

TRAINING NONLINEAR TRANSFORMERS FOR CHAIN-OF-THOUGHT INFERENCE: A THEORETICAL GENERALIZATION ANALYSIS

Hongkang Li¹, Songtao Lu^{2,*}, Pin-Yu Chen³, Xiaodong Cui³, Meng Wang^{1,†}

¹Rensselaer Polytechnic Institute, ²The Chinese University of Hong Kong, ³IBM Research

ABSTRACT

Chain-of-Thought (CoT) is an efficient prompting method that enables the reasoning ability of large language models by augmenting the query using multiple examples with multiple intermediate steps. Despite the empirical success, the theoretical understanding of how to train a Transformer to achieve the CoT ability remains less explored. This is primarily due to the technical challenges involved in analyzing the nonconvex optimization on nonlinear attention models. To the best of our knowledge, this work provides the first theoretical study of training Transformers with nonlinear attention to obtain the CoT generalization capability so that the resulting model can inference on unseen tasks when the input is augmented by examples of the new task. We first quantify the required training samples and iterations to train a Transformer model towards CoT ability. We then prove the success of its CoT generalization on unseen tasks with distribution-shifted testing data. Moreover, we theoretically characterize the conditions for an accurate reasoning output by CoT even when the provided reasoning examples contain noises and are not always accurate. In contrast, in-context learning (ICL), which can be viewed as one-step CoT without intermediate steps, may fail to provide an accurate output when CoT does. These theoretical findings are justified through experiments.

1 INTRODUCTION

Transformer-based large-scale foundation models, such as GPT-3 (Brown et al., 2020), GPT-4 (OpenAI, 2023), LLaMa (Touvron et al., 2023a;b), and Sora (Liu et al., 2024), have demonstrated remarkable success across various tasks, including natural language processing (Brown et al., 2020; Touvron et al., 2023b), multimodal learning (OpenAI, 2023; Radford et al., 2021), and image/video generation (OpenAI, 2023; Liu et al., 2024). What is more surprising is that large language models (LLMs) demonstrate reasoning ability through the so-called “Chain-of-Thought” (CoT) method (Wei et al., 2022). The objective is to let a pre-trained LLM generate K steps of reasoning given input query x_{query} without any fine-tuning. To achieve that, the input x_{query} is augmented with l examples $\{x_i, \{y_{i,j}\}_{j=1}^K\}_{i=1}^l$ of a certain K -step reasoning task, where each x_i is the input with $y_{i,j}$ as the j -th reasoning step, and $y_{i,K}$ is the final output. A pre-trained model then takes the resulting augmented input, referred to as a *prompt*, and outputs the corresponding reasoning steps $\{z_j\}_{j=1}^K$ for x_{query} , or simply outputs z_K . CoT can be viewed as an extended and more intelligent method than the previous in-context learning (ICL) method, where only input-label pairs $\{x_i, y_{i,K}\}_{i=1}^l$ are augmented in the prompt to predict z_K with the pre-trained model.

Inspired by the outstanding empirical performance of CoT in arithmetic reasoning (Wang et al., 2023; Zhang et al., 2023b; Wang & Zhou, 2024), symbolic reasoning (Zhang et al., 2023b; Zhou et al., 2023), and commonsense reasoning (Wang et al., 2023; Wang & Zhou, 2024), there have been some recent works (Li et al., 2023c; Feng et al., 2023; Li et al., 2024d; Yang et al., 2024; Wen et al., 2024) on the theoretical understanding of CoT. These works investigate CoT from the perspective of expressive power, i.e., they construct the Transformer architecture that is proven to have the CoT ability. They also demonstrate empirically that supervised training on pairs of CoT

*This work was done when Prof. Songtao Lu was at IBM Research.

†Corresponding author. Email: wangm7@rpi.edu.

prompts and corresponding outputs can lead to models with CoT ability. However, none of these results theoretically address the question of why a Transformer can obtain generalization-guaranteed CoT ability by training from data with gradient-based methods. Meanwhile, another line of research (Zhang et al., 2023a; Huang et al., 2023; Wu et al., 2023; Li et al., 2024a) aims to unveil the reasons behind the ICL ability of Transformers through characterizing the training dynamics of a Transformer in the supervised setting. These analyses are specifically applicable to ICL. Therefore, a theoretical question still remains less explored, i.e.,

Why can a Transformer be trained to generalize on multi-step reasoning tasks via CoT?

1.1 MAJOR CONTRIBUTIONS

Following Li et al. (2023c); Feng et al. (2023); Li et al. (2024d); Yang et al. (2024); Wen et al. (2024), we train the model in a supervised setting using prompt and label pairs. This paper provides the first theoretical analysis of the training dynamics of nonlinear Transformers to achieve CoT ability. We prove that the learned model has guaranteed CoT ability for new tasks with distribution shifts from the training tasks, even when there exist noisy and erroneous context examples in the prompt. We theoretically characterize the required number of training samples and iterations needed to train a desirable model and the number of context examples required for successful CoT reasoning with a generalization guarantee. Moreover, we provide a theoretical explanation for why CoT outperforms ICL in some cases. Our main technical contributions are as follows:

- 1. A quantitative analysis of how the training can enable the CoT ability:** We theoretically analyze the training dynamics on a one-layer single-head attention-only Transformer and quantify the required number of context examples in each training sample, the total number of training samples, and the number of training iterations needed to acquire CoT ability. We illustrate that the CoT ability results from the property that the attention values of the learned model are concentrated on testing context examples with the same input patterns as the testing query during each reasoning step.
- 2. A quantitative analysis of how context examples affect CoT performance:** We characterize the required number of context examples in the testing prompt for successful CoT reasoning when noise and error exist in contexts. Our quantitative bounds are consistent with the intuition that more accurate context examples and more similar examples to the query improve CoT accuracy.
- 3. A theoretical characterization of why CoT outperforms ICL:** We provide a quantitative analysis of the requirements for successful ICL reasoning with our studied trained model. We show that successful ICL requires an additional condition that the prompt has a dominant number of correct input-label examples, while the success of CoT does not depend on this condition. This can be viewed as one of the possible reasons why CoT outperforms ICL.

1.2 RELATED WORKS

Expressive power of CoT Li et al. (2023c) proves the existence of a Transformer that can learn a multi-layer perceptron (MLP). They interpret CoT as first filtering important tokens and then making predictions by ICL. They also establish the required number of context examples for a desired prediction with the constructed Transformer. Feng et al. (2023); Li et al. (2024d); Merrill & Sabharwal (2024) show that Transformers with CoT are more expressive than Transformers without CoT. Yang et al. (2024); Wen et al. (2024) show the superiority of standard Transformers in some reasoning tasks compared with recurrent neural networks and linear Transformers.

Theoretical analysis of ICL As a simplified one-step version of CoT, ICL has gained much attention from the theoretical community. Garg et al. (2022); Akyürek et al. (2023); Bai et al. (2023); Guo et al. (2023) demonstrate that Transformers are expressive to conduct many machine learning algorithms in context. Akyürek et al. (2023); Von Oswald et al. (2023); Ahn et al. (2023); Cheng et al. (2023); Ding et al. (2024) especially show the existence of Transformers to implement gradient descent and its variants with different input prompts. Zhang et al. (2023a); Huang et al. (2023); Wu et al. (2023); Li et al. (2024a) explore the training dynamics and generalization of ICL on single-attention Transformers. Cui et al. (2024); Chen et al. (2024) provably show the superiority of multi-head attention over single-head attention to achieve ICL ability.

Training and Generalization of Transformers There have been several recent works about the optimization and generalization analysis of Transformers. Jelassi et al. (2022); Li et al. (2023d); Oymak et al. (2023); Li et al. (2023a;b; 2024b); Luo (2023); Huang et al. (2024); Zhang et al. (2024) study

the generalization of one-layer Transformers by assuming spatial association, semantic/contextual structure, or the majority voting of tokens in the data. Oymak et al. (2023); Tarzanagh et al. (2023b;a); Tian et al. (2023a;b); Li et al. (2024c); Ildiz et al. (2024); Nichani et al. (2024); Makkuva et al. (2024b) investigate the training dynamics or loss landscape of Transformers for the next token prediction by assuming infinitely long input sequences, causal structure/Markov Chain of data, or a proper prediction head. Deora et al. (2023); Chen & Li (2024) analyze the optimization and generalization of multi-head attention networks.

2 PROBLEM FORMULATION

We study the problem of learning and generalization of K -steps reasoning tasks. Each task $f = f_K \circ \dots \circ f_2 \circ f_1$ is a composition of functions $\{f_i\}_{i=1}^K$ and outputs labels z_1, z_2, \dots, z_K for the input \mathbf{x}_{query} . During the k -th reasoning step, $k \in [K]$, the label is $z_k = f_k(z_{k-1})$, where $z_0 := \mathbf{x}_{query}$.

2.1 TRAINING TO ACQUIRE THE CHAIN-OF-THOUGHT ABILITY

Following theoretical analysis (Feng et al., 2023; Li et al., 2024d; Wen et al., 2024) and empirical works like process supervision (Lightman et al., 2024), we first investigate the training on a Transformer model to obtain the CoT ability in evaluating new data and tasks. It is a supervised learning setting on pairs of prompts and labels. Different from the testing prompt that includes examples and only \mathbf{x}_{query} , the training prompt includes multiple K -steps reasoning examples and a $(k-1)$ -step reasoning of \mathbf{x}_{query} for any k in $[K]$, and the label for this prompt is z_k . Specifically,

Training Prompt and Label for CoT. For every prompt and output pair from a task $f = f_K \circ \dots \circ f_2 \circ f_1$, we construct a prompt \mathbf{P} that include the query input z_{k-1} by prepending l_{tr} reasoning examples and the first $k-1$ steps of the reasoning query. The prompt \mathbf{P} of the query input z_{k-1} is formulated as:

$$\mathbf{P} = (\mathbf{E}_1, \mathbf{E}_2, \dots, \mathbf{E}_{l_{tr}}, \mathbf{Q}_k) \in \mathbb{R}^{2d_x \times (l_{tr}K + k)},$$

$$\text{where } \mathbf{E}_i = \begin{pmatrix} \mathbf{x}_i & \mathbf{y}_{i,1} & \dots & \mathbf{y}_{i,K-1} \\ \mathbf{y}_{i,1} & \mathbf{y}_{i,2} & \dots & \mathbf{y}_{i,K} \end{pmatrix}, \mathbf{Q}_k = \begin{pmatrix} z_0 & z_1 & \dots & z_{k-2} & z_{k-1} \\ z_1 & z_2 & \dots & z_{k-1} & \mathbf{0} \end{pmatrix}, i \in [l_{tr}], \quad (1)$$

where \mathbf{E}_i is the i -th context example, and \mathbf{Q}_k is the first k steps of the reasoning query for any k in $[K]$. We have $\mathbf{y}_{i,k} = f_k(\mathbf{y}_{i,k-1})$ and $z_k = f_k(z_{k-1})$ for $i \in [l_{tr}]$, $k \in [K]$ with a notation $\mathbf{y}_{i,0} := \mathbf{x}_i$. Let \mathbf{p}_s and \mathbf{p}_{query} be the s -th column and the last column of \mathbf{P} , respectively, for $s \in [l_{tr}K + k - 1]$. $\mathbf{x}_i, \mathbf{y}_{i,k}, \mathbf{z}_j \in \mathbb{R}^{d_x}$ for $i \in [l_{tr}]$ and $j, k \in [K]$. We respectively call \mathbf{x}_i and $\mathbf{y}_{i,k}$ context inputs and outputs of the k -th step of the i th context example. For simplicity of presentation, we denote z as the label of \mathbf{P} , which is indeed z_k for (1). All the notations are summarized in Table 3 in Appendix.

The learning model is a single-head, one-layer attention-only Transformer. We consider positional encoding $\{\mathbf{c}_k\}_{k=1}^K \in \mathbb{R}^{2d_x}$. Following theoretical works (Jelassi et al., 2022; Huang et al., 2024; Ildiz et al., 2024), we add the positional encoding to each \mathbf{p}_i by $\tilde{\mathbf{p}}_i = \mathbf{p}_i + \mathbf{c}_{(i \bmod K)}$, $i \in [K(l_{tr} + 1)]$. $\tilde{\mathbf{p}}_{query}$ is also defined by adding the corresponding \mathbf{c}_k to \mathbf{p}_{query} . Mathematically, given a prompt \mathbf{P} defined in (1) with $\text{len}(\mathbf{P})$ (which is at most $K(l_{tr} + 1)$) denoting the number of columns, it can be written as

$$F(\Psi; \mathbf{P}) = \sum_{i=1}^{\text{len}(\mathbf{P})-1} \mathbf{W}_V \tilde{\mathbf{p}}_i \cdot \text{softmax}((\mathbf{W}_K \tilde{\mathbf{p}}_i)^\top \mathbf{W}_Q \tilde{\mathbf{p}}_{query}), \quad (2)$$

where $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{m \times (2d_x)}$, $\mathbf{W}_V \in \mathbb{R}^{d_x \times (2d_x)}$ are the embedding matrices for queries, keys, and values, respectively. $\Psi := \{\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V\}$ is the set of all model weights¹. Typically, $m > 2d_x$. Here, $\text{softmax}((\mathbf{W}_K \tilde{\mathbf{p}}_i)^\top \mathbf{W}_Q \tilde{\mathbf{p}}_{query}) = e^{(\mathbf{W}_K \tilde{\mathbf{p}}_i)^\top \mathbf{W}_Q \tilde{\mathbf{p}}_{query}} / \sum_{j=1}^{\text{len}(\mathbf{P})-1} e^{(\mathbf{W}_K \tilde{\mathbf{p}}_j)^\top \mathbf{W}_Q \tilde{\mathbf{p}}_{query}}$.

The training problem to enhance the reasoning capability solves the empirical risk minimization,

$$\min_{\Psi} R_N(\Psi) := \frac{1}{N} \sum_{n=1}^N \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n), \quad (3)$$

using N prompt and label pairs $\{\mathbf{P}^n, \mathbf{z}^n\}_{n=1}^N$. For the n -th sample, \mathbf{x}_{query}^n and the context input \mathbf{x}_i^n are all sampled from an unknown distribution \mathcal{D} , the training task f^n is sampled from \mathcal{T} , k is randomly selected from 1 to K , and \mathbf{P}^n is constructed following (1). The loss function is squared loss, i.e., $\ell(\Psi; \mathbf{P}^n, \mathbf{z}^n) = 1/2 \cdot \|\mathbf{z}^n - F(\Psi; \mathbf{P}^n)\|^2$, where $F(\Psi; \mathbf{P}^n)$ is defined in (2).

¹We focus on a one-layer single-head Transformer motivated by recent advancements and current state in Transformer and CoT analysis. Please see Appendix B.1 for discussion.

2.2 TRAINING ALGORITHM

For simplicity of analysis, we let $\mathbf{W} = \mathbf{W}_K^\top \mathbf{W}_Q$ and $\mathbf{W}_V = (\mathbf{0}_{d_x \times d_x} \mathbf{I}_{d_x}) \in \mathbb{R}^{d_x \times (2d_x)}$ as (Jelassi et al., 2022; Huang et al., 2023; Zhang et al., 2023a; Huang et al., 2024). Let $\{\mathbf{c}_k\}_{k=1}^K$ be a set of orthonormal vectors. The model is trained using stochastic gradient descent (SGD) with step size η with batch size B , summarized in Algorithm 1 in Appendix C. Each entry of $\mathbf{W}^{(0)}$ is generated from $\mathcal{N}(0, \xi^2)$ for a tiny $\xi > 0$. \mathbf{W}_V is fixed during the training. The fraction of prompts with \mathbf{z}_{k-1} as the query input is $1/K$ by uniform sampling for any $k \in [K]$ in each batch.

2.3 CHAIN-OF-THOUGHT INFERENCE

We then consider another K -steps reasoning task $f \in \mathcal{T}'$, whose target is to predict labels $\{\mathbf{z}_k\}_{k=1}^K$ given the input query \mathbf{x}_{query} . \mathcal{T}' is the set of testing tasks, and $\mathcal{T}' \neq \mathcal{T}$.

Testing Prompt for CoT. The testing prompt \mathbf{P} is composed of l_{ts} ($\leq l_{tr}$) context examples of K steps plus a query, which is constructed as

$$\mathbf{P} = (\mathbf{E}_1, \mathbf{E}_2, \dots, \mathbf{E}_{l_{ts}}, \mathbf{p}_{query}) \in \mathbb{R}^{(2d_x) \times (l_{ts}K+1)}, \mathbf{p}_{query} = (\mathbf{x}_{query}^\top, \mathbf{0}^\top)^\top, \quad (4)$$

where \mathbf{E}_i follows the form in (1) for $i \in [l_{ts}]$.

We follow the CoT-I/O scheme formulated in (Li et al., 2023c; Feng et al., 2023; Li et al., 2024d; Yang et al., 2024; Park et al., 2024) as the inference method. Specifically, for a K -step CoT with l_{ts} examples on a certain $f \in \mathcal{T}'$, given the testing prompt \mathbf{P} defined in (4), let $\mathbf{P}_1 = \mathbf{P}$ and \mathbf{P}_0 be the first $K \cdot l_{ts}$ columns of \mathbf{P} . When we use CoT prompting for prediction in the k -th step, we first generate the output \mathbf{v}_k , $k \in [K]$ via greedy decoding by feeding the k -th step prompt \mathbf{P}_k to the trained model Ψ obtained from (3). The greedy decoding scheme means outputting the most probable token from the discrete set \mathcal{Y} of all possible outputs, as stated in (5).

$$\mathbf{v}_k = \arg \min_{\mathbf{u} \in \mathcal{Y}} \frac{1}{2} \|F(\Psi; \mathbf{P}_k) - \mathbf{u}\|^2, \text{ (greedy decoding)} \quad (5)$$

Then, we use the output \mathbf{v}_k to update \mathbf{P}_k and use \mathbf{v}_k as the query input to form the input prompt \mathbf{P}_{k+1} for the next step, which is computed as

$$\mathbf{P}_k = (\mathbf{P}_{k-1} \mathbf{q}_k) \in \mathbb{R}^{(2d_x) \times (Kl_{ts}+k)}, \mathbf{P}_{k+1} = (\mathbf{P}_k \mathbf{q}_{k+1}) \in \mathbb{R}^{(2d_x) \times (Kl_{ts}+k+1)}, \quad (6)$$

where $\mathbf{q}_k = (\mathbf{v}_{k-1}^\top \mathbf{v}_k^\top)^\top$, $\mathbf{q}_{k+1} = (\mathbf{v}_k^\top \mathbf{0}^\top)^\top$,

where \mathbf{q}_k is the k -th step reasoning column for the query. The model finally outputs $\mathbf{v}_1, \dots, \mathbf{v}_K$ as CoT result for query \mathbf{x}_{query} by (5). The CoT process is summarized in Algorithm 2 of Appendix C.

When $K \geq 2$, following (Li et al., 2023c; Feng et al., 2023; Li et al., 2024d; Yang et al., 2024), the **CoT generalization error** given the testing query \mathbf{x}_{query} , the testing data distribution \mathcal{D}' , and the labels $\{\mathbf{z}_k\}_{k=1}^K$ on a K -steps testing task $f \in \mathcal{T}'$ is defined as

$$\bar{R}_{CoT, \mathbf{x}_{query} \sim \mathcal{D}', f \in \mathcal{T}'}^f(\Psi) = \mathbb{E}_{\mathbf{x}_{query} \sim \mathcal{D}'} \left[\frac{1}{K} \sum_{k=1}^K \mathbb{1}[\mathbf{z}_k \neq \mathbf{v}_k] \right], \quad (7)$$

which measures the average error between the output and the label of each reasoning step. A zero CoT generalization error indicates correct generations in all K steps.

2.4 IN-CONTEXT LEARNING INFERENCE

The ICL inference on a K -steps reasoning task $f \in \mathcal{T}'$ only predicts the final-step label by prepending examples of input and label pairs before the query. ICL can be viewed as a one-step CoT without intermediate steps. Here, we evaluate the ICL performance of the trained model.

Testing Prompt for ICL. Mathematically, ICL is implemented by constructing a prompt \mathbf{P} as below,

$$\mathbf{P} = (\mathbf{E}_1, \dots, \mathbf{E}_{l_{ts}}, \mathbf{p}_{query}), \text{ where } \mathbf{p}_{query} = \begin{pmatrix} \mathbf{x}_{query} \\ \mathbf{0} \end{pmatrix}, \mathbf{E}_i = \begin{pmatrix} \mathbf{x}_i & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{y}_{i,K} & \mathbf{0} & \dots & \mathbf{0} \end{pmatrix} \quad (8)$$

$\mathbf{P} \in \mathbb{R}^{(2d_x) \times (l_{ts}K+1)}$, $\mathbf{E}_i \in \mathbb{R}^{(2d_x) \times K}$ for $i \in [l_{ts}]$. Note that in the ICL setting, \mathbf{E}_i only has input \mathbf{x}_i and the K -step output $\mathbf{y}_{i,K}$ but does not include any intermediate labels. We pad zeros in \mathbf{E}_i so that its dimension is the same as \mathbf{E}_i in (1) for the inference with the same model as for CoT. The ICL output is $\mathbf{v} = \arg \min_{\mathbf{u} \in \mathcal{Y}} \frac{1}{2} \|F(\Psi; \mathbf{P}) - \mathbf{u}\|^2$, following (5). The **ICL generalization error** is

$$\bar{R}_{ICL, \mathbf{x}_{query} \sim \mathcal{D}', f \in \mathcal{T}'}^f(\Psi) = \mathbb{E}_{\mathbf{x}_{query} \sim \mathcal{D}'} [\mathbb{1}[\mathbf{z}_K \neq \mathbf{v}]], \quad (9)$$

which measures the error between the one-step reasoning output and the final step label.

3 THEORETICAL RESULTS

We first summarize the main theoretical insights in Section 3.1. Then, we introduce the formulation of data and tasks in Section 3.2. Sections 3.3, 3.4, and 3.5, respectively characterize the training analysis of the Transformer and generalization using CoT and ICL with the trained model.

3.1 MAIN THEORETICAL INSIGHTS

We consider the setup that the model is trained using samples generated from tasks in \mathcal{T} that operate on M orthonormal training-relevant (TRR) patterns, while both CoT and ICL are evaluated on tasks in \mathcal{T}' that operate on M' orthonormal testing-relevant (TSR) patterns. We consider the general setup that the context examples in the prompt for CoT and ICL testing are both noisy, i.e., TSR patterns with additive noise, and partially inaccurate, i.e., the reasoning in some examples contains incorrect steps. Our main insights are as follows.

P1. Training Dynamics of Nonlinear Transformer towards CoT. We theoretically analyze the training dynamics on a one-layer single-head attention-only Transformer to acquire the CoT generalization ability and characterize the required number of training samples and iterations. Theorem 1 shows that to learn a model with guaranteed CoT ability, the required number of context examples in each training sample and the total number of training samples/iterations are linear in α^{-1} and α^{-2} , respectively, where α is the fraction of context examples with inputs that share the same TRR patterns as the query. This is consistent with the intuition that the CoT performance is enhanced if more context examples are similar to the query. Moreover, the attention values of the learned model are proved to be concentrated on testing context examples that share similar input TSR patterns as the testing query during each of the reasoning steps (Proposition 1), which is an important property that leads to the success of the CoT generalization.

P2. Guaranteed CoT Generalization. To achieve zero CoT error on tasks in \mathcal{T}' and data based on TSR patterns that contain a non-trivial component in the span of TRR patterns with the learned model, Theorem 2 shows that the required number of context examples, where noise and errors are present, for task f in the testing prompt is proportional to $(\alpha' \tau^f \rho^f)^{-2}$. Here, α' is the fraction of context examples with inputs that share the same TSR patterns as the query. τ^f in $(0, 1)$ measures the fraction of accurate context examples, and a larger constant ρ^f in $(0, 1)$ reflects a higher reasoning accuracy in each step of the examples. This result formally characterizes the intuition that more accurate context examples and more similar examples to the query improve the CoT accuracy.

P3. CoT outperforms ICL. In Theorem 3, We theoretically show that the required number of testing context examples for ICL to be successful has a similar form to that for CoT in Theorem 2, but with an additional requirement (Condition 1) that the fraction of correct input-label examples in the testing prompt must be dominant. Because not all testing cases satisfy this requirement, our result provides one explanation for why CoT sometimes outperforms ICL.

3.2 THE FORMULATION OF DATA AND TASKS

Training data and tasks: Consider M training-relevant (TRR) patterns $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_M$, which form an orthonormal set $\mathcal{M} = \{\boldsymbol{\mu}_i\}_{i=1}^M$. $M = \Theta(d)$, $M \leq d$. $(\boldsymbol{\mu}_i^\top, 0_{d \times}^\top)^\top \perp \mathbf{c}_k$ for $i \in [M']$, $k \in [K]$.

Every training prompt \mathbf{P} in (1) contains the query and training examples from the same training task f in the set of training tasks \mathcal{T} . Specifically, each training task f is a composition of K functions $f = f_K \circ \dots \circ f_2 \circ f_1$ where each function f_k belongs to a function set \mathcal{F} . The k -th step label of the query is $\mathbf{z}_k = f_k(\mathbf{z}_{k-1})$ given the k -th step input \mathbf{z}_{k-1} with $\mathbf{z}_k \in \mathcal{M}$, $k \in [K]$. Moreover, the k -th step label of the i -th ($i \in [l_{tr}]$) context example is $\mathbf{y}_{i,k} = f_k(\mathbf{y}_{i,k-1})$ given the $k-1$ th step input $\mathbf{y}_{i,k-1}$, $k \in [K]$ with $\mathbf{x}_i, \mathbf{y}_{i,k} \in \mathcal{M}$, where $\mathbf{y}_{i,0} := \mathbf{x}_i$ ². We assume that $f_k(\mathbf{x}) \neq f_{k'}(\mathbf{x}')$ if and only if either $\mathbf{x} \neq \mathbf{x}'$ or $f_k \neq f_{k'}$.

Training prompt: Consider a training prompt \mathbf{P} on task $f \in \mathcal{T}$ defined in (1) with the query input \mathbf{z}_{k-1} , $k \in [K]$. Let $\alpha \in (0, 1 - c]$ for some constant $c > 0$ ³ denote the fraction of context examples with input sharing the same TRR pattern as the query input.

²The formulation of f is motivated by recent theoretical works on model training or ICL with Transformers. Please see Appendix B.2 for details.

³This is to prevent the trivial case that the model only learns the positional encoding but not the TRR patterns when α becomes arbitrarily close to 1.

Testing task and query: Consider M' testing-relevant (TSR) patterns $\mu'_1, \mu'_2, \dots, \mu'_{M'}$, which form an orthonormal set $\mathcal{M}' = \{\mu'_i\}_{i=1}^{M'}$. $M' \leq M$. We also have $\mu'_i \perp \mathbf{c}_k$ for $i \in [M'], k \in [K]$. Let \mathcal{T}' denote the set of testing tasks, which all operate on patterns in \mathcal{M}' rather than \mathcal{M} in training tasks in \mathcal{T} . Every testing task $f = f_K \circ \dots \circ f_2 \circ f_1 \in \mathcal{T}'$ is a composition of K functions. The reasoning for the testing query is considered to be *noiseless* and *accurate*. That means,

$$\mathbf{z}_k \in \mathcal{M}' \text{ for all } k \in \{0\} \cup [K], \text{ and } \mathbf{z}_k = f_k(\mathbf{z}_{k-1}), \mathbf{z}_0 = \mathbf{x}_{\text{query}}.$$

Testing prompt: We consider the general setup that testing examples are *noisy* and *erroneous*. By noisy examples, we mean all inputs and outputs of each step are noisy versions of TSR patterns, i.e.,

$$\mathbf{x}_i, \mathbf{y}_{i,k} \in \{\mathbf{b} \in \mathbb{R}^d \mid \mathbf{b} = \mu'_j + \delta, j \in [M'], \delta \perp \mathcal{M}', \|\delta\| \leq \sqrt{2}/2\}, \quad (10)$$

with noise $\delta \neq 0$ for $i \in [Kl_{ts}^f], k \in [K]$. Denote $\text{TSR} : \mathbb{R}^d \mapsto \mathbb{Z}^+$ as a function that outputs the index of the TSR pattern of the noisy input. We consider the case that at least an α' fraction of context examples where the TSR pattern of the input $\mathbf{y}_{s,1}, s \in [l_{ts}^f]$ is the same as $\mathbf{x}_{\text{query}}$.

By erroneous examples, we mean that the reasoning steps in test examples may contain errors. To formally model this, we define the **step-wise transition matrices** $\{A_k^f\}_{k=1}^K \in \mathbb{R}^{M' \times M'}$ such that A_k^f represents the reasoning probabilities of step k in test examples. Specifically, there exists some constant ρ^f in $(0, 1)$ such that for all $s \in [l_{ts}^f], k \in [K]$, the i, j -th entry of A_k^f satisfies

$$A_{k(i,j)}^f = \Pr(\text{TSR}(\mathbf{y}_{s,k}) = j \mid \text{TSR}(\mathbf{y}_{s,k-1}) = i), \quad (11)$$

and $A_{k(i,j^*)}^f \geq 1/(1 - \rho^f) \cdot A_{k(i,j)}^f, \forall j \in [M']$, where $\mu'_{j^*} = f_k(\mu'_i)$,

Note that (11) characterizes a general case in inference that for any given k , in the k -th reasoning step of the test example, the k -th step output is a noisy version of the true label with the highest probability, which guarantees that the examples are overall informative in the k -th step. This requirement is intuitive because otherwise, these examples would overall provide inaccurate information on the k -th step reasoning. Moreover, (11) models the general case that, with some probability, the k -step reasoning is inaccurate in the examples. ρ^f is referred to as the **primacy** of the step-wise transition matrices. ρ^f reflects the difference in the probability of correct reasoning and incorrect reasoning in each step, and a larger ρ^f indicates a larger probability of accurate reasoning.

Let $B^f = \prod_{k=1}^K A_k^f$ be the K -step transition matrix. Then $B_{(i,j)}^f$ is the probability that the K -th step output is a noisy version of μ'_j , when the input is a noisy version of μ'_i in the testing example. We similarly define ρ_o^f in $(0, 1)$ as the primacy of B^f , where

$$B_{(i,j^*)}^f \geq 1/(1 - \rho_o^f) \cdot B_{(i,j)}^f, \forall j \in [M'], j^* = \arg \max_{j \in [M']} B_{(i,j)}^f. \quad (12)$$

Example 1. Consider a simple two-step inference example with $K = 2$, μ'_1, μ'_2 as the TSR pattern, and $\delta = 0$ in inputs and outputs of every step, as shown in Figure 1. The black solid arrows denote the correct inference process, where $f_1(\mu'_1) = \mu'_1, f_1(\mu'_2) = \mu'_2, f_2(\mu'_1) = \mu'_2$, and $f_2(\mu'_2) = \mu'_1$. Hence, $\mu'_1 \rightarrow \mu'_1 \rightarrow \mu'_2$ and $\mu'_2 \rightarrow \mu'_2 \rightarrow \mu'_1$ are two inference **trajectories** under the function f . The testing examples contain errors and follow the transition matrices A_1^f and A_2^f (brown dashed arrows). We let $A_1^f = \begin{pmatrix} 0.6 & 0.4 \\ 0.4 & 0.6 \end{pmatrix}, A_2^f = \begin{pmatrix} 0.4 & 0.6 \\ 0.8 & 0.2 \end{pmatrix}$, which results in $B^f = \begin{pmatrix} 0.56 & 0.44 \\ 0.64 & 0.36 \end{pmatrix}$.

3.3 THE SAMPLE COMPLEXITY ANALYSIS OF THE TRAINING STAGE

We first characterize the convergence and the testing performance of the model during the training stage with sample complexity analysis in Theorem 1.

Theorem 1. For any $\epsilon > 0$, when (i) the number of context examples in every training sample is

$$l_{tr} \geq \Omega(\alpha^{-1}), \quad (13)$$

(ii) the number of iterations satisfies

$$T \geq \Omega(\eta^{-1} \alpha^{-2} K^3 \log \frac{K}{\epsilon} + \eta^{-1} MK(\alpha^{-1} + \epsilon^{-1})), \quad (14)$$

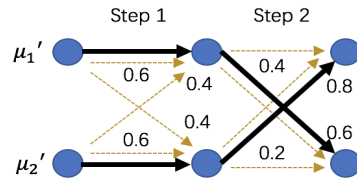


Figure 1: An example of a two-step inference

and (iii) the training tasks and samples are selected such that every TRR pattern is equally likely in every inference step and in each training batch⁴ with batch size $B \geq \Omega(\max\{\epsilon^{-2}, M\} \cdot \log M)$, the step size $\eta < 1$ and $N = BT$ samples, then with a high probability, the returned model guarantees

$$\mathbb{E}_{\mathbf{x}_{\text{query}} \in \mathcal{M}, f \in \mathcal{T}} [\ell(\Psi; \mathbf{P}, \mathbf{z})] \leq \mathcal{O}(\epsilon). \quad (15)$$

Theorem 1 indicates that with long enough training prompts and a sufficient number of iterations and samples for training, a one-layer Transformer can achieve a diminishing loss of $\mathcal{O}(\epsilon)$ on data following the same distribution as training examples. The results indicate that (i) the required number of context examples is proportional to α^{-1} ; (ii) the required number of iterations and samples increases as M and α^{-2} increases. As a sanity check, these bounds are consistent with the intuition that it will make the training stage more time- and sample-consuming if the number of TRR patterns increases or the fraction of prompt examples that share the same TRR pattern as the query decreases.

3.4 CoT GENERALIZATION GUARANTEE

In this section, we first define two quantities, τ^f , and τ_o^f for each testing task $f \in \mathcal{T}'$ based on the formulation of testing data and tasks in Section 3.2. These two quantities are used to characterize the CoT and ICL generalization in Theorems 2 and 3, respectively.

Definition 1. For $f = f_K \circ \dots \circ f_1 \in \mathcal{T}'$, we define the **min-max trajectory transition probability** as:

$$\tau^f = \min_{i \in [M']} \prod_{k=1}^K A_{k(\text{TSR}(f_{k-1} \circ \dots \circ f_0(\mu'_i)), \text{TSR}(f_k \circ \dots \circ f_0(\mu'_i)))}^f, \text{ where } f_0(\mu'_i) := \mu'_i, \forall i \in [M'], \quad (16)$$

which measures the minimum probability, over all the initial TSR patterns, of the K -step reasoning trajectory that has the highest probability over all K -step trajectories. We also define the **min-max input-label transition probability** as

$$\tau_o^f = \min_{i \in [M']} \max_{j \in [M']} B_{i,j}^f, \quad (17)$$

which measures the minimum probability, over all the initial TSR patterns, of the output that has the highest probability over outputs.

For instance, in Example 1 after (12), $\tau^f = \min\{0.36, 0.48\} = 0.36$, $\tau_o^f = \min\{0.56, 0.64\} = 0.56$.

Theorem 2 (CoT generalization). *Given a trained model, the training process of which satisfies conditions (i) to (iii) in Theorem 1, then as long as*

(iv) *each TSR pattern μ'_j in the orthonormal set $\{\mu'_j\}_{j=1}^{M'}$ satisfies*

$$\mu'_j = \lambda_j + \tilde{\mu}_j \quad (18)$$

where $\lambda_j \perp \text{span}(\mu_1, \dots, \mu_M)$, $\tilde{\mu}_j \in \text{span}(\mu_1, \dots, \mu_M)$, and $\|\tilde{\mu}_j\| \geq \Theta((\log \epsilon^{-1})^{-1})$, and (v) the number of testing examples for any $f \in \mathcal{T}'$ is

$$l_{ts}^f \geq \Omega((\alpha' \tau^f \rho^f)^{-2} \log M), \quad (19)$$

we have $\bar{R}_{\text{CoT}, \mathbf{x}_{\text{query}} \in \mathcal{M}', f \in \mathcal{T}'}^f(\Psi) = 0$.

Remark 1. *Theorem 2 proves that a trained one-layer Transformer can generate all K -steps reasoning correctly by CoT for a new task f in \mathcal{T}' with two additional conditions. Condition (iv) means that each TSR pattern in the task set \mathcal{T}' is the summation of a component that belongs to the span of the TRR patterns and a component that is perpendicular to the span.*

Condition (v) indicates that, to achieve the desired CoT accuracy, the number of context examples should be proportional to α'^{-2} , ρ_s^{f-2} , and τ^f , meaning it decreases as α' , ρ_s^f , or τ_s^f increase. It can be interpreted as follows, if the number of context examples remains fixed, an increase in α' , ρ_s^f , or τ_s^f results in improved CoT accuracy. This aligns with intuition, because α' represents the fraction of examples similar to the query, and ρ^f and τ^f reflect the accuracy of the reasoning steps in the context examples.

⁴Our analysis assumes that the whole set of \mathcal{M} is achievable uniformly in each step and training batch. This condition is to ensure a balanced gradient update among all TRR patterns, as used in (Li et al., 2024a) for ICL.

3.5 ICL GENERALIZATION AND COMPARISON WITH CoT

Because only input-label pairs are used as context examples for ICL, the input-label pairs in context examples should be accurate overall to be informative about the task. We formulate this requirement as Condition 1.

Condition 1. For the testing task $f = f_K \circ \dots \circ f_1 \in \mathcal{T}'$, we have that for any $i \in [M']$,

$$\text{TSR}(f(\mu'_i)) = \arg \max_{j \in [M']} B_{(i,j)}^f. \quad (20)$$

Condition 1 requires that in a context example, if the input TSR is μ'_i , then $f(\mu'_i)$ is the output TSR pattern with the highest probability over all TSR patterns. Note that (11) indicates that, for every k and i , when μ'_i is the k -th step input, $f_k(\mu'_i)$ is the step- k output with the highest probability over all TSR patterns. However, (11) does not necessarily imply (20). In Example 1, given the input μ'_1 , although the inference trajectory $\mu'_1 \rightarrow \mu'_1 \rightarrow \mu'_2$ under f has the highest probability over all 2-step trajectories, μ'_1 has the higher probability to be the final output than the correct output μ'_2 by the two-step transition matrix B^f , thus violating Condition 1.

Our result of the ICL generalization is stated as follows.

Theorem 3 (ICL generalization). *Given a trained model, the training process of which satisfies conditions (i) to (iii) of Theorem 1 and (18), for the testing task $f \in \mathcal{T}'$,*

Case A. if Condition 1 does not hold, then $\bar{R}_{ICL, \mathbf{x}_{query} \in \mathcal{M}', f \in \mathcal{T}'}^f(\Psi) \geq \Omega(1)$, no matter how large the number of training samples l_{ts}^f is;

Case B. if Condition 1 holds, then $\bar{R}_{ICL, \mathbf{x}_{query} \in \mathcal{M}', f \in \mathcal{T}'}^f(\Psi) = 0$, provided that

$$l_{ts}^f \geq \Omega((\alpha' \tau_o^f \rho_o^f)^{-2} \log M). \quad (21)$$

Remark 2 (Comparison between CoT and ICL). *Theorem 3(a) formally states that, Condition 1 is necessary for a successful ICL generalization. Because Condition 1 is not required for CoT generalization, CoT performs better than ICL if Condition 1 fails⁵. Theorem 3(b) characterizes that when Condition 1 holds, a desired ICL generalization needs a testing prompt length linear in α'^{-2} , ρ_o^{f-2} , and τ_o^{f-2} for the testing task $f \in \mathcal{T}'$. This result is the counterpart of the requirement (19) for the CoT generalization, indicating that more context examples with the same TSR pattern as the query and more accurate context examples improve ICL generalization.*

Ref. Li et al. (2023c) also shows the advantage of CoT over ICL to learn MLP functions, but in a different setting from ours, where our studied tasks operate on patterns. More importantly, this paper characterizes the CoT and ICL performance theoretically when the testing task has a distribution shift from training tasks (TRR patterns to TSR patterns), and the testing examples contain errors, while Li et al. (2023c) only empirically evaluates the CoT and ICL performance with noisy examples.

4 THE MECHANISM OF CoT AND THE PROOF SKETCH

4.1 TRANSFORMERS IMPLEMENT CoT BY ATTENDING TO THE MOST SIMILAR EXAMPLES EVERY STEP

We characterize the key mechanism of a properly trained one-layer Transformer to implement CoT on a K -steps reasoning task via training dynamics analysis of the attention layer, as demonstrated in Figure 2. This is different from the mechanism study in (Li et al., 2023c; Feng et al., 2023) by constructing a model that can conduct CoT. We have the following proposition for the trained model.

Proposition 1. *Let S_k^* denote the index set of the context columns of the testing prompt \mathbf{P} in (4) that (a) correspond to the k -th step in a context example and (b) share the*

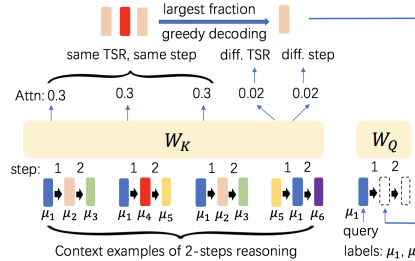


Figure 2: Concentration of attention weights for CoT inference.

⁵Our insight of the comparison between CoT and ICL still holds when we evaluate CoT generalization only by the final step output. This is because a successful CoT generalization in Theorem 2 on all reasoning steps already ensures a satisfactory CoT generalization on the final step.

same TSR pattern in the k -th input as the k -th input \mathbf{v}_{k-1} of the query, $k \in [K]$. Given a trained model that satisfies conditions (i) to (iii) of Theorem 1 and (18) and (19) after T iterations, we have

$$\sum_{i \in \mathcal{S}_k^*} \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(T)} \tilde{\mathbf{q}}_k) \geq 1 - \epsilon, \text{ where } \tilde{\mathbf{p}}_i = \mathbf{p}_i + \mathbf{c}_{(i \bmod K)}, \tilde{\mathbf{q}}_k = \mathbf{q}_k + \mathbf{c}_k, \quad (22)$$

with \mathbf{q}_k defined in (6). Moreover, for any $f \in \mathcal{T}'$, the k -th step output \mathbf{v}_k given $\mathbf{x}_{\text{query}} = \boldsymbol{\mu}_i^f$ satisfies,

$$\mathbf{v}_k = f_k \circ \dots \circ f_1(\boldsymbol{\mu}_i^f). \quad (23)$$

Proposition 1 first illustrates that, when conducting the k -th step reasoning of the query for any $k \in [K']$, the trained model assigns dominant attention weights on the prompt columns that are also the k -th step reasoning of examples and share the same TSR pattern in the k -th step input as the query. Then, given a sufficient number of testing context examples by (19), it is ensured that the fraction of the correct TSR pattern is the largest in the output of each step by (11). Subsequently, the generation by greedy decoding (5) is correct in each step, leading to a successful CoT generalization.

4.2 AN OVERVIEW OF THE PROOF

The technical challenges of the proof are concentrated on Theorem 1, where the property of the trained model is derived. The proof of Theorem 1 is built upon three Lemmas, which characterize the **two stages of the training dynamics**, i.e., Transformers first attend to tokens with the same step as the query and then, among them, further concentrate on tokens that share the same TSR pattern as the query. Specifically, Lemmas 3 and 4 show that if a training prompt \mathbf{P} includes the first k steps of the reasoning query, then the attention weights on columns of \mathbf{P} with a different step from the query decrease to be close to zero in the first stage. Lemma 5 computes the gradient updates in the second stage, where the attention weights on columns in \mathbf{P} that correspond to the same step and have the same TRR pattern as the query gradually become dominant. Theorem 1 unveils this training process by showing the required number of training iterations and sample complexity.

To prove Theorem 2, we first compute the required number of context examples for the new task $f \in \mathcal{T}'$ so that by concentration inequalities, the number of context examples with accurate TSR is larger than examples with inaccurate TSR patterns in all K reasoning steps with high probability. Then, by the correlation between TRR and TSR patterns (18), we also show that the trained Transformer can attend to context columns with the same TSR pattern as the query. Therefore, the model can make the correct generation in each step. Theorem 3 follows a similar proof idea to Theorem 2, with the difference that the trained model predicts output directly from the input query following \mathbf{B}^f instead of $\mathbf{A}_k^f, k \in [K]$ in CoT. Therefore, Condition 1 is required for the success of ICL generalization.

5 NUMERICAL EXPERIMENTS

Data Generation and Model setup. We use synthetic data generated following Sections 2 and 3.2. Let $d_{\mathcal{X}} = 30, M = 20, M' = 10, \alpha = 0.4$. We consider 3-steps tasks for training and testing, i.e., $K = 3$. A reasoning task f is generated by first sampling a set of numbers of permutations $\{p_i\}_{i=1}^M$ with $p_i \in [M]$ and then let $f_k(\boldsymbol{\mu}_{p_i}) = \boldsymbol{\mu}_{p((i+k) \bmod M)}$ for $i \in [M], k, j \in [K]$. The testing noise level is set to be 0.2 for any examples and $f \in \mathcal{T}'$. The learning model is a one-layer single-head Transformer defined in (2) or a three-layer two-head Transformer. We set $\tau^f = 0.5, \rho^f = 0.8, \alpha' = 0.8$ for CoT testing if not otherwise specified.

Experiments on the generalization of CoT. We first verify the required number of context examples for a desired CoT generalization on a one-layer Transformer. We investigate the impact of α', τ^f , and ρ^f by varying one and fixing the other two. Figure 3 illustrates that more testing examples are needed when α', τ^f , or ρ^f is small, which verifies the trend of the lower bound of l_{ts}^f in (19).

Experiments on the generalization of ICL and a comparison with CoT. We then verify the ICL generalization with the trained model. We vary τ_o^f and ρ_o^f by changing τ^f and ρ^f . Figure 3 indicates that more testing examples are required when α', τ_o^f , or ρ_o^f is small, which is consistent with our bound in (21). We then consider the case where $\tau_o^f = 0.4$ and $\rho_o^f = 0.1$ so that the generated testing prompt may not satisfy Condition 1 depending on the specific choices of \mathbf{A}_k^f 's. Figure 5 shows that when Condition 1 holds, the ICL testing error decreases if the number of contexts increases. However, when Condition 1 fails, the ICL testing error remains large, irrespective of the number of contexts.

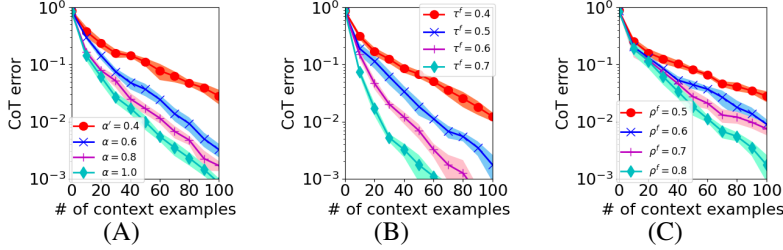


Figure 3: CoT testing error with different (A) α' (B) τ^f (C) ρ^f .

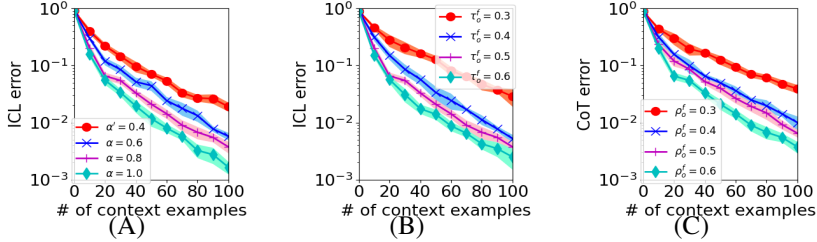


Figure 4: ICL testing error with different (A) α' (B) τ_o^f (C) ρ_o^f .

Experiments on the training dynamics of CoT.

In Figure 6, we compute the total attention weights on four types of testing context columns along the training, which are contexts with the same (or different) TSR pattern and in the same (or different) step as the query. The result shows that the attention weights on contexts that share the same TSR pattern and in the same step as the query increase along the training and converge to around 1. This verifies the mechanism formulated in (22). Meanwhile, Figure 6 also justifies the two-stage training dynamics proposed in Section 4.2, where we add a black vertical dashed line to demonstrate the stage transition boundary. We observe that the attention weights on context columns with a different step, i.e., the red and yellow curves, decrease to zero in the first stage. Then, the attention weights on contexts with the same TSR pattern and the same step as the query, i.e., the blue curve, increase to 1 in the second stage. We also justify the attention mechanism of CoT on a three-layer two-head Transformer with a two-step reasoning task. Figure 7 shows that there exists at least one head in each layer of the Transformer that implements CoT as characterized in Proposition 1. This indicates that the CoT mechanism we characterize on one-layer Transformers can be extended to multi-layer multi-head Transformers.

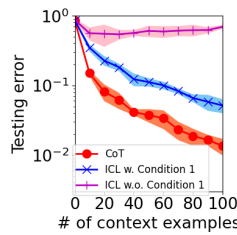


Figure 5: Comparison between CoT and ICL w./w.o. Condition 1

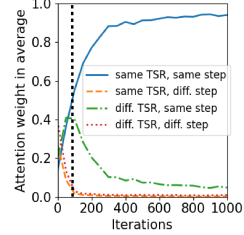


Figure 6: Training dynamics of Transformers for CoT

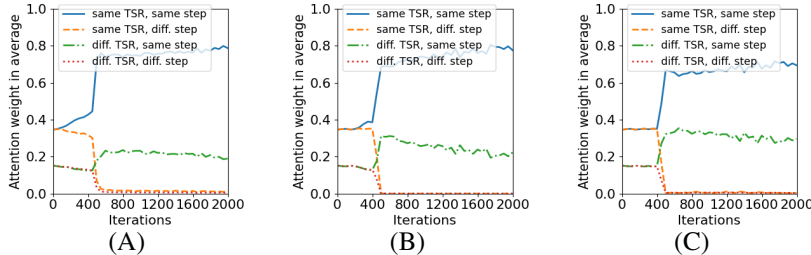


Figure 7: Training dynamics of Transformers. (A) Layer 1, Head 2 (B) Layer 2 Head 2 (C) Layer 3 Head 2.

6 CONCLUSION, LIMITATIONS, AND FUTURE WORKS

This paper theoretically analyzes the training dynamics of Transformers with nonlinear attention, together with the CoT generalization ability of the resulting model on new tasks with noisy and partially inaccurate context examples. We quantitatively characterize and compare the required conditions for the success of CoT and ICL. Although based on a simplified Transformer model and reasoning tasks operating on patterns, this work deepens the theoretical understanding of the CoT mechanism. Future directions include designing efficient prompt-generating methods for CoT and analyzing LLM reasoning on a more complicated data model.

ACKNOWLEDGMENTS

This work was supported by National Science Foundation(NSF) #2430223, Army Research Office (ARO) W911NF-25-1-0020, and the Rensselaer-IBM Future of Computing Research Collaboration (<http://airc.rpi.edu>). We also thank all anonymous reviewers for their constructive comments.

REFERENCES

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *arXiv preprint arXiv:2306.00297*, 2023.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 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.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- Sitan Chen and Yuanzhi Li. Provably learning a multi-head attention layer. *arXiv preprint arXiv:2402.04084*, 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*, 2024.
- Xiang Cheng, Yuxin Chen, and Suvrit Sra. Transformers implement functional gradient descent to learn non-linear functions in context. *arXiv preprint arXiv:2312.06528*, 2023.
- Yingqian Cui, Jie Ren, Pengfei He, Jiliang Tang, and Yue Xing. Superiority of multi-head attention in in-context linear regression. *arXiv preprint arXiv:2401.17426*, 2024.
- Puneesh Deora, Rouzbeh Ghaderi, Hossein Taheri, and Christos Thrampoulidis. On the optimization and generalization of multi-head attention. *arXiv preprint arXiv:2310.12680*, 2023.
- Nan Ding, Tomer Levinboim, Jialin Wu, Sebastian Goodman, and Radu Soricut. CausalLM is not optimal for in-context learning. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=guRNebwZBb>.
- 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, 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.
- Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. How do transformers learn in-context beyond simple functions? a case study on learning with representations. *arXiv preprint arXiv:2310.10616*, 2023.
- Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. In *NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning*, 2023.
- Yu Huang, Zixin Wen, Yuejie Chi, and Yingbin Liang. Transformers provably learn feature-position correlations in masked image modeling. *arXiv preprint arXiv:2403.02233*, 2024.
- M Emrullah Ildiz, Yixiao Huang, Yingcong Li, Ankit Singh Rawat, and Samet Oymak. From self-attention to markov models: Unveiling the dynamics of generative transformers. *arXiv preprint arXiv:2402.13512*, 2024.

- Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. *Advances in Neural Information Processing Systems*, 35:37822–37836, 2022.
- Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. In *The Eleventh International Conference on Learning Representations*, 2023a. URL <https://openreview.net/forum?id=jC1Gv3Qjhb>.
- Hongkang Li, Meng Wang, Tengfei Ma, Sijia Liu, ZAI XI ZHANG, and Pin-Yu Chen. What improves the generalization of graph transformer? a theoretical dive into self-attention and positional encoding. In *NeurIPS 2023 Workshop: New Frontiers in Graph Learning*, 2023b. URL <https://openreview.net/forum?id=BaxFC3z9R6>.
- Hongkang Li, Meng Wang, Songtao Lu, Xiaodong Cui, and Pin-Yu Chen. How do nonlinear transformers learn and generalize in in-context learning? In *Forty-first International Conference on Machine Learning*, 2024a. URL <https://openreview.net/forum?id=I4HTPws9P6>.
- Hongkang Li, Meng Wang, Shuai Zhang, Sijia Liu, and Pin-Yu Chen. Learning on transformers is provable low-rank and sparse: A one-layer analysis. *arXiv preprint arXiv:2406.17167*, 2024b.
- Yingcong Li, Kartik Sreenivasan, Angeliki Giannou, Dimitris Papailiopoulos, and Samet Oymak. Dissecting chain-of-thought: Compositionality through in-context filtering and learning. *Advances in Neural Information Processing Systems*, 36, 2023c.
- 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, 2024c.
- Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: Towards a mechanistic understanding. *arXiv preprint arXiv:2303.04245*, 2023d.
- Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. In *The Twelfth International Conference on Learning Representations*, 2024d. URL <https://openreview.net/forum?id=3EWTEy9MTM>.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=v8L0pN6EOi>.
- Yixin Liu, Kai Zhang, Yuan Li, Zhiling Yan, Chujie Gao, Ruoxi Chen, Zhengqing Yuan, Yue Huang, Hanchi Sun, Jianfeng Gao, et al. Sora: A review on background, technology, limitations, and opportunities of large vision models. *arXiv preprint arXiv:2402.17177*, 2024.
- Yuankai Luo. Transformers for capturing multi-level graph structure using hierarchical distances. *arXiv preprint arXiv:2308.11129*, 2023.
- Ashok Vardhan Makkuva, Marco Bondaschi, Chanakya Ekbote, Adway Girish, Alliot Nagle, Hyeji Kim, and Michael Gastpar. Local to global: Learning dynamics and effect of initialization for transformers. *arXiv preprint arXiv:2406.03072*, 2024a.
- 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*, 2024b.
- William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=NjNGlPh8Wh>.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with gradient descent. *arXiv preprint arXiv:2402.14735*, 2024.

- OpenAI. Gpt-4 technical report. *OpenAI*, 2023.
- Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the role of attention in prompt-tuning. *arXiv preprint arXiv:2306.03435*, 2023.
- Jongho Park, Jaeseung Park, Zheyang Xiong, Nayoung Lee, Jaewoong Cho, Samet Oymak, Kangwook Lee, and Dimitris Papailiopoulos. Can mamba learn how to learn? a comparative study on in-context learning tasks. *arXiv preprint arXiv:2402.04248*, 2024.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines. *arXiv preprint arXiv:2308.16898*, 2023a.
- Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism. *CoRR*, 2023b.
- 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, Zhenyu Zhang, Beidi Chen, and Simon Shaolei Du. Joma: Demystifying multilayer transformers via joint dynamics of mlp and attention. In *Conference on Parsimony and Learning (Recent Spotlight Track)*, 2023b.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- 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.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2609–2634, 2023.
- Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting. *arXiv preprint arXiv:2402.10200*, 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.
- Kaiyue Wen, Xingyu Dang, and Kaifeng Lyu. Rnns are not transformers (yet): The key bottleneck on in-context retrieval. *arXiv preprint arXiv:2402.18510*, 2024.
- Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett. How many pretraining tasks are needed for in-context learning of linear regression? *arXiv preprint arXiv:2310.08391*, 2023.

Kai Yang, Jan Ackermann, Zhenyu He, Guhao Feng, Bohang Zhang, Yunzhen Feng, Qiwei Ye, Di He, and Liwei Wang. Do efficient transformers really save computation? *arXiv preprint arXiv:2402.13934*, 2024.

Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*, 2023a.

Yihua Zhang, Hongkang Li, Yuguang Yao, Aochuan Chen, Shuai Zhang, Pin-Yu Chen, Meng Wang, and Sijia Liu. Visual prompting reimaged: The power of activation prompts, 2024. URL <https://openreview.net/forum?id=0b328CMwn1>.

Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations*, 2023b.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, et al. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*, 2023.

APPENDIX

A EXPERIMENTS ON REAL-WORLD DATA

We consider a simple arithmetic task that outputs $((A_1 o_1 A_2) o_2 A_3) o_3 A_4$ given A_1, A_2, A_3, A_4 chosen from integers from 0 to 9 as the input, where $o_1, o_2, o_3 \in O = \{+, -, \times\}$. The CoT output follows the format of $A_1 o_1 A_2 = S_1, S_1 o_2 A_3 = S_2, S_2 o_3 A_4 = S_3$ and will be evaluated by whether all the three steps are correct for the query as (7). ICL directly outputs S_3 , and the performance is evaluated by the prediction accuracy of S_3 as (9). In the following experimental settings, the accuracy is computed on 50 prompts. Each prompt contains three context examples. The inference model is GPT-4 (OpenAI, 2023).

An increasing number of erroneous examples hurts the CoT generalization. To model the errors in the context examples in the testing prompt, we replace o_3 with one operation \hat{o}_3 from $O \setminus o_3$ in the presentation of some of the context examples in the testing prompt. Note that the output values S_3 are still correctly computed from $S_3 = S_2 o_3 A_4$. Table 1 shows that when the total number of testing examples is fixed to be three, with the increasing number of incorrect examples, the testing accuracy decreases. This is consistent with Remark 1 for Theorem 2.

# of incorrect examples	0	1	2	3
CoT accuracy	100%	100%	56%	0%

Table 1: The accuracy with different numbers of incorrect examples for CoT. Errors in presenting o_3 .

CoT is more robust to erroneous examples with implementation error than ICL. In this setting, the error in a context examples is introduced by replacing o_1 with one operation \hat{o}_1 randomly and independently selected from $O \setminus o_1$. Hence, $S_1 = A_1 \hat{o}_1 A_2$, and the successive computation are based on the wrongly computed S_1 . The results in Table 2 shows that when two incorrect examples exist, CoT performs better than ICL, which justifies Remark 2 for Theorem 3.

# of incorrect examples	0	1	2
CoT accuracy	100%	100%	100%
ICL accuracy	100%	100%	60%

Table 2: The accuracy with different numbers of incorrect examples for CoT and ICL. Errors in implementing o_1 .

B ADDITIONAL DISCUSSIONS

B.1 THE MOTIVATION TO STUDY ONE-LAYER SINGLE-HEAD TRANSFORMERS

The reasons we study one-layer single-head attention-only nonlinear Transformers in this work are as follows.

First, it is much more challenging to theoretically analyze the training dynamics and generalization of multi-layer/head Transformers. This is because the loss landscape for multi-layer/head Transformers is highly nonlinear and non-convex due to the interactions between multiple nonlinear functions. The simplified data helps to characterize the gradient updates in different directions for different patterns and steps. Non-orthogonal data make the updates less separable for different inputs, which is more challenging to analyze.

Second, the state-of-the-art theoretical works (Li et al., 2023a; 2024a; Huang et al., 2023; Makuva et al., 2024a; Ildiz et al., 2024) on optimization and generalization also focus mainly on one-layer Transformers. No existing works study the optimization and generalization of CoT even for one-layer Transformers. Therefore, we plan to focus on the one-layer analysis to obtain more theoretical insights. We leave the theoretical analysis of the multi-layer case as future works.

Third, although we admit the gap between theory and practice, our theory still makes contributions under our settings. Our work is the first one to investigate the optimization and generalization of CoT and characterize the conditions when CoT is better than ICL. We establish the required number of context examples for a successful CoT in terms of how informative and erroneous the prompt is.

We also implement experiments on the attention mechanism for three-layer two-head Transformers on two-step reasoning tasks. Please see Figure 7 for details. The findings of all three layers are generally consistent with Proposition 1 for the single-layer single-head case, which indicates that the CoT mechanism we characterize on one-layer Transformers can be extended to multi-layer multi-head Transformers.

B.2 THE MOTIVATION OF THE DATA AND TASK FORMULATION

There are several reasons for using such data formulation.

First, our data formulation of orthogonal patterns, on which the function is based, is widely used in the state-of-the-art theoretical study of model training or ICL on language and sequential data [(Tian et al., 2023a; Huang et al., 2023; Li et al., 2024a; Chen et al., 2024)]. For example, (Huang et al., 2023; Li et al., 2024a) study ICL on regression or classification tasks, which also use orthogonal patterns as data. Sections 2.1 and 2.2 in (Chen et al., 2024) consider learning n-gram data in ICL by formulating transitions between orthogonal patterns. Section 3 of (Tian et al., 2023a) also assume orthogonal patterns in Transformer model training, and the generation comes from the orthogonal pattern set. The data formulation we use is consistent with the existing theoretical works.

Second, based on this formulation, one can characterize the gradient updates in different directions for different patterns and steps. This enables us to distinguish the impact of different patterns and steps in the convergence analysis of CoT using Transformers. Non-orthogonal data make the model updates less separable for different inputs, which is more challenging to analyze. Moreover, we would like to mention that during the inference, the tokens in testing prompts contain noises as defined in Equation 10. This makes the tokens of different TSR patterns not orthogonal to each other and relaxes our orthogonality condition to some degree.

B.3 THE DISCUSSION OF POSITIONAL ENCODING

The positional encoding (PE) we use is simplified for theoretical analysis. The formulation of PE we use is motivated by (Huang et al., 2024; Nichani et al., 2024), where each token is added with a PE represented by orthogonal vectors. These works formulate the distribution of the PE to be related to the structure of the data, such as patch-wise association (Huang et al., 2024), and sparse token selection (Nichani et al., 2024). Likewise, we follow their intuition to make the PE vary in different steps of our reasoning tasks so that the Transformer can distinguish different steps when making inferences for the query.

Our analysis can be extended to study more general PEs with additional technical work in the future. One possible direction is studying the family of periodic and separable PE. For example, the absolute PE proposed by (Vaswani et al., 2017) considers PE as a sinusoid, which is periodic. Such analysis can be made by relaxing the “orthogonality” of PE vectors to a certain “separability” between PE vectors.

We also conduct experiments on a three-layer single-head Transformer with the standard PE proposed in Section 3.5 of (Vaswani et al., 2017) for our problem. Figure shows that the blue curve increases to be the largest along the training, which means the attention weights on example steps that share the same TSR pattern and the same step as the query. This indicates that the CoT mechanism of using standard PE is the same as the one proposed in Proposition 1 in our paper. One might note that the scores of the blue curve are not as high as Figure 6 in our paper. We guess the reason why the distinction in attention values is more significant in our PE may be the additional orthogonality of our PE and the property that its period is the same as the reasoning length. Nevertheless, the strong similarity between the results on standard PE and our used PE shows the practical significance of our analysis.

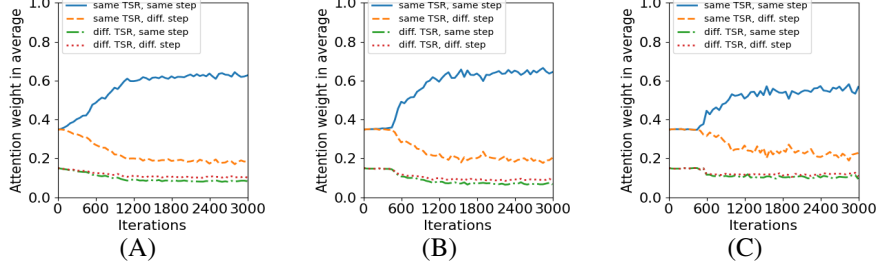


Figure 8: CoT mechanism with standard PE of (A) Layer 1 (B) Layer 2 (C) Layer 3.

C ALGORITHMS

We first present the training algorithm introduced in Section 2.2.

Algorithm 1 Training with Stochastic Gradient Descent (SGD)

- 1: **Hyperparameters:** The step size η , the number of iterations T , batch size B .
- 2: **Initialization:** Let $\mathbf{W} = \mathbf{W}_K^\top \mathbf{W}_Q$ and $\mathbf{W}_V = (\mathbf{0}_{d_X \times d_X} \quad \mathbf{I}_{d_X} \quad \mathbf{0}_{d_X \times d_\varepsilon})$. Each entry of $\mathbf{W}^{(0)}$ is generated from $\mathcal{N}(0, \xi^2)$ for a small constant $\xi > 0$. \mathbf{W}_V and \mathbf{a} are fixed during the training.
- 3: **Training by SGD:** For each iteration, we independently sample $\mathbf{x}_{query} \sim \mathcal{D}$, $f \in \mathcal{T}_{tr}$ to form a batch of training prompt and labels $\{\mathbf{P}^n, z^n\}_{n \in \mathcal{B}_t}$ as introduced in Section 3.2. Each TRR pattern is sampled equally likely in each batch. For each $t = 0, 1, \dots, T - 1$

$$\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_t} \nabla_{\mathbf{W}^{(t)}} \ell(\Psi^{(t)}; \mathbf{P}^n, z^n). \quad (24)$$

- 4: **Output:** $\mathbf{W}^{(T)}$.
-

We then summarize the algorithm of the CoT inference introduced in Section 2.3 as follows.

Algorithm 2 Inference with Chain-of-Thought (CoT)

- 1: **Input:** $z_0 = v_0 = \mathbf{x}_{query}$, \mathbf{P}_0 , and \mathbf{P}_1 .
- 2: **for** $k = 1, \dots, K - 1$, **do**

$$\text{Compute } v_k \text{ by greedy decoding in (5). Then update } \mathbf{P}_k \text{ and } \mathbf{P}_{k+1} \text{ by (6).} \quad (25)$$

- 3: **end for**
 - 4: **Output:** v_1, v_2, \dots, v_{K-1} , and v_K by (5).
-

D PRELIMINARIES

We first summarize the notations we use in this paper in Table 3.

Lemma 1 (Multiplicative Chernoff bounds, Theorem D.4 of (Mohri et al., 2018)). *Let X_1, \dots, X_m be independent random variables drawn according to some distribution \mathcal{D} with mean p and support included in $[0, 1]$. Then, for any $\gamma \in [0, \frac{1}{p} - 1]$, the following inequality holds for $\hat{p} = \frac{1}{m} \sum_{i=1}^m X_i$:*

$$\Pr(\hat{p} \geq (1 + \gamma)p) \leq e^{-\frac{m p \gamma^2}{3}}, \quad (26)$$

$$\Pr(\hat{p} \leq (1 - \gamma)p) \leq e^{-\frac{m p \gamma^2}{2}}. \quad (27)$$

Definition 2 ((Vershynin, 2010)). *We say X is a sub-Gaussian random variable with sub-Gaussian norm $K > 0$, if $(\mathbb{E}|X|^p)^{\frac{1}{p}} \leq K \sqrt{p}$ for all $p \geq 1$. In addition, the sub-Gaussian norm of X , denoted $\|X\|_{\psi_2}$, is defined as $\|X\|_{\psi_2} = \sup_{p \geq 1} p^{-\frac{1}{2}} (\mathbb{E}|X|^p)^{\frac{1}{p}}$.*

Lemma 2 (Vershynin (2010) Proposition 5.1, Hoeffding’s inequality). *Let X_1, X_2, \dots, X_N be independent centered sub-gaussian random variables, and let $K = \max_i \|X_i\|_{\psi_2}$. Then for every $\mathbf{a} = (a_1, \dots, a_N) \in \mathbb{R}^N$ and every $t \geq 0$, we have*

$$\Pr\left(\left|\sum_{i=1}^N a_i X_i\right| \geq t\right) \leq e \cdot \exp\left(-\frac{ct^2}{K^2 \|\mathbf{a}\|^2}\right), \quad (28)$$

Table 3: Summary of Notations

Notations	Annotation
$\mathbf{x}_i, \mathbf{y}_{i,k}, \mathbf{x}_{query}, \mathbf{z}_k$	\mathbf{x}_i is the input to the first step of a reasoning example. $\mathbf{y}_{i,k}$ is the k -th step output label of \mathbf{x}_i . \mathbf{x}_{query} is the query input. \mathbf{z}_k the k -th step output label of \mathbf{x}_{query} . $k \in [K]$.
$\mathbf{P}, \mathbf{p}_{query}, \mathbf{E}_i, \mathbf{Q}_k, \mathbf{v}_k$	\mathbf{P} is a training or testing prompt that consists of multiple training or testing examples and a query. The last column of \mathbf{P} is denoted by \mathbf{p}_{query}^n , which is the query of \mathbf{P} . \mathbf{E}_i is the i -th context example of \mathbf{P} . \mathbf{Q}_k is the first k steps of the reasoning query. $k \in [K]$. \mathbf{v}_k is the k -th step generation by CoT. $k \in [K]$.
$\mathbf{c}_i, \tilde{\mathbf{p}}_i, \tilde{\mathbf{p}}_{query}$	\mathbf{c}_i is the positional encoding for the i -th column of the input sequence. $\tilde{\mathbf{p}}_i = \mathbf{p}_i + \mathbf{c}_i$, where \mathbf{p}_i is the i -th column of \mathbf{P} . $\tilde{\mathbf{p}}_{query}$ is the \mathbf{p}_i of the query column.
$F(\Psi; \mathbf{P}), \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)$	$F(\Psi; \mathbf{P}^n)$ is the Transformer output for \mathbf{P} with Ψ as the parameter. $\ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)$ is the loss function value given \mathbf{P}^n and the corresponding label \mathbf{z}^n .
$\boldsymbol{\mu}_i \in \mathcal{M}, \boldsymbol{\mu}'_i \in \mathcal{M}', \text{TSR}(\cdot)$	$\boldsymbol{\mu}_i$ is the i -th training-relevant (TRR) pattern for $i \in [M]$. $\boldsymbol{\mu}'_i$ is the i -th testing-relevant (TSR) pattern for $i \in [M']$. \mathcal{M} and \mathcal{M}' are the set of TRR and TSR patterns, respectively. $\text{TSR}(\cdot)$ is a function that outputs the index of the TSR pattern of the noisy input.
f_k, f	f is the task function with $f = f_K \circ \dots \circ f_2 \circ f_1$ for a K -steps reasoning. f_k is the k -th step task function.
$\mathcal{T}, \mathcal{T}', \mathcal{D}, \mathcal{D}'$	\mathcal{T} is the distribution of training tasks, while \mathcal{T}' is the distribution of testing tasks. \mathcal{D} is the training data distribution. \mathcal{D}' is the testing data distribution.
α, α'	α (or α') is the fraction of context examples with input sharing the same TRR (or TSR) pattern as the query.
$\mathbf{A}_k^f, \mathbf{B}_k^f$	\mathbf{A}_k^f is the step-wise transition matrix at the k -th step for the task f , $k \in [K]$. \mathbf{B}_k^f is the K -steps transition matrix of the task f .
$\tau^f, \tau_o^f, \rho^f, \rho_o^f$	τ^f is the min-max trajectory transition probability for task f . τ_o^f is the min-max input-label transition probability for task f . ρ^f and ρ_o^f are primacy of the step-wise transition matrices and the K -steps transition matrix, respectively.
S_k^*	The index set of context columns of the prompt that correspond to the k -th step of the example and share the same TSR pattern in the $(k-1)$ -th output as the $(k-1)$ -th output \mathbf{v}_{k-1} of the query.
$p_n(t)$	$p_n(t)$ is the summation of attention weights on context columns that share the same TRR/TSR pattern and in the same step as the query.
B_b	B_b is the SGD batch at the b -th iteration.
l_{tr}	l_{tr} is the universal number of training context examples.
l_{ts}^f	l_{ts}^f is the number of testing context examples of the task f .
$\mathcal{O}(), \Omega(), \Theta()$	We follow the convention that $f(x) = \mathcal{O}(g(x))$ (or $\Omega(g(x))$, $\Theta(g(x))$) means that $f(x)$ increases at most, at least, or in the order of $g(x)$, respectively.
\gtrsim, \lesssim	$f(x) \gtrsim g(x)$ (or $f(x) \lesssim g(x)$) means that $f(x) \geq \Omega(g(x))$ (or $f(x) \lesssim \mathcal{O}(g(x))$).

where $c > 0$ is an absolute constant.

Definition 3. Define that for $\tilde{\mathbf{p}}_i$ that shares the same TRR/TSR pattern and in the same step as the query,

$$p_n(t) = \sum_i \text{softmax}(\tilde{\mathbf{p}}_i^n \top \mathbf{W}^{(t)} \tilde{\mathbf{p}}_{query}^n). \quad (29)$$

Lemma 3. Given the SGD training scheme described in Section 2.2, $B \geq \Omega(M \log M)$, and $l_{tr} \geq \Omega(\alpha^{-1})$, we have the following results. When $\mathcal{O}(\eta^{-1} \alpha^{-2} K^3 \log \frac{K}{\epsilon}) \geq t \geq 1$, for any \mathbf{p} as a

column of context examples in (1), we have

$$\begin{aligned} & \tilde{\mathbf{p}}^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \\ & \leq \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (1 - p_n(t))^2 (-4p_n(t)(1 + \frac{\alpha^2}{K^2}) + \frac{\alpha^2}{K^2} (1 + \frac{2(K-1)}{K})) - \frac{\alpha^2}{K^3} (1 - p_n(t))^2 \right). \end{aligned} \quad (30)$$

For any $\tilde{\mathbf{p}}'$ that shares the same TRR pattern and a different positional encoding as $\tilde{\mathbf{p}}$, we have

$$\begin{aligned} & \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (-4 - (3K-2)(1 - p_n(t))(1 + \frac{\alpha^2}{K^2})) p_n(t)(1 - p_n(t)) + \frac{\alpha^2}{K^3} (1 - p_n(t))^2 \right) \\ & \leq \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \\ & \leq \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (-4 - (3K-2)(1 - p_n(t))(1 + \frac{\alpha^2}{K^2})) p_n(t)(1 - p_n(t)) + \frac{1}{K} p_n(t)(1 - p_n(t))^2 \right. \\ & \quad \left. \cdot (1 + \frac{\alpha^2}{K^2}) \right). \end{aligned} \quad (31)$$

For any $\tilde{\mathbf{p}}'$ that shares a different TRR pattern but the same positional encoding as $\tilde{\mathbf{p}}$, we have

$$\begin{aligned} & \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} \left(-\frac{\alpha^2}{K^2} + (K-1 + \frac{(2K-1)\alpha^2}{K^2}) p_n(t) \right) (1 - p_n(t))^2 - (1 - p_n(t))^2 \frac{\alpha^2}{K^3} \right. \\ & \quad \left. + \frac{1}{K} \cdot (1 - p_n(t))^2 (-p_n(t) + (1 - p_n(t)) \frac{\alpha^2}{K^2}) \right) \\ & \leq \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \\ & \leq \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} \left(-\frac{\alpha^2}{K^2} + (K-1 + \frac{(2K-1)\alpha^2}{K^2}) p_n(t) \right) (1 - p_n(t))^2 - (1 - p_n(t))^2 \frac{\alpha^2}{K^3} \right). \end{aligned} \quad (32)$$

For any $\tilde{\mathbf{p}}'$ that shares a different TRR pattern and a different positional encoding from $\tilde{\mathbf{p}}$, we have

$$\begin{aligned} & \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} p_n(t)(1 - p_n(t))^2 \left(1 + \frac{(2-K)\alpha^2}{K^2} \right) + (1 - p_n(t))^2 \cdot \frac{\alpha^2}{K^3} \right) \\ & \leq \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \\ & \leq \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} p_n(t)(1 - p_n(t))^2 \left(2 - K + \frac{(2-K)\alpha^2}{K^2} \right) + (1 - p_n(t))^2 p_n(t) \left(1 + \frac{\alpha^2}{K^2} \right) \cdot \frac{1}{K} \right). \end{aligned} \quad (33)$$

Lemma 4. Given the SGD training scheme described in Section 2.2, $B \geq \Omega(M \log M)$, and $l_{tr} \geq \Omega(\alpha^{-1})$, and

$$t \gtrsim T_1 := \eta^{-1} \alpha^{-2} K^3 \log \frac{K}{\epsilon}, \quad (34)$$

we have that if \mathbf{p}_{query} is in the k -th step,

$$\sum_{i \in \mathcal{S}_{[K] \setminus k}} \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(t)} \tilde{\mathbf{p}}_{query}) \leq \epsilon \quad (35)$$

where $\mathcal{S}_{[K] \setminus k}$ means the index set of context columns that are not in the k -th step.

Lemma 5. Given the SGD training scheme described in Section 2.2, $B \geq \Omega(M \log M)$, and $l_{tr} \geq \Omega(\alpha^{-1})$, we have the following results. When $t \geq T_1 = \eta^{-1} \alpha^{-2} K^3 \log \frac{K}{\epsilon}$, for any \mathbf{p} as a

column of context examples in (1), we have

$$\tilde{\mathbf{p}}^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \leq -\frac{\eta}{2MB} \sum_{n \in \mathcal{B}_b} 4p_n(t)(1-p_n(t))^2. \quad (36)$$

For any $\tilde{\mathbf{p}}'$ that shares the same TRR pattern and a different positional encoding as $\tilde{\mathbf{p}}$, we have

$$\left| \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \right| \leq \eta \epsilon. \quad (37)$$

For any $\tilde{\mathbf{p}}'$ that shares a different TRR pattern but the same positional encoding as $\tilde{\mathbf{p}}$, we have

$$\left| \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \right| \leq \frac{\eta}{2BM} \sum_{n \in \mathcal{B}_b} p_n(b)(1-p_n(b))^2. \quad (38)$$

For any $\tilde{\mathbf{p}}'$ that shares a different TRR pattern and a different positional encoding from $\tilde{\mathbf{p}}$, we have

$$\left| \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}} \tilde{\mathbf{p}} \right| \leq \eta \epsilon. \quad (39)$$

E PROOF OF MAIN THEOREMS

E.1 PROOF OF THEOREM 1

Proof. By the condition in Lemma 3, we have that

$$B \geq \Omega(M \log M). \quad (40)$$

We know that there exists gradient noise caused by imbalanced TRR patterns in each batch. Then, by Hoeffding's inequality (28),

$$\Pr \left(\left\| \frac{1}{|\mathcal{B}_b|} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} - \mathbb{E} \left[\frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \right] \right\| \geq \left| \mathbb{E} \left[\frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \right] \right| \epsilon \right) \leq e^{-B\epsilon^2} \leq M^{-C}, \quad (41)$$

if $B \gtrsim \epsilon^{-2} \log M$. Therefore, we require

$$B \gtrsim \max\{\epsilon^{-2}, M\} \log M. \quad (42)$$

By Lemma 5 and Definition 3, for $\tilde{\mathbf{p}}_i^n$ that share the same TRR pattern and the same positional encoding of $\tilde{\mathbf{p}}_{query}^n$,

$$\frac{p_n(t+1)}{|\mathcal{S}_1^n|} = \text{softmax}(\tilde{\mathbf{p}}_i^n^\top \mathbf{W}^{(t+1)} \tilde{\mathbf{p}}_{query}^n) \geq \frac{1}{l} \cdot \frac{1}{\frac{\alpha}{K} + (1 - \frac{1}{K}) \cdot \epsilon + (\frac{1}{K} - \frac{\alpha}{K}) e^{-u}}, \quad (43)$$

where by (161),

$$u \gtrsim \frac{\eta}{KM} \sum_{b=0}^t (1-p_n(b))^2 p_n(b). \quad (44)$$

For $\tilde{\mathbf{p}}_i^n$ that only share the same positional encoding of $\tilde{\mathbf{p}}_{query}^n$,

$$\text{softmax}(\tilde{\mathbf{p}}_i^n^\top \mathbf{W}^{(t+1)} \tilde{\mathbf{p}}_{query}^n) \geq \frac{1}{l} \cdot \frac{1}{\frac{\alpha}{K} e^u + (1 - \frac{1}{K}) \cdot \epsilon + (\frac{1}{K} - \frac{\alpha}{K})}. \quad (45)$$

Therefore, to make the attention weights between $\tilde{\mathbf{p}}_{query}^n$ and $\tilde{\mathbf{p}}_i^n$ that share the same TRR pattern and the same positional encoding dominant, we need a large enough u . When $1-p_n(b) \geq \Omega(1)$, we have

$$t \leq T_2 := \eta^{-1} KM \alpha^{-1}. \quad (46)$$

When $1 - p_n(b) \leq O(1)$,

$$p_n(t+1) = \frac{e^u}{e^u + \frac{1-\frac{\alpha}{K}}{\frac{\alpha}{K}}} \gtrsim 1 - \frac{1-\frac{\alpha}{K}}{\frac{\alpha}{K}} e^{-u}, \quad (47)$$

and

$$1 - p_n(t+1) \geq \frac{1-\frac{\alpha}{K}}{\frac{\alpha}{K} e^u + (1-\frac{\alpha}{K})} \gtrsim \frac{1-\frac{\alpha}{K}}{\frac{\alpha}{K}} e^{-u}. \quad (48)$$

Then, we prove that when t is large enough, $u(t) \geq \frac{1}{2} \log \frac{\eta(1-\alpha)^2 t}{\alpha^2 KM}$. We show it by induction. Suppose that the conclusion holds when $t = t_0$, then

$$\begin{aligned} u(t+1) &\geq \frac{\eta}{KM} \sum_{b=0}^{t_0} (1-p_n(b))^2 p_n(b) + \frac{\eta}{KM} (1-p_n(t))^2 p_n(t) \\ &\geq \frac{1}{2} \log \frac{(K-\alpha)^2 t}{2\alpha^2 KM} + \frac{\eta}{KM} (1-p_n(t))^2 p_n(t) \\ &\geq \frac{1}{2} \log \frac{\eta(K-\alpha)^2 (t+1)}{\alpha^2 KM}, \end{aligned} \quad (49)$$

where the last step is by

$$\frac{1}{2} \log\left(1 + \frac{1}{t}\right) \leq \frac{1}{2t} \leq \frac{\eta}{KM} \cdot \left(\frac{K-\alpha}{\alpha}\right)^2 e^{-\log \frac{\eta(K-\alpha)^2 t}{\alpha^2 KM}}. \quad (50)$$

To make $(1 - p_n(t))^2 < \epsilon$, we need

$$\left(\frac{K-\alpha}{\alpha}\right)^2 e^{-2u} \leq \epsilon. \quad (51)$$

Then, we get

$$u \geq \frac{1}{2} \log \frac{1}{\epsilon} + \log \frac{K-\alpha}{\alpha}. \quad (52)$$

Therefore, by

$$\frac{1}{2} \log \frac{\eta t}{KM} + \log \frac{K-\alpha}{\alpha} \geq \frac{1}{2} \log \frac{1}{\epsilon} + \log \frac{K-\alpha}{\alpha}, \quad (53)$$

we finally obtain

$$t \geq T_3 := \eta^{-1} \epsilon^{-1} KM. \quad (54)$$

For $\tilde{\mathbf{p}}_i^n$ that shares the same TSR pattern as the query, we have that when $t = T_1$,

$$\tilde{\mathbf{p}}_i^n \top \mathbf{W}^{(t)} \tilde{\mathbf{p}}_{query}^n \geq \log \frac{1}{\epsilon}. \quad (55)$$

When $t = T_1 + T_2 + T_3$,

$$\tilde{\mathbf{p}}_i^n \top \mathbf{W}^{(t)} \tilde{\mathbf{p}}_{query}^n \geq \Theta(1) \cdot \log \frac{1}{\epsilon} = \Theta(\log \frac{1}{\epsilon}). \quad (56)$$

Then,

$$\begin{aligned} T &:= T_1 + T_2 + T_3 \\ &= \Theta(\eta^{-1} \alpha^{-2} K^3 \log \frac{K}{\epsilon} + \eta^{-1} MK(\alpha^{-1} + \epsilon^{-1})). \end{aligned} \quad (57)$$

Therefore,

$$\mathbb{E}_{\mathbf{x}_{query} \sim \mathcal{D}, f \in \mathcal{T}} [\ell(\Psi; \mathbf{P}, \mathbf{z})] \leq \mathcal{O}(\epsilon). \quad (58)$$

□

E.2 PROOF OF THEOREM 2

Proof. We know that α' is the fraction of examples that share the same TSR pattern as the query. We need that in each step, the number of examples that share the same TSR pattern as the current step of the query is at least 1. Note that the probability of examples where each reasoning step produces the most probable output is

$$\prod_{k=1}^K A_k^f(\text{TSR}(f_{k-1} \circ \dots \circ f_0(\mu'_i)), \text{TSR}(f_k \circ \dots \circ f_0(\mu'_i))), \text{ where } f_0(\mu'_i) := \mu'_i, \forall i \in [M'], \quad (59)$$

where the input to the first step has the TSR pattern μ'_i . Define $m_{k(i)}$ as the TSR pattern in the k -th step output of the i -th context example by the transition matrix defined in 11. Consider that the TSR pattern of the k -th step label of the testing query is μ'_{q_k} , which is also the most probable k -th step output of the k -th step of a certain \mathbf{x}_i with $\text{TSR}(\mathbf{x}_i) = \text{TSR}(\mathbf{x}_{\text{query}}) = q_0$. Let the TSR pattern of another reasoning process, where for a certain first-step input \mathbf{x}_i with $\text{TSR}(\mathbf{x}) = \text{TSR}(\mathbf{x}_{\text{query}}) = q_0$, the k -th step output is the most probable for $k \in [K'] \setminus \{h\}$, while the h -th step output is the second probable. Denote the TSR pattern of the k -th step output of \mathbf{x}_i following this process as μ'_{u_k} with $u_0 = q_0$. By the Chernoff bound of Bernoulli distribution in Lemma 1, we can obtain

$$\begin{aligned} \Pr \left(\frac{1}{l_{ts}} \sum_{i=1}^{l_{ts}} \mathbb{1}[m_{k(i)} = \mu'_{q_k}, \forall k \in [K']] \leq (1 - \rho_s^f/2) \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f \right) \\ \leq e^{-l_{ts} (\rho_s^f)^2 \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f} = M^{-C}, \end{aligned} \quad (60)$$

and by Lemma 2,

$$\begin{aligned} \Pr \left(\frac{1}{l_{ts}} \sum_{i=1}^{l_{ts}} \mathbb{1}[m_{k(i)} = \mu'_{u_k}, \forall k \in [K']] \geq (1 - \rho_s^f/2) \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f \right) \\ \leq \Pr \left(\frac{1}{l_{ts}} \sum_{i=1}^{l_{ts}} \mathbb{1}[m_{k(i)} = \mu'_{u_k}, \forall k \in [K']] \geq \alpha' \prod_{k=1}^{K'} A_{k(u_{k-1}, u_k)}^f + t_0 \right) \\ \leq e^{-l_{ts} t_0^2} = M^{-C}, \end{aligned} \quad (61)$$

for some $c \in (0, 1)$ and $C > 0$, where the first step is by the definition of ρ_s^f in (11), and

$$t_0 \lesssim \rho_s^f \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f. \quad (62)$$

Hence, with a high probability,

$$\begin{aligned} l_{ts} \gtrsim \max \{ (\rho_s^f)^2 \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f \}^{-1} \log M, (\rho_s^f \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f)^{-2} \log M \} \\ \gtrsim (\rho_s^f \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f)^{-2} \log M, \end{aligned} \quad (63)$$

such that the number of examples with the same TSR pattern as the query in each of the total K steps is at least 1. To make the above condition hold for any TSR pattern of the intermediate step of the query, we need

$$\begin{aligned} l_{ts} \gtrsim \max_{q_k \in [M']} (\rho_s^f \alpha' \prod_{k=1}^{K'} A_{k(q_{k-1}, q_k)}^f)^{-2} \log M \\ = \max_{i \in [M']} (\rho_s^f \alpha' \prod_{k=1}^{K'} A_{k(\text{TSR}(f_{k-1} \circ \dots \circ f_0(\mu'_i)), \text{TSR}(f_k \circ \dots \circ f_0(\mu'_i)))})^{-2} \log M \\ = (\rho_s^f \alpha' \tau_s^f)^{-2} \log M. \end{aligned} \quad (64)$$

Then, we show the CoT testing error is zero by induction. In the first step, consider $\mathbf{x}_i = \boldsymbol{\mu}_j + \boldsymbol{\delta}_i$ such that

$$\tilde{\mathbf{p}}_i = \begin{pmatrix} \boldsymbol{\mu}'_j \\ \mathbf{y}_{i,1} \end{pmatrix} + \begin{pmatrix} \boldsymbol{\delta}_i \\ \mathbf{0} \end{pmatrix} + \mathbf{c}_i \pmod K. \quad (65)$$

Since that

$$(\boldsymbol{\delta}_i^\top, \mathbf{0}^\top) \mathbf{W}^{(0)} \tilde{\mathbf{p}}_i \lesssim \xi, \quad (66)$$

by that each entry of $\mathbf{W}^{(0)}$ follows $\mathcal{N}(0, \xi^2)$, and

$$(\boldsymbol{\delta}_i^\top, \mathbf{0}^\top) \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}} \tilde{\mathbf{p}}_{query} = 0, \quad (67)$$

we have that for $\tilde{\mathbf{p}}_i$ that shares the same TSR pattern as the query,

$$\begin{aligned} & \tilde{\mathbf{p}}_i^\top \mathbf{W}^{(T)} \tilde{\mathbf{p}}_{query} \\ &= \tilde{\mathbf{p}}_i^\top (\mathbf{W}^{(0)} + \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}}) \tilde{\mathbf{p}}_{query} \\ &= ((\boldsymbol{\mu}'_j{}^\top, \mathbf{y}_{i,1}^\top) + \mathbf{c}_i^\top \pmod K) (\mathbf{W}^{(0)} + \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}}) \tilde{\mathbf{p}}_{query}. \end{aligned} \quad (68)$$

Since that $\boldsymbol{\lambda}_j$ is orthogonal to all the $\boldsymbol{\mu}_i, i \in [M]$, we have similar conclusion for $\boldsymbol{\lambda}_j$ as $\boldsymbol{\delta}_i$, i.e.,

$$(\boldsymbol{\lambda}_j^\top, \mathbf{0}^\top) \mathbf{W}^{(0)} \tilde{\mathbf{p}}_i \lesssim \xi, \quad (69)$$

and

$$(\boldsymbol{\lambda}_j^\top, \mathbf{0}^\top) \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}} \tilde{\mathbf{p}}_{query} = 0. \quad (70)$$

Let $\boldsymbol{\mu}'_j = \boldsymbol{\lambda}_j + \tilde{\boldsymbol{\mu}}_j = \boldsymbol{\lambda}_j + \sum_{i=1}^{M'} k_{j,i} \boldsymbol{\mu}_i$. Then, we have

$$\begin{aligned} & \tilde{\mathbf{p}}_i^\top \mathbf{W}^{(T)} \tilde{\mathbf{p}}_{query} \\ &= ((\boldsymbol{\lambda}_j^\top + \sum_{i=1}^{M'} k_{j,i} \boldsymbol{\mu}_i^\top, \mathbf{y}_{i,1}^\top) + \mathbf{c}_i^\top \pmod K) (\mathbf{W}^{(0)} + \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}}) ((\boldsymbol{\lambda}_j^\top \\ & \quad + \sum_{i=1}^{M'} k_{j,i} \boldsymbol{\mu}_i^\top, \mathbf{0}^\top) + \mathbf{c}_1)^\top \\ &= \sum_{i=1}^{M'} k_{j,i}^2 ((\boldsymbol{\mu}_i^\top, \mathbf{y}_{i,1}^\top) + \mathbf{c}_i^\top \pmod K) (\mathbf{W}^{(0)} + \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}}) ((\boldsymbol{\mu}_i^\top, \mathbf{0}^\top) + \mathbf{c}_1)^\top \\ & \quad + \sum_{i \neq i'} k_{j,i} k_{j,i'} ((\boldsymbol{\mu}_i^\top, \mathbf{y}_{i,1}^\top) + \mathbf{c}_i^\top \pmod K) (\mathbf{W}^{(0)} + \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}}) ((\boldsymbol{\mu}_{i'}^\top, \mathbf{0}^\top) + \mathbf{c}_1)^\top \\ & \geq C \cdot \Theta(\log \frac{1}{\epsilon}) - \Theta(\xi) \\ & = \Theta(\log \frac{1}{\epsilon}), \end{aligned} \quad (71)$$

where the second to last step is by Theorem 1. The last step holds if $C \geq \Theta(\log^{-1}(1/\epsilon))$. Since the gradient updates for different TRR patterns are very close to each other, we have that $\sum_{i \neq i'} |k_{j,i} k_{j,i'}| \leq 1$ and

$$\begin{aligned} & \sum_{i \neq i'} k_{j,i} k_{j,i'} ((\boldsymbol{\mu}_i^\top, \mathbf{y}_{i,1}^\top) + \mathbf{c}_i^\top \pmod K) (\mathbf{W}^{(0)} + \frac{\eta}{B} \sum_{n \in \mathcal{B}_b} \sum_{b=0}^{T-1} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(b)}}) ((\boldsymbol{\mu}_{i'}^\top, \mathbf{0}^\top) + \mathbf{c}_1)^\top \\ & \lesssim \Theta(1) \cdot \frac{\tilde{\mathbf{p}}_s^\top \mathbf{W}^{(T)} \tilde{\mathbf{p}}_{query}}{\log \frac{1}{\epsilon}}, \end{aligned} \quad (72)$$

where $\tilde{\mathbf{p}}_s$ shares the same TSR pattern and the same step as $\tilde{\mathbf{p}}_{query}$. Hence, for $\tilde{\mathbf{p}}_i$ that shares a different TSR pattern with $\tilde{\mathbf{p}}_{query}$,

$$\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(T)} \tilde{\mathbf{p}}_{query} \lesssim \Theta(1). \quad (73)$$

Therefore, we can derive that

$$\sum_{i \in \mathcal{S}_1^*} \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(T)} \tilde{\mathbf{p}}_{query}) \geq 1 - \epsilon, \quad (74)$$

where \mathcal{S}_1^* is the set of the first step of examples that share the same TSR pattern as the query. Then, the first step leads to a correct prediction with zero testing error, since that $\max_{j \in [M']} A_{k(q_0, j)}$ is the largest to make the correct prediction for \mathbf{x}_{query} if $\mathbf{x}_{query} = \boldsymbol{\mu}'_{q_0}$, i.e.,

$$\mathbf{v}_1 = f_1(\boldsymbol{\mu}'_{q_0}). \quad (75)$$

Suppose that the k -th step generates a zero testing error. Then, for the $k + 1$ -th step, we know that there exists \mathbf{p}_j that shares the same TSR pattern as \mathbf{v}_k . Then, we can also derive that

$$\tilde{\mathbf{p}}_j^\top \mathbf{W}^{(T)} ((\mathbf{v}_k^\top, \mathbf{0}^\top)^\top + \mathbf{c}_k^\top)^\top = \Theta(\log \frac{1}{\epsilon}), \quad (76)$$

and

$$\sum_{j \in \mathcal{S}_k^*} \text{softmax}(\tilde{\mathbf{p}}_j^\top \mathbf{W}^{(T)} ((\mathbf{v}_{k-1}^\top, \mathbf{v}_k^\top)^\top + \mathbf{c}_k^\top)^\top) \geq 1 - \epsilon. \quad (77)$$

Hence, the $k + 1$ -th also makes the correct prediction, i.e.,

$$\mathbf{v}_{k+1} = f_{k+1} \circ \dots \circ f_1(\boldsymbol{\mu}'_{q_0}), \quad (78)$$

where $\boldsymbol{\mu}'_{q_{k+1}}$ is the TSR pattern of the $k + 1$ -th step input. Therefore, we show that CoT makes the correct prediction in each step as well as in the final prediction, such that

$$\bar{R}_{CoT, \mathbf{x} \in \mathcal{M}', f \in \mathcal{T}'}^f(\Psi) = 0. \quad (79)$$

□

E.3 PROOF OF THEOREM 3

Proof. We know that the positional encodings are the same for the ICL inference in all examples. Hence, similar to (74), we can derive that

$$\sum_{i \in \mathcal{S}_K^*} \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(T)} \tilde{\mathbf{p}}_{query}) \geq 1 - \epsilon, \quad (80)$$

where \mathcal{S}_K^* is the set of the last step output of examples that share the same TSR pattern as the last step output of the query. For $\mathbf{x}_{query} = \boldsymbol{\mu}'_q$, $q \in [K']$, we know that the distribution of the corresponding label \mathbf{y} of \mathbf{x} with $\text{TSR}(\mathbf{x}) = q$ follows the q -th row the K -steps transition matrix B^f . Let $F(\Psi; \mathbf{P}) = \sum_{i=1}^{M'} \lambda_i^{\mathbf{P}} \boldsymbol{\mu}'_i$. Hence, based on the output scheme of ICL as stated in Section 2.3, we have that

$$\mathbf{v} = \arg \min_{\mathbf{y} \in \mathcal{M}'} \frac{1}{2} \|F(\Psi; \mathbf{P}) - \mathbf{y}\|^2 = \boldsymbol{\mu}'_{\arg \max_{i \in [M']} \lambda_i^{\mathbf{P}}}. \quad (81)$$

Note that the probability of examples with the most probable final output with $\boldsymbol{\mu}'_q$ as the TSR pattern of the input is

$$B_{(q, \text{TSR}(f(\boldsymbol{\mu}'_q)))}. \quad (82)$$

To ensure that the number of examples with the same TSR pattern as the query that generates the most probable output is at least 1, we compute the following,

$$\begin{aligned} & \Pr \left(\frac{1}{l_{ts}} \sum_{i=1}^{l_{ts}} \mathbb{1}[m_i = \boldsymbol{\mu}'_q] \leq (1 - \rho_o^f/2) \alpha' B_{(q, \text{TSR}(f(\boldsymbol{\mu}'_q)))} \right) \\ & \leq e^{-l_{ts} \rho_o^{f^2} \alpha' B_{(q, \text{TSR}(f(\boldsymbol{\mu}'_q)))}} = M^{-C}, \end{aligned} \quad (83)$$

for some $c \in (0, 1)$ and $C > 0$ by the Chernoff bound of Bernoulli distribution in Lemma 1. Here, m_i is defined as the TSR pattern in the final output of the i -th context example by the K -steps transition matrix defined in 12. The TSR pattern of the most probable output of the testing query is μ'_{q_1} . Similarly, let the TSR pattern of the second most probable output of the testing query be μ'_{q_2} . We also have

$$\begin{aligned} & \Pr \left(\frac{1}{l_{ts}} \sum_{i=1}^{l_{ts}} \mathbb{1}[m_i = \mu'_{q_2}] \geq (1 - \rho_o^f/2) \alpha' B_{(q, q_1)}^f \right) \\ & \leq \Pr \left(\frac{1}{l_{ts}} \sum_{i=1}^{l_{ts}} \mathbb{1}[m_i = \mu'_{q_2}] \geq \alpha' B_{(q, q_2)} + c \cdot \rho_o^f \alpha' B_{(q, q_1)}^f \right) \\ & \leq e^{-l_{ts} \rho_o^{f^2} c^2 \alpha' B_{(q, q_1)}} = M^{-C}, \end{aligned} \quad (84)$$

by Lemma 2 and (12) for some constant $c > 0$. Therefore, to make the number of examples with the same TSR pattern in the output as the label of the query be at least 1 for any TSR pattern of the query and the output be the most probable one, we need

$$\begin{aligned} l_{ts}^f & \gtrsim \max \{ (\rho_o^{f^2} \alpha' \min_{i \in [M']} B_{(i, \text{TSR}(f(\mu'_i)))})^{-1} \log M, (\rho_o^f \alpha' \min_{i \in [M']} B_{(i, \text{TSR}(f(\mu'_i)))})^{-2} \log M \} \\ & = (\rho_o^f \alpha' \tau_o^f)^{-2} \log M. \end{aligned} \quad (85)$$

In addition, if Condition 1 holds such that the most probable output is the actual label, we can derive

$$\bar{R}_{ICL, \mathbf{x} \in \mathcal{M}', f \in \mathcal{T}'}^f(\Psi) = 0. \quad (86)$$

When (85) holds but Condition 1 does not, we know that ICL still always produces the most probable output by the K -steps transition matrix, but such an output is not the label since Condition 1 fails. Hence,

$$\bar{R}_{ICL, \mathbf{x} \in \mathcal{M}', f \in \mathcal{T}'}^f(\Psi) \geq \Omega(1). \quad (87)$$

When both Condition 1 and (85) do not hold, ICL can produce multiple possible outputs with a non-trivial probability, which is decided by the distribution of the prompt instead of the K -steps transition matrix. This can be seen from that (83) and (84) both do not hold since (85) fails. Then, ICL can produce both the most probable and the second most probable output with a constant probability. Let the TSR pattern of the r -th most probable output of the testing query be μ'_r . Recall that $F(\Psi; \mathbf{P}) = \sum_{i=1}^{M'} \lambda_i^{\mathbf{P}} \mu'_i$, we then have that for some small $\epsilon > 0$,

$$\lambda_r^{\mathbf{P}(q)} = \frac{|\{i \in [l_{ts}^f] : \mathbf{y}_i = \mu'_r \text{ in } \mathbf{P}\}|}{l_{ts}^f} \pm \epsilon. \quad (88)$$

Then, by (81), the output of the query is $\mu_{\arg \max_{r \in [M']} \lambda_r}$. Since that (85) does not hold, there exists at least a constant probability of the prompt \mathbf{P}' with the same query as \mathbf{P} such that

$$\lambda_r^{\mathbf{P}'} = \frac{|\{i \in [l_{ts}^f] : \mathbf{y}_i = \mu'_r \text{ in } \mathbf{P}'\}|}{l_{ts}^f} \pm \epsilon \neq \lambda_r^{\mathbf{P}}, \quad (89)$$

for some $r \in [M']$. Therefore, with a constant probability, the output for the same testing query and the same testing task f varies. This leads to

$$\bar{R}_{ICL, \mathbf{x} \in \mathcal{M}', f \in \mathcal{T}'}^f(\Psi) \geq \Omega(1). \quad (90)$$

□

E.4 PROOF OF PROPOSITION 1

Proof. This proposition is derived from the proof of Theorem 2. (22) comes from (77), while (23) comes from (78), both by induction. □

F PROOF OF LEMMAS

F.1 PROOF OF LEMMA 3

Proof.

$$\begin{aligned}
& \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \\
&= \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial F(\Psi; \mathbf{P})} \frac{\partial F(\Psi; \mathbf{P})}{\partial \mathbf{W}} \\
&= \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \\
&\quad \cdot (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_i) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top.
\end{aligned} \tag{91}$$

When $t = 0$, we know that each entry of $\mathbf{W}^{(0)}$ is generated from the Gaussian distribution $\mathcal{N}(0, \xi^2)$. Then,

$$|\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(0)} \tilde{\mathbf{p}}_{query}| = \left| \sum_{k,j} p_{i,k} p_{query,j} W_{k,j}^{(0)} \right| \lesssim \xi. \tag{92}$$

Hence,

$$\text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(0)} \tilde{\mathbf{p}}_{query}) \geq \frac{e^{-\Theta(\xi)}}{l \cdot e^{\Theta(\xi)}} = \frac{1}{l} \cdot e^{-\Theta(\xi)}, \tag{93}$$

$$\text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W}^{(0)} \tilde{\mathbf{p}}_{query}) \leq \frac{e^{-\Theta(\xi)}}{l \cdot e^{\Theta(\xi)}} = \frac{1}{l} \cdot e^{-\Theta(\xi)}. \tag{94}$$

We can obtain

$$F(\Psi; \mathbf{P}) = \sum_{i=1}^l \frac{e^{-\Theta(\xi)}}{l} \mathbf{W}_V \mathbf{p}_i. \tag{95}$$

Since that $\text{PE}(\cdot)$, and $\text{TRR}(\cdot)$ denote the positional encoding, and the TSR pattern of the input, respectively, we have that for \mathbf{p} ,

$$\tilde{\mathbf{p}}^\top \tilde{\mathbf{p}}_{query} = \mathbb{1}[\text{TRR}(\tilde{\mathbf{p}}) = \text{TRR}(\tilde{\mathbf{p}}_{query})] + \mathbb{1}[\text{PE}(\tilde{\mathbf{p}}) = \text{PE}(\tilde{\mathbf{p}}_i)]. \tag{96}$$

Given $\text{lab}(\cdot)$ is the label embedding of the context as the input, we have that for \mathbf{p} ,

$$\tilde{\mathbf{p}}^\top \tilde{\mathbf{p}}_i = \mathbb{1}[\text{TRR}(\tilde{\mathbf{p}}) = \text{TRR}(\tilde{\mathbf{p}}_i)] + \mathbb{1}[\text{lab}(\tilde{\mathbf{p}}) = \text{lab}(\tilde{\mathbf{p}}_i)] + \mathbb{1}[\text{PE}(\tilde{\mathbf{p}}) = \text{PE}(\tilde{\mathbf{p}}_i)], \tag{97}$$

$$(\mathbf{W}_V \tilde{\mathbf{p}})^\top \mathbf{W}_V \tilde{\mathbf{p}}_i = \mathbb{1}[\text{lab}(\tilde{\mathbf{p}}) = \text{lab}(\tilde{\mathbf{p}}_i)]. \tag{98}$$

When $t \geq 1$, we first consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}_{query}$ share the same TRR pattern, label pattern, and the positional encoding,

$$\begin{aligned}
\tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top &\geq 2 \cdot (3 - 3p_n(t) - (1 - p_n(t))) \\
&= 4(1 - p_n(t)),
\end{aligned} \tag{99}$$

and

$$\tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \leq 2 \cdot (3 - 3p_n(t)) = 6(1 - p_n(t)). \tag{100}$$

When $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}_{query}$ only share the same positional encoding or the same TRR pattern,

$$2 - 6p_n(t) \geq \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq -4p_n(t). \tag{101}$$

When $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}_{query}$ share both different positional encodings and TRR patterns,

$$-6p_n(t) \geq \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \mathbf{p}_{query}^\top \tilde{\mathbf{p}} \geq -2 - 4p_n(t). \quad (102)$$

Then, we consider the case where $\tilde{\mathbf{p}}$ only shares the same TRR pattern or the same positional encoding as $\tilde{\mathbf{p}}_i$. If $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}_{query}$ share the same TRR pattern, label pattern, and the positional encoding,

$$\begin{aligned} 3 - p_n(t) &\geq \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq 1 \cdot (3 - p_n(t) - (1 - p_n(t))) \\ &= 2. \end{aligned} \quad (103)$$

When $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}_{query}$ only share the same positional encoding or the same TRR pattern,

$$1 - p_n(t) \geq \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq 0. \quad (104)$$

When $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}_{query}$ only share both different positional encodings and TRR patterns,

$$-p_n(t) \geq \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq -1. \quad (105)$$

Note that $-(1 - p_n(t))p_n(t) + (1 - p_n(t))^2 \alpha^2 / K^2 < 0$ for $p_n(t) \in [\alpha/K, \alpha]$. Then, when $l \geq \Omega(\alpha^{-1})$ and $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_i$,

$$\begin{aligned} &(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ &\cdot \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ &\leq -4p_n(t)(1 - p_n(t))^2 - 4p_n(t)(1 - p_n(t))^2 \cdot \frac{\alpha^2}{K^2} \\ &\quad + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (-4p_n(t)) + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (1 - p_n(t)) (-2 - 4p_n(t)) (K - 1) \\ &= -4p_n(t)(1 - p_n(t))^2 \left(1 + \frac{\alpha^2}{K^2} \right) + \frac{2}{lK} (1 - \alpha) (- (K - 1) - (K + 1)p_n(t) + 2p_n(t)^2 (K - 1)). \end{aligned} \quad (106)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Note that

$$\frac{2}{Kl} \cdot (1 - \alpha) \cdot K(1 - p_n(t)) \lesssim |(-(1 - p_n(t))p_n(t) + (1 - p_n(t))^2 \frac{\alpha^2}{K^2})(1 - p_n(t))|, \quad (107)$$

if $l \geq \Omega(\alpha^{-1})$. Then,

$$\begin{aligned} &(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ &\cdot \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ &\leq -0 \cdot p_n(t)(1 - p_n(t)) + (1 - p_n(t))^2 \frac{\alpha^2}{K^2} \cdot (+2) + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (- (K - 1)) \\ &= 2(1 - p_n(t))^2 \frac{\alpha^2}{K^2} - \frac{K - 1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right). \end{aligned} \quad (108)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_i^\top \tilde{\mathbf{p}} \\ & \leq 0 - (1 - p_n(t))^2 \frac{\alpha^2}{K^2} + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (- (K - 1)) \\ & = - (1 - p_n(t))^2 \frac{\alpha^2}{K^2} - \frac{K - 1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right). \end{aligned} \quad (109)$$

Therefore, as long as

$$l \geq \Omega(\alpha^{-1}), \quad (110)$$

we have

$$\begin{aligned} & \tilde{\mathbf{p}}^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \\ & = \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \\ & \cdot \tilde{\mathbf{p}}^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_i^\top \tilde{\mathbf{p}} \\ & \leq \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (1 - p_n(t))^2 (-4p_n(t)(1 + \frac{\alpha^2}{K^2}) + \frac{2(K-1)\alpha^2}{K^2}) \right. \\ & \quad \left. + \left(\frac{1}{K} - \frac{1}{M} \right) (- (1 - p_n(t))^2 \frac{\alpha^2}{K^2}) \right) \\ & = \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (1 - p_n(t))^2 (-4p_n(t)(1 + \frac{\alpha^2}{K^2}) + \frac{\alpha^2}{K} (1 + \frac{2(K-1)}{K})) - \frac{\alpha^2}{K^3} (1 - p_n(t))^2 \right). \end{aligned} \quad (111)$$

We then consider the case where $\tilde{\mathbf{p}}'$ shares a different positional encoding and the same TRR pattern as $\tilde{\mathbf{p}}$. Let $\tilde{\mathbf{p}}$ share the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$\begin{aligned} 2(3 - p_n(t)) & \geq \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq 2 \cdot (3 - p_n(t) - (1 - p_n(t))) \\ & = 4. \end{aligned} \quad (112)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_{query}$ only share the same positional encoding or the same TRR pattern,

$$2(1 - p_n(t)) \geq \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq 0. \quad (113)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$-2p_n(t) \geq \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq -2. \quad (114)$$

Then, we consider the case where $\tilde{\mathbf{p}}$ only shares the same TRR pattern as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$\begin{aligned} 3 - p_n(t) & \geq \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \geq 1 \cdot (3 - 3p_n(t) - (1 - p_n(t))) = 2(1 - p_n(t)). \end{aligned} \quad (115)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$1 - p_n(t) \geq \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -2p_n(t). \quad (116)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$-p_n(t) \geq \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -1 - 2p_n(t). \quad (117)$$

Next, we consider the case where $\tilde{\mathbf{p}}$ only shares the same positional encoding as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$3 \geq \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 1 \cdot (3 - (1 - p_n(t))) = 2 + p_n(t). \quad (118)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$1 \geq \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq p_n(t). \quad (119)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$0 \geq \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -1 + p_n(t). \quad (120)$$

Then, when $l \geq \Omega(\alpha^{-1})$ and $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^{\top} \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq -4p_n(t)(1 - p_n(t)) + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (-2K). \end{aligned} \quad (121)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Then, by (107),

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^{\top} \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq -2p_n(t)(1 - p_n(t))^2 - 2p_n(t)(1 - p_n(t))^2 \cdot \frac{\alpha^2}{K^2} + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) ((-1 - 2p_n(t))K) \\ & = -2p_n(t)(1 - p_n(t))^2 \left(1 + \frac{\alpha^2}{K^2} \right) + \frac{1}{l} (1 - \alpha) (-1 - 2p_n(t)). \end{aligned} \quad (122)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^{\top} \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^{\top} \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq p_n(t)(1 - p_n(t))^2 + p_n(t)(1 - p_n(t))^2 \frac{\alpha^2}{K^2} + \frac{1}{l} (1 - \alpha) (-1 - 2p_n(t)) \\ & = p_n(t)(1 - p_n(t))^2 \left(1 + \frac{\alpha^2}{K^2} \right) - \frac{1}{l} (1 - \alpha) (1 + 2p_n(t)). \end{aligned} \quad (123)$$

Therefore, as long as

$$l \geq \Omega(\alpha^{-1}), \quad (124)$$

we have

$$\begin{aligned} & \tilde{\mathbf{p}}'^{\top} \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \tilde{\mathbf{p}} \\ &= \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^{\top} \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}^{\top} (\tilde{\mathbf{p}}_i \\ & \quad - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (-4 - 2(K-1)(1-p_n(t))(1 + \frac{\alpha^2}{K^2})) p_n(t)(1-p_n(t)) \right. \\ & \quad \left. + (\frac{1}{K} - \frac{1}{M}) p_n(t)(1-p_n(t))^2 (1 + \frac{\alpha^2}{K^2}) \right) \\ & = \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (-4 - (3K-2)(1-p_n(t))(1 + \frac{\alpha^2}{K^2})) p_n(t)(1-p_n(t)) \right. \\ & \quad \left. + \frac{1}{K} p_n(t)(1-p_n(t))^2 (1 + \frac{\alpha^2}{K^2}) \right), \end{aligned} \quad (125)$$

and

$$\begin{aligned} & \tilde{\mathbf{p}}'^{\top} \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \tilde{\mathbf{p}} \\ & \geq \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (-4 - (3K-2)(1-p_n(t))(1 + \frac{\alpha^2}{K^2})) p_n(t)(1-p_n(t)) \right. \\ & \quad \left. + \frac{1}{K} p_n(t)(1-p_n(t))^2 (1 + \frac{\alpha^2}{K^2}) + \frac{1}{K} \cdot (1-p_n(t))^2 (-p_n(t) + (1-p_n(t)) \frac{\alpha^2}{K^2}) \right) \\ & = \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (-4 - (3K-2)(1-p_n(t))(1 + \frac{\alpha^2}{K^2})) p_n(t)(1-p_n(t)) + \frac{\alpha^2}{K^3} (1-p_n(t))^2 \right). \end{aligned} \quad (126)$$

We next consider the case where $\tilde{\mathbf{p}}'$ shares a different TRR pattern and the same positional encoding as $\tilde{\mathbf{p}}$. Let $\tilde{\mathbf{p}}$ share the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and positional encoding,

$$\begin{aligned} 2(3-p_n(t)) & \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 2 \cdot (3-p_n(t) - (1-p_n(t))) \\ & = 4. \end{aligned} \quad (127)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$2(1-p_n(t)) \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 0. \quad (128)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$-2p_n(t) \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -2. \quad (129)$$

Then, we consider the case where $\tilde{\mathbf{p}}$ only shares the same TRR pattern as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$3 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 1 \cdot (3 - (1-p_n(t))) = 2 + p_n(t). \quad (130)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$1 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq p_n(t). \quad (131)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$0 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -1 + p_n(t). \quad (132)$$

Next, we consider the case where $\tilde{\mathbf{p}}$ only shares the same positional encoding as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$3 - p_n(t) \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 1 \cdot (3 - p_n(t) - (1 - p_n(t))) = 2. \quad (133)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$1 - p_n(t) \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 0. \quad (134)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$-p_n(t) \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -1. \quad (135)$$

Then, when $l \geq \Omega(\alpha^{-1})$, and when $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$, by (107),

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^{\top} \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq 0 - 2(1 - p_n(t))^2 \frac{\alpha^2}{K^2} + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (-2(K - 1)). \end{aligned} \quad (136)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^{\top} \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq -p_n(t)(1 - p_n(t))(-1 + p_n(t)) + p_n(t)(1 - p_n(t))^2 \cdot \frac{\alpha^2}{K^2} + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) K(-1 + p_n(t)) \\ & = p_n(t)(1 - p_n(t))^2 \left(\frac{\alpha^2}{K^2} + 1 \right) + \frac{1}{l} (1 - \alpha)(-1 + p_n(t)). \end{aligned} \quad (137)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned}
& \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\
& \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\
& \leq - (1 - p_n(t))^2 \frac{\alpha^2}{K^2} - 0 + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (-K + 1) \\
& = - (1 - p_n(t))^2 \frac{\alpha^2}{K^2} - \frac{K-1}{Kl} (1 - \alpha).
\end{aligned} \tag{138}$$

Therefore, as long as

$$l \geq \Omega(\alpha^{-1}), \tag{139}$$

we have

$$\begin{aligned}
& \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \\
& = \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \\
& \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\
& \leq \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} \left(-\frac{\alpha^2}{K^2} + (K-1) \left(1 + \frac{\alpha^2}{K^2} \right) p_n(t) \right) (1 - p_n(t))^2 - \left(\frac{1}{K} \right. \right. \\
& \quad \left. \left. - \frac{1}{M} \right) (1 - p_n(t))^2 \frac{\alpha^2}{K^2} \right) \\
& = \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} \left(-\frac{\alpha^2}{K^2} + \left(K-1 + \frac{(2K-1)\alpha^2}{K^2} \right) p_n(t) \right) (1 - p_n(t))^2 - (1 - p_n(t))^2 \frac{\alpha^2}{K^3} \right).
\end{aligned} \tag{140}$$

and

$$\begin{aligned}
& \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \\
& \geq \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} \left(-\frac{\alpha^2}{K^2} + \left(K-1 + \frac{(2K-1)\alpha^2}{K^2} \right) p_n(t) \right) (1 - p_n(t))^2 - (1 - p_n(t))^2 \frac{\alpha^2}{K^3} \right. \\
& \quad \left. + \frac{1}{K} \cdot (1 - p_n(t))^2 \left(-p_n(t) + (1 - p_n(t)) \frac{\alpha^2}{K^2} \right) \right).
\end{aligned} \tag{141}$$

We next consider the case where $\tilde{\mathbf{p}}'$ shares a different TRR pattern and a different positional encoding as $\tilde{\mathbf{p}}$. Let $\tilde{\mathbf{p}}$ share the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$6 \geq \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq 2 \cdot (3 - (1 - p_n(t))) = 4 + 2p_n(t). \tag{142}$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$2 \geq \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \geq 2p_n(t). \tag{143}$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$0 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -2 + 2p_n(t). \quad (144)$$

Then, we consider the case where $\tilde{\mathbf{p}}$ only shares the same TRR pattern as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$3 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 1 \cdot (3 - p_n(t) - (1 - p_n(t))) = 2. \quad (145)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$1 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 0. \quad (146)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$0 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -1. \quad (147)$$

Next, we consider the case where $\tilde{\mathbf{p}}$ only shares the same positional encoding as $\tilde{\mathbf{p}}_{query}$. If $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ share the same TRR pattern, label pattern, and the positional encoding,

$$3 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq 1 \cdot (3 - (1 - p_n(t))) = 2 + p_n(t). \quad (148)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share the same positional encoding or the same TRR pattern,

$$1 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq p_n(t). \quad (149)$$

When $\tilde{\mathbf{p}}'$ and $\tilde{\mathbf{p}}_i$ only share both different positional encodings and TRR patterns,

$$0 \geq \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \geq -1 + p_n(t). \quad (150)$$

Then, when $l \geq \Omega(\alpha^{-1})$, and when $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^{\top} \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq -p_n(t)(1 - p_n(t))(-2 + 2p_n(t)) + (1 - p_n(t))^2 \frac{\alpha^2}{K^2} \cdot 2p_n(t) + \frac{1}{l}(1 - \alpha)(-2 + 2p_n(t)). \end{aligned} \quad (151)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^{\top} \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^{\top} (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^{\top} \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^{\top} \tilde{\mathbf{p}} \\ & \leq 0 + p_n(t)(1 - p_n(t))^2 \cdot \frac{\alpha^2}{K^2} \cdot (-1) + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (-K) \\ & = -p_n(t)(1 - p_n(t))^2 \frac{\alpha^2}{K^2} + \frac{1}{l}(1 - \alpha)(-1). \end{aligned} \quad (152)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned}
& \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\
& \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}} - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\
& \leq - (1 - p_n(t)) p_n(t) (-1 + p_n(t)) + p_n(t) (1 - p_n(t))^2 \frac{\alpha^2}{K^2} + \frac{1}{l} \left(\frac{1}{K} - \frac{\alpha}{K} \right) (-1 + p_n(t)) K \\
& = (1 - p_n(t))^2 p_n(t) \left(1 + \frac{\alpha^2}{K^2} \right) + \frac{1}{l} (1 - \alpha) (-1 + p_n(t)).
\end{aligned} \tag{153}$$

Therefore, as long as

$$l \geq \Omega(\alpha^{-1}), \tag{154}$$

we have

$$\begin{aligned}
& \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \\
& = \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \\
& \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}} - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\
& \leq \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} (-p_n(t)(1 - p_n(t))(-2 + 2p_n(t)) + (3 - K)(1 - p_n(t))^2 \frac{\alpha^2}{K^2} \cdot p_n(t)) \right. \\
& \quad \left. + \left(\frac{1}{K} - \frac{1}{M} \right) (1 - p_n(t))^2 p_n(t) \left(1 + \frac{\alpha^2}{K^2} \right) \right) \\
& = \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} p_n(t)(1 - p_n(t))^2 (2 - K + \frac{(2 - K)\alpha^2}{K^2}) \right. \\
& \quad \left. + (1 - p_n(t))^2 p_n(t) \left(1 + \frac{\alpha^2}{K^2} \right) \cdot \frac{1}{K} \right),
\end{aligned} \tag{155}$$

and

$$\begin{aligned}
& \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \\
& \geq \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} p_n(t)(1 - p_n(t))^2 \left(1 + \frac{(2 - K)\alpha^2}{K^2} \right) + (1 - p_n(t))^2 p_n(t) \left(1 + \frac{\alpha^2}{K^2} \right) \cdot \frac{1}{K} \right. \\
& \quad \left. + \frac{1}{K} \cdot (1 - p_n(t))^2 (-p_n(t) + (1 - p_n(t)) \frac{\alpha^2}{K^2}) \right) \\
& = \eta \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{KM} p_n(t)(1 - p_n(t))^2 \left(1 + \frac{(2 - K)\alpha^2}{K^2} \right) + (1 - p_n(t))^2 \cdot \frac{\alpha^2}{K^3} \right).
\end{aligned} \tag{156}$$

□

F.2 PROOF OF LEMMA 4

Proof. We can derive that when $1 - p_n(t) \geq \Omega(1)$, $\tilde{\mathbf{p}}'^\top \mathbf{W}^{(t)} \tilde{\mathbf{p}}$ increases if $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}'$ share the same positional encoding. Otherwise, $\tilde{\mathbf{p}}'^\top \mathbf{W}^{(t)} \tilde{\mathbf{p}}$ decreases. We know that $p_n(t) \geq \frac{\alpha}{2}$. Combining the

results in Lemma 3, we can derive that when $t \geq 1$,

$$\mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}^{(t)}}. \quad (157)$$

Then, for $\tilde{\mathbf{p}}_i^n$ that share the same TRR pattern and the same positional encoding of $\tilde{\mathbf{p}}_{query}^n$,

$$\begin{aligned} \frac{p_n(t+1)}{|\mathcal{S}_1^n|} &= \text{softmax}(\mathbf{p}_i^{n \top} \mathbf{W}^{(t+1)} \tilde{\mathbf{p}}_{query}^n) \\ &\geq \frac{1}{l} \cdot \frac{1}{\frac{\alpha}{K} + \frac{(K-1)\alpha}{K} \cdot e^{-s_1} + (\frac{1}{K} - \frac{\alpha}{K})((K-1)e^{-s_2} + e^{-s_3})}, \end{aligned} \quad (158)$$

where

$$s_1 \geq \eta \sum_{b=0}^t ((1-p_n(b))^2 \frac{\alpha^2}{K^3} + \frac{\alpha^2}{K^3} (1-p_n(b))^2) = \eta \sum_{b=0}^t (1-p_n(b))^2 \frac{2\alpha^2}{K^3}, \quad (159)$$

$$s_2 \geq \sum_{b=0}^t (1-p_n(b))^2 \cdot \frac{2\eta\alpha^2}{K^3}, \quad (160)$$

$$\begin{aligned} s_3 &\geq -\frac{\eta}{KM} \sum_{b=0}^t (1-p_n(b))^2 (-4p_n(b)(1 + \frac{\alpha^2}{K^2}) + \frac{\alpha^2}{K} (1 + \frac{2(K-1)}{K})) + \frac{\alpha^2}{K^2} \\ &\quad - (K-1 + \frac{2K-1}{K^2} \alpha^2) p_n(b)) \\ &\geq \frac{\eta}{KM} \sum_{b=0}^t (1-p_n(b))^2 (p_n(b)(3 + \frac{\alpha^2}{K^2})(4 + \frac{2K-1}{K^2})), \end{aligned} \quad (161)$$

where the last step is by $Kp_n(b) \geq 4\alpha^2/K^2$ when $p_n(b) \geq \alpha/K$. For $\tilde{\mathbf{p}}_i^n$ that share the same TRR pattern and a different positional encoding of $\tilde{\mathbf{p}}_{query}^n$,

$$\text{softmax}(\tilde{\mathbf{p}}_i^{n \top} \mathbf{W}^{(t+1)} \tilde{\mathbf{p}}_{query}^n) = \frac{1}{l} \cdot \frac{1}{\frac{\alpha}{K} e^{s_1} + \frac{(K-1)\alpha}{K} + (\frac{1}{K} - \frac{\alpha}{K})((K-1)e^{-s_4} + e^{s_5})}, \quad (162)$$

where

$$\begin{aligned} s_4 &\geq -\sum_{b=0}^t \frac{\eta}{M} ((-4 - (3K-2)(1-p_n(b))(1 + \frac{\alpha^2}{K^2})) p_n(b)(1-p_n(b)) \\ &\quad - (2-K)(1 + \frac{\alpha^2}{K^2}) p_n(b)(1-p_n(b))^2) \end{aligned} \quad (163)$$

$$= \sum_{b=0}^t \frac{\eta}{M} (4 + 2K(1-p_n(b))(1 + \frac{\alpha^2}{K^2})) p_n(b)(1-p_n(b)),$$

$$s_5 \geq \sum_{b=0}^t (1-p_n(b))^2 \cdot \frac{2\eta\alpha^2}{K^3}. \quad (164)$$

When $M \geq \Omega(K^4 \alpha^{-1})$ and $t \geq \Omega(\eta^{-1} K^3 \log K \alpha^{-2})$,

$$(K-1)e^{-s_4} + e^{s_5} > K. \quad (165)$$

If $M \geq \Omega(K^4 \alpha^{-1})$ and $t \leq O(\eta^{-1} K^3 \log K \alpha^{-2})$, we cannot ensure

$$(K-1)e^{-s_4} + e^{s_5} > K. \quad (166)$$

For $\tilde{\mathbf{p}}_i^n$ that share a different TRR pattern and the same positional encoding of $\tilde{\mathbf{p}}_{query}^n$,

$$\text{softmax}(\tilde{\mathbf{p}}_i^{n \top} \mathbf{W}^{(t+1)} \tilde{\mathbf{p}}_{query}^n) = \frac{1}{l} \cdot \frac{1}{\frac{\alpha}{K} e^{s_3} + \frac{\alpha}{K} \cdot e^{-s_4} + (\frac{1}{K} - \frac{\alpha}{K})(1 + (K-1)e^{-s_6})}, \quad (167)$$

where

$$s_6 \geq \eta \sum_{b=0}^t \frac{2\alpha^2}{K^3} (1 - p_n(b))^2. \quad (168)$$

For $\tilde{\mathbf{p}}_i^n$ that share a different TRR pattern and a different positional encoding of $\tilde{\mathbf{p}}_{query}^n$,

$$\text{softmax}(\tilde{\mathbf{p}}_i^n \top \mathbf{W}^{(t+1)} \tilde{\mathbf{p}}_{query}^n) = \frac{1}{l} \cdot \frac{1}{\frac{\alpha}{K} e^{s_2} + (\frac{1}{K} - \frac{\alpha}{K})(K - 1 + e^{s_6}) + \frac{\alpha}{K} e^{s_4}}. \quad (169)$$

Note that when $t \lesssim \eta^{-1} \alpha^{-2} K^3$, for \mathbf{p}_{query}^n in the k -th step, we have

$$\sum_{i \in \mathcal{S}_{[K] \setminus \{k\}}} \text{softmax}(\tilde{\mathbf{p}}_i^n \top \mathbf{W}^{(t+1)} \tilde{\mathbf{p}}_{query}^n) \geq \Omega(1), \quad (170)$$

for $\tilde{\mathbf{p}}_i^n$ that share a different positional encoding from $\tilde{\mathbf{p}}_{query}^n$. To make the total softmax values on contexts that share a different positional encoding and a different TRR pattern from the query smaller than ϵ , we need

$$s_1, s_2, s_6 \gtrsim \log \frac{K}{\epsilon}. \quad (171)$$

When t further increases to be larger than $\Omega(\eta^{-1} \alpha^{-2} K^3 \log \frac{K}{\epsilon})$, we also have that the total softmax values on contexts that share a different positional encoding and the same TRR pattern from the query smaller than ϵ . Therefore,

$$t \gtrsim T_1 := \eta^{-1} \alpha^{-2} K^3 \log \frac{K}{\epsilon}. \quad (172)$$

□

F.3 PROOF OF LEMMA 5

Proof. We consider the case when $t \geq T_1$ given Lemma 4. When $l \geq \Omega(\alpha^{-1})$, and when $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \leq -4p_n(t)(1 - p_n(t))^2 + \epsilon \\ & \lesssim -4p_n(t)(1 - p_n(t))^2. \end{aligned} \quad (173)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \lesssim -0 \cdot p_n(t)(1 - p_n(t)) + \epsilon \\ & \lesssim \epsilon. \end{aligned} \quad (174)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r \top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \lesssim \epsilon. \end{aligned} \quad (175)$$

Therefore,

$$\begin{aligned}
& \tilde{\mathbf{p}}^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \\
&= \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \\
&\quad \cdot \tilde{\mathbf{p}}^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\
&\lesssim \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \left(\frac{1}{2M} (-4p_n(t)(1-p_n(t))^2) + \left(\frac{1}{2} - \frac{1}{M} \right) \cdot \epsilon \right) \\
&= -\eta \cdot \frac{1}{2M} \cdot \frac{1}{B} \sum_{n \in \mathcal{B}_b} 4p_n(t)(1-p_n(t))^2.
\end{aligned} \tag{176}$$

We then discuss if $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{p}}'$ only share the same TRR pattern. When $l \geq \Omega(\alpha^{-1})$, and when $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$, we can obtain

$$\begin{aligned}
& \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\
&\quad \cdot \tilde{\mathbf{p}}'^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\
&\gtrsim -2(1-p_n(t))^2 p_n(t).
\end{aligned} \tag{177}$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned}
& \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\
&\quad \cdot \tilde{\mathbf{p}}'^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\
&\gtrsim -(1-p_n(t))(1-p_n(t))p_n(t).
\end{aligned} \tag{178}$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned}
& \left| \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \right. \\
&\quad \left. \cdot \tilde{\mathbf{p}}'^\top (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \right| \\
&\lesssim \epsilon.
\end{aligned} \tag{179}$$

Therefore,

$$\begin{aligned}
& \left| \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \right| \\
&= \left| \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}^\top \right. \\
&\quad \left. \cdot (\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \right| \\
&\leq \eta \epsilon.
\end{aligned} \tag{180}$$

We next discuss when $\tilde{\mathbf{p}}$ only shares the same positional encoding as $\tilde{\mathbf{p}}'$. When $l \geq \Omega(\alpha^{-1})$, and when $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \lesssim \epsilon. \end{aligned} \quad (181)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \lesssim -p_n(t)(1-p_n(t))(-1+p_n(t)) + \frac{1}{M} \\ & \lesssim p_n(t)(1-p_n(t))^2. \end{aligned} \quad (182)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \lesssim \epsilon. \end{aligned} \quad (183)$$

Therefore,

$$\begin{aligned} & \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \\ & = \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}^\top \\ & \cdot \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \lesssim \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{1}{2M} \cdot p_n(b)(1-p_n(b))^2. \end{aligned} \quad (184)$$

We then consider if $\tilde{\mathbf{p}}$ shares a different TRR pattern and a different positional encoding as $\tilde{\mathbf{p}}'$. When $l \geq \Omega(\alpha^{-1})$, and when $\tilde{\mathbf{p}}$ shares the same TRR pattern and the positional encoding as $\tilde{\mathbf{p}}_{query}$,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \gtrsim \epsilon. \end{aligned} \quad (185)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same TRR pattern and the different positional encoding as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \\ & \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \\ & \gtrsim - (1 - p_n(t)) p_n(t). \end{aligned} \quad (186)$$

We next consider the case where $\tilde{\mathbf{p}}$ shares the same positional encoding and the different TRR pattern as $\tilde{\mathbf{p}}_{query}$. Then,

$$\begin{aligned} & \left| \left(\sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i - \mathbf{z}^n \right)^\top \sum_{i=1}^l \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \mathbf{W}_V \tilde{\mathbf{p}}_i \right. \\ & \left. \cdot \tilde{\mathbf{p}}'^\top \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \right| \\ & \lesssim \epsilon. \end{aligned} \quad (187)$$

Therefore,

$$\begin{aligned} & \left| \tilde{\mathbf{p}}'^\top \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} \frac{\partial \ell(\Psi; \mathbf{P}^n, \mathbf{z}^n)}{\partial \mathbf{W}} \mathbf{p} \right| \\ & = \left| \eta \frac{1}{B} \sum_{n \in \mathcal{B}_b} (F(\Psi; \mathbf{P}) - \mathbf{z}^n)^\top \sum_{i=1}^l \mathbf{W}_V \tilde{\mathbf{p}}_i \text{softmax}(\tilde{\mathbf{p}}_i^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}^\top \right. \\ & \quad \left. \left(\tilde{\mathbf{p}}_i - \sum_{r=1}^l \text{softmax}(\tilde{\mathbf{p}}_r^\top \mathbf{W} \tilde{\mathbf{p}}_{query}) \tilde{\mathbf{p}}_r \right) \tilde{\mathbf{p}}_{query}^\top \tilde{\mathbf{p}} \right| \\ & \lesssim \eta \epsilon. \end{aligned} \quad (188)$$

□