



# Learning to Fool the Speaker Recognition

ICASSP 2020

<sup>1,4</sup>Jiguo Li,  
[jiguo.li@vip1.ict.ac.cn](mailto:jiguo.li@vip1.ict.ac.cn)

Joint work with <sup>4</sup>Xinfeng Zhang, <sup>3</sup>Jizheng Xu, <sup>3</sup>Li Zhang, <sup>3</sup>Yue Wang, <sup>2</sup>Siwei Ma, <sup>2</sup>Wen Gao

<sup>1</sup>Institute of Computing Technology(ICT), Chinese Academy of Sciences(CAS)

<sup>2</sup>Institute of Digital Media(IDM), Peking University(PKU)

<sup>3</sup>Bytedance.Inc

<sup>4</sup>University of Chinese Academy of Science(UCAS)

# Biometric Systems

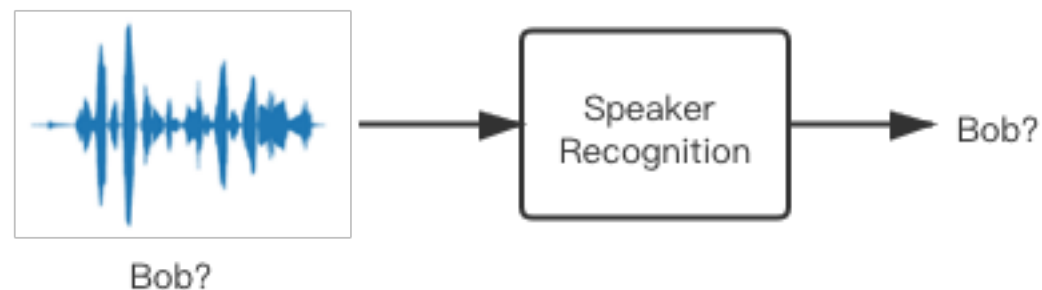
## ◆ Image-based Systems

- ✓ Face, fingerprint



## ◆ Speech-based Systems

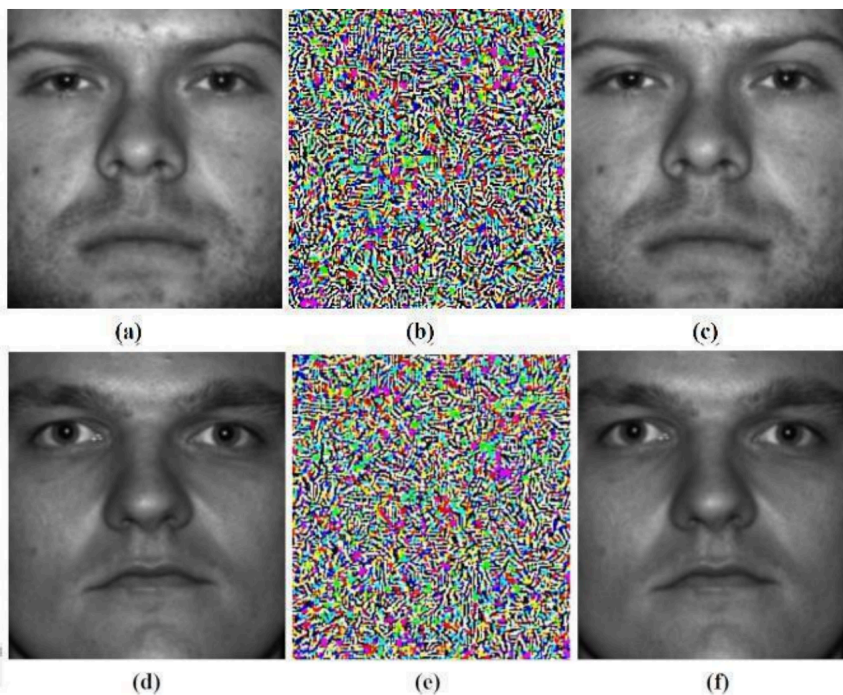
- ✓ Speaker recognition



# Security Risks for Biometric Systems

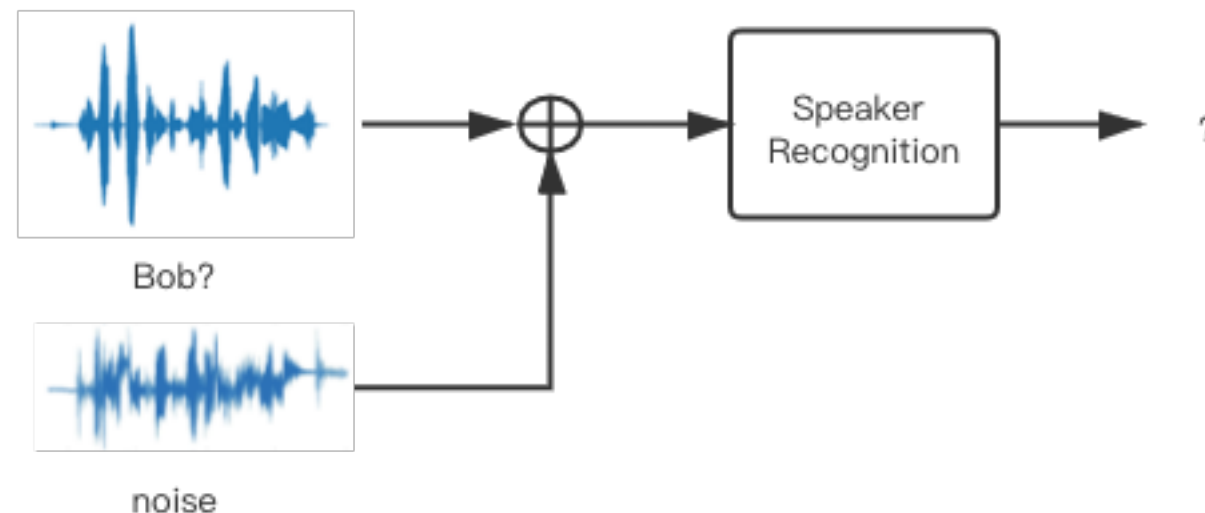
## ◆ Image-based Systems

✓ Face attack



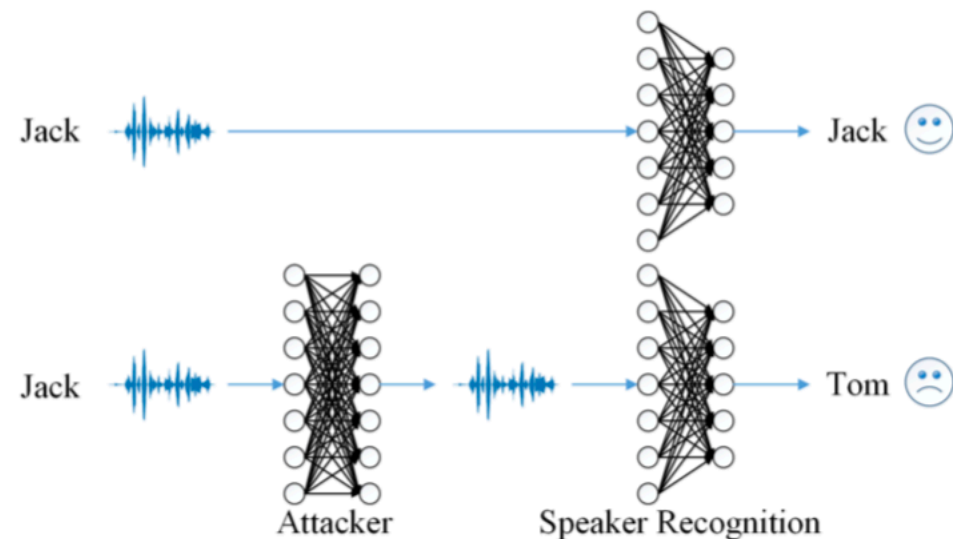
## ◆ Speech-based Systems

✓ Speech attack



# Motivation

- ◆ How to attack the speech-based biometric systems?
- ◆ Is the speech-based biometric systems vulnerable to the adversarial attack?
- ◆ Is it possible to design a biometric systems robust to the adversarial attack?



# Our attack results

## Non-targeted attack

dr1/fcjh0/si1027.ogg  

dr1/fcjh0/sx37.ogg  

dr2/faem0/si762.ogg  

dr8/fbcg1/sx82.ogg  

real

fake

## Targeted attack

dr1/fcjh0/si1027.ogg    

dr1/fdaw0/si1046.ogg    

dr2/faem0/si762.ogg    

dr8/fbcg1/si982.ogg    

real

target0

target100

target200

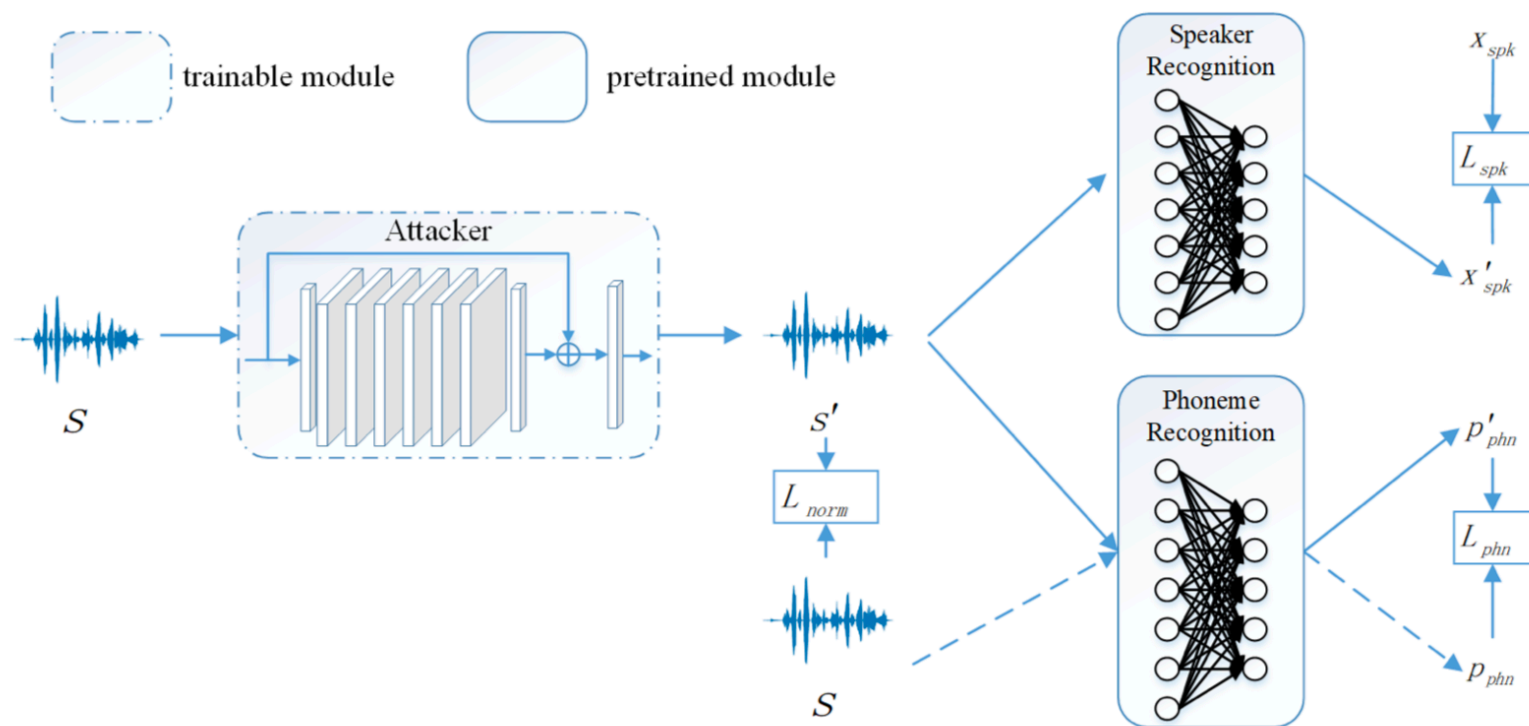
More results: <https://smallflyingpig.github.io/speaker-recognition-attacker/main>



# Proposed Attack Framework

## ◆ Our Framework

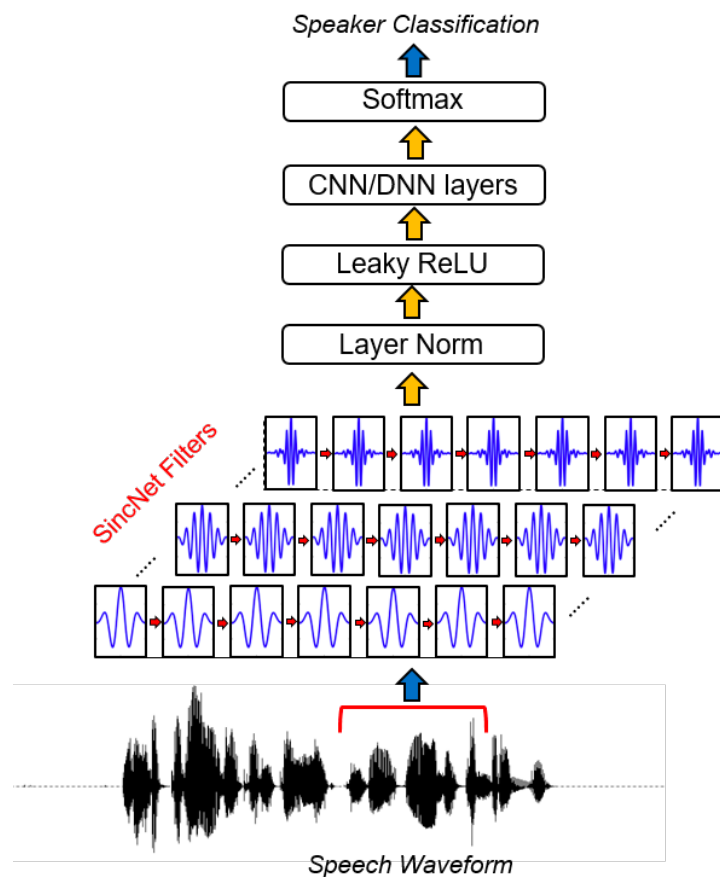
- ✓ An attacker for all samples
- ✓ Optimize the speech via phoneme recognition module



# Proposed Attack Framework

## ◆ Speaker/Phoneme Recognition Model: Sincnet[1]

- ✓ Frequency filters in the first layer
- ✓ Process on the raw waveform
- ✓ More interpretable

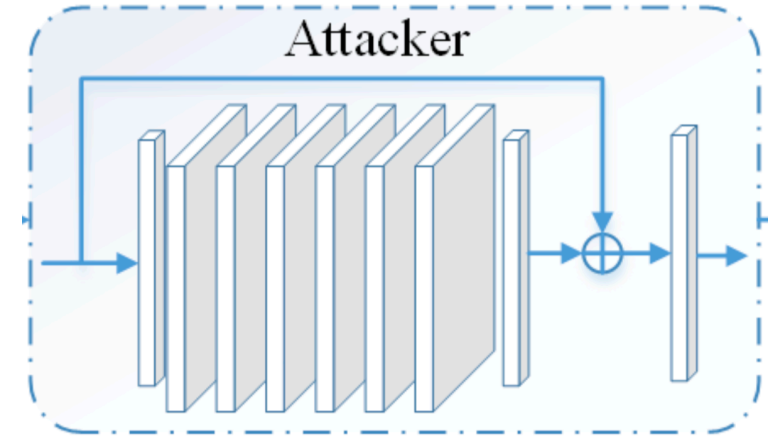




# Proposed Attack Framework

## ◆ Attacker: a Residual Block

- ✓ Referring to Adversarial Transformer Networks (ATNs)[1]
- ✓ Additive perturbations
- ✓ The scale of the perturbation is controllable
- ✓ Training once for all testing samples



[1] Baluja, Shumeet, and Ian Fischer. "Learning to Attack: Adversarial Transformation Networks." *AAAI*. Vol. 1. 2018.



# Proposed Attack Framework

## ◆ Adversarial training/Optimization

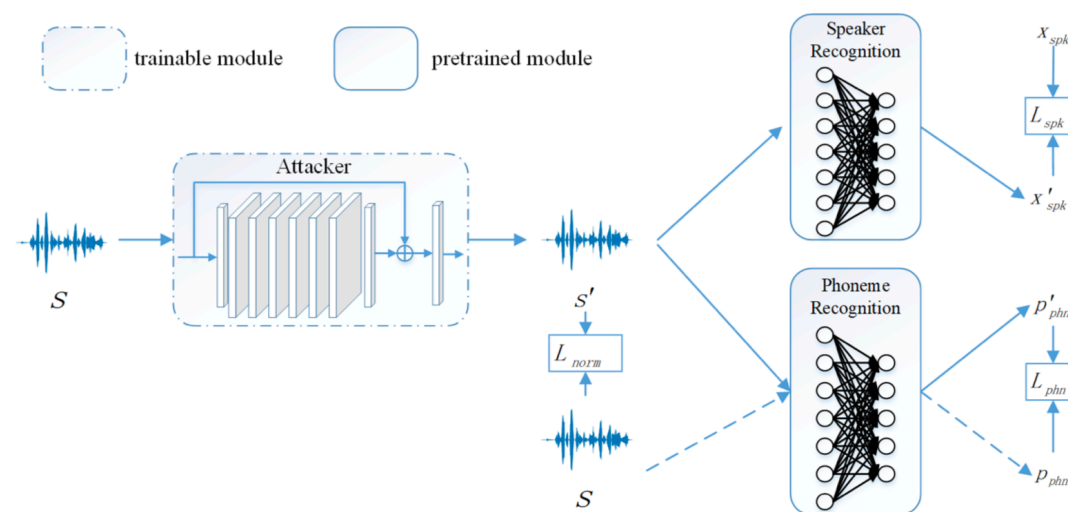
### ✓ Non-targeted attack

$$L_{total} = L_{spk} + \lambda_{phn}L_{phn} + \lambda_{norm}L_{norm}$$

$$L_{spk} = \begin{cases} x'_{spk}[I_{1st}] - x'_{spk}[I_{2nd}], & I_{1st} = y_{spk} \\ 0, & else \end{cases}$$

$$L_{phn} = KL(p_{phn} || p'_{phn})$$

$$L_{norm} = [\max(s - s' - m, 0)]^2$$



### ✓ Targeted attack

$$L_{spk} = \begin{cases} x'_{spk}[I_{1st}] - x'_{spk}[y_{target}], & I_{1st} \neq y_{target} \\ 0, & else \end{cases}$$



# Experimental Results

## ◆ Datasets and Metrics

### ✓ Dataset

Dataset	Label	Speaker number	Samples (train+test)
TIMIT	Speaker+phoneme	462	3694(2309+1385)

### ✓ Metric

- Sentence Error Rate(SER): used for non-targeted attack
- Prediction Target Rate(PTR): used for targeted attack
- Signal-noise Ratio(SNR)
- Perceptual Evaluation of Speech Quality(PESQ): 0.5~4.5



# Experimental Results

◆ Can our proposed model attack the pretrained speaker recognition model?

- ✓ Non-targeted attack
- ✓ SER 90.5% with SNR 59.01 dB
- ✓ SER 90.5% with PESQ 4.28

$\lambda_{phn}$	$\lambda_{norm}$	SER(%)↑	SNR(dB)↑	PESQ↑
-	-	1.52*	-	-
0	0	99.7	18.56	1.09
0	1000	96.5	56.39	3.72
0	2000	86.7	57.79	3.61
1	1000	<b>99.2</b>	57.20	4.20
5	1000	93.9	58.00	4.25
10	1000	90.5	<b>59.01</b>	<b>4.28</b>



# Experimental Results

◆ Can our proposed model attack the pretrained speaker recognition model?

- ✓ Targeted attack
- ✓ Average success rate 72.1%
- ✓ Average SNR 57.64dB
- ✓ Average PESQ 3.48

Target ID	PTR(%)↑	SNR(dB)↑	PESQ↑
0	91.4	57.55	3.36
100	89.3	56.83	3.16
200	63.3	58.42	3.69
300	58.7	56.92	3.52
400	57.6	58.36	3.68
avg	72.1	57.64	3.48



# Experimental Results

◆ Does our design work? (the phoneme recognition model)

✓ With fixed  $\lambda_{norm}$ , larger  $\lambda_{phn}$  results a higher SNR and PESQ

✓ The phoneme brunch works for obtaining a trade-off between SER and SNR/PES

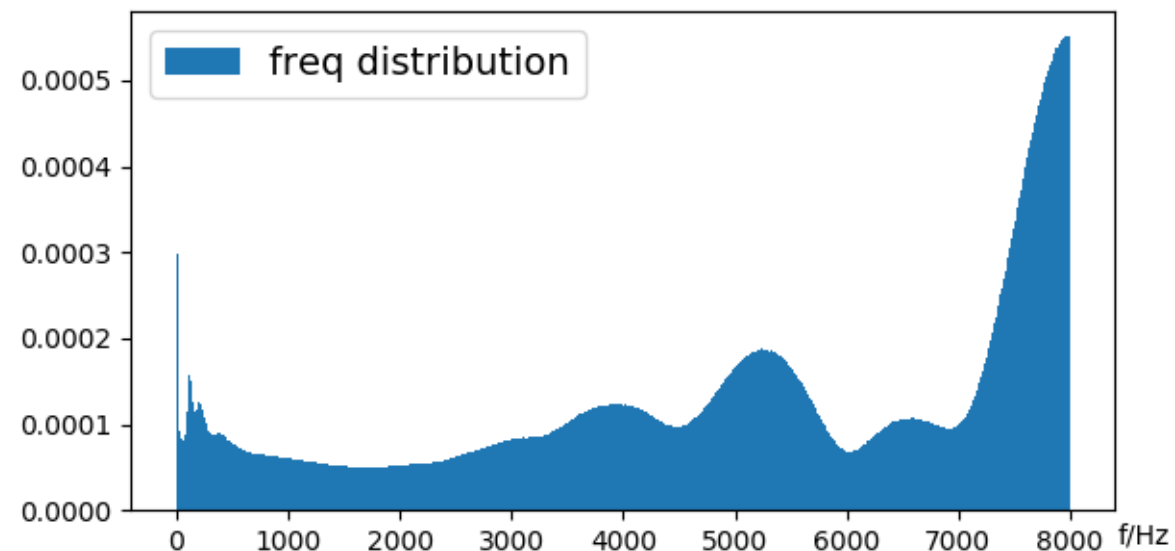
$\lambda_{phn}$	$\lambda_{norm}$	SER(%)↑	SNR(dB)↑	PESQ↑
-	-	1.52*	-	-
0	0	99.7	18.56	1.09
0	1000	96.5	56.39	3.72
0	2000	86.7	57.79	3.61
1	1000	<b>99.2</b>	57.20	4.20
5	1000	93.9	58.00	4.25
10	1000	90.5	<b>59.01</b>	<b>4.28</b>



# Experimental Results

## ◆ Other findings

- ✓ The perturbations concentrate on high frequency
- ✓ Can we design robust speaker recognition models focusing on the low frequency? (future works)



Perturbations distribution



# The questions

- ◆ How to attack the speech-based biometric systems?
  - ✓ Our proposed framework successfully attacked the SOTA speaker recognition model
- ◆ Is the speech-based biometric systems vulnerable to the attacker?
  - ✓ Yes
- ◆ Is it possible to design a biometric systems robust to the adversarial attack?
  - ✓ The future works





# Thanks

## Q & A

Codes, data and more results: <https://smallflyingpig.github.io/speaker-recognition-attacker/main>

Paper early access: <https://ieeexplore.ieee.org/document/9053058>

