

Analyzing Language-Independent Speaker Anonymization Framework under Unseen Conditions

Xiaoxiao Miao¹, Xin Wang¹, Erica Cooper¹, Junichi Yamagishi¹, Natalia Tomashenko²

¹National Institute of Informatics, Japan ²LIA, University of Avignon, France

xiaoxiaomiao@nii.ac.jp

Abstract

In our previous work, we proposed a language-independent speaker anonymization system based on self-supervised learning models. Although the system can anonymize speech data of any language, the anonymization was imperfect, and the speech content of the anonymized speech was distorted. This limitation is more severe when the input speech is from a domain unseen in the training data. This study analyzed the bottleneck of the anonymization system under unseen conditions. It was found that the domain (e.g., language and channel) mismatch between the training and test data affected the neural waveform vocoder and anonymized speaker vectors, which limited the performance of the whole system. Increasing the training data diversity for the vocoder was found to be helpful to reduce its implicit language and channel dependency. Furthermore, a simple correlation-alignment-based domain adaptation strategy was found to be significantly effective to alleviate the mismatch on the anonymized speaker vectors. Audio samples¹ and source code² are available online.

Index Terms: speaker anonymization, self-supervised learning, CORrelation ALignment, multilingual HiFi-GAN.

1. Introduction

The human voice contains a wealth of personal information, such as the speaker's identity and emotion. Since personal information can be revealed by advanced speaker or other types of recognition systems, the demand for privacy-preserving technologies is growing. Although there is no legal definition of privacy [1], through the initiative called VoicePrivacy Challenge (VPC) 2020 [2, 3], the research community has defined a speaker anonymization task, the goal of which is to protect the speaker identity information (privacy) while maintaining the speech intelligibility and naturalness (utility).

While other related methods exist [4], a recently proposed deep neural network (DNN)-based [5] method was used as the VPC 2020 primary baseline to disentangle speaker and other information in the speech data and synthesize speech after anonymizing the speaker information. The effectiveness of this speaker anonymization system (SAS) has been confirmed on English test sets. However, it requires large amounts of text transcriptions for English training data to obtain accurate linguistic representations, which makes it impossible to use for an unknown language. While other digital-signal-processing-based methods [6] require little or no training data, they are less effective than the DNN-based methods at protecting the speaker identity [7, 3].

Aiming at an effective speaker anonymization solution that can be applied to the speech of any language, we have proposed

a self-supervised learning (SSL)-based language-independent SAS [8]. It uses an SSL-based content encoder to extract general context representations regardless of the language of the input speech. The whole SAS requires no text labels or other language-specific resources, allowing the system to anonymize speech data from any language. This SAS has been applied to Mandarin speech data. Even though the Mandarin language was completely unseen to the SAS, the Mandarin speech samples were anonymized reasonably well. However, it was also observed from an increased character error rate (CER) that the speech contents were distorted after anonymization. While the trade-off between speaker anonymization and speech intelligibility is common to many SASs [3], our goal is to push the limit of the language-independent SAS and improve both privacy and utility metrics in unseen conditions.

This study takes one step towards our goal by experimentally analyzing the performance bottleneck of the SAS. Specifically, it analyzes how the components such as the speech generator (i.e., vocoder) are implicitly dependent on a particular language or channel in the training database of that language. While keeping the target language (i.e., Mandarin) unseen, this study finds it beneficial to increase the language diversity of the training data, for example by adding German, Italian and Spanish speech data. Furthermore, this study investigates the language/channel mismatch brought by the impure speaker identity representation. It is found that the mismatch can be alleviated by transforming the anonymized speaker vector using a simple CORrelation ALignment (CORAL)-based domain adaptation strategy [9]. Transforming the anonymized speaker vector from the source domain (English) to a general domain covering German, Italian and Spanish data achieves better performance. When transforming to a more ideally matched domain with a few samples of Mandarin data, the CORAL provides the best utility with remaining high privacy. These findings are expected to be useful to the community for building a better language-independent SAS that can work under unseen conditions.

2. SSL-based Language-Independent Speaker Anonymization System

The baseline SSL-based language-independent SAS [8] disentangles speech into the fundamental frequency (F0), speaker identity representation, and content representation. Figure 1 shows the training and anonymization procedure of the baseline system (black and blue arrows path), respectively. There are two steps in the training stage:

1) *Original F0, speaker identity, and context feature extraction from the original speech recordings.* The YAAPT algorithm [10] is used to extract F0. The emphasized channel attention, propagation and aggregation in a time delay neural network (ECAPA-TDNN) speaker encoder trained on the *VoxCeleb-1* & 2 [11, 12] datasets is used to extract 192-dimensional speaker

¹<https://xiaoxiaomiao39.github.io/IS2022-SAS/>

²<https://github.com/xiaoxiaomiao39/SSL-SAS>

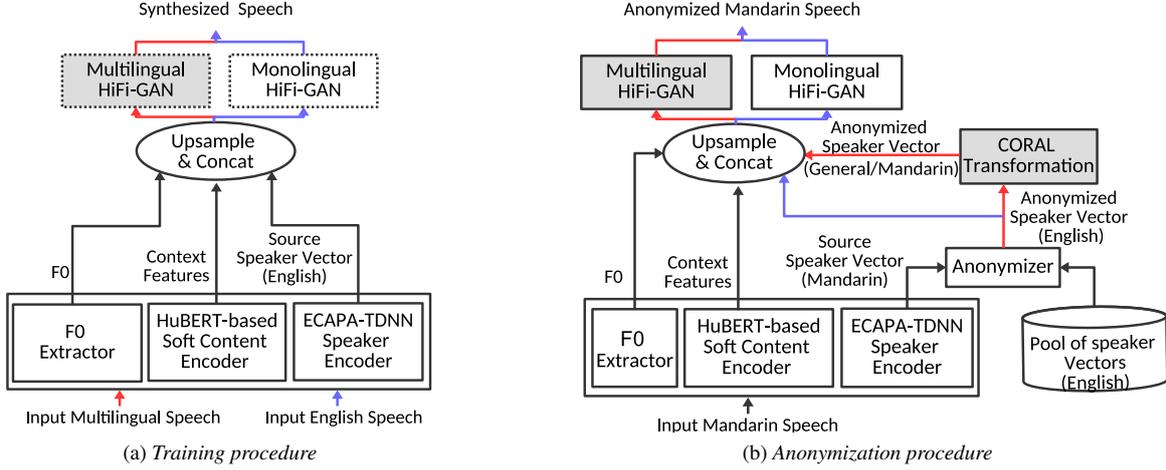


Figure 1: Diagram of language-independent SAS. Blue and red arrows indicate components investigated in this study.

identity vectors. To extract finer-grained context representations, the HuBERT-based soft content encoder [13] downsamples the input speech into a sequence of 768-dimensional continuous representations, which is reduced to 200 dimensions through linear projection. Note that this HuBERT-based soft content encoder took the CNN encoder and the sixth transformer layer from the input of a HuBERT Base model [14] pre-trained on LibriSpeech³ as the backbone. It was fine-tuned on the *LibriTTS-train-clean-100* [15] dataset, and the training criterion is detailed in [8].

2) *Speech synthesis*. The frame-wise content features, F0, and utterance-level source speaker vector are passed to the HiFi-GAN neural vocoder [16] after the operation of upsampling and concatenation to synthesize speech. The HiFi-GAN model was trained on *LibriTTS-train-clean-100* [15], denoted as monolingual HiFi-GAN.

An extra step called *Speaker vector anonymization* is included in the anonymization stage [2, 3]. Given a source speaker vector, cosine distance is used to find the 200 farthest speaker vectors from an external speaker vector pool (*LibriTTS train-other-500*). From these 200 vectors, 100 vectors are randomly selected and their average is used as the anonymized speaker vector [17]. Then, the content features, F0 and anonymized speaker vector instead of the source speaker vector are used to generate the anonymized speech.

3. Analysis of Performance Bottleneck

Although the above SAS performed reasonably well on the unseen Mandarin data, the anonymization performance was below the theoretical optimum. Furthermore, the contents of the anonymized speech were distorted by the SAS. It was shown that the CER of the anonymized speech reached 18.92%, which was much higher than the CER of 10.36% on the original data. Also, when the output speech was produced without actually anonymizing the speaker vector (i.e., resynthesis), its CER was still 14.81%. However, resynthesized speech from an ideal SAS should obtain a similar CER to the original speech.

The above results motivated us to analyze the performance bottleneck of our SAS under unseen circumstances. Here we examine the components in the system one by one and present techniques to alleviate the bottlenecks. Although we only consider Mandarin speech as the unseen test data, the analyses are expected to apply to other unseen languages as well.

³<https://github.com/pytorch/fairseq/tree/main/examples/hubert>

3.1. Robustness of HiFi-GAN

The CER gap between the original (10.36%) and resynthesized speech data (14.81%) suggests that the speech contents in the resynthesized speech may have been distorted. A similar finding has been reported on the language-dependent VPC baseline [18]. Since the system extracted features and re-synthesized the speech waveform without anonymizing the speaker representations, the increased CER is only caused by the content feature extractor and vocoder.

We hypothesize that the HiFi-GAN vocoder is one bottleneck⁴. It is known that to build a neural vocoder that generalizes well to unseen speakers in unseen languages, the training dataset has to cover diverse speakers and languages [19]. However, the HiFi-GAN in our SAS was trained using a subset from LibriTTS with many speakers, but only English-speaking ones. To verify the hypothesis, we used a multilingual database [20] to train the HiFi-GAN and compared its performance with the HiFi-GAN trained on the monolingual LibriTTS data in the experiments. This comparison is illustrated in Figure 1(a). Note that the multilingual database does not contain Mandarin data, and its details are explained in Section 4.2.

3.2. Language and Channel Mismatch on Speaker Vectors

Comparing the CERs in the resynthesized (14.81%) and anonymized (18.92%) cases, the increased CER is likely due to the anonymized speaker vector since it is the only difference between the resynthesized and anonymized data. Many studies have shown that speaker vectors contain speaker-unrelated information from the source domain, e.g., channel conditions and lexical contents [21, 22]. Because the anonymized vector is composed from the pool of English speaker vectors (i.e., those from *LibriTTS train-other-500* [15]), it may carry irrelevant information pertinent to the English database. Therefore, directly using the anonymized vector on Mandarin data may introduce language, channel, or other types of domain mismatch.

To verify this hypothesis, we use a simple but effective unsupervised domain adaption technique called CORAL [9]. The goal of CORAL is to find a transfer matrix \mathbf{A} that can align the feature distributions of the source domain and target domain by minimizing their covariances. Suppose source-domain English speaker vectors $\mathcal{D}_S = \{\mathbf{n}_i\}$, $\mathbf{n}_i \in \mathbb{R}^{192}$, and the target speaker vectors $\mathcal{D}_T = \{\mathbf{m}_i\}$, $\mathbf{m}_i \in \mathbb{R}^{192}$, here \mathbf{n}_i and \mathbf{m}_i are the 192-dimensional vectors extracted from the last projection layer of the ECAPA-TDNN. After the feature normaliza-

⁴We also investigated a multilingual-trained SSL-based soft content encoder but no improvements were observed.

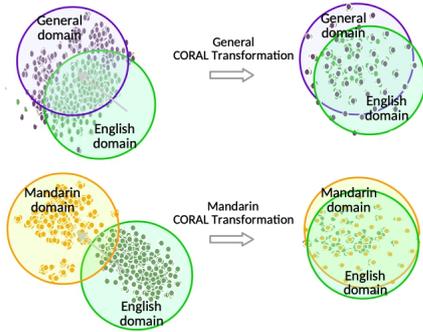


Figure 2: CORAL transformation on speaker vectors.

tion where the speaker vectors \mathcal{D}_S and \mathcal{D}_T are normalized to have zero mean and unit standard deviation in each dimension, the relationship of source- and target-domain statistics is that mean $\mu_S = \mu_T = 0$ and covariance matrices $C_S \neq C_T$. The transfer matrix A can be obtained via:

$$\text{CORAL}(\mathcal{D}_S, \mathcal{D}_T) = \min_A \|A^T C_S A - C_T\|_F^2$$

where $\|\cdot\|_F^2$ represents the matrix Frobenius norm. The detailed procedure to compute the optimal A^* can be found in [9, 23]. Then, the anonymized English speaker vectors transformed by A^* are used as the new anonymized speaker vector.

We consider two cases of CORAL transformations based on the data availability:

General CORAL transformation: we assume that Mandarin data are unavailable, and the target-domain speaker vectors are collected from the German, Italian, and Spanish speech data;

Mandarin CORAL transformation: we assume that a few (no more than 100) Mandarin samples that are completely disjoint (unseen-unheard speakers) with the test set are available as the target-domain speaker vectors.

In the first case, where the domain of the test data to be anonymized remains unseen to the SAS, CORAL using the target speaker vectors from multiple languages and channels is expected to alleviate the mismatch between the English vector pool and the vector to be anonymized. The second case is an oracle scenario where the mismatch caused by the impure speaker vectors can be reduced to a greater extent when they are transferred to the matched domain.

Figure 2 plots the t-SNE visualization of speaker vectors before and after the two cases of CORAL transformations. These speaker vectors are sampled from the different datasets listed at the bottom of Table 1. It is observed that the original distributions of speaker vectors from the different domains are widely spread due to the language and channel mismatch. The distribution between English and general domains has more overlap than between the English and Mandarin domains. Despite the different degrees of overlap, CORAL transformations push the speaker vectors from different domains to move closer.

4. Experiments

4.1. Evaluation Protocols

Our experiments followed the evaluation protocols of VPC 2020 [2, 3]. The SAS anonymized the Mandarin test trials shared by users to protect the speaker identity while preserving the speech contents. To assess how well the speech contents are preserved, CER was computed using a language-matched ASR (ASR_{eval}) as one utility metric. The protection of speaker identity was evaluated via one privacy matrix with equal error rates (EERs) of a language-matched ASV evaluation model (ASV_{eval}) in two setups:

Table 1: The datasets used in the different models.

Model	Dataset	
Training	Monolingual HiFi-GAN	LibriTTS train-clean-100
	Multilingual HiFi-GAN	German & Italian & Spanish LibriTTS train-clean-100
Tran.	General CORAL	German & Italian & Spanish
	Mandarin CORAL	AISHELL-3-test-left

Ignorant: the attackers have access to the ASV_{eval} and a few unanonymized test trials from the users. With ASV_{eval} , the attackers try to recognize the speaker identity by matching the test trials with the unanonymized data (i.e., enrollment data). However, they are unaware that the test trials have been anonymized;

Lazy-informed: the attackers have the same resources as in *Ignorant*. Besides these, they know that the test trials have been anonymized using the anonymization algorithm and CORAL transformation, but are unsure of the detailed parameters (e.g., which speaker vectors are selected to compose the anonymized vector). The attackers anonymize the enrollment data with their knowledge of SAS and use them to recognize the test trials.

This study included another two setups as references: **Unprotected:** the attackers directly recognize the unanonymized test trials using the original enrollment data and ASV_{eval} ; **Resynthesized:** similar to *Unprotected*, but the test trials are resynthesized by the SAS using the original speaker vector. Both setups simulate the case where speaker identity is unprotected, but *Resynthesized* further examines how the feature extractors and vocoder degrade the utility and privacy in the resynthesis process [18].

An ideal SAS should achieve high EERs (close to 50%) in both *Ignorant* and *Lazy-informed* setups. The EERs in *Resynthesized* should be as low as those from *Unprotected*. Meanwhile, the CERs for all setups should be low.

4.2. Databases

In addition to the standard databases to build the SAS (Section 2), this study used external data listed in Table 1 for the SAS training and CORAL transformation. The multilingual dataset for HiFi-GAN consists of *LibriTTS-train-clean-100* and subsets that contains *German*, *Italian*, and *Spanish* data sampled from the Multilingual LibriSpeech corpus [20]. Around 78 hours of clean data were selected for each language. The selection criteria is that the duration of each recording is larger than 10s, and that the signal-to-noise ratio is equal to 100, which is estimated using the WADA SNR algorithm [24].

The Mandarin test trials and enrollment data were sampled from the test set of Mandarin speech corpus *AISHELL-3* [25]. Specifically, 4,179 trials from 44 speakers were randomly sampled as the test trials, and additional 2 utterances of the same speaker were sampled for enrollment. They composed 10,120 enrollment-test pairs for the ASV evaluation, which is denoted as *AISHELL-3-test-veri*. The target-domain vectors for the oracle Mandarin CORAL were extracted from the left data of *AISHELL-3* test set, denoted *AISHELL-3-test-left* in Table 1. Note that there is no speaker or utterance overlap between *AISHELL-3-test-left* and *AISHELL-3-test-veri*. The target-domain vectors for the general CORAL were randomly selected from the *German*, *Italian*, and *Spanish* subsets. The source-domain vectors for both types of CORAL were randomly selected from the *LibriTTS-train-clean-100*. Similar to [8], the ASV_{eval} model was an ECAPA-TDNN trained on the Mandarin *CN-Celeb-1 & 2* [26, 27] datasets. The ASR_{eval} model was an open-source ASR Transformer [28] trained on the Mandarin *AISHELL-1* ASR dataset [29].

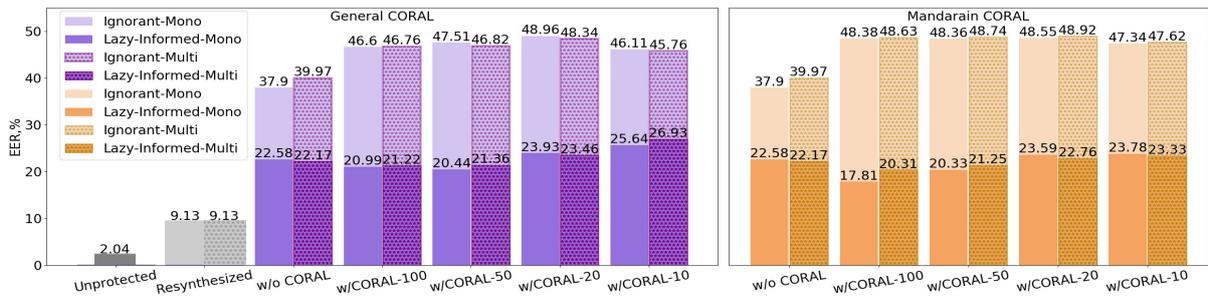


Figure 3: EER values for two anonymization systems using monolingual HiFi-GAN or multilingual HiFi-GAN along with different configurations of the general and Mandarin CORAL transformation.

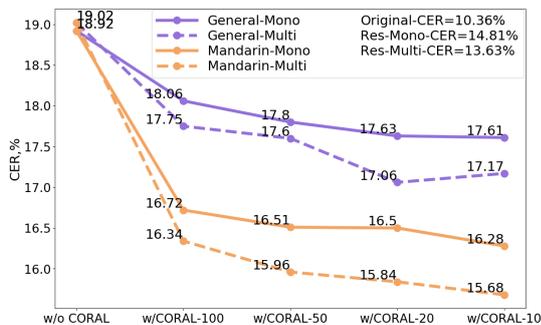


Figure 4: CERs of anonymized speech on AISHELL-3-test-veri.

4.3. Results and analysis

Our experiments compared the SAS performance on the Mandarin test data by varying the two factors illustrated in Figure 1(b): using monolingual or multilingual data to train the HiFi-GAN vocoder, and whether CORAL was applied on the anonymized speaker vector. The EERs from different setups are plotted in Figure 3. The CERs are plotted in Figure 4. Note that, when CORAL was used, we also analyzed the impact of using different amounts of data to estimate the CORAL transformation matrix. The label “w/CORAL- N ” in the figures indicates using N randomly-chosen speaker vectors from the source and target domain respectively. All the results for CORAL transformation are computed over 5 runs of speaker vector selection.

Monolingual vs. Multilingual HiFi-GAN: In the *Resynthesized* setup, the SAS using multilingual training data for HiFi-GAN achieved a CER of 13.63%, which is lower than 14.81% when using monolingual data. Meanwhile, the EERs in both cases are equal to 9.13%. When anonymization was conducted, the SASs using multilingual HiFi-GAN achieved lower CERs than their counterparts using monolingual HiFi-GAN in most of the setups except when no CORAL was used (w/o CORAL, 19.02% > 18.92%). Specifically, when the Mandarin CORAL was used, the multilingual HiFi-GAN outperformed the monolingual case in terms of CER, no matter how much data was used to estimate the CORAL matrix.

These results indicate that the multilingual data is helpful to obtain a robust HiFi-GAN to preserve the speech contents better. The improvement is expected to be larger when the domain mismatch on the anonymized speaker vectors is reduced using CORAL. The EERs were roughly similar no matter which HiFi-GAN was used. Therefore, the benefit of using the multilingual HiFi-GAN is mainly on the better preservation of the speech contents, rather than better protection of the speaker identity.

CORAL: Compared to “w/o CORAL”, all the EERs on the *Ignorant* condition of “w/CORAL- N ” were successfully increased regardless of the choice of HiFi-GAN training data, CORAL types, and amount of data for CORAL matrix estimation. Similarly, the CERs were significantly decreased after applying CORAL. These results suggest that the mismatch from the

Table 2: EER values on Lazy-informed condition for the SAS using Multilingual HiFi-GAN and different CORAL configurations.

EER(%)	w/o CORAL	CORAL-10	CORAL-10-100
General	22.17	26.93	25.93
Mandarin	22.17	23.33	21.66

anonymized speaker vectors severely affected the SAS, and CORAL is effective to reduce various types of mismatches (e.g., language and channel) between different datasets/domains.

For the different configurations for CORAL, we first observed that the oracle Mandarin CORAL performed better on CERs than the general CORAL. However, their differences on EERs are not obvious. This indicates that the SAS performance on speech content preservation is more sensitive to the mismatch of anonymized vectors than the CORAL configurations. For speaker identity protection, using the general CORAL is sufficient. Interestingly, unlike DNN-based methods, using larger N to estimate the CORAL matrix did not constantly improve the results. Using 20 samples (w/CORAL-20) generally performed well in terms of EER and CER. The reason is that users and attackers randomly choose speaker vectors from the target domain individually to approximate CORAL transformation matrices. These matrices can be very different if the number of the speaker vectors N is relatively small, which increases the randomness of the new anonymized speaker vectors used by the users and attackers. Therefore, speaker identity information can be protected better.

Considering that the users may prefer to choose smaller N to protect their privacy, while an attacker may be interested to use larger N to find a more precise CORAL transform matrix on the *Lazy-informed* condition, we then set $N = 10$ for users and $N = 100$ for attackers to compute the CORAL matrix independently, denoted CORAL-10-100. The results from Table 2 show that the EERs of CORAL-10-100 are lower than those of CORAL-10 for both general and Mandarin cases in an acceptable level. Furthermore, “General CORAL-100-10” still performs better than “w/o CORAL”.

5. Conclusions

This paper analyzed the previously proposed SSL-based SAS under unseen conditions. Two hypotheses, which are that the performance bottleneck exists in the HiFi-GAN and in anonymized speaker vectors, were presented and experimentally verified. The results indicate that increasing the language diversity for the HiFi-GAN benefits the preservation of speech contents. The mismatch on the anonymized speaker vectors severely affect the SAS. The SAS using multilingual HiFi-GAN and CORAL strategy easily outperforms the previous SAS using monolingual HiFi-GAN on both privacy and utility.

Acknowledgements This study is supported by JST CREST Grants (JPMJCR18A6 and JPMJCR20D3), MEXT KAKENHI Grants (21K17775, 21H04906, 21K11951, 18H04112), and the VoicePersonal project (ANR-18-JSTS-0001).

6. References

- [1] A. Nautsch, C. Jasserand, E. Kindt, M. Todisco, I. Trancoso, and N. Evans, “The GDPR & Speech Data: Reflections of Legal and Technology Communities, First Steps Towards a Common Understanding,” in *Proc. Interspeech*, 2019, pp. 3695–3699.
- [2] 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 *Proc. Interspeech*, 2020, pp. 1693–1697.
- [3] N. Tomashenko, X. Wang, E. Vincent, J. Patino, B. M. L. Srivastava, P.-G. Noé, A. Nautsch, N. Evans, J. Yamagishi, B. O’Brien *et al.*, “The VoicePrivacy 2020 challenge: Results and findings,” *Computer Speech & Language*, 2022.
- [4] Q. Jin, A. R. Toth, T. Schultz, and A. W. Black, “Speaker deidentification via voice transformation,” in *2009 IEEE Workshop on Automatic Speech Recognition & Understanding*. IEEE, 2009, pp. 529–533.
- [5] F. Fang, X. Wang, J. Yamagishi, I. Echizen, M. Todisco, N. Evans, and J.-F. Bonastre, “Speaker anonymization using x-vector and neural waveform models,” *Proc. 10th ISCA Speech Synthesis Workshop*, pp. 155–160, 9 2019.
- [6] J. Patino, N. Tomashenko, M. Todisco, A. Nautsch, and N. Evans, “Speaker Anonymisation Using the McAdams Coefficient,” in *Proc. Interspeech*, 2021, pp. 1099–1103.
- [7] B. M. L. Srivastava, N. Vauquier, M. Sahidullah, A. Bellet, M. Tommasi, and E. Vincent, “Evaluating voice conversion-based privacy protection against informed attackers,” in *Proc. ICASSP*. IEEE, 2020, pp. 2802–2806.
- [8] X. Miao, X. Wang, E. Cooper, J. Yamagishi, and N. Tomashenko, “Language-independent speaker anonymization approach using self-supervised pre-trained models,” *arXiv preprint arXiv:2202.13097*, 2022.
- [9] B. Sun, J. Feng, and K. Saenko, “Return of frustratingly easy domain adaptation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016.
- [10] K. Kasi and S. A. Zahorian, “Yet another algorithm for pitch tracking,” in *Proc. ICASSP*, vol. 1, 2002, pp. 1–361.
- [11] A. Nagrani, J. S. Chung, and A. Zisserman, “Voxceleb: A large-scale speaker identification dataset,” in *Proc. Interspeech*, 2017, pp. 2616–2620.
- [12] J. S. Chung, A. Nagrani, and A. Zisserman, “Voxceleb2: Deep speaker recognition,” in *Proc. Interspeech*, 2018, pp. 1086–1090.
- [13] B. van Niekerk, M.-A. Carboneau, J. Zaïdi, M. Baas, H. Seuté, and H. Kamper, “A comparison of discrete and soft speech units for improved voice conversion,” *arXiv preprint arXiv:2111.02392*, 2021.
- [14] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, “HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [15] H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, and Y. Wu, “LibriTTS: A corpus derived from LibriSpeech for text-to-speech,” *arXiv preprint arXiv:1904.02882*, 2019.
- [16] J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” 2020.
- [17] B. M. L. Srivastava, N. A. Tomashenko, X. Wang, E. Vincent, J. Yamagishi, M. Maouche, A. Bellet, and M. Tommasi, “Design choices for x-vector based speaker anonymization,” in *Proc. Interspeech*, 2020, pp. 1713–1717.
- [18] P. Champion, D. Jouvét, and A. Larcher, “Speaker information modification in the VoicePrivacy 2020 toolchain,” INRIA Nancy, équipe Multispeech ; LIUM - Laboratoire d’Informatique de l’Université du Mans, Research Report, nov 2020. [Online]. Available: <https://hal.archives-ouvertes.fr/hal-02995855>
- [19] J. Lorenzo-Trueba, T. Drugman, J. Latorre, T. Merritt, B. Putrycz, R. Barra-Chicote, A. Moinet, and V. Aggarwal, “Towards Achieving Robust Universal Neural Vocoding,” in *Proc. Interspeech*, 2019, pp. 181–185.
- [20] V. Pratap, Q. Xu, A. Sriram, G. Synnaeve, and R. Collobert, “MLS: A Large-Scale Multilingual Dataset for Speech Research,” in *Proc. Interspeech 2020*, 2020, pp. 2757–2761.
- [21] D. Raj, D. Snyder, D. Povey, and S. Khudanpur, “Probing the information encoded in x-vectors,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 726–733.
- [22] J. Williams and S. King, “Disentangling style factors from speaker representations,” pp. 3945–3949, 2019.
- [23] J.-F. Cai, E. J. Candès, and Z. Shen, “A singular value thresholding algorithm for matrix completion,” *SIAM Journal on optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.
- [24] C. Kim and R. M. Stern, “Robust signal-to-noise ratio estimation based on waveform amplitude distribution analysis,” in *Ninth Annual Conference of the International Speech Communication Association*, 2008.
- [25] Y. Shi, H. Bu, X. Xu, S. Zhang, and M. Li, “AISHELL-3: A Multi-Speaker Mandarin TTS Corpus,” in *Proc. Interspeech*, 2021, pp. 2756–2760.
- [26] L. Li, R. Liu, J. Kang, Y. Fan, H. Cui, Y. Cai, R. Vippera, T. F. Zheng, and D. Wang, “CN-Celeb: multi-genre speaker recognition,” *Speech Communication*, 2022.
- [27] Y. Fan, J. Kang, L. Li, K. Li, H. Chen, S. Cheng, P. Zhang, Z. Zhou, Y. Cai, and D. Wang, “CN-Celeb: a challenging Chinese speaker recognition dataset,” in *Proc. ICASSP*. IEEE, 2020, pp. 7604–7608.
- [28] M. Ravanelli, T. Parcollet, P. Plantinga, A. Rouhe, S. Cornell, L. Lugosch, C. Subakan, N. Dawalatabad, A. Heba, J. Zhong, J.-C. Chou, S.-L. Yeh, S.-W. Fu, C.-F. Liao, E. Rastorgueva, F. Grondin, W. Aris, H. Na, Y. Gao, R. D. Mori, and Y. Bengio, “SpeechBrain: A general-purpose speech toolkit,” 2021, arXiv:2106.04624.
- [29] H. Bu, J. Du, X. Na, B. Wu, and H. Zheng, “Aishell-1: An open-source Mandarin speech corpus and a speech recognition baseline,” in *2017 20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment (O-COCOSDA)*. IEEE, 2017, pp. 1–5.