

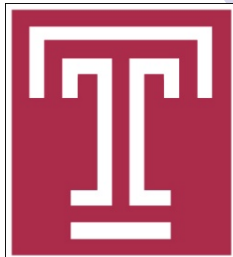
Defending Against Voice Spoofing: A Robust Software-based Liveness Detection System

Jiacheng Shang, Si Chen, and Jie Wu

Center for Networked Computing

Dept. of Computer and Info. Sciences

Temple University

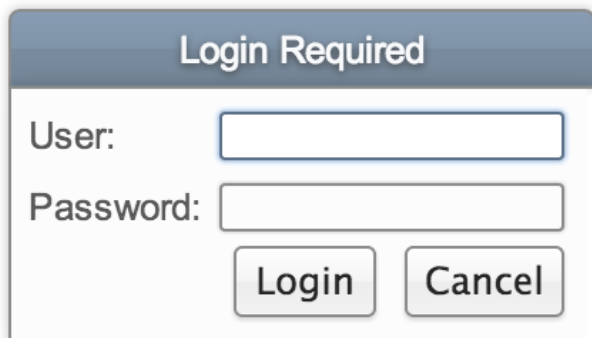


Biometrics: Voiceprint

- Voiceprint

- Promising alternative to password
- Primary way of communication
- Better user experience
- Integration with existing techniques for multi-factor authentication

Applications



Login Required

User:

Password:

Login Cancel



citibank



HSBC

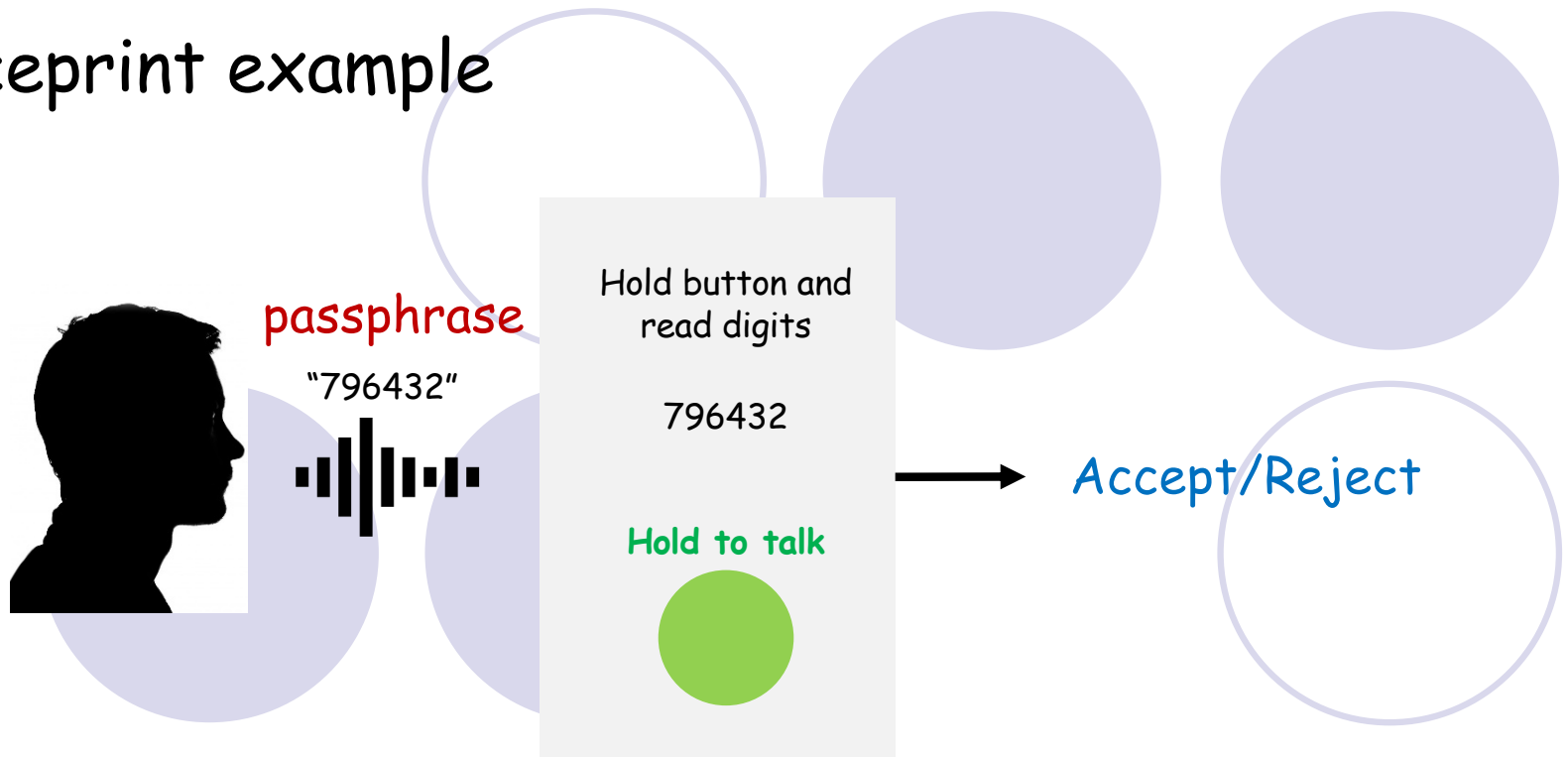
Google

lenovo

SayPay
TECHNOLOGIES, INC.

Biometrics: Voiceprint

- Voiceprint example



Voiceprint-based authentication

Threats

- Human voice is often exposed to the public
- Attackers can "steal" victim's voice with recorders
- Security issues
 - E.g. Adversary could impersonate the victim to spoof the voice-based authentication system



Reverse Turing Test

CAPTCHA

Completely Automated Public Turing test to tell Computers and Humans



voice


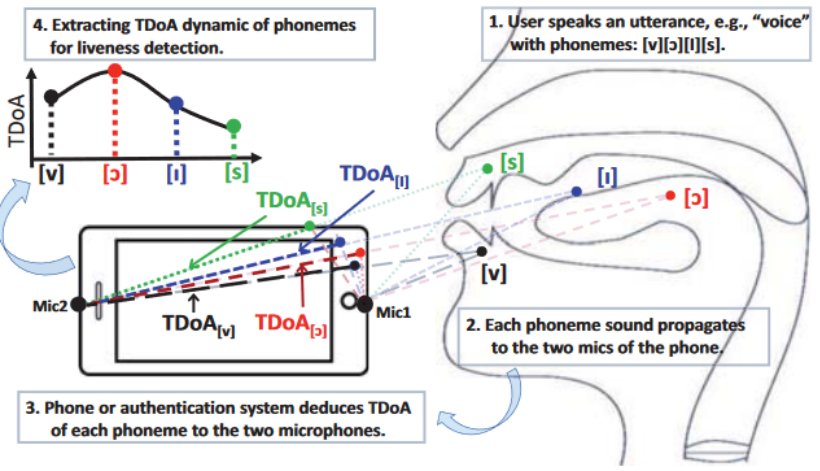


or

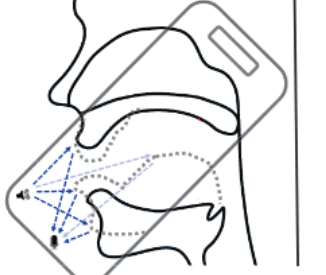
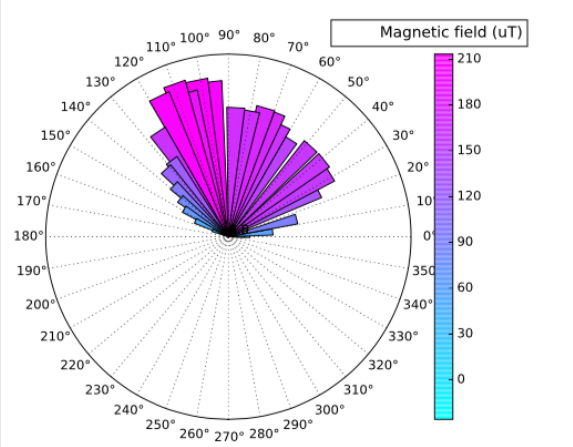


Voiceprint-based authentication

Previous work

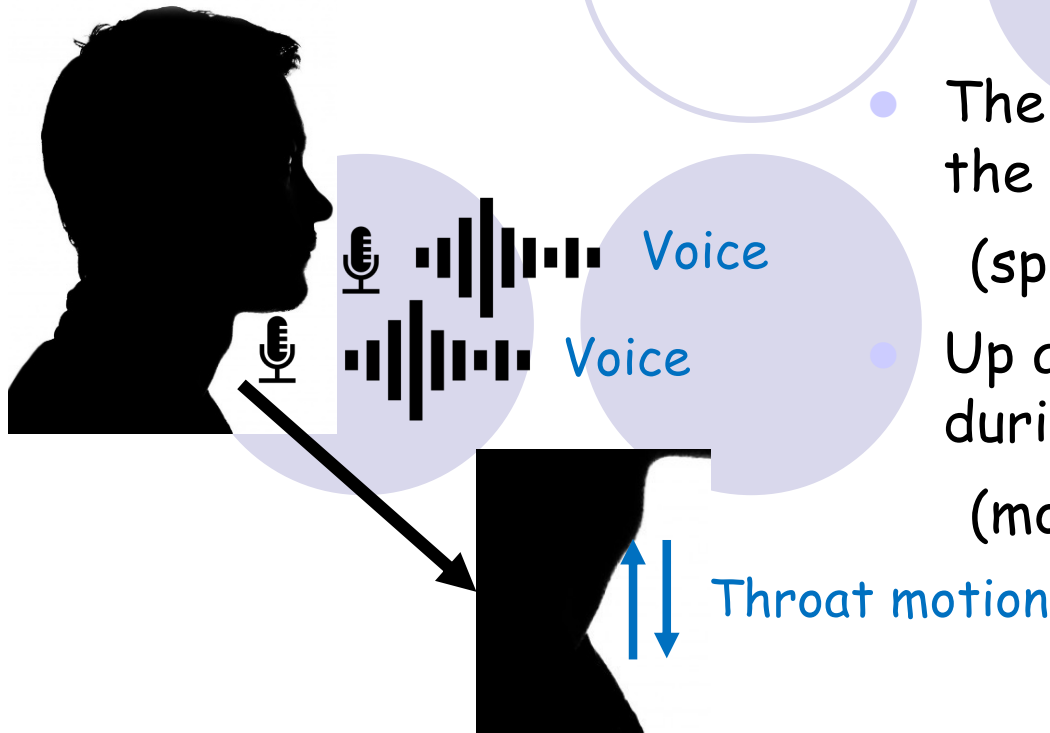
Systems	Limitations
<p data-bbox="170 545 905 597">Automatic speaker verification</p> 	<ul data-bbox="1052 545 1885 699" style="list-style-type: none">• Verifying the speaker's identity (Bob or Alice)• Cannot defend against replay attack
<p data-bbox="170 821 821 930">Phoneme localization-based liveness detection (distance)</p> 	<ul data-bbox="1052 821 1892 971" style="list-style-type: none">• Low true acceptance rate (TAR): the smartphone needs to be static relative to the mouth <p data-bbox="1052 1045 1871 1182">VoiceLive: A Phoneme Localization based Liveness Detection for Voice Authentication on Smartphones (L. Zhang et al. CCS 2016)</p>

Previous work

Systems	Limitations
<p data-bbox="170 488 978 602">Articulatory gesture-based liveness detection (e.g. lip motion)</p>  <p data-bbox="594 688 947 740">(Doppler effect)</p>	<ul data-bbox="1052 488 1892 634" style="list-style-type: none">• Low true acceptance rate (TAR): the smartphone needs to be static relative to the mouth <p data-bbox="1052 708 1818 846">Hearing Your Voice Is Not Enough: An Articulatory Gesture Based Mobile Voice Authentication (L. Zhang et al. CCS 2017)</p>
<p data-bbox="170 964 968 1078">Leveraging the magnetic fields of loudspeakers</p> 	<ul data-bbox="1052 964 1839 1224" style="list-style-type: none">• Low TAR: cannot work if magnetic noise exists• Low true rejection rate (TRR): cannot work if the attacker uses non-conventional loudspeaker <p data-bbox="1052 1292 1776 1471">You Can Hear But You Cannot Steal: Defending against Voice Impersonation Attacks on Smartphones (S. Chen et al. ICDCS 2017)</p>

Basic idea

- Leveraging the structural differences between the vocal systems of human and loudspeakers

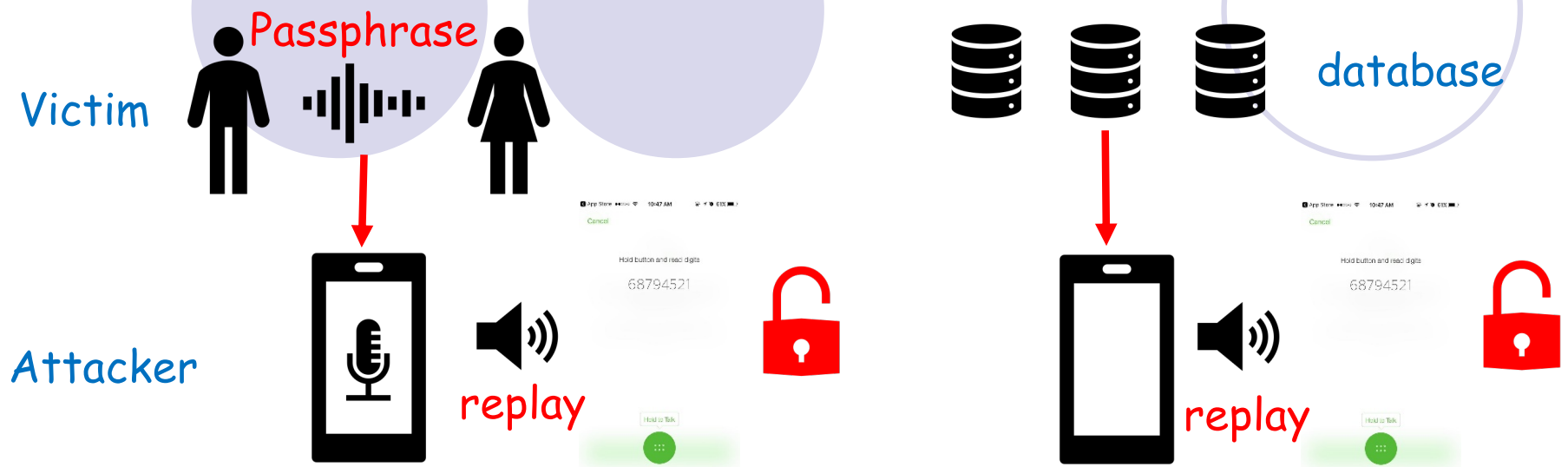


- The **voices** at the mouth and the throat are different (spectrum-based approach)
- Up and down **motions** exist during speaking (motion-based approach)

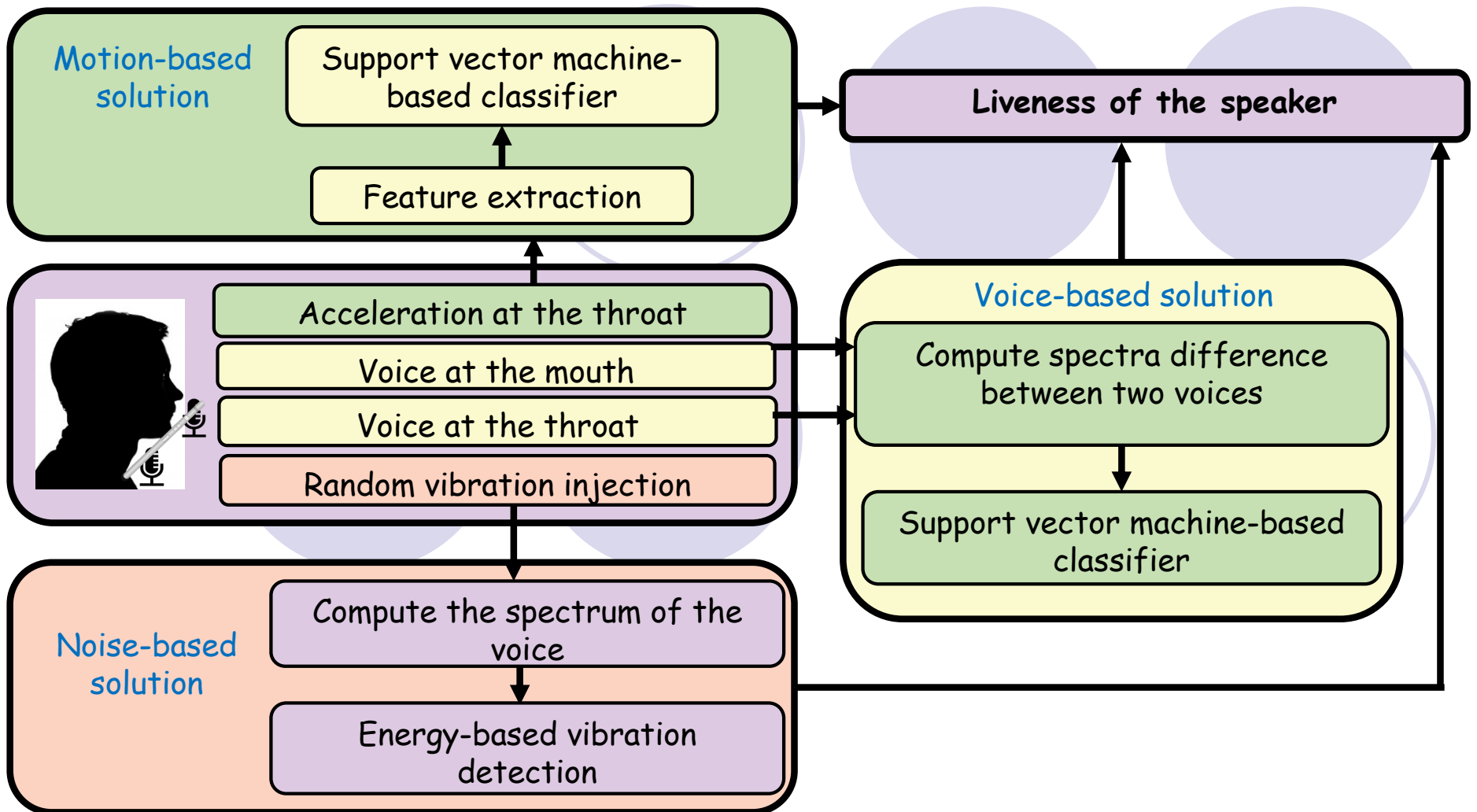
Attack model

- Attack model:

- A simple replay attack: only stealing victim's **voice at the mouth** and replaying it
- A strong replay attack: stealing victim's **throat motions** and **voices at both mouth and throat** from the database and replaying it

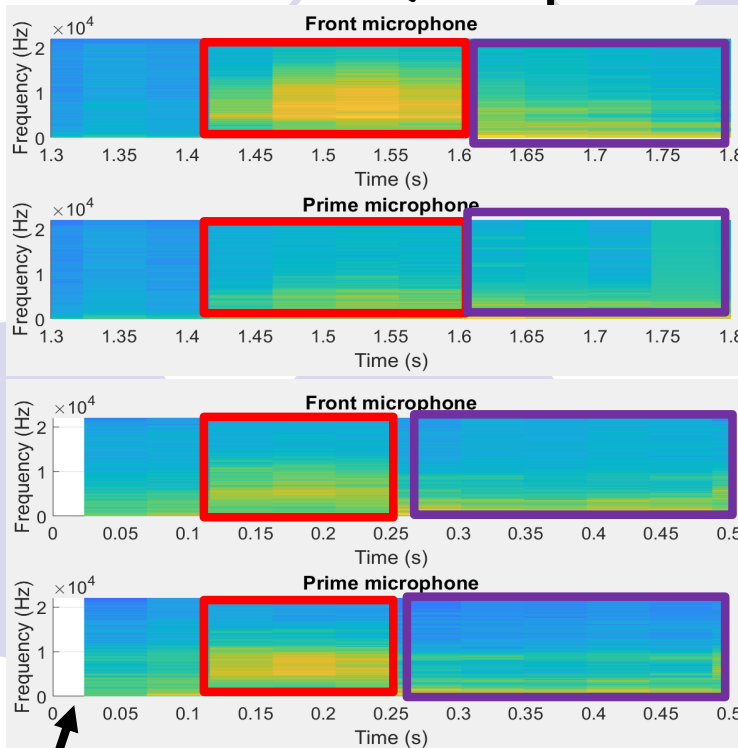
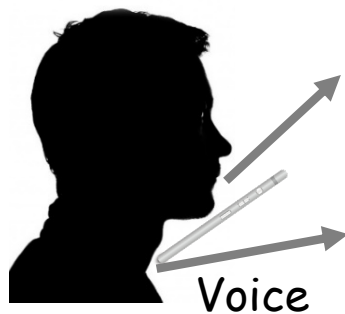


System Architecture

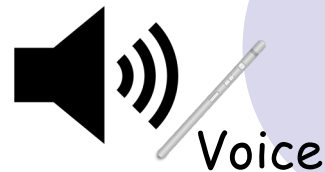
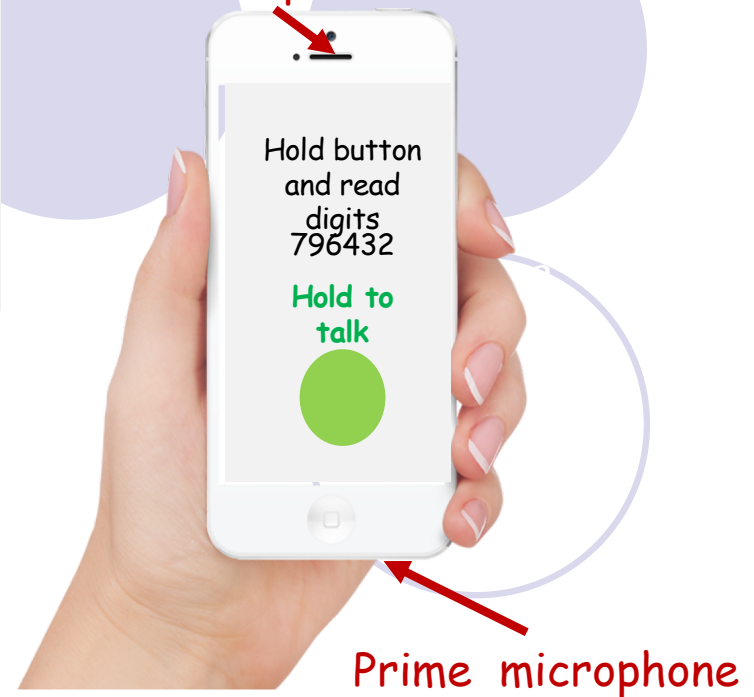


Proposed solutions

- Voice-based solution (Simple attack model)



Front microphone



Computing the spectra using Short-time Fourier transform (STFT)

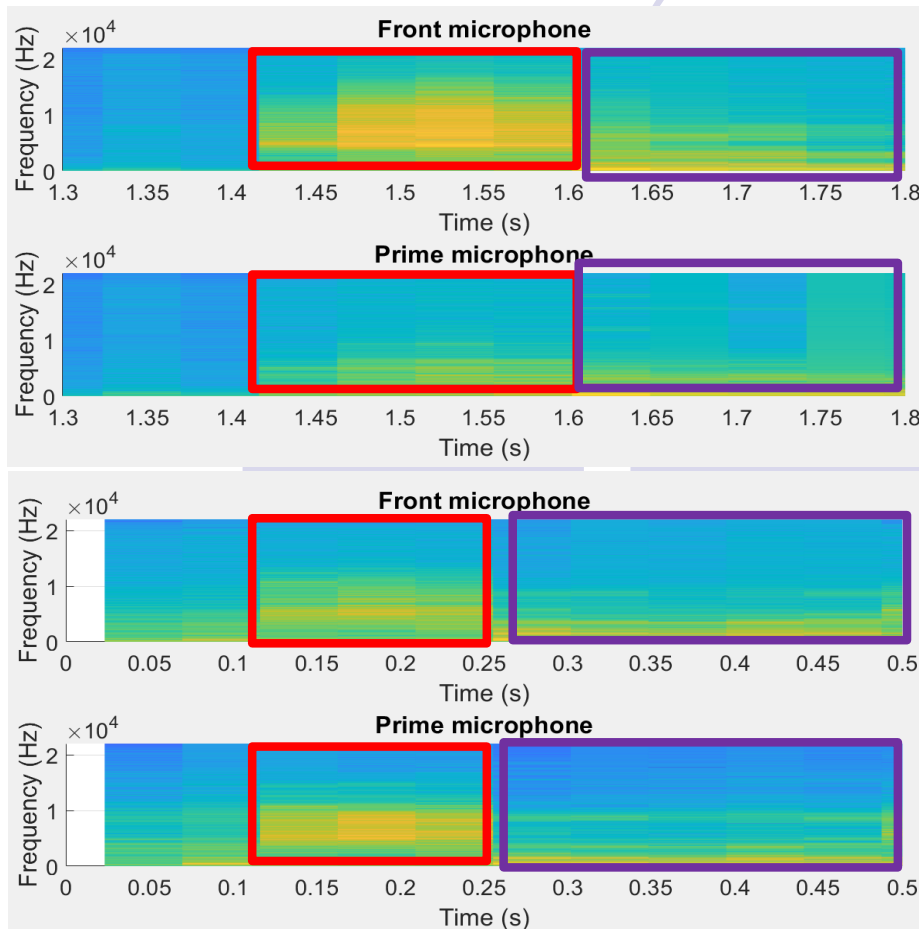
$$\text{spectrogram}\{x[t]\}(m, \omega) = \left| \sum_{n=-\infty}^{\infty} x[n]w[n-m]e^{-j\omega n} \right|^2$$

Convolution
Time domain to frequency domain

$x[n]$: voice $w[n]$: window ω : angular frequency

Proposed solutions

- Voice-based solution for simple attack



- Normal user: two voices are different

- The voice (prime microphone) **does not** contain information of the **unvoiced part**.

- The voice (prime microphone) contains **low-frequency** information of the **voiced part**.

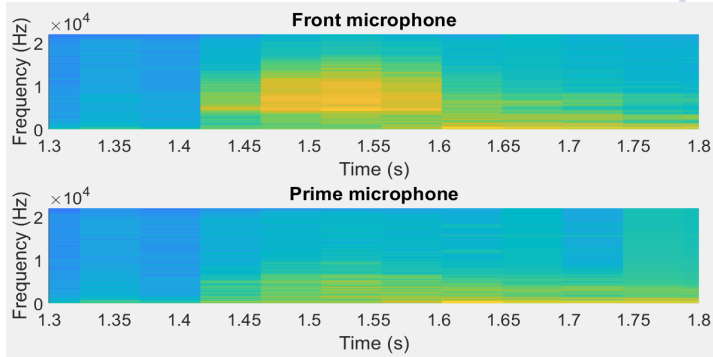
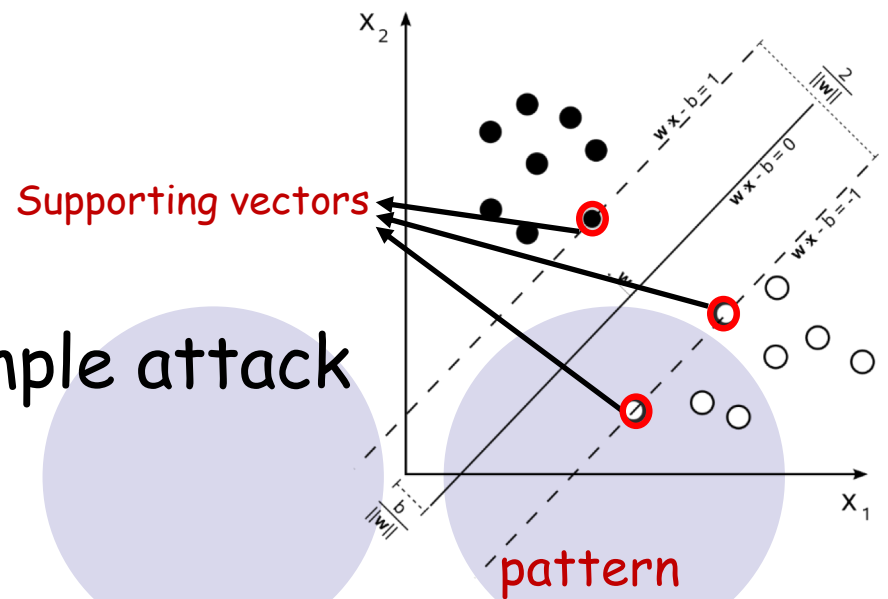
- Attacker: two voices are similar

- The voice (prime microphone) contains information of the **unvoiced part**.

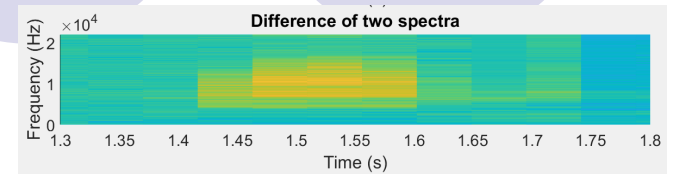
- The voice (prime microphone) contains **most** information of the **voiced part**.

Proposed solutions

- Voice-based solution for simple attack
 - For normal users



Spectra difference

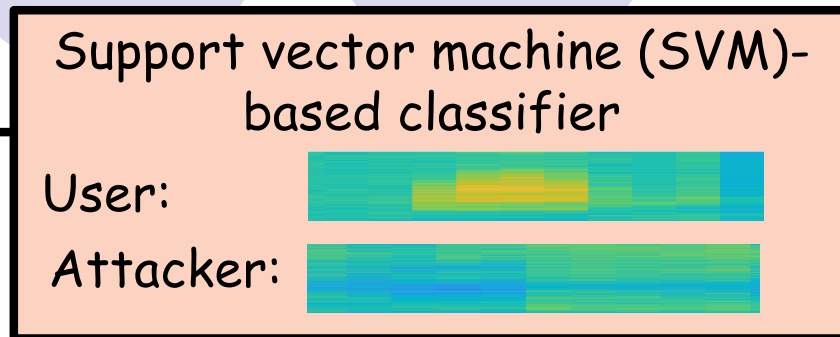


Converting to vector

Input [32,54,3,.....,34,76]

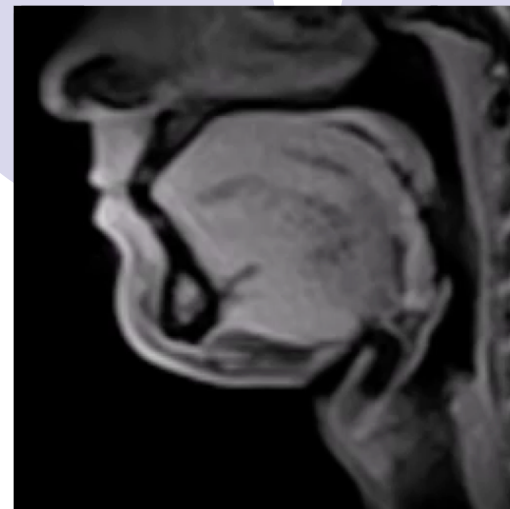
Vector

Accept

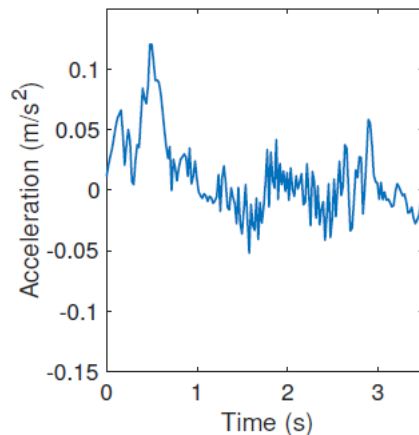


Proposed solutions

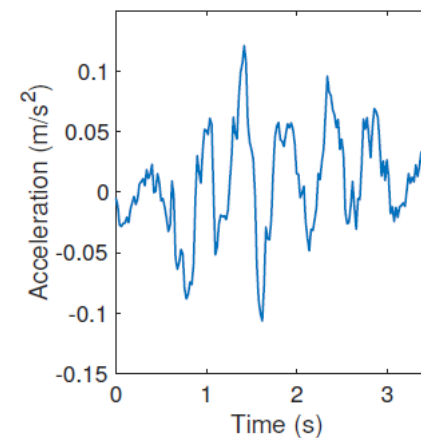
- Motion-based solution for simple attack
 - Using accelerator to capture throat motions
 - 7 features: Variance, minimum, maximum, mean, skewness, kurtosis, standard deviation
 - SVM-based classification model for decision



Normal users

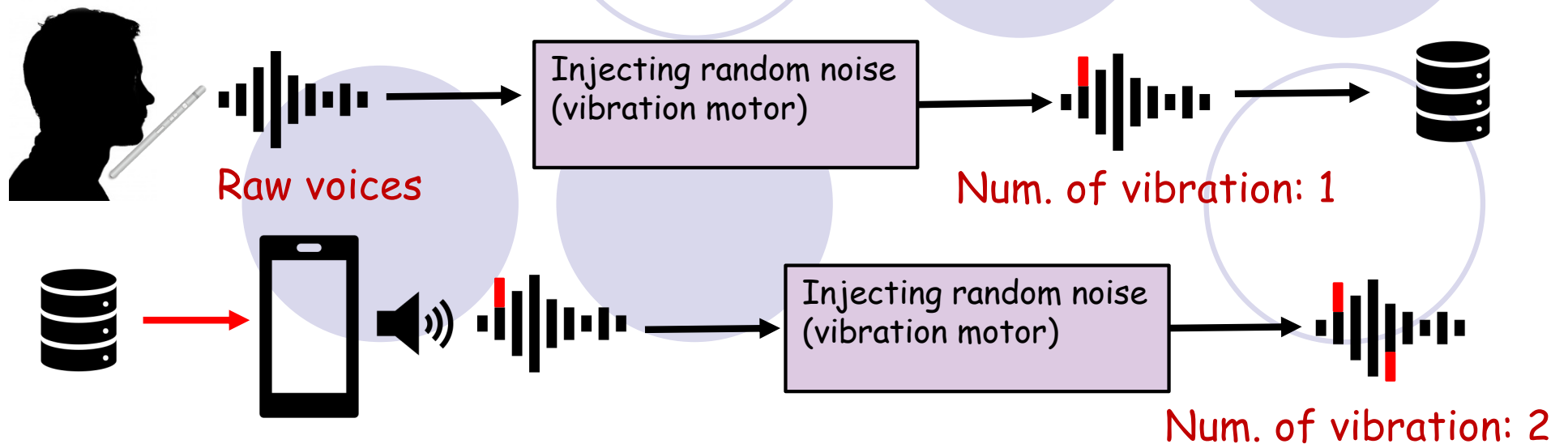


Attacker



Proposed solutions

- Random noise-based solution for strong attack
 - Attackers who can steal victim's voices and throat motions from the database and use multiple loudspeakers to imitate the victim

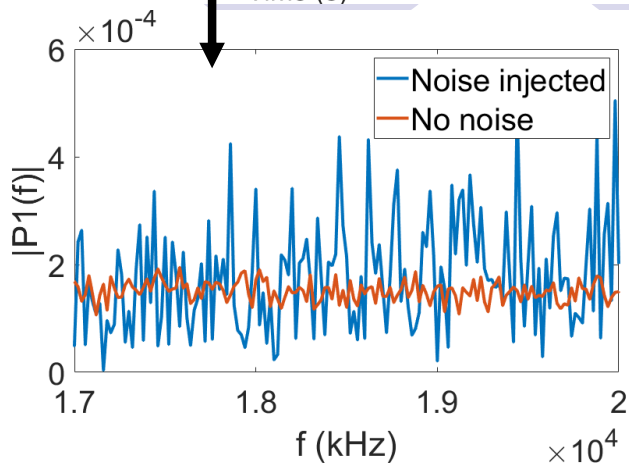
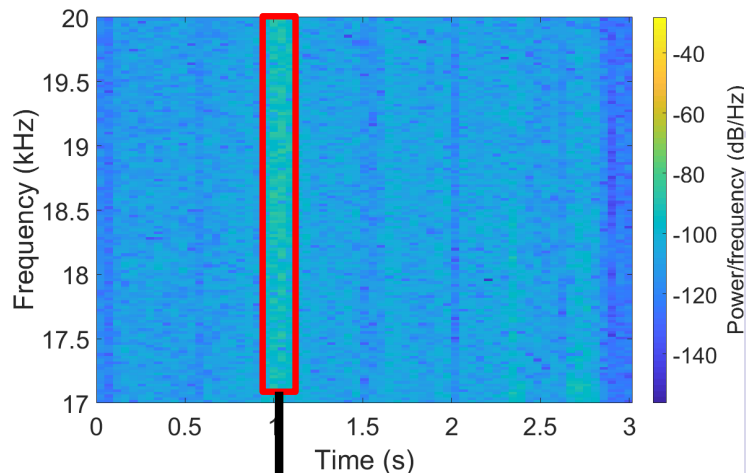


- Our solution:
 - Injecting a random vibration while the user is speaking
 - Checking the number of vibration in the voices

Proposed solutions

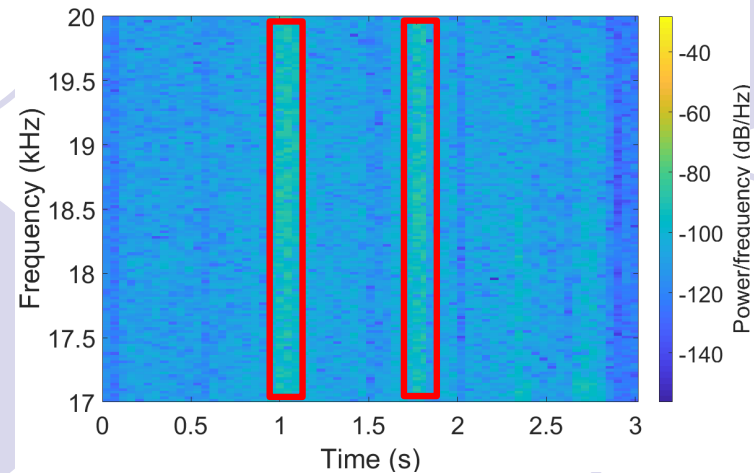
- Random noise-based solution

- For normal users



Computed by STFT

- For the attacker



- The vibration introduces **high energy** to the **high-frequency band**.
- A vibration is detected if the **energy** of a moving window **exceeds a threshold**.

Evaluation

- Methodology

- Implementing our system on real smartphones
- Using two loudspeakers to perform replay attack

Maker	Model	Number of trumpets
Willnorn	SoundPlus	2
Amazon	Echo	2



- Performance metrics

- The standard automatic speaker verification metrics
- True Acceptance Rate (TAR)
- True Rejection Rate (TRR)

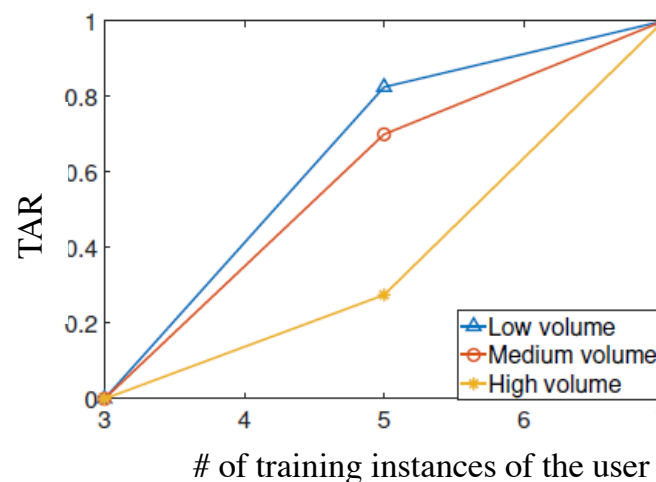
Evaluation

- Influence of locations on random noise-based approach

Locations	TAR	TRR
1	100%	100%
2	100%	100%
3	100%	100%
4	97.5%	100%

- Influence of acoustic noise on spectrum-based approach

7 training instances from the user are sufficient



Evaluation

- Overall performance
 - Simple replay attack

Solutions	TAR	TRR	Computation cost
Voice-based	100%	100%	Medium (SVM+STFT)
Motion-based	93.3%	88.93%	Low (SVM)

- Strong replay attack

Solutions	TAR	TRR	Computation cost
Voice-based & random noise	97.5%	100%	High (SVM+2*STFT)
Motion-based & random noise	91.0%	100%	Medium (SVM+STFT)

Conclusion

- Smartphone-based liveness detection system
 - Leveraging microphones and motion sensors in smartphone - **without additional hardware**
 - Easy to integrate with off-the-shelf mobile phones - **software-based approach**
- Good performance against strong attackers

