Can I Hear Your Face? Pervasive Attack on Voice Authentication Systems with a Single Face Image

Nan Jiang, Bangjie Sun, Terence Sim, and Jun Han*

National University of Singapore * KAIST



Voice Authentication in Daily Life

Login to Account **Activate Voice Assistant** amazon alexa WeChat Hey Siri Google Assistant "Hey, Google" "43257018" Hi Bob, how Voice can I help? Verified! Activate voice assistant Login to WeChat using registered voice using voice

Voice Deepfake Can Bypass Voice Authentication



Sep

🕒 Jul 28, 2023

CATEGORIES

Fake It Until You Make It: Using Deep Fakes to Bypass Voice Biometrics

Home > News > Security

Deepfake Software Fools Voice Authentication With 99% Success Rate

Creating a fake voice to trick authentication systems has never been so easy or effective.

Limitation of Voice Deepfake Attack

• Require **high-quality** recording of the victim's voice

(a) Download from the network



Limitation of Voice Deepfake Attack

• Require **high-quality** recording of the victim's voice

(a) Download from the network



Limited range of victims



Limitation of Voice Deepfake Attack

• Require **high-quality** recording of the victim's voice



How About Replacing Voice with Face Images?

• Easier to obtain face images

(a) Download from the network



Wide availability of face images



Can we launch a novel attack leveraging only a single image of the victim without requiring voice recordings?

Our Work: Foice

• Synthesizes voice recordings from a single face image



Threat Model

- Attacker's Goal: compromise verification
 - Gain unauthorised access to the victim's private account
 - Execute unauthorised commands on the victim's personal voice assistant

Can You Guess What He Sounds Like?



Background: Face and Voice Correlation

• Facial appearance affects how human voice sounds like



Background: Face Offers Limited Information

Human voice is primarily affected by the inner body structure (e.g., vocal cords, chest)

Vibration here causes **facial** resonance

Vibration here causes **chest** resonance



Core Ideas of Foice



 Core idea (1): <u>extract</u> face-dependent features from face image input



Core Ideas of Foice



 Core idea (1): <u>extract</u> face-dependent features from face image input



 Core idea (2): <u>generate</u> candidate supplementary features



Core Ideas of Foice



Goal: To generate candidate supplementary features



Goal: To generate candidate supplementary features



Goal: To generate candidate supplementary features

(1) Minimize distance



Goal: To generate candidate supplementary features

(1) Minimize distance



Dimension is Just Right





Too Narrow



10

10.

10.

10

10

Goal: To generate candidate supplementary features

Minimize distance





Evaluation Setup

- Evaluate *Foice* with **authentication systems** and **voice assistants**
 - o 100 candidate synthetic audios for each subject in the test dataset



Evaluation Setup

• For each subject, a successful attack is defined as:





Summary of Results

 All tested systems vulnerable to Foice attack



- **Comparable** attack performance to voice deepfake attack
- Attack success rate improved by **more than three times** when combing face and voice

• *Foice* outperforming attacks using only age and gender information



• Foice **robust** to various types of image conditions

Occlusion

Resolution



Main Result

• Achieves **comparable** attack performance to the state-of-the-art voice deepfake attack



Discussion

1 Improving Foice

- **Dynamic** face provides more details on the facial bone structure
- Optimize model structure



2 Countermeasures

- Deepfake detection and liveness detection
- Lack of adoption in real-world systems



Conclusion

- Synthesize **voice** recordings using a single **face** image
- Highlight the need to protect voice verification against deepfake attacks

