

EFQAT: AN EFFICIENT FRAMEWORK FOR QUANTIZATION-AWARE TRAINING

Saleh Ashkboos^{1*} Bram Verhoef² Torsten Hoefer¹ Evangelos Eleftheriou² Martino Dazzi²

ABSTRACT

Quantization-aware training (QAT) schemes have been shown to achieve near-full precision accuracy. They accomplish this by training a quantized model for multiple epochs. This is computationally expensive, mainly because of the full precision backward pass. On the other hand, post-training quantization (PTQ) schemes do not involve training and are therefore computationally cheap, but they usually result in a significant accuracy drop. We address these challenges by proposing EfQAT, which generalizes both schemes by optimizing only a subset of the parameters of a quantized model. EfQAT starts by applying a PTQ scheme to a pre-trained model and only updates the most critical network parameters while freezing the rest, accelerating the backward pass. We demonstrate the effectiveness of EfQAT on various CNNs and Transformer-based models using different GPUs. Specifically, we show that EfQAT is significantly more accurate than PTQ with little extra compute. Furthermore, EfQAT can accelerate the QAT backward pass between 1.44-1.64x while retaining most accuracy.

1 INTRODUCTION

Large-scale neural networks are at the core of many advances in Artificial Intelligence. Model scaling is of prime importance and a key enabler in achieving ever-increasing performance on a variety of cognitive tasks. Indeed, state-of-the-art models have moved from 60M parameters in 2012 (Krizhevsky et al., 2017) to 540B in 2022 (Chowdhery et al., 2022), meanwhile showing a strong correlation between the increase in capabilities of the Deep Neural Network (DNNs) models and the number of parameters. Consequently, training these increasingly accurate and general models implies increasing hardware requirements in terms of memory, compute, and energy (OpenAI, 2018; Strubell et al., 2019). Given the current trend, future huge models may be infeasible to train with currently available systems.

Previous work has addressed this resource challenge by introducing methods to reduce the model size. These methods often rely on quantization (Gholami et al., 2021) and sparsification (Hoefer et al., 2021). The former tries to reduce the computational and memory cost of neural network training and inference by allocating a lower number of bits to different network variables (e.g. weights and activations), whereas the latter aims at retaining the most critical parameters and pruning the rest. Quantization has gained more

*Work done during an internship at Axelera AI. ¹ETH Zurich ²Axelera AI. Correspondence to: Saleh Ashkboos <saleh.ashkboos@inf.ethz.ch>.

Quantization Schemes	Near FP Accuracy	High Performance
Post-Training Quantization (PTQ)	✗	✓
Quantization-Aware Training (QAT)	✓	✗
EfQAT	✓	✓

Table 1. Different quantization schemes: EfQAT achieves higher accuracy than PTQ and is faster than QAT. FP refers to the full precision model.

attention as it is easier to be applied in practice with the available hardware. As the dynamic range of activations and weights varies considerably across the network layers, the main challenge of quantizing neural networks lies in mapping real-valued, continuous operands and parameters to a discrete integer-valued domain without significant distortion.

Post-training quantization (PTQ) schemes quantize models in a greedy fashion. Specifically, these algorithms usually optimize an auxiliary loss function representing the distance between the quantized and unquantized operands and parameters in the network. The optimization is applied on a *per-layer* basis using only a few training samples, known as the calibration set. Although these approaches are computationally inexpensive, they often decrease network accuracy.

Quantization-aware training (QAT) is another class of neural network quantization algorithms that tries to make the above transformation towards lower precision networks easier by training the quantized network directly. These

schemes jointly optimize the quantization- and network parameters using standard gradient-based methods (like Stochastic Gradient Descent) on the original network loss. QAT methods have the benefit of higher accuracy compared to PTQ methods. However, QAT methods are computationally expensive as they only accelerate the forward pass, while the backward pass should be performed in high precision, which has twice the number of operations as the forward pass. Thus, QAT methods are not computationally efficient for large models nor are suitable for computationally-constrained environments.

To address these challenges, we generalize the above quantization schemes and introduce Efficient QAT (**EfQAT**). EfQAT starts from a quantized model (output of an arbitrary quantization scheme) and fine-tunes the quantization parameters as well as *the most important* weights of the network. EfQAT aims to combine the benefits of QAT and PTQ as it accelerates the backward pass of the training process up to 2x compared to standard QAT schemes and it also improves the accuracy relative to PTQ schemes (Table 1) We summarize our contributions as follows:

- We propose EfQAT, a general quantization framework that optimizes the quantization-related parameters and the most critical network weights. We define a simple metric to extract such weights and estimate the importance of the weight channels (rows) in the convolution (linear) layers. Our scheme can be used on top of any PTQ and QAT scheme, improving their accuracy (for PTQ schemes) and performance (for QAT schemes).
- We demonstrate that, when EfQAT is applied to networks with 4-bit weights and 8-bit activations, we obtain a 3% improvement over PTQ schemes using ResNet-50 on ImageNet, and 6 F1 on the SQuAD using BERT_{base}) while updating only 5% of the weights.
- We discuss the theoretical speedup of EfQAT and show that our approach can accelerate the backward pass of QAT schemes by up to 2x. We validate the practical performance of EfQAT on various tasks and achieve up to 1.64x speedup for convolution-based models and 1.45x for Transformer-based models compared to the quantization-aware training schemes.

2 RELATED WORK

Quantization methods for accelerating neural network inference can be divided into two main categories (Gholami et al., 2021).

Post-Training Schemes. Given a pre-trained model, the post-training quantization (PTQ) extract the quantization-related parameters (scales and zero points) without (or with

limited) fine-tuning of the network parameters (Frantar et al., 2022b; Dettmers et al., 2023). In the simplest case, according to the MinMax (Krizhevsky et al., 2009) observer, the quantization parameters are computed based on the minimum and maximum values of weights and activations. This approach (Nagel et al., 2021) does not need additional fine-tuning and has successfully been applied to various networks such as ResNet (He et al., 2016) and MobileNet (Sandler et al., 2018) as well as large Transformer models (Vaswani et al., 2017). Moreover, PTQ algorithms have been further enhanced with network-specific considerations such as keeping the outlier features in higher precision (Dettmers et al., 2022; Ashkboos et al., 2023) or smoothing the activations before quantization (Xiao et al., 2022; Ashkboos et al., 2024).

More sophisticated schemes try to overcome the limitations of the post-training methods by minimizing the distance between the output of the quantized and full precision layer. For example, OMSE (Choukroun et al., 2019) uses the layer-wise MSE loss function to find the scale factors of the weights and activations, whereas OBC (Frantar & Alislarh, 2022) uses second-order information to optimize the above loss function. Specifically, it quantizes the weights iteratively and updates the unquantized weights at each iteration to compensate for the quantization error of the quantized weights. Finally, GPTQ (Frantar et al., 2022a) accelerates this process by simultaneously quantizing all the weights in each column and compressing the parameters of large models like OPT-175B (Zhang et al., 2022) and BLOOM (Laurençon et al., 2022) into 3-4 bits.

Quantization-Aware Training Schemes. Unlike post-training schemes, quantization-aware training (QAT) schemes jointly train the model and the quantization parameters by approximating the backward gradient of the rounding operation. Straight Through Estimator (STE) (Bengio et al., 2013) is the most popular approximation for the rounding function, which approximates the gradient by using an identity function. In order to quantize CNNs, DoReFa (Zhou et al., 2016) uses the STE to approximate the gradient of the rounding function and quantizes both weights and activations. PACT (Choi et al., 2018) uses a trainable clipping to compress the dynamic range of the activations and reduce their quantization error. LSQ (Esser et al., 2019) extends the STE approximation to make it sensitive to the distance between the full precision and quantized values. By means of this approximation, LSQ learns a scaling function for quantizing the weights and activations. Finally, TQT (Jain et al., 2020) stabilizes the QAT schemes by proposing an activation function for the quantization parameters. SQuAT (Wang et al., 2022) quantizes the Transformer-based models using sharpness-aware optimization (Foret et al., 2020) to push the quantized weights into flat minima. This approach enables

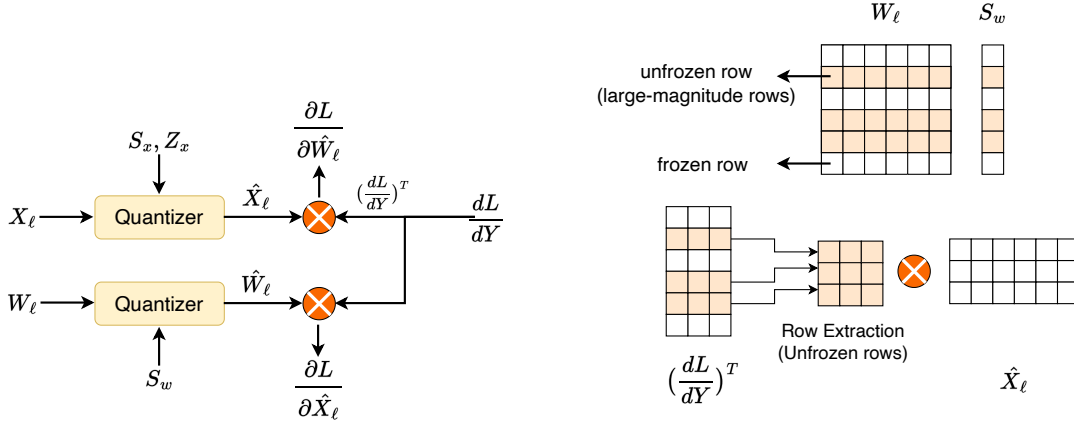


Figure 1. Main backward matrix multiplications of the quantized layer with Symmetric weight (using scale S_w) and Asymmetric input (using scale S_x and zero point Z_x) quantization. **Left:** Quantization-aware training applies both matrix multiplications in full precision. **Right:** EfQAT accelerates the backward pass by performing the matrix multiplication only over the most important/unfrozen rows (the rows with large average magnitude).

the quantization of BERT with 2 to 4 bit representation of parameters.

Our work is a generalization of both quantization schemes, in which we only optimize the most important network parameters alongside with the quantization-related parameters. Using our scheme, we show that with at most one epoch of EfQAT, one can obtain near-full precision accuracy on various models and datasets.

3 EfQAT

Firstly, in this section, we provide a brief introduction to quantization and discuss the main computational bottleneck of the QAT schemes. Then, we present EfQAT and discuss its properties in detail. Finally, we provide some theoretical analysis of the speedup upper bound for our scheme.

3.1 Background

Low-Precision Inference. When the distributions of the weights and activations do not have long tails (i.e., a wide dynamic range is not needed), uniformly spaced quantization points provide a practical choice (Jacob et al., 2018; Wu et al., 2020). Integer quantization provides uniform quantization that can be implemented easily in most hardware. In this work, we use the *round-to-nearest* operation to transform full-precision weights and activations into low-precision values.

Symmetric vs. Asymmetric Quantization. Following (Nagel et al., 2021), we use Asymmetric Quantization for quantizing the input (activations) and Symmetric Quantiza-

tion for the weights. In the former case, we round the input tensor to the b -bit integer using

$$\hat{X} = \max \left(\min \left(\left\lceil \frac{X}{S_x} \right\rceil + Z_x, 2^b - 1 \right), 0 \right), \quad (1)$$

where X represents the activations, S_x and Z_x are the scales and zero points, known as *quantization parameters*, and $\lceil \cdot \rceil$ is the round-to-nearest operation. The quantization parameters are calculated using

$$S_x = \frac{\beta_x - \alpha_x}{2^b - 1}, \quad Z_x = -\left\lceil \frac{\alpha_x}{S_x} \right\rceil, \quad (2)$$

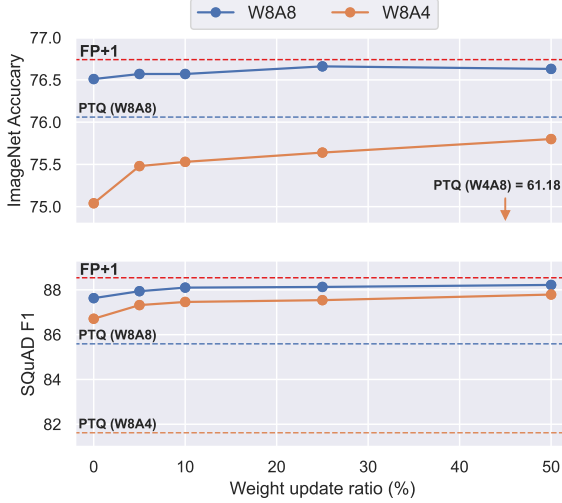
where $[\alpha_x, \beta_x]$ is the quantization range of the activations, estimated by performing the forward pass on the calibration set. We quantize the weights using

$$\hat{W} = \max \left(\min \left(\left\lceil \frac{W}{S_w} \right\rceil, -2^{b-1} + 1 \right), 2^{b-1} - 1 \right). \quad (3)$$

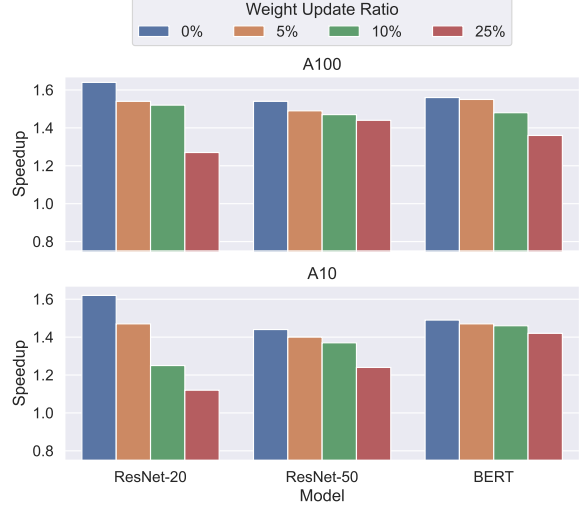
In this case, the quantization range of the weights W is $[\alpha_w, \beta_w]$ and the zero points of the weights are always zero ($Z_w = 0$). The S_w is calculated using

$$S_w = \frac{\max(|\alpha_w|, |\beta_w|)}{2^{b-1} - 1}. \quad (4)$$

Computational Bottleneck of QAT Schemes. Given a linear layer ℓ with weight matrix W_ℓ and input tensor



(a) EfQAT-CWPN improves the accuracy over PTQ in different precisions. **FP+1** is the accuracy of training the full precision checkpoint for one more epoch.



(b) EfQAT-LWPN backward speedup over QAT on A100 and A10 GPUs. The backward runtime is calculated over the total training steps during the EfQAT epoch.

Figure 2. Accuracy and the performance of EfQAT-CWPN/LWPN on the ImageNet (using ResNet-50), SQuAD (using BERT_{base}), and CIFAR-10 (using ResNet-20) datasets.

X_ℓ , a QAT scheme performs the forward pass matrix-multiplication $Y_\ell = \hat{X}_\ell \hat{W}_\ell^T$ in low precision¹, where \hat{X}_ℓ and \hat{W}_ℓ are the quantized input and weights respectively. However, the backward pass requires two full precision matrix multiplications (Figure 1 left)

$$\frac{\partial L}{\partial \hat{X}_\ell} = \frac{\partial L}{\partial Y_\ell} \hat{W}_\ell, \quad \frac{\partial L}{\partial \hat{W}_\ell} = \left(\frac{\partial L}{\partial Y_\ell} \right)^T \hat{X}_\ell, \quad (5)$$

where $\frac{\partial L}{\partial X_\ell}$, $\frac{\partial L}{\partial Y_\ell}$, and $\frac{\partial L}{\partial W_\ell}$ are the gradients with respect to the inputs, outputs, and the weights of the layer. These matrix multiplications dominate the computations in QAT schemes, limiting their practical applications.

3.2 EfQAT Main Idea

Motivation. As mentioned in the previous section, the runtime of QAT schemes is dominated by two full precision matrix multiplications as in Equation (5). The second matrix multiplication can be accelerated by freezing a large portion of the weight matrix and only calculating the gradient for the rest. However, this raises two significant challenges: First, we must define a metric to freeze non-important weights during the backward pass. In addition, unstructured weight freezing does not yield practical speedup, and we need to enforce some structure on the frozen parameters. We address these challenges and define EfQAT in this section.

¹We follow the definition of Linear layer in the PyTorch framework (Paszke et al., 2019).

Weight Importance. Inspired by previous pruning work (Hoeffler et al., 2021), we define the *importance* of a given block B (with size n) of weights from a network layer by the average magnitude of its weights. More precisely, we define

$$\mathcal{I}_B := \frac{1}{n} \sum_{w \in B} |w|, \quad (6)$$

as the metric to decide on freezing a block B of weights.

Structured Weight Freezing. Our goal is to accelerate weight gradient calculation in Equation (5) by only updating a small set of weights during the QAT training step. The most natural case is to consider each channel (row) of the convolution (linear) layer as a block and freeze the less-important weight blocks in each layer (or across the whole network). Figure 3 shows the \mathcal{I}_B when we consider the channels (rows) in the convolution (linear) layers as the weight blocks for ResNet-20 and BERT_{base}. It suggests that there are a few important channels in different network layers.

EfQAT Modes. Based on the granularity structure we use to freeze the network parameters, we define three modes for EfQAT as in Table 2. In the lowest granularity level, Channel-Wise Per-Layer (EfQAT-CWPL), we freeze the less-important *channels* in each *layer*. In the Channel-Wise Per-Network mode (EfQAT-CWPN), we balance the frozen

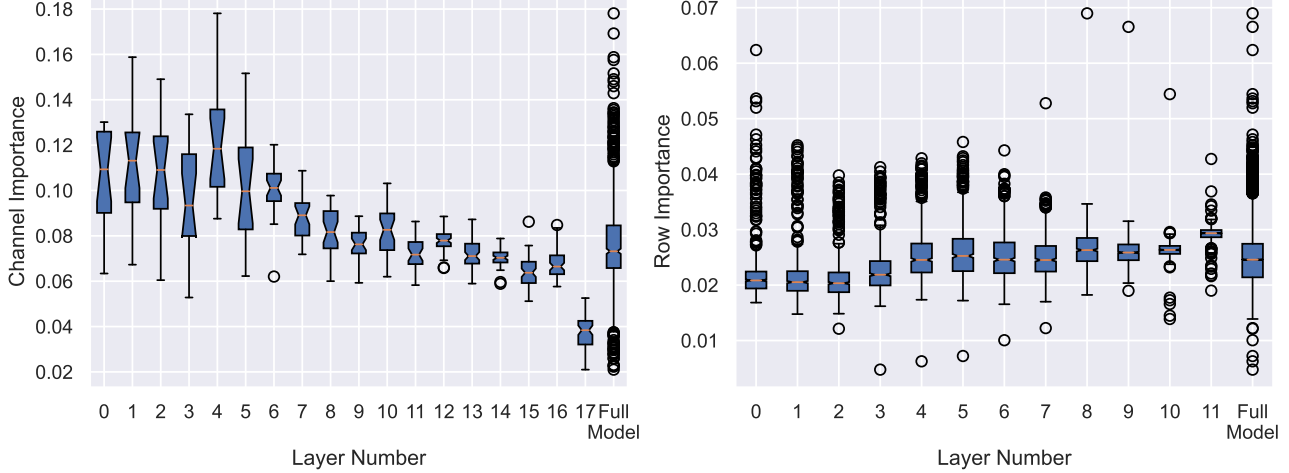


Figure 3. The importance of different channels of convolutions in ResNet-20 (left) and the rows of the output matrix of the self-attention layers in BERT_{base} (right). A significant amount of outliers can be noted for both networks across different layers. In both plots, the last column displays the channel/row importance for the whole network.

Mode	Freezing Granularity	Frozen Weight Selection
EfQAT-CWPL	Channel-Wise	Per-Layer
EfQAT-CWPN	Channel-Wise	Per-Network
EfQAT-LWPN	Layer-Wise	Per-Network

Table 2. Three modes of EfQAT: We gradually increase the granularity of the frozen weights and the selection process.

channels across the whole network. Finally, in the Layer-Wise Per-Network mode (EfQAT-LWPN), we freeze the entire weights of a layer in the network. In all cases, we use the Top-K algorithm to find the important weights.

Freezing Frequency. The importance of the weights will change as long as we update the network parameters. However, updating the \mathcal{L}_B adds overhead to EfQAT (especially in EfQAT-CWPL and EfQAT-CWPN) as we need to iterate over all the unfrozen channels, update their importance, and then select the most important blocks. We found that we can update the frozen channels once after processing every \mathbf{f} data sample without too much accuracy drop. More specifically, we see a negligible accuracy drop if we update the weight importance of every $\mathbf{f} = 4096$ training sample (see Section 4).

Quantization Granularity. We use the per-channel (per-row) quantization for the weights in the convolution (linear) layers. We update the quantization parameters of only if we update the weights of that channel (unfrozen weights).

We use per-tensor quantization for all the activations and apply our freezing scheme over the output channels of the convolution layer and the rows (corresponding to the output features) of the weight matrix in the linear layers.

3.3 Algorithm Pseudocode

Now, we present the full Pseudocode of EfQAT (EfQAT-CWPL and EfQAT-CWPN) in Algorithm 1. The EfQAT-LWPN performs the same and we only skip the weight gradient calculation for the frozen layers.

Algorithm 1 EfQAT on a Neural Network with L layers with ratio k .

```

Start from a PTQ model           {QParam initialization}
for Each Batch  $B$  in Training Set do
  Apply the Forward pass in low precision using  $B$ .
  {Backward Step:}
  for For  $\ell = L, \dots, 1$  do
    Calculate the input gradient using  $\frac{\partial L}{\partial \hat{X}_\ell} = \frac{\partial L}{\partial Y_\ell} \hat{W}_\ell$ .
    Extract the unfrozen channels using  $k$  and save their indices in  $id$ .
    Calculate the unfrozen weight gradients using
       $\frac{\partial L}{\partial \hat{W}_\ell}[id] = \frac{\partial L}{\partial Y_\ell}^T \hat{X}_\ell[id]$ .
    end for
  {Optimizer Step:}
  Update the unfrozen weight channels and their scales  $S_w$ .
  Update the scales and zero points of the activations  $S_x, Z_x$ .
  Update the biases and normalization layer.
end for
    
```

3.4 Theoretical Speed-Up

We consider the theoretical speed-up of EfQAT-CWPL and EfQAT-CWPN with unfrozen ratio $0 \leq r \leq 1$ on two types of layers: linear layers and convolution layers:

Consider a linear layer with the weight matrix $W \in \mathbb{R}^{C_{out} \times C_{in}}$, the input vector to the layer $X \in \mathbb{R}^{C_{in}}$, and the output gradient vector $\frac{\partial L}{\partial W_\ell} \in \mathbb{R}^{C_{out}}$. The total number of operations in both backward matrix multiplications is $2C_{in}C_{out}$.

As EfQAT calculates the gradients of $\lfloor rC_{out} \rfloor$ rows, the total number of operations can be computed as

$$\text{OPS(BWD)}_{\text{Linear}} = \underbrace{C_{in} \lfloor rC_{out} \rfloor}_{\text{grad. w.r.t. weights}} + C_{in}C_{out} \leq (1+r)C_{in}C_{out}. \quad (7)$$

Consider a convolution layer with kernel size k and matrix $\in \mathbb{R}^{C_{out} \times C_{in} \times k \times k}$. Given an input $X \in \mathbb{R}^{C_{in} \times H_{in} \times W_{in}}$, the convolution produces an output $\in \mathbb{R}^{C_{out} \times H_{out} \times W_{out}}$. Formulating the backward pass similarly to the case of the linear layer, the operation count is equal to:

$$\text{OPS(BWD)}_{\text{Conv}} = \underbrace{k^2 C_{in} \lfloor rC_{out} \rfloor H_{out} W_{out}}_{\text{grad. w.r.t. weights}} + k^2 C_{in} C_{out} H_{out} W_{out} \leq (1+r)k^2 C_{in} C_{out} H_{out} W_{out}. \quad (8)$$

The upper bounds in (7) and (8) show that we can achieve up to 2x speedup in the backward operations in the linear and convolution layers with $r = 0$. However, in practice, the QAT backward pass has other operations (e.g., approximating the rounding operator, statistical normalization, and element-wise operations during the quantization parameter gradient calculation), which we do not count in our theoretical analysis.

The minimum and maximum speedups in this section can also be applied to the EfQAT-LWPN. In that case, we have $r \in \{0, 1\}$ as we freeze all weights in a layer. However, in practice, EfQAT-CWPL and EfQAT-CWPN need to extract the rows of the matrices which has memory access (due to data movement) overhead.

4 EXPERIMENTAL RESULTS

Setup. We evaluate EfQAT on a variety of tasks and models. For each task, we start from a full-precision pre-trained checkpoint (**FP**) and quantize it using a **PTQ** scheme. After that, we apply one epoch of **EfQAT** (CWPL, CWPN, or

Model	Full Precision		PTQ	
	FP	FP+1	Bit-Width	Accuracy
ResNet-20	91.71	91.74	W8A8	91.69 ± 0.03
			W4A8	88.17 ± 0.07
			W4A4	80.75 ± 0.41
ResNet-50	76.12	76.74	W8A8	76.06
			W4A8	61.18
			W4A4	19.12
BERT _{base}	88.07	88.54	W8A8	85.59 ± 0.79
			W4A8	81.62 ± 1.17

Table 3. Overview of the baseline models. We use Top-1 accuracy for the ResNet-20 (on CIFAR-10) and ResNet-50 (on ImageNet) and F1 score for BERT_{base} (on SQuAD).

LWPN) to fine-tune the quantized model. As EfQAT applies one training epoch, in addition to the PTQ model, we train the FP model for one more epoch in full precision (**FP+1**) and compare the accuracy against it. To update the network parameters in EfQAT, we use the same optimizer (and hyper-parameters) as FP+1 and always use Adam (Kingma & Ba, 2014) to update the quantization parameters. We use a calibration dataset consisting of 512 samples in all our experiments and STE (Bengio et al., 2013) to approximate the gradient of the rounding function during the EfQAT training epoch. We quantize all convolutions and linear layers (including the input, output, and shortcut layers) of the CNNs but do not quantize the embedding layer in the BERT model. For the CIFAR-10 and SQuAD datasets, we repeat all our experiments with three random seeds and present the mean and standard deviation of the outcomes. However, as we do not observe significant variation (more than 0.01) across different seeds on the ImageNet dataset, we report a single number for those experiments. We implement EfQAT using PyTorch (Paszke et al., 2019) framework. We provide our code and the related instructions in Appendix ??.

CIFAR-10. We evaluate EfQAT on CIFAR-10 using the ResNet-20 network from (He et al., 2016). We train the model for 200 epochs using SGD with momentum 0.9 and 1e-4 weight decay to extract the FP checkpoint. The initial learning rate is 0.1, and we multiply it by 0.1 at epochs 100 and 150. We use the same optimizer as FP for the network parameters (with all states and hyperparameters) during the EfQAT epoch. To update the quantization parameters, we use a learning rate 1e-6.

ImageNet. We evaluate EfQAT on ImageNet using ResNet-50. We use a pre-trained model from torchvision (TorchVision, 2016) and train the parameters using SGD with a 1e-3 learning rate. The quantization parameters

are optimized with a $1e-7$ learning rate.

Question Answering on SQuAD. We evaluate EfQAT on BERT_{base} (Devlin et al., 2018) for question answering task on SQuAD v1.1 (Rajpurkar et al., 2016). To extract the FP checkpoint, we adopt the pre-trained model from the HuggingFace Transformers library (Wolf et al., 2019), and fine-tune it for two epochs using the same settings as (Chen et al., 2020). We exclude 4-bit quantization of the activations and features of BERT (i.e., W4A4) as it leads to a 4.9 F1 score in the QAT experiment. Unlike QAT, we do not update the embedding layer during EfQAT and use a learning rate of $1e-6$ for the quantization parameters.

PTQ Baseline. For the PTQ baseline, we use the Min-Max (Krizhevsky et al., 2009) observer for both the weights and the activations. We apply per-channel (per-row) symmetric quantization for the weights in the convolution (linear) layers and per-tensor asymmetric quantization for the activations. The summary of our baselines is presented in Table 3.

Hyperparameter Exploration. We do not apply any hyperparameter optimization to our experiments. However, we evaluate EfQAT using different learning rates for the quantization parameters. In addition, following (Jain et al., 2020), we train the logarithm of the quantization scales in place of training them directly, to evaluate the stability of the EfQAT. In all cases, we do not observe a significant accuracy drop for EfQAT. See Appendix A.2 for all the results and details.

4.1 Main Results

Accuracy Results. We first study the role of EfQAT on the accuracy of various models. To this end, for a given weight-update ratio, we consider the three modes of EfQAT (CWPL, CWPN, and LWPN). In each case, we apply our scheme with various weight-update ratio and compare it against QAT (where we update all network parameters and quantization parameters) and PTQ. Table 4 summarizes our results for different networks. Interestingly, there is a jump in the accuracy in the 0% case where we only update the quantization parameters and the computationally light-weight parameters (like biases and the normalization layers). In addition, our results show that the changes in accuracy between different EfQAT modes are negligible in almost all the cases, which shows that the acceleration provided by EfQAT-LWPN comes at a minor accuracy cost. Finally, we observe that updating 25% of the network parameters can achieve almost the same accuracy (with $\leq 0.5\%$ drop) as training the full network with QAT.

Freezing Frequency. As mentioned in Section 3.2, we need to update the frozen channels repeatedly during EfQAT

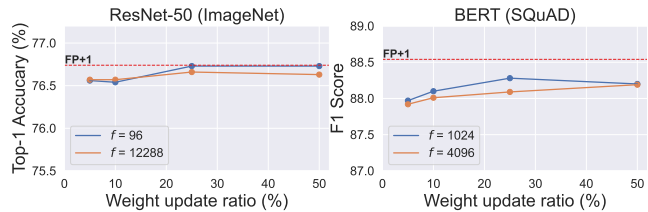


Figure 4. The role of different freezing intervals on the accuracy of EfQAT-CWPN with W8A8. We update the frozen channels every f samples. Larger f does not cause a large accuracy drop during the EfQAT training epoch.

training steps, which causes a computational overhead. To evaluate the impact of the update frequency on the accuracy, we define a frequency ratio f and update the weight importances (and thus the frozen channels) for every $f > 1$ training examples. Figure 4 shows that using large f causes minor accuracy drops. Consequently, we conclude that EfQAT does not rely on frequent updates of the weight importances, making their computational overhead easily amortizable. Additional models and results are in Appendix A.1.

Practical Speedups. To evaluate the practical speedup of EfQAT, we run all the experiments on single A10 and A100 GPUs. We use training batch size 128 for CIFAR-10 dataset and 96 for the ImageNet dataset. For SQuAD, we use batch size 8^2 for the experiments on the A10 GPU and 16 for the ones on the A100. Table 5 compares the running time of EfQAT-CWPN and EfQAT-LWPN. We only report the speedup over the backward pass as the forward pass in QAT may be performed in low precision with hardware-dependent runtime. However, as the backward pass is the main computational bottleneck of QAT, we expect a higher end-to-end speedup when the forward pass is performed in low precision with quantized weights and activations.

EfQAT can achieve up to 1.64x speedup on A100 (1.62x speedup on A10) over the QAT by only updating the quantization parameters of the activations (0% case). Note that we cannot obtain the full 2x theoretical speedup for two reasons: firstly, to calculate the gradients of the unfrozen weights, we need to (1) extract the corresponding rows from the output gradient matrix (Figure 1 right), (2) calculate the matrix multiplication using extracted rows, and (3) finally, save the calculated gradients to the corresponding rows in the weight gradient matrix. However, the above first and last steps cause overhead due to the data movement, limiting the achievable speedup. EfQAT-LWPN solves this issue by freezing the whole weights of a layer, resulting

²The largest batch size on an A10 GPU.

EfQAT: An Efficient Framework for Quantization-Aware Training

Model	Bit-Width	PTQ	Mode	EfQAT Weight Update Ratio (%)					QAT
				0	5	10	25	50	
BERT _{base}	W8A8	85.59 ± 0.79	CWPL CWPN LWPN	87.63 ± 0.04	87.94 ± 0.09 87.94 ± 0.07 87.81 ± 0.05	87.98 ± 0.04 88.10 ± 0.01 87.94 ± 0.19	88.05 ± 0.03 88.13 ± 0.14 88.07 ± 0.16	88.27 ± 0.02 88.22 ± 0.06 88.21 ± 0.02	88.25 ± 0.12
	W4A8	81.62 ± 1.17	CWPL CWPN LWPN	86.71 ± 0.04	87.30 ± 0.04 87.32 ± 0.17 87.22 ± 0.15	87.42 ± 0.16 87.46 ± 0.10 87.26 ± 0.02	87.53 ± 0.13 87.54 ± 0.06 87.56 ± 0.15	87.82 ± 0.03 87.79 ± 0.10 87.72 ± 0.12	87.84 ± 0.04
ResNet-20	W8A8	91.69 ± 0.03	CWPL CWPN LWPN	91.71 ± 0.04	91.72 ± 0.07 91.73 ± 0.01 91.68 ± 0.04	91.74 ± 0.04 91.70 ± 0.07 91.69 ± 0.03	91.75 ± 0.05 91.75 ± 0.06 91.70 ± 0.01	91.71 ± 0.02 91.72 ± 0.02 91.70 ± 0.06	91.72 ± 0.05
	W4A8	88.17 ± 0.07	CWPL CWPN LWPN	91.12 ± 0.04	91.14 ± 0.15 91.19 ± 0.11 91.21 ± 0.04	91.22 ± 0.07 91.16 ± 0.10 91.24 ± 0.03	91.25 ± 0.07 91.35 ± 0.11 91.32 ± 0.07	91.24 ± 0.13 91.32 ± 0.11 91.41 ± 0.10	91.45 ± 0.14
	W4A4	80.75 ± 0.41	CWPL CWPN LWPN	87.23 ± 0.34	87.66 ± 0.21 87.50 ± 0.21 87.24 ± 0.14	87.59 ± 0.23 87.78 ± 0.19 87.60 ± 0.28	87.85 ± 0.31 87.86 ± 0.21 87.89 ± 0.07	87.79 ± 0.13 87.97 ± 0.20 88.35 ± 0.08	88.28 ± 0.37
ResNet-50	W8A8	76.06	CWPL CWPN LWPN	76.51	76.65 76.57 76.53	76.58 76.57 76.51	76.65 76.66 76.63	76.74 76.63 76.68	76.80
	W4A8	61.18	CWPL CWPN LWPN	75.04	75.31 75.48 75.04	75.46 75.53 75.12	75.72 75.64 75.32	75.85 75.80 75.74	76.02
	W4A4	19.12	CWPL CWPN LWPN	66.60	67.41 67.49 67.35	67.77 68.0 67.7	68.32 68.71 68.33	68.75 69.22 68.14	69.31

Table 4. Accuracy of EfQAT on CIFAR-10 (using ResNet-20), ImageNet (using ResNet-50), and SQuAD v1.1 (using BERT_{base}). For each model, we evaluate different EfQAT modes, including Channel-Wise Per-Laye (CWPL), Channel-Wise Per-Network (CWPN), and Layer-Wise Per-Network (LWPN).

in higher speedup. Secondly, the 2x theoretical speedup refers to a single convolution or linear layer. However, the backward pass of the quantized models includes other operations, such as normalizations and element-wise operations, which have been shown to constitute up to 30% of the runtime in BERT (Ivanov et al., 2021). While this is outside of the scope of this work, such operations could be optimized with techniques such as kernel fusion (Wahib & Maruyama, 2014; Ben-Nun et al., 2019).

5 CONCLUSION

We propose EfQAT, a general framework to accelerate quantization-aware training schemes. EfQAT updates only the most essential weights during the backward pass and freezes the rest. We use a simple yet efficient magnitude-based metric to choose the essential weights and freeze the rest using different granularity levels.

We apply EfQAT on different models and datasets and show that EfQAT can accelerate the QAT backward pass up to 1.64x on an A100 GPU (and 1.62x on an A10 GPU) with less than 0.3% accuracy degradation for ResNet-50 on ImageNet and 0.62 F1 score degradation on SQuAD for

Model (GPU)	Mode	Frequency (f)	EfQAT Weight Update Ratio (%)				QAT
			0	5	10	25	
ResNet-20 (A10)	CWPN	16384	2.61	3.46	3.92	4.92	4.23
	LWPN	16384		2.87	3.38	3.77	
ResNet-20 (A100)	CWPN	16384	3.73	5.69	7.1	9.98	6.22
	LWPN	16384		4.02	4.08	4.86	
BERT _{base} (A10)	CWPN	4096	1784	2219	2271	2339	2667
	LWPN	4096		1811	1821	1872	
BERT _{base} (A100)	CWPN	4096	1009	1146	1173	1271	1559
	LWPN	4096		1043	1055	1080	
ResNet-50 (A10)	CWPN	12288	3575	6079	6247	6335	5139
	LWPN	12288		3648	3748	4121	
ResNet-50 (A100)	CWPN	12288	1402	2505	2545	2724	2189
	LWPN	12288		1411	1474	1606	

Table 5. Backward runtime (in second) of EfQAT-CWPN and EfQAT-LWPN on different networks and datasets for both A100 and A10 GPUs.

BERT_{base}. Moreover, on the same datasets, we improve the accuracy of ResNet-50 by 0.45% and of BERT_{base} by 2.03 F1 score compared to PTQ.

REFERENCES

- Ashkboos, S., Markov, I., Frantar, E., Zhong, T., Wang, X., Ren, J., Hoefler, T., and Alistarh, D. Towards end-to-end 4-bit inference on generative large language models. *arXiv preprint arXiv:2310.09259*, 2023.
- Ashkboos, S., Mohtashami, A., Croci, M. L., Li, B., Jaggi, M., Alistarh, D., Hoefler, T., and Hensman, J. Quarot: Outlier-free 4-bit inference in rotated llms. *arXiv preprint arXiv:2404.00456*, 2024.
- Ben-Nun, T., de Fine Licht, J., Ziogas, A. N., Schneider, T., and Hoefler, T. Stateful dataflow multigraphs: A data-centric model for performance portability on heterogeneous architectures. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–14, 2019.
- Bengio, Y., Léonard, N., and Courville, A. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Chen, T., Frankle, J., Chang, S., Liu, S., Zhang, Y., Wang, Z., and Carbin, M. The lottery ticket hypothesis for pre-trained bert networks. *Advances in neural information processing systems*, 33:15834–15846, 2020.
- Choi, J., Wang, Z., Venkataramani, S., Chuang, P. I.-J., Srinivasan, V., and Gopalakrishnan, K. Pact: Parameterized clipping activation for quantized neural networks. *arXiv preprint arXiv:1805.06085*, 2018.
- Choukroun, Y., Kravchik, E., Yang, F., and Kisilev, P. Low-bit quantization of neural networks for efficient inference. In *ICCV Workshops*, pp. 3009–3018, 2019.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Dettmers, T., Lewis, M., Belkada, Y., and Zettlemoyer, L. Llm.int8(): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- Dettmers, T., Svirschevski, R., Egiazarian, V., Kuznedelev, D., Frantar, E., Ashkboos, S., Borzunov, A., Hoefler, T., and Alistarh, D. Spqr: A sparse-quantized representation for near-lossless llm weight compression. *arXiv preprint arXiv:2306.03078*, 2023.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Esser, S. K., McKinstry, J. L., Bablani, D., Appuswamy, R., and Modha, D. S. Learned step size quantization. *arXiv preprint arXiv:1902.08153*, 2019.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Frantar, E. and Alistarh, D. Optimal brain compression: A framework for accurate post-training quantization and pruning. *arXiv preprint arXiv:2208.11580*, 2022.
- Frantar, E., Ashkboos, S., Hoefler, T., and Alistarh, D. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022a.
- Frantar, E., Ashkboos, S., Hoefler, T., and Alistarh, D. GPTQ: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022b.
- Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M. W., and Keutzer, K. A survey of quantization methods for efficient neural network inference. *arXiv preprint arXiv:2103.13630*, 2021.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hoefler, T., Alistarh, D., Ben-Nun, T., Dryden, N., and Peste, A. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *arXiv preprint arXiv:2102.00554*, 2021.
- Ivanov, A., Dryden, N., Ben-Nun, T., Li, S., and Hoefler, T. Data movement is all you need: A case study on optimizing transformers. *Proceedings of Machine Learning and Systems*, 3:711–732, 2021.
- Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H., and Kalenichenko, D. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2704–2713, 2018.
- Jain, S., Gural, A., Wu, M., and Dick, C. Trained quantization thresholds for accurate and efficient fixed-point inference of deep neural networks. *Proceedings of Machine Learning and Systems*, 2:112–128, 2020.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.

- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- Laurençon, H., Saulnier, L., Wang, T., Akiki, C., del Moral, A. V., Le Scao, T., Von Werra, L., Mou, C., Ponferrada, E. G., Nguyen, H., et al. The bigscience roots corpus: A 1.6 tb composite multilingual dataset. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- Nagel, M., Fournarakis, M., Amjad, R. A., Bondarenko, Y., van Baalen, M., and Blankevoort, T. A white paper on neural network quantization. *arXiv preprint arXiv:2106.08295*, 2021.
- OpenAI. AI and compute, 2018. URL <https://openai.com/blog/ai-and-compute/>.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., and Chen, L.-C. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- Strubell, E., Ganesh, A., and McCallum, A. Energy and policy considerations for deep learning in NLP. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- TorchVision, maintainers, c. Torchvision: Pytorch’s computer vision library. <https://github.com/pytorch/vision>, 2016.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Wahib, M. and Maruyama, N. Scalable kernel fusion for memory-bound gpu applications. In *SC’14: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 191–202. IEEE, 2014.
- Wang, Z., Li, J. B., Qu, S., Metze, F., and Strubell, E. Squat: Sharpness-and quantization-aware training for bert. *arXiv preprint arXiv:2210.07171*, 2022.
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- Wu, H., Judd, P., Zhang, X., Isaev, M., and Micikevicius, P. Integer quantization for deep learning inference: Principles and empirical evaluation. *arXiv preprint arXiv:2004.09602*, 2020.
- Xiao, G., Lin, J., Seznec, M., Demouth, J., and Han, S. Smoothquant: Accurate and efficient post-training quantization for large language models. *arXiv preprint arXiv:2211.10438*, 2022.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Zhou, S., Wu, Y., Ni, Z., Zhou, X., Wen, H., and Zou, Y. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160*, 2016.

A SUPPLEMENTARY MATERIAL

A.1 Freezing Frequency

As mentioned in Section 3.2, updating the frozen channels (rows) in the convolution (linear) layers could be performed in large intervals. Figure 4 shows the result of two different frequencies over the EfQAT-CWPN in the BERT_{base} and ResNet-50. Table 6 shows our results on ResNet-20 as well on other frequencies in the BERT_{base} model. We conclude that the accuracy drop is negligible across different frequencies.

Model	Frequency (f)	EfQAT Weight Update Ratio (%)			
		5	10	25	50
BERT _{base}	16	87.94 ± 0.07	88.10 ± 0.01	88.13 ± 0.14	88.22 ± 0.06
	1024	87.97 ± 0.07	88.10 ± 0.16	88.28 ± 0.02	88.20 ± 0.11
	2048	87.94 ± 0.04	88.05 ± 0.06	88.07 ± 0.12	88.21 ± 0.09
	4096	87.92 ± 0.09	88.01 ± 0.02	88.09 ± 0.15	88.19 ± 0.14
ResNet-20	128	91.73 ± 0.01	91.70 ± 0.07	91.75 ± 0.06	91.72 ± 0.02
	8192	91.73 ± 0.03	91.70 ± 0.04	91.72 ± 0.07	91.74 ± 0.03
	16384	91.71 ± 0.05	91.66 ± 0.06	91.68 ± 0.04	91.72 ± 0.01
ResNet-50	96	76.56	76.54	76.73	76.73
	12288	76.57	76.57	76.66	76.63

Table 6. The role of different freezing intervals on the accuracy EfQAT-CWPN with W8A8. We update the frozen channels every f samples.

A.2 Optimization Hyperparameters

In the experiments discussed in Section 4, we train the raw quantization parameters of activations and weights directly. While this is a common practice in QAT algorithms, it could lead to numerical instability. For example, while quantization scales are defined in \mathbb{R}^+ , their direct optimization with a non-optimal learning rate may cause them to turn negative. In an attempt to mitigate the problem, some approaches envision the optimization of a function of the quantization parameters in place of the quantization parameters themselves. (Jain et al., 2020) proposes the training of the logarithm of the quantization scale, showing that its use together with Adam optimizer solves the numerical instability and encourages scale invariance for the updates to the quantization scales. In this section, we compare the optimization of quantization scales with and without the logarithm function. We evaluate ResNet-20 on CIFAR10 and ResNet-50 on ImageNet. In order to assess robustness to learning rate variation, we compare the performance of EfQAT with the nominal learning rate (1e-6 and 1e-7 for ResNet-20 and ResNet-50, respectively) and with a learning rate greater by a factor 1e2 for both log and raw quantization scales. Table 7 reports the results. It can be noted how, for ResNet20 on CIFAR10, the training of log and raw quantization parameters appears to yield accuracy results that are all within statistical error. In the case of ResNet-50 on ImageNet, the training of raw quantization parameters with 1e-7 learning rate outperforms training with log quantization parameters. Although the training of the raw quantization parameters is not numerically stable, the numerical instability was not observed in our experiments even when training with significantly different learning rates. In conclusion, our experiments suggest that the performance of EfQAT is robust to variations in learning rate even when training the quantization parameters directly. Even more so, in the comparison between the two methodologies, optimization of the raw parameters always results in equivalent or higher accuracy compared to optimization of their log counterparts in our experiments.

Model	QParam Func.	LR	EfQAT Weight Update Ratio (%)			
			0	5	10	25
ResNet-20	-	1e-6	91.13 ± 0.08	91.10 ± 0.12	91.17 ± 0.13	91.32 ± 0.19
		1e-4	91.19 ± 0.13	91.26 ± 0.16	91.34 ± 0.06	91.29 ± 0.09
	log ₂	1e-6	91.14 ± 0.10	91.11 ± 0.13	91.17 ± 0.15	91.34 ± 0.18
		1e-4	91.13 ± 0.09	91.26 ± 0.12	91.21 ± 0.11	91.32 ± 0.19
ResNet-50	-	1e-7	76.57	76.57	76.66	76.63
		1e-5	75.07	75.41	75.57	76.20
	log ₂	1e-7	75.08	75.44	75.38	75.97
		1e-5	75.03	75.25	75.44	75.84

Table 7. Accuracy of EfQAT-M2 with different learning rates and activation functions. Following (Jain et al., 2020), we apply the activation functions only over the scales and train the zero points without any activation function.