
From comprehensive study to low rank compensation: Exploring post-training quantization in LLMs

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

19 1 Introduction

20 Large language models (LLMs) like Codex [15] and ChatGPT [24] have demonstrated breakthrough
21 performance across various benchmarks, such as natural language understanding and generation, and
22 are now integrated into everyday applications. However, efficiently serving LLMs has become a
23 pressing concern due to their significant memory consumption and computational demands. Unlike
24 classification or diffusion models, LLMs present unique challenges, as they involve two distinct
25 phases: prompt and generation. The prompt phase is primarily compute-bound, while the generation
26 phase, with low batch size and KV cache, is mainly memory-bound [26].

27 As the progression of hardware bandwidth lags behind that of computational demand [14], the resource
28 demands of extra-large models such as MT-NLG-530B [30]—which necessitates the deployment of
29 multiple nodes for operation—escalate, adding to the complexities of cross-node communication.
30 This has emphasized the urgency to curtail both the size and computational expense of Large Language
31 Models (LLMs). An increasingly effective solution to these issues is post-training quantization (PTQ).
32 This method aids in the reduction of training prerequisites while simultaneously lowering the bit
33 precision of weights and activations to either INT4 or INT8.

34 While the effectiveness of post-training quantization (PTQ) has been underscored in a number of
35 recent studies [36, 12, 35, 7], a comprehensive, systematic investigation into several key dimensions
36 of this technique remains to be undertaken. Specifically, the extant literature falls short in providing

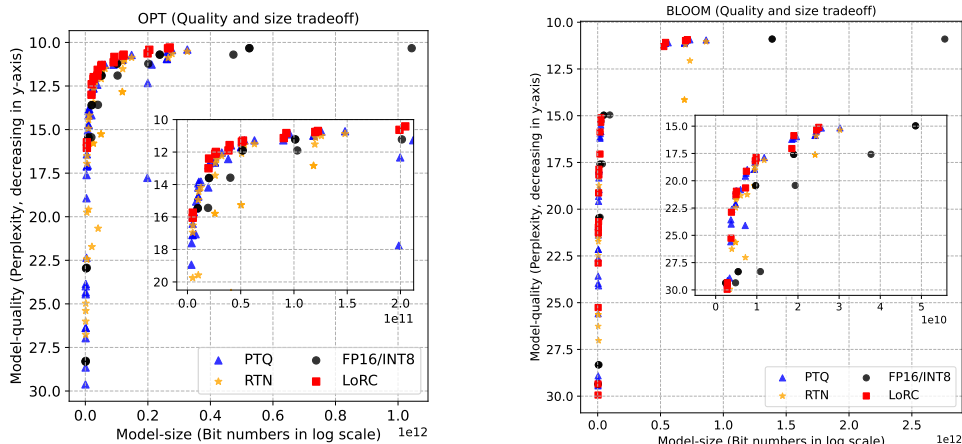


Figure 1: The model size and quality trade-off of different quantization methods on models from OPT and BLOOM families. Here PTQ (with fine-grained quantization) represents the method from [36, 12], RTN means the naive round-to-nearest baseline (with fine-grained quantization as well), and FP16/INT8 is used as the no-accuracy-loss baseline. LoRC is our proposed method that works seamlessly with PTQ. Note that we drop all diverged points for better visualization. For all detailed numbers, please see Appendix E.

37 thorough coverage of the functionality of various PTQ methods or the sensitivity of disparate models.
 38 Moreover, despite current quantization methods demonstrating promising results in the reduction of
 39 model sizes, the question persists as to whether these methods are achieving their optimal potential in
 40 minimizing Large Language Models (LLMs) sizes.

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

44 **Contribution.** To elucidate these queries, we undertake an exhaustive examination of the impact
 45 of PTQ on weight-only, activation-only, and combined weight-and-activation quantization. This
 46 investigation incorporates a range of PTQ methods, including round-to-nearest (RTN), GPTQ [12],
 47 ZeroQuant [36], and their respective variants. To broaden the scope of our analysis, we focus on
 48 two distinct model families, OPT [40] and BLOOM [28], spanning model sizes from 125M to a
 49 massive 176B. Our code will be made available for reproduction. In summary, we make the following
 50 contributions:

51 (1) We provide a thorough **sensitivity analysis** to demonstrate that a) Activation quantization is
 52 generally more sensitive to weight quantization; Smaller models usually have better activation
 53 quantization performance than the relative larger model. b) Different model families show different
 54 INT8 activation quantization behaviors; Particularly for large models, BLOOM-176B has small
 55 accuracy drops (about 1 perplexity or PPL) but OPT-30B and -66B experience worse performance.

56 (2) We carry out a detailed evaluation and comparison of current PTQ methods, utilizing optimal
 57 configurations to maximize model size reduction while minimizing accuracy impact. We found that
 58 the current existing method can barely achieve less than 0.1 PPL points degradation for quantization
 59 with either INT4-weight or INT4-weight-and-INT8-activation (W4A8). To recover the 0.1 PPL, we
 60 strive to push the boundaries of employing **fine-grained quantization** (FGQ) techniques. We observe
 61 FGQ is able to recovered points degradation of <0.1 PPL for large models ($>13B$) for INT4 weight
 62 quantization, but there are still non-negligible model quality drops.

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

67 In Figure 1, we provide model size and quality trade-offs for both OPT and BLOOM families.
 68 As can be seen, using LoRC on top of PTQ methods from [36, 12] and fine-grained quantization,

69 we set a new quantization Pareto frontier for LLMs. Meanwhile, we recommend the following
 70 setting for quantizing LLMs with LoRC (Note that activation quantization should be only applied if
 71 necessary): (1) For larger models (>10B), fine-grained (block size 64–256) 4-bit weight quantization
 72 plus 8-bit activation quantization (block size 64–256) with PTQ can be used for real deployment; (2)
 73 For middle-size models (<10B and >1B), per-row INT8 quantization plus fine-grained (block size
 74 64–256) INT8 activation quantization can be used with PTQ from [12, 36]; (3) For smaller models
 75 (<1B), per-row W8A8 (INT8 weight and INT8 activation) RTN is enough based on [36].

76 2 Related Work

77 Different quantization methods [29, 38, 9, 41, 1, 8, 31, 19] for transformer-based models [32] have
 78 been explored for a while. However, most of those works need quantization-aware finetuning or
 79 even expensive quantization-aware knowledge distillation [17]. Due to the cost of training/finetuning
 80 LLMs [25, 18, 31, 34, 33], it is a challenge for practitioners/researchers to do finetuning/distillation
 81 on those LLMs, particularly for models like GPT-3-175B [4] and BLOOM-176B [28].

82 Post-training quantization (PTQ) [37, 3] is an alternative way to quantize the model with no/minimal
 83 finetuning requirement. Along this line, several recent works focus on LLMs (beyond the million-
 84 parameter scale). [36] proposes vector-based INT8 quantization with layer-by-layer knowledge
 85 distillation to overcome the training cost and quantization error introduced by LLMs. [6] uses similar
 86 vector-based INT8 quantization weight plus mixed-precision (INT8/FP16) quantization for activation
 87 to overcome the sensitivity of activation quantization. However, the inference speed of [6] is generally
 88 even slower than FP16 baseline [2] due to the difficulty of implementing mixed-precision calculation
 89 within a single tensor. More recently, [12] extends OBQ [10, 16, 21] on LLMs for INT4 weight-only
 90 quantization and shows great efficiency on quantization and latency, and [35] shows the outliers
 91 from activations can be smoothed out by migrating the quantization difficulty from activations to its
 92 associated weights. However, [35] can only work for W8A8 quantization as lower weight precision
 93 (INT4) itself already leads to significant accuracy degradation, and the accuracy drop is larger than
 94 0.1 PPL points, which as discussed in the later section is sub-optimal. [7] shows the scaling law of
 95 weight-only quantization with the simplest round-to-nearest baseline, but it does not consider the
 96 weight-and-activation quantization and/or the above PTQ optimization methods. As can be seen
 97 from Figure 1, by using PTQ optimization methods, the model quality can be significantly improved.
 98 Please also see Appendix E for more detailed numbers.

99 Different than existing works, our paper extensively tests the effect of (1) different quantization
 100 schemes, e.g., symmetric and asymmetric quantization, (2) different PTQ methods, e.g., [36, 12],
 101 (3) different model families, e.g., [28, 40], (4) different quantization coverage, e.g., weight-only
 102 and weight-and-activation quantization, and (5) other discussions, e.g., the effect of quantization
 103 granularity. As such, we provide a much more comprehensive understanding of post-training
 104 quantization for large language models compared to the previous works.

105 3 Would different model families behave similarly on quantization?

106 There are mainly two categories of PTQ for LLMs, i.e., weight-only quantization [12] and weight-
 107 and-activation quantization [6, 36, 35]. In the latter, it is uniformly observed across all studies that
 108 activation quantization demonstrates greater sensitivity than weight quantization. However, prior
 109 research tends to concentrate on a single (family) model to emphasize the necessity of their proposed
 110 quantization technique. A comprehensive and systematic evaluation of this PTQ methodology,
 111 particularly the sensitivity of weight/activation quantization for varying model sizes and distinct
 112 model families, has yet to be undertaken. Hence, we conduct an examination on both the OPT [40]
 113 and BLOOM [28] families to elucidate the quantization sensitivity of weight and activation.

114 **Sensitivity setting.** We use the zero-shot validation
 115 perplexity (PPL) differential on three datasets, namely,
 116 Wikitext-2 [23], PTB [22], and C4 [27], before and
 117 after the quantization of these LLMs to illustrate their
 118 sensitivity, as PPL is significantly correlated to zero-
 119 shot/few-shot accuracy measurement [7]. Specifically,
 120 a higher PPL drop indicates enhanced quantization sen-

Table 1: Classification of quantization sensi-
 tivity (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 |

121 sitivity. For simplicity, we also categorize quantization sensitivity (or quantization loss) into three
 122 different classes as depicted in Table 1. Notably, the threshold is chosen because when the model
 123 size approximately doubles (e.g., 13B vs. 30B, and 30B vs. 66B), the PPL improvement is about 0.5
 124 (see Table 2). The sensitivity (or loss) incrementally increases as the class number ascends. From
 125 a practical standpoint, we favor lower quantization sensitivity (accuracy loss), making *Class-1* the
 126 optimal-loss post-training quantization.

127 We employ both symmetric and asymmetric quantization to gauge the quantization sensitivity and
 128 highlight the advantage of asymmetric quantization. Particularly, we implement per-row quantiza-
 129 tion [12] for weight quantization and per-token quantization for activation [36].

130 **Robustness of Weight-only Quantization for Large Models.** The results of weight-only quanti-
 131 zation in OPT and BLOOM models are summarized in Table 2. INT8 weight-only quantization,
 132 either symmetric or asymmetric, results in negligible accuracy loss (less than 0.05, i.e., *Class-1*).
 133 Consequently, for tasks oriented towards generation, FP16 weight can simply be replaced with INT8
 134 weight to reduce memory usage. For INT4 quantization, the asymmetric method outperforms the
 135 symmetric approach in accuracy, attributable to its superior utilization of the quantization range.
 136 Interestingly, larger models exhibit better tolerance to low-precision quantization (i.e., INT4) than
 137 smaller models, with a few exceptions such as OPT-66B.¹ Particularly, BLOOM-176B shows PPL
 138 degradation (around 0.3 points) in *Class-2*, which could explain why the large GLM-130B [39] can
 139 operate with INT4 weight-only quantization out of the box with acceptable accuracy impact.

Table 2: Average PPL of OPT and BLOOM (BLM). See Table E.1 for all results.

| Precision | OPT-6.7b | OPT-13b | OPT-30b | OPT-66b | BLM-1.7b | BLM-3b | BLM-7.1b | BLM-176b |
|-------------------------|----------|---------|---------|---------|----------|--------|----------|----------|
| W16-A16 | 11.90 | 11.22 | 10.70 | 10.33 | 20.43 | 17.58 | 14.96 | 10.90 |
| W8 ^{sym} -A16 | 11.90 | 11.22 | 10.70 | 10.33 | 20.43 | 17.59 | 14.97 | 10.90 |
| W8 ^{asym} -A16 | 11.90 | 11.22 | 10.70 | 10.33 | 20.45 | 17.59 | 14.97 | 10.90 |
| W4 ^{sym} -A16 | 14.36 | 12.73 | 11.77 | 97.05 | 23.18 | 19.36 | 16.27 | 11.28 |
| W4 ^{asym} -A16 | 13.44 | 12.09 | 11.52 | 31.52 | 22.47 | 19.01 | 15.90 | 11.20 |
| W16-A8 ^{sym} | 26.04 | 3171.49 | 2048.21 | 2638.09 | 20.68 | 17.73 | 15.28 | 12.10 |
| W16-A8 ^{asym} | 12.62 | 15.36 | 23.57 | 561.35 | 20.52 | 17.65 | 15.14 | 11.62 |

140 **Challenge Encountered in Activation Quantization for Large Models.** Activation quantization
 141 has consistently proven more difficult than weight quantization [36, 6], as illustrated in Table 2. When
 142 compared to weight-only quantization, activation-only quantization indicates that asymmetric quanti-
 143 zation can significantly improved performance over symmetric quantization. Moreover, contrary to
 144 weight-only quantization, smaller models typically exhibit better tolerance to activation quantization,
 145 as their hidden dimension is smaller and the activation dynamic range is also narrower than larger
 146 models [36]. It should be noted that for models larger than 10B, all fall into *Class-3*, indicating a
 147 degradation of more than 0.5 PPL points.

148 The last two rows of Table 2 show that different model families exhibit significantly different
 149 behaviors. BLOOM does not exhibit divergence issues even up to a model size of 176B, whereas OPT
 150 displays very poor performance from a model size of 6.7B (larger models with INT8 activation have
 151 even worse PPL). This could again be attributed to the Layer Norm issue within the OPT-family¹.

Findings 1 on Sensitivity Analysis. (1) INT8 weight-only quantization can serve as a stan-
 dard method for reducing memory costs in LLMs, with negligible degradation in accuracy. (2) INT4
 weight-only quantization for small models results in substantial accuracy degrada-
 tion (*Class-3*), but this effect lessens as the model size increases (*Class-2*). (3) Contrary
 to (2), INT8 activation results in minimal accuracy drops for small models (*Class-1*) but
 larger models exhibit greater drops (*Class-3*). (4) With INT8 activation, BLOOM shows no
 divergence issues up to a model size of 176B, whereas OPT performs poorly from $\geq 6.7B$
 model sizes.

152

¹[12] discovered that OPT-66B has a high proportion of dead neurons in the early layers, which might
 influence the compression capability. We also identify another potential reason: the Layer Norm of the OPT-
 family is not well trained (except OPT-350M), with the weight and the bias being all 1's and 0's, respectively.

153 **4 Are existing quantization methods optimally harnessing the potential to**
154 **minimize LLMs sizes?**

155 Numerous lightweight optimization-based methods have been proposed, which update the model
156 weights during quantization. These methods such as [36, 12, 35], unlike quantization-aware training,
157 only require a small portion of the training data and a limited training time. Particularly, GPTQ [12]
158 and ZeroQuant [36], have proven to be effective and efficient in terms of GPU resources, time cost,
159 and data usage for INT4 weight quantization.² In this work, we focus on the variants of GPTQ and
160 ZeroQuant as well as the most straightforward baseline, round-to-nearest neighborhood (RTN).

161 **RTN** directly applies PTQ on the trained data and follows the procedure detailed in Section A to
162 perform the quantization. Specifically, for symmetric quantization, we set $S = \max(\text{abs}(x))$ and
163 $Z = 0$; for asymmetric quantization, we set $S = \max(x) - \min(x)$ and $Z = \min(x)$.

164 **GPTQ** extends the OBQ [10]. It tries to optimize the following non-linear least square problem,
165 $\min_{\hat{W}} \|Wx - \hat{W}x\|_2^2$ where W is the weight, x is the activation, and \hat{W} is a quantized weight. GPTQ
166 employs second-order methods to obtain a closed-form solution. In addition, the quantization for each
167 weight matrix is performed column-/row-wisely and the quantization errors from previous columns
168 will be passed to those columns not yet quantized. See [10, 12] for more details.

169 **ZQ-Global** is the original method proposed in [36], where authors treat each layer as a small neural
170 network (a.k.a., subnetwork) and use the FP16 subnetwork as the teacher model to distill the quantized
171 one with a few hundred iterations, i.e., $\min_{\hat{\theta}} |f_{\theta}(x) - f_{\hat{\theta}}(x)|^2$, where θ is a set of weights, $\hat{\theta}$ is
172 the quantized version, f_{θ} is the subnetwork with parameters θ , and x is the input. Thus, it can
173 significantly reduce the GPU resource requirement and time cost.

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

178 **Experimental Setup.** We compare the four methods mentioned above on weight-only and weight-
179 and-activation quantization. As weight quantization is always static (i.e., it does not change during
180 inference), there is virtually no system performance difference between symmetric and asymmetric
181 quantization.³ We use asymmetric quantization for better accuracy, and the conclusions would hold
182 similarly for symmetric quantization. For parameters used for GPTQ, ZQ-Local, and ZQ-Global,
183 please refer to Appendix B. An interesting finding for ZeroQuant is that the hyperparameters (e.g.,
184 learning rate and its scheduler) provided in the original work [36] are sub-optimal. In this work,
185 we find the best configurations for ZQ-Local and ZQ-Global and denote them as ZQ-Local* and
186 ZQ-Global*, respectively, with the best tuned results. To ensure consistent and comparable results,
187 we set a fixed random seed for our experiments. In the context of post-training quantization, varying
188 the random seed has minimal impact on the final results, as indicated in more detail in Table B.1.

189 **Evaluation of Weight-only Quantization.** The results from weight-only quantization using OPT and
190 Bloom are presented in Table 3. The findings indicate that the larger models tend to be less sensitive
191 to INT4 weight-only quantization. This observation holds true across all methods (RTN, GPTQ,
192 ZQ-Local*, and ZQ-Global*) with the exception of OPT-66B, which shows greater degradation than
193 OPT-30B. It is noteworthy that light-weight optimization-based methods significantly outperform the
194 RTN baseline in terms of accuracy. For instance, these methods substantially reduce the degradation
195 in perplexity of OPT-30B/66B compared to baseline. Most quantized models with parameters greater
196 than 6.7B fall under Class II, indicating their potential for real-world applications. For instance, the
197 quality of INT4 OPT-30B (66B) is superior to that of INT8 OPT-13B (30B).

198 Among the optimization-based methods, ZQ-Global* generally performs better on smaller models
199 (those with fewer than 1B parameters), while GPTQ excels on larger models. ZQ-Local* does not
200 outperform GPTQ or ZQ-Global*—a reasonable outcome given that GPTQ employs a closed-form
201 solution to solve the non-linear quadratic problem and ZQ-Global* optimizes a larger subnetwork.
202 The inferior performance of ZQ-Global* compared to GPTQ for larger models is unexpected since
203 ZQ-Global* optimizes an entire transformer layer while GPTQ only optimizes a single linear layer.

²We tested the method proposed by [35] but did not find it better than others for INT4 weight quantization.

³The bias term (a.k.a., the zero point) can be simply fused into the previous activation quantization kernel [36].

Table 3: The evaluation results of different PTQ methods on OPT and BLOOM (BLM) with asymmetric quantization on weight or (and) activation. See more details in Table E.3 and Table E.6.

| 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 |

204 A plausible explanation is that larger models are more sensitive to weight updates, necessitating more
 205 advanced fine-tuning methods.

206 **Evaluation of Weight and Activation Quantization.** The evaluation results for existing methods
 207 using W4A8 quantization are presented in Table 3. The three light-weight optimization-based
 208 methods outperform RTN significantly, underscoring their efficacy. However, all of the results fall
 209 into either *Class-2* or *Class-3*. This suggests that for certain applications, it might be more beneficial
 210 to use smaller models with fewer parameters rather than larger, quantized models.

211 Among quantization-based methods, ZQ-Global* and ZQ-Local* generally outperform GPTQ, which
 212 is anticipated given that GPTQ was originally designed for weight-only quantization. ZQ-Global*
 213 performs better than ZQ-Local* in most cases except for the two largest models, OPT-66B and
 214 Bloom-176B, despite having larger trainable parameters in one step. This again signifies the need for
 215 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.

216

217 4.1 Fine-grained Quantization and Its Evaluation

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

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

228 Although fine-grained quantization can substantially narrow the gap between the quantized tensor
 229 and its floating-point counterpart, the application of RTN still results in a non-trivial accuracy gap.
 230 Consequently, we build upon fine-grained quantization by employing existing optimization-based
 231 methods to further enhance accuracy. Specifically, we utilize GPTQ and ZQ-Global for all models
 232 and settings and apply ZQ-Local to OPT-66B and Bloom-176B. For the hyperparameters used in
 233 ZQ-Global and ZQ-Local, we select the top three identified in Section 4 for all models, except for
 234 Bloom-176B, for which we only use the top-performing hyperparameter to reduce training costs.

235 **4-bit Weight Quantization.** We hereby present the W4A16 results for OPT and BLOOM, as
 236 delineated in Table 4, corresponding to an array of quantization block sizes. The performance
 237 sees a significant improvement with smaller block sizes compared to per-row quantization. The
 238 point of diminishing returns, however, varies for different model sizes. For example, smaller mod-
 239 els (such as OPT-6.7B and BLOOM-1.7b) continue to see substantial gains until the block size
 240 reduces to 32. In contrast, for larger models (those exceeding 10B, with OPT-66B as the excep-

Table 4: Results of **W4^{asym}-A16** quantization with various block-size out of the best result from optimization-based methods on OPT and BLOOM (BLM). See Table E.15 and Table E.16 for full results including RTN. 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 |

Table 5: 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). See Table E.20 for full results including RTN.

| 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 |

241 tion), the benefits derived from smaller block sizes wane rapidly around block-256/512. Most
 242 crucially, for models equal to or larger than 13B, a smaller quantization block size results in quanti-
 243 zation error being classified under *Class-1*, indicating virtually negligible degradation in accuracy.

244
 245 **Activation Quantization (W4A8).** To comprehend the benefits of fine-grained quantization on activation,
 246 we analyze the quantization between per-row and a block size of 128, with INT4 weight, as highlighted in
 247 Table 5. For models of considerable size, specifically those equal to or exceeding 1B, the application of such
 248 fine-grained activation quantization (Case-1) results in a substantial reduction in quantization error compared to per-row activation (Case-2). By implementing
 249 fine-grained activation quantization with weight quantization (Case-3), we are able to almost restore
 250 the performance to the level of their W4A16 counterparts.
 251

252 Furthermore, we detail the impacts of varying activation quantization block sizes in Table 6 on
 253 BLOOM-176B, with INT4 weight. A trend of superior accuracy is observed with smaller block
 254 sizes in contrast to larger ones. However, the enhancement in performance reaches a saturation point
 255 when the size smaller or equal to 256, which corresponds to the range of values INT8 can represent.
 256 Despite INT8’s capability to signify 256 distinct values, activation quantization errors persist due to
 257 the application of uniform quantization.
 258
 259
 260

Table 6: BLOOM-176B with different quantization block sizes on activation. Here weight is asymmetrically quantized with block size 128. See more in Table E.22.

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

261 **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.

262 5 Proposed Method to Further Push the Limit of Post-training Quantization

263 Building on the investigation and conclusions drawn from previous sections, it has become apparent
 264 that there is still a need for an advanced methodology to further refine the existing methods, with
 265 the objective of fully realizing the original fp16 PPL quality. In this section, we introduce a simple
 266 yet effective method called **LoRC** (Low Rank Compensation) to optimize the current existing
 267 quantization error and further bridge the gap between the quality of the original model and its
 268 quantized counterparts.

Table 7: 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 |

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 $\hat{W}_{\text{lorc}} = \hat{W} + \hat{E}$, thereby reducing quantization errors: $\|W - \hat{W}\| \geq \|W - \hat{W}_{\text{lorc}}\|$. LoRC consists of two steps:

Step I: Implement Singular Value Decomposition (SVD) on the error matrix $E = U\Sigma V$, where $U \in \mathbb{R}^{d_{\text{in}} \times d_{\text{in}}}$ and $V \in \mathbb{R}^{d_{\text{out}} \times d_{\text{out}}}$ are unitary matrices, and $\Sigma \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}}$ 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, $U_m = U_{:,1:m} \in \mathbb{R}^{d_{\text{in}} \times m}$, $V_m = V_{1:m,:} \in \mathbb{R}^{m \times d_{\text{out}}}$, 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 [32], 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 $18hml$.⁴ 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 $3m/2h \leq 0.016$.

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 FP16.⁵ 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 7, 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 8). Thus, the two low-rank matrices for LoRC in Table 7 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

Table 8: 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 | | | | Fine-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 |

⁴In 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 $4h \times h$.

⁵For INT8 Activation, please see Table E.23, the observation for FP16 holds similarly for INT8 Activation.

310 combination of fine-grained quantization with LoRC yields the most impressive results, underscoring
 311 the efficacy of LoRC when integrated with FGQ. Overall, the results emphasize the benefits of using
 312 LoRC for enhanced performance in weight quantization and its compatibility with FGQ. Notably,
 313 recovering the last 0.05-0.1 perplexity can be challenging, but with LoRC, we are able to nearly
 314 recover the original model quality for INT4 quantization.

Table 9: 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 |

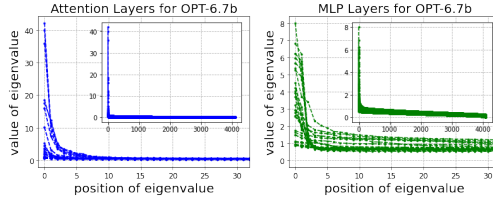


Figure 2: Eigenvalues of the Error matrix E for W4A16

315 **Ablation Study on the Low Rank Dimension m .** An essential aspect of the LoRC method is on the
 316 optimal low-rank dimension, denoted as m , explained in **Step II**. To explore this, we varied m in the
 317 range of 1, 4, 8, 16, and 32 for OPT-1.3b/6.7b/30b models, and applied W4A16 GPTQ quantization.
 318 The outcomes are depicted in Table 9, indicating that the enhancements achieved through LoRC
 319 begin to plateau as the dimension m surpasses 4. The most optimal performance for OPT-6.7b is
 320 realized when $m = 8$.

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

328 6 Discussion

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

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

342 **Future Opportunity.** Throughout the paper, we see several unresolved problems from current
 343 quantization schemes and/or algorithms, and we find potential directions for LLM compression: (1)
 344 Although we use fine-grained quantization schemes in the paper, the real implementation is missing.
 345 How to efficiently implement odd bit precision is also challenging. [12] demonstrated that 3-bit can
 346 achieve better throughput in the generation phase by packing all 3-bit numbers in continuous memory
 347 space. However, this method is sub-optimal as the dequantization step needs to connect bits from
 348 different bytes. One possible way to implement odd bits, e.g., 5 bits, is to use two integer matrices
 349 with INT4 and INT1. During the dequantization stage, we couple the two matrices together. (2) How
 350 to combine PTQ with other lightweight compression techniques, e.g., post-training pruning [20, 11],
 351 is an interesting direction to further reduce the memory consumption and compute cost.

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. For more details, please refer to Appendix D.

352 References

- 353 [1] Haoli Bai, Wei Zhang, Lu Hou, Lifeng Shang, Jing Jin, Xin Jiang, Qun Liu, Michael Lyu, and
354 Irwin King. Binarybert: Pushing the limit of bert quantization. *arXiv preprint arXiv:2012.15701*,
355 2020.
- 356 [2] Big-Science. Bloom inference. <https://github.com/huggingface/transformers-bloom-inference/tree/main/bloom-inference-scripts>, 2022.
- 358 [3] Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Understanding and overcoming
359 the challenges of efficient transformer quantization. *arXiv preprint arXiv:2109.12948*, 2021.
- 360 [4] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
361 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
362 few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 363 [5] Bitu Darvish Rouhani, Daniel Lo, Ritchie Zhao, Ming Liu, Jeremy Fowers, Kalin Ovtcharov,
364 Anna Vinogradsky, Sarah Massengill, Lita Yang, Ray Bittner, et al. Pushing the limits of
365 narrow precision inferencing at cloud scale with microsoft floating point. *Advances in neural
366 information processing systems*, 33:10271–10281, 2020.
- 367 [6] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm. int8 (): 8-bit matrix
368 multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- 369 [7] Tim Dettmers and Luke Zettlemoyer. The case for 4-bit precision: k-bit inference scaling laws.
370 *arXiv preprint arXiv:2212.09720*, 2022.
- 371 [8] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dhar-
372 mendra S Modha. Learned step size quantization. *arXiv preprint arXiv:1902.08153*, 2019.
- 373 [9] Angela Fan, Pierre Stock, Benjamin Graham, Edouard Grave, Remi Gribonval, Herve Jegou,
374 and Armand Joulin. Training with quantization noise for extreme fixed-point compression.
375 *arXiv preprint arXiv:2004.07320*, 2020.
- 376 [10] Elias Frantar and Dan Alistarh. Optimal brain compression: A framework for accurate post-
377 training quantization and pruning. *arXiv preprint arXiv:2208.11580*, 2022.
- 378 [11] Elias Frantar and Dan Alistarh. Massive language models can be accurately pruned in one-shot.
379 *arXiv preprint arXiv:2301.00774*, 2023.
- 380 [12] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training
381 quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- 382 [13] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer.
383 A survey of quantization methods for efficient neural network inference. *arXiv preprint
384 arXiv:2103.13630*, 2021.
- 385 [14] Amir Gholami, Zhewei Yao, Sehoon Kim, Michael W Mahoney, and Kurt Keutzer. Ai and
386 memory wall. *RiseLab Medium Post*, 2021.
- 387 [15] GitHub. Github copilot. <https://github.com/features/copilot/>, 2021.
- 388 [16] Babak Hassibi and David G Stork. Second order derivatives for network pruning: Optimal brain
389 surgeon. In *Advances in neural information processing systems*, pages 164–171, 1993.
- 390 [17] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
391 *Workshop paper in NIPS*, 2014.
- 392 [18] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and
393 Qun Liu. Tinybert: Distilling bert for natural language understanding. *arXiv preprint
394 arXiv:1909.10351*, 2019.
- 395 [19] Sehoon Kim, Amir Gholami, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. I-bert:
396 Integer-only bert quantization. In *International conference on machine learning*, pages 5506–
397 5518. PMLR, 2021.

- 398 [20] Woosuk Kwon, Sehoon Kim, Michael W Mahoney, Joseph Hassoun, Kurt Keutzer, and
399 Amir Gholami. A fast post-training pruning framework for transformers. *arXiv preprint*
400 *arXiv:2204.09656*, 2022.
- 401 [21] Yann LeCun, John S Denker, and Sara A Solla. Optimal brain damage. In *Advances in neural*
402 *information processing systems*, pages 598–605, 1990.
- 403 [22] Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank.
404 *Using Large Corpora*, page 273, 1994.
- 405 [23] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
406 models. In *International Conference on Learning Representations*, 2017.
- 407 [24] OpenAI. Openai chatgpt. <https://openai.com/blog/chatgpt/>, 2022.
- 408 [25] Antonio Polino, Razvan Pascanu, and Dan Alistarh. Model compression via distillation and
409 quantization. *arXiv preprint arXiv:1802.05668*, 2018.
- 410 [26] Reiner Pope, Sholto Douglas, Aakanksha Chowdhery, Jacob Devlin, James Bradbury, Anselm
411 Levskaya, Jonathan Heek, Kefan Xiao, Shivani Agrawal, and Jeff Dean. Efficiently scaling
412 transformer inference. *arXiv preprint arXiv:2211.05102*, 2022.
- 413 [27] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
414 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified
415 text-to-text transformer, 2019.
- 416 [28] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow,
417 Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A
418 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*,
419 2022.
- 420 [29] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W
421 Mahoney, and Kurt Keutzer. Q-BERT: Hessian based ultra low precision quantization of bert.
422 In *AAAI*, pages 8815–8821, 2020.
- 423 [30] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari,
424 Jared Casper, Zhun Liu, Shrimai Prabhunoye, George Zerveas, Vijay Korthikanti, et al. Using
425 deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language
426 model. *arXiv preprint arXiv:2201.11990*, 2022.
- 427 [31] Chaofan Tao, Lu Hou, Wei Zhang, Lifeng Shang, Xin Jiang, Qun Liu, Ping Luo, and Ngai
428 Wong. Compression of generative pre-trained language models via quantization. *arXiv preprint*
429 *arXiv:2203.10705*, 2022.
- 430 [32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
431 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information*
432 *processing systems*, pages 5998–6008, 2017.
- 433 [33] Xiaoxia Wu, Cheng Li, Reza Yazdani Aminabadi, Zhewei Yao, and Yuxiong He. Understanding
434 int4 quantization for transformer models: Latency speedup, composability, and failure cases.
435 *arXiv preprint arXiv:2301.12017*, 2023.
- 436 [34] Xiaoxia Wu, Zhewei Yao, Minjia Zhang, Conglong Li, and Yuxiong He. Extreme compression
437 for pre-trained transformers made simple and efficient. *arXiv preprint arXiv:2206.01859*, 2022.
- 438 [35] Guangxuan Xiao, Ji Lin, Mickael Seznec, Julien Demouth, and Song Han. Smoothquant:
439 Accurate and efficient post-training quantization for large language models. *arXiv preprint*
440 *arXiv:2211.10438*, 2022.
- 441 [36] Zhewei Yao, Reza Yazdani Aminabadi, Minjia Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong
442 He. Zeroquant: Efficient and affordable post-training quantization for large-scale transformers.
443 *arXiv preprint arXiv:2206.01861*, 2022.

- 444 [37] Ali Hadi Zadeh, Isak Edo, Omar Mohamed Awad, and Andreas Moshovos. Gobo: Quantizing
445 attention-based nlp models for low latency and energy efficient inference. In *2020 53rd Annual*
446 *IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 811–824. IEEE,
447 2020.
- 448 [38] Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8BERT: Quantized 8bit bert.
449 *arXiv preprint arXiv:1910.06188*, 2019.
- 450 [39] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang,
451 Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv*
452 *preprint arXiv:2210.02414*, 2022.
- 453 [40] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen,
454 Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained
455 transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- 456 [41] Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao Chen, Xin Jiang, and Qun Liu. Ternarybert:
457 Distillation-aware ultra-low bit bert. *arXiv preprint arXiv:2009.12812*, 2020.