

Push the Limit of Adversarial Example Attack on Speaker Recognition in Physical Domain

Presenter: Qianniu Chen

Qianniu Chen¹, Meng Chen¹, Li Lu¹, Jiadi Yu²,
Yingying Chen³, Zhibo Wang¹, Zhongjie Ba¹, Feng Lin¹, Kui Ren¹

¹ Zhejiang University, Zhejiang, China

² Shanghai Jiao Tong University, Shanghai, China

³ Rutgers University, New Brunswick, NJ, USA



Background

- Speaker Recognition (SR) achieves wide applications in our daily life

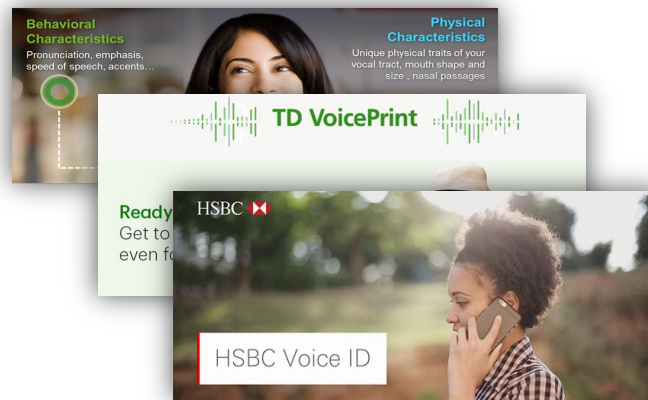


Background

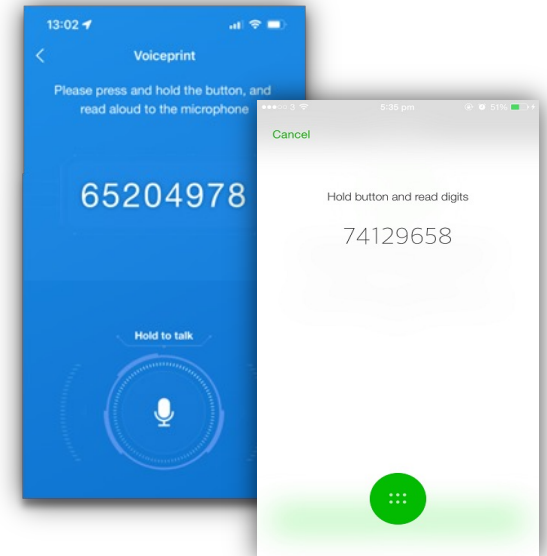
- Speaker Recognition (SR) achieves wide applications in our daily life



Voice Assistant



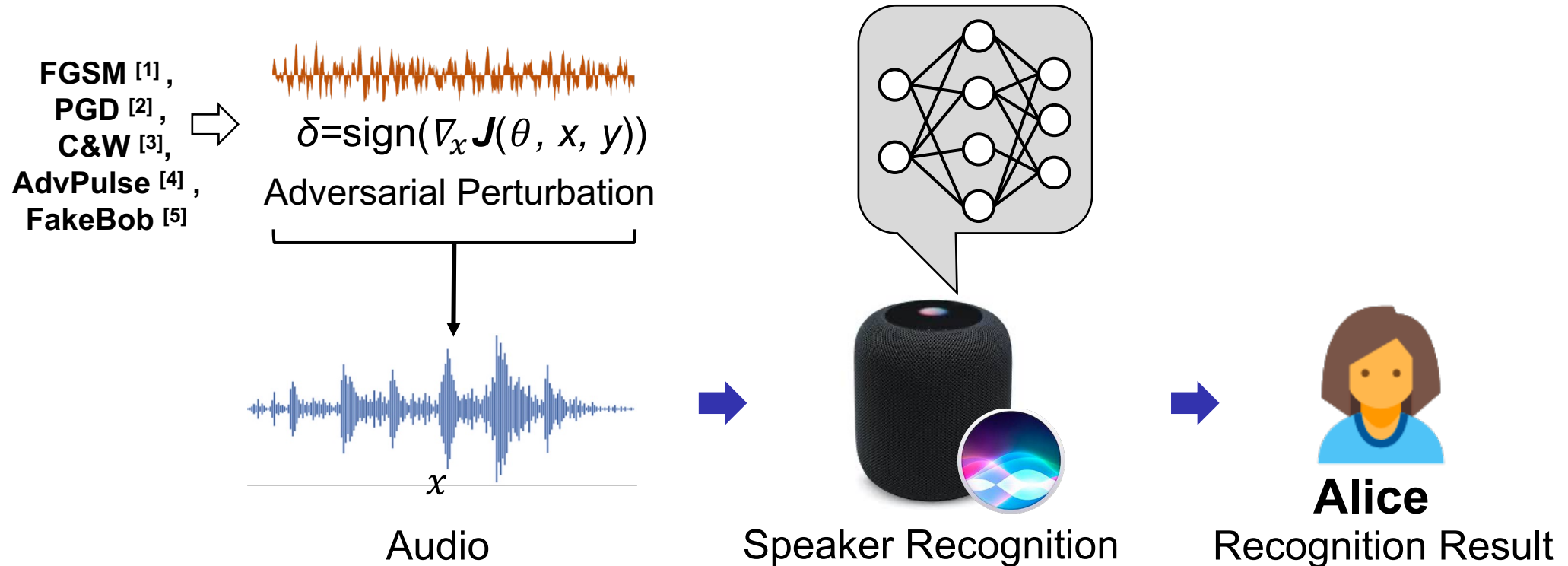
Mobile Bank Voice Password



APP Voice Lock

Background

● Speaker Recognition System & Audio Adversarial Example



[1] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples,” in Proceedings of ICLR. 2015.

[2] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. in Proceedings of ICLR. 2018.

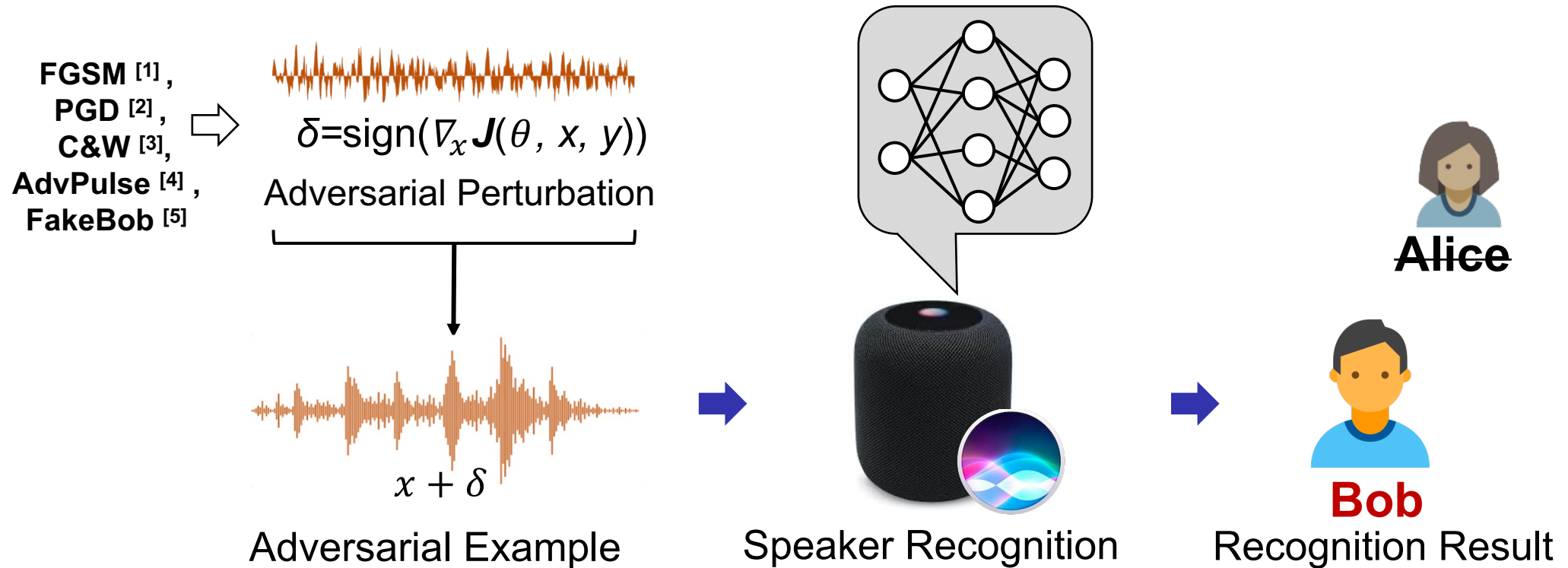
[3] N. Carlini and D. A. Wagner. Towards evaluating the robustness of neural networks. in Proceedings of IEEE S&P. 2017.

[4] Z. Li, Y. Wu, J. Liu, et al. Advpulse: Universal, synchronization-free, and targeted audio adversarial attacks via subsecond perturbations. in Proceedings of ACM CCS. 2020.

[5] G. Chen, S. Chen, L. Fan, et al. Who is real bob? adversarial attacks on speaker recognition systems. in Proceedings of IEEE S&P. 2021.

Background

● Speaker Recognition System & Audio Adversarial Example



[1] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples,” in Proceedings of ICLR. 2015.

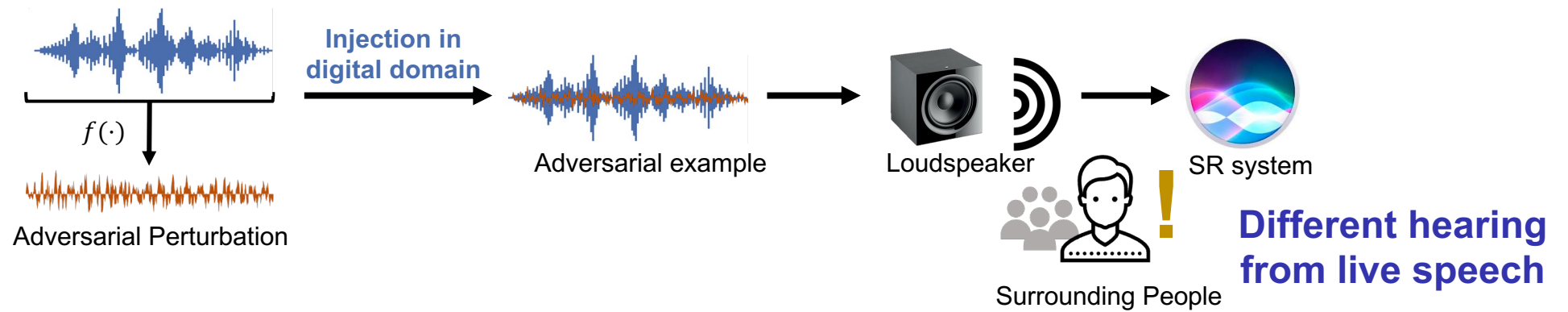
[2] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. in Proceedings of ICLR. 2018.

[3] N. Carlini and D. A. Wagner. Towards evaluating the robustness of neural networks. in Proceedings of IEEE S&P. 2017.

[4] Z. Li, Y. Wu, J. Liu, et al. Advpulse: Universal, synchronization-free, and targeted audio adversarial attacks via subsecond perturbations. in Proceedings of ACM CCS. 2020.

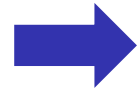
[5] G. Chen, S. Chen, L. Fan, et al. Who is real bob? adversarial attacks on speaker recognition systems. in Proceedings of IEEE S&P. 2021.

Threat Model



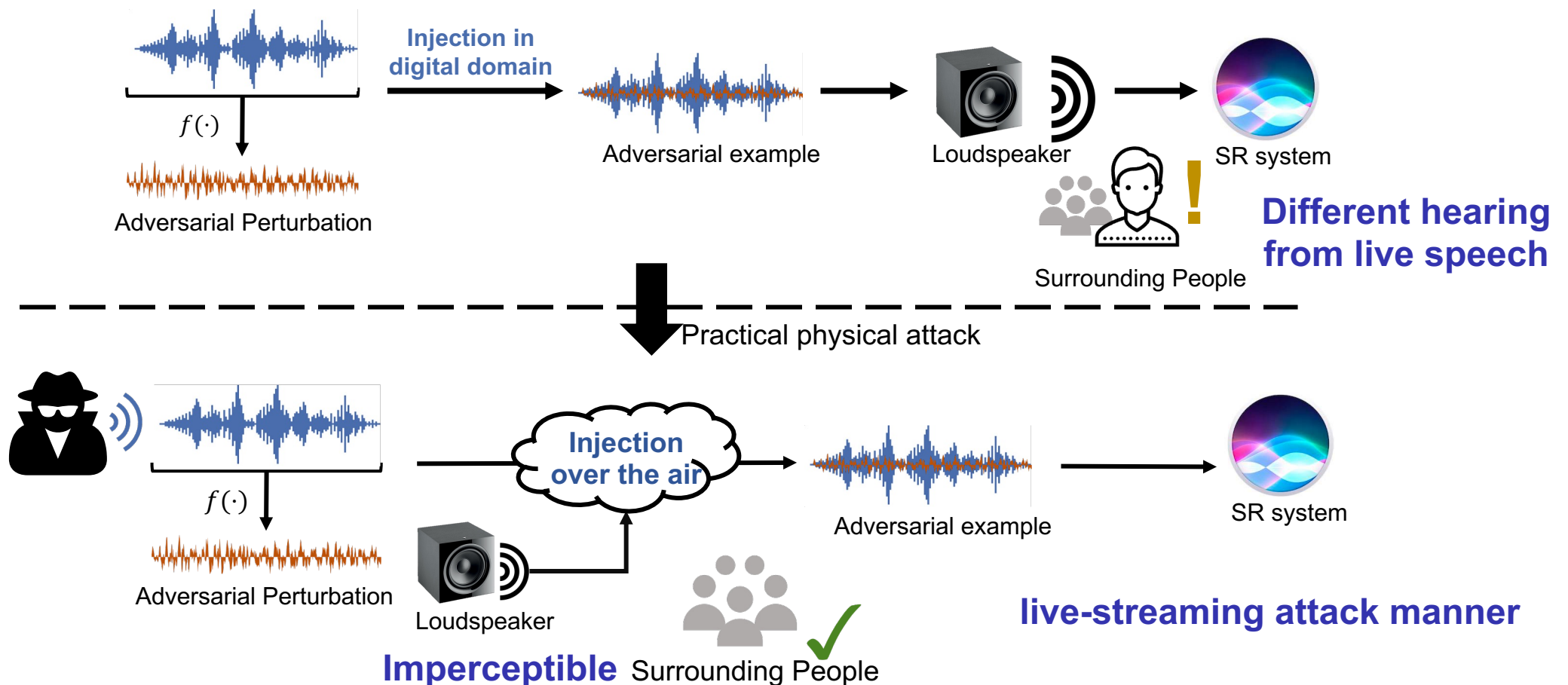
- ➔ Replay **VOICE AUDIO** with adversarial perturbation
- Limited to an ideal attack scenario without others around**

Threat Model



Replay **VOICE AUDIO** with adversarial perturbation

Limited to an ideal attack scenario without others around

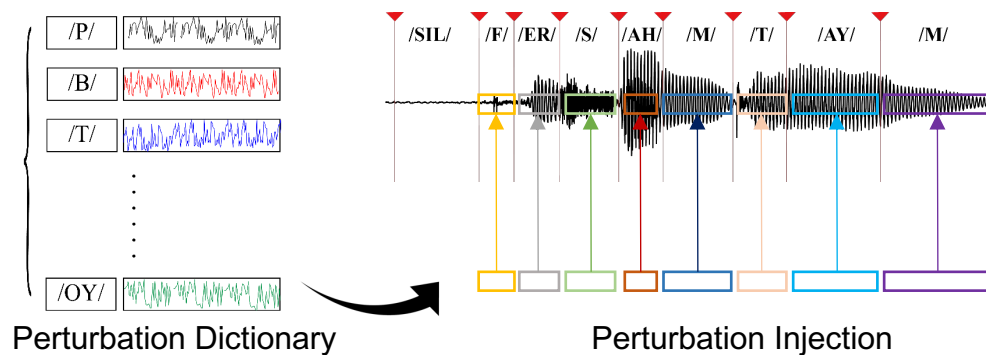


Attack Overview

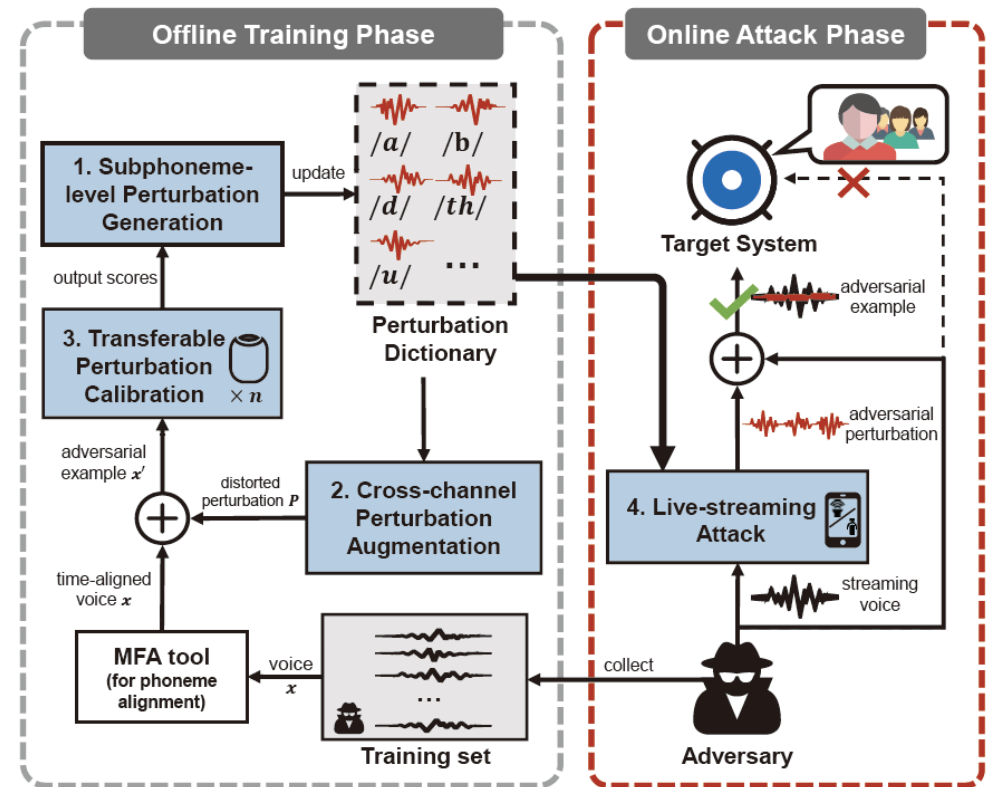
PhyTalker: a live-streaming, channel-robust, transferable audio adversarial example attack

Universal Adversarial Perturbation on Phonemes

- **Combinability:** fast generate perturbation for any speech according to its phoneme sequence.
- **Stability:** cause stable attack effectiveness on the stable acoustic characteristics of phonemes.



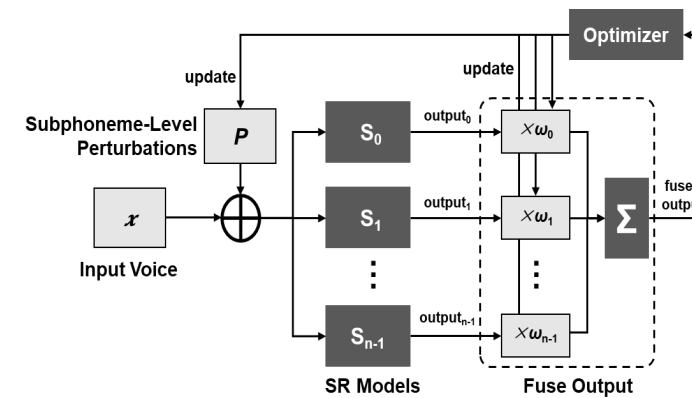
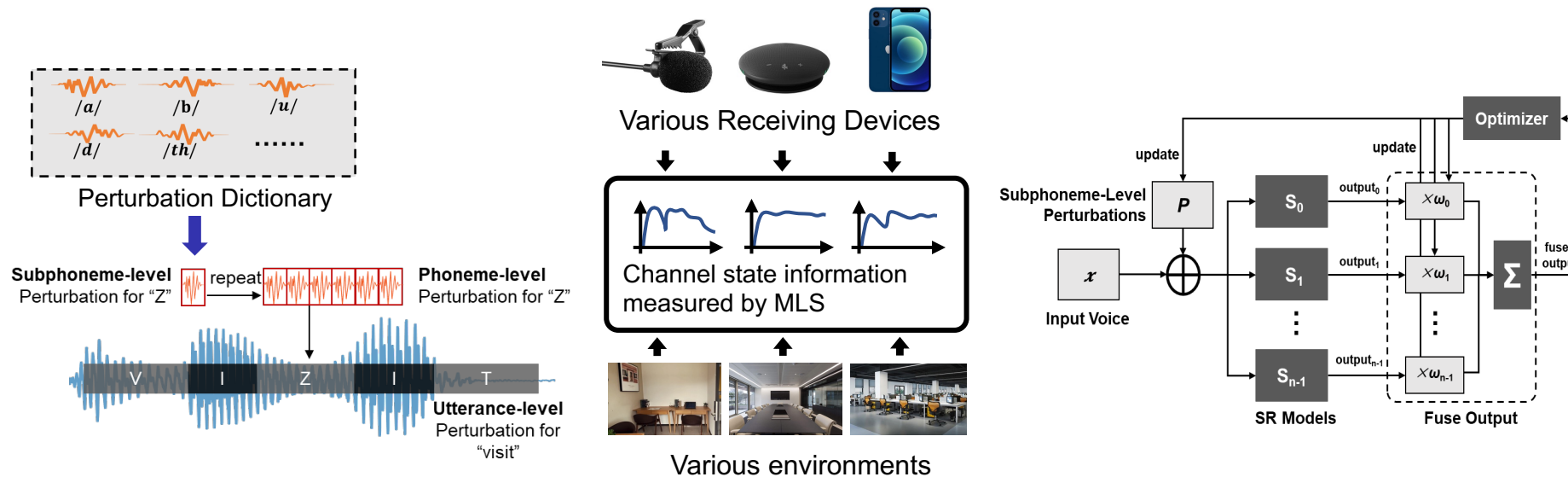
Function of perturbation on phonemes



Attack overview of PhyTalker

Offline Training Phase

- **Subphoneme-level Perturbation:** use fixed short perturbation (<25ms) to form phoneme-level perturbation with various duration(>50ms), repetitively
- **Channel Augmentation:** explore real channel state information for data augmentation with MLS[6]
- **Transferable Calibration:** employ the ensemble learning method to improve transferability
- **Expectation Optimization :** train on a large training set instead of the specific audio



Algorithm 1 Global Optimization Procedure.

Input: Dataset $\chi = \{(x_0, t_0), (x_1, t_1), \dots, (x_{n-1}, t_{n-1})\}$, SR classifiers S_i with threshold θ_i ($i = 0, 1, \dots, k-1$), UIR set C target user y_t , amplitude upper bound ϵ .

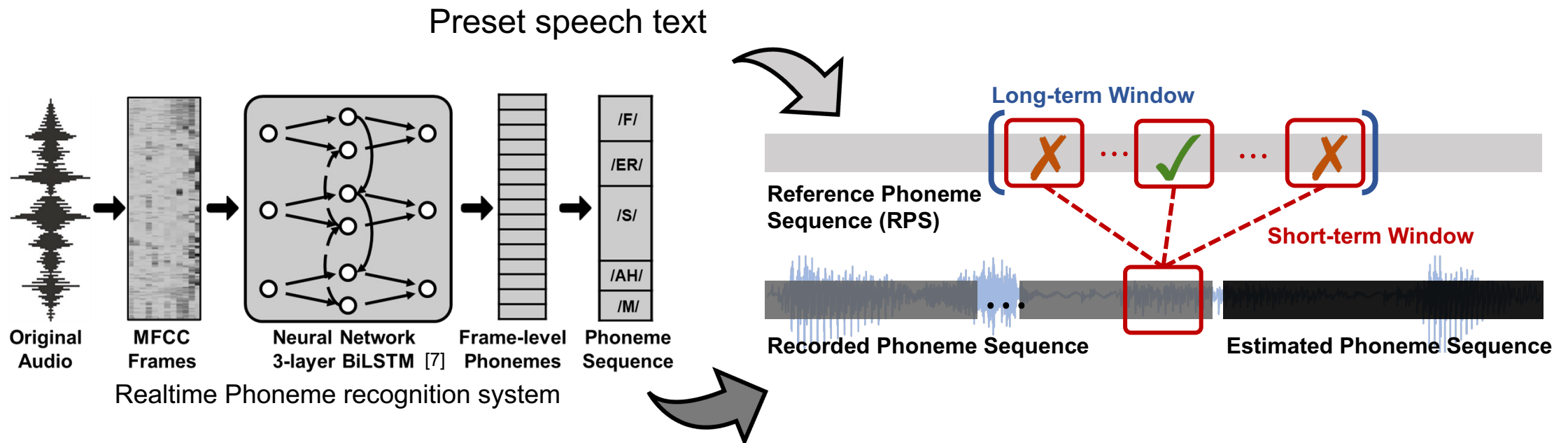
Output:

- 1: Initialize P with normal distribution
- 2: **repeat**
- 3: **for each** (X, T) sampled from χ **do**
- 4: $L_t \leftarrow 0$
- 5: **for each** c in C **do**
- 6: $X_{adv} \leftarrow G(X, P * c)$
- 7: $s \leftarrow \sum_0^{k-1} w_i \cdot L_{S_i, \theta_i}(X_{adv}, y_t)$
- 8: $L \leftarrow \max\{\theta, \max_{i \neq y_t} s_i\} - s_{y_t}$
- 9: $L_t \leftarrow L_t + L$
- 10: **end for**
- 11: Minimize L_t to update P
- 12: **end for**
- 13: **until** Early Stopping
- 14: **return** P ;

[6] Douglas D. Rife and John Vanderkooy. Transfer-function measurement with maximum-length sequences. Journal of the Audio Engineering Society 37, 6 (June). 1989.

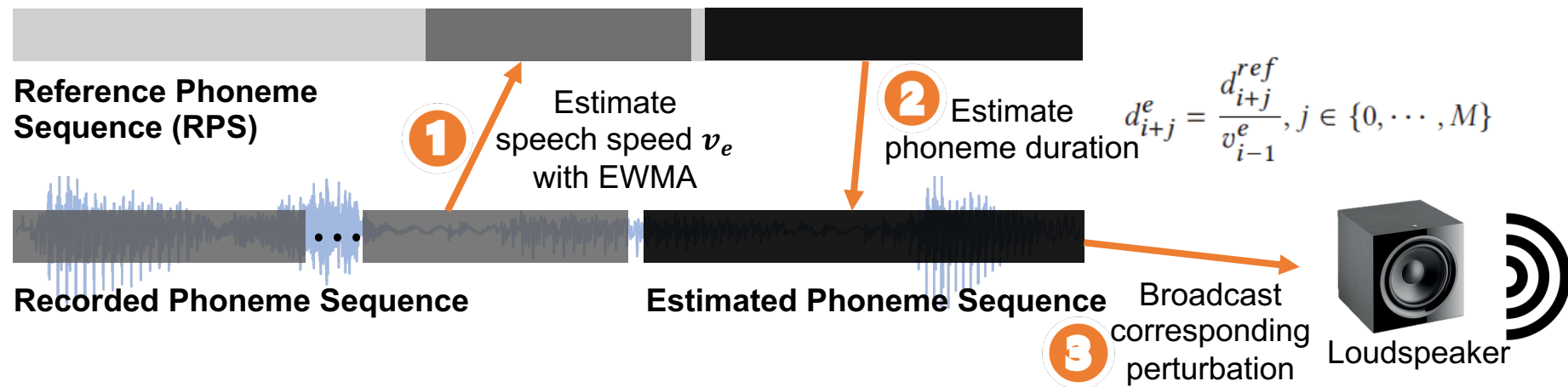
Online Attack Phase

- **Real-time phoneme extraction:** extract current phoneme sequence from the live speech with a fast neural phoneme recognition system
- **Phoneme Alignment :** locate the current phoneme in the RPS with long-short term



Online Attack Phase

- **Real-time phoneme extraction:** extract current phoneme sequence from the live speech with a fast neural phoneme recognition system
- **Phoneme Alignment :** locate the current phoneme in the RPS with long-short term
- **Phoneme Estimation :** Estimate speech voice and patch phoneme durations by referring RPS



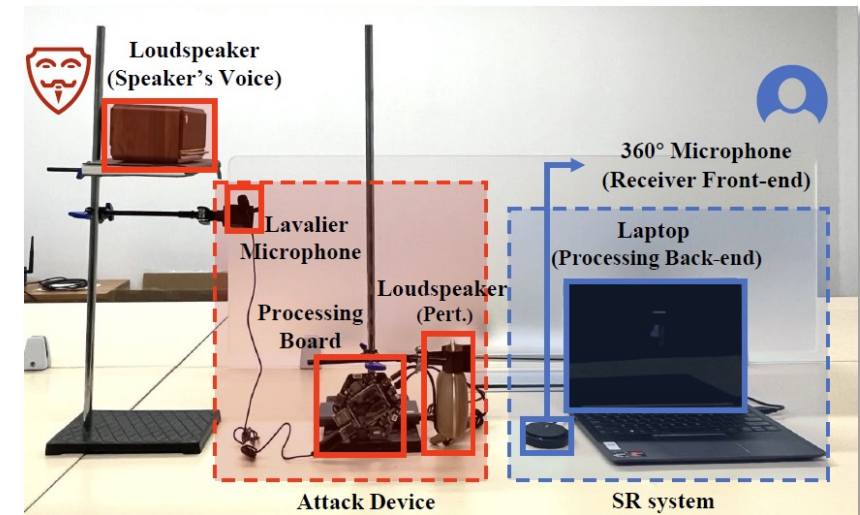
Evaluation

Target Systems Setting

- Architectures: x-vector[8] / d-vector[9] / DeepSpeaker[10]
- Training Set: Voxceleb [11] corpus
- Test Set: LibriSpeech[12] corpus
- Enrollers: 5 speakers(3 males and 2 females)
- Backend: Lenovo Xiaoxin Pro 13

Attack Setting

- Adversaries: 10 speakers(5 males and 5 females)
- Attack Device: ReSpeakerCore v2
- Subphoneme-level perturbation duration: 12.5ms
- Livestreaming synchronization: 0.5s/alignment
- Channel augmentation: 8 CIRs per (receiver, environment)
- Ensemble learning: 4 ensemble models



[8] E. Variani, X. Lei, E. McDermott, I. L. Moreno, and J. Gonzalez-Dominguez. Deep neural networks for small footprint text-dependent speaker verification. in Proceedings of IEEE ICASSP. 2014.

[9] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur. Xvectors: Robust dnn embeddings for speaker recognition. in Proceedings of IEEE ICASSP. 2018.

[10] C. Li, X. Ma, B. Jiang, X. Li, X. Zhang, X. Liu, Y. Cao, A. Kannan, and Z. Zhu. Deep speaker: an end-to-end neural speaker embedding system. CoRR, vol. abs/1705.02304. 2017.

[11] Arsha Nagrani, Joon Son Chung, and Andrew Senior. VoxCeleb: A Large-Scale Speaker Identification Dataset. In Proceedings of ISCA Interspeech. 2017.

[12] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An ASR corpus based on public domain audio books. In Proceedings of IEEE ICASSP. 2015.

Evaluation

- Overall performance
- Human Imperceptibility

Table 3: Overall ASRs, SNR, MCD and RTF of *PhyTalker* and SOTA works under physical attack scenarios.

Attack	ASR(%)			SNR (dB)	MCD (dB)	RTF
	d-vec.	x-vec.	D.S.			
<i>PhyTalker</i>	85.5	80.5	90.5	16.8	2.45	0.5
<i>FakeBob</i> [5]	63.3	77.4	69.8	11.6	4.15	95.3
<i>AdvPulse</i> [4]	N/A	89.9	N/A	4.7	N/A	<1.0

ASR: Attack Success Rate for effectiveness (the higher the better)

MCD: Mel Cepstral Distortion for audibility (the lower the better)

SNR: Signal-to-Noise Ratio for audibility (the higher the better)

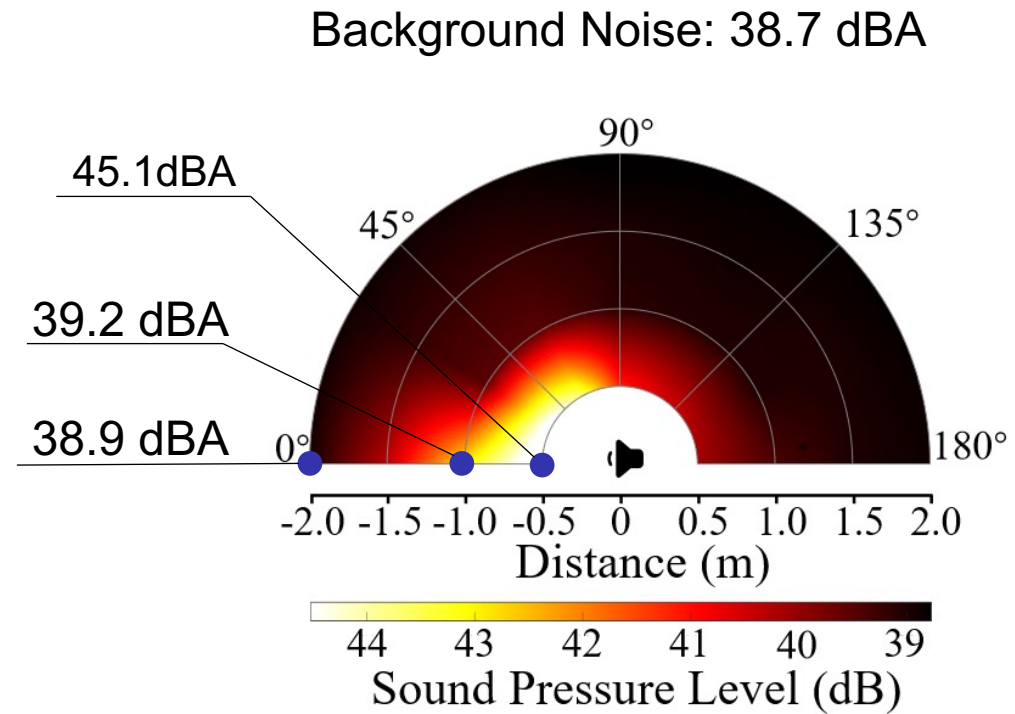
RTF: Real Time Factor for efficiency (the lower the better)

[4] Z. Li, Y. Wu, J. Liu, et al. Advpulse: Universal, synchronization-free, and targeted audio adversarial attacks via subsecond perturbations. in Proceedings of ACM CCS. 2020.

[5] G. Chen, S. Chen, L. Fan, et al. Who is real bob? adversarial attacks on speaker recognition systems. in Proceedings of IEEE S&P. 2021.

Evaluation

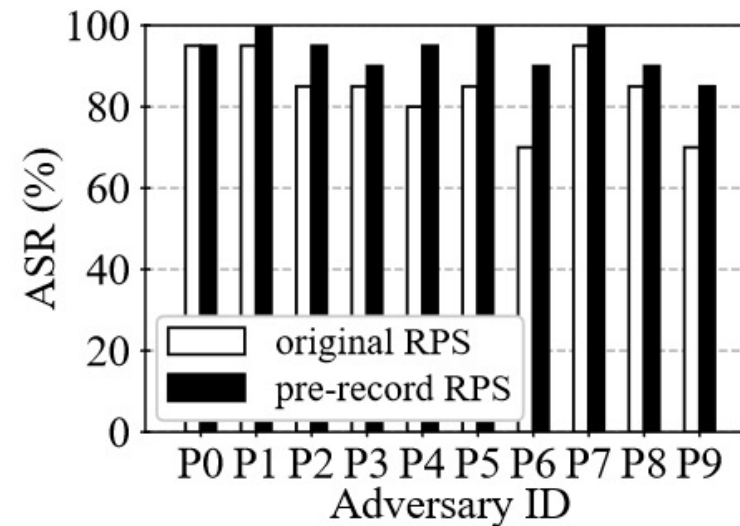
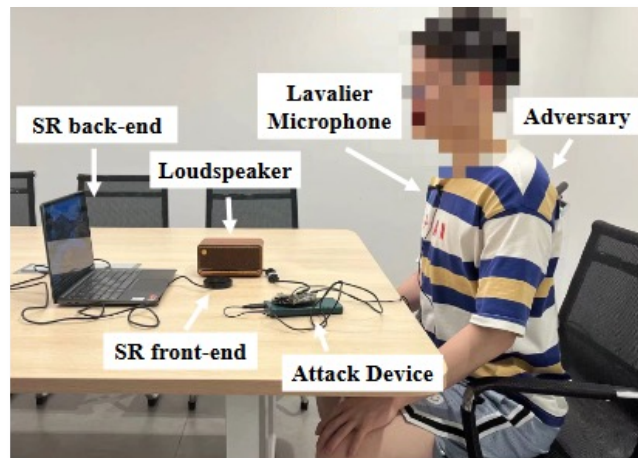
- Overall performance
- **Human Imperceptibility**



Evaluation: in-the-wild evaluation

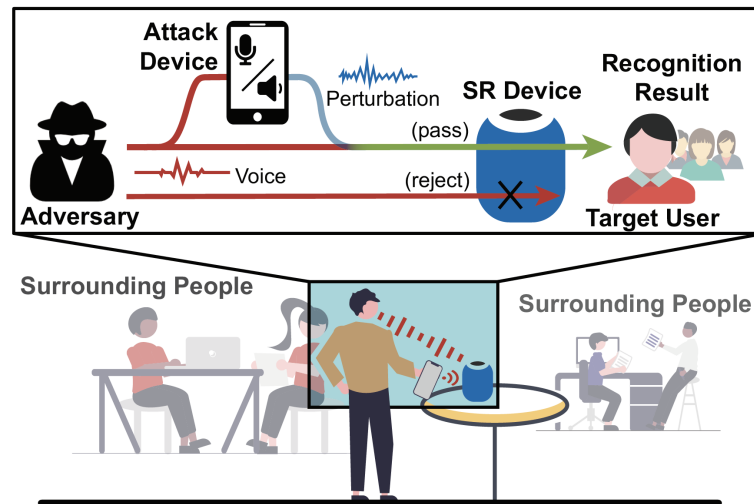
Evaluation Setup

- 10 volunteers as adversaries (4 females and 6 males)
- 20-minute voice record/volunteer for training
- 10 utterances per/volunteer for evaluation



Conclusion

- Explore **three major challenges** underlying a practical physical attack scenario
- Propose a subphoneme-level, channel-robust and transferable adversarial example attack to solve the challenges
- Enables an adversary to conduct a **live-streaming attack manner** in **physical domain**

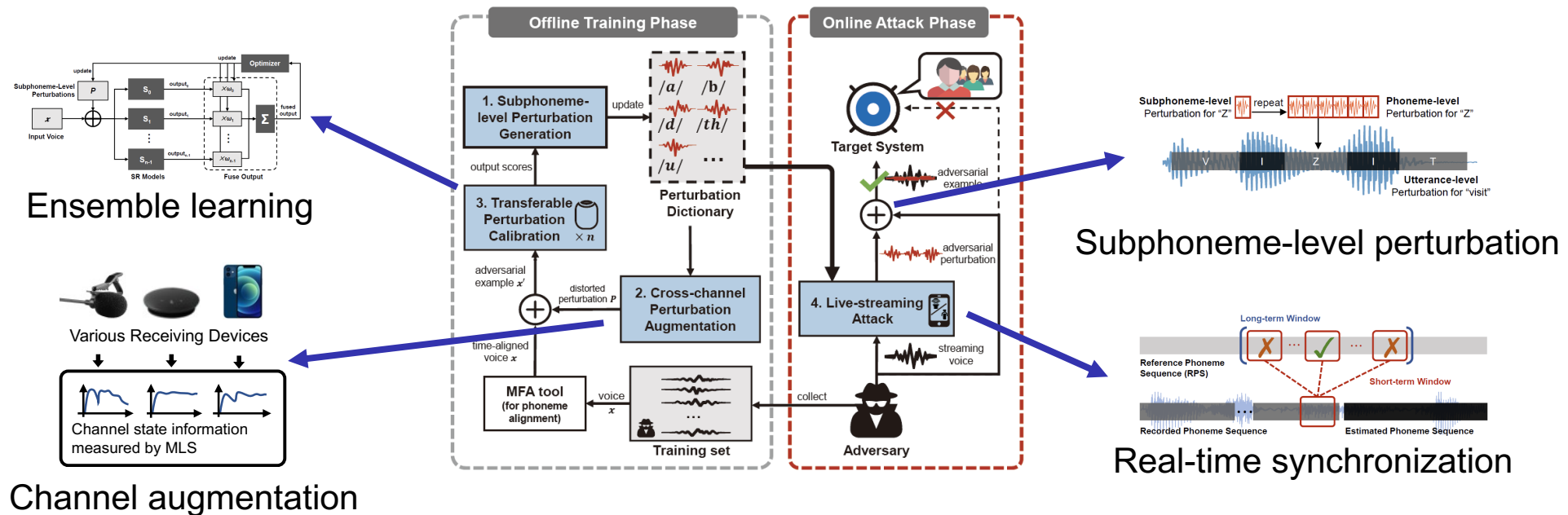


Major challenges

1. Livestreaming manner
2. Channel robustness
3. Black-box optimization

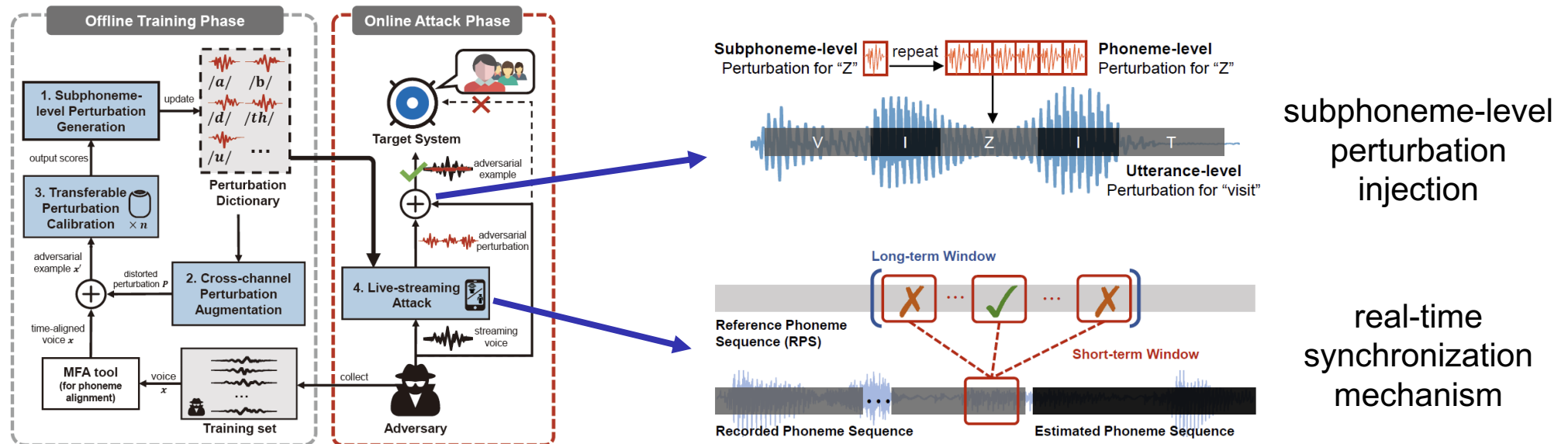
Conclusion

- Explore **three major challenges** underlying a practical physical attack scenario
- Propose a **subphoneme-level, channel-robust** and **transferable** adversarial example attack to solve the challenges
- Enables an adversary to conduct a **live-streaming attack manner** in physical domain

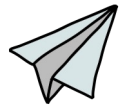


Conclusion

- Explore **three major challenges** underlying a practical physical attack scenario
- Propose a **subphoneme-level, channel-robust** and **transferable** adversarial example attack to solve the challenges
- Enables an adversary to conduct a **live-streaming attack manner** in **physical domain**



Thank you!



qianniuchen@zju.edu.cn

