S²FT: Efficient, Scalable and Generalizable LLM Fine-tuning by Structured Sparsity

Xinyu Yang 1 , Jixuan Leng 1 , Geyang Guo 2 , Jiawei Zhao 3 , Ryumei Nakada 4 , Linjun Zhang 4 , Huaxiu Yao 5 , Beidi Chen 1

¹CMU, ²Georgia Tech, ³Caltech, ⁴Rutgers, ⁵UNC-Chapel Hill xinyuya2, beidic@andrew.cmu.edu

https://infini-ai-lab.github.io/S2FT-Page

Abstract

Current PEFT methods for LLMs can achieve high quality, efficient training, or scalable serving, but not all three simultaneously. To address this limitation, we investigate sparse fine-tuning and observe a remarkable improvement in generalization ability. Utilizing this key insight, we propose a family of Structured Sparse <u>Fine-Tuning</u> (S^2FT) methods for LLMs, which concurrently achieve state-of-theart fine-tuning performance, training efficiency, and inference scalability. S²FT accomplishes this by "selecting sparsely and computing densely". Based on the coupled structures in LLMs, S²FT selects a few attention heads and channels in the MHA and FFN modules for each Transformer block, respectively. Next, it co-permutes the weight matrices on both sides of all coupled structures to connect the selected subsets in each layer into a dense submatrix. Finally, S²FT performs in-place gradient updates on all selected submatrices. Through theoretical analyses and empirical results, our method prevents forgetting while simplifying optimization, delivers SOTA performance on both commonsense and arithmetic reasoning with 4.6% and 1.3% average improvements compared to LoRA, and surpasses full FT by 11.5% when generalizing to various domains after instruction tuning. Using our partial back-propagation algorithm, S^2FT saves training memory up to $3\times$ and improves latency by 1.5-2.7× compared to full FT, while achieving an average 10% improvement over LoRA on both metrics. We further demonstrate that the weight updates in S²FT can be decoupled into adapters, enabling effective fusion, fast switch, and efficient parallelism when serving multiple fine-tuned models.

1 Introduction

Recently, Large Language Models (LLMs) have achieved significant success [16, 1, 66]. With these models being applied in diverse domains, full fine-tuning (FT) is commonly employed to enhance their downstream capabilities [56, 6, 74]. However, retraining all parameters comes with three drawbacks: (i) Full FT suffers from catastrophic forgetting, where a model forgets pre-trained knowledge while acquiring new information [44, 8]. (ii) As the model and dataset sizes grow at scale, full FT becomes increasingly computation-demanding and memory-intensive [70]. (iii) It is impractical to store and serve thousands of fine-tuned LLMs on modern GPUs if each requires full parameter storage [81, 60].

Parameter-efficient fine-tuning (PEFT) methods propose to address these bottlenecks by updating a small fraction of parameters [21]. Rather than merely reducing the number of learnable parameters, an ideal PEFT method should possess three key properties to be practically effective and efficient:

High Quality: It should exhibit both memorization and generalization capabilities, balancing the acquisition of new information from fine-tuning tasks with the retention of pre-trained knowledge. **Efficient Training**: It should minimize the memory footprint for model gradient and optimization states, and further translate such memory efficiency into less computation and fine-tuning speedup. **Scalable Serving**: It should avoid adding inference overhead when serving a single PEFT model. For multiple models, new parameters should be partially stored as adapters to save memory, and allows for effective fusion [78], fast switch [33], and efficient parallelism [60] among thousands of adapters.

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

Step 1: Select sparsely with coupled structures

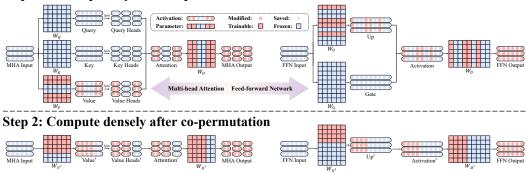


Figure 1: An Overview of the S^2FT Family for LLMs: First, we perform sparse selection of specific attention heads and channels within the coupled structures of the MHA and FFN modules. Next, we apply co-permutation to the weight matrices on both sides of these structures, enabling dense gradient computation only for the selected components. While we demonstrate S^2FT by selecting the same heads/channels on both sides for clarity, our approach also supports asymmetric selection strategies.

However, achieving all the aforementioned goals simultaneously is challenging. Common PEFT approaches, such as LoRA [27], DoRA [38], and Galore [80], project the model's weights or gradients onto a low-rank subspace. While this significantly reduces memory footprint, their performance lags behind full fine-tuning in most large-scale scenarios. Recent state-of-the-art PEFT methods have aimed to improve performance but at the cost of serving efficiency. ReFT operates on a frozen base model and learns task-specific interventions on hidden representations that cannot be merged into the original model, leading to a $2.2\times$ increase in inference latency. LISA [48] employs a coarse-grained selective method by randomly freezing most Transformer blocks during optimization, which requires significantly more trainable parameters. Consequently, in scaled serving settings like S-LoRA [60], LISA can only serve at most $\frac{1}{10}$ as many fine-tuned models as LoRA under the same memory budget.

Prior to the era of LLMs, PEFT methods based on unstructured sparse fine-tuning (SpFT) have shown a strong trade-off between low number of parameters and high model performance without sacrificing serving efficiency [63, 3, 71]. We hypothesize that SpFT, which selectively updates a small subset of model parameters, can outperform LoRA and its variants in generalization capabilities. In Figure 2, our findings across various generalization tasks support this hypothesis. However, the unstructured nature of SpFT necessitates sparse operations in computation, hindering its efficient training and scalable serving on modern hardware. This makes SpFT less practical for adapting LLMs at scale.

In this work, we propose a family of Structured Sparse Fine-Tuning (S²FT) methods to "select sparsely and compute densely" (See Figure 1), thereby closing the efficiency gap in SpFT. Inspired by structured weight pruning techniques [45, 42], we first identify several coupled structures inherent in LLMs that are connected by intermediate activations. For example, in the multi-head attention (MHA) module, each attention head in the query, key, and value projections is linked to only a few rows in the output projection. Similarly, in the feed-forward network (FFN) module, each column in the up and gate projections corresponds to a single row in the down projection. By co-permuting the matrices on both sides of these coupled structures, we can preserve the original output of these structures, with only the order of the intermediate activations changed. Exploiting this property, our S²FT strategically selects a subset of attention heads for the MHA module and a subset of channels for the FFN module. We then permute the coupled structures to connect the selected components within each linear layer into a dense submatrix. Finally, through our partial back-propagation algorithm with only two-line code modification, S²FT performs in-place gradient updates exclusively for all selected submatrices, boosting training efficiency by eliminating redundant forward activations and backward calculation.

Through our theoretical analysis, S^2FT mitigates forgetting under distribution shifts while simplifying optimization. Empirically, S^2FT outperforms other PEFT methods on LLaMA and Mistral family models, improving 1.2-4.1% on commonsense reasoning tasks and 0.6-1.9% on arithmetic reasoning ones. It also surpasses full FT by 11.5% when generalize to various domains after instruction tuning.

Finally, we conduct a comprehensive analysis to verify the training efficiency and serving scalability of S^2FT . Compared to existing PEFT methods, S^2FT not only saves $1.4\text{-}3.0\times$ memory, but also increases latency by 1.5 to $2.7\times$, making LLM fine-tuning more accessible. Additionally, S^2FT 's parameter updates can be decomposed into adapters, enabling adapter fusion with smaller performance drop than LoRA. Our method also results in more scalable and efficient adapter switch and parallelism through reduced matrix multiplications, showcasing strong potential for large-scale LLM serving scenarios.

2 Memorization or Generalization?

In this section, we evaluate the memorization and generalization capabilities of various fine-tuning methods, including full FT, LoRA, and SpFT. We hypothesize that SpFT can generalize better to downstream tasks. To support this hypothesis, we present detailed observations and analyses. Further theoretical analysis about the generalization capabilities of the S²FT family can be found in Section 4.

Hypothesis. We hypothesize that SpFT offers superior generalization than both full FT and LoRA, while maintaining comparable memorization to LoRA with the same number of trainable parameters.

Experimental Setup. We fine-tune the L1ama3-8B on the Math10K data [28] using SpFT, LoRA, and full FT. In addition to training losses, accuracies are measured on downstream tasks in LLM-Adapters, including near out-of-distribution (OOD) generalization on both easy (i.e, MultiArith, AddSub, SingleEq, MAWPS) and hard (i.e, GSM8K, AQuA, SVAMP) arithmetic reasoning tasks, and far OOD generalization on commonsense reasoning ones. For PEFT methods, we set three ratios of trainable parameters (p = 10%, 1%, 0.1%) and search for the optimal hyperparameters on the valid set. In SpFT, trainable parameters are selected randomly with given ratios. See details in Appendix C.

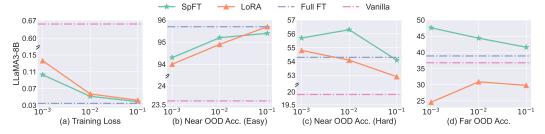


Figure 2: Accuracy comparison of SpFT, LoRA and Full FT at varying ratios of trainable parameters in various settings. SpFT exhibits strong generalization ability while full FT excels in memorization.

Observations. Figure 2 indicates several key findings. First, SpFT achieves lower training losses than LoRA when using the same ratio of trainable parameters, especially at very small ratios. This gap arises from the more complex optimization process in LoRA, which requires the simultaneous updating of two matrices [23]. Second, we observe both elevated training loss and reduced average accuracy on easier math tasks as the ratio decreases, suggesting a positive correlation between memorization abilities and trainable parameters. Notably, with only 10% of the parameters updated, PEFT methods learn comparable memorization abilities to full FT when trained on a 10k-sample dataset.

When generalizing to complex mathematical problems or commonsense reasoning tasks, the performance ranking emerges as: SpFT > Full FT > LoRA. SpFT effectively transfers reasoning abilities to commonsense domains, while LoRA exhibits significant performance drops in far OOD generalization. This indicates (i) freezing a larger fraction of the parameters can retain more pre-trained abilities, and (ii) approximating high-dimensional gradients with low-rank decomposition may overfit fine-tuned data and hinder the model from generalization. Since LLMs are pre-trained on high-quality data, SpFT emerges as the preferred choice for fine-tuning on task-specific data of varying quality.

3 The S^2FT family of methods

While SpFT demonstrates strong generalization ability and good overall performance in Section 2, its unstructured nature poses challenges for efficient training and scalable serving on modern hardware (e.g., GPU). This is because of the need for sparse operations when storing and computing weights, gradients, and optimization states, which are significantly slower than their dense variants on GPU. This motivates our investigation into structured sparsity approaches that utilize only dense operations: Can structured sparsity improve hardware efficiency while preserving performance by selecting sparsely but computing densely? If so, how far can the flexibility of selection be pushed in this context? To answer this question, we design a family of Structured Sparse Fine-Tuning (S²FT) methods with dense-only computations, making PEFT effective, efficient and scalable. We begin by discovering the coupled structure in LLMs in Section 3.1. Leveraging this property, Section 3.2 introduce the selection and permutation strategies of S²FT, with overall pipeline illustrated in Figure 1b. In Section 3.3, we present our partial back-propagation algorithm that enables end-to-end training latency reduction.

3.1 Discover Coupled Structures in LLMs

We initiate our pursuit of flexible structured sparsity by examining the coupled structures in LLMs.

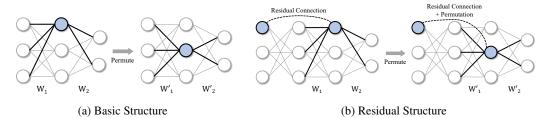


Figure 3: Grouped model weights with basic structure and residual structure. All highlighted weights must be permuted simultaneously. Residual structures require additional permutation during runtime.

Structure Dependency in LLMs. Inspired by prior work on structured pruning [45, 17], our study start by building the dependencies between activations and weights for LLMs. Let A denote an activation and W denote a weight in the model. We define $\operatorname{In}(A)$ as the set of parameters that directly contribute to the computation of A, and $\operatorname{Out}(A)$ as the set of parameters that depend on A in the computation of subsequent activations. The dependency between structures can be defined as follows:

$$W_1 \in \operatorname{In}(A) \wedge \operatorname{Deg}^+(W_1) = 1 \Rightarrow A \text{ is dependent on } W_1$$
 (1)

$$W_2 \in \operatorname{Out}(A) \wedge \operatorname{Deg}^-(W_2) = 1 \Rightarrow W_2 \text{ is dependent on } A$$
 (2)

where $\mathrm{Deg}^+(W_1)$ represents the out-degree of weight W_1 , and $\mathrm{Deg}^-(W_2)$ represents the in-degree of weight W_2 . Each equation represents a unquie directional dependency between activations and weights. When both equations hold simultaneously, a coupled structure exists between W_1 and W_2 . In Figure 3, we employ deep linear networks to illustrate two types of coupled structures in LLMs:

Basic Structures: In Figure 3a, these structures exist in both the multi-head attention (MHA) and feed-forward network (FFN) modules. Taking LLaMA as an example, in the MHA module, we consider the Query (\mathbf{Q}), Key (\mathbf{K}), and Value (\mathbf{V}) projections as W_1 , and the Output (\mathbf{O}) projection as W_2 , while Softmax($\mathbf{Q}\mathbf{K}^{\top}$) $\mathbf{V}(x)$ acting as the activation between weight matrices. Similarly, in the FFN module, the Up (\mathbf{U}) and Gate (\mathbf{G}) projections function as W_1 , with the Down (\mathbf{D}) projection corresponding to W_2 . Here, $\mathbf{U}(x)$ · SwiGLU($\mathbf{G}(x)$) serves as the activations connecting W_1 and W_2 .

Residual Structures: In Figure 3b, this type of coupled structures exists between the MHA and FFN modules. We further consider how residual connections influence the activations in these structures.

Permutation Invariance of Coupled Structures. Figure 3 demonstrates that W_1 and W_2 can be co-permuted using the same order, which only affects the order of activations between them while preserving the original output from the coupled structure. Since residual dependencies require an additional run-time step to permute the residuals, we will focus on basic dependencies in our method.

3.2 Sparse Selection and Permutation

At this point, all coupled structures within the model have been identified. The subsequent sparse selection and permutation processes are straightforward, with overall pipeline illustrated in Figure 1b.

MHA Module: There are four linear layers in a MHA module: $Q, K, V, O \in \mathbb{R}^{d \times d}$. For a model with h attention heads, each head $i \in [h]$ has its own projections denoted as $Q_i \in \mathbb{R}^{d \times d_h}$, $K_i \in \mathbb{R}^{d \times d_h}$, $V_i \in \mathbb{R}^{d \times d_h}$, and $V_i \in \mathbb{R}^{d \times d_h}$, where $V_i \in \mathbb{R}^{d \times d_h}$, to the beginning of each weight matrix, we are able to update these selected heads using dense-only operations, while keeping the other ones frozen.

FFN Module: There are three linear layers in an FFN module: $U, G \in \mathbb{R}^{k \times d}$ and $D \in \mathbb{R}^{d \times k}$. In S²FT, only a few channels require gradient updates. Let $S_{\text{FFN}} \subseteq [d]$ denote the selected channels. We can permute S_{FFN} to the beginning of each weight matrix and only fine-tune this compact subset.

Next, we provide several strategies for identifying and selecting important subsets in each module.

- 1. S²FT-R (S²FT): In this strategy, a subset of channels is randomly selected and set to be trainable.
- 2. S²FT-W: This variant selects subsets based on the magnitude of the weights for linear layers.
- 3. S²FT-A: This variant selects subsets based on the magnitude of activations on a calibration set.
- 4. S²FT-S: Top-K subsets are ranked and selected by the product of weight and activation magnitudes.
- 5. S²FT-G: This variant selects subsets based on the magnitude of gradients on a calibration set.

Here, 1 and 2 can be applied directly without pre-processing. 3 and 4 only require a forward pass on a small calibration dataset. While 5 necessitates a backward pass on this dataset, it does not store optimization states and can mitigate memory footprints for activations through gradient checkpointing [18]. By default, we use S²FT-R for a fair comparison and discuss other variants in Section 5.4.

3.3 Partial Back-propagation Algorithm

Finally, we introduce our partial back-propagation algorithm with only two line modifications in PyTorch, our algorithm stores trainable channels based on their start and end positions, thereby improving training efficiency by eliminating redundant forward activations and backward calculations.

```
def setup_context(ctx, inputs, output):
  activation, weight, bias, start, end = inputs
  # only save partial input tensors for gradient calculation in forward
  ctx.save_for_backward(activation[:, start:end], weight, bias, start, end)
def gradient_update(parameter, gradient, start, end):
    # only modify the assigned positions of weight matrices during optimization
   parameter[:, start:end].add_(gradient)
```

Theoretical Analysis

In this section, we theoretically explain why S²FT demonstrates stronger generalization capabilities compared to LoRA. Following previous work [23, 79, 53, 52], we further show that S²FT is simple and efficient in optimization by maintaining stability in both the magnitude and direction of updates.

4.1 Stronger Generalization Capability

First, we theoretically explore why S²FT demonstrates stronger generalization capabilities compared to LoRA. We consider a pre-trained L-layer deep linear network, which has been widely used to facilitate the theoretical analysis of complex DNNs [59, 30, 43, 22, 34, 5]. Let $f^{\text{pre}}(x) := W_L^{\text{pre}}W_{L-1}^{\text{pre}}\dots W_1^{\text{pre}}x$ be the pre-trained deep linear network, where $W_\ell^{\text{pre}} \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$, with $d_0 = p$ and $d_L = q$. We fine-tune the ℓ -th layer with low-rankness level $r \leq \min\{d_\ell, d_{\ell-1}\}$ or sparsity level $s = \lfloor r \cdot \frac{d_\ell + d_{\ell-1}}{d_{\ell-1}} \rfloor$. Denote a class of adaptation with parameters $U \in \mathbb{R}^{d_\ell \times d}$ and $V \in \mathbb{R}^{d_{\ell-1} \times d}$ as

$$f_{\ell,U,V}(x) := \overline{W}_{\ell+1}^{\text{pre}} (W_{\ell}^{\text{pre}} + UV^{\top}) \underline{W}_{\ell-1}^{\text{pre}} x, \tag{3}$$

 $f_{\ell,U,V}(x) := \overline{W}_{\ell+1}^{\text{pre}}(W_{\ell}^{\text{pre}} + UV^{\top})\underline{W}_{\ell-1}^{\text{pre}}x, \tag{3}$ where $\overline{W}_{\ell}^{\text{pre}} := W_{L}^{\text{pre}}W_{L-1}^{\text{pre}}\dots W_{\ell}^{\text{pre}} \in \mathbb{R}^{d_{L}\times d_{\ell-1}}$ and $\underline{W}_{\ell}^{\text{pre}} := W_{\ell}^{\text{pre}}W_{\ell-1}^{\text{pre}}\dots W_{1}^{\text{pre}} \in \mathbb{R}^{d_{\ell}\times d_{0}}$ with $\underline{W}_{0}^{\text{pre}} = I_{p}$ and $\overline{W}_{L}^{\text{pre}} = I_{q}$. In a transformer-based LLM, each row of W_{ℓ} can represent the parameters in a single attention head for the MHA module or in a single channel for the FFN module.

Given n observations $(x_i^{(i)}, y_i^{(i)}) \subset \mathbb{R}^p \times \mathbb{R}^q$, we fine-tune f^{pre} by minimizing the empirical risk $\mathcal{R}_n^{(i)}(f_{\ell,U,V}) := (1/n) \sum_{i \in [n]} \|y_i^{(i)} - f_{\ell,U,V}(x_i^{(i)})\|^2$ via gradient descent. For LoRA, we train both low-rank matrices (U,V) in Equation (3) with $d \leftarrow r$. For S²FT, we train only V in Equation (3) with $d \leftarrow s$ and fixed $U \leftarrow U_S^{\bar{S}^2FT} := [e_{a_1}; e_{a_2}; \dots; e_{a_s}]$, where $S = \{a_1, \dots, a_s\} \subset [d_\ell]$ and e_a is the a-th standard basis. Similar conclusions hold when we fine-tune only U. Motivated by the implicit regularization in gradient descent [77, 19, 5], we directly consider minimum norm solutions.

We consider a multiple linear regression setting. Assume that the in-distribution training data $(x^{(i)},$ $y^{(i)} \in \mathbb{R}^{p+q}$ and out-of-distribution test data $(x^{(o)}, y^{(o)}) \in \mathbb{R}^{p+q}$ are generated i.i.d. according to

$$y^{(k)} = B^{(k)}x^{(k)} + \epsilon^{(k)}, \ k \in \{i, o\},\$$

where $B^{(k)} \in \mathbb{R}^{q \times p}$ is the coefficient matrix, $x^{(k)}$ and $\epsilon^{(k)}$ are mean zero sub-Gaussian signal and noise with covariance matrices $\Sigma^{(k)}_x$ and $\Sigma^{(k)}_\epsilon$, respectively. The generalization capacity is measured by the fine-tuned model's excess risk $\mathcal{E}(f) := \mathbb{E}[\|y^{(o)} - f(x^{(o)})\|^2] - \inf_{f'} \mathbb{E}[\|y^{(o)} - f'(x^{(o)})\|^2].$

For these OOD data, LoRA suffers from forgetting, while S²FT can maintain pre-training knowledge.

Assumption 4.1 (Distribution Shift). Assume that
$$\Sigma_x^{(i)} = \Sigma_x^{(o)} = \Sigma_x$$
 for some $\Sigma_x \in \mathbb{R}^{p \times p}$, and $\|(\overline{W}_{\ell+1}^{\mathrm{pre}} U_S^{\mathrm{S}^2\mathrm{FT}})(\overline{W}_{\ell+1}^{\mathrm{pre}} U_S^{\mathrm{S}^2\mathrm{FT}})^{\dagger} (B^{(o)} - B^{(i)}) \Sigma_x^{1/2}\|_{\mathrm{F}}^2 \leq \varepsilon^2 \mathcal{E}^{(o)}(f^{\mathrm{pre}})$ for some $\varepsilon > 0$.

Assumption 4.1 states that while the covariate distribution remains unchanged, the label distribution conditioned on covariates may shift, but not exceeding a factor of ϵ^2 of the OOD risk of f^{pre} . This holds for fine-tuning with proper channel selection, where primarily the output distribution is changed.

Theorem 4.2 (Out-of-distribution Excess Risk, Informal). Suppose Assumption 4.1 holds. Consider $n \to \infty$. If $B^{(i)} = \overline{W}_{\ell+1}^{\mathrm{pre}} \tilde{B}^{(i)} \underline{W}_{\ell-1}^{\mathrm{pre}}$ holds for some $\tilde{B}^{(i)} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$, and $s \leq \mathrm{rank}(\Sigma_f^{(i)})$, then,

$$\mathcal{E}^{(\mathrm{o})}(f_{\ell,U_{s}^{\S^{2}\mathrm{FT}},V^{\S^{2}\mathrm{FT}}}) \leq (1+3\varepsilon^{2})\mathcal{E}^{(\mathrm{o})}(f^{\mathrm{pre}}), \ \mathcal{E}^{(\mathrm{o})}(f_{\ell,U^{\mathrm{LORA}},V^{\mathrm{LORA}}}) \geq \|(B^{(\mathrm{o})}-B^{(\mathrm{i})})\Sigma_{x}^{1/2}\|_{\mathrm{F}}^{2}.$$

Theorem 4.2 indicates that the OOD risk of S^2FT is bounded above by that of f^{pre} , while that of LoRA is bounded below by the label shift magnitude. If f^{pre} already has a low risk for OOD tasks, and the label shift is significant, S²FT is expected to outperform LoRA. Essentially, when the OOD task deviates significantly from the FT distribution, LoRA may forget pre-trained knowledge and overfit to the FT data, compromising its generalization capabilities. See formal statements in Theorem F.8.

4.2 Simple and Efficient Optimization

Next, we explain why S²FT is a simple and efficient optimization method. In Equation (3), S²FT can be viewed as a LoRA variant that fixes $U_S^{S^2FT}$ as a combination of multiple orthogonal standard basis vectors while optimizing V^{S^2FT} with zero initialization. The gradient is given by $\frac{\partial \mathcal{L}}{\partial V^{S^2FT}} = (\underline{W}_{\ell-1}^{\text{pre}}x)^{\top}\frac{\partial \mathcal{L}}{\partial \overline{W}_{\ell+1}^{\text{pre}}}U_S^{S^2FT}$. Ignore $\underline{W}_{\ell-1}^{\text{pre}}$, $\overline{W}_{\ell-1}^{\text{pre}}$ and denote $\frac{\partial \mathcal{L}}{\partial \overline{W}_{\ell+1}^{\text{pre}}}$ as \overline{G} , at step t with learning rate η , $\Delta f_{\ell,t}(x) := f_{\ell,t}(x) - f_{\ell,t-1}(x) = U_S^{S^2FT}(V_t^{S^2FT} - V_{t-1}^{S^2FT})^{\top}x = -\eta U_S^{S^2FT}U_S^{S^2FT}^{\top}\overline{G}^{\top}||x||^2$.

$$\Delta f_{\ell,t}(x) := f_{\ell,t}(x) - f_{\ell,t-1}(x) = U_S^{\mathsf{S}^2\mathsf{FT}} (V_t^{\mathsf{S}^2\mathsf{FT}} - V_{t-1}^{\mathsf{S}^2\mathsf{FT}})^\top x = -\eta U_S^{\mathsf{S}^2\mathsf{FT}} U_S^{\mathsf{S}^2\mathsf{FT}^\top} \overline{G}^\top ||x||^2.$$

Since $U_S^{8^2\mathrm{FT}}$ is an orthogonal matrix, the update simplifies to $\Delta f_{\ell,t}(x) = -\eta \overline{G}^\top ||x||^2$. Following LoRA+ [23], assuming that $x = \Theta_n(1)$, where n is the width of the layers in LLMs, we expect $\Delta f_{\ell,t}(x) = \Theta(1)$ to ensure stability and feature learning in the infinite-width limit [72]. S²FT can achieve this when $\eta = \Theta(n^{-1})$ while LoRA requires $\eta_U = \Theta(1)$ and $\eta_V = \Theta(n^{-1})$ for optimal performance. These rates become impractical for modern LLMs with very large n. Therefore, S^2FT aligns with LoRA variants that fix one matrix [52, 79], offering more stable and efficient optimization.

Furthermore, under a given sparsity level as regularization, our model simplifies optimization when approximating the full fine-tuning gradients at non-zero positions. Similar to LoRA-SB [53], let G_V denote the gradient of V^{S^2FT} . The equivalent gradient \tilde{G} , which describes the virtual gradient of the pretrained weight matrices, can be expressed as $U_S^{S^2FT}G_V^{\top}$. Then, the gradient with respect to $V_S^{S^2FT}$ can be expressed in terms of the gradient of the pretrained weight W^{pre} as: $G_V^O = U_S^{\text{S}^2 \text{FT}^\top} G$. Using this relationship, our objective is to minimize the distance between the equivalent gradient and the full gradient as $\min_{G_V} \|\tilde{G} - G\|_F^2$, where the optimal solution is given by $G_V = (U_S^{S^2FT^\top} U_S^{S^2FT})^{-1} G_V^O$. Since $U_S^{S^2FT}$ is orthogonal, we have $G_V = G_V^O$. This shows that S^2FT can keep the optimal update directions throughout the training process, establishing it as an efficient sparse optimization method.

Experiments

In this section, we conduct a series of experiments across three diverse benchmarks covering more than 20 datasets. Our goal is to provide a rich picture of how S²FT performs in different scenarios. Here, we compare our method with different fine-tuning strategies and categories including: (i) Full fine-tuning (FT), (ii) reparameterized fine-tuning: LoRA [27], DoRA [38], and Galore [80], (iii) adapter-based fine-tuning: Series Adapter [26], Parallel Adapter [24], and LoReFT [69], (iv) promptbased fine-tuning: Prefix-Tuning [36], (v) sparse fine-tuning: LISA [48]. For a fair comparison, we keep a comparable number of trainable parameters in S²FT to that of LoRA. The design choices for trainable parameter allocations in S²FT will be detailed in Section 5.4. All other hyperparameters are selected via cross-validation. Detailed setups and dataset descriptions are provided in Appendix E.

5.1 Commonsense Reasoning

The results of eight common sense reasoning tasks in Table 1 show that S²FT consistently outperforms existing PEFT methods in the LLaMA-7B / 13B, LLaMA2-7B and LLaMA3-8B models. Compared to LoRA and DoRA, it achieves average performance gains of 4.6% and 2.8%, respectively. Furthermore, S²FT also shows superior performance against recent approaches, including Galore, LoReFT, and LISA, with improvements of at least 1.0%. Remarkably, despite using less than 1% of trainable parameters, our method surpasses full FT by 0.5%. The 3.0% improvement on the LLaMA3-8B suggests that keeping most pre-trained parameters frozen enables better generalization to test distributions.

5.2 Arithmetic Reasoning

As showcased in Table 2, S²FT consistently outperforms other PEFT methods for different base models. On average, it achieves improvements of 1.3% and 0.9% over LoRA and DoRA, respectively. These results highlight the versatility and effectiveness of our approach across a diverse range of tasks. Additionally, we observe substantial improvements even when compared to Full FT for the LLaMA3-8B model, particularly on complex tasks such as GSM8K and AQuA. This suggests that S²FT better preserves the original reasoning capabilities of this stronger model while acquiring new skills from the fine-tuning data, thereby validating the enhanced generalization ability of our method.

Table 1: Comparison among various fine-tuning methods for the LLaMA-7B/13B, LLaMA2-7B, and LLaMA3-8B models on eight commonsense reasoning tasks. Non-PEFT methods are marked in gray. (1: from DoRA paper, 2: from ReFT paper, 3: reproduced by us, †: projected trainable parameters)

Model	Method	# Param(%)	BoolQ	PIQA	SIQA	HellaSwag	Wino	ARC-e	ARC-c	OBQA	Avg. ↑
ChatGPT ¹	-	-	73.1	85.4	68.5	78.5	66.1	89.8	79.9	74.8	77.0
	Full FT ³	100	70.3	84.2	80.1	92.3	85.4	86.6	72.8	83.4	81.9
	Prefix [36] ¹	0.11	64.3	76.8	73.9	42.1	72.1	72.9	54.0	60.6	64.6
	Series [26] ¹	0.99	63.0	79.2	76.3	67.9	75.7	74.5	57.1	72.4	70.8
	Parallel [24] ¹	3.54	67.9	76.4	78.8	69.8	78.9	73.7	57.3	75.2	72.2
LLaMA-7B	LoRA [27] ³	0.83	69.2	81.7	78.4	83.4	80.8	79.0	62.4	78.4	76.7
	DoRA [38] ¹	0.84	68.5	82.9	79.6	84.8	80.8	81.4	65.8	81.0	78.1
	Galore [80] ³	0.83^{\dagger}	68.6	79.0	78.5	84.7	80.1	80.3	62.1	77.3	76.3
	LoReFT [69] ²	0.03	69.3	84.4	80.3	93.1	84.2	83.2	68.2	78.9	80.2
	LISA [48] ³	9.91	70.4	82.1	78.7	92.4	82.9	84.9	70.2	78.4	80.0
	S ² FT (Ours)	0.81	72.7	83.7	79.6	93.4	83.5	86.1	72.2	83.4	81.8
	Full FT ³	100	74.5	86.3	81.3	94.4	86.9	89.7	77.9	88.8	85.0
	Prefix [36] ¹	0.03	65.3	75.4	72.1	55.2	68.6	79.5	62.9	68.0	68.4
	Series [26] ¹	0.80	71.8	83.0	79.2	88.1	82.4	82.5	67.3	81.8	79.5
LLaMA-13B	Parallel [24] ¹	2.89	72.5	84.9	79.8	92.1	84.7	84.2	71.2	82.4	81.4
	LoRA [27] ¹	0.67	72.1	83.5	80.5	90.5	83.7	82.8	68.3	82.4	80.5
	DoRA [38] ¹	0.68	72.4	84.9	81.5	92.4	84.2	84.2	69.6	82.8	81.5
	LoReFT [69] ²	0.03	72.1	86.3	81.8	95.1	87.2	86.2	73.7	84.2	83.3
	S ² FT (Ours)	0.65	74.2	85.7	80.7	94.9	86.4	88.4	76.3	87.8	84.3
	Full FT ³	100	74.7	84.9	78.7	93.7	84.1	87.5	75.2	85.0	83.0
LLaMA2-7B	LoRA [27] ¹	0.83	69.8	79.9	79.5	83.6	82.6	79.8	64.7	81.0	77.6
LLaWIA2-7B	DoRA [38] ¹	0.84	71.8	83.7	76.0	89.1	82.6	83.7	68.2	82.4	79.7
	S ² FT (Ours)	0.81	72.9	86.1	80.2	94.3	85.5	87.2	74.6	83.4	83.0
	Full FT ³	100	73.9	86.2	79.1	93.1	85.8	88.1	78.2	84.0	83.6
LLaMA3-8B	LoRA [27] ¹	0.70	70.8	85.2	79.7	92.5	84.9	88.9	78.7	84.4	82.5
LLawiA3-6D	DoRA [38] ¹	0.71	74.6	89.3	79.9	95.5	85.6	90.5	80.4	85.8	85.2
	S ² FT (Ours)	0.70	75.0	89.0	80.7	96.5	88.0	92.5	83.4	87.8	86.6

Table 2: Comparison among various fine-tuning methods for different models on seven math reasoning tasks. Non-PEFT methods are marked in gray. (1: from LLM-Adapters paper, 2: reproduced by us)

Model	Method	$\#\operatorname{Param}(\%)$	MultiArith	GSM8K	AddSub	AQuA	SingleEq	SVAMP	MAWPS	Avg. \uparrow
GPT-3.5 ¹	-	-	83.8	56.4	85.3	38.9	88.1	69.9	87.4	72.8
	Full FT ²	100	98.8	43.1	91.1	20.9	94.3	60.6	88.2	71.0
LLaMA-7B	LoRA [27] ²	0.83	98.0	40.0	91.2	21.7	93.1	56.7	85.3	69.7
EEuivii 7 / B	DoRA [38] ²	0.84	97.3	38.9	89.6	22.4	93.9	58.4	85.3	69.4
	S ² FT (Ours)	0.81	98.8	41.3	91.4	21.3	93.5	58.4	86.1	70.1
LLaMA-13B	Full FT ²	100	98.3	47.6	92.9	26.0	95.1	65.7	88.7	73.5
	LoRA [27] ²	0.67	97.5	47.8	89.9	20.5	94.3	61.2	87.4	71.2
	$DoRA [38]^2$	0.68	97.2	48.1	90.6	20.9	93.9	63.8	88.2	71.8
	S ² FT (Ours)	0.65	97.7	48.4	90.4	22.8	95.5	63.9	87.8	72.4
	Full FT ²	100	99.3	47.5	91.1	24.4	96.7	62.5	89.1	72.9
LLaMA2-7B	LoRA [27] ²	0.83	97.5	44.0	91.2	20.9	94.1	59.2	85.7	70.4
LLaWIA2-7D	DoRA [38] ²	0.84	98.2	43.8	90.1	24.4	94.5	59.1	89.1	71.3
	S ² FT (Ours)	0.81	98.5	44.3	91.1	25.2	94.7	61.8	88.2	72.0
LLaMA3-8B	Full FT ²	100	99.2	62.0	93.9	26.8	96.7	74.0	91.2	77.7
	LoRA [27] ²	0.70	99.5	61.6	92.7	25.6	96.3	73.8	90.8	77.2
	$DoRA [38]^2$	0.71	98.8	62.7	92.2	26.8	96.9	74.0	91.2	77.5
	S ² FT (Ours)	0.70	99.7	65.8	93.7	31.5	97.8	76.0	92.4	79.6
			<u> </u>							

5.3 Instruction Following

Table 3 comprehensively compares various methods on eight tasks in the MT-Bench dataset [82]. It is observed that $S^2FT > LISA > Full FT > LoRA/Galore \ge Vanilla for both the Mistral-7B and LLama2-7B model. This is because sparse FT methods like <math>S^2FT$ and LISA retain more pre-trained knowledge while acquiring new skills on the FT dataset, thereby generalizing better to diverse tasks in the MT-Bench dataset. Moreover, our method outperforms LISA due to its fine-grained and flexible selection strategy, enabling all layers to learn to follow instructions on the full fine-tuning set.

Table 3: Performance comparison of LLM fine-tuning methods trained on the Alpaca GPT-4 dataset. We report the MT-Bench score as the evaluation metric. All baseline results are cited from LISA.

Model	Method	Writing	Roleplay	Reasoning	Code	Math	Extraction	STEM	Humanities	Avg.
	Vanilla	5.25	3.20	4.50	1.60	2.70	6.50	6.17	4.65	4.32
	Full FT	5.50	4.45	5.45	2.50	3.25	5.78	4.75	5.45	4.64
Mistral-7B	LoRA	5.30	4.40	4.65	2.35	3.30	5.50	5.55	4.30	4.41
Misuai-/D	Galore	5.05	5.27	4.45	1.70	2.50	5.21	5.52	5.20	4.36
	LISA	6.84	3.65	5.45	2.20	2.75	5.65	5.95	6.35	4.85
	Ours	6.95	4.40	5.50	2.70	3.55	5.95	6.35	6.75	5.27
	Vanilla	2.75	4.40	2.80	1.55	1.80	3.20	5.25	4.60	3.29
	Full FT	5.55	6.45	3.60	1.75	2.00	4.70	6.45	7.50	4.75
II aMA 2.7D	LoRA	6.30	5.65	4.05	1.60	1.45	4.17	6.20	6.20	4.45
LLaMA2-7B	Galore	5.60	6.40	3.20	1.25	1.95	5.05	6.57	7.00	4.63
	LISA	6.55	6.90	3.45	1.60	2.16	4.50	6.75	7.65	4.94
	Ours	6.75	6.60	4.15	1.65	1.85	4.75	7.45	8.38	5.20

5.4 Design Choices for Trainable Parameter Allocations

Finally, we detail how S²FT distribute trainable parameters across layers, modules, and channels.

Uniform across Layers: Following Chen et al. [10], we allocate parameters to each layer uniformly. **Fine-tune Important Modules**: Figure 4 analyzes the effectiveness of different components in a LLaMA-like Transformer Block for fine-tuning, including Query, Key, Value, Output, Up, Gate, and Down projections. To ensure a fair comparison, we maintain a fixed number of trainable parameters when fine-tuning each component. The results show that the effectiveness of components in fine-tuning follows the order: Query/Key & Value/Up/Gate < Output/Down. This is because Query/Key are only used to measure token similarities, while others serve as persistent memories of training data. Based on this finding, we allocate our parameter budget fairly to the Output and Down projections. For the LLama3-8B and Mistral-7B models, we only fine-tune the Down projection due to the inflexible selection in multi-query attention. Further analysis of this setting is left for future research.



Figure 4: The impact of different components in fine-tuning, including Query, Key, Value, Output, Up, Gate, and Down projection. We fix the trainable parameter budget and only fine-tune one component.

Table 4: Comparison of various channel selection strategies on the commonsense and arithmetic reasoning datasets for the LLama3-8B. We report the average accuracy (%) as the evaluation metric.

Task	S ² FT-R	S^2FT-W		S^2FT-A		S^2FT-S		S ² FT-G	
	SIII	Large	Small	Large	Small	Large	Small	Large	Small
Commonsense Arithmetic	86.6 79.6	× /			87.3 _(+0.7) 80.0 _(+0.4)				86.2 _(-0.4) 79.5 _(-0.1)

Selection across Channels: In Section 3.2, we discuss several strategies for channel selection. In our main experiments, we employ random selection to ensure fair comparisons with baseline methods, as these approaches treat all channels with equal importance. However, the sparse structure of S²FT offers controllability during fine-tuning, allowing us to prioritize important channels in the selection process to further boost performance. Table 4 compared nine different strategies, incorporating five varying selection metrics (i.e., random, weight, activation, weight-activation product, and gradient), each choosing either the largest or smallest values. For S²FT-A, S²FT-S, and S²FT-G, we employ 1% of the fine-tuning data as a calibration set, introducing only negligible overhead during inference.

Our results demonstrate that random selection serves as a strong baseline due to its unbiased nature. Among heuristic metrics, selecting channels with the smallest activations (i.e., S²FT-A and S²FT-S) outperforms random selection. This indicates that these channels contain less task-specific information, enabling us to inject new knowledge through fine-tuning while preserving pre-trained capabilities in other channels. In contrast, other strategies introduce bias that compromises model performance. Notably, the counterintuitive accuracy decrease in S²FT-G (Large) suggests that channels with large gradients contain task-related pre-trained knowledge, and modifying them will disrupt these abilities.

6 Analysis

Having demonstrated the strong generalization capability and overall performance of S²FT, we now further explore its training efficiency and serving scalability compared to other fine-tuning techniques.

6.1 Training Efficiency

To evaluate training efficiency, we examine two crucial metrics: peak memory footprint and average training latency. These numbers are measured on a single Nvidia A100 (80G) SXM GPU. We keep a comparable number of parameters for all methods. To obtain the average latency, we fine-tune the model for 50 runs, each run including 200 iterations, with 10 warmup runs excluded in measurement.

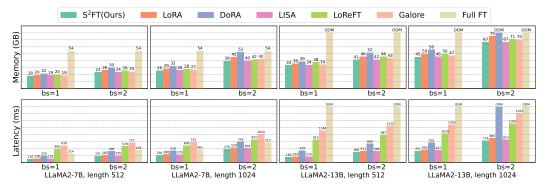


Figure 5: Comparison of memory and computation efficiency during training on the LLaMA2-7B/13B with varying sequence lengths and batch sizes. Average latency and peak memory usage are reported. S²FT significantly improves training latency while reducing memory footprint compared to baselines.

In Figure 5, we thoughtfully profile S^2FT on various model sizes, sequence lengths, and batch sizes. Compared to Full FT, S^2FT saves $1.4\text{-}3.0\times$ memory, and speedups fine-tuning by 1.5-2.7 times. When benchmarking against other PEFT methods, S^2FT establishes new standards for efficiency, offering average reductions of 2% in memory usage and 9% in latency. Notably, S^2FT outperforms the widely adopted LoRA, achieving about 10% improvement in both metrics by avoiding the need to store new parameters and perform additional calculations. Our partial back-propagation algorithm further improves efficiency by saving unnecessary forward activations and backward calculations.

6.2 Serving Scalability

While S²FT avoids additional inference overhead for a single fine-tuned model through in-place gradient updates, we will now discuss its scalability for serving thousands of fine-tuned models. To begin, we introduce the unmerged computation paradigm of S²FT: Given a pre-trained weight matrix $W^{pre} \in \mathbb{R}^{d \times k}$ and its corresponding fine-tuned weight matrix W with sparsity level s, we define the weight difference as $\Delta W = W - W^{pre}$. Similar to Section 4, ΔW can be decomposed into the product of a weight matrix $V \in \mathbb{R}^{k \times s}$ and a permutation matrix $U \in \mathbb{R}^{d \times s}$. This decomposition allows us to "unmerge" an adapter $\Delta W = UV^{\top}$ from W, thereby sharing similarities with other adapters during inference. Following Zhong et al. [83], we consider three different adapter composition scenarios:

Adapter Fusion. To combine knowledge from multiple trained adapters, we employ weighted fusion when fine-tuning is impractical due to limited data access or computational resources. However, this approach degrades performance. In Table 5, we compare the effectiveness of LoRA and S²FT when combining adapters trained separately on commonsense and arithmetic reasoning tasks, where we consider both fine-tuning overlapped and non-overlapped parameters for different adapters in S²FT. Our results show that S²FT with non-overlapped parameters achieves the best performance, while the overlapped variant shows inferior results. This is because S²FT (non-overlap) modifies orthogonal low-rank spaces for different tasks. Similarly, LoRA largely retains task-specific capabilities during adapter fusion by optimizing low-rank projection matrices to create separate spaces for each adapter.

Table 5: Adapter Fusion Results for LoRA and S²FT trained on the commonsense and arithmetic reasoning datasets using the LLama3-8B. We report the average accuracy (%) as the evaluation metric.

Task		LoRA		${f S}^2{f F}{f T}$				
	Commonsense	Arithmetic	Fused	Commonsense	Arithmetic	Fused (overlap)	Fused (non-overlap)	
Commonsense	83.1	32.1	79.8(-3.3)	86.6	42.3	82.0(-4.6)	84.0(-2.6)	
Arithmetic	12.0	77.2	$71.6_{\text{(-5.6)}}$	12.8	79.6	72.2(-7.4)	75.3 _(-4.3)	

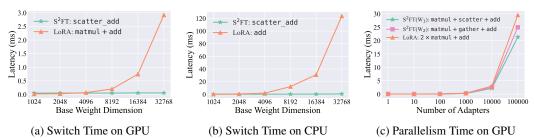


Figure 6: Comparison of latency for adapter switch and parallelism on a single linear layer. S²FT improves scalability for switch on GPU and CPU, while saving 22% time during parallelism on GPU.

Adapter Switch. Another way to leveraging multiple adapters is to dynamically switch between them. This process involves four steps: unfusing the old adapter, unloading it from memory, loading the new adapter, and fusing it into the model. In such setting, LoRA needs two matrix multiplications (matmul) and two additions (add) on GPU whereas S²FT only requires two sparse addition (scatter_add). In Figure 6a, we increase the base weight dimension while maintaining a sparsity of 32 for S²FT and a low-rankness of 16 for LoRA. Notably, we observe that LoRA's switching time scales quadratically, while S²FT remains nearly constant. Moreover, in I/O-constrained scenarios such as deployment on CPU, S²FT further accelerates adapter switch by only updating a small fraction of the original weights, reducing the volume of I/O transfers, as time compared between scatter_add and add in Figure 6b.

Adapter Parallelism. To serve thousands of adapters in parallel, we decompose the computation into separate batched computations for W^{pre} and ΔW following S-LoRA [60]. While LoRA requires two matmul and one add on GPU, S²FT reduces this to a matmul, an add, and either a scatter or gather for W_1 and W_2 in Section 3.1. Figure 6c shows that S²FT achieves up to 22% faster inference than LoRA under the same memory constraints, with more speedup as the number of adapters scales.

7 Related Work

PEFT methods reduce the fine-tuning cost for large models, which can be categorized into 4 groups: Adapter-based Fine-tuning introduces additional trainable module into the original model. Series Adapters insert components between MHA or FFN layers [51, 26], while parallel adapters add modules alongside existing components [24]. Recently, ReFT [69] was introduced to directly learn interventions on hidden representations. However, they introduce additional latency during inference. Prompt-based Fine-tuning adds randomly-initialized soft tokens to the input (usually as a prefix) and train their embeddings while freezing the model weights [36, 40, 35]. These approaches result in poor performance compared to other groups, while come at the cost of significant inference overhead.

Reparameterized Fine-tuning utilizes low-rank projections to reduce trainable parameters while allowing operations with high-dimensional matrices. LoRA[27] and its recent variants like DoRA[38], AsyLoRA [84], and FLoRA [61], use low-rank matrices to approximate additive weight updates during training. To alleviate the limitations of low-rank structure, other work also add or multiply orthogonal matrices to enable high-rank updating, including MoRA [29], OFT [54], and BOFT [39]. These methods require no additional inference cost as the weight updates can be merged into models.

Sparse Fine-tuning aims to reduce the number of fine-tuned parameters by selecting a subset of pre-trained parameters that are critical to downstream tasks while discarding unimportant ones. This kind of methods are commonly used in the pre-LLM era [20, 75, 64]. However, they cannot reduce the memory footprints due to their unstructured nature. Recent approaches address this limitation through three directions: (1) developing structured variants that sacrifice selection flexibility for better hardware efficiency [48, 85], (2) incorporating sparsity into LoRA [68, 15, 41] but yield limited efficiency gains, or (3) using sparse operators for lower memory cost but slow down training [4, 49, 7].

Our work is based on the last category but achieving better performance and efficiency simultaneously. Additionally, we focus on scalable inference of PEFT methods, with S²FT being the only approach that enables effective fusion, rapid switching, and efficient parallelism when serving multiple adapters.

8 Conclusion

This paper introduces S^2FT , a novel PEFT family that simultaneously achieves high quality, efficient training, and scalable serving for LLM fine-tuning. S^2FT accomplishes this by selecting sparsely and compute densely. It selects a subset of heads and channels to be trainable for the MHA and FFN modules, respectively. The weight matrices from the two sides of the coupled structures in LLMs are co-permuted to connect the selected components into dense matrices, and only these parameters are updated using dense operations. We hope S^2FT can be considered as a successor to LoRA for PEFT.

9 Acknowledgement

We would like to thank Songlin Yang, Kaustubh Ponkshe, Raghav Singhal, Jinqi Luo, Tianqi Chen, Hanshi Sun, and Chris De Sa for their helpful discussions, and the authors of LLM-Adapters, ReFT, and DoRA for providing detailed results.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv* preprint arXiv:2303.08774, 2023. 1
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023. 17
- [3] Alan Ansell, Edoardo Maria Ponti, Anna Korhonen, and Ivan Vulić. Composable sparse fine-tuning for cross-lingual transfer. *arXiv* preprint arXiv:2110.07560, 2021. 2
- [4] Alan Ansell, Ivan Vulić, Hannah Sterz, Anna Korhonen, and Edoardo M Ponti. Scaling sparse fine-tuning to large language models. *arXiv preprint arXiv:2401.16405*, 2024. 10
- [5] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. Advances in Neural Information Processing Systems, 32, 2019. 5, 18
- [6] Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for mathematics. arXiv preprint arXiv:2310.10631, 2023. 1
- [7] Kartikeya Bhardwaj, Nilesh Prasad Pandey, Sweta Priyadarshi, Viswanath Ganapathy, Rafael Esteves, Shreya Kadambi, Shubhankar Borse, Paul Whatmough, Risheek Garrepalli, Mart Van Baalen, et al. Rapid switching and multi-adapter fusion via sparse high rank adapters. arXiv preprint arXiv:2407.16712, 2024. 10
- [8] Dan Biderman, Jose Gonzalez Ortiz, Jacob Portes, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, et al. Lora learns less and forgets less. arXiv preprint arXiv:2405.09673, 2024. 1
- [9] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020. 16, 17
- [10] Jiaao Chen, Aston Zhang, Xingjian Shi, Mu Li, Alex Smola, and Diyi Yang. Parameter-efficient fine-tuning design spaces. arXiv preprint arXiv:2301.01821, 2023.
- [11] Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, et al. Spectral methods for data science: A statistical perspective. *Foundations and Trends*® *in Machine Learning*, 14(5):566–806, 2021. 31
- [12] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. arXiv preprint arXiv:1905.10044, 2019. 16, 17
- [13] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. arXiv preprint arXiv:1803.05457, 2018. 16, 17
- [14] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021. 16, 17
- [15] Ning Ding, Xingtai Lv, Qiaosen Wang, Yulin Chen, Bowen Zhou, Zhiyuan Liu, and Maosong Sun. Sparse low-rank adaptation of pre-trained language models. arXiv preprint arXiv:2311.11696, 2023. 10
- [16] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.

- [17] Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards any structural pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16091–16101, 2023. 4
- [18] Jianwei Feng and Dong Huang. Optimal gradient checkpoint search for arbitrary computation graphs. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 11433– 11442, 2021. 4
- [19] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. *Advances in neural information processing systems*, 30, 2017. 5, 18
- [20] Demi Guo, Alexander M Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning. arXiv preprint arXiv:2012.07463, 2020. 10
- [21] Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024. 1
- [22] Moritz Hardt and Tengyu Ma. Identity matters in deep learning. arXiv preprint arXiv:1611.04231, 2016.
 5, 18
- [23] Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. *arXiv* preprint arXiv:2402.12354, 2024. 3, 5, 6
- [24] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *arXiv* preprint arXiv:2110.04366, 2021. 6, 7, 10
- [25] Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 523–533, 2014. 16, 17
- [26] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In *International conference on machine learning*, pages 2790–2799. PMLR, 2019. 6, 7, 10
- [27] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685, 2021. 2, 6, 7, 10
- [28] Zhiqiang Hu, Yihuai Lan, Lei Wang, Wanyu Xu, Ee-Peng Lim, Roy Ka-Wei Lee, Lidong Bing, and Soujanya Poria. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. arXiv preprint arXiv:2304.01933, 2023. 3, 16, 17
- [29] Ting Jiang, Shaohan Huang, Shengyue Luo, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng, Feng Sun, Qi Zhang, Deqing Wang, et al. Mora: High-rank updating for parameter-efficient fine-tuning. arXiv preprint arXiv:2405.12130, 2024. 10
- [30] Kenji Kawaguchi. Deep learning without poor local minima. Advances in neural information processing systems, 29, 2016. 5, 18
- [31] Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. Parsing algebraic word problems into equations. *Transactions of the Association for Computational Linguistics*, 3:585–597, 2015. 16, 17
- [32] Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. Mawps: A math word problem repository. In Proceedings of the 2016 conference of the north american chapter of the association for computational linguistics: human language technologies, pages 1152–1157, 2016. 16, 17
- [33] Rui Kong, Qiyang Li, Xinyu Fang, Qingtian Feng, Qingfeng He, Yazhu Dong, Weijun Wang, Yuanchun Li, Linghe Kong, and Yunxin Liu. Lora-switch: Boosting the efficiency of dynamic llm adapters via system-algorithm co-design. *arXiv preprint arXiv:2405.17741*, 2024. 1
- [34] Thomas Laurent and James Brecht. Deep linear networks with arbitrary loss: All local minima are global. In *International conference on machine learning*, pages 2902–2907. PMLR, 2018. 5, 18
- [35] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. arXiv preprint arXiv:2104.08691, 2021. 10

- [36] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv* preprint arXiv:2101.00190, 2021. 6, 7, 10
- [37] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. *arXiv preprint arXiv:1705.04146*, 2017. 16, 17
- [38] Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. DoRA: Weight-Decomposed low-rank adaptation. arXiv preprint arXiv:2402.09353, 2024. 2, 6, 7, 10
- [39] Weiyang Liu, Zeju Qiu, Yao Feng, Yuliang Xiu, Yuxuan Xue, Longhui Yu, Haiwen Feng, Zhen Liu, Juyeon Heo, Songyou Peng, et al. Parameter-efficient orthogonal finetuning via butterfly factorization. *arXiv* preprint arXiv:2311.06243, 2023. 10
- [40] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. GPT understands, too. AI Open, 2023. 10
- [41] Zequan Liu, Jiawen Lyn, Wei Zhu, Xing Tian, and Yvette Graham. Alora: Allocating low-rank adaptation for fine-tuning large language models. *arXiv* preprint arXiv:2403.16187, 2024. 10
- [42] Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava, Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pages 22137–22176. PMLR, 2023. 2
- [43] Haihao Lu and Kenji Kawaguchi. Depth creates no bad local minima. arXiv preprint arXiv:1702.08580, 2017. 5, 18
- [44] Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *arXiv* preprint arXiv:2308.08747, 2023.
- [45] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. Advances in neural information processing systems, 36:21702–21720, 2023. 2, 4
- [46] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. arXiv preprint arXiv:1809.02789, 2018. 16, 17
- [47] Ryumei Nakada, Halil Ibrahim Gulluk, Zhun Deng, Wenlong Ji, James Zou, and Linjun Zhang. Understanding multimodal contrastive learning and incorporating unpaired data. arXiv preprint arXiv:2302.06232, 2023. 32
- [48] Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. LISA: Layerwise importance sampling for memory-efficient large language model fine-tuning. *arXiv* preprint *arXiv*:2403.17919, 2024. 2, 6, 7, 10
- [49] Ashwinee Panda, Berivan Isik, Xiangyu Qi, Sanmi Koyejo, Tsachy Weissman, and Prateek Mittal. Lottery ticket adaptation: Mitigating destructive interference in llms. arXiv preprint arXiv:2406.16797, 2024. 10
- [50] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple math word problems? In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2080–2094, Online, June 2021. Association for Computational Linguistics.
- [51] Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. Mad-x: An adapter-based framework for multi-task cross-lingual transfer. arXiv preprint arXiv:2005.00052, 2020. 10
- [52] Lai-Man Po, Yuyang Liu, Haoxuan Wu, Tianqi Zhang, Wing-Yin Yu, Zhuohan Wang, Zeyu Jiang, and Kun Li. Sbora: Low-rank adaptation with regional weight updates. arXiv preprint arXiv:2407.05413, 2024. 5, 6
- [53] Kaustubh Ponkshe, Raghav Singhal, Eduard Gorbunov, Alexey Tumanov, Samuel Horvath, and Praneeth Vepakomma. Initialization using update approximation is a silver bullet for extremely efficient low-rank fine-tuning. *arXiv* preprint arXiv:2411.19557, 2024. 5, 6
- [54] Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. *Advances in Neural Information Processing Systems*, 36:79320–79362, 2023. 10

- [55] Subhro Roy and Dan Roth. Solving general arithmetic word problems. *arXiv preprint arXiv:1608.01413*, 2016. 16, 17
- [56] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950, 2023. 1
- [57] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021. 16, 17
- [58] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiqa: Commonsense reasoning about social interactions. arXiv preprint arXiv:1904.09728, 2019. 16, 17
- [59] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. arXiv preprint arXiv:1312.6120, 2013. 5, 18
- [60] Ying Sheng, Shiyi Cao, Dacheng Li, Coleman Hooper, Nicholas Lee, Shuo Yang, Christopher Chou, Banghua Zhu, Lianmin Zheng, Kurt Keutzer, et al. S-lora: Serving thousands of concurrent lora adapters. arXiv preprint arXiv:2311.03285, 2023. 1, 2, 10
- [61] Chongjie Si, Xuehui Wang, Xue Yang, Zhengqin Xu, Qingyun Li, Jifeng Dai, Yu Qiao, Xiaokang Yang, and Wei Shen. Flora: Low-rank core space for n-dimension. *arXiv preprint arXiv:2405.14739*, 2024. 10
- [62] GW Stewart. On the continuity of the generalized inverse. SIAM Journal on Applied Mathematics, 17(1):33–45, 1969. 33
- [63] Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks. Advances in Neural Information Processing Systems, 34:24193–24205, 2021.
- [64] Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks. Advances in Neural Information Processing Systems, 34:24193–24205, 2021. 10
- [65] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca, 2023. 17
- [66] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 2023.
- [67] J Leo van Hemmen and Tsuneya Ando. An inequality for trace ideals. *Communications in Mathematical Physics*, 76:143–148, 1980. 33
- [68] Haoyu Wang, Tianci Liu, Tuo Zhao, and Jing Gao. Roselora: Row and column-wise sparse low-rank adaptation of pre-trained language model for knowledge editing and fine-tuning. arXiv preprint arXiv:2406.10777, 2024. 10
- [69] Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning, and Christopher Potts. ReFT: Representation finetuning for language models. arXiv preprint arXiv:2404.03592, 2024. 6, 7, 10
- [70] Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. arXiv preprint arXiv:2312.12148, 2023.
- [71] Runxin Xu, Fuli Luo, Zhiyuan Zhang, Chuanqi Tan, Baobao Chang, Songfang Huang, and Fei Huang. Raise a child in large language model: Towards effective and generalizable fine-tuning. *arXiv preprint arXiv:2109.05687*, 2021. 2
- [72] Ge Yang, Edward Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tuning large neural networks via zero-shot hyperparameter transfer. Advances in Neural Information Processing Systems, 34:17084–17097, 2021. 6
- [73] Yi Yu, Tengyao Wang, and Richard J Samworth. A useful variant of the davis–kahan theorem for statisticians. *Biometrika*, 102(2):315–323, 2015. 31
- [74] Li Yunxiang, Li Zihan, Zhang Kai, Dan Ruilong, and Zhang You. Chatdoctor: A medical chat model fine-tuned on llama model using medical domain knowledge. arXiv preprint arXiv:2303.14070, 2023. 1

- [75] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021. 10
- [76] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019. 16, 17
- [77] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021. 5, 18
- [78] Jinghan Zhang, Junteng Liu, Junxian He, et al. Composing parameter-efficient modules with arithmetic operation. *Advances in Neural Information Processing Systems*, 36:12589–12610, 2023. 1
- [79] Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning. *arXiv* preprint arXiv:2308.03303, 2023. 5, 6
- [80] Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. arXiv preprint arXiv:2403.03507, 2024. 2, 6, 7, 16
- [81] Justin Zhao, Timothy Wang, Wael Abid, Geoffrey Angus, Arnav Garg, Jeffery Kinnison, Alex Sherstinsky, Piero Molino, Travis Addair, and Devvret Rishi. Lora land: 310 fine-tuned llms that rival gpt-4, a technical report. arXiv preprint arXiv:2405.00732, 2024.
- [82] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. Advances in Neural Information Processing Systems, 36, 2024. 7, 17
- [83] Ming Zhong, Yelong Shen, Shuohang Wang, Yadong Lu, Yizhu Jiao, Siru Ouyang, Donghan Yu, Jiawei Han, and Weizhu Chen. Multi-lora composition for image generation. arXiv preprint arXiv:2402.16843, 2024. 9
- [84] Jiacheng Zhu, Kristjan Greenewald, Kimia Nadjahi, Haitz Sáez de Ocáriz Borde, Rickard Brüel Gabrielsson, Leshem Choshen, Marzyeh Ghassemi, Mikhail Yurochkin, and Justin Solomon. Asymmetry in low-rank adapters of foundation models. *arXiv preprint arXiv:2402.16842*, 2024. 10
- [85] Ligeng Zhu, Lanxiang Hu, Ji Lin, and Song Han. Lift: Efficient layer-wise fine-tuning for large model models. arXiv preprint arXiv:2410.11772, 2023. 10

A Limitations

While our work demonstrates the effectiveness of S^2FT for LLM fine-tuning, several promising directions remain unexplored. First, extending S^2FT to other architectures with coupled structures, such as CNNs and RNNs, can broaden its applicability. Second, verifying our approach beyond language tasks, particularly in large vision/multi-modal models, will enhance its versatility. Third, exploring more selection strategies can provide deeper insights into optimal fine-tuning protocols due to the controllability in S^2FT . Fourth, scaling our method to larger models requires further experiments. Finally, although our work confirms the feasibility of scalable and efficient deployment during inference, developing a practical serving system for S^2FT remains an important next step.

B Broader Impacts

Since our work focuses on PEFT, it leads to a reduction in hardware resource and energy consumption. Given the growing adoption of LLMs across diverse domains and the corresponding surge in fine-tuning demands, S²FT should represent an important step toward more sustainable AI development.

C Detailed Experimental Setups for Section 2

In this study, we used SpFT, LoRA, and Full FT to fine-tune the LLaMA-3-8B model on the Math10K dataset [28]. The Math10K dataset combines training sets from GSM8K [14], MAWPS [32], and AQuA [37], augmented with chain-of-thought steps generated by language models. We conducted training for 3 epochs with a batch size of 64. For both PEFT methods–SpFT and LoRA–we fine-tune with three ratios of trainable parameters for all linear layers: p=10%,1%,0.1%. The model's performance is evaluated on both arithmetic and commonsense reasoning tasks, representing near out-of-distribution (OOD) and far OOD generalization scenarios, respectively. The arithmetic reasoning dataset comprises seven subtasks: MultiArith [55], GSM8K, AddSub [25], AQuA, SingleEq [31], SVAMP [50], and MAWPS. The commonsense reasoning dataset includes eight subtasks: BoolQ [12], PIQA [9], SocialQA [58], HellaSwag [76], WinoGrande [57], ARC-challenge [13], ARC-easy [13], and OpenbookQA [46]. Based on task complexity within arithmetic reasoning (accuracy \geq 90%), we group MultiArith, AddSub, SingleEq, and MAWPS as easy subtasks, while the remaining ones are classified as hard subtasks. This stratification enables us to evaluate whether the model develops advanced reasoning abilities beyond memorizing basic arithmetic operations from the training data.

D Detailed Selection Strategies in Section 3

For the five selection strategies described in Section 3.2, we will detail the methods for identifying and selecting important subsets within each linear layer of both MHA and FFN modules in LLMs.

- 1. S^2FT-R (S^2FT): In this strategy, we will randomly select some heads for the MHA modules and select a few channels for the FFN modules. For the output projection, all channels in the selected heads will be included to enable dense-only computation. In the up and gate projections, we will select a subset of columns, while for the down projection, a few trainable rows will be chosen.
- 2. S^2FT-W : This variant selects subsets based on the weight magnitudes (i.e., $||W||_2$) in the MHA and FFN modules. We will test subsets corresponding to both the largest and smallest weights.
- 3. S^2FT-A : This variant selects subsets based on the magnitude of activations (i.e., $||A||_2$) on a calibration set, using 1% of the fine-tuning data. Since collecting activations requires only forward passes, this approach maintains the same memory footprint as inference and incurs a negligible increase in training time. Similarly, we evaluate both the largest and smallest activation variants.
- 4. **S**²**FT-S**: The Top-K subsets are ranked and selected by the product of the weight and activation magnitudes (i.e, $||W||_2 \cdot ||A||_2$). The activation values are collected in a manner similar to S²FT-A.
- 5. **S**²**FT-G**: This variant selects subsets based on the magnitude of gradients on the calibration set. Since gradients are collected without updating the model, we calculate and discard gradients layer by layer during back-propagation similar to Galore [80], requiring minimal additional memory.

E Detailed Experimental Setups for Section 5

Detailed selection strategies and number of trainable parameters are presented in Section 5.

E.1 Dataset Description

Commonsense Reasoning. The commonsense reasoning dataset comprise eight subsets: BoolQ [12], PIQA [9], SocialQA [58], HellaSwag [76], WinoGrande [57], ARC-challenge [13], ARC-easy [13], and OpenbookQA [46]. Following the experimental setup of LLM-Adapters [28], we split each dataset into training and test sets. Subsequently, we combine the training data from all eight tasks into a single fine-tuning dataset and evaluate performance on the individual test dataset for each task.

Arithmetic Reasoning. We followed Hu et al. [28] and evaluated S²FT on seven math reasoning tasks, including MultiArith [55], GSM8K [14], AddSub [25], AQuA [37], SingleEq [31], SVAMP [50] and MAWPS [32]. Our fine-tuning employed the Math10K dataset [28], which combines training sets from GSM8K, MAWPS, and AQuA, augmented with LM-generated chain-of-thought steps. Therefore, these three tasks are considered ID, while the remaining four are classified as OOD tasks.

Instruction Following. To further showcase S²FT's superior generalization ability, we employ the instruction-following fine-tuning task with Alpaca GPT-4 dataset, which comprises 52k samples generated by GPT-4 [2] based on inputs from Alpaca [65]. Performance is measured on MT-Bench [82], featuring 80 high-quality, multi-turn questions designed to assess LLMs on eight different aspects.

E.2 Hyperparameter Description

Additional hyperparameter configurations for all tasks are provided in Table 6. We maintain the same hyperparameter settings across the LLaMA-7/13B, LLaMA2-7B, LLaMA3-8B, and Mistral-7B models.

Commonsense Reasoning	Arithmetic Reasoning	Instruction Following
AdamW	AdamW	AdamW
2e-4	1e-3	2e-5
linear	linear	cosine
16×4	16×4	16×4
100	100	0
3	3	1
	AdamW 2e-4 linear 16×4	

Table 6: Hyperparameter configurations of S²FT on various base models across three tasks.

F Proofs for Theoretical Results in Section 4

Here we provide proofs for the results in Section 4.

F.1 Notation

For a vector a, let $\|a\|$ be the ℓ_2 norm of a. For $d_1 \geq d_2$, denote a set of orthogonal matrices by $\mathbb{O}_{d_1,d_2} := \{R \in \mathbb{R}^{d_1 \times d_2} : R^\top R = I_{d_2}\}$. For a matrix $A \in \mathbb{R}^{d_1 \times d_2}$, let $\|A\|_F$ and $\|A\|_{op}$ be the Frobenius norm and spectral norm of A, respectively. Denote the condition number of A by $\kappa_*(A) := \|A\|_{op}/\lambda_*(A)$. Let A^\dagger be Moore-Penrose inverse of A. For a symmetric matrix A, denote its effective rank by $r_e(A) := \operatorname{tr}(A)/\|A\|_{op}$. Note that $r_e(A) \leq \operatorname{rank}(A)$ always holds. For $a, b \in \mathbb{R}$, we let $a \vee b := \max(a, b)$ and $a \wedge b := \min(a, b)$. For a matrix $A \in \mathbb{R}^{d_1 \times d_2}$, let $\operatorname{SVD}_r(A) := \Phi_r(A)\Lambda_r(A)\Psi_r^\top(A)$ be the top-r singular value decomposition of A, where $\Phi_r(A) \in \mathbb{O}_{d_1,r}$ and $\Psi_r(A) \in \mathbb{O}_{d_2,r}$ are top-r left and right singular vectors of A, respectively, and $\Lambda_r(A) = \operatorname{diag}(\lambda_1(A),\ldots,\lambda_r(A)) \in \mathbb{R}^{r \times r}$ is a diagonal matrix of singular values of A, where $\lambda_j(A)$ denotes the j-th largest singular value of A. Define $\Phi_*(A) := \Phi_{\operatorname{rank}(A)}(A)$ and $\Psi_*(A) := \Psi_{\operatorname{rank}(A)}(A)$ as the left and right singular vectors of A corresponding to non-zero singular values, respectively. Define the smallest positive singular value of A as $\lambda_*(A) = \lambda_{\operatorname{rank}(A)}(A)$ and let $\Lambda_*(A) = \Lambda_{\operatorname{rank}(A)}(A)$. For a deep learning model fine-tuned on n i.i.d. samples $(x_i^{(i)}, y_i^{(i)}) \subset \mathbb{R}^p \times \mathbb{R}^q$, we say an event \mathcal{F} occurs with high probability when $\mathbb{P}(\mathcal{F}) = 1 - \exp(-\Omega(\log^2(n+p+q)))$.

F.2 Setup

We consider multivariate regression task. Using n i.i.d. samples $(x_i^{(i)}, y_i^{(i)}) \subset \mathbb{R}^p \times \mathbb{R}^q$ from in-distribution task, we fine-tune a pre-trained network $f^{\text{pre}} : \mathbb{R}^p \to \mathbb{R}^q$ for better prediction.

Deep Linear Networks We consider deep linear networks of the form $x\mapsto W_LW_{L-1}\dots W_1x:\mathbb{R}^d\to\mathbb{R}^p$, where $W_\ell\in\mathbb{R}^{d_\ell\times d_{\ell-1}}$, with $d_L=q$ and $d_0=p$. In comparison to multi-head attention transformers, each row of W_ℓ can be viewed as corresponding to the parameters in a single head. Let $f^{\mathrm{pre}}(x)=W_L^{\mathrm{pre}}W_{L-1}^{\mathrm{pre}}\dots W_1^{\mathrm{pre}}x:\mathbb{R}^p\to\mathbb{R}^q$ represent a pre-trained neural network. We denote $\overline{W}_\ell^{\mathrm{pre}}:=W_L^{\mathrm{pre}}W_{L-1}^{\mathrm{pre}}\dots W_\ell^{\mathrm{pre}}\in\mathbb{R}^{d_L\times d_{\ell-1}}$ as the weights up to the ℓ -th layer, and $\underline{W}_\ell^{\mathrm{pre}}:=W_\ell^{\mathrm{pre}}W_{\ell-1}^{\mathrm{pre}}\dots W_1^{\mathrm{pre}}\in\mathbb{R}^{d_\ell\times d_0}$ as the weights above the ℓ -th layer, with the promise that $\underline{W}_0^{\mathrm{pre}}=I$. Deep linear networks have been widely used to facilitate the theoretical analysis of modern complex deep neural networks [59, 30, 43, 22, 34, 5].

Fine-Tuning We employ ℓ_2 distance as the error metric. Given a pre-trained network f^{pre} , we fine-tune its ℓ -th layer by minimizing the empirical in-distribution risk $\mathcal{R}_n^{(i)}(f) := (1/n) \sum_{i \in [n]} \|y_i^{(i)} - f(x_i^{(i)})\|^2$, where $(x_i^{(i)}, y_i^{(i)}) \subset \mathbb{R}^p \times \mathbb{R}^q$ are n i.i.d. observations from in-distribution task. More specifically, we consider a class of rank-d adaptation defined as

$$f_{\ell,U,V}(x) := \overline{W}_{\ell+1}^{\text{pre}} (W_{\ell}^{\text{pre}} + UV^{\top}) \underline{W}_{\ell-1}^{\text{pre}} x, \tag{4}$$

where $U \in \mathbb{R}^{d_\ell \times d}$ and $V \in \mathbb{R}^{d_{\ell-1} \times d}$ are parameters to fine-tune. Note that by regarding multiple consecutive layers as a single layer, our settings can be extended to multi-layer fine-tuning.

We specifically compare two fine-tuning methods: LoRA and S²FT.

• LoRA. For a fixed $\ell \in [L]$, and low-rankness level $1 \le r \le \min\{d_\ell, d_{\ell-1}\}$, we train the low-rank matrices (U, V) in (4) by minimizing the empirical in-distribution risk via gradient descent. Motivated from the previous results that gradient descent has implicit regularization [77, 19, 5], we directly consider the minimum norm solutions:

$$(U^{\text{LoRA}}, V^{\text{LoRA}}) \in \operatorname*{arg\,min}_{U,V} \|(U,V)\|_{\text{F}}^2 \quad \text{s.t. } (U,V) \text{ minimizes } \mathcal{R}_n^{(\text{i})}(f_{\ell,U,V}). \tag{5}$$

• $\mathbf{S}^2\mathbf{FT}$. For a fixed $\ell \in [L]$, and a sparsity level $s = \lfloor r \cdot \frac{d_\ell + d_{\ell-1}}{d_{\ell-1}} \rfloor$, we train only V in (4) with the fixed choice of $U \leftarrow U_S^{\mathbf{S}^2\mathbf{FT}} := [e_{a_1}; e_{a_2}; \dots; e_{a_s}]$, which specifies s channels to fine-tune, where $S = \{a_1, a_2, \dots, a_s\} \subset [d_\ell]$. Here e_a is the standard basis vector with the a-th entry being 1. We minimize the empirical in-distribution risk via gradient descent. Similar to LoRA, we consider the following minimum norm solution:

$$V^{\mathrm{S}^2\mathrm{FT}} = \arg\min_{V} \|V\|_{\mathrm{F}}^2 \quad \text{s.t. } V \text{ minimizes } \mathcal{R}_n^{(\mathrm{i})}(f_{\ell,U_S^{\mathrm{S}^2\mathrm{FT}},V}). \tag{6}$$

Data Generating Process As a simplification of the data generating process, we consider multiple linear regression. Assume that the in-distribution data $(x^{(i)}, y^{(i)}) \in \mathbb{R}^{p+q}$ and out-of-distribution data $(x^{(o)}, y^{(o)}) \in \mathbb{R}^{p+q}$ are generated according to

$$y^{(k)} = B^{(k)}x^{(k)} + \epsilon^{(k)}, \quad k \in \{i, o\},$$
(7)

where $B^{(k)} \in \mathbb{R}^{q \times p}$, and $\epsilon^{(k)} \in \mathbb{R}^q$ is the error term satisfying $\mathbb{E}[\epsilon^{(k)}|x^{(k)}] = 0$. Assume that $\Sigma_{\epsilon}^{(k)} := \mathbb{E}[\epsilon^{(k)}\epsilon^{(k)\top}] \in \mathbb{R}^{q \times q}$ exists and $\mathbb{E}[x^{(k)}] = 0$. The signal covariance matrix is denoted by $\Sigma_x^{(k)} := \mathbb{E}[x^{(k)}x^{(k)\top}] \in \mathbb{R}^{p \times p}$.

We define the in-distribution and out-of-distribution risks of $f: \mathbb{R}^p \to \mathbb{R}^q$ as:

$$\mathcal{R}^{(k)}(f) = \mathbb{E}[\|y^{(k)} - f(x^{(k)})\|], \ k \in \{i, o\}.$$

For notational brevity, we can write $W^{\mathrm{pre}} = \underline{W}^{\mathrm{pre}}_L \in \mathbb{R}^{q \times p}$. Let $X^{(\mathrm{i})} := (x_1^{(\mathrm{i})}, \dots, x_n^{(\mathrm{i})}) \in \mathbb{R}^{p \times n}$, $Y^{(\mathrm{i})} := (y_1^{(\mathrm{i})}, \dots, y_n^{(\mathrm{i})}) \in \mathbb{R}^{q \times n}$, and $E^{(\mathrm{i})} = (\epsilon_1^{(\mathrm{i})}, \dots, \epsilon_n^{(\mathrm{i})}) := Y^{(\mathrm{i})} - B^{(\mathrm{i})}X^{(\mathrm{i})} \in \mathbb{R}^{q \times n}$. Denote the

in-distribution sample covariance matrices by $\hat{\Sigma}_x^{(i)} := (1/n)X^{(i)}X^{(i)}$, $\hat{\Sigma}_{\epsilon}^{(i)} := (1/n)E^{(i)}E^{(i)}$, $\hat{\Sigma}_{\epsilon,x}^{(i)} := \hat{\Sigma}_{x,\epsilon}^{(i)}$. Define $\check{\Sigma}_{x,\epsilon}^{(i)} := (X^{(i)})^{\dagger}E^{(i)}$, $\hat{A} := (\underline{W}_{\ell-1}^{\mathrm{pre}}\hat{\Sigma}_x^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}})^{1/2}$, $\hat{\Delta}_{\epsilon,x}^{(i)} := \hat{\Sigma}_{x,\epsilon}^{(i)}$. Define $\check{\Sigma}_{x,\epsilon}^{(i)} := (X^{(i)})^{\dagger}E^{(i)}$, $\hat{A} := (\underline{W}_{\ell-1}^{\mathrm{pre}}\hat{\Sigma}_x^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}})^{1/2}$, $\hat{\Delta}_{\epsilon,x}^{(i)} := \hat{\Delta}_{\epsilon,x}^{(i)}$, $\hat{\Delta}_{\epsilon,x}^{\mathrm{pre}} := \hat{\Delta}_{\epsilon,x}^{(i)}$, $\hat{\Delta}_{\epsilon,x}^{\mathrm{pre}} := \hat{\Delta}_{\epsilon,x}^{(i)}$. Also define $\hat{\Delta}_{\epsilon,x}^{(i)} := \hat{\Delta}_{\epsilon,x}^{(i)}$. Also define $\hat{\Delta}_{\epsilon,x}^{\mathrm{pre}} := \hat{\Delta}_{\epsilon,x}^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}}\hat{\Delta}_{\epsilon,x}^{(i)}$ and $\hat{\Delta}_{\epsilon,x}^{\mathrm{pre}} := \hat{\Delta}_{\epsilon,x}^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}}\hat{\Delta}_{\epsilon,x}^{(i)}$ be a matrix that captures the covariate shift at the ℓ -th layer.

We consider fine-tuning the ℓ -th ($\ell \in [L]$) layer of the pre-trained deep linear network $f^{\text{pre}}(x) = W_L^{\text{pre}} W_{L-1}^{\text{pre}} \dots W_1^{\text{pre}} x$ using in-distribution observations $(x_i^{(i)}, y_i^{(i)})_{i \in [n]}$.

To measure the performance of models, we define the excess risks of f for the task $k \in \{i, o\}$ as

$$\mathcal{E}^{(k)}(f) := \mathbb{E}[\|y^{(k)} - f(x^{(k)})\|^2] - \inf_{f'} \mathbb{E}[\|y^{(k)} - f'(x^{(k)})\|^2],$$

where the infimum is taken over all square integrable functions.

F.3 Assumptions

We assume that $\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_x^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}\neq 0$, since otherwise $\underline{W}_{\ell-1}^{\mathrm{pre}}x^{(\mathrm{i})}=0$ almost surely and fine-tuning the ℓ -th layer does not improve the performance of the pre-trained model. Define the in-distribution prediction residuals for the pre-trained model f^{pre} by $\Sigma_f^{(\mathrm{i})}:=\mathbb{E}[(B^{(\mathrm{i})}x^{(\mathrm{i})}-W^{\mathrm{pre}}x^{(\mathrm{i})})(B^{(\mathrm{i})}x^{(\mathrm{i})}-W^{\mathrm{pre}}x^{(\mathrm{i})})^{\top}]$. Note that $\mathcal{E}^{(\mathrm{i})}(f^{\mathrm{pre}})=\mathrm{tr}\Big(\Sigma_f^{(\mathrm{i})}\Big)$. We also assume that $\|\Sigma_f^{(\mathrm{i})}\|_{\mathrm{op}}>0$, since otherwise $\mathcal{E}^{(\mathrm{i})}(f^{\mathrm{pre}})=\|\Sigma_f^{(\mathrm{i})}\|_{\mathrm{F}}^2=0$ and there is no room for improvement from the pre-trained model.

Next, we introduce several assumptions.

Assumption F.1 (Sub-Gaussianity). Assume that there exist some constants $c_1, c_2 \in (0, \infty)$ such that $(x^{(i)}, \epsilon^{(i)})$ in the model 7 satisfies

$$\gamma^{\top} \Sigma_x^{(i)} \gamma \geq c_1 \| \gamma^{\top} x^{(i)} \|_{\psi_2}^2$$
, and $\gamma^{\prime \top} \Sigma_{\epsilon}^{(i)} \gamma^{\prime} \geq c_2 \| \gamma^{\prime \top} \epsilon^{(i)} \|_{\psi_2}^2$,

for any $\gamma \in \mathbb{R}^p$ and $\gamma' \in \mathbb{R}^q$, where $\|y\|_{\psi_2}$ is the sub-Gaussian norm defined as

$$||y||_{\psi_2} := \inf\{v > 0 : \mathbb{E}[\exp(y^2/v^2)] \le 2\}$$

for a random variable y taking values in \mathbb{R} .

Assumption F.2 (Sufficiently Many Observations). Assume that

$$n \gg (\kappa_*^4(A)r_e(A^2) + \kappa_*^2(\Sigma_x^{(i)})r_e(\Sigma_x^{(i)}) + r_e(D\Sigma_x^{(i)}D^\top))\log^2(n+p+q),$$

$$n \gg \frac{\|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}}{\|D\Sigma_x^{(i)}D^\top\|_{\text{op}}}(r_e(\Sigma_{\epsilon}^{(i)}) + r_e(A^2))\log^2(n+p+q),$$

and

$$n \gg \kappa_*^4(\Sigma_x^{(i)}) \frac{r_e(\Sigma_x^{(i)}) (r_e(\Sigma_\epsilon^{(i)}) + r_e(\Sigma_x^{(i)}))}{r_e(A^2)} \log^2(n + p + q).$$

Assumption F.3 (Eigengap Condition). Assume that there exists some constant $C_g > 0$ such that

$$\frac{\lambda_s(\Phi'^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger})}{\lambda_s(\Phi'^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger}) - \lambda_{s+1}(\Phi'^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger})} \lesssim C_g$$

holds.

Assumption F.3 is necessary to identify the rank-r approximation of M, which is used to derive the risk of LoRA.

Assumption F.4 (Approximate Sparsity of Channels). Assume that there exists some $S_0 \subset [d_\ell]$ with $|S_0| \leq s$ and $\delta > 0$ such that

$$\sum_{a \in [d_{\ell}] \backslash S_0} \|e_a^\top (\overline{W}_{\ell+1}^{\mathrm{pre}})^\dagger (B^{(\mathrm{i})} - W^{\mathrm{pre}}) \Sigma_x^{(\mathrm{i})1/2} \|^2 \leq \delta^2 \| (\overline{W}_{\ell+1}^{\mathrm{pre}})^\dagger (B^{(\mathrm{i})} - W^{\mathrm{pre}}) \Sigma_x^{(\mathrm{i})1/2} \|_{\mathrm{F}}^2$$

holds.

Assumption F.5 (Distribution Shift). Assume that $\Sigma_x^{(i)} = \Sigma_x^{(o)} = \Sigma_x$ for some $\Sigma_x \in \mathbb{R}^{d \times d}$ and that $\|\Phi_*^{\top}(\overline{W}_{\ell+1}^{\operatorname{pre}}U_S^{2^*\operatorname{FT}})(B^{(o)} - B^{(i)})\Sigma_x^{1/2}\|_{\operatorname{F}}^2 \leq \varepsilon^2 \mathcal{E}^{(o)}(f^{\operatorname{pre}})$ for some $\varepsilon > 0$.

Assumption F.6 (Condition Number). Assume that $\kappa_*(M) \lesssim 1$, $\kappa_*(\overline{W}_{\ell+1}^{\mathrm{pre}}) \lesssim 1$, $\kappa_*(\Sigma_f^{(\mathrm{i})}) \lesssim 1$ and $\kappa_*(\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_x^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}) \lesssim 1$.

Note that Assumption F.6 is not essential to our analysis.

F.4 Main Results

We first demonstrate that LoRA and S²FT exhibit comparable memorization abilities. Next, we present a formal restatement of 4.2 that combine Theorems F.10, F.11, F.13, F.15, and Lemma F.14.

Theorem F.7. Suppose that Assumptions F.1, F.2, F.3, F.4, and F.6 hold. Choose S such that $S \supset S_0$ holds. Let U^{Lora} , V^{Lora} be the Loral adaptation matrices defined in (5). Let V^{S^2FT} be the S^2FT adaptation matrices given U^{S^2FT} defined in (6). Then, for all sufficiently large n, the following holds with probability $1 - \exp(-\Omega(\log^2(n+p+q)))$: for any $\eta > 0$,

$$\begin{split} \mathcal{E}^{(\mathrm{i})}(f_{\ell,U_S^{\mathbf{S}^2\mathrm{FT}},V^{\mathbf{S}^2\mathrm{FT}}}) &\leq (1+\eta)(T_{\mathrm{bias}}^{\mathbf{S}^2\mathrm{FT}})^2 + (1+\eta^{-1})(T_{\mathrm{variance}}^{\mathbf{S}^2\mathrm{FT}})^2, \\ \mathcal{E}^{(\mathrm{i})}(f_{\ell,U_{\mathrm{LORA}},V^{\mathrm{LORA}}}) &\leq (1+\eta)(T_{\mathrm{bias}}^{\mathrm{LoRA}})^2 + (1+\eta^{-1})(T_{\mathrm{variance}}^{\mathrm{LoRA}})^2, \end{split}$$

where

$$\begin{split} 0 & \leq (T_{\text{bias}}^{\text{LoRA}})^2 - \mathcal{E}^{(\text{i})}(f_{\ell}^{\text{full}}) \simeq (T_{\text{bias}}^{\text{S}^2\text{FT}})^2 - \mathcal{E}^{(\text{i})}(f_{\ell}^{\text{full}}) \lesssim \delta^2 \mathcal{E}^{(\text{i})}(f^{\text{pre}}), \\ (T_{\text{variance}}^{\text{S}^2\text{FT}})^2 & \lesssim (\|\Sigma_{\epsilon}^{(\text{i})}\|_{\text{op}} + \|\Sigma_{f}^{(\text{i})}\|_{\text{op}}) \frac{sd_{\ell-1}\log^2(n+p+q)}{n}, \\ (T_{\text{variance}}^{\text{LoRA}})^2 & \lesssim (\|\Sigma_{\epsilon}^{(\text{i})}\|_{\text{op}} + \|\Sigma_{f}^{(\text{i})}\|_{\text{op}}) \frac{r(d_{\ell} + d_{\ell-1})\log^2(n+p+q)}{n}. \end{split}$$

Theorem F.8 (Restatement of Theorem 4.2). Consider the limit $n \to \infty$. Suppose that Assumption F.5 holds. Let $U^{\text{LoRA}}, V^{\text{LoRA}}$ be the LoRA adaptation matrices defined in (15). Let $V^{\text{S}^2\text{FT}}$ be the S^2FT adaptation matrices given $U_S^{\text{S}^2\text{FT}}$ defined in (25). If $B^{(i)} = \overline{W}_{\ell+1}^{\text{pre}} \tilde{B} \underline{W}_{\ell-1}^{\text{pre}}$ holds for some $\tilde{B}^{(i)} \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$, and $s, r \leq \text{rank}(\Sigma_f^{(i)})$, then,

$$\mathcal{E}^{(o)}(f_{\ell, U_S^{S^2FT}, V^{S^2FT}}) \le (1 + 3\varepsilon^2) \mathcal{E}^{(o)}(f^{\text{pre}}),$$

$$\mathcal{E}^{(o)}(f_{\ell, U^{\text{LORA}}, V^{\text{LORA}}}) \ge \|(B^{(o)} - B^{(i)}) \Sigma_{\tau}^{1/2}\|_{\mathcal{F}}^{2}$$

Intuition of the proof of Theorem F.8. LoRA forgets pre-trained tasks due to its model complexity. Consider the simplest low-rank adaptation to a single-layer linear network:

$$\Delta_1 \in \operatorname*{arg\,min}_{\substack{\Delta_1' \in \mathbb{R}^{d_1 \times d_0} \\ \operatorname{rank}(\Delta_1') = r}} \mathbb{E}[\|y^{(i)} - (W_1^{\operatorname{pre}} + \Delta_1')x^{(i)}\|^2].$$

Assume that $\Sigma_x^{(\mathrm{i})}=I$, then we can show that the solution is $\Delta_1=\mathrm{SVD}_r(B^{(\mathrm{i})}-W_1^{\mathrm{pre}})$. Under the condition that the rank of $B^{(\mathrm{i})}-W_1^{\mathrm{pre}}$ is smaller than, or comparable to r, LoRA fine-tuned model can learn the in-distribution best regressor in ℓ_2 sense, since $(W_1^{\mathrm{pre}}+\Delta_1)x\approx B^{(\mathrm{i})}x=\mathbb{E}[y^{(\mathrm{i})}|x^{(\mathrm{i})}=x]$. Hence it makes LoRA fine-tuned model vulunerable to distribution shift.

On the other hand, we model S²FT as fine-tuning only a few channels:

$$\Delta_1 \in \mathop{\arg\min}_{\Delta_1' = \sum_{a \in S} e_a v_a^\top, v_a \in \mathbb{R}^{d_0}} \mathbb{E}[\|y^{(i)} - (W_1^{\mathsf{pre}} + \Delta_1') x^{(i)}\|^2].$$

Although S²FT is a special case of LoRA, the constraint on the direction of low-rank matrix prevents overfitting to the in-distribution task. To see this, note that a sparse fine-tuned model can be written as

$$(W_1^{\text{pre}} + \Delta_1)x = W_1^{\text{pre}}x + \sum_{a \in S} e_a e_a^\top (B^{(\mathbf{i})} - W_1^{\text{pre}})x = \sum_{a \in S^c} e_a e_a^\top W_1^{\text{pre}}x + \sum_{a \in S} e_a e_a^\top B^{(\mathbf{i})}x,$$

where $S \subset [d_1]$ is a set of channels with cardinality s. Since S²FT keeps most of parameters from the pre-trained model, except for rows specified by S, the model forget less pre-training tasks.

F.5 Proofs for LoRA

F.5.1 Excess Risk of LoRA

Lemma F.9 (Excess Risk). Consider the minimum norm solution

$$(U^{\mathrm{LoRA}}, V^{\mathrm{LoRA}}) \in \mathop{\arg\min}_{(U,V) \in \mathbb{R}^{d_{\ell} \times r} \times \mathbb{R}^{d_{\ell-1} \times r}} \|(U,V)\|_{\mathrm{F}}^{2} \quad \textit{s.t.} \ (U,V) \ \textit{minimizes} \ \mathcal{R}_{n}^{(\mathrm{i})}(f_{\ell,U,V}).$$

Then, the low-rank adaptation matrix satisfies

$$U^{\text{Lora}}V^{\text{Lora}\top} = (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger}SVD_r(\overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^{\dagger}\hat{D}\hat{\Sigma}_x^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}\hat{A}^{\dagger})\hat{A}^{\dagger},$$

and

$$\mathcal{E}^{(k)}(f_{\ell,U^{\mathsf{LORA}},V^{\mathsf{LORA}}}) = \operatorname{tr}\left(\left(B^{(k)} - W^{\mathsf{pre}} - SVD_r(\overline{W}_{\ell+1}^{\mathsf{pre}}(\overline{W}_{\ell+1}^{\mathsf{pre}})^{\dagger}\hat{D}\hat{\Sigma}_x^{(i)}\underline{W}_{\ell-1}^{\mathsf{pre}\top}\hat{A}^{\dagger})\hat{A}^{\dagger}\underline{W}_{\ell-1}^{\mathsf{pre}}\right)\Sigma_x^{(k)}$$
$$\cdot \left(B^{(k)} - W^{\mathsf{pre}} - SVD_r(\overline{W}_{\ell+1}^{\mathsf{pre}}(\overline{W}_{\ell+1}^{\mathsf{pre}})^{\dagger}\hat{D}\hat{\Sigma}_x^{(i)}\underline{W}_{\ell-1}^{\mathsf{pre}\top}\hat{A}^{\dagger})\hat{A}^{\dagger}\underline{W}_{\ell-1}^{\mathsf{pre}}\right)^{\top}\right)$$

for $k \in \{i, o\}$.

Proof of Lemma F.9. The empirical risk of $f_{\ell,U,V}$ for the in-distribution task can be written as

$$\begin{split} \mathcal{R}_{n}^{(\mathbf{i})}(f_{\ell,U,V}) &= \frac{1}{n} \sum_{i \in [n]} \| (B^{(\mathbf{i})} - W^{\mathrm{pre}}) x_{i}^{(\mathbf{i})} + \epsilon_{i}^{(\mathbf{i})} - \overline{W}^{\mathrm{pre}}_{\ell+1} U V^{\top} \underline{W}^{\mathrm{pre}}_{\ell-1} x_{i}^{(\mathbf{i})} \|^{2} \\ &= \mathrm{tr} \Big((B^{(\mathbf{i})} - W^{\mathrm{pre}} - \overline{W}^{\mathrm{pre}}_{\ell+1} U V^{\top} \underline{W}^{\mathrm{pre}}_{\ell-1}) \hat{\Sigma}_{x}^{(\mathbf{i})} (B^{(\mathbf{i})} - W^{\mathrm{pre}} - \overline{W}^{\mathrm{pre}}_{\ell+1} U V^{\top} \underline{W}^{\mathrm{pre}}_{\ell-1})^{\top} \Big) \\ &+ 2 \operatorname{tr} \Big((B^{(\mathbf{i})} - W^{\mathrm{pre}} - \overline{W}^{\mathrm{pre}}_{\ell+1} U V^{\top} \underline{W}^{\mathrm{pre}}_{\ell-1}) \hat{\Sigma}_{x,\epsilon}^{(\mathbf{i})} \Big) + \operatorname{tr} \Big(\hat{\Sigma}_{\epsilon}^{(\mathbf{i})} \Big) \\ &= \mathrm{tr} \Big(V^{\top} \underline{W}^{\mathrm{pre}}_{\ell-1} \hat{\Sigma}_{x}^{(\mathbf{i})} \underline{W}^{\mathrm{pre}}_{\ell-1} V U^{\top} \overline{W}^{\mathrm{pre}}_{\ell+1} \overline{W}^{\mathrm{pre}}_{\ell+1} U \Big) \\ &- 2 \operatorname{tr} \Big(\overline{W}^{\mathrm{pre}}_{\ell+1} U V^{\top} \underline{W}^{\mathrm{pre}}_{\ell-1} \Big\{ \hat{\Sigma}_{x}^{(\mathbf{i})} (B^{(\mathbf{i})} - W^{\mathrm{pre}})^{\top} + \hat{\Sigma}_{x,\epsilon}^{(\mathbf{i})} \Big\} \Big) \\ &+ \mathrm{tr} \Big((B^{(\mathbf{i})} - W^{\mathrm{pre}}) \hat{\Sigma}_{x}^{(\mathbf{i})} (B^{(\mathbf{i})} - W^{\mathrm{pre}})^{\top} \Big) + 2 \operatorname{tr} \Big((B^{(\mathbf{i})} - W^{\mathrm{pre}}) \hat{\Sigma}_{x,\epsilon}^{(\mathbf{i})} \Big) + \operatorname{tr} \Big(\hat{\Sigma}_{\epsilon}^{(\mathbf{i})} \Big). \end{split}$$

Since $\hat{\Sigma}_{x,\epsilon}^{(\mathrm{i})} = \hat{\Sigma}_{x}^{(\mathrm{i})} (X^{(\mathrm{i})\top})^{\dagger} E^{(\mathrm{i})\top} = \hat{\Sigma}_{x}^{(\mathrm{i})} \check{\Sigma}_{x,\epsilon}^{(\mathrm{i})}$

$$\mathcal{R}_{n}^{(i)}(f_{\ell,U,V}) = \operatorname{tr}\left(\hat{A}VU^{\top}\overline{W}_{\ell+1}^{\operatorname{pre}}^{\operatorname{pre}}\overline{W}_{\ell+1}^{\operatorname{pre}}UV^{\top}\hat{A}\right) - 2\operatorname{tr}\left(\overline{W}_{\ell+1}^{\operatorname{pre}}UV^{\top}\hat{A}\hat{A}^{\dagger}\underline{W}_{\ell-1}^{\operatorname{pre}}\hat{\Sigma}_{x}^{(i)}\hat{D}^{\top}\right) \\
- 2\operatorname{tr}\left(\overline{W}_{\ell+1}^{\operatorname{pre}}UV^{\top}(I - \hat{A}\hat{A}^{\dagger})\underline{W}_{\ell-1}^{\operatorname{pre}}\hat{\Sigma}_{x}^{(i)}\hat{D}^{\top}\right) \\
+ \operatorname{tr}\left(D\hat{\Sigma}_{x}^{(i)}D^{\top}\right) + 2\operatorname{tr}\left(D\hat{\Sigma}_{x,\epsilon}^{(i)}\right) + \operatorname{tr}\left(\hat{\Sigma}_{\epsilon}^{(i)}\right) \\
= \|\overline{W}_{\ell+1}^{\operatorname{pre}}UV^{\top}\hat{A} - \hat{D}\hat{\Sigma}_{x}^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}}\hat{A}^{\dagger}\|_{F}^{2} - \|\hat{D}\hat{\Sigma}_{x}^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}}\hat{A}^{\dagger}\|_{F}^{2} \\
+ \operatorname{tr}\left(D\hat{\Sigma}_{x}^{(i)}D^{\top}\right) + 2\operatorname{tr}\left(D\hat{\Sigma}_{x,\epsilon}^{(i)}\right) + \operatorname{tr}\left(\hat{\Sigma}_{\epsilon}^{(i)}\right), \tag{9}$$

where we used $(I - \hat{A}\hat{A}^{\dagger})\underline{W}_{\ell-1}^{\mathrm{pre}}\hat{\Sigma}_{x}^{(\mathrm{i})1/2} = 0$. From (9), minimizing $\mathcal{R}_{n}^{(\mathrm{i})}(f_{\ell,U,V})$ is equivalent to minimizing the norm:

$$\begin{split} \|\overline{W}_{\ell+1}^{\text{pre}}UV^{\top}\hat{A} - \hat{D}\hat{\Sigma}_{x}^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}\hat{A}^{\dagger}\|_{\text{F}}^{2} &= \|\overline{W}_{\ell+1}^{\text{pre}}UV^{\top}\hat{A} - \overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^{\dagger}\hat{D}\hat{\Sigma}_{x}^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}\hat{A}^{\dagger}\|_{\text{F}}^{2} \\ &+ \|(I - \overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^{\dagger})\hat{D}\hat{\Sigma}_{x}^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}\hat{A}^{\dagger}\|_{\text{F}}^{2}. \end{split}$$

This is minimized by (U', V') satisfying

$$U'V'^{\top} = (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} \text{SVD}_{r} (\overline{W}_{\ell+1}^{\text{pre}}) (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} \hat{A}^{\dagger}) \hat{A}^{\dagger}$$
$$+ (I - (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} \overline{W}_{\ell+1}^{\text{pre}}) A_{1} + A_{2} (I - \hat{\Psi}' \hat{\Psi}'^{\top}), \tag{10}$$

where $A_1, A_2 \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$ are arbitrary matrices. Since we particularly consider the minimum norm solution, we must have $A_1 = 0$ and $A_2 = 0$. Hence

$$\overline{W}_{\ell+1}^{\mathrm{pre}} U^{\mathrm{LoRA}} V^{\mathrm{LoRA}\top} \underline{W}_{\ell-1}^{\mathrm{pre}} = \mathrm{SVD}_r (\overline{W}_{\ell+1}^{\mathrm{pre}} (\overline{W}_{\ell+1}^{\mathrm{pre}})^\dagger \hat{D} \hat{\Sigma}_x^{(\mathrm{i})} \underline{W}_{\ell-1}^{\mathrm{pre}\top} \hat{A}^\dagger) \hat{A}^\dagger \underline{W}_{\ell-1}^{\mathrm{pre}}.$$

Therefore, the excess risk for $k \in \{i, o\}$ becomes

$$\begin{split} \mathcal{E}^{(k)}(f_{\ell,U^{\mathsf{LORA}},V^{\mathsf{LORA}}}) &= \mathbb{E}\bigg[\Big(B^{(k)} x^{(k)} - \overline{W}^{\mathsf{pre}}_{\ell+1} (W^{\mathsf{pre}}_{\ell} + U^{\mathsf{LORA}} V^{\mathsf{LORA}}^{\mathsf{TORA}}) \underline{W}^{\mathsf{pre}}_{\ell-1} x^{(k)} \Big)^2 \bigg] \\ &= \mathrm{tr}\bigg(\Big(B^{(k)} - W^{\mathsf{pre}} - \mathsf{SVD}_r (\overline{W}^{\mathsf{pre}}_{\ell+1} (\overline{W}^{\mathsf{pre}}_{\ell+1})^\dagger \hat{D} \hat{\Sigma}^{(i)}_x \underline{W}^{\mathsf{pre}}_{\ell-1} \hat{A}^\dagger) \hat{A}^\dagger \underline{W}^{\mathsf{pre}}_{\ell-1} \Big) \Sigma^{(k)}_x \\ & \cdot \Big(B^{(k)} - W^{\mathsf{pre}} - \mathsf{SVD}_r (\overline{W}^{\mathsf{pre}}_{\ell+1} (\overline{W}^{\mathsf{pre}}_{\ell+1})^\dagger \hat{D} \hat{\Sigma}^{(i)}_x \underline{W}^{\mathsf{pre}}_{\ell-1} \hat{A}^\dagger) \hat{A}^\dagger \underline{W}^{\mathsf{pre}}_{\ell-1} \Big)^\top \bigg). \end{split}$$

This concludes the proof.

F.5.2 In-distribution Excess Risk of LoRA

Let $\mathcal{E}^{(i)}(f_{\ell}^{\text{full}})$ denote the excess risk of f^{pre} after fine-tuning all the parameters of the ℓ -th layer under *population* in-distribution risk.

Theorem F.10 (Restatement of Theorem F.7: LoRA Part). Suppose that Assumptions F.1, F.2 and F.3 hold. Then, the following holds with probability $1 - \exp(-\Omega(\log^2(n+p+q)))$. For any $\eta > 0$,

$$\mathcal{E}^{(i)}(f_{\ell,U^{\text{LORA}},V^{\text{LORA}}}) \leq (1+\eta)(T^{\text{LORA}}_{\text{bias}})^2 + (1+\eta^{-1})(T^{\text{LORA}}_{\text{variance}})^2,$$

where

$$(T_{\text{bias}}^{\text{LoRA}})^2 \le \frac{0 \vee (\text{rank}(D\Sigma_x^{(i)}D^\top) - r)}{\text{rank}(D\Sigma_x^{(i)}D^\top)} \kappa_*^2 (D\Sigma_x^{(i)}D^\top) \mathcal{E}^{(i)}(f^{\text{pre}}) + \mathcal{E}^{(i)}(f_\ell^{\text{full}}),$$

$$(11)$$

$$\begin{split} (T_{\text{variance}}^{\text{LoRA}})^2 &\lesssim C^2 \kappa_*^4(M) \| \Sigma_{\epsilon}^{(i)} \|_{\text{op}} \kappa_*^2(A) \frac{r(r_e(\Phi'^{\top} \Sigma_{\epsilon}^{(i)} \Phi') + r_e(A^2)) \log^2(n+p+q)}{n} \\ &+ C^2 \kappa_*^4(M) \| D \Sigma_x^{(i)} D^{\top} \|_{\text{op}} \frac{r(\kappa_*^2(A) r_e(\Phi'^{\top} D \Sigma_x^{(i)} D^{\top} \Phi') + \kappa_*^6(A) r_e(A^2)) \log^2(n+p+q)}{n} \end{split}$$

Note that the first term on the right hand side of (11) depends on the rank of residual matrix $\Sigma_f^{(\mathrm{i})} = D\Sigma_x^{(\mathrm{i})}D^{\top}$. It becomes zero when $\mathrm{rank}(\Sigma_f^{(\mathrm{i})}) \leq r$ and small when $r/\mathrm{rank}(\Sigma_f^{(\mathrm{i})}) \approx 1$.

Proof of Theorem F.10. Let $\overline{W}_{\ell}^{\text{LoRA}} := \overline{W}_{\ell+1}^{\text{pre}} U^{\text{LoRA}} V^{\text{LoRA}}^{\text{LoRA}}$. From Lemma F.9, we have

$$\mathcal{E}^{(i)}(f_{\ell,U^{\text{LoRA}},V^{\text{LoRA}}}) = \text{tr}\Big((D - \overline{W}_{\ell}^{\text{LoRA}} \underline{W}_{\ell-1}^{\text{pre}}) \Sigma_{x}^{(i)} (D - \overline{W}_{\ell}^{\text{LoRA}} \underline{W}_{\ell-1}^{\text{pre}})^{\top} \Big)$$
$$= \| (\overline{W}_{\ell}^{\text{LoRA}} A A^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} - D) \Sigma_{x}^{(i)1/2} \|_{F}^{2},$$

where we used $(I-AA^\dagger)\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_x^{(\mathrm{i})1/2}=0.$ From Lemma F.9

$$\overline{W}_{\ell}^{\text{LoRA}} A = \text{SVD}_r(\overline{W}_{\ell+1}^{\text{pre}} (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} \hat{D} \hat{\Sigma}_x^{(i)} \underline{W}_{\ell-1}^{\text{pre} \top} \hat{A}^{\dagger}) \hat{A}^{\dagger} A.$$

This gives

$$\begin{split} \|(\overline{W}_{\ell}^{\mathsf{LoRA}}AA^{\dagger}\underline{W}_{\ell-1}^{\mathsf{pre}} - D)\Sigma_{x}^{(\mathsf{i})1/2}\|_{\mathsf{F}} &\leq \|(\overline{W}_{\ell}^{\mathsf{LoRA}}A - \mathsf{SVD}_{r}(\overline{W}_{\ell+1}^{\mathsf{pre}}(\overline{W}_{\ell+1}^{\mathsf{pre}})^{\dagger}D\Sigma_{x}^{(\mathsf{i})}\underline{W}_{\ell-1}^{\mathsf{pre}}^{\mathsf{re}}A^{\dagger}))A^{\dagger}\underline{W}_{\ell-1}^{\mathsf{pre}}\Sigma_{x}^{(\mathsf{i})1/2}\|_{\mathsf{F}} \\ &+ \|\mathsf{SVD}_{r}(\overline{W}_{\ell+1}^{\mathsf{pre}}(\overline{W}_{\ell+1}^{\mathsf{pre}})^{\dagger}D\Sigma_{x}^{(\mathsf{i})}\underline{W}_{\ell-1}^{\mathsf{pre}}^{\mathsf{re}}A^{\dagger})A^{\dagger}\underline{W}_{\ell-1}^{\mathsf{pre}}\Sigma_{x}^{(\mathsf{i})1/2} - D\Sigma_{x}^{(\mathsf{i})1/2}\|_{\mathsf{F}} \\ &=: T_{\mathsf{variance}}^{\mathsf{LoRA}} + T_{\mathsf{bias}}^{\mathsf{LoRA}}. \end{split}$$

We bound $T_{\text{variance}}^{\text{LoRA}}$ and $T_{\text{bias}}^{\text{LoRA}}$ separately.

For the term $T_{\mathrm{variance}}^{\mathrm{LoRA}}$, since $A^{\dagger} \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{i})} \underline{W}_{\ell-1}^{\mathrm{pre}\top} A^{\dagger} = A^{\dagger} A^{2} A^{\dagger}$,

$$T_{\text{variance}}^{\text{LoRA}} = \| \text{SVD}_r(\overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^\dagger \hat{D} \hat{\Sigma}_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}\top} \hat{A}^\dagger) \hat{A}^\dagger A - \text{SVD}_r(\overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^\dagger D \Sigma_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}\top} A^\dagger) A^\dagger A \|_{\text{F}}.$$

Therefore,

$$\begin{split} T_{\text{variance}}^{\text{LoRA}} &\leq \|\text{SVD}_r(\overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^{\dagger}D\Sigma_x^{(\text{i})}\underline{W}_{\ell-1}^{\text{pre}}A^{\dagger})A^{\dagger}A - \text{SVD}_r(\overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^{\dagger}\hat{D}\hat{\Sigma}_x^{(\text{i})}\underline{W}_{\ell-1}^{\text{pre}}\hat{A}^{\dagger})A^{\dagger}A\|_{\text{F}} \\ &+ \|\text{SVD}_r(\overline{W}_{\ell+1}^{\text{pre}}(\overline{W}_{\ell+1}^{\text{pre}})^{\dagger}\hat{D}\hat{\Sigma}_x^{(\text{i})}\underline{W}_{\ell-1}^{\text{pre}}^{\text{re}}\hat{A}^{\dagger})(\hat{A}^{\dagger}A - A^{\dagger}A)\|_{\text{F}} \\ &=: T_{\text{variance},1}^{\text{LoRA}} + T_{\text{variance},2}^{\text{LoRA}}, \end{split}$$

We first bound $T_{\text{variance.}1}^{\text{LoRA}}$. From Lemma G.1 and Assumption F.3, we have

$$\begin{split} T_{\text{variance},1}^{\text{LoRA}} &\leq \|\text{SVD}_r(\hat{M}) - \text{SVD}_r(M)\|_{\text{F}} \\ &\leq \kappa_*^2(M) \frac{\lambda_r(M)}{\lambda_r(M) - \lambda_{r+1}(M)} \sqrt{r} \|\hat{M} - M\|_{\text{op}} \\ &\leq \kappa_*^2(M) C \sqrt{r} \|\hat{M} - M\|_{\text{op}}, \end{split}$$

where $\hat{M} = \overline{W}_{\ell+1}^{\text{pre}} (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} \hat{A}^{\dagger}$ and $M = \overline{W}_{\ell+1}^{\text{pre}} (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} D \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} A^{\dagger}$. From Lemma G.3,

$$\begin{split} \|\hat{M} - M\|_{\text{op}} &\leq \|\Phi'^{\top} \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} - \Phi'^{\top} D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top}\|_{\text{op}} \|\hat{A}^{\dagger}\|_{\text{op}} \\ &+ \|D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top}\|_{\text{op}} \|\hat{A}^{\dagger} - A^{\dagger}\|_{\text{op}} \\ &\lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \kappa_{*}(A) \sqrt{\frac{(r_{e}(\Phi'^{\top} \Sigma_{\epsilon}^{(i)} \Phi') + r_{e}(A^{2})) \log^{2}(n+p+q)}{n}} \\ &+ \|D \Sigma_{x}^{(i)} D^{\top}\|_{\text{op}}^{1/2} \kappa_{*}(A) \sqrt{\frac{(r_{e}(\Phi'^{\top} D \Sigma_{x}^{(i)} D^{\top} \Phi') + r_{e}(A^{2})) \log^{2}(n+p+q)}{n}} \\ &+ \|D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top}\|_{\text{op}} \frac{\kappa_{*}(A)}{\lambda_{*}(A)} \sqrt{\frac{r_{e}(A^{2}) \log^{2}(n+p+q)}{n}} \\ &\lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \kappa_{*}(A) \sqrt{\frac{(r_{e}(\Phi'^{\top} \Sigma_{\epsilon}^{(i)} \Phi') + r_{e}(A^{2})) \log^{2}(n+p+q)}{n}} \\ &+ \|D \Sigma_{x}^{(i)} D^{\top}\|_{\text{op}}^{1/2} \sqrt{\frac{(\kappa_{*}^{2}(A) r_{e}(\Phi'^{\top} D \Sigma_{x}^{(i)} D^{\top} \Phi') + \kappa_{*}^{4}(A) r_{e}(A^{2})) \log^{2}(n+p+q)}{n}} \end{split}$$

holds on the event \mathcal{F} , where we used $\|D\Sigma_x^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}} \leq \|D\Sigma_x^{(\mathrm{i})1/2}\|_{\mathrm{op}}\|A\|_{\mathrm{op}}$. Hence

$$\begin{split} T_{\text{variance},1}^{\text{LORA}} &\lesssim C_{\text{g}} \kappa_*^2(M) \| \Sigma_{\epsilon}^{(\text{i})} \|_{\text{op}}^{1/2} \kappa_*(A) \sqrt{\frac{r(r_e(\Phi'^{\top} \Sigma_{\epsilon}^{(\text{i})} \Phi') + r_e(A^2)) \log^2(n+p+q)}{n}} \\ &+ C_{\text{g}} \kappa_*^2(M) \| D \Sigma_x^{(\text{i})} D^{\top} \|_{\text{op}}^{1/2} \sqrt{\frac{r(\kappa_*^2(A) r_e(\Phi'^{\top} D \Sigma_x^{(\text{i})} D^{\top} \Phi') + \kappa_*^4(A) r_e(A^2)) \log^2(n+p+q)}{n}}. \end{split}$$

Next we bound $T_{\text{variance},2}^{\text{LoRA}}$. Again from Lemma G.3,

$$\begin{split} T_{\text{variance},2}^{\text{LoRA}} & \leq \sqrt{r} \| \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} \|_{\text{op}} \| \hat{A}^{\dagger} \|_{\text{op}} \| \hat{A}^{\dagger} - A^{\dagger} \|_{\text{op}} \| A \|_{\text{op}} \\ & \lesssim \| D \Sigma_{x}^{(i)1/2} \|_{\text{op}} \| \Sigma_{x}^{(i)1/2} \underline{W}_{\ell-1}^{\text{pre}\top} \|_{\text{op}} \frac{\kappa_{*}^{2}(A)}{\lambda_{*}(A)} \sqrt{\frac{r \cdot r_{e}(A^{2}) \log^{2}(n+p+q)}{n}} \\ & = \| D \Sigma_{x}^{(i)1/2} \|_{\text{op}} \kappa_{*}^{3}(A) \sqrt{\frac{r \cdot r_{e}(A^{2}) \log^{2}(n+p+q)}{n}} \end{split}$$

holds on the event \mathcal{F} . Therefore,

$$\begin{split} T_{\text{variance}}^{\text{LoRA}} &\lesssim C_{\text{g}} \kappa_{*}^{2}(M) \| \Sigma_{\epsilon}^{(\text{i})} \|_{\text{op}}^{1/2} \kappa_{*}(A) \sqrt{\frac{r(r_{e}(\Phi'^{\top} \Sigma_{\epsilon}^{(\text{i})} \Phi') + r_{e}(A^{2})) \log^{2}(n+p+q)}{n}} \\ &+ C_{\text{g}} \kappa_{*}^{2}(M) \| D \Sigma_{x}^{(\text{i})} D^{\top} \|_{\text{op}}^{1/2} \sqrt{\frac{r(\kappa_{*}^{2}(A) r_{e}(\Phi'^{\top} D \Sigma_{x}^{(\text{i})} D^{\top} \Phi') + \kappa_{*}^{6}(A) r_{e}(A^{2})) \log^{2}(n+p+q)}{n}} \end{split}$$

$$(12)$$

hold with high probability.

Bound $T_{\text{bias}}^{\text{LoRA}}$. Note that

$$\begin{split} (T_{\text{bias}}^{\text{LoRA}})^2 &= \| \text{SVD}_r(M) A^\dagger \underline{W}_{\ell-1}^{\text{pre}} \Sigma_x^{(\text{i})1/2} - D \Sigma_x^{(\text{i})1/2} \|_{\text{F}}^2 \\ &= \| \underbrace{\text{SVD}_r(M) A^\dagger \underline{W}_{\ell-1}^{\text{pre}} \Sigma_x^{(\text{i})1/2} - \Phi' \Phi'^\top D \Sigma_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}} (A^2)^\dagger \underline{W}_{\ell-1}^{\text{pre}} \Sigma_x^{(\text{i})1/2}} \|_{\text{F}}^2 \\ &+ \| \underbrace{D \Sigma_x^{(\text{i})1/2} (I - \Sigma_x^{(\text{i})1/2} \underline{W}_{\ell-1}^{\text{pre}} (A^2)^\dagger \underline{W}_{\ell-1}^{\text{pre}} \Sigma_x^{(\text{i})1/2})}_{=:T_2} \|_{\text{F}}^2 \\ &+ \| \underbrace{(I - \Phi' \Phi'^\top) D \Sigma_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}} (A^2)^\dagger \underline{W}_{\ell-1}^{\text{pre}} \Sigma_x^{(\text{i})1/2}}_{=:T_3} \|_{\text{F}}^2 \end{split}$$

where the second equality follows from the fact that cross terms are zero, i.e., $\operatorname{tr}\big(T_1T_2^\top\big) = \operatorname{tr}\big(T_2T_3^\top\big) = \operatorname{tr}\big(T_3T_1^\top\big) = 0 \text{ since } \Psi_*\big(\underline{W}_{\ell-1}^{\operatorname{pre}}\Sigma_x^{(\mathrm{i})1/2}\big)\Psi_*^\top\big(\underline{W}_{\ell-1}^{\operatorname{pre}}\Sigma_x^{(\mathrm{i})1/2}\big) = \Sigma_x^{(\mathrm{i})1/2}\underline{W}_{\ell-1}^{\operatorname{pre}}(A^2)^\dagger\underline{W}_{\ell-1}^{\operatorname{pre}}\Sigma_x^{(\mathrm{i})1/2} \text{ and }$

$$(I - \Phi' \Phi'^\top) \Phi_*(\mathrm{SVD}_r(M)) = 0, \ \ \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_x^{(\mathrm{i})1/2} (I - \Psi_*(\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_x^{(\mathrm{i})1/2}) \Psi_*^\top(\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_x^{(\mathrm{i})1/2})) = 0$$
 hold. Thus from Lemma F.17,

$$(T_{\text{bias}}^{\text{LoRA}})^2 = \|\text{SVD}_r(\Phi'\Phi'^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger}) - \Phi'\Phi'^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger}\|_F^2 + \mathcal{E}^{(i)}(f_{\ell}^{\text{full}}). \tag{13}$$

Notice that

$$\|\operatorname{SVD}_{r}(\Phi'\Phi'^{\top}D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}\top}A^{\dagger}) - \Phi'\Phi'^{\top}D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}\top}A^{\dagger}\|_{F}^{2}$$

$$\leq \{0 \vee (\operatorname{rank}(\Phi'\Phi'^{\top}D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}\top}A^{\dagger}) - r)\}\|\Phi'\Phi'^{\top}D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}\top}A^{\dagger}\|_{\operatorname{op}}^{2}$$

$$\leq \{0 \vee (\operatorname{rank}(\Phi'\Phi'^{\top}D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}\top}A^{\dagger}) - r)\}\|D\Sigma_{x}^{(i)1/2}\|_{\operatorname{op}}^{2}$$

$$\leq \frac{0 \vee (\operatorname{rank}(D\Sigma_{x}^{(i)}D^{\top}) - r)}{\operatorname{rank}(D\Sigma_{x}^{(i)1/2})}\kappa_{*}^{2}(D\Sigma_{x}^{(i)}D^{\top})\mathcal{E}^{(i)}(f^{\operatorname{pre}}), \tag{14}$$

where the last inequality follows since

$$\|D\Sigma_x^{(i)1/2}\|_{\mathsf{F}}^2 = \|\Lambda_*(D\Sigma_x^{(i)1/2})\|_{\mathsf{F}}^2 \ge \operatorname{rank}(D\Sigma_x^{(i)1/2})\lambda_*^2(D\Sigma_x^{(i)1/2}) = \frac{\operatorname{rank}(D\Sigma_x^{(i)1/2})}{\kappa_*^2(D\Sigma_x^{(i)1/2})} \|D\Sigma_x^{(i)1/2}\|_{\mathsf{op}}^2.$$

Summary Note that for any $\eta > 0$, $(T_{\text{variance}}^{\text{LoRA}} + T_{\text{bias}}^{\text{LoRA}})^2 \le (1+\eta)(T_{\text{bias}}^{\text{LoRA}})^2 + (1+1/\eta)(T_{\text{variance}}^{\text{LoRA}})^2$ holds. Therefore,

$$\mathcal{E}^{(\mathrm{i})}(f_{\ell,U^{\mathrm{LoRA}},V^{\mathrm{LoRA}}}) \leq (1+\eta)(T^{\mathrm{LoRA}}_{\mathrm{bias}})^2 + (1+\eta^{-1})(T^{\mathrm{LoRA}}_{\mathrm{variance}})^2.$$

Combined with (12), (13), and (14), this concludes the proof.

F.5.3 Out-of-distribution Excess Risk of LoRA

We define the low-rank matrix obtained by LoRA under population in-distribution risk as

$$(U_{\infty}^{\text{LoRA}}, V_{\infty}^{\text{LoRA}}) \in \operatorname*{arg\,min}_{U,V} \|(U, V)\|_{\text{F}}^{2} \quad \text{s.t. } (U, V) \text{ minimizes } \mathcal{R}^{(i)}(f_{\ell, U, V}). \tag{15}$$

Theorem F.11 (Restatement of Theorem F.8: LoRA Part). For $(U_{\infty}^{\text{LoRA}}, V_{\infty}^{\text{LoRA}})$, defined in (15)

$$\begin{split} \mathcal{E}^{(\mathrm{o})}(f_{\ell,U_{\infty}^{\mathrm{LoRA}},V_{\infty}^{\mathrm{LoRA}}}) \lesssim & \| (I - \Phi'\Phi'^{\top})B^{(\mathrm{o})}\Sigma_{x}^{(\mathrm{o})1/2} \|_{\mathrm{F}}^{2} + \| (B^{(\mathrm{o})} - B^{(\mathrm{i})})\Sigma_{x}^{(\mathrm{i})1/2} \|_{\mathrm{F}}^{2} \|G_{\ell-1}^{(\mathrm{i},\mathrm{o})}\|_{\mathrm{op}}^{2} \\ & + \| (B^{(\mathrm{o})} - W^{\mathrm{pre}})(\Sigma_{x}^{(\mathrm{o})1/2} - \Sigma_{x}^{(\mathrm{i})1/2}G_{\ell-1}^{(\mathrm{i},\mathrm{o})}) \|_{\mathrm{F}} \\ & + \frac{0 \vee (\mathrm{rank}(D\Sigma_{x}^{(\mathrm{i})}D^{\top}) - r)}{\mathrm{rank}(D\Sigma_{x}^{(\mathrm{i})}D^{\top})} \kappa_{*}^{2} (D\Sigma_{x}^{(\mathrm{i})}D^{\top}) \|G_{\ell-1}^{(\mathrm{i},\mathrm{o})}\|_{\mathrm{op}}^{2} \mathcal{E}^{(\mathrm{i})}(f^{\mathrm{pre}}). \end{split}$$

Furthermore, for any $\eta \in (0, 1)$,

$$\mathcal{E}^{(o)}(f_{\ell,U_{\infty}^{\text{LoRA}},V_{\infty}^{\text{LoRA}}}) \geq (1-\eta) \left\| (B^{(o)} - B^{(i)}) \Sigma_{x}^{(o)1/2} \right\|_{\text{F}}^{2} - 3(\eta^{-1} - 1) \| (I - \Phi'\Phi'^{\top}) B^{(i)} \Sigma_{x}^{(o)1/2} \|_{\text{F}}^{2}$$

$$- 3(\eta^{-1} - 1) \| (B^{(i)} - W^{\text{pre}}) (\Sigma_{x}^{(o)1/2} - \Sigma_{x}^{(i)1/2} G_{\ell-1}^{(i,o)}) \|_{\text{F}}^{2}$$

$$- 3(\eta^{-1} - 1) \frac{0 \vee (\operatorname{rank}(D\Sigma_{x}^{(i)}D^{\top}) - r)}{\operatorname{rank}(D\Sigma_{x}^{(i)}D)} \kappa_{*}^{2} (D\Sigma_{x}^{(i)}D^{\top}) \| G_{\ell-1}^{(i,o)} \|_{\text{op}} \mathcal{E}^{(i)}(f^{\text{pre}}).$$

$$(16)$$

Proof of Theorem F.11. With a slight modification to the proof of Lemma F.9, it follows that

$$\mathcal{E}^{(\mathrm{o})}(f_{\ell,U_{\infty}^{\mathrm{LoRA}},V_{\infty}^{\mathrm{LoRA}}}) = \mathrm{tr}\bigg(\Big(B^{(\mathrm{o})} - W^{\mathrm{pre}} - \mathrm{SVD}_{r}(\overline{W}_{\ell+1}^{\mathrm{pre}}(\overline{W}_{\ell+1}^{\mathrm{pre}})^{\dagger}D\Sigma_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger})A^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Big)\Sigma_{x}^{(\mathrm{o})} \\ \cdot \Big(B^{(\mathrm{o})} - W^{\mathrm{pre}} - \mathrm{SVD}_{r}(\overline{W}_{\ell+1}^{\mathrm{pre}}(\overline{W}_{\ell+1}^{\mathrm{pre}})^{\dagger}D\Sigma_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger})A^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Big)^{\top}\bigg) \\ = \bigg\| (B^{(\mathrm{o})} - W^{\mathrm{pre}})\Sigma_{x}^{(\mathrm{o})1/2} - \mathrm{SVD}_{r}(\Phi'\Phi'^{\top}D\Sigma_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger})A^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_{x}^{(\mathrm{o})1/2} \bigg\|_{\mathrm{F}}^{2}. \tag{17}$$

Recall that $M:=\Phi'\Phi'^{\top}D\Sigma_x^{(\mathbf{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger}.$ Then,

$$\begin{split} & \left\| (B^{(o)} - W^{\text{pre}}) \Sigma_{x}^{(o)1/2} - \text{SVD}_{r} (\Phi' \Phi'^{\top} D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} \mathbf{A}^{\dagger}) A^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} \right\|_{\text{F}} \\ & \leq \left\| (B^{(o)} - W^{\text{pre}}) \Sigma_{x}^{(o)1/2} - \Phi' \Phi'^{\top} D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} \mathbf{A}^{\dagger}) A^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} \right\|_{\text{F}} \\ & + \left\| M A^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} - \Phi' \Phi'^{\top} D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} \right\|_{\text{F}} \\ & = \left\| (B^{(o)} - W^{\text{pre}}) \Sigma_{x}^{(o)1/2} - \Phi' \Phi'^{\top} D \Sigma_{x}^{(i)1/2} (\underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(i)1/2})^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} \right\|_{\text{F}} \\ & + \left\| M A^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} - \text{SVD}_{r} (M) A^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} \right\|_{\text{F}} \\ & \leq \left\| (I - \Phi' \Phi'^{\top}) B^{(o)} \Sigma_{x}^{(o)1/2} \right\|_{\text{F}} + \left\| \Phi' \Phi'^{\top} (B^{(o)} - B^{(i)}) \Sigma_{x}^{(i)1/2} G_{\ell-1}^{(i,o)} \right\|_{\text{F}} \\ & + \left\| \Phi' \Phi'^{\top} (B^{(o)} - W^{\text{pre}}) (\Sigma_{x}^{(o)1/2} - \Sigma_{x}^{(i)1/2} G_{\ell-1}^{(i,o)}) \right\|_{\text{F}} \\ & \leq \left\| (I - \Phi' \Phi'^{\top}) B^{(o)} \Sigma_{x}^{(o)1/2} \right\|_{\text{F}} + \left\| (B^{(o)} - B^{(i)}) \Sigma_{x}^{(i)1/2} \right\|_{\text{F}} \|G_{\ell-1}^{(i,o)} \|_{\text{op}} \\ & + \left\| (B^{(o)} - W^{\text{pre}}) (\Sigma_{x}^{(o)1/2} - \Sigma_{x}^{(i)1/2} G_{\ell-1}^{(i,o)}) \right\|_{\text{F}} + \left\| M - \text{SVD}_{r} (M) \right\|_{\text{F}} \|A^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(o)1/2} \|_{\text{op}}, \end{aligned}$$

where we used $\Phi'\Phi'^{\top}W^{\text{pre}} = W^{\text{pre}}$. From (14), we have

$$\begin{split} \{\mathcal{E}^{(\mathrm{o})}(f_{\ell,U_{\infty}^{\mathrm{LoRA}},V_{\infty}^{\mathrm{LoRA}}})\}^{1/2} &\leq \|(I - \Phi'\Phi'^{\top})B^{(\mathrm{o})}\Sigma_{x}^{(\mathrm{o})1/2}\|_{\mathrm{F}} + \|(B^{(\mathrm{o})} - B^{(\mathrm{i})})\Sigma_{x}^{(\mathrm{i})1/2}\|_{\mathrm{F}} \|G_{\ell-1}^{(\mathrm{i},\mathrm{o})}\|_{\mathrm{op}} \\ &+ \|(B^{(\mathrm{o})} - W^{\mathrm{pre}})(\Sigma_{x}^{(\mathrm{o})1/2} - \Sigma_{x}^{(\mathrm{i})1/2}G_{\ell-1}^{(\mathrm{i},\mathrm{o})})\|_{\mathrm{F}} \\ &+ \|G_{\ell-1}^{(\mathrm{i},\mathrm{o})}\|_{\mathrm{op}} \kappa_{*}(D\Sigma_{x}^{(\mathrm{i})}D^{\top})\sqrt{\frac{0 \vee (\mathrm{rank}(D\Sigma_{x}^{(\mathrm{i})}D^{\top}) - r)}{\mathrm{rank}(D\Sigma_{x}^{(\mathrm{i})1/2})}}\mathcal{E}^{(\mathrm{i})}(f^{\mathrm{pre}}), \end{split}$$

where we used $\|A^\dagger \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_x^{(\mathrm{o})1/2}\|_{\mathrm{op}} = \|G_{\ell-1}^{(\mathrm{i},\mathrm{o})}\|_{\mathrm{op}}$. This gives the first claim.

Using $2\operatorname{tr}(AB^{\top}) \ge -\eta \|A\|_{\mathrm{F}}^2 - (1/\eta)\|B\|_{\mathrm{F}}^2$ for any $\eta > 0$ and any matrices A, B of the same shape, (17) can be rewritten as

$$\mathcal{E}^{(o)}(f_{\ell,U_{\infty}^{LoRA},V_{\infty}^{LoRA}}) = \left\| (B^{(o)} - B^{(i)}) \Sigma_{x}^{(o)1/2} + \underbrace{(I - \Phi' \Phi'^{\top})(B^{(i)} - W^{pre}) \Sigma_{x}^{(o)1/2}}_{=:T_{1}} + \underbrace{\Phi' \Phi'^{\top}(B^{(i)} - W^{pre})(\Sigma_{x}^{(o)1/2} - \Sigma_{x}^{(i)1/2} G_{\ell-1}^{(i,o)})}_{=:T_{2}} + \underbrace{MA^{\dagger} \underline{W}_{\ell-1}^{pre} \Sigma_{x}^{(o)1/2} - SVD_{r}(M)A^{\dagger} \underline{W}_{\ell-1}^{pre} \Sigma_{x}^{(o)1/2}}_{=:T_{3}} \right\|_{F}^{2}$$

$$= \left\| (B^{(o)} - B^{(i)}) \Sigma_{x}^{(o)1/2} \right\|_{F}^{2} + 2 \operatorname{tr} \left((B^{(o)} - B^{(i)}) \Sigma_{x}^{(o)1/2} (T_{1} + T_{2} + T_{3})^{\top} \right) + \|T_{1} + T_{2} + T_{3}\|_{F}^{2}$$

$$\geq (1 - \eta) \| (B^{(o)} - B^{(i)}) \Sigma_{x}^{(o)1/2} \right\|_{F}^{2} + (1 - \eta^{-1}) \|T_{1} + T_{2} + T_{3}\|_{F}^{2}. \tag{18}$$

Choose $\eta \in (0,1)$. By a similar argument as above, and using $\Phi' \Phi'^{\top} W^{\text{pre}} = W^{\text{pre}}$, we can show that $\|T_1 + T_2 + T_3\|_F^2 \le 3\|T_1\|_F^2 + 3\|T_2\|_F^2 + 3\|T_3\|_F^2$

$$\leq 3 \| (I - \Phi' \Phi'^{\top}) B^{(\mathbf{i})} \Sigma_{x}^{(\mathbf{o})1/2} \|_{\mathrm{F}}^{2} + 3 \| (B^{(\mathbf{i})} - W^{\mathrm{pre}}) (\Sigma_{x}^{(\mathbf{o})1/2} - \Sigma_{x}^{(\mathbf{i})1/2} G_{\ell-1}^{(\mathbf{i},\mathbf{o})}) \|_{\mathrm{F}}^{2}$$

$$+ 3 \frac{0 \vee (\mathrm{rank}(D\Sigma_{x}^{(\mathbf{i})}D^{\top}) - r)}{\mathrm{rank}(D\Sigma_{x}^{(\mathbf{i})1/2})} \kappa_{*}^{2} (D\Sigma_{x}^{(\mathbf{i})}D^{\top}) \| G_{\ell-1}^{(\mathbf{i},\mathbf{o})} \|_{\mathrm{op}} \mathcal{E}^{(\mathbf{i})} (f^{\mathrm{pre}}),$$

where we used (14) again. This concludes the proof.

F.6 Proofs for Structured Sparse Fine-tuning

F.6.1 Excess Risk of Structured Sparse Fine-tuning

Lemma F.12 (Excess Risk). Given $S \subset [d_{\ell}]$, consider the minimum norm solution

$$V^{\mathrm{S}^2\mathrm{FT}} \in \mathop{\arg\min}_{V \in \mathbb{R}^{d_{\ell-1} \times s}} \|V\|_{\mathrm{F}}^2 \quad \textit{s.t. } V \textit{ minimizes } \mathcal{R}_n^{(\mathrm{i})}(f_{\ell,U_S^{\mathrm{S}^2\mathrm{FT}},V}).$$

Then, the structured sparse adaptation matrix satisfies

$$U_S^{\mathsf{S}^2\mathsf{FT}}V^{\mathsf{S}^2\mathsf{FT}\top} = U_S^{\mathsf{S}^2\mathsf{FT}}(\overline{W}_{\ell+1}^{\mathsf{pre}}U_S^{\mathsf{S}^2\mathsf{FT}})^{\dagger}\hat{D}\hat{\Sigma}_x^{(\mathsf{i})}\underline{W}_{\ell-1}^{\mathsf{pre}\top}(\hat{A}^{\dagger})^2, \tag{19}$$

and

$$\begin{split} \mathcal{E}^{(k)}(f_{\ell,U_S^{\mathsf{S}^2\mathsf{FT}},V^{\mathsf{S}^2\mathsf{FT}}}) &= \operatorname{tr}\bigg(\Big(B^{(k)} - W^{\mathsf{pre}} - \overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}} (\overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}})^\dagger \hat{D} \hat{\Sigma}_x^{(\mathbf{i})} \underline{W}_{\ell-1}^{\mathsf{pre}\top} (\hat{A}^\dagger)^2 \underline{W}_{\ell-1}^{\mathsf{pre}} \Big) \Sigma_x^{(k)} \\ & \cdot \Big(B^{(k)} - W^{\mathsf{pre}} - \overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}} (\overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}})^\dagger \hat{D} \hat{\Sigma}_x^{(\mathbf{i})} \underline{W}_{\ell-1}^{\mathsf{pre}\top} (\hat{A}^\dagger)^2 \underline{W}_{\ell-1}^{\mathsf{pre}} \Big)^\top \bigg) \end{split}$$

for $k \in \{i, o\}$.

Proof. Since $\hat{\Sigma}_{x,\epsilon}^{(k)} = (1/n)X^{(k)}E^{(k)\top}$ and $\hat{\Sigma}_{x}^{(k)} = (1/n)X^{(k)}X^{(k)\top}$, we have $\hat{\Sigma}_{x,\epsilon}^{(k)} = \hat{\Sigma}_{x}^{(k)}(X^{(k)\top})^{\dagger}E^{(k)\top} =: \hat{\Sigma}_{x}^{(k)}\hat{\Sigma}_{x,\epsilon}^{(k)}$. Similar to (9), we have

$$\begin{split} \mathcal{R}_n^{(\mathrm{i})}(f_{\ell,U_S^{\mathrm{S}^2\mathrm{FT}},V}) &= \|\overline{W}_{\ell+1}^{\mathrm{pre}} U_S^{\mathrm{S}^2\mathrm{FT}} V^\top \hat{A} - \hat{D} \hat{\Sigma}_x^{(\mathrm{i})} \underline{W}_{\ell-1}^{\mathrm{pre}\top} \hat{A}^\dagger \|_{\mathrm{F}}^2 - \|\hat{D} \hat{\Sigma}_x^{(\mathrm{i})} \underline{W}_{\ell-1}^{\mathrm{pre}\top} \hat{A}^\dagger \|_{\mathrm{F}}^2 \\ &+ \mathrm{tr} \Big(D \hat{\Sigma}_x^{(\mathrm{i})} D^\top \Big) + 2 \, \mathrm{tr} \Big(D \hat{\Sigma}_{x,\epsilon}^{(\mathrm{i})} \Big) + \mathrm{tr} \Big(\hat{\Sigma}_{\epsilon}^{(\mathrm{i})} \Big). \end{split}$$

Thus minimizing $\mathcal{R}_n^{(\mathrm{i})}(f_{\ell,U_s^{\mathrm{S}^2\mathrm{FT}},V})$ is equivalent to minimizing the norm

$$\begin{split} \|\overline{W}_{\ell+1}^{\text{pre}} U_{S}^{\text{S}^{2}\text{FT}} V^{\top} \hat{A} - \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} \hat{A}^{\dagger} \|_{\text{F}}^{2} \\ &= \|\overline{W}_{\ell+1}^{\text{pre}} U_{S}^{\text{S}^{2}\text{FT}} V^{\top} \hat{A} - \overline{W}_{\ell+1}^{\text{pre}} U_{S}^{\text{S}^{2}\text{FT}} (\overline{W}_{\ell+1}^{\text{pre}} U_{S}^{\text{S}^{2}\text{FT}})^{\dagger} \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} \hat{A}^{\dagger} \|_{\text{F}}^{2} \\ &+ \| (I - (\overline{W}_{\ell+1}^{\text{pre}} U_{S}^{\text{S}^{2}\text{FT}}) (\overline{W}_{\ell+1}^{\text{pre}} U_{S}^{\text{S}^{2}\text{FT}})^{\dagger}) \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} \hat{A}^{\dagger} \|_{\text{F}}^{2}. \end{split}$$

Using the same argument as in the proof of Lemma F.9, the minimum norm solution V^{S^2FT} is obtained by

$$V^{\mathrm{S}^2\mathrm{FT}} = (\hat{A}^\dagger)^2 \underline{W}_{\ell-1}^{\mathrm{pre}} \hat{\Sigma}_x^{(\mathrm{i})} \hat{D}^\top (U_S^{\mathrm{S}^2\mathrm{FT}} \overline{W}_{\ell+1}^{\mathrm{pre}})^\dagger.$$

The excess risk for $k \in \{i, o\}$ becomes

$$\begin{split} \mathcal{E}^{(k)}(f_{\ell,U_S^{\mathsf{S}^2\mathsf{FT}},V^{\mathsf{S}^2\mathsf{FT}}}) &= \mathbb{E}\bigg[\Big(B^{(k)} x^{(k)} - \overline{W}_{\ell+1}^{\mathsf{pre}} (W_\ell^{\mathsf{pre}} + U_S^{\mathsf{S}^2\mathsf{FT}} V^{\mathsf{S}^2\mathsf{FT}}) \underline{W}_{\ell-1}^{\mathsf{pre}} x^{(k)} \Big)^2 \bigg] \\ &= \mathrm{tr}\bigg(\Big(B^{(k)} - W^{\mathsf{pre}} - \overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}} (\overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}})^\dagger \hat{D} \hat{\Sigma}_x^{(i)} \underline{W}_{\ell-1}^{\mathsf{pre}} (\hat{A}^\dagger)^2 \underline{W}_{\ell-1}^{\mathsf{pre}} \Big) \Sigma_x^{(k)} \\ & \cdot \Big(B^{(k)} - W^{\mathsf{pre}} - \overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}} (\overline{W}_{\ell+1}^{\mathsf{pre}} U_S^{\mathsf{S}^2\mathsf{FT}})^\dagger \hat{D} \hat{\Sigma}_x^{(i)} \underline{W}_{\ell-1}^{\mathsf{pre}} (\hat{A}^\dagger)^2 \underline{W}_{\ell-1}^{\mathsf{pre}} \Big)^\top \bigg). \end{split}$$

This concludes the proof.

F.6.2 In-distribution Excess Risk of Structured Sparse Fine-tuning

Theorem F.13 (Restatement of Theorem F.7: S²FT Part). Suppose that Assumptions F.1 and F.2 hold. Fix $S \subset [d_\ell]$ with |S| = s. Then, the following holds with probability $1 - \exp(-\Omega(\log^2(n+p+q)))$. For any $\eta > 0$,

$$\mathcal{E}^{(\mathrm{i})}(f_{\ell,U_{\mathrm{G}}^{\mathrm{S}^2\mathrm{FT}},V^{\mathrm{S}^2\mathrm{FT}}}) \leq (1+\eta)(T_{\mathrm{bias}}^{\mathrm{S}^2\mathrm{FT}})^2 + (1+\eta^{-1})(T_{\mathrm{variance}}^{\mathrm{S}^2\mathrm{FT}})^2,$$

where

$$(T_{\text{bias}}^{\text{S}^2\text{FT}})^2 \leq \|(\Phi'\Phi'^{\top} - \Phi''_S \Phi''_S^{\top})\Phi_*(D\Sigma_x^{(i)1/2})\|_{\text{op}}^2 \mathcal{E}^{(i)}(f_{\ell}^{\text{pre}}) + \mathcal{E}^{(i)}(f_{\ell}^{\text{full}}), \tag{21}$$

$$(T_{\text{variance}}^{\text{S}^2\text{FT}})^2 \lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}} \kappa_*^2 (A) \frac{s(r_e(\Phi''_S^{\top}\Sigma_{\epsilon}^{(i)}\Phi''_S) + r_e(A^2))\log^2(n+p+q)}{n} + \|D\Sigma_x^{(i)}D^{\top}\|_{\text{op}} \frac{s(\kappa_*^2(A)r_e(\Phi''_S^{\top}D\Sigma_x^{(i)}D^{\top}\Phi''_S) + \kappa_*^8(A)r_e(A^2))\log^2(n+p+q)}{n}.$$

Note that the term $\|(\Phi'\Phi'^{\top} - \Phi''_S\Phi''^{\top})\Phi_*(D\Sigma_x^{(i)1/2})\|_{op}$ in (21) measures the distance between subspaces spanned by Φ' and Φ''_S in a label space, weighted by $\Phi_*(\Sigma_f^{(i)})$. In high level, this quantity shows the closeness between the ℓ -th layer full fine-tuning and S^2FT . It takes small values when the important channels for residual prediction are sparsely distributed among all channels. This aligns with the intuition that S^2FT only selectively fine-tunes small number of coordinates, and thus relying on the information contained in those coordinates.

Proof of Theorem F.13. Using the same argument as in the proof of Theorem F.10 combined with Lemma F.12, we have

$$\mathcal{E}^{(\mathrm{i})}(f_{\ell,U_S^{\mathrm{2FT}},V^{\mathrm{S}^2\mathrm{FT}}}) = \|(\overline{W}_{\ell+1}^{\mathrm{pre}}U_S^{\mathrm{2FT}}V^{\mathrm{S}^2\mathrm{FT}\top}AA^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}} - D)\Sigma_x^{(\mathrm{i})1/2}\|_{\mathrm{F}}^2,$$

and

$$\begin{split} &\|(\overline{W}_{\ell+1}^{\mathrm{pre}}U_{S}^{\mathrm{S}^{2}\mathrm{FT}}V^{\mathrm{S}^{2}\mathrm{FT}\top}AA^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}-D)\Sigma_{x}^{(\mathrm{i})1/2}\|_{\mathrm{F}} \\ &\leq \|\overline{W}_{\ell+1}^{\mathrm{pre}}U_{S}^{\mathrm{S}^{2}\mathrm{FT}}(\overline{W}_{\ell+1}^{\mathrm{pre}}U_{S}^{\mathrm{S}^{2}\mathrm{FT}})^{\dagger}(\hat{D}\hat{\Sigma}_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(\hat{A}^{2})^{\dagger}-D\Sigma_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(A^{2})^{\dagger})A\|_{\mathrm{F}} \\ &+\|\overline{W}_{\ell+1}^{\mathrm{pre}}U_{S}^{\mathrm{S}^{2}\mathrm{FT}}(\overline{W}_{\ell+1}^{\mathrm{pre}}U_{S}^{\mathrm{S}^{2}\mathrm{FT}})^{\dagger}D\Sigma_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(A^{2})^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_{x}^{(\mathrm{i})1/2}-D\Sigma_{x}^{(\mathrm{i})1/2}\|_{\mathrm{F}} \\ &=:T_{\mathrm{variance}}^{\mathrm{S}^{2}\mathrm{FT}}+T_{\mathrm{bias}}^{\mathrm{S}^{2}\mathrm{FT}}. \end{split}$$

We bound $T_{\text{variance}}^{S^2\text{FT}}$ and $T_{\text{bias}}^{S^2\text{FT}}$ separately.

Bound $T_{\text{variance}}^{S^2\text{FT}}$. Note that

$$\begin{split} T_{\text{variance}}^{\text{S}^2\text{FT}} &= \| \overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}} (\overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}})^\dagger \hat{D} \hat{\Sigma}_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}} (\hat{A}^\dagger)^2 A - \overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}} (\overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}})^\dagger D \Sigma_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}} A^\dagger \|_{\text{F}} \\ &\leq \| \overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}} (\overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}})^\dagger D \Sigma_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}} A^\dagger - \overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}} (\overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}})^\dagger \hat{D} \hat{\Sigma}_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}} A^\dagger \|_{\text{F}} \\ &+ \| \overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}} (\overline{W}_{\ell+1}^{\text{pre}} U_S^{\text{S}^2\text{FT}})^\dagger \hat{D} \hat{\Sigma}_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}} ((\hat{A}^\dagger)^2 - (A^\dagger)^2) A \|_{\text{F}} \\ &=: T_{\text{variance},1}^{\text{S}^2\text{FT}} + T_{\text{variance},2}^{\text{S}^2\text{FT}}. \end{split}$$

For the term $T_{\text{variance},1}^{\text{S}^2\text{FT}}$, using Lemma G.3,

$$\begin{split} T_{\text{variance},1}^{\text{S}^2\text{FT}} &\leq 2\sqrt{s} \| \Phi_S''^\top D \Sigma_x^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} - \Phi_S''^\top \hat{D} \hat{\Sigma}_x^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} \|_{\text{op}} \|A^\dagger\|_{\text{op}} \\ &\lesssim \| \Sigma_\epsilon^{(i)} \|_{\text{op}}^{1/2} \kappa_*(A) \sqrt{\frac{s(r_e(\Phi_S''^\top \Sigma_\epsilon^{(i)} \Phi_S'') + r_e(A^2)) \log^2(n+p+q)}{n}} \\ &+ \| D \Sigma_x^{(i)} D^\top \|_{\text{op}}^{1/2} \kappa_*(A) \sqrt{\frac{s(r_e(\Phi_S''^\top D \Sigma_x^{(i)} D^\top \Phi_S'') + r_e(A^2)) \log^2(n+p+q)}{n}} \end{split}$$

holds on the event \mathcal{F} , where the first inequality follows since the term inside the norm is at most rank-2s. Again from Lemma G.3,

$$\begin{split} T_{\text{variance},2}^{\text{S}^2\text{FT}} &\leq \sqrt{s} \| \Phi_S''^\top \hat{D} \hat{\Sigma}_x^{(\text{i})} \underline{W}_{\ell-1}^{\text{pre}\top} \|_{\text{op}} \| (\hat{A}^\dagger)^2 - (A^\dagger)^2 \|_{\text{op}} \| A \|_{\text{op}} \\ &\lesssim \| D \Sigma_x^{(\text{i})1/2} \|_{\text{op}} \| \Sigma_x^{(\text{i})1/2} \underline{W}_{\ell-1}^{\text{pre}\top} \|_{\text{op}} \frac{\kappa_*^3(A)}{\lambda_*(A)} \sqrt{\frac{s r_e(A^2) \log^2(n+d+p)}{n}} \\ &= \| D \Sigma_x^{(\text{i})1/2} \|_{\text{op}} \kappa_*^4(A) \sqrt{\frac{s r_e(A^2) \log^2(n+d+p)}{n}} \end{split}$$

holds on the event \mathcal{F} . Therefore,

$$T_{\text{variance}}^{S^2 \text{FT}} \lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \kappa_*(A) \sqrt{\frac{s(r_e(\Phi_S''^{\top} \Sigma_{\epsilon}^{(i)} \Phi_S'') + r_e(A^2)) \log^2(n+p+q)}{n}} + \|D\Sigma_x^{(i)} D^{\top}\|_{\text{op}}^{1/2} \sqrt{\frac{s(\kappa_*^2(A) r_e(\Phi_S''^{\top} D \Sigma_x^{(i)} D^{\top} \Phi_S'') + \kappa_*^8(A) r_e(A^2)) \log^2(n+p+q)}{n}}.$$
(22)

Bound $T_{\text{bias}}^{\text{S}^2\text{FT}}$. By the same argument as in the proof of Theorem F.10,

$$(T_{\text{bias}}^{S^{2}\text{FT}})^{2} = \|\Phi_{S}''\Phi_{S}''^{\top}D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger} - \Phi'\Phi'^{\top}D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger}\|_{F}^{2} + \mathcal{E}^{(i)}(f_{\ell}^{\text{full}})$$

$$\leq \|(\Phi_{S}''\Phi_{S}''^{\top} - \Phi'\Phi'^{\top})\Phi_{*}(D\Sigma_{x}^{(i)1/2})\|_{\text{op}}^{2}\|D\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\text{pre}\top}A^{\dagger}\|_{F}^{2} + \mathcal{E}^{(i)}(f_{\ell}^{\text{full}})$$

$$= \|(\Phi_{S}''\Phi_{S}''^{\top} - \Phi'\Phi'^{\top})\Phi_{*}(D\Sigma_{x}^{(i)1/2})\|_{\text{op}}^{2}\mathcal{E}^{(i)}(f_{\ell}^{\text{pre}}) + \mathcal{E}^{(i)}(f_{\ell}^{\text{full}}),$$
(24)

where we used $\|D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}^{\top}}A^{\dagger}\|_{\mathrm{F}}^2 \leq \|D\Sigma_x^{(i)1/2}\|_{\mathrm{F}}^2 = \mathcal{E}^{(i)}(f_{\ell}^{\mathrm{pre}})$. We hypothesize that $T_{\mathrm{bias}}^{\mathrm{S}^2\mathrm{FT}} \simeq T_{\mathrm{bias}}^{\mathrm{LoRA}}$ by comparing (13) and (23), Here, $\mathrm{SVD}_s(\Phi'\Phi'^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}^{\top}}A^{\dagger})$ is the best rank-s approximation of $\Phi'\Phi'^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}^{\top}}A^{\dagger}$ and $\Phi''_S\Phi''^{\top}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}^{\top}}A^{\dagger}$ benefits from a rank-r approximation, where r>s.

Summary Note that for any $\eta > 0$, $(T_{\text{variance}}^{\text{S}^2\text{FT}} + T_{\text{bias}}^{\text{S}^2\text{FT}})^2 \le (1+\eta)(T_{\text{bias}}^{\text{S}^2\text{FT}})^2 + (1+1/\eta)(T_{\text{variance}}^{\text{S}^2\text{FT}})^2$ holds. Thus

$$\mathcal{E}^{(\mathrm{i})}(f_{\ell,U_{\mathrm{c}}^{\mathrm{S}^2\mathrm{FT}},V^{\mathrm{S}^2\mathrm{FT}}}) \leq (1+\eta)(T_{\mathrm{bias}}^{\mathrm{S}^2\mathrm{FT}})^2 + (1+\eta^{-1})(T_{\mathrm{variance}}^{\mathrm{S}^2\mathrm{FT}})^2.$$

Combined with (22) and (24), this concludes the proof.

Next we characterize the bias terms $T_{\rm bias}^{\rm LoRA}$ and $T_{\rm bias}^{\rm S^2FT}$ under sparsity assumption.

Lemma F.14. Suppose that Assumption F.4 holds. Then, for a sparse fine-tuned network with the choice $S \supset S_0$, it follows that

$$\mathcal{E}^{(\mathrm{i})}(f_{\ell}^{\mathrm{full}}) \leq (T_{\mathrm{bias}}^{\mathrm{LoRA}})^2 \leq (T_{\mathrm{bias}}^{\mathrm{S}^2\mathrm{FT}})^2 \leq \mathcal{E}^{(\mathrm{i})}(f_{\ell}^{\mathrm{full}}) + \delta^2 \kappa_*^2 (\overline{W}_{\ell+1}^{\mathrm{pre}}) \mathcal{E}^{(\mathrm{i})}(f^{\mathrm{pre}}).$$

Proof. Note that $\Phi_S''\Phi_S''^{\top}$ is a projection into a subspace, which is contained in a subspace projected by $\Phi'\Phi'^{\top}$. Thus

$$\begin{split} \|\Phi_S''\Phi_S''^{\top}D\Sigma_x^{(\mathbf{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger} &- \Phi'\Phi'^{\top}D\Sigma_x^{(\mathbf{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger}\|_{\mathrm{F}}^2 \\ &= \|(\Phi_S''\Phi_S''^{\top} - I)\Phi'\Phi'^{\top}D\Sigma_x^{(\mathbf{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger}\|_{\mathrm{F}}^2 \\ &= \|(\Phi_S''\Phi_S''^{\top} - I)\overline{W}_{\ell+1}^{\mathrm{pre}}(\overline{W}_{\ell+1}^{\mathrm{pre}})^{\dagger}D\Sigma_x^{(\mathbf{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger}\|_{\mathrm{F}}^2 \\ &= \|(\Phi_S''\Phi_S''^{\top} - I)\overline{W}_{\ell+1}^{\mathrm{pre}}((I - U_S^{\mathsf{S}^2\mathrm{FT}}U_S^{\mathsf{S}^2\mathrm{FT}\top}) + U_S^{\mathsf{S}^2\mathrm{FT}}U_S^{\mathsf{S}^2\mathrm{FT}\top})(\overline{W}_{\ell+1}^{\mathrm{pre}})^{\dagger}D\Sigma_x^{(\mathbf{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger}\|_{\mathrm{F}}^2 \\ &= \|(\Phi_S''\Phi_S''^{\top} - I)\overline{W}_{\ell+1}^{\mathrm{pre}}(I - U_S^{\mathsf{S}^2\mathrm{FT}}U_S^{\mathsf{S}^2\mathrm{FT}\top})(\overline{W}_{\ell+1}^{\mathrm{pre}})^{\dagger}D\Sigma_x^{(\mathbf{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}A^{\dagger}\|_{\mathrm{F}}^2, \end{split}$$

where the last equality follows since $(\Phi_S''\Phi_S''^{\top} - I)\overline{W}_{\ell+1}^{\text{pre}}U_S^{\text{22}\text{FT}} = 0$ by definition of $\Phi_S'' = \Phi_*(\overline{W}_{\ell+1}^{\text{pre}}U_S^{\text{22}\text{FT}})$. Thus

$$\begin{split} &\|\Phi_{S}'' \Phi_{S}''^{\top} D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} A^{\dagger} - \Phi' \Phi'^{\top} D \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} A^{\dagger}\|_{\text{F}}^{2} \\ &\leq \|\overline{W}_{\ell+1}^{\text{pre}}\|_{\text{op}}^{2} \|(I - U_{S}^{\text{S}^{2}\text{FT}} U_{S}^{\text{S}^{2}\text{FT}}) (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} D \Sigma_{x}^{(i)1/2} \|_{\text{F}}^{2} \|\Sigma_{x}^{(i)1/2} \underline{W}_{\ell-1}^{\text{pre}} A^{\dagger}\|_{\text{op}}^{2} \\ &= \|\overline{W}_{\ell+1}^{\text{pre}}\|_{\text{op}}^{2} \|\Sigma_{x}^{(i)1/2} \underline{W}_{\ell-1}^{\text{pre}} A^{\dagger}\|_{\text{op}}^{2} \sum_{a \in [d_{\ell}] \setminus S} \|e_{a}^{\top} (\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} D \Sigma_{x}^{(i)1/2} \|^{2} \\ &\leq \delta^{2} \|\overline{W}_{\ell+1}^{\text{pre}}\|_{\text{op}}^{2} \|(\overline{W}_{\ell+1}^{\text{pre}})^{\dagger} D \Sigma_{x}^{(i)1/2} \|_{\text{F}}^{2} \\ &\leq \delta^{2} \kappa_{*}^{2} (\overline{W}_{\ell+1}^{\text{pre}}) \|D \Sigma_{x}^{(i)1/2} \|_{\text{F}}^{2}, \end{split}$$

where the second inequality follows from $\|\Sigma_x^{(i)1/2}\underline{W}_{\ell-1}^{\mathrm{pre}^\top}A^{\dagger}\|_{\mathrm{op}} \leq 1$, Assumption F.4 and $S\supset S_0$. The conclusion follows from (13) and (23).

F.6.3 Out-of-distribution Excess Risk of Structured Sparse Fine-tuning

Given $S \subset [d_\ell]$ with |S| = s, we define the structured sparse adaptation matrix obtained by S²FT under population in-distribution risk as

$$V_{\infty}^{\mathrm{S}^{2}\mathrm{FT}} = \arg\min_{V} \|V\|_{\mathrm{F}}^{2} \quad \text{s.t. } V \text{ minimizes } \mathcal{R}^{(\mathrm{i})}(f_{\ell,U_{S}^{\mathrm{S}^{2}\mathrm{FT}},V}). \tag{25}$$

Theorem F.15 (Restatement of Theorem F.8: S²FT Part). Fix $S \subset [d_{\ell}]$ with |S| = s. For $V_{\infty}^{S^2FT}$ defined in (25),

$$\begin{split} \mathcal{E}^{(\mathrm{o})}(f_{\ell,U_{S}^{\mathsf{S}^{2}\mathsf{FT}},V_{\infty}^{\mathsf{S}^{2}\mathsf{FT}}}) &\leq \mathcal{E}^{(\mathrm{o})}(f^{\mathsf{pre}}) + 3 \big\| \Phi_{S}'' \Phi_{S}''^{\top}(B^{(\mathrm{o})} - B^{(\mathrm{i})}) \Sigma_{x}^{(\mathrm{o})1/2} \big\|_{\mathsf{F}}^{2} \\ &+ 3 \|B^{(\mathrm{i})}(\Sigma_{x}^{(\mathrm{o})1/2} - \Sigma_{x}^{(\mathrm{i})1/2} G_{\ell-1}^{(\mathrm{i},\mathrm{o})}) \|_{\mathsf{F}}^{2} \\ &+ 3 \|\overline{W}_{\ell}^{\mathsf{pre}}\|_{\mathsf{op}}^{2} \|\underline{W}_{\ell-1}^{\mathsf{pre}} \Sigma_{x}^{(\mathrm{o})1/2} - \underline{W}_{\ell-1}^{\mathsf{pre}} \Sigma_{x}^{(\mathrm{i},\mathrm{o})} \|_{\mathsf{F}}^{2}. \end{split}$$

Remark F.16. If there is no covariate shift, i.e., $\Sigma_x^{(i)} = \Sigma_x^{(o)} = \Sigma_x$ for some Σ_x , Theorem F.15 further gives the bound

$$\begin{split} \mathcal{E}^{(\text{o})}(f_{\ell,U_S^{\text{S}^2\text{FT}},V_\infty^{\text{S}^2\text{FT}}}) &\leq \mathcal{E}^{(\text{o})}(f^{\text{pre}}) + 3 \big\| \Phi_S'' \Phi_S''^{\top} (B^{(\text{o})} - B^{(\text{i})}) \Sigma_x^{1/2} \big\|_{\text{F}}^2 \\ &+ 3 \| B^{(\text{i})} \Sigma_x^{1/2} (I - (\underline{W}_{\ell-1}^{\text{pre}} \Sigma_x^{1/2})^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_x^{1/2})) \|_{\text{F}}^2. \end{split}$$

Proof of Theorem F.15. With a slight modification to Lemma F.12, we obtain

$$\begin{split} \mathcal{E}^{(\mathrm{o})}(f_{\ell,U_{S}^{\mathbf{S}^{2}\mathrm{FT}},V_{\infty}^{\mathbf{S}^{2}\mathrm{FT}}}) &= \operatorname{tr} \Bigg(\Big(B^{(\mathrm{o})} - W^{\mathrm{pre}} - \overline{W}_{\ell+1}^{\mathrm{pre}} U_{S}^{\mathbf{S}^{2}\mathrm{FT}} (\overline{W}_{\ell+1}^{\mathrm{pre}} U_{S}^{\mathbf{S}^{2}\mathrm{FT}})^{\dagger} D\Sigma_{x}^{(\mathrm{i})} \underline{W}_{\ell-1}^{\mathrm{pre}\top} (A^{\dagger})^{2} \underline{W}_{\ell-1}^{\mathrm{pre}} \Big) \Sigma_{x}^{(\mathrm{o})} \\ & \cdot \Big(B^{(\mathrm{o})} - W^{\mathrm{pre}} - \overline{W}_{\ell+1}^{\mathrm{pre}} U_{S}^{\mathbf{S}^{2}\mathrm{FT}} (\overline{W}_{\ell+1}^{\mathrm{pre}} U_{S}^{\mathbf{S}^{2}\mathrm{FT}})^{\dagger} D\Sigma_{x}^{(\mathrm{i})} \underline{W}_{\ell-1}^{\mathrm{pre}\top} (A^{\dagger})^{2} \underline{W}_{\ell-1}^{\mathrm{pre}} \Big)^{\top} \Bigg) \\ &= \Big\| (B^{(\mathrm{o})} - W^{\mathrm{pre}}) \Sigma_{x}^{(\mathrm{o})1/2} - \Phi_{S}'' \Phi_{S}''^{\top} D\Sigma_{x}^{(\mathrm{i})} \underline{W}_{\ell-1}^{\mathrm{pre}\top} (A^{\dagger})^{2} \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{o})1/2} \Big\|_{\mathrm{F}}^{2} \\ &= \| (I - \Phi_{S}'' \Phi_{S}''^{\top}) (B^{(\mathrm{o})} - W^{\mathrm{pre}}) \Sigma_{x}^{(\mathrm{o})1/2} \|_{\mathrm{F}}^{2} \\ &+ \Big\| \underline{\Phi}_{S}'' \Phi_{S}''^{\top} \Big\{ (B^{(\mathrm{o})} - W^{\mathrm{pre}}) \Sigma_{x}^{(\mathrm{o})1/2} - D\Sigma_{x}^{(\mathrm{i})1/2} (\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{i})1/2})^{\dagger} \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{o})1/2} \Big\} \Big\|_{\mathrm{F}}^{2}, \end{split}$$

where we used $\Sigma_x^{(i)1/2}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(A^\dagger)^2\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_x^{(\mathrm{o})1/2}=(\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_x^{(i)1/2})^\dagger\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_x^{(\mathrm{o})1/2}.$ Note that

$$\begin{split} \|T\|_{\mathrm{F}} &\leq \left\| \Phi_{S}'' \Phi_{S}''^{\top} \Big\{ B^{(\mathrm{o})} \Sigma_{x}^{(\mathrm{o})1/2} - B^{(\mathrm{i})} \Sigma_{x}^{(\mathrm{i})1/2} (\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{i})1/2})^{\dagger} \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{o})1/2} \Big\} \right\|_{\mathrm{F}} \\ &+ \left\| \Phi_{S}'' \Phi_{S}''^{\top} \overline{W}_{\ell}^{\mathrm{pre}} \Big\{ \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{o})1/2} - \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{i})1/2} (\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{i})1/2})^{\dagger} \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{o})1/2} \Big\} \right\|_{\mathrm{F}} \\ &\leq \left\| \Phi_{S}'' \Phi_{S}''^{\top} (B^{(\mathrm{o})} - B^{(\mathrm{i})}) \Sigma_{x}^{(\mathrm{o})1/2} \right\|_{\mathrm{F}} \\ &+ \left\| \Phi_{S}'' \Phi_{S}''^{\top} \overline{W}_{\ell}^{\mathrm{pre}} (\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{o})1/2} - \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{i},\mathrm{o})}) \right\|_{\mathrm{F}} \\ &+ \left\| \Phi_{S}'' \Phi_{S}''^{\top} \overline{W}_{\ell}^{\mathrm{pre}} (\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{o})1/2} - \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_{x}^{(\mathrm{i},\mathrm{o})} \right\|_{\mathrm{F}}. \end{split}$$

Therefore,

$$\begin{split} \mathcal{E}^{(\mathrm{o})}(f_{\ell,U_S^{\mathrm{S}^2\mathrm{FT}},V_\infty^{\mathrm{S}^2\mathrm{FT}}}) &= \|(I - \Phi_S'' \Phi_S''^\top) (B^{(\mathrm{o})} - W^{\mathrm{pre}}) \Sigma_x^{(\mathrm{o})1/2} \|_{\mathrm{F}}^2 + \|T\|_{\mathrm{F}}^2 \\ &\leq \mathcal{E}^{(\mathrm{o})}(f^{\mathrm{pre}}) + 3 \|\Phi_S'' \Phi_S''^\top (B^{(\mathrm{o})} - B^{(\mathrm{i})}) \Sigma_x^{(\mathrm{o})1/2} \|_{\mathrm{F}}^2 \\ &+ 3 \|B^{(\mathrm{i})} (\Sigma_x^{(\mathrm{o})1/2} - \Sigma_x^{(\mathrm{i})1/2} G_{\ell-1}^{(\mathrm{i},\mathrm{o})}) \|_{\mathrm{F}}^2 \\ &+ 3 \|\overline{W}_\ell^{\mathrm{pre}}\|_{\mathrm{op}}^2 \|\underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_x^{(\mathrm{o})1/2} - \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_x^{(\mathrm{i})1/2} G_{\ell-1}^{(\mathrm{i},\mathrm{o})} \|_{\mathrm{F}}^2 \end{split}$$

where we used $x + y + z \le 3x^2 + 3y^2 + 3z^2$. This concludes the proof.

F.7 Proofs for Full Fine-tuning

Define $f_\ell^{\mathrm{full}}(x) = \overline{W}_{\ell+1}^{\mathrm{pre}}(W_\ell^{\mathrm{pre}} + \Delta_\ell^{\mathrm{full}})\underline{W}_{\ell-1}^{\mathrm{pre}}x$ as a fine-tuned network with full fine-tuning applied to the ℓ -th layer, evaluated under the population in-distribution risk, where $\Delta_\ell^{\mathrm{full}}$ is obtained by

$$\Delta^{\text{full}}_{\ell} \in \mathop{\arg\min}_{\Delta' \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}} \mathbb{E}\bigg[\Big(B^{(\mathrm{i})} x^{(\mathrm{i})} - \overline{W}^{\mathrm{pre}}_{\ell+1} (W^{\mathrm{pre}}_{\ell} + \Delta') \underline{W}^{\mathrm{pre}}_{\ell-1} x^{(\mathrm{i})} \Big)^2 \bigg].$$

Lemma F.17 (In-distribution Excess Risk). For f_{ℓ}^{full} , it holds that

$$\begin{split} \mathcal{E}^{(\mathrm{i})}(f_{\ell}^{\mathrm{full}}) &= \|D\Sigma_{x}^{(\mathrm{i})1/2}(I - \Sigma_{x}^{(\mathrm{i})1/2}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(A^{2})^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_{x}^{(\mathrm{i})1/2})\|_{\mathrm{F}}^{2} \\ &+ \|(I - \Phi'\Phi'^{\top})D\Sigma_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(A^{2})^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_{x}^{(\mathrm{i})1/2}\|_{\mathrm{F}}^{2}. \end{split}$$

Proof of Lemma F.17. Similar to the proof of Theorem F.10, we have

$$\begin{split} \mathcal{E}^{(\mathrm{i})}(f_{\ell}^{\mathrm{full}}) &= \min_{\Delta \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}} \mathbb{E}\bigg[\Big(B^{(\mathrm{i})} x^{(\mathrm{i})} - \overline{W}_{\ell+1}^{\mathrm{pre}}(W_{\ell}^{\mathrm{pre}} + \Delta) \underline{W}_{\ell-1}^{\mathrm{pre}} x^{(\mathrm{i})} \Big)^2 \bigg] \\ &= \min_{\Delta \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}} \| D \Sigma_x^{(\mathrm{i})1/2} - \overline{W}_{\ell+1}^{\mathrm{pre}} \Delta \underline{W}_{\ell-1}^{\mathrm{pre}} \Sigma_x^{(\mathrm{i})1/2} \|_{\mathrm{F}}^2, \end{split}$$

and

$$\|D\Sigma_{x}^{(i)1/2} - \overline{W}_{\ell+1}^{\text{pre}} \Delta \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(i)1/2} \|_{\text{F}}^{2} = \| \underbrace{\overline{W}_{\ell+1}^{\text{pre}} \Delta \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(i)1/2} - \Phi' \Phi'^{\top} D\Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} A^{\dagger}}_{=:T_{1}} \|_{\text{F}}^{2}$$

$$+ \| \underbrace{D\Sigma_{x}^{(i)1/2} (I - \Sigma_{x}^{(i)1/2} \underline{W}_{\ell-1}^{\text{pre}} (A^{2})^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(i)1/2})}_{=:T_{2}} \|_{\text{F}}^{2}$$

$$+ \| \underbrace{(I - \Phi' \Phi'^{\top}) D\Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}} (A^{2})^{\dagger} \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(i)1/2}}_{=:T_{2}} \|_{\text{F}}^{2} ,$$

$$= :T_{2}$$

where we used the fact that the inner products $\operatorname{tr}(T_1T_2^\top) = \operatorname{tr}(T_2T_3^\top) = \operatorname{tr}(T_3T_1^\top) = 0$. By choosing $\Delta = (\overline{W}_{\ell+1}^{\operatorname{pre}})^{\dagger}D\Sigma_x^{(i)}\underline{W}_{\ell-1}^{\operatorname{pre}}A^{\dagger}$ for example, the term T_1 becomes 0. Thus

$$\begin{split} \mathcal{E}^{(\mathrm{i})}(f_{\ell}^{\mathrm{full}}) &= \|D\Sigma_{x}^{(\mathrm{i})1/2}(I - \Sigma_{x}^{(\mathrm{i})1/2}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(A^{2})^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_{x}^{(\mathrm{i})1/2})\|_{\mathrm{F}}^{2} \\ &+ \|(I - \Phi'\Phi'^{\top})D\Sigma_{x}^{(\mathrm{i})}\underline{W}_{\ell-1}^{\mathrm{pre}\top}(A^{2})^{\dagger}\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_{x}^{(\mathrm{i})1/2}\|_{\mathrm{F}}^{2}. \end{split}$$

This gives the desired result.

We obtain the following corollary as a direct consequence of Lemma F.17.

Corollary F.18. For f_{ℓ}^{full} , it holds that

$$\mathcal{E}^{(i)}(f_{\ell}^{\text{full}}) \leq \|\Psi_{*}^{\top}(D\Sigma_{x}^{(i)1/2})(I - \Sigma_{x}^{(i)1/2}\underline{W}_{\ell-1}^{\text{pre}}(A^{2})^{\dagger}\underline{W}_{\ell-1}^{\text{pre}}\Sigma_{x}^{(i)1/2})\|_{\text{op}}\mathcal{E}^{(i)}(f^{\text{pre}}) + \|(I - \Phi'\Phi'^{\top})\Phi_{*}(D\Sigma_{x}^{(i)1/2})\|_{\text{op}}\mathcal{E}^{(i)}(f^{\text{pre}}).$$
(27)

The first term on the right hand side of (27) measures the distance between two subspaces spanned by $\Psi_*(D\Sigma_x^{(\mathrm{i})1/2})$ and $\Psi_*(\underline{W}_{\ell-1}^{\mathrm{pre}}\Sigma_x^{(\mathrm{i})1/2})$. Intuitively, this quantifies the information coded at the ℓ -th layer, and the necessary information to predict residuals. Thus, it bounds the maximum improvement by the ℓ -th layer fine-tuning. The second term measures the subspace distance between the subspace where prediction residuals reside, and the subspace predictable by the ℓ -th layer fine-tuning.

G Auxiliary Results for Proofs

Lemma G.1. Fix $s, d_1, d_2 \in \mathbb{N}^+$. For any $A, B \in \mathbb{R}^{d_1 \times d_2}$, if $\|B - A\|_{op} \leq \|A\|_{op}$ and $\lambda_s(A) > \lambda_{s+1}(A)$ hold, then,

$$||SVD_s(B) - SVD_s(A)||_F \lesssim \kappa_*^2(A) \frac{\lambda_s(A)}{\lambda_s(A) - \lambda_{s+1}(A)} (\sqrt{s}||B - A||_{op} \wedge ||B - A||_F).$$

Proof. By triangle inequality,

$$\|SVD_{s}(B) - SVD_{s}(A)\|_{F} = \|\Phi_{s}(B)\Phi_{s}^{\top}(B)B - \Phi_{s}(A)\Phi_{s}^{\top}(A)A\|_{F}$$

$$\leq \|\Phi_{s}(B)\Phi_{s}^{\top}(B)(B - A)\|_{F} + \|(\Phi_{s}(B)\Phi_{s}^{\top}(B) - \Phi_{s}(A)\Phi_{s}^{\top}(A))A\|_{F}$$

$$\leq \sqrt{s}\|B - A\|_{op} + \|\Phi_{s}(B)\Phi_{s}^{\top}(B) - \Phi_{s}(A)\Phi_{s}^{\top}(A)\|_{F}\|A\|_{op}.$$

Using Davis-Kahan theorem (Theorem 4 from [73]), and Lemma 2.6 from [11],

$$\|\Phi_s(B)\Phi_s^{\top}(B) - \Phi_s(A)\Phi_s^{\top}(A)\|_{\mathsf{F}} \leq \frac{6\sqrt{2}\|A\|_{\mathsf{op}}(\sqrt{s}\|B - A\|_{\mathsf{op}} \wedge \|B - A\|_{\mathsf{F}})}{\lambda_s^2(A) - \lambda_{s+1}^2(A)}.$$

Thus

$$\|SVD_{s}(B) - SVD_{s}(A)\|_{F} \lesssim \frac{\|A\|_{op}^{2}}{\lambda_{s}^{2}(A)} \frac{\lambda_{s}^{2}(A)}{\lambda_{s}^{2}(A) - \lambda_{s+1}^{2}(A)} (\sqrt{s}\|B - A\|_{op} \wedge \|B - A\|_{F})$$
$$\lesssim \frac{\|A\|_{op}^{2}}{\lambda_{s}^{2}(A)} \frac{\lambda_{s}(A)}{\lambda_{s}(A) - \lambda_{s+1}(A)} (\sqrt{s}\|B - A\|_{op} \wedge \|B - A\|_{F}).$$

This concludes the proof.

We cite the concentration inequality for cross-covariance matrices from [47].

Lemma G.2 (Proposition 9.1 from [47]). Let Z and \tilde{Z} be mean zero random vectors taking values in \mathbb{R}^{d_1} and \mathbb{R}^{d_2} , respectively. Denote covariance matrices of Z and \tilde{Z} by Σ_Z and $\Sigma_{\tilde{Z}}$, respectively. Fix any t>0. Assume that there exist constants $c_1,c_2>0$ such that

$$\gamma^{\mathsf{T}} \Sigma_{Z} \gamma \ge c_{1} \| \gamma^{\mathsf{T}} Z \|_{\psi_{2}}^{2} \quad and \quad \gamma'^{\mathsf{T}} \Sigma_{\tilde{Z}} \gamma' \ge c_{2} \| \gamma'^{\mathsf{T}} \tilde{Z} \|_{\psi_{2}}^{2} \tag{28}$$

holds for any $\gamma \in \mathbb{R}^{d_1}$ and $\gamma' \in \mathbb{R}^{d_2}$. Choose $n \gg (r_e(\Sigma_Z) \wedge r_e(\Sigma_{\tilde{Z}}))(t + \log(d_1 + d_2))$. Let $(Z_i, \tilde{Z}_i)_{i \in [n]}$ be n independent copies of (Z, \tilde{Z}) . Then, there exists a constant $C = C(c_1, c_2) > 0$ such that with probability at least $1 - e^{-t}$,

$$\left\| \frac{1}{n} \sum_{i \in [n]} Z_i \tilde{Z}_i^\top - \mathbb{E}[Z \tilde{Z}^\top] \right\|_{\text{op}} \leq C \|\Sigma_Z\|_{\text{op}}^{1/2} \|\Sigma_{\tilde{Z}}\|_{\text{op}}^{1/2} \sqrt{\frac{(r_e(\Sigma_Z) + r_e(\Sigma_{\tilde{Z}})(t + \log(d_1 + d_2))}{n}}$$

hold.

Note that if a random variable Z taking values in \mathbb{R}^d satisfies $\gamma^\top \Sigma_Z \gamma \geq c \| \gamma^\top Z \|_{\psi_2}^2$ for any $\gamma \in \mathbb{R}^d$ with some constant c>0, AZ also satisfies ${\gamma'}^\top \Sigma_{AZ} {\gamma'} \geq c \| {\gamma'}^\top AZ \|_{\psi_2}^2$ for any ${\gamma'} \in \mathbb{R}^{d'}$ and any matrix $A \in \mathbb{R}^{d' \times d}$ and arbitrary $d' \in \mathbb{N}^+$, where $\Sigma_{AZ} = A\Sigma_Z A^\top$.

We then prove the following lemma to show the existance of a 'good' high probability event to bound multiple inequalities.

Lemma G.3. Suppose that Assumptions F.1 and F.2 hold. Fix any $S \subset [d_{\ell}]$. Then, there exists an event \mathcal{F} with $\mathbb{P}(\mathcal{F}) = 1 - \exp(-\Omega(\log^2(n+p+q)))$ such that on the event \mathcal{F} , for $\Phi \in \{\Phi', \Phi''_S\}$,

$$\|\Phi^{\top} \hat{D} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top}\|_{\text{op}} \lesssim \|D \Sigma_{x}^{(i)1/2}\|_{\text{op}} \|A\|_{\text{op}}, \quad \|\hat{A}^{\dagger}\|_{\text{op}} \lesssim \|A^{\dagger}\|_{\text{op}}, \tag{29}$$

and

$$\|(\hat{A}^2)^{\dagger} - (A^2)^{\dagger}\|_{\text{op}} \lesssim \frac{\kappa_*^2(A)}{\lambda_*^2(A)} \sqrt{\frac{r_e(A^2)\log^2(n+p+q)}{n}},$$
 (30)

$$\|\hat{A} - A\|_{\text{op}} \lesssim \kappa_*^2(A) \|A\|_{\text{op}} \sqrt{\frac{r_e(A^2) \log^2(n+p+q)}{n}},$$
 (31)

$$\|\hat{A}^{\dagger} - A^{\dagger}\|_{\text{op}} \lesssim \frac{\kappa_*(A)}{\lambda_*(A)} \sqrt{\frac{r_e(A^2)\log^2(n+p+q)}{n}}$$
(32)

hold. Furthermore,

$$\|\Phi^{\top}(\hat{D}\hat{\Sigma}_{x}^{(i)1/2} - D\Sigma_{x}^{(i)1/2})\underline{W}_{\ell-1}^{\text{pre}^{\top}}\|_{\text{op}}$$

$$\lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2}\|A\|_{\text{op}}\sqrt{\frac{(r_{e}(\Phi^{\top}\Sigma_{\epsilon}^{(i)}\Phi) + r_{e}(A^{2}))\log^{2}(n+p+q)}{n}}$$

$$+ \|D\Sigma_{x}^{(i)}D^{\top}\|_{\text{op}}^{1/2}\|A\|_{\text{op}}\sqrt{\frac{(r_{e}(\Phi^{\top}D\Sigma_{x}^{(i)}D^{\top}\Phi) + r_{e}(A^{2}))\log^{2}(n+p+q)}{n}}$$
(33)

holds on the event \mathcal{F} .

Proof. We only prove for $\Phi = \Phi'$ without loss of generality. Before proving Lemma G.3, we first derive several concentration inequalities. Assumption F.2 implies

$$n \gg r_e(A^2) \log^2(n+p+q),$$

$$n \gg r_e(\Sigma_x^{(i)}) \log^2(n+p+q),$$

$$n \gg (r_e(\Sigma_\epsilon^{(i)}) \wedge r_e(\Sigma_x^{(i)})) \log^2(n+p+q),$$

$$n \gg (r_e(\Phi^{\top} \Sigma_\epsilon^{(i)} \Phi) \wedge r_e(A^2)) \log^2(n+p+q),$$

$$n \gg (r_e(\Phi^{\top} D \Sigma_x^{(i)} D^{\top} \Phi) \wedge r_e(A^2)) \log^2(n+p+q).$$

Using Lemma G.2, we obtain

$$\|\hat{A}^{2} - A^{2}\|_{\text{op}} = \|\underline{W}_{\ell-1}^{\text{pre}} \hat{\Sigma}_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} - \underline{W}_{\ell-1}^{\text{pre}} \Sigma_{x}^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top}\|_{\text{op}}$$

$$\lesssim \|A\|_{\text{op}}^{2} \sqrt{\frac{r_{e}(A^{2}) \log^{2}(n+p+q)}{n}},$$
(34)

and

$$\|\hat{\Sigma}_{\epsilon,x}^{(i)}\|_{\text{op}} \lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \|\Sigma_{x}^{(i)}\|_{\text{op}}^{1/2} \sqrt{\frac{(r_{e}(\Sigma_{\epsilon}^{(i)}) + r_{e}(\Sigma_{x}^{(i)}))\log^{2}(n+p+q)}{n}},$$
(35)

$$\|\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)}\|_{\text{op}} \lesssim \|\Sigma_{x}^{(i)}\|_{\text{op}} \sqrt{\frac{r_{e}(\Sigma_{x}^{(i)})\log^{2}(n+p+q)}{n}},\tag{36}$$

$$\left\| \Phi^{\top} \hat{\Sigma}_{\epsilon,x}^{(i)} (\Sigma_x^{(i)})^{\dagger} \Sigma_x^{(i)} \underline{W}_{\ell-1}^{\text{pre}\top} \right\|_{\text{op}} \lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \|A\|_{\text{op}} \sqrt{\frac{(r_e(\Phi^{\top} \Sigma_{\epsilon}^{(i)} \Phi) + r_e(A^2)) \log^2(n+p+q)}{n}}, \tag{37}$$

$$\left\| \Phi^{\top} D(\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)}) \underline{W}_{\ell-1}^{\mathsf{pre}\top} \right\|_{\mathsf{op}} \lesssim \|D\Sigma_{x}^{(i)} D^{\top}\|_{\mathsf{op}}^{1/2} \|A\|_{\mathsf{op}} \sqrt{\frac{(r_{e}(\Phi^{\top} D\Sigma_{x}^{(i)} D^{\top} \Phi) + r_{e}(A^{2})) \log^{2}(n+p+q)}{n}}, \tag{38}$$

with high probability. Hereafter we only focus on the event \mathcal{F} where these inequalities hold. We divide the proof into 2 parts.

Part 1. In this part we derive (30), (31) and (32). Note that $\|\hat{A}^2 - A^2\|_{op} \le \lambda_*(A^2)/2$ holds on the event \mathcal{F} since $n \gg \kappa_*^4(A)r_e(A^2)\log^2(n+d+p)$ by Assumption F.2, and hence $\operatorname{rank}(\hat{A}^2) = \operatorname{rank}(A^2)$. Using Theorem 5.2 from [62],

$$\frac{\|(\hat{A}^2)^{\dagger} - (A^2)^{\dagger}\|_{\text{op}}}{\|(A^2)^{\dagger}\|_{\text{op}}} \lesssim \left(1 - \frac{\kappa_*(A^2)\|\hat{A}^2 - A^2\|_{\text{op}}}{\|A\|_{\text{op}}^2}\right)^{-1} \frac{\kappa_*(A^2)\|\hat{A}^2 - A^2\|_{\text{op}}}{\|A\|_{\text{op}}^2}.$$

Again from Assumption F.2, (34) gives

$$\|(\hat{A}^2)^{\dagger} - (A^2)^{\dagger}\|_{\text{op}} \lesssim \frac{\kappa_*(A^2)}{\lambda_*(A^2)} \sqrt{\frac{r_e(A^2)\log^2(n+p+q)}{n}}.$$

This yields (30). Proposition 3.2 from [67] and (34) yield,

$$\|(\Phi'''^{\top}\hat{A}^2\Phi''')^{1/2} - (\Phi'''^{\top}A^2\Phi''')^{1/2}\|_{\text{op}} \leq \frac{\|\Phi'''^{\top}(\hat{A}^2 - A^2)\Phi'''\|_{\text{op}}}{\lambda_*^{1/2}(\Phi'''^{\top}A^2\Phi''')} \lesssim \frac{\|A\|_{\text{op}}^2}{\lambda_*(A)}\sqrt{\frac{r_e(A^2)\log^2(n+p+q)}{n}},$$

where $\Phi''' := \Phi_*(A^2)$, and we used $\lambda_*(\Phi'''^\top A^2 \Phi''') \ge \lambda_*(A^2)$. Since $\hat{A} = \Phi'''(\Phi'''^\top \hat{A}^2 \Phi''')^{1/2} \Phi'''^\top$ and $A^{1/2} = \Phi'''(\Phi'''^\top A^2 \Phi''')^{1/2} \Phi'''^\top$, we obtain (31) as

$$\|\hat{A} - A\|_{\text{op}} \lesssim \kappa_*(A) \|A\|_{\text{op}} \sqrt{\frac{r_e(A^2) \log^2(n+p+q)}{n}}.$$
 (39)

Again using Theorem 5.2 from [62] combined with Assumption F.2, we obtain (32) as

$$\|\hat{A}^{\dagger} - A^{\dagger}\|_{\text{op}} \lesssim \frac{\kappa_*^2(A)}{\lambda_*(A)} \sqrt{\frac{r_e(A^2)\log^2(n+p+q)}{n}}.$$

This yields $\|\hat{A}^{\dagger}\|_{\text{op}} \lesssim \|A^{\dagger}\|_{\text{op}}$.

Part 2. Next we derive (33). By a similar argument as Part 1, (36) and Assumption F.2,

$$\|(\hat{\Sigma}_{x}^{(i)})^{\dagger} - (\Sigma_{x}^{(i)})^{\dagger}\|_{\text{op}} \lesssim \frac{\|\Sigma_{x}^{(i)}\|_{\text{op}}}{\lambda_{*}^{2}(\Sigma_{x}^{(i)})} \sqrt{\frac{r_{e}(\Sigma_{x}^{(i)})\log^{2}(n+d+p)}{n}}.$$
(40)

Since
$$\hat{D} - D = \check{\Sigma}_{\epsilon,x}^{(i)} = \hat{\Sigma}_{\epsilon,x}^{(i)} (\hat{\Sigma}_{x}^{(i)})^{\dagger},$$

$$\|\Phi^{\top}(\hat{D}\hat{\Sigma}_{x}^{(i)} - D\Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}}$$

$$\leq \|\Phi^{\top}(\hat{D} - D)\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}} + \|\Phi^{\top}D(\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}} + \|\Phi^{\top}(\hat{D} - D)(\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}}$$

$$= \|\Phi^{\top}\hat{\Sigma}_{\epsilon,x}^{(i)}(\hat{\Sigma}_{x}^{(i)})^{\dagger}\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}} + \|\Phi^{\top}D(\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}} + \|\Phi^{\top}\hat{\Sigma}_{\epsilon,x}^{(i)}(\hat{\Sigma}_{x}^{(i)})^{\dagger}(\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}}$$

$$\leq \|\Phi^{\top}\hat{\Sigma}_{\epsilon,x}^{(i)}(\Sigma_{x}^{(i)})^{\dagger}\Sigma_{x}^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}} + \|\Phi^{\top}\hat{\Sigma}_{\epsilon,x}^{(i)}((\Sigma_{x}^{(i)})^{\dagger}\Sigma_{x}^{(i)} - (\hat{\Sigma}_{x}^{(i)})^{\dagger}\Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}}$$

$$+ \|\Phi^{\top}D(\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}} + \|\Phi^{\top}\hat{\Sigma}_{\epsilon,x}^{(i)}(\hat{\Sigma}_{x}^{(i)})^{\dagger}(\hat{\Sigma}_{x}^{(i)} - \Sigma_{x}^{(i)})\underline{W}_{\ell-1}^{\mathrm{pre}\top}\|_{\mathrm{op}}$$

 $=: Q_1 + R_1 + Q_2 + R_2.$

We bound Q_1 , Q_2 , R_1 and R_2 separately. For the terms Q_1 and Q_2 , (37) and (38) give

$$Q_1 \lesssim \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \|A\|_{\text{op}} \sqrt{\frac{(r_e(\Phi^{\top}\Sigma_{\epsilon}^{(i)}\Phi) + r_e(A^2))\log^2(n+p+q)}{n}},\tag{41}$$

$$Q_2 \lesssim \|D\Sigma_x^{(i)}D^\top\|_{\text{op}}^{1/2}\|A\|_{\text{op}}\sqrt{\frac{(r_e(\Phi^\top D\Sigma_x^{(i)}D^\top\Phi) + r_e(A^2))\log^2(n+p+q)}{n}}.$$
 (42)

For the term R_1 , using (35) and (40),

$$\begin{split} R_{1} &\leq \|\hat{\Sigma}_{\epsilon,x}^{(i)}\|_{\text{op}} \|(\Sigma_{x}^{(i)})^{\dagger} - (\hat{\Sigma}_{x}^{(i)})^{\dagger}\|_{\text{op}} \|\Sigma_{x}^{(i)}\|_{\text{op}}^{1/2} \|\Sigma_{x}^{(i)1/2} \underline{W}_{\ell-1}^{\text{pre}\top}\|_{\text{op}} \\ &\lesssim \frac{\|\Sigma_{x}^{(i)}\|_{\text{op}}^{2} \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \|A\|_{\text{op}}}{\lambda_{*}^{2}(\Sigma_{x}^{(i)})} \sqrt{\frac{(r_{e}(\Sigma_{\epsilon}^{(i)}) + r_{e}(\Sigma_{x}^{(i)})) \log^{2}(n+p+q)}{n}} \sqrt{\frac{r_{e}(\Sigma_{x}^{(i)}) \log^{2}(n+p+q)}{n}} \\ &\lesssim \kappa_{*}^{2}(\Sigma_{x}^{(i)}) \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \|A\|_{\text{op}} \frac{\sqrt{r_{e}(\Sigma_{x}^{(i)})(r_{e}(\Sigma_{\epsilon}^{(i)}) + r_{e}(\Sigma_{x}^{(i)}))} \log^{2}(n+p+q)}{n} \end{split}$$

For the term R_2 , using (35) and (36),

$$\begin{split} R_2 &\leq \|\hat{\Sigma}_{\epsilon,x}^{(i)}\|_{\text{op}} \|(\hat{\Sigma}_x^{(i)})^{\dagger}\|_{\text{op}} \|\hat{\Sigma}_x^{(i)} - \Sigma_x^{(i)}\|_{\text{op}} \|(\Sigma_x^{(i)})^{\dagger}\|_{\text{op}}^{1/2} \|\Sigma_x^{(i)1/2} \underline{W}_{\ell-1}^{\text{pre}\top}\|_{\text{op}} \\ &\lesssim \|(\Sigma_x^{(i)})^{\dagger}\|_{\text{op}}^{3/2} \|A\|_{\text{op}} \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \|\Sigma_x^{(i)}\|_{\text{op}}^{3/2} \sqrt{\frac{(r_e(\Sigma_{\epsilon}^{(i)}) + r_e(\Sigma_x^{(i)})) \log^2(n+p+q)}{n}} \sqrt{\frac{r_e(\Sigma_x^{(i)}) \log^2(n+p+q)}{n}} \\ &\lesssim \kappa_*^{3/2} (\Sigma_x^{(i)}) \|\Sigma_{\epsilon}^{(i)}\|_{\text{op}}^{1/2} \|A\|_{\text{op}} \frac{\sqrt{r_e(\Sigma_x^{(i)})(r_e(\Sigma_{\epsilon}^{(i)}) + r_e(\Sigma_x^{(i)}))} \log^2(n+p+q)}{n}, \end{split}$$

where we used $\|(\hat{\Sigma}_x^{(i)})^{\dagger}\|_{\text{op}} \lesssim \|(\Sigma_x^{(i)})^{\dagger}\|_{\text{op}}$ by Assumption F.2 combined with (40). Again from Assumption F.2, $R_1 + R_2$ is bounded by the right hand side of (41). Therefore,

$$\begin{split} \|\Phi^{\top}(\hat{D}\hat{\Sigma}_{x}^{(\mathbf{i})} - D\Sigma_{x}^{(\mathbf{i})}) & \underline{W}_{\ell-1}^{\mathsf{pre}\top}\|_{\mathsf{op}} \\ & \lesssim \|\Sigma_{\epsilon}^{(\mathbf{i})}\|_{\mathsf{op}}^{1/2} \|A\|_{\mathsf{op}} \sqrt{\frac{(r_{e}(\Phi^{\top}\Sigma_{\epsilon}^{(\mathbf{i})}\Phi) + r_{e}(A^{2}))\log^{2}(n+p+q)}{n}} \\ & + \|D\Sigma_{x}^{(\mathbf{i})}D^{\top}\|_{\mathsf{op}}^{1/2} \|A\|_{\mathsf{op}} \sqrt{\frac{(r_{e}(\Phi^{\top}D\Sigma_{x}^{(\mathbf{i})}D^{\top}\Phi) + r_{e}(A^{2}))\log^{2}(n+p+q)}{n}}. \end{split}$$

Finally, from Assumption F.2, we obtain $\|\Phi^{\top}\hat{D}\hat{\Sigma}_{x}^{(i)}\underline{W}_{\ell-1}^{\mathrm{pre}^{\top}}\|_{\mathrm{op}} \lesssim \|D\Sigma_{x}^{(i)1/2}\|_{\mathrm{op}}\|\Sigma_{x}^{(i)1/2}\underline{W}_{\ell-1}^{\mathrm{pre}^{\top}}\|_{\mathrm{op}}$. This concludes the proof.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly state the contributions of this work.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the limitations of the work in the Conclusion Section (section 8).

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The main assumptions and theorems are provided in Section 4, while additional details and complete proofs can be found in Appendix F.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes].

Justification: The paper has disclosed all the information in the method and experiment part.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA].

Justification: We have the code required to reproduce our experimental results and are working towards making our code available in a public GitHub repository.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes].

Justification: The experimental setting is clearly described in Section 2, Section 5 and Section 6, and we will make our code available in a public GitHub repository.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes].

Justification: All statistics and results included in the paper are accompanied by confidence intervals.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

59946

Answer: [Yes].

Justification: Information for the resources required to reproduce the experiments are included in the oaoer. All experiments are run with 4 x A100 (80G). For the efficiency analysis, a single A100 GPU was used.

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes].

Justification: The research conducted in the paper fully conforms with the NeurIPS Code of Ethics in every respect.

10. **Broader Impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes].

Justification: We discuss the broader impacts of our work in Appendix.

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA].

Justification: Our paper does not introduce any assets that have a high risk for misuse.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes].

Justification: We have explicitly mentioned the citations for the datasets and have ensured that all conditions are fully respected.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA].

Justification: Upon acceptance, we will make our codebase publicly available and complete documentation for our assets.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA].

Justification: We do not include any experiments with human subjects or crowdsourcing.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA].

Justification: We do not include any experiments with human subjects or crowdsourcing.