# THE UNIQUENESS OF LLAMA3-70B SERIES WITH PER-CHANNEL QUANTIZATION

**Minghai Qin**
Western Digital Research
`minghai.qin@wdc.com`

October 2, 2024

## ABSTRACT

We have observed a distinctive quantization-related behavior in the LLaMA3/3.1-70B models that is absent in both the LLaMA2-70B and LLaMA3/3.1/3.2-1B/3B/8B/405B models. Quantization is a crucial technique for deploying large language models (LLMs) efficiently, as it reduces memory consumption, decreases memory transactions, and potentially accelerates inference using lower-precision compute cores. Among various bit widths and representations for weights and activations, the 8-bit integer weight and 8-bit integer activation (W8A8) configuration is particularly popular due to its widespread hardware support. However, the impact of W8A8 post-training quantization on model accuracy, especially on the recently released LLaMA3/3.1 model series, remains contentious. While several studies have suggested calibrating either weights or activations to mitigate accuracy degradation, a comprehensive solution has yet to be identified. In this paper, we explore three key questions: What makes the LLaMA3-70B model series uniquely vulnerable to quantization? Why is this the case? And how can the issue be addressed? We empirically investigate multiple LLMs featured on an open LLM leaderboard, discovering that the LLaMA3-70B model series have a **unique** accuracy degradation behavior with W8A8 per-channel post-training quantization. In contrast, other model series such as LLaMA2, LLaMA3/3.1-8B, LLaMA3.2-1B/3B, Qwen, Mixtral, Mistral, Phi-3, and Falcon demonstrate robust performance with W8A8, sometimes surpassing their FP16 counterparts. Contrary to previous assertions attributing degradation to the large dynamic range of activations, our findings indicate that the weight distribution of the LLaMA3-70B is the primary factor behind the vulnerability. By meticulously analyzing the distinct characteristics of weight distributions across Transformer blocks, we propose two solutions that make different tradeoffs in hardware/software overhead. First, we propose a mixed strategy where less than 3% of the layers employ finer per-group W8A8 quantization granularity, while the remaining 97% retain the per-channel configuration. Second, we introduce a bi-smoothing strategy that balances quantization errors between weights and activations while maintaining per-channel quantization throughout. Experimental results demonstrate that both strategies effectively preserve the accuracy of the entire LLaMA3-70B model series under W8A8 quantization, achieving performance on par with their FP16 counterparts.

***Keywords*** LLaMA models, quantization

## 1 Introduction

The rising popularity of large language models (LLMs) Vaswani et al. [2023], Brown et al. [2020], Devlin et al. [2019], Liu et al. [2019], Touvron et al. [2023] in diverse practical applications underscores their revolutionary impact on natural language processing, data analysis, and artificial intelligence. However, their immense size imposes a substantial burden on GPU memory Kaplan et al. [2020], complicating deployment in resource-constrained environments. Quantization is a pivotal technique to mitigate this issue, reducing the memory footprint of LLMs while maintaining accuracy Gholami
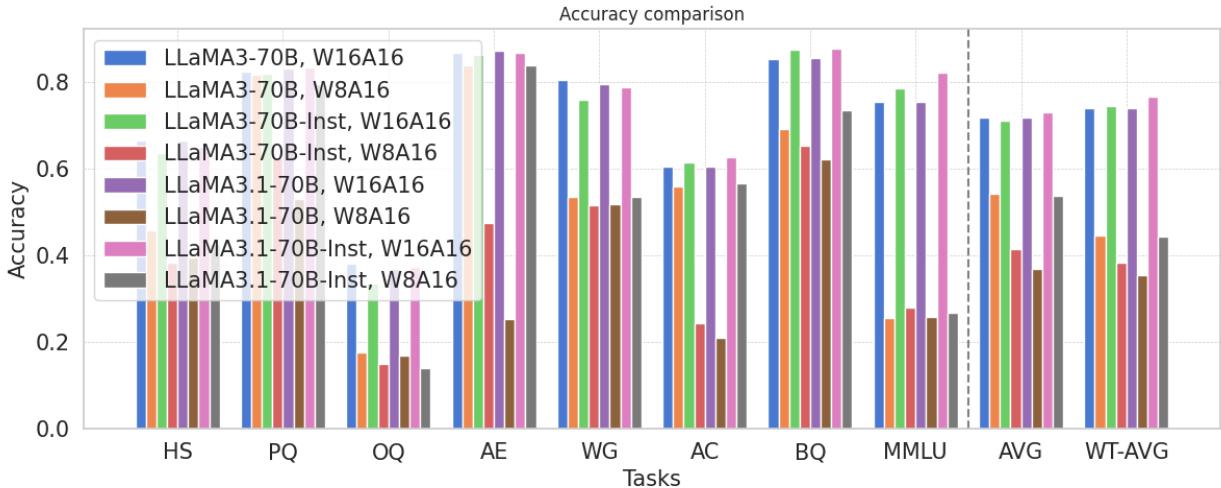
Figure 1: Accuracy of LLaMA3-70B model series with W16A16, and W8A16 per-channel quantization

et al. [2021], Zafrir et al. [2019], Wang et al. [2023]. By converting high-precision weights and activations to lower-precision formats, quantization facilitates more efficient model storage and accelerates inference. As the demand for LLMs grows, the advancement of quantization methods will be crucial for their scalable and cost-effective deployment in real-world applications.

Quantization uses lower precision to represent weights and activations. In this study, we focus on post-training integer-bit quantization (PTQ). The 8-bit weights and 8-bit activations (W8A8) configuration is often regarded as a good balance between efficiency and accuracy. Compared to the FP16 counterpart, W8 can save reduce the memory consumption and transaction by 50%. Although 4-bit weights and 16-bit activations (W4A16) require even less memory, they often lack broad hardware support for the INT4-FP16 matrix-multiplication compute cores. In contrast, existing INT8-INT8 compute cores offer superior acceleration for the inference of W8A8 models. On the other hand, while the W4A4 configuration may reduce memory usage and latency through INT4-INT4 compute cores, it typically suffers from significant accuracy degradation due to the aggressive quantization.

The granularity of a quantization group in large language models (LLMs) is another critical factor influencing the accuracy-acceleration trade-off. Within a quantization group, all values share the same scaling factor and zero point. When the group size matches the input dimension of a layer, this is termed "per-channel" quantization, which maximizes potential speedup. Conversely, "per-group" quantization occurs when an input channel is divided into multiple groups. This requires the intermediate partial sums of matrix-matrix multiplication to be transferred out of the compute cores, significantly reducing system throughput. Therefore, maintaining accuracy with per-channel quantization is essential for practical applications.

In this paper, we investigate the W8A8 quantization for the most popular architectures and models featured on the LLM open leaderboard HuggingFace [2024]. Our study reveals a noteworthy phenomenon: the LLaMA3-70B model series Dubey et al. [2024] exhibits a pronounced vulnerability to W8A8 per-channel quantization, in stark contrast to other models, which demonstrate significantly greater robustness, typically experiencing less than 1% accuracy degradation compared to their FP16 counterparts under the same quantization scheme. Our observation on the "LLaMA3-70B series" also include the LLaMA3.1-70B and various fine-tuned versions, and other robust models include LLaMA2 series, LLaMA3-8B, Mistral-123B, Mixtral-8x7B, Qwen2-72B, Phi3-13B, and Falcon-40B Jiang et al. [2023, 2024], Yang et al. [2024], Abdin et al. [2024], Almazrouei et al. [2023]. Figure 1 provides a comparative analysis of FP16 (labeled as W16A16) and W8A16 quantization for the officially released LLaMA3/3.1-70B and LLaMA3/3.1-70B-Instruct by Meta. This comparison reveals that significant accuracy degradation already occurs with W8A16, indicating that the issue primarily stems from weight quantization rather than activation quantization. Details will be presented in the experiment section. Released in April 2024, the LLaMA3-70B has not been extensively examined in recent research on W8A8 per-channel quantization. Given that this series is currently the most widely used base model on the open leaderboard, further exploration and understanding of its unique characteristics are imperative.

To elucidate the unique vulnerability of the LLaMA3-70B series to W8A8 per-channel quantization, we conducted a comparative analysis of its weight distributions against those of other robust models. Our investigation revealed a
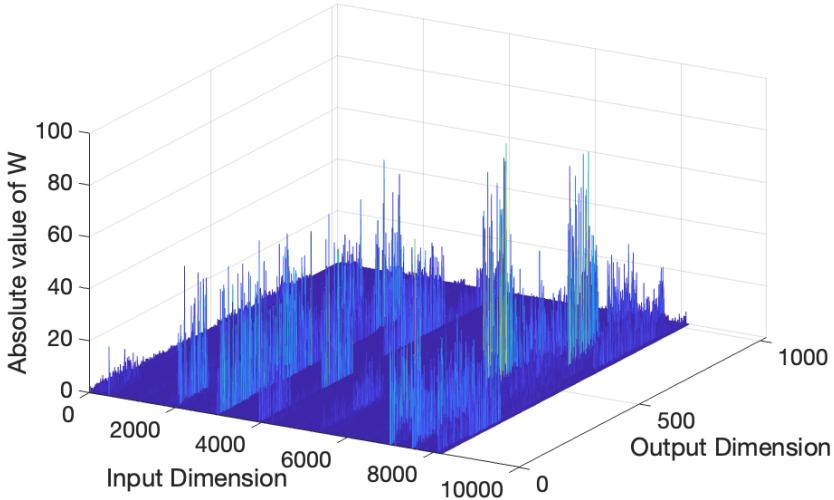
Figure 2: abs_weights of the "V" matrix in the first Transformer block of LLaMA3-70B

significant discrepancy, particularly pronounced in the initial layers. These layers exhibit a distinctive characteristic wherein magnitudes of certain weights exceed others by several orders of magnitude. Figure 2 shows an example of the "V" matrix in the first Transformer block of the LLaMA3-70B model. Weights with large magnitudes are concentrated on specific input dimensions, with the largest values exceeding 90. These outlier weights substantially expand the quantization range, resulting in larger quantization intervals and consequently diminished precision for smaller weight values. Notably, these weight outliers are only prominent in the "Q", "K", "V", "Up", and "Gate" weight matrices of the initial transformer blocks, while the "O" and "Down" weight matrices do not exhibit this behavior. We propose two methods to address this issue. First, to mitigate these quantization errors, we propose implementing per-group quantization specifically for these affected layers. Although these layers constitute merely 2-3% of the total layers in LLaMA3-70B, this targeted approach yields a dramatic enhancement in accuracy, elevating model performance to a level approaching that of FP16 models. Second, we demonstrate a bi-directional smoothing method of the weight and activation matrices that effectively reduces the dynamic range of both, significantly minimizing quantization errors. This technique preserves per-channel quantization while enhancing the model's resilience to W8A8 quantization.

The contribution of this paper is as follows.

- We show that the LLaMA3-70B model series has a unique characteristic that makes it the only model series on the open LLM leaderboard vulnerable to W8A8 per-channel post-training quantization (PTQ). In contrast, all other models can be safely quantized to per-channel-W8A8 with negligible accuracy loss. We believe this observation holds significant implications since LLaMA3-70B serves as one of the most popular base models for numerous models fine-tuned for domain-specific tasks.

- We identify that LLaMA3-70B's sensitivity to per-channel PTQ stems from its significantly different weight distribution, particularly the presence of substantial weight outliers.

- To address this, we first propose a mixed per-channel/per-group quantization strategy. This approach applies per-group quantization to less than 3% of the layers, specifically those with significant weight outliers, while maintaining per-channel quantization for the remaining 97% layers. This mixed strategy effectively restores the accuracy of the LLaMA3-70B model series to levels comparable to its FP16 counterparts.

- The second method introduces bi-directional smoothing to mitigate outliers in both weights and input activations. This approach significantly reduces the dynamic range required for quantization, enabling the LLaMA3-70B model series to retain accuracy comparable to their FP16 counterparts.

## 2  Quantization Basics

In this paper, we use symmetric integer-bit quantization due to its high efficiency and simpler hardware support. In symmetric $n$-bit integer quantization, we map floating-point (FP) values to integers in $[-2^{n-1}, 2^{n-1} - 1]$. This process

involves scaling the FP values based on a scale-factor and an offset. This offset is sometimes called "zero-point" and is set to 0 for symmetric quantization.

Let $x$ be a FP value, and $x_q$ its quantized integer representation. The scale-factor $s$ is determined by the range of the FP values and the desired range of the integers. The symmetric quantization formula is given by:

$$x_q = \text{round}\left(\frac{x}{s}\right)$$

where $\text{round}(\cdot)$ rounds the result to the nearest integer.

The scale-factor $s$ is typically calculated based on the maximum absolute value of the FP values, max_abs, and the maximum representable integer value, $2^{n-1} - 1$. For symmetric quantization, the scale-factor is defined as:

$$s = \frac{\text{max\_abs}}{2^{n-1} - 1}$$

Thus, the quantization and dequantization processes can be expressed as:

$$x_q = \text{round}\left(\frac{x}{s}\right) \quad \text{and} \quad x \approx x_q \cdot s$$

The scale-factor $s$ plays a crucial role in the quantization process. The number of floating-point values used to determine max_abs is typically referred to as the group size. In LLMs, consider a weight matrix $W$ of dimension $M \times N$ and an input activation matrix $A$ of dimension $L \times N$. The resultant output of $A \cdot W^T$ has dimension $L \times M$. Here $L$, $N$, and $M$ represent the sequence length, hidden dimension, and output dimension in a Transformer layer, respectively. For per-channel quantization of $W$, the group size is $N$, implying that each row of $N$ values in $W$ shares a common scale-factor. Consequently, all scale-factors form a vector of length $M$. Similarly, per-channel quantization of $A$ also has a group size of $N$, with each row of $N$ values in $A$ sharing a scale-factor, resulting in a scale-factor vector of length $L$. This configuration enables the computation of floating-point matrix multiplication using integer-bit compute cores (e.g., INT8 tensor cores), followed by dequantization. The dequantization process involves an element-wise multiplication with the scale-factor matrix of dimension $L \times M$, derived as the outer product of the scale-factor vectors of $W$ and $A$. Notably, this dequantization can often be efficiently fused into the subsequent operation after $A \cdot W^T$ in LLMs, which is typically an element-wise operation such as rotary embedding, SiLU activation application, or normalization. In contrast, per-group quantization of $W$ and $A$ employs a finer granularity, where each row of $W$ (and of $A$) is partitioned into groups, each with its own scale-factor. Consequently, the output $A \cdot W^T$ cannot be computed in a single cycle within the INT tensor core, as the scale-factors differ during accumulation in the multiply-accumulation (MAC) unit. This leads to a degradation in the compute efficiency of the quantized model. Given these considerations, per-channel quantization is generally preferred unless the quantization-induced error degrades the model's accuracy to an unacceptable level.

## 3    The Uniqueness of LLaMA3-70B with per-channel W8A8

### 3.1    Test Dataset

We evaluated LLMs across eight reasoning tasks: Hellaswag (HS), PIQA (PQ), OpenbookQA (OQ), ARC-Easy (AE), Winogrande (WG), ARC-Challenge (AC), BoolQ (BQ), and MMLUZellers et al. [2019], Bisk et al. [2019], Mihaylov et al. [2018], Clark et al. [2018], Sakaguchi et al. [2019], Clark et al. [2019], Hendrycks et al. [2021a,b]. We report the "accuracy" for all tasks. Some tasks have "normalized accuracy" and the conclusions derived from "accuracy" remain consistent when applied to "normalized accuracy". We realize that many studies report the average accuracy over these reasoning tasks by simply computing the arithmetic mean of individual task accuracies (denoted by "AVG" in our figures). However, this approach may lack rigor due to the significant variance in the number of questions across tasks. For instance, HS and MMLU each contain over 10,000 questions, while OQ comprises only 500.

To address this disparity, we use an additional metric: the "weighted average accuracy." This metric is calculated by dividing the total number of correct answers across all tasks by the total number of questions. Consequently, this method assigns greater weight to datasets with a larger number of questions, providing a more balanced representation of model performance across tasks of varying sizes. It is denoted by "WT-AVG" in our figures.
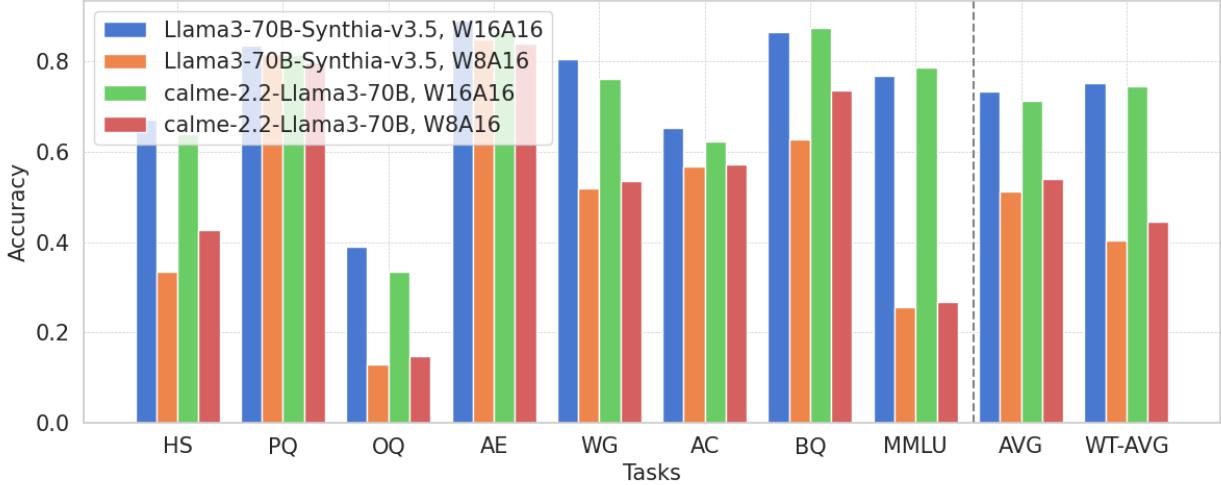
Figure 3: W8A8 per-channel quantization accuracy of models using LLaMA3-70B as the base model

## 3.2 Tested Models

We test the accuracy of the following models: LLaMA3-70B, LLaMA3-70B-Instruct, LLaMA3.1-70B, Llama3-70B-Synthia, calme-2.2-llama3-70b, LLaMA3-8B, LLaMA2-70B, LLaMA2-70B-chat, Qwen2-72B, Mixtral-8x7B, Phi3-14B-Instruct, Mistral-L-Instruct-123B, Falcon-40B. HuggingFace links to them will be provided in the appendix. We also analyze the weight distribution of LLaMA3.1-405B without evaluating its accuracy with FP16 due to GPU memory limitations.

## 3.3 Test Environments

All tests are done within 8xA100 (80GB) GPUs with PyTorch framework using lm-evaluation-harness Gao et al. [2023].

## 3.4 LLaMA3-70B's sensitivity to W8A8

Figure 1 and Figure 3 present a comparative analysis of accuracy for several models in the LLaMA3-70B series subjected to 8-bit per-channel post-training quantization. This analysis encompasses the recently released LLaMA3.1-70B and two additional models fine-tuned from LLaMA3-70B. These figures yield several notable observations:

- The accuracy of models undergoes significant degradation with W8 quantization, even when activations are maintained at FP16 precision. This indicates that the observed accuracy deterioration is not attributable to quantization errors in activations, but rather originates from the 8-bit weight quantization process.

- Models fine-tuned from LLaMA3-70B exhibit the same vulnerability to W8 quantization as their base model. This persistence of vulnerability can be attributed to the fact that fine-tuning typically induces only minor alterations to weight values. In the subsequent section, we will demonstrate that this vulnerability is intrinsically linked to the underlying weight distributions of the model.

## 3.5 LLaMA3-70B is the ONLY ONE sensitive to W8A8

Figure 4 , Figure 6 and Figure 7 compare the accuracy of other models with different architectures with W8A8 quantization. The analysis spans the LLaMA family, including LLaMA3-8B, LLaMA2-70B, and LLaMA2-70B-chat, as well as other prominent Transformer-based designs such as Qwen2-72B, Mixtral-8x7B, Mistral-123B, Phi3-13B, and Falcon-40B. These architectures represent the top-performing architectures in the LLM Open Leaderboard.

The data presented in these figures illustrates a contrast to the behavior observed in the LLaMA3-70B series. Specifically, the models examined here exhibit remarkable resilience to W8A8 per-channel quantization. In the majority of instances, the performance of W8A8 quantized versions closely approximates that of their FP16 counterparts.
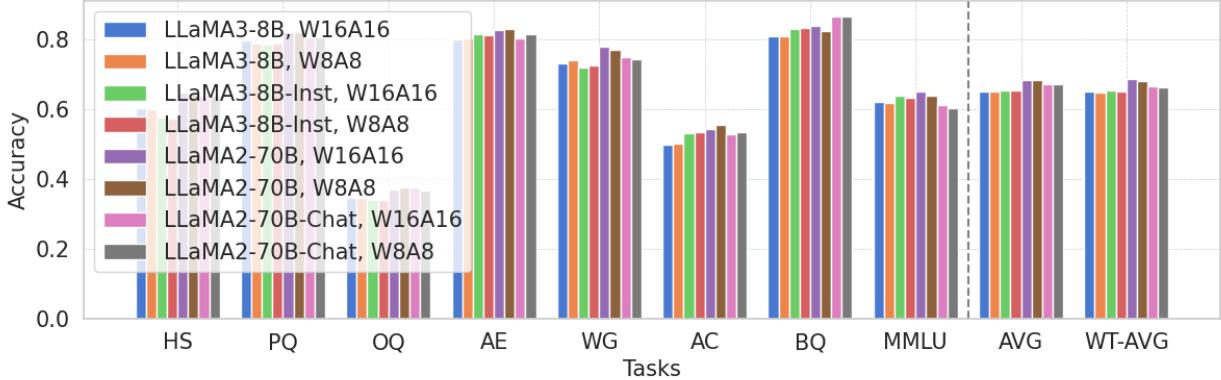
Figure 4: Accuracy of other LLaMA models with W8A8 per-channel quantization
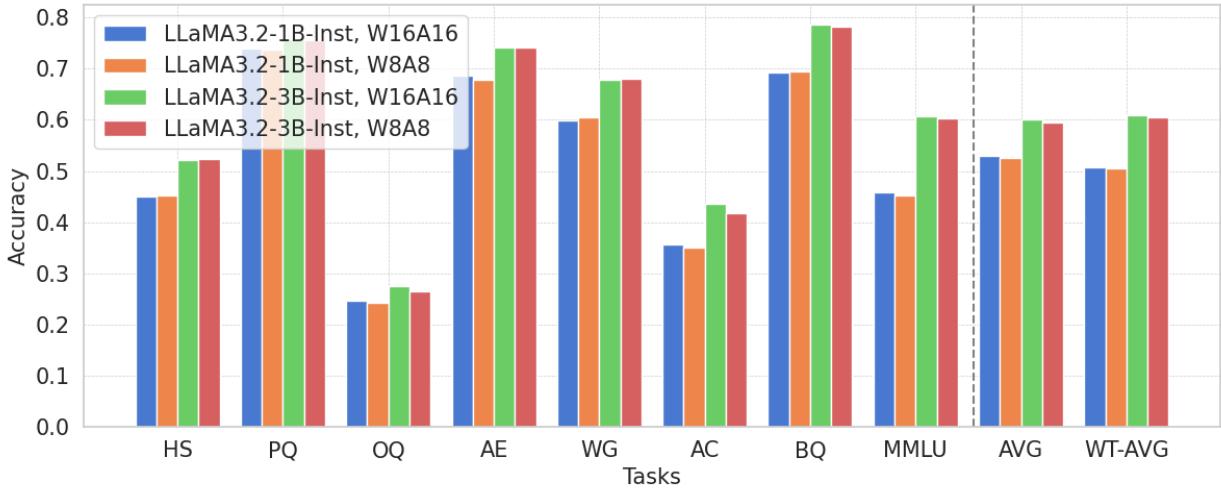


Figure 5: LLaMA3.2 models with W8A8 per-channel quantization

## 3.6 Why LLaMA3-70B is the only one?

In order to understand the uniqueness of LLaMA3-70B models, we explore the weight distributions.

### 3.6.1 A qualitative exploration

Figures 2, 8, 9, and 10, illustrate the absolute values of the V matrix weights in the first Transformer block for LLaMA3-70B, LLaMA3.1-70B, LLaMA3-8B, and LLaMA2-70B models, respectively. Other weight matrices for other models show similar trend and are shown in the Appendix. It is important to note that in per-channel quantization, each group comprises a length-$N$ vector along the input dimension. Our analysis reveals several key observations:

- The maximum absolute weight values (max_abs) in LLaMA3-70B (93) and LLaMA3.1-70B (92.5) surpass those in LLaMA3-8B (0.05) and LLaMA2-70B (0.07) by approximately three orders of magnitude. Furthermore, LLaMA3/3.1-70B exhibits a non-trivial number of weights with large magnitudes along the output dimension, which we term "weight outliers". This extensive range of weight values results in large quantization intervals for LLaMA3/3.1-70B, leading to substantial quantization errors in groups containing any weight outliers.

- LLaMA3-70B displays a distinct pattern where weights with large max_abs values cluster at specific input indices, forming visible "walls" in Figure 2 and 8. In contrast, LLaMA2-70B and LLaMA3-8B lacks such obvious patterning.
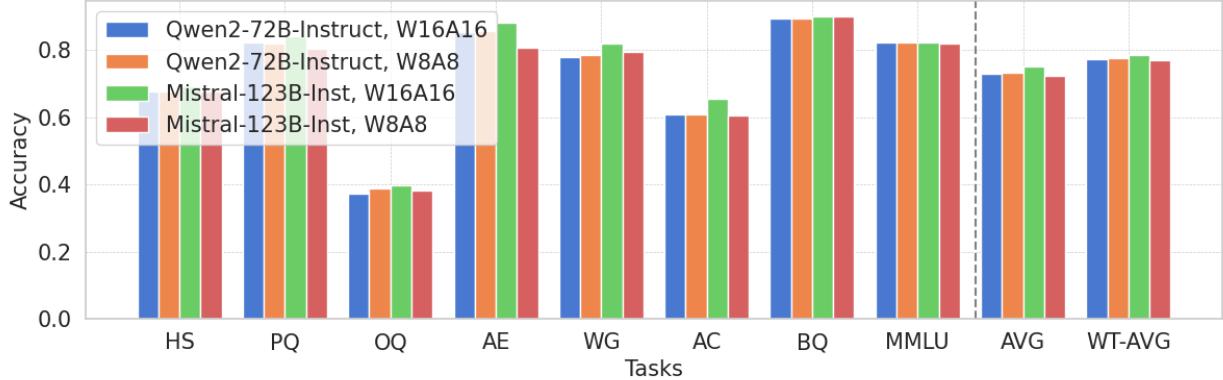
6

Figure 6: Accuracy of other Transformer-based models with W8A8 per-channel quantization
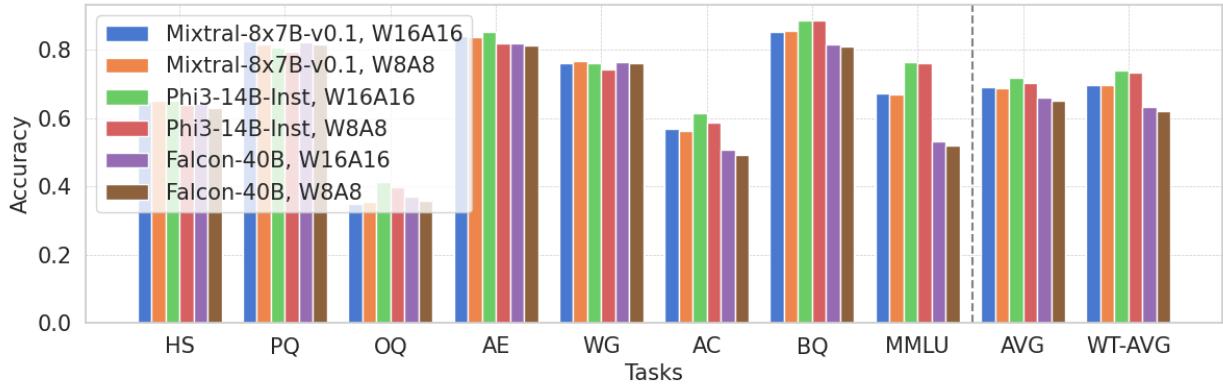


Figure 7: Accuracy of other Transformer-based models with W8A8 per-channel quantization

### 3.6.2 A quantitative exploration

To quantitatively assess the robustness of weights under W8 quantization, we employ two key metrics. The first is the maximum absolute value (max_abs) of weights for each layer in the model, which determines the scale factor and, consequently, the quantization interval. The second metric is the quantization error, measured by the root-mean-square error (RMSE) between the FP16 and 8-bit weights. These metrics are positively correlated, as larger quantization intervals typically result in greater errors.

Figure 11, 12, 13, 14 illustrate the max_abs and quantization error for four LLaMA models: LLaMA3-70B, LLaMA2-70B, LLaMA3-8B, and the recently released LLaMA3.1-405B. The x-axis represents the total number of layers, excluding the initial embedding layer and the final output head. Each Transformer block in a LLaMA model comprises seven matrices (Q, K, V, O, up, gate, down), with 80, 80, 32, and 126 blocks in the four models, respectively. The left y-axis denotes the quantization error, while the right y-axis indicates the max_abs weight of each layer.

Our comprehensive analysis of weight distributions across various layers and different LLM models yields the following conclusions:

- In the initial layers, the LLaMA3-70B model exhibits weight outliers with max_abs values that surpass those of its latter layers and all layers in other models by orders of magnitude. See the subfigure in Figure 11 for a zoom-in view of the first 98 layers (14 blocks). This results in significantly larger quantization intervals, leading to quantization errors that are 1-2 orders of magnitude greater than in other models or layers.

- In the latter layers, the LLaMA3-70B model demonstrates behavior similar to other models in terms of quantization errors and max_abs of weights. This observation suggests that the sensitivity of the LLaMA3-70B model series to W8 quantization may primarily stem from the initial layers. We will provide further confirmation of this hypothesis in the subsequent section.
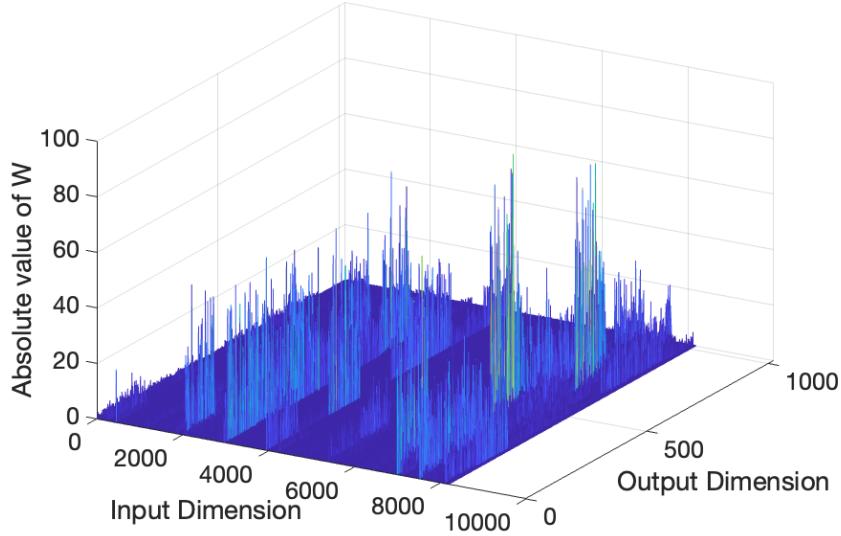
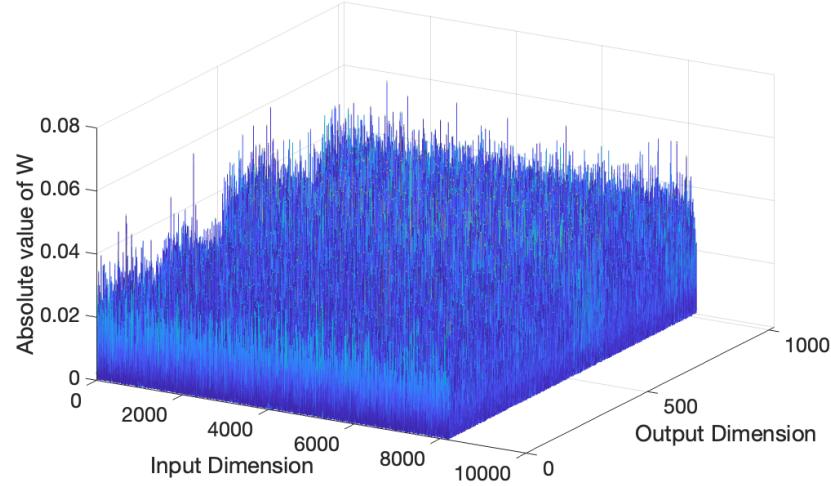Figure 8: abs_weights of the V matrix in the first Transformer block of LLaMA3.1-70B



Figure 9: abs_weights of the V matrix in the first Transformer block of LLaMA2-70B

## 4   A Mixed Grouping Strategy

Based on our analysis of Figure 11 , we have empirically identified that in the LLaMA3-70B model, the Q, K, V, Up, and Gate matrices of Block 0, 1, and 3 exhibit exceptionally high quantization errors and maximum absolute weight values. To mitigate these quantization errors, we propose implementing a finer grouping granularity specifically for these 15 layers.

It is noteworthy that the layers requiring per-group quantization, as opposed to per-channel quantization, constitute only 2.68% (15/560) of the total layers in the model. While many quantization studies employ group sizes of 128 or smaller (e.g., GPTQ, AWQ), our investigation reveals that a group size of 1024 is sufficient to prevent accuracy degradation while maintaining a larger group size for improved hardware efficiency.

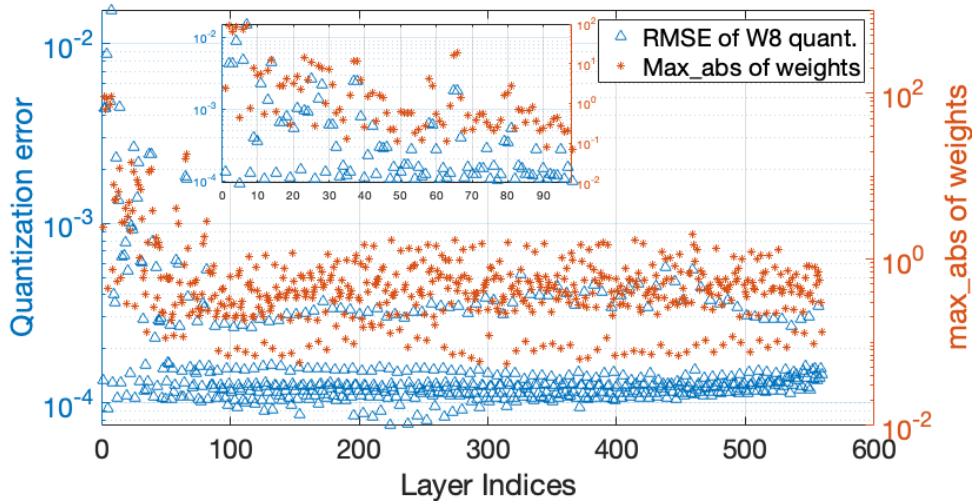Figure 10: abs_weights of the V matrix in the first Transformer block of LLaMA3-8B



Figure 11: Quantization error and max_abs of weights in each layer in LLaMA3-70B. LLaMA3-70B-Instruct, LLaMA3.1-70B, and LLaMA3.1-70B-Instruct have similar patterns.

## 4.1 Main experimental results on mixed grouping

Figure 15 and Figure 16 presents a comparative analysis of the accuracy achieved by our proposed method against that of FP16 models. The presented models includes several LLaMA3-70 and its variants and derivatives in the Open LLM Leaderboard. The results convincingly demonstrate that our mixed quantization approach—combining 2.68% per-group and 97.32% per-channel 8-bit quantization—can match the accuracy of the LLaMA3-70B model with FP16 precision.

This improvement can be attributed to the finer granularity of the quantization process in these critical layers. By reducing the group size from the entire channel to 1024 elements, we effectively create multiple quantization groups within each channel. This approach allows for more precise representation of weight values, particularly in the presence of outliers that previously dominated the quantization scale factors.
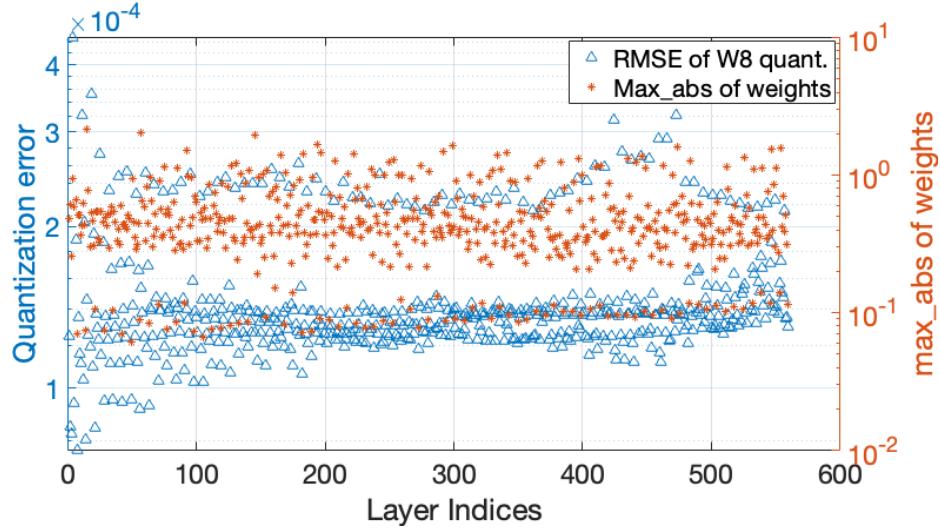
9

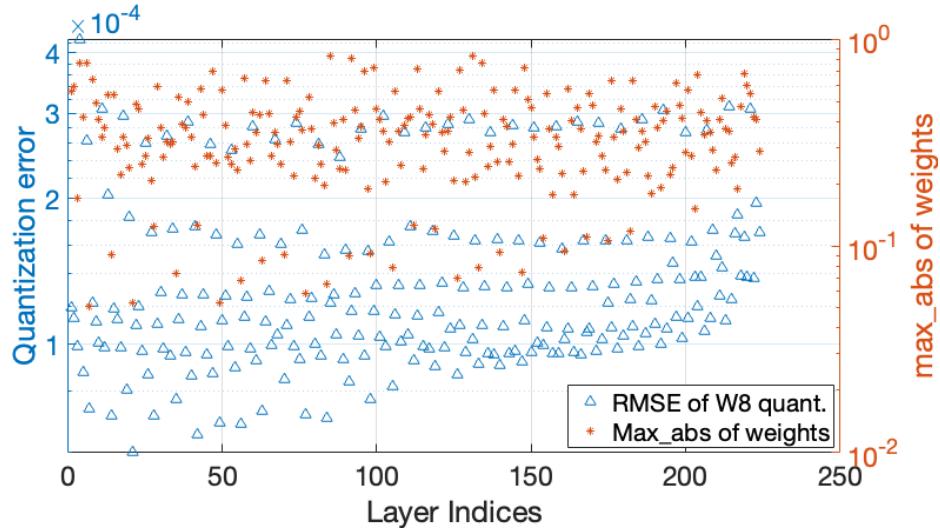Figure 12: Quantization error and max_abs of weights in each layer in LLaMA2-70B



Figure 13: Quantization error and max_abs of weights in each layer in LLaMA3-8B

In summary, this hybrid quantization strategy effectively addresses the unique challenges posed by the LLaMA3-70B model's weight distribution, particularly in its initial layers, while minimizing the hardware complexity overhead associated with per-group quantization across the entire model.

## 4.2 Ablation on the group size

We examine the influence of varying group sizes on the accuracy of W8A8 quantization for the LLaMA3-70B model. Figure 17 illustrates this comparative analysis. Our results reveal that certain datasets, notably AE, AC, and OQ, exhibit particular sensitivity to increased group size applied exclusively to a total of 15 layers within Blocks 0, 1, and 3. From a holistic perspective, per-group quantization employing sizes of 512 or 1024 emerges as a good compromise between model accuracy and hardware efficiency.
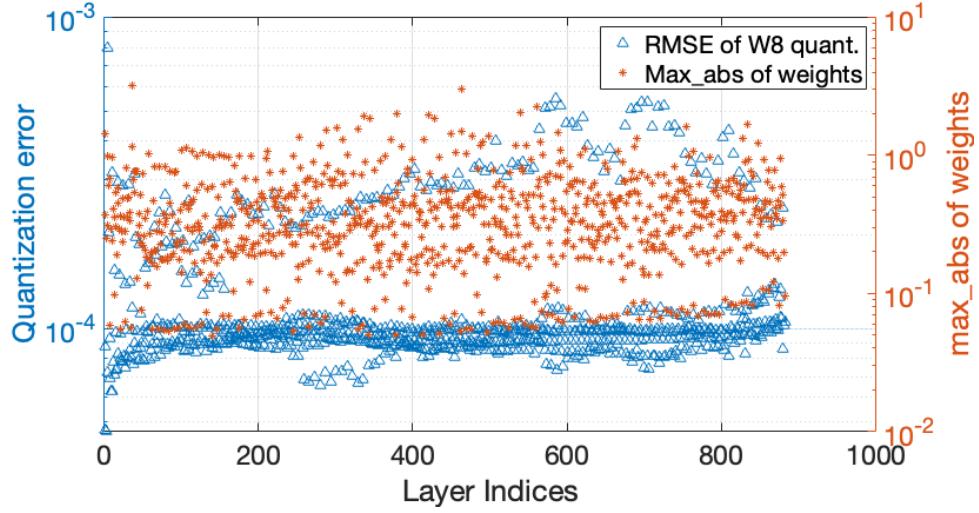
10

Figure 14: Quantization error and max_abs of weights in each layer in LLaMA3.1-405B
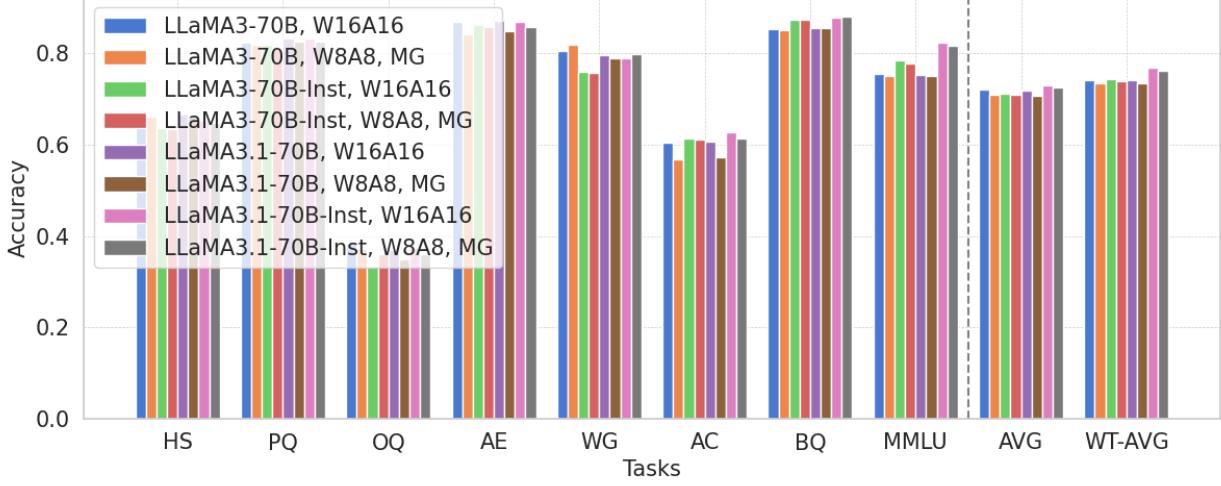


Figure 15: Accuracy comparison of LLaMA3-70 model series with mixed grouping

## 5    A Bi-smoothing Strategy

In Xiao et al. [2024], it was proposed to mitigate activation outliers by reducing the maximum absolute values of activations while amplifying the maximum absolute values of weights. This approach assumes that activation outliers are concentrated in certain channels and are significantly more severe than weight outliers. However, as demonstrated in previous sections, this assumption no longer holds for the LLaMA3-70B and LLaMA3.1-70B model series, including their instruction-tuned versions. In this section, we introduce bi-directional smoothing algorithms that balance the magnitudes of both weights and activations.

### 5.1    Equivalence in bi-directional smoothing

Let $W \in \mathbb{R}^{M \times N}$ be the weights and $A \in \mathbb{R}^{B \times L \times N}$ be the input activation. In the representation of a Transformer linear layer, $M$ is output dimension, $N$ is the hidden dimension, $B$ is the batch size and $L$ is the sequence length. The output of the linear layer $Y \in \mathbb{R}^{B \times L \times M}$ is calculated by

$$Y[i,:,:] = A[i,:,:]W^T, \forall 0 \le i \le B - 1,$$

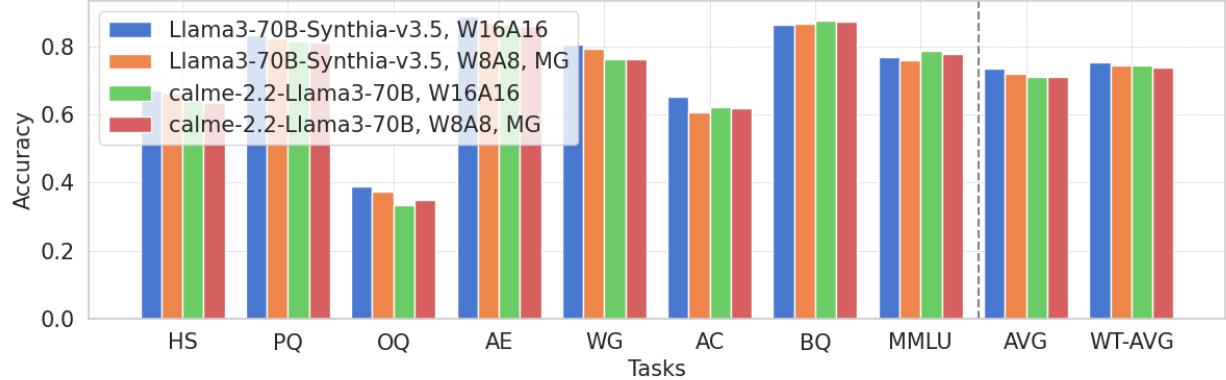where we use Python tensor representation for the 3D tensor $A$, $Y$ and 2D tensor $W$.

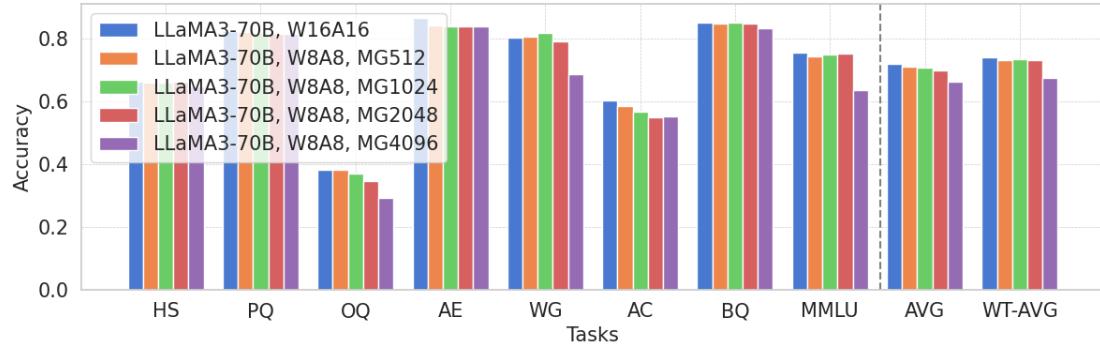Figure 16: Accuracy comparison of LLaMA3-70 derivative models with mixed grouping



Figure 17: Ablation study on the mixed group size for Q, K, V, Up, and Gate of Block 0, 1, 3 in LLaMA3-70B

There is an equivalence when multiplying $A[i, :, :]$ and $W$.

**Theorem 5.1.** *Let $S \in \mathbb{R}^N$ be a non-zero vector. We denote $W_s$ and $A_s$ the column-wise (last dimension) smoothed tensors, respectively, such that*

$$W_s[j, k] = W[j, k]S[k], \forall 0 \leq j \leq M, 0 \leq k \leq N - 1,$$

*and*

$$A_s[i, j, k] = A[i, j, k]/S[k], \forall 0 \leq i \leq B - 1, 0 \leq j \leq L - 1, 0 \leq k \leq N - 1,$$

*Then the following equation holds,*

$$A[i, :, :]W^T = A_s[i, :, :]W_s^T, \forall 0 \leq i \leq B - 1.$$

*Proof.* Let $Y[i, :, :] = A[i, :, :]W^T$ and $Y_s[i, :, :] = A_s[i, :, :]W_s^T$, then $\forall 0 \leq i \leq B - 1, 0 \leq j \leq L - 1, 0 \leq k \leq M - 1,$

$$Y_s[i, j, k] = \sum_z A_s[i, j, z]W_s[z, j]$$
$$= \sum_z A[i, j, z]/S[z] * W[z, j]S[z]$$
$$= \sum_z A[i, j, z]W[z, j]$$
$$= Y[i, j, k].$$

$\square$

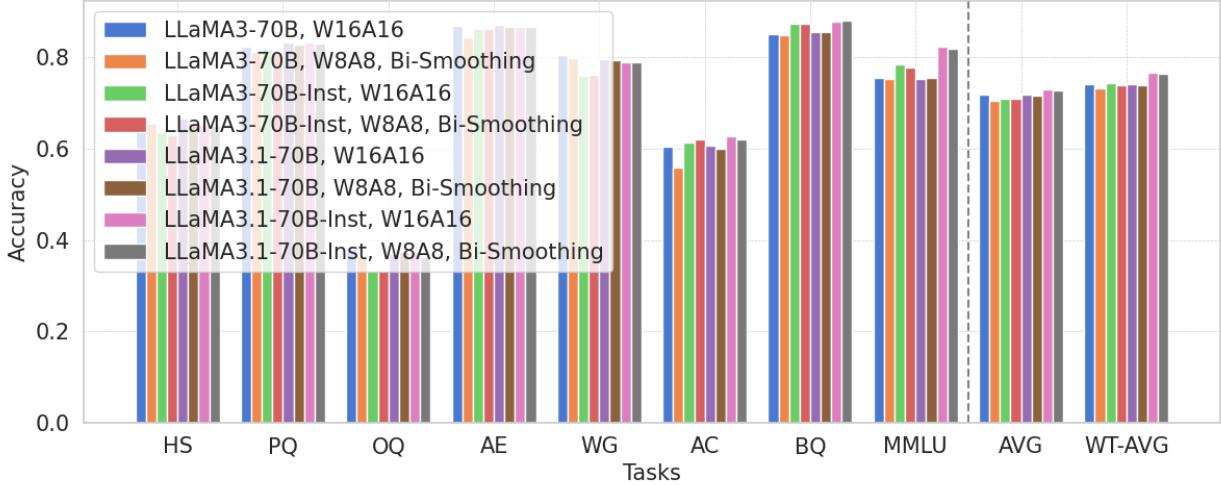We call the vector $S \in \mathbb{R}^N$ the **smooth-factor** for $A$ and $W$.

Figure 18: Accuracy comparison of LLaMA3-70B model series with bi-smoothing

## 5.2 How to select the smooth-factor?

For per-channel quantization, a quantization group consists of $N$ numbers and the max_abs in the quantization group determines the scale-factor and quantization interval. First, assume the batch size is 1, and then $A \in \mathbb{R}^{1 \times L \times N}$. Since it is desirable to minimize the max_abs of both weight and input activation, but the multiplication of these two values, i.e., $\max(|W[:,i]|)$ and $\max(|A[0,:,i]|)$, remains constant to guarantee the output of $AW^T$ is unchanged. Therefore, we choose the smooth-factor as

$$S[k] = \sqrt{\frac{\max_j(|A[0,j,k]|)}{\max_j(|W[j,k]|)}}, \forall 0 \le k \le N-1.$$

Then we have

$$\max_j(|W_s[j,k]|) = \max_j(|W[j,k]S[i]|)$$

$$= \sqrt{\max_j(|A[0,j,k]|)\max_j(|W[j,k]|)}$$

$$= \max_j(|A[0,j,k]S[k]|)$$

$$= \max_j(|A_s[0,j,k]|),$$

i.e., the max_abs after the bi-smoothing for weight and input activation is balanced.

In general when the batch size $B$ is greater than 1, then we can choose $S[k]$ as

$$S[k] = \sqrt{\frac{\text{median}_i \max_j(|A[i,j,k]|)}{\max_j(|W[j,k]|)}}, \forall 0 \le k \le N-1. \tag{1}$$

## 5.3 Main experiment results for bi-smoothing

We show the accuracy comparison of all four LLaMA3-70B series of models (LLaMA3.1-70B-Instruct, LLaMA3.1-70B, LLaMA3-70B-Instruct, and LLaMA3-70B) in the section.

Figure 18 demonstrates the effectiveness of the proposed bi-smoothing method for the entire LLaMA3-70B model series. It is observed that the accuracy of models with per-channel W8A8 quantization are as good as their FP16 counterparts, which shows significant improvement over W8A16 models in Figure 1. Note that to calculate the smooth-factor in Equation (1), we need some calibration data to estimate $\text{median}_i \max_j(|A[i,j,k]|)$. To prevent data contamination, we use only eight questions from the GPQA reasoning task for estimation. The calibration process takes only a few seconds, requiring a single inference pass over these samples. As demonstrated in the later ablation studies, the accuracy of the bi-smoothing method remains stable across varying amounts and sources of calibration data.
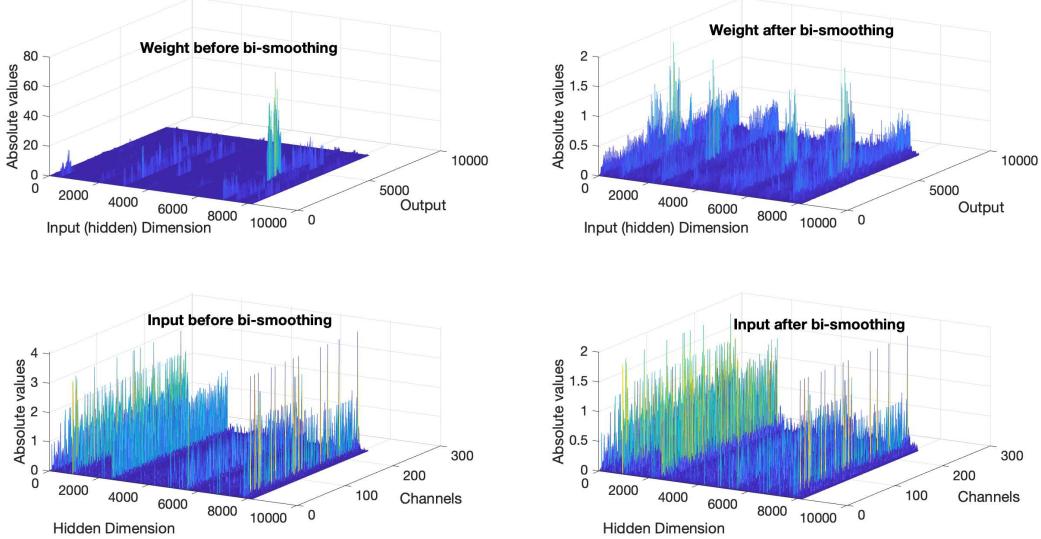
Figure 19: The plot illustrates the absolute values of weights and input activations in the LLaMA3.1-70B-Instruct model, specifically focusing on the first Q matrix in the initial Transformer block. The input activations are multiplied by the Q matrix. **Top-left**: Before bi-smoothing, large weight values (ranging from 60 to 80) form distinct "walls" parallel to the output axis. **Top-right**: After bi-smoothing, these weight "walls" remain but with significantly reduced magnitudes (below 2.0). **Bottom-left**: Input activations before bi-smoothing. **Bottom-right**: After bi-smoothing, the input activations are also suppressed, exhibiting smaller magnitudes.

### 5.3.1 Why bi-smoothing works well?

Figure 19 compares the weights and input activations of the first Q matrix in the LLaMA3.1-70B-Instruct model. The input activations are drawn from a sequence of length 243 in the PIQA dataset. Notably, problematic weight outliers, with magnitudes exceeding 60, are suppressed to below 2.0, alongside activation outliers. This smoothing process makes both weights and activations far more resilient to quantization, achieving W8A8 accuracy nearly identical to that of the W16A16 configuration (as detailed in the next section).

In general, Figure 20 illustrates the quantization error and max_abs across all layers of the LLaMA3.1-70B-Instruct model after bi-smoothing. Compared to Figure 11, depicting the unsmoothed model, there is a substantial reduction in max_abs, yielding an order-of-magnitude decrease in quantization error.

## 5.4 Ablation study on bi-smoothing

From Figure 18, we can see the best performing model is LLaMA3.1-70B-Instruct. So, we use this model for ablation study of bi-smoothing algorithms. We will show the bi-smoothing method is stable in calibration data and the choice of computing the smooth-factor.

### 5.4.1 Calibration data source

Figure 21 presents an accuracy comparison across different calibration data sources. We tested 8 samples from the GPQA dataset, 32 samples from the WG dataset, and 32 samples from the PQ dataset. The results indicate that the choice of calibration data has minimal impact on accuracy; all three datasets enable the W8A8 per-channel quantized model to achieve accuracy comparable to that of the FP16 counterpart.

### 5.4.2 Calibration data size

We evaluated the impact of calibration data volume on smooth-factor computation, comparing results from 8 samples versus a single sample from the GPQA dataset. Figure 22 demonstrates that the difference in performance is negligible; both resulting W8A8 models achieve accuracy comparable to their FP16 counterpart.
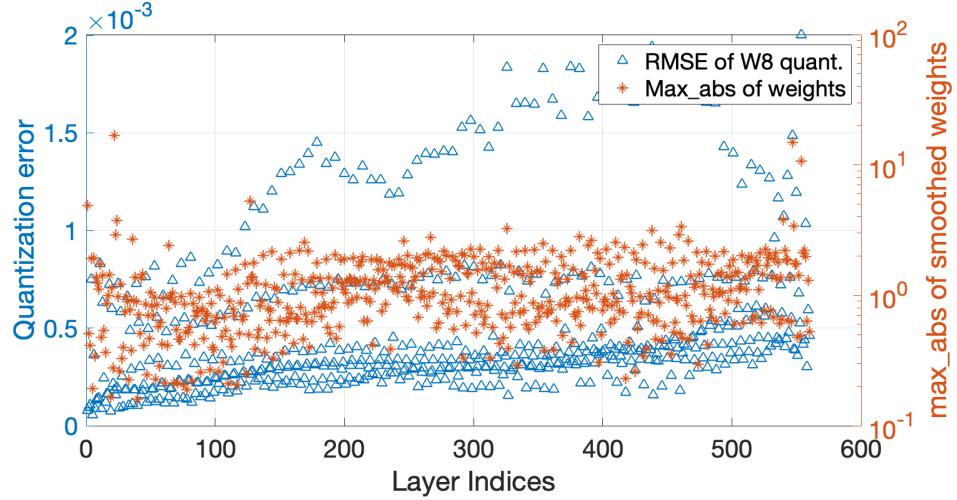
14

Figure 20: Quantization error and max_abs of weights after bi-smoothing of each layer in the LLaMA3.1-70B-Instruct model.
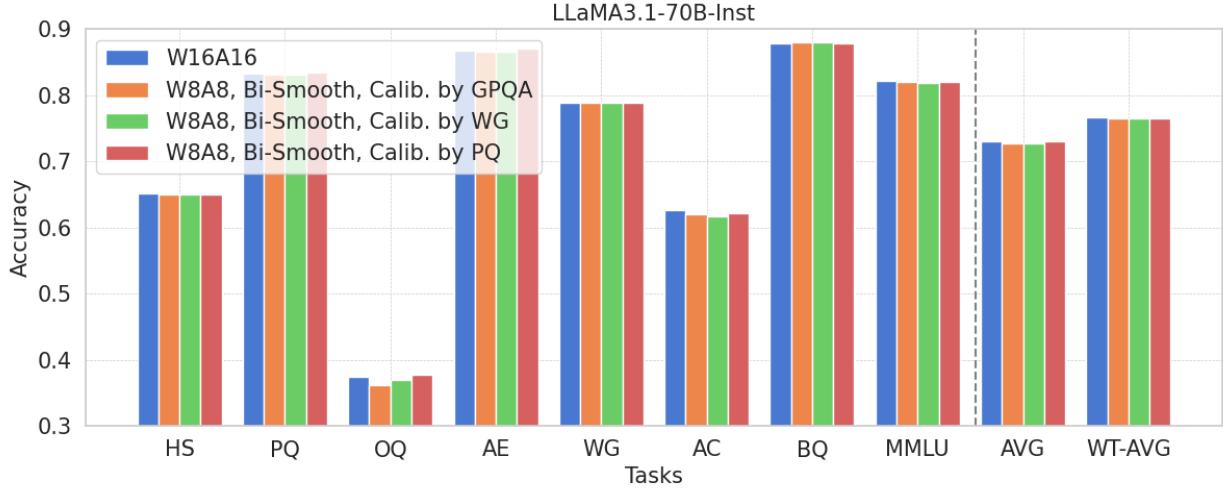


Figure 21: Different calibration data for bi-smoothing

### 5.4.3 Smooth-factor calculation

Note that in Equation (1), we use "median" to calculate the smooth-factor. A natural alternative is to replace "median" by "maximum" in the following equation.

$$S[k] = \sqrt{\frac{\max_i \max_j (|A[i,j,k]|)}{\max_j (|W[j,k]|)}}, \forall 0 \le k \le N - 1. \tag{2}$$

Figure 23 illustrates the comparison, revealing that there is negligible difference between using the "median" or "max" to calculate the smooth factor. This finding underscores the robustness of the bi-smoothing method in recovering the accuracy of W8A8 models.

### 5.4.4 Applying bi-smoothing to a subset of Transformer layers

Our previous analysis revealed that significant weight outliers are predominantly confined to the initial Transformer blocks. An ablation study on the number of blocks with bi-smoothing method was conducted, with results illustrated in Figure 24. Notably, applying bi-smoothing to only the 0th, 1st, and 3rd Transformer blocks yields W8A8 accuracy
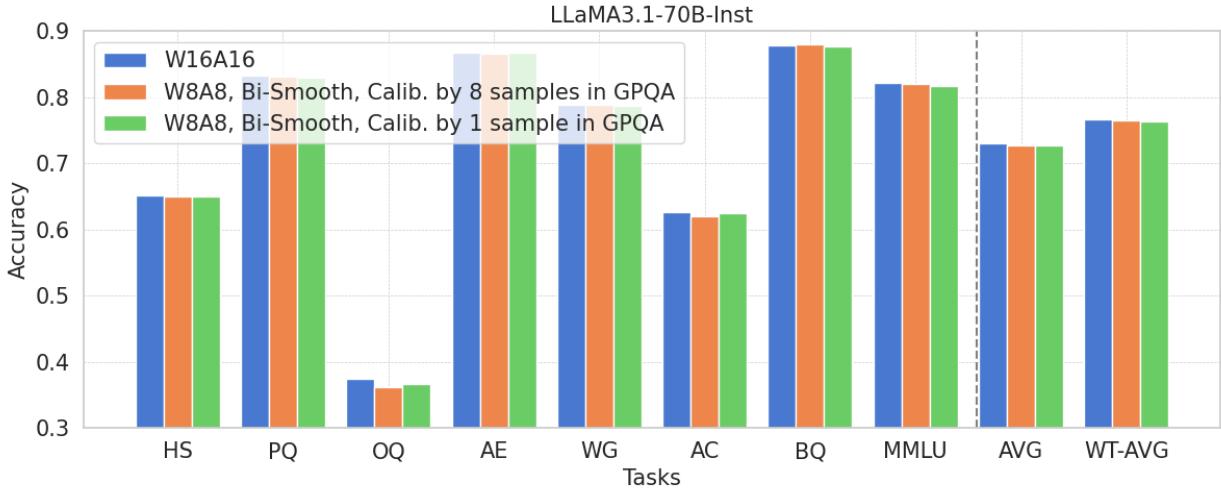
15

Figure 22: Comparison of calibration data using 8 samples and 1 sample from GPQA.
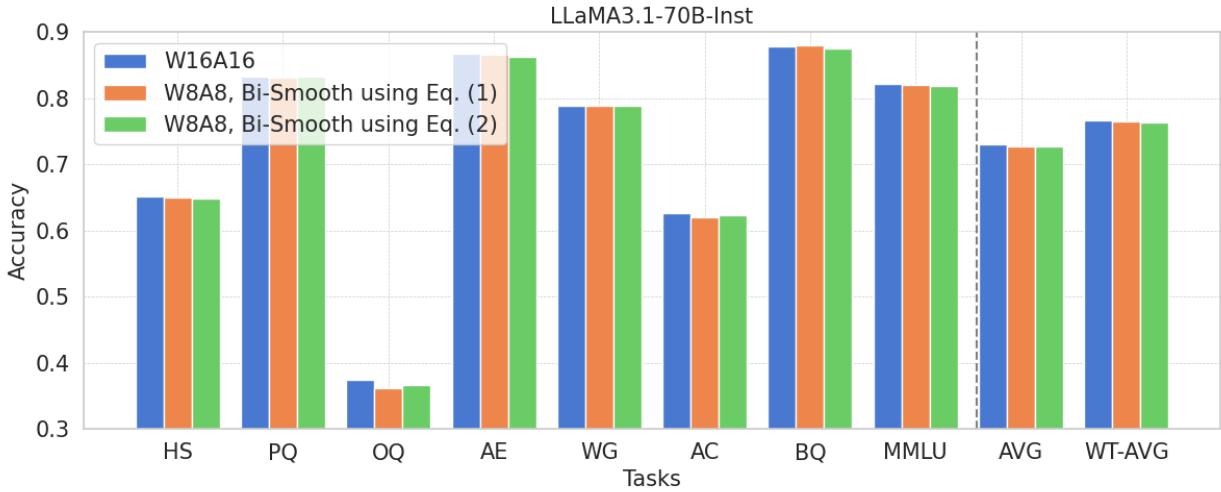


Figure 23: Comparison of using Eq. 1 and Eq. 2 to calculate the smooth-factor

nearly equivalent to that of the full-precision (FP16) model. This finding corroborates our hypothesis that the initial blocks are the primary source of the LLaMA3-70B model series' susceptibility to W8A8 quantization.

# 6 Discussions

According to Meta's technical report Dubey et al. [2024], LLaMA2-70B and LLaMA3-70B share similar model architectures, while LLaMA3-8B/70B and LLaMA3.1-8B/70B/405B are reported to have similar training data and strategies. However, our analysis reveals significantly different weight distribution characteristics in LLaMA3/3.1-70B compared to other LLaMA models. This unexpected divergence in weight patterns, particularly the presence of numerous large outliers exclusively in LLaMA3/3.1-70B models, presents an intriguing peculiarity.

The origins of this phenomenon remain unclear, given that models sharing either architectural design or training methodologies do not exhibit similar characteristics. Unraveling the underlying mechanisms responsible for this unique weight distribution could potentially offer valuable insights into the learning processes of Large Language Models (LLMs), especially in the context of reasoning tasks.
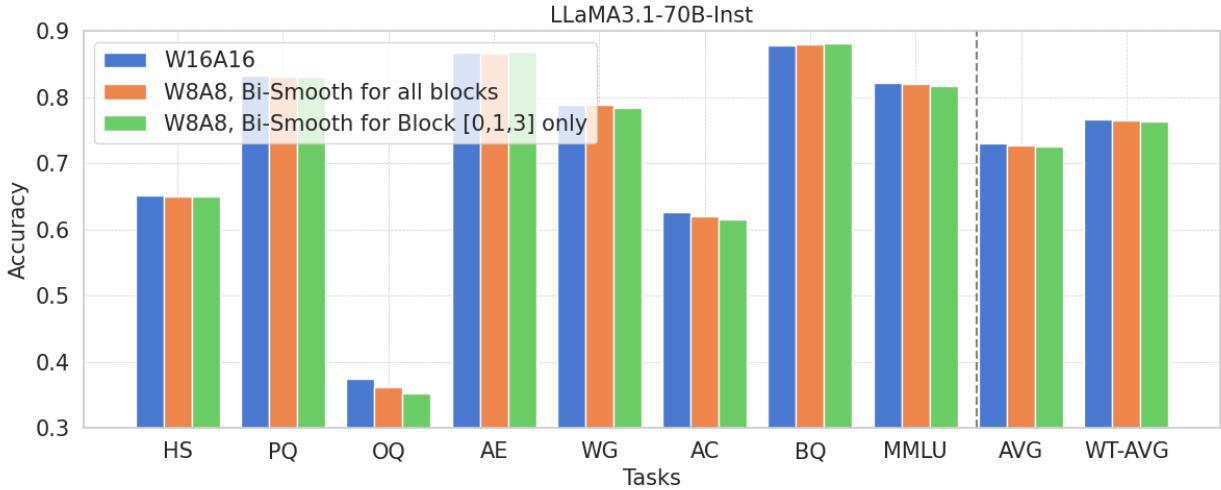
16

Figure 24: Applying bi-smoothing methods to a subset of Transformer blocks

## 7 Related Works

Quantization is useful in reducing model size and potentially increases model inference speed Han et al. [2016], Jacob et al. [2017], Nagel et al. [2019], Bondarenko et al. [2021], Wang et al. [2019], Bengio et al. [2013]. There two main categories of quantization, post-training quantization (PTQ) and quantization-aware training (QAT).

### 7.1 Post-training quantization

PTQ does not use any back-propagation to finetune weights. For **weight-only quantization**, LLM.int8() Dettmers et al. [2022a] uses a mixed precision with 16-bit MAC for outliers. GPTQ Frantar et al. [2023] and AWQ Lin et al. [2024a] quantize the per-group weight-only tensor using as small as 3-4 bits. SqueezeLLM Kim et al. [2024] uses non-uniform quantization and store outliers efficiently. For **weight-activation quantization** that enables faster inference Dettmers et al. [2022b], Wei et al. [2023] SmoothQuant Xiao et al. [2024] transfers the outliers from activation to weights to enable per-tensor W8A8. QServe Lin et al. [2024b] proposes W4A8K4 to serve LLM models. Some other studies enable even smaller number of bits Shao et al. [2024], Zhao et al. [2024], Ashkboos et al. [2024]. Notably, our observations on the weights and activations of the LLaMA3-70B model series contrast with those of SmoothQuant. The proposed bi-smoothing method extends this approach by addressing both weight and activation outliers. Compared to other PTQ methods, we aim to maintain per-channel instead of per-group quantization for maximal hardware efficiency.

### 7.2 Quantization aware training

Quantization-aware training (QAT) requires both calibration data and back-propagation to finetune the quantize weights Bengio et al. [2013], Gholami et al. [2021]. While full finetuning Taori et al. [2023] can recover the accuracy to the maximum extent, Q-LoRA Dettmers et al. [2023] reduces the memory requirement significantly.

## 8 Conclusions

Our investigation reveals a distinctive vulnerability in the LLaMA3-70B model series, including LLaMA3.1-70B, to per-channel quantization. We attribute this susceptibility to significant weight outliers in the initial layers. The contrast in behavior between the LLaMA3-70B and other LLaMA models (LLaMA2, LLaMA3/3.1-8B, LLaMA3.1-405B), as well as other tested LLMs with different architectures, warrants further exploration of the underlying training strategies. To address this issue, we propose a mixed-grouping method that substantially enhances the accuracy of LLaMA3-70B from 45.4% to 73.4%, effectively matching the performance of its FP16 counterpart.

## 9    Appendix

### 9.1    Models

Table 1 shows all the tested models downloaded from Huggingface repository.

| Model name | Huggingface URL |
|---|---|
| LLaMA3-70B | https://huggingface.co/meta-llama/Meta-Llama-3-70B |
| LLaMA3-70B-Instruct | https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct |
| LLaMA3.1-70B | https://huggingface.co/meta-llama/Meta-Llama-3.1-70B |
| LLaMA3.1-70B-Instruct | https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct |
| Llama3-70B-Synthia | https://huggingface.co/migtissera/Llama-3-70B-Synthia-v3.5 |
| calme-2.2-llama3-70b | https://huggingface.co/MaziyarPanahi/calme-2.2-llama3-70b |
| LLaMA3-8B | https://huggingface.co/meta-llama/Meta-Llama-3-8B |
| LLaMA3.2-1B-Instruct | https://huggingface.co/meta-llama/Llama-3.2-1B-Instruct |
| LLaMA3.2-3B-Instruct | https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct |
| LLaMA2-70B | https://huggingface.co/meta-llama/Llama-2-70b-hf |
| LLaMA2-70B-chat | https://huggingface.co/meta-llama/Llama-2-70b-chat-hf |
| Qwen2-72B | https://huggingface.co/Qwen/Qwen2-72B |
| Mixtral-8x7B | https://huggingface.co/mistralai/Mixtral-8x7B-v0.1 |
| Phi3-14B-Instruct | https://huggingface.co/microsoft/Phi-3-medium-128k-instruct |
| Mistral-Large-Instruct-123B | https://huggingface.co/mistralai/Mistral-Large-Instruct-2407 |
| Falcon-40B | https://huggingface.co/tiiuae/falcon-40b |

Table 1: Tested models

### 9.2    Weights of LLaMA3-70B

#### 9.2.1    The first block has more large outliers

We show the distribution of weights in the first block in LLaMA3-70B. Figures 25, 26, 27, 28, 29, 30, 31 shows seven weight matrices in the 0-th Transformer block. The z-axis shows the magnitudes of outliers in these matrices. We can see that the Q, K, V, Up, and Gate matrices have larger outliers and they form a "wall" in some input indices. The other two weights, namely O and Down, do not show such a pattern.

#### 9.2.2    The latter block has less large outliers

Figure 32, 33, 34, 35, 36, and 37 show the Q and Up matrices in the 10th, 40th, and 79th (last) Transformer block in LLaMA3-70B. The magnitude of the weight outliers is significantly smaller than that in the first block, demonstrating more robustness to W8 per-channel quantization.

### 9.3    Weights of Other LLM models

Figure 38, 39, 40, 41 show the weights with maximum quantization error under W8 per-channel quantization in four LLMs, that is, LLaMA2-70B, LLaMA3-8B, LLaMA3.1-405B, and Qwen2-72B. It is worth noting that the maximum quantization error occurs all at the same layer, which is the "K" matrix in the first Transformer block. We can see that the weight outliers in the layer with maximum quantization error is less than 1.0, which is far smaller than that of LLaMA3-70B models, which is usually greater than 90. The difference of 2 orders of magnitude causes the quantization error of the LLaMA3-70Bs to be much more than the rest of the models.
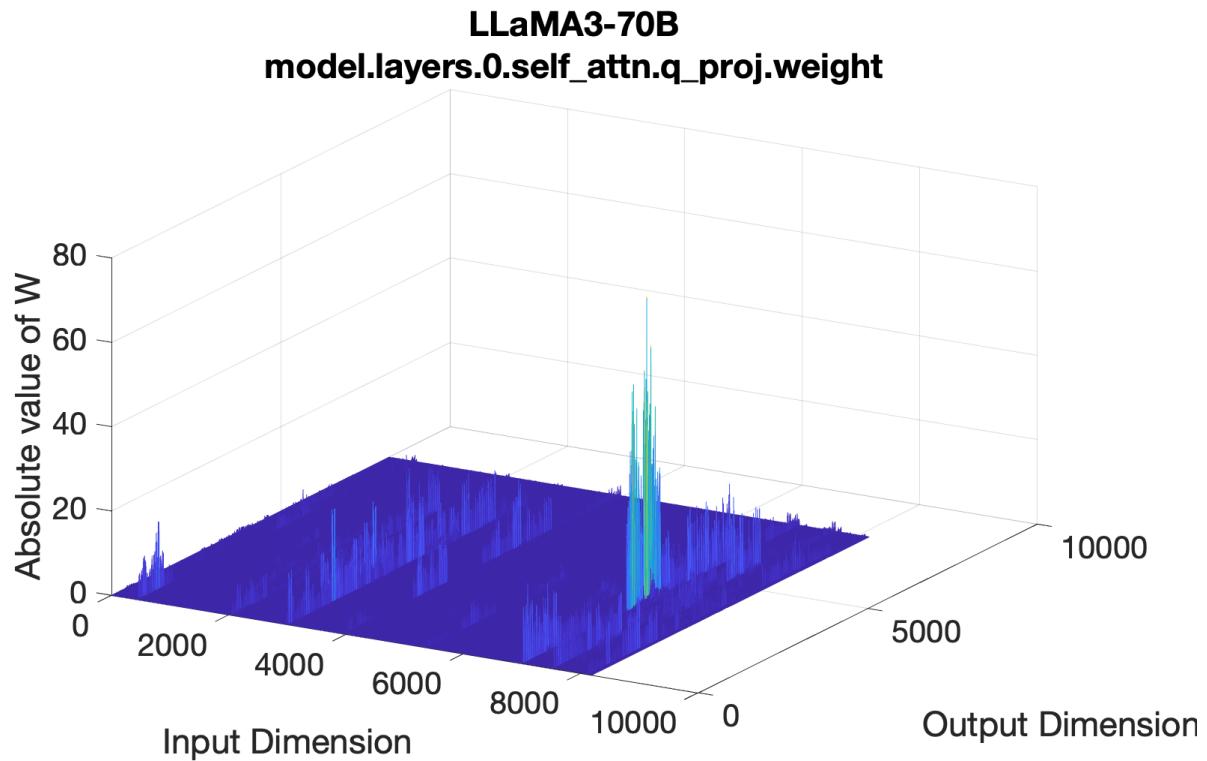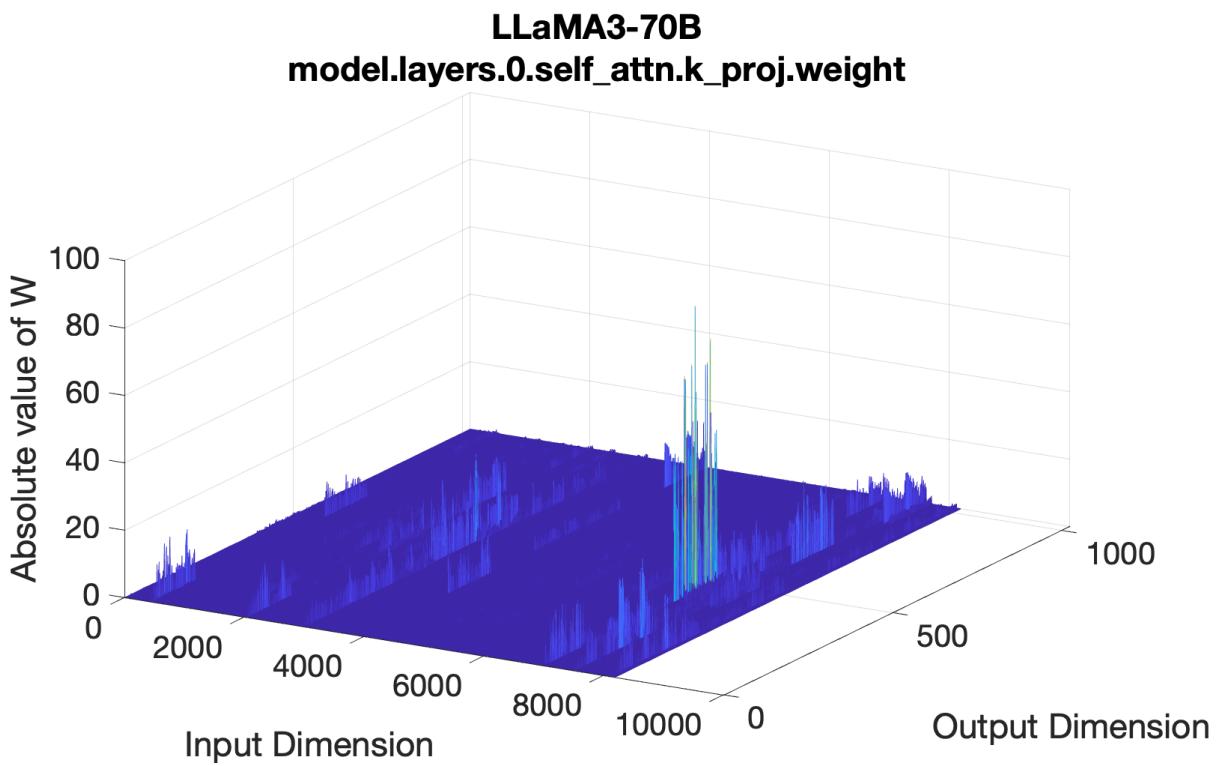
Figure 25: Weights in LLaMA3-70B



Figure 26: Weights in LLaMA3-70B

Figure 27: Weights in LLaMA3-70B



Figure 28: Weights in LLaMA3-70B

**LLaMA3-70B
model.layers.0.mlp.up_proj.weight**



Figure 29: Weights in LLaMA3-70B

**LLaMA3-70B
model.layers.0.mlp.gate_proj.we**
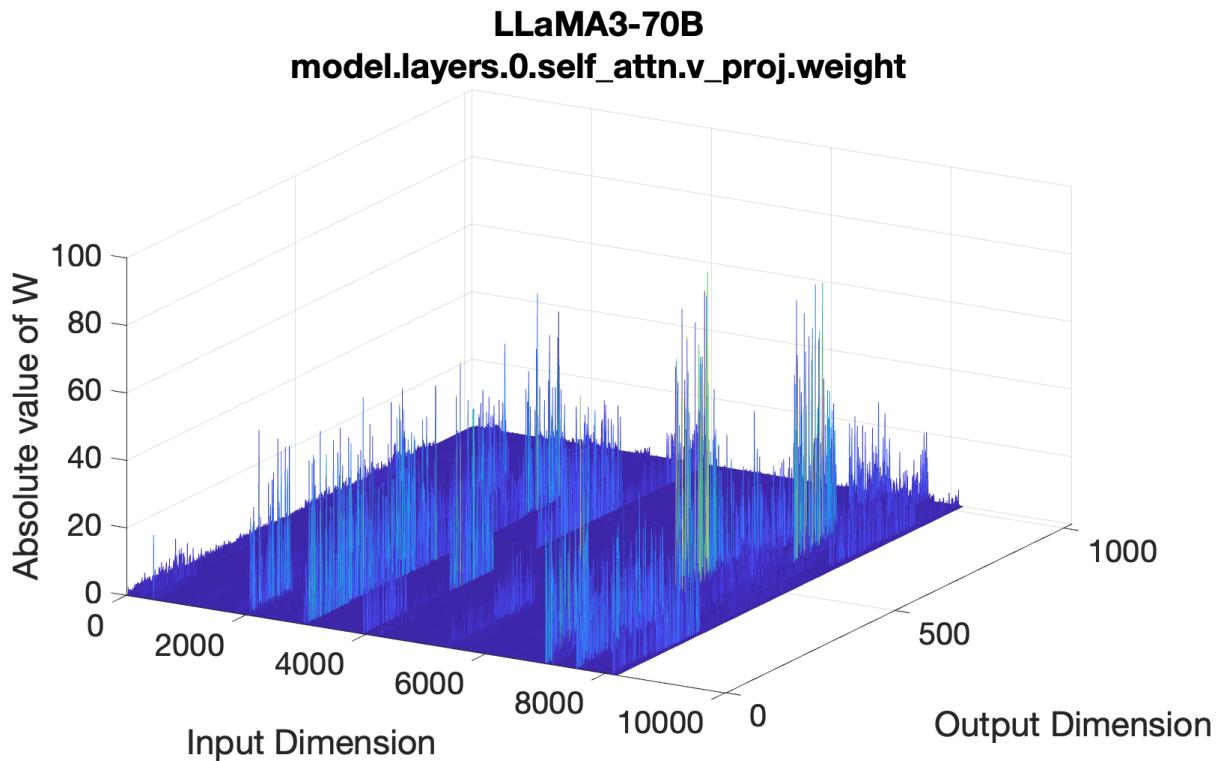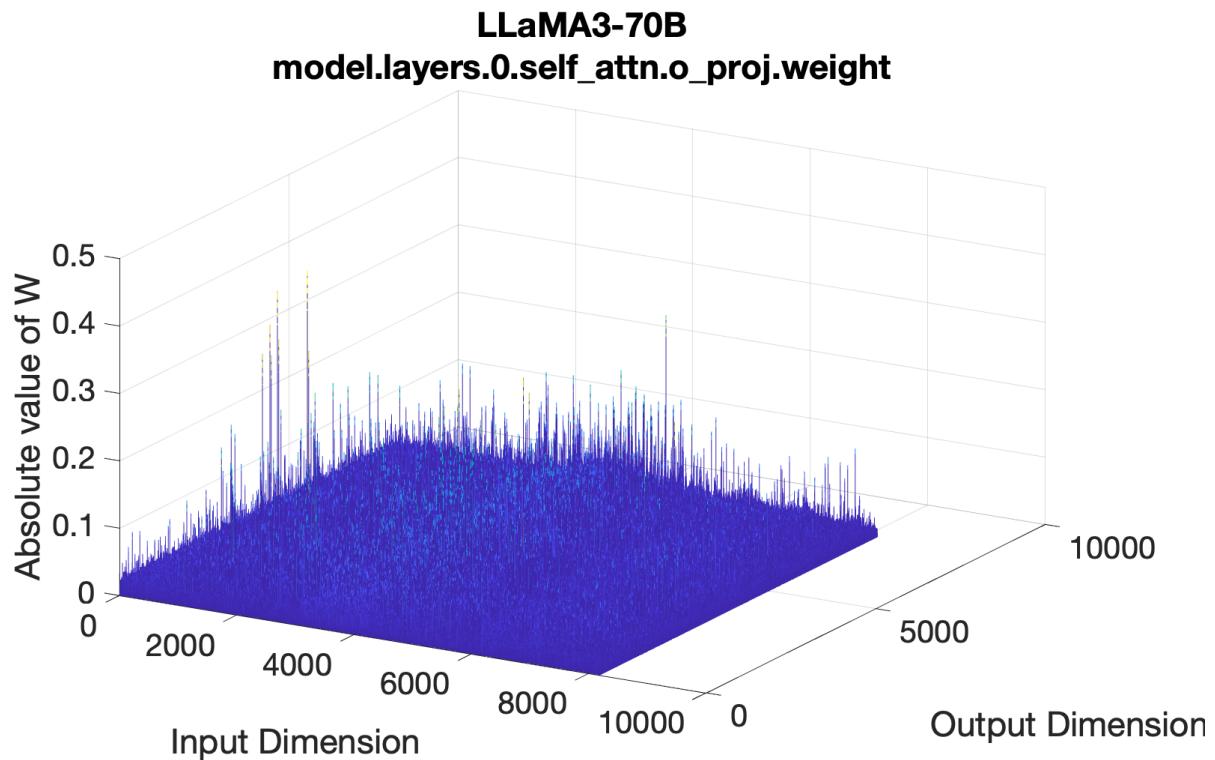


Figure 30: Weights in LLaMA3-70B
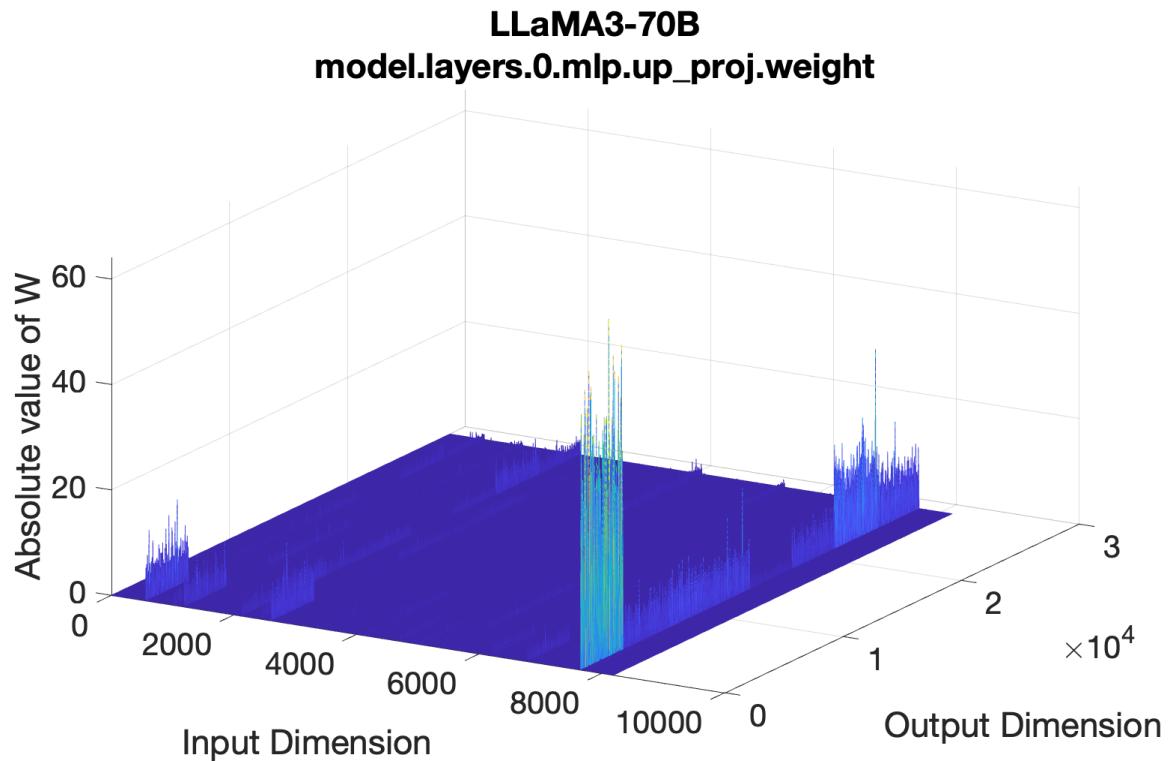
21

Figure 31: Weights in LLaMA3-70B



Figure 32: Weights in LLaMA3-70B

Figure 33: Weights in LLaMA3-70B



Figure 34: Weights in LLaMA3-70B

**LLaMA3-70B**
**model.layers.10.mlp.up_proj.weight**



Figure 35: Weights in LLaMA3-70B

**LLaMA3-70B**
**model.layers.40.mlp.up_proj.weight**



Figure 36: Weights in LLaMA3-70B

**LLaMA3-70B**
**model.layers.79.mlp.up_proj.weight**



Figure 37: Weights in LLaMA3-70B

**LLaMA2-70B**
**model.layers.0.self_attn.k_proj.weight**



Figure 38: Weights with maximum quantization error of W8 in LLaMA2-70B, where max_abs_weight=0.95.

**LLaMA3-8B**
**model.layers.0.self_attn.k_proj.weight**



Figure 39: Weights with maximum quantization error of W8 in LLaMA3-8B, where max_abs_weight=0.77.

**LLaMA3.1-405B**
**model.layers.0.self_attn.k_proj.weight**



Figure 40: Weights with maximum quantization error of W8 in LLaMA3.1-405B, , where max_abs_weight=0.98.

Figure 41: Weights with maximum quantization error of W8 in Qwen2-72B, , where max_abs_weight=0.41.

# References

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL `https://arxiv.org/abs/1706.03762`.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL `https://arxiv.org/abs/2005.14165`.
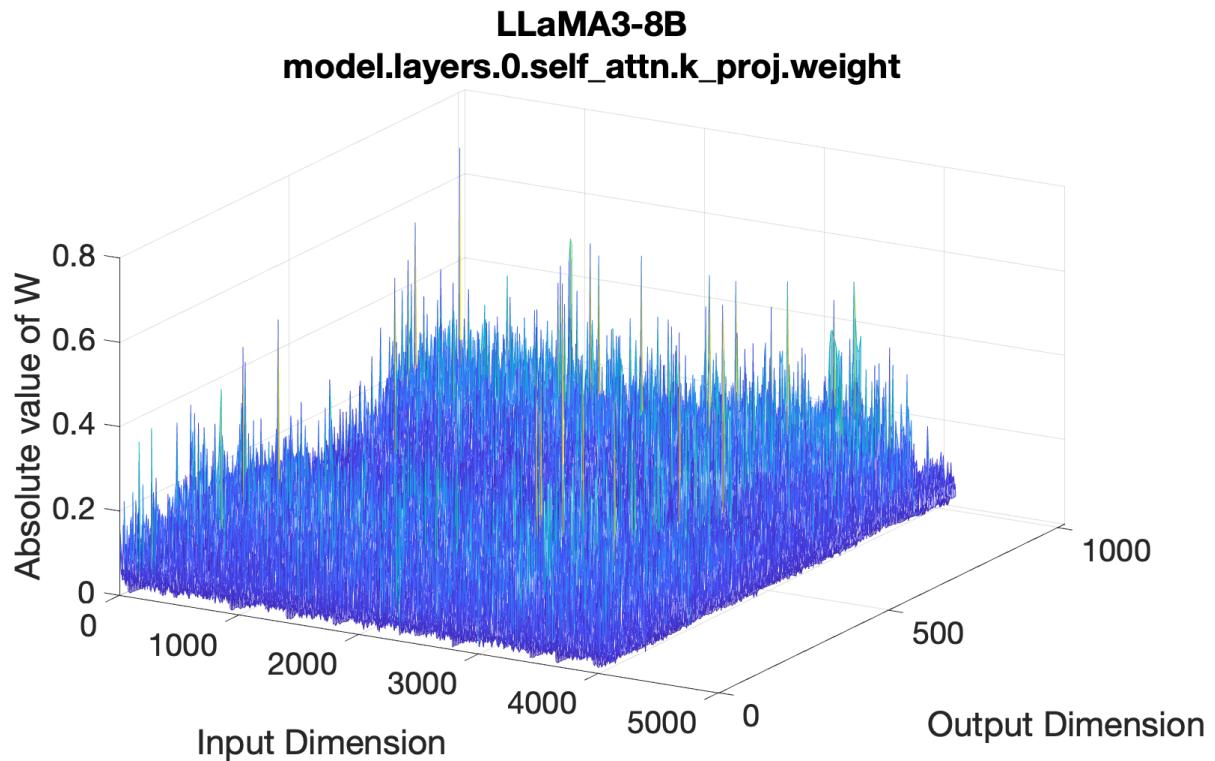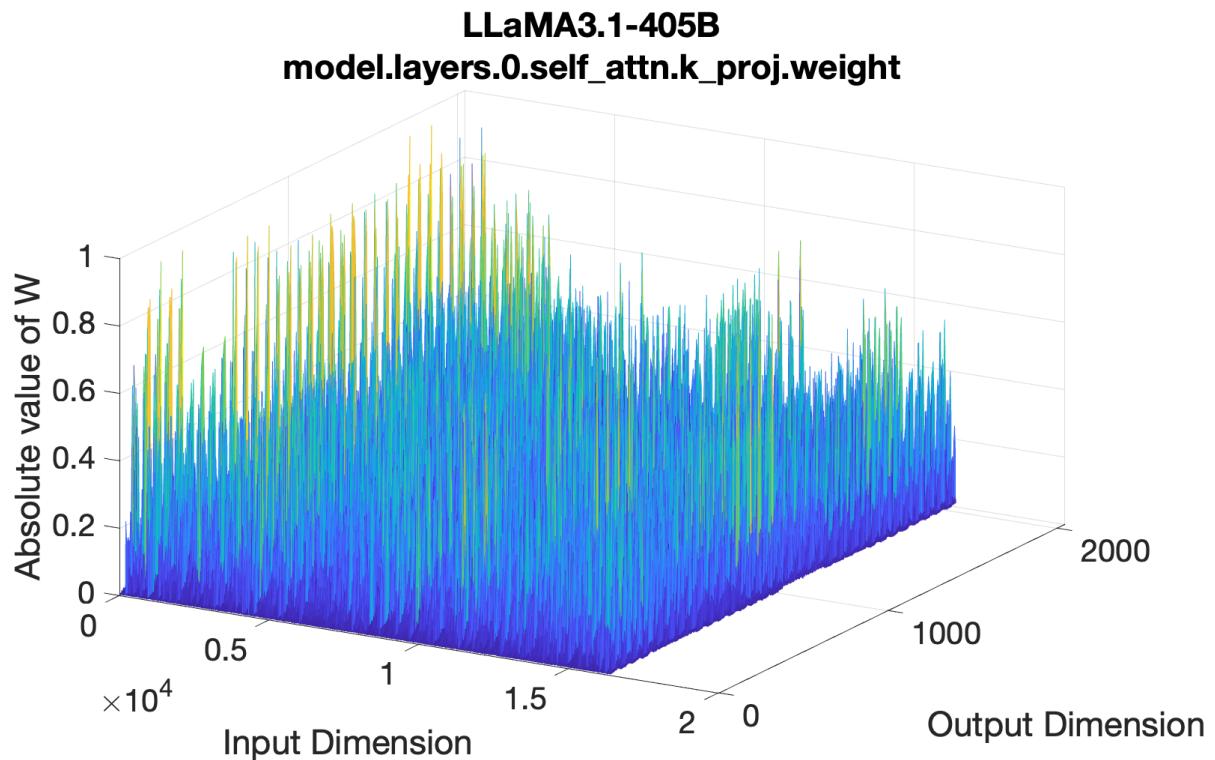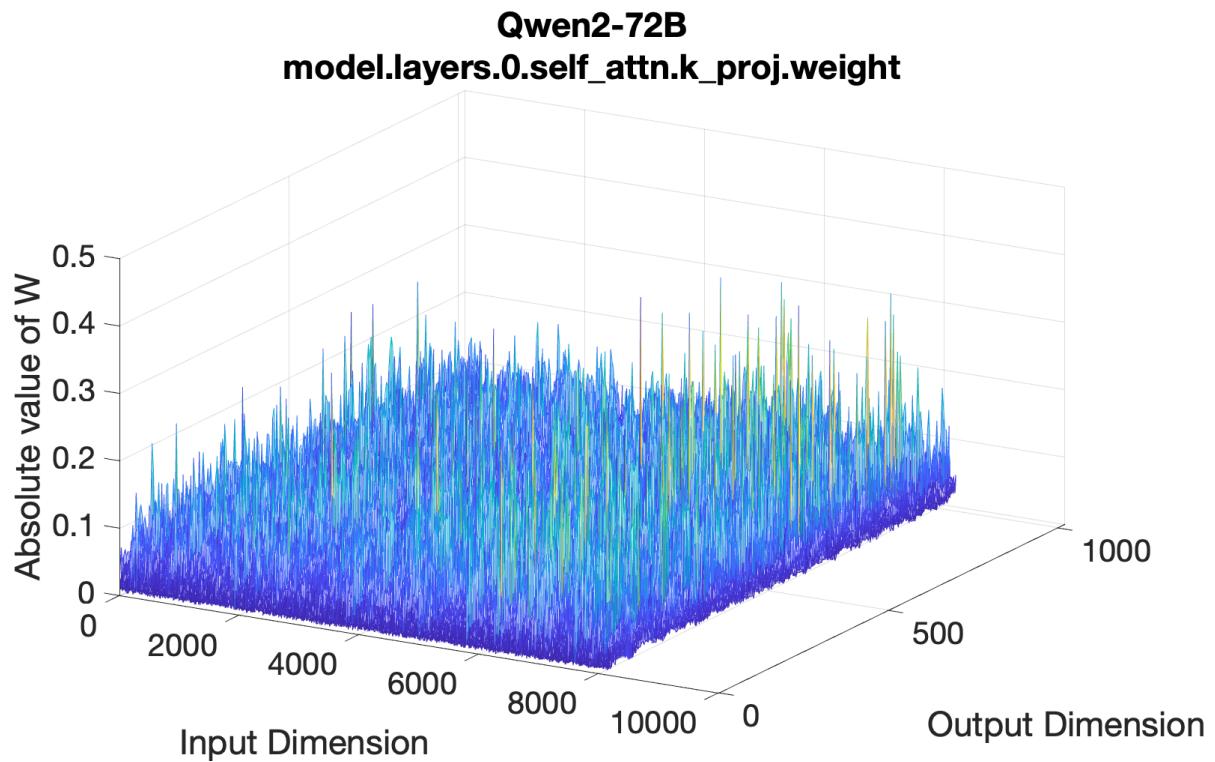
Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. URL `https://arxiv.org/abs/1810.04805`.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019. URL `https://arxiv.org/abs/1907.11692`.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. URL `https://arxiv.org/abs/2302.13971`.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. URL `https://arxiv.org/abs/2001.08361`.

Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W. Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference, 2021. URL `https://arxiv.org/abs/2103.13630`.

Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8bert: Quantized 8bit bert. In *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing - NeurIPS Edition (EMC2-NIPS)*. IEEE, December 2019. doi:10.1109/emc2-nips53020.2019.00016. URL `http://dx.doi.org/10.1109/EMC2-NIPS53020.2019.00016`.

Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Huaijie Wang, Lingxiao Ma, Fan Yang, Ruiping Wang, Yi Wu, and Furu Wei. Bitnet: Scaling 1-bit transformers for large language models, 2023. URL `https://arxiv.org/abs/2310.11453`.

HuggingFace. Open llm leaderboard. `https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard`, 2024.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,

Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian,

Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024. URL `https://arxiv.org/abs/2401.04088`.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024. URL `https://arxiv.org/abs/2407.10671`.

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Qin Cai, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Yen-Chun Chen, Yi-Ling Chen, Parul Chopra, Xiyang Dai, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Victor Fragoso, Dan Iter, Mei Gao, Min Gao, Jianfeng Gao, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Ce Liu, Mengchen Liu, Weishung Liu, Eric Lin, Zeqi Lin, Chong Luo, Piyush Madan, Matt Mazzola, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Xin Wang, Lijuan Wang, Chunyu Wang, Yu Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Haiping Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Sonali Yadav, Fan Yang, Jianwei Yang, Ziyi Yang, Yifan Yang, Donghan Yu, Lu Yuan, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL `https://arxiv.org/abs/2404.14219`.

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. The falcon series of open language models, 2023. URL `https://arxiv.org/abs/2311.16867`.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence?, 2019. URL `https://arxiv.org/abs/1905.07830`.

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language, 2019. URL `https://arxiv.org/abs/1911.11641`.

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019. URL `https://arxiv.org/abs/1907.10641`.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions, 2019. URL `https://arxiv.org/abs/1905.10044`.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021a.

Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021b.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 12 2023. URL `https://zenodo.org/records/10256836`.

Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, 2016. URL `https://arxiv.org/abs/1510.00149`.

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, 2017. URL https://arxiv.org/abs/1712.05877.

Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction, 2019. URL https://arxiv.org/abs/1906.04721.

Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Understanding and overcoming the challenges of efficient transformer quantization, 2021. URL https://arxiv.org/abs/2109.12948.

Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. Haq: Hardware-aware automated quantization with mixed precision, 2019. URL https://arxiv.org/abs/1811.08886.

Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation, 2013. URL https://arxiv.org/abs/1308.3432.

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale, 2022a. URL https://arxiv.org/abs/2208.07339.

Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers, 2023. URL https://arxiv.org/abs/2210.17323.

Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration. In *MLSys*, 2024a.

Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W. Mahoney, and Kurt Keutzer. Squeezellm: Dense-and-sparse quantization, 2024. URL https://arxiv.org/abs/2306.07629.

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. GPT3.int8(): 8-bit matrix multiplication for transformers at scale. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022b. URL https://openreview.net/forum?id=dXiGWqBoxaD.

Xiuying 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, 2023. URL https://arxiv.org/abs/2209.13325.

Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models, 2024. URL https://arxiv.org/abs/2211.10438.

Yujun Lin, Haotian Tang, Shang Yang, Zhekai Zhang, Guangxuan Xiao, Chuang Gan, and Song Han. Qserve: W4a8kv4 quantization and system co-design for efficient llm serving, 2024b. URL https://arxiv.org/abs/2405.04532.

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, 2024. URL https://arxiv.org/abs/2308.13137.

Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. Atom: Low-bit quantization for efficient and accurate llm serving, 2024. URL https://arxiv.org/abs/2310.19102.

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, 2024. URL https://arxiv.org/abs/2404.00456.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms, 2023. URL https://arxiv.org/abs/2305.14314.