approaches proposed in this thesis.

## 2.4 Machine learning based anonymization methods

As described in Section 1.1, personal data such as name, address, gender, ethnicity, health conditions, biometric markers, etc. are sensitive because firstly, they may reveal the true identity of a person that is a direct attack on their privacy, and secondly some of these attributes may be embarrassing for them if published. The goal of anonymization is to transform the personal data of individuals such that it can no longer be linked with their true identity while preserving the utility of the data. In this section, we describe some of the general approaches towards anonymization that have been proposed in the machine learning community, not specifically for speech data.

The earliest ML-based anonymization methods were designed to be applied on large databases [290, 202, 65] that gather sensitive personally identifiable data from users, such as medical records, user surveys, travel history, software usage, browser fingerprints, political opinions, sexual preferences, etc. Before publishing such private information, the database curator must ensure that it is *de-identified* or *pseudonymized*, i.e., sensitive attributes (such as social security number, name, address, etc.) are replaced with a random placeholder which can be reversed using a lookup table, or fully *anonymized*, which implies irreversible removal of attributes. Although it may seem that the database is sanitized by the naive approach of removal of sensitive attributes and can be safely released in public, in fact, the database can often be successfully de-anonymized (i.e., the users can be re-identified) using a combination of attributes leading to unique identifiers and/or auxiliary knowledge from public sources as shown in the case of sanitized movie reviews [205], social networks [20], computer networks [55] and DNA sequences [56]. The common attributes that are present in both the sanitized database and the publicly available sources and are used to shortlist the potential identities of anonymous persons are called quasi-identifiers. It has been shown that 87% of the population of the USA can be uniquely identified using just three quasi-identifiers: zip code, gender, and date of birth [274]. This naive approach is clearly not sufficient to unlink users' data from their identities, hence several formal models of privacy, such as  $k$ -anonymity [245] and differential privacy [78], have been proposed.

**$k$ -Anonymity** The intuition behind these methods is that the machine learning models do not need precise information about each subject. Instead, they aim to infer some aggregate statistics about the

1 population in the database, so releasing the data in its original format is not required. In the re-identification  
 2 scenarios concerning movie reviews, health records, etc., the attacker would generally find a unique set of  
 3 quasi-identifiers for an individual, which is the motivation for imposing  $k$ -anonymity [245] over tabular  
 4 databases. It is defined as the property of the anonymized database such that there are at least  $k$  individuals  
 5 for every set of quasi-identifiers in it. This gives the user the ability to hide in a crowd and the attacker's  
 6 chances are reduced to  $\frac{1}{k}$  to attribute a particular data point to the correct user. It is implemented using  
 7 two basic operations: *suppression*, that is to remove or replace quasi-identifier values with placeholders (e.g.,  
 8 the trailing digits of a zipcode can be replaced by the \* symbol), and *generalization*, that is to club together  
 9 values in a particular column using ranges (e.g., age > 30). The anonymity gets stronger as the value of  $k$   
 10 increases, but of course, it implies a trade-off with the utility of the released data set.

11 There have been efforts to learn efficient machine learning models using a data set that has been  
 12 anonymized using  $k$ -anonymity, by searching for optimal hyperparameters [26], efficient clustering of  
 13 rows [174], and full-domain generalization [165]. Moreover, there are several extensions proposed to address  
 14 issues in  $k$ -anonymity, e.g., the lack of diversity in the anonymized data set is handled by  $l$ -diversity [186], and  
 15 the vulnerability posed by the distinct distribution of sensitive attributes is solved using  $t$ -closeness [170].  
 16 Indeed, such modification of the database implies a significant reduction in utility, and yet there remain  
 17 several vulnerabilities in  $k$ -anonymity due to the lack of randomization in the aggregation and the query  
 18 mechanism. This is the motivation for the differential privacy (DP) [78] paradigm, which provides the  
 19 strongest privacy guarantees at present and is deployed at major corporations [53] including Apple, Google,  
 20 and Microsoft. It was also used to publish population data for the 2020 US census [5].

21 **Differential Privacy** In its simplest form, DP can be explained as a randomized response [308] when a  
 22 binary question is asked to the user with yes/no as the response. In this scenario, a coin is tossed each time  
 23 a response is to be recorded. If the coin says heads, then the true response is recorded, and if the coin says  
 24 tails, then another coin is tossed. Based on the outcome of the second coin, the response may be recorded  
 25 as yes if heads or no if tails. The randomness involved in this mechanism allows the users to have some  
 26 plausible deniability towards their response and hence protects their privacy to some extent. As the number  
 27 of participants increases, the aggregate result for the question becomes more accurate. At this point it is  
 28 to be observed that, while DP can be used to anonymize databases, more generally it is a property of the  
 29 algorithm which executes query functions over databases and generates noisy responses to preserve privacy.  
 30 This mechanism ensures that the output distribution of the DP algorithm is statistically indistinguishable,  
 31 no matter whether the data of a subject is present in the database or not. Hence, for example, the participants  
 32 of a DP-enabled survey can rest assured that their data does not make a significant difference in the aggregate  
 33 response of the survey, and can go on to safely participate in it.

34 DP [75] provides a rigorous probabilistic way to quantify the privacy leakage of an information release  
 35 process. DP also comes with strong mathematical properties and a powerful algorithmic framework [77].  
 36 For these reasons, DP and its variants have become the gold standard notion of privacy in machine learning  
 37 and many other scientific fields. Traditionally, DP is defined using the notion of “neighbouring” databases  
 38 which differ on at most one record. Let  $\mathcal{D}$  be the universal set of databases, and  $d, d' \in \mathcal{D}$  be two databases  
 39 with  $|d|_1$  and  $|d'|_1$  as their respective sizes, then the  $\ell_1$  distance between them is given by  $|d - d'|_1$ , which  
 40 measures how many rows differ between  $d$  and  $d'$ . Now we can define the criteria of response queried from  
 41 these databases which will satisfy the DP guarantee.

**Definition 1 (Differential privacy)** *Let  $\mathcal{A}$  be a randomized algorithm which executes queries over any arbitrary database belonging to  $\mathcal{D}$ , and let  $\epsilon > 0$ . We say that  $\mathcal{A}$  is  $\epsilon$ -differentially private ( $\epsilon$ -DP) if for any*

---

 2.4 Machine learning based anonymization methods

$d, d' \in \mathcal{D}$  such that  $|d - d'|_1 = 1$  and any  $S \subseteq \text{range}(\mathcal{A})$ :

$$\Pr[\mathcal{A}(d) \in S] \leq e^\epsilon \Pr[\mathcal{A}(d') \in S],$$

where the probabilities are taken over the randomness of  $\mathcal{A}$ .

DP essentially requires that the probability of any output does not vary “too much” (as captured by  $\epsilon$ ) when changing the input. The smaller  $\epsilon$ , the stronger the privacy guarantee, hence epsilon is called the privacy budget or privacy leakage.

DP possesses a number of desirable properties. First, any function of an  $\epsilon$ -DP algorithm remains  $\epsilon$ -DP (*robustness to post-processing*). Second, one can easily keep track of the privacy guarantees across multiple analyses (*composition*). In particular, given  $K$  algorithms that satisfy  $\epsilon$ -DP, executing them on the same data and releasing their combined outputs is  $K\epsilon$ -differentially private.

In this thesis, we use a variant of traditional DP, called the local-DP model [68] or the fully distributed model which is more suitable when there is no trusted third party to collect raw data. In this model, it is assumed that the users’ private data is not aggregated in a central database; instead, each user is the sole owner of their data and they can respond to questions in a differentially private manner. This nullifies the possibility of any data theft or malicious use by the database curator, hence it is considered as a stronger privacy model. Local-DP can be defined over individual data points that may originate from features of raw data or intermediate representations of a neural network.

**Definition 2 (Local differential privacy)** Let  $\mathcal{A}$  be a randomized algorithm taking as input a data point in some space  $\mathcal{X}$ , and let  $\epsilon > 0$ . We say that  $\mathcal{A}$  is  $\epsilon$ -local differentially private ( $\epsilon$ -LDP) if for any  $x, x' \in \mathcal{X}$  and any  $S \subseteq \text{range}(\mathcal{A})$ :

$$\Pr[\mathcal{A}(x) \in S] \leq e^\epsilon \Pr[\mathcal{A}(x') \in S],$$

where the probabilities are taken over the randomness of  $\mathcal{A}$ .

It is to be noted that local DP is equivalent to standard DP for databases of size 1, and hence,  $x$  and  $x'$  are always neighboring.

A standard way to design differentially private algorithms is based on output perturbation. A basic approach is to rely on the Laplace mechanism, which consists in adding Laplace noise calibrated to the  $\ell_1$ -sensitivity of the (non-private) function one would like to compute on the data [76].

**Definition 3 (Laplace mechanism)** Let  $f : \mathcal{X} \rightarrow \mathbb{R}^d$  and let the  $\ell_1$ -sensitivity of  $f$  be defined as

$$\Delta_1(f) = \max_{x, x' \in \mathcal{X}} |f(x) - f(x')|_1.$$

Let  $\eta = [\eta_1, \dots, \eta_d] \in \mathbb{R}^d$  be a vector where each  $\eta_i \sim \text{Lap}(\Delta_1(f)/\epsilon)$  is drawn from the centered Laplace distribution with scale  $\Delta_1(f)/\epsilon$ . The algorithm  $\mathcal{A}(\cdot) = f(\cdot) + \eta$  is  $\epsilon$ -local DP.

This approach can also be used to randomize data points directly (*input perturbation*), which corresponds to the case where  $\mathcal{X} \subseteq \mathbb{R}^d$  and  $f$  is the identity function. However, adding noise to raw data often destroys utility. A better strategy is to perturb carefully designed feature representations of the data, as done in existing work on image [185] and text [29, 184, 90]. Our proposed approach to add DP noise in speech features is presented in Section 6.2.

Membership inference attacks [17] are a privacy threat related to trained classifier models where an attacker can deduce if a given data point is present in the training set of the classifier. It has been shown [243]

1 that DP can provably bound the accuracy of such attacks and enhance the privacy of these models considerably.  
2 Although DP provides strong mathematical guarantees, the addition of large random noise with the  
3 assumption of worst-case adversaries significantly reduces the utility and increases the sample complexity of  
4 the model [91, 69].

## 5 2.5 The speaker anonymization task

6 This thesis derives its relevance from Article 25 of the GDPR which mandates the data controller to  
7 implement all possible technical and organisational measures to ensure the protection of personal data  
8 following the principles of privacy by design and by default. More precisely, speech data is protected by  
9 Article 9 of the GDPR which prohibits the processing of biometric personal data that reveals the gender,  
10 racial or ethnic origin, or health indicators of the data subject. This restriction has a direct impact on the  
11 potential scientific and commercial advancements, therefore it is relaxed by Recital 26 of the GDPR which  
12 freely allows the processing of *anonymous* data that can no longer be used to identify the original data subject.  
13 Article 4(5) of the GDPR mentions “pseudonymization” as a possible technical measure to achieve privacy  
14 by design, which is the task of processing personal data such that it can no longer be attributed to the data  
15 subject without using additional information that must be secured separately to safeguard the identity of  
16 the data subject. Such processing may prove to be risky if the so-called additional information is accessed by  
17 a malicious entity, hence a stronger approach, namely anonymization, is explored in this thesis. Although  
18 the GDPR does not explicitly mention anonymization, ISO/IEC 29100:2011 [2] defines it as the process  
19 by which the personally identifiable information is *irreversibly* altered such that the original data subject  
20 cannot be identified directly or indirectly, either by the data controller alone or in collaboration with any  
21 other party.

22 As mentioned before, the goal of this thesis is to completely unlink a speaker’s identity from their speech  
23 utterances while maintaining the usefulness of the signal. In a generalized sense, this task is perceived as  
24 privacy-preserving data publishing [95] in the literature which envisages a similar vision of releasing useful  
25 data sets in the public domain for scientific and commercial progress without compromising individuals’  
26 privacy. Users are the default owners of their data, therefore publishers must anticipate potential hostile  
27 activity against them using their data and take proactive measures to mitigate any such possibility. Fung  
28 et al. [95] present a comprehensive survey of traditional privacy-preserving data publishing approaches  
29 that are applicable to relational databases. They describe the actors involved in data publishing (user, data  
30 publishers, and recipients), the types of attacks and the effectiveness of privacy models against each attack,  
31 anonymization methods and their fundamental operations, and finally metrics to assess the performance of  
32 privacy-preserving data publishing methods in terms of privacy and utility.

33 Due to legal and technological awareness in the society, privacy-preserving data publishing approaches  
34 specific to speech data have gained prominence in the past decade and several different methodologies have  
35 been proposed to achieve anonymization<sup>10</sup> to some extent. Nautsch et al. [208] attempted to describe  
36 the legal and technological advances with respect to privacy-preserving speech processing, but they only  
37 mention the techniques related to cryptographic methods and their evaluation. The methods of speaker  
38 anonymization that are most relevant to this thesis can be divided into five categories based on the type  
39 of technology they use: noise addition, speech transformation, voice conversion, speech synthesis, and  
40 adversarial learning.

---

<sup>10</sup>Legally speaking, the term “anonymization” refers to a method that fully achieves this goal. Following [288], we use it in a broader sense to refer to a method that aims to achieve this goal, even when it has failed to do so.

---

2.5 The speaker anonymization task

43

**Anonymization attempts using noise addition** Ahmed et al. [9] present an end-to-end ASR method that injects DP noise at various levels in the pipeline to protect the user’s identity and publish only the anonymized transcription instead of the speech signal. Publishing the text content of spoken data in a privacy-preserving manner is helpful to safeguard speakers’ identity, but it limits the usage of speech. Hashimoto et al. [116] experiment with a different range of bandpass filtered noise to be added to speech signal to degrade speaker recognition performance in terms of increase in EER and preserve intelligibility. Although this method allows publication of speech data, it is difficult to manually calibrate the noise for unforeseen conditions, e.g., different languages, ambient noise, etc.

**Anonymization attempts using speech transformation** Speech transformation is considered the most convenient anonymization method since it does not require large data sets for training machine learning models; instead, it can be performed using careful signal processing based manipulations of speech parameters. Cohen-Hadria et al. [50] perform anonymization of voice recordings by using a low pass filter to remove formants and inverting the MFCCs to obfuscate speaker information while preserving the acoustic scene. Qian et al. [227] transform the spectrogram frequency scale of original speech by applying a composition of two nonlinear functions with random parameters. The resilience of this technique is dependent on the secrecy of these parameters. They also perform sensitive keyword substitution to completely sanitize the speech for publishing. Patino et al. [218] present the latest speech transformation method as a baseline for the first Voice Privacy Challenge [287] where the pole angles of the linear prediction (LP) spectral envelope are altered using McAdams coefficient ( $\alpha_M$ ) [192]. Specifically, the filter coefficients for each frame are derived using LP source-filter analysis, which are then used to extract the real- and the complex-valued pole positions. Thereafter, the angle of the complex-valued poles is raised to the power of a pre-determined  $\alpha_M$ , causing the associated formant spectrum to expand or contract. The new complex-valued pole positions and the unchanged real-valued poles are then converted back to filter coefficients, and combined with the residual source information to resynthesize an anonymized time-domain speech frame. Gupta et al. [109] significantly improved this work by proposing the modification of both the pole angles as well as the pole radii of the LP spectrum. Although the parameter manipulations performed by speech transformation methods seem perceptually reasonable, they are easy to break using machine learning methods [267], hence they provide weak protection against privacy attacks.

**Anonymization attempts using VC** Voice conversion (VC) methods are the earliest and most obvious choice for speaker anonymization since their goal is to transform the voice of a given speaker into that of another speaker whose voice characteristics are known beforehand. Jin et al. [135] present the first known approach towards speaker anonymization by transforming any given source speaker to a single target voice present in the Festival speech synthesis system [33], and show that it performs well in terms of drop in the speaker identification accuracy. Bahmaninezhad et al. [22] convert a given source speaker into the average of all speakers of the same gender. Pobar and Ipšić [221] pre-train a set of speaker transformations and identify the speaker at test time to select one of the corresponding transformations. These methods are hardly applicable in practice since they require the source speaker to be present in the training set of the VC system and, in the context of anonymization, the amount of speech from the original speaker is often limited to one utterance. To relax this constraint, Magariños et al. [187] find the closest source speaker in the training set and apply one of the corresponding transformations. Yoo et al. [319] presents a many-to-many CycleGAN variational autoencoder-based VC method for speaker anonymization which takes a one-hot vector for the target speaker present in the training set. They experiment with several distributions of training speaker proportions as the target but yet do not allow external speaker identities to be used at run time.

1 **Anonymization attempts using speech synthesis** Techniques based on speech synthesis have also been  
2 proposed to relax the requirement of having source speakers in the training set of the anonymization system.  
3 For instance, Justin et al. [138] transcribe speech into a diphone sequence and re-synthesize it using a single  
4 target. These methods suffer from three limitations. First, they still result in a limited set of target speakers  
5 or speaker transformations, which prevents the original speaker from choosing an arbitrary unseen speaker  
6 as the target. Second, using a real speaker’s voice as the target raises ethical concerns. Third, the conversion  
7 of speech to a sequence of discrete tokens as in [138] is error-prone and destroys all the paralinguistic and  
8 extralinguistic attributes. This motivates the objective of converting the original speaker’s voice into an  
9 arbitrary, imaginary *pseudo-speaker*’s voice without relying on a transcription step. Speaker embeddings such  
10 as x-vectors [262] provide the continuous representation needed to define and generate such pseudo-speakers.  
11 Fang et al. [83] address this objective using a speaker-independent speech synthesis system. They select  
12 x-vectors within an external pool of speakers and average them to obtain a target *pseudo-speaker* x-vector. This  
13 x-vector, along with a representation of the original linguistic and intonation contents, is provided as input to  
14 a neural source-filter (NSF) based speech synthesizer [304] to produce anonymized speech. Han et al. [113]  
15 extend the framework presented in [83] to select a single target x-vector at random within a maximum  
16 distance from the original x-vector that satisfies a privacy metric based on differential privacy. Although these  
17 techniques manage to alleviate the three limitations mentioned before, they use weak evaluation criteria with  
18 an assumption that the attacker does not know about the usage and the parameters of the anonymization  
19 algorithm.

20 **Anonymization attempts using adversarial learning** Recently, there has been some research towards  
21 speaker anonymization using adversarial training which implicitly models speaker information and removes  
22 it from the intermediate representations of a neural network that is optimized for some utility task, thereby  
23 learning privacy-preserving model parameters. One of the first approaches in this direction is investigated  
24 as a part of this thesis which is described in detail in Chapter 4. Espinoza-Cuadros et al. [82] present an  
25 autoencoder-based approach, which reconstructs the speaker representation by adversarially removing  
26 source speaker information from it, and use the new representation as the target pseudo-speaker. In a similar  
27 context, Champion et al. [42] examine the hypothesis that the linguistic features used for speech synthesis  
28 contain speaker information and use speaker adversarial training similar to [265] to mask the identities.

## 29 **2.6 Summary of techniques**

30 In this section, we succinctly recall the relevance of the abovementioned techniques for this thesis and  
31 mention how exactly they are reused or adapted in the following chapters.

32 We first delve into the details of speech generation using the vocal tract, which makes it clear that each  
33 person has a personal physiology that reveals cues to substantiate his/her identity from the speech signal. The  
34 specific configurations of articulators, that generate the phoneme sounds, also exhibit personally identifiable  
35 characteristics. Although in this thesis we do not perform any phoneme-specific analysis with respect to  
36 privacy, such phenomena are crucial to understand that the speaker’s identity is highly entangled with the  
37 linguistic properties of sounds, and may possess identity markers in terms of nativeness or accent [59]. These  
38 properties and identity markers are recognized by examining the speech signal in the time domain or the  
39 frequency domain. The time-domain signal is a useful digital representation of sound and it exhibits several  
40 interesting properties which enable us to define either the local units of speech, i.e., phonemes, or global  
41 speaker-related characteristics like speaking rate [160]. It is noteworthy to mention that the time-domain  
42 signal also contains the effect of ambient noise, the microphone quality, and the subtle variations in the  
43 vocal tract even to produce the same linguistic content. Due to the dominant effect of speaker characteristics

---

## 2.6 Summary of techniques

**45**

among these factors, it is certainly attainable to infer the speaker's identity yet it is difficult to disentangle these variations and capture parameters in the time-domain signal.

Evidently, it is much easier to analyze the speech signal after transforming it into the frequency domain. The spectrogram and the features derived from the frequency domain representation, such as MFCCs, are widely used to identify phonemes, speakers, emotions, etc. These features exhibit numerical patterns that correlate with the perceptual properties of human speech, for example, the fundamental frequency correlates with the pitch, the harmonic structure in the spectrum indicates voicing, and the position of the formant peaks characterize a phoneme. The goal of this thesis is to efficiently identify features corresponding to speaker information and manipulate them such that the speech signal is no longer attributable to the original speaker. We mentioned previously that the speaker information is present throughout the signal as the source characteristic, and shows distinctive cues in the spectral tilt, the nasalization patterns, and the speaking rate of an utterance. Although a rule-based approach, using a combination of acoustic cues and spectral features, can be devised to isolate the factors causing speaker distinction in a controlled setting, it is infeasible to model unknown variations caused by the effect of ambient noise, emotional states, conversational settings, and languages. Hence, strong statistical models, such as deep neural networks, are employed to project speaker or phoneme-related information to a high-dimensional hidden space where discrimination and identification are easier.

TDNNs are well suited for modeling the temporal dependencies in a speech signal, hence, they are used to design effective ASR (except for the end-to-end systems) and ASI models. ASR models are extensively used for two main purposes in this thesis. First, they are used as a key building block of the anonymization pipeline proposed in this thesis. And second, they are used to evaluate the private representations generated by the proposed anonymization techniques in terms of intelligibility. ASI models are used for two main purposes: 1) to model speaker information and generate useful representations, such as the x-vectors, and 2) to test the resilience of the anonymization techniques against re-identification attacks. Chapter 4 investigates whether private representations can be generated by an ASR network when it is trained in a domain-adversarial manner with an ASI adversary. In addition, ASI evaluation metrics are adapted for measuring the degree of privacy protection as explained in the next chapter.

Techniques that allow the manipulation of speaker-related properties of speech and the generation of a new speech signal, such as speech synthesis, speech transformation, and voice conversion, are crucial for the anonymization approaches proposed in this thesis. The remaining chapters, except Chapter 4, leverage these techniques to replace the original speaker's identity in the speech signal with a new identity by manipulating either the inputs or the parameters of these techniques. Since they generate an intelligible speech signal, they can be directly evaluated in terms of their utility for the task of anonymization. The speech data set processed by these techniques can also be published and easily evaluated in terms of speaker re-identification attacks. Finally, we employ differential privacy based anonymization techniques to remove residual speaker information from the speech signal and provide formal guarantees of privacy protection.

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## Chapter 3

# Privacy evaluation using Informed attackers

Facts do not cease to exist because they are ignored.

---

Aldous Huxley

In this chapter, we design the actors involved in the speaker anonymization process, we discuss the potential actions they can take within the encompassing threat model that connects them using their goals. This threat model is then used as a guide to design and evaluate the privacy protection techniques. Specifically, we propose the concept of malicious entities, i.e., attackers who possess a certain degree of knowledge about the anonymization scheme, and use them to rigorously evaluate the level of privacy protection achieved by our schemes. Further, a brief account of the performance metrics that are employed to evaluate these techniques in terms of privacy and utility is presented, followed by a detailed comparison of three privacy metrics to assess their ability to express useful information and vulnerabilities of the proposed techniques. Finally, a preliminary study is conducted using one speech transformation and two voice conversion-based methods to establish the role of the actors in the threat model and validate the claim that ‘knowledgable’ attackers measure the level of privacy protection authentically.

In Section 3.1, we present the threat model, the actors and their goals, and discuss the notion of attackers’ knowledge. Section 3.2 and 3.3 describe the VC techniques and strategies that are used to design the anonymization schemes proposed in this chapter. The performance metrics to evaluate the level of privacy protection and utility achieved by the proposed schemes are mentioned in Section 3.4. Section 3.5 describes the experimental setup including data sets, VC algorithm settings and attacker designs. Different attackers are compared in Section 3.6, while different privacy metrics are compared in Section 3.7. Finally, Section 3.8 summarizes the main findings of this chapter and leads the road for future chapters.

### 3.1 Attack model and the notion of attackers’ knowledge

As explained in Section 1.1, speech data exhibits biometric characteristic of human beings which must be protected to safeguard the identity of individuals. According to the ISO/IEC International Standard 24745 on biometric information protection [1], publicly available biometric references must be *irreversible* and *unlinkable* for full privacy protection. Such protection must be resilient to re-identification attacks that may be strengthened using auxiliary information about the data set or the method of protection.

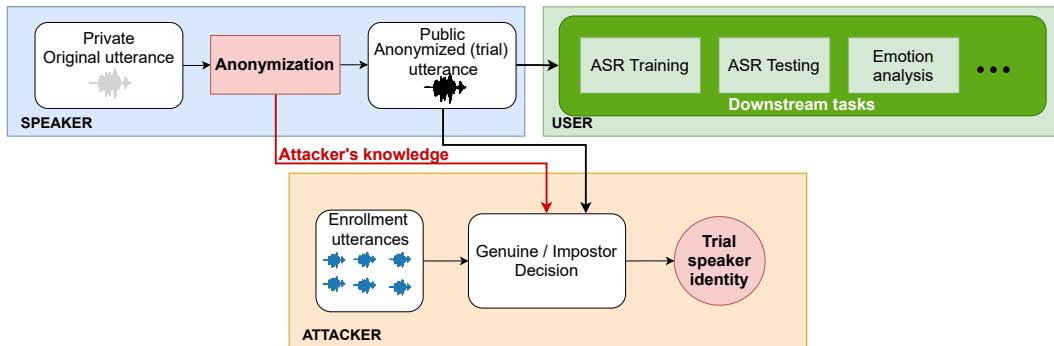


Fig. 3.1 Considered threat model. *Speakers* anonymize speech to conceal their identity before publication; *attackers* use biometric technology and knowledge of the anonymization method to re-identify it; *users* (e.g., speech technology companies) use the published data for downstream tasks such as ASR training.

Throughout this thesis, we consider the following threat model. *Speakers* process their voice through an *anonymization* technique. This anonymization step takes as input one or more *private speech* utterances along with some configuration parameters, and outputs a new speech signal or some kind of derived representation. The transformed utterances from one or more speakers form a *public speech* data set that is processed by a third-party *user* for, e.g., ASR training/decoding or any other downstream task. Given a public data set of anonymized speech (or speech representation in the form of feature vectors) contributed by several speakers, an attacker records/finds a sample of speech of a speaker and attempts to find which utterances in the anonymized data set are spoken by this speaker, possibly leveraging some knowledge about the anonymization method shown as the red arrow in Figure 3.1.

Formally, an attacker has access to two sets of utterances: *A* (*enrollment/found data*) and *B* (*trial/public speech*), but knows the corresponding speakers in *A* only. The attacker designs a linkage function  $LF(a, b)$  that outputs a score for any  $a \in A$  and  $b \in B$ . Typically, this score is a similarity score obtained through a speaker verification system. The attacker then makes a decision (genuine vs. impostor) based on this score. A good speaker anonymization method must defeat such *linkage attacks* by concealing the speaker identity, while preserving the utility of speech for data *users* as measured for instance by the perceived speech naturalness and intelligibility and/or the performance of downstream tasks such as training an automatic speech recognition (ASR) system, thereby achieving a suitable privacy/utility trade-off. Figure 3.1 shows the three actors involved in this model, namely the *speaker*, the *attacker* and the *user*, along with their actions. The goals of the speaker and the user are intimately linked, while the attacker operates independently.

Crucially, all past studies assumed a weak attack scenario where the attacker is unaware that an anonymization method has been applied to the found data [83]. This raises the concern that the privacy protection may entirely rely on the secrecy of the design and implementation of the anonymization scheme, a principle known as “security by obscurity” [197] that has long been rejected by the security community. There is therefore a strong need to evaluate the robustness of the anonymization to the knowledge that the adversary may have about the transformation. In practice, such knowledge may for instance be acquired by inspecting the code embedded in the user’s device or in an open-source implementation.

As opposed to past studies, different linkage attacks are considered in this thesis depending on the attacker’s knowledge of the anonymization method, as illustrated in Figure 3.2. At one end of the continuum, an *Ignorant* attacker is unaware of the speech transformation being applied, while at the other end an *Informed* attacker can leverage complete knowledge of the transformation algorithm. In between, a *Lazy-Informed* attacker may know the voice transformation algorithm but does not exploit it to the full capacity due to computational constraints, hence no model re-training is performed. Re-training the speaker

## 3.2 Voice conversion methods

49

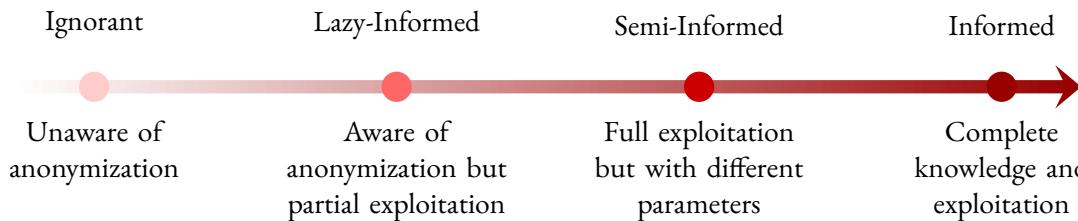


Fig. 3.2 Increasing degree of attacker’s knowledge that determines the strength of re-identification attack. Intermediate points on the continuum can be simulated as attack scenarios.

identification model can reverse the effect of anonymization, as in case of a *Semi-Informed* attacker which fully exploits the voice transformation algorithm by re-training but is unaware about its exact parameter values. In this chapter, we conduct preliminary investigations to assess the strength of *Ignorant*, *Semi-Informed*, and *Informed* attackers, and establish the validity of the proposed attack continuum. Our experiments evaluate three VC methods with different target speaker selection strategies in various attack scenarios to study unlinkability in the spirit of ISO/IEC 30136 standard [3]. In each scenario, we measure how well each VC method protects the speaker identity against attackers that leverage state-of-the-art speaker verification techniques based on i-vectors [61] or x-vectors [262] to design linkage attacks. The *word error rate* (WER) achieved by a state-of-the-art end-to-end automatic speech recognizer [310] is also reported as measure of utility for the data user. While a formal listening test is beyond the scope of this thesis, a few samples of converted speech are available for informal comparison.<sup>1</sup>

## 3.2 Voice conversion methods

**Brij:** maybe use a plot to illustrate the three methods?

The criteria for selecting the VC methods for the experiments performed in this chapter are that they must be **1) non-parallel**, i.e., do not require a parallel corpus of sentences uttered by both the source and target speakers for training — this is important from both privacy and technical perspective. The privacy risk arises due to the scarcity of parallel corpora which limits the data publisher to choose among a few openly available targets, thereby increasing the risk of an inversion attack. Furthermore, it requires the source speaker at test time to be present in the parallel corpora, which limits the usage of such a system to a selected few speakers; **2) many-to-many**, i.e., allow conversion between arbitrary sources and targets so that any speaker in a large corpus can be selected as the target ; **3) source- and language-independent**, i.e., do not require enrollment sentences for the source speaker and do not rely on language-specific ASR or phoneme classification — this is important from a usability perspective as it frees the user from the burden of enrolling and it is applicable to any language (including under-resourced ones), and from a privacy perspective since enrollment translates into the storage of a voiceprint which poses even greater privacy threats.

The third criterion is quite strict: many VC methods, such as StarGAN-VC [140] or the ASR-based method in [83], do not satisfy it. We found that the vocal tract length normalization (VTLN) based methods in [226, 273] and the one-shot method in [45] satisfy all criteria. In this chapter, we use models trained over English speech [215] but do not use any other linguistic resources such as transcriptions, hence in principle they may be applicable to other languages as well.

<sup>1</sup>[https://github.com/brijmohan/adaptive\\_voice\\_conversion/tree/master/samples](https://github.com/brijmohan/adaptive_voice_conversion/tree/master/samples)

### **3.2.1 VoiceMask**

VoiceMask is described in [226] as the frequency warping method based on the composition of a bilinear function, expressed as  $f(\omega, \alpha) = \left| -i \ln \frac{e^{i\omega} - \alpha}{1 - \alpha e^{i\omega}} \right|$ , and a quadratic function, given by  $g(\omega, \beta) = \omega + \beta(\frac{\omega}{\pi} - (\frac{\omega}{\pi})^2)$ . Here  $\omega \in [0, \pi]$  is the normalized frequency,  $\alpha \in [-1, 1]$  is the warping factor for the bilinear function, and  $\beta > 0$  is the warping factor for the quadratic function. Therefore, the warping function is of the form  $g(f(\omega, \alpha), \beta)$ . The two parameters,  $\alpha$  and  $\beta$ , are chosen uniformly at random from a predefined range which is found to produce intelligible speech while perceptually concealing the speaker identity. In the following, we apply this transform to the spectral envelope rather than the pitch-synchronous spectrum as in the original paper. In addition, we apply logarithm Gaussian normalized pitch transformation (see [176]) so as to match the pitch statistics of a target speaker<sup>2</sup>.

The authors claim that this transformation is difficult to inverse when the parameter values are unknown because they are randomly selected from a large interval. However, VoiceMask uses the same parameter values to warp the spectra at each time step of the utterance. This approach is quite limited to conceal the identity of the source speaker and to mimic the target speaker because it warps the entire frequency axis in a single direction.

### **3.2.2 VTLN-based voice conversion**

VTLN-based VC [273] represents each speaker by a set of centroid spectra extracted using the Cheap-Trick [200] algorithm for  $k$  pseudo-phonetic classes. These classes are learned in an unsupervised fashion by clustering all speech frames of all utterances from this speaker. For each class of the source speaker, the procedure finds the class of the target speaker and the warping parameters that minimize the distance between the transformed source centroid spectrum and the target centroid spectrum. All speech frames in that class are then warped using a power function. Similarly to above, we apply this warping to the spectral envelope and also perform Gaussian normalized pitch transformation so as to match the pitch statistics of the target. Compared to VoiceMask, this approach warps the frequency axis in different directions over time. The parameters of this method include the number of classes  $k$  and the chosen target speaker.

Although this algorithm does not require parallel data, it necessitates the storage of  $k$  set of parameters for each speaker who can be used either as the source or the target. Increasing the value of  $k$  might enhance the quality of the output speech because distinct phonetic classes may obtain their own set of parameters rather than a set of averaged parameters over similar classes. It might also have some effect on the strength of anonymization due to an increase in the number of frequency warping parameters within an utterance, i.e., the frequency axis can be warped in more directions than before. Within the scope of this chapter, we fix the value of  $k$  based on the author's recommendations [273] and do not investigate the effect of changing it.

### **3.2.3 Disentangled representation based voice conversion**

The third approach is based on disentangled representation of speech as proposed in [45, 294]. The core idea is that speaker information is statically present throughout the utterance but content information is dynamic. This approach is based on a neural network transformation and uses a *speaker encoder* and a *content encoder* to separate the factors of variation corresponding to speaker and content information. The only parameter of this method is the chosen target speaker.

---

<sup>2</sup>Strictly speaking, VoiceMask is a voice transformation method rather than a VC method: pitch is converted from the source speaker to a target speaker, but the spectral envelope is not related to a particular target speaker.

## 3.3 Target selection strategies and exploitable parameters

51

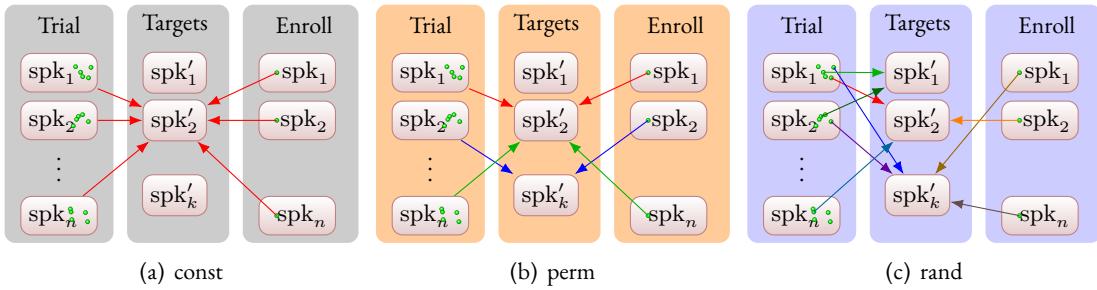


Fig. 3.3 Three target selection strategies: const, perm and rand. Trial utterances are publicly released data set, while enroll utterances are found data used by the attacker. Utterances are shown by small green balls, and the arrows indicate the mapping between original to target speakers.

### 3.3 Target selection strategies and exploitable parameters

In this chapter, we consider that the VC function and the sets of possible parameter values are known to all users. Each user records his/her voice on his/her device and applies a VC scheme locally before sending it to a public database. The threat model (refer Figure 3.1) considers an attacker who performs a linkage attack to try to identify which converted utterances in this public database are spoken by a particular user. To this end, it is assumed that the attacker has access to a small amount of found speech from this user shown as the enrollment utterances in Figure 3.1 (and potentially some additional public resources, such as benchmark speech processing data sets to train generic speaker models).

In the following section, three parameter selection (a.k.a. target selection) strategies are defined for the three VC methods above, which can be seen as key ingredients of a “private-by-design” speech processing system. Thereafter the knowledge of exploitable parameters is described that an attacker trying to compromise the system could have about the VC function and the target selection strategy.

#### 3.3.1 Target selection strategies

Three possible target selection strategies are considered which act as the core part of the anonymization and the methods of defense against the re-identification attack, hence they are also called privacy protection strategies. They are introduced in the increasing order of their randomness. In strategy *const* as shown in Figure 3.3(a), the VC function is constant across all users and all utterances. This means choosing a unique target speaker and, in the case of VoiceMask, fixed values for  $\alpha$  and  $\beta$ . In strategy *perm* depicted in Figure 3.3(b), the conversion parameters are chosen at random once by each user. In other words, when a user downloads the VC module on his/her device, he/she selects a personal target speaker and, in the case of VoiceMask, personal random values for  $\alpha$  and  $\beta$ . Finally, in the *random* strategy illustrated in Figure 3.3(c), each time a user applies VC to an utterance, a random set of parameters is drawn, i.e., a random target speaker is selected and, in the case of VoiceMask, random values are drawn for  $\alpha$  and  $\beta$ .

#### 3.3.2 Exploitable parameters

As explained in Section 3.1, the effectiveness of anonymization is evaluated using three types of attackers based on the extent of their knowledge about the VC function and its exploitable parameters. An *Ignorant* attacker is not aware that VC has been applied at all. In contrast, an *Informed* attacker knows the VC method and its exact parameter values (i.e., the chosen target speaker and the values of  $\alpha$  and  $\beta$ ). One may argue that an *Informed* attacker is not very realistic (except for the *const* strategy), while an *Ignorant* attacker is

1 very weak. Between these two extreme cases, various types of attackers can be defined. For instance, we  
 2 consider a *Semi-Informed* attacker who knows the chosen VC method (VoiceMask, VTLN, or disentangled  
 3 representation) and the target selection strategy (*const*, *perm*, or *random*), but not the actual target (i.e., the  
 4 actual target speaker or the value of  $\alpha$  and  $\beta$ ). This is arguably more realistic since the VC algorithm and the  
 5 target selection strategy may be open-source, while (except for the *const* strategy) the target chosen by the  
 6 user is much less easily accessible.

7 It is important to note that many concrete instances of attackers of the above types can be designed as  
 8 illustrated in Figure 3.2, and finding out the “best” attacker of a particular type is a hard problem. In the  
 9 experiments section, it is proposed how attackers may exploit these different levels of knowledge based on  
 10 the assumptions defined above. A more exhaustive investigation of the design of attackers is left for later  
 11 part of this thesis.

## 12 3.4 Performance metrics

13 Historically, the usual metrics employed in the speaker verification community have been used to assess the  
 14 (in)ability of an attacker to recognize the speaker, which is considered as a proxy for the degree of privacy  
 15 protection. On the other hand, utility is considered a nebulous concept because it depends on the user,  
 16 who may want to use the anonymized speech for transcription, scientific analysis, or broadcast purposes.  
 17 Here, concrete privacy and utility metrics are presented which are used throughout this thesis to evaluate  
 18 the proposed speaker anonymization methods.

### 19 3.4.1 Privacy measures

20 In the context of the attacker’s continuum, privacy may be measured as the deterioration in the attacker’s  
 21 ability to identify the original speaker from the new representation. If the original speaker is present in the  
 22 training set of the system, then an ASI, i.e. *closed-set* identification, can be used over anonymized speech  
 23 to find the decrease in accuracy, precision, and recall for the speaker. On the other hand, if the speaker is  
 24 not present in the training set, then an ASV, i.e. *open-set* authentication, is used and the gain in privacy is  
 25 proportional to the increase in the confusion of the system to identify the original speaker. The open-set  
 26 evaluation is much closer to the realistic scenario where an attacker might obtain a small amount of speech  
 27 data and use it to enroll the speaker under attack. Therefore, ASV systems are often used in this thesis to  
 28 simulate the attackers.

29 The most widely used ASV metric is the *equal error rate* (EER): it considers an attacker that makes a  
 30 decision by comparing speaker similarity scores with a threshold and it assigns the same cost to false alarms  
 31 and misses [132]. The *application-independent log-likelihood-ratio cost function*  $C_{llr}^{\min}$  generalizes the EER by  
 32 considering optimal thresholds over all possible priors and all possible error costs [38]. In the following, we  
 33 consider a third metric called *linkability* which has recently emerged from the biometric template protection  
 34 community but has received little attention in the speech community so far [102]. This metric, denoted as  
 35  $D_{\text{sys}}^{\text{sys}}$ , estimates the distributions of scores for *mated* (same-speaker) vs. *non-mated* (different-speaker) trials  
 36 and computes their overlap. These metrics are formally defined below.

37 **Equal Error Rate (EER)** The EER is the classical metric used in speaker recognition. It assumes a  
 38 threshold-based decision on the score. If  $LF(a, b)$  is greater than a certain threshold  $t$ , the two utterances  $a$   
 39 and  $b$  are considered to be mated. Two types of errors can be made: false alarms with rate  $P_{\text{fa}}(t)$ , and misses  
 40 with rate  $P_{\text{miss}}(t)$ . The EER is the error rate corresponding to the threshold  $\tau$  for which the two types of  
 41 errors are equally likely:

$$42 \quad \text{EER} = P_{\text{miss}}(\tau) = P_{\text{fa}}(\tau). \quad (3.1)$$

**Log-Likelihood-Ratio Cost Function  $C_{\text{llr}}$  and  $C_{\text{llr}}^{\min}$**   $C_{\text{llr}}$  is also a common speaker recognition metric [38]. It is *application-independent* in the sense that it pools across all possible costs for false alarm vs. miss errors, and all possible priors for mated vs. non-mated trials, i.e., a pair of utterances from same vs. different speakers, respectively. Let  $M$  (resp.,  $\bar{M}$ ) be the set of mated (resp., non-mated) trials and  $|M|$  (resp.,  $|\bar{M}|$ ) its cardinality. Denoting by  $\text{llr}(p)$  be the log-likelihood ratio for trial  $p = (a, b)$ ,  $C_{\text{llr}}$  is defined as

$$C_{\text{llr}} = \frac{1}{\log 2} \left[ \frac{1}{|M|} \sum_{p \in M} \log \left( 1 + e^{-\text{llr}(p)} \right) + \frac{1}{|\bar{M}|} \sum_{p \in \bar{M}} \log \left( 1 + e^{\text{llr}(p)} \right) \right]. \quad (3.2)$$

$C_{\text{llr}}$  assesses the overall detection which includes both discrimination and calibration. In practice, discrimination alone is more relevant as a privacy metric. To measure it, a derived metric called  $C_{\text{llr}}^{\min}$  can be computed by optimal calibration of the scores  $LF(p)$  into log-likelihood ratios using a monotonic rising transformation. This transformation is found via the Pool Adjacent Violators algorithm (PAV), see [296] for details.

**Linkability** A linkability metric was proposed in [102] for biometric template protection systems. This metric can be generalized for any two sets of items. Denoting by  $H$  (resp.,  $\bar{H}$ ) the binary variable expressing whether two random utterances  $a$  and  $b$  are mated (resp., non-mated), the local linkability metric for a score  $s = LF(a, b)$  is defined as  $p(H | s) - p(\bar{H} | s)$ . When the local metric is negative, an attacker can deduce with some confidence that the two utterances are from different speakers. The authors of [102] argued that the local metric should estimate the strength of the link described by a score rather than measure how much a score describes non-mated relationships. Therefore they propose a clipped version of the difference:

$$D_{\leftrightarrow}(s) = \max(0, p(H | s) - p(\bar{H} | s)). \quad (3.3)$$

The global linkability metric  $D_{\leftrightarrow}^{\text{sys}}$  is the mean value of  $D_{\leftrightarrow}(s)$  over all mated scores:

$$D_{\leftrightarrow}^{\text{sys}} = \int p(s | H) \cdot D_{\leftrightarrow}(s) ds.$$

In practice,  $D_{\leftrightarrow}(s)$  is rewritten as  $(2 \cdot \omega \cdot \text{lr}(s)) / (1 + \omega \cdot \text{lr}(s)) - 1$  where the likelihood ratio  $\text{lr}(s)$  is  $p(s | H) / p(s | \bar{H})$  and the prior probability ratio  $\omega$  is  $p(H) / p(\bar{H})$ , and  $p(s | H)$  and  $p(s | \bar{H})$  are computed via one-dimensional histograms.

### 3.4.2 Utility measures

Generally, it is assumed that reasonably intelligible, natural-sounding, and good quality audio is sufficient for any kind of utility. Humans are employed to listen and rate these attributes on a perceptual scale similar to the mean opinion score evaluation of TTS systems, but this setup is very costly and time-consuming. Instead, most of the speaker anonymization studies including ours use an ASR as the objective judge of these attributes. It is a good proxy for measuring the utility since the given attributes highly correlate with the performance of ASR. In case the output of anonymization is not an intelligible speech signal, it is tricky to measure the gain or loss in privacy as well as utility, but in practice, similar systems can be deployed with different input features. ASR performance is measured in terms of word error rate (WER) which is derived from the Levenshtein distance. It is computed by first aligning the predicted word sequence with the reference word sequence, and then counting all the number of substitutions ( $S$ ), deletion ( $D$ ), and insertions ( $I$ ) required to convert the predictions to the references, where each reference utterance  $i$  is of length  $L^{(W_i)}$ . The WER for the whole data set is given by the following formula:

$$1 \quad \text{WER} = \frac{S + D + I}{\sum_i L(W_i)}. \quad (3.4)$$

### 2 3.4.3 Comparison of privacy metrics

3 In the later part of this chapter, experiments are performed to assess the suitability of the three metrics: EER,  
 4  $C_{llr}^{\min}$  and  $D_{\leftrightarrow}^{\text{sys}}$ , for the evaluation of speaker anonymization. In addition to comparing the metrics in their  
 5 form and substance, simulated data is generated to exhibit their blindspots. Experiments are also conducted  
 6 on real speech data processed by anonymization techniques proposed in Section 3.3.1 against three different  
 7 attackers defined in Section 3.3.2. Overall, the aim is to understand the complementary factors underlying  
 8 different metrics and ensure that the anonymization techniques being evaluated were not designed to fool  
 9 attackers that follow one specific speaker verification method but would fail with others.

10 Based on the definitions in Section 3.4.1, it is evident that the three metrics do not provide the same  
 11 information. Both the EER and  $C_{llr}^{\min}$  measure the probability of error of an attacker that makes decisions  
 12 based on a threshold on the linkage function (one particular threshold for EER and all possible ones for  
 13  $C_{llr}^{\min}$ ). Linkability measures something different: it evaluates how different the distributions of mated vs.  
 14 non-mated scores are. There is no attacker making a decision and there is no threshold or, from another  
 15 perspective, the best possible *oracle* attacker (not necessarily threshold-based) is assumed. In addition, if we  
 16 consider how general are the metrics, on the one hand  $C_{llr}^{\min}$  is a direct extension of the EER as it does not  
 17 focus on one single threshold. On the other hand,  $D_{\leftrightarrow}^{\text{sys}}$  is evaluated over all the encountered mated scores.  
 18 In Section 3.7, we provide experimental examples that highlight the differences of information provided and  
 19 generality of the metrics.

## 20 3.5 Experimental setup

21 In this section, we describe in detail the data sets, the parameters of the VC methods, and the models used to  
 22 evaluate privacy and utility.

### 23 3.5.1 Data and evaluation setup

Table 3.1 The subsets of the LibriSpeech data set along with their total duration in hours, duration per speaker in minutes, and number of male and female speakers.

	<b>Subset</b>	<b>Duration (h)</b>	<b>per speaker (min)</b>	<b>Male speakers</b>	<b>Female speakers</b>
Evaluation	dev-clean	5.4	8	20	20
	test-clean	5.4	8	20	20
	dev-other	5.3	10	17	16
	test-other	5.1	10	16	17
Training	train-clean-100	100.6	25	126	125
	train-clean-360	363.6	25	482	439
	train-other-500	496.7	30	602	564

24 The majority of experiments in this thesis are performed on the openly available LibriSpeech corpus [215]  
 25 which is a 1000 hours data set containing read English speech derived from a large collection of audiobooks  
 26 in the public domain.<sup>3</sup> The audiobooks are recorded by volunteers in rather clean ambient conditions using

<sup>3</sup><https://librivox.org/>

### 3.5 Experimental setup

55

their own microphone device, hence the recording conditions may differ from speaker to speaker. The audio is sampled at 16kHz and sufficient linguistic resources, such as the lexicon and the language models, are made available for download which makes it suitable for training ASR models. The corpus is gender-balanced in terms of the number of speakers and their individual duration as shown in Table 3.1, therefore it can be well suited for training speaker identification models. The whole data set is divided into evaluation and training subsets with further segregation of utterances incurring lower WER, designated as “clean”, from the ones incurring higher WER, i.e., “other”. Careful selection has been done to maintain the duration of individual utterances to about 10 seconds and to balance the data such that the per-speaker duration within the subset is almost the same.

For the experiments in this chapter, the 460-hour clean training set (*train-clean-100 + train-clean-360*), which contains 1,172 speakers, is used to train the disentanglement transform. Out of the *test-clean* set, an *enrollment* set (438 utterances) and a *trial* set (1,496 utterances) is created with different utterances from the same 29 speakers (13 male and 16 female, not in the training set) considered as source speakers. The details of the trial set are shown in Table 3.2. The target speakers for all three VC methods are randomly picked from the training and *test-clean* sets. As per the threat model considered in Figure 3.1, the trial set is used as the publicly released anonymized data set, the training set is the public data set used for training the evaluation models, i.e., ASR<sub>eval</sub> and ASV<sub>eval</sub> models, and the enrollment set is the found data used by the attacker. The exact method in which the trial, the enrollment and the training set are used to design the attackers is presented in Section 3.5.3.

Table 3.2 Detailed description of the trial set for speaker verification experiments.

	Male	Female
<b>No. of Speakers</b>	13	16
<b>No. of Genuine trials</b>	449	548
<b>No. of Impostor trials</b>	9,457	11,196

For each VC method and target selection strategy, all utterances in the trial set are mapped to possibly different target speakers in the training or trial set. The converted trial set serves as the public database that attackers want to de-anonymize by designing a linkage attack. To this end, attackers have access to the enrollment set which serves as the found data used to model the speakers in the trial set.

The attackers also have access to the 460-hour training set to train state-of-the-art speaker verification methods based on x-vectors [262] and i-vectors [61], which are stronger than the Gaussian mixture model-universal background model (GMM-UBM) based method used in the seminal work of [135]. We adapt the *sre16* Kaldi recipe for training x-vectors and i-vectors to LibriSpeech<sup>4</sup>. The original network architecture is presented in Table 2.1. The recipe is customized to use a smaller network architecture for x-vector computation than the original recipe. Specifically, compared to the given architecture, the *frame4*, *frame5* and *segment7* layers are removed, thereby also reducing the *stats pooling* layer to  $512T \times 1024$  and the *segment6* layer to  $1024 \times 512$ . Here  $T$  refers to the utterance-level context. This reduced architecture performs slightly better on LibriSpeech than the architecture in the original recipe. More details on the different attackers is given below in Section 3.5.3.

Finally, the utility of each VC method is evaluated in terms of the resulting ASR<sub>eval</sub><sup>anon</sup> performance, that is trained and tested on the converted data (Figure 3.4, 3), and this WER is compared with the baseline WER obtained on the clean speech using ASR<sub>eval</sub> (Figure 3.4, 1). A hybrid connectionist temporal classification (CTC) and attention based encoder-decoder [310] is used to build ASR<sub>eval</sub>, that is trained on the converted

<sup>4</sup>[https://github.com/brijmohan/kaldi/tree/master/egs/librispeech\\_spkv/v2](https://github.com/brijmohan/kaldi/tree/master/egs/librispeech_spkv/v2)

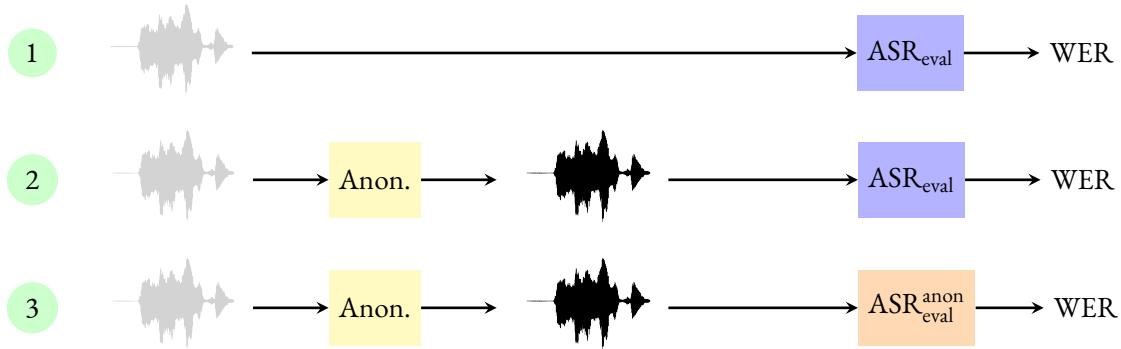


Fig. 3.4 Utility evaluation of (1) original speech data and (2) anonymized speech data using the ASR<sub>eval</sub> model that is trained over original speech. Case (3) indicates that the ASR model ASR<sup>anon</sup><sub>eval</sub> used for decoding is re-trained over anonymized speech data. The yellow block indicates the application of the anonymization algorithm.

1 460-hour training set using the standard recipe for LibriSpeech provided in ESPnet<sup>5</sup>. This model is also  
2 used for the experiments in the next chapter.

### 3 3.5.2 Voice conversion settings

4 **VoiceMask.** Pitch, aperiodicity and spectral envelope are extracted using the pyworld vocoder<sup>6</sup>. Only  
5 *random* strategy is followed for VoiceMask. The value of  $\alpha$  is uniformly sampled such that  $|\alpha| \in [0.08, 0.10]$   
6 then  $\beta$  in  $[-2, 2]$  such that  $0.32 \leq dist_{f_{\alpha,\beta}} \leq 0.40$  where  $dist_{f_{\alpha,\beta}} = \int_0^{\pi} |f_{\alpha,\beta}(\omega) - \omega|$  is the distortion  
7 strength of the warping function. Although our implementation slightly differs from the original method  
8 as described in Section 3.2, we use the ranges for  $\alpha$  and  $\beta$  provided by VoiceMask's authors in [226] since  
9 they produce most intelligible output. A subset of 100 target speakers is randomly selected and, for every  
10 utterance, pitch is transformed so as to match a random speaker within that subset. Other target selections  
11 strategies have not been applied because fixed values for  $\alpha$  and  $\beta$  (whether speaker-dependent or not) are  
12 prone to inversion attacks.

13 **VTLN-based VC.** Pitch, aperiodicity and spectral envelope are extracted using the pyworld vocoder. For  
14 each speaker, speech frames are selected using energy-based voice activity detection (VAD) with a threshold  
15 of 0.06, and their spectral envelopes are clustered via k-means with  $k = 8$ . In strategy *const*, only one target  
16 speaker is selected. In *perm*, a random subset of 100 target speakers is drawn and, for each source speaker, a  
17 random target is selected within the subset. In *random*, a random subset of 100 target speakers is drawn and,  
18 for each source utterance, a random target within the subset is selected.

19 **Disentangled representation based VC.** A publicly available implementation of this method is used.<sup>7</sup>  
20 As per the authors' suggestion in the preprocessing script, the disentanglement models (speaker encoder,  
21 content encoder, decoder) is trained over the *train-clean-100* subset of the LibriTTS corpus (itself a subset of  
22 the 460-hour training set of LibriSpeech), with a batch size of 128 and learning rate of 0.0005 for 500,000  
23 iterations. All three target selection strategies are applied similarly to VTLN-based VC except that only the

<sup>5</sup><https://espnet.github.io/espnet/>

<sup>6</sup><https://github.com/JeremyCCHsu/Python-Wrapper-for-World-Vocoder>

<sup>7</sup>[https://github.com/jjerry2243542/adaptive\\_voice\\_conversion](https://github.com/jjerry2243542/adaptive_voice_conversion)

source utterance and one random utterance from the target speaker are used as inputs to the content and speaker encoders, respectively. Other utterances from the source and targets speakers are unused.

### 3.5.3 Design of attackers using ASV

Several attackers have been implemented depending on the choice of the VC algorithm and the target selection strategy as well as the extent of the attacker's knowledge (*Informed*, *Semi-Informed*, or *Ignorant*). This implementation is illustrated in Figure 3.5, where the first row indicates the baseline case that uses the untransformed trial and enrollment sets along with the  $\text{ASV}_{\text{eval}}$  model trained on the original data. Our *Ignorant* attacker is unaware of the VC step: he/she simply uses the x-vector/i-vector  $\text{ASV}_{\text{eval}}$  model trained on the untransformed training set and applies them to the untransformed enrollment set, while the trial set is anonymized. Our *Semi-Informed* attacker knows the VC algorithm and the target selection strategy (*const*, *random* or *perm*) but not the particular choices of targets. He/she applies this strategy to the training and enrollment sets by drawing random target speakers from the subset of 100 target speakers used by the VC method (we assume that the value of  $k$  in VTLN is known to the attacker). As a result, the training and enrollment data are converted in a similar way as the trial data, but the target speaker associated with every speaker in the enrollment set is typically different from that which is associated with the same speaker in the converted trial set. The training set is then used to train a new  $\text{ASV}_{\text{eval}}^{\text{anon}}$  model to match the testing conditions. Finally, our *Informed* attacker has access to the actual VC models and target choices used to anonymize the trial set, so it converts the training and enrollment sets accordingly.

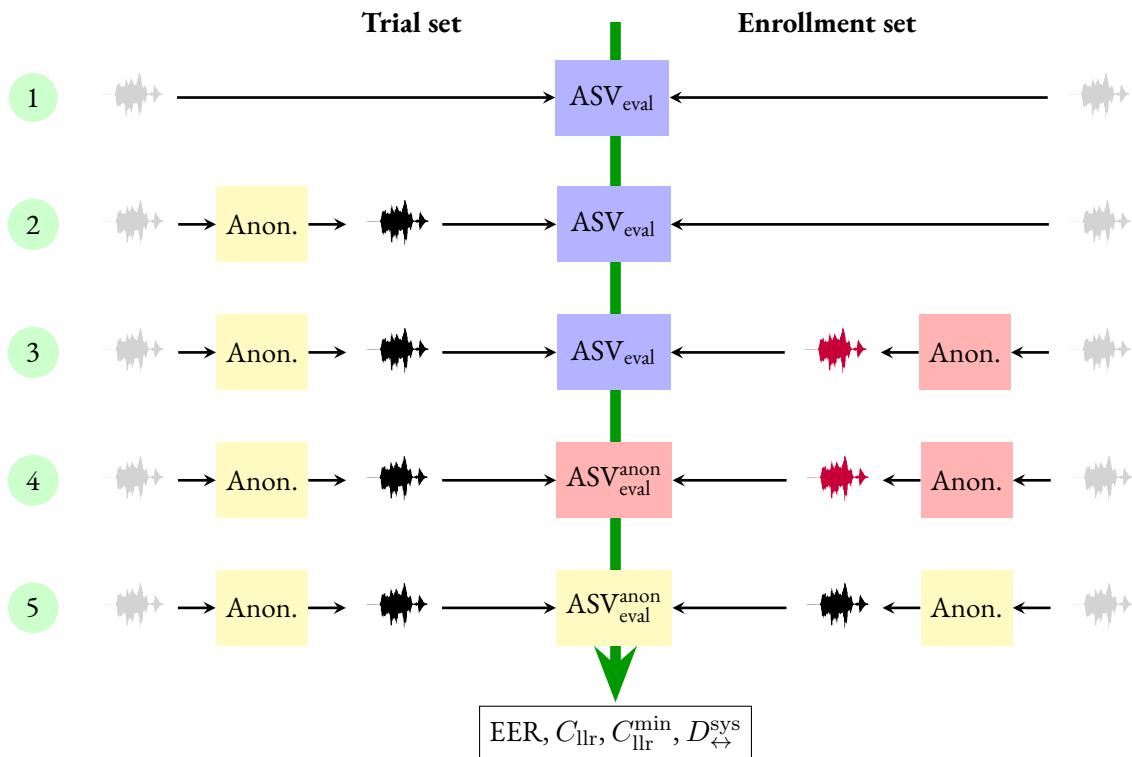


Fig. 3.5 Privacy evaluation in (1) Original, (2) *Ignorant*, (3) *Lazy-Informed*, (4) *Semi-Informed* and (5) *Informed* settings using the  $\text{ASV}_{\text{eval}}$  models. The blue  $\text{ASV}_{\text{eval}}$  block indicates the model trained using original speech, the red  $\text{ASV}_{\text{eval}}^{\text{anon}}$  block is trained using anonymized speech but with different parameters than the trial set, while the yellow  $\text{ASV}_{\text{eval}}^{\text{anon}}$  block is trained using the exact same parameters as the trial set.

In some of the preliminary experiments, attackers who convert the enrollment set only and use ASV<sub>eval</sub> models trained on the untransformed training set were also considered. Unsurprisingly, it was found that this leads to significantly larger *equal error rates* (EER) than re-training the ASV<sub>eval</sub><sup>anon</sup> model (which can easily be done by the attacker using public benchmark data). Although the results for such attackers are not reported in this chapter, later chapters refer to this attacker as *Lazy-Informed* and sometimes also mention the privacy protection obtained in such scenarios.

In the next section, the experiments to assess the strength of different attackers against the privacy protection strategies are presented. In the following, the three privacy metrics are compared in simulation and real data settings to determine their capability to express useful information about the privacy protection mechanism. Results of these experiments are reported and their implications are discussed in detail.

### 3.6 Experimental comparison with different attackers

The three different VC-based anonymization strategies described in Section 3.3.1 are evaluated against the three attackers defined in Section 3.5.3. First and foremost, the ASR<sub>eval</sub> and ASV<sub>eval</sub> systems are trained and applied to the original (untransformed) data for baseline performance. An EER of 4.61% and 4.31% are obtained over the trial set described in Table 3.2 for i-vector and x-vector, respectively, and a WER of 9.4% for ASR<sub>eval</sub> over the test-clean subset of LibriSpeech.

Tables 3.3 and 3.4 present the EER for x-vector and i-vector based speaker verification for the three attackers and the various VC methods and target selection strategies. Interestingly, the *Informed* attacker achieves similar or even slightly lower EER than the baseline in most cases. This indicates that, when the attacker has complete knowledge of the VC scheme and target speaker mapping, none of the VC methods can protect the speaker identity. While an attacker with such complete knowledge is not very realistic in most practical cases, our results show that speaker information has not been totally removed and is somehow still present in the converted speech. They also indicate that privacy protection only relies on the randomization introduced by the target selection strategies.

Table 3.3 EER (%) achieved using x-vector based ASV<sub>eval</sub> for *Ignorant* attacker, and ASV<sub>eval</sub><sup>anon</sup> for *Semi-Informed* and *Informed* attackers. Bold face indicates the best privacy protection strategy against *Informed* attackers.

Attackers ↓ / Strategies →	random	VTLN-based VC			Disentangl.-based VC		
		const	perm	random	const	perm	random
<i>Informed</i>	5.01	4.71	3.91	<b>6.32</b>	4.71	0.20	<b>5.52</b>
<i>Semi-Informed</i>	-	12.84	23.37	6.32	13.64	43.03	5.42
<i>Ignorant</i>	28.69	24.27	30.99	27.38	27.68	32.20	30.59

For the more realistic *Semi-Informed* attacker, it is observed that strategy *perm* is quite effective in protecting privacy and shows the highest gains in EER. This is because the target speaker in the enrolled data may not be the same as the one in the trial, hence greater confusion is induced during inference. It is also important to note that strategy *random* is not much affected by the change of speaker mapping, which is intuitive because in this case the utterances are already being mapped randomly to different speakers. Such mapping would be ineffective due to averaging of randomness. Strategy *const* is also slightly affected by the change of mapping because the training and enrollment speaker is not the same as that of the test speaker, but the effect is not as significant as strategy *perm*. A preliminary *Lazy-Informed* experiment is also

## 3.6 Experimental comparison with different attackers

59

Table 3.4 EER (%) achieved using i-vector based ASV<sub>eval</sub> for *Ignorant* attacker, and ASV<sub>eval</sub><sup>anon</sup> for *Semi-Informed* and *Informed* attackers. Bold face indicates the best privacy protection strategy against *Informed* attackers.

<b>Attackers ↓ / Strategies →</b>	<b>VoiceMask</b>	<b>VTLN-based VC</b>			<b>Disentangl.-based VC</b>		
	<i>random</i>	<i>const</i>	<i>perm</i>	<i>random</i>	<i>const</i>	<i>perm</i>	<i>random</i>
<i>Informed</i>	8.22	6.22	<b>10.23</b>	9.84	4.71	0.20	<b>11.03</b>
<i>Semi-Informed</i>	-	18.25	31.49	18.76	15.65	43.93	10.53
<i>Ignorant</i>	50.55	26.08	49.15	49.15	49.95	47.74	49.85

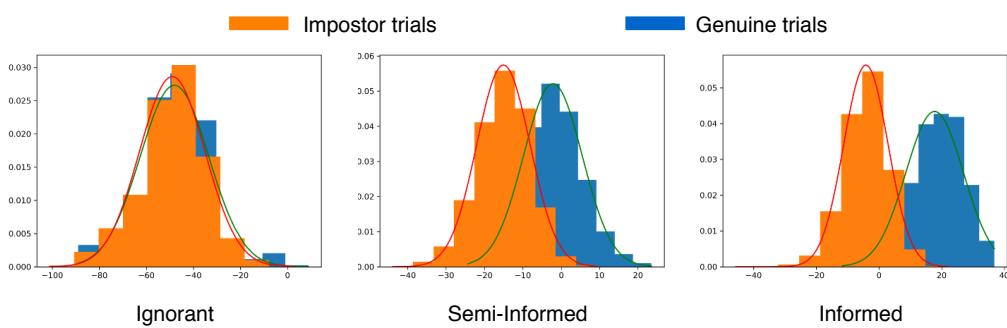


Fig. 3.6 I-vector score distribution for trials conducted on VTLN (strategy *random*) converted data by *Ignorant*, *Semi-Informed*, or *Informed* attackers. The orange distribution indicates impostor scores, while the blue distribution indicates genuine scores. The crossing between the two curves indicates the threshold for EER. More overlap means greater confusion, hence greater privacy protection.

performed with the VoiceMask technique and evaluated using the baseline x-vector ASV<sub>eval</sub> model. The obtained EER was 20.96 which lies in-between the *Informed* (5.01) and the *Ignorant* (28.69) attacker. This scenario is explored in more detail in Chapter 5.

Consistently with past results in the literature, the *Ignorant* attacker performs worst in terms of EER. This confirms that, when the attacker is oblivious to the privacy-preserving mechanism, we can protect speaker identity completely. Figure 3.6 shows the distribution of i-vector PLDA scores for genuine and impostor trials, i.e., the log-likelihood ratios between *same-speaker* and *different-speaker* hypotheses. For full unlinkability, the distributions of genuine and impostor scores must be identical. We observe that the overlap between the two distributions decreases as we move from the *Ignorant* to the *Informed* attacker, hence increasing linkability.

Additional experiments were conducted to investigate that the attacker may not simply use the same protection strategy as the speaker. As observed in Tables 3.3 and 3.4, the *Semi-Informed* attacker sometimes perform slightly better than *Informed*, hence a comparative study was conducted where a *Semi-Informed* attacker, trained using a particular strategy, may try to deduce the speaker's identity from speech samples protected using all the other strategies. Also looking at the *Informed* and *Semi-Informed* attackers trained using *random* strategy performing consistently well, a natural question arises: how well does the *random* attacker perform against *const* and *perm* strategies?

Results, as reported in Table 3.5, indicate that *perm* strategy, as observed before, may not be the best for a speaker because it is not resilient against an attacker trained using *random* strategy. Such an attacker performs decently well against any strategy. It may be the case because having observed several different

Table 3.5 Additional i-vector results with VTLN *Semi-Informed* attackers against the protection strategies. Bold face indicates the best performing attacker’s strategy against a protection strategy.

<b>Semi-Informed attacker (train + enroll)</b>	<b>Strategy (trial)</b>		
	<i>const</i>	<i>perm</i>	<i>random</i>
<i>const</i>	<b>18.25</b>	26.18	25.18
<i>perm</i>	33.00	31.49	33.60
<i>random</i>	20.66	<b>17.35</b>	<b>18.76</b>

- 1 targets, the system has learned to distinguish speaker’s information by ignoring the target selection and  
 2 considering other discriminatory features that are not removed by switching the target identity. Hence it  
 3 is crucial to discover and remove those factors that may be contributing towards the speaker’s identity for  
 4 complete anonymization.

Table 3.6 WER (%) achieved using end-to-end ASR<sub>eval</sub><sup>anon</sup>.

<b>Subset ↓ / Strategies →</b>	<b>Baseline</b>	<b>VoiceMask</b>	<b>VTLN-based VC</b>			<b>Disentangl.-based VC</b>		
			<i>random</i>	<i>const</i>	<i>perm</i>	<i>random</i>	<i>const</i>	<i>perm</i>
dev-clean	9.2	17.7	19.9	17.9	15.5	46.9	23.3	112.9
test-clean	9.4	18.1	19.8	18.4	15.9	41.5	23.7	115.1
dev-other	28.1	37.4	41.2	37.5	34.0	73.9	45.3	113.9
test-other	29.7	39.0	41.4	38.5	35.0	76.6	47.1	111.7

- 5 Table 3.6 gives the WER obtained for each VC method, which is used as a proxy for the usefulness of  
 6 the converted speech. Note that there is no difference between converted data in different attack scenarios,  
 7 hence the WER does not depend on the attacker. VoiceMask and VTLN-based VC achieve reasonable  
 8 WER compared to the untransformed data, while the disentangled representation based VC produces  
 9 unreasonably high WER. Note that these WERs are achieved when ASR is trained solely using converted  
 10 data. In practice, many techniques can be used to optimize the WER, such as using converted data to  
 11 augment clean data.

## 12 3.7 Experimental comparison of privacy metrics

- 13 In this section, first, the privacy metrics are compared in a simulated setting where a variety of data, rep-  
 14 resenting different possible linkage scores and data points, is artificially generated to assess their response  
 15 in different situations. Thereafter, the metrics are compared in a real data setting where the scores were  
 16 obtained using the experimental setting similar to Section 3.6.

### 17 3.7.1 Exhibiting differences and blindspots through simulation

- 18 Two experiments are designed over simulated scores in order to exhibit the differences between the metrics.  
 19 The first experiment relies on discrete scores to highlight the lack of generality of the EER. The second  
 20 experiment relies on Gaussian distributed scores to exhibit the differences between  $C_{llr}^{\min}$  and linkability. All

## 3.7 Experimental comparison of privacy metrics

61

of the metrics are integrated in the Voice Privacy Challenge 2020<sup>8</sup> and an easy-to-use open-source toolkit<sup>9</sup> was developed by Mohamed Maouche.

**Discrete Scores** Let us assume that there are 8 trials, i.e., pairs of utterances  $p_1, \dots, p_8$  and that the score for the  $i$ -th trial is given by the integer  $LF(p_i) = i$  as shown in the header of Table 3.7. The values of EER and  $C_{llr}^{\min}$  vary with the label (mated vs. non-mated) of each trial. In Table 3.7, three particular cases are shown where only the labels of the last three trials (associated with scores 6, 7, and 8) change. It is observed that this has an effect on  $C_{llr}^{\min}$  but not on the EER. This is because the EER searches for a single threshold of the linkage function while  $C_{llr}^{\min}$  averages over all possible thresholds that the attacker might choose. Further, it is also observed that the EER indicates a privacy of 0.25 that is half of the best achievable privacy (0.5), while  $C_{llr}^{\min}$  increases from half of the best achievable privacy (0.5 over 1) to higher values (0.65).

Table 3.7  $C_{llr}^{\min}$  and *EER* with discrete scores in  $\{1, \dots, 8\}$ .  $H$  (resp.  $\bar{H}$ ) denote mated (resp. non-mated) scores.

Score	1	2	3	4	5	6	7	8	$C_{llr}^{\min}$	EER
<b>Case 1</b>	$\bar{H}$	$\bar{H}$	$H$	$\bar{H}$	$H$	$\bar{H}$	$H$	$H$	0.50	0.25
<b>Case 2</b>	$\bar{H}$	$\bar{H}$	$H$	$\bar{H}$	$H$	$H$	$\bar{H}$	$H$	0.59	0.25
<b>Case 3</b>	$\bar{H}$	$\bar{H}$	$H$	$\bar{H}$	$H$	$H$	$H$	$\bar{H}$	0.65	0.25

**Gaussian Scores** Since  $D_{\leftrightarrow}^{\text{sys}}$  relies on density estimation, Gaussian distributed scores are generated to compare  $D_{\leftrightarrow}^{\text{sys}}$  and  $C_{llr}^{\min}$ . Three different Gaussians are considered here:  $G_1 \sim \mathcal{N}(1, \sigma_1)$ ,  $G_2 \sim \mathcal{N}(2, \sigma_2)$  and  $G_3 \sim \mathcal{N}(3, \sigma_3)$ . Each Gaussian  $G_i$  is used to sample either mated or non-mated scores according to a key  $k_i \in \{H, \bar{H}\}$ . In total, four different cases are considered depending on the values of  $(k_1, k_2, k_3)$ : *Mated higher* for  $(\bar{H}, \bar{H}, H)$  or  $(\bar{H}, H, H)$ ; *Mated lower* for  $(H, \bar{H}, \bar{H})$  or  $(H, H, \bar{H})$ ; *Mated in-between* for  $(\bar{H}, H, \bar{H})$ ; *Non-mated in-between* for  $(H, \bar{H}, H)$ . The three given distributions are sampled in order to obtain 5,000 mated and 5,000 non-mated scores. Multiple standard deviations are chosen to obtain different degrees of overlap between the distributions:  $(\sigma_1, \sigma_2, \sigma_3) \in \{0.1, 0.5, 1, 1.5\}^3$ .

The results are presented in Fig. 3.7.  $C_{llr}^{\min}$  and  $D_{\leftrightarrow}^{\text{sys}}$  are considered equivalent when  $C_{llr}^{\min}$  is equal to  $1 - D_{\leftrightarrow}^{\text{sys}}$  (diagonal line). The two metrics agree to a large extent only when the mated scores are higher. When the non-mated scores are higher (mated lower),  $C_{llr}^{\min}$  is always close to 1 while  $D_{\leftrightarrow}^{\text{sys}}$  varies depending on the overlap between the distributions. In the two remaining cases when the mated scores are surrounded by the non-mated scores or vice-versa,  $C_{llr}^{\min}$  is lower-bounded by 0.6 and the two metrics do not agree on the strength of anonymization. This is explained by the fact that threshold-based decision is meaningful in the *mated higher* case and its performance is then strongly related to the overlap between distributions, while it fails partially or totally in the three other cases.

To illustrate why this is an issue and how this may happen in practice, in Figure 3.8, (simulated) x-vectors are drawn for multiple utterances of two speakers, which have all been anonymized by mapping them to another (target) speaker's voice. Each utterance of speaker A has been randomly mapped to the left or the right cluster, while the utterances of speaker B have been mapped to the center cluster. The resulting score distributions match the *non-mated in-between* case above. As expected, the two metrics strongly disagree:  $D_{\leftrightarrow}^{\text{sys}} = 0.99$  (low privacy) and  $C_{llr}^{\min} = 0.81$  (high privacy). While this situation is unlikely to occur with unprocessed data (scores are then expected to match the *mated higher* case), it becomes likely once the

<sup>8</sup><https://www.voiceprivacychallenge.org/>

<sup>9</sup>[https://gitlab.inria.fr/magnet/anonymization\\_metrics](https://gitlab.inria.fr/magnet/anonymization_metrics)

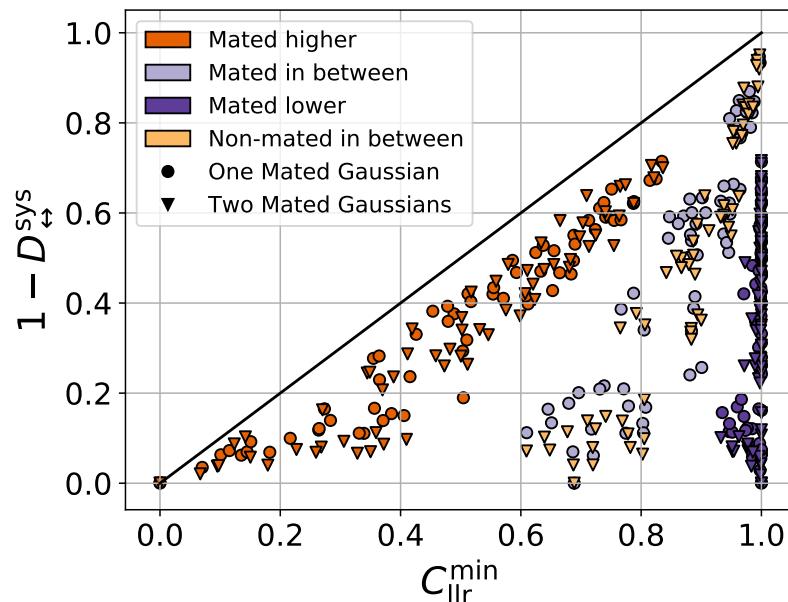


Fig. 3.7  $C_{llr}^{\min}$  vs.  $1 - D_{\leftrightarrow}^{sys}$  on simulated Gaussian scores.

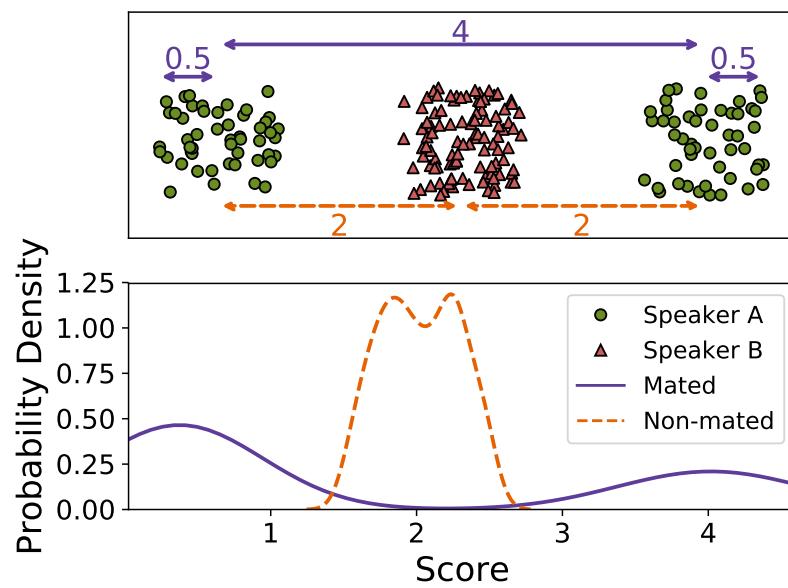


Fig. 3.8 Simulated ‘non-mated in-between’ data. Top: x-vectors visualized in 2D. Bottom: resulting score distributions.

utterances have been anonymized and the anonymization design choices (refer to Section 5.4 for example choices) result in multimodal score distributions.

### 3.7.2 Evaluation on real anonymized speech

In this section, the three privacy metrics are compared under real data settings. The scores and data points used for linkage attacks, along with the observations after the metric comparison are mentioned below.

**Scores under consideration** The linkage scores for this experiment are generated using the VC-based anonymization techniques investigated in Section 3.6. Three target selection strategies, *const*, *perm*, and *random* are used for VTLN and disentanglement approach, while only *random* is used for VoiceMask. For each of them, the three known attackers, namely *Ignorant*, *Semi-Informed*, and *Informed* are used to generate x-vectors. The attacker performs linkage attacks by computing the x-vectors of a trial utterance and an enrollment utterance and comparing them using one of three linkage functions: PLDA affinity, cosine distance, or Euclidean distance. This, along with the baseline, results in a total of 63 combinations of anonymization techniques, target selection strategies, attacker knowledge levels, and linkage functions.

**Results** Figures 3.9 and 3.10 compare the resulting metrics, where each dot corresponds to one of the 63 combinations above. The comparison between the EER and  $C_{llr}^{\min}$  (Fig. 3.9) shows a clear relation between the two metrics. In some cases the EER is stable and  $C_{llr}^{\min}$  varies a little bit but not significantly so. Regarding the comparison between  $D_{\leftrightarrow}^{\text{sys}}$  and  $C_{llr}^{\min}$ , a clear difference can be observed between Fig. 3.10 on real data and Fig. 3.7 on simulated Gaussian scores: on real data, the two metrics follow a clear relation.

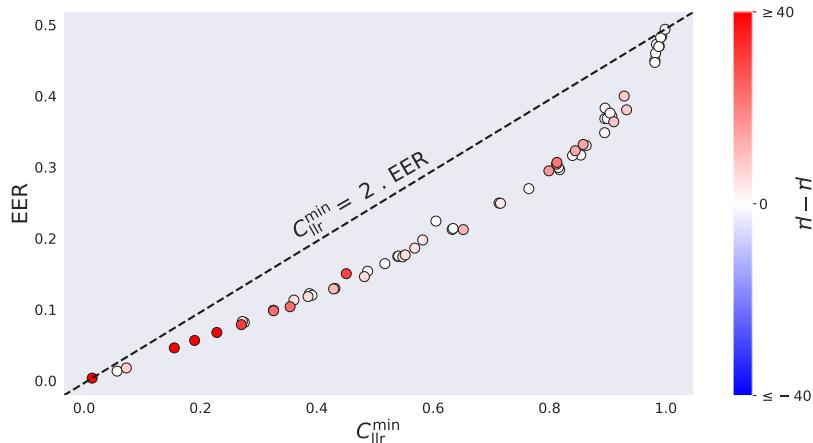


Fig. 3.9  $C_{llr}^{\min}$  vs.  $EER$  on real data. The color scale  $\mu - \bar{\mu}$  is the difference of the means of mated and non-mated scores.

These results can be explained by the fact that, with few exceptions, the score distributions for the specific target selection and attack strategies considered here fall into the *mated higher* case, as can be seen from the colors associated with the dots. It is however likely that advanced target selection strategies aiming for score distributions akin to Fig. 3.7 will be developed in the near future, as these would provide an advantage against attackers making threshold-based decisions. For that reason, this study provides evidence that  $D_{\leftrightarrow}^{\text{sys}}$  should be privileged as a privacy metric, since it provides very similar results to established metrics with current target selection and attack strategies, while being more robust to advanced strategies that will likely be developed soon.

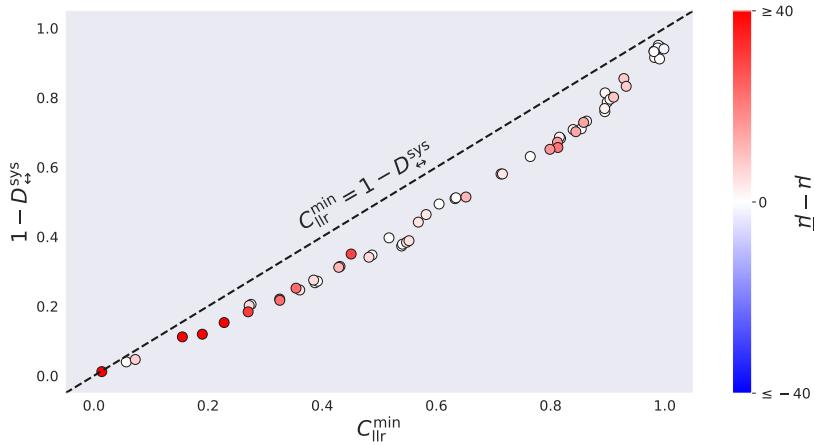


Fig. 3.10  $C_{llr}^{\min}$  vs.  $1 - D_{\leftrightarrow}^{\text{sys}}$  on real data. The color scale  $\mu - \bar{\mu}$  is the difference of the means of mated and non-mated scores.

### 3.8 Summary

In this chapter, We investigated the use of VC methods to protect the privacy of speakers by concealing their identity. Target speaker selection strategies and linkage attack scenarios based on the knowledge of attacker were formally defined. The experimental results indicated that both aspects play an important role in the strength of the protection. Simple methods such as VTLN-based VC with appropriate target selection strategy can provide reasonable protection against linkage attacks with partial knowledge.

The characterization of strategies and attack scenarios in this chapter opens up the avenues for designing better anonymization schemes in the upcoming chapters. Chapter 4 generalizes the idea of *Informed* attacker to measure the amount of speaker-identifiable attributes in the intermediate representations of an ASR network. To increase the naturalness of converted speech, intra-gender VC as well as the use of a supervised phonetic classifier in VTLN can be explored. Although the adversarial learning-based approach proposed in Chapter 4 produces anonymous feature vectors instead of a speech signal as output, these vectors can be used by speech synthesis-based methods to generate a speech signal as shown in Chapter 5. Chapter 5 explores a speech synthesis based approach for generating private speech signal, which has a high-quality and more natural output. Standard local and global unlinkability metrics [102] are used to precisely evaluate the privacy protection in various scenarios. More generally, designing a privacy-preserving transformation which induces a large overlap between genuine and impostor distributions even in the *Informed* attack scenario remains an open question which we will continue to address in the remaining chapters. In the case of disentangled representations, this calls for avoiding any leakage of private attributes into the content embeddings which can be achieved using the technique proposed in Chapter 4.

Furthermore, three metrics to assess the effectiveness of anonymization are compared: the EER, the application-independent log-likelihood-ratio min cost function  $C_{llr}^{\min}$ , and the linkability  $D_{\leftrightarrow}^{\text{sys}}$ . The EER and  $C_{llr}^{\min}$  assume that the attacker makes threshold-based decisions on the linkage score, while  $D_{\leftrightarrow}^{\text{sys}}$  implicitly models a more powerful, non-threshold-based *oracle* attacker. The comparison on real speech data processed via three anonymization techniques with different target selection strategies and with nine attackers suggests that these metrics behave similarly. Yet, experiments on simulated data highlight fundamental differences. Specifically, the EER may yield a fixed value for situations involving different levels of privacy correctly captured by  $C_{llr}^{\min}$ , and  $C_{llr}^{\min}$  becomes less informative than  $D_{\leftrightarrow}^{\text{sys}}$  when the mated scores are lower or interleaved with non-mated scores. While such situations were unlikely to occur in the field of speaker

3.8 Summary**65**

verification, which involves unprocessed speech data, it is expected for them to become frequent in the field of anonymization when more advanced target selection and attack strategies are built. For this reason, this study advocates for the use of  $D_{\leftrightarrow}^{\text{sys}}$  as a robust privacy metric capable of handling both current approaches and future developments in this field.

1  
2  
3  
4

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## Chapter 4

# Adversarial Learning based Anonymization

For self-realization, a rebel demands a strong authority, a worthy opponent, God to his Lucifer.

*Mary McCarthy*

This chapter investigates if there is personally identifiable information present in the bottleneck representation of an ASR network, and to what extent can it be used as features to re-identify the speakers within and outside the training set. Consider a slightly different threat model than the one described in Section 3.1, where individuals use the speech-to-text service provided by digital assistants [177, 148]. In this context, the speech signal is sent from the user device to a cloud-based service, as shown in Figure 4.1, where ASR and natural language understanding are performed in order to address the user request.<sup>1</sup>

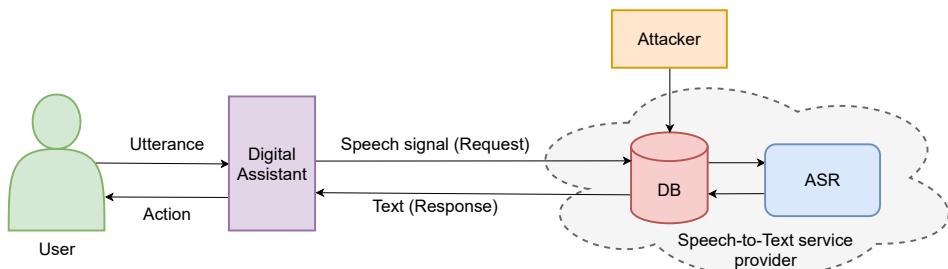


Fig. 4.1 Threat model related to speech-to-text provided by cloud-based services.

While recent studies have identified security vulnerabilities in these devices [166, 47], such studies tend to hide more important privacy risks that can have long-term impact. For instance, if the signal is intercepted by a malicious entity, a re-identification attack can be launched against the user, potentially compromising the person's identity [236], intention [107, 118, 23, 269], gender [326, 156], emotional state [81, 299, 158], pathological condition [67, 295, 249], personality [247, 248] and cultural [250, 301] attributes to a great extent using state-of-the-art speech technologies. These algorithms require just a few tens of hours of

<sup>1</sup>See e.g., <https://cloud.google.com/speech-to-text/>

1 training data to achieve reasonable accuracy, which is easier than ever to collect via virtual assistants. The  
2 dissemination of voice signals in large data centers thereby poses severe privacy threats to the users in the  
3 long run.

4 An alternative software architecture is to pre-process voice data on the device to remove some personal  
5 information before sending it to web services. Although this does not rule out all possible risks, a change  
6 of representation of the voice signal can contribute to limiting unsolicited uses of data. In this chapter, we  
7 investigate how much of a user’s *identity* is encoded in speech representations built for ASR. To this end,  
8 closed- and open-set speaker recognition experiments are conducted. The *closed-set* experiment refers to a  
9 classification setting where all test speakers are known at training time. In contrast, the *open-set* experiment  
10 (a.k.a. speaker verification) aims to measure the capability of an attacker to discriminate between speakers in  
11 a more realistic setting where the test speakers are not known beforehand. The attacker is implemented with  
12 the state-of-the-art x-vector speaker recognition technique [262] as mentioned in the previous chapter.

13 The representations of speech considered in this chapter are given by the encoder output of end-to-end  
14 deep encoder-decoder architectures trained for ASR. Such architectures are natural in our privacy-aware  
15 context, as they correspond to encoding speech on the user device and decoding in the cloud. The baseline  
16 network used here follows the ESPnet architecture [310], with one encoder and two decoders: one based on  
17 connectionist temporal classification (CTC) and the other on an attention mechanism, briefly mentioned  
18 at the end of Section 2.2.3. Inspired by [89], the methods in this chapter propose to extend the baseline  
19 network with a *speaker-adversarial* branch so as to learn representations that perform well in ASR while  
20 hiding the speaker identity.

21 Several papers have recently proposed to use adversarial training for the goal of improving ASR per-  
22 formance by making the learned representations invariant to various conditions. While general form of  
23 acoustic variabilities have been studied [251], there is some work specifically on speaker invariance [292, 195].  
24 Interestingly, there is no general consensus on whether it is more appropriate to use speaker classification in  
25 an adversarial or a multi-task manner, despite the fact that these two strategies implement opposite means  
26 (i.e., encouraging representations to be speaker-invariant or speaker-specific). This question was studied in  
27 [7], in which the authors conclude that both approaches only provide minor improvements in terms of  
28 ASR performance. Their speaker classification experiments also show that the baseline system already tends  
29 to learn speaker-invariant features. However, they did not run speaker verification experiments and hence  
30 did not assess the suitability of these features for the goal of anonymization.

31 In contrast to these studies which aim to increase ASR performance, the goal of this thesis is to assess  
32 the potential benefit of adversarial training for concealing speaker identity in the context of privacy-friendly  
33 ASR. This chapter describes the following contributions of this thesis. First, CTC, attention and adversarial  
34 learning are combined within an end-to-end ASR framework. Second, a rigorous protocol is designed to  
35 quantify speaker identity in ASR representations through a series of closed-set classification and open-set  
36 verification experiments. Third, as per the experiments on the LibriSpeech corpus [215], it is shown that this  
37 framework dramatically reduces speaker classification accuracy, but does not increase speaker verification  
38 error. Several possible reasons are suggested behind this disparity.

39 The rest of the chapter is structured as follows. In Section 4.1.2, the baseline ASR model and our  
40 proposed adversarial model is described. Section 4.2 explains the experimental setup and presents the  
41 obtained results. Finally, Section 4.4 concludes the chapter and briefly discusses future directions.

## 42 4.1 Proposed model

43 This section starts by describing the ASR model used as a baseline, before introducing the proposed speaker-  
44 adversarial network.

### 4.1.1 Baseline end-to-end ASR model

The end-to-end ASR framework presented in [311] is used as the baseline architecture which is also depicted in Figure 2.8. It is composed of three sub-networks: an *encoder* which transforms the input sequence of speech feature vectors into a new representation  $\mathbf{B}$ , and two *decoders* that predict the character sequence from  $\mathbf{B}$ . It is assumed that these networks have already been trained using data previously collected by the service provider (which may be public data, opt-in user data, etc). Then, in the deployment phase of the system that is envisioned in this chapter, the encoder would run on the user device and the resulting representation  $\mathbf{B}$  would be sent to the cloud for decoding.

The first decoder is based on CTC and the second on an attention mechanism. As argued in [311], attention works well in most cases because it does not assume conditional independence between the output labels (unlike CTC). However, it is so flexible that it allows nonsequential alignments which are undesirable in the case of ASR. Hence, CTC acts as a regularizer to prune such misaligned hypotheses. The parameters of the encoder are denoted by  $\theta_{\mathcal{E}}$ , and by  $\theta_{\text{ctc}}$  and  $\theta_{\text{att}}$  the parameters of the CTC and attention decoders respectively. The model is trained in an end-to-end fashion by minimizing an objective function  $\mathcal{L}_{\text{asr}}$  which is a combination of the losses  $\mathcal{L}_{\text{ctc}}$  and  $\mathcal{L}_{\text{att}}$  from both decoder branches:

$$\min_{\theta_{\mathcal{E}}, \theta_{\text{ctc}}, \theta_{\text{att}}} \mathcal{L}_{\text{asr}}(\theta_{\mathcal{E}}, \theta_{\text{ctc}}, \theta_{\text{att}}) = \beta \mathcal{L}_{\text{ctc}}(\theta_{\mathcal{E}}, \theta_{\text{ctc}}) + (1 - \beta) \mathcal{L}_{\text{att}}(\theta_{\mathcal{E}}, \theta_{\text{att}}),$$

with  $\beta \in [0, 1]$  a trade-off parameter between the two decoders.

The form of the two losses  $\mathcal{L}_{\text{ctc}}$  and  $\mathcal{L}_{\text{att}}$  is formally described in Equations (2.20) and (2.21), respectively. Let us briefly recall the notation to eventually describe the speaker-adversarial objective. Each sample in the dataset is denoted as  $\mathbf{S}_i = (\mathbf{O}_i, Y_i, z_i)$ , where  $\mathbf{O}_i = \{\mathbf{o}_1, \dots, \mathbf{o}_T\}$  is the sequence of  $T$  acoustic feature frames,  $Y_i = \{y_1, \dots, y_C\}$  is the sequence of  $C$  characters in the transcription, and  $z_i$  is the speaker label. In the case of CTC, several intermediate label sequences of length  $T$  are created by repeating characters and inserting a speacial *blank* label to mark character boundaries. Let  $\{\mathbf{a}_1^{(i)}, \dots, \mathbf{a}_M^{(i)}\}$  be the set of all such intermediate label sequences or alignments, and  $\mathbf{a}_j^{(i)} = [\bar{y}_1, \dots, \bar{y}_T]$ . The CTC loss  $\mathcal{L}_{\text{ctc}}(\theta_{\mathcal{E}}, \theta_{\text{ctc}})$  is computed as  $\mathcal{L}_{\text{ctc}} = -\log P(Y_i | \mathbf{O}_i; \theta_{\mathcal{E}}, \theta_{\text{ctc}})$  where  $P(Y_i | \mathbf{O}_i; \theta_{\mathcal{E}}, \theta_{\text{ctc}}) = \sum_{j=1}^M P(\mathbf{a}_j^{(i)} | \mathbf{O}_i; \theta_{\mathcal{E}}, \theta_{\text{ctc}})$ . This sum is computed by assuming conditional independence of observing a label  $\bar{y}_t$  over previously observed labels  $\bar{y}_{1:t-1}$ , hence  $P(\mathbf{a}_j^{(i)} | \mathbf{O}_i; \theta_{\mathcal{E}}, \theta_{\text{ctc}}) = \prod_{t=1}^T P(\bar{y}_t | \mathbf{O}_i; \theta_{\mathcal{E}}, \theta_{\text{ctc}})$ . The attention branch does not require an intermediate label representation and conditional independence is not assumed, hence the loss is simply computed as  $\mathcal{L}_{\text{att}}(\theta_{\mathcal{E}}, \theta_{\text{att}}) = -\sum_{c \in C} \ln P(y_c | \mathbf{O}_i, y_{1:c-1}; \theta_{\mathcal{E}}, \theta_{\text{att}})$ .

### 4.1.2 Speaker-adversarial model

In order to encourage the network to learn representations that are not only good at ASR but also hide speaker identity, we propose to extend the above architecture with what we call a *speaker-adversarial* branch. This branch models an adversary which attempts to infer the speaker identity from the encoded representation  $\mathbf{B}$ . We denote by  $\theta_{\text{spk}}$  the parameters of the speaker-adversarial branch. Given the encoder parameters  $\theta_{\mathcal{E}}$ , the goal of the adversary is to find  $\theta_{\text{spk}}$  that minimizes the loss  $\mathcal{L}_{\text{spk}}(\theta_{\mathcal{E}}, \theta_{\text{spk}}) = -\log P(z_i | \mathbf{O}_i; \theta_{\mathcal{E}}, \theta_{\text{spk}})$ . Our new model is then trained in an end-to-end manner by optimizing the following min-max objective:

$$\min_{\theta_{\mathcal{E}}, \theta_{\text{ctc}}, \theta_{\text{att}}} \max_{\theta_s} \mathcal{L}_{\text{asr}}(\theta_{\mathcal{E}}, \theta_{\text{ctc}}, \theta_{\text{att}}) - \lambda \mathcal{L}_{\text{spk}}(\theta_{\mathcal{E}}, \theta_{\text{spk}}), \quad (4.1)$$

where  $\lambda \geq 0$  is a trade-off parameter between the ASR objective and the speaker-adversarial objective. The baseline network can be recovered by setting  $\lambda = 0$ . Note that the max part of the objective corresponds

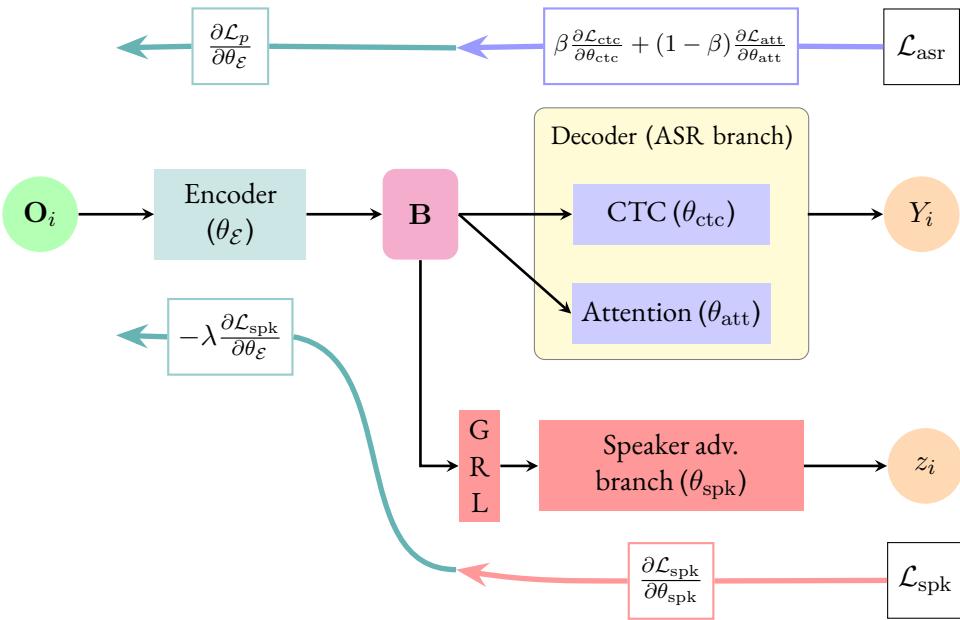


Fig. 4.2 Architecture of the proposed model. The speaker-adversarial branch is shown as a red box. The teal arrow going from GRL to encoder indicates *gradient reversal*. When the model is deployed, the encoder could reside at the client side, while the decoder can be hosted by cloud services.

1 to the adversary, which controls only the speaker-adversarial parameters  $\theta_{spk}$ . The goal of the speaker-  
 2 adversarial branch is to act as a “good adversary” and produce useful gradients to remove the speaker identity  
 3 information from the encoded representation  $B$ . In practice, we use a *gradient reversal layer* [96] between  
 4 the encoder and the speaker-adversarial branch so that the whole network can be trained end-to-end via  
 5 backpropagation. Refer to Figure 4.2, which is an adaptation of Figure 2.14 for an illustration of the full  
 6 architecture.

## 7 4.2 Experimental evaluation

### 8 4.2.1 Datasets

9 We use the Librispeech corpus, described in Table 4.1 for all the experiments in this chapter. Different  
 10 subsets are used for ASR training, adversarial training, and speaker verification. For the sake of clarity we  
 11 refer to them as *data-full*, *data-adv*, and *data-spkv*, respectively. The *data-full* set is almost the original  
 12 Librispeech corpus (Refer Table 3.1), including *train-960* for training, *dev-clean* and *dev-other* for validation,  
 13 and *test-clean* and *test-other* for test, except that utterances with more than 3,000 frames or more than 400  
 14 characters have been removed from *train-960* for faster training.

15 The *data-adv* set is a 100 h subset of *train-960*, which is obtained by removing long utterances from the  
 16 original Librispeech *train-100* set similarly to above. It is split into three subsets in order to perform closed-set  
 17 speaker identification experiments, since the speakers in the original train/dev/test splits are disjoint. There  
 18 are 251 speakers in *data-adv*: we assign 2 utterances per speaker to each *test-adv* and *dev-adv*. The remaining  
 19 utterances are used for training and referred to as *train-adv*.

20 For speaker verification with x-vectors [262], we use *data-spkv*, which is again derived from *data-full*.  
 21 The *train-960* subset was augmented using room impulse responses, isotropic and point-source noises [153]  
 22 as well as music and speech [261] as per the standard *sre16* recipe for training x-vectors [262] from the Kaldi

Table 4.1 Splits of Librispeech used in our experiments.

<b>dataset</b>	<b>data split</b>	<b># utts</b>	<b>duration (h)</b>
<i>data-full</i>	train-960	281,231	960.98
	test-clean	2,620	5.40
	dev-clean	2,703	5.39
	test-other	2,939	5.34
	dev-other	2,864	5.12
<i>data-adv</i>	train-adv	27,535	97.05
	dev-adv	502	1.77
	test-adv	502	1.77
<i>data-spkv</i>	train-spkv	373,985	1,388.79
	train-plda	422,491	1,443.96
	test-clean-enroll	438	0.75
	test-clean-trial	1496	3.60

toolkit [223], which we adapted to Librispeech. This increased the amount of data by a factor of 4. A subset of the augmented data containing 373,985 utterances was used to train the x-vector representation and another subset containing 422,491 utterances to train the probabilistic linear discriminant analysis (PLDA) backend. These subsets are referred to as *train-spkv* and *train-plda*, respectively. For evaluation, we built an enrollment set (*test-clean-enroll*) and a trial set (*test-clean-trial*) from the *test-clean* data. Out of 40, 29 speakers were selected from *test-clean* based on sufficient data availability. For each speaker, we selected a 1 min subset after speech activity detection<sup>2</sup> for enrollment and used the rest for trials. The same evaluation protocol was used in Chapter 3, hence the details of the trials are given in Table 3.2.

#### 4.2.2 Evaluation metrics

For all tested systems, we measure ASR<sub>eval</sub> performance in terms of the word error rate (WER) and we assess the amount of information about speaker identity in the encoded speech representation in terms of both speaker classification accuracy (ACC) and ASV<sub>eval</sub> EER. The WER is reported on the *test-clean* set. The ACC measures how well speakers can be discriminated in a closed-set setting, i.e., speakers are known at training time. It is evaluated over the *test-adv* set using the same classifier architecture as the speaker-adversarial branch of the proposed model (see Section 4.1.2). As opposed to the ACC, the EER measures how well the representations hide the speaker identity for unknown speakers, in an open-set scenario. It reflects the process of confirming whether a person is actually who the attacker thinks it might be. It is evaluated over the trial set (see Table 3.2) using x-vector-PLDA. The open-set evaluation is similar to the *Informed* attacker setting introduced in Section 3.1 because the ASV models are trained using the private representations proposed in this chapter.

The ACC and the EER will be computed for the following representations: the baseline filterbank features, the representations encoded by the network trained for ASR only (corresponding to  $\mathbf{B}_0$ ) as well as those obtained with the speaker-adversarial approach (corresponding to  $\mathbf{B}_\lambda$  for some values of  $\lambda > 0$ ). The baseline measurements are obtained using the ASR<sub>eval</sub> and the ASV<sub>eval</sub> systems, while ASR<sub>eval</sub><sup>anon</sup> and ASV<sub>eval</sub><sup>anon</sup> systems were used to measure the performance of  $\mathbf{B}_\lambda$  representations.

<sup>2</sup>Speech or voice activity detection algorithms predict whether a given time frame has speech or non-speech content [279].

**4.2.3 Network architecture and training**

For all experiments, we use the ESPnet [310] toolkit which implements the hybrid CTC/attention architecture [311]. The input features are 80-dimensional mel-scale filterbank coefficients with pitch and energy features, totalling 84 features per time frame. The *encoder* is composed of a VGG-like CNN layer followed by 5 BLSTM layers with 1,024 units. The VGG layer contains 4 convolutional layers followed by max pooling. The feature maps used in the convolution layers are of dimensions  $(1 \times 64)$ ,  $(64 \times 64)$ ,  $(64 \times 128)$  and  $(128 \times 128)$ . The attention-based decoder consists of location-aware attention [44] with 10 convolutional channels of size 100 each followed by 2 LSTM layers with 1,024 units. The CTC loss is computed over several possible label sequences using dynamic programming. In all experiments, the trade-off parameter  $\beta$  between the two decoder losses is set to 0.5. We train a single-layer recurrent neural network language model (RNNLM) with 1,024 hidden units over the *train-960* transcriptions and use it to rescore the ASR hypotheses. The resulting WER is very close to the state of the art [323] when trained on *train-960*. Finally, we implemented the *speaker-adversarial* branch via a 3 bidirectional LSTM layers with 512 units followed by a softmax layer with 251 outputs corresponding to the 251 speakers in *data-adv*. The adversarial loss  $\mathcal{L}_{\text{spk}}$  is summed across all vectors in the sequence. The speaker label  $z_i$  is duplicated to match the length of the sequence, which is smaller than  $T$  due to the subsampling performed within the encoder. Due to this subsampling as well as to the use of bidirectional LSTM layers within the encoder and the *speaker-adversarial* branch, the frame-level adversarial loss approximates well a utterance-level speaker loss that would be computed from a fixed-sized utterance-level representation, while being easier to train.

In all experiments, we start by pre-training the ASR branch for 10 epochs over *data-full* and then the speaker-adversarial branch for 15 epochs on *data-adv* in order to get a strong adversary on the pre-trained encoded representations. Then, due to time constraints, all networks are fine-tuned on *data-adv*: we run 15 epochs of adversarial training (which corresponds to simple ASR training when  $\lambda = 0$ ). Due to this, the WER is comparable to that typically achieved by end-to-end methods when trained on the *train-100* subset of Librispeech rather than the full *train-960* set. Finally, freezing the resulting encoder, we further fine-tune the speaker-adversarial branch only for 5 epochs to make sure that the reported ACC reflects the performance of a well-trained adversary.

The *encoder* network contains 133.5M parameters. To encode a 10s audio file, it perform 1.1e12 arithmetic operations which can be executed in-parallel on a 40 core CPU in 17.6s and on a single Tesla P100 GPU in 149ms.

**4.3 Results and Discussion**

We train our speaker-adversarial network for  $\lambda \in \{0, 0.5, 2.0\}$ , leading to three encoded representations  $\mathbf{B}_\lambda$ . Recall that  $\lambda = 0$  corresponds to the baseline ASR system as it ignores the speaker-adversarial branch. Table 4.2 summarizes the results.

The first column presents the ACC and EER obtained with the input filerbank features, which are consistent with the numbers reported in the literature. As expected, speaker identification and verification can be addressed to very high accuracy on those features. Using the encoded representation  $\mathbf{B}_0$  trained for ASR only already provides a significant privacy gain: the ACC is divided by 2 and the EER is multiplied by 4, which suggests that a reasonable amount of speaker information is removed during ASR training. Nevertheless,  $\mathbf{B}_0$  still contains some speaker identity information.

More interestingly, our results clearly show that adversarial training drastically reduces the performance in speaker identification but not in verification, which is conducted in an *Informed* setting, i.e., the attacker has complete knowledge of the anonymization and exploits it to train superior models. On the other hand, and counterintuitive to the speaker-invariance claims by several previous studies, we observe that

## 4.4 Summary

73

Table 4.2 ASR and speaker recognition results with different representations. WER (%) is reported on *test-clean* set, ACC (%) on *test-adv* set and EER (%) on *test-clean-trial*.

	Filterbank	$\mathbf{B}_0$	$\mathbf{B}_{0.5}$	$\mathbf{B}_{2.0}$
<b>WER</b>	–	10.9	12.5	12.5
<b>ACC</b>	93.1	46.3	6.4	2.5
<b>EER</b> Pooled	5.72	23.07	21.97	19.56
<b>EER</b> Male	3.34	19.38	18.26	16.26
<b>EER</b> Female	7.48	26.46	24.45	22.45

the verification performance actually improves after adversarial training, which implies that discrimination between the speakers became easier for an attacker. This exhibits a possible limitation in the generalization of adversarial training to unseen speakers and hence establishes the need for further investigation. The reason for the disparity between classification and verification performance might be that the speaker-adversarial branch does not inherently perform verification and hence is not optimized for that task. It might also be attributed to the representation capacity of that branch, to the number of speakers presented during adversarial training, and/or to the exact range of  $\lambda$  needed for generalizable anonymization. These factors of variation open several venues for future experiments.

We also notice that the WER stays reasonably low and stabilizes to the value of 12.5% after increasing  $\lambda$  from 0.5 to 2. In particular, for  $\lambda = 2$  the WER is just 1.6% absolute more than the baseline ( $\lambda = 0$ ).

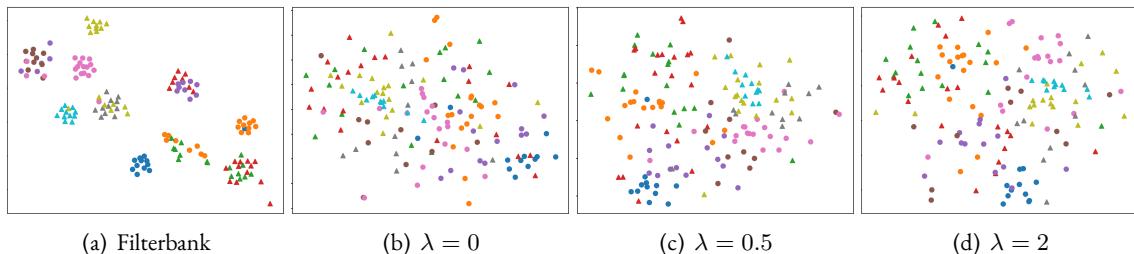


Fig. 4.3 Visualization of x-vector representations of 20 utterances of 10 speakers computed by t-SNE (perplexity equals to 30). Males are represented by circles and females by triangles.

We evaluate whether utterances from the same speaker stay in the same neighborhood or are scattered in the representation space. We compute t-SNE embeddings on the x-vector representations of 20 utterances for 10 speakers (5 male, 5 female), shown in Figure 4.3. When using filterbanks, we can observe well-clustered utterances. The clusters break down when training the x-vectors on  $\mathbf{B}_0$ . For the x-vectors trained on  $\mathbf{B}_{0.5}$  and  $\mathbf{B}_{2.0}$ , the clusters start to re-emerge. The silhouette scores for x-vectors extracted from filterbank,  $\mathbf{B}_0$ ,  $\mathbf{B}_{0.5}$  and  $\mathbf{B}_{2.0}$  representations are 0.14, -0.17, -0.05 and -0.09 respectively, are consistent with the observed EER values.

## 4.4 Summary

We proposed to combine CTC and attention losses with a speaker-adversarial loss within an end-to-end framework with the goal of learning privacy-preserving representations for ASR. Such representations could be safely transmitted to cloud-services for decoding. We investigate the level of speaker identity anonymization

1 achieved by adversarial training through closed-set speaker classification and open-set speaker verification  
2 measures. Adversarial training appears to dramatically reduce the closed-set classification accuracy, seemingly  
3 indicating a high-level of anonymization. However, this observation does not match with the open-set  
4 verification results conducted in an *Informed* setting, which correspond to a strict but real scenario of a  
5 strong adversary trying to confirm the identity of a suspected speaker. Hence we conclude that the adversarial  
6 training does not immediately generalize to produce anonymous representations in speech. We hypothesize  
7 that this disparity might be attributed to the representation capacity of the adversarial branch, the size of  
8 the training set, the formulation of the adversarial loss, and/or the value of the trade-off parameter with the  
9 ASR loss.

10 As a future work, we plan to modify the speaker adversarial branch to inherently optimize for verification  
11 instead of classification and ascertain the impact of these experimental choices over different datasets,  
12 including for languages not seen in training. In Chapter 6, the residual speaker information is removed  
13 from the ASR bottleneck representation by adding differentially private noise, and then a target speech  
14 signal is generated using them. It remains to be seen how well the representations generated by an adversarial  
15 network can perform in terms of privacy when they are used to generate an intelligible speech signal.

# Chapter 5

## X-vector based Anonymization

The measure of intelligence is the ability to change.

---

Albert Einstein

As of now, we have introduced the idea of anonymization techniques that allow speakers to publish their voice data privately. It has also been explained how such techniques can be strictly validated by simulating different attack conditions. In Chapter 3, the proposed methods produce intelligible speech signals as output, while the methods in Chapter 4 estimate private neural representations that can be either used to synthesize a speech signal or directly transmitted to cloud-based services for ASR decoding. The outputs of the techniques proposed in both the chapters were evaluated under *Informed* attack conditions. It is preferred to have a waveform as the output due to the following two reasons: firstly, waveforms are easy to validate in terms of privacy and naturalness; and secondly, due to their wide usability as published speech corpora. Hence, this chapter introduces an anonymization pipeline that replaces the speaker’s identity in an utterance with a fake, private identity and then uses speech synthesis to generate an anonymized utterance.

### 5.1 Fixed-pool voice conversion

Voice conversion methods are a crucial component for designing speaker anonymization techniques as discussed in Chapter 3. We briefly reviewed the different types of VC methods in Section 2.3.3 and then proposed a criteria to choose the method that is most suitable for the task of anonymization, i.e., non-parallel, many-to-many, and source/language independent. Among these, the criteria of being many-to-many is based on the assumption that the potential target speakers must always be present the training set of the VC method, and then the anonymization algorithm would be allowed to choose from them during deployment. The criteria of source-independence obviates the need for the source speaker to be present in the training set of the VC method, hence it is desirable for anonymization. The VC methods that are used to design anonymization strategies in Chapter 3 and many other VC methods described in Section 2.3.3 learn a mapping function  $V$ , either separately for each source-target pair or a single function for all possible pairs, i.e., many-to-many like StarGAN [140]. The limitation of both source and target speakers to be present in the training set severely restricts the capacity of anonymization techniques by fixing the pool of speakers. It is generally hard to quickly scale these techniques to expand the pool size without re-training.

Anonymization algorithms derive their strength from the amount of randomness that can be induced in the output of the algorithm. A large pool of speakers that could be flexibly scaled would allow the techniques

- <sup>1</sup> to select arbitrary unseen speaker as the target or even mix several targets to forge an imaginary sample in  
<sup>2</sup> speaker space, i.e., a *pseudo-speaker*.

### <sup>3</sup> 5.2 Flexible-pool voice conversion

<sup>4</sup> Although the three VC algorithms compared in Chapter 3 satisfy the three abovementioned criteria, they  
<sup>5</sup> did not allow conversion conditioned over a continuous speaker representation, such as x-vectors [262]. We  
<sup>6</sup> briefly discussed in Section 2.5 how the technique proposed by Fang et al. [83] relax this limitation and  
<sup>7</sup> introduce a VC framework based on speech synthesis that does not require source and target speakers to be  
<sup>8</sup> present in its training set. Figure 5.1 shows a schematic diagram of this approach where a source utterance is  
<sup>9</sup> given as the input to block 1 containing the three features extraction algorithms: the  $F_0$  extractor, the ASR  
<sup>10</sup> acoustic model, and the x-vector extractor. Within the scope of this chapter, it is assumed that all the source  
<sup>11</sup> speaker-related information is concentrated in the x-vector extracted from the utterance, and replacing it  
<sup>12</sup> with the target speaker's x-vector is sufficient to remove all the identity markers of the source speaker. This  
<sup>13</sup> assumption may not be completely true as there may be residual speaker information in other features as  
<sup>14</sup> well. We investigate the validity of this assumption in Chapter 6.

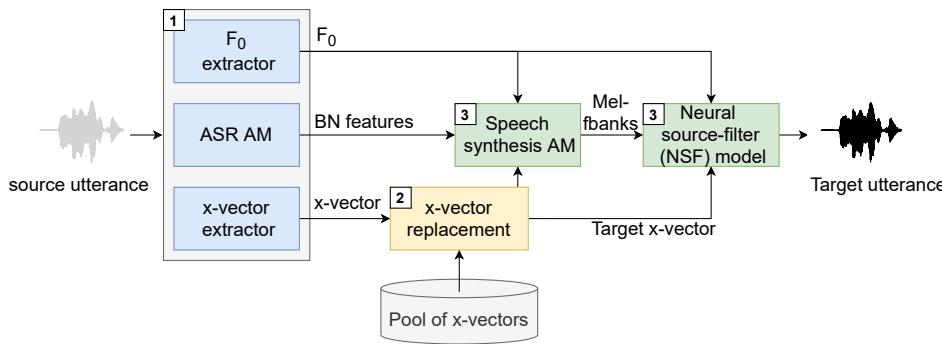


Fig. 5.1 Speech synthesis based VC framework conditioned upon continuous speaker representation that can be replaced by unseen targets.

<sup>15</sup> Nonetheless, in block 2 the new target speaker x-vector is selected from an external pool of speakers for  
<sup>16</sup> identity replacement. This pool can be expanded by simply adding more speaker x-vectors to the pool, given  
<sup>17</sup> that the x-vectors are extracted using the extractor in block 1. The new target speaker x-vector, along with  
<sup>18</sup> the original  $F_0$  and the original BN features extracted from the ASR AM in block 1, is passed to block 3, i.e.,  
<sup>19</sup> the speech synthesis block which first generates the Mel-filterbank features using the acoustic model and  
<sup>20</sup> then the target waveform using the NSF model. The speech synthesis acoustic model and the NSF model  
<sup>21</sup> are described in detail in Section 2.2.4.

### <sup>22</sup> 5.3 The first VoicePrivacy challenge

<sup>23</sup> The first VoicePrivacy challenge<sup>1</sup> was launched in February 2020 to introduce the general public to the  
<sup>24</sup> objectives of the VoicePrivacy initiative [288]. It aims to promote the development of privacy preservation  
<sup>25</sup> tools for speech technology by gathering a new community to define the tasks of interest and the evaluation  
<sup>26</sup> methodology, and benchmarking solutions through a series of challenges. Specifically, the goals of the first  
<sup>27</sup> challenge match the central goals of this thesis as described in Section 1.2, i.e., to develop anonymization

<sup>1</sup><https://www.voiceprivacychallenge.org/>

solutions that suppress personally identifiable information contained within speech signals, and at the same time, they should preserve linguistic content and speech quality/naturalness. In this section, we briefly recall the task, the data sets used for training and testing, the evaluation metrics, the proposed baseline systems, and the obtained results. We limit our description to the selected parts of the challenge that are relevant to understand the contributions made in this thesis.

### 5.3.1 Anonymization task

Recall the threat model depicted in Figure 3.1, and the actors introduced in Section 3.1, i.e., the speakers, the users, and the attackers. Similarly, in the context of the VoicePrivacy initiative, privacy preservation is formulated as a game between speakers who publish some data and attackers who access this data or data derived from it and wish to infer information about the speakers [228, 267]. To protect their privacy, the speakers publish data that contain as little personal information as possible while allowing one or more downstream goals to be achieved. To infer personal information, the attackers may use additional prior knowledge.

Focusing on speech data, a given privacy preservation scenario is specified by: (i) the nature of the data: waveform, features, etc., (ii) the information seen as personal: speaker identity, traits, spoken contents, etc., (iii) the downstream goal(s): human communication, automated processing, model training, etc., (iv) the data accessed by the attackers: one or more utterances, derived data or model, etc., (v) the attackers' prior knowledge: previously published data, privacy preservation method applied, etc. Different specifications lead to different privacy preservation methods from the speakers' point of view and different attacks from the attackers' point of view.

In the context of the VoicePrivacy 2020 challenge, the following scenario is considered where each speaker passes his/her utterances through an anonymization system to hide his/her identity. The resulting anonymized utterances are referred to as *trial* data. They sound as if they had been uttered by another speaker called pseudo-speaker, which may be an artificial voice not corresponding to any real speaker. The task of challenge participants is to design this anonymization system. In order to allow all downstream goals to be achieved, this system should: (a) output a speech waveform, (b) hide speaker identity as much as possible, (c) distort other speech characteristics as little as possible, (d) ensure that all trial utterances from a given speaker appear to be uttered by the same pseudo-speaker, while trial utterances from different speakers appear to be uttered by different pseudo-speakers<sup>2</sup>.

**Attack models and evaluation** The attack model and evaluation metrics for the challenge are very similar to the ones described in Section 3.1 and 3.4, respectively. Recall that the attackers have access to: (a) one or more anonymized trial utterances, (b) possibly, original or anonymized *enrollment* utterances for each speaker. They do not have access to the anonymization system applied by the user<sup>3</sup>. The protection of personal information is assessed via *privacy* metrics, including objective speaker verifiability, and subjective speaker verifiability and linkability. These metrics assume different attack models.

For instance, the objective speaker verifiability metrics assume that the attackers have access to a single anonymized trial utterance and several enrollment utterances. Three sets of metrics are used for original vs. anonymized enrollment data (see Section 5.3.3). In the latter case, it is assumed that the trial and enrollment utterances of a given speaker have been anonymized using the same system, but the corresponding pseudo-speakers are different, i.e., the *Lazy-Informed* scenario. On the other hand, the subjective speaker verifiability

<sup>2</sup>This is akin to “pseudonymization”, which replaces each speaker's identifiers by a unique key. This term is not used here, since it often refers to the distinct case when the identifiers are tabular data and the data controller stores the correspondence table linking speakers and keys.

<sup>3</sup>This case resembles the *Ignorant* and *Lazy-Informed* attacks.

metric assumes that the attackers have access to a single anonymized trial utterance and a single original enrollment utterance. And finally, the subjective speaker linkability metric assumes that the attackers have access to several anonymized trial utterances.

#### 5.3.2 Datasets

Several publicly available corpora are used for the training, development and evaluation of speaker anonymization systems.

**Training set** The training set comprises the 2 800 h *VoxCeleb-1,2* speaker verification corpus [203, 49] and 600 h subsets of the *LibriSpeech* [215] and *LibriTTS* [325] corpora, which were initially designed for ASR and speech synthesis, respectively. The selected subsets are detailed in Table 5.1.

Table 5.1 Statistics of the training data sets.

Subset	Size,h	Number of Speakers			Number of Utterances
		Female	Male	Total	
VoxCeleb-1,2	2 794	2 912	4 451	7 363	1 281 762
LibriSpeech: train-clean-100	100	125	126	251	28 539
LibriSpeech: train-other-500	497	564	602	1 166	148 688
LibriTTS: train-clean-100	54	123	124	247	33 236
LibriTTS: train-other-500	310	560	600	1 160	205 044

**Development set** The development set involves *LibriSpeech dev-clean* and a subset of the VCTK corpus [298] denoted as *VCTK-dev* (see Table 5.2). With the above attack models in mind, we split them into trial and enrollment subsets. For *LibriSpeech dev-clean*, the speakers in the enrollment set are a subset of those in the trial set. For *VCTK-dev*, we use the same speakers for enrollment and trial and we consider two trial subsets: *common* and *different*. The *common* subset comprises utterances #1 – 24 in the VCTK corpus that are identical for all speakers. This is meant for subjective evaluation of speaker verifiability/linkability in a text-dependent manner. The enrollment and *different* subsets comprises distinct utterances for all speakers.

Table 5.2 Statistics of the development data sets.

Subset		Female	Male	Total
LibriSpeech: dev-clean	Speakers in enrollment	15	14	29
	Speakers in trials	20	20	40
	Enrollment utterances	167	176	343
	Trial utterances	1 018	960	1 978
VCTK-dev	Speakers	15	15	30
	Enrollment utterances	300	300	600
	Trial utterances (common)	344	351	695
	Trial utterances (different)	5 422	5 255	10 677

**Evaluation set** Similarly, the evaluation set comprises *LibriSpeech test-clean* and a subset of VCTK called *VCTK-test* (see Table 5.3).

Table 5.3 Statistics of the evaluation data sets.

Subset		Female	Male	Total
LibriSpeech: test-clean	Speakers in enrollment	16	13	29
	Speakers in trials	20	20	40
	Enrollment utterances	254	184	438
	Trial utterances	734	762	1496
VCTK-test	Speakers	15	15	30
	Enrollment utterances	300	300	600
	Trial utterances (common)	346	354	700
	Trial utterances (different)	5 328	5 420	10 748

### 5.3.3 Objective and subjective metrics

Following the attack models in Section 5.3.1, objective and subjective privacy metrics are considered to assess anonymization performance in terms of speaker verifiability and linkability. We also propose objective and subjective utility metrics to assess whether the requirements in Section 5.3.1 are fulfilled. For objective evaluation, an ASV system ( $ASV_{eval}$ ) is trained to assess speaker verifiability and an ASR system ( $ASR_{eval}$ ) is trained to assess ASR decoding error in terms of the WER. Both systems are trained on *LibriSpeech train-clean-360* (Table 5.4) using Kaldi [223].

Table 5.4 Statistics of the training data set for the  $ASV_{eval}$  and  $ASR_{eval}$  evaluation systems.

Subset	Size,h	Number of Speakers			Number of Utterances
		Female	Male	Total	
LibriSpeech: train-clean-360	363.6	439	482	921	104 014

**Objective privacy and utility metrics** The  $ASV_{eval}$  system for *speaker verifiability* evaluation relies on x-vector / PLDA [262] setup as described in detail in Section 2.2.5. Four privacy metrics that are also described in Section 3.4.1 are computed, i.e., the EER, the log-likelihood ratio costs  $C_{llr}$  and  $C_{llr}^{min}$ , and the Linkability  $D_{\leftrightarrow}^{sys}$ . As shown in Fig. 3.5, these metrics are computed for: (1) original trial and enrollment data (original), (2) anonymized trial and original enrollment data (*Ignorant*), (3) anonymized trial and enrollment data (*Lazy-Informed*), and (4) anonymized trial and enrollment data with  $ASV_{eval}^{anon}$  system re-trained using anonymized training data (*Semi-Informed*). The number of target and impostor trials is given in Table 5.5.

The *ASR decoding error* is computed using  $ASR_{eval}$  that is based on the state-of-the-art Kaldi recipe for LibriSpeech involving a TDNN-F acoustic model (Refer Sec. 2.2.3, ‘Conventional approach’) and a trigram language model. As shown in Fig. 3.4, the (1) original and (2) anonymized trial data is decoded using the provided pretrained  $ASR_{eval}$  model and the corresponding WERs are calculated.

**Subjective privacy and utility metrics** Subjective metrics include *speaker verifiability*, *speaker linkability*, *speech intelligibility*, and *speech naturalness*. They are evaluated using listening tests carried out by the organizers.

To evaluate subjective speaker verifiability, listeners are given pairs of one anonymized trial utterance and one distinct original enrollment utterance of the same speaker from the data set described in Table 5.6.

Table 5.5 Number of speaker verification trials in objective evaluation on speaker verifiability.

<b>Subset</b>	<b>Trials</b>	<b>Female</b>	<b>Male</b>	<b>Total</b>
Development	LibriSpeech dev-clean Target	704	644	1 348
	Impostor	14 566	12 796	27 362
VCTK-dev	Target (common)	344	351	695
	Target (different)	1 781	2 015	3 796
Evaluation	Impostor (common)	4 810	4 911	9 721
	Impostor (different)	13 219	12 985	26 204
test-clean	LibriSpeech Target	548	449	997
	Impostor	11 196	9 457	20 653
VCTK-test	Target (common)	346	354	700
	Target (different)	1 944	1 742	3 686
	Impostor (common)	4 838	4 952	9 790
	Impostor (different)	13 056	13 258	26 314

Table 5.6 Number of trials evaluated in subjective evaluation on verifiability, intelligibility, and naturalness. Anonymized trials for subjective evaluation are from 9 anonymized systems (baseline and primary participants' systems as described in [289]). The number of speakers is 30 (15 male and 15 female) in each data set.

<b>Test set</b>	<b>Trials</b>	<b>Female</b>	<b>Male</b>	<b>Total</b>
LibriSpeech test-clean	Original	1330	1330	2660
	Anonymized	1330	1330	2660
VCTK-test (common)	Original	1380	1380	2760
	Anonymized	1380	1380	2760
VCTK-test (different)	Original	1340	1340	2680
	Anonymized	1340	1340	2680

1 Following [180], they are instructed to imagine a scenario in which the anonymized sample is from an  
 2 incoming telephone call, and to rate the similarity between the voice and the original voice using a scale  
 3 of 1 to 10, where 1 denotes ‘different speakers’ and 10 denotes ‘the same speaker’ with highest confidence.  
 4 The second subjective metric assesses speaker linkability, i.e., the ability to cluster several utterances into  
 5 speakers. Listeners are asked to place a set of anonymized trial utterances from different speakers in a 1-  
 6 or 2-dimensional space according to speaker similarity. This relies on a graphical interface, where each  
 7 utterance is represented as a point in space and the distance between two points expresses subjective speaker  
 8 dissimilarity.

9 Listeners are also asked to rate the intelligibility of individual samples (anonymized trial utterances or  
 10 original enrollment utterances) on a scale from 1 (totally unintelligible) to 10 (totally intelligible). The results  
 11 can be visualized through DET curves. Finally, the naturalness of the anonymized speech will be evaluated  
 12 on a scale from 1 (totally unnatural) to 10 (totally natural).

### 5.3.4 Anonymization baselines

Two different baseline systems have been developed for the challenge<sup>4</sup>: (1) anonymization using x-vectors and neural waveform models, and (2) anonymization using McAdams coefficient.

**Baseline-1: Anonymization using x-vectors and neural waveform models** The primary baseline system for the VoicePrivacy 2020 challenge is based on the flexible-pool VC method described in Section 5.2. The anonymization is performed in three steps as indicated by the three blocks in Figure 5.1, i.e., (*Step 1*) extraction of x-vector [262], pitch (F0) and bottleneck (BN) features; (*Step 2*) x-vector anonymization (labeled as x-vector replacement); (*Step 3*) speech synthesis (SS) from the anonymized x-vector and the original  $F_0$ +BN features.

Table 5.7 Baseline-1 system: model architectures, objective functions, output features that are used in the anonymization pipeline, and training corpora are mentioned. Superscript numbers represent feature dimensions.

#	Model	Description	Output features	Training data set
1	ASR AM	TDNN-F Input: MFCC <sup>40</sup> + i-vectors <sup>100</sup> 17 TDNN-F hidden layers Output: 6032 tied states LF-MMI (Eq. (2.17)) and CE criteria (Eq. (2.18))	BN <sup>256</sup> features extracted from the final hidden layer	Librispeech: train-clean-100 train-other-500
2	X-vector extractor	TDNN Input: MFCC <sup>30</sup> 7 hidden layers + 1 stats pooling layer Output: 7232 speaker ids CE criterion	speaker x-vectors <sup>512</sup>	VoxCeleb: 1, 2
3	Speech AM	Autoregressive (AR) network Input: $F_0^1$ + BN <sup>256</sup> + x-vectors <sup>512</sup> FF * 2 + BLSTM + AR + LSTM * 2 + highway-postnet MSE criterion	Mel-filterbanks <sup>80</sup>	LibriTTS: train-clean-100
4	NSF model	h-sinc-NSF in [305] Input: $F_0^1$ + Mel-fbanks <sup>80</sup> + x-vectors <sup>512</sup> STFT criterion (Eq. (2.24))	speech waveform	LibriTTS: train-clean-100
5		Pool of speaker x-vectors		LibriTTS: train-other-500

In *Step 1*, to extract BN features, an ASR acoustic model (AM) is trained (#1 in Table 5.7). It is assumed that these BN features represent the linguistic content of the speech signal. The ASR AM has a TDNN-F model architecture that is described as the ‘Conventional approach’ in Section 2.2.3, and is trained using the Kaldi toolkit [223]. To encode speaker information, an x-vector extractor with a TDNN model topology (#2 in Table 5.7) is also trained using Kaldi. In this step,  $F_0$  is estimated using YAAPT pitch extractor as explained under sub-heading ‘Pitch’ in Section 2.2.

<sup>4</sup>Both baseline systems are available online: <https://github.com/Voice-Privacy-Challenge/Voice-Privacy-Challenge-2020>

In Step 2, for a given source speaker, a new anonymized x-vector is computed by averaging a set of candidate x-vectors from the speaker pool for which the similarity to the x-vector of the source speaker is in the given range. The cosine distance  $\cos(\mathbf{v}_1, \mathbf{v}_2)$  or, optionally, PLDA distance is used as a similarity measure between two x-vectors  $\mathbf{v}_1$  and  $\mathbf{v}_2$ . The candidate x-vectors for averaging are chosen in two steps. First, for a given x-vector,  $N_w$  the most farthest candidates from the speaker pool (#5 in Table 5.7) are selected. Second, a smaller subset of  $\bar{N}_w$  x-vector candidates from this set are chosen randomly<sup>5</sup>. The x-vectors for the speaker pool are extracted from a disjoint data set (*LibriTTS-train-other-500*).

In Step 3, two modules are used to generate the speech waveform: a speech synthesis AM that generates Mel-filterbank features given the F0, the anonymized x-vector, and the BN features, and a NSF waveform model [305] that produces a speech waveform given the F0, the anonymized x-vector, and the generated Mel-filterbanks. These two models are described in detail in Section 2.2.4. Both models (#3 and #4 in Table 5.7) are trained on the same corpus (*LibriTTS-train-clean-100*).

**Baseline-2: Anonymization using McAdams coefficient** A secondary, alternative baseline is proposed based on speech transformation which, in contrast to the primary baseline, does not require any training data. It employs the McAdams coefficient [192] to achieve anonymisation by shifting the pole positions derived from linear predictive coding (LPC) analysis of speech signals. A brief explanation of this method is given in Section 2.5, under the sub-heading ‘Anonymization attempts using speech transformation’. Readers are referred to [218] for more details.

### 5.3.5 Results

The two baseline anonymization systems are evaluated objectively and subjectively in terms of privacy and utility. The results are mentioned below.

**Objective evaluation** Table 5.8 reports the values of objective speaker verifiability metrics obtained before/after anonymization with Baseline-1. The EER,  $C_{llr}^{\min}$  and  $D_{\leftrightarrow}^{\text{sys}}$  metrics behave similarly, while interpretation of  $C_{llr}$  is more challenging due to non-calibration<sup>6</sup>. We hence focus on the EER below. On all data sets, anonymization of the trial data greatly increases the EER. This shows that the anonymization baseline effectively increases the users’ privacy. The EER estimated in *Ignorant* setting (47 to 58%), which is comparable to or above the chance value (50%), suggests that full anonymization has been achieved. However, *Lazy-Informed* scenario result in a much lower EER (26 to 37%), which suggests that  $F_0$  and BN features retain some information about the original speaker. If the attackers have access to anonymized enrollment data, they will be able to re-identify users almost half of the time. Stricter evaluation in *Semi-Informed* setting further reduces the EER (7 to 18%) to a closer value to the baseline and confirms the threat that the attacker can achieve significant performance gain by re-training the re-identification system with the anonymized training set. Note also that the EER is larger for females than males on average. This further demonstrates that failing to define the attack model or assuming a naive attack model leads to a greatly overestimated sense of privacy [267].

Figure 5.2 shows a comparison between the EERs obtained when anonymization is performed using Baseline-1 and Baseline-2. It is clearly observed that for all the cases (*Ignorant*, *Lazy-Informed*, and *Semi-Informed*) case, Baseline-1 outperforms the privacy protection provided by Baseline-2. Table 5.9 reports the WER achieved before/after anonymization with Baseline-1. While the absolute WER stays below 7% on

<sup>5</sup>In the baseline, the following parameter values are used:  $N_w = 200$  and  $\bar{N}_w = 100$ ; and PLDA was used as the distance between x-vectors.

<sup>6</sup>In particular,  $C_{llr} > 1$  is not a problem, since we care more about discrimination metrics than score calibration metrics in the first edition.

## 5.3 The first VoicePrivacy challenge

LibriSpeech and 16% on VCTK, anonymization incurs a large WER increase of 19 to 67% relative. Similarly, Table 5.10 reports the WER achieved with Baseline-2, where an absolute increase of 5-15%, and a relative increase of more than 100% is observed. Hence, the results achieved by Baseline-2 are inferior, both in terms of privacy as well as utility, and are detailed in [289].

Table 5.8 Speaker verifiability achieved by the pretrained  $ASV_{eval}$  model in original, *Ignorant* and *Lazy-Informed* scenarios, and the  $ASV_{eval}^{anon}$  model in *Semi-Informed* case. Baseline-1 is used for anonymization.

Dataset	Gender	Attacker	Development				Test			
			EER (%)	$C_{llr}^{\min}$	$C_{llr}$	$D_{\leftrightarrow}^{\text{sys}}$	EER (%)	$C_{llr}^{\min}$	$C_{llr}$	$D_{\leftrightarrow}^{\text{sys}}$
LibriSpeech	Female	original	8.67	0.304	42.86	0.80	7.66	0.183	26.79	0.89
		<i>Ignorant</i>	50.14	0.996	144.11	0.08	47.26	0.995	151.82	0.07
		<i>Lazy-Informed</i>	36.79	0.894	16.35	0.22	32.12	0.839	16.27	0.29
		<i>Semi-Informed</i>	18.89	0.563	6.9	0.56	12.23	0.384	3.0	0.70
(common)	Male	original	1.24	0.034	14.25	0.97	1.11	0.041	15.30	0.95
		<i>Ignorant</i>	57.76	0.999	168.99	0.10	52.12	0.999	166.66	0.08
		<i>Lazy-Informed</i>	34.16	0.867	24.72	0.24	36.75	0.903	33.93	0.19
		<i>Semi-Informed</i>	7.45	0.241	3.6	0.81	10.69	0.329	5.1	0.68
VCTK	Female	original	2.61	0.088	0.868	0.93	2.89	0.091	0.866	0.92
		<i>Ignorant</i>	49.71	0.995	172.049	0.08	48.27	0.994	162.53	0.07
		<i>Lazy-Informed</i>	27.91	0.741	7.20	0.36	31.21	0.83	9.015	0.27
		<i>Semi-Informed</i>	14.53	0.473	1.6	0.62	18.79	0.552	2.0	0.54
(different)	Male	original	1.425	0.050	1.559	0.96	1.13	0.036	1.041	0.97
		<i>Ignorant</i>	54.99	0.999	192.924	0.09	53.39	1.000	190.136	0.07
		<i>Lazy-Informed</i>	33.33	0.840	23.891	0.24	31.07	0.835	21.68	0.27
		<i>Semi-Informed</i>	16.81	0.518	2.8	0.58	13.28	0.413	1.9	0.65

Table 5.9 ASR decoding error achieved by the pretrained  $ASR_{eval}$  model. Baseline-1 is used for anonymization.

Dataset	Anonymization	Dev. WER (%)	Test WER (%)
LibriSpeech	original	3.83	4.15
	anonymized	6.39	6.73
(comm.+diff.)	original	10.79	12.82
	anonymized	15.38	15.23

**Subjective evaluation** In this section, the subjective analysis results for naturalness, intelligibility, and speaker verifiability of the anonymized speech data are presented and compared with the MOS on target and non-target original speech. There are four types of trials: original or anonymized trials from target or non-target speakers. When displaying the results of naturalness and intelligibility, the anonymized trials

1  
2  
3  
4

5  
6  
7  
8

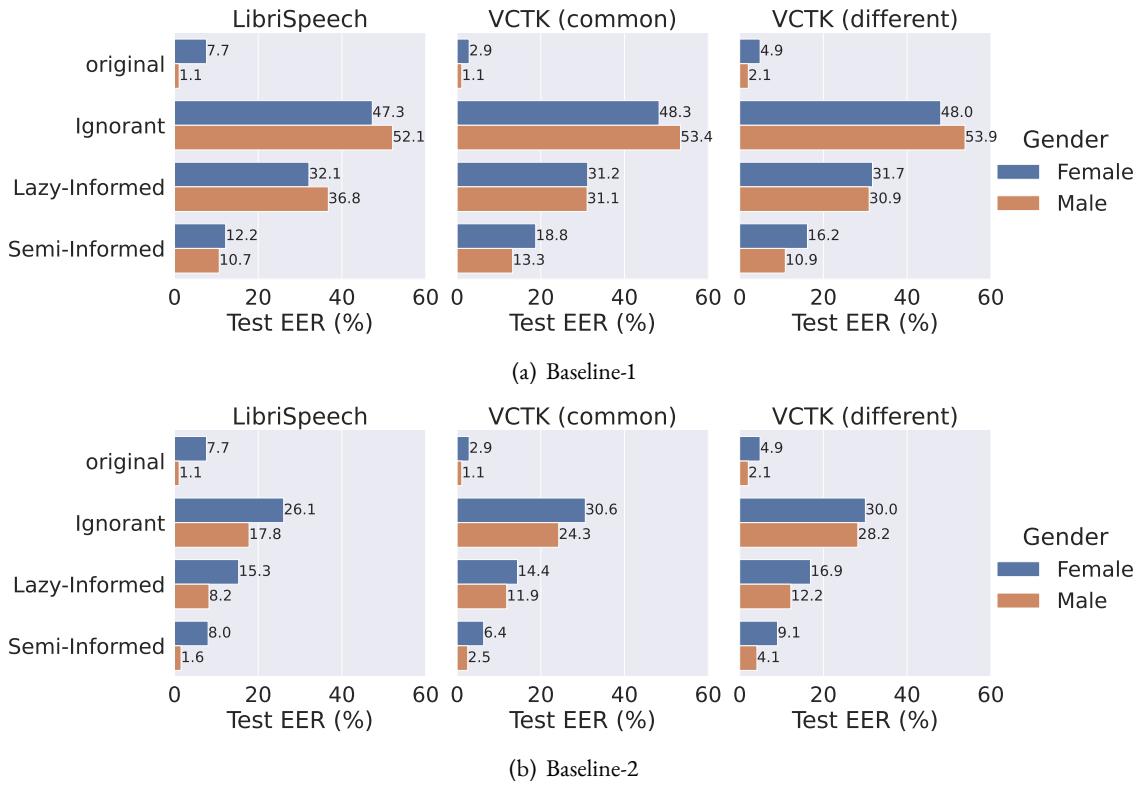


Fig. 5.2 Objective evaluation of privacy protection provided by the two baseline systems in the first VoicePrivacy challenge. Higher EER indicates better protection.

1 of both target and non-target speakers are merged. It is only on speaker verifiability or speaker similarity  
 2 that they need to be separated so that it can be interpreted how well the anonymization system anonymize  
 3 the speech of the target speakers. Hence there are four sub-figures in Figure 5.3. The scores of target and  
 4 non-target original trials are displayed separately for comparison. To reduce perceptual bias of each individual  
 5 evaluator, naturalness, intelligibility, and verifiability scores from the unified subjective test was processed  
 6 using normalized-rank normalization [238]. The processed scores are floating numbers varying from 0 to 1.  
 7 The results of naturalness and intelligibility are as expected. Anonymized samples from both the systems,  
 8 i.e. B1 and B2 are inferior to the target and non-target original data, and the differences are statistically  
 9 significant as per the Mann-Whitney-U tests [238] conducted in [289]. More efforts are necessary to  
 10 address the degradation caused by the two proposed anonymization methods. On similarity, anonymized

Table 5.10 ASR decoding error achieved by the pretrained  $ASR_{eval}$  model. Baseline-2 is used for anonymization.

Dataset	Anonymization	Dev. WER (%)	Test WER (%)
LibriSpeech	original	3.83	4.15
	anonymized	8.77	8.88
VCTK (comm.+diff.)	original	10.79	12.82
	anonymized	25.56	28.22

## 5.4 Design choices in x-vector space

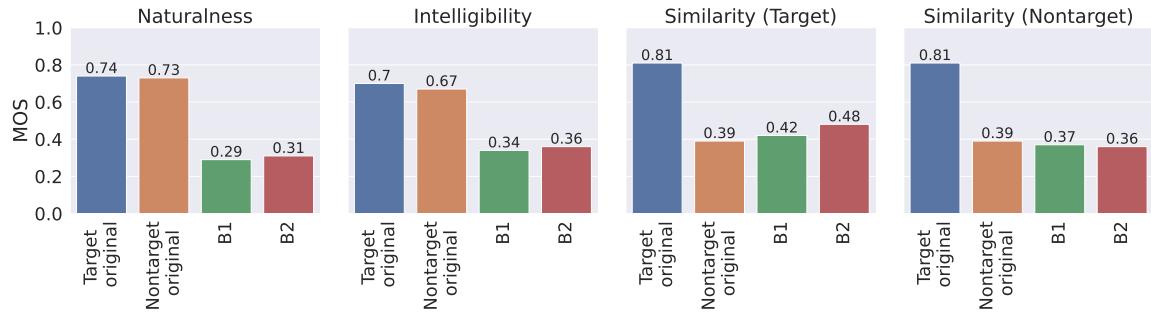


Fig. 5.3 Subjective evaluation of privacy protection provided by the two baseline systems in the first VoicePrivacy challenge. Bar plots of *subjective speech naturalness*, *intelligibility*, and *speaker similarity* obtained from the normalized scores. For naturalness and intelligibility, scores from target and non-target anonymized data are pooled; for similarity, scores for anonymized target and non-target speakers data are separately plotted in 3rd and 4th sub-figures, respectively. Numbers indicate mean values over all the three data sets. Higher values for naturalness and intelligibility correspond to better utility, and lower scores for similarity to target speaker with anonymized data from target speaker — to better privacy

data from target speakers are perceptually less similar to the enrollment data of the target speaker than the unanonymized trial of that target speaker. This performance is welcome because it indicates that both the systems achieve a certain degree of anonymization in human perception. For subjective speaker linkability results, readers are referred to [289].

## 5.4 Design choices in x-vector space

This section focuses on the explanation of the core anonymization logic implemented in the Baseline-1 algorithm, and explores the flexibility of pseudo-speaker selection through four design choices. Experiments are performed to test the hypothesis that a superior degree of privacy protection can be achieved via an optimal configuration of design choices, where each choice essentially exploits a particular property of the target speaker space that is represented by a large pool of speaker representations. We established that Baseline-1, i.e., x-vector based speaker anonymization is the leading approach in the first VoicePrivacy Challenge, which converts the speaker’s voice into that of a random pseudo-speaker. In this section, it is shown that the strength of anonymization varies significantly depending on how the pseudo-speaker is chosen.

Going beyond the experiments conducted in Section 5.3.5, the quality of anonymization is assessed from the perspective of the three actors involved in our threat model (Refer Sec. 3.1), namely the speaker, the user and the attacker. To measure privacy and utility, we use respectively the linkability score  $D_{\leftrightarrow}^{\text{sys}}$  achieved by the attackers and the decoding word error rate achieved by an ASR<sub>eval</sub> model trained on the anonymized data. Experiments on LibriSpeech show that the best combination of design choices yields state-of-the-art performance in terms of both privacy and utility. Experiments on Mozilla Common Voice [14] further show that it guarantees the same anonymization level against re-identification attacks among 50 speakers as original speech among 20,000 speakers.

In order to implement and assess the proposed anonymization method, the following questions arise from the speaker’s and user’s perspectives: Q1: *How to optimally choose and assign the target pseudo-speaker?* Q2: *How well is utility preserved?* Q3: *How much residual speaker information remains?* Furthermore, the attacker must address the following questions: Q4: *Can privacy protection be defeated using some knowledge*

<sup>1</sup> of the anonymization method? Q5: How does the number of possible speakers affect the re-identification  
<sup>2</sup> performance?

<sup>3</sup> In this section, we extend the two target pseudo-speaker generation strategies in [83] (fully random, or  
<sup>4</sup> at a fixed distance from the original as measured by cosine distance between x-vectors) into a whole family  
<sup>5</sup> of strategies based on four design choices: the distance metric between x-vectors, the region of x-vector  
<sup>6</sup> space where the pseudo-speaker is picked, its gender, and whether to assign it to one or all utterances of the  
<sup>7</sup> original speaker. Our experiments suggest an optimal combination of design choices to balance privacy and  
<sup>8</sup> utility (answering Q1). We train and/or evaluate ASR<sub>eval</sub> and ASR<sub>eval<sup>anon</sup></sub> models on original and anonymized  
<sup>9</sup> speech to assess these two forms of utility (answering Q2). We show that some speaker information remains  
<sup>10</sup> in the pitch sequence and apply two different pitch transformation techniques to remove it (answering  
<sup>11</sup> Q3). We conduct these experiments for three types of attackers, i.e., *Ignorant*, *Lazy-Informed* and *Semi-*  
<sup>12</sup> *Informed*, where stronger attackers have more knowledge about the anonymization method (answering Q4).  
<sup>13</sup> Finally, we conduct additional experiments with more than 20,000 possible speakers (answering Q5). These  
<sup>14</sup> contributions significantly extend our preliminary study [266], which provided less detail, did not include  
<sup>15</sup> utterance- vs. speaker-level target assignment and pitch transformation, did not evaluate privacy against the  
<sup>16</sup> strongest (*Semi-Informed*) attacker or with a large number of possible speakers, and did not evaluate utility  
<sup>17</sup> for ASR training.

#### <sup>18</sup> 5.4.1 Anonymization framework

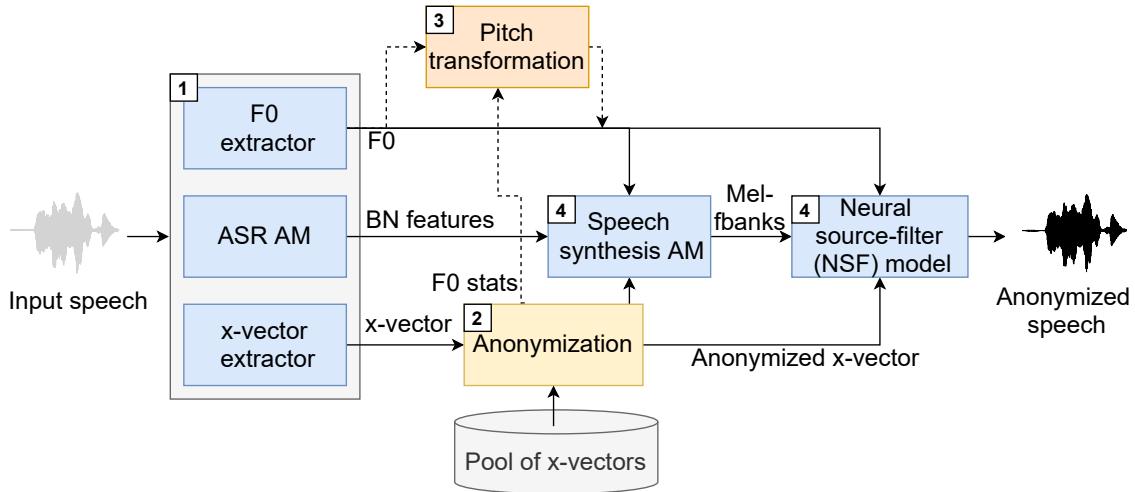


Fig. 5.4 General architecture of the anonymization system.

<sup>19</sup> In the following, we use the anonymization system shown in Fig. 5.4, that is a variant of Baseline-  
<sup>20</sup> 1 architecture described in Section 5.2. This system represents speaker identity, linguistic content and  
<sup>21</sup> intonation using x-vectors,<sup>7</sup> bottleneck (BN) features [321] (a low-dimensional representation extracted  
<sup>22</sup> from an intermediate layer of an ASR model) and pitch contour (F0), respectively. It comprises four steps:  
<sup>23</sup> *Step 1 (Feature extraction)* extracts F0 and BN features and the x-vector from the input signal. *Step 2 (X-*  
<sup>24</sup> *vector anonymization*) generates a target x-vector by averaging  $N^*$  candidate x-vectors from an external pool

<sup>7</sup>Following [83], we use raw x-vectors to represent speaker identity instead of x-vectors compressed and rotated by linear discriminant analysis (LDA), as classically done in the context of ASV. Unless the projected dimension is carefully chosen after several experiments, the impact of the LDA transformation on speaker-specific information cannot be ascertained. Hence we defer experiments with LDA-transformed x-vectors to a future study.

## 5.4 Design choices in x-vector space

87

of speakers.<sup>8</sup> *Step 3 (Pitch transformation)* is an optional step which receives the pseudo-speaker target pitch statistics from the anonymization module and transforms the original pitch. *Step 4 (Speech synthesis)* synthesizes a speech waveform from the anonymized x-vector and the original BN and F0 features using an acoustic model (AM) and the NSF model. With the exception of Step 3 which is new (see Section ??), this system is identical to Baseline-1 proposed for the first VoicePrivacy Challenge. Refer to Table 5.7 for details on the feature dimensions and the architectures of the models in Steps 1 and 4.

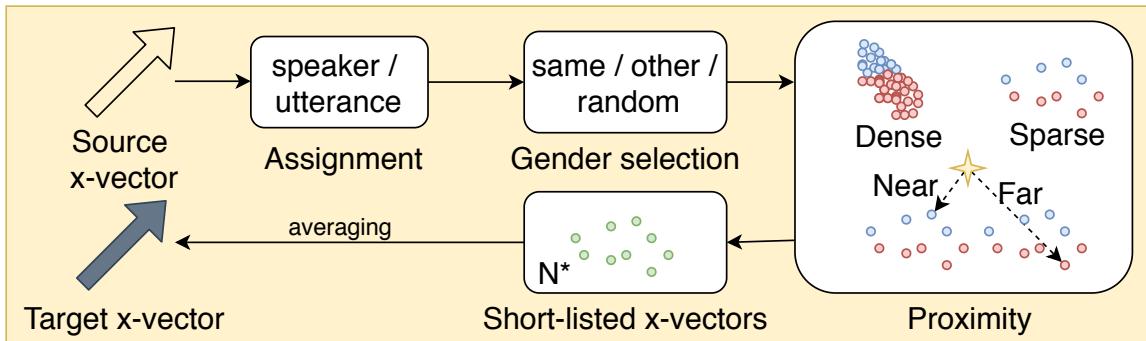


Fig. 5.5 Zoomed-in view of the x-vector anonymization step in Fig. 5.4 showing the design choices for the generation of the target x-vector.

Now given the ability to generate arbitrary external targets in Step 2 (yellow box in Fig. 5.4), the question arises of which strategy the speaker shall employ to select the candidate x-vectors and achieve a suitable privacy-utility tradeoff. Fang et al. [83] select candidate x-vectors at random within the whole pool or within a fixed interval of distances from the original x-vector. Han et al. [113] select a single target x-vector at random within a maximum distance from the original x-vector. In the following, we expand these initial strategies into a broader range of strategies governed by the choice of the distance metric between x-vectors, the region of x-vector space where the candidates are selected, their gender, and the assignment of the resulting target x-vector to one or all utterances of the original speaker. These four design choices, which are illustrated in Fig. 5.5, are detailed below. For the sake of focus, we do not explore other design choices such as the size or the diversity of the anonymization pool.

#### 5.4.2 Distance metric

To design advanced candidate selection strategies, the speaker must first choose a distance metric which dictates the properties of the x-vector space. We compare two such metrics.

The first one is the cosine distance, which was used by [83]. For a pair of x-vectors  $\mathbf{v}_i$  and  $\mathbf{v}_j$ , it is defined as

$$d_{\cos}(\mathbf{v}_i, \mathbf{v}_j) = 1 - \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\|_2 \|\mathbf{v}_j\|_2}. \quad (5.1)$$

The second metric is based on PLDA [131], that is the log-likelihood ratio of the two hypotheses that  $\mathbf{v}_i$  and  $\mathbf{v}_j$  belong to the same speaker ( $\mathcal{H}_s$ ) vs. different speakers ( $\mathcal{H}_d$ ). Previous studies [146] have shown that PLDA yields state-of-the-art performance as the similarity metric between x-vectors in the context of ASV. This is attributed to its formulation which estimates the factorized within-speaker and between-speaker variability in speaker space, making it a superior metric even for short utterances [244]. More specifically, PLDA models x-vectors  $\mathbf{v}$  as  $\mathbf{v} = m + V\mathbf{y} + D\mathbf{z}$ , where  $m$  is the center of the x-vector space, the columns

<sup>8</sup>There is no guarantee that averaging produces a valid x-vector, but all our experiments show that the synthesized anonymized speech is of good quality.

1 of  $V$  capture speaker variability (eigenvoices) with  $y$  depending only on the speaker, and the columns of  $D$   
 2 encode channel variability (eigenchannels) with  $z$  varying from one recording to another. The parameters  
 3  $m$ ,  $V$  and  $D$  are trained on x-vectors extracted using the x-vector extractor in Step 1 from VoxCeleb-1,2 data  
 4 set that is used to train that extractor itself (see Table 5.1 for details on this data set). The log-likelihood ratio  
 5 score

$$6 \quad \text{PLDA} = \log \frac{p(\mathbf{v}_i, \mathbf{v}_j | \mathcal{H}_s)}{p(\mathbf{v}_i, \mathbf{v}_j | \mathcal{H}_d)} \quad (5.2)$$

7 can be computed in closed form [237]. We propose to use minus-PLDA as the “distance” between a pair of  
 8 x-vectors.

### 9 5.4.3 Proximity

10 We propose three alternative criteria resulting in five different “proximity” choices to restrict the region of  
 11 x-vector space from which candidate x-vectors are selected.

#### 12 Random

13 The simplest candidate x-vector selection strategy is to select  $N^*$  x-vectors with a given gender uniformly at  
 14 random from the pool. Note that this strategy does not allow us to choose particular regions of interest in  
 15 the x-vector space.

#### 16 Far/near

17 Alternatively, the chosen distance metric can be used to find candidate x-vectors which resemble most (*near*)  
 18 or least (*far*) the original speaker  $S$ . In essence, we rank all the x-vectors in the pool in increasing order of their  
 19 distance from  $S$  and select either the top  $N$  (*near*) or the bottom  $N$  (*far*). To introduce some randomness,  
 20  $N^* < N$  x-vectors are selected out of these  $N$  uniformly at random.

#### 21 Dense/sparse

22 Another alternative is to identify clusters of x-vectors in the pool and rank them based on their cardinality.  
 23 We construct these clusters using the Affinity Propagation [74] algorithm (see detailed procedure in Sec-  
 24 tion 5.4.6). We filter out the cluster which is closest to the source speaker, then randomly select one cluster  
 25 among those with most (*dense*) or least (*sparse*) members.<sup>9</sup> We then randomly select half of the members of  
 26 that cluster.

27 In all five cases, the selected candidate x-vectors are averaged to obtain the target (pseudo-speaker)  
 28 x-vector.

### 29 5.4.4 Gender selection

30 In practice, instead of applying one of these five proximity choices to the entire speaker pool, we apply it to a  
 31 gender-dependent pool which consists of either all males or all females of the original pool. We propose  
 32 three possible gender selection choices: *same* where all speakers in the pool have the same gender as the  
 33 original speaker; *opposite* where they all have the opposite gender; and *random* where either of the two  
 34 gender-dependent pools is selected at random. This allows us to avoid averaging candidate x-vectors from  
 35 both genders with each other, and to assess the impact of gender selection on privacy and utility.

---

<sup>9</sup>Note that the terms *sparse* and *dense* do not directly reflect the density of x-vectors, since they do not take the diameter of the clusters into account. However, we find that this relation holds in practice.

### 5.4.5 Assignment

The generation of the anonymized waveform is conditioned upon the x-vector sequence, whose length is equal to the number of frames in the original utterance. All the x-vectors in this sequence are identical to each other to indicate a single pseudo-speaker throughout the utterance. In theory, these x-vectors should also be identical across all utterances spoken by this pseudo-speaker but, according to [230], x-vectors also contain channel, duration, and phonetic information, in addition to speaker and gender. Hence, the x-vectors computed for different utterances may exhibit some variations due to utterance-specific properties. To assess the effect of these variations on privacy and utility, we propose two assignment strategies for the target x-vector: speaker-level (*perm*) or utterance-level (*rand*). In the former case, we average the utterance-level x-vectors of all utterances of the original speaker into a single speaker-level x-vector, we generate a corresponding target x-vector, and we use it to anonymize all utterances of that speaker. In the latter case, we consider the utterance-level x-vector for each utterance of the original speaker, we generate a corresponding target x-vector (using the same distance metric, proximity, and gender across all utterances), and we use it to anonymize that utterance only.

### 5.4.6 Experimental setup

Along with the privacy evaluation using *Ignorant*, *Lazy-Informed* and *Semi-Informed* attackers that is relevant from the speakers' and attackers' perspective, the utility of ASR training is also evaluated which is relevant from the users' perspective. Unlike Section 5.3.5, no subjective analysis is conducted and the design choices are only objectively evaluated.

## Data

The experiments in Section 5.4.7 follow the VoicePrivacy Challenge setup. The training data sets for the components of the anonymization system, i.e., the ASR AM, the x-vector extractor, and the speech synthesis AM and NSF model are described in Table 5.7. The *train-other-500* subset of LibriTTS is used as the external pool of speakers for x-vector anonymization. The development and test sets are built from the *dev-clean* and *test-clean* subsets of LibriSpeech, respectively as described in Table 5.5.<sup>10</sup> Each of these two sets consists of *trial* utterances from 40 speakers and *enrollment* utterances from a subset of 29 speakers (see Section 5.4.6).

In Section 5.5, we employ the same trained models and the same external pool of speakers but we build multiple test sets from the Mozilla Common Voice [14] English corpus, in order to study the attacker's success against anonymization with a larger number of possible speakers. This corpus contains more than 52,000 speakers, out of which we select up to 24,616 male speakers (see Section 5.5.3).

## Algorithm settings

The *dense* and *sparse* anonymization choices are implemented as follows. We use Affinity Propagation [74] to cluster the speakers in the external pool. This non-parametric clustering method determines the number of clusters automatically through a message passing protocol. Two parameters govern the final number of clusters: *preference* assigns prior weights to samples which may be likely candidates for centroids, and *damping factor* is a floating-point multiplier to responsibility and availability messages. In our experiments, equal *preference* is assigned to each sample and the *damping factor* is set to 0.5. Out of 1,160 speakers in the pool, 80 clusters are found, including 46 male and 34 female. The number of speakers per cluster ranges from 6 to 36. Candidate x-vector selection is achieved by picking either the 10 clusters with least

<sup>10</sup>The VoicePrivacy Challenge involves development and evaluation sets built from LibriSpeech and VCTK. Due to space limitations, we focus on LibriSpeech.

members (*sparse*) or the 10 clusters with most members (*dense*). The remaining clusters are ignored. During anonymization, one of the 10 clusters is selected at random and 50% of its members are averaged to produce the target x-vector.

#### 4 Privacy evaluation

As explained in Section 5.3.1, privacy protection can be seen as a contest between two entities: a *speaker* who publishes anonymized utterances, and an *attacker* who attempts to uncover the speaker’s identity by comparing these utterances with utterances whose speaker is known. Following the classical ASV terminology adopted in the VoicePrivacy Challenge, these are called *trial* and *enrollment* utterances, respectively, and each such comparison is called a *trial*. The attacker has full control over the enrollment set and the speaker identities within it. Hence he/she may use some knowledge about the anonymization scheme to transform the enrollment data and reduce the mismatch with the trial data. To assess the strength of anonymization against attackers with increasing knowledge, we perform the evaluation in four scenarios similar to the ones presented in Section 5.3.5:

- *Original*: The speaker does not perform any anonymization. The attacker uses original speech for enrollment and an ASV system trained on original speech. This offers the lowest possible privacy protection.
- *Ignorant*: The speaker anonymizes his/her speech, unbeknownst to the attacker who still uses original speech for enrollment and an ASV system trained on original speech.
- *Lazy-Informed*: The speaker anonymizes his/her speech. The attacker anonymizes the enrollment data using the same anonymization system and the same design choices. However, he/she is not aware of the random numbers drawn by the speaker to obtain the *random* target gender (Section 5.4.4) or the candidate x-vectors (Section 5.4.3). Hence, different pseudo-speakers are assigned to the trial and enrollment utterances of a given speaker.
- *Semi-Informed*: The speaker anonymizes his/her speech. The attacker anonymizes the enrollment data using the same system and design choices. In addition, he/she anonymizes the training data set for the ASV<sub>eval</sub> system and re-trains it to get ASV<sub>eval</sub><sup>anon</sup>. This scenario is the one in which the speaker is most “vulnerable” despite anonymization, hence we consider it as the most trustworthy assessment of privacy.<sup>11</sup>

In Section 5.4.7, privacy is assessed in terms of the *linkability* [103, 190] achieved by an x-vector-PLDA ASV system trained on the *train-clean-360* subset of LibriSpeech (anonymized ASV<sub>eval</sub><sup>anon</sup> in the *Semi-Informed* scenario, original ASV<sub>eval</sub> otherwise). Recall that this metric computes the overlap between the distributions of PLDA scores of same-speaker and different-speaker trials as described in Section 3.4.1. It behaves similarly to the EER and  $C_{llr}^{\min}$  [37], but it does not rely on any restrictive assumption (e.g., threshold-based decision) which makes it a more trustworthy metric [190]. For the sake of reproducibility, we use the same set of trials as in Table 5.5.<sup>12</sup> Lower linkability means higher privacy.

We observe a high correlation between EER, linkability and  $C_{llr}^{\min}$  in *Lazy-Informed* setting as shown in Fig. 5.6. While the diagonal plots represent the marginal distribution of a particular metric for gender-specific

<sup>11</sup>The *Informed* scenario described in Section 3.1 where the attacker is aware of the random numbers drawn by the speaker is not part of our study, since it falls into a security problem rather than just a privacy problem.

<sup>12</sup>As classically assumed in the speaker verification literature, the two speakers in each trial have the same original gender. In practice though, the gender of the original speaker may be unknown to the attacker. Hence, the resulting linkability values can be seen as worst-case values from the speaker’s point of view and best-case values from the attacker’s point of view.

## 5.4 Design choices in x-vector space

91

trials, others indicate the correlation of two given metrics on x- and y-axis. As expected, EER and  $C_{llr}^{\min}$  are positively correlated, while both of them are negatively correlated with linkability.

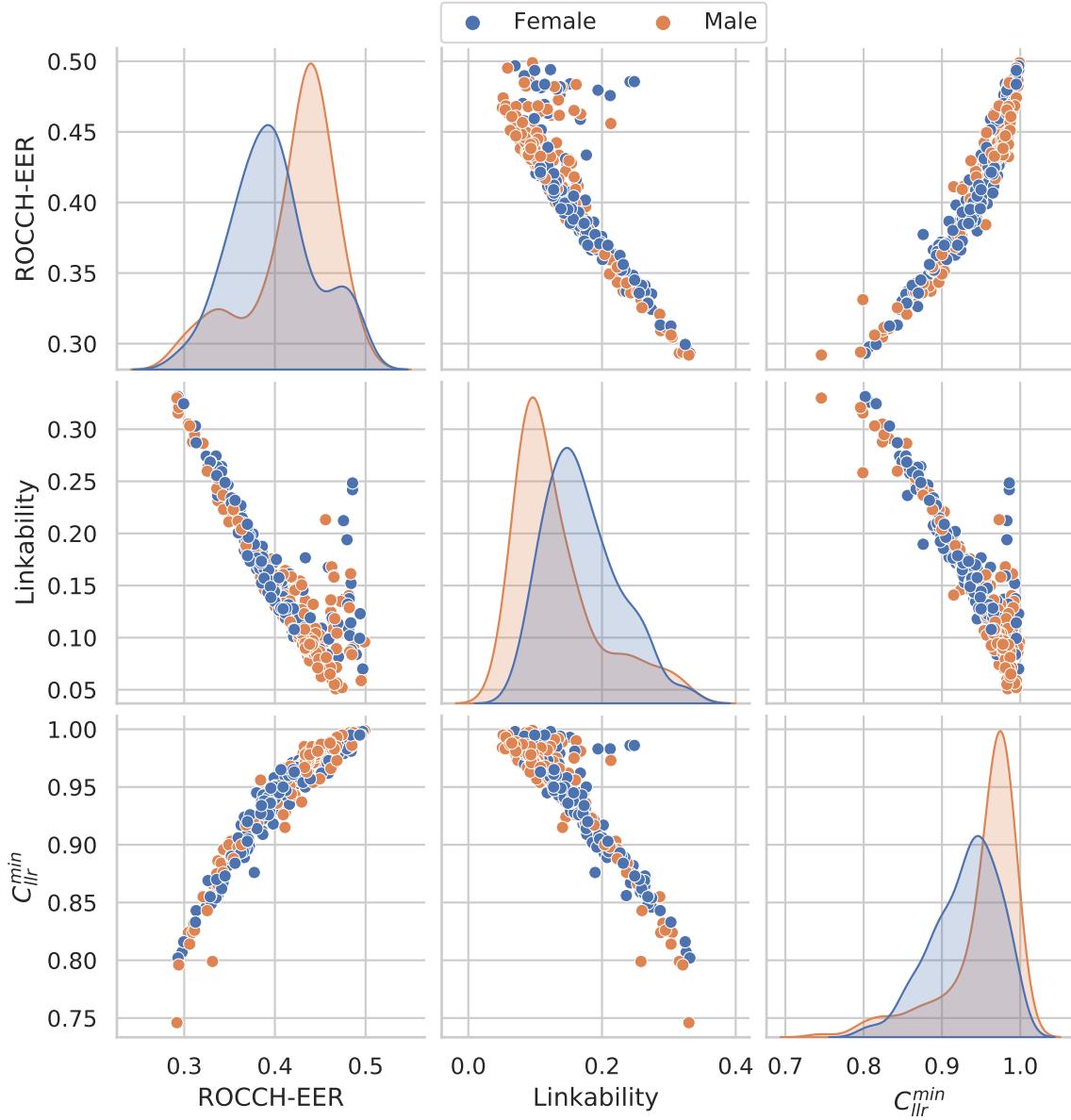


Fig. 5.6 Relationship between EER, linkability and  $C_{llr}^{\min}$  in *Lazy-Informed* setting for various combination of design choices.

In Section 5.5, we also evaluate the average *rank* of the true speaker and the *top-k precision* achieved for closed-set ASI. Instead of training speaker classification systems on subsets of Common Voice, which would overfit the speakers therein, we compute the PLDA scores between each trial utterance and all enrollment utterances (one per speaker, including the true speaker) using the same x-vector and PLDA models as in Section 5.4.7 and sort them in decreasing order. The higher the rank and the lower the top-k precision, the higher the privacy.

1  
23  
4  
5  
6  
7  
8

**1 Utility evaluation**

**2** In Section 5.4.7 (“Speaker’s perspective”), we evaluate the utility for ASR decoding in terms of the WER  
**3** achieved by an ASR<sub>eval</sub> system trained on the *train-clean-360* subset of LibriSpeech and applied to the  
**4** anonymized utterances. In Section 5.4.7 (“User’s perspective”), we evaluate the utility both for ASR decoding  
**5** and training in terms of the WER achieved by an ASR system trained either on the original (ASR<sub>eval</sub>) or the  
**6** anonymized *train-clean-360* data set (ASR<sub>eval</sub><sup>anon</sup>) and used to decode either original or anonymized speech.  
**7** For more details on the ASR system architecture, see Section 5.3.3. A lower WER indicates higher utility.

**8 5.4.7 Results and Discussion**

**9** In this section, we evaluate the design choices introduced in Section 5.4 under the VoicePrivacy Challenge  
**10** setup. Our experiments are organized according to the three actors in our threat model. First, the speaker  
**11** finds the two most promising combinations of design choices on the development set in terms of privacy  
**12** in the *Ignorant* and *Lazy-Informed* scenarios and utility for ASR decoding. This is motivated by the high  
**13** computational cost of anonymizing the *train-clean-360* subset of LibriSpeech and retraining ASV and  
**14** ASR systems on it, which prevents the evaluation of privacy in the *Semi-Informed* scenario and utility for  
**15** ASR training for all 54 combinations. Second, the user assesses the utility of these two combinations for  
**16** both ASR training and decoding. Third, the attacker quantifies the resulting privacy in the *Semi-Informed*  
**17** scenario, which leads us to identify the best combination among them. Finally, we show how the proposed  
**18** pitch transformation further improves privacy in this scenario without loss of utility.

**19 Speaker’s perspective**

**20** We first evaluate the design choices from the speaker’s perspective in terms of privacy in the *Ignorant* and  
**21** *Lazy-Informed* scenarios and utility for ASR decoding on the development set. The results are displayed in  
**22** the form of swarm plots, i.e., scatter plots where each dot represents the privacy or utility value associated  
**23** with one combination of design choices. In order to avoid overlapping dots with similar values, the dots are  
**24** spread horizontally.

**25 Distance** Figure 5.7 evaluates the effect of the chosen distance metric on privacy. We observe that both  
**26** cosine distance and PLDA result in similarly low linkability in the *Ignorant* case but PLDA marginally  
**27** outperforms cosine distance (i.e., it results in a lower linkability) in the *Lazy-Informed* case. Since both  
**28** distance measures perform similarly in terms of utility (see Fig. 5.12(a)), PLDA has an advantage. Therefore  
**29** we consider only PLDA as the distance metric in the following experiments.

**30 Proximity** Next, we assess the five choices of target *proximity* described in Section 5.4.3, namely *random*,  
**31** *near*, *far*, *sparse* and *dense*. The distance metric is fixed to PLDA and the values of  $N$  and  $N^*$  are fixed to  
**32** 200 and 100 respectively.<sup>13</sup> We discover the clusters in x-vector space and select *pseudo-speakers* from *sparse*  
**33** and *dense* clusters using the procedure described in Section 5.4.6.

**34** We observe in Fig. 5.8 that, although selecting candidate x-vectors *far* from the original speaker achieves  
**35** the lowest linkability in the *Ignorant* case together with the *random* strategy, it is largely outperformed in  
**36** the *Lazy-Informed* case by selection from *sparse* or *dense* clusters and by the *random* strategy. This shows  
**37** that clustering based pseudo-speaker mapping results in more robust anonymization as compared to simple  
**38** distance-based mapping.

---

<sup>13</sup>We noticed a sharp decline in utility for smaller values of  $N^*$ .

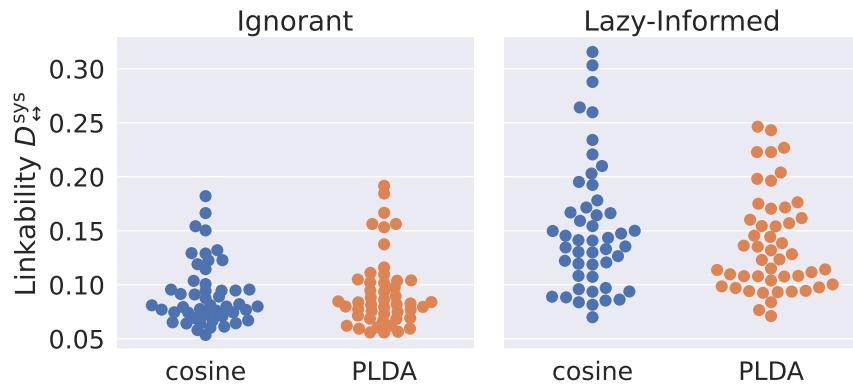


Fig. 5.7 Privacy against *Ignorant* and *Lazy-Informed* attackers depending on the distance choice and the gender of the original speaker. Each swarm plot shows the 24 linkability values for each gender on the development set resulting from all combinations of proximity (excluding *random*), gender selection, and assignment choices.



Fig. 5.8 Privacy against *Ignorant* and *Lazy-Informed* attackers depending on the proximity choice and the gender of the original speaker. Distance is fixed to PLDA. Each swarm plot shows the 6 linkability values for each gender on the development set resulting from all combinations of gender selection and assignment choices.

1      Compared to the *sparse* selection strategy, the *dense* strategy provides comparable privacy protection in  
 2      the *Lazy-Informed* case, but much higher utility (see Fig. 5.12(b)). This can be attributed to the fact that  
 3      speakers in *sparse* clusters stand out more from the crowd than those in *dense* clusters, therefore they are  
 4      more likely to suffer from poor ASR performance.

5      Finally, *random* target selection yields similar privacy protection in the *Lazy-Informed* case and slightly  
 6      better utility as compared to *dense*. Hence we consider the *random* and *dense* strategies to be the best choices  
 7      for proximity.

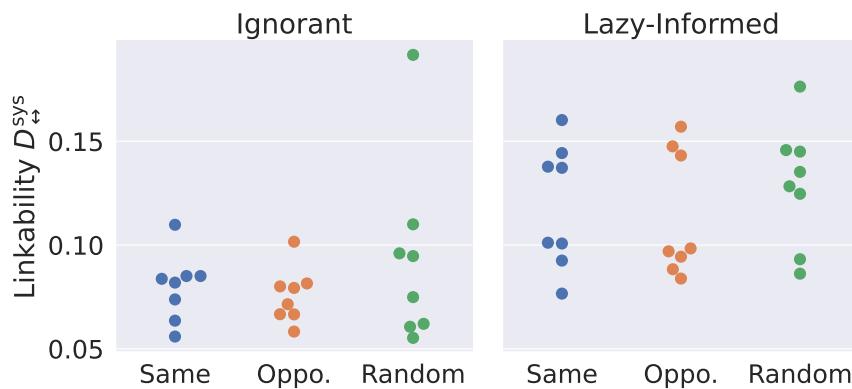


Fig. 5.9 Privacy against *Ignorant* and *Lazy-Informed* attackers depending on the gender selection choice and the gender of the original speaker. Distance is fixed to PLDA and proximity to *dense* or *random*. Each swarm plot shows the 4 linkability values for each gender on the development set resulting from the assignment choice and the 2 proximity choices.

8      **Gender selection**    We now investigate the gender selection strategy described in Section 5.4.4. The distance  
 9      is fixed to PLDA and proximity to *dense* or *random*. As per the results shown in Fig. 5.9 it is hard to find  
 10     the best choice for gender selection in terms of privacy since the linkability is not consistently lower for any  
 11     specific choice.

12     As dictated by intuition, the *same* gender pseudo-speaker mapping results in slightly better utility (see  
 13     Fig. 5.12(c)) than *opposite* or *random* gender.

14     We report in Fig. 5.10 the effect of *proximity* and *gender selection* on the distance of the target pseudo-  
 15     speaker from the original speaker. The x-axis represents the average cosine distance between the original  
 16     x-vector and the selected pseudo-speaker x-vector over all the trial utterances, while the y-axis indicates  
 17     the average cosine distance between the original speaker and the x-vector computed from the anonymized  
 18     utterance after the speech synthesis. The diagonal indicates the  $x = y$  line, which is the desired effect. With  
 19     *same* gender selection, the cosine distance between the source and the target (before and after anonymization)  
 20     is naturally lower than with the two other gender selection choices and the actual distance (after anonymiza-  
 21     tion) is larger than the desired distance (before anonymization). With *other* gender selection, the actual  
 22     distance is typically lower than the desired distance instead. *Random* gender selection lies in the middle and  
 23     most of the points adhere to the diagonal, which indicates that the desired distance is preserved after speech  
 24     synthesis.

25     Overall, since *random* gender selection produces robust privacy and utility results across both genders,  
 26     we consider it to be the best choice for gender selection.

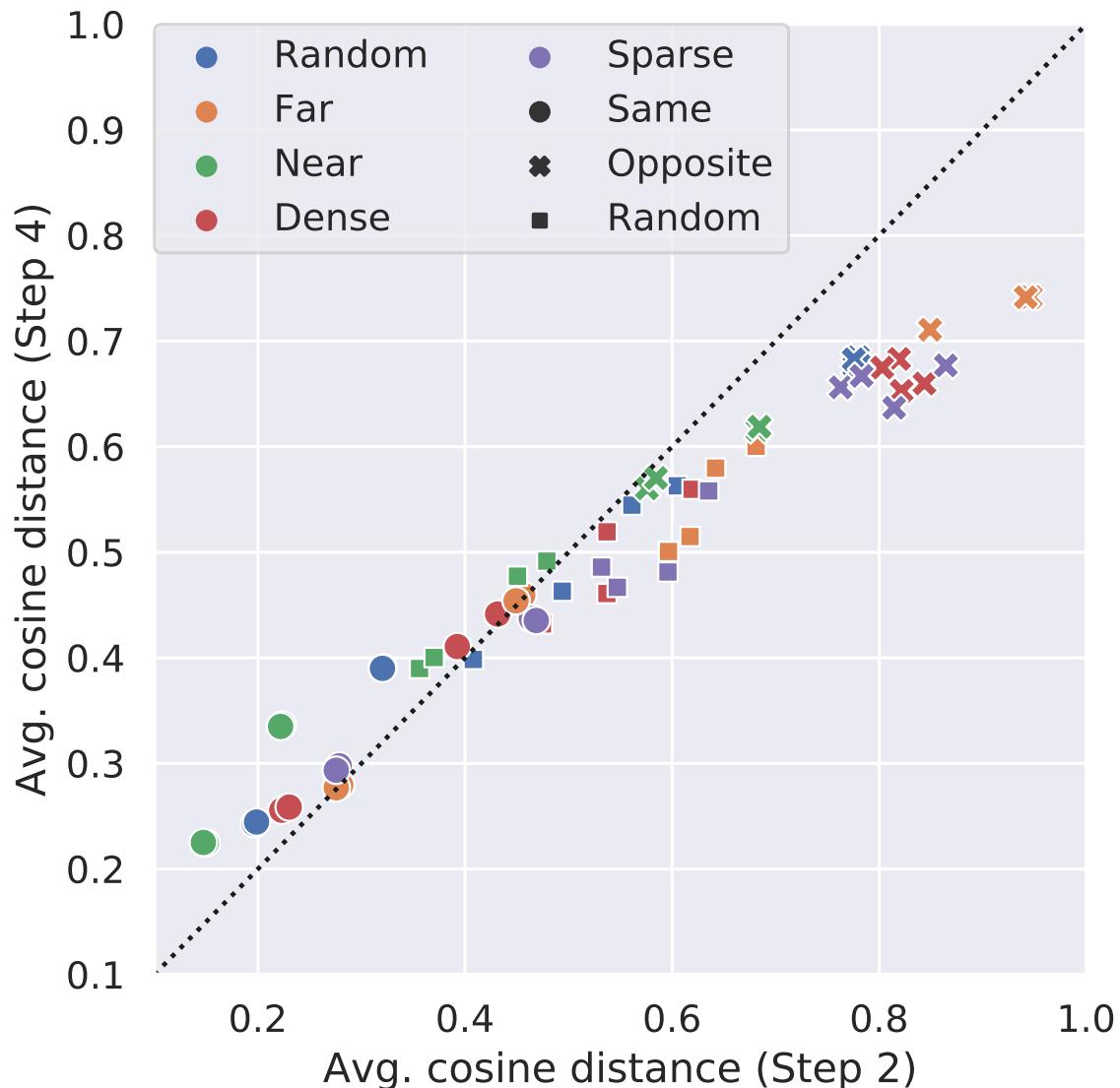


Fig. 5.10 Average cosine distance on the development set between the x-vectors of the original utterance and the anonymized utterance in Step 4 plotted against the average cosine distance between the x-vector of the original utterance and the target x-vector in Step 2. Distance is fixed to PLDA. For each choice of proximity (color) and gender selection (marker shape), four values are shown depending on the gender of the original speaker and the assignment choice.

**Emmanuel:** I'm not convinced by this argument. Unless there are other application constraints, I think we should care about privacy and utility, not about Fig. 6. All methods are comparable in terms of privacy (perhaps 'random' is a little worse than the two others) while 'same' is clearly better in terms of utility (and 'random' is clearly worse). Therefore I believe 'same' is the best method and 'random' is the worst.

1

2 **Assignment** Finally the design choice of *assignment* is examined from the speaker's perspective as de-  
 3 scribed in Section 5.4.5. The distance is fixed to PLDA, proximity to dense and gender-selection to random.  
 4 The results reported in Fig. 5.11 show that *utterance-level* pseudo-speaker assignment results in lower linka-  
 5 bility and also exhibits lower variability across genders as compared to *speaker-level* assignment.

6 The WER resulting from *utterance-level* assignment is higher than from speaker-level assignment (see  
 7 Fig. 5.11). However, in order to conform with the four requirements of the anonymization task  
 8 mentioned in Section 3.2 of the VoicePrivacy Challenge evaluation plan [287], we propose to use *speaker-*  
 9 *level* assignment. This ensures that all utterances from a given original speaker appear to be uttered by the  
 10 same pseudo-speaker.

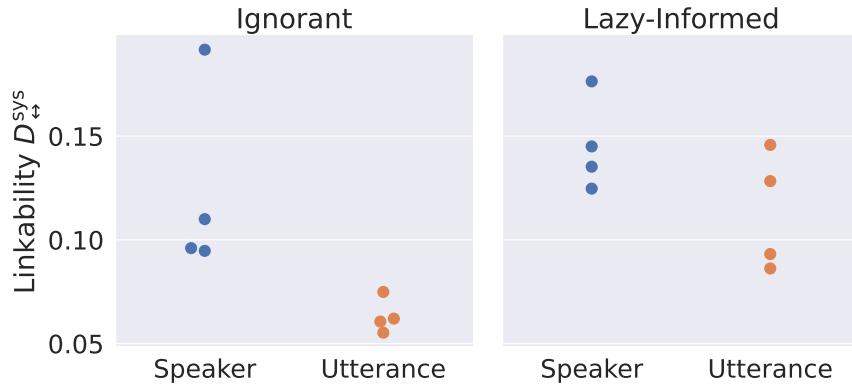


Fig. 5.11 Privacy against *Ignorant* and *Lazy-Informed* attackers depending on the assignment choice and the gender of the original speaker. Distance is fixed to PLDA, proximity to *dense* or *random*, and gender selection to *random*. Each swarm plot shows the 2 linkability values for each gender on the development set resulting from the 2 proximity choices.

11 Based on these indications, the speaker may choose specific parameters according to their application  
 12 needs. For the sake of further experimentation, we choose distance as **PLDA**, proximity as **random** or  
 13 **dense**, gender selection as **random** and assignment as **speaker-level** to be the best combinations of design  
 14 choices based on our observations. The different anonymization techniques can be visually compared in  
 15 Fig. 5.13. The random proximity mapping (Fig. 5.13(a)) causes male and female clusters to move apart  
 16 since the mapping is within the *same* gender selection. It causes compression of speaker space thereby  
 17 inducing ambiguity in speaker's identity, but the highly separated gender clusters may not be desirable  
 18 from anonymization perspective. The *dense* mapping of pseudo-speakers with *random* gender-selection  
 19 (Fig. 5.13(d)) causes speakers to further accumulate in selected clusters and non-separable boundary between  
 20 genders causing lower discrimination between speakers. Such case results in significantly high and robust  
 21 privacy protection.

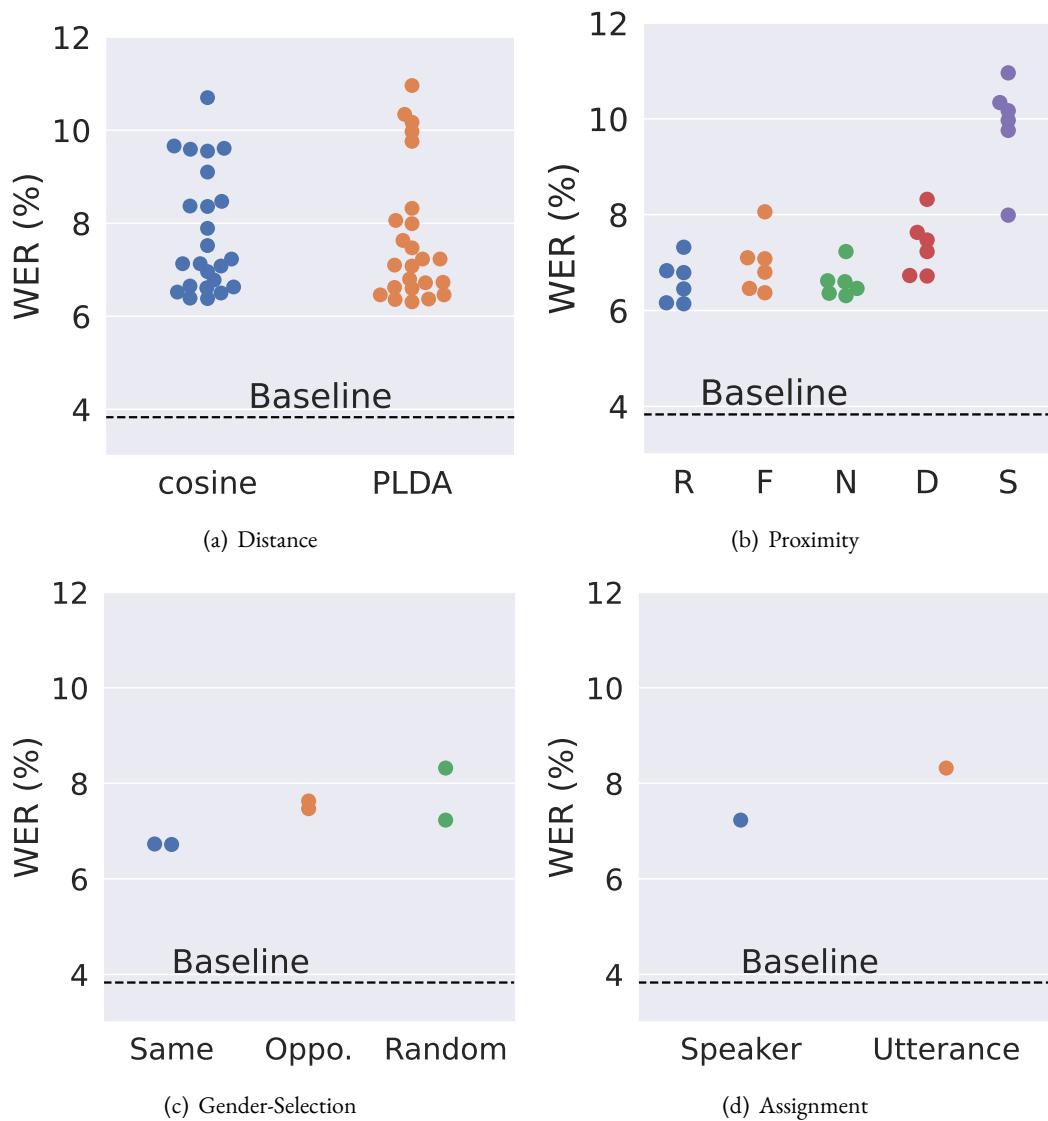


Fig. 5.12 Utility of anonymized speech in terms of WER compared to the original (baseline) speech depending on the different design choices and the gender of the original speaker. Each swarm plot shows the WER values on the development set for each gender and for a given design choice. The remaining design choices are fixed in the same way as in Figs. 5.7, 5.8, 5.9 and 5.11.

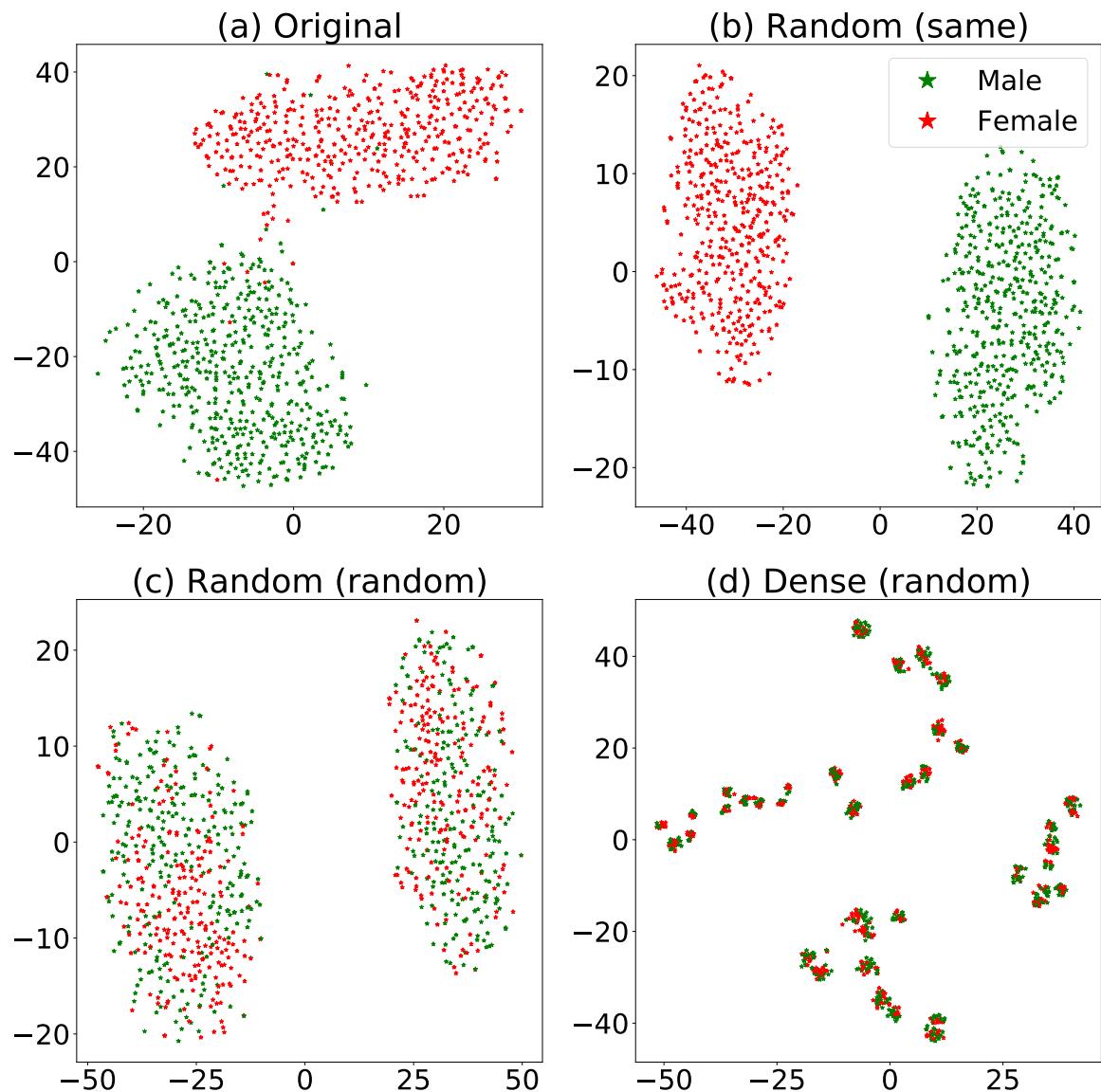


Fig. 5.13 t-SNE visualization of speaker-level x-vectors from the LibriSpeech *train-clean-360* data set transformed using different anonymization schemes.

## 5.4 Design choices in x-vector space

99

**User's perspective**

We now present some complementary results from user's perspective where we measure the feasibility of using the anonymized speech corpus in downstream tasks such as ASR training. Recall that in our threat model, the user is an actor who consumes the anonymized speech corpus for some specific application. The primary concern of any user is the quality of speech in terms of naturalness and intelligibility and its usefulness in downstream tasks. We specifically exhibit the quality of anonymized speech in terms of its viability to train a good  $ASR_{eval}^{anon}$  model. In Section 5.4.7 we will investigate how naturalness and intelligibility can be increased using pitch interpolation techniques.

In order to re-train a similar  $ASR_{eval}^{anon}$  model which is being used hitherto, we anonymize the entire speech corpus which was used for the previous model.

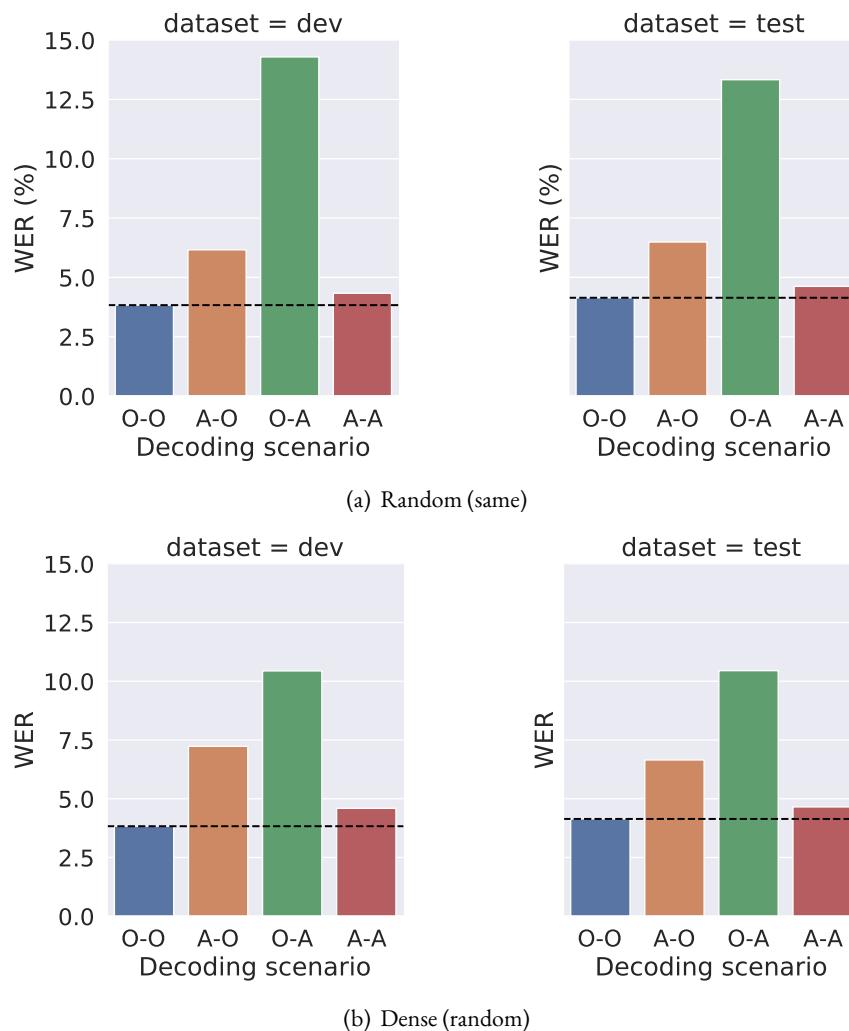


Fig. 5.14 Performance of  $ASR_{eval}^{anon}$  models re-trained using anonymized speech corpus obtained by two methods, namely Random proximity and Dense proximity.

We choose two fairly distinct kinds of anonymization methods to compare their capacity to produce good quality corpus usable for training the  $ASR_{eval}$  model. *Random* anonymization is agnostic to distance metric in x-vector space, while *Dense* anonymization is where pseudo-speakers are selected from a dense

1 cluster in x-vector space. The results for both are marginalised over all the *gender-selection* and *assignment*  
 2 criteria.

3 Figure 5.14 shows that results of the two anonymization methods. The four bars in each plot represent  
 4 the decoding scenarios: O-O indicates original (un-anonymized) speech being decoded by original model  
 5 (trained on un-anonymized speech), A-O indicates that anonymized speech is being decoded by the original  
 6 model, O-A is when original speech is decoded using re-trained model using anonymized speech, while A-A  
 7 is the case where anonymized speech is decoded using the model trained using anonymized speech corpus.

8 We observe that the A-A (red) bar is almost always equal to the O-O (blue) bar indicating that each  
 9 method produces viable speech corpora for training the ASR model. Such models can decode the speech  
 10 with WER as low as the baseline. The middle two bars indicate mismatch between training and decoding  
 11 data. The WER degradation for *random* and *dense* case is much higher when original speech is decoded  
 12 using re-trained model and reasonably lower when anonymized samples are decoded using original model.

13 Such asymmetry indicates “loss of generalization” when ASR is trained using anonymized speech,  
 14 thereby exclusion of certain crucial factors which are present in the original speech. We hypothesize that  
 15 such a loss can be recovered by mixing original and anonymized speech while re-training the ASR model  
 16 so that all the factors may be included. This type of re-training is out of scope for this study and can be  
 17 conducted as part of future experiments.

18 In conclusion, we have shown that the anonymized speech data is suitable for training a viable ASR  
 19 acoustic model with very little loss of generalisation.

## 20 Attacker’s perspective

21 The primary objective of the attacker is to deduce the original speaker’s identity from the anonymized speech,  
 22 i.e. to achieve high Linkability value while conducting speaker authentication trials. We have thus far seen  
 23 two types of attackers, namely *Ignorant* and *Lazy-Informed* who utilize a pre-trained  $ASV_{eval}$  to conduct  
 24 trials. Such a model face an intrinsic limitation of poor discrimination between anonymized speakers in  
 25 x-vector space.

26 In this section we present results with respect to *Semi-Informed* attacker who is aware of the anonymization  
 27 method being used to generate the published corpus. The attacker then uses this method to generate  
 28 the anonymized training data set to be able to better discriminate between anonymized pseudo-speakers in  
 29 x-vector space.

30 We use the same recipe and data set as used by the original  $ASV_{eval}$  model. Similar to Section 5.4.7, we  
 31 use the two diverse methods, namely *Random* and *Dense* to anonymize the training data set.

32 The results for the two methods are shown in Fig. 5.15. We observe that the value of Linkability  
 33 rises gradually as we move from *Ignorant* to *Semi-Informed* attacker in case of both *Random* and *Dense*  
 34 anonymization. It stays below 0.4 for *Random* and below 0.2 for *Dense* even in the strongest scenario. This  
 35 indicates the robustness of *Dense* proximity over *Random*. The results of the two methods are marginalised  
 36 over all the values of *gender-selection* and *assignment*, hence we can further infer from Fig. 5.15(a) and  
 37 Fig. 5.15(c) that the *Dense* method incurs smaller variance of results than *Random*, thereby proving to be  
 38 more stable and superior anonymization.

39 Here we exhibit the robustness of the selected design choices against an attacker with complete knowledge  
 40 about the anonymization method and parameters but only lacks the knowledge of exact pseudo-speaker  
 41 targets. As opposed to signal processing based method, which provides no protection against a strong  
 42 attacker, the design choices selected for x-vector based method are capable of cutting attackers’ Linkability  
 43 by up to half, or even quarter of the baseline value.

## 5.4 Design choices in x-vector space

101

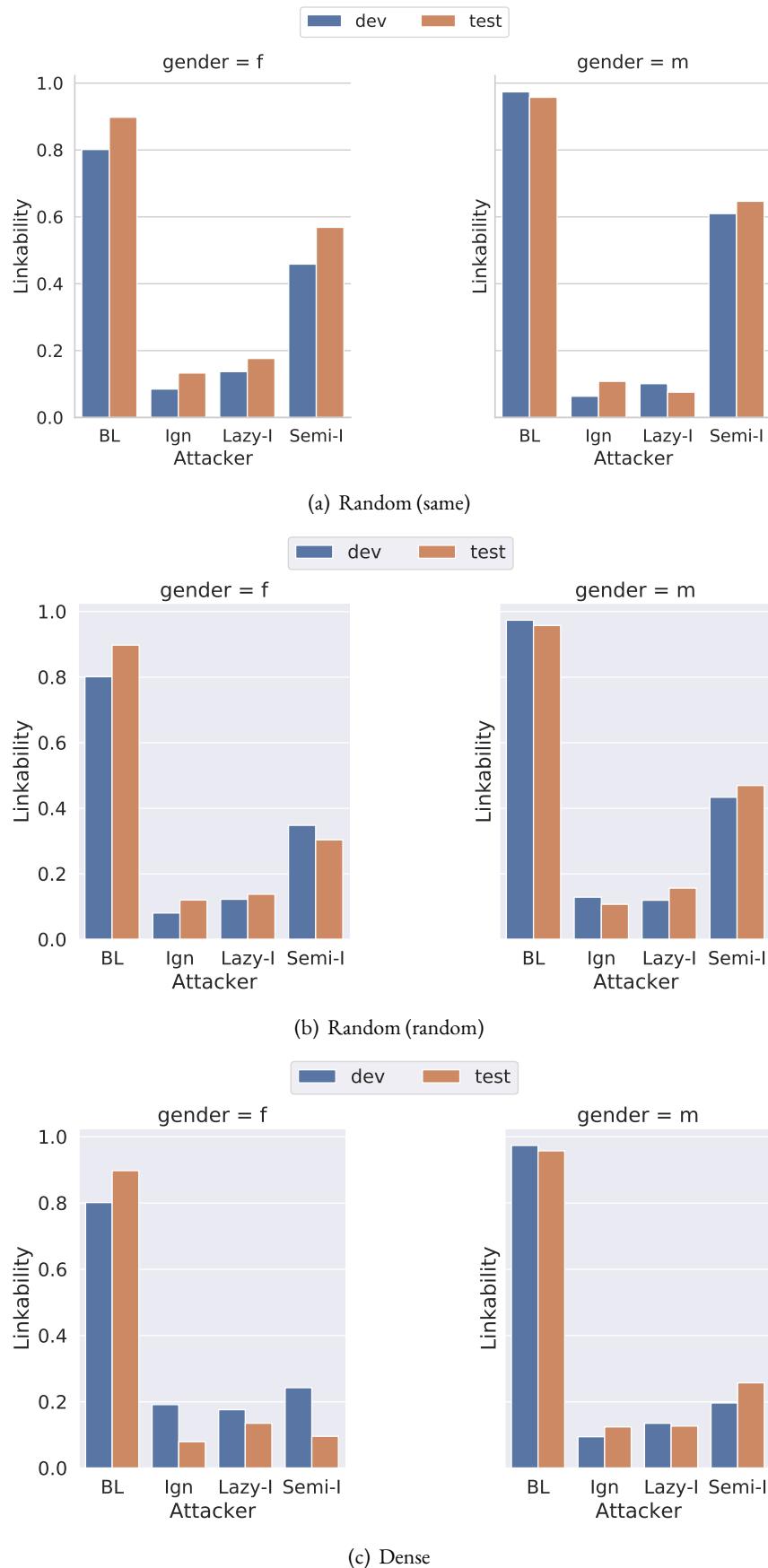


Fig. 5.15 Performance of informed ASV<sub>eval</sub><sup>anon</sup> models re-trained using anonymized speech corpus obtained by two methods, namely Random proximity and Dense proximity. BL = Original (baseline), Ign = *Ignorant*, Lazy-I = *Lazy-Informed* and Semi-I = *Semi-Informed* attacker.

## <sup>1</sup> Pitch transformation

- <sup>2</sup> In this section we investigate whether we can remove residual speaker information in pitch contour and  
<sup>3</sup> possibly achieve better quality speech by transforming the original pitch of the anonymized utterance to the  
<sup>4</sup> pseudo-speaker's pitch.

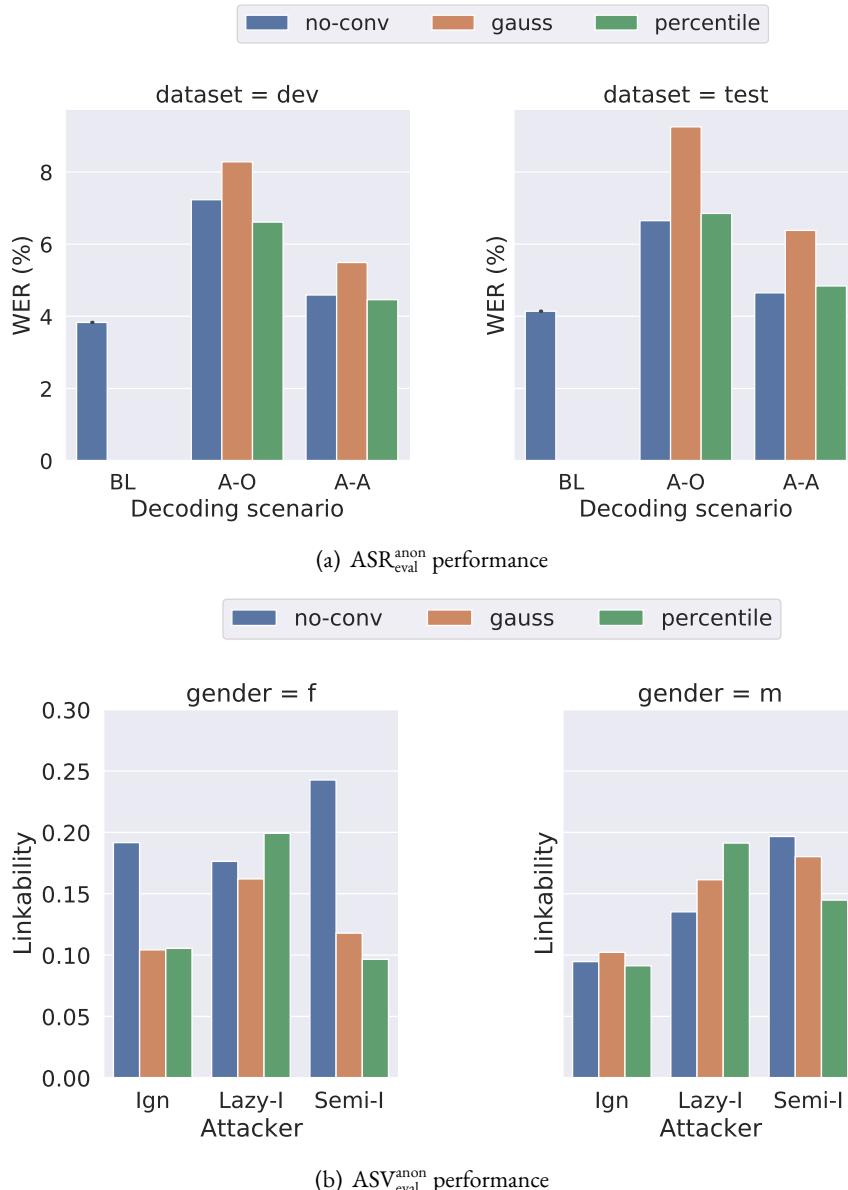


Fig. 5.16 Performance of ASR<sub>eval</sub><sup>anon</sup> and ASV<sub>eval</sub><sup>anon</sup> after the two types of target pitch transformation as compared to original pitch.

- <sup>5</sup> As described in Section 5.4.1, the anonymization framework takes as input the BN features, the original  
<sup>6</sup> x-vector and the pitch contour. Only the original x-vector is transformed into the target pseudo-speaker  
<sup>7</sup> while the other two are left unchanged. Yet, prosodical features of an utterance contribute prominently  
<sup>8</sup> towards the speaker's identity and the presence of original F0 might reveal some information about the

speaker [87]. Hence, we employ two methods of pitch transformation to better conceal the identity as well as to enhance the naturalness of output speech.

The first method is logarithm Gaussian F0 normalisation [176], where the value of original F0 is linearly interpolated to the target F0 in log domain using mean and variance of original and target as given in the following equation:

$$\log(p_{tgt}) = \frac{\log(p_{src}) - \mu_{src}}{\sigma_{src}^2} \sigma_{tgt}^2 + \mu_{tgt}, \quad (5.3)$$

where  $\mu_{src}$ ,  $\sigma_{src}^2$ ,  $p_{src}$  and  $\mu_{tgt}$ ,  $\sigma_{tgt}^2$ ,  $p_{tgt}$  are the mean, variance and pitch of the original and target speaker respectively. This method is also employed by Champion et al. [?] in the same setting to study the effect of F0 modification on privacy and utility for different gender.

Additionally, we propose a second method of transformation which is based on mapping the specific percentile of original pitch to the corresponding percentile of target pitch. The percentile for a specific pitch value is computed as follows:

$$\rho[i] = \frac{p[i] - \min(p)}{\max(p) - \min(p)} 100.0, \quad (5.4)$$

where  $p$  represents the array of pitch values excluding zeros and  $\rho$  is the corresponding percentile value. The advantage of percentile-based mapping is that all the resulting values come from the set of valid pitch values, while in case of Gaussian normalisation, the computed pitch might not be within the valid range of pitch values. To the best of our knowledge, the percentile based pitch transformation is a novel approach and has not been reported in previous literature. Note that the unvoiced regions are not taken into account when performing both transformation methods.

It is observed in Fig. 5.16(a) that Gaussian normalisation significantly increases the WER (thereby causing loss of intelligibility), while percentile-based method maintains the original WER. Figure 5.16(b) shows that both the methods reduce the Linkability significantly, especially in the *Semi-Informed* attack. It can be implied that conversion of pitch removes some of the residual speaker information in the anonymized speech, thereby improving privacy protection. It is worth mentioning that the naturalness of cross-gender voice conversion has also shown noticeable improvement after percentile based pitch transformation as compared to no-transformation when listened in an informal setting.

## 5.5 Large-scale speaker study

In this section we analyse the attacker's performance in presence of increasingly large population of speakers. Specifically, we ask from the attacker's perspective: does the anonymization sustain the desired level of privacy in the presence of thousands of enrollment speakers? If yes, then how is it affected in terms of ASI performance?

We employ Mozilla's Common Voice data set for investigating the effect of speaker population on the attacker's performance. The number of speakers are considered as the attacker's prior knowledge since smaller number of speakers reflect the ability of the attacker to narrow down the original speaker to a smaller number of possible suspects possibly using contextual information. Our main goal is to study if the speaker's identity can be hidden in the crowd or gets revealed when subjected to ASV or ASI in presence of large speaker population. We increase the speaker population exponentially and measure the attacker's performance at each step.

Previous research by Sholokhov et al. [256, 255] studies a similar phenomenon from voice spoofing perspective where an attacker desires to be accepted through an ASV authentication system by finding the "closest impostor" who would be accounted as a false alarm. The attacker has access to a speech sample of a

1 target speaker and the scoring mechanism of the ASV system. They show that the chance of acceptance of  
 2 the impostor may reach up to 50% in the worst case as the population approaches  $10^5$  impostors.

3 Another similar problem is posed by Multi-target speaker detection challenge [257] where membership  
 4 (TOP-S) and identification (TOP-1) of a speaker must be assessed from a large set of blacklisted speakers.  
 5 They show that the performance in both the cases gradually degrade as the number of speakers in the blacklist  
 6 increase. In the scope of this article, we would like to observe the performance of a *Semi-Informed* attacker  
 7 when it has access to the speech samples of large number of speakers. We do not strengthen the attack  
 8 by selecting only the “closest impostor” like [255], and test only the overall linkability of speakers as the  
 9 potential non-mated trials increase multifold.

### 10 5.5.1 Remark on utility

11 We notice that the WER obtained over Common Voice evaluation set using ASR<sub>eval</sub> model increased from  
 12 the baseline 4.71 to 12.64 after anonymization. This looks reasonable given that the ASR<sub>eval</sub> is trained using  
 13 LibriSpeech training set.

### 14 5.5.2 Privacy evaluation metrics

15 To assess the performance of a re-identification attack, we first report the results in terms of binary decision  
 16 metrics such as EER and Linkability as reported in previous sections. Thereafter we perform closet-set ASI  
 17 and report metrics such as *Rank* and *top-k* membership performance. We compute the PLDA score for a  
 18 particular speaker with all the speakers in the selected population (which usually contains the true speaker)  
 19 and sort it in descending order. Ideally the true speaker must attain highest score and appear as the first  
 20 rank, but in a more private setting the given speaker will not resemble like itself and the rank of the true  
 21 speaker will increase manyfolds. We report the rank of the true speaker along with the top- $k$  precision which  
 22 indicates the binary presence of the true speaker within initial  $k$  ranks. The rank and top- $k$  help us assess the  
 23 relative difficulty of ASI faced by different attackers in the presence of large number of speakers. Adding  
 24 more speakers to the population might certainly increase the rank of true speaker, hence we also report  
 25 normalized rank, which is the absolute rank divided by the number of speakers in the population, to assess  
 26 the bona fide improvement in privacy as a function of number of speakers.

27 We plot the *chance-level* rank as the reference of absolute privacy which can be defined as the expectation  
 28 of rank for a speaker when the attacker is selecting a rank at random among all the speakers with no knowledge  
 29 of anonymization. Such a rank is given by  $\mathbb{E}(R) = \frac{N_{\text{spk}}+1}{2}$ , where  $N_{\text{spk}}$  is the total number of speakers in  
 30 the population. Similarly we define the normalized *chance-level* rank as  $\frac{N_{\text{spk}}+1}{2N_{\text{spk}}} \approx 0.5$ . The derivation for  
 31 chance-level rank is given as follows. Let  $R_i \in \{1, \dots, N_{\text{spk}}\}$  be the set of all possible ranks for speaker  $i$   
 32 that can be obtained during identification with probability  $P(R_i)$ . Then the expected rank is given by:

$$33 \quad \mathbb{E}(R) = \sum_i R_i P(R_i) \quad (5.5)$$

34 To obtain the chance-level rank, we set the value of  $P(R_i) = \frac{1}{n}$  and substitute the value of  $R_i$ . Hence:

$$\begin{aligned}
 \mathbb{E}(R) &= \frac{1}{N_{\text{spk}}} \sum_i R_i \\
 &= \frac{1}{N_{\text{spk}}} \frac{N_{\text{spk}}(N_{\text{spk}} + 1)}{2} \\
 &= \frac{N_{\text{spk}} + 1}{2}
 \end{aligned} \tag{5.6} \quad 1$$

When the rank is normalized, we can divide the chance-level rank by  $N_{\text{spk}}$ , therefore the expected normalized rank is given by:

$$\mathbb{E}(R_{\text{norm}}) = \frac{N_{\text{spk}} + 1}{2N_{\text{spk}}} \approx 0.5 \tag{5.7} \quad 4$$

### 5.5.3 Gender identification

We observed that a large number of speakers in the common voice data set did not specify their gender, which is crucial for conducting trials pertaining to the experiments described in this section. Before training a gender identification method over the training set composed of speakers who specified their gender, we visualized the x-vector space of the training set as presented in Fig. 5.17. We notice that there is a significant overlap between the male and female clusters which is uncommon in x-vector space as observed in other data sets. After manually listening to some of the outlier audio samples we discovered that this overlap is due to large amount of *children speakers* and *gender mislabelling* in the data set. Previous studies [220, 232] have shown that it is especially challenging to identify gender in presence of children’s voice in the speech data set.

We also observe that the labelled part of the data set consists of large number of male speakers (22.1%) and only a small amount of female speakers (5.45%). Such unbalanced data set would result in huge bias against female speakers. We employ the gender identification technique suggested by Kanervisto et al. [142], where the x-vectors are first projected into 1-D space using linear discriminant analysis (LDA) and then clustered using Gaussian mixture models (GMM). We obtain the F-1 score of 92% for female speakers and 91% for male speakers in the test set. The speaker gender distribution before and after the gender identification is presented in Fig. 5.18.

We notice that out of 72.46% speakers with unspecified gender, 48.3% speakers are predicted to be male and only 24.16% as female. Moreover, many children speakers are classified as female hence the number of female speakers is further reduced. Since there are large number of confirmed male speakers in the overall data set, we chose to conduct the large-scale speaker study only with the male speakers.

### 5.5.4 Experimental setup

We select 24,616 speakers as the total population possessed by the attacker from the set of male speakers whose total duration is more than 10 seconds after removing silent frames using voice activity detection (VAD). The maximum duration for each speaker in the population is limited to 2 minutes. Only the utterances with Signal-to-Noise ratio (SNR) more than 75 dB are selected. We compute the SNR using WADA-SNR [149] Algorithm.<sup>14</sup> Out of these, 20 speakers whose total duration is greater than 5 minutes, are selected for testing which represents the publicly released data subjected to re-identification attack. We manually listened to each speaker in the test set to confirm that it is a distinct male speaker due to low gender

<sup>14</sup><https://gist.github.com/johnmeade/d8d2c67b87cda95cd253f55c21387e75>

TSNE for 52973 speakers in cv51-en-all. One vector per speaker.

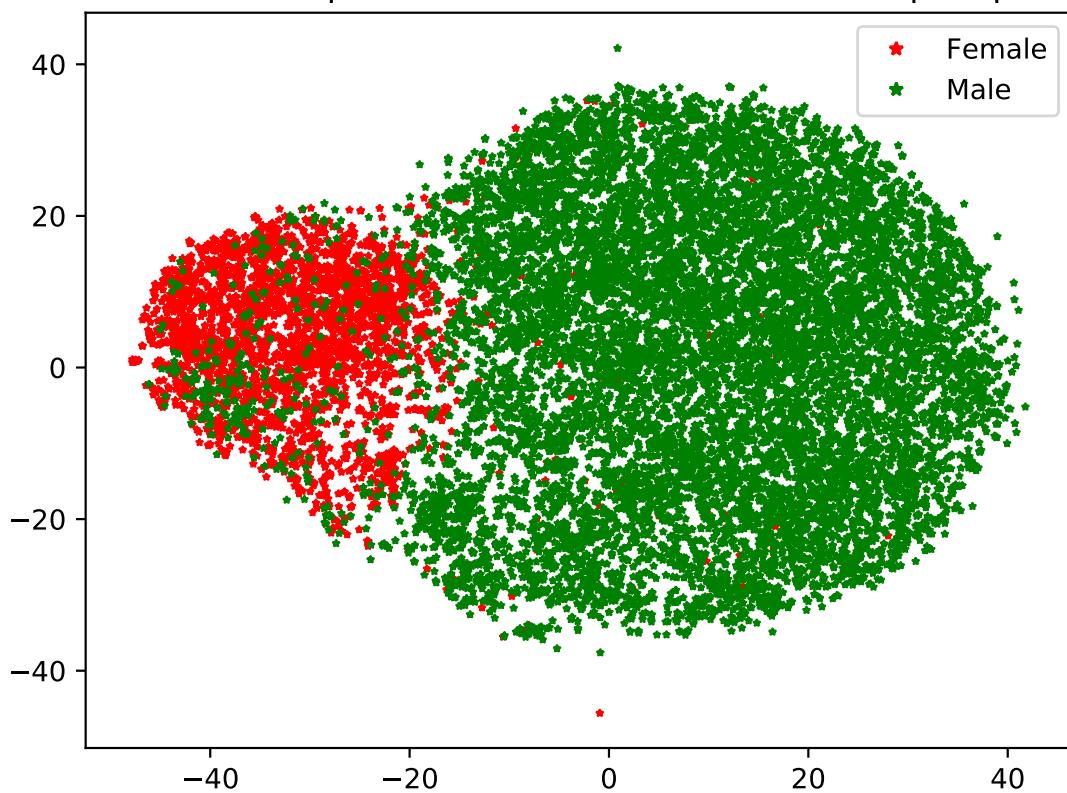


Fig. 5.17 t-SNE representation of speaker x-vectors in common voice data set.

1 identification performance on this data set as shown previously. After computing PLDA scores between the  
 2 speaker population and the test speakers, we get 4696 same-speaker scores and 115,563,864 different-speaker  
 3 scores.

4 Initially we select the subset of scores which are computed only among 20 speakers present in the test set.  
 5 Thereafter we double the speakers in the population at each step and include the scores corresponding to  
 6 these speakers. The newly added speakers are randomly sampled 5 times from the entire speaker population  
 7 to ascertain if the sampling process induces a bias due to proximity with speakers.

### 8 5.5.5 Experiments and results

9 Let us first discuss the ASV measures as shown in Fig. 5.19. We notice in Fig. 5.19(a) that the baseline EER  
 10 starts with a value of 7% and slightly increases to converge at 12% as we increase the number of speakers.  
 11 The three attackers perform much poorly as we notice that *Ignorant* and *Lazy-Informed* cases start from  
 12 37% EER and quickly ascend to converge at 50%. The *Semi-Informed* attacker performs in a consistent  
 13 manner even after addition of large number of speakers and the EER converges below the *Ignorant* and  
 14 *Lazy-Informed* attackers. This observation might indicate the precise measurement of EER when large  
 15 number of scores are present. Hence the previous experiments can be considered as the lower bound of  
 16 privacy while the privacy improves with more speakers in the population (see Fig. 5.20). The Linkability  
 17 performance is consistent with the observations in EER except that the *Semi-Informed* attacker displays  
 18 lowest value at all steps.

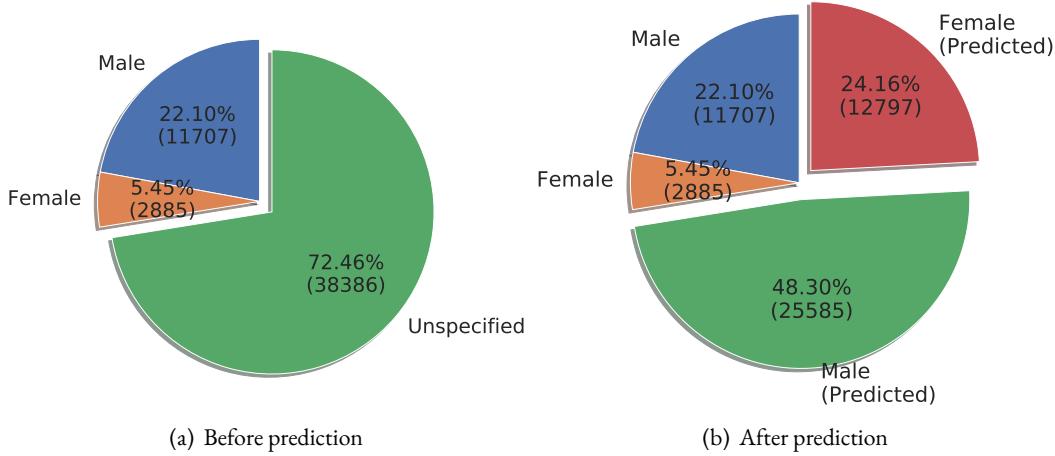


Fig. 5.18 Speaker gender distribution observed in the common voice data set. The exact number of speakers for each gender are indicated in the brackets. Four speakers are discarded due to lack of data.

Figure 5.20 indicates the un-normalized ( $\log_{10}(Rank)$ ) and normalized rank of true speaker obtained by different attackers before and after anonymization. In Fig. 5.20(a), we notice a steep rise in the value of absolute rank, i.e. decline in the ASI performance, in original as well as anonymized case. However the baseline performance is well below the chance level rank even when there are thousands of speakers in the enrollment set indicating the distinctive characteristics of speakers in the population. All the attackers start with better performance than chance-level but soon converge very close to the chance-level rank as the speakers in the population increase. We also plot the normalized rank and the chance-level normalized rank in Fig. 5.20(b). We observe that this plot resembles the EER performance depicted in Fig. 5.19(a). The ASI performance obtained by *Ignorant* and *Lazy-Informed* attacker quickly degrades and converge to a value worse than the chance-level normalized rank, while the *Semi-Informed* attacker maintains a consistent performance, which is better than chance-level. Although this performance is not significantly better than chance-level performance, it warns the speaker that more information to attacker might pose a re-identification threat.

We further study the top- $k$  precision obtained by different attackers and compare them to their corresponding baseline performance. In Fig. 5.21 we only show the results for top-20 precision which is a binary measure indicating the presence of a particular same-speaker score within top 20 of the sorted rank list averaged over all the test speakers in the population. This value is always within the range [0, 1], where 0 indicates that no test speaker could be matched with their own utterance within 20 ranks and 1 indicates that all the speakers match their own utterance within 20 ranks. It exhibits the general observation noticed in all the top- $k$  results and realistically an attacker might look at top-20 results when it wants to shortlist the probable speaker identities. We observe that the precision drops much faster after anonymization as compared to the baseline i.e. no anonymization. This indicates that after anonymization it is much difficult for an attacker to identify a speaker as compared to the baseline. In other words, hiding the identity of an *anonymized* speaker in a crowd of  $n$  speakers is equivalent to hiding the *original* speaker in a crowd of  $N$  speakers, and  $N$  increases with much faster rate as compared to  $n$ .

Concretely, we can infer from Fig. 5.21 that while the attacker's ability to re-identify the speaker naturally reduces with the number of candidate speakers, the best instance of our anonymization scheme with 50

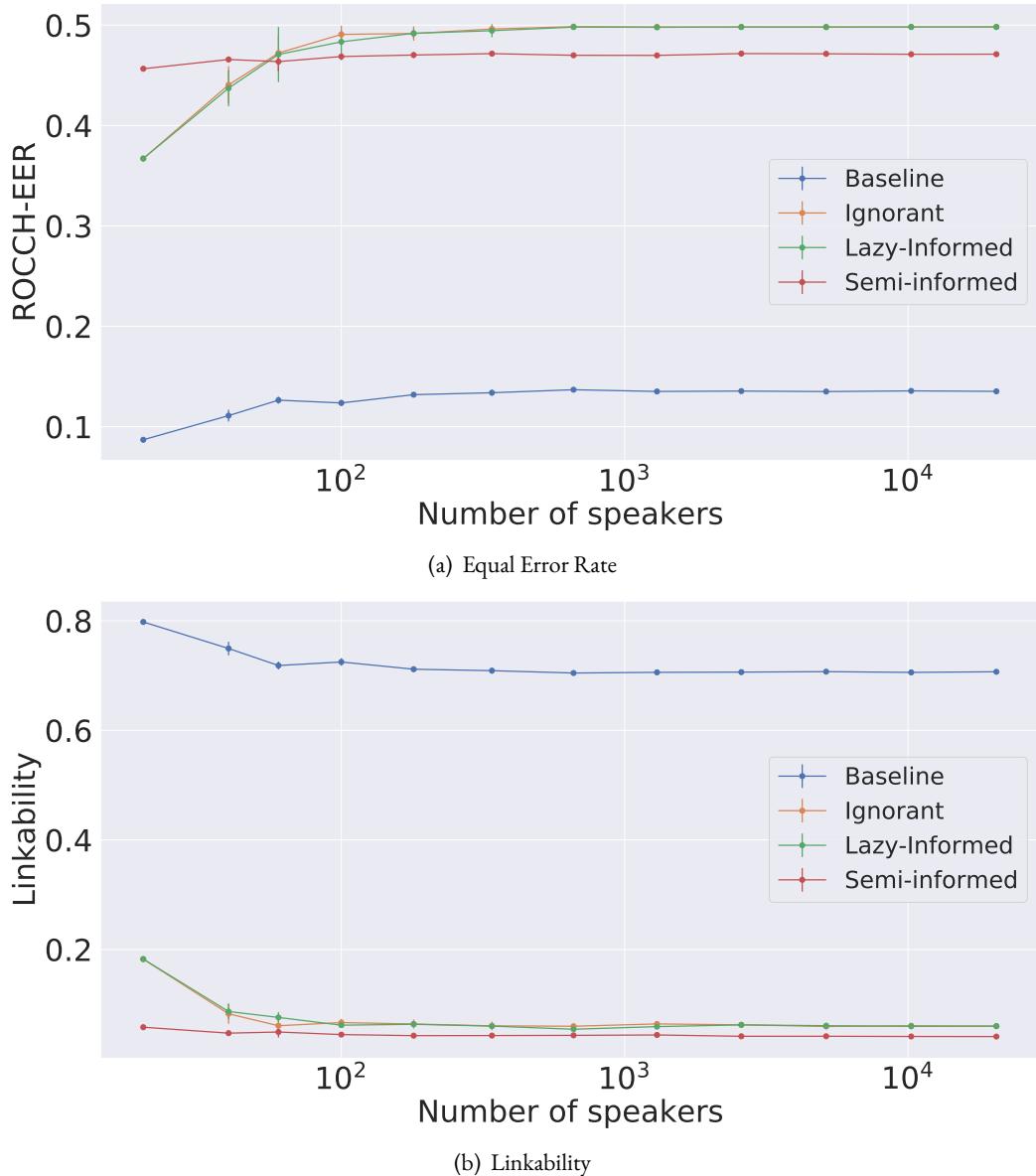


Fig. 5.19 Open-set ASV performance of different attackers in terms of EER and Linkability as the speakers in the population double at each step.

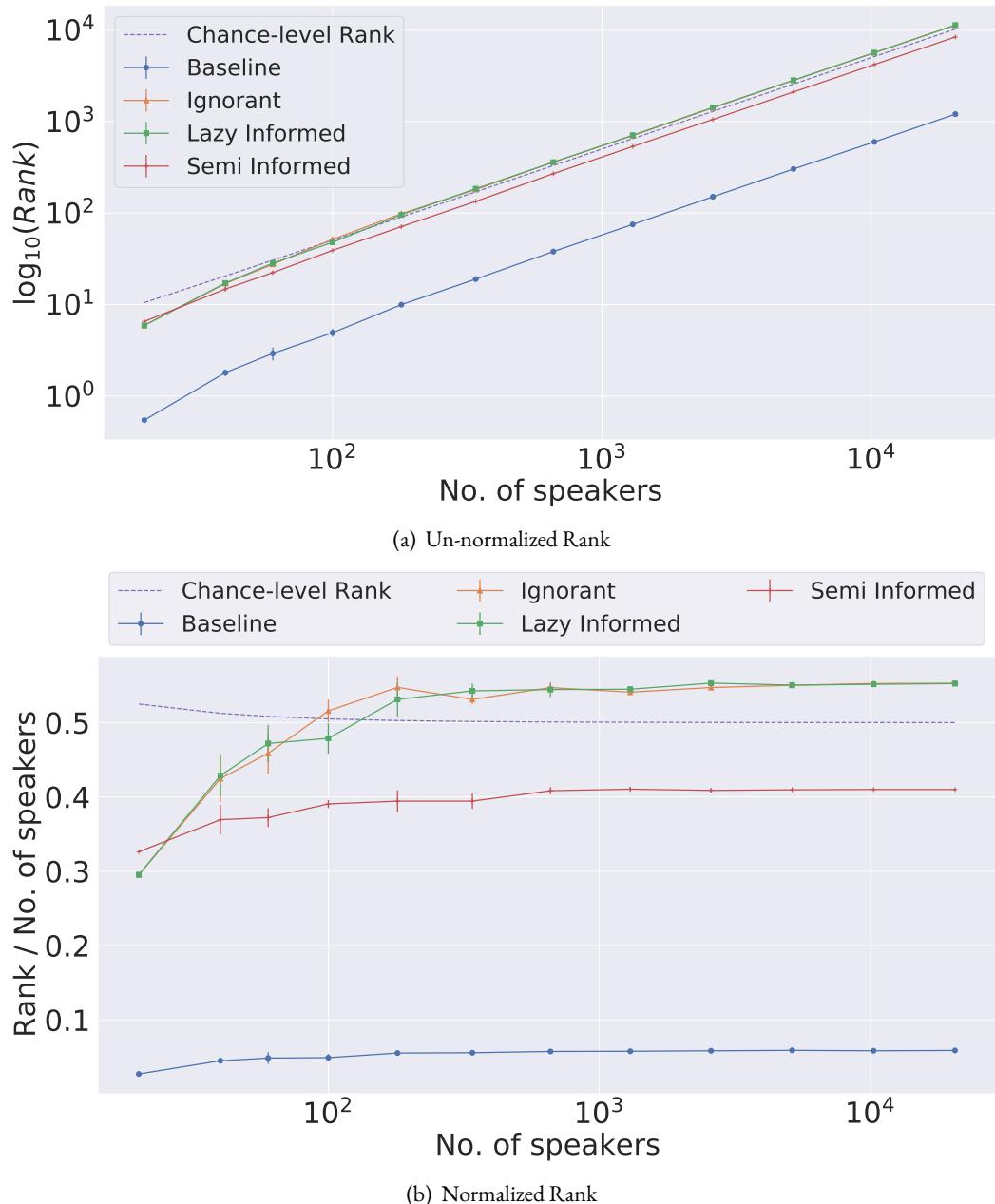


Fig. 5.20 Closed-set ASI performance in terms of un-normalized and normalized rank obtained by different attackers as the number of speakers in the population double at each step.

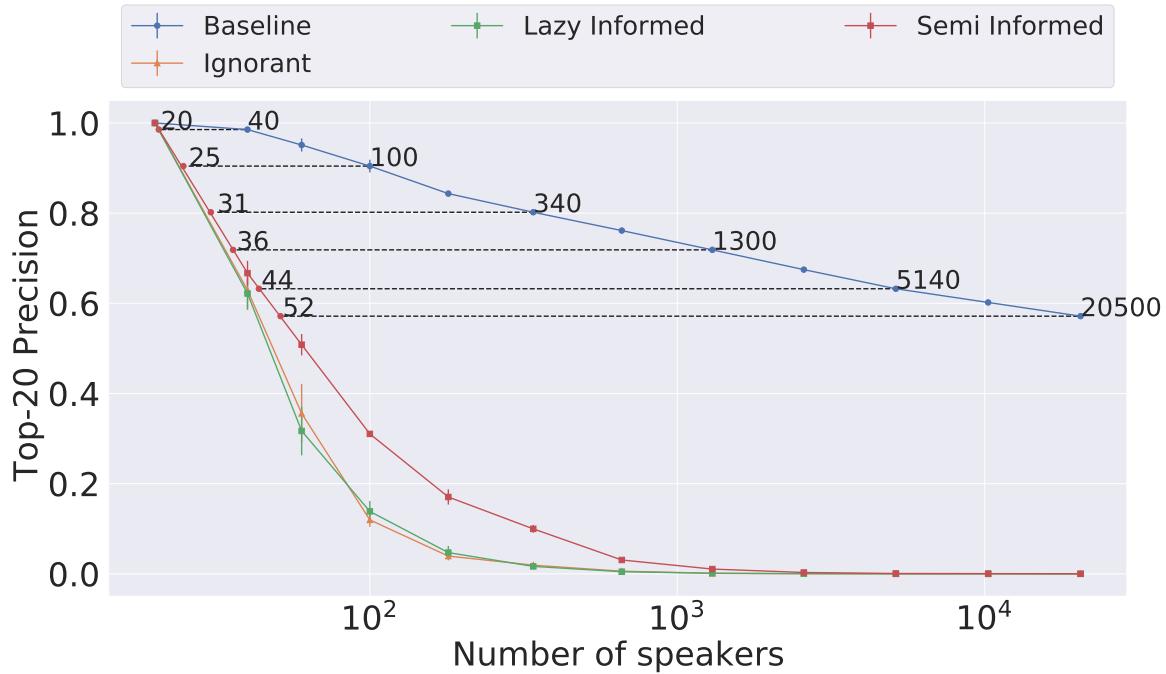


Fig. 5.21 Top-20 precision of ASI for different attackers as the speaker population is doubled at each step. The number of speakers needed before anonymization ( $N$  on blue curve) and after anonymization ( $n$  on red curve) to achieve equivalent drop in precision are highlighted.

<sup>1</sup> candidates speakers guarantees the same anonymization level as raw speech with 20,000 speakers. Refer  
<sup>2</sup> Appendix A for top-1, top-10, and top-50 precision plots.

### <sup>3</sup> 5.5.6 Worst-case analysis

<sup>4</sup> In the previous section, especially in Fig. 5.20(b) we observed the normalized rank of original and anonymized  
<sup>5</sup> utterances as the enrollment speaker population increases exponentially. A lower normalized rank escalates  
<sup>6</sup> the threat of re-identification, thereby compromising the privacy of individuals participating in the data  
<sup>7</sup> collection. To assure optimal protection, data publishers would be interested to know the lower bound of  
<sup>8</sup> privacy achieved by an anonymization scheme and the properties of maximum-risk individuals to prevent  
<sup>9</sup> such conditions. To this end, we conduct the worst-case analysis of the results obtained in the previous  
<sup>10</sup> section over the Mozilla CommonVoice data set. Since Fig. 5.20(b) shows the average normalized rank of all  
<sup>11</sup> the utterances of the selected 20 speakers across all the trials at a given step on the x-axis which corresponds to  
<sup>12</sup> the number of speakers in the enrollment set. This aggregate normalized rank does not answer the questions  
<sup>13</sup> related to the worst-case, which are: (1) *What is the distribution of normalized ranks in case of Semi-Informed*  
<sup>14</sup> *attack as compared to the original speech?* (2) *What are the properties of the worst-performing utterance in the*  
<sup>15</sup> *data set, the worst-performing speaker, and the worst-performing utterance for each speaker?*

<sup>16</sup> To answer the questions above, we analyze the results in terms of the overall normalized rank of the  
<sup>17</sup> worst-performing utterance,  $U^{worst}$ , the normalized rank of the utterance set of the worst-performing speaker,  
<sup>18</sup>  $S^{worst}$ , and the normalized rank of the set of single worst-performing utterance from each speaker,  $U_S^{worst}$ .  
<sup>19</sup> We only consider the ranks obtained using the original speech (baseline) and the *Semi-Informed* attacker  
<sup>20</sup> here. The remaining results with *Ignorant* and *Lazy-Informed* attackers can be found in Appendix B.

**Normalized rank distribution** First, we plot the distribution of normalized ranks of all the trials in the baseline case in Figure 5.22, where each subplot represents the number of speakers in the enrollment set. As expected, it is observed that a majority of utterances exhibit the normalized rank very close to zero and hence are extremely vulnerable to re-identification attacks even in the presence of thousands of enrollment speakers. The value of  $U^{\text{worst}}$ ,  $\bar{S}^{\text{worst}}$ , i.e., the mean of set  $S^{\text{worst}}$ , and  $\bar{U}_S^{\text{worst}}$ , i.e., the mean of set  $U_S^{\text{worst}}$ , represented by the red, green and black dashed lines, respectively, are all close to zero, in fact, they are superimposed upon each other. The goal of anonymization is to increase the value of these variables and shift the utterance density to higher ranks so that the lower bound on privacy protection can be raised.

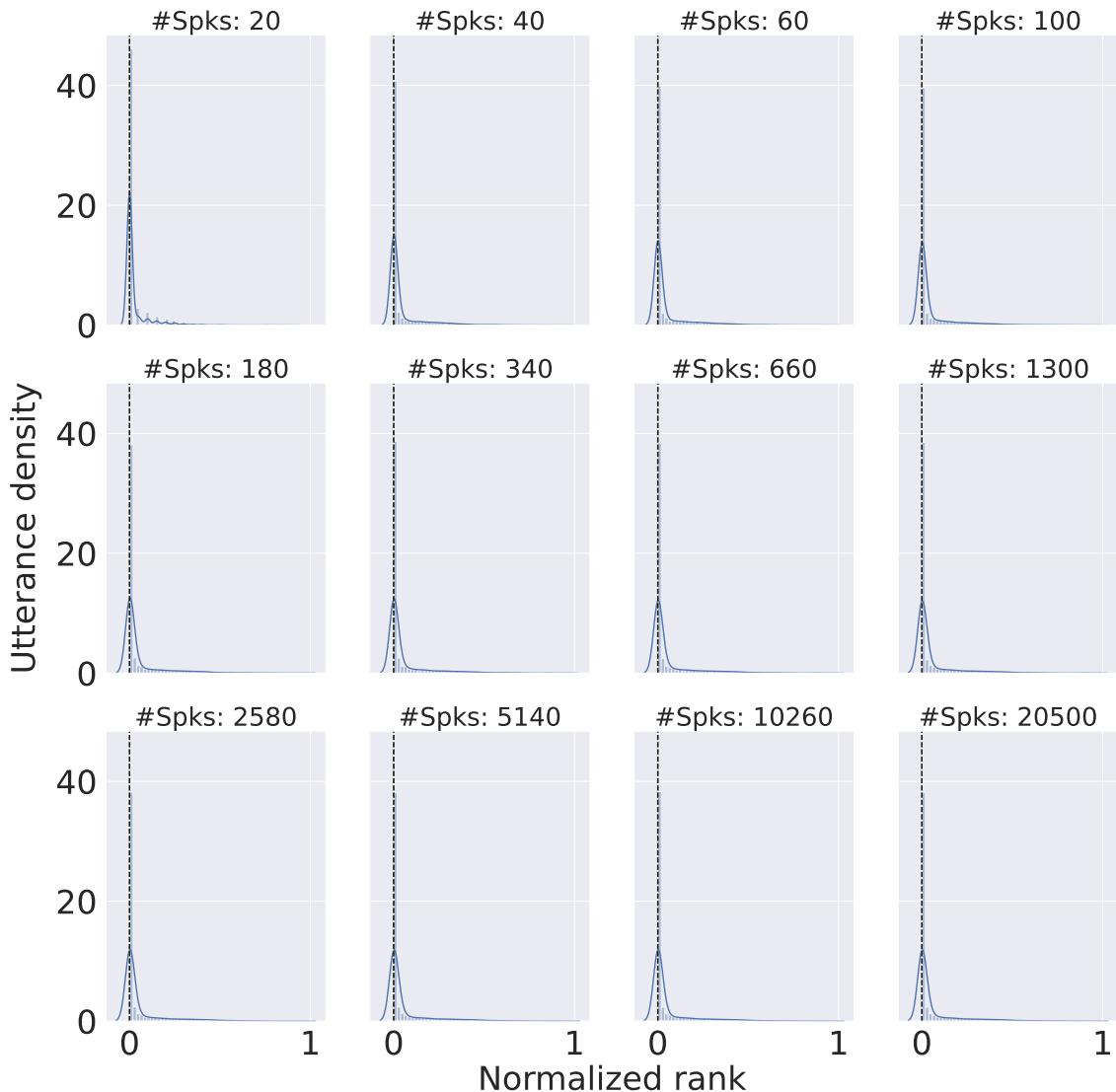


Fig. 5.22 The normalized rank distribution for the baseline case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of  $U^{\text{worst}}$  (red),  $\bar{S}^{\text{worst}}$  (green),  $\bar{U}_S^{\text{worst}}$  (black).

1      Compare Figure 5.22 with Figure 5.23, which shows the normalized rank distribution in the *Semi-*  
 2 *Informed* setting. Although several utterances still tend to exhibit lower ranks, the utterance density becomes  
 3 more evenly distributed over the whole range of normalized rank, thereby protecting several utterances  
 4 from re-identification attacks. Moreover,  $U^{\text{worst}}$  is much worse than  $\bar{S}^{\text{worst}}$  and  $\bar{U}_S^{\text{worst}}$  indicating that not  
 5 all utterances of a single speaker are vulnerable to a strong attack. The gap between  $\bar{S}^{\text{worst}}$  and  $\bar{U}_S^{\text{worst}}$   
 6 indicates that a majority of speakers have their worst-performing utterances better protected than the overall  
 7 worst-performing utterance and the worst-performing speaker.

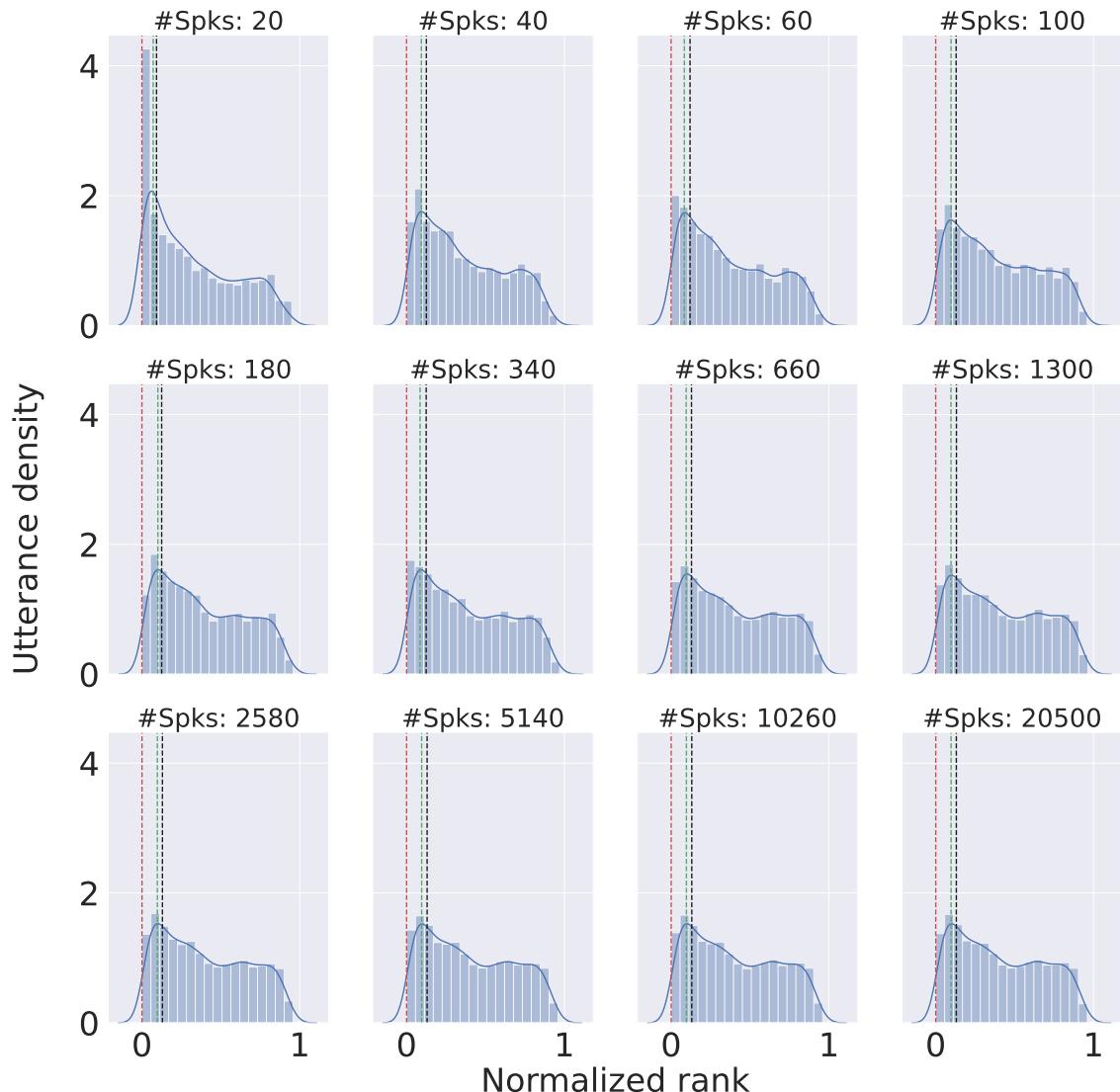


Fig. 5.23 The normalized rank distribution for the *Semi-Informed* case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of  $U^{\text{worst}}$  (red),  $\bar{S}^{\text{worst}}$  (green),  $\bar{U}_S^{\text{worst}}$  (black).

**Duration-specific normalized rank** Next, we investigate if there is a relationship between the duration of an utterance and its vulnerability to attacks. Figure 5.24 shows that for original speech, shorter duration utterances exhibit higher normalized rank, hence they are more well-protected than the longer ones. Similar to Figure ??, the value of  $U^{\text{worst}}$ ,  $\bar{S}^{\text{worst}}$  and  $\bar{U}_S^{\text{worst}}$  are all close to zero.

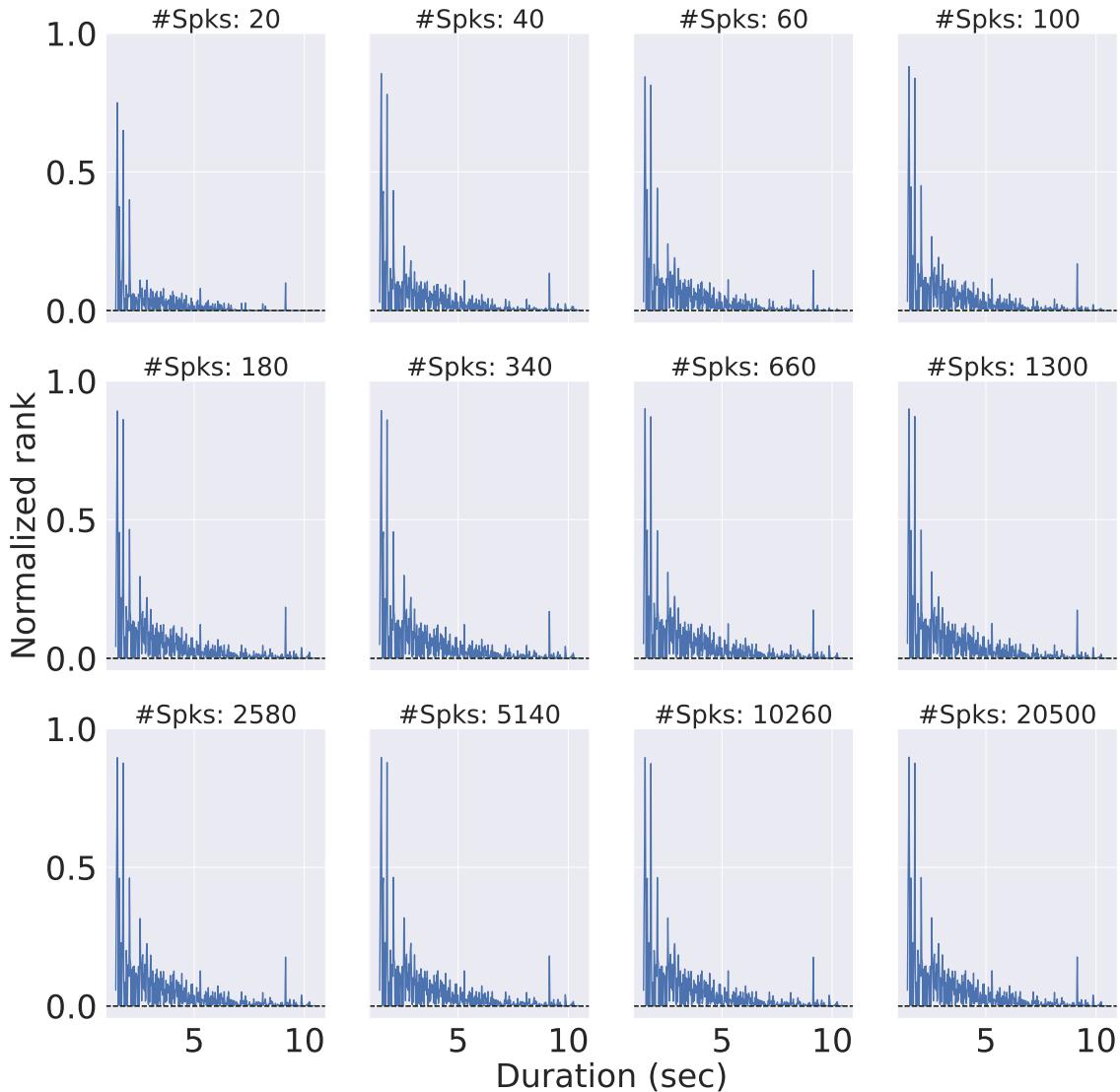


Fig. 5.24 Line plot in original case showing the normalized rank (y-axis) against the duration of utterances (x-axis). The dashed horizontal lines show the value of  $U^{\text{worst}}$  (red),  $\bar{S}^{\text{worst}}$  (green),  $\bar{U}_S^{\text{worst}}$  (black).

Contrary to the above results, Figure 5.25 shows that for the anonymized utterances in *Semi-Informed* case, the privacy protection gradually ascends as the duration increases till approximately 7-8 seconds, and then starts to descend. Such an observation is a crucial indicator for a data publisher to choose an optimal duration for the utterances undergoing anonymization.

Refer to Figures B.3, B.4 and B.5, which show the trend of normalized ranks in all the attack conditions for  $U^{\text{worst}}$ ,  $\bar{S}^{\text{worst}}$  and  $\bar{U}_S^{\text{worst}}$ , respectively. The large standard deviation of privacy protection exhibited by the utterances in *Semi-Informed* case corroborate the conclusions drawn from the duration plots.

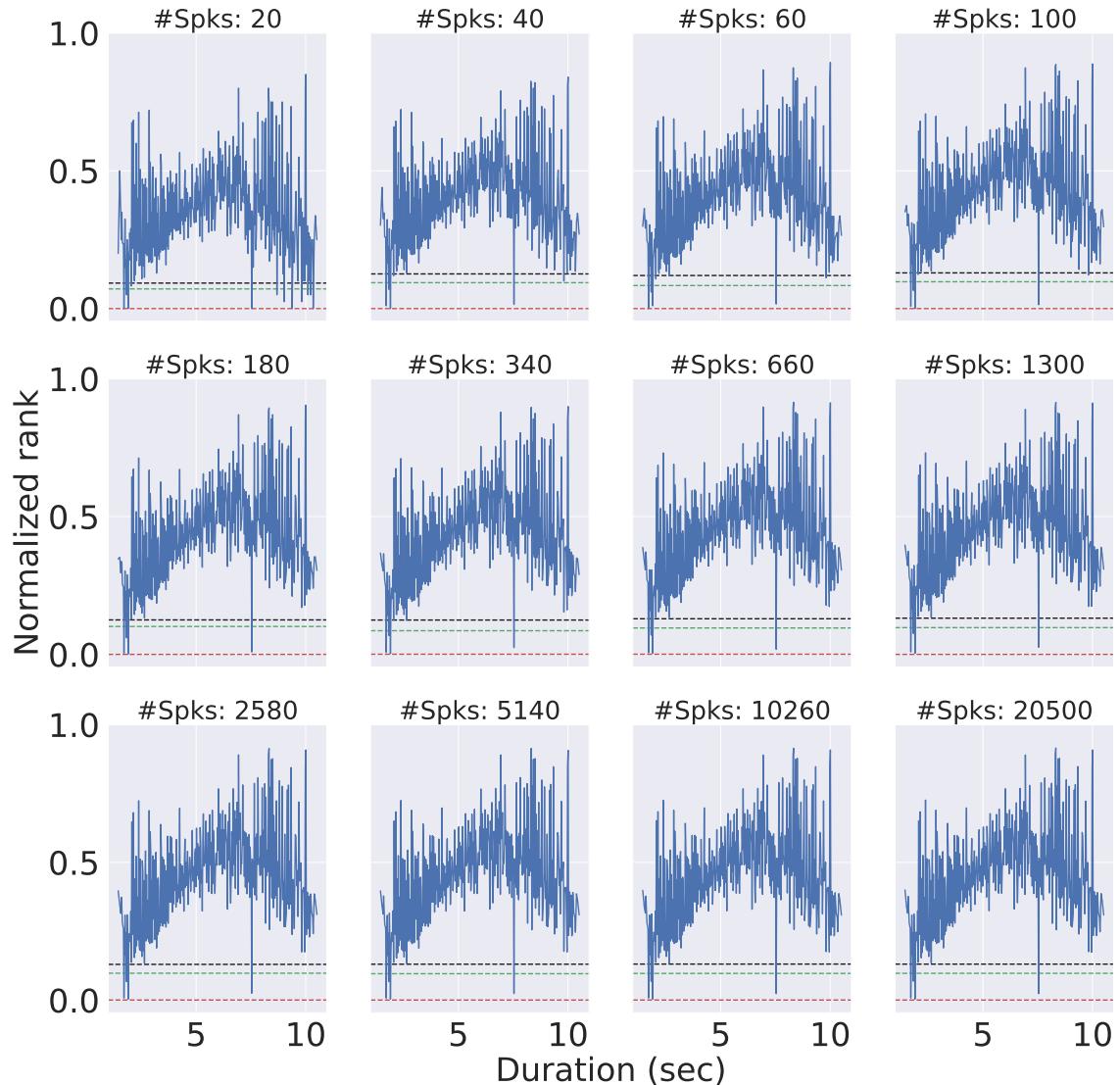


Fig. 5.25 Line plot in *Semi-Informed* case showing the normalized rank (y-axis) against the duration of utterances (x-axis). The dashed horizontal lines show the value of  $U^{\text{worst}}$  (red),  $\bar{S}^{\text{worst}}$  (green),  $\bar{U}_S^{\text{worst}}$  (black).

## 5.6 Summary

We present a detailed comparison between four design choices for x-vector based speaker anonymization. The effect of each design choice was studied with respect to privacy and utility metrics of the anonymized speech and the optimal combination of choices was recommended. Experiments showed that the anonymized speech corpus is suitable to train a viable ASR model and that reasonable amount of privacy protection is achieved even if a *Semi-Informed* attacker endeavors to re-train the ASV model with anonymized speech corpus.

We further investigated two pitch transformation methods to remove the residual speaker information present in the pitch contour of the anonymized speech and possibly to enhance the naturalness of the synthesized voice. The newly proposed percentile based pitch transformation method outperforms the conventional Gaussian normalization method in terms of privacy as well as utility.

We assessed the robustness of the proposed anonymization method assuming that the attacker possesses the data for large number of speakers. We conducted closed-set ASI by incrementally adding thousands of speakers in the population and observed that the rank of the true speaker quickly ascends and converges close to chance-level performance after anonymization. Another interesting observation can be noted using the results of top- $k$  membership analysis where the loss of precision before anonymization that is seen after adding thousands of speaker in the enrollment set is equivalent to adding only a couple of speakers after anonymization.

Finally, we performed the worst-case privacy protection assessment of the proposed anonymization scheme and conclude that a majority of utterances are well protected after anonymization, raising the lower bound of privacy significantly. Additionally, the data publishers can expect maximum protection if they choose an optimal duration of speech signal for anonymization.

A fundamental assumption that we made in this chapter was that, out of the three features extracted from the speech signal, i.e., pitch, BN features, and x-vector, speaker-related information is concentrated only in the x-vector and replacing it with a new pseudo-speaker will delete all the identity markers of the speaker in the synthesized speech. We investigate this assumption in Chapter 6, and propose techniques to remove the residual speakers' identity from the pitch and BN features.

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## Chapter 6

# Towards removing residual speaker information and provable guarantee

There is nothing like looking, if you want to find something.

---

*J.R.R. Tolkien*

In the last chapter, we introduced the state-of-the-art anonymization scheme which was used as the primary baseline for the first VoicePrivacy challenge. It extracts the pitch, the BN features, and the x-vector from the source utterance, replaces the x-vector by a random target pseudo-speaker and re-synthesizes the speech signal using the new x-vector and the original pitch and BN features. This method is based on the assumption that there is perfect disentanglement of linguistic, prosodic and speaker-related attributes in speech, thereby all the speaker-related information is concentrated in the x-vector. Furthermore, investigations in the last chapter have shown that anonymization provides brittle privacy protection, even less so any provable guarantee.

In this chapter, we show that disentanglement is indeed not perfect: linguistic and prosodic attributes still contain speaker information. We remove speaker information from these attributes by introducing differentially private extractors based on an autoencoder and an automatic speech recognizer, respectively, trained using noise layers. We plug these extractors in our anonymization pipeline and generate, for the first time, differentially private utterances with a provable upper bound on the speaker information they contain. We evaluate empirically the privacy and utility resulting from our differentially private anonymization scheme on the LibriSpeech data set. Experimental results show that the generated utterances are intelligible while protected against strong attackers who have full knowledge of the anonymization process.

### 6.1 Individual impact on privacy and utility

comparative study of pitch, bottleneck features and x-vector.

How each contribute towards privacy and utility in VPC settings.

### 6.2 Adding Differentially-Private Noise in F0 and BN

Work with Ali

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## Chapter 7

# Usability of anonymized speech for training ASR (20 pages)

Deadline: May 15, 2021

7.1 Impact of re-training ASR

7.2 Data augmentation

7.3 Model adaptation

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

## **Chapter 8**

### **Conclusion and Perspectives (4 pages)**

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11

# References

- [1] 24745:2011, D. I. (2011). Information Technology—Security techniques—Biometric Information Protection. *ISO/IEC JTC1 SC27 Security Techniques*. 1  
2  
3
- [2] 29100:2011, D. I. (2011). Information Technology—Security techniques—Privacy framework. *ISO/IEC JTC1 SC27 Information security, cybersecurity and privacy protection*. 4  
5
- [3] 30136, I. F. (2017). Information Technology—Performance Testing of Biometric Protection Schemes. *ISO/IEC JTC1 SC37 Biometrics*. 6  
7
- [4] Abdi, N., Ramokapane, K. M., and Such, J. M. (2019). More than smart speakers: security and privacy perceptions of smart home personal assistants. In *Fifteenth Symposium on Usable Privacy and Security ({SOUPS} 2019)*. 8  
9  
10
- [5] Abowd, J. M. (2018). Protecting the confidentiality of america’s statistics: Adopting modern disclosure avoidance methods at the census bureau. *Census Blogs: Research Matters*. 11  
12
- [6] Adams, O., Wiesner, M., Watanabe, S., and Yarowsky, D. (2019). Massively multilingual adversarial speech recognition. *arXiv preprint arXiv:1904.02210*. 13  
14
- [7] Adi, Y., Zeghidour, N., Collobert, R., Usunier, N., Liptchinsky, V., and Synnaeve, G. (2018). To reverse the gradient or not: An empirical comparison of adversarial and multi-task learning in speech recognition. *arXiv preprint arXiv:1812.03483*. 15  
16  
17
- [8] Adi, Y., Zeghidour, N., Collobert, R., Usunier, N., Liptchinsky, V., and Synnaeve, G. (2019). To reverse the gradient or not: An empirical comparison of adversarial and multi-task learning in speech recognition. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3742–3746. IEEE. 18  
19  
20  
21
- [9] Ahmed, S., Chowdhury, A. R., Fawaz, K., and Ramanathan, P. (2020). Preech: A system for privacy-preserving speech transcription. In *29th {USENIX} Security Symposium ({USENIX} Security 20)*, pages 2703–2720. 22  
23  
24
- [10] Alexander, A., Botti, F., Dessimoz, D., and Drygajlo, A. (2004). The effect of mismatched recording conditions on human and automatic speaker recognition in forensic applications. *Forensic science international*, 146:S95–S99. 25  
26  
27
- [11] Algabri, M., Mathkour, H., Bencherif, M. A., Alsulaiman, M., and Mekhtiche, M. A. (2017). Automatic speaker recognition for mobile forensic applications. *Mobile Information Systems*, 2017. 28  
29
- [12] Amodei, D., Ananthanarayanan, S., Anubhai, R., Bai, J., Battenberg, E., Case, C., Casper, J., Catanzaro, B., Cheng, Q., Chen, G., et al. (2016). Deep speech 2: End-to-end speech recognition in english and mandarin. In *International conference on machine learning*, pages 173–182. PMLR. 30  
31  
32
- [13] Aono, Y., Hayashi, T., Wang, L., Moriai, S., et al. (2017). Privacy-preserving deep learning via additively homomorphic encryption. *IEEE Transactions on Information Forensics and Security*, 13(5):1333–1345. 33  
34

- <sup>1</sup> [14] Ardila, R., Branson, M., Davis, K., Henretty, M., Kohler, M., Meyer, J., Morais, R., Saunders, L., Tyers, F. M., and Weber, G. (2019). Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670*.
- <sup>4</sup> [15] Arık, S. Ö., Chrzanowski, M., Coates, A., Diamos, G., Gibiansky, A., Kang, Y., Li, X., Miller, J., Ng, A., Raiman, J., et al. (2017). Deep voice: Real-time neural text-to-speech. In *International Conference on Machine Learning*, pages 195–204. PMLR.
- <sup>7</sup> [16] Association, I. P., Staff, I. P. A., et al. (1999). *Handbook of the International Phonetic Association: A guide to the use of the International Phonetic Alphabet*. Cambridge University Press.
- <sup>9</sup> [17] Ateniese, G., Mancini, L. V., Spognardi, A., Villani, A., Vitali, D., and Felici, G. (2015). Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers. *International Journal of Security and Networks*, 10(3):137–150.
- <sup>12</sup> [18] Ayala-Rivera, V. and Pasquale, L. (2018). The grace period has ended: An approach to operationalize gdpr requirements. In *2018 IEEE 26th International Requirements Engineering Conference (RE)*, pages 136–146. IEEE.
- <sup>15</sup> [19] Bachan, J., Kuczmarski, T., and Francuzik, P. (2012). Evaluation of synthetic speech using automatic speech recognition. In *XIV International PhD Workshop (OWD 2012). Conference Archives PTETiS*, volume 30, pages 500–505.
- <sup>18</sup> [20] Backstrom, L., Dwork, C., and Kleinberg, J. (2007). Wherefore art thou r3579x? anonymized social networks, hidden patterns, and structural steganography. In *Proceedings of the 16th international conference on World Wide Web*, pages 181–190.
- <sup>21</sup> [21] Bahdanau, D., Chorowski, J., Serdyuk, D., Brakel, P., and Bengio, Y. (2016). End-to-end attention-based large vocabulary speech recognition. In *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4945–4949. IEEE.
- <sup>24</sup> [22] Bahmaninezhad, F., Zhang, C., and Hansen, J. H. (2018). Convolutional neural network based speaker de-identification. In *Odyssey*, pages 255–260.
- <sup>26</sup> [23] Ballmer, T. and Brennstuhl, W. (2013). *Speech act classification: A study in the lexical analysis of English speech activity verbs*, volume 8. Springer Science & Business Media.
- <sup>28</sup> [24] Banisar, D. and Davies, S. (1999). Global trends in privacy protection: An international survey of privacy, data protection, and surveillance laws and developments. *J. Marshall J. Computer & Info. L.*, 18:1.
- <sup>31</sup> [25] Battenberg, E., Chen, J., Child, R., Coates, A., Li, Y. G. Y., Liu, H., Satheesh, S., Sriram, A., and Zhu, Z. (2017). Exploring neural transducers for end-to-end speech recognition. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 206–213.
- <sup>34</sup> [26] Bayardo, R. J. and Agrawal, R. (2005). Data privacy through optimal k-anonymization. In *21st International conference on data engineering (ICDE'05)*, pages 217–228. IEEE.
- <sup>36</sup> [27] Bayerl, S. P., Brasser, F., Busch, C., Frassetto, T., Jauernig, P., Kolberg, J., Nautsch, A., Riedhammer, K., Sadeghi, A.-R., Schneider, T., et al. (2019). Privacy-preserving speech processing via stpc and tees.
- <sup>38</sup> [28] Beenau, B. W., Bonalle, D. S., Fields, S. W., Gray, W. J., Larkin, C., Montgomery, J. L., and Saunders, P. D. (2010). Voiceprint biometrics on a payment device. US Patent 7,814,332.
- <sup>40</sup> [29] Beigi, G., Shu, K., Guo, R., Wang, S., and Liu, H. (2019). I am not what i write: Privacy preserving text representation learning. *arXiv preprint arXiv:1907.03189*.

- [30] Biemans, M. (1998). The effect of biological gender (sex) and social gender (gender identity) on three pitch measures. *Linguistics in the Netherlands*, 15(1):41–52. 1  
2
- [31] Bimbot, F., Bonastre, J.-F., Fredouille, C., Gravier, G., Magrin-Chagnolleau, I., Meignier, S., Merlin, T., Ortega-García, J., Petrovska-Delacrétaz, D., and Reynolds, D. A. (2004). A tutorial on text-independent 3  
4 speaker verification. *EURASIP Journal on Advances in Signal Processing*, 2004(4):1–22. 5
- [32] Bispham, M. K., Agrafiotis, I., and Goldsmith, M. (2018). A taxonomy of attacks via the speech 6  
7 interface.
- [33] Black, A., Taylor, P., Caley, R., and Clark, R. (1998). The festival speech synthesis system. 8
- [34] Board, E. D. P. (2021). Guidelines 02/2021 on virtual voice assistants. *EDPB Website*. 9
- [35] Bohn, D. (2019). Amazon says 100 million alexa devices have been sold — what's next? *The Verge*. 10
- [36] Brasser, F., Frassetto, T., Riedhammer, K., Sadeghi, A.-R., Schneider, T., and Weinert, C. (2018). 11  
Voiceguard: Secure and private speech processing. In *Interspeech*, volume 18, pages 1303–1307. 12
- [37] Brümmer, N. and Du Preez, J. (2006). Application-independent evaluation of speaker detection. 13  
*Computer Speech and Language*, 20(2-3):230–275. 14
- [38] Brümmer, N. and du Preez, J. A. (2006). Application-independent evaluation of speaker detection. 15  
*Computer Speech and Language*, 20(2-3):230–275. 16
- [39] Burkhardt, F. (2005). Emofilt: the simulation of emotional speech by prosody-transformation. In 17  
*Ninth European Conference on Speech Communication and Technology*. 18
- [40] Campbell, J. P. (1997). Speaker recognition: A tutorial. *Proceedings of the IEEE*, 85(9):1437–1462. 19
- [41] Carr, P. (2019). *English phonetics and phonology: An introduction*. John Wiley & Sons. 20
- [42] Champion, P., Jouvet, D., and Larcher, A. (2020). Speaker information modification in the VoicePrivacy 2020 toolchain. Research report, INRIA Nancy, équipe Multispeech ; LIUM - Laboratoire 21  
d’Informatique de l’Université du Mans. 22  
23
- [43] Chatterjee, S. (2019). Is data privacy a fundamental right in india? *International Journal of Law and 24  
Management*. 25
- [44] Chorowski, J. K., Bahdanau, D., Serdyuk, D., Cho, K., and Bengio, Y. (2015). Attention-based models 26  
for speech recognition. In *NIPS*, pages 577–585. 27
- [45] Chou, J.-c., Yeh, C.-c., and Lee, H.-y. (2019). One-shot voice conversion by separating speaker and 28  
content representations with instance normalization. *arXiv preprint arXiv:1904.05742*. 29
- [46] Chou, J.-c., Yeh, C.-c., Lee, H.-y., and Lee, L.-s. (2018). Multi-target voice conversion without parallel 30  
data by adversarially learning disentangled audio representations. *arXiv preprint arXiv:1804.02812*. 31
- [47] Chung, H., Iorga, M., Voas, J., and Lee, S. (2017). Alexa, can i trust you? *Computer*, 50(9):100–104. 32
- [48] Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. (2014). Empirical evaluation of gated recurrent 33  
neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*. 34
- [49] Chung, J. S., Nagrani, A., and Zisserman, A. (2018). Voxceleb2: Deep speaker recognition. *Proc. 35  
Interspeech 2018*, pages 1086–1090. 36

- 1 [50] Cohen-Hadria, A., Cartwright, M., McFee, B., and Bello, J. P. (2019). Voice anonymization in urban  
2 sound recordings. In *2019 IEEE 29th International Workshop on Machine Learning for Signal Processing*  
3 (*MLSP*), pages 1–6. IEEE.
- 4 [51] Commission, E. (2002). The eprivacy directive. *Official Journal of the European Union*.
- 5 [52] Commission, E. (2016). General data protection regulation. *Official Journal of the European Union*.
- 6 [53] Cormode, G., Jha, S., Kulkarni, T., Li, N., Srivastava, D., and Wang, T. (2018). Privacy at scale: Local  
7 differential privacy in practice. In *Proceedings of the 2018 International Conference on Management of*  
8 *Data*, pages 1655–1658.
- 9 [54] Costan, V. and Devadas, S. (2016). Intel sgx explained. *IACR Cryptol. ePrint Arch.*, 2016(86):1–118.
- 10 [55] Coull, S. E., Wright, C. V., Monrose, F., Collins, M. P., Reiter, M. K., et al. (2007). Playing devil’s  
11 advocate: Inferring sensitive information from anonymized network traces. In *Ndss*, volume 7, pages  
12 35–47.
- 13 [56] Craig, D. W., Pearson, J. V., Szelinger, S., Sekar, A., Redman, M., Corneveaux, J. J., Pawlowski, T. L.,  
14 Laub, T., Nunn, G., Stephan, D. A., et al. (2008). Identification of genetic variants using bar-coded  
15 multiplexed sequencing. *Nature methods*, 5(10):887–893.
- 16 [57] Csáji, B. C. et al. (2001). Approximation with artificial neural networks. *Faculty of Sciences, Etvs Lornd*  
17 *University, Hungary*, 24(48):7.
- 18 [58] Dathathri, R., Saarikivi, O., Chen, H., Laine, K., Lauter, K., Maleki, S., Musuvathi, M., and Mytkowicz,  
19 T. (2019). Chet: an optimizing compiler for fully-homomorphic neural-network inferencing. In *Proceed-  
20 ings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*,  
21 pages 142–156.
- 22 [59] de Jong, G., McDougall, K., and Nolan, F. (2007). Sound change and speaker identity: an acoustic  
23 study. In *Speaker Classification II*, pages 130–141. Springer.
- 24 [60] de l’Informatique et des Libertés, C. N. (2020). "on the record": Cnil publishes a white paper on voice  
25 assistants. *CNIL Website*.
- 26 [61] Dehak, N., Kenny, P. J., Dehak, R., Dumouchel, P., and Ouellet, P. (2010). Front-end factor analysis  
27 for speaker verification. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(4):788–798.
- 28 [62] Denisov, P., Vu, N. T., and Font, M. F. (2018). Unsupervised domain adaptation by adversarial  
29 learning for robust speech recognition. In *Speech Communication; 13th ITG-Symposium*, pages 1–5.  
30 VDE.
- 31 [63] Desai, S., Raghavendra, E. V., Yegnanarayana, B., Black, A. W., and Prahallad, K. (2009). Voice  
32 conversion using artificial neural networks. In *2009 IEEE International Conference on Acoustics, Speech  
33 and Signal Processing*, pages 3893–3896. IEEE.
- 34 [64] Deutsch, D., Henthorn, T., and Lapidis, R. (2011). Illusory transformation from speech to song. *The  
35 Journal of the Acoustical Society of America*, 129(4):2245–2252.
- 36 [65] Di Cerbo, F. and Trabelsi, S. (2018). Towards personal data identification and anonymization using  
37 machine learning techniques. In *European Conference on Advances in Databases and Information Systems*,  
38 pages 118–126. Springer.
- 39 [66] Dias, M., Abad, A., and Trancoso, I. (2018). Exploring hashing and cryptonet based approaches for  
40 privacy-preserving speech emotion recognition. In *2018 IEEE International Conference on Acoustics,  
41 Speech and Signal Processing (ICASSP)*, pages 2057–2061. IEEE.

- [67] Dibazar, A. A., Narayanan, S., and Berger, T. W. (2002). Feature analysis for automatic detection of pathological speech. In *2nd Joint EMBS-BMES Conference*, volume 1, pages 182–183. 1  
2
- [68] Duchi, J. C., Jordan, M. I., and Wainwright, M. J. (2013). Local privacy and statistical minimax rates. 3  
In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 429–438. IEEE. 4
- [69] Duchi, J. C., Jordan, M. I., and Wainwright, M. J. (2018). Minimax optimal procedures for locally 5  
private estimation. *Journal of the American Statistical Association*, 113(521):182–201. 6
- [70] Dudley, H. (1939a). The automatic synthesis of speech. *Proceedings of the National Academy of 7  
Sciences of the United States of America*, 25(7):377. 8
- [71] Dudley, H. (1939b). Remaking speech. *The Journal of the Acoustical Society of America*, 11(2):169–177. 9
- [72] Dudley, H. (1940a). The carrier nature of speech. *Bell System Technical Journal*, 19(4):495–515. 10
- [73] Dudley, H. (1940b). The vocoder—electrical re-creation of speech. *Journal of the Society of Motion 11  
Picture Engineers*, 34(3):272–278. 12
- [74] Dueck, D. (2009). *Affinity Propagation: Clustering Data by Passing Messages*. PhD thesis, University 13  
of Toronto. 14
- [75] Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006a). Calibrating noise to sensitivity in private 15  
data analysis. In *Theory of Cryptography (TCC)*. 16
- [76] Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006b). Calibrating noise to sensitivity in private 17  
data analysis. In *Theory of cryptography conference*, pages 265–284. Springer. 18
- [77] Dwork, C. and Roth, A. (2014). The Algorithmic Foundations of Differential Privacy. *Foundations 19  
and Trends in Theoretical Computer Science*, 9(3–4):211–407. 20
- [78] Dwork, C., Roth, A., et al. (2014). The algorithmic foundations of differential privacy. *Foundations 21  
and Trends in Theoretical Computer Science*, 9(3–4):211–407. 22
- [79] d’Alessandro, C. (2012). Voice source parameters and prosodic analysis. In *Methods in empirical 23  
prosody research*, pages 63–88. De Gruyter. 24
- [80] Edu, J. S., Such, J. M., and Suarez-Tangil, G. (2019). Smart home personal assistants: a security and 25  
privacy review. *arXiv preprint arXiv:1903.05593*. 26
- [81] El Ayadi, M., Kamel, M. S., and Karray, F. (2011). Survey on speech emotion recognition: Features, 27  
classification schemes, and databases. *Pattern Recognition*, 44(3):572–587. 28
- [82] Espinoza-Cuadros, F. M., Perero-Codosero, J. M., Antón-Martín, J., and Hernández-Gómez, L. A. 29  
(2020). Speaker de-identification system using autoencoders and adversarial training. *arXiv preprint 30  
arXiv:2011.04696*. 31
- [83] Fang, F., Wang, X., Yamagishi, J., Echizen, I., Todisco, M., Evans, N., and Bonastre, J.-F. (2019). 32  
Speaker anonymization using x-vector and neural waveform models. In *Proc. 10th ISCA Speech Synthesis 33  
Workshop*, pages 155–160. 34
- [84] Fant, G. (1993). Some problems in voice source analysis. *Speech Communication*, 13(1–2):7–22. 35
- [85] Fant, G. (2004). *Speech acoustics and phonetics: Selected writings*, volume 24. Springer Science & 36  
Business Media. 37

- <sup>1</sup> [86] Fant, G., Kruckenberg, A., Liljencrants, J., and Bavegdrd, M. (1994). Voice source parameters in continuous speech, transformation of lf-parameters. In *Third International Conference on Spoken Language Processing*.
- <sup>4</sup> [87] Farrús, M., Wagner, M., Anguita, J., and Hernando, J. (2008). How vulnerable are prosodic features to professional imitators? In *Odyssey*.
- <sup>6</sup> [88] Fernandes, E., Jung, J., and Prakash, A. (2016). Security analysis of emerging smart home applications. In *2016 IEEE symposium on security and privacy (SP)*, pages 636–654. IEEE.
- <sup>8</sup> [89] Feutry, C., Piantanida, P., Bengio, Y., and Duhamel, P. (2018). Learning anonymized representations with adversarial neural networks. *arXiv preprint arXiv:1802.09386*.
- <sup>10</sup> [90] Feyisetan, O., Balle, B., Drake, T., and Diethé, T. (2020). Privacy- and utility-preserving textual analysis via calibrated multivariate perturbations. In *WSDM*.
- <sup>12</sup> [91] Fienberg, S. E., Rinaldo, A., and Yang, X. (2010). Differential privacy and the risk-utility tradeoff for multi-dimensional contingency tables. In *International Conference on Privacy in Statistical Databases*, pages 187–199. Springer.
- <sup>15</sup> [92] Fontan, L., Ferrané, I., Farinas, J., Pinquier, J., Tardieu, J., Magnen, C., Gaillard, P., Aumont, X., and Füllgrabe, C. (2017). Automatic speech recognition predicts speech intelligibility and comprehension for listeners with simulated age-related hearing loss. *Journal of Speech, Language, and Hearing Research*, 60(9):2394–2405.
- <sup>19</sup> [93] Forney, G. (1973). The viterbi algorithm. *Proceedings of the IEEE*, 61(3):268–278.
- <sup>20</sup> [94] Fu, S.-W., Li, P.-C., Lai, Y.-H., Yang, C.-C., Hsieh, L.-C., and Tsao, Y. (2016). Joint dictionary learning-based non-negative matrix factorization for voice conversion to improve speech intelligibility after oral surgery. *IEEE Transactions on Biomedical Engineering*, 64(11):2584–2594.
- <sup>23</sup> [95] Fung, B. C., Wang, K., Chen, R., and Yu, P. S. (2010). Privacy-preserving data publishing: A survey of recent developments. *ACM Computing Surveys (Csur)*, 42(4):1–53.
- <sup>25</sup> [96] Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. (2016). Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030.
- <sup>28</sup> [97] Garson, J. (1997). Connectionism.
- <sup>29</sup> [98] Geiping, J., Bauermeister, H., Dröge, H., and Moeller, M. (2020). Inverting gradients—how easy is it to break privacy in federated learning? *arXiv preprint arXiv:2003.14053*.
- <sup>31</sup> [99] Ghahremani, P., Manohar, V., Hadian, H., Povey, D., and Khudanpur, S. (2017). Investigation of transfer learning for asr using lf-mmi trained neural networks. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 279–286. IEEE.
- <sup>34</sup> [100] Glackin, C., Chollet, G., Dugan, N., Cannings, N., Wall, J., Tahir, S., Ray, I. G., and Rajarajan, M. (2017). Privacy preserving encrypted phonetic search of speech data. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6414–6418. IEEE.
- <sup>37</sup> [101] Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings.

- [102] Gomez-Barrero, M., Galbally, J., Rathgeb, C., and Busch, C. (2017a). General framework to evaluate unlinkability in biometric template protection systems. *IEEE Transactions on Information Forensics and Security*, 13(6):1406–1420.
- [103] Gomez-Barrero, M., Galbally, J., Rathgeb, C., and Busch, C. (2017b). General framework to evaluate unlinkability in biometric template protection systems. *IEEE Transactions on Information Forensics and Security*, 13(6):1406–1420.
- [104] Gontier, F., Lagrange, M., Lavandier, C., and Petiot, J.-F. (2020). Privacy aware acoustic scene synthesis using deep spectral feature inversion. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 886–890. IEEE.
- [105] Graves, A. and Jaitly, N. (2014). Towards end-to-end speech recognition with recurrent neural networks. In *International conference on machine learning*, pages 1764–1772. PMLR.
- [106] Graves, A. and Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural networks*, 18(5-6):602–610.
- [107] Gu, Y., Li, X., Chen, S., Zhang, J., and Marsic, I. (2017). Speech intention classification with multimodal deep learning. In *Canadian Conference on Artificial Intelligence*, pages 260–271.
- [108] Gulati, A., Qin, J., Chiu, C.-C., Parmar, N., Zhang, Y., Yu, J., Han, W., Wang, S., Zhang, Z., Wu, Y., et al. (2020). Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*.
- [109] Gupta, P., Prajapati, G. P., Singh, S., Kamble, M. R., and Patil, H. A. (2020). Design of voice privacy system using linear prediction. In *2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 543–549. IEEE.
- [110] Haderlein, T., Moers, C., Möbius, B., Rosanowski, F., and Nöth, E. (2011). Intelligibility rating with automatic speech recognition, prosodic, and cepstral evaluation. In *International Conference on Text, Speech and Dialogue*, pages 195–202. Springer.
- [111] Hadian, H., Sameti, H., Povey, D., and Khudanpur, S. (2018). End-to-end speech recognition using lattice-free mmi. In *Interspeech*, pages 12–16.
- [112] Han, W., Zhang, Z., Zhang, Y., Yu, J., Chiu, C.-C., Qin, J., Gulati, A., Pang, R., and Wu, Y. (2020a). Contextnet: Improving convolutional neural networks for automatic speech recognition with global context. *arXiv preprint arXiv:2005.03191*.
- [113] Han, Y., Li, S., Cao, Y., Ma, Q., and Yoshikawa, M. (2020b). Voice-indistinguishability: Protecting voiceprint in privacy-preserving speech data release. In *2020 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6.
- [114] Hannun, A., Case, C., Casper, J., Catanzaro, B., Diamos, G., Elsen, E., Prenger, R., Satheesh, S., Sengupta, S., Coates, A., et al. (2014). Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*.
- [115] Hardcastle, W. J., Laver, J., and Gibbon, F. E. (2012). *The handbook of phonetic sciences*, volume 119. John Wiley & Sons.
- [116] Hashimoto, K., Yamagishi, J., and Echizen, I. (2016). Privacy-preserving sound to degrade automatic speaker verification performance. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5500–5504. IEEE.

- <sup>1</sup> [117] Hedelin, P. (1981). A tone oriented voice excited vocoder. In *ICASSP'81. IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 6, pages 205–208. IEEE.
- <sup>3</sup> [118] Hellbernd, N. and Sammler, D. (2016). Prosody conveys speaker's intentions: Acoustic cues for speech act perception. *Journal of Memory and Language*, 88:70–86.
- <sup>5</sup> [119] Hermansky, H. (2011). Speech recognition from spectral dynamics. *Sadhana*, 36(5):729–744.
- <sup>6</sup> [120] Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A.-r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T. N., et al. (2012). Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal processing magazine*, 29(6):82–97.
- <sup>9</sup> [121] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- <sup>11</sup> [122] Hsu, C.-C., Hwang, H.-T., Wu, Y.-C., Tsao, Y., and Wang, H.-M. (2016). Voice conversion from non-parallel corpora using variational auto-encoder. In *2016 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*, pages 1–6. IEEE.
- <sup>14</sup> [123] Hsu, C.-C., Hwang, H.-T., Wu, Y.-C., Tsao, Y., and Wang, H.-M. (2017a). Voice conversion from unaligned corpora using variational autoencoding wasserstein generative adversarial networks. *arXiv preprint arXiv:1704.00849*.
- <sup>17</sup> [124] Hsu, W.-N., Zhang, Y., and Glass, J. (2017b). Unsupervised learning of disentangled and interpretable representations from sequential data. *arXiv preprint arXiv:1709.07902*.
- <sup>19</sup> [125] Hu, S., Xie, X., Liu, S., Lam, M. W., Yu, J., Wu, X., Liu, X., and Meng, H. (2019). Lf-mmi training of bayesian and gaussian process time delay neural networks for speech recognition. In *Interspeech*, pages 2793–2797.
- <sup>22</sup> [126] Huang, W.-C., Hayashi, T., Watanabe, S., and Toda, T. (2020a). The sequence-to-sequence baseline for the voice conversion challenge 2020: Cascading asr and tts. *arXiv preprint arXiv:2010.02434*.
- <sup>24</sup> [127] Huang, W.-C., Hayashi, T., Wu, Y.-C., Kameoka, H., and Toda, T. (2021). Pretraining techniques for sequence-to-sequence voice conversion. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:745–755.
- <sup>27</sup> [128] Huang, W.-C., Luo, H., Hwang, H.-T., Lo, C.-C., Peng, Y.-H., Tsao, Y., and Wang, H.-M. (2020b). Unsupervised representation disentanglement using cross domain features and adversarial learning in variational autoencoder based voice conversion. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 4(4):468–479.
- <sup>31</sup> [129] Huang, X., Baker, J., and Reddy, R. (2014). A historical perspective of speech recognition. *Communications of the ACM*, 57(1):94–103.
- <sup>33</sup> [130] Huang, Y., Obada-Obieh, B., and Beznosov, K. (2020c). Amazon vs. my brother: How users of shared smart speakers perceive and cope with privacy risks. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–13.
- <sup>36</sup> [131] Ioffe, S. (2006). Probabilistic linear discriminant analysis. In *9th European Conference on Computer Vision (ECCV)*, pages 531–542.
- <sup>38</sup> [132] ISO/IEC 19795-1:2006 (2006). Information Technology — Biometric performance testing and reporting — Part 1: Principles and framework.
- <sup>40</sup> [133] Jackson, C. and Orebaugh, A. (2018). A study of security and privacy issues associated with the amazon echo. *International Journal of Internet of Things and Cyber-Assurance*, 1(1):91–100.

- [134] Jain, A., Hong, L., and Pankanti, S. (2000). Biometric identification. *Communications of the ACM*, 43(2):90–98.
- [135] Jin, Q., Toth, A. R., Schultz, T., and Black, A. W. (2009). Speaker de-identification via voice transformation. In *2009 IEEE Workshop on Automatic Speech Recognition & Understanding*, pages 529–533. IEEE.
- [136] Juang, B. H. and Chen, T. (1998). The past, present, and future of speech processing. *IEEE signal processing magazine*, 15(3):24–48.
- [137] Jurafsky, D. and Martin, J. H. (2000). Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition.
- [138] Justin, T., Štruc, V., Dobrišek, S., Vesnicić, B., Ipšić, I., and Mihelič, F. (2015). Speaker de-identification using diphone recognition and speech synthesis. In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, volume 4, pages 1–7.
- [139] Juvela, L., Bollepalli, B., Yamagishi, J., and Alku, P. (2019). Gelp: Gan-excited linear prediction for speech synthesis from mel-spectrogram. *arXiv preprint arXiv:1904.03976*.
- [140] Kameoka, H., Kaneko, T., Tanaka, K., and Hojo, N. (2018). Stargan-vc: Non-parallel many-to-many voice conversion using star generative adversarial networks. In *2018 IEEE Spoken Language Technology Workshop (SLT)*, pages 266–273. IEEE.
- [141] Kaneko, T. and Kameoka, H. (2018). Cyclegan-vc: Non-parallel voice conversion using cycle-consistent adversarial networks. In *2018 26th European Signal Processing Conference (EUSIPCO)*, pages 2100–2104. IEEE.
- [142] Kanervisto, A., Vestman, V., Sahidullah, M., Hautamäki, V., and Kinnunen, T. (2017). Effects of gender information in text-independent and text-dependent speaker verification. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5360–5364.
- [143] Karmakar, P., Teng, S. W., and Lu, G. (2021). Thank you for attention: A survey on attention-based artificial neural networks for automatic speech recognition. *arXiv preprint arXiv:2102.07259*.
- [144] Kasi, K. (2002). *Yet Another Algorithm for Pitch Tracking:(YAAPT)*. PhD thesis, Citeseer.
- [145] Kawahara, H. (2006). Straight, exploitation of the other aspect of vocoder: Perceptually isomorphic decomposition of speech sounds. *Acoustical science and technology*, 27(6):349–353.
- [146] Kenny, P. (2010a). Bayesian speaker verification with heavy-tailed priors. In *Odyssey*, page 14.
- [147] Kenny, P. (2010b). Bayesian speaker verification with heavy-tailed priors. In *Odyssey*.
- [148] Kepuska, V. and Bohouta, G. (2018). Next-generation of virtual personal assistants (Microsoft Cortana, Apple Siri, Amazon Alexa and Google Home). In *IEEE CCWC*, pages 99–103.
- [149] Kim, C. and Stern, R. M. (2008). Robust signal-to-noise ratio estimation based on waveform amplitude distribution analysis. In *Interspeech*, pages 2598–2601.
- [150] Kinnunen, T. and Li, H. (2010). An overview of text-independent speaker recognition: From features to supervectors. *Speech communication*, 52(1):12–40.
- [151] Kiparsky, P. (2003). The phonological basis of sound change. *The handbook of historical linguistics*, 313:342.

- 1 [152] Kniesza, V. (1988). Sound substitution, sound change, spelling in french loanwords in middle english.  
2 In *Luick Revisited: Papers Read at the Luick-Symposium at Schloss Liechtenstein, 15.-18.9. 1985*, volume  
3 288, page 205. Gunter Narr Verlag.
- 4 [153] Ko, T., Peddinti, V., Povey, D., Seltzer, M. L., and Khudanpur, S. (2017). A study on data augmen-  
5 tation of reverberant speech for robust speech recognition. In *2017 IEEE International Conference on*  
6 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 5220–5224. IEEE.
- 7 [154] Kopparapu, S. K. and Laxminarayana, M. (2010). Choice of mel filter bank in computing mfcc of a  
8 resampled speech. In *10th International Conference on Information Science, Signal Processing and their*  
9 *Applications (ISSPA 2010)*, pages 121–124. IEEE.
- 10 [155] Kottasová, I. (2018). These companies are getting killed by gdpr. *CNN Business*.
- 11 [156] Kotti, M. and Kotropoulos, C. (2008). Gender classification in two emotional speech databases. In  
12 *ICPR*, pages 1–4.
- 13 [157] Kratzenstein, C. G. (1782). Sur la naissance de la formation des voyelles. *Journal de Physique*,  
14 21:358–380.
- 15 [158] Kwon, O.-W., Chan, K., Hao, J., and Lee, T.-W. (2003). Emotion recognition by speech signals. In  
16 *EuroSpeech*.
- 17 [159] Labov, W. (1963). The social motivation of a sound change. *Word*, 19(3):273–309.
- 18 [160] Ladefoged, P. (1996). *Elements of acoustic phonetics*. University of Chicago Press.
- 19 [161] Latif, S., Khalifa, S., Rana, R., and Jurdak, R. (2020). Federated learning for speech emotion  
20 recognition applications. In *2020 19th ACM/IEEE International Conference on Information Processing*  
21 *in Sensor Networks (IPSN)*, pages 341–342. IEEE.
- 22 [162] Lau, J., Zimmerman, B., and Schaub, F. (2018). Alexa, are you listening? privacy perceptions, concerns  
23 and privacy-seeking behaviors with smart speakers. *Proceedings of the ACM on Human-Computer*  
24 *Interaction*, 2(CSCW):1–31.
- 25 [163] LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *nature*, 521(7553):436–444.
- 26 [164] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to  
27 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.
- 28 [165] Lee, H., Kim, S., Kim, J. W., and Chung, Y. D. (2017). Utility-preserving anonymization for health  
29 data publishing. *BMC medical informatics and decision making*, 17(1):1–12.
- 30 [166] Lei, X., Tu, G.-H., Liu, A. X., Li, C.-Y., and Xie, T. (2017). The insecurity of home digital voice  
31 assistants — Amazon Alexa as a case study. *arXiv preprint arXiv:1712.03327*.
- 32 [167] Leong, R. (2018). Analyzing the privacy attack landscape for amazon alexa devices.
- 33 [168] Leroy, D., Coucke, A., Lavril, T., Gisselbrecht, T., and Dureau, J. (2019). Federated learning for  
34 keyword spotting. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal*  
35 *Processing (ICASSP)*, pages 6341–6345. IEEE.
- 36 [169] Levin, A. and Nicholson, M. J. (2005). Privacy law in the united states, the eu and canada: The allure  
37 of the middle ground. *U. Ottawa L. & Tech.J.*, 2:357.
- 38 [170] Li, N., Li, T., and Venkatasubramanian, S. (2007). t-closeness: Privacy beyond k-anonymity and  
39 l-diversity. In *2007 IEEE 23rd International Conference on Data Engineering*, pages 106–115. IEEE.

- [171] Li, Q., Wen, Z., Wu, Z., Hu, S., Wang, N., and He, B. (2019). A survey on federated learning systems: vision, hype and reality for data privacy and protection. *arXiv preprint arXiv:1907.09693*. 1  
2
- [172] Liao, C.-F., Tsao, Y., Lee, H.-Y., and Wang, H.-M. (2018). Noise adaptive speech enhancement using 3  
domain adversarial training. *arXiv preprint arXiv:1807.07501*. 4
- [173] Lieberman, P., Laitman, J. T., Reidenberg, J. S., and Gannon, P. J. (1992). The anatomy, physiology, 5  
acoustics and perception of speech: essential elements in analysis of the evolution of human speech. 6  
*Journal of Human Evolution*, 23(6):447–467. 7
- [174] Lin, J.-L. and Wei, M.-C. (2008). An efficient clustering method for k-anonymization. In *Proceedings 8  
of the 2008 international workshop on Privacy and anonymity in information society*, pages 46–50. 9
- [175] Lippmann, R. P. (1989). Review of neural networks for speech recognition. *Neural computation*, 10  
1(1):1–38. 11
- [176] Liu, K., Zhang, J., and Yan, Y. (2007). High quality voice conversion through phoneme-based linear 12  
mapping functions with STRAIGHT for Mandarin. In *Proc. Fourth International Conference on Fuzzy 13  
Systems and Knowledge Discovery (FSKD 2007)*, volume 4, pages 410–414. 14
- [177] López, G., Quesada, L., and Guerrero, L. A. (2017). Alexa vs. Siri vs. Cortana vs. Google Assistant: 15  
a comparison of speech-based natural user interfaces. In *International Conference on Applied Human 16  
Factors and Ergonomics*, pages 241–250. 17
- [178] Lorenzo-Trueba, J., Fang, F., Wang, X., Echizen, I., Yamagishi, J., and Kinnunen, T. (2018a). Can we 18  
steal your vocal identity from the internet?: Initial investigation of cloning Obama’s voice using GAN, 19  
WaveNet and low-quality found data. In *Proc. Odyssey: The Speaker and Language Recognition Workshop*, 20  
pages 240–247. 21
- [179] Lorenzo-Trueba, J., Fang, F., Wang, X., Echizen, I., Yamagishi, J., and Kinnunen, T. (2018b). Can 22  
we steal your vocal identity from the internet?: Initial investigation of cloning obama’s voice using gan, 23  
wavenet and low-quality found data. In *Proc. Odyssey 2018 The Speaker and Language Recognition 24  
Workshop*, pages 240–247. 25
- [180] Lorenzo-Trueba, J., Yamagishi, J., Toda, T., Saito, D., Villavicencio, F., Kinnunen, T., and Ling, Z. 26  
(2018c). The Voice Conversion Challenge 2018: Promoting development of parallel and nonparallel 27  
methods. In *Odyssey*, pages 195–202. 28
- [181] Luger, E. and Sellen, A. (2016). “like having a really bad pa” the gulf between user expectation 29  
and experience of conversational agents. In *Proceedings of the 2016 CHI conference on human factors in 30  
computing systems*, pages 5286–5297. 31
- [182] Lukács, A. (2016). what is privacy? the history and definition of privacy. *University of Szeged*. 32
- [183] Luong, H.-T. and Yamagishi, J. (2019). Bootstrapping non-parallel voice conversion from speaker- 33  
adaptive text-to-speech. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop 34  
(ASRU)*, pages 200–207. IEEE. 35
- [184] Lyu, L., He, X., and Li, Y. (2020). Differentially private representation for nlp: Formal guarantee 36  
and an empirical study on privacy and fairness. In *EMNLP (Findings)*. 37
- [185] Lécuyer, M., Atlidakis, V., Geambasu, R., Hsu, D., and Jana, S. (2019). Certified robustness to 38  
adversarial examples with differential privacy. In *S&P*. 39
- [186] Machanavajjhala, A., Kifer, D., Gehrke, J., and Venkatasubramaniam, M. (2007). l-diversity: Privacy 40  
beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):3–es. 41

- <sup>1</sup> [187] Magariños, C., Lopez-Otero, P., Docio-Fernandez, L., Rodriguez-Banga, E., Erro, D., and Garcia-Mateo, C. (2017). Reversible speaker de-identification using pre-trained transformation functions. *Computer Speech and Language*, 46:36–52.
- <sup>4</sup> [188] Malkin, N., Deatrick, J., Tong, A., Wijesekera, P., Egelman, S., and Wagner, D. (2019). Privacy attitudes of smart speaker users. *Proceedings on Privacy Enhancing Technologies*, 2019(4):250–271.
- <sup>6</sup> [189] Manohar, V. et al. (2019). *Semi-supervised training for automatic speech recognition*. PhD thesis, Johns Hopkins University.
- <sup>8</sup> [190] Maouche, M., Srivastava, B. M. L., Vauquier, N., Bellet, A., Tommasi, M., and Vincent, E. (2020). A comparative study of speech anonymization metrics. In *Interspeech*, pages 1708–1712.
- <sup>10</sup> [191] Matrouf, D., Bonastre, J.-F., and Fredouille, C. (2006). Effect of speech transformation on impostor acceptance. In *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*, volume 1, pages I–I. IEEE.
- <sup>13</sup> [192] McAdams, S. E. (1984). *Spectral fusion, spectral parsing and the formation of auditory images*. Stanford university.
- <sup>15</sup> [193] McAulay, R. and Quatieri, T. (1986). Speech analysis/synthesis based on a sinusoidal representation. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 34(4):744–754.
- <sup>17</sup> [194] Mehri, S., Kumar, K., Gulrajani, I., Kumar, R., Jain, S., Sotelo, J., Courville, A., and Bengio, Y. (2016). Samplernn: An unconditional end-to-end neural audio generation model. *arXiv preprint arXiv:1612.07837*.
- <sup>20</sup> [195] Meng, Z., Li, J., Chen, Z., Zhao, Y., Mazalov, V., Gong, Y., and Juang, B.-H. (2018a). Speaker-invariant training via adversarial learning. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5969–5973. IEEE.
- <sup>23</sup> [196] Meng, Z., Li, J., Gong, Y., et al. (2018b). Adversarial feature-mapping for speech enhancement. *arXiv preprint arXiv:1809.02251*.
- <sup>25</sup> [197] Mercuri, R. T. and Neumann, P. G. (2003). Security by obscurity. *Communications of the ACM*, 46(11):160.
- <sup>27</sup> [198] Miao, Y., Gowayyed, M., and Metze, F. (2015). Eesen: End-to-end speech recognition using deep rnn models and wfst-based decoding. In *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, pages 167–174. IEEE.
- <sup>30</sup> [199] Morales, N., Gu, L., and Gao, Y. (2007). Adding noise to improve noise robustness in speech recognition. In *Eighth Annual Conference of the International Speech Communication Association*.
- <sup>32</sup> [200] Morise, M. (2015). CheapTrick, a spectral envelope estimator for high-quality speech synthesis. *Speech Communication*, 67:1–7.
- <sup>34</sup> [201] Morise, M., Yokomori, F., and Ozawa, K. (2016). World: a vocoder-based high-quality speech synthesis system for real-time applications. *IEICE TRANSACTIONS on Information and Systems*, 99(7):1877–1884.
- <sup>37</sup> [202] Muntés-Mulero, V. and Nin, J. (2009). Privacy and anonymization for very large datasets. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pages 2117–2118.
- <sup>39</sup> [203] Nagrani, A., Chung, J. S., and Zisserman, A. (2017). Voxceleb: A large-scale speaker identification dataset. *Proc. Interspeech 2017*, pages 2616–2620.

- [204] Nakashika, T., Takiguchi, T., and Ariki, Y. (2014). Voice conversion using rnn pre-trained by recurrent temporal restricted boltzmann machines. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(3):580–587.
- [205] Narayanan, A. and Shmatikov, V. (2008). Robust de-anonymization of large sparse datasets. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*, pages 111–125. IEEE.
- [206] Nasr, M., Shokri, R., and Houmansadr, A. (2019). Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *2019 IEEE symposium on security and privacy (SP)*, pages 739–753. IEEE.
- [207] Nautsch, A., Jasserand, C., Kindt, E., Todisco, M., Trancoso, I., and Evans, N. (2019a). The gdpr & speech data: Reflections of legal and technology communities, first steps towards a common understanding. *arXiv preprint arXiv:1907.03458*.
- [208] Nautsch, A., Jiménez, A., Treiber, A., Kolberg, J., Jasserand, C., Kindt, E., Delgado, H., Todisco, M., Hmani, M. A., Mtibaa, A., et al. (2019b). Preserving privacy in speaker and speech characterisation. *Computer Speech & Language*, 58:441–480.
- [209] Neekhara, P., Hussain, S., Dubnov, S., Koushanfar, F., and McAuley, J. (2021). Expressive neural voice cloning.
- [210] Ning, Y., He, S., Wu, Z., Xing, C., and Zhang, L.-J. (2019). A review of deep learning based speech synthesis. *Applied Sciences*, 9(19):4050.
- [211] Oord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. (2016). Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*.
- [212] Oppenheim, A. V., Buck, J. R., and Schafer, R. W. (2001). *Discrete-time signal processing. Vol. 2.* Upper Saddle River, NJ: Prentice Hall.
- [213] Orlandi, C., Piva, A., and Barni, M. (2007). Oblivious neural network computing via homomorphic encryption. *EURASIP Journal on Information Security*, 2007:1–11.
- [214] Paliwal, K. and Wójcicki, K. (2008). Effect of analysis window duration on speech intelligibility. *IEEE signal processing letters*, 15:785–788.
- [215] Panayotov, V., Chen, G., Povey, D., and Khudanpur, S. (2015). LibriSpeech: an ASR corpus based on public domain audio books. In *Proc. ICASSP*, pages 5206–5210.
- [216] Park, S.-w., Kim, D.-y., and Joe, M.-c. (2020). Cotatron: Transcription-guided speech encoder for any-to-many voice conversion without parallel data. *arXiv preprint arXiv:2005.03295*.
- [217] Pathak, M. A. (2012). *Privacy-preserving machine learning for speech processing*. Springer Science & Business Media.
- [218] Patino, J., Tomashenko, N., Todisco, M., Nautsch, A., and Evans, N. (2020). Speaker anonymisation using the mcdams coefficient. *arXiv preprint arXiv:2011.01130*.
- [219] Peddinti, V., Povey, D., and Khudanpur, S. (2015). A time delay neural network architecture for efficient modeling of long temporal contexts. In *Sixteenth Annual Conference of the International Speech Communication Association*.
- [220] Perry, T. L., Ohde, R. N., and Ashmead, D. H. (2001). The acoustic bases for gender identification from children’s voices. *The Journal of the Acoustical Society of America*, 109(6):2988–2998.

- [221] Pobar, M. and Ipšić, I. (2014). Online speaker de-identification using voice transformation. In *37th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*, pages 1264–1267.
- [222] Povey, D., Cheng, G., Wang, Y., Li, K., Xu, H., Yarmohammadi, M., and Khudanpur, S. (2018). Semi-orthogonal low-rank matrix factorization for deep neural networks. In *Interspeech*, pages 3743–3747.
- [223] Povey, D., Ghoshal, A., Boulian, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y., Schwarz, P., et al. (2011). The Kaldi speech recognition toolkit. Technical report.
- [224] Povey, D., Peddinti, V., Galvez, D., Ghahremani, P., Manohar, V., Na, X., Wang, Y., and Khudanpur, S. (2016). Purely sequence-trained neural networks for asr based on lattice-free MMI. In *Interspeech*, pages 2751–2755.
- [225] PrivacyGuard (2019). Smart speaker technology: Privacy risks and solutions. *PrivacyGuard*.
- [226] Qian, J., Du, H., Hou, J., Chen, L., Jung, T., and Li, X.-Y. (2018a). Hidebehind: Enjoy voice input with voiceprint unclonability and anonymity. In *Proc. the 16th ACM Conference on Embedded Networked Sensor Systems*, pages 82–94. ACM.
- [227] Qian, J., Du, H., Hou, J., Chen, L., Jung, T., Li, X.-Y., Wang, Y., and Deng, Y. (2017). Voicemask: Anonymize and sanitize voice input on mobile devices. *arXiv preprint arXiv:1711.11460*.
- [228] Qian, J., Han, F., Hou, J., Zhang, C., Wang, Y., and Li, X.-Y. (2018b). Towards privacy-preserving speech data publishing. In *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*, pages 1079–1087. IEEE.
- [229] Rabiner, L. (1989). A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286.
- [230] Raj, D., Snyder, D., Povey, D., and Khudanpur, S. (2019). Probing the information encoded in x-vectors. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 726–733.
- [231] Ramos, D. (2018). Voice assistants: How artificial intelligence assistants are changing our lives every day. *smartsheet*.
- [232] Ramteke, P. B., Dixit, A. A., Supanekar, S., Dharwadkar, N. V., and Koolagudi, S. G. (2018). Gender identification from children’s speech. In *2018 International Conference on Contemporary Computing (IC3)*, pages 1–6.
- [233] Rane, S. and Boufounos, P. T. (2013). Privacy-preserving nearest neighbor methods: Comparing signals without revealing them. *IEEE Signal Processing Magazine*, 30(2):18–28.
- [234] Rasmussen, D. J. (2013). Voice print identification for identifying speakers. US Patent 8,606,579.
- [235] Recommendation, I. (1994). Telephone transmission quality subjective opinion tests. a method for subjective performance assessment of the quality of speech voice output devices. page 85.
- [236] Reynolds, D. A. (1995). Speaker identification and verification using Gaussian mixture speaker models. *Speech Communication*, 17(1-2):91–108.
- [237] Rohdin, J., Biswas, S., and Shinoda, K. (2014). Constrained discriminative PLDA training for speaker verification. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1670–1674.

- [238] Rosenberg, A. and Ramabhadran, B. (2017). Bias and statistical significance in evaluating speech synthesis with mean opinion scores. In *Interspeech*, pages 3976–3980. 1  
2
- [239] Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386. 3  
4
- [240] Ruder, S. (2016). An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*. 5  
6
- [241] Ruggiero, G., Zovato, E., Caro, L. D., and Pollet, V. (2021). Voice cloning: a multi-speaker text-to-speech synthesis approach based on transfer learning. 7  
8
- [242] Rumelhart, D. E., McClelland, J. L., Group, P. R., et al. (1988). *Parallel distributed processing*, volume 1. IEEE Massachusetts. 9  
10
- [243] Sablayrolles, A., Douze, M., Schmid, C., Ollivier, Y., and Jégou, H. (2019). White-box vs black-box: Bayes optimal strategies for membership inference. In *International Conference on Machine Learning*, pages 5558–5567. PMLR. 11  
12  
13
- [244] Salmun, I., Opher, I., and Lapidot, I. (2016). On the use of PLDA i-vector scoring for clustering short segments. In *Odyssey*, pages 407–414. 14  
15
- [245] Samarati, P. and Sweeney, L. (1998). Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and suppression. 16  
17
- [246] Saon, G., Soltau, H., Nahamoo, D., and Picheny, M. (2013). Speaker adaptation of neural network acoustic models using i-vectors. In *2013 IEEE Workshop on Automatic Speech Recognition and Understanding*, pages 55–59. IEEE. 18  
19  
20
- [247] Schuller, B. and Batliner, A. (2013). *Computational paralinguistics: emotion, affect and personality in speech and language processing*. John Wiley & Sons. 21  
22
- [248] Schuller, B., Steidl, S., Batliner, A., Nöth, E., Vinciarelli, A., Burkhardt, F., van Son, R., Weninger, F., Eyben, F., Bocklet, T., et al. (2015). A survey on perceived speaker traits: Personality, likability, pathology, and the first challenge. *Computer Speech & Language*, 29(1):100–131. 23  
24  
25
- [249] Schuller, B., Steidl, S., Batliner, A., Vinciarelli, A., Scherer, K., Ringeval, F., Chetouani, M., Weninger, F., Eyben, F., Marchi, E., et al. (2013). The INTERSPEECH 2013 computational paralinguistics challenge: social signals, conflict, emotion, autism. In *Interspeech*, pages 148–152. 26  
27  
28
- [250] Sekiyama, K. (1997). Cultural and linguistic factors in audiovisual speech processing: The McGurk effect in Chinese subjects. *Perception & Psychophysics*, 59(1):73–80. 29  
30
- [251] Serdyuk, D., Audhkhasi, K., Brakel, P., Ramabhadran, B., Thomas, S., and Bengio, Y. (2016). Invariant representations for noisy speech recognition. *arXiv preprint arXiv:1612.01928*. 31  
32
- [252] Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N., Yang, Z., Chen, Z., Zhang, Y., Wang, Y., Skerry-Ryan, R., et al. (2018). Natural tts synthesis by conditioning wavenet on mel spectrogram predictions. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4779–4783. IEEE. 33  
34  
35  
36
- [253] Shi, J., Amith, J. D., García, R. C., Sierra, E. G., Duh, K., and Watanabe, S. (2021). Leveraging end-to-end asr for endangered language documentation: An empirical study on yolox\ochitl mixtec. *arXiv preprint arXiv:2101.10877*. 37  
38  
39
- [254] Shinohara, Y. (2016). Adversarial multi-task learning of deep neural networks for robust speech recognition. In *Interspeech*, pages 2369–2372. San Francisco, CA, USA. 40  
41

- <sup>1</sup> [255] Sholokhov, A., Kinnunen, T., Vestman, V., and Lee, K. A. (2020a). Extrapolating false alarm rates in automatic speaker verification. In *Interspeech*, pages 4218–4222.
- <sup>3</sup> [256] Sholokhov, A., Kinnunen, T., Vestman, V., and Lee, K. A. (2020b). Voice biometrics security: Extrapolating false alarm rate via hierarchical bayesian modeling of speaker verification scores. *Computer Speech and Language*, 60:101024.
- <sup>6</sup> [257] Shon, S., Dehak, N., Reynolds, D., and Glass, J. (2019). MCE 2018: The 1st multi-target speaker detection and identification challenge evaluation. In *Interspeech*, pages 356–360.
- <sup>8</sup> [258] Signol, F., Barras, C., and Liénard, J.-S. (2008). Evaluation of the pitch estimation algorithms in the monopitch and multipitch cases. In *Acoustics'08*.
- <sup>10</sup> [259] Sisman, B., Yamagishi, J., King, S., and Li, H. (2020). An overview of voice conversion and its challenges: From statistical modeling to deep learning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.
- <sup>13</sup> [260] Smaragdis, P. and Shashanka, M. (2007). A framework for secure speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(4):1404–1413.
- <sup>15</sup> [261] Snyder, D., Chen, G., and Povey, D. (2015). MUSAN: A music, speech, and noise corpus. *arXiv preprint arXiv:1510.08484v1*.
- <sup>17</sup> [262] Snyder, D., Garcia-Romero, D., Sell, G., Povey, D., and Khudanpur, S. (2018). X-vectors: Robust dnn embeddings for speaker recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5329–5333. IEEE.
- <sup>20</sup> [263] Spanias, A. S. (1994). Speech coding: A tutorial review. *Proceedings of the IEEE*, 82(10):1541–1582.
- <sup>21</sup> [264] Sriram, A., Jun, H., Gaur, Y., and Satheesh, S. (2018). Robust speech recognition using generative adversarial networks. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5639–5643. IEEE.
- <sup>24</sup> [265] Srivastava, B. M. L., Bellet, A., Tommasi, M., and Vincent, E. (2019). Privacy-preserving adversarial representation learning in ASR: Reality or illusion? In *Proc. INTERPSPEECH*, pages 3700–3704.
- <sup>26</sup> [266] Srivastava, B. M. L., Tomashenko, N., Wang, X., Vincent, E., Yamagishi, J., Maouche, M., Bellet, A., and Tommasi, M. (2020a). Design choices for x-vector based speaker anonymization. In *Interspeech*, pages 1713–1717.
- <sup>29</sup> [267] Srivastava, B. M. L., Vauquier, N., Sahidullah, M., Bellet, A., Tommasi, M., and Vincent, E. (2020b). Evaluating voice conversion-based privacy protection against informed attackers. In *2020 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 2802–2806.
- <sup>32</sup> [268] Stevens, K. N. (2001). Acoustic phonetics.
- <sup>33</sup> [269] Stolcke, A., Shriberg, E., Bates, R., Coccaro, N., Jurafsky, D., Martin, R., Meteer, M., Ries, K., Taylor, P., Van Ess-Dykema, C., et al. (1998). Dialog act modeling for conversational speech. In *AAAI Spring Symposium on Applying Machine Learning to Discourse Processing*, pages 98–105.
- <sup>36</sup> [270] Stylianou, Y. (2009). Voice transformation: a survey. In *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 3585–3588. IEEE.
- <sup>38</sup> [271] Sun, L., Li, K., Wang, H., Kang, S., and Meng, H. (2016). Phonetic posteriograms for many-to-one voice conversion without parallel data training. In *2016 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE.

- [272] Sun, S., Yeh, C.-F., Hwang, M.-Y., Ostendorf, M., and Xie, L. (2018). Domain adversarial training for accented speech recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4854–4858. IEEE.
- [273] Sundermann, D. and Ney, H. (2003). VTLN-based voice conversion. In *Proc. 3rd IEEE International Symposium on Signal Processing and Information Technology*, pages 556–559.
- [274] Sweeney, L. (2000). Simple demographics often identify people uniquely. *Health (San Francisco)*, 671(2000):1–34.
- [275] Synnaeve, G., Xu, Q., Kahn, J., Likhomanenko, T., Grave, E., Pratap, V., Sriram, A., Liptchinsky, V., and Collobert, R. (2019). End-to-end asr: from supervised to semi-supervised learning with modern architectures. *arXiv preprint arXiv:1911.08460*.
- [276] Székely, É., Henter, G. E., Beskow, J., and Gustafson, J. (2019). Spontaneous conversational speech synthesis from found data. In *Proc. INTERSPEECH*, pages 4435–4439.
- [277] Talkin, D. and Kleijn, W. B. (1995). A robust algorithm for pitch tracking (rapt). *Speech coding and synthesis*, 495:518.
- [278] Tan, H. H. and Lim, K. H. (2019). Vanishing gradient mitigation with deep learning neural network optimization. In *2019 7th International Conference on Smart Computing & Communications (ICSCC)*, pages 1–4. IEEE.
- [279] Tan, Z.-H., Dehak, N., et al. (2020). rvad: An unsupervised segment-based robust voice activity detection method. *Computer speech & language*, 59:1–21.
- [280] Tanaka, K., Kameoka, H., Kaneko, T., and Hojo, N. (2019). Atts2s-vc: Sequence-to-sequence voice conversion with attention and context preservation mechanisms. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6805–6809. IEEE.
- [281] Taylor, J. and Richmond, K. (2020). Enhancing sequence-to-sequence text-to-speech with morphology. *Submitted to IEEE ICASSP*.
- [282] Teixeira, F., Abad, A., and Trancoso, I. (2019). Privacy-preserving paralinguistic tasks. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6575–6579. IEEE.
- [283] Thiruvaran, T., Ambikairajah, E., and Epps, J. (2008). Fm features for automatic forensic speaker recognition. In *Ninth Annual Conference of the International Speech Communication Association*.
- [284] Tian, X., Chng, E. S., and Li, H. (2019). A speaker-dependent wavenet for voice conversion with non-parallel data. In *Interspeech*, pages 201–205.
- [285] Tian, X., Wang, J., Xu, H., Chng, E. S., and Li, H. (2018). Average modeling approach to voice conversion with non-parallel data. In *Odyssey*, volume 2018, pages 227–232.
- [286] Toda, T., Black, A. W., and Tokuda, K. (2007). Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(8):2222–2235.
- [287] Tomashenko, N., Srivastava, B. M. L., Wang, X., Vincent, E., Nautsch, A., Yamagishi, J., Evans, N., Bonastre, J.-F., Noé, P.-G., Todisco, M., and Patino, J. (2020a). The VoicePrivacy 2020 Challenge evaluation plan.

- [288] Tomashenko, N., Srivastava, B. M. L., Wang, X., Vincent, E., Nautsch, A., Yamagishi, J., Evans, N., Patino, J., Bonastre, J.-F., Noé, P.-G., and Todisco, M. (2020b). Introducing the VoicePrivacy initiative. In *Interspeech*, pages 1693–1697.
- [289] Tomashenko, N., Wang, X., Vincent, E., Patino, J., Srivastava, B. M. L., Noé, P.-G., Nautsch, A., Yamagishi, J., O’Brien, B., Chanclu, A., Bonastre, J.-F., Todisco, M., and Maouche, M. (2021). The VoicePrivacy 2020 Challenge: Results and findings. *Submitted to Computer Speech and Language*.
- [290] Torra, V. and Navarro-Arribas, G. (2016). Big data privacy and anonymization. In *IFIP International Summer School on Privacy and Identity Management*, pages 15–26. Springer.
- [291] Tripathi, A., Mohan, A., Anand, S., and Singh, M. (2018). Adversarial learning of raw speech features for domain invariant speech recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5959–5963. IEEE.
- [292] Tsuchiya, T., Tawara, N., Ogawa, T., and Kobayashi, T. (2018). Speaker invariant feature extraction for zero-resource languages with adversarial learning. In *IEEE ICASSP*, pages 2381–2385.
- [293] Tu, M., Tang, Y., Huang, J., He, X., and Zhou, B. (2019). Towards adversarial learning of speaker-invariant representation for speech emotion recognition. *arXiv preprint arXiv:1903.09606*.
- [294] Ulyanov, D., Vedaldi, A., and Lempitsky, V. (2017). Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *Proc. CVPR*, pages 6924–6932.
- [295] Umapathy, K. and Krishnan, S. (2005). Feature analysis of pathological speech signals using local discriminant bases technique. *Medical and Biological Engineering and Computing*, 43(4):457–464.
- [296] van Leeuwen, D. A. and Brümmer, N. (2007). An introduction to application-independent evaluation of speaker recognition systems. In *Speaker Classification I: Fundamentals, Features, and Methods*, pages 330–353.
- [297] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. *arXiv preprint arXiv:1706.03762*.
- [298] Veaux, C., Yamagishi, J., and MacDonald, K. (2019). CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit (version 0.92).
- [299] Ververidis, D. and Kotropoulos, C. (2004). Automatic speech classification to five emotional states based on gender information. In *EUSIPCO*, pages 341–344.
- [300] Vestman, V., Kinnunen, T., Hautamäki, R. G., and Sahidullah, M. (2020). Voice mimicry attacks assisted by automatic speaker verification. *Computer Speech & Language*, 59:36–54.
- [301] Vinciarelli, A., Pantic, M., and Bourlard, H. (2009). Social signal processing: Survey of an emerging domain. *Image and Vision Computing*, 27(12):1743–1759.
- [302] von Kempelen, W. (1791). *Le mécanisme de la parole, suivi de la description d'une machine parlante*. Imprimé chez B. Bauer.
- [303] Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., and Lang, K. J. (1989). Phoneme recognition using time-delay neural networks. *IEEE transactions on acoustics, speech, and signal processing*, 37(3):328–339.
- [304] Wang, X., Takaki, S., and Yamagishi, J. (2019a). Neural source-filter-based waveform model for statistical parametric speech synthesis. In *2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5916–5920.

- [305] Wang, X., Takaki, S., and Yamagishi, J. (2019b). Neural source-filter waveform models for statistical parametric speech synthesis. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:402–415. 1  
2  
3
- [306] Wang, Y., Skerry-Ryan, R., Stanton, D., Wu, Y., Weiss, R. J., Jaityl, N., Yang, Z., Xiao, Y., Chen, Z., Bengio, S., et al. (2017). Tacotron: A fully end-to-end text-to-speech synthesis model. *arXiv preprint arXiv:1703.10135*, 164. 4  
5  
6
- [307] Wang, Y., Stanton, D., Zhang, Y., Ryan, R.-S., Battenberg, E., Shor, J., Xiao, Y., Jia, Y., Ren, F., and Saurous, R. A. (2018). Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis. In *International Conference on Machine Learning*, pages 5180–5189. PMLR. 7  
8  
9
- [308] Warner, S. L. (1965). Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69. 10  
11
- [309] Warren, S. D. and Brandeis, L. D. (1890). Right to privacy. *Harv. L. Rev.*, 4:193. 12
- [310] Watanabe, S., Hori, T., Karita, S., Hayashi, T., Nishitoba, J., Unno, Y., Enrique Yalta Soplin, N., Heymann, J., Wiesner, M., Chen, N., Renduchintala, A., and Ochiai, T. (2018). ESPnet: End-to-end speech processing toolkit. In *Proc. INTERSPEECH*, pages 2207–2211. 13  
14  
15
- [311] Watanabe, S., Hori, T., Kim, S., Hershey, J. R., and Hayashi, T. (2017). Hybrid ctc/attention architecture for end-to-end speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 11(8):1240–1253. 16  
17  
18
- [312] Werbos, P. J. (1990). Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560. 19  
20
- [313] Wheatstone, C. (2011). *The Scientific Papers of Sir Charles Wheatstone*. Cambridge University Press. 21
- [314] Wiesner, M., Renduchintala, A., Watanabe, S., Liu, C., Dehak, N., and Khudanpur, S. (2018). Pre-training by backtranslation for end-to-end asr in low-resource settings. *arXiv preprint arXiv:1812.03919*. 22  
23
- [315] Wu, Y.-C., Hwang, H.-T., Hsu, C.-C., Tsao, Y., and Wang, H.-M. (2016). Locally linear embedding for exemplar-based spectral conversion. In *INTERSPEECH*, pages 1652–1656. 24  
25
- [316] Wu, Z., Evans, N., Kinnunen, T., Yamagishi, J., Alegre, F., and Li, H. (2015). Spoofing and countermeasures for speaker verification: A survey. *speech communication*, 66:130–153. 26  
27
- [317] Wu, Z., Virtanen, T., Chng, E. S., and Li, H. (2014). Exemplar-based sparse representation with residual compensation for voice conversion. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22(10):1506–1521. 28  
29  
30
- [318] Yi, J., Tao, J., Wen, Z., and Bai, Y. (2018). Language-adversarial transfer learning for low-resource speech recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(3):621–630. 31  
32
- [319] Yoo, I.-C., Lee, K., Leem, S., Oh, H., Ko, B., and Yook, D. (2020). Speaker anonymization for personal information protection using voice conversion techniques. *IEEE Access*, 8:198637–198645. 33  
34
- [320] Yu, D. and Deng, L. (2016). *AUTOMATIC SPEECH RECOGNITION*. Springer. 35
- [321] Yu, D. and Seltzer, M. (2011). Improved bottleneck features using pretrained deep neural networks. In *Interspeech*, pages 237–240. 36  
37
- [322] Zahorian, S. A. and Hu, H. (2008). A spectral/temporal method for robust fundamental frequency tracking. *The Journal of the Acoustical Society of America*, 123(6):4559–4571. 38  
39

- 1 [323] Zeghidour, N., Xu, Q., Liptchinsky, V., Usunier, N., Synnaeve, G., and Collobert, R. (2018). Fully  
2 convolutional speech recognition. *arXiv preprint arXiv:1812.06864*.
- 3 [324] Zeiler, M. D., Ranzato, M., Monga, R., Mao, M., Yang, K., Le, Q. V., Nguyen, P., Senior, A.,  
4 Vanhoucke, V., Dean, J., et al. (2013). On rectified linear units for speech processing. In *2013 IEEE*  
5 *International Conference on Acoustics, Speech and Signal Processing*, pages 3517–3521. IEEE.
- 6 [325] Zen, H., Dang, V., Clark, R., Zhang, Y., Weiss, R. J., Jia, Y., Chen, Z., and Wu, Y. (2019). LibriTTS:  
7 A Corpus Derived from LibriSpeech for Text-to-Speech. *Proc. Interspeech 2019*, pages 1526–1530.
- 8 [326] Zeng, Y.-M., Wu, Z.-Y., Falk, T., and Chan, W.-Y. (2006). Robust GMM based gender classification  
9 using pitch and RASTA-PLP parameters of speech. In *International Conference on Machine Learning*  
10 and *Cybernetics*, pages 3376–3379.
- 11 [327] Zhang, J.-X., Ling, Z.-H., Jiang, Y., Liu, L.-J., Liang, C., and Dai, L.-R. (2019a). Improving sequence-  
12 to-sequence voice conversion by adding text-supervision. In *ICASSP 2019-2019 IEEE International*  
13 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6785–6789. IEEE.
- 14 [328] Zhang, M., Sisman, B., Zhao, L., and Li, H. (2020). Deepconversion: Voice conversion with limited  
15 parallel training data. *Speech Communication*, 122:31–43.
- 16 [329] Zhang, M., Wang, X., Fang, F., Li, H., and Yamagishi, J. (2019b). Joint training framework for text-to-  
17 speech and voice conversion using multi-source tacotron and wavenet. *arXiv preprint arXiv:1903.12389*.
- 18 [330] Zhang, M., Zhou, Y., Zhao, L., and Li, H. (2021). Transfer learning from speech synthesis to voice  
19 conversion with non-parallel training data. *IEEE/ACM Transactions on Audio, Speech, and Language*  
20 *Processing*, 29:1290–1302.
- 21 [331] Zhang, S.-X., Gong, Y., and Yu, D. (2019c). Encrypted speech recognition using deep polynomial  
22 networks. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing*  
23 (*ICASSP*), pages 5691–5695. IEEE.
- 24 [332] Zhang, Y., Chan, W., and Jaitly, N. (2017). Very deep convolutional networks for end-to-end speech  
25 recognition. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*,  
26 pages 4845–4849.
- 27 [333] Zhao, L., Mammadov, M., and Yearwood, J. (2010). From convex to nonconvex: a loss function  
28 analysis for binary classification. In *2010 IEEE International Conference on Data Mining Workshops*,  
29 pages 1281–1288. IEEE.
- 30 [334] Zhou, Y., Tian, X., Xu, H., Das, R. K., and Li, H. (2019). Cross-lingual voice conversion with  
31 bilingual phonetic posteriogram and average modeling. In *ICASSP 2019-2019 IEEE International*  
32 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6790–6794. IEEE.

## Appendix A

### Extra top- $k$ results

This appendix presents the top-1, top-10 and top-50 precision plots that are computed as part of the large-scale speaker study described in Section 5.5.

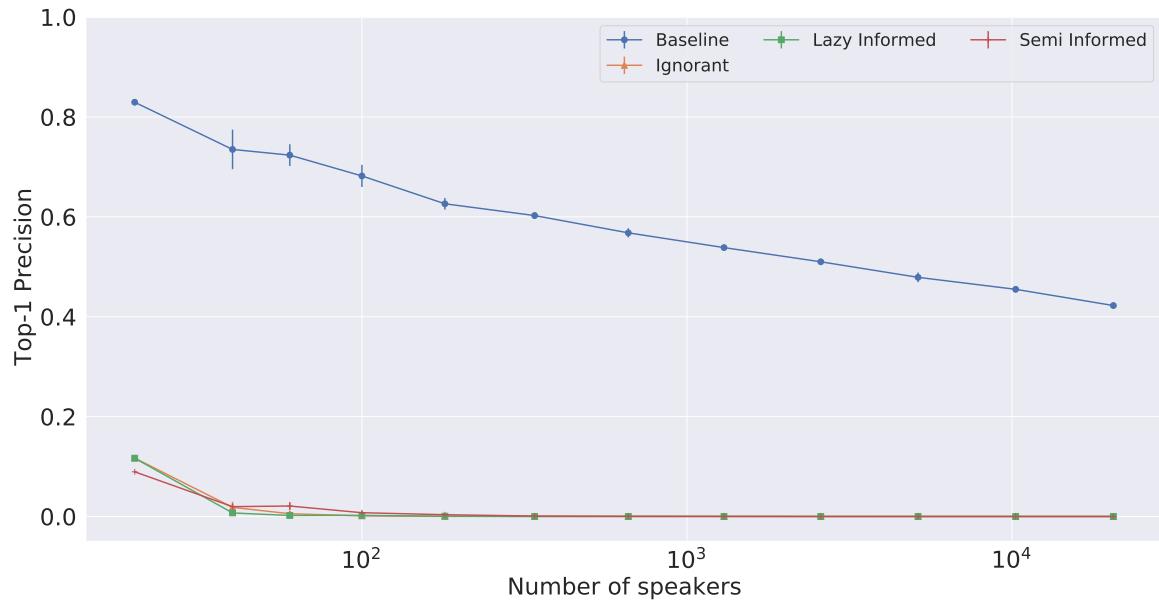


Fig. A.1 Speaker identification top-1 performance as enrollment speakers increase.

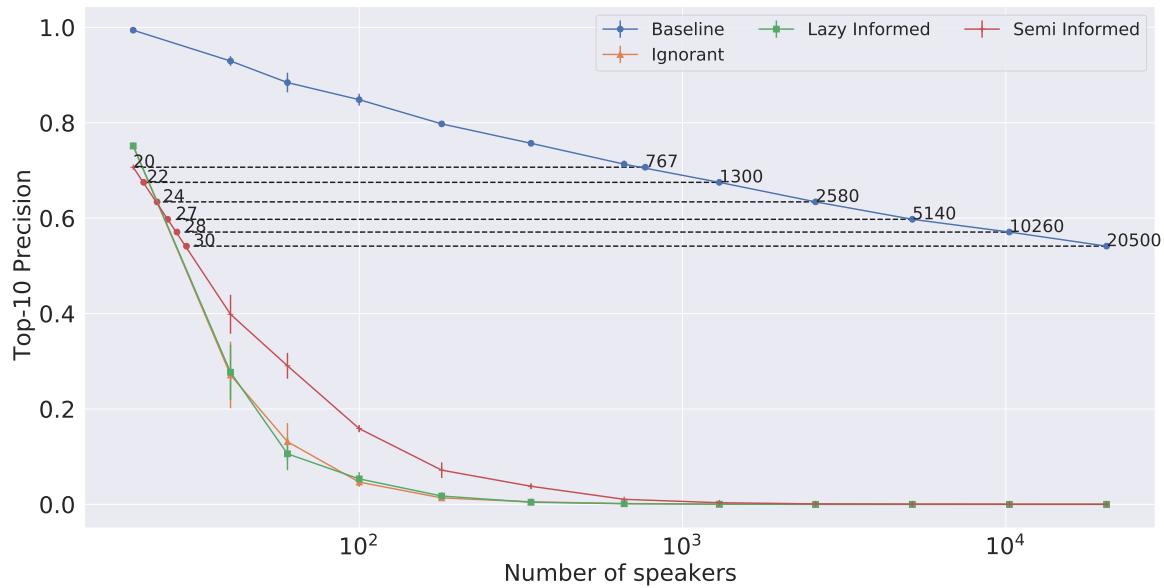


Fig. A.2 Speaker identification top-10 performance as enrollment speakers increase.

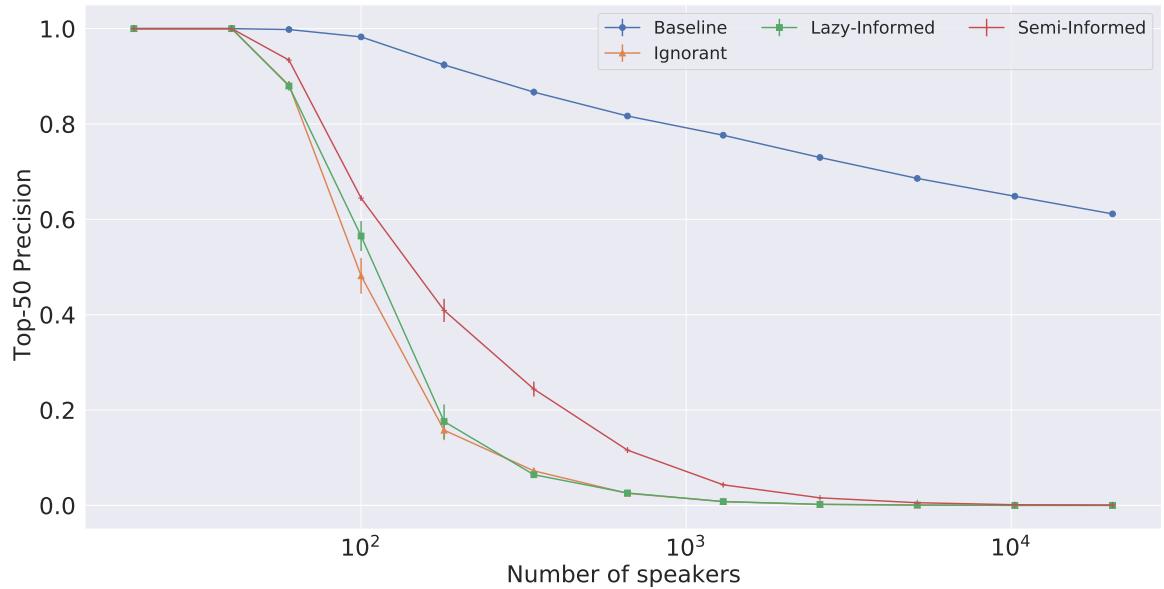


Fig. A.3 Speaker identification top-50 performance as enrollment speakers increase.

## Appendix B

# Worst-case analysis of the anonymization scheme

This appendix presents the top-1, top-10 and top-50 precision plots that are computed as part of the large-scale speaker study described in Section 5.5.

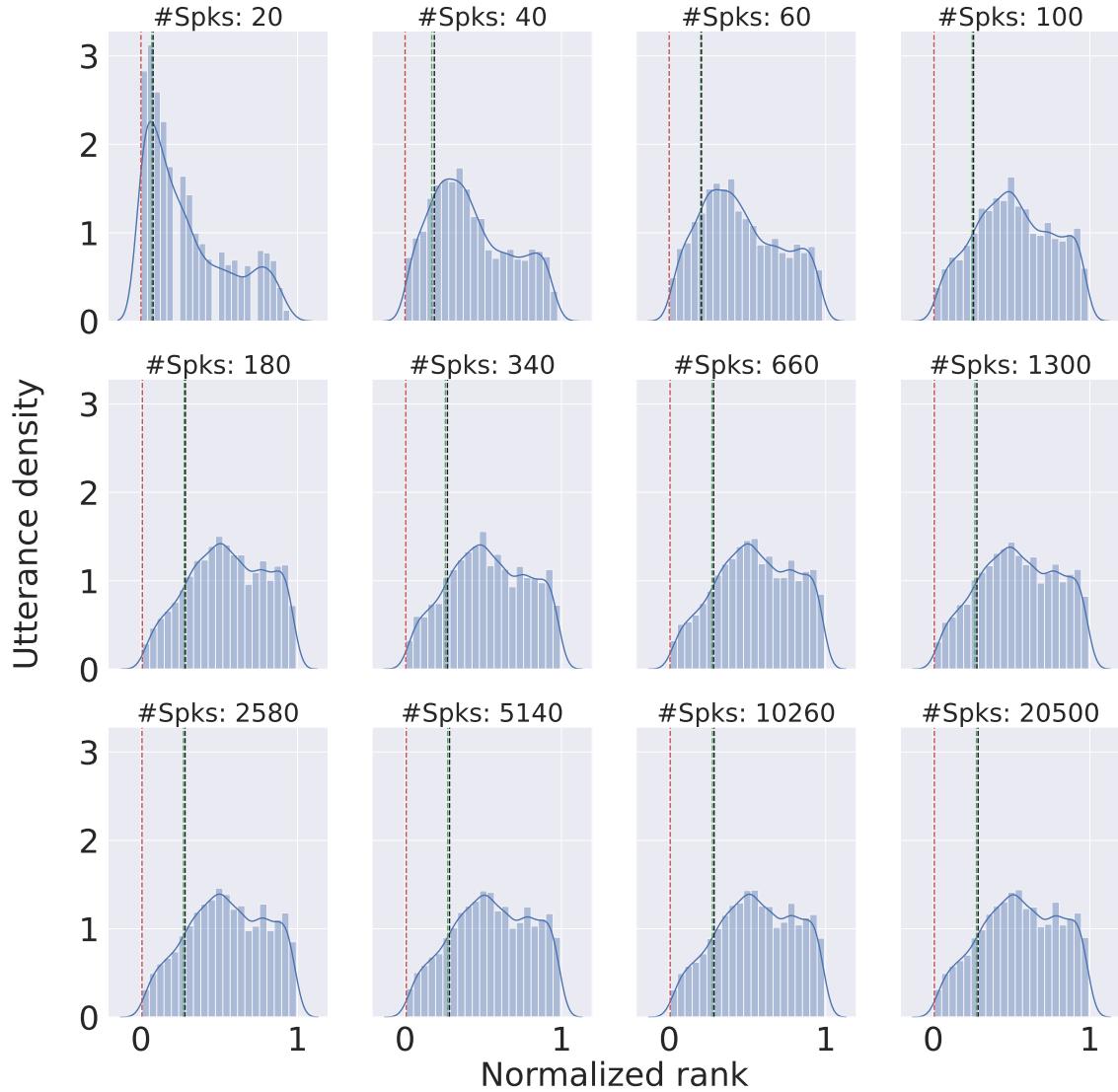


Fig. B.1 The normalized rank distribution for the *Ignorant* case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of  $U^{\text{worst}}$  (red),  $S^{\text{worst}}$  (green),  $U_S^{\text{worst}}$  (black).

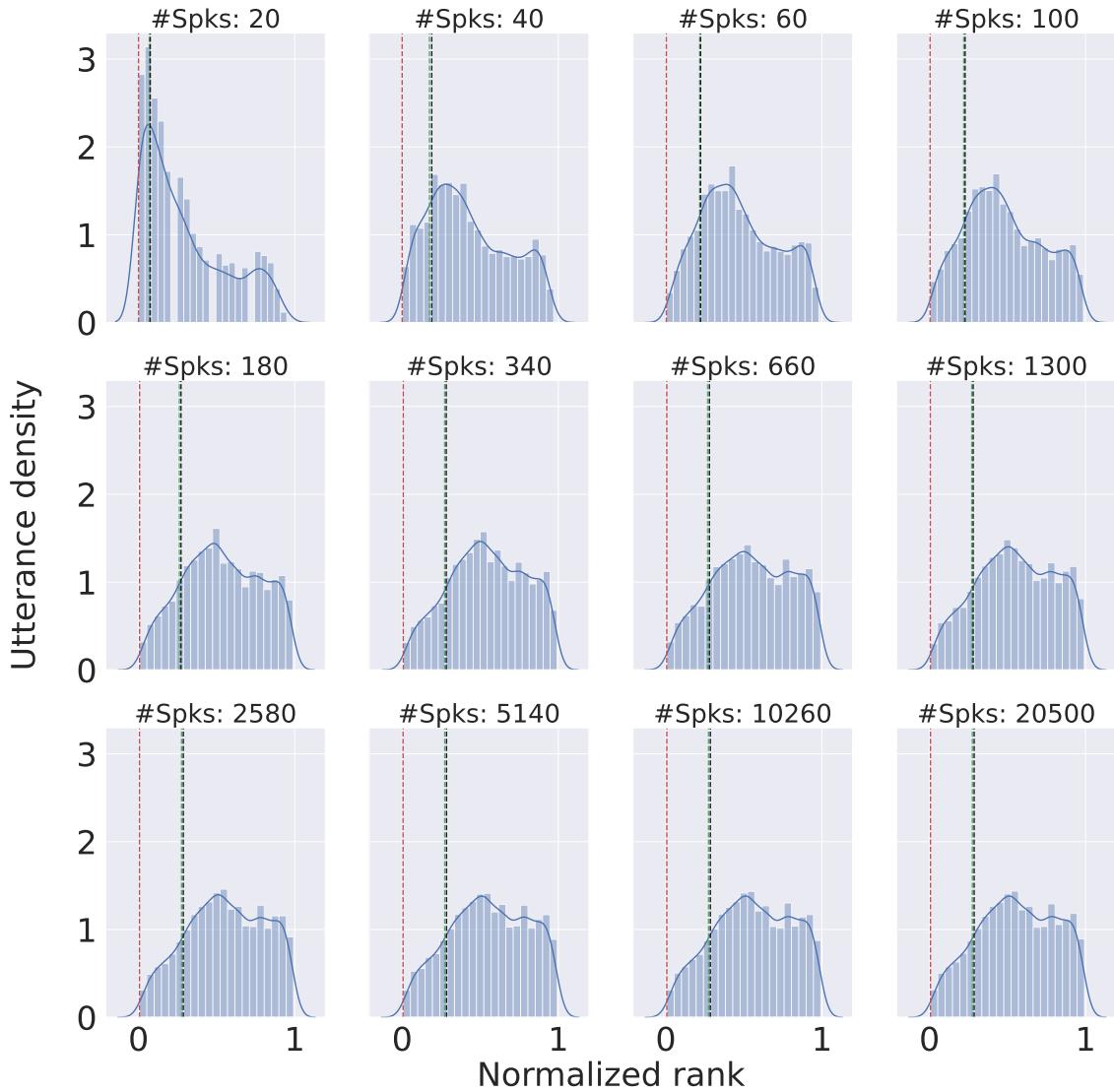


Fig. B.2 The normalized rank distribution for the *Lazy-Informed* case. X-axis represents the bins for normalized rank, and y-axis represents the density of utterances for each normalized rank bin. The number of enrollment speakers considered for a subplot are mentioned as the title of the subplot. The dashed vertical lines show the value of  $U^{\text{worst}}$  (red),  $S^{\text{worst}}$  (green),  $U_S^{\text{worst}}$  (black).

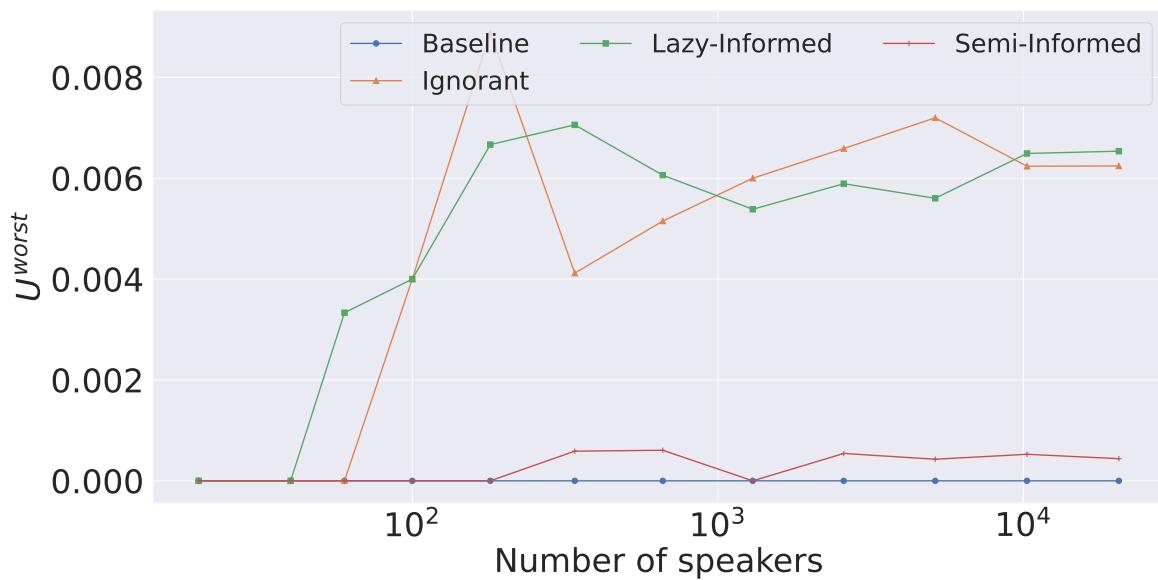


Fig. B.3 Normalized rank for the worst-performing utterances, i.e.,  $U^{\text{worst}}$ .

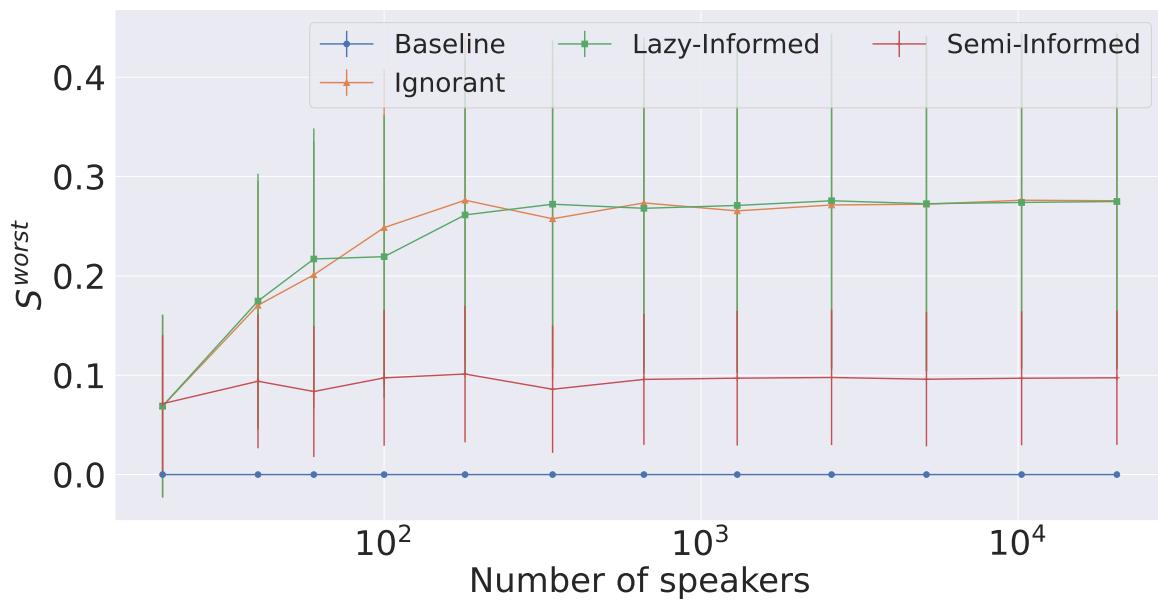


Fig. B.4 Normalized rank for the worst-performing speaker, i.e.,  $S^{\text{worst}}$ . Whiskers indicate the standard deviation of the normalized rank for the utterance of the worst-performing speaker.

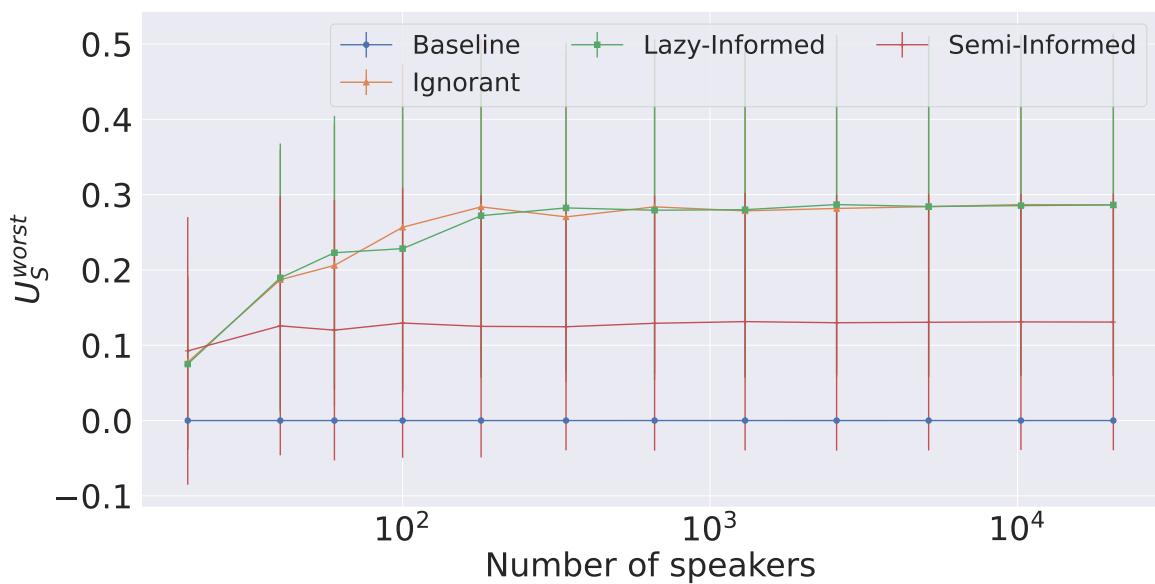


Fig. B.5 Normalized rank for the worst-performing utterances of each speaker, i.e.,  $U_S^{\text{worst}}$ . Whiskers indicate the standard deviation of the normalized rank for the worst-performing utterance of each speaker.

Draft - v1.0

Tuesday 7<sup>th</sup> September, 2021 – 12:11