

# META-LEARNING ADAPTABLE FOUNDATION MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The power of *foundation models* (FMs) lies in their capacity to learn highly expressive representations that can be adapted to a broad spectrum of tasks. However, these pretrained models require multiple stages of fine-tuning to become effective for downstream applications. Conventionally, the model is first retrained on the aggregate of a diverse set of tasks of interest and then adapted to specific low-resource downstream tasks by utilizing a parameter-efficient fine-tuning (PEFT) scheme. While this procedure seems reasonable, the independence of the retraining and fine-tuning stages causes a major issue, as there is no guarantee the retrained model will achieve good performance post-fine-tuning. To explicitly address this issue, we introduce a meta-learning framework infused with PEFT in this intermediate retraining stage to learn a model that can be easily adapted to unseen tasks. For our theoretical results, we focus on linear models using low-rank adaptations. In this setting, we demonstrate the suboptimality of standard retraining for finding an adaptable set of parameters. Further, we prove that our method recovers the optimally adaptable parameters. We then apply these theoretical insights to retraining the RoBERTa model to predict the continuation of conversations between different personas within the ConvAI2 dataset. Empirically, we observe significant performance benefits using our proposed meta-learning scheme during retraining relative to the conventional approach.

## 1 INTRODUCTION

*Foundation Models* (FMs) learn rich representations that are useful for a variety of downstream tasks. FMs are trained in three general stages to fit user-specific tasks like context-specific language generation and personalized image synthesis, among others. The first stage is commonly referred to as pretraining, where FMs are trained from scratch on a combination of massive public, propriety, and synthetic sources of data to learn a general-purpose model (Devlin et al., 2019; Brown et al., 2020; Abdin et al., 2024; Radford et al., 2021). This stage is largely inaccessible to most due to the enormous cost of training state-of-the-art models on such large datasets.

Thus, the most popular and viable way to utilize FMs for individual tasks is to take a pretrained model and *retrain* it for a specific objective. In this second training stage, we refine the pretrained model and retrain it on a large set of tasks of interest. For clarity, we generally refer to this intermediate stage as retraining. Other works have referred to this stage as pre-finetuning (Aghajanyan et al., 2021) or supervised fine-tuning (Dong et al., 2024). In the third stage, referred to as fine-tuning, the model is ultimately trained on an individual low-resource task. For example, a pretrained large language model (LLM) can be downloaded and retrained on a large multi-lingual corpus to perform English-Spanish and English-Italian translations. Then, one may adapt the model to translate English to French using a small English-French translation dataset. For this last stage, the model is typically fine-tuned using parameter efficient fine-tuning (PEFT) methods – training heuristics which sacrifice learning expressiveness for improved computational efficiency (Hu et al., 2021; Li & Liang, 2021). PEFT is especially useful in the low-resource setting, as running full fine-tuning of the model’s parameters on a small number of samples is expensive and potentially unnecessary.

Conventional retraining updates either a subset or all of the model parameters to fit the aggregation of the different retraining tasks. While this approach seems reasonable and has been successful in improving downstream task performance (Khashabi et al., 2020; Raffel et al., 2020), it does not leverage knowledge of the downstream fine-tuning procedure to cater the retrained model to perform well after such adaptation. Rather, it retrains the model to minimize the average loss across

054 the retraining tasks regardless of the PEFT method to be employed later. This raises two key issues.  
055 Firstly, there may not exist a single set of model parameters that simultaneously fits the various  
056 retraining tasks. Secondly, even if the model is sufficiently over-parameterized, there is no assurance  
057 the recovered retrained solution is indeed adaptable to future unseen tasks relative to other possible  
058 solutions during retraining, as the retraining and fine-tuning are performed independently.

059 We address these issues by drawing upon ideas from *meta-learning*, a framework designed to explic-  
060 itly train models for future adaptation. Meta-learning is a common method to improve model per-  
061 formance after fine-tuning, typically in low-resource, few-shot settings using gradient-based adap-  
062 tations (Finn et al., 2017; Lee & Choi, 2018). Moreover, it has begun to be applied to FM retraining  
063 to prepare models for downstream fine-tuning (Hou et al., 2022; Hong & Jang, 2022; Bansal et al.,  
064 2022; Gheini et al., 2022; Hu et al., 2023). However, it is not yet understood whether meta-learning  
065 how to fine-tune can provably confer performance benefits over standard retraining followed by  
066 PEFT. In this work, we provide rigorous theoretical and empirical evidence that this is indeed the  
067 case. We first study a stylized linear model where the ground truth parameters for both the retrain-  
068 ing and fine-tuning tasks are realizable by low-rank adaptations. We validate our theory through  
069 synthetic data and show that our insights improve performance on real language tasks using large  
070 language models (LLMs). Specifically, our contributions are as follows:

- 071 • We develop a generalized framework to model standard retraining and propose the Meta-  
072 Adapters objective for retraining, a meta-learning-inspired objective function for infusing  
073 PEFT in foundation model retraining. Our framework can be implemented with any PEFT  
074 algorithm, but we emphasize the incorporation of LoRA (Hu et al., 2021).
- 075 • For a linear model applied to multiple tasks whose ground truth parameters are realizable by  
076 LoRA, we show standard retraining does not recover an adaptable set of model parameters  
077 (Theorem 1) and thus incurs significant loss on unseen tasks after fine-tuning (Corollary 1).  
078 We prove two key results for the Meta-Adapter’s objective function:
  - 079 – Any model that globally minimizes this objective can be exactly fine-tuned to unseen  
080 tasks (Theorem 2), and when retraining on three or more tasks, the ground truth param-  
081 eters are the unique global minimum up to orthogonal symmetry (Theorem 3).  
082 This uniqueness property holds as long as the data dimension is sufficiently large,  
083 which is counterintuitive to previous work on multi-task learning theory that requires  
084 the number of tasks to be larger than the effective task dimension (Du et al., 2021;  
085 Collins et al., 2022).
  - 086 – For two retraining tasks, second-order stationarity is sufficient to guarantee global  
087 minimization for our Meta-Adapters loss (Theorem 4). In this case, our Meta-  
088 Adapters objective function is provably amenable to local optimization methods.
- 089 • To test our theoretical insights, we compare the performance of the standard retraining and  
090 Meta-Adapters objectives for linear models using LoRA while relaxing the assumptions  
091 from our theory. We show clear improvements using the Meta-Adapters objective for all  
092 data generation parameter settings and for different numbers of tasks. Then, we apply our  
093 meta-learning method to the RoBERTa (Liu et al., 2019) large language model (LLM) on  
094 the ConvAI2 dataset (Dinan et al., 2019), a real-world multi-task dataset for generating  
095 continuations of conversations between different personas. Again, we show improvements  
096 using the Meta-Adapters relative to retraining then fine-tuning.

## 097 1.1 RELATED WORK

098 Meta-learning is a framework for learning models that can be rapidly adapted to new unseen tasks  
099 by leveraging access to prior tasks during training. For example, Model-Agnostic Meta-Learning  
100 (MAML) (Finn et al., 2017) is a popular, flexible method that aims to find a model that can be  
101 adapted to a new unseen task after a small number of steps of gradient descent on the unseen task’s  
102 loss function. Further, other works have proposed methods specific to low-dimensional linear mod-  
103 els and have shown strong results and connections between meta-learning and representation learn-  
104 ing (Collins et al., 2022; Thekumparampil et al., 2021).

105  
106 In the case of FMs, other lines of work have proposed meta-learning approaches where the task-  
107 specific adaptation incorporates PEFT methods rather than few-shot gradient updates of all model  
parameters. (Hong & Jang, 2022; Bansal et al., 2022; Gheini et al., 2022) apply meta-learning with

architecture adaptations that inject small task-specific trainable layers within the FM architecture. (Hou et al., 2022) further combines architecture adaptations with parameter perturbation adaptations similar to LoRA. They consider a complicated meta-learning loss that separates the available training tasks data into training and testing tasks, and they update the adapters and FM weights over different splits of the data. Using combinations of architecture and parameter adaptation methods, they show empirical gains over retraining, then fine-tuning, and other gradient-based MAML-style algorithms. (Aghajanyan et al., 2021) similarly proposes a multi-task objective that trains an FM on different tasks simultaneously to encourage learning a universally applicable representation. They force the FM to learn a common shared data representation and apply a different prediction head for each retraining task. They run extensive empirical studies and observe performance improvements in a large-scale setting when 15 or more tasks are used in the retraining stage.

These works propose some kind of meta-learning or multi-task objective and show empirical gains over standard retraining strategies on natural language datasets, yet none explain when standard retraining is insufficient relative to meta-learning and multi-task approaches, how many tasks are needed to learn a rich representation, and how to best adapt to tasks unseen in the training stage.

Lastly, although we focus on LoRA, different PEFT methods have been proposed, including variants of LoRA (Liu et al., 2024; Dettmers et al., 2023; Zhang et al., 2023) and architecture adaptations (Houlsby et al., 2019) among others. Further, recent works have begun analyzing theoretical aspects of LoRA in the fine-tuning stage (Jang et al., 2024; Zeng & Lee, 2023). These works have started advancing the theory of LoRA, but they explore orthogonal directions to the analysis of meta-learning infused with LoRA. We include an extended discussion of these works in Appendix A.

**Notation.** We use bold capital letters for matrices and bold lowercase letters for vectors.  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  refers to the multivariate Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ .  $\|\cdot\|_F$  refers to the Frobenius norm.  $S_d$  refers to the set of  $d \times d$  symmetric matrices, and  $S_d^+ = \{\mathbf{X} \in S_d \mid \mathbf{X} \succcurlyeq \mathbf{0}\}$  is the set of  $d \times d$  symmetric positive semi-definite matrices.  $O_d$  refers to the set of  $d \times d$  orthogonal matrices.  $[n]$  refers to the set  $\{1, \dots, n\}$ . For a matrix  $\mathbf{X}$ ,  $\text{im}(\mathbf{X})$  and  $\text{ker}(\mathbf{X})$  refer to the image and kernel of  $\mathbf{X}$ , respectively. For subspaces  $M, N$ ,  $\dim(M)$  refers to the dimension of  $M$  and  $M + N = \{\mathbf{x} + \mathbf{y} \mid \mathbf{x} \in M, \mathbf{y} \in N\}$ . If  $M \cap N = \{\mathbf{0}\}$ , then we write the direct sum  $M \oplus N$ .

## 2 FOUNDATION MODEL RETRAINING AND FINE-TUNING

In this section, we first briefly recap the optimization process for conventional retraining of a foundation model (FM) across multiple tasks, followed by its fine-tuning on a downstream task. We then introduce our meta-learning-based approach which adjusts the retraining phase to incorporate insights from the final fine-tuning procedure.

### 2.1 STANDARD RETRAINING THEN FINE-TUNING

Consider a collection of  $T$  tasks of interest  $\mathcal{T} = \{\mathcal{T}_t\}_{t=1}^T$  where each task  $\mathcal{T}_t$  is drawn from task distribution  $\mathcal{D}$  and consists of  $n_t$  labeled examples  $\mathcal{T}_t = \{(\mathbf{x}_{t,j}, \mathbf{y}_{t,j})\}_{j=1}^{n_t}$ , where  $(\mathbf{x}_{t,j}, \mathbf{y}_{t,j})$  are i.i.d. from the  $t$ th task’s data distribution  $\mathcal{D}_{\mathcal{T}_t}$ . Without loss of generality we assume that for all tasks  $\mathcal{T}_t$  drawn from  $\mathcal{D}$ ,  $\mathcal{D}_{\mathcal{T}_t}$  generates samples  $\mathbf{x}_{t,j} \in \mathbb{R}^{d_x}$ ,  $\mathbf{y}_{t,j} \in \mathbb{R}^{d_y}$  for all  $t \in [T]$ ,  $j \in [n_t]$ . Consider a model  $\Phi(\cdot; \mathbf{W}) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$  parameterized by weights  $\mathbf{W}$  that maps feature vectors  $\mathbf{x} \in \mathbb{R}^{d_x}$  to predicted labels  $\hat{\mathbf{y}} \in \mathbb{R}^{d_y}$ . Typically  $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_m)$  is a list of matrices where  $\mathbf{W}_i \in \mathbb{R}^{d \times d}$  parameterize the layers of a neural network. We assume each  $\mathbf{W}_i$  is square for convenience.

**Retraining Phase.** Given a loss function  $\mathcal{L}$ , standard retraining attempts to minimize the aggregated loss over a collection of training tasks (Liu et al., 2019; Brown et al., 2020). This amounts to solving

$$\min_{\mathbf{W}} \sum_{t=1}^T \sum_{j=1}^{n_t} \mathcal{L}(\Phi(\mathbf{x}_{t,j}; \mathbf{W}), \mathbf{y}_{t,j}). \quad (1)$$

In other words, the above optimization problem seeks a set of universal parameters that define a unique mapping function capable of translating inputs to outputs across all tasks involved in the retraining phase. We denote the set of weights obtained by solving (1) as  $\hat{\mathbf{W}}_{\text{SR}}$ , and the corresponding input-output mapping function as  $\Phi(\cdot; \hat{\mathbf{W}}_{\text{SR}})$ , where SR stands for Standard Retraining.

**Fine-Tuning Phase.** In the subsequent fine-tuning step, we refine either the retrained weights, the model’s feature map, or both to fit a downstream task with fewer labeled samples. More precisely, consider a downstream task  $\mathcal{T}_{T+1}$  drawn from the same distribution  $\mathcal{D}$  where  $\mathcal{T}_{T+1} = \{(\mathbf{x}_{T+1,j}, \mathbf{y}_{T+1,j})\}_{j=1}^{n_{T+1}}$ . To fit the model to task  $\mathcal{T}_{T+1}$  we do not retrain the retrained weights  $\hat{\mathbf{W}}_{\text{SR}}$ , but instead fine-tune the mapping  $\Phi(\cdot; \hat{\mathbf{W}}_{\text{SR}})$  using additional parameters  $\theta$ . For example,  $\theta$  could parameterize transformations of  $\hat{\mathbf{W}}_{\text{SR}}$  that adapt the retrained weights or new trainable layers inserted into the architecture of the retrained model (Hu et al., 2021; Liu et al., 2024; Aghajanyan et al., 2021). We denote the fine-tuned model’s mapping as  $\Phi_{\text{FT}}(\cdot; \hat{\mathbf{W}}_{\text{SR}}, \theta) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ . During the *fine-tuning stage*, the goal is to find the optimal additional parameters,  $\theta$ , that minimize the loss for the downstream task  $\mathcal{T}_{T+1}$ , solving:

$$\min_{\theta} \sum_{j=1}^{n_{T+1}} \mathcal{L}(\Phi_{\text{FT}}(\mathbf{x}_{T+1,j}; \hat{\mathbf{W}}_{\text{SR}}, \theta), \mathbf{y}_{T+1,j}). \quad (2)$$

In particular, when the LoRA PEFT method is used for fine-tuning, the model is adapted to task  $\mathcal{T}_{T+1}$  by fixing the model architecture and the retrained weights  $\hat{\mathbf{W}}_{\text{SR}}$  and only training low-rank perturbations for each of the matrices  $\hat{\mathbf{W}}_{\text{SR},1}, \dots, \hat{\mathbf{W}}_{\text{SR},m}$ . For rank- $k$  adaptations, we parameterize  $\theta = ((\mathbf{U}_1, \mathbf{V}_1), \dots, (\mathbf{U}_m, \mathbf{V}_m))$ , where  $\mathbf{U}_i, \mathbf{V}_i \in \mathbb{R}^{d \times k}$  are the factors of the low-rank adaptation of the  $i$ th matrix in  $\hat{\mathbf{W}}_{\text{SR}}$ . The fine-tuned model is just the original model where the  $i$ th weight matrix  $\mathbf{W}_i$  is now perturbed to be  $\mathbf{W}_i + \mathbf{U}_i \mathbf{V}_i^\top$ . Then the LoRA fine-tuning optimization problem is:

$$\min_{\{\mathbf{U}_i, \mathbf{V}_i\}_{i=1}^m} \sum_{j=1}^{n_{T+1}} \mathcal{L}(\Phi(\mathbf{x}_{T+1,j}; (\hat{\mathbf{W}}_{\text{SR},1} + \mathbf{U}_1 \mathbf{V}_1^\top, \dots, \hat{\mathbf{W}}_{\text{SR},m} + \mathbf{U}_m \mathbf{V}_m^\top)), \mathbf{y}_{T+1,j}). \quad (3)$$

This pipeline seems reasonable as we first fit the model to the aggregation of the retraining tasks which we hope will promote learning the general structure of the tasks drawn from  $\mathcal{D}$ . However, there may not exist a single model that can model each retraining task simultaneously, so retraining the model on the aggregation of the retraining tasks does not align with our implicit assumption that each task is realizable after task-specific adaptations from a common model. Further, even if the model is sufficiently overparameterized where many possible solutions fit the retraining tasks, standard retraining finds a solution independent of the subsequent PEFT method to be used for fine-tuning. Nothing about standard retraining promotes learning an adaptable solution relative to other candidate solutions that fit the retraining tasks.

## 2.2 META-ADAPTERS

Since the ultimate goal of our model is to perform well on a variety of unseen downstream tasks, we propose the Meta-Adapters objective that explicitly fits weights and adapter parameters to the training tasks. Intuitively, this objective promotes sets of parameters that can be adapted to future unseen tasks drawn from the same distribution as those seen in retraining.

Rather than training a single model on the aggregation of the retraining tasks, we instead incorporate the adapters during the retraining process and learn adapted models for each task. Let  $\theta^{(t)}$  be the set of adapter parameters for the  $t_{th}$  training task  $\mathcal{T}_t$ . The Meta-Adapters method searches for a single set of base weights  $\hat{\mathbf{W}}_{\text{Meta}}$  such that for all  $t \in [T]$ , the  $t_{th}$  adapted model  $\Phi_{\text{FT}}(\cdot; \hat{\mathbf{W}}_{\text{Meta}}, \theta^{(t)})$  minimizes the loss over the training task  $\mathcal{T}_t$ . We define the Meta-Adapters objective as:

$$\min_{\mathbf{W}} \sum_{t=1}^T \min_{\theta^{(t)}} \sum_{j=1}^{n_t} \mathcal{L}(\Phi_{\text{FT}}(\mathbf{x}_{t,j}; \mathbf{W}, \theta^{(t)}), \mathbf{y}_{t,j}). \quad (4)$$

When we use LoRA as the adaptation method, we define  $\mathbf{U}_i^{(t)} (\mathbf{V}_i^{(t)})^\top \in \mathbb{R}^{d \times d}$  as the factorization of the low-rank adapter for the  $i$ th weight matrix for the  $t_{th}$  task. Then the objective reduces to:

$$\min_{\mathbf{W}} \sum_{t=1}^T \min_{\{\mathbf{U}_i^{(t)}, \mathbf{V}_i^{(t)}\}_{i=1}^m} \sum_{j=1}^{n_t} \mathcal{L}(\Phi(\mathbf{x}_{t,j} (\mathbf{W}_1 + \mathbf{U}_1^{(t)} \mathbf{V}_1^{(t)\top}, \dots, \mathbf{W}_m + \mathbf{U}_m^{(t)} \mathbf{V}_m^{(t)\top})), \mathbf{y}_{t,j}). \quad (5)$$

In this case, we refer to the objective function as Meta-LoRA. This proposed optimization problem is designed to replace the standard retraining objective in (1). After minimizing (4) we recover base parameters  $\hat{\mathbf{W}}_{\text{Meta}}$  that are explicitly designed to be adaptable downstream. To perform finetuning, we then run the exact same minimization in (2) but using retrained weights  $\hat{\mathbf{W}}_{\text{Meta}}$  instead of  $\hat{\mathbf{W}}_{\text{SR}}$ .

### 3 THEORETICAL RESULTS

To establish our theoretical results, we consider  $T$  multi-output linear regression retraining tasks and one test task, with the caveat that the ground-truth regressor for each task is a low-rank modification of a common single matrix. More precisely, consider the matrix  $\mathbf{A}^* \in \mathbb{R}^{d \times d}$ , which is a common parameter shared across all tasks, and task-specific adapters  $\mathbf{U}_t^* \mathbf{U}_t^{*\top}$  for  $t \in [T+1]$ , where  $\mathbf{U}_t^* \in \mathbb{R}^{d \times k}$  and the entries of  $\mathbf{U}_t^*$  are i.i.d. from  $\mathcal{N}(0, 1)$ . We work in the setting where  $k \ll d$  and  $k(T+1) < d$ . Assume the data generation for task  $\mathcal{T}_t \sim \mathcal{D}$  is given by  $\mathbf{y}_{t,j} = (\mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top}) \mathbf{x}_{t,j} + \epsilon_{t,j}$ . Here,  $\mathbf{x}_{t,j}$  is the  $j$ -th sample of task  $t$  which is i.i.d. from  $\mathcal{N}(\mathbf{0}, \sigma_x^2 \mathbf{I}_d)$ , and  $\epsilon_{t,j}$  is i.i.d.  $\mathcal{N}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I}_d)$  noise sampled independently of the data  $\mathbf{x}_{t,j}$ . As mentioned above,  $\mathbf{A}^*$  can be considered as the common parameter which is close to the ground truth of each task up to a low-rank adaptation.

For each task  $t$ , the learner uses the linear predictor  $\Phi(\mathbf{x}; \mathbf{A}_t) = \mathbf{A}_t \mathbf{x}$  for  $\mathbf{A}_t \in \mathbb{R}^{d \times d}$ ,  $\mathbf{x} \in \mathbb{R}^d$ . In the ideal case, we hope to recover parameter value  $\hat{\mathbf{A}} = \mathbf{A}^*$  in the retraining phase so that the fine-tuned model  $\Phi_{\text{FT}}(\mathbf{x}; \hat{\mathbf{A}}, \mathbf{U}, \mathbf{V}) = (\hat{\mathbf{A}} + \mathbf{U}\mathbf{V}^\top) \mathbf{x}$  with a proper low-rank adapter  $\mathbf{U}\mathbf{V}^\top$  can fit the data distribution of any downstream task also drawn from  $\mathcal{D}$ .

Given  $N$  samples for each task, the loss for each task is  $\mathcal{L}_t^N(\mathbf{A}_t) = \frac{1}{2N} \sum_{j=1}^N \|\mathbf{y}_{t,j} - \mathbf{A}_t \mathbf{x}_{t,j}\|_2^2 = \frac{1}{2N} \sum_{j=1}^N \|(\mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{A}_t) \mathbf{x}_{t,j} + \epsilon_{t,j}\|_2^2$ . We define  $\mathcal{L}_t(\mathbf{A}_t)$  as the shifted and scaled infinite sample loss:

$$\mathcal{L}_t(\mathbf{A}_t) = \frac{1}{2} \left\| \mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{A}_t \right\|_F^2 = \frac{1}{\sigma_x^2} \left( \mathbb{E} [\mathcal{L}_t^N(\mathbf{A}_t)] - \frac{\sigma_\epsilon^2}{2} \right) \quad (6)$$

We assume access to infinite samples during the retraining process, as in practice, we have access to large retraining tasks relative to the low-resource downstream tasks to be used for fine-tuning.

**Remark 1.** For convenience, we require a mild sense of task diversity and assume that the aggregated columns from all  $\mathbf{U}_t^*$ ,  $t \in [T+1]$ , form a linearly independent set. Precisely, we assume  $\dim(\text{im}(\mathbf{U}_1^*) \oplus \dots \oplus \text{im}(\mathbf{U}_{T+1}^*)) = k(T+1)$ . Since  $k(T+1) < d$ , the nature of the generation process of each  $\mathbf{U}_t^*$  ensures that this assumption holds almost surely.

Given access to the loss functions defined in 6, the goal of the learner is to find an  $\hat{\mathbf{A}}$  that can be adapted to the unseen task  $\mathcal{T}_{T+1}$ . The infinite sample test loss for adapter factors  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1}$  and fixed  $\hat{\mathbf{A}}$  is the LoRA loss on  $\mathcal{T}_{T+1}$  which reduces to the low-rank matrix factorization problem:

$$\mathcal{L}_{\text{Test}}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}) = \frac{1}{2} \left\| \mathbf{A}^* + \mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{*\top} - \hat{\mathbf{A}} - \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top \right\|_F^2. \quad (7)$$

We compare the standard retraining and Meta-LoRA objectives for utilizing each  $\mathcal{L}_t$  to recover a common set of base parameters  $\hat{\mathbf{A}}$  whose low-rank adaptation  $\hat{\mathbf{A}} + \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top$  minimizes the test loss  $\mathcal{L}_{\text{Test}}$  for some  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1}$ . We include complete proofs for all theorems in Appendix B.

#### 3.1 STANDARD RETRAINING THEN FINE-TUNING

First, consider the standard retraining then fine-tuning setup as a candidate for ultimately minimizing 7. Here, the learner first finds a single matrix  $\hat{\mathbf{A}}_{\text{SR}}$  that minimizes the sum of losses  $\sum_{t=1}^T \mathcal{L}_t$ :

$$\hat{\mathbf{A}}_{\text{SR}} = \arg \min_{\mathbf{A}} \frac{1}{2} \sum_{t=1}^T \left\| \mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{A} \right\|_F^2. \quad (8)$$

Then when given a new task  $\mathcal{T}_{T+1}$ , the learner runs LoRA to minimize the loss over the unseen task in 7. However, this strategy suffers substantial loss on the test task.

**Theorem 1.** For standard retraining,  $\text{rank}(\hat{\mathbf{A}}_{\text{SR}} - \mathbf{A}^*) = kT$ .

The above theorem demonstrates that the standard retraining process is unable to recover the ground truth shared matrix  $\mathbf{A}^*$ . Specifically, it shows that the discrepancy between the obtained solution and the ground truth  $\mathbf{A}^*$  has a rank of  $kT$ . Consequently, any fine-tuning method constrained to a rank lower than  $kT$  will fail to recover the correct model for the downstream task. This result follows from the fact that the obtained model from the standard retraining scheme can be written as

$$\hat{\mathbf{A}}_{\text{SR}} = \mathbf{A}^* + \frac{1}{T} \sum_{t=1}^T \mathbf{U}_t^* \mathbf{U}_t^{*\top}. \quad (9)$$

Now, given the fact that  $\mathbf{U}_t^*$  are linearly independent, it follows that  $\text{rank}(\sum_{t=1}^T \mathbf{U}_t^* \mathbf{U}_t^{*\top}) = kT$ . Hence,  $\hat{\mathbf{A}}_{\text{SR}}$  is far from both  $\mathbf{A}^*$  and the test task ground truth parameters.

**Corollary 1.** For number of retraining tasks  $T \geq 1$ , if test task adaptation rank  $k' < k(T+1)$ , then  $\mathcal{L}_{\text{Test}}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}_{\text{SR}}) > 0$  for all rank- $k'$  adapters  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1}^\top$  where  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1} \in \mathbb{R}^{d \times k'}$ .

**Corollary 2.** For a large number of retraining tasks  $T$  and test task adaptation rank  $k' < k(T+1)$ ,  $\mathcal{L}_{\text{Test}}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}_{\text{SR}}) = \Omega((d - k')k^2)$  for all  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1} \in \mathbb{R}^{d \times k'}$ .

Both corollaries follow from the classic result of (Mirsky, 1960). In the infinite sample setting, the LoRA rank needed to fit the test task after standard retraining is  $k(T+1)$ , and using anything smaller results in test error that scales with  $d$ . Thus, **standard retraining recovers parameters that cannot be low-rank adapted to any relevant task**. To address these issues, we employ the Meta-LoRA objective which explicitly searches for a low-rank adaptable solution.

### 3.2 META-LORA

Although we have shown that standard retraining can lead to large losses on downstream tasks after LoRA, it is not yet clear whether any other retraining method can do better in this setting. We next explore whether minimizing the Meta-LoRA objective results in a matrix  $\hat{\mathbf{A}}_{\text{Meta}}$  that indeed leads to a smaller test loss  $\mathcal{L}_{\text{Test}}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}_{\text{Meta}})$  for some values of  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1}$ .

As in (5), we introduce low-rank adapters during the retraining phase to model the different training tasks. We search for a value of  $\mathbf{A}$  such for all  $\mathcal{T}_t$ , the loss  $\mathcal{L}_t$  after running LoRA on  $\mathcal{T}_t$  is minimized. This promotes values of  $\mathbf{A}$  that can be easily adapted to unseen tasks downstream. We use the Meta-LoRA loss but with symmetric low-rank adapters  $\mathbf{U}_t \mathbf{U}_t^\top$  for the  $t_{th}$  task  $\mathcal{T}_t$  in retraining. We allow asymmetric adapters at test time. The infinite sample Meta-LoRA loss is then

$$\mathcal{L}_{\text{Meta}}(\mathbf{A}) = \sum_{t=1}^T \min_{\mathbf{U}_t} \mathcal{L}_t(\mathbf{A} + \mathbf{U}_t \mathbf{U}_t^\top). \quad (10)$$

Define the concatenation of each  $\mathbf{U}_t$  as  $\mathbf{U} = (\mathbf{U}_1, \dots, \mathbf{U}_T) \in (\mathbb{R}^{d \times k})^T$ . Then minimizing (10) is equivalent to solving  $\min_{\mathbf{A}, \mathbf{U}} \mathcal{L}(\mathbf{A}, \mathbf{U})$  where

$$\mathcal{L}(\mathbf{A}, \mathbf{U}) = \frac{1}{2} \sum_{t=1}^T \left\| \mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{A} - \mathbf{U}_t \mathbf{U}_t^\top \right\|_F^2. \quad (11)$$

We have seen that standard retraining does not recover an optimal solution, but it is unclear what the global minima of this new objective function are and if they can be easily found. Note that by fixing  $\mathbf{A}$ , (11) is  $T$  independent symmetric matrix factorization problems, and by fixing  $\mathbf{U}$ , (11) is a convex quadratic problem over  $\mathbf{A}$ . Despite these well-understood sub-problems, joint minimization over  $\mathbf{A}$  and  $\mathbf{U}$  presents challenging variable interactions that complicate the analysis. Nevertheless, we employ a careful landscape analysis of (11) to address these questions.

#### 3.2.1 LANDSCAPE OF GLOBAL MINIMA OF (11)

First, we show that the objective is well-posed, i.e., minimization of  $\mathcal{L}$  leads to an adaptable solution.

**Theorem 2.** For any  $T \geq 1$ , if  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$ , then  $\hat{\mathbf{A}} = \mathbf{A}^* + \mathbf{C}$  where  $\text{rank}(\mathbf{C}) \leq 2k$

Clearly, any point is a global minimum of (11) if and only if it achieves zero loss. Theorem 2 guarantees that the values of  $\mathbf{A}$  that induce global minima of (11) are at most rank- $2k$  away from the ground truth parameter  $\mathbf{A}^*$ .

**Corollary 3.** For any  $T \geq 1$ , if  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$ , then there exists a rank- $3k$  adapter  $\mathbf{U}_{T+1}\mathbf{V}_{T+1}^\top$  where  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1} \in \mathbb{R}^{d \times 3k}$  such that  $\mathcal{L}_{Test}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}) = 0$ .

This again follows from classic low-rank factorization results, as  $\mathcal{L}_{Test}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \mathbf{A}^* + \mathbf{C}) = \frac{1}{2} \|\mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{\top} - \mathbf{C} - \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top\|_F^2$  and  $\text{rank}(\mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{\top} - \mathbf{C}) \leq 3k$ . Note that  $3k$  is still much smaller than  $d$  as  $k \ll d$ .

*Proof sketch of Theorem 2.* Notice that any set of parameters  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  such that  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$  must be a critical point as  $\mathcal{L} \geq 0$ . This directly implies that  $\hat{\mathbf{A}} = \mathbf{A}^* + \frac{1}{T} \sum_{t=1}^T \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{U}_t \mathbf{U}_t^\top$  and  $\mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{U}_t \mathbf{U}_t^\top = \mathbf{U}_j^* \mathbf{U}_j^{*\top} - \mathbf{U}_j \mathbf{U}_j^\top$  for all  $1 \leq i, j \leq T$ . It then follows that  $\hat{\mathbf{A}} = \mathbf{A}^* + \mathbf{U}_1^* \mathbf{U}_1^{*\top} - \mathbf{U}_1 \mathbf{U}_1^\top$ , and  $\text{rank}(\mathbf{U}_1^* \mathbf{U}_1^{*\top} - \mathbf{U}_1 \mathbf{U}_1^\top) \leq \text{rank}(\mathbf{U}_1^* \mathbf{U}_1^{*\top}) + \text{rank}(\mathbf{U}_1 \mathbf{U}_1^\top) \leq 2k$ .  $\square$

This result shows that for any  $T \geq 1$ , any global minimum of (11) recovers  $\mathbf{A}^*$  with an error up to rank- $2k$ . Consequently, it can perform well on a downstream task after fine-tuning with a rank- $3k$  adaptor. Furthermore, we demonstrate that when the number of tasks satisfies  $T \geq 3$ , a stronger result can be established. Specifically, in this case, we can prove the exact recovery of the ground truth parameter  $\mathbf{A}^*$  is possible.

**Theorem 3.** If  $T \geq 3$ , then  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$  implies  $\hat{\mathbf{A}} = \mathbf{A}^*$  and  $\mathbf{U}_t \mathbf{U}_t^\top = \mathbf{U}_t^* \mathbf{U}_t^{*\top}$  for all  $t \in [T]$

Theorem 3 guarantees that the ground truth parameters are the unique global minimum up to orthogonal symmetry when there are three or more tasks, regardless of the ambient dimension or the number of columns  $k$ . This result is surprising, as most theoretical results for multi-task learning require higher task diversity, typically where the number of tasks  $T$  is required to be larger than the effective task dimension  $k$  (Du et al., 2021; Collins et al., 2022). However, we establish this uniqueness result for the absolute condition  $T \geq 3$ . This implies that exact test task fine-tuning can be achieved with a rank  $k$ -adaptation.

**Corollary 4.** For any  $T \geq 3$ , if  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$ , then there exists a rank- $k$  adapter  $\mathbf{U}_{T+1}\mathbf{V}_{T+1}^\top$  where  $\mathbf{U}_{T+1}, \mathbf{V}_{T+1} \in \mathbb{R}^{d \times k}$  such that  $\mathcal{L}_{Test}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}) = 0$ .

This follows directly from  $\mathcal{L}_{Test}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \mathbf{A}^*) = \frac{1}{2} \|\mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{\top} - \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top\|_F^2$  and  $\text{rank}(\mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{\top}) = k$ .

*Proof sketch of Theorem 3.* We again rely on the fact that a set of parameters that achieves zero loss must satisfy  $\mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top = \mathbf{U}_s^* \mathbf{U}_s^{*\top} - \hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top$  for all  $t, s \in [T]$ . Then

$$\mathbf{U}_1^* \mathbf{U}_1^{*T} = \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^T + \mathbf{U}_2^* \mathbf{U}_2^{*T} - \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^T = \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^T + \mathbf{U}_3^* \mathbf{U}_3^{*T} - \hat{\mathbf{U}}_3 \hat{\mathbf{U}}_3^T.$$

Since  $\mathbf{U}_1^* \mathbf{U}_1^{*T} \succcurlyeq 0$ , both  $\text{im}(\hat{\mathbf{U}}_2) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_2^*)$  and  $\text{im}(\hat{\mathbf{U}}_3) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_3^*)$ . This then implies that the image of  $\hat{\mathbf{U}}_1$  is a subset of two key subspaces:

$$\text{im}(\mathbf{U}_1^*) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_2^*) \quad \text{and} \quad \text{im}(\mathbf{U}_1^*) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_3^*). \quad (12)$$

We then make use of a key lemma to prove the result. The proof can be found in Appendix B.3.

**Lemma 1.**  $([\text{im}(\hat{\mathbf{U}}_1) \oplus \text{im}(\mathbf{U}_2^*)] \cap [\text{im}(\hat{\mathbf{U}}_1) \oplus \text{im}(\mathbf{U}_3^*)]) = \text{im}(\hat{\mathbf{U}}_1)$

Combining Lemma 1 with (12) implies that  $\text{im}(\hat{\mathbf{U}}_1) = \text{im}(\mathbf{U}_1^*)$ . Then applying the same argument for the other indices shows that  $\text{im}(\hat{\mathbf{U}}_t) = \text{im}(\mathbf{U}_t^*)$  for all  $t \in [T]$ . Since  $\mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top = \mathbf{U}_s^* \mathbf{U}_s^{*\top} - \hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top$  for all  $t, s \in [T]$  and  $\text{im}(\mathbf{U}_1^*) \cap \text{im}(\mathbf{U}_2^*) = \{\mathbf{0}\}$ , it follows that  $\mathbf{U}_t^* \mathbf{U}_t^{*\top} = \hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top$  for all  $t \in [T]$ . Then since  $\nabla_{\mathbf{A}} \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$ ,  $\hat{\mathbf{A}} = \mathbf{A}^* + \frac{1}{T} \sum_{t=1}^T \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top = \mathbf{A}^*$ .  $\square$

The proof of Theorem 3 relies on the assumption that there are at least three training tasks. This is necessary to some degree as if there are only two tasks, we can construct ground truth parameters that have infinite solutions as in the example in Appendix D.1.

**Summary.** The previous two theorems show that for any  $T \geq 1$ , the set of global minima of the meta objective is always adaptable to the downstream task. Furthermore, if  $T \geq 3$ , the global minima of

the meta-objective are the unique ground truth parameters  $(\mathbf{A}^*, \mathbf{U}_t^*)$  up to orthogonal symmetry of  $\mathbf{U}_t$ . In other words, minimizing (11) guarantees the recovery of the ground truth parameters.

### 3.2.2 ALGORITHMS FOR MINIMIZING (11)

The above results establish that minimizing the meta-objective (11) leads to recovery of the ground truth parameters, with a small error term when  $T = 2$ . However, it is unclear if this minimization problem can always be solved by local optimization methods.

**Theorem 4.** *If  $T = 2$ , then  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$  if and only if  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is a second order stationary point (SOSP) of  $\mathcal{L}$ .*

Thus, local optimization algorithms for finding SOSPs, such as perturbed gradient descent and cubic-regularized Newton method, can efficiently find the minima of the meta-learning objective.

*Proof sketch.* Clearly if  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$ , then  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is an SOSP. The reverse direction is the challenging part of the proof. We equivalently prove that if  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is a critical point and  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \neq 0$ , then  $\nabla^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  has a negative eigenvalue.

Assume for the sake of contradiction that  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is an SOSP and  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \neq 0$ . Considering  $\mathcal{L}$  as a function of the flattened vector  $[\text{vec}(\mathbf{A}); \text{vec}(\mathbf{U}_1); \text{vec}(\mathbf{U}_2)]$ , the idea of the proof is to contradict the assumption that  $\nabla^2 \mathcal{L} \succcurlyeq \mathbf{0}$ .

Since  $\nabla_{\hat{\mathbf{A}}}^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = T\mathbf{I} \succ \mathbf{0}$ , we can work with the Schur complement  $\mathbf{Q} = (\nabla^2 \mathcal{L} / \nabla_{\hat{\mathbf{A}}}^2 \mathcal{L})(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  as  $\nabla^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \succcurlyeq \mathbf{0}$  if and only if  $\mathbf{Q} \succcurlyeq \mathbf{0}$ . Inspection of the condition  $\nabla \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = \mathbf{0}$  along with the assumptions  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \neq 0$  and  $T = 2$  gives three key properties:

$$\left( \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top \right) \mathbf{x} = \left( \mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top} \right) \mathbf{x} \quad \forall \mathbf{x} \text{ s.t. } \left( \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top \right) \mathbf{x} \neq \mathbf{0} \quad (13)$$

$$\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top \neq \mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top} \quad (14)$$

$$\dim \ker \left( \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top \right) > 0 \quad (15)$$

Thus, there is an eigenvector  $\mathbf{z}$  of  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$  with eigenvalue  $\lambda \neq 0$  such that  $\mathbf{z} \in \ker(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top)$ . Assume without loss of generality  $\lambda > 0$ , and consider  $\boldsymbol{\alpha} \in \mathbb{R}^{2k}$ . Define  $g(\cdot; \mathbf{z}) : \mathbb{R}^{2k} \rightarrow \mathbb{R}$  such that  $g(\boldsymbol{\alpha}; \mathbf{z}) = (\boldsymbol{\alpha} \otimes \mathbf{z})^\top \mathbf{Q} (\boldsymbol{\alpha} \otimes \mathbf{z})$ , where  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1; \boldsymbol{\alpha}_2]$  with  $\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2 \in \mathbb{R}^k$ . Then,

$$g(\boldsymbol{\alpha}; \mathbf{z}) = \left\| \hat{\mathbf{U}}_1 \boldsymbol{\alpha}_1 + \hat{\mathbf{U}}_2 \boldsymbol{\alpha}_2 \right\|_2^2 + \lambda \left( \|\boldsymbol{\alpha}_1\|_2^2 - \|\boldsymbol{\alpha}_2\|_2^2 \right). \quad (16)$$

We prove the existence of  $\boldsymbol{\alpha} \in \mathbb{R}^{2k}, \mathbf{x} \in \mathbb{R}^d$  such that  $g(\boldsymbol{\alpha}; \mathbf{x}) < 0$  considering two different cases. Define  $N^- : S_d \rightarrow \mathbb{Z}$  as the function that returns the number of negative eigenvalues of its input.

**Case 1:**  $N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) < k$ : Then there exists  $\mathbf{z}^- \in \mathbb{R}^d$  that is a  $\lambda^-$ -eigenvector of  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$ ,  $\lambda^- < 0$ , where  $\mathbf{z}^- \in \ker(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top)$ . By (15), we can pick  $\boldsymbol{\alpha}$  such that  $\hat{\mathbf{U}}_1 \boldsymbol{\alpha}_1 + \hat{\mathbf{U}}_2 \boldsymbol{\alpha}_2 = \mathbf{0}, \boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2 \neq \mathbf{0}$ . Then  $g(\boldsymbol{\alpha}; \mathbf{z}^-) = -g(\boldsymbol{\alpha}; \mathbf{z}^-) = \|\boldsymbol{\alpha}_1\|_2^2 - \|\boldsymbol{\alpha}_2\|_2^2$ .

If  $\|\bar{\boldsymbol{\alpha}}_1\|_2 \neq \|\bar{\boldsymbol{\alpha}}_2\|_2$ ,  $\min\{g(\boldsymbol{\alpha}; \mathbf{z}), g(\boldsymbol{\alpha}; \mathbf{z}^-)\} < 0$ . Else,  $g(\bar{\boldsymbol{\alpha}}; \mathbf{z}) = 0$ , but  $\nabla_{\boldsymbol{\alpha}_1} g(\bar{\boldsymbol{\alpha}}; \mathbf{z}) = \hat{\mathbf{U}}_1^\top (\hat{\mathbf{U}}_1 \bar{\boldsymbol{\alpha}}_1 + \hat{\mathbf{U}}_2 \bar{\boldsymbol{\alpha}}_2) - 2\lambda \bar{\boldsymbol{\alpha}}_2 = -2\lambda \bar{\boldsymbol{\alpha}}_2 \neq \mathbf{0}$ . Thus there exists  $\bar{\boldsymbol{\alpha}}$  in an infinitesimal neighborhood around  $\boldsymbol{\alpha}$  where  $g(\bar{\boldsymbol{\alpha}}; \mathbf{z}) < 0$ .

**Case 2:**  $N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k$ : By (15),  $\exists \boldsymbol{\Gamma} \in O_k$  such that  $\hat{\mathbf{U}}_2 \boldsymbol{\Gamma} \mathbf{e}_1 \in (\text{im}(\hat{\mathbf{U}}_1) \cap \text{im}(\hat{\mathbf{U}}_2))$ . Define  $\mathbf{y} = \hat{\mathbf{U}}_2 \boldsymbol{\Gamma} \mathbf{e}_1$ . Then

$$k = N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) \geq N^-(\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) \geq N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k.$$

Thus,  $N^-(\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k$  and  $\text{rank}(\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) \leq k$ , so  $\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top \preceq \mathbf{0}$ . Take  $\boldsymbol{\alpha}$  such that  $\hat{\mathbf{U}}_1 \boldsymbol{\alpha}_1 = -\mathbf{y}$  and  $\boldsymbol{\alpha}_2 = \boldsymbol{\Gamma} \mathbf{e}_1$ . Then  $\mathbf{y}_1 \mathbf{y}_1^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top = \hat{\mathbf{U}}_1 (\boldsymbol{\alpha}_1 \boldsymbol{\alpha}_1^\top - \mathbf{I}) \hat{\mathbf{U}}_1^\top \preceq \mathbf{0}$ . Therefore  $\|\boldsymbol{\alpha}_1\|_2 \leq 1$ . Then  $g(\boldsymbol{\alpha}; \mathbf{z}) = \|\hat{\mathbf{U}}_1 \boldsymbol{\alpha}_1 + \hat{\mathbf{U}}_2 \boldsymbol{\alpha}_2\|_2^2 + \lambda(\|\boldsymbol{\alpha}_1\|_2^2 - \|\boldsymbol{\alpha}_2\|_2^2) = \lambda(\|\boldsymbol{\alpha}_1\|_2^2 - 1) \leq 0$ .



If  $g(\alpha; \mathbf{z}) < 0$  we are done. Else, the same analysis from Case 1 will show that  $\nabla g(\bar{\alpha}; \mathbf{z}) \neq \mathbf{0}$ , so there exists  $\bar{\alpha}$  in an infinitesimal neighborhood around  $\alpha$  where  $g(\bar{\alpha}; \mathbf{z}) < 0$ .  $\square$

**Summary.** We have shown that when  $T = 2$ , any optimization algorithm for finding an SOSP will find a global minimum of the meta-objective (11). Surprisingly, when there are three or more tasks, numerical experiments (see Appendix D.2) show that adversarially picking  $U_t^*$  can result in specific instantiations of (11) with spurious local minima. In the next section, we perform extensive numerical experiments for various values of  $T$  which show that these spurious minima are almost never found in practice and vanilla gradient descent is sufficient to minimize (11).

## 4 EXPERIMENTS

### 4.1 LINEAR EXPERIMENTS

To test our algorithm, we perform experiments on a synthetic dataset. We generate  $\mathbf{A}^* \in \mathbb{R}^{d \times d}$  and  $U_t^* \in \mathbb{R}^{d \times k}$  for all tasks  $t \in [T + 1]$ , where the entries of  $\mathbf{A}^*$  and each  $U_t^*$  are i.i.d.  $\mathcal{N}(0, 1)$  random variables. Then we generate  $\frac{N}{T}$  samples for each retraining task  $t \in [T]$  as  $\mathbf{y}_{t,j} = (\mathbf{A}^* + U_t^* U_t^{*\top}) \mathbf{x}_{t,j} + \epsilon_{t,j}$ ,  $j \in [\frac{N}{T}]$  and  $N'$  samples for the held-out task as  $\mathbf{y}_{T+1,j} = (\mathbf{A}^* + U_{T+1}^* U_{T+1}^{*\top}) \mathbf{x}_{T+1,j} + \epsilon_{T+1,j}$ ,  $j \in [N']$ , where  $\mathbf{x}_{t,j} \sim \mathcal{N}(0, \mathbf{I}_d)$  and  $\epsilon_{t,j} \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I}_d)$  are i.i.d. feature and noise vectors respectively.

We apply gradient descent to the Meta-LoRA and standard retraining objectives on the  $T$  retraining tasks and then fine-tune to the  $(T + 1)$ -th task using LoRA. We use symmetric adapters for the Meta-LoRA retraining objective and asymmetric adapters during fine-tuning for each retraining method. We conduct experiments by varying one hyperparameter at a time from the fixed values of  $d = 10$ ,  $T = 3$ ,  $N = 5000$ ,  $N' = 100$ ,  $k = 1$  and  $\sigma = 0.1$ . When  $T = 2$ , we use a rank-3 adaptation during fine-tuning and use a rank-1 adaptation otherwise for both retraining schemes.

We plot the population loss on the test task after training and fine-tuning with Meta-LoRA and SR+LoRA, respectively, in Figure 1. Meta-LoRA significantly outperforms SR+LoRA for all data generation parameter settings. We observe from Figure 1b that with more retraining data, Meta-LoRA performance first improves and then stagnates because of the finite sample noise floor during the fine-tuning stage. We observe a similar phenomenon in Figure 1c. Figure 1d shows that the performance of Meta-LoRA improves for  $T > 2$  relative to  $T = 2$  but is agnostic to  $T$  once in the  $T > 2$  regime. Lastly, Figure 1a shows how performance worsens with increasing dimension.

### 4.2 LLM EXPERIMENTS

To test the Meta-LoRA objective beyond linear models, we perform experiments using the pre-trained 355 million parameter RoBERTa-Large model on the ConvAI2 dataset. ConvAI2 consists of conversations between two personas, i.e. people with different personalities. Each persona is associated with a short list of factual information that guides the content of their responses. We model learning the dialogue continuations of each individual persona as a different task. A training sample for a given persona consists of the previous conversation as input and 20 candidate dialogue continuations, where one of the 20 candidates is the true continuation. We consider the supervised learning task of selecting the correct continuation. During training, we maximize the log-likelihood of the correct continuation and minimize the log-likelihood of each of the incorrect continuations conditioned on the observed conversation history. To run inference given the past conversation and the 20 possible continuations, we select the continuation with the highest conditional likelihood.

For both the standard retraining and Meta-LoRA objectives, we retrain the model using the  $T = 10$  largest retraining tasks, with an average of 117.4 training samples and 36.5 heldout samples per retraining task. We select the model from the epoch with the best average accuracy on the heldout samples and then fine-tune to each of the 10 largest test tasks. For each test task, we take the accuracy on the heldout data from the best performing epoch. We run 5 random trials for this entire retraining and fine-tuning process and report the median best heldout accuracy for each task. All training was done on a single Nvidia A40 GPU, and we report our training hyperparameters in Appendix C.

We compare performance across the test tasks in Table 1. We first minimize the Meta-LoRA objective using rank-8 adapters on the retraining tasks and denote this model Meta-LoRA-8. In Table 1a,

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

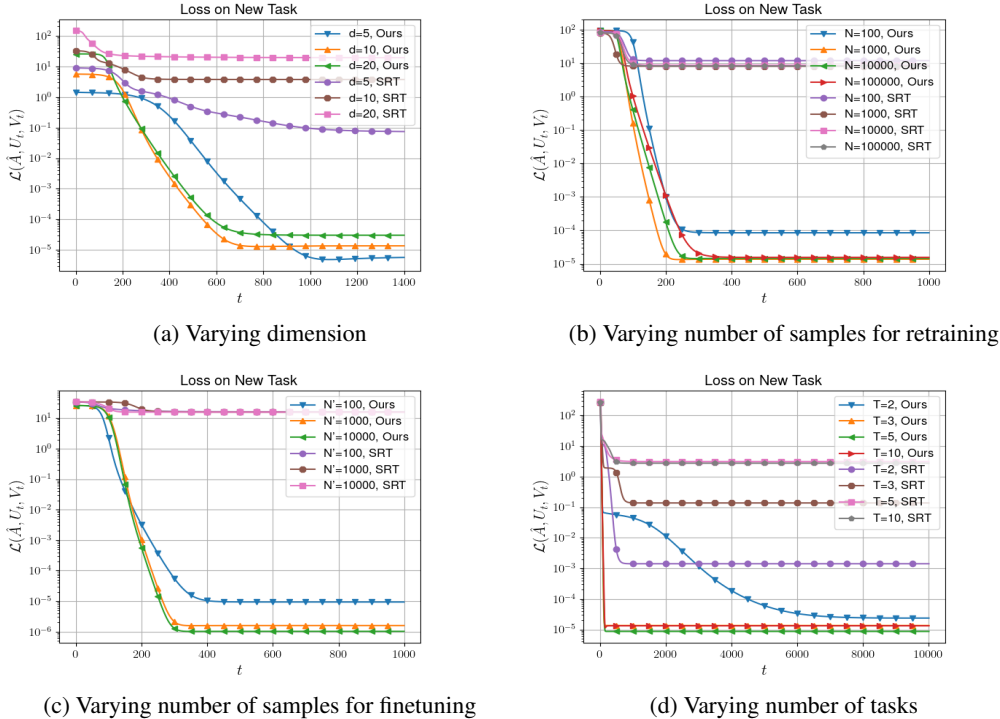


Figure 1: Evaluating the linear Meta-LoRA algorithm in different settings.

we show this improves performance over standard retraining followed by rank-8 LoRA on the test set, denoted SR+LoRA. As suggested by Theorem 2, we test if we can improve performance by increasing the LoRA rank during fine-tuning relative to the rank of the adapters in retraining with Meta-LoRA. Table 1b shows that the retrained Meta-LoRA-8 model fine-tuned with rank-16 adaptations outperforms both standard retraining followed by rank-16 LoRA as well as the Meta-LoRA-16 model which was retrained and fine-tuned with rank-16 adaptations.

Table 1: Comparison of Meta-LoRA and the SRT+LoRA algorithms on the ConvAI2 dataset

(a) Rank-8 fine-tuning adaptations

Algorithm	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6	Task 7	Task 8	Task 9	Task 10	Average
SR+LoRA	43.75	40.00	43.48	41.94	41.03	37.23	42.73	43.20	41.13	40.76	41.52
Meta-LoRA-8	<b>50.00</b>	<b>50.00</b>	<b>47.82</b>	<b>48.39</b>	<b>46.15</b>	<b>41.49</b>	<b>44.55</b>	<b>44.00</b>	<b>42.55</b>	<b>42.68</b>	<b>45.76</b>

(b) Rank-16 fine-tuning adaptations

Algorithm	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6	Task 7	Task 8	Task 9	Task 10	Average
SR+LoRA	43.75	43.33	39.13	38.71	39.74	35.11	38.18	39.20	39.72	38.85	39.57
Meta-LoRA-8	<b>50.0</b>	<b>53.33</b>	<b>50.0</b>	<b>50.0</b>	<b>48.72</b>	<b>42.55</b>	<b>45.45</b>	<b>44.80</b>	<b>45.39</b>	<b>44.59</b>	<b>47.48</b>
Meta-LoRA-16	43.75	33.33	36.96	40.32	43.59	39.36	42.73	41.60	40.43	40.13	40.22

## 5 CONCLUSION

We introduced the Meta-Adapters objective function for retraining an FM on a collection of tasks in a way that prepares the model for subsequent downstream fine-tuning. We provide theoretical justifications on the shortcomings of standard retraining as well as where the Meta-Adapters objective using LoRA (Meta-LoRA) can provably improve performance. Empirically, our Meta-LoRA objective outperforms standard retraining for adapting to unseen downstream tasks. Future avenues include extending our theoretical analysis to finite sample settings and to more general adapters.

## REFERENCES

- 540  
541  
542 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen  
543 Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko,  
544 Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dong-  
545 dong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang  
546 Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit  
547 Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao,  
548 Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin  
549 Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim,  
550 Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden,  
551 Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong  
552 Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro  
553 Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-  
554 Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo  
555 de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim,  
556 Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla,  
557 Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua  
558 Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp  
559 Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Ji-  
560 long Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan,  
561 Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan  
562 Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your  
563 phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- 564 Armen Aghajanyan, Anchit Gupta, Akshat Shrivastava, Xilun Chen, Luke Zettlemoyer, and Sonal  
565 Gupta. Muppet: Massive multi-task representations with pre-finetuning. In Marie-Francine  
566 Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021*  
567 *Conference on Empirical Methods in Natural Language Processing*, pp. 5799–5811, Online  
568 and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguis-  
569 tics. doi: 10.18653/v1/2021.emnlp-main.468. URL <https://aclanthology.org/2021.emnlp-main.468>.
- 570 Trapit Bansal, Salaheddin Alzubi, Tong Wang, Jay-Yoon Lee, and Andrew McCallum. Meta-  
571 adapters: Parameter efficient few-shot fine-tuning through meta-learning. In *First Conference*  
572 *on Automated Machine Learning (Main Track)*, 2022. URL [https://openreview.net/](https://openreview.net/forum?id=BCGNf-prLg5)  
573 [forum?id=BCGNf-prLg5](https://openreview.net/forum?id=BCGNf-prLg5).
- 574 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-  
575 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
576 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
577 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz  
578 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec  
579 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In  
580 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-  
581 ral Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc.,  
582 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)  
583 [file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- 584 Liam Collins, Aryan Mokhtari, Sewoong Oh, and Sanjay Shakkottai. MAML and ANIL provably  
585 learn representations. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,  
586 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine*  
587 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 4238–4310. PMLR, 7  
588 2022. URL <https://proceedings.mlr.press/v162/collins22a.html>.
- 589 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning  
590 of quantized LLMs. In *Thirty-seventh Conference on Neural Information Processing Systems*,  
591 2023. URL <https://openreview.net/forum?id=OUIFPHEgJU>.
- 592  
593 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and

- 594 Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of*  
595 *the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*  
596 *and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Com-  
597 putational Linguistics. doi: 10.18653/v1/N19-1423. URL [https://aclanthology.org/](https://aclanthology.org/N19-1423)  
598 N19-1423.
- 599 Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander H. Miller, Kurt Shuster, Jack Ur-  
600 banek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W.  
601 Black, Alexander I. Rudnicky, Jason Williams, Joelle Pineau, Mikhail S. Burtsev, and Jason We-  
602 ston. The second conversational intelligence challenge (convai2). *ArXiv*, abs/1902.00098, 2019.  
603 URL <https://api.semanticscholar.org/CorpusID:59553505>.
- 604 Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei  
605 Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models  
606 are affected by supervised fine-tuning data composition, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2310.05492)  
607 [abs/2310.05492](https://arxiv.org/abs/2310.05492).
- 608 Simon Shaolei Du, Wei Hu, Sham M. Kakade, Jason D. Lee, and Qi Lei. Few-shot learning via  
609 learning the representation, provably. In *International Conference on Learning Representations*,  
610 2021. URL <https://openreview.net/forum?id=pW2Q2xLwIMD>.
- 611 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of  
612 deep networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International*  
613 *Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp.  
614 1126–1135. PMLR, 8 2017. URL [https://proceedings.mlr.press/v70/finn17a.](https://proceedings.mlr.press/v70/finn17a.html)  
615 [html](https://proceedings.mlr.press/v70/finn17a.html).
- 616 Mozhdeh Gheini, Xuezhe Ma, and Jonathan May. Know where you’re going: Meta-learning for  
617 parameter-efficient fine-tuning, 2022.
- 618 S. K. Hong and Tae Young Jang. AMAL: Meta knowledge-driven few-shot adapter learning. In  
619 Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference*  
620 *on Empirical Methods in Natural Language Processing*, pp. 10381–10389, Abu Dhabi, United  
621 Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/  
622 2022.emnlp-main.709. URL <https://aclanthology.org/2022.emnlp-main.709>.
- 623 Zejiang Hou, Julian Salazar, and George Polovets. Meta-Learning the Difference: Preparing Large  
624 Language Models for Efficient Adaptation. *Transactions of the Association for Computational*  
625 *Linguistics*, 10:1249–1265, 11 2022. ISSN 2307-387X. doi: 10.1162/tacl\_a\_00517. URL  
626 [https://doi.org/10.1162/tacl\\_a\\_00517](https://doi.org/10.1162/tacl_a_00517).
- 627 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An-  
628 drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for  
629 NLP. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th Inter-*  
630 *national Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Re-*  
631 *search*, pp. 2790–2799. PMLR, 6 2019. URL [https://proceedings.mlr.press/v97/](https://proceedings.mlr.press/v97/houlsby19a.html)  
632 [houlsby19a.html](https://proceedings.mlr.press/v97/houlsby19a.html).
- 633 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,  
634 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*  
635 *arXiv:2106.09685*, 2021.
- 636 Nathan Zixia Hu, Eric Mitchell, Christopher D Manning, and Chelsea Finn. Meta-learning online  
637 adaptation of language models. In *The 2023 Conference on Empirical Methods in Natural Lan-*  
638 *guage Processing*, 2023. URL <https://openreview.net/forum?id=jPr118r4RA>.
- 639 Uijeong Jang, Jason D. Lee, and Ernest K. Ryu. Lora training in the ntk regime has no spurious  
640 local minima, 2024.
- 641 Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and  
642 Hannaneh Hajishirzi. UNIFIEDQA: Crossing format boundaries with a single QA system. In  
643 Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Lin-*  
644 *guistics: EMNLP 2020*, pp. 1896–1907, Online, November 2020. Association for Computational

- 648 Linguistics. doi: 10.18653/v1/2020.findings-emnlp.171. URL <https://aclanthology.org/2020.findings-emnlp.171>.
- 649
- 650
- 651 Yoonho Lee and Seungjin Choi. Gradient-based meta-learning with learned layerwise metric and  
652 subspace. In *International Conference on Machine Learning*, 2018. URL <https://api.semanticscholar.org/CorpusID:3350728>.
- 653
- 654 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.  
655 In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th*  
656 *Annual Meeting of the Association for Computational Linguistics and the 11th International Joint*  
657 *Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4582–4597, Online,  
658 August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353.  
659 URL <https://aclanthology.org/2021.acl-long.353>.
- 660
- 661 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-  
662 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv*  
663 *preprint arXiv:2402.09353*, 2024.
- 664
- 665 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike  
666 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining  
667 approach, 2019. URL <http://arxiv.org/abs/1907.11692>. cite arxiv:1907.11692.
- 668
- 669 L. Mirsky. SYMMETRIC GAUGE FUNCTIONS AND UNITARILY INVARIANT NORMS. *The*  
670 *Quarterly Journal of Mathematics*, 11(1):50–59, 01 1960. ISSN 0033-5606. doi: 10.1093/qmath/  
671 11.1.50. URL <https://doi.org/10.1093/qmath/11.1.50>.
- 672
- 673 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
674 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
675 Sutskever. Learning transferable visual models from natural language supervision. In *Internat-*  
676 *ional Conference on Machine Learning*, 2021. URL <https://api.semanticscholar.org/CorpusID:231591445>.
- 677
- 678 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
679 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
680 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 681
- 682 Kiran K Thekumparampil, Prateek Jain, Praneeth Netrapalli, and Sewoong Oh. Statisti-  
683 cally and computationally efficient linear meta-representation learning. In M. Ranzato,  
684 A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neu-*  
685 *ral Information Processing Systems*, volume 34, pp. 18487–18500. Curran Associates, Inc.,  
686 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/  
file/99e7e6ce097324aceb45f98299ceb621-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/99e7e6ce097324aceb45f98299ceb621-Paper.pdf).
- 687
- 688 Yuchen Zeng and Kangwook Lee. The expressive power of low-rank adaptation, 2023.
- 689
- 690 Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and  
691 Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh In-*  
692 *ternational Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=lq62uWRJjiY>.
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701 Yuqian Zhang, Qing Qu, and John Wright. From symmetry to geometry: Tractable nonconvex  
problems, 07 2020.

## A RELATED WORK ON LORA-STYLE PEFT

There is a vast amount of work in developing PEFT methods for FMs. The LoRA algorithm (Hu et al., 2021) has established itself as a popular and successful PEFT strategy and has inspired various extensions such as QLoRA, DoRA, and others (Dettmers et al., 2023; Liu et al., 2024; Zhang et al., 2023). These algorithms are heuristics for mimicking the full finetuning of an FM to a specific downstream task and have proven to be empirically successful in various settings. However, there is a lack of theoretical analysis on the adaptability of PFMs under LoRA-style adaptations, the ability to efficiently optimize LoRA-style objectives, and the kinds of solutions they recover. Some recent works have attempted to analyze different parts of these theoretical questions.

**Convergence of LoRA.** (Jang et al., 2024) analyzes the optimization landscape for LoRA for the Neural Tangent Kernel regime. The authors show that LoRA finetuning converges in this setting as they prove that the objective function satisfies a strict saddle property, ensuring that there are no spurious local minima. However, this focuses on the actual ability of LoRA to converge to the optimal low-rank adapter given an FM, and does not consider the adaptability of the FM in the first place.

**Expressivity of LoRA.** (Zeng & Lee, 2023) derives the expressive power of LoRA as a function of model depth. This work shows that under some mild conditions, fully connected and transformer networks when respectively adapted with LoRA can closely approximate arbitrary smaller networks. They quantify the required LoRA rank to achieve this approximation as well as the resulting approximation error.

## B PROOFS

### B.1 PROOF OF THEOREM 1 AND COROLLARIES 1,2

By definition,

$$\hat{\mathbf{A}}_{\text{SR}} = \arg \min_{\mathbf{A}} \frac{1}{2} \sum_{t=1}^T \left\| \mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{A} \right\|_F^2$$

This optimization problem is just a quadratic function of  $\mathbf{A}$ , so we can simply solve for the point at which the gradient is  $\mathbf{0}$ . Thus,  $\hat{\mathbf{A}}_{\text{SR}}$  must satisfy:

$$\sum_{t=1}^T \left( \mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{A}}_{\text{SR}} \right) = \mathbf{0}$$

Thus,

$$\hat{\mathbf{A}}_{\text{SR}} = \mathbf{A}^* + \frac{1}{T} \sum_{t=1}^T \mathbf{U}_t^* \mathbf{U}_t^{*\top}$$

Therefore,  $\text{rank}(\hat{\mathbf{A}}_{\text{SR}} - \mathbf{A}^*) = \text{rank}\left(\frac{1}{T} \sum_{t=1}^T \mathbf{U}_t^* \mathbf{U}_t^{*\top}\right) = kT$ . Further,

$$\begin{aligned} \mathcal{L}_{\text{Test}}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}_{\text{SR}}) &= \frac{1}{2} \left\| \mathbf{A}^* + \mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{*\top} - \hat{\mathbf{A}}_{\text{SR}} - \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top \right\|_F^2 \\ &= \frac{1}{2} \left\| \mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{*\top} - \frac{1}{T} \sum_{t=1}^T \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top \right\|_F^2 \\ &\approx \frac{1}{2} \left\| \mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{*\top} - k\mathbf{I} - \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top \right\|_F^2 \quad \text{for large } T \end{aligned}$$

$\mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{*\top} - k\mathbf{I}$  has  $d - k$  eigenvalues of magnitude  $k$ , and the rank- $k'$  factorization  $\mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top$  can only capture  $k'$  of them, so  $\mathbf{U}_{T+1}^* \mathbf{U}_{T+1}^{*\top} - k\mathbf{I} - \mathbf{U}_{T+1} \mathbf{V}_{T+1}^\top$  has at least  $d - k' - k$  eigenvalues of magnitude  $k$ . Thus,  $\mathcal{L}_{\text{Test}}(\mathbf{U}_{T+1}, \mathbf{V}_{T+1}; \hat{\mathbf{A}}_{\text{SR}})$  scales as  $(d - k' - k)k^2 \approx (d - k')k^2$  since  $k \ll d$ .

## B.2 PROOF OF THEOREM 2

*Proof.* Since  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$  and  $\mathcal{L} \geq 0$  we must have that  $\nabla_{\mathbf{A}} \mathcal{L} = 0$ .

Thus,  $\hat{\mathbf{A}} = \mathbf{A}^* - \frac{1}{T} \sum_{j=1}^T (\hat{\mathbf{U}}_j \hat{\mathbf{U}}_j^\top - \mathbf{U}_j^* \mathbf{U}_j^{*\top})$ . Plugging this into  $\mathcal{L}$  gives

$$\begin{aligned} 0 = \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) &= \frac{1}{2} \sum_{t=1}^T \left\| \mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \left( \mathbf{A}^* - \frac{1}{T} \sum_{s=1}^T (\hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top - \mathbf{U}_s^* \mathbf{U}_s^{*\top}) \right) - \mathbf{U}_t \mathbf{U}_t^\top \right\|_F^2 \\ &= \frac{1}{2} \sum_{t=1}^T \left\| \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{U}_t \mathbf{U}_t^\top - \frac{1}{T} \sum_{s=1}^T (\hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top - \mathbf{U}_s^* \mathbf{U}_s^{*\top}) \right\|_F^2. \end{aligned}$$

Thus each term of the summation is zero, so for all  $t, s \in [T]$ ,

$$\hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top - \mathbf{U}_t^* \mathbf{U}_t^{*\top} = \hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top - \mathbf{U}_s^* \mathbf{U}_s^{*\top}.$$

Combining these results gives that

$$\begin{aligned} \hat{\mathbf{A}} &= \mathbf{A}^* - \frac{1}{T} \sum_{s=1}^T (\hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top - \mathbf{U}_s^* \mathbf{U}_s^{*\top}) \\ &= \mathbf{A}^* - (\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top - \mathbf{U}_1^* \mathbf{U}_1^{*\top}) \end{aligned}$$

Let  $\mathbf{C} = -(\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top - \mathbf{U}_1^* \mathbf{U}_1^{*\top})$ . Then  $\hat{\mathbf{A}} = \mathbf{A}^* + \mathbf{C}$  and  $\text{rank}(\mathbf{C}) \leq \text{rank}(\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) + \text{rank}(\mathbf{U}_1^* \mathbf{U}_1^{*\top}) \leq 2k$

□

## B.3 PROOF OF THEOREM 3

*Proof.* Since  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$ , we have that for all  $t, s \in [T]$ ,

$$\hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top - \mathbf{U}_t^* \mathbf{U}_t^{*\top} = \hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top - \mathbf{U}_s^* \mathbf{U}_s^{*\top} \quad (17)$$

Applying this to the first three tasks and rearranging gives that

$$\mathbf{U}_1^* \mathbf{U}_1^{*\top} = \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top + \mathbf{U}_2^* \mathbf{U}_2^{*\top} - \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top \quad (18)$$

$$= \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top + \mathbf{U}_3^* \mathbf{U}_3^{*\top} - \hat{\mathbf{U}}_3 \hat{\mathbf{U}}_3^\top. \quad (19)$$

We first show that  $\text{im}(\hat{\mathbf{U}}_1) = \text{im}(\mathbf{U}_1^*)$ .

Since  $\mathbf{U}_1^* \mathbf{U}_1^{*\top} \succcurlyeq 0$ , we must have that  $\text{im}(\hat{\mathbf{U}}_2) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_2^*)$  and  $\text{im}(\hat{\mathbf{U}}_3) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_3^*)$ , as otherwise there would exist a vector on  $\ker(\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top + \mathbf{U}_2^* \mathbf{U}_2^{*\top}) \cap \ker(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top)^\perp$  whose existence contradicts the positive semi-definiteness of  $\mathbf{U}_1^* \mathbf{U}_1^{*\top}$ .

Thus,

$$\text{im}(\mathbf{U}_1^*) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_2^*) \quad (20)$$

$$\text{im}(\mathbf{U}_1^*) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_3^*) \quad (21)$$

Using that fact that for subspaces  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ ,  $\mathbf{X} \subseteq \mathbf{Y} \implies \mathbf{X} + \mathbf{Z} \subseteq \mathbf{Y} + \mathbf{Z}$ , we can add  $\text{im}(\mathbf{U}_2^*)$  and  $\text{im}(\mathbf{U}_3^*)$  to both sides of 20 and 21 respectively. This gives that

$$\text{im}(\mathbf{U}_1^*) \oplus \text{im}(\mathbf{U}_2^*) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_2^*) \quad (22)$$

$$\text{im}(\mathbf{U}_1^*) \oplus \text{im}(\mathbf{U}_3^*) \subseteq \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\mathbf{U}_3^*). \quad (23)$$

For  $t \in \{2, 3\}$ , we clearly have that  $\dim(\text{im}(\hat{U}_1) + \text{im}(U_t^*)) \leq \dim \text{im}(\hat{U}_1) + \dim \text{im}(U_t^*) \leq 2k$ , and  $\dim(\text{im}(U_1^*) + \text{im}(U_t^*)) = 2k$ . Thus,

$$(\text{im}(U_1^*) \oplus \text{im}(U_2^*)) = (\text{im}(\hat{U}_1) \oplus \text{im}(U_2^*)) \quad (24)$$

$$(\text{im}(U_1^*) \oplus \text{im}(U_3^*)) = (\text{im}(\hat{U}_1) \oplus \text{im}(U_3^*)) \quad (25)$$

**Lemma 2.**  $([\text{im}(\hat{U}_1) \oplus \text{im}(U_2^*)] \cap [\text{im}(\hat{U}_1) \oplus \text{im}(U_3^*)]) = \text{im}(\hat{U}_1)$

*Proof.* Clearly,  $\text{im}(\hat{U}_1) \subseteq ([\text{im}(\hat{U}_1) \oplus \text{im}(U_2^*)] \cap [\text{im}(\hat{U}_1) \oplus \text{im}(U_3^*)])$ . To show the converse, consider  $x \in ([\text{im}(\hat{U}_1) \oplus \text{im}(U_2^*)] \cap [\text{im}(\hat{U}_1) \oplus \text{im}(U_3^*)])$ .

By assumption there exists some  $a, b, c, d \in \mathbb{R}^k$  such that

$$x = \hat{U}_1 a + U_2^* b \quad (26)$$

$$= \hat{U}_1 c + U_3^* d \quad (27)$$

Thus,

$$\hat{U}_1(a - c) + U_2^* b - U_3^* d = \mathbf{0}. \quad (28)$$

By Equation 24, we can write

$$\begin{aligned} \text{im}(U_2^*) &= ([\text{im}(U_1^*) \oplus \text{im}(U_2^*)] \cap [\text{im}(U_2^*) \oplus \text{im}(U_3^*)]) \\ &= ([\text{im}(\hat{U}_1) \oplus \text{im}(U_2^*)] \cap [\text{im}(U_2^*) \oplus \text{im}(U_3^*)]) \end{aligned}$$

Thus,  $\text{im}(\hat{U}_1) \cap [\text{im}(U_2^*) \oplus \text{im}(U_3^*)] \subseteq \text{im}(\hat{U}_1) \cap \text{im}(U_2^*) = \{\mathbf{0}\}$ , so

$$\text{im}(\hat{U}_1) \cap [\text{im}(U_2^*) \oplus \text{im}(U_3^*)] = \{\mathbf{0}\} \quad (29)$$

Applying Equation (29) to Equation (28) implies that  $a = c$  and  $b = d = \mathbf{0}$ . Thus  $x = \hat{U}_1 a \in \text{im}(\hat{U}_1)$ , so  $([\text{im}(\hat{U}_1) \oplus \text{im}(U_2^*)] \cap [\text{im}(\hat{U}_1) \oplus \text{im}(U_3^*)]) \subseteq \text{im}(\hat{U}_1)$ .  $\square$

Then Equations (20) and (21) combined with Lemma (2) implies that  $\text{im}(U_1^*) \subseteq \text{im}(\hat{U}_1)$  but  $\dim(\text{im}(U_1^*)) = \dim(\text{im}(\hat{U}_1)) = k$ , so  $\text{im}(U_1^*) = \text{im}(\hat{U}_1)$ .

Since the initial assumptions about  $\hat{U}_1$  and  $U_1^*$  analogously hold for the corresponding matrices for tasks 2 and 3, by the exact same argument we can show that

$$\text{im}(U_t^*) = \text{im}(\hat{U}_t) \quad \forall t \in [T]. \quad (30)$$

Then by equation (17),  $\text{im}(U_1^*) \supseteq \text{im}(\hat{U}_1 \hat{U}_1^T - U_1^* U_1^{*T}) = \text{im}(\hat{U}_2 \hat{U}_2^T - U_2^* U_2^{*T}) \subseteq \text{im}(U_2^*)$ . Thus,

$$\begin{aligned} \text{im}(\hat{U}_1 \hat{U}_1^T - U_1^* U_1^{*T}) &\subseteq \text{im}(U_1^*) \cap \text{im}(U_2^*) \\ &= \{\mathbf{0}\}. \end{aligned}$$

Thus  $\hat{U}_1 \hat{U}_1^T = U_1^* U_1^{*T}$ . Then by Equation (17),  $\hat{U}_t \hat{U}_t^T = U_t^* U_t^{*T}$  for all  $t \in [T]$ . Lastly, since  $\mathcal{L}(\hat{A}, \hat{U}) = 0$ , we have that  $\nabla_{\mathcal{A}} \mathcal{L}(\hat{A}, \hat{U}) = 0$ , so

$$\hat{A} = A^* + \frac{1}{T} \sum_{t=1}^T U_t^* U_t^{*T} - U_t U_t^T = A^*$$

$\square$



## B.4 PROOF OF THEOREM 4

Clearly if  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$ , then  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is an SOSP. The reverse direction is the challenging part of the proof. We equivalently prove that if  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is a critical point and  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \neq 0$ , then  $\nabla^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  has a negative eigenvalue.

Assume for the sake of contradiction that  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is a critical point and  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \neq 0$ . Then,

$$\nabla_{\mathbf{A}} \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = T(\hat{\mathbf{A}} - \mathbf{A}^*) + \sum_{t=1}^T (\hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top - \mathbf{U}_t^* \mathbf{U}_t^{*\top}) = \mathbf{0} \quad (31)$$

$$\nabla_{\mathbf{U}_t} \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 2(\hat{\mathbf{A}} - \mathbf{A}^* + \hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top - \mathbf{U}_t^* \mathbf{U}_t^{*\top}) \hat{\mathbf{U}}_t = \mathbf{0} \quad (32)$$

Thus,

$$\hat{\mathbf{A}} = \mathbf{A}^* - \frac{1}{T} \sum_{t=1}^T (\hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top - \mathbf{U}_t^* \mathbf{U}_t^{*\top}). \quad (33)$$

Define  $\mathbf{B}_t(\hat{\mathbf{U}}) = \hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top - \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \frac{1}{T} \sum_{s=1}^T (\hat{\mathbf{U}}_s \hat{\mathbf{U}}_s^\top - \mathbf{U}_s^* \mathbf{U}_s^{*\top})$ . Despite being a slight abuse of notation, we refer to  $\mathbf{B}_t(\hat{\mathbf{U}})$  as just  $\mathbf{B}_t$  for the remainder of the proof.

Then (32) equivalently states:

$$\mathbf{B}_t \hat{\mathbf{U}}_t = \mathbf{0}. \quad (34)$$

Note that by construction,  $\sum_{t=1}^T \mathbf{B}_t = \mathbf{0}$ .

Considering  $\mathcal{L}$  as a function of the flattened vector  $[\text{vec}(\mathbf{A}); \text{vec}(\mathbf{U}_1); \text{vec}(\mathbf{U}_2)]$ , and let  $\mathbf{U}_1 = [\mathbf{x}_1 \dots \mathbf{x}_k]$ ,  $\mathbf{U}_2 = [\mathbf{y}_1 \dots \mathbf{y}_k]$ , we compute the Hessian

$$\nabla^2 \mathcal{L} = \begin{bmatrix} \nabla_{\mathbf{A}}^2 \mathcal{L} & \nabla_{\mathbf{U}_1} \nabla_{\mathbf{A}} \mathcal{L} & \nabla_{\mathbf{U}_2} \nabla_{\mathbf{A}} \mathcal{L} \\ (\nabla_{\mathbf{U}_1} \nabla_{\mathbf{A}} \mathcal{L})^\top & \nabla_{\mathbf{U}_1}^2 \mathcal{L} & \mathbf{0} \\ (\nabla_{\mathbf{U}_2} \nabla_{\mathbf{A}} \mathcal{L})^\top & \mathbf{0} & \nabla_{\mathbf{U}_2}^2 \mathcal{L} \end{bmatrix} \quad (35)$$

where

$$\nabla_{\mathbf{A}}^2 \mathcal{L} = 2\mathbf{I}_{d^2}$$

$$\nabla_{\mathbf{U}_1} \nabla_{\mathbf{A}} \mathcal{L} = [(\mathbf{x}_1 \oplus \mathbf{x}_1) \dots (\mathbf{x}_k \oplus \mathbf{x}_k)] \in \mathbb{R}^{d^2 \times dk}$$

$$\nabla_{\mathbf{U}_2} \nabla_{\mathbf{A}} \mathcal{L} = [(\mathbf{y}_1 \oplus \mathbf{y}_1) \dots (\mathbf{y}_k \oplus \mathbf{y}_k)] \in \mathbb{R}^{d^2 \times dk}$$

$$\nabla_{\mathbf{U}_1}^2 \mathcal{L} = 2(\mathbf{A} + \mathbf{U}_1 \mathbf{U}_1^\top - \mathbf{A}^* - \mathbf{U}_1^* \mathbf{U}_1^{*\top}) \otimes \mathbf{I}_k$$

$$+ 2 \begin{bmatrix} \mathbf{x}_1 \mathbf{x}_1^\top + \|\mathbf{x}_1\|_2^2 \mathbf{I} & \mathbf{x}_1^\top \mathbf{x}_2 \mathbf{I} + \mathbf{x}_2 \mathbf{x}_1^\top & \dots & \mathbf{x}_1^\top \mathbf{x}_k \mathbf{I} + \mathbf{x}_k \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \mathbf{x}_1 \mathbf{I} + \mathbf{x}_1 \mathbf{x}_2^\top & \mathbf{x}_2 \mathbf{x}_2^\top + \|\mathbf{x}_2\|_2^2 \mathbf{I} & \dots & \mathbf{x}_2^\top \mathbf{x}_k \mathbf{I} + \mathbf{x}_k \mathbf{x}_2^\top \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_k^\top \mathbf{x}_1 \mathbf{I} + \mathbf{x}_1 \mathbf{x}_k^\top & \dots & \dots & \mathbf{x}_k \mathbf{x}_k^\top + \|\mathbf{x}_k\|_2^2 \mathbf{I} \end{bmatrix}$$

$$\nabla_{\mathbf{U}_2}^2 \mathcal{L} = 2(\mathbf{A} + \mathbf{U}_2 \mathbf{U}_2^\top - \mathbf{A}^* - \mathbf{U}_2^* \mathbf{U}_2^{*\top}) \otimes \mathbf{I}_k$$

$$+ 2 \begin{bmatrix} \mathbf{y}_1 \mathbf{y}_1^\top + \|\mathbf{y}_1\|_2^2 \mathbf{I} & \mathbf{y}_1^\top \mathbf{y}_2 \mathbf{I} + \mathbf{y}_2 \mathbf{y}_1^\top & \dots & \mathbf{y}_1^\top \mathbf{y}_k \mathbf{I} + \mathbf{y}_k \mathbf{y}_1^\top \\ \mathbf{y}_2^\top \mathbf{y}_1 \mathbf{I} + \mathbf{y}_1 \mathbf{y}_2^\top & \mathbf{y}_2 \mathbf{y}_2^\top + \|\mathbf{y}_2\|_2^2 \mathbf{I} & \dots & \mathbf{y}_2^\top \mathbf{y}_k \mathbf{I} + \mathbf{y}_k \mathbf{y}_2^\top \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{y}_k^\top \mathbf{y}_1 \mathbf{I} + \mathbf{y}_1 \mathbf{y}_k^\top & \dots & \dots & \mathbf{y}_k \mathbf{y}_k^\top + \|\mathbf{y}_k\|_2^2 \mathbf{I} \end{bmatrix}$$

Note that  $\oplus$  denotes the Kronecker sum defined as  $\mathbf{X} \oplus \mathbf{Y} = \mathbf{I} \otimes \mathbf{X} + \mathbf{Y} \otimes \mathbf{I}$  where  $\otimes$  is the Kronecker product.

**Lemma 3.**  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$  if and only if  $\mathbf{B}_t = \mathbf{0}$  for each  $t \in [T]$ .

*Proof.* Since  $(\hat{\mathbf{A}}, \hat{\mathbf{U}})$  is a critical point, then plugging Equation (33) into the definition of  $\mathcal{L}$  gives that

$$\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = \frac{1}{2} \sum_{t=1}^T \|\mathbf{B}_t\|_F^2.$$

Thus  $\mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) = 0$  if and only if  $\mathbf{B}_t = \mathbf{0} \quad \forall t$ .  $\square$

**Lemma 4.** If  $\nabla_{\mathbf{U}}^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \succcurlyeq \mathbf{0}$ , then the eigenvectors corresponding to the non-zero eigenvalues of  $\hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top$  are the leading non-negative eigenvectors of  $\mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{A}}$  for all  $t \in [T]$ .

*Proof.* Consider the function  $\bar{f}_t(\mathbf{U}_t; \hat{\mathbf{A}}) = \frac{1}{2} \left\| \mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{A}} - \mathbf{U}_t \mathbf{U}_t^\top \right\|_F^2$ .  $\bar{f}_t$  is simply the  $t$ th summand in  $\mathcal{L}$  where  $\mathbf{A} = \hat{\mathbf{A}}$  is fixed and we only consider the variable  $\mathbf{U}_t$ . Minimising  $\bar{f}_t$  is identical to the problem of symmetric matrix factorization.

Using well-known properties of symmetric matrix factorization, since  $\nabla \bar{f}_t(\hat{\mathbf{U}}_t) = \mathbf{0}$ , we must have that  $\hat{\mathbf{U}}_t = \mathbf{V}_t \mathbf{\Gamma}$  where the columns of  $\mathbf{V}_t$  are the properly scaled eigenvectors of  $\mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{A}}$  with non-negative eigenvalues where each column has norm equal to the square root of its corresponding eigenvalue, and  $\mathbf{\Gamma} \in O_k$  is some orthogonal matrix. Further, if the eigenvectors corresponding to the non-zero eigenvalues of  $\hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top$  are not the leading non-negative eigenvectors, then  $\nabla^2 \bar{f}_t(\hat{\mathbf{U}}_t) \not\succeq \mathbf{0}$  by (Zhang et al., 2020). Since  $\nabla^2 \bar{f}_t(\hat{\mathbf{U}}_t)$  is a diagonal block of  $\nabla^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}})$ ,  $\nabla^2 \bar{f}_t(\hat{\mathbf{U}}_t) \not\succeq \mathbf{0}$  would imply  $\nabla^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \not\succeq \mathbf{0}$ .  $\square$

**Remark 2.** Without loss of generality, we can assume that the eigenvectors corresponding to the non-zero eigenvalues of  $\hat{\mathbf{U}}_t \hat{\mathbf{U}}_t^\top$  are the leading non-negative eigenvectors of  $\mathbf{A}^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \hat{\mathbf{A}}$  for all  $i$ .

**Lemma 5.**  $(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) \mathbf{x} = (\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}) \mathbf{x}$  for all  $\mathbf{x} \in \text{im}(\hat{\mathbf{U}}_1) + \text{im}(\hat{\mathbf{U}}_2)$ .

*Proof.* Recall  $\mathbf{B}_1 = \frac{1}{2} (\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top - \mathbf{U}_1^* \mathbf{U}_1^{*\top} - \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top + \mathbf{U}_2^* \mathbf{U}_2^{*\top})$ . Then applying first-order stationarity and the fact that  $\mathbf{B}_2 = -\mathbf{B}_1$ , we have

$$\begin{aligned} (\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) \hat{\mathbf{U}}_1 &= (\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}) \hat{\mathbf{U}}_1 \\ (\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) \hat{\mathbf{U}}_2 &= (\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}) \hat{\mathbf{U}}_2. \end{aligned}$$

$\square$

**Corollary 5.**  $\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top$  and  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$  share an eigenbasis.

*Proof.* Using the lemma, any non-zero eigenvector-eigenvalue pair of  $\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top$  is also an eigenvector-eigenvalue pair of  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$ . Denote the space defined by the span of these eigenvectors as  $\mathcal{S}$ . Then all other eigenvectors of  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$  are orthogonal to  $\mathcal{S}$ , so they are also 0-eigenvectors of  $\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top$ . Thus the two matrices share an eigenbasis.  $\square$

**Corollary 6.**  $\dim(\text{im} \hat{\mathbf{U}}_1 + \text{im} \hat{\mathbf{U}}_2) \leq 2k - 1$ , i.e., the set of columns of  $\hat{\mathbf{U}}_1$  and  $\hat{\mathbf{U}}_2$  are not linearly independent.

*Proof.* Assume for contradiction that the vectors in the set  $\mathcal{S} = \{\hat{\mathbf{U}}_1 \mathbf{e}_i \mid i = 1, \dots, k\} \cup \{\hat{\mathbf{U}}_2 \mathbf{e}_i \mid i = 1, \dots, k\}$  are linearly independent, where  $\mathbf{e}_i$  is the  $i$ th standard basis vector in  $\mathbb{R}^k$ .

Then note that  $(\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top - \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top) \mathbf{x} \neq \mathbf{0}$  and  $(\mathbf{U}_1^* \mathbf{U}_1^{*\top} - \mathbf{U}_2^* \mathbf{U}_2^{*\top}) \mathbf{x} \neq \mathbf{0}$  for all  $\mathbf{x} \in \mathcal{S}$ . By Lemma (5),  $\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top - \hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top$  and  $\mathbf{U}_1^* \mathbf{U}_1^{*\top} - \mathbf{U}_2^* \mathbf{U}_2^{*\top}$  agree for each vector on the  $2k$ -dimensional

space  $\text{span}(\mathcal{S})$ . But, both  $\text{rank}(\hat{U}_1\hat{U}_1^\top - \hat{U}_2\hat{U}_2^\top)$ ,  $\text{rank}(U_1^*U_1^{*\top} - U_2^*U_2^{*\top}) \leq 2k$  by construction. Then by dimension counting,  $\hat{U}_1\hat{U}_1^\top - \hat{U}_2\hat{U}_2^\top$  and  $U_1^*U_1^{*\top} - U_2^*U_2^{*\top}$  must send  $\text{span}\{\mathcal{S}\}^\perp$  to  $\mathbf{0}$ . Thus,  $\hat{U}_1\hat{U}_1^\top - \hat{U}_2\hat{U}_2^\top$  and  $U_1^*U_1^{*\top} - U_2^*U_2^{*\top}$  agree on the entire basis formed by concatenating basis vectors of  $\text{span}\{\mathcal{S}\}^\perp$  with those of  $\text{span}(\mathcal{S})$ . This implies that  $\hat{U}_1\hat{U}_1^\top - \hat{U}_2\hat{U}_2^\top = U_1^*U_1^{*\top} - U_2^*U_2^{*\top}$  and thus  $B_1 = \hat{U}_1\hat{U}_1^\top - \hat{U}_2\hat{U}_2^\top - U_1^*U_1^{*\top} + U_2^*U_2^{*\top} = \mathbf{0}$ . Then  $B_2 = -B_1 = \mathbf{0}$  so by Lemma 3,  $\mathcal{L}(\hat{A}, \hat{U}) = 0$  which is a contradiction.  $\square$

**Lemma 6.**  $U_2^*U_2^{*\top} - U_1^*U_1^{*\top}$  has exactly  $k$  positive and  $k$  negative eigenvalues.

*Proof.* First, note that  $U_2^*U_2^{*\top}$  has exactly  $k$  positive eigenvalues and  $k - d$  eigenvalues of  $\mathbf{0}$ . Then  $U_2^*U_2^{*\top} - (U_1^*e_1)(U_1^*e_1)^\top$  has rank  $k + 1$  because of the linear independence of the columns of the combined set of columns  $U_1^*$  and  $U_2^*$ . Further, since we subtract  $(U_1^*e_1)(U_1^*e_1)^\top$ , we must be accumulating an additional negative eigenvalue relative to  $U_2^*U_2^{*\top}$ . Continuing this process shows that subtracting  $(U_1^*e_{j+1})(U_1^*e_{j+1})^\top$  from  $U_2^*U_2^{*\top} - \sum_{i=1}^j (U_1^*e_i)(U_1^*e_i)^\top$  contributes exactly one more negative eigenvalue, since  $U_1^*e_{j+1}$  can never be written as a linear combination of  $\{U_1^*e_1, \dots, U_1^*e_k, U_2^*e_1, \dots, U_2^*e_j\}$  for  $0 < j < k$ . The result then follows from induction.  $\square$

**Lemma 7.**  $\text{rank}(\hat{U}_1) = \text{rank}(\hat{U}_2) = k$ .

*Proof.* Assume for contradiction that  $\text{rank}(\hat{U}_1) = m < k$  without loss of generality. Since by Remark (2) we assume the columns of  $\hat{U}_1$  are the leading  $k$  non-negative eigenvectors of  $A^* + U_1^*U_1^{*\top} - \hat{A} = \hat{U}_1\hat{U}_1^\top - B_1$ , this must imply that  $A^* + U_1^*U_1^{*\top} - \hat{A} - \hat{U}_1\hat{U}_1^\top = -B_1 \preceq \mathbf{0}$ .

Plugging in the definition of  $B_1$  gives that  $\frac{1}{2}(\hat{U}_1\hat{U}_1^\top - U_1^*U_1^{*\top} - \hat{U}_2\hat{U}_2^\top + U_2^*U_2^{*\top}) \succcurlyeq \mathbf{0}$ . Thus,  $\hat{U}_1\hat{U}_1^\top \succcurlyeq U_1^*U_1^{*\top} + \hat{U}_2\hat{U}_2^\top - U_2^*U_2^{*\top} \succcurlyeq U_1^*U_1^{*\top} - U_2^*U_2^{*\top}$ . This contradicts the fact from Lemma (6) that  $U_1^*U_1^{*\top} - U_2^*U_2^{*\top}$  has  $k$  positive eigenvalues.  $\square$

With this lemma, we will prove the existence of a direction of  $\nabla^2\mathcal{L}$  with negative curvature. Instead of directly working with this matrix, we instead use the Schur complement to work with a different form.

**Theorem 5. (Schur Complement)** Since  $\nabla_A^2\mathcal{L}(\hat{A}, \hat{U}) = 2I \succ \mathbf{0}$ ,  $\nabla^2\mathcal{L}(\hat{A}, \hat{U}) \succcurlyeq \mathbf{0}$  if and only if  $\nabla_U^2\mathcal{L}(\hat{A}, \hat{U}) - (\nabla_A\nabla_U\mathcal{L}(\hat{A}, \hat{U}))(\nabla_A^2\mathcal{L}(\hat{A}, \hat{U}))^{-1}(\nabla_U\nabla_A\mathcal{L}(\hat{A}, \hat{U})) \succcurlyeq \mathbf{0}$ .

Define  $Q = \nabla_U^2\mathcal{L}(\hat{A}, \hat{U}) - (\nabla_A\nabla_U\mathcal{L}(\hat{A}, \hat{U}))(\nabla_A^2\mathcal{L}(\hat{A}, \hat{U}))^{-1}(\nabla_U\nabla_A\mathcal{L}(\hat{A}, \hat{U}))$ .

For example, when  $k = 2$  and letting  $U_1 = [x_1 \ x_2]$ ,  $U_2 = [y_1 \ y_2]$ , we have

$$Q = \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{12}^\top & Q_{22} \end{bmatrix},$$

where

$$Q_{11} = \begin{bmatrix} 2B_1 + x_1x_1^\top + \|x_1\|_2^2 & x_1^\top x_2 I + x_2x_1^\top \\ x_2^\top x_1 I + x_1x_2^\top & 2B_1 + x_2x_2^\top + \|x_2\|_2^2 \end{bmatrix}$$

$$Q_{12} = \begin{bmatrix} -x_1^\top y_1 I - y_1x_1^\top & -x_1^\top y_2 I - y_2x_1^\top \\ -x_2^\top y_1 I - y_1x_2^\top & x_2^\top y_2 I - y_2x_2^\top \end{bmatrix}$$

$$Q_{22} = \begin{bmatrix} 2B_2 + y_1y_1^\top + \|y_1\|_2^2 & y_1^\top y_2 I + y_2y_1^\top \\ y_2^\top y_1 I + y_1y_2^\top & 2B_2 + y_2y_2^\top + \|y_2\|_2^2 \end{bmatrix}$$

For brevity, we do not include the full form of  $Q$  for general  $k$ . However, we can make an easy simplification that will allow for a much cleaner expression.

Using Corollaries (5) and (6), there is an eigenvector  $\mathbf{z}$  of  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$  with eigenvalue  $\lambda \neq 0$  such that  $\mathbf{z} \in \ker(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top)$ . Assume without loss of generality that  $\lambda > 0$ , and consider  $\boldsymbol{\alpha} \in \mathbb{R}^{2k}$ . Define the function  $g(\cdot; \mathbf{z}) : \mathbb{R}^{2k} \rightarrow \mathbb{R}$  parameterized by  $\mathbf{z}$  such that  $g(\boldsymbol{\alpha}; \mathbf{z}) = (\boldsymbol{\alpha} \otimes \mathbf{z})^\top \mathbf{Q}(\boldsymbol{\alpha} \otimes \mathbf{z})$ , where we partition  $\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1; \boldsymbol{\alpha}_2]$ ,  $\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2 \in \mathbb{R}^k$ . Then after some algebra,

$$g(\boldsymbol{\alpha}; \mathbf{z}) = \left\| \hat{\mathbf{U}}_1 \boldsymbol{\alpha}_1 + \hat{\mathbf{U}}_2 \boldsymbol{\alpha}_2 \right\|_2^2 + \lambda \left( \|\boldsymbol{\alpha}_1\|_2^2 - \|\boldsymbol{\alpha}_2\|_2^2 \right). \quad (36)$$

We prove the existence of  $\boldsymbol{\alpha} \in \mathbb{R}^{2k}$ ,  $\mathbf{x} \in \mathbb{R}^d$  such that  $g(\boldsymbol{\alpha}; \mathbf{x}) < 0$  considering two different cases. Define  $N^- : S_d \rightarrow \mathbb{Z}$  as the function that returns the number of negative eigenvalues of its input.

**Case 1:**  $N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) < k$ .

Using Corollary (6), we can pick  $\boldsymbol{\alpha}$  such that  $\hat{\mathbf{U}}_1 \boldsymbol{\alpha}_1 + \hat{\mathbf{U}}_2 \boldsymbol{\alpha}_2 = 0$ ,  $\boldsymbol{\alpha}_1, \boldsymbol{\alpha}_2 \neq 0$ .

Because  $N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) < k$ ,  $N^-(\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}) = k$ , and  $\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top$  and  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$  share an eigenbasis by Corollary 5, there exists  $\mathbf{z}^- \in \mathbb{R}^d$  that is a  $\lambda^-$ -eigenvector of  $\mathbf{U}_2^* \mathbf{U}_2^{*\top} - \mathbf{U}_1^* \mathbf{U}_1^{*\top}$ ,  $\lambda^- < 0$ , where  $\mathbf{z} \in \ker(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top)$

Then for the same choice of  $\boldsymbol{\alpha}$ ,

$$\begin{aligned} \text{sign}(g(\boldsymbol{\alpha}; \mathbf{z})) &= \text{sign}\left(\|\boldsymbol{\alpha}_1\|_2^2 - \|\boldsymbol{\alpha}_2\|_2^2\right) \\ \text{sign}(g(\boldsymbol{\alpha}; \mathbf{z}^-)) &= \text{sign}\left(\|\boldsymbol{\alpha}_2\|_2^2 - \|\boldsymbol{\alpha}_1\|_2^2\right). \end{aligned}$$

Then if  $\|\boldsymbol{\alpha}_1\|_2 \neq \|\boldsymbol{\alpha}_2\|_2$ , one of the above expressions is negative and thus  $\mathbf{Q}$  has a negative eigenvalue. This then implies  $\nabla^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \not\preceq 0$ .

Otherwise  $\|\boldsymbol{\alpha}_1\|_2 = \|\boldsymbol{\alpha}_2\|_2$ . Then  $g(\boldsymbol{\alpha}; \mathbf{z}) = 0$ , but  $\nabla_{\boldsymbol{\alpha}_1} g(\boldsymbol{\alpha}; \mathbf{z}) = \hat{\mathbf{U}}_1^\top (\hat{\mathbf{U}}_1 \bar{\boldsymbol{\alpha}}_1 + \hat{\mathbf{U}}_2 \boldsymbol{\alpha}_2) - 2\lambda \boldsymbol{\alpha}_2 = -2\lambda \boldsymbol{\alpha}_2 \neq 0$ . Thus  $g(\boldsymbol{\alpha}; \mathbf{z}) = 0$  and  $\nabla g(\boldsymbol{\alpha}; \mathbf{z}) \neq 0$  so there exists  $\bar{\boldsymbol{\alpha}}$  in an infinitesimal neighborhood around  $\boldsymbol{\alpha}$  where  $g(\bar{\boldsymbol{\alpha}}; \mathbf{z}) < 0$ . Thus  $\mathbf{Q}$  has a negative eigenvalue so  $\nabla^2 \mathcal{L}(\hat{\mathbf{A}}, \hat{\mathbf{U}}) \not\preceq 0$ .

**Case 2:**  $N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k$ .

Define  $m = \dim(\text{im}(\hat{\mathbf{U}}_1) \cap \text{im}(\hat{\mathbf{U}}_2))$ . By Corollary 6,  $m \geq 1$ , so we can select orthogonal matrix  $\boldsymbol{\Gamma} \in O_k$  such that  $\hat{\mathbf{U}}_2 \boldsymbol{\Gamma} \mathbf{e}_1 \in (\text{im}(\hat{\mathbf{U}}_1) \cap \text{im}(\hat{\mathbf{U}}_2))$ . Define  $\mathbf{y} = \hat{\mathbf{U}}_2 \boldsymbol{\Gamma} \mathbf{e}_1$ .

Clearly for any  $\mathbf{B} \in S_d$  and  $\mathbf{R} \in S_d^+$ ,  $N^-(\mathbf{B}) \geq N^-(\mathbf{B} + \mathbf{R})$ . Then since  $N^-(\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k$  by Lemma (7), we have that

$$\begin{aligned} k &= N^-(\hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) \geq N^-(\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = N^-\left(\left(\hat{\mathbf{U}}_2 \boldsymbol{\Gamma} \mathbf{e}_1\right) \left(\hat{\mathbf{U}}_2 \boldsymbol{\Gamma} \mathbf{e}_1\right)^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top\right) \\ &\geq N^-\left(\left(\hat{\mathbf{U}}_2 \boldsymbol{\Gamma}\right) \left(\hat{\mathbf{U}}_2 \boldsymbol{\Gamma}\right)^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top\right) = N^-(\hat{\mathbf{U}}_2 \hat{\mathbf{U}}_2^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k, \end{aligned}$$

Thus,  $N^-(\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k$ . But, since  $\mathbf{y} \in \text{im}(\hat{\mathbf{U}}_1)$ ,  $\text{rank}(\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top) = k$ . Thus,

$$\mathbf{y} \mathbf{y}^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top \preceq 0. \quad (37)$$

Take  $\boldsymbol{\alpha}$  such that  $\hat{\mathbf{U}}_1 \boldsymbol{\alpha}_1 = -\mathbf{y}$  and  $\boldsymbol{\alpha}_2 = \boldsymbol{\Gamma} \mathbf{e}_1$ . Then

$$\mathbf{y}_1 \mathbf{y}_1^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top = \left(\hat{\mathbf{U}}_1 \boldsymbol{\alpha}\right) \left(\hat{\mathbf{U}}_1 \boldsymbol{\alpha}\right)^\top - \hat{\mathbf{U}}_1 \hat{\mathbf{U}}_1^\top \quad (38)$$

$$= \hat{\mathbf{U}}_1 \left(\boldsymbol{\alpha}_1 \boldsymbol{\alpha}_1^\top - \mathbf{I}\right) \hat{\mathbf{U}}_1^\top \preceq 0. \quad (39)$$

1080 Therefore  $\|\alpha_1\|_2 \leq 1$ .

1081  
1082 Then  $g(\alpha; \mathbf{z}) = \left\| \hat{U}_1 \alpha_1 + \hat{U}_2 \alpha_2 \right\|_2^2 + \lambda \left( \|\alpha_1\|_2^2 - \|\alpha_2\|_2^2 \right) = \lambda \left( \|\alpha_1\|_2^2 - 1 \right) \leq 0$ .

1083  
1084 If  $g(\alpha; \mathbf{z}) < 0$ , then we are done. Otherwise,  $g(\alpha; \mathbf{z}) = 0$ . Then the same analysis from Case 1  
1085 will show that  $\nabla g(\alpha; \mathbf{z}) \neq \mathbf{0}$ , so there exists  $\bar{\alpha}$  in an infinitesimal neighborhood around  $\alpha$  where  
1086  $g(\bar{\alpha}; \mathbf{z})$  is strictly negative. This then implies our desired result.

## 1087 1088 B.5 DERIVATION OF EQUATION (6)

1089 Recall our generative model  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma_x^2 \mathbf{I})$ ,  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I})$ , and  $\mathbf{y} = \mathbf{A}_t^* \mathbf{x} + \epsilon$ , where  $\mathbf{x}$  and  $\epsilon$   
1090 are independent. Then,

$$\begin{aligned}
1091 & 2\mathbb{E}[\mathcal{L}_t^1(\mathbf{A}_t)] = \mathbb{E} \left[ \|\mathbf{y} - \mathbf{A}_t \mathbf{x}\|_2^2 \right] \\
1092 & = \mathbb{E} \left[ \|\mathbf{A}_t^* \mathbf{x} + \epsilon - \mathbf{A}_t \mathbf{x}\|_2^2 \right] \\
1093 & = \mathbb{E} \left[ \|(\mathbf{A}_t^* - \mathbf{A}_t) \mathbf{x} + \epsilon\|_2^2 \right] \\
1094 & = \mathbb{E} \left[ \left( \|(\mathbf{A}_t^* - \mathbf{A}_t) \mathbf{x}\|_2^2 + \|\epsilon\|_2^2 + 2\epsilon^\top (\mathbf{A}_t^* - \mathbf{A}_t) \mathbf{x} \right) \right] \\
1095 & = \mathbb{E} \left[ \mathbf{x}^\top (\mathbf{A}_t^* - \mathbf{A}_t)^\top (\mathbf{A}_t^* - \mathbf{A}_t) \mathbf{x} \right] + \mathbb{E} \left[ \|\epsilon\|_2^2 \right] + 2\mathbb{E} \left[ \epsilon^\top (\mathbf{A}_t^* - \mathbf{A}_t) \mathbf{x} \right] \\
1096 & = \mathbb{E} \left[ \text{tr} \left( \mathbf{x}^\top (\mathbf{A}_t^* - \mathbf{A}_t)^\top (\mathbf{A}_t^* - \mathbf{A}_t) \mathbf{x} \right) \right] + \sigma_\epsilon^2 + 2\mathbb{E} \left[ \epsilon^\top (\mathbf{A}_t^* - \mathbf{A}_t) \mathbb{E}[\mathbf{x}] \right] \quad (\epsilon, \mathbf{x} \text{ are independent}) \\
1097 & = \mathbb{E} \left[ \text{tr} \left( (\mathbf{A}_t^* - \mathbf{A}_t)^\top (\mathbf{A}_t^* - \mathbf{A}_t) \mathbf{x} \mathbf{x}^\top \right) \right] + \sigma_\epsilon^2 \quad (\text{by cyclic property of trace and since } \mathbb{E}[\mathbf{x}] = 0) \\
1098 & = \text{tr} \left( (\mathbf{A}_t^* - \mathbf{A}_t)^\top (\mathbf{A}_t^* - \mathbf{A}_t) \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \right) + \sigma_\epsilon^2 \\
1099 & = \sigma_x^2 \text{tr} \left( (\mathbf{A}_t^* - \mathbf{A}_t)^\top (\mathbf{A}_t^* - \mathbf{A}_t) \right) + \sigma_\epsilon^2 \\
1100 & = \sigma_x^2 \|\mathbf{A}_t^* - \mathbf{A}_t\|_F^2 + \sigma_\epsilon^2
\end{aligned}$$

1101 Thus,  $\mathbb{E}[\mathcal{L}_t^1(\mathbf{A}_t)] = \frac{1}{2} \left( \sigma_x^2 \|\mathbf{A}_t^* - \mathbf{A}_t\|_F^2 + \sigma_\epsilon^2 \right)$ . Then  $\mathbb{E}[\mathcal{L}_t^N(\mathbf{A}_t)] = \mathbb{E}[\mathcal{L}_t^1(\mathbf{A}_t)]$  by linearity of  
1102 expectation, so

$$\begin{aligned}
1103 & \frac{1}{2} \left\| \mathbf{A}_t^* + \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \mathbf{A}_t \right\|_F^2 = \frac{1}{\sigma_x^2} \left( \mathbb{E}[\mathcal{L}_t^N(\mathbf{A}_t)] - \frac{\sigma_\epsilon^2}{2} \right) \\
1104 & \\
1105 & \\
1106 & \\
1107 & \\
1108 & \\
1109 & \\
1110 & \\
1111 & \\
1112 & \\
1113 & \\
1114 & \\
1115 & \\
1116 & \\
1117 & \\
1118 & \\
1119 & \\
1120 & \\
1121 & \\
1122 & \\
1123 & \\
1124 & \\
1125 & \\
1126 & \\
1127 & \\
1128 & \\
1129 & \\
1130 & \\
1131 & \\
1132 & \\
1133 &
\end{aligned}$$

## C LLM TRAINING HYPERPARAMETERS

Hyperparameter	Standard Retraining	Meta-LoRA-8	Meta-LoRA-16
Learning Rate	5e-5	3e-5	5e-5
Learning Rate Schedule	Linear	Linear	Linear
Batch Size	6	4	4
Epochs	30	30	30
Optimizer	AdamW	AdamW	AdamW
LoRA Rank	N/A	8	16
LoRA Dropout	N/A	0.1	.1
LoRA Alpha	N/A	16	16

Table 2: Retraining Hyperparameters

Hyperparameter	Rank- $k$ LoRA Fine-Tuning
Learning Rate	3e-5
Learning Rate Schedule	Linear
Batch Size	6
Epochs	30
Optimizer	AdamW
LoRA Rank	$k$
LoRA Dropout	.1
LoRA Alpha	16

Table 3: Rank- $k$  LoRA Fine-Tuning Hyperparameters,  $k \in \{8, 16\}$ 

### C.1 NOTE ON NUMBER OF TRAINABLE PARAMETERS

For simplicity assume our model architecture consisted of  $m$  layers, where each layer was parameterized by a  $d \times d$  matrix, and we use rank- $k$  adaptations for each layer for our Meta-LoRA objective, where  $k \ll d$ . Then the standard retraining method uses  $md^2$  trainable parameters, while minimizing the Meta-LoRA objective uses  $m(d^2 + 2kdT)$  trainable parameters. Although Meta-LoRA uses some additional parameters, since  $k$  is small relative to  $d$  and we work in the setting where  $k(T + 1) < d$ , asymptotically  $m(d^2 + 2kdT) = O(md^2)$  so the increase in trainable parameters is minor. After running either of these retraining procedures, the fine-tuning stages are identical and require the same number of trainable parameters no matter which retraining procedure was run.

## D THEORY NOTES

### D.1 NON-UNIQUENESS OF GLOBAL MIN FOR $T = 2$

Consider  $T = 2, k = 1, d = 2, \mathbf{A}^* = \mathbf{0}$ , and  $\mathbf{u}_t^* = \mathbf{e}_t$  for  $t = 1, 2$ , where  $\mathbf{e}_t$  is the  $t_{th}$  standard basis vector. Clearly the ground truth perturbations  $\mathbf{u}_i^*$  are orthonormal and thus linearly independent. The set of global minima of  $\mathcal{L}$  are  $(\mathbf{A}, \mathbf{U})$  such that  $\mathbf{A} = \frac{1}{T} \sum_{t=1}^T (\mathbf{u}_t^* \mathbf{u}_t^{*\top} - \mathbf{u}_t \mathbf{u}_t^\top)$  and  $\mathbf{u}_t \mathbf{u}_t^\top - \mathbf{u}_t^* \mathbf{u}_t^{*\top} - \frac{1}{T} \sum_{s=1}^T (\mathbf{u}_s \mathbf{u}_s^\top - \mathbf{u}_s^* \mathbf{u}_s^{*\top}) = \mathbf{0}$ . It is not hard to see that a global minimum follows from any set values of  $\mathbf{u}_1, \mathbf{u}_2$  such that  $\mathbf{u}_1 \mathbf{u}_1^\top - \mathbf{u}_2 \mathbf{u}_2^\top = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$ . When properly parameterized, this system of equations defines a hyperbola where each point corresponds to a global minimum of  $\mathcal{L}$ .

### D.2 SPURIOUS LOCAL MINIMA

We observe that for  $T \geq 3$ , for certain tasks  $\mathbf{U}^* = (\mathbf{U}_1^*, \mathbf{U}_2^*, \mathbf{U}_3^*)$ , it is possible to find points  $\mathbf{U}$  that are local minima, but not global minima. To find these points, we sample true tasks  $\mathbf{U}^*$  from a

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

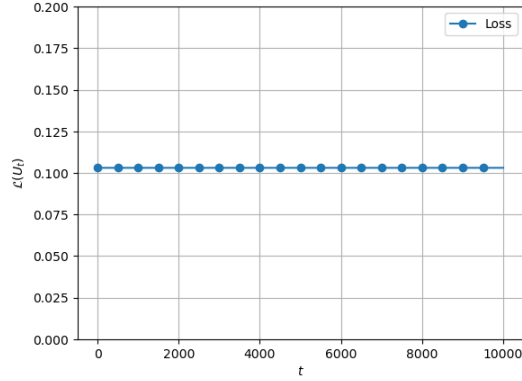


Figure 2: Loss does not decrease near these spurious local minima

normal distribution and use a numerical solver to find zeros of the gradient of the reduced loss

$$\hat{\mathcal{L}}(\mathbf{U}) = \sum_{t=1}^T \left\| \mathbf{U}_t \mathbf{U}_t^\top - \mathbf{U}_t^* \mathbf{U}_t^{*\top} - \frac{1}{T} \sum_{s=1}^T (\mathbf{U}_s \mathbf{U}_s^\top - \mathbf{U}_s^* \mathbf{U}_s^{*\top}) \right\|_F^2.$$

Through the Schur complement argument used to prove Theorem 4, we can see that  $\hat{\mathcal{L}}$  has a spurious local minimum only if  $\mathcal{L}$  has a spurious local minimum.

Typically, these zeros are close to the global minimum. Occasionally, it is possible to find a point  $\hat{\mathbf{U}}$  with gradients close to 0 and with positive definite Hessians. We then confirm that these are close to the spurious local minimum through the following argument.

Consider the function

$$r(\mathbf{U}) = \text{vec}(\mathbf{U} - \hat{\mathbf{U}})^\top \text{vec}(\nabla \hat{\mathcal{L}}(\mathbf{U})).$$

Clearly, there is a minimum of  $\hat{\mathcal{L}}$  in the  $\delta$ -ball of  $\hat{\mathbf{U}}$  if  $r(\mathbf{U}) > 0$  for all  $\mathbf{U}$  on the boundary of the  $\delta$ -ball. As  $r$  is continuous, if for some small enough  $\epsilon, \gamma > 0$  if  $r(\mathbf{U}) > \gamma > 0$  for all  $\mathbf{U}$  on the  $\epsilon$ -net of the boundary of the  $\delta$ -ball, then there exists a spurious local minimum in the  $\delta$ -ball around  $\hat{\mathbf{U}}$ . Numerically, such points and  $\epsilon, \delta$ , and  $\gamma$  can be found which would imply that spurious local minima exist, barring any errors due to numerical computation. To confirm, we run gradient descent from this point and observe that the loss stays constant.

## E EXAMPLE PSEUDOCODE FOR MINIMIZING (4)

---

### Algorithm 1 Meta-Adapter Training

---

- 1: **Input:** Tasks  $\mathcal{T}_t, t \in [T]$ , learning rate  $\eta$ , number of epochs  $N_e$ , batches per epoch  $N_b$
  - 2: **Initialize:** Model parameters  $\mathbf{W}_0, \boldsymbol{\theta}_0^{(t)}$  for all  $t = 1, \dots, T$
  - 3: **for** epoch  $e = 1$  to  $N_e$  **do**
  - 4:   **for**  $b = 1, \dots, N_b$  **do**
  - 5:     **for**  $t = 1, \dots, T$  **do**
  - 6:       Load next batch  $\beta_{t,b}$  from  $\mathcal{T}_t$
  - 7:       Compute gradient  $\mathbf{g}^{(t)} = \nabla_{\mathbf{W}, \boldsymbol{\theta}^{(t)}} \left( \sum_{(\mathbf{x}, \mathbf{y}) \in \beta_{t,b}} \mathcal{L}(\Phi_{\text{FT}}(\mathbf{x}; \mathbf{W}, \boldsymbol{\theta}^{(t)}), \mathbf{y}) \right)$
  - 8:       Update adapter parameters:  $\boldsymbol{\theta}_{e+1}^{(t)} \leftarrow \boldsymbol{\theta}_e^{(t)} - \eta_e \mathbf{g}^{(t)}$
  - 9:     **end for**
  - 10:    Update base parameters:  $\mathbf{W}_{e+1} \leftarrow \mathbf{W}_e - \eta_e \sum_{t=1}^T \mathbf{g}^{(t)}$
  - 11:   **end for**
  - 12: **end for**
-