Interspeech ICASSP                     2
AI

The development of deep neural networks has been progressing rapidly and the evolution of speech recognition systems has been incredibly fast. The study aims to provide researchers with ideas on improving system security in light of the increasingly severe security issues.

In this project, we carefully studied the principles of speech recognition systems and researched all possible attack details. We summarized our findings in a review and proposed methods for improving the front-end and back-end of speech recognition systems. We expanded our research scope with a universal point of view Similar attacks can co-exist in speech-related systems, not just speech recognition systems. We also consider adversarial attacks as particular noise, then combining traditional speech enhancement, modeling, and post-processing methods in system development can sufficiently deal with this attack.
Top journals and conferences in the speech field accepted our achievements, such as Interspeech and ICASSP. Above two years of research achievement have been introduced into two books (ISBN 978-4-904020-26-5, ISBN 978-4-904020-28-9) by NICT and stored in the national library Kansai. These efforts are our contribution to ensuring the security and reliability of AI systems.

speech recognition adversarial attack privacy perserving deepfake detection spoken dialogue federated learning security

１．研究開始当初の背景

Nowadays, spoken dialogue systems have been widely used in our society. The automatic speech recognition (ASR) module, as these dialogue systems' most natural human-machine interface, is based on deep neural networks (DNNs). The development of deep neural networks has been progressing rapidly, from basic DNNs to end-to-end models, to the popular self-supervised learning models in recent years, and now to the revolutionary large language models (LLMs). However, DNN is known to be vulnerable to adversarial examples (or attacks). This is a severe problem.

２．研究の目的

The study of adversarial attacks on DNN-based speech recognition systems serves as a starting point for our research, aiming to provide speech development researchers with ideas on improving system security in light of the increasingly severe security issues related to speech-based systems.

３．研究の方法

**Current academia is focusing on designing individual purposed perfect adversarial attacks** on Automatic Speech Recognition (ASR) systems while still sounding normal to human listeners with the following focuses:

1. Black-box attacks: In black-box attacks, the attacker has limited or no access to the internals of the target model, such as its architecture or weights. They craft adversarial examples by exploiting the transferability of adversarial perturbations across models or by querying the target model and using its output to estimate its behavior.

2. Untargeted attacks: Adversarial attacks on speech systems can be categorized into targeted and untargeted attacks. Targeted attacks aim to make the system produce a specific, desired output, while untargeted attacks cause incorrect output without a specific target.

3. Real-world applicability: Researchers are increasingly focusing on adversarial attacks that can be effective in real-world settings, such as over-the-air attacks or attacks robust to various environmental conditions like background noise.

These works studying individual threats are very important contributions to the community. However, **other substantial threats co-exist in real-world scenarios and adversarial attacks**. These attacks can compromise the integrity of machine learning models and seriously affect critical applications, such as autonomous vehicles, healthcare, and security systems.

In our proposed method, we propose **a universal point of view on dealing with these threats**.

1. **Defense of the whole system pipeline and taking adversarial audio as noise**: The speech recognition module is not isolated in the industry system pipeline. Taking adversarial audio as a type of noise, we confirm that three traditional methods can effectively deal with adversarial attacks: 1. Changing/detecting front-end input signal (denoising, adding noise, or detect-then-reject), 2. Changing models (tuning structure or retraining), and 3. Correct by back-end language models. All these methods exist in current speech systems.

2. **Extend to concerning the robustness of other speech-related models**: Similar to the adversarial attack, other substantial threats exist. **(1) Deepfakes** are manipulated media, particularly images and videos, created using advanced

machine-learning techniques. Deepfakes can spread misinformation, manipulate public opinion, and harm reputations, posing a severe threat to individuals, organizations, and society. **(2) Privacy leakage** refers to unauthorized access or disclosure of sensitive personal information, often resulting from inadequate data protection measures or malicious activities. Privacy leakage can lead to identity theft, financial loss, and damage to an individual's reputation. **(3) Adversarial attacks to other speech models:** the same immediate impact as deepfakes or privacy leakage, adversarial attacks on speaker recognition systems can still be a significant threat.

４．研究成果

In the first year of this project, we carefully studied the principles of speech recognition systems and researched all possible attack details. We summarized our findings in a review and proposed methods for improving the front-end and back-end of speech recognition systems. Our achievements were accepted by top conferences in the speech field, such as Interspeech and ICASSP, as followings:

- We construct speech recognition systems with recent popular training toolkits and neural network types (accepted in Journals and conferences, e.g., ICASSP2022).
- We did surveys for the current attack methods (accepted in Journal).
- We implement robust adversarial attacks using the Kaldi-based ASR systems (accepted in SLT2020).
- We are also happy to see that this framework can be used to protect sensitive speech content (accepted in LREC 2022).
- To defend against attacks, we find that adversarial audios are very sensitive. Moreover, the feature of its spectrogram is very different from the human voice, and it can be treated as a special kind of noise. We constructed speech enhancement systems and studied their mechanism this year (accepted in Journals and conferences, e.g., ICASSP2022).

In the project's second year, we discovered that attacks on speech recognition systems are not only limited to the content of speech recognition but also involve more attacks related to speaker attributes. Therefore, we expanded our research scope with a universal point of view.

- We continued research on the front-end method about synthesized speech quality estimation, signal-to-noise ratio estimation framework (accepted in conferences, such as APSIPA2022, EUSIPCO2022, and Interspeech2022), and back-end of the language model, speech recognition/translation systems and other down-streaming systems (accepted in Journals and conferences, e.g., JNLP, International Journal of Asian Language Processing, Interspeech2022, ICASSP2023).
- We participated in several privacy attack competitions, proposing our deepfake detection (accepted by Interspeech2022). Moreover, we published our first database following security protection law (O-COCOSDA2022).
- At the same time, we expanded our research from speech recognition systems to spoken dialogue systems (accepted in Sigdial2022).
- We also studied the state-of-the-art modeling training method, self-supervised training (accepted in Interspeech2022), and federated learning (accepted in ICASSP2023).

Above two years of research achievement have been introduced into two books (ISBN: 978-4-904020-26-5, ISBN: 978-4-904020-28-9) by NICT and stored in the national library Kansai.

These efforts have laid the foundation for further research to ensure the security and reliability of AI systems.

Kai Li, Xugang Lu, Masato Akagi, Jianwu Dang, Sheng Li, Masashi Unoki

Relationship Between Speakers' Physiological Structure and Acoustic Speech Signals: Data-Driven Study Based on Frequency-Wise Attentional Neural Network

30th European Signal Processing Conference (EUSIPCO)

2022

Kak Soky, Sheng Li, Masato Mimura, Chenhui Chu, Tatsuya Kawahara

Leveraging Simultaneous Translation for Enhancing Transcription of Low-resource Language via Cross Attention Mechanism

INTERSPEECH 2022

2022

Longfei Yang, Wenqing Wei, Sheng Li, Jiyi Li, Takahiro Shinozaki

Augmented Adversarial Self-Supervised Learning for Early-Stage Alzheimer's Speech Detection

INTERSPEECH 2022

2022

Kai Li, Sheng Li, Xugang Lu, Masato Akagi, Meng Liu, Lin Zhang, Chang Zeng, Longbiao Wang, Jianwu Dang, Masashi Unoki

Data Augmentation Using McAdams-Coefficient-Based Speaker Anonymization for Fake Audio Detection

INTERSPEECH 2022

2022

| |
|---|
| Zhengdong Yang, Wangjin Zhou, Chenhui Chu, Sheng Li, Raj Dabre, Raphael Rubino, Yi Zhao |
| Fusion of Self-supervised Learned Models for MOS Prediction |
| INTERSPEECH 2022 |
| 2022 |

| |
|---|
| Hao Shi, Longbiao Wang, Sheng Li, Jianwu Dang, Tatsuya Kawahara |
| Monaural Speech Enhancement Based on Spectrogram Decomposition for Convolutional Neural Network-sensitive Feature Extraction |
| INTERSPEECH 2022 |
| 2022 |

| |
|---|
| Longfei Yang, Jiyi Li, Sheng Li, Takahiro Shinozaki |
| Multi-Domain Dialogue State Tracking with Top-k Slot Self Attention |
| SIGdial Meeting Discourse and Dialogue 2022 |
| 2022 |

| |
|---|
| Kak Soky, Zhuo Gong, Sheng Li |
| Nict-Tib1: A Public Speech Corpus Of Lhasa Dialect For Benchmarking Tibetan Language Speech Recognition Systems |
| 25th Conference of the Oriental COCOSDA International Committee for the Co-ordination and Standardisation of Speech Databases and Assessment Techniques (O-COCOSDA) |
| 2022 |

| |
|---|
| Hao Shi, Longbiao Wang, Sheng Li, Jianwu Dang, Tatsuya Kawahara |
| Subband-based Spectrogram Fusion for Speech Enhancement by Combining Mapping and Masking Approaches |
| Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC) |
| 2022 |

| |
|---|
| Zhuo Gong, Saito Daisuke, Sheng Li, Hisashi Kawai, Minematsu Nobuaki |
| Can We Train a Language Model Inside an End-to-End ASR Model? - Investigating Effective Implicit Language Modeling |
| Proceedings of the Second Workshop on When Creative AI Meets Conversational AI |
| 2022 |

| |
|---|
| Chao Tan, Yang Cao, Sheng Li and Masatoshi Yoshikawa |
| GENERAL OR SPECIFIC? INVESTIGATING EFFECTIVE PRIVACY PROTECTION IN FEDERATED LEARNING FOR SPEECH EMOTION RECOGNITION |
| ICASSP |
| 2023 |

| |
|---|
| Kak Soky, Sheng Li, Chenhui Chu, Tatsuya Kawahara |
| DOMAIN AND LANGUAGE ADAPTATION USING HETEROGENEOUS DATASETS FOR WAV2VEC2.0-BASED SPEECH RECOGNITION OF LOW-RESOURCE LANGUAGE |
| ICASSP |
| 2023 |

Z. Gong, D. Saito, L. Yang, T. Shinozaki, S. Li, H. Kawai and N. Minematsu

Self-Adaptive Multilingual ASR Rescoring with Language Identification and Unified Language Model

ISCA-Odyssey (The Speaker and Language Recognition Workshop)

2022

S. Li, J. Li, Q. Liu and Z. Gong

Adversarial Speech Generation and Natural Speech Recovery for Speech Content Protection

LREC (Language Resources and Evaluation Conference)

2022

Y. Lv, L. Wang, M. Ge, S. Li, C. Ding, L. Pan, Y. Wang, J. Dang, K. Honda

Compressing Transformer-based ASR Model by Task-driven Loss and Attention-based Multi-level Feature Distillation

in Proc. IEEE-ICASSP, pp. 7992--7996, 2022.

2022

K. Wang, Y. Peng, H. Huang, Y. Hu, and S. Li

Mining Hard Samples Locally and Globally for Improved Speech Separation

in Proc. IEEE-ICASSP, pp. 6037--6041, 2022.

2022

H. Shi, L. Wang, S. Li, C. Fan, J. Dang, and T. Kawahara

Spectrograms Fusion-based End-to-End Robust Automatic Speech Recognition

In Proc. APSIPA ASC, pp. 438--442, 2021.

2021

Y. Peng, J. Zhang, H. Zhang, H. Xu, H. Huang, S. Li, and E.S. Chng

Multilingual Approach to Joint Speech and Accent Recognition with DNN-HMM Framework

In Proc. APSIPA ASC, pp. 1043--1048, 2021.

2021

K. Soky, S. Li, M. Mimura, C. Chu, and T. Kawahara

On the Use of Speaker Information for Automatic Speech Recognition in Speaker-imbalanced Corpora

In Proc. APSIPA ASC, pp. 433--437, 2021.

2021

D. Wang, S. Ye, X. Hu, S. Li, and X. Xu

An End-to-End Dialect Identification System with Transfer Learning from a Multilingual Automatic Speech Recognition Model

in Proc. INTERSPEECH, pp. 3266--3270, 2021.

2021

K. Wang, H. Huang, Y. Hu, Z. Huang, and S. Li

End-to-End Speech Separation Using Orthogonal Representation in Complex and Real Time- Frequency Domain

in Proc. INTERSPEECH, pp. 3046--3050, 2021.

2021

S. Li, R. Dabre, R. Raphael, W. Zhou, Z. Yang, C. Chu, Y. Zhao

The System Description for VoiceMOS Challenge 2022 (KK team, main/ood tasks)

VoiceMOS Challenge 2022

2022

D. Wang, S. Ye, X. Hu, S. Li

The Royal Flush-NICT System Description for AP21-OLR Challenge (Silk-road team, full tasks)

OLR2021 (oriental language recognition challenge)

2021

W. Wei, R. Wang, S. Li, Y. Guo and H. Huang

System description of Alzheimer's disease early detection (Silk-road team, short speech track)

In special session of NCMMSC2021 (Alzheimer's disease detection challenge), 2021

2021

X. Chen, H. Huang, and S. Li

Adversarial Attack and Defense on Deep Neural Network-based Voice Processing Systems: An Overview

National Conference on Man-Machine Speech Communication (NCMMSC), 2021. (report is selected to publish in Applied Sciences, Special Issues of Machine Speech Communication)

2021

L. Qiang, H. Shi, M. Ge, H. Yin, N. Li, L. Wang, S. Li and J. Dang

Speech Dereverberation Based on Scale-aware Mean Square Error Loss

International Conference on Neural Information Processing (ICONIP2021), pp 55-63, Springer, 2021.

2021

H. Yin, L. Qiang, H. Shi, L. Wang, S. Li, M. Ge, G. Zhang and J. Dang

Simultaneous Progressive Filtering-based Monaural Speech Enhancement

International Conference on Neural Information Processing (ICONIP2021), pp 213-221, Springer, 2021.

2021

D. Liu, L. Wang, S. Li, H. Li, C. Ding, J. Zhang and J. Dang

Exploring Effective Speech Representation via ASR for High-Quality End-to-End Multispeaker TTS

International Conference on Neural Information Processing (ICONIP2021), pp 110-118, Springer, 2021.

2021

| Sheng Li | 2022 |
|---|---|
| NICT | 112 |
| Voices of the Himalayas: Investigation of Speech Recognition Technology for the Tibetan Language | |

| Sheng Li | 2022 |
|---|---|
| NICT | 110 |
| Phantom in the Opera: The Vulnerabilities of Speech-based Artificial Intelligence Systems | |

The international collaboration with China is based on NICT's international collaboration activity in 2021. In 2022, we turn to Singapore according to NICT's rule.

https://www.nict.go.jp/outcome/journals/journals_2021_j.html

https://www.nict.go.jp/outcome/proceedings/proceedings_2021_j.html
google scholar of Sheng Li
https://scholar.google.com/citations?user=zHAhsOIAAAAJ&hl=en
Lab homepage of Sheng Li
https://ast-astrec.nict.go.jp/member/sheng-li/index.html

| | | |
|---|---|---|
| | | |

0

| | |
|---|---|
| | |

| | Tianjin University | Xinjiang University | Royal Flush AI Research Inc. | |
| --- | --- | --- | --- | --- |
| | Nanyang Technological University | | | |