# TARGETED LOW-RANK REFINEMENT: ENHANCING SPARSE NEURAL NETWORKS WITH PRECISION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Pruning is a widely used technique for compressing large neural networks that eliminates weights that have minimal impact on the model's performance. Current pruning methods, exemplified by magnitude pruning, assign an importance score to each weight based on its magnitude and remove weights with scores below a certain threshold. Nonetheless, these methods often create a gap between the original dense and the pruned sparse model, potentially impairing performance. Especially when the sparsity ratio is high, the gap becomes more pronounced. To mitigate this issue, we introduce a method to bridge the gap left by pruning by utilizing a low-rank approximation of the difference between the dense and sparse matrices. Our method entails the iterative refinement of the sparse weight matrix augmented by a low-rank adjustment. This technique captures and retains the essential information often lost during pruning, thereby improving the performance of the pruned model. Furthermore, we offer a comprehensive theoretical analysis of our approach, emphasizing its convergence properties and establishing a solid basis for its efficacy. Experimental results on LLaMa models validate its effectiveness on large language models across various pruning techniques and sparsity levels. Our method shows significant improvements: at 50% sparsity, it reduces perplexity by 53.9% compared to conventional magnitude pruning on LLaMa-7B. Furthermore, to achieve a specific performance target, our approach enables an 8.6% reduction in model parameters while maintaining a sparsity ratio of about 50%.

029 030 031

032

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

### 1 INTRODUCTION

Pruning is a crucial technique in the field of model compression, particularly for large language models (LLMs), which have become the cornerstone of natural language processing tasks (Devlin, 2018; Brown, 2020; Hoffmann et al., 2022). Pruning involves the removal of specific weights or parameters from the neural network that are considered to have minimal impact on the overall performance of the model (LeCun et al., 1989; Hassibi et al., 1993; Han et al., 2015b; Frankle & Carbin, 2018; Frankle et al., 2020). This technique is especially valuable for LLMs, which often contain billions of parameters and require substantial computational resources for both training and inference (Han, 2017; Touvron et al., 2023; Minaee et al., 2024).

041 One of the most representative pruning techniques is magnitude pruning, which removes weights that 042 have the smallest absolute values. This method is based on the assumption that smaller weights have 043 less effect on the network's overall performance (Han et al., 2015a). Furthermore, as model sizes 044 continue to grow, the number of redundant parameters also increases. For LLMs with billions of parameters, even half of the layers can be dropped without significantly affecting performance (Men et al., 2024; Fan et al., 2024). However, recent research shows that pruning can cause irreparable loss 046 of knowledge and performance drops, especially for difficult tasks, a phenomenon known as the Junk 047 DNA Hypothesis (Yin et al., 2024). This consistent degradation in performance is observed across a 048 spectrum of pruning methods, including magnitude pruning, SparseGPT (Frantar & Alistarh, 2023), and Wanda (Sun et al., 2023), and applies to both unstructured pruning and structured N:M pruning. 050

Conventional approaches to post-pruning recovery face a three-fold challenge: (1) Computational
 burden: Re-training (Frankle & Carbin, 2018; Xia et al., 2023; Kim et al., 2024a) and knowledge
 distillation (Hinton, 2015; Wan et al., 2024a; Muralidharan et al., 2024) methods are computationally
 expensive and time-consuming. (2) Data and model dependency: These techniques typically require

either extensive datasets or access to a high-performing teacher model, which may not always be
feasible. (3) Sparsity inconsistency: Recent low-rank approximation methods (Li et al., 2023;
Mozaffari et al., 2024; Zhang & Papyan, 2024) often fail to maintain a consistent sparsity pattern,
making them unsuitable for structured pruning, which is crucial for hardware efficiency.

058 In this study, we address the challenge of post-pruning recovery by approximating the dense matrix as the sum of an updated sparse matrix with a maintained sparsity pattern and a low-rank matrix. We 060 propose an iterative refinement process that concurrently updates the sparse matrix and the low-rank 061 component in a data-free manner. This approach effectively recovers crucial information typically 062 lost during pruning. Our method features an adaptive low-rank approximation that dynamically 063 adjusts to complement the sparse matrix, enabling efficient information recovery. Unlike traditional 064 techniques that rely on large datasets or high-performing teacher models, our approach operates directly on model weights, offering a computationally efficient and broadly applicable solution. This 065 approach aims to improve the performance of the pruned model without significantly increasing the 066 parameter count. Our method combines the advantages of both sparse and low-rank approximation, 067 ensuring the model maintains its efficiency while enhancing its accuracy. In addition to our empirical 068 findings, we provide a comprehensive theoretical analysis of the iterative refinement process, which 069 rigorously demonstrates the favorable convergence properties of our method.

Experimental results on the LLaMa models demonstrate the effectiveness of our approach. When applying 50% sparsity, our method achieves a 53.9% reduction in perplexity compared to conventional magnitude pruning. The benefits of our approach become increasingly evident as sparsity increases: at 60% sparsity, we observe a 92.0% decrease in perplexity, while at 70% sparsity, an impressive 99.6% reduction is achieved. These findings highlight the effectiveness of our method, especially in scenarios of high sparsity where traditional approaches typically face considerable challenges.

077

078

079

081

082

084

085

090

092

093 094

095

To summarize, the main contributions of this paper are:

- In this work, we bridge the gap between the original dense and pruned sparse model by leveraging a low-rank component. This approach effectively fills the gap left by pruning, enhancing the model's performance with minimal parameter increase.
- We develop an iterative algorithm that incrementally refines the sparse weight matrix and incorporates the low-rank approximation. By prioritizing the preservation of weight components associated with larger singular values, our method allows for a more aggressive reduction of less important components, leading to a more precise approximation.
  - We provide a thorough theoretical analysis of our proposed method, which offers a rigorous foundation understanding of the effectiveness, convergence, and stability of our approach.
- We validate the effectiveness of our method on LLMs across various sparsity levels, achieving substantial perplexity reductions compared to baselines. This is particularly noteworthy at high sparsity levels, maintaining robust improvements even as high as 70% sparsity.
- 2 PRELIMINARY

### 2.1 THE GENERAL FRAMEWORK OF LAYER-WISE PRUNING

Neural network pruning is a crucial technique for model compression, aiming to reduce parameters 096 while preserving performance. In this section, we first present the general framework for (local) layer-wise pruning, outlining the key steps and components involved and it is illustrated in Fig. 1(a). 098 Let  $W \in \mathbb{R}^{m \times n}$  represent a weight matrix, and  $M = \mathcal{M}(W, \mathcal{D}) \in \mathbb{R}^{m \times n}$  denote its corresponding importance score matrix. Here,  $\mathcal{M}: \mathbb{R}^{m \times n} \times \mathcal{D} \to \mathbb{R}^{m \times n}$  is a metric function that computes the 100 importance scores based on the weight matrix W and an optional dataset  $\mathcal{D}$ . The most straightforward 101 approach is to use the magnitude of the weights as the importance score, i.e.,  $M_{ij} = |W_{ij}|$ . More 102 sophisticated methods can be employed to capture the importance of each weight more accurately, 103 such as the SparseGPT (Frantar & Alistarh, 2023) and Wanda (Sun et al., 2023). Given a pruning 104 ratio  $p \in [0, 1]$ , we set threshold h as the p-th percentile of M as the decision boundary for pruning. 105 Then the binary pruning mask P is obtained by  $P_{ij} = \mathbb{I}(M_{ij} > h)$ . The pruned sparse weight matrix S is obtained by  $S = W \odot P$ , where  $\odot$  denotes element-wise multiplication. This zeroes out 106 less important weights while retaining significant ones. An optional weight update procedure can be 107 applied to refine the pruned weights further. For example, after the pruning process, a re-training



(a) The general framework of neural network pruning. (b) Refine with a low-rank patch.

Figure 1: (a) An overview of the general framework of (local) layer-wise pruning, where W, M, P, and S are the original dense weight matrix, the importance score matrix, the binary pruning mask, and the pruned sparse weight matrix, respectively. Here we show the case of pruning using the magnitude of the weights as the importance score. (b) Our proposed method, where S is updated and a low-rank matrix  $L_k$  is introduced to minimize the approximation error.

phase may be implemented to fine-tune the remaining weights and recover some lost performance. During this process, the pruning structure is maintained by enforcing the constraint  $S^{(t)} = S^{(t)} \odot P$ at each iteration t. This ensures that the pruned weights remain zero throughout the re-training process, preserving the sparsity achieved through pruning.

### 126 2.2 SINGULAR VALUE DECOMPOSITION

Here we briefly review the Singular Value Decomposition (SVD) of a matrix. Given a matrix  $W \in \mathbb{R}^{m \times n}$ , its reduced SVD is given by (Olver & Shakiban, 2018):

114

120

125

127

128

131 132  $\boldsymbol{W} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\top} = \sum_{i=1}^{r} \sigma_{i}\boldsymbol{u}_{i}\boldsymbol{v}_{i}^{\top}, \qquad (1)$ 

133 where  $U \in \mathbb{R}^{m \times r}$  and  $V \in \mathbb{R}^{r \times n}$  are orthogonal matrices, and  $\Sigma \in \mathbb{R}^{r \times r}$  is a diagonal matrix 134 containing the singular values of W. The *i*-th columns of matrices U and V are represented by  $u_i$ 135 and  $v_i$ , respectively. Additionally, we use  $\sigma_i$  to denote the *i*-th diagonal element of the matrix  $\Sigma$ . 136 For convenience, we extend this notation to represent functions that map W to its SVD components. 137 This allows us to express the decomposition as  $W = U(W)\Sigma(W)V^{\top}(W)$ , or alternatively as a 138 sum of outer products:  $W = \sum_{i=1}^{r} \sigma_i(W)u_i(W)v_i^{\top}(W)$ . While this slightly abuses the original 139 notation, it provides a concise way to refer to the SVD components of any given matrix.

### 3 Method

141 142

140

143 In this study, we focus on the local pruning method and aim to fill the gap between the original dense 144 weight matrix W and the pruned weight S by introducing a low-rank matrix  $L_k$  to approximate the 145 difference matrix  $\Delta W = W - S$ . We illustrate the main idea in Fig. 1(b), where the dense matrix 146 W is factorized into the superposition of a low-rank matrix L and a updated sparse matrix S':

147 148

$$W = S' + L \approx S' + L_k. \tag{2}$$

In this formulation, S' maintains the same sparsity pattern as S, ensuring that  $S' = S' \odot P$  and  $S = S \odot P$  are satisfied. It's important to note that both S' and L can have elements of any magnitude without restrictions.  $L_k$  is the best rank-k approximation of W, which can be obtained using the SVD of L as  $L_k = U_{:k}(L)\Sigma_{:k}(L)V_{:k}^{\top}(L)$ . Begin with fixed S, we have the baseline:

**Baseline 1 (Zero-shot SVD)** The most straightforward method to obtain the low-rank matrix  $L_k$ is to directly perform the SVD on W - S without updating S. This can be expressed as  $L_k = U_{:k}(W - S)\Sigma_{:k}(W - S)V_{:k}^{\top}(W - S)$ . Here,  $U_{:k}(\cdot)$ ,  $\Sigma_{:k}(\cdot)$ , and  $V_{:k}(\cdot)$  represent the first kcolumns of  $U(\cdot)$ , the top-left  $k \times k$  submatrix of  $\Sigma(\cdot)$ , and the first k columns of  $V(\cdot)$ , respectively.

**Parameter efficiency analysis.** We analyze the computational efficiency of low-rank refinement by examining its parameter count and FLOPs (floating-point operations). We compare these metrics with the dense model and a pruned model, considering the impact of sparsity. For a weight matrix  $W \in \mathbb{R}^{m \times n}$ , the dense model has mn parameters and requires 2mn FLOPs for a forward pass. With pruning at sparsity ratio p, these reduce to (1 - p)mn parameters and 2(1 - p)mn FLOPs.

187

188

199 200

203

204

205 206 207



Figure 2: Analysis of the residual matrix L = W - S' and its low-rank approximation  $L_k$  using different methods. Results are shown for zero-shot SVD, PCP baseline (T = 5000), and our proposed method with varying k (64, 128, 512) and T = 50. (a) Singular value spectrum of L. (b) Proportion of total energy captured by the top k singular values, calculated as  $\sum_{j=1}^{i} \sigma_j^2(L_k) / \sum_{j=1}^{r} \sigma_j^2(L)$ .

Our method introduces a low-rank matrix  $L_k = BA$ , which introduces k(m + n) parameters and 2k(m + n) FLOPs, slightly reducing the overall sparsity ratio by  $k(\frac{1}{m} + \frac{1}{n})$ . The choice of k and p presents a trade-off between model size, computational complexity, and performance. In practice, using a rank k of 128 results in only a 4.9% increase in parameters, while possibly reducing perplexity by half on a LLaMa-7B model with both unstructured 50% sparsity and structured 4:8 sparsity.

This optimization problem is generally NP-hard and similar challenges have been addressed in the
fields of matrix completion (Chandrasekaran et al., 2011) and robust Principal Component Analysis
(PCA) (Candès et al., 2011; Peng et al., 2020) by solving a Principal Component Pursuit (PCP) as:

$$\min_{\boldsymbol{L},\boldsymbol{S}'} \|\boldsymbol{L}\|_* + \lambda \|\boldsymbol{S}'\|_1, \quad \text{s.t. } \boldsymbol{L} + \boldsymbol{S}' = \boldsymbol{W}.$$
(3)

Where  $\|\cdot\|_*$  denotes the nuclear norm and  $\|\cdot\|_1$  denotes the  $\ell_1$ -norm,  $\lambda$  is a hyperparameter that 189 controls the trade-off between the rank of L and the sparsity of S'. The nuclear norm serves as a 190 convex approximation for the rank of a matrix, while the  $\ell_1$ -norm acts as a convex proxy for the 191  $\ell_0$ -norm (which represents the count of non-zero elements in a matrix). In this convex optimization 192 problem, the  $||L||_*$  term promotes a low-rank solution for L, whereas the  $\lambda ||S'||_1$  term promotes 193 sparsity in S'. Nevertheless, it's important to note that the solution to Eq.(3) does not necessarily 194 preserve the fixed sparsity pattern of S' that is required to match S in Eq.(2). To address this issue, 195 we propose to incorporate the binary mask P into the optimization process to ensure that the sparsity 196 pattern of S' is fixed as S. This is achieved by rewriting Eq.(2) as follows: 197

$$W = \underbrace{(\overline{W \odot P} - Q \odot P)}_{\text{sparse part } S'} + \underbrace{(W \odot (1 - P) + Q \odot P)}_{\text{low-rank part } L}.$$
(4)

Where Q is a learnable matrix in the same shape as W.

**Baseline 2 (PCP with mask)** By directly substituting the decomposition from Eq.(4) into the optimization problem presented in Eq.(3) yields:

$$\min_{\boldsymbol{Q}} \|\boldsymbol{W} \odot (1-\boldsymbol{P}) + \boldsymbol{Q} \odot \boldsymbol{P}\|_{*} + \lambda \|\boldsymbol{W} \odot \boldsymbol{P} - \boldsymbol{Q} \odot \boldsymbol{P}\|_{1}, and \ \boldsymbol{L} = \boldsymbol{W} \odot (1-\boldsymbol{P}) + \boldsymbol{Q} \odot \boldsymbol{P}.$$
(5)

This equation represents a constrained optimization problem where we seek to find the optimal matrix Q. The binary mask P plays a crucial role in maintaining the desired sparsity pattern. To solve this optimization problem, we can employ iterative methods such as gradient descent or its variants. Following Candès et al. (2011), we set  $\lambda = 1/\sqrt{\max(m, n)}$ . However, the nuclear norm minimization approach has a limitation: it applies equal shrinkage to all rank components (Zha et al., 2019). This uniform treatment may not be optimal for our objective in Eq.(2), where we aim to approximate the residual matrix W - S using a low-rank matrix  $L_k$  with a rank lower than that of L.

To address this limitation, we propose a forward-only method that prioritizes the preservation of low-rank components associated with larger singular values, while allowing for a more aggressive

216 reduction of components with smaller singular values. This approach offers greater flexibility and 217 precision in low-rank approximation compared to the PCP baseline, aligning more closely with our 218 goal of achieving an efficient low-rank approximation. We begin by setting  $S^{(0)} = W \odot P$ , and 219 then iteratively refine  $S^{(t+1)}$  using the following update rule: 220

$$\boldsymbol{S}^{(t+1)} = \boldsymbol{S}^{(t)} + \boldsymbol{P} \odot \left\{ \sum_{i=r^{(t)}+1}^{r} \sigma_i \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{u}_i \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{v}_i^{\top} \left( \boldsymbol{L}^{(t)} \right) \right\}$$
(6)

$$= \boldsymbol{S}^{(t)} + \boldsymbol{P} \odot \left\{ \boldsymbol{U}_{r^{(t)}} \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{\Sigma}_{r^{(t)}} \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{V}_{r^{(t)}}^{\top} \left( \boldsymbol{L}^{(t)} \right) \right\}.$$
(7)

In this equation,  $L^{(t)}$  represents the difference  $W - S^{(t)}$ , and  $r^{(t)}$  is defined as  $\lfloor 1 + \frac{k-1}{T-1}t \rfloor$ , where 227 T denotes the total number of iterations, k is the target rank of  $L_k$ , and t ranges from 0 to T-1. 228 We summarize the above algorithm in Algorithm. 1. In each iteration, our algorithm focuses on 229 eliminating the least significant components from the residual matrix  $L^{(t)}$ , which are associated with 230 smaller singular values, while preserving those with larger singular values. This approach allows for a 231 progressive refinement of the sparse matrix S. As we advance through the iterations, we gradually in-232 corporate more subtle details from the original dense weight matrix W into the updated sparse term S. 233

Incrementally increasing the rank  $r^{(t)}$  from 2 234 to k over T iterations enables a more nuanced 235 exploration of the weight space. Furthermore, 236 the consistent application of the binary mask P237 throughout this optimization process ensures the 238 preservation of the desired sparsity pattern.

239 In Figure 2, we show the singular value spec-240 trum of L and the cumulative energy reten-241 tion in  $L_k$  for different methods using a 512  $\times$ 242 512 weight block extracted from a fine-tuned 243 LLaMa-7B model. Subfigure 2a demonstrates 244 that the proposed method with different target 245 rank values k (64, 128, and 512), consistently produces a more pronounced decay in singular 246 values compared to both the zero-shot SVD and 247 the PCP baseline. Subfigure 2b shows that the 248

### Algorithm 1 Iterative Weight Update

- 1: **Inputs:** Dense weight matrix W, binary mask P, target rank k, number of iterations T
- 2: Initialize  $S^{(0)} \leftarrow W \odot P$
- 3: for t = 0 to T 1 do
- $oldsymbol{L}^{(t)} \leftarrow oldsymbol{W} oldsymbol{S}^{(t)}$ 4: Compute SVD:  $\boldsymbol{L}^{(t)} = \boldsymbol{U}^{(t)} \boldsymbol{\Sigma}^{(t)} \boldsymbol{V}^{(t)^{\top}}$
- 5: 6:
- $r^{(t)} \leftarrow \lfloor 1 + \frac{k-1}{T-1}t \rfloor$
- $m{S}^{(t+1)} \leftarrow m{S}^{(t)} + m{P} \odot \Big\{ m{U}^{(t)}_{r^{(t)}:} m{\Sigma}^{(t)}_{r^{(t)}:} m{V}^{(t)^{ op}}_{r^{(t)}:} \Big\}$ 7:

8: end for **T***TT*  $\mathbf{T}(T)$ 

9: 
$$\boldsymbol{L}^{(T)} \leftarrow \boldsymbol{W} - \boldsymbol{S}^{(T)}$$
  
10: Returns:  $\boldsymbol{S}^{(T)}, \boldsymbol{L}^{(T)}$ 

cumulative energy retention in L, calculated as  $E(i) = \sum_{j=1}^{i} \sigma_j^2(L) / \sum_{j=1}^{i} \sigma_j^2(L)$  or  $\|L_i\|_F^2 / \|L\|_F^2$ , 249 increases more rapidly with truncated rank i for the proposed method compared to the PCP baseline 2. 250

251 From Algorithm 1 and observations from Figure 2, we can summarize the technical contributions 252 of our proposed method as follows: (1) Iterative refinement with adaptive rank increase. Our 253 method iteratively refines the sparse matrix S while gradually increasing the rank  $r^{(t)}$  from 2 to 254 k over T iterations. This enables progressive elimination of less significant components from the 255 residual matrix  $L^{(t)}$ . (2) Sparsity-preserving. Throughout the optimization process, our method 256 consistently applies the binary mask P, ensuring the preservation of the desired sparsity pattern. 257 Sparsity-preserving enables compatibility with both unstructured and structured pruning techniques, facilitating hardware-friendly implementations, particularly for structured sparsity patterns. (3) 258 **Efficient information capture.** Our approach creates low-rank approximations that more efficiently 259 capture the essential information of the original weight matrix compared to zero-shot SVD. 260

261 262

225 226

#### 3.1 THEORETICAL ANALYSIS

In this section, we delve into the theoretical underpinnings of our proposed method, offering a more 264 rigorous analysis of its properties and performance. We present key theorems that elucidate the 265 behavior of our algorithm, focusing on two critical aspects: sparsity preservation and convergence. 266

267 First, we demonstrate that our method maintains the desired sparsity pattern throughout the iterative 268 process, ensuring that the final solution adheres to the specified binary mask. This property is crucial for applications where specific weight connections must remain zero.

Theorem 1 (Sparsity Preservation) For all iterations t, the sparsity pattern of  $S^{(t)}$  is preserved and matches the binary mask P, i.e.,  $S^{(t)} = S^{(t)} \odot P$  for all t. Refer to Proof. 5.

Second, we prove the convergence of our algorithm, showing that it approaches a well-defined
solution as the number of iterations increases. This convergence guarantee provides theoretical
justification for the stability and reliability of our method.

Theorem 2 (Convergence) For any weight matrix  $W \in \mathbb{R}^{m \times n}$  and binary mask  $P \in \{0, 1\}^{m \times n}$ , the iterative weight update algorithm converges to a solution  $(S^*, L^*)$  as  $T \to \infty$ , such that  $W = S^* + L^*$ ,  $S^* = S^* \odot P$ . Moreover,  $\lim_{T\to\infty} ||S^{(T)} - S^*||_F = 0$  and  $\lim_{T\to\infty} ||L^{(T)} - L^*||_F = 0$ .

**Theorem 3 (Asymptotic Convergence)** There exists a time step  $T_0$  such that for all  $t > T_0$ , the Frobenius norm of the error decreases monotonically, i.e., for some

$$\left\| \boldsymbol{W} - \left( \boldsymbol{S}^{(t+1)} + \boldsymbol{L}_{k}^{(t+1)} \right) \right\|_{F} \leq \left\| \boldsymbol{W} - \left( \boldsymbol{S}^{(t)} + \boldsymbol{L}_{k}^{(t)} \right) \right\|_{F}.$$
(8)

A corollary of Theorem. 3 is that if we choose to fix the  $r^{(t)}$  to be k for all t, the Frobenius norm of the error decreases monotonically. Because in this case, the term  $\|\boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)}\|_{F}^{2}$  vanishes, and Eq.(43) becomes  $\|\boldsymbol{E}^{(t+1)}\|_{F}^{2} \leq \|\boldsymbol{E}^{(t)}\|_{F}^{2} - \|\boldsymbol{P} \odot \boldsymbol{E}^{(t)}\|_{F}^{2}$ . Formally, we have the following corollary:

**Corollary 1 (Monotonic Improvement)** If we choose to fix the  $r^{(t)}$  to be k for all t, the Frobenius norm of the error decreases monotonically, i.e.,

$$\left\|\boldsymbol{E}^{(t+1)}\right\|_{F} - \left\|\boldsymbol{E}^{(t)}\right\|_{F} \leq -\left\|\boldsymbol{P} \odot \boldsymbol{E}^{(t)}\right\|_{F} \leq 0,$$
(9)

the equality holds if and only if  $E^{(t)} = 0$ .

However, our empirical observations suggest that fixing  $r^{(t)}$  to be k throughout the process is not optimal in practice, and gradually increasing  $r^{(t)}$  as the iterations progress leads to faster convergence.

**Theorem 4 (Error Bound)** At each iteration t, the Frobenius norm of the difference between the original weight matrix W and its approximation  $S^{(t)} + L_k^{(t)}$  is bounded by:

$$\left\| \boldsymbol{W} - \left( \boldsymbol{S}^{(t)} + \boldsymbol{L}_{k}^{(t)} \right) \right\|_{F} = \sqrt{\sum_{i=k+1}^{r} \sigma_{i}^{2} \left( \boldsymbol{L}^{(t)} \right)} \leq \sqrt{(r-k)} \sigma_{k+1} \left( \boldsymbol{L}^{(t)} \right), \tag{10}$$

where  $\sigma_i(\mathbf{L}^{(t)})$  are the singular values of  $\mathbf{L}^{(t)} = \mathbf{W} - \mathbf{S}^{(t)}$ , and  $r = \operatorname{rank}(\mathbf{L}^{(t)})$ .

For a more comprehensive treatment, including detailed proofs of these theorems and additional supporting lemmas, we direct the reader to Appendix A. This appendix contains the full mathematical derivations and supplementary results that underpin our theoretical analysis.

4 EXPERIMENT

### 4.1 EXPERIMENTAL SETUP

We conducted our experiments using LLaMa models, evaluating their performance on the WikiText-2. Our primary metric for assessment was perplexity, where lower values indicate better performance. We tested various sparsity levels, including unstructured sparsity and structured N:M sparsity.

# 4.2 EFFECTIVENESS OF LOW-RANK REFINEMENT

First, we examine the benefits of low-rank refinement on the LLaMa-7B model by evaluating the WikiText-2 validation perplexity. Figure 3 illustrates the perplexity evaluation for various sparsity levels using magnitude pruning and low-rank refinement methods, with the target rank k set to 128. In subfigure 3a, we compare the perplexity of sparsity-only pruning and zero-shot SVD



Figure 3: Perplexity evaluation for different sparsity levels and methods, k = 128 for low-rank refinement. (a) Low-rank refinement advantage over sparsity-only pruning. (b) Comparison of our proposed iterative weight update method with other baseline methods.

Table 1: Comparison of WikiText validation perplexity ( $\downarrow$  is better) across various sparsity levels on LLaMa-7B. All methods and sparsity levels use a target rank of k = 128 (4.9% more parameters).

|                  | SPARSITY LEVEL (LLaMa-7B) |                |                 |                |                |  |
|------------------|---------------------------|----------------|-----------------|----------------|----------------|--|
| METHOD           | 50%                       | 60%            | 70%             | 4:8            | 2:4            |  |
| Dense            |                           |                | 5.68            |                |                |  |
| Magnitude        | 17.29 (0%)                | 152.36 (0%)    | 48427.85 (0%)   | 16.83 (0%)     | 42.53 (0%)     |  |
| w/ Zero-shot SVD | 8.06 (-53.4%)             | 13.59 (-91.1%) | 283.74 (-99.4%) | 9.29 (-44.8%)  | 12.72 (-70.0%) |  |
| w/ PCP with mask | 8.70 (-49.7%)             | 16.67 (-89.1%) | 727.54 (-98.5%) | 10.60 (-37.0%) | 16.62 (-60.9%) |  |
| w/ Ours          | 7.97 (-53.9%)             | 12.14 (-92.0%) | 200.09 (-99.6%) | 8.86 (-47.4%)  | 10.74 (-74.7%) |  |
| Wanda            | 7.26 (0%)                 | 10.69 (0%)     | 84.69 (0%)      | 8.57 (0%)      | 11.53 (0%)     |  |
| w/ Zero-shot SVD | 7.09 (-2.3%)              | 9.60 (-10.2%)  | 35.65 (-57.9%)  | 8.14 (-5.0%)   | 10.48 (-9.1%)  |  |
| w/ PCP with mask | 7.28 (+0.3%)              | 10.19 (-4.7%)  | 47.11 (-44.4%)  | 8.63 (-0.7%)   | 11.22 (-2.7%)  |  |
| w/ Ours          | 6.92 (-4.7%)              | 8.97 (-16.1%)  | 32.90 (-61.2%)  | 7.74 (-9.7%)   | 9.18 (-20.4%)  |  |

refinement to highlight the benefits of the low-rank refinement strategy. By incorporating a low-rank structure into the pruned weights, the parameter count increases, so the x-axis represents the parameter reduction relative to the dense model, and we use a dashed black arrow to indicate the correspondence between the pruned model and the model with refinement. The results indicate that we gain 8.6% parameter reduction at approximately 50% sparsity, highlighting the superiority of low-rank refinement compared to sparsity-only pruning techniques. In subfigure 3b, we compare our proposed method with other baseline methods. Our proposed method consistently outperforms the other two baseline methods, with the benefits becoming more pronounced at higher sparsity levels. It is evident that PCP with mask underperforms compared to zero-shot SVD and our iterative weight update method, although it still surpasses sparsity-only pruning. This can be attributed to PCP's uniform shrinkage of all singular values, including smaller ones, along with our choice of a very low target rank k = 128. As illustrated in Figure 4a, PCP also yields larger singular values for indices ranging from  $10^2$  to nearly  $10^3$ , demonstrating its less discriminative approach to value reduction.

To further analyze the performance of low-rank refinement with different sparsity types and levels, we provide a more detailed comparison of WikiText-2 validation perplexity across various sparsity levels and pruning methods in Tables 1 and 2. The tables include results for different sparsity levels, ranging from 50% to 70%, as well as structured sparsity patterns like 4:8 and 2:4. For each sparsity level, we compare the performance of low-rank refinement incorporated with magnitude pruning and Wanda pruning (Sun et al., 2023). At 50% sparsity, our method reduces perplexity by 53.9% compared to sparse-only magnitude pruning. This improvement becomes more significant at higher sparsity levels, reaching a 92.0% reduction at 60% sparsity and a 99.6% reduction at 70% sparsity. Our proposed method consistently achieves lower perplexity values across all sparsity levels. This detailed comparison highlights the effectiveness of our method in maintaining low perplexity even at higher sparsity levels, demonstrating its robustness and superiority over other baseline methods.

Table 2: Comparison of WikiText validation perplexity ( $\downarrow$  is better) across various sparsity levels on LLaMa-13B. All methods and sparsity levels use a target rank of k = 128 (3.8% more parameters).

|   | SPARSITY LEVEL (LLaMa-13B)                       |   |   |   |  |  |
|---|--|---|---|---|--|--|
| METHOD  | 50%  | 60%   | 70%   | 4:8   | 2:4  |  |
| Dense   |  |   | 4.57  |   |  |  |
| Magnitude<br>w/ Zero-shot SVD<br>w/ <b>Ours</b> | 5.98 (0%)<br>5.73 (-4.2%)<br><b>5.65 (-5.5%)</b> | 9.91 (0%)<br>8.83 (-10.9%)<br>8.83 (-10.9%) | 408.75 (0%)<br>163.96 (-59.9%)<br><b>99.27 (-75.7%)</b> | 6.76 (0%)<br>6.58 (-2.7%)<br><b>6.40 (-5.3%</b> ) | 8.32 (0%)<br>8.86 (-6.5%)<br><b>7.76 (-6.7%)</b> |  |

Table 3: Performance comparison on LLaMa models across several benchmark datasets (k = 128).

|     | METHOD   | TruthfulQA                  | GSM8K                    | ARC-c                       | MMLU                        | AVG.  |
|-----|--|-----------------------------|--------------------------|-----------------------------|-----------------------------|---|
|     | Dense  | 34.1                        | 10.3                     | 44.7                        | 32.1                        | 30.3 (0%)   |
| 7B  | <i>Magnitude 50%</i><br>w/ Zero-shot SVD<br>w/ <b>Ours</b> | <b>35.3</b><br>34.3<br>34.2 | 1.0<br>1.5<br><b>3.4</b> | 33.5<br>36.9<br><b>41.5</b> | 24.6<br>26.0<br>26.0        | 23.6 (-22.1%)<br>24.7 (-18.5%)<br>26.3 (-13.2%)         |
|     | Dense  | 36.9                        | 23.4                     | 49.1                        | 52.1                        | 40.4 (0%)   |
| 13B | Magnitude 2:4<br>w/ Zero-shot SVD<br>w/ <b>Ours</b>        | <b>38.4</b><br>37.6<br>36.9 | 1.7<br>1.8<br><b>9.4</b> | 34.9<br>32.2<br><b>36.7</b> | 27.8<br>27.0<br><b>41.9</b> | 25.7 (-36.4%)<br>24.7 (-38.9%)<br><b>31.2 (-22.8%</b> ) |

In Table 3, we evaluate the performance of our method on several benchmark datasets. We compare the performance of our method with that of dense models, sparsity-only pruned models, and low-rank refined using zero-shot SVD. Our method consistently outperforms both magnitude pruning and zero-shot SVD across most tasks and model sizes. For the 7B model, it achieves a 13.2% reduction in average performance compared to 22.1% for magnitude pruning and 18.5% for zero-shot SVD.

### 409 4.3 ITERATIVE WEIGHT UPDATE ANALYSIS

Here we analyze the behavior of our proposed iterative weight update algorithm (Algorithm 1). Our goal is to empirically validate the theoretical properties established in Section 3, specifically regarding convergence, error reduction, and emperically investigate the lower bound of rank(L).

Convergence and Error Analysis. In Figure 4, we visualize several key properties of the residual matrix L = W - S' and its low-rank approximation  $L_k$ , we use a log scale for the x-axis to better visualize the decay of singular values at low ranks. Throughout the analysis, we use magnitude pruning with a sparsity level of 50%. (a) Singular Value Spectrum. We first visualize the singular value spectrum of L for different methods and hyperparameter configurations for our method. It is clear that the iterative weight update method exhibits a more pronounced decay in singular values compared to both zero-shot SVD and the PCP baseline. (b) Energy Retention. Here we show the cumulative energy retention of L, calculated as  $\sum_{j=1}^{i} \sigma_j^2(L_k) / \sum_{j=1}^{r} \sigma_j^2(L) = \|L_i\|_F^2 / \|L\|_F^2$ . Across various target rank values, our method captures a higher proportion of energy within the top k singular values, suggesting that our low-rank approximations preserve more information from the original matrix, thereby enhancing performance. (c) Error Analysis. This subfigure presents the Frobenius norm of the error  $E_i = W - (S^{(t)} + L_i^{(t)})$ , noting that  $E_i = L - L_i$  as well. It is observed that at the desired target rank, our method achieves a lower Frobenius norm error than both zero-shot SVD and the PCP baseline. (d) Convergence Analysis. We illustrate how the approximation error between the dense matrix W and the combined approximation  $S^{(t)} + L_k^{(t)}$  diminishes over iterations, as predicted by Theorem 2 and Theorem 3. Specifically, it shows the Frobenius norm  $\|E_i^{(t)}\|_F = \|L^{(t)} - L_i^{(t)}\|_F$  at various iterations for k = 512 and T = 50. 

431 Investigating the Lower Bound of rank(L) An intriguing question arises: Is it possible to fully compress the information contained in  $W \odot (1 - P)$  into S'? In other words, does the rank of L



Figure 4: Analysis of the residual matrix L = W - S' and its low-rank approximation  $L_k$  using different methods. Results are shown for zero-shot SVD, PCP baseline (T = 100), and our proposed method with varying k (64, 128, 512) and T = 50. We show the x-axis in log scale and vertical dashed lines at i = 64, 128, 512 for better visualization. Subfigures (a), (b), and (c) have a shared legend. (a) Singular value spectrum of L. (b) Proportion of total energy captured by the top ksingular values, calculated as  $\sum_{i=1}^{k} \sigma_i^2(L_k) / \sum_{i=1}^{r} \sigma_i^2(L)$ . (c) Frobenius norm of the dropped matrix  $\|L - L_i\|_F$ . (d)  $\|L - L_i\|_F$  at different number of iterations t. Here T = 50 and k = 512 and the vertical dashed line at i = 512 indicates the target rank.



(a) Spectrum of L across various ranks of W.

457

458

459

460

461 462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477 478

(b) Spectrum of L across various sparsity levels.

Figure 5: Singular value spectrum analysis of L across different conditions. We compare our method with the zero-shot SVD method  $L = W \odot (1 - P)$ . (a) This subfigure shows the spectrum of Lacross different ranks of the original weight matrix W (64, 128, 256, 512). (b) This subfigure shows the spectrum of L across different sparsity levels (0.3, 0.5, 0.7) of magnitude pruning.

have a lower bound? To address this, we conduct an empirical investigation by applying our method to various target ranks k, using both synthetic matrices and real-world weight matrices extracted from the LLaMa-7B model. We first construct a series of synthetic matrices W of known rank r by multiplying two Gaussian random matrices:  $U \in \mathbb{R}^{m \times r}$  and  $V \in \mathbb{R}^{n \times r}$ , such that  $W = UV^T$ . Here, m = n = 512. We subsequently apply our method along with zero-shot SVD to W, varying the rank r and sparsity level p, and visualize the spectrum of L = W - S' in Figure 5.

In both subfigures in Figure 5, the solid lines represent the zero-shot SVD method ( $L = W \odot (1-P)$ ), 489 while the dashed lines represent the proposed method. The proposed method consistently shows a 490 steeper decay in singular values compared to the zero-shot SVD method, indicating better compression 491 of information. In subfigure 5a, we observe that for diverse ranks of W and a constant sparsity level 492 p = 0.5, the rank of L remains invariant. Conversely, subfigure 5b shows that fixing the rank of W 493 at r = 512 and varying the sparsity level p from 0.3 to 0.7, reveals a positive correlation between 494 the rank of L and increasing sparsity. More specifically, the rank of L is around  $p \times \min(m, n)$ . 495 Empirically, we come to the conclusion that the rank of L is independent of rank(W) but dependent 496 on the sparsity level p. This is consistent with the intuitive explanation that as the sparsity level increases, the pruned matrix contains more information, and thus the rank of L is higher. Similar 497 observations can be made in the LLaMa-7B model, as shown in Figure 2. 498

499 500

501

### 5 RELATED WORK

LLM Compression. LLMs have become essential in natural language processing tasks, but their 502 substantial size poses challenges in terms of computational resources and efficiency. Various tech-503 niques have been proposed to compress these models while maintaining their performance, including 504 pruning, low-rank compression, quantization, and knowledge distillation (Cheng et al., 2017; Choud-505 hary et al., 2020; Haroush et al., 2020). In this study, we focus on the pruning technique, and the 506 low-rank approximations are utilized to recover the lost performance. Pruning LLMs with billions of 507 parameters differs significantly from pruning smaller models (Gale et al., 2019; Frankle et al., 2020; 508 Kurtic & Alistarh, 2022), as current pruning techniques often necessitate extensive re-training after 509 the pruning process, which is prohibitively costly for LLMs (Komatsuzaki et al., 2022; Chung et al., 510 2024; Snell et al., 2024). Pruning methods are often categorized into unstructured and structured 511 pruning (Liu et al., 2017; Fan et al., 2019; He & Xiao, 2023). Structured pruning techniques, including layer pruning (Chen & Zhao, 2018; Kim et al., 2024b), channel pruning (He et al., 2017; Zhuang et al., 512 2018), and N:M pruning (Sun et al., 2021), focus on eliminating entire neurons, layers, or N out of M 513 elements in a regular pattern. Unstructured pruning removes individual weights without considering 514 the underlying structure, which can often lead to better performance but is less hardware-friendly. 515 Low-rank compression reduces the dimensions of the weight matrix by focusing on larger singular 516 values in both the column and row spaces (Cheng et al., 2005; Idelbayev & Carreira-Perpinán, 2020). 517 While quantization (Xiao et al., 2023; Lin et al., 2024) and knowledge distillation (Wan et al., 2024b) 518 have been used to compress LLMs, they are orthogonal to our approach. 519

Post-Pruning Performance Recovery. To recover the performance after pruning step, several post-pruning recovery techniques have been explored such as re-training the pruned model and using knowledge distillation (Muralidharan et al., 2024). However, these approaches can be computationally expensive and time-consuming. Recent research has focused on incorporating low-rank approximations to recover the lost performance. By adding a low-rank component to the pruned model, it is possible to approximate the original dense model more closely with minimal computational overhead (Li et al., 2023; Mozaffari et al., 2024; Zhang & Papyan, 2024).

526 527 528

### 6 CONCLUSION

In this work, we present an iterative weight update algorithm for low-rank refinement of sparse-529 pruned models. Our approach aims to bridge the performance gap between dense and pruned sparse 530 models in large language models. The proposed method offers a computationally efficient solution 531 that does not rely on extensive datasets or high-performing teacher models, making it a practical 532 choice for improving sparse model performance. A notable advantage of our method is its sparsity-533 preserving property, which allows for the concurrent update of the sparse matrix while maintaining its 534 sparsity pattern and incorporating a low-rank component. This approach effectively recovers crucial 535 information lost during pruning, leading to performance recovery, especially at high sparsity ratios. 536

These experimental results highlight the potential of combining low-rank and sparsity in LLMs. Future work could explore methods for automatically determining the optimal rank for the low-rank component based on the specific characteristics of each layer or the overall model architecture.

## 540 REFERENCES

556

563

565

574

575

576

577

| 542 | Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.   |
|-----|--|
| 543 | Emmanuel I Candès, Xiaodong Li, Yi Ma, and John Wright, Robust principal component analysis? |
| 544 | <i>Journal of the ACM (JACM)</i> 58(3):1–37, 2011  |
| 545 | <i>bournal of the right (brown), 50(5),1 57, 2011.</i>                                       |
| 546 | Venkat Chandrasekaran, Sujay Sanghavi, Pablo A Parrilo, and Alan S Willsky. Rank-sparsity    |
| 547 | incoherence for matrix decomposition. SIAM Journal on Optimization, 21(2):572-596, 2011.     |

Shi Chen and Qi Zhao. Shallowing deep networks: Layer-wise pruning based on feature representations. *IEEE transactions on pattern analysis and machine intelligence*, 41(12):3048–3056, 2018.

- Hongwei Cheng, Zydrunas Gimbutas, Per-Gunnar Martinsson, and Vladimir Rokhlin. On the compression of low rank matrices. *SIAM Journal on Scientific Computing*, 26(4):1389–1404, 2005.
- Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. A survey of model compression and acceleration
   for deep neural networks. *arXiv preprint arXiv:1710.09282*, 2017.
- Tejalal Choudhary, Vipul Mishra, Anurag Goswami, and Jagannathan Sarangapani. A comprehensive
   survey on model compression and acceleration. *Artificial Intelligence Review*, 53:5113–5155, 2020.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,
   Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language
   models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
  - Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv* preprint arXiv:1810.04805, 2018.
- Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with
   structured dropout. *arXiv preprint arXiv:1909.11556*, 2019.
- Siqi Fan, Xin Jiang, Xiang Li, Xuying Meng, Peng Han, Shuo Shang, Aixin Sun, Yequan Wang, and Zhongyuan Wang. Not all layers of llms are necessary during inference. *arXiv preprint arXiv:2403.02181*, 2024.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural
   networks. *arXiv preprint arXiv:1803.03635*, 2018.
  - Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pp. 3259–3269. PMLR, 2020.
- Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pp. 10323–10337. PMLR, 2023.
- Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv* preprint arXiv:1902.09574, 2019.
- 583 Song Han. *Efficient methods and hardware for deep learning*. PhD thesis, Stanford University, 2017.
- Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015a.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for
   efficient neural network. *Advances in neural information processing systems*, 28, 2015b.
- Matan Haroush, Itay Hubara, Elad Hoffer, and Daniel Soudry. The knowledge within: Methods for
   data-free model compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8494–8502, 2020.
- 593 Babak Hassibi, David G Stork, and Gregory J Wolff. Optimal brain surgeon and general network pruning. In *IEEE international conference on neural networks*, pp. 293–299. IEEE, 1993.

| 594<br>595<br>596        | Yang He and Lingao Xiao. Structured pruning for deep convolutional neural networks: A survey. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 2023.  |
|--------------------------|--|
| 597<br>598               | Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks.<br>In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 1389–1397, 2017.   |
| 599<br>600<br>601        | Geoffrey Hinton. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> , 2015.  |
| 602<br>603<br>604        | Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. <i>arXiv preprint arXiv:2203.15556</i> , 2022.                     |
| 605<br>606<br>607<br>608 | Yerlan Idelbayev and Miguel A Carreira-Perpinán. Low-rank compression of neural nets: Learning the rank of each layer. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 8049–8059, 2020.  |
| 609<br>610<br>611        | Bo-Kyeong Kim, Geonmin Kim, Tae-Ho Kim, Thibault Castells, Shinkook Choi, Junho Shin, and Hyoung-Kyu Song. Shortened llama: A simple depth pruning for large language models. <i>arXiv</i> preprint arXiv:2402.02834, 2024a.   |
| 612<br>613<br>614        | Jinuk Kim, Marwa El Halabi, Mingi Ji, and Hyun Oh Song. Layermerge: Neural network depth compression through layer pruning and merging. <i>arXiv preprint arXiv:2406.12837</i> , 2024b.  |
| 615<br>616<br>617        | Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. Sparse upcycling: Training mixture-of-experts from dense checkpoints. <i>arXiv preprint arXiv:2212.05055</i> , 2022.                            |
| 618<br>619<br>620        | Eldar Kurtic and Dan Alistarh. Gmp*: Well-tuned gradual magnitude pruning can outperform most bert-pruning methods. <i>arXiv preprint arXiv:2210.06384</i> , 2022.   |
| 621<br>622<br>623        | Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. Advances in neural information processing systems, 2, 1989.   |
| 624<br>625<br>626        | Yixiao Li, Yifan Yu, Qingru Zhang, Chen Liang, Pengcheng He, Weizhu Chen, and Tuo Zhao.<br>Losparse: Structured compression of large language models based on low-rank and sparse approxi-<br>mation. In <i>International Conference on Machine Learning</i> , pp. 20336–20350. PMLR, 2023.          |
| 627<br>628<br>629<br>630 | Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device Ilm compression and acceleration. <i>Proceedings of Machine Learning and Systems</i> , 6: 87–100, 2024. |
| 632<br>633<br>634        | Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learn-<br>ing efficient convolutional networks through network slimming. In <i>Proceedings of the IEEE</i><br><i>international conference on computer vision</i> , pp. 2736–2744, 2017.                         |
| 635<br>636<br>637<br>638 | Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect. <i>arXiv preprint arXiv:2403.03853</i> , 2024.  |
| 639<br>640<br>641        | Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language models: A survey. <i>arXiv preprint arXiv:2402.06196</i> , 2024.  |
| 642<br>643<br>644        | Mohammad Mozaffari, Amir Yazdanbakhsh, Zhao Zhang, and Maryam Mehri Dehnavi. Slope:<br>Double-pruned sparse plus lazy low-rank adapter pretraining of llms. <i>arXiv preprint</i><br><i>arXiv:2405.16325</i> , 2024.   |
| 646<br>647               | Saurav Muralidharan, Sharath Turuvekere Sreenivas, Raviraj Joshi, Marcin Chochowski, Mostofa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Jan Kautz, and Pavlo Molchanov. Compact language models via pruning and knowledge distillation. <i>arXiv preprint arXiv:2407.14679</i> , 2024.              |

| 648<br>649<br>650<br>651 | Peter J. Olver and Chehrzad Shakiban. Applied Linear Algebra. Undergraduate Texts in Mathematics.<br>Springer International Publishing, Cham, 2018. ISBN 978-3-319-91040-6 978-3-319-91041-<br>3. doi: 10.1007/978-3-319-91041-3. URL http://link.springer.com/10.1007/<br>978-3-319-91041-3. |  |  |  |  |  |
|--------------------------|---|--|--|--|--|--|
| 652<br>653<br>654<br>655 | Chong Peng, Yongyong Chen, Zhao Kang, Chenglizhao Chen, and Qiang Cheng. Robust principal component analysis: A factorization-based approach with linear complexity. <i>Information Sciences</i> , 513:581–599, 2020.   |  |  |  |  |  |
| 656<br>657               | Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. <i>arXiv preprint arXiv:2408.03314</i> , 2024.  |  |  |  |  |  |
| 658<br>659<br>660        | Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach large language models. <i>arXiv preprint arXiv:2306.11695</i> , 2023.  |  |  |  |  |  |
| 661<br>662<br>663        | Wei Sun, Aojun Zhou, Sander Stuijk, Rob Wijnhoven, Andrew O Nelson, Henk Corporaal, et al. Dominosearch: Find layer-wise fine-grained n: M sparse schemes from dense neural networks. <i>Advances in neural information processing systems</i> , 34:20721–20732, 2021.                        |  |  |  |  |  |
| 664<br>665<br>666        | Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.                        |  |  |  |  |  |
| 668<br>669               | Fanqi Wan, Xinting Huang, Deng Cai, Xiaojun Quan, Wei Bi, and Shuming Shi. Knowledge fusion of large language models. <i>arXiv preprint arXiv:2401.10491</i> , 2024a.   |  |  |  |  |  |
| 670<br>671<br>672        | Fanqi Wan, Ziyi Yang, Longguang Zhong, Xiaojun Quan, Xinting Huang, and Wei Bi. Fusechat:<br>Knowledge fusion of chat models. <i>arXiv preprint arXiv:2402.16107</i> , 2024b.   |  |  |  |  |  |
| 673<br>674               | Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language model pre-training via structured pruning. <i>arXiv preprint arXiv:2310.06694</i> , 2023.  |  |  |  |  |  |
| 675<br>676<br>677        | Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant:<br>Accurate and efficient post-training quantization for large language models. In <i>International</i><br><i>Conference on Machine Learning</i> , pp. 38087–38099. PMLR, 2023.                    |  |  |  |  |  |
| 679<br>680<br>681        | Lu Yin, Ajay Kumar Jaiswal, Shiwei Liu, Souvik Kundu, and Zhangyang Wang. Junk DNA hypothesis: Pruning small pre-trained weights irreversibly and monotonically impairs "difficult" downstream tasks in llms. In <i>Forty-first International Conference on Machine Learning</i> , 2024.      |  |  |  |  |  |
| 682<br>683<br>684        | Zhiyuan Zha, Bihan Wen, Jiachao Zhang, Jiantao Zhou, and Ce Zhu. A comparative study for the nuclear norms minimization methods. In <i>2019 IEEE International Conference on Image Processing (ICIP)</i> , pp. 2050–2054. IEEE, 2019.   |  |  |  |  |  |
| 685<br>686<br>687<br>688 | Stephen Zhang and Vardan Papyan. OATS: Outlier-Aware Pruning Through Sparse and Low Rank Decomposition, September 2024. URL http://arxiv.org/abs/2409.13652. arXiv:2409.13652 [cs].   |  |  |  |  |  |
| 689<br>690<br>691<br>692 | Zhuangwei Zhuang, Mingkui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and Jinhui Zhu. Discrimination-aware channel pruning for deep neural networks. <i>Advances in neural information processing systems</i> , 31, 2018.   |  |  |  |  |  |
| 693<br>694<br>695<br>696 |   |  |  |  |  |  |
| 697<br>698<br>699<br>700 |   |  |  |  |  |  |

The appendix is organized into multiple sections, with each section offering additional details.

| A | Theoretical Analysis   | 14 |
|---|------------------------|----|
|   | A.1 Additional Lemmas  | 14 |
|   | A.2 Proofs of Theorems | 15 |
|   |                        |    |

### A THEORETICAL ANALYSIS

In this section, we provide the proofs for the theorems presented in the main text and additional lemmas.

### 715 A.1 ADDITIONAL LEMMAS

In this subsection, we present additional lemmas that are useful for the proofs of the theorems in the main text.

**Lemma 1 (Mask Norm Inequality I)** For any matrix  $M \in \mathbb{R}^{m \times n}$  and binary mask  $P \in \{0,1\}^{m \times n}$ , the following inequality holds:

$$\|\boldsymbol{P} \odot \boldsymbol{M}\|_F \le \|\boldsymbol{M}\|_F \tag{11}$$

**Proof 1 (Proof of Lemma 1)** Let  $m_{ij}$  and  $p_{ij}$  be the elements of M and P respectively. By definition of the Frobenius norm and element-wise multiplication:

$$\|\boldsymbol{P} \odot \boldsymbol{M}\|_{F}^{2} = \sum_{i=1}^{m} \sum_{j=1}^{n} (p_{ij}m_{ij})^{2} = \sum_{i=1}^{m} \sum_{j=1}^{n} p_{ij}^{2}m_{ij}^{2}$$
(12)

Since P is a binary mask,  $p_{ij} \in \{0, 1\}$ , which means  $p_{ij}^2 = p_{ij}$ . Therefore:

$$\|\boldsymbol{P} \odot \boldsymbol{M}\|_{F}^{2} = \sum_{i=1}^{m} \sum_{j=1}^{n} p_{ij} m_{ij}^{2} \le \sum_{i=1}^{m} \sum_{j=1}^{n} m_{ij}^{2} = \|\boldsymbol{M}\|_{F}^{2}$$
(13)

Taking the square root of both sides preserves the inequality:

$$\|\boldsymbol{P} \odot \boldsymbol{M}\|_F \le \|\boldsymbol{M}\|_F \tag{14}$$

**Proof 2 (Proof of Lemma 2)** Let  $m_{ij}$  and  $p_{ij}$  be the elements of M and P respectively. By definition of the Frobenius inner product and element-wise multiplication:

$$\langle \mathbf{M}, \mathbf{P} \odot \mathbf{M} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n m_{ij}(p_{ij}m_{ij}) = \sum_{i=1}^m \sum_{j=1}^n p_{ij}m_{ij}^2 \le \sum_{i=1}^m \sum_{j=1}^n m_{ij}^2 = \langle \mathbf{M}, \mathbf{M} \rangle_F$$
 (15)

**Lemma 2 (Mask Norm Inequality II)** For any matrix  $M \in \mathbb{R}^{m \times n}$  and binary mask  $P \in \{0,1\}^{m \times n}$ , the following inequality holds:

$$\langle \boldsymbol{M}, \boldsymbol{P} \odot \boldsymbol{M} \rangle_F \le \langle \boldsymbol{M}, \boldsymbol{M} \rangle_F$$
 (16)

**Lemma 3 (Mask Norm Equality I)** For any matrix  $M \in \mathbb{R}^{m \times n}$  and binary mask  $P \in \{0, 1\}^{m \times n}$ , the following equality holds:

$$\langle \boldsymbol{M}, \boldsymbol{P} \odot \boldsymbol{M} \rangle_F = \| \boldsymbol{P} \odot \boldsymbol{M} \|_F^2$$
 (17)

**Proof 3 (Proof of Lemma 3)** Because P is a binary mask, we have  $p_{ij}^2 = p_{ij}$ . Therefore, we can rewrite the Frobenius inner product as:

$$\langle \boldsymbol{M}, \boldsymbol{P} \odot \boldsymbol{M} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n m_{ij}(p_{ij}m_{ij}) = \sum_{i=1}^m \sum_{j=1}^n p_{ij}^2 m_{ij}^2 = \| \boldsymbol{P} \odot \boldsymbol{M} \|_F^2$$
 (18)

**Lemma 4** (Mask Norm Equality II) For any matrix  $A, B \in \mathbb{R}^{m \times n}$  and binary mask  $P \in$  $\{0,1\}^{m \times n}$ , the following equality holds: 

$$\langle \boldsymbol{P} \odot \boldsymbol{A}, \boldsymbol{P} \odot \boldsymbol{B} \rangle_F = \langle \boldsymbol{A}, \boldsymbol{P} \odot \boldsymbol{B} \rangle_F = \langle \boldsymbol{P} \odot \boldsymbol{A}, \boldsymbol{B} \rangle_F$$
 (19)

**Proof 4 (Proof of Lemma 4)** By definitions of Frobenius inner product and element-wise multiplication, we have:

$$\langle \boldsymbol{P} \odot \boldsymbol{A}, \boldsymbol{P} \odot \boldsymbol{B} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n (p_{ij} a_{ij}) (p_{ij} b_{ij}) = \sum_{i=1}^m \sum_{j=1}^n p_{ij}^2 a_{ij} b_{ij}$$
(20)

$$\langle \mathbf{A}, \mathbf{P} \odot \mathbf{B} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n a_{ij}(p_{ij}b_{ij}) = \sum_{i=1}^m \sum_{j=1}^n p_{ij}a_{ij}b_{ij}$$
 (21)

$$\mathbf{P} \odot \mathbf{A}, \mathbf{B} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n (p_{ij} a_{ij}) b_{ij} = \sum_{i=1}^m \sum_{j=1}^n p_{ij} a_{ij} b_{ij}$$
(22)

For binary mask,  $p_{ij}^2 = p_{ij}$ , so the three expressions are equivalent. 

#### A.2 PROOFS OF THEOREMS

**Proof 5 (Proof of Theorem 1)** We prove this by induction. For t = 0,  $S^{(0)} = W \odot P$ , so the property holds. Assume the property holds for iteration t, i.e.,  $S^{(t)} = S^{(t)} \odot P$ . We need to prove that it holds for iteration t + 1, i.e.,  $S^{(t+1)} = S^{(t+1)} \odot P$ . Therefore, we have: 

$$\begin{split} \boldsymbol{S}^{(t+1)} &= \boldsymbol{S}^{(t)} + \boldsymbol{P} \odot \left\{ \boldsymbol{U}_{r^{(t)}:} \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{\Sigma}_{r^{(t)}:} \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{V}_{r^{(t)}:}^{\top} \left( \boldsymbol{L}^{(t)} \right) \right\} \\ &= (\boldsymbol{S}^{(t)} \odot \boldsymbol{P}) + \boldsymbol{P} \odot (SVD \ terms) \quad (applying \ induction \ hypothesis) \\ &= \boldsymbol{P} \odot (\boldsymbol{S}^{(t)} + SVD \ terms) = \boldsymbol{P} \odot \boldsymbol{S}^{(t+1)} \end{split}$$

Thus, we have shown that if the property holds for t, it also holds for t + 1. Combined with the base case, this completes the induction proof, showing that  $S^{(t)} = S^{(t)} \odot P$  for all t.  $\square$ 

**Proof 6 (Proof of Theorem 2)** Let  $S^*$  and  $L^*$  be the limits of  $S^{(T)}$  and  $L^{(T)}$  as  $T \to \infty$ , respec-tively. We will show that these limits exist and satisfy the stated properties. First, note that for any t,  $S^{(t)} = S^{(t)} \odot P$  by construction of the algorithm. Let  $\epsilon_t = \|S^{(t+1)} - S^{(t)}\|_F$ . We need to show that  $\lim_{T\to\infty} \epsilon_{T-1} = 0$ . At each iteration t, we have: 

$$\boldsymbol{S}^{(t+1)} = \boldsymbol{S}^{(t)} + \boldsymbol{P} \odot \left\{ \boldsymbol{U}_{r^{(t)}} \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{\Sigma}_{r^{(t)}} \left( \boldsymbol{L}^{(t)} \right) \boldsymbol{V}_{r^{(t)}}^{\top} \left( \boldsymbol{L}^{(t)} \right) \right\}$$
(23)

$$= \boldsymbol{S}^{(t)} + \boldsymbol{P} \odot \left( \boldsymbol{W} - \boldsymbol{S}^{(t)} \right)_{r^{(t)}}.$$
(24)

where  $(\cdot)_{r^{(t)}}$  denotes the truncated SVD reconstruction from  $r^{(t)}$  onwards. From the properties of SVD, we can express the difference between the truncated SVDs as:

$$\left(\boldsymbol{W} - \boldsymbol{S}^{(t)}\right)_{r^{(t)}:} - \left(\boldsymbol{W} - \boldsymbol{S}^{(t)}\right)_{k:} = \sum_{i=r^{(t)}+1}^{k} \sigma_{i}^{(t)} \boldsymbol{u}_{i}^{(t)} \boldsymbol{v}_{i}^{(t)^{\top}},$$
(25)

where  $\sigma_i^{(t)}$ ,  $u_i^{(t)}$ , and  $v_i^{(t)}$  are the *i*-th singular value and corresponding left and right singular vectors of  $W - S^{(t)}$ . The Frobenius norm of this difference is: 

$$\left\|\sum_{i=r^{(t)}+1}^{k} \sigma_i^{(t)} \boldsymbol{u}_i^{(t)} \boldsymbol{v}_i^{(t)^{\top}}\right\|_F^2 = \sum_{i=r^{(t)}+1}^{k} \sigma_i^{(t)^2}.$$
(26)

While we cannot directly guarantee that this sum approaches zero as t increases, we can bound it as: 

808  
809 
$$\sum_{i=r^{(t)}+1}^{k} \sigma_i^{(t)^2} \le \left(k - r^{(t)}\right) \left(\sigma_{r^{(t)}+1}^{(t)}\right)^2.$$
(27)

$$\left\| \boldsymbol{P} \odot \left( \left( \boldsymbol{W} - \boldsymbol{S}^{(t)} \right)_{r^{(t)}:} - \left( \boldsymbol{W} - \boldsymbol{S}^{(t)} \right)_{k:} \right) \right\|_{F} \leq \left\| \left( \boldsymbol{W} - \boldsymbol{S}^{(t)} \right)_{r^{(t)}:} - \left( \boldsymbol{W} - \boldsymbol{S}^{(t)} \right)_{k:} \right\|_{F} < \epsilon$$
(28)

Setting  $\delta = \epsilon$ , we have shown that for any  $\delta > 0$ , there exists a  $T_0$  such that for all  $t > T_0$ :

$$\left\| \boldsymbol{P} \odot \left( \boldsymbol{W} - \boldsymbol{S}^{(t)} \right)_{r^{(t)}:} - \boldsymbol{P} \odot \left( \boldsymbol{W} - \boldsymbol{S}^{(t)} \right)_{k:} \right\|_{F} < \delta$$

$$t > T_{0}:$$
(29)

This implies that for 
$$t > T_0$$
:

$$\epsilon_t = \|\boldsymbol{S}^{(t+1)} - \boldsymbol{S}^{(t)}\|_F < \delta \tag{30}$$

Since  $\delta$  can be arbitrarily small, we conclude that  $\lim_{T\to\infty} \epsilon_{T-1} = 0$ . This shows that the sequence  $\{S^{(T)}\}$  is Cauchy and therefore converges to some limit  $S^*$ . Since  $W = S^{(t)} + L^{(t)}$  for all t, and  $S^{(T)}$  converges to  $S^*$ , it follows that  $L^{(T)}$  must converge to  $L^* = W - S^*$ . Finally, since  $S^{(t)} = S^{(t)} \odot P$  for all t, we have  $S^* = S^* \odot P$ . Thus, the algorithm converges to the solution  $(S^*, L^*)$  as  $T \to \infty$ , satisfying all the stated properties.

**Proof 7 (Proof of Theorem 3)** Let  $E^{(t)} = W - (S^{(t)} + L_k^{(t)}) = (L^{(t)})_{k:}$  be the error at iteration *t*. From the update rule in Eq.(7), we have:

$$\boldsymbol{E}^{(t+1)} = \boldsymbol{W} - \left(\boldsymbol{S}^{(t+1)} + \boldsymbol{L}_{k}^{(t+1)}\right)$$
(31)

$$= \boldsymbol{W} - \left(\boldsymbol{S}^{(t)} + \boldsymbol{P} \odot \left(\boldsymbol{L}^{(t)}\right)_{r^{(t)}:} + \boldsymbol{L}_{k}^{(t+1)}\right)$$
(32)

$$= \boldsymbol{L}^{(t)} - \boldsymbol{P} \odot \left( \boldsymbol{L}^{(t)} \right)_{r^{(t)}:} - \boldsymbol{L}_{k}^{(t+1)}$$
(33)

836 Let's denote  $\mathbf{A} = \mathbf{P} \odot (\mathbf{L}^{(t)})_{r^{(t)}}$ . Then  $\mathbf{E}^{(t+1)} = (\mathbf{L}^{(t)} - \mathbf{A}) - \mathbf{L}_{k}^{(t+1)}$ . By construction,  $\mathbf{L}_{k}^{(t+1)}$  is 837 the best rank-k approximation of  $\mathbf{L}^{(t+1)} = \mathbf{L}^{(t)} - \mathbf{A}$ . Therefore:

$$\left\|\boldsymbol{E}^{(t+1)}\right\|_{F} = \left\|\left(\boldsymbol{L}^{(t)} - \boldsymbol{A}\right) - \boldsymbol{L}_{k}^{(t+1)}\right\|_{F} \le \left\|\left(\boldsymbol{L}^{(t)} - \boldsymbol{A}\right) - \boldsymbol{L}_{k}^{(t)}\right\|_{F}$$
(34)

Now,  $oldsymbol{L}^{(t)} - oldsymbol{L}^{(t)}_k = oldsymbol{E}^{(t)}$ , so we have:

$$\left\| \boldsymbol{E}^{(t+1)} \right\|_{F}^{2} \leq \left\| \left( \boldsymbol{L}^{(t)} - \boldsymbol{A} \right) - \boldsymbol{L}_{k}^{(t)} \right\|_{F}^{2} = \left\| \left( \boldsymbol{L}^{(t)} - \boldsymbol{L}_{k}^{(t)} \right) - \boldsymbol{A} \right\|_{F}^{2}$$
(35)

$$= \left\| \boldsymbol{E}^{(t)} - \boldsymbol{A} \right\|_{F}^{2} = \left\| \boldsymbol{E}^{(t)} \right\|_{F}^{2} + \left\| \boldsymbol{A} \right\|_{F}^{2} - 2\left\langle \boldsymbol{E}^{(t)}, \boldsymbol{A} \right\rangle_{F}$$
(36)

According to Lemma 2 and  $A = P \odot \left( L_{r^{(t)}:k}^{(t)} + E^{(t)} \right)$ , we have:

$$\left\langle \boldsymbol{E}^{(t)}, \boldsymbol{A} \right\rangle_{F} = \left\langle \boldsymbol{E}^{(t)}, \boldsymbol{P} \odot \left( \boldsymbol{L}_{r^{(t)}:k}^{(t)} + \boldsymbol{E}^{(t)} \right) \right\rangle_{F}$$
(37)

$$= \left\langle \boldsymbol{E}^{(t)}, \boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)} \right\rangle_{F} + \left\langle \boldsymbol{E}^{(t)}, \boldsymbol{P} \odot \boldsymbol{E}^{(t)} \right\rangle_{F}$$
(38)

$$= \left\langle \boldsymbol{E}^{(t)}, \boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)} \right\rangle_{F} + \left\| \boldsymbol{P} \odot \boldsymbol{E}^{(t)} \right\|_{F}^{2}$$
(39)

The last equality follows from Lemma 3. Substituting this back into our earlier inequality Eq.(36):

$$\left\| \boldsymbol{E}^{(t+1)} \right\|_{F}^{2} \leq \left\| \boldsymbol{E}^{(t)} \right\|_{F}^{2} + \left\| \boldsymbol{A} \right\|_{F}^{2} - 2\left\langle \boldsymbol{E}^{(t)}, \boldsymbol{A} \right\rangle_{F}$$
(40)

$$\leq \left\| \boldsymbol{E}^{(t)} \right\|_{F}^{2} + \left\| \boldsymbol{A} \right\|_{F}^{2} - 2 \left\langle \boldsymbol{E}^{(t)}, \boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)} \right\rangle_{F}^{2} - 2 \left\| \boldsymbol{P} \odot \boldsymbol{E}^{(t)} \right\|_{F}^{2}$$
(41)

859  
860
$$= \left\| \boldsymbol{E}^{(t)} \right\|_{F}^{2} + \left\| \boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)} \right\|_{F}^{2} + \left\| \boldsymbol{P} \odot \boldsymbol{E}^{(t)} \right\|_{F}^{2} + 2 \left\langle \boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)}, \boldsymbol{P} \odot \boldsymbol{E}^{(t)} \right\rangle_{F}$$
860

861  

$$-2 \left\langle \boldsymbol{E}^{(t)}, \boldsymbol{P} \odot \boldsymbol{L}^{(t)}_{r(t):k} \right\rangle_{F} - 2 \left\| \boldsymbol{P} \odot \boldsymbol{E}^{(t)} \right\|_{F}^{2} \quad (by \ Lemma. \ 4) \tag{42}$$
862

863 
$$= \left\| \boldsymbol{E}^{(t)} \right\|_{F}^{2} + \left\| \boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)} \right\|_{F}^{2} - \left\| \boldsymbol{P} \odot \boldsymbol{E}^{(t)} \right\|_{F}^{2}$$
(43)

As  $t \to T-1$ ,  $r^{(t)} \to k$ , so  $\left\| \boldsymbol{P} \odot \boldsymbol{L}_{r^{(t)}:k}^{(t)} \right\|_{F}^{2} \to 0$ . Therefore, there exists a time step  $T_{0}$  such that for all  $t > T_0$ ,  $\|\boldsymbol{E}^{(t+1)}\|_F^2 - \|\boldsymbol{E}^{(t)}\|_F^2 \le -\|\boldsymbol{P} \odot \boldsymbol{E}^{(t)}\|_F^2 \le 0.$ **Proof 8 (Proof of Theorem 4)** At iteration t, we have  $W = S^{(t)} + L^{(t)}$  by construction. Let  $L^{(t)} = U\Sigma V^{\top}$  be the SVD of  $L^{(t)}$ . Then  $L_k^{(t)} = U_{:k}\Sigma_{:k}V_{\cdot k}^{\top}$  is the best rank-k approximation of  $L^{(t)}$ . The error can be expressed as:  $\left\|\boldsymbol{W} - \left(\boldsymbol{S}^{(t)} + \boldsymbol{L}_k^{(t)}\right)\right\|_{\mathrm{F}} = \left\|\left(\boldsymbol{S}^{(t)} + \boldsymbol{L}^{(t)}\right) - \left(\boldsymbol{S}^{(t)} + \boldsymbol{L}_k^{(t)}\right)\right\|_{\mathrm{F}}$  $= \left\| \boldsymbol{L}^{(t)} - \boldsymbol{L}_k^{(t)} \right\|_F = \| \boldsymbol{U}_{k:} \boldsymbol{\Sigma}_{k:} \boldsymbol{V}_{k:}^\top \|_F$  $= \sqrt{\sum_{i=k+1}^{r} \sigma_i^2 \left( \boldsymbol{L}^{(t)} \right)} \leq \sqrt{(r-k)} \sigma_{k+1} \left( \boldsymbol{L}^{(t)} \right)$ 

(44)

(45)

(46)