
LoftQ: LoRA-Fine-Tuning-Aware Quantization for Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Quantization is an indispensable technique for serving Large Language Models
2 (LLMs) and has recently found its way into LoRA fine-tuning (Dettmers et al.,
3 2023). In this work we focus on the scenario where quantization and LoRA fine-
4 tuning are applied together on a pre-trained model. In such cases it is common
5 to observe a consistent gap in the performance on downstream tasks between full
6 fine-tuning and quantization plus LoRA fine-tuning approach. In response, we
7 propose LoftQ (**LoRA-Fine-Tuning-aware Quantization**), a novel quantization
8 framework that simultaneously quantizes an LLM and finds a proper low-rank
9 initialization for LoRA fine-tuning. Such an initialization alleviates the discrep-
10 ancy between the quantized and full-precision model and significantly improves
11 generalization in downstream tasks. We evaluate our method on natural language
12 understanding, question answering, summarization, and natural language genera-
13 tion tasks. Experiments show that our method is highly effective and outperforms
14 existing quantization methods, especially in the challenging 2-bit and 2/4-bit mixed
15 precision regimes. We will release our code.

16 1 Introduction

17 The advent of Pre-trained Language Models (PLMs) has marked a transformative shift in the field
18 of Natural Language Processing (NLP), offering versatile solutions across various applications (He
19 et al., 2021b; Lewis et al., 2019; Touvron et al., 2023). However, the extensive computational and
20 memory demands of these models pose significant challenges, especially in real-world deployments
21 where resources are often constrained and need to be shared among many users.

22 To mitigate the extensive storage requirements of pre-trained models, quantization serves as a pivotal
23 compression technique (Zafir et al., 2019; Shen et al., 2020; Bai et al., 2022; Dettmers et al.,
24 2022), converting high-precision numerical values into a discrete set of values. Typically, model
25 parameters, originally stored in a 16-bit float format, are transformed into a 4-bit integer format
26 through quantization, resulting in a substantial 75% reduction in storage overhead.

27 When quantizing pre-trained models, however, practitioners often concentrate primarily on the
28 quantization technique, inadvertently neglecting the importance of subsequent task adaption, e.g.,
29 LoRA fine-tuning (Dettmers et al., 2023; Diao et al., 2023). For example, QLoRA inherits the fixup
30 initialization (Zhang et al., 2019) used in LoRA, which (Dettmers et al., 2023) attaches zero initialized
31 low-rank adapters to the quantized pre-trained model. The inevitable discrepancy introduced by
32 quantization during the approximation of the original high-precision numbers, a scenario particularly
33 pronounced in low-bit situations such as the 2-bit regime, can adversely impact the initialization of
34 LoRA fine-tuning. This deviation often results in an inferior fine-tuning performance.

35 In this paper, we introduce a novel quantization framework, called **LoRA-Fine-Tuning-aware**
36 **Quantization (LoftQ)**. It is designed specifically for pre-trained models that require quantization

37 and LoRA fine-tuning. This framework actively integrates low-rank approximation, working in
 38 tandem with quantization to jointly approximate the original high-precision pre-trained weights. This
 39 synergy significantly enhances alignment with the original pre-trained weights. Consequently, our
 40 method provides an advantageous initialization point for subsequent LoRA fine-tuning, leading to
 41 improvements in downstream tasks.

42 We evaluate our quantization framework by conducting extensive experiments on downstream tasks.
 43 Experiments show that LoftQ consistently outperforms QLoRA across all precision levels. LoftQ is
 44 also competitive with different quantization methods, e.g., NormalFloat (Dettmers et al., 2023) and
 45 the uniform quantization.

46 2 Method

47 We propose **LoRA-Fine-Tuning-aware Quantization (LoftQ)**, a quantization framework for LLMs. It
 48 alternatively applies quantization and low-rank approximation to approximate original pre-trained
 49 weights. This quantization framework provides a promising initialization for LoRA fine-tuning,
 50 which alleviates the quantization discrepancy in QLoRA and improves generalization in downstream
 51 tasks significantly.

52 2.1 LoRA-Aware Quantization

53 We use an N -bit quantized weight $Q \in \mathbb{R}_N^{d_1 \times d_2}$ and low-rank approximations $A \in \mathbb{R}^{d_1 \times r}, B \in$
 54 $\mathbb{R}^{d_2 \times r}$ to approximate the original high-precision pre-trained weight $W \in \mathbb{R}^{d_1 \times d_2}$ as the initialization
 55 of LoRA fine-tuning. Specifically, before fine-tuning, we initialize the network by minimizing the
 56 following objective:

$$\min_{Q,A,B} \|W - Q - AB^\top\|_F, \quad (1)$$

57 where $\|\cdot\|_F$ denotes the Frobenious norm. This objective in (1) takes LoRA fine-tuning into consid-
 58 eration by jointly optimizing the initial values of the quantized backbone Q and low-rank adapters
 59 A, B . Contrarily, practitioners typically convert the pre-trained weight W into a quantized weight Q
 60 outright and initialize the low-rank adapters by $A \sim \mathcal{N}(0, \sigma^2), B = 0$, neglecting the subsequent
 61 LoRA fine-tuning process. This oversight leads to notable performance degradation in downstream
 62 tasks arising from the quantization discrepancy.

63 2.2 Alternating Optimization

64 We solve the minimization problem in (1) by alternating between quantization and singular value
 65 decomposition (SVD). Initially, we set A_0 , and B_0 equal to 0. The total alternating steps are T .

66 **Quantization.** At the t -th step, we quantize the difference between the original pre-trained weight W
 67 and the low-rank approximation $A_{t-1}B_{t-1}^\top$ from the last step to obtain the quantized weight Q_t by

$$Q_t = q_N(W - A_{t-1}B_{t-1}^\top), \quad (2)$$

68 where $q_N(\cdot): \mathbb{R}^{m \times n} \mapsto \mathbb{R}_N^{m \times n}$ is a quantization function that maps a high-precision matrix into a
 69 quantized matrix, and $\mathbb{R}_N: \{\mathcal{T}[i] \in \mathbb{R} | 0 \leq i < 2^N\}$.

70 We remark that our algorithm is compatible with different quantization functions $q_N(\cdot)$. We apply
 71 NF4 and the uniform quantization in Section 3 as examples. We also remark that Q_t is not an exact
 72 solution of the minimization in (1), given the fixed $A_{t-1}B_{t-1}^\top$, but it is an efficient approximation.

73 **SVD.** After obtaining the t -th quantized weight Q_t , SVD is applied to the residual of the quantization
 74 denoted by $R_t = W - Q_t$ by

$$R_t = \sum_{i=1}^d \sigma_{t,i} u_{t,i} v_{t,i}^\top, \quad (3)$$

75 where $d = \min\{d_1, d_2\}$, $\sigma_{t,1} \geq \sigma_{t,2} \geq \dots \geq \sigma_{t,d}$ are the singular values of R_t , $u_{t,i}$'s and $v_{t,i}$'s are
 76 the associated left and right singular vectors of R_t . We then obtain a rank- r approximation of R_t by
 77 $A_t B_t^\top$, where

$$\begin{aligned} A_t &= [\sqrt{\sigma_{t,1}} u_{t,1}, \dots, \sqrt{\sigma_{t,r}} u_{t,r}], \\ B_t &= [\sqrt{\sigma_{t,1}} v_{t,1}, \dots, \sqrt{\sigma_{t,r}} v_{t,r}]. \end{aligned} \quad (4)$$

78 It is worth noting that $T = 1$ is a special case where Q_1 is the exact quantized weight obtained
 79 by QLoRA, and low-rank approximations A_1, B_1 are obtained by the SVD of the quantization

80 residual $W - Q_1$. $T = 1$ is sufficient to mitigate the quantization discrepancy, and alternating
 81 optimization helps to find a closer initialization to the pre-trained weight W , which further improves
 82 the performance.

83 We remark that the computational cost of LoftQ is negligible because it is applied to individual
 84 weight matrices and therefore can be executed in parallel. Also, it only requires being applied once to
 85 a pre-trained model, and the initialization obtained can be reused for various downstream tasks

86 We then use Q_T, A_T, B_T obtained by the alternating optimization above as LoRA fine-tuning
 87 initialization. We freeze the quantized backbone Q_T and optimize the low-rank adapters, starting
 88 from A_T, B_T , with an efficient optimization algorithm, e.g., AdamW (Loshchilov & Hutter, 2017).

89 3 Experiments

90 We evaluate our method on NLU and NLG tasks. We apply LoftQ for quantizing DeBERTaV3-base
 91 (He et al., 2021b), BART-large (Lewis et al., 2019), and LLAMA-2 series (Touvron et al., 2023).

92 **Implementation Details.** Following the prior works of LoRA variants (Zhang et al., 2023; He et al.,
 93 2021a), we freeze all the backbone weight matrices and add low-rank adapters to weight matrices in
 94 MHA and FFN of all layers. We quantize the weight matrices that are attached by low-rank adapters.

95 **Quantization Methods.** We apply two quantization methods to demonstrate LoftQ is compatible
 96 with different quantization functions. We apply uniform quantization and NF4 quantization¹ in our
 97 experiments. We perform 2-bit and 4-bit quantization on all models, achieving compression ratios of
 98 25-30% and 15-20% at the 4-bit and 2-bit levels, respectively. The compression ratios and trainable
 99 parameter ratios for all models are detailed in the Appendix A.

100 **Baselines.** We compare LoftQ with baseline methods of Full fine-tuning, Full precision LoRA
 101 (LoRA), and QLoRA. Specific introductions of baseline methods are given in Appendix B

102 3.1 Encoder-only Model: DeBERTaV3

103 **Models, Datasets, Implementations, and Results** We quantize the DeBERTaV3-base (He et al.,
 104 2021b) with LoftQ, then finetune the model on the General Language Understanding Evaluation
 105 (GLUE) benchmark (Wang et al., 2019), SQuADv1.1 (Rajpurkar et al., 2016), and ANLI (Nie et al.,
 106 2019). We show the implementation details in Appendix E.3. Table 1 summarize the results for
 107 2-bit quantization on the GLUE, SQuADv1.1, and ANLI datasets, by NF2 quantization. Our method
 108 consistently outperforms QLoRA on all settings with respect to different ranks, and datasets. The
 109 4-bit quantization experiment results are presented in Appendix E.1 as both LoftQ and QLoRA
 110 achieve performance close to full fine-tuning. We also show our method excels over baseline using
 111 uniform quantization shown in Table 8, indicating our method is applicable to different methods.

Table 1: Results with 2-bit LoftQ of DeBERTaV3-base models on GLUE development set, SQuADv1.1 development set, ANLI test set using **NF2 quantization**. *N.A.* indicates the model does not converge. The best results on each dataset are shown in **bold**.

Rank	Method	MNLI m / mm	QNLI Acc	RTE Acc	SST Acc	MRPC Acc	CoLA Matt	QQP Acc	STS P/S Corr	SQuAD EM/F1	ANLI Acc
	Full FT	90.4/90.5	94.6	85.1	95.1	89.9/93.6	69.9	92.0/89.4	91.7/91.1	87.3/93.1	59.8
16	LoRA	90.5/90.6	94.8	85.2	95.0	89.9/93.6	69.8	92.0/89.4	91.6/91.0	87.0/93.1	60.2
16	QLoRA	75.4/75.6	82.4	55.9	86.5	73.8/82.8	N.A.	86.8/82.3	83.0/82.8	61.5 / 71.2	N.A.
	LoftQ	84.7/85.1	86.6	61.4	90.2	83.8/88.6	37.4	90.3/86.9	87.1/86.9	81.5/88.6	47.1
32	QLoRA	78.5/78.7	80.4	56.7	86.9	73.8/82.7	N.A.	87.1/82.7	83.6/83.3	64.6/73.8	N.A.
	LoftQ	86.0/86.1	89.9	61.7	92.0	83.6/87.2	47.5	91.0/87.9	87.5/87.0	82.9/89.8	49.0

112 3.2 Encoder-Decoder Model: BART

113 **Models, Datasets, Implementations, and Results** We quantize BART-large model (Lewis et al.,
 114 2020) with LoftQ, then finetune and evaluate the model on two commonly used summarization
 115 datasets: XSum (Narayan et al., 2018) and CNN/DailyMail(Hermann et al., 2015). The implemen-
 116 tation details are given in Appendix F. The 2-bit quantization results are shown in Table 2. Our
 117 observation is consistent with the NLU experiments, that LoftQ demonstrates the convergence to
 118 reasonable results, while QLoRA does not converge. This indicates our method is robust by nar-

¹We abbreviate NF4 for NormalFloat quantization used in QLoRA. NF2 is its 2-bit variant

119 rowing the initialization gap. We remark that our method is also successful within 4-bit quantization
 120 precision. Results shown in Table 13

Table 2: Results with 2-bit LoftQ of BART-large on XSum and CNN/DailyMail using **NF2 quantization**. *N.A.* indicates the model does not converge. We report ROUGE-1/2/L.

Rank	Method	XSum	CNN/DailyMail
-	Lead-3	16.30/1.60/11.95	40.42/17.62/36.67
	Full FT	45.14/22.27/37.25	44.16/21.28/40.90
8	FP LoRA	43.40/20.20/35.20	44.72/21.58/41.84
	QLoRA	N.A.	N.A.
	LoftQ	39.63/16.65/31.62	42.24/19.44/29.04
16	FP LoRA	43.95/20.72/35.68	45.03/21.84/42.15
	QLoRA	N.A.	N.A.
	LoftQ	40.81/17.85/32.80	42.52/19.81/39.51

121 3.3 Decoder-only Model: LLAMA-2

122 **Models, Datasets, Implementations, and Results** We quantize LLAMA-2-7b and LLAMA-2-13b
 123 (Touvron et al., 2023) with LoftQ. We then fine-tune and evaluate the models on two NLG datasets:
 124 GSM8K (Cobbe et al., 2021) and WikiText-2 (Merity et al., 2016). Please see Appendix G for more
 125 details about the datasets. The implementation details are given in Appendix G. Table 3 presents a
 126 summary of our experiments on LLAMA-2-7b and LLAMA-2-13b using 2-bit, 4-bit, and mixed-
 127 precision NormalFloat quantization methods on WikiText-2 and GSM8K datasets. In WikiText-2, our
 128 method consistently outperforms QLoRA across all quantization precision settings on both models.
 129 When dealing with the challenging 2-bit precision, where QLoRA fails to converge, LoftQ manages
 130 to achieve a perplexity of 7.85. In GSM8K, our method also achieves better or on par performance
 131 compared to QLoRA across different model sizes and quantization precision levels.

132 We also explore mixed-precision quantization where matrices in the first 4 layers are quantized
 133 using 4 bits, and the rest matrices remain 2 bits. We witness a remarkable 5.9% accuracy boost
 134 on the GSM8K dataset using LLAMA-2-7b and a 12.7% boost using LLAMA-2-13b. This result
 135 underscores the potential of LoftQ for complex mixed-precision quantization scenarios.

Table 3: Results of LoftQ using 4-bit, 2.25 bit and 2-bit for LLAMA-2 series on WikiText-2 and GSM8K. 2.25-bit indicates mixed-precision NormalFloat quantization: 4-bit precision for the first 4 layers and 2-bit precision for the rest of layers. We report the perplexity (the smaller the better) for WikiText-2 and accuracy for GSM8K. The rank of low-rank adapters is 64. *N.A.* indicates the model does not converge.

Method	Bit	LLAMA-2-7b		LLAMA-2-13b	
		WikiText-2↓	GSM8K↑	WikiText-2↓	GSM8K↑
LoRA	16	5.08	36.9	5.12	43.1
QLoRA	4	7.41	35.1	5.22	39.9
LoftQ	4	5.24	35.0	5.16	45.0
QLoRA	2.25	N.A.	N.A.	N.A.	N.A.
LoftQ	2.25	6.13	26.5	5.45	38.1
QLoRA	2	N.A.	N.A.	N.A.	N.A.
LoftQ	2	7.85	20.9	7.69	25.4

136 4 Conclusion

137 We propose LoftQ, a quantization framework for LLMs, which alternatively applies quantization
 138 and low-rank approximation to the original high-precision pre-trained weights, to obtain an ini-
 139 tialization for the subsequent LoRA fine-tuning. Experiments on natural language understanding,
 140 question answering, summarization, and natural language generation show that our framework remark-
 141 ably surpasses existing methods, e.g., QLoRA, for quantizing encoder-only, encoder-decoder, and
 142 decoder-only models. We have not observed our method exhibiting worse performance over QLoRA.
 143 Moreover, our quantization framework demonstrates effectiveness and robustness particularly in
 144 low-bit quantization regimes, e.g., the 2-bit level.

145 References

- 146 Haoli Bai, Lu Hou, Lifeng Shang, Xin Jiang, Irwin King, and Michael R Lyu. Towards efficient post-
147 training quantization of pre-trained language models. *Advances in Neural Information Processing*
148 *Systems*, 35:1405–1418, 2022.
- 149 Roy Bar-Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan
150 Szpektor. The second pascal recognising textual entailment challenge. 2006.
- 151 Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth pascal recognizing
152 textual entailment challenge. In *TAC*, 2009.
- 153 Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. SemEval-2017 task 1:
154 Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of*
155 *the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pp. 1–14, Vancouver,
156 Canada, August 2017. Association for Computational Linguistics. doi: 10.18653/v1/S17-2001.
- 157 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
158 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
159 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 160 Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment
161 challenge. In *Machine Learning Challenges Workshop*, 2007.
- 162 Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix
163 multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- 164 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning
165 of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- 166 Shizhe Diao, Rui Pan, Hanze Dong, Ka Shun Shum, Jipeng Zhang, Wei Xiong, and Tong Zhang.
167 Lmflow: An extensible toolkit for finetuning and inference of large foundation models. *arXiv*
168 *preprint arXiv:2306.12420*, 2023.
- 169 William B. Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases.
170 In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005.
- 171 Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third PASCAL recognizing
172 textual entailment challenge. In *Proceedings of the ACL-PASCAL Workshop on Textual Entailment*
173 *and Paraphrasing*, pp. 1–9, Prague, June 2007. Association for Computational Linguistics.
- 174 Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a
175 unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2021a.
- 176 Pengcheng He, Jianfeng Gao, and Weizhu Chen. Debertav3: Improving deberta using electra-style
177 pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*,
178 2021b.
- 179 Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa
180 Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. *Advances in neural*
181 *information processing systems*, 28, 2015.
- 182 Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In
183 *Thirteenth international conference on the principles of knowledge representation and reasoning*,
184 2012.
- 185 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer
186 Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for
187 natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*,
188 2019.
- 189 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
190 Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence pre-training for
191 natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual*
192 *Meeting of the Association for Computational Linguistics*, pp. 7871–7880, Online, July 2020.
193 Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.703.

- 194 Yixiao Li, Yifan Yu, Qingru Zhang, Chen Liang, Pengcheng He, Weizhu Chen, and Tuo Zhao.
195 Lospars: Structured compression of large language models based on low-rank and sparse approxi-
196 mation. *arXiv preprint arXiv:2306.11222*, 2023.
- 197 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
198 *arXiv:1711.05101*, 2017.
- 199 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
200 models, 2016.
- 201 Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don’t give me the details, just the summary!
202 topic-aware convolutional neural networks for extreme summarization. *ArXiv*, abs/1808.08745,
203 2018.
- 204 Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial
205 nli: A new benchmark for natural language understanding. *ArXiv*, abs/1910.14599, 2019. URL
206 <https://api.semanticscholar.org/CorpusID:207756753>.
- 207 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
208 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward
209 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,
210 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep
211 learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran
212 Associates, Inc., 2019.
- 213 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions
214 for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods*
215 *in Natural Language Processing*, pp. 2383–2392, Austin, Texas, November 2016. Association for
216 Computational Linguistics. doi: 10.18653/v1/D16-1264.
- 217 Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney,
218 and Kurt Keutzer. Q-bert: Hessian based ultra low precision quantization of bert. In *Proceedings*
219 *of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 8815–8821, 2020.
- 220 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and
221 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank.
222 In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp.
223 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.
- 224 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
225 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
226 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 227 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.
228 GLUE: A multi-task benchmark and analysis platform for natural language understanding. In
229 *International Conference on Learning Representations*, 2019.
- 230 Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. Neural network acceptability judgments.
231 *Transactions of the Association for Computational Linguistics*, 7:625–641, 2019. doi: 10.1162/
232 tacl_a_00290.
- 233 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for
234 sentence understanding through inference. In *Proceedings of the 2018 Conference of the North*
235 *American Chapter of the Association for Computational Linguistics: Human Language Technolo-*
236 *gies, Volume 1 (Long Papers)*, pp. 1112–1122, New Orleans, Louisiana, June 2018. Association
237 for Computational Linguistics. doi: 10.18653/v1/N18-1101.
- 238 Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8bert: Quantized 8bit bert. In *2019*
239 *Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS Edition*
240 *(EMC2-NIPS)*, pp. 36–39. IEEE, 2019.
- 241 Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without
242 normalization. *arXiv preprint arXiv:1901.09321*, 2019.

243 Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen,
244 and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. *arXiv preprint*
245 *arXiv:2303.10512*, 2023.

246 **A Model Compression Ratio**

Table 4: Compression ratios of backbones.

Model	Compression ratio (%)	Trainable ratio (%)	Rank	Bits	Quantization method
DeBERTaV3-base	15.6	3.1	16	2	uniform
DeBERTaV3-base	18.8	6.3	32	2	uniform
DeBERTaV3-base	17.2	3.1	16	2	nf2
DeBERTaV3-base	20.4	6.3	32	2	nf2
BART-large	15.3	1.2	8	2	nf2
BART-large	16.7	2.5	16	2	nf2
BART-large	27.8	1.2	8	4	nf4
BART-large	29.0	2.5	16	4	nf4
BART-large	26.2	1.2	8	4	uniform
BART-large	27.5	2.5	16	4	uniform
LLAMA-2-7b	16.6	2.4	64	2	nf2
LLAMA-2-7b	29.0	2.4	64	4	nf4
LLAMA-2-13b	16.0	1.9	64	2	nf2
LLAMA-2-13b	28.5	1.9	64	4	nf4

247 **B Baseline Methods**

- 248 • *Full fine-tuning* is the most common approach for adapting a pre-trained model to down-
249 stream tasks. The model is initialized with pre-trained weights and all parameters are
250 updated through an SGD-type optimization method.
- 251 • *Full precision LoRA (LoRA)* is a lightweight method for task adaptation, where it stores the
252 backbone using 16-bit numbers and optimizes the low-rank adaptors only. The adaptors are
253 applied to the same matrices as in LoftQ.
- 254 • *QLoRA* is similar to *LoRA* except the backbone is quantized into low-bit regime. The
255 low-rank adapters are zero initialized using and are applied to the same matrices as in LoftQ.

256 **C GLUE Dataset Statistics**

We present the dataset statistics of GLUE Wang et al. (2019) in the following table.

Corpus	Task	#Train	#Dev	#Test	#Label	Metrics
Single-Sentence Classification (GLUE)						
CoLA	Acceptability	8.5k	1k	1k	2	Matthews corr
SST	Sentiment	67k	872	1.8k	2	Accuracy
Pairwise Text Classification (GLUE)						
MNLI	NLI	393k	20k	20k	3	Accuracy
RTE	NLI	2.5k	276	3k	2	Accuracy
QQP	Paraphrase	364k	40k	391k	2	Accuracy/F1
MRPC	Paraphrase	3.7k	408	1.7k	2	Accuracy/F1
QNLI	QA/NLI	108k	5.7k	5.7k	2	Accuracy
Text Similarity (GLUE)						
STS-B	Similarity	7k	1.5k	1.4k	1	Pearson/Spearman corr

Table 5: Summary of the GLUE benchmark.

257
258 GLUE includes two single-sentence classification tasks: SST-2 (Socher et al., 2013) and CoLA
259 (Warstadt et al., 2019), and three similarity and paraphrase tasks: MRPC (Dolan & Brockett, 2005),
260 STS-B (Cer et al., 2017), and QQP. GLUE also includes four natural language inference tasks in
261 GLUE: MNLI (Williams et al., 2018), QNLI (Rajpurkar et al., 2016), RTE (Dagan et al., 2007;
262 Bar-Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009), and WNLI (Levesque et al.,
263 2012).

264 **D Decompose Time**

265 We report the execution time of LoftQ applying to a single weight matrix. The time is tested on
266 Intel(R) Xeon(R) CPU E5-2650 v4 @ 2.20GHz.

Table 6: Execution time of LoftQ applying to different weight matrices.

Model	Size	Step T	Quantization method	Time
DeBERTaV3-base	768*768	5	Uniform	1s
BART-large	1024*1024	5	NF4	1s
LLAMA-2-7b	4096*4096	5	NF4	21s
LLAMA-2-13b	5120*5120	5	NF4	43s

267 E Natural Language Understanding

268 E.1 GLUE with 4-bit

269 We show the 4-bits results in the Table 7. Both methods can achieve performance close to full-finetuning.

Table 7: Results with 4-bit LoftQ of DeBERTaV3-base models on GLUE development set using NF4 quantization. We report the median over four seeds. Results with N.A. indicate the model does not converge. The best results on each dataset are shown in bold

Method	Rank	MNLI m / mm	SST-2 Acc	QNLI Acc	ANLI Acc
Full FT	-	90.5/90.6	95.3	94.0	59.8
QLoRA	32	89.9/89.9	95.3	94.2	59.4
LoftQ	32	89.9/90.0	95.3	94.1	59.9

270

271 E.2 GLUE with uniform 2-bit quantization

Table 8: Results with 2-bit LoftQ of DeBERTaV3-base models on GLUE development set, SQuADv1.1 development set using **Uniform quantization**. We report the median over four seeds. N.A. indicates the model does not converge. The best results on each task are shown in **bold**.

Rank	Method	MNLI m / mm	QNLI Acc	RTE Acc	SST Acc	MRPC Acc	CoLA Acc	QQP Mcc	STSB P/S Corr	SQuAD Em/F1
-	Full FT	90.5/90.6	94.0	82.0	95.3	89.5/93.3	69.2	92.4/89.8	91.6/91.1	88.5/92.8
16	LoRA	90.4/90.5	94.6	85.1	95.1	89.9/93.6	69.9	92.0/89.4	91.7/91.1	87.3/93.1
16	QLoRA	76.5/76.3	83.8	56.7	86.6	75.7/84.7	N.A.	87.1/82.6	83.5/83.4	69.5/77.6
	LoftQ	87.3/87.1	90.6	61.1	94.0	87.0/90.6	59.1	90.9/88.0	87.9/87.6	84.4/91.2
32	QLoRA	79.9/79.5	83.7	57.8	86.9	76.5/84.5	N.A.	88.6/84.7	84.1/84.0	71.6/80.2
	LoftQ	88.0/88.1	92.2	63.2	94.7	87.5/91.2	60.5	91.3/88.3	89.5/89.2	85.2/91.6

272 E.3 Training Details

273 **Implementation Details.** The implementation of LoftQ is based on publicly available Huggingface
274 (Paszke et al., 2019) code-base ².

275 **Hyper-parameter Details.** We select the learning rate of $\{1 \times 10^{-5}, 5 \times 10^{-5}, 1 \times 10^{-4}, 5 \times 10^{-4}\}$,
276 and use the selected learning rate for both uniform quantization experiments and nf2 quantization
277 experiments. We use batch size of 32 for all GLUE tasks and ANLI. We use batch size of 16 for
278 SQuADv1.1. We use LoftQ of 5 iterations for all GLUE tasks.

279 Table 9 summarizes the detailed hyperparameters for each task used in training DeBERTaV3-base
280 using uniform quantization. Table 10 summarizes the detailed hyperparameters for each task used in
281 training DeBERTaV3-base using nf2 quantization.

²<https://github.com/huggingface/transformers/tree/main/examples/pytorch>

Table 9: Hyper-parameter setup of LoftQ for GLUE benchmark for training DeBERTaV3-base using uniform quantization.

Hyper-parameter	MNLI	RTE	QNLI	MRPC	QQP	SST-2	CoLA	STS-B	SQuADv1.1	ANLI
# epochs	5	20	10	60	10	10	60	60	10	12
Learning rate	1×10^{-4}	5×10^{-4}	5×10^{-5}	1×10^{-4}	5×10^{-5}	5×10^{-5}	5×10^{-5}	5×10^{-5}	5×10^{-5}	5×10^{-5}

Table 10: Hyper-parameter setup of LoftQ for GLUE benchmark for training DeBERTaV3-base using nf2 quantization.

Hyper-parameter	MNLI	RTE	QNLI	MRPC	QQP	SST-2	CoLA	STS-B	SQuADv1.1	ANLI
# epochs	5	20	10	60	10	10	60	60	10	12
Learning rate	1×10^{-4}	5×10^{-5}	5×10^{-5}	1×10^{-4}	5×10^{-5}	5×10^{-5}	5×10^{-5}	1×10^{-4}	5×10^{-5}	5×10^{-5}

282 F Summarization

283 F.1 Training Details

284 We set the batch size as 32, the number of training epoch as 10. We choose Adam as the optimizer
 285 and try learning rate from $\{1 \times 10^{-5}, 5 \times 10^{-5}, 7 \times 10^{-5}, 2 \times 10^{-4}, 3 \times 10^{-4}, 4 \times 10^{-4}\}$. We show
 286 the optimal learning rate for different settings in Table We use LoftQ of 1 iteration for all BART-large
 287 experiments. Table 11 and Table 12 summarize the learning rate for CNN/DailyMail and XSum

Table 11: Hyper-parameter setup of LoftQ BART-large on CNN/DailyMail

Hyperparameter	NF4		4-bit Uniform		NF2	
	rank8	rank16	rank8	rank16	rank8	rank16
Learning rate	2e-4	2e-4	2e-4	3e-4	2e-4	2e-4

Table 12: Hyper-parameter setup of LoftQ BART-large on XSum

Hyperparameter	NF4		4-bit Uniform		NF2	
	rank8	rank16	rank8	rank16	rank8	rank16
Learning rate	2e-4	2e-4	2e-4	2e-4	2e-4	2e-4

288 F.2 BART-large experiments with NF4 quantization

Table 13: Results with 4-bit LoftQ of BART-large on XSum and CNN/DailyMail. We report ROUGE-1/2/L, the higher the better. *Lead-3* means choosing the first 3 sentences as the summary. *N.A.* indicates the model does not converge. *Full FT* refers to the full fine-tuning where all parameters are tuned. We report the median over five seeds.

Quantization	Rank	Method	XSum	CNN/DailyMail
-	-	Lead-3	16.30/1.60/11.95	40.42/17.62/36.67
		Full FT	45.14/22.27/37.25	44.16/21.28/40.90
	8	LoRA	43.40/20.20/35.20	44.72/21.58/41.84
		LoRA	43.95/20.72/35.68	45.03/21.84/42.15
NF4	8	QLoRA	42.91/19.72/34.82	43.10/20.22/40.06
		LoftQ	44.08/20.72/35.89	43.81/20.95/40.84
	16	QLoRA	43.29/20.05/35.15	43.42/20.62/40.44
		LoftQ	44.51/21.14/36.18	43.96/21.06/40.96

289 G Natural Language Generation

290 We set the batch size as 32 for WikiText-2 and 16 for GSM8K. We train 2 epochs on WikiText-2 and
 291 6 epochs on GSM8K. We select learning rate from $\{1 \times 10^{-5}, 5 \times 10^{-5}, 7 \times 10^{-5}, 1 \times 10^{-4}, 3 \times$

292 $10^{-4}, 4 \times 10^{-4}$. We use five iterations for all experiments. Specific settings are summarized as
 293 below

Table 14: Hyper-parameter setup of LoftQ LLAMA-2-series on GSM8K

Model	Hyperparameter	4-bit NF4	2-bit NF2	Mixed-precision
LLAMA-2-7b	learning rate	1e-4	1e-4	3e-4
LLAMA-2-13b	learning rate	1e-4	1e-4	3e-4

293

Table 15: Hyper-parameter setup of LoftQ LLAMA-2-series on WikiText-2

Model	Hyperparameter	4-bit NF4	2-bit NF2	Mixed-precision
LLAMA-2-7b	learning rate	1e-4	1e-4	3e-4
LLAMA-2-13b	learning rate	1e-4	1e-4	3e-4

294 H Comparison to Pruning

295 Pruning is also a widely used compression method. Here we compare LoftQ with the state-of-the-art
 296 pruning method Li et al. (2023). We show the comparison in Table 16. We can see our method
 297 significantly outperforms the pruning methods on DeBERTaV3-base model. We also remark that
 298 LoftQ can consistently reduce the memory of both training and storage. In contrast, pruning requires
 299 training the entire full-precision matrix, which implies that it can not achieve any memory savings
 during the training stage.

Table 16: Results of LoftQ using 2-bits uniform quantization compared with LoSparse with DeBERTaV3-base models on some of GLUE development sets. Here *Ratio* is the proportion of total remaining weights. Results with *N.A.* indicate the model does not converge.

Method	Ratio	MNLI m / mm	SST-2 Acc	QNLI Acc
Full FT	100%	90.5 / 90.6	95.3	94.0
LoSparse	15%	83.3/82.9	87.6	90.4
	20%	84.5/83.8	91.7	88.6
LoftQ	15.6%	87.3/87.1	94.0	90.6
	18.8%	88.0/88.1	94.7	92.4

300