

Wireless Hearables With Programmable Speech AI Accelerators

Malek Itani¹, Tuochao Chen¹, Arun Raghavan², Gavriel Kohlberg², Shyamnath Gollakota¹

¹Paul G. Allen School of Computer Science and Engineering, University of Washington

²Department of Otolaryngology—Head and Neck Surgery, University of Washington

Abstract

The conventional wisdom has been that designing ultra-compact, battery-constrained wireless hearables with on-device speech AI models is challenging due to the high computational demands of streaming deep learning models. Speech AI models require continuous, real-time audio processing, imposing strict computational and I/O constraints.

We present *NeuralAids*, a fully on-device speech AI system for wireless hearables, enabling real-time speech enhancement and denoising on compact, battery-constrained devices. Our system bridges the gap between state-of-the-art deep learning for speech enhancement and low-power AI hardware by making three key technical contributions: 1) a wireless hearable platform integrating a speech AI accelerator for efficient on-device streaming inference, 2) an optimized dual-path neural network designed for low-latency, high-quality speech enhancement, and 3) a hardware-software co-design that uses mixed-precision quantization and quantization-aware training to achieve real-time performance under strict power constraints. Our system processes 6 ms audio chunks in real-time, achieving an inference time of 5.54 ms while consuming 71.6 mW. In real-world evaluations, including a user study with 28 participants, our system outperforms prior on-device models in speech quality and noise suppression, paving the way for next-generation intelligent wireless hearables that can enhance hearing entirely on-device.

1 Introduction

In recent years, intelligent hearables have made significant advances in enhanced hearing, leveraging deep learning to program acoustic scenes in real time [7, 10, 60]. These advancements enable capabilities such as speech enhancement [52], noise suppression [7], and even target speech hearing [61]. However, current implementations rely on wired headsets and computationally demanding platforms like smartphones or high-power embedded systems [10, 60, 61]—devices with far greater processing power, memory, and energy budgets than what can realistically fit within ultra-compact, battery-constrained hearables like earbuds and hearing aids.

The fundamental challenge is that streaming deep learning models are traditionally associated with high computational demands and power consumption. This raises key questions about whether real-time speech AI models can

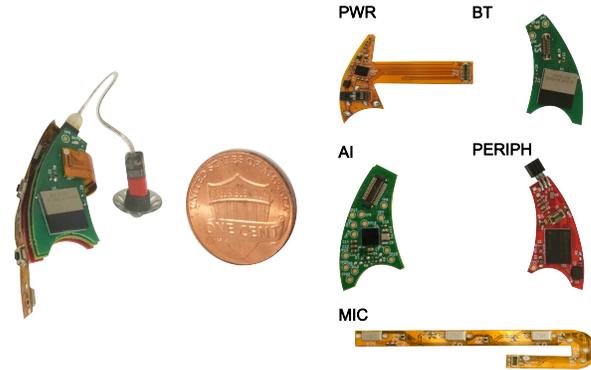


Figure 1: NeuralAid hardware. The device has five interconnected flexible and rigid circuit boards that together form an AI-enabled hearable. PWR (2-layer flexible PCB): manages power, charging, and programming, BT (2-layer rigid PCB): houses the BLE SoC, AI (6-layer rigid PCB): contains a low-power AI accelerator for real-time speech AI, PERIPH (4-layer rigid PCB): hosts peripherals including RAM, NOR flash, IMU, and I2S DAC, MIC (2-layer flexible PCB): has a microphone array with three mics and two push buttons.

be deployed on small, battery-powered hearables. Offloading computations to a smartphone is unreliable due to the stringent sub-10 ms algorithmic latency requirements of enhanced and augmented audio applications, which are highly sensitive to wireless network variability, I/O delays, and operating system overhead on smartphones. Meanwhile, running speech AI models on-device has so far been constrained by the limited availability of AI accelerators that are both powerful enough for real-time speech AI and small enough to fit within these miniature wireless devices. Overcoming these barriers is essential to unlocking the next generation of truly intelligent and self-sufficient wireless hearables.

In this paper, we explore whether it is possible to design wireless hearables that can run real-time speech AI models entirely on-device while understanding the trade-offs between power consumption and efficiency. Achieving this requires addressing three key challenges.

- Existing research platforms for hearables lack the computational capabilities required to run deep learning models for augmented and enhanced audio applications. In contrast to classification tasks, these models must continuously operate in a streaming manner, processing audio input at a minimum sampling rate of 16 kHz. To do this, they should process audio in small chunks of around 6 ms, maintain an inference



Figure 2: A person wearing a NeuralAid, which rests behind the ear, with a receiver inserted into their ear canal.

time of under 6 ms to ensure real-time performance, and continuously output audio at the same sampling rate. Meeting these requirements imposes significant computational and I/O constraints. Further, due to the limited non-volatile storage available on these compact devices, the model size must not exceed 1.5 MB and to enable more than six hours of continuous operation on a 675 hearing aid battery, deep learning power consumption must be below 100 mW.

- Tiny, low-power AI accelerator hardware [1, 3] is inherently far more constrained than GPUs and general-purpose embedded CPUs like the Raspberry Pi. On the other hand, effective neural architectures for speech enhancement rely on components such as convolutions [30], LSTMs [28], transformers [55], and state-space layers [27]. These models prioritize speech quality over real-time, on-device, or low-power constraints — indeed, transformers and state-space models exceed our target hardware’s runtime and memory limits. While efforts like TinyLSTM [20] and its successor TinyDenoiser [46] aim for on-device deployment, their 25 ms algorithmic latency falls short of the low-latency requirements for enhanced hearing applications.

- While smartphones and general-purpose embedded devices can run floating point operations, this increases the runtime and power consumption on low-power hardware. In general, low-power AI accelerators function as a compute cluster with standard processors for general-purpose tasks and a dedicated neural accelerator optimized for low-power, highly quantized integer computations. However, quantizing all neural network weights, activations, and inputs can significantly degrade audio quality. Further, different neural network layers are affected differently by quantization errors, necessitating mixed-precision quantization and compensation for non-linear errors at each layer. Thus, achieving optimal performance requires a low-power hardware-software co-design approach.

We introduce *NeuralAids*, which makes three key contributions across wireless wearable hardware, embedded systems, deep learning, and hardware-software co-design.

- **Wireless wearables with speech AI accelerators.** We design NeuralAids, a real-time speech AI system built on a modular hardware architecture with five stacked PCBs (Fig. 1). It integrates a low-power AI accelerator (GreenWaves GAP9) with an nRF53-based System-on-Chip (SoC) for Bluetooth Low Energy (BLE) connectivity, audio processing, sensors, and both volatile and non-volatile memory. The system runs speech AI models, using GAP9’s compute cluster for real-time neural network inference. Optimized power management, high-speed audio interfaces, and flexible storage enable continuous operation in a behind-the-ear compact form factor (Fig. 2), making NeuralAids a unique platform for real-time speech AI applications.

- **Real-time on-device efficient neural network.** Prior on-device models like TinyDenoiser [46] are constrained by their limited complexity and inability to achieve effective performance with sub-10ms latencies. Instead, we begin with dual-path models that process audio in the time-frequency (TF) domain, treating time and frequency components as sequences using recurrent networks [6, 69]. These architectures achieve state-of-the-art performance by modeling both temporal and frequency relationships. However, they are computationally expensive, as running recurrent networks across each frequency prevents real-time operation. Thus, they cannot run in real time on our hardware. To address this, our proposed architecture compresses frequencies using a linear layer and replaces LSTMs with Gated Recurrent Units (GRUs). These architectural changes and others listed in §2.2.2 achieve real-time performance by reducing runtime while outperforming TinyDenoiser for the speech denoising task.

- **Hardware-software co-design.** Floating-point speech denoising models cannot achieve real-time operation on our target accelerator hardware. To address this, we explore various quantized model configurations and use quantization-aware training (QAT) to simulate non-linear quantization errors and fine-tune the network to mitigate error propagation across layers. This further narrows the performance gap between our mixed-precision and fully floating-point speech models.

We evaluated our system for speech enhancement in noisy scenarios. Our results are as follows:

- Before quantization, our floating-point dual-path network architecture improves speech quality by 2.41 dB over TinyDenoiser. Applying quantization-aware training reduces the performance gap between our mixed-precision and floating-point networks from 7.86 dB to 0.57 dB.

- Our proposed mixed-precision network can process 6 ms audio chunks within 5.54 ms on the GAP9 processor, achieving real-time operation. Furthermore, the model uses only 299 kB of memory and consumes 71.64 mW on our hardware.
- Our system generalizes to six real-world indoor and outdoor environments, as well as to wearer head motion. Our training relies solely on synthetic data and does not require any training data collection using our hearable hardware.
- In a user study with 28 participants, our system achieved a higher mean opinion score and better noise removal compared to both the raw, unprocessed input and TinyDenoisier.

This paper demonstrates that effective speech AI models can indeed run on low-power wireless hearables. We believe that this paves the way for integrating on-device AI models for enhanced hearing into billions of wireless earbuds and hearing aids, unlocking exciting new opportunities.

2 System Design

We present our hearable hardware, an efficient streaming neural network and a hardware-software co-design.

2.1 NeuralAids Hardware

The hardware has a low-power BLE SoC that controls a tiny AI accelerator to process incoming audio in real-time.

2.1.1 Stacked printed circuits boards. The systems are distributed across the 5 printed circuit boards (PCBs) shown in Fig. 1, referred to as PWR, BT, AI, PERIPH and MIC. These circuits are stacked on top of each other and placed in the hearing aid case. The board functionalities are as follows:

- *PWR:* A 2-layer flex board that houses components for power management, battery charging, and monitoring. It also exposes programming pads for the BLE SoC and AI accelerator. The board is powered by a 3.85V battery pack consisting of four CP1254 batteries (75 mAh each, 300 mAh total). A Texas Instruments BQ25120 power management and charging IC is integrated to charge the battery and generate a regulated 1.8V supply required to power the device components. This chip can also produce a 3.3V voltage domain on demand via I2C to power indicator LEDs on other circuits. A switch is included to manually turn the device on and off. The programming pads are designed for interfacing with a 2×7 array of 1.27 mm pogo pins, with one pad also serving as a charging interface for the battery via a 5V source. The board has a long flexible section with a Molex 2167010209 board-to-board connector at the tip that connects to the BT board. Since the board exposes programming pads, it is placed along the edge of the case.
- *BT:* A 2-layer rigid board containing an ISP2053, which is a BLE SoC based on nRF5340 that integrates additional

components such as capacitors, oscillators and an antenna. A resistor is soldered in one of two possible positions to hardcode the device’s side (i.e. left or right device). A board-to-board connector (Molex 513382674) attaches this circuit to the AI board, enabling power and data exchange between the two boards.

- *AI:* A 6-layer rigid board houses the low-power AI accelerator (GreenWaves GAP9). An FPF1204UCX load switch is used to enable the BLE SoC to control power to the AI chip, allowing it to connect and disconnect power as needed. This board connects to the PERIPH circuit.
- *PERIPH:* A 4-layer rigid board contains peripheral devices, including an AP Memory APS256XXN 256 Mbit RAM chip and an SSM6515 ultra-low-power I2S DAC and amplifier. The AI accelerator can control power to the RAM chip using a load switch (FPF1204UCX) when the RAM is not needed. Additionally, this board includes a NOR flash (Macronix MX25UW12845G) and an IMU (BMI323), which are reserved for future research. A right-angled 1 mm pitch connector is used to connect a Phonak Audeo Marvel M Receiver-in-canal (RIC), allowing audio playback into the ear canal. The board connects to the MIC circuit using a Molex 5050661022 connector.
- *MIC:* A 2-layer flexible microphone array with 3 PDM microphones (TDK T5837) and 2 buttons (Omron B3U-1000P).

2.1.2 Hardware subsystems. The NeuralAid system consists of three main subsystems: Bluetooth, AI, and Audio.

- *Bluetooth subsystem:* This is built around the BLE SoC, which is the first chip to boot when the device powers on and manages power delivery to the AI chip. It has direct access to the buttons on the MIC board and receives system-level interrupts such as wake-up and reset signals. There are two channels between the BLE and AI chips: 1) A 115200 baud UART interface for low-speed, spontaneous communication, and 2) an I2S interface for high-speed, continuous communication. The BLE SoC enables wireless connectivity with external devices. After waking up the AI chip, it transmits advertising packets and accepts connections, allowing remote control and data exchange for audio streaming or playback. The BLE SoC is dual-core; we run time-critical BLE-related tasks on one core, while tasks like interfacing or GPIO control run on another core.
- *Audio subsystem:* This manages audio from microphone capture to playback through the RIC. Although audio interfaces directly with GAP9, all its Serial Audio Interfaces (SAIs) derive their clocks from the BLE chip. GAP9 manages three SAIs for different purposes: 1) reading audio from microphones, 2) writing audio to speakers, and 3) exchanging audio with the BLE SoC (e.g., for streaming). The BLE SoC

Table 1: Power consumption of various NeuralAid components. Speaker amplification was calibrated by placing a NeuralAid in a silicone ear model and measuring the sound level 8 mm from the RIC tip. The amplification was adjusted to increase the sound level by 20 dB. The ambient noise level was 33 dBA in both the calibration and testing environments.

Component	Power (mW)
BLE chip	6.75
Microphone array	2.02
Speaker	1.49

provides a 3.072 MHz audio clock, enabling high-quality microphone capture. GAP9’s Smart Filter Unit (SFU) converts incoming PDM samples to PCM using an 8th-order cascaded integrator-comb filter with a $64\times$ decimation ratio and 2 samples per stage. The output is shifted by 24 bits to generate 32-bit PCM samples at 48 kHz, which are then downsampled to 16 kHz via the SFU’s resampler block for processing. For playback, the speaker SAI interface transmits 32-bit PCM samples at 48 kHz via I2S to the DAC. Since internal processing occurs at 16 kHz, the SFU’s resampler upsamples audio back to 48 kHz. The BLE I2S interface relays recorded audio and external playback between the BLE and AI chips.

- *AI subsystem:* The AI accelerator processes incoming audio in real-time using a neural network. Every 96 audio samples (6 ms at 16 kHz), it performs speech enhancement. First, the system converts 32-bit PCM audio to 16-bit float using a dedicated fixed-to-floating point converter. A pre-emphasis filter with a coefficient of 0.97 removes DC components. The audio is then placed in a ring buffer and an FFT converts it into the frequency domain. The neural network processes the frequency-domain audio, after which an inverse FFT and overlap-add reconstruct the time-domain signal. The output is converted back to 32-bit integers for playback through the RIC. Although the microphone array has three microphones, only one is used in this paper. For efficiency, FFT and inverse FFT run on GAP9’s compute cluster, while other operations execute on the fabric controller. All GAP9 components, including the fabric controller, compute cluster, and peripherals, are clocked at at 370 MHz. The AI chip also communicates with onboard RAM via an Octal SPI interface.

2.1.3 Component power consumption. Table 1 presents the power consumption of the BLE SoC, microphone array, and speaker while running real-time AI speech enhancement. The combined power consumption of these components is 10.26 mW. The power consumption of the AI accelerator, which varies based on the complexity of the neural network, is analyzed in detail in the following sections.

2.2 Efficient Streaming Neural Network

We have three key requirements for achieving a streaming speech enhancement network: (1) achieving a low end-to-end latency of less than 20 ms between the audio and visual scenes, with an algorithmic latency of 10 ms or less, (2) completing the processing of the current chunk before the next chunk arrives, and (3) generalizing effectively to unseen reverberant and noisy environments, as well as diverse wearers, without requiring training data collected from hardware.

As shown in Fig. 3A, end-to-end latency is defined as the time for a single audio sample to travel from the microphone input buffer, through the speech enhancement network, and into the speaker output buffer. This latency consists of two components: algorithmic latency, introduced by factors such as chunk size, lookahead, and overlap-add; and hardware latency, which accounts for the computation and I/O time required to process each chunk [63].

2.2.1 Problem Formulation. Given an acoustic environment, let $s(t)$ represent the clean target speech received at the microphone, and let $n(t)$ denote all background noise. The mixture audio, $x(t)$, can be expressed as:

$$x(t) = s(t) + n(t) \tag{1}$$

Since our enhancement network \mathcal{N} processes audio chunks in a streaming manner, we divide the continuous waveform into smaller audio chunks: $x(t) = [x_0, x_1, \dots, x_N]$, $s(t) = [s_0, s_1, \dots, s_N]$, and $n(t) = [n_0, n_1, \dots, n_N]$. The process of streaming speech enhancement can be formulated as,

$$\hat{s}_i, h_i = \mathcal{N}(x_i, h_{i-1}) \tag{2}$$

Here, \hat{s}_i represents an estimation of the clean speech s_i , while h_{i-1} denotes the cache or memory state from previous chunks, which can be reused as the previous intermediate output, and h_i represents the current cache or memory state.

2.2.2 Neural Network Architecture. Existing low-power, on-device speech enhancement networks, like TinyLSTM [20] and its successor TinyDenoyer [46], have two drawbacks: (1) their simple architecture significantly limits enhancement performance, and (2) they exhibit high algorithmic latency of around 25 ms, restricting real-time applications. While several advanced speech enhancement models have been proposed in recent years [8, 48, 62], they fail to meet our real-time, on-device, and low-power constraints. Our goal is to significantly improve speech enhancement performance while maintaining low-latency and real-time constraints.

The overall architecture of our network is shown in Fig.3C. It is a dual-path Time-Frequency (TF) domain model based on TF-GridNet[62], which achieves superior performance for speech enhancement tasks. First, the small audio chunk x_i is transformed into its time-frequency representation X_i using a Short-Time Fourier Transform (STFT) followed by a causal

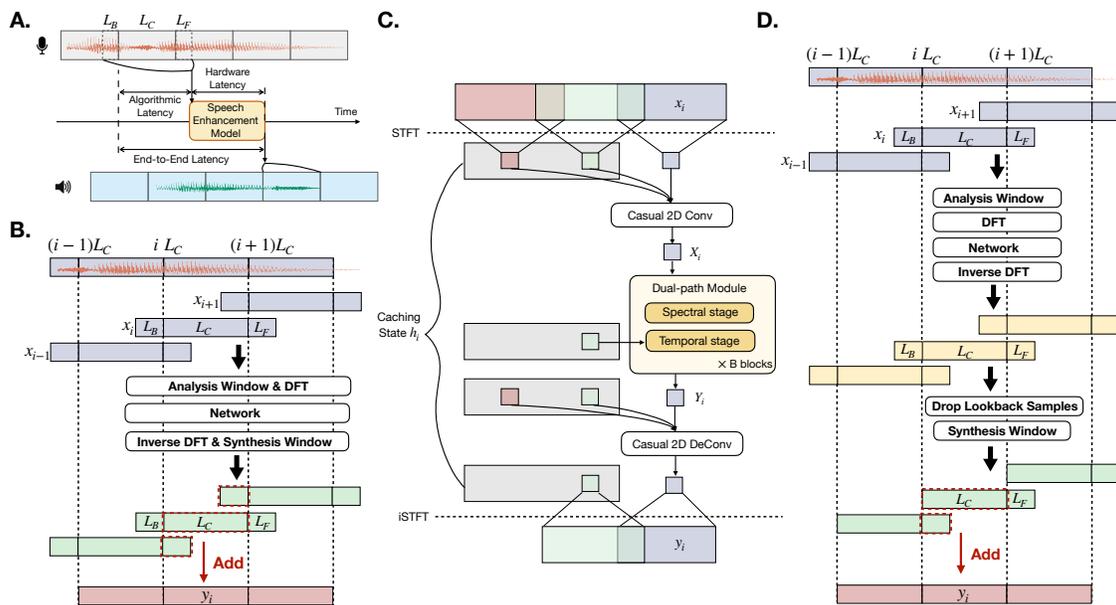


Figure 3: Efficient streaming neural network. (A) Decomposition of the end-to-end latency in streaming speech enhancement. (B) The normal overlap-add operation introduces additional algorithmic latency with lookback padding. (C) The overall architecture of our streaming speech enhancement network with history caching states for computing reuse. (D) Dual-window approach during overlap-add can reduce the additional algorithmic latency introduced by the lookback padding.

2D convolutional layer. The transformed representation X_i is then processed by multiple dual-path modules—the most computationally intensive components of the network. Each dual-path module has two stages: a spectral stage that processes frequency information and a temporal stage that handles time dependencies. Finally, the enhanced time-frequency representation Y_i is passed through a 2D causal deconvolution layer and reconstructed back into a time-domain signal y_i using the inverse STFT and overlap-add operations.

To enable real-time operation on low-power devices while maintaining low latency, we highlight three modifications.

Low-complexity dual-path architecture. In the original TF-GridNet, the spectral stage consists of a bi-directional LSTM that processes the frequency domain of X_i , while the temporal stage consists of a causal uni-directional LSTM that processes the time domain of X_i . We observed that the primary contributor to inference time is the bi-directional LSTM in the spectral stage, as it must process the long sequence of frequency bins sequentially. Using recurrent models to process frequency sequences sequentially slows computation.

To reduce its complexity, we first apply a strided convolution layer to compress the frequency dimension before using the frequency-domain LSTM. This reduces the sequence length processed by the LSTM, thereby decreasing its execution time. Additionally, compared to LSTMs, GRUs offer

lower computational cost while maintaining similar performance [14]. To further optimize the runtime of the spectral module, we replace the LSTM module with a GRU module.

Dual-window approach for low-algorithmic latency.

As shown in Fig. 3A, to achieve low algorithmic latency, we process the input audio x_i in small chunks of size L_C (ms), with an additional lookahead of L_F (ms) and a lookback of L_B (ms); we set $L_C = 6ms$, $L_B = 6ms$ and $L_F = 4ms$. Lookahead and lookback padding improve frequency resolution, which is important for enhancement performance [63]. However, increasing the lookahead and lookback padding also introduce higher algorithmic latency.

With lookahead padding L_F , the network must wait for an additional L_F ms of audio samples before applying the STFT. While lookback padding does not introduce additional waiting time at input, it causes algorithmic latency during the iSTFT process. Specifically, the overlap-add processing step of the current chunk y_i requires the lookback padding from the next chunk y_{i+1} , as shown in Fig. 3B. Thus, the overlap-add process must wait for the next chunk (an additional $L_B + L_C$ ms) to reconstruct the current output chunk.

To eliminate the algorithmic latency introduced by lookback padding, we adopt a dual-window-size approach [32, 63], where the synthesis window is shorter than the analysis window, as shown in Fig. 3D. By discarding the lookback samples after the inverse DFT and before applying the synthesis window, the overlap-add process for the current chunk

no longer depends on the next chunk’s output. We use a rectangular window for analysis and use a specific function for the synthesis window s , described below, valid for $L_C > L_F$, to achieve perfect signal reconstruction [22]:

$$s[i] = \begin{cases} 1 & \text{if } i \in [L_F, L_C] \\ \frac{1}{\lfloor \frac{L_C + L_B}{L_C} \rfloor + 1} & \text{otherwise} \end{cases}$$

With this approach, our lookback padding can increase frequency resolution without introducing additional algorithmic latency. As a result, we achieve a low algorithmic latency of $L_C + L_F = 10$ ms.

Cache state management. When processing a stream of continuous chunks, many values can be reused to avoid redundant computations. To achieve this, we maintain and update a cache state h_i at every inference step. As shown in Fig. 3C, our cache state h_i consists of four components: 1) Since the kernel of the first 2D causal convolution layer requires prior chunks, we can avoid recomputing the STFT frames of these prior chunks by caching and reusing them, as illustrated in the figure. 2) Similarly, we store the output of the sequence of dual-path modules from previous chunks, allowing us to compute the causal 2D deconvolution more efficiently. 3) Since computing the ISTFT for the current chunk depends on previous chunks, we maintain a buffer for the intermediate outputs of the 2D deconvolution layer. 4) We store the hidden and cell states of the temporal unidirectional LSTM for each dual-path module. This allows us to fully leverage the long-term receptive field of the recurrent network.

2.2.3 Training for Real-world Generalization. Our speech enhancement system must generalize to complex real-world acoustics, where variations arise from multipath propagation, head-related transfer functions (HRTFs), diverse noise profiles, and motion. To ensure real-world generalization, we adopt a two-step training strategy.

Training. In the first stage, we train our model on the LibriSpeech [36] dataset (360 hours of clean speech). To enhance generalization across reverberation and human wearers, we convolve each speech sample and background noise with a binaural room-impulse-response (BRIR) that captures the acoustic transformations caused by a room and the human head. We use four real-world BRIR datasets—CIPIC[5], RRBRIR [23], ASH-Listening-Set [53], and CATTRIR [24]—and split them into training, validation, and test sets without overlap. We also use WHAM! (58.03 hours of diverse background noise) [66]. For each training sample, we randomly select a 5-second clean speech clip $a(t)$ from LibriSpeech, a BRIR $h_{\theta,\phi}(t)$ from BRIR datasets recorded at azimuthal angle θ and polar angle ϕ , and a 5-second background noise clip $n(t)$ from WHAM. Then we mix them as,

$x(t) = h_{\theta,\phi}(t) * a(t) + n(t)$. We do not convolute $n(t)$ with $h_{\theta,\phi}(t)$, as WHAM! noise is recorded in a binaural format.

During training, our target is $s(t) = h_{\theta,\phi}(t) * a(t)$ to preserve the binaural characteristic of clean speech. We optimize the model using the SNR loss: $L_{SNR}(\hat{s}, s) = \frac{\|s\|^2}{\|s - \hat{s}\|^2}$. Training runs for 200 epochs, each with 20k samples, using AdamW with gradient clipping (0.1). The learning rate follows three sequential schedulers: 1) linearly increases from 1e-4 to 1e-3 over 10 epochs, 2) maintains 1e-3 for 140 epochs, and 3) halves every 15 epochs for the final 50 epochs.

In practice, each device runs its own network, processing audio independently. Thus, we train two separate networks — one for the left channel and one for the right.

Fine-tuning with motion and colored noise. In the previous stage, the sound source and wearer’s head were assumed to be static. However, real-world scenarios involve motion. To address this, we fine-tune the model with time-varying θ and ϕ , following the motion simulation method from [61].

Specifically, the source’s position updates every 25 ms, with a 2.5% probability of triggering a random motion event. When triggered, the angular velocity in the azimuthal and polar directions is sampled from $[\pi/6, \pi/2]$ rad/s, and the source moves at this velocity for a random duration between 0.1 and 1 s, during which no other motion event is triggered. The Steam Audio SDK [51] is used for motion trajectory simulation, creating a mix of stationary and moving segments within the same audio clip. Since BRIR datasets contain recordings at discrete positions, we approximate intermediate positions using nearest-neighbor BRIR selection based on the current $[\theta(t), \phi(t)]$.

We also augment the mixture signal with white, pink, and brown noise, mimicking real-world noise sources like microphone thermal noise and HVAC systems. White noise is generated with a standard deviation sampled from $[0, 0.002)$, while pink and brown noise are created using the Python colorednoise library and scaled by a random factor from $[0, 0.05]$. We fine-tune the model for 100 additional epochs, incorporating motion simulation and noise augmentation. Instead of SNR loss, we use a multi-resolution spectrogram loss [54, 68]. Optimization uses AdamW with a gradient clipping of 0.1, and a ReduceLROnPlateau scheduler (patience = 5, factor = 0.5) starting at 1e-3 learning rate.

2.3 Hardware-Software Co-Design

2.3.1 Hardware Constraints. GAP9’s limited memory and compute resources impose strict constraints on neural networks that it can run efficiently. It has only 1.5 MB of L2 memory and 128 KB of L1 memory. Models exceeding this require external L3 memory, which is slow and power-intensive.

Efficient DSP and neural network inference rely on GAP9’s compute cluster — 9 RISC-V cores and the NE16 accelerator.

The NE16 is optimized for streamed multiply-accumulate operations but only supports up to 8-bit quantized weights and 8-/16-bit activations. The RISC-V cores handle both integer and floating-point operations but are less efficient for neural networks. Thus, the network must be carefully designed by identifying the parts of the network that can be efficiently run on the NE16 with quantization without degrading the output speech quality.

2.3.2 Quantizing Neural Networks. Quantization is the process of converting tensors in the computational graph from a full-precision representation (e.g., FP32) to a lower-precision fixed-point representation (e.g., INT8). This helps reduce memory usage, runtime and power consumption.

We denote the parameters of the network $\mathcal{N}(\cdot)$ as θ , which are originally stored in floating-point precision. Quantization can be classified into two types:

- *Weight quantization:* Applied to network parameters θ .
- *Activation quantization:* Applied to intermediate activation maps generated during inference.

We use uniform INT8 quantization, with asymmetric thresholds for activations and symmetric thresholds for weights.

A uniform quantizer is defined as follows: let $r \in \mathbb{R}^n$ be a vector to be quantized, $S \in \mathbb{R}^+$ be the quantizer scaling factor, $Z \in \mathbb{R}$ be the zero-point and let b be the bit width. The quantization process can be formulated as:

$$Q(r) = \lfloor r/S \rfloor + Z \quad (3)$$

where $Q(r)$ is the fixed point representation of r in quantization space. $\lfloor \cdot \rfloor$ denotes the rounding of the input to the nearest integer value.

Dequantization maps $Q(r)$ back to r , and is given by:

$$\hat{r} = S(Q(r) - Z) \quad (4)$$

Since the recovered values \hat{r} may not exactly match r due to rounding, a **quantization error** is introduced.

A key aspect of uniform quantization is choosing the scaling factor S . It defines how a given range $[\alpha, \beta]$ of real values r is divided into discrete partitions [21]: $S = \frac{\beta - \alpha}{2^b - 1}$. In our implementation, we use the Min-Max Moving Average Observer [16] to determine $[\alpha, \beta]$.

In asymmetric quantization, the range $[\alpha, \beta]$ is not necessarily symmetric with respect to the origin, i.e., $-\alpha \neq \beta$. In symmetric quantization, in contrast, a symmetric clipping range ($-\alpha = \beta$) is used and its scaling factor S is computed as, $S = \frac{\max(-\alpha, \beta)}{2^{b-1} - 1}$ [21], with the zero-point set to $Z = 0$.

Furthermore, in neural networks, quantization is usually applied to a high-dimension tensor (dimension ≥ 2). There are two methods for quantizing higher dimension tensors:

- *Per-tensor quantization:* The entire tensor is quantized using the same S and Z .

- *Per-channel quantization:* Each channel is quantized independently with different S and Z .

Given the constraints of the GAP9 hardware, we deploy per-tensor asymmetric quantization for activations and per-channel symmetric quantization for weights.

2.3.3 Mixed-Precision Configuration. One common method for quantizing a trained floating-point neural network is Post-Training Quantization (PTQ). This involves: 1) calibrating the quantization range using a representative dataset, and 2) quantizing all tensors in the network to a lower-precision format.

We start PTQ by applying full INT8 quantization, but as shown in Table 3, this results in severe performance degradation due to quantization errors. Instead, inspired by [47], we adopt mixed-precision quantization to mitigate performance loss. The key idea is that different network submodules have varying sensitivity to quantization errors. By selectively quantizing certain modules to BFLOAT16 instead of INT8, we can preserve critical information and recover performance from quantization errors. According to [17, 47], input and output quantization is sensitive to the quantization error and leads to performance degradation. The rationale is that the first convolution layer processes raw input data, so maintaining high precision helps retain essential features. The last deconvolution layer reconstructs the final output, where high precision is crucial for generating high-quality audio. Hence, we quantized first input convolution and last deconvolution layers to BFLOAT16 instead of INT8.

2.3.4 Quantization-Aware Training for Better Performance. While mixed-precision PTQ helps mitigate some performance loss, there remains a significant performance drop compared to the floating-point network. This degradation is even more pronounced (Table. 3) in model architectures with a large number of sub-components like ours, which restricts the deployment of advanced models on our accelerator.

To further reduce quantization errors, we apply Quantization Aware Training (QAT). The key idea behind QAT is to simulate quantization errors during training, allowing the model to adjust and converge to a more optimal solution under quantization constraints.

During QAT, floating-point weights and activations are rounded to their quantized equivalents. Since quantization is a non-differentiable operation, we use the Straight-Through Estimator to approximate gradients, where: $\partial \lfloor x \rfloor / x = 1$.

QAT implementation. We describe the details below.

Initialization: We begin with a pretrained floating-point model with parameters θ and apply the mixed-precision quantization described in §2.3.3.

Learned Step-Size Quantization (LSQ): We enhance vanilla QAT with LSQ [18], which allows the network to learn the

quantizer scale factor S using gradients of the task loss. This helps reduce activation quantization errors.

Fine-tuning: The mixed-precision model is fine-tuned for 30 epochs using the FQSE framework [16].

Training Setup: Each epoch processes only 4,000 mixtures due to QAT’s high computational cost. We use an initial learning rate of $1e-3$, with a ReduceLRonPlateau scheduler.

3 Experiments and Results

3.1 Benchmark results

3.1.1 Evaluation Metrics. Our evaluation consists of two main components: speech enhancement evaluation and system evaluation. The former assesses the network’s noise suppression capability and the quality of the enhanced speech.

- *SISDRi:* Scale-Invariant Signal-to-Distortion Ratio (SISDR) is a commonly used metric in speech enhancement tasks to assess the quality of enhanced speech relative to the reference clean speech [26]. SISDR Improvement (SISDRi) can be computed between the input and output speech to assess how much the network improves the quality of the speech compared to the clean reference speech.

- *PESQ:* Perceptual Evaluation of Speech Quality (PESQ) is an objective metric that assesses how closely an enhanced signal matches a reference clean speech signal [43], mimicking human perception.

- *DNSMOS:* Deep Noise Suppression MOS (DNSMOS) is a neural network based objective metric designed to predict Mean Opinion Score (MOS) for speech quality in speech enhancement. We use the ITU P.835 personalized DNSMOS OVRL implementation as our score [42].

For system evaluation, we measure memory consumption, runtime, and power consumption on the AI accelerator.

- *Memory:* GreenWaves’ conversion process reports where neural network parameters are stored. We sum L1 and L2 memory usage to determine total memory consumption.

- *Runtime:* Models run on GAP9 by sending tasks to its compute cluster. Runtime is measured by sending a model execution task, counting elapsed clock cycles until the task completes, and dividing the elapsed cycles by clock frequency. The reported runtime is averaged over 100 iterations.

- *Power:* To measure GAP9’s power consumption during continuous speech enhancement, we power it with a 1.8V source, run inference on each incoming audio chunk every 6 ms, and measure the current over several minutes. Power is calculated as the product of current and voltage.

3.1.2 Model comparison. We compare multiple FP32 models that when fully quantized meet our real-time requirements:

Table 2: Floating-point (FP32) network results.

Model	SISDRi (dB)	PESQ	DNSMOS
TinyDenoiser	6.24 ± 3.43	1.37 ± 0.27	2.06 ± 0.64
TFGridNet-6F	8.43 ± 3.46	1.74 ± 0.49	2.48 ± 0.68
Our model	8.65 ± 3.41	1.76 ± 0.49	2.50 ± 0.66

- *TinyDenoiser:* Prior state-of-the-art TinyDenoiser model [47] performs on-device speech denoising. For fair comparison, we modify its algorithmic latency to 10ms with the same chunk size and padding as our model.

- *TFGridNet-6F:* We use the causal implementation from [61] without self-attention and LayerNormalization modules. Its hyperparameters are $L_B = 6ms$, $L_C = 6ms$, $L_F = 4ms$, $B = 6$, $D = 32$ and $H = 32$. Since TF-GridNet is not real-time even at $4\times$ compression, we apply $6\times$ frequency compression to the spectral module to meet real-time requirements.

- *Our Model:* We set model hyperparameters to $L_B = 6ms$, $L_C = 6ms$, $L_F = 4ms$, $B = 6$, $D = 32$ and $H = 32$. We apply $4\times$ frequency compression on the spectral module.

As shown in Table. 2, our FP32 model achieves the best SISDRi, PESQ and DNSMOS scores at floating model resolution.

3.1.3 Quantization evaluation. Next, we evaluate how different quantization configurations affect performance, runtime, power, and memory consumption. We compare our model with TinyDenoiser [47] under three configurations:

- *BFLOAT16:* We quantize both TinyDenoiser and our model into BFLOAT16 resolution.

- *INT8 PTQ:* We quantize both TinyDenoiser and our model into INT8 resolution, using Post Training Quantization.

- *Mix PTQ:* For TinyDenoiser, we follow the mix-precision configuration as the original paper [47], where we quantize the linear layers into BFLOAT16 and quantize the LSTM layers into INT8. For our model, we quantize the input 2D Conv and output 2D DeConv into BFLOAT16 and quantize other parts into INT8. Then we apply Post Training Quantization on these mixed precision models.

- *Mix QAT:* We follow the same mixed precision quantization configuration for TinyDenoiser and our model, then we apply Quantization-Aware Training.

The results of different quantization configurations are shown in Table 3. While full INT8 quantization achieves the lowest runtime, memory usage, and power consumption for both TinyDenoiser and our model, it significantly degrades speech enhancement performance. TinyDenoiser experiences a 2.76 dB drop, while our model suffers a 10.46 dB decline. The larger performance drop in our model is due to more severe quantization noise accumulation in its deeper and more complex architecture.

Table 3: Power consumption and quantization results. We measure the SISDRi, model size, runtime and power consumption for different Post-Training Quantization (PTQ) and Quantization-Aware Training (QAT) strategies and different models. The BFLOAT16 network configurations cannot run in real-time on the target AI accelerator.

Model	Quantization config	SISDRi (dB)	Memory (kB)	Runtime (ms)	Power (mW)
TinyDenoiser	BFLOAT16	6.30 ± 3.52	–	–	–
TinyDenoiser	INT8 PTQ	3.54 ± 3.34	1135.2	0.53	23.08
TinyDenoiser	Mix PTQ	3.90 ± 3.33	1195.9	0.58	24.12
TinyDenoiser	Mix QAT	5.97 ± 3.35	1195.9	0.58	24.12
Our model	BFLOAT16	8.76 ± 3.41	–	–	–
Our model	INT8 PTQ	-1.70 ± 7.5	280.4	5.19	58.57
Our model	Mix PTQ	0.90 ± 5.40	298.8	5.54	71.64
Our model	Mix QAT	8.19 ± 3.38	298.8	5.54	71.64

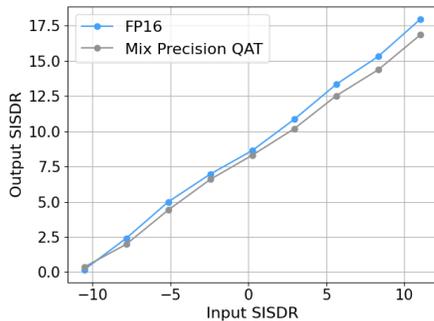


Figure 4: QAT reduces the gap between our floating-point and quantized models across input noise levels.

Applying mixed-precision quantization helps recover 0.36 dB of SISDRi for TinyDenoiser and 2.6 dB for our model. However, the performance loss compared to the floating-point model remains non-negligible. After QAT, TinyDenoiser achieves an SISDRi of 5.97 dB, while our model reaches 8.19 dB, demonstrating the effectiveness of QAT, especially for deeper and more complex networks. Our mixed-precision model with QAT achieves the best performance within hardware constraints. Considering runtime, memory, and power consumption, our model balances performance and memory efficiency while maintaining real-time requirements (< 6 ms) and power constraints (71.6 mW), trading some runtime and power efficiency for improved speech enhancement.

In Fig. 4, we plot the output SISDR as a function of input SISDR for both the floating-point model and the mixed-precision QAT model. The results show that QAT reduces the performance gap between our floating-point and quantized models across all input noise levels. An interesting observation is that the performance gap increases with higher input SISDR. As input SISDR increases, the output speech quality becomes clearer, making quantization noise more dominant and its effects more pronounced compared to scenarios with lower input SISDR [16].

3.1.4 End-to-end hardware runtime evaluation. Fig. 5 shows CDF plots of the hardware runtime for key subcomponents of our audio processing pipeline. In addition to AI inference with our speech enhancement model, the pipeline performs FFT and inverse FFT for time-frequency conversion. Runtime is measured by recording clock cycles before and after the execution of each component and dividing the difference by the clock frequency. The plots are generated by consolidating runtime measurements over 100 iterations. Since we use rectangular synthesis windows (i.e., containing all ones), there is no need to apply a window before the FFT. However, for the inverse FFT, we apply an analysis window (not necessarily all ones) and so include the overlap-add operation in the runtime measurements. Figs. 5a-5c show consistent runtimes, with AI inference being the most time-consuming step. Fig. 5d shows that the full pipeline, including additional tasks such as data copying and type conversion, executes in under 6 ms, ensuring real-time performance.

3.1.5 Wireless throughput. Fig. 6 plots the wireless throughput of our NeuralAid device at different distances from the receiver. To measure throughput, we stream 10,000 packets, each 196 bytes in size, over BLE from the hearing aid to a laptop in a large conference room. The total transmission time is recorded and used to compute the throughput. As expected, throughput decreases as the distance increases due to higher packet loss and retransmissions. Beyond 3 meters, throughput drops below 200 kbps, primarily due to the limited range of the device’s low-profile integrated antenna.

3.2 In-the-Wild Evaluation

We evaluated our system in previously unseen indoor and outdoor environments, using participants who were not included in the training data. It is important to emphasize that our training data consists solely of publicly available datasets, as detailed in §2.2.3, and does not include any data

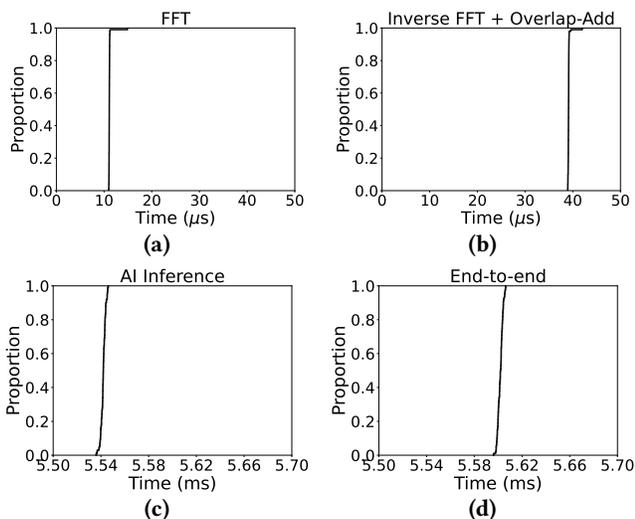


Figure 5: End-to-end hardware run-time evaluation.

collected with our own hardware. The various acoustic environments evaluated are shown in Fig. 7, covering indoor settings such as noisy office spaces, as well as outdoor locations like busy streets with traffic noise and natural environments like parks with multiple noise sources. In indoor spaces, the noise sources included chatter from a nearby event with many people, as well as sounds of coughing, table knocking, and other typical office noises. In outdoor settings, in addition to traffic noise, we also encountered airplane sounds. In all these settings, we had a speaker read a different text in the presence of uncontrolled environmental sounds. The wearer and the speaker were free to move and/or rotate their head and adopted different postures like sitting and standing.

Evaluation procedure. Since the speaker is speaking in the presence of unknown noise, it is difficult to obtain the ground truth audio signal for our speakers in the real world. So, we cannot rely on objective metrics to evaluate the system performance. Instead, we use subjective metrics to allow human participants to rate the audio quality.

For this user study, we recruited 28 participants (19 male and 9 female) ranging in age from 18 to over 70 years, with an average age of 43. The only inclusion criteria was that the participants were adults and could follow English instructions. Each participant evaluated the system in three modes, selected in a random order, across 15 scenarios.

- (1) **No AI:** In this mode, there is no noise suppression and thus the participants hear the unfiltered mixture of the speech signal and background noise.
- (2) **TinyDenoiser:** Prior state-of-the-art TinyDenoiser model.
- (3) **Our method:** In this mode, we use our NeuralAid model fine-tuned with motion and colored noise as described in §2.2.3.

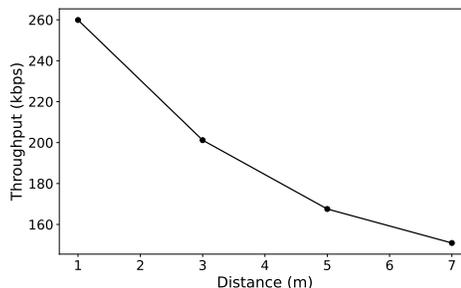


Figure 6: Wireless throughput from the NeuralAid device to a nearby receiver as a function of distance.

For each of these modes, we ask the participants to rate the speech quality by asking them the following questions [61]:

- (1) **Noise suppression:** *How INTRUSIVE/NOTICEABLE were the INTERFERING SPEAKERS and BACKGROUND NOISES? 1 - Very intrusive, 2 - Somewhat intrusive, 3 - Noticeable, but not intrusive, 4 - Slightly noticeable, 5 - Not noticeable*
- (2) **Overall MOS:** *If the goal is to focus on this target speaker, how was your OVERALL experience? 1 - Bad, 2 - Poor, 3 - Fair, 4 - Good, 5 - Excellent*

Results. As shown in Fig. 8, our system significantly reduces background noise, as demonstrated by the increase in the mean opinion score for the noise suppression task from 2.15 in the “no AI” setting to 3.57. Additionally, our NeuralAids framework improved the overall mean opinion score (MOS) from 2.96 to 3.38.

Notably, while TinyDenoiser improves the average noise suppression score from 2.15 in the “no AI” setting to 2.38, it reduces the overall mean opinion score (MOS) from 2.96 to 1.96. This is because, although TinyDenoiser effectively reduces noise, it also significantly degrades speech quality by introducing distortions and choppiness. As a result, understanding the extracted speech becomes more difficult. This outcome aligns with our quantitative benchmarking results from §3.1. The likely reason is that TinyDenoiser, designed to minimize computational complexity, incorporates only a minimal number of components, making it challenging to both suppress noise and preserve speech quality effectively.

An important observation was that the system effectively adapts to sudden and rapid changes in the speaker’s position caused by the wearer rotating their head. Many participants in the study frequently turned their heads to look at different sound sources and objects in their surroundings. To assess this quantitatively, we evaluated the performance of a fine-tuned model across angular velocities ranging from 10 deg/s to 90 deg/s using our synthetic test set. As shown in Table. 4, the fine-tuned system demonstrates slightly better performance than the non-fine-tuned version and exhibits robustness to angular velocities of up to 90 deg/s. Notably, the

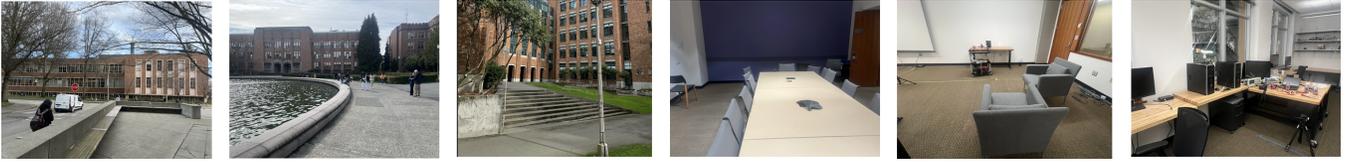


Figure 7: Different outdoor and indoor in-the-wild scenarios. Indoor pictures taken without humans just for this figure.

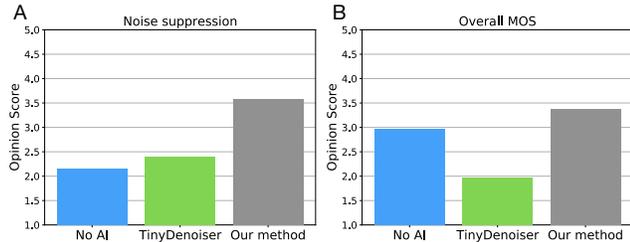


Figure 8: Subjective in-the-wild evaluations. (a) Mean opinion score for the noise suppression quality reported for the three modes of operation, and (b) overall reported mean opinion score on the speech quality.

non-fine-tuned version still performs well on samples with moving speakers despite not being trained on them. This could be due to the fact that this network operates on left and right audio channels independently and it does not use time-difference features which vary drastically with motion.

4 Related Work

AI-based enhanced hearing systems. Recent systems, such as Clearbuds [7], enhance speech of the user wearing the earbuds. However, the processing is not on-device but instead occurs on a smartphone. Further, the target application is telephony, with delay constraints of 100–200 ms, which is an order of magnitude higher than our system.

Enhanced hearing systems like Semantic Hearing [60] allow users to pick and choose which classes of sounds they want to hear (e.g., car honks). This relies on a wired headset, with processing performed on an attached smartphone that has significantly higher processing capabilities. Similarly, prior work on target speech hearing [61] and sound bubbles [10] enable users to hear target speakers based on user-selected characteristics or distance. These require multiple microphones distributed across a wired over-the-head headset, with processing performed on an external platform like an Orange Pi or a Raspberry Pi.

Unlike all these prior works, our system operates on fully wireless hearables and addresses the fundamental question of whether one can design wireless hearables with programmable low-power AI accelerators and if deep learning models for speech processing can be designed to run on these wireless hearables in real time while meeting the size, power, and compute constraints of these wireless hearable platforms.

Table 4: Motion results. Enhancement quality (SISDRi) with and without motion fine-tuning (FT). Results are shown for different ranges of angular speeds (deg/s).

Speed	(-90,-60)	(-60,-30)	(-30,0)	(0,30)	(30,60)	(60,90)
FT	10.71	10.28	10.46	10.44	10.24	10.49
W/o FT	10.42	9.95	10.19	10.18	9.99	10.20

Commercial hearables. Conventional hearables use statistical signal processing for speech enhancement. However, deep neural networks have been demonstrated to achieve superior source separation, outperforming traditional methods by up to 9 dB [7, 29, 56]. Companies like Google [37] have begun integrating AI accelerators, such as the Tensor A1 chip in Pixel earbuds, though specific details remain undisclosed, and its primary use appears to be for active noise cancellation. Similarly, in late 2024, Phonak [2] introduced AI-powered hearing aids, but their deep learning models appear to be heavily duty-cycled to conserve battery life. These commercial solutions remain proprietary, making it unclear whether they use programmable AI or fixed-function integrated circuits (ICs) [57]. This highlights the need for academic exploration into the design space of programmable, low-power speech AI for hearables and a better understanding of the challenges of running streaming deep learning models on hearables—a gap our work addresses.

Hearable platforms. Existing platforms have been designed to enable earable research [9, 11, 19, 41], but none support on-device speech AI acceleration. The Nokia Lab eSense platform [25] introduced sensor-integrated earbuds, enabling data collection and application development for physiological sensing and tracking. OpenEarable [44], OpenEarable 2.0 [45], and ClearBuds [7] further advanced open-source earable platforms. However, the OpenEarable series lacks AI acceleration, relies on a DSP chip and focuses on physiological sensing [31]. ClearBuds, meanwhile, uses a compute-limited microcontroller, making real-time on-device deep learning challenging. Similarly, the OpenMHA platform [39] lacks a hearable or hearing aid form factor and does not include any AI accelerator.

OmniBuds [33] expanded eSense’s bio-sensing capabilities and introduced on-board machine learning for physiological signal processing and classification tasks. However, it has not demonstrated speech AI, which is computationally more demanding. Additionally, it remains closed-source, and

even the specific chip used for on-device ML has not been disclosed – the authors contacted the OmniBuds team, who could not provide this information. Speech AI requires significantly more computational power than classification tasks due to its high sampling rate, extensive I/O requirements, real-time constraints, and the need for causal processing within sub-10 ms latency. Prior to our work, achieving this on low-power hearables was considered challenging.

Low-latency speech processing. Applications have different latency constraints in that they either can wait to process an entire audio file or require much quicker responses. In augmented hearing, minimizing latency between input and processed output is crucial. However, this can degrade performance due to less available information for predicting the output [64]. Prior work has proposed architectures for low-latency speech tasks [49, 50, 58, 59, 61], and some speech challenges focus on these low-latency models [4]. But, these are designed for devices with significantly higher clock frequencies, power budgets, and memory footprints than our target hardware and none of these neural networks meet the computational constraints of wireless hearables.

Deep learning with computational constraints. Common methods include quantization [40], pruning, and knowledge distillation [12, 34, 35]. In the audio domain, these techniques have been applied to tasks such as keyword spotting [70], speaker verification [38], and sound event detection [67]. Prior work has also explored these techniques for speech enhancement and denoising [12, 13, 15, 34, 65]. However, these models are not designed to operate in real-time on wearable hardware. Moreover, they do not process time and frequency components as individual sequences—an essential feature of state-of-the-art enhancement models. [17] proposes a quantized network, but it is neither causal nor real-time.

The closest works to ours are TinyLSTM [20] and TinyDenoiser [46], recurrent neural network (RNN)-based methods for speech enhancement. TinyDenoiser has been shown to run in real-time on GAP9. However, both models have an algorithmic latency of 25 ms. Furthermore, as shown in §3.1.3, when modified for lower latency, performance degrades significantly. In contrast, we present the first dual-path speech enhancement network capable of real-time operation on low-power AI accelerators. We also explore mixed-precision quantization and quantization-aware training, evaluating various trade-offs concerning run-time performance on the target accelerator platform.

5 Limitations and Discussion

The power consumption figures reported in this paper assume that the AI accelerator is continuously active, processing audio at all times. However, in real-world usage, energy consumption and battery life depend on how often the AI

accelerator is actually in use. A more efficient design could activate the AI accelerator only when significant background noise is detected, allowing it to operate in a duty-cycled manner and significantly reducing average energy consumption. The effectiveness of this approach would depend on how frequently the wearer encounters noisy environments throughout the day, and the noise threshold, which is user dependent. Exploring such an adaptive system would be an interesting direction for future research.

Our goal was to demonstrate that effective speech AI models can run on low-power wireless hearables. We address this fundamental challenge using the GAP9 accelerator. However, other low-power AI accelerators are becoming increasingly available including the Analog Devices MAX78000 and MAX78002, Arm Ethos U-55 and U-85, and Kendryte K210, K230, and K510. Investigating their capabilities for running speech AI models would be important for future research.

Recent research has explored techniques with wired headsets like semantic hearing [60], sound bubbles [10], and target speech hearing [61]. While we do not demonstrate the feasibility of implementing these capabilities on wireless hearables, it is worth noting that some of these systems [10, 61] rely on the dual-path model architecture that we optimize in this paper. This suggests a potential path for expanding to these capabilities in the near future.

Our hardware currently utilizes passive noise cancellation, where the earbuds physically block external sounds. However, we can enhance isolation by incorporating active noise cancellation (ANC). This would involve implementing ANC signal processing algorithms, which must meet strict delay requirements. The GAP9 hardware supports ANC algorithms, making integration feasible for future iterations.

Future work could also integrate additional sensors, such as PPG and temperature sensors, to support non-speech applications like physiological sensing [33]. Although our hardware supports multiple microphones per device, our neural networks use only one. Exploring multi-microphone processing on a single wearable and across both ears could enhance performance and warrants further investigation.

Since the inclusion criteria in our user study does not focus on hearing loss, future research is needed to incorporate signal-processing-based personalization algorithms tailored to individuals with hearing loss. This would involve adapting the system based on medical hearing loss prescriptions and systematically evaluating its effectiveness across different levels of hearing impairment.

Finally, our paper focuses on behind-the-ear hearing aid form factor devices. The next step in this research would be integrating these speech AI models into earbud-form-factor devices, which is likely feasible due to the compact size of the target low-power AI accelerator.

6 Conclusion

The emergence of low-power programmable AI accelerators presents an opportunity to bridge state-of-the-art speech AI with wireless earbuds and hearing aids. In this paper, we introduce NeuralAids, a fully on-device speech AI system for real-time speech enhancement and denoising on wireless hearables. Our real-world evaluations highlight the feasibility of deploying advanced speech AI models on low-power wireless hearables, paving the way for next-gen intelligent audio devices that achieve on-device enhanced hearing.

Acknowledgments

The researchers are partly supported by the Moore Inventor Fellow award #10617, UW WE-REACH grant, Thomas J. Cable Endowed Professorship, and a UW CoMotion innovation gap fund. This work was facilitated through the use of computational, storage, and networking infrastructure provided by the UW HYAK Consortium.

References

- [1] 2023. GAP9 processor | GreenWaves Technologies. https://github.com/GreenWaves-Technologies/nn_menu_gap9.
- [2] 2024. AI noise cancelling hearing aid: Phonak audeo sphere™ infinio, Phonak. <https://www.phonak.com/en-us/hearing-devices/hearing-aids/audeo-sphere>
- [3] 2025. NDP120 – Syntiant. <https://www.syntiant.com/ndp120>.
- [4] Michael A. Akeroyd, Will Bailey, Jon Barker, Trevor J. Cox, John F. Culling, Simone Graetzer, Graham Naylor, Zuzanna Podwińska, and Zehai Tu. 2023. The 2nd Clarity Enhancement Challenge for Hearing Aid Speech Intelligibility Enhancement: Overview and Outcomes. In *ICASSP*.
- [5] V.R. Algazi, R.O. Duda, D.M. Thompson, and C. Avendano. 2001. The CIPIC HRTF database. 99–102 pages. doi:10.1109/ASPAA.2001.969552
- [6] Rong Chao, Wen-Huang Cheng, Moreno La Quatra, Sabato Marco Siniscalchi, Chao-Han Huck Yang, Szu-Wei Fu, and Yu Tsao. 2024. An Investigation of Incorporating Mamba for Speech Enhancement. *arXiv* (2024).
- [7] Ishan Chatterjee, Maruchi Kim, Vivek Jayaram, Shyamnath Gollakota, Ira Kemelmacher, Shwetak Patel, and Steven M Seitz. 2022. ClearBuds: wireless binaural earbuds for learning-based speech enhancement. In *MobiSys*.
- [8] Jingjing Chen, Qirong Mao, and Dong Liu. 2020. Dual-path transformer network: Direct context-aware modeling for end-to-end monaural speech separation. *arXiv preprint arXiv:2007.13975* (2020).
- [9] Tao Chen, Xiaoran Fan, Yongjie Yang, and Longfei Shangguan. 2023. Towards Remote Auscultation with Commodity Earphones. In *Proceedings of the 20th ACM Conference on Embedded Networked Sensor Systems* (Boston, Massachusetts) (*SenSys '22*). Association for Computing Machinery, New York, NY, USA, 853–854. doi:10.1145/3560905.3568084
- [10] Tuochao Chen, Malek Itani, Sefik Eskimez, Takuya Yoshioka, and Shyamnath Gollakota. 2024. Hearable devices with sound bubbles. *Nature Electronics* (2024).
- [11] Tao Chen, Yongjie Yang, Xiaoran Fan, Xiuzhen Guo, Jie Xiong, and Longfei Shangguan. 2024. Exploring the Feasibility of Remote Cardiac Auscultation Using Earphones. In *Proceedings of the 30th Annual International Conference on Mobile Computing and Networking* (Washington D.C., DC, USA) (*ACM MobiCom '24*). Association for Computing

Machinery, New York, NY, USA, 357–372. doi:10.1145/3636534.3649366

- [12] Xiuyi Chen, Guangcan Liu, Jing Shi, Jiaming Xu, and Bo Xu. 2018. Distilled Binary Neural Network for Monaural Speech Separation. In *IJCNN*. doi:10.1109/IJCNN.2018.8489456
- [13] Hyeong-Seok Choi, Sungjin Park, Jie Hwan Lee, Hoon Heo, Dongsuk Jeon, and Kyogu Lee. 2021. Real-Time Denoising and Dereverberation with Tiny Recurrent U-Net. In *ICASSP*.
- [14] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014).
- [15] Elad Cohen, Hai Victor Habi, and Arnon Netzer. [n. d.]. Towards Fully Quantized Neural Networks For Speech Enhancement. In *Interspeech 2023*.
- [16] Elad Cohen, Hai Victor Habi, and Arnon Netzer. 2023. Towards fully quantized neural networks for speech enhancement.
- [17] Elad Cohen, Hai Victor Habi, Reuven Peretz, and Arnon Netzer. 2024. Fully Quantized Neural Networks for Audio Source Separation. *IEEE Open Journal of Signal Processing* 5 (2024), 926–933. doi:10.1109/OJSP.2024.3425287
- [18] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. 2019. Learned step size quantization. *arXiv preprint arXiv:1902.08153* (2019).
- [19] Xiaoran Fan, Longfei Shangguan, Siddharth Rupavatharam, Yanyong Zhang, Jie Xiong, Yunfei Ma, and Richard Howard. 2021. HeadFi: bringing intelligence to all headphones. In *Proceedings of the 27th Annual International Conference on Mobile Computing and Networking* (New Orleans, Louisiana) (*MobiCom '21*). Association for Computing Machinery, New York, NY, USA, 147–159. doi:10.1145/3447993.3448624
- [20] Igor Fedorov, Marko Stamenovic, Carl Jensen, Li-Chia Yang, Ari Mandell, Yiming Gan, Matthew Mattina, and Paul N Whatmough. 2020. TinyLSTMs: Efficient neural speech enhancement for hearing aids. *InterSpeech* (2020).
- [21] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. 2022. A survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*.
- [22] Daniel Griffin and Jae Lim. 1984. Signal estimation from modified short-time Fourier transform. *IEEE Transactions on acoustics, speech, and signal processing* 32, 2 (1984), 236–243.
- [23] IoSR-Surrey. 2016. IoSR-surrey/realroombrirs: Binaural impulse responses captured in real rooms. <https://github.com/IoSR-Surrey/RealRoomBRIRs>.
- [24] IoSR-Surrey. 2023. Simulated Room Impulse Responses. <https://iosr.uk/software/index.php>.
- [25] Fahim Kawsar, Chulhong Min, Akhil Mathur, and Alessandro Montanari. 2018. Earables for Personal-Scale Behavior Analytics. *IEEE Pervasive Computing* 17, 3 (2018), 83–89. doi:10.1109/MPRV.2018.03367740
- [26] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R Hershey. 2019. SDR—half-baked or well done?. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 626–630.
- [27] Kai Li and Guo Chen. 2024. SPMamba: State-space model is all you need in speech separation. In *arXiv*.
- [28] Yi Luo, Zhuo Chen, and Takuya Yoshioka. 2020. Dual-Path RNN: Efficient Long Sequence Modeling for Time-Domain Single-Channel Speech Separation. In *ICASSP*.
- [29] Yi Luo and Nima Mesgarani. 2019. Conv-tasnet: Surpassing ideal time–frequency magnitude masking for speech separation. *IEEE/ACM transactions on audio, speech, and language processing* 27, 8 (2019).
- [30] Yi Luo and Nima Mesgarani. 2019. Conv-TasNet: Surpassing Ideal Time–Frequency Magnitude Masking for Speech Separation. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.* (2019).

- [31] Dong Ma, Andrea Ferlini, and Cecilia Mascolo. 2021. OESense: employing occlusion effect for in-ear human sensing. In *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services (Virtual Event, Wisconsin) (MobiSys '21)*. Association for Computing Machinery, New York, NY, USA, 175–187. doi:10.1145/3458864.3467680
- [32] Dirk Mauler and Rainer Martin. 2007. A low delay, variable resolution, perfect reconstruction spectral analysis-synthesis system for speech enhancement. In *2007 15th European Signal Processing Conference*. IEEE, 222–226.
- [33] Alessandro Montanari, Ashok Thangarajan, Khaldoon Al-Naimi, Andrea Ferlini, Yang Liu, Ananta Narayanan Balaji, and Fahim Kawsar. 2024. OmniBuds: A Sensory Earable Platform for Advanced Bio-Sensing and On-Device Machine Learning. arXiv:2410.04775 [cs.ET] <https://arxiv.org/abs/2410.04775>
- [34] Rayan Daod Nathoo, Mikolaj Kegler, and Marko Stamenovic. 2024. Two-Step Knowledge Distillation for Tiny Speech Enhancement. In *ICASSP*.
- [35] Rayan Daod Nathoo, Mikolaj Kegler, and Marko Stamenovic. 2024. Two-Step Knowledge Distillation for Tiny Speech Enhancement. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 10141–10145. doi:10.1109/ICASSP48485.2024.10446796
- [36] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: An ASR corpus based on public domain audio books. In *ICASSP*.
- [37] Bill Park. 2024. Meet Pixel Buds Pro 2, the first buds built for Gemini, Google. <https://blog.google/products/pixel/google-pixel-buds-pro-2/>
- [38] Massimo Pavan, Gioele Mombelli, Francesco Sinacori, and Manuel Roveri. 2024. TinySV: Speaker Verification in TinyML with On-device Learning. (2024).
- [39] Chaslav Pavlovic, Volker Hohmann, Hendrik Kayser, Louis Wong, Tobias Herzke, S. Prakash, zezhang Hou, and Paul Maanen. 2018. Open portable platform for hearing aid research. *The Journal of the Acoustical Society of America* 143 (03 2018), 1738–1738. doi:10.1121/1.5035670
- [40] Antonio Polino, Razvan Pascanu, and Dan Alistarh. 2018. Model compression via distillation and quantization. In *ICLR*.
- [41] Jay Prakash, Zhijian Yang, Yu-Lin Wei, Haitham Hassanieh, and Romit Roy Choudhury. 2020. EarSense: earphones as a teeth activity sensor. In *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking (London, United Kingdom) (MobiCom '20)*. Association for Computing Machinery, New York, NY, USA, Article 40, 13 pages. doi:10.1145/3372224.3419197
- [42] Chandan K A Reddy, Vishak Gopal, and Ross Cutler. 2022. DNSMOS P.835: A Non-Intrusive Perceptual Objective Speech Quality Metric to Evaluate Noise Suppressors. arXiv:2110.01763 [eess.AS] <https://arxiv.org/abs/2110.01763>
- [43] Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra. 2001. Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs. In *2001 IEEE international conference on acoustics, speech, and signal processing. Proceedings (Cat. No. 01CH37221)*, Vol. 2. IEEE, 749–752.
- [44] Tobias Röddiger, Tobias King, Dylan Ray Roodt, Christopher Clarke, and Michael Beigl. 2023. OpenEarable: Open Hardware Earable Sensing Platform (*UbiComp/ISWC '22 Adjunct*).
- [45] Tobias Röddiger, Michael Küttner, Philipp Lepold, Tobias King, Dennis Moschina, Oliver Bagge, Joseph A. Paradiso, Christopher Clarke, and Michael Beigl. 2025. OpenEarable 2.0: Open-Source Earphone Platform for Physiological Ear Sensing. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* (2025).
- [46] Manuele Rusci, Marco Fariselli, Martin Croome, Francesco Paci, and Eric Flamand. 2022. Accelerating RNN-based Speech Enhancement on a Multi-Core MCU with Mixed FP16-INT8 Post-Training Quantization. In *arXiv*.
- [47] Manuele Rusci, Marco Fariselli, Martin Croome, Francesco Paci, and Eric Flamand. 2023. Accelerating RNN-Based Speech Enhancement on a Multi-core MCU with Mixed FP16-INT8 Post-training Quantization. In *Machine Learning and Principles and Practice of Knowledge Discovery in Databases*. Springer Nature Switzerland, Cham, 606–617.
- [48] Kohei Saijo, Gordon Wichern, François G Germain, Zexu Pan, and Jonathan Le Roux. 2024. TF-LoCoformer: Transformer with local modeling by convolution for speech separation and enhancement. In *IWAENC*. IEEE, 205–209.
- [49] Hiroshi Sato, Takafumi Moriya, Masato Mimura, Shota Horiguchi, Tsubasa Ochiai, Takanori Ashihara, Atsushi Ando, Kentaro Shinayama, and Marc Delcroix. 2024. SpeakerBeam-SS: Real-time Target Speaker Extraction with Lightweight Conv-TasNet and State Space Modeling. arXiv:2407.01857 [eess.AS] <https://arxiv.org/abs/2407.01857>
- [50] Hendrik Schröter, Alberto N. Escalante-B., Tobias Rosenkranz, and Andreas Maier. 2022. DeepFilterNet2: Towards Real-Time Speech Enhancement on Embedded Devices for Full-Band Audio. arXiv:2205.05474 [eess.AS] <https://arxiv.org/abs/2205.05474>
- [51] SDK. 2023. Steam Audio. <https://valvesoftware.github.io/steam-audio/>.
- [52] Irtaza Shahid, Yang Bai, Nakul Garg, and Nirupam Roy. 2022. VoiceFind: Noise-Resilient Speech Recovery in Commodity Headphones (*IASA '22*).
- [53] ShanonPearce. 2022. Shanonpearce/ash-listening-set: A dataset of filters for headphone correction and binaural synthesis of spatial audio systems on headphones. <https://github.com/ShanonPearce/ASH-Listening-Set/tree/main>
- [54] Christian J Steinmetz and Joshua D Reiss. 2020. auraloss: Audio focused loss functions in PyTorch. In *Digital music research network one-day workshop (DMRN+ 15)*.
- [55] Cem Subakan, Mirco Ravanelli, Samuele Cornell, Mirko Bronzi, and Jianyuan Zhong. 2021. Attention is All You Need in Speech Separation. In *ICASSP*.
- [56] Cem Subakan, Mirco Ravanelli, Samuele Cornell, Mirko Bronzi, and Jianyuan Zhong. 2021. Attention is all you need in speech separation. In *ICASSP*. IEEE.
- [57] Thierry Tamba, En-Yu Yang, Glenn G. Ko, Yuji Chai, Coleman Hooper, Marco Donato, Paul N. Whatmough, Alexander M. Rush, David Brooks, and Gu-Yeon Wei. 2023. A 16-nm SoC for Noise-Robust Speech and NLP Edge AI Inference With Bayesian Sound Source Separation and Attention-Based DNNs. *IEEE Journal of Solid-State Circuits* (2023).
- [58] Jean-Marc Valin, Umut Isik, Neeraj Phansalkar, Ritwik Giri, Karim Helwani, and Arvindh Krishnaswamy. 2020. A Perceptually-Motivated Approach for Low-Complexity, Real-Time Enhancement of Fullband Speech. arXiv:2008.04259 [eess.AS] <https://arxiv.org/abs/2008.04259>
- [59] Bandhav Veluri, Justin Chan, Malek Itani, Tuochao Chen, Takuya Yoshioka, and Shyamnath Gollakota. 2023. Real-Time Target Sound Extraction. In *ICASSP*.
- [60] Bandhav Veluri, Malek Itani, Justin Chan, Takuya Yoshioka, and Shyamnath Gollakota. 2023. Semantic Hearing: Programming Acoustic Scenes with Binaural Hearables. In *ACM UIST*.
- [61] Bandhav Veluri, Malek Itani, Tuochao Chen, Takuya Yoshioka, and Shyamnath Gollakota. 2024. Look Once to Hear: Target Speech Hearing with Noisy Examples. In *ACM CHI*.
- [62] Zhong-Qiu Wang, Samuele Cornell, Shukjae Choi, Younglo Lee, Byeong-Yeol Kim, and Shinji Watanabe. 2023. TF-GRIDNET: Making Time-Frequency Domain Models Great Again for Monaural Speaker Separation. In *ICASSP*. doi:10.1109/ICASSP49357.2023.10094992

- [63] Zhong-Qiu Wang, Gordon Wichern, Shinji Watanabe, and Jonathan Le Roux. 2022. STFT-domain neural speech enhancement with very low algorithmic latency. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 31 (2022), 397–410.
- [64] Zhong-Qiu Wang, Gordon Wichern, Shinji Watanabe, and Jonathan Le Roux. 2022. STFT-domain neural speech enhancement with very low algorithmic latency. *Trans. on Audio, Speech, and Language Processing* (2022).
- [65] Nils L. Westhausen and Bernd T. Meyer. 2023. Low Bit Rate Binaural Link for Improved Ultra Low-Latency Low-Complexity Multichannel Speech Enhancement in Hearing Aids. In *WASPAA*.
- [66] Gordon Wichern, Joe Antognini, Michael Flynn, Licheng Richard Zhu, Emmett McQuinn, Dwight Crow, Ethan Manilow, and Jonathan Le Roux. 2019. WHAM!:: Extending Speech Separation to Noisy Environments. arXiv:1907.01160 [cs.SD]
- [67] Yushu Wu, Xiao Quan, Mohammad Rasool Izadi, and Chuan-Che Jeff Huang. 2024. “It os Okay to be Uncommon”: Quantizing Sound Event Detection Networks on Hardware Accelerators with Uncommon Sub-Byte Support. In *ICASSP*. 281–285.
- [68] Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. 2020. Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram. In *ICASSP*.
- [69] Lei Yang, Wei Liu, Ruijie Meng, Gunwoo Lee, Soonho Baek, and Han-Gil Moon. 2024. Fspen: an Ultra-Lightweight Network for Real Time Speech Enhancement. In *ICASSP*.
- [70] Yundong Zhang, Naveen Suda, Liangzhen Lai, and Vikas Chandra. 2018. Hello Edge: Keyword Spotting on Microcontrollers. In *arXiv*.