

TRANSFORMERS LEARN TO IMPLEMENT MULTI-STEP GRADIENT DESCENT WITH CHAIN OF THOUGHT

Anonymous authors

Paper under double-blind review

ABSTRACT

Chain of Thought (CoT) prompting has been shown to significantly improve the performance of large language models (LLMs), particularly in arithmetic and reasoning tasks, by instructing the model to produce intermediate reasoning steps. Despite the remarkable empirical success of CoT and its theoretical advantages in enhancing expressivity, the mechanisms underlying CoT training remain largely unexplored. In this paper, we study the training dynamics of transformers over a CoT objective on an in-context weight prediction task for linear regression. We prove that while a one-layer linear transformer without CoT can only implement a single step of gradient descent (GD) and fails to recover the ground-truth weight vector, a transformer with CoT prompting can learn to perform multi-step GD autoregressively, achieving near-exact recovery. Furthermore, we show that the trained transformer effectively generalizes on the unseen data. Empirically, we demonstrate that CoT prompting yields substantial performance improvements.

1 INTRODUCTION

Transformer-based Large Language Models (LLMs) have demonstrated significant success across various language modeling tasks, achieving state-of-the-art performance in numerous domains (OpenAI, 2023). Remarkably, these models have also unlocked complex reasoning abilities, particularly in mathematical problem-solving and coding tasks (Chowdhery et al., 2023; Anil et al., 2022; Achiam et al., 2023). A key method driving this advancement is the Chain of Thought (CoT), which enables LLMs to generate intermediate reasoning steps autoregressively rather than providing a direct answer. This process effectively improves the model’s capacity to solve complex problems. In practice, CoT reasoning can be elicited either by providing few-shot CoT examples or by appending prompts like “let’s think step by step” to bootstrap the model’s response (Kojima et al., 2022; Wei et al., 2022; Suzgun et al., 2022; Nye et al., 2021).

Theoretically, CoT enables LLMs to perform multi-step sequential computations by generating intermediate results, thereby significantly improving the expressive power of transformers (Li et al., 2024b; Feng et al., 2024; Merrill & Sabharwal, 2023a) compared to standard decoder transformers that generate direct outputs without intermediate reasoning (Liu et al., 2022; Merrill & Sabharwal, 2023b). Despite these theoretical insights, it remains unclear how transformers are **trained** on CoT data to effectively execute multi-step reasoning. Furthermore, it is unknown whether a transformer trained specifically with an auto-regressive objective with multi-step CoT can substantially outperform one trained to directly output answers without CoT.

This paper takes an initial step beyond expressiveness to study the training dynamics of transformers when trained on CoT data. Specifically, following the modified in-context learning (ICL) setting on linear regression proposed by Ahn et al. (2023); Zhang et al. (2023), we use it as a testbed to analyze the training process with the CoT framework implemented. We name the task **in-context weight prediction** where the goal is to predict the linear weight vector from the sequence of input prompts. Instead of performing direct ICL and outputting a prediction, the transformer with CoT prompting is allowed to generate multiple intermediate steps before arriving at the final answer. We theoretically investigate the transformer’s training trajectory on the CoT objective and show the expressiveness gap between transformers trained with CoT and those without. Our main results show this separation is **learnable**: gradient-based algorithm can learn the constructed transformer with CoT.

We summarize our contributions as follows:

- **Expressiveness Gap.** We characterize the global optimum of the population loss for the in-context weight prediction task on linear regression using a one-layer transformer without CoT prompting. Our results show that, without CoT, the transformer at the global minimizer effectively performs a single step of gradient descent (GD) (Theorem 3.1), leading to significant errors in predicting the d -dimensional weight vector $w^* \in \mathbb{R}^d$ when the number of examples for ICL is $n = \tilde{\Theta}(d)$ (Corollary 3.1). In contrast, we demonstrate that a one-layer transformer with CoT prompting can achieve near-exact recovery by executing multi-step GD (Theorem 3.2).
- **Global convergence.** We prove the convergence results of running gradient flow on the population CoT loss under mild assumptions (Theorem 4.1). Our analysis uses a novel stage-wise approach combining dynamics analysis and landscape properties: the parameters initially approach the global minimizer, followed by local convergence toward the final solution. Our proof technique involves a novel characterization of the complicated population gradient. Furthermore, we prove that the trained transformer can exhibit both in-distribution and out-of-distribution generalization (Theorem 4.2) at inference time. We are the first to establish the learnable separation between transformers with and without CoT under the in-context linear regression setting. We empirically validate that the trained transformer converges to the minimizer predicted by our theory, with a distinct performance gap between models trained with and without CoT prompting.

Outline. In Section 2, we formalize the problem setting including the data model, the one-layer transformer architecture, and the CoT prompting format. In Section 3, we theoretically show the performance gap between the transformer with and without CoT. Section 4 consists of our main results, including our dynamics analysis and out-of-distribution (OOD) generalization result. Section 5 empirically validates the advantage of CoT.

1.1 RELATED WORKS

Training dynamics of transformers. Several works have studied the training process of specific transformer architectures. Jelassi et al. (2022); Li et al. (2023) examined the training process and sample complexity of Vision Transformer (Dosovitskiy et al., 2020). Tarzanagh et al. (2023); Ataee Tarzanagh et al. (2023); Li et al. (2024a) explored the connection between the optimization landscape of self-attention mechanisms and the Support Vector Machine problem. Tian et al. (2023a;c) provided insights into the training dynamics of the self-attention and MLP layers during the training process respectively.

A related line of research focuses on Markov-like data models. Bietti et al. (2024) studied the *induction head* mechanism from the perspective of associative memory. Nichani et al. (2024) demonstrated that a simplified two-layer transformer provably learns a generalized induction head on latent causal graphs. Chen et al. (2024b) further proved that a modified two-layer multi-head transformer can learn in-context generalized n -gram. Edelman et al. (2024) investigated the multi-stage phase transitions during training on bigram and n -gram ($n \geq 3$). Additionally, Makkuva et al. (2024) studied the loss landscape of transformers trained on sequences from a Markov Chain.

Another growing body of literature aims to understand the training dynamics of in-context learning (ICL). Garg et al. (2022) first empirically studied the ICL capabilities of transformers over a variety of function classes. Akyürek et al. (2022); Von Oswald et al. (2023) investigated the behavior of transformers on random ICL instances of linear regression. Several works have also established the existence of deep transformers capable of implementing multi-step gradient descent (GD) across different domains (Fu et al., 2023; Bai et al., 2023; Giannou et al., 2023). Mahankali et al. (2023); Ahn et al. (2024) analyzed the loss landscape of the linear regression ICL task and Zhang et al. (2023) proved global convergence on a one-layer linear self-attention layer using gradient flow. Gatmiry et al. (2024) demonstrated that a linear looped transformer with specific update procedures can learn to implement multi-step GD for linear regression. Further analyses of training dynamics under more realistic assumptions about data models and architectures have been conducted by Huang et al. (2023); Kim & Suzuki (2024); Chen et al. (2024a). For a detailed discussion see Appendix A.1.

Compared to prior works, our study and Huang et al. (2023); Ahn et al. (2024); Zhang et al. (2023); Tarzanagh et al. (2023); Nichani et al. (2024); Kim & Suzuki (2024); Wang et al. (2024); Chen et al. (2024b) all use similar reparameterizations that combine key and query matrices to simplify the training dynamics. Moreover, many previous studies (Tian et al., 2023a; Zhang et al., 2023; Huang et al., 2023; Nichani et al., 2024; Kim & Suzuki, 2024; Chen et al., 2024a; Gatmiry et al., 2024) adopted the population loss to facilitate the analysis of these dynamics.

A closely related work is Gatmiry et al. (2024), which shows that a looped transformer can implement multi-step GD on the ICL linear regression task to directly predict the query answer in context. In comparison, the goal of our setting is to predict the weight vector from the input examples using a realistic CoT autoregressive generation process. Theoretically, we also establish a performance gap between transformers with CoT and those without. See Appendix A.2 for a more detailed discussion.

Chain of Thought and Scratchpad The CoT prompting method was first introduced by Wei et al. (2022) to enhance the multi-step reasoning capability of LLMs. Before the formalization of CoT, Nye et al. (2021) demonstrated that allowing language models to generate intermediate results on “*scratchpads*” dramatically boosts the multi-step computation ability of LLMs. Wang et al. (2022b); Yao et al. (2024); Creswell et al. (2022); Zhou et al. (2022) further proposed variants of the CoT/scratchpad method to improve the efficiency and reliability of generation.

Recently, several works have attempted to understand CoT from both experimental and theoretical perspectives. Wang et al. (2022a); Saparov & He (2022); Shi et al. (2022); Paul et al. (2023) empirically studied the capability of CoT, providing valuable insights on its reasoning processes. Meanwhile, Wu et al. (2023); Tutunov et al. (2023); Hou et al. (2023); Cabannes et al. (2024) investigated CoT through the lens of mechanistic interpretability. On the theoretical side, Liu et al. (2022); Merrill & Sabharwal (2023a); Li et al. (2024b); Feng et al. (2024) explored the expressive power of transformers with CoT, showing that CoT can significantly extend the expressivity of transformers in the context of circuit complexity. Hu et al. (2024) investigated the statistical foundations of CoT. However, the training dynamics of CoT remain largely unexplored. To the best of our knowledge, this work is among the first theoretical analyses of training dynamics on CoT/scratchpad objectives.

2 PRELIMINARIES

In this section, we describe the modified in-context learning linear regression task, i.e. **in-context weight prediction**, the one-layer linear self-attention architecture, and the Chain of Thought (CoT) prompting formulation.

Notation We use $[T]$ to denote the set $\{1, 2, \dots, T\}$. Scalars are in lower-case unbolded letters (y, α , etc.). Matrices and vectors are denoted in upper-case bold letters (\mathbf{W}, \mathbf{V} , etc.) and lower-case bold letters (\mathbf{x}, \mathbf{w} , etc.), respectively. $\mathbf{W}_{[i,j]}$, $\mathbf{W}_{[i,:]}$, $\mathbf{W}_{[:,j]}$ respectively denotes the (i, j) -th entry, i -th row, and j -th column of the matrix \mathbf{W} . $\mathbf{W}_{[:, -1]}$ means the last column of the matrix \mathbf{W} . The notation \mathbf{W}_{ij} denotes block matrices/vectors on the i -th row and j -th column according to context. For norm, $\|\cdot\|$ denotes ℓ_2 norm and $\|\cdot\|_F$ denotes the Frobenius norm. We use $\mathbb{1}\{\cdot\}$ to denote the indicator function. We use $\tilde{O}(\cdot)$ to hide logarithmic factors in the asymptotic notations.

2.1 IN-CONTEXT WEIGHT PREDICTION

Previous works (Zhang et al., 2023; Ahn et al., 2023; 2024; Akyürek et al., 2022; Mahankali et al., 2023) focus on the in-context learning (ICL) task on linear regression. We suppose the data sequence is sampled from a linear regression task where the ground-truth

$$\mathbf{w}^* \sim \mathcal{N}(0, \mathbf{I}_d) \quad \mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d) \quad y_i = \mathbf{w}^{*\top} \mathbf{x}_i \text{ for all } i \in [n]. \quad (1)$$

The goal of in-context learning is to predict the correct label $\mathbf{w}^{*\top} \mathbf{x}_{\text{query}}$ given a query $\mathbf{x}_{\text{query}}$ and the previous example pairs (\mathbf{x}_i, y_i) . Most previous works (Zhang et al., 2023; Ahn et al., 2024; Mahankali et al., 2023) show the transformer predicts the query label y_{query} by implicitly doing a one-step gradient descent without predicting the linear classifier \mathbf{w}^* .

In this work, we go one step further: instead of directly outputting the query label, we require the transformers to implement gradient descent to learn the ground-truth weight vector \mathbf{w}^* . We call this task **in-context weight prediction** for linear regression. Specifically, the data sequence is in the following format:

$$\mathbf{Z}_0 = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n & 0 \\ y_1 & \cdots & y_n & 0 \\ 0 & \cdots & 0 & \mathbf{w}_0 \\ 0 & \cdots & 0 & 1 \end{bmatrix} := \begin{bmatrix} \mathbf{X} & 0 \\ \mathbf{y} & 0 \\ \mathbf{0}_{d \times n} & \mathbf{w}_0 \\ \mathbf{0}_{1 \times n} & 1 \end{bmatrix} \in \mathbb{R}^{d_e \times (n+1)}, \quad (2)$$

where $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n]$ is the data matrix and \mathbf{w}_0 is the initialization of the linear parameter $\hat{\mathbf{w}}$. We assume $\mathbf{w}_0 = \mathbf{0}_d$ for simplicity, and define $d_e = 2d + 2$. Our setting is similar to the setting in Bai et al. (2023) where multi-layer transformers are constructed to do explicit multi-step GD on the

weight vector \hat{w} . We separate the input example space and the weight vector space as in Bai et al. (2023) (the $\{\mathbf{p}_i\}_{i \in [N+1]}$) in order to facilitate training. Moreover, we add a dummy token (an extra 1) at the end of each token similar to what Bai et al. (2023) did in their input sequence format.

2.2 LINEAR SELF-ATTENTION LAYER

We consider a one-layer linear self-attention (LSA) module with residual connection, following the setting in Zhang et al. (2023); Ahn et al. (2023); Gatmiry et al. (2024): we remove the $\text{softmax}(\cdot)$ non-linearity, consolidate the projection and value matrix into a single matrix $\mathbf{V} \in \mathbb{R}^{d_e \times d_e}$, and merge the key and query matrices into $\mathbf{W} \in \mathbb{R}^{d_e \times d_e}$. We denote

$$f_{\text{LSA}}(\mathbf{Z}; \mathbf{V}, \mathbf{W}) = \mathbf{Z} + \mathbf{V} \mathbf{Z} \cdot \frac{\mathbf{Z}^\top \mathbf{W} \mathbf{Z}}{n} \quad (3)$$

The prediction of the transformer will be the last token of the output sequence, namely

$$f_{\text{LSA}}(\mathbf{Z}; \mathbf{V}, \mathbf{W})_{[:, -1]} = \mathbf{Z}_{[:, -1]} + \mathbf{V} \mathbf{Z} \cdot \frac{\mathbf{Z}^\top \mathbf{W} \mathbf{Z}_{[:, -1]}}{n} \quad (4)$$

Since the first $(d+1)$ entries of the full weight tokens $(0, 0, \mathbf{w}, 1)$ are zero, only part of the \mathbf{W} and \mathbf{V} affect the prediction. We can rewrite the parameter \mathbf{V}, \mathbf{W} into block matrices

$$\mathbf{V} = \begin{bmatrix} \mathbf{V}_{11} & \mathbf{V}_{12} & \mathbf{V}_{13} & \mathbf{V}_{14} \\ \mathbf{V}_{21} & v_{22} & \mathbf{V}_{23} & v_{24} \\ \mathbf{V}_{31} & \mathbf{V}_{32} & \mathbf{V}_{33} & \mathbf{V}_{34} \\ \mathbf{V}_{41} & v_{42} & \mathbf{V}_{43} & v_{44} \end{bmatrix}, \mathbf{W} = \begin{bmatrix} \mathbf{W}_{11} & \mathbf{W}_{12} & \mathbf{W}_{13} & \mathbf{W}_{14} \\ \mathbf{W}_{21} & w_{22} & \mathbf{W}_{23} & w_{24} \\ \mathbf{W}_{31} & \mathbf{W}_{32} & \mathbf{W}_{33} & \mathbf{W}_{34} \\ \mathbf{W}_{41} & w_{42} & \mathbf{W}_{43} & w_{44} \end{bmatrix} \in \mathbb{R}^{(2d+2) \times (2d+2)}$$

where the block matrices are in the following shape $(i, j \in \{1, 2\})$:

$$\mathbf{V}_{2i-1, 2j-1}, \mathbf{W}_{2i-1, 2j-1} \in \mathbb{R}^{d \times d}; \mathbf{V}_{2i-1, 2j}, \mathbf{W}_{2i-1, 2j}, \mathbf{V}_{2i, 2j-1}^\top, \mathbf{W}_{2i, 2j-1}^\top \in \mathbb{R}^{d \times 1}; v_{2i, 2j}, w_{2i, 2j} \in \mathbb{R}.$$

In the following sections, we will show only \mathbf{V}_{31} , \mathbf{W}_{13} , and w_{24} affects the prediction. We will further prove that all other entries are always zero along the training trajectory if initialized at zero.

2.3 CHAIN-OF-THOUGHT PROMPTING

In language modeling tasks, transformers have been proven to be versatile in various downstream tasks. However, transformers struggle to solve mathematical or scientific problems with one single generation, where several reasoning steps are required. CoT was then proposed to make transformers learn to generate intermediate results auto-regressively before reaching the answer.

With CoT, we allow the transformer to generate k steps before it outputs the final prediction \hat{w}_k for the ground-truth \mathbf{w}^* . Specifically, given the generated input sequence $\hat{\mathbf{Z}}_i$ at the i -th step of generation, we have $f_{\text{LSA}}(\hat{\mathbf{Z}}_i)_{[:, -1]}$ as the prediction of the next token $((i+1)$ -th token), and append it to the end of the current sequence s.t. $\hat{\mathbf{Z}}_{i+1} = [\hat{\mathbf{Z}}_i, f_{\text{LSA}}(\hat{\mathbf{Z}}_i)_{[:, -1]}]$. After k generation steps, the CoT process induces k intermediate sequences $\{\hat{\mathbf{Z}}_i\}_{i=1}^k$ in the following form:

$$\hat{\mathbf{Z}}_i = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n & 0 & \star & \cdots & \star \\ y_1 & \cdots & y_n & 0 & \star & \cdots & \star \\ 0 & \cdots & 0 & \mathbf{w}_0 & \hat{\mathbf{w}}_1 & \cdots & \hat{\mathbf{w}}_i \\ 0 & \cdots & 0 & 1 & 1 & \cdots & 1 \end{bmatrix} \in \mathbb{R}^{d_e \times (n+i+1)}, i \in [k] \quad (\text{Inference})$$

Here, we define $\hat{\mathbf{w}}_i := f_{\text{LSA}}(\hat{\mathbf{Z}}_{i-1})_{[d+2:2d+1, -1]}$ as the i -th step prediction for the weight vector. The other entries in the same column are irrelevant and we denote them as \star . Finally, the transformer inputs the last generated sequence $\hat{\mathbf{Z}}_k$ back to the transformer once again to generate the final output $\hat{\mathbf{w}}_{k+1} := f_{\text{LSA}}(\hat{\mathbf{Z}}_k)_{[d+2:2d+1, -1]}$ as the prediction of the weight vector \mathbf{w}^* .

Different from the inference time generation, the training process is similar to pre-training on the ground-truth sequence to predict the next token. Specifically, we input the transformer with CoT ground-truth sequences \mathbf{Z}_i :

$$\mathbf{Z}_i = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n & 0 & 0 & \cdots & 0 \\ y_1 & \cdots & y_n & 0 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & \mathbf{w}_0 & \mathbf{w}_1 & \cdots & \mathbf{w}_i \\ 0 & \cdots & 0 & 1 & 1 & \cdots & 1 \end{bmatrix} \in \mathbb{R}^{d_e \times (n+i+1)}, i \in [k] \quad (\text{Training})$$

where $\mathbf{w}_i = \mathbf{w}_{i-1} - \eta \cdot \frac{\mathbf{X}(\mathbf{X}^\top \mathbf{w}_{i-1} - \mathbf{y}^\top)}{n}$ is the ground-truth intermediate weight vector after i gradient steps on the linear regression objective. Each gradient step adopts a fixed learning rate η for all possible training instances $\{\mathbf{X}, \mathbf{w}\}$ when generating the ground-truth sequence \mathbf{Z}_i . Note that \mathbf{Z}_i is the corresponding ground-truth sequence of $\hat{\mathbf{Z}}_i$.

In the training objective for the i -th step, the transformer is required to predict the next token $\mathbf{Z}_{i+1}[:, -1] := (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1)$ given the i -th ground-truth intermediate sequence \mathbf{Z}_i . Finally, we predict the final ground-truth weight vector \mathbf{w}^* with the final intermediate sequence \mathbf{Z}_k . The CoT **training** objective given a sample prompt \mathbf{X}, \mathbf{y} then becomes:

$$\ell^{\text{CoT}}(\mathbf{X}, \mathbf{w}^*; \mathbf{V}, \mathbf{W}) = \frac{1}{2} \sum_{i=0}^k \left\| f_{\text{LSA}}(\mathbf{Z}_i)[:, -1] - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right\|^2 \quad (5)$$

Here we denote $\mathbf{w}_{k+1} := \mathbf{w}^*$ for clarity. Following Zhang et al. (2023); Nichani et al. (2024); Kim & Suzuki (2024); Tian et al. (2023b); Chen et al. (2024a); Gatmiry et al. (2024), we consider the gradient flow dynamics over the population loss of the CoT objective:

$$\mathcal{L}^{\text{CoT}}(\mathbf{V}, \mathbf{W}) = \mathbb{E}_{\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_d), \mathbf{w}^* \sim \mathcal{N}(0, \mathbf{I}_d)} [\ell^{\text{CoT}}(\mathbf{X}, \mathbf{w}^*; \mathbf{V}, \mathbf{W})] \quad (6)$$

For clarity, we write the expectation as $\mathbb{E}_{\mathbf{X}, \mathbf{w}^*}[\cdot]$. The following differential equation gives the gradient flow dynamics of the parameters:

$$\frac{d\boldsymbol{\theta}}{dt} = -\nabla \mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}), \quad \boldsymbol{\theta} := (\mathbf{V}, \mathbf{W}).$$

When measuring the performance after training, we apply the CoT **inference** procedure to generate k intermediate sequences $\{\hat{\mathbf{Z}}_i\}_{i=1}^k$ and consider the final output token $f(\hat{\mathbf{Z}}_k)[:, -1]$ by inputting the last generated sequence $\hat{\mathbf{Z}}_k$. The performance evaluation is measured on the error between the final output $f(\hat{\mathbf{Z}}_k)[:, -1]$ and the ground-truth \mathbf{w}^* :

$$\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\left\| f_{\text{LSA}}(\hat{\mathbf{Z}}_k)[:, -1] - (\mathbf{0}_d, 0, \mathbf{w}^*, 1) \right\|^2 \right] \quad (7)$$

When CoT prompting is not used ($k = 0$), the evaluation loss $\mathcal{L}^{\text{Eval}}$ is equivalent to \mathcal{L}^{CoT} .

3 EXPRESSIVENESS IMPROVEMENT WITH CHAIN OF THOUGHT

In this section, we theoretically explore the performance gap on our data model between transformers with CoT and those without. We first prove that a one-layer transformer without CoT can only implement a one-step GD and cannot recover the ground-truth, while it can near-exactly predict the ground-truth parameter with CoT by implementing multi-step GD.

3.1 ONE-LAYER TRANSFORMER CANNOT RECOVER GROUND-TRUTH

For the ICL linear regression task, the optimal prediction given by a one-layer linear transformer is equivalent to a single step of GD on the MSE objective of linear regression (Mahankali et al., 2023). What about our task on predicting the ground-truth weight vector \mathbf{w}^* in context? The following theorem proves that the optimal solution is still a one-step GD solution.

Theorem 3.1 (Lower bound without CoT). *If the global minimizer of $\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W})$ is $(\mathbf{V}^*, \mathbf{W}^*)$, the corresponding one-layer transformer $f_{\text{LSA}}(\mathbf{Z}_0)[:, -1]$ implements one step GD on a linear model with some learning rate $\eta^* = \frac{n}{n+d+1}$ and the transformer outputs $(\mathbf{0}_d, 0, \frac{\eta^*}{n} \mathbf{X} \mathbf{y}^\top, 1)$.*

We briefly present the high-level intuitions in the proof and the detailed proof is deferred to Appendix B.1. We use a similar technique in Mahankali et al. (2023) when proving the optimality of one-step GD in the ICL task. The key strategy of the proof is to replace $(\mathbf{0}_d, 0, \mathbf{w}^*, 1)$ in the evaluation loss $\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W})$ (Equation (7)) with $(\mathbf{0}_d, 0, \frac{\eta^*}{n} \mathbf{X} \mathbf{y}^\top, 1)$ in the following form.

$$\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) = \frac{1}{2} \mathbb{E} \left[\left\| f_{\text{LSA}}(\mathbf{Z}_0)[:, -1] - \left(\mathbf{0}_d, 0, \frac{\eta^*}{n} \mathbf{X} \mathbf{y}^\top, 1 \right) \right\|^2 \right] + C$$

In order to prove this equation above, we show the gradient of the original loss Equation (7) and this formula are identical. We first obtain the closed-form formula of the expected gradient for both sides with regard to \mathbf{X}, \mathbf{w}^* . Then we use the symmetric property of the distribution of \mathbf{X}, \mathbf{w}^* to simplify the gradient expressions, and eventually prove them equal.

The equivalent form of loss indicates that the evaluation loss only depends on the ℓ_2 distance between the output of the linear self-attention module and $(\mathbf{0}_d, 0, \frac{\eta^*}{n} \mathbf{X} \mathbf{y}^\top, 1)$. Therefore, any (\mathbf{V}, \mathbf{W}) is a global minimizer of this loss function if and only if the output of $f_{\text{LSA}}(\mathbf{Z}_k)_{[:, -1]}$ is $(\mathbf{0}_d, 0, \frac{\eta^*}{n} \mathbf{X} \mathbf{y}^\top, 1)$. Meanwhile, one can assign

$$\mathbf{V}^* = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ -\eta^* \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \mathbf{W}^* = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & -1 \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (8)$$

and the one-layer transformer achieves the optimal solution, which concludes the proof.

Is a one-step gradient solution good enough? Most of the previous ICL work Zhang et al. (2023); Ahn et al. (2023); Gatmiry et al. (2024) consider the number of examples $n \rightarrow +\infty$ when d is fixed. In this case, the one-step GD solution can perfectly find the ground-truth weight vector \mathbf{w}^* . However, a simple corollary of this theorem indicates that the one-step solution has a non-negligible error when there are limited samples, e.g. $n = \tilde{\Theta}(d)$. This number of examples n is required to guarantee the reconstruction of $\mathbf{w}^* \in \mathbb{R}^d$.

Corollary 3.1. *For any parameters (\mathbf{V}, \mathbf{W}) in the one-layer transformer, $\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) \geq \Theta\left(\frac{d^2}{n}\right)$.*

Moreover, if $n = \tilde{\Theta}(d)$, $\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) = \tilde{\Theta}(d) \xrightarrow{d \rightarrow +\infty} +\infty$.

Proof. By Theorem 3.1, we directly calculate the evaluation loss on the global optimum:

$$\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) \geq \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| \frac{\eta^*}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* - \mathbf{w}^* \right\|^2 = \frac{1}{2} \mathbb{E}_{\mathbf{X}} \text{tr} \left(\mathbf{I} - \frac{\eta^*}{n} \mathbf{X} \mathbf{X}^\top \right)^2$$

since $\mathbb{E}_{\mathbf{w}^*} [\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}$. Apply $\mathbb{E}[\mathbf{X} \mathbf{X}^\top] = n\mathbf{I}$ and $\mathbb{E}[(\mathbf{X} \mathbf{X}^\top)^2] = n(n+d+1)\mathbf{I}$,

$$\frac{1}{2} \mathbb{E}_{\mathbf{X}} \text{tr} \left(\mathbf{I} - \frac{\eta^*}{n} \mathbf{X} \mathbf{X}^\top \right)^2 = \frac{1}{2} \left(d - 2\eta^* + \frac{\eta^{*2}}{n} (n+d+1)d \right) = \Theta\left(\frac{d^2}{n}\right)$$

and we finish the proof by substituting n with $\tilde{\Theta}(d)$. \square

3.2 ONE-LAYER TRANSFORMER WITH CoT CAN IMPLEMENT MULTI-STEP GD

The previous subsection shows that the one-step solution by the one-layer transformer without CoT is not the endgame. Nevertheless, CoT can become the savior for this simple transformer because it enables the transformer to generate several intermediate computation steps to improve the final performance. The following theorem shows that with the reinforcement of CoT, there exists a one-layer transformer that can perform multi-step GD using intermediate generations. We show that $\Theta(\log d)$ steps of CoT can remarkably improve the performance, reducing the error from $\Theta(\frac{d}{\text{poly } \log d})$ to $O(1/\text{poly } d)$. With constant learning rate, $\Theta(\log d)$ steps of GD is also necessary to reconstruct \mathbf{w}^* accurately. The proof is deferred to Appendix B.2.

Theorem 3.2 (Informal). *There exists \mathbf{V}^* and \mathbf{W}^* s.t. $f_{\text{LSA}}(\mathbf{Z}_k)_{[:, -1]}$ outputs $(\mathbf{0}_d, 0, \mathbf{w}_k, 1)$ where $\mathbf{w}_k := (\mathbf{I} - (\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top)^k) \mathbf{w}^*$ is the k -step GD solution with learning rate η on a linear regression model. Moreover, if $n = \tilde{\Omega}(d)$, $k = \Omega(\log d)$, $\eta \in (0.1, 1)$, then the evaluation loss*

$$\mathcal{L}^{\text{Eval}}(\mathbf{V}^*, \mathbf{W}^*) = \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\left\| \left(\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \right)^{k+1} \mathbf{w}^* \right\|^2 \right] \leq O\left(\frac{1}{\text{poly}(d)}\right) \quad (9)$$

With the one-step GD solution in Theorem 3.1, the proof is straightforward: we assign the parameters (\mathbf{V}, \mathbf{W}) in the same form of Equation (8), with the η^* replaced by η . However, now the

transformer is allowed to generate k steps before reaching the final output. We can inductively calculate the i -th step of generation, showing that the output is exactly the i -th gradient step:

$$f_{\text{LSA}}(\mathbf{Z}_{i-1})_{[:, -1]} = (\mathbf{0}_d, 0, \mathbf{w}_i, 1), \quad i = 1, 2, \dots, k+1$$

After $k+1$ steps, we have the final output $(\mathbf{I} - (\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top)^{k+1}) \mathbf{w}^*$ by induction and the evaluation loss becomes Equation (9). By Lemma D.4, the final loss is upper bounded by $O\left(\frac{1}{\text{poly}(d)}\right)$. This is strictly better than a one-step GD solution by comparing with Corollary 3.1.

Now we theoretically display the expressivity improvement of transformers brought by CoT. In the following sections, we will further prove that **this separation is learnable** simply by gradient flow.

4 GRADIENT DYNAMICS OVER CHAIN OF THOUGHT

In this section, we go beyond the construction and prove our convergence result on the CoT objective. We show that the final solution found by gradient flow is approximately our construction in Theorem 3.2, which is significantly better than the one-step gradient descent solution without CoT.

4.1 MAIN RESULTS

According to our construction in Theorem 3.2, we use the following specific initialization to zero out the irrelevant blocks while keeping the essential blocks \mathbf{W}_{13} , \mathbf{V}_{31} , and w_{24} .

Assumption 4.1 (Initialization). *Let $\sigma > 0$ be a parameter. We assume the initialization of the parameters satisfies that*

$$\mathbf{V} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{V}_{31}(0) & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \mathbf{W} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{W}_{13}(0) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & w_{24} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}$$

Here $\mathbf{W}_{13}(0) = \sum_{i=1}^d \lambda_i^{\mathbf{W}} \mathbf{u}_i \mathbf{u}_i^\top$ and $\mathbf{V}_{31}(0) = \sum_{i=1}^d \lambda_i^{\mathbf{V}} \mathbf{u}_i \mathbf{u}_i^\top$ are symmetric and simultaneously diagonalizable, $\lambda_i^{\mathbf{V}} \leq -\sigma$, $\lambda_i^{\mathbf{W}} \in [\sigma, \frac{1}{2}]$. Further, we fix $w_{24} = -1$ for all $t > 0$.

This initialization follows Chen et al. (2024a) by assuming \mathbf{V}_{31} and \mathbf{W}_{13} share the same set of eigenvectors. It is close to the particular symmetric random initialization schemes discussed in Zhang et al. (2023) with a scaling factor σ . We use this specific initialization to zero out the irrelevant blocks along the training trajectory and facilitate the analysis in the early stages. To simplify the analysis of the complex dynamical system, we fix $w_{24} = -1$ to break the homogeneity of the model and avoid the occurrence of multiple global minimizers.

Now we prove that under appropriate initialization, gradient flow will nearly converge to the global minimizer. We provide a proof sketch in the next subsection. See Appendix C.3 for details.

Theorem 4.1 (Informal, Global Convergence). *Suppose $n = \tilde{\Omega}(d)$, $\eta \in (0.1, 0.9)$, $k = \Theta(\log d)$. Under Assumption 4.1 with $\sigma = \Theta(1)$, if we run gradient flow on the population loss in Equation (5), then after time $t = O(\log d + \log \frac{1}{\epsilon})$, we have $\mathcal{L}^{\text{CoT}}(t) \leq \epsilon$ for any $\epsilon \in (\frac{1}{\text{poly}(d)}, 1)$.*

4.2 PROOF IDEAS

In this subsection, we briefly outline the proof of Theorem 4.1.

Before analyzing the training dynamics, we will first prove that under Assumption 4.1, the gradient dynamics will only depend on the parameter blocks $\mathbf{W}_{13}(t)$, $\mathbf{V}_{31}(t)$, w_{24} , while other blocks stay zero (Lemma C.2). This is because our Gaussian data assumption makes sure the gradients on all the blocks are zero once they are initialized at zero, except for $\mathbf{W}_{13}(t)$, $\mathbf{V}_{31}(t)$, w_{24} . By this lemma, we can simplify the linear self-attention formula and consider the following equivalent yet simplified loss (we denote $\tilde{\mathbf{W}} := \mathbf{W}_{13}$, $\tilde{\mathbf{V}} := \mathbf{V}_{31}$, and w_{24} is fixed as -1):

$$\begin{aligned} \mathcal{L}^{\text{CoT}}(\theta) = & \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \sum_{i=0}^{k-1} \left\| \frac{1}{n} (\tilde{\mathbf{V}} \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{W}} + \eta \mathbf{X} \mathbf{X}^\top) \mathbf{w}_i - \frac{1}{n} (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right\|_2^2 \\ & + \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| \left(\mathbf{I} + \frac{1}{n} \tilde{\mathbf{V}} \mathbf{X} \mathbf{X}^\top \tilde{\mathbf{W}} \right) \mathbf{w}_k - \left(\frac{1}{n} \tilde{\mathbf{V}} \mathbf{X} \mathbf{X}^\top + \mathbf{I} \right) \mathbf{w}^* \right\|_2^2 \end{aligned}$$

For ease of presentation, we denote $\mathbf{S} := \frac{1}{n} \mathbf{X} \mathbf{X}^\top$. To analyze the gradient dynamics, we first need to compute the exact closed-form gradient instead of keeping the expectation. However, there exists difficulty calculating the closed form of the gradient: the formula involves the i -th step weight vector $\mathbf{w}_i = (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \mathbf{w}^*$, involving the higher order moments of the Wishart matrix \mathbf{S}^1 whose closed form is hard to obtain. Here, we provide a tighter estimate compared to previous work (Gatmiry et al., 2024) using the concentration of the Wishart matrix \mathbf{S} (Vershynin, 2018) when $n = \Theta(d \text{ poly } \log d)$ to estimate the expectation. In particular, we use the exponential decaying tail probability bound for the operator norm of the error $\delta \mathbf{S} := \mathbf{S} - \mathbf{I}$. For example, when estimating the expectation $\mathbb{E}[(\mathbf{I} - \eta \mathbf{S})^i]$, we can decompose the expectation into two cases: when $\|\delta \mathbf{S}\|_{op}$ is small, $(\mathbf{I} - \eta \mathbf{S})^i \approx (1 - \eta)^i \mathbf{I}$; when $\|\delta \mathbf{S}\|$ is larger than a threshold, the rest part of the expectation can be controlled by integrating the exponential decaying tail probability.² The concentration lemmas are provided in Appendix D.

The motivation behind a better concentration estimation is to ensure nearly independent dynamics along different eigenspaces $\{\mathbf{u}_i\}_{i=1}^d$ of $\tilde{\mathbf{W}}$ and $\tilde{\mathbf{V}}$. As an extreme case, we consider $n \rightarrow \infty$ and \mathbf{S} converges to \mathbf{I} almost surely. Now the gradient component on the $\mathbf{u}_i \mathbf{u}_i^\top$ subspace is only dependent on $\lambda_i^{\tilde{\mathbf{V}}}$ and $\lambda_j^{\tilde{\mathbf{W}}}$ without any other $\lambda_j^{\tilde{\mathbf{V}}}, \lambda_j^{\tilde{\mathbf{W}}}, j \neq i$ involved. That means there is no interaction between two different subspaces, i.e. the dynamics are independent. However, some interactions are introduced since the concentration error $\delta \mathbf{S} \neq 0$ when n is finite. Therefore, the improved characterization of the expected gradient is essential to upper bound the interaction between the dynamics of different eigenspaces $\{\mathbf{u}_i\}_{i=1}^d$, leading to a nearly independent evolution at initialization.

This independence property motivates us to conduct a stage-wise analysis. We first analyze the dynamics in **Stage 1** when the distance between the parameters $\tilde{\mathbf{V}}, \tilde{\mathbf{W}}$ and the ground-truth is larger than $O(1/\text{poly } \log d)$. In this stage, the bounded error can be dominated by the signal terms in the gradient, maintaining the nearly independent dynamics along each direction \mathbf{u}_i . After this stage, we enter **Stage 2** as a local convergence phase. We describe the dynamics below in detail.

Stage 1: $\tilde{\mathbf{W}}, \tilde{\mathbf{V}}$ converges to near-optimal. In this stage, the dynamics along each direction \mathbf{u}_i stay nearly independent. Specifically, we can expand the gradient flow dynamics for $\tilde{\mathbf{V}}, \tilde{\mathbf{W}}$ and project them into the eigenspaces $\mathbf{u}_i \mathbf{u}_i^\top$ to get the dynamics of the eigenvalues $\lambda_i^{\tilde{\mathbf{V}}} := \mathbf{u}_i^\top \tilde{\mathbf{V}} \mathbf{u}_i$, $\lambda_i^{\tilde{\mathbf{W}}} := \mathbf{u}_i^\top \tilde{\mathbf{W}} \mathbf{u}_i$. The dynamics of eigenvalues are characterized by the following Lemma 4.1 where we can prove that the interaction terms between different subspaces are bounded by $O(1/\log^2 d)$.

Lemma 4.1 (Informal version of Lemma C.6). *The dynamics of $\lambda_i^{\tilde{\mathbf{V}}}$ and $\lambda_i^{\tilde{\mathbf{W}}}$ are given by the following equations with $|\delta_j^{\tilde{\mathbf{V}}}| \leq O\left(\frac{1}{\log^2 d}\right), |\delta_j^{\tilde{\mathbf{W}}}| \leq O\left(\frac{1}{\log^2 d}\right)$:*

$$\begin{aligned} \frac{d\lambda_j^{\tilde{\mathbf{V}}}}{dt} &= - \left[(k+1) \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right)^2 + \frac{2}{\eta} \lambda_j^{\tilde{\mathbf{W}}} \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right) + \frac{\lambda_j^{\tilde{\mathbf{W}}^2}}{\eta(2-\eta)} \right] \lambda_j^{\tilde{\mathbf{V}}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{W}}} - 1 + \delta_j^{\tilde{\mathbf{V}}} \\ \frac{d\lambda_j^{\tilde{\mathbf{W}}}}{dt} &= \left[k + 1 - \frac{1}{\eta} \right] \lambda_j^{\tilde{\mathbf{V}}^2} \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right) + \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{V}}^2} \lambda_j^{\tilde{\mathbf{W}}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{V}}} - \delta_j^{\tilde{\mathbf{W}}}. \end{aligned}$$

This nearly independent evolution along each eigenvector \mathbf{u}_i enables us to analyze the individual dynamics of $\lambda_i^{\tilde{\mathbf{V}}}$ and $\lambda_i^{\tilde{\mathbf{W}}}$ at the beginning of training. Under Assumption 4.1, $\lambda_j^{\tilde{\mathbf{V}}}, \lambda_j^{\tilde{\mathbf{W}}}$ are initialized $\Theta(1)$. By Lemma 4.1, we prove by induction that the eigenvalues will go through two phases: (1) $\lambda_j^{\tilde{\mathbf{V}}}$ increases yet stay smaller than $-O\left(\frac{1}{k(1-\lambda_j^{\tilde{\mathbf{W}}})}\right)$, while $\lambda_j^{\tilde{\mathbf{W}}}$ increases to $1 - o(1)$. (2) $\lambda_j^{\tilde{\mathbf{W}}}$ stays $o(1)$ -close to 1, and $\lambda_j^{\tilde{\mathbf{V}}}$ also converges to $o(1)$ -close to $-\eta$. Here all $o(1)$ terms are some $O(1/\log^c d)$

¹To deal with the similar problem, Gatmiry et al. (2024) proposed a simple combinatorial method to estimate the expectation. We use the same technique to get a certain form of the expectation (see Appendix D), but the bound is not tight enough to get the desired results. See discussion in Appendix A.2.

²This method can keep the $(1-\eta)^i$ factor to prevent introducing unwanted estimation errors when i is large.

terms for some constant $c > 0$. That indicates that the distance between the eigenvalues and the target $|\lambda_j^{\tilde{V}} + \eta|, |\lambda_j^{\tilde{W}} - 1|$ converge to $O(1/\log^c d)$ for all $j \in [d]$ at the end of Stage 1.

Stage 2: Local convergence. One may expect that after Stage 1, the transformer can approximate gradient steps quite accurately since the parameter \tilde{V}, \tilde{W} are both $o(1)$ -close to ground-truth along each direction u_i . Unfortunately, the sum of error in d directions can still be $\tilde{\Theta}(d)$ since we can only reduce the error to $O(1/\text{poly } \log d)$ in each direction. Therefore, the solution still cannot recover the weight vector w^* at this stage. To address this issue, we further look into the exact form of the interaction terms $\delta_j^{\tilde{W}}, \delta_j^{\tilde{V}}$ and analyze the local convergence. By fine-grained expansion of the error terms, we notice that $\delta_j^{\tilde{W}}$ and $\delta_j^{\tilde{V}}$ are always coupled with some individual residual like $(1 - \lambda_j^{\tilde{W}})$, $(\eta + \lambda_j^{\tilde{V}})$, or some weighted average or those individual residuals. Meanwhile, the coefficient of the residual in the interaction terms is still upper bounded by $O(1/\text{poly } \log d)$. That enables us to derive some gradient lower bound similar to PL-conditions (Lemma C.12) when \tilde{V}, \tilde{W} are close to the ground-truth, leading to local convergence to near-optimal at a linear rate.

The final training error is some $O(\frac{1}{\text{poly } d})$, which depends on the inference step k and ground-truth η . Note that the optimal loss value is also at least polynomially small in d given $\Theta(\log d)$ CoT steps. Therefore, now we can conclude that the transformer can learn to implement multi-step GD when given intermediate ground-truth states after optimizing the CoT loss with gradient flow.

4.3 OUT-OF-DISTRIBUTION GENERALIZATION AT INFERENCE

In this section, we prove that after training, the transformer not only correctly predicts the weight vector in context with CoT generation, but also can generalize out-of-distribution (OOD). The following theorem shows that the trained transformer obtained from Theorem 4.1 with CoT generalizes over other problem instances when the input example sequence has an OOD covariance, as long as the covariance is not too ill-conditioned. Here $\mathcal{L}_{\Sigma}^{\text{Eval}}$ is defined as the OOD evaluation loss in eq. (7) with the in-context examples $x_i \sim \mathcal{N}(0, \Sigma)$ and weight vector $w^* \sim \mathcal{N}(0, I)$:

$$\mathcal{L}_{\Sigma}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) = \frac{1}{2} \mathbb{E}_{x_i \sim \mathcal{N}(0, \Sigma), w^*} \left[\left\| f_{\text{LSA}}(\hat{\mathbf{Z}}_k)_{[:, -1]} - (\mathbf{0}_d, 0, w^*, 1) \right\|^2 \right]$$

Theorem 4.2 (Informal, Theorem C.2). *Suppose $n = \tilde{\Omega}(d)$, $\eta \in (0.1, 0.9)$, $k' = \Theta(\log d)$. Assume the out-of-distribution covariance is well-conditioned: $\frac{\delta}{\eta} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq \frac{2-\delta}{\eta}$ for some constant $\delta > 0$. Then after training in Theorem 4.1, we have $\mathcal{L}_{\Sigma}^{\text{Eval}}(t) \leq \epsilon$ for any $\epsilon \in (\frac{1}{\text{poly}(d)}, 1)$.*

Note that this theorem covers both in-distribution (when $\eta = \delta$) and OOD tasks at evaluation, indicating that the transformer is trained to implement a general iterative optimization algorithm. Moreover, the inference step number k' in this theorem can go beyond the training CoT steps k , achieving better estimation for w^* .

One may think once the next-token-prediction training loss \mathcal{L}^{CoT} converges to the global minimizer based on ground-truth CoT data, the transformer naturally learns to do multi-step reasoning at inference, i.e. $\mathcal{L}^{\text{Eval}}$ is small. However, at the i -th generation step, the transformer is predicting the next weight token \hat{w}_{i+1} based on the previous generation \hat{w}_i instead of the ground-truth intermediate step w_i . It is possible that prediction error for each step accumulates or even increases exponentially.

Fortunately, the trained transformer guarantees a converging series of errors throughout the inference process, and we can expand and upper bound the sum of all the prediction errors at each step. That also ensures we can achieve any $O(\frac{1}{\text{poly}(d)})$ -small evaluation loss when we have $k' = \Theta(\log d)$ reasoning steps. The detailed proof is provided in Appendix C.4.

5 EXPERIMENTS

In this section, we introduce our experimental setup on our in-context weight vector prediction task to numerically validate our theoretical results. Specifically, we show that parameters of the transformer match the prediction of our theory when optimized over the CoT loss. Furthermore, we present the gap of evaluation loss $\mathcal{L}^{\text{Eval}}$ in Equation (7) between transformers with and without CoT.

Experimental Setup We train the transformer architecture in Equation (3) on the synthetic data. The data distribution follows our in-context weight prediction task in Equation (1). In particular, we choose the token dimensions $d = 10$, number of in-context examples $n = 20$, and GD learning rate $\eta = 0.4$ for generating the ground-truth intermediate states. We use a batch size $B = 1000$ and run Adam with learning rate $\alpha = 0.001$ for $\tau = 750$ iterations. More details refer to Appendix E.

Global convergence Our experiments show that the structure that weights of the full model exhibit is consistent with Theorem 3.2. At final convergence, all of the entries of \mathbf{W} converge to zero except the elements on the diagonal in the top-right corner block (the red box in the heatmap of \mathbf{W} , Figure 1), while all the entries of \mathbf{V} are near zero except elements on the diagonal in the bottom-left corner (the red box in the heatmap of \mathbf{V} , Figure 1). Also, the pattern shows $\mathbf{W}_{13} = \alpha \mathbf{I}$, $w_{24} = -\alpha$, and $\mathbf{V}_{31} = -\frac{\eta}{\alpha} \mathbf{I}$ with some scaling factor α ,³ which is equivalent to the construction stated in Theorem 3.2 and Theorem 4.1. That means the transformer implements one step of gradient descent $(\mathbf{0}_d, 0, -\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top (\mathbf{w}_i - \mathbf{w}^*), 0)$ before the residual connection, and the autoregressive CoT process enables model to perform multi-step GD.

Performance improvement We empirically verify the evaluation loss gap between transformers with and without CoT shown by Theorem 3.1 and Theorem 3.2. Our experiments in Figure 2 demonstrate that the evaluation loss of transformers with CoT converges to near zero even when $k = 10$. In comparison, the optimal expected loss that the one-layer linear transformer can achieve (the pink dashed line, from Corollary 3.1) is much larger than any of the model that applies multiple steps of computation. We also observe that evaluation loss at convergence keeps decreasing when the number of reasoning steps k increases from 10 to 40, which is consistent with Theorem C.1 where larger k allows for smaller error ϵ .

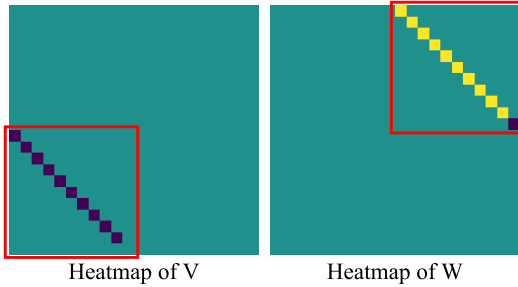


Figure 1: **Model weights:** We present the heatmap of the weights of the trained transformer. We initialize \mathbf{V}, \mathbf{W} randomly at $t = 0$, where $n = 20$, $d = 10$ and $k = 20$. After training, all entries of \mathbf{V} and \mathbf{W} converge to zero except the two blocks highlighted in the red box. Moreover, the pattern matches the theoretical results.

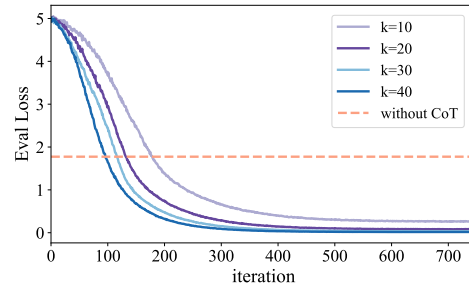


Figure 2: **k -step v.s. 1-step:** We plot the evaluation loss $\mathcal{L}^{\text{Eval}}$ when $n = 20$, $d = 10$. We randomly initialize the transformer. For transformers with CoT, loss converges to near zero while transformers without CoT cannot. Moreover, the loss at convergence decreases when k increases.

6 CONCLUSION

This paper investigates the training dynamics of transformers when the Chain of Thought (CoT) prompting is introduced. By focusing on the in-context weight prediction task, our theoretical results demonstrate that transformers can learn to implement iterative algorithms like multi-step GD with the enhancement of CoT, highlighting the essential role of CoT in multi-step reasoning tasks. Our empirical findings corroborate these theoretical insights, indicating that CoT prompting provides significant performance benefits.

There are still many open problems. Can we move beyond population loss on the in-context weight prediction task and show a sample complexity guarantee? Can CoT empower the transformer to acquire compositional reasoning capability instead of doing the same iterative steps?

³In Figure 1, $\alpha > 0$ while all $\alpha \neq 0$ works for the construction. Empirically, the sign of α depends on the random initialization, and both positive and negative solutions exist.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Linear attention is (maybe) all you need (to understand transformer optimization). *arXiv preprint arXiv:2310.01082*, 2023.
- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556, 2022.
- Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism. *Advances in Neural Information Processing Systems*, 36:48314–48362, 2023.
- Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*, 2023.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36, 2024.
- Vivien Cabannes, Charles Arnal, Wassim Bouaziz, Alice Yang, Francois Charton, and Julia Kempe. Iteration head: A mechanistic study of chain-of-thought. *arXiv preprint arXiv:2406.02128*, 2024.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv preprint arXiv:2402.19442*, 2024a.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Unveiling induction heads: Provable training dynamics and feature learning in transformers. *arXiv preprint arXiv:2409.10559*, 2024b.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.
- Nan Ding, Tomer Levinboim, Jialin Wu, Sebastian Goodman, and Radu Soricut. Causallm is not optimal for in-context learning. *arXiv preprint arXiv:2308.06912*, 2023.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Benjamin L Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution of statistical induction heads: In-context learning markov chains. *arXiv preprint arXiv:2402.11004*, 2024.

- Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. Towards revealing the mystery behind chain of thought: a theoretical perspective. *Advances in Neural Information Processing Systems*, 36, 2024.
- Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order optimization methods for in-context learning: A study with linear models. *arXiv preprint arXiv:2310.17086*, 2023.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Khashayar Gatmiry, Nikunj Saunshi, Sashank J. Reddi, Stefanie Jegelka, and Sanjiv Kumar. Can looped transformers learn to implement multi-step gradient descent for in-context learning? In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=o8AaRKbP9K>.
- Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. *arXiv preprint arXiv:2301.13196*, 2023.
- Yifan Hou, Jiaoda Li, Yu Fei, Alessandro Stolfo, Wangchunshu Zhou, Guangtao Zeng, Antoine Bosselut, and Mrinmaya Sachan. Towards a mechanistic interpretation of multi-step reasoning capabilities of language models. *arXiv preprint arXiv:2310.14491*, 2023.
- Xinyang Hu, Fengzhuo Zhang, Siyu Chen, and Zhuoran Yang. Unveiling the statistical foundations of chain-of-thought prompting methods. *arXiv preprint arXiv:2408.14511*, 2024.
- Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. *arXiv preprint arXiv:2310.05249*, 2023.
- Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. *Advances in Neural Information Processing Systems*, 35:37822–37836, 2022.
- Juno Kim and Taiji Suzuki. Transformers learn nonlinear features in context: Nonconvex mean-field dynamics on the attention landscape. *arXiv preprint arXiv:2402.01258*, 2024.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. *arXiv preprint arXiv:2302.06015*, 2023.
- Yingcong Li, Yixiao Huang, Muhammed E Ildiz, Ankit Singh Rawat, and Samet Oymak. Mechanics of next token prediction with self-attention. In *International Conference on Artificial Intelligence and Statistics*, pp. 685–693. PMLR, 2024a.
- Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024b.
- Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022.
- Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. *arXiv preprint arXiv:2307.03576*, 2023.
- Ashok Vardhan Makkuva, Marco Bondaschi, Adway Girish, Alliot Nagle, Martin Jaggi, Hyeji Kim, and Michael Gastpar. Attention with markov: A framework for principled analysis of transformers via markov chains. *arXiv preprint arXiv:2402.04161*, 2024.

- William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought. *arXiv preprint arXiv:2310.07923*, 2023a.
- William Merrill and Ashish Sabharwal. The parallelism tradeoff: Limitations of log-precision transformers. *Transactions of the Association for Computational Linguistics*, 11:531–545, 2023b.
- Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with gradient descent. *arXiv preprint arXiv:2402.14735*, 2024.
- Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. Show your work: Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*, 2021.
- OpenAI. Gpt-4 technical report, 2023.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.
- Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. Refiner: Reasoning feedback on intermediate representations. *arXiv preprint arXiv:2304.01904*, 2023.
- Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. Language models are multilingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057*, 2022.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines. *arXiv preprint arXiv:2308.16898*, 2023.
- Yuangdong Tian, Yiping Wang, Beidi Chen, and Simon Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. *arXiv preprint arXiv:2305.16380*, 2023a.
- Yuangdong Tian, Yiping Wang, Beidi Chen, and Simon Du. Scan and Snap: Understanding Training Dynamics and Token Composition in 1-layer Transformer, July 2023b. URL <http://arxiv.org/abs/2305.16380>. arXiv:2305.16380 [cs].
- Yuangdong Tian, Yiping Wang, Zhenyu Zhang, Beidi Chen, and Simon Du. Joma: Demystifying multilayer transformers via joint dynamics of mlp and attention. *arXiv preprint arXiv:2310.00535*, 2023c.
- Rasul Tutunov, Antoine Grosnit, Juliusz Ziomek, Jun Wang, and Haitham Bou-Ammar. Why can large language models generate correct chain-of-thoughts? *arXiv preprint arXiv:2310.13571*, 2023.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. Towards understanding chain-of-thought prompting: An empirical study of what matters. *arXiv preprint arXiv:2212.10001*, 2022a.

- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022b.
- Zixuan Wang, Stanley Wei, Daniel Hsu, and Jason D Lee. Transformers provably learn sparse token selection while fully-connected nets cannot. In *Forty-first International Conference on Machine Learning*, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.
- Skyler Wu, Eric Meng Shen, Charumathi Badrinath, Jiaqi Ma, and Himabindu Lakkaraju. Analyzing chain-of-thought prompting in large language models via gradient-based feature attributions. *arXiv preprint arXiv:2307.13339*, 2023.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*, 2023.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.

A DISCUSSION AND LIMITATION

A.1 RELATED WORKS ON EXPRESSIVENESS

Our work is closely related to the previous works in multi-step GD using multi-layer attention layers, including Bai et al. (2023); Fu et al. (2023); Ding et al. (2023); Ahn et al. (2024); Giannou et al. (2023); Gatmiry et al. (2024). These works guarantee that transformers are **expressive enough** to do in-context learning by implementing gradient descent, and they serve as the foundation of our work which focuses on **optimization**. Most of them focus on the in-context learning setup as the testbed so we naturally follow the setup to understand the advantage of CoT.

Most of the above works on **expressiveness** focus on those iterative algorithms, e.g. (pre-conditioned) gradient descent on various objectives (Bai et al., 2023; Ahn et al., 2024; Ding et al., 2023), Newton methods/matrix inverse (Giannou et al., 2023), etc. Those papers have similar constructive proof techniques using multi-layer transformers: they construct a basic block(s) to represent one step of some iterative algorithm and stack them up to do multi-steps of that algorithm. Sometimes the blocks can be even the same, which means a “looped” transformer, i.e. implementing the same transformer blocks several times as a loop, can express those algorithms. In our warm-up construction for a better understanding of the setup, we use similar techniques to construct the linear transformer that allows auto-regressive generation to iteratively implement the block. However, we require the practical auto-regressive setting, which is novel in the literature.

Most importantly, despite the close relation between our work and those previous expressiveness papers, our work mainly focuses on the **optimization** perspective. It is a big step beyond expressiveness because there is no guarantee that one can algorithmically find the constructed solutions in the previous work. Ahn et al. (2024); Gatmiry et al. (2024) are the only two papers related to optimization of multi-layer transformers over in-context linear regression setup. Ahn et al. (2024) analyzed the global optimizer/critical points for multi-layer transformers, but they didn’t prove that any gradient-based algorithm can reach those solutions. Compared to all the works above, our proof techniques for the main theorems are completely orthogonal and **not** straightforward extensions of the previous papers like Bai et al. (2023).

Gatmiry et al. (2024) is the most related work to us. They also proved some results on **learning** to implement multi-step GD by looped transformer. We will highlight the differences and **our novel contributions** of our work in the next section.

A.2 DISCUSSION ON GATMIRY ET AL. (2024)

In this section, we compare our work with Gatmiry et al. (2024). We begin by outlining the similarities and connections between the two works before highlighting our theoretical contributions in contrast to Gatmiry et al. (2024).

Both Gatmiry et al. (2024) and our study analyze the dynamics of a one-layer linear transformer in the context of a linear regression task, demonstrating that transformers can implement multi-step gradient descent. We adopt similar architectural frameworks to those in Zhang et al. (2023); Ahn et al. (2024; 2023); Mahankali et al. (2023), as well as several other works. The key connection between our work and Gatmiry et al. (2024) lies in the observation that both looped transformers and transformers with CoT prompting through autoregressive generation are capable of naturally implementing iterative algorithms like gradient descent.

However, our data model and training objective are intrinsically different from those in Gatmiry et al. (2024), leading to distinct insights. While Gatmiry et al. (2024) focuses on an ICL setting for linear regression tasks involving examples and a query, our task is centered on predicting the ground-truth weight vector w^* within context, i.e. in-context weight prediction. The final converging solutions are totally different, even though they both are equivalent to some type of GD. From the perspective of the training objective, Gatmiry et al. (2024) uses a standard squared loss over the ICL objective. In contrast, we use a sum of squared losses across all intermediate steps, corresponding to the CoT loss defined in Equation (6). Therefore, we highlight the effectiveness in improving the performance of the CoT prompting on a shallow transformer, while Gatmiry et al. (2024) stress a multi-layer transformer with shared weights (looped transformer) can do multi-step GD through the layers.

From a technical perspective, Gatmiry et al. (2024) fix the outer layer and train only the matrix \mathbf{A} , which is analogous to our matrix \mathbf{W} . In contrast, our work allows for training both layers of the transformer, providing a stronger analysis of training dynamics. Our proof strategy is also novel, given that our training dynamics are more complicated: obtaining our final solution requires solving a challenging d -dimensional dynamical system, whereas prior work in ICL reduces the outer layer to a scalar.

As a more profound theoretical contribution, we rigorously establish a clear performance gap between the one-layer transformer without CoT and the ones with CoT. Specifically, the one-layer transformer without CoT is restricted to a single step of GD, with the final error $\Theta(d/\text{poly } \log d)$, while a one-layer transformer with CoT can achieve a $O(1/\text{poly } d)$ loss with only $\Theta(\log d)$ steps. On the other hand, Gatmiry et al. (2024) do not show their transformer implementing the multi-step GD can outperform the transformer with one-step GD. According to their Theorem 4.2, their looped transformer can only provably get the final loss down to $\frac{d^{5/2}L \cdot 4^L}{\sqrt{n}}$, where L is the number of the loops. However, a one-layer transformer can achieve $\Theta(d^2/n)$ loss by implementing one-step of GD, **which is asymptotically better than the multi-step solution in Gatmiry et al. (2024).**

We conjecture the gap between our analysis lies in our different methods of calculating the terms in the gradient concerning Wishart matrices. For intuition, we introduce **the novel expectation calculation method** in Section 4, which asymptotically improves the estimation of higher moments of Wishart matrices in Gatmiry et al. (2024). We adopt the combinatorial technique in Gatmiry et al. (2024) to compute the form of $\mathbb{E}[\mathbf{S}\mathbf{A}\mathbf{S}^k\mathbf{T}\mathbf{S}^{k'}]$, but when we calculate the expected gradient we use the concentration tail bound technique to calculate the expectation. That enables us to better approximate the expectation. We hypothesize that applying our techniques could potentially demonstrate that their looped transformers outperform those without loops in the ICL setting.

A.3 LIMITATION AND FUTURE DIRECTIONS

Architecture and parameterization In this work, we use the single-layer linear transformer to analyze the training dynamics. Moreover, we adopt the same reparameterization and similar initialization in previous works (Zhang et al., 2023; Tian et al., 2023a; Chen et al., 2024a; Mahankali et al., 2023; Ahn et al., 2024). It deviates from the practical softmax attention with $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ parameterization and random initialization, which is a limitation of this work.

However, analyzing the linear counterpart of the model before targeting the more difficult practical models is common in the development of learning theory. As for linear attention, the connection between linear attention and softmax attention is also partially justified by the empirical observations in Ahn et al. (2023). Analyzing the dynamics using more practical architectures will be a very important and fundamental future direction.

Population loss and sample complexity Following most of the previous work, we use population loss when analyzing the training trajectory instead of using finite sample loss. This modification is to simplify the analysis and focus on the population dynamics without noise. A possible future step is to generalize this analysis to a finite sample setting and train the model with online SGD.

CoT on iterative tasks In this work, we mainly focus on **iterative** tasks, one of the simplest forms where multi-step CoT can help yield better performance. That serves as the initial step towards understanding why CoT helps reasoning following the first principle. As a limitation, though CoT can empower the transformer to acquire compositional reasoning capability instead of doing the same iterative step, it is a much harder question beyond our paper’s scope. It is a very important future direction and definitely worth further exploring.

B PROOFS OF THEOREMS IN SECTION 3

In this section, we prove the expressiveness results of the linear transformers with and without CoT. In Appendix B.1, we prove that a one-layer linear transformer without CoT can only obtain the one-step gradient descent solution. In Appendix B.2, we prove that there exists a one-layer linear

transformer that implements multi-step gradient descent with the CoT prompting. As corollaries, there exists a separation between the one-step and multi-step solutions.

B.1 PROOF OF THEOREM 3.1

We first restate the theorem:

Theorem B.1 (Lower bound without CoT). *If the global minimizer of $\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W})$ is $(\mathbf{V}^*, \mathbf{W}^*)$, the corresponding one-layer transformer $f_{\text{LSA}}(\mathbf{Z}_0)_{[:, -1]}$ implements one step GD on a linear model with some learning rate $\eta = \frac{n}{n+d+1}$ and the transformer outputs $\frac{\eta}{n} \mathbf{X} \mathbf{y}^\top$.*

Proof. Recall the loss expression in Equation (5) when $k = 0$,

$$\begin{aligned} \mathcal{L}(\mathbf{V}, \mathbf{W}) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| f_{\text{LSA}}(\mathbf{Z}_0)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}^*, 1) \right\|^2 \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| \mathbf{V} \mathbf{Z}_0 \cdot \frac{\mathbf{Z}_0^\top \mathbf{W} \mathbf{Z}_0_{[:, -1]}}{n} - (\mathbf{0}_d, 0, \mathbf{w}^*, 0)^\top \right\|^2 \quad (\text{since } \mathbf{w}_0 = \mathbf{0}_d.) \end{aligned}$$

The key insight of the proof is to replace the \mathbf{w}^* with the one-step GD solution $\frac{\eta}{n} \mathbf{X} \mathbf{y}^\top$,

$$\mathcal{L}(\mathbf{V}, \mathbf{W}) = \frac{1}{2} \mathbb{E} \left[\left\| \mathbf{V} \mathbf{Z}_0 \cdot \frac{\mathbf{Z}_0^\top \mathbf{W} \mathbf{Z}_0_{[:, -1]}}{n} - \left(\mathbf{0}_d, 0, \frac{\eta}{n} \mathbf{X} \mathbf{y}^\top, 0 \right)^\top \right\|^2 \right] + C$$

After proving this property, we can conclude that the optimal solution without CoT is exactly the one-step solution $\frac{\eta}{n} \mathbf{X} \mathbf{y}^\top$. We prove this result by showing the gradient of those two loss functions are the same.

First, before calculating the gradient, we extract the identical parts of the loss. Notice that the ground-truth entries are all zero in $i = 1, 2, \dots, d, d+1, 2d+2$ positions in both expressions. Therefore, that part of error is the norm of the output $f_{\text{LSA}}(\mathbf{Z}_0)_{[:, -1]}$ in those corresponding entries:

$$\frac{1}{2} \mathbb{E} \left[\left\| \mathbf{V} \mathbf{Z}_0 \cdot \frac{\mathbf{Z}_0^\top \mathbf{W} \mathbf{Z}_0_{[1:d+1, -1]}}{n} \right\|^2 \right] + \frac{1}{2} \mathbb{E} \left[\left\| \mathbf{V} \mathbf{Z}_0 \cdot \frac{\mathbf{Z}_0^\top \mathbf{W} \mathbf{Z}_0_{[2d+2, -1]}}{n} \right\|^2 \right]$$

which is the same for both expressions. Therefore, we just need to consider

$$f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} = \mathbf{V}_{[d+2:2d+1, :]} \mathbf{Z}_0 \cdot \frac{\mathbf{Z}_0^\top \mathbf{W} \mathbf{Z}_0_{[:, -1]}}{n},$$

which corresponds to the ground-truth signals. Here $\mathbf{V}_{[d+2:2d+1, :]} = [\mathbf{V}_{31}, \mathbf{V}_{32}, \mathbf{V}_{33}, \mathbf{V}_{34}]$. We only need to prove that

$$\mathbb{E} \left\| f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} - \mathbf{w}^* \right\|^2 = \mathbb{E} \left\| f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right\|^2 + C$$

for some constant C .

We show the gradients of both sides are the same, and equivalently the differential of both sides should be the same. The differential of L.H.S. is

$$\begin{aligned} & d \left(\mathbb{E} \left\| f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} - \mathbf{w}^* \right\|^2 \right) \\ &= 2 \mathbb{E} \left[(f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} - \mathbf{w}^*)^\top d f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} \right] \end{aligned}$$

and the differential of R.H.S. is

$$\begin{aligned} & d \left(\mathbb{E} \left\| f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right\|^2 \right) \\ &= 2 \mathbb{E} \left[(f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^*)^\top d f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} \right] \end{aligned}$$

Therefore, we only need to prove that

$$\mathbb{E} \left[\mathbf{w}^{*\top} \mathrm{d}f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} \right] = \mathbb{E} \left[\left(\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right)^\top \mathrm{d}f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} \right] \quad (10)$$

We expand this expression $f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]}$ (Note that now we don't have the assumption of initialization):

$$\begin{aligned} & \mathbf{V}_{[d+2:2d+1, :]} \mathbf{Z}_0 \cdot \frac{\mathbf{Z}_0^\top \mathbf{W} \mathbf{Z}_0[:, -1]}{n} \\ &= \frac{1}{n} [\mathbf{V}_{31} \quad \mathbf{V}_{32} \quad \mathbf{V}_{33} \quad \mathbf{V}_{34}] \begin{bmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{y} & 0 \\ \mathbf{0}_{d \times n} & \mathbf{w}_0 \\ \mathbf{0}_{1 \times n} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{X}^\top & \mathbf{y}^\top & \mathbf{0}_{n \times d} & \mathbf{0}_n \\ \mathbf{0}_{1 \times d} & 0 & \mathbf{w}_0^\top & 1 \end{bmatrix} \mathbf{W} \begin{bmatrix} \mathbf{0} \\ 0 \\ \mathbf{w}_0 \\ 1 \end{bmatrix} \\ &= \frac{1}{n} [\mathbf{V}_{31} \quad \mathbf{V}_{32} \quad \mathbf{V}_{33} \quad \mathbf{V}_{34}] \begin{bmatrix} \mathbf{X} \mathbf{X}^\top & \mathbf{X} \mathbf{y}^\top & \mathbf{0}_{d \times d} & \mathbf{0}_d \\ \mathbf{y} \mathbf{X}^\top & \mathbf{y} \mathbf{y}^\top & \mathbf{0}_{1 \times d} & 0 \\ \mathbf{0}_{d \times d} & \mathbf{0}_d & \mathbf{0}_{d \times d} & \mathbf{0}_d \\ \mathbf{0}_{1 \times d} & 0 & \mathbf{0}_{1 \times d} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{W}_{14} \\ w_{24} \\ \mathbf{W}_{34} \\ w_{44} \end{bmatrix} \quad (\text{since } \mathbf{w}_0 = \mathbf{0}_d) \\ &= \frac{1}{n} [\mathbf{V}_{31} \quad \mathbf{V}_{32} \quad \mathbf{V}_{33} \quad \mathbf{V}_{34}] \begin{bmatrix} \mathbf{X} \mathbf{X}^\top \mathbf{W}_{14} + w_{24} \mathbf{X} \mathbf{y}^\top \\ \mathbf{y} \mathbf{X}^\top \mathbf{W}_{14} + w_{24} \mathbf{y} \mathbf{y}^\top \\ \mathbf{0}_d \\ w_{44} \end{bmatrix} \\ &= \frac{1}{n} \left(\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top} \right) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) + \frac{\mathbf{V}_{34} w_{44}}{n} \quad (\mathbf{y} = \mathbf{X}^\top \mathbf{w}^*) \end{aligned}$$

and the differential of $f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]}$ is

$$\begin{aligned} & \mathrm{d}f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]} \\ &= \mathrm{d} \left(\frac{1}{n} \left(\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top} \right) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \right) + \mathrm{d} \frac{\mathbf{V}_{34} w_{44}}{n} \\ &= \frac{1}{n} \left(\mathrm{d}\mathbf{V}_{31} + \mathrm{d}\mathbf{V}_{32} \mathbf{w}^{*\top} \right) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) + \frac{1}{n} (\mathrm{d}\mathbf{V}_{34} \cdot w_{44} + \mathbf{V}_{34} \mathrm{d}w_{44}) \\ & \quad + \frac{1}{n} \left(\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top} \right) \mathbf{X} \mathbf{X}^\top (\mathrm{d}\mathbf{W}_{14} + \mathrm{d}w_{24} \mathbf{w}^*) \end{aligned}$$

Now, to prove Equation (10), we compare the differential for each parameter on both sides. For all parameter, we start from the left side and prove it equal to the right.

V₃₁: The \mathbf{V}_{31} term of differential in $\mathrm{d}f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]}$ is $\frac{1}{n} \mathrm{d}\mathbf{V}_{31} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*)$,

$$\begin{aligned} & \mathbb{E} \left[\mathbf{w}^{*\top} \cdot \frac{1}{n} \mathrm{d}\mathbf{V}_{31} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \right] \\ &= \mathbb{E} \left[\text{tr} \left(\mathbf{w}^{*\top} \cdot \frac{1}{n} \mathrm{d}\mathbf{V}_{31} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \right) \right] \quad (\text{It is a scalar in the trace.}) \\ &= \mathbb{E} \left[\text{tr} \left(\frac{1}{n} \mathrm{d}\mathbf{V}_{31} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \mathbf{w}^{*\top} \right) \right] \\ &= \mathbb{E} [\text{tr} (\mathrm{d}\mathbf{V}_{31} w_{24})] \quad (\mathbb{E} [\mathbf{X} \mathbf{X}^\top] = n \mathbf{I}_d, \mathbb{E} [\mathbf{w}^*] = 0, \mathbb{E} [\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}_d.) \\ &= \mathbb{E} \left[\text{tr} \left(\frac{\eta}{n^2} \cdot \mathrm{d}\mathbf{V}_{31} w_{24} \mathbf{X} \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \right) \right] \quad (\mathbb{E} [(\mathbf{X} \mathbf{X}^\top)^2] = n(n+d+1) \mathbf{I}_d, \eta = \frac{n}{n+d+1}.) \\ &= \mathbb{E} \left[\text{tr} \left(\frac{\eta}{n^2} \cdot \mathbf{X} \mathbf{X}^\top \mathrm{d}\mathbf{V}_{31} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \mathbf{w}^{*\top} \right) \right] \quad (\mathbb{E} [\mathbf{w}^*] = 0, \mathbb{E} [\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}_d.) \\ &= \mathbb{E} \left[\left(\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right)^\top \cdot \frac{1}{n} \mathrm{d}\mathbf{V}_{31} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \right] \end{aligned}$$

So those two $\mathrm{d}\mathbf{V}_{31}$ terms are identical.

V₃₂: The \mathbf{V}_{32} term of differential in $\mathrm{d}f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1, -1]}$ is $\frac{\mathrm{d}\mathbf{V}_{32}}{n} \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*)$,

$$\mathbb{E} \left[\mathbf{w}^{*\top} \cdot \frac{\mathrm{d}\mathbf{V}_{32}}{n} \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \right]$$

$$\begin{aligned}
&= \mathbb{E} \left[\text{tr} \left(\mathbf{w}^{*\top} \cdot \frac{d\mathbf{V}_{32}}{n} \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \right) \right] && \text{(It is a scalar in the trace.)} \\
&= \mathbb{E} \left[\text{tr} \left(\frac{d\mathbf{V}_{32}}{n} \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \mathbf{w}^{*\top} \right) \right] \\
&= \mathbb{E} \left[\text{tr} \left(\frac{d\mathbf{V}_{32}}{n} \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top \mathbf{W}_{14} \mathbf{w}^{*\top} \right) \right] && (\mathbb{E}[\mathbf{w}^*] = \mathbf{0} \text{ and } \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \mathbf{w}^{*\top} \text{ is odd}) \\
&= \mathbb{E} \left[\text{tr} \left(\frac{d\mathbf{V}_{32}}{n} \mathbf{W}_{14}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \mathbf{w}^{*\top} \right) \right] && (\mathbf{W}_{14}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \text{ is a scalar.}) \\
&= \mathbb{E} [\text{tr} (d\mathbf{V}_{32} \mathbf{W}_{14}^\top)] && (\mathbb{E}[\mathbf{X} \mathbf{X}^\top] = n\mathbf{I}_d, \mathbb{E}[\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}_d.) \\
&= \mathbb{E} \left[\text{tr} \left(\frac{\eta}{n^2} \cdot d\mathbf{V}_{32} \mathbf{W}_{14}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \right) \right] && (\mathbb{E}[(\mathbf{X} \mathbf{X}^\top)^2] = n(n+d+1)\mathbf{I}_d, \eta = \frac{n}{n+d+1}.) \\
&= \mathbb{E} \left[\text{tr} \left(\frac{\eta}{n^2} \cdot \mathbf{X} \mathbf{X}^\top d\mathbf{V}_{32} \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \mathbf{w}^{*\top} \right) \right] && (\mathbb{E}[\mathbf{w}^*] = \mathbf{0}, \mathbb{E}[\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}_d.) \\
&= \mathbb{E} \left[\left(\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right)^\top \cdot \frac{1}{n} d\mathbf{V}_{32} \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{14} + w_{24} \mathbf{w}^*) \right]
\end{aligned}$$

So those two $d\mathbf{V}_{32}$ terms are identical.

\mathbf{V}_{34} : The \mathbf{V}_{34} term of differential in $df_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1,-1]}$ is $\frac{1}{n} d\mathbf{V}_{34} w_{44}$,

$$\mathbb{E} \left[\mathbf{w}^{*\top} \frac{1}{n} d\mathbf{V}_{34} w_{44} \right] = 0 = \mathbb{E} \left[\left(\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right)^\top \frac{1}{n} d\mathbf{V}_{34} w_{44} \right]$$

since $\mathbb{E}[\mathbf{w}^*] = \mathbf{0}_d$. Therefore those two are equal.

\mathbf{W}_{14} : The \mathbf{W}_{14} term of differential in $df_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1,-1]}$ is $\frac{1}{n} (\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top d\mathbf{W}_{14}$,

$$\begin{aligned}
&\mathbb{E} \left[\mathbf{w}^{*\top} \cdot \frac{1}{n} (\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top d\mathbf{W}_{14} \right] \\
&= \mathbb{E} \left[\text{tr} \left(\mathbf{w}^{*\top} \cdot \frac{1}{n} (\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top d\mathbf{W}_{14} \right) \right] && \text{(It is a scalar in the trace.)} \\
&= \mathbb{E} \left[\text{tr} \left(\frac{1}{n} (\mathbf{w}^{*\top} \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top d\mathbf{W}_{14} \right) \right] && (\mathbb{E}[\mathbf{w}^*] = \mathbf{0}_d.) \\
&= \mathbb{E} \left[\text{tr} \left(\frac{1}{n} (\mathbf{V}_{32}^\top \mathbf{w}^* \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top d\mathbf{W}_{14} \right) \right] && (\mathbf{V}_{32}^\top \mathbf{w}^* \text{ is a scalar.}) \\
&= \mathbb{E} [\text{tr} (\mathbf{V}_{32}^\top d\mathbf{W}_{14})] && (\mathbb{E}[\mathbf{X} \mathbf{X}^\top] = n\mathbf{I}_d, \mathbb{E}[\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}_d.) \\
&= \mathbb{E} \left[\text{tr} \left(\frac{\eta}{n^2} \cdot \mathbf{X} \mathbf{X}^\top \mathbf{V}_{32}^\top \mathbf{X} \mathbf{X}^\top d\mathbf{W}_{14} \right) \right] && (\mathbb{E}[(\mathbf{X} \mathbf{X}^\top)^2] = n(n+d+1)\mathbf{I}_d, \eta = \frac{n}{n+d+1}.) \\
&= \mathbb{E} \left[\left(\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right)^\top \cdot \frac{1}{n} (\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top d\mathbf{W}_{14} \right]
\end{aligned}$$

Thus the two $d\mathbf{W}_{14}$ terms are the same.

w_{24} : The w_{24} term in $df_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1,-1]}$ is $\frac{1}{n} (\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top dw_{24} \mathbf{w}^*$,

$$\begin{aligned}
&\mathbb{E} \left[\mathbf{w}^{*\top} \cdot \frac{1}{n} (\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top \mathbf{w}^* dw_{24} \right] \\
&= \mathbb{E} \left[\text{tr} \left(\mathbf{w}^{*\top} \cdot \frac{1}{n} (\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top}) \mathbf{X} \mathbf{X}^\top \mathbf{w}^* dw_{24} \right) \right] && \text{(It is a scalar in the trace.)} \\
&= \mathbb{E} \left[\text{tr} \left(\frac{1}{n} (\mathbf{w}^{*\top} \mathbf{V}_{31}) \mathbf{X} \mathbf{X}^\top \mathbf{w}^* dw_{24} \right) \right] && (\mathbb{E}[\mathbf{w}^*] = \mathbf{0}_d.) \\
&= \mathbb{E} [\text{tr} (\mathbf{V}_{31} dw_{24})] && (\mathbb{E}[\mathbf{X} \mathbf{X}^\top] = n\mathbf{I}_d, \mathbb{E}[\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}_d.) \\
&= \mathbb{E} \left[\text{tr} \left(\frac{\eta}{n^2} \cdot \mathbf{X} \mathbf{X}^\top \mathbf{V}_{31} \mathbf{X} \mathbf{X}^\top dw_{24} \right) \right] && (\mathbb{E}[(\mathbf{X} \mathbf{X}^\top)^2] = n(n+d+1)\mathbf{I}_d, \eta = \frac{n}{n+d+1}.)
\end{aligned}$$

$$= \mathbb{E} \left[\left(\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right)^\top \cdot \frac{1}{n} \left(\mathbf{V}_{31} + \mathbf{V}_{32} \mathbf{w}^{*\top} \right) \mathbf{X} \mathbf{X}^\top \mathbf{w}^* d\mathbf{w}_{24} \right]$$

Therefore the differential for w_{24} are the same.

w_{44} : The w_{44} term of differential in $d f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1,-1]}$ is $\frac{1}{n} \mathbf{V}_{34} d\mathbf{w}_{44}$,

$$\mathbb{E} \left[\mathbf{w}^{*\top} \frac{1}{n} \mathbf{V}_{34} d\mathbf{w}_{44} \right] = 0 = \mathbb{E} \left[\left(\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right)^\top \frac{1}{n} \mathbf{V}_{34} d\mathbf{w}_{44} \right]$$

since $\mathbb{E}[\mathbf{w}^*] = \mathbf{0}_d$. Therefore those two are also equal.

In conclusion, Equation (10) holds since all the differential terms are equal. Therefore, $\exists C$

$$\mathbb{E} \left\| f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1,-1]} - \mathbf{w}^* \right\|^2 = \mathbb{E} \left\| f_{\text{LSA}}(\mathbf{Z}_0)_{[d+2:2d+1,-1]} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{w}^* \right\|^2 + C$$

which finishes our proof. \square

B.2 PROOF OF THEOREM 3.2

Here we restate the Theorem 3.2 and provide the detailed proof.

Theorem B.2. Suppose $n = \Theta(d \log^5 d)$, $k \geq C \log d$, $\eta \in (0.1, 0.9)$. There exists \mathbf{V}^* and \mathbf{W}^* s.t. $f_{\text{LSA}}(\mathbf{Z}_k)_{[:, -1]}$ outputs $(\mathbf{0}_d, 0, \mathbf{w}_{k+1}, 1)$ where $\mathbf{w}_i := (\mathbf{I} - (\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top)^i) \mathbf{w}^*$ is the k -step GD solution with learning rate η on a linear regression model. Moreover, the evaluation loss

$$\mathcal{L}^{\text{Eval}}(\mathbf{V}^*, \mathbf{W}^*) = \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\left\| \left(\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \right)^{k+1} \mathbf{w}^* \right\|^2 \right] \leq \frac{1}{d^{C \log(\frac{1}{1-\eta})}} \quad (11)$$

Proof. We construct \mathbf{V}^* and \mathbf{W}^* in the following way,

$$\mathbf{V}^* = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ -\eta \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \mathbf{W}^* = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & -1 \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (12)$$

Now the transformer is allowed to generate k steps before reaching the final output. We can inductively calculate the i -th step of generation, showing that the output is exactly the parameter after i -th gradient step ($i = 1, 2, \dots, k+1$):

$$\begin{aligned} f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} &= (\mathbf{0}_d, 0, \mathbf{w}_i, 1) + \mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_i_{[:, -1]}}{n} \\ &= (\mathbf{0}_d, 0, \mathbf{w}_i, 1) + \frac{1}{n} (\mathbf{0}_d, 0, \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*), 0) \\ &= (\mathbf{0}_d, 0, \mathbf{w}_i, 1) + (\mathbf{0}_d, 0, -\frac{\eta}{n} \mathbf{X} \mathbf{X}^\top (\mathbf{w}_i - \mathbf{w}^*), 0) \\ &= (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \end{aligned}$$

After $k+1$ steps, we have the final output $(\mathbf{I} - (\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top)^{k+1}) \mathbf{w}^*$ by induction and the evaluation loss becomes Equation (9). By Lemma D.4, the final loss is

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\left\| \left(\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \right)^{k+1} \mathbf{w}^* \right\|^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\text{tr} \left(\left(\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \right)^{2k+2} \right) \right] \quad (\mathbb{E}[\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}.) \\ &= \frac{1}{2} \text{tr} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\left(\mathbf{I} - \frac{\eta}{n} \mathbf{X} \mathbf{X}^\top \right)^{2k+2} \right] \\ &= \frac{1}{2} \text{tr}((1-\eta)^k (1+\delta) \mathbf{I}) \quad (\text{By Lemma D.4}) \\ &\leq d(1-\eta)^k \leq d^{-C \log(\frac{1}{1-\eta})}. \end{aligned}$$

\square

C PROOF OF THEOREM 4.1

C.1 GRADIENT COMPUTATION OF THE FULL MODEL OVER THE CoT OBJECTIVE

In this appendix, we compute the gradient of the full model given the Assumption 4.1 and prove the equivalence between the dynamics of the full model and a simplified model. Throughout the appendix, we denote the $\mathbf{S} = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$ for simplicity. And recall the i -th step of the linear classifier is $\mathbf{w}_i = (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \mathbf{w}^*$.

In Section 2.2, we have the full attention model

$$f_{\text{LSA}}(\mathbf{Z}; \mathbf{V}, \mathbf{W})_{[:, -1]} = \mathbf{Z}_{[:, -1]} + \mathbf{V} \mathbf{Z} \cdot \frac{\mathbf{Z}^\top \mathbf{W} \mathbf{Z}_{[:, -1]}}{n}$$

and the Chain of Thought (CoT) objective

$$\mathcal{L}^{\text{CoT}}(\mathbf{V}, \mathbf{W}) = \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\frac{1}{2} \sum_{i=0}^k \left\| f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right\|^2 \right]$$

We define the error for the i -th step

$$\mathcal{L}_i := \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right\|^2$$

By linearity of expectation, we know the gradient of the CoT objective is the sum of gradients of all CoT steps: $\nabla \mathcal{L}^{\text{CoT}} = \sum_{i=1}^k \nabla \mathcal{L}_i$. Now we can calculate the gradients of \mathbf{V}, \mathbf{W} based on the loss of each CoT step:

Lemma C.1 (Gradients of the full model). *The gradient of \mathbf{V}, \mathbf{W} are given by the following equations:*

$$\begin{aligned} \nabla_{\mathbf{V}} \mathcal{L} &= \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \sum_{i=0}^k \left(\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:, -1]}}{n} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1} - \mathbf{w}_i, 0)^\top \right) \mathbf{Z}_{i[:, -1]}^\top \mathbf{W}^\top \mathbf{Z}_i \mathbf{Z}_i^\top \\ \nabla_{\mathbf{W}} \mathcal{L} &= \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \sum_{i=0}^k \mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{V}^\top \left(\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:, -1]}}{n} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1} - \mathbf{w}_i, 0)^\top \right) \mathbf{Z}_{i[:, -1]}^\top \end{aligned}$$

Proof. The loss is given by eq. (6):

$$\mathcal{L}^{\text{CoT}}(\mathbf{V}, \mathbf{W}) = \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\frac{1}{2} \sum_{i=0}^k \left\| f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right\|^2 \right] = \sum_{i=1}^k \mathcal{L}_i$$

Take differential of the loss for the i -th step \mathcal{L}_i and we have

$$\begin{aligned} d\mathcal{L}_i &= \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left(f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right)^\top d f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} \\ &= \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left(f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right)^\top d \left(\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:, -1]}}{n} \right) \\ &= \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left(f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right)^\top d(\mathbf{V}) \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:, -1]}}{n} \\ &\quad + \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left(f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right)^\top \mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top d\mathbf{W} \mathbf{Z}_{i[:, -1]}}{n} \end{aligned}$$

Then the gradients of \mathbf{W}, \mathbf{V} of the \mathcal{L}_i are:

$$\nabla_{\mathbf{V}} \mathcal{L}_i = \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left(f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) \right) \mathbf{Z}_{i[:, -1]}^\top \mathbf{W}^\top \mathbf{Z}_i \mathbf{Z}_i^\top$$

$$\begin{aligned}
&= \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left(\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:,-1]}}{n} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1} - \mathbf{w}_i, 0)^\top \right) \mathbf{Z}_{i[:,-1]}^\top \mathbf{W}^\top \mathbf{Z}_i \mathbf{Z}_i^\top \\
\nabla_{\mathbf{W}} \mathcal{L}_i &= \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{V}^\top (f_{\text{LSA}}(\mathbf{Z}_i)_{[:,-1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1)) \mathbf{Z}_{i[:,-1]}^\top \\
&= \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{V}^\top \left(\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:,-1]}}{n} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1} - \mathbf{w}_i, 0)^\top \right) \mathbf{Z}_{i[:,-1]}^\top
\end{aligned}$$

Take the sum of the two equations above from $i = 0$ to k , and we finish the proof. \square

Now we consider the gradient flow (GF) trajectory (note that \mathbf{w}_{24} is fixed under Assumption 4.1):

$$\frac{d\boldsymbol{\theta}}{dt} = -\nabla \mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}), \quad \boldsymbol{\theta} := (\mathbf{V}, \mathbf{W} \setminus \{\mathbf{w}_{24}\}).$$

Recall the block matrix form of \mathbf{V}, \mathbf{W} :

$$\mathbf{V} = \begin{bmatrix} \mathbf{V}_{11} & \mathbf{V}_{12} & \mathbf{V}_{13} & \mathbf{V}_{14} \\ \mathbf{V}_{21} & v_{22} & \mathbf{V}_{23} & v_{24} \\ \mathbf{V}_{31} & \mathbf{V}_{32} & \mathbf{V}_{33} & \mathbf{V}_{34} \\ \mathbf{V}_{41} & v_{42} & \mathbf{V}_{43} & v_{44} \end{bmatrix}, \quad \mathbf{W} = \begin{bmatrix} \mathbf{W}_{11} & \mathbf{W}_{12} & \mathbf{W}_{13} & \mathbf{W}_{14} \\ \mathbf{W}_{21} & w_{22} & \mathbf{W}_{23} & w_{24} \\ \mathbf{W}_{31} & \mathbf{W}_{32} & \mathbf{W}_{33} & \mathbf{W}_{34} \\ \mathbf{W}_{41} & w_{42} & \mathbf{W}_{43} & w_{44} \end{bmatrix}$$

According to the construction in Theorem 3.2, the blocks $\mathbf{W}_{13}, \mathbf{V}_{31}, w_{24}$ are the only relevant parameter blocks, while the others should be zeroed out. Next, we prove that if we initialize those irrelevant blocks to 0, then they will stay at 0 along the gradient descent trajectory.

Lemma C.2. *Under the Assumption 4.1, when the linear transformer is trained under GF, we have for all $t > 0$, the parameters $\mathbf{V}(t), \mathbf{W}(t)$ have the following form:*

$$\mathbf{V}(t) = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 0 & \mathbf{0} & 0 \\ \mathbf{V}_{31}(t) & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 0 & \mathbf{0} & 0 \end{bmatrix}, \quad \mathbf{W}(t) = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{W}_{13}(t) & \mathbf{0} \\ \mathbf{0} & 0 & \mathbf{0} & -1 \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 0 & \mathbf{0} & 0 \end{bmatrix}$$

Proof. To prove this lemma, we prove that when the irrelevant blocks are 0, the gradients $\nabla_{\mathbf{V}} \mathcal{L}_i, \nabla_{\mathbf{W}} \mathcal{L}_i$ for those blocks are always 0 and they never update the corresponding parameter block. Also, note that $w_{24} = -1$ for all $t > 0$.

First, we calculate the output of the linear self-attention $\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:,-1]}}{n}$:

$$\begin{aligned}
&\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:,-1]}}{n} \\
&= \frac{1}{n} \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 0 & \mathbf{0} & 0 \\ \mathbf{V}_{31}(t) & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 0 & \mathbf{0} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{X} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{y} & 0 & 0 & \cdots & 0 \\ \mathbf{0}_{d \times n} & \mathbf{w}_0 & \mathbf{w}_1 & \cdots & \mathbf{w}_i \\ \mathbf{0}_{1 \times n} & 1 & 1 & \cdots & 1 \end{bmatrix} \mathbf{Z}_i^\top \mathbf{W} \begin{bmatrix} \mathbf{0} \\ 0 \\ \mathbf{w}_i \\ 1 \end{bmatrix} \\
&= \frac{1}{n} \begin{bmatrix} \mathbf{0}_{d \times n} & \mathbf{0}_d & \cdots & \mathbf{0}_d \\ \mathbf{0}_{1 \times n} & 0 & \cdots & 0 \\ \mathbf{V}_{31}(t) \mathbf{X} & \mathbf{0}_d & \cdots & \mathbf{0}_d \\ \mathbf{0}_{1 \times n} & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{X} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{y} & 0 & 0 & \cdots & 0 \\ \mathbf{0}_{d \times n} & \mathbf{w}_0 & \mathbf{w}_1 & \cdots & \mathbf{w}_i \\ \mathbf{0}_{1 \times n} & 1 & 1 & \cdots & 1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{W}_{13}(t) \mathbf{w}_i \\ -1 \\ \mathbf{0}_d \\ 0 \end{bmatrix} \\
&= \frac{1}{n} \begin{bmatrix} \mathbf{0}_{d \times n} & \mathbf{0}_d & \cdots & \mathbf{0}_d \\ \mathbf{0}_{1 \times n} & 0 & \cdots & 0 \\ \mathbf{V}_{31}(t) \mathbf{X} & \mathbf{0}_d & \cdots & \mathbf{0}_d \\ \mathbf{0}_{1 \times n} & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{X}^\top \mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{y}^\top \\ \mathbf{0}_{i+1} \end{bmatrix} = \frac{1}{n} \begin{bmatrix} \mathbf{0}_d \\ 0 \\ \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) \\ 0 \end{bmatrix}
\end{aligned}$$

The last line is because $\mathbf{y}^\top = \mathbf{X}^\top \mathbf{w}^*$. Now, we consider the gradient for \mathbf{V} :

$$\nabla_{\mathbf{V}} \mathcal{L}_i = \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\left(\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_{i[:,-1]}}{n} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1} - \mathbf{w}_i, 0)^\top \right) \mathbf{Z}_{i[:,-1]}^\top \mathbf{W}^\top \mathbf{Z}_i \mathbf{Z}_i^\top \right]$$

$$\begin{aligned}
&= \frac{1}{n^2} \begin{bmatrix} \mathbf{0}_d \\ 0 \\ \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i) \\ 0 \end{bmatrix} \mathbf{Z}_i^\top \mathbf{W}^\top \mathbf{Z}_i^\top \\
&= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \begin{bmatrix} \mathbf{0}_d \\ 0 \\ \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i) \\ 0 \end{bmatrix} \begin{bmatrix} \mathbf{w}_i^\top \mathbf{W}_{13}^\top(t) \\ -1 \\ \mathbf{0}_d \\ 0 \end{bmatrix}^\top \mathbf{Z}_i \mathbf{Z}_i^\top \\
&= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \begin{bmatrix} \mathbf{0}_d \\ 0 \\ \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i) \\ 0 \end{bmatrix} \begin{bmatrix} \mathbf{w}_i^\top \mathbf{W}_{13}^\top(t) \mathbf{X} \mathbf{X}^\top - \mathbf{y} \mathbf{X}^\top \\ \mathbf{w}_i^\top \mathbf{W}_{13}^\top(t) \mathbf{X} \mathbf{y}^\top - \mathbf{y} \mathbf{y}^\top \\ \mathbf{0}_d \\ 0 \end{bmatrix}^\top \\
&= \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \nabla_{\mathbf{V}_{31}} \mathcal{L}_i(t) & \nabla_{\mathbf{V}_{32}} \mathcal{L}_i(t) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}
\end{aligned}$$

Therefore, we know all blocks of the gradient are zero except the positions of \mathbf{V}_{31} and \mathbf{V}_{32} .

Now look at $\nabla_{\mathbf{V}_{32}} \mathcal{L}_i$:

$$\begin{aligned}
\nabla_{\mathbf{V}_{32}} \mathcal{L}_i &= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} [(\mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i)) \\
&\quad (\mathbf{w}_i^\top \mathbf{W}_{13}^\top(t) \mathbf{X} \mathbf{y}^\top - \mathbf{y} \mathbf{y}^\top)] \\
&= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} [(\mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i)) \\
&\quad (\mathbf{w}_i^\top \mathbf{W}_{13}^\top(t) \mathbf{X} \mathbf{X}^\top \mathbf{w}^* - \mathbf{w}^{*\top} \mathbf{X} \mathbf{X}^\top \mathbf{w}^*)]
\end{aligned}$$

Note that $\mathbf{w}_i = (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \mathbf{w}^*$ for all $i \in [k]$, and $\mathbf{w}_{k+1} = \mathbf{w}^*$. Therefore, for all $i \in \{0, 1, \dots, k+1\}$ the formula inside the expectation is an odd function of \mathbf{w}^* . Since $\mathbf{w}^* \sim \mathcal{N}(0, \mathbf{I}_d)$, the expectation should be $\mathbf{0}_d$.

Similarly, we calculate the gradient of the \mathbf{W} :

$$\begin{aligned}
\nabla_{\mathbf{W}} \mathcal{L}_i &= \frac{1}{n} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{V}^\top \left(\mathbf{V} \mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W} \mathbf{Z}_i^\top}{n} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1} - \mathbf{w}_i, 0)^\top \right) \mathbf{Z}_i^\top \right] \\
&= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{V}^\top \begin{bmatrix} \mathbf{0}_d \\ 0 \\ \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i) \\ 0 \end{bmatrix} \mathbf{Z}_i^\top \right] \\
&= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \begin{bmatrix} \mathbf{0}_{d \times d} & 0 & \mathbf{X} \mathbf{X}^\top \mathbf{V}_{31}(t)^\top & 0 \\ \mathbf{0}_{d \times d} & 0 & \mathbf{y} \mathbf{X}^\top \mathbf{V}_{31}(t)^\top & 0 \\ \mathbf{0}_{d \times d} & 0 & \mathbf{0}_{d \times d} & 0 \\ \mathbf{0}_{d \times d} & 0 & \mathbf{0}_{d \times d} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{0}_d \\ 0 \\ \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i) \\ 0 \end{bmatrix} \mathbf{Z}_i^\top \right] \\
&= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\begin{bmatrix} \mathbf{X} \mathbf{X}^\top \mathbf{V}_{31}(t)^\top \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i) \\ \mathbf{y} \mathbf{X}^\top \mathbf{V}_{31}(t)^\top \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i) \\ \mathbf{0}_d \\ 0 \end{bmatrix} \begin{bmatrix} \mathbf{0}_d \\ 0 \\ \mathbf{w}_i \\ 1 \end{bmatrix}^\top \right] \\
&= \begin{bmatrix} \mathbf{0} & \mathbf{0} & \nabla_{\mathbf{W}_{13}} \mathcal{L}_i(t) & \nabla_{\mathbf{W}_{14}} \mathcal{L}_i(t) \\ \mathbf{0} & \mathbf{0} & \nabla_{\mathbf{W}_{23}} \mathcal{L}_i(t) & \nabla_{\mathbf{W}_{24}} \mathcal{L}_i(t) \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}
\end{aligned}$$

Since we fix \mathbf{w}_{24} , we only consider the remaining three blocks. First, we consider the gradient of the vector block \mathbf{W}_{14} :

$$\nabla_{\mathbf{W}_{14}} \mathcal{L}_i(t) = \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} [\mathbf{X} \mathbf{X}^\top \mathbf{V}_{31}(t)^\top \mathbf{V}_{31}(t) \mathbf{X} \mathbf{X}^\top (\mathbf{W}_{13}(t) \mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i)].$$

Notice that the $\mathbf{X}\mathbf{X}^\top \mathbf{V}_{31}(t)^\top \mathbf{V}_{31}(t) \mathbf{X}\mathbf{X}^\top (\mathbf{W}_{13}(t)\mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i)$ is odd in \mathbf{w}^* . Therefore the expectation is $\mathbf{0}_d$. Similarly, we consider the other block \mathbf{W}_{23} :

$$\begin{aligned}\nabla_{\mathbf{W}_{23}} \mathcal{L}_i(t) &= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[(\mathbf{y}\mathbf{X}^\top \mathbf{V}_{31}(t)^\top \mathbf{V}_{31}(t) \mathbf{X}\mathbf{X}^\top (\mathbf{W}_{13}(t)\mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i)) \mathbf{w}_i^\top \right] \\ &= \frac{1}{n^2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[(\mathbf{w}^{*\top} \mathbf{X}\mathbf{X}^\top \mathbf{V}_{31}(t)^\top \mathbf{V}_{31}(t) \mathbf{X}\mathbf{X}^\top (\mathbf{W}_{13}(t)\mathbf{w}_i - \mathbf{w}^*) - n(\mathbf{w}_{i+1} - \mathbf{w}_i)) \mathbf{w}_i^\top \right] \\ &= \mathbf{0}_{1 \times d}.\end{aligned}$$

In conclusion, all the blocks have zero gradient except \mathbf{V}_{31} , \mathbf{W}_{13} given that they are all zero matrices. Under Assumption 4.1, all the irrelevant blocks remain zero matrices for all $t \geq 0$. \square

By Lemma C.2, we prove that along the gradient flow trajectory under Assumption 4.1, the objective of the linear self-attention model with residual connection can be equivalently transform to the following simplified form.

Lemma C.3. *Under Assumption 4.1, we have the training objective*

$$\mathcal{L}^{\text{CoT}}(\mathbf{V}, \mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\sum_{i=0}^k \|\mathbf{V}_{31}(\mathbf{S}\mathbf{W}_{13}\mathbf{w}_i - \mathbf{S}\mathbf{w}^*) - \Delta\mathbf{w}_i\|^2 \right]$$

where $\mathbf{S} = \frac{1}{n} \mathbf{X}\mathbf{X}^\top$ and $\Delta\mathbf{w}_i := \mathbf{w}_{i+1} - \mathbf{w}_i$, $i = 0, 1, \dots, k$ is the residual for each step i .

Proof. Given the following CoT objective,

$$\mathcal{L}^{\text{CoT}}(\mathbf{V}, \mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\sum_{i=0}^k \|f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1)\|^2 \right]$$

By Lemma C.2, we plug in the \mathbf{V} , \mathbf{W} expressions and get:

$$\begin{aligned}f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) &= \mathbf{V}\mathbf{Z}_i \cdot \frac{\mathbf{Z}_i^\top \mathbf{W}\mathbf{Z}_i_{[:, -1]}}{n} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1} - \mathbf{w}_i, 0)^\top \\ &= \left(\mathbf{0}_d, 0, \frac{1}{n} \mathbf{V}_{31}(\mathbf{X}\mathbf{X}^\top \mathbf{W}_{13}\mathbf{w}_i - \mathbf{X}\mathbf{y}^\top) - \Delta\mathbf{w}_i, 0 \right)\end{aligned}$$

Since $\mathbf{y}^\top = \mathbf{X}^\top \mathbf{w}^*$, we have

$$f_{\text{LSA}}(\mathbf{Z}_i)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}_{i+1}, 1) = (\mathbf{0}_d, 0, \mathbf{V}_{31}(\mathbf{S}^\top \mathbf{W}_{13}\mathbf{w}_i - \mathbf{S}\mathbf{w}^*) - \Delta\mathbf{w}_i, 0)^\top$$

Put it back to the loss expression and we complete the proof. \square

Now the chain of thought loss can be rewritten into the form by Lemma C.3, we can directly calculate the gradient update using the simplified loss for clarity. We denote the only relevant blocks $\tilde{\mathbf{W}} := \mathbf{W}_{13}$ and $\tilde{\mathbf{V}} := \mathbf{V}_{31}$. Moreover, we can further expand the CoT loss with $\Delta\mathbf{w}_i = -\eta \cdot \frac{\mathbf{X}\mathbf{X}^\top}{n}(\mathbf{w}_i - \mathbf{w}^*)$ for $i \in \{0, 1, \dots, k-1\}$, and $\Delta\mathbf{w}_k = \mathbf{w}^* - \mathbf{w}_k$. That leads to the following expression of the CoT loss:

$$\begin{aligned}\mathcal{L}^{\text{CoT}}(\theta) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \sum_{i=0}^{k-1} \left\| \mathbf{w}_i + \tilde{\mathbf{V}}\mathbf{S}(\tilde{\mathbf{W}}\mathbf{w}_i - \mathbf{w}^*) - \mathbf{w}_{i+1} \right\|_2^2 \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| \mathbf{w}_k + \tilde{\mathbf{V}}\mathbf{S}(\tilde{\mathbf{W}}\mathbf{w}_k - \mathbf{w}^*) - \mathbf{w}^* \right\|_2^2\end{aligned}\tag{13}$$

Observe that the final loss only depends on the $(d+2)$ to $(2d+2)$ entries of the transformer's output, indicating we can simplify the model a bit and prune out the irrelevant part. We can define a simplified one-layer transformer to get the loss form above, where the dynamics of the equivalent model is exactly the same with the original dynamics of \mathbf{W}_{13} and \mathbf{V}_{31} . Accordingly, the last token input of the transformer for i -th step becomes \mathbf{w}_i and the label becomes \mathbf{w}_{i+1} since the other entries in the original input/label $(\mathbf{0}, 0, \mathbf{w}_i, 1)$ do not affect prediction.

Definition C.1 (Reduced transformer). Let $\theta = (\tilde{V}, \tilde{W})$. Define

$$f_\theta(\mathbf{X}, \mathbf{Z}_i) = \mathbf{w}_i + \tilde{V} S(\tilde{W} \mathbf{w}_i - \mathbf{w}^*)$$

to be the reduced model of the one-layer transformer in Equation (3). For ease of presentation, we denote $f_\theta(\mathbf{w}_i) := f_\theta(\mathbf{X}, \mathbf{Z}_i)$.

In the following sections, we will consider the equivalent form of transformer. Here we present the gradient with regard to the reduced model. For clarification, throughout this section we will denote $\mathbf{w}_{k+1} := (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k+1}) \mathbf{w}^*$ as the $(k+1)$ -th update, and \mathbf{w}^* is the ground-truth.

Lemma C.4. The gradient of \tilde{V} and \tilde{W} are given by the following expectations:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{V}} &= \sum_{i=0}^k \mathbb{E} \left[(f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) (\mathbf{w}_i^\top \tilde{W}^\top - \mathbf{w}^{*T}) \mathbf{S} \right] + \mathbb{E} \left[(\mathbf{w}_{k+1} - \mathbf{w}^*) (\mathbf{w}_k^\top \tilde{W}^\top - \mathbf{w}^{*T}) \mathbf{S} \right], \\ \frac{\partial \mathcal{L}}{\partial \tilde{W}} &= \sum_{i=0}^k \mathbb{E} \left[\tilde{V}^\top (f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) \mathbf{w}_i^\top \right] + \mathbb{E} \left[\tilde{V}^\top (\mathbf{w}_{k+1} - \mathbf{w}^*) \mathbf{w}_k^\top \right]. \end{aligned}$$

Proof. Given the equivalent CoT loss in Equation (13), we take the gradient with regard to \tilde{V} of the loss and we have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{V}} &= \sum_{i=0}^{k-1} \mathbb{E} \left[(f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) (\mathbf{w}_i^\top \tilde{W}^\top - \mathbf{w}^{*T}) \mathbf{S} \right] + \mathbb{E} \left[(f_\theta(\mathbf{w}_k) - \mathbf{w}^*) (\mathbf{w}_k^\top \tilde{W}^\top - \mathbf{w}^{*T}) \mathbf{S} \right] \\ &= \sum_{i=0}^k \mathbb{E} \left[(f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) (\mathbf{w}_i^\top \tilde{W}^\top - \mathbf{w}^{*T}) \mathbf{S} \right] + \mathbb{E} \left[(\mathbf{w}_{k+1} - \mathbf{w}^*) (\mathbf{w}_k^\top \tilde{W}^\top - \mathbf{w}^{*T}) \mathbf{S} \right] \end{aligned}$$

The second step is because we subtract $\mathbb{E} \left[(f_\theta(\mathbf{w}_k) - \mathbf{w}_{k+1}) (\mathbf{w}_k^\top \tilde{W}^\top - \mathbf{w}^{*T}) \mathbf{S} \right]$ from the second term and put it into the summation. Similarly, the partial derivative of \tilde{W} should be:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{W}} &= \sum_{i=0}^{k-1} \mathbb{E} \left[\tilde{V}^\top (f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) \mathbf{w}_i^\top \right] + \mathbb{E} \left[\tilde{V}^\top (f_\theta(\mathbf{w}_k) - \mathbf{w}^*) \mathbf{w}_k^\top \right] \\ &= \sum_{i=0}^k \mathbb{E} \left[\tilde{V}^\top (f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) \mathbf{w}_i^\top \right] + \mathbb{E} \left[\tilde{V}^\top (\mathbf{w}_{k+1} - \mathbf{w}^*) \mathbf{w}_k^\top \right] \end{aligned}$$

Therefore we complete the proof. \square

C.2 GRADIENT CHARACTERIZATION OVER THE CoT OBJECTIVE

In this section, we compute the exact gradient for the reduced model parameters to facilitate analysis on the dynamics. For clarification, throughout this section we will denote $\mathbf{w}_{k+1} := (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k+1}) \mathbf{w}^*$ as the $(k+1)$ -th update, and \mathbf{w}^* is the ground-truth.

We first compute our gradients for the simplified model defined in Definition C.1, which is equivalent to the full model's dynamics. Recall that under assumption 4.1, we have \tilde{V}, \tilde{W} are simultaneously diagonalizable, with the orthonormal basis $\{\mathbf{u}_i\}_{i=1}^d$. We denote the orthogonal matrix formed by the basis as \mathbf{U} . We will observe that \mathbf{u}_i are always the eigenvector of \tilde{V}, \tilde{W} , so we denote $\tilde{V} = \mathbf{U} \Lambda^{\tilde{V}} \mathbf{U}^\top, \tilde{W} = \mathbf{U} \Lambda^{\tilde{W}} \mathbf{U}^\top$. For clarity, we ignore the timestamp when calculating the gradients and dynamics.

We present an accurate estimate of the gradient in the following Lemma C.5. We intensively use the concentration lemma in Appendix D to separate the main terms dominating the gradient flow dynamics, and some bounded error terms that may complicate the analysis. We also call the error terms as ‘interaction terms’, since they contain the interactions between two subspaces $\mathbf{u}_i \mathbf{u}_i^\top$ and $\mathbf{u}_j \mathbf{u}_j^\top$. The structure of the interaction terms $\Delta^{\tilde{V}}, \Delta^{\tilde{W}}$ are further characterized in this lemma, which is essential for the final local convergence analysis.

Lemma C.5. Suppose $n = \Theta(d \log^5 d)$, $\eta \in (0.1, 0.9)$, $k = \lceil c \log d \rceil$. Under Assumption 4.1, if we run gradient flow on the population loss in Equation (6), then the gradient of $\tilde{\mathbf{V}}$ and $\tilde{\mathbf{W}}$ are characterized by the following equations:

$$\begin{aligned} \mathbf{U}^\top \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{V}}} \mathbf{U} &= \left[\left(k+1 - \frac{2}{\eta} + \frac{1}{\eta(2-\eta)} \right) \mathbf{\Lambda} \tilde{\mathbf{W}}^2 - 2 \left(k+1 - \frac{1}{\eta} \right) \mathbf{\Lambda} \tilde{\mathbf{W}} + (k+1) \mathbf{I} \right] \mathbf{\Lambda} \tilde{\mathbf{V}} \\ &\quad - \frac{1-\eta}{2-\eta} \mathbf{\Lambda} \tilde{\mathbf{W}} + \mathbf{I} + \Delta \tilde{\mathbf{V}}, \\ \mathbf{U}^\top \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{W}}} \mathbf{U} &= \left(k+1 - \frac{2}{\eta} + \frac{1}{\eta(2-\eta)} \right) \mathbf{\Lambda} \tilde{\mathbf{V}}^2 \mathbf{\Lambda} \tilde{\mathbf{W}} - \left(k+1 - \frac{1}{\eta} \right) \mathbf{\Lambda} \tilde{\mathbf{V}}^2 - \frac{1-\eta}{2-\eta} \mathbf{\Lambda} \tilde{\mathbf{V}} + \Delta \tilde{\mathbf{W}}. \end{aligned}$$

where the error terms (interaction terms) $\|\Delta \tilde{\mathbf{V}}\|_{op} \leq O\left(\frac{1}{\log^2 d}\right)$, $\|\Delta \tilde{\mathbf{W}}\|_{op} \leq O\left(\frac{1}{\log^2 d}\right)$. Moreover, there exist diagonal matrices $\mathbf{A}^{\tilde{\mathbf{V}}}, \mathbf{B}^{\tilde{\mathbf{V}}}, \mathbf{A}^{\tilde{\mathbf{W}}}, \mathbf{B}^{\tilde{\mathbf{W}}}$ with $O\left(\frac{1}{\log^2 d}\right)$ -operator norm, $\mathbf{C}^{\tilde{\mathbf{V}}}, \mathbf{D}^{\tilde{\mathbf{V}}}, \mathbf{C}^{\tilde{\mathbf{W}}}, \mathbf{D}^{\tilde{\mathbf{W}}}, \mathbf{E}^{\tilde{\mathbf{W}}}$ with $O\left(\frac{1}{d \log^2 d}\right)$ -operator norm and $\mathbf{E}^{\tilde{\mathbf{V}}}, \mathbf{F}^{\tilde{\mathbf{W}}}$ with $O\left((1-\eta)^k\right)$ -operator norm s.t. the error terms $\Delta \tilde{\mathbf{V}}, \Delta \tilde{\mathbf{W}}$ can be written as

$$\begin{aligned} \Delta \tilde{\mathbf{V}} &= (\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{A}^{\tilde{\mathbf{V}}} + (\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{B}^{\tilde{\mathbf{V}}} + \text{tr}((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{C}^{\tilde{\mathbf{V}}} + \text{tr}(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{D}^{\tilde{\mathbf{V}}} + \mathbf{E}^{\tilde{\mathbf{V}}}, \\ \Delta \tilde{\mathbf{W}} &= (\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{A}^{\tilde{\mathbf{W}}} + (\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{B}^{\tilde{\mathbf{W}}} + \text{tr}(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{C}^{\tilde{\mathbf{W}}} + \text{tr}((\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{\Lambda} \tilde{\mathbf{V}}) \mathbf{D}^{\tilde{\mathbf{W}}} \\ &\quad + \text{tr}((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{V}}^2) \mathbf{E}^{\tilde{\mathbf{W}}} + \mathbf{F}^{\tilde{\mathbf{W}}}. \end{aligned}$$

Proof. Recall the gradients formula of $\tilde{\mathbf{V}}$ and $\tilde{\mathbf{W}}$ by Lemma C.4:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{V}}} &= \sum_{i=0}^k \mathbb{E} \left[(f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) (\mathbf{w}_i^\top \tilde{\mathbf{W}}^\top - \mathbf{w}^{*T}) \mathbf{S} \right] + \mathbb{E} \left[(\mathbf{w}_{k+1} - \mathbf{w}^*) (\mathbf{w}_k^\top \tilde{\mathbf{W}}^\top - \mathbf{w}^{*T}) \mathbf{S} \right] \\ \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{W}}} &= \sum_{i=0}^k \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}) \mathbf{w}_i^\top \right] + \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (\mathbf{w}_{k+1} - \mathbf{w}^*) \mathbf{w}_k^\top \right] \end{aligned}$$

We expand the reduced model $f_\theta(\mathbf{w}_i)$ in Definition C.1, and get the residual term

$$\begin{aligned} f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1} &= \mathbf{w}_i + \tilde{\mathbf{V}} \mathbf{S} (\tilde{\mathbf{W}} \mathbf{w}_i - \mathbf{w}^*) - \mathbf{w}_{i+1} \\ &= \tilde{\mathbf{V}} \mathbf{S} (\tilde{\mathbf{W}} \mathbf{w}_i - \mathbf{w}^*) + \eta \mathbf{S} (\mathbf{w}_i - \mathbf{w}^*) \\ &= (\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) \mathbf{w}_i - (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} \mathbf{w}^* \end{aligned}$$

Substitute $f_\theta(\mathbf{w}_i) - \mathbf{w}_{i+1}$ term in the dynamics by the equation above, we have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{V}}} &= \sum_{i=0}^k \mathbb{E} \left[(\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i)^2 \tilde{\mathbf{W}}^\top \mathbf{S} \right] - \sum_{i=0}^k \mathbb{E} \left[(\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \mathbf{S} \right] \\ &\quad - \sum_{i=0}^k \mathbb{E} \left[(\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \tilde{\mathbf{W}}^\top \mathbf{S} \right] + \sum_{i=0}^k \mathbb{E} \left[(\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S}^2 \right] \\ &\quad - \mathbb{E} \left[(\mathbf{I} - \eta \mathbf{S})^{k+1} ((\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k) \tilde{\mathbf{W}}^\top - \mathbf{I}) \right] \\ &= \sum_{i=0}^k (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{W}} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i)^2 \tilde{\mathbf{W}}^\top \mathbf{S} \right] \tag{Term 1} \\ &\quad + \eta \sum_{i=0}^k \mathbb{E} \left[\mathbf{S} (\mathbf{I} - \tilde{\mathbf{W}}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i)^2 \tilde{\mathbf{W}}^\top \mathbf{S} \right] \tag{Term 2} \end{aligned}$$

$$- \sum_{i=0}^k (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{W}} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) \mathbf{S} \right] \quad (\text{Term 3})$$

$$- \eta \sum_{i=0}^k \mathbb{E} \left[\mathbf{S} \left(\mathbf{I} - \tilde{\mathbf{W}} \right) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) \mathbf{S} \right] \quad (\text{Term 4})$$

$$- \sum_{i=0}^k (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbb{E} \left[\mathbf{S} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) \tilde{\mathbf{W}}^\top \mathbf{S} \right] \quad (\text{Term 5})$$

$$+ \sum_{i=0}^k (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbb{E} [\mathbf{S}^2] \quad (\text{Term 6})$$

$$- \mathbb{E} \left[(\mathbf{I} - \eta \mathbf{S})^{k+1} \left(\left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \tilde{\mathbf{W}}^\top - \mathbf{I} \right) \right]. \quad (\text{Term 7})$$

To get an accurate estimate of the gradient, we apply Lemma C.14, Lemma C.15 respectively to each of the terms (Term 1 to Term 7) and separate the interaction terms introduced by the moments of Wishart matrix, which is bounded by $O\left(\frac{1}{\log^3 d}\right)$.

Consider Term 7 and the i -th term in the summation of Term 1 to Term 6. By Lemma C.14 and Lemma C.15, there exist diagonal matrices $\xi_j, j \in [6]$ satisfying $\|\xi_j\|_{op} \leq O\left(\frac{1}{\log^3 d}\right)$ such that

$$\begin{aligned} \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{W}} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 \tilde{\mathbf{W}}^\top \mathbf{S} \right] &= \mathbf{U} \left[\left(1 - (1 - \eta)^k \right)^2 \mathbf{\Lambda} \tilde{\mathbf{W}}^2 + \xi_1 \right] \mathbf{U}^\top \\ \mathbb{E} \left[\mathbf{S} \left(\mathbf{I} - \tilde{\mathbf{W}} \right) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 \tilde{\mathbf{W}}^\top \mathbf{S} \right] &= \mathbf{U} \left[\left(1 - (1 - \eta)^k \right)^2 \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{\Lambda} \tilde{\mathbf{W}} + \xi_2 \right] \mathbf{U}^\top \\ \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{W}} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) \mathbf{S} \right] &= \mathbf{U} \left[\left(1 - (1 - \eta)^k \right) \mathbf{\Lambda} \tilde{\mathbf{W}} + \xi_3 \right] \mathbf{U}^\top \\ \mathbb{E} \left[\mathbf{S} \left(\mathbf{I} - \tilde{\mathbf{W}} \right) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) \mathbf{S} \right] &= \mathbf{U} \left[\left(1 - (1 - \eta)^k \right) \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) + \xi_4 \right] \mathbf{U}^\top \\ \mathbb{E} \left[\mathbf{S} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) \tilde{\mathbf{W}}^\top \mathbf{S} \right] &= \mathbf{U} \left[\left(1 - (1 - \eta)^k \right) \mathbf{\Lambda} \tilde{\mathbf{W}} + \xi_5 \right] \mathbf{U}^\top \\ \mathbb{E} [\mathbf{S}^2] &= \mathbf{U} (\mathbf{I} + \xi_6) \mathbf{U}^\top \end{aligned}$$

By Lemma D.4, there exists diagonal matrix ξ_7 satisfying $\|\xi_7\|_{op} \leq O\left((1 - \eta)^k\right)$ such that

$$\mathbb{E} \left[(\mathbf{I} - \eta \mathbf{S})^{k+1} \left(\left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \tilde{\mathbf{W}}^\top - \mathbf{I} \right) \right] = \mathbf{U} \xi_7 \mathbf{U}^\top.$$

Moreover, there exist $\alpha_1, \alpha_2 \leq O\left(\frac{1}{\log^3 d}\right), \alpha_3, \alpha_4, \alpha_5 \leq O\left(\frac{1}{d \log^3 d}\right)$ such that

$$\xi_2 = \left(\alpha_1 \mathbf{\Lambda} \tilde{\mathbf{W}} + \alpha_2 \mathbf{I} \right) \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) + \text{tr} \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \left(\alpha_3 \mathbf{\Lambda} \tilde{\mathbf{W}} + \alpha_4 \mathbf{I} \right) + \alpha_5 \text{tr} \left(\left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{I},$$

and exist $\beta_1 \leq O\left(\frac{1}{\log^3 d}\right), \beta_2 \leq O\left(\frac{1}{d \log^3 d}\right)$ such that

$$\xi_4 = \beta_1 \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) + \beta_2 \text{tr} \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{I}.$$

We define $\Delta_i^{\tilde{\mathbf{V}}}$ as the sum of all the interaction terms $(\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I})(\xi_1 - \xi_3 - \xi_5 + \xi_6) + \eta(\xi_2 - \xi_4)$ for the i -th term in the summation of dynamics of $\tilde{\mathbf{V}}$. From the analysis above, there exist diagonal matrices $\mathbf{A}_i^{\tilde{\mathbf{V}}}, \mathbf{B}_i^{\tilde{\mathbf{V}}}, \mathbf{C}_i^{\tilde{\mathbf{V}}}, \mathbf{D}_i^{\tilde{\mathbf{V}}}$ with their operator norm $O\left(\frac{1}{\log^3 d}\right)$, such that (note every matrix is diagonal, so they commute)

$$\begin{aligned} \Delta_i^{\tilde{\mathbf{V}}} &= \left(\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I} \right) \mathbf{A}_i^{\tilde{\mathbf{V}}} + O\left(\frac{1}{d}\right) \text{tr} \left(\left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{B}_i^{\tilde{\mathbf{V}}} \\ &\quad + \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{C}_i^{\tilde{\mathbf{V}}} + O\left(\frac{1}{d}\right) \text{tr} \left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right) \mathbf{D}_i^{\tilde{\mathbf{V}}}. \end{aligned}$$

We define $\Delta_{-1}^{\tilde{V}}$ as the interaction term brought by Term 7 since there is no summation in Term 7. It is obvious that $\|\Delta_{-1}^{\tilde{V}}\|_{op} \leq O((1-\eta)^k)$.

Now we denote

$$\hat{\Delta}^{\tilde{V}} = \sum_{i=0}^k \Delta_i^{\tilde{V}} - \Delta_{-1}^{\tilde{V}}$$

to be the sum of all interaction term of the dynamics of $\Lambda^{\tilde{V}}$. From the definition of $\Delta_i^{\tilde{V}}$ and $\Delta_{-1}^{\tilde{V}}$ above, there exist diagonal matrices $A^{\tilde{V}}, B^{\tilde{V}}, C^{\tilde{V}}, D^{\tilde{V}}$ and $E_0^{\tilde{V}}$ satisfying $\|A^{\tilde{V}}\|, \|C^{\tilde{V}}\| \leq O\left(\frac{1}{\log^2 d}\right)$, $\|B^{\tilde{V}}\|, \|D^{\tilde{V}}\| \leq O\left(\frac{1}{d \log^2 d}\right)$ and $\|E_0^{\tilde{V}}\| \leq O((1-\eta)^k)$ such that (because $k = \Theta(\log d)$)

$$\hat{\Delta}^{\tilde{V}} = (\Lambda^{\tilde{V}} + \eta I)A^{\tilde{V}} + \text{tr}\left((I - \Lambda^{\tilde{W}})\Lambda^{\tilde{W}}\right)B^{\tilde{V}} + (I - \Lambda^{\tilde{W}})C^{\tilde{V}} + \text{tr}\left((I - \Lambda^{\tilde{W}})D^{\tilde{V}}\right) + E_0^{\tilde{V}}$$

Sum up all the seven terms together and we have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{V}} = & U \left[\left(k+1 - \frac{2(1-(1-\eta)^{k+1})}{\eta} + \frac{1-(1-\eta)^{2k+2}}{\eta(2-\eta)} \right) \Lambda^{\tilde{W}} (\Lambda^{\tilde{V}} \Lambda^{\tilde{W}} + \eta I) \right] U^\top \\ & - U \left[\left(k+1 - \frac{1-(1-\eta)^{k+1}}{\eta} \right) (\Lambda^{\tilde{V}} \Lambda^{\tilde{W}} + \eta I) \right] U^\top \\ & - U \left[\left(k+1 - \frac{1-(1-\eta)^{k+1}}{\eta} \right) \Lambda^{\tilde{W}} (\Lambda^{\tilde{V}} + \eta I) \right] U^\top \\ & + U \left[(k+1) (\Lambda^{\tilde{V}} + \eta I) \right] U^\top + U \hat{\Delta}^{\tilde{V}} U^\top \end{aligned}$$

Denote $E_1^{\tilde{V}}$ to be the sum of all $O((1-\eta)^k)$ terms in the dynamics of \tilde{V} :

$$E_1^{\tilde{V}} = \left(\frac{2(1-\eta)^{k+1}}{\eta} - \frac{(1-\eta)^{2k+2}}{\eta(2-\eta)} \right) \Lambda^{\tilde{W}} (\Lambda^{\tilde{V}} \Lambda^{\tilde{W}} + \eta I) - \frac{(1-\eta)^{k+1}}{\eta} (2\Lambda^{\tilde{V}} \Lambda^{\tilde{W}} + \eta \Lambda^{\tilde{W}} + \eta I)$$

Denote $E^{\tilde{V}} = E_0^{\tilde{V}} + E_1^{\tilde{V}}$ and denote $\Delta^{\tilde{V}} = \hat{\Delta}^{\tilde{V}} + E_1^{\tilde{V}}$, we have

$$\begin{aligned} U^\top \frac{\partial \mathcal{L}}{\partial \tilde{V}} U = & \left[\left(k+1 - \frac{2}{\eta} + \frac{1}{\eta(2-\eta)} \right) \Lambda^{\tilde{W}^2} - 2 \left(k+1 - \frac{1}{\eta} \right) \Lambda^{\tilde{W}} + (k+1) I \right] \Lambda^{\tilde{V}} \\ & - \frac{1-\eta}{2-\eta} \Lambda^{\tilde{W}} + I + \Delta^{\tilde{V}} \end{aligned}$$

Moreover, $\Delta^{\tilde{V}}$ has the form

$$\Delta^{\tilde{V}} = (\Lambda^{\tilde{V}} + \eta I)A^{\tilde{V}} + \text{tr}\left((I - \Lambda^{\tilde{W}})\Lambda^{\tilde{W}}\right)B^{\tilde{V}} + (I - \Lambda^{\tilde{W}})C^{\tilde{V}} + \text{tr}\left((I - \Lambda^{\tilde{W}})D^{\tilde{V}}\right) + E^{\tilde{V}}$$

Similar to the calculation of the dynamics of \tilde{V} , we can also have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{W}} = & \sum_{i=0}^k \mathbb{E} \left[S \tilde{V}^\top (f_\theta(w_i) - w_{i+1}) w_i^\top \right] + \mathbb{E} \left[S \tilde{V}^\top (w_{k+1} - w^*) w_k^\top \right] \\ = & \sum_{i=0}^k \mathbb{E} \left[S \tilde{V}^\top (\tilde{V} S \tilde{W} + \eta S) (I - (I - \eta S)^i)^2 \right] - \sum_{i=0}^k \mathbb{E} \left[S \tilde{V}^\top (\tilde{V} + \eta I) S (I - (I - \eta S)^i) \right] \\ & - \mathbb{E} \left[S \tilde{V}^\top (I - \eta S)^{k+1} (I - (I - \eta S)^k) \right] \\ = & \sum_{i=0}^k \mathbb{E} \left[S \tilde{V}^\top (\tilde{V} S (\tilde{W} - I) + (V + \eta I) S) (I - (I - \eta S)^i)^2 \right] \end{aligned}$$

$$\begin{aligned}
& - \sum_{i=0}^k \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \right] - \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (\mathbf{I} - \eta \mathbf{S})^{k+1} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k) \right] \\
& = \sum_{i=0}^k \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top \tilde{\mathbf{V}} \mathbf{S} (\tilde{\mathbf{W}} - \mathbf{I}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i)^2 \right] + \sum_{i=0}^k \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i)^2 \right] \\
& - \sum_{i=0}^k \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \right] - \mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (\mathbf{I} - \eta \mathbf{S})^{k+1} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k) \right]
\end{aligned}$$

We apply Lemma C.16 and Lemma C.17 to each term, similarly define $\Delta_i^{\tilde{\mathbf{W}}}$ for $i \in [k] \cup \{0\}$ as the sum of all interaction terms for the i -th term in the summation of dynamics of $\tilde{\mathbf{W}}$. There exists diagonal matrices $\mathbf{A}_i^{\tilde{\mathbf{W}}}, \mathbf{B}_i^{\tilde{\mathbf{W}}}, \mathbf{C}_i^{\tilde{\mathbf{W}}}, \mathbf{D}_i^{\tilde{\mathbf{W}}}, \mathbf{E}_i^{\tilde{\mathbf{W}}}$ with their operator norm $O\left(\frac{1}{\log^3 d}\right)$, such that

$$\begin{aligned}
\Delta_i^{\tilde{\mathbf{W}}} & = (\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \mathbf{A}_i^{\tilde{\mathbf{W}}} + (\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \mathbf{B}_i^{\tilde{\mathbf{W}}} + O\left(\frac{1}{d}\right) \text{tr}(\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \mathbf{C}_i^{\tilde{\mathbf{W}}} \\
& + O\left(\frac{1}{d}\right) \text{tr}((\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \Lambda^{\tilde{\mathbf{V}}}) \mathbf{D}_i^{\tilde{\mathbf{W}}} + O\left(\frac{1}{d}\right) \text{tr}((\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \Lambda^{\tilde{\mathbf{V}^2}}) \mathbf{E}_i^{\tilde{\mathbf{W}}}
\end{aligned}$$

We define $\Delta_{-1}^{\tilde{\mathbf{W}}}$ as the interaction term brought by the last term

$$\mathbb{E} \left[\mathbf{S} \tilde{\mathbf{V}}^\top (\mathbf{I} - \eta \mathbf{S})^{k+1} (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k) \right].$$

It is clear that $\|\Delta_{-1}^{\tilde{\mathbf{W}}}\|_{op} \leq O((1 - \eta)^k)$. Similarly denote

$$\hat{\Delta}^{\tilde{\mathbf{W}}} = \sum_{i=0}^k \Delta_i^{\tilde{\mathbf{W}}} - \Delta_{-1}^{\tilde{\mathbf{W}}},$$

then there exist diagonal matrices $\mathbf{A}^{\tilde{\mathbf{W}}}, \mathbf{B}^{\tilde{\mathbf{W}}}, \mathbf{C}^{\tilde{\mathbf{W}}}, \mathbf{D}^{\tilde{\mathbf{W}}}, \mathbf{E}^{\tilde{\mathbf{W}}}, \mathbf{F}_0^{\tilde{\mathbf{W}}}$ satisfying $\|\mathbf{A}^{\tilde{\mathbf{W}}}\|, \|\mathbf{B}^{\tilde{\mathbf{W}}}\| \leq O\left(\frac{1}{\log^2 d}\right), \|\mathbf{C}^{\tilde{\mathbf{W}}}\|, \|\mathbf{D}^{\tilde{\mathbf{W}}}\|, \|\mathbf{E}^{\tilde{\mathbf{W}}}\| \leq O\left(\frac{1}{d \log^2 d}\right), \|\mathbf{F}_0^{\tilde{\mathbf{W}}}\| \leq O((1 - \eta)^k)$ such that

$$\begin{aligned}
\hat{\Delta}^{\tilde{\mathbf{W}}} & = (\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \mathbf{A}^{\tilde{\mathbf{W}}} + (\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \mathbf{B}^{\tilde{\mathbf{W}}} + \text{tr}(\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \mathbf{C}^{\tilde{\mathbf{W}}} \\
& + \text{tr}((\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \Lambda^{\tilde{\mathbf{V}}}) \mathbf{D}^{\tilde{\mathbf{W}}} + \text{tr}((\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \Lambda^{\tilde{\mathbf{V}^2}}) \mathbf{E}^{\tilde{\mathbf{W}}} + \mathbf{F}_0^{\tilde{\mathbf{W}}}.
\end{aligned}$$

Denote $\mathbf{F}_1^{\tilde{\mathbf{W}}}$ to be the sum of all $O((1 - \eta)^k)$ terms in the dynamics of $\tilde{\mathbf{W}}$, $\mathbf{F}^{\tilde{\mathbf{W}}} = \mathbf{F}_0^{\tilde{\mathbf{W}}} + \mathbf{F}_1^{\tilde{\mathbf{W}}}$ and $\Delta^{\tilde{\mathbf{W}}} = \hat{\Delta}^{\tilde{\mathbf{W}}} + \mathbf{F}_1^{\tilde{\mathbf{W}}}$. Thus we have

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \tilde{\mathbf{W}}} & = \sum_{i=0}^k \mathbf{U} \left(1 - (1 - \eta)^i \right)^2 \Lambda^{\tilde{\mathbf{V}}} (\Lambda^{\tilde{\mathbf{V}}} \Lambda^{\tilde{\mathbf{W}}} + \eta \mathbf{I}) \mathbf{U}^\top - \sum_{i=0}^k \mathbf{U} \left(1 - (1 - \eta)^i \right) \Lambda^{\tilde{\mathbf{V}}} (\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \mathbf{U}^\top + \mathbf{U} \Delta^{\tilde{\mathbf{W}}} \mathbf{U}^\top \\
& = \mathbf{U} \left[\left(k + 1 - \frac{2(1 - (1 - \eta)^{k+1})}{\eta} + \frac{1 - (1 - \eta)^{2k+2}}{\eta(2 - \eta)} \right) \Lambda^{\tilde{\mathbf{V}}} (\Lambda^{\tilde{\mathbf{V}}} \Lambda^{\tilde{\mathbf{W}}} + \eta \mathbf{I}) \right] \mathbf{U}^\top \\
& - \mathbf{U} \left[\left(k + 1 - \frac{1 - (1 - \eta)^{k+1}}{\eta} \right) \Lambda^{\tilde{\mathbf{V}}} (\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \right] \mathbf{U}^\top + \mathbf{U} \hat{\Delta}^{\tilde{\mathbf{W}}} \mathbf{U}^\top \\
& = \mathbf{U} \left[\left(k + 1 - \frac{2}{\eta} + \frac{1}{\eta(2 - \eta)} \right) \Lambda^{\tilde{\mathbf{V}^2}} \Lambda^{\tilde{\mathbf{W}}} - \left(k + 1 - \frac{1}{\eta} \right) \Lambda^{\tilde{\mathbf{V}^2}} - \frac{1 - \eta}{2 - \eta} \Lambda^{\tilde{\mathbf{V}}} + \Delta^{\tilde{\mathbf{W}}} \right] \mathbf{U}^\top
\end{aligned}$$

Moreover, $\Delta^{\tilde{\mathbf{W}}}$ has the form

$$\hat{\Delta}^{\tilde{\mathbf{W}}} = (\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \mathbf{A}^{\tilde{\mathbf{W}}} + (\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \mathbf{B}^{\tilde{\mathbf{W}}} + \text{tr}(\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \mathbf{C}^{\tilde{\mathbf{W}}}$$

$$+ \text{tr} \left((\Lambda^{\tilde{\mathbf{V}}} + \eta \mathbf{I}) \Lambda^{\tilde{\mathbf{V}}} \right) \mathbf{D}^{\tilde{\mathbf{W}}} + \text{tr} \left((\mathbf{I} - \Lambda^{\tilde{\mathbf{W}}}) \Lambda^{\tilde{\mathbf{V}^2}} \right) \mathbf{E}^{\tilde{\mathbf{W}}} + \mathbf{F}^{\tilde{\mathbf{W}}}.$$

Since $\|\mathbf{A} + \mathbf{B}\|_{op} \leq \|\mathbf{A}\|_{op} + \|\mathbf{B}\|_{op}$ and $\|\mathbf{AB}\|_{op} \leq \|\mathbf{A}\|_{op} \|\mathbf{B}\|_{op}$, it is obvious that

$$\|\Delta^{\tilde{\mathbf{V}}}\|_{op} \leq O\left(\frac{1}{\log^2 d}\right), \quad \|\Delta^{\tilde{\mathbf{W}}}\|_{op} \leq O\left(\frac{1}{\log^2 d}\right)$$

□

After obtaining the estimation of the gradient by lemma C.5, we can decompose the gradient updates into the dynamics along each eigenspace \mathbf{u}_i , which can be characterized by the following lemma.

Lemma C.6. Suppose $\tilde{\mathbf{V}} = \sum_{j=1}^d \lambda_j^{\tilde{\mathbf{V}}} \mathbf{u}_j \mathbf{u}_j^\top$, $\tilde{\mathbf{W}} = \sum_{j=1}^d \lambda_j^{\tilde{\mathbf{W}}} \mathbf{u}_j \mathbf{u}_j^\top$. The dynamics of the eigenvalues of $\tilde{\mathbf{V}}$ and $\tilde{\mathbf{W}}$ are given by the following equations:

$$\begin{aligned} \frac{d\lambda_j^{\tilde{\mathbf{V}}}}{dt} &= - \left[(k+1) \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right)^2 + \frac{2}{\eta} \lambda_j^{\tilde{\mathbf{W}}} \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right) + \frac{1}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{W}^2}} \right] \lambda_j^{\tilde{\mathbf{V}}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{W}}} - 1 + \delta_j^{\tilde{\mathbf{V}}} \\ \frac{d\lambda_j^{\tilde{\mathbf{W}}}}{dt} &= \left(k+1 - \frac{1}{\eta} \right) \lambda_j^{\tilde{\mathbf{V}^2}} \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right) + \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{V}^2}} \lambda_j^{\tilde{\mathbf{W}}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{V}}} - \delta_j^{\tilde{\mathbf{W}}} \end{aligned}$$

where $|\delta_j^{\tilde{\mathbf{V}}}| \leq O\left(\frac{1}{\log^2 d}\right)$, $|\delta_j^{\tilde{\mathbf{W}}}| \leq O\left(\frac{1}{\log^2 d}\right)$.

Proof. This is directly obtained from Lemma C.5. □

C.3 PROOF OF THE MAIN THEOREM 4.1

In this section, we prove Theorem 4.1, which characterizes the CoT loss of the trained transformer. First, we restate the theorem.

Theorem C.1 (Global Convergence). Suppose $n = \Theta(d \log^5 d)$, $\eta \in (0.1, 0.9)$, $k = \lceil c \log d \rceil$, $c \log\left(\frac{1}{1-\eta}\right) > 2$. Under Assumption 4.1 with some constant $\sigma > \frac{3(1-\eta)}{(2-\eta)} \frac{1}{k+1}$, if we run gradient flow on the population loss in Equation (6), then after time $t = O(\log d + \log \frac{1}{\epsilon})$, we have $\mathcal{L}^{\text{CoT}}(t) \leq \epsilon$ for any $\epsilon \geq \Theta\left(\frac{\log d}{d^{c \log\left(\frac{1}{1-\eta}\right)-2}}\right)$.

Proof. According to the previous sections, we can reduce the original optimization problem to Equation (13), and consider the equivalent reduced model (Definition C.1). By Lemma C.5, we fully characterized the gradient expression, which decomposes the gradient of $\tilde{\mathbf{V}}$ and $\tilde{\mathbf{W}}$ into main signal terms with large norm at initialization (terms before $\Delta^{\tilde{\mathbf{V}}}, \Delta^{\tilde{\mathbf{W}}}$) and interaction terms ($\Delta^{\tilde{\mathbf{V}}}, \Delta^{\tilde{\mathbf{W}}}$) with bounded norm $O(\frac{1}{\log^2 d})$ for all $t > 0$.

The decomposition motivates us to conduct a stage-wise analysis. We first analyze the dynamics in **Stage 1** when the distance between the parameters $\tilde{\mathbf{V}}, \tilde{\mathbf{W}}$ and the ground-truth is larger than $O(\frac{1}{\log^2 d})$. In this stage, the bounded error can be dominated by the signal terms in the gradient, leading to nearly independent dynamics along each direction \mathbf{u}_i . After this stage, we enter **Stage 2** as a local convergence phase. We describe the dynamics below in detail.

Stage 1 In the first stage, the dynamics are dominated by the main terms, and the interaction terms $\Delta^{\tilde{\mathbf{V}}}, \Delta^{\tilde{\mathbf{W}}}$ can be somehow be ignored. Specifically, by Lemma C.6, given the dynamics of $\lambda_j^{\tilde{\mathbf{V}}}, \lambda_j^{\tilde{\mathbf{W}}}$:

$$\begin{aligned} \frac{d\lambda_j^{\tilde{\mathbf{V}}}}{dt} &= - \left[(k+1) \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right)^2 + \frac{2}{\eta} \lambda_j^{\tilde{\mathbf{W}}} \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right) + \frac{1}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{W}^2}} \right] \lambda_j^{\tilde{\mathbf{V}}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{W}}} - 1 + \delta_j^{\tilde{\mathbf{V}}} \\ \frac{d\lambda_j^{\tilde{\mathbf{W}}}}{dt} &= \left(k+1 - \frac{1}{\eta} \right) \lambda_j^{\tilde{\mathbf{V}^2}} \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right) + \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{V}^2}} \lambda_j^{\tilde{\mathbf{W}}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{V}}} - \delta_j^{\tilde{\mathbf{W}}} \end{aligned}$$

we can conclude that the dynamics of the eigenvalue $\lambda_j^{\tilde{V}}, \lambda_j^{\tilde{W}}$ mainly depend on themselves when the main term (terms before $\delta_j^{\tilde{W}}, \delta_j^{\tilde{V}}$) are larger than $O(\frac{1}{\log^2 d})$, which is within the stage 1. That is, the dynamics within the subspace $\mathbf{u}_i \mathbf{u}_i^\top$ for \tilde{V}, \tilde{W} are almost independent with other subspaces. In this stage, we focus on the analysis of $\lambda_j^{\tilde{V}}, \lambda_j^{\tilde{W}}$ depending on their own value.

The first stage can be further divided into two phases.

Stage 1, Phase 1. At the beginning of training, we have

$$\lambda_j^{\tilde{V}}(0) + \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}}(0))} < -\sigma + \frac{3(1-\eta)}{(2-\eta)} \frac{1}{k+1} < 0$$

then by Lemma C.8, we can prove an upper bound of $\lambda_j^{\tilde{V}}$ when $\lambda_j^{\tilde{W}} \leq 1 - (k+1)^{-\frac{7}{12}}$,

$$\lambda_j^{\tilde{V}} < -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}$$

according to the dynamics for both sides. With this upper bound, we prove $\frac{d\lambda_j^{\tilde{W}}}{dt} \geq O(\frac{1}{k})$. Therefore, $\lambda_j^{\tilde{W}}$ will converge to $1 - (k+1)^{-\frac{7}{12}}$ in $t_1 = O(\log d)$ time (Lemma C.9).

Stage 1, Phase 2. After time t_1 , we have $\lambda_j^{\tilde{W}}$ very close to the ground-truth value 1. Meanwhile, the lower bound for $\lambda_j^{\tilde{V}}$ still holds, and it will further decrease. Specifically,

$$\lambda_j^{\tilde{W}}(t_1) = 1 - (k+1)^{-\frac{7}{12}} \quad \lambda_j^{\tilde{V}}(t_1) < -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}}(t_1))}$$

By Lemma C.10, we can prove that $\lambda_j^{\tilde{W}}$ will stay close to $1 - o(1)$:

$$1 - 2(k+1)^{-\frac{7}{12}} < \lambda_j^{\tilde{W}}(t) < 1 + (k+1)^{-\frac{7}{12}}$$

for any $t \geq t_1$. With this condition, a converging condition for $(\lambda_j^{\tilde{V}} + \eta)$ can be deduced from Lemma C.11:

$$\frac{d(\lambda_j^{\tilde{V}} + \eta)^2}{dt} \leq -\frac{1}{2\eta(2-\eta)} (\lambda_j^{\tilde{V}} + \eta)^2$$

Lemma C.11 shows that $|\lambda_j^{\tilde{V}} + \eta|$ converges to $(k+1)^{-\frac{1}{12}}$ in $t_2 = O(\log \log d)$ time.

Stage 2. Now the eigenvalues are already close to ground-truth:

$$|\lambda_j^{\tilde{V}}(t_1 + t_2) + \eta| = O((k+1)^{-\frac{1}{12}}), \quad |\lambda_j^{\tilde{W}}(t_1 + t_2) - 1| \leq 2(k+1)^{-\frac{7}{12}}.$$

According to the expansion of the error terms in Lemma C.5, we notice that $\delta_j^{\tilde{W}}$ and $\delta_j^{\tilde{V}}$ are always coupled with some individual residual like $(\Lambda^{\tilde{V}} + \eta \mathbf{I})$, $(\Lambda^{\tilde{W}} - \mathbf{I})$, or some weighted average $\frac{1}{d} \text{tr}((\Lambda^{\tilde{V}} + \eta \mathbf{I}) \Lambda^{\tilde{V}})$. Meanwhile, the coefficient of this kind of residual in the interaction terms is still upper bounded by $O(1/\log^2 d)$. That helps us to derive the PL-condition like gradient lower bound (Lemma C.12):

$$\begin{aligned} & \frac{d \text{tr}[(\Lambda^{\tilde{V}} + \eta \mathbf{I})^2]}{dt} + \frac{d \text{tr}[(\mathbf{I} - \Lambda^{\tilde{W}})^2]}{dt} \\ & \leq -\frac{1}{2\eta(2-\eta)} \text{tr}[(\Lambda^{\tilde{V}} + \eta \mathbf{I})^2] - \frac{\eta^2}{2}(k+1) \text{tr}[(\mathbf{I} - \Lambda^{\tilde{W}})^2] + \alpha \end{aligned}$$

where $\alpha = O((1 - \eta)^k) \geq 0$.

By Lemma C.12, we know $|\lambda_j^{\tilde{V}} + \eta|$ and $|1 - \lambda_j^{\tilde{W}}|$ converge to $\delta \in \left(\Theta\left(d^{\frac{\epsilon}{2} \log(1-\eta) + \frac{1}{2}}\right), 1\right)$ in $t_3 = O(\log \frac{1}{\delta})$ time. At this time, there exist diagonal matrices \mathbf{A} and \mathbf{B} satisfying $\|\mathbf{A}\|_{op} \leq \Theta(1)$ and $\|\mathbf{B}\|_{op} \leq \Theta(1)$ such that

$$\mathbf{\Lambda}^{\tilde{V}} = -\eta \mathbf{I} + \delta \cdot \mathbf{A} \quad \mathbf{\Lambda}^{\tilde{W}} = \mathbf{I} + \delta \cdot \mathbf{B}.$$

Now we consider the CoT loss given by Lemma C.3

$$\begin{aligned} \mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \sum_{i=0}^{k-1} \left\| (\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) \mathbf{w}_i - (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} \mathbf{w}^* \right\|_2^2 \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}}) \mathbf{w}_k - (\tilde{\mathbf{V}} \mathbf{S} + \mathbf{I}) \mathbf{w}^* \right\|_2^2. \end{aligned}$$

Apply Lemma C.13, we directly obtain that

$$\mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}) = O(\delta^2 d \log d).$$

Since $\delta \in \left(\Theta\left(d^{\frac{\epsilon}{2} \log(1-\eta) + \frac{1}{2}}\right), 1\right)$, the CoT loss is smaller than $\epsilon = \Theta(d^{c \log(1-\eta) + 2} \log d)$. The local convergence takes $t_3 = O(\log \frac{1}{\delta}) = O(\log \frac{1}{\epsilon})$. Considering all stages, at time $t = t_1 + t_2 + t_3 = O(\log d) + O(\frac{1}{\epsilon})$, we have

$$\mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}) \leq \epsilon.$$

□

C.3.1 TECHNICAL LEMMA IN APPENDIX C.3

Lemma C.7. Assume $\lambda_j^{\tilde{W}} \leq 1 - (k+1)^{-\frac{7}{12}}$, if $-\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})} \leq \lambda_j^{\tilde{V}} < 0$, it holds that

$$\frac{d\left(\lambda_j^{\tilde{V}} + \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right)}{dt} < 0 \quad (14)$$

Proof. Directly consider the derivative

$$\frac{d\left(\lambda_j^{\tilde{V}} + \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right)}{dt} = \frac{d\lambda_j^{\tilde{V}}}{dt} + \frac{3(1-\eta)}{2(k+1)(2-\eta)} \frac{1}{(1-\lambda_j^{\tilde{W}})^2} \frac{d\lambda_j^{\tilde{W}}}{dt}$$

Substitute the derivatives with the equations in Lemma C.6, we have

$$\begin{aligned} &\frac{d\lambda_j^{\tilde{V}}}{dt} + \frac{3(1-\eta)}{2(k+1)(2-\eta)} \frac{1}{(1-\lambda_j^{\tilde{W}})^2} \frac{d\lambda_j^{\tilde{W}}}{dt} \\ &= - \left[(k+1) \left(1 - \lambda_j^{\tilde{W}}\right)^2 + \frac{2}{\eta} \lambda_j^{\tilde{W}} \left(1 - \lambda_j^{\tilde{W}}\right) + \frac{1}{\eta(2-\eta)} \lambda_j^{\tilde{W}^2} \right] \lambda_j^{\tilde{V}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{W}} - \left(1 + \delta_j^{\tilde{V}}\right) \\ &\quad + \frac{3(1-\eta)}{2(k+1)(2-\eta)} \frac{1}{(1-\lambda_j^{\tilde{W}})^2} \left[\left(k+1 - \frac{1}{\eta}\right) \lambda_j^{\tilde{V}^2} \left(1 - \lambda_j^{\tilde{W}}\right) + \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{V}^2} \lambda_j^{\tilde{W}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{V}} - \delta_j^{\tilde{W}} \right] \end{aligned}$$

Since $-\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})} \leq \lambda_j^{\tilde{V}} < 0$, we have

$$\frac{d\lambda_j^{\tilde{V}}}{dt} + \frac{3(1-\eta)}{2(k+1)(2-\eta)} \frac{1}{(1-\lambda_j^{\tilde{W}})^2} \frac{d\lambda_j^{\tilde{W}}}{dt}$$

$$\begin{aligned}
&\leq \left[(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^2 + \frac{2}{\eta} \lambda_j^{\widetilde{\mathbf{W}}} \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right) + \frac{1}{\eta(2-\eta)} \lambda_j^{\widetilde{\mathbf{W}^2}} \right] \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)} \\
&\quad + \frac{1-\eta}{2-\eta} \lambda_j^{\widetilde{\mathbf{W}}} - \left(1 + \delta_j^{\widetilde{\mathbf{V}}}\right) \\
&\quad + \frac{3(1-\eta)}{2(k+1)(2-\eta)} \frac{1}{\left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^2} \left[(k+1) \left(\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)} \right)^2 \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right) \right. \\
&\quad \left. + \frac{1-\eta}{\eta(2-\eta)} \left(\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)} \right)^2 \lambda_j^{\widetilde{\mathbf{W}}} - \delta_j^{\widetilde{\mathbf{W}}} \right] \\
&= \frac{3(1-\eta)}{2(2-\eta)} \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right) + \frac{1}{k+1} \frac{3(1-\eta)}{\eta(2-\eta)} \lambda_j^{\widetilde{\mathbf{W}}} + \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)} \frac{3(1-\eta)}{2\eta(2-\eta)^2} \lambda_j^{\widetilde{\mathbf{W}^2}} \\
&\quad + \frac{1-\eta}{2-\eta} \lambda_j^{\widetilde{\mathbf{W}}} - \left(1 + \delta_j^{\widetilde{\mathbf{V}}}\right) + \left[\frac{3(1-\eta)}{2(2-\eta)} \right]^3 \frac{1}{(k+1)^2 \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^3} \\
&\quad + \left[\frac{3(1-\eta)}{2(2-\eta)} \right]^3 \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\widetilde{\mathbf{W}}} \frac{1}{(k+1)^3 \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^4} - \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^2} \delta_j^{\widetilde{\mathbf{W}}} \\
&= \left[\frac{1-\eta}{2(2-\eta)} \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right) - \frac{1}{2-\eta} \right] + \frac{1}{k+1} \frac{3(1-\eta)}{2-\eta} \lambda_j^{\widetilde{\mathbf{W}}} \\
&\quad + \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)} \frac{3(1-\eta)}{2\eta(2-\eta)^2} \lambda_j^{\widetilde{\mathbf{W}^2}} - \delta_j^{\widetilde{\mathbf{V}}} + \left[\frac{3(1-\eta)}{2(2-\eta)} \right]^3 \frac{1}{(k+1)^2 \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^3} \\
&\quad + \left[\frac{3(1-\eta)}{2(2-\eta)} \right]^3 \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\widetilde{\mathbf{W}}} \frac{1}{(k+1)^3 \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^4} - \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^2} \delta_j^{\widetilde{\mathbf{W}}}
\end{aligned}$$

Put in the assumption on $\lambda_j^{\widetilde{\mathbf{W}}}$ that $\lambda_j^{\widetilde{\mathbf{W}}} \leq 1 - (k+1)^{-\frac{7}{12}}$, we have

$$\begin{aligned}
&\frac{d\lambda_j^{\widetilde{\mathbf{V}}}}{dt} + \frac{3(1-\eta)}{2(k+1)(2-\eta)} \frac{1}{\left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)^2} \frac{d\lambda_j^{\widetilde{\mathbf{W}}}}{dt} \\
&\leq -\frac{1+\eta}{2(2-\eta)} + \frac{1}{k+1} \frac{3(1-\eta)}{\eta(2-\eta)} + \frac{1}{(k+1)^{\frac{5}{12}}} \frac{3(1-\eta)}{2\eta(2-\eta)^2} + \left| \delta_j^{\widetilde{\mathbf{V}}} \right| + \left[\frac{3(1-\eta)}{2(2-\eta)} \right]^3 \frac{1}{(k+1)^{\frac{1}{4}}} \\
&\quad + \left[\frac{3(1-\eta)}{2(2-\eta)} \right]^3 \frac{1-\eta}{\eta(2-\eta)} \frac{1}{(k+1)^{\frac{2}{3}}} + \frac{3(1-\eta)}{2(2-\eta)} \left| \delta_j^{\widetilde{\mathbf{W}}} \right| \\
&= -\frac{1+\eta}{2(2-\eta)} + O\left(\frac{1}{\log^{\frac{1}{4}} d} \right)
\end{aligned}$$

□

Lemma C.8 (Upper bound of $\lambda_j^{\widetilde{\mathbf{V}}}$). *Under Assumption 4.1, if $\lambda_j^{\widetilde{\mathbf{W}}} \leq 1 - (k+1)^{-\frac{7}{12}}$, it holds that*

$$\lambda_j^{\widetilde{\mathbf{V}}} < -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1) \left(1 - \lambda_j^{\widetilde{\mathbf{W}}}\right)} \quad (15)$$

Proof. We prove by induction. First, check the initialization $\lambda_j^{\tilde{V}}(0) \leq -\sigma$, $\sigma \leq \lambda_j^{\tilde{W}}(0) \leq \frac{1}{2}$. If $\sigma \geq \frac{3(1-\eta)}{2-\eta} \frac{1}{k+1}$, then we have

$$\lambda_j^{\tilde{V}}(0) + \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}}(0))} < -\sigma + \frac{3(1-\eta)}{(2-\eta)} \frac{1}{k+1} \leq 0$$

If the inequality holds until some time t_1 , that is for any $t < t_1$, we have

$$\lambda_j^{\tilde{V}}(t) < -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}}(t))}$$

but

$$\lambda_j^{\tilde{V}}(t_1) \geq -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}}(t_1))}$$

By Lemma C.7, we have

$$\left. \frac{d\left(\lambda_j^{\tilde{V}} + \frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right)}{dt} \right|_{t=t_1} < 0$$

Therefore, there exists some time $t' < t_1$ such that

$$\lambda_j^{\tilde{V}}(t') \geq -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}}(t'))}$$

which is a contradiction. Hence, the proof is complete. \square

Lemma C.9 ($\lambda_j^{\tilde{W}}$ converges to near optimal). *Under Assumption 4.1, it takes $O(\log d)$ time for $\lambda_j^{\tilde{W}}$ to converge to $1 - (k+1)^{-\frac{7}{12}}$.*

Proof. Recall the gradient of $\lambda_j^{\tilde{W}}$ in Lemma C.6

$$\frac{d\lambda_j^{\tilde{W}}}{dt} = \left(k+1 - \frac{1}{\eta}\right) \lambda_j^{\tilde{V}^2} (1 - \lambda_j^{\tilde{W}}) + \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{V}^2} \lambda_j^{\tilde{W}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{V}} - \delta_j^{\tilde{W}}$$

Substitute $\lambda_j^{\tilde{V}}$ with Lemma C.8, we have

$$\begin{aligned} \frac{d\lambda_j^{\tilde{W}}}{dt} &\geq \left(k+1 - \frac{1}{\eta}\right) \left(\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right)^2 (1 - \lambda_j^{\tilde{W}}) \\ &\quad + \frac{1-\eta}{\eta(2-\eta)} \left(\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right)^2 \lambda_j^{\tilde{W}} \\ &\quad - \frac{1-\eta}{2-\eta} \left(\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right) - \delta_j^{\tilde{W}} \\ &\geq \frac{4}{5}(k+1) \left(\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right)^2 (1 - \lambda_j^{\tilde{W}}) \\ &\quad - \frac{1-\eta}{2-\eta} \left(\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{W}})}\right) - |\delta_j^{\tilde{W}}| \end{aligned}$$

$$\begin{aligned}
&= \frac{3}{10} \frac{(1-\eta)^2}{(2-\eta)^2} \frac{1}{(k+1)(1-\lambda_j^{\tilde{\mathbf{W}}})} - |\delta_j^{\tilde{\mathbf{W}}}| \\
&\geq \frac{1}{5} \frac{(1-\eta)^2}{(2-\eta)^2} \frac{1}{(k+1)(1-\lambda_j^{\tilde{\mathbf{W}}})} \\
&\geq \frac{1}{5} \frac{(1-\eta)^2}{(2-\eta)^2} \frac{1}{k+1}
\end{aligned}$$

In $O(\log d)$ time, $\lambda_j^{\tilde{\mathbf{W}}}$ can converge to $1 - (k+1)^{-\frac{7}{12}}$. \square

Lemma C.10. Assume $\lambda_j^{\tilde{\mathbf{W}}}(t_1) = 1 - (k+1)^{-\frac{7}{12}}$ and $\lambda_j^{\tilde{\mathbf{V}}}(t_1) < -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{\mathbf{W}}}(t_1))}$, for any $t \geq t_1$ it holds that

$$1 - 2(k+1)^{-\frac{7}{12}} < \lambda_j^{\tilde{\mathbf{W}}}(t) < 1 + (k+1)^{-\frac{7}{12}}.$$

Proof. First, it is clear that the inequality holds at time t_1 . If the inequality doesn't hold, then there exists $t' > t_1$ such that

$$\begin{aligned}
1 - 2(k+1)^{-\frac{7}{12}} &< \lambda_j^{\tilde{\mathbf{W}}}(t) < 1 + (k+1)^{-\frac{7}{12}} && \text{for any } t_1 \leq t < t' \\
\lambda_j^{\tilde{\mathbf{W}}}(t') &= 1 - 2(k+1)^{-\frac{7}{12}}
\end{aligned}$$

or

$$\begin{aligned}
1 - 2(k+1)^{-\frac{7}{12}} &< \lambda_j^{\tilde{\mathbf{W}}}(t) < 1 + (k+1)^{-\frac{7}{12}} && \text{for any } t_1 \leq t < t' \\
\lambda_j^{\tilde{\mathbf{W}}}(t') &= 1 + (k+1)^{-\frac{7}{12}}
\end{aligned}$$

In the first case, it suffices to prove

$$\lambda_j^{\tilde{\mathbf{V}}}(t') \leq -\frac{3(1-\eta)}{2(2-\eta)} (k+1)^{-\frac{5}{12}} < -\frac{3(1-\eta)}{2(2-\eta)} \frac{1}{(k+1)(1-\lambda_j^{\tilde{\mathbf{W}}}(t'))}$$

to show

$$\left. \frac{d\lambda_j^{\tilde{\mathbf{W}}}}{dt} \right|_{t=t'} > 0$$

which says there exists $t_1 \leq t'' < t'$ such that

$$\lambda_j^{\tilde{\mathbf{W}}}(t'') \leq 1 - 2(k+1)^{-\frac{7}{12}}$$

and leads to a contradiction. Recall the gradient of $\lambda_j^{\tilde{\mathbf{V}}}$ in Lemma C.6

$$\begin{aligned}
\frac{d\lambda_j^{\tilde{\mathbf{V}}}}{dt} &= - \left[(k+1)(1-\lambda_j^{\tilde{\mathbf{W}}})^2 + \frac{2}{\eta} \lambda_j^{\tilde{\mathbf{W}}} (1-\lambda_j^{\tilde{\mathbf{W}}}) + \frac{1}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{W}^2}} \right] \lambda_j^{\tilde{\mathbf{V}}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{W}}} - (1+\delta_j^{\tilde{\mathbf{V}}}) \\
&\leq - \left[4(k+1)^{-\frac{1}{6}} + \frac{4}{\eta} \left[(k+1)^{-\frac{7}{12}} + (k+1)^{-\frac{7}{6}} \right] + \frac{1}{\eta(2-\eta)} \left[1 + 2(k+1)^{-\frac{7}{12}} + (k+1)^{-\frac{7}{6}} \right] \right] \lambda_j^{\tilde{\mathbf{V}}} \\
&\quad - \frac{1}{2-\eta} + \frac{1-\eta}{2-\eta} (k+1)^{-\frac{7}{12}} + |\delta_j^{\tilde{\mathbf{V}}}| \\
&\leq - \frac{2}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{V}}} - \frac{1}{2(2-\eta)}
\end{aligned}$$

and thus we have

$$\lambda_j^{\tilde{\mathbf{V}}}(t) \leq C e^{-\frac{2}{\eta(2-\eta)}(t-t_1)} - \frac{\eta}{4}$$

If $C \leq 0$, then

$$\lambda_j^{\tilde{\mathbf{V}}}(t') \leq -\frac{\eta}{4} \leq -\frac{3(1-\eta)}{2(2-\eta)} (k+1)^{-\frac{5}{12}}$$

else

$$\lambda_j^{\tilde{V}}(t') \leq \lambda_j^{\tilde{V}}(t_1) = -\frac{3(1-\eta)}{2(2-\eta)}(k+1)^{-\frac{5}{12}}$$

In the second case,

$$\lambda_j^{\tilde{V}}(t) \leq -\frac{3(1-\eta)}{2(2-\eta)}(k+1)^{-\frac{5}{12}}$$

still holds for any $t_1 \leq t \leq t'$. Recall the gradient of $\lambda_j^{\tilde{W}}$ in Lemma C.6

$$\begin{aligned} \left. \frac{d\lambda_j^{\tilde{W}}}{dt} \right|_{t=t'} &= \left(k+1 - \frac{1}{\eta} \right) \lambda_j^{\tilde{V}^2} (1 - \lambda_j^{\tilde{W}}) + \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{V}^2} \lambda_j^{\tilde{W}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{V}} - \delta_j^{\tilde{W}} \\ &= - \left(k+1 - \frac{1}{\eta} \right) \lambda_j^{\tilde{V}^2} (k+1)^{-\frac{7}{12}} + \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{V}^2} \left[1 + (k+1)^{-\frac{7}{12}} \right] + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{V}} - \delta_j^{\tilde{W}} \\ &\leq - \frac{(k+1)^{\frac{5}{12}}}{2} \lambda_j^{\tilde{V}^2} + \frac{1-\eta}{2(2-\eta)} \lambda_j^{\tilde{V}} + \left| \delta_j^{\tilde{W}} \right| \\ &\leq - \frac{9(1-\eta)^2}{8(2-\eta)^2} (k+1)^{-\frac{5}{12}} - \frac{3(1-\eta)^2}{4(2-\eta)^2} (k+1)^{-\frac{5}{12}} + \left| \delta_j^{\tilde{W}} \right| \\ &\leq - \frac{(1-\eta)^2}{(2-\eta)^2} (k+1)^{-\frac{5}{12}} \end{aligned}$$

There exists $t_1 \leq t'' < t'$ such that

$$\lambda_j^{\tilde{W}}(t'') \geq 1 + (k+1)^{-\frac{7}{12}}$$

which is a contradiction. Hence, the proof is complete. \square

Lemma C.11 ($\lambda_j^{\tilde{V}}$ converges to near optimal). *Assume*

$$1 - 2(k+1)^{-\frac{7}{12}} < \lambda_j^{\tilde{W}}(t) < 1 + (k+1)^{-\frac{7}{12}}$$

then it takes $O(\log \log d)$ time for $\left| \lambda_j^{\tilde{V}} + \eta \right|$ to converge to $(k+1)^{-\frac{1}{12}}$.

Proof. From Lemma C.10, we know

$$\begin{aligned} \frac{d\lambda_j^{\tilde{V}}}{dt} &= - \left[(k+1) \left(1 - \lambda_j^{\tilde{W}} \right)^2 + \frac{2}{\eta} \lambda_j^{\tilde{W}} \left(1 - \lambda_j^{\tilde{W}} \right) + \frac{1}{\eta(2-\eta)} \lambda_j^{\tilde{W}^2} \right] \lambda_j^{\tilde{V}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{W}} - \left(1 + \delta_j^{\tilde{V}} \right) \\ &\leq - \left[4(k+1)^{-\frac{1}{6}} + \frac{4}{\eta} \left[(k+1)^{-\frac{7}{12}} + (k+1)^{-\frac{7}{6}} \right] + \frac{1}{\eta(2-\eta)} \left[1 + 2(k+1)^{-\frac{7}{12}} + (k+1)^{-\frac{7}{6}} \right] \right] \lambda_j^{\tilde{V}} \\ &\quad - \frac{1}{2-\eta} + \frac{1-\eta}{2-\eta} (k+1)^{-\frac{7}{12}} + \left| \delta_j^{\tilde{V}} \right| \\ &= - \left[\frac{1}{\eta(2-\eta)} + O\left((k+1)^{-\frac{1}{6}} \right) \right] \left(\lambda_j^{\tilde{V}} + \eta \right) + O\left((k+1)^{-\frac{1}{6}} \right) \end{aligned}$$

and

$$\begin{aligned} \frac{d\lambda_j^{\tilde{V}}}{dt} &= - \left[(k+1) \left(1 - \lambda_j^{\tilde{W}} \right)^2 + \frac{2}{\eta} \lambda_j^{\tilde{W}} \left(1 - \lambda_j^{\tilde{W}} \right) + \frac{1}{\eta(2-\eta)} \lambda_j^{\tilde{W}^2} \right] \lambda_j^{\tilde{V}} + \frac{1-\eta}{2-\eta} \lambda_j^{\tilde{W}} - \left(1 + \delta_j^{\tilde{V}} \right) \\ &\geq - \left[-\frac{2}{\eta} \left[(k+1)^{-\frac{7}{12}} + (k+1)^{-\frac{7}{6}} \right] + \frac{1}{\eta(2-\eta)} \left[1 - 4(k+1)^{-\frac{7}{12}} + 4(k+1)^{-\frac{7}{6}} \right] \right] \lambda_j^{\tilde{V}} \\ &\quad - \frac{1}{2-\eta} - \frac{2(1-\eta)}{2-\eta} (k+1)^{-\frac{7}{12}} - \left| \delta_j^{\tilde{V}} \right| \\ &= - \left[\frac{1}{\eta(2-\eta)} + O\left((k+1)^{-\frac{1}{6}} \right) \right] \left(\lambda_j^{\tilde{V}} + \eta \right) + O\left((k+1)^{-\frac{1}{6}} \right) \end{aligned}$$

Therefore,

$$\begin{aligned} \frac{d(\lambda_j^{\tilde{V}} + \eta)^2}{dt} &= 2(\lambda_j^{\tilde{V}} + \eta) \frac{d\lambda_j^{\tilde{V}}}{dt} \\ &\leq -\left[\frac{1}{\eta(2-\eta)} + O((k+1)^{-\frac{1}{6}})\right](\lambda_j^{\tilde{V}} + \eta)^2 + O((k+1)^{-\frac{1}{6}})(\lambda_j^{\tilde{V}} + \eta) \end{aligned}$$

If $|\lambda_j^{\tilde{V}} + \eta|$ converges to $\epsilon = (k+1)^{-\frac{1}{12}}$, then

$$\frac{d(\lambda_j^{\tilde{V}} + \eta)^2}{dt} \leq -\frac{1}{2\eta(2-\eta)}(\lambda_j^{\tilde{V}} + \eta)^2$$

Thus, there exists $c \leq \Theta(1)$ such that

$$\epsilon^2 = (\lambda_j^{\tilde{V}} + \eta)^2 \leq c^2 \exp\left(-\frac{1}{2\eta(2-\eta)}t\right)$$

In $O(\log(\frac{1}{\epsilon})) = O(\log \log d)$ time, $|\lambda_j^{\tilde{V}} + \eta|$ can converge to ϵ . \square

Lemma C.12 (Local convergence). *Suppose $k = \lceil c \log d \rceil$. Assume*

$$|\lambda_j^{\tilde{W}}(t) - 1| \leq 2(k+1)^{-\frac{7}{12}} \quad |\lambda_j^{\tilde{V}}(t) + \eta| = O((k+1)^{-\frac{1}{12}}),$$

then there exists $\alpha = O((1-\eta)^k) \geq 0$ such that $\Lambda^{\tilde{V}}$ and $\Lambda^{\tilde{W}}$ comply with

$$\begin{aligned} &\frac{d \operatorname{tr}[(\Lambda^{\tilde{V}} + \eta \mathbf{I})^2]}{dt} + \frac{d \operatorname{tr}[(\mathbf{I} - \Lambda^{\tilde{W}})^2]}{dt} \\ &\leq -\frac{1}{2\eta(2-\eta)} \operatorname{tr}[(\Lambda^{\tilde{V}} + \eta \mathbf{I})^2] - \frac{\eta^2}{2}(k+1) \operatorname{tr}[(\mathbf{I} - \Lambda^{\tilde{W}})^2] + \alpha, \end{aligned}$$

thus $|\lambda_j^{\tilde{V}} + \eta|$ and $|1 - \lambda_j^{\tilde{W}}|$ can converge to $\epsilon \in (d^{-\frac{\epsilon}{2} \log(\frac{1}{1-\eta}) + \frac{1}{2}}, 1)$ in $O(\log \frac{1}{\epsilon})$ time.

Proof. Consider the error term in Lemma C.6 more carefully, we have

$$\begin{aligned} \frac{d\lambda_j^{\tilde{V}}}{dt} &= -\left[(k+1)(1 - \lambda_j^{\tilde{W}})^2 + \frac{2}{\eta}\lambda_j^{\tilde{W}}(1 - \lambda_j^{\tilde{W}}) + \frac{1}{\eta(2-\eta)}\lambda_j^{\tilde{W}^2}\right](\lambda_j^{\tilde{V}} + \eta) \\ &\quad + \eta(k+1)(1 - \lambda_j^{\tilde{W}})^2 + \left(\frac{3-2\eta}{2-\eta}\lambda_j^{\tilde{W}} - 1\right)(1 - \lambda_j^{\tilde{W}}) \\ &\quad + (\lambda_j^{\tilde{V}} + \eta)O\left(\frac{1}{\log^2 d}\right) + (1 - \lambda_j^{\tilde{W}})O\left(\frac{1}{\log^2 d}\right) \\ &\quad + \operatorname{tr}(\mathbf{I} - \Lambda^{\tilde{W}})O\left(\frac{1}{d \log^2 d}\right) + \operatorname{tr}((\mathbf{I} - \Lambda^{\tilde{W}})\Lambda^{\tilde{W}})O\left(\frac{1}{d \log^2 d}\right) \\ &\quad + O((1-\eta)^k) \end{aligned}$$

and

$$\begin{aligned} \frac{d\lambda_j^{\tilde{W}}}{dt} &= \left(k+1 - \frac{2}{\eta} + \frac{1}{\eta(2-\eta)}\right)\lambda_j^{\tilde{V}^2}(1 - \lambda_j^{\tilde{W}}) + \frac{1-\eta}{\eta(2-\eta)}\lambda_j^{\tilde{V}}(\lambda_j^{\tilde{V}} + \eta) \\ &\quad + (\lambda_j^{\tilde{V}} + \eta)O\left(\frac{1}{\log^2 d}\right) + (1 - \lambda_j^{\tilde{W}})O\left(\frac{1}{\log^2 d}\right) \\ &\quad + \operatorname{tr}(\mathbf{I} - \Lambda^{\tilde{W}})O\left(\frac{1}{d \log^2 d}\right) + \operatorname{tr}((\Lambda^{\tilde{V}} + \eta \mathbf{I})\Lambda^{\tilde{V}})O\left(\frac{1}{d \log^2 d}\right) \end{aligned}$$

$$+ \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{V}}^2 \right) O\left(\frac{1}{d \log^2 d}\right) + O\left((1 - \eta)^k\right)$$

Now we consider the decay rate of the distance between $\lambda_j^{\tilde{\mathbf{V}}}$, $\lambda_j^{\tilde{\mathbf{W}}}$ and their ground truth.

$$\begin{aligned} & \frac{d(\lambda_j^{\tilde{\mathbf{V}}} + \eta)^2}{dt} + \frac{d(\lambda_j^{\tilde{\mathbf{W}}} - 1)^2}{dt} \\ &= - \left[(k+1) \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right)^2 + \frac{2}{\eta} \lambda_j^{\tilde{\mathbf{W}}} \left(1 - \lambda_j^{\tilde{\mathbf{W}}}\right) + \frac{1}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{W}}^2} + O\left(\frac{1}{\log^2 d}\right) \right] (\lambda_j^{\tilde{\mathbf{V}}} + \eta)^2 \\ & \quad + \left(\frac{3-2\eta}{2-\eta} \lambda_j^{\tilde{\mathbf{W}}} - \frac{1-\eta}{\eta(2-\eta)} \lambda_j^{\tilde{\mathbf{V}}} - 1 + O\left(\frac{1}{\log^2 d}\right) \right) (1 - \lambda_j^{\tilde{\mathbf{W}}}) (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \\ & \quad - \left[\left(k+1 - \frac{2}{\eta} + \frac{1}{\eta(2-\eta)} \right) \lambda_j^{\tilde{\mathbf{V}}^2} - \eta(k+1) (\lambda_j^{\tilde{\mathbf{V}}} + \eta) + O\left(\frac{1}{\log^2 d}\right) \right] (1 - \lambda_j^{\tilde{\mathbf{W}}})^2 \\ & \quad + (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) O\left(\frac{1}{d \log^2 d}\right) \right) + (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{W}} \right) O\left(\frac{1}{d \log^2 d}\right) \\ & \quad + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) O\left(\frac{1}{d \log^2 d}\right) \right) + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{\Lambda} \tilde{\mathbf{V}} \right) O\left(\frac{1}{d \log^2 d}\right) \\ & \quad + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{V}}^2 \right) O\left(\frac{1}{d \log^2 d}\right) + O\left((1 - \eta)^k\right) \\ &= - \left[\frac{1}{\eta(2-\eta)} + O\left(\frac{1}{\log^{\frac{1}{6}} d}\right) \right] (\lambda_j^{\tilde{\mathbf{V}}} + \eta)^2 \\ & \quad + \left(\frac{2(1-\eta)}{2-\eta} + O\left(\frac{1}{\log^{\frac{1}{12}} d}\right) \right) (1 - \lambda_j^{\tilde{\mathbf{W}}}) (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \\ & \quad - \left[\eta^2(k+1) + O\left(k^{\frac{11}{12}}\right) \right] (1 - \lambda_j^{\tilde{\mathbf{W}}})^2 \\ & \quad + (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) O\left(\frac{1}{d \log^2 d}\right) \right) + (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{W}} \right) O\left(\frac{1}{d \log^2 d}\right) \\ & \quad + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) O\left(\frac{1}{d \log^2 d}\right) \right) + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{\Lambda} \tilde{\mathbf{V}} \right) O\left(\frac{1}{d \log^2 d}\right) \\ & \quad + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{V}}^2 \right) O\left(\frac{1}{d \log^2 d}\right) + O\left((1 - \eta)^k\right) \end{aligned}$$

Utilizing Mean Inequality, we have

$$\left| (1 - \lambda_j^{\tilde{\mathbf{W}}}) (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \right| \leq \frac{1}{2\sqrt{k}} (\lambda_j^{\tilde{\mathbf{V}}} + \eta)^2 + \frac{\sqrt{k}}{2} (1 - \lambda_j^{\tilde{\mathbf{W}}})^2$$

Insert the inequality into the equation and we have

$$\begin{aligned} & \frac{d(\lambda_j^{\tilde{\mathbf{V}}} + \eta)^2}{dt} + \frac{d(\lambda_j^{\tilde{\mathbf{W}}} - 1)^2}{dt} \\ & \leq - \left[\frac{1}{\eta(2-\eta)} + O\left(\frac{1}{\log^{\frac{1}{6}} d}\right) \right] (\lambda_j^{\tilde{\mathbf{V}}} + \eta)^2 - \left[\eta^2(k+1) + O\left(k^{\frac{11}{12}}\right) \right] (1 - \lambda_j^{\tilde{\mathbf{W}}})^2 \\ & \quad + (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) O\left(\frac{1}{d \log^2 d}\right) \right) + (\lambda_j^{\tilde{\mathbf{V}}} + \eta) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{W}} \right) O\left(\frac{1}{d \log^2 d}\right) \\ & \quad + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) O\left(\frac{1}{d \log^2 d}\right) \right) + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{\Lambda} \tilde{\mathbf{V}} \right) O\left(\frac{1}{d \log^2 d}\right) \\ & \quad + (1 - \lambda_j^{\tilde{\mathbf{W}}}) \text{tr} \left((\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}}) \mathbf{\Lambda} \tilde{\mathbf{V}}^2 \right) O\left(\frac{1}{d \log^2 d}\right) + O\left((1 - \eta)^k\right) \end{aligned}$$

There exist $\alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5 = O\left(\frac{1}{d \log^2 d}\right) \geq 0$ and $\alpha_6 = O\left((1 - \eta)^k\right) \geq 0$ such that

$$\begin{aligned} \frac{d(\lambda_j^{\tilde{V}} + \eta)^2}{dt} + \frac{d(\lambda_j^{\tilde{W}} - 1)^2}{dt} &\leq - \left[\frac{1}{\eta(2 - \eta)} + O\left(\frac{1}{\log^{\frac{1}{6}} d}\right) \right] (\lambda_j^{\tilde{V}} + \eta)^2 \\ &\quad - \left[\eta^2(k + 1) + O\left(k^{\frac{11}{12}}\right) \right] (1 - \lambda_j^{\tilde{W}})^2 \\ &\quad + \alpha_1 |\lambda_j^{\tilde{V}} + \eta| \operatorname{tr}(|I - \Lambda \tilde{W}|) + \alpha_2 |\lambda_j^{\tilde{V}} + \eta| \cdot \left| \operatorname{tr}((I - \Lambda \tilde{W}) \Lambda \tilde{W}) \right| \\ &\quad + \alpha_3 |1 - \lambda_j^{\tilde{W}}| \operatorname{tr}(|I - \Lambda \tilde{W}|) + \alpha_4 |1 - \lambda_j^{\tilde{W}}| \cdot \left| \operatorname{tr}((\Lambda \tilde{V} + \eta I) \Lambda \tilde{V}) \right| \\ &\quad + \alpha_5 |1 - \lambda_j^{\tilde{W}}| \cdot \left| \operatorname{tr}((I - \Lambda \tilde{W}) \Lambda \tilde{V}^2) \right| + \alpha_6. \end{aligned}$$

Notice that for diagonal matrices A and B , we have

$$\operatorname{tr}(AB) \leq |\operatorname{tr}(AB)| = \left| \sum_i a_{ii} b_{ii} \right| \leq \sum_i |a_{ii}| |b_{ii}| \leq \sum_i |a_{ii}| \|B\| = \operatorname{tr}(|A|) \|B\|$$

Plug in the inequality and we have

$$\begin{aligned} &\frac{d(\lambda_j^{\tilde{V}} + \eta)^2}{dt} + \frac{d(\lambda_j^{\tilde{W}} - 1)^2}{dt} \\ &\leq - \left[\frac{1}{\eta(2 - \eta)} + O\left(\frac{1}{\log^{\frac{1}{6}} d}\right) \right] (\lambda_j^{\tilde{V}} + \eta)^2 - \left[\eta^2(k + 1) + O\left(k^{\frac{11}{12}}\right) \right] (1 - \lambda_j^{\tilde{W}})^2 \\ &\quad + \alpha_1 |\lambda_j^{\tilde{V}} + \eta| \operatorname{tr}(|I - \Lambda \tilde{W}|) + \alpha_2 |\lambda_j^{\tilde{V}} + \eta| \operatorname{tr}(|I - \Lambda \tilde{W}|) \cdot \|\Lambda \tilde{W}\| \\ &\quad + \alpha_3 |1 - \lambda_j^{\tilde{W}}| \operatorname{tr}(|I - \Lambda \tilde{W}|) + \alpha_4 |1 - \lambda_j^{\tilde{W}}| \operatorname{tr}(|\Lambda \tilde{V} + \eta I|) \cdot \|\Lambda \tilde{V}\| \\ &\quad + \alpha_5 |1 - \lambda_j^{\tilde{W}}| \operatorname{tr}(|I - \Lambda \tilde{W}|) \cdot \|\Lambda \tilde{V}^2\| + \alpha_6 \\ &= - \left[\frac{1}{\eta(2 - \eta)} + O\left(\frac{1}{\log^{\frac{1}{6}} d}\right) \right] (\lambda_j^{\tilde{V}} + \eta)^2 - \left[\eta^2(k + 1) + O\left(k^{\frac{11}{12}}\right) \right] (1 - \lambda_j^{\tilde{W}})^2 \\ &\quad + (\alpha_1 + \alpha_2 \|\Lambda \tilde{W}\|) \cdot |\lambda_j^{\tilde{V}} + \eta| \cdot \operatorname{tr}(|I - \Lambda \tilde{W}|) + \alpha_4 \|\Lambda \tilde{V}\| \cdot |1 - \lambda_j^{\tilde{W}}| \cdot \operatorname{tr}(|\Lambda \tilde{V} + \eta I|) \\ &\quad + (\alpha_3 + \alpha_5 \|\Lambda \tilde{V}^2\|) \cdot |1 - \lambda_j^{\tilde{W}}| \cdot \operatorname{tr}(|I - \Lambda \tilde{W}|) + \alpha_6 \end{aligned}$$

Take the sum of both sides separately, we have

$$\begin{aligned} &\frac{d \operatorname{tr}[(\Lambda \tilde{V} + \eta I)^2]}{dt} + \frac{d \operatorname{tr}[(I - \Lambda \tilde{W})^2]}{dt} \\ &\leq - \left[\frac{1}{\eta(2 - \eta)} + O\left(\frac{1}{\log^{\frac{1}{6}} d}\right) \right] \operatorname{tr}[(\Lambda \tilde{V} + \eta I)^2] - \left[\eta^2(k + 1) + O\left(k^{\frac{11}{12}}\right) \right] \operatorname{tr}[(I - \Lambda \tilde{W})^2] \\ &\quad + (\alpha_1 + \alpha_2 \|\Lambda \tilde{W}\| + \alpha_4 \|\Lambda \tilde{V}\|) \cdot \operatorname{tr}(|\Lambda \tilde{V} + \eta I|) \cdot \operatorname{tr}(|I - \Lambda \tilde{W}|) \\ &\quad + (\alpha_3 + \alpha_5 \|\Lambda \tilde{V}^2\|) \cdot \operatorname{tr}^2(|I - \Lambda \tilde{W}|) + \alpha_6 \end{aligned}$$

From Jensen's Inequality with $f(x) = x^2$, we have

$$\left(\frac{\sum_{i=1}^d \lambda_i}{d} \right)^2 \leq \frac{\sum_{i=1}^d \lambda_i^2}{d}.$$

Therefore, it holds for diagonal matrix $\mathbf{\Lambda} \in \mathbb{R}^{d \times d}$ that

$$\text{tr}^2(\mathbf{\Lambda}) \leq d \text{tr}(\mathbf{\Lambda}^2)$$

Plug in the inequality and we have

$$\begin{aligned} & \frac{d \text{tr} \left[\left(\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I} \right)^2 \right]}{dt} + \frac{d \text{tr} \left[\left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right)^2 \right]}{dt} \\ & \leq -\frac{1}{2\eta(2-\eta)} \text{tr} \left[\left(\mathbf{\Lambda} \tilde{\mathbf{V}} + \eta \mathbf{I} \right)^2 \right] - \frac{\eta^2}{2} (k+1) \text{tr} \left[\left(\mathbf{I} - \mathbf{\Lambda} \tilde{\mathbf{W}} \right)^2 \right] + \alpha_6 \end{aligned}$$

Because $k = \lceil c \log d \rceil$, we have $O(d(1-\eta)^k) = d^{-c \log(\frac{1}{1-\eta})+1}$. So in $O(\log \frac{1}{\epsilon})$ time, $|\lambda_j^{\tilde{\mathbf{V}}} + \eta|$ and $|1 - \lambda_j^{\tilde{\mathbf{W}}}|$ converge to $\epsilon \in \left(\Theta\left(d^{\frac{\epsilon}{2} \log(1-\eta)+\frac{1}{2}}\right), 1 \right)$. \square

Lemma C.13. Suppose $\delta \in \left(\Theta\left(d^{\frac{\epsilon}{2} \log(1-\eta)+\frac{1}{2}}\right), 1 \right)$ and there exist diagonal matrices \mathbf{A} and \mathbf{B} satisfying $\|\mathbf{A}\|_{op} \leq \Theta(1)$ and $\|\mathbf{B}\|_{op} \leq \Theta(1)$ such that

$$\mathbf{\Lambda} \tilde{\mathbf{V}} = -\eta \mathbf{I} + \delta \cdot \mathbf{A} \quad \mathbf{\Lambda} \tilde{\mathbf{W}} = \mathbf{I} + \delta \cdot \mathbf{B},$$

then it holds that

$$\mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}) = O(\delta^2 d \log d).$$

Proof. Now we consider the CoT loss given by Lemma C.3

$$\begin{aligned} \mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \sum_{i=0}^{k-1} \left\| (\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) \mathbf{w}_i - (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} \mathbf{w}^* \right\|_2^2 \\ &\quad + \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left\| (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}}) \mathbf{w}_k - (\tilde{\mathbf{V}} \mathbf{S} + \mathbf{I}) \mathbf{w}^* \right\|_2^2. \end{aligned}$$

Plug in the expression of $\mathbf{\Lambda} \tilde{\mathbf{V}}$ and $\mathbf{\Lambda} \tilde{\mathbf{W}}$, we get

$$\mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}) = \frac{\delta^2}{2} \mathbb{E} \sum_{i=0}^{k-1} \left\| (\mathbf{A} \mathbf{S} - \eta \mathbf{S} \mathbf{B} + \delta \mathbf{A} \mathbf{S} \mathbf{B}) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) - \mathbf{A} \mathbf{S} \right\|_F^2 \quad (16)$$

$$+ \frac{1}{2} \mathbb{E} \left\| -(\mathbf{I} - \eta \mathbf{S})^k + \mathbf{\Lambda} \tilde{\mathbf{V}} \mathbf{S} \left[-(\mathbf{I} - \eta \mathbf{S})^k + \delta \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \right] \right\|_F^2. \quad (17)$$

We first consider the term in the summation:

$$\begin{aligned} & \mathbb{E} \left\| (\mathbf{A} \mathbf{S} - \eta \mathbf{S} \mathbf{B} + \delta \mathbf{A} \mathbf{S} \mathbf{B}) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) - \mathbf{A} \mathbf{S} \right\|_F^2 \\ &= \mathbb{E} \left\| (-\eta \mathbf{S} \mathbf{B} + \delta \mathbf{A} \mathbf{S} \mathbf{B}) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) - \mathbf{A} \mathbf{S} (\mathbf{I} - \eta \mathbf{S})^i \right\|_F^2 \\ &= \text{tr} \mathbb{E} \left[(-\eta \mathbf{S} \mathbf{B} + \delta \mathbf{A} \mathbf{S} \mathbf{B}) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 (-\eta \mathbf{B} \mathbf{S} + \delta \mathbf{B} \mathbf{S} \mathbf{A}) \right] \\ &\quad - 2 \text{tr} \mathbb{E} \left[(-\eta \mathbf{S} \mathbf{B} + \delta \mathbf{A} \mathbf{S} \mathbf{B}) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) (\mathbf{I} - \eta \mathbf{S})^i \mathbf{S} \mathbf{A} \right] + \text{tr} \mathbb{E} \left[\mathbf{A} \mathbf{S} (\mathbf{I} - \eta \mathbf{S})^{2i} \mathbf{S} \mathbf{A} \right] \\ &= \text{tr} \left((-\eta \mathbf{I} + \delta \mathbf{A}) \mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 \mathbf{B} \mathbf{S} \right] (-\eta \mathbf{I} + \delta \mathbf{A}) \right) \quad (\text{Term 1}) \\ &\quad - 2 \text{tr} \left((-\eta \mathbf{I} + \delta \mathbf{A}) \mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) (\mathbf{I} - \eta \mathbf{S})^i \mathbf{S} \right] \mathbf{A} \right) \quad (\text{Term 2}) \\ &\quad + \text{tr} \left(\mathbf{A} \mathbb{E} \left[\mathbf{S} (\mathbf{I} - \eta \mathbf{S})^{2i} \mathbf{S} \right] \mathbf{A} \right) \quad (\text{Term 3}) \end{aligned}$$

Apply Lemma C.15 to the expectation in Term 1, we have

$$\mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 \mathbf{B} \mathbf{S} \right]$$

$$= \left(1 - (1 - \eta)^i\right)^2 \mathbf{B}^2 + O\left(\frac{1}{\log^3 d}\right) \left[\mathbf{B}^2 + O\left(\frac{1}{d}\right) \text{tr}(\mathbf{B})\mathbf{B} + O\left(\frac{1}{d}\right) \text{tr}(\mathbf{B}^2)\mathbf{I} + O\left(\frac{1}{d^2}\right) \text{tr}^2(\mathbf{B})\mathbf{I} \right].$$

It is obvious that

$$\left\| \mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 \mathbf{B} \mathbf{S} \right] \right\|_{op} \leq \Theta(1).$$

Therefore, for Term 1 we have

$$\begin{aligned} & \text{tr} \left((-\eta \mathbf{I} + \delta \mathbf{A}) \mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 \mathbf{B} \mathbf{S} \right] (-\eta \mathbf{I} + \delta \mathbf{A}) \right) \\ & \leq d \| -\eta \mathbf{I} + \delta \mathbf{A} \|_{op}^2 \cdot \left\| \mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right)^2 \mathbf{B} \mathbf{S} \right] \right\|_{op} \leq O(d). \end{aligned}$$

(all matrices in the inequality are diagonal matrices.)

Similarly, for Term 2 and Term 3, we have

$$\begin{aligned} & \left| \text{tr} \left((-\eta \mathbf{I} + \delta \mathbf{A}) \mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) (\mathbf{I} - \eta \mathbf{S})^i \mathbf{S} \right] \mathbf{A} \right) \right| \leq O(d) \\ & \text{tr} \left(\mathbf{A} \mathbb{E} \left[\mathbf{S} (\mathbf{I} - \eta \mathbf{S})^{2i} \mathbf{S} \right] \mathbf{A} \right) \leq O(d). \end{aligned}$$

Add Term 1, 2, 3 together and we have

$$\mathbb{E} \left\| (\mathbf{A} \mathbf{S} - \eta \mathbf{S} \mathbf{B} + \delta \mathbf{A} \mathbf{S} \mathbf{B}) \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i \right) - \mathbf{A} \mathbf{S} \right\|_F^2 \leq O(d).$$

We then consider the second term in Equation (17):

$$\begin{aligned} & \mathbb{E} \left\| -(\mathbf{I} - \eta \mathbf{S})^k + \mathbf{\Lambda} \tilde{\mathbf{V}} \mathbf{S} \left[-(\mathbf{I} - \eta \mathbf{S})^k + \delta \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \right] \right\|_F^2 \\ & = \mathbb{E} \left\| -\left(\mathbf{I} + \mathbf{\Lambda} \tilde{\mathbf{V}} \mathbf{S} \right) (\mathbf{I} - \eta \mathbf{S})^k + \delta \mathbf{\Lambda} \tilde{\mathbf{V}} \mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \right\|_F^2 \\ & = \text{tr} \left(\mathbb{E} \left[\left(\mathbf{I} + \mathbf{\Lambda} \tilde{\mathbf{V}} \mathbf{S} \right) (\mathbf{I} - \eta \mathbf{S})^{2k} \left(\mathbf{I} + \mathbf{S} \mathbf{\Lambda} \tilde{\mathbf{V}} \right) \right] \right) \\ & \quad - 2\delta \text{tr} \left(\mathbb{E} \left[\left(\mathbf{I} + \mathbf{\Lambda} \tilde{\mathbf{V}} \mathbf{S} \right) (\mathbf{I} - \eta \mathbf{S})^k \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \mathbf{B} \mathbf{S} \right] \mathbf{\Lambda} \tilde{\mathbf{V}} \right) \\ & \quad + \delta^2 \text{tr} \left(\mathbf{\Lambda} \tilde{\mathbf{V}} \mathbb{E} \left[\mathbf{S} \mathbf{B} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right)^2 \mathbf{B} \mathbf{S} \right] \mathbf{\Lambda} \tilde{\mathbf{V}} \right) \\ & \leq O(\delta^2 d). \end{aligned}$$

((1 - \eta)^k \leq \delta)

Recall the CoT loss in Equation (16) and Equation (17). By the analysis above, we directly obtain that

$$\mathcal{L}^{\text{CoT}}(\boldsymbol{\theta}) = O(\delta^2 d \log d).$$

Hence, the proof is complete. \square

Lemma C.14. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$, $\|\mathbf{\Gamma}\|_{op} \leq \Theta(1)$. Then the expectation

$$\mathbb{E} \left[\mathbf{S} \mathbf{\Lambda} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \mathbf{\Gamma} \mathbf{S} \right] = \left(1 - (1 - \eta)^k \right) \mathbf{\Lambda} \mathbf{\Gamma} + \Delta,$$

where $\|\Delta\|_{op} = O\left(\frac{k^2 d}{n}\right) \leq O\left(\frac{1}{\log^3 d}\right)$. Moreover, the error is in the form

$$\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} + \alpha_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}$$

where $\alpha_1 = O\left(\frac{k^2 d}{n}\right)$, $\alpha_2, \alpha_3, \alpha_5 = O\left(\frac{k^2}{n}\right)$, $\alpha_4 = O\left(\frac{k^2}{nd}\right)$.

Proof. We can directly get the lemma by applying Lemma D.2 to $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{\Gamma} \mathbf{S}]$, $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} (\mathbf{I} - \eta \mathbf{S})^k \mathbf{\Gamma} \mathbf{S}]$. \square

Lemma C.15. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$, $\|\mathbf{\Gamma}\|_{op} \leq \Theta(1)$. Then the expectation

$$\mathbb{E} \left[\mathbf{S} \mathbf{\Lambda} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right)^2 \mathbf{\Gamma} \mathbf{S} \right] = \left(1 - (1 - \eta)^k \right)^2 \mathbf{\Lambda} \mathbf{\Gamma} + \Delta,$$

where $\|\Delta\|_{op} = O\left(\frac{k^2 d}{n}\right) \leq O\left(\frac{1}{\log^3 d}\right)$. Moreover, the error is in the form

$$\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} + \alpha_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}$$

where $\alpha_1 = O\left(\frac{k^2 d}{n}\right)$, $\alpha_2, \alpha_3, \alpha_5 = O\left(\frac{k^2}{n}\right)$, $\alpha_4 = O\left(\frac{k^2}{nd}\right)$.

Proof. We can directly get the lemma by applying Lemma D.2 to $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{\Gamma} \mathbf{S}]$, $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} (\mathbf{I} - \eta \mathbf{S})^k \mathbf{\Gamma} \mathbf{S}]$ and $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} (\mathbf{I} - \eta \mathbf{S})^{2k} \mathbf{\Gamma} \mathbf{S}]$. \square

Lemma C.16. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$, $\|\mathbf{\Gamma}\|_{op} \leq \Theta(1)$. Then the expectation

$$\mathbb{E} \left[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right) \right] = \left(1 - (1 - \eta)^k \right) \mathbf{\Lambda} \mathbf{\Gamma} + \Delta,$$

where $\|\Delta\|_{op} = O\left(\frac{k^2 d}{n}\right) \leq O\left(\frac{1}{\log^3 d}\right)$. Moreover, the error is in the form

$$\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} + \alpha_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}$$

where $\alpha_1 = O\left(\frac{k^2 d}{n}\right)$, $\alpha_2, \alpha_3, \alpha_5 = O\left(\frac{k^2}{n}\right)$, $\alpha_4 = O\left(\frac{k^2}{nd}\right)$.

Proof. We can directly get the lemma by applying Lemma D.3 to $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma}]$, $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} (\mathbf{I} - \eta \mathbf{S})^k]$. \square

Lemma C.17. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$, $\|\mathbf{\Gamma}\|_{op} \leq \Theta(1)$. Then the expectation

$$\mathbb{E} \left[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} \left(\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^k \right)^2 \right] = \left(1 - (1 - \eta)^k \right)^2 \mathbf{\Lambda} \mathbf{\Gamma} + \Delta,$$

where $\|\Delta\|_{op} = O\left(\frac{k^2 d}{n}\right) \leq O\left(\frac{1}{\log^3 d}\right)$. Moreover, the error is in the form

$$\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} + \alpha_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}$$

where $\alpha_1 = O\left(\frac{k^2 d}{n}\right)$, $\alpha_2, \alpha_3, \alpha_5 = O\left(\frac{k^2}{n}\right)$, $\alpha_4 = O\left(\frac{k^2}{nd}\right)$.

Proof. We can directly get the lemma by applying Lemma D.3 to $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma}]$, $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} (\mathbf{I} - \eta \mathbf{S})^k]$ and $\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} (\mathbf{I} - \eta \mathbf{S})^{2k}]$. \square

C.4 OUT-OF-DISTRIBUTION GENERALIZATION

We restate the formal theorem here. We still denote $\mathbf{S} := \frac{1}{n} \mathbf{X} \mathbf{X}^\top$ for simplicity. Note that the number of steps k can be different/larger compared to the step number in the previous training theorem.

Theorem C.2. Suppose $n = \Theta(d \log^5 d)$, $\eta \in (0.1, 0.9)$, $k = C \log d$. Assume the out-of-distribution input data $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}_d, \mathbf{\Sigma})$, $i \in [n]$ where $\frac{\delta}{\eta} \leq \lambda_{\min}(\mathbf{\Sigma}) \leq \lambda_{\max}(\mathbf{\Sigma}) \leq \frac{2-\delta}{\eta}$ for some constant $\delta > 0.1$, and $\mathbf{w}^* \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I})$. Then the trained transformer in Theorem 4.1 satisfies that $\mathcal{L}_{\mathbf{\Sigma}}^{\text{Eval}}(t) \leq \epsilon$ for any $\epsilon \in \left(d^{-C \log(\min\{\frac{1}{1-\eta}, \frac{1}{1-\delta}\})+1} \log^2 d, 1\right)$.

Proof. Recall the definition of the evaluation loss and our reduced transformer (Definition C.1)

$$\begin{aligned}\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) &= \frac{1}{2} \mathbb{E}_{\mathbf{X}, \mathbf{w}^*} \left[\left\| f_{\text{LSA}}(\hat{\mathbf{Z}}_k)_{[:, -1]} - (\mathbf{0}_d, 0, \mathbf{w}^*, 1) \right\|^2 \right] \\ &= \frac{1}{2} \mathbb{E} \left[\left\| f_{\theta}(\hat{\mathbf{w}}_k) - \mathbf{w}^* \right\|^2 \right]\end{aligned}$$

where $\hat{\mathbf{Z}}_k$ is the generated sequence after k steps and $\hat{\mathbf{w}}_k := f_{\theta}(\hat{\mathbf{w}}_{k-1})$ is the k -th generated intermediate weight vector. Note that each step the transformer is inputted with the last step prediction. We define the prediction error at each step i is $\Delta \mathbf{w}_i := \hat{\mathbf{w}}_i - \mathbf{w}_i = f_{\theta}(\hat{\mathbf{w}}_{i-1}) - \mathbf{w}_i$. We expand the term $f_{\theta}(\hat{\mathbf{w}}_k) - \mathbf{w}^*$ and sum up the error accumulation as follows:

$$\begin{aligned}f_{\theta}(\hat{\mathbf{w}}_k) - \mathbf{w}^* &\leq (\mathbf{w}_{k+1} - \mathbf{w}^*) + (f_{\theta}(\hat{\mathbf{w}}_k) - \mathbf{w}_{k+1}) \\ &= (\mathbf{w}_{k+1} - \mathbf{w}^*) + \hat{\mathbf{w}}_k + \tilde{\mathbf{V}}\mathbf{S}(\tilde{\mathbf{W}}\hat{\mathbf{w}}_k - \mathbf{w}^*) - \mathbf{w}_{k+1} \\ &\leq (\mathbf{w}_{k+1} - \mathbf{w}^*) + \left(\mathbf{w}_k + \tilde{\mathbf{V}}\mathbf{S}(\tilde{\mathbf{W}}\mathbf{w}_k - \mathbf{w}^*) - \mathbf{w}_{k+1} \right) + \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right) \Delta \mathbf{w}_k.\end{aligned}$$

After one step of decomposition, we notice that the error $\Delta \mathbf{w}_{k+1}$ can be decomposed into two parts: (1) The approximation error predicting \mathbf{w}_{k+1} with ground-truth input \mathbf{w}_k . We define it

$$\Delta_{k+1}^{\text{pred}} := \mathbf{w}_k + \tilde{\mathbf{V}}\mathbf{S}(\tilde{\mathbf{W}}\mathbf{w}_k - \mathbf{w}^*) - \mathbf{w}_{k+1}$$

(2) The accumulated error from the last inference step: $(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}})\Delta \mathbf{w}_k$. Therefore, we can inductively calculate the sum of the error:

$$\begin{aligned}f_{\theta}(\hat{\mathbf{w}}_k) - \mathbf{w}^* &\leq (\mathbf{w}_{k+1} - \mathbf{w}^*) + \left(\mathbf{w}_k + \tilde{\mathbf{V}}\mathbf{S}(\tilde{\mathbf{W}}\mathbf{w}_k - \mathbf{w}^*) - \mathbf{w}_{k+1} \right) + \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right) \Delta \mathbf{w}_k \\ &= (\mathbf{w}_{k+1} - \mathbf{w}^*) + \Delta_{k+1}^{\text{pred}} + \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right) \Delta \mathbf{w}_k \\ &= (\mathbf{w}_{k+1} - \mathbf{w}^*) + \Delta_{k+1}^{\text{pred}} + \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right) \Delta_k^{\text{pred}} + \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right)^2 \Delta \mathbf{w}_{k-1} \\ &= (\mathbf{w}_{k+1} - \mathbf{w}^*) + \sum_{i=0}^k \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right)^i \Delta_{k-i+1}^{\text{pred}} \quad (\Delta \mathbf{w}_0 = 0 \text{ by definition.})\end{aligned}$$

Then we have our evaluation loss upper bounded:

$$\begin{aligned}\frac{1}{2} \mathbb{E} \left[\left\| f_{\theta}(\hat{\mathbf{w}}_k) - \mathbf{w}^* \right\|^2 \right] &= \frac{1}{2} \mathbb{E} \left\| (\mathbf{w}_{k+1} - \mathbf{w}^*) + \sum_{i=0}^k \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right)^i \Delta_{k-i+1}^{\text{pred}} \right\|^2 \\ &\leq \frac{k+2}{2} \left(\mathbb{E} \left\| (\mathbf{w}_{k+1} - \mathbf{w}^*) \right\|^2 + \sum_{i=0}^k \mathbb{E} \left\| \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right)^i \Delta_{k-i+1}^{\text{pred}} \right\|^2 \right) \quad (*)\end{aligned}$$

We first consider the first term: $\mathbb{E} \left\| (\mathbf{w}_{k+1} - \mathbf{w}^*) \right\|^2$:

$$\begin{aligned}\mathbb{E} \left\| \mathbf{w}_{k+1} - \mathbf{w}^* \right\|^2 &= \mathbb{E} \left\| (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k+1}) \mathbf{w}^* - \mathbf{w}^* \right\|^2 = \text{tr} \left(\mathbb{E} (\mathbf{I} - \eta \mathbf{S})^{2k+2} \right) \\ &\leq 2d(1 - \delta)^{2k+2} \leq 2d^{-2c \log(\frac{1}{1-\delta})+1}. \quad (\text{Lemma D.5})\end{aligned}$$

Then we consider the second summation term. Since the parameters of the reduced model $\tilde{\mathbf{V}} = -\eta \mathbf{I} + \mathbf{A}$, $\tilde{\mathbf{W}} = \mathbf{I} + \mathbf{B}$, where $\|\mathbf{A}\|_{\text{op}}, \|\mathbf{B}\|_{\text{op}} \leq d^{-\frac{1}{2}C \log(\frac{1}{1-\eta}) + \frac{1}{2}}$ for some constant $c > 0$, we want to bound the prediction error given the ground-truth input. By Lemma D.6, we have

$$\begin{aligned}&\mathbb{E} \sum_{i=0}^k \left\| \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right)^i \Delta_{k-i+1}^{\text{pred}} \right\|^2 \\ &= \mathbb{E} \sum_{i=0}^k \left\| \left(\mathbf{I} + \tilde{\mathbf{V}}\mathbf{S}\tilde{\mathbf{W}} \right)^i (\mathbf{w}_{k-i} + \tilde{\mathbf{V}}\mathbf{S}(\tilde{\mathbf{W}}\mathbf{w}_{k-i} - \mathbf{w}^*) - \mathbf{w}_{k-i+1}) \right\|^2 \\ &\leq O \left(d^{-C \log(\frac{1}{1-\eta})+1} \cdot k \right).\end{aligned}$$

Therefore, plug those back to Equation (*), the total evaluation loss should be upper bounded by

$$\mathcal{L}^{\text{Eval}}(\theta) \leq O \left(d^{-C \log(\min\{\frac{1}{1-\eta}, \frac{1}{1-\delta}\})+1} \cdot k^2 \right) = O(d^{-C \log(\min\{\frac{1}{1-\eta}, \frac{1}{1-\delta}\})+1} \log^2 d)$$

□

D SUPPLEMENTARY LEMMAS

D.1 CONCENTRATION LEMMAS

In this appendix, we prove some concentration lemmas to estimate the expected gradient more accurately. Throughout the proof, $\mathbf{\Lambda}, \mathbf{\Gamma}$ are both symmetric matrices with orthonormal eigenbasis $\{\mathbf{u}_i\}_{i=1}^d$.

Lemma D.1. *Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$. Then the expectation*

$$\mathbb{E}[\mathbf{S}\mathbf{\Lambda}(\mathbf{I} - \eta\mathbf{S})^k \mathbf{S}] = (1 - \eta)^k (\mathbf{\Lambda} + \Delta),$$

where $\|\Delta\|_{op} \leq O(\frac{k^2 d}{n}) = O(\frac{1}{\log^3 d})$. Moreover, the error is in the form $\Delta = \alpha_1 \mathbf{\Lambda} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{I}$, where $\alpha_1 = O(\frac{k^2 d}{n})$, $\alpha_2 = O(\frac{k^2}{n})$.

Proof. Denote $\delta\mathbf{S} := \mathbf{S} - \mathbf{I}$. Then we expand the term $\mathbf{S}\mathbf{\Lambda}(\mathbf{I} - \eta\mathbf{S})^k \mathbf{S}$:

$$\begin{aligned} & \mathbf{S}\mathbf{\Lambda}(\mathbf{I} - \eta\mathbf{S})^k \mathbf{S} \\ &= (\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda}((1 - \eta)\mathbf{I} - \eta\delta\mathbf{S})^k (\mathbf{I} + \delta\mathbf{S}) \\ &= (1 - \eta)^k (\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda} \left(\mathbf{I} - \frac{\eta}{(1 - \eta)} \delta\mathbf{S} \right)^k (\mathbf{I} + \delta\mathbf{S}) \\ &= (1 - \eta)^k (\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda} \left(\mathbf{I} - \frac{k\eta}{(1 - \eta)} \delta\mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \\ & \quad + (1 - \eta)^k (\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda} \left(\mathbf{I} - \frac{k\eta}{(1 - \eta)} \delta\mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \delta\mathbf{S} \end{aligned}$$

Now take expectation to both sides. Note that $\mathbb{E}[\delta\mathbf{S}] = 0$, so all the terms only contain first order $\delta\mathbf{S}$ vanish. We denote

$$(1 - \eta)^k \tilde{\Delta} = \mathbf{S}\mathbf{\Lambda}(\mathbf{I} - \eta\mathbf{S})^k \mathbf{S} - (1 - \eta)^k \left(\mathbf{\Lambda} + \delta\mathbf{S}\mathbf{\Lambda} + \mathbf{\Lambda}\delta\mathbf{S} - \frac{k\eta}{1 - \eta} \mathbf{\Lambda}\delta\mathbf{S} \right),$$

which denotes all the higher order terms (the degree of $\delta\mathbf{S} \geq 2$.)

Since we have the tail bound for $\delta\mathbf{S}$ in Theorem 4.6.1 Vershynin (2018) (In this lemma $\|\cdot\|$ is operator norm if without specification):

$$\Pr(\|\delta\mathbf{S}\| > \max(\delta, \delta^2)) \leq 2 \exp\{-s^2\}, \text{ where } \delta = C \left(\sqrt{\frac{d}{n}} + \frac{s}{\sqrt{n}} \right) \quad (18)$$

We can estimate the expectation using this property. First, given $s = \sqrt{d}$ and $\|\delta\mathbf{S}\| \leq \max(\delta, \delta^2) = C\sqrt{\frac{d}{n}}$ (since $n = \Theta(d \log^5 d)$), we can upper bound the operator norm of $\tilde{\Delta}$:

$$\begin{aligned} \|\tilde{\Delta}\|_{op} &\leq \left\| \mathbf{\Lambda} \left(\binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \right\|_{op} \\ & \quad + \left\| \delta\mathbf{S}\mathbf{\Lambda} \left(-\frac{k\eta}{(1 - \eta)} \delta\mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \right\|_{op} \\ & \quad + \left\| \delta\mathbf{S}\mathbf{\Lambda} \left(\mathbf{I} - \frac{k\eta}{(1 - \eta)} \delta\mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \delta\mathbf{S} \right\|_{op} \end{aligned}$$

$$+ \left\| \Lambda \left(-\frac{k\eta}{(1-\eta)} \delta \mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1-\eta} \right)^2 \delta \mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1-\eta} \right)^j \delta \mathbf{S}^j \right) \delta \mathbf{S} \right\|_{op}$$

Now upper bound all matrices with their operator norm and combine all terms with the same degree of $\delta \mathbf{S}$. We have

$$\begin{aligned} \|\tilde{\Delta}\|_{op} &\leq \sum_{j=2}^{k+2} \|\Lambda\| \left(\binom{k}{j} \left(\frac{\eta}{1-\eta} \right)^j + 2 \binom{k}{j-1} \left(\frac{\eta}{1-\eta} \right)^{j-1} + \binom{k}{j-2} \left(\frac{\eta}{1-\eta} \right)^{j-2} \right) \|\delta \mathbf{S}\|^j \\ &\leq \sum_{j=2}^{k+2} \|\Lambda\| ((9k)^j + 2(9k)^{j-1} + (9k)^{j-2}) \|\delta \mathbf{S}\|^j \quad \left(\frac{\eta}{1-\eta} \leq 9, \binom{k}{j} \leq k^j \right) \\ &\leq 4 \sum_{j=2}^{k+2} \|\Lambda\| \cdot (9k)^j \left(C \sqrt{\frac{d}{n}} \right)^j \quad (\|\delta \mathbf{S}\| \leq C \sqrt{\frac{d}{n}}) \\ &\leq 4 \|\Lambda\| \cdot \frac{81C^2 k^2 d}{n} \cdot \frac{1}{1 - (9kd^{1/2}/n^{1/2})} \leq C' \frac{k^2 d}{n} \leq O\left(\frac{1}{\log^3 d}\right). \quad (*) \end{aligned}$$

Given this upper bound, we can now upper bound the operator norm of the error term $\Delta := \mathbb{E}[\tilde{\Delta}]$.

Suppose $\mathbf{u} := \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \frac{\|\Delta \mathbf{u}\|}{\|\mathbf{u}\|}$, then the operator norm becomes:

$$\begin{aligned} \|\Delta\| &= \left| \mathbf{u}^\top \mathbb{E}[\tilde{\Delta}] \mathbf{u} \right| \\ &= \mathbb{E} \left[\left| \mathbf{u}^\top \tilde{\Delta} \mathbf{u} \right| \left(\mathbb{1} \left\{ \|\tilde{\Delta}\| \leq C' \frac{k^2 d}{n} \right\} + \mathbb{1} \left\{ \|\tilde{\Delta}\| > C' \frac{k^2 d}{n} \right\} \right) \right] \\ &\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds \end{aligned}$$

When $\|\tilde{\Delta}\| \geq s$ where $s \geq \frac{C' k^2 d}{n}$, we can first upper bound the $\|\tilde{\Delta}\|$ with $\|\delta \mathbf{S}\|$ using the second row of eq. (*1): there exists some constant $C_1 > 0$ s.t.

$$\|\tilde{\Delta}\| \leq 4 \sum_{j=2}^{k+2} \|\Lambda\| \cdot (9k \|\delta \mathbf{S}\|)^j \leq \max \left((C_1 k \|\delta \mathbf{S}\|)^2, (C_1 k \|\delta \mathbf{S}\|)^{k+2} \right).$$

Therefore, when $\|\tilde{\Delta}\| \geq s$, $\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\}$. To apply the tail bound, we need to make sure we pick some s' such that $\max(\delta, \delta^2) \leq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\}$ to upper bound the integral of probability, where $\delta = C \left(\sqrt{\frac{d}{n}} + \frac{s'}{\sqrt{n}} \right)$. Now since $s > \frac{C' k^2 d}{n}$, $\min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \geq C_\alpha \sqrt{\frac{d}{n}}$ for some constant C_α . Therefore, we just need $\max \left\{ \frac{s'}{\sqrt{n}}, \frac{s'^2}{n} \right\} \leq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\}$, i.e. $s' \leq \min \left\{ C_2 \frac{s^{1/(k+2)} \sqrt{n}}{k}, C_3 \frac{s^{1/(2k+4)} \sqrt{n}}{\sqrt{k}}, C_4 \frac{\sqrt{sn}}{k}, C_5 \frac{s^{1/4} \sqrt{n}}{\sqrt{k}} \right\}$.

Applying the tail bound (18) with $s' = \min \left\{ C_2 \frac{s^{1/(k+2)} \sqrt{n}}{k}, C_3 \frac{s^{1/(2k+4)} \sqrt{n}}{\sqrt{k}}, C_4 \frac{\sqrt{sn}}{k}, C_5 \frac{s^{1/4} \sqrt{n}}{\sqrt{k}} \right\}$ where C_2, C_3, C_4, C_5 are some constant, we have the error term for the tail expectation,

$$\begin{aligned} \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds &\leq \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds \\ &\leq 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp \{-s'^2\} ds. \end{aligned}$$

Now we estimate the upper bound of error with

$$s'^2 = \min \left\{ C_2^2 \cdot \frac{s^{2/(k+2)}}{k^2} n, C_3^2 \cdot \frac{s^{1/(k+2)}}{k} n, C_4^2 \cdot \frac{sn}{k^2}, C_5^2 \cdot \frac{\sqrt{sn}}{k} \right\}.$$

For the first term, let $x = \frac{C_2^2 n}{k^2} s^{2/(k+2)}$:

$$\begin{aligned} & 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp\left\{-C_2^2 \cdot \frac{s^{2/(k+2)}}{k^2} n\right\} ds \\ &= (k+2) \int_{\frac{C_2^2 n}{k^2} \left(\frac{C' k^2 d}{n}\right)^{\frac{2}{k+2}}}^{\infty} \left(\frac{k^2}{C_2^2 n}\right)^{(k+2)/2} \exp\{-x\} x^{k/2} dx \\ &\leq (k+2) \cdot \left(\frac{k^2}{C_2^2 n}\right)^{(k+2)/2} \cdot \left(\frac{C_2^2 n}{k^2} \left(\frac{C' k^2 d}{n}\right)^{\frac{2}{k+2}}\right)^{k/2} \exp\left\{-\frac{C_2^2 n}{k^2} \left(\frac{C' k^2 d}{n}\right)^{\frac{2}{k+2}}\right\} \leq \frac{k^2 d}{n}. \end{aligned}$$

The second term, let $x = C_3^2 \cdot \frac{s^{1/(k+2)}}{k} n$:

$$\begin{aligned} & 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp\left\{-C_3^2 \cdot \frac{s^{1/(k+2)}}{k} n\right\} ds \\ &= 2(k+2) \int_{\frac{C_3^2 n}{k} \left(\frac{C' k^2 d}{n}\right)^{\frac{1}{k+2}}}^{\infty} \left(\frac{k}{C_3^2 n}\right)^{k+2} \exp\{-x\} x^{k+1} dx \\ &\leq 2(k+2) \cdot \left(\frac{k}{C_3^2 n}\right)^{k+2} \cdot \left(\frac{C_3^2 n}{k} \left(\frac{C' k^2 d}{n}\right)^{\frac{1}{k+2}}\right)^{k+1} \exp\left\{-\frac{C_3^2 n}{k} \left(\frac{C' k^2 d}{n}\right)^{\frac{1}{k+2}}\right\} \leq \frac{k^2 d}{n}. \end{aligned}$$

For the third term, let $x = \frac{C_4^2 s n}{k^2}$:

$$\begin{aligned} & 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp\left\{-\frac{C_4^2 s n}{k^2}\right\} ds = \int_{\frac{C' k^2 d}{n} \cdot \frac{C_4^2 n}{k^2}}^{\infty} \frac{k^2}{C_4^2 n} \exp\{-x\} dx \\ &\leq \frac{k^2}{C_4^2 n} \exp\left\{-\frac{C' k^2 d}{n} \cdot \frac{C_4^2 n}{k^2}\right\} \leq \frac{k^2 d}{n}. \end{aligned}$$

The fourth term, let $x = C_5^2 \cdot \frac{s^{1/2}}{k} n$:

$$\begin{aligned} & 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp\left\{-C_5^2 \cdot \frac{s^{1/2}}{k} n\right\} ds \\ &= \frac{4k^2}{n^2 C_5^4} \int_{C_5^2 \frac{n}{k} \left(\frac{C' k^2 d}{n}\right)^{1/2}}^{\infty} \exp\{-x\} x dx \\ &\leq \frac{4k^2}{n^2 C_5^4} \cdot C_5^2 \frac{n}{k} \left(\frac{C' k^2 d}{n}\right)^{1/2} \exp\left\{-C_5^2 \frac{n}{k} \left(\frac{C' k^2 d}{n}\right)^{1/2}\right\} \leq \frac{k^2 d}{n}. \end{aligned}$$

Therefore, we plug this error back to the upper bound of $\|\Delta\|$:

$$\begin{aligned} \|\Delta\| &\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr\left[\|\tilde{\Delta}\| \geq s\right] ds \\ &\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr\left[\|\delta \mathbf{S}\| \geq \min\left\{\frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k}\right\}\right] ds = O\left(\frac{k^2 d}{n}\right) \leq O\left(\frac{1}{\log^3 d}\right). \end{aligned}$$

Finally, since by Lemma D.8, we know the error is in the form $\Delta = \alpha_1 \mathbf{\Lambda} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{I}$ for all $\mathbf{\Lambda}$.

Therefore $\alpha_1 = O\left(\frac{k^2 d}{n}\right)$, $\alpha_2 = O\left(\frac{k^2}{n}\right)$. \square

Lemma D.2. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$, $\|\mathbf{\Gamma}\|_{op} \leq \Theta(1)$. Then the expectation

$$\mathbb{E}[\mathbf{S} \mathbf{\Lambda} (\mathbf{I} - \eta \mathbf{S})^k \mathbf{\Gamma} \mathbf{S}] = (1 - \eta)^k (\mathbf{\Lambda} \mathbf{\Gamma} + \Delta),$$

where $\|\Delta\|_{op} = O\left(\frac{k^2 d}{n}\right) \leq O\left(\frac{1}{\log^3 d}\right)$. Moreover, the error is in the form

$$\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} + \alpha_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}$$

where $\alpha_1 = O\left(\frac{k^2 d}{n}\right)$, $\alpha_2, \alpha_3, \alpha_5 = O\left(\frac{k^2}{n}\right)$, $\alpha_4 = O\left(\frac{k^2}{nd}\right)$.

Proof. Denote $\delta \mathbf{S} := \mathbf{S} - \mathbf{I}$. Then we expand the term $\mathbf{S}\mathbf{\Lambda}(\mathbf{I} - \eta\mathbf{S})^k\mathbf{\Gamma}\mathbf{S}$:

$$\begin{aligned}
& \mathbf{S}\mathbf{\Lambda}(\mathbf{I} - \eta\mathbf{S})^k\mathbf{\Gamma}\mathbf{S} \\
&= (\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda}((1 - \eta)\mathbf{I} - \eta\delta\mathbf{S})^k\mathbf{\Gamma}(\mathbf{I} + \delta\mathbf{S}) \\
&= (1 - \eta)^k(\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda}\left(\mathbf{I} - \frac{\eta}{(1 - \eta)}\delta\mathbf{S}\right)^k\mathbf{\Gamma}(\mathbf{I} + \delta\mathbf{S}) \\
&= (1 - \eta)^k(\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda}\left(\mathbf{I} - \frac{k\eta}{(1 - \eta)}\delta\mathbf{S} + \binom{k}{2}\left(\frac{\eta}{1 - \eta}\right)^2\delta\mathbf{S}^2 + \sum_{j=3}^k\binom{k}{j}\left(\frac{-\eta}{1 - \eta}\right)^j\delta\mathbf{S}^j\right)\mathbf{\Gamma} \\
&\quad + (1 - \eta)^k(\mathbf{I} + \delta\mathbf{S})\mathbf{\Lambda}\left(\mathbf{I} - \frac{k\eta}{(1 - \eta)}\delta\mathbf{S} + \binom{k}{2}\left(\frac{\eta}{1 - \eta}\right)^2\delta\mathbf{S}^2 + \sum_{j=3}^k\binom{k}{j}\left(\frac{-\eta}{1 - \eta}\right)^j\delta\mathbf{S}^j\right)\mathbf{\Gamma}\delta\mathbf{S}
\end{aligned}$$

Take expectation to both sides. Note that $\mathbb{E}[\delta\mathbf{S}] = 0$, so all the first order term vanish. We denote

$$(1 - \eta)^k\tilde{\Delta} = \mathbf{S}\mathbf{\Lambda}(\mathbf{I} - \eta\mathbf{S})^k\mathbf{\Gamma}\mathbf{S} - (1 - \eta)^k\left(\mathbf{\Lambda} + \delta\mathbf{S} \cdot \mathbf{\Lambda}\mathbf{\Gamma} + \mathbf{\Lambda}\mathbf{\Gamma} \cdot \delta\mathbf{S} - \frac{k\eta}{1 - \eta}\mathbf{\Lambda} \cdot \delta\mathbf{S}\mathbf{\Gamma}\right),$$

which denotes all the higher order terms (the degree of $\delta\mathbf{S} \geq 2$.)

We can estimate the expectation using similar technique as in Lemma D.1. First, given $s = \sqrt{d}$ and $\|\delta\mathbf{S}\| \leq \max(\delta, \delta^2) = C\sqrt{\frac{d}{n}}$ (since $n = \Theta(d \log^5 d)$), we upper bound the operator norm of $\tilde{\Delta}$:

$$\begin{aligned}
\|\tilde{\Delta}\|_{op} &\leq \left\| \mathbf{\Lambda} \left(\binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \mathbf{\Gamma} \right\|_{op} \\
&\quad + \left\| \delta\mathbf{S} \mathbf{\Lambda} \left(-\frac{k\eta}{(1 - \eta)} \delta\mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \mathbf{\Gamma} \right\|_{op} \\
&\quad + \left\| \delta\mathbf{S} \mathbf{\Lambda} \left(\mathbf{I} - \frac{k\eta}{(1 - \eta)} \delta\mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \mathbf{\Gamma} \delta\mathbf{S} \right\|_{op} \\
&\quad + \left\| \mathbf{\Lambda} \left(-\frac{k\eta}{(1 - \eta)} \delta\mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta\mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta\mathbf{S}^j \right) \mathbf{\Gamma} \delta\mathbf{S} \right\|_{op}
\end{aligned}$$

Now upper bound all matrices with their operator norm and combine all terms with the same degree of $\delta\mathbf{S}$. We have

$$\begin{aligned}
\|\tilde{\Delta}\|_{op} &\leq \sum_{j=2}^{k+2} \|\mathbf{\Gamma}\| \|\mathbf{\Lambda}\| \left(\binom{k}{j} \left(\frac{\eta}{1 - \eta} \right)^j + 2 \binom{k}{j-1} \left(\frac{\eta}{1 - \eta} \right)^{j-1} + \binom{k}{j-2} \left(\frac{\eta}{1 - \eta} \right)^{j-2} \right) \|\delta\mathbf{S}\|^j \\
&\leq \sum_{j=2}^{k+2} \|\mathbf{\Gamma}\| \|\mathbf{\Lambda}\| ((9k)^j + 2(9k)^{j-1} + (9k)^{j-2}) \|\delta\mathbf{S}\|^j \quad \left(\frac{\eta}{1 - \eta} \leq 9, \binom{k}{j} \leq k^j \right) \\
&\leq 4 \sum_{j=2}^{k+2} \|\mathbf{\Gamma}\| \|\mathbf{\Lambda}\| \cdot (9k)^j \left(C\sqrt{\frac{d}{n}} \right)^j \quad (\|\delta\mathbf{S}\| \leq C\sqrt{\frac{d}{n}}) \\
&\leq 4 \|\mathbf{\Lambda}\| \|\mathbf{\Gamma}\| \cdot \frac{81C^2k^2d}{n} \cdot \frac{1}{1 - \left(\frac{9kd^{1/2}}{n^{1/2}} \right)} \leq C' \frac{k^2d}{n} \leq O\left(\frac{1}{\log^3 d} \right)
\end{aligned}$$

Now upper bound the operator norm of the error term $\Delta := \mathbb{E}[\tilde{\Delta}]$. Suppose $\mathbf{u} := \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \frac{\|\Delta\mathbf{u}\|}{\|\mathbf{u}\|}$, then the operator norm becomes:

$$\|\Delta\| = \left| \mathbf{u}^\top \mathbb{E}[\tilde{\Delta}] \mathbf{u} \right|$$

$$\begin{aligned}
&= \mathbb{E} \left[\left| \mathbf{u}^\top \tilde{\Delta} \mathbf{u} \right| \left(\mathbb{1} \left\{ \|\tilde{\Delta}\| \leq C' \frac{k^2 d}{n} \right\} + \mathbb{1} \left\{ \|\tilde{\Delta}\| > C' \frac{k^2 d}{n} \right\} \right) \right] \\
&\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds
\end{aligned}$$

When $\|\tilde{\Delta}\| \geq s$ where $s \geq \frac{C' k^2 d}{n}$, there exists some constant $C_1 > 0$ s.t.

$$\|\tilde{\Delta}\| \leq \max \left((C_1 k \|\delta \mathbf{S}\|)^2, (C_1 k \|\delta \mathbf{S}\|)^{k+2} \right).$$

Therefore, when $\|\tilde{\Delta}\| \geq s$, $\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\}$. Like Lemma D.1, applying the tail bound (18) with $s' \leq \min \left\{ C_2 \frac{s^{1/(k+2)} \sqrt{n}}{k}, C_3 \frac{s^{1/(2k+4)} \sqrt{n}}{\sqrt{k}}, C_4 \frac{\sqrt{s n}}{k}, C_5 \frac{s^{1/4} \sqrt{n}}{\sqrt{k}} \right\}$ where C_2, C_3, C_4, C_5 are some constant, we have the error term for the tail expectation

$$\begin{aligned}
\int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds &\leq \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds \\
&\leq 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp \{-s'^2\} ds.
\end{aligned}$$

Use the exact same argument, $2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp \{-s'^2\} ds \leq \frac{k^2 d}{n}$. Thus, the upper bound of $\|\Delta\|$ is:

$$\begin{aligned}
\|\Delta\| &\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds \\
&\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds = O \left(\frac{k^2 d}{n} \right) \leq O \left(\frac{1}{\log^3 d} \right).
\end{aligned}$$

Finally by Lemma D.8, we know the error is in the form $\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I}$ for all $\mathbf{\Lambda}, \mathbf{\Gamma}$. Therefore $\alpha_1 = O \left(\frac{k^2 d}{n} \right), \alpha_2, \alpha_3 = O \left(\frac{k^2}{n} \right), \alpha_4 = O \left(\frac{k^2}{nd} \right)$. \square

Lemma D.3. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$, $\|\mathbf{\Gamma}\|_{op} \leq \Theta(1)$. Then the expectation

$$\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} (\mathbf{I} - \eta \mathbf{S})^k] = (1 - \eta)^k (\mathbf{\Lambda} \mathbf{\Gamma} + \Delta),$$

where $\|\Delta\|_{op} \leq O \left(\frac{1}{\log^3 d} \right)$. Moreover, the error is in the form

$$\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} + \alpha_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}$$

where $\alpha_1 = O \left(\frac{k^2 d}{n} \right), \alpha_2, \alpha_3, \alpha_5 = O \left(\frac{k^2}{n} \right), \alpha_4 = O \left(\frac{k^2}{nd} \right)$.

Proof. Denote $\delta \mathbf{S} := \mathbf{S} - \mathbf{I}$. Then we expand the term $\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} (\mathbf{I} - \eta \mathbf{S})^k$:

$$\begin{aligned}
&\mathbf{S} \mathbf{\Lambda} \mathbf{S} \mathbf{\Gamma} (\mathbf{I} - \eta \mathbf{S})^k \\
&= (1 - \eta)^k (\mathbf{I} + \delta \mathbf{S}) \mathbf{\Lambda} (\mathbf{I} + \delta \mathbf{S}) \mathbf{\Gamma} \left(\mathbf{I} - \frac{\eta}{(1 - \eta)} \delta \mathbf{S} \right)^k \\
&= (1 - \eta)^k (\mathbf{I} + \delta \mathbf{S}) \mathbf{\Lambda} \mathbf{\Gamma} \left(\mathbf{I} - \frac{k \eta}{(1 - \eta)} \delta \mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta \mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta \mathbf{S}^j \right) \\
&\quad + (1 - \eta)^k (\mathbf{I} + \delta \mathbf{S}) \mathbf{\Lambda} \delta \mathbf{S} \mathbf{\Gamma} \left(\mathbf{I} - \frac{k \eta}{(1 - \eta)} \delta \mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta \mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta \mathbf{S}^j \right)
\end{aligned}$$

Take expectation to both sides. Note that $\mathbb{E}[\delta \mathbf{S}] = 0$, so all the first order term vanish. We denote

$$(1 - \eta)^k \tilde{\Delta} = \mathbf{S} \mathbf{\Lambda} (\mathbf{I} - \eta \mathbf{S})^k \mathbf{\Gamma} \mathbf{S} - (1 - \eta)^k \left(\mathbf{\Lambda} + \delta \mathbf{S} \cdot \mathbf{\Lambda} \mathbf{\Gamma} + \mathbf{\Lambda} \delta \mathbf{S} \mathbf{\Gamma} - \frac{k \eta}{1 - \eta} \mathbf{\Lambda} \mathbf{\Gamma} \cdot \delta \mathbf{S} \right),$$

which denotes all the higher order terms (the degree of $\delta \mathbf{S} \geq 2$.)

We can estimate the expectation using similar technique as in Lemma D.1. Given $s = \sqrt{d}$ and $\|\delta \mathbf{S}\| \leq \max(\delta, \delta^2) = C\sqrt{\frac{d}{n}}$ (since $n = \Theta(d \log^5 d)$), we upper bound the operator norm of $\tilde{\Delta}$. We directly expand the formula and upper bound all matrices with their operator norm and combine all terms with the same degree of $\delta \mathbf{S}$. We have

$$\begin{aligned} \|\tilde{\Delta}\|_{op} &\leq \sum_{j=2}^{k+2} \|\mathbf{\Gamma}\| \|\mathbf{\Lambda}\| \left(\binom{k}{j} \left(\frac{\eta}{1-\eta} \right)^j + 2 \binom{k}{j-1} \left(\frac{\eta}{1-\eta} \right)^{j-1} + \binom{k}{j-2} \left(\frac{\eta}{1-\eta} \right)^{j-2} \right) \|\delta \mathbf{S}\|^j \\ &\leq \sum_{j=2}^{k+2} \|\mathbf{\Gamma}\| \|\mathbf{\Lambda}\| ((9k)^j + 2(9k)^{j-1} + (9k)^{j-2}) \|\delta \mathbf{S}\|^j \quad \left(\frac{\eta}{1-\eta} \leq 9, \binom{k}{j} \leq k^j \right) \\ &\leq 4 \sum_{j=2}^{k+2} \|\mathbf{\Gamma}\| \|\mathbf{\Lambda}\| \cdot (9k)^j \left(C\sqrt{\frac{d}{n}} \right)^j \leq C' \frac{k^2 d}{n} \leq O\left(\frac{1}{\log^3 d} \right) \quad (\|\delta \mathbf{S}\| \leq C\sqrt{\frac{d}{n}}) \end{aligned}$$

Now upper bound the operator norm of $\Delta := \mathbb{E}[\tilde{\Delta}]$. Suppose $\mathbf{u} := \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \frac{\|\Delta \mathbf{u}\|}{\|\mathbf{u}\|}$, then

$$\begin{aligned} \|\Delta\| &= \mathbb{E} \left[\left| \mathbf{u}^\top \tilde{\Delta} \mathbf{u} \right| \left(\mathbb{1} \left\{ \|\tilde{\Delta}\| \leq C' \frac{k^2 d}{n} \right\} + \mathbb{1} \left\{ \|\tilde{\Delta}\| > C' \frac{k^2 d}{n} \right\} \right) \right] \\ &\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds \end{aligned}$$

When $\|\tilde{\Delta}\| \geq s$ where $s \geq \frac{C' k^2 d}{n}$, there exists some constant $C_1 > 0$ s.t.

$$\|\tilde{\Delta}\| \leq \max \left((C_1 k \|\delta \mathbf{S}\|)^2, (C_1 k \|\delta \mathbf{S}\|)^{k+2} \right).$$

Therefore, when $\|\tilde{\Delta}\| \geq s$, $\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\}$. Like Lemma D.1, applying the tail bound (18) with $s' \leq \min \left\{ C_2 \frac{s^{1/(k+2)} \sqrt{n}}{k}, C_3 \frac{s^{1/(2k+4)} \sqrt{n}}{\sqrt{k}}, C_4 \frac{\sqrt{sn}}{k}, C_5 \frac{s^{1/4} \sqrt{n}}{\sqrt{k}} \right\}$ where C_2, C_3, C_4, C_5 are some constant, we have

$$\int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds \leq \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds \leq 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp \{-s'^2\} ds.$$

Use the exact same argument, $2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp \{-s'^2\} ds \leq \frac{k^2 d}{n}$. Thus, the upper bound of $\|\Delta\|$ is:

$$\|\Delta\| \leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds = O\left(\frac{k^2 d}{n} \right) \leq O\left(\frac{1}{\log^3 d} \right).$$

Finally by Lemma D.8, we know the error is in the form $\Delta = \alpha_1 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I}$ for all $\mathbf{\Lambda}, \mathbf{\Gamma}$. Therefore $\alpha_1 = O\left(\frac{k^2 d}{n} \right)$, $\alpha_2, \alpha_3 = O\left(\frac{k^2}{n} \right)$, $\alpha_4 = O\left(\frac{k^2}{nd} \right)$. \square

Lemma D.4. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$, $\|\mathbf{\Lambda}\|_{op} \leq \Theta(1)$. Then there exists $\delta = O\left(\frac{k^2 d}{n} \right) \leq O\left(\frac{1}{\log^3 d} \right)$, the expectation is

$$\mathbb{E}[\mathbf{\Lambda}(\mathbf{I} - \eta \mathbf{S})^k] = (1 - \eta)^k (1 + \delta) \mathbf{\Lambda},$$

Proof. Denote $\delta \mathbf{S} := \mathbf{S} - \mathbf{I}$. Then we expand the term $\mathbf{\Lambda}(\mathbf{I} - \eta \mathbf{S})^k$:

$$\begin{aligned} \mathbf{\Lambda}(\mathbf{I} - \eta \mathbf{S})^k &= (1 - \eta)^k \mathbf{\Lambda} \left(\mathbf{I} - \frac{\eta}{(1 - \eta)} \delta \mathbf{S} \right)^k \\ &= (1 - \eta)^k \mathbf{\Lambda} \left(\mathbf{I} - \frac{k\eta}{(1 - \eta)} \delta \mathbf{S} + \binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta \mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta \mathbf{S}^j \right) \end{aligned}$$

Take expectation to both sides. Note that $\mathbb{E}[\delta \mathbf{S}] = 0$, so all the first order term vanish. We denote

$$(1 - \eta)^k \tilde{\Delta} = \mathbf{\Lambda}(\mathbf{I} - \eta \mathbf{S})^k - (1 - \eta)^k \left(\mathbf{\Lambda} - \frac{k\eta}{1 - \eta} \mathbf{\Lambda} \cdot \delta \mathbf{S} \right),$$

which denotes all the higher order terms (the degree of $\delta \mathbf{S} \geq 2$.)

We can estimate the expectation using similar technique as in Lemma D.1. First, given $s = \sqrt{d}$ and $\|\delta \mathbf{S}\| \leq \max(\delta, \delta^2) = C\sqrt{\frac{d}{n}}$ (since $n = \Theta(d \log^5 d)$), we upper bound the operator norm of $\tilde{\Delta}$:

$$\|\tilde{\Delta}\|_{op} \leq \left\| \mathbf{\Lambda} \left(\binom{k}{2} \left(\frac{\eta}{1 - \eta} \right)^2 \delta \mathbf{S}^2 + \sum_{j=3}^k \binom{k}{j} \left(\frac{-\eta}{1 - \eta} \right)^j \delta \mathbf{S}^j \right) \right\|_{op}$$

Now upper bound all matrices by operator norm and combine all terms with the same degree of $\delta \mathbf{S}$:

$$\begin{aligned} \|\tilde{\Delta}\|_{op} &\leq \sum_{j=2}^k \|\mathbf{\Lambda}\| \left(\binom{k}{j} \left(\frac{\eta}{1 - \eta} \right)^j \right) \|\delta \mathbf{S}\|^j \leq \sum_{j=2}^{k+2} \|\mathbf{\Lambda}\| (9k)^j \|\delta \mathbf{S}\|^j \quad \left(\frac{\eta}{1 - \eta} \leq 9, \binom{k}{j} \leq k^j \right) \\ &\leq \|\mathbf{\Lambda}\| \cdot \frac{81C^2k^2d}{n} \cdot \frac{1}{1 - \left(\frac{9kd^{1/2}}{n^{1/2}} \right)} \leq C' \frac{k^2d}{n} \leq O\left(\frac{1}{\log^3 d} \right) \quad (\|\delta \mathbf{S}\| \leq C\sqrt{\frac{d}{n}}) \end{aligned}$$

Now upper bound the operator norm of the error. Suppose $\mathbf{u} := \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \frac{\|\Delta \mathbf{u}\|}{\|\mathbf{u}\|}$, we have

$$\begin{aligned} \|\Delta\| &= \left| \mathbf{u}^\top \mathbb{E}[\tilde{\Delta}] \mathbf{u} \right| = \mathbb{E} \left[\left| \mathbf{u}^\top \tilde{\Delta} \mathbf{u} \right| \left(\mathbb{1} \left\{ \|\tilde{\Delta}\| \leq C' \frac{k^2d}{n} \right\} + \mathbb{1} \left\{ \|\tilde{\Delta}\| > C' \frac{k^2d}{n} \right\} \right) \right] \\ &\leq C' \frac{k^2d}{n} + \int_{\frac{C'k^2d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds \end{aligned}$$

When $\|\tilde{\Delta}\| \geq s$ where $s \geq \frac{C'k^2d}{n}$, there exists some constant $C_1 > 0$ s.t.

$$\|\tilde{\Delta}\| \leq \max \left((C_1 k \|\delta \mathbf{S}\|)^2, (C_1 k \|\delta \mathbf{S}\|)^{k+2} \right).$$

Therefore, when $\|\tilde{\Delta}\| \geq s$, $\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\}$. Like Lemma D.1, applying the tail bound (18) with $s' \leq \min \left\{ C_2 \frac{s^{1/(k+2)} \sqrt{n}}{k}, C_3 \frac{s^{1/(2k+4)} \sqrt{n}}{\sqrt{k}}, C_4 \frac{\sqrt{sn}}{k}, C_5 \frac{s^{1/4} \sqrt{n}}{\sqrt{k}} \right\}$ where C_2, C_3, C_4, C_5 are some constant, we have the error term for the tail expectation

$$\int_{\frac{C'k^2d}{n}}^{\infty} \Pr \left[\|\tilde{\Delta}\| \geq s \right] ds \leq \int_{\frac{C'k^2d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds \leq 2 \int_{\frac{C'k^2d}{n}}^{\infty} \exp \{-s'^2\} ds.$$

Use the exact same argument, $2 \int_{\frac{C'k^2d}{n}}^{\infty} \exp \{-s'^2\} ds \leq \frac{k^2d}{n}$. Thus, the upper bound of $\|\Delta\|$ is:

$$\|\Delta\| \leq C' \frac{k^2d}{n} + \int_{\frac{C'k^2d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds = O\left(\frac{k^2d}{n} \right) \leq O\left(\frac{1}{\log^3 d} \right).$$

Finally by Lemma D.8, we know the error is in the form $\Delta = \alpha_1 \mathbf{\Lambda}$ for all $\mathbf{\Lambda}$. So $\alpha_1 = O\left(\frac{k^2d}{n} \right)$. \square

D.2 CONCENTRATION LEMMAS FOR OUT-OF-DISTRIBUTION DATA

For non-isotropic covariance Gaussian data input, we also have the concentration around the covariance Σ when $n = \Theta(d \log^c d)$ for $c > 0$. We still denote $\mathbf{S} = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$. The following lemmas are involved in the calculation for the evaluation process, for in-distribution and out-of-distribution input examples \mathbf{X} .

Lemma D.5. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ where $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}_d, \Sigma)$, $\frac{\delta}{\eta} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq \frac{2-\delta}{\eta}$ for some constant $\delta > 0.1$, $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$. Then the expectation

$$\text{tr}(\mathbb{E}(\mathbf{I} - \eta \mathbf{S})^k) \leq 2d(1 - \delta)^k.$$

Proof. Denote $\delta \mathbf{S} := \mathbf{S} - \Sigma$. Then we expand the term $\Lambda(\mathbf{I} - \eta \mathbf{S})^k$:

$$\begin{aligned} (\mathbf{I} - \eta \mathbf{S})^k &= (1 - \delta)^k \left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} - \frac{\eta}{1 - \delta} \delta \mathbf{S} \right)^k \\ &= (1 - \delta)^k \left(\left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right)^k - \frac{k\eta}{(1 - \delta)} \left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right)^{k-1} \delta \mathbf{S} + \sum_{j=2}^k \binom{k}{j} \left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right)^{k-j} \left(\frac{-\eta}{1 - \delta} \right)^j \delta \mathbf{S}^j \right) \end{aligned}$$

Take expectation to both sides. Note that $\mathbb{E}[\delta \mathbf{S}] = 0$, so all the first order term vanish. We denote

$$(1 - \delta)^k \tilde{\Delta} = (\mathbf{I} - \eta \mathbf{S})^k - (1 - \delta)^k \left(\left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right)^k - \frac{k\eta}{(1 - \delta)} \left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right)^{k-1} \delta \mathbf{S} \right),$$

which denotes all the higher order terms (the degree of $\delta \mathbf{S} \geq 2$). Note $\left\| \frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right\|_{op} \leq 1$.

We can estimate the expectation using similar technique as in Lemma D.1. First, given $s = \sqrt{d}$ and $\|\delta \mathbf{S}\| \leq \max(\delta, \delta^2) = C\sqrt{\frac{d}{n}}$ (since $n = \Theta(d \log^5 d)$), we upper bound the operator norm of $\tilde{\Delta}$:

$$\|\tilde{\Delta}\|_{op} \leq \left\| \sum_{j=2}^k \binom{k}{j} \left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right)^{k-j} \left(\frac{-\eta}{1 - \delta} \right)^j \delta \mathbf{S}^j \right\|_{op}$$

Now upper bound all matrices by operator norm and combine all terms with the same degree of $\delta \mathbf{S}$:

$$\begin{aligned} \|\tilde{\Delta}\|_{op} &\leq \sum_{j=2}^k \binom{k}{j} \left(\frac{\eta}{1 - \delta} \right)^j \|\delta \mathbf{S}\|^j \leq \sum_{j=2}^{k+2} (9k)^j \|\delta \mathbf{S}\|^j \quad \left(\frac{\eta}{1 - \delta} \leq 9, \binom{k}{j} \leq k^j \right) \\ &\leq \frac{81C^2 k^2 d}{n} \cdot \frac{1}{1 - \left(\frac{9kd^{1/2}}{n^{1/2}} \right)} \leq C' \frac{k^2 d}{n} \leq O\left(\frac{1}{\log^3 d} \right) \quad (\|\delta \mathbf{S}\| \leq C\sqrt{\frac{d}{n}}) \end{aligned}$$

Now upper bound the operator norm of the error. Suppose $\mathbf{u} := \arg \max_{\mathbf{u}: \|\mathbf{u}\|=1} \frac{\|\Delta \mathbf{u}\|}{\|\mathbf{u}\|}$, we have

$$\begin{aligned} \|\Delta\| &= |\mathbf{u}^\top \mathbb{E}[\tilde{\Delta}] \mathbf{u}| = \mathbb{E} \left[|\mathbf{u}^\top \tilde{\Delta} \mathbf{u}| \left(\mathbb{1} \left\{ \|\tilde{\Delta}\| \leq C' \frac{k^2 d}{n} \right\} + \mathbb{1} \left\{ \|\tilde{\Delta}\| > C' \frac{k^2 d}{n} \right\} \right) \right] \\ &\leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr[\|\tilde{\Delta}\| \geq s] ds \end{aligned}$$

When $\|\tilde{\Delta}\| \geq s$ where $s \geq \frac{C' k^2 d}{n}$, there exists some constant $C_1 > 0$ s.t.

$$\|\tilde{\Delta}\| \leq \max \left((C_1 k \|\delta \mathbf{S}\|)^2, (C_1 k \|\delta \mathbf{S}\|)^{k+2} \right).$$

Therefore, when $\|\tilde{\Delta}\| \geq s$, $\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\}$. Like Lemma D.1, applying the tail bound (18) with $s' \leq \min \left\{ C_2 \frac{s^{1/(k+2)} \sqrt{n}}{k}, C_3 \frac{s^{1/(2k+4)} \sqrt{n}}{\sqrt{k}}, C_4 \frac{\sqrt{sn}}{k}, C_5 \frac{s^{1/4} \sqrt{n}}{\sqrt{k}} \right\}$ where C_2, C_3, C_4, C_5 are some constant, we have the error term for the tail expectation

$$\int_{\frac{C' k^2 d}{n}}^{\infty} \Pr[\|\tilde{\Delta}\| \geq s] ds \leq \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds \leq 2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp\{-s'^2\} ds.$$

Use the exact same argument, $2 \int_{\frac{C' k^2 d}{n}}^{\infty} \exp\{-s'^2\} ds \leq \frac{k^2 d}{n}$. Thus, the upper bound of $\|\Delta\|$ is:

$$\|\Delta\| \leq C' \frac{k^2 d}{n} + \int_{\frac{C' k^2 d}{n}}^{\infty} \Pr \left[\|\delta \mathbf{S}\| \geq \min \left\{ \frac{s^{1/2}}{C_1 k}, \frac{s^{1/(k+2)}}{C_1 k} \right\} \right] ds = O\left(\frac{k^2 d}{n} \right) \leq O\left(\frac{1}{\log^3 d} \right) < \frac{1}{2}.$$

Finally, the absolute value of the trace should be upper bounded by

$$\text{tr} \left((1 - \delta)^k \left(\left(\frac{\mathbf{I} - \eta \Sigma}{1 - \delta} \right)^k + \Delta \right) \right) \leq 2d(1 - \delta)^k.$$

□

The next lemma deals with the prediction error.

Lemma D.6. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ where $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}_d, \Sigma)$, and the covariance matrix satisfies $\frac{\delta}{\eta} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq \frac{2-\delta}{\eta}$ for some constant $\delta > 0$. Assume $n = \Theta(d \log^5 d)$, $k = O(\log d)$, $\eta = \Theta(1) \in (0.1, 0.9)$. Denote that $\mathbf{A} := \tilde{\mathbf{V}} + \eta \mathbf{I}$, $\mathbf{B} := \tilde{\mathbf{W}} - \mathbf{I}$, $\|\mathbf{A}\|_{op}, \|\mathbf{B}\|_{op} \leq \Theta(d^{-c})$. Then for any $i < k$,

$$\mathbb{E} \left\| (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}})^i (\mathbf{w}_{k-i} + \tilde{\mathbf{V}} \mathbf{S} (\tilde{\mathbf{W}} \mathbf{w}_{k-i} - \mathbf{w}^*) - \mathbf{w}_{k-i+1}) \right\|^2 \leq O\left(\frac{(1-\delta)^{2i}}{d^{-2c+1}}\right)$$

Proof. We will adopt a similar method as we did throughout Lemma D.1 to Lemma D.4.

First, we expand the left hand side loss:

$$\begin{aligned} & \mathbb{E} \left\| (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}})^i (\mathbf{w}_{k-i} + \tilde{\mathbf{V}} \mathbf{S} (\tilde{\mathbf{W}} \mathbf{w}_{k-i} - \mathbf{w}^*) - \mathbf{w}_{k-i+1}) \right\|^2 \\ &= \mathbb{E} \left\| (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}})^i \left((\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k-i}) \mathbf{w}^* - (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} \mathbf{w}^* \right) \right\|^2 \\ &= \mathbb{E} \left\| (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}})^i \left((\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k-i}) - (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} \right) \right\|_F^2 \\ &\leq d \cdot \mathbb{E} \left\| (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}})^i \left((\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k-i}) - (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} \right) \right\|_{op}^2 \end{aligned}$$

The second equation is due to $\mathbf{w}_i = (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^i) \mathbf{w}^*$, and we arranged to stress the error terms. The third line is because $\mathbb{E}[\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}$. The last line is $\|\cdot\|_F \leq \sqrt{d} \|\cdot\|_{op}$.

Now we expand each term of the expression within the operator norm into $\mathbf{A}, \mathbf{B}, \mathbf{I}$, and \mathbf{S} :

$$\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} = \mathbf{I} - \eta \mathbf{S} + \mathbf{A} \mathbf{B} - \eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{S}.$$

$$\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S} = \mathbf{A} \mathbf{B} - \eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{S}, \tilde{\mathbf{V}} + \eta \mathbf{I} = \mathbf{A}.$$

Therefore the formula becomes (consider each term separately)

$$\begin{aligned} & (\mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}})^i = (\mathbf{I} - \eta \mathbf{S} + \mathbf{A} \mathbf{B} - \eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{S})^i \\ & \left((\tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} + \eta \mathbf{S}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k-i}) - (\tilde{\mathbf{V}} + \eta \mathbf{I}) \mathbf{S} \right) \\ &= (-\eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{B}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k-i}) - \mathbf{A} \mathbf{S} (\mathbf{I} - \eta \mathbf{S})^{k-i} \end{aligned}$$

We still denote $\delta \mathbf{S} = \mathbf{S} - \Sigma$. We first consider when the concentration holds, a.k.a $\|\delta \mathbf{S}\| \leq C \sqrt{\frac{d}{n}}$.

Since $\|\mathbf{A}\|, \|\mathbf{B}\| \leq O(d^{-c})$, their error are dominated by $C \sqrt{\frac{d}{n}}$. We reduce this case to the previous Lemma D.5. Therefore we can upper bound the expression by

$$\begin{aligned} & \left\| \mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} \right\|^i = \left\| \mathbf{I} - \eta \mathbf{S} + \mathbf{A} \mathbf{B} - \eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{S} \right\|^i \leq \frac{3}{2} (1 - \delta)^i \\ & \left\| (-\eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{B}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k-i}) - \mathbf{A} \mathbf{S} (\mathbf{I} - \eta \mathbf{S})^{k-i} \right\| \leq O(d^{-c}). \end{aligned}$$

That means this part of the expectation is upper bounded by $d \cdot \frac{9}{4} (1 - \delta)^{2i} \cdot O(d^{-2c}) = O\left(\frac{(1-\delta)^{2i}}{d^{-2c+1}}\right)$

Then we estimate the tail expectation. We first upper bound the above formula by $\|\delta \mathbf{S}\|$:

$$\begin{aligned} & \left\| \mathbf{I} + \tilde{\mathbf{V}} \mathbf{S} \tilde{\mathbf{W}} \right\|^i = \left\| \mathbf{I} - \eta \mathbf{S} + \mathbf{A} \mathbf{B} - \eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{S} \right\|^i \leq O(k(1 - \delta)^i \min\{\|\delta \mathbf{S}\|, 1\}^i) \\ & \left\| (-\eta \mathbf{S} \mathbf{B} + \mathbf{A} \mathbf{B}) (\mathbf{I} - (\mathbf{I} - \eta \mathbf{S})^{k-i}) - \mathbf{A} \mathbf{S} (\mathbf{I} - \eta \mathbf{S})^{k-i} \right\| \leq O(kd^{-c} \min\{1, \|\delta \mathbf{S}\|^{k-i}\}). \end{aligned}$$

Use the same argument as in Lemma D.1 to calculate the integral of tail bound, the tail expectation can also be upper bounded by $O\left(\frac{(1-\delta)^{2i}}{d^{-2c+1}}\right)$. Combine those two part and we finish the proof.

□

D.3 THE FORM OF EXPECTATION

Lemma D.7. Suppose $S = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, then the expectation is in the following form for any k :

$$\begin{aligned}\mathbb{E}[\mathbf{S} \mathbf{u}_s \mathbf{u}_s^\top \mathbf{S}^k \mathbf{u}_t \mathbf{u}_t^\top \mathbf{S}^{k'}] &= \alpha_1 \mathbf{u}_s \mathbf{u}_s^\top + \alpha_2 \mathbf{u}_t \mathbf{u}_t^\top + \alpha_3 \mathbf{I} \text{ for any } s \neq t. \\ \mathbb{E}[\mathbf{S} \mathbf{u}_s \mathbf{u}_s^\top \mathbf{S}^k \mathbf{u}_s \mathbf{u}_s^\top \mathbf{S}^{k'}] &= \alpha_4 \mathbf{u}_s \mathbf{u}_s^\top + \alpha_5 \mathbf{I}.\end{aligned}$$

Proof. We notice that by changing the basis to $\{\mathbf{u}_s\}_{s=1}^d$,

$$\mathbb{E}[\mathbf{S} \mathbf{u}_s \mathbf{u}_s^\top \mathbf{S}^k \mathbf{u}_t \mathbf{u}_t^\top \mathbf{S}^{k'}] = \mathbf{U} \mathbb{E}[(\mathbf{U}^\top \mathbf{S} \mathbf{U}) \mathbf{e}_s \mathbf{e}_s^\top (\mathbf{U}^\top \mathbf{S} \mathbf{U})^k \mathbf{e}_t \mathbf{e}_t^\top (\mathbf{U}^\top \mathbf{S} \mathbf{U})^{k'}] \mathbf{U}^\top. \quad (19)$$

Define $\hat{\mathbf{x}}_i = \mathbf{U}^\top \mathbf{x}_i$. Since gaussian is isotropic, we have $\mathbb{E}[\hat{\mathbf{x}}_i] = \mathbf{U}^\top \mathbb{E}[\mathbf{x}_i] = \mathbf{0}$. After we change the basis, the covariance matrix of $\hat{\mathbf{x}}_i$ should also be the same:

$$\text{Cov}(\hat{\mathbf{x}}_i) = \mathbf{U}^\top \text{Cov}(\mathbf{x}_i) \mathbf{U} = \mathbf{I}.$$

Therefore $\hat{\mathbf{x}}_i$ has the same distribution as \mathbf{x}_i and we have

$$\mathbf{U} \mathbb{E}[(\mathbf{U}^\top \mathbf{S} \mathbf{U}) \mathbf{e}_s \mathbf{e}_s^\top (\mathbf{U}^\top \mathbf{S} \mathbf{U})^k \mathbf{e}_t \mathbf{e}_t^\top (\mathbf{U}^\top \mathbf{S} \mathbf{U})^{k'}] \mathbf{U}^\top = \mathbf{U} \mathbb{E}[\mathbf{S} \mathbf{e}_s \mathbf{e}_s^\top \mathbf{S}^k \mathbf{e}_t \mathbf{e}_t^\top \mathbf{S}^{k'}] \mathbf{U}^\top.$$

Subsequently, we only need to consider the expectation of $\mathbf{S} \mathbf{e}_s \mathbf{e}_s^\top \mathbf{S}^k \mathbf{e}_t \mathbf{e}_t^\top \mathbf{S}^{k'}$. Decompose \mathbf{x}_i into the sum of basis vectors and we get $\mathbf{x}_i = \sum_{j=1}^d x_{ij} \mathbf{e}_j$.

Plug in the decomposition into the expectation and we have

$$\begin{aligned}& n^{k+2} \mathbb{E}[\mathbf{S} \mathbf{e}_s \mathbf{e}_s^\top \mathbf{S}^k \mathbf{e}_t \mathbf{e}_t^\top \mathbf{S}^{k'}] \\&= \mathbb{E}\left[\left(\sum_{i_0=1}^n \sum_{j_0, j_1 \in [d]} x_{i_0 j_0} x_{i_0 j_1} \mathbf{e}_{j_0} \mathbf{e}_{j_1}^\top\right) \mathbf{e}_s \mathbf{e}_s^\top \prod_{l=1}^{k'} \left(\sum_{i_l=1}^n \sum_{j_{2l}, j_{2l+1} \in [d]} x_{i_l j_{2l}} x_{i_l j_{2l+1}} \mathbf{e}_{j_{2l}} \mathbf{e}_{j_{2l+1}}^\top\right)\right. \\&\quad \left. \mathbf{e}_t \mathbf{e}_t^\top \prod_{l=1}^k \left(\sum_{i_l=1}^n \sum_{j_{2l}, j_{2l+1} \in [d]} x_{i_l j_{2l}} x_{i_l j_{2l+1}} \mathbf{e}_{j_{2l}} \mathbf{e}_{j_{2l+1}}^\top\right)\right] \\&= \mathbb{E}\left[\sum_{i_0, \dots, i_{k+k'} \in [n]} \sum_{j_0, \dots, j_{2(k+k')+1} \in [d]} x_{i_0 j_0} x_{i_0 j_1} \cdots x_{i_{k+k'} j_{2(k+k')}} x_{i_{k+k'} j_{2(k+k')+1}} \right. \\&\quad \left. \mathbf{e}_{j_0} \mathbf{e}_{j_1}^\top \mathbf{e}_s \mathbf{e}_s^\top \mathbf{e}_{j_2} \mathbf{e}_{j_3}^\top \cdots \mathbf{e}_{j_{2k}} \mathbf{e}_{j_{2k+1}}^\top \mathbf{e}_t \mathbf{e}_t^\top \mathbf{e}_{j_{2k+2}} \mathbf{e}_{j_{2k+3}}^\top \cdots \mathbf{e}_{j_{2(k+k')}} \mathbf{e}_{j_{2(k+k')+1}}^\top\right] \\&= \sum_{i_0, \dots, i_{k+k'} \in [n]} \sum_{j_0, \dots, j_{2(k+k')+1} \in [d]} \mathbb{E}[x_{i_0 j_0} x_{i_0 j_1} \cdots x_{i_{k+k'} j_{2(k+k')}} x_{i_{k+k'} j_{2(k+k')+1}}] \\&\quad \mathbf{e}_{j_0} \mathbf{e}_{j_1}^\top \mathbf{e}_s \mathbf{e}_s^\top \mathbf{e}_{j_2} \mathbf{e}_{j_3}^\top \cdots \mathbf{e}_{j_{2k}} \mathbf{e}_{j_{2k+1}}^\top \mathbf{e}_t \mathbf{e}_t^\top \mathbf{e}_{j_{2k+2}} \mathbf{e}_{j_{2k+3}}^\top \cdots \mathbf{e}_{j_{2(k+k')}} \mathbf{e}_{j_{2(k+k')+1}}^\top.\end{aligned}$$

Note that $\mathbf{e}_a^\top \mathbf{e}_b \neq 0$ only when $a = b$, so $\mathbf{e}_a^\top \mathbf{e}_s \mathbf{e}_s^\top \mathbf{e}_b \neq 0$ only when $a = b = s$. Therefore, we only need to consider the case where $j_{2q-1} = j_{2q}$ for any $q \in [1, k+k']$. By symmetry, we know $\mathbb{E}[\mathbf{S} \mathbf{e}_s \mathbf{e}_s^\top \mathbf{S}^k \mathbf{e}_t \mathbf{e}_t^\top \mathbf{S}^{k'}]$ is a diagonal matrix, so we have $j_0 = j_{2(k+k')+1}$. We denote

$$\mathbf{E}_{j_0} = \mathbf{e}_{j_0} \mathbf{e}_{j_1}^\top \mathbf{e}_s \mathbf{e}_s^\top \mathbf{e}_{j_1} \mathbf{e}_{j_2}^\top \cdots \mathbf{e}_{j_k} \mathbf{e}_{j_{k+1}}^\top \mathbf{e}_t \mathbf{e}_t^\top \mathbf{e}_{j_{k+1}} \mathbf{e}_{j_{k+2}}^\top \cdots \mathbf{e}_{j_{k+k'}} \mathbf{e}_{j_0}^\top$$

to be one of the standard basis in $\mathbb{R}^{d \times d}$ space. It is a non-zero matrix when $j_1 = s$ and $j_{k+1} = t$. By the analysis above, we have

$$n^{k+2} \mathbb{E}[\mathbf{S} \mathbf{e}_s \mathbf{e}_s^\top \mathbf{S}^k \mathbf{e}_t \mathbf{e}_t^\top \mathbf{S}^{k'}] = \sum_{i_0, \dots, i_{k+k'} \in [n]} \sum_{j_0, \dots, j_{k+k'} \in [d]} \mathbb{E}[x_{i_0 j_0} x_{i_0 j_1} \cdots x_{i_{k+k'} j_{k+k'}} x_{i_{k+k'} j_0}] \mathbf{E}_{j_0}.$$

Let $\mathcal{P}(2k+2)$ be the set of all distinct ways of partitioning $\{i_0j_0, i_0j_1, \dots, i_{k+k'}j_{k+k'}, i_{k+k'}j_0\}$ into $k+1$ unordered pairs $p = ((p_1, p_2), \dots, (p_{2k+1}, p_{2k+2}))$. From Isserlis' theorem, we have

$$\mathbb{E}\left[x_{i_0j_0}x_{i_0j_1}\cdots x_{i_{k+k'}j_{k+k'}}x_{i_{k+k'}j_0}\right] = \sum_{p \in \mathcal{P}(2k+2)} \prod_{i=0}^{k+k'} \mathbb{E}[x_{p_{2i}}x_{p_{2i+1}}].$$

Plug it in the expectation and we have

$$\begin{aligned} n^{k+2} \mathbb{E}\left[\mathbf{S}e_s e_s^\top \mathbf{S}^k e_t e_t^\top \mathbf{S}^{k'}\right] &= \sum_{i_0, \dots, i_{k+k'} \in [n]} \sum_{j_0, \dots, j_{k+k'} \in [d]} \sum_{p \in \mathcal{P}(2k+2)} \prod_{i=0}^{k+k'} \mathbb{E}[x_{p_{2i}}x_{p_{2i+1}}] \mathbf{E}_{j_0} \\ &= \sum_{p \in \mathcal{P}(2k+2)} \sum_{i_0, \dots, i_k \in [n]} \sum_{j_0, \dots, j_k \in [d]} \prod_{i=0}^{k+k'} \mathbb{E}[x_{p_{2i}}x_{p_{2i+1}}] \mathbf{E}_{j_0}. \end{aligned}$$

To make sure the term in the summation is non-zero, $p_{2q-1} = p_{2q}$ should hold for any $1 \leq q \leq k+1$. Now consider the graph \mathcal{G}_p and \mathcal{G}'_p with vertices $\{0, 1, \dots, k+k'\}$. If $i_{u_1}j_{v_1}$ is paired with $i_{u_2}j_{v_2}$, then we put an edge between u_1 and u_2 into \mathcal{G}_p and put an edge between v_1 and v_2 into \mathcal{G}'_p , which means $i_{u_1} = i_{u_2}$ and $j_{v_1} = j_{v_2}$. Therefore, for a cycle $C = (u_1, u_2, \dots, u_r)$ in \mathcal{G}_p or \mathcal{G}'_p , we have $i_{u_1} = i_{u_2} = \dots = i_{u_r}$ or $j_{u_1} = j_{u_2} = \dots = j_{u_r}$. Note that we have n or d choices for the value of the circle. Here we use $C(\cdot)$ to denote the set of circles in the graph and use $|C(\cdot)|$ to denote the number of circles in the graph. Let c^* be the cycle in \mathcal{G}'_p which includes the vertex j_0 .

Case 1: $s \neq t$. For the partition p where $j_1 \in c^*$ and $j_{k+1} \in c \neq c^*$, there is only one choice for c and c^* to take. So the term in the summation should be $n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_s e_s^\top$. Similarly, for the partition p where $j_{k+1} \in c^*$ and $j_1 \in c \neq c^*$, the term in the summation should be $n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_t e_t^\top$. For the partition p where $j_1 \in c' \neq c^*$ and $j_{k+1} \in c'' \neq c^*$, there is only one choice for c' and c'' to take. Therefore, the expectation should be

$$\begin{aligned} &n^{k+2} \mathbb{E}\left[\mathbf{S}e_s e_s^\top \mathbf{S}^k e_t e_t^\top \mathbf{S}^{k'}\right] \\ &= \sum_{\mathcal{P}: j_1 \in c^*, j_{k+1} \notin c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_s e_s^\top + \sum_{\mathcal{P}: j_{k+1} \in c^*, j_1 \notin c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_t e_t^\top \\ &\quad + \sum_{\mathcal{P}: j_1, j_{k+1} \notin c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_{j_0} e_{j_0}^\top \\ &= \sum_{\mathcal{P}: j_1 \in c^*, j_{k+1} \notin c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_s e_s^\top + \sum_{\mathcal{P}: j_{k+1} \in c^*, j_1 \notin c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_t e_t^\top \\ &\quad + \sum_{\mathcal{P}: j_1, j_{k+1} \notin c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-3} \mathbf{I}. \end{aligned}$$

Recall Equation (19), we prove that

$$\mathbb{E}\left[\mathbf{S}u_s u_s^\top \mathbf{S}^k u_t u_t^\top \mathbf{S}^{k'}\right] = \alpha_1 u_s u_s^\top + \alpha_2 u_t u_t^\top + \alpha_3 \mathbf{I}.$$

Case 2: $s = t$. For the partition p where $j_1, j_{k+1} \in c^*$, there is only one choice for c^* to take. So the term in the summation should be $n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-1} e_s e_s^\top$. For the partition p where $j_1 \in c^*$ and $j_{k+1} \in c \neq c^*$, there is only one choice for c and c^* to take. So the term in the summation should be $n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_s e_s^\top$. Similarly, for the partition p where $j_{k+1} \in c^*$ and $j_1 \in c \neq c^*$, the term in the summation should be $n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_s e_s^\top$. For the partition p where $j_1 \in c' \neq c^*$ and $j_{k+1} \in c'' \neq c^*$, there is only one choice for c' and c'' to take. Therefore, the expectation should be

$$\begin{aligned} &n^{k+2} \mathbb{E}\left[\mathbf{S}e_s e_s^\top \mathbf{S}^k e_s e_s^\top \mathbf{S}^{k'}\right] \\ &= \sum_{\mathcal{P}: j_1, j_{k+1} \in c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-1} e_s e_s^\top + \sum_{\mathcal{P}: j_1 \in c^*, j_{k+1} \notin c^*} n^{|C(\mathcal{G}_p)|} d^{|C(\mathcal{G}'_p)|-2} e_s e_s^\top \end{aligned}$$

$$\begin{aligned}
& + \sum_{\mathcal{P}: j_{k+1} \in c^*, j_1 \notin c^*} n^{|\mathcal{C}(\mathcal{G}_p)|} d^{|\mathcal{C}(\mathcal{G}'_p)|-2} \mathbf{e}_s \mathbf{e}_s^\top + \sum_{\mathcal{P}: j_1, j_{k+1} \notin c^*} n^{|\mathcal{C}(\mathcal{G}_p)|} d^{|\mathcal{C}(\mathcal{G}'_p)|-2} \mathbf{e}_{j_0} \mathbf{e}_{j_0}^\top \\
& = \left[\sum_{\mathcal{P}: j_1, j_{k+1} \in c^*} n^{|\mathcal{C}(\mathcal{G}_p)|} d^{|\mathcal{C}(\mathcal{G}'_p)|-1} + \sum_{\mathcal{P}: j_1 \in c^*, j_{k+1} \notin c^*} n^{|\mathcal{C}(\mathcal{G}_p)|} d^{|\mathcal{C}(\mathcal{G}'_p)|-2} \right. \\
& \quad \left. + \sum_{\mathcal{P}: j_{k+1} \in c^*, j_1 \notin c^*} n^{|\mathcal{C}(\mathcal{G}_p)|} d^{|\mathcal{C}(\mathcal{G}'_p)|-2} \right] \mathbf{e}_s \mathbf{e}_s^\top + \sum_{\mathcal{P}: j_1, j_{k+1} \notin c^*} n^{|\mathcal{C}(\mathcal{G}_p)|} d^{|\mathcal{C}(\mathcal{G}'_p)|-3} \mathbf{I}.
\end{aligned}$$

Recall Equation (19), we prove that

$$\mathbb{E}[\mathbf{S} \mathbf{u}_s \mathbf{u}_s^\top \mathbf{S}^k \mathbf{u}_s \mathbf{u}_s^\top \mathbf{S}^{k'}] = \alpha_4 \mathbf{u}_s \mathbf{u}_s^\top + \alpha_5 \mathbf{I}.$$

Hence, the proof is complete. \square

Lemma D.8. Suppose $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, then the expectation is in the following form for any k :

$$\mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S}^k \mathbf{\Gamma} \mathbf{S}^{k'}] = \beta_1 \mathbf{\Lambda} \mathbf{\Gamma} + \beta_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \beta_3 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \beta_4 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} + \beta_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}.$$

where $\mathbf{\Lambda} = \sum_{j=1}^d \lambda_j^\mathbf{\Lambda} \mathbf{u}_j \mathbf{u}_j^\top$, $\mathbf{\Gamma} = \sum_{j=1}^d \lambda_j^\mathbf{\Gamma} \mathbf{u}_j \mathbf{u}_j^\top$.

Proof. By lemma D.7, we have:

$$\begin{aligned}
& \mathbb{E}[\mathbf{S} \mathbf{\Lambda} \mathbf{S}^k \mathbf{\Gamma} \mathbf{S}^{k'}] \\
& = \sum_{j=1}^d \sum_{i \neq j}^d \lambda_i^\mathbf{\Lambda} \lambda_j^\mathbf{\Gamma} (\alpha_1 \mathbf{u}_i \mathbf{u}_i^\top + \alpha_2 \mathbf{u}_j \mathbf{u}_j^\top + \alpha_3 \mathbf{I}) + \sum_{i=1}^d \lambda_i^\mathbf{\Lambda} \lambda_i^\mathbf{\Gamma} (\alpha_4 \mathbf{u}_i \mathbf{u}_i^\top + \alpha_5 \mathbf{I})
\end{aligned}$$

The first term here can be expand into the following form:

$$\begin{aligned}
& \sum_{j=1}^d \sum_{i \neq j}^d \lambda_i^\mathbf{\Lambda} \lambda_j^\mathbf{\Gamma} (\alpha_1 \mathbf{u}_i \mathbf{u}_i^\top + \alpha_2 \mathbf{u}_j \mathbf{u}_j^\top + \alpha_3 \mathbf{I}) \\
& = \alpha_1 \text{tr}(\mathbf{\Gamma}) \mathbf{\Lambda} + \alpha_2 \text{tr}(\mathbf{\Lambda}) \mathbf{\Gamma} + \alpha_3 \text{tr}(\mathbf{\Lambda}) \text{tr}(\mathbf{\Gamma}) \mathbf{I} - (\alpha_1 + \alpha_2) \mathbf{\Lambda} \mathbf{\Gamma} - \alpha_3 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}
\end{aligned}$$

Meanwhile, the second term is directly $\alpha_4 \mathbf{\Lambda} \mathbf{\Gamma} + \alpha_5 \text{tr}(\mathbf{\Lambda} \mathbf{\Gamma}) \mathbf{I}$. We pick $\beta_2 = \alpha_1, \beta_3 = \alpha_2, \beta_4 = \alpha_3, \beta_1 = \alpha_4 - \alpha_1 - \alpha_2, \beta_5 = \alpha_5 - \alpha_3$, and we complete the proof. \square

E EXPERIMENTAL DETAILS

For all our experiments, we use pytorch Paszke et al. (2019) and models are trained on an NVIDIA RTX A6000s. Each experiment takes about 1 hour.

Setup In all our experiments, we choose $d = 10$, $n = 20$ and $\eta = 0.4$. The architecture is

$$f_{\text{LSA}}(\mathbf{Z}; \mathbf{V}, \mathbf{W})_{[:, -1]} = \mathbf{Z}_{[:, -1]} + \mathbf{V} \mathbf{Z} \cdot \frac{\mathbf{Z}^\top \mathbf{W} \mathbf{Z}_{[:, -1]}}{n}$$

and data is drawn from the distribution in Equation (1). The batch size B is 1000 and the learning rate α is 0.001. The total time is $\tau = 750$ iterations. In the first experiment, k is chosen as 20 while $k = 10, 20, 30, 40$ in the second experiment. The baseline (evaluation loss of transformers without CoT) is given by Corollary 3.1 where $\eta^* = \frac{n}{n+d+1}$:

$$\mathcal{L}^{\text{Eval}}(\mathbf{V}, \mathbf{W}) \geq \frac{1}{2} \left(d - 2\eta^* d + \frac{\eta^{*2}}{n} (n + d + 1) d \right)$$

In-distribution Generalization We empirically verify the evaluation loss gap between transformers with and without CoT shown by Theorem 3.1 and Theorem 3.2. Our experiments in Figure 2 demonstrate that the evaluation loss of transformers with CoT converges to near zero even when $k = 10$. See Section 5 for details.

Out-of-distribution Generalization In addition, we empirically verify the OOD generalization result shown by Theorem 4.2. We sample 10 different covariance matrices from the distribution which complies to

$$\frac{\delta}{\eta} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq \frac{2 - \delta}{\eta}$$

where $\eta = 0.4$ and $\delta = 0.4$. 10 experiments are taken to show the generality of our results for each set of experiment. Our experiment in Figure 3 exhibits that the OOD loss of transformers with CoT converges to near zero when $k = 10, 20, 30, 40$ as the training loss/in-distribution loss converges to zero. The final loss also drops when the number of reasoning steps increases.

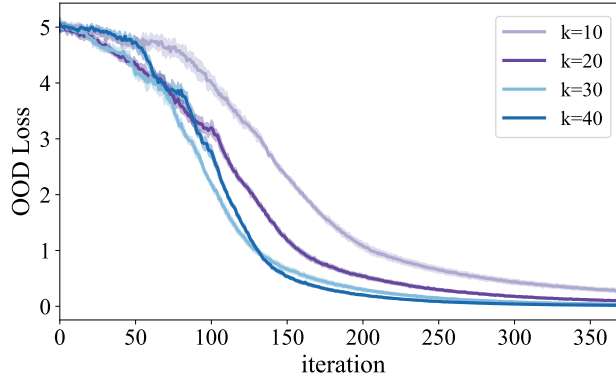


Figure 3: **OOD Generalization:** We plot the OOD loss $\mathcal{L}_{\Sigma}^{\text{Eval}}$ when $n = 20, d = 10$. Each set of experiments sampled 10 different Σ . The mean results are presented as line charts, with variance represented by shaded areas. As shown, OOD loss will converge to near zero.

Given all experiments above, we conclude that transformers with CoT can converge to our construction (Theorem 4.1), surpass those without CoT (Corollary 3.1, Theorem 3.2) and generalize well to unseen data (Theorem 4.2).