

TRAMBA: A Hybrid Transformer and Mamba Architecture for Practical Audio and Bone Conduction Speech Super Resolution and Enhancement for Mobile and Wearable Platforms

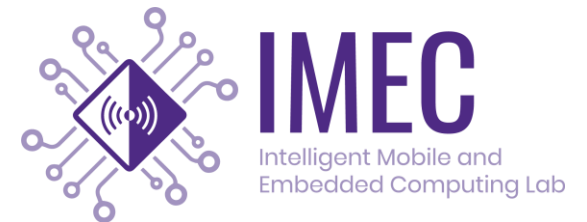
October 14, 2025

UbiComp/ISWC 2025

Yueyuan Sui, Minghui Zhao, Junxi Xia, Xiaofan Jiang, **Stephen Xia**

Department of Electrical and Computer Engineering

Northwestern University



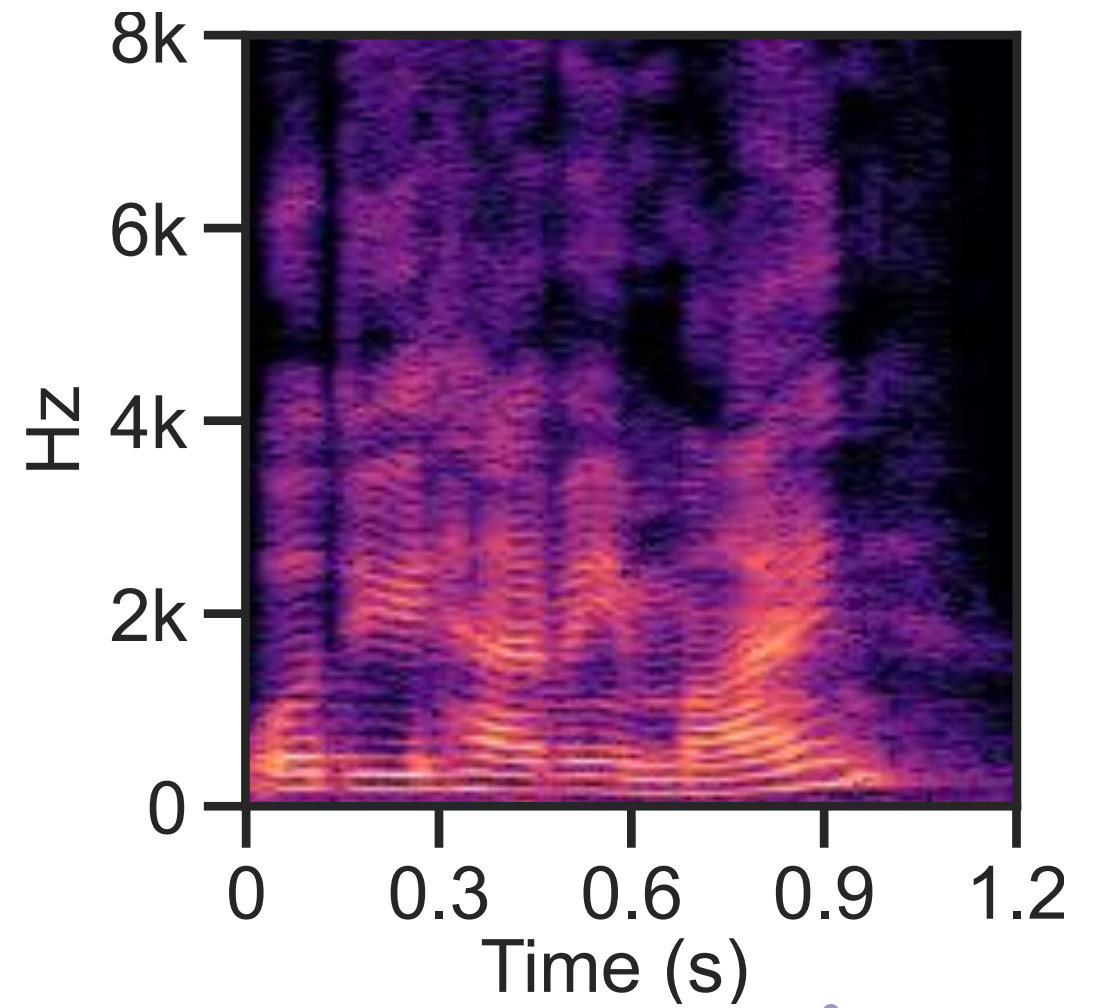
Bone and Vibration-based Microphones

- Other sounds
- Speech is attenuated

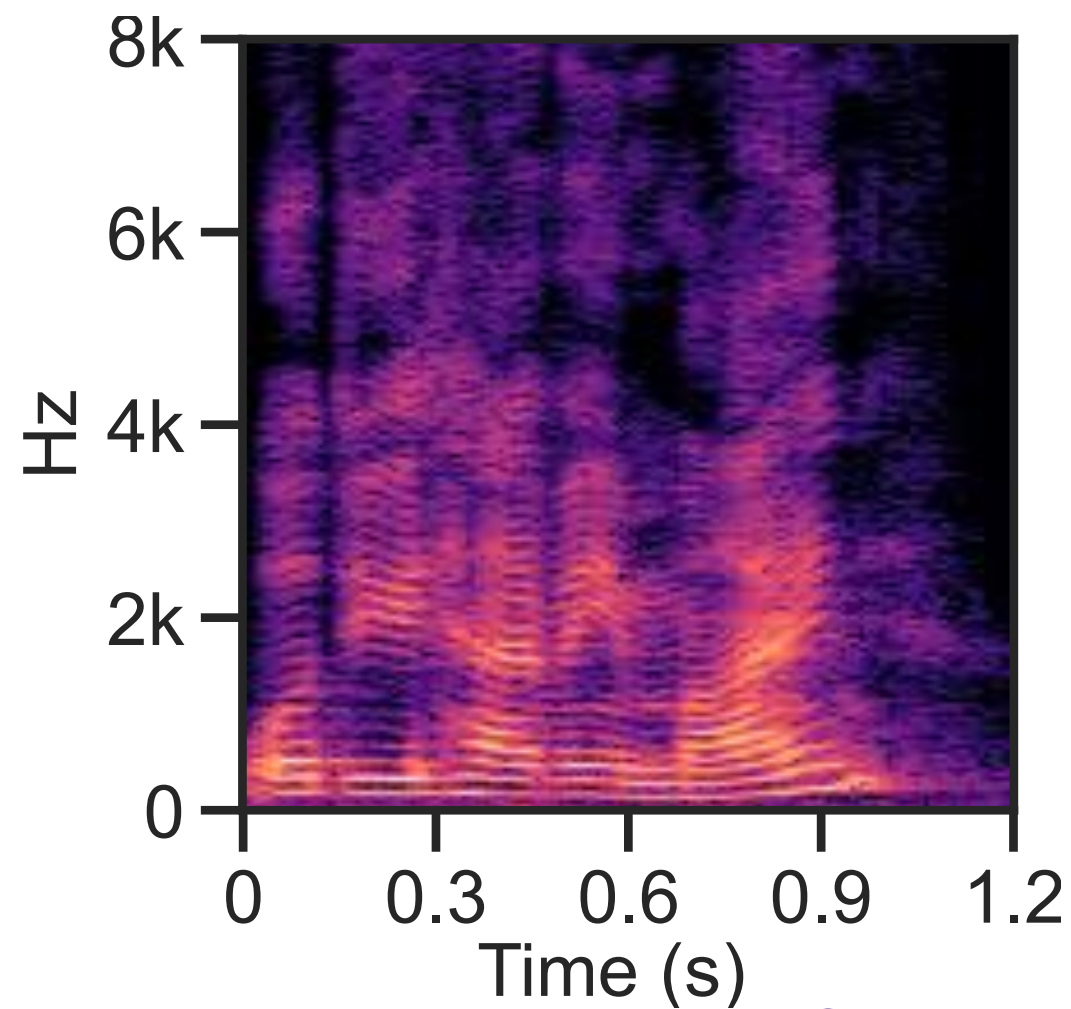
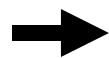
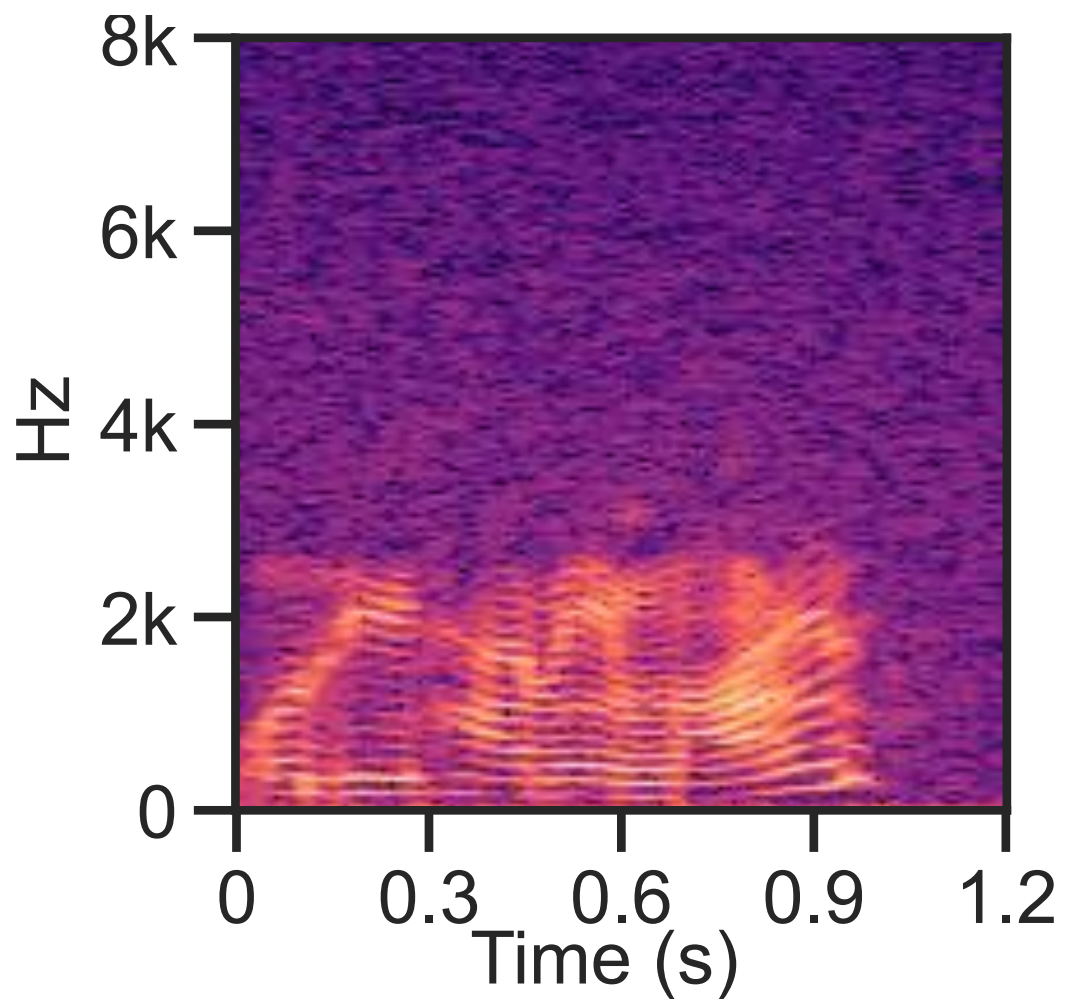
Bone and Vibration-based Microphones

- Other sounds
 - Mitigated with vibration-based sensing (BCM or IMU)
- **Speech is attenuated**

Standard Speech

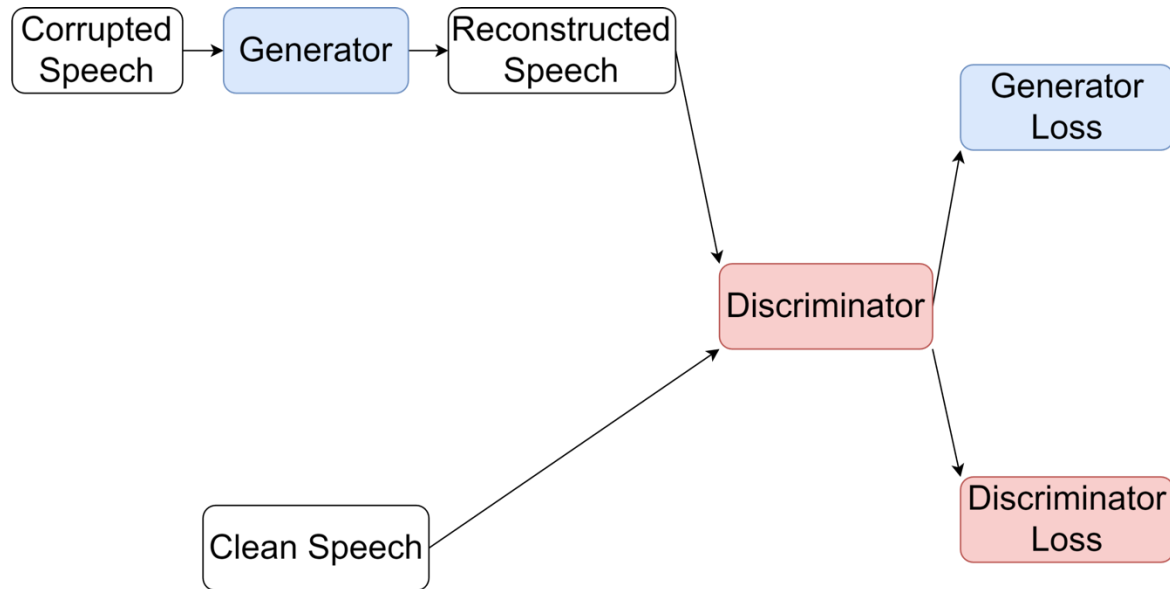


Corrupted Speech



Current Solutions:

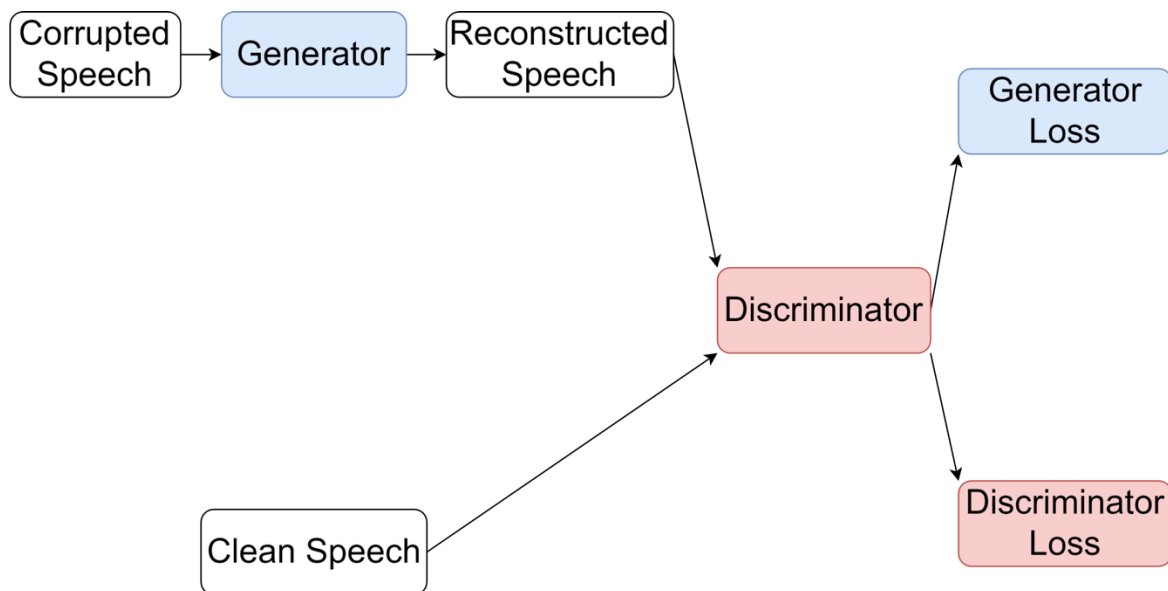
Generative Adversarial Networks



Performance: good, Compute: heavy

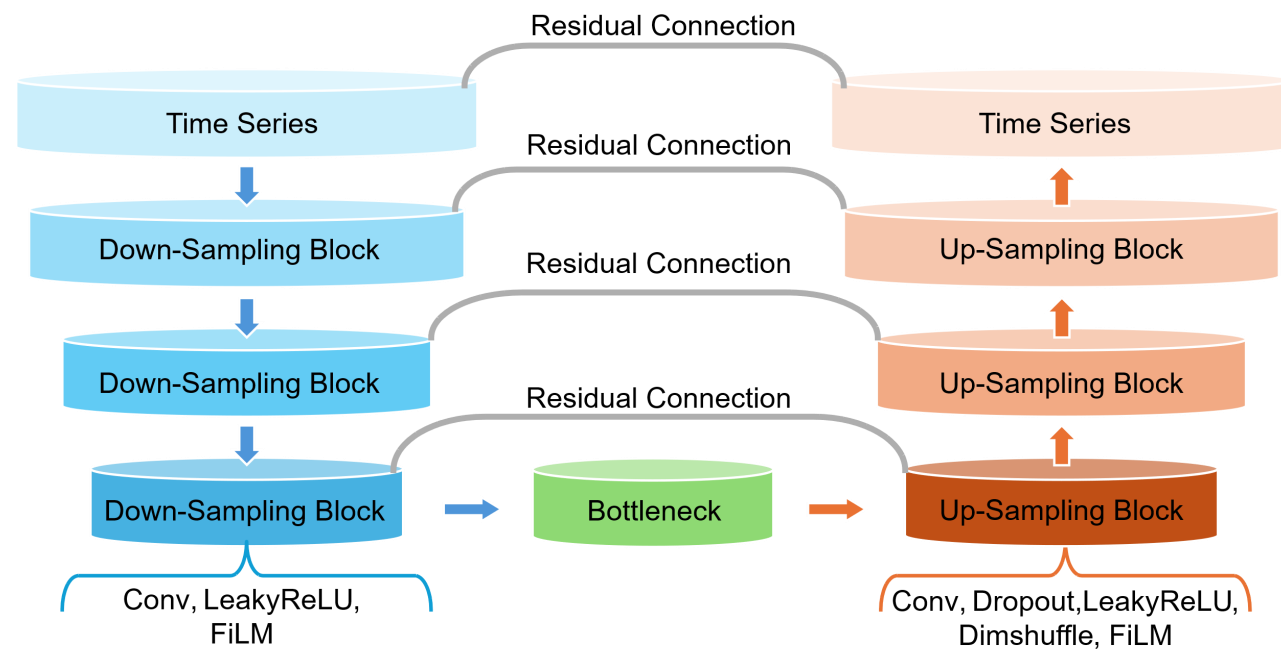
Current Solutions:

Generative Adversarial Networks



Performance: good, Compute: heavy

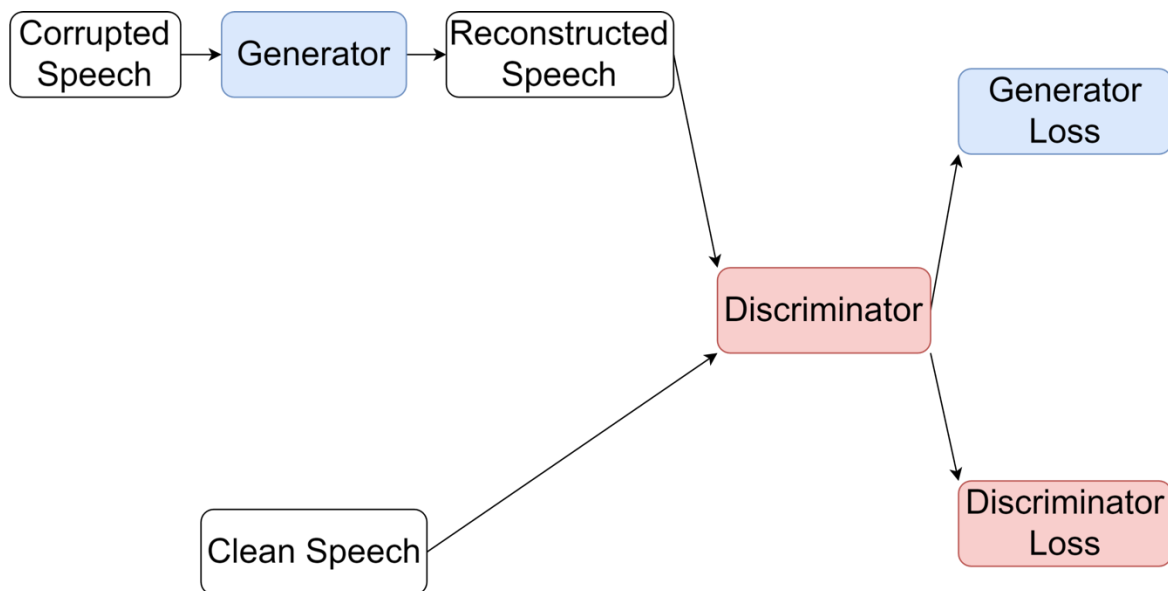
U-Net



Performance: less good, Compute: light

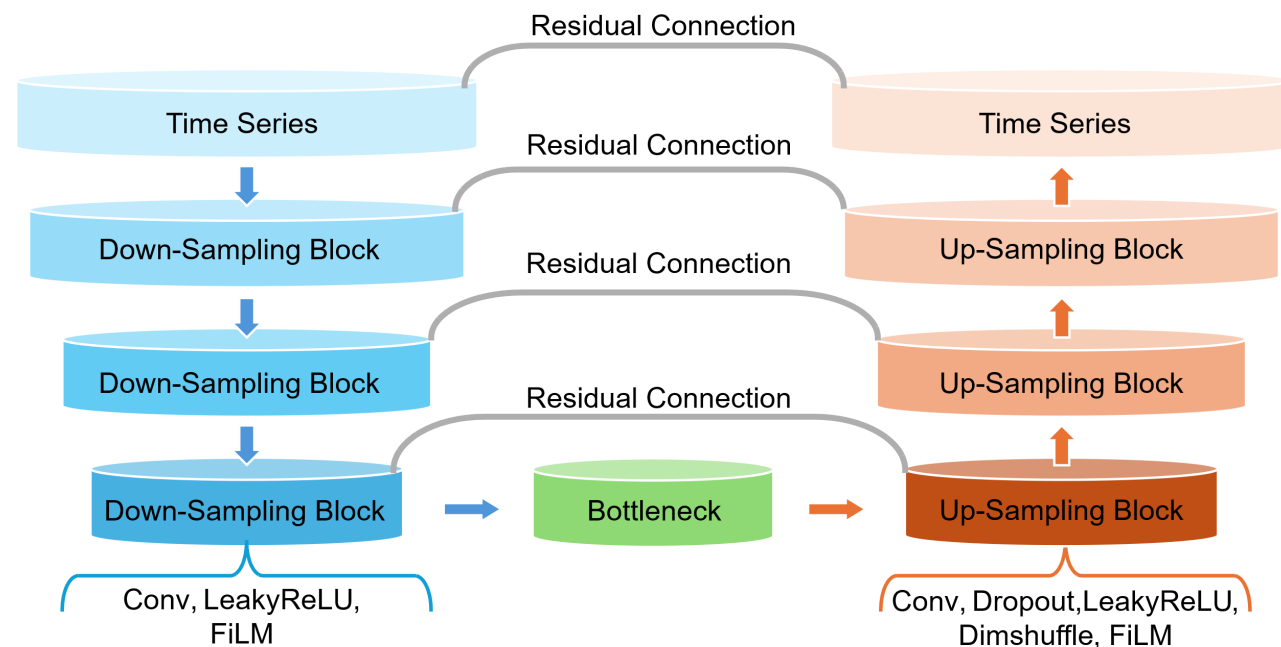
Current Solutions:

Generative Adversarial Networks



Performance: good, Compute: heavy

U-Net



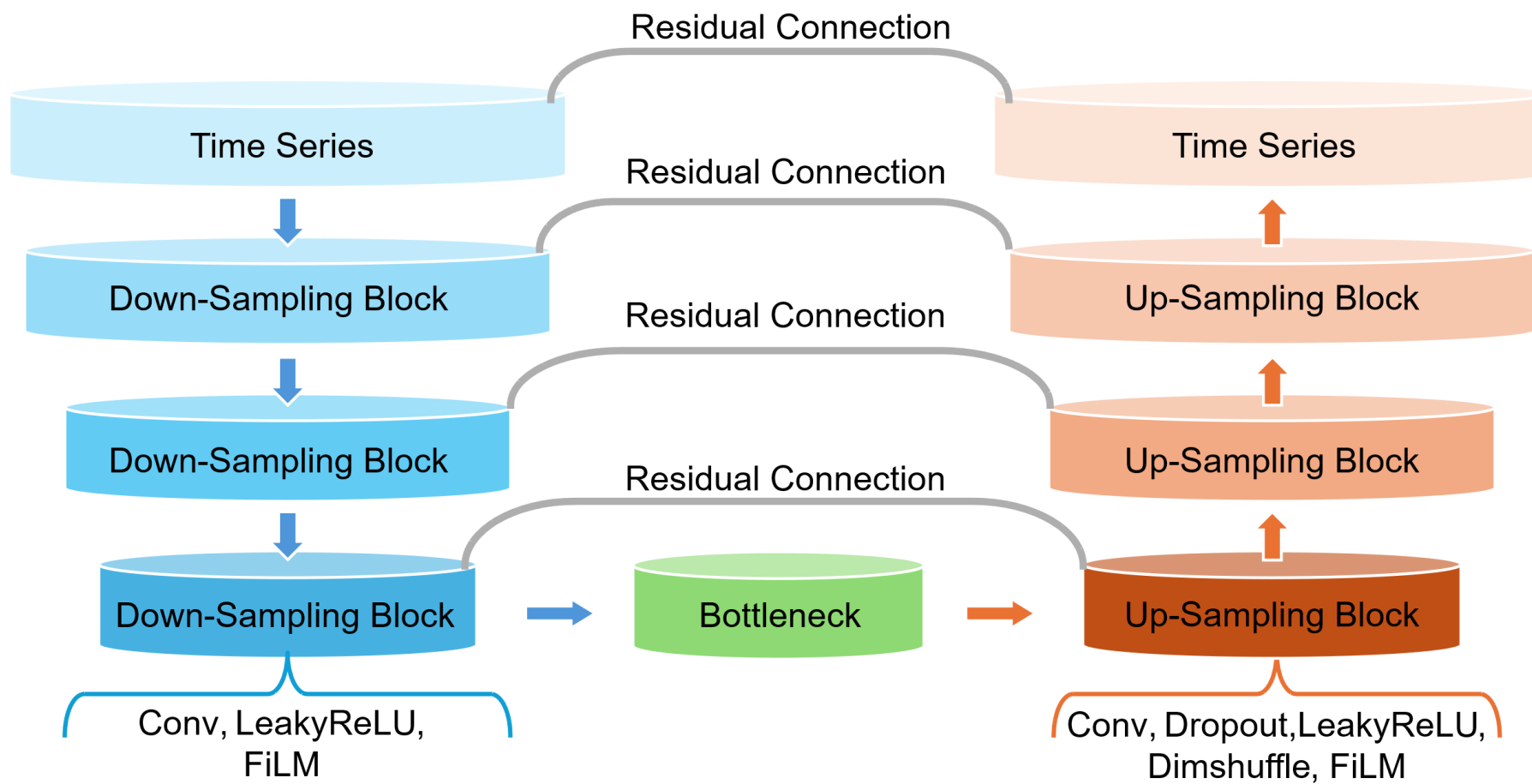
Performance: less good, Compute: light

Goal: bridge performance, speed, and efficiency

TRAMBA: good, compute: light

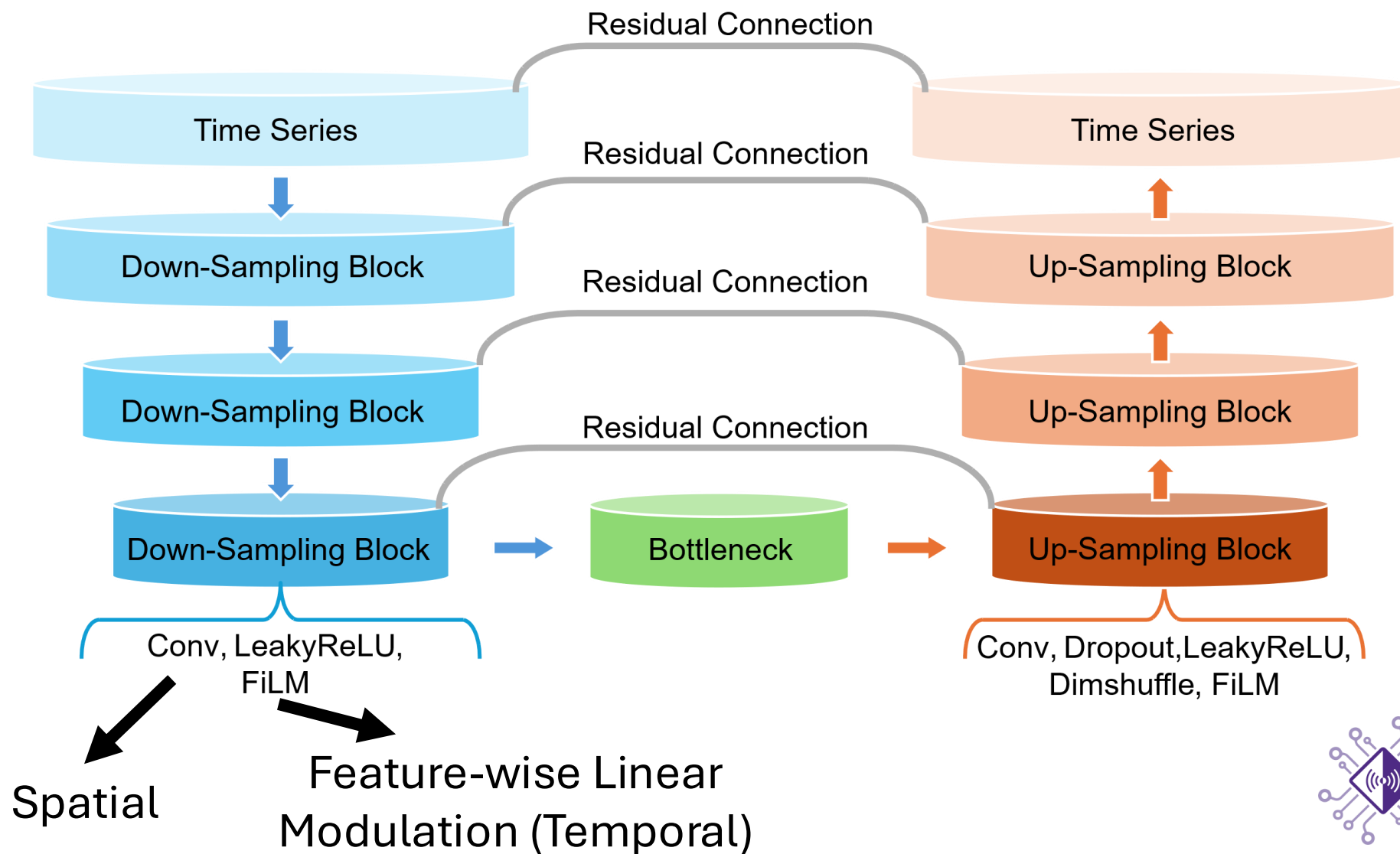
Design

Traditional Speech Enhancement (U-Net)



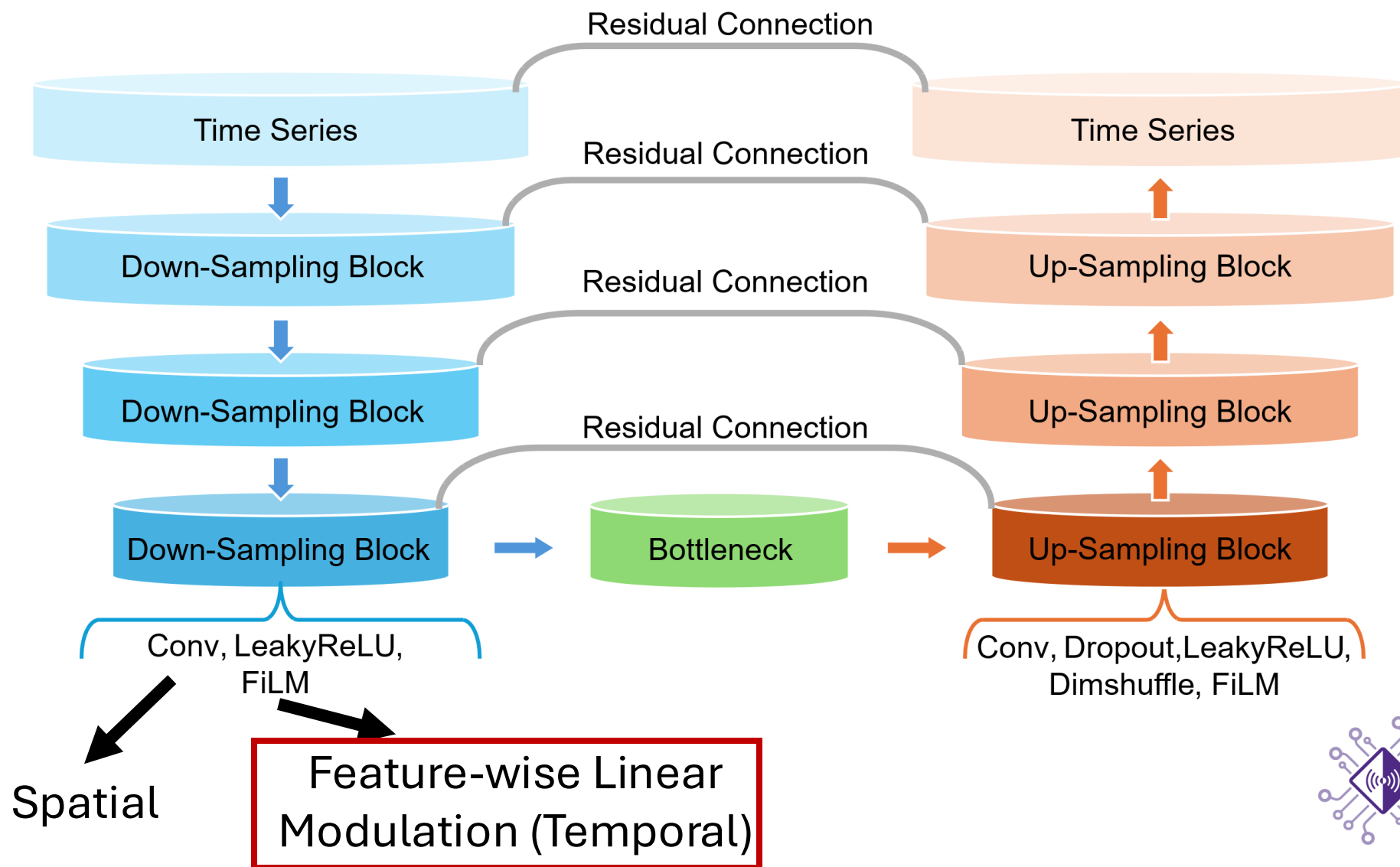
Design

Traditional Speech Enhancement (U-Net)

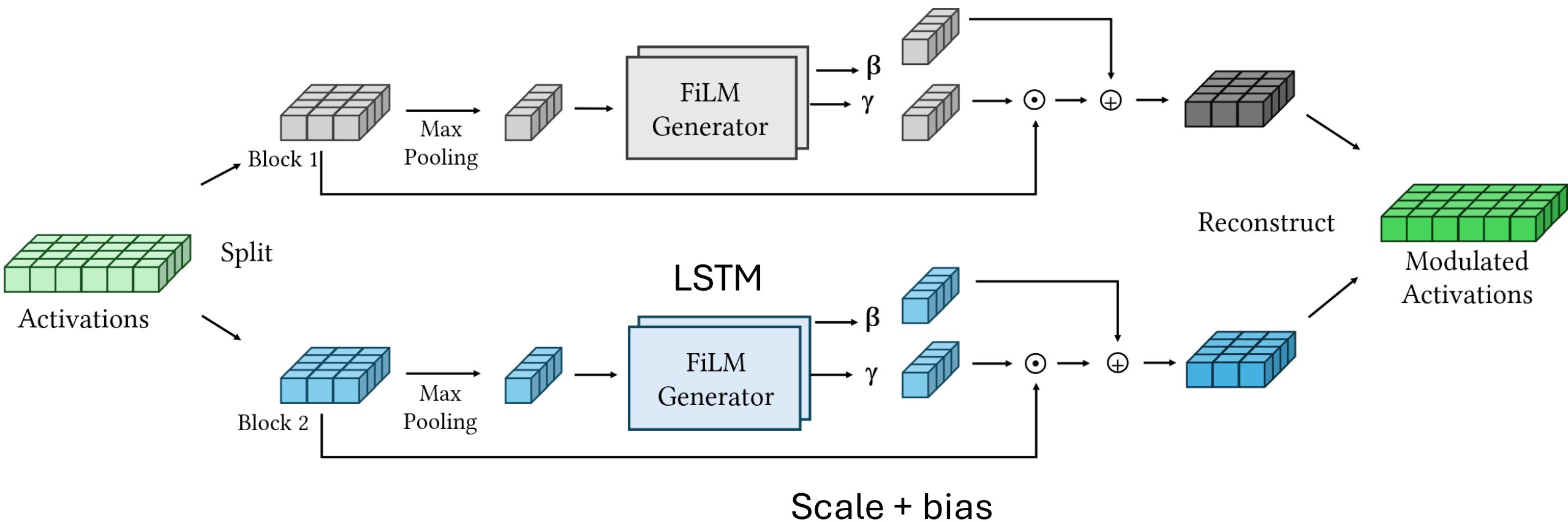


Design

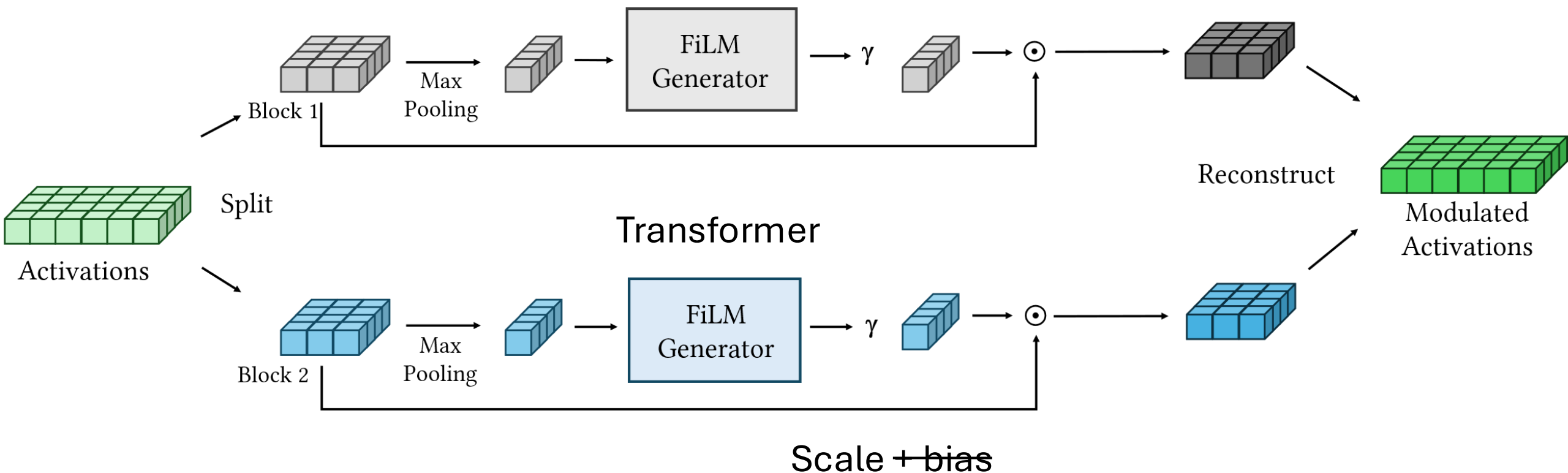
Traditional Speech Enhancement (U-Net)



Temporal Feature-wise Linear Modulation (TFiLM)

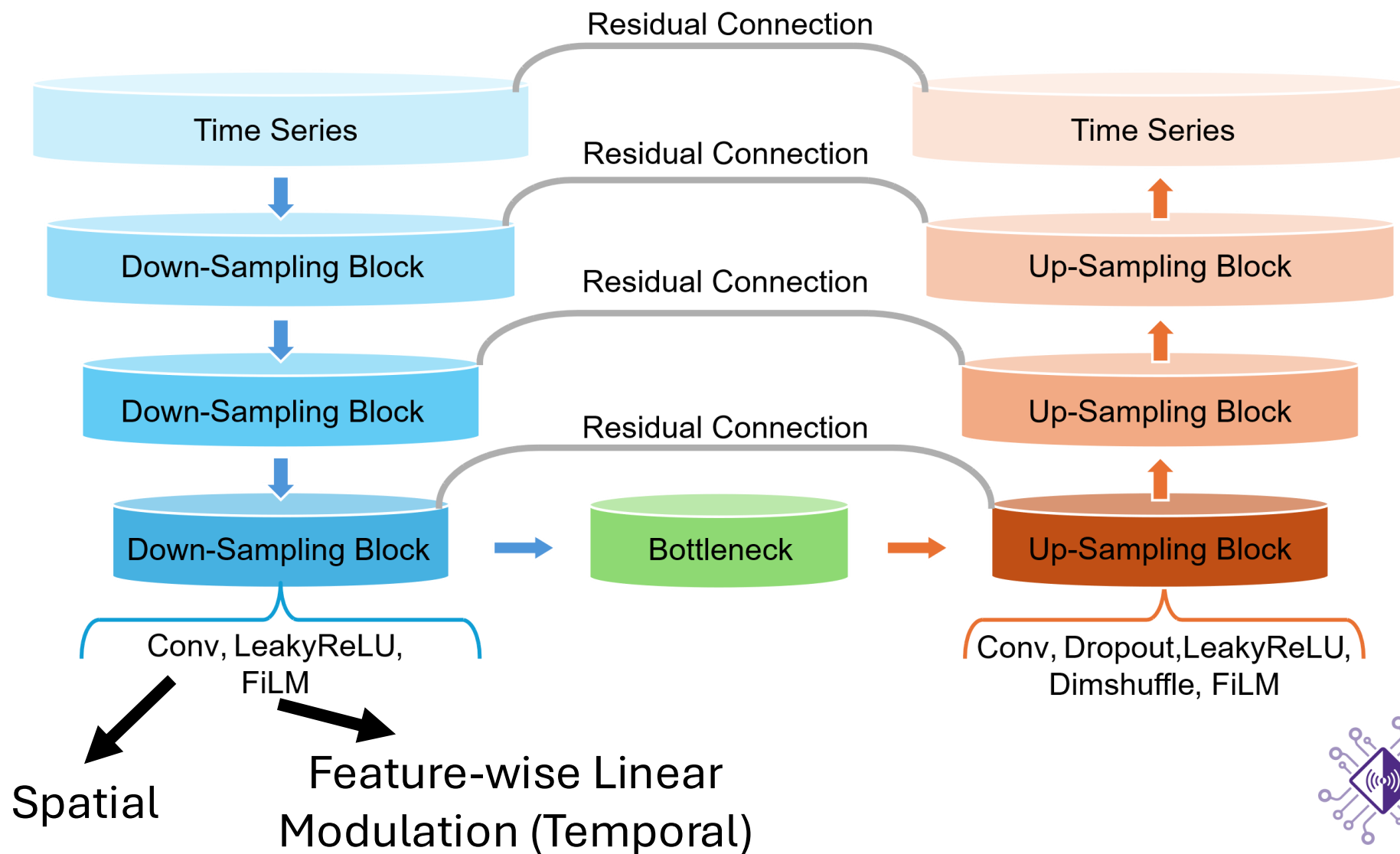


Scale-Only Attention-based Feature-wise Linear Modulation (SAFiLM)



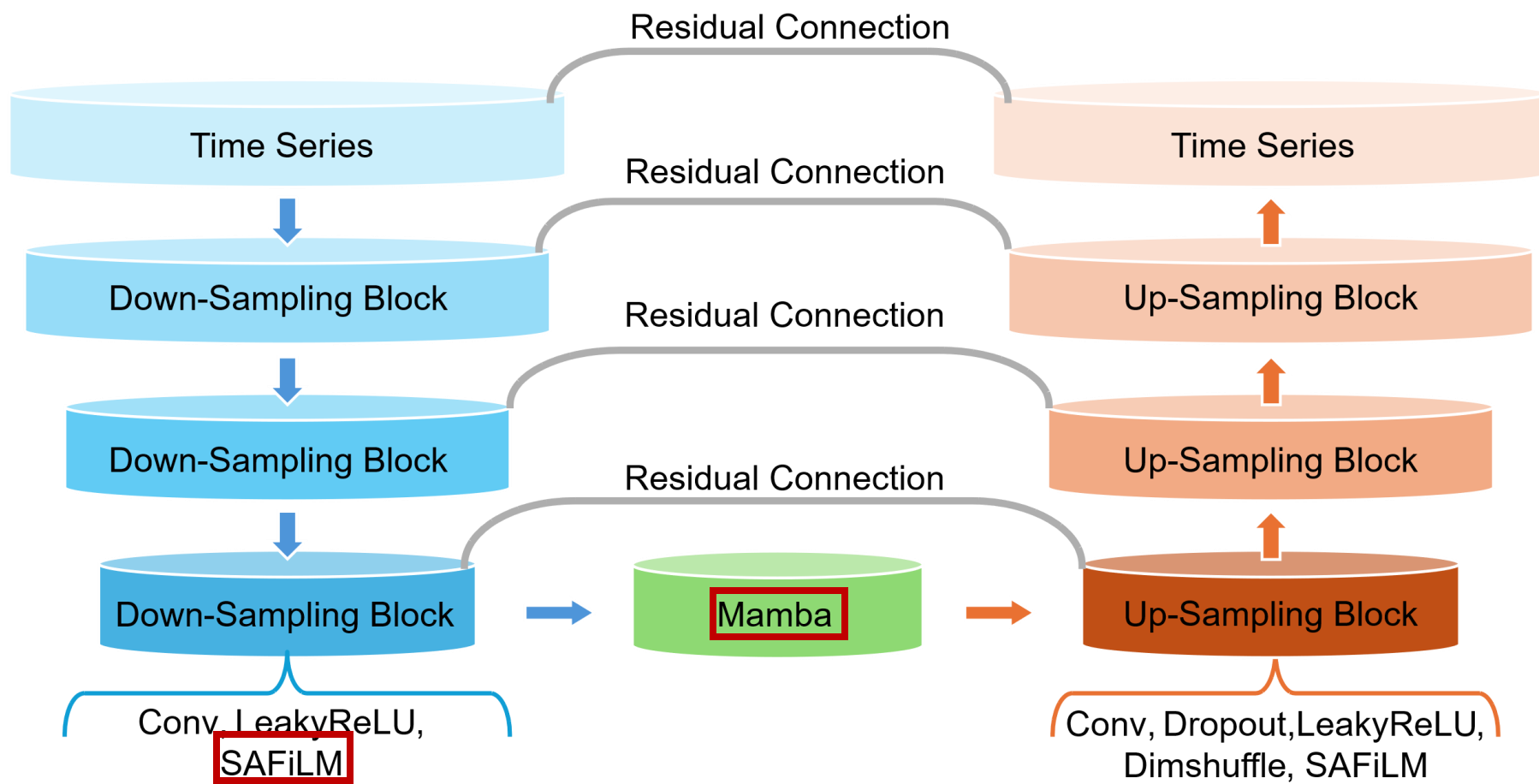
Design

Traditional Speech Enhancement (U-Net)

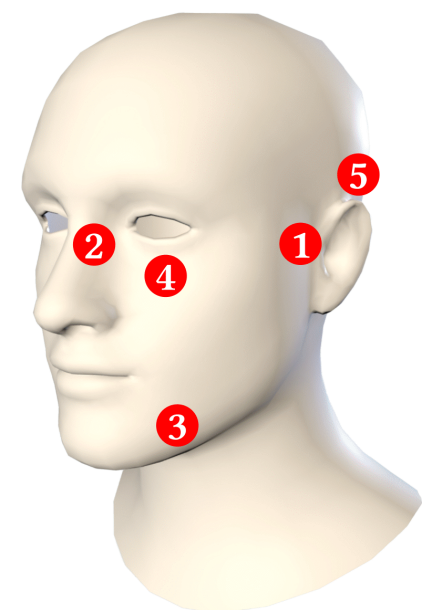


Design

Our Solution (TRAMBA)



Deployment and Evaluation



- 1 **Temporal Bone**
(Bone-conduction Headsets)
- 2 **Nasal Bone**
(Smart Glasses)
- 3 **Mandible**
(Face Masks)
- 4 **Zygomatic Bone**
(VR Headsets)
- 5 **Parietal Bone**
(Hats and Bone-conduction Headsets)

(a) Attachment positions tested



A. Quiet Office



B. Cafeteria

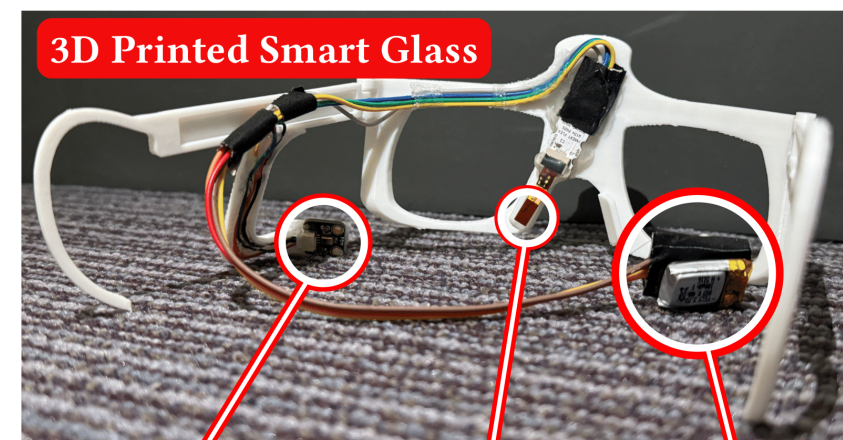


C. Loud Music



D. Construction

(b) Environments tested



3D Printed Smart Glass

OTA Mic.
Collects ground-truth voice for fine-tuning

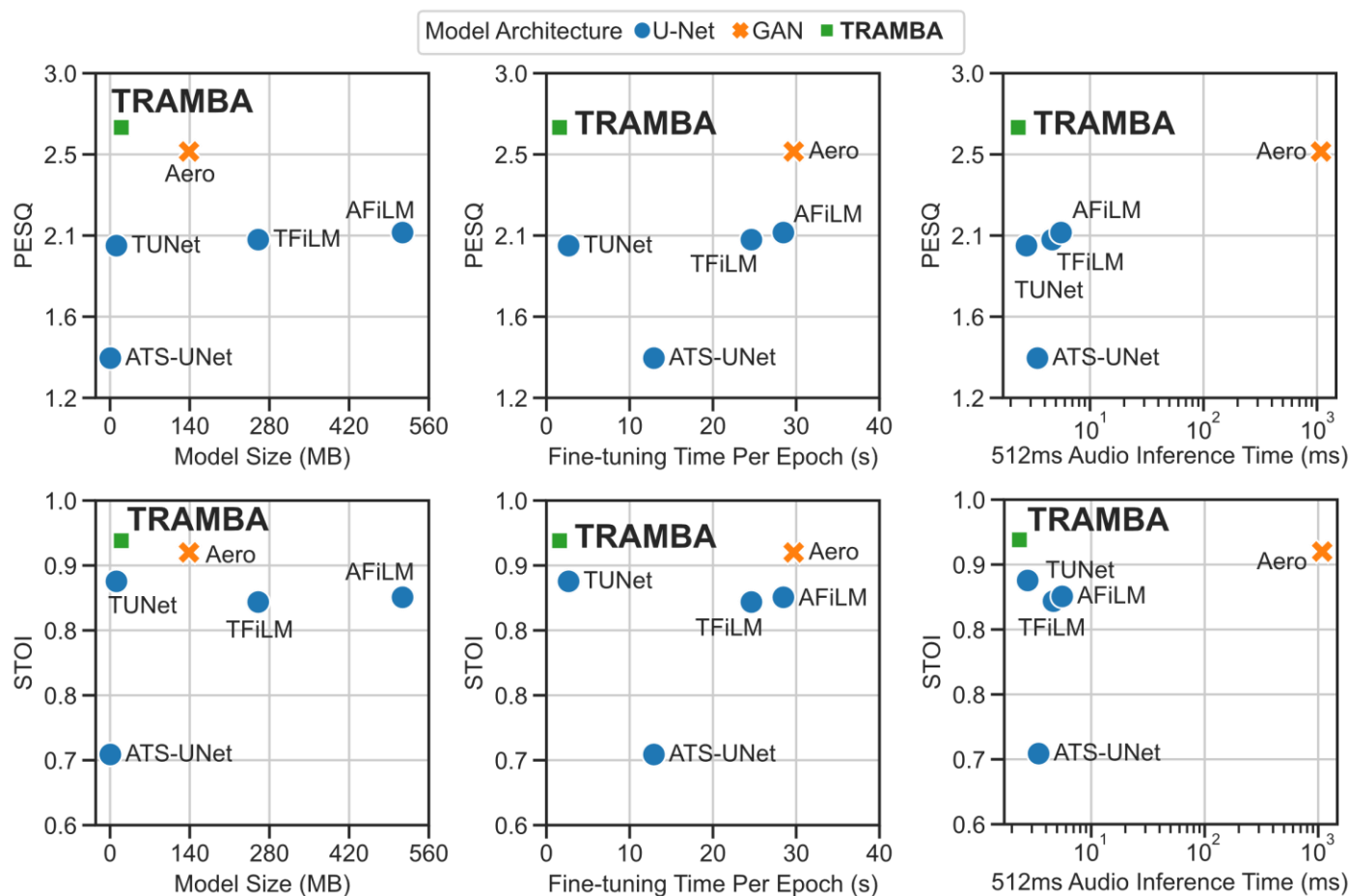
BCM/ACCEL
Captures env. noise-free voice data through vibration

nRF52840 BLE SoC
Streams data to smartphone

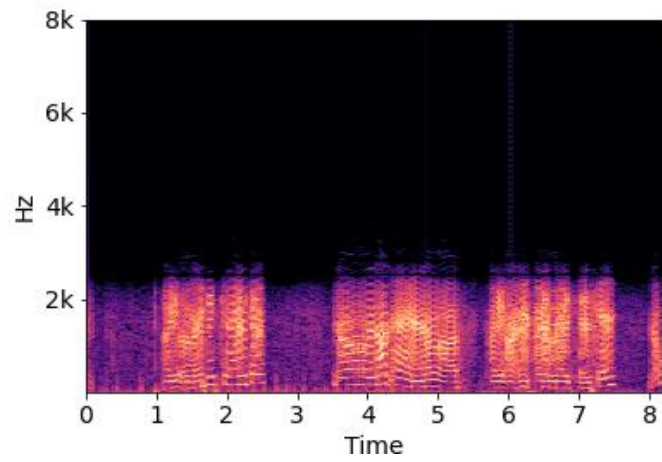
(c) Mobile-TRAMBA prototype

Overall Performance

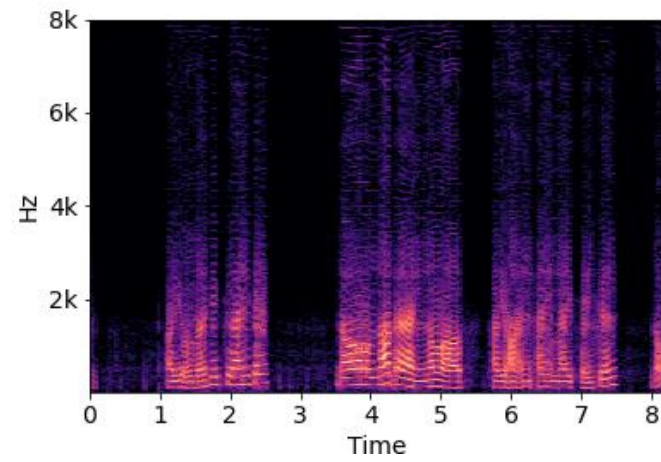
- Highest Perceptual Quality and Intelligibility
- Only ~20 MB
- Real-time
- Reduced power



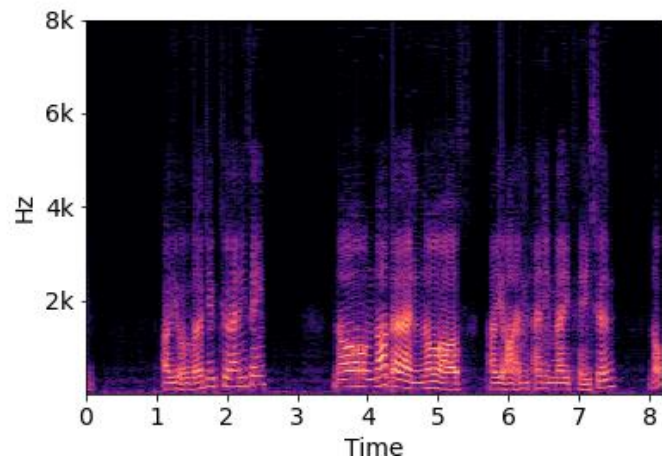
Demo




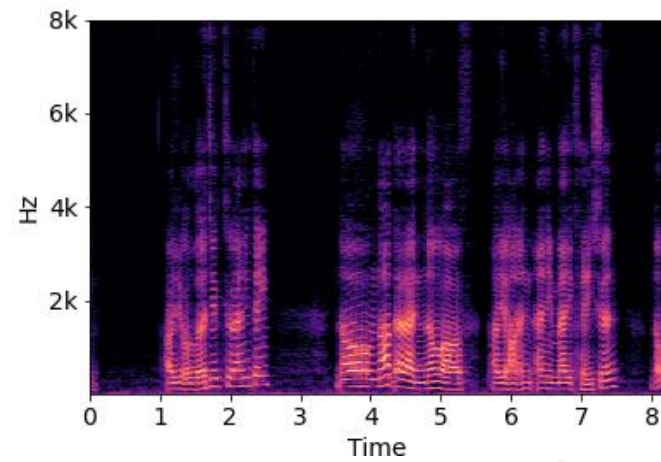
Raw BCM: 



TFILM: 



TRAMBA (Ours): 



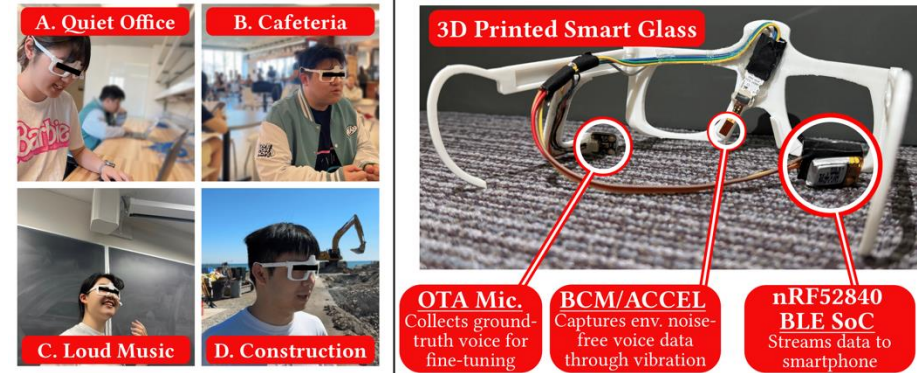
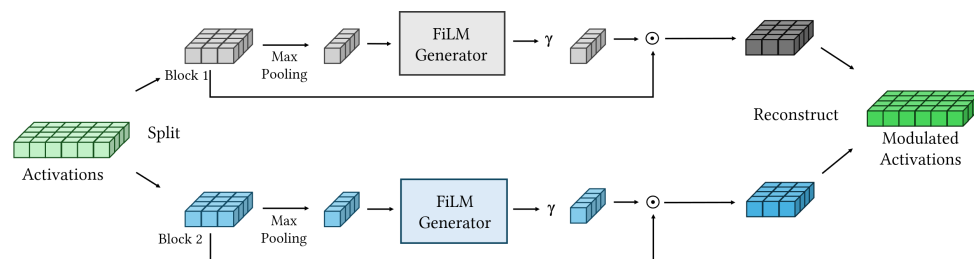
Ground Truth: 

TRAMBA: Practical Speech Enhancement for Mobile and Wearable Systems

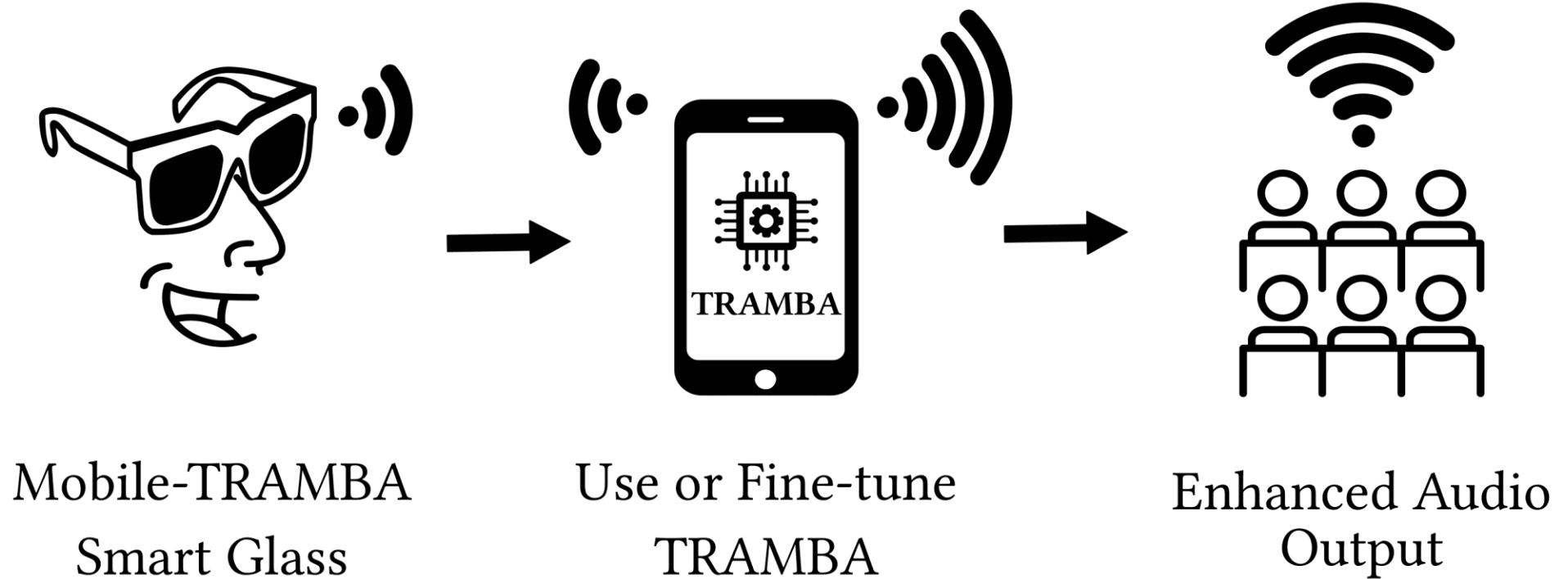
- Efficient Temporal Modeling: Scale-only Attention-based Feature-wise Linear Modulation (SAFiLM)
- Efficient Bottleneck: Mamba
- Bridges performance + efficiency

Feel free to reach out!

- Yueyuan Sui: yueyuansui2024@u.northwestern.edu
- Stephen Xia: stephen.xia@northwestern.edu

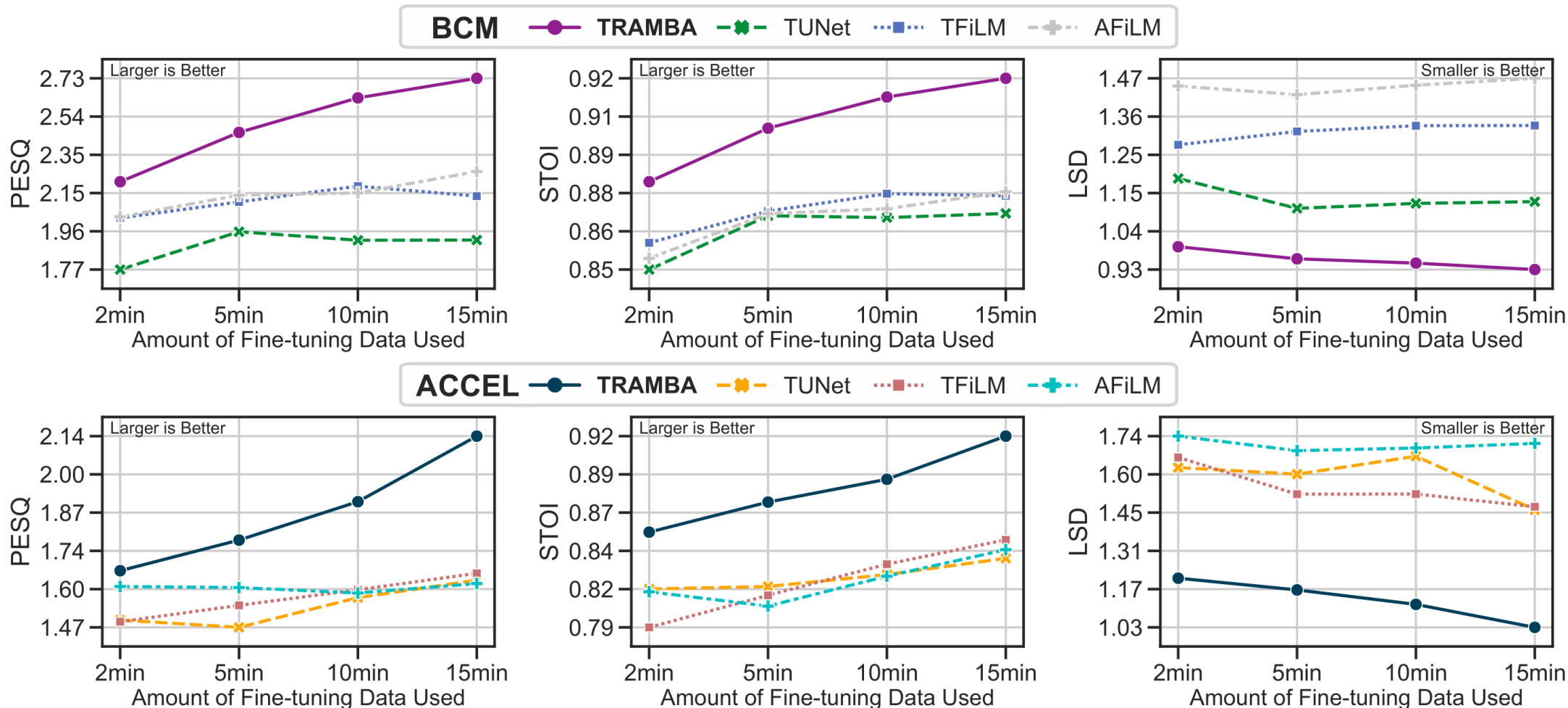


Training and Fine-tuning: Dealing with Lack of Data

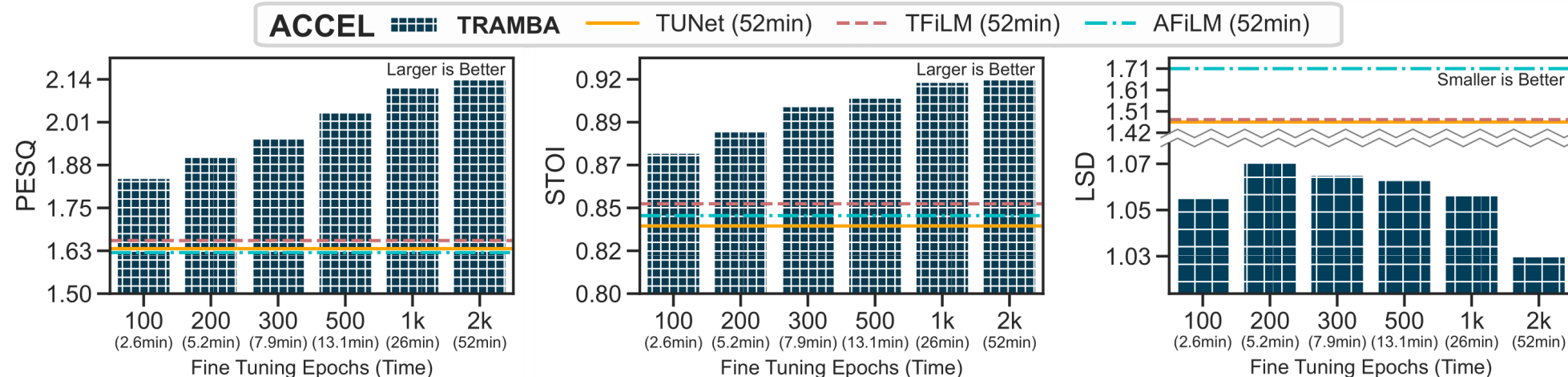
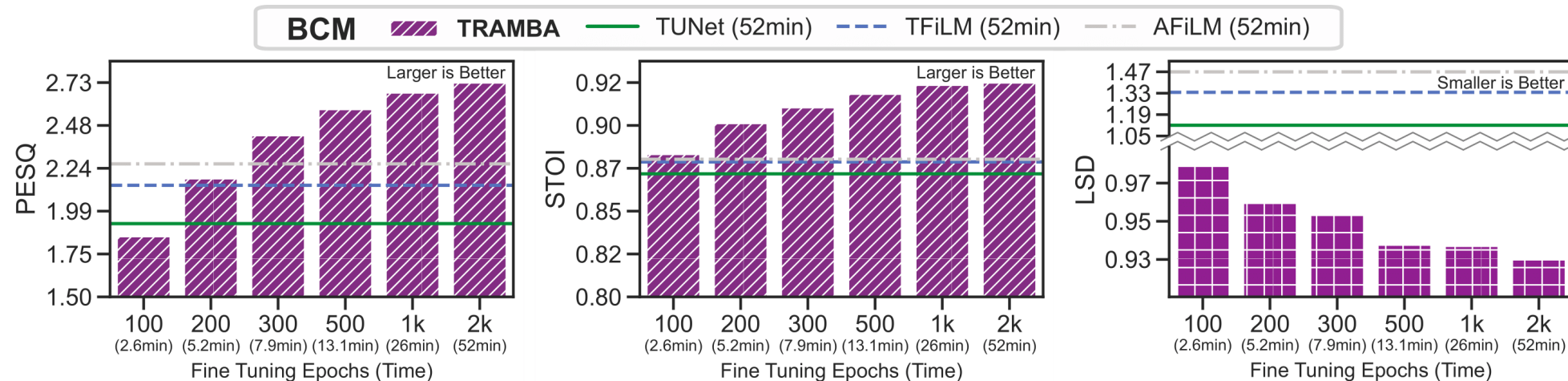


- Pretrain with Over-the-Air Audio
 - Subsample and Decimate
- Fine-tune with small amount of user speech

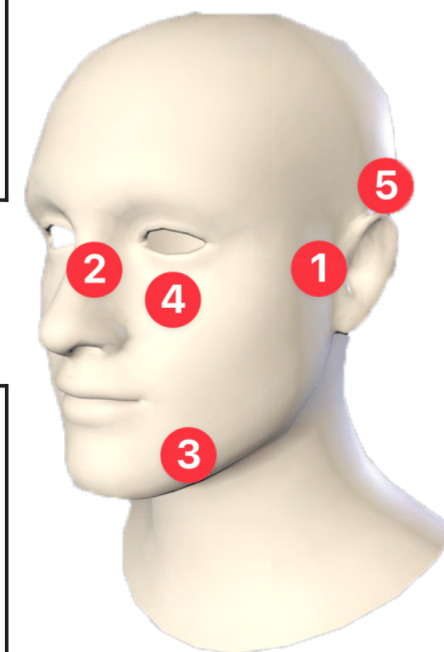
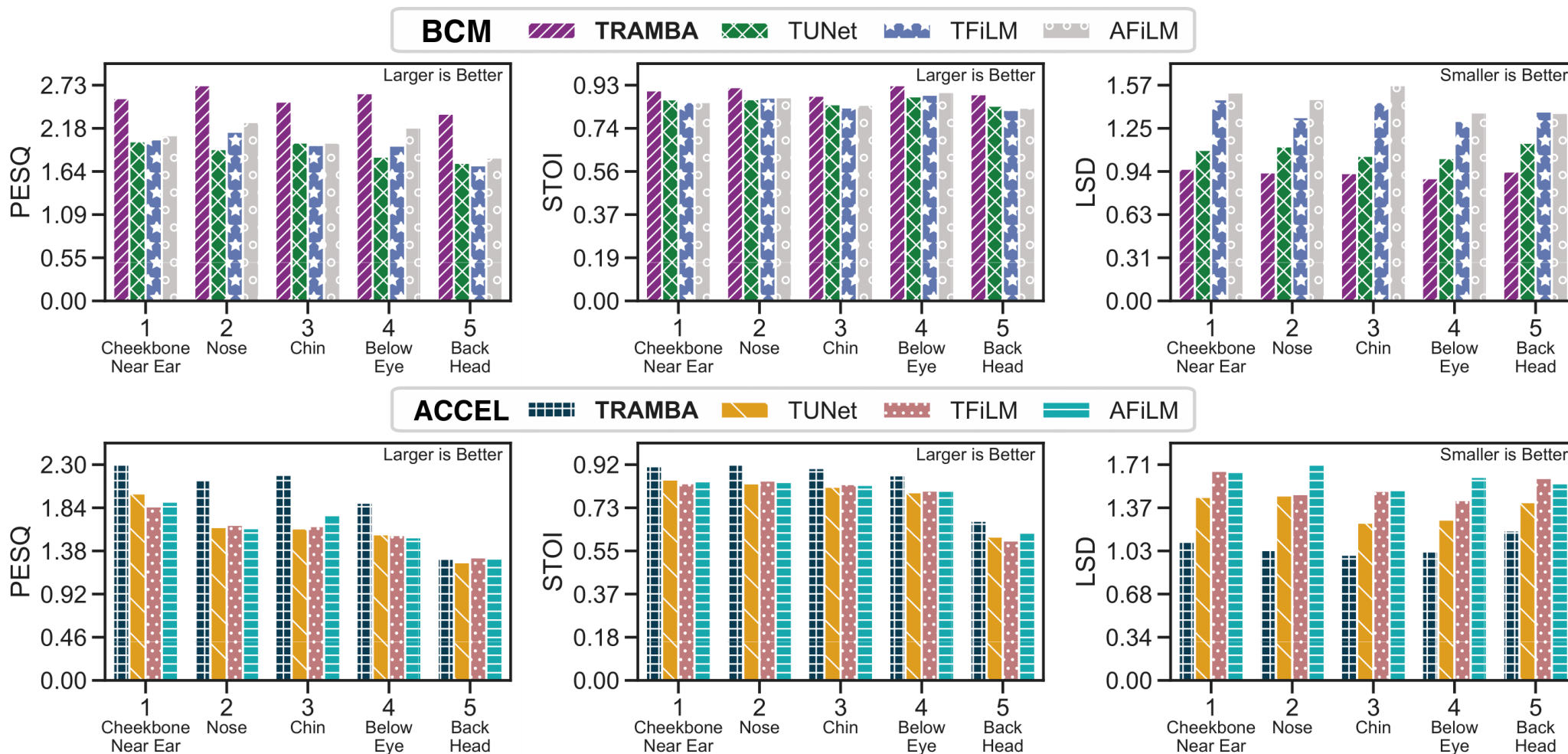
Fine-tuning Performance



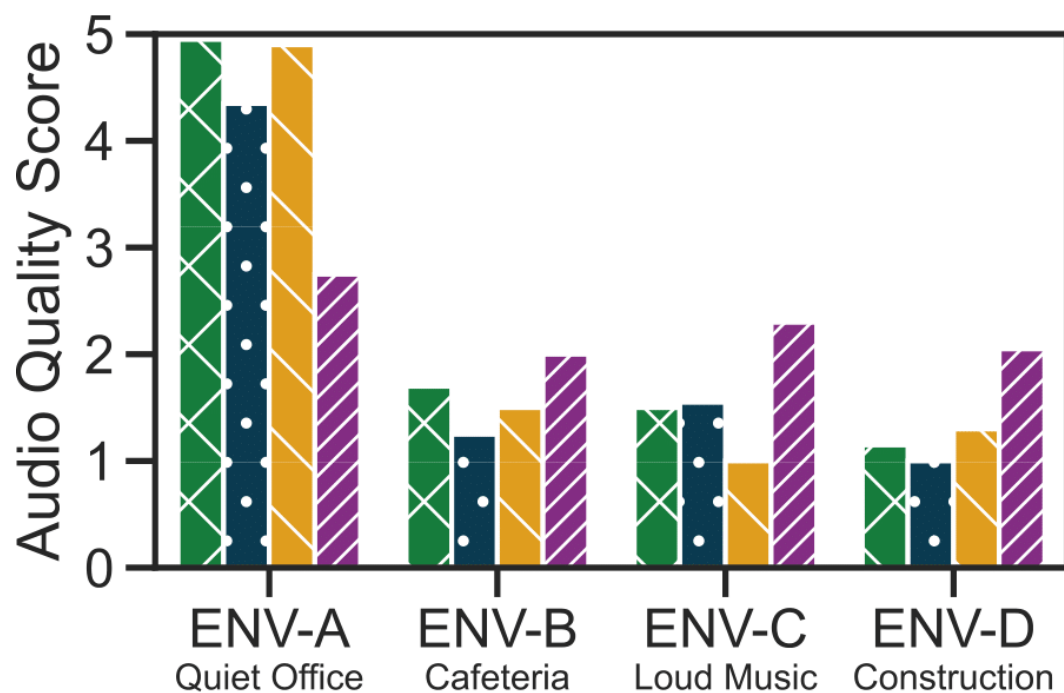
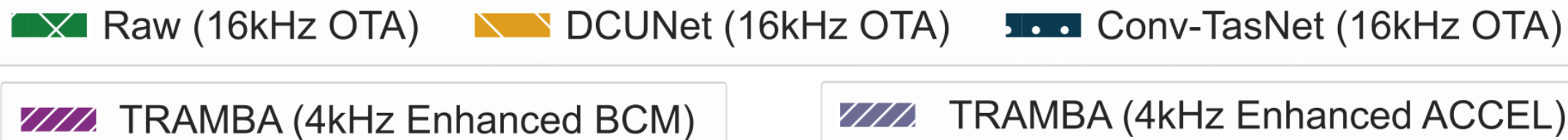
Fine-tuning Performance



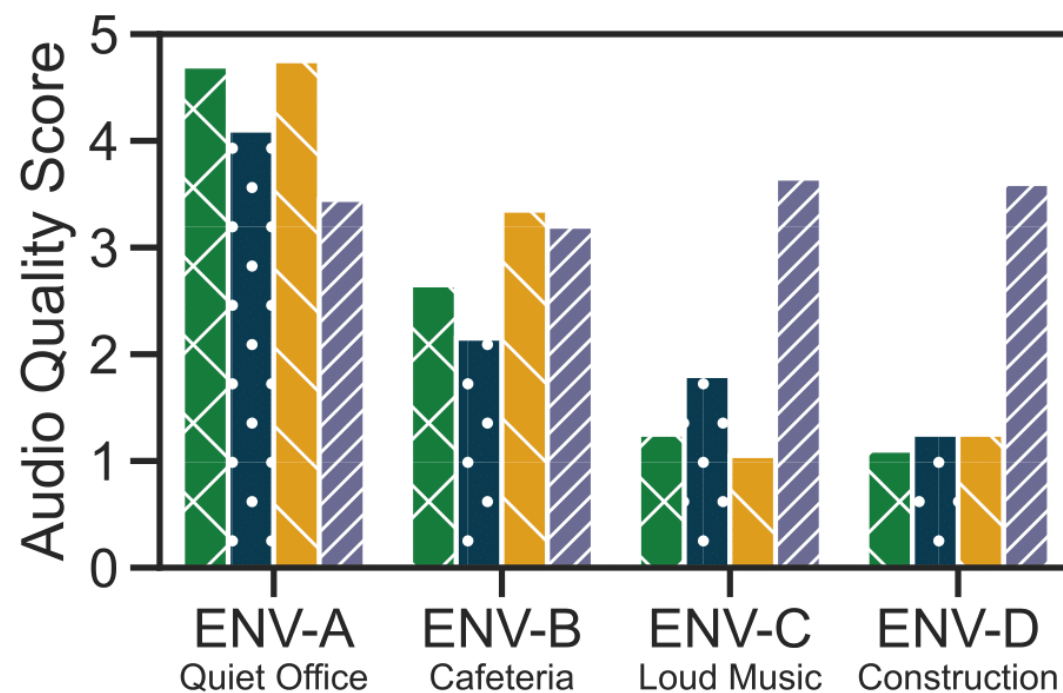
Microphone Placement



Different Environments



(a) TRAMBA-enhanced BCM Data



(b) TRAMBA-enhanced ACCEL Data

Under Motion

