

Metamorph: Injecting Inaudible Commands into Over-the-air Voice Controlled Systems

Tao Chen¹ Longfei Shangguan² Zhenjiang Li¹ Kyle Jamieson³

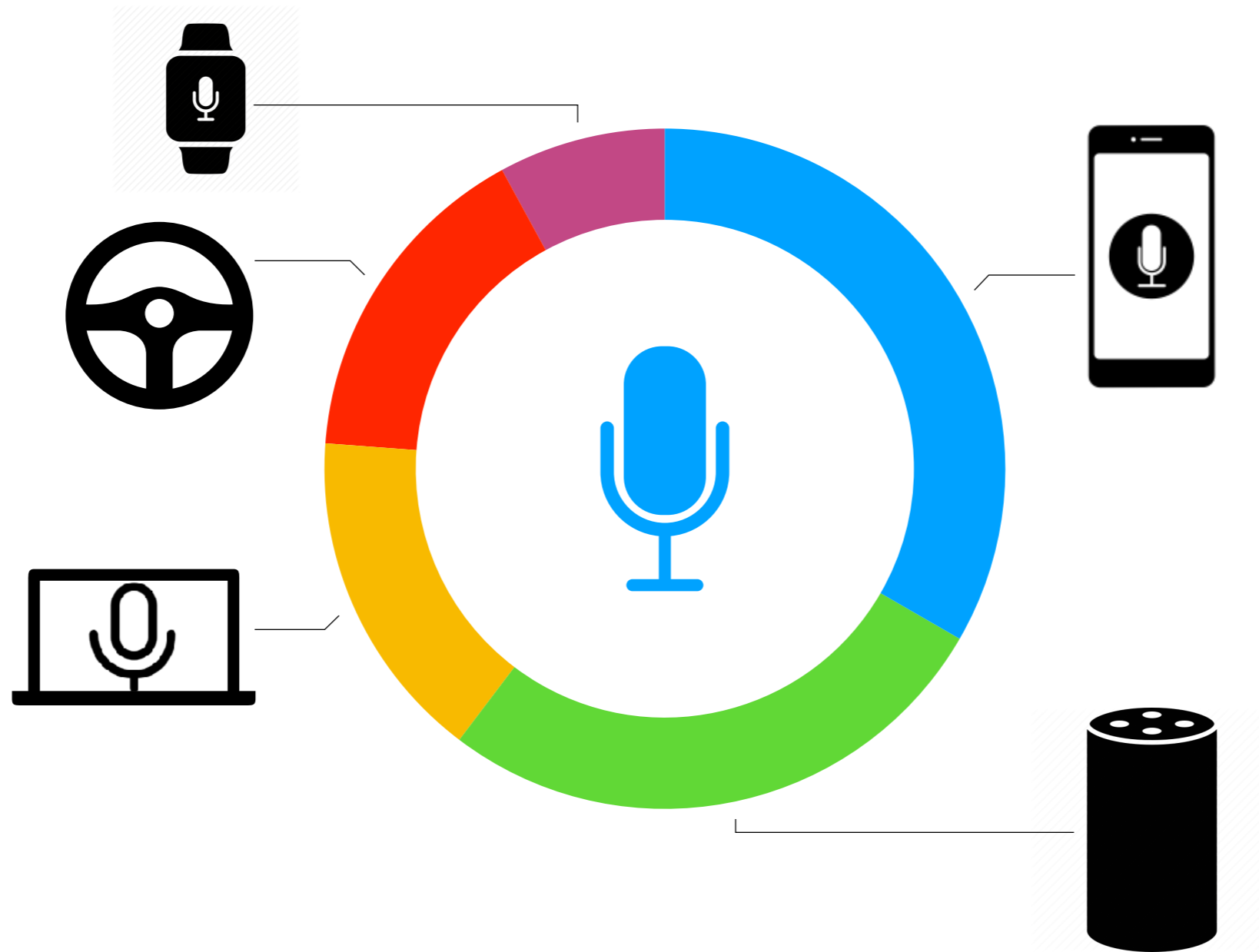
¹City University of Hong Kong, ²Microsoft, ³Princeton University



香港城市大學
City University of Hong Kong



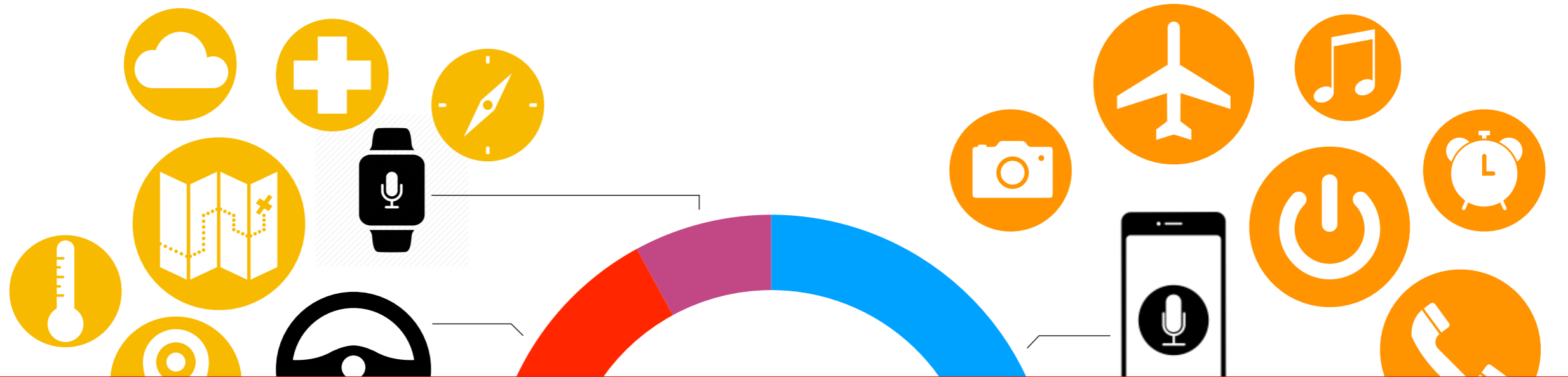
Voice Assistants in Smart Home



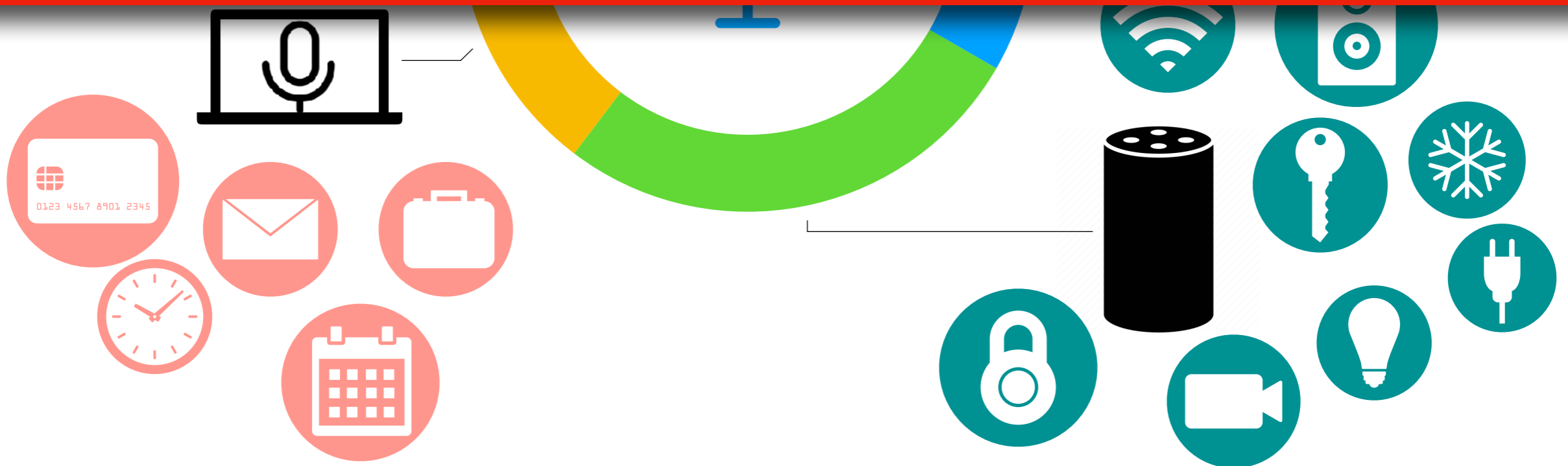
Voice Assistants in Smart Home



Voice Assistants in Smart Home



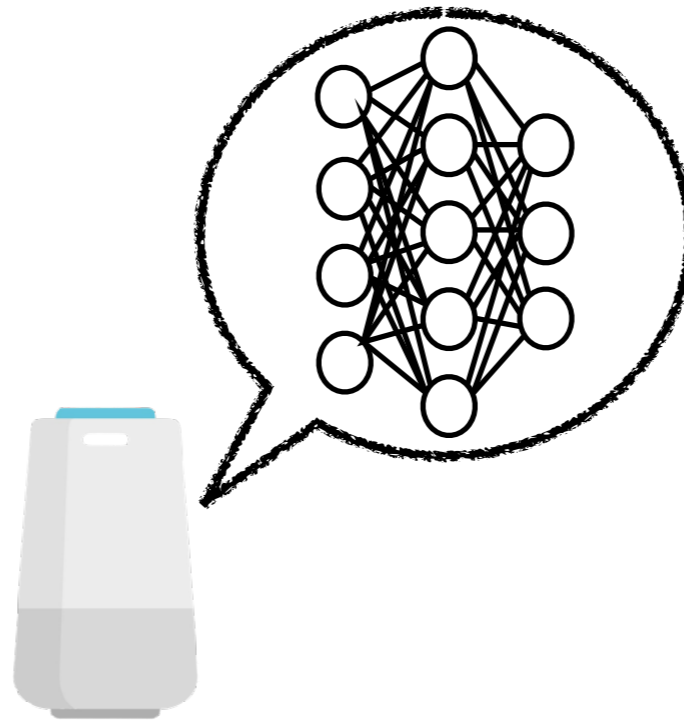
111.8 million people in U.S. use voice assistants and related services!



Are they safe enough?

How to attack the voice assistant?

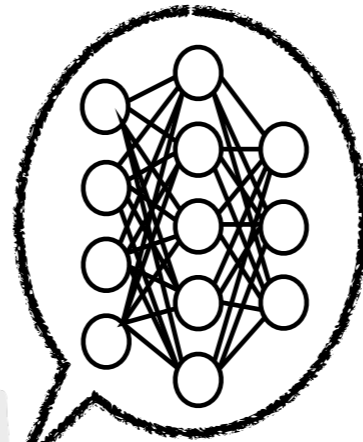
Neural networks



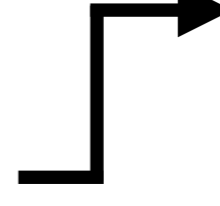
Speech Recognition Models (SR)

How to attack the voice assistant?

Audio Clip: I

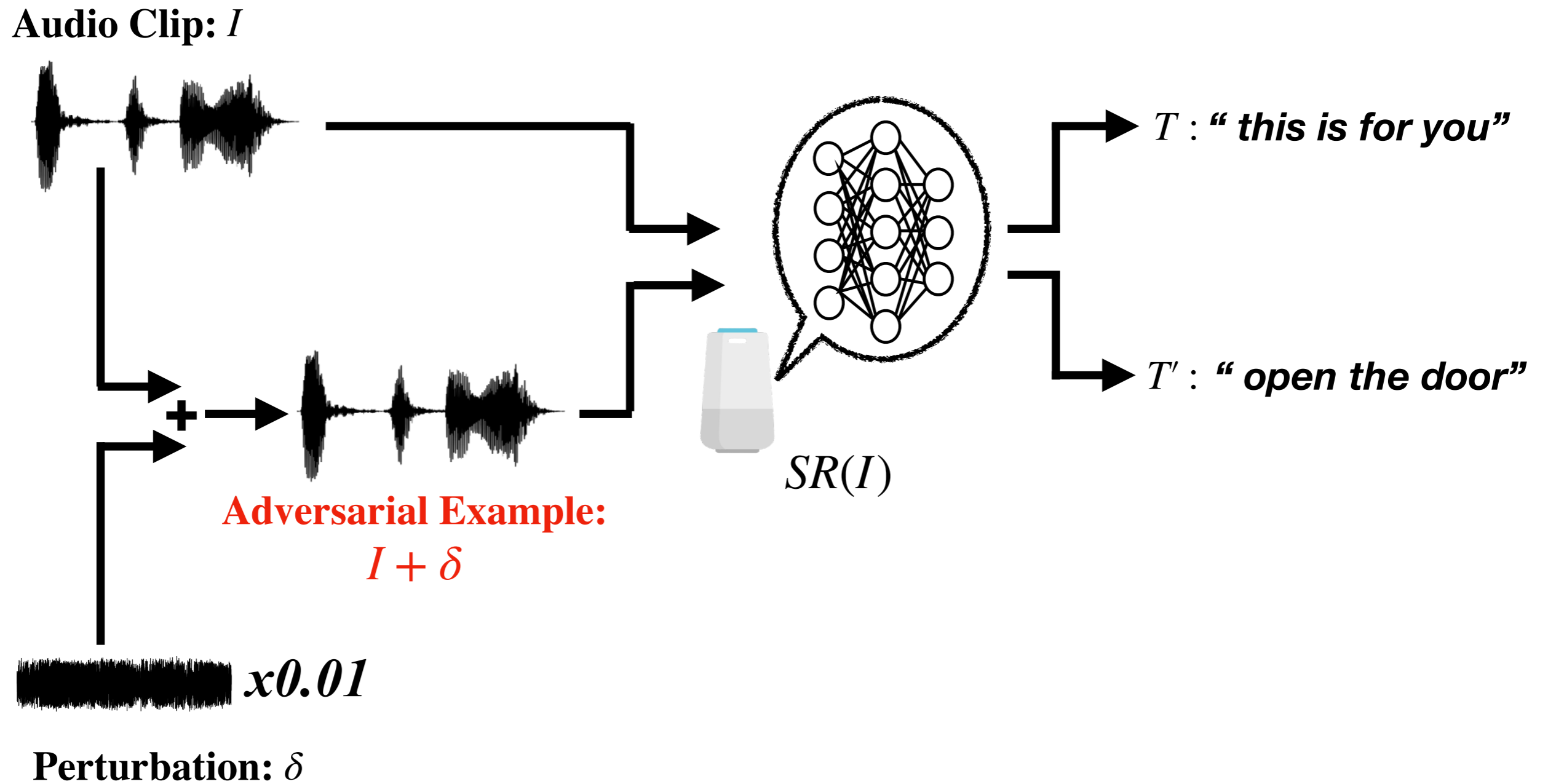


$SR(I)$

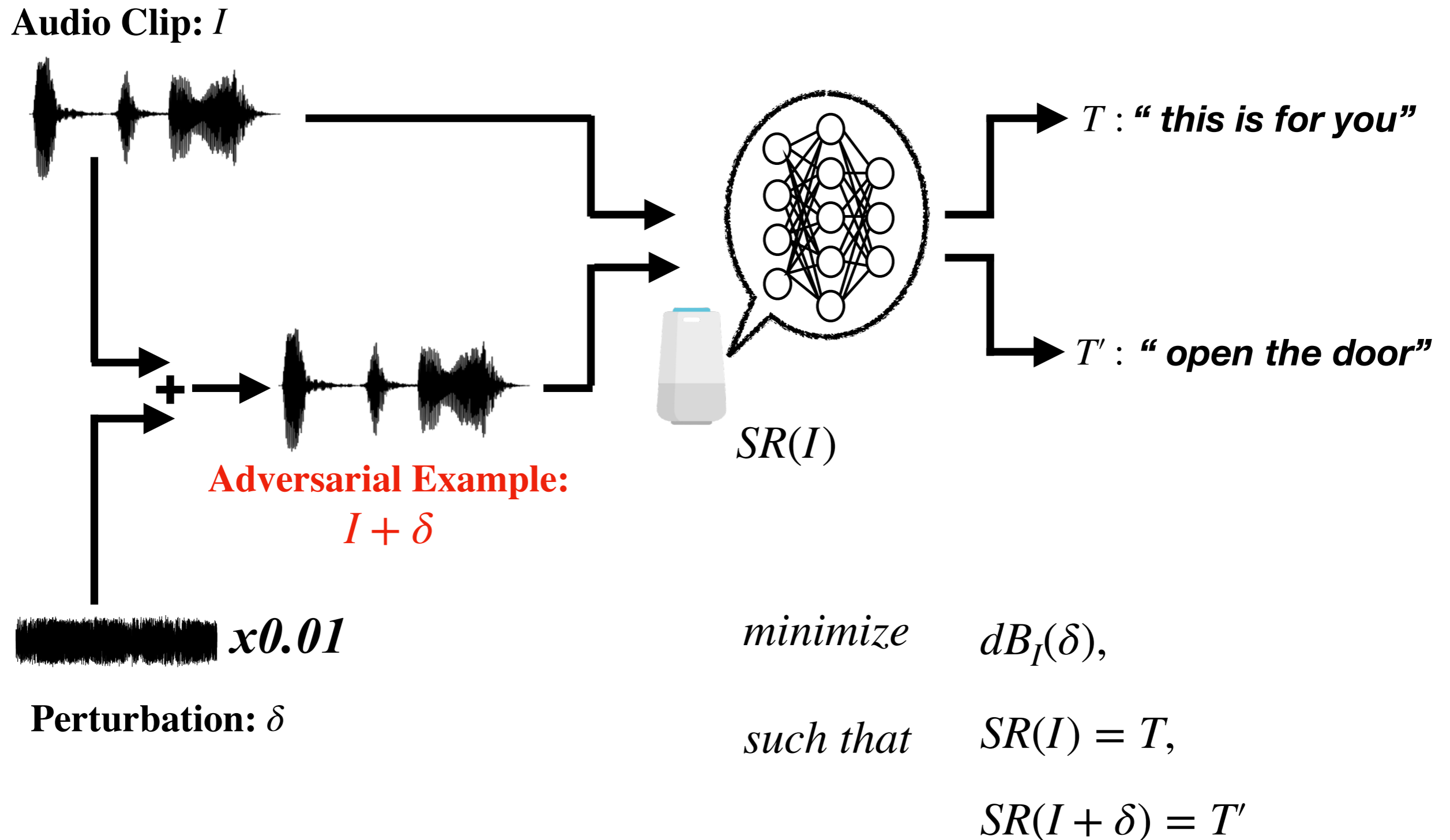


T : "this is for you"

How to attack the voice assistant?

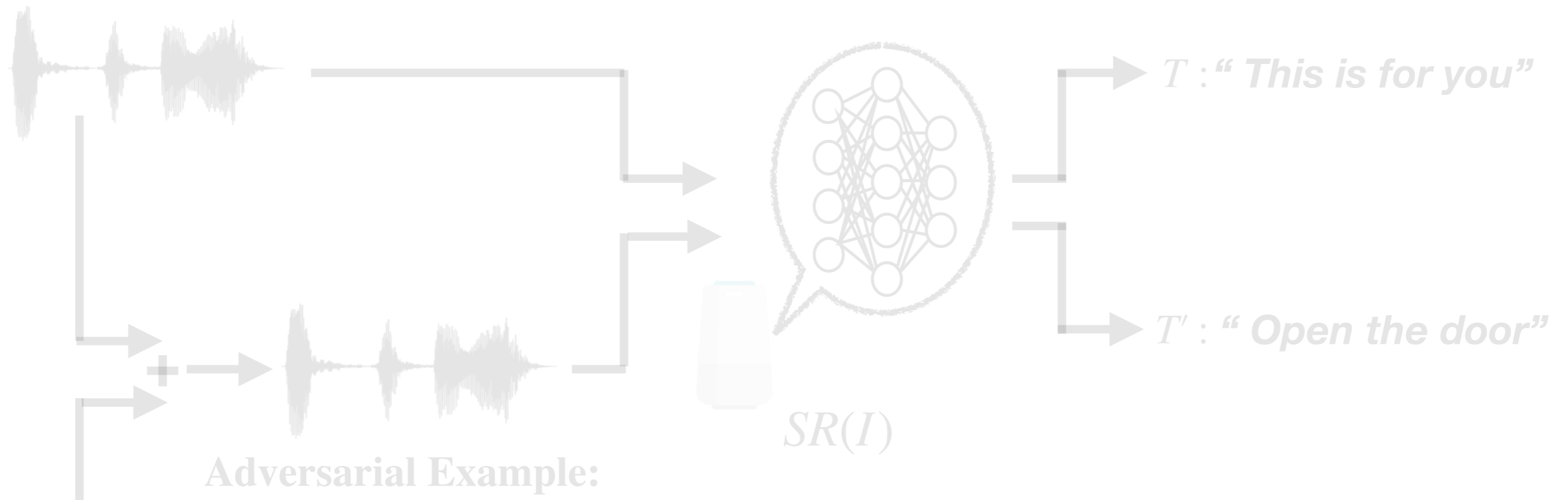


How to attack the voice assistant?



How to attack the voice assistant?

Audio Clip: I



Audio Adversarial Attack

 $\times 0.01$

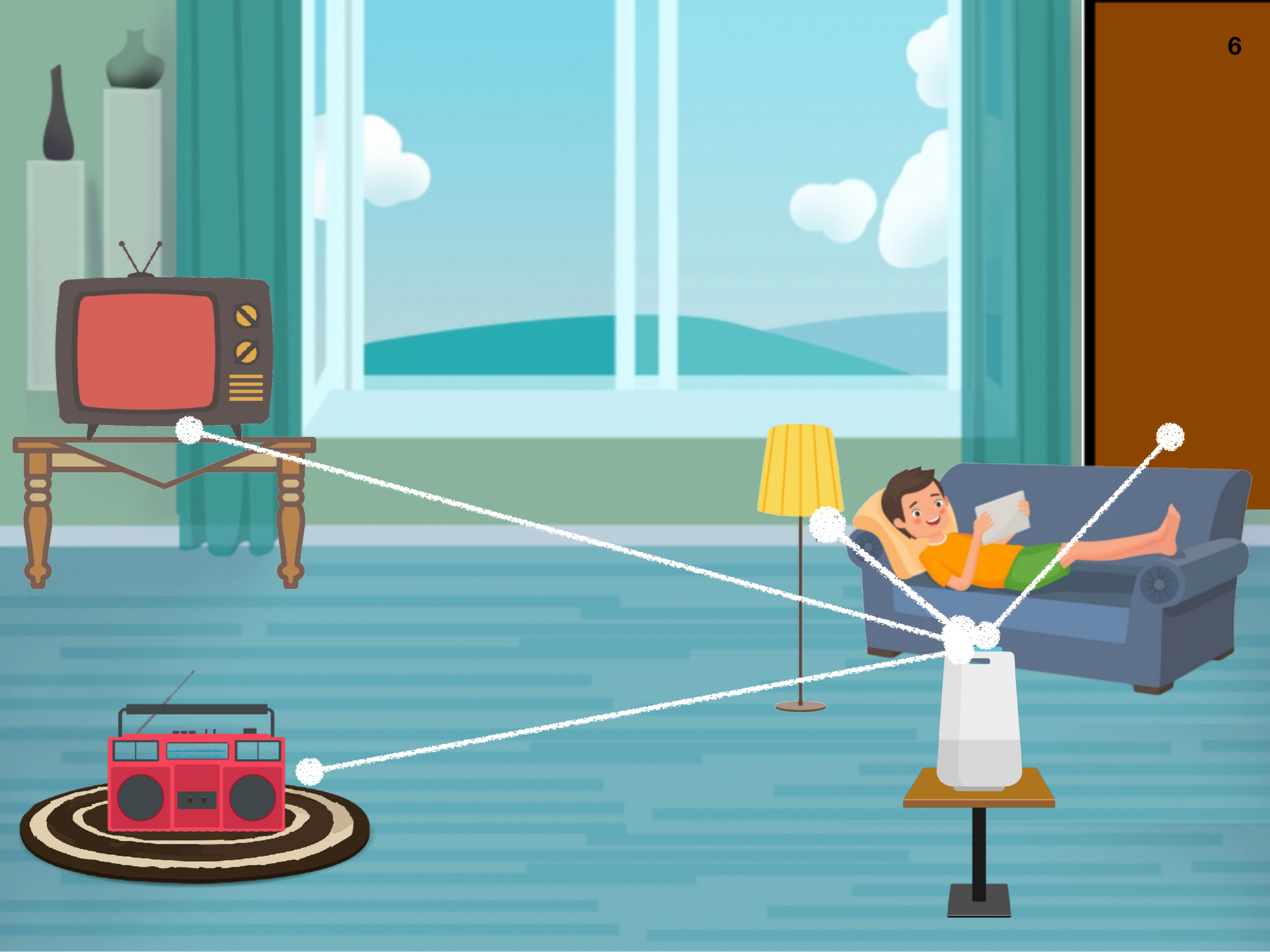
Perturbation: δ

minimize $dB_I(\delta)$,

such that $SR(I) = T$,

$SR(I + \delta) = T'$











Is it a real threat? Yes!



Adversarial Example



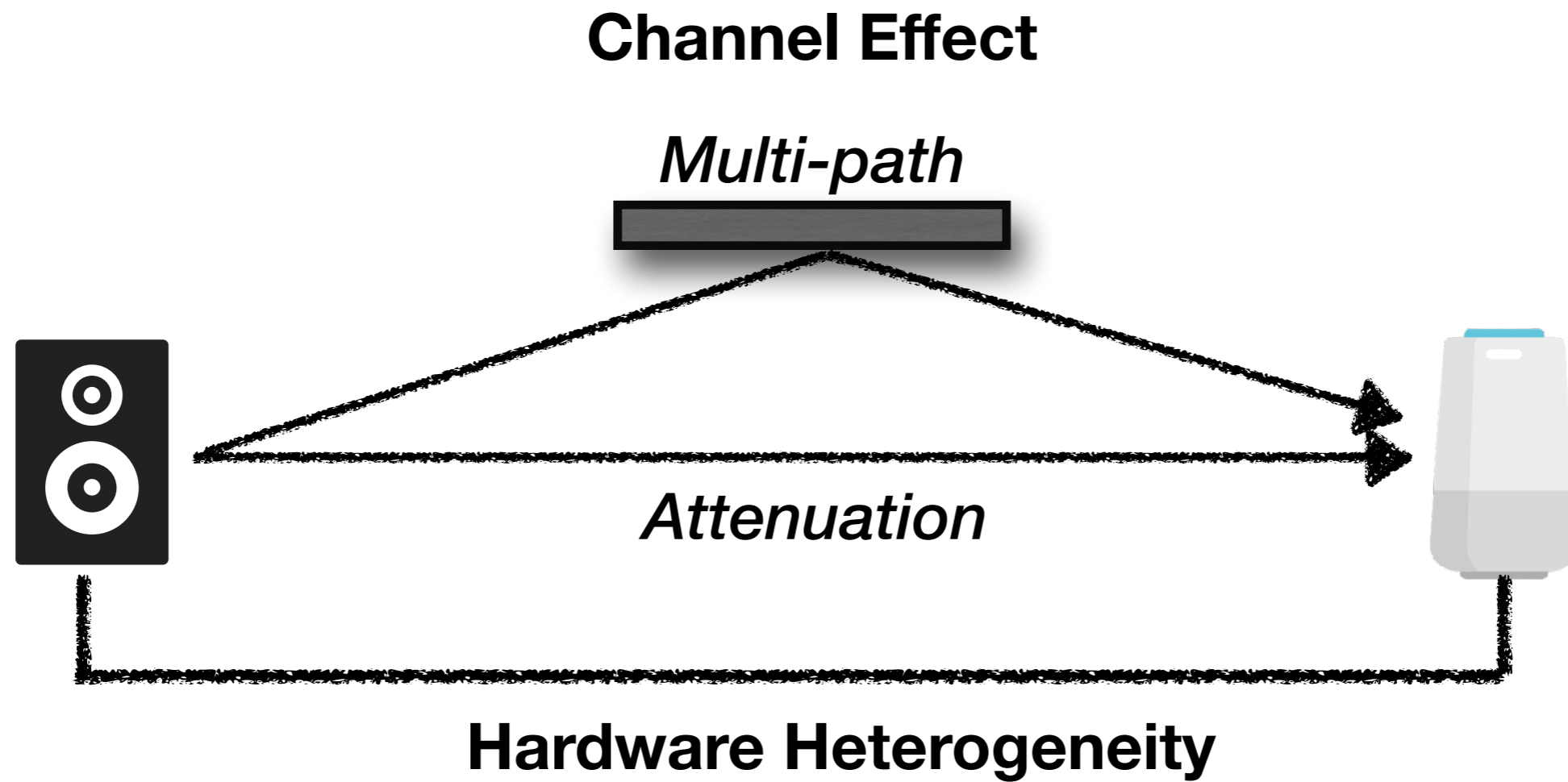


Adversarial Example

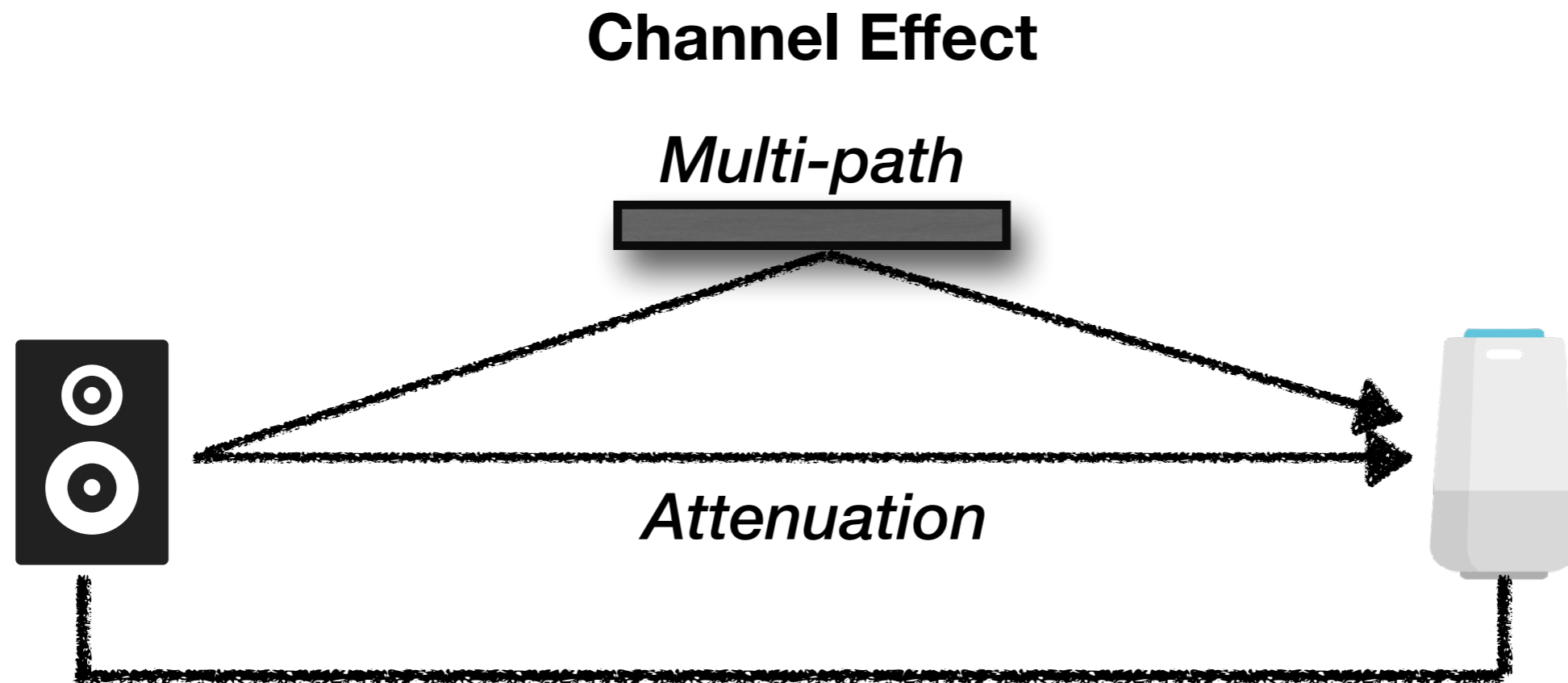
But, failed Over-the-air!



Challenge



Challenge

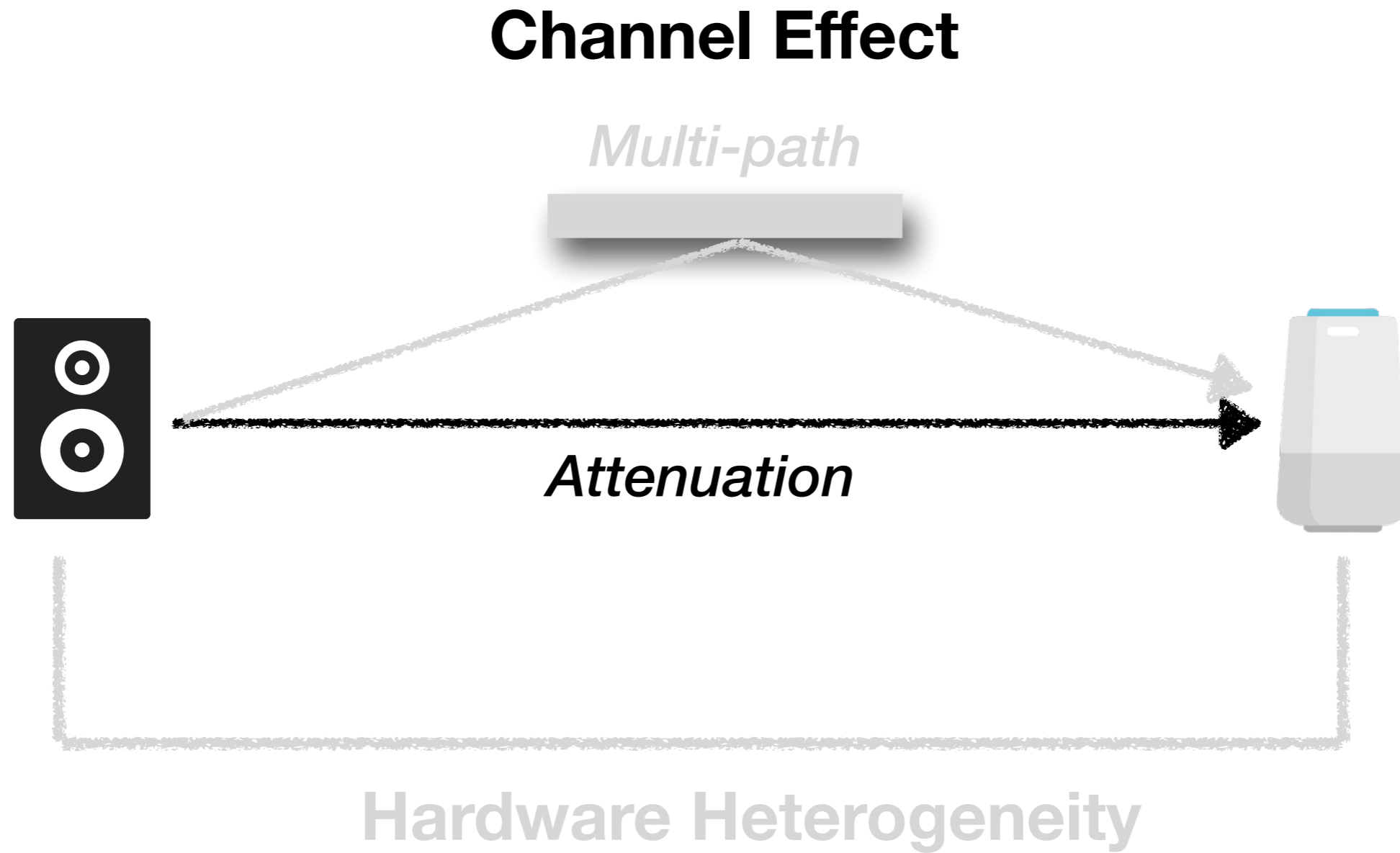


Hardware Heterogeneity

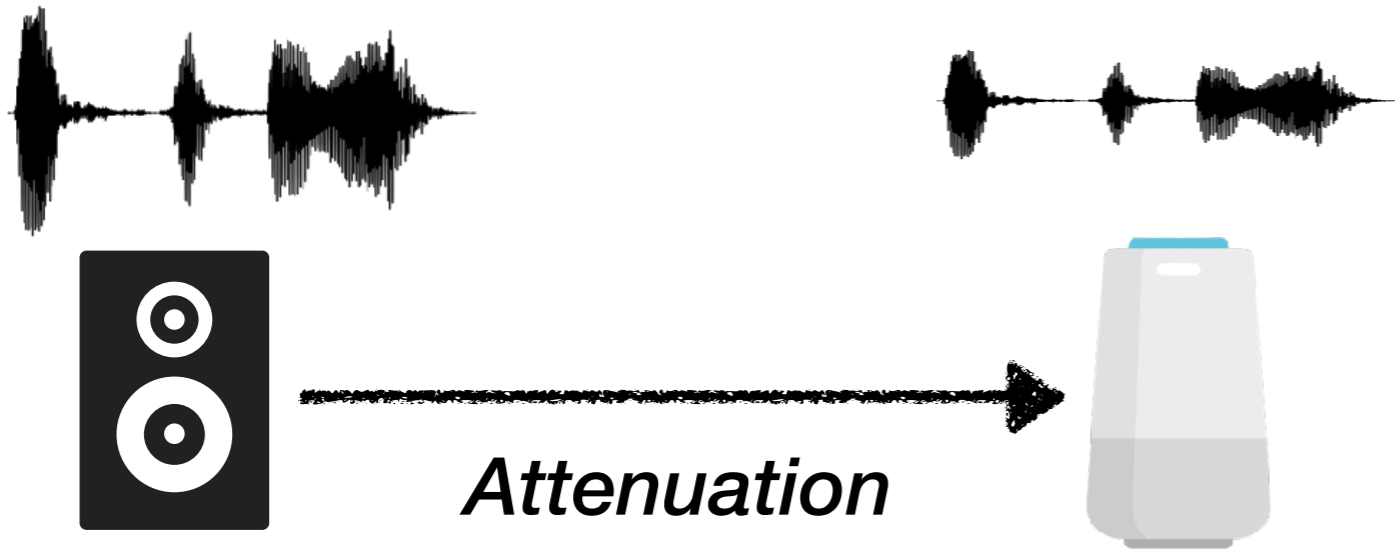
$$SR(I + \delta) \text{ VS } SR(H(I + \delta))$$

H is unknown in advance!

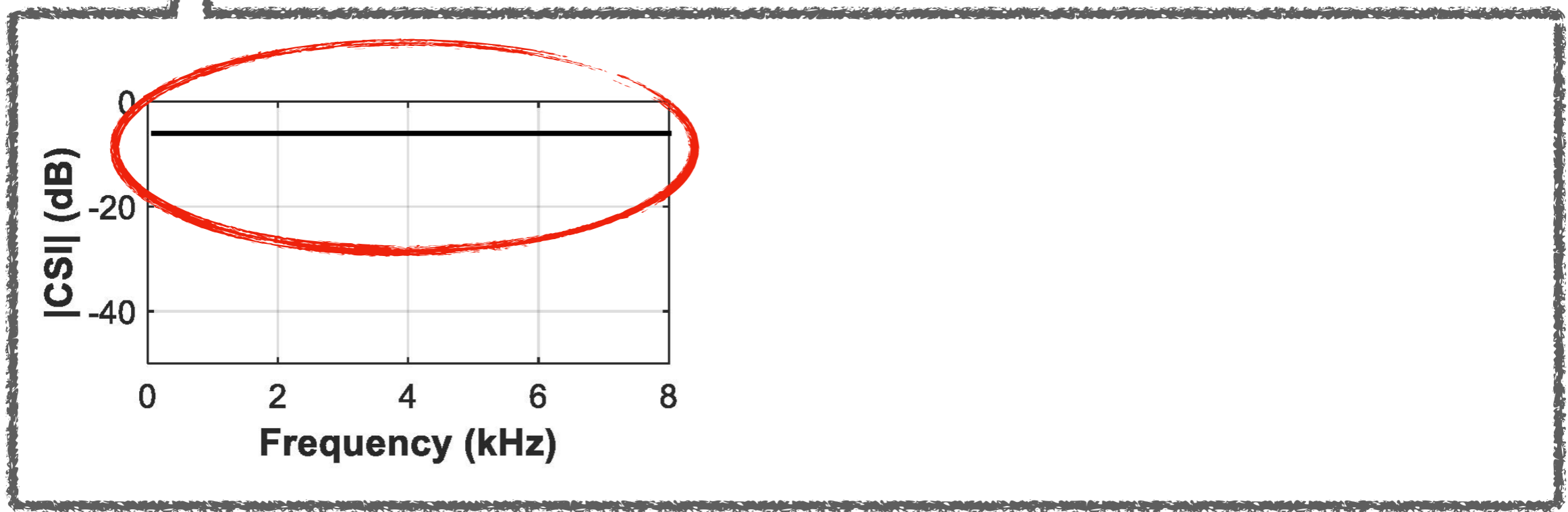
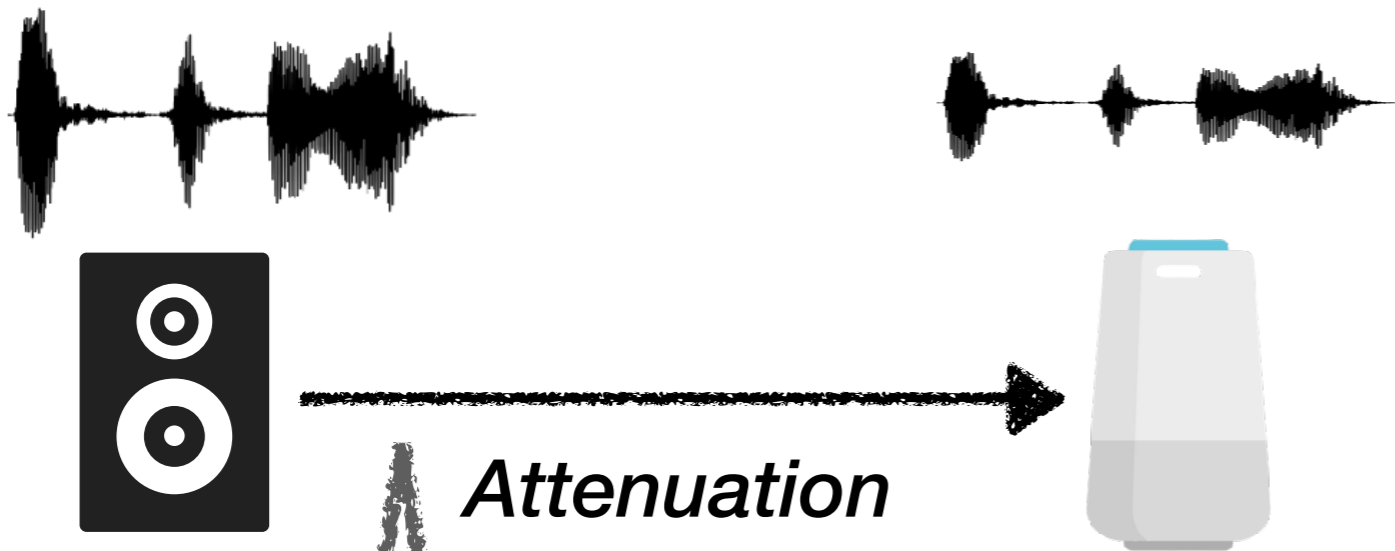
Understand Over-the-air Attack



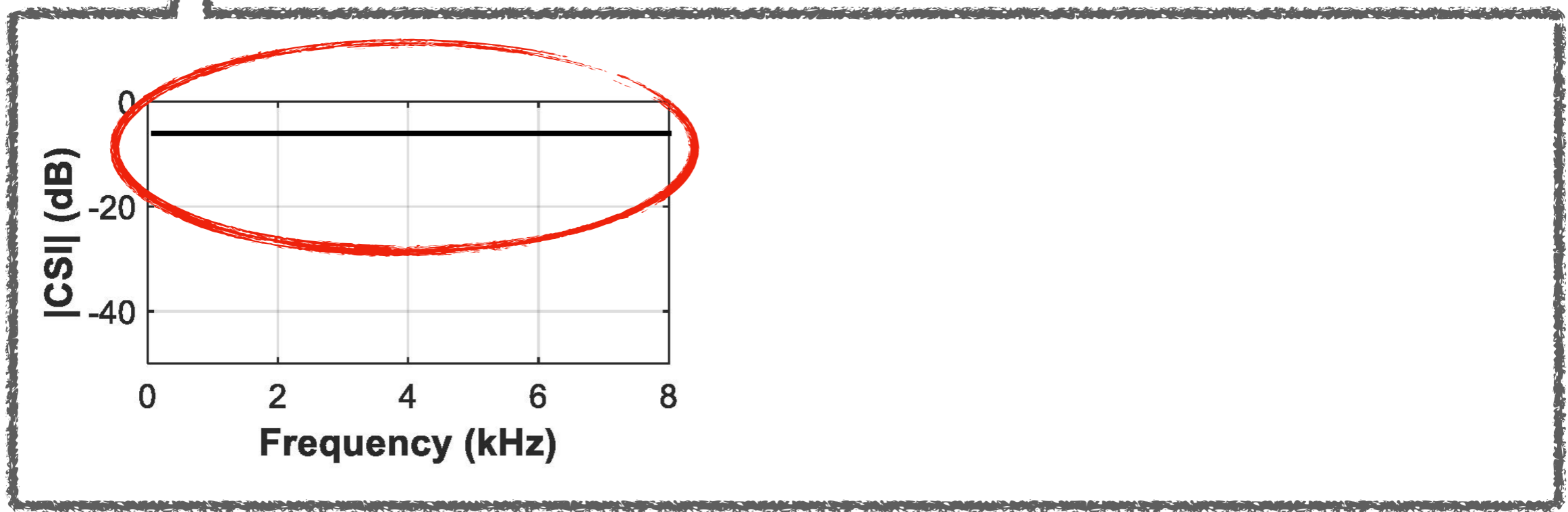
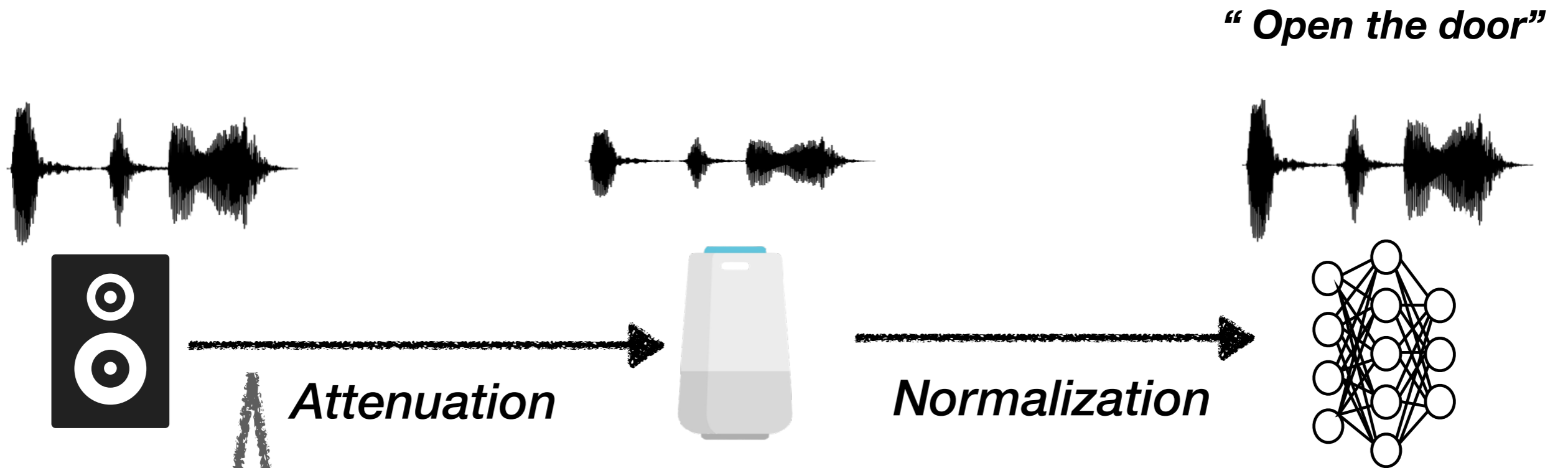
Attenuation



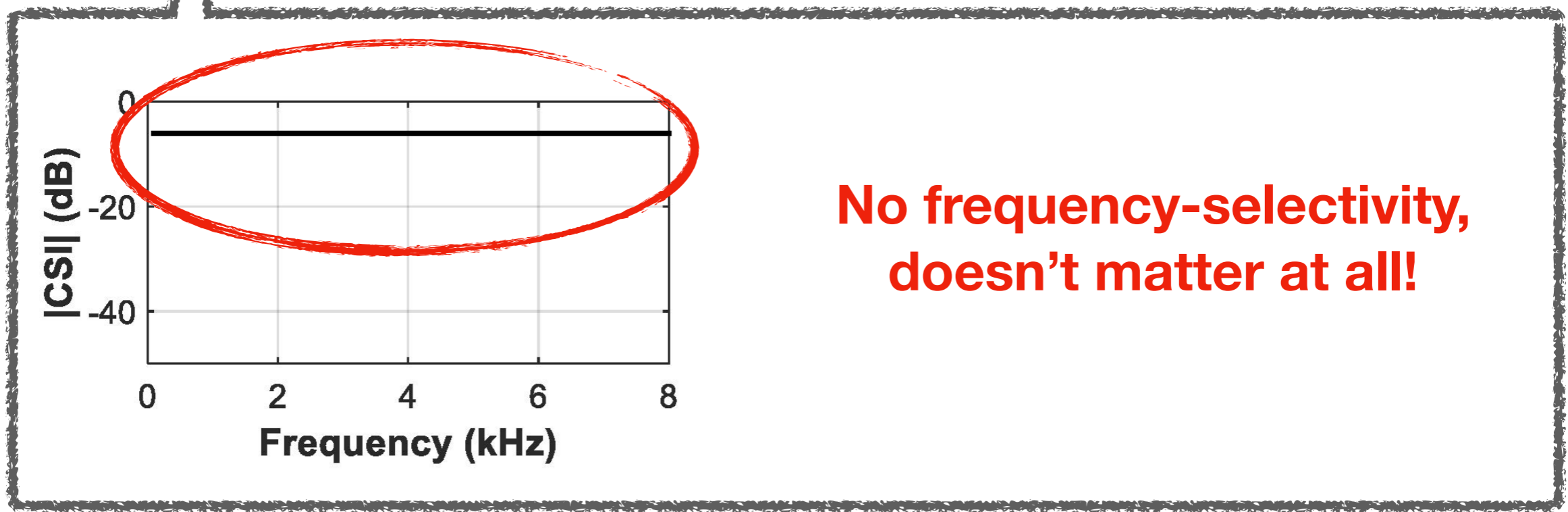
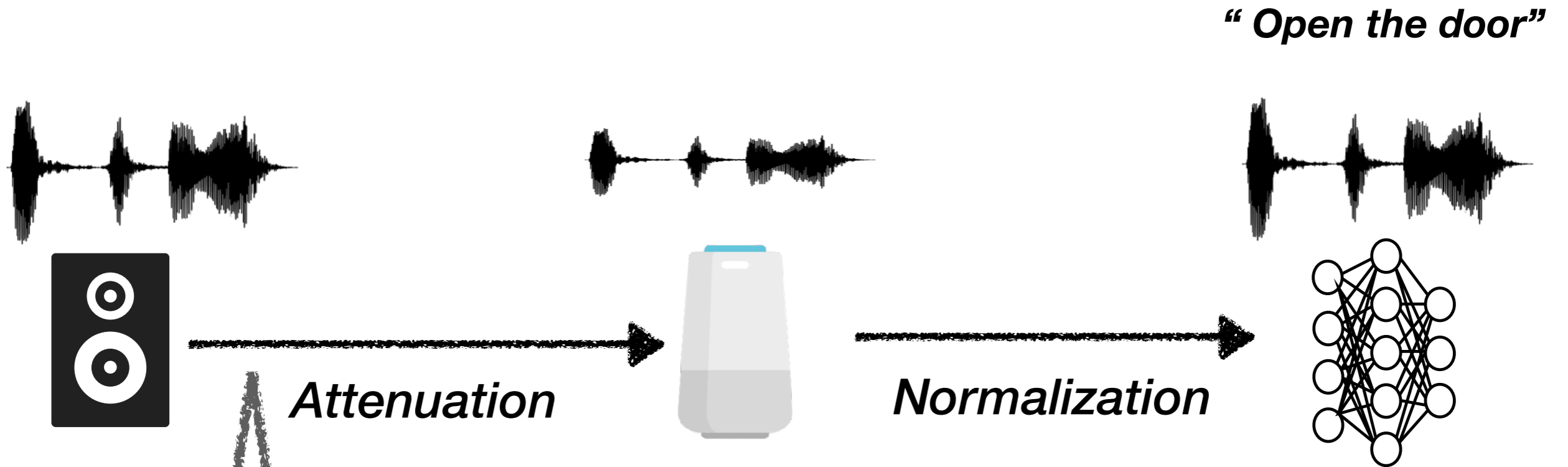
Attenuation



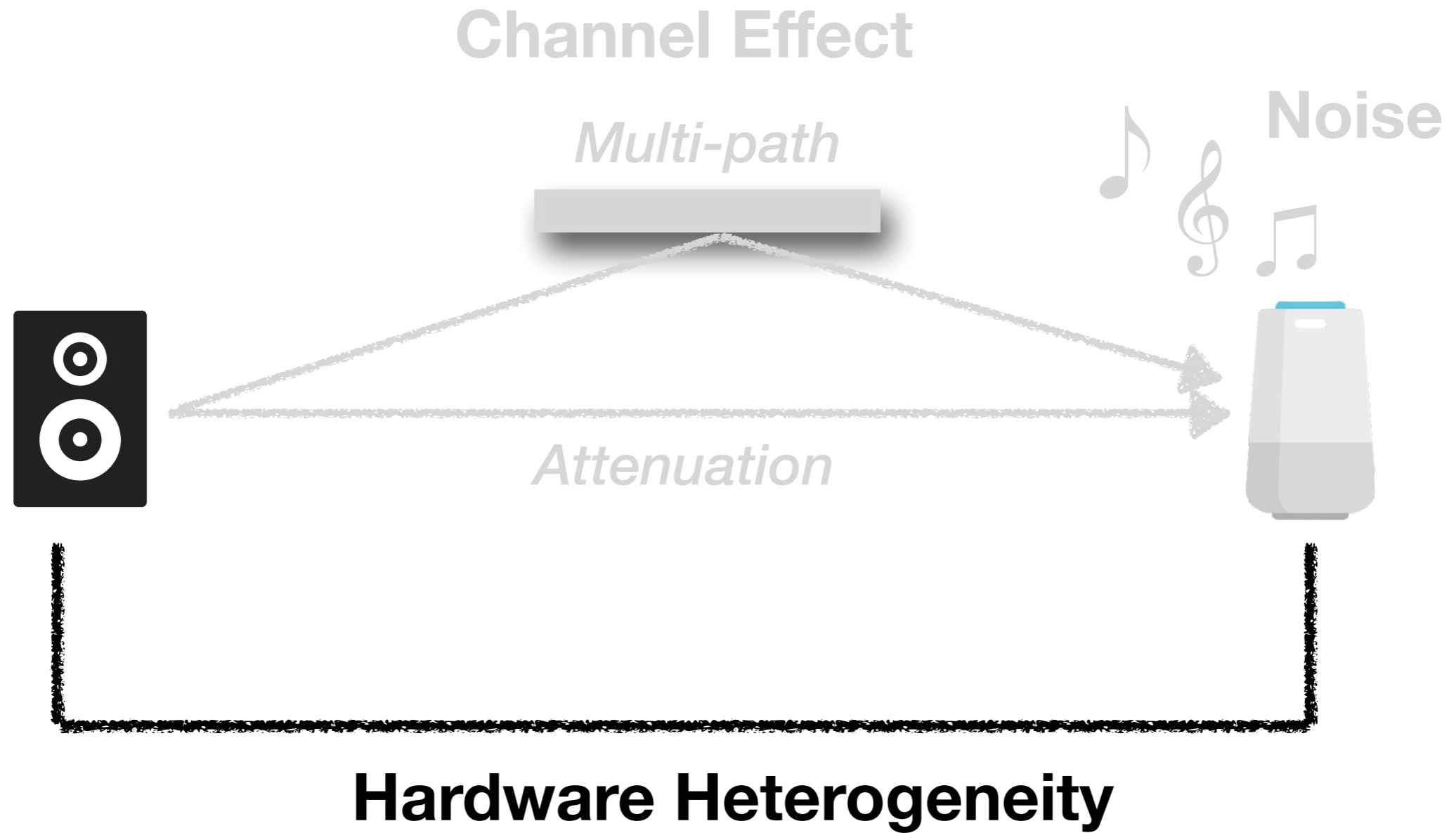
Attenuation



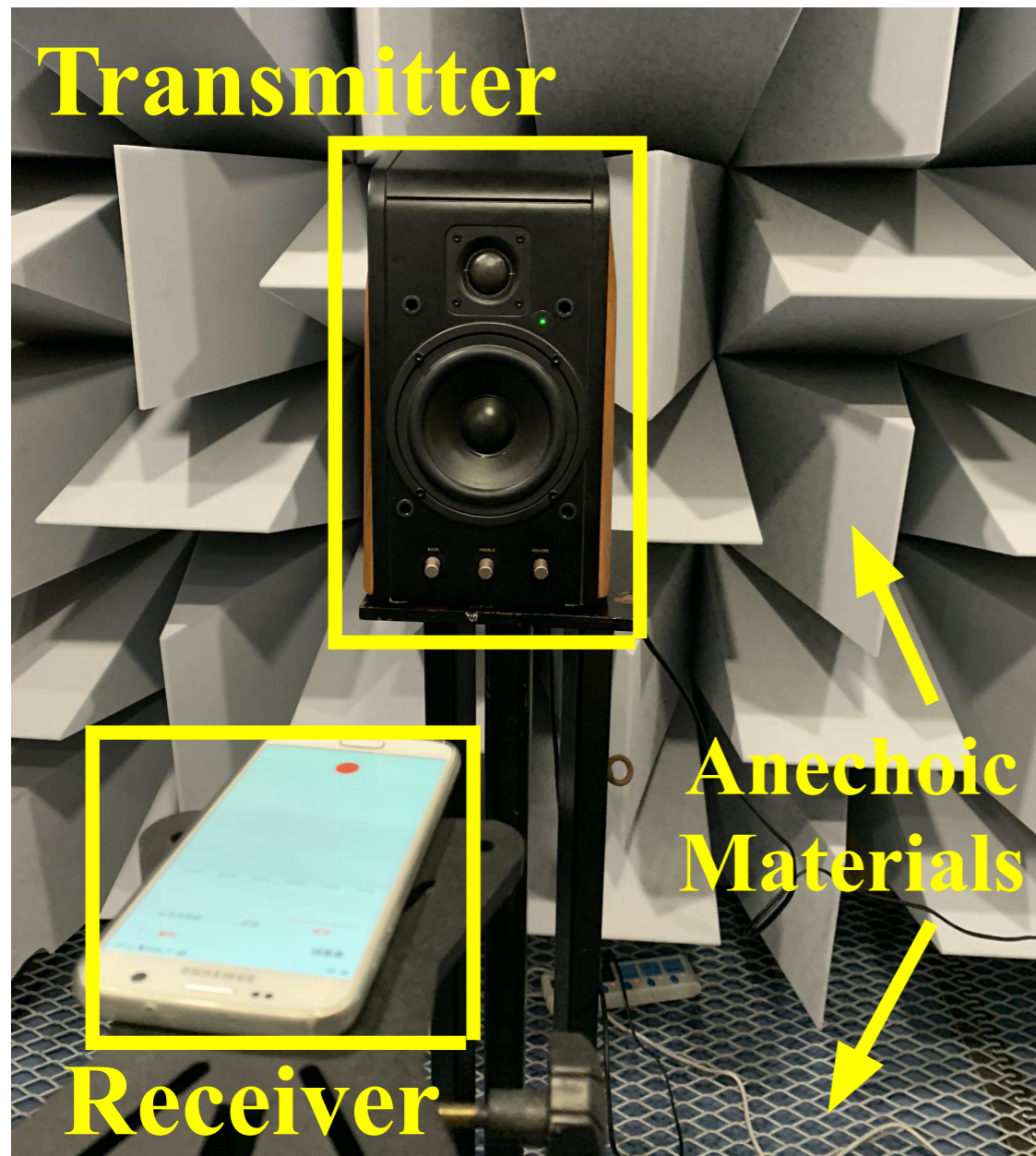
Attenuation



Understand Over-the-air Attack

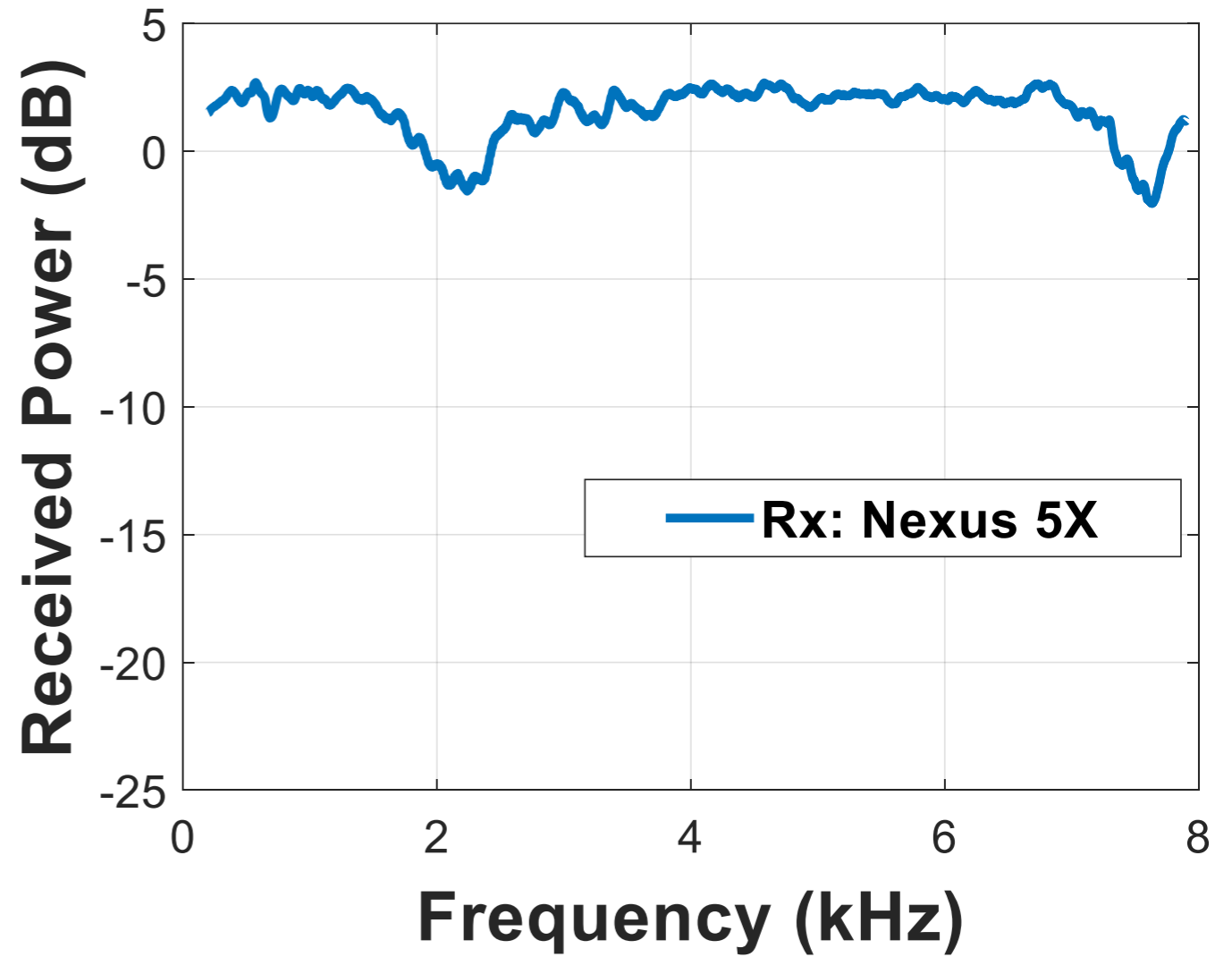
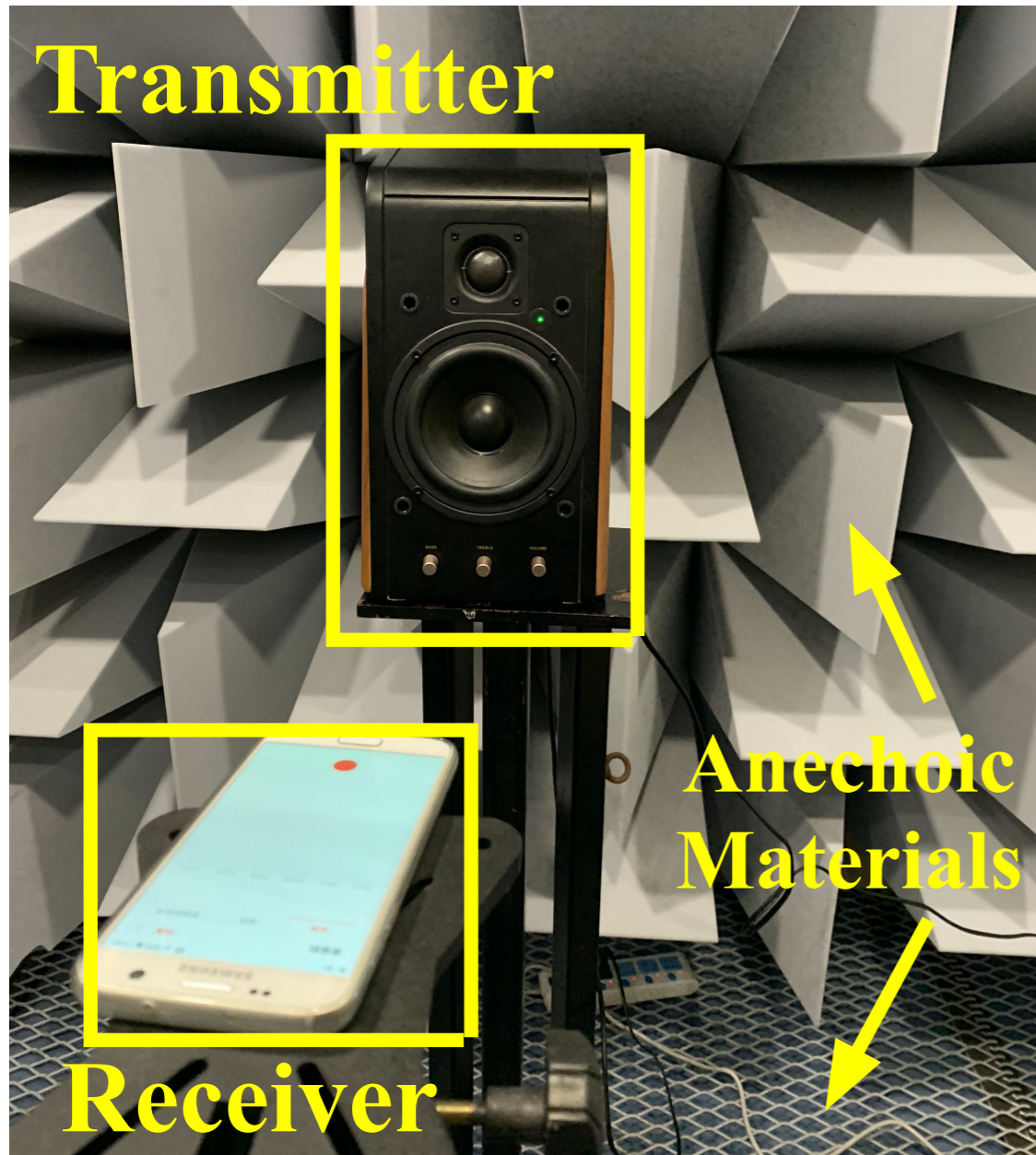


Hardware Heterogeneity



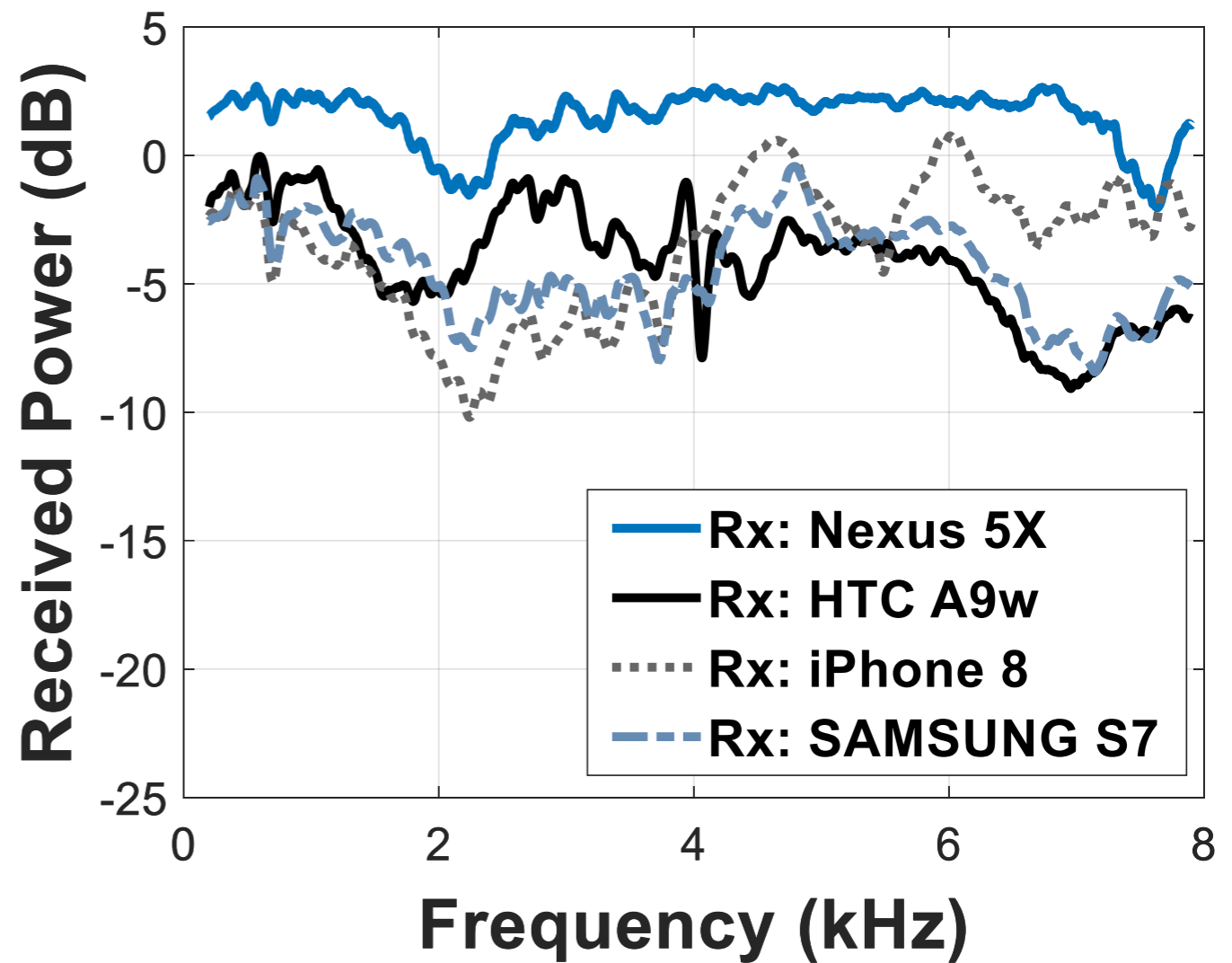
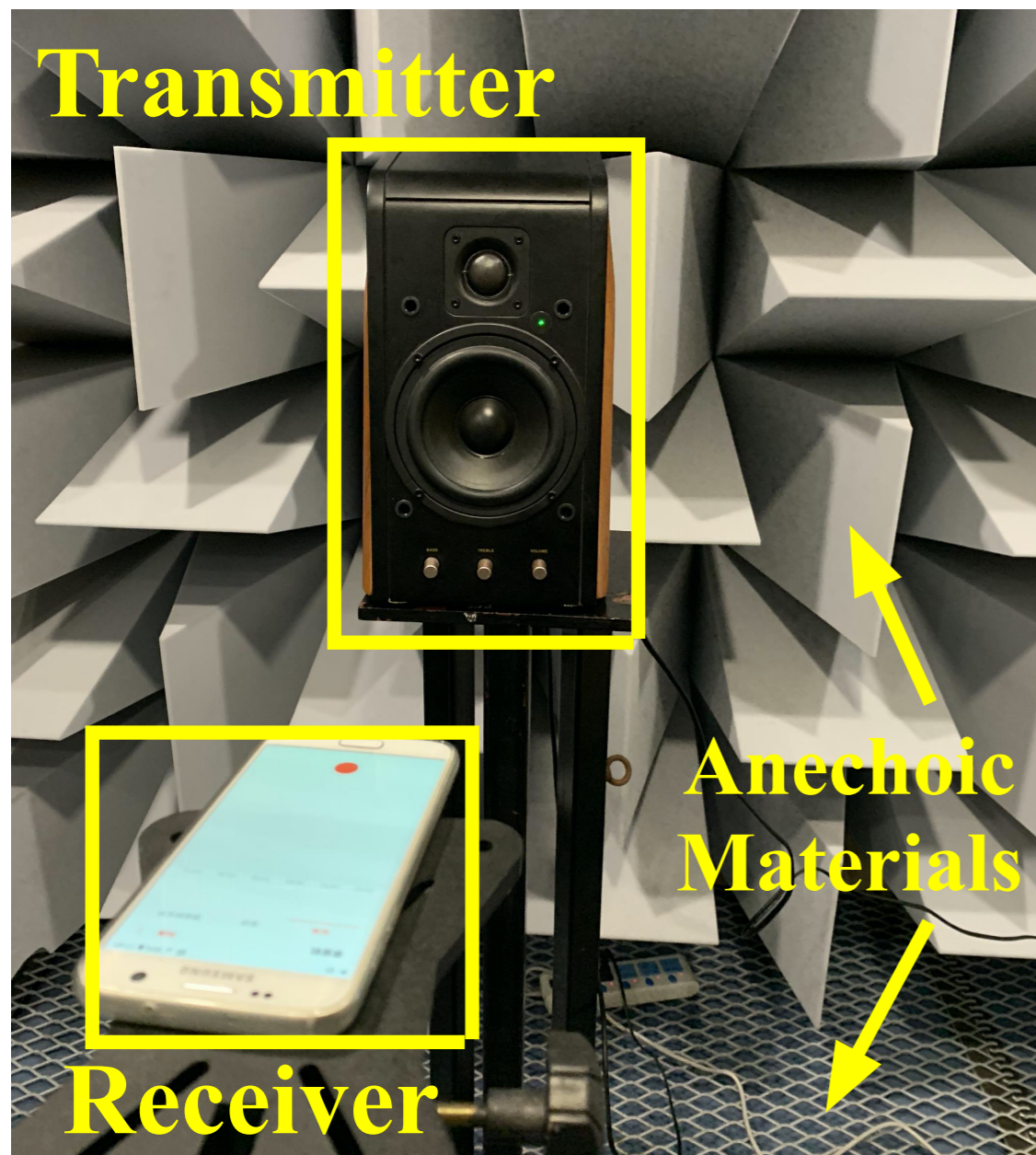
Anechoic Chamber Testing

Hardware Heterogeneity



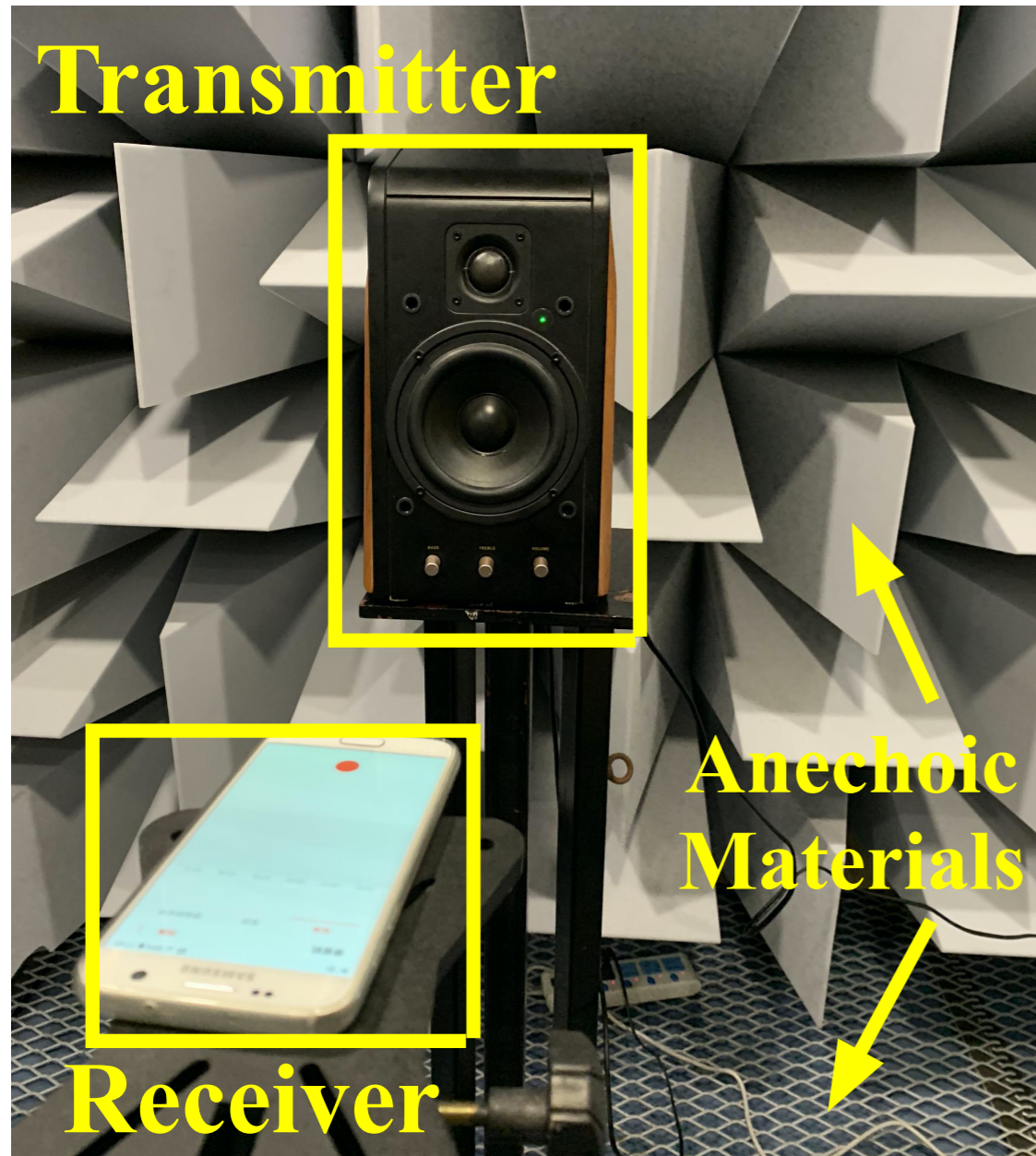
Anechoic Chamber Testing

Hardware Heterogeneity

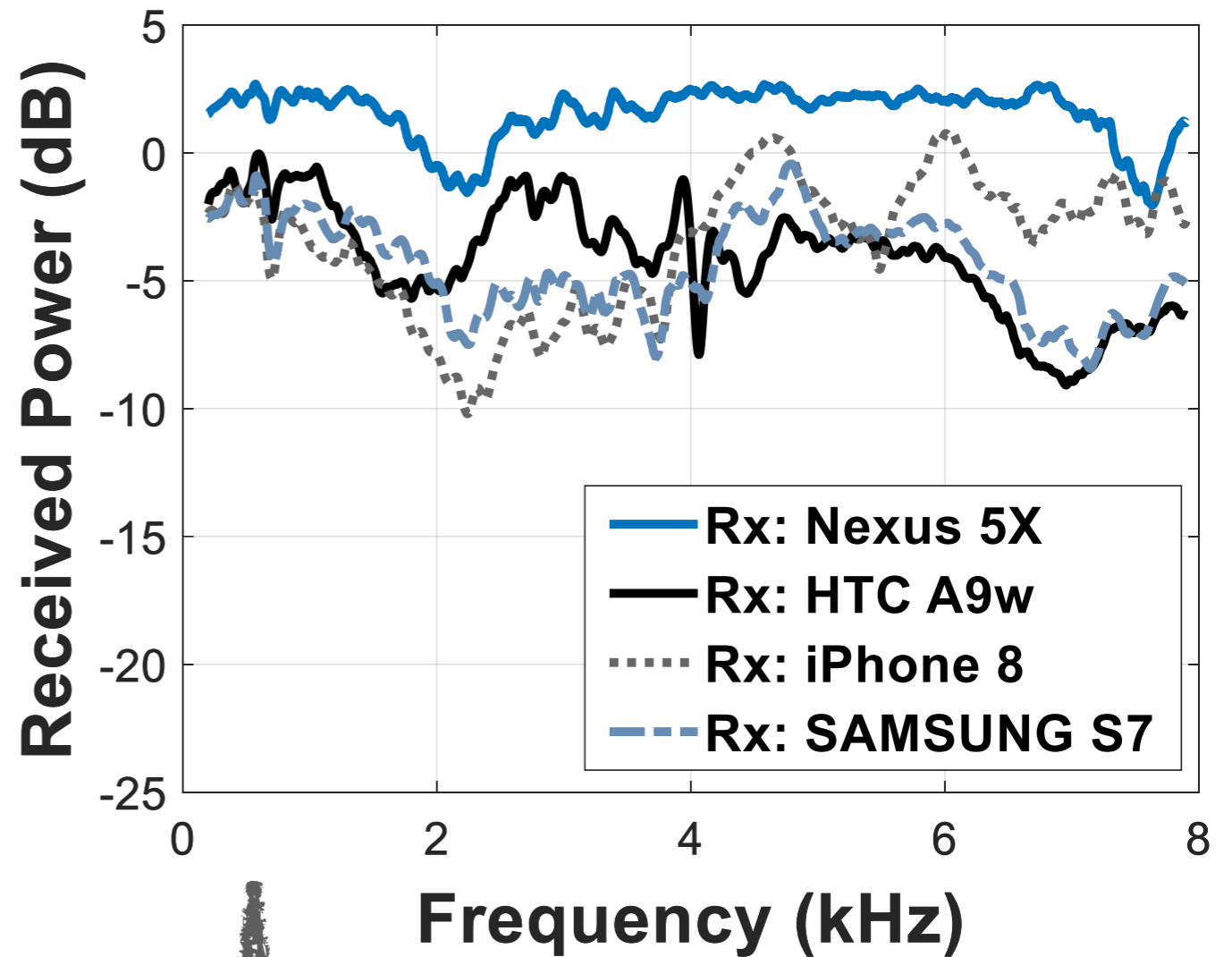


Anechoic Chamber Testing

Hardware Heterogeneity

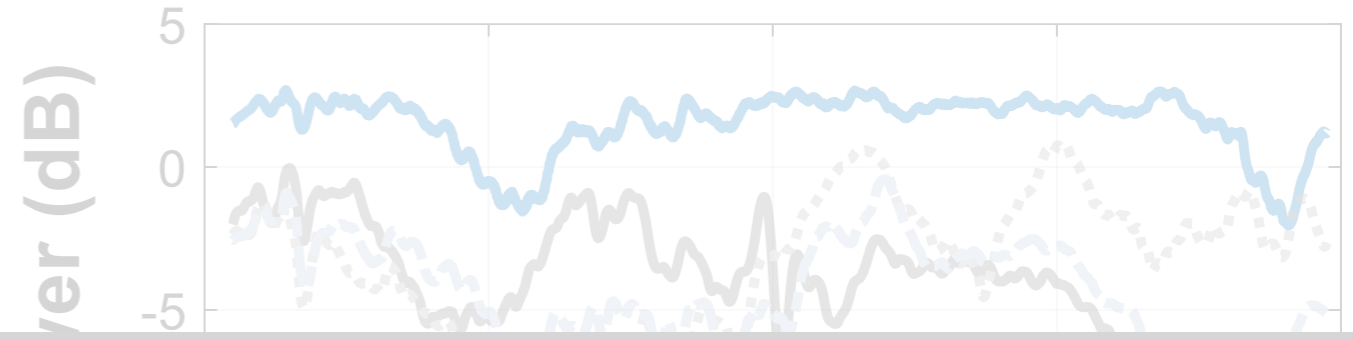


Anechoic Chamber Testing

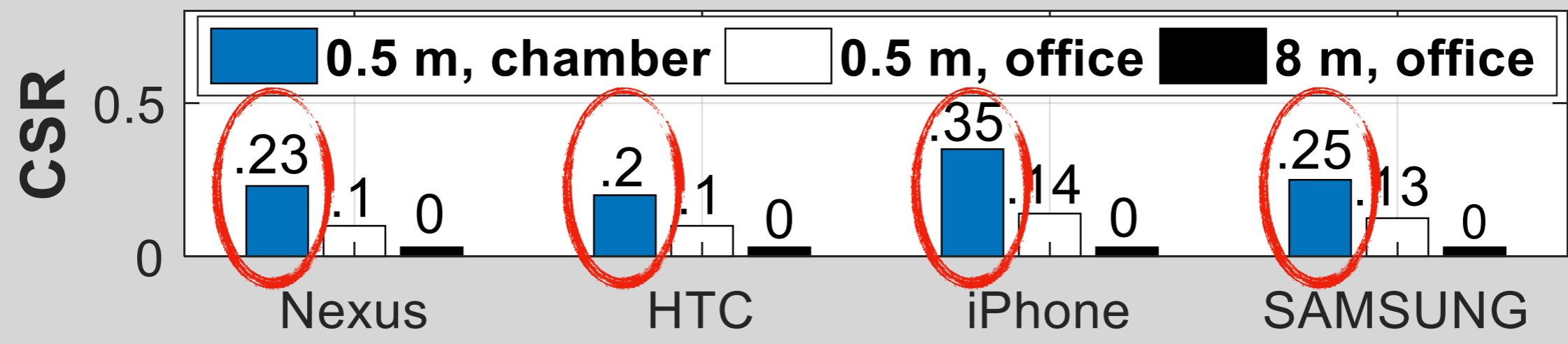


Not strong, device's inherent feature, compensable!

Hardware Heterogeneity



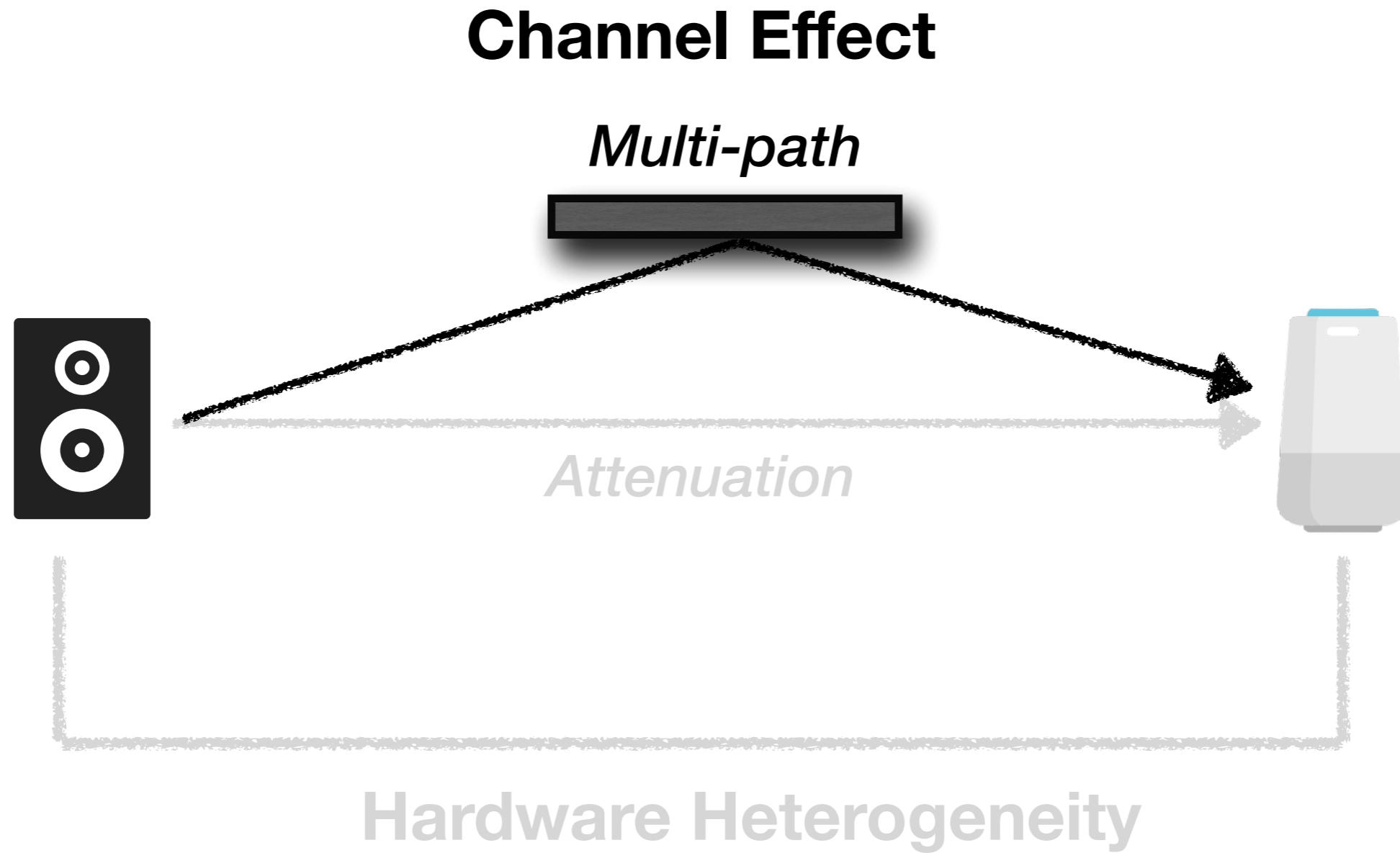
Character Successful Rate (CSR):



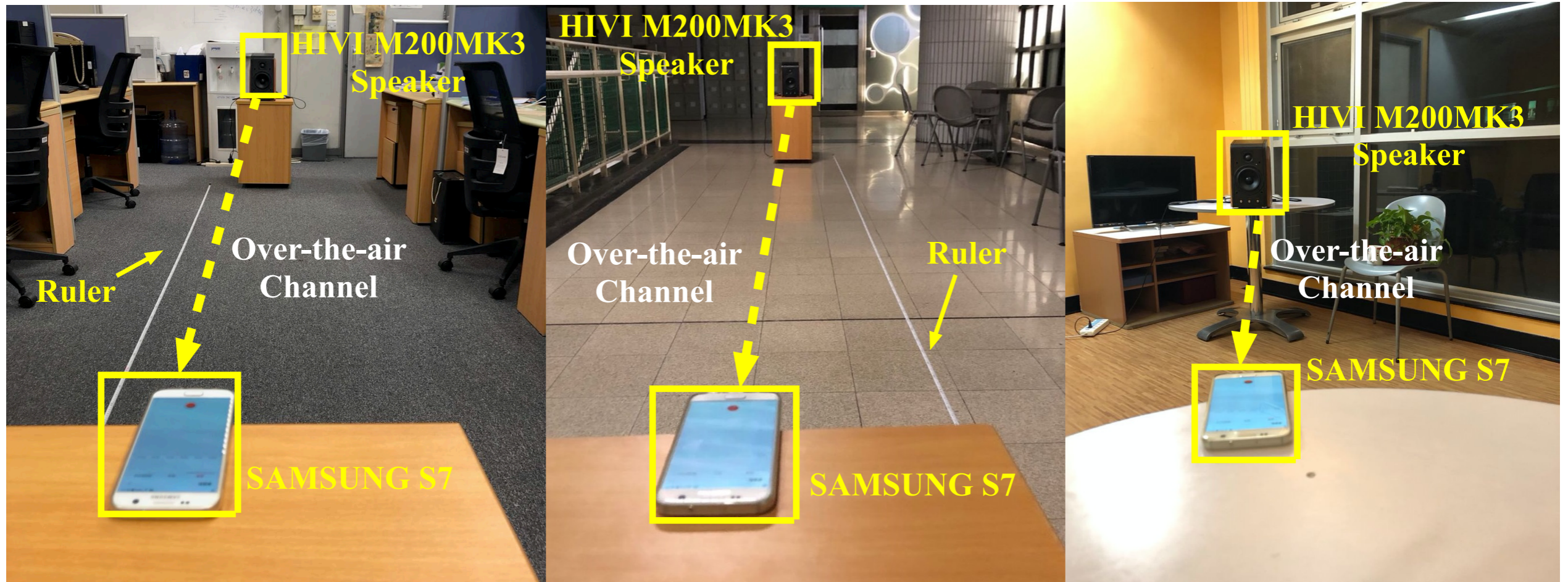
Static, predictable and compensable!

Anechoic Chamber Testing

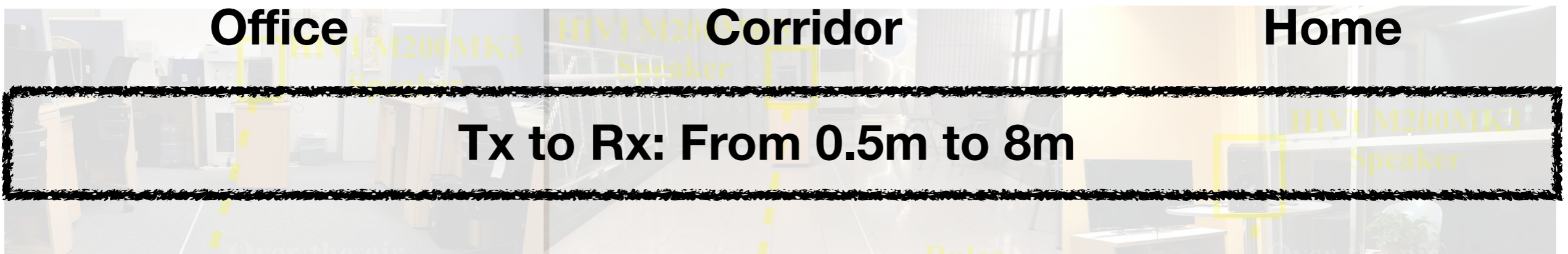
Understand Over-the-air Attack



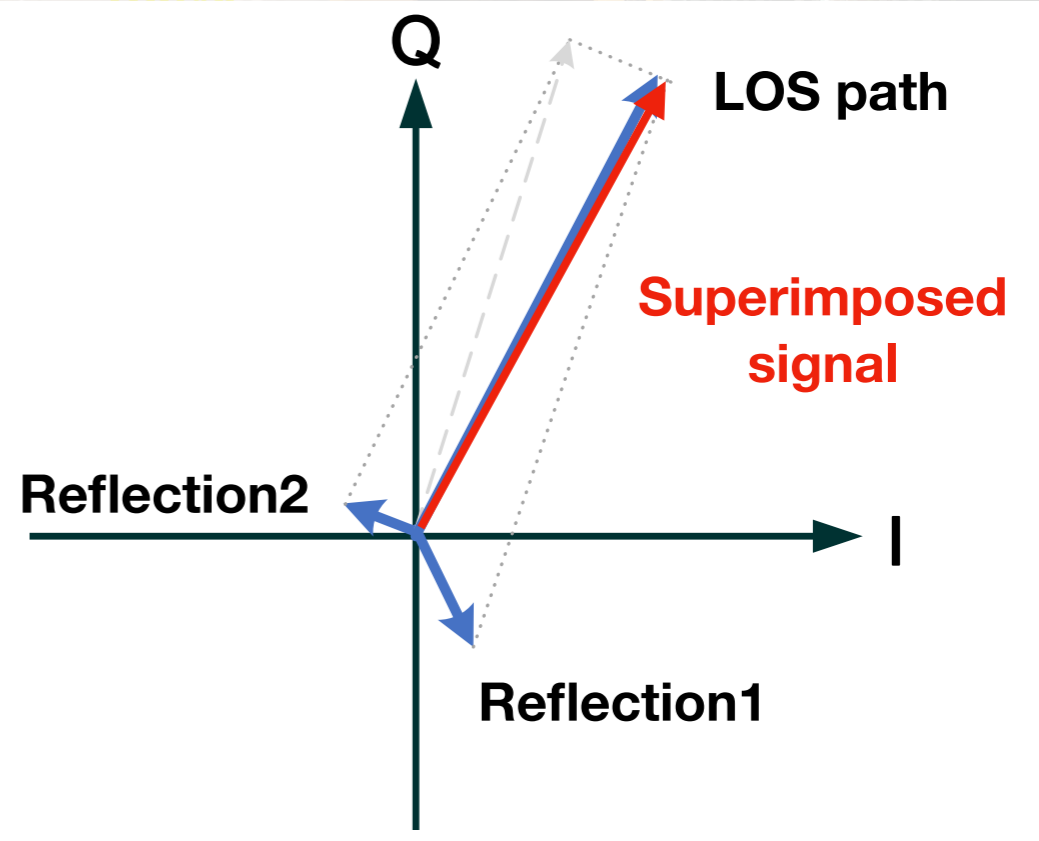
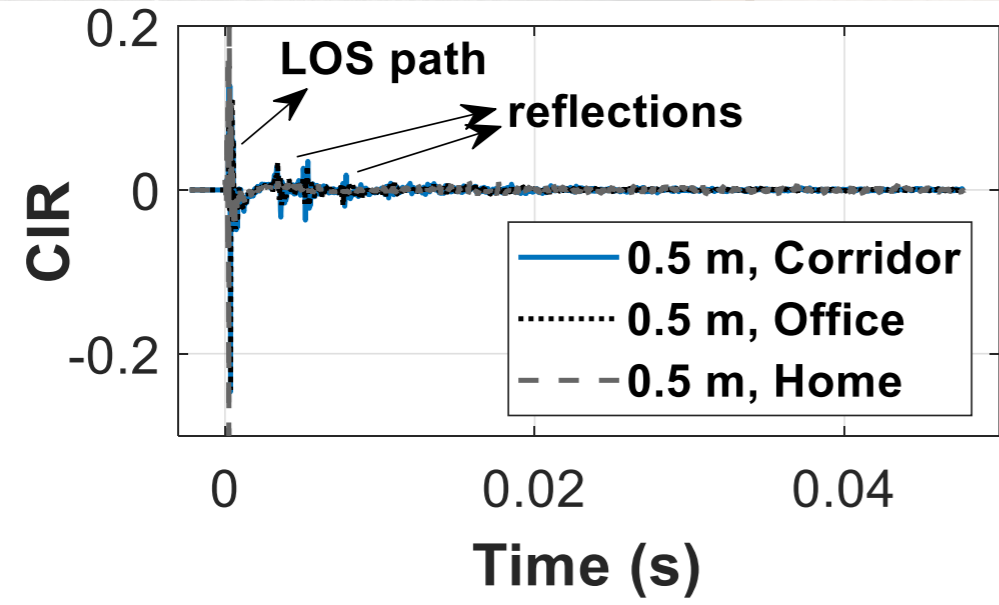
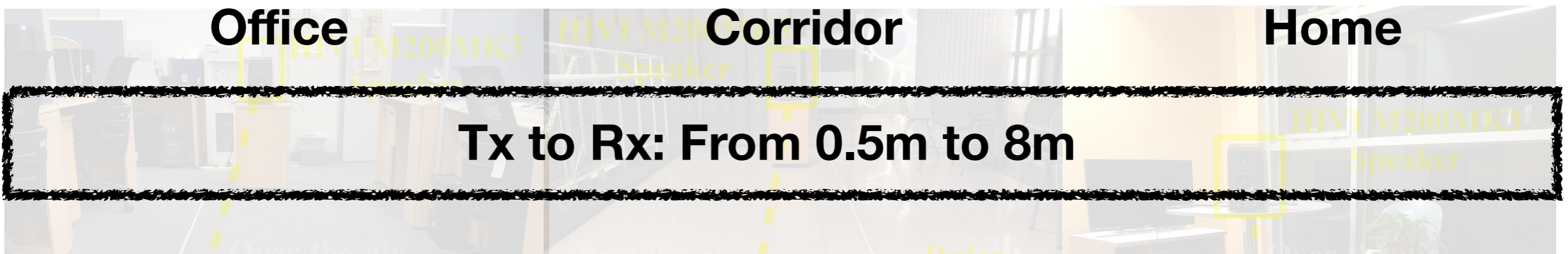
Multi-path



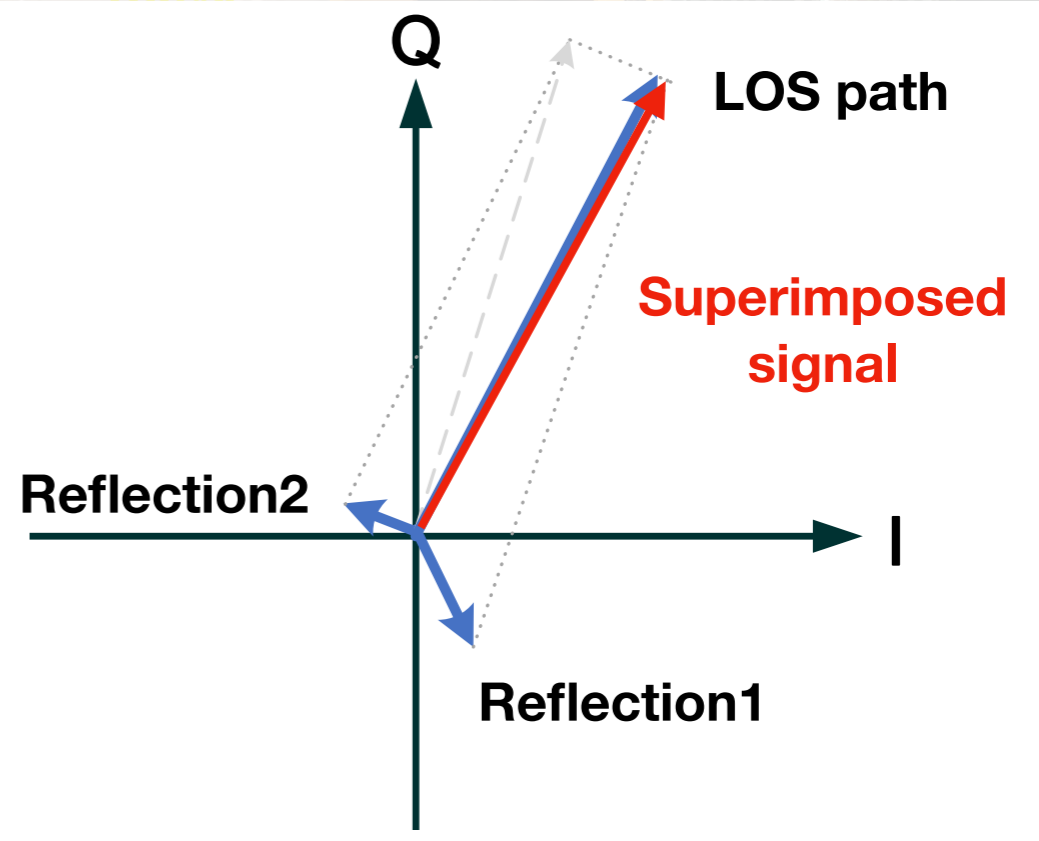
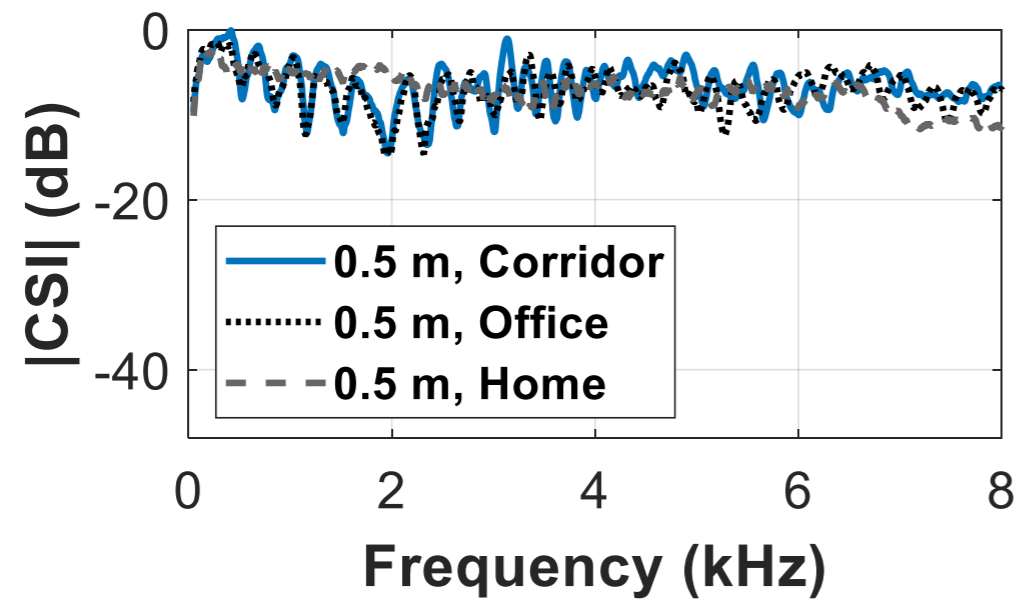
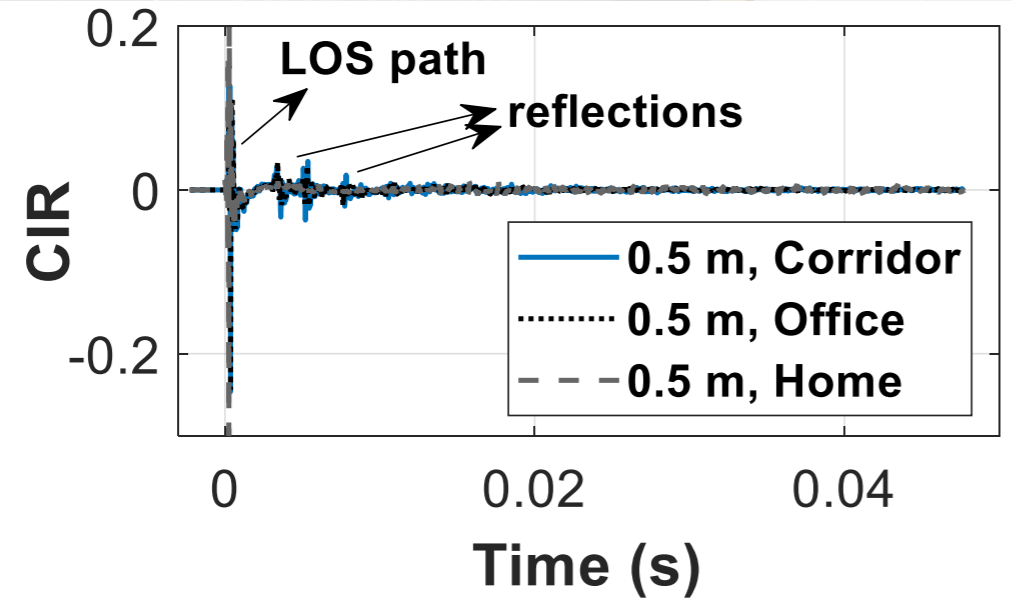
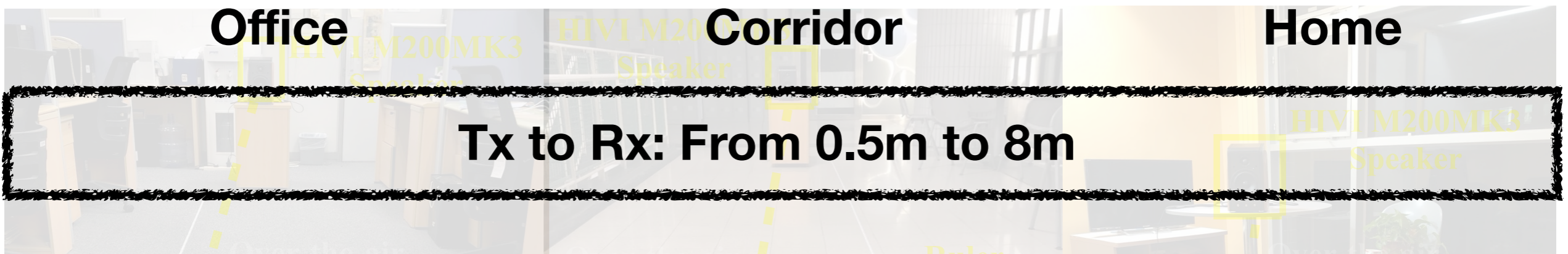
Multi-path: Near range



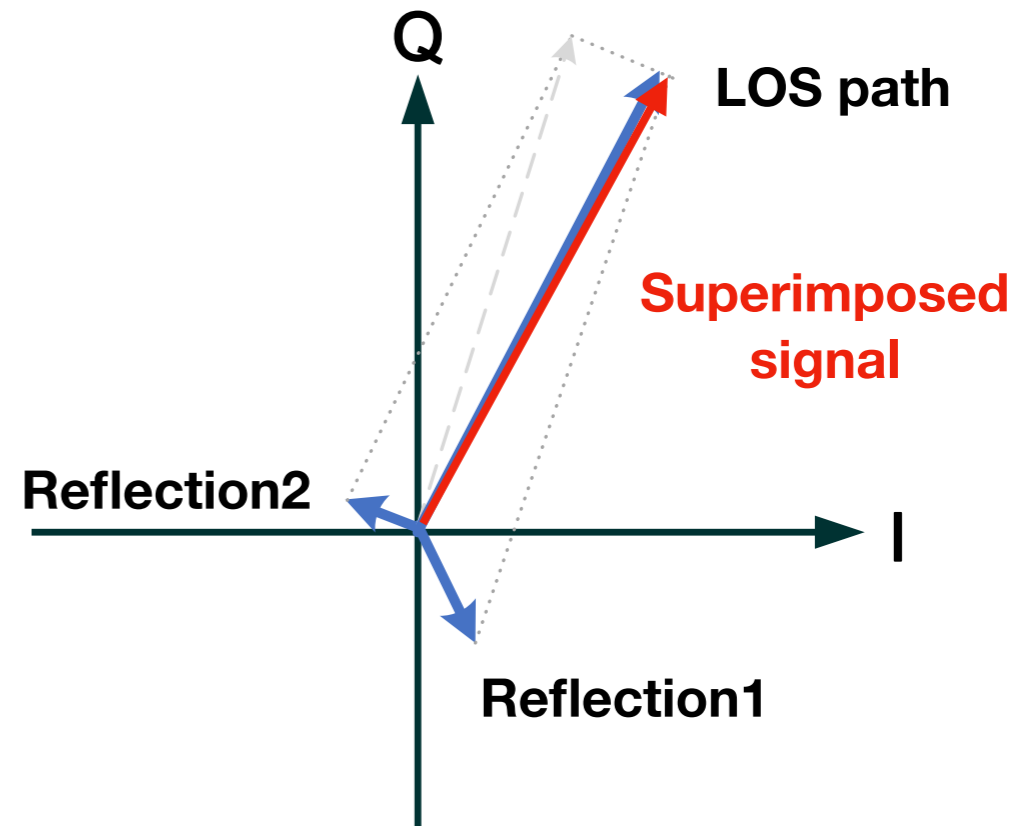
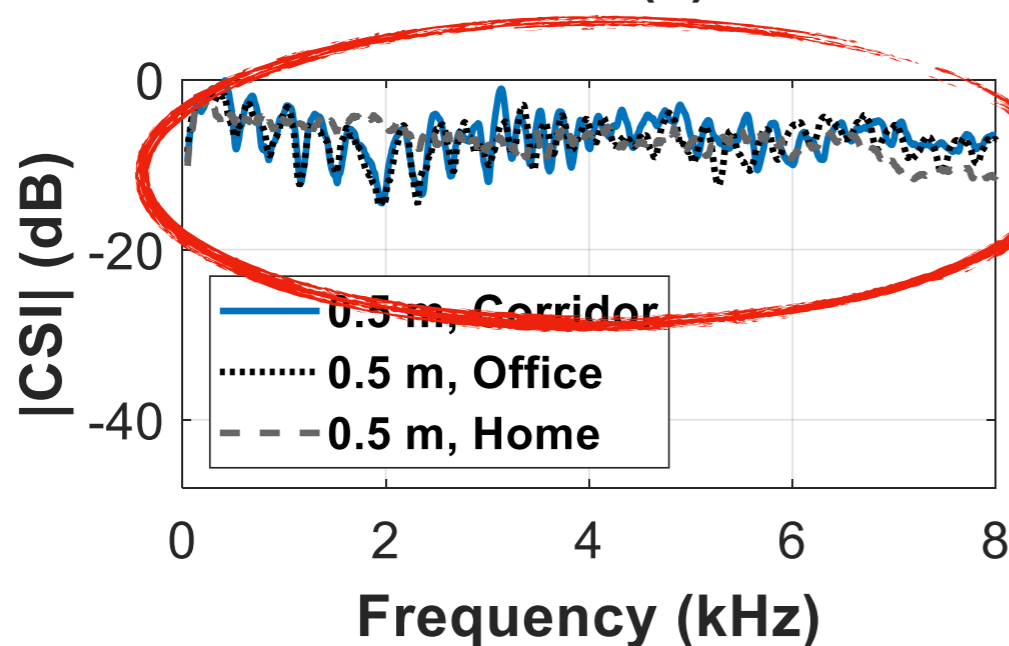
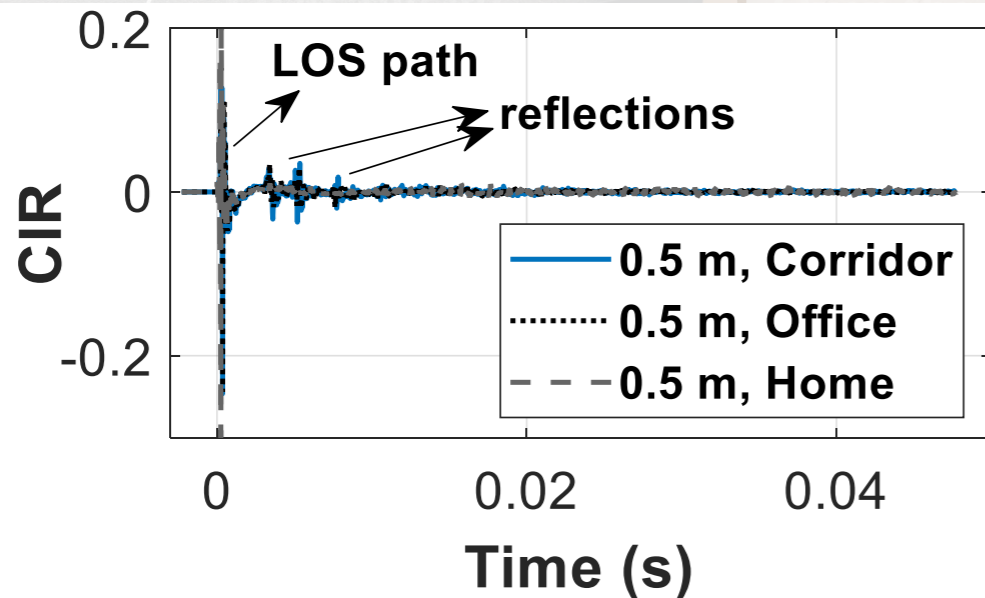
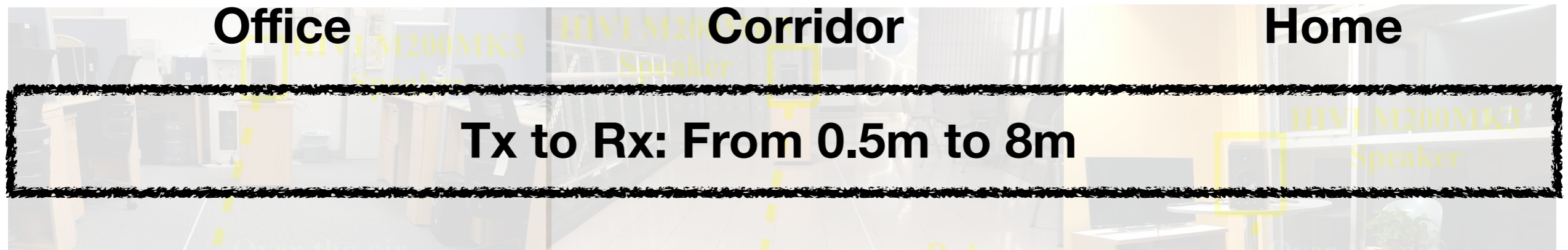
Multi-path: Near range



Multi-path: Near range



Multi-path: Near range

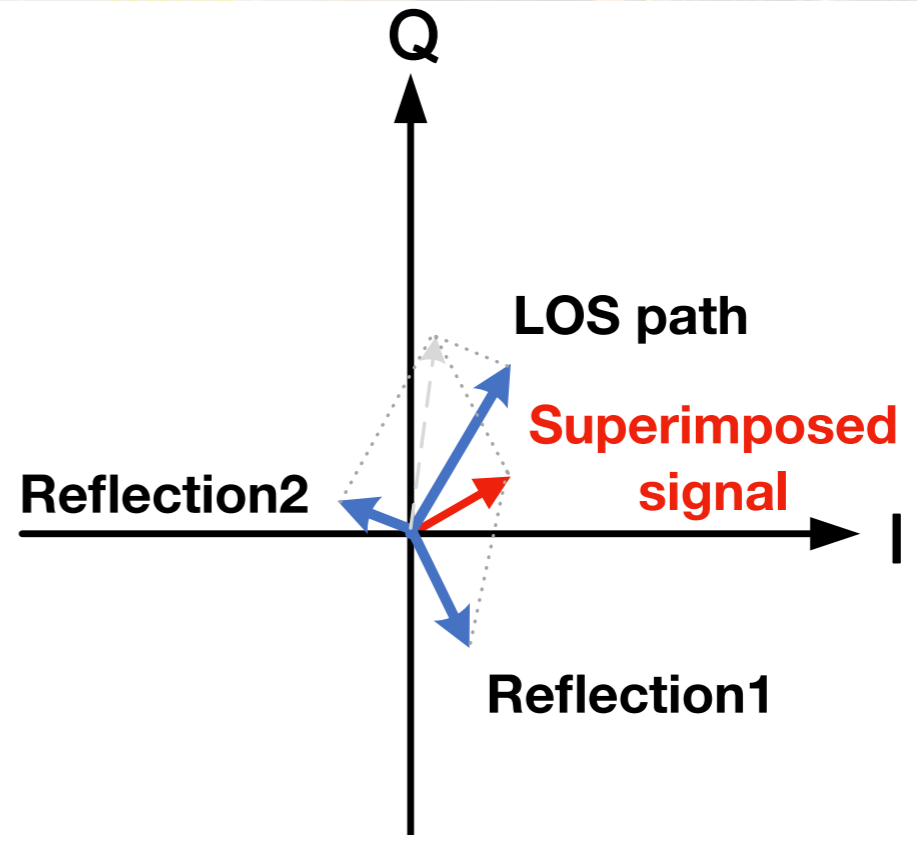
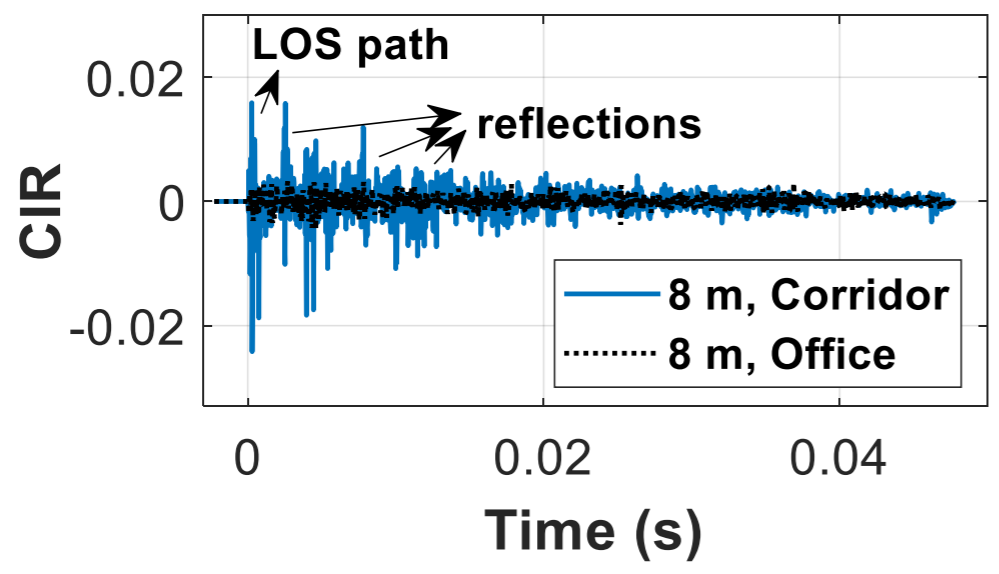
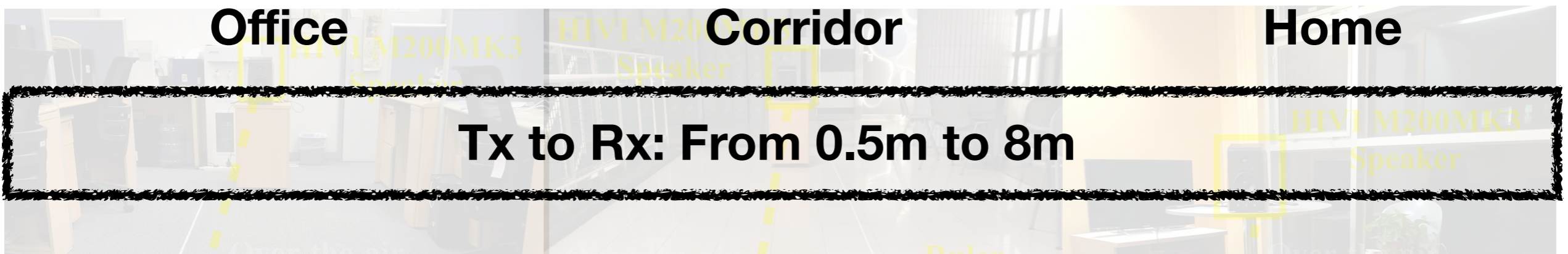


Also not strong and similar!

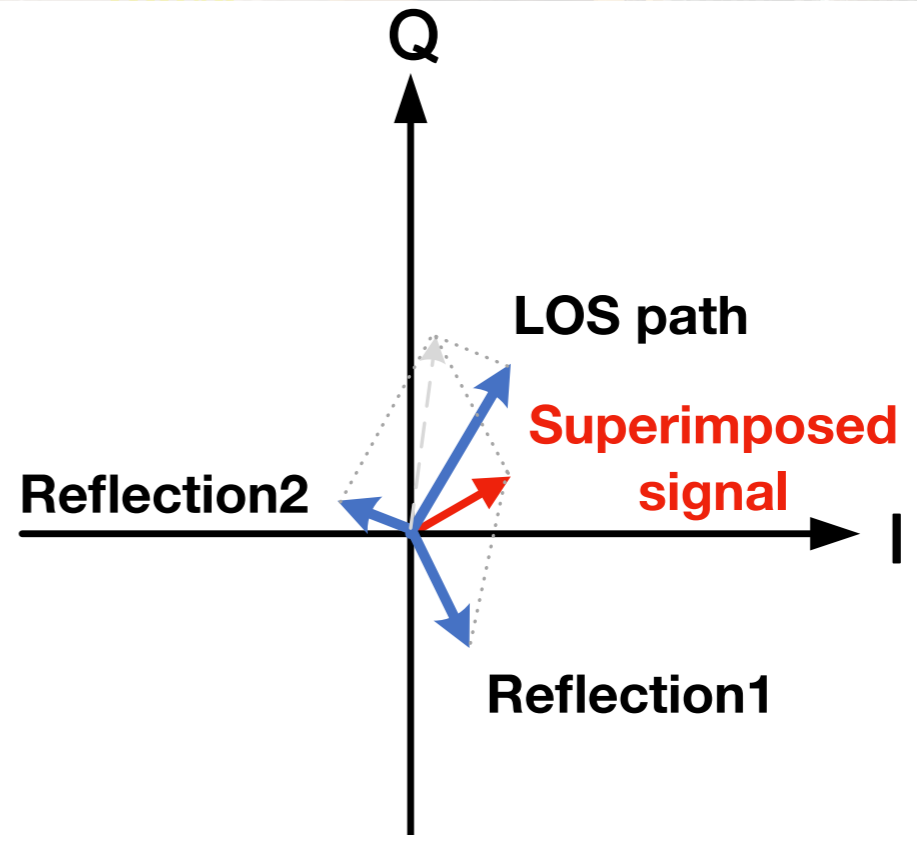
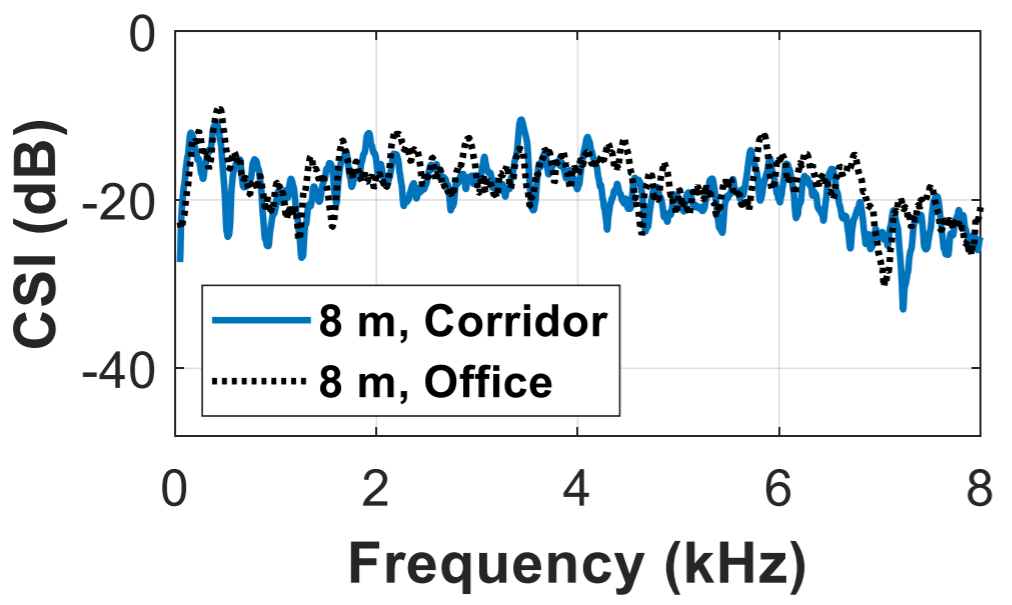
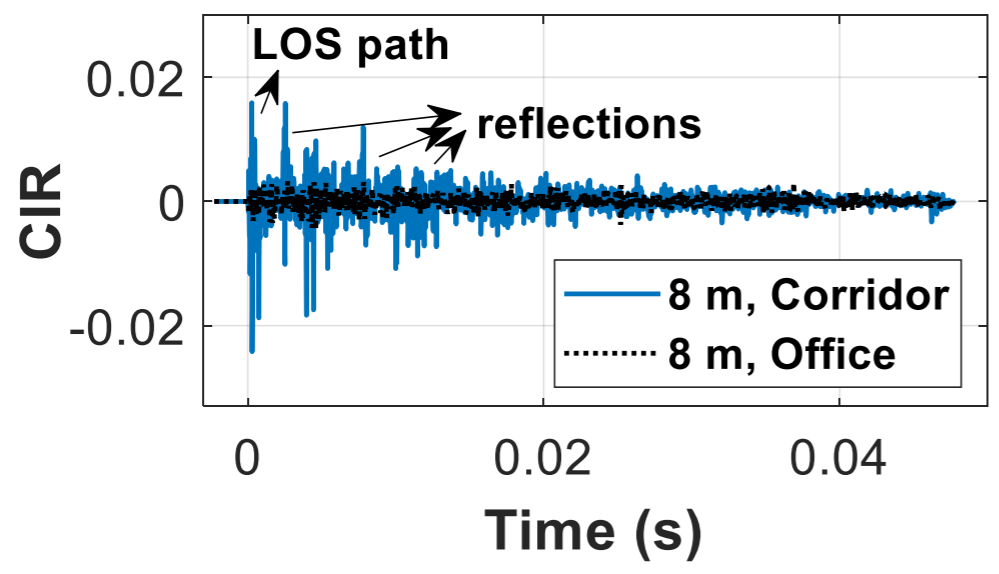
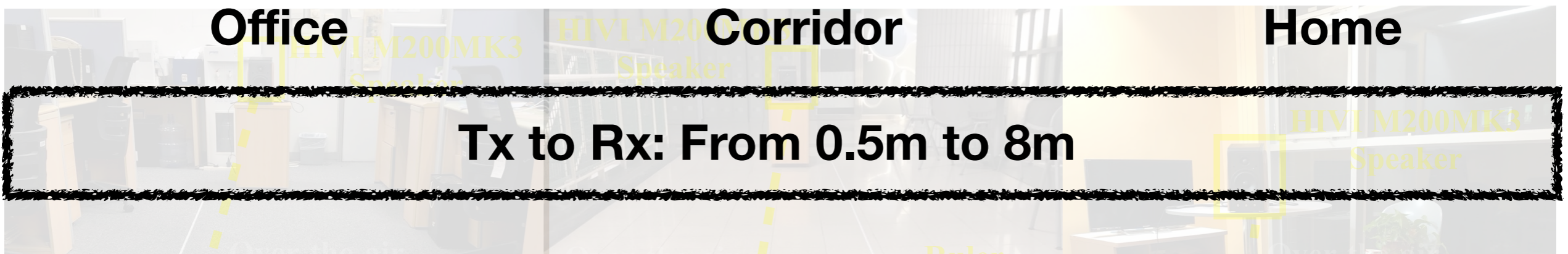
Multi-path: Long range



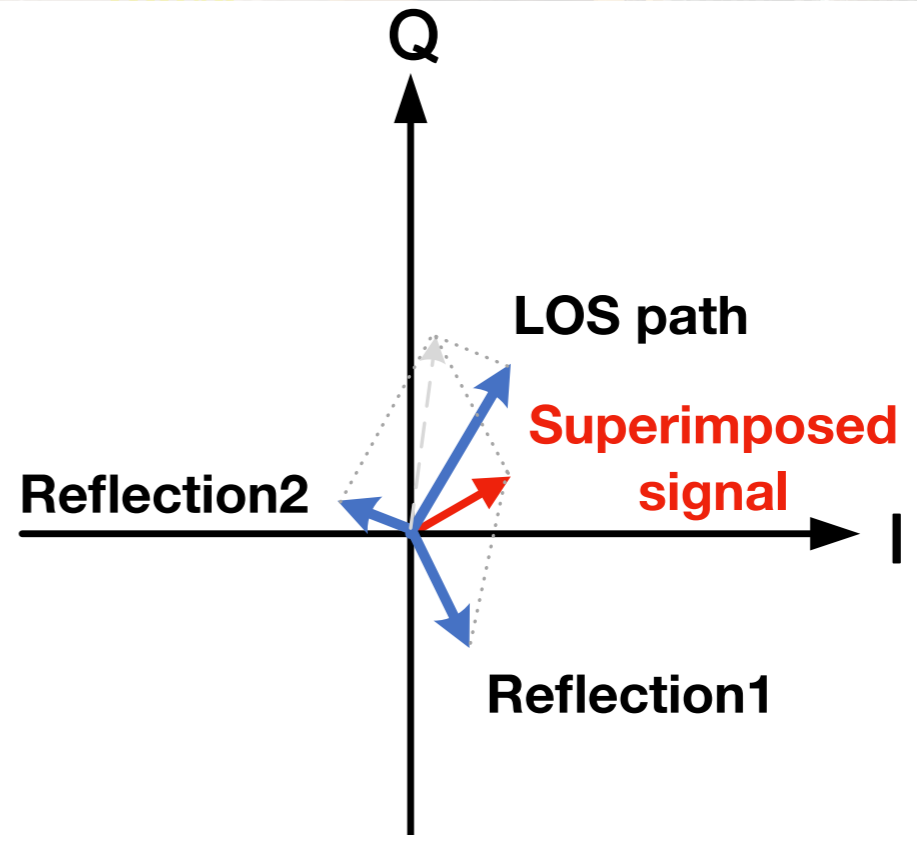
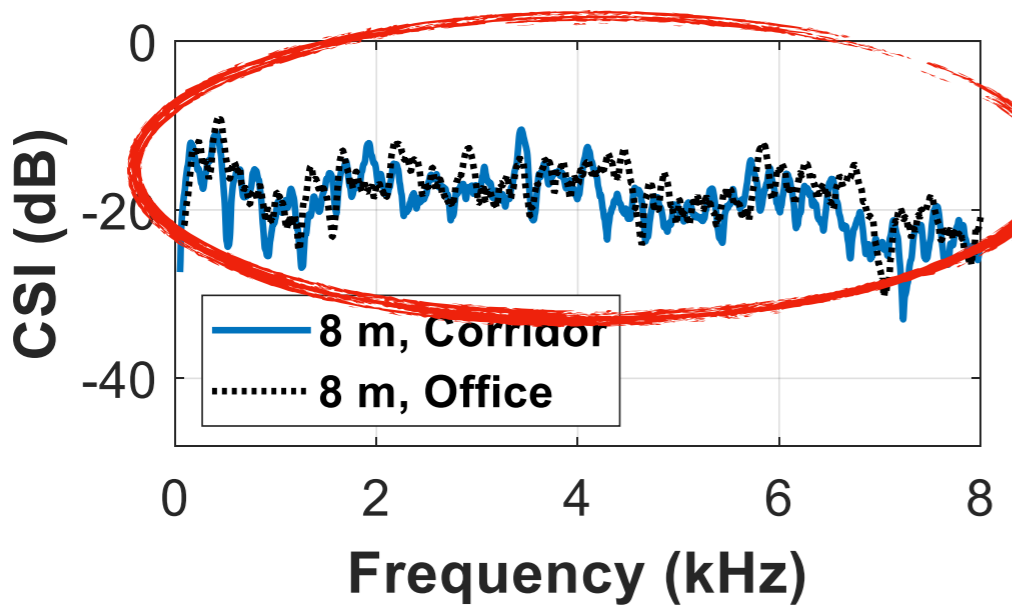
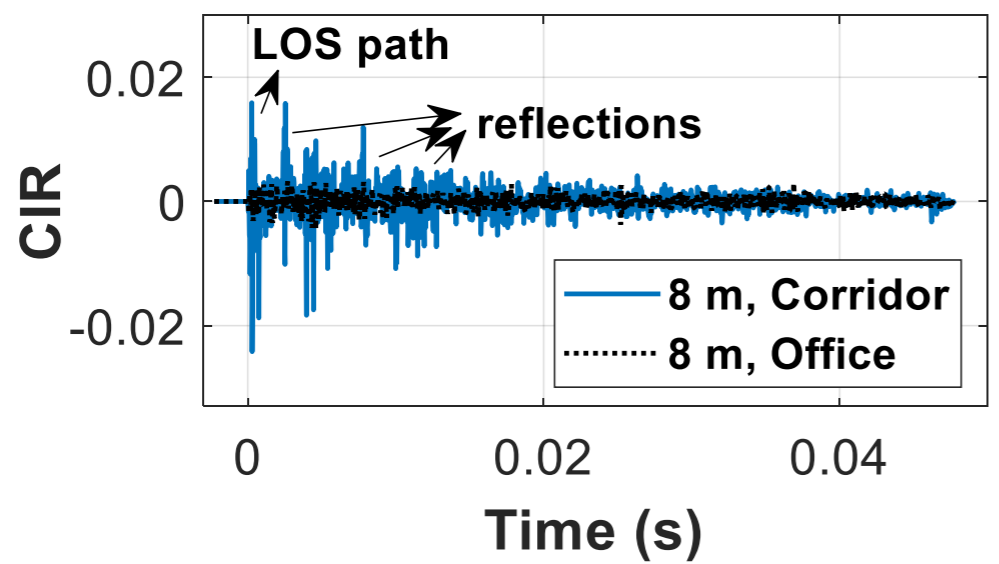
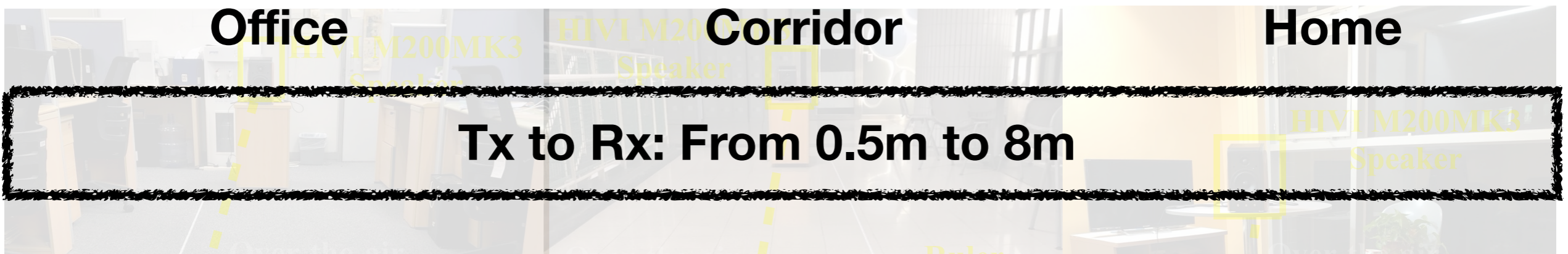
Multi-path: Long range



Multi-path: Long range



Multi-path: Long range

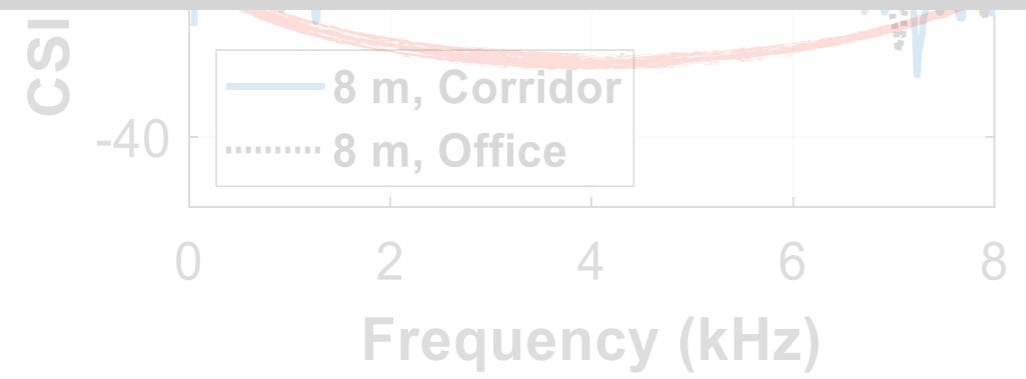
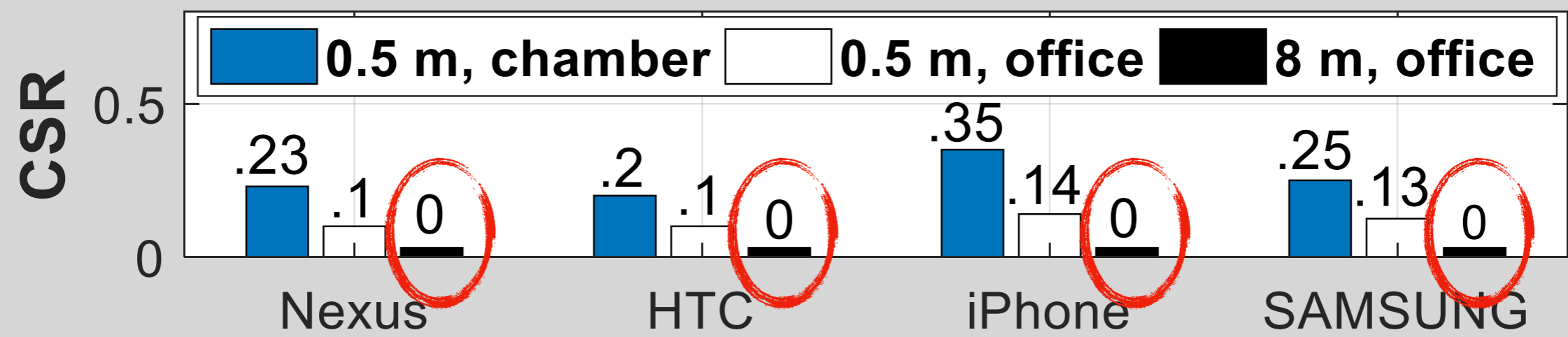


Stronger and unpredictable!

Multi-path: Long range

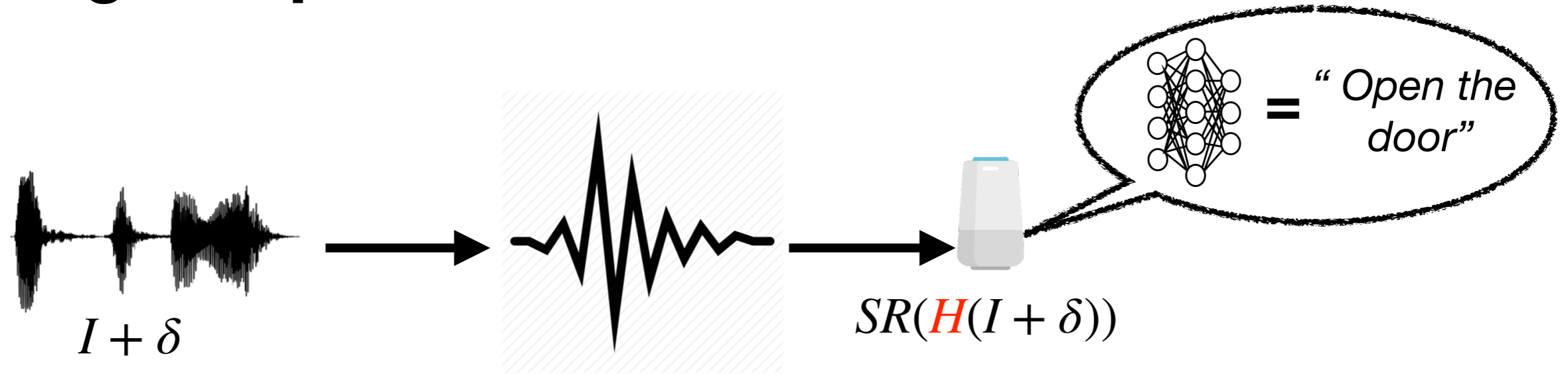


Character Successful Rate (CSR):



Highly unpredictable!

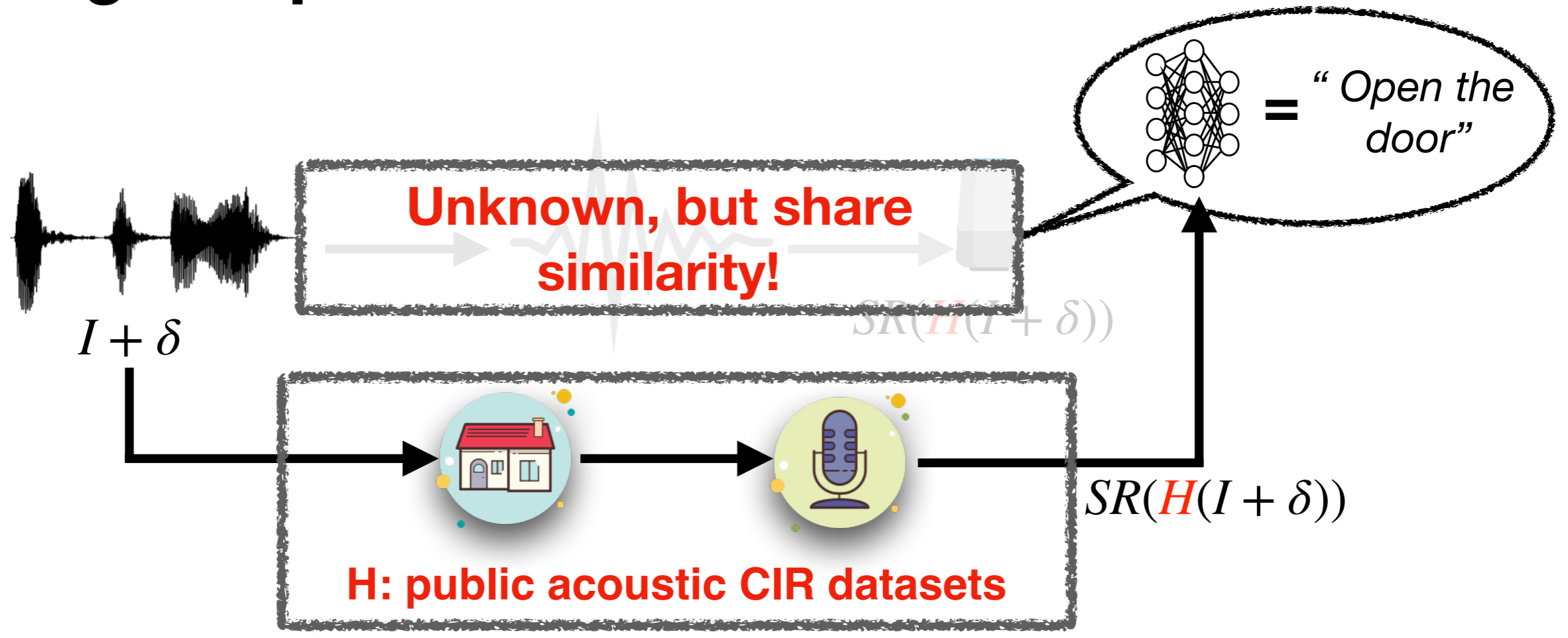
Design Inspiration



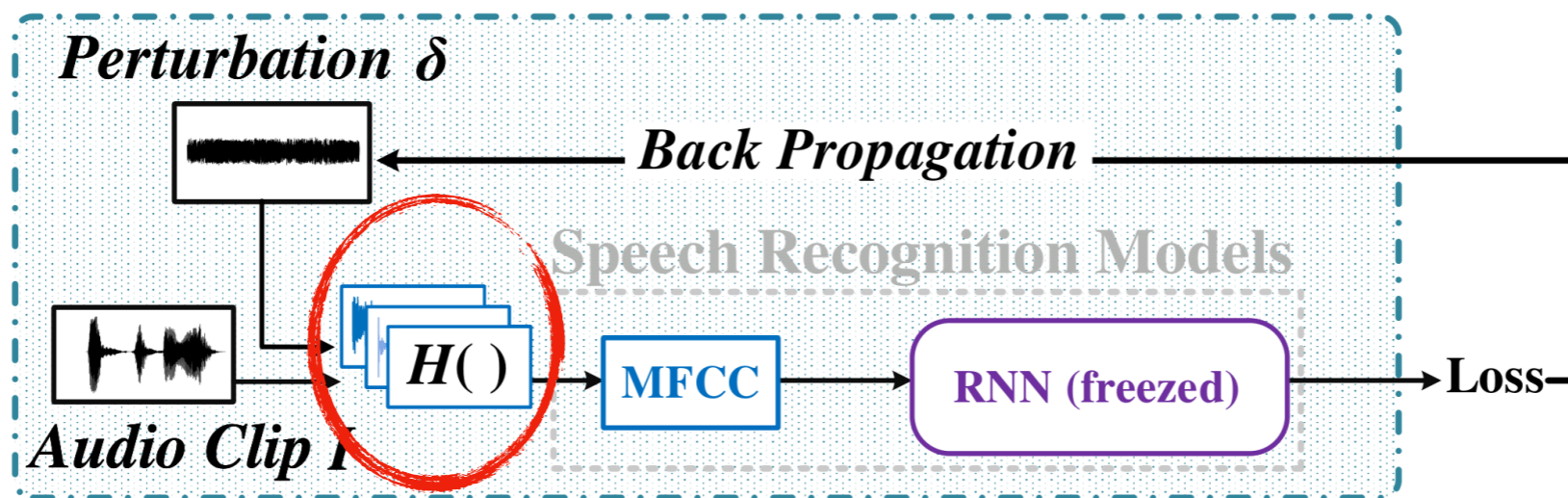
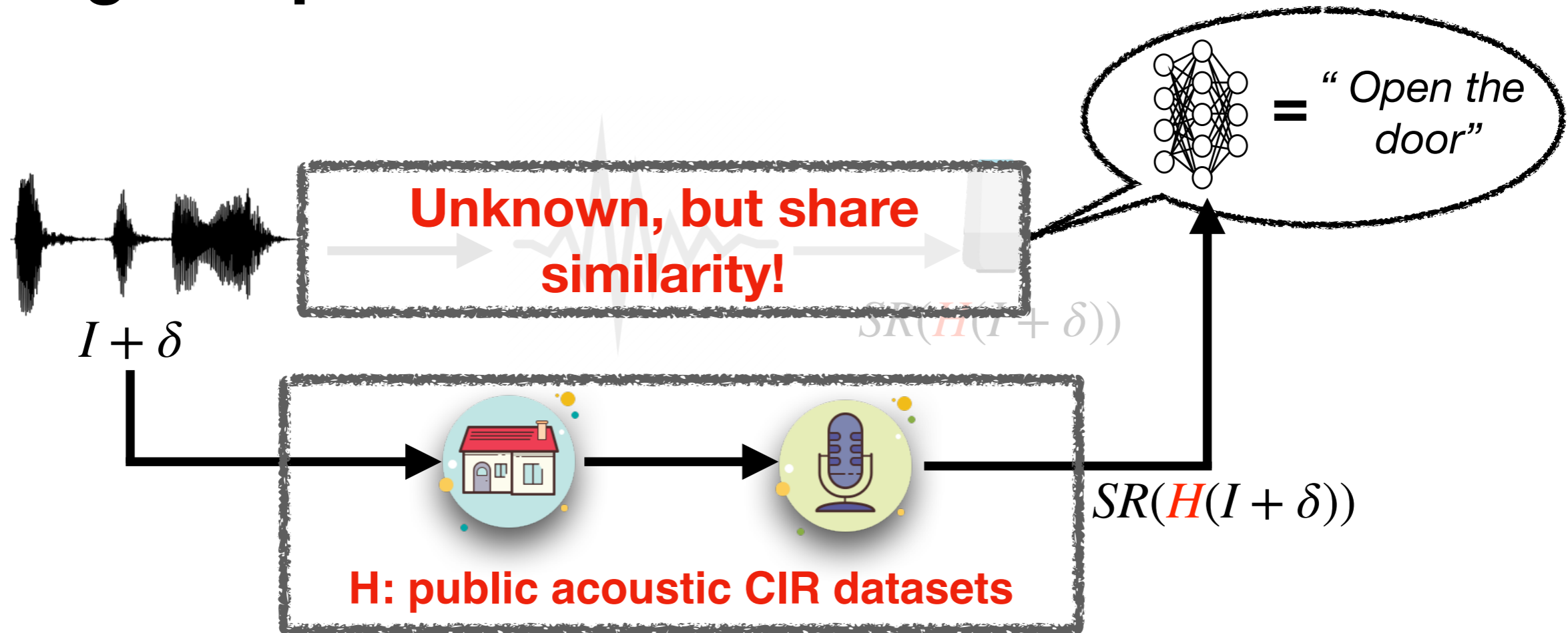
Design Inspiration



Design Inspiration

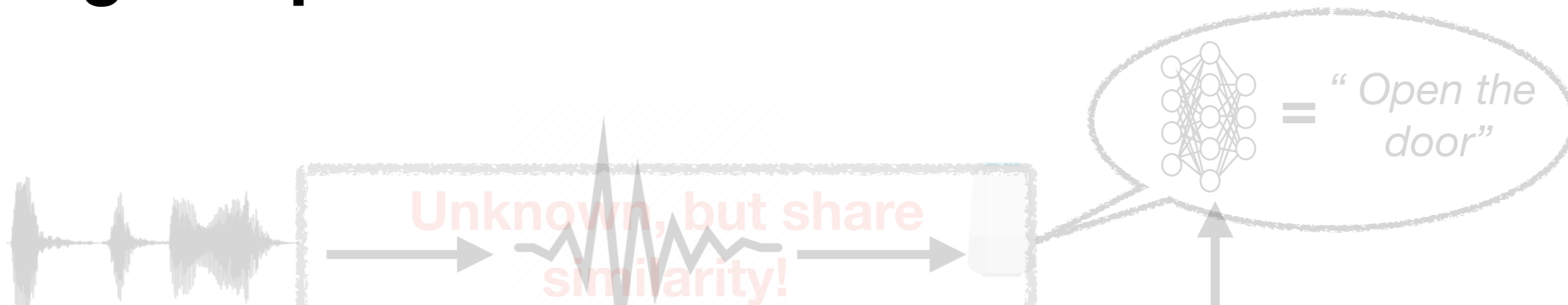


Design Inspiration

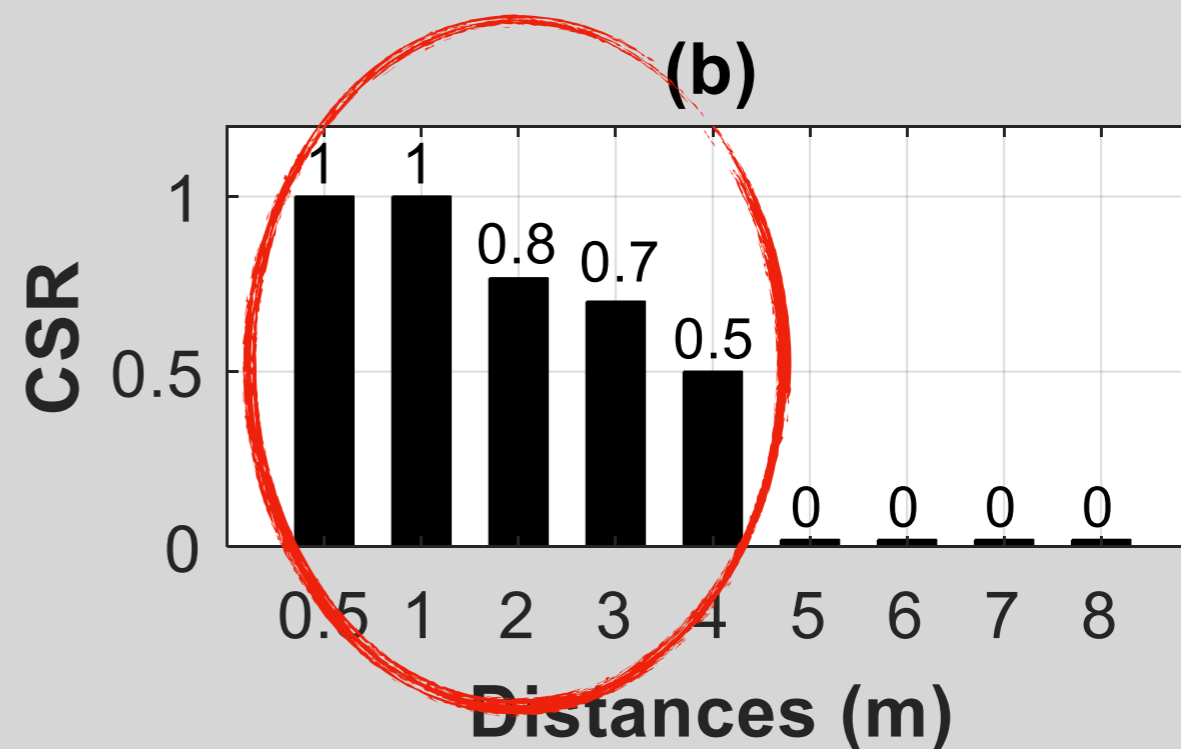
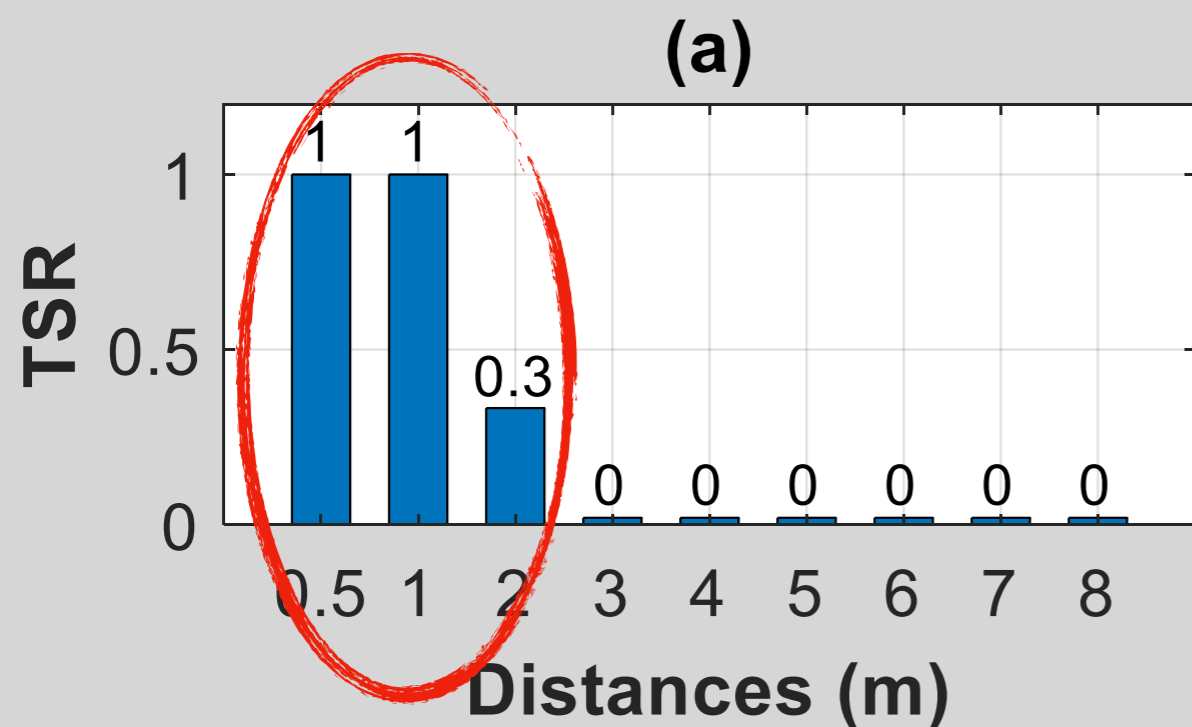


$$\arg \min_{\delta} \alpha \cdot dB_I(\delta) + \frac{1}{M} \sum_i Loss(SR(H_i(I + \delta)), T')$$

Design Inspiration

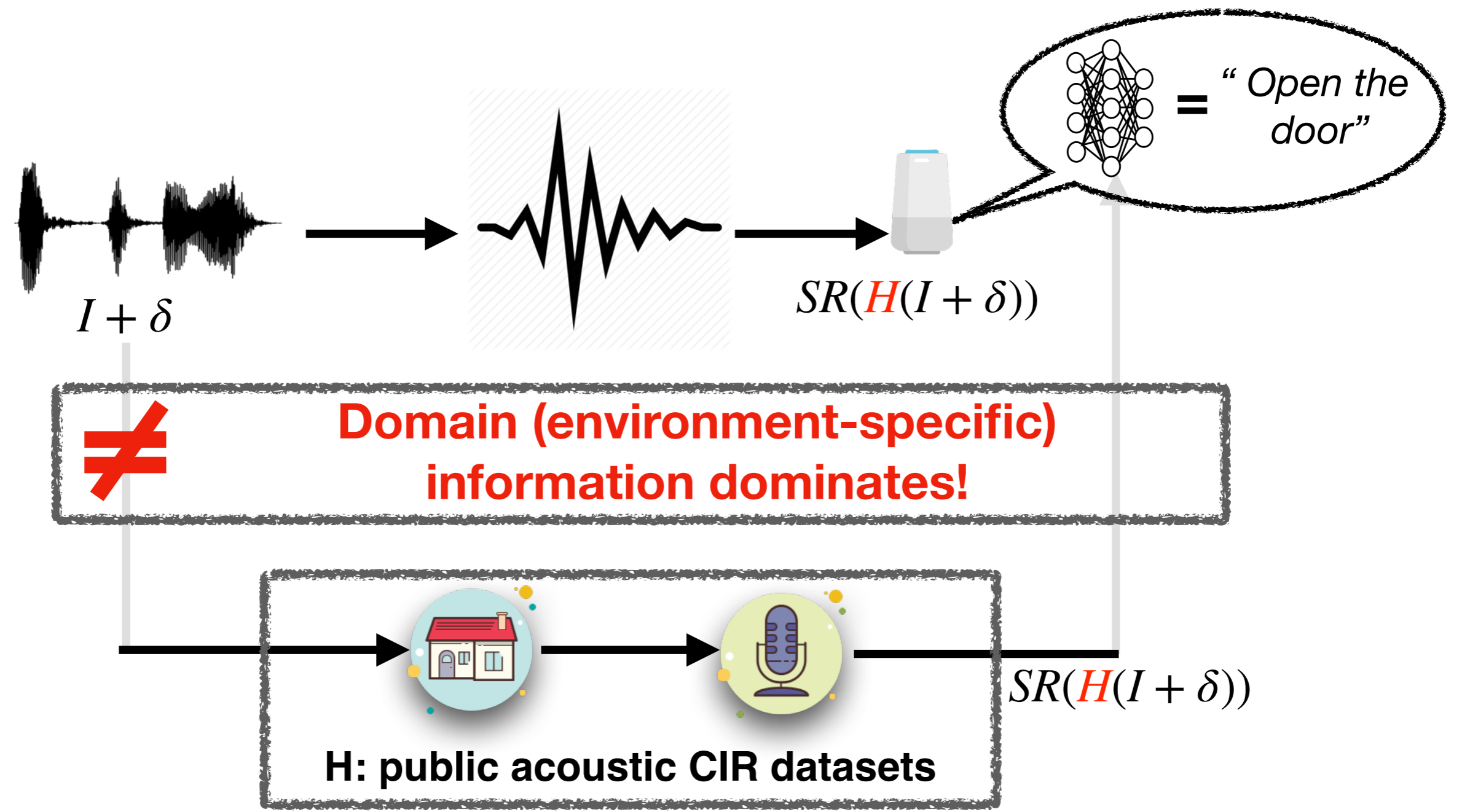


Transcript and Character Successful Rate:



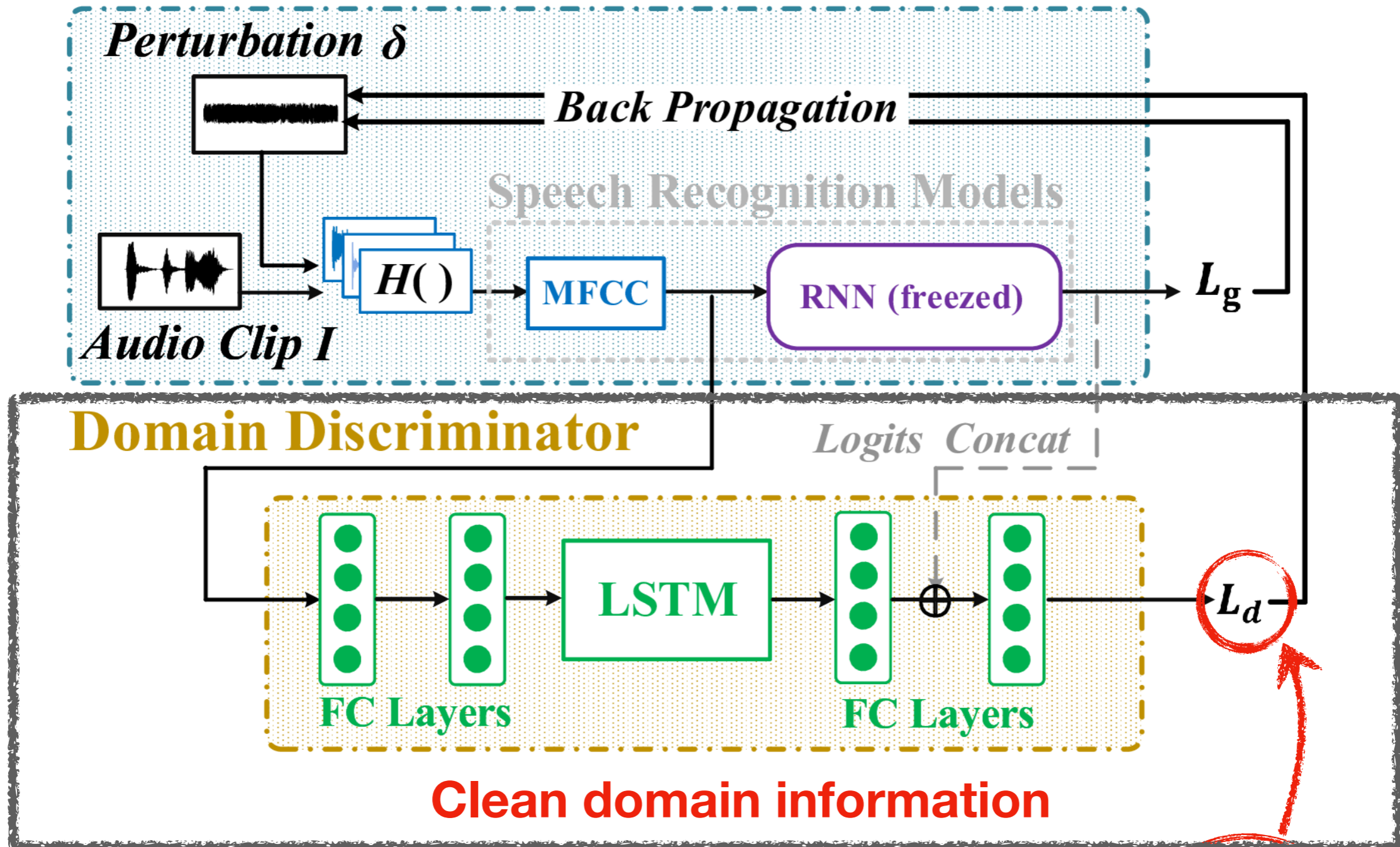
$$\arg \min_{\delta} \alpha \cdot dB_I(\delta) + \frac{1}{M} \sum_i \text{Loss}(SR(H_i(I + \delta)), T')$$

Design Inspiration



Metamorph: Meta-Enha


Adversarial Example Generator



$$\arg \min_{\delta} \alpha \cdot dB_I(\delta) + \frac{1}{M} \sum_i \text{Loss}(SR(H_i(I + \delta)), T') - \beta \cdot L_d$$

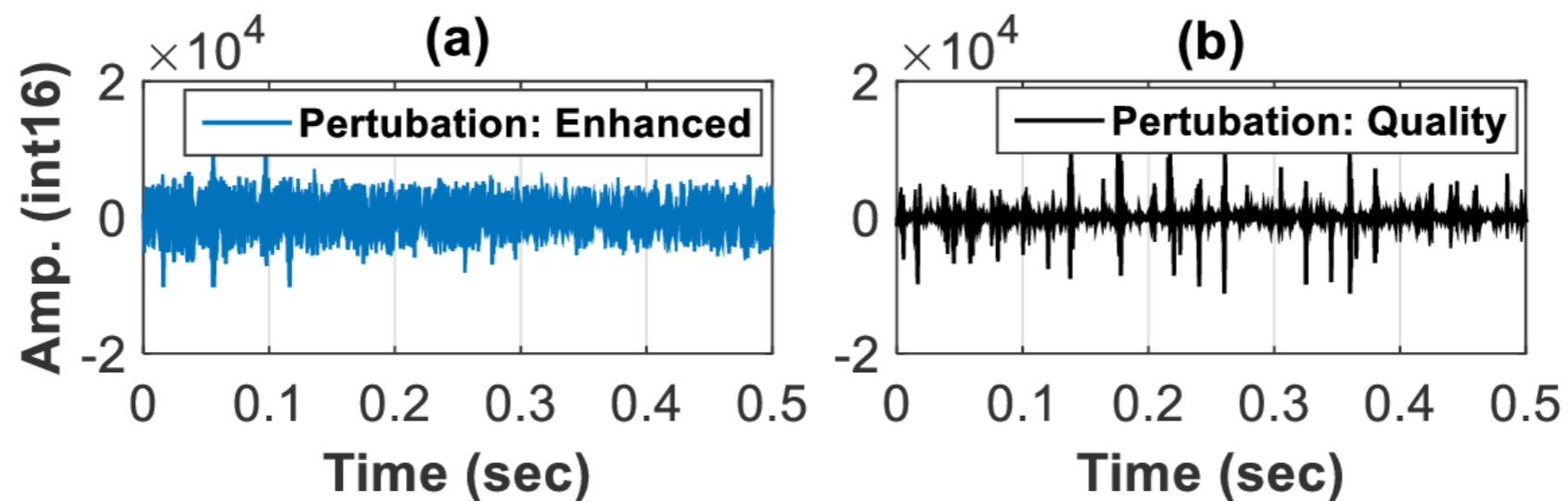
Metamorph: Meta-Qual

- Acoustic Graffiti:

$$\text{distance}(\delta, \hat{N})$$


- Reducing Perturbation's Coverage:

L1/L2 regularization



Evaluation: Audio Quality

- Examples

Classical music

Original:
[no transcription]

Meta-Enha:
“hello world”

Meta-Qual:
“hello world”

Human speech

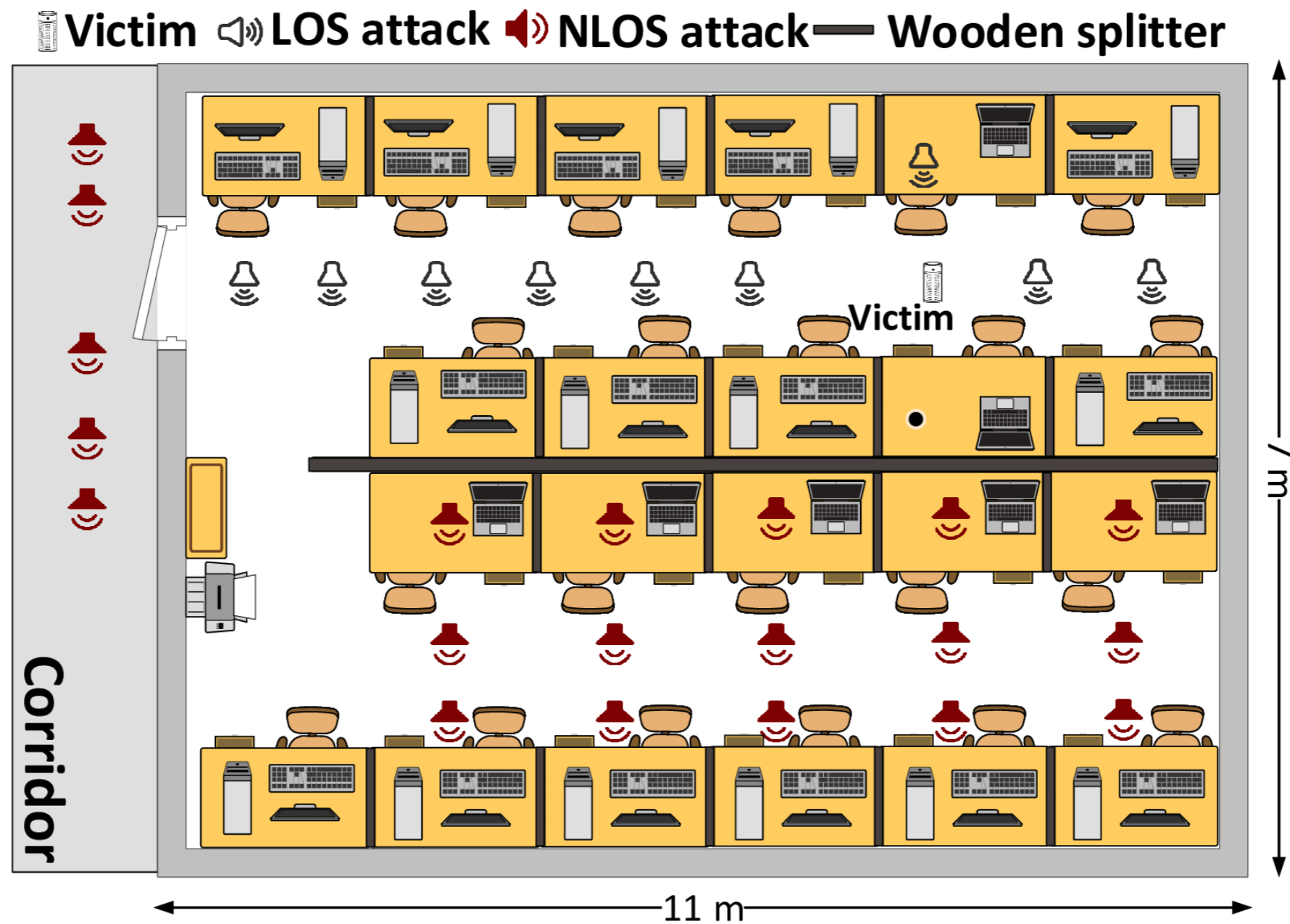
Original:
“your son went to
serve at a distant
place and became
a centurion”

Meta-Enha:
“open the door”

Meta-Qual:
“open the door”

Evaluation: Attack Successful Rate

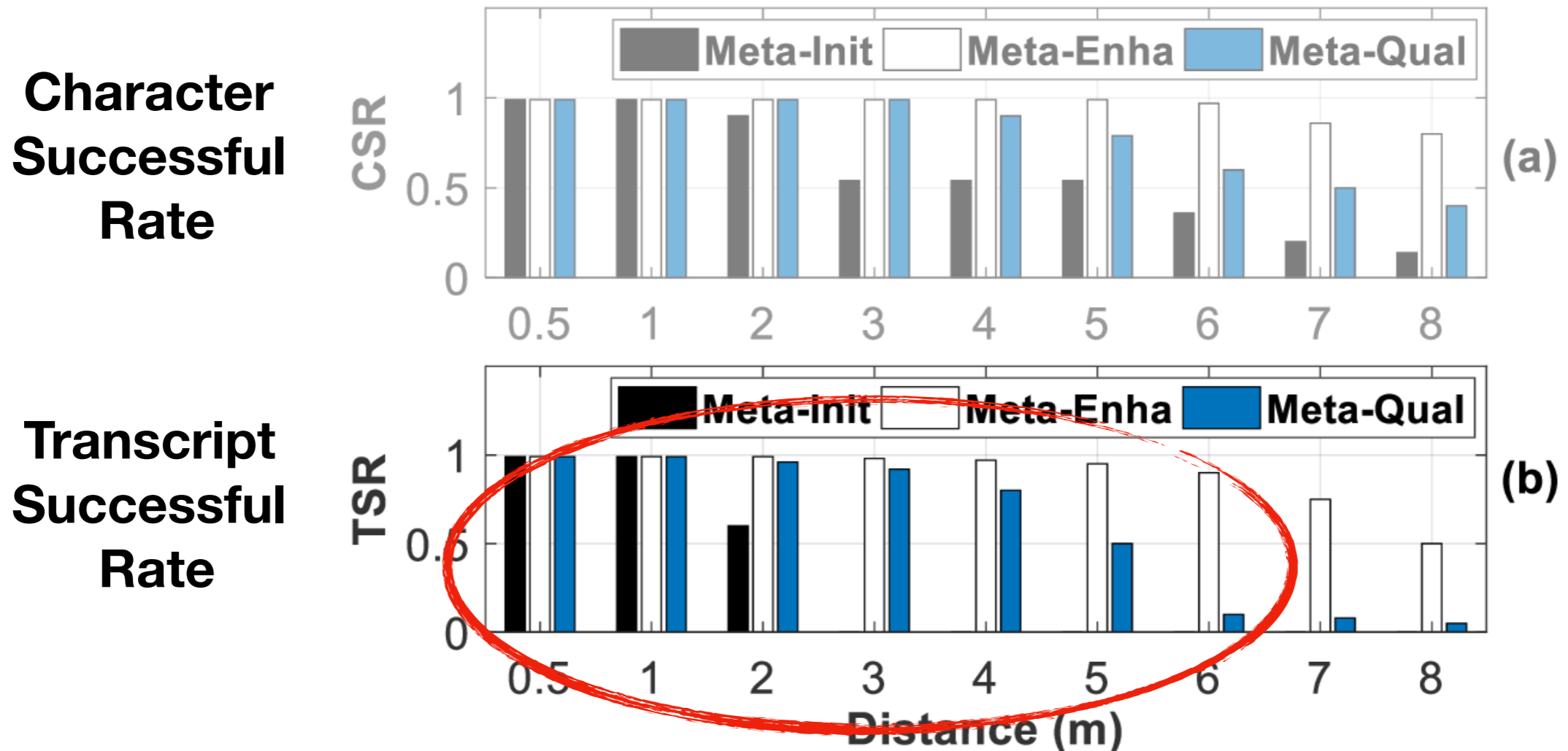
- Attack Target: “DeepSpeech” (White-Box)



A multi-path prevalent office

Evaluation: Attack Successful Rate

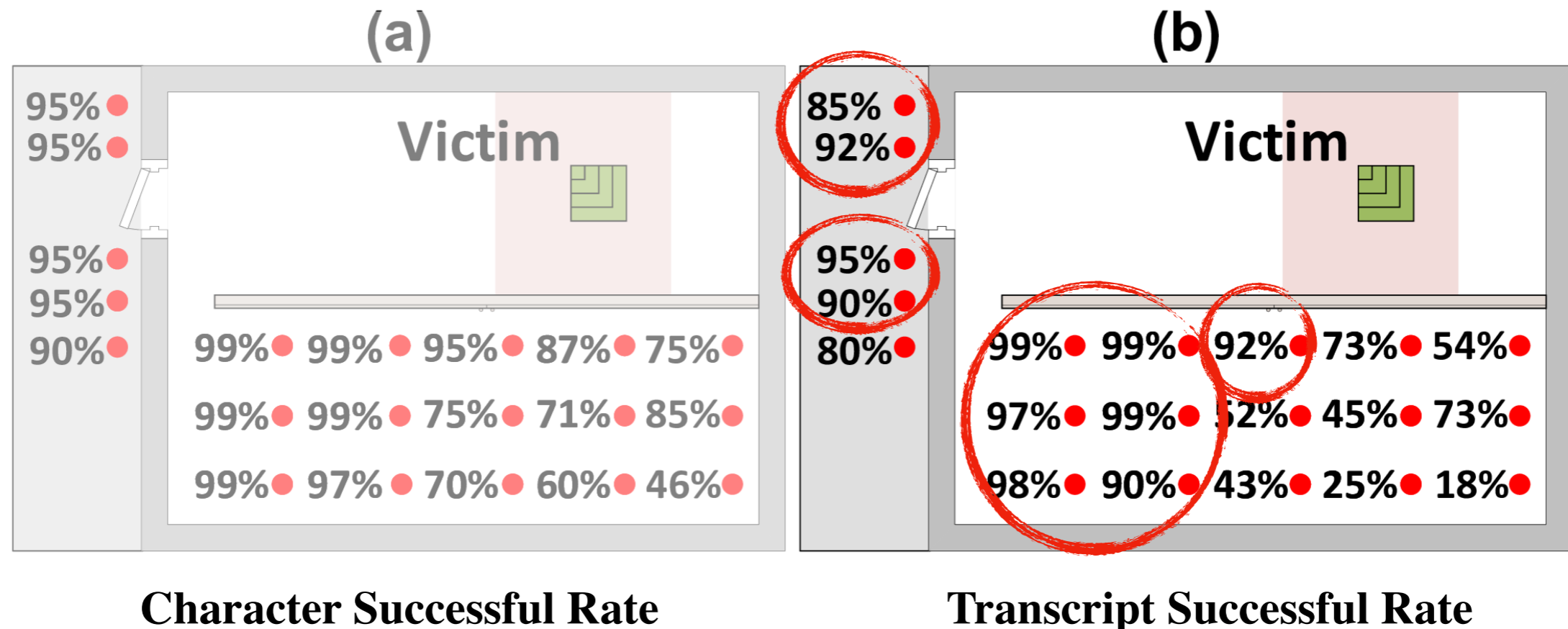
- Line-of-Sight (LOS) Attack



Meta-Enha: > 90% attack successful rate

Evaluation: Attack Successful Rate

- No-Line-of-Sight (NLOS) Attack



**Meta-Enha: over 85% attack successful rate
across 11/20 NLOS location!**

Conclusion

1. Investigate over-the-air audio adversarial attacks systematically.
2. Propose a “generate-and-clean” two-phase design and improve the audio quality.
3. Develop a prototype and conduct extensive evaluations.

**Visit [acoustic-metamorph-system.github.io](https://github.com/acoustic-metamorph-system)
for more information!**