# FIMA-Q: Post-Training Quantization for Vision Transformers by Fisher Information Matrix Approximation

Zhuguanyu Wu[1,2*], Shihe Wang[1,2*], Jiayi Zhang[1,2], Jiaxin Chen[1,2✉], Yunhong Wang[1,2✉]

[1]State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, China
[2]School of Computer Science and Engineering, Beihang University, Beijing, China

{goatwu, shihewang, zhangjyi, jiaxinchen, yhwang}@buaa.edu.cn

## Abstract

*Post-training quantization (PTQ) has stood out as a cost-effective and promising model compression paradigm in recent years, as it avoids computationally intensive model retraining. Nevertheless, current PTQ methods for Vision Transformers (ViTs) still suffer from significant accuracy degradation, especially under low-bit quantization. To address these shortcomings, we analyze the prevailing Hessian-guided quantization loss, and uncover certain limitations of conventional Hessian approximations. By following the block-wise reconstruction framework, we propose a novel PTQ method for ViTs, dubbed FIMA-Q. Specifically, we firstly establish the connection between KL divergence and FIM, which enables fast computation of the quantization loss during reconstruction. We further propose an efficient FIM approximation method, namely DPLR-FIM, by employing the diagonal plus low-rank principle, and formulate the ultimate quantization loss. Our extensive experiments, conducted across various vision tasks with representative ViT-based architectures on public datasets, demonstrate that our method substantially promotes the accuracy compared to the state-of-the-art approaches, especially in the case of low-bit quantization. The source code is available at* [https://github.com/ShiheWang/FIMA-Q](https://github.com/ShiheWang/FIMA-Q).

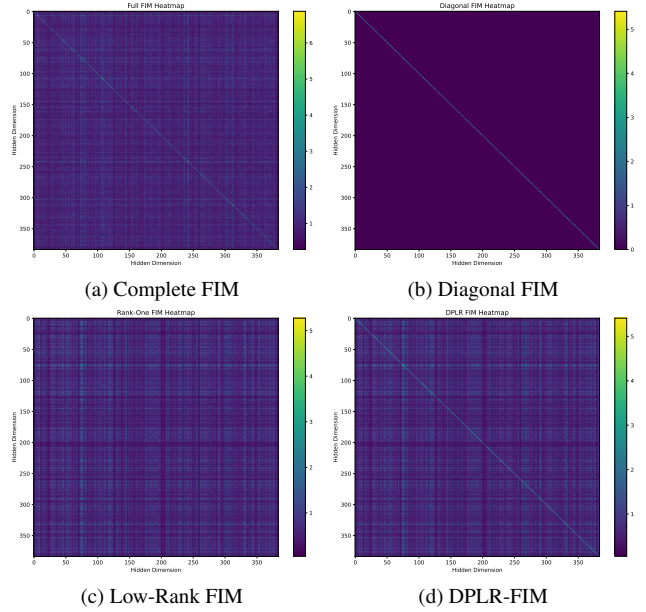(a) Complete FIM  (b) Diagonal FIM

(c) Low-Rank FIM  (d) DPLR-FIM

Figure 1. Illustration on the heatmap of FIM for the class token. (a) shows the complete FIM, where the diagonal elements are notably prominent, and the off-diagonal elements are also non-negligible. (b) and (c) display the diagonal approximation and low-rank approximation on FIM, respectively, both of which omit critical information from the complete FIM. (d) shows the proposed diagonal plus low-rank approximation on FIM, which clearly achieves improvement over compared approaches.

## 1. Introduction

By leveraging the powerful self-attention mechanism [2, 5], Vision Transformers (ViTs) have recently established themselves as a fundamental architecture in the computer vision community [9, 34, 35, 38, 39], challenging the longstanding dominance of Convolutional Neural Networks (CNNs) [12, 15, 28]. However, the improved performance comes at the cost of larger model sizes, more parameters, and increased computational overhead. In order to facilitate the

deployment of ViTs on resource-constrained hardware, substantial research efforts have been devoted to developing model compression techniques.

Model quantization, as one of the most effective model compression methods, aims to convert floating-point weights or activations to low-bit integers, thereby reducing memory consumption and computational cost [14]. However, this process inevitably introduces quantization loss, leading to accuracy degradation. Conventional approaches employ Quantization-Aware Training (QAT) [8, 10, 17] to re-

---

* Equal contribution.
✉ Corresponding author.

cover performance through end-to-end retraining. Despite the significantly improved accuracy, QAT methods often requires training on the entire pretraining dataset, incurring prohibitive computational costs. Consequently, Post-Training Quantization (PTQ) methods, which calibrate quantization parameters on a small unlabeled data, have emerged as a promising alternative [16, 26, 30].

While current PTQ methods have delivered promising performance for CNNs, they suffer from significant performance degradation when applied to ViTs, primarily due to the imbalanced and asymmetric activation distributions of ViTs. Although numerous studies have developed specialized quantizers to accommodate the unique structure and distribution of ViTs, these approaches still underperform at low bit-widths [19, 31, 33]. This propels us to examine the fundamental aspects of existing PTQ methods, and observe that most of them rely on a loss function to guide parameter tuning for minimizing quantization error. Consequently, the formulation and measurement of quantization loss substantially influence the performance of quantized models.

Motivated by our observation that many existing methods essentially employ Hessian-guided quantization loss [6, 7, 16, 33], we investigate the construction of the foundational quantization loss through Fisher Information Matrix (FIM) approximation in this work. Considering the computational intractability of exact Hessian computation, we explore viable approximation strategies. The existing diagonal approximation method, introduced by BRECQ [16], serves as a candidate, initially replacing the Hessian with FIM, followed by a diagonal approximation.

However, our in-depth analysis of FIM approximation reveal two critical findings: 1) As validated in Table 3, the diagonal approximation proposed by BRECQ achieves inferior performance compared to simple mean square error (MSE) loss on ViTs. This is particularly surprising since the MSE loss, being agnostic to task-specific objectives, should theoretically underperform Hessian-based methods in case of block-wise quantization. 2) As shown in Fig. 1 (a), the visualization of the complete FIM for the class token implies that both the diagonal elements and the off-diagonal inter-token correlations of FIM are important. This finding suggests that the prevailing diagonal approximation approach discards potentially valuable off-diagonal components that may be crucial for maintaining the performance.

To address these issues, we conduct a rigorous analysis of FIM approximation. Our investigation challenges the hypothesis of BRECQ, revealing that FIM is linearly proportional to the gradient of the Kullback-Leibler (KL) divergence, instead of the square of the gradient as BRECQ posits. Building upon this insight, we firstly propose a rank-one approximation method that preserves critical off-diagonal correlations and extend it to a low-rank approximation. Subsequently, inspired by the diagonal plus low-rank approx-

imation, we further incorporate the diagnal and low-rank components, formulating a more concise estimation on FIM, as displayed in Fig. 1 (d).

The main contributions of this paper are summarized in the following three aspects:

1) We conduct a thorough analysis on FIM approximation adopted by the prevailing Hessian-guided loss for posttraining quantization. By exploring the relationship between KL divergence and FIM, we reveal that FIM is linearly proportional to the gradient of KL divergence.

2) We propose a novel rank-one approximation method for FIM and further extend it to low-rank approximation. By combining with the diagonal approximation, we formulate a new loss for quantization reconstruction.

3) We extensively evaluate the performance of our proposed method on image classification and object detection across various ViT-based network architectures. The experimental results demonstrate that our method, with a simple uniform quantizer, significantly outperforms the state-of-the-art approaches that adopt specialized quantizers, especially in the case of low-bit quantization.

## 2. Related Work

Current quantization approaches are primarily divided into two categories: Quantization Aware Training (QAT) and Post Training Quantization (PTQ). QAT [8, 10, 17, 18] integrates quantization directly into the training process, which requires a tremendous amount of training data, computational resources, and time. In contrast, PTQ requires only a smaller dataset and lower costs for the fine-tuning process. Currently, numerous PTQ methods have demonstrated promising performance when applied to CNNs. AdaRound [26] pioneers the adaptive quantitative rounding strategy, bringing revolutionary progress in low-bit quantization. BRECQ [16] introduces the block reconstruction framework with a Hessian quantization loss. Successively, QDrop [30] develops the randomly dropping quantization.

Due to the unique structure and properties of ViTs, the straightforward application of PTQ methods originally designed for CNNs [16, 26] does not perform well on ViTs.The problem can be summed up in two aspects: 1) inter-channel variation in layernorm and 2) non-uniform distribution in softmax and GELU. FQ-ViT [21] designed Power-of-Two Factor and $\log 2$ quantizer to solve both problems first. To solve the inter-channel variation, many studies are dedicated to identifying suitable group quantization techniques, such as PEG [1] and IGQ-ViT [25]. RepQ-ViT [19] introduced a scale reparameterization technique, enabling the application of layer-wise quantizers. Meanwhile, to further match distributions of activation values after softmax and GELU, PTQ4ViT [33] proposed twin uniform quantization. Building on PTQ4ViT, APQ-ViT [6] worked on preserving the Matthew effect of post-Softmax activations and presented
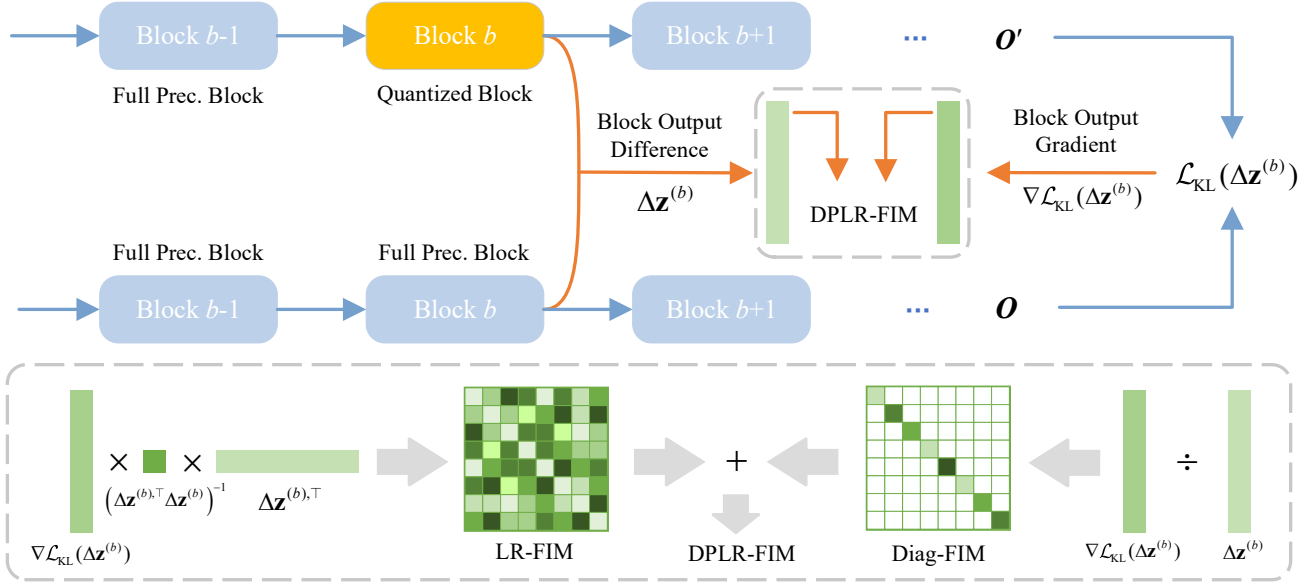
Figure 2. Framework overview of the proposed FIMA-Q method. We follow the block-wise quantization pipeline. For each block, we compute the difference between outputs before and after quantization, perform forward propagation through the rest of the networks, calculate the KL divergence, and conduct backward propagation to obtain the gradients. Based on the output difference and gradients, we compute the low-rank approximation and diagonal approximation on FIM respectively, and combine them to formulate the DPLR-FIM loss for quantization reconstruction.

a unified Bottom-elimination Blockwise Calibration. Furthermore, AdaLog [31] developed a non-uniform quantizer to adaptively select the logarithm base, further improving the accuracy. From an alternative perspective, we observe that many works [6, 7, 16, 33] have employed the hessian matrix as an importance metric. In this paper, we focus on a more accurate approximation of the Hessian matrix and more effective Hessian-guided quantization loss.

## 3. The Proposed Method

Fig. 2 illustrates the framework of the proposed FIMA-Q method. Basically, FIMA-Q is built upon the QDrop [30] framework. Specifically, we reconstruct each block in the ViTs individually. Apart from replacing the MSE loss in QDrop with the DPLR-FIM based quantization loss proposed in Sec. 3.3, we maintain the other steps as utilized in QDrop. The overall pipeline of the FIMA-Q for a certain Transformer block is summarized in Algorithm 1.

### 3.1. Preliminaries

BRECQ introduces a Hessian-guided quantization loss that models quantization as output perturbation during block-wise quantization. Through Taylor expansion, it specifically minimizes the second-order term as below:

$$\min \mathbb{E}\left[-\Delta\mathbf{z}^{(b)\top}\mathbf{H}^{(\mathbf{z}^{(b)})}\Delta\mathbf{z}^{(b)}\right], \quad (1)$$

where $\mathbf{z}^{(b)}$ is the output of the block, $\Delta\mathbf{z}^{(b)}$ is the difference between the output before and after quantization, $\mathbf{H}^{(\mathbf{z}^{(b)})}$ is the Hessian matrix of the task loss $\mathcal{L}_{\text{task}}$ w.r.t. $\mathbf{z}^{(b)}$.

Due to the large size of the Hessian matrix, BRECQ employs a diagonal approximation. Specifically, BRECQ estimates the diagonal of the negative Hessian using squared gradient values:

$$-\mathbf{H}^{(\mathbf{z}^{(b)})} \approx \text{Diag}\left((\frac{\partial\mathcal{L}_{\text{task}}}{\partial\mathbf{z}_1^{(b)}})^2, \cdots, (\frac{\partial\mathcal{L}_{\text{task}}}{\partial\mathbf{z}_a^{(b)}})^2\right). \quad (2)$$

However, within the PTQ framework, the absence of labeled data prevents task loss from being directly computed. Therefore, BRECQ adopts the KL divergence between pre- and post-quantization distributions to approximate the task loss. When calculating the negative Hessian matrix, BRECQ utilizes several approximations: 1) utilizing the Fisher Information matrix (FIM) to approximate the negative Hessian matrix. 2) using the diagonal FIM to approximate the full FIM. 3) adopting squared gradients of the task loss to approximate the diagonal of the FIM. 4) applying the gradient of KL divergence instead of the task loss.

We analyze the validity of these approximations as follows. First, we demonstrate the soundness of approximating the negative Hessian with the Fisher information matrix (FIM) in Sec. 3.2. Second, while the diagonal approximation on FIM is generally acceptable, we argue it introduces un-

**Algorithm 1** Pipeline of Block-wise FIMA-Q.

**Input:** Full-precision model $\mathcal{M}$, full-precision block $\mathcal{B}_{\text{full}}$, quantized block $\mathcal{B}_{\text{quant}}$, calibration data $\mathcal{D}_{\text{calib}}$, rank $r$, maximal iteration number $\text{max\_iter}$, and iteration interval $x$.
**Output:** The optimized quantized block $\mathcal{B}_{\text{quant}}$.
# Calculate Rank-One FIM:
1: Generate the raw input $\boldsymbol{X}_{\text{raw}}$ and output $\mathbf{z}^{(b)}$ of $\mathcal{B}_{\text{full}}$, the quantized input $\boldsymbol{X}_{\text{quant}}$ of $\mathcal{B}_{\text{quant}}$ based on $\mathcal{D}_{\text{calib}}$.
2: Generate the model output $O_{\text{raw}}$ and $O_{\text{quant}}$ by performing forward propagations through the network starting from $\mathcal{B}_{\text{raw}}$ and $\mathcal{B}_{\text{quant}}$ based on $\boldsymbol{X}_{\text{raw}}$, respectively.
3: Calculate the perturbation $\Delta\mathbf{z}^{(b)}$ and the loss $\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)})$ based on Eq. (9), and perform backward propagation to compute the gradient $\nabla\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)})$.
# Perform Quantization Reconstruction:
4: Initialize the current rank $k = 1$, and the next data iteration $y = x$.
5: **for** $i = 1, \cdots, \text{max\_iter}$ **do**
6:     **if** $k < r$ and $i = y$ **then**
7:         Calculate $\Delta\mathbf{z}^{(b)'}$ and $\nabla\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)'})$.
8:         Calculate $B$ in Eq. (19).
9:         Set $k := k + 1$ and $y := y + x$.
10:     **end if**
11:     Calculate $\mathcal{L}_{\text{DPLR}}$ with Eq. (21) and perform BP.
12:     Update all the AdaRound weights by [26] and activation scaling factors in $\mathcal{B}_{\text{quant}}$.
13: **end for**

---

necessary limitations and develops an improved approach in Sec. 3.3. Third, we identify the squared gradient approximation utilized in BRECQ is inaccurate and provide elaborate analysis in Sec. 3.2. Finally, we show the KL-divergence gradient substitution is well-justified, and provide the proof in Sec. A.1 of the *supplementary material*.

### 3.2. Relationship Between FIM and KL-Divergence

Our method fundamentally relies on replacing the Hessian in Eq. (1) with FIM for optimization. We therefore first establish the validity of this replacement.

Specifically, on the calibration set, the expected value of reconstruction loss is written as:

$$\mathbb{E}\left[\mathcal{L}_{\text{calib}}\right] = -\mathbb{E}\left[\Delta\mathbf{z}^{(b)\top}\mathbf{H}^{(\mathbf{z}^{(b)})}\Delta\mathbf{z}^{(b)}\right]. \quad (3)$$

Since the Hessian matrix is a function of the sample and the full-precision block output $\mathbf{z}^{(b)}$, we assume that it is independent of the perturbation $\Delta\mathbf{z}^{(b)}$ caused by quantization. Therefore, we can deduce the following:

$$\mathbb{E}\left[\mathcal{L}_{\text{calib}}\right] = -\mathbb{E}\left[\Delta\mathbf{z}^{(b),\top}\mathbb{E}\left[\mathbf{H}^{(\mathbf{z}^{(b)})}\right]\Delta\mathbf{z}^{(b)}\right]. \quad (4)$$

---

**Theorem 3.1.** *When the expected gradient of the log-likelihood of model outputs becomes zero, FIM $\boldsymbol{F}^{(\boldsymbol{z}^{(b)})}$ equals to the expected negative Hessian, i.e.,*

$$\mathbf{F}^{(\mathbf{z}^{(b)})} = \mathbb{E}\left[\left(\nabla_{\mathbf{z}^{(b)}}\log p(y;\mathbf{z}^{(b)})\right)\left(\nabla_{\mathbf{z}^{(b)}}\log p(y;\mathbf{z}^{(b)})\right)^{\top}\right],$$
$$(5)$$

*where $\log p(y;\mathbf{z}^{(b)})$ is the log-likelihood function.*

We provide the detailed proof in Sec. A.2 of the *supplementary material*. Based on Theorem 3.1 and Eq. (3), the target to optimize can be written as:

$$\min\mathbb{E}\left[\Delta\mathbf{z}^{(b),\top}\mathbf{F}^{(\mathbf{z}^{(b)})}\Delta\mathbf{z}^{(b)}\right]. \quad (6)$$

Note that FIM is an inherent attribute of the network block. In BRECQ, directly approximating the diagonal elements of FIM as the squared gradient is inaccurate. Actually, by definition, the following equation holds

$$\text{Diag}\left(\mathbf{F}^{(\mathbf{z}^{(b)})}\right) = \mathbb{E}\left[\left(\nabla_{\mathbf{z}^{(b)}}\log p(y;\mathbf{z}^{(b)})\right)^{2}\right]$$
$$= \mathbb{E}\left[\left(\nabla_{\mathbf{z}^{(b)}}\log p(y;\mathbf{z}^{(b)})\right)\right]^{2} \quad (7)$$
$$+ \text{Var}\left(\nabla_{\mathbf{z}^{(b)}}\log p(y;\mathbf{z}^{(b)})\right).$$

In BRECQ, FIM is regarded as a sample-dependent matrix, and the second variance term in Eq. (7) is omitted, resulting in larger approximation errors.

Thereafter, we derive the relationship between the KL-divergence and FIM to propose an improved approximation of FIM. The KL-divergence based loss is defined as:

$$\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)}) = D_{\text{KL}}(p(y;\mathbf{z}^{(b)})\|p(y;\mathbf{z}^{(b)} + \Delta\mathbf{z}^{(b)}))$$
$$= \sum_{x} p(y;\mathbf{z}^{(b)})\log\frac{p(y;\mathbf{z}^{(b)})}{p(y;\mathbf{z}^{(b)} + \Delta\mathbf{z}^{(b)})}. \quad (8)$$

Accordingly, $\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)})$ satisfies the following property.

**Theorem 3.2.** *When adopting the KL-divergence as the task loss, the following relationship holds:*

$$\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)}) = \frac{1}{2}\Delta\mathbf{z}^{(b),\top}\mathbf{F}^{(\mathbf{z}^{(b)})}\Delta\mathbf{z}^{(b)}. \quad (9)$$

We provide the detailed proof in Sec. A.3 of the *supplementary material*. Differentiating both sides of Eq. (9) w.r.t. $\Delta\mathbf{z}^{(b)}$ yields the following:

$$\nabla\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)}) = \frac{\partial\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)})}{\partial\Delta\mathbf{z}^{(b)}} = \mathbf{F}^{(\mathbf{z}^{(b)})}\Delta\mathbf{z}^{(b)}. \quad (10)$$

By adopting the diagonal form approximation on FIM employed in prior works, we can derive that

$$\mathbf{F}_{\text{Diag}}^{(\mathbf{z}^{(b)})} = \text{Diag}\left(\frac{\nabla\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)})_1}{\Delta\mathbf{z}_1^{(b)}}, \cdots, \frac{\nabla\mathcal{L}_{\text{KL}}(\Delta\mathbf{z}^{(b)})_a}{\Delta\mathbf{z}_a^{(b)}}\right),$$
$$(11)$$

which indicates that the diagonal elements of FIM is linearly correlated with the gradient of KL-divergence, rather than being related to the squared gradient as claimed in BRECQ.

### 3.3. Improved Approximations of FIM

In this section, we establish improved estimations on FIM. Concretely, for the $i-$th sample, we denote that

$$\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b,i)}) = \mathbf{F}^{(\mathbf{z}^{(b)})} \Delta \mathbf{z}^{(b,i)}. \tag{12}$$

Therefore, $\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})$ and $\Delta \mathbf{z}^{(b)}$ are rewritten as

$$\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) = \sum_{i=1}^{n} \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b,i)}),$$
$$\Delta \mathbf{z}^{(b)} = \sum_{i=1}^{n} \Delta \mathbf{z}^{(b,i)}. \tag{13}$$

By assuming the independence between samples, Eq. (10) holds when using Eq. (13).

Our optimization problem now reduces to optimizing Eq. (6) under the constraints that FIM satisfies Eq. (10) while being symmetry and positive definite. As the solution is not unique, we present four distinct quantization losses based on different approximation forms of FIM.

**1. Diagonal approximation.** By assuming FIM is diagonal and based on Eq. (11), the loss is formulated as:

$$\mathcal{L}_{\mathrm{diag}} = \left( \frac{\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})}{\Delta \mathbf{z}^{(b)}} \right)^{\top} \left( \Delta \mathbf{z}^{(b,i)} \right)^2. \tag{14}$$

**2. Rank-one approximation.** By assuming FIM $\mathbf{F}_{\mathrm{rank}-1}^{(\mathbf{z}^{(b)})} = \boldsymbol{u}\boldsymbol{u}^{\top}$ where $\boldsymbol{u} \in \mathbb{R}^{a \times 1}$, we have

$$\boldsymbol{u} = \frac{\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})}{\sqrt{\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})^{\top} \Delta \mathbf{z}^{(b,i)}}}. \tag{15}$$

We refer to Sec. A.4 of the *supplementary material* for details. Given the $i-$th sample, the loss is written as

$$\mathcal{L}_{\mathrm{rank}-1} = \frac{\left( \Delta \mathbf{z}^{(b,i)\top} \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) \right)^2}{\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})^{\top} \Delta \mathbf{z}^{(b,i)}}, \tag{16}$$

which maintains an $O(a)$ computational complexity. Unlike the diagonal approximation that treats block outputs as independent contributors to the task loss, our method captures collective dependencies among elements of the output, without introducing extra computational overhead.

**3. Low-rank approximation.** Before extending the rank-one approximation to rank-$k$, we first establish the following corollary.

**Corollary 3.1.** *When $k > 1$, it is typically difficult to find a low-rank matrix $\boldsymbol{u} \in \mathbb{R}^{a \times k}$ such that $\mathbf{F}^{(\mathbf{z}^{(b)})} = \boldsymbol{u}\boldsymbol{u}^{\top}$ satisfying Eq. (10).*

We refer to Sec. A.5 of the *supplementary material* for a detailed description. Accordingly, we consider relaxing the constraint of symmetry to compute the low-rank FIM.

Specifically, by taking the Moore-Penrose inverse of $\Delta \mathbf{z}^{(b)}$, the rank-$k$ FIM is computer as

$$\Delta \mathbf{z}^{(b),+} = \left( \Delta \mathbf{z}^{(b)\top} \Delta \mathbf{z}^{(b)} \right)^{-1} \Delta \mathbf{z}^{(b)\top},$$
$$\mathbf{F}_{\mathrm{rank}-k}^{(\mathbf{z}^{(b)})} = \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) \Delta \mathbf{z}^{(b)+}. \tag{17}$$

Consequently, the optimization objective for the $i-$th sample is formed as

$$\mathcal{L}_{\mathrm{rank}-k} = \Delta \mathbf{z}^{(bi)\top} \mathbf{F}^{(\mathbf{z}^{(b)})} \Delta \mathbf{z}^{(b,i)} = A \cdot B \cdot C, \tag{18}$$

where

$$A = \Delta \mathbf{z}^{(b,i)\top} \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}),$$
$$B = \left( \Delta \mathbf{z}^{(b)\top} \Delta \mathbf{z}^{(b)} \right)^{-1},$$
$$C = \Delta \mathbf{z}^{(b)\top} \Delta \mathbf{z}^{(b,i)}. \tag{19}$$

The computational complexity of computing $A$ and $C$ is $O(ak)$, and $B$ can be preprocessed. Therefore, the overall complexity of computing $\mathcal{L}_{\mathrm{recon}}$ is $O(ak)$, which is acceptable for small rank $k$.

It is worth noting that the low-rank approximation requires that the matrix $\Delta \mathbf{z}^{(b)\top} \Delta \mathbf{z}^{(b)}$ is invertible, implying that the rank of $\Delta \mathbf{z}^{(b)}$ should be $k$. Therefore, we need to obtain $k$ linearly independent perturbations and their corresponding gradients. We employ a progressive strategy to gradually increase $k$ in practice. Specifically, we initialize with a rank of 1 and increase 1 rank by $x$ iterations. For every $x$ iteration, we recompute the averaged quantization perturbations and averaged gradients based on Eq. (13), and concatenate them with the previously computed perturbations $\Delta \mathbf{z}^{(b)}$ and gradients $\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})$.

**4. DPLR approximation.** As shown in Fig. 1, the low-rank approximation of the FIM may degrade the impact of individual outputs on task loss. Thus, we combine the diagonal and low-rank approximate and obtain the following diagonal plus low-rank (DPLR) form:

$$\mathbf{F}_{\mathrm{DPLR}}^{(\mathbf{z}^{(b)})} = \alpha \cdot \mathbf{F}_{\mathrm{rank}-k}^{(\mathbf{z}^{(b)})} + (1 - \alpha) \cdot \mathbf{F}_{\mathrm{diag}}^{(\mathbf{z}^{(b)})}. \tag{20}$$

Accordingly, the optimization objective is finally formulated as below:

$$\mathcal{L}_{\mathrm{DPLR}} = \alpha \cdot \mathcal{L}_{\mathrm{rank}-k} + (1 - \alpha) \cdot \mathcal{L}_{\mathrm{diag}}. \tag{21}$$

In this work, we adopt Eq. (21) based on DPLR-FIM approximation in Eq. (20) as the ultimate quantization loss.

Table 1. Comparison of the top-1 accuracy across various ViT-based models on ImageNet. "*" denotes the results are based on our re-implementation as QDrop is originally designed for CNNs, "OP" indicates whether the PTQ method is an optimization-based one. "SQ" signifies that the method requires a specific quantizer rather than a standard uniform quantizer. The best results are highlighted in **bold**.

| Method | OP | SQ | W/A | ViT-S | ViT-B | DeiT-T | DeiT-S | DeiT-B | Swin-S | Swin-B |
|---|---|---|---|---|---|---|---|---|---|---|
| Full-Prec | - | - | 32/32 | 81.39 | 84.54 | 72.71 | 79.85 | 81.80 | 83.23 | 85.27 |
| PTQ4ViT [33] | × | ✓ | 3/3 | 0.10 | 0.10 | 3.50 | 0.10 | 31.06 | 28.69 | 20.13 |
| RepQ-ViT [19] | × | ✓ | 3/3 | 0.10 | 0.10 | 0.10 | 0.10 | 0.10 | 0.10 | 0.10 |
| AdaLog [31] | × | ✓ | 3/3 | 13.88 | 37.91 | 31.56 | 24.47 | 57.47 | 64.41 | 69.75 |
| I&S-ViT [36] | ✓ | ✓ | 3/3 | 45.16 | 63.77 | 41.52 | 55.78 | 73.30 | 74.20 | 69.30 |
| DopQ-ViT [32] | ✓ | ✓ | 3/3 | 54.72 | 65.76 | 44.71 | 59.26 | 74.91 | 74.77 | 69.63 |
| QDrop* [30] | ✓ | × | 3/3 | 41.05 | 74.75 | 46.88 | 50.95 | 72.97 | 74.67 | 76.57 |
| **FIMA-Q (Ours)** | ✓ | × | 3/3 | **64.09** | **77.63** | **55.55** | **69.13** | **76.54** | **77.26** | **78.82** |
| PTQ4ViT [33] | × | ✓ | 4/4 | 42.57 | 30.69 | 36.96 | 34.08 | 64.39 | 76.09 | 74.02 |
| APQ-ViT [6] | × | ✓ | 4/4 | 47.95 | 41.41 | 47.94 | 43.55 | 67.48 | 77.15 | 76.48 |
| RepQ-ViT [19] | × | ✓ | 4/4 | 65.05 | 68.48 | 57.43 | 69.03 | 75.61 | 79.45 | 78.32 |
| ERQ [37] | × | ✓ | 4/4 | 68.91 | 76.63 | 60.29 | 72.56 | 78.23 | 80.74 | 82.44 |
| IGQ-ViT [25] | × | ✓ | 4/4 | 73.61 | 79.32 | 62.45 | 74.66 | 79.23 | 80.98 | 83.14 |
| AdaLog [31] | × | ✓ | 4/4 | 72.75 | 79.68 | 63.52 | 72.06 | 78.03 | 80.77 | 82.47 |
| I&S-ViT [36] | ✓ | ✓ | 4/4 | 74.87 | 80.07 | 65.21 | 75.81 | 79.97 | 81.17 | 82.60 |
| DopQ-ViT [32] | ✓ | ✓ | 4/4 | 75.69 | 80.95 | 65.54 | 75.84 | 80.13 | 81.71 | 83.34 |
| QDrop* [30] | ✓ | × | 4/4 | 71.84 | 82.63 | 65.27 | 72.64 | 79.96 | 81.21 | 82.99 |
| OASQ [24] | ✓ | × | 4/4 | 72.88 | 76.59 | 66.31 | 76.00 | 78.83 | 81.02 | 82.46 |
| **FIMA-Q (Ours)** | ✓ | × | 4/4 | **76.68** | **83.04** | **66.84** | **76.87** | **80.33** | **81.82** | **83.60** |
| PTQ4ViT [33] | × | ✓ | 6/6 | 78.63 | 81.65 | 69.68 | 76.28 | 80.25 | 82.38 | 84.01 |
| APQ-ViT [6] | × | ✓ | 6/6 | 79.10 | 82.21 | 70.49 | 77.76 | 80.42 | 82.67 | 84.18 |
| NoisyQuant [22] | × | ✓ | 6/6 | 78.65 | 82.32 | - | 77.43 | 80.70 | 82.86 | 84.68 |
| RepQ-ViT [19] | × | ✓ | 6/6 | 80.43 | 83.62 | 70.76 | 78.90 | 81.27 | 82.79 | 84.57 |
| IGQ-ViT [25] | × | ✓ | 6/6 | 80.76 | 83.77 | 71.15 | 79.28 | 81.71 | 82.86 | 84.82 |
| AdaLog [31] | × | ✓ | 6/6 | **80.91** | 84.80 | 71.38 | 79.39 | 81.55 | **83.19** | **85.09** |
| I&S-ViT [36] | ✓ | ✓ | 6/6 | 80.43 | 83.82 | 70.85 | 79.15 | 81.68 | 82.89 | 84.94 |
| DopQ-ViT [32] | ✓ | ✓ | 6/6 | 80.52 | 84.02 | 71.17 | 79.30 | 81.69 | 82.95 | 84.97 |
| QDrop* [30] | ✓ | × | 6/6 | 79.59 | 84.68 | 71.48 | 79.15 | 81.69 | 83.01 | 84.94 |
| OASQ [24] | ✓ | × | 6/6 | 80.60 | 83.81 | 71.52 | 79.50 | 81.72 | 82.76 | 84.91 |
| **FIMA-Q (Ours)** | ✓ | × | 6/6 | 80.64 | **84.82** | **71.53** | **79.52** | **81.74** | **83.19** | 85.01 |

## 4. Experimental Results and Analysis

In this section, we evaluate the performance of our method across distinct vision tasks including image classification, object detection and instance segmentation with various ViT-based architectures, by comparing to the state-of-the-art PTQ approaches. We also perform extensive ablation studies on the proposed components.

### 4.1. Experimental setup

**Datasets and Models.** For the image classification task, we adopt the ImageNet [27] dataset for validation, based on various vision transformer architectures, including ViT [9], DeiT [29], and Swin [23]. For object detection and instance segmentation, we evaluate on the COCO [20] dataset based on

the representative Mask R-CNN [13] and Cascade R-CNN [3] models that utilize Swin Transformer as the backbone.

**Implementation Details.** Similar to [19, 33], all the pretrained full-precision Vision Transformers for classification are obtained from the timm library[1]. The pretrained detection and segmentation models are obtained from MMDetection [4]. By following the reconstruction-based PTQ methods [16, 24, 30, 36], we randomly select 1024 unlabeled images from ImageNet and 256 unlabeled images from COCO as the calibration sets for classification and object detection, respectively. We adopt the vanilla channel-wise uniform quantizers for weight quantization and layer-wise uniform quantizers for activation quantization, including post-Softmax activa-

---

[1]https://github.com/huggingface/pytorch-image-models

Table 2. Quantization results (%) on COCO for the object detection and instance segmentation tasks. "*" and "†" indicate the results are based on our re-implementation or re-production using the officially released source code. The best results are highlighted in **bold**.

| Method | Opt | SQ | W/A | Mask R-CNN | | | | Cascade Mask R-CNN | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Swin-T | | Swin-S | | Swin-T | | Swin-S | |
| | | | | $AP^b$ | $AP^m$ | $AP^b$ | $AP^m$ | $AP^b$ | $AP^m$ | $AP^b$ | $AP^m$ |
| Full-Precision | - | - | 32/32 | 46.0 | 41.6 | 48.5 | 43.3 | 50.4 | 43.7 | 51.9 | 45.0 |
| Baseline | ✗ | ✗ | 4/4 | 34.6 | 34.2 | 40.8 | 38.6 | 45.9 | 40.2 | 47.9 | 41.6 |
| PTQ4ViT [33] | ✗ | ✓ | 4/4 | 6.9 | 7.0 | 26.7 | 26.6 | 14.7 | 13.5 | 0.5 | 0.5 |
| APQ-ViT [6] | ✗ | ✓ | 4/4 | 23.7 | 22.6 | **44.7** | 40.1 | 27.2 | 24.4 | 47.7 | 41.1 |
| RepQ-ViT [19] | ✗ | ✓ | 4/4 | 36.1 | 36.0 | $44.2_{42.7}$† | 40.2 | 47.0 | 41.1 | 49.3 | 43.1 |
| ERQ [37] | ✗ | ✓ | 4/4 | 36.8 | 36.6 | 43.4 | 40.7 | 47.9 | 42.1 | 50.0 | 43.6 |
| I&S-ViT [36] | ✓ | ✓ | 4/4 | 37.5 | 36.6 | 43.4 | 40.3 | 48.2 | 42.0 | 50.3 | 43.6 |
| DopQ-ViT [32] | ✓ | ✓ | 4/4 | 37.5 | 36.5 | 43.5 | 40.4 | 48.2 | 42.1 | 50.3 | **43.7** |
| QDrop* [30] | ✓ | ✗ | 4/4 | 36.2 | 35.4 | 41.6 | 39.2 | 47.0 | 41.3 | 49.0 | 42.5 |
| **FIMA-Q (Ours)** | ✓ | ✗ | 4/4 | **38.7** | **37.8** | 44.2 | **41.1** | **48.7** | **42.5** | 50.4 | **43.7** |

tions. We fix the rank for the proposed DPLR-FIM module as $k = 15$.

As established in Theorem 3.1, the FIM can substitute for the Hessian matrix when the loss function corresponds to the log-likelihood, specifically the Cross-Entropy loss. In order to adapt this framework to the detection and segmentation tasks, we employ the classification head from the regional proposal network (RPN) as the task output. Specifically, we compute FIM using the KL divergence between the classification outputs from RPN by the original and quantized models. We apply block reconstruction exclusively to the backbone network and Feature Pyramid Network (FPN), while conducting calibration-only quantization for the RPN and Region of Interest (ROI) modules.

## 4.2. Comparison to the State-of-the-art Approaches

**Quantization Results for Classification on ImageNet.** We first evaluate the performance of our method on the classification task on ImageNet in terms of top-1 accuracy, compared to the state-of-the-art PTQ approaches. To highlight the advantages of our method, we report results across various representative Transformer architectures, including ViT-S/B, DeiT-T/S and Swin-S/B, under 6, 4 and 3 bits.

As displayed in Tab. 1, FIMA-Q consistently promotes accuracy across different settings, with particularly significant gains in case of low-bit quantization. Concretely, for 6-bit quantization, our method achieves either superior or competitive accuracy compared to the second-best approaches. For 4-bit quantization, while most existing methods suffer substantial accuracy degradation, our method maintains robust performance in most cases. As for the challenging 3-bit quantization, the performance of competing methods degrades dramatically, some of which even exhibit extremely poor results, such as PTQ4ViT [33] and RepQ-ViT [19]. In com-

parison, our proposed FIMA-Q achieves minimal accuracy loss, surpassing the second-best approach by 5.31% on average. It is worth noting that many compared approaches such as PTQ4ViT, RepQ-ViT, AdaLog and DopQ-ViT attempt to boost the performance by designing specific quantizers (SQ), which however are generally difficult to implement on hardware in practice. In contrast, our method only utilizes a standard uniform quantizer, making it hardware-friendly while reaching superior accuracy.

**Quantization Results for Object Detection and Instance Segmentation on COCO.** As shown in Tab. 2, our method achieves the best results under W4/A4 in most cases for object detection and instance segmentation. To establish a fair comparison baseline, we adapt RepQ-ViT by replacing its specialized quantizer with a standard uniform quantizer. Notably, the methods that exclusively adopt uniform quantizers including Baseline and QDrop consistently underperform approaches employing specific quantizers. These results imply the inherent limitations of uniform quantization in handling activation distributions. Despite this challenge, our method can significantly mitigate the quantization error by leveraging the proposed FIM-based loss, reaching state-of-the-art performance even when using a standard uniform quantizer.

## 4.3. Ablation Study

As we mainly analyze the FIM-based reconstruction and propose a quantization loss based on DPLR-FIM approximation in this paper, we evaluate its effectiveness by comparing to distinct quantization losses and studying the influence of different ranks on the DPLR-FIM component.

**On distinct quantization losses.** We choose the MSE loss as the baseline method, and compare with the conventional representative Hessian-guided loss adopted by BRECQ [16], denoted by BRECQ-FIM. We also report the results of our

Table 3. Comparison of the top-1 accuracy (%) by using different quantization losses across distinct Transformer architectures on ImageNet.

| #Bits (W/A) | Method | ViT-S | ViT-B | DeiT-T | DeiT-S | DeiT-B | Swin-S | Swin-B |
|---|---|---|---|---|---|---|---|---|
| | MSE [30] | 41.05 | 74.75 | 46.88 | 50.95 | 72.97 | 74.67 | 76.57 |
| | BRECQ-FIM [16] | 14.65 | 11.61 | 36.57 | 49.20 | 58.76 | 66.26 | 70.15 |
| 3/3 | Diag-FIM **(Ours)** | 60.02 | 76.29 | 55.54 | 68.68 | 76.32 | 75.08 | 77.87 |
| | LR-FIM **(Ours)** | **64.09** | 77.46 | 55.25 | 68.91 | 76.33 | 76.03 | 77.59 |
| | **DPLR-FIM (Ours)** | **64.09** | **77.63** | **55.55** | **69.13** | **76.54** | **77.26** | **78.82** |
| | MSE [30] | 71.84 | 82.63 | 65.27 | 72.64 | 79.96 | 81.21 | 82.99 |
| | BRECQ-FIM [16] | 63.70 | 76.26 | 61.99 | 72.52 | 76.59 | 80.52 | 81.80 |
| 4/4 | Diag-FIM **(Ours)** | 75.88 | 83.02 | 66.81 | 76.79 | 80.19 | 81.18 | 83.35 |
| | LR-FIM **(Ours)** | 76.47 | **83.04** | 66.78 | 76.66 | 80.30 | 81.60 | 83.15 |
| | **DPLR-FIM (Ours)** | **76.65** | **83.04** | **66.84** | **76.87** | **80.33** | **81.82** | **83.60** |



Figure 3. Influence of the rank $k$ on the accuracy of DPLR-FIM on ImageNet.

Table 4. The training time cost (GPU Minutes) using different ranks $k$ on a single Nvidia RTX 4090 GPU.

| Model | Qdrop* | Ours | | | |
|---|---|---|---|---|---|
| | | k=1 | k=5 | k=10 | k=15 |
| DeiT-T | 100 | 100 | 115 | 150 | 180 |
| DeiT-S | 105 | 105 | 120 | 160 | 225 |
| DeiT-B | 145 | 150 | 160 | 240 | 310 |
| Swin-S | 160 | 165 | 235 | 360 | 420 |
| Swin-B | 170 | 180 | 250 | 420 | 480 |

method using the diagonal estimation, low-rank estimation as well as their combination, denoted by Diag-FIM, LR-FIM, and DPLR-FIM, respectively.

As shown in Tab. 3, despite incorporating second-order Hessian information, BRECQ-FIM generally performs worse than the conventional MSE loss, due to the inaccurate estimation on FIM. In contrast, based on our finding that FIM is linearly proportional to the KL divergence gradient rather than its square, both the proposed Diag-FIM and LR-FIM unleash the potential of Hessian-guided losses, substantially surpassing MSE and BRECQ-FIM. Their combination dubbed DPLR-FIM further boosts the performance, promoting the accuracy of MSE by 23.04% and 18.18% on ViT-S and DeiT-S under W3/A3 respectively, with an average improvement of 8.74%.

**Influence of rank $k$ on DPLR-FIM.** As the rank $k$ plays an important role in DPLR-FIM, we study its influence on both accuracy and efficiency of DPLR-FIM. As illustrated in Fig. 3, our method performs steadily when using various values of $k$ at 4 bits. Under 3-bit quantization, the accuracy of DPLR-FIM improves with higher ranks, showing a gradual upward trend despite minor fluctuations.

In regards of efficiency, as described in Sec. 3.3, the computational complexity of the loss $\mathcal{L}_{\text{recon}}$ is $O(ak)$, indicating that higher ranks linearly increase the reconstruction

time cost. This necessitates balancing between accuracy and computational efficiency. To explore the trade-off, Table 4 illustrates the impact of the rank $k$ on the reconstruction time for 3-bit quantization. The results reveal that the training time cost slightly increases as $k$ becomes large for small models but grows substantially for large network architectures such as Swin-B. However, the overall training cost remains affordable, requiring less than 480 GPU minutes on a single Nvidia RX 4090 GPU.

## 5. Conclusion

In this paper, we propose a novel approach dubbed FIMA-Q for post-training quantization of Vision Transformers. Specifically, we propose a more accurate approximation method for FIM. Specifically, we first demonstrate that FIM is proportional to the gradient of KL-divergence, based on which we develop a novel estimation method dubbed DPLR-FIM by integrating both the diagonal and off-diagonal information. Extensive experimental results across distinct vision transformer architectures validate the effectiveness of FIMA-Q for various visual tasks. The results reveal that our method has achieved significant performance improvements in the cases of low-bit quantization, especially with an average improvement of 5.31%, compared to the current state-of-the-art approaches in 3-bit quantization.

## Acknowledgments

## References

[1] Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Understanding and overcoming the challenges of efficient transformer quantization. In *EMNLP*, 2021. 2

[2] 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. In *NeurIPS*, 2020. 1

[3] Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: delving into high quality object detection. In *CVPR*, pages 6154–6162, 2018. 6

[4] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 6

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, pages 4171–4186, 2019. 1

[6] Yifu Ding, Haotong Qin, Qinghua Yan, Zhenhua Chai, Junjie Liu, Xiaolin Wei, and Xianglong Liu. Towards accurate post-training quantization for vision transformer. In *ACM MM*, pages 5380–5388, 2022. 2, 3, 6, 7

[7] Yifu Ding, Weilun Feng, Chuyan Chen, Jinyang Guo, and Xianglong Liu. Reg-ptq: Regression-specialized post-training quantization for fully quantized object detector. In *CVPR*, 2024. 2, 3

[8] Peiyan Dong, Lei Lu, Chao Wu, Cheng Lyu, Geng Yuan, Hao Tang, and Yanzhi Wang. Packqvit: Faster sub-8-bit vision transformers via full and packed quantization on the mobile. In *NeurIPS*, 2023. 1, 2

[9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 1, 6

[10] Steven K. Esser, Jeffrey L. McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S. Modha. Learned step size quantization. In *ICLR*, 2020. 1, 2

[11] Fisher and Ronald Aylmer. On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 222:309–368, 1922. 11

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 1

[13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. In *ICCV*, pages 2980–2988, 2017. 6

[14] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018. 1

[15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1106–1114, 2012. 1

[16] Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. BRECQ: pushing the limit of post-training quantization by block reconstruction. In *ICLR*, 2021. 2, 3, 6, 7, 8

[17] Yanjing Li, Sheng Xu, Baochang Zhang, Xianbin Cao, Peng Gao, and Guodong Guo. Q-vit: Accurate and fully quantized low-bit vision transformer. In *NeurIPS*, 2022. 1, 2

[18] Zhikai Li and Qingyi Gu. I-vit: Integer-only quantization for efficient vision transformer inference. In *ICCV*, pages 17019–17029, 2023. 2

[19] Zhikai Li, Junrui Xiao, Lianwei Yang, and Qingyi Gu. Repq-vit: Scale reparameterization for post-training quantization of vision transformers. In *ICCV*, pages 17227–17236, 2023. 2, 6, 7

[20] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, pages 740–755, 2014. 6

[21] Yang Lin, Tianyu Zhang, Peiqin Sun, Zheng Li, and Shuchang Zhou. Fq-vit: Post-training quantization for fully quantized vision transformer. In *IJCAI*, pages 1173–1179, 2022. 2

[22] Yijiang Liu, Huanrui Yang, Zhen Dong, Kurt Keutzer, Li Du, and Shanghang Zhang. Noisyquant: Noisy bias-enhanced post-training activation quantization for vision transformers. In *CVPR*, pages 20321–20330, 2023. 6

[23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 9992–10002, 2021. 6

[24] Yuexiao Ma, Huixia Li, Xiawu Zheng, Feng Ling, Xuefeng Xiao, Rui Wang, Shilei Wen, Fei Chao, and Rongrong Ji. Outlier-aware slicing for post-training quantization in vision transformer. In *ICML*, 2024. 6

[25] Jaehyeon Moon, Dohyung Kim, Junyong Cheon, and Bumsub Ham. Instance-aware group quantization for vision transformers. In *CVPR*, 2024. 2, 6

[26] Markus Nagel, Rana Ali Amjad, Mart van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? adaptive rounding for post-training quantization. In *ICML*, pages 7197–7206, 2020. 2, 4

[27] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 6

[28] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 1

[29] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357, 2021. 6

[30] Xiuying Wei, Ruihao Gong, Yuhang Li, Xianglong Liu, and Fengwei Yu. Qdrop: Randomly dropping quantization for extremely low-bit post-training quantization. In *ICLR*, 2022. 2, 3, 6, 7, 8

[31] Zhuguanyu Wu, Jiaxin Chen, Hanwen Zhong, Di Huang, and Yunhong Wang. Adalog: Post-training quantization for vision transformers with adaptive logarithm quantizer. In *ECCV*, 2024. 2, 3, 6

[32] Lianwei Yang, Haisong Gong, and Qingyi Gu. Dopq-vit: Towards distribution-friendly and outlier-aware post-training quantization for vision transformers. *arXiv preprint arXiv:2408.03291*, 2024. 6, 7

[33] Zhihang Yuan, Chenhao Xue, Yiqi Chen, Qiang Wu, and Guangyu Sun. Ptq4vit: Post-training quantization for vision transformers with twin uniform quantization. In *ECCV*, pages 191–207, 2022. 2, 3, 6, 7

[34] Yanan Zhang, Jiaxin Chen, and Di Huang. Cat-det: Contrastively augmented transformer for multi-modal 3d object detection. In *CVPR*, 2022. 1

[35] Hanwen Zhong, Jiaxin Chen, Yutong Zhang, Di Huang, and Yunhong Wang. Transforming vision transformer: Towards efficient multi-task asynchronous learner. In *NeurIPS*, 2024. 1

[36] Yunshan Zhong, Jiawei Hu, Mingbao Lin, Mengzhao Chen, and Rongrong Ji. I&s-vit: An inclusive & stable method for pushing the limit of post-training vits quantization. *arXiv preprint arXiv:2311.10126*, 2023. 6, 7

[37] Yunshan Zhong, Jiawei Hu, You Huang, Yuxin Zhang, and Rongrong Ji. ERQ: Error reduction for post-training quantization of vision transformers. In *ICML*, 2024. 6, 7

[38] Chao Zhou, Yanan Zhang, Jiaxin Chen, and Di Huang. Octr: Octree-based transformer for 3d object detection. In *CVPR*, 2023. 1

[39] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: deformable transformers for end-to-end object detection. In *ICLR*, 2021. 1

# FIMA-Q: Post-Training Quantization for Vision Transformers by Fisher Information Matrix Approximation

## Supplementary Material

## A. Main Proofs

The proofs of this section are based on a assumption that are generally adopted in Fisher Information Matrix (FIM), which is called the regularity condition [11].

**Regularity Condition:** Supposing that $f$ is the likelihood with parameter $\theta$, and $T(x)$ is a function independent of the differentiation parameter $\theta$, then :

$$\frac{\partial}{\partial \theta} \int_{\mathbb{R}} T(x) f(x; \theta) \mathrm{d}x = \int_{\mathbb{R}} T(x) \frac{\partial}{\partial \theta} f(x; \theta) \mathrm{d}x. \quad (22)$$

### A.1. Proof of Approximation 4

We believe that BRECQ adopts the gradient of KL divergence instead of the task loss gradient is based on the following theorems.

**Theorem A.1.** *When the model's output distribution matches the true data distribution, the Hessian matrix of the KL divergence after a small perturbation of the model is exactly equal to the expectation of the Hessian matrix of the model's likelihood function.*

*Proof.* The Hessian matrix of the model's likelihood function is defined as:

$$\boldsymbol{H}(\theta) \triangleq \frac{\partial^2}{\partial \theta^2} \log f(X; \theta). \quad (23)$$

As mentioned in Theorem 3.1, when the assumption that the model's output distribution matches the true data distribution is satisfied, the expectation of the Hessian matrix is equal to the negative Fisher Information Matrix.

We adopt the integral form of the KL divergence to derive the KL divergence after a small perturbation of the model. Assume the output distribution of the model is $p(x) = f(x; \theta)$, the output after perturbation is $q(x) = f(x; \theta')$, where $\theta'$ is a small perturbation w.r.t $\theta$:

$$D_{\mathrm{KL}}(p \| q) = \int_{\mathbb{R}} f(x; \theta) \log \frac{f(x; \theta)}{f(x; \theta')} \, \mathrm{d}x. \quad (24)$$

Similar to the definition of FIM [11], under the regularity condition, the Hessian matrix of KL divergence can be written as:

$$\frac{\partial^2}{\partial \theta'_i \partial \theta'_j} D_{\mathrm{KL}}(p \| q) = - \int_{\mathbb{R}} f(x; \theta) \left( \frac{\partial^2 \log f(x; \theta')}{\partial \theta'_i \partial \theta'_j} \right) \mathrm{d}x. \quad (25)$$

It can be observed that when $f(x; \theta)$ matches the true data distribution, it can be regarded as the probability density function of the true data distribution. Thus, the Hessian matrix of KL divergence is equal to the expectation of the Hessian matrix of the log-likelihood of $f(x, \theta')$. When $\theta$ and $\theta'$ are sufficiently close, the Hessian matrix of the KL divergence is the expectation of the Hessian matrix of the model's log-likelihood function. $\square$

### A.2. Proof of Theorem 3.1

In order to prove Theorem 3.1, we begin with the definition of the score function.

**Definition 1.** *The Fisher Information Matrix is defined as the variance of the score function, where the score function is the gradient of the log-likelihood function.*

**Theorem A.2.** *When the model's output distribution matches the true distribution, the expected value of the score function becomes $0$.*

*Proof.* According to the definition of the score function [11], we have:

$$
\begin{aligned}
\mathbb{E}\left[\frac{\partial}{\partial \theta} \log f(X; \theta) \middle| \theta\right] &= \int_{\mathbb{R}} \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} f(x; \theta) \mathrm{d}x \\
&= \frac{\partial}{\partial \theta} \int_{\mathbb{R}} f(x; \theta) \mathrm{d}x \\
&= \frac{\partial}{\partial \theta} 1 \\
&= 0,
\end{aligned} \quad (26)
$$

where the likelihood function $f(X; \theta)$ denotes the probability of the model output random variable $X$, and $f(x; \theta)$ denotes the probability density of $X$ at $x$. This equation holds if and only if the output distribution matches the true distribution, allowing the use of $f(x; \theta)$ as the probability density function for integration. $\square$

Based on the conclusion above, we prove Theorem 3.1 as follows.

*Proof.* Based on the definition of the Fisher Information Matrix and the definition of variance $D(X) = E(X^2) - E^2(X)$, when the expected gradient of the log-likelihood is $0$, we have:

$$\mathbf{F}(\theta) = \mathbb{D}\left[\frac{\partial}{\partial \theta} \log f(X; \theta) \middle| \theta\right] = \mathbb{E}\left[\left(\frac{\partial}{\partial \theta} \log f(X; \theta)\right)^2 \middle| \theta\right]. \quad (27)$$

The second derivative of the log-likelihood function with respect to the parameters (*i.e.*, the Hessian matrix) is:

$$
\begin{aligned}
\frac{\partial^2}{\partial \theta^2} \log f(X;\theta) &= \frac{\partial}{\partial \theta} \frac{\frac{\partial}{\partial \theta} f(X;\theta)}{f(X;\theta)} \\
&= \frac{\left(\frac{\partial^2}{\partial \theta^2} f(X;\theta)\right) \cdot f(X;\theta) - \left(\frac{\partial}{\partial \theta} f(X;\theta)\right)^2}{f(X;\theta)^2} \\
&= \frac{\frac{\partial^2}{\partial \theta^2} f(X;\theta)}{f(X;\theta)} - \left(\frac{\frac{\partial}{\partial \theta} f(X;\theta)}{f(X;\theta)}\right)^2 \\
&= \frac{\frac{\partial^2}{\partial \theta^2} f(X;\theta)}{f(X;\theta)} - \left(\frac{\partial}{\partial \theta} \log f(X;\theta)\right)^2,
\end{aligned}
\tag{28}
$$

where the second term is the definition of FIM, and under the regularity condition, the expectation of the first term is 0:

$$
\begin{aligned}
\mathbb{E}\left[\frac{\frac{\partial^2}{\partial \theta^2} f(X;\theta)}{f(X;\theta)}\,\middle|\,\theta\right] &= \int_{\mathbb{R}} \frac{\frac{\partial^2}{\partial \theta^2} f(x;\theta)}{f(x;\theta)} f(x;\theta)\mathrm{d}x \\
&= \frac{\partial^2}{\partial \theta^2}\int_{\mathbb{R}} f(x;\theta)\mathrm{d}x = 0.
\end{aligned}
\tag{29}
$$

Therefore, when the model's output distribution matches the true distribution, the Fisher Information Matrix is equivalent to the expectation of the negative second derivative of the log-likelihood function, i.e., the expectation of the Hessian matrix of the negative log-likelihood function. □

### A.3. Proof of Theorem 3.2

*Proof.* In the context of block-wise post-training quantization, we regard the KL divergence as a function of the perturbation to the block output:

$$
\begin{aligned}
\mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) &= D_{\mathrm{KL}}(p(x;\mathbf{z}^{(b)})\|p(x;\mathbf{z}^{(b)}+\Delta\mathbf{z}^{(b)})) \\
&= \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)}) \log \frac{p(x;\mathbf{z}^{(b)})}{p(x;\mathbf{z}^{(b)}+\Delta\mathbf{z}^{(b)})} \\
&= \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)}) \log p(x;\mathbf{z}^{(b)}) \\
&\quad - \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)}) \log p(x;\mathbf{z}^{(b)}+\Delta\mathbf{z}^{(b)}).
\end{aligned}
\tag{30}
$$

We perform a second order Taylor expansion to $\log p(x;\mathbf{z}^{(b)}+\Delta\mathbf{z}^{(b)})$ as below:

$$
\begin{aligned}
&\log p(x;\mathbf{z}^{(b)}+\Delta\mathbf{z}^{(b)}) \\
&= \log p(x;\mathbf{z}^{(b)}) + \nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})^\top \Delta\mathbf{z}^{(b)} \\
&\quad + \frac{1}{2}\Delta\mathbf{z}^{(b)\top}\nabla^2_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})\Delta\mathbf{z}^{(b)}
\end{aligned}
\tag{31}
$$

Thus, we have

$$
\mathcal{L}_{\mathrm{KL}}(\Delta\mathbf{z}^{(b)}) = -\Delta\mathbf{z}^{(b)\top}\cdot S_1 - \frac{1}{2}\Delta\mathbf{z}^{(b)\top}\cdot S_2 \cdot \Delta\mathbf{z}^{(b)}, \tag{32}
$$

where

$$
\begin{aligned}
S_1 &= \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)}) \\
S_2 &= \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)})\nabla^2_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)}).
\end{aligned}
\tag{33}
$$

According to the properties of logarithmic function differentiation, we can deduce the following:

$$
\nabla_{\mathbf{z}^{(b)}} p(x;\mathbf{z}^{(b)}) = p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)}). \tag{34}
$$

Thus, under the regularity condition, we have:

$$
\begin{aligned}
S_1 &= \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)}) \\
&= \int_{\mathbb{R}} \nabla_{\mathbf{z}^{(b)}} p(x;\mathbf{z}^{(b)}) \\
&= \nabla_{\mathbf{z}^{(b)}} \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)}) \\
&= \nabla_{\mathbf{z}^{(b)}} 1 \\
&= 0.
\end{aligned}
\tag{35}
$$

For $S_2$, according to Eq. (34), the following equations hold:

$$
\begin{aligned}
&\nabla^2_{\mathbf{z}^{(b)}} p(x;\mathbf{z}^{(b)}) \\
&= \nabla_{\mathbf{z}^{(b)}}(p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})) \\
&= \nabla_{\mathbf{z}^{(b)}} p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})^\top \\
&\quad + p(x;\mathbf{z}^{(b)})\nabla^2_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)}) \\
&= p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})^\top \\
&\quad + p(x;\mathbf{z}^{(b)})\nabla^2_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)}).
\end{aligned}
\tag{36}
$$

Therefore, $S_2$ can be written as

$$
\begin{aligned}
S_2 &= \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)})\nabla^2_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)}) \\
&= \int_{\mathbb{R}} \nabla^2_{\mathbf{z}^{(b)}} p(x;\mathbf{z}^{(b)}) \\
&\quad - \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})^\top.
\end{aligned}
\tag{37}
$$

Under the regularity condition, we can derive that

$$
\begin{aligned}
\int_{\mathbb{R}} \nabla^2_{\mathbf{z}^{(b)}} p(x;\mathbf{z}^{(b)}) &= \nabla^2_{\mathbf{z}^{(b)}} \int_{\mathbb{R}} p(x;\mathbf{z}^{(b)}) \\
&= \nabla^2_{\mathbf{z}^{(b)}} 1 \\
&= 0.
\end{aligned}
\tag{38}
$$

Thus,

$$
\begin{aligned}
S_2 &= -\int_{\mathbb{R}} p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})\nabla_{\mathbf{z}^{(b)}} \log p(x;\mathbf{z}^{(b)})^\top \\
&= -\mathbf{F}^{(\mathbf{z}^{(b)})}.
\end{aligned}
\tag{39}
$$

By substituting $S_1$ and $S_2$ into Eq. (32), we have:

$$\mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) = \frac{1}{2} \Delta \mathbf{z}^{(b)\top} \mathbf{F}^{(\mathbf{z}^{(b)})} \Delta \mathbf{z}^{(b)}. \qquad (40)$$

$\square$

### A.4. Derivation of Eq. (15)

Given

$$\mathbf{F}^{(\mathbf{z}^{(b)})} = \boldsymbol{u} \boldsymbol{u}^\top, \qquad (41)$$

$$\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) = \mathbf{F}^{(\mathbf{z}^{(b)})} \Delta \mathbf{z}^{(b)}, \qquad (42)$$

where $\boldsymbol{u}, \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}), \Delta \mathbf{z}^{(b)} \in \mathbb{R}^{a \times 1}$, we define a scalar $\alpha = \boldsymbol{u}^\top \cdot \Delta \mathbf{z}^{(b)}$ such that

$$\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) = \alpha \boldsymbol{u}. \qquad (43)$$

Then, we can deduce the following

$$\boldsymbol{u}^\top = \frac{\left( \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) \right)^\top}{\alpha}. \qquad (44)$$

Thus,

$$\alpha = \sqrt{\left( \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) \right)^\top \Delta \mathbf{z}^{(b)}}, \qquad (45)$$

$$\boldsymbol{u} = \frac{\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})}{\sqrt{\left( \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) \right)^\top \Delta \mathbf{z}^{(b)}}}. \qquad (46)$$

### A.5. Proof of Corollary 3.1

*Proof.* Given

$$\mathbf{F}^{(\mathbf{z}^{(b)})} = \boldsymbol{u} \boldsymbol{u}^\top, \qquad (47)$$

$$\nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) = \mathbf{F}^{(\mathbf{z}^{(b)})} \Delta \mathbf{z}^{(b)}, \qquad (48)$$

we can deduce the following

$$\Delta \mathbf{z}^{(b)\top} \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) = \Delta \mathbf{z}^{(b)\top} \boldsymbol{u} \boldsymbol{u}^\top \Delta \mathbf{z}^{(b)}. \qquad (49)$$

Since the right-hand side of Eq. (49) is a symmetric matrix, the left-hand side should also be symmetric:

$$\Delta \mathbf{z}^{(b)\top} \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) = \left( \nabla \mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)}) \right)^\top \Delta \mathbf{z}^{(b)}. \quad (50)$$

When $k = 1$, both sides of Eq. (50) are scalars, implying that Eq. (50) naturally holds. When $k > 1$, since $\Delta \mathbf{z}^{(b)}$ and $\mathcal{L}_{\mathrm{KL}}(\Delta \mathbf{z}^{(b)})$ are not directly related, we cannot guarantee their symmetry.

As a consequence, it is generally difficult to find a $\boldsymbol{u}$ such that $\mathbf{F}^{(\mathbf{z}^{(b)})} = \boldsymbol{u} \boldsymbol{u}^\top$ satisfying Eq. (10) in most cases. $\square$

## B. More Experiments

As both FIM approximation (FIMA) and reconstruction steps depend on calibration data, we separately evaluate their performance utilizing different numbers of samples. As shown in Table A, the accuracy using FIMA increases as sample size grows, but is generally robust to the sample size. However, the reconstruction step is more sensitive to the number of calibration samples.

Table A. Ablation results (%) w.r.t. the samples size with W3/A3 on ImageNet.

| Sample Size | In FIMA Step | | | In Reconstruction Step | | |
|---|---|---|---|---|---|---|
| | ViT-S | DeiT-S | Swin-S | ViT-S | DeiT-S | Swin-S |
| **128** | 63.52 | 68.99 | 77.10 | 49.64 | 64.12 | 71.71 |
| **256** | 63.18 | 69.10 | 77.00 | 56.02 | 66.21 | 74.12 |
| **512** | 63.61 | 69.14 | 77.18 | 60.45 | 67.87 | 75.8 |
| **1024** | 64.09 | 69.13 | 77.26 | 64.09 | 69.13 | 77.26 |

Since the Fisher Information Matrix (FIM) captures global information, its computation involves averaging over the sample dimension. Theoretically, a larger sample size leads to a more accurate approximation due to reduced sampling error. However, since the averaging process mitigates the impact of individual sample variations, the difference is not particularly significant. In fact, even using a single sample for approximation can still yield an acceptable level of accuracy. However, as shown in Tab. A, directly altering the overall sample size leads to a more substantial accuracy change, as the reconstruction process in Adaround is more sensitive to the number of samples.