

# An Empirical Study of Low Precision Quantization for TinyML

Shaojie Zhuo\*  
Qualcomm AI Research  
Qualcomm Canada ULC  
shaojie@qti.qualcomm.com

Hongyu Chen\*  
University of Toronto  
hy.chen@mail.utoronto.ca

Ramchalam Kinattinkara Ramakrishnan  
Qualcomm AI Research  
Qualcomm Canada ULC  
rkinatti@qti.qualcomm.com

Tommy Chen  
Qualcomm AI Research  
Qualcomm Canada ULC  
tommchen@qti.qualcomm.com

Chen Feng  
Qualcomm AI Research  
Qualcomm Canada ULC  
chenf@qti.qualcomm.com

Yicheng Lin  
Qualcomm AI Research  
Qualcomm Canada ULC  
yichengl@qti.qualcomm.com

Parker Zhang  
Qualcomm AI Research  
Qualcomm Canada ULC  
xiaopeng@qti.qualcomm.com

Liang Shen  
Qualcomm AI Research  
Qualcomm Canada ULC  
liangs@qti.qualcomm.com

## ABSTRACT

Tiny machine learning (tinyML) has emerged during the past few years aiming to deploy machine learning models to embedded AI processors with highly constrained memory and computation capacity. Low precision quantization is an important model compression technique that can greatly reduce both memory consumption and computation cost of model inference. In this study, we focus on post-training quantization (PTQ) algorithms that quantize a model to low-bit (less than 8-bit) precision with only a small set of calibration data and benchmark them on different tinyML use cases. To achieve a fair comparison, we build a simulated quantization framework to investigate recent PTQ algorithms. Furthermore, we break down those algorithms into essential components and reassembled a generic PTQ pipeline. With ablation study on different alternatives of components in the pipeline, we reveal key design choices when performing low precision quantization. We hope this work could provide useful data points and shed lights on the future research of low precision quantization.

## KEYWORDS

tinyML, neural networks, quantization, benchamrk

### ACM Reference Format:

Shaojie Zhuo, Hongyu Chen, Ramchalam Kinattinkara Ramakrishnan, Tommy Chen, Chen Feng, Yicheng Lin, Parker Zhang, and Liang Shen. 2022. An Empirical Study of Low Precision Quantization for TinyML. In *Proceedings of tinyML Research Symposium (tinyML Research Symposium'22)*. ACM, New York, NY, USA, 9 pages.

## 1 INTRODUCTION

Embedding AI directly on edge devices is becoming a key to revolutionize the internet of things (IoT), where billions of tiny devices are leveraged to gain great productivity and efficiency in areas including consumer, medical, automotive and industrial. As a result, tiny machine learning (tinyML) has emerged over the past few years

\*Equal contribution

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

*tinyML Research Symposium'22, March 2022, San Jose, CA*

© 2022 Copyright held by the owner/author(s).

that aims to deploy machine learning algorithms, especially neural networks, on embedded AI processors such as micro-controllers or embedded NPUs (e.g. Qualcomm's eNPU [23] and ARM's MicroNPU [1]) with very low power consumption in the level of a few milli-watts. By performing inferences on the devices, tinyML brings the advantage of high energy efficiency, fast responsiveness, high privacy and strong autonomy of edge devices. However, the embedded AI processors are usually highly resource-constrained with limited memory and computation capability. A large body of literature has focused on addressing these issues by making neural networks more efficient, including efficient neural architecture search [4, 18, 29], neural networks and system co-design [16, 17] and model compression such as pruning [10], knowledge distillation [11] and quantization [10, 14].

Neural network quantization is one of the most important model compression techniques. It compresses models by using low-bit precision representation for weight and activation tensors, instead of the 32-bit (or 16-bit for half-precision) precision that is commonly used during model training. The memory consumption and computation of a low precision model are greatly reduced. And quantization compresses a model without changing its architecture. This is particularly useful when the model architecture is already optimally designed for a specific inference system or device. Besides, quantization can be applied along with other model compression techniques, such as pruning and knowledge distillation.

Neural networks have been shown to be quite robust to quantization with 8-bit precision without sacrificing the accuracy [20]. Quantization with lower precision can further reduce the memory consumption and computation. And ultra-low precision (1 or 2-bit) operations can often be computed efficiently with bit-wise arithmetic and thus achieving significant computation acceleration [28]. However, due to the large quantization noise, the benefits of low precision quantization often come at the cost of significant accuracy degradation. Among approaches to mitigate the accuracy drop caused by quantization, post-training quantization (PTQ) directly performs quantization on a pre-trained full precision model to convert it to a corresponding low precision model. Unlike quantization-aware training (QAT), PTQ does not require the full training pipeline of the original model and can be applied with a small amount of calibration data (or even data-free [21]), thus becoming a popular model quantization routine.

In this paper, we aim to study and benchmark recent PTQ algorithms on neural network models designed for tinyML use cases. By doing that, we try to answer the following questions.

First of all, *what is the performance of PTQ algorithms on tinyML models?* Existing work [5, 15, 19, 21, 25] has shown successful applications of low precision PTQ on big neural network models for large-scale tasks such as image classification on ImageNet. However, tinyML models are usually very compact and are designed for relatively simple tasks. It is unknown whether those models can be quantized to low precision while preserving their accuracy.

Secondly, *what are the design choices making a PTQ algorithm superior to others?* Comparing the performance of different PTQ algorithms helps us to select a proper quantization algorithm. And further study on their algorithm design choices can reveal the key reason of an algorithm's success and provide insights for further improvement of the algorithm.

Last but not least, *what is the trade-off between accuracy and memory/computation when applying low precision quantization?* This is an important question when deploying a tinyML model. By knowing the sensitivity of different models to quantization, we are able to seek the best trade-off between accuracy and memory/computation considering the very limited resources on tinyML devices.

By conducting extensive experiments with various tinyML models, we find answers to the previous questions and many insightful observations. To summarize our contributions:

- We create a quantization simulation framework and an unified PTQ pipeline to encourage fair comparison of different PTQ algorithms.
- We conduct ablations to reveal the key design choices of the PTQ pipeline using a set of models with varieties of architectures in tinyML.
- we demonstrate that low precision quantization is useful for tinyML by compressing a model to greatly save memory and computations while preserving accuracy.
- We find key gaps of low precision quantization for tinyML and point out the directions of potential improvements.

We believe that low precision quantization is important to tinyML. And the observations, challenges and open questions in this study are good data points to understand the existing PTQ algorithms and the characteristic of low precision quantization, and shed lights on the direction of future work to improve the performance.

## 2 RELATED WORK

*Quantization.* We refer to the readers a good overview of quantization basics and algorithms in [7]. Quantization algorithms can be classified into quantization-aware training (QAT) and post-training quantization (PTQ). QAT re-trains a model with quantized parameters so that it converges to a point with a better loss. An important challenge of the backpropagation during re-training is how to treat the non-differentiable rounding operator in the quantization function. Straight Through Estimator (STE) [3] is one of the main techniques used to approximate the gradient. Built on top of that, several recent papers [5, 13] proposed to learn the quantization parameters together with the network parameters and achieve better results. Although STE-based methods enable the training of

quantized networks with gradient-based optimization, the gradient mismatch between forward and backward passes affects the training effectiveness, especially for low precision quantization. To alleviate the gradient mismatch, soft quantizers with sigmoid or tanh function [8] are used to approximate the discrete rounding in both forward and backward passes. These approaches, however, cause the accuracy drop due to the discrepancy between the soft rounding during training and the hard rounding at inference time. Despite the promising results given by QAT methods, they usually need the whole training pipeline and much more computation. In contrast, PTQ is attractive in the sense that it requires only a small set of (unlabeled) calibration data and thus become the main focus of this study. A fundamental problem in PTQ is finding good quantization parameters for both weight and activation tensors. It is particular challenging when there is imbalanced distribution among different channels of a tensor. To resolve that, cross layer equalization (CLE) [21] and outlier channel splitting (OCS) [27] are proposed to re-balance the distribution of weights. Although good performance are achieved for 8-bit quantization, PTQ suffers from severe performance degradation as precision goes lower than 8-bit. Recently, a new paradigm makes huge progress for low precision PTQ. The approach utilizes layerwise optimization to minimize the error between the output of a quantized layer and that of its corresponding full-precision layer, with the help of a small set of calibration data. A series of work leveraging the layer-wise calibration paradigm, namely AdaRound [19], AdaQuant [12], BitSplit [25], BRECC [15], has push the limit of PTQ down to 2 bit. We thus focus on the study of those methods and compare their performances.

*Quantization for tinyML.* There are works [14, 26] focusing on the study of quantization for general applications with big models in the study. In the context of tinyML, [22] proposed a new QAT method for low precision quantization. [2, 6] focused on deployment of the 8-bit quantized models and their latency on actual tinyML devices. By contrast, Our study aims on benchmarking the performance of the state of the art quantization methods on low precision (less than 8-bit) quantization for tinyML. We selected several representative applications in tinyML to study the effects of different algorithm design choices and the trade-off among accuracy, memory consumption and computational cost.

## 3 QUANTIZATION FRAMEWORK

The performances of existing PTQ algorithms are usually reported based on separate implementations without a common framework. The implementation differences (e.g. the way how quantization is simulated) makes it difficult to do a fair comparison of different algorithms. And existing work usually focuses on improvement of different aspects in the quantization process without optimizing the whole pipeline. To this end, we implement a quantization framework to facilitate the quantization study. Built on top of that, we propose a unified post-training quantization pipeline, so that existing PTQ algorithms can be deconstructed and fitted into the pipeline as improvement of parts of its components. It does not only make fair comparisons of PTQ algorithms possible, but also has the flexibility to ablate on different choices of each components

in the pipeline and find the optimal one that generate the best performance. In this section, we introduce our quantization simulation framework and the proposed unified PTQ pipeline.

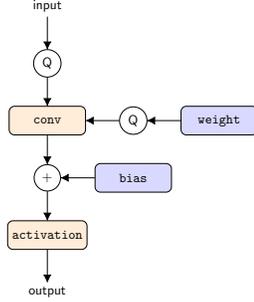


Figure 1: Simulated quantization for the forward pass of a convolution layer.

### 3.1 Simulated quantization

We use simulated quantization [14] in our framework. A tensor is quantized and then de-quantized back to floating precision to simulate quantization loss while the operations are carried out with floating point arithmetic. Figure 1 shows an example of the quantized forward pass of a convolution layer. Both input and weight tensor go through the quantization function in the quantizer to simulate quantization and then are fed to the convolution computation in full precision. The bias remains in full precision. The quantization of other operators is simulated in a similar way.

We adopt the asymmetric quantization scheme in our framework for quantization simulation. The asymmetric quantization function is defined by three quantization parameters: the scale factor  $s$ , the zero-point  $z$  and the bit-width  $b$ ,

$$\hat{x} = q(x; s, z, b) = s \left[ \text{clamp} \left( \left\lfloor \frac{x}{s} \right\rfloor + z, n, p \right) - z \right] \quad (1)$$

where  $\hat{x}$  is the quantized approximation of  $x$ ,  $\lfloor \cdot \rfloor$  indicates the round function and  $\text{clamp}(\cdot)$  clamps values between  $n$  and  $p$ , which are determined by the bit-width  $b$  and define the range of the integer grids. As a simplified version of the asymmetric quantization, symmetric quantization restricts the zero-point to be 0 to reduce the additional computation introduced by the zero-point offset. We follow [20] to use *asymmetric quantization for activation tensors* and *symmetric quantization for parameters*, considering that the distributions of activation tensors are often asymmetric about zero while those for parameters are roughly symmetric about zero.

Quantization granularity is another factor affecting the accuracy of a quantized model. Per-tensor quantization has a single set of quantization parameter for the whole tensor while per-channel quantization has a separate set of quantization parameters for each channel of the tensor. Per-channel quantization usually introduces lower quantization noise due to its finer granularity. In our framework, we use *per-channel quantization for parameters* and *per-tensor quantization for activation tensors*, considering the readiness of hardware support and to better preserve the accuracy of a model [20].

### 3.2 PTQ Pipeline

As illustrated in Figure 2, the proposed unified PTQ pipeline takes a pre-trained full precision model with a set of calibration data as

input and produces the optimized low precision model through four sequential steps.

**Step 1: FP model preprocessing.** In this step, a full precision (FP) model is pre-processed to make it more quantization friendly. One common issue for quantization is the imbalanced distribution among various channels of the tensor, which makes it hard to find a suitable quantization parameters for the tensor. This is especially prevalent in the weights of efficient models with depth-wise separable convolutions [21]. Fine-grained quantization (e.g., per-channel quantization) is able to alleviate the issue. Without it, cross-layer range equalization (CLE) [21] and Outlier channel splitting (OCS) [27] are introduced to re-balance the distribution of a tensor. Although we use per-channel quantization for parameters, the activation tensors are quantized using per-tensor quantization. It is worthwhile to show the impact of FP model preprocessing on quantization.

**Step 2: Add quantizers and initialization.** After the preprocessing, we add quantizers to a model to simulate quantization for both parameters and activation tensors. To initialize the quantization parameters (the scale  $s$  and zero-point  $z$ ) of each quantizer, one common approach is to find the clipping range  $[\alpha, \beta]$  first and then use it to calculate the scale and zero-point, i.e.,

$$s = \frac{\beta - \alpha}{2^b - 1}, z = \lfloor \frac{-\alpha}{s} \rfloor + n \quad (2)$$

And there are multiple ways to determine the clipping range,

- *MinMax* initialization sets the clipping range using the min/max of a tensor. It is simple but usually suffers from outliers.
- *MSE* initialization finds the range that minimize the mean squared error (MSE) between the original and quantized tensor using grid search or analytical approximations  $\alpha^*, \beta^* = \arg \min_{\alpha, \beta} \|x - \hat{x}\|_F^2$ , where  $\hat{x}$  is the quantized version of  $x$  and  $\|\cdot\|_F$  is the Frobenius norm.

**Step 3: Quantization optimization.** As shown in Equation 1, the quantization function maps a floating point value to a low-precision value and then de-quantize it back to a floating point value. During the process, clipping error is introduced by the  $\text{clamp}(\cdot)$  function and rounding error is introduced by the round function  $\lfloor \cdot \rfloor$ . Quantization optimization is the process to seek the right trade-off between the two errors. We adopt a general layerwise optimization process, as illustrated in Figure 3. A low precision model is optimized layer by layer to minimize the loss (usually the mean squared error) between the outputs of the full precision layer and the quantized layer. As shown in equation 1, multiple variables can potentially be optimized,

- Optimizing scale and zero-point (*opt\_qparam*). Learnt step-size quantization [5] treats the scale  $s$  and zero-point  $z$  as trainable parameter and optimizes them to minimize the layerwise MSE.
- Optimizing weights (*opt\_weights*). AdaQuant [12] extends the flexibility and optimizes the weights directly to minimize the layerwise MSE.
- Optimizing bits (*opt\_bits*). Instead of working on the floating point, BitSplit [25] performs the optimization on the quantized values directly and splits the integer into multiple bits and optimizes them to minimize the layerwise MSE.

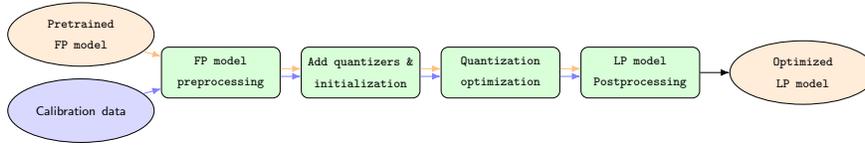


Figure 2: The unified pipeline for PTQ. It takes a pre-trained full precision model and a small set of calibration data as input and output a optimized quantized model through multiple steps.

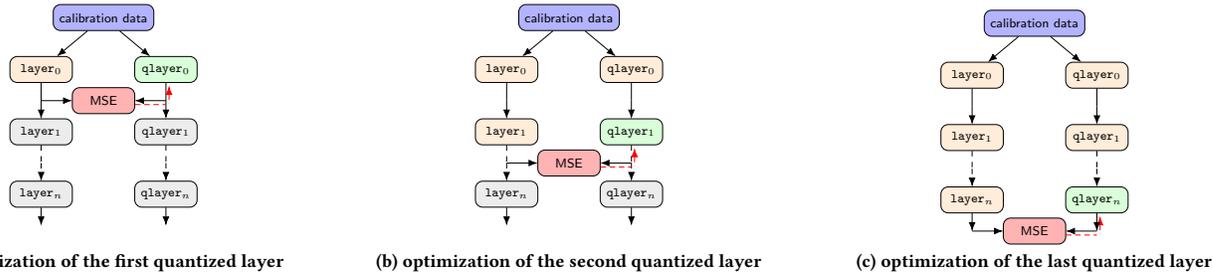


Figure 3: Layerwise quantization optimization. With a small set of calibration data, (a) the first quantized layer (in green) is activated and optimized, while the rest of the quantized layers stay inactive (in gray). The first quantized is optimized to minimize the mean squared error(MSE) between its output and that of the first FP layer; (b) the first quantized layer is optimized and frozen (in orange) and the second quantized layer is activated and optimized; (c) The quantized model is optimized layer by layer until the last one.

- Optimizing rounding (*opt\_round*). AdaRound [19] optimizes the rounding by introducing a auxiliary parameter for each weight and constrains it to be either 0 or 1 to achieve rounding up or down.

In the quantization optimization step, instead of layer-wise optimization, block-wise optimization [15] is also proposed to capture the cross-layer dependency within the block and it is claimed that block-wise optimization could be a better bias-variance trade-off choice with limited calibration data for low precision quantization.

*Step 4: LP model post-processing.* After quantization optimization, bias tuning [12] can be applied on the low-precision (LP) model to update biases in a model to further compensate for quantization loss. The biases are updated to minimize the MSE of the final outputs between FP model and the LP model. In theory, the bias tuning can be done during the quantization optimization, but the rationale to put it in the post-processing stage is that, by updating a limited part of the model parameters in an end-to-end fashion, it can provide global guidance to optimize for the final output without over-fitting.

With the unified pipeline, existing PTQ algorithms can be deconstructed and fitted into the pipeline. We vary each component in the pipeline to construct new PTQ algorithms. We do ablations on different options for each component to reveal the best options and then use them to construct the best PTQ algorithm.

## 4 EXPERIMENTS

In this section, we benchmark and ablate the performance of recent PTQ algorithms on multiple tinyML models. We first describe the tinyML tasks and models selected for the experiments and then show the experiment results of ablations on different components of the PTQ pipeline and results of end-to-end study. In our experiments, for all the models, we report the mean and standard deviation of the accuracy, calculated using 5 runs with different initial seeds. And we use 1024 unlabeled data extracted from the original training data as the calibration dataset.

### 4.1 TinyML tasks and models

We select four use cases for our study: image classification, visual wake words, Keyword spotting and human activity recognition. The first three are from the MLPerf Tiny benchmark [2]. The selection is based on the relevance to real world applications, the availability of open-sourced dataset and models, and the coverage of a wide variety of applications and a wide variety of model architectures in tinyML. We train each model based on the recommended setting in [2] to get the full precision models. A summary of the tasks and models are shown in Table 1.

### 4.2 Ablations of the PTQ pipeline

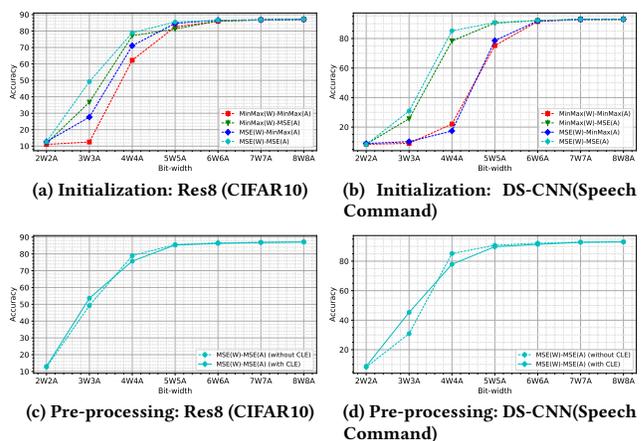


Figure 4: Ablation study on FP model preprocessing and quantizer initialization for both weights and activations. (a) and (b) show the accuracy after quantizer initialization. (c) and (d) show the accuracy after MSE-baed quantizer initialization with and without CLE as the FP model preprocessing.

Use case	Dataset	Model	Float Accuracy	MAC	Num Param.	Peak activation
Image Classification	CIFAR10	Res8	86.75	12,591,808	77,706	36,608
Keyword Spotting	Speech Commands	DS-CNN	93.03	2,736,832	22,604	36,864
Visual Wake Words	VWV	MobileNetV1	85.10	7,723,776	210,850	32,768
Human Activity Recognition	UCI-HAR	CNN	90.36	2,298,368	523,462	8,064

Table 1: TinyML tasks and models selected for the evaluation of quantization algorithms. The detailed architecture of each model can be found in the appendix

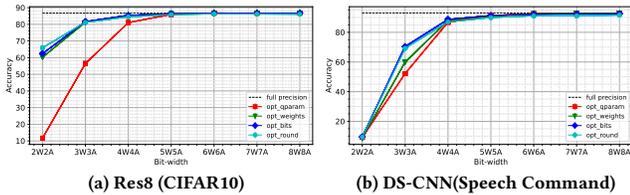


Figure 5: Ablation study on optimization variables during quantization optimization. Better zoom in to view the results.

*FP model preprocessing and quantizer initialization.* We consider FP model pre-processing and quantizer initialization together in this study to see how the combination of the two affects the accuracy. Figure 4 (a) and (b) show the comparison of model accuracy with different initialization methods. We need to ignore the results here for 2W2A because the quantized model after initialization performs as a random classifier caused by large quantization noise. From the results, we can see that MSE-based quantizer initialization usually outperforms MinMax-based one. And MSE-based initialization for both weights and activations performs the best consistently, this is aligned with previous study on quantizer initialization [20]. For FP model preprocessing, we consider cross layer equalization (CLE) because it is effective without changing the architecture of the model. Figure 4 (c) and (d) compare the model accuracy with and without CLE after MSE-based initialization. We see that there is no accuracy discrepancy for higher bitwidth. Although we use per-channel quantization for weights, CLE still helps to improve the accuracy for low precision cases (3W3A), while it has negative impact for the precision in between(4W4A or 5W5A). This is because activation tensors have per-tensor quantization and CLE also changes the distribution of activation tensors.

*Optimization variable.* As discussed in section 3.2, different variables can be optimized to minimize the layer-wise loss. We investigate the effect of optimizing different variables. Please note that only the weights are quantized while activation are in full precision for this study. As shown in Figure 5, we can see that *opt\_qparam* is clearly less effective compared with the other three, especially when the bit-width is low. It is because that only the quantization parameters (zero-points and offsets) are trainable, which restricted its capability to to compensate for the large quantization noise introduced by low precision quantization. *opt\_round* and *opt\_bits* usually outperforms *opt\_weights*. We believe that the STE [3] used in *opt\_weights* causes biased gradient estimation under low-precision cases that affects the efficacy of optimization. And the constrained optimization space in *opt\_bits* and *opt\_round* makes them less prone to over-fitting with limited calibration data.

*Optimization granularity.* We conduct the study on optimization granularity and the results are shown in Figure 6. The solid lines represent the layerwise optimization and the dotted lines represent the blockwise optimization. *opt\_bits* is excluded here because it

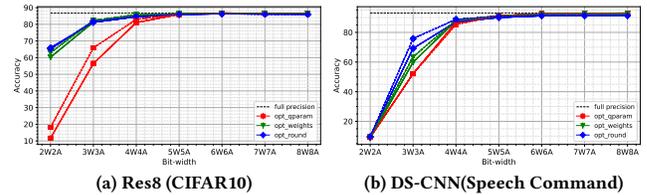


Figure 6: Ablation study on optimization granularity during quantization optimization. The solid lines represent the layerwise optimization and the dotted lines represent the blockwise optimization. Better zoom in to view the results.

doesn't support blockwise optimization by design. The block is defined as the basic residual block in Res8 and the depth separable convolution in DS-CNN. We observe that both layerwise optimization and blockwise optimization have very similar performance when the models are quantized to 4W4A and higher, regardless of the different optimization variables. For quantization with precision lower than 4W4A, blockwise optimization provides quite significant improvement on the accuracy for Res8 quantized to 3W3A and 2W2A, and for DS-CNN quantized to 3W3A.

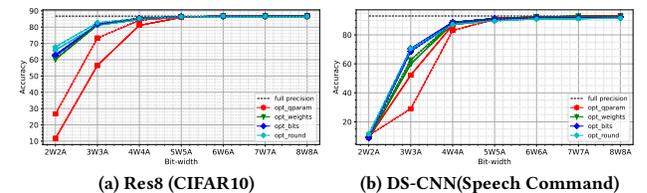


Figure 7: Ablation study on bias tuning for LP model post-processing. The solid lines represent optimization without post-processing and the dotted lines represent the optimization with post-processing. Better zoom in to view the results.

*LP model Post-processing.* We also conduct experiment to show the impact of bias tuning as the low precision model post-processing. As shown in Figure 7, bias tuning consistently improve the accuracy for all the methods (except for the case of DS-CNN using *opt\_qparam*). For Res8 using *opt\_qparam*, it has more than 10% improvement for both 3W3A and 2W2A quantization. This empirically shows the effectiveness of end-to-end bias tuning.

### 4.3 Results

To compare end-to-end performance with both weight and activation quantization enabled, we fix all the options in the pipeline using the best ones based on ablation study results, except for the optimization variables. We use symmetric per-channel quantization for weights and asymmetric per-tensor quantization for activations. The FP model pre-processing is enabled and MSE-based initialization is applied on both weight quantizers and activation quantizers. And we use layerwise optimization with bias tuning enabled as

Bitwidth (W/A)	Method	Res8 (CIFAR10)	DS-CNN (KWS)	MobileNetv1 (VWW)	CNN (UCI-HAR)
32/32	full precision	86.75	93.03	85.10	90.36
8/8	opt_qparam	<b>86.81±0.06</b>	<b>92.63±0.14</b>	85.13±0.12	<b>90.39±0.17</b>
	opt_weights	86.76±0.04	92.55±0.06	85.14±0.05	90.81±0.24
	opt_bits	86.70±0.10	91.93±0.10	<b>85.28±0.15</b>	90.04±0.13
	opt_round	86.65±0.14	91.65±0.20	84.98±0.21	90.37±0.26
7/7	opt_qparam	<b>86.79±0.11</b>	92.57±0.10	<b>85.18±0.16</b>	90.54±0.24
	opt_weights	86.69±0.08	<b>92.62±0.14</b>	85.15±0.17	<b>90.88±0.31</b>
	opt_bits	86.64±0.08	91.92±0.15	85.16±0.17	89.98±0.15
	opt_round	86.70±0.14	91.54±0.21	84.87±0.11	90.37±0.35
6/6	opt_qparam	86.67±0.15	92.01±0.19	84.90±0.14	90.24±0.11
	opt_weights	<b>86.71±0.14</b>	<b>92.19±0.19</b>	85.06±0.10	<b>90.86±0.22</b>
	opt_bits	86.63±0.06	91.80±0.14	<b>85.11±0.10</b>	90.03±0.31
	opt_round	86.64±0.06	90.56±0.31	84.81±0.18	90.20±0.23
5/5	opt_qparam	86.35±0.03	<b>90.86±0.28</b>	84.33±0.19	90.61±0.41
	opt_weights	86.47±0.12	90.82±0.20	<b>84.82±0.15</b>	<b>91.00±0.47</b>
	opt_bits	86.37±0.15	90.42±0.36	84.39±0.16	90.09±0.21
	opt_round	<b>86.64±0.10</b>	90.13±0.93	84.21±0.24	90.41±0.30
4/4	opt_qparam	84.61±0.15	83.31±1.02	78.40±1.52	91.12±0.22
	opt_weights	<b>85.91±0.18</b>	<b>88.41±0.25</b>	<b>82.96±0.23</b>	<b>90.76±0.13</b>
	opt_bits	85.33±0.09	87.88±0.36	81.72±0.62	90.38±0.32
	opt_round	85.34±0.32	88.33±0.95	81.89±0.50	90.57±0.35
3/3	opt_qparam	74.30±0.53	29.58±4.60	56.79±1.43	90.57±0.23
	opt_weights	82.62±0.13	63.99±2.98	76.03±1.36	90.83±0.31
	opt_bits	81.81±0.14	63.95±2.66	74.45±0.71	90.32±0.27
	opt_round	<b>83.21±0.41</b>	<b>75.22±1.57</b>	<b>76.71±1.21</b>	<b>91.38±1.15</b>
2/2	opt_qparam	25.26±2.2	9.16±1.18	47.15±0.00	82.71±0.76
	opt_weights	62.81±0.53	8.83±0.85	62.24±0.58	<b>90.46±0.54</b>
	opt_bits	63.72±0.47	8.77±0.23	<b>62.48±0.34</b>	86.74±0.59
	opt_round	<b>69.40±0.74</b>	<b>9.42±0.51</b>	60.58±0.96	89.63±0.83

Table 2: Comparison of PTQ pipelines with different optimization variables while fixing other options the best in the pipeline.

LP model post-processing. The results are shown in Table 2. For a higher bit-width (5W5A to 8W8A), different PTQ pipelines perform similarly with very small accuracy drop compared to the full precision accuracy. However, for a lower bit-width, different pipelines show big performance discrepancy. The pipelines using *opt\_round* usually perform better. However, even for the best quantization algorithm, there is still big accuracy drop when going to ultra-low precision. The accuracy for 2W2A is 69.40% for CIFAR10, with accuracy drop as large as 17.5%.

To better show the benefits of low precision quantization, we view the results on two axes: the accuracy drop of the best PTQ pipeline, and corresponding computation saving (or memory saving) with different bit-widths. We use bit operation (BOP) as a system-agnostic proxy to the computation complexity. The BOP count of a layer  $l$  is computed as:  $BOP(l) = MACs(l)b_w b_a$ , where  $MACs(l)$  is the multiply-accumulate count of the layer,  $b_w$  is the bit-width of the weight and  $b_a$  is the bit-width of the input. We use peak memory as the measurement of the memory consumption, which is computed as  $peak\_memory = total\_weight * b_w + peak\_activation * b_a$ . it measures the maximum memory required to store the weights and activations during inferences, assuming that previous intermediate activations can be effectively released if they are not used by the later layers. The results are shown in Figure 8. We can see that 4W4A is a sweet point that has 50% memory saving and 75% computation saving with small accuracy drop (< 5%). This is mismatching with most inference frameworks that support only 8-bit quantization [6]. Some models can even tolerant lower precision quantization with negligible accuracy drop. For example, the CNN model for UCI-HAR can be quantized to 2W2A with less than 2% accuracy drop. However, this could be a model/task dependent behavior. For most of the models, we have to sacrifice the accuracy if going for ultra-low precision such as 3W3A or 2W2A, and thus facing the accuracy and memory/computation trade-off. We can

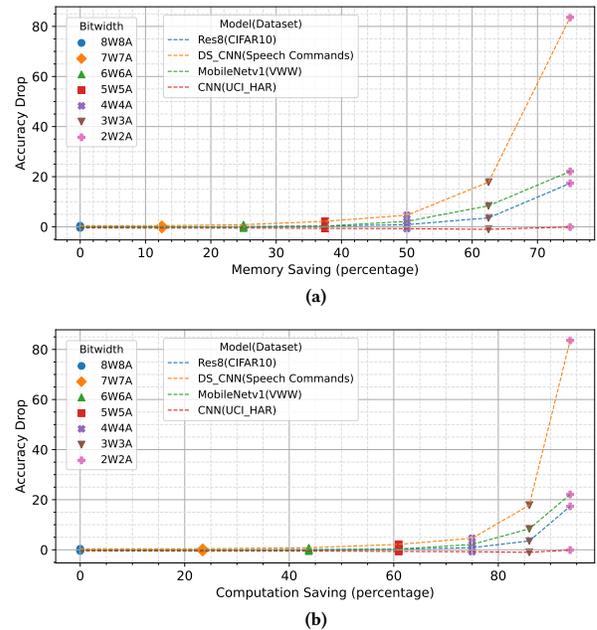


Figure 8: Tradeoff between accuracy and memory/computation saving. 4W4A is a sweet point having large memory/computation saving with small accuracy drop. Some models can go lower precision with negligible accuracy drop.

clearly see the benefits of low precision quantization in terms of computation and memory saving (and eventually power saving).

## 5 CONCLUSION AND FUTURE DIRECTIONS

In this paper, with a unified quantization pipeline, we benchmark PTQ algorithms on carefully selected models for tinyML. We conduct ablations to reveal the key design choices of the pipeline. We show that low precision quantization with the best PTQ algorithm can have huge memory and computation savings with negligible accuracy drop, although there is still big accuracy drop if quantized models to ultra-low precision. Our empirical study provides useful data points on low precision quantization and points out potential future directions for improvement on algorithm design and system support for tinyML.

How to further improve the accuracy for ultra-low precision quantization is an important direction for future research. Based on our unified PTQ pipeline, new improvement on different components can further reduce the accuracy degradation. And mix-precision quantization and quantization aware training could be potential solutions. Besides, existing embedded AI frameworks, such as TensorFlow Lite Micro [9] and STM32Cube.AI [24], can only support computation in floating point or 8-bit integers. The framework with lower precision quantization support is not widely available. It first leaves big room to squeeze the model by combining quantization with other model compression techniques. And it is important to investigate the best combination in practice. At the meanwhile, developing efficient embedded AI frameworks with low precision quantization support will be particularly useful for deployment of tinyML models.

REFERENCES

- [1] Arm. 2021. Ethos-U55 Machine Learning Processor. <https://www.arm.com/products/silicon-ip-cpu/ethos/ethos-u55>. last visited 05/03/2022.
- [2] Colby Banbury, Vijay Janapa Reddi, Peter Torelli, Jeremy Holleman, Nat Jeffries, Csaba Kiraly, Pietro Montino, David Kanter, Sebastian Ahmed, Danilo Pau, et al. 2021. MLPerf Tiny Benchmark. *Neurips* (2021).
- [3] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation. *arXiv:1308.3432* (2013).
- [4] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. 2020. Once-for-All: Train One Network and Specialize it for Efficient Deployment. In *ICLR*.
- [5] Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Aparambath, and Dharmendra Modha. 2020. Learned Step Size Quantization. *ICLR*.
- [6] Sedigh Ghamari, Koray Ozcan, Thu Dinh, Andrey Melnikov, Juan Carvajal, Jan Ernst, and Sek Chai. 2021. Quantization-Guided Training for Compact TinyML Models. *TinyML Research Symposium 2021* (2021).
- [7] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W. Mahoney, and Kurt Keutzer. 2021. A Survey of Quantization Methods for Efficient Neural Network Inference. *arXiv:2103.13630* (2021).
- [8] Ruihao Gong, Xianglong Liu, Shenghu Jiang, Tianxiang Li, Peng Hu, Jiazhen Lin, Fengwei Yu, and Junjie Yan. 2019. Differentiable Soft Quantization: Bridging Full-Precision and Low-Bit Neural Networks.
- [9] Google. 2021. TensorFlow Lite for Microcontrollers. <https://www.tensorflow.org/lite/microcontrollers>. last visited 05/03/2022.
- [10] Song Han, Huizi Mao, and William J. Dally. 2016. Deep Compression: Compressing Deep Neural Network with Pruning, Trained Quantization and Huffman Coding. In *ICLR*.
- [11] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the Knowledge in a Neural Network. *arXiv:1503.02531* (2015).
- [12] Itay Hubara, Yury Nahshan, Yair Hanani, Ron Banner, and Daniel Soudry. 2020. Improving Post Training Neural Quantization: Layer-wise Calibration and Integer Programming. *arXiv:2006.10518* (2020).
- [13] Sambhav R. Jain, Albert Gural, Michael Wu, and Chris H. Dick. 2020. Trained Quantization Thresholds for Accurate and Efficient Fixed-Point Inference of Deep Neural Networks. *arXiv:1903.08066* (2020).
- [14] Raghuraman Krishnamoorthi. 2018. Quantizing deep convolutional networks for efficient inference: A whitepaper.
- [15] Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. 2021. BRECO: Pushing the Limit of Post-Training Quantization by Block Reconstruction. In *ICLR*.
- [16] Ji Lin, Wei-Ming Chen, Han Cai, Chuang Gan, and Song Han. 2021. MCUNetV2: Memory-Efficient Patch-based Inference for Tiny Deep Learning. In *NeurIPS*.
- [17] Ji Lin, Wei-Ming Chen, John Cohn, Chuang Gan, and Song Han. 2020. MCUNet: Tiny Deep Learning on IoT Devices. In *NeurIPS*.
- [18] Hanxiao Liu, Karen Simonyan, and Yiming Yang. 2019. DARTS: Differentiable Architecture Search. In *ICLR*.
- [19] Markus Nagel, Rana Ali Amjad, Mart van Baalen, Christos Louizos, and Tijmen Blankevoort. 2020. Up or Down? Adaptive Rounding for Post-Training Quantization. In *ICML*.
- [20] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tijmen Blankevoort. 2021. A White Paper on Neural Network Quantization. *arXiv:2106.08295* (2021).
- [21] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. 2019. Data-Free Quantization Through Weight Equalization and Bias Correction. In *ICCV*.
- [22] Pierre-Emmanuel Novac, Ghouthi Boukli Hacene, Alain Pegatoquet, Benoit Miramond, and Vincent Gripon. 2021. Quantization and Deployment of Deep Neural Networks on Microcontrollers. *Sensors* (2021).
- [23] Qualcomm. 2021. Hexagon DSP SDK Tools and Resources. <https://developer.qualcomm.com/software/hexagon-dsp-sdk/tools>. last visited 05/03/2022.
- [24] STMicroelectronics. 2021. STM32Cube.AI. <https://www.st.com/en/embedded-software/x-cube-ai.html>. last visited 05/03/2022.
- [25] Peisong Wang, Qiang Chen, Xiangyu He, and Jian Cheng. 2020. Towards Accurate Post-training Network Quantization via Bit-Split and Stitching. In *ICML*.
- [26] Hao Wu, Patrick Judd, Xiaojie Zhang, Mikhail Isaev, and Paulius Micikevicius. 2020. Integer Quantization for Deep Learning Inference: Principles and Empirical Evaluation. *arXiv:2004.09602* (2020).
- [27] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Christopher De Sa, and Zhiru Zhang. 2019. Improving Neural Network Quantization without Retraining using Outlier Channel Splitting. *arXiv:1901.09504* (2019).
- [28] Shien Zhu, Luan H. K. Duong, and Weichen Liu. 2022. TAB: Unified and Optimized Ternary, Binary and Mixed-Precision Neural Network Inference on the Edge. *ACM Transactions on Embedded Computing Systems* (2022).
- [29] Barret Zoph and Quoc V. Le. 2016. Neural Architecture Search with Reinforcement Learning. *CoRR* (2016).

A MODEL ARCHITECTURES

Here we show the basic blocks used to define the models and the detailed architecture of each TinyML model.

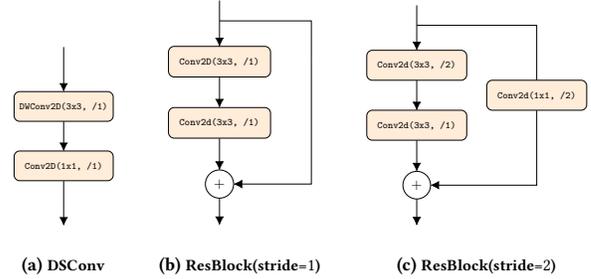


Figure 9: The basic blocks used in the model architectures.

Input	Operator	#out	stride
32 × 32 × 3	Conv2d, 3 × 3	16	1
32 × 32 × 16	ResBlock	16	1
32 × 32 × 16	ResBlock	32	2
16 × 16 × 32	ResBlock	64	2
8 × 8 × 64	AvgPool, 8 × 8	-	1
1 × 64	Fully Connected	10	-

Table 3: Res8 (CIFAR10) architecture definition

Input	Operator	#out	stride
10 × 49 × 1	Conv2d, 4 × 10	64	2
5 × 25 × 64	DSCConv	64	1
5 × 25 × 64	DSCConv	64	1
5 × 25 × 64	DSCConv	64	1
5 × 25 × 64	DSCConv	64	1
5 × 25 × 64	AvgPool, 5x25	-	1
1 × 64	Fully Connected	12	-

Table 4: DSCNN(Speech Commands) architecture definition

Input	Operator	#out	stride
96 × 96 × 3	Conv2d, 3 × 3	8	2
48 × 48 × 8	DSCConv	16	1
48 × 48 × 16	DSCConv	32	2
24 × 24 × 32	DSCConv	32	1
24 × 24 × 32	DSCConv	64	2
12 × 12 × 64	DSCConv	64	1
12 × 12 × 64	DSCConv	128	2
6 × 6 × 128	DSCConv	128	1
6 × 6 × 128	DSCConv	128	1
6 × 6 × 128	DSCConv	128	1
6 × 6 × 128	DSCConv	256	2
3 × 3 × 256	DSCConv	256	1
3 × 3 × 256	AvgPool, 3 × 3	-	1
1 × 256	Fully Connected	2	-

Table 5: MobileNetV1(VWW) architecture definition

Input	Operator	#out	stride
128 × 9	Conv1d, 3	64	1
126 × 64 × 16	Conv1d, 3	64	1
124 × 64	MaxPool	-	2
62 × 64	Flatten	-	-
1 × 3960	Fully Connected	6	-

Table 6: CNN(UCI-HAR) architecture definition