
Exploring Post-training Quantization in LLMs from Comprehensive Study to Low Rank Compensation

Zhewei Yao*, Xiaoxia Wu*, Cheng Li, Stephen Youn, Yuxiong He
Microsoft
{zhewei Yao, xiaoxia Wu, cheng Li, stephen.youn, yuxiong He}@microsoft.com

Abstract

Post-training quantization (PTQ) has emerged as a promising technique for mitigating memory consumption and computational costs in large language models (LLMs). However, a systematic examination of various quantization schemes, model families, and quantization bit precision has been absent from the literature. In this paper, we conduct a comprehensive analysis of these factors by investigating the effects of PTQ on weight-only, activation-only, and weight-and-activation quantization using diverse methods such as round-to-nearest (RTN), GPTQ, ZeroQuant, and their variants. We apply these methods to two distinct model families with parameters ranging from 125M to 176B. Our contributions include: (1) a sensitivity analysis revealing that activation quantization is generally more susceptible to weight quantization, with smaller models often outperforming larger models in terms of activation quantization; (2) an evaluation and comparison of existing PTQ methods to optimize model size reduction while minimizing the impact on accuracy, revealing that none of the current methods can achieve the original model quality for quantization with either INT4-weight or INT4-weight-and-INT8-activation; (3) based on these insights, we propose an optimized method called Low-Rank Compensation (LoRC), which employs low-rank matrices to enhance model quality recovery with a minimal increase in model size.

1 Introduction

As the progression of hardware bandwidth lags behind that of computational demand [13], the resource demands of extra-large models such as MT-NLG-530B [25]—which necessitates the deployment of multiple nodes for operation—escalate, adding to the complexities of cross-node communication. This has emphasized the urgency to curtail both the size and computational expense of Large Language Models (LLMs). An increasingly effective solution to these issues is post-training quantization (PTQ). This method aids in the reduction of training prerequisites while simultaneously lowering the bit precision of weights and activations to either INT4 or INT8.

While the effectiveness of post-training quantization (PTQ) has been underscored in a number of recent studies [31, 11, 30, 7], a comprehensive, systematic investigation into several key dimensions of this technique remains to be undertaken. Specifically, the extant literature falls short in providing thorough coverage of the functionality of various PTQ methods or the sensitivity of disparate models. Moreover, despite current quantization methods demonstrating promising results in the reduction of model sizes, the question persists as to whether these methods are achieving their optimal potential in minimizing Large Language Models (LLMs) sizes.

*Equal Contribution. Code will be released as a part of <https://github.com/microsoft/DeepSpeed>

With these observations in mind, our study sets forth to address two salient questions: (1) When subjected to quantization, do LLMs of varying sizes and pretraining data exhibit similar behavior? (2) Are existing quantization methods truly leveraging their full potential in reducing the sizes of LLMs?

Contribution. To elucidate these queries, we undertake an exhaustive examination of the impact of PTQ on weight-only, activation-only, and combined weight-and-activation quantization. This investigation incorporates a range of PTQ methods, including round-to-nearest (RTN), GPTQ [11], ZeroQuant [31], and their respective variants. To broaden the scope of our analysis, we focus on two distinct model families, OPT [35] and BLOOM [23], spanning model sizes from 125M to a massive 176B. Our code will be made available for reproduction. In summary, we make the following contributions:

(1) We provide a thorough **sensitivity analysis** to demonstrate that a) Activation quantization is generally more sensitive to weight quantization; We carry out a detailed evaluation and comparison of current PTQ methods, utilizing optimal configurations to maximize model size reduction while minimizing accuracy impact. We found that the current existing method can barely achieve less than 0.1 PPL points degradation for quantization with either INT4-weight or INT4-weight-and-INT8-activation (W4A8). To recover the 0.1 PPL, we strive to push the boundaries of employing **fine-grained quantization** (FGQ) techniques. We observe FGQ is able to recovered points degradation of <0.1 PPL for large models (>13B) for INT4 weight quantization.

(2) Based on the above understanding, we further optimize existing methods and introduce a technique called **Low Rank Compensation** (LoRC), which employs low-rank matrix factorization on the quantization error matrix. Complementary to FGQ, LoRC plays a crucial role in enhancing the full model quality recovery, while there is little increase of the model size.

In Figure 1, we provide model size and quality trade-offs for both OPT and BLOOM families. As can be seen, using LoRC on top of PTQ methods from [31, 11] and fine-grained quantization, we set a new quantization Pareto frontier for LLMs. Meanwhile, we recommend the following setting for quantizing LLMs with LoRC (Note that activation quantization should be only applied if necessary): (1) For larger models (>10B), fine-grained (block size 64–256) 4-bit weight quantization plus 8-bit activation quantization (block size 64–256) with PTQ can be used for real deployment; (2) For middle-size models (<10B and >1B), per-row INT8 quantization plus fine-grained (block size 64–256) INT8 activation quantization can be used with PTQ from [11, 31]; (3) For smaller models (<1B), per-row W8A8 (INT8 weight and INT8 activation) RTN is enough based on [31].

2 Sensitivity Analysis and Comparison Between Existing Methods

There are mainly two categories of PTQ for LLMs, i.e., weight-only quantization [11] and weight-and-activation quantization [6, 31, 30]. In the latter, it is uniformly observed across all studies that activation quantization demonstrates greater sensitivity than weight quantization. However, prior research tends to concentrate on a single (family) model to emphasize the necessity of their proposed quantization technique. A comprehensive and systematic evaluation of this PTQ methodology, particularly the sensitivity of weight/activation quantization for varying model sizes and distinct model families, has not yet to be undertaken. Hence, we conduct an examination on both the OPT [35] and BLOOM [23] families.

Sensitivity setting. We use the zero-shot validation perplexity (PPL) differential on three datasets, namely, Wikitext-2 [20], PTB [19], and C4 [22], before and after the quantization of these LLMs to

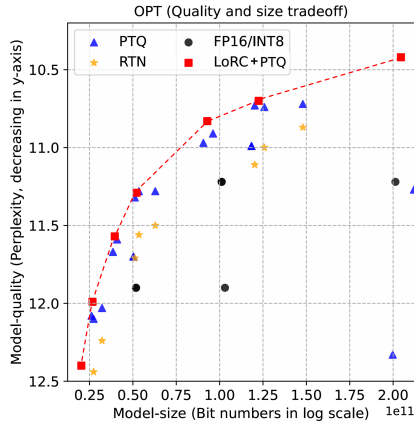


Figure 1: The model size and quality trade-off of different quantization methods on models from OPT families. Here PTQ (with fine-grained quantization) represents the method from [31, 11], RTN means the naive round-to-nearest baseline (with fine-grained quantization), and FP16/INT8 is used as the no-accuracy-loss baseline. LoRC is our proposed method that works seamless with PTQ.

Table 1: Classification of quantization sensitivity (or quantization loss). The sensitivity increases from *Class-1* to *Class-3*.

| Class | <i>Class-1</i> | <i>Class-2</i> | <i>Class-3</i> |
|-----------------|----------------|----------------|----------------|
| PPL Degradation | ≤0.1 | >0.1 & ≤0.5 | >0.5 |

illustrate their sensitivity, as PPL is significantly correlated to zero-shot/few-shot accuracy measurement [7]. Specifically, a higher PPL drop indicates enhanced quantization sensitivity. For simplicity, we also categorize quantization sensitivity (or quantization loss) into three different classes as depicted in Table 1. Notably, the threshold is chosen because when the model size approximately doubles (e.g., 13B vs. 30B, and 30B vs. 66B), the PPL improvement is about 0.5 (see Table 2). The sensitivity (or loss) incrementally increases as the class number ascends. From a practical standpoint, we favor lower quantization sensitivity (accuracy loss), making **Class-1** the optimal-loss PTQ.

Robustness of Weight-only Quantization for Large Models. The results of weight-only quantization in OPT and BLOOM models are summarized in Table 2. INT8 weight-only quantization, either symmetric or asymmetric, results in negligible accuracy loss (less than 0.05, i.e., **Class-1**). Consequently, for tasks oriented towards generation, FP16 weight can simply be replaced with INT8 weight to reduce memory usage. For INT4 quantization, the asymmetric method outperforms the symmetric approach in accuracy, attributable to its superior utilization of the quantization range. Interestingly, larger models exhibit better tolerance to low-precision quantization (i.e., INT4) than smaller models, with a few exceptions such as OPT-66B. Particularly, BLOOM-176B shows PPL degradation (around 0.3 points) in **Class-2**, which could explain why the large GLM-130B [34] can operate with INT4 weight-only quantization out of the box with acceptable accuracy impact.

Challenge Encountered in Activation Quantization for Large Models.

Activation quantization has consistently proven more difficult than weight quantization [31, 6], as illustrated in Table 2. When compared to weight-only quantization, activation-only quantization indicates that asymmetric quantization can significantly improved performance over symmetric quantization. Moreover, contrary to weight-only quantization, smaller models typically exhibit better tolerance to activation quantization, as their hidden dimension is smaller and the activation dynamic range is also narrower than larger models [31]. It should be noted that for models larger than 10B, all fall into **Class-3**, indicating a degradation of >0.5 PPL.

Comparison between GPTQ and ZeroQuant. Numerous lightweight optimization-based methods have been proposed, which update the model weights during quantization. These methods such as [31, 11, 30], unlike quantization-aware training, only require a small portion of the training data and a limited training time. Particularly, GPTQ [11] and ZeroQuant [31], have proven to be effective and efficient in terms of GPU resources, time cost, and data usage for INT4 weight quantization.² We focus on the variants of GPTQ and ZeroQuant as well as the most straightforward baseline, round-to-nearest neighborhood (RTN). Due to space limit, we defer the experimental results in Appendix C. The main conclusions are: (1) GPTQ typically performs better for weight-only quantization, while ZeroQuant yields superior results for weight and activation quantization. (2) The tested optimization-based methods cannot achieve **Class-1** quantization error for either INT4 weight-only or W4A8 quantization with the exception of GPTQ on OPT-30B with weight-only quantization.

3 Proposed Method to Further Push the Limit of Post-training Quantization

Building on the investigation and conclusions drawn from previous sections, it has become apparent that there is still a need for an advanced methodology to further refine the existing methods, with the objective of fully realizing the original fp16 PPL quality. In this section, we introduce a simple yet effective method called **LoRC** (Low Rank Compensation) to optimize the current existing quantization error and further bridge the gap between the quality of the original model and its quantized counterparts. LoRC is inspired by the employment of low-rank matrix factorization on the quantization error matrix $E := W - \hat{W}$, where W represents the original weight and \hat{W} is the quantized weight. LoRC approximates the error E with $\hat{E} = \hat{U}\hat{V}$ by using two low-rank matrices \hat{U} and \hat{V} . This results in a more accurate approximation of the original weight matrix W by $\hat{W}_{\text{lorc}} = \hat{W} + \hat{E}$, reducing errors: $\|W - \hat{W}\| \geq \|W - \hat{W}_{\text{lorc}}\|$. LoRC consists of two steps:

²We tested the method proposed by [30] but did not find it better than others for INT4 weight quantization.

Table 3: W#^{asym}-A16 quantization with # being 4-bit, 3-bit and 2-bit on OPT and BLOOM (BLM).

| Bits | LoRC | Coarse-grained weight quantization (per-row block-size) | | | | | Fine-grained quantization on weight (256 block-size) | | | | |
|-------|------|---|---------|---------|---------|----------|--|---------|---------|---------|----------|
| | | OPT-6.7b | OPT-13b | OPT-30b | OPT-66b | BLM-176b | OPT-6.7b | OPT-13b | OPT-30b | OPT-66b | BLM-176b |
| W8A16 | | 11.90 | 11.22 | 10.70 | 10.33 | 10.90 | 11.90 | 11.22 | 10.70 | 10.33 | 10.90 |
| W4A16 | ✗ | 12.28 | 11.42 | 10.78 | 10.78 | 11.02 | 12.05 | 11.28 | 10.74 | 10.50 | 10.95 |
| | ✓ | 12.10 | 11.36 | 10.76 | 10.34 | 10.98 | 11.99 | 11.29 | 10.70 | 10.29 | 10.93 |
| W3A16 | ✗ | 14.18 | 12.43 | 11.28 | 17.77 | 49.46 | 12.79 | 11.63 | 10.9 | 11.34 | 11.13 |
| | ✓ | 13.00 | 11.90 | 11.14 | 10.63 | 11.30 | 12.40 | 11.57 | 10.83 | 10.42 | 11.08 |
| W2A16 | ✗ | 120.56 | 40.17 | 25.74 | 225.45 | Explode | 23.13 | 15.55 | 12.68 | 308.49 | 12.64 |
| | ✓ | 24.17 | 18.53 | 14.39 | 13.01 | 14.15 | 16.27 | 14.30 | 12.37 | 11.54 | 12.21 |

Step I: Implement Singular Value Decomposition (SVD) on the error matrix $E = U\Sigma V$, where $U \in \mathbb{R}^{d_{in} \times d_{in}}$ and $V \in \mathbb{R}^{d_{out} \times d_{out}}$ are unitary matrices, and $\Sigma \in \mathbb{R}^{d_{in} \times d_{out}}$ is a diagonal matrix with its diagonal elements ordered in a descending manner.

Step II: We formulate the matrix $\hat{E} = \hat{U}\hat{V}$ where $\hat{U} = U_m(\Sigma_m)^{\frac{1}{2}}$ and $\hat{V} = (\Sigma_m)^{\frac{1}{2}}V_m$. Here, $U_m = U_{:,1:m} \in \mathbb{R}^{d_{in} \times m}$, $V_m = V_{1:m,:} \in \mathbb{R}^{m \times d_{out}}$, and $\Sigma_m = \Sigma_{1:m,1:m} \in \mathbb{R}^{m \times m}$.

The objective of LoRC is to achieve a good approximation of the error matrix E using low-rank matrices, with minimal impact on the increase in model size. For instance, consider the standard transformer models [27], where each layer is comprised of a multi-headed attention (MHA) module and a multi-linear perception (MLP) module.

Significantly, LoRC can be viewed as a supplementary feature to existing quantization methodologies such as RTN, GPTQ, and ZeroQuant-Local/Global, and can be seamlessly integrated with FGQ. We have conducted experiments to evaluate the performance of LoRC on both OPT and BLOOM, applying 4-bit, 3-bit, and 2-bit weights by setting the activation to FP16. Based on the discoveries in the preceding sections, we utilize the GPTQ quantization strategy. To gain a comprehensive understanding of LoRC, we include the results with and without the application of FGQ. The datasets and hyperparameters are consistent with those detailed in earlier sections.

Evaluation Results. The findings are showcased in Table 3, split into two sections: coarse-grained weight quantization (per-row) and fine-grained quantization (block-size 256). Notably, we observe that the two low-rank matrices, \hat{U} and \hat{V} , can be quantized to 8-bit without any performance discrepancy (Table 9). Thus, the two low-rank matrices for LoRC in Table 3 are INT8 with a low-rank dimension of $m = 8$. Several key observations can be made. Firstly, LoRC consistently boosts performance across all bit sizes and block sizes, as indicated by the lower perplexity scores when LoRC is activated. Secondly, the enhancement brought about by LoRC becomes more substantial as the bit size diminishes, especially noticeable for W2A16, which displays a markedly greater impact compared to W4A16 and W3A16 in most scenarios. Lastly, the combination of fine-grained quantization with LoRC yields the most impressive results, underscoring the efficacy of LoRC when integrated with FGQ. Overall, the results emphasize the benefits of using LoRC for enhanced performance in weight quantization and its compatibility with FGQ. Notably, recovering the last 0.05-0.1 perplexity can be challenging, but with LoRC, we nearly recover the original model quality.

We provide ablation study for the low-rank dimension n in Appendix D.

4 Discussion

Conclusion. In this work, we provide a comprehensive study of post-training quantization (PTQ) on large language models with different PTQ methods (e.g., RTN, GPTQ, ZeroQuant), and with different quantization coverage (weight-only and weight-and-activation quantization), etc. We find that PTQ methods are critical to improving the quantized model quality, and that fine-grained quantization (FGQ) can bring acceptable accuracy and model size trade-off. Finally, we introduced an optimization technique called Low Rank Compensation (LoRC), which works synergistically with PTQ and FGQ, playing a crucial role in enhancing full model quality recovery with a minimal increase in model size.

Limitation. Despite quantizing over 10,000 experiments, our study was constrained by our computing resources. This restriction made us choose between diversifying the model sizes and varying the tasks. We strategically limited our datasets to WikiText, PTB, and C4 to concentrate on a broad range of quantization methods. Consequently, our general findings are more robust concerning the two model families and three datasets examined in this paper. However, caution should be exercised when generalizing these findings to tasks that are dissimilar to those covered in this study.

References

- [1] Haoli Bai, Wei Zhang, Lu Hou, Lifeng Shang, Jing Jin, Xin Jiang, Qun Liu, Michael Lyu, and Irwin King. Binarybert: Pushing the limit of bert quantization. *arXiv preprint arXiv:2012.15701*, 2020.
- [2] Big-Science. Bloom inference. <https://github.com/huggingface/transformers-bloom-inference/tree/main/bloom-inference-scripts>, 2022.
- [3] Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Understanding and overcoming the challenges of efficient transformer quantization. *arXiv preprint arXiv:2109.12948*, 2021.
- [4] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [5] Bitu Darvish Rouhani, Daniel Lo, Ritchie Zhao, Ming Liu, Jeremy Fowers, Kalin Ovtcharov, Anna Vinogradsky, Sarah Massengill, Lita Yang, Ray Bittner, et al. Pushing the limits of narrow precision inferencing at cloud scale with microsoft floating point. *Advances in neural information processing systems*, 33:10271–10281, 2020.
- [6] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm. int8 (): 8-bit matrix multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- [7] Tim Dettmers and Luke Zettlemoyer. The case for 4-bit precision: k-bit inference scaling laws. *arXiv preprint arXiv:2212.09720*, 2022.
- [8] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. Learned step size quantization. *arXiv preprint arXiv:1902.08153*, 2019.
- [9] Angela Fan, Pierre Stock, Benjamin Graham, Edouard Grave, Remi Gribonval, Herve Jegou, and Armand Joulin. Training with quantization noise for extreme fixed-point compression. *arXiv preprint arXiv:2004.07320*, 2020.
- [10] Elias Frantar and Dan Alistarh. Optimal brain compression: A framework for accurate post-training quantization and pruning. *arXiv preprint arXiv:2208.11580*, 2022.
- [11] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- [12] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. *arXiv preprint arXiv:2103.13630*, 2021.
- [13] Amir Gholami, Zhewei Yao, Sehoon Kim, Michael W Mahoney, and Kurt Keutzer. Ai and memory wall. *RiseLab Medium Post*, 2021.
- [14] Babak Hassibi and David G Stork. Second order derivatives for network pruning: Optimal brain surgeon. In *Advances in neural information processing systems*, pages 164–171, 1993.
- [15] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *Workshop paper in NIPS*, 2014.
- [16] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*, 2019.
- [17] Sehoon Kim, Amir Gholami, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. I-bert: Integer-only bert quantization. In *International conference on machine learning*, pages 5506–5518. PMLR, 2021.
- [18] Yann LeCun, John S Denker, and Sara A Solla. Optimal brain damage. In *Advances in neural information processing systems*, pages 598–605, 1990.

- [19] Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. *Using Large Corpora*, page 273, 1994.
- [20] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *International Conference on Learning Representations*, 2017.
- [21] Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*, 2018.
- [22] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2019.
- [23] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- [24] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Q-BERT: Hessian based ultra low precision quantization of bert. In *AAAI*, pages 8815–8821, 2020.
- [25] Shaden Smith, Mostofa Patwary, Brandon Norrick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, et al. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *arXiv preprint arXiv:2201.11990*, 2022.
- [26] Chaofan Tao, Lu Hou, Wei Zhang, Lifeng Shang, Xin Jiang, Qun Liu, Ping Luo, and Ngai Wong. Compression of generative pre-trained language models via quantization. *arXiv preprint arXiv:2203.10705*, 2022.
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [28] Xiaoxia Wu, Cheng Li, Reza Yazdani Aminabadi, Zhewei Yao, and Yuxiong He. Understanding int4 quantization for transformer models: Latency speedup, composability, and failure cases. *arXiv preprint arXiv:2301.12017*, 2023.
- [29] Xiaoxia Wu, Zhewei Yao, Minjia Zhang, Conglong Li, and Yuxiong He. Extreme compression for pre-trained transformers made simple and efficient. *arXiv preprint arXiv:2206.01859*, 2022.
- [30] Guangxuan Xiao, Ji Lin, Mickael Seznec, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. *arXiv preprint arXiv:2211.10438*, 2022.
- [31] Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers. *arXiv preprint arXiv:2206.01861*, 2022.
- [32] Ali Hadi Zadeh, Isak Edo, Omar Mohamed Awad, and Andreas Moshovos. Gobo: Quantizing attention-based nlp models for low latency and energy efficient inference. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 811–824. IEEE, 2020.
- [33] Ofir Zafir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8BERT: Quantized 8bit bert. *arXiv preprint arXiv:1910.06188*, 2019.
- [34] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022.
- [35] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.

- [36] Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao Chen, Xin Jiang, and Qun Liu. Ternarybert: Distillation-aware ultra-low bit bert. *arXiv preprint arXiv:2009.12812*, 2020.

A Related Work

Different quantization methods [24, 33, 9, 36, 1, 8, 26, 17] for transformer-based models [27] have been explored for a while. However, most of those works need quantization-aware finetuning or even expensive quantization-aware knowledge distillation [15]. Due to the cost of training/finetuning LLMs [21, 16, 26, 29, 28], it is a challenge for practitioners/researchers to do finetuning/distillation on those LLMs, particularly for models like GPT-3-175B [4] and BLOOM-176B [23].

Post-training quantization (PTQ) [32, 3] is an alternative way to quantize the model with no/minimal finetuning requirement. Along this line, several recent works focus on LLMs (beyond the million-parameter scale). [31] proposes vector-based INT8 quantization with layer-by-layer knowledge distillation to overcome the training cost and quantization error introduced by LLMs. [6] uses similar vector-based INT8 quantization weight plus mixed-precision (INT8/FP16) quantization for activation to overcome the sensitivity of activation quantization. However, the inference speed of [6] is generally even slower than FP16 baseline [2] due to the difficulty of implementing mixed-precision calculation within a single tensor. More recently, [11] extends OBQ [10, 14, 18] on LLMs for INT4 weight-only quantization and shows great efficiency on quantization and latency, and [30] shows the outliers from activations can be smoothed out by migrating the quantization difficulty from activations to its associated weights. However, [30] can only work for W8A8 quantization as lower weight precision (INT4) itself already leads to significant accuracy degradation, and the accuracy drop is larger than 0.1 PPL points, which as discussed in the later section is sub-optimal. [7] shows the scaling law of weight-only quantization with the simplest round-to-nearest baseline, but it does not consider the weight-and-activation quantization and/or the above PTQ optimization methods. As can be seen from Figure 1, by using PTQ optimization methods, the model quality can be significantly improved.

Different than existing works, our paper extensively tests the effect of (1) different quantization schemes, e.g., symmetric and asymmetric quantization, (2) different PTQ methods, e.g., [31, 11], (3) different model families, e.g., [23, 35], (4) different quantization coverage, e.g., weight-only and weight-and-activation quantization, and (5) other discussions, e.g., the effect of quantization granularity. As such, we provide a much more comprehensive understanding of post-training quantization for large language models compared to the previous works.

B Background of Quantization

Quantization maps floating point (e.g., FP16/FP32) numbers to integer numbers (e.g., INT4/INT8) so that lower memory usage (weight quantization) and faster integer arithmetic (weight-and-activation quantization) can be achieved compared to the floating point format. In this work, we are focusing on uniform quantization, i.e.,

$$Q(x) = \text{INT}((x - Z)/S) - Z, \tag{1}$$

where Q is the quantization function, x is a floating point input vector/tensor, S is a real valued scaling factor, and Z is an integer zero point. Based on different settings, the quantization method can be viewed as (1) symmetric vs. asymmetric quantization ($Z = 0$ or not), (2) fine-grained vs. coarse-grained quantization (how to partition the input x and get its associated scaling factor, e.g., matrix wise or row wise). See [12] for more details.

Throughout this work, we focus on post-training quantization (PTQ), i.e., no or minimal training effort is applied after quantization, for which large accuracy degradation usually exhibits for coarse-grained quantization (per matrix/tensor) due to their large quantization error. As such, we focus on fine-grained quantization. Particularly, we use the per-row quantization (one row of the weight matrix or one token for the activation) from [31] as our coarsest-grained quantization method, and we use block-k quantization (for every k elements, they have their own scaling factor and/or zero point) as our finer-grained quantization scheme.

C Are existing quantization methods optimally harnessing the potential to minimize LLMs sizes?

Numerous lightweight optimization-based methods have been proposed, which update the model weights during quantization. These methods such as [31, 11, 30], unlike quantization-aware training, only require a small portion of the training data and a limited training time. Particularly, GPTQ [11]

Table 4: The evaluation results of different PTQ methods on OPT and BLOOM (BLM) with asymmetric quantization on weight or (and) activation.

| Precision | Method | OPT-6.7b | OPT-13b | OPT-30b | OPT-66b | BLM-1.7b | BLM-3b | BLM-7.1b | BLM-176b |
|-----------|------------|----------|---------|---------|---------|----------|--------|----------|----------|
| W16A16 | | 11.90 | 11.22 | 10.70 | 10.33 | 20.43 | 17.58 | 14.96 | 10.90 |
| W4A16 | RTN | 13.44 | 12.09 | 11.52 | 31.52 | 22.47 | 19.01 | 15.90 | 11.20 |
| | GPTQ | 12.28 | 11.42 | 10.78 | 10.52 | 21.58 | 18.33 | 15.50 | 11.02 |
| | ZQ-Local* | 12.46 | 11.64 | 11.05 | 10.79 | 21.70 | 18.50 | 15.55 | 11.11 |
| | ZQ-Global* | 12.38 | 11.62 | 11.04 | 10.68 | 21.38 | 18.33 | 15.52 | 11.05 |
| W4A8 | RTN | 14.80 | 26.36 | 86.26 | 815.00 | 22.75 | 19.17 | 16.19 | 12.22 |
| | GPTQ | 13.88 | 17.28 | 20.71 | 648.69 | 21.71 | 18.44 | 15.75 | 11.86 |
| | ZQ-Local* | 13.24 | 14.23 | 18.53 | 16.32 | 21.86 | 18.66 | 15.75 | 11.19 |
| | ZQ-Global* | 13.17 | 13.07 | 14.65 | 37.82 | 21.43 | 18.39 | 15.58 | 11.49 |

and ZeroQuant [31], have proven to be effective and efficient in terms of GPU resources, time cost, and data usage for INT4 weight quantization.³ In this work, we focus on the variants of GPTQ and ZeroQuant as well as the most straightforward baseline, round-to-nearest neighborhood (RTN).

RTN directly applies PTQ on the trained data and follows the procedure detailed in Section B to perform the quantization. Specifically, for asymmetric quantization, we set $S = \max(x) - \min(x)$ and $Z = \min(x)$.

GPTQ extends the OBQ [10]. It tries to optimize the following non-linear least square problem, $\min_{\hat{W}} \|Wx - \hat{W}x\|_2^2$ where W is the weight, x is the activation, and \hat{W} is a quantized weight. GPTQ employs second-order methods to obtain a closed-form solution. In addition, the quantization for each weight matrix is performed column-/row-wisely and the quantization errors from previous columns will be passed to those columns not yet quantized. See [10, 11] for more details.

ZQ-Global is the original method proposed in [31], where authors treat each layer as a small neural network (a.k.a., subnetwork) and use the FP16 subnetwork as the teacher model to distill the quantized one with a few hundred iterations, i.e., $\min_{\hat{\theta}} \|f_{\theta}(x) - f_{\hat{\theta}}(x)\|_2^2$, where θ is a set of weights, $\hat{\theta}$ is the quantized version, f_{θ} is the subnetwork with parameters θ , and x is the input. Thus, it can significantly reduce the GPU resource requirement and time cost.

ZQ-Local is an extension mode of ZQ-Global for further GPU requirement reduction and training cost reduction. Particularly, instead of using each transformer layer as the subnetwork, we treat each linear layer as the subnetwork. This method can be viewed as an iterative first-order optimization method (e.g., SGD) to solve $\min_{\hat{W}} \|Wx - \hat{W}x\|_2^2$.

Experimental Setup. We compare the four methods mentioned above on weight-only and weight-and-activation quantization. As weight quantization is always static (i.e., it does not change during inference), there is virtually no system performance difference between symmetric and asymmetric quantization.⁴ We use asymmetric quantization for better accuracy, and the conclusions would hold similarly for symmetric quantization. An interesting finding for ZeroQuant is that the hyperparameters (e.g., learning rate and its scheduler) provided in the original work [31] are sub-optimal. In this work, we find the best configurations for ZQ-Local and ZQ-Global and denote them as ZQ-Local* and ZQ-Global*, respectively, with the best tuned results. To ensure consistent and comparable results, we set a fixed random seed for our experiments. In the context of post-training quantization, varying the random seed has minimal impact on the final results.

Evaluation of Weight-only Quantization. The results from weight-only quantization using OPT and Bloom are presented in Table 4. The findings indicate that the larger models tend to be less sensitive to INT4 weight-only quantization. This observation holds true across all methods (RTN, GPTQ, ZQ-Local*, and ZQ-Global*) with the exception of OPT-66B, which shows greater degradation than OPT-30B. It is noteworthy that light-weight optimization-based methods significantly outperform the RTN baseline in terms of accuracy. For instance, these methods substantially reduce the degradation in perplexity of OPT-30B/66B compared to baseline. Most quantized models with parameters greater than 6.7B fall under Class II, indicating their potential for real-world applications. For instance, the quality of INT4 OPT-30B (66B) is superior to that of INT8 OPT-13B (30B).

Evaluation of Weight and Activation Quantization. The evaluation results for existing methods using W4A8 quantization are presented in Table 4. The three light-weight optimization-based

³We tested the method proposed by [30] but did not find it better than others for INT4 weight quantization.

⁴The bias term (a.k.a., the zero point) can be simply fused into the previous activation quantization kernel [31].

methods outperform RTN significantly, underscoring their efficacy. However, all of the results fall into either *Class-2* or *Class-3*. This suggests that for certain applications, it might be more beneficial to use smaller models with fewer parameters rather than larger, quantized models.

Among quantization-based methods, ZQ-Global* and ZQ-Local* generally outperform GPTQ, which is anticipated given that GPTQ was originally designed for weight-only quantization. ZQ-Global* performs better than ZQ-Local* in most cases except for the two largest models, OPT-66B and Bloom-176B, despite having larger trainable parameters in one step. This again signifies the need for a more suitable and advanced optimization method for large language models (LLMs).

Finding 2 on Comparisons. (1) GPTQ typically performs better for weight-only quantization, while ZeroQuant (including both ZQ-Global* and ZQ-Local*) yields superior results for weight and activation quantization. (2) The tested optimization-based methods cannot achieve *Class-1* quantization error for either INT4 weight-only or W4A8 quantization with the exception of GPTQ on OPT-30B with weight-only quantization.

C.1 Fine-grained Quantization and Its Evaluation

With PTQ and row-wise quantization, achieving *Class-1* quantization error is challenging for both weight-only and weight-and-activation quantization. Generally, utilizing a smaller model with INT8 weight is more advantageous than employing a model that is twice as large with INT4 weight.

One potential solution to this issue is the implementation of finer-grained quantization schemes [5], where every k elements possess their own scaling factor and/or zero point. This approach can significantly reduce quantization error. In the extreme case, where every single element has its own scaling factor, the original FP16 number can be precisely recovered. Importantly, block-k quantization can be implemented on modern GPUs, one of the most prevalent deep learning architectures, since the compute unit (streaming multiprocessor) of GPUs processes tiles of data (e.g., 128 by 128 tiling size) for matrix computation.

Although fine-grained quantization can substantially narrow the gap between the quantized tensor and its floating-point counterpart, the application of RTN still results in a non-trivial accuracy gap. Consequently, we build upon fine-grained quantization by employing existing optimization-based methods to further enhance accuracy. Specifically, we utilize GPTQ and ZQ-Global for all models and settings and apply ZQ-Local to OPT-66B and Bloom-176B. For the hyperparameters used in ZQ-Global and ZQ-Local, we select the top three identified in Section C for all models, except for Bloom-176B, for which we only use the top-performing hyperparameter to reduce training costs.

4-bit Weight Quantization. We hereby present the W4A16 results for OPT and BLOOM, as delineated in Table 5, corresponding to an array of quantization block sizes. The performance sees a significant improvement with smaller block sizes compared to per-row quantization. The point of diminishing returns, however, varies for different model sizes. For example, smaller models (such as OPT-6.7B and BLOOM-1.7b) continue to see substantial gains until the block size reduces to 32. In contrast, for larger models (those exceeding 10B, with OPT-66B as the exception), the benefits derived from smaller block sizes wane rapidly around block-256/512. Most crucially, for models equal to or larger than 13B, a smaller quantization block size results in quantization error being classified under *Class-1*, indicating virtually negligible degradation in accuracy.

Table 6: OPT W4^{asym}-A8 with various block-size out of the best result from GPTQ, ZQ-Local, and ZQ-Global on OPT and BLOOM (BLM).

| Precision | block-size (W/A) | OPT-6.7b | OPT-13b | OPT-30b | OPT-66b | BLM-1.7b | BLM-3b | BLM-7.1b | BLM-176b |
|-----------|---------------------------|----------|---------|---------|---------|----------|--------|----------|----------|
| W4A16 | 128 NA | 12.10 | 11.28 | 10.74 | 10.44 | 20.92 | 17.90 | 15.17 | 10.94 |
| W4A8 | Case-1: per-row per-row | 13.17 | 13.07 | 14.65 | 16.32 | 21.43 | 18.39 | 15.58 | 11.19 |
| | Case-2: per-row 128 | 12.29 | 11.45 | 10.80 | 10.61 | 21.59 | 18.31 | 15.52 | 11.03 |
| | Case-3: 128 128 | 12.04 | 11.31 | 10.75 | 10.45 | 21.27 | 17.86 | 15.19 | 10.96 |

Activation Quantization (W4A8). To comprehend the benefits of fine-grained

Table 7: BLOOM-176B with different quantization block sizes on activation. Here weight is asymmetrically quantized with block size 128.

| A8 Block Size | 1024 | 512 | 256 | 128 | 32 |
|---------------|-------|-------|-------|-------|-------|
| PPL | 10.98 | 10.97 | 10.95 | 10.95 | 10.95 |

Table 5: Results of **W4^{asym}-A16** quantization with various block-size out of the best result from optimization-based methods on OPT and BLOOM (BLM). N/A means that the block size is not divisible by the hidden size.

| Block-size | OPT-6.7b | OPT-13b | OPT-30b | OPT-66b | BLM-1.7b | BLM-3b | BLM-7.1b | BLM-176b |
|------------|----------|---------|---------|---------|----------|--------|----------|----------|
| W16A16 | 11.90 | 11.22 | 10.70 | 10.33 | 20.43 | 17.58 | 14.96 | 10.90 |
| Per-row | 12.28 | 11.42 | 10.78 | 10.52 | 21.38 | 18.33 | 15.50 | 11.02 |
| 1024 | 12.16 | 11.36 | 10.75 | 10.52 | 31.03 | N/A | 15.24 | 10.96 |
| 512 | 12.08 | 11.32 | 10.73 | 10.52 | 20.93 | 17.99 | 15.20 | 10.95 |
| 256 | 12.05 | 11.28 | 10.74 | 10.50 | 20.95 | 17.97 | 15.18 | 10.95 |
| 128 | 12.10 | 11.28 | 10.74 | 10.44 | 20.92 | 17.90 | 15.17 | 10.94 |
| 32 | 12.03 | 11.28 | 10.72 | 10.41 | 20.82 | 17.88 | 15.16 | 10.95 |

quantization on activation, we analyze the quantization between per-row and a block size of 128, with INT4 weight, as highlighted in Table 6. For models of considerable size, specifically those equal to or exceeding 1B, the application of such fine-grained activation quantization (Case-1) results in a substantial reduction in quantization error compared to per-row activation (Case-2). By implementing fine-grained activation quantization with weight quantization (Case-3), we are able to almost restore the performance to the level of their W4A16 counterparts.

Furthermore, we detail the impacts of varying activation quantization block sizes in Table 7 on BLOOM-176B, with INT4 weight. A trend of superior accuracy is observed with smaller block sizes in contrast to larger ones. However, the enhancement in performance reaches a saturation point when the size smaller or equal to 256, which corresponds to the range of values INT8 can represent. Despite INT8’s capability to signify 256 distinct values, activation quantization errors persist due to the application of uniform quantization.

Finding 3 on FGQ. (1) Larger models ($\geq 10B$) are capable of attaining *Class-1* error for 4-bit quantization. These models can leverage low-precision quantization as the model size with INT4 is similar to an INT8 model that is half its size, with improved accuracy. On the other hand, smaller models ($\leq 10B$) typically reach only *Class-2* or *Class-3* error levels. (2) For larger models ($> 10B$), the difference between fine-grained weight-and-activation quantization and fine-grained weight-only quantization is insignificant. (3) The advantage of fine-grained activation quantization fades for larger models when the block size reaches 256.

D Additional Results for LORC

The objective of LoRC is to achieve a good approximation of the error matrix E using low-rank matrices, with minimal impact on the increase in model size. For instance, consider the standard transformer models [27], where each layer is comprised of a multi-headed attention (MHA) module and a multi-linear perception (MLP) module. Let h represent the hidden dimension and l the number of layers. The total number of parameters is $12lh^2$ as each layer contains $4h^2$ for MHA (for key, query, value, and projection matrices), and $8h^2$ for MLP (two matrices of sizes $h \times 4h$ and $4h \times h$). With the addition of low-rank LoRC to the six matrices in each layer, the total number of parameters for l layers would amount to $18hml$.⁵ Consequently, the ratio of parameters added to the existing model is $3m/2h$. It’s important to note that the low-rank dimension m can be as small as 4 or 8 (which we will discuss in detail in a later section) while the standard hidden dimension $h \geq 768$, making the number $3m/2h \leq 0.016$.

⁵In the MHA module, LoRC contributes $2hm$ to each of key, query, value, and the projection matrices. In the MLP module, LoRC contributes $8hm$ and $2hm$ respectively to the matrices of dimensions $h \times 4h$ and $4h \times h$.

Table 9: Results of W4^{asym} A16 quantization with LoRC approximating $\hat{E} = \hat{U}\hat{V}$ on OPT model family. \hat{U} and \hat{V} can be represented with FP16 or INT8, of which the performance are represented below. There is hardly any difference between FP16 and INT8.

| LoRC \hat{U}, \hat{V} | Coarse-grained weight quantization | | | | Fain-grained weight Quantization | | |
|----------------------------|------------------------------------|-------|-------|-------|----------------------------------|--------|--------|
| | 6.7b | 13b | 30b | 66b | 6.7b | 13b | 30b |
| FP16 | 12.08 | 11.35 | 10.76 | 10.31 | 11.993 | 11.290 | 10.703 |
| INT8 | 12.10 | 11.36 | 10.76 | 10.34 | 11.987 | 11.290 | 10.700 |

Table 8: W4A16 quantization with LoRC by varying the low-rank dimension m .

| LoRC-dim m | OPT-1.3b | OPT-6.7b | OPT-30b |
|------------------|----------|----------|---------|
| $m = 0$ baseline | 15.95 | 12.06 | 10.73 |
| $m = 1$ | 15.93 | 12.01 | 10.73 |
| $m = 4$ | 15.73 | 12.00 | 10.72 |
| $m = 8$ | 15.76 | 11.99 | 10.70 |
| $m = 16$ | 15.74 | 12.00 | 10.69 |
| $m = 32$ | 15.71 | 12.01 | 10.69 |

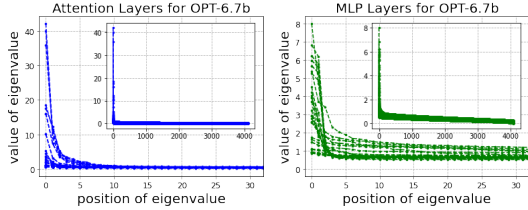


Figure 2: Eigenvalues of the Error matrix E for W4A16

Ablation Study on the Low Rank Dimension m . An essential aspect of the LoRC method is on the optimal low-rank dimension, denoted as m , explained in **Step II**. To explore this, we varied m in the range of 1, 4, 8, 16, and 32 for OPT-1.3b/6.7b/30b models, and applied W4A16 GPTQ quantization. The outcomes are depicted in Table 8, indicating that the enhancements achieved through LoRC begin to plateau as the dimension m surpasses 4. The most optimal performance for OPT-6.7b is realized when $m = 8$.

This observation may seem counterintuitive initially, as one might anticipate that larger LoRC dimensions would yield more significant improvements. To gain a more comprehensive understanding, we conducted an analysis of the eigenvalues of the actual error matrix $E = W - \hat{W}$ for each matrix. By randomly selecting 20 matrices from MHA and MLP layers, we plotted the eigenvalues of E as a curve, depicted in Figure 2. The two plots reveal a rapid flattening of eigenvalues after index 8, which elucidates why increasing the LoRC dimension does not considerably enhance performance. Hence, a sensible dimension for \hat{U} and \hat{V} in the LoRC methodology could be 8.⁶

⁶Please note that this observation is only true for PTQ. If one uses quantize-aware training (QAT) and let \hat{U} and \hat{V} updated during QAT, we arrive at contrasting conclusions.