

“IT IS OKAY TO BE UNCOMMON”

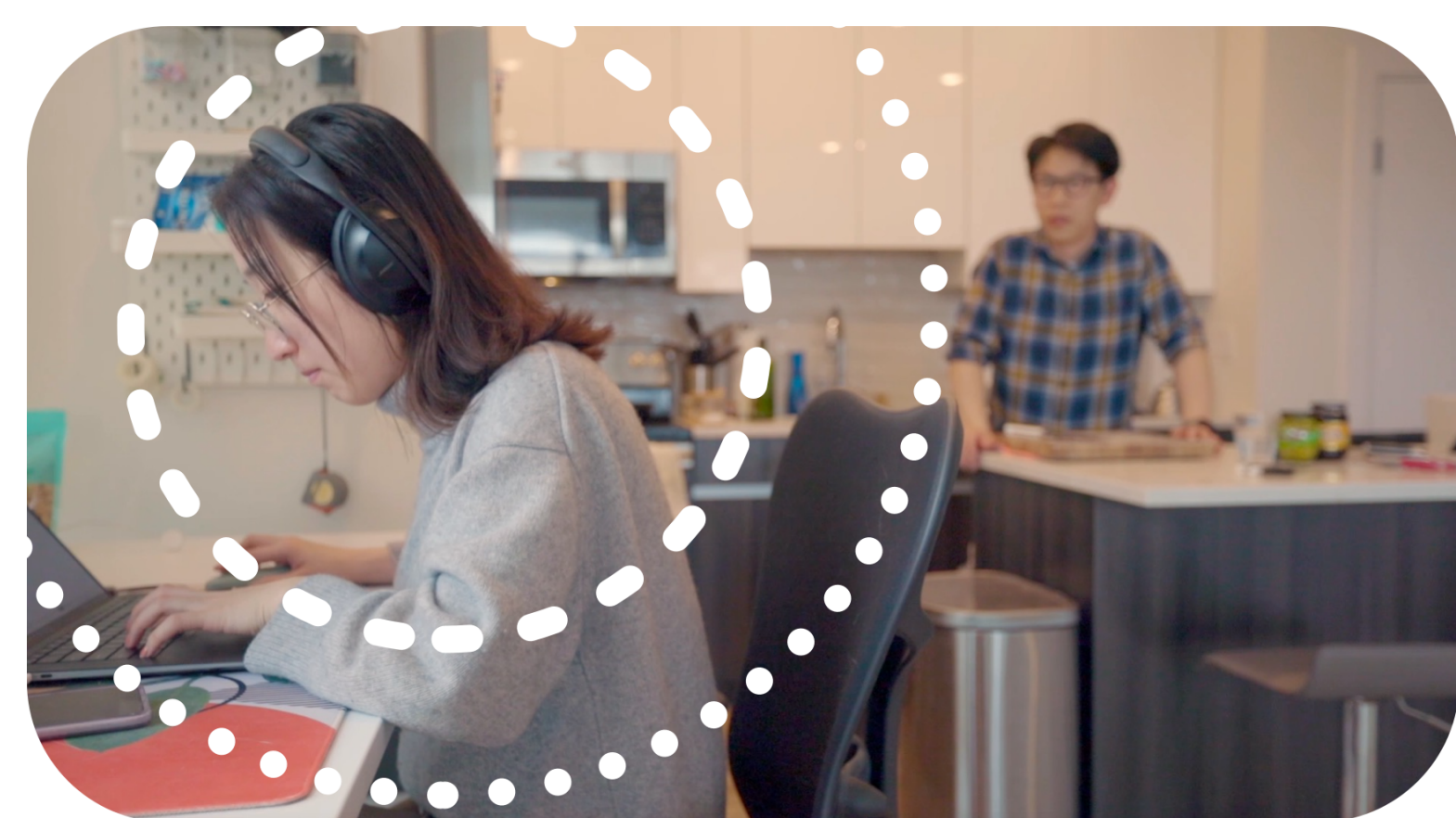
Quantizing Sound Event Detection Networks on Hardware Accelerators with Uncommon Sub-byte Support

Yushu Wu^{1,2}, Xiao Quan¹, Mohammad Rasool Izadi¹, Chuan-Che (Jeff) Huang^{1*}

*chuan-che_huang@bose.com

¹Bose Corporation, USA

²Northeastern University, USA



Context-Aware Headphones

Headphones that understand our audio environments can enable several new user experiences (e.g., inform us of important sound events and adjust audio rendering based on content).

TinyML Challenges

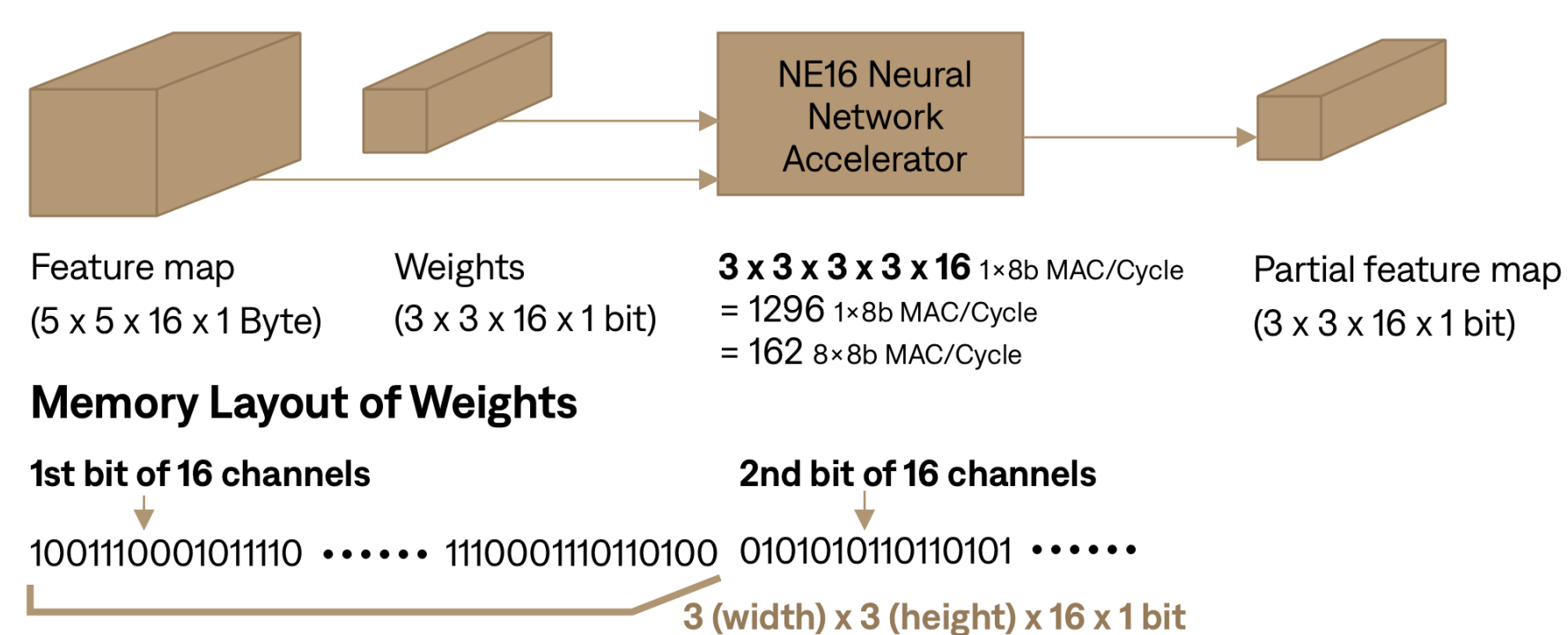
However, running multiple neural networks to understand our audio environments on-device remains a challenging task due to energy and memory constraints.

Contributions

Identify New NN Accelerators

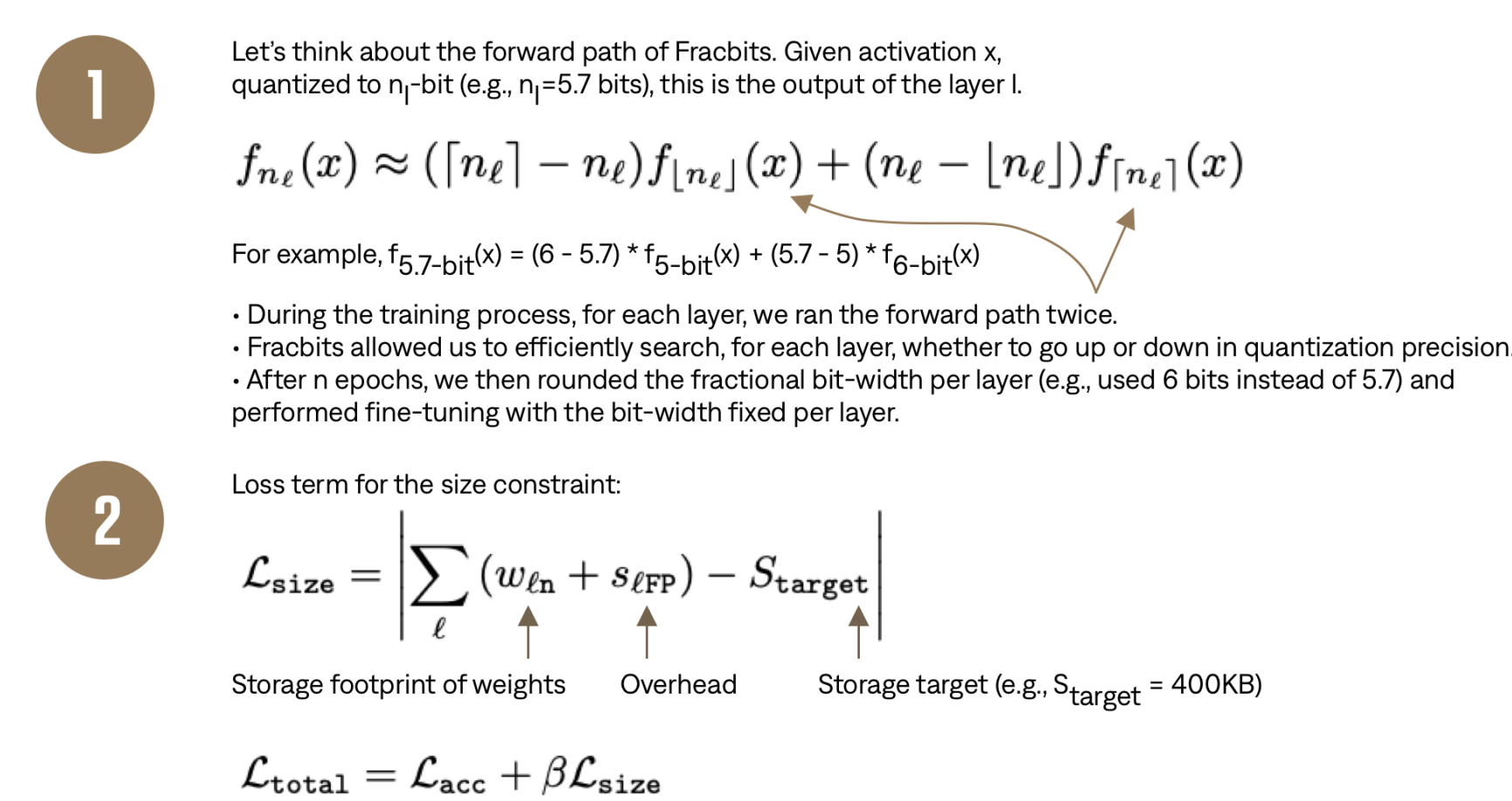
(e.g., NE16 on GAP9 [1]) that support both common (e.g., 4 bit) and uncommon (e.g., 3, 5 bit) sub-byte operations.

Conv2D Example



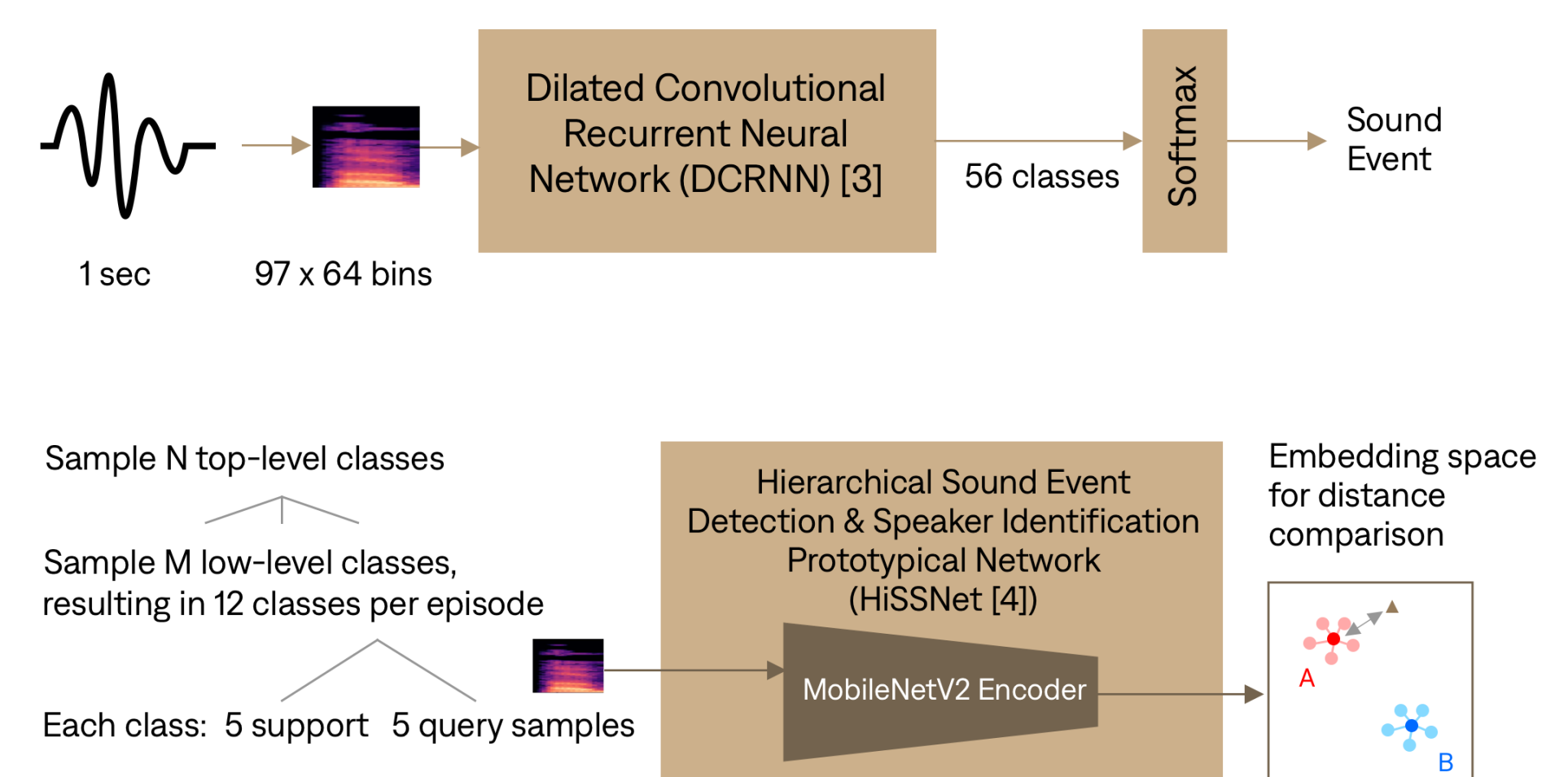
Apply DNAs Over Bit-Width

Apply an efficient differential neural architecture search technique (i.e., Fracbits [2]) to search over the optimal bit-width per layer of a network and evaluate the impact on actual hardware.



Evaluate On Two SED Tasks

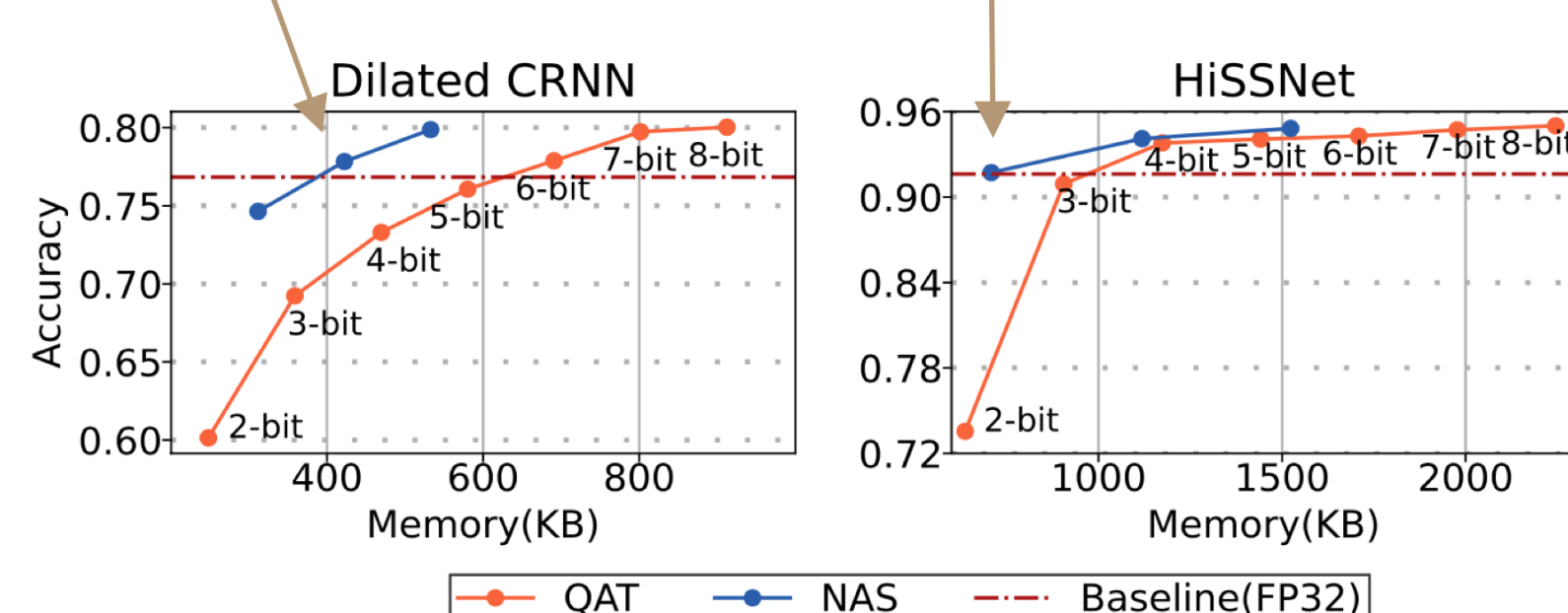
I.e., generic classification and few-shot learning, which potentially have different requirements on quantization granularity.



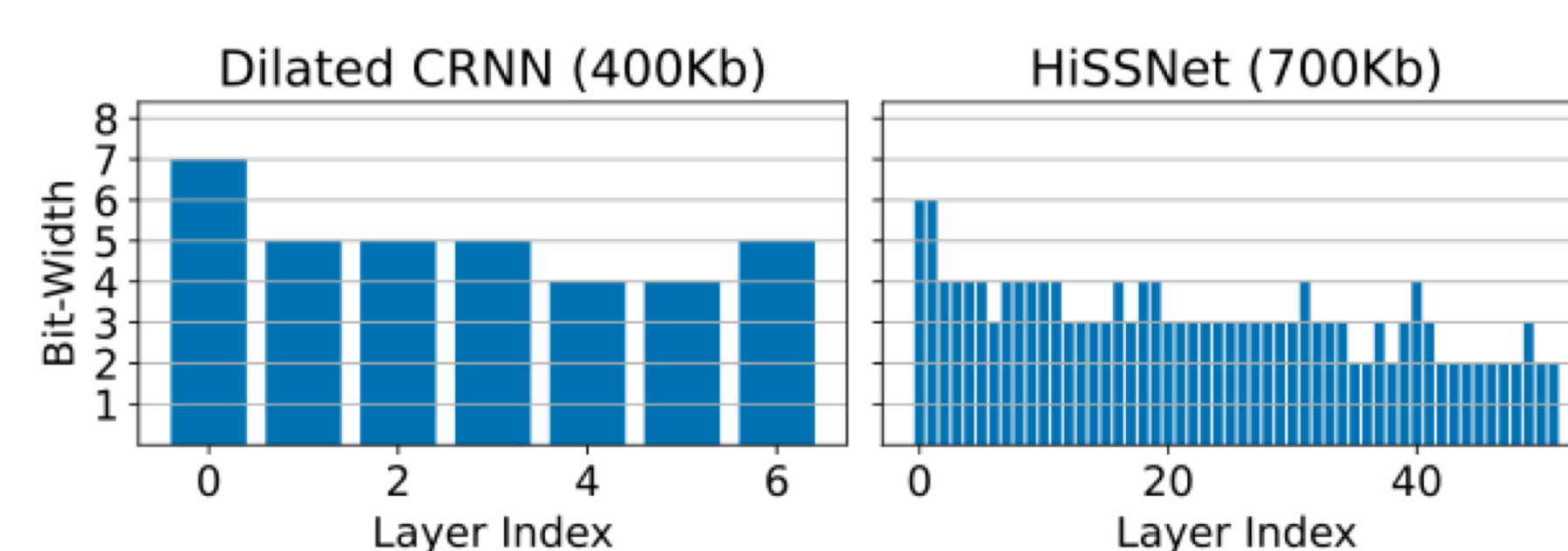
Results

We achieved an average of **62% memory reduction**, **46% latency reduction**, and **61% energy reduction** compared to 8-bit models (trained with quantization-aware training) while maintaining floating point performance.

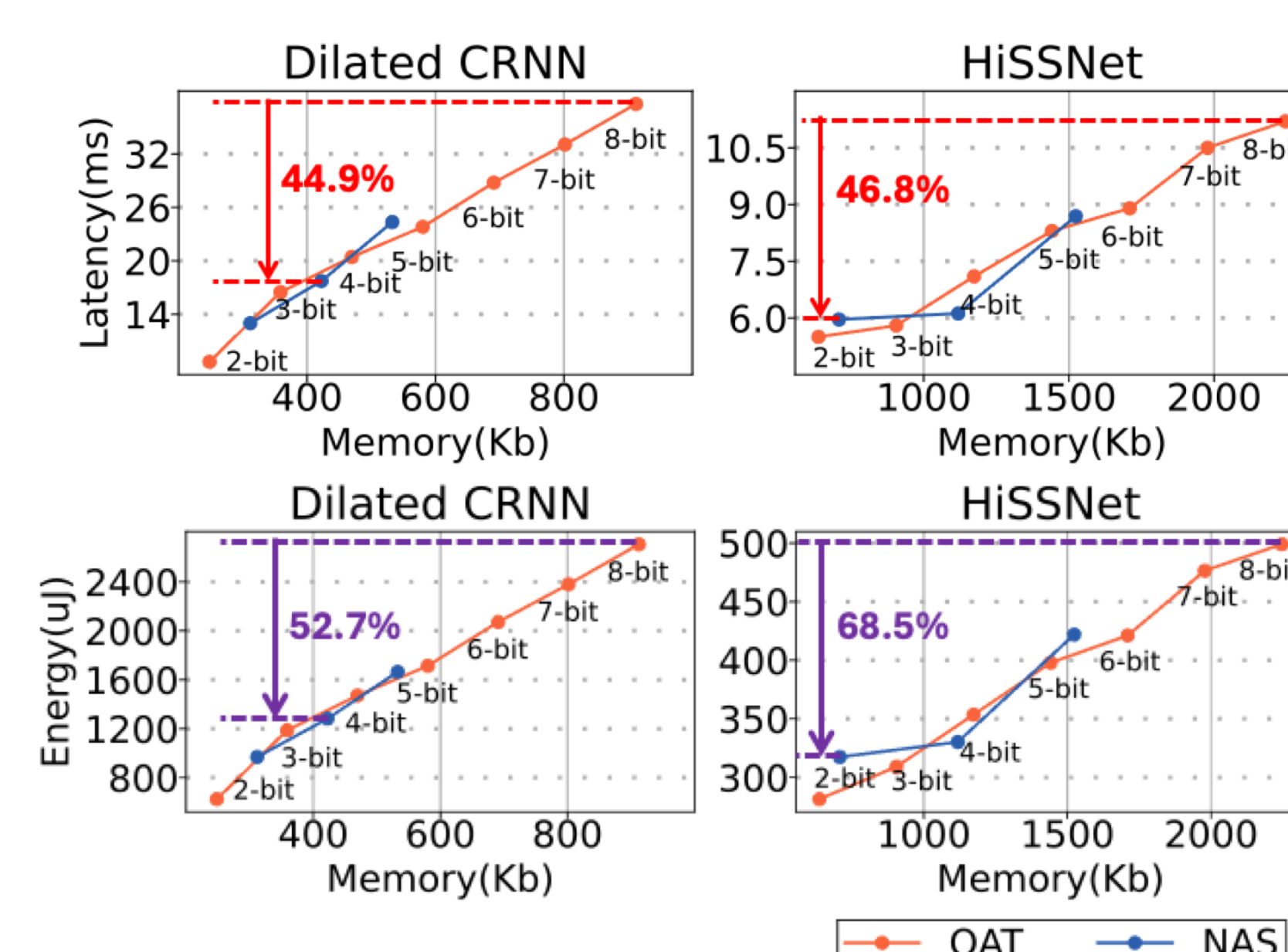
1 Smallest models found that still maintain floating point performance ($S_{\text{target}}=400\text{KB}$ for DCRNN, and $S_{\text{target}}=700\text{KB}$ for HiSSNet)



2 These best models use a variety of bit-widths for different layers -- 7, 5, and 4 for DCRNN and 6, 4, 3, and 2 for HiSSNet.



3 We evaluated these models' energy consumption and latency on an actual hardware (GAP9).



Experimental Setup

Datasets Aggregated monophonic audio recordings from 7 different datasets: ESC50, TUT, TAU, FSD50K, BBC, VCTK and LibriSpeech. 211.6K files, 673.4 hours of audio. 90% for training, 10% evaluation.

Classes 56 classes for generic SED (e.g., emergency alarm) 1263 unique voices of speakers.

We used all of the classes when training few-shot learning models for sound event detection and speaker identification (using HiSSNet) and removed the speaker subset when training for generic SED (using Dilated CRNN).

Training Generic SED We compared quantization-aware training (QAT) using predefined bit-widths and DNAs in this work.

QAT: First, 100 epochs in float, then 10 with QAT
NAS: First, 100 epochs in float, then 5 epochs of Fracbits, then 5 epochs of fine-tuning with the bit-widths fixed.

Training Few-Shot Learning We used a 100-episode, 12-way, 5-shot setup
QAT: First, 1000 epochs in float, then 100 with QAT
NAS: First, 1000 epochs in float, then 50 epochs of Fracbits, then 50 epochs of fine-tuning with the bit-widths fixed.

References

- [1] "GAP9 Product Brief," https://greenwaves-technologies.com/wp-content/uploads/2023/02/GAP9-Product-Brief-V1_14_non_NDA.pdf
- [2] Yang and Jin, "Fracbits: Mixed Precision Quantization Via Fractional Bit-Widths," AACL 2021
- [3] Li et al., "Sound Event Detection Via Dilated Convolutional Recurrent Neural Networks," ICASSP 2020.
- [4] Shashaank et al., "HiSSNet: Sound Event Detection and Speaker Identification via Hierarchical Prototypical Networks for Low-Resource Headphones," ICASSP 2023