# VibOmni: Towards Scalable Bone-conduction Speech Enhancement on Earables

Lixing He, Yunqi Guo, *Member, IEEE*, Haozheng Hou, Zhenyu Yan, *Member, IEEE*

*Abstract*—Earables, such as True Wireless Stereo earphones and VR/AR headsets, are increasingly popular, yet their compact design poses challenges for robust voice-related applications like telecommunication and voice assistant interactions in noisy environments. Existing speech enhancement systems, reliant solely on omnidirectional microphones, struggle with ambient noise like competing speakers. To address these issues, we propose VibOmni, a lightweight, end-to-end multi-modal speech enhancement system for earables that leverages bone-conducted vibrations captured by widely available Inertial Measurement Units (IMUs). VibOmni integrates a two-branch encoder-decoder deep neural network to fuse audio and vibration features. To overcome the scarcity of paired audio-vibration datasets, we introduce a novel data augmentation technique that models Bone Conduction Functions (BCFs) from limited recordings, enabling synthetic vibration data generation with only 4.5% spectrogram similarity error. Additionally, a multi-modal SNR estimator facilitates continual learning and adaptive inference, optimizing performance in dynamic, noisy settings without on-device back-propagation. Evaluated on real-world datasets from 32 volunteers with different devices, VibOmni achieves up to 21% improvement in Perceptual Evaluation of Speech Quality (PESQ), 26% in Signal-to-Noise Ratio (SNR) and about 40% WER reduction with much less latency on mobile devices. A user study with 35 participants showed 87% preferred VibOmni over baselines, demonstrating its effectiveness for depolyment in diverse acoustic environments.

*Index Terms*—Speech enhancement, earables, bone-conduction vibration.

## I. INTRODUCTION

Earables are smart devices designed to be worn on users' heads or ears, including products such as True Wireless Stereo (TWS) earphones, VR/AR headsets, and smart glasses. These devices are equipped with various sensors and support a range of applications, including virtual and augmented reality (VR/AR), motion recognition, and voice assistants. TWS earphones, like the Apple AirPods series, have significantly contributed to the growth of the earables market, which has become the largest category among all wearable devices, with an estimated shipment of over 273 million units worldwide in 2023 [2]. Manufacturers are continuously enhancing earables by adding new functionalities. For instance, many earphones and headphones now feature active noise cancellation (ANC) to improve the listening experience. Additionally, voice-related applications that utilize the microphones in earables are among the most commonly used features.

L. He, Y. Guo, H. Hou, and Z. Yan are with the department of information engineering, The Chinese University of Hong Kong, Hong Kong.
This paper is an extended version of our prior work appeared in the proceeding of ACM MobiSys 2023 [1].

One of the key applications of earables is making phone calls, which an increasing number of people are utilizing. Users can also interact with voice assistants, such as Siri and Alexa, through voice commands. However, the speech quality on earables often falls short due to several challenges: First, most earables use omnidirectional microphones that capture sound from all directions, which can result in unwanted environmental noise. Second, while many earables are equipped with multiple microphones to create an array for noise reduction via beamforming, the proximity of the microphones to one another can hinder their performance in effectively isolating the user's voice. Third, the speech audio is significantly attenuated by the time it reaches the earables, as they are typically positioned far away from the user's mouth.

Various approaches have been developed for speech enhancement. Signal processing-based methods [3] remove noises based on their statistical models. However, these approaches cannot handle complex environments. Microphone beamforming [4], [5] removes noises based on their directions. But they fail to distinguish the speaker's voice and noise when the microphones are placed too close. Several research [6], [7], [8] adopt deep neural networks (DNNs) to improve speech quality. However, their performance varies when the domain changes. Except for the audio-only solution, several works leverage other modalities like contact sensors [9], vibration sensor [10], camera [11], mm-Wave Radar [12], [13], [14], ultrasonic sensing [15], [16], and Lidar [17], which introduce additional hardware requirements or user overhead to existing earables.

Motivated by the above works, we envision finding a modality that: 1) can be obtained with existing earables without significant modification. 2) has a close correlation with the user's speech and may not be influenced by noises. Since earables are well connected to the user's head, it becomes a desired position to obtain the vibration transmitted to the user's head. The vibration, which is also known as bone conduction vibration, is less influenced by the ambient sound and mainly depends on the clear speech from the user's vocal tract. Fortunately, the existing sensors of commercial earables include an Inertial Measurement Unit (IMU), whose original function is to track the head pose. At the same time, the IMU (including accelerometer and gyroscope) can also capture the subtle vibration on the head, which corresponds to the bone conduction vibration. Considering the mainstream sampling rate of IMU is about 1.6 kHz, the bandwidth (800Hz) overlaps with the lower part of the frequency range of human speech.

To enhance speech quality in earables, we propose Vi-bOmni, a low-latency, end-to-end multi-modal speech en-
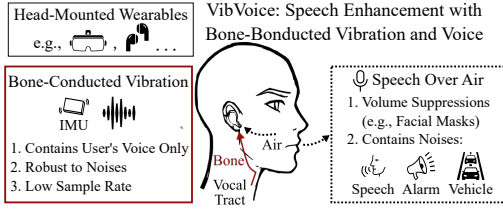
Fig. 1: VibOmni enhances speech quality of head-mounted wearables by extracting the user's clear voice from the bone-conducted vibrations.

hancement system using bone-conducted vibrations and audio, as shown in Fig. 1. Optimized for mobile devices, VibOmni leverages vibration's noise-robustness and audio's rich information, addressing three key challenges:

- **Multi-Modal Data Fusion**: Microphone audio, with its higher sample rate, provides richer information than vibration data, risking DNN overfitting and neglect of vibration input. A specialized training strategy with dual loss functions is needed to balance both modalities.
- **Scarcity of Paired Data**: Collecting paired audio-vibration datasets for head-mounted wearables is labor-intensive. Existing public datasets lack paired data, and large-scale collection involves significant overhead from diverse volunteers and extensive annotated recordings.
- **Real-World Deployment**: Dynamic interferences complicate deploying deep learning speech enhancement. Limited and variable-quality user-collected data may hinder training a robust model, even with rapid data collection.

First, unlike previous work like audio-only speech enhancement or target speaker extraction that relies on time-invariant speaker embeddings, our multi-modal speech enhancement network (§IV-B) employs two encoders to extract and concatenate features from both modalities, with a projection layer aligning frequency dimensions due to differing sampling rates. A DPRNN module separates speech at the feature level and utilizes two decoders to estimate the full-band and low-band speech (to avoid the collapse of multi-modal learning), which is trained by SISNR loss. The whole model is constructed by causal convolutions and unidirectional RNNs for frame-by-frame inference.

Second, to overcome the scarcity of paired audio-vibration datasets, VibOmni employs an innovative pre-training strategy using Bone Conduction Functions (BCFs) (§IV-C). By estimating BCFs from limited recordings and applying them to public datasets like LibriSpeech [18], we generate synthetic vibration data, achieving a low 4.5% error in spectrogram similarity. This approach, enhanced by cubic interpolation and Gaussian modeling, ensures diverse training data and improves model generalization across users and conditions, reducing the need for extensive user-specific fine-tuning. These advancements enable VibOmni to bridge the gap between training and real-world testing, particularly in noisy, dynamic environments.

Lastly, to tackle real-world speech enhancement challenges, VibOmni integrates advanced continual learning and adaptive inference mechanisms, as detailed in Sections IV-D3 and

IV-D4. We propose a multi-modal SNR estimator to distinguish clean audio from noise, enabling a continual self-supervised learning approach that leverages only noisy, in-the-wild mixtures. Additionally, an adaptive inference mechanism dynamically adjusts the model's depth based on estimated noise levels, optimizing computational efficiency. These strategies ensure effective, privacy-preserving adaptation without requiring on-device back-propagation, making VibOmni ideal for resource-constrained mobile devices in diverse acoustic environments.

We evaluate VibOmni on both add-ons of earables and the development board, with paired vibration and audio from 32 volunteers in total. We evaluate VibOmni's performance with both a synthetic noise dataset and in-the-wild noise. Specifically, VibOmni achieves the best performance with 31 times less latency on mobile devices than the two strong baselines, which is promising to be deployed on mobile devices. In addition, data augmentation can reduce the requirement of paired data by more than 72 times. Besides, our proposed SNR estimator obtains 3dB errors in average, which is much better than the audio-only baseline. With the estimated SNR, the adaptive training and testing achieve up to 3dB boost of SNR improvement and prune the unnecessary computation. For the perceptual evaluation, we recruit 35 volunteers for a user study in which VibOmni is preferred by 87% users compared to the baseline.

We summarize our contributions as follows:

- We develop a multi-modal deep neural network that extracts clean speech from noisy audio with assistance from the vibration signal, and employs a lightweight architecture suitable for low-latency execution on mobile devices.
- We introduce a novel data augmentation approach for modeling the Bone Conduction Function, enabling the augmentation of paired vibration and audio data by leveraging a large public dataset.
- We propose a multi-modal SNR estimator to enhance adaptation in training and testing phases, enabling continual learning and adaptive inference, which boosts the effectiveness and efficiency.

## II. RELATED WORK

### A. Speech Enhancement

*a) Audio-only enhancement:* Traditional speech enhancement assumes signal stationarity, speech-noise independence in the time-frequency domain [3], or uses microphone arrays for beamforming to enhance audio quality by leveraging arrival time differences [4], [5]. These methods struggle with dynamic noises without prior knowledge. Recent DNN-based approaches [6], [7], [8] enhance speech by capturing voice and noise features from large datasets, supporting either single-channel [6] or multi-channel audio with [7]. Importantly, ClearBuds [8] uses DNNs for stereo audio from Earbuds but needs dedicated devices.

*b) Multi-modal enhancement:* Except for using audio only, any other modalities that are correlated to the speech can be leveraged for speech enhancement. Previous work

[11] uses camera videos to correlate audio-visual data for speech enhancement and separation via deep learning and cross-modal embeddings. However, cameras are not always available in daily scenarios. Differently, works like [15], [16] use smartphones to emit inaudible acoustic signals ($>17$ kHz) and capture lip-reflected echoes to enhance noisy audio. Others works [14] combine mmWave radar with microphones for speech recognition. Recent studies use bone conduction sensors or accelerometers with microphones for multi-modal speech enhancement [19], [20], [21]. However, their scalability is questionable without an explicit understanding of bone-conduction.

*c) Target speech enhancement:* Different from introducing a new modality, a predefined target can be applied to speech enhancement as well. Compared to incorporating a new modality (e.g., vision, vibration) as the target, the above target can be obtained in a one-time effort. Previous work considers features that correlate well with the target sound as the condition, such as the speaker embeddings [22], sound class [23], or proximity to the user [24]. However, all of them need manual involvement of the user, such as setting a key parameter (the distance to define proximity).

### B. Acoustic Sensing on Earables

*a) In-ear sensing:* In-ear sensing in ANC earphones leverages the built-in microphone and speaker to detect vital signs and enable various applications. The occlusion effect, where the ear canal amplifies low-frequency sounds, supports step counting, activity recognition, and gesture recognition, as demonstrated in studies like Oesense [25]. Additionally, active sensing (using both the microphone and the speaker) using ultrasonic waves emitted by the speaker and analyzing their reflections enables applications such as authentication [26], silent speech interfaces [27].

*b) Out-ear sensing:* Sensing the details of the face by earables primarily relies on active sensing due to the absence of the occlusion effect present in in-ear sensing. Research focuses on analyzing facial expressions [28], enhancing speech [29], authentication [30]. Smart glasses, with their distinct microphone and speaker placement, support applications like authentication [31], facial expression analysis [32], and pose estimation [33].

## III. BACKGROUND AND MOTIVATION STUDY

### A. Bone Conduction Vibration

In this paper, we focus on the vibration transmission from the vocal cord to the skull and refer to it as bone conduction vibration as follows:

$$s_{vib} = f(s_{speech}) + \epsilon_{vib} \quad s_{mic} = s_{speech} + \epsilon_{mic} \quad (1)$$

where $s_{vib}$ and $s_{mic}$ are the raw data captured by the accelerometer and the microphone, respectively; $s_{speech}$ denotes the ground-truth (clean) speech audio; $\epsilon_{vib}$ and $\epsilon_{mic}$ are environmental noises captured by the accelerometer and the microphone, respectively; and $f$ is the Bone Conduction Function (BCF). The noiseless feature of BCF has been discovered



(a) *EarSense* with an Airpods Pro.  (b) Test Setup.

Fig. 2: *EarSense* is an open-sourced data collector attachable to commercial earabless for vibration sensing.



(a) The user is talking with no noise as the reference recording.  (b) The user is talking with a competing speaker.
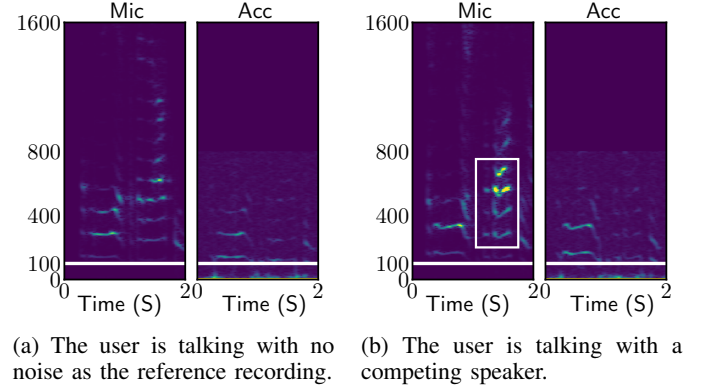
Fig. 3: Microphone and accelerometer recording.

[21], which has also been applied in extremely harsh environments like battlefields or underwater [34]. However, there are still challenges to leveraging it for speech enhancement on earables: 1) the propagation through the head is complicated, resulting in an unstable response even for the same location [35], [36]. 2) They are too expensive (e.g., 60 US dollars) and not available on commercial earables.

### B. Vibration Sensing Platform

Most commercial earables do not provide APIs to collect raw acceleration data. Recently, Apple provided APIs [37] to collect motion data from AirPods earphones at $100\,\text{Hz}$, which is too low for speech recording and doesn't fully utilize commercial IMU. Hence, it is desirable to develop a new sensing platform *EarSense* that can collect the acceleration and acoustic data synchronously on commercial earables. Fig. 2 presents the prototype, which is a sensing platform with a 3D-printed enclosure and a Bosch BMI-160 IMU sensor [38]. We can attach our platform on commercial earables like AirPods Pro to capture bone-conducted vibration, as shown in Fig. 2a. Specifically, two EarSense units connect to a battery-powered Raspberry Pi (RPi) for data collection, and the accelerometer data is streamed by I2C. The microphone and accelerometer sample at 16 kHz and 1.6 kHz, respectively. Specifically, we only access mono audio recording due to AirPods Pro limitations.

### C. Measurement Study

We conducted an experiment where the user (with *EarSense*) was in a meeting room ($10m^2$), speaking when we recorded both the audio and vibration. Note that we apply L2-Norm to the three axes of acceleration to extract vibration
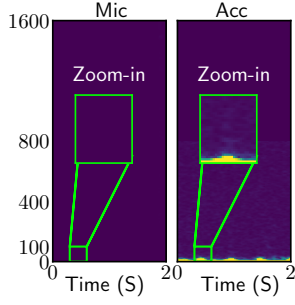
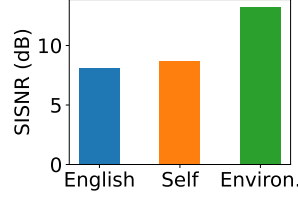Fig. 4: The user is walking with no noise as the motion reference.



Fig. 5: Enhancement performance on different kinds of noises.



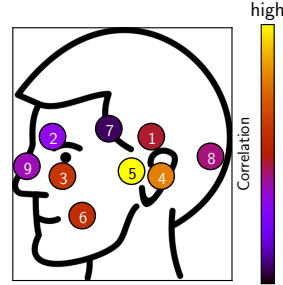Fig. 6: Potential placements on devices.



Fig. 7: Intensities of received bone-conducted vibrations on ten locations of the head.
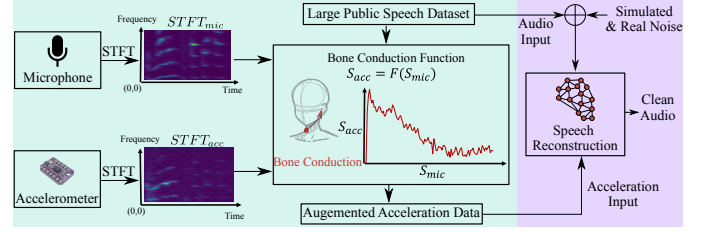


Fig. 8: The overview of VibOmni system.

7): upper ear, eyebrow, cheekbones, ear, temporomandibular joint, cheek, temple, back of the head, nose, and headphone (EarSense taped to an over-ear headphone pad). Some locations suit commercial earables like glasses and headsets (Fig. 6). Pearson correlation coefficients between audio and acceleration, visualized in Fig. 7, show vibrations at most locations, with higher correlations near the mouth, enhancing clean speech extraction.

**Summary:** a) Bone-conducted vibration is robust to environmental voice, b) The user's motion only generates vibrations lower than $85\,\mathrm{Hz}$, c) BCF has a unique frequency response, and d) We can receive bone-conducted vibration from multiple locations on the head.

## IV. METHODOLOGY

### A. Overview of VibOmni

We introduce VibOmni, a novel speech enhancement system designed for earables, as illustrated in Fig. 8. The system comprises three core components: a multi-modal speech enhancement network, pre-training with BCFs (BCFs), and adaptation strategies for speech enhancement.

First, the multi-modal speech enhancement network leverages the complementary strengths of audio and vibration signals, as detailed in Section IV-B. Unlike traditional methods that rely on time-invariant speaker embeddings, VibOmni employs time-dependent features from vibration signals to preserve fine-grained information, balancing compression and performance for robust speech extraction.

Second, to address the scarcity of paired audio-vibration datasets, we propose a pre-training strategy using BCFs, as described in Section IV-C. By estimating BCFs from limited recordings and applying them to public audio datasets like LibriSpeech, we generate synthetic vibration data for training. This approach ensures diverse audio-vibration pairs, enhancing the model's generalization across users and conditions.

Lastly, VibOmni is expected to perform reliably in diverse, noisy environments where clean paired data is scarce. Traditional supervised learning struggles in such conditions, requiring robust algorithms adaptable to varied noises. We propose an adaptive speech enhancement that boost VibOmni in both training and testing. For the adaptive training, we propose a continual self-supervised learning framework (Section IV-D3) that leverages noisy audio via multi-modal SNR estimation and SNR-aware training, enabling effective use of noisy data for robust performance. For the adaptive testing, we propose an adaptive inference framework (Section IV-D4) that adjusts computational resources based on estimated noise levels, modulating separator block depth to balance quality

intensity while reducing effects from wearing position and motion.

**Competing speaker** refers to the scenario when there is another speaker around the user. We simulate it with a loudspeaker 1 meter away playing pre-recorded speech at 3 dB SNR, matching the user's voice volume. This poses a challenge for earables due to the similarity to the user's speech. Spectrograms (Fig. 3a, 3b) show that microphone audio is distorted by the competing speaker (highlighted in Fig. 3b), while accelerometer data, capturing bone-conducted vibrations, remains unaffected, attenuating high frequencies. Thus, accelerometers effectively isolate the user's voice from environmental noise, enhancing speech clarity.

**User motion** may impact the reading of the accelerometer. In Figure 4, we illustrate that walking introduces low-frequency noise in the accelerometer data. The green boxes highlight a zoomed-in view of signals below 100 Hz, revealing clear periodic fluctuations at low frequencies (i.e., $< 50 Hz$), which correspond to the volunteer's steps. Consequently, user motion primarily affects the accelerometer at lower frequencies, leaving speech on higher frequencies unaffected.

**Frequency Response** is the key property of BCF, as we can observe in Fig. 3a and 3b where the vibration signal lacks high-frequency components, consistent with the low-pass filter effect described earlier. We leave the details of the modeling of BCF for the later section.

**Locations on the head** can impact the BCF property since it indicates a different propagation path. We measured bone-conducted vibration intensities at ten head locations (Fig.
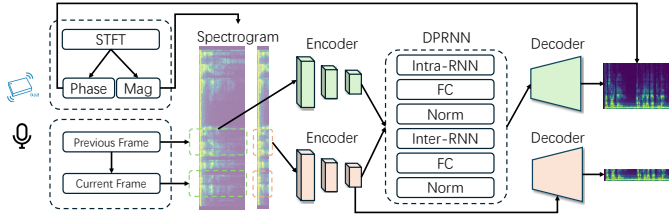
Fig. 9: Network architecture of multi-modal speech enhancement network.

and efficiency. We re-trained the speech enhancement model by multi-level loss function, enabling arbitrary number of blocks conditioned on the noise level for resource-constrained devices.

### B. Multi-Modal Speech Enhancement

*1) Problem formulation:* Suppose the recordings of the microphone and vibration sensor (e.g., accelerometer) are $S_m$ and $S_v$, respectively. Generally speaking, our design goal is to extract the desired conversation from a mixture of speech under the vibration sensor as a condition. Compared to other conditional speech extraction models like voicefilter [39], the vibration provides fine-grained information but lacks the high-frequency component, as illustrated before.

*2) Model architecture:* We illustrate the neural network architecture as follows, which is shown in Fig. 9.

**Features**. Our speech enhancement approach operates in the time-frequency domain. We transform both waveform audio and vibration signals into spectrograms using the Short-Time Fourier Transform (STFT). Note that the three axes acceleration is first normalized (L2) and converted to a spectrogram. The STFT output is separated into phase and magnitude components, with our primary neural network processing only the magnitude. Due to the typically higher sampling rate of audio compared to vibration, we adjust the STFT parameters accordingly to make both of them have the same time dimension. For instance, we set the window size for audio to 640, while the vibration window size is set to 64 to account for the difference.

**Encoder**. We employ convolutional neural networks (CNNs) to extract high-level features from the two modalities, respectively. The features from both encoders are concatenated along the channel dimension. Each encoder is constructed by stacking basic blocks, where each block comprises a 2D convolutional layer, batch normalization, ReLU activation, and max-pooling. To accelerate training, we incorporate a residual shortcut from the block input to the layer before the final deconvolution. To enhance the receptive field and capture harmonic patterns across the entire spectrogram, we use dilated convolutions instead of standard ones. The audio data has a sampling rate of 16 kHz, which is ten times higher than the acceleration data sampling rate of 1.6 kHz. As a result, we use different window lengths for the audio and vibration, which leads to a ten times larger frequency dimension than that of the vibration data. To merge the differences, we apply three more layers to the audio branch and a projection layer at the end of the vibration encoder to ensure that the dimensions are aligned.

**DPRNN**. After the encoder, we aim to separate speech and noise at the feature level, drawing on concepts from sound separation models. To achieve this, we incorporate the DPRNN [40] neural network, known for its lightweight and effective performance in sound separation. Specifically, DPRNN employs two RNNs: one for inter-block modeling (time dimension) and another for intra-block modeling (frequency dimension). We configure DPRNN to focus on a single source, as our goal is to isolate speech without estimating both noise and speech.

**Decoder**. We designed two decoders: a fusion decoder and an auxiliary decoder. Both share the same block structure as the encoder but feature a decreasing number of filters and increasing output tensor sizes. The fusion decoder processes concatenated features from both modalities to generate a spectrogram mask, which is applied to the original noisy audio spectrogram through element-wise multiplication to produce a clean spectrogram. We then incorporate the noisy phase and apply inverse STFT to reconstruct the waveform audio. The auxiliary decoder, which processes only accelerometer features, predicts the low-frequency component of the clean audio. We add the auxiliary decoder since the audio branch contains much more information than the vibration branch. During training, the model may discard the vibration branch since audio-only speech enhancement is also valid when the noise is different from the target speech. To compensate for the collapse, forcing the vibration branch to reconstruct itself can make sure the features are valid.

**Real-time inference**. Real-time speech enhancement is critical to minimize user experience disruptions. When processing audio at the frame level, the processing time per frame must be shorter than the frame duration to ensure the end-to-end latency equals the frame length. However, models like encoders, decoders, and separators typically depend on the current frame, past frames, and sometimes future frames. To enable frame-by-frame inference, previous hidden states must be retained and utilized in subsequent frames. To further reduce latency, all convolutional layers should employ causal convolution, which relies solely on past frames. Additionally, the RNN module in DPRNN should be configured as unidirectional to support streaming processing.

**Loss**. The fusion decoder uses the SISNR loss, where we use $s$ and $\hat{s}$ to represent the clean and enhanced signals as follows:

$$L = 20 \log_{10} \left( \frac{\|s\|^2}{\|s - \hat{s}\|^2} \right) \qquad (2)$$

The above loss applies to full-band audio, which may omit the vibration data that only exists in the low-band. Consequently, we have another auxiliary loss that applies to the extra decoder for the vibration data. Different from the above loss, we use the low-passed clean audio as the target instead. We set a weight of 0.05 for the auxiliary loss to balance the scales of the two losses.

### C. Pre-training with BCF

*1) Problem formulation:* Motivated by our findings in Section III-C, bone-conducted vibration is a promising complementary sensing modality to microphone recording for speech
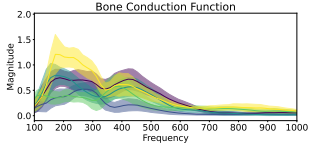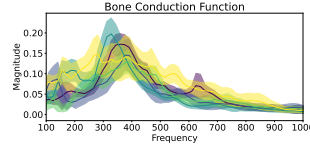
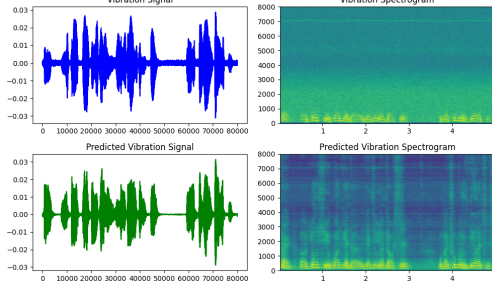Fig. 10: BCF of dataset 1.    Fig. 11: BCF of dataset 2.



Fig. 12: Vibration and the reconstructed vibration.

enhancement under environmental noises. However, there is no large dataset with paired acceleration and audio on earabless available. As a result, we consider using the BCF, which depicts the transfer function from audio to vibration, to generate a virtual vibration dataset. We note that this estimation does not resort to black-box deep learning approaches, which require a huge amount of training data. Instead, we take advantage of prior knowledge that a frequency response exists between the acceleration and audio spectrograms.

*2) Function estimation:* To estimate the BCF, we split the paired data into 5-second windows, and each window can contribute one sample of BCF. Specifically, the power spectral density (PSD) is considered the frequency response between the two signals, which can be estimated using Welch's method [41]. This method involves dividing the data into overlapping segments, computing a modified periodogram for each segment, and then averaging these periodograms. Since both audio and vibration signals are sparse in the frequency domain, the estimated PSD may not always be reliable. To address this, we model the BCF as a Gaussian distribution in the frequency domain because it exhibits and similar pattern with non-trivial variance due to the complex structure of the head skeleton [35], as illustrated in Fig. 10 and 11. Specifically, we denote this function as $f \sim N(\mu, \sigma^2)$, in which $\mu$ and variance $\sigma$ contribute to the contour and fluctuation of frequency response, estimated from the recording of a group of users.

*3) Data Augmentation with BCFs:* We develop a data augmentation approach with BCFs described in Section IV-C2. Note that we cannot apply the inverse BCF to turn the acceleration back to audio due to the significantly limited sample rate of the acceleration data. In addition, the frequency band (larger than $500\,\mathrm{Hz}$) has very slim energy, which can cause a large energy after the inversion. We utilize these functions to generate acceleration signals using a large-scale audio dataset, i.e., LibriSpeech [18]. To be specific, for each audio clip, we first select a BCF (i.e., a list of means and variances over different frequencies) from the pool randomly. Then, we restore the frequency response from the Gaussian distribution of given parameters. Lastly, we can augment the

audio with synthetic acceleration data by directly multiplying the frequency response. Fig. 12 shows the spectrograms of augmented acceleration and real acceleration signals, respectively. The augmented spectrogram is close to the real acceleration spectrogram. We compute the similarity by calculating the mean of the absolute distance of all pixels in the whole spectrogram and dividing it by the largest value of the real acceleration spectrogram. The average error for all volunteers is only 4.5%, which indicates that our proposed acceleration augmentation is reliable and reduce the data collection overhead.

### D. Adaptive Speech Enhancement

*1) Problem formulation:* In speech enhancement, performance depends not only on the quality of bone-conducted vibrations—which are influenced by personal physiological factors and hardware quality, as discussed earlier—but also on the characteristics of interfering noise. For instance, louder noise typically presents a greater challenge for suppression.

To improve the performance of VibOmni, we propose making the system noise-aware, allowing for adaptive speech enhancement. By leveraging noise characteristics, we can dynamically adjust enhancement strategies, unlocking new capabilities in real-world scenarios. This approach involves two key components:

- Noise-aware training – Training a model exclusively on noisy data to improve robustness.
- Noise-aware inference – Optimizing inference under limited computational resources while maintaining performance.

Before implementing these strategies, however, we must first accurately estimate the noise profile.

*2) Noise estimation:* Based on the analysis above, it is essential to first obtain knowledge of noise. We conclude that there are two impact factors of noise: noise type and noise levels. The former can be represented by the audio classification of the noise, which is relatively mature [42]. On the other hand, the noise level can be represented by the SNR.

We have observed that there is a strong correlation between bone-conducted vibration and the user's speech. The above property also indicates that the combination of audio and vibration can infer the SNR. Specifically, the correlation between audio and vibration decreases when the audio is contaminated by noise. Therefore, we propose a multi-modal SNR estimator that builds upon the audio-only SNR estimator described in [43]. Specifically, we transform both audio and vibration by STFT and keep the magnitude only. We conduct zero-padding for the vibration spectrogram is necessary. Then, the spectrograms are concatenated. The estimator consists of five convolutional layers, each with a kernel size of 4, 128 channels, and a stride of 1. This is followed by a statistical pooling layer and two fully connected layers, each containing 256 neurons. We use ReLU (Rectified Linear Unit) activation functions for all the layers. To ensure the estimated SI-SNR values range between -20 and 20 dB, we normalize the network's output to fall between 0 and 1. In this normalization, a value of 0 corresponds to -20 dB, while a value of 1
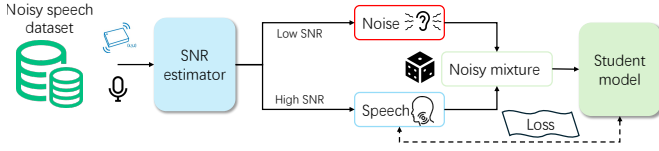
Fig. 13: Continual learning (adaptive training) of VibOmni.



Fig. 14: Adaptive inference of VibOmni.

corresponds to 20 dB. This range compression is achieved using a sigmoid function applied to the output of the network.

*3) Noise-aware training:* Given the obtained noise information—including noise type and noise level—we first focus on addressing noise type. Similar to target speech enhancement [22] or target sound enhancement [23], conditioning the enhancement process on noise type is a promising approach. However, noise type is highly dynamic and cannot be sufficiently represented by either a single class or a static feature. Instead, we characterize noise type through the concept of noise domain, which captures the distribution of noise sources that vary spatiotemporally (e.g., by location and time). To explore this, we trained our multi-modal network using in-domain noise samples: Mandarin speech from AI-SHELL [44] and general environmental noise from FSD50K [45]. When evaluating the model on out-of-domain noise—such as English speech (TIMIT [46]), self-noise from the speaker, or environmental sounds (DEMAND [47])—we observed a significant performance degradation, as shown in Fig. 5. The red line (in-domain performance) consistently surpasses results for English noise and self-noise, highlighting the challenge of domain generalization.

Consequently, to enable VibOmni's capability to work on the noise type, equivalent to domain adaptation to a new noise domain in a continuous manner. A straightforward solution is to fine-tune the model with out-of-domain noise, effectively making it in-domain. Considering it is not practical to ask the user to record clean speech and noise as we do in the lab, self-supervised or unsupervised learning becomes necessary. Specifically, our problem can be formulated as an unsupervised continual learning where we only have access to noisy audio that contains in-domain noise. As a naive solution, we can consider the noisy audio as the clean audio; there is no doubt that it can not perform as well as clean audio, or even collapse when the data is extremely noisy. However, we observe that there are also relatively clean audio samples in the noisy data, so it is critical to find out those samples (high-quality data) effectively. Suppose we can select the high-quality data, then we also need an effective algorithm to utilize it.

Suppose the noisy audio is $S_m$ and vibration is $S_v$, the estimated SNR is $E = Estimator(S_m, S_v)$. In our initial approach, we only select samples with an estimated SNR above a certain threshold, categorizing these as positive samples. Conversely, samples with an estimated SNR below another threshold are classified as negative samples. We can then mix these positive and negative samples offline to fine-tune the model. While this method is effective, it discards a large portion of the data, limiting performance due to the reduced availability of high-quality samples. To enhance the quality of our dataset, we can preprocess both the noisy audio and
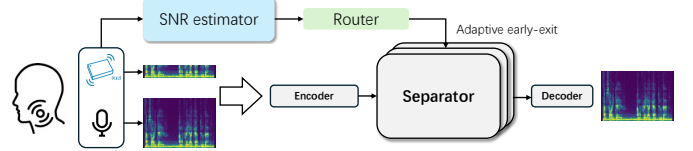
vibration using a pretrained speech enhancement model. The pretrained model, which plays a similar role as the teacher model in RemixIT [48], can be updated at the end of each epoch to ensure continuous enhancement of our data quality.

---

**Algorithm 1** Continual learning for the noisy dataset $D_m$

---

1: $\{\mathbf{D}_{\text{clean}}, \mathbf{D}_{\text{noise}}, \mathbf{D}_{\text{noisy}}\} = \{\}, threshold = \beta$
2: **for** each batch $\mathbf{m} \in \mathcal{D}_m, \mathbf{m} \in \mathbf{R}^{B \times T}$ **do**
3: $\quad$ SNR $\leftarrow$ ESTIMATE_SNR($\mathbf{m}$)
4: $\quad$ **if** $|SNR| > threshold$ **then**
5: $\quad\quad \{\mathbf{D}_{\text{clean}}, \mathbf{D}_{\text{noise}}\} \leftarrow (\mathbf{m}, \text{SNR})$
6: $\quad$ **end if**
7: $\quad$ **if** $\mathcal{D}_{\text{clean}}, \mathcal{D}_{\text{noise}}$ is enough **then**
8: $\quad\quad (\mathbf{D}_{\text{noisy}}) \leftarrow$ REMIX_DATASET($\mathbf{D}_{\text{clean}}, \mathbf{D}_{\text{noise}}$)
9: $\quad$ **end if**
10: **end for**

---

*4) Noise-aware inference:* Beyond noise type, noise level also critically impacts speech enhancement—but in a distinct way. Unlike noise type, it is uncommon for a model to perform well only on strong noise while failing on weak noise. Instead, noise level exhibits a clear correlation with task difficulty: higher noise levels generally demand more aggressive enhancement, while lower levels may require minimal processing. This property opens opportunities to optimize inference efficiency. During deployment of our multi-modal speech enhancement system, we observed dynamic variations in both user speech presence and noise intensity. Running the same enhancement model continuously—even on relatively clean signals—wastes computational resources and introduces unnecessary latency. A straightforward optimization is to integrate voice activity detection (VAD) using head vibration signals [49], activating enhancement only when speech is detected.

However, even during active speech, noise levels can fluctuate substantially (e.g., from loud to near-absent), making on-off enhancement inefficient. To address this, we propose an adaptive VibOmni that dynamically adjusts its processing based on input characteristics. The key is identifying a conditioning factor that reliably reflects enhancement difficulty (e.g., real-time noise level). By feeding this factor to the model as an input, we can selectively allocate computational resources—minimizing overhead without sacrificing performance.

Based on prior analysis, noise level serves as the key factor for enabling adaptive speech enhancement by leveraging the previously proposed SNR estimator. Specifically, we can set an SNR threshold and dynamically adjust the depth of the speech enhancement model until the output meets the desired threshold. The number of separator blocks is an ideal control

| Dataset name | Content | Duration | Role |
|---|---|---|---|
| LibriSpeech-train | English | 1000 hours | Pre-train |
| LibriSpeech-dev | English | 1000 hours | Noise |
| Ai-shell | Mandarin | hours | Noise |
| VibVoice | English | 3 hours | Fine-tune |
| FSD50K | General sound | hours | Noise |

TABLE I: Information about the dataset used.

factor, as it significantly impacts computation without altering the model's pipeline. To support adaptive inference, the training process for speech enhancement must also be modified. This involves defining multiple loss functions corresponding to different numbers of modules, which are averaged to compute the final loss as follows:

$$L = \sum_{i=0}^{N} w_i 20 \log_{10}\left(\frac{\|s\|^2}{\|s - \hat{s}_i\|^2}\right)$$

, where $\hat{s}_i$ refers to the output after the $i^{th}$ separator modules and $w_i$ refers to the weight.

## V. EXPERIMENT SETUP

### A. Dataset

For pre-training with Bone Conduction Functions (BCFs) in Sec. IV-C, we use LibriSpeech [18] as the source dataset. Besides, we recruited 15 volunteers to collect audio-vibration datasets in a clean lab and noisy real-world settings, reading English content from LibriSpeech [18]. Each volunteer contributed 10 minutes of data. Noise Datasets are assumed to be clean by default. For training and evaluation, we add noise with SNRs from -5 dB to 15 dB (average input SNR of 5 dB). Noise types, in equal proportions, include: 1) FSD50K [45] general noise, 2) speech noise from Ai-shell [44] or LibriSpeech [18], and 3) self-noise from the same user (different utterance). Audio is convolved with room impulse responses from [50] to mimic real-world conditions.

### B. Metrics

**Signal-to-Noise Ratio (SNR)**: Measures signal quality relative to noise, computed as $\text{SNR}(x, y) = 20 \log_{10}\left(\left(\frac{y}{x-y}\right)^2\right)$, where $x$ is the estimated audio and $y$ is the clean audio. Higher values indicate better quality. Scale-invariant SNR is used by default. **Perceptual Evaluation of Speech Quality (PESQ)**: Per ITU-T P.862, assesses speech quality with scores from 1 (poor) to 5 (excellent). Wide-band version evaluates full-band speech. **Log-Spectral Distance (LSD)**: Measures frequency-domain quality between reconstructed and ground truth audio: $\text{LSD}(x, y) = \frac{1}{L}\sum_{l=1}^{L}\sqrt{\frac{1}{K}\sum_{k=1}^{K}\left(X(l,k) - \hat{X}(l,k)\right)^2}$, where $l$ and $k$ are time and frequency indices, $X = \log(|\text{STFT}(y)|^2)$, and $\hat{X} = \log(|\text{STFT}(x)|^2)$. Lower values indicate higher quality.

### C. Baselines

We deploy two baselines, i.e., FullSubNet (FSN) [51] and SEANet (SN) [21]. FSN and SN are two state-of-the-art speech enhancement approaches using audio-only and audio-vibration
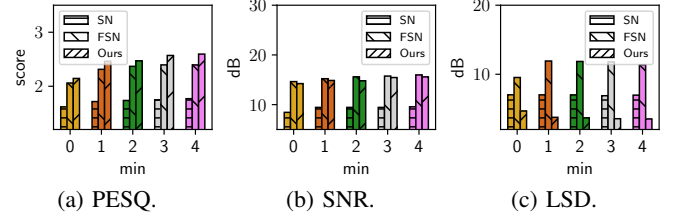


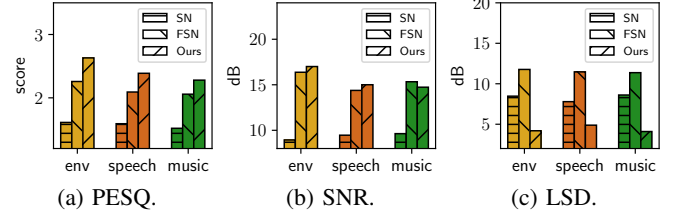Fig. 15: Impact of calibration time.



Fig. 16: Impact of noise types.

inputs, respectively. We train SN using our dataset since the one used in its paper is not available to the public. Specifically, we train SN using vibration data with a sample rate of 1.6kHz to ensure it is the same as VibOmni.

## VI. EVALUATION

### A. Overall Performance

**Calibration**. VibOmni operates out of the box but benefits from target-user data to enhance performance. Fig. 15 compares VibOmni with two baselines using varying amounts of target-user data (zero indicates no user data during training). Results show that more calibration data improves performance across all methods, with VibOmni outperforming baselines at equivalent data levels. VibOmni achieves the highest PESQ for perceptual quality and the best LSD for spectrogram reconstruction, with SNR comparable to FSN due to similar time-domain signal outputs.

**Noise type**. We evaluate the impact of different types of noises in Fig. 16, i.e., environmental noises, competing speakers, and music. The result shows that VibOmni performs better under all noises and metrics except for the SNR with music noise. This is because the complex and dynamic spectrum of the user's speech and music's vocals introduce minor fluctuations in high frequencies, reflecting large fluctuations in SNR.

**Noise level**. We test VibOmni's performance under noise levels of low (10 dB), medium (5 dB), and high (0 dB with only speech noise). The results in Fig. 17 show that VibOmni has better performance improvements, especially when the noise is challenging, i.e., 21% improvement on PESQ and 26% improvement on SNR. This is because the vibration can more robustly identify the target speech, whereas the audio-only
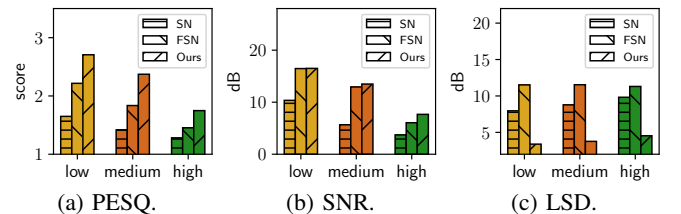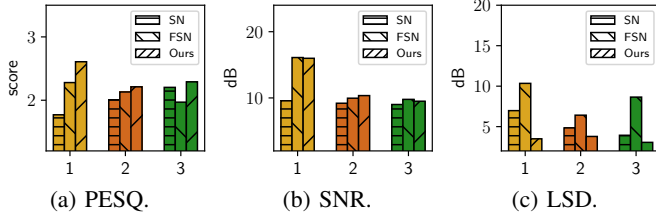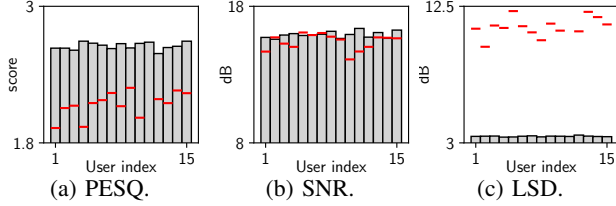


Fig. 17: Impact of noise levels.

(a) PESQ.  (b) SNR.  (c) LSD.

Fig. 18: Number of noise sources.



(a) PESQ.  (b) SNR.  (c) LSD.

Fig. 19: VibVoice on different users. Red line: the performance of the baseline, i.e., FullSubNet.



(a) PESQ.  (b) SNR.  (c) LSD.

Fig. 20: VibVoice on different head locations. Red line: the performance of the baseline, i.e., FullSubNet.

|  | PESQ | SNR | LSD |
|---|---|---|---|
| VibVoice | 2.6 | 15.6 | 3.5 |
| w/o auxiliary decoder | 2.5 | 15.1 | 4.4 |
| w/o augmentation | 1.9 | 14 | 5 |
| w/o Gaussian approx | 2.4 | 15.2 | 4.2 |
| Accelerometer sample rate: 1200 Hz | 2.47 | 14.4 | 4.2 |
| Accelerometer sample rate: 800 Hz | 2.45 | 14.3 | 4.5 |
| Accelerometer sample rate: 400 Hz | 2.4 | 14.2 | 4.4 |

TABLE II: Ablation study.

solution is difficult to differentiate sound with a similar pattern (e.g., strong speech noise).

**The number of noise sources**. We evaluate the impact of different numbers of noise sources by repeatedly mixing clean audio with random audio clips. Fig. 18 shows that VibOmni's performance degrades as the number of noise sources increases, but still outperforms all the baselines.

**Temporal stability**. We further examine how VibOmni performs for the same user over time. Note that the offset of sensor placement and minor changes in speech can cause a slight change. We collect ten-minute data from three volunteers twice, six months apart. The results show that the performance of VibOmni has negligible changes, from 2.6 to 2.5 for PESQ, 15.7 to 15.5 for SNR, and 4.3 to 4.6 for LSD. Besides, VibOmni outperforms FSN, whose performance is 2.1 for PESQ, 15.2 for SNR, and 11 for LSD. The results affirm that VibOmni is robust to temporal changes.

**Airway blockage**. Facial masks, which block air transmission, can reduce speech volume. We tested VibOmni with three volunteers speaking identical content with and without masks. VibOmni's performance degrades minimally: 0.05 ($< 2\%$) for PESQ, 0.4 ($< 3\%$) for SNR, and 0.1 ($< 3\%$) for LSD with masks. Compared to FSN (PESQ: 2.15, SNR: 13.8, LSD: 10), VibOmni shows superior resilience, as expected, due to its use of bone-conducted vibration, unaffected by air transmission.

**Variances among users**. Speech and bone-conducted vibration can differ across users due to vocal features, head and skull shapes, body fat, etc. Fig. 19 shows the performance of VibOmni across 15 users. VibOmni shows stable and significantly better performance in PESQ compared to baseline and comparable performance in SNR.

**Sensor positions**. We test VibOmni when EarSense is placed in ten locations on the head as defined in Fig. 7, validating VibOmni's effectiveness for different HMW devices. The bars in Fig. 20 show VibOmni's performance at each location. The red line represents the performance of the baseline at #4 *ear*. The results show that VibOmni achieves satisfactory performance at all locations in PESQ. Note that locations like #1 *upper ear*, #2 *eyebrow*, #5 *temporomandibular joint*, #7 *temple*, and #10 *interior of the pad of the headphone* show
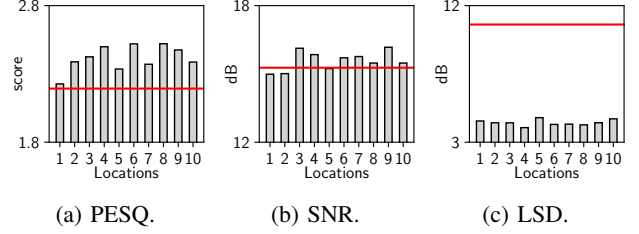
similar or slightly lower SNR than the baseline, which is because the vibration intensity is slim due to their far distance to the audio source.

**Summary**. VibOmni outperforms FSN and SN by up to 21% for PESQ when the noise volume is low, where the PESQ for VibOmni and FSN are 2.7 and 2.21, respectively. VibOmni outperforms the baselines up to 26% for SNR when the noise is speech with high volume, where the SNR of VibOmni and FSN are 2.0 and 1.6, respectively. In addition, VibOmni outperforms the baselines 50~80% in LSD under most impact factors, indicating the efficiency of our multi-modal design and novel data augmentation. VibOmni has slightly lower performance in SNR for some cases, as it can be biased due to the similar spectrum of the user's speech and music's vocals. The dynamic also introduces minor fluctuations. In comparison, LSD evaluates the whole band without preference, so VibOmni outperforms the two baselines by a large margin.

### B. Ablation Study

We conduct an ablation study to understand the performance of different design components in VibOmni. The performance of VibOmni without different components is listed in Table II.

**No auxiliary decoder**. First, we remove the self-supervise loss, meaning the audio may dominate the model. The results indicate that the variant slightly degrades by 0.1 in PESQ, 0.3 in SNR, and 0.8 in LSD, respectively.

**No data augmentation**. Second, we remove the data augmentation based on the Bone Conduction Function. The performance significantly degrades to 1.9 in PESQ, 14 in SNR, and 5 in LSD. According to the definition of ITU and mean opinion score (MOS), the audio quality is poor when the score drops from 2.5 to 1.9 [52]. This confirms that small-scale self-collected datasets are insufficient for robust neural network training. In addition, we compare VibOmni trained on: a) 18–180 hours of paired audio-vibration data [19] (5 kHz vibration bandwidth) and b) three hours of paired data with augmentation. Results in Fig. 21 show that data augmentation
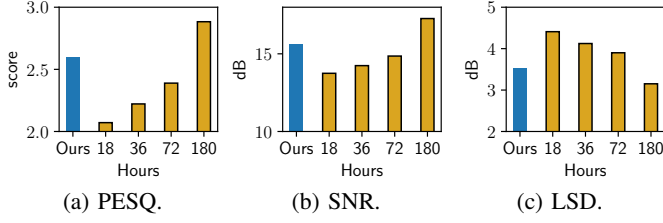
(a) PESQ.　　　　(b) SNR.　　　　(c) LSD.

Fig. 21: Effectiveness of data augmentation. Blue bars: Vib-Voice using less than three hours of paired data with augmentation. Yellow bars: VibVoice using 18- to 180-hour paired data without augmentation.

|  | VibOmni | FSN | SN |
|---|---|---|---|
| Desktop CPU | 0.05 | 0.27 | 0.5 |
| Desktop GPU | 0.016 | 0.034 | 0.07 |
| P30 | 0.16 | 5 | 1.9 |
| Mate20 | 0.29 | 4.6 | 1.7 |
| Pixel7 | 0.31 | 5.2 | 1.4 |

TABLE III: Runtime analysis (second/instance).

achieves comparable performance with $\sim 24\times$ less paired data.

**No Gaussian approximation**. Third, the Bone Conduction Function is modeled by only the mean, while the variance is zero. The performance drops 0.2 in PESQ, 0.4 in SNR, and 0.5 for LSD, indicating that our Gaussian approximation is close to the nature of the Bone Conduction Function.

**Lower sample rate**. Lastly, we evaluate the performance of VibOmni with a lower sample rate by downsampling the vibration data to 1200 Hz, 800 Hz, and 400 Hz. The results show that VibOmni is robust to various sample rates, which consume less power in processing and communication. When the sample rate is 800 Hz, the output's PESQ only degrades 10%.

### C. Runtime Evaluation

We test the execution latency of VibOmni and the two baselines (i.e., FSN [51] and SN [21]) on a desktop PC (i.e., i7-11700k CPU and RTX 3060 GPU) and three smartphones (i.e., Huawei P30, Huawei Mate20, and Google Pixel 7). We run the inference of a 5-second clip 100 times and record the mean latency. The results in Table III show that VibOmni reduces up to $31\times$ and $12\times$ less latency on average than FSN and SN, respectively. On the other hand, the results show that VibOmni exhibits a greater advantage in runtime for low-end devices. In conclusion, VibOmni can support real-time voice applications with a delay of fewer than 0.6 seconds (i.e., two times the inference latency) for processing 5 5-second audio clip. The real-time factor is only 0.12, significantly less than the minimal requirement, which means less energy consumption and space to process other tasks.

### D. Extension Evaluation

Except for the dataset collected from our prototype, we observe that there are a few public dataset that contains both vibration and synchronized audio, including EMSB [19] and ABCS [20]. Specifically, both are collected in quiet places

| SNR (dB) | ABCS | EMSB |
|---|---|---|
| VibOmni (TFGridNet) | 10.76 | 11.46 |
| VibOmni (DPRNN) | 7.96 | 8.03 |
| TFGridNet | 9.25 | 10.54 |
| DPRNN | 3.04 | 3.08 |

TABLE IV: Performance comparison on public datasets: ABCS and EMSB.
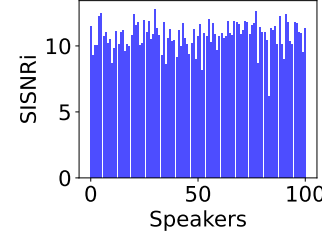


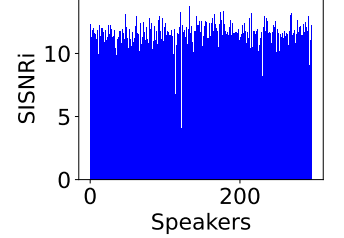Fig. 22: Performance for different speakers (ABCS).



Fig. 23: Performance for different speakers (EMSB).

by customized earphones with a vibration sensor (without detailed documentation). Since both datasets are Mandarin, we consider the Mandarin dataset Ai-shell as the speech noise. We evaluate VibOmni on the two datasets, compared to the latest baselines: DPRNN [40], TFGridNet [53]. We report the average performance of the two datasets in the Table. IV. Besides, we present the SNR improvement, against the input noisy audio (5dB) in Fig. 22 and 23 for each speaker of the dataset. We observe that our model works well for most of the speakers and obtains around 10 dB improvement, indicating a similar performance on our self-collected dataset.

### E. Adaptive Speech Enhancement

To evaluate the continual learning, we evaluate it from two perspectives: 1) the performance of SNR estimation, and 2) the performance of continual learning.

**SNR Estimation** is assessed using the mean average error (MAE) between the estimated SNR and the ground truth. It is important to note that the output is constrained to the range of possible SNR values (-20dB to 20dB). In Fig. 24, we compare the proposed SNR estimation method with the audio-only SNR estimation [43]. Our observations indicate that the benefits of multi-modal learning are not as pronounced when the input SNR is relatively low ($SNR < 0dB$); however, the estimation accuracy significantly improves when the input SNR is high. In contrast, the audio-only SNR estimation struggles in high SNR conditions, as it fails to distinguish between the user's speech
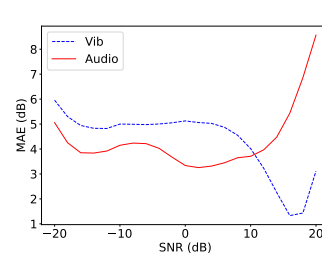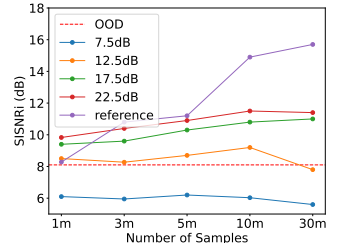


Fig. 24: SNR estimation errors.



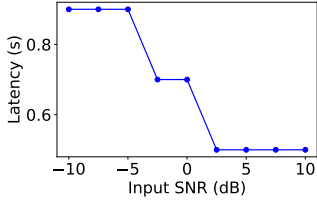Fig. 25: Continual learning performance vs. SNR threshold

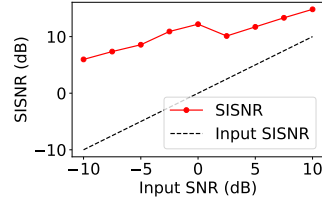Fig. 26: Adaptive speech enhancement latency.



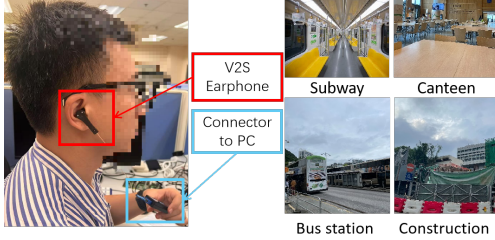Fig. 27: Adaptive speech enhancement performance.



Fig. 28: Volunteers with Knowles V2S200D (left) and locations where we collect the dataset (right).

and similar background interferences. Since the continual learning is interested in finding the clean audio (high SNR), our proposed SNR estimator is much better than the baseline. **Adaptive training** is evaluated at various SNR thresholds. Specifically, we classify data samples with an estimated SNR above the threshold as "clean" data. It's important to note that while a higher threshold indicates more reliable data, it also results in a lower proportion of effective data. As illustrated in Fig. 25, when we set the threshold at 7.5dB or 12.5dB, the training process fails due to the presence of noisy data. In contrast, with thresholds of 17.5dB and 22.5dB, continual learning achieves an approximate 3dB improvement compared to the out-of-domain model, all without requiring any clean data.

**Adaptive Inference** is evaluated in terms of latency and corresponding effectiveness, as shown in Fig. 14. As the input SNR increases, VibOmni dynamically adjusts the model's depth, which helps to reduce computational latency. By default, we set the SNR threshold at 15 dB, meaning that latency begins to decrease once this threshold is surpassed, as indicated in Figs. 26 and 27. In Fig. 26, it is demonstrated that VibOmni automatically changes the model's depth twice, successfully reducing the latency from 0.9 seconds to 0.5 seconds while maintaining the same output SNR.

### F. In-the-wild Evaluation

We have developed an advanced prototype of our earables using Knowles' V2S200D Voice Vibration Sensor development kit [54], which enables dual-channel audio recording (capturing both audio and vibration signals) at a 48 kHz sampling rate. Compared to the prototype described in Sec. III, this version offers greater user-friendliness, making it suitable for conducting user studies.

**Dataset.** Specifically, we conduct extensive experiments to validate the performance of VibOmni in the presence of ongoing noise in real-world environments. Volunteers are asked to read the same content as in the experiments described in Sec. VI-A. In total, our dataset includes 22 speakers. We collect 10
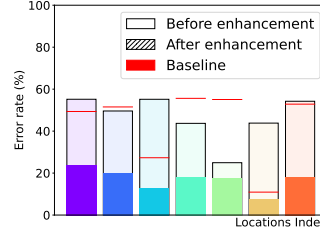


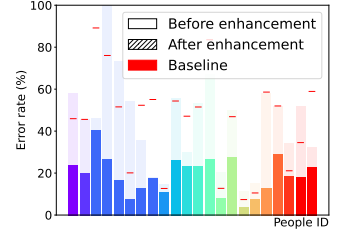Fig. 29: In-the-wild evaluation per location.



Fig. 30: In-the-wild evaluation per speaker.

minutes of data from each volunteer at each location. The data is gathered from seven different locations, where each speaker may appear in multiple settings. Specifically, the locations include: restaurant, roadside, coffee bar, office, park, subway, and bus.

**Metric.** Unlike synthetic noisy speech, we can neither capture the ground truth of clean speech nor evaluate the metrics like SNR and PESQ, nor train the model. Instead, we use the result of Automatic Speech Recognition [1] to evaluate the quality of the speech. We use Word Error Rate (WER $= \frac{S+D+I}{N}$) as the evaluation metric, where S is the number of substitutions, D is the number of deletions, I is the number of insertions, and N is the number of words in the reference. The higher the value of WER, the lower the audio quality. Considering the dataset is in Mandarin, we calculate the Character-Error-Rate in the same way as WER.

**Results.** As shown in Fig. 30 and Fig. 29, VibOmni outperforms both the baseline and the input noisy audio by a large margin in most cases. Specifically, only in one place VibOmni performs similarly to the baseline, which is the roadside, where the noise is less significant. And we do observe that the improvement from baseline is not obvious for 3 people because the audio is not so noisy, so it is not necessary to have good noise removal. In total, VibOmni gains 44% less word error rate compared to the baseline.

### G. User Study

**Questionnaire Design**. We recruited 35 volunteers to evaluate the perceived performance of VibOmni. Volunteers listened to 5-second audio clips from the dataset in Section 5.1. In the first part, they transcribed audio enhanced by VibOmni to assess intelligibility using Word Error Rate (WER). In the second part, they compared pairs of audio clips with identical content, selecting the higher-quality audio in two scenarios: VibOmni versus original noisy audio and VibOmni versus the baseline (FSN). Each comparison was repeated five times. We measured VibOmni's improvement using the correct ratio $\frac{P}{P+N}$, where $P$ is the number of times VibOmni was preferred and $N$ is the alternative.

**Study results**. According to the results of the first part, VibOmni achieves an overall WER of $21.5\%$, which is acceptable for understanding the audio content and confirms the effectiveness of VibOmni. According to the answers to the second question, the survey results show that $87\%$ of the participants choose VibOmni over the baseline, and $72\%$

---

[1] https://github.com/speechbrain/speechbrain

of them choose VibOmni over the original audio without any enhancements. In addition, we discuss with participants why they prefer the original audio sounds over the baseline. The baseline can produce acoustic artifacts and sometimes wrongly suppress the sound of the target speaker. We note that some participants observed that the impact of artifacts and suppression is hindered after knowing the content or listening repeatedly. However, the speech generated by the baseline causes lots of misunderstanding for the first-time listener. In conclusion, the user study results show that VibOmni can enhance speech quality and improve user experience compared with the original audio and the baseline.

## VII. DISCUSSION

- The current wireless communication in earables does not support two-channel audio recording by default. VibOmni requires a bit rate of 153.6 kbps ($6 \times 16$ bits $\times 1.6$ kHz) to transmit acceleration data to a mobile device without compression, which is below the maximal bandwidth for Bluetooth 5.0. However, the compatibility with existing profiles and IMU data compression is necessary.

- To save the computation resource, we can offload neural network inference on a smartphone, but this still brings additional energy consumption from the ADC (0.54 mW). In comparison, each AirPods Pro earbud has a 43 mAh battery, where VibVoice adds only approximately 1.5% to the power consumption of earphones.

- The quality of vibration data in earables depends on the hardware form factor and sensor placement. Based on extensive evaluations, we recommend two guidelines for optimal IMU sensor placement: 1) position sensors close to the vibrating organ, and 2) ensure tight contact with the head. Additionally, considerations for user comfort and compatibility with existing devices are crucial.

- We envision integrating VibOmni with on-device processing to enhance compatibility (eliminating the need for smartphone software installation) and reduce overhead. Devices like OmniBuds [2] already support on-device machine learning, making this a promising avenue to explore.

## REFERENCES

[1] L. He, H. Hou, S. Shi, X. Shuai, and Z. Yan, "Towards bone-conducted vibration speech enhancement on head-mounted wearables," in *Proceedings of the 21st Annual International Conference on Mobile Systems, Applications and Services (MobiSys)*, 2023, pp. 14–27.

[2] Statista, "Wearable shipments by category 2024," https://www.statista.com/statistics/690731/wearables-worldwide-shipments-by-product-category/, 2020, (Accessed on 02/15/2022).

[3] Y. Ephraim and H. L. Van Trees, "A signal subspace approach for speech enhancement," *IEEE Transactions on speech and audio processing*, vol. 3, no. 4, pp. 251–266, 1995.

[4] N. Yousefian, J. H. Hansen, and P. C. Loizou, "A hybrid coherence model for noise reduction in reverberant environments," *IEEE Signal Processing Letters*, vol. 22, no. 3, pp. 279–282, 2014.

[5] Y. Ji, J. Byun, and Y.-c. Park, "Coherence-based dual-channel noise reduction algorithm in a complex noisy environment." in *INTERSPEECH*, 2017, pp. 2670–2674.

[6] C. Subakan, M. Ravanelli, S. Cornell, M. Bronzi, and J. Zhong, "Attention is all you need in speech separation," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 21–25.

[7] S. Zhang and X. Li, "Microphone array generalization for multichannel narrowband deep speech enhancement," *arXiv preprint arXiv:2107.12601*, 2021.

[8] I. Chatterjee, M. Kim, V. Jayaram, S. Gollakota, I. Kemelmacher, S. Patel, and S. M. Seitz, "Clearbuds: wireless binaural earbuds for learning-based speech enhancement," in *Proceedings of the 20th Annual International Conference on Mobile Systems, Applications and Services*, 2022, pp. 384–396.

[9] P. Gupta, Y. Jeong, J. Choi, M. Faingold, and F. Ayazi, "Precision high-bandwidth out-of-plane accelerometer as contact microphone for body-worn auscultation devices," in *2018 Hilton Head Workshop*, 2018, pp. 30–33.

[10] H. A. C. Maruri, P. Lopez-Meyer, J. Huang, W. M. Beltman, L. Nachman, and H. Lu, "V-speech: Noise-robust speech capturing glasses using vibration sensors," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 2, no. 4, pp. 1–23, 2018.

[11] R. Gao and K. Grauman, "Visualvoice: Audio-visual speech separation with cross-modal consistency," in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2021, pp. 15 490–15 500.

[12] M. Z. Ozturk, C. Wu, B. Wang, and K. Liu, "Radiomic: Sound sensing via mmwave signals," *arXiv preprint arXiv:2108.03164*, 2021.

[13] H. Li, C. Xu, A. S. Rathore, Z. Li, H. Zhang, C. Song, K. Wang, L. Su, F. Lin, K. Ren *et al.*, "Vocalprint: exploring a resilient and secure voice authentication via mmwave biometric interrogation," in *Proceedings of the 18th Conference on Embedded Networked Sensor Systems*, 2020, pp. 312–325.

[14] T. Liu, M. Gao, F. Lin, C. Wang, Z. Ba, J. Han, W. Xu, and K. Ren, "Wavoice: A noise-resistant multi-modal speech recognition system fusing mmwave and audio signals," in *Proceedings of the 19th ACM Conference on Embedded Networked Sensor Systems*, 2021, pp. 97–110.

[15] K. Sun and X. Zhang, "Ultrase: single-channel speech enhancement using ultrasound," in *Proceedings of the 27th Annual International Conference on Mobile Computing and Networking*, 2021, pp. 160–173.

[16] Q. Zhang, D. Wang, R. Zhao, Y. Yu, and J. Shen, "Sensing to hear: Speech enhancement for mobile devices using acoustic signals," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 5, no. 3, pp. 1–30, 2021.

[17] S. Sami, Y. Dai, S. R. X. Tan, N. Roy, and J. Han, "Spying with your robot vacuum cleaner: eavesdropping via lidar sensors," in *Proceedings of the 18th Conference on Embedded Networked Sensor Systems*, 2020, pp. 354–367.

[18] 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)*. IEEE, 2015, pp. 5206–5210.

[19] H. Wang, X. Zhang, and D. Wang, "Fusing bone-conduction and air-conduction sensors for complex-domain speech enhancement," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 3134–3143, 2022.

[20] M. Wang, J. Chen, X. Zhang, Z. Huang, and S. Rahardja, "Multi-modal speech enhancement with bone-conducted speech in time domain," *Applied Acoustics*, vol. 200, p. 109058, 2022.

[21] M. Tagliasacchi, Y. Li, K. Misiunas, and D. Roblek, "Seanet: A multi-modal speech enhancement network," *arXiv preprint arXiv:2009.02095*, 2020.

[22] B. Veluri, M. Itani, T. Chen, T. Yoshioka, and S. Gollakota, "Look once to hear: Target speech hearing with noisy examples," in *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 2024, pp. 1–16.

[23] B. Veluri, M. Itani, J. Chan, T. Yoshioka, and S. Gollakota, "Semantic hearing: Programming acoustic scenes with binaural hearables," in *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, 2023, pp. 1–15.

[24] T. Chen, M. Itani, S. E. Eskimez, T. Yoshioka, and S. Gollakota, "Hearable devices with sound bubbles," *Nature Electronics*, pp. 1–12, 2024.

[25] D. Ma, A. Ferlini, and C. Mascolo, "Oesense: employing occlusion effect for in-ear human sensing," in *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services*, 2021, pp. 175–187.

[26] Z. Wang, S. Tan, L. Zhang, Y. Ren, Z. Wang, and J. Yang, "Eardynamic: An ear canal deformation based continuous user authentication using

---

[2] https://www.omnibuds.tech/

in-ear wearables," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 5, no. 1, pp. 1–27, 2021.

[27] X. Dong, Y. Chen, Y. Nishiyama, K. Sezaki, Y. Wang, K. Christofferson, and A. Mariakakis, "Rehearsse: Recognizing hidden-in-the-ear silently spelled expressions," in *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 2024, pp. 1–16.

[28] X. Song, K. Huang, and W. Gao, "Facelistener: Recognizing human facial expressions via acoustic sensing on commodity headphones," in *2022 21st ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN)*. IEEE, 2022, pp. 145–157.

[29] D. Duan, Y. Chen, W. Xu, and T. Li, "Ease: Bringing robust speech enhancement to cots headphones," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 7, no. 4, pp. 1–33, 2024.

[30] D. Duan, Z. Sun, T. Ni, S. Li, X. Jia, W. Xu, and T. Li, "F2key: Dynamically converting your face into a private key based on cots headphones for reliable voice interaction," in *Proceedings of the 22nd Annual International Conference on Mobile Systems, Applications and Services*, 2024, pp. 127–140.

[31] K. Li, D. Agarwal, R. Zhang, V. Gunda, T. Mo, S. Mahmud, B. Chen, F. Guimbretiĕre, and C. Zhang, "Sonicid: User identification on smart glasses with acoustic sensing," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 8, no. 4, pp. 1–27, 2024.

[32] K. Li, R. Zhang, S. Chen, B. Chen, M. Sakashita, F. Guimbretiĕre, and C. Zhang, "Eyeecho: Continuous and low-power facial expression tracking on glasses," in *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 2024, pp. 1–24.

[33] S. Mahmud, K. Li, G. Hu, H. Chen, R. Jin, R. Zhang, F. Guimbretiĕre, and C. Zhang, "Posesonic: 3d upper body pose estimation through egocentric acoustic sensing on smartglasses," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 7, no. 3, pp. 1–28, 2023.

[34] "Bone-conduction ear microphone — shop motorola solutions," https://shop.motorolasolutions.com/bone-conduction-ear-microphone-system/product/PMLN5464A, 2022, (Accessed on 08/19/2022).

[35] Y. Chang, N. Kim, and S. Stenfelt, "The development of a whole-head human finite-element model for simulation of the transmission of bone-conducted sound," *The Journal of the Acoustical Society of America*, vol. 140, no. 3, pp. 1635–1651, 2016.

[36] S. Y. Won and J. Berger, "Estimating transfer function from air to bone conduction using singing voice," in *ICMC*, 2005.

[37] Apple, "Cmheadphonemotionmanager — apple developer documentation," https://developer.apple.com/documentation/coremotion/cmheadphonemotionmanager, 5 2023, (Accessed on 05/16/2023).

[38] B. Sensortec, "Inertial measurement unit bmi160," https://www.bosch-sensortec.com/products/motion-sensors/imus/bmi160/, 2021.

[39] Q. Wang, H. Muckenhirn, K. Wilson, P. Sridhar, Z. Wu, J. Hershey, R. A. Saurous, R. J. Weiss, Y. Jia, and I. L. Moreno, "Voicefilter: Targeted voice separation by speaker-conditioned spectrogram masking," *arXiv preprint arXiv:1810.04826*, 2018.

[40] Y. Luo, Z. Chen, and T. Yoshioka, "Dual-path rnn: efficient long sequence modeling for time-domain single-channel speech separation," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 46–50.

[41] P. Welch, "The use of fast fourier transform for the estimation of power spectra: A method based on time averaging over short, modified periodograms," *IEEE Transactions on audio and electroacoustics*, vol. 15, no. 2, pp. 70–73, 2003.

[42] F. Schmid, K. Koutini, and G. Widmer, "Efficient large-scale audio tagging via transformer-to-cnn knowledge distillation," in *ICASSP 2023-2023 IEEE international Conference on acoustics, Speech and signal processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[43] C. Subakan, M. Ravanelli, S. Cornell, and F. Grondin, "Real-m: Towards speech separation on real mixtures," 2021.

[44] H. Bu, J. Du, X. Na, B. Wu, and H. Zheng, "Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline," in *2017 20th conference of the oriental chapter of the international coordinating committee on speech databases and speech I/O systems and assessment (O-COCOSDA)*. IEEE, 2017, pp. 1–5.

[45] E. Fonseca, X. Favory, J. Pons, F. Font, and X. Serra, "Fsd50k: an open dataset of human-labeled sound events," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 829–852, 2021.

[46] W. Fisher, "The darpa speech recognition research database: specifications and status," in *Proc. DARPA Workshop on speech recognition*, 1986, pp. 93–99.

[47] J. Thiemann, N. Ito, and E. Vincent, "Demand: a collection of multi-channel recordings of acoustic noise in diverse environments," *(No Title)*, 2013.

[48] E. Tzinis, Y. Adi, V. K. Ithapu, B. Xu, P. Smaragdis, and A. Kumar, "Remixit: Continual self-training of speech enhancement models via bootstrapped remixing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1329–1341, 2022.

[49] P. Schilk, N. Polvani, A. Ronco, M. Cernak, and M. Magno, "In-ear-voice: Towards milli-watt audio enhancement with bone-conduction microphones for in-ear sensing platforms," in *Proceedings of the 8th ACM/IEEE Conference on Internet of Things Design and Implementation*, 2023, pp. 1–12.

[50] E. Hadad, F. Heese, P. Vary, and S. Gannot, "Multichannel audio database in various acoustic environments," in *2014 14th international workshop on acoustic signal enhancement (IWAENC)*. IEEE, 2014, pp. 313–317.

[51] X. Hao, X. Su, R. Horaud, and X. Li, "Fullsubnet: A full-band and sub-band fusion model for real-time single-channel speech enhancement," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6633–6637.

[52] Wikipedia, "Perceptual evaluation of speech quality - wikipedia," https://en.wikipedia.org/wiki/Perceptual_Evaluation_of_Speech_Quality, 2020, (Accessed on 12/02/2022).

[53] Z.-Q. Wang, S. Cornell, S. Choi, Y. Lee, B.-Y. Kim, and S. Watanabe, "Tf-gridnet: Making time-frequency domain models great again for monaural speaker separation," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

[54] [Online]. Available: https://www.digikey.com/en/product-highlight/k/knowles/v2s200d-voice-vibration-sensor
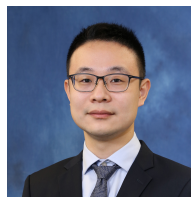
**Lixing He** received the B.S. degree in automation from UESTC, Chengdu, China, in 2021, and is a Ph.D. student at Embedded AI and IoT Lab (AIoT Lab), Department of Information Engineering, The Chinese University of Hong Kong. His research interests include audio, wearables technology, and human-centric sensing,

**Yunqi Guo** received the B.S. degree in computer science from Shanghai Jiao Tong University, Shanghai, China, in 2016, and the M.S. and Ph.D. degrees in computer science from the University of California, Los Angeles, Los Angeles, CA, USA, in 2018 and 2023, advised by Prof. Songwu Lu. He is currently a Postdoctoral Fellow at The Chinese University of Hong Kong, working with Prof. Guoliang Xing. His research interests lie at the intersection of augmented reality, mobile systems, visual–language interaction, and intelligent sensing.

**Haozheng Hou** is a Ph.D. student at Embedded AI and IoT Lab (AIoT Lab), Department of Information Engineering, The Chinese University of Hong Kong. He received his B.S. degree (2021) from Nanjing University. His research interests include underwater acoustic sensing and human activity recognition.

**Zhenyu Yan** is an Assistant Professor at The Chinese University of Hong Kong. Dr. Yan has extensive experience in sensing systems, signal and information processing, cyber-physical systems, and machine learning in IoT systems. His works have been published in top international conferences and journals, such as MobiCom, SenSys, IPSN, IEEE Transactions on Mobile Computing, and ACM Transactions on Sensor Networks. He is the recipient of the Rising Star Award from ACM SIGBED China. His papers also received the Best Community Contributions Award at ACM MobiCom 2023, the Best Paper Award Runner-up at ACM MobiCom 2022, and the Best Artifact Award Runner-up at ACM/IEEE IPSN 2021.