

---



# DuQuant: Distributing Outliers via Dual Transformation Makes Stronger Quantized LLMs

---

Haokun Lin<sup>\*1,3,4</sup>, Haobo Xu<sup>\*2</sup>, Yichen Wu<sup>\*4</sup>, Jingzhi Cui<sup>2</sup>, Yingtao Zhang<sup>2</sup>,  
Linzhan Mou<sup>5</sup>, Linqi Song<sup>4</sup>, Zhenan Sun<sup>†1,3</sup>, Ying Wei<sup>†4,5</sup>

<sup>1</sup> School of Artificial Intelligence, University of Chinese Academy of Sciences

<sup>2</sup> Tsinghua University    <sup>3</sup> NLPR & MAIS, Institute of Automation, CAS

<sup>4</sup> City University of Hong Kong    <sup>5</sup> Zhejiang University

haokun.lin@cripac.ia.ac.cn    xuhb20@mails.tsinghua.edu.cn  
wuyichen.am97@gmail.com    znsun@nlpr.ia.ac.cn    ying.wei@zju.edu.cn

## Abstract

Quantization of large language models (LLMs) faces significant challenges, particularly due to the presence of outlier activations that impede efficient low-bit representation. Traditional approaches predominantly address *Normal Outliers*, which are activations across all tokens with relatively large magnitudes. However, these methods struggle with smoothing *Massive Outliers* that display significantly larger values, which leads to significant performance degradation in low-bit quantization. In this paper, we introduce DuQuant, a novel approach that utilizes rotation and permutation transformations to more effectively mitigate both massive and normal outliers. First, DuQuant starts by constructing the rotation matrix, using specific outlier dimensions as prior knowledge, to redistribute outliers to adjacent channels by block-wise rotation. Second, We further employ a zigzag permutation to balance the distribution of outliers across blocks, thereby reducing block-wise variance. A subsequent rotation further smooths the activation landscape, enhancing model performance. DuQuant simplifies the quantization process and excels in managing outliers, outperforming the state-of-the-art baselines across various sizes and types of LLMs on multiple tasks, even with 4-bit weight-activation quantization. Our code is available at <https://github.com/Hsu1023/DuQuant>.

## 1 Introduction

Large language models (LLMs) [51, 7, 50] have demonstrated exceptional performance across a wide range of natural language processing tasks. However, their billions of parameters present considerable deployment challenges on resource-constrained edge devices, particularly in terms of memory usage and inference speed [22, 15, 56]. In response to these challenges, network quantization methods [20, 23] have been extensively explored to minimize memory usage by converting floating-point parameters into low-bit formats [18, 32, 8], and to expedite inference by quantizing both activations and weights for accelerating the matrix multiplication process [64, 34, 74].

Among LLM quantization methods, a primary issue is the presence of activation outliers, which enlarge the quantization step sizes and subsequently cause significant accuracy loss [59]. To mitigate this problem, current research has developed various methods to address **Normal Outliers** in activations, which are *persistent in several channels across all tokens* [13, 64]. However, besides Normal Outliers, there exists another type of activation outlier [48, 35], termed **Massive Outliers**.

---

\*Equal contribution.

†Corresponding authors.

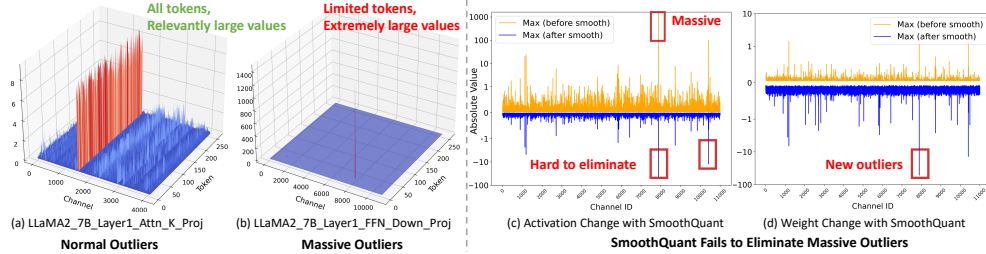


Figure 1: Visualizations of Outliers in LLaMA2-7B. (a) Input activation of Layer1 attention key projection shows Normal Outliers with relatively high magnitudes across all token sequences. (b) Input activation of Layer1 FFN down projection reveals Massive Outliers, presenting extremely high magnitudes (around 1400) at very few tokens. (c) Application of SmoothQuant on FFN down projection, illustrating its struggle with massive outliers in the Activation matrix. (d) Corresponding weight changes with SmoothQuant, highlighting the emergence of new outliers.

These outliers are characterized by their *exceedingly high values and limited occurrence in a subset of tokens*, as depicted in Figure 1(b). Unfortunately, existing LLM quantization methods struggle to effectively address these Massive Outliers. For instance, SmoothQuant [64], despite using a smooth factor to shift some of the activation outliers to the weight part, still cannot effectively handle Massive Outliers with extremely large values, as shown in Figure 1(c)(d). OmniQuant [47] and AffineQuant [39], on the other hand, exhibit training instability issues [34] due to the presence of Massive Outliers. Consequently, there is a pressing need for an LLM quantization approach that effectively addresses both Normal and Massive Outliers.

To tackle this challenge, we propose the **Dual transformations Quantization (DuQuant)** method. Our motivation is to **redistribute the activation outlier values across different channels**, facilitating easier quantization. Specifically, we construct the orthogonal rotation matrix and the orthogonal permutation matrix. By multiplying these matrices with the activations, we can effectively perform column transformations on the activations, which in turn allows for the redistribution of outliers. For the **rotation transformation** aspect, we first identify specific dimensions of outliers as the prior knowledge and employ a greedy algorithm to construct the rotation matrix. To enhance the multiplication efficiency, we utilize diagonal block-wise rotation matrices, with each matrix responsible for a small portion of the activations. However, this approach may result in uneven outlier magnitudes across different blocks. Therefore, we propose the **zigzag permutation** for reordering the activation channels, which promotes a more uniform distribution across different blocks. Concretely, we distribute the channels with the highest activations across the blocks in a back-and-forth pattern. After establishing blocks with uniformly distributed outlier magnitudes, we employ another rotation transformation to further redistribute the outliers within each block. Note that we multiply the weight matrix with the transpose of the rotation and permutation matrices at the same time, preserving the linear layer equivalence and smoothing weights. Theoretical analysis confirms that the rotation and permutation transformations greatly mitigate quantization challenges induced by outliers.

As a result, DuQuant offers several clear advantages over QuaRot [2]: (1) DuQuant’s optimal rotation matrix, derived through a greedy search guided by prior knowledge, surpasses QuaRot’s Hadamard rotation in managing outliers; (2) our unique zigzag permutation significantly reduces activation variance across blocks, providing a distinct advantage for handling massive outliers; and (3) by jointly smoothing weights and activations, DuQuant avoids time-consuming GPTQ algorithm in QuaRot. Extensive evaluations demonstrate that our DuQuant approach significantly outperforms existing 4-bit weight-activation quantization baselines across various benchmarks. Notably, DuQuant achieves a 5% improvement in Commonsense QA tasks across all LLaMA model sizes and a 10% increase in zero-shot MMLU benchmarks for the Vicuna-v1.5-13B. Moreover, in practical applications with the LLaMA2-7B model, DuQuant not only accelerates pre-filling phase by up to  $2.08\times$  but also reduces memory usage during decoding phase by  $3.50\times$ , with minimal impact on performance: only a 0.61 increase in perplexity and a 2.71% drop in accuracy compared to the FP16 model. These results highlight the effectiveness of DuQuant in enhancing the efficiency and capacity of quantized LLMs.

## 2 Motivation

**Normal Outliers and Massive Outliers.** Previous works [13, 72, 32] have highlighted the challenge posed by activation outliers in LLMs for model compression. These outlier features consistently

manifest large values across specific feature dimensions and are present in all token sequences [64], which we refer to as *Normal Outliers*. Recently, a distinct type of outlier [48, 35], termed *Massive Outliers*, has been observed in LLMs. The primary distinctions between normal and massive outliers are: 1) Normal outliers persist across all token sequences, whereas massive outliers are confined to a limited number of tokens. 2) Massive outliers exhibit significantly larger magnitudes, often surpassing 100 and being approximately 1000 times greater than the median of other activations [48]. In our study, we delve deeper into the impact of these two distinct types of outliers on quantization.

**Massive Outliers Exist at the Second Linear Layer of FFN Module.** In contrast to previous studies [48, 35] that observe massive outliers at the output of Transformer blocks, we **first** discover that these extremely large activations exist at the input of the down-projection layer within the FFN module. As depicted in Figure 1, the input of the down-projection layer in the LLaMA2-7B model Layer 1 contains a single activation of significant magnitude (approximately 1400). This activation is isolated to one token and therefore classified as one of massive activations. This phenomenon is consistently observed across different layers and sizes of models, as illustrated in Appendix I.

**Massive Outliers Enlarge Quantization Difficulty.** Although previous studies [64, 47, 39, 1] have proposed various approaches to eliminate outlier features, they still face challenges in effectively managing massive outliers. SmoothQuant [64], for instance, attempts to shift the quantization difficulty from activations to weights by dividing the activation by a per-channel smoothing factor and multiplying it to the weight matrix. Nevertheless, we observe that this transfer at the input of the down-projection layer can cause the weights of the down-projection to display noticeable outliers, as demonstrated in Figure 1. This issue arises because massive outliers cause the smoothing factor to become significantly large. Moreover, extremely large outliers can lead optimization-based methods to encounter problems with loss explosion. Both OmniQuant [47] and AffineQuant [39] have had to exclude their learnable parameters for the down projection layer due to unstable gradients. Given the poor accuracy observed with 4-bit quantization, QUIK [1] opts to use INT8 quantization for the down projection layer and Atom [74] applies INT8 quantization for 128 outlier channels. Consequently, massive outliers introduce new challenges to the quantization process that existing methods cannot fully address. This observation has motivated us to develop rotation and permutation transformations, which effectively handles both massive and normal outliers and achieves state-of-the-art performance.

### 3 Method

In this section, we delve into the distribution of outliers and introduce our proposed DuQuant method. The DuQuant method is built on two key components: 1) the block-diagonal rotation matrix, tasked with the local redistribution of feature outliers, and 2) the zigzag permutation, responsible for the global reordering of outliers across different blocks.

#### 3.1 Preliminaries

As the common modules within each transformer block of LLMs, both Multi-head Self-Attention (MSA) and Feed-Forward Network (FFN) fundamentally consist of basic linear layers, which can be represented as,  $\mathbf{Y} = \mathbf{X} \cdot \mathbf{W} \in \mathbb{R}^{T \times C_{out}}$ . Here,  $\mathbf{X} \in \mathbb{R}^{T \times C_{in}}$  is the activation input and  $\mathbf{W} \in \mathbb{R}^{C_{in} \times C_{out}}$  denotes the weight matrix. In this paper, we focus on integer uniform quantization [25] of both activation and weight, aiming to achieve better hardware support. Specifically, the  $b$ -bit quantization process maps the FP16 tensor  $\mathbf{X}$  to low-bit integer  $\mathbf{X}_q$ :

$$\mathbf{X}_q = \text{clamp} \left( \left\lfloor \frac{\mathbf{X}}{\Delta} \right\rfloor + z, 0, 2^b - 1 \right), \text{ where } \Delta = \frac{\max(\mathbf{X}) - \min(\mathbf{X})}{2^b - 1}, z = - \left\lfloor \frac{\min(\mathbf{X})}{\Delta} \right\rfloor. \quad (1)$$

The notation  $\lfloor \cdot \rfloor$  means the nearest rounding operation,  $\Delta$  is the quantization step size and  $z$  represents the zero point. Following [64, 47, 34, 39], we employ per-token quantization for activation and per-channel quantization for weight, which entails assigning different step sizes to individual tokens of activations ( $\Delta_{\mathbf{X}} \in \mathbb{R}^{T \times 1}$ ) and different output channels of weights ( $\Delta_{\mathbf{W}} \in \mathbb{R}^{1 \times C_{out}}$ ).

#### 3.2 The proposed DuQuant Method

To address the Normal Outliers issue stated in Section 2, current quantization methods, such as SmoothQuant [64] and OmniQuant [64], usually adopt the smooth technique. Concretely, it involves the utilization of a per-channel smoothing diagonal matrix, denoted as  $\mathbf{\Lambda}$ , to scale the input activation and weight matrix. The adjustment allows us to rewrite the original linear layer as  $\mathbf{Y} = \mathbf{X} \cdot \mathbf{W} = (\mathbf{X} \cdot$

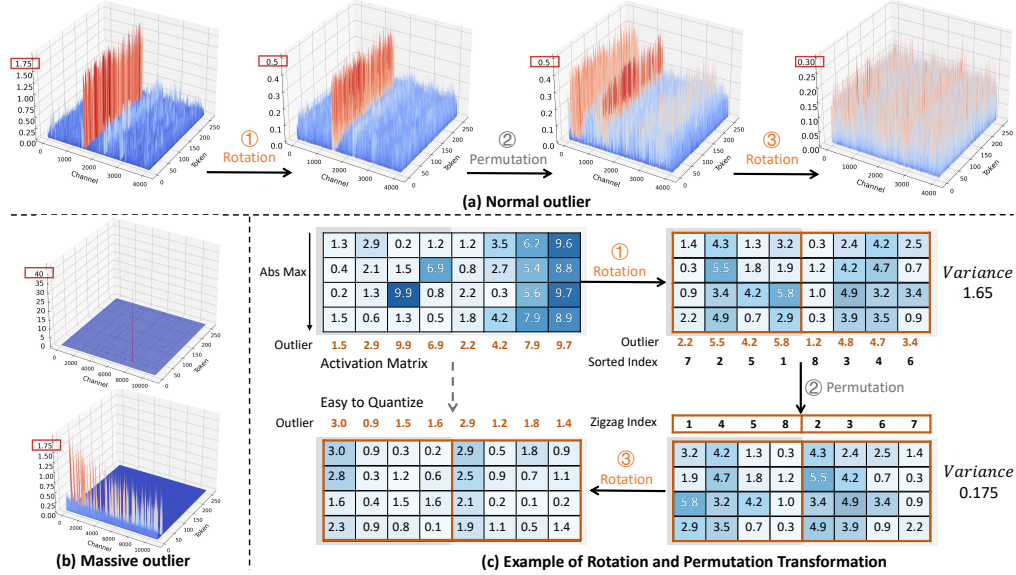


Figure 2: Transformation Steps for Activation Matrices after smooth technique. (a) Sequential transformations on Normal Outliers: ① initial rotation to reduce outliers within blocks, ② permutation to evenly distribute outliers across blocks, and ③ a second rotation for further smoothing. (b) Activation changes for Massive Outliers before and after DuQuant. (c) A sample matrix for highlighting the continual reduction of outliers through rotation and permutation, with outliers marked in dark blue.

$\Lambda^{-1})(\Lambda \cdot \mathbf{W})$ . The diagonal element  $\Lambda_j$  within  $\Lambda$  is computed as  $\Lambda_j = \max(|\mathbf{X}_j|)^\alpha / \max(|\mathbf{W}_j|)^{1-\alpha}$ , where  $\alpha$  is a hyper-parameter representing the migration strength. However, despite the ability of this smoothing technique to shift the quantization challenge from activations to weights, it still faces difficulties in effectively managing Massive Outliers, as depicted in Figure 1. This challenge stems from the extremely large massive outliers inducing large scaling factors  $\Lambda_j$ , which in turn introduce new outliers in the weight matrix and result in significant performance declines in 4-bit quantization.

According to these findings, we propose the DuQuant method, which includes the Rotation and Permutation transformations based on the smooth technique. By combining rotation transformation and channel permutation, our DuQuant method aims to redistribute these features within the activation space, thereby mitigating the effects of both Normal and Massive Outliers.

**The Rotation Transformation.** In contrast to the smooth technique, our aim is to apply a rotation matrix for row or column transformations, mitigating the impact of both Normal and Massive outliers. The ideal rotation matrix, denoted as  $\mathbf{R}$ , should possess the following properties: 1)  $\mathbf{R}$  is an orthogonal matrix satisfying  $\mathbf{R}\mathbf{R}^\top = \mathbf{I}$  and  $|\mathbf{R}| = \pm 1$ . This allows us to reformulate the linear layer within the transformer as  $\mathbf{Y} = \mathbf{X} \cdot \mathbf{W} = (\mathbf{X}\mathbf{R})(\mathbf{R}^\top\mathbf{W})$ ; 2)  $\mathbf{R}$  should be capable of effectively target the positions of outliers and effectively mitigating them through matrix multiplication. However, due to the Massive Outliers are usually randomly distributed within the activation space, it is challenging to directly identify the optimal rotation matrix  $\mathbf{R}$  capable of mitigating outliers through a single rotation transformation. To address this problem, we employ a greedy search with prior knowledge to compute a rotation matrix  $\hat{\mathbf{R}}$ , thereby approximating the ideal rotation matrix  $\mathbf{R}$ . Specifically, the calculation of  $\hat{\mathbf{R}}$  involves the following steps,

- Identify the feature dimension  $d^{(1)}$  where the outlier are primarily concentrated, i.e.,  $d^{(1)} = \arg \max_j (\max_i |\mathbf{X}_{ij}|)$ . Here,  $\mathbf{X}_{ij}$  represents the element in the  $i$ -th row and  $j$ -th column of  $\mathbf{X}$ .
- Based on the searched dimensions  $d^{(1)}$ , we construct the rotation matrix as follows,

$$\mathbf{R}^1 = \mathbf{E}_{d^{(1)}} \hat{\mathbf{R}} \mathbf{Q} \mathbf{E}_{d^{(1)}}, \quad \mathbf{Q} = \begin{bmatrix} 1 & \mathbf{O} \\ \mathbf{O} & \mathbf{Q}' \end{bmatrix}. \quad (2)$$

Here,  $\mathbf{E}_{d^{(1)}}$  is the switching matrix used to swap the first and the  $d^{(1)}$ -th columns of the activation, and  $\hat{\mathbf{R}}$  represents an orthogonal initialized rotation matrix, in which the first row is specifically uniformly distributed. The motivation behind this is to mitigate outliers in the first column after the transformation by  $\mathbf{E}_{d^{(1)}}$ . To further increase the randomness of the rotation operation,

we retain the first column, where outliers have been mitigated, and randomly rotate the other columns by multiplying them with a random orthogonal matrix  $\mathbf{Q}'$ .

- Let  $N$  denote the greedy search steps, then the approximated rotation matrix  $\hat{\mathbf{R}} = \mathbf{R}^1 \mathbf{R}^2 \cdots \mathbf{R}^n$ , where  $n = \arg \min_{k \in [1:N]} (\max_{i,j} |(\mathbf{X}\mathbf{R}^1 \cdots \mathbf{R}^k)_{ij}|)$ . Each  $\mathbf{R}^i$  is constructed according to Eqn. (2) and the identified feature dimension  $d^{(i)}$ . Appendix G provides detailed pseudo code.

Through this construction manner, we can ensure that the approximated optimal rotation matrix  $\hat{\mathbf{R}}$  can effectively mitigate outliers with large magnitudes, as opposed to merely using a randomly selected orthogonal rotation matrix. Nevertheless, directly constructing the entire rotation matrix is time-consuming and results in substantial memory overhead. For fast matrix multiplication, following [63], we approximate the rotation matrix  $\hat{\mathbf{R}} \in \mathbb{R}^{C_{in} \times C_{in}}$  in a block-wise manner,

$$\hat{\mathbf{R}} = \text{BlockDiag}(\hat{\mathbf{R}}_{b_1}, \dots, \hat{\mathbf{R}}_{b_K}), \quad (3)$$

where  $\hat{\mathbf{R}}_{b_i} \in \mathbb{R}^{2^n \times 2^n}$  denotes a square matrix of the  $i$ -th block, which is constructed following the three steps mentioned above. And the block numbers  $K$  is calculated by  $K = C_{in}/2^n$ .

**The Permutation Transformation.** Despite adopting the block-diagonal rotation matrix  $\hat{\mathbf{R}}$  for its time and storage efficiency, its focus on local information introduces a potential limitation in further reducing the outliers. This is because the rotation transformation, conducted within each small block, cannot integrate the information across different blocks to further minimize outliers. Consequently, one block may have relatively larger outliers while another block has smaller outliers, resulting in high variance among different blocks, as shown in Figure 2. This limitation explains that merely utilizing the block-diagonal rotation matrix is insufficient to effectively reduce the overall outliers.

To effectively mitigate the overall outliers, it is essential to balance the outliers' magnitudes among various blocks. Specifically, within each small block, we denote the largest outlier in dimension  $d_j$  as  $O_j$ . Meanwhile,  $M_{b_i}$  represents the mean value of all  $O_j$  in the  $i$ -th block, where  $i = 1, 2, \dots, K$ . Then the variance in activation magnitudes across various blocks can be expressed as,

$$\text{Var}([M_{b_1}, M_{b_2}, \dots, M_{b_K}]). \quad (4)$$

To minimize this variance and further reduce the overall outliers, we introduce the **zigzag permutation**. Concretely, we generate a zigzag sequence that starts by assigning channels with the highest activations to the first block. The process continues by assigning channels with the next highest activations to the subsequent blocks in descending order until the end of block  $K$ . Upon reaching the final block, the order reverses, starting from the channel with the next highest activation and proceeding in ascending order. This back-and-forth patterning continues throughout all the blocks, ensuring that no single block consistently receives either the highest or lowest activation channels. It is worth noting that the constructed permutation is an orthogonal matrix, which we denote as  $\mathbf{P}$ , satisfying the conditions  $\mathbf{P}\mathbf{P}^\top = \mathbf{I}$  and  $|\mathbf{P}| = \pm 1$ . By employing the zigzag permutation, we achieve a balanced distribution of outliers across different blocks. This allows us to use an additional rotation transformation to further smooth the outliers. Figure 2 provides an illustration of outlier mitigation.

**The Overall DuQuant Method.** To effectively mitigate both Normal and Massive Outliers, we first employ the smooth technique to shift the quantization challenge from activations to weights. Next, we introduce the block-diagonal rotation matrix  $\hat{\mathbf{R}}$  to locally redistribute feature outliers within the activation space. We then propose the zigzag permutation matrix for globally balancing the outliers across different blocks, followed by another application of the block-diagonal rotation transformation. To sum up, the linear layers within the transformer can be rewrite as,

$$\mathbf{Y} = \mathbf{X} \cdot \mathbf{W} = \underbrace{[(\mathbf{X} \cdot \mathbf{\Lambda}^{-1}) \hat{\mathbf{R}}_{(1)} \cdot \mathbf{P} \cdot \hat{\mathbf{R}}_{(2)}]}_{\mathbf{G}} \cdot \underbrace{[\hat{\mathbf{R}}_{(2)}^\top \cdot \mathbf{P}^\top \cdot \hat{\mathbf{R}}_{(1)}^\top (\mathbf{\Lambda} \cdot \mathbf{W})]}_{\mathbf{G}^{-1}}, \quad (5)$$

where the notation  $\mathbf{P}$  denotes the orthogonal permutation matrix learned via the zigzag manner, the  $\hat{\mathbf{R}}_{(1)}$  and  $\hat{\mathbf{R}}_{(2)}$  represent the first and second block-diagonal rotation matrix, respectively.

**Remark 1.** It is worth noting that the proposed DuQuant method can simultaneously smooth the weight matrix. While the commonly adopted smooth technique is effective, it can cause the weight matrix of the down-projection layer to exhibit pronounced outliers, leading to performance degradation. However, in the proposed DuQuant method, the rotation transformation we designed is applied to not only the activation input but also the weight matrix. As a result, the outliers induced by the smooth technique can be mitigated through our approximated rotation matrix  $\hat{\mathbf{R}}$ , yielding a

smoother, more quantization-friendly weight matrix. Moreover, this approach eliminates the reliance on complex weight quantization techniques, such as GPTQ [18] used in Atom [74] and QuaRot [2].

**Remark 2.** To further decrease the computation and memory costs, we initially construct the  $k$ -th block rotation matrix  $\hat{\mathbf{R}}_{b_k}$ , with the  $k$ -th block containing the largest outlier. We then assign  $\hat{\mathbf{R}}_{b_i} = \hat{\mathbf{R}}_{b_k}$  for all  $1 \leq i \leq K$ . This strategy not only effectively mitigates the impact of outliers, but also reduces the number of block rotation matrices from  $K$  to 1, significantly reducing computation and memory requirements. Importantly, incorporating the invertible matrix  $\mathbf{G}$  from Eqn. (5) significantly eases the quantization challenges for  $\mathbf{X}$  and  $\mathbf{W}$ . Consequently, the quantization process acts as  $\mathbf{Y} = (\mathbf{X}\mathbf{G})(\mathbf{G}^{-1}\mathbf{W}) = \hat{\mathbf{X}} \cdot \hat{\mathbf{W}} \approx \Delta_{\hat{\mathbf{X}}}\Delta_{\hat{\mathbf{W}}}(\hat{\mathbf{X}}_q - z_{\hat{\mathbf{X}}})(\hat{\mathbf{W}}_q - z_{\hat{\mathbf{W}}})$ .

### 3.3 Theoretical Analysis

To further demonstrate the effectiveness of the proposed DuQuant method, we conduct a theoretical analysis of the rotation and permutation transformations. Theorem 1 shows that within each block, the constructed rotation matrix effectively mitigates the maximum outlier, thereby reducing the outlier magnitude through a greedy search. Theorem 2 reveals that the employed zigzag permutation ensures a balanced upper bound shared among different blocks. This suggests that the zigzag permutation effectively reduces the variance shown in Eqn. (4) and thus assists the rotation matrix in further decreasing the outliers. Please refer to Appendix B for detailed proofs.

**Theorem 1 (Rotation).** *For the activation input  $\mathbf{X} \in \mathbb{R}^{T \times C_{in}}$ ,  $\hat{\mathbf{R}} \in \mathbb{R}^{2^n \times 2^n}$  is a diagonal block matrix constructed as per Eqn. (3). For a specific block  $b_i$ , let  $O_j(\cdot)$  represent the maximum outlier of the  $j$ -th dimension  $d_j$  within the input. Then, we can deduce that,*

$$\max_{1 \leq j \leq 2^n} O_j(\mathbf{X}_{b_i} \hat{\mathbf{R}}_{b_i}) \leq \max_{1 \leq j \leq 2^n} O_j(\mathbf{X}_{b_i}). \quad (6)$$

**Theorem 2 (Zigzag Permutation).** *For the activation input  $\mathbf{X} \in \mathbb{R}^{T \times C_{in}}$ , it can be divided into  $K$  blocks, where  $K = C_{in}/2^n$ . Let  $O_j$  denote the max outlier of the dimension  $d_j$  in  $\mathbf{X}$ , the reordered outliers from large to small is expressed as  $O^{(1)}, O^{(2)}, \dots, O^{(C_{in})}$ . Moreover, the  $M_{b_i}$  represents the mean value of all  $O_j$  in the  $i$ -th block,  $i = 1, 2, \dots, K$ . Let  $\delta := \max\{|O^{(i+1)} - O^{(i)}|\}$ ,  $i = 1, 2, \dots, C_{in}-1$ . Then, following the zigzag permutation described in Section 3.2, the mean value  $M_{b_i}$  within each  $i$ -th block consistently satisfies,*

$$M_{b_i} \leq O^{(1)} + \frac{(2^n K - 1)(2^{n-1} - 1)}{2^n} \delta, \quad i = 1, 2, 3, \dots, K. \quad (7)$$

## 4 Experiment

**Models and Evaluations.** We apply our DuQuant on pre-trained LLMs: LLaMA (7B-65B) [51], LLaMA2 (7B-70B) [52], LLaMA3 (8B, 70B), Mistral, Phi2 and instruction-tuned LLMs: Vicuna-v1.5 (7B-13B) [10]. We evaluate quantized pre-trained LLMs on language generation tasks and commonsense QA tasks. Specifically, we assess the perplexity on WikiText2 [40] and C4 [44] datasets, as well as the zero-shot accuracy on PIQA [6], ARC [12], BoolQ [11], HellaSwag [70], and WinoGrande [45] datasets. Moreover, we evaluate quantized Vicuna models on MMLU [21] and MT-Bench [76] benchmarks, as well as their long-form generative capabilities on LongBench [4].

**Implementation Details.** In line with prior studies [34, 47, 39], we apply per-token activation quantization and per-channel weight quantization. Given that W8A8 quantization has been established as lossless in precision by SmoothQuant [64], our primary evaluation in this paper focuses on 4-bit and 6-bit quantization for weights and activations. As for details, we quantize all intermediate activations, excluding the SoftMax output. Moreover, we have developed two types of quantized models, denoted as **DuQuant** and **DuQuant<sub>+LWC</sub>**. For DuQuant, we employ round-to-nearest quantization, using a clipping ratio of 0.9 for activations and 0.8 for weights. To improve weight matrix quantization, DuQuant<sub>+LWC</sub> integrates the learnable weight clipping (LWC) technique from OmniQuant. Concretely, LWC adjusts weights by training parameters  $\gamma, \beta \in [0, 1]$  to compute step size  $\Delta = \frac{\gamma \max(\mathbf{X}) - \beta \min(\mathbf{X})}{2^b - 1}$  in Eqn. (1). Notably, the smoothing diagonal matrix and the learned weight clipping factor can be integrated into the quantized weights, introducing no additional computational or memory costs. More details and hyperparameters are left in Appendix C.

**Baselines.** We compare with state-of-the-art (SOTA) weight-activation PTQ methods, including SmoothQuant [64], Outlier Supression+ [59], OmniQuant [47], QLLM [34], AffineQuant [39], and Atom [74]. For Atom, we reproduce the results with no group-wise asymmetric quantization.

Table 1: Perplexity ( $\downarrow$ ) results under 4-bit weight-activation quantization. The results for W6A6 can be found in Table D8. Atom and OmniQuant unprocessed group-query attention for LLaMA2-70B.

Dataset	#Bit	Method	1-7B	1-13B	1-30B	1-65B	2-7B	2-13B	2-70B
WikiText2	FP16	-	5.68	5.09	4.10	3.53	5.47	4.88	3.31
	W4A4	SmoothQuant	25.25	40.05	192.40	275.53	83.12	35.88	26.01
		OmniQuant	11.26	10.87	10.33	9.17	14.26	12.30	NaN
		AffineQuant	10.28	10.32	9.35	-	12.69	11.45	-
		QLLM	9.65	8.41	8.37	6.87	11.75	9.09	7.00
		Atom	8.15	7.43	6.52	5.14	8.40	6.96	NaN
		<b>DuQuant</b>	6.40	5.65	4.72	4.13	6.28	5.42	3.79
		<b>DuQuant+LWC</b>	<b>6.18</b>	<b>5.47</b>	<b>4.55</b>	<b>3.93</b>	<b>6.08</b>	<b>5.33</b>	<b>3.76</b>
C4	FP16	-	7.08	6.61	5.98	5.62	6.97	6.46	5.52
	W4A4	SmoothQuant	32.32	47.18	122.38	244.35	77.27	43.19	34.61
		OmniQuant	14.51	13.78	12.49	11.28	18.02	14.55	NaN
		AffineQuant	13.64	13.44	11.58	-	15.76	13.97	-
		QLLM	12.29	10.58	11.51	8.98	13.26	11.13	8.89
		Atom	10.34	9.57	8.56	8.17	10.96	9.12	NaN
		<b>DuQuant</b>	7.84	7.16	6.45	6.03	7.90	7.05	5.87
		<b>DuQuant+LWC</b>	<b>7.73</b>	<b>7.07</b>	<b>6.37</b>	<b>5.93</b>	<b>7.79</b>	<b>7.02</b>	<b>5.85</b>

Table 2: Zero-shot QA ( $\uparrow$ ) results of LLaMA1 models under 4-bit weight-activation quantization. The results for LLaMA2 models and W6A6 quantization can be found in Table D1 D9, and D10.

Model	Method	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg.
LLaMA1-7B W4A4	FP16	77.47	52.48	41.46	73.08	73.00	67.07	64.09
	SmoothQuant	49.80	30.40	25.80	49.10	27.40	48.00	38.41
	OS+	62.73	39.98	30.29	60.21	44.39	52.96	48.43
	OmniQuant	66.15	45.20	31.14	63.51	56.44	53.43	52.65
	AffineQuant	69.37	42.55	31.91	63.73	57.65	55.33	53.42
	QLLM	68.77	45.20	31.14	-	57.43	56.67	51.84
	Atom	71.44	47.74	35.49	67.71	63.89	55.01	56.88
	<b>DuQuant</b>	<b>76.44</b>	<b>50.04</b>	<b>38.99</b>	<b>70.98</b>	<b>69.39</b>	<b>64.72</b>	<b>61.76</b>
<b>DuQuant+LWC</b>	76.22	<b>50.04</b>	38.31	70.09	<b>69.82</b>	62.59	61.18	
LLaMA1-13B W4A4	FP16	79.10	59.89	44.45	68.01	76.21	70.31	66.33
	SmoothQuant	61.04	39.18	30.80	61.80	52.29	51.06	49.36
	OS+	63.00	40.32	30.38	60.34	53.61	51.54	49.86
	OmniQuant	69.69	47.39	33.10	62.84	58.96	55.80	54.37
	AffineQuant	66.32	43.90	29.61	64.10	56.88	54.70	52.58
	QLLM	71.38	47.60	34.30	-	63.70	59.43	55.28
	Atom	71.38	49.07	36.69	64.53	68.00	58.56	58.04
	<b>DuQuant</b>	77.26	<b>58.04</b>	<b>41.55</b>	<b>67.55</b>	73.62	<b>66.69</b>	<b>64.12</b>
<b>DuQuant+LWC</b>	<b>77.64</b>	57.32	41.21	66.79	<b>74.12</b>	65.98	63.84	
LLaMA1-30B W4A4	FP16	80.08	58.92	45.47	68.44	79.21	72.53	67.44
	SmoothQuant	58.65	35.53	27.73	60.42	35.56	48.06	44.83
	OS+	67.63	46.17	34.40	60.70	54.32	52.64	52.62
	OmniQuant	71.21	49.45	34.47	65.33	64.65	59.19	56.63
	AffineQuant	70.84	49.41	37.12	70.12	65.53	58.64	58.61
	QLLM	73.83	50.67	38.40	-	67.91	58.56	57.87
	Atom	71.98	49.07	40.02	66.85	70.45	58.64	59.50
	<b>DuQuant</b>	78.56	<b>56.99</b>	42.32	66.73	76.70	69.61	65.15
<b>DuQuant+LWC</b>	<b>78.73</b>	56.52	<b>43.17</b>	<b>68.84</b>	<b>77.53</b>	<b>70.96</b>	<b>65.96</b>	
LLaMA1-65B W4A4	FP16	80.79	58.71	46.24	82.29	80.72	77.50	71.04
	SmoothQuant	64.47	40.44	29.82	59.38	39.90	52.24	47.71
	OS+	68.06	43.98	35.32	62.75	50.73	54.30	52.52
	OmniQuant	71.81	48.02	35.92	73.27	66.81	59.51	59.22
	QLLM	73.56	52.06	39.68	-	70.94	62.90	59.83
	Atom	74.48	51.60	40.61	73.76	73.78	62.12	62.73
	<b>DuQuant</b>	79.71	57.95	<b>45.05</b>	<b>79.82</b>	78.66	<b>72.29</b>	<b>68.91</b>
	<b>DuQuant+LWC</b>	<b>79.98</b>	<b>58.29</b>	44.80	77.89	<b>79.22</b>	72.21	68.73

## 4.1 Main Results

**Quantization of LLaMA1 and LLaMA2 Models.** We conduct a comprehensive comparison of our DuQuant with several SOTA baselines on LLaMA1 and LLaMA2 models. Results for W4A4 quantization are presented in this Section, while results for W6A6 quantization are provided in

Table 3: Zero-shot and five-shot results on the MMLU benchmark for Vicuna-v1.5-13B under 4-bit weight-activation quantization. The results for Vicuna-v1.5-7b can be found in Table D2.

Model	Method	MMLU (0 shot) ↑					MMLU (5 shot) ↑				
		STEM	Hums	Social	Others	Avg.	STEM	Hums	Social	Others	Avg.
Vicuna-v1.5-13B W4A4	FP16	43.70	50.48	62.72	62.74	54.54	44.96	51.97	65.26	62.40	55.78
	SmoothQuant	21.70	24.29	22.13	23.16	22.82	25.31	24.97	26.00	27.08	25.84
	OmniQuant	26.81	26.57	30.35	28.75	28.12	28.79	27.29	31.13	28.99	29.05
	Atom	32.54	39.60	46.02	46.11	41.07	35.35	39.21	59.72	45.77	45.01
	<b>DuQuant</b>	<b>40.82</b>	46.61	58.73	57.59	50.94	40.92	<b>48.78</b>	<b>60.42</b>	57.71	<b>51.96</b>
	<b>DuQuant<sub>+LWC</sub></b>	40.13	<b>47.48</b>	<b>58.86</b>	<b>57.83</b>	<b>51.08</b>	<b>41.42</b>	48.52	58.73	<b>57.74</b>	51.61

Table 4: Long-context generation results for 4-bit Vicuna models on the LongBench benchmark.

Vicuna	Setting	Qasper	QMSum	MultiNews	TREC	TriviaQA	SAMSum	DuReader	RepoBench-P	Avg
Vicuna-v1.5-7B W4A4	FP16	23.27	21.07	26.91	66.00	82.59	41.06	25.53	48.23	41.83
	SmoothQuant	4.11	2.00	6.05	15.00	1.62	1.55	4.24	25.92	7.56
	OmniQuant	1.62	3.93	2.64	1.00	0.81	0.61	1.87	14.97	3.43
	Atom	17.97	20.24	24.60	58.00	67.20	37.94	19.41	29.34	34.34
	<b>DuQuant</b>	<b>19.98</b>	<b>21.15</b>	<b>25.85</b>	<b>64.00</b>	<b>78.91</b>	<b>42.24</b>	<b>23.15</b>	<b>47.66</b>	<b>40.37</b>
Vicuna-v1.5-13B W4A4	FP16	24.41	21.24	26.53	68.00	86.81	41.97	27.57	43.08	42.45
	SmoothQuant	2.18	2.95	3.54	1.50	1.83	0.35	6.71	11.57	3.83
	OmniQuant	0.68	1.78	2.83	9.00	1.13	0.45	13.83	8.46	4.77
	Atom	17.67	20.23	23.39	59.00	80.75	38.72	21.79	37.31	37.36
	<b>DuQuant</b>	<b>18.93</b>	<b>20.72</b>	<b>26.59</b>	<b>66.50</b>	<b>83.04</b>	<b>42.67</b>	<b>26.02</b>	<b>38.09</b>	<b>40.32</b>

Table 5: Perplexity and QA results of LLaMA3-8B under 4-bit/6-bit weight-activation quantization.

#Bits	Method	WikiText2 ↓	C4 ↓	PTB ↓	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg. ↑
FP16	-	6.14	8.88	9.91	80.85	77.78	53.41	81.28	79.16	72.84	74.22
LLaMA3-8B W6A6	SmoothQuant	7.07	9.57	11.69	78.94	75.88	49.49	77.58	77.39	70.8	71.68
	OmniQuant	7.24	9.82	11.90	78.90	73.95	47.35	74.95	76.77	70.56	70.41
	AffineQuant	7.35	9.99	12.30	78.73	73.32	46.08	74.59	77.08	70.88	70.11
	<b>DuQuant</b>	<b>6.27</b>	<b>8.38</b>	<b>10.77</b>	<b>80.20</b>	<b>77.27</b>	<b>52.05</b>	<b>80.12</b>	<b>79.14</b>	<b>72.77</b>	<b>73.59</b>
	<b>DuQuant<sub>+LWC</sub></b>	<b>6.27</b>	<b>8.38</b>	10.78	79.71	<b>77.57</b>	<b>53.07</b>	80.00	78.70	<b>73.09</b>	<b>73.69</b>
LLaMA3-8B W4A4	SmoothQuant	210.19	187.93	278.02	54.57	31.9	24.23	52.72	31.26	51.14	40.97
	OmniQuant	3.64e3	2.80e3	3.09e3	50.22	26.94	24.57	37.98	26.55	50.20	36.08
	AffineQuant	21.21e3	34.60e3	16.72e3	50.71	25.93	26.02	40.55	26.07	48.46	36.29
	Atom	22.14	31.83	40.04	62.95	49.45	30.12	60.31	53.75	56.04	52.10
	<b>DuQuant</b>	8.56	11.98	13.66	75.68	68.48	41.81	71.99	73.07	66.22	66.21
	<b>DuQuant<sub>+LWC</sub></b>	<b>8.06</b>	<b>11.29</b>	<b>13.19</b>	<b>76.22</b>	<b>70.41</b>	<b>43.69</b>	<b>74.34</b>	<b>73.87</b>	<b>67.80</b>	<b>67.72</b>

Appendix D. Table 1 indicates that our DuQuant quantized models notably outperform other baselines on both the WikiText2 and C4 datasets. Notably, LWC technique further enhances model capacity, with our DuQuant<sub>+LWC</sub> achieving comparable performance with FP16 models. Table 2 and Table D1 showcase the zero-shot accuracy of W4A4 quantization on Commonsense QA tasks, where DuQuant significantly improves the average accuracy. Our method surpasses QLLM by +9%, and Atom by +5% for all model sizes. These results demonstrate the superiority of our rotation and permutation transformation, which establishes new SOTA performance by effectively eliminating outlier features.

**Quantization of Instruction-tuned Models.** We quantize Vicuna-v1.5 [10] models to assess the generalizability of our DuQuant. Table 3 illustrates that our quantized models surpass the baselines across all task categories on MMLU benchmark. For Vicuna-13B, our DuQuant<sub>+LWC</sub> surpasses Atom by 10.01% under zero-shot settings and 6.95% under five-shot settings. Moreover, we compare our DuQuant with Atom and OmniQuant using MT-Bench and utilize GPT-4 to evaluate the answers from quantized models. As shown in Figure 3, DuQuant quantized models significantly outperform both Atom and OmniQuant in win rates. Specifically, for Vicuna-7B, DuQuant only lost 16 and 1 times to Atom and OmniQuant, respectively, while achieving 68 and 155 wins against them.

**Evaluation of Long-context Generation.** To further evaluate the long-text generative capabilities, we follow [37, 33, 55] and conduct a comprehensive comparison of DuQuant against state-of-the-art baselines on the LongBench [4], which includes a variety of generative tasks to provide a broader evaluation. We set the maximum sequence length to 3500 for Vicuna models, with results presented in Table 4. DuQuant achieves performance comparable to FP16 models, demonstrating the effectiveness of our dual transformations. More detailed results on different subtasks are listed in Table D3, D4.

**Quantization of LLaMA3 Models.** LLaMA3, known for its superior performance in various tasks, faces significant degradation in low-bit quantization [24]. To address this, we apply our DuQuant to quantize LLaMA3-8B. Table 5 displays the perplexity and zero-shot accuracy results. Notably, under W6A6 setting, our DuQuant achieves performance comparable to FP16 model. Furthermore, unlike other methods that show weaker results under W4A4 setting, our DuQuant maintains competitive performance, indicating its robustness with LLaMA3. We attribute this success to the advanced handling of outliers achieved through dual transformations, which is not restricted to specific models.



## 4.2 Ablation Study

**Module-wise Impact.** We ablate four distinct operations within DuQuant: 1) only the smoothing technology like SmoothQuant; 2) one rotation following the smoothing operation; 3) a sequence of rotation, permutation, and another rotation without smoothing; and k4) full DuQuant approach. Table 6 shows that the smoothing operation plays a basic role in our DuQuant by shifting activation outliers to weight. The initial rotation significantly enhances model performance, yielding competitive PPL results. Finally, permutation combined with a second rotation further enhances the quantized model.

Table 6: Influence of different components in DuQuant under 4-bit weight-activation quantization.

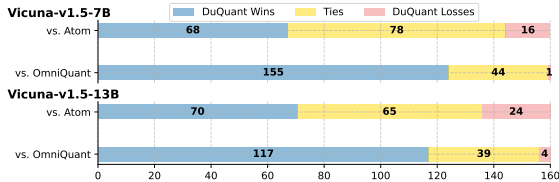
Modules				LLaMA2-7B		LLaMA2-13B	
Smooth	Rotation 1	Permutation	Rotation 2	WikiText2 ↓	C4 ↓	WikiText2 ↓	C4 ↓
✓				NaN	1379.46	160.30	203.87
	✓			8.48	10.63	14.32	21.73
✓				7.92	10.64	5.96	7.94
	✓			6.79	8.51	6.06	8.03
✓		✓	✓	<b>6.28</b>	<b>7.90</b>	<b>5.42</b>	<b>7.05</b>

**Influence of Normal/Massive Outliers.** In this section, we comprehensively explore the influence of massive and normal outliers on quantization. Notably, we observe that massive outliers primarily occur at the down-projection of the FFN module. To isolate their effect, we remove the rotation and permutation transformations, applying only the smoothing technique to all down-projection inputs. The resulting perplexity for LLaMA2-7B and LLaMA-13B showed significant degradation, presented in Table 7. Conversely, when we eliminate the rotation and permutation transformations for normal outliers, the performance decrease was noticeable but less severe compared to massive outliers. These findings indicate that: 1) massive outliers exert a more substantial impact on quantization, corroborating our claims in Section 2; 2) the smoothing technique alone struggles to fully mitigate the influence of outliers, particularly massive ones; and 3) our rotation and permutation methods prove highly effective against both types of outliers, leading to superior performance.

Table 7: Outliers impact on quantization. We only apply the smooth technique on Normal and Massive outliers for W4A4 quantization.

Outlier Type		LLaMA2-7B		LLaMA2-13B	
Normal	Massive	WikiText2 ↓	C4 ↓	WikiText2 ↓	C4 ↓
	✓	18.16	26.42	10.51	16.01
✓		10.88	13.89	7.87	10.52
✓	✓	<b>6.28</b>	<b>7.90</b>	<b>5.42</b>	<b>7.05</b>

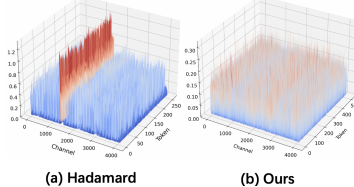
Figure 3: GPT-4 evaluation on the MT-Bench.



**Comparison with QuaRot [2]** In light of the recent introduction of Hadamard rotations by QuaRot[2] to eliminate outlier features, we have undertaken a detailed analysis to highlight the key differences between our DuQuant and QuaRot. To ensure a balanced evaluation, we have re-implemented QuaRot in accordance with our quantization settings. The results demonstrate that 1) the rotation matrix constructed by DuQuant outperforms QuaRot’s approach of simply selecting a randomly initialized Hadamard matrix. As depicted in Figure 10, our DuQuant more effectively smooths activations than QuaRot. This is attributed to the prior knowledge utilized by DuQuant to accurately target the outliers; 2) As demonstrated by the perplexity in Table 8, QuaRot employs GPTQ for their weight quantization method, whereas our DuQuant, with its sophisticated outlier management, attains competitive results using RTN quantization. The superiority proves the effectiveness of our zigzag permutation to enhance capacity. For a more comprehensive comparison, please refer to Appendix F.

Table 8: PPL (↓) comparison under W4A4 setting. Figure 4: LLaMA2-7B Attention key\_proj.

Method	1-7B	1-13B	1-30B	2-7B	2-13B
FP16	5.68	5.09	4.10	5.47	4.88
QuaRot-RTN	7.08	6.57	5.44	9.66	6.73
QuaRot-GPTQ	6.44	5.63	4.73	6.39	5.75
<b>DuQuant</b>	6.40	5.65	4.72	6.28	5.42
<b>DuQuant+LWC</b>	<b>6.18</b>	<b>5.47</b>	<b>4.55</b>	<b>6.08</b>	<b>5.33</b>



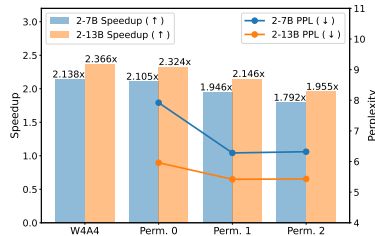
**Permutation Frequency.** We conduct ablations on the rotation and permutation frequencies in DuQuant. As shown in Figure 5, “Perm 1” (two rotations with one permutation) achieves stronger

performance compared with “Perm 0” (no permutation), while incurring an additional 8.9% computational cost on LLaMA2-7B and 9.3% on LLaMA2-13B compared to the W4A4 setup. Considering the approximately  $2\times$  speedup and the impressive performance, these additional costs are deemed acceptable. Further permutations, like “Perm 2,” do not improve performance and reduce inference efficiency. Consequently, “Perm 1” strikes the best balance between perplexity and inference speed, making it the optimal configuration for DuQuant.

**Inference Speedup.** To assess the inference speedup delivered by our DuQuant, we adopt the measurement strategy and W4A4 kernel from [2]. We evaluate the layer-wise speedup of LLaMA2 models on one NVIDIA 3090 GPU, with results detailed in Table 9 and 10. We set the pre-filling sequence length at 2048 and decode for 128 steps. In the pre-filling stage, DuQuant achieves a  $2.08\times$  speedup over FP16 for LLaMA2-7B and a  $2.34\times$  speedup for LLaMA2-13B, with slight variations across different batch sizes. In the decoding stage, batching the token generation phase yields high throughput without any downside [43]. Consequently, we enlarge the batch size to 64 and the results for LLaMA2-7B in Table 10 prove DuQuant achieves speedup comparable to QuaRot. More detailed analyses and end-to-end speedup are available in Appendix E.1.

Table 9: Layer-wise speedup during pre-filling stage for 4-bit weight-activation quantization.

Model	Batch Size	Speedup
LLaMA2-7B	1	1.95 $\times$
	4	2.03 $\times$
	16	2.08 $\times$
LLaMA2-13B	1	2.15 $\times$
	4	2.30 $\times$
	16	2.34 $\times$



**Memory Consumption.** We measure the peak memory usage of DuQuant with the W4A4 kernel on LLaMA2-7B using a single NVIDIA 3090 GPU. We process 2048 tokens for pre-filing and run 128 decoding steps, with the results listed in Table 11. In the pre-filling stage, DuQuant, SmoothQuant, and QuaRot achieve up to  $3.2\times$  memory reduction, while QLLM performs worse. In the decoding stage, DuQuant maintains strong memory efficiency, with superior performance.

Table 10: Decoding stage.

INT4, BS=64	Speedup
FP16	-
SmoothQuant	1.508 $\times$
QLLM	OOM
QuaRot	1.442 $\times$
DuQuant	1.321 $\times$

Table 11: Peak memory usage with a batch size of 1.

LLaMA2-7B	pre-filling (GB)	Saving	Decoding (GB)	Saving
FP16	15.282	-	13.638	-
SmoothQuant	4.782	3.196 $\times$	3.890	3.506 $\times$
QLLM	5.349	2.857 $\times$	3.894	3.502 $\times$
QuaRot	4.784	3.194 $\times$	3.891	3.505 $\times$
DuQuant	4.786	3.193 $\times$	3.893	3.503 $\times$

**Runtime.** Our DuQuant stands out for its efficiency, surpassing other baselines [47, 34, 39, 74]. The block-wise rotation ensures fast multiplication between the rotation and activation matrices. Zigzag permutation,

involving simple channel swaps, is much faster than complex algorithms like Simulated Annealing, as discussed in Appendix E.3. Moreover, the advanced management of outliers makes DuQuant not rely on GPTQ or gradient-based training. Hence, DuQuant enables a rapid quantization process shown in Table F24, e.g., we successfully quantize LLaMA2-13B in just 71s with superior performance.

Table 12: Quantization runtime on one NVIDIA A100.

Model	Omni.	Affine.	QLLM	Atom	DuQuant
LLaMA2-7B	2.0h	9.1h	1.1h	20min	50s
LLaMA2-13B	3.2h	16.0h	1.7h	36min	71s
LLaMA2-70B	14.6h	18.6h	9.3h	3.5h	270s

## 5 Conclusion

In conclusion, this paper presents DuQuant, an innovative quantization strategy for large language models (LLMs) that effectively addresses the challenge of outlier activations. By integrating rotation and permutation transformations, DuQuant effectively mitigates the impacts of both massive and normal outliers. This strategic redistribution of outliers not only simplifies the quantization process but also leads to substantial improvements in model performance. Consequently, DuQuant establishes new state-of-the-art results in 4-bit weight-activation quantization scenarios. This advancement enhances the deployment of efficient LLMs in resource-constrained environments.

## Acknowledgements

This work is supported in part by the National Natural Science Foundation of China (Grant No. U23B2054, 62276263, and 62371411), the Research Grants Council of the Hong Kong SAR under Grant GRF 11217823 and Collaborative Research Fund C1042-23GF, InnoHK initiative, the Government of the HKSAR, Laboratory for AI-Powered Financial Technologies.

## References

- [1] Saleh Ashkboos, Iliia Markov, Elias Frantar, Tingxuan Zhong, Xincheng Wang, Jie Ren, Torsten Hoefler, and Dan Alistarh. Towards end-to-end 4-bit inference on generative large language models. *arXiv preprint arXiv:2310.09259*, 2023.
- [2] Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L Croci, Bo Li, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. *arXiv preprint arXiv:2404.00456*, 2024.
- [3] 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.
- [4] Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*, 2023.
- [5] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- [6] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [8] Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization of large language models with guarantees. *Advances in Neural Information Processing Systems*, 36, 2024.
- [9] Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. Efficientqat: Efficient quantization-aware training for large language models. *arXiv preprint arXiv:2407.11062*, 2024.
- [10] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023.
- [11] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2924–2936, 2019.
- [12] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [13] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. In *Conference on Neural Information Processing Systems*, 2022.

- [14] Tim Dettmers, Ruslan A. Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. SpQR: A sparse-quantized representation for near-lossless LLM weight compression. In *The Twelfth International Conference on Learning Representations*, 2024.
- [15] Peijie Dong, Lujun Li, Zhenheng Tang, Xiang Liu, Xinglin Pan, Qiang Wang, and Xiaowen Chu. Pruner-zero: Evolving symbolic pruning metric from scratch for large language models. In *Proceedings of the 41st International Conference on Machine Learning*. PMLR, 2024. URL <https://arxiv.org/abs/2406.02924>. [arXiv: 2406.02924].
- [16] Dayou Du, Yijia Zhang, Shijie Cao, Jiaqi Guo, Ting Cao, Xiaowen Chu, and Ningyi Xu. Bitdistiller: Unleashing the potential of sub-4-bit llms via self-distillation. *arXiv preprint arXiv:2402.10631*, 2024.
- [17] Haojie Duanmu, Zhihang Yuan, Xiuhong Li, Jiangfei Duan, Xingcheng Zhang, and Dahua Lin. Skvq: Sliding-window key and value cache quantization for large language models. *arXiv preprint arXiv:2405.06219*, 2024.
- [18] 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.
- [19] Yuxin Guo, Siyang Sun, Shuailei Ma, Kecheng Zheng, Xiaoyi Bao, Shijie Ma, Wei Zou, and Yun Zheng. Crossmae: Cross-modality masked autoencoders for region-aware audio-visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26721–26731, 2024.
- [20] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- [21] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2020.
- [22] Jung Hwan Heo, Jeonghoon Kim, Beomseok Kwon, Byeongwook Kim, Se Jung Kwon, and Dongsoo Lee. Rethinking channel dimensions to isolate outliers for low-bit weight quantization of large language models. *arXiv preprint arXiv:2309.15531*, 2023.
- [23] Lu Hou, Quanming Yao, and James T Kwok. Loss-aware binarization of deep networks. *arXiv preprint arXiv:1611.01600*, 2016.
- [24] Wei Huang, Xudong Ma, Haotong Qin, Xingyu Zheng, Chengtao Lv, Hong Chen, Jie Luo, Xiaojuan Qi, Xianglong Liu, and Michele Magno. How good are low-bit quantized llama3 models? an empirical study. *arXiv preprint arXiv:2404.14047*, 2024.
- [25] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2704–2713, 2018.
- [26] Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 2023.
- [27] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [28] Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W Mahoney, and Kurt Keutzer. Squeezellm: Dense-and-sparse quantization. *arXiv preprint arXiv:2306.07629*, 2023.

- [29] Liang Li, Qingyuan Li, Bo Zhang, and Xiangxiang Chu. Norm tweaking: High-performance low-bit quantization of large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18536–18544, 2024.
- [30] Shiyao Li, Xuefei Ning, Luning Wang, Tengxuan Liu, Xiangsheng Shi, Shengen Yan, Guohao Dai, Huazhong Yang, and Yu Wang. Evaluating quantized large language models. *arXiv preprint arXiv:2402.18158*, 2024.
- [31] Haokun Lin, Haoli Bai, Zhili Liu, Lu Hou, Muye Sun, Linqi Song, Ying Wei, and Zhenan Sun. Mope-clip: Structured pruning for efficient vision-language models with module-wise pruning error metric. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27370–27380, 2024.
- [32] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration. *arXiv preprint arXiv:2306.00978*, 2023.
- [33] Akide Liu, Jing Liu, Zizheng Pan, Yefei He, Gholamreza Haffari, and Bohan Zhuang. Mini-cache: Kv cache compression in depth dimension for large language models. *arXiv preprint arXiv:2405.14366*, 2024.
- [34] Jing Liu, Ruihao Gong, Xiuying Wei, Zhiwei Dong, Jianfei Cai, and Bohan Zhuang. QLLM: Accurate and efficient low-bitwidth quantization for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=FIplmUWdm3>.
- [35] Ruikang Liu, Haoli Bai, Haokun Lin, Yuening Li, Han Gao, Zhengzhuo Xu, Lu Hou, Jun Yao, and Chun Yuan. IntactKV: Improving large language model quantization by keeping pivot tokens intact. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics ACL 2024*, pages 7716–7741, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.460. URL <https://aclanthology.org/2024.findings-acl.460>.
- [36] Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models. *arXiv preprint arXiv:2305.17888*, 2023.
- [37] Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint arXiv:2402.02750*, 2024.
- [38] Shijie Ma, Fei Zhu, Zhun Zhong, Wenzhuo Liu, Xu-Yao Zhang, and Cheng-Lin Liu. Happy: A debiased learning framework for continual generalized category discovery. *arXiv preprint arXiv:2410.06535*, 2024.
- [39] Yuexiao Ma, Huixia Li, Xiawu Zheng, Feng Ling, Xuefeng Xiao, Rui Wang, Shilei Wen, Fei Chao, and Rongrong Ji. Affinequant: Affine transformation quantization for large language models. *arXiv preprint arXiv:2403.12544*, 2024.
- [40] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *International Conference on Learning Representations*, 2016.
- [41] Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization. In *International Conference on Machine Learning*, pages 7197–7206. PMLR, 2020.
- [42] Davide Paglieri, Saurabh Dash, Tim Rocktäschel, and Jack Parker-Holder. Outliers and calibration sets have diminishing effect on quantization of modern llms. *arXiv preprint arXiv:2405.20835*, 2024.
- [43] Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Íñigo Goiri, Saeed Maleki, and Ricardo Bianchini. Splitwise: Efficient generative llm inference using phase splitting. In *2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA)*, pages 118–132. IEEE, 2024.

- [44] 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. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [45] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- [46] Yuzhang Shang, Zhihang Yuan, and Zhen Dong. PB-LLM: Partially binarized large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [47] Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- [48] Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. Massive activations in large language models. *arXiv preprint arXiv:2402.17762*, 2024.
- [49] 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.
- [50] Chaofan Tao, Qian Liu, Longxu Dou, Niklas Muennighoff, Zhongwei Wan, Ping Luo, Min Lin, and Ngai Wong. Scaling laws with vocabulary: Larger models deserve larger vocabularies. *arXiv preprint arXiv:2407.13623*, 2024.
- [51] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [52] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [53] Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks. *arXiv preprint arXiv:2402.04396*, 2024.
- [54] Zhongwei Wan, Xin Wang, Che Liu, Samiul Alam, Yu Zheng, Zhongnan Qu, Shen Yan, Yi Zhu, Quanlu Zhang, Mosharaf Chowdhury, et al. Efficient large language models: A survey. *arXiv preprint arXiv:2312.03863*, 1, 2023.
- [55] Zhongwei Wan, Xinjian Wu, Yu Zhang, Yi Xin, Chaofan Tao, Zhihong Zhu, Xin Wang, Siqi Luo, Jing Xiong, and Mi Zhang. D2o: Dynamic discriminative operations for efficient generative inference of large language models. *arXiv preprint arXiv:2406.13035*, 2024.
- [56] Zhongwei Wan, Ziang Wu, Che Liu, Jinfa Huang, Zhihong Zhu, Peng Jin, Longyue Wang, and Li Yuan. Look-m: Look-once optimization in kv cache for efficient multimodal long-context inference. *arXiv preprint arXiv:2406.18139*, 2024.
- [57] Haoxuan Wang, Yuzhang Shang, Zhihang Yuan, Junyi Wu, and Yan Yan. Quest: Low-bit diffusion model quantization via efficient selective finetuning. *arXiv preprint arXiv:2402.03666*, 2024.
- [58] Xin Wang, Yu Zheng, Zhongwei Wan, and Mi Zhang. Svd-llm: Truncation-aware singular value decomposition for large language model compression. *arXiv preprint arXiv:2403.07378*, 2024.
- [59] Xiuying Wei, Yunchen Zhang, Yuhang Li, Xiangguo Zhang, Ruihao Gong, Jinyang Guo, and Xianglong Liu. Outlier suppression+: Accurate quantization of large language models by equivalent and effective shifting and scaling. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1648–1665, 2023.

- [60] Junyi Wu, Haoxuan Wang, Yuzhang Shang, Mubarak Shah, and Yan Yan. Ptq4dit: Post-training quantization for diffusion transformers. *arXiv preprint arXiv:2405.16005*, 2024.
- [61] Yichen Wu, Long-Kai Huang, and Ying Wei. Adversarial task up-sampling for meta-learning. *Advances in Neural Information Processing Systems*, 35:31102–31115, 2022.
- [62] Yichen Wu, Long-Kai Huang, Renzhen Wang, Deyu Meng, and Ying Wei. Meta continual learning revisited: Implicitly enhancing online hessian approximation via variance reduction. In *The Twelfth International Conference on Learning Representations*, 2024.
- [63] Haocheng Xi, Changhao Li, Jianfei Chen, and Jun Zhu. Training transformers with 4-bit integers. *Advances in Neural Information Processing Systems*, 36:49146–49168, 2023.
- [64] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023.
- [65] Haobo Xu, Yuchen Yan, Dingsu Wang, Zhe Xu, Zhichen Zeng, Tarek F Abdelzaher, Jiawei Han, and Hanghang Tong. Slog: An inductive spectral graph neural network beyond polynomial filter. In *Forty-first International Conference on Machine Learning*, 2024.
- [66] Jaewoo Yang, Hayun Kim, and Younghoon Kim. Mitigating quantization errors due to activation spikes in glu-based llms. *arXiv preprint arXiv:2405.14428*, 2024.
- [67] Lianwei Yang, Zhikai Li, Junrui Xiao, Haisong Gong, and Qingyi Gu. Mgrq: Post-training quantization for vision transformer with mixed granularity reconstruction. *arXiv preprint arXiv:2406.09229*, 2024.
- [68] 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. In *Conference on Neural Information Processing Systems*, 2022.
- [69] Zhihang Yuan, Lin Niu, Jiawei Liu, Wenyu Liu, Xinggong Wang, Yuzhang Shang, Guangyu Sun, Qiang Wu, Jiayang Wu, and Bingzhe Wu. Rptq: Reorder-based post-training quantization for large language models. *arXiv preprint arXiv:2304.01089*, 2023.
- [70] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, 2019.
- [71] 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.
- [72] Yingtao Zhang, Haoli Bai, Haokun Lin, Jialin Zhao, Lu Hou, and Carlo Vittorio Cannistraci. Plug-and-play: An efficient post-training pruning method for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Tr01Px9woF>.
- [73] Tianchen Zhao, Xuefei Ning, Tongcheng Fang, Enshu Liu, Guyue Huang, Zinan Lin, Shengen Yan, Guohao Dai, and Yu Wang. Mixdq: Memory-efficient few-step text-to-image diffusion models with metric-decoupled mixed precision quantization. *arXiv preprint arXiv:2405.17873*, 2024.
- [74] Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. Atom: Low-bit quantization for efficient and accurate llm serving. *arXiv preprint arXiv:2310.19102*, 2023.
- [75] Zhenghao Zhao, Yuzhang Shang, Junyi Wu, and Yan Yan. Dataset quantization with active learning based adaptive sampling. *arXiv preprint arXiv:2407.07268*, 2024.
- [76] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023.
- [77] Yinan Zhou, Yuxin Chen, Haokun Lin, Shuyu Yang, Li Zhu, Zhongang Qi, Chen Ma, and Ying Shan. Doge: Towards versatile visual document grounding and referring. *arXiv preprint arXiv:2411.17125*, 2024.

## Appendix Overview

- Section A: Related work.
- Section B: Theory proofs.
- Section C: Additional implementation details.
- Section D: More empirical results.
- Section E: More detailed ablation studies.
- Section F: Detailed comparison with QuaRot.
- Section G: Algorithm for rotation matrix.
- Section H: Limitations and broader impacts.
- Section I: More visualization examples.

## A Related Work

### A.1 Network Quantization

Neural networks achieve great success in recent years [65, 77, 38, 61, 31, 19, 62], quantization [1, 63, 67, 74, 57, 75] is a widely utilized technique aimed at reducing model size and memory usage. Research in this area generally falls into two main categories: quantization-aware training (QAT) [71, 3, 49] and post-training quantization (PTQ) [41, 60, 73]. QAT involves training quantized model weights using additional data, often with the assistance of a straight-through estimator (STE) [5]. However, the computational cost associated with QAT poses challenges, particularly for large language models (LLMs) with millions of parameters, which necessitate significant amounts of data for retraining [36, 16, 9]. In contrast, PTQ has gained popularity for LLMs [54, 22, 69, 58, 37] due to its efficient approach, involving the training of quantized models using a small amount of data, known as calibration data [18]. However, PTQ often leads to significant performance degradation, especially when employing low-bit settings [18, 53, 46, 30]. Consequently, our work focuses on enhancing the performance of low-bit PTQ quantized models.

### A.2 Post Training Quantization of LLM

Post-training quantization for LLMs can be categorized into weight-only quantization [32, 14, 28, 29] and weight-activation quantization [68, 64, 74]. We focus on 4-bit weight-activation quantization due to the actual speedup it provides with low-bit quantization kernels [1]. Quantizing LLMs faces challenges due to activation outlier features persisting across different tokens and layers [13, 64]. Some approaches [13, 74] retain a small portion of crucial outlier channels at high precision (e.g., INT8), which poses challenges to hardware compatibility and leads to additional memory footprint. Other methods [64, 59, 47] attempt to shift quantization difficulty from activation to weight channels. However, the learnable equivalent transformation in OmniQuant [47] and the affine transform matrix in AffineQuant [39] exhibit instability as discussed in Section 2. The channel disassembly and assembly in QLLM [34], coupled with LoRA-tuning, incur significant time costs. Notably, these methods demonstrate poor performance under W4A4 quantization. We attribute this degradation to the ineffective handling of outlier features, especially massive outliers. Hence, we propose DuQuant to effectively eliminate outlier features through rotation matrices and channel permutation, achieving state-of-the-art performance. In contrast with QuaRot [2] also utilizing Hadamard matrices to enhance weight-activation quantization, our approach uniquely incorporates knowledge about the actual outlier channels. Furthermore, unlike QuaRot, which relies on GPTQ [18] for weight quantization, our permutation transformation has been proven helpful and efficient, facilitating a faster quantization process. The more detailed analysis and comparison with QuaRot are left in Appendix F. In addition, unlike RPTQ [69] and SKVQ [17], which use channel reordering to cluster similar activations, our method employs Permutation transformations with a fundamentally different goal: to evenly distribute outliers across blocks. This balanced distribution is crucial for enabling effective secondary rotations, ultimately leading to smoother activations that facilitate easier quantization.



## B Proofs

**Theorem 1** (Rotation). *For the activation input  $\mathbf{X} \in \mathbb{R}^{T \times C_{in}}$ ,  $\hat{\mathbf{R}} \in \mathbb{R}^{2^n \times 2^n}$  is a diagonal block matrix constructed as per Eqn. (3). For a specific block  $b_i$ , let  $O_j(\cdot)$  represent the maximum outlier of the  $j$ -th dimension  $d_j$  within the input. Then, we can deduce that,*

$$\max_{1 \leq j \leq 2^n} O_j(\mathbf{X}_{b_i} \hat{\mathbf{R}}_{b_i}) \leq \max_{1 \leq j \leq 2^n} O_j(\mathbf{X}_{b_i}). \quad (8)$$

*Proof.* In the case of a specific block  $b_i$ , the potential maximum value  $O_j(\mathbf{X}_{b_i} \hat{\mathbf{R}}_{b_i})$ , where  $\mathbf{X}_{b_i} \in \mathbb{R}^{T \times 2^n}$ , can be achieved by ensuring that the different  $j$ -th column outliers  $O_j(\mathbf{X}_{b_i})$  are located in the same  $t$ -th row of the activation input  $\mathbf{X}_{b_i}$ . This can be formally defined as follows,

$$\max_{1 \leq j \leq 2^n} O_j(\mathbf{X}_{b_i} \hat{\mathbf{R}}_{b_i}) = O_1(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \cdot \frac{1}{\sqrt{M}} + O_2(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \cdot \delta_1 + \dots + O_M(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \delta_M,$$

$$\text{where } M = 2^n \text{ and } \left\| \frac{1}{2^n} + \delta_2^2 + \dots + \delta_M^2 \right\| = 1, \quad (\hat{\mathbf{R}}_{b_i} \text{ is orthogonal})$$

Without loss of generality, let's assume that  $\delta_i \geq 0$  and define  $m := \arg \max_i \delta_i^2$ . Then we can derive that  $\delta_m \geq \frac{1}{\sqrt{2^n}}$ . Consequently, we can obtain the following inequality,

$$\begin{aligned} \max_{1 \leq j \leq 2^n} O_j(\mathbf{X}_{b_i} \hat{\mathbf{R}}_{b_i}) &= O_1(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \cdot \frac{1}{\sqrt{2^n}} + O_2(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \cdot \delta_1 + \dots + O_M(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \delta_M \\ &\leq O_1(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \cdot \frac{1}{\sqrt{2^n}} + (O_2(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) + \dots + O_M(\mathbf{X}_{b_i} \mathbf{E}_{d(1)})) \delta_m \\ &\stackrel{(1)}{\leq} O_1(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \cdot \frac{1}{\sqrt{2^n}} + O_1(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \cdot \frac{\sqrt{2^n} - 1}{\sqrt{2^n}} = O_1(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) \\ &= \max_{1 \leq j \leq 2^n} O_j(\mathbf{X}_{b_i}). \end{aligned} \quad (9)$$

The inequality (1) holds because the switch matrix  $\mathbf{E}_{d(1)}$  has swap the largest outliers  $O_1(\mathbf{X}_{b_i})$  in the first column, i.e.,  $O_1(\mathbf{X}_{b_i} \mathbf{E}_{d(1)}) > \max_{2 \leq m \leq M} O_m(\mathbf{X}_{b_i} \mathbf{E}_{d(1)})$ . □

**Theorem 2** (Zigzag Permutation). *For the activation input  $\mathbf{X} \in \mathbb{R}^{T \times C_{in}}$ , it can be divided into  $K$  blocks, where  $K = C_{in}/2^n$ . Let  $O_j$  denote the max outlier of the dimension  $d_j$  in  $\mathbf{X}$ , the reordered outliers from large to small is expressed as  $O^{(1)}, O^{(2)}, \dots, O^{(C_{in})}$ . Moreover, the  $M_{b_i}$  represents the mean value of all  $O_j$  in the  $i$ -th block,  $i = 1, 2, \dots, K$ . Let  $\delta := \max\{|O^{(i+1)} - O^{(i)}|\}, i = 1, 2, \dots, C_{in}-1$ . Then, following the zigzag permutation described in Section 3.2, the mean value  $M_{b_i}$  within each  $i$ -th block consistently satisfies,*

$$M_{b_i} \leq O^{(1)} + \frac{(2^n K - 1)(2^{n-1} - 1)}{2^n} \delta, \quad i = 1, 2, 3, \dots, K. \quad (10)$$

*Proof.* According to the zigzag permutation described in Section 3.2, and considering the reordered outliers  $O^{(1)}, O^{(2)}, \dots, O^{(C_{in})}$ , we can redistribute the channels they occupy (i.e.,  $O_c^{(i)}$ ) across different blocks. Specifically, for the  $i$ -th block, it contains the following channels,

$$\begin{aligned} b_i &= \{O_c^{2mK+i}, O_c^{2(m+1)K-i+1} | m = 0, 1, \dots, 2^{n-1} - 1\} \\ &= \{O_c^{(i)}, O_c^{(2K-i+1)}, \dots, O_c^{(2^n K+K)}, O_c^{(2^n K+K+1)}\} \end{aligned}$$

Since  $\delta = \max\{|O^{(i+1)} - O^{(i)}|\}, i = 1, 2, \dots, C_{in}-1$ , then we can get

$$\begin{aligned} M_{b_i} &= \frac{1}{2^n} \{O^{(1)} + O^{(2K)} + \dots + O^{2^n - 2K + 1} + O^{2^n K}\} \\ &\leq \frac{1}{2^n} \{O^{(1)} + (4K - 1)\delta + (8K - 1)\delta + [(2^{n+1} - 4)K - 1]\delta + (2^n K - 1)\delta\} \\ &\leq O^{(1)} + \frac{(2^n K - 1)(2^{n-1} - 1)}{2^n} \delta. \end{aligned} \quad (11)$$

Similarly, we can deduce that all  $M_{b_i}$  ( $i = 1, 2, \dots, K$ ) share the same upper bound after applying our zigzag permutation. □

## C Additional Implementation Details

In this work, all experiments are done on NVIDIA RTX 3090 GPUs for small-scale models and NVIDIA A100 GPUs for large-scale models. We set sequence length to 2048 for all evaluation tasks.

For calibration data, following [47, 39, 34], we randomly select 128 sampled sequences from the WikiText2 dataset, with the sequence length of 2048. For rotation and permutation transformations, the rotation block size  $2^n$  is set to 128, and maximum greedy search steps  $N$  equals 256. We adopt once permutation times for efficiency. We conduct detailed ablation studies in Appendix E.2, E.3, E.4.

Regarding quantization details, for multiplications between activations in MSA, such as Query and Key, attention outputs and Value, we apply a Hadamard rotation matrix for rapid and straightforward processing. A Hadamard matrix is an orthogonal and symmetric matrix filled with elements  $\pm 1/\sqrt{2^n}$ . For smooth parameter  $\alpha$ , we set it to 0.6 for DuQuant and 0.5 DuQuant-LWC. We clip the maximum activation values in all projection blocks, and the clipping ratio is set to 0.9. For DuQuant we also clip the maximum values in weight matrices, with a clipping ratio of 0.8. For LWC, we keep the same default epoch numbers of 20, batch size as 1, learning rate as 5e-3, and zero weight decay, as [47].

## D More Empirical Results

Table D1: Zero-shot common-sense QA ( $\uparrow$ ) results of LLaMA2 models under 4-bit WA quantization.

Model	Method	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg.
LLaMA2-7B W4A4	FP16	76.88	53.54	40.53	71.13	72.96	67.25	63.72
	SmoothQuant	60.17	35.23	27.13	57.92	37.08	49.57	44.52
	OS+	63.11	39.10	28.84	-	51.30	45.93	45.66
	OmniQuant	65.61	44.28	30.38	62.66	53.51	51.85	51.38
	AffineQuant	67.36	44.23	31.91	62.75	54.38	55.18	52.64
	QLLM	67.68	44.40	30.89	-	58.45	56.59	51.60
	Atom	69.75	47.35	34.22	62.42	63.21	56.51	55.58
	<b>DuQuant</b>	<b>75.24</b>	<b>51.89</b>	<b>36.77</b>	<b>67.86</b>	<b>69.54</b>	<b>62.12</b>	<b>60.57</b>
	<b>DuQuant-LWC</b>	<b>75.68</b>	<b>50.00</b>	<b>37.46</b>	<b>69.24</b>	<b>69.74</b>	<b>63.93</b>	<b>61.01</b>
LLaMA2-13B W4A4	FP16	79.05	57.91	44.20	69.02	76.60	69.69	66.08
	SmoothQuant	62.30	40.28	30.72	60.49	42.24	49.96	47.67
	OS+	64.47	41.46	32.17	-	59.30	51.38	49.76
	OmniQuant	69.80	47.22	33.79	65.47	59.34	55.49	55.19
	AffineQuant	68.55	47.64	32.34	66.97	59.97	55.07	55.09
	QLLM	70.46	48.48	34.39	-	62.80	55.41	54.31
	Atom	71.16	50.89	37.88	63.91	67.51	58.40	58.29
	<b>DuQuant</b>	<b>77.31</b>	<b>55.60</b>	<b>41.55</b>	<b>66.61</b>	<b>73.68</b>	<b>66.06</b>	<b>63.47</b>
	<b>DuQuant-LWC</b>	<b>77.26</b>	<b>56.23</b>	<b>42.15</b>	<b>65.78</b>	<b>73.68</b>	<b>65.43</b>	<b>63.42</b>
LLaMA2-70B W4A4	FP16	81.01	59.68	47.95	75.87	80.87	76.95	70.39
	SmoothQuant	64.09	41.84	32.00	58.56	54.21	51.07	50.30
	OS+	66.16	42.72	34.90	-	56.93	52.96	50.73
	QLLM	74.27	50.59	37.20	-	71.62	59.43	58.62
	<b>DuQuant</b>	<b>79.27</b>	<b>58.16</b>	<b>46.07</b>	<b>70.46</b>	<b>79.21</b>	<b>74.19</b>	<b>67.89</b>
	<b>DuQuant-LWC</b>	<b>79.82</b>	<b>59.76</b>	<b>46.76</b>	<b>73.12</b>	<b>79.38</b>	<b>74.11</b>	<b>68.83</b>

**Zero-shot QA Results for 4-bit LLaMA2 Models.** Table D1 showcases the zero-shot common-sense QA results for INT4 quantized LLaMA2 models. Our DuQuant method excels across various model sizes and datasets, demonstrating state-of-the-art performance in commonsense reasoning tasks. For example, DuQuant outperforms Atom by 5.43% for the LLaMA2-7B model and by 5.18% for the LLaMA2-13B model. In contrast to Atom [74], which relies on GPTQ for weight quantization and maintains 128 channels at INT8, thereby increasing memory usage, our method offers a rapid and more efficient weight-activation quantization solution through Rotation and Permutation.

**MMLU Results for 4-bit Vicuna-v1.5-7B.** Vicuna-v1.5 models [10], fine-tuned from LLaMA-2 models using high-quality user-shared conversations, are considered state-of-the-art chatbots. Table D2 displays the INT4 quantization results for Vicuna-v1.5-7B on the MMLU benchmarks. In comparison to SmoothQuant, OmniQuant, and Atom, our DuQuant method exhibits the smallest

performance decline and maintains competitive capacities in both zero-shot and five-shot settings. These results demonstrate the effectiveness of DuQuant in generalizing to instruction-tuned models.

Table D2: Zero-shot and five-shot results on the MMLU benchmark for quantized Vicuna-v1.5-7B.

Model	Method	MMLU (0 shot) ↑					MMLU (5 shot) ↑				
		STEM	Hums	Social	Others	Avg.	STEM	Hums	Social	Others	Avg.
Vicuna-v1.5-7B W4A4	FP16	38.70	45.42	56.13	56.01	49.07	39.56	45.76	58.14	57.43	50.22
	SmoothQuant	27.10	25.16	27.40	26.71	26.59	25.22	25.06	24.99	26.68	25.49
	OmniQuant	27.20	24.00	27.14	25.08	25.86	29.39	24.95	27.30	24.80	26.39
	Atom	30.28	34.73	38.97	40.56	36.14	31.97	35.37	40.46	40.81	37.15
	<b>DuQuant</b>	<b>35.85</b>	<b>42.66</b>	<b>52.03</b>	<b>51.23</b>	<b>45.44</b>	<b>38.90</b>	<b>42.57</b>	51.80	51.23	<b>46.13</b>
	<b>DuQuant+LWC</b>	35.18	41.91	51.28	50.52	44.72	37.34	42.21	<b>53.07</b>	<b>51.76</b>	46.10

**Long-context Evaluation Results.** LongBench [4] is proposed to access the long-context generation ability of LLMs, which covers several key long-text application scenarios. We evaluate the 4-bit quantized Vicuna models on five different tasks. Specifically, Qasper, MultiFieldQA, and NarrativeQA (F1 score) are Single-Document QA tasks; DuReader (Rouge-L score) and 2WikiMulti-hopQA (F1 score) are Multi-Document QA tasks; QMSum, GovReport (F1 score) and MultiNews (Rouge-L score) are Summarization tasks; TREC (Accuracy CLS), TriviaQA (F1 score), and SAMSum (Rouge-L score) are Few-shot Learning tasks; and RepoBench-P (similarity score) is Code Completion task. Table D3 and Table D4 show that our DuQuant outperforms other baselines by a clear margin, maintaining the ability for long context generation tasks compared with FP16 models.

Table D3: Long-context generation results on the LongBench benchmark for 4-bit Vicuna-v1.5-7B.

Vicuna-v1.5-7B	RepoBench-P	MultiFieldQA-en	GovReport	MultiNews	DuReader	2WikiMQA	TriviaQA
FP16	48.23	38.30	27.93	26.91	25.53	18.02	82.59
SmoothQuant	25.92	4.66	2.62	6.05	4.24	2.02	1.62
OmniQuant	14.97	2.30	2.51	2.64	1.87	0.48	0.81
Atom	29.34	31.15	23.60	24.60	19.41	<b>17.10</b>	67.20
<b>DuQuant</b>	<b>47.66</b>	<b>35.62</b>	<b>25.66</b>	<b>25.85</b>	<b>23.15</b>	15.09	<b>78.91</b>
Vicuna-v1.5-7B	QMSum	MultiFieldQA-zh	NarrativeQA	Qasper	SAMSum	TREC	Avg
FP16	21.07	32.56	14.96	23.27	41.06	66.00	35.88
SmoothQuant	2.00	0.88	1.75	4.11	1.55	15.00	5.57
OmniQuant	3.93	1.40	1.10	1.62	0.61	1.00	2.71
Atom	20.24	21.55	11.57	17.97	37.94	58.00	29.21
<b>DuQuant</b>	<b>21.15</b>	<b>29.56</b>	<b>11.31</b>	<b>19.98</b>	<b>42.24</b>	<b>64.00</b>	<b>33.86</b>

Table D4: Long-context generation results on the LongBench benchmark for 4-bit Vicuna-v1.5-13B.

Vicuna-v1.5-13B	RepoBench-P	MultiFieldQA-en	GovReport	MultiNews	DuReader	2WikiMQA	TriviaQA
FP16	43.08	42.69	28.43	26.53	27.57	29.40	86.81
SmoothQuant	11.57	1.64	2.81	3.54	6.71	1.39	1.83
OmniQuant	8.46	4.32	0.74	2.83	13.83	0.75	1.13
Atom	37.31	37.31	19.34	23.39	21.79	15.16	80.75
<b>DuQuant</b>	<b>38.09</b>	<b>44.12</b>	<b>26.97</b>	<b>26.59</b>	<b>26.02</b>	<b>22.07</b>	<b>83.04</b>
Vicuna-v1.5-13B	QMSum	MultiFieldQA-zh	NarrativeQA	Qasper	SAMSum	TREC	Avg
FP16	21.24	40.44	15.41	24.41	41.97	68.00	40.64
SmoothQuant	2.95	0.82	0.97	2.18	0.35	1.50	4.21
OmniQuant	1.78	1.06	0.62	0.68	0.45	9.00	4.58
Atom	20.23	28.02	8.81	17.67	38.72	59.00	33.58
<b>DuQuant</b>	<b>20.72</b>	<b>30.85</b>	<b>13.36</b>	<b>18.93</b>	<b>42.67</b>	<b>66.50</b>	<b>38.13</b>

**Comparison with FP16 models on MT-Bench.**

We conducted additional comparisons using the MT-Bench between our INT4 quantized models and the FP16 models. As shown in Table D5, for both 7B and 13B models, our DuQuant performs comparably to FP16, which further underscores the effectiveness of dual transformations in maintaining high accuracy even with reduced precision.

Table D5: More results on MT-Bench.

DuQuant v.s. FP16	Former Win	Tie	Former Loss
Vicuna-v1.5-7B	36	56	68
Vicuna-v1.5-13B	43	53	64

**Results for 4-bit Mistral-7B and Phi2-2.8B.** We have extended the application of DuQuant to include Mistral [27] and Phi2 [26] under 4-bit WA quantization. From Table D6, we can observe that

DuQuant largely surpasses other baselines, particularly with Mistral-7B. Regarding the Phi2-2.8B model, it often experiences instability in matrix multiplication between queries and values, leading to overflow issues and posing great challenges to quantization. However, while DuQuant may not perform as well as FP models, it still significantly outperforms other baselines. In addition, we have visualized the massive outliers in the down projection layer of the Mistral-7B model and the feature space after our dual transformations. These visualizations are shown in Figure I9. It can be observed that our DuQuant perfectly eliminates these outliers. These results underscore the effectiveness of our dual transformation approach in addressing massive outliers across various types of LLMs.

Table D6: Perplexity results of Mistral-7B and Phi2-2.8B under 4-bit weight-activation quantization.

Model	Method	WikiText2	C4	Model	Method	WikiText2	C4
Mistral-7B W4A4	FP16	5.25	7.75	Phi2-2.8B W4A4	FP16	9.71	12.76
	RTN	306.26	300.07		RTN	230.59	253.79
	SmoothQuant	100.59	158.02		SmoothQuant	63.84	83.24
	OmniQuant	5490.31	6094.82		OmniQuant	NaN	NaN
	Atom	8.65	12.43		Atom	35.72	41.26
	<b>DuQuant</b>	<b>5.86</b>	<b>8.48</b>		<b>DuQuant</b>	<b>20.65</b>	<b>22.49</b>

**Results for 4-bit LLaMA3-70B.** As LLaMA3 models have proven to be sensitive to quantization, we apply our DuQuant to the LLaMA3-70B and present the results in Table D7. Due to time constraints, we do not add learnable weight clipping. The results demonstrate that our DuQuant-quantized models outperform SmoothQuant by 12.9% on Commonsense QA tasks and significantly reduce perplexity across the WikiText2, C4, and PTB datasets. These improvements underscore the robustness of our DuQuant method when applied to the LLaMA3-70B model.

Table D7: Perplexity and QA results of LLaMA3-70B under 4-bit weight-activation quantization.

#Bits	Method	WikiText2 ↓	C4 ↓	PTB ↓	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg. ↑
FP16	-	2.9	6.9	8.2	82.4	86.9	60.3	85.2	84.9	80.6	80.1
LLaMA3-70B W4A4	SmoothQuant	9.6	16.9	17.7	76.9	75.8	43.5	64.4	62.9	58.9	63.7
	<b>DuQuant</b>	<b>4.9</b>	<b>8.3</b>	<b>8.7</b>	<b>81.1</b>	<b>80.8</b>	<b>57.3</b>	<b>81.3</b>	<b>82.1</b>	<b>77.0</b>	<b>76.6</b>

**W6A6 Quantization Results.** To thoroughly evaluate the effectiveness of our DuQuant models, we conduct comprehensive assessments under the W6A6 quantization setting. The perplexity results for language generation tasks are displayed in Table D8, while the zero-shot accuracy for Commonsense QA tasks is detailed in Tables D9 and D10. Our findings reveal that DuQuant not only surpasses other baselines but also achieves nearly lossless performance with FP16 models in these tasks. Interestingly, in several instances, DuQuant slightly outperforms DuQuant+LWC. This suggests that the Rotation and Permutation transformations alone are sufficient to create highly competitive quantized models under W6A6 settings, without the need for additional enhancements such as the learnable weight clipping (LWC) technique. These outcomes highlight the exceptional versatility and robustness of DuQuant across various quantization scenarios, confirming its potential as a leading solution in post-training quantization for large language models.

Table D8: Perplexity (↓) results on the WikiText2 and C4 datasets under 6-bit WA quantization.

Dataset	#Bit	Method	1-7B	1-13B	1-30B	1-65B	2-7B	2-13B	2-70B
WikiText2	FP16	-	5.68	5.09	4.10	3.53	5.47	4.88	3.31
	W6A6	SmoothQuant	6.03	5.42	4.55	3.88	6.20	5.18	3.69
		OmniQuant	5.96	5.28	4.38	3.75	5.87	5.14	3.71
		QLLM	5.89	5.28	4.30	3.73	5.91	5.08	3.55
		<b>DuQuant</b>	<b>5.73</b>	<b>5.13</b>	<b>4.14</b>	<b>3.57</b>	<b>5.53</b>	<b>4.92</b>	<b>3.35</b>
		<b>DuQuant+LWC</b>	5.74	<b>5.13</b>	4.15	3.60	<b>5.53</b>	<b>4.92</b>	<b>3.35</b>
C4	FP16	-	7.08	6.61	5.98	5.62	6.97	6.46	5.52
	W6A6	SmoothQuant	7.47	6.97	6.34	5.99	7.76	6.76	5.88
		OmniQuant	7.43	6.84	6.22	5.82	7.48	6.74	5.91
		QLLM	7.34	6.82	6.17	5.80	7.31	6.71	5.76
		<b>DuQuant</b>	<b>7.12</b>	<b>6.64</b>	<b>6.00</b>	<b>5.64</b>	<b>7.03</b>	<b>6.50</b>	<b>5.54</b>
		<b>DuQuant+LWC</b>	7.13	<b>6.64</b>	6.01	<b>5.64</b>	<b>7.03</b>	<b>6.50</b>	<b>5.54</b>

Table D9: Zero-shot common-sense QA ( $\uparrow$ ) results of LLaMA1 models under 6-bit WA quantization.

Model	Method	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg
LLaMA1-7B W6A6	FP16	77.47	52.48	41.46	73.08	73.00	67.07	64.09
	SmoothQuant	76.75	51.64	39.88	71.75	71.67	65.03	62.81
	OS+	76.82	51.35	41.13	72.08	71.42	65.98	61.13
	OmniQuant	77.09	51.89	40.87	72.53	71.61	65.03	63.17
	AffineQuant	76.60	52.29	40.63	72.65	71.29	63.85	62.89
	QLLM	77.26	52.02	41.04	-	71.40	65.19	61.38
	<b>DuQuant</b>	<b>77.53</b>	<b>51.47</b>	<b>41.13</b>	<b>72.78</b>	<b>72.76</b>	<b>66.69</b>	<b>63.73</b>
	<b>DuQuant-LWC</b>	<b>77.42</b>	<b>52.65</b>	<b>40.53</b>	<b>71.53</b>	<b>72.64</b>	<b>67.72</b>	<b>63.75</b>
LLaMA1-13B W6A6	FP16	79.10	59.89	44.45	68.01	76.21	70.31	66.33
	SmoothQuant	77.91	56.60	42.40	64.95	75.36	69.36	64.43
	OS+	78.29	56.90	43.09	66.98	75.09	69.22	64.92
	OmniQuant	78.40	57.28	42.91	67.00	75.82	68.27	64.95
	QLLM	77.91	57.70	42.92	-	75.02	69.14	64.54
	<b>DuQuant</b>	<b>78.62</b>	<b>59.51</b>	<b>44.03</b>	<b>68.44</b>	<b>75.98</b>	<b>70.08</b>	<b>66.11</b>
	<b>DuQuant-LWC</b>	<b>79.16</b>	<b>59.39</b>	<b>43.69</b>	<b>68.10</b>	<b>75.81</b>	<b>69.06</b>	<b>65.87</b>
	LLaMA1-30B W6A6	FP16	80.08	58.92	45.47	68.44	79.21	72.53
SmoothQuant		77.14	57.61	42.91	65.56	78.07	69.92	65.20
OS+		80.14	58.92	45.05	68.02	77.96	71.98	67.01
OmniQuant		79.81	58.79	45.22	68.38	78.95	72.21	67.23
QLLM		79.65	58.08	44.11	-	78.38	<b>73.24</b>	66.69
<b>DuQuant</b>		<b>79.43</b>	<b>59.34</b>	<b>44.54</b>	<b>70.15</b>	<b>78.89</b>	<b>72.77</b>	<b>67.52</b>
<b>DuQuant-LWC</b>		<b>80.09</b>	<b>57.95</b>	<b>45.05</b>	<b>68.72</b>	<b>79.17</b>	<b>73.09</b>	<b>67.35</b>
LLaMA1-65B W6A6		FP16	80.79	58.71	46.24	82.29	80.72	77.50
	SmoothQuant	80.25	57.92	45.50	80.22	80.18	74.76	69.80
	OS+	79.67	55.68	45.22	80.02	78.03	73.95	68.76
	OmniQuant	81.01	58.12	46.33	80.64	79.91	75.69	70.28
	QLLM	80.14	57.79	45.05	-	79.74	74.59	67.46
	<b>DuQuant</b>	<b>80.96</b>	<b>59.09</b>	<b>46.76</b>	<b>82.20</b>	<b>80.68</b>	<b>77.27</b>	<b>71.16</b>
	<b>DuQuant-LWC</b>	<b>80.63</b>	<b>58.00</b>	<b>46.50</b>	<b>82.08</b>	<b>80.49</b>	<b>76.87</b>	<b>70.76</b>

Table D10: Zero-shot common-sense QA ( $\uparrow$ ) results of LLaMA2 models under 6-bit WA quantization.

Model	Method	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg
LLaMA2-7B W6A6	FP16	76.88	53.54	40.53	71.13	72.96	67.25	63.72
	SmoothQuant	75.57	53.62	39.93	69.54	71.76	66.14	62.76
	OS+	76.22	52.74	40.70	-	71.89	65.19	61.35
	OmniQuant	76.55	53.83	40.96	68.75	55.89	65.59	60.26
	QLLM	77.48	52.99	39.33	-	71.38	65.98	61.43
	<b>DuQuant</b>	<b>76.99</b>	<b>52.99</b>	<b>40.87</b>	<b>70.40</b>	<b>72.49</b>	<b>67.32</b>	<b>63.51</b>
	<b>DuQuant-LWC</b>	<b>76.88</b>	<b>52.31</b>	<b>40.44</b>	<b>69.72</b>	<b>72.60</b>	<b>66.93</b>	<b>63.15</b>
	LLaMA2-13B W6A6	FP16	79.05	57.91	44.20	69.02	76.60	69.69
SmoothQuant		78.29	57.41	43.86	69.50	75.02	66.93	65.17
OS+		78.29	59.13	43.34	-	75.37	67.56	64.74
OmniQuant		78.24	57.58	43.86	71.10	75.52	68.35	65.78
AffineQuant		78.35	57.58	43.34	66.73	74.71	68.59	64.88
QLLM		78.78	58.29	43.77	-	75.10	68.43	64.87
<b>DuQuant</b>		<b>78.62</b>	<b>56.94</b>	<b>43.43</b>	<b>68.35</b>	<b>76.19</b>	<b>69.22</b>	<b>65.46</b>
<b>DuQuant-LWC</b>		<b>78.94</b>	<b>57.95</b>	<b>44.11</b>	<b>68.81</b>	<b>76.17</b>	<b>68.98</b>	<b>65.83</b>
LLaMA2-70B W6A6	FP16	81.01	59.68	47.95	75.87	80.87	76.95	70.39
	SmoothQuant	79.87	57.32	45.65	77.13	79.01	74.03	68.84
	OS+	79.33	59.09	47.18	-	79.46	75.06	68.02
	OmniQuant	80.20	60.27	46.84	-	80.55	76.01	68.77
	QLLM	80.63	59.01	45.99	-	79.64	75.37	68.13
	<b>DuQuant</b>	<b>80.96</b>	<b>59.39</b>	<b>47.27</b>	<b>77.34</b>	<b>80.70</b>	<b>76.40</b>	<b>70.34</b>
	<b>DuQuant-LWC</b>	<b>81.18</b>	<b>59.26</b>	<b>47.78</b>	<b>77.86</b>	<b>80.68</b>	<b>76.95</b>	<b>70.62</b>

## E More Ablation Studies

### E.1 Time Speedup and Memory Saving

The current generation of LLMs usually splits into pre-filling and decoding phases and deploys on two separate machines [43]. Here, we present more speedup and memory-saving results for these two phases achieved with the LLaMA2-7B model on a single NVIDIA RTX 3090 GPU. We set the input sequence length to 2048 and the decoding steps to 256. End-to-end results for time speedup and memory savings during the pre-filling stage are shown in Table E11 and E12. We can observe that DuQuant achieves a maximum speedup of  $2.01\times$  during the pre-filling phase, with speedup increasing as the batch size grows. From Table E12, DuQuant demonstrates significant memory savings, effectively reducing memory usage by up to  $3.20\times$  through quantization. For the decoding phase, we enlarge the batch size to 64 and measure speedup along with memory usage for one LLaMA2-7B layer, constrained by the 24 GB memory of the GPU. As shown in Table E13, DuQuant maintains speedup and memory usage comparable to QuaRot. These results underscore the efficiency of DuQuant in optimizing resource utilization, highlighting its potential to enhance performance and reduce costs in deploying large language models, particularly in resource-constrained environments.

Table E11: End-to-end pre-filling speedup on LLaMA2-7B model.

Batch Size	FP16 Time	DuQuant Time	Speedup
1	568ms	294ms	$1.93\times$
2	1003ms	509ms	$1.97\times$
3	1449ms	720ms	$2.01\times$

Table E12: Peak memory usage during pre-filling phase of LLaMA2-7B model.

Batch Size	FP16 Mem.	DuQuant Mem.	Saving Factor
1	15.28GB	4.79GB	$3.20\times$
2	17.94GB	5.94GB	$3.02\times$
3	20.56GB	7.10GB	$2.90\times$

Table E13: Decoding phase results of one LLaMA2-7B layer with a batch size of 64.

Method	Time (ms)	Saving Factor	Memory (GB)	Saving Factor
FP16	659	-	3.550	-
SmoothQuant	437	$1.508x$	1.669	$2.127x$
QLLM	OOM	-	OOM	-
QuaRot	457	$1.442x$	1.678	$2.116x$
DuQuant	499	$1.321x$	1.677	$2.117x$

### E.2 Effects of Rotation Matrix

**Ablation of Rotation Block Size.** To further explore the impact of rotation block size, we apply varying block sizes in the rotation matrices to both LLaMA2-7B and LLaMA2-13B models and evaluate the perplexity of the quantized models. The results, presented in Table E14, indicate that increasing block sizes generally improves model performance. This improvement occurs because larger block sizes allow outliers to be distributed across more channels, evening out values throughout the activation/weight matrix thereby enhancing quantization accuracy and performance. Additionally, quantization runtime decreases with larger block sizes, likely due to more efficient transformations during the reshaping of original activation/weight matrices. Consequently, we adopt 128 as our rotation block size for all experiments for efficiency and effectiveness.

**Ablation of Rotation Times.** Identifying the optimal rotation matrix  $\mathbf{R}$  is a complex challenge, so we employ a greedy search algorithm to approximate the matrix as  $\hat{\mathbf{R}}$ . We conduct an ablation study on the number of greedy steps  $N$  and summarize the results in Table E15. Initially, as  $N$  increases, the model performance improves, reflecting our ability to determine  $\hat{\mathbf{R}}$  more effectively. However, when  $N$  reaches 1024, the model begins to overfit. Consequently, we have chosen  $N = 256$  for all our experiments, as it offers the optimal balance between model performance and time usage.

Table E14: Impact of rotation block size.

Block Size	LLaMA2-7B			LLaMA2-13B		
	WikiText2 ↓	C4 ↓	Time/s	WikiText2 ↓	C4 ↓	Time/s
4	18.69	26.48	64.4	8.81	13.03	97.7
8	10.77	15.04	53.8	7.02	9.68	80.8
16	8.69	11.46	48.2	6.12	8.12	75.2
32	6.96	8.85	48.3	5.61	7.35	76.2
64	6.38	8.07	50.1	5.45	7.13	74.0
128	6.28	7.90	48.6	5.42	7.05	74.0

Table E15: Impact of rotation times.

Rotation Times	LLaMA2-7B			LLaMA2-13B		
	WikiText2 ↓	C4 ↓	Time/s	WikiText2 ↓	C4 ↓	Time/s
1	6.60	8.41	22.9	5.48	7.12	37.7
4	6.34	8.04	22.6	5.41	7.06	38.7
16	6.32	7.98	28.8	5.43	7.05	41.8
64	6.34	7.98	29.0	5.43	7.06	47.0
256	6.28	7.90	48.6	5.42	7.05	74.0
1024	6.31	8.01	129.7	5.46	7.12	179.8

### E.3 Effects of Permutation Algorithm.

As discussed in Section 3.2, rotation transformations within each block are limited and unable to redistribute outliers across different blocks. To address this, we introduce a permutation transformation aimed at balancing outliers more comprehensively. Our primary goal is to minimize the variance among different blocks, as outlined in Eqn. (4). We explore several optimization algorithms, with the results detailed in Table E16. Note that the variance values are measured on activation values of the query project in the first layer of each model, and the time in the table represents the runtime of calibration. The Zigzag permutation notably reduces the variance to  $3.0e-4$ , achieving this with minimal time expenditure and yielding competitive perplexity results. While Simulated Annealing slightly outperforms Zigzag in terms of perplexity for the LLaMA2-7B model, it was significantly more time-consuming, and the marginal gains did not justify the additional complexity. Therefore, we select Zigzag permutation as our preferred method, leading to smoother outlier distribution and more effective quantized models.

Table E16: Impact of channel permutation algorithm.

Permutation Method	LLaMA2-7B				LLaMA2-13B			
	WikiText2 ↓	C4 ↓	Variance	Time/s	WikiText2 ↓	C4 ↓	Variance	Time/s
w.o. Permutation	7.92	10.64	$3.9e-2$	27.5	5.96	7.94	$3.1e-2$	44.7
Random	6.40	8.08	$4.9e-3$	89.5	5.43	7.07	$3.9e-3$	148.6
Simulated Annealing	6.26	7.89	$1.7e-4$	769.6	5.42	7.06	$1.5e-4$	1257.8
Zigzag	6.28	7.90	$3.0e-4$	48.6	5.42	7.05	$2.5e-4$	74.0

### E.4 Effects of Calibration Datasets

**Ablation of Different Calibration Datasets.** We apply our DuQuant to quantize the LLaMA2-7B model using different calibration datasets, with results presented in Table E17. It can be observed that the selection of calibration datasets has a relatively minor impact on quantization performance. This is because our method uses the calibration data solely to identify outlier channels, rather than for gradient-based parameter learning as seen in methods like OmniQuant [47] and AffineQuant [39]. This ablation study underscores the robustness of our DuQuant method.

Table E17: Ablation of calibration datasets.

LLaMA2-7B		Eval.	
		WikiText2 ↓	C4 ↓
Calib.	WikiText2	6.28	7.90
	C4	6.25	7.87

**Calibration-free Quantization.** To further explore the robustness of DuQuant under varying calibration conditions, we generate random calibration data within the vocabulary range of the model, setting the sample count to 256. The results, shown in Table E18, indicate that even in calibration-free settings, our method continues to perform well, achieving results that are competitively close to those obtained with actual calibration data. This demonstrates that DuQuant could provide a viable solution in real-world scenarios where obtaining specific calibration data is challenging or impossible. In addition, our findings suggest that outliers are inherent to certain model layers, reflecting characteristics of the model weights or modules, especially in recent LLMs. This aligns with concurrent research: [66] identified consistent massive outliers at the FFN down projection layer in GLU-based LLMs, such as LLaMA, Mistral, Mixtral, SOLAR, and Gemma, while [42] reported that, although OPT models are sensitive to different calibration sets, newer models demonstrate robustness to outliers and maintain stable activations. These insights reinforce the idea that outliers are more linked to the internal structure of model weights and modules than to calibration data. DuQuant’s ability to deliver high performance without relying on traditional calibration data creates opportunities for deploying quantized models in environments with stringent privacy requirements or limited data availability. This highlights a promising direction for future research, focusing on improving model adaptability and deployment flexibility.

Table E18: Calibration-free quantization, where we generate random data within vocabulary range.

LLaMA2-7B		Eval.	
		WikiText2 ↓	C4 ↓
Calib.	Randomly Generated	6.25	7.86
	WikiText2	6.25	7.87

LLaMA2-13B		Eval.	
		WikiText2 ↓	C4 ↓
Calib.	Randomly Generated	5.45	7.05
	WikiText2	5.44	7.05

**Ablation of Different Numbers of Calibration Samples.** We utilize our DuQuant to quantize the LLaMA2-7B model using varying numbers of calibration samples from the WikiText2 dataset, with results detailed in Table E19. Interestingly, the quantization performance shows a low correlation with the number of samples, demonstrating the robustness of DuQuant. This stability arises because we utilize the mean activation values from these samples to construct our rotation matrices. Since we average the activations, the influence of any single, potentially non-representative sample is minimized, ensuring consistent performance. Notably, as we use mean values, the time cost of our quantization process remains constant regardless of the number of samples, enhancing the efficiency.

Table E19: Ablation of different numbers in the calibration dataset.

# of Samples	WikiText2 ↓	C4 ↓
16	6.29	7.88
32	6.31	7.99
64	6.29	7.88
128	6.28	7.90
256	6.23	7.88



## F Detailed Comparison with QuaRot

Table F20: Comparison of quantization settings between QuaRot and DuQuant.

Setting	Weight	Activation	Query
QuaRot	per-channel symmetric	per-token symmetric	FP16
DuQuant	per-channel asymmetric	per-token asymmetric	per-token asymmetric

Table F21: Evaluation results between QuaRot and DuQuant under DuQuant settings.

Model	Method	WikiText2 ↓	c4 ↓	PIQA	ARC-E	ARC-C	BoolQ	HellaSwag	WinoGrande	Avg ↑
LLaMA1-7B W4A4	FP16	5.68	7.08	77.47	52.48	41.46	73.08	73.00	67.07	64.09
	QuaRot-RTN	7.08	8.73	74.59	48.57	36.01	68.99	65.69	58.56	46.03
	QuaRot-GPTQ	6.44	7.87	76.17	49.96	38.23	70.80	69.29	63.06	61.25
	<b>DuQuant</b>	6.40	7.84	<b>76.44</b>	<b>50.04</b>	<b>38.99</b>	<b>70.98</b>	69.39	<b>64.72</b>	<b>61.76</b>
	<b>DuQuant-LWC</b>	<b>6.18</b>	<b>7.73</b>	76.22	<b>50.04</b>	38.31	70.09	<b>69.82</b>	62.59	61.18
LLaMA2-7B W4A4	FP16	5.47	6.97	76.88	53.54	40.53	71.13	72.96	67.25	63.72
	QuaRot-RTN	9.66	11.98	69.48	46.25	32.76	64.80	60.75	56.67	44.04
	QuaRot-GPTQ	6.39	8.15	75.15	49.15	36.68	67.89	68.87	61.33	59.85
	<b>DuQuant</b>	6.28	7.90	75.24	<b>51.89</b>	36.77	67.86	69.54	62.12	60.57
	<b>DuQuant-LWC</b>	<b>6.08</b>	<b>7.79</b>	<b>75.68</b>	50.00	<b>37.46</b>	<b>69.24</b>	<b>69.74</b>	<b>63.93</b>	<b>61.01</b>
LLaMA3-8B W4A4	FP16	6.14	8.88	80.85	77.78	53.41	81.28	79.16	72.84	74.22
	QuaRot-RTN	13.89	17.59	69.64	57.58	34.56	66.76	63.46	62.75	59.13
	QuaRot-GPTQ	8.69	12.40	74.54	67.38	40.61	70.43	70.47	65.11	64.76
	<b>DuQuant</b>	8.53	12.01	76.93	<b>70.88</b>	<b>45.05</b>	<b>74.59</b>	73.17	66.14	<b>67.79</b>
	<b>DuQuant-LWC</b>	<b>8.06</b>	<b>11.29</b>	<b>76.22</b>	70.41	43.69	74.34	<b>73.87</b>	<b>67.80</b>	67.72

Table F22: Evaluation results between QuaRot and DuQuant under QuaRot settings.

Model	Method	WikiText2 ↓	C4 ↓	PIQA	WinoGrande	HellaSwag	ARC-E	ARC-C	LAMBADA	Avg ↑
LLaMA2-7B W4A4 QuaRot Setting	FP16	5.47	6.97	79.11	69.06	75.99	74.58	46.25	73.90	69.82
	QuaRot-RTN	8.37	-	72.09	60.69	65.4	58.88	35.24	57.27	58.26
	QuaRot-GPTQ	6.1	-	76.77	63.77	72.16	69.87	40.87	<b>70.39</b>	65.64
	<b>DuQuant</b>	6.23	7.91	76.28	66.93	72.96	69.99	40.53	69.61	66.05
	<b>DuQuant-LWC</b>	<b>6.01</b>	<b>7.67</b>	<b>77.64</b>	<b>67.8</b>	<b>72.97</b>	<b>70.37</b>	<b>41.81</b>	69.53	<b>66.69</b>
LLaMA2-13B W4A4 QuaRot Setting	FP16	4.88	6.46	80.47	72.22	79.39	77.48	49.23	76.75	72.59
	QuaRot-RTN	6.09	-	77.37	67.32	73.11	70.83	43.69	70.66	67.16
	QuaRot-GPTQ	5.4	-	78.89	70.24	76.37	72.98	46.59	73.67	69.79
	<b>DuQuant</b>	5.39	7.05	78.51	<b>70.88</b>	76.80	<b>74.62</b>	<b>48.21</b>	73.92	<b>70.49</b>
	<b>DuQuant-LWC</b>	<b>5.27</b>	<b>6.93</b>	<b>78.73</b>	<b>70.88</b>	<b>77.20</b>	74.07	47.27	<b>73.96</b>	70.35

Table F23: Matrices comparison between DuQuant and QuaRot under W4A4 quantization.

Model	LLaMA2-7B		LLaMA2-13B	
Dataset	WikiText2 ↓	C4 ↓	WikiText2 ↓	C4 ↓
QuaRot	9.66	11.98	6.73	8.69
<b>DuQuant</b>	<b>7.92</b>	<b>10.64</b>	<b>5.96</b>	<b>7.94</b>

Table F24: Quantization runtime comparison on a single NVIDIA A100 80G GPU.

Model	LLaMA2-7B	LLaMA2-13B	LLaMA2-70B
QuaRot	20min	36min	5.1h
<b>DuQuant</b>	50s	71s	270s

In this section, we present a detailed comparison between our DuQuant and QuaRot [2]. QuaRot employs Hadamard matrices to mitigate outliers in activations and utilizes the GPTQ algorithm for weight quantization to achieve competitive performance. However, our DuQuant method demonstrates several distinct advantages:

- **Effective Use of Prior Knowledge:** DuQuant leverages prior knowledge to accurately target and eliminate outliers through multiple rotations, achieving a smoother activation distribution compared to the Hadamard transformation, as demonstrated in Figure 10.

- **Efficient Channel Permutation:** Our channel permutation not only further smooths outlier features but also benefits from rapid implementation, enhancing overall performance.
- **Simultaneous Weight Matrix Smoothing:** Unlike QuaRot, DuQuant directly and efficiently smooths the weight matrix, avoiding the time-consuming GPTQ algorithm and accelerating the quantization process, as demonstrated high quantization efficiency in Table F24.

Experimental results underscore the superiority of DuQuant over QuaRot. We first summarize the experimental setting differences between the original paper of QuaRot with ours in Table F20\*. For a fair comparison, we reproduce the QuaRot under 4-bit per-channel weight and per-token activation asymmetric quantization. Table F21 displays the perplexity (PPL) and zero-shot accuracy for models LLaMA1-7B, LLaMA2-7B, and LLaMA3-8B. Our DuQuant method consistently outperforms QuaRot-RTN across all benchmarks, showcasing our advanced weight matrix management. Furthermore, compared to QuaRot-GPTQ, DuQuant and DuQuant+LWC achieve better average accuracy across six QA tasks and demonstrate superior performance on the WikiText and C4 datasets, particularly for LLaMA3-8B. Moreover, we further provide the results of DuQuant under the setting utilized in the original paper of QuaRot in Table F22. DuQuant still surpasses QuaRot by a large margin.

Additionally, we assess the effectiveness of the rotation matrices utilized in DuQuant which incorporate prior knowledge against the Hadamard matrices used in QuaRot. We omit the permutation step in DuQuant and directly contrast it with QuaRot-RTN. Results in Table F23 show that our DuQuant without permutation outperforms QuaRot by a clear margin, which confirms that our rotation transformation is more effective than Hadamard by leveraging prior knowledge. It is worth noting that because a Hadamard matrix is orthogonal and symmetric, it multiplies by itself to yield the identity matrix. In other words, the Hadamard matrix is not suitable for greedy searches aimed at finding smaller outliers. These findings differentiate DuQuant from QuaRot and highlight the effectiveness of our approach in managing outliers for post-training quantization of large language models.

## G Algorithm for Rotation Matrix

---

### Algorithm 1 Construction of the Rotation Matrix

---

**Input:** Pre-initialized rotation matrix  $\tilde{\mathbf{R}}$ , greedy search steps  $N$ , activation matrix  $\mathbf{X}$  with shape of  $[T, C_{in}]$   
**Output:** Rotation matrix  $\hat{\mathbf{R}}$   
**function** get\_rotation\_matrix ( $\mathbf{X}, \tilde{\mathbf{R}}, N$ )  
1:  $T, C_{in} = \mathbf{X}.shape$   
2:  $\mathbf{R} = \text{eye}(C_{in})$   
3:  $a = \max_{i,j} |\mathbf{X}_{ij}|$   
4: **for**  $k$  in  $1, \dots, N$  **do**  
5:     channel\_max =  $\mathbf{X}.abs().\max(\text{dim} = 0).values$   
6:     outlier\_channel =  $\arg \max(\text{channel\_max})$   
7:     Obtain randomly initialized orthogonal matrix  $\mathbf{Q}'$  with shape of  $[C_{in} - 1, C_{in} - 1]$   
8:      $\mathbf{Q}' = \text{concat}([\text{zeros}(C_{in} - 1, 1), \mathbf{Q}'], \text{dim}=1)$   
9:      $\mathbf{Q} = \text{concat}([\text{zeros}(1, C_{in}), \mathbf{Q}'], \text{dim}=0)$   
10:      $\mathbf{Q}[0, 0] = 1$   
11:      $\mathbf{R}' = \text{matmul}(\tilde{\mathbf{R}}, \mathbf{Q})$   
12:      $\mathbf{R}'[:, \text{outlier\_channel}], \mathbf{R}'[:, 0] = \mathbf{R}'[:, 0], \mathbf{R}'[:, \text{outlier\_channel}]$   
13:      $\mathbf{R}'[\text{outlier\_channel}, :], \mathbf{R}'[0, :] = \mathbf{R}'[0, :], \mathbf{R}'[\text{outlier\_channel}]$   
14:      $\mathbf{R} = \text{matmul}(\mathbf{R}, \mathbf{R}')$   
15:      $\mathbf{X} = \text{matmul}(\mathbf{X}, \mathbf{R}')$   
16:     **if**  $\max_{i,j} |\mathbf{X}_{ij}| < a$  **then**  
17:          $\hat{\mathbf{R}} = \mathbf{R}$   
18:          $a = \max_{i,j} |\mathbf{X}_{ij}|$   
19:     **end if**  
20: **end for**  
21: **return**  $\hat{\mathbf{R}}$

---

\*Following prior works [47, 34, 39], we describe operations on Query, Key, and Value states as per-tensor quantization for consistency. To be precise, these operations are effectively applied on a per-head basis.

## H Limitations and Broader Impacts

**Limitations.** The primary limitation of our method is the lack of a specialized strategy for calibration data selection. We adhere to established practices [47, 39, 34, 74, 2, 1] by randomly selecting 128 samples from the WikiText2 dataset to compute the mean embeddings that inform our rotation matrix and zigzag permutation order. We also explore the possibility of calibration-free quantization and show some promising results. However, further investigating more tailored choices for calibration data can potentially enhance the performance of our quantized models. We leave this for future study.

**Broader Impacts.** Our work identifies the presence of massive outliers in down-projection layer of FFN modules, which significantly complicates low-bit weight-activation quantization. To address this challenge, we implement a combination of rotation and permutation matrices to effectively smooth both massive and uniform outliers, proving both fast and effective. Consequently, we establish a new state-of-the-art for INT4 weight-activation post-training quantization methods. Our approach aims to accelerate large language models and reduce memory usage during deployment, offering substantial benefits to the field of LLM research. These advancements could lead to more efficient and accessible LLM applications, facilitating broader usage and enabling more sustainable AI implementations.

## I More Visualizations

We provide additional visualizations of normal and massive outliers in various models (LLaMA1, LLaMA2, Vicuna-v1.5, Mistral) from Figure I1 to Figure I9. In each figure except Mistral, the left side illustrates changes in normal outliers before and after applying our rotation and permutation transformations, while the right side shows the changes in massive outliers before and after transformations. It is evident that massive outliers consistently occur in the down-projection layer of the FFN module across all models, supporting our findings discussed in Section 2. Conversely, normal outliers appear in different modules within the transformation block. For instance, Figure I3 shows normal outliers at the up-projection layer of the FFN module in LLaMA1-13B. Significantly, both massive and normal outliers are reduced markedly after our rotation and permutation transformations, leading to easier quantization of activations. This underscores the effectiveness of our DuQuant in managing outlier features across diverse LLM models.

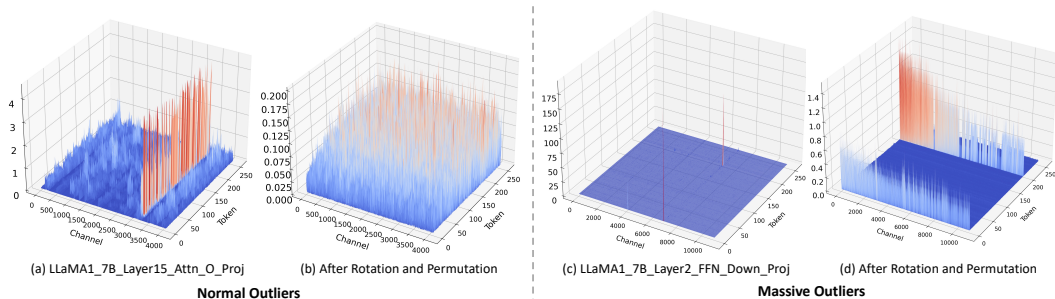


Figure I1: Activation change with the use of our DuQuant for LLaMA1-7B.

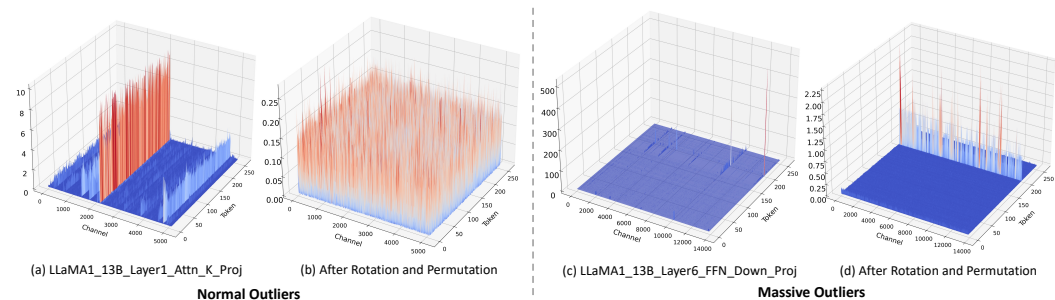


Figure I2: Activation change with the use of our DuQuant for LLaMA1-13B.

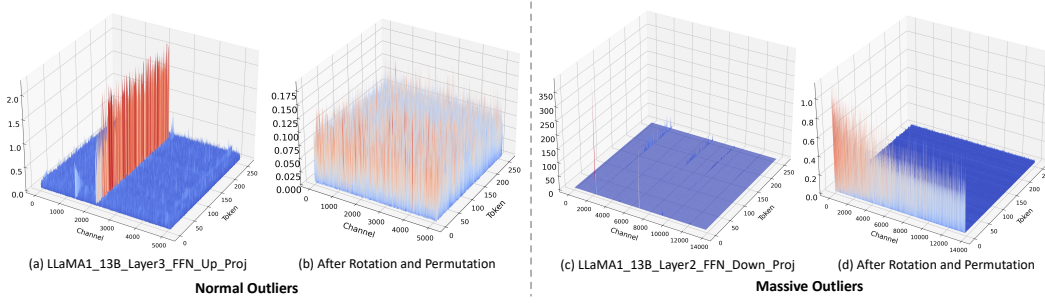


Figure I3: More examples of activation change with the use of our DuQuant for LLaMA1-13B.

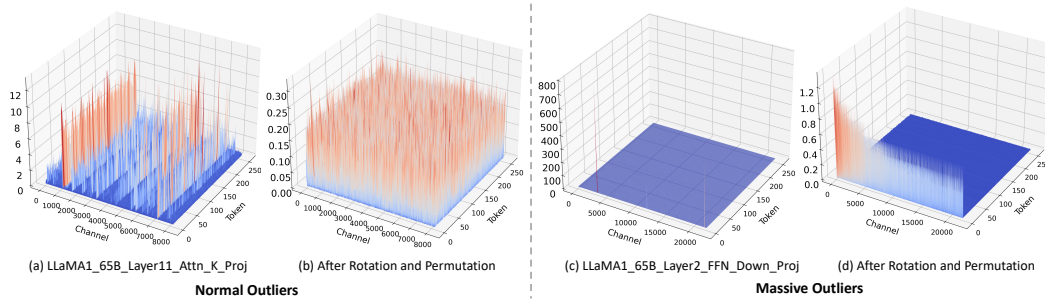


Figure I4: Activation change with the use of our DuQuant for LLaMA1-65B.

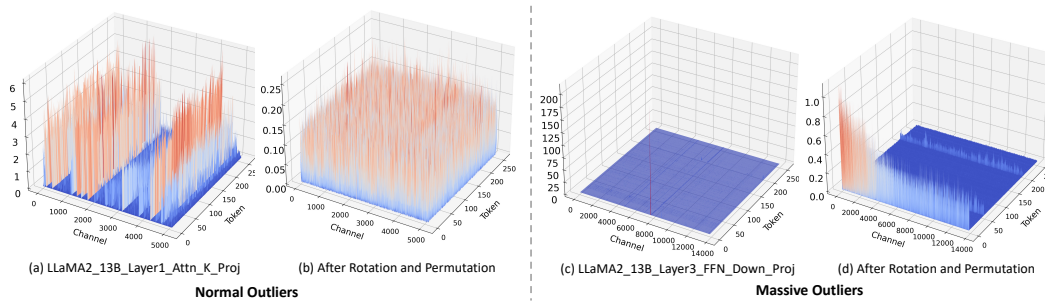


Figure I5: Activation change with the use of our DuQuant for LLaMA2-13B.

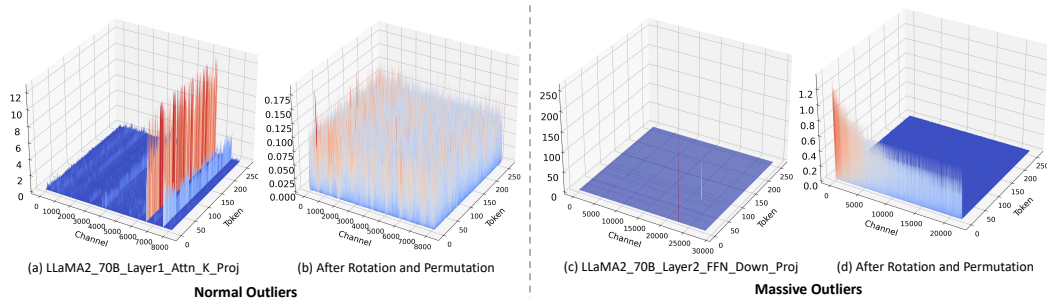


Figure I6: Activation change with the use of our DuQuant for LLaMA2-70B.

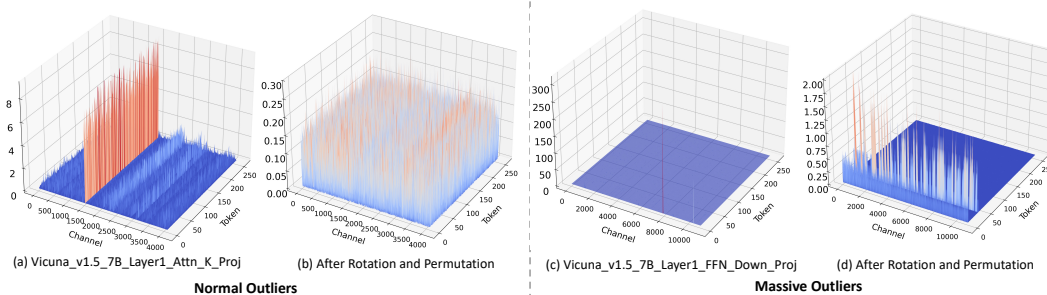


Figure I7: Activation change with the use of our DuQuant for Vicuna-v1.5-7B.

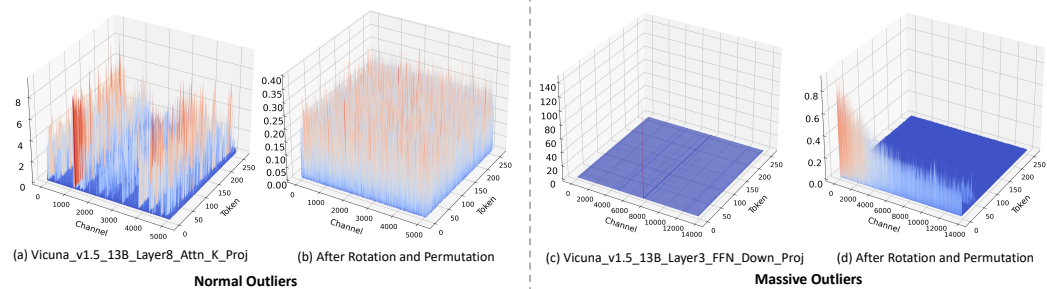


Figure I8: Activation change with the use of our DuQuant for Vicuna-v1.5-13B.

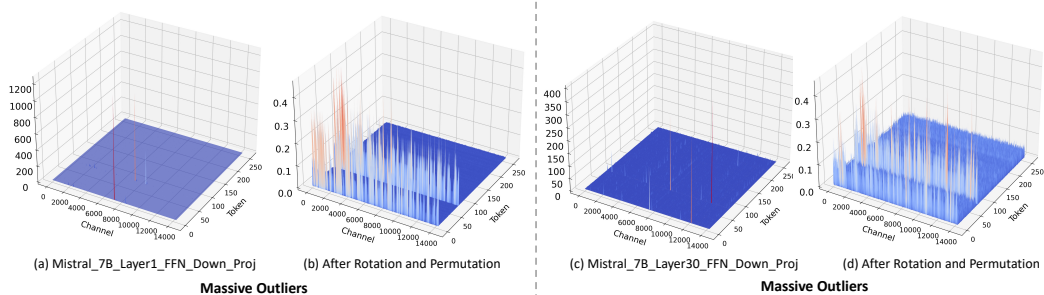


Figure I9: Massive activation change with the use of our DuQuant for Mistral7B.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We state the contributions in the abstract and introduction sections.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the main limitations in Appendix H.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Please refer to Appendix B for related assumptions and proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Implementation details are provided in Appendix C. We have released our code at <https://github.com/Hsu1023/DuQuant>.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We have released our code at <https://github.com/Hsu1023/DuQuant>.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Please refer to Section 4 and Appendix C for experimental settings.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: Our study used a fixed seed for all quantization operations, following standards in post-training quantization, and thus did not report statistical significance.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).



- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: All experiments are done on NVIDIA RTX 3090 GPUs and A100 GPUs.

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

Answer: [Yes]

Justification: We confirm the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We discuss the broader impacts of our work in Appendix H.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This work poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cite all datasets and models utilized in our experiments.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This work does not release new assets.

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This work does not involve crowdsourcing nor research with human subjects.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

#### 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This work does not involve crowdsourcing nor research with human subjects.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.