



# First Steps Towards Voice Anonymization for Code-Switching Speech

Sarina Meyer, Ekaterina Kolos, Ngoc Thang Vu

Institute for Natural Language Processing, University of Stuttgart, Germany

sarina.meyer@ims.uni-stuttgart.de

## Abstract

The goal of voice anonymization is to modify an audio such that the true identity of its speaker is hidden. Research on this task is typically limited to the same English read speech datasets, thus the efficacy of current methods for other types of speech data remains unknown. In this paper, we present the first investigation of voice anonymization for the multilingual phenomenon of code-switching speech. We prepare two corpora for this task and propose adaptations to a multilingual anonymization model to make it applicable for code-switching speech. By testing the anonymization performance of this and two language-independent methods on the datasets, we find that only the multilingual system performs well in terms of privacy and utility preservation. Furthermore, we observe challenges in performing utility evaluations on this data because of its spontaneous character and the limited code-switching support by the multilingual speech recognition model.

**Index Terms:** voice anonymization, privacy, code-switching

## 1. Introduction

Voice anonymization (VA) refers to the task of hiding the identity of a speaker in an audio. Common approaches modify the audio using signal processing [1], voice conversion [2, 3, 4], or cascade systems of speech recognition (ASR) and text-to-speech (TTS) [5]. The main goal, as defined by the Voice Privacy Challenges (VPC) [6, 7, 8], is to protect the speaker's privacy from being identified by a speaker verification (ASV) system while keeping the utility for a downstream ASR task.

Most VA systems are evaluated using the procedures proposed by the VPC. While the exact test conditions and metrics differ in each edition, the evaluation datasets remain constant at subsets of LibriSpeech [9] and VCTK [10].<sup>1</sup> These datasets consist of single-speaker read speech, are monolingual in English, and recorded in noise-free environments. Thus, VA approaches are mainly evaluated on data that differs greatly from speech found in the wild such as spontaneous conversations. One aspect of many conversations across the world is that they are not, or not only, in English. While there have been previous approaches to explore VA for different languages, by developing multilingual [12, 13] and language-independent [2, 14] methods, to this date, the presence of several languages within one recording has not yet been investigated. This commonly occurs in multilingual communities as *code-switching*, in which a speaker alternates between two or more languages within a conversation, sentence or word. Code-switching is a phenomenon that has been found challenging for speech processing topics in the past [15], yet its effects on VA are so far unknown.

<sup>1</sup>VCTK was replaced by IEMOCAP [11] in VPC 2024 [8] to measure emotional preservation but not privacy or linguistic utility.

Thus, in this paper, we present to our knowledge the first investigation of VA in the context of code-switching. Specifically, we start by analyzing how well current VA methods deal with code-switching speech. While the way speakers code-switch differs between individuals [16] and might be preferred to be changed during anonymization, as a starting point, we assume that the code-switching behavior of the original speaker should be preserved. We prepare two existing code-switching datasets from Mandarin-English and Spanish-English communities for their applicability in anonymization and evaluate the performances of two language-independent VA models and one multilingual one that we adapt for code-switching. We find the anonymization of the multilingual model to be generally successful, achieving high privacy scores and only little degradation in utility, whereas the language-independent models offer almost no privacy protection already on monolingual audios. At the same time, the experiments reveal special challenges for the utility evaluation of this kind of data. While Whisper, being a powerful multilingual ASR model, is capable of keeping the speakers' code-switching behavior to some extent, its performance on the monolingual utterances of these datasets with 20-30% mixed error rate (MER) is substantially worse than the 1.5-3% achieved on the standard VPC datasets. Our analysis reveals that this is mainly due to the code-switching typical spontaneous nature and audio quality of the data. We observe a further drop to  $\sim 39\%$  MER for code-switched utterances and find that the original code-switching behavior is significantly altered and reduced during speech recognition towards more monolingual transcriptions, due to omissions and translations. We conclude that the evaluation of VA needs to be adapted to the nature of the data that is being anonymized. Thus, there is a need for a more diverse evaluation framework that supports a variety of speech phenomena, including code-switching. We publish our code and data online to encourage more research in this field<sup>2</sup>.

## 2. Methods

### 2.1. Data

We perform all experiments on two datasets, SEAME (Mandarin-English) and MIAMI (Spanish-English). They contain two types of code-switching: Inter-sentential, in which the language switch happens between utterances, and intra-sentential with language switches within an utterance. We split the data into single utterances and denote only intra-sentential cases as *code-switching* and everything else as *monolingual*. Thus, we have three *language settings* per dataset: *English* (EN), *Mandarin* (ZH) and *code-switching* (CS) in SEAME, and *English*, *Spanish* (ES) and *code-switching* in MIAMI.

<sup>2</sup><https://github.com/DigitalPhonetics/speaker-anonymization>

Table 1: Number of speakers and utterances of the SEAME and MIAMI datasets, divided into female (F) and male (M) subsets, language settings (EN, ZH/ES, CS) and ASV sets (enroll, trial).

SEAME		#spk	#utts EN		#utts ZH		#utts CS	
			Enroll	Trial	Enroll	Trial	Enroll	Trial
Train	F	28	1,632		3,792		8,646	
	M	26	1,326		2,824		9,942	
Dev	F	10	62	311	66	994	91	3,261
	M	10	68	1,165	70	685	90	3,211
Test	F	12	86	472	86	983	114	5,246
	M	8	54	1,121	53	572	75	2,317

  

MIAMI		#spk	#utts EN		#utts ES		#utts CS	
			Enroll	Trial	Enroll	Trial	Enroll	Trial
Test	F	17	100	2,442	101	884	139	425
	M	8	55	942	56	362	68	201

SEAME [17] consists of 297 conversational and interview recordings of 156 Singaporean and Malaysian speakers using Mandarin-English code-switching. We select only the interview part of the data, comprising 94 speakers and 210 recordings. 20 speakers each are selected for the development and testing of the anonymization, while we reserve the remaining 54 speakers for training or finetuning of an ASV attacker.

Bangor Miami<sup>3</sup> (shortened as MIAMI) consists of Spanish-English code-switching utterances of 56 conversational recordings of 84 speakers in Miami (USA). We exclude 15 recordings that contain only the speech of one specific speaker (María) and further exclude all speakers that have less than 20 utterances in each language setting. The remaining dataset comprises 25 speakers which we all use for testing.

Each dataset is split into three parts according to the annotated language setting, with their statistics shown in Table 1. We exclude utterances that overlap between speakers and remove annotations and discourse particles from transcriptions. Short utterances of less than two seconds in duration are concatenated with other short utterances of the same speaker and similar loudness levels. In the end, all utterances in SEAME have a duration between 2 and 42 seconds, with 130 to 1550 utterances per speaker. In MIAMI, utterances range between 2 to 21 seconds, with 110 to 400 utterances per speaker. We further split the development and test data into enrollment and trial subsets for ASV computation, with 4 to 10 utterances per speaker and language setting for enrollment and the remaining for trial.

These datasets are in several ways different to the datasets usually used for evaluating VA, namely LibriSpeech and VCTK. First, they are spontaneous and conversational speech. Utterances can be incomplete and ungrammatical, and contain repetitions, hesitations, or discourse particles. Moreover, being in a multilingual context, speakers might use words of other languages than the two primary ones and can have stronger accents. Finally, MIAMI has not been recorded in labs, thus recordings differ in audio quality. This includes noise and background speakers, as well as microphone differences.

## 2.2. Voice anonymization

We select three VA systems based on their ability to process several languages and their availability of open-source code.

B2 [1, 8] is a parameter-free technique based on signal processing and is known as baseline B2 of the VPC. We use the

<sup>3</sup><https://biling.talkbank.org/access/Bangor/Miami.html>

code of the VPC 2024<sup>4</sup>. It extracts pole positions from the input signal using linear predictive coding and then shifts these positions using the McAdams coefficient [18]. This method is generally seen as a weak baseline with low privacy protection but has the advantage of being naturally language-independent.

SSL<sup>5</sup> [2, 14] is a neural voice anonymization method specifically designed to be language-independent. From an input signal, it extracts the linguistic content as soft content units using HuBERT [19], F0 information using YAAPT [20], and the speaker embedding using an ECAPA-TDNN model [21]. The latter is anonymized by replacing it with an average of several speaker embeddings sampled from an external pool of speakers. The extracted and modified information is converted back into speech using a HiFiGAN vocoder [22].

GAN<sub>multi</sub><sup>6</sup> is a multilingual version of baseline B3 of the VPC 2024, as proposed in our previous work [12]. It is an ASR-TTS cascade system with a GAN-based [23] anonymization technique. Whisper [24] is used to recognize the linguistic content as a transcript, and to identify the language of the input if a language has not been specified. The speaker information is anonymized by sampling an artificial ECAPA-TDNN-like embedding from a GAN. An audio is synthesized using a multilingual FastSpeech2-based TTS [25] and a HiFiGAN vocoder, as provided in IMS Toucan [26]. The system has the option to preserve prosodic information during anonymization, however, we disabled this functionality because it would produce unreliable pitch estimates for these datasets. Instead, we rely on the prosody estimation of the TTS, as it has been done in a previous version of this model [27]. Since the system expects only one language per audio, we extend this multilingual version to a code-switching variant. For this, we update the TTS and HiFiGAN models to their latest versions<sup>7</sup> because of their improved multilingual support [28]. During synthesis, a phonemizer and language embedding is chosen depending on the input language. For code-switching, we first detect the language of each word in the recognized transcript using either the dragonmapper tool<sup>8</sup> to recognize Chinese characters (SEAME) or a BERT-based language identifier<sup>9</sup> for Spanish-English code-switching (MIAMI). We then phonemize each word according to its language and concatenate the phonemized representations of all words. Lastly, we adapt the code such that the language embedding is chosen per phone instead of once per utterance.

## 2.3. Evaluation

The anonymization is evaluated using the framework of the VPC 2024, which is based on [29]. Each utterance is anonymized towards a different target speaker. For privacy, an ASV model is trained as *semi-informed*, i.e., on data that has been anonymized by the same model that is being evaluated. The training data is LibriSpeech train-clean-360.<sup>10</sup> The privacy protection is then measured as the equal error rate (EER) of the ASV, in which a higher EER signifies better privacy. We perform ASV for female and male speakers separately and report

<sup>4</sup><https://github.com/Voice-Privacy-Challenge/Voice-Privacy-Challenge-2024>

<sup>5</sup><https://github.com/nii-yamagishilab/SSL-SAS>

<sup>6</sup><https://github.com/DigitalPhonetics/speaker-anonymization>

<sup>7</sup><https://github.com/DigitalPhonetics/IMS-Toucan/releases/tag/v3.0>

<sup>8</sup><https://github.com/tsroten/dragonmapper>

<sup>9</sup><https://huggingface.co/sagorsarker/codeswitch-spaeng-lid-lince>

<sup>10</sup>We experimented with finetuning this ASV model on the anonymized SEAME train data but did not observe an improvement compared to the LibriSpeech-only ASV model.

Table 2: Anonymization on VPC datasets. For WER and original EER, lower is better. For anonymized EER, higher is better.

Anon Model	EER (%)		WER (%)	
	Libri	VCTK	Libri	VCTK
<i>Original</i>	4.59	2.37	2.77	1.57
B2	7.81	4.69	4.13	9.83
SSL	2.05	3.46	<b>3.39</b>	3.11
GAN <sub>multi</sub>	<b>46.60</b>	<b>50.01</b>	3.64	<b>2.24</b>

their average. The performance is also assessed on the original, non-anonymized data, for which the ASV model is trained on the non-anonymized train data. The utility is typically measured as the word error rate (WER) of an English ASR system. We extend this to the mixed error rate (MER), which corresponds to WER for English and Spanish words, and to the character error rate (CER) for Mandarin ones. Following [12], we use Whisper-large-v3<sup>11</sup> for this purpose, similar to its use in the GAN<sub>multi</sub> model. Before computing the MER, we convert Chinese characters to Pinyin. Additionally, we measure the phone error rate (PER) using the phonemizer of the IMS Toucan toolkit [26].

#### 2.4. Recognizing code-switching speech with Whisper

Whisper is used for ASR in the anonymization of GAN<sub>multi</sub> as well as the utility evaluation. It is currently one of the most powerful multilingual ASR systems but is not specifically trained for code-switching. In fact, the model expects a single language per audio. However, we observe in our experiments that Whisper is able to transcribe code-switched speech, and is thus suitable for preliminary investigations of this topic, given that other code-switching specific ASR models like [30] are not publicly accessible. We experiment with different language prompts to find the optimal settings for the two datasets. For SEAME, we achieve the best results by prompting Whisper with the language *Mandarin*. For MIAMI, it is best to let Whisper recognize the language itself, but we rerun the recognition using the prompt *Spanish* if the recognized language is different from Spanish and English. Moreover, we set the output probabilities of special characters (e.g., digits, currency signs) to a low value such that Whisper would transcribe these words as pronounced. We use the same settings for the Whisper version in GAN<sub>multi</sub> and the utility evaluation.

### 3. Experiments

#### 3.1. Anonymization of VPC datasets

In order to understand the general anonymization abilities of each VA model, we first test them on the English datasets provided by the VPC [7]. The results are shown in Table 2. All methods increase privacy protection except SSL on LibriSpeech. Generally, the anonymization of B2 and SSL result in only weak protection of privacy, with EER scores below 10%. In contrast, the results of GAN<sub>multi</sub> are close to 50% EER, marking a strong anonymization. It also achieves the best utility preservation, having on average the lowest MER scores. We note that the performance of SSL is lower than reported in [2, 14] because they used a less informed evaluation strategy of earlier VPC editions. On the other hand, our privacy results for GAN<sub>multi</sub> are higher than in [12], which is mainly due to dis-

abling the prosody cloning mechanism, leading to less leakage of speaker information during VA.

#### 3.2. Anonymization of code-switching data

The results of the VA on SEAME and MIAMI are shown in Table 3. In terms of privacy, we see the same trends across anonymization systems as for LibriSpeech and VCTK. The data of B2 and SSL can be identified almost as well as the original data, with EER scores below 10%. Interestingly, this trend is broken for B2 on the Miami dataset: There, the privacy protection is considerably higher, with EER between 20 and 26%, which is due to increased noise levels after VA. For GAN<sub>multi</sub>, the results are equally good across all datasets, with EER close to 50% suggesting an overall successful anonymization.

Similarly, the GAN<sub>multi</sub> model is better at preserving the utility of the speech. There is an absolute degradation of 1-6% MER across datasets, but less pronounced than the 20-41% MER of B2 and 13-42% MER of SSL. Both of them distort the speech in such a way that utility is heavily affected. For B2, the degradation remains relatively constant across datasets and language settings. For SSL, though being considered language-independent, the anonymization has a smaller effect on English utterances than Spanish and code-switching speech. The PER scores are generally lower but show the same differences between the models, thus confirming these trends.

The differences between code-switching and monolingual utterances can be seen by comparing their results on original speech. Speaker verification seems to be facilitated by code-switching, leading to EER scores that are 2-4 points lower than on monolingual data. For utility, on the other hand, the ASR has more difficulties in recognizing the speech in CS than the monolingual subsets, with an absolute difference of up to 20% MER. However, given that Whisper expects audios to be monolingual, the error rate is with less than 40% MER still relatively small. We will analyze this finding in the following.

### 4. Analysis

Among the VA systems we tested, only GAN<sub>multi</sub> could achieve a high level of privacy and keep the utility degradation comparably low. However, the utility scores on the original data are substantially higher than on the standard VPC data. Thus, we examine the outcome of the ASR before and after anonymization with GAN<sub>multi</sub> more closely to understand to what extent the code-switching behavior is kept and reflected in the results.

#### 4.1. Preservation of code-switching during anonymization

We compare the number of code-switching points of an utterance in CS before (CSP(O)) and after (CSP(A)) anonymization with GAN<sub>multi</sub> in order to estimate how much the model changes the code-switching behavior. Although this does not show if the code-switching is kept for the same words but only the frequency, it gives some indication of how the model deals with code-switching. The transcription-based language identification tools of GAN<sub>multi</sub> are reused to map each recognized token to a language and count the number of language changes in an utterance. We notice  $CSP(O) = 0$  for 0.3% of utterances in SEAME and 14% in MIAMI, even though they had been annotated as code-switching. Thus, we only use utterances with  $CSP(O) > 0$  in our analysis, shown in Table 4. On average, code-switching in SEAME is reduced from  $CSP(O) = 3.47$  to  $CSP(A) = 1.75$ , with less code-switching points for 65% and no code-switching for 30% of utterances after VA. In a similar

<sup>11</sup><https://huggingface.co/openai/whisper-large-v3>

Table 3: Anonymization on code-switching datasets. For reference, all scores are also computed for the original data (Orig).

		Privacy				Utility							
		Orig↓	EER (%)↑		GAN <sub>multi</sub>	Orig	MER (%)↓			Orig	PER (%)↓		
			B2	SSL	GAN <sub>multi</sub>		B2	SSL	GAN <sub>multi</sub>		B2	SSL	GAN <sub>multi</sub>
SEAME	EN	7.54	9.46	7.41	<b>48.93</b>	27.39	57.05	41.41	<b>28.26</b>	16.86	40.62	27.02	<b>17.29</b>
	ZH	6.23	9.45	7.59	<b>48.40</b>	19.36	44.05	40.28	<b>25.31</b>	14.55	33.75	29.57	<b>17.61</b>
	CS	3.18	4.49	3.72	<b>50.61</b>	38.72	56.97	57.22	<b>44.27</b>	30.64	44.58	43.53	<b>35.05</b>
MIAMI	EN	8.48	21.77	6.29	<b>49.22</b>	20.12	52.13	33.16	<b>21.28</b>	15.03	41.08	25.31	<b>15.26</b>
	ES	8.78	19.81	7.08	<b>48.61</b>	28.86	69.37	70.89	<b>30.56</b>	20.71	52.37	51.16	<b>21.62</b>
	CS	6.75	20.81	4.19	<b>46.94</b>	39.29	74.28	78.59	<b>43.02</b>	30.39	55.59	58.41	<b>32.39</b>

Table 4: Code-switching points before (CSP(O)) and after (CSP(A)) anonymization. \* denotes a significant difference compared to the total MER.

	# utterances		MER	
	SEAME	MIAMI	SEAME	MIAMI
Total	7,533	536	44.24	43.90
CSP(A) < CSP(O)	4,926 (65%)	382 (71%)	50.90*	50.58*
CSP(A) = 0	2,281 (30%)	342 (64%)	59.36*	51.95*
CSP(A) = CSP(O)	2,232 (30%)	137 (26%)	27.81*	24.38*

30% of cases, the code-switching remains stable (CSP(A) = CSP(O)). For MIAMI, the reduction is more severe: 71% of utterances have less and 64% no code-switching after anonymization, it stays the same in 26% of cases. On average, code-switching is reduced from CSP(O) = 1.47 to CSP(A) = 0.52. This affects the MER: For both SEAME and MIAMI, the MER is significantly<sup>12</sup> higher if CSP(O) = 0 and significantly lower if CSP(A) = CSP(O), as compared to the total MER. The fact that the anonymized utterances have in most cases less code-switching than before might explain the higher total MER for CS compared to the monolingual subsets.

#### 4.2. Subjective analysis of anonymized speech content

To understand how the speech content is changed during VA, we perform a small user study with 6 subjects each for SEAME and MIAMI, selected based on their knowledge of the respective languages. For each dataset, we present ten CS audios as anonymized with GAN<sub>multi</sub> and their gold transcription to the subjects and ask them to compare the content of the audio to the text. The answers reveal how the errors made by the ASR in GAN<sub>multi</sub> are reflected in the anonymized speech. In both datasets, parts of one language were regularly translated into the more dominant language or recognized as similar-sounding words (near-homophones) of the other language. For MIAMI, the recognition would ignore one language altogether in some cases. Several words would be recognized as near-homophones in the same language. This same-language change was not always perceived as such by the listeners, e.g., 4 out of 6 did not hear the difference between “jobs” and “jumps” (MIAMI). Interestingly, for two utterances in SEAME, the content of the anonymized audio almost completely matched with the gold transcript but was mistranscribed during MER computation, resulting in an unfairly high error rate. The outcome of the user study raises the question of which kind of utterances might be

<sup>12</sup>All statistical significance tests were performed with the one-sided Mann-Whitney U test and  $\alpha = 0.025$ .

more prone to transcription errors than others.

#### 4.3. ASR performance based on data characteristics

Since SEAME and MIAMI come with rich annotations, we can divide the dataset into subsets depending on the presence of certain phenomena in an utterance. By comparing the MER achieved on each dataset after removing this subset to the scores in Table 3, we can estimate the effect of this characteristic on the utility computation. We report the results on the original data but could observe the same trends also after VA. In this analysis, we find that utterances containing abbreviations and acronyms (e.g., USA), foreign words in languages other than the two dataset languages, and filled pauses and discourse particles (e.g., “oh”), have a significant impact on the MER in SEAME. Without the respective samples, the dataset has only 29% (CS) to 38% (EN) of its previous size left but the MER is also reduced by 24% (EN, 21% MER), 14% (ZH, 17% MER) and 11% (CS, 35% MER). This shows a higher influence of these characteristics for monolingual than code-switching speech. For MIAMI, the most influential phenomena are the presence of words tagged as *unknown* by the annotator, repeated words or phrases, trailing off and incomplete utterances, and audios with a sub-average loudness. Their removal leads to a dataset of 37% (CS) to 47% (ES) of its original size and reduces the MER by 17% for EN (17% MER) and 15% for ES (24% MER). In contrast, the MER increases by 8% for CS (42% MER). We conclude that the characteristics that lead to errors in speech recognition are different for code-switching than for monolingual utterances. To improve the performance for CS, code-switching-specific processing is necessary.

## 5. Conclusion

In this paper, we present the first work of applying voice anonymization to code-switching speech. We prepare two Mandarin-English and Spanish-English datasets for the evaluation and experiment with three anonymization methods. Two of them are designed as language-independent but fail to achieve sufficient privacy preservation even on monolingual data. The third method is a multilingual model which we adapt for code-switching and which performs well in terms of privacy and utility preservation. We find, however, that Whisper, used for utility evaluation and in the multilingual model, has several difficulties in recognizing the speech in these datasets. In future work, further investigations should be made with other evaluation models and aspects such as speech emotion recognition, as well as more code-switching pairs and datasets.

## 6. Acknowledgements

This work is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project: Multilingual Controllable Voice Privacy (VoiPy) - Project number 533241795.

## 7. References

- [1] J. Patino, N. Tomashenko, M. Todisco, A. Nautsch, and N. Evans, “Speaker anonymisation using the mcadams coefficient,” in *Interspeech 2021*, 2021, pp. 1099–1103.
- [2] X. Miao, X. Wang, E. Cooper, J. Yamagishi, and N. Tomashenko, “Language-Independent Speaker Anonymization Approach Using Self-Supervised Pre-Trained Models,” in *Proc. The Speaker and Language Recognition Workshop (Odyssey 2022)*, 2022, pp. 279–286.
- [3] M. Panariello, F. Nespoli, M. Todisco, and N. Evans, “Speaker Anonymization Using Neural Audio Codec Language Models,” in *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024, pp. 4725–4729.
- [4] A. S. Shamsabadi, B. M. L. Srivastava, A. Bellet, N. Vauquier, E. Vincent, M. Maouche, M. Tommasi, and N. Papernot, “Differentially private speaker anonymization,” *Proc. Privacy Enhancing Technologies*, vol. 2023, no. 1, pp. 98–114, Jan. 2023.
- [5] S. Meyer, F. Lux, J. Koch, P. Denisov, P. Tilli, and N. T. Vu, “Prosody Is Not Identity: A Speaker Anonymization Approach Using Prosody Cloning,” in *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [6] N. Tomashenko, B. M. L. Srivastava, X. Wang, E. Vincent, A. Nautsch, J. Yamagishi, N. Evans, J. Patino, J.-F. Bonastre, P.-G. Noé, and M. Todisco, “Introducing the VoicePrivacy Initiative,” in *Interspeech 2020*, 2020, pp. 1693–1697.
- [7] M. Panariello, N. Tomashenko, X. Wang, X. Miao, P. Champion, H. Nourtel, M. Todisco, N. Evans, E. Vincent, and J. Yamagishi, “The VoicePrivacy 2022 Challenge: Progress and Perspectives in Voice Anonymisation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 32, pp. 3477–3491, 2024.
- [8] N. Tomashenko, X. Miao, P. Champion, S. Meyer, X. Wang, E. Vincent, M. Panariello, N. Evans, J. Yamagishi, and M. Todisco, “The VoicePrivacy 2024 Challenge Evaluation Plan,” 2024. [Online]. Available: <https://arxiv.org/abs/2404.02677>
- [9] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An ASR corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.
- [10] J. Yamagishi, C. Veaux, and K. MacDonald, “CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit (version 0.92),” 2019.
- [11] C. Busso, A. Kazemzadeh, and C.-C. Lee, “IEMOCAP: interactive emotional dyadic motion capture database,” 2008.
- [12] S. Meyer, F. Lux, and N. T. Vu, “Probing the feasibility of multilingual speaker anonymization,” in *Interspeech 2024*, 2024, pp. 4448–4452.
- [13] J. Yao, Q. Wang, P. Guo, Z. Ning, Y. Yang, Y. Pan, and L. Xie, “Musa: Multi-lingual speaker anonymization via serial disentanglement,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 33, pp. 1664–1674, 2025.
- [14] X. Miao, X. Wang, E. Cooper, J. Yamagishi, and N. Tomashenko, “Analyzing Language-Independent Speaker Anonymization Framework under Unseen Conditions,” in *Interspeech 2022*, 2022, pp. 4426–4430.
- [15] M. B. Mustafa, M. A. Yusoo, H. K. Khalaf, A. A. Rahman Mahmoud Abushariah, M. L. M. Kiah, H. N. Ting, and S. Muthaiyah, “Code-switching in automatic speech recognition: The issues and future directions,” *Applied Sciences*, vol. 12, no. 19, p. 9541, 2022.
- [16] P. Auer, “From codeswitching via language mixing to fused lects: Toward a dynamic typology of bilingual speech,” *International Journal of Bilingualism*, vol. 3, no. 4, pp. 309–332, 1999.
- [17] D.-C. Lyu, T.-P. Tan, E. S. Chng, and H. Li, “Seame: a mandarin-english code-switching speech corpus in south-east asia,” in *Interspeech 2010*, 2010, pp. 1986–1989.
- [18] S. McAdams, “Spectral fusion, spectral parsing and the formation of auditory images,” Ph.D. dissertation, Stanford University, Stanford, California, 05/1984 1984.
- [19] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhota, R. Salakhutdinov, and A. Mohamed, “Hubert: Self-supervised speech representation learning by masked prediction of hidden units,” *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, vol. 29, p. 3451–3460, Oct. 2021.
- [20] K. Kasi and S. A. Zahorian, “Yet another algorithm for pitch tracking,” in *2002 IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 1, 2002, pp. 1–361–364.
- [21] B. Desplanques, J. Thienpondt, and K. Demuynck, “Ecapa-tddn: Emphasized channel attention, propagation and aggregation in tddn based speaker verification,” in *Interspeech 2020*, 2020, pp. 3830–3834.
- [22] J. Kong, J. Kim, and J. Bae, “HiFi-GAN: Generative Adversarial Networks for Efficient and High Fidelity Speech Synthesis,” in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., pp. 17 022–17 033.
- [23] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, Eds., vol. 27. Curran Associates, Inc., 2014.
- [24] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” in *Proceedings of the 40th International Conference on Machine Learning*, ser. ICML’23. JMLR.org, 2023.
- [25] Y. Ren, C. Hu, X. Tan, T. Qin, S. Zhao, Z. Zhao, and T.-Y. Liu, “Fastspeech 2: Fast and high-quality end-to-end text to speech,” in *International Conference on Learning Representations*, 2021.
- [26] F. Lux, J. Koch, A. Schweitzer, and N. Thang Vu, “The ims toucan system for the blizzard challenge 2021,” in *The Blizzard Challenge 2021*, 2021, pp. 14–19.
- [27] S. Meyer, P. Tilli, P. Denisov, F. Lux, J. Koch, and N. T. Vu, “Anonymizing speech with generative adversarial networks to preserve speaker privacy,” in *2022 IEEE Spoken Language Technology Workshop (SLT)*, 2023, pp. 912–919.
- [28] F. Lux, S. Meyer, L. Behringer, F. Zalkow, P. Do, M. Coler, E. A. P. Habets, and N. T. Vu, “Meta Learning Text-to-Speech Synthesis in over 7000 Languages,” in *Interspeech 2024*. ISCA, 2024.
- [29] S. Meyer, X. Miao, and N. T. Vu, “VoicePAT: An Efficient Open-Source Evaluation Toolkit for Voice Privacy Research,” *IEEE Open Journal of Signal Processing*, vol. 5, pp. 257–265, 2024.
- [30] Y. Yang, Y. Peng, H. Huang, E. S. Chng, and X. Zhong, “Adapting openai’s whisper for speech recognition on code-switch mandarin-english seame and asru2019 datasets,” in *2024 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2024, pp. 1–6.