# GPLQ: A General, Practical, and Lightning QAT Method for Vision Transformers

**Guang Liang**[1,2]    **Xinyao Liu**[3]    **Jianxin Wu**[1,2*]

[1]State Key Laboratory for Novel Software Technology, Nanjing University, China
[2]School of Artificial Intelligence, Nanjing University, China
[3]University of Science and Technology of China, Hefei, China

liangg@lamda.nju.edu.cn, liuxinyao@mail.ustc.edu.cn, wujx2001@nju.edu.cn

## Abstract

Vision Transformers (ViTs) are essential in computer vision but are computationally intensive, too. Model quantization, particularly to low bit-widths like 4-bit, aims to alleviate this difficulty, yet existing Post-Training Quantization (PTQ) and Quantization-Aware Training (QAT) methods exhibit significant limitations. PTQ often incurs substantial accuracy drop, while QAT achieves high accuracy but suffers from prohibitive computational costs, limited generalization to downstream tasks, training instability, and lacking of open-source codebase. To address these challenges, this paper introduces General, Practical, and Lightning Quantization (GPLQ), a novel framework designed for efficient and effective ViT quantization. GPLQ is founded on two key empirical insights: the paramount importance of activation quantization and the necessity of preserving the model's original optimization "basin" to maintain generalization. Consequently, GPLQ employs a sequential "activation-first, weights-later" strategy. Stage 1 keeps weights in FP32 while quantizing activations with a feature mimicking loss in only 1 epoch to keep it stay in the same "basin", thereby preserving generalization. Stage 2 quantizes weights using a PTQ method. As a result, GPLQ is 100x faster than existing QAT methods, lowers memory footprint to levels even below FP32 training, and achieves 4-bit model performance that is highly competitive with FP32 models in terms of both accuracy on ImageNet and generalization to diverse downstream tasks, including fine-grained visual classification and object detection. We will release an easy-to-use open-source toolkit supporting multiple vision tasks.

## 1 Introduction

Vision Transformer (ViT) [6, 33] has emerged as the mainstream backbone network in computer vision, but it demands substantial computational and memory resources. Model quantization is one of the key techniques to address this challenge by reducing the numerical precision of model parameters and/or activation values[18, 19]. However, existing quantization methods still faces challenges, especially in low-bit (e.g., 4-bit) quantization.

Mainstream methods include Post-Training Quantization (PTQ) [25] and Quantization-Aware Training (QAT) [7]. PTQ has fast speed and low resource consumption, but often leads to large accuracy drop under 4-bit quantization [20]. On the other hand, QAT simulates quantization operations during training and enables higher accuracy than that of PTQ, or *even higher than that of floating-point models*. Nevertheless, in this paper we will show that existing QAT methods have inherent limitations:
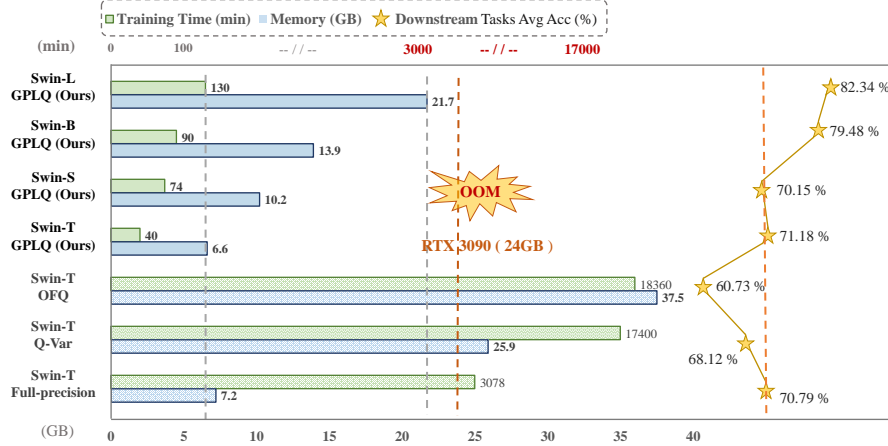
---

*Corresponding author.

Figure 1: Core advantages of our GPLQ: Generality, Practicality, and Lightning efficiency.

- **High Computational Costs.** QAT requires lengthy fine-tuning of the entire model. Training time and GPU memory required in QAT often *far exceed those for training the FP32 model* [18]. This makes QAT *cumbersome and very slow* for deployment in real-world applications.

- **Limited Generalization Ability.** QAT methods often boast higher accuracy than their FP32 counterparts. However, in this paper we will show that such models are *generalizing worse* than FP32 or PTQ quantized models in downstream tasks. That is, they are likely *non-generalizable* beyond ImageNet [4], the dataset on which they were trained.

- **Training Instability and Complexity.** QAT is prone to training instability [15], and complex Knowledge Distillation (KD) techniques [19, 15] severely increase memory footprint. Some also rely on external, extremely powerful teacher models, which are *not* available in practical scenarios. In short, existing QAT methods are *not practical*.

- **Classification Only and Code Missing.** Open-source code for QAT is rare, and is only for classification when it exists. This further makes QAT *impractical* for real-world applications.

To this end, we propose GPLQ (General, Practical, and Lightning Quantization). The core objective of GPLQ is to provide a quantization solution that is far more training-efficient than traditional QAT, superior to PTQ in accuracy and generalization, easy to use, and highly practical. As a result, Figure 1 demonstrates 3 core advantages of GPLQ.

- **General.** GPLQ exhibits excellent average accuracy on multiple *downstream tasks*: close to or even surpassing FP32 models, and significantly outperforming existing QAT methods.

- **Practical.** GPLQ has *very small training memory footprint* (far lower than existing QAT methods), which avoids out-of-memory (OOM) issues in many applications and enables quantization of larger models. GPLQ's design allows it to be conveniently applied to *other tasks such as object detection*.

- **Lightning.** GPLQ is blazingly fast: *hundreds of times faster* than existing QAT methods.

GPLQ is based on our empirical findings. First, *activations are far more important than weights* in low-bit quantization. Second, quantization should not change its optimization "basin" (i.e., *avoid jumping out of the current local minimum*) in order to keep the generalization ability.

Based on these findings, GPLQ adopts a sequential quantization paradigm. First, activations are quantized with weights kept at FP32. To maintain generalization, we draw inspiration from TCS [37] and employ a PCA-based feature mimicking loss to guide the quantized model's feature outputs to approximate those of the original FP32 model (i.e., stay in the same basin). Second, after activations are quantized, existing efficient PTQ methods are used to quantize the weights. This "activation-first, weights later" strategy not only drastically reduces QAT *training time from days to 1-2 hours* and

with *memory footprint even lower than FP32 training*, but also allows a 4-bit model to achieve both accuracy and generalization nearly identical to the original FP32 model. The main contributions are:

1. **Insights.** We reveal that activation quantization is the main bottleneck in QAT, and staying in the original optimization basin is crucial for generalization.

2. **GPLQ.** We propose "activation-first" sequential quantization: first optimize activations then quantize weights via PTQ.

3. **Code.** GPLQ provides an easy-to-use quantization tool supporting classification, detection and other downstream tasks. We will open-source GPLQ upon paper acceptance.

## 2 Related Work

Model quantization aims to enhance model efficiency by reducing the numerical precision of weights and activations in neural networks [29].

**Post-Training Quantization (PTQ)**. PTQ operates without retraining, requires only a small calibration set, and is very fast. Various techniques have been proposed: AIQViT [16], GPTQ [8], PTQ4ViT [36]), SmoothQuant [35], AWQ [21]. More methods like RepQ-ViT [20] and QwT [9] perform optimization through scale reparameterization and lightweight compensation modules, respectively. Accuracy degradation remains a severe challenge in low-bit scenarios.

The second stage of GPLQ employs PTQ to quantize weights. Since activations have been quantized via QAT in the first stage, PTQ's duty changes from W32A32 → W4A4 to W32A4 → W4A4.

**Quantization-Aware Training (QAT)**. QAT introduces simulated quantization during training or fine-tuning, and achieves higher accuracy than PTQ. It often uses a Straight-Through Estimator (STE) to handle gradient [7]. Research directions include learning quantization scales [7], improving training stability) (OFQ [23], Quantization-Variation [15]), enhancing efficiency (EfficientQAT [3]), and specific optimizations for ViTs (e.g., Q-ViT [19], PackQViT [5]).

The bottlenecks of QAT are high computational cost, training instability, and potential degradation in generalization ability. GPLQ, with an extremely short QAT stage (*only 1 epoch*) focused solely on activations, effectively alleviates the cost, stability, and generalization issues of traditional QAT.

**Knowledge Distillation (KD) in QAT**. KD [14, 11, 34, 38] is often used in QAT to learn from the FP32 model or even a much stronger, external teacher. Researchers have proposed methods (DGD [36], MCKD [15]) that are both heavy and complex. TCS [37] offers an efficient approach by capturing the linear subspace of the teacher model's features through Principal Component Analysis (PCA) for knowledge transfer. GPLQ draws inspiration from TCS, aiming at efficiently transferring knowledge to maintain generalization and avoiding the high costs of complex distillation.

## 3 Methodology

GPLQ is directly derived from two empirical insights in the ViT quantization process.

### 3.1 Activations are Crucial & Stay in the Same Basin

Our first empirical insight is that *activations are more critical than weights* in quantization. Taking a ViT network such as DeiT [31] pre-trained on ImageNet-1k [4], we independently applied 4-bit PTQ (using a percentile-based per-channel quantization calibration method) to either weights (with activations kept at FP32) or activations (with weights kept at FP32). Figure 2 shows the results, which consistently indicates that quantizing activations to 4 bits (weights at FP32) leads to larger Top-1 accuracy drop than quantizing weights to 4 bits (activations at FP32). That is, activations face more severe challenges under low-bit quantization compared to weights. This finding leads to our "activation-first" quantization strategy in GPLQ.

QAT methods heavily adjust their weights, and make them stay in dramatically different basins (local minima) in the loss landscape before and after QAT learning. That is, they are significantly different from the initial FP32 model. Although they achieve high accuracy on the pre-training ImageNet data, we observe that this aggressive retraining weakens the transferability of the learned representations
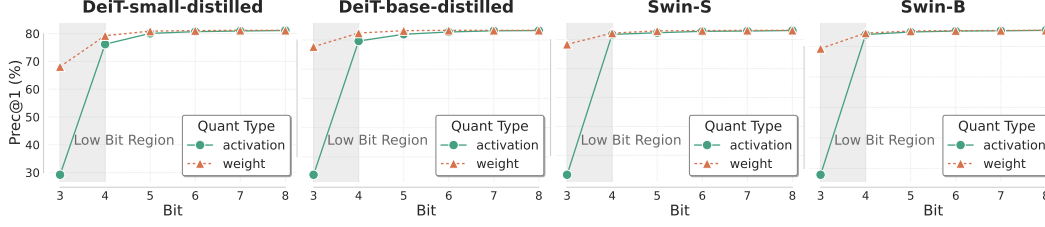
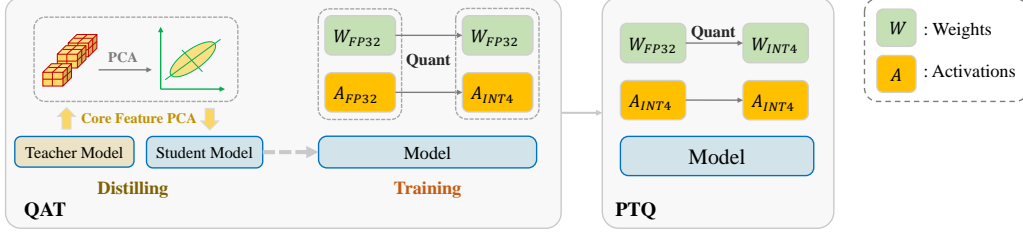Figure 2: Impact of quantizing weights and activations separately.



Figure 3: Overview of GPLQ: QAT stage first only for activations, then PTQ stage only for weights.

to downstream tasks (cf. Table 1 and Figure 1). Our second insight is to *make the quantized model stay in the same basin*. This is achieved in GPLQ by i) restricting QAT to optimize only activation quantizers; ii) use a low learning rate and only 1 epoch training in this QAT stage; iii) a feature mimicking loss that encourages the quantized model to retain the key feature structures of the original FP32 model.

## 3.2 The GPLQ Framework

### 3.2.1 Stage 1: Activation-only Quantization-Aware Training, or Act-QAT

This stage only quantizes activations. The key is that all model weights are *learnable but kept at FP32 precision*, thereby decoupling activation quantizer learning from weight updates and effectively circumventing weight oscillation that occurs in QAT which quantizes the two simultaneously.

Activations employ the uniform symmetric quantization implemented with per-channel granularity. A comparison of per-tensor versus per-group quantization is detailed in the appendix. Activations are quantized to 4 bits ($b = 4$) in our experiments, but our method and code work for lower-bit quantization, too. The quantization scaling factor $s_a$ is learned using LSQ [7]. One *novel proposal in GPLQ* is that we initialize $s_a$ and calibrate it on a small subset of training data using a percentile-based min-max PTQ method. For internal per-token quantized operations (if applicable), the initial quantization range is set based on the 1st and 99th percentiles of observed activation values to mitigate outlier effects; for per-channel quantization, min-max values are used directly. This initialization process is fast (typically in seconds) and provides a better starting point for subsequent LSQ-like optimization. The quantized activation $\hat{x}$ for an input $x$ is calculated as follows:

$$\hat{x} = \text{clamp}(\text{round}(x/s_a), -2^{b-1}, 2^{b-1} - 1) \times s_a \,, \tag{1}$$

where $b = 4$ in our experiments. The gradient of $s_a$ is estimated using Straight-Through Estimator (STE).

To preserve the rich representational power of pre-trained models (i.e., keep the generalization power), we follow a lightweight feature mimicking loss inspired by TCS [37]. We first extract features $f_t$ from the penultimate layer of the original FP32 teacher model. Then, PCA is performed on a set of teacher features $F_t = \{f_{t_i}\}$ obtained from a subset of the training data. The resulting principal components $V$ define a low-dimensional subspace that captures the main information in the teacher's feature manifold. We select principal components that explain a majority of the variance in $F_t$ (e.g., approximately 60%), and for hardware friendliness, adjust the number of selected principal components to be a multiple of 32. Specifically, for Swin-T [24] with 768 dimensions, 256 are selected; for DeiT-T [31] with 192 dimensions, 64 are selected. The corresponding features $f_s$ of the

4

student model (the model undergoing Act-QAT) are projected onto this PCA-defined subspace. The loss $L_{PCA}$ is defined as the Mean Squared Error (MSE) between the student's projected features and the teacher's projected features, i.e., matching projections in the PCA space:

$$L_{PCA} = \frac{1}{N} \sum_{i=1}^{N} ||(f_s^i - \mu_t)V_{sel} - (f_t^i - \mu_t)V_{sel}||_2^2, \tag{2}$$

where $\mu_t$ is the mean of the teacher features in $F_t$, $V_{sel}$ are the selected principal components, and $N$ is the number of samples (e.g., batch size). This loss encourages the post-quantization activations to compactly retain salient features of the FP32 teacher model, thereby enhancing generalization.

### 3.2.2 Stage 2: Post-Training Quantization of Weights, or Weight-PTQ

Now the model has FP32 weights and 4-bit quantized activations (W32A4). The second stage rapidly quantizes the weights to 4 bits using mature PTQ techniques to generate the final W4A4 model. We leverage existing efficient PTQ methods in this stage. Specifically, after quantizing the weights using RepQ-ViT [20], we further apply the QwT (Quantization without Tears) [9] method to compensate for the accuracy loss introduced by weight quantization. The entire calibration process is completed on a small, randomly selected subset of the ImageNet training data.

With activation quantizers *frozen*, we significantly simplify PTQ of weights. Since activations are already fixed and quantized to 4 bits (A4), error coupling that hurts weight PTQ algorithm is greatly reduced. The primary source of error is weight quantization itself (W32 $\rightarrow$ W4), rather than the compound error in simultaneous W32 $\rightarrow$ W4 and A32 $\rightarrow$ A4 conversions in traditional PTQ. This simplified objective makes the compensation technique (QwT) easier to perform. QwT corrects quantization error by introducing a lightweight linear compensation layer, whose parameters $W^*$ are determined by the following closed-form solution:

$$W^* = (Y - Y_Z)X_Z^T(X_Z X_Z^T + \lambda I)^{-1}, \tag{3}$$

which includes a regularization term $\lambda I$ for stability. In our W32A4 setting, $X_Z$ represents the input activations that have already been quantized to 4 bits, $Y$ is the output of the layer with FP32 weights and A4 inputs ($Y_{W32A4}$), while $Y_Z$ is the output with 4-bit weights and A4 inputs ($Y_{W4A4}$).

### 3.3 Advantages of Our GPLQ

GPLQ is designed with practical ease of use in mind, and we provide an implementation that encapsulates the two-stage process into simple and easy-to-use code.

The principles and effectiveness of GPLQ are not limited to image classification. We extended the code to object detection on MS-COCO [22]. For COCO tasks, Act-QAT (Stage 1) is also performed for 1 epoch, but due to larger input resolutions and higher model complexity the batch size per GPU is adjusted to 1. The subsequent weight PTQ (Stage 2) follows a similar procedure to that for image classification. This demonstrates the good generality of our framework across diverse vision tasks.

Compared to traditional QAT and PTQ methods, GPLQ offers an attractive alternative. Requiring only 1 epoch of activation QAT, its training duration and required computational resources are far less than those of a typical full QAT process (e.g., *hundreds of times faster*). Furthermore, pre-quantized activations create a more tractable optimization problem for subsequent weight PTQ, thereby enhancing the effectiveness of methods like QwT and the weight quantization in RepQ-ViT. Finally, our design avoids jumping out of the FP32 model's local minima, thus is useful for preserving the generalization ability of the pre-trained FP32 model.

## 4 Experiments

We conducted a comprehensive evaluation of GPLQ on multiple benchmark datasets and vision tasks. For image classification, we use ImageNet-1k [4] for pre-training and primary performance evaluation. For object detection and instance segmentation tasks, we employ the COCO 2017 [22] dataset, by training models on the 'train2017' set and reporting performance on the 'val2017' set.

To evaluate the model's generalization ability, we also selected five commonly used Fine-Grained Visual Classification (FGVC) datasets, including Aircraft [27], Food101 [2], Flowers102 [28],

Table 1: Comparison (Top-1 accuracy in percentage) between GPLQ and SOTA QAT methods on ImageNet-1k and 5 downstream FGVC tasks. 'Avg Task' is the average accuracy on 5 FGVC tasks.

| Network | Method | Mode | ImageNet | Aircraft | Food101 | Flowers102 | Pets | Cars | Avg Task |
|---|---|---|---|---|---|---|---|---|---|
| Swin-T | FP32 | W32A32 | 81.2 | 39.72 | 73.85 | 91.10 | 93.21 | 56.06 | 70.79 |
| | OFQ [23] | W4A4 | 81.9 | 26.58 | 64.79 | 84.40 | 91.74 | 36.13 | 60.73 |
| | Q-Var [15] | W4A4 | **82.4** | 37.02 | 70.98 | 87.15 | **92.86** | 52.57 | 68.12 |
| | RepQ-ViT [20] | W4A4 | 73.0 | 35.46 | 60.59 | 86.83 | 88.74 | 42.33 | 62.79 |
| | **GPLQ** | W4A4 | 79.8 | **43.50** | **71.02** | **92.05** | 92.31 | **57.02** | **71.18** |
| Swin-S | FP32 | W32A32 | 83.2 | 38.13 | 72.63 | 91.32 | 92.80 | 55.48 | 70.07 |
| | RepQ-ViT [20] | W4A4 | 71.9 | 32.16 | 69.44 | 90.40 | 93.08 | 49.27 | 66.87 |
| | **GPLQ** | W4A4 | **81.9** | **39.03** | **71.81** | 90.13 | **93.24** | **56.54** | **70.15** |
| Swin-B | FP32 | W32A32 | 85.3 | 49.71 | 85.12 | 99.48 | 94.49 | 65.86 | 78.93 |
| | RepQ-ViT [20] | W4A4 | 69.0 | 44.85 | 63.44 | 95.74 | 89.13 | 59.25 | 70.48 |
| | **GPLQ** | W4A4 | **84.2** | **52.93** | **83.81** | 99.40 | **93.87** | **67.39** | **79.48** |
| Swin-L | FP32 | W32A32 | 86.3 | 51.52 | 87.10 | 99.63 | 94.93 | 71.25 | 80.89 |
| | RepQ-ViT [20] | W4A4 | 83.2 | 52.66 | 84.84 | 99.46 | 94.36 | 67.31 | 79.73 |
| | **GPLQ** | W4A4 | **85.5** | **57.52** | **86.21** | 99.66 | **94.55** | **73.77** | **82.34** |
| DeiT-S | FP32 | W32A32 | 81.2 | 34.41 | 64.92 | 87.36 | 91.85 | 50.50 | 65.81 |
| | OFQ [23] | W4A4 | **81.1** | 29.07 | **65.88** | 83.61 | **92.10** | 42.31 | 62.59 |
| | RepQ-ViT [20] | W4A4 | 72.7 | 26.55 | 57.68 | 85.28 | 89.67 | 40.59 | 59.95 |
| | QwT [9] | W4A4 | 74.8 | 35.61 | 61.23 | 87.75 | 88.93 | 48.09 | 64.32 |
| | **GPLQ** | W4A4 | 78.8 | **39.81** | 64.95 | **89.77** | 91.44 | **50.55** | **67.30** |
| DeiT-B | FP32 | W32A32 | 83.3 | 45.06 | 72.96 | 91.84 | 93.35 | 63.93 | 73.43 |
| | RepQ-ViT [20] | W4A4 | 76.3 | 48.90 | 69.47 | 93.35 | 92.56 | 62.39 | 73.33 |
| | QwT [9] | W4A4 | 78.5 | 49.23 | **73.37** | **93.85** | 92.59 | **65.96** | **75.00** |
| | **GPLQ** | W4A4 | **82.0** | **49.84** | 71.69 | 93.35 | **93.27** | 65.15 | 74.66 |

Pets [30], and Cars [17]. To fairly compare the feature extraction capabilities and downstream generalization performance, we train models using linear probing. Furthermore, for a fair comparison, the hyperparameters used for all methods were kept consistent following DTL [10]: 100 epochs, learning rate 0.001, batch size of 64, and drop path rate 0.1.

Stage 1 (Act-QAT) trained for 1 epoch on ImageNet-1k (classification) or COCO (detection) using AdamW [26] with a fixed learning rate of $5 \times 10^{-6}$ and no decay. Activations use per-channel symmetric 4-bit quantization. The subspace dimension used for PCA feature mimicking is dynamically selected based on the model's feature dimension, with the selection primarily based on accumulated variance when the accumulated variance is around 60%. Specifically, Swin-T uses a 256-dimensional PCA subspace, and DeiT-T (with 192 dimensions) uses a 64-dimensional subspace. Training was conducted on 8 GPUs, with a batch size of 16 per GPU. This configuration allows quantizing of the entire Swin Transformer series even on consumer-grade GPUs.

Stage 2 (Weight-PTQ) employed percentile-based per-channel symmetric 4-bit quantization for weights, combined with QwT [9] for compensation. The calibration set consists of 512 randomly selected images from the ImageNet training set.

## 4.1 Image Classification Performance

We evaluated GPLQ on ImageNet-1k and five downstream fine-grained classification tasks. The results are shown in Table 1, from which we can observe that:

1. **On ImageNet Itself.** GPLQ significantly outperforms PTQ methods (RepQ-ViT and QwT). On Swin-T, GPLQ is 6.8% higher than RepQ-ViT and only 1.4% lower than FP32. For DeiT-S, GPLQ also far surpasses RepQ-ViT and QwT. Compared to QAT methods (OFQ, Q-Var), GPLQ is slightly inferior to these QAT models. But, these QAT methods not only require *unacceptable training time* but *may also lead to overfitting* (which we discuss next).

2. **Downstream Task Generalization:** In terms of average downstream accuracy, GPLQ almost always surpasses other quantization methods. And it consistently achieves better

Table 2: Object detection and instance segmentation results ($AP^{box}$ / $AP^{mask}$).

| Method | Bits (W/A) | Swin-T (1x) | Swin-T (3x) | Swin-S (3x) |
|---|---|---|---|---|
| Full-Precision | 32/32 | 0.426 / 0.393 | 0.460 / 0.416 | 0.485 / 0.433 |
| PTQ4ViT [36]) | 4/4 | — / — | 0.069 / 0.070 | 0.267 / 0.266 |
| APQ-ViT [7] | 4/4 | — / — | 0.237 / 0.226 | 0.447 / 0.401 |
| RepQ-ViT [20] | 4/4 | 0.135 / 0.137 | 0.361 / 0.360 | 0.426 / 0.400 |
| GPLQ (Act-QAT only) | 32/4 | 0.397 / 0.381 | 0.430 / 0.402 | 0.457 / 0.421 |
| **GPLQ** | 4/4 | **0.379 / 0.368** | **0.401 / 0.389** | **0.434 / 0.413** |

generalization than FP32. Swin-T GPLQ's average downstream accuracy is 71.18%, higher than FP32's 70.79%, and far exceeding Q-Var (68.12%) and OFQ (60.73%). We want to emphasize that *QAT models, although exhibiting highest accuracy on ImageNet, lag behind PTQ methods in terms of generalization* (downstream accuracy). On the other hand, our *GPLQ has clearly better generalization than PTQ methods*, despite being a QAT method.

3. Due to the limited open-source availability of QAT methods, the number of compared QAT methods is small. And, QAT methods compared in this paper are only for small models, because they are out-of-memory even for medium size models. Our GPLQ *scales to large models and will be open-source*.

4. Q-Var performs better than OFQ in generalization, because it uses an external EfficientNet-L2 (88.2% accuracy on ImageNet) pre-trained on JFT-300M as a teacher. Finding an equally powerful teacher in other tasks is impossible or difficult and limits its application scenarios.

In short, our GPLQ quantizes super-fast, has both high accuracy and excellent generalization.

## 4.2 Object Detection Performance

We also evaluated GPLQ for object detection using the Mask R-CNN [12]framework on COCO 2017. Here GPLQ did not employ PCA feature mimicking. Results are shown in Table 2.

For a fair comparison, our weight quantization uses the same weight quantization method as in RepQ-ViT to complete the second stage. Although there is some degradation compared to W32A4 (our Stage 1 model), GPLQ still significantly outperforms other W4A4 methods such as RepQ-ViT, PTQ4ViT, and APQ-ViT. On Swin-T (3x), GPLQ achieves 0.401 $AP^{box}$, while RepQ-ViT only achieves 0.361. Even without using PCA feature mimicking, the core two-stage idea of GPLQ still demonstrates strong competitiveness in object detection tasks, which shows its generality.

## 4.3 Ablation Studies

We first investigate the impact of different activation quantization granularities (per-channel vs. per-layer) during the Act-QAT stage. As shown in Table 3, using per-channel activation quantization consistently outperforms per-layer quantization in both ImageNet and average downstream accuracy. Notably, even with per-layer quantization, our GPLQ can still achieve generalization performance close to the original floating-point model in most cases. For example, Swin-T (Layer-wise) achieved an 'Avg Task' of 69.2%, only 1.6% lower than FP32's 70.8%. Per-channel quantization even surpassed FP32 in all cases.

**Contribution of different components in the second stage of GPLQ.** We first use the first stage Act-QAT to obtain a W32A4 model. Then, we apply different PTQ strategies to the weights of this model. Results are shown in Table 4. It can be observed that even when only using the basic percentile-based PTQ method to directly quantize weights to 4 bits (PTQ weight), the resulting W4A4 model is clearly worse than the W32A4 model, but is already within an acceptable range, indicating that our Act-QAT lays a good foundation for subsequent weight quantization. Furthermore, when QwT is applied for compensation after weight quantization ('+QwT'), the model's accuracy loss is significantly mitigated. The W4A4 model with QwT achieved ImageNet accuracy close to W32A4, and its average downstream accuracy was even slightly higher than the original FP32 model. This demonstrates the effectiveness of our two-stage method, where activations are first stabilized, then weights are quantized and supplemented with lightweight compensation.

Table 3: Ablation on Act-QAT granularities.

| Model (Swin) | Prec. (W/A) | ImageNet Top-1 Acc. (%) | Avg Task (%) |
|---|---|---|---|
| Swin-T FP32 | W32A32 | 81.2 | 70.8 |
|   Channel wise | W4A4 | 79.8 (-1.4) | 71.2 (+0.4) |
|   Layer wise | W4A4 | 78.6 (-2.6) | 69.2 (-1.6) |
| Swin-S FP32 | W32A32 | 83.2 | 70.1 |
|   Channel wise | W4A4 | 81.9 (-1.3) | 70.2 (+0.1) |
|   Layer wise | W4A4 | 81.3 (-1.9) | 68.9 (-1.2) |
| Swin-B FP32 | W32A32 | 85.3 | 78.9 |
|   Channel wise | W4A4 | 84.2 (-1.1) | 79.5 (+0.6) |
|   Layer wise | W4A4 | 83.1 (-2.2) | 75.2 (-3.7) |
| Swin-L FP32 | W32A32 | 86.3 | 80.9 |
|   Channel wise | W4A4 | 85.5 (-0.8) | 82.3 (+1.4) |
|   Layer wise | W4A4 | 84.9 (-1.4) | 81.4 (+0.5) |

Table 4: Ablation on weight-PTQ components.

| Model (Swin) | Operation | Prec. (W/A) | ImageNet Top-1 Acc. (%) | Avg Task (%) |
|---|---|---|---|---|
| Swin-T | FP32 | W32A32 | 81.2 | 70.8 |
| | GPLQ Act-QAT | W32A4 | 80.4 (-0.8) | 72.2 (+1.4) |
| | + PTQ weight | W4A4 | 79.3 (-1.9) | 69.2 (-1.6) |
| | + QwT [9] | W4A4 | 79.8 (-1.4) | 71.2 (+0.4) |
| Swin-S | FP32 | W32A32 | 83.2 | 70.1 |
| | GPLQ Act-QAT | W32A4 | 82.3 (-0.9) | 70.9 (+0.8) |
| | + PTQ weight | W4A4 | 81.6 (-1.6) | 68.9 (-1.2) |
| | + QwT [9] | W4A4 | 81.9 (-1.3) | 70.2 (+0.1) |
| Swin-B | FP32 | W32A32 | 85.3 | 78.9 |
| | GPLQ Act-QAT | W32A4 | 84.6 (-0.7) | 80.4 (+1.5) |
| | + PTQ weight | W4A4 | 83.9 (-1.4) | 78.0 (-0.9) |
| | + QwT [9] | W4A4 | 84.2 (-1.1) | 79.5 (+0.6) |

Table 5: Optimizing basin retention.

| Medthod | Quantized model ImageNet(%) | Internal FP32 ImageNet(%) | Downstream Avg(%) |
|---|---|---|---|
| Raw FP32 | 81.2 | 81.2 | 70.79 |
| Q-Variation [15] | 82.4 | 68.9 | 68.12 |
| **GPLQ** | **79.8** | **81.1** | **71.18** |

Table 6: PCA Dimensionality.

| Method | PCA dim | ImageNet (%) | Avg Task Top-1 (%) |
|---|---|---|---|
| w/o PCA | – | 80.3 | 69.55 |
| w/ PCA | 64 | 80.4 | 70.42 |
| W/ PCA | 256 | 80.4 | 71.22 |
| W/ PCA | 512 | 80.3 | 70.80 |

**Importance of Preserving Optimization Basin for Generalization.** We designed an experiment to validate our hypothesis: if QAT methods significantly deviate from the original FP32 optimization basin, generalization ability may be impaired. We compared the FP32 model, a QAT method (using Quantization-Variation as an example), and our GPLQ. For QAT and GPLQ, we extracted their trained FP32 weights (i.e., remove quantization nodes and use the learned "latent" FP32 weights, denoted as 'Internal FP32') and evaluated their performance on ImageNet.

As shown in Table 5, although Quantization-Variation achieved a high quantized model accuracy on ImageNet (82.4%), the performance of its internal FP32 weights dropped significantly to 68.9%, far below FP32's 81.2%. This indicates that, the quantized weights have significantly deviated from the original optimization basin, potentially led to overfitting to the ImageNet train and val set, and ultimately impaired its generalization ability on downstream tasks (Avg Acc. 68.12%). Figure 4a compares the loss curves of GPLQ and traditional QAT. GPLQ exhibits a smooth convergence process without severe oscillations, indicating that it remains within the original optimization basin of the FP32 model. In contrast, QAT's loss soars to the level for a randomly initialized network (above 7.0) at the very beginning of training—it jumped out of the original local minimum and started to fit the specifics of ImageNet, thereby affecting final generalization adversely.

In contrast, GPLQ quantized model has an accuracy of 79.8%, and the accuracy of its extracted FP32 weights is 81.1%, which is very close to the original FP32 model. This fact strongly suggests that GPLQ successfully stays near its original optimization basin. This strategy not only achieves high accuracy on ImageNet, more critically, it maintains excellent downstream task generalization (Avg Acc. 71.18%), even slightly outperforming the FP32 model.

**Impact of PCA Dimensionality in Feature Distillation.** To further investigate the effectiveness of the PCA feature distillation module in GPLQ, we conducted ablation experiments specifically on the choice of PCA projection dimensionality on the Swin-T model. Except for whether PCA projection learning was used, all other settings followed GPLQ's default settings. We report its Top-1 accuracy on the ImageNet validation set and the average accuracy from linear probing on 5 downstream FGVC datasets. Results are shown in Table 6.

We observed that without PCA feature mimicking ('Without PCA'), the model's average downstream task accuracy was 69.55%. When PCA feature mimicking was introduced, even with few dimensions (64), downstream task performance improved (70.42%). As the PCA dimensionality increased (256, at which point the cumulative explained variance was about 60%), the average downstream accuracy reached 71.22%. Further increasing the dimensionality to 512 resulted in a slight drop to 70.80%, but it is still better than not using PCA. These results indicate that feature mimicking can effectively guide

(a) Comparison of Training Loss Curves      (b) Swin-T PCA Cumulative Explained Variance
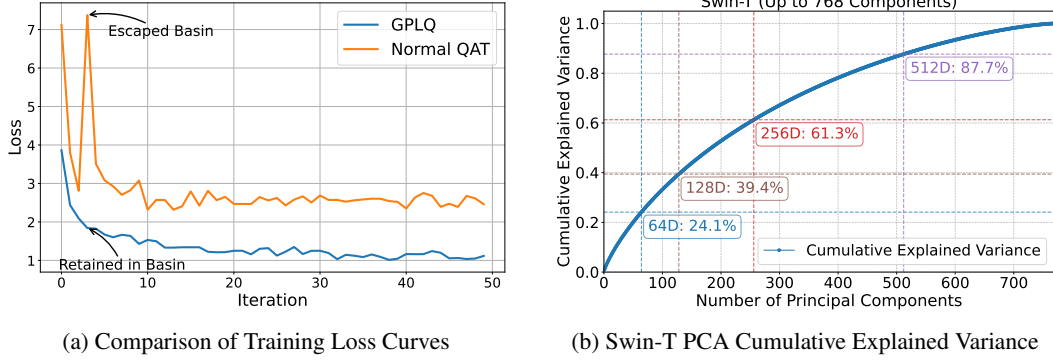
Figure 4: Training loss curves and Percentage of explained variance in GPLQ.

the quantized model to learn key features from the FP32 model, thereby enhancing its generalization ability. Selecting a dimension that captures around 60% of the original feature variance represents a good trade-off.

Figure 4b shows Swin-T model's PCA cumulative explained variance with varying dimensions. At about 256 dimensions, the cumulative explained variance reaches approximately 61.3%. Even with 64 dimensions where only 24.1% variance explained, PCA feature mimicking already shows better generalization ability than without using PCA feature mimicking, as shown in Table 6.

## 5 Conclusion

We presented GPLQ, a novel quantization framework for Vision Transformers that significantly improves efficiency and generalization over existing PTQ and QAT methods. Our approach is grounded in the empirical finding that activation quantization is critical and that preserving the original model's optimization basin is key to maintaining generalization. GPLQ's "activation-first, weights-later" strategy, featuring single-epoch activation quantization with PCA-based feature mimicking, followed by PTQ for weights, achieves 4-bit performance competitive with, and sometimes superior to, FP32 models in terms of generalization. This methodology not only drastically cuts training overhead, making advanced quantization more accessible, but also consistently outperforms prior art. GPLQ thus offers a practical and robust path for deploying low-bit ViTs in resource-constrained scenarios, facilitated by our forthcoming open-source toolkit.

## 6 Limitations and Future Work

Limitations include:

- Broader QAT comparison: Limited open-source availability of advanced QAT methods restricted a comprehensive comparative analysis.
- Dependence on PTQ techniques: GPLQ's second stage performance is tied to the capabilities and limitations of existing PTQ methods.
- Low bit-width exploration: Our work primarily focused on 4-bit quantization, mainly because 4-bit is where the hardware support ends at. Deeper investigation into even lower bit-widths is desired, too.

Future directions include:

- Extending PCA-based feature mimicking to further enhance model generalization across a broader range of vision tasks, such as object detection and semantic segmentation.
- Conducting more comprehensive and rigorous QAT benchmarks against a wider array of contemporary methods as they become publicly accessible and resources permitting.
- Adapting and evaluating the GPLQ framework for its applicability and effectiveness on diverse neural network architectures, including CNNs [13] and emerging LLMs [1, 32].

9

- Integrating next-generation PTQ advancements to continually improve GPLQ's performance.

## Acknowledgments and Disclosure of Funding

## A  Appendix

This appendix provides supplementary experimental results and analyses that further support the findings presented in the main paper.

### A.1   Detailed Performance with Varying PCA Dimensions

Due to space constraints in the main text, detailed results for downstream tasks across various datasets with different PCA dimensions were not fully elaborated. Table 7 presents these detailed outcomes for the Swin-T model under W4A4 quantization. Optimal results are highlighted in bold. As discussed in the main paper (Section 4.3, Table 6), the model generally achieves its best performance when the PCA dimension is set to 256, corresponding to a cumulative explained variance of approximately 61.3%.

Table 7: Detailed results of Swin-T(W4A4) on downstream tasks with different PCA dimensions used in the Act-QAT stage. ImageNet Top-1 accuracy (%) and cumulative explained variance (%) by the selected principal components are also shown. Best results on downstream tasks are highlighted in **bold**.

| Setting | PCA dim | ImageNet (%) | Var. (%) | Downstream Task Accuracy (%) | | | | | Avg Task Acc (%) |
|---------|---------|--------------|----------|---------|---------|------------|------|-------|------------------|
| | | | | Aircraft | Food101 | Flowers102 | Pets | Cars | |
| w/o PCA | - | 80.3 | - | 42.27 | 69.52 | 91.14 | 91.11 | 53.67 | 69.54 |
| w/ PCA | 64 | 80.4 | 24.1 | 42.57 | 70.37 | **92.36** | 91.99 | 54.74 | 70.41 |
| w/ PCA | 256 | **80.4** | 61.3 | **43.50** | **71.02** | 92.05 | **92.31** | **57.02** | **71.18** |
| w/ PCA | 512 | 80.3 | 87.7 | 42.51 | 70.25 | 92.00 | 92.67 | 56.01 | 70.69 |

Consistent with the conclusions in the main text, the Swin-T model demonstrates peak performance when the PCA dimension is 256, capturing around 60% of the cumulative variance. This configuration yields the best balance for generalization across the evaluated downstream tasks.

### A.2   Quantization under Constrained Computational Resources

To assess the robustness and efficiency of our GPLQ method under more restrictive computational environments, we investigated the impact of varying the number of available GPUs during the activation quantization stage (Act-QAT). For these experiments, only activations were quantized for 1 epoch. The batch size per GPU was maintained at 16. Consequently, a reduction in the number of GPUs corresponds to a proportional decrease in the effective batch size and the learning rate was adjusted accordingly. Other training parameters remained consistent with the settings described in the main paper.

Table 8 details the performance of Swin-T on ImageNet and the associated training times.

The results indicate that GPLQ maintains high accuracy on ImageNet even as the number of GPUs and, consequently, the effective batch size and learning rate, are significantly reduced. The performance remains remarkably stable (80.3-80.4% top-1 accuracy) across all tested configurations. This resilience suggests that our method can effectively train quantized models even with minimal training resources, a capability not typically demonstrated by traditional QAT methods. Traditional QAT

Table 8: Impact of varying GPU counts on Swin-T (W4A4 Act-QAT only) ImageNet accuracy and training time. Batch size per GPU is 16.

| Number of GPUs | Equivalent Batch Size | Learning Rate | ImageNet Acc (%) | Time (min) |
|---|---|---|---|---|
| 8 | 128 (16×8) | 5e-6 | 80.4 | 35 |
| 4 | 64 (16×4) | 2.5e-6 | 80.3 | 71 |
| 2 | 32 (16×2) | 1e-6 | 80.3 | 139 |
| 1 | 16 (16×1) | 1e-6 | 80.3 | 275 |

approaches often rely on large batch sizes and learning rates for stable convergence, and their performance is expected to degrade under such resource-constrained conditions. Our findings underscore the practical advantage of GPLQ in scenarios with limited hardware availability.

## A.3 Impact of Training Data Volume on Model Performance

We further analyzed the influence of the training data volume on the Swin-T model's performance during the 1-epoch activation quantization stage. All training settings were kept consistent with the GPLQ defaults, except for the number of training images used, which was varied by controlling the number of training iterations.

Table 9: Impact of training data volume (number of images used in 1 epoch of Act-QAT) on Swin-T ImageNet Top-1 accuracy (%) for W32A4 (only activations QAT) and W4A4 (activations QAT, weights PTQ) settings.

| Training Iterations | Images Used | W32A4 Acc (%) | W4A4 Acc (%) |
|---|---|---|---|
| 1 | 128 | 70.9 | 68.5 |
| 10 | 1,280 | 74.1 | 72.5 |
| 100 | 12,800 | 77.9 | 77.0 |
| 1000 | 128,000 | 79.9 | 79.2 |
| 10009 | Full Dataset (approx. 1.28M) | 80.4 | 79.8 |

As expected, increasing the volume of training data generally leads to improved model accuracy. Even with a relatively small number of images (e.g., 128,000, corresponding to about 10% of the full ImageNet training set for 1 epoch), the model achieves a respectable 79.2% accuracy. Training on the full dataset for 1 epoch yields 79.8% accuracy for the W32A4 model (output of Act-QAT stage), which forms a strong basis for the subsequent weight PTQ stage.

## A.4 Impact of Direct W4A4 Training versus Sequential Quantization

To further highlight the benefits of our proposed sequential "activation-first, weights-later" (W32A4 → W4A4) strategy, we compare it against a more direct approach where both weights and activations are quantized simultaneously from the start of the 1-epoch QAT process (direct W4A4). All other training hyperparameters for this direct W4A4 baseline were kept identical to those used in the Act-QAT stage of our GPLQ method.

Table 10 presents the ImageNet top-1 accuracy for various Swin Transformer models.

Table 10: Comparison of ImageNet top-1 accuracy (%) for Swin Transformers using direct W4A4 QAT versus GPLQ's sequential (W32A4 → W4A4) approach. Optimal results are in **bold**.

| Training Method | Swin-T | Swin-S | Swin-B | Swin-L |
|---|---|---|---|---|
| Direct W4A4 (1 epoch QAT) | 78.9 | 81.4 | 83.9 | 85.2 |
| GPLQ (W32A4 → W4A4) | **79.8** | **81.9** | **84.2** | **85.5** |

The results clearly demonstrate that our proposed sequential quantization strategy (GPLQ) consistently outperforms the direct W4A4 QAT approach across all Swin Transformer variants. As discussed in the main paper, quantizing activations first (while keeping weights in FP32) and then applying PTQ to

the weights offers several advantages. Beyond the slight improvements in training speed and reduced memory footprint during the Act-QAT stage (as pseudo-quantization of weights is not performed), this sequential approach helps to avoid the weight oscillations often encountered in traditional QAT. This leads to a smoother optimization process for the model weights, ultimately resulting in improved final model performance, as evidenced by the higher accuracies in Table 10.

# References

[1] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

[2] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.

[3] Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang, Peng Gao, Kaipeng Zhang, and Ping Luo. Efficientqat: Efficient quantization-aware training for large language models. *arXiv preprint arXiv:2407.11062*, 2024.

[4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[5] 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. *Advances in Neural Information Processing Systems*, 36:9015–9028, 2023.

[6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[7] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. Learned step size quantization. *arXiv preprint arXiv:1902.08153*, 2019.

[8] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.

[9] Minghao Fu, Hao Yu, Jie Shao, Junjie Zhou, Ke Zhu, and Jianxin Wu. Quantization without tears. *arXiv preprint arXiv:2411.13918*, 2024.

[10] Minghao Fu, Ke Zhu, and Jianxin Wu. Dtl: Disentangled transfer learning for visual recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 12082–12090, 2024.

[11] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.

[12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[14] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[15] Xijie Huang, Zhiqiang Shen, Pingcheng Dong, and Kwang-Ting Cheng. Quantization variation: A new perspective on training transformers with low-bit precision. *arXiv preprint arXiv:2307.00331*, 2023.

[16] Runqing Jiang, Ye Zhang, Longguang Wang, Pengpeng Yu, and Yulan Guo. Aiqvit: Architecture-informed post-training quantization for vision transformers. *arXiv preprint arXiv:2502.04628*, 2025.

[17] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013.

[18] Jiedong Lang, Zhehao Guo, and Shuyu Huang. A comprehensive study on quantization techniques for large language models. In *2024 4th International Conference on Artificial Intelligence, Robotics, and Communication (ICAIRC)*, pages 224–231. IEEE, 2024.

[19] Yanjing Li, Sheng Xu, Baochang Zhang, Xianbin Cao, Peng Gao, and Guodong Guo. Q-vit: Accurate and fully quantized low-bit vision transformer. *Advances in neural information processing systems*, 35:34451–34463, 2022.

[20] Zhikai Li, Junrui Xiao, Lianwei Yang, and Qingyi Gu. Repq-vit: Scale reparameterization for post-training quantization of vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17227–17236, 2023.

[21] 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 on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024.

[22] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014.

[23] Shih-Yang Liu, Zechun Liu, and Kwang-Ting Cheng. Oscillation-free quantization for low-bit vision transformers. In *International conference on machine learning*, pages 21813–21824. PMLR, 2023.

[24] 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.

[25] Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, and Wen Gao. Post-training quantization for vision transformer. *Advances in Neural Information Processing Systems*, 34:28092–28103, 2021.

[26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[27] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.

[28] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pages 722–729. IEEE, 2008.

[29] Lorenzo Papa, Paolo Russo, Irene Amerini, and Luping Zhou. A survey on efficient vision transformers: algorithms, techniques, and performance benchmarking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[30] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012.

[31] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021.

[32] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[34] Guo-Hua Wang, Yifan Ge, and Jianxin Wu. Distilling knowledge by mimicking features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):8183–8195, 2021.

[35] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023.

[36] Zhihang Yuan, Chenhao Xue, Yiqi Chen, Qiang Wu, and Guangyu Sun. Ptq4vit: Post-training quantization for vision transformers with twin uniform quantization. In *European conference on computer vision*, pages 191–207. Springer, 2022.

[37] Junjie Zhou, Ke Zhu, and Jianxin Wu. All you need in knowledge distillation is a tailored coordinate system. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 22946–22954, 2025.

[38] Ke Zhu, Yin-Yin He, and Jianxin Wu. Quantized feature distillation for network quantization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 11452–11460, 2023.