SLIM-LLM: SALIENCE-DRIVEN MIXED-PRECISION QUANTIZATION FOR LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models (LLMs) have achieved remarkable progress, but their extensive number of parameters results in high memory usage, significant loading latency, and substantial computational demands. To address these challenges, posttraining quantization (PTQ) has emerged as an effective technique for compressing model weights. In the context of PTQ for LLMs, existing uniform quantization methods, though efficient in terms of memory and computational requirements, often struggle to maintain performance. In this paper, we propose **SliM-LLM**, a Salience-Driven Mixed-Precision Quantization scheme that achieves group-wise bit-width allocation with mixed precisions for efficient LLMs with high accuracy. Building on our observation that salient/important weights often follow a structured distribution, we incorporate two core components to preserve post-quantization performance in LLMs while maintaining efficiency: 1) Salience-Determined Bit Allocation adaptively assigns bit widths to groups within each layer based on their group-level salience, aiming to minimize the reconstruction error of activations; and 2) Salience-Weighted Quantizer Calibration optimizes quantizer parameters by incorporating element-level salience, ensuring that the most critical weights are preserved, further preserving important weights information. With its structured group partitioning, SliM-LLM offers a hardware-friendly quantization approach, maintaining computational and memory efficiency comparable to highly optimized uniform quantization methods. Extensive experiments demonstrate that SliM-LLM significantly improves the accuracy of various LLMs when quantized to ultra-low bit widths. For instance, a 2-bit quantized LLaMA-7B model achieves nearly 6x memory reduction compared to its floating-point counterpart, alongside a 48% reduction in perplexity compared to the leading gradient-free PTQ method, all while maintaining GPU inference speed. Furthermore, SliM-LLM+, which incorporates gradient-based quantizers, reduces perplexity by an additional 35.1%.

034 035

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

1 INTRODUCTION

038 039 Large language models (LLMs) have exhibited exceptional performance across a wide array of natural 040 language benchmarks (Brown et al., 2020; Hendrycks et al., 2020). Notably, LLaMA (Touvron et al., 041 2023a) and GPT (Brown et al., 2020) series have significantly contributed to the ongoing evolution of 042 LLMs towards universal language intelligence. The powerful language understanding capabilities of 043 LLMs have been transferred to multi-modal domains (Li et al., 2024b; Achiam et al., 2023; Team 044 et al., 2023; Zhang et al., 2023; Huang et al., 2024b), laying the foundation for artificial general intelligence (AGI) (Bubeck et al., 2023). Despite these significant achievements, the substantial computational and memory requirements of LLMs pose huge challenges for real-world applications, 046 particularly in resource-constrained environments. 047

To address resource constraints of LLMs, post-training quantization (PTQ) has emerged as an efficient yet effective compression technique (Dettmers et al., 2022), showing success in quantizing the weights of pre-trained LLMs (Frantar et al., 2022; Lin et al., 2023; Shao et al., 2023; Lee et al., 2023; Chee et al., 2024). As LLMs continue to scale, the demand for more aggressive low-bit compression becomes critical due to limited computational and storage resources in application (Huang et al., 2024a; Tseng et al., 2024). However, significant performance degradation remains a challenge in low-bit scenarios (≤ 3-bit). To mitigate this, unstructured mixed-precision quantization (Shang et al.,



Figure 1: (a) The perplexity (↓) of existing low-bit PTQ methods of LLaMA at 2-bit. Solid-line indicates methods with structured quantization group. (b) Compare PTQ methods with gradient quantizer at 3-bit. (c) Features of current low-bit quantization methods. C denotes codebook-based, S is statistic-based, and G represents gradient-based quantizers.

2023; Huang et al., 2024a; Dettmers et al., 2023) and vector quantization (Chee et al., 2024; Tseng et al., 2024; Egiazarian et al., 2024) methods have been developed to preserve performance. While these approaches have advanced the field, they are often hardware-unfriendly, introducing extra storage requirements such as storing bitmaps or code indices, along with additional computations for vector decoding. This creates a bottleneck, limiting further reductions in memory and computational demands during deployments. In sum, ensuring the accuracy of LLMs while maintaining efficiency during deployment remains a significant challenge for current PTQ approaches.

075 This paper presents a Salience-Driven Mixed-Group LLM (SliM-LLM) framework, an accurate 076 and inference-efficient PTQ method for LLMs (\leq 3-bit). Our approach is grounded in the key 077 observation that salient or important weights, which are critical to model performance, exhibit a structured distribution, often clustering within certain channels (see Sec. 3.2.1 and Fig. 3). This insight, largely overlooked by prior research (Frantar et al., 2022), forms the basis for designing SliM-079 LLM as a structured, hardware-friendly mixed-precision low-bit quantization method. It preserves 080 performance through two key designs that retain important weights at both the global group and local 081 element levels. First, we develop a novel Salience-Determined Bit Allocation (SBA) method, which adaptively assigns bit-widths to each quantization group based on their group-level salience ranking. 083 The allocation strategy is optimized to reduce activation reconstruction errors. By applying higher 084 precision to more important groups and reducing the bit-width for less critical ones, SBA achieves 085 a low average bit-width while enhancing the overall performance of LLMs. Next, we introduce the Salience-Weighted Quantizer Calibration (SOC), which enhances sensitivity to locally salient 087 weights, ensuring that critical information within groups is preserved. SQC works collaboratively 880 with SBA, exploiting the local and global salience of weights to preserve the performance of LLMs after quantization. Unlike element-wise mixed-precision methods (Shang et al., 2023; Dettmers 089 et al., 2023; Huang et al., 2024a), SliM-LLM is inherently structured, eliminating additional bit or 090 computational overhead while preserving high performance. This is further demonstrated through 091 our deployment of SliM-LLM in an application-level inference tool¹ for LLMs, enabling efficient 092 mixed-precision inference on GPUs with consistently strong performance. 093

094 Experiments show that for various LLM families, SliM-LLM surpasses existing training-free PTQ methods on diverse benchmarks, particularly in low-bit scenarios. Using GPTQ as the backbone, SliM-LLM improves the perplexity scores of 2-bit LLaMA-13B and LLaMA2-13B on WikiText2 (Merity 096 et al., 2016) from 20.44 and 28.14 to 8.87 and 9.41, denoting performance improvements of over 56%, respectively. SliM-LLM even outperforms other element-wise mixed-precision PTO methods. 098 such as PB-LLM (Shang et al., 2023), APTQ (Guan et al., 2024) and LLM-MQ (Li et al., 2024a), in a deployment-friendly manner, showcasing its superior low-bit accuracy and efficiency. We also 100 integrate SliM-LLM into OmniQuant (Shao et al., 2023) and obtain SliM-LLM⁺ through gradient 101 optimization to further improve quantization quality. Moreover, the group-wise mixed-precision 102 strategy can smoothly be adapted to existing quantization-aware training (QAT) (Liu et al., 2023), fine-103 tuning based (Guo et al., 2023; Liao & Monz, 2024; Dettmers et al., 2024), or codebook-based (Chee 104 et al., 2024; Egiazarian et al., 2024; Tseng et al., 2024) LLMs compression methodologies. The 105 structure of weight salience we theoretically identify introduces a new practical view of the weight 106 quantization of LLMs.

https://github.com/AutoGPTQ/AutoGPTQ



Figure 2: Illustration of our proposed SliM-LLM. The Salience-Determined Bit Allocation (SBA) optimizes activation-aware structured precision, optimizing the global information distribution in quantization. Salience-Weighted Quantizer Calibration (SQC) detects discretely distributed salient weights, enhancing the local important information in LLMs.

120 **RELATED WORK** 2

121 122

108

110 111

112

113 114

115

116

117

118

119

Large Language Models (LLMs) have been significantly developed in diverse natural language 123 processing domains, establishing a prominent paradigm in these fields (Bubeck et al., 2023; Chang 124 et al., 2024; Zhao et al., 2023; Brown et al., 2020; Touvron et al., 2023a). Nevertheless, the exceptional 125 success of LLMs depends on massive parameters and computations, posing significant challenges for deployment in resource-constrained environments. Consequently, research into the compression 126 of LLMs has emerged as a promising field. Existing compression techniques for LLMs primarily 127 include low-bit quantization, pruning, distillation, and low-rank decomposition (Xu et al., 2023; 128 Ganesh et al., 2021; Frantar et al., 2022; Xiao et al., 2023a; Shao et al., 2023; Chee et al., 2024; 129 Zhu et al., 2023; Frantar & Alistarh, 2023; Huang et al., 2024a; Qin et al., 2024). Among these 130 technologies, low-bit quantization gains remarkable attention, for efficiently reducing the model size 131 without change of network structure(Zhu et al., 2023; Zhao et al., 2023; Chang et al., 2024). 132

Quantization of LLMs can be generally divided into QAT and PTQ. QAT, by employing a retraining 133 strategy based on quantized perception, better preserves the performance of quantized models. LLM-134 QAT (Liu et al., 2023; Ma et al., 2024a) addresses the data obstacle issue in QAT through data-free 135 distillation. However, for LLMs with huge size of parameters, the cost of retraining is extremely 136 inefficient(Chang et al., 2024). Therefore, PTQ has become a more efficient choice for LLMs. For 137 instance, LLM.int8() (Liu et al., 2023) and ZeroQuant (Yao et al., 2022) explore the quantization 138 strategies for LLMs in block-wise, which is a low-cost grouping approach that reduces hardware 139 burden. Subsequently, AWQ (Lin et al., 2023) and OWQ (Lee et al., 2023) also propose scaling 140 transformations on outlier channels of weight to preserve their information representation capacity. 141 GPTQ (Frantar et al., 2022) reduces the group quantization error of LLMs through Hessian-based 142 error compensation (Frantar & Alistarh, 2022), achieving commendable quantization performance at 3-bit. OmniQuant (Shao et al., 2023) introduces a learnable scaling quantizer to reduce quantization 143 errors in an output-aware manner. To achieve LLM quantization at ultra-low bit-width, recent 144 novel efforts such as QuIP (Chee et al., 2024), QuIP# (Tseng et al., 2024), and AQLM (Egiazarian 145 et al., 2024) promote quantization performance at 2-bit through matrix decomposition with learnable 146 codebooks and fine-tuning. Recent studies (Qin et al., 2024; Liao & Monz, 2024; Dettmers et al., 147 2024; Guo et al., 2023) have further refined compression techniques by integrating post-training 148 quantization (PTQ) with parameter-efficient fine-tuning (PEFT) to enhance model performance via 149 additional parameter learning. 150

Mixed-Precision Quantization exploits variations in the importance and redundancy of model 151 parameters, assigning different bit-widths to each component. HAWQ V2 (Dong et al., 2020) and 152 V3 (Yao et al., 2021) optimize bit-width allocation layer-wise in traditional visual networks through 153 Hessian analysis and Integer Linear Programming (ILP). Alternatively, OMPQ (Ma et al., 2023) 154 employs network orthogonality instead of Hessian for similar purposes. In LLMs, APTQ (Guan 155 et al., 2024) extends HAWQ's strategy, allocating varied bit-widths to different transformer blocks 156 based on Hessian-trace, thus improving the accuracy of 3-bit LLMs. However, such block-wise or 157 layer-wise mixed-precision allocation at 2-bit still fails to maintain post-compression performance 158 in LLMs. Recent studies such as SpQR (Dettmers et al., 2023), PB-LLM (Shang et al., 2023), and 159 LLM-MQ (Li et al., 2023) have introduced finer-grained partitioning for grouped quantization with element-wise mixed-precision for accurate weight quantization. Nevertheless, these low-bit methods 160 still rely on special structures and fine-grained grouping to ensure accuracy, which brings the huge 161 burden of real hardware deployment and inference speed.

¹⁶² 3 SLIM-LLM

This section introduces a group-wise mixed-precision quantization method, SliM-LLM, designed to
overcome the accuracy and efficiency bottlenecks of mixed-precision LLMs. We devise two novel
strategies for LLMs, including the use of *Salience-Determined Bit Allocation* (SBA) based on global
salience distribution to determine group bit-widths, and *Salience-Weighted Quantizer Calibration*(SQC) to enhance the perception of locally important weight information. We introduce SBA and
SQC in Sec. 3.2 and Sec. 3.3, respectively.

170 171 3.1 PRELIMINARIES

177 178

183

188 189

172 173 174 175 176 Quantization Framework. We first present the general uniform quantization process of LLMs 174 according to common practice (Liu et al., 2023; Shao et al., 2023; Achiam et al., 2023). The 174 quantization process requires mapping float-point weights distributed within the interval $[w_{\min}, w_{\max}]$ 175 to an integer range of 2^N , where N is the target bit-width. The quantization function for weight 176 $w_f \in \mathbb{R}^{n \times m}$ follows:

$$\hat{\boldsymbol{w}}_q = \operatorname{clamp}(\lfloor \frac{\boldsymbol{w}_f}{s} \rceil + z, 0, 2^N - 1), \ s = \frac{w_{\max} - w_{\min}}{2^N - 1}, \ z = -\lfloor \frac{w_{\min}}{s} \rceil$$
(1)

where \hat{w}_q indicates quantized weight which is integer, $\lfloor \cdot \rceil$ is round operation and $\operatorname{clamp}(\cdot)$ constrains the value within integer range (e.g. [0, 1, 2, 3], N = 2). Δ is scale factor and z is quantization zero point, respectively. When converted to 1-bit quantization, the calculation follows:

$$\hat{\boldsymbol{w}}_b = \operatorname{sign}(\boldsymbol{w}_f), \ \operatorname{sign}(\boldsymbol{w}) = \begin{cases} 1 & \text{if } \boldsymbol{w} \ge 0, \\ -1 & \text{others.} \end{cases}, \ \alpha = \frac{1}{l} ||\boldsymbol{w}_f||_{\ell 1}$$
(2)

185 where \hat{w}_b is binary result. α denots binarization scales and l is the number of elements in weight (Qin 186 et al., 2023), used for dequantization through $\alpha \hat{w}_b$. We can formalize the per-layer loss in PTQ, 187 following the common practice (Nagel et al., 2020; Frantar et al., 2022):

$$\mathcal{L}(\hat{\boldsymbol{w}}_f) = ||\boldsymbol{x}\boldsymbol{w}_f^{\top} - \boldsymbol{x}\hat{\boldsymbol{w}}_f^{\top}||^2 \approx \operatorname{tr}((\hat{\boldsymbol{w}}_f - \boldsymbol{w})\boldsymbol{H}(\hat{\boldsymbol{w}}_f - \boldsymbol{w})^{\top})$$
(3)

where $\boldsymbol{x} \in \mathbb{R}^{t \times m}$ denotes the input vectors from calibration dataset, $\hat{\boldsymbol{w}}_f \in \mathbb{R}^{n \times m}$ is dequantized weight from quantization result in Eq. (1) or Eq. (2), and $\boldsymbol{H} = \frac{1}{P} \sum_{k=1}^{P} \boldsymbol{x}^{[k]^{\top}} \boldsymbol{x}^{[k]}$ is proxy Hessian matrix by Levenberg-Marquardt approximation (Marquardt, 1963; Frantar & Alistarh, 2022) from a set of input activations.

Parameter Salience. In LLMs, the importance of each element in the weight matrix is various (Dettmers et al., 2023; Frantar & Alistarh, 2023). According to Eq. (3), quantizing different elements causes different impacts on the model's output loss. Elements that significantly influence the loss are termed salient weights. Consequently, we follow the SparseGPT (Frantar & Alistarh, 2023) to define the salience of each element as:

Definition 1. In the quadratic approximation of the loss as expressed in Eq. (3), we give the Hessian matrix $H \in \mathbb{R}^{m \times m}$ generated by $\frac{1}{P} \sum_{k=1}^{P} \boldsymbol{x}^{[k]^{\top}} \boldsymbol{x}^{[k]}$ for a weight matrix, the removal of the element at (i, j) induces an error $\delta_{i,j} = \frac{w_{i,j}^2}{[H^{-1}]_{j,j}^2}$ to the output matrix for linear projection in LLMs.

203 where $[H^{-1}]_{ij}$ denotes the j^{th} diagonal entry for the inverse Hessian, and H^{-1} can be efficiently 204 calculated through Cholesky decomposition (Krishnamoorthy & Menon, 2013). According to 205 Definition. 1, we map the elimination error δ_{ij} to the salience measure of each weight element in 206 LLMs, representing the impact of different weights on the output loss and the language capabilities, 207 which also leads the generation of mixed-precision quantization strategies (Dettmers et al., 2023; 208 Shang et al., 2023; Huang et al., 2024a; Li et al., 2024a) for LLMs. However, existing mixed-precision 209 solutions require the discrete allocation of bit-widths across the entire weight matrix, which imposes 210 a significant burden on hardware computations, thereby affecting the inference efficiency.

211 212

213

3.2 SALIENCE-DETERMINED BIT ALLOCATION

214 We reveal the phenomenon of spatial clustering in the distribution of weight salience, which inspires 215 our proposed concept of group-wise mixed-precision quantization for LLMs, and then introduce the *Salience-Determined Bit Allocation* (SBA) technique to allocate the optimal precision to each group.

216 3.2.1 SPATIAL DISTRIBUTION OF GLOBAL SALIENCE

218 We first conduct an empirical investigation into the weight salience distribution. The results reveal that 219 certain channels exhibit higher salience and show ten-220 dencies for spatial clustering. As illustrated in Fig. 3, 221 salient clustering are identified around the 2100^{th} , 222 3218^{th} and 3853^{rd} channels within the 2^{nd} layer's at-223 tention projection of the LLaMA-7B model. A similar 224 structured pattern is observed near the 600^{th} , 2200^{th} 225 and 3992^{nd} channels in the 10^{th} layer. Also, clustered 226 salience is detected in other layers (as shown in Fig. 3). 227 More examples of spatial clustering of salience are 228 provided in Appendix G.



256

266 267 268



(c) Salience of layer-10-Out (d) Salience of layer-10-Down Figure 3: Salience weight distribution in layer-2 and layer-10 of LLaMA-7B.

LLMs, activations exhibit extreme outlier channels, while the numerical differences in weights are
 relatively slight (Xiao et al., 2023a; Nrusimha et al., 2024). Therefore, we propose an analysis of how
 the outlier channels in activations influence the distribution of weight salience:

Theorem 1. Given the input calibration activation $x \in \mathbb{R}^{t \times m}$ with an outlier channel $x_{:,p}^* \gg x_{:,j}, \forall j \in [0,m], j \neq p$ at the position of channel-p. The trace elements of $H = x^{\top}x$ will show great outlier value at (p,p), where $H_{p,p} \gg H_{j,j}, \forall j \in [0,m], j \neq p$, as $H_{p,p}$ is produced by $[x_{:,p}^{*\top}x_{:,p}^*] = \sum_{i=0}^{t} x_{i,p}^{*2}$, which further leads to the parameter salience larger at the p^{th} channel of weight, where $\delta_{:,p} > \delta_{:,k}, \delta_{:,k} = \frac{w_{i,k}^2}{|H^{-1}|_{k,p}^2}, \forall k \in [0,t], k \neq p$.

Theorem 1 elucidates the influence of outlier activation on the distribution of channel salient weights
(detailed proof in the Appendix G.1). Furthermore, recent research indicates that outlier channels
in LLMs activations consistently appear in fixed yet clustered patterns (Nrusimha et al., 2024).
According to Theorem 1, these consistently occurring anomalous activations result in the distribution
of salient weights, as depicted in Fig. 3. Then, during group-wise quantization, the average salience
of each group shows different features.

Meanwhile, previous unstructured mixed-precision, incurred additional storage requirements and computational overheads, affecting the real-time inference. However, the strong spatial structured characteristics observed in the salient of weights in this section strongly inspire us to first develop a group-wise mixed-precision strategy within the weight matrix while maintaining inference efficiency. Therefore, we aim to allocate bit-widths based on intra-group salient disparities, which not only enhances quantization accuracy but also ensures the inference efficiency of LLMs with structured bit-widths saving and dequantization.

257 3.2.2 SALIENCE-DETERMINED BIT ALLOCATION FOR STRUCTURED GROUP

To allocate optimal bit-widths to each group, we introduce a *Salience-Determined Bit Allocation* (SBA) technique for mixed-precision LLMs, as depicted in Fig. 2. This technique, predicated on the differences in group salience, determines the optimal bit-width allocation for different groups by minimizing the distance of information entropy with the original weight output.

Specifically, we first utilize the average salience as the importance indicator for each weight group
 and rank them accordingly. The proposed SBA optimizes the following formula to determine the
 optimal number of salient-unsalient quantization groups of LLMs:

Objective : argmin
$$\mathcal{D}_{kl} (\boldsymbol{x} \boldsymbol{w}_{f}^{\top} || \boldsymbol{x} (\hat{\boldsymbol{w}}_{sba})^{\top}), \ \hat{\boldsymbol{w}}_{sba} = [\hat{\boldsymbol{w}}_{0,b_{0}}, \hat{\boldsymbol{w}}_{1,b_{1}}...\hat{\boldsymbol{w}}_{k-1,b_{k-1}}, \hat{\boldsymbol{w}}_{k,b_{k}}]$$

Constrain : $|\mathcal{G}_{N-1}| = |\mathcal{G}_{N+1}|, \ \mathcal{G}_{N-1} = \{b_{i}|b_{i} = N-1\}, \ \mathcal{G}_{N+1} = \{b_{j}|b_{j} = N+1\},$ (4)

where $\mathcal{D}_{kl}(\cdot||\cdot)$ denotes the Kullback-Leibler (KL) divergence between two outputs, \hat{w}_{f}^{sba} generally represents the de-quantization results of weight, employing group-wise mixed-precision designated

as $[\hat{w}_{0,b_0}, \hat{w}_{1,b_1}...\hat{w}_{k-1,b_{k-1}}, \hat{w}_{k,b_k}]$, where b_i represents the bit-width for the i^{th} group and \mathcal{G} is a set of groups with the same bit-width, N is the targeted average bit-width. We apply a compensation constraints strategy to maintain a consistent average bit-width for our SBA. For example, in 2-bit quantization, the groups with the highest salience are quantized to 3-bit. To offset the additional bits, we quantize an equal number of groups with the lowest salience to 1-bit ($|\mathcal{G}_{N-1}| = |\mathcal{G}_{N+1}|$), while the remaining groups are set to 2-bit.

276 We utilize an effective double-pointer search (more detailed examples in Appendix C) to optimize 277 our objective in Eq. (4). When the weight output channel size is m and group size is 128, $k = \frac{m}{198}$, 278 the search region for weight is limited to $[0, \frac{k}{2}]$, which is highly efficient with limited searching 279 space, e.g., only 16 iterations are needed in LLaMA-7B. We also provide detailed searching error 280 examples in Appendix C. Notably, SBA diverges from traditional quantization with mean squared 281 error (MSE) in Eq. (3) by instead utilizing the KL divergence as its measure of loss. Compared to 282 using the mean squared error (MSE) for weights, SBA leverages the KL divergence of block outputs 283 as a precision allocating metric, aiming to maximize the similarity between the distribution of the LLM's output activation matrix and the quantized activation distribution. This approach enhances the 284 model's information representation capacity under low-bit quantization, facilitating optimal bit-width 285 allocation. We note that HAWQ v2 (Dong et al., 2019) employs ILP to allocate bit-width for layers, which can also be adapted to our group-wise target. However, unlike the allocation of precision 287 based solely on the weight matrix loss of each group, SBA can accurately perceive the impact of 288 different precisions within each block on the model's output information, allowing for a more optimal 289 bit-width allocation. More experiments comparing SBA and ILP are shown in Section 4.2. 290

- 291
- 292 293

295

296

297

298

299

3.3 SALIENCE-WEIGHTED QUANTIZER CALIBRATION

In addition to the global group-wise distribution of salience, we notice that salience within the group still shows local differences in discrete distribution. Common existing quantizers apply uniform consideration across all weights to minimize the effect (error) of quantization, lacking the capability to perceive differences in local salience. Therefore, in this section, we introduce a *Salience-Weighted Quantizer Calibration* (SQC) to enhance the information of significant weights within the group by amplifying the quantizer awareness of salient weight.

300 301 302

303 304

3.3.1 DISCRETE DISTRIBUTION OF LOCAL SALIENCE

In the aforementioned section, we group-wisely allo-305 cate the bit-width for each group based on the global 306 salience. To maintain the efficiency of quantized in-307 ference, we employ a commonly used sequential struc-308 tured grouping (Frantar et al., 2022; Lin et al., 2023; 309 Shao et al., 2023). However, this group-wise mixed-310 precision also leads to differences in salience among the 311 various elements within the same group. Specifically, 312 as the salience distribution in Fig. 4, within the 10^{th} 313 attention output layer of LLaMA-7b, a subset of sparse 314 weights within the comparatively less salient Group-2 315 (Fig. 4) still maintains a high level of importance. In LLMs, a small number of weight elements with outliers 316 affect the local distribution of salience. These discrete 317 weights typically account for only approximately 1% 318



Figure 4: Local salience distribution of the 10^{th} MHA output layer in LLaMA-7B.

within the group but play a crucial role in the modeling capability of LLMs.

The existing vanilla quantizers face the challenge of representing significant weight information, by
 only considering the mean error of all elements within a group. When quantizing weights according
 to Eq. (1) in group-wise format, a large number of non-salient weights at the intra-group statistical
 level tend to dominate the parameters generated by the quantizer. This leads to a degradation of
 salient information within the group, thereby affecting the model performance of LLMs.

324 3.3.2 SALIENCE-WEIGHTED QUANTIZER CALIBRATION FOR LOCAL SALIENCE AWARENESS

To prevent the degradation of local salient weight information in each group, we propose the *Salience-Weighted Quantizer Calibration* (SQC), which enhances the expression of salient weights through local salience awareness, thereby reducing the quantization error of these significant elements and improving the compressed performance of LLMs.

330 Based on a common observation (Dettmers et al., 2023; Huang et al., 2024a), the proportion of 331 relatively salient weights in each group is only 1-5%. Therefore, we employ the 3- σ rule for a mask 332 to select the salience part ($w < (\mu - 3\sigma) \cup w > (\mu + 3\sigma)$) in each group (Fig. 2), which accounts 333 for about 1% elements. After the selection, we get $w_i = w_i^s \cup w_i^{us}$, where w_i^s is the salient part and w_i^{us} represents the non-salient elements within group *i*. To effectively keep the information of local 334 salient weights, SQC first introduces the calibration parameter τ to the SQC quantizer, liberating the 335 perception interval during quantization. Then we define the local salience awareness loss of the SQC 336 quantizer through calibration: 337

338 339

340

 $\underset{\tau}{\operatorname{argmin}} ||\boldsymbol{w}_{i}^{s} - \tau \cdot s\{\mathcal{Q}(\boldsymbol{w}_{i}^{s}, \tau \cdot s, \tau \cdot z) - \tau \cdot z\}||_{2}^{2} + ||\boldsymbol{w}_{i}^{us} - \tau \cdot s(\mathcal{Q}(\boldsymbol{w}_{i}^{us}, \tau \cdot s, \tau \cdot z) - \tau \cdot z)||_{2}^{2} (5)$

where $Q(\cdot)$ denotes the quantization process in Eq. (1), $|| \cdot ||_2^2$ represents the ℓ_2 loss, aligned with Eq. (3). w_i^s and w_i^{us} denotes the salient and less salient part of group *i*, respectively, generated from a mask operation. In Eq. (5), τ expands the solution space of *s* and *z*, flexibly adjusts *s* and *z* to search the optimal loss under τ^* , without bringing additional parameters, as w_i^s and w_i^{us} share the same quantizer. The search space for τ by linearly dividing the interval $[1-\lambda, 1+\lambda]$ into 2n candidates. We empirically set λ at 0.1 and *n* at 50 to achieve a balance between efficiency and accuracy.

347 Compared to traditional quantizer calibration methods, SQC effectively mitigates the degradation of 348 intra-group local salient weights caused by general average loss by enhancing the loss sensitivity to salient elements during the calibration (more experiments are detailed in Appendix E). Moreover, the 349 SQC process allows w_i^s and w_i^{us} to share a set of parameters τ^*s and τ^*z , eliminating the need to 350 differentiate intra-group weights during storage and inference. This facilitates straightforward group-351 wise dequantization calculations, thereby avoiding the hardware overhead associated with element-352 wise bitmap and unstructured grouping. SQC and SBA each capture local salient weight information 353 within groups and global salient weight combinations across groups, effectively enhancing the 354 protection of critical information during quantization, thereby accurately preserving the overall 355 performance of LLMs at extremely low bit-widths. 356

357 358

359

3.4 IMPLEMENTATION PIPELINE OF SLIM-LLM

We integrate our mixed-precision framework into advanced PTQ methods, such as GPTQ (Frantar et al., 2022) and OmniQuant (Shao et al., 2023), all of which are inference-friendly with group-wise quantization. We primarily integrate SBA and SQC into GPTQ to get SliM-LLM. For SliM-LLM⁺, the SBA is plugged into OmniQuant with a learnable quantizer. The plugging pipeline of SliM-LLM is provided in Algorithm 1 (line 4 and line 9), detailed functions are shown in Appendix B.1.

Algorithm 1 Main Framework of SliM-LLM.	
func SliM-LLM($\boldsymbol{w}, \boldsymbol{x}_F, \beta, \lambda, N$)	4: $\mathcal{G}{\cdot} \coloneqq \mathrm{SBA}(\boldsymbol{w}, \boldsymbol{x}_F, \boldsymbol{H}^{\mathrm{in}}, \beta, N)$
Input: $\boldsymbol{w} \in \mathbb{R}^{n imes m}$ - FP16 weight	5: for $b = 0, \beta, 2\beta,$ do
$oldsymbol{x}_F \in \mathbb{R}^{t imes m}$ - calibration data	6: $oldsymbol{w}^b\coloneqqoldsymbol{w}_{:,b:b+eta}$
β - group size	7: $g_b \coloneqq \mathcal{G}[b]$
λ - hessian regularizer	8: $\boldsymbol{w}_{s}^{b}, \boldsymbol{w}_{us}^{b} \coloneqq \operatorname{sal}_{\max}(\boldsymbol{w}^{b})$
N - average bit-width	9: $\hat{\boldsymbol{w}}_{a}^{b} \coloneqq \operatorname{SQC}(\boldsymbol{w}_{s}^{b}, \boldsymbol{w}_{us}^{b}, g_{b})$
Output: \hat{w}_q - quantized weight	10: GPTQ-error compensation:
	11: $\boldsymbol{E} \coloneqq (\boldsymbol{w}_{:,b:b+\beta} - \hat{\boldsymbol{w}}_a^b) / \boldsymbol{H}_{bb:b+\beta b+\beta}^{\text{in}}$
1: $\boldsymbol{H} \coloneqq \frac{1}{P} \sum_{k=1}^{P} \boldsymbol{x}_{F}^{[k]} \boldsymbol{x}_{F}^{[k]T}$ hessian matrix	12: $\boldsymbol{w}_{:,b+\beta:} \coloneqq \boldsymbol{w}_{:,b+\beta:} - \boldsymbol{E} \cdot \boldsymbol{H}_{b:b+\beta,b+\beta}^{\text{in}}$
2: $\boldsymbol{H}^{\text{in}} \coloneqq \text{Cholesky}((\boldsymbol{H} + \lambda \boldsymbol{I})^{-1})$	13: end for
3: $\hat{\boldsymbol{w}}_a \coloneqq 0^{n \times m}$	14: return \hat{w}_{q}
Ч	1

not give the	icsuits on iis	icu mou	cis of the	coues						
#W PPL↓	Method	1-7B	1-13B	1-30B	1-65B	2-7B	2-13B	2-70B	3-8B	3-70B
16-bit	-	5.68	5.09	4.10	3.53	5.47	4.88	3.31	5.75	2.9
	APTQ	6.76	-	-	-	-	-	-	-	-
	LLM-MQ	-	-	-	-	-	8.54	-	-	-
2 1.4	RTN	7.01	5.88	4.87	4.24	6.66	5.51	3.97	27.91	11.84
3-011	AWQ	6.46	5.51	4.63	3.99	6.24	5.32	-	8.22	4.81
	GPTQ	6.55	5.62	4.80	4.17	6.29	5.42	3.85	8.19	5.22
	SliM-LLM	6.40	5.48	4.61	3.99	6.24	5.26	3.67	7.16	4.08
	LLM-MQ	-	-	-	-	-	12.17	-	-	-
	RTN	1.9e3	781.20	68.04	15.08	4.2e3	122.08	27.27	1.9e3	4.6e5
	AWQ	2.6e5	2.8e5	2.4e5	7.4e4	2.2e5	1.2e5	-	1.7e6	1.7e6
2-bit	GPTQ	152.31	20.44	13.01	9.51	60.45	28.14	8.78	210.00	11.90
	QuIP	29.74	12.48	11.57	7.83	39.73	13.48	6.64	84.97	13.03
	PB-LLM	24.61	17.73	12.65	7.85	25.37	49.81	NAN	44.12	11.68
	SliM-LLM	14.58	8.87	7.33	5.90	16.01	9.41	6.28	39.66	9.46

Table 1: Quantization results of LLaMA family with statistic quantizer. We report the WikiText2 perplexity in this table, C4 results are shown in Appendix H. '-' denotes that the selected works did not give the results on listed models or the codes

Table 2: Quantization results of LLaMA-1 and LLaMA-2 models with learnable quantizer. We report the WikiText2 perplexity in this Table, C4 results are shown in Appendix H. '-' denotes that the selected works have not reported the results on listed models or published the codes

#W PPL↓	Method	1-7B	1-13B	1-30B	1-65B	2-7B	2-13B	2-70B
16-bit	-	5.68	5.09	4.10	3.53	5.47	4.88	3.31
3-bit	OmniQuant AffineQuant SliM-LLM ⁺	6.15 6.14 6.07	5.44 5.45 5.37	4.56 4.59 4.34	3.94 	6.03 6.08 5.94	5.28 5.28 5.11	3.78
2-bit	OmniQuant AffineQuant SliM-LLM ⁺	9.72 13.51 9.68	7.93 7.22 7.17	7.12 6.49 6.41	5.95 - 5.74	11.06 10.87 10.87	8.26 7.64 7.59	6.55 - 6.44

EXPERIMENTS

We evaluated SliM-LLM and SliM-LLM⁺ under weight-only conditions, focusing on 2/3-bit pre-cisions. Per-channel group quantization is utilized in our framework with groupsize = 128 in experiments. Since no back-propagation in SliM-LLM, the quantization is carried out on a sin-gle NVIDIA A800 GPU. For SliM-LLM⁺, we employ the AdamW optimizer, following Omni-Quant (Shao et al., 2023), which is also feasible on a single A800. We randomly select 128 samples from WikiText2 (Merity et al., 2016) as calibration data, each with 2048 tokens.

Models and Evaluation. To comprehensively demonstrate the low-bit performance advantages of SliM-LLM and SliM-LLM⁺, we conduct experiments across OPT (Zhang et al., 2022), LLaMA (Tou-vron et al., 2023a), LLaMA-2 (Touvron et al., 2023b) and LLaMA-3. We employ the perplexity as our evaluation metric, which is widely recognized as a stable measure of language generation capabilities (Frantar et al., 2022; Lin et al., 2023; Huang et al., 2024a; Shang et al., 2023; Shao et al., 2023; Chee et al., 2024; Egiazarian et al., 2024; Huang et al., 2024b), particularly in compression scenarios. Experiments are carried out on the WikiText2 (Merity et al., 2016) and C4 (Raffel et al., 2020) datasets. Furthermore, to assess the practical application capabilities of quantized LLMs, we also evaluate their accuracy on zero-shot benchmarks, including PIQA (Bisk et al., 2020), ARC (Clark et al., 2018), BoolQ (Clark et al., 2019), and HellaSwag (Clark et al., 2018).

Baseline. Since SliM-LLM and SliM-LLM⁺ are efficient PTQ approaches without additional training or fine-tuning, QAT and re-training methods are not within the comparison range of our work. The experiments evaluate existing advanced quantization methods and GPU-friendly computations, including vanilla round-to-nearest (RTN), GPTQ (Frantar et al., 2022), AWQ (Lin et al., 2023). And mixed-precision quantization techniques, including PB-LLM (Shang et al., 2023) ($\frac{1}{7} \times 8$ -bit+ $\frac{6}{7} \times 1$ -bit), LLM-MQ (Li et al., 2024a), and APTQ (Guan et al., 2024), as well as the codebook-based method QuIP (Chee et al., 2024) are also compared in this work. We compare SliM-LLM⁺ with

Model / Acc↑	#W	Method	PIQA	ARC-e	ARC-c	BoolQ	HellaSwag	Winogrande	Avg.
	16-bit	-	77.47	52.48	41.46	73.08	73.00	67.07	64.09
-	2-bit	GPTQ	55.49	31.02	$2\overline{2}.\overline{17}$	53.49	33.84	41.91	39.65
II MA 7B	2-bit	AWQ	47.78	28.77	21.31	31.19	24.47	40.03	32.26
LLawA-/D	2-bit	SliM-LLM	57.83	33.46	25.09	56.05	36.70	52.64	43.84
	2-bit	_ OmniQuant	63.63	43.91	27.32^{-}	$5\overline{8}.0\overline{2}$	48.78	52.97	49.11
	2-bit	SliM-LLM ⁺	64.96	45.66	28.67	64.59	48.86	53.35	51.02
	16-bit	-	79.10	59.89	44.45	68.01	76.21	70.31	66.33
	2-bit	<u>G</u> PTQ	70.37	47.74	35.88	51.57	61.39	60.84	54.63
II aMA_13B	2-bit	AWQ	49.23	30.01	29.49	30.88	26.72	46.30	35.44
LLawA-15D	2-bit	SliM-LLM	73.19	47.95	36.27	55.92	63.04	61.79	56.36
-	2-bit	_ OmniQuant_	73.14	49.38	-36.93	63.34	62.19	61.77	57.64
	2-bit	SliM-LLM ⁺	74.15	50.26	37.04	64.31	63.57	63.11	58.74
	16-bit	-	80.08	58.92	45.47	68.44	79.21	72.53	67.44
-	2-bit	GPTQ	71.92	48.27	36.20	61.27	65.76	63.11	57.76
LL 2MA-30B	2-bit	AWQ	49.17	28.56	25.97	34.73	24.97	46.99	35.07
LLawiA-JOD	2-bit	SliM-LLM	75.52	51.29	39.29	62.01	66.10	64.07	59.71
-	2-bit	_ OmniQuant_	$^{-}7\overline{6}.\overline{2}3^{-}$	53.23	39.52	63.34	65.57	64.82	60.22
	2-bit	SliM-LLM ⁺	76.31	54.07	39.79	63.35	67.14	64.93	60.91
	16-bit	-	80.79	58.71	46.24	82.29	80.72	77.50	71.04
LLaMA-30B	2-bit	GPTQ	76.16	52.48	40.14	77.23	71.96	70.22	64.70
II aMA_65P	2-bit	SliM-LLM	77.09	53.72	40.25	77.51	72.05	70.91	65.26
LLawiA-05D	2-bit	OmniQuant	77.78	53.71	$4\bar{0}.\bar{90}$	$7\bar{8}.0\bar{4}$	74.55	68.85	65.64
		CHARTER AT	=0 0C	F2 00	41 10	70.22	777	(0.00	((10

Table 3: Performance comparisons of different quantization methods for zero-shot tasks.

gradient optimizer-based OmniQuant (Shao et al., 2023) and AffineQuant (Ma et al., 2024b). When applying SliM-LLM, the quantization process for a 7B model takes only about 50 minutes.

4.1 MAIN RESULTS

432

456

457

458 459

460

476

461 We show experiments within the LLaMA family in this section and detailed results for the OPT 462 models are available in Appendix H. For language generation tasks, as depicted in Tab. 1, SliM-LLM 463 markedly outperforms its backbone GPTQ, particularly under the 2-bit. Specifically, on LLaMA-7B, SliM-LLM achieves a 90% decrease in perplexity, while on LLaMA-3-8B, it improves by 81%. In 464 comparison with the element-wise mixed-precision PB-LLM and the codebook-based QuIP method, 465 SliM-LLM further reduces the perplexity by 41%~51%. As shown in Tab. 1, the performance of SliM-466 LLM⁺ is still ahead compared to OmniQuant and AffineQuant, further proving the effectiveness and 467 of the mixed-precision framework. We also provide dialogue examples of 2-bit instruction fine-tuning 468 Vicuna-13B (Chiang et al., 2023) and LLaMA-13B in Appeandix I. 469

Morever, our method exhibits zero-shot advantages at 2-bit, as shown in Tab. 3, where SliM-LLM and SliM-LLM⁺ still outperforms other methods. For instance, compared with GPTQ and OmniQuant, our approach achieves an average improvement of 4.19% and 1.91% on LLaMA-7B. Meanwhile, for LLaMA-65B, 2-bit SliM-LLM and SliM-LLM⁺ is close to FP16 results (less than 6% degradaion in accuracy). Overall, our proposed mixed-precision framwork demonstrates superior performance across different model sizes, with its advantages becoming increasingly significant at lower bit-width.

477 4.2 ABLATION RESULTS

Abliation of SBA and SQC. We conduct a detailed ablation study to illustrate the benefits of bit-width allocation and the impact of each component. Fig. 5(a) compares three strategies for allocating bit-widths across groups, including random allocation, head-tail allocation by spatial order, and our proposed SBA. When the

Table 4: WikiText2↓	performance of SBA
and ILP on LLaMA.	-

Method	#W	7B	13B	30B	65B
ILP	2-bit	17.55	9.51	9.27	7.46
SBA	2-bit	14.58	8.87	7.33	5.90

average bit-width remains constant, random and head-tail mixed-precision allocation prove ineffective
 and even result in performance degradation, as shown in Fig. 5(a). In contrast, SBA consistently
 delivers significant improvements in post-quantization performance, validating the efficacy of our



Figure 5: Ablation results on OPT models. Random means randomly selecting the same number of lower/higher-bit groups; head-tail denotes using the head groups as the lower-bit and the same number of tails as the higher-bit on the original sequence of group.

Table 5: Deployment results of GPTQ and Slim-LLM on GPU. Group size is set to 128.

#W	LLaMA-*	1-7B					1-13B				2-7B			
		WM	RM	PPL↓	Token/s	WM	RM	PPL↓	Token/s	WM	RM	PPL↓	Token/s	
FP16	-	12.6G	14.4G	5.68	69.2	24.3G	27.1G	5.09	52.5	12.7G	14.6G	5.47	69.3	
3-bit	GPTQ SliM-LLM	3.2G 3.2G	5.1G 5.2G	6.55 6.40	83.4 79.1	5.8G 5.8G	8.7G 8.8G	5.62 5.48	57.6 48.5	3.2G 3.2G	5.2G 5.4G	6.29 6.26	56.3 55.9	
2-bit	GPTQ SliM-LLM	2.2G 2.3G	4.1G 4.4G	152.31 14.58	83.9 61.2	4.0G 4.1G	7.5G 7.8G	20.44 8.87	92.6 73.7	2.2G 2.3G	4.1G 4.1G	60.45 16.01	83.6 64.4	

507 mixed-precision approach. Fig. 5(b) presents the ablation effects of SBA and SQC, demonstrating that 508 both methods, based on the perception of global and local salience, enhance quantization performance. 509 SBA is particularly effective in smaller models, and combining these two methods can further boost capabilities of LLMs. We also provide the detailed ablation results on group size in Appendix F. 510

511 Compare of SBA and ILP. We compare the performance between the ILP model in HAWQ v2 (Dong 512 et al., 2019) and SBA on the LLaMA model. Tab. 4 shows that SBA achieves comprehensive 513 performance superiority on LLaMA. We observed that under a 2-bit scenario, ILP ensures an equal 514 number of 1-bit and 3-bit groups within the search space {1-bit, 2-bit, 3-bit}. The advantage of ILP 515 lies in a broader selection range for target bit-widths, but under commonly used fixed integer bitwidths (e.g. 2-bit, 3-bit), SBA's double-pointer search strategy based on output feature KL proposed 516 by SBA can achieve a more optimal matching strategy. 517

518 519

521

522

523

495

496 497

504 505 506

4.3 EFFICIENT INFERENCE ON DEVICE

520 We utilize the open-source AutoGPTQ to extend CUDA kernel supporting experimental mixedprecision inference, with detailed process in Appendix B.2. We evaluate the deployment performance of LLaMA-7/13B and LLaMA-2-7B under 2/3-bit settings in Tab. 5. The results indicate that our mixed-precision approach maintains a good compression rate on GPUs and significantly enhances 524 model accuracy, only with a slight decrease in inference speed on the A800 (due to the inference 525 alignment of different bit-width). Since current 1-bit operations lack well hardware support, additional 526 consumption of storage and computation is required on device. There remains considerable scope for optimization in mixed-precision computing, and we aim to further improve this in future work.

5 CONCLUSION

531 In this work, we introduce **SliM-LLM**, a group-wise mixed-precision PTQ framework tailored for 532 LLMs, designed to enhance performance with low-bit weights in a deployment-friendly manner. The essence of SliM-LLM lies in employing the Salience-Determined Bit Allocation to dynamically 534 allocate bit widths, thereby improving the preservation of global salience information. Within groups, the Salience-Weighted Quantizer Calibration is designed to enhance local information 536 perception, further minimizing the loss associated with locally salient weights. Experiments validate the effectiveness of SliM-LLM, showing notable accuracy improvements across various LLMs, and ensuring efficiency in inference. In conclusion, SliM-LLM is versatile and can be seamlessly 538 integrated with different quantization frameworks and successfully improves the performance of LLMs supporting practical deployment in resource-constrained environments.

540 REFERENCES 541

542 543 544	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> , 2023.
545 546 547	Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pp. 7432–7439, 2020.
548 549 550 551	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. <i>Advances in Neural Information Processing Systems</i> , 33:1877–1901, 2020.
552 553 554	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with GPT-4. <i>arXiv preprint arXiv:2303.12712</i> , 2023.
555 556 557 558	Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. <i>ACM Transactions on Intelligent Systems and Technology</i> , 15(3):1–45, 2024.
559 560 561	Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization of large language models with guarantees. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
562 563 564 565	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing GPT-4 with 90%* chatgpt quality. See https://vicuna. lmsys. org (accessed 14 April 2023), 2023.
566 567 568	Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. <i>arXiv preprint arXiv:1905.10044</i> , 2019.
569 570 571	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. <i>arXiv preprint arXiv:1803.05457</i> , 2018.
573 574	Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. LLM.int8 (): 8-bit matrix multiplication for transformers at scale. <i>arXiv preprint arXiv:2208.07339</i> , 2022.
575 576 577 578	Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashk- boos, Alexander Borzunov, Torsten Hoefler, and Dan Alistarh. SpQR: A sparse-quantized repre- sentation for near-lossless LLM weight compression. <i>arXiv preprint arXiv:2306.03078</i> , 2023.
579 580	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
581 582 583 584	Zhen Dong, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Hawq: Hessian aware quantization of neural networks with mixed-precision. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 293–302, 2019.
585 586 587	Zhen Dong, Zhewei Yao, Daiyaan Arfeen, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Hawq-v2: Hessian aware trace-weighted quantization of neural networks. <i>Advances in neural</i> <i>information processing systems</i> , 33:18518–18529, 2020.
588 589 590 591	Vage Egiazarian, Andrei Panferov, Denis Kuznedelev, Elias Frantar, Artem Babenko, and Dan Alistarh. Extreme compression of large language models via additive quantization. <i>arXiv preprint arXiv:2401.06118</i> , 2024.
592 593	Elias Frantar and Dan Alistarh. Optimal brain compression: A framework for accurate post-training quantization and pruning. <i>Advances in Neural Information Processing Systems</i> , 35:4475–4488, 2022.

594 Elias Frantar and Dan Alistarh. SparseGPT: Massive language models can be accurately pruned in 595 one-shot. In International Conference on Machine Learning, pp. 10323–10337. PMLR, 2023. 596 597 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. 598 Prakhar Ganesh, Yao Chen, Xin Lou, Mohammad Ali Khan, Yin Yang, Hassan Sajjad, Preslav Nakov, 600 Deming Chen, and Marianne Winslett. Compressing large-scale transformer-based models: A case 601 study on bert. Transactions of the Association for Computational Linguistics, 9:1061–1080, 2021. 602 603 Ziyi Guan, Hantao Huang, Yupeng Su, Hong Huang, Ngai Wong, and Hao Yu. APTQ: Attention-604 aware post-training mixed-precision quantization for large language models. arXiv preprint 605 arXiv:2402.14866, 2024. 606 Han Guo, Philip Greengard, Eric P Xing, and Yoon Kim. Lq-lora: Low-rank plus quantized matrix 607 decomposition for efficient language model finetuning. arXiv preprint arXiv:2311.12023, 2023. 608 609 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and 610 Jacob Steinhardt. Measuring massive multitask language understanding. arXiv preprint 611 arXiv:2009.03300, 2020. 612 613 Wei Huang, Yangdong Liu, Haotong Qin, Ying Li, Shiming Zhang, Xianglong Liu, Michele Magno, and Xiaojuan Qi. BiLLM: Pushing the limit of post-training quantization for llms. arXiv preprint 614 arXiv:2402.04291, 2024a. 615 616 Wei Huang, Xudong Ma, Haotong Oin, Xingyu Zheng, Chengtao Ly, Hong Chen, Jie Luo, Xiaojuan 617 Qi, Xianglong Liu, and Michele Magno. How Good Are Low-bit Quantized LLaMA3 Models? 618 An Empirical Study. arXiv preprint arXiv:2404.14047, 2024b. 619 620 Aravindh Krishnamoorthy and Deepak Menon. Matrix inversion using cholesky decomposition. In 2013 signal processing: Algorithms, architectures, arrangements, and applications (SPA), pp. 621 70-72. IEEE, 2013. 622 623 Changhun Lee, Jungyu Jin, Taesu Kim, Hyungjun Kim, and Eunhyeok Park. OWQ: Lessons 624 learned from activation outliers for weight quantization in large language models. arXiv preprint 625 arXiv:2306.02272, 2023. 626 627 Shiyao Li, Xuefei Ning, Ke Hong, Tengxuan Liu, Luning Wang, Xiuhong Li, Kai Zhong, Guohao 628 Dai, Huazhong Yang, and Yu Wang. Llm-mq: Mixed-precision quantization for efficient llm 629 deployment. In The Efficient Natural Language and Speech Processing Workshop with NeurIPS, volume 9, 2023. 630 631 Shiyao Li, Xuefei Ning, Ke Hong, Tengxuan Liu, Luning Wang, Xiuhong Li, Kai Zhong, Guohao 632 Dai, Huazhong Yang, and Yu Wang. LLM-MQ: Mixed-precision quantization for efficient LLM 633 deployment. In Advances in Neural Information Processing Systems (NeurIPS) ENLSP Workshop, 634 2024a. 635 636 Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng 637 Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. arXiv preprint arXiv:2403.18814, 2024b. 638 639 Baohao Liao and Christof Monz. Apiq: Finetuning of 2-bit quantized large language model. arXiv 640 preprint arXiv:2402.05147, 2024. 641 642 Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. AWO: 643 Activation-aware weight quantization for LLM compression and acceleration. arXiv preprint 644 arXiv:2306.00978, 2023. 645 Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang 646 Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. LLM-QAT: Data-Free Quantization Aware 647 Training for Large Language Models. arXiv preprint arXiv:2305.17888, 2023.

648 649 650	Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong, Ruiping Wang, Jilong Xue, and Furu Wei. The era of 1-bit llms: All large language models are in 1.58 bits. <i>arXiv preprint arXiv:2402.17764</i> , 2024a.
652 653 654	Yuexiao Ma, Taisong Jin, Xiawu Zheng, Yan Wang, Huixia Li, Yongjian Wu, Guannan Jiang, Wei Zhang, and Rongrong Ji. Ompq: Orthogonal mixed precision quantization. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 37, pp. 9029–9037, 2023.
655 656 657	Yuexiao Ma, Huixia Li, Xiawu Zheng, Feng Ling, Xuefeng Xiao, Rui Wang, Shilei Wen, Fei Chao, and Rongrong Ji. Affinequant: Affine transformation quantization for large language models. <i>arXiv preprint arXiv:2403.12544</i> , 2024b.
659 660	Donald W Marquardt. An algorithm for least-squares estimation of nonlinear parameters. <i>Journal of the society for Industrial and Applied Mathematics</i> , 11(2):431–441, 1963.
661 662	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. <i>arXiv preprint arXiv:1609.07843</i> , 2016.
664 665 666	Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization. In <i>International Conference on Machine Learning</i> , pp. 7197–7206. PMLR, 2020.
667 668 669	Aniruddha Nrusimha, Mayank Mishra, Naigang Wang, Dan Alistarh, Rameswar Panda, and Yoon Kim. Mitigating the impact of outlier channels for language model quantization with activation regularization. <i>arXiv preprint arXiv:2404.03605</i> , 2024.
670 671 672 673	Haotong Qin, Mingyuan Zhang, Yifu Ding, Aoyu Li, Zhongang Cai, Ziwei Liu, Fisher Yu, and Xianglong Liu. Bibench: Benchmarking and analyzing network binarization. <i>arXiv preprint arXiv:2301.11233</i> , 2023.
674 675 676	Haotong Qin, Xudong Ma, Xingyu Zheng, Xiaoyang Li, Yang Zhang, Shouda Liu, Jie Luo, Xianglong Liu, and Michele Magno. Accurate LoRA-Finetuning Quantization of LLMs via Information Retention. <i>arXiv preprint arXiv:2402.05445</i> , 2024.
677 678 679 680	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>The Journal of Machine Learning Research</i> , 21(1):5485–5551, 2020.
681 682	Yuzhang Shang, Zhihang Yuan, Qiang Wu, and Zhen Dong. PB-LLM: Partially binarized large language models. <i>arXiv preprint arXiv:2310.00034</i> , 2023.
683 684 685 686	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. <i>arXiv preprint arXiv:2308.13137</i> , 2023.
687 688	Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. <i>arXiv preprint arXiv:2306.11695</i> , 2023.
689 690 691 692	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> , 2023.
693 694 695	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. <i>arXiv preprint arXiv:2302.13971</i> , 2023a.
696 697 698 699	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. <i>arXiv preprint arXiv:2307.09288</i> , 2023b.
700 701	Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better LLM quantization with hadamard incoherence and lattice codebooks. <i>arXiv preprint</i> <i>arXiv:2402.04396</i> , 2024.

702 703 704	Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In <i>International</i> <i>Conference on Machine Learning</i> , pp. 38087–38099. PMLR, 2023a.
705 706 707	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. <i>arXiv preprint arXiv:2309.17453</i> , 2023b.
708 709 710	Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhensu Chen, Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models. arXiv preprint arXiv:2309.14717, 2023.
711 712 713 714	Z Yao, RY Aminabadi, M Zhang, X Wu, C Li, and Y Zeroquant He. Efficient and affordable post-training quantization for large-scale transformers. <i>URL https://arxiv. org/abs/2206.01861</i> , 2022.
715 716 717	Zhewei Yao, Zhen Dong, Zhangcheng Zheng, Amir Gholami, Jiali Yu, Eric Tan, Leyuan Wang, Qijing Huang, Yida Wang, Michael Mahoney, et al. Hawq-v3: Dyadic neural network quantization. In <i>International Conference on Machine Learning</i> , pp. 11875–11886. PMLR, 2021.
718 719 720 721	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022.
722 723 724	Yiyuan Zhang, Kaixiong Gong, Kaipeng Zhang, Hongsheng Li, Yu Qiao, Wanli Ouyang, and Xiangyu Yue. Meta-transformer: A unified framework for multimodal learning. <i>arXiv preprint arXiv:2307.10802</i> , 2023.
725 726 727 728	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. <i>arXiv</i> preprint arXiv:2303.18223, 2023.
729 730 731	Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. <i>arXiv preprint arXiv:2308.07633</i> , 2023.
732 733	
734 735 736	
737 738 739	
740 741	
742 743 744	
745 746	
747 748 749	
750 751 752	
753 754	
755	

756 A LIMITATIONS

Though the mixed-precision framework significantly improves the quantization performance of
LLMs, the current out-of-the-box deployment tools still cannot well support efficient mixed-precision
computing. Meanwhile, the support for 1/2/3-bit inference on GPUs remains limited, which affects
the inferencing advantages of low-bit models. We believe there is significant room for improvement
in the hardware efficiency of mixed-precision LLMs in the future.

763 764

765 766

767

B SLIM-LLM IMPLEMENTATION

B.1 DETAILED IMPLEMENTATION

In this section, we present the specific implementation details of SliM-LLM, which utilizes
GPTQ (Frantar et al., 2022) as its backbone for mixed-precision quantization and incorporates
both SBA and SQC. SliM-LLM⁺ is consistent with SliM-LLM in SBA computations but does
not include the SQC component, instead retaining learnable weight clipping (LWC) approach in
OmniQuant (Shao et al., 2023) for gradient optimization.

773 Algorithm 2 Detailed functions in SliM-LLM. 774 func SQC($\boldsymbol{w}_{s}^{b}, \boldsymbol{w}_{us}^{\overline{b}}, g_{b}$) 775 func SBA $(\boldsymbol{w}, \boldsymbol{x}_F, \boldsymbol{H}^{\text{in}}, \beta, N)$ 776 1: $w_{\max} \coloneqq \max(\boldsymbol{w}_s^b \cup \boldsymbol{w}_{us}^b)$ 2: $w_{\min} \coloneqq \min(\boldsymbol{w}_s^b \cup \boldsymbol{w}_{us}^b)$ 1: $\mathcal{G}\{\cdot\} := \{0\} //$ initialize group bit-width 777 2: $e \coloneqq \inf //$ bit-width searching error 778 3: $p^* \coloneqq 0 //$ number of (N-1)-bit and (N+1)-3: $\lambda \coloneqq 0.1$ 779 4: $n \coloneqq 50$ bit 5: $e \coloneqq \inf H$ scale searching error 4: $l \coloneqq N - 1 //$ lower bit-width 781 6: $\Delta^* \in \mathbb{R}^{n \times 1}$ // per-channel scale 5: $h \coloneqq N + 1 //$ higher bit-width 782 7: $z^* \in \mathbb{R}^{n \times 1}$ // per-channel zero point 6: $S\{\cdot\} \coloneqq \operatorname{average}\left(\frac{\boldsymbol{w}^2}{[\boldsymbol{H}^{\operatorname{in}}]^2_{\operatorname{diag}}}\right)$ 783 8: for $\tau \in [1 - \lambda, 1 + \lambda]$ with 2n slices do 7: for $p = 1, 2, ..., \left[\frac{m}{2\beta}\right]$ do 784 9: $\Delta \coloneqq \tau (w_{\max} - w_{\min}) / (2^{g_s} - 1)$ 785 $\hat{\boldsymbol{w}}_{l}^{b}\coloneqq \mathrm{fakequant}(\boldsymbol{w}_{b\in\mathrm{top}_k_\min(\mathrm{p})}^{b},l,)$ 8: 10: $z \coloneqq -|(\tau w_{\min})/\Delta]$ $\begin{aligned} \hat{\boldsymbol{w}}_{s}^{\text{c}} &:= \operatorname{fakequant}(\boldsymbol{w}_{s}^{b}, g_{b}, \Delta, z) \\ \hat{\boldsymbol{w}}_{s}^{b} &:= \operatorname{fakequant}(\boldsymbol{w}_{us}^{b}, g_{b}, \Delta, z) \\ \hat{\boldsymbol{w}}_{us}^{b} &:= \operatorname{fakequant}(\boldsymbol{w}_{us}^{b}, g_{b}, \Delta, z) \\ \mathcal{L}_{s} &:= ||\boldsymbol{w}_{s}^{b} - \hat{\boldsymbol{w}}_{s}^{b}||^{2} \\ \mathcal{L}_{us} &:= ||\boldsymbol{w}_{us}^{b} - \hat{\boldsymbol{w}}_{us}^{b}||^{2} \end{aligned}$ $\hat{\boldsymbol{w}}_{h}^{b} \coloneqq \text{fakequant}(\boldsymbol{w}_{b \in \text{top}_k_max(p)}^{b}, h,)$ 786 11: 9: 787 $\hat{\boldsymbol{w}}_{h}^{b} \coloneqq \text{fakequant}(\boldsymbol{w}_{b \in \text{others}}^{b}, N,) \\ \hat{\boldsymbol{w}}_{q} \coloneqq \hat{\boldsymbol{w}}_{l}^{b} \cup \hat{\boldsymbol{w}}_{l}^{b} \cup \hat{\boldsymbol{w}}_{h}^{b}$ 12: 10: 788 13: 11: 14: 789 if $\mathcal{D}_{kl} (\boldsymbol{x} \boldsymbol{w}^\top || \boldsymbol{x} \hat{\boldsymbol{w}}_a^\top) < e$ then 12: if $\mathcal{L}_s + \mathcal{L}_{us} < e$ then 15: 790 $e \coloneqq \mathcal{L}_s + \mathcal{L}_{us}$ $e \coloneqq \mathcal{D}_{kl} \left(\boldsymbol{x} \boldsymbol{w}^\top \mid \mid \boldsymbol{x} \hat{\boldsymbol{w}}_a^\top \right)$ 16: 13: 791 $z^* \coloneqq z$ $p^* \coloneqq p$ 17: 14: 792 $\Delta^*\coloneqq\Delta$ 18: 15: end if 793 16: end for 19: end if 794 20: end for 17: $\mathcal{G}{l} \coloneqq S{\operatorname{top}_k_{\min}(p^*) = l}$ 21: $\hat{\boldsymbol{w}}_q^b \coloneqq \text{fakequant}(\boldsymbol{w}^b, g_b, \Delta^*, z^*)$ 18: $\mathcal{G}{h} := S{\text{top}_k_max}(p^*) = h}$ 19: $\mathcal{G}{N} := S{\text{middle}_k}([\frac{m}{2}] - 2p^*) = N}$ 796 22: return \hat{w}_a^b 797 20: return $\mathcal{G}\{\cdot\}$ 798

798 799

Algorithm 2 primarily encompasses the core details of both SBA and SQC. In SBA, the importance of each group is determined by sorting the average salience of groups, followed by a bi-pointer search that increases the number of (N - 1)-bit and (N + 1)-bit groups to maintain their quantity equilibrium. The optimization function then utilizes the KL divergence from Eq. (4) to determine the optimal mixed-precision ratio. SQC, on the other hand, enhances its information by amplifying the quantization error of unstructured weight groups. When the last two parameters, scale and zero point, in the fakequant(\cdot) function are omitted, the default values from Eq. (1) are used.

806 807

808

B.2 MIXED BIT STORAGE AND COMPUTING

809 We developed a framework for storage and inference deployment supporting mixed-precision quantization based on AutoGPTQ. The deployment process is as follows. After completing mixed-precision



Figure 6: The memory layout shown in the figure is modified based on AutoGPTQ. The transposed original weights $w^{\top} \in \mathbb{R}^{m \times n}$ are still divided into multiple groups along the row direction after quantization. The elements within each group are vertically packed into integers and then reassembled into \hat{w}_{int} . The figure employs corresponding colors to indicate how each original number is mapped to a specific position within the packed integers after quantization, which finally generates $\hat{w}_{int} \in \mathbb{R}^{m^* \times n}$, where m^* is compressed from m by packing several low-bit number. Similarly, \hat{z}_{int} is also packed into integers to save memory.

836 837

quantization with SliM-LLM, it outputs scales, zeros, and group-wise bit-width generated during the quantization process to identify the quantization parameters and precision of each group in the Linear Projection weights. AutoGPTQ then packs the weights and zeros into integer-compressed representations (denoted by \hat{w}_{int} and \hat{z}_{int} respectively) based on the precision of different groups, significantly reducing storage and operational bit-width. After the quantized weights are packed, AutoGPTQ loads the model onto the GPU, where the mixed precision quantization kernel on the GPU performs dequantization on the weights and zeros of different groups and calculation with input activation, ultimately producing the final output.

845 In the mixed-precision deployment of AutoGPTQ, the weight memory layout is organized by group, 846 with each group sharing the same precision, which is shown in Fig. 6. Within each group, elements 847 with the same precision are packed as integers, eliminating the need for additional padding, which 848 saves space. Given that the bit-width of integers is a power of 2, this is compatible with group size 849 that is also a power of 2. For instance, even with the odd-bit such as 3-bit storage, integers can store these numbers without padding, as the commonly used group size is 128, a multiple of almost all 850 definition of integer type. This ensures that elements within a group fully utilize the space provided 851 by integers, without storing numbers of different precision within the same integer. \hat{z}_{int} follow the 852 original logic of AutoGPTQ but are packed with a uniform precision along the channel direction 853 for ease of use. Other tensors, like scales, remain in the same floating-point format to ensure the 854 correctness of dequantization calculations. 855

To indicate the precision of each group, we also introduce an additional array to store bit-width of each group, where each number is represented as a 2-bit value aggregated into integers, marking the quantization precision of each group for accurate reconstruction. We use cumulative calculations to determine the starting index of each group, ensuring correctness despite changes in \hat{w}_{int} height and starting indices caused by varying precision. Using the above methods to store the quantized weights, zeros, and additional bit arrays effectively reduces memory usage during model storage and loading, thereby lowering the resource overhead required for model deployment.

863 Once the weights are packed, we follow the modified AutoGPTQ logic for GPU inference. The GPU processes and dequantizes the weights group by group for computation. During GPU computation,



a thread dequantizes a segment of continuous memory data in one column of \hat{w}_{int} and performs vector dot product calculations with the input activation shared within the block, accumulating the results in the corresponding result matrix. When threads form a logical block, the block handles the computation and reduction of a continuous channel region. We complete the linear layer computation by iterating through all logical blocks. Leveraging AutoGPTQ's initial logic and CUDA Warp's 32-thread units, we ensure similar code structure and data access logic for threads within each warp when group size is 128. This method was primarily conducted to validate feasibility os SliM-LLM, demonstrating that the mixed precision quantization with integer packing does not cause additional computational overhead, indicating the efficiency and accuracy advantage of SliM-LLM. In summary, by dividing weight into several structured groups with mixed precision and employing a reasonable GPU utilization strategy, Slim-LLM balances performance and efficiency.

C SEARCHING DETAILS OF GROUP-WISE SALIENCE-DETERMINED BIT ALLOCATION

We optimize the mixed-precision configuration based on the output information entropy (KL-divergence), searching for the optimal compensation bit-width ratio as shown in Eq. (4).

Initially, we rank each group by their average salience, a metric for quantization, and employ a double-pointer that moves simultaneously from both the beginning (lowest salience) and end (highest salience) of the sorted list. This ensures an equal number of groups at low and high bit-widths, effectively balancing the global average bit-width compensation. We then calculate the relative entropy under the corresponding precision ratio and search for the optimal ratio. Fig 7 displays the search error curves related to the 2^{nd} , 10^{th} , and 15^{th} Transformer layers in the OPT1.3B model, showcasing the search curves for certain self-attention layers (Query, Key, Value, FC2).

Due to the limited range of the search, extreme scenarios involve either a half (N - 1)-bit and half (N + 1)-bit without N-bit or all groups being N-bit (uniform precision). Fig 7 demonstrates that lower quantization errors can be achieved under mixed-precision compared to quantization at the uniform bit-width. We also find that multiple low-error precision combinations are possible within a group of weights, allowing SBA to flexibly select the optimal ratio through its versatile search.

911 912

868

870 871 872

873

874

875 876 877

878

879

880

882

883

884

885

887

889

890

891

892 893 894

895

896 897

898

899

D EVLUATION FUNCTION OF SBA

913 914

In Tab. 6, we employ various objective functions and compare their performance in SBA across different models. Compared to the commonly used Mean Squared Error (MSE) loss, Kullback-Leibler (KL) divergence ensures the distribution of critical activation positions within the model from an information entropy perspective, making it a superior choice for the bit-width allocation strategy in

SBA for the OPT and LLaMA models. When computing KL divergence in this context, we first transform the layer outputs into probability distributions using softmax.

	Table 6: Comparison of MSE and KL Divergence in SBA.												
Method	# W	OPT-1.3B	OPT-2.7B	OPT-6.7B	OPT-13B	LLaMA-7B	LLaMA2-7B						
MSE KL Divergence	2-bit 2-bit	32.50 30.71	27.58 13.26	15.14 11.27	13.28 10.12	21.94 14.58	16.86 16.01						

E EXTENSION ABLATION ON SQC

In this section, we visualize the effectiveness of SQC in mitigating the degradation of information in locally salient weights. We observed the absolute error of weights in a randomly selected channel of the quantized OPT-1.3B model. As shown in Fig. 8, the overall absolute error of the weights post-quantization with a standard quantizer was 0.0055, while with SQC it was reduced to 0.0039. This further demonstrates that the search parameter τ , as applied in Eq. (5), effectively optimizes the quantizer parameters, thereby reducing quantization errors.

More importantly, SQC effectively perceives the information of locally salient weights, as indicated
by the red regions in Fig. 8. Compared to the vanilla quantizer, SQC significantly reduces the
error of salient weights. Specifically, the prominent weights at indices 375 in Fig. 8(a) show higher
quantization errors, while in Fig. 8(b), this error is effectively reduced. This confirms SQC's ability
to perceive locally salient weights, effectively preventing the degradation of critical information.

F EXTENSION ABLATION ON QUANTIZATION GROUP-SIZE

To investigate the impact of different group sizes on the quantization effectiveness of SliM-LLM, we evaluated performance with 256 and 512 columns at a 3-bit level, observing that larger group sizes enhance GPU efficiency during inference. The findings suggest that increased group granularity does not substantially elevate perplexity across four models, indicating that SliM-LLM is robust and conducive to more efficient deployment methods. In contrast, at 2-bit, we assessed group sizes of 64 and 32 columns. With finer group granularity, the models displayed reduced perplexity. This is attributed to smaller groups providing more detailed data representation and utilizing additional quantization parameters, although they also raise computational and storage demands. A group size of 128 strikes a better balance between efficiency and quantization performance.

G EXTENSION ON SALIENCE CHANNEL CLUSTERING

G.1 DISCUSSION OF THEOREM 1

Theorem 1. Given the input calibration activation $\mathbf{x} \in \mathbb{R}^{t \times m}$ with an outlier channel $\mathbf{x}_{:,p}^* \gg \mathbf{x}_{:,j}, \forall j \in [0,m], j \neq p$ at the position of channel-p. The trace elements of $\mathbf{H} = \mathbf{x}^\top \mathbf{x}$ will show great outlier value at (p, p), where $\mathbf{H}_{p,p} \gg \mathbf{H}_{j,j}, \forall j \in [0,m], j \neq p$, as $\mathbf{H}_{p,p}$ is produced by $[\mathbf{x}_{:,p}^{*\top}\mathbf{x}_{:,p}^*] = \sum_{i=0}^t x_{i,p}^{*2}$, which further leads to the parameter salience larger at the p^{th} channel of weight, where $\delta_{:,p} > \delta_{:,k}, \delta_{:,k} = \frac{w_{i,k}^2}{[\mathbf{H}^{-1}]_{k,k}^2}, \forall k \in [0,t], k \neq p.$

Proof. Given $x \in \mathbb{R}^{t \times m}$ with outlier channel $x_{i,p}^*$, $p \in [0, m]$, and other elements with small magnitude $x_{i,j}$, where $x_{q,p}^* \gg x_{i,j}$ and $i, j \neq q, p$. We can get the Hessian matrix with Levenberg-



Figure 8: Absolute channel error of the weight of the OPT-1.3B model. The red line represents the quantization error for the locally salient weights, and the lightmauve represents other weights. (a) Vanilla quantizer error on the 794^{th} channel of OPT-1.3B. (b) SQC error on the 794^{th} channel of OPT-1.3B

1005Table 7: Ablation results on OPT-6.7B, LLaMA-7B, LLaMA-2-7B, LLaMA-3-8B with SliM-LLM1007under different group size (#g denotes the group size).

$\textbf{Precision / PPL}{\downarrow}$	#g	OPT-6.7B	LLaMA-7B	LLaMA-2-7B	LLaMA-3-8B
3-bit	512	11.65	6.96	6.69	8.87
	256	11.33	6.92	6.94	8.14
	128	11.27	6.40	6.24	7.62
2-bit	128	14.41	14.58	16.01	39.66
	64	13.95	13.41	15.02	29.84
	32	12.47	11.91	11.95	16.93

1017 Marquardt (Marquardt, 1963) approximation in Eq. (3):

018	$\binom{x_{11}}{x_{21}}$	$x_{12} \\ x_{22}$	$x_{13} \\ x_{23}$	 	$\begin{pmatrix} x_{1m} \\ x_{2m} \end{pmatrix}$	$\binom{x_{11}}{x_{21}}$	$x_{21} \\ x_{22}$	 	$\begin{pmatrix} x_{t1} \\ x_{t2} \end{pmatrix}$		$(x_{11}^2 +)$			`
020 021	:	÷	÷	·	:		÷	·	÷		:	·		÷
022	:	÷	÷	$x_{a,n}^*$:	·	÷	$x^*_{a \ n}$	÷	=	:	÷	$x_{1,n}^{*}^{2} +$	÷
023 024		:	÷	·.	:		÷	·.	÷				1,p 	·
)25	$\backslash x_{t1}$	x_{t2}	x_{t3}		x_{tm} /	$\backslash x_{1m}$	x_{2m}		x_{tm}		X			(6)

1038 1039

1046

1075 1076

where $[x_{:,p}^{*\top}x_{:,p}^{*}]$ will appears at position $H_{p,p}$. And following SparseGPT (Frantar & Alistarh, 2023), the inverse matrix of H can be formulated as:

$$\delta_{i,j} = \frac{w_{i,j}^2}{[\operatorname{diag}((\boldsymbol{x}^{\top}\boldsymbol{x} + \lambda \boldsymbol{I})^{-1})]^2}$$
(7)

where $(\boldsymbol{x}^{\top}\boldsymbol{x} + \lambda \boldsymbol{I})^{-1}$ is the new representation of Hessian matrix \boldsymbol{H} for the layer-wise reconstruction problem, and λ is the dampening factor for the Hessian to prevent the collapse of the inverse computation. Additionally, in accordance with the configuration in LLMs (Frantar & Alistarh, 2023; Frantar et al., 2022; Sun et al., 2023), the value of λ set is extremely small ($\lambda \leq e^{-1}$), while the values located at the diagonal of Hessian are large. Therefore, only considering the influence of diagonal elements (Sun et al., 2023), we can further approximate salience as:

$$\delta_{i,j} = \frac{w_{i,j}^2}{[\operatorname{diag}((\boldsymbol{x}^{\top}\boldsymbol{x} + \lambda \boldsymbol{I})^{-1})]^2} \approx \frac{w_{i,j}^2}{[(\operatorname{diag}(\boldsymbol{x}^{\top}\boldsymbol{x}))^{-1}]^2} = (w_{i,j} \cdot ||\boldsymbol{x}_j||_2^2)^2 \tag{8}$$

Here the diagonal of $x^{\top}x$ is diag $(||x_j||_2^2)$, and $||x_j||_2$ evaluates the ℓ_2 norm of j^{th} channel across different tokens. Consequently, it can be summarized that when there is an outlier channel-p, the value of $||x_p||_2$ is primarily influenced by $[x_{:,p}^{*\top}x_{:,p}^*]$. Additionally, since the activation values are relatively large and the differences in weight values are comparatively small, the p^{th} channel of weights will also exhibit salience.

G.2 DISTRIBUTION OF SALIENCE, ACTIVATION AND WEIGHT MAGNITUDE

Fig. 9 illustrates the distribution of salience among certain weights in LLMs. This section provides additional examples to demonstrate how the distribution of weights and input activation characteristics influence the salience of parameters in LLMs. The figure captures seven linear projections in the multihead self-attention (MHA) and feed-forward block (FFB) layers of the 2^{nd} and 10^{th} Transformer modules in the LLaMA-7B model.



Figure 9: Salience, activation and weight distribution in the 2nd and 10th layers of LLaMA-7B

In line with previous findings (Nrusimha et al., 2024; Xiao et al., 2023a), activations demonstrate particularly marked outlier phenomena on anomalous tokens and channels, with extremes differing by more than two orders of magnitude. Notably, distinct anomalous channels are present in the MHA's Query, Key, and Value layers, where outliers vary significantly across different tokens. This

pattern is consistent in the FFB layers. We observe that disparities in weight magnitudes are less pronounced than those in activation, thus exerting a reduced impact on outlier channels. Moreover, weights distribute structurally along rows or columns (Dettmers et al., 2023; Huang et al., 2024a), affecting the overall distribution of salience from a row-wise perspective (Fig. 9). However, the most prominent salience is predominantly driven by activation across channels (column-wise).

1086 G.3 HESSIAN DIAGONAL CLUSTERING

Sec. 3.2.1 demonstrates that outlier tokens in input activations result in significant values at the 1088 corresponding positions along the diagonal of the weight Hessian matrix. Additionally, due to 1089 the token sink phenomenon (Xiao et al., 2023b; Nrusimha et al., 2024), areas around significantly 1090 activated key tokens exhibit increased salience, creating clusters of salient regions along the Hessian 1091 matrix diagonal. To further elucidate this phenomenon, Fig. 10 shows the values along the diagonal of 1092 the Hessian matrix for selected weights in the 2^{nd} and $10^{\bar{t}h}$ layers of the LLaMA-7B model. Within 1093 this diagonal, certain positions display pronounced values (indicated in red), whereas others are 1094 relatively moderate. In the attention aggregation layer of the 10^{th} layer, the token sink phenomenon 1095 results in a pronounced convergence of significant values along the Hessian matrix diagonal, with deep 1096 red areas indicating regional clustering. These findings reinforce the influence of input activations on the diagonal of the Hessian matrix, subsequently leading to a clustering phenomenon in the salience 1098 distribution of weights across channels.





1121

1087

1099

Figure 10: Hessian diagonal magnitude in attention layers of 2^{nd} and 10^{th} layers of LLaMA-7B

1122 H MORE COMPARISONS

1124 In this section, we provide supplementary experiments for SliM-LLM. Tab. 8 displays the comparative 1125 results of SliM-LLM and SliM-LLM+ with other methods on the OPT series models. Tab. 9 shows 1126 the performance of SliM-LLM when quantizing the LLaMA family models on the C4 dataset, 1127 while Tab. 10 also compares the results of SliM-LLM+ on the C4 dataset. In Tab. 11, we compared 1128 the quantization results of GPTQ, AWQ, and SliM-LLM at 2-bit on the Gemma2 and Mixtral 1129 models, demonstrating the greater stability of SliM-LLM across a wider range of model structures. Additionally, in Tab. 12, we supplemented the 4-bit results of different quantization methods in the 1130 LLaMA series models, showing that SliM-LLM and SliM-LLM+ exhibit the smallest quantization 1131 errors at practical 4-bit levels. To provide a comprehensive evaluation across a broader set of 1132 benchmarks, we further compared the quantization results on MMLU and MathQA in Tab. 13. 1133

1	1	3	4	

#W PPL↓	Method	1.3B	2.7B	6.7B	13B	30B
16-bit	-	14.63	12.47	10.86	10.12	9.56
	RTN	1.2e2	3.0e2	23.54	46.03	18.80
	GPTQ	16.47	13.69	11.65	10.35	9.73
	AWQ	16.32	13.58	11.41	10.68	9.85
3 hit	QuIP	16.21	13.79	11.51	10.50	9.75
5-011	SliM-LLM	15.91	13.26	11.27	10.26	9.70
	ÖmniQuant	15.72	13.18	11.27	- 10.47 -	9.79 -
	AffineQuant	15.61	12.98	11.18	10.51	9.81
	SliM-LLM ⁺	15.58	12.84	11.18	10.44	9.67
	RTN	1.3e4	5.7e4	7.8e3	7.6e4	1.3e4
	GPTQ	1.1e2	61.59	20.18	21.36	12.71
	AWQ	47.97	28.50	16.20	14.32	12.31
2-bit	QuIP	41.64	28.98	18.57	16.02	11.48
	PB-LLM	45.92	39.71	20.37	19.11	17.01
	SliM-LLM	30.71	24.08	14.41	13.68	11.34
	OmniQuant	23.95	18.13	14.43	12.94	11.39
	SliM-LLM ⁺	24.57	17.98	14.22	12.16	11.27

Table 9: Quantization results of LLaMA Family with statistic quantizer on C4 (group size is 128).

1160	#W PPL↓	Method	1-7B	1-13B	1-30B	1-65B	2-7B	2-13B	2-70B	3-8B	3-70B
1161	16-bit	-	7.08	6.61	5.98	5.62	6.97	6.46	5.52	9.22	6.85
1163		APTQ PTN	6.24	-	-	-	-	-	-	-	-
1164 1165	3-bit	AWQ	8.02 7.92 7.85	7.07	6.37 6.47	5.94	7.84 7.80	6.94 7.00	- 5 85	11.62 11.62	8.03
1166		SliM-LLM	6.14	6.05	6.33	5.94	7.89 7.74	5.26	5.85 5.09	13.07	8.64
1167		RTN	1.0e3	4.5e2	99.45	17.15	4.9e3	1.4e2	42.13	2.5e4	4.6e5
1100		AWQ	1.9e5	2.3e5	2.4e5	7.5e4	1.7e5	9.4e4	-	2.1e6	1.4e6
1169	2-bit	GPTQ	34.63	15.29	11.93	11.99	33.70	20.97	NAN	4.1e4	21.82
1170	2-011	QuIP	33.74	21.94	10.95	13.99	31.94	16.16	8.17	1.3e2	22.24
1171		PB-LLM	49.73	26.93	17.93	11.85	29.84	19.82	8.95	79.21	33.91
1172		SliM-LLM	32.91	13.85	11.27	10.95	16.00	9.41	7.01	1.1e2	15.92

Table 10: Quantization results of LLaMA-1 and LLaMA-2 models with learnable quantizer on C4.

Method	1-7B	1-13B	1-30B	1-65B	2-7B	2-13B	2-70B
-	7.08	6.61	5.98	5.62	6.97	6.46	5.52
OmniQuant	7.75	7.05	6.37	5.93	7.75	6.98	5.85
AffineQuant	7.75	7.04	6.40	-	7.83	6.99	-
SliM-LLM ⁺	7.75	6.91	6.36	5.96	7.71	6.90	5.85
OmniQuant	12.97	10.36	9.36	8.00	15.02	11.05	8.52
AffineQuant	14.92	12.64	9.66	-	16.02	10.98	-
SliM-LLM ⁺	14.99	10.22	9.33	7.52	18.18	10.24	8.40
	Method - OmniQuant AffineQuant SliM-LLM ⁺ OmniQuant AffineQuant SliM-LLM ⁺	Method 1-7B - 7.08 OmniQuant 7.75 AffineQuant 7.75 SliM-LLM ⁺ 7.75 OmniQuant 12.97 AffineQuant 14.92 SliM-LLM ⁺ 14.99	Method 1-7B 1-13B - 7.08 6.61 OmniQuant 7.75 7.05 AffineQuant 7.75 7.04 SliM-LLM ⁺ 7.75 6.91 OmniQuant 12.97 10.36 AffineQuant 14.92 12.64 SliM-LLM ⁺ 14.99 10.22	Method 1-7B 1-13B 1-30B - 7.08 6.61 5.98 OmniQuant 7.75 7.05 6.37 AffineQuant 7.75 7.04 6.40 SliM-LLM ⁺ 7.75 6.91 6.36 OmniQuant 12.97 10.36 9.36 AffineQuant 14.92 12.64 9.66 SliM-LLM ⁺ 14.99 10.22 9.33	Method 1-7B 1-13B 1-30B 1-65B - 7.08 6.61 5.98 5.62 OmniQuant 7.75 7.05 6.37 5.93 AffineQuant 7.75 7.04 6.40 - SliM-LLM ⁺ 7.75 6.91 6.36 5.96 OmniQuant 12.97 10.36 9.36 8.00 AffineQuant 14.92 12.64 9.66 - SliM-LLM ⁺ 14.99 10.22 9.33 7.52	Method 1-7B 1-13B 1-30B 1-65B 2-7B - 7.08 6.61 5.98 5.62 6.97 OmniQuant 7.75 7.05 6.37 5.93 7.75 AffineQuant 7.75 7.04 6.40 - 7.83 SliM-LLM ⁺ 7.75 6.91 6.36 5.96 7.71 OmniQuant 12.97 10.36 9.36 8.00 15.02 AffineQuant 14.92 12.64 9.66 - 16.02 SliM-LLM ⁺ 14.99 10.22 9.33 7.52 18.18	Method 1-7B 1-13B 1-30B 1-65B 2-7B 2-13B - 7.08 6.61 5.98 5.62 6.97 6.46 OmniQuant 7.75 7.05 6.37 5.93 7.75 6.98 AffineQuant 7.75 7.04 6.40 - 7.83 6.99 SliM-LLM ⁺ 7.75 6.91 6.36 5.96 7.71 6.90 OmniQuant 12.97 10.36 9.36 8.00 15.02 11.05 AffineQuant 14.92 12.64 9.66 - 16.02 10.98 SliM-LLM ⁺ 14.99 10.22 9.33 7.52 18.18 10.24

Table 11: PPL Comparison on Gemma2 and Mixtral.						
Model/Evaluation	Method	PPL (wikitext2)				
	GPTQ 2-bit	186.77				
Gemma2-9B	AWQ 2-bit	217.83				
	SliM-LLM 2bit	26.30				
	GPTQ 2-bit	16.38				
Mixtral 8x7B	AWQ 2-bit	3.2e5				
	SliM-LLM 2bit	7.44				

Table 12: The PPL results of our proposed method and other methods under 4bit quantization.

Method	LLaMA-7B	LLaMA-13B	LLaMA2-7B	LLaMA2-13B	LLaMA3-8B
FP16	5.68	5.09	5.47	4.88	5.75
AWQ	5.81	5.30	5.62	4.97	6.63
GPTQ	5.85	5.20	5.61	4.98	6.50
SliM-LLM	5.83	5.16	5.59	4.95	6.42
Omniquant	5.77	-	5.58	-	-
SliM-LLM ⁺	5.75	-	5.57	-	-

_

Table 13: The results(%) on MMLU and MathQA for multiple quantized LLaM	A models.
---	-----------

Model	Method	Humanities	Social Sciences	STEM	Other	MMLU	MathQ
	GPTQ 2-bit	24.87	21.84	21.79	24.01	23.32	21.11
LLaMA-7B	AWQ 2-bit	24.21	21.71	21.25	23.98	22.95	22.21
	SliM-LLM 2bit	24.94	23.60	23.40	25.50	25.10	23.74
	GPTQ 2-bit	24.23	23.20	22.99	24.78	23.85	21.68
LLaMA-13B	AWQ 2-bit	24.17	31.07	28.61	25.14	26.89	21.98
	SliM-LLM 2bit	25.12	31.74	29.19	26.17	27.05	23.17
	GPTQ 2-bit	25.02	22.13	22.61	23.17	23.44	21.07
LLaMA2-7B	AWQ 2-bit	25.12	22.79	24.26	24.01	24.51	19.06
	SliM-LLM 2bit	26.60	23.23	25.70	25.70	25.81	22.55
	GPTQ 2-bit	23.91	27.17	26.10	25.78	25.53	20.87
LLaMA2-13B	AWQ 2-bit	24.17	31.07	28.61	25.14	26.89	19.53
	SliM-LLM 2bit	26.27	32.20	29.98	26.46	27.34	23.48

1242 I REAL DIALOG EXAMPLES

In this section, we show some dialogue examples of LLaMA-2-13B and Vicuna-13B with SliM-LLM-2bit and GPTQ-2bit in Fig. 11.

Prompt:	The eight planets in Solar System are Mercury,
GPTQ-2bit:	, 10.2. The eight pill o free games of p. 62.2. The way we do this.
Slim-LLM-2bit:	Venus, Earth, Mars, Jupiter, Saturn, Uranus and Neptune.
Prompt:	The tallest and most famous tower in France is
GPTQ-2bit:	is 350 and most recently.
Slim-LLM-2bit:	320 meters high.
Prompt:	The 2012 Summer Olympic was held in
GPTQ-2bit:	The 2012 Summer Olympic was held in 1016 (a.185) and the 775,000 km from each.
Slim-LLM-2bit:	2012 in London, the United Kingdom.
Vicuna-13B	
Prompt:	What are the capitals of Brazil, Mexico and France?
GPTQ-2bit:	What is the RER score?
Slim-LLM-2bit:	The capital of Brazil is Brasília, the capital of Mexico is Mexico City, and the capital of France is Paris.
Prompt:	What are the ancient civilizations in Europe?
GPTQ-2bit:	I am trying to find a book to nook over the life of to 300,000 or so. The ancient civilization is the one in which history books be pop and I have a be, oh, do you know the one the Mayan, the one of the total for Dake.
Slim-LLM-2bit:	The ancient civilizations in Europe include the Roman Empire, the Byzantine Empire and the Ottoman Empire.
Prompt:	What are the logical circuits?
GPTQ-2bit:	What is the order of the logical cumulated? What is the best strategy to do in a single? What is the most popular to be eaten? What are the most sought?
Slim-LLM-2bit:	The logical circuits are the basic components for building the logical functions, which are required for digital computing. These circuits use logic structures and devices that allow the use of Boolean and multiple logic to implement the functional requirements as defined by a problem solving method of computer. In digital computing, circuits perform the logical tasks required for solving problems.

Figure 11: Some examples of conversations. LLaMA-2-13B and Vicuna-13B are chosen to show the case of language supplementary and Q&A ability. And GPTQ-2bit is selected as the comparison. We color the text to show the reasonable or inappropriate responses.

1286 1287

1288

1289

1290

1291

1292

1293