

# BASE-Q: BIAS AND ASYMMETRIC SCALING ENHANCED ROTATIONAL QUANTIZATION FOR LARGE LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Rotation-based methods have become essential for state-of-the-art LLM quantization by effectively mitigating outliers in weights and activations. Current approaches predominantly focus on optimizing the global rotation matrix to achieve marginal accuracy improvements—a strategy that incurs prohibitive computational costs through full-model backpropagation while offering limited practical utility. We fundamentally reassess this optimization paradigm and identify two critical error sources that persist even under optimal rotation conditions: (i) channel mean misalignment, which amplifies rounding errors during quantization, and (ii) clipping-induced energy loss, which is exacerbated by the rotation-induced Gaussian-like distributions. Our analysis reveals that directly addressing these issues offers a more effective path to achieving high quantization accuracy. Based on these insights, we introduce **BASE-Q**, a lightweight quantization framework that circumvents expensive global rotation learning. **BASE-Q** employs simple yet powerful transformer-block-wise correction strategy: **bias correction** to eliminate channel mean variance and **asymmetric scaling** to compensate for clipping-induced energy loss. This blockwise strategy drastically reduces optimization overhead, enabling efficient quantization of 70B parameter models on a single GPU. Extensive experiments across diverse LLMs and benchmarks validate the effectiveness of BASE-Q, narrowing the accuracy gap to full-precision models by 50.5%, 42.9%, and 29.2% compared to previous rotation method QuaRot, SpinQuant, and OSTQuant respectively, demonstrating the superiority of our lightweight paradigm.

## 1 INTRODUCTION

Large language models (LLMs) (Touvron et al., 2023; Bai & et al., 2023; Dubey & et al., 2024; Touvron & et al., 2023; DeepSeek-AI, 2024) drive advances across diverse natural language tasks, but their ever-increasing scales present significant challenges for efficient deployment, particularly regarding inference latency and memory consumption. Low-precision quantization (Nagel et al., 2020; Li et al., 2021; Liu et al., 2022; 2025) (e.g., 4 bits or lower) represents a crucial technique for addressing these computational bottlenecks. However, aggressively reducing numerical precision to such low bit-widths results in substantial accuracy degradation. The fundamental challenge underlying this degradation stems from the presence of significant outliers in both activations and weights (Wei et al., 2023; Dettmers et al., 2022), which necessitate wider quantization ranges that consequently amplify quantization errors.

To mitigate the adverse effects of outliers in LLM quantization, recent approaches have leveraged linear equivalent transformations, particularly scaling and rotation operations. SmoothQuant (Xiao et al., 2022) pioneered this direction by introducing channel-wise scaling to reduce activation variance, effectively smoothing the distribution and achieving robust improvements in INT8 post-training quantization (PTQ) across various LLMs. Building on this direction, QuaRot (Ashkboos et al., 2024) introduces Hadamard rotations as another form of equivalent transformation to redistribute outliers and reduce channel-wise disparities, enabling robust INT4 PTQ performance. Since RMSNorm in LLMs normalizes activation energy across residual blocks while preserving this en-

054 energy under rotations, identical rotation matrix can be shared across all transformer blocks and fused  
055 into adjacent linear layers without incurring additional computation overhead.

056 Building upon the effectiveness of fixed rotations, a prevailing trend in subsequent research (e.g.,  
057 SpinQuant (Liu et al., 2024), OSTQuant (Hu et al., 2025)) has pursued further accuracy improve-  
058 ments by learning a single, globally shared rotation matrix. We argue that this paradigm, while  
059 seemingly intuitive, is fundamentally flawed for two critical reasons. Firstly, it incurs prohibitive  
060 optimization costs due to full-model backpropagation, rendering it impractical for truly large-scale  
061 models. More importantly, we find that the core assumption—that an ‘optimal’ global rotation is the  
062 key to minimizing quantization error—is misguided. A single transformation shared across dozens  
063 of diverse transformer blocks inherently lacks the expressiveness to simultaneously suppress out-  
064 liers and accommodate block-specific distributional variations. This theoretical limitation explains  
065 the observed marginal improvements and necessitates a paradigm shift away from global optimiza-  
066 tion.

067 Motivated by this insight, we shift the focus from optimizing the rotation itself to directly correcting  
068 the residual errors it leaves behind. Our analysis identifies two dominant, yet previously overlooked,  
069 error sources: **(i) Rounding Error from Channel Mean Misalignment:** We reveal that significant  
070 variance persists among channel-wise means even after rotation, as a single global rotation can-  
071 not simultaneously nullify this variance across all layers. We address this through a lightweight,  
072 block-specific bias correction that precisely realigns channel means with negligible computational  
073 overhead. **(ii) Clipping Error from Distributional Shift:** We uncover that rotation, by render-  
074 ing activations more Gaussian-like, paradoxically increases the energy loss from quantization clip-  
075 ping, thereby violating the theoretical equivalence that rotation aims to preserve. We counteract this  
076 through a block-aware asymmetric scaling strategy that restores signal energy post-clipping. By cir-  
077 cumventing the costly and ineffective pursuit of a perfect global rotation, our framework, BASE-Q,  
078 employs these simple yet targeted corrections at the block level, achieving superior accuracy while  
079 enabling efficient and memory-friendly optimization.

080 Our main contributions are:

- 081 • **Theoretical Insight.** We provide a systematic analysis revealing key components of quan-  
082 tization error after rotation, elucidating why existing approaches plateau and identifying  
083 opportunities for further optimization.
- 084 • **Efficient Transformer-Block-wise Correction.** We propose BASE-Q, a novel quantiza-  
085 tion framework that employs blockwise channel bias correction and asymmetric scaling  
086 under fixed rotations, achieving significant performance enhancement without the expen-  
087 sive full-model optimization.
- 088 • **Extensive Empirical Validation.** Extensive experiments across diverse LLMs and bench-  
089 marks validate the effectiveness of BASE-Q, narrowing the accuracy gap to full-precision  
090 models by 50.5%, 42.9%, and 29.2% compared to previous rotation method QuaRot (Ashk-  
091 boos et al., 2024), SpinQuant (Liu et al., 2024), and OSTQuant (Hu et al., 2025) respec-  
092 tively.

## 093 2 RELATED WORK

094 **Equivalent Transformations in LLM Quantization.** Post-training quantization (PTQ) has at-  
095 tracted considerable attention for the efficient deployment of LLMs, yet remains challenging due to  
096 frequent outliers in both activations and weights. AWQ (Lin et al., 2024b) introduced channel-wise  
097 scaling for weight-only PTQ, while SmoothQuant (Xiao et al., 2022) employed rescaling for activa-  
098 tions and weights to suppress outlier effects and enable robust INT8 quantization. OmniQuant (Shao  
099 et al., 2023) extended this concept by introducing learnable scaling coefficients for each submodule,  
100 allowing for finer-grained adaptation across network components. AffineQuant (Ma et al., 2024)  
101 further generalized these ideas by employing learnable affine transformations to jointly align mean  
102 and variance before quantization. Beyond scaling-based methods, QuIP (Chee et al., 2023; Tseng  
103 et al., 2024) first applied rotation transformations for weight-only PTQ. QuaRot (Ashkboos et al.,  
104 2024) proposed applying Hadamard rotations to both activations and weights, rendering distribu-  
105 tions more Gaussian and further suppressing outliers, thus simplifying the quantization process.  
106 DuQuant (Lin et al., 2024a) employed rotation and permutation to more effectively eliminate out-  
107

108 liers. DFRot (Xiang & Zhang, 2024) attributes the success of Hadamard transforms to their han-  
 109 dling of rare tokens with massive activations and proposes a weighted loss function combined with  
 110 an alternating Procrustes-based optimization to learn a globally improved rotation matrix. Spin-  
 111 Quant (Liu et al., 2024) advanced this direction by learning optimal rotation matrices from cali-  
 112 bration data, achieving lower quantization errors at the cost of greater computational and memory  
 113 requirements. OSTQuant (Hu et al., 2025) unified learnable rotations and scaling within a single  
 114 framework, providing additional flexibility and consistently outperforming previous methods on  
 115 various LLM benchmarks. FlatQuant (Sun et al., 2024) employed layer-wise learned online matrix  
 116 transforms to improve quantized linear layers, although at the cost of increased inference overhead  
 117 and parameter count.

### 118 3 ERROR ANALYSIS OF ROTATION QUANTIZATION

119 In this section, we theoretically analyze the errors derived by rotation-based quantization, specif-  
 120 ically focusing on the Hadamard transformation. We begin by investigating the error reduction  
 121 mechanism inherent in this transformation, deriving quantitative expressions for both rounding and  
 122 clipping errors in Section 3.1. Subsequently, in Sections 3.2 and 3.3, we provide an in-depth analysis  
 123 of these error sources, uncovering potential components that are amenable to further optimization  
 124 and providing key insights for methodological improvements. Finally, drawing inspiration from  
 125 OSTQuant (Hu et al., 2025), which demonstrates complementary effects of scaling and rotation, we  
 126 interpret and substantiate this synergy through our error decomposition framework in Section 3.4.

127 **Notation.** Throughout this paper, we follow a consistent notational convention. We denote matrices  
 128 by uppercase bold letters (e.g.,  $\mathbf{W}$ ), vectors by lowercase bold letters (e.g.,  $\mathbf{x}$ ), and scalars by regular  
 129 lowercase letters (e.g.,  $s$ ). All vectors are considered column vectors unless specified otherwise.

#### 130 3.1 QUANTIZATION ERROR ANALYSIS FOR HADAMARD ROTATIONS

131 The key challenge in LLM quantization is the presence of extreme outliers in the activation channels,  
 132 which significantly increases the dynamic range and severely impedes quantization performance.  
 133 To elucidate how Hadamard rotations mitigate this issue, consider a token activation  $\mathbf{x} \in \mathbb{R}^n$  with  
 134 outlier values structured as follows:

$$135 \mathbf{x} = \mathbf{g} + \sum_{i \in \mathcal{O}} a_i \delta \mathbf{e}_i, \quad (1)$$

136 where  $\mathbf{g} \sim \mathcal{N}(\mu, \delta^2 \mathbf{I})$  represents the ‘main mass’ as a Gaussian component, and each outlier channel  
 137  $i \in \mathcal{O}$  is modeled as a one-hot vector  $\mathbf{e}_i$  of amplitude  $a_i \delta$  with  $a_i \gg 1$  and  $|\mathcal{O}| \ll n$ . When an  
 138 orthogonal Hadamard matrix  $\mathbf{H}$  is applied (where each entry is  $\pm n^{-\frac{1}{2}}$ ), the transformed activation  
 139 becomes:

$$140 \mathbf{Hx} = \mathbf{Hg} + \sum_{i \in \mathcal{O}} a_i \delta \mathbf{H} \mathbf{e}_i. \quad (2)$$

141 Since Gaussian distributions are invariant under orthogonal transformations, the main mass  $\mathbf{Hg}$   
 142 remains Gaussian. Critically, each outlier term  $a_i \delta \mathbf{e}_i$  is now assigned to the  $i$ -th column of  $\mathbf{H}$ , a  
 143 dense vector whose large amplitude  $a_i \delta$  is evenly distributed across all channels, with each channel  
 144 receiving only  $a_i \delta / \sqrt{n}$ . For large  $n$ , the per-channel impact of any outlier is significantly diminished  
 145 effectively, outliers are ‘absorbed’ into the Gaussian bulk, yielding a distribution that is far more  
 146 conducive to quantization.

147 This mechanism extends beyond simple Gaussian activations to more general cases where acti-  
 148 vations exhibit correlated structures or multiple modes (e.g., Gaussian mixtures or channel-wise  
 149 mean/variance discrepancies). In such cases, applying a principal component transformation to di-  
 150 agonalize the covariance matrix, followed by a Hadamard rotation, further balances the marginal  
 151 variances and minimizes the prominence of outliers:

$$152 \mathbf{x}' = \mathbf{HU}^T (\mathbf{x} - \mathbb{E}[\mathbf{x}]), \quad (3)$$

153 where  $\mathbf{U}$  is the matrix whose columns are the principal component vectors. To quantify the effect  
 154 on quantization, consider uniform quantization with per-channel clipping:

$$155 \mathbf{x}_q = F_{clip}(F_{round}(\frac{\mathbf{x} - z}{s}), 0, 2^b - 1) * s + z, \quad (4)$$

where  $z$  and  $s$  denote the quantization lower bound and step size, respectively, and  $b$  represents the target bit-width. The overall quantization error comprises two components:

- **Rounding error**, with expected  $\ell_2$  energy per channel:  $\mathbb{E}[\|\varepsilon_{\text{round}}\|_2^2] = \frac{\Delta^2}{12}$ , where  $\Delta$  denotes the quantization step size (bin width). This classical result stems from modeling the rounding error as uniform noise, a concept we will elaborate on in Section 3.4 (see Eq. 10).
- **Clipping error**, defined by the mass outside quantization bounds:  $\mathbb{E}[\|\varepsilon_{\text{clip}}\|_2^2] = \int_{-\infty}^z x^2 P(x) dx + \int_{z+\Delta}^{+\infty} x^2 P(x) dx$ , where  $\Delta = s(2^b - 1)$  and  $P(x)$  is the empirical channel distribution.

By dispersing outlier energy and compressing the activation dynamic range, Hadamard transformations significantly reduce rounding errors by enabling smaller quantization steps  $s$ .

### 3.2 ROUNDING ERRORS FROM MISALIGNED CHANNEL MEANS

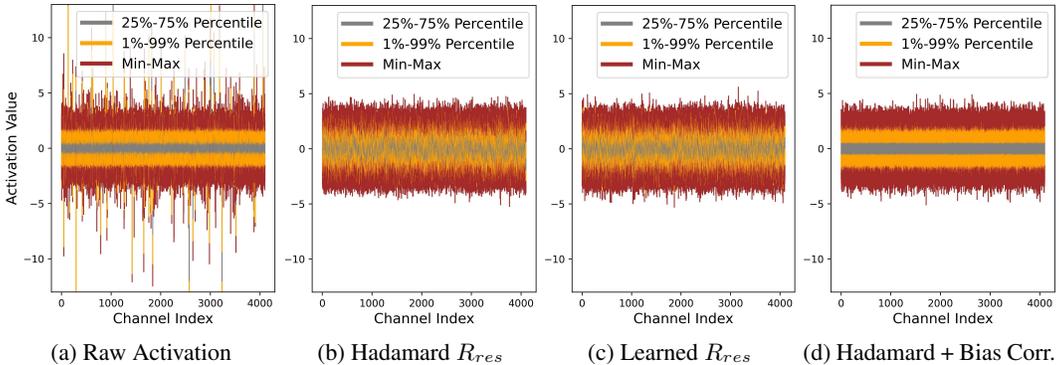


Figure 1: (a) Raw activation distribution of the first MLP block in Llama3-8B. (b) Hadamard rotation suppresses outliers but leaves residual inter-channel misalignment. (c) Learned rotations also fail to address this misalignment. (d) Our bias correction eliminates most inter-channel mean variability, thereby reducing  $Var(\mu_j)$ , which is the main rounding error component as formalized in Equations (5) and (6). Additional visualizations demonstrating this effect are provided in Section D.

As demonstrated in the previous section, the energy of the rounding error is proportional to the quantization scale,  $s = \frac{\Delta}{2^b - 1}$ . As rotations gaussianize the activation distribution, the quantization range can be considered proportional to the standard deviation of activations,  $\sigma$ . Therefore, the expected rounding error satisfies:

$$\mathbb{E}[\|\varepsilon_{\text{round}}\|_2^2] \propto s^2 \propto \sigma^2. \tag{5}$$

The total variance  $\sigma^2$  of activations across all channels can be formally decomposed using the Law of Total Variance, which states:

$$\sigma^2 = \frac{1}{n} \sum_{j=1}^n \sigma_j^2 + Var(\mu_j), \tag{6}$$

where  $\sigma_j^2$  and  $\mu_j$  denote the variance and mean of the  $j$ -th channel, respectively. The first term, representing the average channel variance, remains invariant under orthogonal transformations and is determined by the data distribution. Ideally, an optimal rotation would align all channel means, reducing  $Var(\mu_j)$  to zero. However, we observe that even globally learned rotations (e.g., in SpinQuant) fail to achieve this alignment, yielding only marginal improvements over fixed rotations in this regard (see Figure 1(c)). This limitation motivates our direct approach. As demonstrated in Figure 1(d) and Table 5, our blockwise bias correction effectively eliminates  $Var(\mu_j)$ , which renders the choice of the global rotation matrix  $R_{res}$  insignificant. Consequently, models with different  $R_{res}$  choices converge to consistently high accuracy. This stands in stark contrast to prior work, where learning  $R_{res}$  is crucial to performance stability. Ultimately, our framework fully subsumes

the optimization potential of a learnable global rotation, allowing its computationally expensive optimization to be entirely removed from the quantization pipeline.

Ideally, optimizing rotation should align all channel means to achieve  $Var(\mu_j) = 0$ . In practice, however, the LLM shares a single rotation  $\mathbf{R}_{res}$  across all transformer blocks, despite blockwise variation in bias. The limited expressiveness of a global rotation renders it infeasible to simultaneously eliminate outliers and achieve mean alignment. As a result, globally learned rotations (e.g., SpinQuant) can only minimize  $Var(\mu_j)$  in a least-loss sense but cannot completely eliminate this error arising from mean discrepancies. This limitation motivates our explicit bias correction strategy, which precisely cancels the variance-of-means term post-rotation. Moreover, this approach supports blockwise optimization, circumventing the computational burden of full-model optimization.

### 3.3 CLIPPING ELIMINATES NON-NEGLIGIBLE ENERGY

Orthogonal rotation, such as the Hadamard transformation, significantly alters the distribution of activations. While raw activations typically exhibit a heavy-tailed distribution with sparse outliers, post-rotation activations approximate a Gaussian distribution. This transformation has critical implications: unlike heavy-tailed distributions where extreme values are rare, the Gaussian shape concentrates a larger proportion of activation energy towards the distribution tails. Consequently, applying a fixed clipping threshold removes a non-negligible amount of energy from values exceeding the threshold.

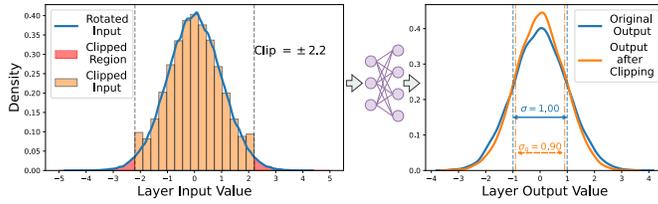
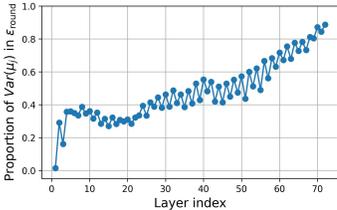


Figure 2: Misalignment causes up to 85% layer-wise rounding error in Qwen2.5-3B. Figure 3: (Left) MSE-optimal clipping on rotated activations results in a 18.4% loss of layer’s input energy. (Right) This induces significant discrepancies in the layer output.

Formally, consider activations  $x \sim \mathcal{N}(\mu, \sigma^2)$  quantized to  $b = 4$  bits (INT4), the minimum MSE-optimal clipping threshold is commonly determined as:

$$\theta^* = \arg \min_{\theta} \mathbb{E}[\|x - (F_{clip}(F_{round}(\frac{x - \theta}{2\theta/(2^b - 1)}), 0, 2^b - 1) + \theta)\|_2] \approx 2.2\sigma \quad (7)$$

Here,  $\sigma$  denotes the standard deviation of the post-rotation activations, which are approximately Gaussian. The factor  $\theta^* \approx 2.2\sigma$  is the empirically found optimal clipping threshold coefficient for minimizing. At this threshold, we can quantify the proportion of total activation energy contained within the clipped regions (i.e., the energy loss). For a zero-mean Gaussian distribution, the expected  $L_2$  energy of the clipped values constitutes approximately 18.4% of the total energy:

$$\mathbb{E}[\|\varepsilon_{clip}\|_2] = \int_{-\infty}^{-2.2\sigma} x^2 P(x) dx + \int_{+2.2\sigma}^{+\infty} x^2 P(x) dx \approx 18.4\%, \quad (8)$$

where  $P(x)$  denotes the probability density function. 18.4% of the energy is **non-negligible** in the context of deep neural networks. When this clipping is applied layer throughout a deep network, the accumulated energy loss becomes substantial. The empirical consequence of this breakdown is illustrated in Figure 3, where the distribution of the layer output significantly mismatches the original. Moreover, in practical LLM architectures, neither globally nor per-layer optimized clipping bounds can fully resolve this energy loss, as joint optimization across all layers is intractable. This limitation explains the diminishing returns observed in aggressive clipping strategies for rotation-based quantization schemes, even with parameter tuning.

To address this issue, we introduce asymmetric scaling for each activation quantization step, which restores the energy loss and maintains accurate signal magnitude, thereby improving both theoretical fidelity and empirical performance.

### 3.4 ROLE OF SCALING IN ROTATIONAL QUANTIZATION

Combining scaling with orthogonal transformations (rotations) has been shown to be crucial for enhancing quantization performance. To understand this synergy, we analyze the effect of quantization noise from the perspective of the original, pre-rotation space. In the rotated space, the rounding noise introduced by a uniform quantizer can be modeled as additive noise  $\epsilon$  where each component is an independent random variable uniformly distributed over the quantization bin:

$$\epsilon \sim \mathcal{U}\left(-\frac{\Delta}{2}, \frac{\Delta}{2}\right), \quad (9)$$

where  $\Delta$  denote the rounding noise and quantization bin width. The variance of this noise is  $\frac{\Delta^2}{12}$ . Since orthogonal transform preserves the variance, this noise manifests as additive isotropic Gaussian noise in the original space, which becomes:

$$\epsilon \sim \mathcal{N}\left(0, \frac{\Delta^2}{12}\right). \quad (10)$$

Referring to the matrix multiplication of weights and activations in the specific layer, the layer output includes terms arising from error propagation:

$$(w + \epsilon_w) \cdot (a + \epsilon_a) = w \cdot a + (w \cdot \epsilon_a + \epsilon_w \cdot a) + \epsilon_w \cdot \epsilon_a, \quad (11)$$

where  $w$  and  $a$  represent the original weights and activations respectively, while  $\epsilon_a$  and  $\epsilon_w$  denote their corresponding quantization errors.

Introducing per-channel symmetric scaling factors  $s$  (applied to  $w$ ) and  $1/s$  (applied to  $a$ ) allows balancing the propagated error energy between these terms. The total variance of the noisy product, dominated by  $w \cdot \epsilon_a$  and  $\epsilon_w \cdot a$ , achieves its minimum when their variances are equal. By the arithmetic-geometric mean (AM-GM) inequality, the optimal scaling factor satisfies:

$$s^2 = \frac{\mathbb{E}[|w|_2]}{\mathbb{E}[|a|_2]}. \quad (12)$$

This choice of symmetric scaling enables optimal allocation of quantization noise energy. Unfortunately, symmetric scaling breaks this per-token energy, making it incompatible with fusible rotations for layers immediately following an RMSNorm (i.e.,  $Q$ ,  $K$ ,  $V$ ,  $Up$ ,  $Gate$  projections). It can only be effectively paired with fusible rotation in layers like the  $O$  and  $Down$  projections.

In summary, our analysis in this section has deconstructed the residual errors in rotational quantization, identifying two key challenges that persist even after outlier smoothing: 1) significant rounding errors stemming from the variance of channel means (Sec. 3.2), and 2) substantial activation energy loss due to clipping (Sec. 3.3). These findings suggest that solely optimizing the rotation matrix offers diminishing returns. To overcome these fundamental limitations, in the following section, we introduce BASE-Q, a framework designed to directly target and correct these two specific error sources with high efficiency.

## 4 BASE-Q

As discussed in Sections 3.2 and 3.3, while rotation-based quantization methods effectively smooth outliers, they exhibit three notable limitations: (a) Channel-mean variance cannot be eliminated through rotation-only optimization, leading to a persistent rounding error term (as formalized in Equations (5) and (6)). (b) Implementing flexible and optimal activation clipping schemes within standard rotation-based PTQ is challenging, resulting in cumulative energy loss across layers. (c) Optimizing full-network rotation parameters requires prohibitive GPU memory and computational resources, especially for large models. To address these issues, we propose **BASE-Q** (**B**ias and **A**symmetric **S**caling **E**nhanced **Q**uantization), a lightweight yet powerful quantization framework that strengthens rotation-based methods through two key innovations: explicit **bias correction** to eliminate channel mean variance, and **asymmetric scaling** to improve clipping and quantizer fit at each block. The blockwise optimization capability circumvents the heavy cost and complexity of full-model rotation learning, thereby achieving efficiency and superior quantization performance.

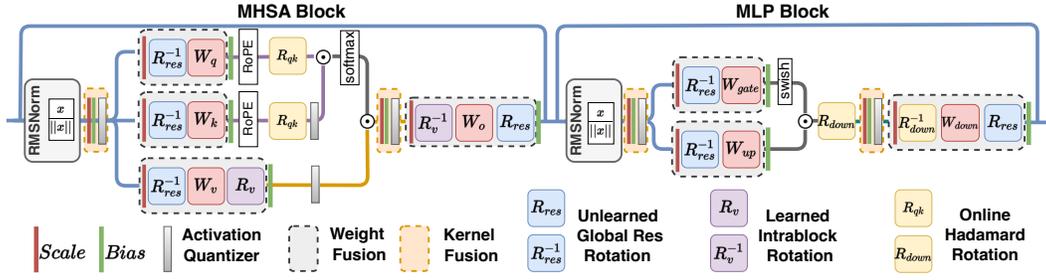


Figure 4: An overview of **BASE-Q**, highlighting three key design features: (a) channel-wise bias correction to reduce rounding error in activation quantization. (b) asymmetric scaling to compensate for the loss of computational equivalence caused by the clipping; (c) elimination of  $R_{res}$  learning, thereby avoiding full model optimization. This mechanism matches the parameter count and hardware fusion efficiency of typical scaling strategies.

**Fundamental Rotation Settings.** Figure 4 provides a schematic of our BASE-Q. The core optimized transformations within each transformer block include the global shared residual rotation ( $R_{res}$ ), the blockwise value rotation ( $R_v$ ), and the online rotations ( $R_{qk}$  and  $R_{down}$ ).  $R_{res}$  is applied to the residual pathway and can be fused into linear weights without additional inference overhead. Since the residual path connects all blocks, a single shared rotation is used throughout. As our method explicitly addresses channel mean variance and optimizes quantization for activations, we configure  $R_{res}$  to primarily optimize weight statistics. Following Section 3.1, we compute the PCA transform  $U$  across all related linear weights, and compose it with the Hadamard transformation:

$$R_{res} = U^T H, \quad (13)$$

where  $U$  is derived from the covariance of all associated linear weights.  $R_v$  can be exactly fused layer-wise into linear parameters and supports independent, block-level optimization. We employ PyTorch’s built-in Cayley orthogonal mapping for rotation learning, initialized as a Hadamard matrix.  $R_{qk}$  and  $R_{down}$  require on-the-fly computation during inference. To reduce memory and computation overhead, we adopt Hadamard rotations, which contain binary elements and can be efficiently implemented through fast algorithms.

**Bias Correction.** To directly address the rounding error component caused by  $Var(\mu_j)$ , we introduce a learnable bias prior to activation quantization. After quantized inference and the corresponding linear projection, this bias is subtracted, yielding:

$$y = Q_w(WsR) Q_a([R^{-1}x - b^c]) + \underbrace{WRb^c + b}_{fused\ bias} \quad (14)$$

where  $WR$  denotes rotated weights and  $b$  represents the usual layer bias. The learnable bias correction terms—including  $b_{qkv}^c$ ,  $b_o^c$ ,  $b_{up}^c$  and  $b_{down}^c$ —incur negligible parameter overhead (less than 0.1% of original model size) and are highly efficient to optimize.

**Asymmetric Scaling.** To counteract the energy loss from clipping in post-rotation Gaussian-like distributions, we introduce an per-quantizer asymmetric scaling  $s^a$ . This asymmetric scaling  $s^a$  adjusts the quantization range at inference time to better accommodate activation statistics and clipping requirements:

$$Y = Q_w(WsR) Q_a(s^a R^{-1} s^{-1} x - b^c) + (WRb^c + b) \quad (15)$$

**Blockwise Optimization.** Combining these techniques, our blockwise optimization strategy jointly optimizes  $R_v$ ,  $b_i^c$ ,  $s_j$ ,  $s_i^a$ , and activation clipping threshold  $\alpha_i$  per block via minimizing the MSE between the floating-point and quantized outputs:

$$\underset{R_v, b_i^c, s_j, s_i^a, \alpha_k}{\operatorname{argmin}} \mathcal{L}_{mse}(y_{FP}, y_Q; R_v, b_i^c, s_j, s_i^a, \alpha_i, \theta) \quad (16)$$

where  $y_{FP}$  and  $y_Q$  denote the reference and quantized outputs respectively, while  $\theta$  represents the set of frozen model and rotation parameters. This framework achieves state-of-the-art quantization accuracy with negligible memory and computational overhead during the quantization process.

## 5 EXPERIMENT

**Models and Tasks.** We evaluate our method on 12 open-source LLMs spanning various sizes: Llama-2-7/13/70B (Touvron et al., 2023), Llama-3-8/70B (Dubey & et al., 2024), Llama-3.1-8/70B, Llama-3.2-1/3B, and Qwen2.5-3/14/32B (Yang et al., 2024). We report perplexity on the Wikitext-2 dataset (Merity et al., 2016) and zero-shot accuracy on nine downstream tasks, including ARC-Easy and ARC-Challenge (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), LAMBADA (Radford et al., 2019), OpenBookQA (Mihaylov et al., 2018), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), and WinoGrande (Sakaguchi et al., 2021). We compare against baseline methods: QuaRot (Ashkboos et al., 2024), which employs fixed Hadamard rotations; SpinQuant (Liu et al., 2024), which leverages learned rotations; and OSTQuant (Hu et al., 2025), which incorporates both learned rotations and scaling.

**Deployment Details.** Our quantization framework is implemented in PyTorch and evaluated using lm-eval (Gao et al., 2024). We apply per-token asymmetric dynamic quantization for activations, per-head asymmetric dynamic quantization for KV-cache, and per-channel symmetric quantization for weights. Calibration is performed using 128 samples from Wikitext-2. For each transformer block, we first train symmetric scaling for 3 epochs, then apply weight quantization using GPTQ (Frantar et al., 2022), and finally train bias terms, asymmetric scaling, and learnable clipping factors for 5 epochs. Learning rates are initialized as  $1e-2$  for scaling and clipping,  $1e-3$  for bias, with cosine decay applied throughout training. The quantization process requires approximately 0.7 hours for 3B models and 10 hours for 70B models on a single A800 GPU. [This highlights the efficiency of our blockwise approach, which avoids the substantial resource requirements of full-model optimization methods like SpinQuant, which necessitates at least 5 A800 GPUs and 30 GPU-hours for a 70B model.](#)

### 5.1 MAIN RESULTS

Table 1 presents the overall performance of BASE-Q and baseline methods under W4A4KV4 quantization, reporting perplexity on Wikitext-2 and the average accuracy across nine zero-shot tasks. Since BASE-Q is explicitly designed to tackle the most challenging 4-bit activation quantization, we focus solely on the W4A4KV4 configuration. Across all evaluated models, BASE-Q consistently achieves superior results, achieving an average accuracy degradation of only 3.40% from full precision. In comparison, the average accuracy drops for QuaRot, SpinQuant, and OSTQuant are 6.87%, 5.95%, and 4.80%, respectively (averaged over eight supported models). This corresponds to BASE-Q reducing the performance gap to full precision by 50.5%, 42.9%, and 29.2% relative to QuaRot, SpinQuant, and OSTQuant, respectively. Notably, BASE-Q performs exceptionally well on Qwen2.5-3B, where all existing methods suffer severe degradation (perplexity increases by more than 100%). These results strongly validate our theoretical insights and demonstrate the practical effectiveness of our proposed bias correction and asymmetric scaling strategies. To provide a more comprehensive evaluation, we additionally test our method on the complex, multi-subject MMLU benchmark in Table 2.

Additional evaluations across different bit-width settings (W6A6KV6, W3A4KV4, W2A4KV4) are presented in Tables 3, 10 and 11, including comparisons against FlatQuant. FlatQuant improves quantization performance via additional learned online transforms for each layer. While this approach proves effective at certain bit-widths, it introduces substantial computational overhead compared to methods based on fusible rotation. Furthermore, we observe that these complex transforms are susceptible to numerical instability. Notably, FlatQuant fails on the Qwen2.5-3B model under our W3A4 setting and struggles to converge under W2A4 across most models, highlighting a practical robustness challenge at ultra-low precision. In contrast, BASE-Q’s method of combining a fusible global rotation with lightweight corrections proves to be both more stable and more efficient across all tested bit-widths.

### 5.2 PREFILL ACCELERATION ON GPU

We benchmark INT4 quantization during the compute-bound prefill stage on an NVIDIA 3090 GPU. Int4 matrix multiplications are implemented using NVIDIA’s Cutlass library, while other custom operators are written in Triton to ensure flexibility and speed. To minimize computational overhead,

Table 1: Comparison of Wikitext-2 perplexity and accuracy on 9 zero-shot benchmark tasks. All baseline results for QuaRot, SpinQuant, and OSTQuant are reproduced using their official open-source implementations, with necessary modifications to support Qwen models, which differ from LLaMA primarily through the inclusion of attention bias. Due to the absence of FSDP (Fully Sharded Data Parallel) support in the official OSTQuant repository, its results are limited to models with up to 14B parameters. Complete results are provided in Section B.

W4A4KV4	Qwen-2.5 3B		Qwen-2.5 14B		Qwen-2.5 32B		LLaMA-3.1 8B		LLaMA-3.1 70B		LLaMA-3.2 1B	
	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)
Full-Precision	64.17	8.03	70.95	5.29	71.11	5.02	68.70	6.23	73.76	2.81	55.89	9.75
QuaRot	44.30	69.33	67.23	6.77	68.14	6.04	63.74	7.82	69.56	5.31	48.66	14.44
SpinQuant	46.86	46.35	67.29	6.55	68.51	5.88	64.58	7.51	70.69	4.74	49.41	13.46
OSTQuant	50.81	20.09	67.81	6.37	OOM	OOM	64.91	7.40	OOM	OOM	<b>50.85</b>	12.84
<b>BASE-Q</b>	<b>58.93</b>	<b>10.43</b>	<b>69.17</b>	<b>6.28</b>	<b>70.18</b>	<b>5.65</b>	<b>65.36</b>	<b>7.17</b>	<b>71.54</b>	<b>4.17</b>	50.61	<b>12.66</b>

W4A4KV4	Llama-2 7B		Llama-2 13B		Llama-2 70B		LLaMA-3 8B		LLaMA-3 70B		LLaMA-3.2 3B	
	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)										
Full-Precision	65.22	5.47	67.62	4.88	71.57	3.32	68.11	6.14	73.82	2.86	63.59	7.81
QuaRot	61.59	6.12	65.03	5.39	70.28	3.76	62.89	7.82	68.76	5.62	55.86	10.07
<b>DFRot</b>	61.80	6.25	64.95	5.43	68.78	4.02	62.94	7.91	69.62	5.03	-	-
SpinQuant	61.38	5.99	65.63	5.30	70.22	3.71	63.80	7.49	69.93	5.11	57.74	9.32
OSTQuant	62.08	5.92	65.43	5.24	OOM	OOM	64.72	7.36	OOM	OOM	59.29	9.16
<b>BASE-Q</b>	<b>62.50</b>	<b>5.85</b>	<b>65.95</b>	<b>5.19</b>	<b>70.74</b>	<b>3.59</b>	<b>65.60</b>	<b>7.12</b>	<b>71.83</b>	<b>4.06</b>	<b>60.02</b>	<b>9.01</b>

Table 2: Summary of MMLU average accuracy (5-shot) comparison. BASE-Q demonstrates consistent performance gains, especially on challenging models like Qwen2.5-3B.

Method	Llama-2-7B	Llama-2-13B	Llama-3-8B	Qwen2.5-3B	Qwen2.5-14B
FP16	45.87	55.23	65.30	65.70	79.73
QuaRot	39.28	50.10	55.96	28.68	72.67
SpinQuant	<b>40.69</b>	50.05	56.25	29.92	73.15
<b>BASE-Q</b>	40.65	<b>51.57</b>	<b>57.82</b>	<b>53.80</b>	<b>75.49</b>

we fuse the online bias and scaling operations introduced by BASE-Q with quantization and de-quantization into a single Triton kernel (Figure 6), rendering the additional costs of bias and scaling computation negligible. For query and key projections (with small dimensions, e.g., 128×128), we implement the Hadamard transformation directly as a matrix multiplication. For the much larger *down* layer, we apply a single Cooley-Tukey step, ensuring that all compute-intensive kernels operate on manageable, high-parallelism sub-matrices. As shown in Figure 5, our approach yields 2.1× to 2.4× acceleration across all batch sizes compared to standard FP16. Notably, our optimizations incur only minimal overhead versus pure INT4 quantization without any online operations. During the decoding phase, which is typically characterized by a small batch size, the limited parallelism makes it challenging to gain throughput improvements from standard INT4 GEMM kernels. We found that a more effective strategy in this memory-bandwidth-bound scenario is to perform an INT4-FP16 GEMV (Wang et al., 2024) kernel, which yields a significant 1.8x to 2.3x throughput

Table 3: Performance comparison for W3A4KV4 quantization. PPL is measured on WikiText2, and Avg.acc represents the average accuracy on Zero-shot Common Sense Reasoning tasks.

Method	Qwen2.5-3B		Qwen2.5-14B		Llama2-7B		Llama2-13B		Llama3-8B	
	PPL(↓)	Avg.acc(↑)	PPL(↓)	Avg.acc(↑)	PPL(↓)	Avg.acc(↑)	PPL(↓)	Avg.acc(↑)	PPL(↓)	Avg.acc(↑)
FP16	8.03	65.80	5.29	72.68	5.47	66.70	4.88	69.98	6.14	69.97
Quarot	89.33	45.81	7.33	67.70	6.81	59.50	5.87	64.02	9.87	61.37
Spinquant	56.36	47.40	7.15	68.47	6.72	60.13	5.78	64.83	8.70	61.61
Flatquant	10045	30.62	6.95	<b>69.94</b>	6.48	61.49	5.52	<b>66.60</b>	8.49	63.03
<b>Base-Q</b>	<b>11.77</b>	<b>56.26</b>	<b>6.83</b>	67.35	<b>6.21</b>	<b>61.71</b>	<b>5.45</b>	65.55	<b>7.955</b>	<b>63.98</b>

increase. Furthermore, across all decoding scenarios, our method consistently achieves a 3.4x to 3.6x reduction in memory footprint through quantized weights and KV-cache (Table 6).

### 5.3 ABLATION STUDY

We perform a systematic ablation study on 3B to 8B models to assess the effects of different quantization strategies in BASE-Q, alongside comparisons with Quarot, SpinQuant, and OSTQuant. Our experimental results, presented in Table 4 indicate that bias correction yields substantial perplexity reductions for Qwen2.5-3B and Llama3-8B, highlighting the importance of addressing inter-channel bias in these models. In contrast, its effect on Llama2-7B remains marginal. Notably, asymmetric scaling delivers consistent improvements across all three models. A detailed analysis of bias correction is provided in Section C.

Table 4: Ablation Study on WikiText2 (word perplexity↓)

Method	Fixed Rotation	Learned Rotation	Bias Corect.	Asym. Scale	Scale	Qwen2.5-3B	Llama-2-7B	Llama-3-8B
QuaRot	$R_{res} R_v R_{qk} R_{down}$					69.33	6.12	7.82
	$R_{res} R_{qk} R_{down}$	$R_v$				54.38 -5.95	6.18 +0.06	7.63 -0.19
	$R_{res} R_{qk} R_{down}$	$R_v$	✓			13.59 -40.79	6.12 -0.06	7.42 -0.21
	$R_{res} R_{qk} R_{down}$	$R_v$	✓	✓		10.82 -2.77	5.92 -0.20	7.22 -0.20
<b>BASE-Q</b>	$R_{res} R_{qk} R_{down}$	$R_v$	✓	✓	✓	10.83 +0.01	5.85 -0.07	7.14 -0.08
SpinQuant	$R_{qk} R_{down}$	$R_{res} R_v$				46.35	5.99	7.49
OSTQuant	$R_{qk} R_{down}$	$R_{res} R_v$			✓	20.09	5.92	7.36

We conducted additional ablation studies with global rotations using three distinct approaches: standard Hadamard matrices, random Hadamard matrices, and learnable rotations (following SpinQuant’s official codebase). In SpinQuant (cf. Fig. 4 in (Liu et al., 2024)), the authors observed that using random Hadamard rotations could cause large fluctuations and subpar quantization performance. Based on this observation, SpinQuant advocates for learnable global rotations to stabilize and improve quantization results under their framework.

However, according to Table 5, quantization metrics showed negligible differences between different rotation types within our BASE-Q framework across five benchmark LLMs. This suggests that any potential improvements from learnable global rotations are effectively subsumed by BASE-Q’s bias correction component. Therefore, the additional memory and computational cost introduced by learnable global rotations become unnecessary for maintaining quantization quality in our method.

Table 5: Ablation study on  $R_{res}$  choice.

W4A4KV4	$R_{res}$ Setting	Qwen-2.5 3B		Qwen-2.5 14B		LLaMA-2 7B		LLaMA-3 8B		LLaMA-2 13B	
		0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)	0-shot <sup>9</sup> Avg.(↑)	Wiki (↓)
Full-Precision		64.17	8.03	70.95	5.29	65.22	5.47	68.11	6.14	67.62	4.88
QuaRot	Random Hadamard	44.30	69.33	67.23	6.77	61.59	6.12	62.89	7.82	65.03	5.39
SpinQuant	Learned Rotation	46.86	46.35	67.29	6.55	61.38	5.99	63.80	7.49	65.63	5.30
<b>BASE-Q</b>	Standard Hadamard	<b>58.93</b>	<b>10.43</b>	<b>69.17</b>	<b>6.28</b>	<b>62.08</b>	<b>5.85</b>	<b>65.33</b>	<b>7.13</b>	<b>65.89</b>	<b>5.20</b>
	Random Hadamard	<b>58.60</b>	<b>10.43</b>	<b>69.22</b>	<b>6.28</b>	<b>62.24</b>	<b>5.87</b>	<b>65.21</b>	<b>7.14</b>	<b>65.46</b>	<b>5.20</b>
	Learned Hadamard	<b>57.95</b>	<b>10.44</b>	<b>68.43</b>	<b>6.26</b>	<b>62.68</b>	<b>5.86</b>	<b>65.01</b>	<b>7.13</b>	<b>65.60</b>	<b>5.20</b>

## 6 CONCLUSION

In this work, we analyze key challenges in rotation-based quantization for large language models, identifying channel mean discrepancies and cumulative clipping-induced energy loss as major sources of quantization errors. To address these issues, we propose BASE-Q, which introduces blockwise bias correction and per-channel asymmetric scaling to achieve accurate and efficient quantization with minimal resource overhead. Comprehensive experiments on established LLM benchmarks demonstrate that BASE-Q narrows the accuracy gap to floating-point baselines while substantially reducing memory usage, enabling single-GPU quantization even for large-scale models. Our results highlight BASE-Q as a practical and scalable approach to quantizing LLMs.

## 540 REPRODUCIBILITY STATEMENT

541  
542 To ensure the reproducibility of our results, we have included comprehensive details of our method-  
543 ology, experimental setup, and all hyperparameters in the main paper and its appendices. We will  
544 release our source code and quantized model checkpoints to facilitate verification and future work.  
545 An anonymized version of the code and checkpoints will be made available during the rebuttal pe-  
546 riod, and a public release will follow upon acceptance of the paper.

## 547 REFERENCES

- 548  
549 Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L. Croci, Bo Li, Martin Jaggi, Dan Al-  
550 listarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in ro-  
551 tated llms. *ArXiv*, abs/2404.00456, 2024. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:268819214)  
552 [CorpusID:268819214](https://api.semanticscholar.org/CorpusID:268819214).
- 553  
554 Jinze Bai and et al. Qwen technical report. *ArXiv*, abs/2309.16609, 2023. URL [https://api.](https://api.semanticscholar.org/CorpusID:263134555)  
555 [semanticscholar.org/CorpusID:263134555](https://api.semanticscholar.org/CorpusID:263134555).
- 556  
557 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-  
558 monsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*,  
559 volume 34, pp. 7432–7439, 2020.
- 560  
561 Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization  
562 of large language models with guarantees. *Advances in Neural Information Processing Systems*,  
563 36:4396–4429, 2023.
- 564  
565 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
566 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint*  
*arXiv:1905.10044*, 2019.
- 567  
568 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
569 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
*arXiv preprint arXiv:1803.05457*, 2018.
- 570  
571 Tri Dao. fast-hadamard-transform, 2023. [https://github.com/Dao-AILab/](https://github.com/Dao-AILab/fast-hadamard-transform)  
572 [fast-hadamard-transform](https://github.com/Dao-AILab/fast-hadamard-transform).
- 573  
574 DeepSeek-AI. Deepseek-v3 technical report. *ArXiv*, abs/2412.19437, 2024. URL [https://api.](https://api.semanticscholar.org/CorpusID:275118643)  
575 [semanticscholar.org/CorpusID:275118643](https://api.semanticscholar.org/CorpusID:275118643).
- 576  
577 Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit ma-  
578 trix multiplication for transformers at scale. *ArXiv*, abs/2208.07339, 2022. URL [https://](https://api.semanticscholar.org/CorpusID:251564521)  
[api.semanticscholar.org/CorpusID:251564521](https://api.semanticscholar.org/CorpusID:251564521).
- 579  
580 Abhimanyu Dubey and et al. The llama 3 herd of models. *ArXiv*, abs/2407.21783, 2024. URL  
581 <https://api.semanticscholar.org/CorpusID:271571434>.
- 582  
583 Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training  
584 quantization for generative pre-trained transformers. *ArXiv*, abs/2210.17323, 2022. URL  
<https://api.semanticscholar.org/CorpusID:253237200>.
- 585  
586 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Fos-  
587 ter, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muen-  
588 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang  
589 Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model  
590 evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- 591  
592 Xing Hu, Yuan Cheng, Dawei Yang, Zukang Xu, Zhihang Yuan, Jiangyong Yu, Chen Xu, Zhe  
593 Jiang, and Sifan Zhou. Ostquant: Refining large language model quantization with orthogonal  
and scaling transformations for better distribution fitting. *ArXiv*, abs/2501.13987, 2025. URL  
<https://api.semanticscholar.org/CorpusID:275907083>.

- 594 Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang,  
595 and Shi Gu. Brecq: Pushing the limit of post-training quantization by block reconstruc-  
596 tion. *ArXiv*, abs/2102.05426, 2021. URL [https://api.semanticscholar.org/  
597 CorpusID:231861390](https://api.semanticscholar.org/CorpusID:231861390).
- 598  
599 Haokun Lin, Haobo Xu, Yichen Wu, Jingzhi Cui, Yingtao Zhang, Linzhan Mou, Linqi Song, Zhenan  
600 Sun, and Ying Wei. Duquant: Distributing outliers via dual transformation makes stronger quan-  
601 tized llms. *Advances in Neural Information Processing Systems*, 37:87766–87800, 2024a.
- 602 Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan  
603 Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for  
604 on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:  
605 87–100, 2024b.
- 606  
607 Yijiang Liu, Huanrui Yang, Zhen Dong, Kurt Keutzer, Li Du, and Shanghang Zhang. Noisyquant:  
608 Noisy bias-enhanced post-training activation quantization for vision transformers. *2023  
609 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 20321–20330,  
610 2022. URL <https://api.semanticscholar.org/CorpusID:254069623>.
- 611 Yijiang Liu, Hengyu Fang, Liulu He, Rongyu Zhang, Yichuan Bai, Yuan Du, and Li Du. Fbquant:  
612 Feedback quantization for large language models. *ArXiv*, abs/2501.16385, 2025. URL [https://  
613 //api.semanticscholar.org/CorpusID:275932200](https://api.semanticscholar.org/CorpusID:275932200).
- 614  
615 Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Kr-  
616 ishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. Spinqant: Llm  
617 quantization with learned rotations. *ArXiv*, abs/2405.16406, 2024. URL [https://api.  
618 semanticscholar.org/CorpusID:270062819](https://api.semanticscholar.org/CorpusID:270062819).
- 619 Yuexiao Ma, Huixia Li, Xiawu Zheng, Feng Ling, Xuefeng Xiao, Rui Wang, Shilei Wen, Fei  
620 Chao, and Rongrong Ji. Affinequant: Affine transformation quantization for large language  
621 models. *ArXiv*, abs/2403.12544, 2024. URL [https://api.semanticscholar.org/  
622 CorpusID:268531127](https://api.semanticscholar.org/CorpusID:268531127).
- 623  
624 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
625 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 626  
627 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
628 electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*,  
629 2018.
- 630 Markus Nagel, Rana Ali Amjad, Mart van Baalen, Christos Louizos, and Tijmen Blankevoort. Up  
631 or down? adaptive rounding for post-training quantization. *ArXiv*, abs/2004.10568, 2020. URL  
632 <https://api.semanticscholar.org/CorpusID:216056295>.
- 633  
634 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
635 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 636  
637 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-  
638 sarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 639  
640 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialliqa: Common-  
641 sense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 642  
643 Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqiang Li, Kaipeng  
644 Zhang, Peng Gao, Yu Jiao Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated  
645 quantization for large language models. *ArXiv*, abs/2308.13137, 2023. URL [https://api.  
646 semanticscholar.org/CorpusID:261214575](https://api.semanticscholar.org/CorpusID:261214575).
- 647  
648 Yuxuan Sun, Ruikang Liu, Haoli Bai, Han Bao, Kang Zhao, Yuening Li, Jiabin Hu, Xianzhi Yu,  
649 Lu Hou, Chun Yuan, et al. Flatquant: Flatness matters for llm quantization. *arXiv preprint  
650 arXiv:2410.09426*, 2024.

648 Hugo Touvron and et al. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*,  
649 abs/2307.09288, 2023. URL [https://api.semanticscholar.org/CorpusID:  
650 259950998](https://api.semanticscholar.org/CorpusID:259950998).  
651

652 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
653 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
654 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

655 Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#:  
656 Even better llm quantization with hadamard incoherence and lattice codebooks. *arXiv preprint  
657 arXiv:2402.04396*, 2024.

658 Lei Wang, Lingxiao Ma, Shijie Cao, Quanlu Zhang, Jilong Xue, Yining Shi, Ningxin Zheng, Ziming  
659 Miao, Fan Yang, Ting Cao, et al. Ladder: Enabling efficient {Low-Precision} deep learning com-  
660 puting through hardware-aware tensor transformation. In *18th USENIX Symposium on Operating  
661 Systems Design and Implementation (OSDI 24)*, pp. 307–323, 2024.

662 Xiuying Wei, Yunchen Zhang, Yuhang Li, Xiangguo Zhang, Ruihao Gong, Jinyang Guo, and  
663 Xianglong Liu. Outlier suppression+: Accurate quantization of large language models by  
664 equivalent and optimal shifting and scaling. *ArXiv*, abs/2304.09145, 2023. URL [https:  
665 //api.semanticscholar.org/CorpusID:258187503](https://api.semanticscholar.org/CorpusID:258187503).  
666

667 Jingyang Xiang and Sai Qian Zhang. Dfrot: Achieving outlier-free and massive activation-free for  
668 rotated llms with refined rotation. *arXiv preprint arXiv:2412.00648*, 2024.

669 Guangxuan Xiao, Ji Lin, Mickael Seznec, Julien Demouth, and Song Han. Smoothquant: Accurate  
670 and efficient post-training quantization for large language models. *ArXiv*, abs/2211.10438, 2022.  
671 URL <https://api.semanticscholar.org/CorpusID:253708271>.  
672

673 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,  
674 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint  
675 arXiv:2412.15115*, 2024.

676 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
677 chine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A SPEED-UP WITH KERNEL FUSION

We implement kernel fusion to merge the bias, scaling, quantization, and dequantization operations into a single computational step, which reduces on-chip memory access. For the online Hadamard transform, we find that the fast algorithm (Dao, 2023) adopted in Quarot, which uses the recursive Cooley-Tukey method to achieve  $O(n \log n)$  complexity, does not fully leverage tensor core parallelism on modern GPUs. To address this, we implement the Hadamard transform as a matrix multiplication in Triton, better utilizing available compute. For the prefill stage, we accelerate int4 GEMM using the CUTLASS library. In the decoding stage, especially with low batch sizes, we leverage the BitBLAS (Wang et al., 2024) library to accelerate int4-fp16 GEMV computations.

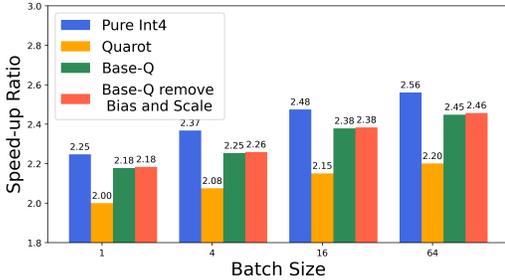


Figure 5: Prefill speedup for Llama2-7B with seqLens 2048 as batch size scales from 1 to 64.

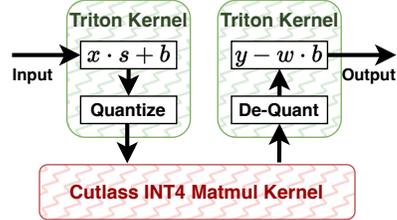


Figure 6: Illustration of kernel fusion.

Table 6: Memory saving on decoding stage for Llama-2-7B.

bsz × seqLens	1 × 256	1 × 1024	1 × 4096	16 × 256	16 × 1024	16 × 4096
Fp16	12.82	13.13	14.74	14.71	20.82	45.26
W4A4KV4	3.70	3.81	4.25	4.26	5.88	12.44
Saving Ratio	71.1%	71.0%	71.2%	71.0%	71.8%	72.5%

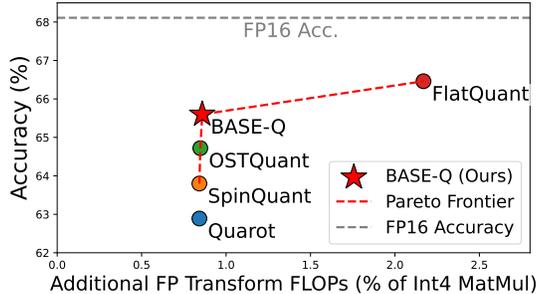


Figure 7: Pareto frontier plot between accuracy and inference cost.

## B FULL QUANTIZATION RESULTS

We present comprehensive quantization results in this section, including perplexity on WikiText2 and zero-shot accuracy on nine evaluation datasets. All baseline results, including QuaRot, SpinQuant, and OSTQuant, are reproduced using their official open-source implementations, with necessary modifications to support the Qwen model, which differs from LLaMA mainly by including attention bias. As the official OSTQuant repository does not support FSDP (Fully Sharded Data Parallel) training, its results are limited to models with up to 14B parameters, while larger models, such as the 70B model, encounter out-of-memory (OOM) issues. Results for the Llama2 series are summarized in Table 7, those for the Llama3 series are in Table 8, and for the Qwen2.5 series in Table 9. We also perform additional bit-widths (W6A6, W3A4, W2A4) to better demonstrate the robustness of BASE-Q.

Table 7: Complete comparison of the perplexity score on WikiText2 and accuracy on Zero-shot Common Sense Reasoning tasks for **Llama-2 models**.

Model	#Bits W-A-KV	Method	ARC-c (↑)	ARC-e (↑)	BoolQ (↑)	HellaS. (↑)	Lam. (↑)	OBQA (↑)	PIQA (↑)	SIQA (↑)	WinoG. (↑)	Avg. (↑)	Wiki2 (↓)
2-7B	16-16-16	Full Precision	46.33	74.54	77.74	76.02	73.90	44.20	79.05	46.16	69.06	65.22	5.47
		Quarot	42.15	70.37	73.00	73.09	70.68	39.40	77.26	43.30	65.04	61.59	6.12
	4-4-4	SpinQuant	40.78	70.45	73.79	72.40	71.30	38.40	75.63	43.45	66.22	61.38	5.99
		OSTQuant	41.64	68.94	74.43	73.17	71.61	<b>42.20</b>	77.09	<b>43.71</b>	65.90	62.08	5.92
		<b>BASE-Q</b>	<b>42.24</b>	<b>71.42</b>	<b>74.74</b>	<b>73.69</b>	<b>71.32</b>	41.80	<b>77.53</b>	43.35	<b>66.38</b>	<b>62.50</b>	<b>5.85</b>
2-13B	16-16-16	Full Precision	49.15	77.44	80.61	79.38	76.73	45.20	80.52	47.39	72.14	67.62	4.88
		Quarot	46.76	75.08	76.97	75.84	74.36	42.80	78.84	45.19	69.46	65.03	5.39
	4-4-4	SpinQuant	<b>48.46</b>	74.33	77.28	76.27	74.83	<b>44.60</b>	78.67	<b>46.01</b>	70.24	65.63	5.30
		OSTQuant	47.10	75.20	77.46	<b>77.71</b>	<b>75.14</b>	<b>44.60</b>	78.67	45.75	68.03	65.41	5.24
		<b>BASE-Q</b>	47.01	<b>75.67</b>	<b>78.90</b>	77.42	74.87	<b>44.60</b>	<b>79.27</b>	45.34	<b>70.48</b>	<b>65.95</b>	<b>5.19</b>
2-70B	16-16-16	Full Precision	57.25	81.06	83.76	83.82	79.62	48.80	82.70	49.18	77.98	71.57	3.32
		Quarot	55.97	<b>80.18</b>	81.87	82.25	78.73	48.00	81.39	47.49	75.69	70.28	3.76
	4-4-4	SpinQuant	54.78	79.76	81.90	82.78	79.20	47.40	81.77	48.46	<b>76.80</b>	70.22	3.71
		OSTQuant	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM
		<b>BASE-Q</b>	<b>56.40</b>	80.09	<b>82.17</b>	<b>82.79</b>	<b>79.37</b>	<b>48.20</b>	<b>82.32</b>	<b>48.62</b>	76.72	<b>70.74</b>	<b>3.59</b>

Table 8: Complete comparison of the perplexity score on WikiText2 and accuracy on Zero-shot Common Sense Reasoning tasks for **Llama-3 models**.

Model	#Bits W-A-KV	Method	ARC-c (↑)	ARC-e (↑)	BoolQ (↑)	HellaS. (↑)	Lam. (↑)	OBQA (↑)	PIQA (↑)	SIQA (↑)	WinoG. (↑)	Avg. (↑)	Wiki2 (↓)
3-8B	16-16-16	Full Precision	53.33	77.74	81.35	79.15	76.01	45.00	80.79	47.13	72.53	68.11	6.14
		Quarot	46.08	70.50	74.50	74.47	70.58	40.60	76.50	44.93	67.88	62.89	7.82
	4-4-4	SpinQuant	47.35	73.36	75.75	74.74	70.70	41.20	77.04	44.93	69.14	63.80	7.49
		OSTQuant	48.21	72.69	<b>79.02</b>	75.69	70.52	<b>44.00</b>	77.86	44.98	69.53	64.72	7.36
		<b>BASE-Q</b>	<b>50.51</b>	<b>76.05</b>	78.26	<b>76.48</b>	<b>71.22</b>	43.60	<b>78.89</b>	<b>45.75</b>	<b>69.61</b>	<b>65.60</b>	<b>7.12</b>
3-70B	16-16-16	Full Precision	64.33	85.90	85.23	84.89	79.82	48.60	84.55	50.72	80.35	73.82	2.86
		Quarot	53.16	78.11	82.91	80.80	75.49	44.40	79.65	47.49	76.87	68.76	5.62
	4-4-4	SpinQuant	57.25	80.22	83.06	81.05	75.04	46.20	81.88	47.44	77.27	69.93	5.11
		OSTQuant	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM
		<b>BASE-Q</b>	<b>59.98</b>	<b>82.70</b>	<b>84.65</b>	<b>83.92</b>	<b>78.42</b>	<b>47.20</b>	<b>83.08</b>	<b>48.41</b>	<b>78.14</b>	<b>71.83</b>	<b>4.06</b>
3.1-8B	16-16-16	Full Precision	53.50	81.19	82.08	78.90	75.82	44.80	81.23	47.19	73.56	68.70	6.23
		Quarot	46.08	72.31	77.34	74.46	71.36	42.40	76.82	44.37	68.51	63.74	7.82
	4-4-4	SpinQuant	47.87	76.05	76.76	74.63	70.54	<b>43.20</b>	<b>79.27</b>	45.70	67.17	64.58	7.51
		OSTQuant	47.44	75.34	<b>78.84</b>	75.48	<b>71.38</b>	41.80	78.24	47.13	68.51	64.91	7.40
		<b>BASE-Q</b>	<b>49.15</b>	<b>77.19</b>	78.56	<b>76.41</b>	70.72	42.40	79.16	<b>45.85</b>	<b>69.22</b>	<b>65.36</b>	<b>7.17</b>
3.1-70B	16-16-16	Full Precision	64.93	86.66	85.41	85.02	79.16	48.00	84.28	50.56	79.79	73.76	2.81
		Quarot	57.51	80.35	83.27	81.45	75.65	44.80	81.66	45.60	75.77	69.56	5.31
	4-4-4	SpinQuant	59.13	82.28	<b>84.37</b>	82.24	76.32	46.00	82.21	47.24	76.40	70.69	4.74
		OSTQuant	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM
		<b>BASE-Q</b>	<b>60.49</b>	<b>84.13</b>	83.12	<b>83.10</b>	<b>77.06</b>	<b>47.40</b>	<b>83.24</b>	<b>47.85</b>	<b>77.43</b>	<b>71.54</b>	<b>4.17</b>
3.2-1B	16-16-16	Full Precision	36.35	60.48	64.07	63.65	62.93	37.20	74.54	42.99	60.77	55.89	9.75
		Quarot	31.23	51.01	59.42	54.74	43.90	34.40	66.92	40.48	55.88	48.66	14.44
	4-4-4	SpinQuant	32.68	51.68	58.47	<b>56.68</b>	47.22	<b>34.80</b>	67.79	40.23	55.17	49.41	13.46
		OSTQuant	<b>33.45</b>	<b>55.26</b>	<b>61.80</b>	56.02	<b>48.90</b>	33.80	70.02	<b>41.61</b>	56.83	<b>50.85</b>	12.84
		<b>BASE-Q</b>	31.23	53.91	59.79	56.61	46.73	31.60	<b>70.57</b>	40.74	<b>56.91</b>	49.79	<b>12.63</b>
3.2-3B	16-16-16	Full Precision	45.90	71.63	73.33	73.60	70.48	43.00	77.58	46.98	69.85	63.59	7.81
		Quarot	38.40	59.13	64.56	66.49	60.00	37.20	72.25	42.68	62.04	55.86	10.07
	4-4-4	SpinQuant	37.12	60.90	68.72	68.83	61.67	39.60	73.88	44.58	64.40	57.74	9.32
		OSTQuant	41.04	<b>68.06</b>	68.84	68.75	<b>63.50</b>	40.40	74.21	<b>44.78</b>	64.01	59.29	9.16
		<b>BASE-Q</b>	<b>41.64</b>	66.86	<b>72.45</b>	<b>69.83</b>	62.57	<b>40.80</b>	<b>75.84</b>	44.52	<b>65.67</b>	<b>60.02</b>	<b>9.01</b>

Table 9: Complete comparison of the perplexity score on WikiText2 and accuracy on Zero-shot Common Sense Reasoning tasks for **Qwen2.5 Models**.

Model	#Bits W-A-KV	Method	ARC-c ( $\uparrow$ )	ARC-e ( $\uparrow$ )	BoolQ ( $\uparrow$ )	HellaS. ( $\uparrow$ )	Lam. ( $\uparrow$ )	OBQA ( $\uparrow$ )	PIQA ( $\uparrow$ )	SIQA ( $\uparrow$ )	WinoG. ( $\uparrow$ )	Avg. ( $\uparrow$ )	Wiki2 ( $\downarrow$ )
2.5-3B	16-16-16	Full Precision	47.53	73.06	77.22	73.52	67.11	42.00	78.84	49.80	64.14	64.17	8.03
	4-4-4	Quarot	30.72	52.69	51.16	49.52	20.63	34.60	66.70	38.84	53.83	44.30	69.33
		SpinQuant	34.47	57.58	54.04	50.91	24.14	32.20	68.01	40.48	57.93	46.86	46.35
		OSTQuant	38.48	58.08	61.90	56.19	36.33	37.60	67.90	42.48	58.33	50.81	20.09
		<b>BASE-Q</b>	<b>42.75</b>	<b>67.26</b>	<b>71.96</b>	<b>66.21</b>	<b>52.24</b>	<b>38.40</b>	<b>74.59</b>	<b>46.01</b>	<b>63.77</b>	<b>58.13</b>	<b>10.83</b>
2.5-14B	16-16-16	Full Precision	58.87	79.17	85.26	82.91	74.62	45.20	82.05	55.17	75.30	70.95	5.29
	4-4-4	Quarot	55.80	79.46	80.24	78.40	70.06	41.60	78.45	50.20	70.88	67.23	6.77
		SpinQuant	53.75	79.97	78.87	79.18	<b>70.97</b>	44.00	79.22	49.23	70.40	67.29	6.55
		OSTQuant	54.78	78.91	80.73	<b>79.81</b>	69.18	<b>44.40</b>	79.71	50.15	72.61	67.81	6.37
		<b>BASE-Q</b>	<b>56.14</b>	<b>82.41</b>	<b>82.23</b>	79.75	70.89	42.80	<b>80.20</b>	<b>52.97</b>	<b>72.69</b>	<b>68.90</b>	<b>6.28</b>
2.5-32B	16-16-16	Full Precision	55.63	80.93	87.19	84.06	76.96	44.00	82.32	56.29	75.22	71.11	5.02
	4-4-4	Quarot	52.05	74.66	85.41	81.70	73.14	42.40	79.98	52.81	71.11	68.14	6.04
		SpinQuant	52.82	76.05	85.11	80.99	72.71	44.40	80.03	52.15	72.30	68.51	5.88
		OSTQuant	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM	OOM
		<b>BASE-Q</b>	<b>54.61</b>	<b>78.16</b>	<b>87.06</b>	<b>82.49</b>	<b>75.61</b>	<b>44.60</b>	<b>81.12</b>	<b>53.63</b>	<b>74.35</b>	<b>70.18</b>	<b>5.65</b>

Table 10: Performance comparison for W6A6KV6 quantization. PPL is measured on WikiText2, and Avg.acc represents the average accuracy on Zero-shot Common Sense Reasoning tasks.

Method	Qwen2.5-3B		Qwen2.5-14B		Llama2-7B		Llama2-13B		Llama3-8B	
	PPL( $\downarrow$ )	Avg.acc( $\uparrow$ )								
FP16	8.03	65.80	5.29	72.68	5.47	66.70	4.88	69.98	6.14	69.97
Quarot	8.13	<b>66.09</b>	5.39	72.53	5.64	65.22	<b>4.90</b>	69.29	6.23	69.78
Spinquant	8.13	65.82	5.39	72.54	5.50	66.58	<b>4.90</b>	69.11	6.23	<b>69.92</b>
RoLoRa	-	-	-	-	-	<b>67.10</b>	-	68.80	-	68.10
Omniquant	-	-	-	-	7.48	58.65	6.74	61.02	-	-
<b>Base-Q</b>	<b>8.12</b>	65.45	<b>5.38</b>	<b>72.59</b>	<b>5.49</b>	66.68	<b>4.90</b>	<b>69.40</b>	<b>6.21</b>	69.48

Table 11: Performance comparison for W2A4 quantization. PPL is measured on WikiText2, and Avg.acc represents the average accuracy on Zero-shot Common Sense Reasoning tasks.

Method	Qwen2.5-3B		Qwen2.5-14B		Llama2-7B		Llama2-13B		Llama3-8B	
	PPL( $\downarrow$ )	Avg.acc( $\uparrow$ )								
FP16	8.03	65.80	5.29	72.68	5.47	66.70	4.88	69.98	6.14	69.97
Quarot	549.48	36.64	18.74	49.96	65.22	39.14	16.58	43.09	73.08	38.78
Spinquant	430.54	36.62	<b>17.81</b>	<b>52.29</b>	60.82	37.13	15.56	44.32	61.68	40.57
FlatQuant	NaN	NaN								
<b>Base-Q</b>	<b>22.86</b>	<b>44.09</b>	49.78	47.13	<b>14.01</b>	<b>48.62</b>	<b>14.20</b>	<b>57.95</b>	<b>21.92</b>	<b>45.93</b>

## C THE ABLATION STUDY OF BIAS CORRECTION

To clarify the effectiveness of bias correction, we analyzed the behavior of bias terms during quantization by examining the cosine similarity between the initialized channel means and the optimized bias correction coefficients on Qwen2.5-3B. Except for the first block—where channel mean deviation is not yet present, as shown in Figure 2 and table 12—the optimized bias values in all subsequent layers remain highly correlated with the initial channel means (cosine similarity  $>97.5\%$ ), which fully aligns with our theoretical analysis.

Table 12: Cosine similarity between blocks after quantization.

<b>Block index</b>	0	1	2	3	4	5	6	7
Cosine Similarity	0.556	0.976	0.979	0.997	0.988	0.988	0.988	0.989

## D VISUALIZATION OF ACTIVATION DISTRIBUTION

We provide visualization results for the input activations of the multi-head attention blocks and MLP blocks at the 1st, 11th, and 31st layers of Qwen2.5-3B, Llama2-7B, and Llama3-8B. For each selected layer, we illustrate the distributions of activations under various rotation strategies. We observe that inter-channel mean misalignment consistently occurs in different layers and across all evaluated models, and that learned rotations are insufficient to eliminate this issue. As discussed in Section 3.2, this misalignment is a significant source of quantization error, which can be effectively mitigated by our proposed bias correction technique.

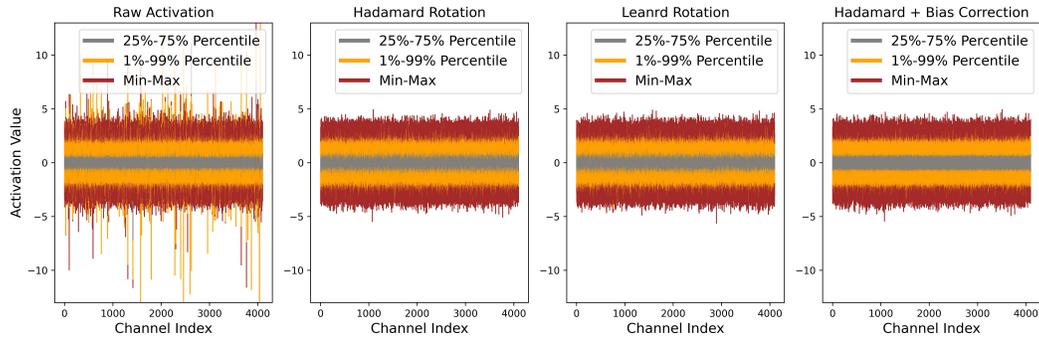


Figure 8: Visualizations comparing the activation distributions from the 1st MHA block in Llama2-7B.

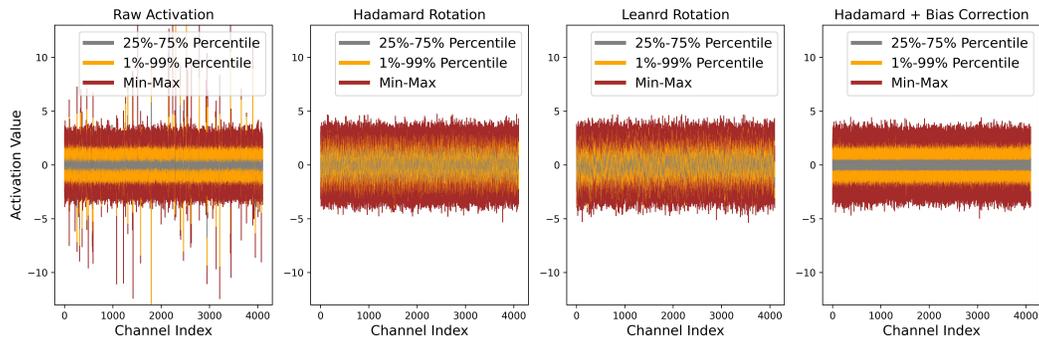


Figure 9: Visualizations comparing the activation distributions from the 1st MLP block in Llama2-7B.

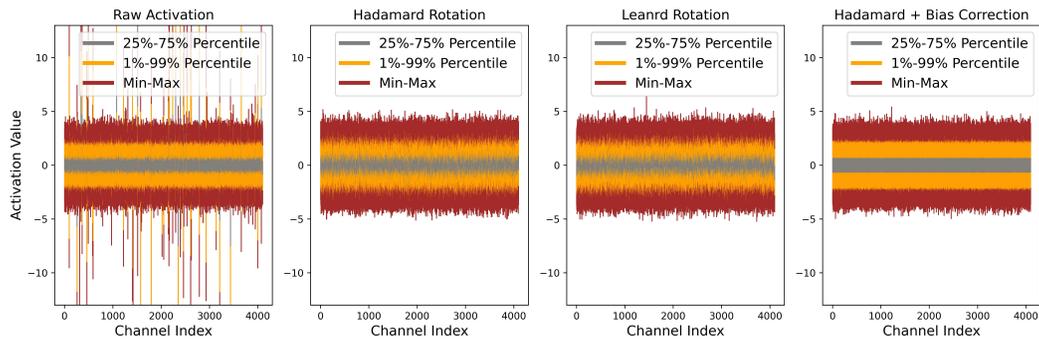


Figure 10: Visualizations comparing the activation distributions from the 11th MHA block in Llama2-7B.

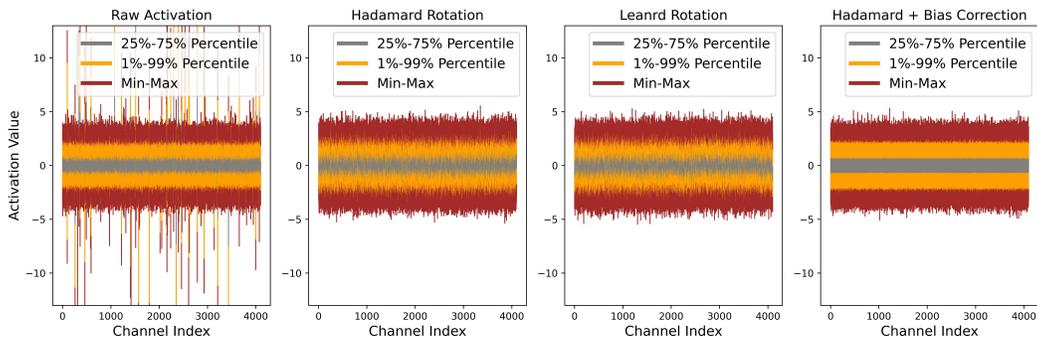


Figure 11: Visualizations comparing the activation distributions from the 11th MLP block in Llama2-7B.

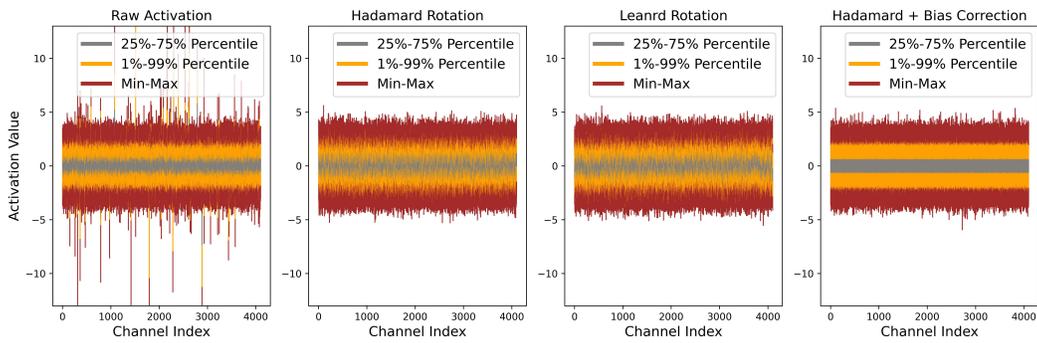


Figure 12: Visualizations comparing the activation distributions from the 31st MHSA block in Llama2-7B.

## E THE USE OF LARGE LANGUAGE MODELS (LLMs)

This paper was partially created with the assistance of a Large Language Model (LLM), which was used for tasks such as sentence polishing, brainstorming, and content organization. All content has been finally reviewed and confirmed by the author.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

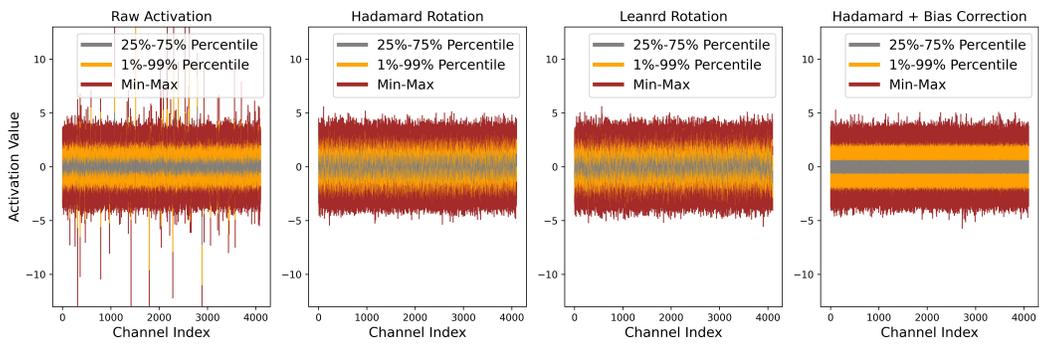


Figure 13: Visualizations comparing the activation distributions from the 31st MLP block in Llama2-7B.