

---

# Memory-Efficient Vision Transformers: An Activation-Aware Mixed-Rank Compression Strategy

---

Seyedarmin Azizi<sup>1</sup> Mahdi Nazemi<sup>1</sup> Massoud Pedram<sup>1</sup>

## Abstract

As Vision Transformers (ViTs) increasingly set new benchmarks in computer vision, their practical deployment on inference engines is often hindered by their significant memory bandwidth and (on-chip) memory footprint requirements. This paper addresses this memory limitation by introducing an activation-aware model compression methodology that uses selective low-rank weight tensor approximations of different layers to reduce the parameter count of ViTs. The key idea is to decompose the weight tensors into a sum of two parameter-efficient tensors while minimizing the error between the product of the input activations with the original weight tensor and the product of the input activations with the approximate tensor sum. This approximation is further refined by adopting an efficient layer-wise error compensation technique that uses the gradient of the layer’s output loss. The combination of these techniques achieves excellent results while it avoids being trapped in a shallow local minimum early in the optimization process and strikes a good balance between the model compression and output accuracy. Notably, the presented method significantly reduces the parameter count of DeiT-B by 60% with less than 1% accuracy drop on the ImageNet dataset, overcoming the usual accuracy degradation seen in low-rank approximations. In addition to this, the presented compression technique can compress large DeiT/ViT models to have about the same model size as smaller DeiT/ViT variants while yielding up to 1.8% accuracy gain. These results highlight the efficacy of our approach, presenting a viable solution for embedding ViTs in memory-constrained environments without compromising their performance.

---

<sup>1</sup>Department of Electrical and Computer Engineering, University of Southern California, Los Angeles, United States. Correspondence to: Seyedarmin Azizi <seyedarm@usc.edu>.

## 1. Introduction

Vision Transformers (ViTs), highlighted in key studies like (Dosovitskiy et al., 2021), are recognized for their strong performance in various computer vision tasks. These tasks include image classification (Dosovitskiy et al., 2021; Touvron et al., 2021; Liu et al., 2021a), object detection (Carion et al., 2020; Fang et al., 2021), semantic segmentation (Chen et al., 2021a; Strudel et al., 2021), and multi-modal virtual assistance (Anonymous, 2024; Liu et al., 2023; Shayegani et al., 2023). However, the broader implementation of these transformers is notably hampered by their extensive parameter requirements, resulting in significant memory footprints.

In response to this challenge, model compression has emerged as the quintessential strategy for facilitating the deployment of models characterized by a high parameter count. Predominant techniques in this domain encompass model pruning (Zhu & Gupta, 2018; Zhu et al., 2021; Liu et al., 2019; Yu et al., 2022; Chen et al., 2021b), token pruning (Kong et al., 2022), quantization (Liu et al., 2021b; Yuan et al., 2022), and knowledge distillation (Touvron et al., 2021; Wu et al., 2022). These methodologies collectively aim to reduce the computational and storage burden, thereby enabling the efficient deployment of these advanced neural network architectures in resource-constrained environments.

Within the array of model compression techniques, low-rank approximation stands out as a particularly effective strategy for model compression, due to two reasons: (1) It has a solid theoretical foundation with proven optimality, as shown in (Eckart & Young, 1936), and (2) Its structured application directly translates to hardware efficiency and implementation ease. The effectiveness of this method is demonstrated in significant studies like (Jaderberg et al., 2014; Noach & Goldberg, 2020; Hsu et al., 2022).

Despite its potential, the naïve application of low-rank decomposition to the weights of ViTs often leads to a significant decline in performance, specifically when targeting higher compression rates. This issue arises primarily because the parameters of transformer-based models are typically not inherently suited to low-rank structures, as detailed in (Hsu et al., 2022). This underscores the need for a more nuanced and tailored approach in the application of

low-rank approximation techniques to ViTs.

To rectify the challenges associated with the application of low-rank decomposition to ViT weights, this paper introduces an innovative method for the decomposition of a pre-trained weight matrix. Our approach involves the decomposition of the weight matrix into a sum of two low-rank matrices, each contributing distinctively to the accurate reconstruction of the original matrix. Given a pre-trained weight matrix  $\mathbf{W}$ , we aim to approximate it as:

$$\mathbf{W} \approx \mathbf{UV}^T + \mathbf{Z} \quad (1)$$

where  $\mathbf{U}$  and  $\mathbf{V}$  are low-rank matrices, and  $\mathbf{Z}$  is also a matrix with a constrained number of parameters. Importantly,  $\mathbf{Z}$  must be designed for efficient hardware implementation and effective in approximating  $\mathbf{W} - \mathbf{UV}^T$ . This configuration is strategically designed to maintain the aggregate parameter count of approximation matrices significantly lower than that of the original, pre-trained weight matrix.

The proposed method, which will be detailed later, ensures that the product  $\mathbf{UV}^T$  and  $\mathbf{Z}$  provide distinct yet complementary contributions to the reconstruction of  $\mathbf{W}$ . Specifically,  $\mathbf{UV}^T$  captures the **principal energy** of the matrix through singular value decomposition (SVD), while  $\mathbf{Z}$  aims to offset the **residual error** from SVD using a layer-wise gradient-based optimization process.

Our approach for low-rank decomposition of weight matrices of ViTs is supported by two important observations: (1) statistics of the input feature map for each layer play a key role in influencing the approximation error associated with the parameters of the layer (as also discussed in (Yu & Wu, 2023)), and (2) layers of ViTs exhibit different susceptibilities to low-rank approximation, that is, aggressive rank reduction in some layers results in notable performance degradation at the model output. These insights form the foundation of our strategy, described as **activation-aware mixed-rank compression**, allowing for a smoother and tailored reduction in the number of parameters, thereby preserving the principal energy of the original ViT model’s weight matrices. The contributions of this paper may be summarized as follows:

- We formulate model compression as a general optimization problem that aims to find a low-rank approximation for each weight matrix in the ViT while minimizing the aggregate energy loss across all weight matrices.
- Our investigation delves into the impact of activation awareness in the application of singular value decomposition (SVD). We present a practical and highly effective methodology that incorporates input activations for the approximation of weight matrices, enhancing the approximation quality and capturing the principal energy contents of each layer.

- In terms of SVD implementation, we employ a strategic, gradual rank reduction approach, which judiciously assigns varying ranks to different layers within the model. The method is based on the fitness of the layers to low-rank approximation, and picks a layer which if its approximate tensor undergoes a parameter count reduction, the amount of information loss is minimum.
- To address the approximation error that is inherent in activation-aware SVD, we formulate a layer-specific gradient-based optimization problem. This approach aims to minimize the reconstruction error at each layer by decomposing the original weight matrix into a combination of the SVD result and a low parameter-count matrix, denoted as  $\mathbf{Z}$ . This crucial step serves to recuperate the energy loss encountered as a result of compression.
- We extend our methodology to various ViTs, conducting comprehensive experiments. These experiments yield compelling results in both accuracy and compression, demonstrating significant parameter count reduction without compromising model performance. Although our primary goal is reducing the memory footprint, we show in appendix A that adopting low-rank approximation not only does not hurt the computational efficiency but also improves it. Thus, our method introduces no memory-computation trade-off.

## 2. Background

A ViT architecture comprises a collection of identical blocks, each block comprising four layers, including the attention Query-Key-Value (QKV) layer, Attention Projection (AttnProj) layer, and two feed-forward Up Projection and Down Projection layers realized as Multi-Layer Perceptrons (MLP1 and MLP2).

Given a layer’s pre-trained weight matrix  $\mathbf{W} \in \mathbb{R}^{n \times d}$ , singular value decomposition (SVD) may be used to factorize it into  $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Here,  $\mathbf{U} \in \mathbb{R}^{n \times n}$  and  $\mathbf{V} \in \mathbb{R}^{d \times d}$  represent unitary matrices comprising the left and right singular vectors, respectively, while  $\mathbf{\Sigma} \in \mathbb{R}^{n \times d}$  is a diagonal matrix containing the singular values. To approximate  $\mathbf{W}$  with a rank  $r$  using SVD, the process involves retaining only the top  $r$  largest singular values and their corresponding singular vectors, resulting in an approximation  $\hat{\mathbf{W}} \approx \mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}_r^T$ . As established in (Eckart & Young, 1936), this specific rank- $r$  approximation is optimal, yielding minimal error amongst all potential rank  $r$  matrices.

Aligned with previous studies on implementing low-rank structures in neural networks (Jaderberg et al., 2014; Tai et al., 2016; Yu et al., 2017), the linear transformation of a layer may be approximated as follows:

$$\mathbf{O} = \mathbf{XW} + \mathbf{b} \approx \mathbf{X}(\mathbf{U}_r \mathbf{\Sigma}_r \mathbf{V}_r^T) + \mathbf{b} \quad (2)$$

where  $\mathbf{X}$  is the layer’s input activation,  $\mathbf{b}$  is the bias, and  $\mathbf{O}$  is the output.

Following the methodology in (Hsu et al., 2022), the singular values can be integrated into the left and right singular vectors. This integration results in a total post-approximation parameter count of  $n \times r + d \times r$ . **Throughout the paper we use the  $n_l$  and  $d_l$  to refer to the dimensions of the original weight matrix of layer  $l$ .**

The exploration of low-rank approximation, particularly within the realms of transformers and ViTs, has been an active area of research. For instance, (Hsu et al., 2022) effectively utilized Fisher Information as a means to evaluate the significance of the weight matrices in transformers, subsequently refining the objective of SVD to incorporate gradient awareness, thereby boosting its efficacy in model compression. Further advancements in this field are evident in the work of (Yu & Wu, 2023), where eigendecomposition is applied to the covariance of the layer’s output, revealing its enhanced suitability for low-rank approximation. Lastly, (Kumar, 2022) introduced a hybrid methodology that combines pruning techniques for feed-forward layers with low-rank decomposition for attention blocks, showcasing an innovative approach to model compression.

Despite these advancements, a common challenge persists across these methods: a significant drop in accuracy when high levels of compression are applied through low-rank structures, such as a 50% reduction in parameter count. This accuracy decline is largely attributed to the inability of these approaches to adequately compensate for the perturbations and energy loss induced by low-rank approximation. This highlights the necessity for more refined techniques that can effectively balance the trade-off between model compression and performance retention.

Another related work (Li et al., 2023) approximates the pre-trained weight matrix as a summation of a low-rank and sparse matrix, where the low-rank matrix captures the coherent part of the matrix, and a sparse matrix approximates the remaining incoherent residual; they demonstrated that this decoupling makes the matrix easy to prune. Despite this, since the SVD decomposition is directly applied to the weight matrix, they suffer from considerable energy losses (removals of large singular values), which results in noticeable performance degradation in high compression rates, highlighting the need for more sophisticated methods that can preserve the essential characteristics of the weight matrix.

### 3. Problem Formulation and Solution Methodology

In our setting, given a pre-trained neural network model  $\mathcal{M}$ , which has  $L$  layers, each represented by a weight

matrix  $\mathbf{W}_l$ , we are looking for a configuration of ranks  $\mathbf{r} = [r_1, r_2, \dots, r_L]$  such that when each layer  $l$  is approximated by rank  $r_l$ , the summation of the layers’ normalized errors across model  $\mathcal{M}$  (denoted by  $\text{Err}(\mathcal{M})$ ) is minimized:

$$\begin{aligned} \text{Err}(\mathcal{M}) = \min_{\mathbf{r}} \sum_l \min_{\hat{\mathbf{W}}_l} E_{\mathbf{X}_l} \left[ \frac{\|\mathbf{X}_l \hat{\mathbf{W}}_l - \mathbf{X}_l \mathbf{W}_l\|_F^2}{\|\mathbf{X}_l \mathbf{W}_l\|_F^2} \right] \\ \text{s.t. } \rho(\hat{\mathbf{W}}_l) \leq r_l \quad \forall l, \\ \psi(\mathcal{M}, \mathbf{r}) \geq \alpha \end{aligned} \quad (3)$$

Here,  $\mathbf{X}_l$  denotes the input activation for a layer  $l$ ,  $\mathbf{W}_l$  denotes the layer’s pre-trained weight matrix,  $\hat{\mathbf{W}}_l$  is the low-rank approximation of the weight matrix, constrained from above by rank  $r_l$ ,  $\rho(\cdot)$  denotes a function that returns the rank of the input matrix,  $E_{\mathbf{X}_l}(\cdot)$  denotes the expectation value over the set of all  $\mathbf{X}_l$  activations, and the function  $\psi(\mathcal{M}, \mathbf{r})$  computes the ratio of parameter count reduction of the model  $\mathcal{M}$  after approximating each original weight matrix  $\mathbf{W}_l$  by a low-rank matrix  $\hat{\mathbf{W}}_l$  of rank  $r_l$ . The main differentiating point between this formulation and prior work that deal with layer-wise output reconstruction (e.g., (Hubara et al., 2021; Frantar & Alistarh, 2022; Nagel et al., 2020)) is the incorporation of the sum of total normalized reconstruction errors across all layers, which effectively accounts for the interaction of layers in terms of compressibility.

If we define the quantity inside inner minimization as  $\varepsilon_l = \min_{\hat{\mathbf{W}}_l} E_{\mathbf{X}_l} \left[ \frac{(\|\mathbf{X}_l \hat{\mathbf{W}}_l - \mathbf{X}_l \mathbf{W}_l\|_F^2)}{(\|\mathbf{X}_l \mathbf{W}_l\|_F^2)} \right]$ , then  $\varepsilon_l$  captures at layer  $l$  the Frobenius norm of the output difference between the compressed representation and its uncompressed counterpart **in a normalized sense**. The Frobenius norm  $\|\mathbf{A}\|_F^2$  is defined as  $\sum_{i,j} a_{ij}^2$  and is equivalent to the sum of the squares of the singular values  $\sigma_i$  of  $\mathbf{A}$  (Horn & Johnson, 2012), representing the **energy of the matrix**.

This optimization seeks a rank configuration vector  $\mathbf{r}$  assigning ranks to layers, which minimizes the summation of normalized energy losses across all layers. The Min-Sum-Min nature of this optimization problem aims to achieve a balanced, minimal energy loss landscape across all layers, thus preserving the model’s inherent characteristics. This is a challenging optimization problem because the objective function is non-convex. Moreover, since the rank of each layer can be any integer number, the total number of possible configurations  $\mathbf{r}$  that achieve the desired compression ratio  $\alpha$  is exponentially large. Essentially, any configuration meeting the following condition is a feasible candidate:

$$\psi(\mathcal{M}, \mathbf{r}) = \frac{\sum_l n_l \times d_l}{\sum_l (n_l \times r_l + d_l \times r_l)} \geq \alpha. \quad (4)$$

This problem, in its general form, is NP-hard, indicating a high level of complexity and computational challenge.

To efficiently solve this problem, we have developed a multi-step heuristic flow, as explained below. Section 3.1 presents an **Activation-Aware Low-Rank Approximation** technique where, for a given rank configuration vector  $\mathbf{r}$  and layer  $l$ , the inner minimization of (3) is solved. That is, we compute the minimum energy loss  $\varepsilon_l$  and the corresponding optimized low-rank weight matrices  $\mathbf{U}_l$  and  $\mathbf{V}_l$ . This aspect of the framework leverages the characteristics of each layer’s input activations to determine the most effective low-rank approximation.

Section 3.2 introduces a **Mixed-Rank Compression** technique. To solve the outer minimization in (3), this strategy defines the rank configuration vector  $\mathbf{r}$  through a methodical, greedy local neighborhood search. It continuously integrates feedback from the activation-aware SVD process, ensuring that the rank distribution across layers is optimally aligned with their remaining compression potential and performance constraints.

Finally, section 3.3 proposes a **Layer-wise Error Compensation** technique, whose goal is to balance the residual errors introduced by the SVD. By employing a layer-wise gradient-based optimization technique, this technique finds a new matrix  $\mathbf{Z}_l$  to be added to the product of  $\mathbf{U}_l$  and  $\mathbf{V}_l$  for each layer  $l$ , further refining the low-rank approximation and enhancing the overall quality of the approximation.

### 3.1. Activation-Aware Low-Rank Approximation

First and foremost, we reformulate the reconstruction error based on our proposed approach:

$$\varepsilon_l = \min_{\mathbf{U}_l, \mathbf{V}_l, \mathbf{Z}_l} E_{\mathbf{X}_l} \left[ \frac{(\|\mathbf{X}_l(\mathbf{U}_l \mathbf{V}_l^T + \mathbf{Z}_l) - \mathbf{X}_l \mathbf{W}_l\|_F^2)}{(\|\mathbf{X}_l \mathbf{W}_l\|_F^2)} \right] \quad (5)$$

Given the challenges in directly computing the expectation, we follow a methodology similar to (Frantar & Alistarh, 2022) by constructing a proxy dataset. This dataset comprises samples of input activations for each layer, and we approximate the expectation with the average over these samples in the proxy dataset. A critical insight we’ve observed is the importance of considering the input activation during the application of singular value decomposition (SVD) to maximally preserve a layer’s output, as targeted in (5).

Re-examining Problem 5, if we initially set  $\mathbf{Z}_l$  to zero and the value of  $r_l$  is predetermined, then for a single input image  $\mathbf{X}_l^i$ , the problem can be reformulated as  $\min_{\mathbf{U}_l, \mathbf{V}_l} \|\mathbf{X}_l^i(\mathbf{U}_l \mathbf{V}_l^T) - \mathbf{X}_l^i \mathbf{W}_l\|_F^2$  (notice that the denominator is not a function of  $\mathbf{U}_l$  or  $\mathbf{V}_l$ ). For this special case, we can find the optimal  $\mathbf{U}_l$  and  $\mathbf{V}_l$  that minimize the reconstruction error as follows:

$$\mathbf{X}_l^i \mathbf{U}_l \mathbf{V}_l^T = \text{SVD}_{r_l}(\mathbf{X}_l^i \mathbf{W}_l) = \mathbf{U}_l^* \boldsymbol{\Sigma}_l^* \mathbf{V}_l^{*T} \quad (6)$$

In this context,  $\text{SVD}_k$  denotes the application of SVD while retaining the largest  $k$  singular values  $\boldsymbol{\Sigma}_l^*$  and their associated left and right singular vectors  $(\mathbf{U}_l^*, \mathbf{V}_l^*)$ , respectively. A feasible solution for  $\mathbf{U}$  and  $\mathbf{V}$  could be as follows:

$$\mathbf{U}_l = \mathbf{X}_l^{i\dagger} \mathbf{U}_l^* \sqrt{\boldsymbol{\Sigma}_l}, \quad \mathbf{V}_l = \sqrt{\boldsymbol{\Sigma}_l} \mathbf{V}_l^*, \quad (7)$$

where  $\mathbf{X}_l^{i\dagger}$  denotes the pseudo-inverse of the input activation.

This formulation considers the influence of single input activation. When utilizing a proxy dataset comprising  $N$  samples extracted from the original dataset rather than a single image, it becomes necessary to generate a representative input  $\mathbf{X}_{\text{rep}}$  for the application of (7). To simplify this process, one practical approach is to use the average of the samples in the proxy dataset. This representative sample  $\mathbf{X}_{\text{rep}}$  can be calculated as  $\mathbf{X}_{\text{rep}} = \frac{1}{N} \sum_{i=1}^N \mathbf{X}^i$ . This method ensures that  $\mathbf{X}_{\text{rep}}$  effectively captures the characteristics of the proxy dataset, providing a more holistic representation for the application of SVD. Then, we can use  $\mathbf{X}_{\text{rep}}$  in (7) to obtain  $\mathbf{U}$  and  $\mathbf{V}$ .

The approach we have taken in applying SVD aligns with the findings of (Yu & Wu, 2023), particularly regarding the layer’s output being more amenable to low-rank approximation compared to its weights. The matrix energy, as defined above, serves as a critical metric in this context. In Fig. 1, we analyze the energy loss in various layers of the DeiT-B model, comparing our activation-aware SVD approach to the direct application of SVD on weight matrices. Results indicate that activation-aware SVD is more effective in preserving the matrix energy, especially for higher rank reduction factors.

### 3.2. Mixed-Rank Model Compression

To achieve a specific model compression factor  $\alpha$ , regardless of the type of low-rank approximation employed (i.e., whether it is activation-aware or conventional SVD), it is possible to assign the ranks non-uniformly across the layers. This approach allows for a more flexible and potentially more effective compression strategy.

In Fig. 2, we present the relative decrease in energy of the weight matrices as the number of retained ranks in their spectrum is reduced. This figure points to the necessity of assigning different ranks to various transformer blocks and even to individual layers within each block to preserve the model’s performance effectively. This differentiation in rank allocation is crucial for maintaining a balance between model compression and performance retention. This phenomenon also aligns with the findings of the other domains of compression, including mixed-precision quantization, where layers of the model are differently sensitive to low bit-width quantization (Wang et al., 2019; Azizi et al., 2023; Dong et al., 2019).

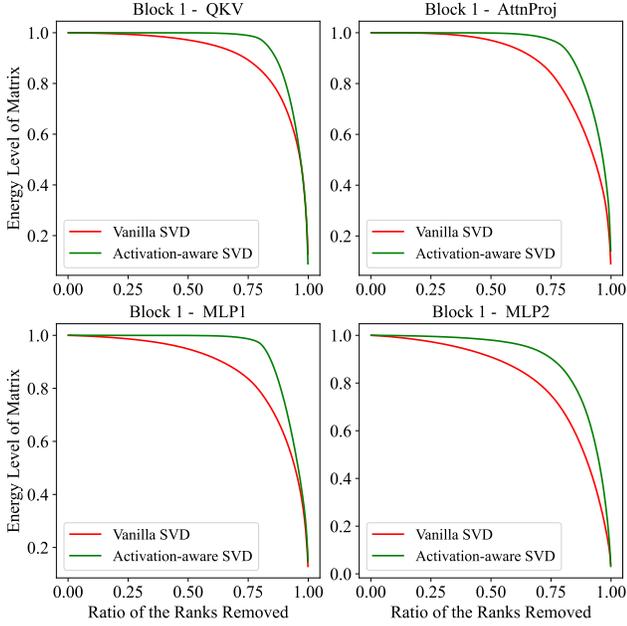


Figure 1. Impact of SVD-based rank reduction on energy level of different matrices in the first block of DeiT-B.

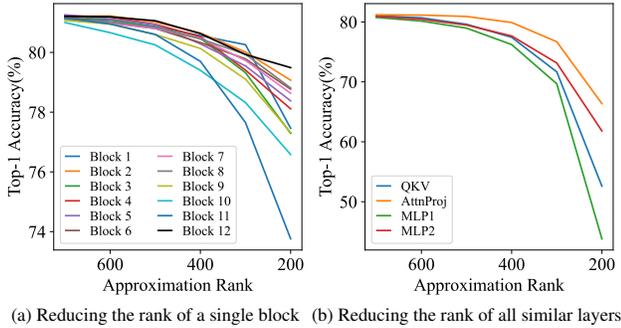


Figure 2. Impact of rank reduction on top-1 ImageNet accuracy

Considering Problem (3), our goal is to identify an optimal rank configuration  $\mathbf{r}$ , which not only reduces the model’s parameter count by a factor of  $\alpha$  but also minimizes the sum of energy losses observed across all layers of the model. As discussed in section 3.1, we have already presented a technique that, for a given rank  $r_l$ , calculates the optimal energy loss  $\varepsilon_l$  at the output of each layer  $l$ . So, we can simplify the formulation of (3) as follows:

$$\begin{aligned} \text{Err}(\mathcal{M}) &= \min_{\mathbf{r}} \sum_l \varepsilon_l \\ \text{s.t. } & \psi(\mathcal{M}, \mathbf{r}) \geq \alpha \end{aligned} \quad (8)$$

As we discussed earlier in section 3, this is still an NP-hard problem. To address this, we propose an efficient greedy solution based on local neighborhood search. Instead of a drastic, single-step reduction in ranks and parameters, we

adopt a gradual optimization strategy. This method involves incremental steps, each of which slightly reduces the number of parameters in line with a pre-defined **compression rate scheduling policy**. For each step of parameter reduction, we select the layer whose (further) rank reduction yields the **minimum normalized energy loss** compared to its uncompressed counterpart, and reduce its rank to meet the target parameter count reduction. This iterative process avoids disproportionate normalized energy loss at any individual layer, thus greedily minimizing the sum of normalized energy loss across all layers (see (3)).

Algorithm 1 details this process where function **ActivationAwareSVD()** is what we developed in section 3.1 for approximating  $\varepsilon_l$ , and  $\tau$  is the rate scheduling function that sets the rate of rank reduction. We adopt an exponential of the form below:

$$\tau(\mathcal{M}, \text{iter}, \alpha) = N_{\text{target}} + (N_0 - N_{\text{target}}) \exp\left(-\frac{\text{iter}}{\gamma}\right) \quad (9)$$

in which  $N_0$  is the initial number of parameters,  $N_{\text{target}}$  is our desired number of parameters, and  $\gamma$  is the decaying rate of the schedule. The exponential nature of the schedule encourages significant parameter reduction in the initial stages, primarily targeting layers that are less sensitive to rank reduction. As the process progresses, the required rate of parameter reduction is decreased like an annealing schedule. Consequently, this approach minimizes the impact on layers that are more sensitive to rank reduction. The function **ComputeEnergyLoss()**, which is called for each layer  $l$  and in each iteration  $t$  of algorithm, computes the ratio of the lost energy for layer  $l$  and its original total energy if it only retains  $m^t[l]$  of its singular values.

A key aspect of Algorithm 1 is its efficiency in handling the Singular Value Decomposition (SVD) of each layer. Notably, the SVD for each layer needs to be computed only once. After this initial computation, we store the singular values obtained from the decomposition. Subsequently, all optimization processes leverage these stored singular values. This approach significantly reduces computational overhead, as it eliminates the need for repeated SVD computations for each layer during the optimization steps.

### 3.3. Layer-Wise Error Compensation

Using the combination of the techniques developed in sections 3.1 and 3.2 we are able to come up with the low-rank approximation matrices  $\mathbf{U}$  and  $\mathbf{V}$  for each layer, with the ranks specific to each layer. However, applying SVD and removing some of the singular vectors would unavoidably introduce some energy loss. In mixed-rank compression we managed to create a balanced distribution of energy losses across all layers. In this section, we present a novel and efficient method to compensate for the energy loss experienced

**Algorithm 1** Mixed-Rank Compression

---

**Input:** Model  $\mathcal{M}$ , proxy dataset  $\mathcal{D}$ , Compression ratio  $\alpha$   
**Output:** List of ranks  $\mathbf{r}$

- 1: Singular Values  $\mathbf{S} = \text{ActivationAwareSVD}(\mathcal{M}, \mathcal{D})$
- 2: Initialize  $p_{\text{tot}} = p_{\text{cur}} = \sum_l n_l \times d_l$  // # of param.
- 3: Initialize  $\mathbf{r}^0 = \min(n_l, d_l) \forall l$  // hidden dim. in ViT
- 4:  $t = 1$
- 5: **repeat**
- 6:    $p^t = \tau(\mathcal{M}, t, \alpha)$  // # of param. to remove at step  $t$
- 7:    $\ell = []$  // vector for keeping energy losses
- 8:    $\mathbf{m}^t = []$  // vector for keeping temporary ranks
- 9:    $\mathbf{r}^t = \mathbf{r}^{t-1}$
- 10:   **for**  $l = 1$  to  $L$  **do**
- 11:      $\mathbf{m}^t[l] = \mathbf{r}^{t-1}[l] - \frac{p^t}{n_l + d_l}$  // # of ranks to keep
- 12:      $\ell[l] = \text{ComputeEnergyLoss}(\mathbf{S}[l], \mathbf{m}^t[l])$
- 13:   **end for**
- 14:    $l^* = \arg \min_l \ell$  // layer with minimum energy loss
- 15:    $\mathbf{r}^t[l^*] = \mathbf{m}^t[l^*]$
- 16:    $p_{\text{cur}} = p_{\text{cur}} - p^t$
- 17:    $t = t + 1$
- 18: **until**  $p_{\text{cur}} \leq \frac{p_{\text{tot}}}{\alpha}$
- 19: **return**  $\mathbf{r}^t$

---

by all layers effectively. Revisiting (5), we are looking for an auxiliary matrix  $\mathbf{Z}$  that is efficient in terms of the number of parameters and can compensate for the residual error introduced by activation-aware SVD.

We investigate four distinct configurations for the  $\mathbf{Z}$  matrix: (1, 2) Sparse matrix (both structured and unstructured) in a manner similar to (Li et al., 2023), where sparse matrices were utilized for capturing the uncorrelated components of the approximation’s residual; (3) Diagonal matrix; (4) Low-rank matrix, with a rank substantially smaller than that of  $\mathbf{U}$  and  $\mathbf{V}$ . Table 1 compares these strategies in terms of the parameter count overhead and normalized remaining residual error after the layer-wise error compensation. Low-rank matrix is a compelling option regarding both metrics, thus we introduce  $\mathbf{G}\mathbf{Y}^T$  as the matrix for capturing the residual error, where  $\mathbf{G} \in \mathbb{R}^{n \times q}$  and  $\mathbf{Y} \in \mathbb{R}^{d \times q}$  have a rank much smaller than that of  $\mathbf{U}$  and  $\mathbf{V}$  ( $q \ll r_l \forall l$ ).

Table 1. Comparison of different approaches for residual error compensation via  $\mathbf{Z}$ .

↑ and ↓ indicates relatively high and low values, respectively.

Approach	Parameter Count Overhead	Average Residual Error
Sparse (unstructured)	↑	0.11 ↓
Sparse (structured)	↓	0.23 ↑
Diagonal	↓	0.27 ↑
<b>Low-Rank</b>	↓	0.13 ↓

The introduction of  $\mathbf{G}$  and  $\mathbf{Y}$  serves as a key innovation for addressing the energy loss inherent in the SVD process. By starting with these matrices as zero and progressively updating them, we are able to fine-tune the approximation in a way that specifically targets the reconstruction errors and energy deficiencies resulting from the initial SVD. Firstly, let’s revisit (5) with a slight modification. In this iteration, we maintain  $\mathbf{U}$  and  $\mathbf{V}$  as fixed, based on the values obtained from the activation-aware low-rank approximation (section 3.1). Additionally, we substitute the expectation term in the equation with a summation over the samples in the proxy dataset. This adjustment aligns with our earlier discussion about the computational challenges of directly computing the expectation and the practical solution of using a proxy dataset  $\mathcal{D}$  for approximation:

$$\varepsilon_l = \min_{\mathbf{G}_l, \mathbf{Y}_l} \frac{1}{|\mathcal{D}|} \sum_i \left[ \frac{\| \mathbf{X}_l^i (\mathbf{G}_l \mathbf{Y}_l^T) - \mathbf{X}_l^i (\mathbf{W}_l - \mathbf{U}_l \mathbf{V}_l^T) \|_F^2}{(\| \mathbf{X}_l^i \mathbf{W}_l \|_F^2)} \right] \quad (10)$$

To simplify the implementation process, we choose to fix the rank  $q$  for each layer statically. This approach entails setting the total number of parameters in  $\mathbf{G}_l$  and  $\mathbf{Y}_l$  to a pre-defined percentage of the parameters in  $\mathbf{W}_l$ . We opt for this percentage to be 5%.

With the dimensions of these low-rank matrices  $\mathbf{G}_l$  and  $\mathbf{Y}_l$  fixed, the problem essentially transforms into a regression task. For each layer  $l$ , we define  $\mathbf{A}_l^i = \frac{\mathbf{X}_l^i}{(\| \mathbf{X}_l^i \mathbf{W}_l \|_F^2)}$ ,  $\mathbf{B}_l^i = \frac{\mathbf{X}_l^i (\mathbf{W}_l - \mathbf{U}_l \mathbf{V}_l^T)}{(\| \mathbf{X}_l^i \mathbf{W}_l \|_F^2)}$ , and the loss function  $\mathcal{L}_l = \frac{1}{|\mathcal{D}|} \sum_i \| \mathbf{A}_l^i (\mathbf{G}_l \mathbf{Y}_l^T) - \mathbf{B}_l^i \|_F^2$ . This setting forms the basis of our optimization problem. A crucial observation is that this optimization process can be conducted for each layer of the model independently. As a result, the optimization tasks for different layers can be executed **in parallel**. This parallelization greatly enhances the efficiency of the optimization process, allowing for simultaneous adjustments across multiple layers. Although there is no closed-form solution, this problem can be solved using gradient descent since the gradient of the loss function with respect to  $\mathbf{G}$  and  $\mathbf{Y}$  can be easily computed as below:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{G}} = \frac{1}{|\mathcal{D}|} \sum_i 2 \mathbf{A}^{iT} (\mathbf{A}^i \mathbf{G} \mathbf{Y}^T - \mathbf{B}^i) \mathbf{Y} \quad (11)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{Y}} = \frac{1}{|\mathcal{D}|} \sum_i 2 (\mathbf{A}^i \mathbf{G} \mathbf{Y}^T - \mathbf{B}^i)^T \mathbf{A} \mathbf{G} \quad (12)$$

Intuitively, the  $\mathbf{X}\mathbf{U}\mathbf{V}^T$  component is designed to capture the main components and a substantial portion of the energy of  $\mathbf{X}\mathbf{W}$  via activation-aware SVD. Consequently,  $\mathbf{X}\mathbf{G}\mathbf{Y}^T$  endeavors to minimize the residual error left by the SVD, utilizing the relatively small matrices  $\mathbf{G}$  and  $\mathbf{Y}$ . This dual

approach aims to ensure a thorough and efficient approximation of the original weight matrices.

To demonstrate the effectiveness of this addition, we analyze the Frobenius norm error of the attention projection module in DeiT-B model. This analysis is conducted post-activation-aware SVD, following the gradient-based optimization on an unseen (validation) proxy dataset. The results, illustrated in Fig. 3, highlight the impact of our approach on these layers, showcasing the reduction in reconstruction error achieved through our optimization process. The activation-aware SVD achieves an almost uniform distribution of error across the layers. This outcome aligns perfectly with the primary objective of our problem; on top of that, the gradient-based optimization reduces the layer-wise error as much as possible. If we omit the gradient-based error compensation and instead allocate the parameter budget of  $\mathbf{G}$  and  $\mathbf{Y}$  to  $\mathbf{U}$  and  $\mathbf{V}$ , the error reduction would be a lot less than the residual error reduction through  $\mathbf{G}$  and  $\mathbf{Y}$ . This scenario highlights the critical role of error compensation, as evidenced by the comparative results. An important aspect of this optimization is that  $\mathbf{U}_i \mathbf{V}_i^T$  remains constant while optimizing  $\mathbf{G}_i \mathbf{Y}_i^T$ . Given that the dimensions of  $\mathbf{G}$  and  $\mathbf{Y}$  are significantly smaller than those of  $\mathbf{U}$  and  $\mathbf{V}$ , there is a minimal risk of overfitting to the calibration dataset, thus providing more generalizability and robustness.

## 4. Results and Discussions

In this section, we thoroughly evaluate our compression framework.

### 4.1. Experimental Setup

In our experimental setup, we have carefully selected specific hyperparameters to optimize the performance of our model. Here is an overview of the key settings:

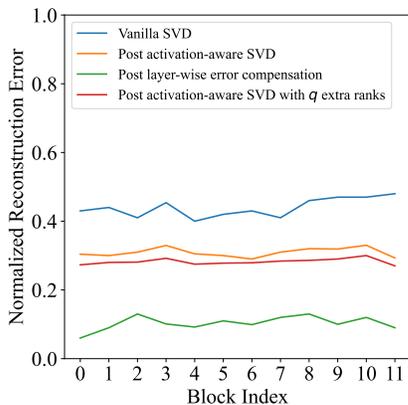


Figure 3. Normalized Frobenius norm of the error at the output of the AttnProj layer in DeiT-B

1. Proxy Dataset Creation: For both activation-aware SVD and gradient-based optimization, we construct the proxy dataset using 1024 samples randomly selected from the dataset. This sample size is chosen to provide a representative subset of the original data.

2. Rank Setting for  $\mathbf{G}$  and  $\mathbf{Y}$ : The rank  $q$  for the matrices  $\mathbf{G}$  and  $\mathbf{Y}$  is determined so that these matrices together account for 5% of the parameters in each layer.

3. Mixed-Rank Gradual Compression: In this approach, we set the decay rate  $\gamma$  at 80 and the iterative algorithm for this compression is run for 500 iterations.

4. Layer-wise Error Compensation Settings (Section 3.3): During the gradient-based optimization phase, we set the learning rate for updating  $\mathbf{G}$  and  $\mathbf{Y}$  to  $10^{-3}$ . This process involves running Mini-Batch gradient descent for 2000 iterations with a batch size of 64. Since this optimization is done in a layer-wise manner, no backpropagation is involved, and layers can be processed in parallel; thus, it is very fast.

5. As the final step, we fine-tune the uncompressed parameters of the model, including the **LayNorm** parameters, **head**, **biases**, and the **patch embedding** module on the standard ImageNet dataset. This **partial fine-tuning** is done for 20 epochs with a learning rate of  $10^{-4}$  and a cosine scheduling. Since almost 2% of the network parameters are being fine-tuned and the weight matrices are frozen (they are already optimized using our flow), this step is very fast and efficient.

6. Library and Hardware: We utilize pre-trained models from the timm library (Wightman, 2019) and implement our model optimization using PyTorch (Paszke et al., 2019). All experiments are conducted on NVIDIA A6000 GPUs.

### 4.2. ImageNet Classification

We assess our framework’s effectiveness in compressing various Vision Transformer architectures, including ViT (Dosovitskiy et al., 2021), DeiT (Touvron et al., 2021), and Swin Transformer (Liu et al., 2021a), on the ImageNet dataset. The detailed results are presented in table 2. The table includes Params, indicating both the absolute number of parameters and the percentage of reduction relative to the baseline model. We evaluate the models at three different levels of parameter count reduction: -40%, -50%, and -60%, and correspondingly report the Top-1 validation accuracy.

Our results demonstrate that we can achieve a 50% reduction in the parameter count of DeiT-B without any loss in accuracy. Furthermore, we accomplish a 60% reduction in parameters with less than a 1% drop in accuracy, showcasing the efficacy of our compression framework across different levels of parameter reduction. It can be seen that our model outperforms the previous compression approaches in accu-

Table 2. Comparison of the parameter count and top-1 ImageNet accuracy for different compression methods and ViT architectures.

FT stands for fine-tuning.

 The results of other methods are directly sourced from their respective papers (where available): AAFM (Yu & Wu, 2023), WDPPruning (Yu et al., 2022), SPViT (Kong et al., 2022), S<sup>2</sup>ViTE (Chen et al., 2021b), and VTP (Zhu et al., 2021);

Architecture	Method	# Params (Millions)	Top-1 Accuracy (%)
DeiT-B	Baseline	86.6	81.80
	AAFM	51.9(-40%)	81.28
	WDPPruning	60.6(-30%)	81.10
	S <sup>2</sup> ViTE	56.8(-35%)	82.20
	SPViT	62.3(-28%)	81.60
	Vanilla SVD + FT	52.0(-40%)	77.30
	<b>Ours</b>	<b>52.0(-40%)</b>	<b>81.80</b>
	<b>Ours</b>	<b>43.3(-50%)</b>	<b>81.35</b>
	<b>Ours</b>	<b>34.6(-60%)</b>	<b>81.10</b>
DeiT-S	Baseline	22.1	79.80
	WDPPruning	15.0(-32%)	78.60
	S <sup>2</sup> ViTE	14.6(-34%)	79.20
	SPViT	16.4(-26%)	78.30
	Vanilla SVD + FT	13.2(-40%)	75.20
	<b>Ours</b>	<b>13.2(-40%)</b>	<b>79.60</b>
	<b>Ours</b>	<b>11.1(-50%)</b>	<b>79.34</b>
	<b>Ours</b>	<b>8.9(-60%)</b>	<b>78.60</b>
	ViT-B	Baseline	86.5
VTP		48.0(-45%)	80.70
Vanilla SVD + FT		42.3(-50%)	78.50
<b>Ours</b>		<b>52.9(-40%)</b>	<b>84.20</b>
<b>Ours</b>		<b>42.3(-50%)</b>	<b>83.70</b>
Swin-B	Baseline	88.1	85.45
	AAFM	60.2(-33%)	82.68
	SPViT	68.0(-24%)	83.20
	Vanilla SVD + FT	52.9(-40%)	79.10
	<b>Ours</b>	<b>52.9(-40%)</b>	<b>83.90</b>
	<b>Ours</b>	<b>44.1(-50%)</b>	<b>83.65</b>
<b>Ours</b>	<b>35.3(-60%)</b>	<b>83.14</b>	

racy, compression, or both.

To demonstrate our method’s capabilities, we undertake extreme compression cases: compressing ViT-L to match DeiT-B’s size, DeiT-B to DeiT-S’s size, and DeiT-S to DeiT-T’s size. This tests our method’s efficacy in significantly reducing model sizes while maintaining performance. The outcomes of these compression experiments are detailed in table 3. The data shows that our compressed models, despite having a similar number of parameters, outperform the pre-trained models they are compared against, highlighting the effectiveness of our compression method.

### 4.3. Compatibility with Weight Quantization

To illustrate the compatibility of the presented approach with weight quantization, we apply post-training quantization (PTQ) to weight matrices  $U$ ,  $V$ ,  $G$ , and  $Y$ . For this

Table 3. Comparison of the top-1 ImageNet accuracy on models with approximately the same size.

Architecture	# Params (Millions)	Top-1 Accuracy (%)
ViT-B	86.6	81.8
<b>Compressed ViT-L</b>	<b>86.0 (-72%)</b>	<b>83.1</b>
DeiT-S	22.1	79.8
<b>Compressed DeiT-B</b>	<b>21.7 (-75%)</b>	<b>80.3</b>
DeiT-T	5.0	72.2
<b>Compressed DeiT-S</b>	<b>5.0 (-77%)</b>	<b>74.0</b>

Table 4. Compatibility of the presented compression approach with weight quantization.

Method	Model Size (MiB)	Top-1 Accuracy (%)
DeiT-B baseline (FP16)	165.2	81.80
DeiT-B 50% compressed (FP16)	82.6	81.35
<b>DeiT-B 50% compressed + 8-bit PTQ</b>	<b>41.7</b>	<b>81.15</b>
DeiT-B baseline + 4-bit PTQ	42.2	80.72

purpose, we employ 8-bit, channel-wise, round-to-nearest quantization, targeting only the weights (the details of the quantization function applied can be found in appendix). We then compare this mixed compression strategy with applying 4-bit PTQ to the baseline model, which yields about the same level of compression.

As summarized in table 4, the 8-bit version of our compressed model surpasses the 4-bit uncompressed model in accuracy by 0.43%. This result indicates the compatibility of quantization with our low-rank approximation and its superiority compared with quantization-only compression.

## 5. Conclusions and Future Work

In this work, we proposed a novel methodology for compressing ViTs. We adopted activation-aware SVD to approximate the outputs of the layers within the model while maintaining the principal energy components of the matrices. This approximation was refined by developing a greedy strategy for assigning various ranks to different layers. In the end, we also proposed layer-wise error compensation for reducing the error introduced by SVD as much as possible. Overall, our methodology significantly reduces the parameter count of ViTs, facilitating their efficient deployment in inference engines.

While our current work has been focused on the compression of ViTs, the presented approach exhibits promising characteristics that suggest its potential applicability to other transformer-based architectures like large language models. Investigating and adapting our methodology to other transformer variants is an exciting direction for future work.

## References

- Anonymous. Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=plmBsXHxgR>.
- Azizi, S., Nazemi, M., Fayyazi, A., and Pedram, M. Sensitivity-aware mixed-precision quantization and width optimization of deep neural networks through cluster-based tree-structured parzen estimation. *arXiv preprint arXiv:2308.06422*, 2023.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. In Vedaldi, A., Bischof, H., Brox, T., and Frahm, J. (eds.), *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part I*, volume 12346 of *Lecture Notes in Computer Science*, pp. 213–229. Springer, 2020. doi: 10.1007/978-3-030-58452-8\_13. URL [https://doi.org/10.1007/978-3-030-58452-8\\_13](https://doi.org/10.1007/978-3-030-58452-8_13).
- Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, A. L., and Zhou, Y. Transunet: Transformers make strong encoders for medical image segmentation. *CoRR*, abs/2102.04306, 2021a. URL <https://arxiv.org/abs/2102.04306>.
- Chen, T., Cheng, Y., Gan, Z., Yuan, L., Zhang, L., and Wang, Z. Chasing sparsity in vision transformers: An end-to-end exploration. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 19974–19988, 2021b. URL <https://proceedings.neurips.cc/paper/2021/hash/a61f27ab2165df0e18cc9433bd7f27c5-Abstract.html>.
- Dong, Z., Yao, Z., Gholami, A., Mahoney, M. W., and Keutzer, K. HAWQ: hessian aware quantization of neural networks with mixed-precision. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pp. 293–302. IEEE, 2019. doi: 10.1109/ICCV.2019.00038. URL <https://doi.org/10.1109/ICCV.2019.00038>.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Eckart, C. and Young, G. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- Fang, Y., Liao, B., Wang, X., Fang, J., Qi, J., Wu, R., Niu, J., and Liu, W. You only look at one sequence: Rethinking transformer in vision through object detection. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 26183–26197, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/dc912a253d1e9ba40e2c597ed2376640-Abstract.html>.
- Frantar, E. and Alistarh, D. Optimal brain compression: A framework for accurate post-training quantization and pruning. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/1caf09c9f4e6b0150b06a07e77f2710c-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/1caf09c9f4e6b0150b06a07e77f2710c-Abstract-Conference.html).
- Horn, R. A. and Johnson, C. R. *Matrix Analysis, 2nd Ed.* Cambridge University Press, 2012. ISBN 9780521548236. doi: 10.1017/CBO9781139020411. URL <https://doi.org/10.1017/CBO9781139020411>.
- Hsu, Y., Hua, T., Chang, S., Lou, Q., Shen, Y., and Jin, H. Language model compression with weighted low-rank factorization. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=uPv9Y3gmAI5>.
- Hubara, I., Nahshan, Y., Hanani, Y., Banner, R., and Soudry, D. Accurate post training quantization with small calibration sets. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 4466–4475. PMLR, 2021. URL <http://proceedings.mlr.press/v139/hubara21a.html>.
- Jaderberg, M., Vedaldi, A., and Zisserman, A. Speeding up convolutional neural networks with low rank expansions. In Valstar, M. F., French, A. P., and Pridmore, T. P. (eds.), *British Machine Vision Conference, BMVC 2014, Nottingham, UK, September 1-5, 2014*. BMVA Press, 2014. URL <http://www.bmva.org/bmvc/2014/papers/paper073/index.html>.

- Kong, Z., Dong, P., Ma, X., Meng, X., Niu, W., Sun, M., Shen, X., Yuan, G., Ren, B., Tang, H., Qin, M., and Wang, Y. Spvit: Enabling faster vision transformers via latency-aware soft token pruning. In Avidan, S., Brostow, G. J., Cissé, M., Farinella, G. M., and Hassner, T. (eds.), *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XI*, volume 13671 of *Lecture Notes in Computer Science*, pp. 620–640. Springer, 2022. doi: 10.1007/978-3-031-20083-0\\_37. URL [https://doi.org/10.1007/978-3-031-20083-0\\\_37](https://doi.org/10.1007/978-3-031-20083-0\_37).
- Kumar, A. Vision transformer compression with structured pruning and low rank approximation. *arXiv preprint arXiv:2203.13444*, 2022.
- Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.
- Li, Y., Yu, Y., Zhang, Q., Liang, C., He, P., Chen, W., and Zhao, T. Lospase: Structured compression of large language models based on low-rank and sparse approximation. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 20336–20350. PMLR, 2023. URL <https://proceedings.mlr.press/v202/li23ap.html>.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. In *NeurIPS*, 2023.
- Liu, Z., Sun, M., Zhou, T., Huang, G., and Darrell, T. Rethinking the value of network pruning. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=rJlnB3C5Ym>.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. Swin transformer: Hierarchical vision transformer using shifted windows. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 9992–10002. IEEE, 2021a. doi: 10.1109/ICCV48922.2021.00986. URL <https://doi.org/10.1109/ICCV48922.2021.00986>.
- Liu, Z., Wang, Y., Han, K., Zhang, W., Ma, S., and Gao, W. Post-training quantization for vision transformer. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 28092–28103, 2021b. URL <https://proceedings.neurips.cc/paper/2021/hash/ec8956637a99787bd197eacd77acce5e-Abstract.html>.
- Nagel, M., Amjad, R. A., van Baalen, M., Louizos, C., and Blankevoort, T. Up or down? adaptive rounding for post-training quantization. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 7197–7206. PMLR, 2020. URL <http://proceedings.mlr.press/v119/nagel20a.html>.
- Noach, M. B. and Goldberg, Y. Compressing pre-trained language models by matrix decomposition. In Wong, K., Knight, K., and Wu, H. (eds.), *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing, AACL/IJCNLP 2020, Suzhou, China, December 4-7, 2020*, pp. 884–889. Association for Computational Linguistics, 2020. URL <https://aclanthology.org/2020.acl-main.88/>.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E. Z., DeVito, Z., Raiison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 2019.
- Shayegani, E., Mamun, M. A. A., Fu, Y., Zaree, P., Dong, Y., and Abu-Ghazaleh, N. Survey of vulnerabilities in large language models revealed by adversarial attacks, 2023.
- Strudel, R., Pinel, R. G., Laptev, I., and Schmid, C. Seg-menter: Transformer for semantic segmentation. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 7242–7252. IEEE, 2021. doi: 10.1109/ICCV48922.2021.00717. URL <https://doi.org/10.1109/ICCV48922.2021.00717>.
- Tai, C., Xiao, T., Wang, X., and E, W. Convolutional neural networks with low-rank regularization. In Bengio, Y. and LeCun, Y. (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1511.06067>.
- Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., and Jégou, H. Training data-efficient image transformers & distillation through attention. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021,*

- 18-24 July 2021, Virtual Event, volume 139 of *Proceedings of Machine Learning Research*, pp. 10347–10357. PMLR, 2021. URL <http://proceedings.mlr.press/v139/touvron21a.html>.
- Wang, K., Liu, Z., Lin, Y., Lin, J., and Han, S. HAQ: hardware-aware automated quantization with mixed precision. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 8612–8620. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00881. URL [http://openaccess.thecvf.com/content\\_CVPR\\_2019/html/Wang\\_HAQ\\_Hardware-Aware\\_Automated\\_Quantization\\_With\\_Mixed\\_Precision\\_CVPR\\_2019\\_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Wang_HAQ_Hardware-Aware_Automated_Quantization_With_Mixed_Precision_CVPR_2019_paper.html).
- Wightman, R. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- Wu, K., Zhang, J., Peng, H., Liu, M., Xiao, B., Fu, J., and Yuan, L. Tinyvit: Fast pretraining distillation for small vision transformers. In Avidan, S., Brostow, G. J., Cissé, M., Farinella, G. M., and Hassner, T. (eds.), *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXI*, volume 13681 of *Lecture Notes in Computer Science*, pp. 68–85. Springer, 2022. doi: 10.1007/978-3-031-19803-8\5. URL <https://doi.org/10.1007/978-3-031-19803-8.5>.
- Yu, F., Huang, K., Wang, M., Cheng, Y., Chu, W., and Cui, L. Width & depth pruning for vision transformers. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pp. 3143–3151. AAAI Press, 2022. doi: 10.1609/AAAI.V36I3.20222. URL <https://doi.org/10.1609/aaai.v36i3.20222>.
- Yu, H. and Wu, J. Compressing transformers: Features are low-rank, but weights are not! In Williams, B., Chen, Y., and Neville, J. (eds.), *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pp. 11007–11015. AAAI Press, 2023. doi: 10.1609/AAAI.V37I9.26304. URL <https://doi.org/10.1609/aaai.v37i9.26304>.
- Yu, X., Liu, T., Wang, X., and Tao, D. On compressing deep models by low rank and sparse decomposition. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 67–76. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.15. URL <https://doi.org/10.1109/CVPR.2017.15>.
- Yuan, Z., Xue, C., Chen, Y., Wu, Q., and Sun, G. Ptq4vit: Post-training quantization for vision transformers with twin uniform quantization. In Avidan, S., Brostow, G. J., Cissé, M., Farinella, G. M., and Hassner, T. (eds.), *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XII*, volume 13672 of *Lecture Notes in Computer Science*, pp. 191–207. Springer, 2022. doi: 10.1007/978-3-031-19775-8\12. URL <https://doi.org/10.1007/978-3-031-19775-8.12>.
- Zhu, M. and Gupta, S. To prune, or not to prune: Exploring the efficacy of pruning for model compression. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Workshop Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=Syl1iIDkPM>.
- Zhu, M., Tang, Y., and Han, K. Vision transformer pruning. *arXiv preprint arXiv:2104.08500*, 2021.

## A. Impact of Low-Rank Approximation on Compute Efficiency

Given a matrix multiplication  $O = XW$ , where  $X$  and  $W$  are  $n \times k$  and  $k \times m$  matrices respectively, we approximate the output as  $O \approx X(UV^T + GY^T)$ . Here,  $U$ ,  $V$ ,  $G$ , and  $Y$  are  $k \times r$ ,  $m \times r$ ,  $k \times q$ , and  $m \times q$  matrices, respectively. The original matrix multiplication requires  $nk m$  individual multiplication operations.

If we compute  $UV^T$  and  $GY^T$  directly, add them, and then multiply by  $X$ , the total number of multiplications is:

$$krm \text{ (for } UV^T) + kqm \text{ (for } GY^T) + nkm \text{ (when multiplying by } X) \quad (13)$$

This increases the multiplication count. However, a more efficient approach is to first multiply  $X$  by  $U$  and then the result by  $V^T$ , and similarly for  $G$  and  $Y$ . Compute  $XG$  first, then multiply it by  $Y^T$ . The total multiplications become:

$$nkr \text{ (for } XU) + nrm \text{ (for } XUV^T) + nkq \text{ (for } XG) + nqm \text{ (for } XGY^T) \quad (14)$$

The ratio of these multiplications to the original is:

$$\frac{nkr + nrm + nkq + nqm}{nkm} = \frac{nk(r+q) + mn(r+q)}{nkm} \quad (15)$$

Since  $r+q$  is smaller than the original matrix dimensions  $k$  and  $m$ , this ratio is less than 1. Thus, our method, while not primarily aimed at computational efficiency, inadvertently achieves a reduction in multiplication count, leading to a speedup.

## B. Quantization

In our experiments, post-training quantization (PTQ) was implemented to demonstrate how our weight compression technique can be effectively combined with other methods. The specific quantization function used is the basic round-to-nearest linear quantization, defined as follows:

$$U_q = \text{clamp}(\lfloor \frac{U}{s} + z \rceil, 0, 2^N - 1), \quad \hat{U} = s \times (U_q - z) \quad (16)$$

Here,  $s = \frac{\max(U)}{2^N - 1}$  and  $z = -\frac{\min(U)}{2^N - 1}$ , where  $N$  is the number of bit-widths used. The operation  $\lfloor \cdot \rceil$  represents the rounding process.  $U_q$  is the quantized version, and  $\hat{U}$  is the de-quantized version of the original matrix  $U$ .