InfoMasker: Preventing Eavesdropping Using Phoneme-Based Noise

Peng Huang, Yao Wei, Peng Cheng, Zhongjie Ba, Li Lu, Feng Lin, Fan Zhang, Kui Ren



Eavesdropping with Smart Devices

• Widespread of smart devices equipped with microphone



• Developers are committed for privacy protection



$\left(\right)$	
	Record calls and audio
For ser	providing audio recording vices
	Deny
	Once
AI	low only while using the app



Eavesdropping with Smart Devices

- Still an unsolved problem
 - Third-party operating systems
 - Malicious fake applications
 - Uncontrolled legal recordings
 - Hidden Recorders



- Need to physically block voice eavesdroppers
 - Makes the voice privacy controllable to the users.

Problem Setup

• Application scenario



- Design goals
 - Effectiveness
 - Successfully mislead human ears
 - Successfully mislead automaticspeech-recognition tools
 - Robustness
 - Could not be removed by noise reduction methods
 - User-friendly
 - Should not disturb users

Existing Methods to Jam Microphone

- Electromagnetic interference-based jamming
 - Pros: No disturbance to users
 - Cons: Limited coverage & Affect other devices
- Adversarial example-based jamming
 - Pros: No need for special hardware
 - Cons: No effect to human ear& generalization ability
- Ultrasound-based jamming
 - **Pros**: No disturbance to users & Reasonable coverage



Principle of Ultrasound-Based Microphone Jamming

- Nonlinearity in microphone will cause self-demodulation of input signals.
 - Zhang et al. (2017) inject inaudible voice commands to microphone via ultrasound^[1]
- Nonlinearity in microphone



Principle of Ultrasound-Based Microphone Jamming

- Nonlinearity in microphone will cause self-demodulation of input signals.
 - Zhang et al. (2017) inject inaudible voice commands to microphone via ultrasound^[13]
- Inject audible noise n(t) with inaudible ultrasound



1. High demand for noise energy vs. Limited transmission energy



- 2. Target speech recognition tools (human and ASR) have strong denoising ablility
 - Common noises with limited energy will be easily removed
- Cocktail party effect^[2] in human ear





Human brain can easily **focus on the target speech** in a noisy environment

- 2. Target speech recognition tools (human and ASR) have strong denoising ablility
 - Common noises with limited energy will be easily removed
- Noise reduction methods in ASR



- 2. Target speech recognition tools (human and ASR) have strong denoising ablility
 - Common noises with limited energy will be easily removed
- Both methods rely on the differences between the noise and the speech



Jamming Strategy: Energetic v.s. Informational



Energetic masking: Covering

Origin wave Masked wave

Characteristics

Pros: No need for prior knowledge

Cons: High energy requirement & Easily to remove

- Informational Masking: Disturbing

 Origin Word: desk
 Phonogram: / desk /
 Inject / I / de I sk / → desk? disk?
 - Characteristics

Pros: Low energy requirement & Hard to remove

Cons: Needs prior knowledge

Informational Masking for Human Speech Jamming

- Prior knowledge for jamming human speech
 - Signal structure: a series of phonemes
 - Frequency domain properties: User dependent
 - Fundamental frequency (F0)
 - Timbre
 - Time domain properties : Varying and uncertain

Main idea: Inject phonemes similar to the target speech to disturb it



Phoneme-Based Jamming Noise Design

• Noise structure





System Workflow

- User Registration
 - Get the user's voice features
- Data Augmentation
 - Get enough data for noise generation
- Noise Generation
 - Get the noise
- Jamming
 - Inject the noise to microphone





User Registration

Data Augmentation

Noise Generation



User Registration

- Purpose: Obtain enough phoneme data with similar timbre as the user.
- Extracting from the user's speech is time consuming, and so not practical
- Extract user's voice feature from short registration audios and match speech data from public corpus



Data Augmentation

- Increase the amount of phonemes while retaining similarity with original data
- *Method:* Fine-tune the emotional-related speech properties^[3].

Phonetical	Modification	Emotional Impact			
Properties	Range	\uparrow	\downarrow		
Speech Rate	0.3-1.8	Fear or Disgust	Sadness		
F0 Mean	0.9-1.1	Anger or Happiness	Disgust or Sadness		
F0 Contour	0.7-1.3	Anger or Happiness	Sadness		
Energy	0.5-2.0	-	-		
Sequential Order	-	-	-		

Data Augmentation

- Increase the amount of phonemes while retaining similarity with original data
- *Method:* Fine-tune the emotional-related speech properties^[3].



Noise Transmission

• Lower-sideband modulation to achieve higher transmission energy



Noise Transmission

• Pre-compensation to reduce distortion during transmission



• Estimate $h_1(t)$ and $h_2(t)$, pre-compensate s(t) with $h_1(t) \circledast h_2^{-1}(t)$



System Overview







Evaluation: Experimental Setting

- Speech recognition tools
 - 4 Commercial ASR tools
 - 2 Open-Source ASR tools
 - Human recognition
- Datasets
 - LibriSpeech^[4] for most experiments
 - TIMIT^[5] for training targeted ASRs
 - Harvard Sentences^[6] for human recognition

- Evaluate aspects
 - Effectiveness
 - Robustness
- Scenarios
 - Digital domain
 - Real-world jamming
 - Case study: A common office

Evaluation: Effectiveness

- Digital domain
 - 27000 words for each ASR
 - Compared with [0, 8] kHz bandlimited white noise.
- Real-world jamming



• 70 hours data

SNR	<-4	[-4,-2]	[-2, 0]	[0,2]	[2,4]	>4	Clear
Avg WER(%) Min WER(%)	85.8 68.6	81.6 77.0	77.6 62.4	70.2 62.2	56.4 45.3	42.3 30.3	11.5 -
Digital WER(%)	88.6	85.4	68.8	48.67	28.9	17.0	4.1

Evaluation: Effectiveness

- Comparisions with existing works
 - Two previous works and one commercial device.
 - With the presence of noise reduction methods
- Real-world end-to-end scenario





Evaluation: Robustness

- Speech enhancement method^[10]
 - Makes the distrubed speech harder to be recognized



• Speech Separation^[11]



• Specialized ASR



Evluation: Case Study

• Setting





• Results

Tupos		WER(%	6)	
Types	Phone A	Phone B	Laptop	iPad
Α	98.0	98.2	95.7	99.3
В	98.8	98.4	88.1	93.8
С	98.5	56.4	95.8	98.6
D	95.7	97.7	97.9	95.3
Amplifiers On	25.8	26.3	32.5	32.0
Clear	16.0	7.1	19.9	15.5

Summary

- We propose a new type of noise based on the idea of informational masking
- We conduct an in-depth study of the ultrasound transmission method, and then optimize several aspects to make it practical
- We design and implement a system which is proved to have high effectiveness and high security in real-world scenarios

Thank You!

Peng Huang, Yao Wei, Peng Cheng, Zhongjie Ba, Li Lu, Feng Lin, Fan Zhang, Kui Ren



References

[1] G. Zhang, C. Yan, X. Ji, T. Zhang, T. Zhang, and W. Xu, "DolphinAttack: Inaudible Voice Commands," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, Dallas, Texas, USA, Oct. 2017, pp. 103-117.

[2] D. S. Brungart, "Informational and energetic masking effects in the perception of two simultaneous talkers," *The Journal of the Acoustical Society of America*, vol. 109, no. 3, pp. 1101–1109, Mar. 2001.

[3] R. Cowie *et al.*, "Emotion recognition in human-computer interaction," in *IEEE Signal Processing Magazine*, vol. 18, no. 1, pp. 32-80, Jan 2001.

[4] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An asr corpus based on public domain audio books," in 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2015, pp. 5206-5210.

[5] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, D. S. Pallett, and N. L. Dahlgren, "Timit acoustic-phonetic continuous speech corpus ldc93s1," 1993.

[6] "IEEE recommended practice for speech quality measurements," *IEEE Transactions on Audio and Electroacoustics*, vol. 17, no. 3, pp. 225–246, 1969.