# Can I Hear Your Face? Pervasive Attack on Voice Authentication Systems with a Single Face Image
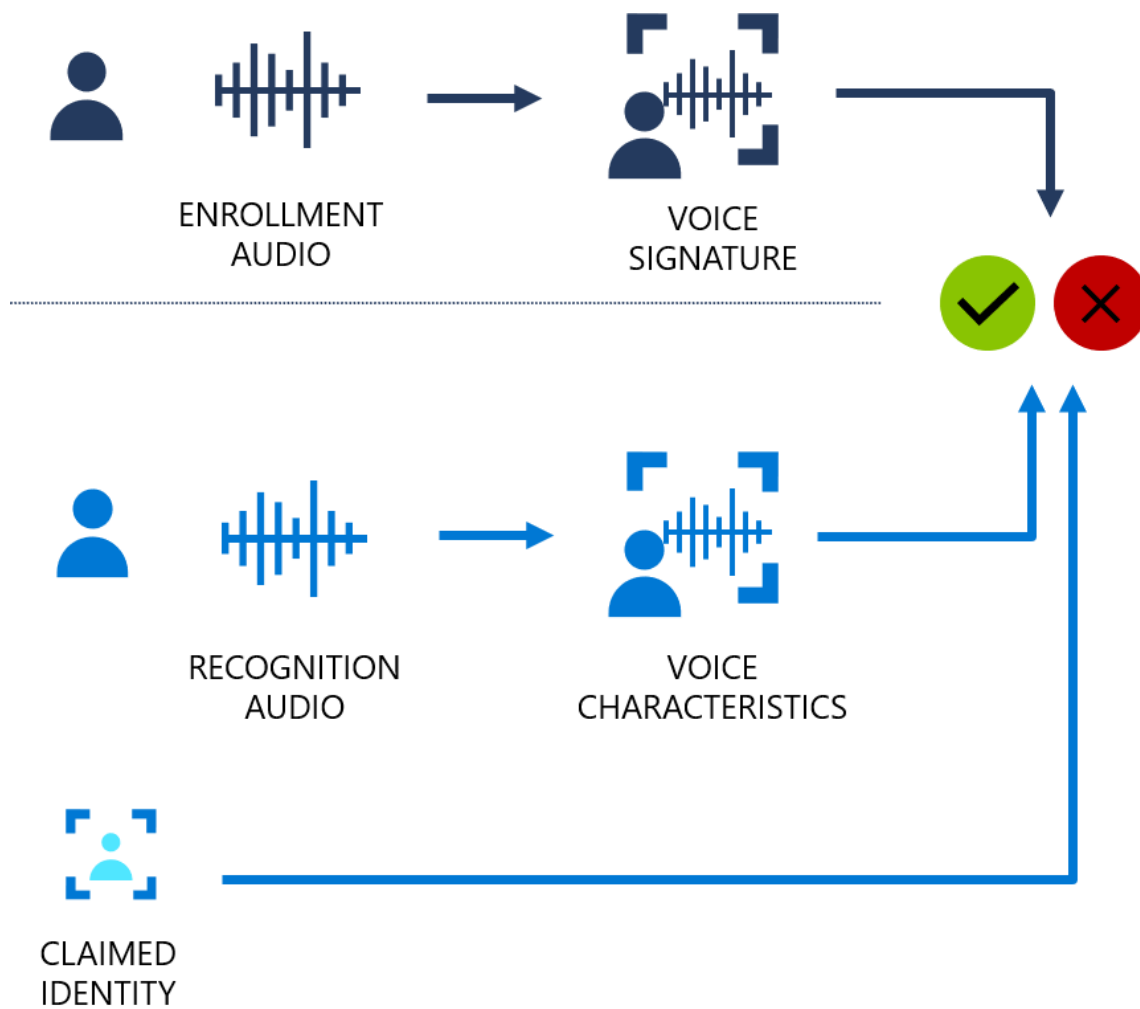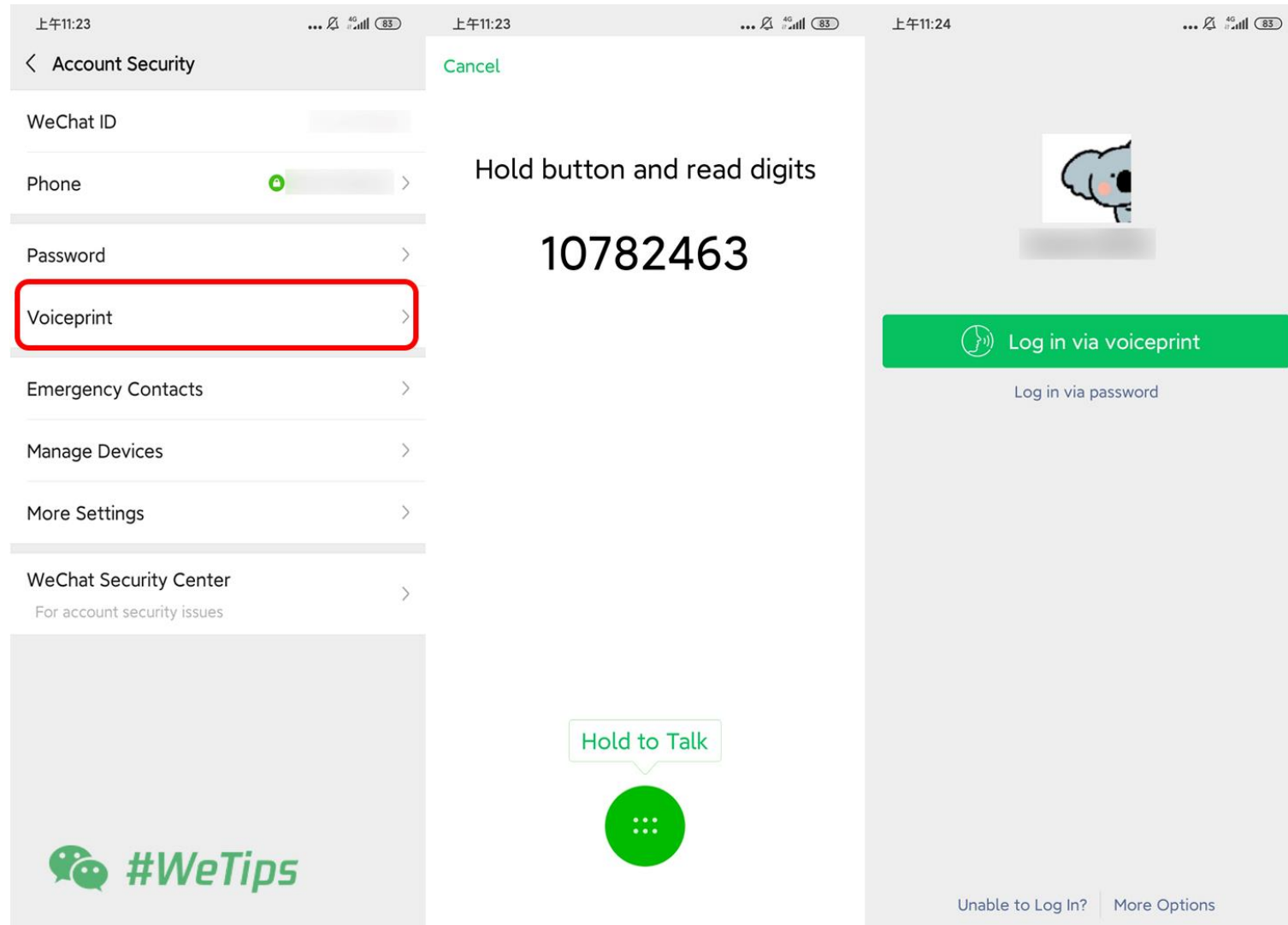
Jiang et al.
USENIX '24

# Voice Authentication?
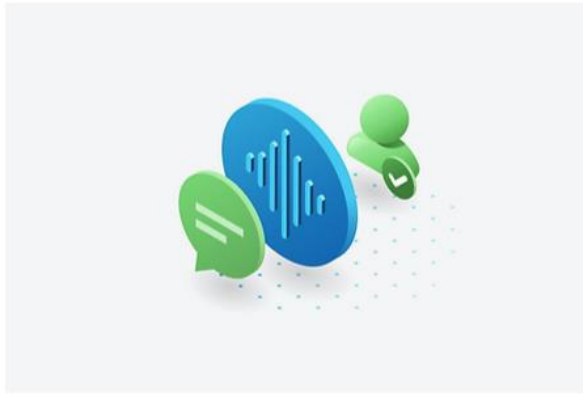


ENROLLMENT
AUDIO

VOICE
SIGNATURE

RECOGNITION
AUDIO

VOICE
CHARACTERISTICS

CLAIMED
IDENTITY

https://learn.microsoft.com/en-us/azure/ai-
services/speech-service/speaker-recognition-overview

# Voice Authentication is widely used

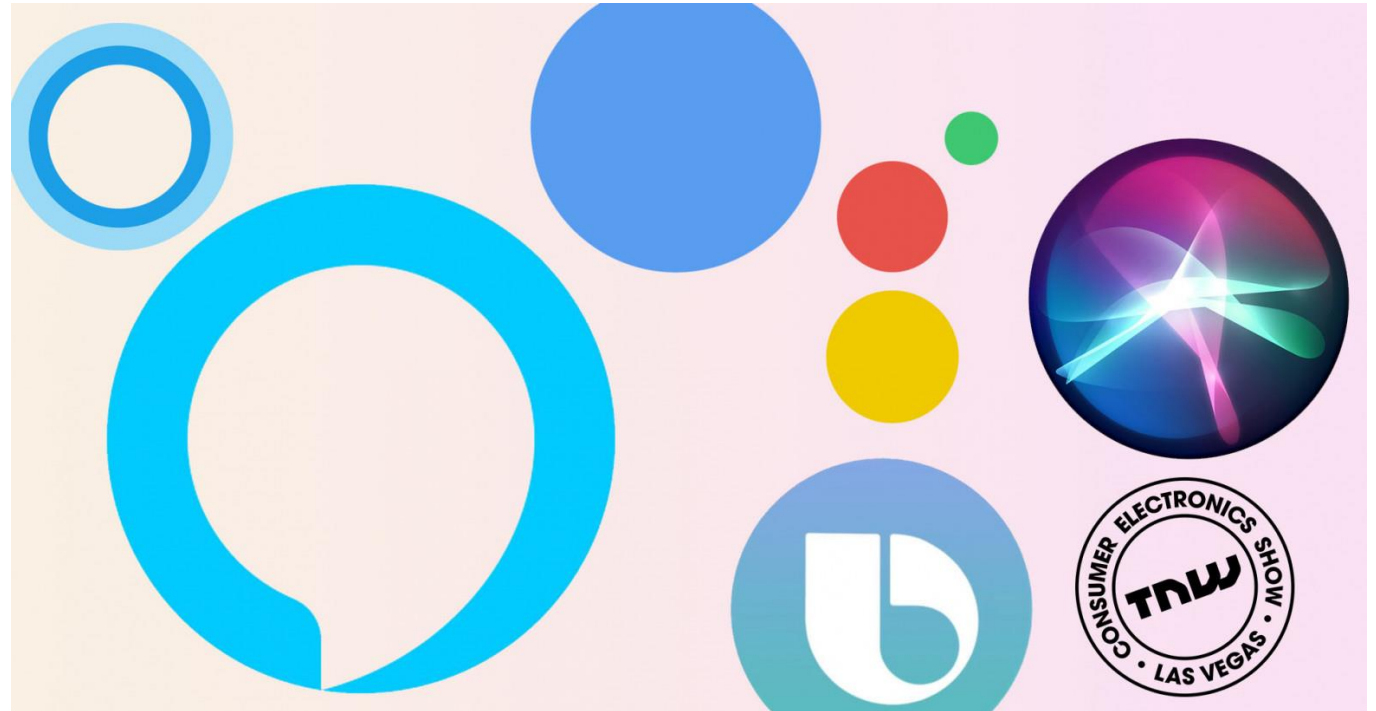# Voice Authentication is widely used



### Verify and recognize speakers

Confirm a person's identity or recognize who's speaking in a meeting by adding speaker verification and identification to your app.

Learn more about Speaker Recognition >

Learn how to recognize speakers in your app >

https://thenextweb.com/news/its-about-time-alexa-and-google-assistant-learn-to-play-nice

# Audio Deepfake



https://www.youtube.com/playlist?list=PLpBfg7XImIEgHy
uStO0g3F7t8xnbnqv8C

# What if..?

Generate
Synthetic Voice

Tries to
unlock App

Deepfake model
trained with
IU's voice

# Okay, but if there's no voice samples?

Can't been trained..

Generate
Synthetic Voice

Tries to
unlock App

## Deepfake model
## trained..?

# Can I hear your face?

Facial structure

Generate
Synthetic Voice

Tries to
unlock App

Foice,
suggested
model

# Brief Overview



Jiang et al., Can I Hear Your Face? Pervasive Attack on
Voice Authentication Systems with a Single Face Image,
USENIX'24

# How to generate synthetic voice?



(a) Voice Feature Extraction

Hello, how are you?

Input: Victim's Voice Recording → Speaker Encoder → Output: Voice Feature Vector

(b) Voice Synthesis

Input: Voice Feature Vector + Hey, Google! → Synthesizer → Output: Synthesized Voice Recording (Hey, Google!)

Jiang et al., Can I Hear Your Face? Pervasive Attack on
Voice Authentication Systems with a Single Face Image,
USENIX'24

# How to generate synthetic voice?

| English | Korean | Chinese |
|---------|--------|---------|
| Nasal cavity | 비강 | 鼻腔 |
| Oral cavity | 구강 | 口腔 |
| Vocal cords | 성대 | 声带 |
| Larynx | 후두 | 喉 |



Jiang et al., Can I Hear Your Face? Pervasive Attack on
Voice Authentication Systems with a Single Face Image,
USENIX'24

# Few things to consider

- Voice Authentication
  - Enrollment Phase(i.e. system displays ###)
  - Authentication Phase(i.e. Speak displayed #)
- System set <span style="color:red">low threshold</span> to ensure user can authenticate in noisy environments.
- Many voice authentication system <span style="color:red">do not restrict</span> the number of authentication attempts.

# System Design



Jiang et al., Can I Hear Your Face? Pervasive Attack on
Voice Authentication Systems with a Single Face Image,
USENIX'24

# Face & Voice Processing



Jiang et al., Can I Hear Your Face? Pervasive Attack on
Voice Authentication Systems with a Single Face Image,
USENIX'24

# Face-Dependent Feature Generator

- F: Feature
  - GT: Ground Truth
  - dep: Face-Dependent
- E: Encoder
  - face: Encoder trained with facial image
- C: Converter
  - f–⟩v: facial feature to F_dep



**Face-dependent Voice Feature Extractor** (§4.3)

**(a) Training Phase**

Input — $Img_{face}$ → $E_{face}$ → $F_{face}$ → $C_{f \to v}$ → Output $F_{dep}$

$F_{GT}$

Minimize Distance

**(b) Attack Phase**

Input — $Img_{face}$ → $E_{face}$ (Trained Encoder) → $C_{f \to v}$ (Trained Converter) → Output $F_{dep}$

$$F_{face} = E_{face}(Img_{face}), \qquad F_{dep} = C_{f \to v}(F_{face}),$$

$$\min_{E_{face}(\cdot), C_{f \to v}(\cdot)} Err(F_{dep}, F_{GT})$$

Jiang et al., Can I Hear Your Face? Pervasive Attack on Voice Authentication Systems with a Single Face Image, USENIX'24

# Face-Independent Feature Generator



Face-independent Voice Feature Generator (§4.4)

**(a) Training Phase**

Input / Output

$F_{dep}$

$F_{GT}$

B

Search Space

$F_{indep} \sim N(0, I)$

R

$F_{recon}$

Minimize Distance

**(b) Attack Phase**

Input / Output

$F_{dep}$

$\{F_{indep_i}\}$

N Samples from Gaussian Distr.

$\sim N(0, I)$

R

Trained Reconstructor

$\{F_{recon_i}\}$

xN

- F: Feature
  - GT: Ground Truth
  - dep: Face-Dependent
  - recon: reconstructed
  - indep: Face-Independent
- B: Bottleneck
  - Project F to smaller indep. search space, tries F to follow gaussian distribution
- R: Reconstructor
  - Search space to reconstructed F

$$F_{indep} = B(F_{GT}), \qquad F_{recon} = R(F_{indep}, F_{dep})$$

$$\min_{B(\cdot), R(\cdot, \cdot)} Err(F_{GT}, F_{recon}) + KL[P_{F_{indep}}(\cdot) \| \mathcal{N}(0, I)],$$

Jiang et al., Can I Hear Your Face? Pervasive Attack on Voice Authentication Systems with a Single Face Image, USENIX'24

# Experimental Setup

- Voice Authentication Systems
  - On-device: Commercial System installed on smartphones
  - Cloud: Microsoft, iFlytek, VGGVox, DeepSpeaker

| Category | System | System Type | Commercial/ Academic |
|---|---|---|---|
| On-Device System | WeChat | Authentication | Commercial |
| | Siri | Voice Assistant | Commercial |
| | Google Assistant | Voice Assistant | Commercial |
| | Bixby | Voice Assistant | Commercial |
| Cloud Service | Microsoft API | Authentication | Commercial |
| | iFlytek API | Authentication | Commercial |
| | VggVox | Authentication | Academic |
| | DeepSpeaker | Authentication | Academic |

# Experimental Setup

- Benchmark Voice Deepfake System:
  - SV2TTS: SoTA TTS System can produce voice of the un-seen speaker in training phase "naturally".
- Speaker Dataset:
  - VoxCeleb1(100K Videos, 1251 celebs) – For evaluation
  - VoxCeleb2(1M Videos, 6112 celebs) – For Training
  - 10 Participants recorded in quiet environment.
- Performance Metrics
  - Overall Success Rate: Percentage of speakers with at least one successful voice attack
  - Individual SR: Percentage of successful synthetic voice attacks for specific person
  - Foice Individual SR: Fraction of voice cloned from single face image passes verification

- Is Foice attack effective against diverse modern implementations of speaker authentication systems and voice assistants?
- Can Foice provide more voice information other than age and gender?
- Can we combine Foice and the existing voice deepfake system to improve the attack's effectiveness?
- How do different experimental conditions affect the effectiveness of Foice?

# RQ1: Is Foice attack effective against diverse modern implementations of speaker authentication systems and voice assistants?

- Method
  - Used custom dataset(10)
  - Foice: 100 synthetic voice recordings per participants
  - SV2TTS: 1 synthetic voice recording per participant
  - Laptop Speaker played synthetic recording to cell phone

| Category | System | System Type | Commercial/ Academic | Eval. Param. | | Overall Success Rate | | | Average Individual Success Rate (Foice) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | #Spk. | Threshold | SV2TTS [31] | Foice | Augmentation Attack (Foice + SV2TTS) | |
| On-Device System | WeChat | Authentication | Commercial | 10 | — | 50.0% | **100%** | — | 29.7% |
| | Siri | Voice Assistant | Commercial | 10 | — | 50.0% | **70.0%** | — | 40.9% |
| | Google Assistant | Voice Assistant | Commercial | 10 | — | 50.0% | **60.0%** | — | 10.3% |
| | Bixby | Voice Assistant | Commercial | 10 | — | 30.0% | **50.0%** | — | 3.6% |

RQ1: Is Foice attack effective against diverse modern implementations of speaker authentication systems and voice assistants?: On-device

- Method
  - Used custom dataset(10)
  - Foice: 100 synthetic voice recordings per participants w\ img
  - SV2TTS: 1 synthetic voice recording per participant
  - Laptop Speaker played synthetic recording to cell phone
- Analysis
  - SV2TTS struggles, input voice is not noise-free(i.e. echo).

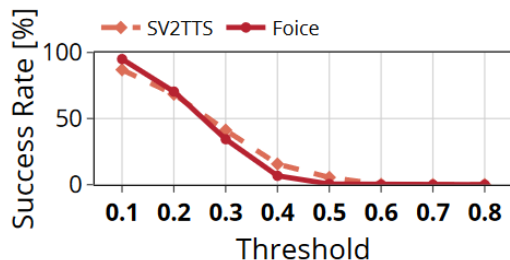| Category | System | System Type | Commercial/ Academic | Eval. Param. | | Overall Success Rate | | | Average Individual Success Rate (Foice) |
|---|---|---|---|---|---|---|---|---|---|
| | | | | #Spk. | Threshold | SV2TTS [31] | Foice | Augmentation Attack (Foice + SV2TTS) | |
| On-Device System | WeChat | Authentication | Commercial | 10 | — | 50.0% | 100% | — | 29.7% |
| | Siri | Voice Assistant | Commercial | 10 | — | 50.0% | 70.0% | — | 40.9% |
| | Google Assistant | Voice Assistant | Commercial | 10 | — | 50.0% | 60.0% | — | 10.3% |
| | Bixby | Voice Assistant | Commercial | 10 | — | 30.0% | 50.0% | — | 3.6% |

# RQ1: Is Foice attack effective against diverse modern implementations of speaker authentication systems and voice assistants?: Cloud
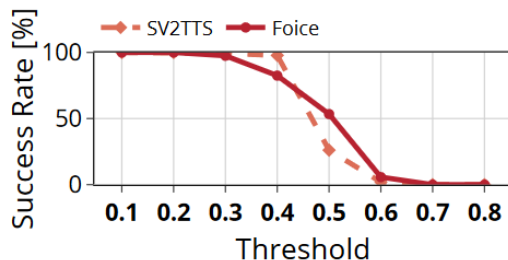
- Method
    - Used VoxCeleb1 dataset(1029)
    - Foice: Synthetic voice recordings w\ celeb's img
    - SV2TTS: Synthetic voice recording from new voice recording
- Analysis
    - SV2TTS ~= Foice to commercial APIs
    - SV2TTS ⟨ Foice on Academic models

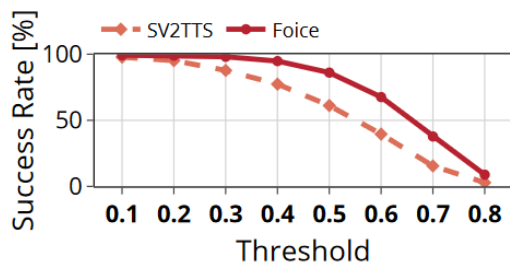| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Microsoft API | Authentication | Commercial | 597 | 0.1 - 0.8 | 0% - 86.9% | 0% - 95.0% | **0% - 99.6%** | 0% - 29.5% |
| | iFlytek API | Authentication | Commercial | 1021 | 0.1 - 0.8 | 0% - 100% | 0% - 100% | **0% - 100%** | 0% - 99.5% |
| | VggVox | Authentication | Academic | 1029 | 0.1 - 0.8 | 2.9% - 97.8% | 8.9% - 99.3% | **26.1% - 99.9%** | 2.3% - 84.6% |
| Cloud Service | DeepSpeaker | Authentication | Academic | 1029 | 0.1 - 0.8 | 0.2% - 99.5% | 0.5% - 100% | **2.4% - 100%** | 1.0% - 99.4% |

# RQ1: Is Foice attack effective against diverse modern implementations of speaker authentication systems and voice assistants?: Cloud
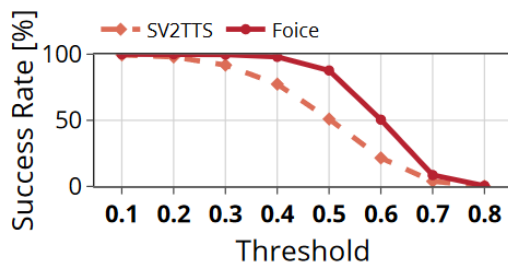


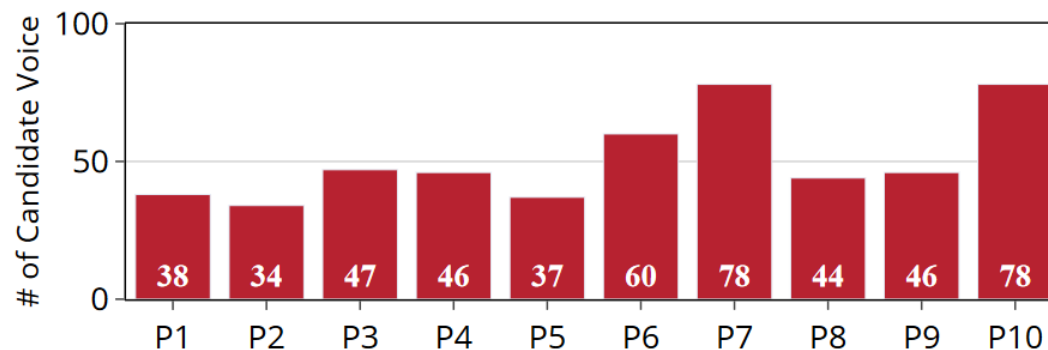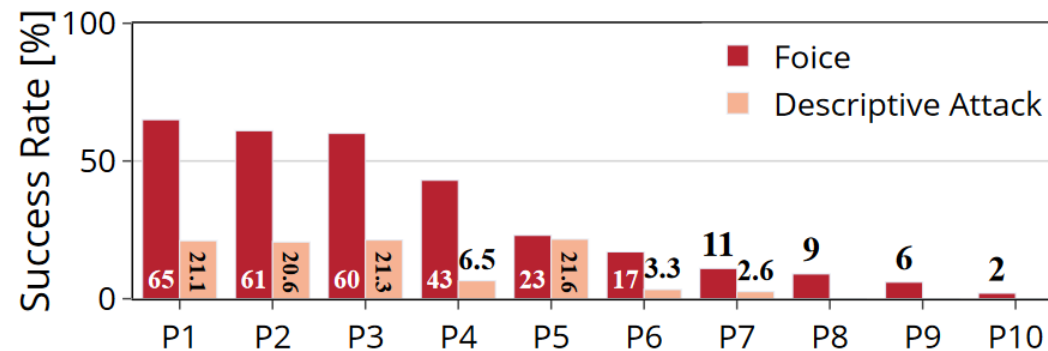(a) Microsoft

(b) iFlytek

(c) VGGVox

(d) DeepSpeaker

| System | Default/Optimal Threshold | Overall Success Rate | | Average Individual Success Rate |
|---|---|---|---|---|
| | | SV2TTS [31] | Foice | |
| Microsoft | 0.5 | **5.5%** | 0.5% | 1% |
| iFlytek | 0.6 | 2.0% | **5.7%** | 3.3% |
| VGGVox | 0.6 | 39.6% | **67.6%** | 15.4% |
| DeepSpeaker | 0.5 | 51% | **87.7%** | 32.7% |

# RQ2: Can Foice provide more voice information other than age and gender?

- Method: Target system: WeChat
  - P1..10: Cluster VoxCeleb1(1029) to 10, using age and gender



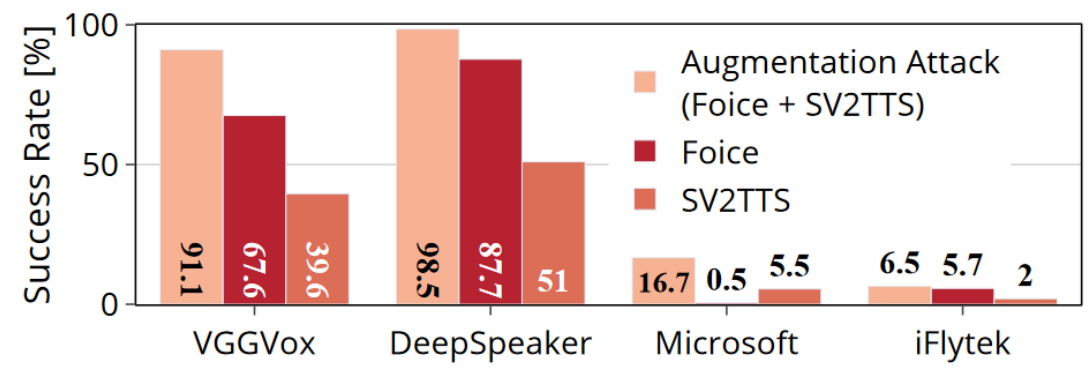(a) Number of candidate voice recordings



(b) *Foice vs. Descriptive Attack*

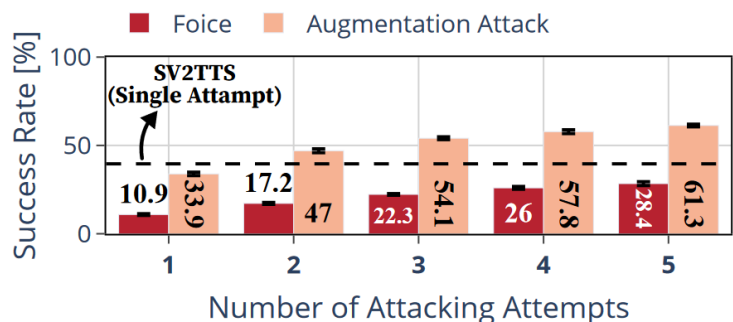# RQ3: Can we combine Foice and the existing voice deepfake system to improve the attack's effectiveness?

- ## Method:
    - Average 100 voice feature vector of SV2TTS and Foice
    - Generate 100 synthetic voice from averaged vector.

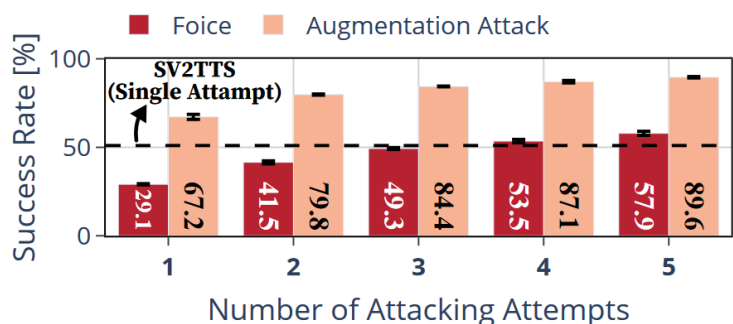| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Microsoft API | Authentication | Commercial | 597 | 0.1 - 0.8 | 0% - 86.9% | 0% - 95.0% | **0% - 99.6%** | 0% - 29.5% |
| iFlytek API | Authentication | Commercial | 1021 | 0.1 - 0.8 | 0% - 100% | 0% - 100% | **0% - 100%** | 0% - 99.5% |
| VggVox | Authentication | Academic | 1029 | 0.1 - 0.8 | 2.9% - 97.8% | 8.9% - 99.3% | **26.1% - 99.9%** | 2.3% - 84.6% |
| DeepSpeaker | Authentication | Academic | 1029 | 0.1 - 0.8 | 0.2% - 99.5% | 0.5% - 100% | **2.4% - 100%** | 1.0% - 99.4% |

Cloud Service

# RQ4: How do different experimental conditions affect the effectiveness of Foice?

- Varying number of attacking attempts, image occlusion, image resolution
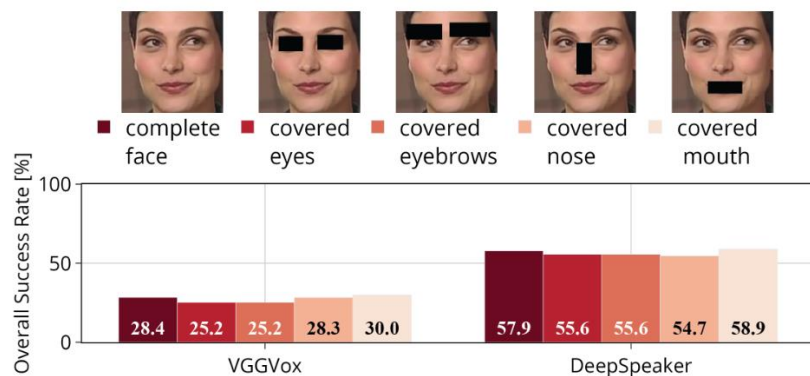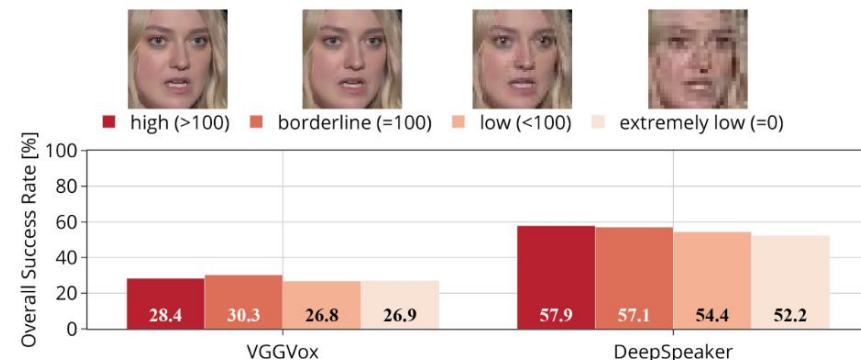


Attacking attempts

Image occlusion

Image resolution

# Discussion

- Voice Authentication is trending, performance will increase with 3D Face photos + Voice (i.e. Face ID + Voice)
- Safety
  - Restricting Login Attempts
  - Deepfake detection algorithm
- Foice shows how "<span style="color:red">vulnerable</span>" current voice authentication systems are; safety enhancement is needed.

Thank You! Any Questions?

# Dimension of B − finding optimal is necessary



Jiang et al., Can I Hear Your Face? Pervasive Attack on
Voice Authentication Systems with a Single Face Image,
USENIX'24