

Improving Generalization Ability for Audio Deepfake Detection 提高音频鉴伪的泛化能力

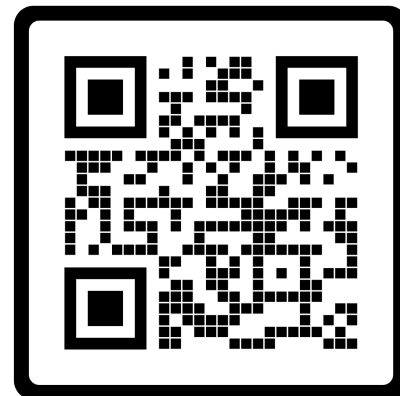
You Zhang 张优, University of Rochester



Dec 28, 2023

Self-introduction

- PhD candidate at University of Rochester, working with Prof. Zhiyao Duan
- B.Eng from UESTC, Exchange studies at UC Berkeley
- Research interests:
 - Computer Audition
 - Audio Deepfake Detection
 - Spatial Audio Personalization
 - Audio-Visual Rendering and Analysis



SCAN ME

Outline

- Introduction to speech anti-spoofing
- Generalization ability to unseen synthetic attacks
 - One-class learning: OC-Softmax;
 - Multi-center one-class learning: SAMO
- Beyond speech anti-spoofing: **Singing voice** deepfake detection
- Future directions

Demo of Speech Deepfakes

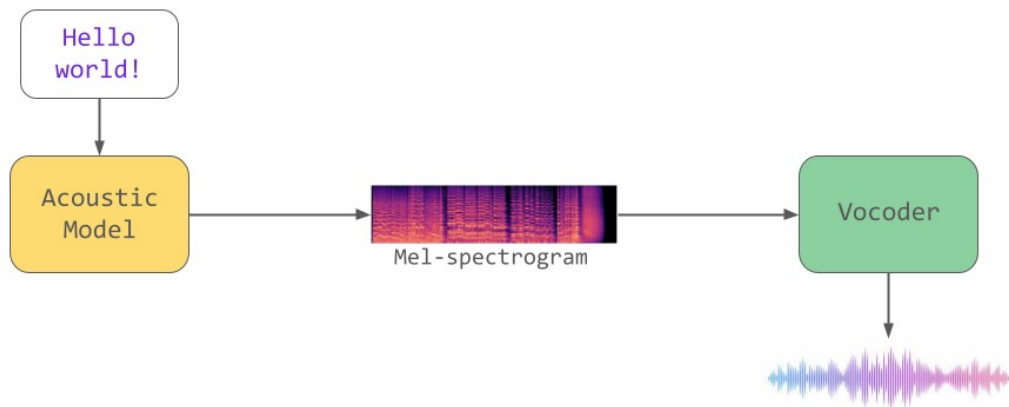
- Tacotron2



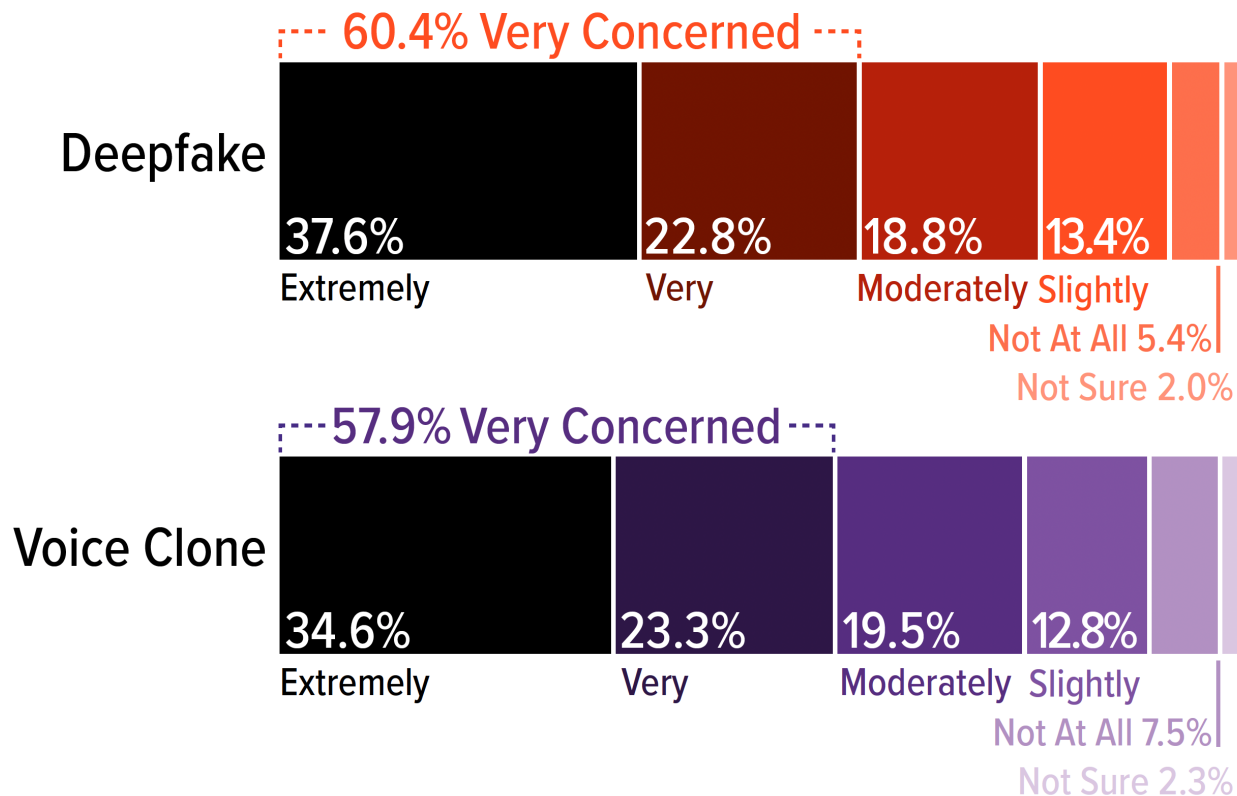
- Fastpitch



- BigVGAN



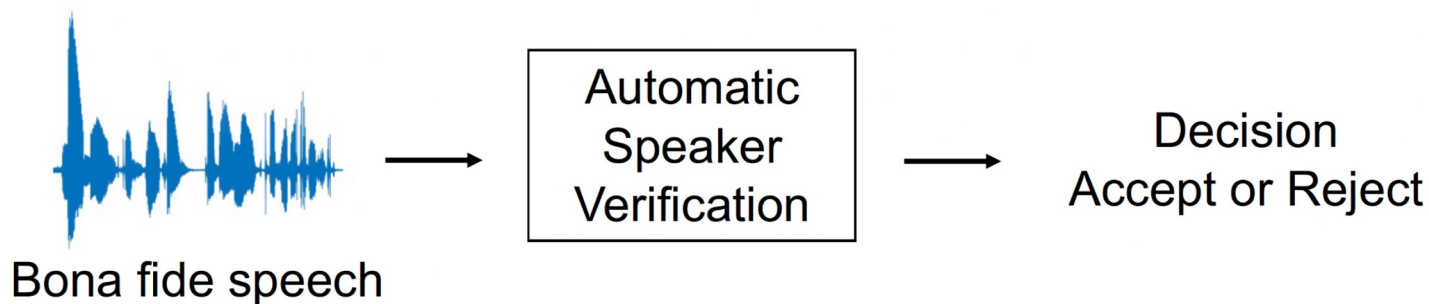
Deepfake and Voice Clone Concern Among Consumers



Source: Voicebot

Voice Biometrics

Verify the identity of a speaker



We expect that the input (bona fide speech) is from a real person.

TECHNEWSWORLD

COMPUTING | INTERNET | IT | MOBILE TECH | REVIEWS | SECURITY

ARTIFICIAL INTELLIGENCE

Microsoft's New AI Can Simulate Anyone's Voice From a 3-Second Sample

By John P. Mello Jr. • January 11, 2023 8:06 AM PT • [Email Article](#)

RESEMBLE.AI

PRODUCTS | USE CASES | PRICING | SIGN IN

Your Complete Generative Voice AI Toolkit

community built voices

- Text-to-Speech
- Speech-to-Speech
- Neural Audio Editing
- Language Dubbing

Resemble's AI voice generator lets you create human-like voice overs in seconds.

Forbes

FORBES > INNOVATION > CYBERSECURITY

EDITORS' PICK

Fraudsters Cloned Company Director's Voice In \$35 Million Bank Heist, Police Find

DIGITAL MUSIC NEWS

CATEGORIES + SYNC NEWS JOBS + PODCASTS

Home > Music Industry News

AI Voice Tool Abused to Make Celebrity Deepfake Audio Clips

Ashley King • February 1, 2023

[top left](#)
[top right](#)
[bottom left](#)
[bottom right](#)

Spoofing attacks

Impersonation

- twins and professional mimics

Replay

- reuse pre-recorded audio, most accessible

Text-to-speech (TTS)

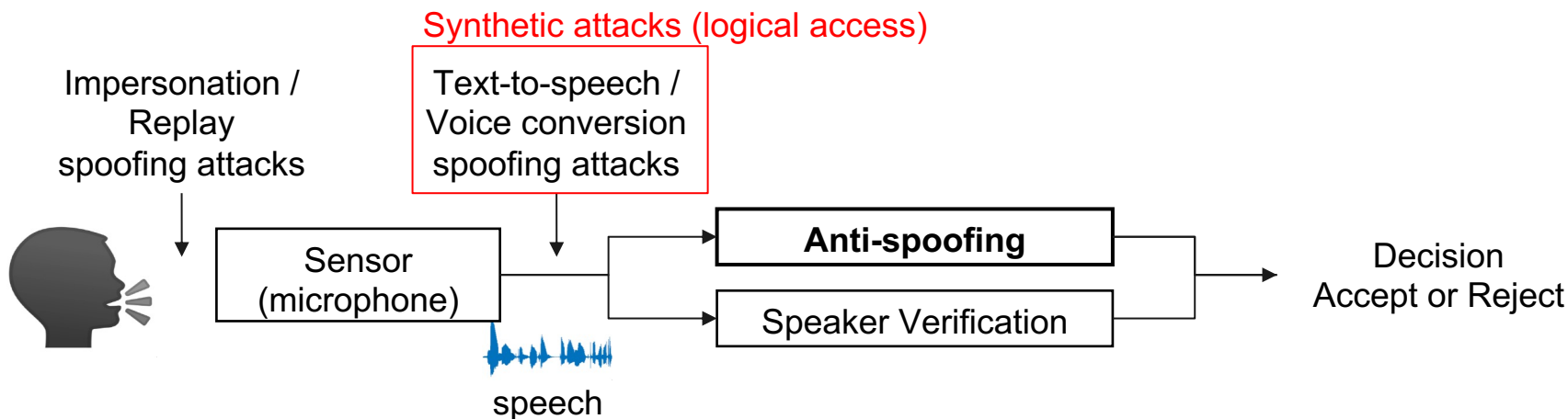
- convert written text into spoken words with speech synthesis

Voice conversion (VC)

- convert speech from source speaker to target speaker's voice

Speech Anti-Spoofing (Speech Deepfake Detection)

A voice anti-spoofing system is desired to distinguish synthetic **speech** from **bona fide speech**.



ASVspoof challenge series

- LA: Robust to channel variability
- PA: Involve real replayed samples
- DF: a new speech deepfake task

2015

Replay spoofing
attacks detection

2019

Text-to-speech
(TTS) and voice
conversion (VC)
spoofing attacks
detection

2017

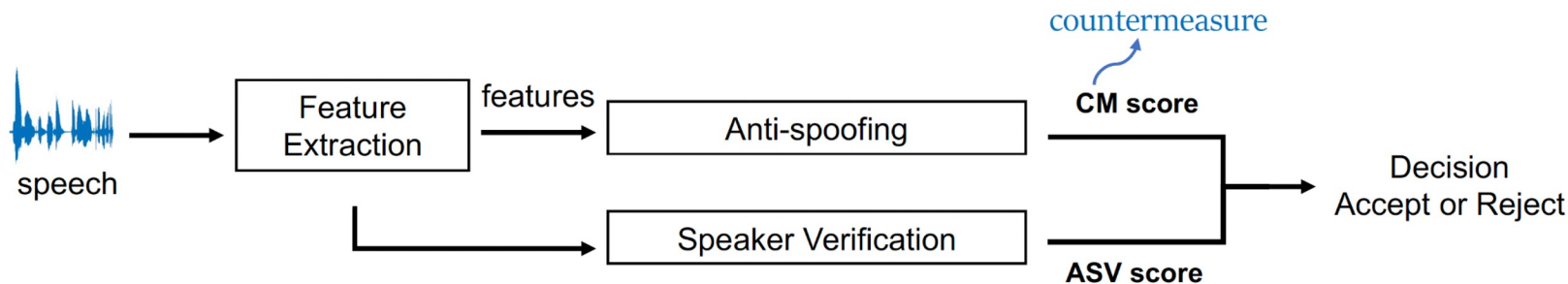
- LA: Advanced TTS and VC attacks
- PA: More controlled setup for replay attacks
- A new evaluation metric

2021

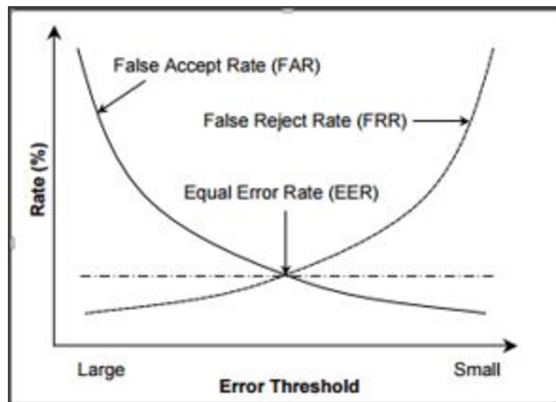
Logical Access (LA): algorithm-related artifacts

Physical Access (PA): device-related artifacts

Evaluation metric



- EER (Equal Error rate)



$$P_{fa}(\theta) = \frac{\#\{\text{spooof trials with score} > \theta\}}{\#\{\text{total spooof trials}\}},$$

$$P_{miss}(\theta) = \frac{\#\{\text{human trials with score} \leq \theta\}}{\#\{\text{total human trials}\}}$$

$$P_{fa}(\theta_{EER}) = P_{miss}(\theta_{EER})$$

Dataset

ASVspoof 2019 Logical Access (TTS + VC)

- Bona fide speech (VCTK dataset)
- 6 known attacks (appear in training set)
- 11 unknown attacks (only appear in eval set)
- 2 attacks (use known algorithms but trained with more data)

	Bona fide	Spoofed	
	# utterance	# utterance	attacks
Training	2,580	22,800	A01 - A06
Development	2,548	22,296	A01 - A06
Evaluation	7,355	63,882	A07 - A19

Research question

Motivation:

- The fast development of speech synthesis are posing increasingly more threat.
- The **distribution mismatch** between the training set and test set for the spoofing attacks class.

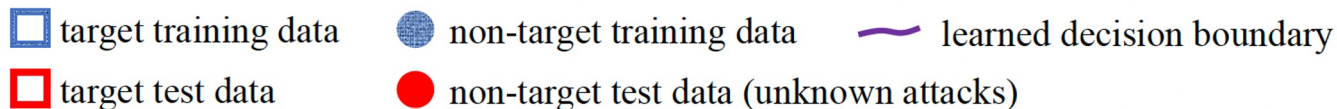
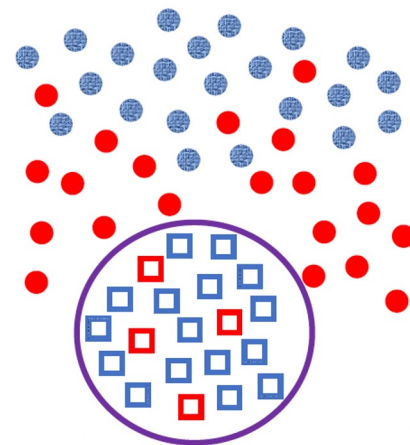
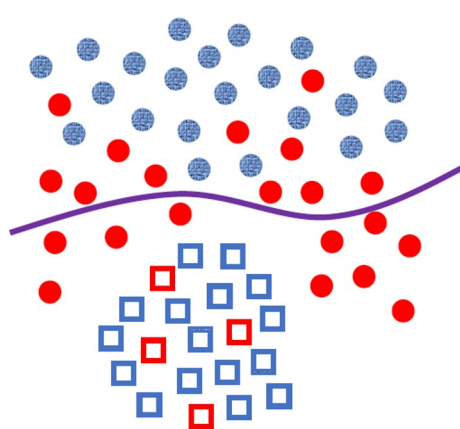
How can the anti-spoofing system defend against **unseen** spoofing attacks?

Generalization ability!

Definition of one-class classification

- “One of the classes (referred to as the positive class or **target** class)
 - is **well characterized** by instances in the training data.
- For the other class (**nontarget**),
 - it has either **no instances** at all,
 - **very few** of them,
 - or they do **not form a statistically-representative** sample of the negative concept.”

Illustration of comparison



(a) Binary classification

(b) One-class classification

One-class learning

- Compact the bona fide speech representation
- Isolate the spoofing attacks

Training: OC-Softmax loss (Proposed)

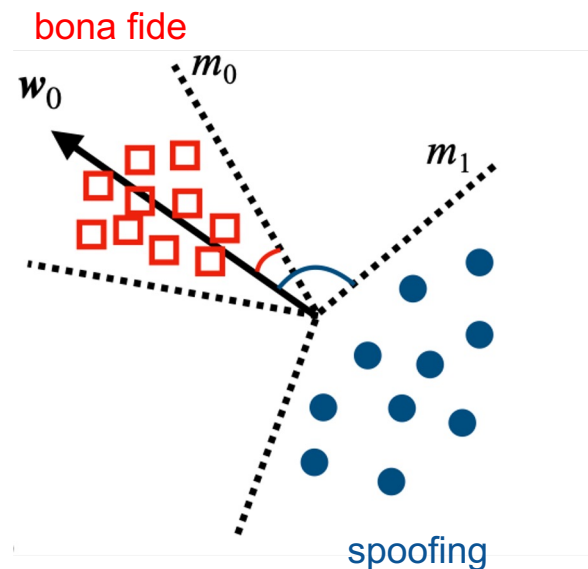
$$\mathcal{L}_{OCS} = \frac{1}{N} \sum_{i=1}^N \log \left(1 + e^{\alpha(m_{y_i} - \hat{w}_0 \hat{x}_i)(-1)^{y_i}} \right).$$

Diagram illustrating the OC-Softmax loss formula with annotations:

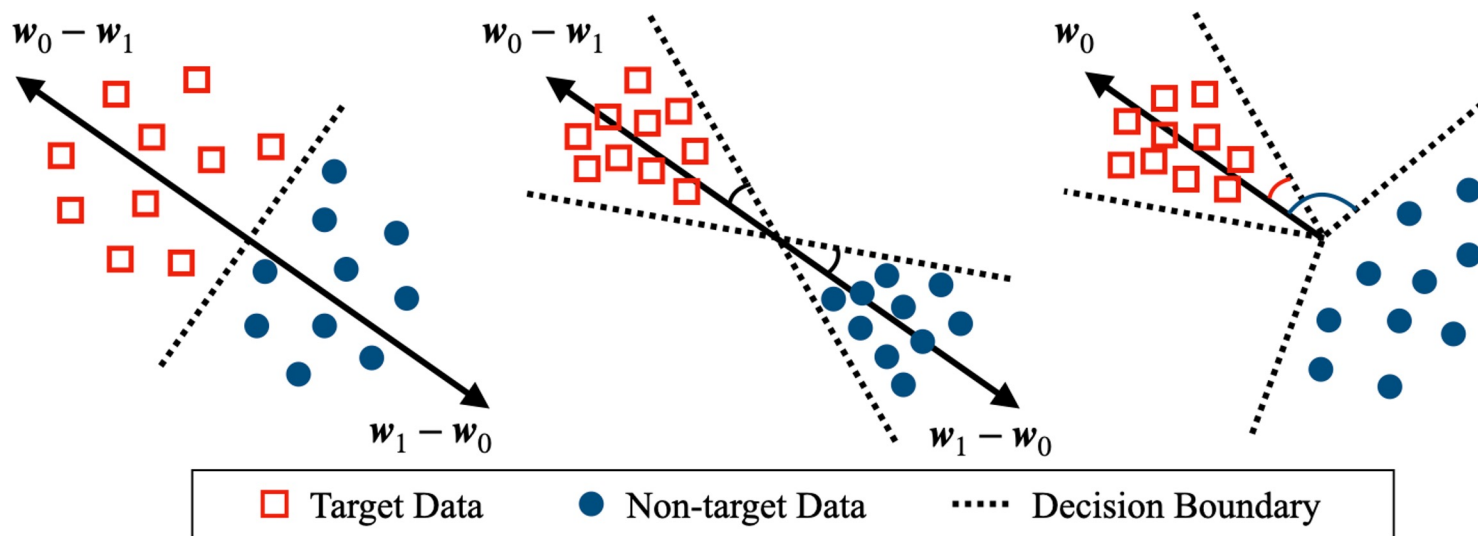
- $\frac{1}{N}$: # samples
- α : scale factor
- m_{y_i} : margin
- \hat{w}_0 : center vector
- \hat{x}_i : embedding
- $(-1)^{y_i}$: label

Inference: cosine similarity

$$S_{OCS} = \hat{w}_0 \hat{x}_i.$$



Comparing OC-Softmax with binary classification



(a) Original Softmax

(b) AM-Softmax

(c) OC-Softmax (Proposed)

Comparing OC-Softmax with binary classification

Softmax:

$$\begin{aligned}\mathcal{L}_S &= -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i}}{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i} + e^{\mathbf{w}_{1-y_i}^T \mathbf{x}_i}} \\ &= \frac{1}{N} \sum_{i=1}^N \log (1 + e^{(\mathbf{w}_{1-y_i} - \mathbf{w}_{y_i})^T \mathbf{x}_i}),\end{aligned}$$

AM-Softmax:

$$\begin{aligned}\mathcal{L}_{AMS} &= -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\alpha(\hat{\mathbf{w}}_{y_i}^T \hat{\mathbf{x}}_i - m)}}{e^{\alpha(\hat{\mathbf{w}}_{y_i}^T \hat{\mathbf{x}}_i - m)} + e^{\alpha \hat{\mathbf{w}}_{1-y_i}^T \hat{\mathbf{x}}_i}} \\ &= \frac{1}{N} \sum_{i=1}^N \log \left(1 + e^{\alpha(m - (\hat{\mathbf{w}}_{y_i} - \hat{\mathbf{w}}_{1-y_i})^T \hat{\mathbf{x}}_i)} \right),\end{aligned}$$

OC-Softmax:

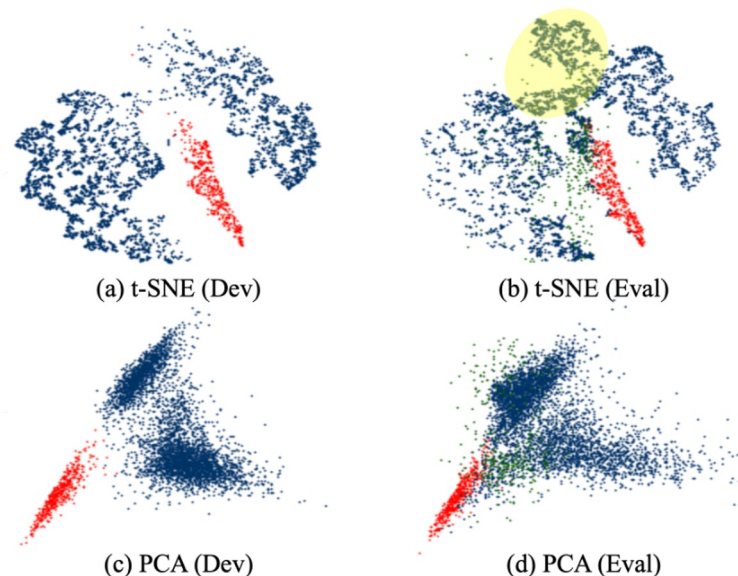
$$\mathcal{L}_{OCS} = \frac{1}{N} \sum_{i=1}^N \log (1 + e^{\alpha(m_{y_i} - \hat{\mathbf{w}}_0 \hat{\mathbf{x}}_i)(-1)^{y_i}}).$$

Evaluation of OC-Softmax

Results on the development and evaluation sets of ASVspoof 2019 LA using different losses

Loss	Dev Set		Eval Set	
	EER (%)	min t-DCF	EER (%)	min t-DCF
Softmax	0.35	0.010	4.69	0.125
AM-Softmax	0.43	0.013	3.26	0.082
OC-Softmax	0.20	0.006	2.19	0.059

- OC-Softmax performs the best on unseen attacks.
- Achieved the state-of-the-art single-system performance.



Feature Embedding Visualization
(red: bona fide, green: A17 attack, blue: spoofing attacks)

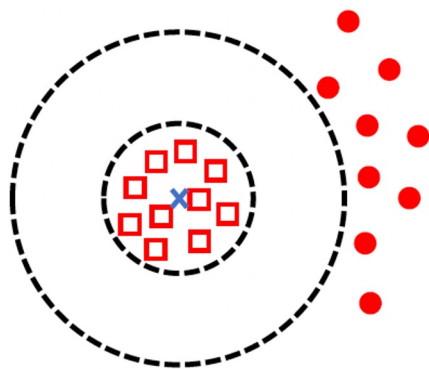
Comparison with single systems

System	EER (%)	min t-DCF
CQCC + GMM [3]	9.57	0.237
LFCC + GMM [3]	8.09	0.212
Chettri et al. [22]	7.66	0.179
Monterio et al. [14]	6.38	0.142
Gomez-Alanis et al. [16]	6.28	-
Aravind et al. [18]	5.32	0.151
Lavrentyeva et al. [21]	4.53	0.103
ResNet + OC-SVM	4.44	0.115
Wu et al. [17]	4.07	0.102
Tak et al. [19]	3.50	0.090
Chen et al. [15]	3.49	0.092
Proposed	2.19	0.059

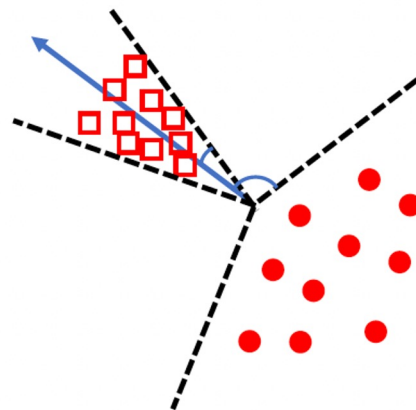
Other one-class loss functions

Euclidean distance-based one-class loss (isolate loss, single-center loss)

Cosine similarity-based one-class loss (OC-Softmax, angular isolate loss)



(a) Euclidean distance-based



(b) Cosine similarity-based

Other one-class loss functions

Isolate loss:

$$L_{ISO} = \frac{1}{|\Omega|} \sum_{\Omega} \max(0, \|\mathbf{x}_i - \mathbf{c}\| - r_0) + \frac{1}{|\Omega|} \sum_{\Omega} \max(0, r_1 - \|\mathbf{x}_i - \mathbf{c}\|),$$

Single-center loss:

$$L_{SCL} = \frac{1}{|\Omega|} \sum_{\Omega} \|\mathbf{x}_i - \mathbf{c}\| + \max\left(0, \frac{1}{|\Omega|} \sum_{\Omega} \|\mathbf{x}_i - \mathbf{c}\| - \frac{1}{|\Omega|} \sum_{\Omega} \|\mathbf{x}_i - \mathbf{c}\| + \beta \sqrt{D}\right)$$

Angular isolate loss:

$$L_{AISO} = \frac{1}{|\Omega|} \sum_{\Omega} \log\left(1 + e^{\alpha(m_0 - \hat{\mathbf{w}}^T \hat{\mathbf{x}}_i)}\right) + \frac{1}{|\Omega|} \sum_{\Omega} \log\left(1 + e^{\alpha(\hat{\mathbf{w}}^T \hat{\mathbf{x}}_i - m_1)}\right)$$

Follow-up works

ONE-CLASS KNOWLEDGE DISTILLATION FOR SPOOFING

Jingze Lu^{1,2}, Yuxiang Zhang^{1,2}, Wenchao Wang¹, Zengqiang Sha

¹Key Laboratory of Speech Acoustics and Content Understanding
Chinese Academy of Sciences, Beijing, China

²University of Chinese Academy of Sciences, Beijing, China

{ Interspeech 2022
18-22 September 2022, Incheon, Korea

A Deep One-Class Learning Method for Replay Attack Detection

Yijie Lou, Shiliang Pu, Jianfeng Zhou, Xin Qi, Qinbo Dong, Hongwei Zhou

Hikvision Research Institute, Hangzhou

louyijie@hikvision.com, pushiliang.hri@hikvision.com

2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)

2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)

Improved One-class Learning for Voice Spoofing Detection

Lixiang Li*^{† §}, Xiaopeng Xue^{§ †}, Haipeng Peng^{§ †}, Yeqing Ren^{§ †} and Mengmeng Zhao^{§ † ‡}



[†]State Key Laboratory of Networking and Switching Technology,

[‡]Institute of Information and Communication Engineering,

[§]Engineering Laboratory for Disaster Backup and Recovery,

[¶]Department of Posts and Telecommunications, Beijing, China

^{||}Department of Science and Engineering, Zaozhuang University, Zaozhuang, China

*Corresponding author. Email: lixiang@bupt.edu.cn

IEEE/ACM TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING, VOL. 31, 2023

Generalized Voice Spoofing Detection via Integral Knowledge Amalgamation

Yeqing Ren^{||}, Haipeng Peng^{||}, Lixiang Li^{||}, Xiaopeng Xue, Yang Lan, and Yixian Yang^{||}

Research question

Motivation:

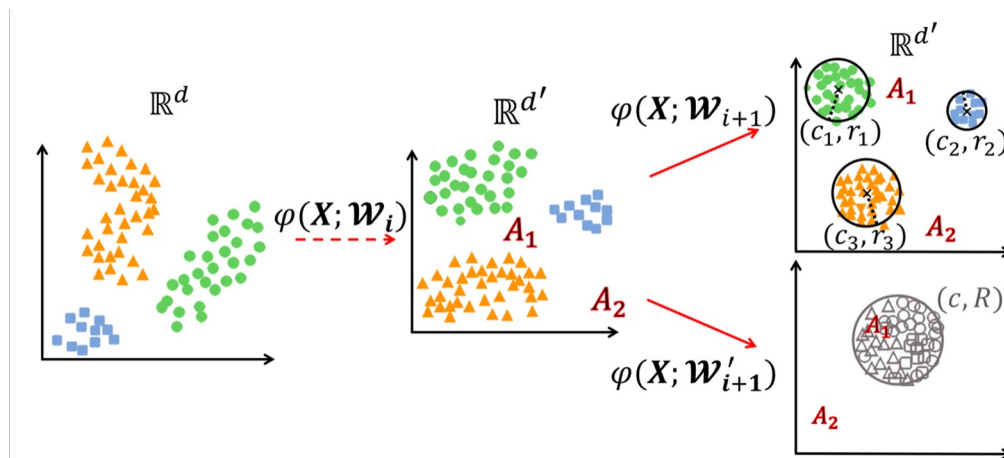
- In our previous work, we compact the embedding space of the bona fide speech into one cluster.
- However, due to **the variety of timbre and speaking traits of different speakers**, the bona fide speech of different speakers naturally forms multiple clusters in the embedding space.

How to improve the **generalization** ability while **maintaining the variation** of bona fide speech?

Inspiring work

Deep Multi-sphere Support Vector (SDM'20)

If data is naturally multi-cluster, merging them into one cluster could be harmful for detecting anomaly.



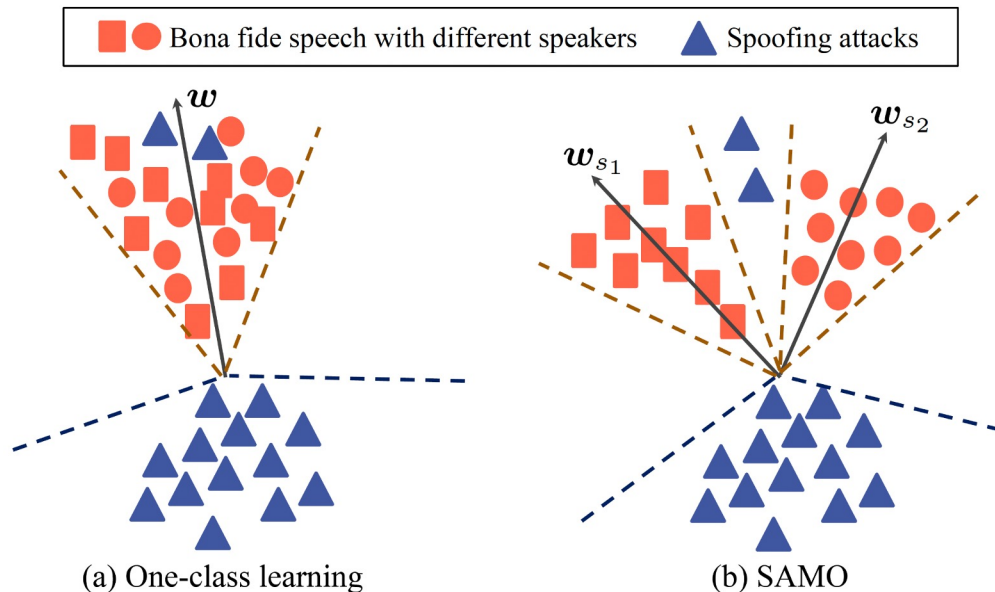
(Figure 2 in Ghafoori et al.)

Ghafoori, Z., & Leckie, C. (2020). Deep multi-sphere support vector data description. In *Proceedings of the 2020 SIAM International Conference on Data Mining* (pp. 109-117). Society for Industrial and Applied Mathematics.

Speaker attractor multi-center one-class learning

Model speaker diversity while maintaining the generalization ability brought by one-class learning

- Discriminate bona fide vs. spoofing attacks
- Cluster bona fide speech according to speakers



Siwen Ding, You Zhang, and Zhiyao Duan. "SAMO: Speaker Attractor Multi-Center One-Class Learning for Voice Anti-Spoofing", *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023.

Speaker attractors

Define: a speaker-specific **anchor** in the embedding space

Compute: average the embeddings of each speaker's bona fide speech

- **Training:** **attract** bona fide speech embeddings of the same speaker
- **Inference:** **cosine similarity** between test utterance and enrolled utterance or attractors of training speakers

$$S_{SAMO} = \begin{cases} \hat{\mathbf{w}}_{s_i} \hat{\mathbf{x}}_i & \text{if } s_i \text{ is enrolled} \\ \max_s (\hat{\mathbf{w}}_s \hat{\mathbf{x}}_i), s \in \mathcal{S}_{train} & \text{otherwise} \end{cases},$$

The diagram includes the following callouts:

- speaker attractor**: points to $\hat{\mathbf{w}}_{s_i}$
- embedding**: points to $\hat{\mathbf{x}}_i$
- speaker label of i -th utterance**: points to s_i
- speakers in the training set**: points to \mathcal{S}_{train}

Loss function for multi-center one-class learning

- Compact the bona fide speech representation belonging to the same speaker
- Push away the spoofing attacks from all speaker attractors

$$\mathcal{L}_{SAMO} = \frac{1}{N} \sum_{i=1}^N \log \left(1 + e^{\alpha(m_{y_i} - d_i)(-1)^{y_i}} \right),$$

Annotations for the equation:

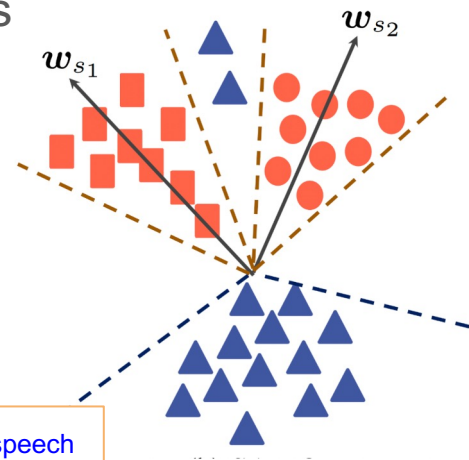
- # samples (points to N)
- scale factor (points to α)
- margin (points to m_{y_i})
- label (points to y_i)

where d_i is calculated by

$$d_i = \begin{cases} \hat{w}_{s_i} \hat{x}_i & \text{if } y_i = 0 \\ \max_s (\hat{w}_s \hat{x}_i), s \in \mathcal{S}_{train} & \text{if } y_i = 1 \end{cases}$$

Annotations for the equation:

- bona fide speech (points to the $y_i = 0$ case)
- spoofing attacks (points to the $y_i = 1$ case)



SAMO training algorithm

- Compact the bona fide utterances spoken by the same speaker
- Push away spoofing utterances from all speaker attractors

Algorithm 1: SAMO Training Algorithm

Require: T : Total number of epochs

M : speaker attractor update interval (# epochs)

```
1 Initialize network  $F$  with random weights
2 Initialize speaker attractors  $w_s$  as one-hot vectors
3 for  $i \leftarrow 1$  to  $T$  do
4   | if  $i \bmod M = 0$  then
5   |   | Update  $w_s$  as the average bona fide embedding for
6   |   |   each speaker  $s \in \mathcal{S}_{train}$ 
7   |   | end if
8   |   | Update  $F$  by  $\mathcal{L}_{SAMO}$  with mini-batches ▷ Eq. (3)
9   | end for
10 return Optimized network  $F$  and speaker attractors  $w_s$ 
```

Comparison with state-of-the-art methods

Table 2. Comparison of our proposed SAMO with Softmax and OC-Softmax on the target-only portion of the ASVspoof2019 LA evaluation set. All the systems use AASIST [9] backbone. The average (best) results across 3 trials with random training seeds are shown.

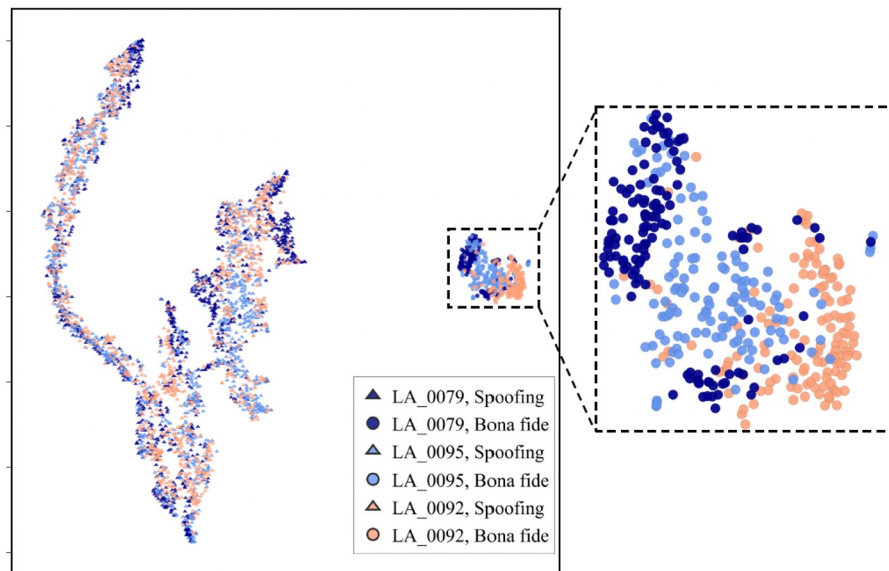
Method	EER(%)	min t-DCF
Softmax	1.74 (1.25)	0.0583 (0.0425)
OC-Softmax	1.25 (1.17)	0.0415 (0.0393)
SAMO (test w/o enrollment)	1.09 (0.91)	0.0363 (0.0306)
SAMO (test w/ enrollment)	1.08 (0.88)	0.0356 (0.0291)

SAMO further improves the performance, indicating the advantage brought by the multi-center modeling of bona fide speech.

Embedding visualization

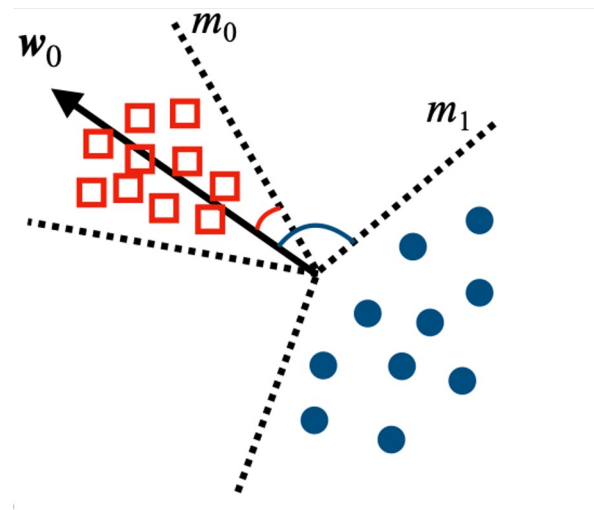
2D t-SNE visualization of SAMO feature embeddings of bona fide and spoofed speech of three speakers

- Bona fide utterances are grouped in a small region.
- Utterances of the three speakers are generally clustered according to speaker identity.



Takeaways

- One-class learning aims to **compact the target** class representation in the embedding space, and **push away non-target**.
- The proposed OC-Softmax and SAMO could improve the **generalization ability** of anti-spoofing system against **unseen spoofing attacks**.



Demo of singing voice deepfakes

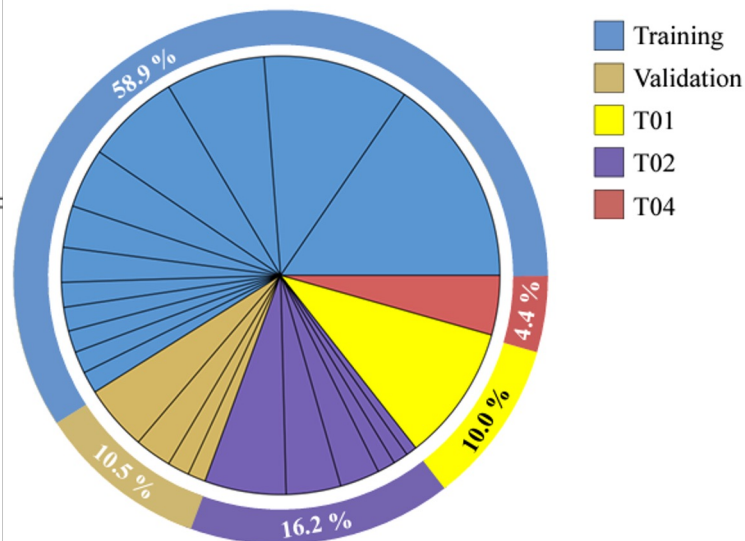


<https://www.bilibili.com/video/BV1Sb4y1V76A>

SingFake dataset: singing voice deepfake detection (SVDD)

Table 1. SingFake statistics for each split.

Splits	Description	# Singers	Languages (Sorted by percentages in the splits)	# Clips (Real / Fake)
Train	Training set	12	Mandarin, Cantonese, Japanese, English, Others	5251 / 4519
Val	Validation set (unseen singers)	4		1089 / 543
T01	Test set for seen singer Stefanie Sun	1		370 / 1208
T02	Test set for unseen singers	6		1685 / 1006
T03	T02 over 4 communication codecs	6		6740 / 4024
T04	Test set for Persian musical context	17		353 / 166

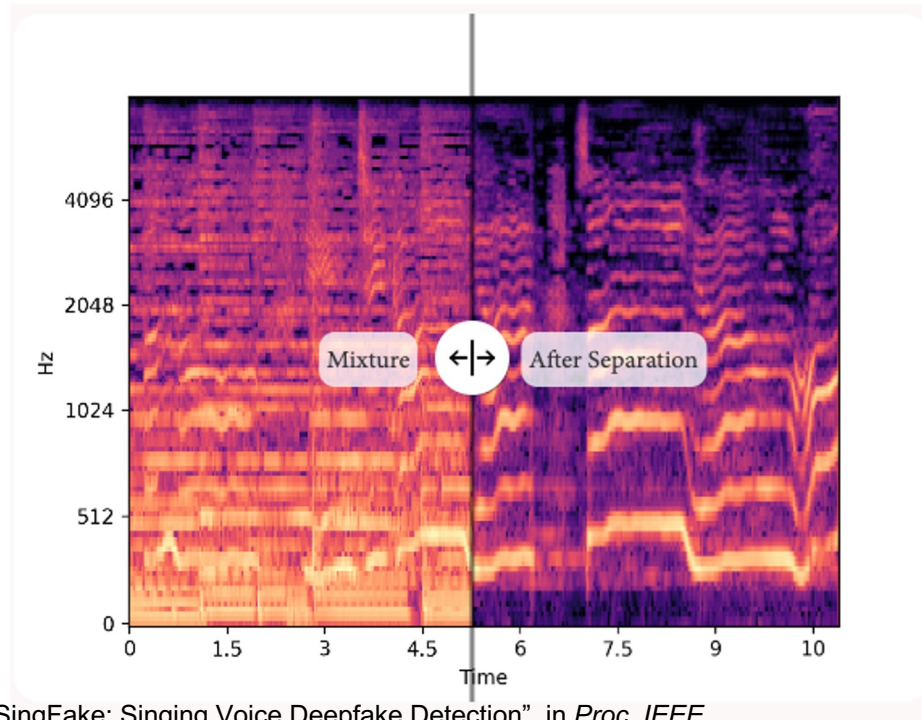


Yongyi Zang*, You Zhang*, Mojtaba Heydari, and Zhiyao Duan. "SingFake: Singing Voice Deepfake Detection", in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024. (* equal contribution)

Speech anti-spoofing heavily degrades on SVDD task

Table 2. Test results on speech and singing voice with CM systems trained on speech utterance from ASVspoof2019LA (EER (%)).

Method	ASVspoof2019	SingFake-T02	
	LA - Eval	Mixture	Vocals
AASIST	0.83	58.12	37.91
Spectrogram+ResNet	4.57	51.87	37.65
LFCC+ResNet	2.41	45.12	54.88
Wav2Vec2+AASIST	7.03	56.75	57.26



Yongyi Zang*, You Zhang*, Mojtaba Heydari, and Zhiyao Duan. "SingFake: Singing Voice Deepfake Detection", in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024. (* equal contribution)

Performance of training on the SingFake data

Table 3. Evaluation results for SVDD systems on all testing conditions in our SingFake dataset (EER (%))

Method	Setting	Train	T01	T02	T03	T04
AASIST	Mixture	4.10	7.29	11.54	17.29	38.54
	Vocals	3.39	8.37	10.65	13.07	43.94
Spectrogram+ResNet	Mixture	4.97	14.88	22.59	24.15	48.76
	Vocals	5.31	11.86	19.69	21.54	43.94
LFCC+ResNet	Mixture	10.55	21.35	32.40	31.85	50.07
	Vocals	2.90	15.88	22.56	23.62	39.27
Wav2Vec2+AASIST (Joint-finetune)	Mixture	1.57	4.62	8.23	13.62	42.77
	Vocals	1.70	5.39	9.10	10.03	42.19

Training on singing voices improves SVDD performance

SVDD systems show limited robustness to unseen scenarios

Yongyi Zang*, You Zhang*, Mojtaba Heydari, and Zhiyao Duan. "SingFake: Singing Voice Deepfake Detection", in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024. (* equal contribution)

Singing voice deepfake detection (SVDD) challenge

- Controlled setting
 - Clean vocals generated by state-of-the-art singing voice synthesis (SVS) and singing voice conversion (SVC) systems based on open source pop song datasets.
- In-the-wild setting
 - Extended SingFake dataset.



Future directions

Generalizing to diversified spoofing attacks

- Replay + TTS + VC + Adversarial + PartialSpoof

Robustness

- Additive noise, channel variation, quality of TTS/VC systems (In-the-wild)

Explainability

- The artifacts or the cues that distinguish bona fide from spoofed speech

Visually-informed speech anti-spoofing

- Audio-visual deepfake detection

Takeaways

Introduction to speech anti-spoofing

Generalization ability to unseen synthetic attacks

- One-class learning: OC-Softmax,
- Multi-center one-class learning: SAMO

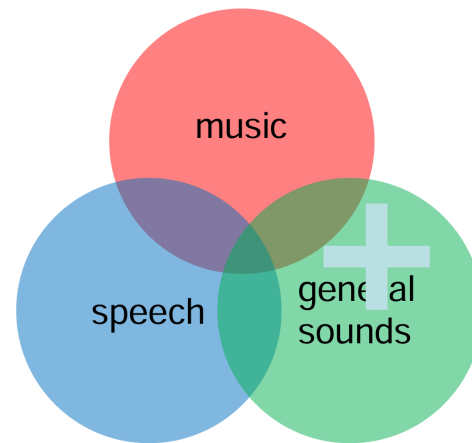
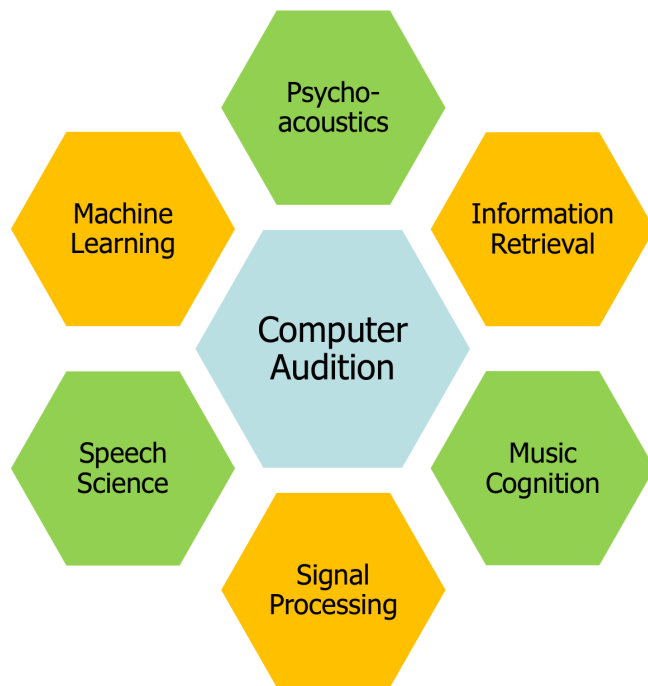
Beyond speech anti-spoofing: **Singing voice** deepfake detection

Thank you! Questions?

References

- [1] **You Zhang**, Fei Jiang, and Zhiyao Duan, "One-class Learning Towards Synthetic Voice Spoofing Detection", *IEEE Signal Processing Letters*, vol. 28, pp. 937-941, 2021. [\[link\]](#)[\[code\]](#)[\[video\]](#)
- [2] **You Zhang**, Ge Zhu, Fei Jiang, and Zhiyao Duan, "An Empirical Study on Channel Effects for Synthetic Voice Spoofing Countermeasure Systems", in *Proc. Interspeech*, pp. 4309-4313, 2021. [\[link\]](#)[\[code\]](#)[\[video\]](#)
- [3] Xinhui Chen*, **You Zhang***, Ge Zhu*, and Zhiyao Duan, "UR Channel-Robust Synthetic Speech Detection System for ASVspoof 2021", in *Proc. ASVspoof 2021 Workshop*, pp. 75-82, 2021. (* equal contribution) [\[link\]](#)[\[code\]](#)[\[video\]](#)
- [4] **You Zhang**, Fei Jiang, Ge Zhu, Xinhui Chen, and Zhiyao Duan, "Generalizing Voice Presentation Attack Detection to Unseen Synthetic Attacks and Channel Variation", *Handbook of Biometric Anti-spoofing (3rd Ed.)*, Springer, 2023. [\[link\]](#)[\[code\]](#)
- [5] **You Zhang**, Ge Zhu, and Zhiyao Duan, "A Probabilistic Fusion Framework for Spoofing Aware Speaker Verification", in *Proc. Odyssey*, 2022. [\[link\]](#)[\[code\]](#)
- [6] Siwen Ding, **You Zhang**, and Zhiyao Duan. "SAMO: Speaker Attractor Multi-Center One-Class Learning for Voice Anti-Spoofing", in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023. [\[link\]](#)[\[code\]](#)
- [7] Yongyi Zang, **You Zhang**, and Zhiyao Duan. "Phase Perturbation Improves Channel Robustness for Speech Spoofing Countermeasures", in *Proc. Interspeech*, pp. 3162-3166, 2023. [\[link\]](#)[\[code\]](#)
- [8] Yongyi Zang*, **You Zhang***, Mojtaba Heydari, and Zhiyao Duan. "SingFake: Singing Voice Deepfake Detection", in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024. (* equal contribution) [\[link\]](#)[\[code\]](#)[\[webpage\]](#)

Computer Audition Research Areas



vision

text

EEG

other modalities

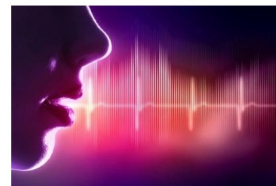
Audio Information Research (AIR) Lab

Machine Understanding of Sounds



MUSIC INFORMATION RETRIEVAL

- Music transcription, alignment
- Source separation
- Generation
- Interactive performance



SPEECH PROCESSING

- Separation and enhancement
- Verification and anti-spoofing
- Emotion analysis
- Diarization
- Text-to-speech
- Voice conversion



ENVIRONMENTAL SOUND UNDERSTANDING

- Sound search by vocal imitation
- Sound event detection
- Source localization
- HRTF modeling
- Smart acoustics



AUDIO-VISUAL PROCESSING

- Talking face generation
- Music performance analysis and generation
- Audio-visual source separation