MQUANT: UNLEASHING THE INFERENCE POTENTIAL OF MULTIMODAL LARGE LANGUAGE MODELS VIA FULL STATIC QUANTIZATION

Anonymous authors

006

012 013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

Paper under double-blind review

ABSTRACT

Recently, multimodal large language models (MLLMs) have garnered widespread attention due to their ability to perceive and understand multimodal signals. However, their large parameter sizes and substantial computational demands severely hinder their practical deployment and application. While quantization is an effective way to reduce model size and inference latency, its application to MLLMs remains underexplored. In this paper, we conduct an in-depth analysis of MLLMs quantization and identify several challenges: slow inference speed of the visual tokens, distributional differences across modalities, and visual outlier clipping degrades performance. To address these challenges, we propose MQuant, a quantization framework tailored for MLLMs. Specifically, 1) we design Modality-specific Quantization (MSQ) and Attention-Invariant Flexible Switching (AIFS) to support per-tensor static quantization and facilitate efficient inference. 2) we introduce an unified LayerNorm-to-RMSNorm transformation, achieving seamless integration of the MLLM vision encoder with Hadamard rotation. 3) we propose Rotation Magnitude Suppression (RMS) to mitigate outliers introduced by Hadamard rotation. Experiments conducted on five mainstream MLLMs demonstrate the superior performance and broad applicability of MQuant. For example, it maintains around 98% of the floating-point accuracy under the W4A8 setting. To the best of our knowledge, *MQuant* is the first quantization solution for MLLMs, paving the way for future advancements in their application.

032 1 INTRODUCTION

Recently, large language models (LLMs) such as GPT (Brown et al., 2020) and the LLaMA series (Touvron et al., 2023a;b; Dubey et al., 2024) have shown remarkable success across various domains, but they often struggle when handling non-textual data. Multimodal large language models (MLLMs) (Reid et al., 2024; Achiam et al., 2023b; Wang et al., 2023) overcome this by incorporating various modalities such as images and video, enabling a comprehensive understanding of diverse data types. However, the vast number of parameters and substantial computational demands due to the processing of multimodal inputs, which severely hinders their practical application. This is especially true in resource-constrained and privacy-sensitive environments.

Quantization is an effective technique for compressing and accelerating neural networks by converting
 high-precision values (e.g., FP32) to low-precision values (e.g., INT8), significantly reduces memory
 usage and inference latency. While quantization has been widely applied to LLMs (Yuan et al.,
 2023a; Shao et al., 2023), its application to MLLMs is essential but remains relatively under-explored.

045 When quantizing MLLMs, it is essential to consider the differences in inference efficiency that arise 046 from their multi-modal nature. A typical difference is the Time to First Token (TTFT) in MLLMs 047 increases significantly with the resolution and aspect ratio of input images or videos. As shown 048 in Fig 1(a), in models such as Qwen2-VL (Wang et al., 2024), the number of prefill visual tokens grows as image resolution increases (detailed figure in Fig. 7). This rapid expansion in token count exacerbates the inference latency, particularly in per-token dynamic quantization, which requires 051 separate processing for each token, such as memory access and scale computation. As a result, the TTFT in MLLMs increases drastically, severely impacting overall inference latency. Moreover, in 052 higher-demand scenarios such as video-based tasks and multi-image dialogues, the accumulation of visual tokens becomes even more pronounced, further exacerbating the increase in TTFT.

066

067 068



Figure 1: (a) The number of prefill visual tokens across different MLLMs as the image splits or resolution increases. (b) The activation values of visual tokens range from -20 to 10, whereas textual tokens are centered around 0, with only a few having absolute magnitudes exceeding 0.5.

One intuitive approach to address this issue is to employ per-tensor static quantization for prefill 069 multimodal tokens, which eliminates the need for token-wise scale computation. While this approach achieves up to 66% acceleration across all LLM linear layers (Fig. 4), it still results in a serious 071 accuracy degradation. We investigate the root causes and identify the following key challenges 072 in MLLM quantization: **1** Balancing efficiency and accuracy when adopting per-tensor static 073 quantization for multi-modal tokens. In MLLMs, input multi-modal tokens include visual and 074 textual with distinct value distributions. Unlike per-token dynamic quantization, which can adaptively 075 compute token-wise scales in an online manner, static quantization applies few scale factors across the 076 mixed multi-modal tokens in an offline manner, leading to a difficult trade-off between efficiency and 077 accuracy. **2** The data distribution of visual and textual tokens exhibits significant differences. Unlike LLMs, which solely handle textual features, MLLMs face the challenge of processing input with distinct distributions across multiple modalities. As shown in Fig 1 (b), the activation 079 distributions between visual and textual tokens reveal substantial numerical discrepancies. Joint quantization of these tokens may cause the larger values of visual tokens to overshadow the smaller 081 textual values, resulting in a loss of textual information. Therefore, MLLMs require modalityspecific approaches to account for the heterogeneous nature of visual and textual token distributions. 083

O Clipping outliers in visual tokens leads to perfor-084 mance degradation. As shown in the Tab. 1, applying 085 a stricter clipping range causes a large drop in accuracy for visual tokens, while textual tokens exhibit only a 087 minor decrease. In contrast, using a full clipping range 088 slightly mitigates the performance drop for visual tokens. This indicates that careful handling of outliers 089 in visual tokens is crucial-indiscriminate clipping of 090 outliers leads to severe accuracy loss. 091

Table 1: Results of the Qwen-VL on TextVQA using different clipping ranges for visual and textual tokens. Activations are quantized using per-tensor static setting.

Clip Range	Bits (W/A)	Visual Tokens	Textual Tokens
-	BF16 / BF16	61	.40
(0-0.99999) (0-1.0)	BF16 / INT16 BF16 / INT16	18.92 (↓42.48) 61.2 (↓0.20)	$60.09 (\downarrow 1.31)$ $61.25 (\downarrow 0.15)$

092 To this end, we propose *MQuant*, an accurate and efficient post-training quantization (PTQ) solution specifically designed for multimodal large language models (MLLMs). First, to reduce the TTFT 094 while maintaining accuracy, we introduce Modality-Specific Quantization (MSQ), which handles the distinct distribution differences between visual and textual tokens. For further acceleration, 095 we design an Attention-Invariant Flexible Switching (AIFS) scheme. AIFS transforms mixed 096 multimodal tokens into unified, modality-decoupled tokens without the need for dynamic position 097 vectors, maintaining computational equivalence and memory efficiency by avoiding the overhead 098 associated with dynamic processing. Second, based on the *computation invariance* (Ashkboos et al., 099 2024a), we further develop a **Post-LN + Rotate** scheme to extend and accommodate various MLLMs 100 architecture, enabling seamless integration with rotation to smooth outliers. Third, we reveal the 101 weight outliers caused by the online Hadamard rotations through theoretical analysis and propose 102 Rotation Magnitude Suppression (RMS) to mitigate them, further improving quantization accuracy. 103

We evaluate *MQuant* on various datasets and five mainstream MLLMs, including InternVL (Chen et al., 2024a), Qwen-VL (Bai et al., 2023b), MiniCPM-V (Yao et al., 2024), CogVLM2 (Hong et al., 2024) and Qwen2-VL (Wang et al., 2024). The results demonstrate that *MQuant* achieves less than 1% accuracy loss on all MLLMs under the W4A8 setting, highlighting its superior performance and wide applicability. We publish the code, hoping to provide new insights and advance the deployment and application of MLLMs. Our main contributions are summarized as follows:

- We present the first comprehensive analysis of quantization issues in MLLMs, unveiling the root causes of performance collapse and identifying inference speed bottlenecks as well as the quantization challenges posed by modality differences.
- We design Modality-specific Quantization (MSQ) and Attention-Invariant Flexible Switching (AIFS) to support static quantization. Besides, we propose a Post-LN + Rotate scheme tailored to accommodate various MLLMs.
- 113 114 115

110

111

112

• We analyze the root issues of weight outliers caused by online rotation, and propose Rotation Magnitude Suppression (RMS) to effectively mitigate it.

- 116 117
- We propose *MQuant*, to the best of our knowledge, the first PTQ solution for MLLMs.

118 2 RELATED WORK

119 Multimodal Large Language Models. The rapid advancement of LLMs has spurred significant 120 progress in MLLMs. Flamingo (Alayrac et al., 2022) pioneered connecting pre-trained visual 121 encoders to LLMs, demonstrating strong generalization across visual-language tasks. Following the 122 emergence of ChatGPT Achiam et al. (2023a), numerous open-source models based on pre-trained 123 LLMs like LLaMA (Touvron et al., 2023a) and its variants (Touvron et al., 2023b; Zheng et al., 2024) 124 have been proposed (Li et al., 2023; Huang et al., 2024; Zhu et al., 2023; Liu et al., 2024). Subsequent 125 efforts like Qwen-VL (Bai et al., 2023b), InternVL (Chen et al., 2024b), and CogVLMV2 Hong et al. (2024) enhanced MLLMs from perspectives such as high-resolution input and larger-scale training 126 data. However, the substantial parameters of MLLMs lead to high computational costs, limiting 127 broader application. Recently, smaller MLLMs like Mini-Gemini (Li et al., 2024), MobileVLM Chu 128 et al. (2023; 2024), and MiniCPM-V (Yao et al., 2024) have emerged. Despite these developments, 129 dedicated quantization methods for MLLMs to further reduce memory usage and accelerate inference 130 remain under-explored. 131

132 Post-Training Quantization for LLMs. Existing post-training quantization (PTQ) methods for 133 LLMs are categorized into weight-only and weight-activation quantization (Zhao et al., 2023; Yuan et al., 2024). Weight-only methods like GPTQ (Frantar et al., 2022), QuIP (Chee et al., 2024), and 134 AWQ (Lin et al., 2023) achieve high compression rates but offer limited inference acceleration. In con-135 trast, weight-activation quantization methods (Xiao et al., 2022; Wei et al., 2022; Yuan et al., 2023b; 136 Zhang et al., 2024) quantize both weights and activations, improving memory usage and latency. The 137 main challenge is activation outliers causing quantization errors. Techniques like SmoothQuant (Xiao 138 et al., 2022) shift quantization difficulty from activations to weights, while OmniQuant (Shao et al., 139 2023) optimizes performance by training quantization parameters. SliceGPT (Ashkboos et al., 2024a) 140 reduces memory demands by designing a Pre-LN + Rotate Scheme for LLMs sparsification based on 141 computational invariance. They achieve this by adding a linear layer in the residual connection (see 142 Appendix A.14). Unlike SliceGPT, we further develop a Post-LN + Rotate scheme to accommodate 143 more vision encoder and extends its applicability to various MLLMs. This enhancement broadens the 144 the LayerNorm + Rotate approach, making it suitable for both Pre-LN and Post-LN configurations across various MLLMs. Recently, Quarot (Ashkboos et al., 2024b) introduces rotations to eliminate 145 outliers; however, this solution is not applicable to MLLMs due to inherent modality differences. 146

147 3 PRILIMINARY

159

149 3.1 MLLM ARCHITECTURE

The existing MLLM framework (Bai et al., 2023b; Chen et al., 2024b) consists of three primary modules: a vision encoder E for processing visual inputs and feature extraction, a visual-language projector P that serves as a bridge to align the two modalities and a large language model (*LLM*) that handles the multi-modal tokens and performs reasoning.

Vision Encoder Taking the input image or video X_v as input, Existing vision encoders typically adhere to the Vision Transformer (ViT) (Dosovitskiy, 2021) architecture, such as CLIP (Radford et al., 2021), OpenCLIP (Ilharco et al., 2021) and SigLIP (Zhai et al., 2023). the vision encoder compresses the original vision information into more compact patch features F_v , which can be formulated as:

$$\boldsymbol{F}_{v} = E(\boldsymbol{X}_{v}) \tag{1}$$

160 Vision-Language Projector The task of the vision-language projector P is to map the visual patch 161 features F_v into the textual feature space:

$$\boldsymbol{E}_{v} = P(\boldsymbol{F}_{v}) \tag{2}$$

162 Large Language Model The pre-trained large language model serves as the core component of 163 MLLMs, endowing the framework with exceptional capabilities, such as zero-shot generalization, 164 instruction following, and in-context learning.

165 Typically, a text tokenizer is integrated with the LLM, mapping text 166 prompts X_t to the text tokens E_t . The text tokens E_t and the visual 167 tokens E_v are then concatenated to form the input for LLMs, which 168 outputs the final response sequence O_a in an autoregressive manner: 169

$$LLM(\boldsymbol{O}_a|\boldsymbol{E}_v, \boldsymbol{E}_t) = \prod_{i=1}^{r} LLM(y_i|\boldsymbol{E}_v, \boldsymbol{E}_t, y_{< i})$$
(3)

where l denotes the length of O_a . The parameter sizes of large lan-172 guage models (LLMs) range from 3 billion to tens of billions. Com-173 monly used open-source LLMs include the Llama series (Touvron 174 et al., 2023a;b), Qwen (Bai et al., 2023a), InternLM (Cai et al., 2024), 175 MiniCPM (Hu et al., 2024), ChatGLM (GLM et al., 2024). 176



178 Transformer networks (Vaswani et al., 2017) have proven highly 179 effective across a wide range of tasks, including language modeling



(5)

and visual perception. A transformer consists of multiple layers, each containing a multi-head self-181 attention (MHSA) block followed by a feed-forward network (FFN) block, with a pre-normalization 182 or post-normalization layer between them. The details of whole Transformer are in Appendix A.15.

183 **Attention Block in MLLMs** Given an input prompt, the generation process of LLMs can be broadly categorized into two distinct phases: the prefill phase, which computes and stores the KV 185 cache for input tokens, and the decoding phase, where new tokens are generated through a next-token-186 prediction scheme. Given input data X after LayerNorm (LN), and an attention block with its weight 187 matrices \mathbf{W}_q , \mathbf{W}_k and \mathbf{W}_v , attention scores is A, the prefill phase is formulated as: 188

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}_k, \quad \mathbf{V} = \mathbf{X}\mathbf{W}_v \tag{4}$$

189 190

191 192

193

194 195

170

171

 $\mathbf{A} = \operatorname{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{T}}{\sqrt{D}}\right), \quad \mathbf{O} = \mathbf{A}\mathbf{V}$ **Quantization** Quantization maps high-precision value into discrete levels, we adopt uniform quantization (Jacob et al., 2018) in our study. Given a floating-point (FP) tensor x (weights or activations), it can be uniformly quantized to *b*-bits in signed quantization as follows:

$$\hat{\mathbf{x}} = \mathbf{Q}_U(\mathbf{x}, b) = (clamp(\lfloor \frac{\mathbf{x}}{s} \rceil + z, q_{min}, q_{max}) - z) \cdot s, \quad s = \frac{\max(|\mathbf{x}|)}{2^{b-1} - 1}$$
(6)

196 where $|\cdot|$ is the rounding-to-nearest operator, and the function $clamp(\cdot)$ clips values outside the 197 integer range $[q_{min}, q_{max}]$. z is zero-point. s denotes the quantization scale factor, which reflects the proportional relationship between FP values and integers. $[q_{min}, q_{max}]$ is the quantization 199 range determined by the bit-width b. Generally, when we quantize the network's weight with 4-bit 200 and activations with 8-bit, called it as W4A8. We can calculate s offline using the activations 201 from calibration samples, known as **static quantization**. We can also use the runtime statistics of 202 activations to get s, referred to as **dynamic quantization**. More details are in Appendix A.11. 203

4 METHODOLOGY

205 In this section, we present *MQuant*, a post-training quantization solution specifically designed for 206 MLLMs. In Sec. 4.1, we describe modality-specific quantization (MSQ) and attention-invariant 207 flexible switching (AIFS). In Sec. 4.2, we introduce the Post LayerNorm-to-RMSNorm transforma-208 tion. In Sec. 4.3, we identify the weight outliers caused by the online Hadamard rotations and state 209 Rotation Magnitude Suppression (RMS). We provide the detailed MQuant algorithm for FP MLLMs 210 in Appendix A.5 Algorithm 1.

211

204

MODALITY-SPECIFIC QUANTIZATION AND ATTENTION-INVARIANT FLEXIBLE 4.1 212 SWITCHING 213

As discussed in Sec 1, increasing the resolution of input images or videos significantly raises the 214 number of visual tokens (Fig 1 (a)), resulting in a sharp increase in inference latency for traditional 215 per-token dynamic quantization methods. Although static per-tensor quantization is more efficient, it degrades performance compared to per-token dynamic quantization. Given the substantial differences
 between visual and textual tokens (Fig 1 (b)), we propose modality-specific quantization (MSQ) to
 address the heterogeneous nature of their distributions. However, the interleaved arrangement of
 input tokens complicates the efficient implementation of MSQ. To address this, we further introduce
 the attention-invariant flexible switching (AIFS) scheme, which reorganizes tokens by placing all
 image tokens before text tokens. This preserves the attention equivalence while eliminating the
 input-dependent dynamic position overhead, thus improving computational efficiency.

Modality-Specific Quantization. As shown in Figure 2, our proposed modality-specific quantization (MSQ) apply different static per-tensor scaling factors to visual and textual tokens respectively. For an input sequence E of length L containing mixed text and visual multimodal tokens. Here, without loss of generality, we define $E = \{e_1^t, ..., e_m^v, ..., e_n^t\} \in (E_v, E_t)\}$, where m and ndenote the start and end indices of the visual tokens. Then, we have:

228 229 230

258 259

$$E = \underbrace{(e_1^t, ..., e_{m-1}^t)}_{(e_1^t, ..., e_{m-1}^t)}, \underbrace{e_m^v, ..., e_n^v}_{e_m^v, ..., e_n^v} \underbrace{e_{n+1}^t, ..., e_L^t}_{(e_1^t, ..., e_L^t)}$$
(7)

MSQ effectively addresses the substantial distributional differences between visual and textual tokens

231 by applying distinct per-tensor scales factor for each modality. This separa-232 tion not only captures the unique dis-233 tribution of each modality but also pre-234 vents the significant time increase asso-235 ciated with token number, which grows 236 as image and video resolutions increase. 237 More discussions about differences be-238 tween MSQ and per-tensor static and 239 per-token dynamic quantization are in 240 Appendix A.11, highlighting its advan-241 tages. 242

Attention-Invariant Flexible Switching. Due to the interleaved and nonfixed positions of visual and textual to-



Figure 2: The depiction of modality-specific quantization (MSQ) and attention-invariant flexible switching (AIFS).

245 kens in the input sequence, directly applying MSQ introduces extra and irregular data process and 246 memory operation (e.g., slice, concat, pad), which reduces the computational efficiency of all the 247 GEMM-based¹ layers including QK and FC. To address this issue, we propose the attention-invariant 248 flexible switching (AIFS) scheme. AIFS rearranges the visual tokens to the front of the sequence and 249 places the textual tokens at the end, while simultaneously adjusting the causal mask in the attention 250 mechanism. This rearrangement preserves the computational consistency of both image and text tokens, enhancing computational efficiency. As shown in Figure 2, regardless of the input order of 251 visual and textual tokens, the AIFS scheme consistently reorganizes the tokens in this manner. 252

Incorporating the AIFS scheme necessitates an adaptation of the naive causal attention. Causal attention, also known as masked self-attention (Vaswani et al., 2017), is a fundamental structure in LLMs and MLLMs due to its auto-regressive properties. It ensures that each token can only attend to or be influenced by previous tokens, preventing any influence from future tokens. The naive casual mask for token sequence E can be formulated as follows:

$$\mathbf{A} = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{D}} + M_{i,j} \right), M_{i,j} = \begin{cases} 0 & \text{if } j \le i \\ -\infty & \text{if } j > i \end{cases}$$
(8)

After reordering token in AIFS, it is essential to maintain the causal relationships for the unified modality-decoupled token sequence $E^u = \{e_m^v, ..., e_n^v, e_1^t, ..., e_{m-1}^t, ..., e_{n+1}^t, ..., e_L^t \in (E_v, E_t)\}$. Consequently, the corresponding causal mask is modified, resulting in a unified causal mask:

$$M_{i,j}^{u} = \begin{cases} 0 & \text{if one of the following conditions is met:} \\ (i \le (n-m), j \le i \text{ or } (n-m) < j \le n) \\ \text{or } ((n-m) < i \le n, (n-m) < j \le i) \\ \text{or } (i > n, j \le i) \\ -\infty & \text{otherwise} \end{cases}$$
(9)

¹GEMM (General Matrix Multiplication) is a common operation in linear algebra for matrix multiplication.

This unified causal mask preserves the inherent causal relationships among tokens, ensuring numerical equivalence during attention computations even after token reordering. Additionally, we apply corresponding changes to the position embeddings to ensure that they align with the new token indices after AIFS. This adjustment is crucial for maintaining the computation equivalence of the attention computations, as it ensures that the position information accurately reflects the revised ordering of the tokens. More details can be found in the Appendix A.

Furthermore, MSQ can be naturally integrated with AIFS to form a per-modality static quantization approach. Extensive experiments demonstrate three key advantages brought by MSQ and AIFS: (1)
Computational Equivalence and Strong Compatibility, which can achieve SOTA quantization performance across 5 MLLMs (Table 2). (2) Reduced Inference Latency, achieving the speed of the per-tensor static quantization (Table 4). (3) Enhanced Memory and Computational Efficiency, achieving up to 24.7% speedup and 152.9% memory savings (Table 10.)

282 4.2 POST LAYERNORM-TO-RMSNORM TRANSFORMATION 283 Unlike the common Pre-LN structure in LLMs, MLLMs incorpo-284 rate both Pre-LN and Post-LN structure in their vision encoder. 285 As discussed in Sec 2, SliceGPT only designed a Pre-LN + Rotate 286 Scheme for LLMs. we further develop a Post-LN + Rotate scheme 287 to accommodate a wider range of vision encoders, enabling our 288 MQuant to be applicable to various MLLMs in both Pre-LN and 289 Post-LN configurations. The difference of LN style from 5 mainstream MLLMs are in Appendix 9. Here, we present how to 290 transform post-LN into RMSNorm layer while ensuring compu-291 tational invariance. Due to the recentering operation, LN exhibits 292 invariance to shifts, such that $LN(\mathbf{X}_k - a\mathbf{1}) = LN(\mathbf{X}_k), \forall a \in \mathbb{R}$. 293 Therefore, as shown in Figure 3. we can replace LN as RMSNorm 294 layer through adjusting the weights A_2 and bias b_2 of the the 295 linear ℓ_2 as follows: 1 296



Figure 3: Post-LN to RM-SNorm.

$$\hat{A}_2 = A_2 - \frac{1}{D} \mathbf{1} \mathbf{1}^T A_2, \hat{b}_2 = b_2 - \mu(b_2) \mathbf{1},$$
 (10)

Eq10 is the recenter operation in Figure 3 (b). Therefore, based on Eq10, the Post-LN can be replaced with an RMSNorm with the same arithmetic functionality and $X_{k+1} = X_k + RMSNorm(\ell_2(g(\ell_1(X_k))))$. Ultimately, we establish the equivalence of Post-LN and Post-RMSNorm Transformers. Now that every LN in MLLMs vision encoder has been converted to RMSNorm, where we can seamlessly integrate any orthogonal matrices for handling weights and activation outliers, while ensuring that the model's output remains arithmetically equivalent. 4.3 ROTATION MAGNITUDE SUPPRESSION

304

Theoretical Analysis of Hadamard Transformation on Weight Quantization. (Chee et al., 2024; Tseng et al., 2024) introduced incoherence (Eq. 11) to measure quantization difficulty. Lower incoherence means easier quantization. They also showed that applying a Hadamard transform to weights and activations effectively reduces incoherence. Quarot (Ashkboos et al., 2024b) employs offline Hadamard transforms on weights and activations, along with partial online Hadamard transfor LLMs, it does not meet expectations for MLLMs, as demonstrated in Table 2.

312 313

> 319 320

$$|V_{ij}| = |e_i^{\top} W e_j| \le \mu \frac{\|W\|_F}{\sqrt{mn}}.$$
 (11)

To address this, we analyze the incoherence of weight input channels before and after transformation, since per-channel quantization isn't applicable to input channels. For a weight matrix W_{ℓ_2} , according to (Eq. 11), we can compute an incoherence coefficient μ_{ℓ_2} satisfying: $\max(|W_{\ell_2}|) = \mu_{\ell_2} \frac{|W_{\ell_2}||_F}{\sqrt{mn}}$. Since the Hadamard transform preserves the Frobenius norm, $||HW_{\ell_2}||_F = ||W_{\ell_2}||_F$. Let the incoherence coefficient of the transformed weight HW_{ℓ_2} be $\mu_{H\ell_2}$. Then,

$$\frac{\mu_{H\ell_2}}{\mu_{\ell_2}} = \frac{\max(|HW_{\ell_2}|)}{\max(|W_{\ell_2}|)}.$$
(12)

For Hadamard transform matrices, the elements in the first row and first column are all $\sqrt{\frac{1}{n}}$, and the sums of the other rows and columns are zero. For the first input channel after transformation, it always holds that $HW_{\ell_20j} = \sum_{r=1}^{n} \frac{1}{\sqrt{n}} w_{rj} = \sqrt{n} \cdot \text{mean}(w_{i,j})$. For the other input channels 324 after transformation, as-325 suming the weights follow 326 a normal distribution with 327 mean μ and variance σ^2 , 328 the maximum value before transformation is μ + 329 $\sigma\sqrt{2\ln n}$, and the maxi-330 mum value after transfor-331 mation is $\sigma \sqrt{2 \ln n}$. From 332 the above equations, it is 333 evident that only the first 334 transformed input channel 335 may cause the maximum 336 value to increase. There-337 fore, when 338





(13)



, the Hadamard transform

339

340

341

342

343 344

365

will degrade the weight quantization performance. In Figure 4 (b), we show a weight from MLLMs identified using Equation 13, and its transformed distribution matches our derivations.

Rotation Magnitude Suppression 345 for Weight Outliers. To address 346 the weight outliers introduced by the 347 Hadamard transformation, we propose 348 rotation magnitude suppression RMS 349 mechanism. As shown in Figure 5, 350 we present the quantization process 351 for W4A8 setting. Since outliers only 352 appear in the first row of the input 353 channel weights, we separate this 354 channel weights and the corresponding 355 activations for the GEMV kernels. For the remaining weights, we perform a 356 GEMM kernels with the original activation 357 values. To ensure computational equiv-358 alence, we set the separated positional 359 weights to 0. Finally, we add the results 360 of the separated operations to obtain the



Figure 5: Schematic of RMS.

output. RMS is a plug-and-play module with minimal overhead and, when combined with Quarot or
 other Hadamard-based transformations, significantly enhances quantization performance. Algorithm
 2 outlines the specific procedure.

5 EXPERIMENTS

Models and Datasets. We evaluate our *MQUANT* on five MLLMs: InternVL2-8B (Chen et al., 2024a), Qwen-VL-Chat-9.6B (Bai et al., 2023b), MiniCPM-V 2.6-8B (Yao et al., 2024), Qwen2-VL-7B (Wang et al., 2024), and GLM-4V-9B (Hong et al., 2024). Evaluations are conducted on four benchmarks covering OCR and general question answering: TextVQA (Singh et al., 2019), DocVQA (Mathew et al., 2021), OCRBench (Liu et al., 2023), and MME (Fu et al., 2023), which assesses perception and cognition across 14 subtasks. These MLLMs' details are in Appendix A.9.

Baselines and Implementation Details. We test W8A8 and W4A8 quantization settings for both visual encoders and LLMs, comparing RTN, SmoothQuant (Xiao et al., 2022), and Quarot (Ashkboos et al., 2024b). Notably, we apply static per-tensor activation quantization for both components, unlike the dynamic per-token quantization typically used in existing MLLMs. The calibration dataset consists of 256 randomly selected samples from the corresponding benchmark training sets (Singh et al., 2019; Mathew et al., 2021; Liu et al., 2023). The batch size for latency evaluation is 1.

MLLMs	Method	Bits S Visual	Setting LLM	TextVQA Val↑	DocVQA Val↑	OCRBench ↑	MN
	-	BF16	BF16	77.65	90.97	794	22
	RTN			52.02	59.04	542	15
	SmoothQuant	WOAD	WAAS	59.88	59.75	544	15
InternVI 2-8B	Quarot	WORD	W 4A0	73.34	84.07	715	20
InternVL2-8B	MQuant (Ours)			77.49	90.27	785	2
	RTN			40.06	31.58	302	14
	SmoothQuant	W4A8	W4A8	46.48	31.21	310	1:
	Quarot			49.10	33.62	361	19
	MQuant (Ours)			76.62	88.42	725	2
	-	BF16	BF16	61.40	60.36	493	1
	RTN			0.45	0.03	189	6
	SmoothQuant	W8A8	W4A8	7.45	7.70	160	
	Quarot		W +A0	45.32	42.44	286	9
Qwen-VL-Chat-9.6B	MQuant (Ours)	-		61.16	59.31	483	1
	RIN			1.02	0.02	193	-
	SmoothQuant	W4A8	W4A8	8.59	4.28	188	ç
	Quarot			46.77	37.35	289	1
	MQuant (Ours)			60.50	58.72	473	1
	-	BF16	BF16	79.10	89.18	847	2
MiniCPM-V 2.6-8B	RTN		W4A8	61.00	65.16	332	1
	SmoothQuant	W848		62.40	65.76	424	1
	Quarot	WORO	W4A0	73.71	80.04	736	1
	MQuant (Ours)			80.41	89.15	844	2
	RTN	W448	W448	60.70	62.23	351	1
	SmoothQuant			65.67	60.02	455	1
	Quarot W4A8	10 W4A0	68.96	79.63	685	1	
	MQuant (Ours)			81.14	89.75	839	2
	-	BF16	BF16	84.43	93.87	842	2
	RTN			33.92	52.61	442	1
	SmoothQuant	W848	W448	49.11	53.97	444	1
	Quarot	0000	11-1-10	79.36	89.57	754	2
Qwen2-VL-7B	MQuant (Ours)			84.43	93.61	830	2
	RTN			40.20	38.82	422	1
	SmoothQuant	W4A8	W4A8	46.25	52.36	411	1
	Quarot			71.44	83.96	670	1
	MQuant (Ours)			84.32	93.58	824	2
	-	BF16	BF16	82.82	81.16	782	2
	RTN			7.05	3.70	0.00	1
	SmoothQuant	W/8 A 9	WAAP	9.05	4.10	0.00	1
	Quarot	woAd	w+A0	82.00	80.17	782	2
GLM-4V-9B	MQuant (Ours)			82.06	80.53	782	2
	RTN			7.61	3.60	0.00	1
	SmoothQuant	W4A8	W4A8	9.85	4.40	0.00	1
	Quarot		,, 1110	64.16	45.52	516	2
	MQuant (Ours)			81.58	79.67	754	2

378	Table 2: Comprehensive quantization results of different MLLMs across various evaluation datasets.
379	

416 5.1 OVERALL RESULTS

417 Weight-activation quantization results of various MLLMs. MOuant can be applied to the 418 quantization of various MLLMs. As shown in Table 2, our MQuant demonstrates significant improvements over several representative quantization methods. In W8A8 setting, MQuant achieves 419 performance nearly equivalent to that of FP models across all evaluation datasets. Notably, even in the 420 more challenging W4A8 setting, MQuant maintains comparable performance with FP models, while 421 other advanced quantization methods exhibit significant performance degradation. These results 422 indicate that our MQuant provide a more general and effective PTQ solution with strong compatibility 423 for maintaining high accuracy in MLLMs under various quantization settings. 424

Speedup and Memory Savings. We fixed the input sequence as "text-image-text" with 15 textual tokens, varying the image resolution from 280 × 280 to 5600 × 5600. Notably, the "text-image-text" sequence setting is not arbitrarily chosen; instead, it is a common setting in existing evaluation datasets (Duan et al., 2024). We evaluate speedup and memory savings by comparing PyTorch's BF16, AWQ (W4-only), and our MQuant (W4A8). O Speedup: As shown in Table 10, MQuant consistently achieves speedups over both PyTorch and AWQ across all resolutions, with a maximum of 24.76% over PyTorch at 840 × 840. Notably, MQuant outperforms AWQ, which is slower than PyTorch at most resolutions due to negative speedups. This significant speedup highlights the advantage of our

Table 3: Comparison of latency and memory saving with Pytorch (BF16), AWQ (W4-only) and ours MQuant (W4A8) on Qwen2-VL-7B. Pytorch and AWQ using the Qwen2-VL-7B official implementation. ↓ means lower values are better, ↑ means larger values are better.

Image size			Latency (s)		Memory (G)					
$H \times W$	Pytorch	AWQ↓	Speedup↑	MQuant↓	Speedup↑	Pytorch	AWQ↓	Improve↑	MQuant↓	Improve↑
280^{2}	0.257	0.286 (+0.029)	-10.14%	0.220 (-0.037)	+16.82%	16.45	7.45 (-9.00)	+120.67%	6.50 (-9.95)	+152.92%
840 ²	0.261	0.304 (+0.043)	-14.14%	0.210 (-0.051)	+24.76%	16.45	7.45 (-9.00)	+120.67%	6.50 (-9.95)	+152.92%
1400^{2}	0.559	0.652 (+0.093)	-14.29%	0.432 (-0.127)	+20.24%	16.90	7.90 (-9.00)	+113.92%	6.97 (-9.93)	+142.71%
1960^{2}	1.369	1.598 (+0.229)	-14.26%	1.112 (-0.257)	+16.63%	17.82	8.82 (-9.00)	+100.00%	7.85 (-9.97)	+119.93%
2520^{2}	2.820	3.175 (+0.355)	-11.14%	2.357 (-0.436)	+19.63%	19.04	10.05 (-8.99)	+89.59%	9.10 (-9.94)	+109.92%
3080^{2}	5.208	5.872 (+0.664)	-11.27%	4.488 (-0.720)	+16.02%	20.58	11.58 (-9.00)	+77.60%	10.61 (-9.97)	+96.45%
5600 ²	8.380	9.393 (+1.013)	-10.78%	7.469 (-0.911)	+12.19%	22.22	13.22 (-9.00)	+57.54%	12.25 (-9.97)	+61.65%

per-tensor static quantization, eliminating the overhead of token-wise scale computation. Even at higher resolutions (e.g., 5600²), MQuant maintains a 12.19% latency improvement, demonstrating scalability across various image sizes. **2** Memory Savings: MQuant offers substantial memory reductions compared to both PyTorch and AWQ. It consistently reduces memory usage by over 100% compared to PyTorch (e.g., 152.92% at 840²) and significantly outperforms AWQ's memory efficiency, achieving up to 101.07% savings at higher resolutions. These experiments demonstrate MQuant's strengths in both latency and memory savings, achieving up to 24.76% faster inference and reducing memory consumption by over 100% compared to baselines. This significant gain is attributed to MQuant's conversion of mixed input tokens into modality-decoupled tokens, eliminating irregular memory operations (e.g., slice, concat, pad) introduced by directly applying MSQ. This transformation reduces memory consumption and mitigates the computational overhead and increased latency caused by the surge in token counts from higher image and video resolutions.

454 5.2 ABLATION STUDY

In this section, we select Qwen2-VL-7B (Wang et al., 2024), currently one of the most powerful
 open-source MLLMs, to ablate the effectiveness of our proposed designs.

Table 4: The accuracy and speedup of our MSQ and AIFS scheme on the linear layer during prefill stage. The input sequence is structured as "text-image-text-image-text-image-text" with an image resolution of 2240×2240 and 50 textual tokens. Latency were tested on an NVIDIA RTX 6000 Ada Generation. \downarrow means lower values are better, \uparrow means larger values are better.

Activation	Weight	TextVQA Val↑	DocVQA Val↑	OCRBench ↑	MME↑	Latency \downarrow (s)	Speedup ↑
BF16	BF16	84.43	93.87	842	2319	1.690	-
W4-g128(AWQ)	BF16	83.93 (-0.50)	93.13 (-0.74)	828 (-14)	2252 (-67)	2.057 (+0.367)	-17.8%
A8-per-token dyn		84.32 (-0.11)	93.61 (-0.26)	830 (-12)	2269 (-50)	1.253 (-0.437)	+34.9%
A8-per-tensor sta	W4 par abannal sta	40.20 (-44.12)	38.82 (-54.79)	422 (-408)	1082 (-1187)	1.016 (-0.674)	+66.3%
A8-MSQ	W4-per-channel sta	84.32 (-0.11)	93.61 (-0.26)	830 (-12)	2269 (-50)	1.085 (-0.605)	+55.8%
A8-MSQ+AIFS		84.32 (-0.11)	93.61 (-0.26)	830 (-12)	2269 (-50)	<u>1.017</u> (-0.673)	+66.2%

Accuracy and Speedup Efficiency of the AIFS Scheme. We compare the accuracy and speedup results of our MSQ and AIFS. As previously discussed, the number of visual tokens in MLLMs increases quadratically with image resolution, and existing dynamic per-token quantization introduces additional inference latency due to the computation of token-wise scale factors.

Our AIFS transforms mixed visual and textual tokens into a unified multi-modal representation, enabling static per-tensor quantization for acti-vations and addressing the quadratic time complexity. As shown in Table 4, MSQ+AIFS achieves the same speedup efficiency as per-tensor static quantization while maintaining near-lossless accuracy comparable to the original Float model across all lin-ear layers in MLLMs. Furthermore,



Figure 6: The Speedup of AIFS+MSQ on Qwen2-VL-7B.

we evaluate the speedip of AIFS for visual tokens under varying image resolutions, comparing it
 with per-token setting and FP model. As shown in Figure 6, results indicate that as image resolution
 increases, AIFS achieves speedups around 20% and 80%. This substantiates the effectiveness and
 efficiency of our MSQ and AIFS design for MLLM quantization. Moreover, since AIFS requires

486 only a one-time rearrangement of the input data (adjusting the causal mask and token index offline), it does not alter the overall computation graph. We also report the acceleration results both in prefill 488 and decode stage in Table 7 of Appendix.

Table 5: Ablation study of proposed quantization designs on the OCR and MME benchmarks.

MLLMs	Methods		Bits Setting		T	D VOLVIA			• · · · · · ·						
	Static	AIFS + MSQ	LN2RN	RMS	Visual	LLM	lext v QA val	DocvQA var	OCRBench	MINE	Lat (IIIS)				
					BF16	BF16	84.43	93.87	842	2319	6523				
	 ✓ 	×	×	X			71.44	83.96	670	1911	5479				
Qwen2-VL-7B	 Image: A second s	\checkmark	×	X			111440	11/110	111440	111440	78.95	87.55	721	2095	5484
	\checkmark	\checkmark	\checkmark	X	W4A8	W4A8	82.48	91.91	803	2174	5452				
	\checkmark	\checkmark	\checkmark	\checkmark			84.32	93.58	824	2255	5471				

Ablation Study of Proposed Quantization Methods. We conduct ablation studies on Qwen2-VL-7B (Table 5) to evaluate our quantization methods. Starting from the 'Static' baseline (both LLM and visual components quantized using GPTQ with online Hadamard transformation), we progressively incorporate our designs. Applying AIFS and MDQ to the LLM part, while keeping the visual part quantized with GPTQ, significantly improves performance, demonstrating their effectiveness for the LLM. Adding LN2RN to the visual part allows us to apply offline and online Hadamard rotation to the vision encoder, further enhancing performance. Adding RMS further achieves performance close to the FP model, showing the effectiveness of our proposed designs while maintaining high accuracy.

Table 6: Comparative quantized results under different quantization settings on the OCR and MME benchmarks. † means re-implementation based on the official weight-only quantization setting with a group size of 128 (Wang et al., 2024).

MLLMa	MLI Ma Mathad		etting	ToxtVOA Valt	DeeVOA Valt	OCDBonab	мме
IVILLIVIS	Method	Visual	LLM	Iext vQA val	DocvQA val	OCKDENCI	IVIIVIT.
	-	BF16	BF16	84.43	93.87	842	2319
	GPTQ (g128) [†]	BF16	W8	84.33	93.97	842	2313
	GPTQ (g128) [†]	BF16	W4	84.18	93.25	831	2285
Qwen2-VL-7	$\mathbf{B} \mid AWQ \ (g128)^{\dagger}$	BF16	W4	83.93	93.13	828	2252
	MQuant (g128)	BF16	W4	84.55	93.18	832	2304
	MQuant (g128)	W4	W4	84.70	93.57	828	2292
	MQuant (per-channel)	W4A8	W4A8	84.32	93.58	824	2255

517 Weight-only Quantization. We compare our weight-only quantization results with the official Qwen2-VL-7B, which uses GPTQ (Frantar et al., 2022) and AWQ (Lin et al., 2023) to quantize only 518 the LLM while keeping the visual encoder in BF16. For fairness, we adopt the same group size of 128. 519 As shown in Table 6, our W4 LLM quantization aligns with their settings, achieving nearly lossless 520 accuracy compared to other methods. Extending W4 quantization to the visual encoder enables 521 4-bit quantization of the entire MLLM; MQuant still performs comparably to the BF16 model and 522 even surpasses other methods. Under the W4A8 weight-activation quantization setting, our results 523 remain consistent, with some metrics surpassing advanced weight-only methods. These experiments 524 confirm MQuant's effectiveness and robustness across various quantization configurations, whether 525 for weight-only or weight-activation quantization, and for partial or full quantization of MLLMs. 526

6 CONCLUSION 527

528 In this paper, we propose *MQuant*, the first accurate and efficient post-training quantization solution for multimodal large language models (MLLMs). Addressing their unique challenges, MQuant 529 reduces the time to first token (TTFT) with per-tensor static quantization and introduces modality-530 specific quantization (MSQ) to handle distribution discrepancies between visual adn textual tokens. 531 Besides, we propose attention-invariant flexible switching (AIFS) to ensure consistent positioning of 532 multimodal tokens and computational invariance, enhancing quantization efficiency. Furthermore, 533 we propose a unified LayerNorm to RMSNorm transformation, enabling seamless integration of the 534 MLLM vision encoder with Hadamard rotation. Finally, we reveal the weight outliers caused by the 535 online Hadamard rotations and propose Rotation Magnitude Suppression (RMS) to mitigate them. 536 Extensive experiments on five mainstream MLLMs demonstrate that MQuant attains state-of-the-art 537 PTQ performance, with the W4A8 model matching FP model accuracy. These findings highlight 538 MQuant's potential to advance MLLM quantization and enhance practical applicability.

487

489

497

498

499

500

501

502

503

504

505

506

507

540 REFERENCES 541

548

555

556

558

566

567

568

569

577

579

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 542 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. 543 *arXiv preprint arXiv:2303.08774*, 2023a. 3 544
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 546 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. 547 arXiv preprint arXiv:2303.08774, 2023b. 1
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel 549 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: A visual language 550 model for few-shot learning. NeurIPS, 35:23716-23736, 2022. 3 551
- 552 Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James 553 Hensman. Slicegpt: Compress large language models by deleting rows and columns. *International* Conference on Learning Representations (ICLR), 2024a. 2, 3, 21, 22 554
 - Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L Croci, Bo Li, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. arXiv preprint arXiv:2404.00456, 2024b. 3, 6, 7, 15, 17, 21
- JL Ba. Layer normalization. arXiv preprint arXiv:1607.06450, 2016. 20, 21 559
- 560 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, et al. Qwen technical report. arXiv preprint 561 arXiv:2309.16609, 2023a. 4 562
- 563 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. 564 arXiv preprint arXiv:2308.12966, 2023b. 2, 3, 7, 16, 18 565
 - Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in Neural Information Processing Systems, 2020. 1
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, et al. InternIm2 technical report, 2024. 4 570
- 571 Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher De Sa. Quip: 2-bit quantization of 572 large language models with guarantees. arXiv preprint arXiv:2307.13304, 2023. 17 573
- Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization of 574 large language models with guarantees. Advances in Neural Information Processing Systems, 36, 575 2024.3,6 576
- Zhe Chen, Weivun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi 578 Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. arXiv preprint arXiv:2404.16821, 2024a. 2, 7, 18
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong 581 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning 582 for generic visual-linguistic tasks. In Proceedings of the IEEE/CVF Conference on Computer 583 Vision and Pattern Recognition, pp. 24185–24198, 2024b. 3, 16 584
- 585 Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, et al. Mobilevlm: A fast, reproducible and strong vision language 586 assistant for mobile devices. arXiv preprint arXiv:2312.16886, 2023. 3
- 588 Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yiming 589 Hu, Xinyang Lin, Bo Zhang, et al. Mobilevlm v2: Faster and stronger baseline for vision language 590 model. arXiv preprint arXiv:2402.03766, 2024. 3
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-592 efficient exact attention with io-awareness. Advances in Neural Information Processing Systems, 35:16344–16359, 2022. 20

594	Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale.
595	<i>ICLR</i> , 2021. 3
596	

- Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang
 Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating large
 multi-modality models. In *Proceedings of the 32nd ACM International Conference on Multimedia*,
 pp. 11198–11201, 2024. 8, 19
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
 arXiv preprint arXiv:2407.21783, 2024. 1, 15, 21
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022. 3, 10, 19
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. MME: A comprehensive evaluation
 benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 7
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu
 Feng, Hanlin Zhao, Hanyu Lai, et al. Chatglm: A family of large language models from glm-130b
 to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024. 4
- Wenyi Hong, Weihan Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvlm2: Visual language models for image and video understanding. *arXiv preprint arXiv:2408.16500*, 2024. 2, 3, 7, 18
- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang,
 Yuxiang Huang, Weilin Zhao, et al. MiniCPM: Unveiling the potential of small language models
 with scalable training strategies. *arXiv preprint arXiv:2404.06395*, 2024. 4
- Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Barun Patra, et al. Language is not all you need: Aligning perception with language models. *NeurIPS*, 36, 2024. 3
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori,
 Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali
 Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL https://doi.org/10.5281/zenodo.
 5143773. If you use this software, please cite it as below. 3
- Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig
 Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient
 integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 4
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ICML*, pp. 19730–19742, 2023. 3
- Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-Gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024. 3
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation aware weight quantization for llm compression and acceleration. *arXiv preprint arXiv:2306.00978*, 2023. 3, 10
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 36, 2024. 3
- Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui
 Chen, Chunyuan Li, Lianwen Jin, et al. On the hidden mystery of OCR in large multimodal models. arXiv preprint arXiv:2305.07895, 2023. 7

684

- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. DocVQA: A dataset for VQA on document images. In *WACV*, pp. 2200–2209, 2021. 7
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pp. 8748–8763. PMLR, 2021. 3
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini
 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024. 1
- Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large language models. *CoRR*, abs/2308.13137, 2023. 1, 3
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,
 and Marcus Rohrbach. Towards VQA models that can read. In *CVPR*, pp. 8317–8326, 2019. 7
- Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with
 rotary position embedding, 2021. 15
- Fuwen Tan, Royson Lee, Łukasz Dudziak, Shell Xu Hu, Sourav Bhattacharya, Timothy Hospedales, Georgios Tzimiropoulos, and Brais Martinezs. Mobilequant: Mobile-friendly quantization for on-device language models. In *The 2024 Conference on Empirical Methods in Natural Language Processing*, 2024. URL https://openreview.net/forum?id=48ptWWA54E. 20
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a. 1, 3, 4, 15, 21
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b. 1, 3, 4, 15
- Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even
 better llm quantization with hadamard incoherence and lattice codebooks. *Forty-first International Conference on Machine Learning*, 2024. 6, 15, 17
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 2017. 4, 5, 22
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1, 2, 7, 9, 10, 16, 18
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,
 Lei Zhao, Xixuan Song, et al. CogVLM: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023. 1
- Kiuying Wei, Yunchen Zhang, Xiangguo Zhang, Ruihao Gong, Shanghang Zhang, Qi Zhang, Fengwei
 Yu, and Xianglong Liu. Outlier suppression: Pushing the limit of low-bit transformer language
 models. Advances in Neural Information Processing Systems, 2022. 3
- Guangxuan Xiao, Ji Lin, Mickael Seznec, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. *arXiv preprint arXiv:2211.10438*, 2022. 3, 7, 19
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li,
 Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024. 2, 3, 7, 16, 18

702	Zhewei Yao, Reza Yazdani Aminabadi, Minija Zhang, Xiaoxia Wu, Conglong Li, and Yuxiong He.
703	Zeroquant: Efficient and affordable post-training quantization for large-scale transformers, $arXiv$
704	preprint arXiv:2206.01861, 2022. 19
705	

- Zhihang Yuan, Lin Niu, Jiawei Liu, Wenyu Liu, Xinggang Wang, Yuzhang Shang, Guangyu Sun, Qiang Wu, Jiaxiang Wu, and Bingzhe Wu. Rptq: Reorder-based post-training quantization for large language models. *arXiv preprint arXiv:2304.01089*, 2023a. 1
- Zhihang Yuan, Yuzhang Shang, Yue Song, Qiang Wu, Yan Yan, and Guangyu Sun. Asvd: Activation aware singular value decomposition for compressing large language models. *arXiv preprint arXiv:2312.05821*, 2023b. 3
- Zhihang Yuan, Yuzhang Shang, Yang Zhou, Zhen Dong, Chenhao Xue, Bingzhe Wu, Zhikai Li, Qingyi Gu, Yong Jae Lee, Yan Yan, et al. Llm inference unveiled: Survey and roofline model insights. *arXiv preprint arXiv:2402.16363*, 2024. 3
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
 pp. 11975–11986, 2023. 3
- Biao Zhang et al. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019. 20, 21
- Ying Zhang, Peng Zhang, Mincong Huang, Jingyang Xiang, Yujie Wang, Chao Wang, Yineng Zhang,
 Lei Yu, Chuan Liu, and Wei Lin. Qqq: Quality quattuor-bit quantization for large language models.
 arXiv preprint arXiv:2406.09904, 2024. 3
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv* preprint arXiv:2303.18223, 2023. 3
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging LLM-as-a-judge with MT-Bench and
 Chatbot Arena. *NeurIPS*, 36, 2024. 3
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 3

1	4

756 A APPENDIX

A.1 ROTARY POSITION EMBEDDING FOR ATTENTION-INVARIANT FLEXIBLE SWITCHING

Many modern LLMs (Touvron et al., 2023a;b; Dubey et al., 2024) use rotary position embedding (RoPE) (Su et al., 2021) to encode information about the order of tokens in the input sequence. Rotary position embeddings are linear transformations applied to keys and queries defined as:

$$\mathbb{R}^{d_h}_{\Theta,m} = \begin{pmatrix} \cos i\theta_1 & -\sin i\theta_1 & 0 & 0 & \cdots & 0 & 0\\ \sin i\theta_1 & \cos i\theta_1 & 0 & 0 & \cdots & 0 & 0\\ 0 & 0 & \cos i\theta_2 & -\sin i\theta_2 & \cdots & 0 & 0\\ 0 & 0 & \sin i\theta_2 & \cos i\theta_2 & \cdots & 0 & 0\\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots\\ 0 & 0 & 0 & 0 & \cdots & \cos i\theta_{d_h/2} & -\sin i\theta_{d_h/2}\\ 0 & 0 & 0 & 0 & \cdots & \sin i\theta_{d_h/2} & \cos i\theta_{d_h/2} \end{pmatrix}$$
(14)

774 where $i \in [1, L]$ is the token index, $\Theta = \{\theta_i = 10000^{-2(i-1)/D}, i \in [1, 2, ..., D/2]\}$, and $\theta_i, i \in 1..D/2$ are predefined constants.

In the proposed Attention-Invariant Flexible Switching (AIFS) mechanism, we also apply the rearrangement for position embedding to maintain the computation equivalence. For a mixed input token $E = \{e_1^t, ..., e_m^v, ..., e_L^v\} \in (\mathbf{E}_v, \mathbf{E}_t)\}$ (in Eq 7), where m and n denote the start and end indices of the visual tokens. Specifically, after AIFS, the unified token can formulated as: $E = \{e_m^v, ..., e_m^v, e_1^t, ..., e_L^t\} \in (\mathbf{E}_v, \mathbf{E}_t)\}$. Therefore, the unified token indices after AIFS can be represented as:

$$(m, ..., n, 1, ..., m - 1, n + 1, ..., L) = AIFS(1, ..., m, ..., n, ..., L)$$
(15)

Due to we are aware of the final indices for input token after AIFS, than we can utilize the reorder token indices for Eq 14 to get the corresponding position embedding.

A.2 SPEEDUP OF MSQ+ AIFS.

Stage	BF16	Per-token Dynamic	Ours	Ours + GEMV	Speedup
Prefill	1690	1253	1017	-	+23%
Decode	17.5	16.4	13.06	8.2	+100%

Table 7: Latency (ms) comparison with Float and per-token dynamic quantization in W4A8 setting.

Here, we provide a more comprehensive comparison, we further report the acceleration results on the decode stage. Configuration is aligned with Table 4, measuring mean latency (ms) of linear layer for decoding 2,000 tokens. A custom kernel was implemented for W4A8 kernel GEMV operations. As shown in Table 7, compared to per-token dynamic quantization, in addition to achieving 23%speed improvement during the prefill stage, our method achieves an 100% speed up in decode stage. Overall, our AIFS+MSQ transforms the time-consuming online dynamic quantization into offline static quantization, achieving significant acceleration with almost no loss in accuracy, especially in long sequences. Notably, in practical applications, using **OpenAI's** token pricing as an example, our method can save $\approx 30\%$ in costs, and this effect is even more pronounced in other MLLMs, as visual tokens are more expensive than textual tokens.

A.3 WEIGHTS OUTLIERS AFTER ONLINE HADAMARD TRANSFORMATION

Following (Tseng et al., 2024; Ashkboos et al., 2024b) we make use of fast Hadamard transformation where convenient. Hadamard matrix is an orthogonal matrix with entries proportional to $\{+1, -1\}$.

A Walsh-Hadamard matrix is a square matrix of size 2^n with For a Hadamard matrix:

$$\boldsymbol{H}_{2} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1\\ 1 & -1 \end{bmatrix} \quad \text{and} \quad \boldsymbol{H}_{2^{n}} = \boldsymbol{H}_{2} \otimes \boldsymbol{H}_{2^{n-1}}.$$
(16)

$$(AB)_{ij} = \sum_{r=1}^{N} a_{ir}b_{rj} = a_{i1}b_{1j} + a_{i2}b_{2j} + \dots + a_{in}b_{nj}$$
(17)

Thereby the HW_{ℓ_2} can be formulated as:

$$(\boldsymbol{H}W_{\ell_2})_{ij} = \sum_{r=1}^{\infty} h_{ir} w_{rj} = h_{i1} w_{1j} + h_{i2} w_{2j} + \dots + h_{in} w_{nj}$$
(18)

where $\boldsymbol{H} \in \mathbb{R}^{d_{in} \times d_{in}}$ and $W_{\ell_2} \in \mathbb{R}^{d_{in} \times d_{out}}$, d_{in} and d_{out} is the dimension of input and output of weight W_{ℓ_2} . Due to the properties of the Hadamard matrix H, whose first row consists entirely of 1, for the first row in (HW_{ℓ_2}) , $(HW_{\ell_2})_{0j} = \sum_{r=1}^n w_{rj}$, due to the property of Hadamard matrix H. So, the values in W_{ℓ_2} are subject to continuous accumulation and summation, resulting in the exists of outliers in the first row of the output matrix HW_{ℓ_2} . Notably, the matrix $HW_{\ell_2}Q$ still has the same problem, for simplicity, we omit the matrix Q in the main paper.

IMAGE TOKENS IN VARIOUS MLLMS A.4





As shown in Fig 7, for different MLLMs (Bai et al., 2023b; Wang et al., 2024; Yao et al., 2024; Chen et al., 2024b), the number of prefill visual tokens grows as image resolution increases. This rapid expansion in token count exacerbates the inference latency, particularly in per-token dynamic quantization, which requires separate processing for each token, such as memory access and scale computation. As a result, the TTFT in MLLMs increases drastically, severely impacting overall inference latency. Moreover, in higher-demand scenarios such as video-based tasks and multi-image dialogues, the accumulation of visual tokens becomes even more pronounced, further exacerbating the increase in TTFT.

A.5 MQUANT ALGORITHM

Here, we present our *MQuant* algorithm for MLLMs in Algorithm 1.

 Algorithm 1 MQuant Quantization Algorithm Input: Full-precision (FP) MLLM model with a vision encoder E, visual-language projector P, a a large language model LLM; Calibration dataset D^c. Output: 	ind
Input: Full-precision (FP) MLLM model with a vision encoder E , visual-language projector P , a large language model LLM; Calibration dataset D^c . Output :	and
a large language model LLM; Calibration dataset D^c . Output:	
⁸⁶⁷ Output:	
• For E and P : per-channel weight scale s_w^E , per-channel weight zero-point z_w^E , per-ten activation scales s_a^E , per-tensor activation zero-point z_a^E .	SOI
• For LLM: per-channel weight scale s_w^{llm} , per-channel weight zero-point z_w^{llm} , per-ten activation scales $s_{a^v}^{llm}$ for visual tokens and $s_{a^t}^{llm}$ for textual tokens, per-tensor activat zero-points z^{llm} for visual tokens and z^{llm} for textual tokens	sor ion
2210 points z_{a^v} for visual tokens and z_{a^s} for textual tokens.	
 Apply Hadamard Rotation to the LLM Part as described: Apply the offline and online Hadamard rotations to all the weights and activations in LLM. Ownetics Weights of the LLM: 	
5. Qualifize weights of the LLM. 4. Input the calibration dataset D^c to the FP MI I M	
5: Use GPTO to quantize the weights for LLM, obtaining per-channel weight quantization para	me-
ters s_w^{llm} and z_w^{llm} .	
6: For the LLM Part:	
(a) Input D^c to the FP MLLM.	
(b) Compute per-tensor activation quantization parameters $s_{a^v}^{llm}$ and $s_{a^t}^{llm}$, $z_{a^v}^{llm}$ and z	\lim_{a^t}
for visual and textual tokens respectively, based on the proposed Modality-Spec	ific
884 Quantization (MSQ) in Sec 4.1.	
(c) Reorder the input mixed token sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^t$ to a to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^t$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^t$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^t$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to be a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_L^v$ to be a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_m^v$ to be a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_m^v$ to be a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_m^v$ to be a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_m^v$ to be a sequence from $E = e_1^t, \dots, e_m^v, \dots, e_m^v$ to be a sequence from $E = e_1^t, \dots, e_m^v$ to be a sequence from $E =$	ıni-
1886 field modality-decoupled sequence $E^{u} = e_{m}^{v}, \dots, e_{n}^{v}, e_{1}^{v}, \dots, e_{m-1}^{v}, \dots, e_{n+1}^{v}, \dots, e_{n+1}^{v}, \dots, e_{m-1}^{v}, \dots, \dots, e_{m-1}^{v}, \dots, \dots, e_{m-1}^{v}, \dots, \dots, e_{m-1}^{v}, \dots, \dots,$	e_L^{ι}
⁸⁸⁷ 7 The formula is a second secon	•••
 ⁽²⁾ Iransform all the LayerNorm to RMISNOrm in MILLM Vision encoder <i>E</i> and Visual-Langui ⁸⁸⁸ Projector <i>D</i> using the proposed Post I overNorm to PMSNorm transformation in Sec. 4.2 	ige
⁸⁸⁹ 8. Apply the offline and online Hadamard rotations to all the weights and activations in <i>E</i> and	Р
890 9: Quantize Weights of E and P :	
⁸⁹¹ 10: Input D^c to the transformed FP vision encoder E and Visual-Language Projector P .	
11: Use GPTQ to quantize E and P , obtaining per-channel weight quantization parameters s_w^E a z_w^E .	ınd
 Address the weight outliers using caused by online Hadamard based on the proposed Rotat Magnitude Suppression (RMS) in Sec 4.3. 	ion
895 Wagintude Suppression (KWS) in Sec 4.5.	

A.6 EFFECTIVENESS OF ROTATIONAL MAGNITUDE SUPPRESSION FOR WEIGHT OUTLIERS IN LLMs.

902 For LLMs, compared 903 to the original Quarot, T 904 integrating RMS with w 905 Quarot leads to perfor-906 mance improvements 907 across LLaMA2 models 908 with 7B, 13B, and 70B 909 parameters, as detailed in Table 8. For MLLMs, 910 ablation studies presented 911

Table 8: WikiText-2 PPL on 4-bit quantization of LLAMA-2 models
with 2048 sequence length.

Method	Weight Quantization	LLAMA-27B	LLAMA-2 13B	LLAMA-2 70B
Baseline	-	5.47	4.88	3.32
QuaRot	GPTQ	6.10	5.40	3.79
QuaRot	RMS + GPTQ	6.04	5.32	3.67

in Table 5 demonstrate that the RMS method significantly enhances quantization performance. 912

913 A.7 WEIGHT OUTLIERS FROM ON ROTATION

899 900

901

914

Previous study Quarot (Ashkboos et al., 2024b) leveraged incoherence processing (Chee 915 et al., 2023; Tseng et al., 2024) to make the weights matrix and activations easier to quan-916 tize by multiplying randomized Hadamard matrix. Here, we briefly summarize the pipeline 917 of Quarot for LLMs. As shown in Figure 8 (a), we select a randomized Hadamard ma918 trix with size that matches the hidden dimension of the model, denoted as Q. (1) the 919 input X is incoherence process to enable improved quantization, becoming $X \leftarrow XQ$; 920 (2) we remove the scaling

921 operation from RMSNorm 922 $(\operatorname{diag}(\boldsymbol{\alpha}))$ and absorbing into the subsequent weight ma-923 trices W_{ℓ_1} of input linear; 924 (3) to cancel the impact of 925 Q, Q^{\perp} is applyed to weight 926 matrices W_{ℓ_1} , then $W_{\ell_1} \leftarrow$ 927 $Q^{\top}(\alpha)W_{\ell_1}$; (4) to improve 928 the quantize-ability of the acti-929 vations within each block, in-930 sert online Hadamard opera-931 tions before the output linear 932 ℓ_2 ; (5) to retain the computa-933 tional invariance, Hadamard matrix H is implicitly re-934 versed by fusing into the 935



Figure 8: (a) The pipeline of Quarot, (b) transformed weights outliers. output linear matrix W_{ℓ_2} : $W_{\ell_2} \leftarrow HW_{\ell_2}Q$. Without loss of generality, the weight matrices W_{ℓ_1} can also be query/key/valueprojection matrix, and the weight matrices W_{ℓ_2} can also be output weight matrix, please refer to Quarot for more quantization details.

A.8 RMS ALGORITHM

936

937

938

939 940

941 942

943

960

961

Here, we present our *RMS* design with Quarot in Algorithm 2.

Re	quire: An LLM or MLLM model
Er	sure: Quantized model with RMS
1	Initialize an empty list marks
2	for each linear layer in the model do
3	if Layer satisfies Equation 13 then
4	Mark the layer and append its ID to marks
5	end if
6	end for
7	Apply Quarot or other Hadamard-based transformations to the model
8	for each layer ID in marks do
9	Modify the layer's implementation using the RMS method
10	end for
11	Quantize the model
12	return model

COMPARISON OF DIFFERENT MLLMS: INPUT PRE-PROCESS, LAYERNORM A.9 ARCHITECTURE IN VISION ENCODER, MODEL PARAMETERS AND FLOPS.

In this section, we compare the visual input pre-process methods, LayerNorm structures in MLLM 962 vision encoder, model parameters, and Flops across five mainstream MLLMs: InternVL2-8B (Chen 963 et al., 2024a), Qwen-VL-Chat-9.6B (Bai et al., 2023b), MiniCPM-V 2.6-8B (Yao et al., 2024), 964 Qwen2-VL-7B (Wang et al., 2024), and GLM-4V-9B (Hong et al., 2024). This comparison highlights 965 the architectural differences between these models, particularly the TTFT sensitivity to input image 966 resolution, illustrating the unique challenges these variations present for quantization and efficient 967 inference. 968

969 A.10 SPEEDUP AND MEMORY SAVINGS WITH SCALING IMAGE RESOLUTION 970

We fixed the input sequence as "text-image-text" with 15 textual tokens and presente the detailed 971 changes of speed and memory, varying the image resolution from 280×280 to 5600×5600 . Notably, 972 Table 9: Comparison of TTFT sensitivity to image resolution, model Parameters and flops in 973 mainstream MLLMs.[†] means the Flops values are measured with the number of visual tokens is 256.

MITM.	TTET:	T	Params (B)		FLOPs (T) [†]	
MLLMS	I I F I's sensitivity to input image resolution	LayerNorm	Visual	LLM	Visual	LLM
nternVL2-8B	TTFT increases with input image aspect ratio	Pre-LN	0.34	7.74	1.28	7.54
Qwen-VL-Chat-9.6B	Fixed input resolution (448×448)	Pre-LN	1.94	7.73	4.23	3.70
MiniCPM-V 2.6-8B	TTFT increases with input image aspect ratio	Pre-LN	0.49	7.61	4.16	3.64
Qwen2-VL-7B	TTFT increases quadratically with input image resolution	Pre-LN	0.68	7.61	1.31	3.61
GLM-4V-9B	Fixed input resolution (1120×1120)	Post-LN	4.51	8.78	12.10	4.70

Table 10: Comparison of latency and memory saving with Pytorch and AWQ on Qwen2-VL-7B.[‡] means the Qwen2-VL-7B official implementation.

I	mage size	Pytorch [‡] (BF16)		AWQ [‡] (W4-only)		MQuant (W4A8)	
	$H \times W$	Latency(s)	Memory(G)	Latency(s)	Memory(G)	Latency(s)	Memory(G)
	280^{2}	0.257	16.45	0.286 (-10.14%)	7.45 (+120.67%)	0.220 (+16.82%)	6.50 (+152.92%)
	560^{2}	0.252	16.45	0.292 (-13.70%)	7.45 (+120.67%)	0.211 (+19.48%)	6.50 (+152.92%)
	840^{2}	0.261	16.45	0.304 (-14.14%)	7.45 (+120.67%)	0.210 (+ 24.76%)	6.50 (+152.92%)
	1120^{2}	0.326	16.58	0.384 (-15.10%)	7.56 (+119.51%)	0.262 (+24.48%)	6.61 (+151.59%)
	1400^{2}	0.559	16.90	0.652 (-14.29%)	7.90 (+113.92%)	0.432 (+ 20.24%)	6.97 (+142.71%)
	1680^{2}	0.881	17.33	1.066 (-17.39%)	8.33 (+108.57%)	0.705 (+18.23%)	7.40 (+130.27%)
	1960^{2}	1.369	17.82	1.598 (-14.26%)	8.82 (+100.00%)	1.112 (+ 16.63 %)	7.85 (+119.93%)
	2240^{2}	2.013	18.40	2.294 (-12.24%)	9.40 (+95.74%)	1.653 (+ 17.83%)	8.44 (+117.84%)
	2520^{2}	2.820	19.04	3.175 (-11.14%)	10.05 (+89.59%)	2.357 (+ 19.63%)	9.10 (+109.92%)
	2880^{2}	3.883	19.77	4.345 (-10.64%)	10.77 (+83.81%)	3.297 (+17.69%)	9.82 (+101.07%)
	3080^{2}	5.208	20.58	5.872 (-11.27%)	11.58 (+77.60%)	4.488 (+16.02%)	10.61 (+96.45%)
	3360^{2}	6.814	21.46	7.548 (-9.73%)	12.46 (+70.65%)	6.004 (+ 13.41 %)	11.50 (+83.01%)
	3640^{2}	8.360	22.22	9.377 (-10.84%)	13.22 (+57.54%)	7.469 (+11.91%)	12.25 (+61.65%)
	4480^{2}	8.349	22.22	9.379 (-10.97%)	13.22 (+57.54%)	7.469 (+11.71%)	12.25 (+61.65%)
	5600^{2}	8.380	22.22	9.393 (-10.78%)	13.22 (+57.54%)	7.469 (+12.19%)	12.25 (+61.65%)

997 998

982

999 the "text-image-text" sequence setting is not arbitrarily chosen; instead, it is a common setting 1000 in existing evaluation datasets (Duan et al., 2024). We evaluate speedup and memory savings by 1001 comparing PyTorch's BF16, AWQ (W4-only), and our MQuant (W4A8). Speedup: As shown in 1002 Table 10, MQuant consistently achieves speedups over both PyTorch and AWQ across all resolutions, 1003 with a maximum of 24.76% over PyTorch at 840×840 . Notably, MQuant outperforms AWQ, which is slower than PyTorch at most resolutions due to negative speedups. This significant speedup 1004 highlights the advantage of our per-tensor static quantization, eliminating the overhead of token-wise 1005 scale computation. Even at higher resolutions (e.g., 5600^2), MQuant maintains a 12.19% latency 1006 improvement, demonstrating scalability across various image sizes. Memory Savings: MQuant offers 1007 substantial memory reductions compared to both PyTorch and AWQ. It consistently reduces memory 1008 usage by over 100% compared to PyTorch (e.g., 152.92% at 840²) and significantly outperforms 1009 AWQ's memory efficiency, achieving up to 101.07% savings at higher resolutions. These experiments 1010 demonstrate MQuant's strengths in both latency and memory savings, achieving up to 24.76% faster 1011 inference and reducing memory consumption by over 100% compared to baseline methods. This 1012 makes MQuant a more efficient solution for deploying MLLMs in resource-constrained environments. 1013

1014

1017

1015 1016

A.11 QUANTIZATION GRANULARITY

Furthermore, as mentioned in SoomthQuant (Xiao et al., 2022), there are different different granular-1018 ity levels. The per-tensor static quantization uses a single step size for the entire matrix. Per-token 1019 **dynamic** quantization employs different s for the activations associated with each token, being a 1020 common granularity for activations quantization of existing LLMs. For weights, per-channel quanti-1021 zation applies distinct s for each output channel of the weights, while group-wise quantization utilizes 1022 a coarse-grained s for different channel groups. Notably, group-wise quantization is a prevalent 1023 granularity for weight quantization in LLMs (Frantar et al., 2022; Yao et al., 2022). Please refer to 1024 appendix for more quantization basics. 1025

1026 A.12 ADVANTAGE OF MSQ AND AIFS

In per-tensor static quantization, the quantization parameters (i.e., scale and zero-point) are precomputed for an entire tensor (e.g., weights or activations) and remain fixed throughout inference. While
 efficient, this approach often leads to large and unacceptable accuracy loss in MLLMs due to their
 diverse activation distributions across varying inputs.

1032 In contrast, per-token dynamic quantization computes quantization parameters on-the-fly for each 1033 input token during inference. This approach incurs significantly higher computational overhead, 1034 as the quantization parameters must be recalculated for every input token, along with multiple 1035 additional memory traversals. Such requirements make per-token dynamic quantization unfriendly or impractical for edge devices and some AI accelerators, which struggle with fine-grained dynamic 1036 operations Tan et al. (2024). This issue is especially severe in MLLMs, where the token count 1037 increases significantly with higher image resolution or more video frames. The Modality-Specific 1038 Quantization (MSQ) in MQuant is a novel per-modality quantization approach specifically designed 1039 to address the unique challenges of MLLMs quantization. 1040

Furthermore, MSQ can be naturally applied to the unified modality-decoupled tokens generated by 1041 AIFS. By integrating MSQ and AIFS, our designs yields three key advantages: (1) Computational 1042 Equivalence and Strong Compatibility: The unified causal mask and token index introduced by 1043 AIFS preserves the inherent causal relationships among tokens, ensuring numerical equivalence 1044 during attention computations. Moreover, since AIFS requires only a one-time rearrangement of the 1045 input data (adjust causal mask and token index in offline), it does not alter the overall computation 1046 graph. This characteristic allows for seamless integration with other LLM inference acceleration 1047 methods, such as FlashAttention (Dao et al., 2022), ensuring both computational equivalence and 1048 strong compatibility. As shown in Table 2, MQuant achieves SOTA quantization performance across 1049 5 mainstream MLLMs. (2) Reduced Inference Latency: MSQ not only addresses the substantial 1050 distributional differences between modalities but also mitigates the significant computational overhead 1051 and increased inference latency caused by the surge in token counts from higher image and video resolutions. As shown in Table 4, MSQ+AIFS significantly reduces latency from 2.057s to 1.1017s, 1052 closely matching the speed of the per-tensor static setting while maintaining near-lossless accuracy 1053 comparable to the original Float model. (3) Enhanced Memory and Computational Efficiency: By 1054 combining MSQ and AIFS, we convert mixed input tokens into unified, modality-decoupled tokens, 1055 eliminating the irregular memory operations (e.g., slice, concat, pad) introduced by directly applying 1056 MSQ. This transformation reduces memory consumption and improves efficiency of GEMM kernel, 1057 which would otherwise be compromised by the interleaved and non-fixed positions of visual and 1058 textual tokens. As shown in Table 10, MQuant can achieve up to 24.7% speedup and 152.9% memory 1059 savings. 1060

1061 A.13 LAYERNORM, RMSNORM AND COMPUTATIONAL INVARIANCE

1067 1068

1077 1078

We introduce LayerNorm, RMSNorm, the computational Invariance, and their usage in Transformers.

Layer Normalization (LayerNorm, LN) (Ba, 2016) is a technique to normalize the activations of intermediate layers of neural networks. Given a vector $x \in \mathbb{R}^d$, LayerNorm normalizes it to obtain a zero-mean unit-variance vector,

LayerNorm
$$(\boldsymbol{x}) = \frac{\boldsymbol{x} - \mu(\boldsymbol{x})\mathbf{1}}{\sqrt{\boldsymbol{x}_2^2/d - \mu^2(\boldsymbol{x}) + \epsilon}},$$
 where $\mu(\boldsymbol{x}) = \frac{\mathbf{1}^T \boldsymbol{x}}{d}, \epsilon > 0.$ (19)

LayerNorm recenters and rescales the activations and gradients in the forward and backward computations, which enables fast and robust training of neural networks.

Root Mean Square Normalization (RMSNorm) (Zhang et al., 2019) is another technique used for normalizing the activations. It is similar to LayerNorm in that it aims to accelerate and stabilize the training but uses a different normalization approach. Instead of normalizing the inputs based on their mean and variance, RMSNorm normalizes them based on their root mean square (RMS) value. It is defined in the following equation,

$$\text{RMSNorm}(\boldsymbol{x}) = \frac{\boldsymbol{x}}{\sqrt{\boldsymbol{x}_2^2/d + \epsilon}}, \text{ where } \epsilon > 0. \tag{20}$$

1079 RMSNorm only rescales the input vector and the corresponding gradients, discarding the recentering process. As shown in their definitions, RMSNorm is computationally simpler and more efficient

1080 than LayerNorm. It is reported that replacing LayerNorm with RMSNorm can achieve comparable performance and save training and inference time by 7% - 64% (Zhang et al., 2019). 1082

Given a zero-mean vector \boldsymbol{x} , these two kinds of normalization are equivalent. Formally, if $\mu(\boldsymbol{x}) = 0$, 1083 then LayerNorm(x) = RMSNorm(x). We may optionally introduce learnable parameters and apply 1084 an element-wise affine transformation on the output of LayerNorm and RMSNorm.

1086 **LayerNorm** (LN) and RMSNorm Given the input concated token X after embeddings with the 1087 shape $L \times D$, the X is passed through a LayerNorm (Ba, 2016) operation, which subtracts the mean 1088 from each row of the matrix, divides the row by its standard deviation, rescales (columnwise), and 1089 adds an offset. Follow (Ashkboos et al., 2024a), we write the LayerNorm block as

1090 1091

1100

1110

1132

$$LayerNorm(X) = RMSNorm(XM)diag(\alpha)\sqrt{D} + \mathbf{1}_N\beta^{\top}$$
(21)

1092 where RMSNorm(X) applies $x \leftarrow x/||x||$ to each row of X, and $X = concat(E_v, E_t))$ is the 1093 concatenation between text tokens E_t and the visual tokens E_n . The vector parameter α and offset 1094 (vector) parameter β are learned independently at each LayerNorm instance. The constant matrix 1095 $M = I - \frac{1}{D} \mathbf{1} \mathbf{1}^{\top}$ is a $D \times D$ matrix which subtracts the mean from each row of X, called recentering 1096 operation. Formally, if M = I, the input X has a zero-mean, the Eq 21 is equivalent to RMSNorm. 1097 Specifically, LayerNorm is widely employed in visual encoders E, whereas RMSNorm (Zhang et al., 2019) is commonly used in LLMs (Touvron et al., 2023a; Dubey et al., 2024) and has been shown to 1099 accelerate training and inference time with similar performance (Zhang et al., 2019).

Computational Invariance in RMSNorm. Based on the *computational invariance*, recent stud-1101 ies (Ashkboos et al., 2024a;b) have shown that orthogonal transformations can effectively smooth 1102 outliers and improve the quantize-ability of both weights and activations. In particular, for 1103 transformers, inserting linear layers with an orthogonal matrices Q before and after the RM-1104 SNorm (Zhang et al., 2019) layer in a transformer, the network remains unchanged. In detail, 1105 given the input X and orthogonal matrix Q for RMSNorm layer, the *computational invariance* 1106 means: RMSNorm $(XQ)Q^{\top}$ = RMSNorm. Here, $Q^{\top}Q = QQ^{\top} = I$ and a rotation matrix is an 1107 orthogonal matrix with |Q| = 1. Note that multiplying a vector x by Q does not change the norm of 1108 the vector, since $||Qx|| = \sqrt{x^\top Q^\top Qx} = \sqrt{x^\top x} = ||x||$. 1109

A.14 LAYERNORM TO RMSNORM TRANSFORMATION.





Pre-LayerNorm to RMSNorm Transformation. Here, we propose unified LayerNorm-to-1133 RMSNorm transformation, aiming to synthesize the transformer architecture of MLLMs' vision

1134 encoders and LLMs, endowing them with rotation-friendly characteristics that facilitate the effective 1135 removal of outliers. We take Pre-LN transformer as an example to show that how transform Pre-LN 1136 into RMSNorm layer while guaranting arithmetic equivalence. As shown in Figure 10 (a), for the input X_k of the k-th block, the main block in pre-LN transformer is $X_{k+1} = X_k + \ell_2(g(\ell_1(LN(X_k)))))$, 1137 1138 where $k \in [1, N]$, and N is the block number. If g is an activation function, such as GELU, this block is a multi-layer perceptron (MLP) module. If g is a multi-head attention, then this 1139 block is the casual attention module (Vaswani et al., 2017). Due to the recentering operation, 1140 LN exhibits invariance to shifts, such that $LN(\mathbf{X}_k - a\mathbf{1}) = LN(\mathbf{X}_k), \forall a \in \mathbb{R}$. Therefore, 1141 as shown in Figure 10 (b), we can replace LN as RMSNorm layer through two modifications: 1142

0 recenter the input X_k to $X_k - \mu(X_k)\mathbf{1}$, ensuring 1143 that the input to norm layer maintain a zero mean. 1144 adjust the weights A_2 and bias b_2 of the linear ℓ_2 1145 to $\hat{A}_2 = A_2 - \frac{1}{D} \mathbf{1} \mathbf{1}^T A_2, \hat{b}_2 = b_2 - \mu(b_2) \mathbf{1}$. Con-1146 sequently, the LN can be replaced with an RMSNorm 1147 layer with the same arithmetic functionality. The first op-1148 eration is to recenter X_k , while the second operation is 1149 to recenter the output of main branches. Notably, since 1150 $X_{k+1} = X_k + \ell_2(g(\ell_1(LN(X_k)))))$, after applying **1** and 1151 2, the input and the output of main branch are re-centered with zero-mean, while the input of residual branches also 1152 maintain a zero mean. Therefore, the output after current 1153 blocks, X_{k+1} (which serves as the input for next block), 1154 still maintain zero-mean. A detailed proof is provided in 1155 the Appendix. Ultimately, we establish the equivalence of 1156 Pre-LN and Pre-RMSNorm Transformers. Now that every 1157 LayerNorm in the transformer has been converted to RM-1158 SNorm in MLLMs, we can use any orthogonal matrices Q



Figure 10: The illustration of transformation from Pre-LN to RMSNorm.

to the model. Therefore, the visual encoder and LLMs are in a rotation-friendly RMSNorm-only transformer architecture.

1162 A.15 TRANSFORMER FORWARD PASS

1161

1181

Here, we refer from SliceGPT (Ashkboos et al., 2024a) to describe the transformer forward pass
process. As described in Figure 10 illustrates part of a transformer network: an attention block
connected to a Feed Forward Network (FFN) block through a LayerNorm or RMSNorm block, with
residual connections.

Embeddings Let *D* be the embedding dimension of our transformer, *N* be the sequence length. The transformer model takes as input a sequence of token IDs and position IDs, and uses them to index the embedding matrices, producing the initial signal X with shape $N \times D$. In what follows we consider, without loss of generality, a single embedding matrix W_{embd} indexed by input sequence *s*.

Attention Blocks in LLMs The attention block has four matrices: W_k, W_q, W_v and W_o , each of 1172 dimension $D \times D$. The input signal arriving into the block is projected into the Key (**XW**_k), Query 1173 (XW_q) , and Value (XW_v) matrices, which are then split into multiple heads. A nonlinear operation 1174 is applied at each head before the signals are combined and multiplied by the output weight matrix 1175 W_o . Since the first three weight matrices are applied separately to the inputs, we can concatenate 1176 them and perform a single matrix multiplication. We can consider the concatenation of these matrices 1177 to be a single linear layer, which we denote as Linear l_1 , as shown in Figure 10. The MHSA in 1178 Figure 10 means the attention computation. We also refer to the output linear as Linear2 ℓ_2 , as 1179 shown in Figure 10. We treat the attention block as MHSA(XW_1) W_2 , where MHSA represents the 1180 multi-head attention operation.

FFN Blocks The other type of block that appears in transformer architectures is a Feed Forward Network (FFN) block. In many cases, this is a Multi-layer Perceptron (MLP), which also can be represented with Figure 10. A Linear1 layer ℓ_1 with weight matrix W_1 , followed by an non-linear operation (i.e., GELU function), and Linear2 layer: ℓ_2 with weight matrix W_2 . We can therefore denote the operation of MLP or gated FFN layers as GELU(XW_1) W_2 .

- **Language Modelling (LM) Head** All of the transformer networks to which we apply in this paper have a decoder-only structure: after multiple layers applying alternating attention and FFN blocks, a
 - 22

head block computes logits which are used to compute the loss during training and token prediction on deployment. The head operation is $XW_{head} + b_{head}$.

¹¹⁹⁰Once the model is trained and all of the parameters are set, the computations required in a transformer network to produce predictions involve passing signal matrices from one block to the next until the head node is reached. Since we are able to define both FFN and attention blocks in the form $\sigma(XW_1)W_2$, where we understand that σ represents either a point-wise (GELU function) or multihead-attention nonlinearity (MHSA), we are able to describe the forward pass using Algorithm 3.

Algorithm 3 Pre-LN Transformer	
Require: $\{W_1^{\ell}, W_2^{\ell}\}_{\ell=1}^L$ Require: $\{\sigma_{\ell}\}_{\ell=1}^L$ Require: $\{Norm_{\ell}\}_{\ell=0}^L$ Require: $W_{embd}, W_{head}, b_{head}$ Require: s	weights and biases of FFN and attention blocks nonlinearity associated with each block LayerNorm or RMSNorm instances to perform between blocks embedding and head matrices input sequence
$oldsymbol{X} \leftarrow oldsymbol{W}_{ ext{embd}}[oldsymbol{s},:] \ oldsymbol{X} \leftarrow ext{Norm}_0(oldsymbol{X})$	index embeddings normalize
: for $\ell = 1 \dots L$ do : $Z \leftarrow \sigma_{\ell} (X W_{in}^{\ell}) W_{2}^{\ell}$: $X \leftarrow \operatorname{Norm}_{\ell} (X + Z)$: end for	apply FFN or attention normalize and apply residual connection
7: return $XW_{head} + b_{head}$	apply model head