

Hiding Sensitive Information in Desensitized Voice Sequences*

*Note: Sub-titles are not captured in Xplore and should not be used

1st Given Name Surname, 2nd Given Name Surname, 3rd Given Name Surname
dept. name of organization (of Aff.), City, Country
dept. name of organization (of Aff.), City, Country if needed
 email address or ORCID of corresponding author(s)

Abstract—Voice data is acquired and utilized by a variety of consumer services. During the computation of these services, most of the user’s raw records are sent to the web server associated with them. In this way, serious privacy problems are bound to be caused, because the raw voice contains voiceprint and content information, which is considered sufficient to identify the speaker or various sensitive information related to the speaker. In order to mitigate such problem, Desensitization methods are employed as an intermediary between user and the cloud services. However, if these methods are provided by a third party as a black box, it may not be safe enough. In this paper, we discuss and demonstrate the possibility of hiding information sufficient to recover original voice from in seemingly desensitized voices that may be used for online services with deep learning technologies, utilizing StarGAN-based voice transformation and voice-optimized audio stenography technologies.

Index Terms—privacy, voice, desensitization, stenography

I. INTRODUCTION

Since the last century, voice has become one of the important means of human-computer interaction. Voice recognition technology enables this interaction.

The number of speech recognition technologies in existence is very large. Among them, voiceprint extraction technology enables computers to use only certain features in human voice for authorization, without involving traditional password-based methods, thereby avoiding security concerns with weak passwords or those been in use for a long period. Also the need for users to remember a robust and secure sequence of password is eliminated. (描述一下传统的方法有什么问题). Some voice assistant apps on smartphones support voiceprint-based security measures to prevent accidental (or malicious) activation by anyone other than the smartphone owner, such as Siri from Apple Inc. and Xiao Ai from Xiaomi Inc. [19], [20]. Also, there are many voice-typing enabled input methods for various kinds of devices, such as iFly Input Method from iFlyTek and GBoard from Google, Inc. [21], [22].

In recent years, deep learning and deep learning-based speech analysis techniques have made great progress. This also leads to the fact that humans cannot easily distinguish the speech sequences produced by the generation-oriented services, and even the human ear cannot recognize the most subtle features in the speech sequences provided by the

discrimination-oriented services. [12]. At the same time, the advancement of this speech analysis technology also warns that people will be vulnerable to high-precision recognition [6], [7]. Because voice contains a lot of sensitive personal information, if users’ voice data is processed by companies with insufficient security infrastructure or even malicious companies, it is likely to lead to a privacy crisis.

If the original recording of the user’s voice, which contains too much private information, is completely processed locally on the user’s device (i.e., from recording to recognition), the security risk is not great because only the unrecognized results are exported or sent to the remote server.

However, these recognition services consume a lot of computing power, and it is almost impossible to execute them locally on the user. Therefore, service providers process this data on their servers. Although some vendors claim to process data locally, it is still unavoidable that part of the data will be processed in the server, which obviously does not meet the security requirements. For example, Xiaomi claims that Xiao Ai can “do most of the training and evaluation locally” [20].

Ideally, an honest service provider would not use the raw recordings for anything other than the intended purpose, as depicted in Fig. 1. However, in reality this imposes a serious privacy concern since the raw recording from user is known to contain much more data than merely for content recognition. Such additional data may be used by malicious service providers or attackers to identify the speaker and cause a data breach incident, as depicted in Fig. 2. To address this type of problem, many speech desensitization algorithms are being developed, the general workflow of which is shown in Figure 3(图片描述不够细致).

(这里开始到72行这些脱敏方法的介绍尽量简短一些因为介绍他们主要是为了引出本文的攻击方法) Conventional voice desensitization methods are believed to focus on two radically different areas: Content and Voiceprint.

Those which focus on content desensitization employ certain methodologies to remove or replace voice segments that are detected to contain sensitive information. There are publicly available APIs [13] and dedicated softwares for this purpose [14]. As a safe method or the last resort, some may employ manual audio editing to achieve this goal, and there exists a number of softwares sufficient for this use case [15].

Other methods may apply distortion on entire voice se-

quences. Deep learning based methodologies are employed in both types of methods along with traditional algorithms [9]–[11].

It's notable that many approaches to bring privacy to remotely handled voices combine these 2 ideas to provide better performance. There exists a number of active researches on this subject [8]. (下面可以简短一些: 脱敏语音中没有敏感数据是不现实的, 本文就针对脱敏模型的弱点...做出了...攻击) (将别人的工作、算法放到related work去介绍)

It's unlikely to guarantee the absence of sensitive data in desensitized voices, primarily due to the weakness of current desensitization. However, the desensitization method itself could also be modified to contain adversarial features without drastically degradation of on-paper performance [1], as depicted in Fig. 4. Targeting this particular case, we performs attack on a desensitization model in this paper.

In summary, this paper makes the following contributions.

- We demonstrate the adversarial case: Hiding sensitive data in desensitized voice sequences.
- We present its concrete workflow: a novel, adversarial exploitation of voice desensitization frameworks that.
- We conduct experiments on a particular voice transformation model with well-known voice samples. From them, we analyzed the performance and potential of this attack framework.

(将下面这段分为几点填到上面item里面概括一下主要贡献) The rest of the paper consists related works, an detailed explanation of the proposed methodology, including the discussion of potential use cases and opportunities of future works.

II. RELATED WORK

(查阅一些对于脱敏算法相关的攻击方法补充一些)

A. Adversary against desensitization

[PPRL-VGAN, Subverting] There are researches on attacking privacy-preserving data transformation models. Some employ similar techniques that attempt to embed certain amount of data in sanitized data with modification to original deep-learning models and recover the original data after the exposure of sanitized data in public by victims.

An notable instance among them works with images, targeting a privacy-preserving facial expression recognition algorithm, PPRL-VGAN. It sets up the attack with weak assumptions of user, who have white-box access to the attacked model. In order to achieve the adversary, its adversarial parts are embedded in the original model as additional layers or modification of existing layers, thus avoiding user's discovery.

B. Speech anonymization

III. PROBLEM STATEMENT

将Methodology的一些内容放到这个章节(攻击、攻击者的设定是怎么样的? 黑盒白盒? 攻击者能获得什么模型什么信息? 为了完成攻击做出了什么假设? 考虑的博弈模型的是怎么样的)

In this adversary, the attacker would like to distribute a modified black-box desensitization service that process the raw voice sequences solely on user's device, which implies that the attacker is not able to obtain the original voice sequences. With the fact that only the desensitized voice is available, attacker would like the service to output sequences that contain adequate information for it to recover original ones from.

Our proposed adversary framework primarily works as an add-on to a benign (non-adversarial) voice desensitization model. The use case, as mentioned before, is depicted by Fig. 4, in which the user chose to use a black-box that disguised itself to be a normal desensitization service, but embedded additional data on top of the desensitized voice sequence for adversarial purposes.

For users, before or after the adversary, the desensitization model remains as a black box, which implies that users are not able to distinguish them by their internal structure.

For the attacker, there are slightly different cases.

It's possible that the attacker is the owner of the desensitization model, or obtained a publicly available model. In this case, the model is a white-box, which indicates that the attacker is able to better integrate the stenography part into the black-box presented to users by either training or fine-tuning the stenography model along with the benign model or embedding the model into the benign one, as additional layers.

It's also possible that the attacker has only black-box access to the benign model. In this case, the workflow remain unchanged except for that the attacker have no further control of the integration.

To conclude, the user treating the adversarial workflow as a desensitization service wants to minimize the possibility of identification after the desensitization process while the attacker wants to maximize the quality of extracted data while fulfilling the user's original demands.

IV. METHODOLOGY

(这个章节中只需要描述攻击方法、组件、定义、算法等等信息)

In order to complete such workflow, three distinct models are needed as depicted in Fig. 5. They are explained in detail below, but they can also be summarized as following:

- A **Benign model** to be attacked by. This model generates desensitized voice sequences that users would normally want to get from this black-box.
- A **Stenography model** to embed features of original voice sequences on the desensitized ones. The output sequences should trick users into believing it's clean desensitized ones.
- A **Extraction model** to extract the features that **Stenography model** hid in the voice sequences and recover original ones from them.

(我觉得Methodology一直到这个位置可以拆分放到上一个章节只留下下面的三个模型具体的介绍和攻击方法1)

A. Benign model

As two slightly different scenarios, the benign model could be either a in-house solution designed by the attacker itself

(and may be even trained with adversarial models) or a production model with predefined weights acquired from a third-party organization or the public. The source of benign model theoretically shows no significant difference in terms of the overall structure of the adversary workflow we designed, but the desensitization performance of different models may have an impact on the quality of stenography-added voices and recovered ones. For our experiments we used an StarGAN-VC [3] based model, the StarGanVCDialectConversion project hosted by user "Didnelpsun" on GitHub [2] in particular.

Specifically, this specific StarGAN-VC implementation employs a transformer-based, multi-layered architecture for its generation model, which closely resemble the approach mentioned in the original StarGAN-VC publication. There are 5 sets of layers in both the encoder and decoder with each set of layer consisting of two convolution layers and two instance normalization layers. It's notable that this particular model is not strictly designed for desensitization, but we consider the fact that this model is capable to erase the original identity from the voice sequence and give it another one to be the ability to desensitize.

B. Stenography model

The stenography model works as an intermediate between the user and the benign model to embed information of original voice into clean product of the benign model. It's clear that, users of this adversarial product are not intended to notice the presence of such stenography model that attempts to tamper the clean output.

Details are discussed in the following sections. However, the characteristics a stenography model needed to be adequate for this purpose can be summarized as follows:

- **Integration with benign model:** The stenography model should not alienate the adversary workflow by requiring additional data, other than what is provided inside the black-box (the original voice sample and the desensitized one) or be able to be detected despite being inside the black-box.
- **Sufficient performance:** The stenography model should embed sufficient information for the extraction model to recover the original voice sequence, thus ensuring the possibility of the adversary. But it should also avoid excess modifications applied to the samples to reduce suspicion.

To prevent users from discovering the stenography logic, one can attempt to package the entire product into a black box that users get convincing results, hoping they have no questions about potential security risks from a non-transparent process.

However, such actions may rise suspicion among users, which may not be desired in some cases. Another robust solution to this problem is to embed the stenography model into the victim, as an additional set of layers for "post-processing", thus integrated with it. In this way, the presence of stenography model could no longer be easily detected without dedicated research of behaviour or, even, source code.

Also, the benign desensitization model attempts to replace a certain amount of features of original voices with features that does not belong to the owner of them while the stenography model attempts to embed a certain amount of information of original features into the already-desensitized voices.

It's very likely that during this process the textual content of the voice or the overall audio quality will experience degradation to a degree. For this reason, the stenography model chosen for this purpose should be high quality enough to avoid excessively modifying the voice output.

There exist a number of solutions for this purpose [5]. In order to achieve such a task to hide sufficient information in desensitized ones while minimizing the differences made, we employed the Hide and Speak model [4], an voice-centralized stenography model that could be considered as state of the art as it handles different lengths of carrier data and target data and produce high quality results.

C. Extraction model

The extraction model accepts the final product from the combination of benign and stenography model, which is likely to be the voices user trusted to be desensitized, and attempts recovery of original voice from it. As a model trained in conjunction with the stenography model, the one we employed in our experiment is also from the Hide and Speak mentioned before.

V. EXPERIMENT

A. Setup

As mentioned before, we use StarGanVCDialectConversion, a StarGAN-VC implementation as the benign model to produce desensitized voices, Hide and Speak for stenography and extraction process. All our experiments are performed on a x86-64 based container-enabled Linux server with a NVIDIA Tesla V100 GPU. We used software packages required by each models, but with Python 1.10 and PyTorch 1.13, which are slightly newer.

For convenience of further processes, we used voice samples from TM1, TM2, SF1 and SF2, speakers in the VCC2016 data set [16], a well-known data set that is used by numerous projects and is versatile [17], [18], as preferred by the configuration of benign model. The voices of TM1 is considered the original voices that contain sensitive information and ones of SF1, on the other hand, is considered desensitized. With this setup in mind, the experiment can be described as: Voices of TM1 are transformed into ones of SF1 by the benign model, and the stenography model takes the output and embed voice data of TM1, generating the final output. Extraction model then use the final output to recover the voices of TM1.

B. Training

We trained the models locally on the server. In order to streamline the experiment process while preserving the most accurate possible results, we avoid excess modifications to the models. Particularly, the benign model is trained to 200000 steps (60 epochs), as the default settings. Similar approach is applied on the stenography model, that it is trained to preferred settings by the model authors.

C. Generation

It's necessary to clarify that, in order to achieve the maximum quantity of samples and prevent the quality of benign model from having excess impact on our overall process, we use the whole VCC2016 training set as the original voices. Each identity contains 162 samples, and we get 648 samples in total.

As a StarGAN-based model, our benign model is capable to transform a voice sequence from any known speakers to another. We consider the fact that the transformed voice belongs to other identity than the original speaker to be a form of desensitization, as mentioned before. Since each sample can be used to generate 3 sequences targeting different speakers, we get 1944 in total.

The remaining steps are straightforward. We used each sample generated by benign model to get a stenography one with the original sample associated with it. Finally we used the extraction model to recover each voices.

D. Evaluation

In order to generate creditable numbers for our voice samples that correctly reflect the amount of sensitive data, the identity of original speaker, we employed cloud-based voiceprint analysis solution provided by iFlyTek. After learning about the identity of four speakers in our domain, this service is able to generate 4 numbers for each voice sample, denoting the probability of the voice to belong to a particular speaker. iFlyTek suggests in the official manual that a score that is higher than 0.6 meant that the identity of a sample can be confirmed. We consider the identity with the highest score for a particular sample to be the classification made by this service.

VI. RESULTS

A. Explanation

Table I to IV presents the essential statistics of our results. All of the statistics are based on the score of each sample being classified into the original identity of their own. A higher score mean that our evaluation service, as mentioned before, deduces that the sample has a higher probability to belong to its original identity.

Each row contains statistics targeting the original identity of samples where their original identity is the speaker denoted in the header. Each column consists different types of results of statistical computation according to the following description:

- **Mean** - Average value of scores. Higher values mean generally closer to the original speaker.
- **Definiteness** - Ratio of scores being greater than 0.6, which means the sample can be confirmed to bear the same speaker as the original one. Higher values mean values are more definite.
- **Best** - Best value of scores.
- **Worst** - Worst value of scores. Closer value with **Best** means better stability.
- **Class Ratio** - Ratio of samples being classified into its original speaker, i.e. have the highest score being

the original speaker, as mentioned before. Higher values mean higher probability an generic classification model will think the samples have the same identity of the original speakers.

B. Analysis

What is depicted in Fig. 6 is 4 sets of voice samples in different stages of processing, original(**ori**), desensitized(**des**), stenography(**msg**) and extracted(**rec**). From which it's noticeable that the overall performance of this adversary is reasonable.

1) *Original samples*: According to Table I, all of the score statistics, except for the **Worst**, are close to 1, which is intended behavior for original samples. As the worst case, **Worst** is also close or greater than 0.6, which means that even this type of cases are confirmed to have the same identity as the original speaker. As a result, it's safe to confirm that both the benign model and the validation service are of desirable performance.

2) *Desensitized samples*: According to Table II, the overall score dropped drastically from over 0.8 to over 0.4. It can be argued that from the **Best** and **Class Ratio** numbers that there still exhibits a portion of samples that are classified into the original speaker. However, these type of classification results can not be trusted because they are vague, according to the close-to-zero **Definiteness** value.

3) *Stenography samples*: Scores of stenography samples are similar to the desensitized ones, given the numerical changes of statistical numbers are mostly less than 0.1. However, the changes stenography model made to the samples did not cause any form of degradation of desensitization performance, but upgraded it instead.

Table V is the statistics of direct comparison between these two sets of samples. According to the **Minimum** value, it's still possible that the desensitization performance would suffer significant degradation from the stenography process, but from the **Mean** and **Variance** value we can see that the performance changes are even-spread.

It's highly likely that the changes are caused by the slight content degradation, i.e. perturbation, made by stenography model and would not strongly affect the overall performance of such adversary. We also performed manual audio quality tests on a random subset of this set of samples and confirmed that the degradation is not audible.

4) *Extracted samples*: As mentioned before, extracted samples are meant to be as close to the original samples as possible. As presented in Table IV, the values are slightly inferior than the original samples with the decrease of score within 0.1 to 0.2. However, despite the worsen results, the **Definiteness** and **Class Ratio** are still well desirable, suggesting that the usability of these samples are comparable to the original ones, which declares the success of adversary.

VII. DISCUSSION

According to our experiments, it is safe to consider conventional acoustical-based or NN-based audio transforming solutions not sufficient for voice desensitization. Besides the

TABLE I
STATISTICS OF ORIGINAL SAMPLES

	SF1	SF2	TM1	TM2
Mean	0.853889	0.837284	0.84537	0.837963
Definiteness	1.0	0.987654	1.0	1.0
Best	0.95	0.94	0.94	0.94
Worst	0.6	0.56	0.6	0.61
Class Ratio	1.0	1.0	1.0	1.0

TABLE II
STATISTICS OF DESENSITIZED SAMPLES

	SF1	SF2	TM1	TM2
Mean	0.407654	0.440206	0.460556	0.417078
Definiteness	0	0.022634	0.047325	0.004115
Best	0.58	0.64	0.62	0.62
Worst	0.22	0.23	0.27	0.21
Class Ratio	0.125514	0.236626	0.195473	0.012346

StarGAN-VC solution we used, there exists many more such "voice changer" services on the Internet available for public use. One would consider these solutions secure because of the vast audible differences they made on its voice sequences. However, these solution exhibits potential of adversary with such method we demonstrated in this paper, which is not negligible. It's not likely that human ears can pick up subtle changes a stenography program made to certain parts of a voice sequence.

It's possible to avoid or mitigate risks of being attacked by such method. The most straightforward way to go is to avoid desensitization models from unknown or unsound sources. Due to the black-box nature of various proprietary services, it is not likely possible for users to have practical means to test these services for potential adversaries. When possible, users could train their own desensitization models. If the models are acquired from third-party, users should pay attention to the behaviour of model and the choice of training data set.

VIII. FUTURE WORK

In our experiment, we simply directed the voice sequence generated by the benign model to the stenography model. As mentioned before, this approach may not be sufficient in

TABLE III
STATISTICS OF STENOGRAPHY SAMPLES

	SF1	SF2	TM1	TM2
Mean	0.387695	0.366070	0.412119	0.394115
Definiteness	0	0.002058	0	0
Best	0.52	0.61	0.56	0.56
Worst	0.21	0.15	0.17	0.20
Class Ratio	0.119342	0.183128	0.1893	0.014403

TABLE IV
STATISTICS OF EXTRACTED SAMPLES

	SF1	SF2	TM1	TM2
Mean	0.681975	0.681728	0.695123	0.702593
Definiteness	0.901235	0.864198	0.950617	0.938272
Best	0.81	0.81	0.81	0.85
Worst	0.44	0.47	0.53	0.46
Class Ratio	0.993827	1	1	0.993827

TABLE V
STATISTICS OF DIFFERENCES BETWEEN DESENSITIZED AND STENOGRAPHY SAMPLES

Mean	0.041373
Variance	0.010478
Maximum	0.37
Minimum	-0.29
Maximum Absolute Value	0.37
Minimum Absolute Value	0

terms of stealthiness. Also, the overall storage consumption of this black box will increase significantly and the processing performance may not be ideal.

We believed that turning the stenography model into extra layers of benign model, eliminating the redundant audio encoding and decoding processes, can mitigate such problems while potentially increase the overall adversarial quality. In this way, the benign model can train in conjunction with the adversary model, taking advantage of intermediate representations of it. Moreover, this creates the potential of reducing the overall storage consumption and be less suspicious.

As another way to mitigate such problem, users can use various traditional or ML-based methods to apply inaudible perturbation on processed voices to attempt erasure of potential stenography while preserving high audio quality.

IX. CONCLUSION

Targeting voice desensitization models based on generative NN models, we designed an adversary scheme that attempts to recover original voices from desensitized ones via stenography means, thus cause a privacy bleach. Our experiments prove this idea to be viable and the conventional solutions to be vulnerable to this type of adversary.

REFERENCES

- [1] N. Subramanian, O. Elharrouss, S. Al-Maadeed and A. Bouridane, "Image Steganography: A Review of the Recent Advances," in IEEE Access, vol. 9, pp. 23409-23423, 2021, doi: 10.1109/ACCESS.2021.3053998.
- [2] <https://github.com/DidneIpsun/StarGanVCDialectConversion>
- [3] Y. Li, X. Qiu, P. Cao, Y. Zhang, and B. Bao, "Non-parallel Voice Conversion Based on Perceptual Star Generative Adversarial Network," Circuits Syst Signal Process, vol. 41, no. 8, pp. 4632-4648, Aug. 2022, doi: 10.1007/s00034-022-01998-5.
- [4] F. Kreuk, Y. Adi, B. Raj, R. Singh, and J. Keshet, "Hide and Speak: Towards Deep Neural Networks for Speech Steganography," arXiv, Jul. 27, 2020. doi: 10.48550/arXiv.1902.03083.
- [5] N. Takahashi, M. K. Singh, and Y. Mitsufuji, "Source Mixing and Separation Robust Audio Steganography," arXiv, Feb. 17, 2022. doi: 10.48550/arXiv.2110.05054.
- [6] Andreas Nautsch, Abelino Jiménez, Amos Treiber, Jascha Kolberg, Catherine Jasserand, Els Kindt, Héctor Delgado, Massimiliano Todisco, Mohamed Amine Hmani, Aymen Mtibaa, Mohammed Ahmed Abdelraheem, Alberto Abad, Francisco Teixeira, Driss Matrouf, Marta Gomez-Barrero, Dijana Petrovska-Delacrétaz, Gérard Chollet, Nicholas Evans, Thomas Schneider, Jean-François Bonastre, Bhiksha Raj, Isabel Trancoso, and Christoph Busch. 2019. Preserving privacy in speaker and speech characterisation. Comput. Speech Lang. 58, C (Nov 2019), 441-480. <https://doi.org/10.1016/j.csl.2019.06.001>
- [7] Kröger, J.L., Lutz, O.H.M., Raschke, P. (2020). Privacy Implications of Voice and Speech Analysis – Information Disclosure by Inference. In: Friedewald, M., Önen, M., Lievens, E., Krenn, S., Fricker, S. (eds) Privacy and Identity Management. Data for Better Living: AI and Privacy. Privacy and Identity 2019. IFIP Advances in Information and Communication Technology(), vol 576. Springer, Cham.

- [8] Jaemin Lim, Kiyeon Kim, Hyunwoo Yu, and Suk-Bok Lee. 2022. Overo: Sharing Private Audio Recordings. In Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security (CCS '22). Association for Computing Machinery, New York, NY, USA, 1933–1946. <https://doi.org/10.1145/3548606.3560572>
- [9] J. Qian, H. Du, J. Hou, L. Chen, T. Jung and X. -Y. Li, "Speech Sanitizer: Speech Content Desensitization and Voice Anonymization," in IEEE Transactions on Dependable and Secure Computing, vol. 18, no. 6, pp. 2631–2642, 1 Nov.-Dec. 2021, doi: 10.1109/TDSC.2019.2960239.
- [10] Jianwei Qian, Haohua Du, Jiahui Hou, Linlin Chen, Taeho Jung, and Xiang-Yang Li. 2018. Hidebehind: Enjoy Voice Input with Voiceprint Unclonability and Anonymity. In Proceedings of the 16th ACM Conference on Embedded Networked Sensor Systems (SenSys '18). Association for Computing Machinery, New York, NY, USA, 82–94. <https://doi.org/10.1145/3274783.3274855>
- [11] J. Qian, F. Han, J. Hou, C. Zhang, Y. Wang and X. -Y. Li, "Towards Privacy-Preserving Speech Data Publishing," IEEE INFOCOM 2018 - IEEE Conference on Computer Communications, Honolulu, HI, USA, 2018, pp. 1079–1087, doi: 10.1109/INFOCOM.2018.8486250.
- [12] Nandwana, Mahesh Kumar, Julien van Hout, Mitchell McLaren, Allen R. Stauffer, Colleen Richey, Aaron D. Lawson and Martin Graciarena. "Robust Speaker Recognition from Distant Speech under Real Reverberant Environments Using Speaker Embeddings." Interspeech (2018).
- [13] <https://docs.aws.amazon.com/transcribe/latest/dg/pii-redaction.html>
- [14] Vidizmo - Automatic audio redaction software. <https://www.vidizmo.com/vidizmo-artificial-intelligence-solutions/redaction/>
- [15] Audacity - Open source audio software. <https://www.audacityteam.org/>
- [16] T. Toda, L.-H. Chen, D. Saito, F. Villavicencio, M. Wester, Z. Wu, J. Yamagishi, "The Voice Conversion Challenge 2016," Proc. INTERSPEECH, pp. 1632–1636, 2016.
- [17] M. Wester, Z. Wu, J. Yamagishi, "Analysis of the Voice Conversion Challenge 2016 Evaluation Results," Proc. INTERSPEECH, pp. 1637–1641, 2016.
- [18] M. Wester, Z. Wu, J. Yamagishi, "Multidimensional scaling of systems in the Voice Conversion Challenge 2016," Proc. SSW9, pp. 40–45, 2016.
- [19] Siri - Voice assistant software. <https://www.apple.com/siri/>
- [20] Xiao Ai - Voice assistant software. <https://xiaoi.mi.com/>
- [21] iFly Input Method - Chinese/English input method software. <https://srf.xunfei.cn>
- [22] GBoard - Multilingual input method software. <https://play.google.com/store/apps/details?id=com.google.android.inputmethod>