

SingleQuant: Efficient Quantization of Large Language Models in a Single Pass

Jinying Xiao, Bin Ji, Shasha Li, Xiaodong Liu, Jun Ma, Ye Zhong, Wei Li, Xuan Xie, Qingbo Wu, Jie Yu*

National University of Defense Technology
Changsha, Hunan 410073 China
jinyingxiao@nudt.edu.cn

Abstract

Large Language Models (LLMs) quantization facilitates deploying LLMs in resource-limited settings, but existing methods that combine incompatible gradient optimization and quantization truncation lead to serious convergence pathology. This prolongs quantization time and degrades LLMs' task performance. Our studies confirm that Straight-Through Estimator (STE) on Stiefel manifolds introduce non-smoothness and gradient noise, obstructing optimization convergence and blocking high-fidelity quantized LLM development despite extensive training. To tackle the above limitations, we propose SingleQuant, a single-pass quantization framework that decouples from quantization truncation, thereby eliminating the above non-smoothness and gradient noise factors. Specifically, SingleQuant constructs Alignment Rotation Transformation (ART) and Uniformity Rotation Transformation (URT) targeting distinct activation outliers, where ART achieves smoothing of outlier values via closed-form optimal rotations, and URT reshapes distributions through geometric mapping. Both matrices comprise strictly formulated Givens rotations with predetermined dimensions and rotation angles, enabling promising LLMs task performance within a short time. Experimental results demonstrate SingleQuant's superiority over the selected baselines across diverse tasks on 7B-70B LLMs. To be more precise, SingleQuant enables quantized LLMs to achieve higher task performance while necessitating less time for quantization. For example, when quantizing LLaMA-2-13B, SingleQuant achieves $1,400\times$ quantization speedup and increases +0.57% average task performance compared to the selected best baseline.

Introduction

Large Language Models (LLMs) have shown exceptional capabilities in addressing various scientific tasks (Tang et al. 2025; Wang et al. 2025b). However, running LLMs for inference demands substantial computation and storage resources. Consequently, reducing resources required by running LLMs has attracted much research attention, with post-training quantization (PTQ) (Ashkboos et al. 2023; Lin et al. 2024; Zhao et al. 2024; Ma et al. 2024b; Liu et al. 2025) being one of the potential solutions.

In PTQ scenarios, outliers remain an open challenge, including massive outliers (MO) and normal outliers (NO) (Lin et al. 2024; Jin et al. 2025; Ramachandran, Kundu, and Krishna 2025). These outliers reduce the effective bit allocation for the majority of values, resulting in low quantization space utilization (Hu et al. 2025). Prior research leverages rotational invariance to mitigate outliers and enhance model quantizability through gradient-based optimization of rotation matrices (Sun et al. 2025; Liu et al. 2024b; Hu et al. 2025).

However, we observe that gradient optimization and quantization operations are not compatible, manifesting as pathological convergence (Fig. 2). Furthermore, gradient-based rotation matrix adjustment incurs substantial time overhead (Fig. 1a) and remains dependent on GPTQ support. Our theoretical analysis reveals that Cayley SGD with Straight-Through Estimator (STE) (Li, Li, and Todorovic 2020) manifests non-smooth and oscillatory convergence properties. Specifically, STE introduces fixed noise relative to true gradients (Bengio, Léonard, and Courville 2013), inducing optimization oscillations that impede convergence. Additionally, STE creates non-smoothness across quantization boundaries, violating the Lipschitz smoothness prerequisite for Cayley SGD convergence, thereby preventing gradient stabilization. Consequently, SpinQuant suffers from pathological optimization that consistently yields suboptimal results as empirically validated.

These limitations critically hinder the broader adoption of quantization techniques. To address this, we propose SingleQuant (see Fig. 1c), a mathematically constructed framework decoupled from quantization operations and avoiding optimization—thereby rendering quantization a deterministic computation. The core innovation leverages geometric properties of Givens rotations (Press 2007) to directly address both outlier categories: (1) Alignment Rotation Transformation (ART) targets sparse yet extreme-valued MO by smoothing extreme values via closed-form optimal rotation angles in a single operation. As shown in Fig. 1b, ART aligns outliers within the quantization domain, significantly boosting quantization space utilization; (2) Uniformity Rotation Transformation (URT) addresses widely-distributed NO with moderate values by reshaping distributions to break quantization bottlenecks, specifically leveraging Givens mapping feasibility (Ma et al. 2024a).

*Corresponding Author.

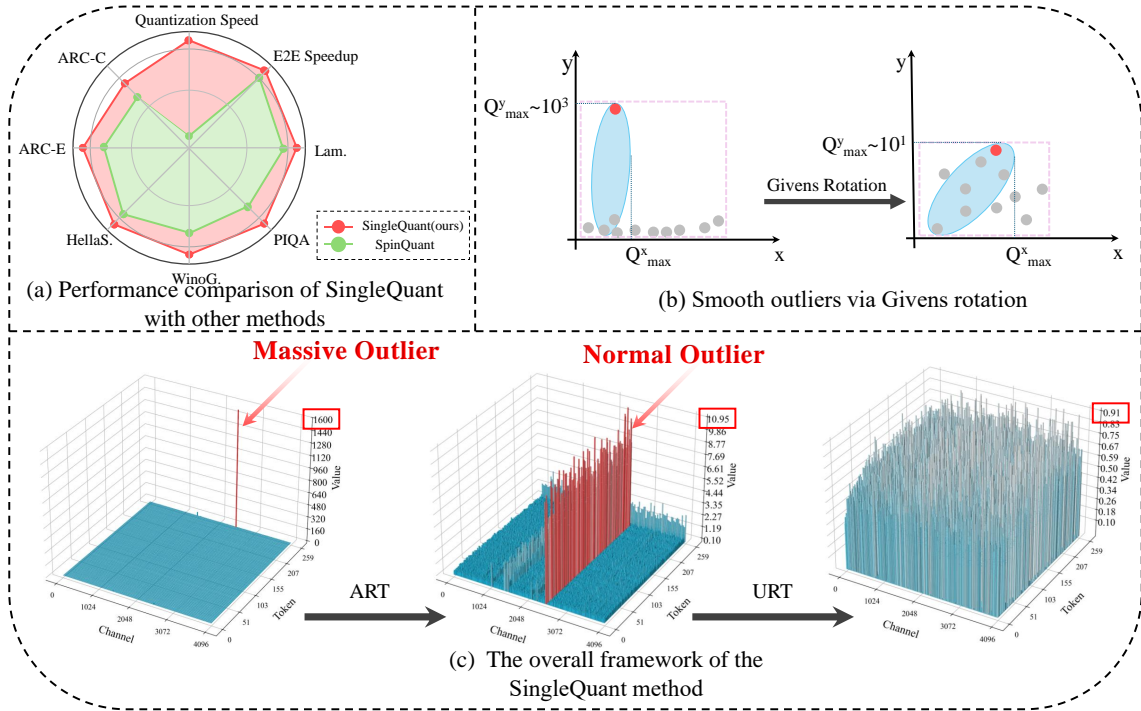


Figure 1: The sub-figure (a) compares SingleQuant and SpinQuant in quantization speed (LLMs quantization per hour), end-to-end speedup, and performance under various QA tasks. The sub-figure (b) illustrates deterministic outlier smoothing via Givens rotation on 2D data containing MO. Grey circles represent data points, with red ones indicating MO. Blue ellipses depict quantization space (size determined by bit-width), where greater coverage of data points within this fixed space indicates higher quantization space utilization. The sub-figure (c) presents SingleQuant’s framework comprising two components: ART smooths outlier magnitudes targeting prominent/scattered outliers, while URT performs secondary smoothing through distribution optimization. The diagram demonstrates ART/URT operations against distinct outlier types.

to achieve uniform distribution mapping through linearly-scaling Givens rotations, thereby neutralizing NO-induced quantization interference.

The core components of both ART and URT consist of rigorously formalized Givens rotations. The novel design lies in the fact that the dimension and rotation angle of the Givens rotation are based on the magnitude and distribution of outliers. Thus, all critical parameters of the quantization process are determined through a mathematically derived closed-form solution, rather than requiring pathological convergence on Stiefel manifolds via gradient feedback as in conventional methods. As demonstrated in Fig. 1a, SingleQuant outperforms the selected baselines in terms of task performance, inference latency, and quantization speed. Extensive experiments conducted across diverse tasks and LLMs confirm that SingleQuant establishes new state-of-the-art quantization results from both task performance and quantization speed.

The contributions of this work are summarized below:

- (1) Our research demonstrates that prior methods suffer from pathological convergence due to gradient noise induced by the STE and non-smoothness in Cayley SGD, which fundamentally conflicts with quantization truncation. Consequently, this defective optimization process generates substantial overhead and yields subopti-

mal quantization results.

- (2) To address the above limitations, we propose the SingleQuant quantization framework that decouples from quantization truncation. This approach first processes outlier magnitudes via the ART, then optimizes outlier distributions through the URT, and integrates Kronecker matrix decomposition. Extensive ablation studies confirm that this integrated approach enables more efficient outlier adjustment and superior quantization performance.
- (3) Experimental results demonstrate that SingleQuant achieves new state-of-the-art performance. For example, under extreme W4A4 quantization settings on LLaMA-2-70B, it attains 76.30% average zero-shot task accuracy—surpassing prior SOTA methods by 5.08%. Regarding quantization time, SingleQuant drastically compresses quantization time to 37 seconds for LLaMA-2-13B, a $1400\times$ reduction compared to the 14-hour requirement of baseline methods.

Related Work

LLMs Quantization. Quantization is a critical technique for reducing memory footprint and accelerating inference by using fewer bits for storage and computation. LLMs have been shown to exhibit outliers in activations (Shen et al.

2020; Bai et al. 2021) and MO in pivot tokens (Wei et al. 2022; Barbero et al. 2025; Liu et al. 2024a; Lin et al. 2024; Sun et al. 2024; Jin et al. 2025), which severely degrade quantization precision. To mitigate the impact of these outliers, pre-quantization transformations are widely adopted for weight and activation quantization. These methods employ specific invariances to scale or orthogonally transform activations and weights before quantization, smoothing and redistributing outliers (Xiao et al. 2023; Ashkboos et al. 2024; Ma et al. 2024b).

Learnable Transformation. Recent research has focused on exploring learnable transformations and clipping thresholds. For example, SpinQuant (Liu et al. 2024b) designs learnable orthogonal transformations while enforcing the orthogonality constraint via Cayley SGD. DuQuant (Lin et al. 2024) reduces quantization time via greedy rotation learning but achieves noncompetitive performance. FlatQuant (Sun et al. 2025) employs online matrix transformations to improve the distributions of weights and activations, at the cost of increased inference overhead and parameter count. OS-TQuant (Hu et al. 2025) outperforms prior methods across LLM benchmarks using learnable rotations and scalings.

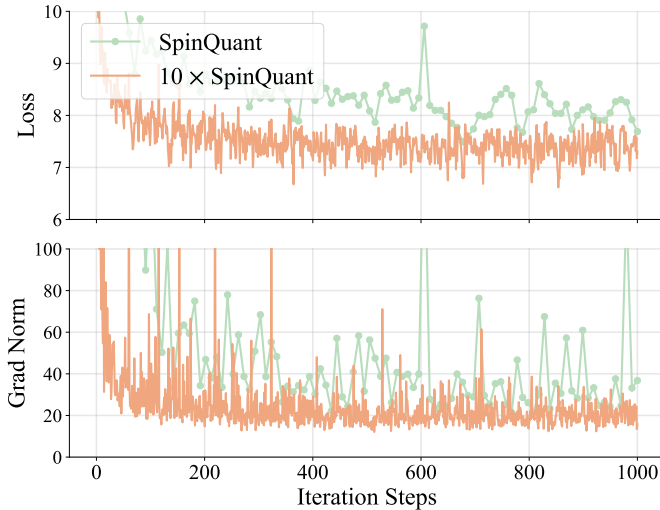


Figure 2: SpinQuant applies W4A4 quantization to LLaMA-2-7B with linearly decaying LR. The orange curve uses $10\times$ SpinQuant’s claimed iterations; adjacent green points are spaced 10 iterations apart. The figure shows optimization loss and gradient norm. More model results can be seen in Appendix C.

Discussions

While prior work (Liu et al. 2024b) employs Cayley SGD for optimization on the Stiefel manifold, we observe that gradient descent exhibits incompatibility with quantization operations, manifesting specifically as convergence pathologies during optimization. As evidenced in Fig. 2, SpinQuant exhibits persistent oscillations in loss and gradient norms with no convergence trend even after $10\times$ more iterations

(1,000 steps). We establish a theoretical basis for this pathology in this section.

Specifically, they optimize the rotation matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$ using the following formula:

$$\underset{\mathbf{R} \in \mathcal{M}}{\operatorname{argmin}} \mathcal{L}(\mathbf{R} \mid \mathbf{W}, \mathbf{X}). \quad (1)$$

Among them, \mathcal{M} denotes the Stiefel manifold, specifically the set of orthogonal matrices, and \mathbf{R} is constrained to update on this surface. $\mathcal{L}(\cdot)$ denotes the objective loss.

During the quantization optimization process, the non-differentiable quantization function $\operatorname{quant}(\cdot)$ is introduced, making it impossible to directly obtain its true gradient. To address this challenge, the STE (Yin et al. 2019) is commonly employed to approximate the gradient. According to prior research (Bengio, Léonard, and Courville 2013), the gradient of the STE can be expressed as:

$$\nabla \operatorname{STE}(\mathbf{R}) \approx \nabla \mathcal{L}(\mathbf{R}) + \mathbf{e}, \quad (2)$$

where \mathbf{e} denotes the noise error introduced by STE on gradient estimation. This noise numerically tends to exhibit substantial errors in both expectation and variance. We introduce the following lower bound for this expectation:

$$\mathbb{E}[\|\mathbf{e}\|] \geq c. \quad (3)$$

Furthermore, when the iterative trajectory of the parameter \mathbf{R} approaches a quantization transition point, the deviation between $\nabla \operatorname{STE}(\mathbf{R})$ and the theoretical true gradient $\nabla \mathcal{L}(\mathbf{R})$ increases significantly, thereby introducing additional stochastic perturbations during optimization. Further, prior studies on the convergence of Cayley SGD on the Stiefel manifold (Li, Li, and Todorovic 2020) rely on the critical Lipschitz assumption: that the gradient $\nabla \mathcal{L}(\mathbf{R})$ satisfies Lipschitz continuity across the entire Stiefel manifold. However, the quantization operation introduces non-smoothness at the transition points, violating the Lipschitz condition. More specifically, if the parameters are denoted as \mathbf{R}^- and \mathbf{R}^+ from either side of a quantization boundary, a noticeable difference in the gradient is observable at this boundary:

$$\|\nabla \operatorname{STE}(\mathbf{R}^+) - \nabla \operatorname{STE}(\mathbf{R}^-)\| \geq \frac{C}{\delta}, \quad (4)$$

where C is a constant, and δ denotes the quantization resolution (such as in 8-bit quantization, $\delta \sim 10^{-3}$). When $\delta \rightarrow 0$, the right-hand side of the above expression becomes boundless, indicating that $\nabla \operatorname{STE}$ cannot be bounded by any global Lipschitz constant at this point. Based on the foregoing analysis, we propose:

Assumption 1. *The gradient $\nabla \operatorname{STE}(\mathbf{R})$ does not satisfy Lipschitz continuity on the Stiefel manifold, i.e., there exists no global constant $L > 0$ such that:*

$$\forall \mathbf{R}_1, \mathbf{R}_2, \quad |\nabla \operatorname{STE}(\mathbf{R}_1) - \nabla \operatorname{STE}(\mathbf{R}_2)| \leq L \|\mathbf{R}_1 - \mathbf{R}_2\|_F. \quad (5)$$

(Proof and detailed discussion can be found in Appendix A)

Building upon the aforementioned assumption, we focus on the gradient variation of Cayley SGD under finite-step size conditions.

Theorem 1. *Under a finite step size, the Stiefel gradient norm of Cayley SGD satisfies:*

$$\liminf_{t \rightarrow \infty} \mathbb{E}[|\nabla_{\mathcal{M}} \mathcal{L}(\mathbf{R}_t)|] \geq \frac{c^2}{2} > 0. \quad (6)$$

(Proof and detailed discussion can be found in Appendix A)

Our findings directly demonstrate that since the STE noise term maintains a persistent lower bound during gradient estimation, the gradient norm relied upon by the algorithm cannot tend to zero as iterations progress. This implies that the optimization struggles to converge to an exact solution and is more inclined to oscillate within the solution space with constant amplitude.

By introducing the lower bound of STE noise and the assumption of non-smoothness induced by quantization, we systematically reveal that the synergistic effect between the noise lower bound and non-smoothness causes SpinQuant to oscillate within the solution space. This prevents convergence to the ideal scenario within a finite number of steps, incurring significant optimization overhead and suboptimal performance.

Method

Following the above analysis, we reveal fundamental convergence issues in quantized model optimization. The direct consequence of these issues is that even with substantial optimization costs, the resulting quantized models remain suboptimal. Therefore, we aim to employ a single rotation to smooth outliers, for which we design two components, both composed of strictly Givens orthogonal rotations.

Givens Rotation

In numerical linear algebra, the Givens rotation (Press 2007) represents a selective rotation within a two-dimensional plane. Formally, given a Givens orthogonal rotation denoted as $\mathbf{G}(i, j; \theta)$, the elements $\cos(\theta)$ and $\sin(\theta)$ occupy the intersection of rows i and j and columns i and j in \mathbf{G} , while the remaining non-zero elements lie along the diagonal. Geometrically, for a vector $x \in \mathbb{R}^d$, $x\mathbf{G}(i, j; \theta)$ signifies that vector x is rotated clockwise by angle θ in the subspace plane spanned by the i -th and j -th coordinate axes.

Due to the limitations of quantization space, previous works (Liu et al. 2024b; Lin et al. 2024) primarily focus on smoothing MO. Specifically, it can be expressed as:

$$\mathbf{X}\mathbf{W}^T = (\mathbf{X}\mathbf{R})(\mathbf{R}^T\mathbf{W}^T), \quad (7)$$

where the input activations are $\mathbf{X} \in \mathbb{R}^{T \times n}$, the corresponding weights are $\mathbf{W}^T \in \mathbb{R}^{n \times C_{out}}$, and \mathbf{R} matrices are designed algorithmically. However, owing to the characteristic of MO being sparse yet exhibiting large magnitudes, Givens rotation matrices can specifically target and smooth the MO.

Theorem 2. *Locally optimal smoothing rotation in two-dimensional subproblems.*

Let $V = (a, b)$ be a two-dimensional row vector. Among all matrices $\mathbf{G} \in \mathcal{O}(2)$, the matrix:

$$\mathbf{G}^* = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}, \quad \theta = \arctan \frac{b-a}{a+b}, \quad (8)$$

is the unique minimization function for V :

$$\Phi(\mathbf{G}) = \|\mathbf{V}\mathbf{G}\|_{\infty} = \max\{|x_1|, |x_2|\}, \quad (9)$$

where $\mathbf{V}\mathbf{G} = (x_1, x_2)$. Therefore, for vector V , the rotated value is:

$$\min_{\mathbf{G} \in \mathcal{O}(2)} \Phi(\mathbf{G}) = \sqrt{\frac{a^2 + b^2}{2}}. \quad (10)$$

(Proof and detailed discussion can be found in Appendix A)

Theorem 2 demonstrates that by specifying rotation axes (i, j) and angle θ^* , the target vector is rotated within the two-dimensional subspace spanned by the i -th and j -th coordinate axes. This rotation balances the energy of the two rotated components, achieving smoothing of MO along specific components.

SingleQuant

In this subsection, we introduce a single-optimization method—SingleQuant. This approach comprises two rotation components: For sparsely MO, we design an Alignment Rotation Transformation (ART) that precisely mitigates these outliers; For densely NO, we devise a Uniformity Rotation Transformation (URT) that adjusts their distribution through homogeneous mapping. Remarkably, SingleQuant requires only a single forward propagation pass to complete quantization without any additional optimization.

Kronecker Product. SingleQuant relies on using the Kronecker product to construct large rotation matrices $\mathbf{R} \in \mathbb{R}^{n \times n}$. Specifically, \mathbf{R} is composed of two lightweight matrices $\mathbf{R}_1 \in \mathbb{R}^{n_1 \times n_1}$ and $\mathbf{R}_2 \in \mathbb{R}^{n_2 \times n_2}$, i.e., $\mathbf{R} = \mathbf{R}_1 \otimes \mathbf{R}_2$. For a row vector V , we obtain:

$$V(\mathbf{R}_1 \otimes \mathbf{R}_2) = \text{vec}(\mathbf{R}_2^T V_{mat} \mathbf{R}_1)^T, \quad (11)$$

where $V_{mat} \in \mathbb{R}^{n_1 \times n_2}$ is the matrix reshaped from V , thus Equation 7 can be expressed as:

$$\mathbf{X}\mathbf{W}^T = \mathbf{R}_1^T \times_1 \tilde{\mathbf{X}} \times_2 \mathbf{R}_1 \times (\mathbf{R}_1^{-1} \times_1 \tilde{\mathbf{W}} \times_2 (\mathbf{R}_2^{-1})^T)^T, \quad (12)$$

where $\tilde{\mathbf{X}} \in \mathbb{R}^{k \times n_1 \times n_2}$ and $\tilde{\mathbf{W}} \in \mathbb{R}^{m \times n_1 \times n_2}$ are reshaped from \mathbf{X} and \mathbf{W} respectively, and \times_i denotes matrix multiplication along the i -th dimension.

The computational load of Equation (11) is dominated by the matrix multiplication $\mathbf{R}_2^T V_{mat} \mathbf{R}_1$. Notably, the pseudocode for factorizing the dimension n into n_1 and n_2 is provided in Appendix B. When n_1 and n_2 are appropriately chosen on the order of \sqrt{n} , this method achieves a significant dimensionality reduction from $O(n^2)$ to $O(n^{3/2})$, substantially reducing quantization and inference time.

In SingleQuant, \mathbf{R}_1 and \mathbf{R}_2 are both composed of ART and URT. Next, we introduce how to construct these two rotations.

Alignment Rotation Transformation. Conventional global rotation or progressive iterative methods (Shao et al. 2024; Liu et al. 2024b; Lin et al. 2024) often struggle to simultaneously achieve precision localization and efficient suppression of outliers: the former lacks targeted handling of MO, while the latter introduces additional overhead. Directly smoothing outliers through a single rotational transformation presents significant challenges. Therefore, we design the ART. ART leverages the conclusion of Theorem 2,

combining rapid MO detection with closed-form optimal rotation to smooth the most salient abnormal activations in a single transformation.

The construction expression for the ART matrix $\mathbf{R}^A \in \mathbb{R}^n$ is as follows:

$$\mathbf{R}^A = \begin{pmatrix} \mathbf{G}(\theta^*) & \mathbf{0} \\ \mathbf{0} & \mathbf{O} \end{pmatrix} \mathbf{P}_{ij}, \quad (13)$$

where i and j are the dimensions where the outlier and minimum value of the current activation \mathbf{X} reside, \mathbf{P}_{ij} is a permutation matrix that locates the outlier into the subspace $\mathbb{R}^{2 \times 2}$. $\mathbf{O} \in \mathbb{R}^{(n-2) \times (n-2)}$ is a randomly orthogonalized matrix that preserves metric invariance in high-dimensional subspaces, preventing additional noise from rotational propagation while ensuring Givens rotation acts solely on target dimensions.

ART smooths the outlier by computing the locally optimal rotation angle θ^* (see Theorem 2). Crucially, it operates independently of STE and its construction remains fully decoupled from quantization boundary operations, fundamentally avoiding the convergence issues to Cayley SGD. Moreover, ablation studies confirm this single-pass closed-form solution achieves theoretical optimality.

Uniformity Rotation Transformation. Although ART smooths MO, a substantial number of NO (Lin et al. 2024) persist in LLM activations. These outliers exhibit consistent median values across specific feature dimensions and exist in all token sequences (Xiao et al. 2023). We design the URT, which is implemented through two mappings while strictly preserving norm invariance. Specifically, for a row vector V , we aim to construct a rotation matrix \mathbf{R}^U such that:

$$V\mathbf{R}^U = U \in \mathbb{R}^n, \quad (14)$$

where the elements in U follow a uniform distribution and satisfy $\|V\|_2 = \|U\|_2$. Next, we explain how to construct \mathbf{R}^U . By respectively mapping V and U , we obtain:

$$V\mathbf{R}_{map} = U\mathbf{R}'_{map} = \|V\|_2 e_1^T, \quad (15)$$

where e_1 is the standard basis vector. Since Ma et al. demonstrated that two vectors with identical norms can be mutually transformed by $n-1$ Givens rotations (Ma et al. 2024a), thus through $O(n)$ complexity mapping of U and V onto e_1^T , we obtain \mathbf{R}_{map} and \mathbf{R}'_{map} . Thus the uniformity rotation matrix \mathbf{R}^U can be expressed as:

$$\mathbf{R}^U = \mathbf{R}_{map}(\mathbf{R}'_{map})^T. \quad (16)$$

Similar to the ART, the URT is entirely based on closed-form rotation matrix construction without additional gradient calculation or iterative tuning. Moreover, \mathbf{R}^U flattens the original activations into a uniform distribution and reduces quantization noise shift caused by outliers.

The Overall SingleQuant Method. To simultaneously address the smoothing of MO and NO, we first employ ART to smooth MO, eliminating quantization obstacles caused by extreme values. Next, we introduce the URT to flatten and uniformly distribute activations for quantization adaptation. Overall, \mathbf{R} in Equation 7 can be expressed as:

$$\mathbf{R} = (\mathbf{R}_1^U \mathbf{R}^A)^T \otimes (\mathbf{H} \mathbf{R}_2^U), \quad (17)$$

where \mathbf{H} denotes the Hadamard matrix, and \mathbf{R}_1^U and \mathbf{R}_2^U represent uniformity rotation matrices built according to different dimensions of the activation \mathbf{X} .

In summary, SingleQuant employs a unified transformation framework to seamlessly integrate ART and URT, achieving synchronous smoothing of both activations and weights. Notably, for weight processing, owing to the orthogonal property of rotation matrices, we apply structurally equivalent rotations to the original weight matrix. Formally, the SingleQuant methodology not only simultaneously mitigates long-tail risks in the distributions of both activations and weights, significantly reducing precision degradation caused by extreme outliers, but also achieves low-overhead transformation integration through rotation fusion techniques.

Experiments

Experimental Settings

Evaluation and Baselines. We evaluated the performance of SingleQuant on multiple models (Touvron et al. 2023; Chiang et al. 2023; Grattafiori et al. 2024). Quantized LLMs were assessed across three task categories: Language Generation Tasks (Merity et al. 2017), Zero-shot QA Tasks (Clark et al. 2018; Zellers et al. 2019; Paperno et al. 2016; Bisk et al. 2020; Sakaguchi et al. 2021) and MMLU (Hendrycks et al. 2021). SingleQuant was compared against popular INT4 post-training quantization (Wang et al. 2025a; Shao et al. 2024; Lin et al. 2024), including the two state-of-the-art approaches: SpinQuant (Liu et al. 2024b) and OS-TQuant (Hu et al. 2025).

Implementation Details. We implement the quantization and evaluation of SingleQuant building upon the Huggingface (Wolf et al. 2019), PyTorch (Paszke 2019), and lm-evaluation-harness (Gao et al. 2024) frameworks. To ensure experimental accuracy and mitigate randomness, each SingleQuant experiment was executed across 10 distinct random seeds, with the results reported representing the averaged metrics.

Main Results

Results on Language Generation Tasks. Table 1 reports SingleQuant’s perplexity results on WikiText-2 and C4 datasets. Notably, SingleQuant with RTN weight quantizer consistently outperforms prior SOTA quantization methods across most benchmarks, exhibiting particularly superior performance on larger models. Despite significant degradation in low-bit quantization for LLaMA-3 (Huang et al. 2024), SingleQuant on LLaMA3-70B/C4 exceeds the FP16 baseline by merely 0.84 while outperforming other SOTA baselines by 1.62. Remarkably, RTN-based SingleQuant surpasses GPTQ-based alternatives, demonstrating that our ART/URT components effectively manage outliers and compensate for RTN’s limitations.

Results on Zero-shot QA Tasks. On QA tasks, we evaluate six zero-shot tasks as shown in Table 2. SingleQuant further narrows the performance gap between the quantized model and the FP16 baseline. For instance, on the LLaMA-2-70B model, SingleQuant is only 0.91% lower in aver-

Table 1: WikiText-2 and C4 perplexity of 4-bit weight & activation quantized LLaMA models. ↓ denotes that the smaller the score, the better the performance. Bold values denote the best performance scores.

Method	W Quant.	WikiText-2↓					C4↓				
		2-7B	2-13B	2-70B	3-8B	3-70B	2-7B	2-13B	2-70B	3-8B	3-70B
FP16	-	5.47	4.88	3.32	6.14	2.86	7.26	6.73	5.71	9.45	7.17
SmoothQuant	RTN	83.12	35.88	26.01	210.19	9.60	77.27	43.19	34.61	187.93	16.90
OmniQuant	RTN	14.74	12.28	-	-	-	21.40	16.24	-	-	-
AffineQuant	RTN	12.69	11.45	-	-	-	15.76	13.97	-	-	-
QuaRot	RTN	8.56	6.10	4.14	10.60	55.44	11.86	8.67	6.42	17.19	79.48
QuaRot	GPTQ	6.10	5.40	3.79	8.16	6.60	8.32	7.54	6.12	13.38	12.87
QUIK-4B	GPTQ	8.87	7.78	6.91	-	-	-	-	-	-	-
SpinQuant	RTN	6.14	5.44	3.82	7.96	7.58	9.19	8.11	6.26	13.45	15.39
SpinQuant	GPTQ	5.96	5.24	3.70	7.39	6.21	8.28	7.48	6.07	12.19	12.82
MergeQuant	GPTQ	6.09	5.29	3.78	7.92	6.86	7.87	6.98	5.89	11.71	10.52
DuQuant	RTN	6.28	5.42	3.79	8.56	6.06	7.90	7.05	5.87	11.98	9.63
SingleQuant (Ours)	RTN	6.12	5.22	3.65	7.86	4.71	7.60	6.82	5.80	11.36	8.01

Table 2: Zero-shot⁶ AVG.¹ results of 4-bit weight & activation quantized LLaMA models. * denotes the usage of GPTQ weight quantization. The official OSTQuant repository lacks support for parallel training and inevitably encounters out-of-memory (OOM) errors when processing 70B models.

Method	Zero-shot ⁶ AVG. ↑				
	2-7B	2-13B	2-70B	3-8B	3-70B
FP16	69.87	72.55	77.05	73.23	79.95
QuaRot	57.73	66.25	73.47	61.34	35.36
QuaRot*	65.01	68.91	75.68	65.79	65.37
SpinQuant	63.52	68.56	75.09	66.98	65.66
SpinQuant*	66.23	70.93	76.06	68.70	71.33
DuQuant	61.34	64.98	69.39	65.76	72.97
OSTQuant*	66.39	70.27	OOM	69.08	OOM
SingleQuant	67.18	71.50	76.14	68.96	76.30

age accuracy than FP16, while other baseline methods suffer greater losses exceeding 1%. More detailed results are reported in Appendix C. These substantial improvements demonstrate that across diverse tasks, the two proposed components exhibit significant efficacy in smoothing weight and activation outliers.

Results on Instruction-tuned Models. To evaluate SingleQuant’s multi-architecture generalizability, we quantized the Vicuna-v1.5 model (Chiang et al. 2023). As shown in Table 3, our quantized model surpasses baselines on the MMLU benchmark, achieving an average improvement of 0.65% in 0-shot settings and 1.62% in 5-shot settings. Extended experiments are detailed in the Appendix C. These results confirm that SingleQuant maintains consistent efficacy across diverse architectures without requiring model-specific structural adjustments.

¹It denotes the average zero-shot performance across 6 tasks.

Table 3: Zero-shot and five-shot results on the MMLU benchmark for Vicuna-v1.5-7B under 4-bit weight & activation quantization.

LLM	Method	MMLU (0 shot) Avg. ↑	MMLU (5 shot) Avg. ↑
Vicuna v1.5-7B	FP16	48.75	49.85
	SmoothQuant	26.59	25.49
	OmniQuant	25.86	26.39
	Atom	36.14	37.15
	DuQuant	43.65	43.63
	SingleQuant	44.30	45.25

Ablation Study

Influence of ART/URT. We conducted an ablation study on SingleQuant, focusing on the independent contributions and synergistic effects of the proposed ART and URT components. Experiments on LLaMA-2 and LLaMA-3 (Table 4) demonstrate that: ART significantly enhances quantized model accuracy, validating its efficacy in smoothing MO to preserve performance—consistent with prior conclusions (Jin et al. 2025); furthermore, while URT alone yields limited gains, its integration with ART produces marked synergy. We attribute this to complementary mechanisms: ART reduces MO magnitudes, increasing NO prevalence with non-uniform distributions, while URT remaps these outliers into flatter distributions, optimizing quantization space utilization. These findings corroborate Fig. 1c: ART resolves MO magnitude issues and URT adjusts NO distribution structures, jointly enabling comprehensive outlier processing.

Quantization Time. Compared to existing baselines, SingleQuant demonstrates significant quantization speed advantages. Through the Kronecker decomposition in Equation (11) ensuring $O(n^{3/2})$ complexity order, and ART/URT

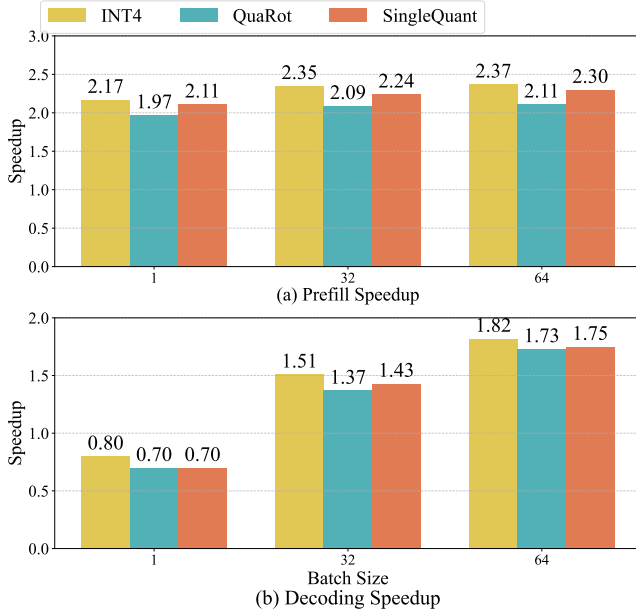


Figure 3: Prefill and decoding speedup of LLaMA-2-7B model across different batch sizes. We decode 256 tokens after the prefill on a sequence length of 2048.

matrices are single-pass optimized, SingleQuant is substantially faster than gradient-optimized or greedily algorithms. As shown in Table 5, SingleQuant quantizes a 13B model in 37 seconds, achieving up to $1420\times$ acceleration over baseline methods.

Inference Latency. To evaluate the inference speedup provided by SingleQuant, we adopted the measurement strategy and acceleration kernel from (Sun et al. 2025). As shown in Fig. 3, during the Prefill phase, SingleQuant achieves a $2.3\times$ speedup at batch size 64, merely $0.07\times$ lower than INT4. Notably, its acceleration exhibits strong robustness to batch size variations, significantly outperforming existing baselines across different batch scales. In the Decode phase, SingleQuant attains a $1.43\times$ speedup over FP16 at batch size 32, surpassing the baseline QuaRot by $0.06\times$. While a minor gap persists versus INT4, SingleQuant maintains significant advantages in preserving model accuracy, making it more suitable for practical deployment scenarios

Table 4: Effect of component ablation on LLaMA-2/3: PPL AVG. denotes mean WikiText-2 and C4 perplexity; 0-shot⁶ AVG. denotes average across 6 zero-shot tasks.

Rotation		2-13B		3-8B	
ART (R^A)	URT (R^U)	PPL AVG.↓	0-shot ⁶ AVG.↑	PPL AVG.↓	0-shot ⁶ AVG.↑
		7.15	68.91	11.82	63.44
	✓	6.89	69.13	10.52	64.17
✓		6.10	71.23	9.93	68.24
✓	✓	6.02	71.50	9.74	68.96

Table 5: Comparisons of quantization time cost by quantizing various LLMs. We run each experiment 30 times on a single NVIDIA A800 80GB GPU and report the averaged time.

LLM	OST Quant	Spin Quant	Du Quant	Single Quant
LLaMA2-7B	2,184s	22,334s	58s	24s
LLaMA2-13B	4,598s	52,561s	95s	37s
LLaMA3-8B	2,226s	24,849s	60s	27s

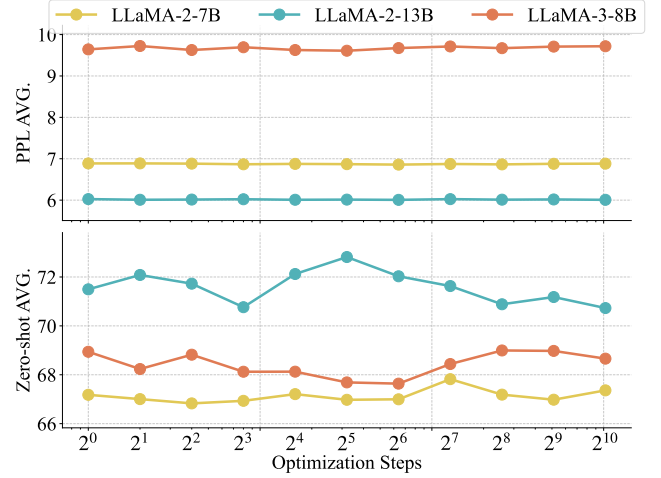


Figure 4: Performance comparisons of ART on SingleQuant through multiple optimization steps.

given its superior accuracy-performance trade-off.

Optimization Analysis. Previous sections demonstrate that MO can be effectively smoothed via single-pass Givens rotations. Here, we conduct an ablation study examining the impact of multi-pass optimization on SingleQuant—iteratively applying R^A transformations—with results reported in Fig. 4. The curves reveal SingleQuant’s robustness to optimization steps: perplexity remains essentially stable across LLaMA-2/3, while zero-shot metrics exhibit only minor fluctuations. Our findings confirm that ART achieves theoretical optimum via a single-step closed-form solution, effectively smoothing MO while ensuring no new outliers are introduced in other dimensions.

Conclusion

In this study, to address the incompatibility and prohibitive costs of LLM quantization, we propose SingleQuant, a single-optimization framework that effectively bridges the performance gap between full-precision and 4-bit weight-activation quantization. By leveraging rotation invariance, SingleQuant strategically inserts rotation matrices by constructing two specialized rotation components targeting distinct outlier types, while accelerating quantization and inference via Kronecker product decomposition. Consequently, SingleQuant establishes new state-of-the-art performance

for 4-bit weight-activation quantization, significantly facilitating efficient LLM deployment in resource-constrained environments.

References

- Ashkboos, S.; Markov, I.; Frantar, E.; Zhong, T.; Wang, X.; Ren, J.; Hoefler, T.; and Alistarh, D. 2023. Towards end-to-end 4-bit inference on generative large language models. *arXiv preprint arXiv:2310.09259*.
- Ashkboos, S.; Mohtashami, A.; Croci, M. L.; Li, B.; Cameron, P.; Jaggi, M.; Alistarh, D.; Hoefler, T.; and Hensman, J. 2024. Quarot: Outlier-free 4-bit inference in rotated llms. *Advances in Neural Information Processing Systems*, 37: 100213–100240.
- Bai, H.; Zhang, W.; Hou, L.; Shang, L.; Jin, J.; Jiang, X.; Liu, Q.; Lyu, M.; and King, I. 2021. BinaryBERT: Pushing the Limit of BERT Quantization. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 4334–4348.
- Barbero, F.; Vitvitskyi, A.; Perivolaropoulos, C.; Pascanu, R.; and Veličković, P. 2025. Round and Round We Go! What makes Rotary Positional Encodings useful? In *The Thirteenth International Conference on Learning Representations*.
- Bengio, Y.; Léonard, N.; and Courville, A. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*.
- Bisk, Y.; Zellers, R.; Gao, J.; Choi, Y.; et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 7432–7439.
- Chiang, W.-L.; Li, Z.; Lin, Z.; Sheng, Y.; Wu, Z.; Zhang, H.; Zheng, L.; Zhuang, S.; Zhuang, Y.; Gonzalez, J. E.; et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2(3): 6.
- Clark, P.; Cowhey, I.; Etzioni, O.; Khot, T.; Sabharwal, A.; Schoenick, C.; and Tafford, O. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- Gao, L.; Tow, J.; Abbasi, B.; Biderman, S.; Black, S.; DiPofi, A.; Foster, C.; Golding, L.; Hsu, J.; Le Noac’h, A.; Li, H.; McDonell, K.; Muennighoff, N.; Ociepa, C.; Phang, J.; Reynolds, L.; Schoelkopf, H.; Skowron, A.; Sutawika, L.; Tang, E.; Thite, A.; Wang, B.; Wang, K.; and Zou, A. 2024. The Language Model Evaluation Harness.
- Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Hendrycks, D.; Burns, C.; Basart, S.; Zou, A.; Mazeika, M.; Song, D.; and Steinhardt, J. 2021. Measuring Massive Multitask Language Understanding. In *International Conference on Learning Representations*.
- Hu, X.; Cheng, Y.; Yang, D.; Chen, Z.; Xu, Z.; Yuan, Z.; Zhou, S.; et al. 2025. OSTQuant: Refining Large Language Model Quantization with Orthogonal and Scaling Transformations for Better Distribution Fitting. In *The Thirteenth International Conference on Learning Representations*.
- Huang, W.; Ma, X.; Qin, H.; Zheng, X.; Lv, C.; Chen, H.; Luo, J.; Qi, X.; Liu, X.; and Magno, M. 2024. How good are low-bit quantized llama3 models? an empirical study. *arXiv e-prints*, arXiv–2404.
- Jin, M.; Mei, K.; Xu, W.; Sun, M.; Tang, R.; Du, M.; Liu, Z.; and Zhang, Y. 2025. Massive Values in Self-Attention Modules are the Key to Contextual Knowledge Understanding. In *Forty-second International Conference on Machine Learning*.
- Li, J.; Li, F.; and Todorovic, S. 2020. Efficient Riemannian Optimization on the Stiefel Manifold via the Cayley Transform. In *International Conference on Learning Representations*.
- Lin, H.; Xu, H.; Wu, Y.; Cui, J.; Zhang, Y.; Mou, L.; Song, L.; Sun, Z.; and Wei, Y. 2024. Duquant: Distributing outliers via dual transformation makes stronger quantized llms. *Advances in Neural Information Processing Systems*, 37: 87766–87800.
- Liu, R.; Bai, H.; Lin, H.; Li, Y.; Gao, H.; Xu, Z.; Hou, L.; Yao, J.; and Yuan, C. 2024a. IntactKV: Improving Large Language Model Quantization by Keeping Pivot Tokens Intact. In *Findings of the Association for Computational Linguistics ACL 2024*, 7716–7741.
- Liu, Y.; Fang, H.; He, L.; Zhang, R.; Bai, Y.; Du, Y.; and Du, L. 2025. Fbquant: Feedback quantization for large language models. *arXiv preprint arXiv:2501.16385*.
- Liu, Z.; Zhao, C.; Fedorov, I.; Soran, B.; Choudhary, D.; Krishnamoorthi, R.; Chandra, V.; Tian, Y.; and Blankevoort, T. 2024b. Spinquant: Llm quantization with learned rotations. *arXiv preprint arXiv:2405.16406*.
- Ma, X.; Chu, X.; Yang, Z.; Lin, Y.; Gao, X.; and Zhao, J. 2024a. Parameter efficient quasi-orthogonal fine-tuning via givens rotation. In *Proceedings of the 41st International Conference on Machine Learning*, 33686–33729.
- Ma, Y.; Li, H.; Zheng, X.; Ling, F.; Xiao, X.; Wang, R.; Wen, S.; Chao, F.; and Ji, R. 2024b. AffineQuant: Affine Transformation Quantization for Large Language Models. In *The Twelfth International Conference on Learning Representations*.
- Merity, S.; Xiong, C.; Bradbury, J.; and Socher, R. 2017. Pointer Sentinel Mixture Models. In *International Conference on Learning Representations*.
- Paperno, D.; Kruszewski, G.; Lazaridou, A.; Pham, N.-Q.; Bernardi, R.; Pezzelle, S.; Baroni, M.; Boleda, G.; and Fernández, R. 2016. The LAMBADA dataset: Word prediction requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1525–1534.
- Paszke, A. 2019. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*.

- Press, W. H. 2007. *Numerical recipes 3rd edition: The art of scientific computing*. Cambridge university press.
- Ramachandran, A.; Kundu, S.; and Krishna, T. 2025. Microscopiq: Accelerating foundational models through outlier-aware microscaling quantization. In *Proceedings of the 52nd Annual International Symposium on Computer Architecture*, 1193–1209.
- Sakaguchi, K.; Bras, R. L.; Bhagavatula, C.; and Choi, Y. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9): 99–106.
- Shao, W.; Chen, M.; Zhang, Z.; Xu, P.; Zhao, L.; Li, Z.; Zhang, K.; Gao, P.; Qiao, Y.; and Luo, P. 2024. OmniQuant: Omnidirectionally Calibrated Quantization for Large Language Models. In *The Twelfth International Conference on Learning Representations*.
- Shen, S.; Dong, Z.; Ye, J.; Ma, L.; Yao, Z.; Gholami, A.; Mahoney, M. W.; and Keutzer, K. 2020. Q-bert: Hessian based ultra low precision quantization of bert. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 8815–8821.
- Sun, M.; Chen, X.; Kolter, J. Z.; and Liu, Z. 2024. Massive activations in large language models. *arXiv preprint arXiv:2402.17762*.
- Sun, Y.; Liu, R.; Bai, H.; Bao, H.; Zhao, K.; Li, Y.; Yu, X.; Hou, L.; Yuan, C.; Jiang, X.; et al. 2025. FlatQuant: Flatness Matters for LLM Quantization. In *Forty-second International Conference on Machine Learning*.
- Tang, H.; Zhang, C.; Jin, M.; Yu, Q.; Wang, Z.; Jin, X.; Zhang, Y.; and Du, M. 2025. Time series forecasting with llms: Understanding and enhancing model capabilities. *ACM SIGKDD Explorations Newsletter*, 26(2): 109–118.
- Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Wang, J.; Wang, J.; Sun, H.; Yang, T.; Zhuang, Z.; Ning, W.; Yin, Y.; Qi, Q.; and Liao, J. 2025a. MergeQuant: Accurate 4-bit Static Quantization of Large Language Models by Channel-wise Calibration. *arXiv preprint arXiv:2503.07654*.
- Wang, X.; Tan, S.; Jin, M.; Wang, W. Y.; Panda, R.; and Shen, Y. 2025b. Do larger language models imply better reasoning? a pretraining scaling law for reasoning. *arXiv preprint arXiv:2504.03635*.
- Wei, X.; Zhang, Y.; Zhang, X.; Gong, R.; Zhang, S.; Zhang, Q.; Yu, F.; and Liu, X. 2022. Outlier suppression: Pushing the limit of low-bit transformer language models. *Advances in Neural Information Processing Systems*, 35: 17402–17414.
- Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Xiao, G.; Lin, J.; Seznec, M.; Wu, H.; Demouth, J.; and Han, S. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, 38087–38099. PMLR.
- Yin, P.; Lyu, J.; Zhang, S.; Osher, S.; Qi, Y.; and Xin, J. 2019. Understanding straight-through estimator in training activation quantized neural nets. *arXiv preprint arXiv:1903.05662*.
- Zellers, R.; Holtzman, A.; Bisk, Y.; Farhadi, A.; and Choi, Y. 2019. HellaSwag: Can a Machine Really Finish Your Sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4791–4800.
- Zhao, Y.; Lin, C.-Y.; Zhu, K.; Ye, Z.; Chen, L.; Zheng, S.; Ceze, L.; Krishnamurthy, A.; Chen, T.; and Kasikci, B. 2024. Atom: Low-bit quantization for efficient and accurate llm serving. *Proceedings of Machine Learning and Systems*, 6: 196–209.