# Pruning vs XNOR-Net: A Comprehensive Study of Deep Learning for Audio Classification on Edge-devices

**MD MOHAIMENUZZAMAN[1], CHRISTOPH BERGMEIR[1] AND BERND MEYER[1]**

[1]Department of Data Science and AI, Monash University, Australia (e-mail: md.mohaimen, christoph.bergmeir, bernd.meyer@monash.edu)

Corresponding author: Md Mohaimenuzzaman(e-mail: md.mohaimen@monash.edu).

**ABSTRACT** Deep learning has celebrated resounding successes in many application areas of relevance to the Internet of Things (IoT), such as computer vision and machine listening. These technologies must ultimately be brought directly to the edge to fully harness the power of deep leaning for the IoT. The obvious challenge is that deep learning techniques can only be implemented on strictly resource-constrained edge devices if the models are radically downsized. This task relies on different model compression techniques, such as network pruning, quantization, and the recent advancement of XNOR-Net. This study examines the suitability of these techniques for audio classification on microcontrollers. We present an application of XNOR-Net for end-to-end raw audio classification and a comprehensive empirical study comparing this approach with pruning-and-quantization methods. We show that raw audio classification with XNOR yields comparable performance to regular full precision networks for small numbers of classes while reducing memory requirements 32-fold and computation requirements 58-fold. However, as the number of classes increases significantly, performance degrades, and pruning-and-quantization based compression techniques take over as the preferred technique being able to satisfy the same space constraints but requiring approximately 8x more computation. We show that these insights are consistent between raw audio classification and image classification using standard benchmark sets. To the best of our knowledge, this is the first study to apply XNOR to end-to-end audio classification and evaluate it in the context of alternative techniques. All codes are publicly available on GitHub.

**INDEX TERMS** Sound Classification, Audio Classification, Deep Learning, Model Compression, Filter Pruning, Channel Pruning, XNOR-Net, Edge-AI, Microcontroller, and Image Classification

## I. INTRODUCTION

AUDIO classification is a fundamental building block of many smart Internet of Things (IoT) applications such as predictive maintenance [1]–[3], surveillance [4], [5], and ecosystem monitoring [6], [7]. Smart sensors driven by microcontroller units (MCUs) are at the core of these applications. MCUs, installed at the edge of the networks, sense data and send them to the cloud for classification and detection. However, the energy requirements for transmitting high volumes of data are a burden for battery-powered MCUs. Furthermore, this increases latency, and data transmission may further lead to privacy concerns. The increased latency makes real-time or near real-time analytics infeasible. One way to solve these challenges is to move the analysis and

recognition directly to the edge. In practice, this means that all the processing must take place on resource-impoverished MCUs.

MCUs are low-powered resource-constrained devices typically based on system-on-a-chip (SoC) hardware with less than a megabyte (1 MB) of RAM and below 200 MHz clock speeds. All the recent state-of-the-art audio classification models, however, are based on deep learning (DL) [8]–[16] requiring very resource-intensive computation. Usually, the memory size of such models varies from several MBs to even gigabytes (GB). To run such models on MCUs requires extreme minimization of the model's size and computation requirements with minimum loss of accuracy. Recent studies have applied model compression techniques such

as pruning connections and neurons from fully connected neural networks (FCNNs) [17], filter or channel pruning from convolutional neural networks (CNNs) [18]–[20], knowledge distillation [21] and low-precision quantization [17], [22]. The most recent advancement in extreme downsizing of DL models and their computation requirements is XNOR-Net [23], where a model's activations and inputs are fully binarized.

There are many state-of-the-art XNOR-Net models for different computer vision tasks, such as Rastegari et al. [23] for MNIST, Cong [24] for CIFAR-10, and Bulat et al. [25] for CIFAR-100 and ImageNet datasets. However, for audio classification, the only work we are aware of is presented by Cerutti et al. [26]. This study uses XNOR-Net in combination with spectrogram-based input to effectively perform audio classification via image classification. As the literature on full-precision audio networks shows, this approach does not usually deliver the best results for difficult classification problems with conventional CNN architectures [27], [28] (see Section II-B for further details). To the best of our knowledge, the current literature has not yet considered XNOR-Net for raw audio classification.

Most of the work on XNOR has been performed for computer vision. The leaderboards of benchmark image datasets show considerable differences in the classification accuracy of state-of-the-art full precision nets and XNOR-Nets (see Table 9 for the CIFAR-100 and ImageNet datasets). Furthermore, models produced by XNOR-Net generally yield a up to $32x$ reduction in memory requirements and up to $58x$ reduction in computation cost compared to their full-precision counterparts [23]. This may not result in sufficient reduction for MCUs. The memory requirements of XNOR-Net-based DL models producing comparable accuracy [23]–[25] typically reach several MBs, while typical MCUs offer only 128KB to 1MB of memory.

In this paper, we seek to understand the comparison of XNOR model minimization with pruning-and-quantization approaches and the potential of XNOR in the context of audio classification. We present XNOR-Nets for raw audio classification followed by a comprehensive study that compares this approach with traditional model compression techniques (pruning-and-quantization). Our extensive experimental study reveals that XNOR-Net may be preferred for scenarios comprising a small number of classes (e.g., 10 classes) along with extremely tight resource constraints, specifically where computation speed is concerned. In contrast, for complex scenarios with a larger number of classes, pruning-and-quantization-based compression techniques would still be the choice because they produce sufficiently small models for off-the-shelf MCUs that have higher performance in terms of accuracy.

Experiments show that, when two models are generated by pruning-and-quantization and XNOR-Net for the same memory constraint, XNOR-Net requires $\sim 8x$ less computation while both produce comparable classification accuracy for small problem sizes (number of classes). While

XNOR-Nets require extremely little computation compared to their pruning-and-quantization based counterparts, their classification performance on datasets degrades rapidly with an increasing number of classes. The performance loss of pruning-and-quantization based models for the same audio benchmarks is much more graceful, so that this approach still appears to be preferable for problems with larger class numbers.

To harden the above findings, we have conducted a similar study on image classification datasets and analyzed this in the context of current state-of-the-art full precision and XNOR-Net leaderboards for various image benchmarks. The behavior is found to be consistent in both the audio and the image domains.

Thus, the contribution of this paper is twofold: 1) it presents the first application of XNOR-Net for raw audio classification as a benchmark for future research. 2) It presents the first comprehensive empirical comparison of pruning, quantization, and XNOR-Net-based model compression techniques and derives guidelines for when to use pruning, quantization, and XNOR-Net, respectively.

## II. BACKGROUND AND RELATED WORK

### A. AUDIO REPRESENTATIONS

Audio can be represented in the time domain as a wave form of amplitude changes over time (Figure 1). If this time series is directly used as an input to the network, we speak of raw audio classification.
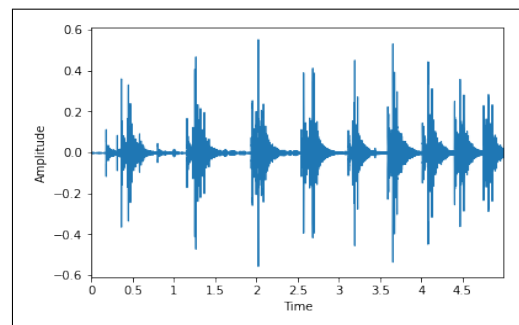


**FIGURE 1.** 1-D representation of audio as wave

An alternative is to use spectrograms. For this, the audio is first transformed into the frequency domain using Fourier transforms of short overlapping windows and presented with respect to time, frequency, and amplitude (Figure 2). The spectrogram captures the intensity of the different frequency component of the signal against time (see Figure 3).
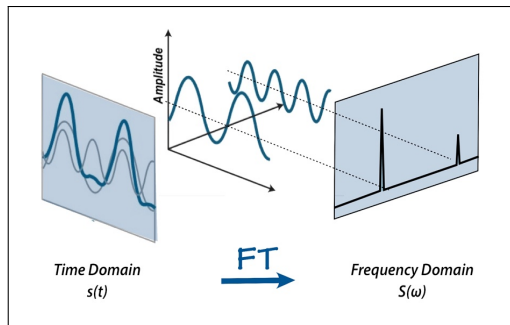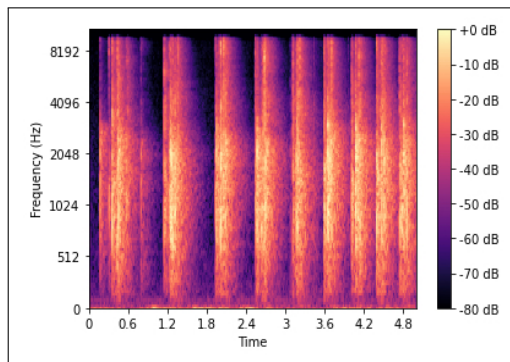
**FIGURE 2.** 2-D representation of audio [29].



**FIGURE 3.** Example spectrogram

## B. AUDIO FEATURE PROCESSING

Input representation is a fundamental decision when applying deep learning to any problem. Owing to the tremendous success of deep learning in the image processing domain, researchers have used spectrograms directly as the input representations [28], [30]–[32]. Somewhat surprisingly, the results achieved so far have not matched the performance that might have been expected based on the state-of-the-art in visual image processing [27], [28].

While being a visual structure, spectrograms have properties different from those of natural images. A pixel of a certain color or the similar neighboring pixels of an image may often belong to the same visual object. On the other hand, although frequencies move together according to a common relationship of the sound, a particular frequency or several frequencies do not belong to a single sound [27], [28]. Furthermore, in images, both axes carry spatial information, but the axes of the spectrograms have different meanings. Sounds are not static two-dimensional objects such as images; they genuinely are time series. Thus, the invariances in the natural images and spectrograms are fundamentally different. For example, moving a face in any direction does not change the fact that it is the same face. However, moving frequencies upwards may not only change an adult voice to a child's voice or something completely different, it may also change the spatial information of the sound [28].

Hence, this research considers audio classification using raw audio time series, an approach that has proven to be successful with traditional full-precision networks [10].

## C. AUDIO CLASSIFICATION AT THE EDGE

The recent state-of-the-art performances for audio classification are all produced by different resource-intensive DL models [10], [13]–[16]. These models must be extremely compressed to be run on the MCUs.

### 1) Model Compression

There are different DL model compression techniques such as unstructured compression [17] (i.e. weight pruning), structured compression [19], [20] (i.e. channel/filter/neuron pruning), knowledge distillation [21], and quantization [22].

Unstructured pruning produces sparse weight matrices that requires sparse computation to fully utilize its benefits. This is not yet supported by the MCUs [33]–[37].

In structured pruning, the neurons and the channels are pruned to find a tiny version of the baseline network. This is an iterative process in which a single iteration performs a global ranking of filters/neurons (using methods such as L2-Norm, Taylor criteria [19] and binary index-based ranking [38]), followed by the removal of the lowest-ranked neuron/channel from the network and retraining of the network for one or two epochs to recover the loss of accuracy due to the pruning [19]. This is repeated until the target number of filters/neurons is removed from the network to find a model of the desired size. This process is illustrated in Figure 4.



**FIGURE 4.** Iterative process of structured model compression

There has been a significant amount of research on structured channel pruning, such as [10], [19], [20], [33], [34], [36], [37]. However, all of this except [10] is proposed for computer vision. The hybrid channel pruning technique proposed in [10] is the only study to date for audio analysis. This method incorporates the Taylor expansion criteria (TE) [19] in the ranking process. In the beginning of an iteration, a forward pass with the entire training dataset takes place and the gradient is applied to the activations to determine the least

affected channels for ranking. The change in gradient can be expressed as:

$$\Theta_{TE}(Z_l^i) = \mid \Delta C Z_l^i \mid \qquad (1)$$

where $Z_l^i$ is the $i$th feature map of layer $l$, $\Delta C$ is the change in loss denoted by $\frac{\delta C}{\delta Z_l^i}$ and $\Theta_{TE}(Z_l^i)$ denotes the change determined by TE. Now, the gradient is applied to the activation as:

$$Z_l^i = Z_l^i + \Theta_{TE}(Z_l^i) \qquad (2)$$

For the ranking of the channels, all the feature maps are normalized layer-wise. Thus, the normalization for a layer $l$ of a network can be expressed as:

$$\bar{Z}_l = \frac{|Z_l^{(i)}|}{\sqrt{\sum |Z_l|^2}} \qquad (3)$$

where $Z_l$ is the list of activations for all the channels $c$ of a layer $l$ in a CNN. Now, the ranking of the channels across all layers is performed and the index of the lowest-ranked channel is determined as follows:

$$i_{lc} = \kappa(\bar{Z}) \qquad (4)$$

where $\bar{Z}$ are the normalized activations of all channels of all layers, $\kappa$ is the function that takes $\bar{Z}$ and returns the information of the channel with the lowest magnitude $i_{lc}$, where $i$ is the index of channel $c$ of layer $l$.

Once this iterative process produces the final compressed and fine-tuned model, it is further compressed using low-precision quantization. Quantization is an independent process, and this study uses 8-bit quantization to achieve a further $4x$ compression.

Knowledge distillation is a very different approach to pruning where the knowledge of a large teacher network is gradually transferred to a smaller student network. However, [10] showed that structured pruning produces superior performance to knowledge distillation in audio classification tasks.

According to the current literature, the best pruning-based approach to derive a tiny model from a state-of-the-art deep neural network (DNN) model is structured pruning. Furthermore, the only work that successfully deploys a DNN model for audio classification to MCUs uses this technique followed by 8-bit post-training quantization to compress the state-of-the-art model [10]. In that study, the compressed model still produces close to that of the state-of-the-art classification accuracy. Thus, for DNN models using pruning-based techniques, this study focuses on structured pruning and quantization. For simplicity, we refer to this as "pruning-and-quantization".

### 2) XNOR-Net
In an XNOR-Net [23], all layers except the first and the last are binary. The input, activations, and the weights of the binary layers are represented using either +1 or -1 and are stored efficiently with single bits. Figure 5 shows the

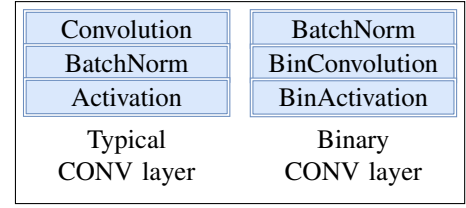construction of a typical convolution layer and an equivalent binary convolution layer.



| Convolution | BatchNorm |
|---|---|
| BatchNorm | BinConvolution |
| Activation | BinActivation |
| Typical CONV layer | Binary CONV layer |

**FIGURE 5.** Typical convolution layer vs binary convolution layer

The convolutions between the matrices (input/activations and weights) are implemented using exclusive-NOR (XNOR) and bit-counting operations. For this, the convolution between two vectors $\in \mathbb{R}^n$ is approximated by the dot products between two vectors $\in \{-1, +1\}^n$ [23]. Thus, convolution between the input $X$ and weight $W$ can be written as:

$$Z \approx (sign(X) \odot sign(W)) \odot \alpha\beta \qquad (5)$$

where $\alpha$ and $\beta$ denote the scaling factors for all the sub tensors in input $X$ and weight $W$ respectively and $\odot$ denotes element-wise multiplication. Due to the binary activations the dot product between $sign(X)$ and $sign(W)$ can be replaced by XNOR and pop-count operations which require extremely little computation. Hence, Equation 5 is reduced to:

$$Z \approx (X' \circledast W') \odot \alpha\beta \qquad (6)$$

where $X' = sign(X)$, $W' = sign(W)$ with all zeros replaced by -1, and $\circledast$ denotes the XNOR and pop-count operations between $X'$ and $W'$. This process is extremely efficient in terms of memory and energy usage. According to Rastegari et al. [23] XNOR-Net requires $32x$ less memory and reduces $58x$ computation requirements. Figure 6 provides an example of how the dot product between $sign(X)$ and $sign(W)$ is replaced by XNOR and pop-count operations between $X'$ and $W'$.
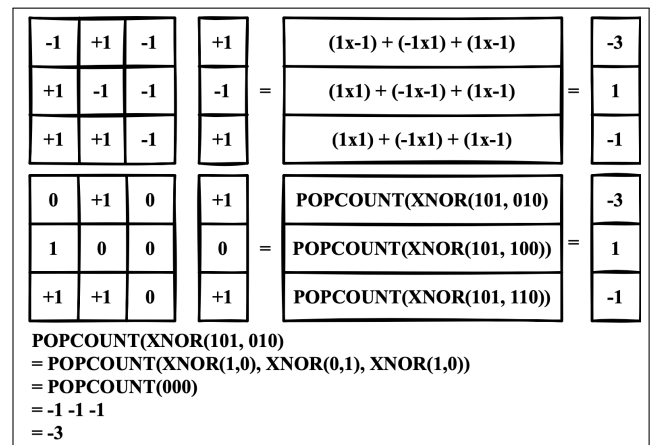


**FIGURE 6.** XNOR and POPCOUNT in XNOR-Net

Like structured pruning based compression works, almost all the works based on XNOR-Net are proposed for computer vision. The current state-of-the-art XNOR models for the benchmark image datasets are: [23] for MNIST, [24] for CIFAR-10 and CIFAR-100 and [25] for ImageNet. In contrast, we have found [26] to be the only XNOR-Net based work for audio classification. However, that work uses spectrograms as the input to the model which is similar to image classification using XNOR-Net. According to the discussion presented in Section II-B, using raw audio as input is the appropriate method for audio classification in MCUs. According to our knowledge, there is no such XNOR-Net based work present in the current literature.

## III. EXPERIMENTAL DETAILS

All experiments are conducted using Python version 3.7.4, and the GPU versions of Torch 1.8.1. All experimental codes are available at: https://github.com/mohaimenz/pruning-vs-xnor.

### A. DATASETS

The experiments are conducted on three widely used audio benchmark datasets: Environmental Sound, namely ESC-50 and ESC-10 [39], and UrbanSound8k [40]. For the image classification experiments, we use the CIFAR-10 and CIFAR-100 benchmark image datasets.

ESC-50 contains 2000 samples that are equally distributed over 50 disjoint classes. The length of the audio samples is 5s recorded at 16kHz and 44.1kHz. Furthermore, the dataset is partitioned into five folds for cross validation to help researchers obtain directly comparable results. ESC-10 is a subset of ESC-50 with 400 audio samples distributed equally over 10 classes. The subsets of the ESC datasets are as follows:

$S_{10}$ = The whole ESC-10 dataset
$S_{20} = S_{10} \cup \{5, 31, 18, 27, 48, 8, 15, 45, 25, 34\}$
$S_{30} = S_{20} \cup \{3, 14, 23, 36, 43, 7, 22, 28, 30, 49\}$
$S_{40} = S_{30} \cup \{6, 9, 16, 17, 24, 29, 32, 35, 37, 44\}$
$S_{50}$ = The whole ESC-50 dataset

UrbanSound8k contains 8732 labeled audio samples of approx. 4s each, recorded at 22.05kHz. The data are pre-sorted into 10 folds and distributed over 10 classes for easy reproduction and comparison of the performances of different algorithms.

The AudioEvent dataset has 28 classes and 5223 samples unevenly distributed across 28 classes. The data is recorded at 16kHz with a bit depth of 16bit. The subsets of the AudioEvent datasets are as follows:

$S_{10} = \{0, 2, 5, 9, 11, 12, 17, 21, 25, 26\}$
$S_{20} = S_{10} \cup \{1, 4, 7, 8, 14, 15, 19, 20, 23, 27\}$
$S_{28}$ = The whole AudioEvent dataset

The CIFAR-10 dataset contains 60,000 image samples equally distributed over 10 classes. 50,000 of the samples are used for training and 10,000 for testing. CIFAR-100 has 100 classes, and 60,000 samples are equally distributed across the classes. For each class, there are 500 training samples and 100 test samples. We create the subsets of CIFAR-100 using Equation 8. The subsets are defined as follows:

$$S_{x_i} = S_{x_{i-1}} \cup \left\{ \frac{x}{10} + m \right\} \quad (7)$$

where $x \in \{10, 20, \ldots, \lfloor n \rfloor\} \mid n \in [1, 100]$ and $m \in [1, 10]$. This gives us the following subsets such as:
$S_{10}$ = The whole CIFAR-10 dataset.
$s_{10} = \{0, 10, 20, 30, 40, 50, 60, 70, 80, 90\}$
$S_{20} = s_{10} \cup \{1, 11, 21, 31, 41, 51, 61, 71, 81, 91\}$
$S_{30} = S_{20} \cup \{1, 11, 21, 31, 41, 51, 61, 71, 81, 91\}$
...
...
$S_{90} = S_{80} \cup \{9, 19, 29, 39, 49, 59, 69, 70, 89, 99\}$
$S_{100}$ = The whole CIFAR-100 dataset.

### B. DATA PREPROCESSING

For the audio datasets (ESC-10, ..., 50 and UrbanSound8k), we train the DL models with samples of length 30,225, i.e., $\approx 1.51s$ audio at 20kHz (for the AudioEvent dataset the length is 51,215, i.e., $\approx 3.2s$ audio at 16kHz). We use data augmentation as described in [15] and [10] for the audio datasets. For the image datasets, we use random cropping, horizontal flipping, and rotation available in the Transforms module of the PyTorch TorchVision library.

All implementations of the data augmentation procedures are available in our GitHub repository.

### C. MODELS AND HYPERPARAMETERS

The model configuration is available in the GitHub repository, and the hyperparameters for all the experiments conducted on the audio and image datasets are listed in Table 1.

**TABLE 1.** Hyperparameter settings for experiments conducted on Audio Datasets (ESC10, ...,50, UrbanSound8k and AudioEvent) and Image datasets (CIFAR-10,...,100). In the loss functions row, KLD stands for KL Divergence, CE for Cross Entropy. In the Optimizer row, SGD stands for Stochastic Gradient Descent and ADAM for Adaptive Momentum Estimation. In the LR Scheduler row, CosAnLR stands for CosineAnnealingLR.

| Datasets→ Networks→ | Audio Datasets | | Image Datasets | |
|---|---|---|---|---|
| | ACDNet, AclNet & their derivatives | | RESNET-18 | |
| Hyperparams↓ | Full Precision | XNOR | Full Precision | XNOR |
| Input shape (ch, h, w) | (1, 1, 30225) | | (3, 32, 32) | |
| Loss function | KLD | | CE | |
| Optimizer | SGD | ADAM | SGD | ADAM |
| Weight decay | 5e-4 | 1e-4 | 5e-4 | 1e-4 |
| Momentum | 0.9 | - | 0.9 | - |
| Initial LR | 0.1 | 0.001 | 0.1 | 0.001 |
| Epochs | 2000 (ESC-10,...,50), 1000 (UrbanSound8k) & 1500 (AudioEvent) | | 400 | |
| LR Scheduler | (0.3, 0.6, 0.9) | CosAnLR | CosAnLR | |
| Warmup epochs | 10 | - | 10 | - |
| LR decay | **0.1x** | - | - | |
| Batch size | 64 | | | |

## IV. ANALYSIS: PRUNING VS XNOR-NET

In this section, we compare pruning-and-quantization techniques with XNOR-Net for raw audio classification. Our

analysis through extensive experiments on different standard benchmark datasets for audio classification shows that XNOR is more effective for relatively simple problems and very tight constraints on the computational resources; however, the higher classification accuracy achieved by pruning-and-quantization outweighs the benefit of XNOR-Net for problems with complex data characteristics, large number of classes, etc.

As our starting point for pruning-and-quantization techniques, we use two of the recent state-of-the-art DL networks for raw audio classification called ACDNet [10] and AclNet [41] as the baseline networks and build their derivatives using pruning-and-quantization as well as the XNOR-Net technique. We test these models on three standard benchmark sound classification datasets - ESC-10 [39], ESC-50 [39], UrbanSound8k [40], and on subsets of them for an in-depth analysis of the effect of increasing the number of classes. We create subsets of a dataset $S_n$ with $n$ classes as:

$$S_x \subseteq S_n \mid x \in \{10, 20, \ldots, \lfloor n \rfloor\} \qquad (8)$$

We measure classification accuracy through 5-fold cross validation. The experimental details are provided in Section III.

### A. PRUNING-AND-QUANTIZATION

We have used the hybrid structured pruning technique proposed in [10] to derive Micro-ACDNet by pruning 80% of the channels from ACDNet. The trained model is first sparsified using the $l0$ norm. It can be expressed as:

$$\hat{Z} = W[\chi(W)] \qquad (9)$$

where $W$ denotes the weights of all the layers of the network, and $\chi$ is the function that sorts $W$ and returns the indices of the bottom 95% of the weights. The channels are then ranked, and the lowest-ranked channel is removed from the network using Equations 1, 2, 3, and 4. The resulting model using this iterative pruning and fine-tuning process is called Micro-ACDNet [10]. Micro-ACDNet is quantized using a post-training 8-bit quantization technique. We refer to this quantized model as QMicro-ACDNet. The memory and computation requirements of QMicro-ACDNet make it suitable for current off-the-shelf MCUs (see Table 2). We follow the same procedure to derive the respective derivatives of AclNet.

**TABLE 2.** Size and computation requirements for ACDNet, AclNet including their Micro and QMicro versions

| Models | Params (M) | | Size (KB) | | FLOPs (M) | |
|---|---|---|---|---|---|---|
| | ACDNet | AclNet | ACDNet | AclNet | ACDNet | AclNet |
| Baseline | 4.74 | 10.63 | 18498.00 | 41512.96 | 544.00 | 1806.57 |
| Micro | 0.13 | 0.13 | 501.00 | 502.00 | 14.82 | 21.50 |
| QMicro | 0.13 | 0.13 | 133.12 | 128.50 | 14.82 | 21.50 |

Micro-ACDNet and QMicro-ACDNet require $36x$ and $144x$ less memory, respectively, than ACDNet that requires 18.06MB to store its 4.74M parameters. Furthermore, both smaller versions require $37x$ less floating point operations (FLOPs) compared to the base model. For AclNet, the memory reductions for the same derivatives are $86x$ and $324x$

with $84x$ less FLOPs. We note that QMicro version requires the same number of FLOPs as the Micro version but only at quarter precision. If a full 32-bit floating-point unit (FPU) is used, this does not make a practical difference, but it allows us to exploit significant speed gains using 16-bit FPUs, 16-bit FPU modes on full-precision FPUs, and software emulations when no hardware FPU is available.

Although the smaller models require fewer resources, they lose classification accuracy because they have less capacity to learn. According to [10], Micro-ACDNet has 80% less capacity than ACDNet and QMicro-ACDNet is a quarter precision version (8-bit) of Micro-ACDNet. Table 3 and Figure 7 present a comparison of the accuracy achieved by the three versions of the both the baseline networks on ESC-10,...,50 datasets (derived using Equation 8. For further details, see Section III-A). The table and the figure show that all the three versions of both the networks produce state-of-the-art and near state-of-the-art accuracy on ESC-10, however, the accuracy drops continuously with an increase in the number of classes, as may be expected.

**TABLE 3.** Prediction accuracy (%) of Micro versions of ACDNet and AclNet including their quantized versions (QMicro) on ESC-10,...,50 datasets. The column headers '#Cls' stands for 'No. of Classes' and 'Baseline' stands for the base model (i.e., ACDNet or AclNet).

| #Cls. | ACDNet | | | AclNet | | |
|---|---|---|---|---|---|---|
| | Baseline | Micro | QMicro | Baseline | Micro | QMicro |
| 10 | 96.75 | 96.25 | 92.75 | 95.75 | 94.00 | 90.75 |
| 20 | 92.38 | 90.13 | 82.55 | 91.12 | 88.50 | 80.25 |
| 30 | 89.25 | 86.83 | 81.67 | 87.81 | 84.75 | 79.34 |
| 40 | 85.94 | 81.56 | 75.71 | 85.00 | 79.10 | 73.81 |
| 50 | 87.05 | 83.25 | 75.50 | 85.70 | 80.05 | 72.25 |

The base ACDNet achieves an accuracy of 96.75% and 87.05% on ESC-10 and ESC-50, respectively, whereas Micro-ACDNet produces 96.25% and 83.25% on ESC-10 and ESC-50, respectively. The quantized version experiences a larger drop in accuracy as the number of classes increases, starting at 92.75% for ESC-10 and ending with 75.50% for ESC-50. For AclNet and its different smaller versions, the trend is similar if not the same. Figure 7 clearly shows these trends.

**TABLE 4.** Prediction accuracy (%) of Micro versions of ACDNet and AclNet including their quantized versions (QMicro) on UrbanSound8k dataset. Here, 'Baseline' stands for the base model (i.e., ACDNet or AclNet)

| Models | Accuracy (%) | |
|---|---|---|
| | ACDNet | AclNet |
| Baseline | 84.45 | 79.17 |
| Micro | 78.28 | 75.80 |
| QMicro | 70.93 | 67.47 |

The scenario is not different when we apply the same networks on the UrbanSound8k dataset. Table 4 shows that QMicro-ACDNet has lost almost 13.5% accuracy compared to the base network. For QMicro-AclNet, the loss in accuracy is 11.7% compared to the base network. We note that specialized quantization targeted to a particular model can potentially reduce the loss of the accuracy that the
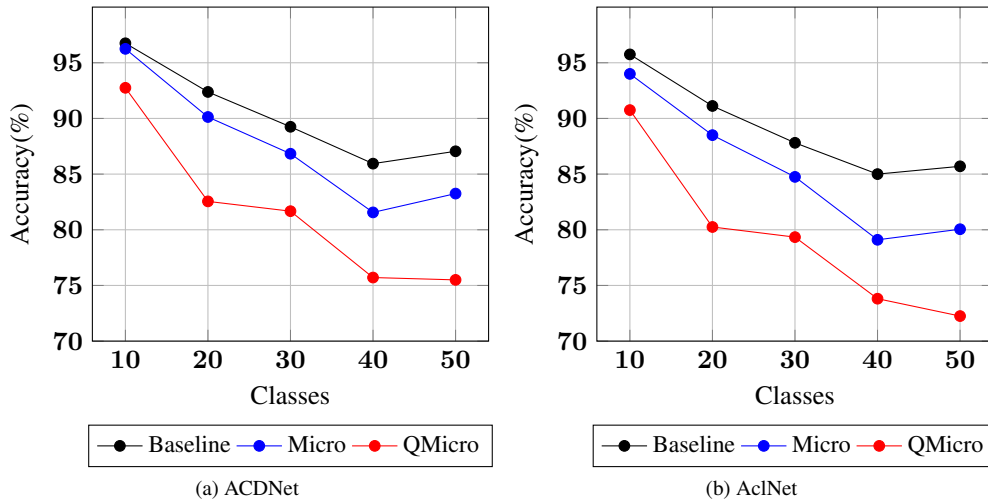
**FIGURE 7.** Comparison of prediction accuracy between the baseline models (i.e., ACDNet and AclNet) and the micro versions including their quantized versions (QMicro) on ESC-10,. . . ,50 datasets.

models experience during the quantization process. However, improving any technique used to demonstrate the process is beyond the scope of this paper and while a change in quantization techniques may shift the results, the general trends are expected to remain the same.

### B. XNOR-NET

We have applied XNOR-Net technique to derive the XNOR-Net versions of ACDNet (i.e., XACDNet) and AclNet (i.e., XAclNet). Table 5 lists the memory and the computation requirements of the full precision networks and their XNOR counterparts. The table shows that the memory required by the XACDNet is 578KB and XAclNet is 1300KB. These figures are still too much for typical MCUs offering less than 1 MB of RAM, since the parameter spaces alone already leave hardly any room for the actual computations. Hence, we need smaller versions of the original networks whose resource requirements do not exceed the resources available in the MCUs. To compare the network performance, memory, and computation requirements with QMicro-ACDNet, we create Mini-ACDNet such that its XNOR-Net version has similar requirements to QMicro-ACDNet (see Tables 2 and 5). To derive Mini-ACDNet, we use the same technique as that used to derive Micro-ACDNet, summarized above and fully detailed in [10]. The same procedure is used to derive Mini-AclNet from AclNet.

**TABLE 5.** Size and computation requirements for the baseline, Mini and their XNOR-Net versions XBaseline and XMini respectively. Here 'Baseline' stands for the base networks (i.e., ACDNet or AclNet)

| Model | Params (M) | | Size (KB) | | FLOPs (M) | |
|---|---|---|---|---|---|---|
| | ACDNet | AclNet | ACDNet | AclNet | ACDNet | AclNet |
| Baseline | 4.74 | 10.63 | 18498 | 41512.96 | 544 | 1806.57 |
| XBaseline | 4.74 | 10.63 | 578.00 | 1300.48 | 8.88 | 30.02 |
| Mini | 1.05 | 1.05 | 4096.00 | 4096.00 | 112.00 | 119.00 |
| XMini | 1.05 | 1.05 | 128.5 | 133.12 | 2.79 | 3.20 |

The sizes of the XNOR-Net versions of ACDNet and

AclNet are calculated according to [23]. The number of binary operations required for the networks is calculated according to [42]. We express the computation (Binary operation + FLOPS) required for an XNOR network as FLOPs for simplicity. If $a$ and $b$ are the FLOPs required for the first and the last full precision layers of the XNOR network and $x$ is the total number of FLOPs of the full precision version of the network, then the calculation of FLOPs of an XNOR network can be expressed as $FLOPs = a + (x - (a+b))/64 + b$.

Table 5 also shows that XMini versions are extremely small (128.5KB and 133.12KB) and require considerably less computation. This clearly allows the network to be deployed on the current off-the-shelf MCUs. Furthermore, from Table 6, we can see that the XNOR networks produce reasonable accuracy for a smaller number of classes (ESC-10 with 10 classes). However, as the number of classes increases, the network performance of the XNOR versions decreases rapidly.

**TABLE 6.** Prediction accuracy (%) of Baselines and Mini versions including their XNOR-Net versions on ESC-10,...,50 datasets. The column header '#Cls' stands for 'No. of Classes' and 'Baseline' stands for the base network (i.e., ACDNet or AclNet)

| #Cls. | ACDNet | | | | AclNet | | | |
|---|---|---|---|---|---|---|---|---|
| | Baseline | XBaseline | Mini | XMini | Base | XBase | Mini | XMini |
| 10 | 96.75 | 91.25 | 96.75 | 82.25 | 95.75 | 89.70 | 93.50 | 80.50 |
| 20 | 92.38 | 77.25 | 91.12 | 54.75 | 91.12 | 66.00 | 88.38 | 52.75 |
| 30 | 89.25 | 70.42 | 88.58 | 46.83 | 87.81 | 57.33 | 85.6 | 48.00 |
| 40 | 85.94 | 61.87 | 84.50 | 38.56 | 85.00 | 48.00 | 81.44 | 35.60 |
| 50 | 87.05 | 56.40 | 85.60 | 31.70 | 85.70 | 39.85 | 82.80 | 31.55 |

Figure 8 shows the performance graph of the Baseline models (i.e., ACDNet and AclNet) and their mini, pruning-and-quantization and XNOR counterparts. The line at the bottom (XMini) shows how extremely the smallest XNOR-Net is affected when the number of classes increases, while QMicro manages degrade much more gracefully. For all the XNOR-Net networks (XBaseline and XMini versions of the

Baselines), the slope is much steeper than that for the other compression methods.

Figure 9 shows the performance of the same networks on the UrbanSound8k dataset, another widely used audio benchmark. From the graph, we observe that the results for UrbanSound8k confirm our findings for ESC-10,...,50.

In summary, we observe that the memory requirements of QMicro and XMini versions are essentially the same. XMini-ACDNet requires $7.97x$ less FLOPs than QMicro-ACDNet, on the other hand, XMini-AclNet requires $6.72x$ less FLOPs that QMicro-AclNet. However, the XNOR-based models do not reach comparable classification accuracies when the number of classes is larger. Although the pruning-and-quantization-based QMicro-ACDNet also sees a loss in accuracy, the accuracy is still reasonable for larger problems given its minuscule model size. This observation is consistent with both the baseline models and their derivatives.

### C. EXTENDED ANALYSIS

#### 1) Comparing with Existing Work

To the best of our knowledge, Cerutti et al. [26] have presented the only XNOR network for audio classification so far, termed BNN-GAP8. They have used a different benchmark, namely, the less commonly used AudioEvent dataset [43]. We cannot test BNN-GAP8 on the most widely used benchmarks, ESC50 and Urbansound-8k, since BNN-GAP8 has not been described in sufficient detail in [26] to make it reproducible.

To facilitate a direct comparison, we extend our analysis to this dataset, which has also been used in a several other studies [43], [44]. These results are consistent with what we have seen above for the most widely used standard benchmarks. Table 7 and Figure 10 show the performance of our base nets on the AudioEvent dataset and its smaller subsets.

**TABLE 7.** Prediction accuracy (%) of the Baseline, Mini including their XNOR-Net versions and Micro including its quantized version (QMicro) on AudioEvent-10,20,28 datasets. The column header '#Cls' stands for 'No. of Classes' and the Baseline model is ACDNet or AclNet. XMini and QMicro versions have approximately the same memory requirements.

| #Cls.→ | 10 | | 20 | | 28 | |
|---|---|---|---|---|---|---|
| Models↓ | ACDNet | AclNet | ACDNet | AclNet | ACDNet | AclNet |
| Baseline | 96.25 | 94.05 | 92.82 | 90.46 | 92.57 | 90.15 |
| XBaseline | 78.05 | 77.08 | 53.27 | 50.84 | 49.83 | 41.34 |
| Mini | 95.86 | 93.45 | 92.93 | 90.82 | 92.49 | 89.99 |
| XMini | 74.16 | 73.33 | 53.48 | 50.09 | 43.40 | 44.61 |
| Micro | 95.66 | 94.06 | 90.03 | 88.35 | 89.69 | 87.51 |
| QMicro | 94.08 | 92.48 | 86.71 | 84.28 | 84.98 | 81.57 |

However, the larger resource requirements of our base networks do not allow us a direct and fair comparison to the BNN-GAP8 network presented in [26]. Therefore, we derive two additional smaller networks (QNano-ACDNet and XGAP8-ACDNet) with memory requirements equivalent to BNN-GAP8. QNano-ACDNet is an 8-bit quantized version of the associated full-precision network Nano-ACDNet that is derived by pruning the channels from ACDNet. XGAP8-ACDNet is an XNOR network derived from the full precision network GAP8-ACDNet. GAP8-ACDNet is also derived by

pruning channels from ACDNet. The equivalent derivatives are also derived from AclNet. The new models are constructed using the same procedures as described above. The tests of these networks on the AudioEvent dataset used in [26] are summarized in Table 8.

**TABLE 8.** Prediction accuracy (%) of BNN-GAP8 [26], Nano and GAP8 versions of both the baselines including the XNOR-Net for GAP8 (i.e., XGAP8) and the quantized version of Nano (i.e., QNano) on AudioEvent (28 classes) dataset

| Models | Accuracy (%) | | Memory (KB) | | FLOPs (M) | |
|---|---|---|---|---|---|---|
| CNN-CNP [44] | 85.1 | | 1766 | | 2478 | |
| BNN-GAP8 [26] | 77.9 | | 58 | | 1768 | |
| | ACDNet | AclNet | ACDNet | AclNet | ACDNet | AclNet |
| Nano | 86.22 | 84.75 | 245 | 245 | 10.54 | 11.37 |
| QNano | 82.34 | 80.15 | 61 | 61 | 10.54 | 11.37 |
| GAP8 | 90.35 | 87.57 | 1960 | 1960 | 40.51 | 54.02 |
| XGAP8 | 35.48 | 35.93 | 61 | 61 | 1.60 | 2.18 |

BNN-GAP8 shows good performance but the quantized networks QNano-ACDNet and QNano-AclNet clearly outperforms it for equivalent storage requirements. The computational requirements for BNN-GAP8 are given in MACs rather than FLOPs in [26]. Using the same re-scaling as detailed above, we arrive at 1768M FLOPs for BNN-GAP8. The quantized QNano networks compare very favorably, using only 10.54M FLOPs for the ACDNet version and 11.37M FLOPs for the AclNet version. In a nutshell, the pruned-and-quantized QNano versions deliver better performance than BNN-GAP8 XNOR-Net with smaller resource requirements.

The performance of BNN-GAP8 on the 28-class problem is much better than that of the the XNOR networks derived from ACDNet and AclNet (i.e., XGAP8-ACDNet and XGAP8-AclNet). The results for the latter are consistent with the performance drop established for the slightly larger XMini-ACDNet on ESC30 (ESC50 reduced to 30 classes, cf. Table 6).

How can BNN-GAP8 achieve such good performance on a 28-class problem? The likely reason is that BNN-GAP8 benefits from using spectrograms as input to the network. This means that in the BNN-GAP8 implementation much of the necessary full precision computation required, for which a binary net is not suitable, is encapsulated outside of the network in the conversion of raw audio to spectrograms. The other XNOR networks in this comparison are deprived of this possibility because they perform end-to-end classification for raw audio input and thus must perform all required computations within the network.

From a pragmatic perspective, using spectrograms may be a useful approach for handling simple audio classification problems on MCUs that feature a DSP with hardware support for FFT computation, which is the basis of generating spectrograms. However, as has been shown in Table 8 and Figure 11, even where such support exists, forgoing it and using end-to-end classification with a pruned-and-quantized network instead still offers performance and resource benefits over the XNOR network. Furthermore, from the broader literature on DNN audio classification, it is evident that purely spectrogram-based classification does not always allow us to achieve state-of-the-art performance with conventional
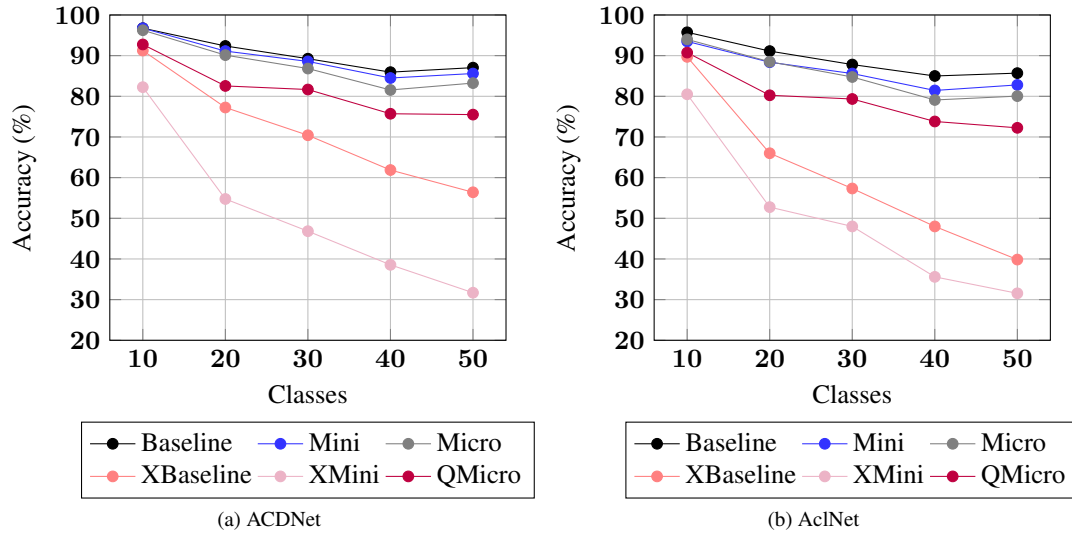
**FIGURE 8.** Comparison of prediction accuracy between Baseline, Mini including their XNOR-Net versions and Micro including its quantized version (QMicro) on ESC-10,...,50 datasets. The Baseline network is AcdNet or AclNet. XMini and QMicro versions have approximately the same memory requirements.
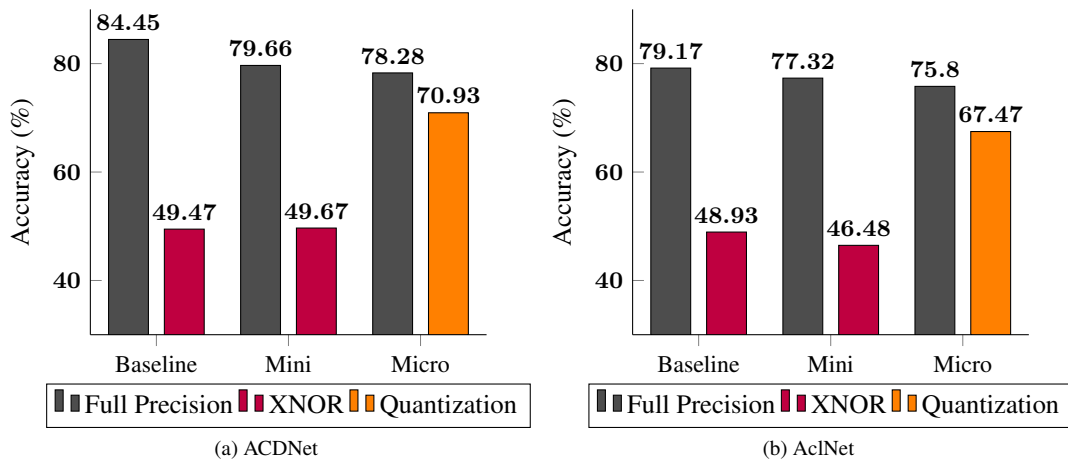


**FIGURE 9.** Comparison of prediction accuracy between Baseline, Mini including their XNOR-Net versions and Micro including its quantized version (QMicro) on the UrbanSound8k dataset. The Baseline network is ACDNet or AclNet. XMini and QMicro versions have approximately the same memory requirements.

CNN architectures and that features learned from raw audio or multi-channel features are preferable for more difficult benchmarks [8], [9], [11], [12], [16].

### 2) Analysis with Image Datasets

Most of the work on XNOR nets so far has taken place in the image domain. A comparison with this body of work is instructive. We summarize the state-of-the-art for the most widely used datasets in Table 9. The accuracy achieved by XNOR nets is impressive (second last column), but it has to be noted that these nets are significantly larger than our target size. None of these models fit on the relevant MCUs with the single exception of the one for MNIST. This is a very simple dataset with only 10 classes.

**TABLE 9.** State-of-the-art accuracy (%) for various image datasets. '#Cls' stands for 'number of Classes'.

| Datasets (#cls.) | Full Precision | | XNOR | |
|---|---|---|---|---|
| | Accuracy | Size (MB) | Accuracy | Size (MB) |
| MNIST (10) | 99.87 [45] | 5.8 | 99.23 [23] | 0.10 |
| CIFAR-10 (10) | 99.50 [46] | 2411 | 88.74 [24] | 1.3 |
| CIFAR-100 (100) | 96.08 [47] | - | 77.80 [25] | 7.8 |
| Imagenet (1000) | 90.45 [48] | 7030 | 71.20 [25] | 7.8 |

To verify whether our own results for audio classification are consistent with the image domain, we conducted additional experiments on image classification using XNOR. We have used the XNOR version of RESNET-18 [49] to classify the widely used benchmark image datasets CIFAR-10 and CIFAR-100. We have created variably sized subsets of CIFAR-100 to determine whether the trend of the loss in accuracy is similar to audio classification. The experimental details are provided in Section III. Table 10 and Figure 12
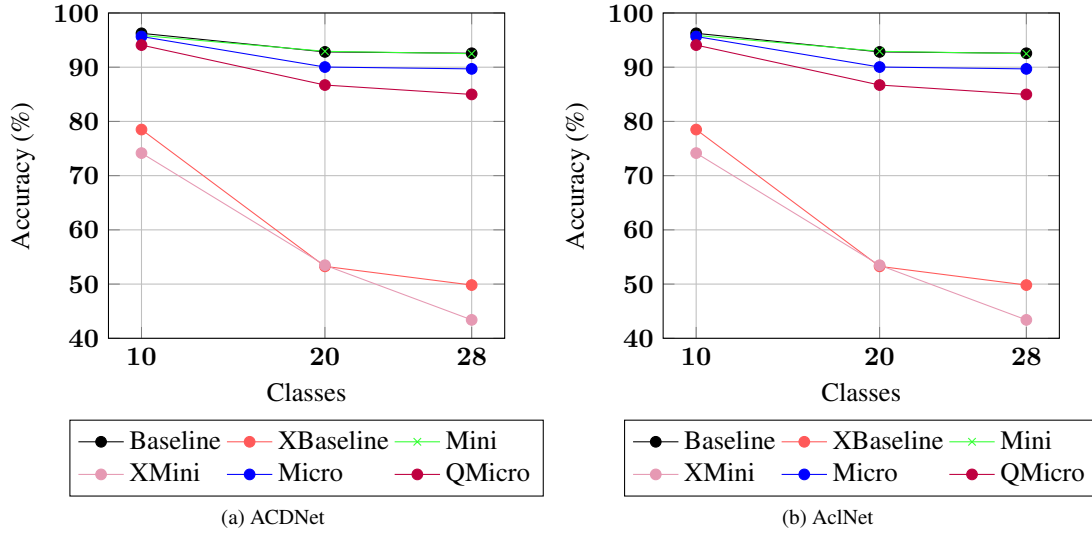
**FIGURE 10.** Comparison of prediction accuracy between the Baseline, Mini including their XNOR-Net versions and Micro including its quantized version (QMicro) on AudioEvent datasets and its subsets (10, 20 and 28 classes). The Baseline model is ACDNet or AclNet. XMini and QMicro versions have approximately the same memory requirements.
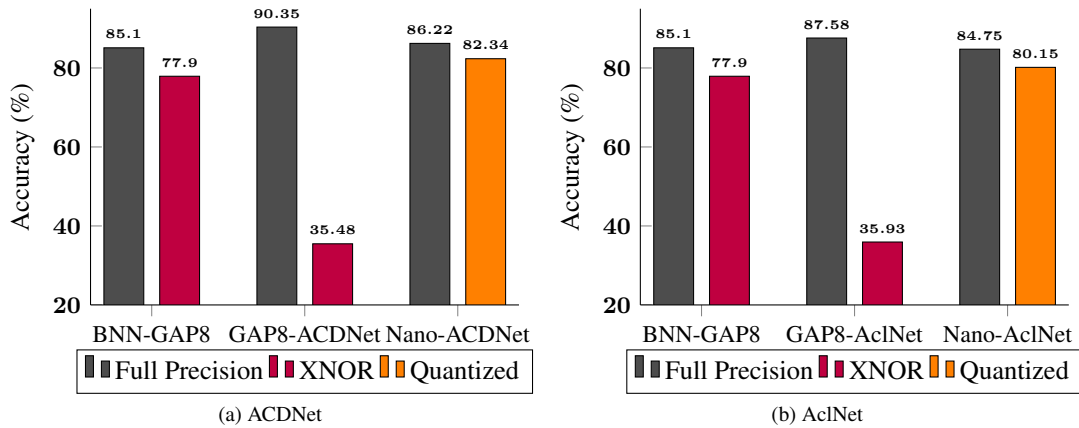


**FIGURE 11.** Comparison of prediction accuracy between BNN-GAP8 [26], GAP8-ACDNet and Nano-ACDNet on AudioEvent dataset for 28 classes. BNN-GAP8 runs on spectrogram, in contrast, GAP8-ACDNet and Nano-ACDNet run on raw audio. GAP8-ACDNet has approximately the same resource requirements of BNN-GAP8. QNano-ACDNet has equivalent memory requirement of the XNOR versions of BNN-GAP8 and XGAP8-ACDNet

summarize the results. The sizes of the full precision models are between 42.63MB and 42.80MB and the computation requirements between 95.17M FLOPs and 95.21M FLOPs. For the XNOR-Net version, the model size is approximately 1.34MB using 2.06M FLOPs. The results of this experiment are consistent with those for the audio domain, and it is clearly visible that a similar performance drop for XNOR occurs as the number of classes increases.

**TABLE 10.** RESNET-18 and its XNOR-Net version on CIFAR-10,20,30,....,100 datasets.

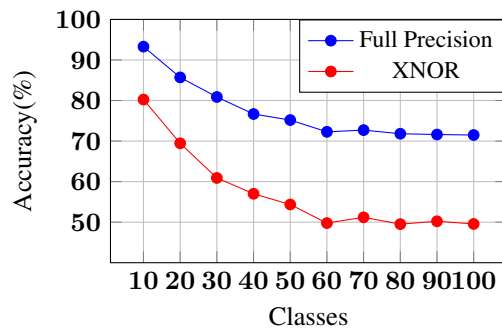| Datasets | Full Precision | XNOR |
|---|---|---|
| | Accuracy (%) | Accuracy (%) |
| CIFAR-10 | 93.29 | 80.24 |
| CIFAR-20 | 85.70 | 69.45 |
| CIFAR-30 | 80.87 | 60.90 |
| CIFAR-40 | 76.67 | 57.00 |
| CIFAR-50 | 75.18 | 54.36 |
| CIFAR-60 | 72.28 | 49.77 |
| CIFAR-70 | 72.70 | 51.20 |
| CIFAR-80 | 71.81 | 49.51 |
| CIFAR-90 | 71.61 | 50.22 |
| CIFAR-100 | 71.49 | 49.56 |

**FIGURE 12.** Prediction accuracy of RESNET-18 and its XNOR-Net version on CIFAR-10,20,30,....,100 datasets

## V. CONCLUSION

This paper presents the first study of XNOR-Net for end-to-end classification of raw audio. We emphasize state-of-the-art MCUs as practically important target architectures for realistic Edge-AI applications.

For small problem sizes, measured in the number of classes, our comprehensive experimental analysis shows that XNOR-Net can produce models with reasonable performance that are sufficiently small to fit the target architectures. For our test scenarios with small problem sizes, XNOR networks require significantly less computation than comparably sized models generated by pruning-and-quantization alone. For relatively simple problems, the performance of XNOR-Nets on MCUs can be further increased by using spectrograms as input to the net. This enables us to capitalize on special DSP hardware so that the full-precision computations required to handle raw audio can be confined outside of the net in the spectrogram generation.

The picture changes as the problem complexity grows and larger numbers of classes need to be distinguished. Our analysis shows that XNOR-Nets face large drops in classification accuracy for datasets with many classes. This performance degradation is much more rapid than that of their pruning-and-quantization counterparts. This performance degradation can be handled by growing the network size. Although growing the network size may still allow the model to be run on CPUs, it takes away the ability to fit the model in MCUs because XNOR-Net can only reduce the memory requirement by a maximum of 32-fold comparing to its full-precision version. Hence, for larger models, it is necessary to use other model compression techniques to find a model that is at most 32 times bigger than the target size before XNOR is applied.

In contrast, comparably sized models generated by pruning-and-quantization show significantly better performance while still satisfying the constraints of the target architectures. The advantages of spectrogram-based network input can no longer be exploited with complex data characteristics as feature-learning from raw audio or multi-channel input are required to reach state-of-the-art performance. This renders pruning-and-quantization the preferred approach from a certain problem complexity unless computation speed rather than model size dominates the decision.

Furthermore, to the best of our knowledge, there is no off-the-shelf computation kernel for XNOR-Nets yet. This means that it is difficult to realize the theoretical advantage of XNOR-Net on existing add-multiplication-based hardware and that custom hardware is required to achieve the full benefit of faster computation [50]. However, given the popularity of XNOR-Nets, we hope that such support is not too far away.

Our study provides an experimental analysis of the current state of XNOR-Net and alternative compression methods as a basis to evaluate possible paths towards competitive deep learning architectures in edge-AI applications. The creation of specialized computation kernels for XNOR-Net and theoretical studies of XNOR-Net performance degradation for larger class sizes and how to avoid it are still open research problems and important areas for future research.

## REFERENCES

[1] Matthias Auf der Mauer, Tristan Behrens, Mahdi Derakhshanmanesh, Christopher Hansen, and Stefan Muderack. Applying sound-based analysis at porsche production: Towards predictive maintenance of production machines using deep learning and internet-of-things technology. In Digitalization Cases, pages 79–97. Springer, 2019.

[2] Feng Jia, Yaguo Lei, Liang Guo, Jing Lin, and Saibo Xing. A neural network constructed by deep learning technique and its application to intelligent fault diagnosis of machines. Neurocomputing, 272:619–628, 2018.

[3] Huitaek Yun, Hanjun Kim, Eunseob Kim, and Martin BG Jun. Development of internal sound sensor using stethoscope and its applications for machine monitoring. Procedia Manufacturing, 48:1072–1078, 2020.

[4] Roneel V Sharan and Tom J Moir. An overview of applications and advancements in automatic sound recognition. Neurocomputing, 200:22–34, 2016.

[5] Weitao Xu, Xiang Zhang, Lina Yao, Wanli Xue, and Bo Wei. A multi-view cnn-based acoustic classification system for automatic animal species identification. Ad Hoc Networks, 102:102115, 2020.

[6] Dan Stowell, Tereza Petrusková, Martin Šálek, and Pavel Linhart. Automatic acoustic identification of individuals in multiple species: improving identification across recording conditions. Journal of the Royal Society Interface, 16(153):20180940, 2019.

[7] Xiao Yan, Hemin Zhang, Desheng Li, Daifu Wu, Shiqiang Zhou, Mengmeng Sun, Haiping Hu, Xiaoqiang Liu, Shijie Mou, Shengshan He, et al. Acoustic recordings provide detailed information regarding the behavior of cryptic wildlife to support conservation translocations. Scientific reports, 9(1):1–11, 2019.

[8] Yuan Gong, Yu-An Chung, and James Glass. AST: Audio Spectrogram Transformer. arXiv preprint arXiv:2104.01778, 2021.

[9] Loris Nanni, Gianluca Maguolo, Sheryl Brahnam, and Michelangelo Paci. An ensemble of convolutional neural networks for audio classification. Applied Sciences, 11(13):5796, 2021.

[10] Md Mohaimenuzzaman, Christoph Bergmeir, Ian Thomas West, and Bernd Meyer. Environmental sound classification on the edge: Deep acoustic networks for extremely resource-constrained devices. arXiv e-prints, pages arXiv–2103, 2021.

[11] Jaehun Kim. Urban sound tagging using multi-channel audio feature with convolutional neural networks. Proceedings of the Detection and Classification of Acoustic Scenes and Events, 2020, page http://dcase.community/documents/challenge2020/technical_reports/DCASE2020_JHKim_21_t5.pdf, 2020.

[12] Anurag Kumar and Vamsi Ithapu. A sequential self teaching approach for improving generalization in sound event recognition. In Proceedings of the International Conference on Machine Learning, ICML 2020, pages 5447–5457. PMLR, 2020.

[13] Zhichao Zhang, Shugong Xu, Shunqing Zhang, Tianhao Qiao, and Shan Cao. Learning attentive representations for environmental sound classification. IEEE Access, 7:130327–130339, 2019.

[14] Yu Su, Ke Zhang, Jingyu Wang, and Kurosh Madani. Environment sound classification using a two-stream CNN based on decision-level fusion. Sensors, 19(7):1733, 2019.

[15] Yuji Tokozume, Yoshitaka Ushiku, and Tatsuya Harada. Learning from between-class examples for deep sound recognition. In Proceedings of the 6th International Conference on Learning Representations, ICLR 2018. https://openreview.net/forum?id=B1Gi6LeRZ, 2018.

[16] Hardik B Sailor, Dharmesh M Agrawal, and Hemant A Patil. Unsupervised filterbank learning using Convolutional Restricted Boltzmann Machine for environmental sound classification. In Proceedings of the 18th Annual Conference of the International Speech Communication Association, Interspeech 2017, pages 3107–3111, 2017.

[17] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In Proceedings of the 4th International Conference on Learning Representations, ICLR 2016. https://arxiv.org/abs/1510.00149, 2016.

[18] Xiaolong Ma, Geng Yuan, Sheng Lin, Zhengang Li, Hao Sun, and Yanzhi Wang. ResNet can be pruned 60×: Introducing network purification and unused path removal (P-RM) after weight pruning. In Proceedings of the IEEE/ACM International Symposium on Nanoscale Architectures, NANOARCH 2019, pages 1–2. IEEE, 2019.

[19] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. In Proceedings of the 5th International Conference on Learning Representations, ICLR 2017. https://openreview.net/forum?id=SJGCiw5gl, 2017.

[20] Oyebade Oyedotun, Djamila Aouada, and Bjorn Ottersten. Structured compression of deep neural networks with debiased elastic group LASSO. In Proceedings of the IEEE Winter Conference on Applications of Computer Vision, 2020, pages 2277–2286, 2020.

[21] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. stat, 1050:9, 2015.

[22] Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and quantization. In Proceedings of the 6th International Conference on Learning Representations, ICLR 2018. https://openreview.net/forum?id=S1XolQbRW, 2018.

[23] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. XNOR-Net: Imagenet classification using binary convolutional neural networks. In Proceedings of the European conference on computer vision, ECCV 2016, pages 525–542. Springer, 2016.

[24] Cong Wang. Pytorch XNOR-Net: XNOR-Net, with binary gemm and binary conv2d kernels, support both CPU and GPU. https://github.com/cooooorn/Pytorch-XNOR-Net. Accessed: Jun 15, 2021.

[25] Adrian Bulat, Brais Martinez, and Georgios Tzimiropoulos. High-capacity expert binary networks. In Proceedings of the 9th International Conference on Learning Representations, ICLR 2021. https://openreview.net/forum?id=MxaY4FzOTa, 2021.

[26] Gianmarco Cerutti, Renzo Andri, Lukas Cavigelli, Elisabetta Farella, Michele Magno, and Luca Benini. Sound event detection with binary neural networks on tightly power-constrained IoT devices. In Proceedings of the ACM/IEEE International Symposium on Low Power Electronics, ISLPED 2020, pages 19–24, 2020.

[27] Daniel Rothmann. What's wrong with spectrograms and CNNs for audio processing? https://towardsdatascience.com/whats-wrong-with-spectrograms-and-cnns-for-audio-processing-311377d7ccd, Mar 2018.

[28] L Wyse. Audio spectrogram representations for processing with convolutional neural networks. In Proceedings of the First International Conference on Deep Learning and Music, 2017, pages 37–41, 2017.

[29] Kartik Chaudhary. Understanding audio data, fourier transform, FFT, spectrogram and speech recognition. https://towardsdatascience.com/understanding-audio-data-fourier-transform-fft-spectrogram-and-speech-recognition-a4072d228520, Jun 2021.

[30] Prateek Verma and Julius O Smith. Neural style transfer for audio spectograms. arXiv preprint arXiv:1801.01589, 2018.

[31] Muhammad Irfan, Zheng Jiangbin, Shahid Ali, Muhammad Iqbal, Zafar Masood, and Umar Hamid. Deepship: An underwater acoustic benchmark dataset and a separable convolution based autoencoder for classification. Expert Systems with Applications, 183:115270, 2021.

[32] Akon O Ekpezu, Isaac Wiafe, Ferdinand Katsriku, and Winfred Yaokumah. Using deep learning for acoustic event classification: The case of natural disasters. The Journal of the Acoustical Society of America, 149(4):2926–2935, 2021.

[33] Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung. Structured pruning of deep convolutional neural networks. ACM Journal on Emerging Technologies in Computing Systems (JETC), 13(3):32, 2017.

[34] Ariel Gordon, Elad Eban, Ofir Nachum, Bo Chen, Hao Wu, Tien-Ju Yang, and Edward Choi. MorphNet: Fast & simple resource-constrained structure learning of deep networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, pages 1586–1595, 2018.

[35] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In Proceedings of the 5th International Conference on Learning Representations, ICLR 2017. https://openreview.net/forum?id=rJqFGTslg, 2017.

[36] Jian-Hao Luo, Hao Zhang, Hong-Yu Zhou, Chen-Wei Xie, Jianxin Wu, and Weiyao Lin. ThiNet: pruning CNN filters for a thinner net. IEEE transactions on pattern analysis and machine intelligence, 41(10):2525–2538, 2019.

[37] Pravendra Singh, Vinay Kumar Verma, Piyush Rai, and Vinay Namboodiri. Leveraging filter correlations for deep model compression. In Proceedings of the IEEE Winter Conference on Applications of Computer Vision, 2020, pages 835–844, 2020.

[38] Jian-Hao Luo and Jianxin Wu. AutoPruner: An end-to-end trainable filter pruning method for efficient deep model inference. Pattern Recognition, 107:107461, 2020.

[39] Karol J. Piczak. ESC: Dataset for Environmental Sound Classification. In Proceedings of the 23rd Annual ACM Conference on Multimedia, 2015, pages 1015–1018. ACM Press, 2015.

[40] Justin Salamon, Christopher Jacoby, and Juan Pablo Bello. A dataset and taxonomy for urban sound research. In Proceedings of the 22nd ACM international conference on Multimedia, 2014, pages 1041–1044. ACM, 2014.

[41] Jonathan J Huang and Juan Jose Alvarado Leanos. AclNet: efficient end-to-end audio classification CNN. arXiv preprint arXiv:1811.06669, 2018.

[42] Zechun Liu, Baoyuan Wu, Wenhan Luo, Xin Yang, Wei Liu, and Kwang-Ting Cheng. Bi-Real Net: Enhancing the performance of 1-bit CNNs with improved representational capability and advanced training algorithm. In Proceedings of the European conference on computer vision, ECCV 2018, pages 722–737, 2018.

[43] Naoya Takahashi, Michael Gygli, Beat Pfister, and Luc Van Gool. Deep convolutional neural networks and data augmentation for acoustic event detection. arXiv preprint arXiv:1604.07160, 2016.

[44] Matthias Meyer, Lukas Cavigelli, and Lothar Thiele. Efficient convolutional neural network for audio event detection. arXiv preprint arXiv:1709.09888, 2017.

[45] Adam Byerly, Tatiana Kalganova, and Ian Dear. No Routing Needed Between Capsules. arXiv preprint arXiv:2001.09136v6, 2021.

[46] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.

[47] Mingxing Tan and Quoc Le. EfficientNet: Rethinking model scaling for convolutional neural networks. In Proceedings of the International Conference on Machine Learning, ICML 2019, pages 6105–6114. PMLR, 2019.

[48] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. arXiv preprint arXiv:2106.04560, 2021.

[49] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, CVPR 2016, pages 770–778, 2016.

[50] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In Proceedings of the IEEE conference on computer vision and pattern recognition, CVPR 2018, pages 2704–2713, 2018.

**MD MOHAIMENUZZAMAN** received the B.Sc. and M.Sc. degrees in computer science and engineering, in 2007 and 2013, respectively. He is currently pursuing the Ph.D. degree in data science with the Department of Data Science and AI, Faculty of Information Technology, Monash University, Australia.

Before starting Ph.D., he developed software applications for international clients for about a decade. He is currently working as a Sessional Teaching Associate with the Faculty of Information Technology, where he teaches data science and software engineering related courses. He received the "2020 Faculty Teaching Excellence" award for teaching "Introduction to Data Science" to graduate students.

**CHRISTOPH BERGMEIR** received the Ph.D. degree in computer science from the University of Granada, Spain, and the M.Sc. degree in computer science from the University of Ulm, Germany.

He is currently a Senior Lecturer in data science and artificial intelligence with Monash University. He is a 2019 ARC DECRA Fellow with the Department of Data Science and AI, Monash University, where he develops "efficient and effective analytics for real-world time series forecasting." He works as a Data Scientist in a variety of projects with external partners in diverse sectors, e.g., in healthcare or infrastructure maintenance. He has published on time series prediction using machine learning methods, recurrent neural networks and long short-term memory neural networks (LSTM), time series predictor evaluation, and on medical applications and software packages in the R programming language, in journals, such as IEEE TRANSACTIONS ON NEURAL NETWORKS and LEARNING SYSTEMS, Journal of Statistical Software, Computational Statistics & Data Analysis, and Information Sciences.

**BERND MEYER** received the Ph.D. degree in computer science from the University of Hagen, Germany, in 1994. He is currently a Professor with the Department of Data Science and AI, Faculty of Information Technology, Monash University, Australia.

He works on data-intensive computational ecology, develops mathematical and computational models for the interactions of organisms with their environment, mostly focusing on the collective behaviour of social insects, such as bees and ants. How these self-organised "super-organisms" coordinate their actions remains a fascinating enigma. He also works on AI-based methods for monitoring animal activity as the basis for ecosystem monitoring and for automating experiments.

• • •