# L4Q: PARAMETER EFFICIENT QUANTIZATION-AWARE FINE-TUNING ON LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

# Abstract

Due to the high memory and computational costs associated with large language models (LLMs), model compression techniques such as quantization, which reduces inference costs, and parameter-efficient fine-tuning (PEFT) methods like Low-Rank Adaptation (LoRA), which reduce training costs, have gained significant popularity. This trend has spurred active research into quantization-aware PEFT techniques, aimed at maintaining model accuracy while minimizing memory overhead during both inference and training. Previous quantization-aware PEFT methods typically follow a two-step approach: first, applying post-training quantization (PTQ) to model weights, followed by PEFT on the quantized model. However, recovering from the quantization error introduced by PTQ through finetuning has proven challenging. Additionally, most PTQ-based PEFT methods result in a mixture of low-precision quantized weights and high-precision adapter weights, limiting the efficiency of full quantization during inference. While a previous method attempted to address these issues, it still suffers from limited adaptability due to the constrained LoRA parameter structure required to produce fullyquantized models. To overcome these challenges, we propose L4Q, a method that integrates Quantization-Aware Training (QAT) with LoRA to effectively reduce quantization error. By employing a memory-optimized layer design, L4Q significantly reduces QAT's memory overhead while producing fully-quantized weights, enabling effective adaptation to downstream tasks. Our experiments demonstrate that this combined approach to quantization and fine-tuning achieves superior accuracy compared to decoupled fine-tuning schemes, particularly in sub-4-bit quantization, positioning L4Q as an efficient QAT solution. Using the LLaMA model families and instructional datasets, we showcase L4Q's capabilities in language tasks and few-shot learning.

034

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

# 1 INTRODUCTION

036 037 I INTRODUCTION

Given their impressive scalability, Large Language Models (LLMs) such as GPT, OPT, PaLM, and LLaMA (Brown et al., 2020; Zhang et al., 2022; Chowdhery et al., 2024; Touvron et al., 2023a;b) are popular in natural language processing. However, their substantial memory and computational 040 demands pose challenges for practical deployment, making model compression (Han et al., 2015) 041 crucial for LLM deployment. Quantization is a prominent method that reduces model size by low-042 ering the bit precision of model parameters (Hubara et al., 2017), so LLM quantization has been 043 actively studied (Liu et al., 2024; Xiao et al., 2023; Frantar et al., 2023; Dettmers & Zettlemoyer, 044 2023). These quantization methods are generally divided into two categories: post-training quantization (PTQ) and quantization-aware training (QAT). QAT effectively reduces quantization error by integrating quantization into the training process, where both model weights and quantization 046 parameters, are trained together (Esser et al., 2020; Bhalgat et al., 2020). However, applying QAT 047 to LLM quantization is challenging due to the significant memory overhead it incurs during train-048 ing. As a result, PTQ, which applies quantization without retraining the pre-trained model weights and with minimal calibration of the quantization parameters, is widely adopted for LLM quantization (Xiao et al., 2023; Lin et al., 2024; Heo et al., 2024). 051

Concurrently, to enhance the problem-solving abilities of LLMs for specific applications, fine-tuning
 pre-trained LLMs on downstream tasks is crucial as it improves accuracy on target tasks and related
 tasks (Wei et al., 2022; Scialom et al., 2022). However, fine-tuning is a resource-intensive pro-

cess due to the large number of model weight parameters involved. Parameter-efficient fine-tuning (PEFT) addresses this issue (Hu et al., 2022; Li & Liang, 2021; Liu et al., 2022a; Wang et al., 2023) by training a small subset of parameters while freezing the majority of pre-trained weights.
One of the most prominent techniques within PEFT is Low-Rank Adaptation (LoRA) (Hu et al., 2022), which inserts trainable rank decomposition matrices into each layer to represent updates to the frozen weights.

060 The integration of quantization and PEFT 061 holds significant potential for developing 062 highly efficient and accurate LLMs for 063 downstream tasks. Recent research has 064 introduced quantization-aware PEFT approaches to achieve high-quality quantized 065 models (Dettmers et al., 2024b; Kim et al., 066 2024; Xu et al., 2023; Li et al., 2024). 067 Previous works involve a two-stage opti-068 mization strategy: first, a PTQ technique, 069 such as GPTQ (Frantar et al., 2023), is 070 applied to pre-trained LLMs for compres-071 sion. Then, these quantized LLMs un-072 dergo PEFT, such as LoRA, where quan-073 tized weights are kept fixed and only the 074 LoRA parameters are fine-tuned. While 075 fine-tuning can mitigate the effects of quantization errors, separating quantiza-076 tion and fine-tuning into distinct state hin-077



Figure 1: A diagram of model fine-tuning and quantization represented with an example of LLaMA-1 7B model. L4Q produces fast and accurate quantized model.

ders the models from achieving the best accuracy. Furthermore, as high-precision LoRA parameters are adopted alongside the quantized weight matrix, these methods eventually produce mixed-precision models, which limits the efficiency of full quantization during inference. Recently, QA-LoRA (Xu et al., 2023) addresses this issue by strictly constraining the LoRA parameter structure to integrate with quantization parameters, but this constraint can limit the fine-tuning capability.

In this paper, we propose a novel quantization-aware fine-tuning technique, named L4Q (Low-rank 083 adaptive Learning quantization for LLMs). L4Q addresses the limitations of PTQ-based PEFT 084 methods by introducing QAT alongside LoRA. While QAT and LoRA have advantages in reduc-085 ing quantization error and enabling memory-efficient training, respectively, their straightforward 086 integration diminishes the benefits of each approach. Therefore, L4Q carefully integrates these two 087 approaches to properly leverage their advantages. First, L4Q applies the quantization process after 088 fully combining the model weights and LoRA parameters in the linear layer. This approach produces 089 a fully-quantized model that enables memory-efficient and fast inference without limiting the training capabilities of either QAT or LoRA. Moreover, to preserve the memory-efficient nature of LoRA 091 during training, we design the backpropagation path of L4Q to eliminate the need to store weight 092 gradients required for QAT. Finally, the full integration of QAT and LoRA in the proposed L4Q allows for the joint optimization of both the quantization and LoRA parameters, thereby improving 093 of quantized LLMs. As a result, L4Q significantly improves the accuracy of quantized models while maintaining low memory costs during both inference and training, and achieves inference speed comparable to state-of-the-art approaches, making it a more effective solution compared to previous 096 works, as illustrated in Figure 1. 097

098 099

100

# 2 BACKGROUNDS

101 2.1 PEFT WITH LORA

103 LoRA inserts the rank-decomposition matrices composed of a pair of parameter matrices  $A \in R^{r \times i}$ 104 and  $B \in R^{o \times r}$ . Here, *i* and *o* represent the size of input and output dimensions of the original 105 weight matrix, respectively.  $r \ll i, o$  is the rank of the LoRA matrices, and  $\alpha$  is a constant that 106 adjusts the influence of the LoRA matrices. During the fine-tuning process, the pre-trained weight 107 matrix  $W_0 \in R^{o \times i}$  is frozen, preserving the pre-trained features. For a given input activation 108  $X \in R^{i \times s \times b}$  (s: sequence length, b: batch size), the output  $Y \in R^{o \times s \times b}$  of a layer utilizing LoRA is computed as follows:

110 111

117

118

119 120

121

131 132 133

$$Y = (W_0 + \alpha BA)X = W_0X + \alpha BAX \tag{1}$$

The fine-tuning of the LoRA parameters is guided by the gradient of a loss function L, which is calculated with respect to each parameter matrix. The gradients are derived as follows:

$$\frac{\partial L}{\partial A} = \alpha \frac{\partial L}{\partial \tilde{X}} X^{\top}, \quad \frac{\partial L}{\partial B} = \alpha \frac{\partial L}{\partial Y} \tilde{X}^{\top}$$
(2)

Here,  $\tilde{X} := AX$  represents the intermediate input activation of B, which is the transformation of X by A. These gradients guide the adjustment of the LoRA parameters to minimize the loss and more accurately approximate the necessary updates to the original model weights.

2.2 QUANTIZATION

Uniform quantization is a widely used quantization scheme due to its simplicity and broad compatibility with various computing kernels and hardware units (Liu et al., 2022b). Therefore, considering the adaptability on the broad quantization scheme, including both weight-only and activation quantization scenarios, we refer to the term 'quantization' specifically as uniform quantization throughout this paper. A common practice is to organize a quantization group consisting of a certain number of consecutive weight elements that share the same quantization scale s and bias b, which control the quantization step size and zero point, respectively.

The weights W within the quantization group are quantized according to the following equation: 130

$$\tilde{w} = R(clamp(\frac{W-b}{s}, Q_n, Q_p)), \quad W_q = \tilde{w} \times s + b$$
(3)

Here,  $\tilde{w}$  denotes the quantized integer value obtained through a sequence of division, clamping, and rounding. Clamping is applied within the range  $Q_n = -2^{n-1}$  to  $Q_p = 2^{n-1} - 1$ , where *n* is the bit-width, followed by the rounding function *R*.  $W_q$  represents the dequantized version of the quantized weight, which is adjusted using *s* and *b* from  $\tilde{w}$  to approximate the original weight.

During QAT, the straight-through estimator (STE) approximates the derivative of the rounding func-138 tion with an identity function (Bengio et al., 2013; Choi et al., 2018; Esser et al., 2020), enabling 139 gradients to propagate through non-differentiable rounding operations and allowing effective weight 140 parameter training. Building on conventional QAT, LSQ (Esser et al., 2020) and LSQ+ (Bhal-141 gat et al., 2020) extend this process by training quantization parameters s and b, alongside model 142 weights. This tuning scheme provides finer control over quantization, improving model accuracy. 143 The quantization parameters tuning during backpropagation proceeds as follows, where w denotes the value  $\frac{W-b}{c}$ : 144 145

$$\frac{\partial L}{\partial s} = \frac{\partial W_q}{\partial s} \frac{\partial L}{\partial W_q} \quad s.t. \quad \frac{\partial W_q}{\partial s} = -w + \tilde{w}, \quad \text{if } Q_n \le w \le Q_p \tag{4}$$

147 148

146

149 150

151 152

153

$$\frac{\partial L}{\partial b} = \frac{\partial W_q}{\partial b} \frac{\partial L}{\partial W_q} \quad s.t. \quad \frac{\partial W_q}{\partial b} = 1, \qquad \text{if } w < Q_n \quad or \quad w > Q_p \tag{5}$$

More details on QAT with LSQ and LSQ+ are provided in Appendix A.1.

# 154 2.3 LLM QUANTIZATION

Quantization compress LLMs by lowering the bit precision of model parameters (Hubara et al., 2017). However, a key challenge with quantization is the introduction of errors that reduce model accuracy, leading to extensive research aimed at mitigating these losses through calibration or training. A notable examples of PTQ for LLM compression are GPTQ(Frantar et al., 2023) and OmniQuant (Shao et al., 2023). In contrast, QAT integrates quantization into the training process, adaptively training model parameters to account for quantization errors during training, ensuring that the quantized model retains much of its accuracy and functionality through training.



Figure 2: A categorization of training scheme and inference strategy of QA-LoRA, QLoRA, QAT-LoRA, and L4Q. Compared to QA-LoRA, L4Q utilizes higher optimization ability with non-constrained LoRA parameters and quantization parameters. Additionally, compared to QLoRA and QAT-LoRA, L4Q exploit fully-quantized weights rather than the mixed-precision weights during inference and perform a solid co-optimization of parameters.

Despite its advantages, QAT faces challenges, primarily due to its high training overhead, which 189 limits its use in resource-intensive models like LLMs. The memory overhead of QAT stems from 190 storing weight gradients and their optimizer states, each requiring multiple times the memory of the 191 weights. For example, applying QAT to a 7B model that fine-tunes all weight parameters requires at 192 least 14GB for the model weights, 14GB for the gradients, 28GB for optimizer states, and additional 193 dozens of GB for activations. This often exceeds the memory capacity of a single GPU, which 194 typically offers a maximum of 80GB of memory. Therefore, previous works have primarily focused 195 on the application of PTQ for LLM compression, while research on QAT remains in its early stages, 196 especially in terms of improving training efficiency.

197 198 199

### 2.4 QUANTIZATION-AWARE PEFT

200 To improve the accuracy of quantized LLMs, recent research has introduced quantization-aware 201 PEFT approaches (Dettmers et al., 2024b; Kim et al., 2024; Xu et al., 2023; Li et al., 2024). Among 202 these, QLoRA (Dettmers et al., 2024b), QA-LoRA (Xu et al., 2023), and LoftQ (Li et al., 2024) 203 stand out as notable methods. As illustrated in Figure 2, QLoRA begins by applying PTQ to a 204 pre-trained model. After this initial quantization, LoRA fine-tuning is performed, with the quan-205 tized weight parameters kept frozen. The LoRA parameters are fine-tuned using higher precision 206 formats such as bfloat16 or float16, allowing the method to correct quantization errors during the 207 fine-tuning. However, QLoRA introduces computational inefficiencies during inference due to the additional forward path on LoRA parameters. This inefficiency arises because the high-precision 208 LoRA parameters and low-precision quantized weights cannot be merged into low-precision values. 209 We further examine the impact of this unmerged LoRA path on inference efficiency by comparing 210 the speed of fully-quantized models with mixed-precision models in Section 4. 211

QA-LoRA (Xu et al., 2023) addresses the issue of high-precision LoRA parameters by modifying
 the structure of the LoRA matrix, allowing these parameters to be integrated with the quantized
 weights after training (Figure 2). The input dimension of the LoRA matrix A is set to the number of
 weight quantization groups. This adjustment ensures that each element of *BA* corresponds directly
 with individual quantization groups, enabling the LoRA parameters to be seamlessly integrated into

the quantization bias as  $b' = b - \alpha BA$  at the end of training. However, this solution presents a new challenge: the constrained LoRA structure in this setup limits the model's ability to achieve optimal accuracy during the PEFT stage.

219 A broader issue with previous quantization-aware PEFT approaches is that quantization and fine-220 tuning are performed sequentially rather than simultaneously. Starting the fine-tuning process from a 221 pre-quantized model with inherent quantization errors is less optimal than starting from a pre-trained 222 model. While LoftQ attempts to address these quantization errors by approximating them using 223 LoRA with iterative singular-value decomposition (SVD), the limitations on subsequent adaptation 224 persist, making it challenging to fully recover and optimize model performance during fine-tuning. 225 These challenges highlight the need for continued research to improve quantization-aware PEFT 226 techniques, aiming to enhance both the quantization and PEFT processes for better accuracy and efficiency in LLMs. 227

228 229

# 3 Methods

230 231 232

233

# 3.1 STRAIGHTFORWARD INTEGRATION OF QAT AND LORA

**QAT-LoRA** One of the key principles in our proposed L4Q scheme is the integration of the QAT 234 and LoRA to facilitate the simultaneous calibration for quantization and fine-tuning on downstream 235 tasks. To achieve this, we begin with a straightforward integration of QAT and LoRA, referred to as 236 QAT-LoRA, which serves as our baseline approach for combining QAT and PEFT. In QAT-LoRA, 237 pre-trained weights are frozen, while LoRA parameters are added to the linear layers (Figure 2). 238 Additionally, quantization scales and bias parameters are introduced, similar to advanced QAT tech-239 niques like LSQ, which are crucial for calibrating the quantization function. Freezing the weights 240 reduces the need for optimizer states, while a small number of LoRA and quantization parameters 241 are introduced to approximate updates to the weight matrix and to update the quantization func-242 tion, respectively. This results in more efficient memory usage compared to standard QAT. Detailed 243 analysis results of the memory efficiency of QAT-LoRA is further discussed in Section 4.

244 **Limitations of QAT-LoRA** However, several issues arise with this straightforward integration of 245 QAT and LoRA, where quantized weights and LoRA parameters remain distinct. First, although 246 freezing the pre-trained weights eliminates the need for optimizer states, weight gradients  $\frac{\partial L}{\partial W_{-}}$  must 247 still be stored to update the quantization parameters, as shown in Equations 4-5. As a result, QAT-248 LoRA still incurs memory overhead from weight gradients, undermining the memory efficiency 249 benefits of LoRA fine-tuning. Secondly, QAT-LoRA produces a mixed-precision model with both 250 quantized weights and high-precision LoRA parameters. This mixed-precision approach negates the 251 advantages of LLM quantization, similar to previous methods such as QLoRA and LoftQ discussed 252 in Section 2.4. Lastly, the gradient updates for quantization and LoRA parameters are decoupled, limiting the potential for comprehensive optimization across the model. As outlined in Equations 253 4-5, updates to the quantization parameters rely on the quantized weight matrix  $W_q$ , while updates 254 to the LoRA parameters depend on weights A and B. This limits the effectiveness of model train-255 ing, as it prevents holistic adjustments where changes in LoRA parameters could directly influence 256 quantization adjustments and vice versa. 257

To address these challenges, we introduce L4Q, an enhanced integration of QAT and LoRA. L4Q features an advanced layer design that seamlessly integrates QAT and LoRA. By applying quantization after merging the weights and LoRA parameters, along with a custom backpropagation path that reduces the memory overhead from the complex quantization and LoRA processes, L4Q effectively overcomes the primary challenges encountered with QAT-LoRA.

263 264

265

# 3.2 L4Q: LOW-RANK ADAPTIVE LEARNING QUANTIZATION

Fully-Quantized Linear Layer As high-precision LoRA weights introduces inference overhead, it is crucial to design a fully-quantized linear layer. In this context, L4Q first combines the original weights  $W_0$  and the LoRA parameters BA into a unified parameter matrix:

$$W_{comb} = W_0 + \alpha B A \tag{6}$$

273

277

283

284

285

287

307 308

270 Then, quantization is applied to the fully combined weight  $W_{comb}$  as follows: 271

$$W_q = R(clamp(\frac{W_{comb} - b}{s}, Q_n, Q_p)) \times s + b$$
(7)

274 In this way, during inference, L4Q only uses quantized weights, simplifying the forward path of the 275 linear layer as follows: 276

$$Y = W_q X \tag{8}$$

While QA-LoRA also achieves fully-quantized linear layers by introducing constraints on the LoRA 278 parameter structure, the proposed L4Q imposes no such restrictions. This flexibility allows L4Q to 279 fully leverage the benefits of LoRA-based fine-tuning, all while achieving fully-quantized linear 280 layers. 281

**Memory Efficient QAT** As discussed in the previous section, QAT requires weight gradients for 282 training the quantization parameters s and b. Since weight gradients are a major source of memory overhead during training, we compute these gradients locally in the backpropagation path of the linear layer, utilizing the input and output of the layer. The weight gradients are calculated as follows:

$$\frac{\partial L}{\partial W_a} = \frac{\partial L}{\partial Y} X^\top \tag{9}$$

We use these weight gradients to calculate gradients of s and b with Equation 4 and 5. Once the 289 gradient computation for the linear layer is complete, the weight gradients are immediately flushed 290 to conserve memory. 291

292 **Efficient LoRA Training** To compute the gradients of the LoRA parameters in the L4Q linear 293 layer, we must trace back from Equation 8 to 6. Due to the full integration of quantization and LoRA parameters (Equation 6 and 7), the gradients of the LoRA parameters also rely on the weight gradients. This then be described as follows: 295

$$\frac{\partial L}{\partial A} = \frac{\partial L}{\partial W_a} \frac{\partial W_q}{\partial A}, \quad \frac{\partial L}{\partial B} = \frac{\partial L}{\partial W_a} \frac{\partial W_q}{\partial B}$$
(10)

Since we reuse the weight gradient  $\frac{\partial L}{\partial W_a}$  that have already been computed for QAT, we only need 300 to compute  $\frac{\partial W_q}{\partial A}$  and  $\frac{\partial W_q}{\partial B}$  to obtain the gradients of LoRA parameters. Both terms are derived by applying the chain rule to Equation 7 and Equation 6. Since Equation 7 contains a rounding function, 301 302 we apply STE and clamping with conditional gradient propagation to  $\frac{\partial W_q}{\partial A}$  and  $\frac{\partial W_q}{\partial B}$ . Moreover, in Equation 6, the LoRA parameters A and B are simply multiplied, so its gradients are expressed directly as  $\alpha B^{\top}$  and  $\alpha A^{\top}$ , respectively. This leads to the following expressions for  $\frac{\partial W_q}{\partial A}$  and  $\frac{\partial W_q}{\partial B}$ : 303 304 305 306

$$\frac{\partial W_q}{\partial A} = \begin{cases} \alpha B^\top, & \text{if } Q_n \le w \le Q_p \\ 0, & otherwise \end{cases}, \quad \frac{\partial W_q}{\partial B} = \begin{cases} \alpha A^\top, & \text{if } Q_n \le w \le Q_p \\ 0, & otherwise \end{cases}$$
(11)

310 Therefore, the proposed L4O efficiently processes LoRA training by simply reusing the weight 311 gradients computed for QAT parameter training. For more detailed explanations of the gradient 312 calculation in L4Q, please refer to Appendix A.2, and the memory efficiency of L4Q will be further 313 examined in Section 4. 314

Joint Optimization of Quantization and LoRA parameters Since  $\frac{\partial L}{\partial W_q}$  is involved in the gradi-315 ent calculation for the LoRA parameters (Equation 10), the proposed L4Q ensures that the impact 316 of quantization is directly reflected in the updates to the LoRA parameters. This enables the joint 317 optimization of LoRA parameters and the quantization process, enhancing the accuracy of the fully-318 quantized LLMs. 319

320 In summary, the proposed L4Q produces a fully-quantized model for memory-efficient and fast LLM 321 inference by fully integrating the model weights and LoRA parameters prior to the quantization process. Additionally, the training process of L4Q is memory-efficient due to careful handling of 322 gradient computation for quantization. Finally, L4Q can improve the accuracy of quantized LLMs 323 through the joint optimization of the quantization and LoRA parameters.



Figure 3: The MMLU 5-shot benchmark results for LLaMA-2 7B models with 100-step training, along with the sum of quantization and clipping errors, are presented for the initialization methods: LSQ+, symmetric, asymmetric, and L4Q. Errors are reported at both the initial point (Initial) and after training (Post-train).

## 3.3 QUANTIZATION PARAMETER INITIALIZATION

341 In L4Q, we address the outlier-sensitive nature of LLMs by improving the quantization parameter 342 initialization at the start of fine-tuning. While the initialization scheme in LSQ+ (LSQ+ $_{init}$ ) has 343 been shown to be effective by using multiples of standard deviation of weights as the quantiza-344 tion range and setting the quantization scale accordingly, this approach was originally designed for 345 CNNs, not transformer-based LLMs. Studies on LLM quantization have shown that outliers in acti-346 vation and their corresponding salient weights significantly impact model performance (Xiao et al., 347 2023; Dettmers et al., 2024a; Lin et al., 2024). In LLM, the standard deviation of the small number of weight groups is often much smaller than the absolute maximum value, causing large weights to 348 be more prone to clipping, resulting in significant clipping errors. To address this, we introduce a 349 quantization initialization scheme, L4Q<sub>init</sub>, which minimizes clipping errors by adopting a conser-350 vative quantization scale that captures both minimum and maximum outliers during training. The 351 detailed equation is as follows: 352

$$s = Max(\left|\frac{Min(W)}{Qn}\right|, \left|\frac{Max(W)}{Qp}\right|), \quad b = 0$$
(1)

2)

Note that conventional symmetric initialization involves division of the absolute maximum value of
 the numbers, and asymmetric initialization involves division of the minimum to maximum range
 of the number for the quantization scale initialization. Detailed explanation on LSQ+, symmetric,
 asymmetric, and L4Q quantization initialization is presented in Appendix B.

We evaluate the models trained with L4Q using various initialization methods, comparing the accu-360 racy, the sum of quantization error  $|W - W_q|$ , and the sum of clipping error of the overflowed outliers 361 at both the initialization point and the end of training. As shown in Figure 3,  $L4Q_{init}$  achieves the 362 highest accuracy, while LSQ+*init* struggles to recover from quantization errors. Although the overall quantization error remains similar across methods, clipping errors are notably higher in  $LSQ+_{init}$ 364 and symmetric initialization. Additionally, while asymmetric initialization avoids outliers at the 365 start, it fails to handle emerging outliers during fine-tuning as the weight distribution shifts. Its tight 366 initialization, focused on the minimum to maximum range, becomes insufficient as weights change, 367 leading to higher clipping error. In contrast, L4Q<sub>init</sub> considers the broader weight distribution in 368 LLMs, reducing both quantization and clipping errors effectively.

369 370 371

334

335

336

337 338 339

340

353

354 355

### 3.4 OVERALL L4Q ALGORITHM

The overall L4Q algorithm is outlined in Algorithm 1. In the initial phase, L4Q undergoes a warmup stage with LoRA fine-tuning only. Since the quantization function includes a clamping mechanism, applying it to LoRA at the beginning of training can suppress gradient updates for outlier values, which negatively affects LLM accuracy. Hence, L4Q warms up the LoRA parameters before applying quantization. The LoRA warm-up consists of  $T_{wp}$  iterations, where  $T_{wp} \ll T$ , with T representing the total number of training steps. An analysis of the effect of LoRA warm-up can be found in Appendix C. After the warm-up phase, L4Q combines the pre-trained weights and

Algorithm 1 L4Q Algorithm	
<b>Require:</b> Pre-trained weight $W_0$ , Fine-tuned weight $W_{co}$	$_{mb}$ , Quantized weight $W_q$
<b>Require:</b> LoRA parameter $A, B$ , LoRA configuration $r, d$	α
<b>Require:</b> Quantization parameter $s, b$ , Quantization bit-w	vidth n
<b>Require:</b> Training step T, LoRA warm-up step $T_{wp}$ , Stat	e I <sub>init</sub>
Initialize $A \sim N(0, \sigma^2), B \leftarrow 0, \alpha \leftarrow \frac{1}{2\pi}, I_{init} \leftarrow Fallet$	se
for $t = 1$ to $T$ do	
if $t \leq T_{wp}$ then	⊳ LoRA warm-up stage
Forward $W_0X + \alpha BAX$ and Update $(A, B)$	
else	▷ L4Q fine-tuning stage
if not $I_{init}$ then	
Set $W_{comb} \leftarrow W_0 + \alpha BA$	
Initialize $s, b \leftarrow L4Q_{init}$	$\triangleright$ L4Q <sub><i>init</i></sub> described in the equation 12
Set $I_{init} \leftarrow$ True	
end if	
Forward $W_q X$ and Update $(A, B, s, b)$	$\triangleright W_q$ described in the equation 8
end if	1
end for	

LoRA parameters (Equation 6), and then initializes the quantization parameters *s* and *b*, which determine the quantization range, using these combined weights. While previous QAT schemes use the standard deviation of weights to set the initial quantization range, and prior LLM quantization approaches rely on the min/max of the weights, L4Q employs a relaxed min/max-based initialization as discussed in Section 3.3. Then, the proposed L4Q conducts both QAT and LoRA, with gradient computation as described in Section 3.2, to achieve high-accuracy fully-quantized models.

402 403 404

405

396 397

398

399

400

401

4 EXPERIMENTS

406 4.1 EXPERIMENTAL SETTINGS

Target Foundation Models OpenLLaMA<sup>1</sup>, LLaMA-1 (Touvron et al., 2023a), and LLaMA-2 (Touvron et al., 2023b) model families are used for the evaluation. Specifically, we assess the OpenLLaMA model with 3B parameters, LLaMA-1 models with 7B, 13B, and 33B parameters, and LLaMA-2 models with 7B and 13B parameters. When applying quantization, the quantization group size is set to 128 for the LLaMA-1 and LLaMA-2 models, and 64 for the OpenLLaMA models.

Baselines We compare the proposed L4Q with previous PTQ and QA-PEFT methods. The baseline PTQ methods include GPTQ and OmniQuant, while the baseline QA-PEFT methods include
QLoRA, QA-LoRA, and LoftQ. Although the original QLoRA and LoftQ employ non-uniform
quantization, we apply uniform quantization in our experiments for consistency across methods.
These uniformly quantized versions are referred to as QLoRA\* and LoftQ\*.

418 **Evaluation Setups** We establish the L4Q framework based on the Lit-GPT<sup>2</sup> and huggingface trans-419 formers<sup>3</sup>, both of which are open-source LLM frameworks. The rank of LoRA parameter r is set to 420 4 by default, as our experiments showed that r = 4 is sufficient to maintain performance. Further details on the effect of rank size can be found in Appendix D. Meanwhile, as QA-LoRA reduces 421 the input dimension of the LoRA parameters to integrate them with the quantization bias parameter, 422 we double the rank r to 8 for a fair comparison in terms of the number of LoRA parameters. For 423 the fine-tuning process, we use the Stanford-Alpaca (Taori et al., 2023) dataset that consists of 52k 424 instruction-following samples generated from the instruction-tuned GPT 3.5 model (Brown et al., 425 2020). This dataset is split into a training set with 50k samples and a validation set with 2k samples. 426 In line with standard LLM training conventions, we used bfloat16 for numerically stable fine-tuning. 427 The batch size for fine-tuning is 128. All experiments were conducted on an NVIDIA A100 80GB 428 GPU. Detailed hyperparameter settings for fine-tuning are provided in Appendix E.

```
<sup>2</sup>https://github.com/Lightning-AI/lit-gpt.git
```

```
<sup>3</sup>https://github.com/huggingface/transformers.git
```

<sup>&</sup>lt;sup>1</sup>https://github.com/openlm-research/open\_llama

	OpenLLaMA		LLaMA-1	
Methods	3B	7B	13B	33B
LoRA	15.1	25.1	43.8	71.9
QLoRA	5.2	7.9	19.6	31.9
LoftQ	5.2	7.9	19.6	31.9
QA-LoRA	7.8	14.8	27.8	67.2
QAT	44.2	79.5	OOM	OOM
QAT-LoRA	22.6	41.9	70.6	OOM
140	15.3	25.4	44.3	73.2





Figure 4: The average inference speedup of quantized models compared to pre-trained models.

**Evaluation Metrics** We evaluate the accuracy of LLMs on Commonsense QA (CSQA) (Gao et al., 2021) and MMLU (Hendrycks et al., 2021) benchmarks. The CSQA benchmark includes tasks from Hellaswag (Zellers et al., 2019), PIQA (Bisk et al., 2020), ARC-challenge and ARC-448 easy (Clark et al., 2018), Winogrande (Sakaguchi et al., 2020), BoolQ (Clark et al., 2019), and 449 Openbook QA (Talmor et al., 2019). The MMLU benchmark spans four subject categories: Hu-450 manities, STEM, Social Sciences, and others made up of 57 subcategories of language tasks.

# 4.2 EVALUATION RESULTS

445

446

447

451 452

453

454 Memory Cost for Fine-Tuning We measure the peak memory usage during fine-tuning of 4-bit 455 LLMs using both the baseline methods and L4Q, as reported in Table 1. Since PTQ-based PEFT 456 methods apply PTQ before fine-tuning, the size of linear layers in the model is reduced by a factor 457 of 1/4, resulting in lower memory usage compared to naïve LoRA, which fine-tunes full-precision 458 models. While both QAT and QAT-LoRA result in 2-3 times higher memory costs compared to 459 LoRA, the memory usage of L4Q is similar to that of LoRA. Although L4Q cannot reduce memory 460 costs as much as PTQ-based PEFT methods due to the use of 16-bit full-precision models during fine-tuning, this result highlights that L4Q is well-designed to leverage the memory-efficient nature 461 of LoRA in training. Further analysis on the training efficiency of L4Q and baseline methods can be 462 found in Appendix F. 463

464 We measure the inference speed of 16-bit pre-trained models and quantized Inference Speedup 465 models using LLaMA-1 models. The quantized models include fully-quantized 4-bit models (L4Q and QA-LoRA), which contain only quantized parameters, and mixed-precision 4&16-bit models 466 (LoftQ\*, QLoRA\*, and QAT-LoRA), which use additional 16-bit LoRA parameters. The inference 467 speed was measured with input batch sizes ranging from 1 to 64. The average speedup of quantized 468 models compared to full-precision 16-bit models is reported in Figure 4. The 4-bit models achieve 469 a speedup of  $1.8 \times$  to  $2.3 \times$  over the 16-bit models. More importantly, these 4-bit models achieve 470 a  $1.4 \times$  to  $1.6 \times$  speedup compared to mixed-precision 4&16-bit models, which are also quantized 471 versions of LLMs. This demonstrates that the full integration of QAT and LoRA in L4Q plays 472 a crucial role in achieving the best inference speed. More analysis on speedup can be found in 473 Appendix G. 474

Accuracy Results We compare the accuracy of baselines and L4Q. Specifically, we evaluate zero-475 shot accuracy on the CSQA benchmark, and evaluate both zero-shot and few-shot (5-shot) accuracy 476 on the MMLU benchmark. Table 2 presents a comprehensive comparison between baselines and 477 the proposed L4Q after 4-bit LLM quantization. Since previous QA-PEFT methods involve a fine-478 tuning stage after PTQ, they generally achieve higher accuracy compared to PTQ methods. L4Q 479 further improves accuracy of 4-bit models by adopting QAT strategy, and it achieves accuracy com-480 parable to 16-bit models. Moreover, L4Q consistently outperforms QAT-LoRA, which adopts QAT 481 but keeps quantization and LoRA parameters decoupled. This highlights the advantage of L4Q in 482 achieving superior accuracy through the joint optimization of quantization and LoRA parameters. 483 We also evaluate accuracy with 3-bit LLM quantization. As presented in Table 3. L4Q achieves the best accuracy in this case as well. Notably, the impact of applying QAT on accuracy is more pro-484 nounced in 3-bit quantization, as previous approaches experience significant accuracy degradation 485 with 3-bit quantization.

		Pre- trained	LoRA	GPTQ	OmniQ.	LoftQ*	QLoRA*	QA- LoRA	QAT- LoRA
Model	Benchmark	16	16	4	4	4+16	4+16	4	4+16
OpenLLaMA 3B	CSQA	54.8	55.9	50.7	54.1	54.2	54.4	54.5	54.6
LLaMA-1 7B	CSQA	61.7	63.4	59.4	58.1	62.6	61.3	61.3	62.4
	MMLU <sub>0-shot</sub>	32.5	36.3	28.3	30.9	33.0	32.8	34.5	33.8
	MMLU5-shot	35.1	36.7	32.7	33.3	35.1	33.6	35.6	34.8
LLaMA-1 13B	CSQA	63.8	65.2	63.5	60.4	64.2	63.8	62.5	64.4
	MMLU <sub>0-shot</sub>	43.6	44.3	40.1	42.6	42.4	42.1	42.4	42.0
	MMLU5-shot	46.3	47.0	45.7	45.7	45.4	45.9	45.8	45.5
LLaMA-1 33B	CSQA	67.4	68.3	65.7	62.9	67.4	66.2	65.3	67.3
	MMLU <sub>0-shot</sub>	53.0	54.4	51.4	52.0	51.8	51.0	48.9	52.3
	MMLU <sub>5-shot</sub>	56.4	57.6	55.7	55.8	56.4	55.6	55.0	56.7
LLaMA-2 7B	CSQA	61.9	63.3	60.7	59.5	61.7	61.3	61.0	61.9
	MMLU <sub>0-shot</sub>	41.6	43.9	37.1	41.0	38.5	38.6	38.9	37.9
	MMLU <sub>5-shot</sub>	45.4	46.0	42.9	45.4	43.7	44.6	44.4	43.8
LLaMA-2 13B	CSQA	65.0	66.5	64.4	59.9	64.9	64.0	64.5	64.7
	MMLU <sub>0-shot</sub>	52.1	52.5	50.0	51.8	51.7	50.7	50.4	50.7
	MMLU <sub>5-shot</sub>	54.8	55.7	54.7	54.7	54.5	54.2	54.1	53.8

Table 2: Accuracy (%) evaluation results with 4-bit quantization. The bit precision of weight parameters is indicated under the method name. The notation '4+16' refers to the requirement of 16-bit LoRA parameters alongside 4-bit weights for inference.

Table 3: Accuracy (%) evaluation results with 3-bit quantization. The bit precision of weight parameters is indicated under the method name. The notation '3+16' refers to the requirement of 16-bit LoRA parameters alongside 3-bit weights for inference.

		Pre- trained	LoRA	GPTQ	OmniQ	LoftQ*	QLoRA*	QA- LoRA	QAT- LoRA	L4Q
Model	Benchmark	16	16	3	3	3+16	3+16	3	3+16	3
OpenLLaMA 3B	CSQA	54.8	55.9	52.2	50.0	38.1	51.0	51.5	53.2	54.0
LLaMA-1 7B	CSQA	61.7	63.4	53.4	56.5	49.8	59.1	58.7	60.7	61.2
	MMLU <sub>0-shot</sub>	32.5	36.3	23.7	29.0	23.4	27.7	28.0	30.6	30.6
	MMLU5-shot	35.1	36.7	27.3	31.6	23.1	31.5	29.1	31.5	31.8
LLaMA-1 13B	CSQA	63.8	65.2	61.0	58.9	54.0	61.3	61.1	63.2	63.4
	MMLU <sub>0-shot</sub>	43.6	44.3	33.1	34.8	25.0	36.1	37.5	38.8	40.7
	MMLU5-shot	46.3	47.0	38.2	41.6	25.3	40.4	38.2	40.9	41.8
LLaMA-2 7B	CSQA	61.9	63.3	57.6	57.9	34.7	57.6	56.3	57.4	61.3
	MMLU <sub>0-shot</sub>	41.6	43.9	31.3	34.3	22.9	32.5	31.0	31.5	34.9
	MMLU5-shot	45.4	46.0	37.5	37.7	24.2	37.6	37.5	36.8	38.0
LLaMA-2 13B	CSQA	65.0	66.5	61.7	59.9	39.3	62.5	61.7	64.3	65.1
	MMLU <sub>0-shot</sub>	52.1	52.5	46.3	46.3	23.5	46.8	46.4	45.9	47.1
	MMLU5-shot	54.8	55.7	50.4	50.2	26.0	50.6	49.9	48.9	50.0

# 5 CONCLUSION

In this work, we introduce L4Q, a parameter-efficient quantization-aware fine-tuning method for large language models. L4Q enables element-wise adaptation of model weights for downstream tasks while simultaneously optimizing quantization parameters. This concurrent optimization en-sures that the adaptation parameters effectively account for quantization errors. We demonstrate the efficiency of L4Q, which significantly reduces training resource requirements compared to traditional QAT. Moreover, since the L4Q layer is designed to produce fully quantized low-bit model weights, it maintains inference efficiency, unlike QLoRA, LoftQ, or QAT-LoRA, which result in mixed-precision models. The effectiveness of L4Q as a QAT framework is further supported by experimental results across various task evaluations. L4Q consistently achieves superior quality maintenance on language tasks, demonstrating its enhanced adaptability compared to QAT-LoRA and PTQ-based PEFT methods. 

# 540 ETHICS STATEMENTS

541 542 543

544

546 547

548

556

558

This paper presents work aimed at advancing the field of machine learning, particularly in the context of language models. While our research contributes to the development and exploitation of LLMs, we acknowledge that there are potential societal consequences associated with this work. However, we do not delve into the potential harms or biases inherent in these models.

# Reproducibility Statement

This paper introduces a novel algorithm for quantization layer design and the training process. To ensure reproducibility, the source code is provided as an anonymous, downloadable link in the supplementary materials, along with detailed explanations of the environmental setup and execution process. The theoretical results are thoroughly explained with clear assumptions, and complete proofs are included in the appendix. Additionally, detailed descriptions of the data processing steps for all datasets used in the experiments are available in the supplementary materials. These resources are intended to fully enable the reproduction of our results.

- References
- 559 Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients 560 through stochastic neurons for conditional computation, 2013.
- Yash Bhalgat, Jinwon Lee, Markus Nagel, Tijmen Blankevoort, and Nojun Kwak. Lsq+: Improving
   low-bit quantization through learnable offsets and better initialization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2020.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about
   physical commonsense in natural language. In *The Association for the Advancement of Artificial Intelligence (AAAI)*, 2020.
- 568 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, 569 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. 570 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, 571 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, 572 Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Proceedings of the 573 34th International Conference on Neural Information Processing Systems (NeurIPS), Vancouver, 574 BC, Canada, 2020. Curran Associates Inc. ISBN 9781713829546. 575
- Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srini vasan, and Kailash Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural
   networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- 579 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam 580 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, 581 Kensen Shi, Sashank Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, 582 Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, 583 James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, 584 Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Bar-585 ret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, 586 Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Bren-588 nan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas 589 Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. 590 Journal of Machine Learning Research (JMLR), 24(1):240:1–240:113, March 2024. ISSN 1532-4435. 592
- 593 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Annual*

596

607

619

629

636

637

638

639

640

Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), 2019.

- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
   Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge,
   2018.
- Tim Dettmers and Luke Zettlemoyer. The case for 4-bit precision: k-bit inference scaling laws. In
   *Proceedings of the 40th International Conference on Machine Learning (ICML)*, 2023.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2024a. Curran Associates Inc. ISBN 9781713871088.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning
   of quantized llms. In *Proceedings of the 37th International Conference on Neural Information Processing Systems (NeurIPS)*, New Orleans, LA, USA, 2024b.
- Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S. Modha. Learned step size quantization. In *International Conference on Learning Representations (ICLR)*, 2020.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training
   quantization for generative pre-trained transformers. *International Conference on Learning Representations (ICLR)*, abs/2210.17323, 2023. URL https://api.semanticscholar.
   org/CorpusID:253237200.
- Leo Gao, Jonathan Tow, Stella Biderman, Shawn Black, Anthony DiPofi, Charlie Foster, Leo Golding, Jasmine Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Lucy Reynolds, Eric Tang, Alex Thite, Ben Wang, Kevin Wang, and Amanda Zou. A framework for few-shot language model evaluation, 2021. URL https://doi.org/10.5281/zenodo.5371629. Zenodo Repository.
- Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. arXiv: Computer Vision and Pattern Recognition, 2015. URL https://api.semanticscholar.org/CorpusID: 2134321.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
   Steinhardt. Measuring massive multitask language understanding. In *International Conference* on Learning Representations (ICLR), 2021.
- Jung Hwan Heo, Jeonghoon Kim, Beomseok Kwon, Byeongwook Kim, Se Jung Kwon, and Dong soo Lee. Rethinking channel dimensions to isolate outliers for low-bit weight quantization of
   large language models. In *International Conference on Learning Representations (ICLR)*, 2024.
  - J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *International Conference on Learning Representations (ICLR)*, abs/2106.09685, 2022. URL https://api. semanticscholar.org/CorpusID:235458009.
- Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *Journal* of Machine Learning Research (JMLR), 18(1):6869–6898, January 2017.
- Jeonghoon Kim, Jung Hyun Lee, Sungdong Kim, Joonsuk Park, Kang Min Yoo, Se Jung Kwon, and
   Dongsoo Lee. Memory-efficient fine-tuning of compressed large language models via sub-4-bit
   integer quantization. In *Proceedings of the 37th International Conference on Neural Information Processing Systems (NeurIPS)*, New Orleans, LA, USA, 2024.

- Kiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL) and the 11th International Joint Conference on Natural Language Processing (IJCNLP), pp. 4582–4597, Online, 2021. Association for Computational Linguistics. doi: 10.18653/v1/ 2021.acl-long.353. URL https://aclanthology.org/2021.acl-long.353.
- 4653
   4654
   4655
   4656
   471
   472
   473
   474
   474
   474
   474
   474
   474
   474
   474
   474
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475
   475</l
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration. In *Annual Conference on Machine Learning and Systems (MLSys)*, 2024.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and
   Colin A. Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
   learning. In Advances in Neural Information Processing Systems (NeurIPS), volume 35, pp. 1950–
   1965, 2022a. URL https://proceedings.neurips.cc/paper\_files/paper/
   2022/file/0cde695b83bd186c1fd456302888454c-Paper-Conference.pdf.
- Z. Liu, K. Cheng, D. Huang, E. Xing, and Z. Shen. Nonuniform-to-uniform quantization: Towards accurate quantization via generalized straight-through estimation. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4932–4942, Los Alamitos, CA, USA, jun 2022b. IEEE Computer Society. doi: 10.1109/CVPR52688.2022.00489. URL https://doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.00489.
- Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. LLM-QAT: Data-free quantization aware training for large language models. In *Findings of the Association for Computational Linguistics (ACL 2024)*, pp. 467–484, Bangkok, Thailand and Virtual Meeting, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.26. URL https://aclanthology.org/2024.findings-acl.26.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. In *The Association for the Advancement of Artificial Intelligence (AAAI)*, 2020.
- Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. Fine-tuned language models are
   continual learners. In *Proceedings of the 2022 Conference on Empirical Methods in Natu- ral Language Processing (EMNLP)*, pp. 6107–6122, Abu Dhabi, United Arab Emirates, 2022.
   Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.410. URL
   https://aclanthology.org/2022.emnlp-main.410.

689

- Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large language models. arXiv preprint arXiv:2308.13137, 2023.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In *North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xinyun Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model, 2023. URL https://github.com/tatsu-lab/stanford\_alpaca. GitHub Repository.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
   Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
   language models, 2023a. URL https://arxiv.org/abs/2302.13971.

- 702 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-703 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, 704 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy 705 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, 706 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, 708 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, 709 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh 710 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen 711 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, 712 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 713 2023b. URL https://arxiv.org/abs/2307.09288. 714
- Zhen Wang, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, Huan Sun, and Yoon Kim. Multitask prompt tuning enables parameter-efficient transfer learning, 2023. URL https://arxiv. org/abs/2303.02861.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. *International Conference on Learning Representations (ICLR)*, abs/2109.01652, 2022. URL https://api. semanticscholar.org/CorpusID:237416585.
  - Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *Proceedings of the* 40th International Conference on Machine Learning (ICML), 2023.
- Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models, 2023. URL https://arxiv.org/abs/2309.14717.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer,
  Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models, 2022. URL https://arxiv.
  org/abs/2205.01068.
- 740 741 742 743 744 745 746 747 748
- 749

722

723

724

- 750
- 751
- 752 753
- 754
- 755

### DETAILS ON QUANTIZATION SCALE LEARNING PROCEDURE А

### QUANTIZATION SCALE UPDATE ON QAT A.1

. .

. .

From the conditions and notations in Equation 3, Equations 4 and 5 are derived as follows. First, the derivative of s is presented as follows.

$$\frac{\partial W_q}{\partial s} = \frac{\partial}{\partial s} (\tilde{w} \times s + b) = s \frac{\partial}{\partial s} (\tilde{w}) + \tilde{w} = s \frac{\partial}{\partial s} (R \cdot clamp(w)) + \tilde{w}$$
(13)

By applying the STE, the rounding function R is considered as an identity function. Therefore, the rounding function, combined with a clamp function  $R := R \cdot clamp$ , and its derivative is induced as follows. Note that  $w = \frac{W-b}{s}$ . 

$$\tilde{R}(w) = \begin{cases} Qn, & if \ w < Qn \\ w, & if \ Q_n \le w \le Q_p \\ Qp, & if \ w > Q_p \end{cases} \quad \frac{\partial}{\partial w} \tilde{R}(w) = \begin{cases} 1, & if \ Q_n \le w \le Q_p \\ 0, & otherwise \end{cases}$$
(14)

By applying the chain rule, the derivation of term  $R \cdot clamp(w) = \tilde{R}((W-b)/s)$  is expressed as below.

  $\frac{\partial}{\partial s}(\tilde{R}(w)) = \frac{\partial \tilde{R}}{\partial w} \frac{\partial w}{\partial s} = \frac{\partial \tilde{R}}{\partial w} \frac{\partial}{\partial s} (\frac{W-b}{s}) = \frac{\partial \tilde{R}}{\partial w} (-\frac{W-b}{s^2})$ 

Therefore, Equation 13 can be represented with a value w and quantized value  $\tilde{w}$  as follows.

$$\frac{\partial W_q}{\partial s} = s \frac{\partial \tilde{R}}{\partial w} \left(-\frac{W-b}{s^2}\right) + \tilde{w} = \frac{\partial \tilde{R}}{\partial w} \left(-\frac{W-b}{s}\right) + \tilde{w} = \begin{cases} Qn, & \text{if } w < Qn \\ -w + \tilde{w}, & \text{if } Q_n \le w \le Q_p \\ Qp, & \text{if } w > Q_p \end{cases}$$
(16)

Secondly, with a similar context above, the derivative of b is presented as follows.

$$\frac{\partial W_q}{\partial b} = \frac{\partial}{\partial b} (\tilde{w} \times s + b) = s \frac{\partial}{\partial b} (\tilde{R}(w)) + 1 = s \frac{\partial \tilde{R}}{\partial w} (\frac{\partial}{\partial b} (\frac{W - b}{s})) + 1$$
(17)

$$=\frac{\partial \tilde{R}}{\partial w}(-1)+1 = \begin{cases} 0, & \text{if } Q_n \le w \le Q_p \\ 1, & \text{otherwise} \end{cases}$$
(18)

We also note that the gradient of  $W_q$  is presented as follows.

$$\frac{\partial L}{\partial W_a} = \frac{\partial L}{\partial Y} X^\top \tag{19}$$

(15)

As a result, the updates on the quantization scale and bias are calculated as multiplication of Equation 19 with Equation 13 and with Equation 17, respectively. This update helps calibrate the quantization function, effectively reducing quantization errors.

824 825 826

837 838

# A.2 QUANTIZATION SCALE AND LORA PARAMETER UPDATE ON L4Q

In L4Q, as described in Equation 8, the quantized weight  $W_q$  is obtained as follows. First, the pretrained model weight  $W_0$  and LoRA parameters are integrated to  $W_{comb} = W_0 + \alpha BA$ . Next, the integrated weight is quantized by the quantization parameters s, b.

Here, the LoRA parameters A and B are independent of the quantization parameters, scale s and bias b. Therefore, the derivatives of s, b follow the same process as in Equation A.1, but with the term  $w, \tilde{w}$  defined as follows. Note that  $W_q = \tilde{w} \times s + b$ .

$$w = \frac{W_0 + \alpha BA - b}{s}, \quad \tilde{w} = \tilde{R}(w) \quad s.t. \quad \tilde{R} = R \cdot clamp \tag{20}$$

Seen from the L4Q layer that integrates the LoRA parameters and quantization parameters together, A, B are now considered as variables of  $W_q$ . Therefore, from the conditions in Equation 8, the derivative of A, B is presented as follows.

$$\frac{\partial L}{\partial A} = \frac{\partial W_q}{\partial A} \frac{\partial L}{\partial W_q}, \quad \frac{\partial L}{\partial B} = \frac{\partial L}{\partial W_q} \frac{\partial W_q}{\partial B}$$
(21)

The derivatives  $\frac{\partial w}{\partial A}$  and  $\frac{\partial w}{\partial B}$  are then can be computed by applying the chain rule with w, as follows:

$$\frac{\partial W_q}{\partial A} = \frac{\partial w}{\partial A} \frac{\partial W_q}{\partial w}, \quad \frac{\partial W_q}{\partial B} = \frac{\partial W_q}{\partial w} \frac{\partial w}{\partial B}$$
(22)

From Equation 20, the terms  $\frac{\partial w}{\partial A}$ ,  $\frac{\partial w}{\partial A}$ , and  $\frac{\partial W}{\partial w}$  can be expressed as follows:

$$\frac{\partial w}{\partial A} = \frac{\alpha B^{\top}}{s}, \quad \frac{\partial w}{\partial B} = \frac{\alpha A^{\top}}{s}, \quad \frac{\partial W}{\partial w} = \frac{\partial}{\partial w} (\tilde{R}(w)s + b) = s \frac{\partial \tilde{R}}{\partial w}$$
(23)

Therefore, by substitution of Equation 23 and applying STE on  $\frac{\partial \tilde{R}}{\partial w}$  from Equation 14 on Equation 22, the equation is simplified by the crossed-out products between the terms. As a result, the partial derivatives presented in Equation 10 can be derived as follows.

$$\frac{\partial W_q}{\partial A} = \frac{\partial w}{\partial A} \frac{\partial W_q}{\partial w} = \left(\frac{\alpha B^{\top}}{s}\right) \left(\frac{s\partial \tilde{R}}{\partial w}\right) = \begin{cases} \alpha B^{\top}, & \text{if } Q_n \le w \le Q_p \\ 0, & \text{otherwise} \end{cases}$$
(24)

$$\frac{\partial W_q}{\partial B} = \frac{\partial W}{\partial w} \frac{\partial w}{\partial B} = (s \frac{\partial \tilde{R}}{\partial w}) (\frac{\alpha A^{\top}}{s}) = \begin{cases} \alpha A^{\top}, & \text{if } Q_n \le w \le Q_p \\ 0, & \text{otherwise} \end{cases}$$
(25)

Finally, substitution of Equation 24 and 25 to Equation 21 derives the Equation 26 and 27.

$$\frac{\partial L}{\partial A} = \begin{cases} \alpha B^{\top} (\frac{\partial L}{\partial Y} X^{\top}), & \text{if } Q_n \le w \le Q_p \\ 0, & otherwise \end{cases}$$
(26)

$$\frac{\partial L}{\partial B} = \begin{cases} \alpha(\frac{\partial L}{\partial Y}X^{\top})A^{\top}, & \text{if } Q_n \le w \le Q_p \\ 0, & otherwise \end{cases}$$
(27)

This form closely resembles the original backpropagation structure of the LoRA parameters A, Bas shown in Equation 2, where the updates are expressed as  $\frac{\partial L}{\partial A} = \alpha \frac{\partial L}{\partial \tilde{X}} X^{\top} = \alpha (B^{\top} \frac{\partial L}{\partial Y}) X^{\top}$ , and  $\frac{\partial L}{\partial B} = \alpha \frac{\partial L}{\partial Y} \tilde{X}^{\top} = \alpha \frac{\partial L}{\partial Y} (AX)^{\top}$ , respectively. However, in L4Q, this process includes an added gating condition on the quantized weights, which accounts for the integration of quantization into the LoRA parameters. As a result, we conclude that the backward process of the L4Q layer, which integrates both quantization parameter learning and LoRA parameter adaptation, is designed to account for the impact of quantization on the LoRA parameter updates.

### В QUANTIZATION INITIALIZATION

We evaluate various quantization initialization schemes within L4Q, including method introduced in Section 3.3, LSQ+ (Bhalgat et al., 2020), and conventional symmetric and asymmetric quanti-zation parameter initialization. The methods are depicted as L4Q<sub>init</sub>, LSQ+<sub>init</sub>, Symm, Asymm, respectively. In specific, each methods can be represented as the equations below, with quantiza-tion scale s and bias b and group-wise aligned model weight W with quantization bit-width n and  $Q_n = -2^{n-1}, Q_p = 2^{n-1} - 1.$ 

b = 0

b =

 $LSQ+_{init}: \quad s = \frac{Max(|\mu - 3\sigma(W)|, |\mu + 3\sigma(W)|)}{2^{n-1}}$ (28)

$$Symm: \quad s = \frac{Max(Abs(W))}{2^{n-1}} \tag{30}$$

$$Asymm: \quad s = \frac{Max(W) - Min(W)}{Q_p - Q_n} \tag{32}$$

$$b = Max(W) - s \times Q_p = Min(W) - s \times Q_n$$
(33)

$$L4Q_{init}: \quad s = Max(|\frac{Min(W)}{Qn}|, |\frac{Max(W)}{Qp}|)$$
(34)

$$b = 0 \tag{35}$$

We report the detailed model accuracy evaluation results of L4Q fine-tuning across different initial-ization methods, along with the quantization error and clipping error for each method, measured both at the initialization point and the end of the training, in Table 4. The LLaMA-27B model was trained for a total of 12,800 iterations with batch size 128, using the same hyperparameters as those in the main evaluation.

Table 4: MMLU 5-shot benchmark and the sum of quantization errors for various quantization parameter initialization methods within L4Q on the LLaMA-2 7B model. Quantization errors are represented in order of  $10^6$  and clipping errors are represented in order of  $10^3$ .

				MMLU 5-shot					tial	Post-train	
Model	Method	#Bits	Human.	STEM	Social.	Others	Average	Equant	$E_{clip}$	Equant	$E_{clip}$
LLaMA-2 7B	LSQ+	4	26.7	26.8	26.2	22.9	25.7	11.8	278.0	11.8	360.6
	Symm	4	40.8	35.9	48.2	50.1	43.5	11.1	260.0	11.0	282.1
	Asymm	4	41.0	37.1	49.7	50.2	44.2	10.5	0.0	10.5	64.7
	L4Q	4	42.9	37.7	50.5	51.9	45.3	11.4	0.0	11.6	36.1

### С EFFECT OF LORA WARM-UP ON L4Q TRAINING

We investigated the effect of LoRA warm-up on L4Q training and found that a small number of LoRA warm-up steps is often beneficial for quickly restoring model accuracy. Using LLaMA-1 and LLaMA-2 7B models, we trained the models for a total of 12,800 iterations with 128 batches. The remaining training conditions, such as the quantization target layers and the rank of LoRA parameters, were kept consistent with the main experimental setup. The Commonsense QA evaluation results are presented in Table 5.

Table 5: Commonsense QA benchmark result with the variations of LoRA warm-up on LLaMA-1 and LLaMA-2 7B models. The numbers represent the accuracy(%) of each task.

Model	#Bits	Warmup	Hella.	PIQA	ARC-c	ARC-e	Winog.	BoolQ	OBQA	Avg.
LLaMA-1 7B	4	0	58.0	78.9	45.1	76.6	70.1	76.5	37.0	63.2
		10	58.0	78.5	44.9	76.4	69.7	76.5	36.4	62.9
		20	57.9	78.8	45.2	76.4	70.1	76.1	36.2	62.9
		40	57.5	78.6	45.6	76.0	69.8	75.8	35.4	62.7
		80	57.2	78.6	45.1	75.3	69.6	75.1	34.8	62.2
	3	0	55.8	77.2	41.9	73.6	68.1	76.1	32.8	60.8
		10	55.7	77.2	42.0	74.0	68.2	75.9	33.0	60.8
		20	55.0	77.5	41.5	74.0	68.0	74.9	32.2	60.4
		40	54.3	77.3	41.1	73.3	66.5	73.7	31.8	59.7
		80	52.6	76.8	39.5	71.7	65.4	71.0	31.8	58.4
LLaMA-27B	4	0	57.5	78.1	46.9	75.8	69.8	76.0	35.0	62.7
		10	57.1	78.2	46.0	76.0	69.5	75.8	35.8	62.6
		20	57.2	78.5	46.1	76.5	69.9	77.7	34.2	62.8
		40	57.1	78.2	45.7	76.5	69.5	77.3	34.8	62.7
		80	56.5	78.4	46.2	76.4	70.7	74.5	34.0	62.4
	3	0	55.6	77.9	42.1	74.3	69.5	72.0	34.0	60.8
		10	55.7	78.6	44.0	74.2	68.4	72.5	34.8	61.2
		20	55.5	77.8	43.0	74.6	69.0	71.7	34.2	60.8
		40	55.4	77.8	43.4	73.6	68.9	71.7	35.0	60.8
		80	53.9	76.8	42.1	73.7	67.3	71.8	34.4	60.0

Models with fewer than 20 LoRA warm-up steps performed best, as initializing LoRA parameters with minimal disturbance from quantization errors facilitated better convergence. Extending the LoRA warm-up beyond the quantization-aware fine-tuning phase, however, resulted in performance degradation, likely due to an insufficient number of training steps to fully compensate for quanti-zation errors. Considering the variability in impact across different models and configurations, we generally applied 10 LoRA warm-up steps in our experiments. 

### D ABLATIVE STUDY ON LORA RANK

We investigated the effect of LoRA rank size on L4Q training and found that a rank size of around 4 is sufficient for effective training. Using the LLaMA-1 7B model, we conducted training over 12,800 iterations with 128 batches. The remaining training conditions are consistent with the main experiments. The evaluation results for Commonsense QA and MMLU are presented in Table 6 and Table 7, respectively.

Table 6: Commonsense QA benchmark result on LLaMA-1 7B model. The numbers represent accuracy (%) for each task. 

Model	Rank	HellaSwag	PIQA	ARC-c	ARC-e	Winogrande	BoolQ	OBQA	Average
LLaMA-17B	1	56.9	79.1	43.7	76.0	69.9	76.6	35.4	62.5
	2	58.6	78.6	45.4	76.8	69.9	74.1	36.4	62.8
	4	57.8	79.1	45.3	76.0	69.5	76.1	34.8	62.7
	8	58.5	79.1	45.8	76.5	70.2	75.8	37.0	62.7
	16	58.3	78.8	44.5	75.8	69.8	78.3	35.2	63.0
	32	58.6	79.0	45.2	76.3	70.2	76.7	37.6	63.4
	64	58.4	79.0	45.5	76.5	70.0	78.7	37.8	63.2
	128	58.6	79.1	45.9	76.7	70.6	75.8	36.4	63.3

Table 7: MMLU benchmark result on LLaMA-1 7B model. The numbers represent accuracy (%) for each category.

				0-shot					5-shot		
Model	Rank	Hums.	STEM	Social	Others	Avg.	Hums.	STEM	Social	Others	Avg.
LLaMA-1 7B	1	31.8	29.2	32.9	37.0	32.7	32.3	30.7	35.2	39.9	34.4
	2	30.7	28.3	32.7	37.5	32.2	33.0	29.3	34.3	37.2	33.5
	4	32.4	32.1	36.7	39.4	34.9	32.9	31.4	38.7	39.5	35.4
	8	33.5	31.9	37.7	39.4	35.5	34.2	30.7	38.4	39.8	35.7
	16	30.0	31.6	35.1	38.1	33.5	32.2	31.2	35.3	38.1	34.8
	32	32.1	30.8	33.2	36.5	33.1	32.5	31.7	37.2	39.5	35.0
	64	33.4	31.2	36.1	38.8	34.8	33.6	31.1	35.9	39.5	35.0
	128	31.8	29.7	34.5	37.6	33.3	33.3	31.1	35.7	38.3	34.5

Increasing the rank beyond 32 or 4 does not lead to further performance improvements, which aligns with the observations in the original LoRA paper (Hu et al., 2022). Therefore, we generally applied a rank size of 4, considering that higher rank sizes introduce memory and computational overhead during training. 

### Ε **EXPERIMENTAL SETTINGS**

The baselines and L4Q are trained with AdamW optimizer Loshchilov & Hutter (2019) with a weight decay of 0.01. For the learning rate scheduler, a cosine decay scheduler with a linear warm-up through 10% of the total training steps. Learning rates are presented in Table 8.

Table 8: Learning rate conditions used to fine-tuning on each models for L4Q and baselines: QLoRA\*, QA-LoRA, and QAT-LoRA.

		Met	hods	
Model	QLoRA*	QA-LoRA	QAT-LoRA	L4Q
OpenLLaMA 3B	$1 \times 10^{-5}$	$2 \times 10^{-5}$	$5 \times 10^{-5}$	$5 \times 10^{-5}$
LLaMA-1 7B	$1 \times 10^{-5}$	$2 \times 10^{-5}$	$5 \times 10^{-5}$	$5 \times 10^{-5}$
LLaMA-1 13B	$1 \times 10^{-5}$	$5 \times 10^{-5}$	$4 \times 10^{-5}$	$4 \times 10^{-5}$
LLaMA-1 33B	$1 \times 10^{-5}$	$5 \times 10^{-5}$	$2 \times 10^{-4}$	$2 \times 10^{-4}$
LLaMA-2 7B	$2 \times 10^{-5}$	-	$2 \times 10^{-4}$	$2 \times 10^{-4}$
LLaMA-2 13B	$2 \times 10^{-5}$	-	$2 \times 10^{-4}$	$2 \times 10^{-4}$

Batch size is set to 128. 50K iterations for the baselines that utilize PTQ applied schemes, such as QLoRA\* and QA-LoRA, and 25K iterations for QAT involving schemes, such as QAT, QAT-LoRA, and L4Q are utilized to match the training latency. The training sequence length of training is set to above the maximum sequence length of the dataset, which is 2048, except for a 33B model with L4Q that is set to 128.

### F **TRAINING TIME**

We report the total training time of L4Q and the quantization-aware PEFT baselines, QLoRA\* and QA-LoRA, in Table 9. L4Q demonstrates similar time performance to the baselines, highlighting its scalability and time efficiency for larger model sizes, comparable to that of the baseline methods. 

Table 9: Training time (in hours) spent on fine-tuning on OpenLLaMA and LLaMA-1 models with a A100 GPU.

	OpenLLaMA		LLaMA-1	
Methods	3B	7B	13B	33B
QLoRA*	4.5	9.9	18.0	38.4
QA-LoRA	5.0	11.2	19.8	36.2
L4Q	4.4	10.1	16.9	39.6

### 1080 G THROUGHPUT AND SPEEDUP OF FULLY-QUANTIZED MODELS AND 1081 MIXED-PRECISION MODELS 1082

1083 We investigate the throughput and speedup of fully-quantized models and mixed-precision models, 1084 demonstrating that, although the number of LoRA parameters is negligible, it causes a noticeable drop in throughput when its forward path is not merged with that of the base linear layer. Using the 1086 LLaMA-1 7B, 13B, and 33B models, we conducted experiments to measure throughput (tokens per 1087 second) and compare speedup. In both fully-quantized and mixed-precision models, uniform quan-1088 tization is applied to the linear layers, except for the head layer (lm\_head), using the EXLLaMA2 kernel<sup>4</sup>, which is designed for 4-bit weight-only quantized inference. For fp16 computations in 1089 LoRA within mixed-precision models or the baseline, default GEMM kernels are used. We mea-1090 sured the elapsed time for inferencing 512 tokens over 2000 data points with batch sizes ranging 1091 from 1 to 64, calculating throughput by dividing the number of tokens, which is set to be 512, by 1092 the elapsed time. The results are presented in Table 10. 1093

1094 Table 10: Throughput (tokens/sec) and Speedup for LLaMA models. L4Q represents fully-1095 quantized models, and QLoRA\* represents mixed-precision models. 'OOM' indicates out-of-1096 memory cases. 1097

					Batch size	ze			
Model	Method	1	2	4	8	16	32	64	Speedu
LLaMA-1 7B	L4Q	38.04	75.12	148.63	216.51	255.64	276.43	318.43	1.81
	QLoRA*	24.80	47.88	96.51	184.06	234.79	247.79	299.94	1.33
	None	17.04	33.81	67.27	124.13	199.19	241.31	OOM	1.00
LLaMA-1 13B	L4Q	30.68	59.53	115.78	144.44	160.95	191.20	OOM	1.92
	QLoRA*	19.71	38.90	77.83	128.67	150.72	156.16	OOM	1.41
	None	13.67	26.97	53.23	85.47	124.17	OOM	OOM	1.00
LLaMA-1 30B	L4Q	20.43	40.05	64.00	73.29	79.77	OOM	OOM	2.25
	QLoRA*	13.22	25.48	50.81	66.89	75.11	OOM	OOM	1.44
	None	9.13	17.68	OOM	OOM	OOM	OOM	OOM	1.00

Fully-quantized models demonstrate a speedup of over 1.8x, while mixed-precision models achieve 1110 a maximum speedup of 1.4x, despite using the same quantization scheme and execution kernel, 1111 compared to the fp16 baselines. As a result, fully-quantized models achieve a 30% to 50% greater 1112 speedup compared to mixed-precision models. This demonstrates that L4O, which produces fully-1113 quantized models, offers higher inference efficiency and better hardware utilization than conven-1114 tional quantization-aware PEFT methods, such as QLoRA and LoftQ, which retain unmerged for-1115 ward paths for LoRA.

1116 1117

1119

1109

#### DETAILED RESULT ON MAIN EVALUATIONS Η 1118

We present the Commonsense QA and MMLU benchmark results with averaged accuracy score 1120 on Section 4. We present the detailed results of each benchmarks composed of several categories 1121 of tasks in Table 11 and Table 12 below. Through evaluation, we demonstrate that L4Q generally 1122 achieves higher accuracy in low-bit quantized models compared to both PTQ methods and PTQ-1123 based fine-tuning methods. Notably, L4Q surpasses the pre-trained models on the Commonsense 1124 QA benchmarks and on the MMLU benchmarks with LLaMA-1 7B and 33B models. In contrast, 1125 PTQ-based fine-tuning methods, including those that incorporate high-precision LoRA weights, 1126 show lower performance compared to both L4Q and the pre-trained models. These results emphasize 1127 the challenges of recovering from quantization errors with PTQ alone and highlight the effectiveness 1128 of L4Q's joint quantization and fine-tuning scheme.

1129

- 1131 1132
- 1133

<sup>&</sup>lt;sup>4</sup>https://github.com/turboderp/exllamav2.git

1136	Model	Method	#Bits	Hella.	PIQA	ARC-c	ARC-e	Winogr.	BoolQ	OBQA	Avg.
1137	OpenLLaMA 3B	None	16	48.8	75.0	33.9	69.2	61.6	66.9	28.2	54.8
1138		GPTQ	4	49.8	75.6	37.0	58.8	63.1	68.0 57.9	27.2	50.7
1139		OmniQ LoftO*	4	48.2	73.8	33.1	69.5 68.6	60.1	67.5	26.6	54.1
1140		QLoRA* OA-LoRA	4+16 4	48.4	74.3 74.9	33.0 33.8	69.4 69.2	61.5 61.9	67.1 66.7	26.8 26.2	54.4 54.5
1141		QAT-LoRA L4Q	4+16 4	48.8 49.1	74.5 74.9	35.0 35.2	70.1 69.8	61.9 61.1	65.2 67.7	27.0 27.4	54.6 55.0
1142		GPTQ	3	46.3	72.6	31.8	64.7 56.6	58.1	66.5	25.6	52.2
1143		LoftQ*	3+16	27.9	57.3	19.5	37.3	51.0	61.9	12.0	38.1
1144		QLoRA* QA-LoRA	3+16	45.6 46.3	72.6	29.3 28.9	61.6 66.0	59.7 59.5	64.2 63.4	24.4 23.8	51.0 51.5
1145		QAI-LOKA L4Q	3+16	46.7 47.2	74.1 75.0	33.2 32.3	67.2 68.3	60.5 60.9	64.1 67.2	26.4 27.0	53.2 54.0
1146	LLaMA 7B	None LoRA	16 16	57.0 58.3	78.7 78.8	41.9 45.7	75.3 76.1	69.9 70.6	75.1 78.7	34.4 35.4	61.7 63.4
1147		GPTQ OmniQ	4 4	53.9 55.7	77.7 77.7	40.3 38.8	73.5 67.5	67.9 65.3	72.9 72.5	30.0 29.2	59.4 58.1
1148		LoftQ*	4+16	57.8	79.2	43.1	76.9	69.8 70.0	75.8 74.6	35.4	62.6
1149		QA-LoRA OAT-LoRA	4 4+16	57.2 57.7	78.9 78.9	41.2 44.7	74.9 75.3	70.6 68.9	73.6 75.8	32.6 35.6	61.3 62.4
1150		LAQ	4	57.8	79.1	45.3	76.0	69.5	76.1	34.8	62.7
1151		OmniQ	3	40.0 54.0	71.9	32.4 35.6	63.4 64.9	63.0 64.7	71.2	24.0	56.5
1152		LoftQ* QLoRA*	3+16 3+16	43.4 53.9	68.9 76.2	33.0 39.3	65.5 71.5	56.5 68.9	58.5 72.8	23.0 31.0	49.8 59.1
1153		QA-LoRA QAT-LoRA	$3^{3}_{2}+16$	55.4 56.1	76.3	39.8 41.6	72.5	69.5 68.0	67.1 76.0	30.6 33.0	58.7 60.7
1154	LLaMA 13B	None	16	59.9	79.2	46.5	74.1	72.8	78.0	33.2	63.8
1155		LoRA GPTQ	16	60.8 58.9	79.7	50.3 46.5	78.6	72.3	80.2	34.8	65.2
1156		OmniQ LoftO*	4	58.6	79.7	43.8	73.5	70.5	68.7	28.4	60.4
1157		QLoRA*	4+16 4	59.6 60.1	79.0 79.2 79.0	46.5	77.1 77.0	72.5 71.4	78.1 67.1	33.4 36.2	63.8 62.5
1158		QAT-LoRA L4O	4+16 4	60.9 60.9	79.2 79.8	48.2 48.2	78.6 78.5	71.5 71.7	77.0 76.7	35.6 35.4	64.4 64.5
1159		GPTQ	3	57.3	77.3	42.6	73.0	71.0	74.6	31.4	61.0
1160		LoftQ*	3+16	47.8	72.1	37.6	70.8	58.3	65.5	25.8	54.0
1161		QLoRA* QA-LoRA	3+16	56.6 57.7	77.8 78.0	43.9 44.7	75.1 75.3	70.8 71.2	73.5 68.6	31.6 32.4	61.3 61.1
1162		L4Q	3+16	59.1 58.9	78.1 78.4	46.3 45.8	77.4	70.8	74.7 77.7	36.2 35.2	63.4
1163	LLaMA 33B	None LoRA	16 16	63.3 64.1	81.0 81.3	52.8 53.7	80.4 81.6	75.9 75.5	82.6 84.0	36.0 37.6	67.4 68.3
1104		GPTQ OmniQ	4 4	61.8 62.3	80.5 80.0	49.1 48.5	78.9 75.8	73.6 73.9	82.2 69.1	33.6 31.0	65.7 62.9
1166		LoftQ* OLoRA*	4+16 4+16	63.3 62.3	80.3 80.2	51.8 50.2	81.4 79.5	75.3 74 9	82.9 81.0	37.0 35.4	67.4 66.2
1167		QA-LoRA QAT-LoRA	4 4+16	62.8 62.3	80.3 81.3	50.1 53.0	79.5 81.6	75.1 74.9	73.2 82.7	36.4 35.4	65.3 67.3
1168		L4Q LoftO*	4	63.9	81.0	53.0 38.9	81.3	75.0	82.8 63.5	35.8	67.5
1169		QLoRA* OA-LoRA	3+16	59.6 61.1	78.7 79.6	46.0 47.8	76.8 78.0	72.5 73.8	81.6 79.3	34.6 33.0	64.3 64.6
1170		QAT-LoRA L4Q	3+16 3	63.3 63.0	81.0 81.0	52.8 52.6	81.3 81.4	75.5 75.5	82.8 82.8	35.4 35.4	67.4 67.4
1171	LLaMA 2 7B	None LoRA	16 16	57.1 57.9	78.1 78.9	43.4 48.0	76.3 77.4	69.1 70.3	77.7 75.8	31.4 34.8	61.9 63.3
1172		GPTQ	4	56.0	77.5	42.2	75.0	68.2 67.8	76.4	29.8	60.7
1173		LoftQ*	4+16	57.0	78.0	43.3	76.3	69.2	76.8	31.4	61.7
1174		QLORA* QA-LoRA	4+16	56.6 56.4	77.8 79.3	43.3	75.2 39.2	69.1 71.8	75.3 75.5	31.8 31.4	61.3 61.0
1175		L4Q	4+10	57.2	78.8	47.1	76.9	70.2	80.4	34.8	63.6
1176		GPTQ OmniQ	3 3	53.1 54.6	76.2 76.4	35.8 37.5	70.3 67.6	67.7 66.1	72.4 71.9	27.6 31.0	57.6 57.9
1177		LoftQ* OLoRA*	3+16 3+16	27.1 52.4	55.7 75.9	19.0 37.6	31.1 69.9	48.8 65.6	48.1 74.1	12.8 27.4	34.7 57.6
1178		QA-LoRA QAT-LoRA	3 3+16	56.5 52.0	77.8 75.2	42.3 39.3	74.7 71.1	68.0 65.1	30.8 69.9	43.8 29.3	56.3 57.4
1179	LL MA 2 12D	L4Q None	3	55.5 60.1	77.3	42.8	73.8	68.8	77.2	34.0	61.3
1180	Elawin 2 19b	LoRA	16	61.2	79.4	53.0	79.8	73.2	81.4	37.4	66.5
1181		ÖmniQ	4	59.0	78.1	43.7	71.3	68.7	66.6	32.0	59.9
1182		LottQ* QLoRA*	4+16 4+16	60.0 59.6	79.3 78.4	48.1 46.6	79.7 77.9	71.9 72.2	80.7 79.2	34.8 33.8	64.9 64.0
1183		QA-LOKA QAT-LoRA	4 4+16 4	59.4 59.5	78.8 80.1	48.4 51.2	40.9 79.2 70 7	72.3 71.5 71.0	80.7 80.9 82.2	34.4 34.4 35 e	64.5 65.8
1184		GPTQ	3	57.3	77.2	43.5	76.1	69.9	74.0	34.0	61.7
1185		UmniQ LoftO*	3	57.8 28.7	60.6	42.0	45.3	68.0 50.7	69.9 55.1	31.2	59.9 39.3
1186		QLoŘA* QA-LoRA	3+16 3	57.8 57.3	77.9 77.2	44.3 76.0	76.7 43.4	70.0 70.1	78.1 73.7	32.6 34.0	62.5 61.7
1187		QAT-LoRA L4Q	3+16 3	55.8 59.3	77.1 78.7	67.6 51.2	$76.0 \\ 78.5$	67.6 70.6	75.1 79.9	30.8 37.4	64.3 65.1

1134Table 11: Commonsense QA benchmark result. The numbers represent accuracy (%) of each task.

Madal	Mal	# <b>D</b> *4		OTEM	0-shot	04		11	CTEN 4	5-shot	04	
LLaMA 7B	None	<b>#Bits</b>	32.9	26.9	32.1	37.3	Avg.	33.9	30.6	38.2	38.2	AV:
LLaMA 13B	LoRA	16	36.1	31.5	36.9	40.6	36.3	34.4	30.3	39.9	43.1	36.
	OmniQ	4	31.1	26.7	29.8	35.5	30.9	31.1	29.8	35.5	37.5	33.
	LoftQ* QLoRA*	4+16 4+16	32.3 33.1	30.0 27.1 20.5	32.7 33.1	37.3 37.5	33.0 32.8	33.6 32.3	30.7 29.0	37.2 35.4	39.0 38.0	35. 33.
	QAT-LORA L4O	4 4 4	33.5 32.4	29.5 29.5 32.1	37.5 36.7	37.9 39.4	34.5 34.9	34.1 34.2	31.2 30.7	38.5 38.4	39.9 39.0 39.8	35.
	GPTQ OmniO	3 3	25.0 27.8	22.5 29.7	22.0 26.8	24.5 32.2	23.7 29.0	25.9 31.6	25.7 32.1	28.2 33.7	29.7 29.7	27. 31.
	LoftQ*	3+16	24.3	21.7	22.4	24.5	23.4	23.4	23.2	21.9	23.5	23.
	QA-LoRA OAT-LoRA	3+16 3	28.2 28.9	29.7 27.1	32.0 25.8	33.7 29.6	30.6 28.0	29.8 29.1	29.2 26.6	32.2 29.7	34.6 31.1	31. 29.
	L4Q None	3	29.5	27.8	32.1	33.3	30.6	29.3	31.0	33.5	30.4	31.
	LoRA	16	42.4	34.0	49.4	51.9	44.3	45.0	36.4	54.1	53.1	47.
	OmniQ	4	39.8	35.1	44.9	44.9	40.1	43.1	35.9	52.5	52.3	45.
	LoftQ* QLoRA*	4+16 4+16	39.0 39.0	34.8 35.7	47.8 47.5	48.5 47.2	42.4 42.1	43.4 43.8	34.3 35.3	52.3 52.0	52.1 52.8	45.
	QAT-LORA L4O	4+10 4 4	35.8 35.3 40.5	35.9 35.2 35.3	40.9 47.7 49.5	48.5 47.9 48.5	42.0 42.4 43.2	43.2 43.0 43.4	34.6 36.1	53.0 52.3	53.6 53.2	45. 45. 46.
	GPTQ	3	32.1	27.2	36.3 37.1	36.8 41.8	33.1 34.8	36.0 39.1	29.8 33.1	42.2	45.6 47.6	38. 41
	LoftQ*	3+16	25.1	24.7	24.0	26.1	25.0	24.9	27.4	24.1	25.0	25.
	QA-LoRA OAT-LoRA	3+16 3	37.6 35.9	31.6 28.4	41.9 42.4	44.3 43.5	38.8 37.5	38.0 34.8	34.1 31.5	44.5 43.0	47.9 44.8	40.
LLaMA 33B	L4Q None	3	38.5	33.2	44.7 62.2	47.2	40.7	39.3 54.4	34.0	46.6	48.4	41.
	LoRA	16	49.2	41.3	61.4	58.7	54.4	55.2	46.1	66.4	63.3	57.
	OmniQ	4	49.4	40.3	61.3	59.1	52.0	53.4	44.8	64.7	61.1	55.
	QLoRA*	4+16 4+16 4+16	49.2 48.5 50.2	40.2 39.0 39.8	60.8 59.7 60.9	58.0 57.8 58.9	51.8 51.0 52.3	54.6 54.5 55.0	44.5 44.2 45.5	65.2 63.4 65.1	61.5 60.5 61.8	56. 55. 56
	QAT-LoRA L4Q	4 4	45.2 50.8	39.7 42.1	56.6 61.5	55.5 59.4	48.9 53.3	52.7 53.5	43.5 46.6	63.4 66.1	61.0 61.8	55. 56.
	LoftQ* OLoRA*	3+16 3+16	24.7 41.8	24.0 34.6	23.2	26.4 52.3	24.6 45.6	24.3 46.4	23.2 40.9	22.9 57.9	25.6 56.7	24.
	QA-LoRA QAT-LoRA	3+16 3	47.7 41.5	38.9 37.2	58.0 54.4	56.6 53.1	50.1 46.1	46.8 45.4	41.2 39.9	58.2 55.6	57.5 55.0	50 48
LLaMA 2 7B	L4Q None	3	46.4 39.3	38.7	57.6 47.9	<u> </u>	49.4	46.4	41.6	58.2	58.8	50 45
	LoRA GPTQ	4	41.0	34.6	50.8 41.3	50.2 41.1	43.9	43.4	37.0	51.8 48.9	52.4 48.6	46.
	LoftQ*	4 4+16	37.7	34.6	47.2	45.7	41.0 38.5	42.4	37.7	49.5	49.8	45.
	QLoŘA* QA-LoRA	4+16 4+16	37.3 36.5	31.5 32.4	43.3 40.6	42.7 42.6	38.6 37.9	42.1 41.8	35.9 35.2	50.2 48.6	51.1 50.2	44. 43.
	QAI-LORA L4Q	4	37.3 38.7	32.3 33.8	43.5 45.6	43.0 46.4	38.9 40.9	42.0 42.9	35.7 37.7	49.6 50.5	50.8 51.9	44.
	GPTQ OmniQ	3 3	28.9 33.1	28.5 30.4	35.3 39.1	33.7 35.5	31.3 34.3	36.0 34.1	31.7 32.4	39.3 41.6	43.5 44.4	37. 37.
	LoftQ* QLoRA*	3+16 3+16	24.1 30.2	21.3 29.1	21.8 36.0	23.8 35.5	22.9 32.5	23.7 35.4	26.1 32.5	22.4 40.5	24.9 42.7	24 37
	QA-LORA QAT-LoRA	3+16 3	31.1 28.9 31.0	27.2 27.8 32.7	33.9 34.7 38.6	33.8 33.7 39.2	31.5 31.0 34.9	34.2 36.0 34.3	31.2 31.6 32.3	39.9 39.5 42.2	42.7 43.4 44.9	36 37 38
LLaMA 2 13B	None	16	47.8	42.3	60.5	59.2 59.4	52.1 52.1	51.5 52.0 54.4	43.8	63.0	61.2	54 54
	GPTQ	4	46.5	40.2	57.7	56.8	50.0	52.3	43.1	62.7	61.5	54
	LoftQ*	4 4+16	47.8	41.9	60.1	58.9	51.8	53.0	43.0	62.5	60.5	54 54
	QLoRA* QA-LoRA	4+16 4+16	46.9 47.5	40.9 41.0	58.8 58.8	57.6 56.8	50.7 50.7	51.3 50.3	43.1 42.9	62.5 62.3	60.8 60.7	54 53
	L4Q	4	40.5	40.8	58.5 60.4	58.4	51.9	53.6	44.3	62.5	60.5	55
	GPTQ OmniQ	3	43.5 42.3	37.3 38.9	53.6 54.5	51.8 51.3	46.3 46.3	46.3 43.4	42.7 43.0	57.3 58.8	56.2 56.5	50 50
	LoftQ* QLoRA*	3+16 3+16	24.2 43.9	21.7 38.3	23.8 53.9	23.7 52.2	23.5 46.8	24.6 48.5	28.5 41.1	24.0 57.7	27.4 55.9	26 50
	QA-LoRA QAT-LoRA	3+16 3	42.5 43.4	37.3 37.4	53.0 53.7	52.1 52.1	45.9 46.4	44.8 47.5	40.5 41.5	56.3 56.3	55.7 55.3	48 49

Table 12: MMLU benchmark result. The numbers represent accuracy(%) of each task.