# Low-Rank Prehab: Preparing Neural Networks for SVD Compression

Haoran Qin, Shansita Sharma, Ali Abbasi, Chayne Thrash, and Soheil Kolouri

Department of Computer Science, Vanderbilt University, TN, USA

*Abstract*— Low-rank approximation methods such as singular value decomposition (SVD) and its variants (e.g., Fisher-weighted SVD, Activation SVD) have recently emerged as effective tools for neural network compression. In this setting, decomposition acts as a "surgical" intervention, followed by fine-tuning that serves as "rehab" to recover accuracy. Inspired by prehabilitation in surgery, we introduce a pre-compression fine-tuning stage, *Low-Rank Prehab*, that explicitly encourages low-rank structure in weight matrices while preserving task performance. By conditioning the model before SVD, Prehab steers weights toward spectrally compact regions of the parameter space, enabling smoother low-rank approximation and improved recovery. Experiments on large language models (LLMs) and other Transformer-based architectures, including Vision Transformers (ViTs), show that Prehab substantially reduces the immediate accuracy drop after compression and consistently improves post-finetuning performance. Across a wide range of compression ratios, our method outperforms state-of-the-art SVD-based techniques such as SVD-LLM, highlighting the importance of preparing models for compression rather than only improving the compression and recovery stages. Source code is available at **https://github.com/niqretnuh/PREHAB-SVD**.

## I. INTRODUCTION

Recent advances in computer vision, language, and multimodal reasoning have been propelled by increasingly large Transformer-based architectures such as Vision Transformers (ViTs) and Large Language Models (LLMs). While scaling has delivered remarkable accuracy gains, it also brings heavy computational, storage, and deployment costs. Models with billions of parameters remain impractical for on-device inference or cross-agent communication. Consequently, efficient compression methods that shrink model size while preserving task fidelity are essential for real-world deployment.

Among available strategies, *low-rank approximation* via Singular Value Decomposition (SVD) stands out for its simplicity and compatibility with existing architectures. By decomposing a weight matrix $W \in \mathbb{R}^{m \times n}$ into two low-rank factors $U_r \Sigma_r V_r^\top$, SVD reduces both parameter count and multiply–accumulate cost. However, directly truncating singular values often incurs a sharp performance drop, prompting extensive research into post-training refinements that account for loss geometry or parameter sensitivity.

Most existing SVD-based approaches operate *after* training. Fisher-Weighted SVD (FWSVD) [1] and Generalized Fisher-Weighted SVD (GFWSVD) [2] apply weighting of reconstruction error based on parameter importance using diagonal or Kronecker-factored approximation of Fisher Information Matrix (FIM). Activation-SVD (ASVD) [3] and SVD-LLM [4] improve truncation by normalizing input activations through data whitening, aligning the singular

spectrum with the true loss (i.e., the activation subspaces). More recently, Dobi-SVD [5] learns heterogeneous per-layer ranks through differentiable truncation and reconstructs compressed weights using incremental PCA (IPCA) with a memory–rank bijection for remapping. Despite their differences, all these methods share a reactive design: compression is performed on a converged model, followed by fine-tuning ("rehab") to restore accuracy. This reactive pipeline limits attainable fidelity because the pre-compression weights were never optimized to be low-rank friendly.

We re-examine compression from a geometric perspective. Let $\mathcal{M}_{\mathcal{L}}$ denote the manifold of optimal solutions, i.e., weight configurations that achieve low training loss. As illustrated in Figure 1, standard SVD projects a converged weight matrix $W_{\text{orig}}$ onto a low-rank subspace, but the resulting point $W'_{\text{SVD}}$ may lie far from $\mathcal{M}_{\mathcal{L}}$. This misalignment explains the immediate degradation of compressed models.

Most existing approaches focus on optimizing the projection/truncation step, ignoring how the original weights, $W_{\text{orig}}$, are obtained. Our key insight is that compression loss often arises because training drifts away from the manifold of low-rank, near-optimal solutions. We address this by conditioning the model *before* compression, guiding its learning dynamics toward regions of the loss manifold, $\mathcal{M}_{\mathcal{L}}$, that are inherently low-rank. This encourages weights that naturally align with these directions to be more robust to truncation, reducing the gap between original and compressed models.

We propose **Low-Rank Prehab**, a pre-compression fine-tuning stage that prepares the model for subsequent SVD-based compression. Drawing analogy to medicine, SVD acts as the "surgery," post-compression fine-tuning as "rehab," and our pre-conditioning phase as "prehab." Prehab co-optimizes task loss with an auxiliary term for smooth rank surrogate, gradually steering weights along $\mathcal{M}_{\mathcal{L}}$ toward low-rank configurations while maintaining task performance.

**Contributions.** This work makes three main contributions:

1) **Pre-compression conditioning framework.** We introduce **Low-Rank Prehab**, a lightweight, architecture-agnostic stage that shapes the model's geometry for smoother SVD compression.

2) **Geometric interpretation.** We formalize model compression as projection from the manifold of high-performing solutions onto a low-rank subspace and show that pre-conditioning reduces this manifold–subspace gap.

3) **Empirical validation.** Experiments on Llama, ViT, and BERT demonstrate that Prehab reduces immediate post-compression accuracy loss and improves final re-
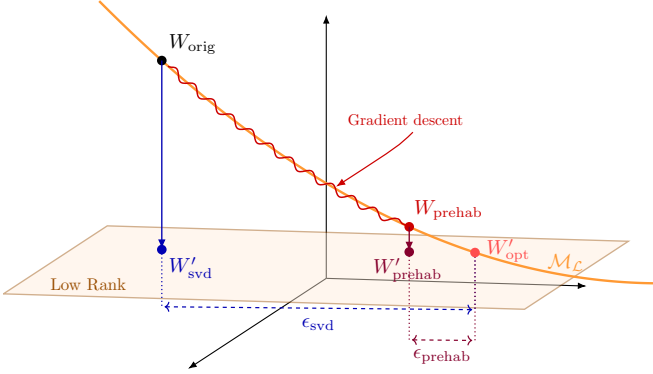
Fig. 1: Geometric view of the proposed **Low-Rank Prehab**. The original weights $W_{\mathrm{orig}}$ lie on the manifold $\mathcal{M}_{\mathcal{L}}$ of optimal solutions. Standard SVD projects $W_{\mathrm{orig}}$ directly onto the low-rank plane, producing $W'_{\mathrm{SVD}}$ far from the optimal intersection $W'_{\mathrm{opt}}$. Prehab first moves the weights along $\mathcal{M}_{\mathcal{L}}$ via joint optimization of task loss with a smooth rank regularizer, yielding $W_{\mathrm{prehab}}$. Its SVD projection $W'_{\mathrm{prehab}}$ lies closer to $W'_{\mathrm{opt}}$, reducing the geometric deviation $\epsilon_{\mathrm{prehab}} \ll \epsilon_{\mathrm{SVD}}$ and consequently minimizing post-compression loss.

covery, outperforming SVD-LLM across compression ratios.

## II. RELATED WORK

### A. Fisher-Aware SVD Methods

**Fisher-Weighted SVD (FWSVD)** replaces the Frobenius reconstruction norm with a Fisher-weighted metric, emphasizing parameters with high task sensitivity. A diagonal Fisher approximation, estimated from gradients on a calibration set, provides importance weighting but neglects correlations. **Generalized Fisher-Weighted SVD (GFWSVD)** extends this formulation using a Kronecker-factored (K-FAC) Fisher approximation, preserving structured dependencies between input and output dimensions. GFWSVD achieves higher fidelity under aggressive compression (up to 20×) but, like FWSVD, acts purely after convergence and cannot promote compressibility during training.

### B. Truncation- and Loss-Aligned SVD Methods

**SVD-LLM** links compression loss to the singular spectrum of each layer. It first whitens activations via a Cholesky factor of their covariance so that singular values directly reflect expected loss contributions, i.e., the compression objective is aligned with the task loss. SVD is then applied in the whitened space, followed by a closed-form layerwise correction or optional LoRA fine-tuning.

**Dobi-SVD** advances post-training SVD through *heterogeneous rank selection*, *incremental PCA (IPCA) recovery*, and *post-decomposition remapping*. Instead of a uniform rank, each layer learns its optimal compression ratio via differentiable truncation stabilized by truncated Taylor expansions. Compressed weights are reconstructed analytically using Eckart–Young–Mirsky–optimal updates with IPCA, while a memory–rank bijection "remapping" ensures information-preserving remapping under strict storage constraints. Despite strong empirical results, Dobi-SVD remains

TABLE I: Core objectives for SVD-based compression methods in related work. Here, $F$ denotes the diagonal Fisher information matrix used in FWSVD; $S = X^{\top} X$ is the activation covariance matrix (whitening in SVD-LLM); and $F \approx G \otimes A$ is the Kronecker-factored Fisher approximation in GFWSVD. $A_i = x_i W$, where $V_{A_i}$ are the right singular vectors of $A_i$, and $G_k = \mathrm{diag}(\mathbf{1}_k, \mathbf{0}_{d-k})$.

| Method | Optimization (rank-$r$ approximation) |
|---|---|
| **SVD** | $\widehat{W}_{\mathrm{SVD}} = \underset{\mathrm{rank}(\widehat{W}) \leq r}{\arg\min} \ \| W - \widehat{W} \|_{\mathsf{F}}^2 = U_r \Sigma_r V_r^{\top}$ |
| **FWSVD** | $\widehat{W}_{\mathrm{FWSVD}} = \underset{\mathrm{rank}(\widehat{W}) \leq r}{\arg\min} \ \| F^{1/2} \odot (W - \widehat{W}) \|_{\mathsf{F}}^2$ |
| **SVD-LLM** | $\widehat{W}_{\mathrm{SVD\text{-}LLM}} = \underset{\mathrm{rank}(\widehat{W}) \leq r}{\arg\min} \ \| (W - \widehat{W}) S^{1/2} \|_{\mathsf{F}}^2$ |
| **GFWSVD** | $\widehat{W}_{\mathrm{GFWSVD}} = \underset{\mathrm{rank}(\widehat{W}) \leq r}{\arg\min} \ \| A^{1/2} (W - \widehat{W}) G^{1/2} \|_{\mathsf{F}}^2$ |
| **Dobi-SVD** | $\widehat{W}_{\mathrm{Dobi}} = \underset{\mathrm{rank}(\widehat{W}) \leq k}{\arg\min} \ \sum_{i=1}^{n} \| x_i W V_{A_i} G_k V_{A_i}^{\top} - x_i \widehat{W} \|_2^2$ *Differentiable truncation of activations* learns $k$; *IPCA* provides an optimal Eckart–Young–Mirsky recovery of $\widehat{W}$. |

a post-training method and does not alter the model's pre-compression geometry.

### C. From Post-Projection to Pre-Conditioning

All current SVD-based approaches focus on improving the "surgery"—the decomposition and post-hoc recovery—yet leave the model unprepared for compression. **Low-Rank Prehab** fills this gap by conditioning the model before compression through a brief optimization phase that jointly minimizes task loss and a rank surrogate. This shifts weights toward low-rank-compatible regions of $\mathcal{M}_{\mathcal{L}}$, yielding smaller projection error and greater stability after truncation. Prehab is orthogonal to Fisher- or truncation-aware techniques and can serve as a universal pre-processing stage preceding any SVD variant. By introducing this "pre-surgical" preparation, we show that preparing a model for compression can be as crucial as the compression itself.

## III. METHOD

Previous SVD-based methods treat each trained weight matrix $W$ as static and project it directly onto the low-rank manifold. This purely algebraic step ignores the loss landscape, so large projection errors can induce substantial performance degradation. In contrast, and building on the geometric intuition of Section I, we introduce **Low-Rank Prehab**, a lightweight *pre-conditioning* stage applied before compression. Instead of reacting to post-SVD compression loss, Prehab proactively steers network parameters toward the low-rank manifold while keeping the loss low.

Formally, given pretrained weights $\{W_\ell\}_{\ell=1}^{L}$, the goal is to obtain perturbed weights $\{W_\ell^{\mathrm{prehab}}\}$ that (i) maintain task accuracy and (ii) lie closer to the manifold of low-rank matrices under the same Fisher-weighted geometry used in SVD-LLM. Recall that the rank of $W$ is the number of nonzero singular values, $\mathrm{rank}(W) = \|\sigma(W)\|_0$, where $\sigma(W)$ denotes the vector of singular values. Since the $\ell_0$ norm is non-differentiable, we optimize a smooth surrogate $\mathcal{R}_{\mathrm{rank}}(W_\ell X_\ell)$ defined over the Fisher-whitened matrices

$W_\ell X_\ell$, where $X_\ell$ is the Cholesky factor of the (uncentered) activation covariance used in SVD-LLM. The overall Prehab objective combines task fidelity and rank regularization:

$$\mathcal{L}_{\text{prehab}} = \mathcal{L}_{\text{task}}(\{W_\ell X_\ell X_\ell^{-1}\}_{\ell=1}^L) + \lambda \sum_{\ell=1}^L \mathcal{R}_{\text{rank}}(W_\ell X_\ell), \tag{2}$$

where $\lambda > 0$ controls the trade-off between compressibility and accuracy, and layer-specific coefficients $\lambda_\ell$ may optionally be used. The factor $X_\ell X_\ell^{-1}$ preserves functional equivalence while enabling gradients to act on the Fisher-weighted representation $W_\ell X_\ell$.

### A. Prehab Loss Function

The Prehab objective consists of a standard task loss, evaluated on the uncompressed model, and an auxiliary spectral regularizer applied to the whitened weights:

$$\mathcal{L}_{\text{task}} = \mathbb{E}_{(x,y)\sim\mathcal{D}}\left[\ell(f(x; \{W_\ell X_\ell X_\ell^{-1}\}_{\ell=1}^L), y)\right]. \tag{3}$$

Here, $\mathcal{L}_{\text{task}}$ measures predictive fidelity, while $\mathcal{R}_{\text{rank}}$ penalizes high-rank structure in the Fisher-weighted space. In practice, we optimize $\mathcal{L}_{\text{prehab}}$ for a few fine-tuning epochs using AdamW, freezing $\{X_\ell\}$ and letting gradients propagate through the round-trip $X_\ell X_\ell^{-1}$ transformation.

### B. Rank Surrogates: $\ell_1$-Norm and Stable Rank

Because $\text{rank}(W)$ is not differentiable, we consider two smooth surrogates defined on the whitened matrix $WX$:

*a) $\ell_1$-norm of singular values.:*

$$\mathcal{R}_{\ell_1}(WX) = \sum_i |\sigma_i(WX)|, \tag{4}$$

which imposes an $\ell_1$ penalty on the singular spectrum, promoting sparsity but requiring explicit SVD and scaling poorly for large models.

*b) Stable rank [6]:*

$$\text{SRank}(WX) = \frac{\|WX\|_*^2}{\|WX\|_F^2}, \tag{5}$$

a smooth, scale-invariant proxy for true rank that depends only on matrix norms and can be efficiently approximated using randomized trace estimators or sketching methods, often in nearly linear time in the matrix size. In our experiments, we rely on such approximations, making stable-rank regularization practical even for large models. Its gradient with respect to $W$ is

$$\nabla_W \text{SRank}(WX) = 2(WX)\left(\frac{\Sigma_X}{\|WX\|_F^2} - \text{SRank}(WX)\,I\right)X^\top, \tag{6}$$

where $\Sigma_X$ contains the singular values of $WX$. Minimizing $\text{SRank}(WX)$ directly shapes the Fisher-weighted spectrum that SVD-LLM will later truncate.

---

**Algorithm 1** Fisher-Aligned Low-Rank Prehab

---

**Require:** Pretrained weights $\{W_\ell^{(0)}\}$, whitening matrices $\{X_\ell\}$, dataset $\mathcal{D}$, rank coefficient $\lambda$, epochs $T_{\text{prehab}}$, learning rate $\eta$

1: **procedure** PREPROCESSING (PREHAB)
2:     **for** $t = 1$ **to** $T_{\text{prehab}}$ **do**
3:         Sample minibatch $(x, y) \sim \mathcal{D}$
4:         $\mathcal{L}_{\text{task}} \leftarrow \ell(f(x; \{W_\ell^{(t)} X_\ell X_\ell^{-1}\}_{\ell=1}^L), y)$
5:         $\mathcal{R}_{\text{rank}} \leftarrow \sum_\ell \text{SRank}(W_\ell^{(t)} X_\ell)$
6:         $W^{(t+1)} \leftarrow W^{(t)} - \eta \nabla_{W^{(t)}}(\mathcal{L}_{\text{task}} + \lambda \mathcal{R}_{\text{rank}})$
7:     **end for**
8: **end procedure**
9: **procedure** POSTPROCESSING (SURGERY & REHAB)
10:     **Surgery:** SVD-based compression (e.g., SVD-LLM)
11:     **Rehab:** Fine-tune to recover residual loss
12: **end procedure**

---

### C. Tunable Task–Compression Trade-Off via $\lambda$

The hyperparameter $\lambda$ governs the trade-off between accuracy and compressibility. Small $\lambda$ values preserve task fidelity, while larger ones promote stronger low-rank alignment by enhancing singular-value decay in $WX$. Empirically, logarithmic sweeps of $\lambda$ produce stable post-compression results. From a geometric standpoint, $\lambda$ controls how far Prehab moves $W$ along the solution manifold $\mathcal{M}_\mathcal{L}$ toward the Fisher-weighted low-rank manifold $\mathcal{M}_r^X$.

### D. Algorithmic Procedure

The complete Prehab process is outlined in Algorithm 1. It alternates between minimizing task loss and applying Fisher-weighted rank regularization. After pre-conditioning, SVD-LLM truncation operates on a smoother spectrum with minimal degradation.

### E. Summary of the Prehab Framework

Low-Rank Prehab serves as a Fisher-aligned pre-conditioning mechanism that harmonizes training geometry with compression geometry. By minimizing the stable rank of $WX$ while preserving the functional mapping $f(x; WXX^{-1})$, it teaches the model to inhabit a subspace where subsequent SVD truncation naturally aligns with task-relevant directions. The procedure requires no architectural change, adds negligible cost, and integrates seamlessly with post-training compressors such as SVD-LLM.

## IV. EXPERIMENTS

### A. Vision Transformer Evaluation

We evaluate Prehab-SVD on the ViT-B/16 [7] architecture to examine its applicability beyond language models. Training follows the standard ImageNet-1K [8] setup with $\ell_1$-norm regularization on the singular values and cross-entropy loss. Each model is fine-tuned for 500 steps with batch size 64 and $\lambda = 10^{-1}$. For LoRA fine-tuning, we follow the SVD-LLM protocol of alternating $U$, $V$ matrix updates using LoRA of rank $r = 10$, batch size 64, and 100 steps per factor.

**Results without LoRA.** Table II shows the Top-1 accuracy of ViT-B under varying compression ratios. Prehab-SVD consistently improves over SVD-LLM across all compression levels, demonstrating its effectiveness as a spectral pre-conditioning step. At moderate, 50% compression, Prehab-SVD improves accuracy from 58.42% (SVD-LLM) to 64.06%, yielding a +5.64% improvement. The benefit becomes more pronounced as compression increases: at 60% compression, Prehab-SVD reaches 48.12% accuracy compared to 28.38% for SVD-LLM, corresponding to a substantial +19.74% improvement. These results show that pre-conditioning becomes stronger in higher spectral sparsity regimes.

**Results with LoRA fine-tuning.** When applying the lightweight LoRA fine-tuning stage (denoted by $^\dagger$), Prehab-SVD retains its advantage across all configurations. Although the relative gains narrow at lower compression levels, Prehab-SVD$^\dagger$ continues to outperform SVD-LLM$^\dagger$. At 50% compression, Prehab-SVD$^\dagger$ improves accuracy from 66.41% (SVD-LLM$^\dagger$) to 68.59%, yielding a +2.18% gain. At 60% compression, Prehab-SVD$^\dagger$ reaches 58.97%, compared to 50.06% for SVD-LLM$^\dagger$, achieving a +8.61% improvement. This indicates that pre-conditioning helps LoRA recover lost performance more effectively, yielding better convergence and spectral alignment.

**Prehab-SVD significantly improves performance at high compression rates.** Figure 2 reports the percentage gain of Prehab-SVD over SVD-LLM, computed as (Prehab − SVD-LLM)/SVD-LLM to reveal the true magnitude of improvement. We observe that performance gains grow exponentially with compression rate, reflecting improved robustness to aggressive compression and highlighting the importance of *prehab* at such magnitudes of compression.

$\lambda$ **Ablation.** We further validate the effect of differing $\lambda$ on ViT performance. To investigate this, we set $\lambda = 10^{-1}$, $\lambda = 5.0$, and $\lambda = 10.0$ for multiple compression ratios and showcase the results in Table V. We observe that $\lambda = 10.0$ notably enhances performance under the highest compression, while introducing trade-offs at lower
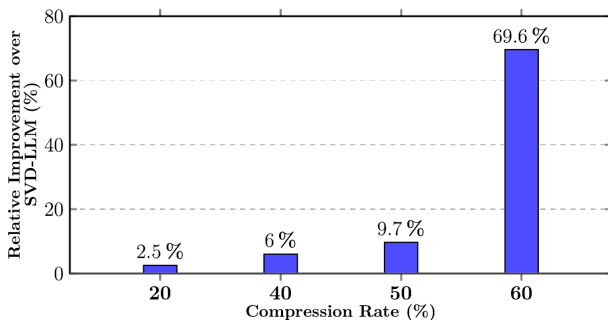


Fig. 2: **Relative improvement of Prehab-SVD over SVD-LLM on ViT-B.** Gains are computed as (Prehab − SVD-LLM)/SVD-LLM for each compression rate. Performance gains increase with compression, indicating enhanced robustness of Prehab-SVD at higher compression rates.

TABLE II: ImageNet Top-1 (%) for ViT-B across truncation levels (higher is better).

| Ratio | Method | Top-1 Accuracy (%) |
|---|---|---|
| **0% (Baseline)** | Original ViT-B | 81.10 |
| 20% | SVD | 70.24 |
| | FWSVD | 71.90 |
| | GFWSVD | 76.75 |
| | SVD-LLM | 76.56 |
| | Prehab-SVD | **78.49** (↑1.93) |
| | SVD-LLM$^\dagger$ | 77.51 |
| | Prehab-SVD$^\dagger$ | **78.68** (↑1.17) |
| 40% | SVD | 39.60 |
| | FWSVD | 53.93 |
| | GFWSVD | 68.95 |
| | SVD-LLM | 67.83 |
| | Prehab-SVD | **71.92** (↑4.09) |
| | SVD-LLM$^\dagger$ | 72.95 |
| | Prehab-SVD$^\dagger$ | **74.41** (↑1.46) |
| 50% | SVD | 10.06 |
| | FWSVD | 32.61 |
| | GFWSVD | 60.48 |
| | SVD-LLM | 58.42 |
| | Prehab-SVD | **64.06** (↑5.64) |
| | SVD-LLM$^\dagger$ | 66.41 |
| | Prehab-SVD$^\dagger$ | **68.59** (↑2.18) |
| 60% | SVD | 0.52 |
| | FWSVD | 11.95 |
| | GFWSVD | 43.18 |
| | SVD-LLM | 28.38 |
| | Prehab-SVD | **48.12** (↑19.74) |
| | SVD-LLM$^\dagger$ | 50.06 |
| | Prehab-SVD$^\dagger$ | **58.97** (↑8.91) |

$^\dagger$ Lightweight LoRA fine-tuning ("+LoRA"). Green ↑ indicates Top-1 accuracy improvement vs. SVD-LLM at the same ratio. "Ratio" = fraction of parameters removed in compressed blocks; Groups: *Baselines* (SVD/FWSVD/GFWSVD), *Loss-aligned* (SVD-LLM, Prehab-SVD), and *+LoRA variants*.

compression ratios. This behavior reflects the importance of identifying a *sweet spot* for the prehab regularization strength $\lambda$, balancing compression-target alignment and task fidelity for optimal performance.

**Runtime.** All ViT-B experiments were trained on a single NVIDIA A6000 GPU. The Prehab stage required only 7.5 minutes in total, averaging 0.9 seconds per step due to the small additional computational overhead of the rank-regularization term, confirming that Prehab introduces negligible additional cost relative to baseline SVD-based pipelines while achieving significant gains in performance.

*B. BERT-Base Evaluation*

We evaluate Prehab-SVD on BERT-Base [9] over the GLUE benchmark [10] to examine its effectiveness on smaller transformer models for natural language understanding. For each GLUE task, prior to applying any compression, the BERT-Base model was finetuned on the corresponding dataset for 3 epochs with a cosine scheduler with an initial learning rate $5 \times 10^{-5}$ for CoLA, and $2 \times 10^{-5}$ for the remaining tasks. Prehab was applied as a short pre-conditioning

TABLE III: Performance of BERT-Base on GLUE tasks across truncation levels. Best results per block are in **bold**.

| Ratio | Method | CoLA | SST-2 | MRPC | RTE | STS-B | QNLI | QQP | MNLI-m / MNLI-mm | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| 0% | BERT-Base | 0.567 | 0.925 | 0.898 | 0.658 | 0.857 | 0.903 | 0.876 | 0.827 / 0.832 | 0.816 |
| 20% | SVD | 0.021 | 0.841 | 0.684 | 0.520 | 0.616 | 0.717 | 0.727 | 0.579 / 0.583 | 0.588 |
| | FWSVD | 0.287 | 0.890 | 0.825 | 0.552 | 0.821 | 0.862 | 0.838 | 0.705 / 0.698 | 0.722 |
| | GFWSVD | 0.399 | 0.888 | 0.855 | 0.523 | 0.644 | 0.828 | 0.828 | 0.718 / 0.723 | 0.711 |
| | SVD-LLM | 0.452 | 0.898 | 0.884 | 0.527 | 0.847 | 0.871 | 0.847 | 0.808 / 0.812 | 0.767 |
| | Prehab-SVD | **0.492** | **0.914** | **0.901** | **0.614** | **0.852** | **0.876** | **0.853** | **0.812 / 0.818** | **0.790** (↑0.023) |
| 40% | SVD | 0.000 | 0.789 | 0.075 | 0.509 | 0.437 | 0.645 | 0.596 | 0.451 / 0.456 | 0.438 |
| | FWSVD | 0.000 | 0.828 | 0.823 | 0.502 | 0.810 | 0.817 | 0.801 | 0.546 / 0.549 | 0.641 |
| | GFWSVD | 0.219 | 0.852 | 0.724 | 0.462 | 0.601 | 0.809 | 0.809 | 0.632 / 0.633 | 0.638 |
| | SVD-LLM | 0.410 | 0.891 | 0.869 | 0.486 | 0.809 | 0.784 | 0.829 | 0.781 / 0.772 | 0.732 |
| | Prehab-SVD | **0.421** | **0.905** | **0.887** | **0.581** | **0.821** | **0.818** | **0.838** | **0.790 / 0.789** | **0.758** (↑0.026) |
| 60% | SVD | 0.000 | 0.627 | 0.000 | 0.523 | 0.114 | 0.561 | 0.498 | 0.356 / 0.352 | 0.335 |
| | FWSVD | 0.000 | 0.799 | 0.811 | 0.516 | 0.680 | 0.667 | 0.710 | 0.445 / 0.447 | 0.584 |
| | GFWSVD | 0.072 | 0.737 | 0.007 | 0.469 | 0.677 | 0.670 | 0.680 | 0.472 / 0.485 | 0.474 |
| | SVD-LLM | 0.259 | 0.836 | 0.830 | 0.520 | 0.662 | 0.651 | 0.763 | 0.658 / 0.659 | 0.647 |
| | Prehab-SVD | **0.309** | **0.877** | **0.847** | **0.567** | **0.681** | **0.671** | **0.771** | **0.698 / 0.695** | **0.677** (↑0.030) |

We report F1 for QQP/MRPC, Spearman for STS-B, and Acc. for the others; green ↑ marks performance relative to SVD-LLM at the same compression.

phase prior to SVD-LLM based compression for 1 epoch with batch size 64 and learning rates identical to the initial finetuning phase, jointly optimizing cross-entropy loss with the stable-rank regularizer with regularization coefficient $\lambda = 10^{-1}$.

**Effect without LoRA Fine-tuning.** Table III presents the performance gain of applying Prehab prior to compression under different compression ratios. Without any LoRA adaptation, Prehab-SVD consistently improves over the post-training baseline across most scenarios, confirming that pre-conditioning enhances low-rank compatibility and reduces performance loss after compression. Notably, our method yields substantial gains at 60% compression, especially on tasks such as RTE and CoLA, relative to the baseline SVD-LLM that we build upon. On average, Prehab-SVD offers strong performance gains across compression rates on BERT.

**Effect with LoRA Fine-tuning.** After applying lightweight LoRA recovery, the benefits of Prehab-SVD become less pronounced on BERT, as LoRA's parameter-efficient updates largely replicate prehab's low-rank regularization effects, reducing the relative impact of pre-conditioning. Results for this setting are omitted due to space limitations.

**Runtime.** All BERT experiments were trained on a single NVIDIA A6000 GPU. The Prehab stage was applied with the fast stable rank surrogate and required varying durations dependent on the size of the dataset, at the average rate of 1.1 seconds for each batch, confirming the negligible overhead.

### C. LLaMA-7B Evaluation

We further evaluate **Prehab-SVD** on **LLaMA-7B** to assess its effectiveness on large-scale transformer architectures. Following the experimental design of SVD-LLM, all methods are tested under identical compression ratios on LlaMa-7B [11] and evaluated on both **WikiText-2** [12] and **C4** [13]. Prehab was applied as a short pre-conditioning phase using $\lambda = 10^{-3}$ and learning rate $10^{-5}$, jointly optimizing cross-entropy loss and the stable-rank regularizer for 3 epochs.

Each prehab run used 256 WikiText-2 samples; LoRA fine-tuning (when applicable) was performed on 50k Alpaca [14] samples.

**Effect without LoRA Fine-tuning.** Table IV summarizes the performance of various SVD-based compression algorithms across compression ratios. Without any LoRA adaptation, Prehab-SVD consistently improves over all post-training baselines on WikiText-2. At 20% compression, Prehab-SVD achieves an 11.5% perplexity reduction versus SVD-LLM (7.03 vs. 7.94), and at 60% compression, a 29.8% reduction (46.75 vs. 66.62). However, on the larger and more diverse C4 dataset, Prehab-SVD without LoRA shows mild overfitting to WikiText-2, with degradation at lower compression levels (e.g., +131.3% at 20% ratio).

**Effect with Lightweight LoRA Recovery.** After applying a lightweight LoRA fine-tuning ("rehab") stage on Alpaca, Prehab-SVD preserves its performance advantage while regaining generalization on C4. For instance, Prehab-SVD† improves over SVD-LLM† by 6.3% on WikiText-2 and 10.1% on C4 at 20% compression, and by 12.1% and 10.1% respectively at 60%. These results demonstrate that pre-conditioning yields benefits beyond immediate compression fidelity, enabling faster and more stable recovery during downstream adaptation.

**Summary.** Across all ratios, Prehab-SVD narrows the perplexity gap between compressed and original models while improving robustness to data distribution shift. The combination of spectral pre-conditioning and lightweight LoRA recovery achieves state-of-the-art efficiency–performance trade-offs for large-scale low-rank adaptation.

**Runtime and Computational Cost.** All experiments were conducted on a single NVIDIA H200 GPU. The *Prehab* stage was lightweight, requiring only one epoch of fine-tuning on WikiText-2, which completed in approximately 10 minutes. LoRA fine-tuning on the Alpaca dataset required 30 minutes per epoch and was trained for 6 epochs in total—three epochs

TABLE IV: Perplexity (↓) of LLaMA-7B compressed by various SVD-based methods (lower is better).

| Ratio | Method | WikiText-2 ↓ | C4 ↓ |
|---|---|---|---|
| **0%** | Original | 5.68 | 7.34 |
| 20% (10.2 GB) | SVD | 20061 | 18800 |
| | FWSVD | 1727 | 1511 |
| | ASVD | 11.14 | 15.93 |
| | SVD-LLM | 7.94 | 15.84 |
| | Prehab-SVD | **7.03** (↓11.5%) | 36.64 (↑131.3%) |
| | SVD-LLM$^\dagger$ | 7.73 | 12.23 |
| | Prehab-SVD$^\dagger$ | 7.24 (↓6.3%) | **10.99** (↓10.1%) |
| 40% (7.76 GB) | SVD | 52489 | 47774 |
| | FWSVD | 18156 | 12847 |
| | ASVD | 1407 | 1109 |
| | SVD-LLM | 13.73 | 75.42 |
| | Prehab-SVD | 11.12 (↓19.0%) | 118.24 (↑56.7%) |
| | SVD-LLM$^\dagger$ | 9.27 | 15.63 |
| | Prehab-SVD$^\dagger$ | **9.15** (↓1.3%) | **13.67** (↓12.5%) |
| 60% (5.35 GB) | SVD | 105474 | 106976 |
| | FWSVD | 32194 | 29292 |
| | ASVD | 57057 | 43036 |
| | SVD-LLM | 66.62 | 471.83 |
| | Prehab-SVD | 46.75 (↓29.8%) | 298.59 (↓36.7%) |
| | SVD-LLM$^\dagger$ | 15.00 | 26.26 |
| | Prehab-SVD$^\dagger$ | **13.19** (↓12.1%) | **23.60** (↓10.1%) |

$^\dagger$ Indicates lightweight LoRA fine-tuning ("+LoRA"). Green ↓ and red ↑ denote relative perplexity change vs. SVD-LLM at the same ratio, computed as (Prehab − SVD-LLM)/SVD-LLM. Groups: *Baselines* (SVD/FWSVD/ASVD), *Loss-aligned* (SVD-LLM, Prehab-SVD), and *+LoRA variants*.

TABLE V: **ImageNet Top-1 (%) for ViT-B/16 across truncation levels (higher is better).**

| Ratio | Method | Top-1 Accuracy (%) |
|---|---|---|
| **0% (Baseline)** | Original ViT-B/16 | 81.10 |
| 40% | SVD-LLM | 67.83 |
| | Prehab-SVD ($\lambda=10^{-1}$) | **71.92** (↑4.09) |
| | Prehab-SVD ($\lambda=5.0$) | 71.76 (↑3.93) |
| | Prehab-SVD ($\lambda=10.0$) | 55.66 (↓12.17) |
| 50% | SVD-LLM | 58.42 |
| | Prehab-SVD ($\lambda=10^{-1}$) | **64.06** (↑5.64) |
| | Prehab-SVD ($\lambda=5.0$) | 63.72 (↑5.3) |
| | Prehab-SVD ($\lambda=10.0$) | 51.61 (↓6.81) |
| 60% | SVD-LLM | 28.38 |
| | Prehab-SVD ($\lambda=10^{-1}$) | 48.12 (↑19.74) |
| | Prehab-SVD ($\lambda=5.0$) | **49.25** (↑20.87) |
| | Prehab-SVD ($\lambda=10.0$) | 43.59 (↑15.21) |
| 80% | SVD-LLM | 0.72 |
| | Prehab-SVD ($\lambda=10^{-1}$) | 2.90 (↑2.18) |
| | Prehab-SVD ($\lambda=5.0$) | 2.47 (↑1.75) |
| | Prehab-SVD ($\lambda=10.0$) | **13.74** (↑13.02) |

Green ↑/red ↓ show Top-1 accuracy change vs. SVD-LLM at the same truncation. "Ratio" = fraction of parameters removed in compressed blocks.

jointly stabilize multiple compression operators.

each for the $U$ and $V$ projection matrices—resulting in an overall runtime of roughly 3 hours. Compared to the multi-day fine-tuning schedules typical for large-scale compression pipelines, this demonstrates that Prehab introduces minimal overhead while providing substantial compression readiness and downstream generalization improvements.

## V. CONCLUSION

We introduced *Low-Rank Prehab*, a lightweight pre-conditioning stage that aligns training and compression geometries by minimizing a smooth, loss-aligned stable-rank surrogate on activation-weighted parameter weights while preserving task loss. Across ViT-B/16 (ImageNet-1K), BERT-Base (GLUE), and LLaMA-7B (WikiText-2/C4), Prehab consistently improves performance—with or without subsequent LoRA—at negligible additional cost.

**Future work.** Our results suggest several natural extensions: (i) *Heterogeneous-rank Prehab:* extending the framework to layer-specific rank targets as in Dobi-SVD, where $\lambda_\ell$ adapts to each layer's curvature or Fisher energy. This would couple Prehab with per-layer compression planning and better exploit cross-layer redundancy. (ii) *Adaptive regularization schedules:* learning $\lambda$ through bilevel or meta-optimization that directly minimizes post-compression validation loss, enabling automatic trade-off discovery between compressibility and task fidelity. (iii) *Unified pre-conditioning across modalities:* exploring Prehab's compatibility with other efficiency paradigms such as quantization, pruning, and KV-cache compression, where shared low-rank alignment could

## REFERENCES

[1] Y.-C. Hsu, T. Hua, S. Chang, Q. Lou, Y. Shen, and H. Jin, "Language model compression with weighted low-rank factorization," 2022. [Online]. Available: https://arxiv.org/abs/2207.00112

[2] V. Chekalina, D. Moskovskiy, D. Cherniuk, M. Kurkin, A. Kuznetsov, and E. Frolov, "Generalized fisher-weighted svd: Scalable kronecker-factored fisher approximation for compressing large language models," 2025. [Online]. Available: https://arxiv.org/abs/2505.17974

[3] Z. Yuan, Y. Shang, Y. Song, D. Yang, Q. Wu, Y. Yan, and G. Sun, "Asvd: Activation-aware singular value decomposition for compressing large language models," 2025. [Online]. Available: https://arxiv.org/abs/2312.05821

[4] X. Wang, Y. Zheng, Z. Wan, and M. Zhang, "Svd-llm: Truncation-aware singular value decomposition for large language model compression," 2025. [Online]. Available: https://arxiv.org/abs/2403.07378

[5] Q. Wang, J. Ke, M. Tomizuka, Y. Chen, K. Keutzer, and C. Xu, "Dobi-svd: Differentiable svd for llm compression and some new perspectives," 2025. [Online]. Available: https://arxiv.org/abs/2502.02723

[6] M. Rudelson and R. Vershynin, "Sampling from large matrices: An approach through geometric functional analysis," *J. ACM*, vol. 54, no. 4, p. 21–es, Jul. 2007. [Online]. Available: https://doi.org/10.1145/1255443.1255449

[7] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," 2021. [Online]. Available: https://arxiv.org/abs/2010.11929

[8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

[9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," 2019. [Online]. Available: https://arxiv.org/abs/1810.04805

[10] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, "Glue: A multi-task benchmark and analysis platform for natural language understanding," 2019. [Online]. Available: https://arxiv.org/abs/1804.07461

[11] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "Llama: Open and efficient foundation language models," 2023. [Online]. Available: https://arxiv.org/abs/2302.13971

[12] S. Merity, C. Xiong, J. Bradbury, and R. Socher, "Pointer sentinel mixture models," 2016. [Online]. Available: https://arxiv.org/abs/1609.07843

[13] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," 2023. [Online]. Available: https://arxiv.org/abs/1910.10683

[14] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto, "Stanford alpaca: An instruction-following llama model," https://github.com/tatsu-lab/stanford_alpaca, 2023.