

HOW TRANSFORMERS IMPLEMENT INDUCTION HEADS: APPROXIMATION AND OPTIMIZATION ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformers have demonstrated exceptional in-context learning capabilities, yet the theoretical understanding of the underlying mechanisms remains limited. A recent work (Elhage et al., 2021) identified a “rich” in-context mechanism known as induction head, contrasting with “lazy” n -gram models that overlook long-range dependencies. In this work, we provide both approximation and dynamics analyses of how transformers implement induction heads. In the approximation analysis, we formalize both standard and generalized induction head mechanisms, and examine how transformers can efficiently implement them, with an emphasis on the distinct role of each transformer submodule. For the dynamics analysis, we study the training dynamics on a synthetic mixed target, composed of a 4-gram and an in-context 2-gram component. This controlled setting allows us to precisely characterize the entire training process and uncover an *abrupt transition* from lazy (4-gram) to rich (induction head) mechanisms as training progresses.

1 INTRODUCTION

Transformer, introduced by Vaswani et al. (2017), have achieved remarkable success across various domains, including natural language processing, computer vision, and scientific computing. An emergent observation is that transformers, trained on trillions of tokens, can perform (few-shot) in-context learning (ICL), which makes prediction based on the contextual information without needing model retraining (Brown et al., 2020). This ICL ability is widely regarded as crucial for enabling large language models (LLMs) to solve reasoning tasks, representing a key step toward more advanced artificial intelligence.

To understand how transformers implement ICL, Elhage et al. (2021) and Olsson et al. (2022) identified a simple yet powerful mechanism known as **induction head**. Specifically, given an input sequence $[\dots ab\dots a]$, an induction head predicts b as the next token by leveraging the prior occurrence of the pattern ab in the context, effectively modeling an in-context bi-gram. In contrast, traditional n -gram model (Shannon, 1948) (with a small n) utilizes only a limited number of recent tokens to predict the next token, which is context-independent and inevitably overlooks long-range dependence. Based on the extent of context utilization, we categorize n -gram model as a “*lazy*” *mechanism*, whereas the induction head represents a more “*rich*” *mechanism*.

Practically, induction heads have been demonstrated to play a critical role in enabling LLMs’ ICL capabilities (Song et al., 2024; Crosbie and Shutova, 2024), and even used to test new LLM architectures (Gu and Dao, 2023). Theoretically, induction heads also serve as a controllable tool for understanding various aspects of LLMs, such as multi-step reasoning (Sanford et al., 2024b) and inductive biases of different architectures (Jelassi et al., 2024).

In this paper, we aim to provide a theoretical analysis of how transformers can efficiently implement induction heads. The first key problem is to rigorously formalize induction heads and evaluate the efficiency of transformers in representing them. According to Elhage et al. (2021), the original induction head can be implemented using a two-layer, twelve-head transformer without feed-forward networks (FFNs). However, practical scenarios demand more powerful induction heads. Thus, it is crucial to generalize the mechanism behind and explore how different transformer submodules, such as varying the number of attention heads or incorporating FFNs, impact the transformer’s ability to implement them. This forms our first research objective:

(Approximation). Investigate how two-layer transformers express the induction head mechanism and its potential variants.

The next problem is to investigate the dynamics of transformers in learning induction heads. The pioneering works by Elhage et al. (2021) and Olsson et al. (2022) demonstrated that transformers undergo an abrupt phase transition to learning induction heads. A recent empirical study on synthetic datasets replicate this behavior, further showing that 2-gram is always learned prior to induction heads (Bietti et al., 2024). However, a rigorous theoretical analysis of this learning progression is still lacking. Closing this gap forms our second research objective:

(Optimization). Understand how transformers transition from relying on n -gram patterns to employing the induction head mechanism as training progresses.

Focusing on these two key problems, in this paper, we make the following contributions:

- **Approximation analysis: how transformers express induction heads.** We consider three types of induction heads with varying complexities. First, we show that two-layer, single-head transformers without FFNs can efficiently approximate the vanilla induction head (Elhage et al., 2021). We then introduce two generalized induction heads, which leverage richer in-context n -gram information and incorporate a general similarity function. Our analysis clarifies the distinct roles of multihead attention, positional encoding, dot-product structure, and FFNs in implementing these generalized induction heads.
- **Optimization analysis: how learning undergoes a sharp transition from n -gram to induction head.** We study the learning dynamics of a two-layer transformer without FFNs for a mixed target, composed of a 4-gram and an in-context 2-gram component. This toy setting allows us to capture the entire training process precisely. Specifically, we show that learning progresses through four phases: partial learning of the 4-gram, plateau of induction head learning, emergence of the induction head, and final convergence, showcasing a sharp transition from 4-gram to induction head. Our analysis identifies two key drivers of the transition: 1) time-scale separation due to low- and high-order parameter dependencies in self-attention, and 2) speed differences caused by the relative proportions of the two components in the mixed target.

2 RELATED WORKS

Empirical observations of induction head. The induction head mechanism was first identified by Elhage et al. (2021) in studying how two-layer transformers perform language modeling. Subsequently, Olsson et al. (2022) conducted a more systematic investigation, revealing two key findings: 1) induction head emerges abruptly during training, and 2) induction head plays a critical role in the development of in-context learning capabilities. To obtain a fine-grained understanding of how induction head emerges during training, recent studies have developed several synthetic settings (Reddy, 2024; Edelman et al., 2024; Bietti et al., 2024). Particularly, Bietti et al. (2024) successfully reproduced the fast learning of (global) bigrams and the slower development of induction head. Despite these efforts, a comprehensive theoretical understanding of how the induction head operates in two-layer transformers and how it is learned during training remains elusive.

Expressiveness of transformers. Theoretically, Dehghani et al. (2019); Pérez et al. (2021); Wei et al. (2022) explored the Turing-completeness of transformers; Yun et al. (2019) established the universal approximation property of transformers. Subsequent studies examined the efficiency of transformers in representing specific functions or tasks, such as sparse functions (Edelman et al., 2022), targets with nonlinear temporal kernels (Jiang and Li, 2023), practical computer programs (Giannou et al., 2023), long but sparse memories (Wang et al., 2024), induction head (Sanford et al., 2024a;b; Rajaraman et al., 2024), and memorization and reasoning (Chen and Zou, 2024). Besides, many studies suggest that transformers achieve in-context learning by approximating gradient-based iterations across various layers (Garg et al., 2022; Akyürek et al., 2022; Von Oswald et al., 2023; Mahankali et al., 2023; Bai et al., 2023; Shen et al., 2023). Besides, several studies explored the limitation of transformer’s expressivity, particularly in modeling formal languages or simulating circuits (Hahn, 2020; Weiss et al., 2021; Bhattamishra et al., 2020; Merrill et al., 2022; Merrill and Sabharwal, 2023). Among all these works, the most closely related to ours are Rajaraman et al. (2024), which examined a generalized induction head similar to our Eq. (6). Specifically, they showed that multi-

layer transformers with single-head attention can implement this mechanism. In contrast, we prove that two-layer transformers are sufficient if multihead attention is used.

Training dynamics of transformers. To gain insights into the dynamics of training transformers, several studies have analyzed simplified transformers on toy tasks. These tasks include learning distinct/common tokens (Tian et al., 2023), leaning balance/inbalanced features (Huang et al., 2023), linear regression task (Zhang et al., 2023; Ahn et al., 2024), multi-task linear regression (Chen et al., 2024a), binary classification (Li et al., 2024), transformer with diagonal weights (Abbe et al., 2024), learning causal structure (Nichani et al., 2024), sparse token selection task (Wang et al., 2024), and learning n -gram Markov chain (Chen et al., 2024b). Additionally, studies such as those by Ataee Tarzanagh et al. (2023), Tarzanagh et al. (2023) and Vasudeva et al. (2024) have analyzed scenarios where transformers converge to max-margin solutions. Furthermore, Thrampoulidis (2024) has examined the implicit bias of next-token prediction. Among these works, the most closely related to ours are Nichani et al. (2024) and Chen et al. (2024b), which proved that two-layer transformers can converge to induction head solutions. In this work, we explore a setting where the target is a mixture of 4-gram and induction head. We show that two-layer transformers can effectively converge to this mixed target and provide a precise description of the learning process associated with each component. Importantly, we are able to capture the *abrupt transition* from learning 4-gram patterns to mastering the induction head mechanism—a critical phase in the learning of induction heads, as highlighted in the seminal works (Elhage et al., 2021; Olsson et al., 2022).

3 PRELIMINARIES

Notations. For $k \in \mathbb{N}^+$, let $[k] = \{1, 2, \dots, k\}$. For a vector v and $1 \leq p \leq \infty$, we denote by $\|v\|_p$ the ℓ_p norm of v . For a matrix $A = (a_{i,j})$, we denote by $\|A\|$, $\|A\|_F$ the spectral and Frobenius norms, respectively; let $\|A\|_{1,1} = \sum_{i,j} |a_{i,j}|$. For an event S , we define $\mathbb{I}\{S\} = 1$ if S is true, and 0 otherwise. We use standard big-O notations \mathcal{O} , Ω , Θ to hide absolute positive constants, and use $\tilde{\mathcal{O}}$, $\tilde{\Omega}$, $\tilde{\Theta}$ to further hide logarithmic constants.

Sequence modeling. Given a sequence of tokens (x_1, x_2, x_3, \dots) with each token lying in \mathbb{R}^d , let $X_L = (x_1, x_2, \dots, x_L) \in \mathbb{R}^{d \times L}$ and $X_{m:n} = (x_m^\top, x_{m+1}^\top, \dots, x_n^\top)^\top \in \mathbb{R}^{(n-m+1)d}$. Given $A = (a_1, \dots, a_n) \in \mathbb{R}^{m \times n}$, we denote $(a_s)_{s=i}^j = (a_i, \dots, a_j) \in \mathbb{R}^{m \times (j-i+1)}$. Then, we consider the next-token prediction task: predict x_{L+1} using $X_L = (x_1, x_2, \dots, x_L)$.

In a n -gram model (Shannon, 1948), the conditional probability of predicting the next token is given by $p(x_{L+1}|X_L) = p(x_{L+1}|X_{L-n+2:L})$, meaning that the prediction depends only on the most recent $n-1$ tokens. In practice, the value of n is typically small (e.g., 2, 3, or 4), as the computational cost of n -gram models grows exponentially with n . However, n -gram models with small n cannot capture long-range interactions, leading to inferior performance in sequence modeling.

Transformer is designed to more efficiently capture long-range dependencies in sequence modeling (Vaswani et al., 2017). Specifically, given an L -token input sequence $X = (x_1, \dots, x_L) \in \mathbb{R}^{d \times L}$, an U -layer transformer TF processes it as follows. First, each input token is embedded into a higher-dimensional space through an *embedding layer*:

$$x_s^{(0)} = W_E x_s + b_E, \quad s \in [L], \quad \text{with } W_E \in \mathbb{R}^{D \times d}, b_E \in \mathbb{R}^D.$$

Next, the U -layer attention blocks process the embedded sequence $X^{(0)} = (x_1^{(0)}, \dots, x_L^{(0)})$ as follows, and the output of the final layer is taken as the output sequence $\text{TF}(X) = X^{(L)} \in \mathbb{R}^{D \times L}$:

$$\begin{aligned} X^{(u-\frac{1}{2})} &= X^{(u-1)} + \text{SA}^{(u)}(X^{(u-1)}), \quad u \in [U]; \\ X^{(u)} &= X^{(u-\frac{1}{2})} + \text{FFN}^{(u)}(X^{(u-\frac{1}{2})}), \quad u \in [U]. \end{aligned} \tag{1}$$

Here, $\text{FFN}^{(u)}$ denotes a (token-wise) two-layer FFN of width M , and $\text{SA}^{(u)}$ represents the multi-head self-attention operation. Specifically, when applied to a sequence $Z = (z_1, \dots, z_L) \in \mathbb{R}^{D \times L}$, $\text{SA}^{(l)}$ operates it as follows:

$$\begin{aligned} \text{SA}^{(u)}(Z) &= W_O^{(u)} \sum_{h=1}^{H_u} \text{SA}^{(u,h)}(Z), \\ \text{SA}^{(u,h)}(Z) &= \left(W_V^{(u,h)} Z \right) \text{softmax} \left(\left\langle W_Q^{(u,h)} Z, W_K^{(u,h)} Z \right\rangle + R^{(u,h)} \right), \end{aligned} \tag{2}$$

where $W_Q^{(u,h)}, W_K^{(u,h)}, W_V^{(u,h)}, W_O^{(u)}$ $\in \mathbb{R}^{D \times D}$ correspond to the query, key, value and output matrices of the (u, h) -th head, respectively. softmax represents taking softmax normalization across columns. $\langle W_Q^{(u,h)} X, W_K^{(u,h)} X \rangle$ is called the dot-product (DP) structure. Furthermore, $R^{(u,h)} = (R_{i,j}^{(u,h)}) \in \mathbb{R}^{L \times L}$ denotes the additive relative positional encoding matrix, which satisfies $R_{i,j}^{(u,h)} = -\infty$ if $i \leq j$ for the next-token prediction task.

Relative positional encoding (RPE). Throughout this paper, we focus on the Alibi RPE (Press et al., 2022), where $R_{i,j}^{(u,h)}$ follows a Toeplitz structure, i.e., $R_{i,j}^{(u,h)} = \phi(i - j; p^{(u,h)})$ for $i, j \in [L]$. Here, $p^{(u,h)}$'s are learnable parameters and we consider $\phi(\cdot; p)$ of the following form:

$$\phi(z; p) = \begin{cases} -p \cdot (z - 1) & \text{if } z \geq 1 \\ -\infty & \text{otherwise} \end{cases}. \quad (3)$$

Note that we adopt the Alibi RPE only for simplicity and our results can be extended to other additive RPEs, such as T5 (Raffel et al., 2020). However, extending our analysis to the popular rotary RPE (Su et al., 2024) may be nontrivial, and we leave this for future work.

4 FORMULATION AND APPROXIMATION OF INDUCTION HEAD

In this section, we formalize three types of induction head mechanisms with varying levels of complexity. We then theoretically investigate how two-layer single- or multi-head transformers, with or without FFNs, can efficiently implement these mechanisms, highlighting the distinct roles of different transformer submodules

4.1 VANILLA INDUCTION HEADS

The original induction head, proposed in Elhage et al. (2021) and Olsson et al. (2022), is regarded as one of the key mechanisms to implement ICL and reasoning. This induction head suggests that two-layer multi-head transformers without FFNs can execute a simple in-context algorithm to predict the next token b from a context $[\dots ab \dots a]$ through retrieval, copying, and pasting, based on in-context bi-gram pairs, as illustrated in Figure 1.

<START>Mr and Mrs Dursley, of number four, Privet Drive, were proud to say that they were perfectly normal, thank you very much. They were the last people you'd expect to be involved in anything strange or mysterious, because they just didn't hold with such nonsense. Mr Dursley was the director of a firm called Grunnings, which made drills. He was a big, beefy man with hardly any neck, although he did have a very large moustache. Mrs Dursley was thin and blonde and had nearly twice the usual amount of neck, which came in very useful as she spent so much of her time craning over garden fences, spying on the neighbours. The Dursleys had a small son called Dudley and in

Figure 1: An illustration of the original induction head (taken from Elhage et al. (2021)). The induction head proceeds the context $[\dots \text{The D}]$ by retrieving the preceding information most relevant to the current token (D), then copying and pasting the subsequent token (the green `urs`) as the current prediction. Notably, the two self-attention layers focus on the highlighted red and green tokens respectively. For further details, refer to the description below Theorem 4.1.

Formulation of IH_2 . Based on the phenomenon illustrated in Figure 1, we define the vanilla induction head $\text{IH}_2 : \cup_{L \in \mathbb{N}^+} \mathbb{R}^{d \times L} \mapsto \mathbb{R}^d$ as follows:

$$\text{IH}_2(X_L) = (x_s)_{s=2}^{L-1} \text{softmax} \left((x_L^\top W^* x_{s-1})_{s=2}^{L-1} \right)^\top \quad (4)$$

Specifically, IH_2 retrieves in-context information based on the similarities of in-context bi-gram pairs $\{(x_s, x_L)\}_{s=1}^{L-2}$. Note that the magnitude of matrix W^* controls the sparsity of retrieval, since increasing $\|W^*\|$ causes the softmax output to concentrate as a delta measure over the preceding tokens. Additionally, IH_2 can handle input sequences of arbitrary length.

This model retrieves previous tokens x_{s-1} 's that are similar to the current token x_L based on a dot-product similarity, and then copies and pastes x_{s-1} 's subsequent token x_s as the current prediction

x_{L+1} . For example, in Figure 1, the current token x_L is `D`, and the model retrieves previous tokens similar to `D`, copying and pasting its subsequent token `urs` as the prediction.

Comparison with previous formulations. As shown in Figure 1, the current token `D` appears multiple times in the preceding context, and the induction head detects all occurrences of `D`. Our formulation (4) captures this behavior, as the softmax scores for all preceding `D` are identical. In contrast, previous formulations, such as Sanford et al. (2024a) and Sanford et al. (2024b), focus solely on the most recent occurrence of `D`, neglecting this multi-occurrence aspect.

Measure of approximation. Consider a target function $H : \cup_{L \in \mathbb{N}^+} \mathbb{R}^{d \times L} \mapsto \mathbb{R}^d$, where d is the token dimension and L denotes the sequence length. Given an input sequence $X \in \mathbb{R}^{d \times L}$, transformer TF approximates $H(X)$ using its last output token, i.e., $\text{TF}_{-1}(X) \in \mathbb{R}^d$. To quantify the approximation error, we define the following metric: for $1 \leq p \leq +\infty$,

$$\|H - \text{TF}\|_{L,p} := (\mathbb{E}_{X_L} [\|H(X_L) - \text{TF}_{-1}(X_L)\|_{\infty}^p])^{1/p}. \quad (5)$$

The next theorem shows that a two-layer *single-head* transformer *without FFNs* suffices to implement vanilla induction heads.

Theorem 4.1 (two-layer single-head TF w/o FFNs). *Let IH_2 satisfy Eq. (4). Then exists a constant $C > 0$ and a two-layer single-head transformer TF (without FFNs), with $D = 2d$, $W_K^{(1,1)} = W_Q^{(1,1)} = 0$, $p^{(2,1)} = 0$, and $\|W_K^{(2,1)}\|, \|W_Q^{(2,1)}\| \leq \mathcal{O}(1, \|W^*\|_F)$, such that*

$$\sup_{L \in \mathbb{N}^+} \|\text{IH}_2 - \text{TF}\|_{L,\infty} \leq \frac{C}{e^{p^{(1,1)}}}.$$

This theorem shows that single head suffices to approximate the vanilla induction head and moreover, the approximation efficiency is independent of the sequence length. The proof is provided in Appendix A.1, offering the following insights into how two-layer single-head transformers without FFNs implement vanilla induction heads:

- **The first layer** aggregates local tokens and outputs $(z_s = [x_{s-1}, x_s])_{2 \leq s \leq L}$ for the s -th token. This is achieved by using SA with only RPE (no DP). Specifically, RPE allows SA to capture the *preceding token* via $x_{s-1} = \sum_{j \geq 1} x_{s-j} \rho(j)$ for each token x_s , where $\rho(\cdot) = \mathbb{I}\{\cdot = 1\}$. Hence, DP in this layer is not essential and can be omitted.
- **The second layer** extracts the relevant tokens using DP similarity. First, DP computes the similarity $\langle W_Q z_L, W_K z_s \rangle = x_L^\top W^* x_{s-1}$, where $z_L = [x_{L-1}, x_L]$ and $z_s = [x_{s-1}, x_s]$ represent the hidden tokens outputted by the first layer. This similarity measure enables SA to identify tokens that match x_L . Subsequently, the value/output component extracts x_s in z_s , effectively copying the subsequent token and using it as the current prediction. In this layer, RPE is not necessary and can be omitted.

Remark 4.2 (Alignment with experimental findings). Our theoretical analysis is consistent with the experimental observations reported in Elhage et al. (2021). Specifically, the experiments there demonstrate that SA in the first layer attends to adjacent tokens, while SA in the second layer retrieves information related to the current token. Our analysis identifies components responsible for these two operations, and reveals that *single-head* transformers suffice to perform them efficiently.

4.2 GENERALIZED INDUCTION HEADS: IN-CONTEXT n -GRAM AND GENERIC SIMILARITY

Although the standard induction head defined in Eq. (4) is intuitive, it exhibits notable limitations: **1**) it retrieves only a *single token*, potentially missing *complete local information* and leading to false retrievals; **2**) it relies solely on the *dot-product* to measure the similarity between two tokens, which is not sufficiently general.

Formulation of IH_n . Motivated by the limitation **1**) above, we define a generalized induction head:

$$\text{IH}_n(X_L) = (x_s)_{s=n}^{L-1} \text{softmax} \left((X_{L-n+2:L}^\top W^* X_{s-n+1:s-1})_{s=n}^{L-1} \right)^\top, \quad (6)$$

where the patch $X_{s-n+1:s-1}$ incorporates richer local information near x_s and $X_{L-n+2:L}$ denotes the current patches. This formulation is more general than Eq. (4), which only focuses on x_{s-1} . This

induction head operates based on the similarity between the n -gram pairs: $(X_{s-n+1:s-1}; X_{L-n+2:L})$ for $s = n, \dots, L - 1$.

Integrating richer local information facilitates more accurate information retrieval. The model (6) retrieves previous $(n - 1)$ -token patch that are similar to the current $(n - 1)$ -token patch, thereby generalizing the vanilla induction head (4), which considers only single-token retrieval. For example, as depicted in Figure 1, if the current local information is The D (comprising two tokens), and prior local information such as Mr D and Mrs D is identified as similar to The D, transformer would copy and paste their subsequent token, urs, as the prediction.

Theorem 4.3 (two-layer multi-head TF w/o FFNs). *Let IH_n satisfy Eq. (6). Then, for any $q \in \mathbb{N}^+$, there exists an absolute constant $C > 0$ and a two-layer H -head transformer $\text{TF}(\cdot)$ (without FFNs), with $D = nd$, such that:*

$$\sup_{L \in \mathbb{N}^+} \|\text{IH}_n - \text{TF}\|_{L, \infty} \leq C \left(\frac{ne^{1+0.01n}}{H} \right)^q.$$

This theorem demonstrates that two-layer multi-head transformers, even without FFNs, can efficiently implement the generalized induction head (6). Notably, the approximation error scales as $\mathcal{O}(H^{-q})$, where q can be arbitrarily large, and $H \gtrsim ne^{1+0.01n}$ is sufficient to ensure a good approximation. Furthermore, n is typically small when extracting local semantics. For example, in the vanilla induction head, $n = 2$. The proof of this theorem is provided in Appendix A.2.

The role of multiple heads. In Theorem 4.3, multiple heads are employed in the first layer to approximate the n -gram interaction, represented by the $n - 1$ memory kernels $\{\rho_j := \mathbb{I}\{\cdot = j\}\}_{j=1}^{n-1}$. Thus, TF can capture $n - 1$ preceding tokens via $x_{s-j} = \sum_{k \geq 1} x_{s-k} \rho_j(k)$ for $j \in [n - 1]$. Intuitively, as n increases, more memory kernels are required for accurate approximation, necessitating more attention heads. In contrast, Theorem 4.1 only requires approximating a single memory kernel $\mathbb{I}\{\cdot = 1\}$, which can be efficiently achieved using a single attention head.

Recently, Rajaraman et al. (2024) explored a generalized induction head similar to Eq. (6) and showed that multi-layer single-head transformers can implement it. In contrast, our Theorem 4.3 demonstrates that two layers suffice if multi-head self-attention is adopted.

Formulation of GIH_n . Building on the formulation (6), and motivated by the limitation 2) above, we further consider the following generalized induction head:

$$\text{GIH}_n(X_L) = (x_s)_{s=n}^{L-1} \text{softmax} \left(\left(g(X_{L-n+2:L}; X_{s-n+1:s-1}) \right)_{s=n}^{L-1} \right)^\top, \quad (7)$$

where $g : \mathbb{R}^{D \times (n-1)} \times \mathbb{R}^{D \times (n-1)} \rightarrow \mathbb{R}$ denotes a generic function measuring the similarity between two $(n - 1)$ -length patches.

This model retrieves previous relevant multi-token patch $X_{s-n+1:s-1}$ that is similar to the current multi-token patch $X_{L-n+2:L}$, utilizing the generalized similarity function $g(\cdot, \cdot)$. This mechanism is more general than Eq. (6), which is limited to dot-product similarities. For instance, the use of general similarity g enables the model to recognize not only synonymous but also antonymic semantics, thereby improving both the accuracy and diversity of in-context retrievals.

Theorem 4.4 (two-layer multi-head TF with FFNs). *Let GIH_n satisfy Eq. (7). Suppose the similarity function g is α -well-behaved (see Definition A.7). Then, for any $q \in \mathbb{N}^+$, there exist constants $A_{g,q,n}, B_{g,\alpha} > 0$ and a two-layer H -head transformer $\text{TF}(\cdot)$ with FFNs of width M , such that*

$$\|\text{GIH}_n - \text{TF}\|_{L,2} \leq A_{g,q,n} H^{-q} + B_{g,\alpha} L^{1/(1+2\alpha)} M^{-\alpha/(1+3\alpha)}.$$

This theorem establishes that if the similarity function g is well-behaved, two-layer multi-head transformers with FFNs can efficiently implement the generalized induction head (7).

The role of FFNs. In contrast to Theorem 4.3, transformer models in Theorem 4.4 include FFNs. These FFN layers are used to approximate the similarity function g . Specifically, we consider the proper orthogonal decomposition (POD) of g , which can be viewed as an extension of the matrix singular value decomposition (SVD) applied to functions of two variables. For $g : \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$, its POD is $g(u, v) = \sum_{k=1}^{\infty} \sigma_k \phi_k(u) \psi_k(v)$, where ϕ_k, ψ_k are orthonormal bases for $L^2(\mathcal{I})$ (see Appendix D for details). Intuitively, the FFN in the first layer is used to efficiently approximate K bases (ϕ_i 's and ψ_i 's). Then, in the second layer, DP in SA can approximately reconstruct g by using the truncated sum $g(u, v) \approx \sum_{k=1}^K \sigma_k \phi_k(u) \psi_k(v)$. The complete proof is deferred to Appendix A.3.

5 THE TRANSITION FROM LAZY TO RICH MECHANISMS IN LEARNING INDUCTION HEADS

In this section, we investigate the dynamics of learning induction heads using a transformer, particularly focusing on how this differs from n -gram learning. To facilitate the analysis, we consider a mixed target function that comprises a 4-gram component and a vanilla induction head component as defined in Eq. (4). Specifically, we study the gradient flow dynamics of a two-layer multi-head transformer without FFNs on this task.

5.1 SETUPS

5.1.1 MIXED TARGET FUNCTION

Mixed target function. Let the input sequence be $X = (x_1, \dots, x_L) \in \mathbb{R}^{1 \times L}$. Our mixed target function f^* contains both a 4-gram component $f_{\mathbb{G}_4}^*$ and an in-context 2-gram component $f_{\mathbb{H}_2}^*$:

$$f^*(X) := \left(\frac{\alpha^*}{1 + \alpha^*} f_{\mathbb{G}_4}^*(X), \frac{1}{1 + \alpha^*} f_{\mathbb{H}_2}^*(X) \right)^\top \in \mathbb{R}^2, \quad (8)$$

where $\alpha^* > 0$ represents the relative weight between the two components: $f_{\mathbb{G}_4}^*(X)$ and $f_{\mathbb{H}_2}^*(X)$. Here, $f_{\mathbb{G}_4}^*$ represents a 4-gram component and $f_{\mathbb{H}_2}^*$ is given by the vanilla induction head (4) to represent a type of in-context 2-gram information:

$$f_{\mathbb{G}_4}^*(X) := x_{L-2}, \quad f_{\mathbb{H}_2}^*(X) := (x_s)_{s=2}^{L-1} \text{softmax} \left((x_L w^{*2} x_{s-1})_{s=2}^{L-1} \right)^\top.$$

Note that $f_{\mathbb{G}_4}^*$ denotes a ‘‘simplest’’ 4-gram target, where the next token is predicted according to the conditional probability $p(z|X) = p(z|x_L, x_{L-1}, x_{L-2}) = \mathbb{I}\{z = x_{L-2}\}$.

Remark 5.1 (The reason for considering 4-gram). Note that our target includes a 4-gram component rather than simpler 2- or 3-gram components. As suggested by the experimental results in Elhage et al. (2021), for a learned two-layer transformer that implements vanilla induction head \mathbb{H}_2 , the first layer has extracted both x_L and x_{L-1} , which can be outputted using the residual block. Thus, the 2- and 3-gram targets: $p(z|X) = \mathbb{I}\{z = x_L\}$ and $p(z|X) = \mathbb{I}\{z = x_{L-1}\}$ must be learned prior to the induction head. Hence we focus on the more challenging 4-gram target to avoid trivializing the learning process, though our analysis extends straightforwardly to the 2- or 3-gram scenarios.

Remark 5.2 (Extension). Since the transformer studied in this section does not have FFNs, its expressive power is limited. Consequently, we only consider the simple but representative mixed target (8). However, (8) can be generalized to $f^*(X) = F(f_{\mathbb{G}_4}^*(X); f_{\mathbb{H}_2}^*(X))$, where F is general nonlinear function. Such a form can be efficiently approximated by transformers with FFNs. We leave the optimization analysis under this general setting for future work.

5.1.2 TWO-LAYER MULTI-HEAD TRANSFORMER WITH REPARAMETERIZATION

Two-layer multi-head transformer w/o FFNs. We consider a simple two-layer multi-head transformer TF, where the first layer contains a single head $\text{SA}^{(1,1)}$, and the second layer contain two heads $\text{SA}^{(2,1)}$, $\text{SA}^{(2,2)}$. Given an input sequence $X = (x_1, \dots, x_L) \in \mathbb{R}^{1 \times L}$, it is first embedded as $X^{(0)} := (X^\top, 0^\top) \in \mathbb{R}^{2 \times L}$. The model then processes the sequence as follows:

$$\begin{aligned} X^{(1)} &= X^{(0)} + \text{SA}^{(1,1)}(X^{(0)}), \\ \text{TF}(X) &= \text{SA}^{(2,1)}(X^{(1)}) + \text{SA}^{(2,2)}(X^{(1)}). \end{aligned}$$

Reparameterization. Despite the simplification, the transformer above is still too complicated for dynamics analysis. To overcome this challenge, we adopt the reparameterization trick used in previous works (Tian et al., 2023; Huang et al., 2023; Chen et al., 2024b). Specifically, by Theorem 4.1 and its proof, *the first layer does not require DP, and the second layer does not require RPE*. Moreover, to express the 4-gram component $f_{\mathbb{G}_4}^*$, we only need an additional head without DP in the second layer. Therefore, we can reparameterize the model as follows:

- **The first layer.** This layer has only one trainable parameter $p^{(1,1)}$. In the unique head $\text{SA}^{(1,1)}$, DP is removed by setting $W_Q^{(1,1)} = W_K^{(1,1)} = 0$, and we let $W_V^{(1,1)} = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}$. The output

sequence of this layer given by $X^{(1)} = X^{(0)} + \mathbf{SA}^{(1,1)}(X^{(0)}) = \begin{pmatrix} x_1, \dots, x_L \\ y_1, \dots, y_L \end{pmatrix}$, where

$$y_s = (x_\tau)_{\tau=1}^{s-1} \text{softmax} \left(\left(-p^{(1,1)}(s-1-\tau) \right)_{\tau=1}^{s-1} \right)^\top, \quad s \in [L], \quad (9)$$

where $p^{(1,1)}$, used in RPE (3), is the unique trainable parameter in this layer.

- **The second layer.** This layer has 5 trainable parameters: $w_V^{(2,1)}, w_V^{(2,2)}, p^{(2,1)}, w_K^{(2,2)}, w_Q^{(2,2)}$ for parametrizing the two heads. The first head $\mathbf{SA}^{(2,1)}$ without DP is responsible to fit $f_{\mathbf{G}_4}^*$, while the second head $\mathbf{SA}^{(2,2)}$ without RPE is responsible to fit $f_{\mathbf{IH}_2}^*$. Specifically,

$$W_Q^{(2,1)} = W_K^{(2,1)} = 0, W_V^{(2,1)} = \begin{pmatrix} 0 & w_V^{(2,1)} \\ 0 & 0 \end{pmatrix}, p^{(2,2)} = 0, W_V^{(2,2)} = \begin{pmatrix} w_V^{(2,2)} & 0 \\ 0 & 0 \end{pmatrix}.$$

Then the second layer processes $X^{(1)}$ and outputs the last token:

$$\mathbf{TF}_{-1}(X; \theta) = \begin{pmatrix} (w_V^{(2,1)} y_s)_{s=2}^{L-2} \text{softmax} \left(\left(-p^{(2,1)}(L-1-s) \right)_{s=2}^{L-2} \right)^\top \\ (w_V^{(2,2)} x_s)_{s=2}^{L-2} \text{softmax} \left((x_L w_Q^{(2,2)} w_K^{(2,2)} x_{s-1})_{s=2}^{L-2} \right)^\top \end{pmatrix}, \quad (10)$$

where y_s is given by (9). $p^{(2,1)}, w_V^{(2,1)}$ are trainable parameters in $\mathbf{SA}^{(2,1)}$, while $w_Q^{(2,2)}, w_K^{(2,2)}, w_V^{(2,2)}$ are trainable parameters in $\mathbf{SA}^{(2,2)}$.

The set of all six trainable parameters across both layers is denoted by θ .

5.1.3 GRADIENT FLOW ON SQUARE LOSS

We consider the Gaussian input and square loss, both of which are commonly used in analyzing transformer dynamics and ICL (Akyürek et al., 2022; Huang et al., 2023; Wang et al., 2024). The loss is defined as:

$$\mathcal{L}(\theta) = \frac{1}{2} \mathbb{E}_{X \sim \mathcal{N}(0, I_{L \times L})} \left[\|\mathbf{TF}_{-1}(X; \theta) - f^*(X)\|_2^2 \right], \quad (11)$$

To characterize the learning of \mathbf{G}_4 and \mathbf{IH}_2 , we introduce the following two partial losses:

$$\mathcal{L}_{\mathbf{G}_4}(\theta) = \frac{1}{2} \mathbb{E}_X (\mathbf{TF}_{-1,1}(X; \theta) - f_1^*(X))^2, \quad \mathcal{L}_{\mathbf{IH}_2}(\theta) = \frac{1}{2} \mathbb{E}_X (\mathbf{TF}_{-1,2}(X; \theta) - f_2^*(X))^2,$$

which correspond to the two dimensions in $\mathbf{TF}_{-1}(X; \theta) - f^*(X) \in \mathbb{R}^2$, respectively. It follows that $\mathcal{L}(\theta) = \mathcal{L}_{\mathbf{G}_4}(\theta) + \mathcal{L}_{\mathbf{IH}_2}(\theta)$.

Gradient flow (GF). We analyze the GF for minimizing the objective (11):

$$\frac{d\theta(t)}{dt} = -\nabla \mathcal{L}(\theta(t)), \text{ starting with } \theta(0) = (\sigma_{\text{init}}, \dots, \sigma_{\text{init}})^\top, \quad (12)$$

where $0 < \sigma_{\text{init}} \ll 1$ is sufficiently small. Note that $\sigma_{\text{init}} \neq 0$ prevents $\nabla \mathcal{L}(\theta(0)) = 0$.

Layerwise training paradigm. We consider a layerwise training paradigm in which, during each stage, only one layer is trained by GF. Specifically,

- **Training Stage I:** In this phase, only the parameter in the first layer, i.e., $p^{(1,1)}$, is trained.
- **Training Stage II:** In this phase, the first layer parameter $p^{(1,1)}$ keeps fixed and only parameters in the second layer are trained: $w_V^{(2,1)}, w_V^{(2,2)}, p^{(2,1)}, w_Q^{(2,2)}, w_K^{(2,2)}$.

This type of layerwise training has been widely used to study the training dynamics of neural networks, including FFN networks (Safra and Lee, 2022; Bietti et al., 2023; Wang et al., 2023) and transformers (Tian et al., 2023; Nichani et al., 2024; Chen et al., 2024b).

Lemma 5.3 (Training Stage I). *For the Training Stage I, $\lim_{t \rightarrow +\infty} p^{(1,1)}(t) = +\infty$.*

According to (9), this lemma implies that, at the end of Training Stage I, the first layer captures the preceding token x_{s-1} for each token x_s , i.e., $y_s = x_{s-1}$. This property is crucial for transformers to implement induction heads and aligns with our approximation result in Theorem 4.1. The proof of Lemma 5.3 is deferred to Appendix B.

5.2 TRAINING STAGE II: TRANSITION FROM 4-GRAM TO INDUCTION HEAD

In this section, we analyze the dynamics in Training Stage II. We start from the following lemma:

Lemma 5.4 (Parameter balance). *In Training Stage II, it holds that $|w_Q^{(2,2)}(t)|^2 \equiv |w_K^{(2,2)}(t)|^2$.*

Lemma 5.4 is similar to the balance result for homogeneous networks (Du et al., 2018), and its proof can be found at the start of Appendix C. By this lemma, we can define $w_{KQ}^{(2,2)} := w_Q \equiv w_K$. Additionally, Lemma 5.3 ensures that $p^{(1,1)} = +\infty$ holds during Stage II. For simplicity, we denote $w_{V_1} := w_V^{(2,1)}$, $w_{V_2} := w_V^{(2,2)}$, $p := p^{(2,1)}$, $w_{KQ} := w_{KQ}^{(2,2)}$. Consequently, the training dynamics are reduced to four parameters

$$\theta = (w_{V_1}, w_{V_2}, p, w_{KQ}),$$

where we still denote the set of parameters as θ **without introducing ambiguity**. It is important to note that the problem remains **highly non-convex** due to the joint optimization of both inner parameters (p, w_{KQ}) and outer parameters (w_{V_1}, w_{V_2}) in the two heads. At this training stage, GF has a **unique fixed point**:

$$w_{V_1} = \frac{\alpha^*}{1 + \alpha^*}, w_{V_2} = \frac{1}{1 + \alpha^*}, p = +\infty, w_{KQ} = w^*,$$

which corresponds to a global minimizer of the objective (11).

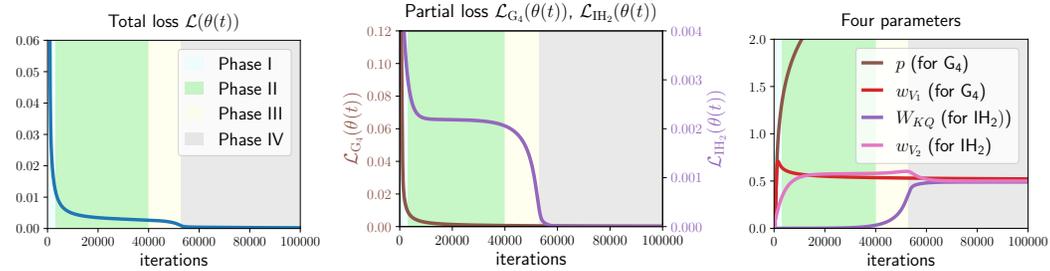


Figure 2: Visualize the dynamical behavior of Training Stage II with total loss, partial loss, and the parameter evolution. Here, $\alpha^* = 1$, $w^* = 0.49$, $\sigma_{\text{init}} = 0.01$, $L = 40$. The is clearly shown that transformer learns the 4-gram component first and then, starts to learn the induction head mechanism. Notably, the entire dynamics unfold in four distinct phases, consistent with our theoretical results (Theorem 5.5). For more experimental details, we refer to Appendix E.1.

As shown in Figure 2, a learning transition from the 4-gram mechanism to the induction head mechanism does occur in our setting. Moreover, the learning process exhibits a four-phase dynamics. The next theorem provides a precise characterization of the four phases, whose proof can be found in Appendix C.

Theorem 5.5 (Learning transition and 4-phase dynamics). *Let $\alpha^* = \Omega(1)$ and $w^* = \mathcal{O}(1)$, and we consider the regime of small initialization ($0 < \sigma_{\text{init}} \ll 1$) and long input sequences ($L \gg 1$). Then we have the following results:*

- **Phase I (partial learning).** *In this phase, most of the 4-gram component in the mixed target is learned, while a considerable number of induction head component have not yet been learned. Specifically, let $T_I = \mathcal{O}(1)$, then we have the following estimates:*

$$\mathcal{L}_{G_4}(\theta(T_I)) \leq 0.01 \cdot \mathcal{L}_{G_4}(\theta(0)), \quad \mathcal{L}_{IH_2}(\theta(T_I)) \geq 0.99 \cdot \mathcal{L}_{IH_2}(\theta(0)).$$

- **Phase II (plateau) + Phase III (emergence).** *In these two phases, the learning of the induction head first gets stuck in a plateau for T_{II} time, then is learned suddenly. Specifically, denoted by an observation time $T_o = \Theta(L)$, we have the following tight estimate of the duration:*

$$T_{II} := \inf \left\{ t > T_o : \mathcal{L}_{IH_2}(\theta(t)) \leq 0.99 \cdot \mathcal{L}_{IH_2}(\theta(T_o)) \right\} = \Theta \left((\alpha^* + 1)^2 L \log(1/\sigma_{\text{init}}) / w^{*2} \right);$$

$$T_{III} := \inf \left\{ t > T_o : \mathcal{L}_{IH_2}(\theta(t)) \leq 0.01 \cdot \mathcal{L}_{IH_2}(\theta(T_o)) \right\} = \Theta \left((\alpha^* + 1)^2 L \log(1/\sigma_{\text{init}}) / w^{*2} \right).$$

486 During these phases, the parameter w_{KQ} (for learning w^* in IH_2) increases exponentially:
 487

$$488 w_{KQ}(t) = \sigma_{\text{init}} \cdot \exp\left(\Theta\left(\frac{w^{*2}t}{(1 + \alpha^*)^2L}\right)\right), t < T_{\text{III}}.$$

- 489
 490
 491 • **Phase IV (convergence).** In this phase, the *loss converges toward zero*. Specifically, the
 492 following convergence rates hold for all $t > T_{\text{III}}$:

$$493 \mathcal{L}_{\text{G}_4}(\theta(t)) = \mathcal{O}\left(\frac{1}{t}\right), \quad \mathcal{L}_{\text{IH}_2}(\theta(t)) = \mathcal{O}\left(\exp\left(-\Omega\left(\frac{w^{*2}t}{(1 + \alpha^*)^2L}\right)\right)\right),$$

494 and $\mathcal{L}(\theta(t)) = \mathcal{L}_{\text{G}_4}(\theta(t)) + \mathcal{L}_{\text{IH}_2}(\theta(t))$.
 495
 496
 497
 498

499 By this theorem, the 4-gram mechanism is first learned, taking time T_{I} . Then, the learning of the
 500 induction head mechanism enters a plateau, taking time T_{II} , followed by a sudden emergence of
 501 learning, taking time $T_{\text{III}} - T_{\text{II}}$. Finally, the loss for both components converges to zero.

502 **The clear learning transition.** When any one of $L, \alpha^*, 1/\sigma_{\text{init}}, 1/w^*$ is sufficiently large, Phase II
 503 lasts for $T_{\text{II}} \gg 1$. During this phase, the 4-gram component has been learned well but the induction
 504 head component remains underdeveloped, demonstrating a distinct learning transition. Moreover,
 505 Theorem 5.5 and its proof reveal two key factors that drive this transition:
 506

- 507 • **Time-scale separation due to high- and low-order parameter dependence in self attention.**
 508 The learning of DP and RPE components differ in their parameter dependencies. DP component
 509 exhibits a quadratic dependence on the parameter w_{KQ} , while RPE component shows linear
 510 dependence on the parameter p . With small initialization $\sigma_{\text{init}} \ll 1$, a clear time-scale separation
 511 emerges: $|\dot{w}_{KQ}| \sim w_{KQ} \ll 1$ (DP, slow dynamics) and $|\dot{p}| \sim 1$ (RPE, fast dynamics).
 512 Consequently, the induction head (fitted by DP) is learned much slower than the 4-gram
 513 component (fitted by RPE). This time-scale separation accounts for the term $\log(1/\epsilon_{\text{init}})$ in the
 514 plateau time T_{II} .
 515 • **Speed difference due to component proportions in the mixed target.** The 4-gram target
 516 component and the induction-head component have differing proportions in the mixed target.
 517 A simple calculation shows: $\mathcal{L}_{\text{G}_4}(0) \sim \alpha^{*2}/(1 + \alpha^*)^2$; If $w^* = \mathcal{O}(1)$, then $\mathcal{L}_{\text{IH}_2}(0) \sim$
 518 $1/[(1 + \alpha^*)^2L]$. Notably, $\mathcal{L}_{\text{IH}_2}(0)$ is significantly smaller than $\mathcal{L}_{\text{G}_4}(0)$. This proportion disparity
 519 accounts for the $(1 + \alpha^*)^2L$ term in the plateau time T_{II} .
 520

521 **Proof idea.** We highlight that our fine-grained analysis of entire learning process is guided by two key
 522 observations: 1) the dynamics of the two heads can be decoupled; 2) there exist a distinct transition
 523 point in the dynamics of each head, as shown in Figure 2 (right). These insights lead us to divide the
 524 analysis of each head into two phases: a monotonic phase and a convergence phase.
 525

526 6 EXPERIMENTAL VALIDATION

527
 528 To further support both our approximation results and optimization dynamics, we conduct a series of
 529 experiments ranging from simple toy models to real-world natural language training tasks. Due to
 530 space constraints, the detailed experimental setups and results are presented in Appendix E.
 531

532 7 CONCLUSION

533
 534 In this work, we present a comprehensive theoretical analysis of how transformers implement
 535 induction heads, examining both the approximation and optimization aspects. From the approximation
 536 standpoint, we identify the distinct roles of each transformer component in implementing induction
 537 heads of varying complexity. On the optimization side, we analyze a toy setting, where we clearly
 538 characterize how learning transitions from n -grams to induction heads. Looking forward, an important
 539 direction for future research is to investigate the dynamics of learning general induction heads, which
 are crucial for realizing stronger ICL capabilities.

REFERENCES

- 540
541
542 Emmanuel Abbe, Samy Bengio, Enric Boix-Adsera, Etai Littwin, and Joshua Susskind. Transformers
543 learn through gradual rank increase. *Advances in Neural Information Processing Systems*, 36,
544 2024. 3
- 545 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement
546 preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing
547 Systems*, 36, 2024. 3
- 548 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-
549 rithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*,
550 2022. 2, 8
- 551 Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes.
552 *arXiv preprint arXiv:1610.01644*, 2016. 42, 43
- 554 Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token
555 selection in attention mechanism. *Advances in Neural Information Processing Systems*, 36:48314–
556 48362, 2023. 3
- 557 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:
558 Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*,
559 2023. 2, 32
- 560 Andrew R Barron. Neural net approximation. In *Proc. 7th Yale Workshop on Adaptive and Learning
561 Systems*, volume 1, pages 69–72, 1992. 19
- 562 Andrew R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE
563 Transactions on Information theory*, 39(3):930–945, 1993. 19
- 564 Andrew R Barron. Approximation and estimation bounds for artificial neural networks. *Machine
565 Learning*, 14(1):115–133, 1994. 19
- 566 Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of transformers
567 to recognize formal languages. *Proceedings of the 2020 Conference on Empirical Methods in
568 Natural Language Processing (EMNLP)*, 2020. 2
- 569 Alberto Bietti, Joan Bruna, and Loucas Pillaud-Vivien. On learning gaussian multi-index models
570 with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023. 8
- 571 Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a
572 transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36, 2024.
573 2, 43
- 574 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
575 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
576 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1
- 577 Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head
578 softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv preprint
579 arXiv:2402.19442*, 2024a. 3
- 580 Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Unveiling induction heads: Provable
581 training dynamics and feature learning in transformers. *arXiv preprint arXiv:2409.10559*, 2024b.
582 3, 7, 8
- 583 Xingwu Chen and Difan Zou. What can transformer learn with varying depth? case studies on
584 sequence learning tasks. *arXiv preprint arXiv:2404.01601*, 2024. 2
- 585 Joy Crosbie and Ekaterina Shutova. Induction heads as an essential mechanism for pattern matching
586 in in-context learning. *arXiv preprint arXiv:2407.07011*, 2024. 1
- 587 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal
588 transformers. *International Conference on Learning Representations*, 2019. 2
- 590
591
592
593

- 594 Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous
595 models: Layers are automatically balanced. *Advances in neural information processing systems*,
596 31, 2018. 9
- 597 Weinan E, Chao Ma, and Lei Wu. A priori estimates of the population risk for two-layer neural
598 networks. *Communications in Mathematical Sciences*, 17(5):1407–1425, 2019. 19
- 600 Weinan E, Chao Ma, and Lei Wu. The barron space and the flow-induced function spaces for neural
601 network models. *Constructive Approximation*, pages 1–38, 2021. 19
- 602 Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable
603 creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages
604 5793–5831. PMLR, 2022. 2, 39
- 606 Benjamin L Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution
607 of statistical induction heads: In-context learning markov chains. *arXiv preprint arXiv:2402.11004*,
608 2024. 2, 43, 44
- 609 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
610 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli,
611 Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal
612 Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris
613 Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
614 <https://transformer-circuits.pub/2021/framework/index.html>. 1, 2, 3, 4, 5, 7
- 615 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn
616 in-context? a case study of simple function classes. *Advances in Neural Information Processing*
617 *Systems*, 35:30583–30598, 2022. 2
- 619 Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dimitris
620 Papailiopoulos. Looped transformers as programmable computers. *International Conference on*
621 *Machine Learning*, 2023. 2
- 622 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
623 *preprint arXiv:2312.00752*, 2023. 1
- 625 Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions of*
626 *the Association for Computational Linguistics*, 8:156–171, 2020. 2
- 627 Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. *arXiv preprint*
628 *arXiv:2310.05249*, 2023. 3, 7, 8
- 630 Samy Jelassi, David Brandfonbrener, Sham M Kakade, and Eran Malach. Repeat after me: Trans-
631 formers are better than state space models at copying. *arXiv preprint arXiv:2402.01032*, 2024.
632 1
- 633 Haotian Jiang and Qianxiao Li. Approximation theory of transformer networks for sequence modeling.
634 *arXiv preprint arXiv:2305.18475*, 2023. 2, 19
- 636 Hongkang Li, Meng Wang, Songtao Lu, Xiaodong Cui, and Pin-Yu Chen. Training nonlinear
637 transformers for efficient in-context learning: A theoretical learning and generalization analysis.
638 *arXiv preprint arXiv:2402.15607*, 2024. 3
- 639 Chao Ma, Stephan Wojtowytsch, Lei Wu, and Weinan E. Towards a mathematical understanding
640 of neural network-based machine learning: what we know and what we don’t. *arXiv preprint*
641 *arXiv:2009.10713*, 2020. 19
- 642 Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is
643 provably the optimal in-context learner with one layer of linear self-attention. *arXiv preprint*
644 *arXiv:2307.03576*, 2023. 2
- 645 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
646 models. *arXiv preprint arXiv:1609.07843*, 2016. 40, 41

- 648 William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought.
649 *arXiv preprint arXiv:2310.07923*, 2023. 2
- 650
- 651 William Merrill, Ashish Sabharwal, and Noah A Smith. Saturated transformers are constant-depth
652 threshold circuits. *Transactions of the Association for Computational Linguistics*, 10:843–856,
653 2022. 2
- 654 Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with
655 gradient descent. *arXiv preprint arXiv:2402.14735*, 2024. 3, 8
- 656
- 657 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
658 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads.
659 *arXiv preprint arXiv:2209.11895*, 2022. 1, 2, 3, 4
- 660 Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing complete. *The Journal of*
661 *Machine Learning Research*, 22(1):3463–3497, 2021. 2
- 662
- 663 Ofir Press, Noah A Smith, and Mike Lewis. Train short, test long: Attention with linear biases
664 enables input length extrapolation. *International Conference on Learning Representations*, 2022. 4
- 665
- 666 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
667 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
668 transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 4
- 669 Nived Rajaraman, Marco Bondaschi, Kannan Ramchandran, Michael Gastpar, and Ashok Vardhan
670 Makkuva. Transformers on markov data: Constant depth suffices. *arXiv preprint arXiv:2407.17686*,
671 2024. 2, 6
- 672
- 673 Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context
674 classification task. *International Conference on Learning Representations*, 2024. 2
- 675 Michael Reed and Barry Simon. *Methods of modern mathematical physics: Functional analysis*,
676 volume 1. Gulf Professional Publishing, 1980. 19
- 677
- 678 Itay Safran and Jason Lee. Optimization-based separations for neural networks. In *Conference on*
679 *Learning Theory*, pages 3–64. PMLR, 2022. 8
- 680 Clayton Sanford, Daniel Hsu, and Matus Telgarsky. One-layer transformers fail to solve the induction
681 heads task. *arXiv preprint arXiv:2408.14332*, 2024a. 2, 5
- 682
- 683 Clayton Sanford, Daniel Hsu, and Matus Telgarsky. Transformers, parallel computation, and logarithmic
684 depth. *arXiv preprint arXiv:2402.09268*, 2024b. 1, 2, 5
- 685 Johannes Schmidt-Hieber et al. Nonparametric regression using deep neural networks with ReLU
686 activation function. *Annals of Statistics*, 48(4):1875–1897, 2020. 20
- 687
- 688 Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical*
689 *journal*, 27(3):379–423, 1948. 1, 3
- 690
- 691 Lingfeng Shen, Aayush Mishra, and Daniel Khashabi. Do pretrained transformers really learn
692 in-context by gradient descent? *arXiv preprint arXiv:2310.08540*, 2023. 2
- 693 Jonathan W Siegel and Jinchao Xu. Approximation rates for neural networks with general activation
694 functions. *Neural Networks*, 128:313–321, 2020. 19
- 695
- 696 Jiajun Song, Zhuoyan Xu, and Yiqiao Zhong. Out-of-distribution generalization via composition: a
697 lens through induction heads in transformers. *arXiv preprint arXiv:2408.09503*, 2024. 1
- 698
- 699 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced
700 transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024. 4
- 701 Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as
support vector machines. *arXiv preprint arXiv:2308.16898*, 2023. 3

- 702 Christos Thrampoulidis. Implicit bias of next-token prediction. *arXiv preprint arXiv:2402.18551*,
703 2024. 3
- 704
- 705 Yuandong Tian, Yiping Wang, Beidi Chen, and Simon S Du. Scan and snap: Understanding
706 training dynamics and token composition in 1-layer transformer. *Advances in Neural Information*
707 *Processing Systems*, 36:71911–71947, 2023. 3, 7, 8
- 708 Bhavya Vasudeva, Puneesh Deora, and Christos Thrampoulidis. Implicit bias and fast convergence
709 rates for self-attention. *arXiv preprint arXiv:2402.05738*, 2024. 3
- 710
- 711 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
712 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
713 *systems*, 30, 2017. 1, 3
- 714 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
715 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In
716 *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023. 2
- 717
- 718 Mingze Wang and Weinan E. Understanding the expressive power and mechanisms of transformer
719 for sequence modeling. *Advances in Neural Information Processing Systems*, 2024. 39
- 720 Zihao Wang, Eshaan Nichani, and Jason D Lee. Learning hierarchical polynomials with three-layer
721 neural networks. *arXiv preprint arXiv:2311.13774*, 2023. 8
- 722
- 723 Zixuan Wang, Stanley Wei, Daniel Hsu, and Jason D Lee. Transformers provably learn sparse token
724 selection while fully-connected nets cannot. *arXiv preprint arXiv:2406.06893*, 2024. 2, 3, 8
- 725 Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on
726 approximating turing machines with transformers. *Advances in Neural Information Processing*
727 *Systems*, 35:12071–12083, 2022. 2
- 728
- 729 Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Conference*
730 *on Machine Learning*, pages 11080–11090. PMLR, 2021. 2
- 731 Norman Yarvin and Vladimir Rokhlin. Generalized gaussian quadratures and singular value de-
732 compositions of integral operators. *SIAM Journal on Scientific Computing*, 20(2):699–718, 1998.
733 19
- 734 Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar.
735 Are transformers universal approximators of sequence-to-sequence functions? *arXiv preprint*
736 *arXiv:1912.10077*, 2019. 2
- 737
- 738 Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.
739 *arXiv preprint arXiv:2306.09927*, 2023. 3
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

Appendix

756		
757		
758		
759		
760		
761	A Proofs in Section 4	15
762		
763	A.1 Proof of Theorem 4.1	15
764	A.2 Proof of Theorem 4.3	16
765	A.3 Proof of Theorem 4.4	19
766		
767		
768	B Proofs of Optimization Dynamics: Training Stage I	24
769		
770		
771	C Proofs of Optimization Dynamics: Training Stage II	26
772	C.1 Dynamics of the parameters for 4-gram	27
773	C.2 Dynamics of the parameters for induction head	32
774	C.3 Proof of Theorem 5.5	37
775		
776		
777	D Useful Inequalities	39
778		
779	E Experiments	40
780		
781	E.1 Experimental details for Figure 2	40
782	E.2 Additional experiments supporting optimization dynamics	40
783	E.3 Experiments supporting approximation results	42
784		
785		
786	F Detailed Comparison with Related Works	43
787		
788		
789		
790		
791	A PROOFS IN SECTION 4	
792		
793	A.1 PROOF OF THEOREM 4.1	
794		
795	$\text{IH}_2(X_L) = (x_s)_{s=2}^{L-1} \text{softmax} \left((x_L^\top W^* x_{s-1})_{s=2}^{L-1} \right)^\top, \quad (13)$	
796	Theorem A.1 (Restatement of Theorem 4.1). <i>Let IH_2 satisfy Eq. (13). Then, there exists a constant</i>	
797	<i>$C > 0$ and a two-layer single-head transformer TF (without FFNs), with $D = 2d$, $W_K^{(1,1)} =$</i>	
798	<i>$W_Q^{(1,1)} = 0$, $p^{(2,1)} = 0$, and $\ W_K^{(2,1)}\ , \ W_Q^{(2,1)}\ \leq \mathcal{O}(1, \ W^*\ _F)$, such that</i>	
799		
800	$\sup_{L \in \mathbb{N}^+} \ \text{IH}_2 - \text{TF}\ _{L, \infty} \leq \frac{C}{e^{p^{(1,1)}}}.$	
801		
802		
803	<i>Proof.</i> We consider two-layer single-head transformer without FFN, where the first layer has the	
804	residual block, while the second layer does not have the residual block.	
805		
806	We first embed each token into \mathbb{R}^D as $\begin{pmatrix} x_s \\ 0 \end{pmatrix}$ and take $W_V^{(1)} = \begin{pmatrix} 0 & 0 \\ I_{d \times d} & 0 \end{pmatrix}$, then the s -th output	
807	token of the first layer is	
808		
809	$\begin{pmatrix} x_s \\ y_s \end{pmatrix} = \begin{pmatrix} x_s \\ (x_\tau)_{\tau=1}^{s-1} \text{softmax} \left((-p^{(1,1)}(s-1-\tau))_{\tau=1}^{s-1} \right) \end{pmatrix}.$	

Then for the second layer, we choose $p^{(2,1)} = 0$,

$$W_Q^{(2,1)} = \begin{pmatrix} 0 & 0 \\ I_{d \times d} & 0 \end{pmatrix}, W_K^{(2,1)} = \begin{pmatrix} 0 & 0 \\ 0 & W^* \end{pmatrix}, W_V^{(2,1)} = \begin{pmatrix} I_{d \times d} & 0 \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{D \times D},$$

and the projection $W_O^{(2)} = (I_{d \times d} \ 0_{d \times d}) \in \mathbb{R}^{d \times D}$.

Then the last output token of the second layer is

$$(x_s)_{s=2}^{L-1} \text{softmax} \left((x_L^\top W^* y_s)_{s=2}^{L-1} \right)^\top.$$

By Lemma D.1, for any $L \in \mathbb{N}^+$

$$\begin{aligned} & \| \mathbf{H}_2 - \mathbf{TF} \|_{L, \infty} \\ &= \sup_{X_L} \| \mathbf{H}(X_L) - \mathbf{TF}_{-1}(X_L) \|_\infty \\ &= \left\| (x_s)_{s=2}^{L-1} \text{softmax} \left((x_L^\top W^* y_s)_{s=2}^{L-1} \right)^\top - (x_s)_{s=2}^{L-1} \text{softmax} \left((x_L^\top W^* x_{s-1})_{s=2}^{L-1} \right)^\top \right\|_\infty \\ &\leq \| (x_s)_{s=2}^{L-1} \|_{\infty, \infty} \left\| \text{softmax} \left((x_L^\top W^* y_s)_{s=2}^{L-1} \right) - \text{softmax} \left((x_L^\top W^* x_{s-1})_{s=2}^{L-1} \right) \right\|_1 \\ &\leq 2 \sup_{2 \leq s \leq L-1} |x_L^\top W^* y_s - x_L^\top W^* x_{s-1}| \\ &\leq 2 \|x_L^\top W^*\|_1 \sup_s \|y_s - x_{s-1}\|_\infty \\ &\leq 2 \sum_{i,j} |W_{i,j}^*| \sup_s \left\| (x_\tau)_{\tau=1}^{s-1} \text{softmax} \left(\left(-p^{(1,1)}(s-1-\tau) \right)_{\tau=1}^{s-1} \right)^\top - x_{s-1} \right\|_\infty \\ &\leq 2 \|W^*\|_{1,1} \sup_s \left\| \text{softmax} \left(\left(-p^{(1,1)}(s-1-\tau) \right)_{\tau=1}^{s-1} \right) - \mathbf{e}_{s-1} \right\|_1 \\ &\leq 4 \|W^*\|_{1,1} \frac{e^{-p^{(1,1)}}}{1 - e^{-p^{(1,1)}}} \leq \mathcal{O} \left(e^{-p^{(1,1)}} \right). \end{aligned}$$

□

A.2 PROOF OF THEOREM 4.3

$$\mathbf{H}_n(X_L) = (x_s)_{s=n}^{L-1} \text{softmax} \left((X_{L-n+2:L}^\top W^* X_{s-n+1:s-1})_{s=n}^{L-1} \right)^\top, \quad (14)$$

Theorem A.2 (Restatement of Theorem 4.3). *Let \mathbf{H}_n satisfy Eq. (14). Then for any $q \in \mathbb{N}^+$, there exists a constant $C_{q,n} > 0$ and a two-layer H -head transformer $\mathbf{TF}(\cdot)$ (without FFNs), with $D = nd$, such that:*

$$\sup_{L \in \mathbb{N}^+} \| \mathbf{H}_n - \mathbf{TF} \|_{L, \infty} \leq C \left(\frac{ne^{1+0.01n}}{H} \right)^q.$$

Proof. We consider two-layer multi-head transformer without FFN, where the first layer has the residual block, while the second layer does not have the residual block.

First, we choose the embedding dimension $D = nd$, and parameters in the embedding map

$$W_E = \begin{pmatrix} I_{d \times d} \\ 0_{(D-d) \times d} \end{pmatrix} \in \mathbb{R}^{D \times d}, \quad b_E = 0 \in \mathbb{R}^D,$$

then each token $x_s^{(0)}$ after embedding is

$$x_s^{(0)} = W_E x_s + b_E = \begin{pmatrix} x_s \\ 0 \end{pmatrix} \in \mathbb{R}^D.$$

This proof can be summarized as the following process for TF_{-1} :

$$\begin{aligned} & (x_s)_{s=n}^{L-1} \text{softmax} \left((X_{L-n+2:L}^\top W^* X_{s-n+1:s-1})_{s=n}^{L-1} \right)^\top \\ & \quad \text{Step II. 2-st Attn } \uparrow \\ & \quad \quad X_{L-n+2:L} \\ & \quad \text{Step I. 1-st Attn } \uparrow \\ & \quad \quad (x_L^\top, 0_{D-d}^\top)^\top \end{aligned}$$

Step I. The first layer. We use 1-st Attn with residual to copy the previous tokens $(x_{s-n+1}, \dots, x_{s-1})$ of each token x_s . We use $H = \sum_{i=1}^{n-1} H_i$ attention heads to realize this step, and the following projection matrices are needed:

$$P_i := (0_{d \times id}, I_{d \times d}, 0_{d \times (D-(i+1)d)}) \in \mathbb{R}^{d \times D}, \quad i = 1, \dots, n-1.$$

By lemma D.2, there exist a constant $C > 0$ such that: for any rate $q \in \mathbb{N}^+$, there exists a function

$$\phi_i^{\text{exp}}(t) = \sum_{h=1}^{H_i} \alpha_{h,i} e^{-\beta_{h,i}(t-1)}$$

such that $\beta_h > 0$ and

$$\|\mathbb{I}\{\cdot = i\} - \phi_i^{\text{exp}}(\cdot)\|_{\ell_1(\mathbb{N})} = \sum_{s=i}^{+\infty} |\mathbb{I}\{s = 1\} - \phi^{\text{exp}}(s)| \leq \frac{C e^{q+0.01(q+1)i}}{H_i^q}$$

For $h = \sum_{j=1}^{i-1} H_j, 1 + \sum_{j=1}^{i-1} H_j, \dots, \sum_{j=1}^i H_j$, we choose parameters as follows

$$\begin{aligned} p^{(1,h)} &= \beta_{h,i}, \quad W_V^{(1,h)} = \alpha_{h,i} \left(\sum_{j=0}^{H_i} \exp(-\beta_{h,i}(j-1)) \right) S_i, \\ W_K^{(1,h)} &= W_Q^{(1,h)} = 0, \quad W_O^{(1)} = I_{D \times D} \end{aligned}$$

where $S_i \in \mathbb{R}^{D \times D}$ is a shift matrix that takes out the first d elements of a vector and shifts it backward to the $(id+1)$ -th to $(i+1)d$ -th elements. Then

$$\left(P_i \sum_{h=\sum_{j=1}^{i-1} H_j}^{\sum_{j=1}^i H_j} \text{SA}^{(1,h)}(X_L^{(0)}) \right)_{-1} = \sum_{h=\sum_{j=1}^{i-1} H_j}^{\sum_{j=1}^i H_j} \alpha_{h,i} \sum_{s=1}^{L-1} e^{-\beta_{h,i}(s-1)} x_{L-s}.$$

We denote $x_L^{(1)} := \text{SA}^{(1)}(X_L^{(0)})_{-1}$, then the approximation error of this step is

$$\begin{aligned} \varepsilon_{\text{SA}}^{(1)} &:= \sup_s \left\| x_s^{(1)} - \begin{pmatrix} x_s \\ \vdots \\ x_{s-n+1} \end{pmatrix} \right\|_\infty \leq \sup_s \sum_{i=1}^{n-1} \|P_i x_s^{(1)} - x_{s-i}\|_\infty \\ &\leq \sup_s \sum_{i=1}^{n-1} \|\mathbb{I}\{\cdot = i\} - \phi_i^{\text{exp}}(\cdot)\|_{\ell_1(\mathbb{N})} \leq C e^q \sum_{i=1}^{n-1} \frac{e^{0.01(q+1)i}}{H_i^q}. \end{aligned}$$

Consequently, one detail is to assign the head number $\{H_i\}_{i=1}^n$ such that the error's sum $\sum_{i=1}^{n-1} \frac{e^{0.01(q+1)i}}{H_i^q}$ is as small as possible. Our way is solving the minimization problem

$$\begin{aligned} \min : & \sum_{i=1}^{n-1} \frac{e^{0.01(q+1)i}}{H_i^q} \\ \text{s.t.} & \sum_{i=1}^{n-1} H_i = H, \end{aligned}$$

which suggests that we should choose the head number:

$$H_i = \frac{e^{0.01i}}{\sum_{j=1}^{n-1} e^{0.01j}}, i \in [n-1].$$

Thus, we obtain the bound

$$\varepsilon_{\text{SA}}^{(1)} \leq \frac{Ce^q}{H^q} \left(\sum_{i=1}^{n-1} e^{0.01i} \right)^q \leq C \left(\frac{ene^{0.01n}}{H} \right)^q.$$

Additionally, we denote the output of this step as

$$x_s^{(1)} := \begin{pmatrix} x_s \\ \hat{x}_{s-1} \\ \vdots \\ \hat{x}_{s-n+1} \end{pmatrix} := \hat{X}_{s-n+1:s-1}.$$

We choose H large enough so that $x_s^{(1)} \in [-2, 2]^D$.

Step II. The second layer. For the second Attn, we only need use the first head (by setting $W_V^{(2,h)} = 0$ for $h \geq 1$). Specifically, we choose $p^{(2,1)} = 0$,

$$W_Q^{(2,1)} = \begin{pmatrix} 0 & 0 \\ I_{(D-d) \times (D-d)} & 0 \end{pmatrix}, W_K^{(2,1)} = \begin{pmatrix} 0 & 0 \\ 0 & W^* \end{pmatrix}, W_V^{(2)} = \begin{pmatrix} I_{d \times d} & 0 \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{D \times D},$$

and the projection $W_O^{(2)} = (I_{d \times d} \quad 0_{(D-d) \times d}) \in \mathbb{R}^{d \times D}$.

Then the output of this layer is

$$x_L^{(2)} = (x_s)_{s=n}^{L-1} \text{softmax} \left(\left(\hat{X}_{L-n+2:L}^\top W^* \hat{X}_{s-n+1:s-1} \right)_{s=n}^{L-1} \right)^\top$$

According to Lemma D.1,

$$\begin{aligned} & \left\| x_L^{(2)} - (x_s)_{s=n}^{L-1} \text{softmax} \left(\left(X_{L-n+2:L}^\top W^* X_{s-n+1:s-1} \right)_{s=n}^{L-1} \right)^\top \right\|_\infty \\ & \leq \sum_{s=n}^{L-1} \left| \text{softmax} \left(\left(\hat{X}_{L-n+2:L}^\top W^* \hat{X}_{s-n+1:s-1} \right)_{s=n}^{L-1} \right)^\top \right. \\ & \quad \left. - \text{softmax} \left(\left(X_{L-n+2:L}^\top W^* X_{s-n+1:s-1} \right)_{s=n}^{L-1} \right)^\top \right| \\ & \leq 2 \max_s \left| \hat{X}_{L-n+2:L}^\top W^* \hat{X}_{s-n+1:s-1} - X_{L-n+2:L}^\top W^* X_{s-n+1:s-1} \right| \\ & \leq 2 \|W^*\|_{(1,1)} \cdot \varepsilon_{\text{SA}}^{(1)}. \end{aligned}$$

Since the above inequality holds for any L and X_L , we have:

$$\sup_{L \in \mathbb{N}^+} \|\text{IH}_n - \text{TF}\|_{L,\infty} \leq C \left(\frac{ne^{1+0.01n}}{H} \right)^q.$$

Additionally, our proof primarily focuses on the case of $H \geq n$. For the case of $H < n$, the approximation error can be trivially bounded by:

$$\sup_{L \in \mathbb{N}^+} \|\text{IH}_n - \text{TF}\|_{L,\infty} \leq \sup_{L \in \mathbb{N}^+} \|\text{IH}_n - 0\|_{L,\infty} \leq 1 \leq C \left(\frac{ne^{1+0.01n}}{H} \right)^q.$$

Then, the two cases can be unified. □

972 A.3 PROOF OF THEOREM 4.4

973 A.3.1 APPROXIMATION RESULTS FOR FFNS

974 Since the setting in this subsection includes FFNs, we introduce the following preliminary results
975 about the approximation of FFNs.

976 The well-known universal approximation result for two-layer FNNs asserts that two-layer FNNs
977 can approximate any continuous function (Barron, 1992; 1993; 1994). Nonetheless, this result lacks
978 a characterization of the approximation efficiency, i.e., how many neurons are needed to achieve
979 a certain approximation accuracy? Extensive pre-existing studies aimed to address this gap by
980 establishing approximation rates for two-layer FFNs. A representative result is the Barron theory (E
981 et al., 2019; 2021; Ma et al., 2020): any function f in Barron space \mathcal{B} can be approximated by a
982 two-layer FFN with M hidden neurons can approximate f efficiently, at a rate of $\mathcal{O}(\|f\|_{\mathcal{B}}/\sqrt{M})$.
983 This rate is remarkably independent of the input dimension, thus avoiding the Curse of Dimensionality.
984 Specifically, Barron space is defined in as follows:

985 **Definition A.3** (Barron space (E et al., 2019; 2021; Ma et al., 2020)). Consider functions $f : X \rightarrow \mathbb{R}$
986 that admit the following representation: $f(x) = \int_{\Omega} a\sigma(b^{\top}x + c)\rho(\mathrm{d}a, \mathrm{d}b, \mathrm{d}c)$, $x \in X$. For any
987 $p \in [1, +\infty]$, we define the Barron norm as $\|f\|_{\mathcal{B}_p} := \inf_{\rho} \left(\mathbb{E}_{\rho} [|a|^p (\|b\|_1 + |c|)^p] \right)^{1/p}$. Then the
988 Barron space are defined as: $\mathcal{B}_p := \{f \in \mathcal{C} : \|f\|_{\mathcal{B}_p} < +\infty\}$.

989 **Proposition A.4** (E et al. (2019)). For any $p \in [1, +\infty]$, $\mathcal{B}_p = \mathcal{B}_{\infty}$ and $\|f\|_{\mathcal{B}_p} = \|f\|_{\mathcal{B}_{\infty}}$.

990 **Remark A.5.** From the Proposition above, the Barron spaces \mathcal{B}_p are equivalent for any $p \in [1, +\infty]$.
991 Consequently, in this paper, we use \mathcal{B} and $\|\cdot\|_{\mathcal{B}}$ to denote the Barron space and Barron norm.

992 The next lemma illustrates the approximation rate of two-layer FFNs for Barron functions.

993 **Lemma A.6** (Ma et al. (2020)). For any $f \in \mathcal{B}$, there exists a two-layer ReLU neural network
994 FFN(x) = $\sum_{k=1}^M a_w \sigma(b_k^{\top} x + c_k)$ with M neurons such that

$$995 \|f - \text{FFN}\|_{L^{\infty}([0,1]^d)} \leq \tilde{\mathcal{O}} \left(\frac{\|f\|_{\mathcal{B}}}{\sqrt{M}} \right).$$

1000 A.3.2 PROPER ORTHOGONAL DECOMPOSITION

1001 Proper orthogonal decomposition (POD) can be viewed as an extension of the matrix singular value
1002 decomposition (SVD) applied to functions of two variables. Specifically, for a square integrable
1003 function $g : \mathcal{I} \times \mathcal{I} \rightarrow \mathbb{R}$, it has the following decomposition (Theorem 3.4 in Yarvin and Rokhlin
1004 (1998), Theorem VI.17 in Reed and Simon (1980)):

$$1005 g(u, v) = \sum_{k=1}^{\infty} \sigma_k \phi_k(u) \psi_k(v). \quad (15)$$

1006 Here, ϕ_k, ψ_k are orthonormal bases for $L^2(\mathcal{I})$, and $\sigma_k \geq 0$ are the singular values, arranged in
1007 descending order.

1008 Recently, Jiang and Li (2023) also used POD to study the approximation rate of single-layer single-
1009 head Transformer for the targets with nonlinear temporal kernels.

1010 Given that two-layer FFNs can efficiently approximate Barron functions (Ma et al., 2020), which is
1011 dense in $L^2([0, 1]^d)$ (Siegel and Xu, 2020), we introduce the following technical definition regarding
1012 the well-behavior POD, which is used for our theoretical analysis.

1013 **Definition A.7** (Well-behaved POD). Let the POD of $g : [-2, 2]^D \times [-2, 2]^D \mapsto \mathbb{R}$ be $g(u, v) =$
1014 $\sum_{k=1}^{\infty} \sigma_k \phi_k(u) \psi_k(v)$. We call the function g has α -well-behaved POD ($\alpha > 0$) if:

- 1015 • The decay rate of singular values satisfies $\sigma_k = \mathcal{O}(1/k^{1+\alpha})$;
- 1016 • The L^{∞} norms, Barron norms, and Lipschitz norms of the POD bases are all uniformly bounded:
1017 $\sup_k \left(\|\phi_k\|_{L^{\infty}} \vee \|\psi_k\|_{L^{\infty}} \vee \|\phi_k\|_{\mathcal{B}} \vee \|\psi_k\|_{\mathcal{B}} \vee \|\phi_k\|_{\text{Lip}} \vee \|\psi_k\|_{\text{Lip}} \right) < \infty$.

A.3.3 PROOF OF THEOREM 4.4

$$\text{GIH}_n(X_L) = (x_s)_{s=n}^{L-1} \text{softmax} \left((g(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right)^\top, \quad (16)$$

Theorem A.8 (Restatement of Theorem 4.4). *Let GIH_n satisfy Eq. (16). Suppose the similarity function g is α -well-behaved (see Definition A.7). Then, for any $q > 0$, there exist constants $A_{g,q,n}, B_{g,\alpha} > 0$ and a two-layer H -head transformer $\text{TF}(\cdot)$ with FFN of width M , such that the following approximation rate holds:*

$$\|\text{GIH}_n - \text{TF}\|_{L,2} \leq \frac{A_{g,q,n}}{H^n} + \frac{B_{g,\alpha} L^{1/(1+2\alpha)}}{M^{\alpha/(1+3\alpha)}}.$$

Proof. We consider two-layer multi-head transformer with FFN, where the first layer has the residual block.

First, we set an constant $K \in \mathbb{N}^+$, and we will optimize it finally. We choose the embedding dimension $D = nd + 2(n-1)K$, and the flowchart of the theorem proof is as follows:

$$\begin{aligned} & (x_s)_{s=n}^{L-1} \text{softmax} \left((g(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right)^\top \\ & \quad \text{Step III. 2-st Attn } \uparrow \\ & (x_L^\top, \dots, \hat{x}_{L-n+1}, \hat{\phi}_1(\hat{X}_{L-n+2:L}), \dots, \hat{\phi}_K(\hat{X}_{L-n+2:L}), \\ & \quad \hat{\psi}_1(\hat{X}_{L-n+1:L-1}), \dots, \hat{\psi}_K(\hat{X}_{L-n+1:L-1}))^\top \\ & \quad \text{Step II. 1-st FFN } \uparrow \\ & (x_L^\top, \hat{x}_{L-1}, \dots, \hat{x}_{L-n+1}, \mathbf{0}^\top)^\top \\ & \quad \text{Step I. 1-st Attn } \uparrow \\ & (x_t^\top, \mathbf{0}^\top)^\top \end{aligned}$$

Recalling Definition A.7, there exists constants $C_g^\infty, C_g^\mathcal{B}, C_g^{\text{Lip}} > 0$ such that:

$$\sup_k (\|\phi_k\|_\infty \vee \|\psi_k\|_\infty) \leq C_g^\infty, \sup_k (\|\phi_k\|_\mathcal{B} \vee \|\psi_k\|_\mathcal{B}) \leq C_g^\mathcal{B}, \sup_k (\|\phi_k\|_{\text{Lip}} \vee \|\psi_k\|_{\text{Lip}}) \leq C_g^{\text{Lip}}.$$

Additionally, $\sigma_k = \mathcal{O}(1/k^{1+\alpha})$ implies that there exists a $C_\alpha > 0$ such that:

$$\sum_{k=K}^{\infty} \sigma_k < \frac{C_\alpha}{K^\alpha}, \quad \forall K \geq 1.$$

Step I: Error in 1-st Attn layer. This step is essentially the same as Step I in the proof of Theorem 4.3, so we write down the error of the first Attn layer directly:

$$\epsilon_{\text{SA}}^{(1)} \leq \frac{C_{q,n}}{H^n}.$$

Moreover, due to $\|\hat{X}_{s-n+2:s} - X_{s-n+2:s}\| \leq \epsilon_{\text{SA}}^{(1)}$, for all s , we have:

$$\hat{X}_{s-n+2:s} \in [-2, 2]^D.$$

Step II: Error in 1-st FFN layer. The 1-st FFN is used to approximate ϕ_k, ψ_k ($k = 1, \dots, K$). Each function is approximated by a 2-layer neural networks with $\frac{M}{2K}$ neurons defined on \mathbb{R}^D , and the FFNs are concatenated together (refer to section 7.1 "Parallelization" in Schmidt-Hieber et al. (2020)) as FFN⁽¹⁾. We denote them as

$$\hat{\phi}_k(y) = \sum_{m=1}^{\frac{M}{2K}} a_m^k \sigma(b_m^{k^\top} y + c_m^k)$$

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

$$\hat{\psi}_k(y) = \sum_{m=1}^{\frac{M}{2K}} \tilde{a}_m^k \sigma(\tilde{b}_m^{k\top} y + \tilde{c}_m^k)$$

Then according to lemma A.6, such FFNs exist and satisfy the following properties hold for all $1 \leq k \leq K$:

$$\begin{aligned} \|\hat{\phi}_k - \phi_k\|_{L^\infty([-2,2]^D)} &\leq \tilde{O} \left(\|\phi_k\|_{\mathcal{B}} \sqrt{\frac{K}{M}} \right) \leq \epsilon_{\text{FFN}}^{(1)}, \\ \|\hat{\psi}_k - \psi_k\|_{L^\infty([-2,2]^D)} &\leq \tilde{O} \left(\|\psi_k\|_{\mathcal{B}} \sqrt{\frac{K}{M}} \right) \leq \epsilon_{\text{FFN}}^{(1)}, \end{aligned}$$

where

$$\epsilon_{\text{FFN}}^{(1)} := \cdot \tilde{O} \left(C_g^{\mathcal{B}} \sqrt{\frac{K}{M}} \right).$$

Step III: Error in 2nd Attn layer.

We use matrices in the second layer to take out elements needed

$$\begin{aligned} W_V^{(2)} &= (I_{d \times d}, 0_{d \times D}) \in \mathbb{R}^{d \times D}, \\ W_K^{(2,1)} &= \sum_{i=k}^K \sqrt{\sigma_k} e_{k, (n-1)d+k} \in \mathbb{R}^{D \times D}, \\ W_Q^{(2,1)} &= \sum_{k=1}^K \sqrt{\sigma_k} e_{k, (n-1)d+K+k} \in \mathbb{R}^{D \times D}. \end{aligned}$$

We denote the rank- K truncation of g as

$$g_K := \sum_{k=1}^K \sigma_k \phi_k \psi_k,$$

and its approximation as

$$\hat{g}_K := \sum_{k=1}^K \sigma_k \hat{\phi}_k \hat{\psi}_k$$

The second FFN is set to be identity map and we denote the final output as

$$x_L^{(2)} := (x_s)_{s=n}^{L-1} \text{softmax} \left((\hat{g}_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}))_{s=n}^{L-1} \right)^\top.$$

First, we consider the error under the first norm, $\|\cdot\|_\infty$, which can be divided the total error into three components:

$$\begin{aligned} &\left\| x_L^{(2)} - (x_s)_{s=n}^{L-1} \text{softmax} \left((g(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right)^\top \right\|_\infty \\ &\leq \left\| \text{softmax} \left((\hat{g}_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}))_{s=n}^{L-1} \right) - \text{softmax} \left((g(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right) \right\|_\infty \\ &\leq \left\| \text{softmax} \left((\hat{g}_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}))_{s=n}^{L-1} \right) - \text{softmax} \left((g_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}))_{s=n}^{L-1} \right) \right\|_\infty \\ &\quad + \left\| \text{softmax} \left((g_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}))_{s=n}^{L-1} \right) - \text{softmax} \left((g_K(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right) \right\|_\infty \\ &\quad + \left\| \text{softmax} \left((g_K(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right) - \text{softmax} \left((g(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right) \right\|_\infty \\ &\leq \max_s \left| \hat{g}_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}) - g_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}) \right| \\ &\quad + \max_s \left| g_K(\hat{X}_{L-n+2:L}; \hat{X}_{s-n+1:s-1}) - g_K(X_{L-n+2:L}; X_{s-n+1:s-1}) \right| \\ &\quad + \sum_{s=n}^{L-1} \left| g_K(X_{L-n+2:L}; X_{s-n+1:s-1}) - g(X_{L-n+2:L}; X_{s-n+1:s-1}) \right| \end{aligned} \tag{17}$$

For the first term in RHS of (17), it holds that:

$$\begin{aligned}
& \max_s \left| \hat{g}_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) - g_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) \right| \\
& \leq \max_s \sum_{k=1}^K \sigma_k \left| \hat{\phi}_k(\hat{X}_{L-n+2:L}) \hat{\psi}_k(\hat{X}_{s-n+1:s-1}) - \phi_k(\hat{X}_{L-n+2:L}) \psi_k(\hat{X}_{s-n+1:s-1}) \right| \\
& \leq \sum_{k=1}^K \sigma_k \left(\|\hat{\phi}_k\|_{L^\infty} \|\hat{\psi}_k - \psi_k\|_{L^\infty} + \|\psi_k\|_{L^\infty} \|\hat{\phi}_k - \phi_k\|_{L^\infty} \right) \\
& \leq \epsilon_{\text{FFN}}^{(1)} \cdot \sum_{k=1}^K \sigma_k \left(\|\hat{\phi}_k\|_{L^\infty} + \|\psi_k\|_{L^\infty} \right) \\
& \leq \epsilon_{\text{FFN}}^{(1)} \cdot \sum_{k=1}^K \sigma_k \left(\|\phi_k\|_{L^\infty} + \|\hat{\phi}_k - \phi_k\|_{L^\infty} + \|\psi_k\|_{L^\infty} \right) \\
& \leq \epsilon_{\text{FFN}}^{(1)} \cdot (2C_g^\infty + 1) \sum_{k=1}^K \sigma_k \leq (2C_g^\infty + 1) C_\alpha \epsilon_{\text{FFN}}^{(1)}.
\end{aligned}$$

For the second term in RHS of (17), we have:

$$\begin{aligned}
& \max_s \left| g_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) - g_K(X_{L-n+2:L}, X_{s-n+1:s-1}) \right| \\
& \leq \max_s \sum_{k=1}^K \sigma_k \left(\|\phi_k\|_{L^\infty} |\psi_k(\hat{X}_{s-n+1:s-1}) - \hat{\psi}_k(X_{s-n+1:s-1})| \right. \\
& \quad \left. + \|\psi_k\|_{L^\infty} |\phi_k(\hat{X}_{L-n+1:L-1}) - \phi_k(X_{L-n+1:L-1})| \right) \\
& \leq \max_s \sum_{k=1}^K \sigma_k \left(\|\phi_k\|_{L^\infty} \|\psi_k\|_{\text{Lip}} \|\hat{X}_{s-n+1:s-1} - X_{s-n+1:s-1}\| \right. \\
& \quad \left. + \|\psi_k\|_{L^\infty} \|\phi_k\|_{\text{Lip}} \|\hat{X}_{L-n+1:L-1} - X_{L-n+1:L-1}\| \right) \\
& \leq 2C_g^\infty C_g^{\text{Lip}} \epsilon_{\text{SA}}^{(1)} \cdot \left(\max_s \sum_{k=1}^K \sigma_k \right) \leq 2C_g^\infty C_g^{\text{Lip}} C_\alpha \epsilon_{\text{SA}}^{(1)}.
\end{aligned}$$

Additionally, the third term in RHS of (17), its L^2 holds that:

$$\begin{aligned}
& \int_{[0,1]^{d \times L}} \left(\sum_{s=n}^{L-1} |g_K(X_{L-n+2:L}, X_{s-n+1:s-1}) - g(X_{L-n+2:L}, X_{s-n+1:s-1})| \right)^2 dX \\
& \leq (L-1-n) \sum_{s=n}^{L-1} \int_{[0,1]^{D \times L}} |g_K(X_{-n+2:t}, X_{s-n+1:s-1}) - g(X_{L-n+2:L}, X_{s-n+1:s-1})|^2 dX \\
& = (L-1-n)^2 \int_{[0,1]^D \times [0,1]^D} |g(u, v) - g_K(u, v)|^2 du dv \\
& = (L-1-n)^2 \int \left(\sum_{k=K+1}^{+\infty} \sigma_k \phi_k(u) \psi_k(v) \right)^2 du dv \\
& \leq \int \left(\sum_{k=K+1}^{+\infty} \sigma_k \phi_k^2(u) \right) \left(\sum_{k=K+1}^{+\infty} \sigma_k \psi_k^2(v) \right) du dv \\
& \leq (L-1-n)^2 \left(\sum_{k=K+1}^{\infty} \sigma_k \right)^2 \leq \frac{(L-1-n)^2 C_\alpha^2}{K^{2\alpha}}.
\end{aligned}$$

Now we combine three error terms together to obtain the total L^2 error for the output of this layer:

$$\begin{aligned}
& \int_{X \in [0,1]^{d \times L}} \left\| x_L^{(2)} - (x_s)_{s=n}^{L-1} \text{softmax} \left((g(X_{L-n+2:L}; X_{s-n+1:s-1}))_{s=n}^{L-1} \right)^\top \right\|_\infty^2 dX \\
& \leq 3 \int_{X \in [0,1]^{d \times L}} \max_s \left| \hat{g}_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) - g_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) \right|^2 dX \\
& \quad + 3 \int_{X \in [0,1]^{d \times L}} \max_s \left| g_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) - g_K(X_{L-n+2:L}, X_{s-n+1:s-1}) \right|^2 dX \\
& \quad + 3 \int_{X \in [0,1]^{d \times L}} \left(\sum_{s=n}^{L-1} |g_K(X_{L-n+2:L}, X_{s-n+1:s-1}) - g(X_{L-n+2:L}, X_{s-n+1:s-1})| \right)^2 dX \\
& \leq 3 \max_s \left| \hat{g}_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) - g_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) \right|^2 \\
& \quad + 3 \max_s \left| g_K(\hat{X}_{L-n+2:L}, \hat{X}_{s-n+1:s-1}) - g_K(X_{L-n+2:L}, X_{s-n+1:s-1}) \right|^2 \\
& \quad + 3 \left(\frac{(L-1-n)C_\alpha}{K^\alpha} \right)^2 \\
& \leq 3 \left((2C_g^\infty + 1)C_\alpha \epsilon_{\text{FFN}}^{(1)} \right)^2 + 3 \left(2C_g^\infty C_g^{\text{Lip}} C_\alpha \epsilon_{\text{SA}}^{(1)} \right)^2 + 3 \left(\frac{(L-1-n)C_\alpha}{K^\alpha} \right)^2 \\
& \leq 3 \left((2C_g^\infty + 1)C_\alpha \epsilon_{\text{FFN}}^{(1)} + 2C_g^\infty C_g^{\text{Lip}} C_\alpha \epsilon_{\text{SA}}^{(1)} + \frac{(L-1-n)C_\alpha}{K^\alpha} \right)^2.
\end{aligned}$$

This estimate implies that

$$\begin{aligned}
& \|\text{GIH}_n - \text{TF}\|_{L,2} \\
& \leq \sqrt{3} \left(2C_g^\infty C_g^{\text{Lip}} C_\alpha \epsilon_{\text{SA}}^{(1)} + (2C_g^\infty + 1)C_\alpha \epsilon_{\text{FFN}}^{(1)} + \frac{(L-1-n)C_\alpha}{K^\alpha} \right) \\
& \leq \mathcal{O} \left(\frac{C_{g,q,n}}{H^n} \right) + \tilde{\mathcal{O}} \left(\frac{C_{g,\alpha} \sqrt{K}}{\sqrt{M}} \right) + \mathcal{O} \left(\frac{tC_\alpha}{K^\alpha} \right)
\end{aligned} \tag{18}$$

Step IV. Optimizing K in (18).

Notice that in RHS of (18), only $\tilde{\mathcal{O}} \left(\frac{C_{g,\alpha} \sqrt{K}}{\sqrt{M}} \right)$ and $\mathcal{O} \left(\frac{LC_\alpha}{K^\alpha} \right)$ depend on K .

By Young's inequality, with $p = \frac{\alpha + \frac{1}{2}}{\alpha}$ and $q = 2(\alpha + \frac{1}{2})$, we have:

$$\begin{aligned}
& \min_K : \frac{\alpha}{\frac{1}{2} + \alpha} \frac{C_{g,\alpha} \sqrt{K}}{\sqrt{M}} + \frac{\frac{1}{2}}{\frac{1}{2} + \alpha} \frac{LC_\alpha}{K^\alpha} \\
& = \min_K : \frac{\alpha}{\frac{1}{2} + \alpha} \left(\left(\frac{C_{g,\alpha} \sqrt{K}}{\sqrt{M}} \right)^{\frac{\alpha}{\frac{1}{2} + \alpha}} \right)^{\frac{\frac{1}{2} + \alpha}{\alpha}} + \frac{\frac{1}{2}}{\frac{1}{2} + \alpha} \left(\left(\frac{LC_\alpha}{K^\alpha} \right)^{\frac{1}{\frac{1}{2} + \alpha}} \right)^{2(\frac{1}{2} + \alpha)} \\
& = \frac{C'_{g,\alpha} L^{1/(1+2\alpha)}}{M^{\alpha/(1+2\alpha)}}.
\end{aligned}$$

Thus, we obtain our final bound:

$$\|\text{GIH}_n - \text{TF}\|_{L,2} \leq \mathcal{O} \left(\frac{C_{g,q,n}}{H^n} \right) + \left\{ \tilde{\mathcal{O}} \left(\frac{C_{g,\alpha} \sqrt{K}}{\sqrt{M}} \right) + \mathcal{O} \left(\frac{LC_\alpha}{K^\alpha} \right) \right\}_{\min:K}$$

$$\leq \mathcal{O}\left(\frac{C_{g,q,n}}{H^n}\right) + \tilde{\mathcal{O}}\left(\frac{C'_{g,\alpha} L^{1/(1+2\alpha)}}{M^{\alpha/(1+2\alpha)}}\right) \leq \frac{A_{g,q,n}}{H^n} + \frac{B_{g,\alpha} L^{1/(1+2\alpha)}}{M^{\alpha/(1+3\alpha)}}.$$

□

B PROOFS OF OPTIMIZATION DYNAMICS: TRAINING STAGE I

In this subsection we focus on training the first layer of Transformer model to capture the token ahead. For simplicity, we introduce some notations:

$$\tilde{p} := p^{(1,1)}, \quad p := p^{(2,1)}, \quad g := w_V^{(2,1)}, \quad h := w_V^{(2,2)}, \quad w_K := w_K^{(2,2)}, \quad w_Q := w_Q^{(2,2)},$$

and denote the initialization of each parameter as $\tilde{p}(0), p(0), g(0), w_Q(0), w_K(0), h(0)$ respectively.

We initialize $p(0), w_k(0), w_Q(0) = 0$ while the other parameters are all initialized at σ_{init} . In this training stage, we only train \tilde{p} . And our goal, **the proof of Lemma 5.3** can be deduced from which, is to prove:

$$\lim_{t \rightarrow +\infty} \tilde{p}(t) = +\infty.$$

In this stage, the s -th output token of the first layer is represented as

$$\begin{pmatrix} x_s \\ (x_\tau)_{\tau=1}^{s-1} \text{softmax} \left((-\tilde{p}(s-1-\tau))_{\tau=1}^{s-1} \right)^\top \end{pmatrix},$$

and the target function and output of transformer are as follows

$$f^*(X) = \begin{pmatrix} \frac{\alpha^*}{1+\alpha^*} x_{L-2} \\ \frac{1}{1+\alpha^*} (x_s)_{s=2}^{L-1} \text{softmax} \left((x_L w^{*2} x_{s-1})_{s=2}^{L-1} \right)^\top \end{pmatrix},$$

$$\begin{aligned} f_\theta(X) &= \begin{pmatrix} g(0) \left(\sum_{\tau=1}^{s-1} \text{softmax}_s(-\tilde{p}(s-1-\tau)) x_\tau \right)_{s=2}^{L-1} \text{softmax} \left(-p(0)(L-1-s)_{s=2}^{L-1} \right) \\ h(0) (x_s)_{s=2}^{L-2} \cdot \text{softmax} \left(\left(w_K(0) w_Q(0) x_L \cdot (x_\tau)_{\tau=1}^{s-1} \text{softmax} \left((-\tilde{p}(s-1-\tau))_{\tau=1}^{s-1} \right)^\top \right)_{s=2}^{L-2} \right) \end{pmatrix} \\ &= \begin{pmatrix} g(0) \frac{1}{L-2} \sum_{\tau=1}^{L-2} \left(\sum_{s=\tau+1}^{L-1} \text{softmax} \left(-\tilde{p}(s-1-t)_{t=1}^{s-1} \right)_{t=\tau} \right) x_\tau \\ h(0) \frac{1}{L-2} \sum_{s=2}^{L-2} x_s \end{pmatrix}. \end{aligned}$$

Since we only focus on \tilde{p} and the other parameters remain the initialization value, the loss function can be simplified as

$$\begin{aligned} \mathcal{L}(\theta) &= \mathbb{E}_{X \sim \mathbb{N}(0,1)^L} \left[\frac{\alpha^{*2}}{(1+\alpha^*)^2} x_{L-2}^2 + \frac{g(0)^2}{(L-2)^2} \sum_{\tau=1}^{L-2} \left(\sum_{s=\tau+1}^{L-1} \text{softmax} \left(-\tilde{p}(s-1-t)_{t=1}^{s-1} \right)_{t=\tau} \right)^2 x_\tau^2 \right. \\ &\quad \left. + \frac{2g(0)}{L-2} \frac{\alpha^*}{1+\alpha^*} \text{softmax} \left(-p(0)(L-1-s)_{s=2}^{L-1} \right)_{s=L-1} x_{L-2}^2 \right] + C(w^*, \alpha^*, w(0), h(0)) \end{aligned}$$

where the second term $C(w^*, \alpha^*, w(0), h(0))$ is a constant depends on $w^*, \alpha^*, w(0)$ and $h(0)$, produced by calculating the error of the second head, i.e., loss of induction head, while the first term is 4-gram loss.

We first define several functions that will be useful for calculation in this stage and the second one:

1296 *Function I.* This function is purely defined for the calculation of $\frac{dp}{dt}$. Denoted by $q(\tilde{p}) :=$
 1297 $\sum_{\tau=1}^{L-2} \left(\sum_{s=\tau+1}^{L-1} \frac{e^{\tilde{p}(s-1-\tau)}}{\sum_{k=0}^{s-2} e^{-\tilde{p}k}} \right)^2$, we first prove $q'(\tilde{p}) \leq 0$.

$$\begin{aligned}
 1300 \quad q(\tilde{p}) &:= \sum_{\tau=1}^{L-2} \left(\sum_{s=\tau+1}^{L-1} \frac{e^{\tilde{p}(s-1-\tau)}}{\sum_{k=0}^{s-2} e^{-\tilde{p}k}} \right)^2 \\
 1301 & \\
 1302 & \\
 1303 &= \sum_{\tau=1}^{L-2} \left(\sum_{s=\tau+1}^{L-1} \frac{e^{-\tilde{p}(s-1-\tau)}}{1 - e^{-\tilde{p}(s-1)}} (1 - e^{-\tilde{p}}) \right)^2 \\
 1304 & \\
 1305 &= (1 - e^{-\tilde{p}})^2 \sum_{\tau=1}^{L-2} \left(\sum_{s=\tau+1}^{L-1} \frac{e^{-\tilde{p}(s-1-\tau)}}{1 - e^{-\tilde{p}(s-1)}} \right)^2 \\
 1306 & \\
 1307 &= (1 - e^{-\tilde{p}})^2 \sum_{\tau=1}^{L-2} e^{2\tilde{p}\tau} \left(\sum_{s=\tau+1}^{L-1} \frac{e^{-\tilde{p}(s-1)}}{1 - e^{-\tilde{p}(s-1)}} \right)^2 \\
 1308 & \\
 1309 &= (1 - e^{-\tilde{p}})^2 \sum_{\tau=1}^{L-2} e^{2\tilde{p}\tau} \left(\sum_{s=\tau+1}^{L-1} \frac{1}{e^{\tilde{p}(s-1)} - 1} \right)^2 \\
 1310 & \\
 1311 & \\
 1312 & \\
 1313 & \\
 1314 & \\
 1315 &
 \end{aligned}$$

1316 Then we take its derivative of \tilde{p}

$$\begin{aligned}
 1317 \quad q'(\tilde{p}) &= 2(1 - e^{-\tilde{p}})e^{-\tilde{p}} \sum_{\tau=1}^{L-2} e^{2\tilde{p}\tau} \left(\sum_{s=\tau+1}^{L-1} \frac{1}{e^{\tilde{p}(s-1)} - 1} \right)^2 \\
 1318 & \\
 1319 &+ (1 - e^{-\tilde{p}})^2 \sum_{\tau=1}^{L-2} 2\tau e^{2\tilde{p}\tau} \left(\sum_{s=\tau+1}^{L-1} \frac{1}{e^{\tilde{p}(s-1)} - 1} \right)^2 \\
 1320 & \\
 1321 &+ (1 - e^{-\tilde{p}})^2 \sum_{\tau=1}^{L-2} 2e^{2\tilde{p}\tau} \left(\sum_{s=\tau+1}^{L-1} \frac{1}{e^{\tilde{p}(s-1)} - 1} \right) \left(\sum_{s=\tau+1}^{L-1} \frac{-(s-1)e^{\tilde{p}(s-1)}}{(e^{\tilde{p}(s-1)} - 1)^2} \right) \\
 1322 & \\
 1323 &= 2(1 - e^{-\tilde{p}}) \sum_{\tau=1}^{L-2} e^{2\tilde{p}\tau} \left(\sum_{s=\tau+1}^{L-1} \frac{1}{e^{\tilde{p}(s-1)} - 1} \right) \left(\sum_{s=\tau+1}^{L-1} \frac{e^{-\tilde{p}} + \tau(1 - e^{-\tilde{p}})}{e^{\tilde{p}(s-1)} - 1} - \frac{(s-1)e^{\tilde{p}(s-1)}}{(e^{\tilde{p}(s-1)} - 1)^2} \right) \\
 1324 & \\
 1325 & \\
 1326 & \\
 1327 & \\
 1328 & \\
 1329 &
 \end{aligned}$$

1330 $q'(\tilde{p})$'s last factor can be formed as

$$\begin{aligned}
 1331 & \\
 1332 & \frac{(\tau - (\tau - 1)e^{-\tilde{p}}) (e^{\tilde{p}(s-1)} - 1) - (s - 1)e^{\tilde{p}(s-1)}}{e^{\tilde{p}(s-1)} - 1)^2} \\
 1333 & \\
 1334 &= \frac{(\tau + 1 - s)t^{s-1} - (\tau - 1)t^{s-2} - \tau + \frac{\tau-1}{t}}{e^{\tilde{p}(s-1)} - 1)^2} \\
 1335 & \\
 1336 &
 \end{aligned}$$

1337 where $t = e^{-\tilde{p}} \geq 1$. Since $s \geq \tau + 1$, $q'(\tilde{p}) \leq 0$.

1338 *Function II.* For simplicity, we define $M(p)$ and its derivative $m(p)$:

$$\begin{aligned}
 1340 \quad M(p) &:= \sum_{s=2}^{L-1} \exp(-p(L-1-s)) = \sum_{s=0}^{L-3} \exp(-ps) = \frac{1 - e^{-p(L-2)}}{1 - e^{-p}}, \\
 1341 & \\
 1342 & \\
 1343 & \\
 1344 \quad m(p) &:= \sum_{s=1}^{L-3} s \exp(-ps) = \frac{e^{-p} - (L-2)e^{-p(L-2)} + (L-3)e^{-p(L-1)}}{(1 - e^{-p})^2}. \\
 1345 & \\
 1346 &
 \end{aligned}$$

1347 *Function III.* The third function is derivative of softmax. By straightforward calculation, we obtain:

$$1348 \quad \frac{d}{dp} \text{softmax} \left(-p(L-1-t)_{t=2}^{L-1} \right)_{t=L-1-s} = \frac{d}{dp} \frac{\exp(-ps)}{\sum_{\tau=0}^{L-3} \exp(-p\tau)} = \frac{-s \exp(-ps)M(p) + \exp(-ps)m(p)}{M(p)^2}.$$

Through the quantities and their properties above, we obtain the dynamic of \tilde{p}

$$\begin{aligned} \frac{d\tilde{p}}{dt} &= -\frac{g(0)^2}{(L-2)^2} q'(\tilde{p}) + \frac{2\alpha^* g(0)}{(1+\alpha^*)(L-2)} \frac{m(p)}{M(p)^2} \\ &\geq \frac{2\alpha^* g(0)}{(1+\alpha^*)(L-2)} e^{-\tilde{p}}, \end{aligned}$$

which implies:

$$\lim_{t \rightarrow +\infty} \tilde{p}(t) = +\infty.$$

C PROOFS OF OPTIMIZATION DYNAMICS: TRAINING STAGE II

In this training stage, the first layer is already capable of capturing the token ahead i.e. $y_s = x_{s-1}$. And we train the parameters $w_{V_1}, w_{V_2}, p, w_{KQ}$ in the second layer.

We start from proving the parameter balance lemma:

Lemma C.1 (Restate of Lemma 5.4). *In Training Stage II, it holds that $w_Q^{(2,2)^2}(t) \equiv w_K^{(2,2)^2}(t)$.*

Proof. Notice that

$$\begin{aligned} \frac{d}{dt} \left(w_Q^{(2,2)^2}(t) - w_K^{(2,2)^2}(t) \right) &= -w_Q^{(2,2)} \frac{\partial \mathcal{L}}{\partial w_Q^{(2,2)}} + w_K^{(2,2)} \frac{\partial \mathcal{L}}{\partial w_K^{(2,2)}} \\ &= -w_Q^{(2,2)} w_K^{(2,2)} \frac{\partial \mathcal{L}}{\partial \left(w_Q^{(2,2)} w_K^{(2,2)} \right)} + w_K^{(2,2)} w_Q^{(2,2)} \frac{\partial \mathcal{L}}{\partial \left(w_Q^{(2,2)} w_K^{(2,2)} \right)} \equiv 0. \end{aligned}$$

Thus, we have:

$$w_Q^{(2,2)^2}(t) - w_K^{(2,2)^2}(t) \equiv w_Q^{(2,2)^2}(0) - w_K^{(2,2)^2}(0) = 0.$$

□

For simplicity, we still use the following notations:

$$p := p_1, \quad g := w_{V_1}, \quad w := w_{KQ}, \quad h := w_{V_2}.$$

and notations for initialization $p(0), g(0), w(0), h(0)$. Then the target function and output of Transformer can be formed as follows

$$\begin{aligned} f^*(X) &= \left(\frac{\frac{\alpha^*}{1+\alpha^*} x_{L-2}}{\frac{1}{1+\alpha^*} (x_s)_{s=2}^{L-1} \cdot \text{softmax} \left((w^{*2} x_L x_{s-1})_{s=2}^{L-1} \right)} \right), \\ \text{TF}(X; \theta) &= \left(\begin{array}{c} g \cdot (x_{s-1})_{s=2}^{L-2} \cdot \text{softmax} \left((-p(L-1-s))_{s=2}^{L-2} \right) \\ h \cdot (x_s)_{s=2}^{L-2} \cdot \text{softmax} \left((w^2 x_L x_{s-1})_{s=2}^{L-2} \right) \end{array} \right). \end{aligned}$$

And the loss function is expressed as:

$$\begin{aligned} \mathcal{L}(\theta) &= \frac{1}{2} \mathbb{E}_{X \sim \mathcal{N}(0,1)^L} [\|f^*(x) - \text{TF}(x; \theta)\|^2] \\ &= \frac{1}{2} \mathbb{E}_X \left[\left(\frac{\alpha^*}{1+\alpha^*} x_{L-2} - g \cdot (x_{s-1})_{s=2}^{L-2} \cdot \text{softmax} \left((-p(L-1-s))_{s=2}^{L-2} \right) \right)^2 \right] \\ &\quad + \frac{1}{2} \mathbb{E}_X \left[\left(\frac{1}{1+\alpha^*} (x_s)_{s=2}^{L-1} \cdot \text{softmax} \left((w^{*2} x_L x_{s-1})_{s=2}^{L-1} \right) - h \cdot (x_s)_{s=2}^{L-2} \cdot \text{softmax} \left((w^2 x_L x_{s-1})_{s=2}^{L-2} \right) \right)^2 \right]. \end{aligned}$$

The total loss can naturally be divided into two parts:

$$\mathcal{L}(\theta) = \mathcal{L}_{\mathbb{G}_4}(\theta) + \mathcal{L}_{\mathbb{H}_2}(\theta),$$

where

$$\begin{aligned} \mathcal{L}_{\mathbb{G}_4}(\theta) &= \mathcal{L}_{\mathbb{G}_4}(p, g) \\ &= \frac{1}{2} \mathbb{E}_X \left[\left(\frac{\alpha^*}{1 + \alpha^*} x_{L-2} - g \cdot (x_{s-1})_{s=2}^{L-2} \cdot \text{softmax} \left((-p(L-1-s))_{s=2}^{L-2} \right) \right)^2 \right], \end{aligned}$$

$$\begin{aligned} \mathcal{L}_{\mathbb{H}_2}(\theta) &= \mathcal{L}_{\mathbb{H}_2}(w, h) \\ &= \frac{1}{2} \mathbb{E}_X \left[\left(\frac{1}{1 + \alpha^*} (x_s)_{s=2}^{L-1} \cdot \text{softmax} \left((w^* x_L x_{s-1})_{s=2}^{L-1} \right) - h \cdot (x_s)_{s=2}^{L-2} \cdot \text{softmax} \left((w^* x_L x_{s-1})_{s=2}^{L-2} \right) \right)^2 \right]. \end{aligned}$$

Notably, the dynamics of (p, g) and (w, h) are **decoupled**, which allows us to analyze them separately.

Additionally, we denote the optimal values of the parameters as:

$$p^* = +\infty, \quad g^* = \frac{\alpha^*}{1 + \alpha^*}, \quad w^* := w^*, \quad h^* = \frac{1}{1 + \alpha^*}.$$

For the initialization scale and the sequence length, we consider the case:

$$\sigma_{\text{init}} = \mathcal{O}(1) \ll 1, \quad L = \Omega(1/\sigma_{\text{init}}) \gg 1.$$

C.1 DYNAMICS OF THE PARAMETERS FOR 4-GRAM

First, we define two useful auxiliary functions:

$$\begin{aligned} M(p) &:= \frac{1 - e^{-p(L-2)}}{1 - e^{-p}}, \\ m(p) &:= \frac{e^{-p} - (L-2)e^{-p(L-2)} + (L-3)e^{-p(L-1)}}{(1 - e^{-p})^2}. \end{aligned}$$

Then, a straightforward calculation, combined with Lemma D.3 and Lemma D.4, yields the explicit formulation of $\mathcal{L}_{\mathbb{G}_4}(\theta)$ and the GF dynamics of p and g :

$$\mathcal{L}_{\mathbb{G}_4}(\theta) = \frac{1}{2} \left(\frac{\alpha^*}{1 + \alpha^*} \right)^2 + \frac{1}{2} g^2 \frac{M(2p)}{M(p)^2} - \frac{\alpha^* g}{1 + \alpha^*} \frac{1}{M(p)}. \quad (19)$$

$$\begin{aligned} \frac{dp}{dt} &= -\frac{\partial \mathcal{L}}{\partial p} = -\frac{\partial \mathcal{L}_{\mathbb{G}_4}}{\partial p} = \frac{m(p)}{M(p)^2} \left[g^2 \frac{m(2p)}{m(p)} - g^2 \frac{M(2p)}{M(p)} + \frac{\alpha^* g}{1 + \alpha^*} \right], \\ \frac{dg}{dt} &= -\frac{\partial \mathcal{L}}{\partial g} = -\frac{\partial \mathcal{L}_{\mathbb{G}_4}}{\partial g} = \frac{\alpha^*}{1 + \alpha^*} \frac{1}{M(p)} - g \frac{M(2p)}{M(p)^2}, \end{aligned}$$

Equivalently, the dynamics can be written as:

$$\begin{aligned} \frac{dp}{dt} &= \frac{m(p)g}{M(p)^2} \left(g^* - g \frac{M(2p)}{M(p)} + g \frac{m(2p)}{m(p)} \right), \\ \frac{dg}{dt} &= \frac{1}{M(p)} \left(g^* - g \frac{M(2p)}{M(p)} \right). \end{aligned}$$

Notice that at the initialization, it holds that $\frac{dp}{dt}|_{t=0} > 0$ and $\frac{dg}{dt}|_{t=0} > 0$. Then we first define a hitting time:

$$T_1^g := \inf\{t > 0 : g(t) > g^*\}.$$

1458 Noticing $g(0) = \sigma_{\text{init}} \ll g^*$ and the continuity, $T_1^g > 0$.

1459 Our subsequent proof can be divided into **two phases**: a monotonic phase $t < T_1^g$, and a stable
1460 convergence phase $t > T_1^g$.

1462 **Part I. Analysis for the monotonic phase $t < T_1^g$.**

1463

1464

1465
$$\frac{dp}{dt} = \frac{m(p)g}{M(p)^2} \left(g^* - g \frac{M(2p)}{M(p)} + g \frac{m(2p)}{m(p)} \right) = \frac{m(p)g}{M(p)^2} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} + g \frac{m(2p)}{m(p)} \right),$$

1466

1467

1468

1469

$$\frac{dg}{dt} = \frac{1}{M(p)} \left(g^* - g \frac{M(2p)}{M(p)} \right) = \frac{1}{M(p)} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} \right).$$

1470 It is easy to see that p, g are monotonically increasing for $t < T_1^g$. We can choose sufficiently large

1471

1472

$$L = \Omega(1/p(0)) = \Omega(1/\sigma_{\text{init}})$$

1473

1474 such that:

1475

1476

$$(L-3)e^{-(L-3)p(t)}, e^{-(L-5)p(t)} < 0.0001, \quad \forall p > \sigma_{\text{init}}.$$

1477

1478 Then we can calculate the following three terms in the dynamics:

1479

1480

1481

1482

1483

$$\frac{m(p)}{M^2(p)} = \frac{e^{-p} (1 - (L-2)e^{-p(L-3)} + (L-3)e^{-p(L-2)})}{1 - e^{-p(L-2)}} = \frac{e^{-p}(1 + \xi_1(p))}{1 + \xi_2(p)},$$

$$\frac{1}{M(p)} = \frac{1 - e^{-p(L-2)}}{1 - e^{-p}} = \frac{1 + \xi_3(p)}{1 - e^{-p}},$$

1484

1485

1486

1487

1488

1489

$$\begin{aligned} \frac{m(2p)}{m(p)} &= \frac{e^{-p} (1 - (L-2)e^{-2p(L-3)} + (L-3)e^{-2p(L-2)})}{(1 + e^{-p})^2 (1 - (L-2)e^{-p(L-3)} + (L-3)e^{-p(L-2)})} \\ &= \frac{e^{-p}(1 + \xi_4(p))}{(1 + e^{-p})^2(1 + \xi_5(p))}, \end{aligned}$$

1490

1491 where the error functions satisfy:

1492

1493

1494

$$|\xi_1(p)|, \dots, |\xi_5(p)| \leq 0.0001, \quad \forall t > T_1^g.$$

1495

1496 Then the dynamics satisfy:

1497

1498

1499

$$\frac{dp}{dt} = \frac{e^{-p}g(1 + \xi_1(p))}{1 + \xi_2(p)} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} + \frac{ge^{-p}(1 + \xi_3(t))}{(1 + e^{-p})^2(1 + \xi_5(t))} \right),$$

$$\frac{dg}{dt} = \frac{1 + \xi_3(p)}{1 - e^{-p}} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} \right).$$

1500

1501

1502 When $g < \frac{1}{2} \frac{\alpha^*}{1 + \alpha^*}$, we have

1503

1504

$$\frac{dp}{dg} \leq 2(e^{-p} - e^{-2p})g.$$

1505

1506 By define $T_{1/2}^g := \inf\{t > 0 : g(t) > g^*/2\}$ and $\tilde{p} := p(T_{1/2}^g)$, we have

1507

1508

1509

$$\ln(e^{\tilde{p}} - 1) \leq \frac{1}{4}g^{*2} - g(0)^2 + e^{p(0)} - 1 + \ln(e^{p(0)} - 1)$$

1510

1511 then $\tilde{p} \leq \mathcal{O}(\sqrt{p(0)})$, from which we infer that p barely increases when $t \leq T_{1/2}^g$.

1512

For $0 \leq t \leq T_{1/2}^g$,

1513

$$\frac{dg}{dt} \geq \frac{1}{1 - e^{-p(0)}} \left[g^* - \frac{g}{1 + e^{-p(0)}} \right]$$

$$g \geq g^*(1 + e^{-p(0)}) + \left[g(0) - g^*(1 + e^{-p(0)}) \right] \exp\left(\frac{-t}{1 - e^{-2p(0)}}\right)$$

so

$$T_{1/2}^g \leq (1 - e^{-2p(0)}) \ln\left(\frac{g^*(1 + e^{-p(0)}) - g(0)}{g^* \left((1 + e^{-p(0)}) - \frac{1}{2} \right)}\right) = \mathcal{O}(2p(0))$$

For $T_{1/2}^g \leq t \leq T_1^g$, let $p_1 := p(T_1^g)$,

$$\begin{aligned} \frac{dp}{dg} &\leq 1.01e^{-p}(1 - e^{-p})g \left(1 + \frac{\frac{g}{1+e^{-p}} - \frac{g}{(1+e^{-p})^2}}{\frac{\alpha^*}{1+\alpha^*} - \frac{g}{1+e^{-p}}} \right) \\ &\leq \frac{1.01}{4} \frac{\alpha^*}{1+\alpha^*} (1 + e^{-p_1}) \end{aligned}$$

then

$$\begin{aligned} p_1 - p(0) &\leq \frac{1.01}{4} \left(\frac{\alpha^*}{1+\alpha^*} \right)^2 (1 + e^{p_1}), \\ p_1 &\leq \frac{1}{2 \left(\frac{\alpha^*}{1+\alpha^*} \right)^2 - 1}, \end{aligned}$$

and we take $\alpha^* > 1$.

Since for $T_{1/2}^g \leq t \leq T_1^g$,

$$\begin{aligned} \frac{dp}{dt} &\leq 2e^{-p}g^* \left(g^* - \frac{1}{8}g^* \right), \\ \frac{dp}{dt} &\geq \frac{1}{2}e^{-p}g^* \left(g^* - \frac{1}{1+e^{-p_1}}g^* \right), \end{aligned}$$

we have

$$T_1^g - t_1 \leq \mathcal{O}\left((e^{2p_1} - 1) \left(\frac{1 + \alpha^*}{\alpha^*} \right)^2 \right).$$

Hence, putting the two part of time together we have

$$\begin{aligned} T_1^g &\leq \mathcal{O}\left(p(0) + (e^{2p_1} - 1) \left(\frac{1 + \alpha^*}{\alpha^*} \right)^2 \right) \\ &= \mathcal{O}\left(\sigma_{\text{init}} + (e^{2p_1} - 1) \left(\frac{1 + \alpha^*}{\alpha^*} \right)^2 \right) = \mathcal{O}(1). \end{aligned} \tag{20}$$

Part II. Analysis for the convergence phase $t > T_1^g$.

We will prove that, in this phase, (p, g) keep in a stable region, and the convergence occurs.

Recall the dynamics:

$$\begin{aligned} \frac{dp}{dt} &= \frac{m(p)g}{M(p)^2} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} + g \frac{m(2p)}{m(p)} \right), \\ \frac{dg}{dt} &= \frac{1}{M(p)} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} \right). \end{aligned}$$

Using contradiction, it is easy to verify that for all $t > T_1^g$,

$$g^* < g(t) < 2g^*, \quad \frac{dp(t)}{dt} > 0,$$

1566 which means g has entered a stable region (although it is possible that g is non-monotonic), while p
 1567 keeps increase. In fact, if $\hat{t} := \inf\{t > 0 : g(t) = 2g^*\}$, then $\frac{dg}{dt}|_{\hat{t}} < 0$, which leads to a contradiction.
 1568 If $\hat{t} := \inf\{t > 0 : \frac{dp(t)}{dt} = 0\}$, then
 1569

$$1570 \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} + g \frac{m(2p)}{m(p)} \right) \Big|_{\hat{t}} = 0, \quad \frac{dg}{dt} < 0,$$

$$1571 \frac{d}{dt} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} + g \frac{m(2p)}{m(p)} \right) = -g' \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} + g' \frac{m(2p)}{m(p)} > 0,$$

1572 where the last inequality leads to a contradiction.

1573 Thus, $p(t) > p(T_1^g) > p(0) = \sigma_{\text{init}}$ holds in this phase. Therefore, the dynamics
 1574

$$1575 \frac{dp}{dt} = \frac{e^{-p}g(1 + \xi_1(p))}{1 + \xi_2(p)} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} + \frac{ge^{-p}(1 + \xi_3(t))}{(1 + e^{-p})^2(1 + \xi_5(t))} \right),$$

$$1576 \frac{dg}{dt} = \frac{1 + \xi_3(p)}{1 - e^{-p}} \left(g^* - g \frac{1 + e^{-p(L-2)}}{1 + e^{-p}} \right),$$

1577 also satisfy

$$1578 |\xi_1(p)|, \dots, |\xi_5(p)| \leq 0.0001, \quad \forall t > T_1^g.$$

1580 For simplicity, we consider the transform:

$$1581 u := e^{-p}.$$

1582 Then the dynamics of u and g can be written as:

$$1583 \frac{du}{dt} = -\frac{(1 + \xi_1(p))u^2g}{1 + \xi_2(p)} \left(g^* - g \frac{1 + u^{L-2}}{1 + u} + \frac{gu(1 + \xi_4(p))}{(1 + u)^2(1 + \xi_5(p))} \right),$$

$$1584 \frac{dg}{dt} = \frac{1 + \xi_3(p)}{1 - u} \left(g^* - g \frac{1 + u^{L-2}}{1 + u} \right).$$

1585 Notice that this dynamics are controlled by high-order terms. Consequently, we construct a variable
 1586 to reflect the dynamics of high-order term:

$$1587 v := ug^* + (g^* - g).$$

1588 Then the dynamics of u and v satisfy:

$$1589 \frac{du}{dt} = -\frac{(1 + \xi_1(p))u^2g}{1 + \xi_2(p)} \left(\frac{v - u^{L-2}g}{1 + u} + \frac{gu(1 + \xi_4(p))}{(1 + u)^2(1 + \xi_5(p))} \right),$$

$$1590 \frac{dv}{dt} = -\frac{(1 + \xi_1(p))u^2gg^*}{1 + \xi_2(p)} \left(\frac{v - u^{L-2}g}{1 + u} + \frac{gu(1 + \xi_4(p))}{(1 + u)^2(1 + \xi_5(p))} \right) - \frac{1 + \xi_3(p)}{1 - u^2} (v - u^{L-2}g).$$

1591 Now we consider the Lyapunov function about u, v :

$$1592 G(u, v) := \frac{1}{2} (u^2 + v^2).$$

1593 Then it is straightforward:

$$1594 \frac{dG}{dt} = u \frac{du}{dt} + v \frac{dv}{dt}$$

$$1595 = -\frac{u^3g(1 + \xi_1(p))}{1 + \xi_2(p)} \left(\frac{v - u^{L-2}g}{1 + u} + \frac{gu(1 + \xi_4(p))}{(1 + u)^2(1 + \xi_5(p))} \right)$$

$$1596 - \frac{(1 + \xi_1(p))u^2vgg^*}{1 + \xi_2(p)} \left(\frac{v - u^{L-2}g}{1 + u} + \frac{gu(1 + \xi_4(p))}{(1 + u)^2(1 + \xi_5(p))} \right)$$

1620

1621

1622

$$- \frac{1 + \xi_3(p)}{1 - u^2} (v - u^{L-2}g) v.$$

1623

By $|\xi_1|, \dots, |\xi_5| \leq 0.0001$, we have the following estimate for the Lyapunov dynamics:

1624

1625

1626

1627

1628

1629

1630

1631

1632

By $u^{L-5} = e^{-p(L-5)} < 0.0001$ and $0 < u < e^{-p(T_1^g)}$, we further have:

1633

1634

1635

1636

1637

1638

$$\begin{aligned} \frac{dG}{dt} &\leq \frac{1.001g}{1+u} |u^3v| + \frac{1.0001g^2}{1+u} u^{L+1} - \frac{0.999g^2}{(1+u^2)} u^4 \\ &\quad - \frac{0.999gg^*}{1+u} u^2v^2 + \frac{1.001g^2g^*}{1+u} |u^L v| + \frac{1.001g^2g^*}{(1+u^2)} |u^3v| \\ &\quad - \frac{0.999}{1-u^2} v^2 + \frac{1.001g}{1-u^2} |u^{L-2}v| \end{aligned}$$

1639

By using the following inequalities:

1640

1641

1642

1643

1644

1645

1646

1647

we have

1648

1649

1650

1651

1652

1653

Since $g^* < g < 2g^*$, $u > 0$ for $t > T_1^g$, and $\frac{u^2}{1+u} \leq \frac{1}{2}$ for $0 \leq u \leq 1$, we have:

1654

1655

1656

1657

1658

1659

1660

then

1661

1662

1663

1664

1665

1666

1667

which implies:

1668

1669

1670

1671

1672

1673

$$G(u(t), v(t)) \leq \frac{1}{G(u(t_1), v(t_1)) + \frac{g^{*2}}{64}(t - t_1)}, \quad \forall t > T_1^g.$$

Hence,

$$u^2(t), \quad v^2(t) = \mathcal{O}\left(\frac{1}{g^{*2}t}\right) = \mathcal{O}\left(\frac{1}{t}\right), \quad \forall t > T_1^g = \mathcal{O}(1)$$

1674 which implies:

$$1675 e^{-p(t)} = u(t) = \mathcal{O}\left(\frac{1}{\sqrt{t}}\right), \quad \forall t > T_1^g = \mathcal{O}(1);$$

$$1676$$

$$1677 g(t) - g^* = g^*u(t) - v(t) \leq \mathcal{O}\left(\frac{g^*}{\sqrt{t}}\right) + \mathcal{O}\left(\frac{1}{\sqrt{t}}\right) = \mathcal{O}\left(\frac{1}{\sqrt{t}}\right), \quad \forall t > T_1^g = \mathcal{O}(1). \quad (21)$$

$$1678$$

$$1679$$

$$1680$$

1681 **Notably**, these proofs capture the **entire** training dynamics of p, g , from $t = 0$ to $t = T_1^g$, and finally
 1682 to $t \rightarrow +\infty$, providing a fine-gained analysis for each phase.

1684 C.2 DYNAMICS OF THE PARAMETERS FOR INDUCTION HEAD

1685 Recall the partial loss about the induction head:

$$1686 \mathcal{L}_{\text{IH}_2}(\theta) = \frac{1}{2} \mathbb{E}_X \left[\left(\frac{1}{1 + \alpha^*} (x_s)_{s=2}^{L-1} \cdot \text{softmax}\left((w^{*2} x_L x_{s-1})_{s=2}^{L-1}\right) - h \cdot (x_s)_{s=2}^{L-2} \cdot \text{softmax}\left((w^2 x_L x_{s-1})_{s=2}^{L-2}\right) \right)^2 \right].$$

$$1687$$

$$1688$$

$$1689$$

1690 **Technical simplification.** Unlike $\mathcal{L}_{\text{G}_4}(\theta)$, the denominators of the softmax terms
 1691 $\text{softmax}\left((w^{*2} x_L x_{s-1})_{s=2}^{L-1}\right)$ and $\text{softmax}\left((w^2 x_L x_{s-1})_{s=2}^{L-2}\right)$ in $\mathcal{L}_{\text{IH}_2}(\theta)$ depend on the input
 1692 tokens X , making it hard to derive a closed-form expression for $\mathcal{L}_{\text{IH}_2}(\theta)$. In Bai et al. (2023), the
 1693 authors consider a simplified transformer model, which replaces the softmax $\text{softmax}(z_1, \dots, z_L)$
 1694 with $\frac{1}{L} \exp(z_1, \dots, z_L)$. This approximation is nearly tight when $z_1, \dots, z_L \approx 0$. Notice that 1)
 1695 $w^2 x_L x_{s-1} \approx 0$ holds near the small initialization, i.e., for $w \approx \sigma_{\text{init}} \ll 1$. In fact, our analysis shows
 1696 that $w \approx \sigma_{\text{init}}$ is maintained over a long period. 2) $w^* = \mathcal{O}(1)$, which implies that $w^2 x_L x_{s-1} \approx 0$
 1697 for most input sequence. Thus, we adopt the simplification used in Bai et al. (2023), resulting in the
 1698 following approximation of the loss function:

$$1699 \mathcal{L}_{\text{IH}_2}(\theta) := \frac{1}{2} \mathbb{E}_X \left[\left(\frac{1}{1 + \alpha^*} \frac{1}{L-2} \sum_{s=2}^{L-1} \exp(w^{*2} x_L x_{s-1}) x_s - h \frac{1}{L-2} \sum_{s=2}^{L-2} \exp(w^2 x_L x_{s-1}) x_s \right)^2 \right].$$

$$1700$$

$$1701$$

$$1702$$

1703 Then by a straightforward calculation with Lemma D.3, we can derive its explicit formulation:

$$1704 \mathcal{L}_{\text{IH}_2}(\theta) = \frac{(1 - 4w^{*4})^{-\frac{1}{2}}}{2(1 + \alpha^*)^2(L-2)} + \frac{1}{2} \frac{h^2}{L-2} (1 - 4w^4)^{-\frac{1}{2}} - \frac{h(1 - (w^2 + w^{*2})^2)^{-\frac{1}{2}}}{(1 + \alpha^*)(L-2)}. \quad (22)$$

$$1705$$

$$1706$$

$$1707$$

1708 Furthermore, we can calculate GF dynamics as follows:

$$1709 \frac{dw}{dt} = \frac{h}{(1 + \alpha^*)(L-2)} (1 - (w^2 + w^{*2})^2)^{-\frac{3}{2}} \cdot (w^2 + w^{*2}) \cdot 2w - \frac{h^2}{L-2} (1 - 4w^4)^{-\frac{3}{2}} \cdot 4w^3,$$

$$1710$$

$$1711 \frac{dh}{dt} = \frac{1}{(1 + \alpha^*)(L-2)} (1 - (w^2 + w^{*2})^2)^{-\frac{1}{2}} - \frac{h}{L-2} (1 - 4w^4)^{-\frac{1}{2}}.$$

$$1712$$

$$1713$$

$$1714$$

1715 For simplicity, we denote:

$$1716 w^* := w^*, \quad h^* := \frac{1}{1 + \alpha^*}.$$

$$1717$$

$$1718$$

1719 **Part I. The trend and monotonicity of w, h .**

1720 For simplicity, we denote the tuning time point of h :

$$1721 T_2^h := \inf \left\{ t > 0 : \frac{dh(t)}{dt} = 0 \right\}.$$

$$1722$$

$$1723$$

$$1724$$

1725 In this step, we will prove the following three claims regarding the trend and monotonicity of w, h ,
 1726 which are essential for our subsequent analysis:

- 1727 • **(P1.1)** h initially increases beyond h^* , and then remains above this value.

- **(P1.2)** w keeps increasing but always stays below w^* .
- **(P1.3)** h increases before T_2^h , but decreases after T_2^h .

(P1.1) h initially increases beyond h^* , and then remains above this value.

We will prove that initially, h increases beyond h^* , and keeps growing beyond h^* . Define

$$T_1^h := \inf\{t > 0 : h(t) > h^*\},$$

we will prove that h remains above h^* thereafter.

For simplicity, we denote

$$\psi(x) = (1 - x^2)^{-\frac{1}{2}}, \quad \phi(x) = (1 - x^2)^{-\frac{3}{2}} \cdot x,$$

then the dynamics holds:

$$\begin{aligned} \frac{dh}{dt} &= \frac{h}{L-2} \psi(w^2 + w^{*2}) \left[\frac{h^*}{h} - \frac{\psi(2w^2)}{\psi(w^2 + w^{*2})} \right], \\ \frac{dw}{dt} &= \frac{2h^2 w}{L-2} \cdot \phi(w^2 + w^{*2}) \cdot \left[\frac{h^*}{h} - \frac{\phi(2w^2)}{\phi(w^2 + w^{*2})} \right]. \end{aligned}$$

Notice that $\frac{\phi(2w^2)}{\phi(w^2 + w^{*2})} < \frac{\psi(2w^2)}{\psi(w^2 + w^{*2})}$, $w < w^*$, while $\frac{\phi(2w^2)}{\phi(w^2 + w^{*2})} > \frac{\psi(2w^2)}{\psi(w^2 + w^{*2})}$, $w > w^*$.

We denote the first hitting time of h decreasing to h^* as \hat{t} :

$$\hat{t} := \inf\{t > T_2^h : h(t) < h^*\}.$$

If $w(\hat{t}) \geq w^*$, then at the first hitting time of w increasing to w^* , $\frac{dw}{dt} < 0$, which leads to a contradiction. If $w(\hat{t}) < w^*$, then $\frac{dh}{dt}|_{\hat{t}} > 0$, which also leads to a contradiction. Hence, $\hat{t} = +\infty$, which means that h always remains above h^* for $t > T_2^h$.

(P1.2) w keeps increasing but always below w^* .

We first prove that w always remains below w^* . We denote the first hitting time of w increasing to w^* as t' , then it is not difficult to see $\frac{dw}{dt}|_{t'} < 0$, which leads to a contradiction.

Next we prove that w keeps increasing throughout. We define the following functions

$$\begin{aligned} H &:= \frac{1}{1 + \alpha^*} \left(1 - (w^2 + w^{*2})^2\right)^{-\frac{3}{2}} (w^2 + w^{*2}) - h(1 - 4w^4)^{-\frac{3}{2}} \cdot 2w^2 \\ Q &:= \frac{1}{1 + \alpha^*} \left(1 - (w^2 + w^{*2})^2\right)^{-\frac{1}{2}} - h(1 - 4w^4)^{-\frac{1}{2}} \end{aligned}$$

If at some \bar{t} , $\frac{dw}{dt}$ reaches its zero point at the first time, then

$$\left. \frac{dH}{dt} \right|_{\bar{t}} = -h'(\bar{t})(1 - 4w^{*4})^{-\frac{3}{2}} \cdot 2w(\bar{t}) > 0,$$

which leads to a contradiction. Hence \bar{t} does not exist and w keeps increasing.

(P1.3) After the tuning point $t > T_2^h$, h will be monotonically decreasing.

The first sign-changing zero point of $\frac{dh}{dt}$ is T_2^h , then $Q(T_2^h) = 0$. $H(T_2^h) > 0$,

$$\begin{aligned} \left. \frac{dQ}{dt} \right|_{T_2^h} &= \frac{1}{1 + \alpha^*} (1 - (w(T_2^h)^2 + w^{*2})^2)^{-\frac{1}{2}} \cdot 2w(T_2^h) \cdot w'(T_2^h) \\ &\quad \cdot \left[(1 - (w(T_2^h)^2 + w^{*2})^2)^{-1} \cdot (w(T_2^h)^2 + w^{*2}) - (1 - 4w(T_2^h)^4)^{-1} \cdot 4w(T_2^h)^2 \right]. \end{aligned}$$

We can see that T_2^h is a sign-changing zero point only if

$$\frac{(1 - 4w(T_2^h)^4) \cdot (w(T_2^h)^2 + w^{*2})}{(1 - (w(T_2^h)^2 + w^{*2})^2) \cdot 4w(T_2^h)^2} < 1,$$

1782 i.e. we have:

$$1783 w(T_2^h) > w^\circ := \sqrt{\frac{3 - 4w^{*4} - \sqrt{(4w^{*4} - 3)^2 - 16w^{*4}}}{8w^{*2}}} \geq \frac{w^*}{2}, \quad (23)$$

1786 when $w^* = \mathcal{O}(1)$.

1788 Next we show that h keeps decreasing after T_2^h . We denote the first zero point of $\frac{dh}{dt}$ as t° , then
1789 $Q(t^\circ) = 0$. Since $\frac{dw}{dt}|_{t^\circ} > 0$, we have $\frac{dQ}{dt}|_{t^\circ} > 0$ which leads to a contradiction. Hence t° does not
1790 exist and h keeps decreasing after T_2^h .

1792 **Part II. Estimation of T_1^h, T_2^h , and the tight estimate of $w(t)$ before T_2^h .**

1793 At the first stage, we prove that h grows first and w barely increases. If $w \leq 0.01w^*$ and $h \leq$

$$1794 \frac{1}{1+\alpha^*} \frac{(1-w^{*4})^{-\frac{1}{2}}}{(1-0.01^4w^{*4})^{-\frac{1}{2}}},$$

$$1795 \frac{dh}{dt} \geq \frac{-1}{L-2} \left[h(1-0.01^4w^{*4})^{-\frac{1}{2}} - \frac{1}{1+\alpha^*} (1-w^{*4})^{-\frac{1}{2}} \right],$$

$$1796 h \geq \frac{1}{1+\alpha^*} \frac{(1-w^{*4})^{-\frac{1}{2}}}{(1-0.01^4w^{*4})^{-\frac{1}{2}}} - \left[\frac{1}{1+\alpha^*} \frac{(1-w^{*4})^{-\frac{1}{2}}}{(1-0.01^4w^{*4})^{-\frac{1}{2}}} - h(0) \right] \exp\left(\frac{-t}{(L-2)(1-0.01w^{*4})^{\frac{1}{2}}}\right). \quad (24)$$

1800 For h increasing from $h(0)$ to $\frac{1}{1+\alpha^*}$, it takes

$$1801 T_1^h \leq (1-0.01w^{*4})^{\frac{1}{2}}(L-2) \ln\left(\frac{1}{1 - \frac{(1-w^{*4})^{\frac{1}{2}}}{(1-0.01^4w^{*4})^{\frac{1}{2}}}}\right)$$

$$1802 \leq 2(L-2)(1 - \frac{1}{2}w^{*4}) = \mathcal{O}(L). \quad (25)$$

1803 For $0 \leq t \leq T_1^h$,

$$1804 \frac{dw}{dt} \leq \frac{1}{L-2} (1-4w^{*4})^{-\frac{3}{2}} \cdot w^{*2} \cdot 4w.$$

1805 Hence, it take $\mathcal{O}(L \log(1/\sigma_{\text{init}}))$ for w to reach $0.01w^*$, which allows sufficient time for h to reach
1806 $\frac{1}{1+\alpha^*}$ beforehand.

1807 Therefore, there exists a small constant $\varepsilon(w(0), w^*)$ only depends on $w(0)$ and w^* such that h is
1808 dominated by $1 + \varepsilon(w(0), w^*)$ times right hand side of (24), from which we deduce that (25) is a
1809 tight estimation of T_1^h instead of an upper bound, i.e. $T_1^h = \Theta(L)$.

1810 We then give a bound for $h(T_2^h)$. By $\frac{dh}{dt} = 0$,

$$1811 h(T_2^h)/h^* \leq \frac{(1-4w^4)^{\frac{1}{2}}}{(1-(w^2+w^{*2})^2)^{\frac{1}{2}}} := r(w).$$

1812 Moreover, $r(w)$ is an decreasing function of w for $w > w^\circ$, and w° is a function of w^* , we have

$$1813 h(T_2^h)/h^* \leq r(w^\circ) := R(w^*),$$

1814 where w° is a function about w^* , defined in Eq. (23). It is clear that

$$1815 R(w^* = 0) = 1, \quad R'(w^* = 0) = 0.$$

1816 Then using the continuity of $R'(\cdot)$ (in $[0, 0.4]$), there exists $c > 0$ such that $|R'(w^*)| < 0.04$ holds
1817 for all $0 < w^* < c$, which implies:

$$1818 R(w^*) = R(0) + \int_0^{w^*} R'(v)dv < 1 + 0.04w^*, \quad 0 < w^* < c.$$

1836 i.e., if $w^* = O(1)$, then $R(w^*) < 1 + 0.04w^*$. This implies:

$$1837 \quad h^* \leq h(t) \leq (1 + 0.04375w^*)h^*, \quad \forall t \geq T_1^h. \quad (26)$$

1840 By some computation, we can prove that $w^\circ(w^*)$ is an increasing function of w^* , and is always
1841 above $\frac{1}{2}w^*$. Thus we obtain a lower bound of w° for the estimation of lower bound of T_2^h :

1842 For the second stage, h barely changes and w starts to grow exponentially fast, and we use the tight
1843 estimation of $T_{1/2}^w := \inf \{t > 0 : w(t) > \frac{1}{2}w^*\}$ to give a lower bound of T_2^h . During this stage,
1844

$$1845 \quad \frac{dw}{dt} \leq \frac{2w}{(1 + \alpha^*)^2(L - 2)} \left[(1 - (w^2 + w^{*2})^2)^{-\frac{3}{2}} \cdot (w^2 + w^{*2}) - (1 - 4w^4)^{\frac{3}{2}} \right]$$

$$1846 \quad \leq \frac{2w}{(1 + \alpha^*)^2(L - 2)} (1 - 4w^{*4})^{\frac{3}{2}} \cdot 2w^{*2},$$

1849 and w has upper bound

$$1851 \quad w \leq w(0) \exp \left(\frac{4w^{*2}(1 - 4w^{*4})^{\frac{3}{2}}}{(1 + \alpha^*)(L - 2)} t \right). \quad (27)$$

1853 Hence, the lower bound of time for w to reach $\frac{1}{2}w^*$ is

$$1855 \quad T_{1/2}^w - T_1^h = \frac{(1 + \alpha^*)^2(L - 2)}{4w^{*2}(1 - 4w^{*4})^{\frac{3}{2}}} \ln \left(\frac{w^*}{2w(0)} \right),$$

1857 and lower bound for $T_{1/2}^w$ is

$$1859 \quad T_{1/2}^w \geq (L - 2) \left[\frac{(1 + \alpha^*)^2 \ln \left(\frac{w^*}{2w(0)} \right)}{4w^{*2}(1 - 4w^{*4})^{\frac{3}{2}}} - \ln \left(1 - (1 - w^{*4})^{\frac{1}{2}} \right) \right]$$

$$1860 \quad \geq \frac{(L - 2)(1 + \alpha^*)^2}{16w^{*2}} \ln \left(\frac{1}{w(0)} \right) = \Omega \left(\frac{(1 + \alpha^*)^2 L}{w^{*2}} \log \left(\frac{1}{\sigma_{\text{init}}} \right) \right). \quad (28)$$

1866 On the other hand, we estimate the lower bound of w . Let

$$1867 \quad C(x) = (1 - x^2)^{-\frac{3}{2}} \cdot x,$$

1869 then

$$1870 \quad C'(x) = 3(1 - x^2)^{-\frac{5}{2}}x^2 + (1 - x^2)^{-\frac{3}{2}} > 1, \quad 0 < x < 1,$$

$$1871 \quad C''(x) = 15x^3(1 - x^2)^{-\frac{7}{2}} + 6x(1 - x^2)^{-\frac{5}{2}} + 3x(1 - x^2)^{-\frac{5}{2}} > 0, \quad 0 < x < 1.$$

1872 $C(x)$ is a monotonically increasing convex function on $(0, 1)$ and $C(x) \geq x$.

1874 Using conclusions above, before w^2 increases to $\frac{1}{2\gamma(w^*) + \beta - 1}w^{*2}$ for some $\beta > 0$,

$$1875 \quad C(w^2 + w^{*2})$$

$$1876 \quad \geq C((2\gamma(w^*) + \beta)w^2)$$

$$1877 \quad \geq C(2\gamma(w^*) \cdot w^2) + C(\beta w^2) \quad (\text{Lemma D.6})$$

$$1878 \quad \geq \gamma(w^*) \cdot C(2w^2) + \beta w^2 \quad (C(ax) \geq aC(x), \text{ for } a > 1)$$

1882 then we have

$$1883 \quad \frac{dw}{dt} \geq \frac{2w}{(1 + \alpha^*)^2(L - 2)} (C(w^2 + w^{*2}) - \gamma(w^*) \cdot C(2w^2))$$

$$1884 \quad \geq \frac{2w}{(1 + \alpha^*)^2(L - 2)} \frac{\beta}{\gamma(w^*) + \beta} w^{*2}$$

1888 and

$$1889 \quad w \geq w(0) \exp \left(\frac{2\beta}{\gamma(w^*) + \beta} \frac{1}{(1 + \alpha^*)^2(L - 2)} w^{*2} t \right).$$

1890 Take $\beta = 2$, then

$$1891 \quad w \geq w(0) \exp\left(\frac{w^{*2}t}{(1+\alpha^*)^2(L-2)}\right), \forall t \in [0, T_{1/2}^w]. \quad (29)$$

1892 From the above inequality, (28) is not only an upper bound, but a tight estimation of $T_{1/2}^w$, i.e.

$$1893 \quad T_{1/2}^w = \Theta\left(\frac{(1+\alpha^*)^2L}{w^{*2}} \log\left(\frac{1}{\sigma_{\text{init}}}\right)\right).$$

1894 **Part II. Dynamics after the critical point $T_{1/2}^w$.**

1895 For simplicity, we consider:

$$1896 \quad v := w^2,$$

1897 and denote $v^* := w^{*2}, h^* := \frac{1}{1+\alpha^*}$. Then we focus on the dynamics of v and h .

1898 Additionally, we introduce a few notations used in this part:

$$1899 \quad \phi(x) := \frac{x}{(1-x^2)^{3/2}}, \quad \psi(x) := \frac{1}{(1-x^2)^{1/2}}.$$

1900 Then the dynamics of v and g are:

$$1901 \quad \frac{dv}{dt} = \frac{4vh}{L-2} (h^* \phi(v+v^*) - h \phi(2v)),$$

$$1902 \quad \frac{dh}{dt} = \frac{1}{L-2} (h^* \psi(v+v^*) - h \psi(2v)).$$

1903 **Step II.1. A coarse estimate of the relationship between v and h .**

1904 It is easy to verify the monotonicity that $\frac{dv}{dt} > 0$ and $\frac{dh}{dt} < 0$ for $t > t_2$. Additionally, we have

$$1905 \quad \frac{\psi(v+v^*)}{\psi(2v)} < \frac{h}{h^*} < \frac{\phi(v+v^*)}{\phi(2v)}.$$

1906 Then by Monotone convergence theorem, we obtain:

$$1907 \quad \lim_{t \rightarrow +\infty} v = v^*, \quad \lim_{t \rightarrow +\infty} h = h^*.$$

1908 **Step II.2. Convergence analysis by Lyapunov function.**

1909 This step aims to establish the convergence rate of v and h .

1910 In fact, the dynamics of v, h can be approximately characterized by their linearized dynamics. In contrast, the dynamics of p, g are controlled by high-order terms. Therefore, the proof for v and h is significantly simpler than the corresponding proof for p and g . We only need to consider the simplest Lyapunov function:

$$1911 \quad G(v, h) := \frac{1}{2} \left((v - v^*)^2 + (h - h^*)^2 \right).$$

1912 It is easy to verify that

$$1913 \quad (L-2) \frac{dG(v, h)}{dt} = (v - v^*) \frac{dv}{dt} + (h - h^*) \frac{dh}{dt}$$

$$1914 \quad = 4vh(v - v^*) (h^* \phi(v + v^*) - h \phi(2v)) + (h - h^*) (h^* \psi(v + v^*) - h \psi(2v))$$

$$1915 \quad = 4vh(v - v^*) \left(\phi(v + v^*) (h^* - h) - h (\phi(v + v^*) - \phi(2v)) \right)$$

$$1916 \quad \quad + (h - h^*) \left((h^* - h) \psi(v + v^*) + h (\psi(v + v^*) - \psi(2v)) \right)$$

$$1917 \quad = -4vh^2(v^* - v) (\phi(v + v^*) - \phi(2v)) - \psi(v + v^*) (h - h^*)^2$$

$$+ 4vh\phi(v + v^*)(v - v^*)(h^* - h) + h(h - h^*)(\psi(v + v^*) - \psi(2v)).$$

Let $v^* \leq 0.3 = \mathcal{O}(1)$. Recalling (23) and (26), as well as the monotonicity about p and w , we have:

$$\frac{v^*}{4} < v(t) < v^*; \quad h^* < h(t) < 1.02h^*, \quad \forall t > T_2^h.$$

Combining these estimates with the properties of ϕ and ψ , we have the following straight-forward estimates:

$$\begin{aligned} \phi(v + v^*) - \phi(2v) &= \phi'(\xi)(v^* - v) = \frac{1 + 2\xi^2}{(1 - \xi^2)^{5/2}}(v^* - v) \geq v^* - v; \\ \phi(v + v^*) &\leq \phi(2v^*) \leq 1; \\ \psi(v + v^*) &= \frac{1}{(1 - (v + v^*)^2)^{1/2}} \geq 1; \\ \psi(v + v^*) - \psi(2v) &= \psi'(\xi)(v^* - v) = \frac{\xi}{(1 - \xi^2)^{3/2}}(v^* - v) \leq 1.3v^*(v^* - v). \end{aligned}$$

Thus, we have the following estimate for the Lyapunov function:

$$\begin{aligned} &(L - 2) \frac{dG(v, h)}{dt} \\ &\leq -\frac{4}{1.02}v^*h^{*2}(v - v^*)^2 - (h - h^*)^2 \\ &\quad + 4.08v^*h^*(v - v^*)(h^* - h) + 1.3 \cdot 1.02v^*h^*(v^* - v)(h - h^*) \\ &= -\frac{4}{1.02}v^*h^{*2}(v - v^*)^2 - (h - h^*)^2 + 5.41v^*h^*(v^* - v)(h - h^*) \\ &\leq -3.92v^*h^{*2}(v - v^*)^2 - (h - h^*)^2 + \left(9.6v^*h^{*2}(v - v^*)^2 + \frac{3}{4}(h - h^*)^2\right) \\ &\leq -(3.92 - 9.6 \cdot 0.3)v^*h^{*2}(v - v^*)^2 - 0.25(h - h^*)^2 \leq -\frac{1}{4}v^*h^{*2}G(v, h). \end{aligned}$$

Consequently, we have the exponential bound for all $t > T_2^h$:

$$G(v(t), h(t)) \leq G(v(T_2^h), h(T_2^h)) \exp\left(-\frac{v^*h^{*2}}{4(L-2)}(t - T_2^h)\right), \quad \forall t > T_2^h,$$

This can imply:

$$\begin{aligned} (h(t) - h^*)^2 &= (h(T_2^h) - h^*)^2 \exp\left(-\Omega\left(\frac{w^{*2}(t - T_2^h)}{L(1 + \alpha^*)^2}\right)\right) \\ &= \mathcal{O}\left(h^{*2} \exp\left(-\Omega\left(\frac{w^{*2}(t - T_2^h)}{L(1 + \alpha^*)^2}\right)\right)\right), \quad \forall t > T_2^h; \\ (w(t) - w^*)^2 &= (w(T_2^h) - w^*)^2 \exp\left(-\Omega\left(\frac{w^{*2}(t - T_2^h)}{L(1 + \alpha^*)^2}\right)\right) \\ &= \mathcal{O}\left(w^{*2} \exp\left(-\Omega\left(\frac{w^{*2}(t - T_2^h)}{L(1 + \alpha^*)^2}\right)\right)\right), \quad \forall t > T_2^h. \end{aligned} \tag{30}$$

Notably, these proofs capture the **entire** training dynamics of w, h , from $t = 0$ to $t = T_1^h$, to $t = T_{1/2}^w \leq T_2^h$, and finally to $t \rightarrow +\infty$, providing a fine-gained analysis for each phase.

C.3 PROOF OF THEOREM 5.5

This theorem is a direct corollary of our analysis of the entire training dynamics in Appendix C.1 and C.2, leveraging the relationship between the parameters and the loss.

Proof of Phase I (partial learning).

By combining (19) and (21), it follows that: $\mathcal{L}_{\mathcal{G}_4}(\theta(0)) = \Theta(1)$. Moreover,

$$\mathcal{L}_{\mathcal{G}_4}(\theta(t)) = \mathcal{O}\left(\frac{1}{t}\right), \quad t > T_1^g = \mathcal{O}(1).$$

Thus, there exists a sufficiently large $T_1 = \Theta(1)$, such that:

$$\mathcal{L}_{\mathcal{G}_4}(\theta(T_1)) \leq 0.01\mathcal{L}_{\mathcal{G}_4}(\theta(0)).$$

Recalling our proof in Appendix C.2, for $t < T_1^h = \mathcal{O}(L)$, it holds that $h(t) < \sigma_{\text{init}} + \mathcal{O}(t/((1 + \alpha^*)L))$, $w(t) < \sigma_{\text{init}} + o(t/((1 + \alpha^*)L))$. Additionally, since $T_1 = \Theta(1) \ll \Theta(L)$, it follows that

$$w(T_1) = \mathcal{O}(\sigma_{\text{init}} + 1/L) < 2\sigma_{\text{init}} \ll w^*, \quad h(T_1) = \mathcal{O}(\sigma_{\text{init}} + 1/L) < 2\sigma_{\text{init}} \ll h^*.$$

Substituting these estimates into (22), we obtain by Lipschitz continuity of $\mathcal{L}_{\mathbb{H}_2}$:

$$\begin{aligned} |\mathcal{L}_{\mathbb{H}_2}(\theta(T_1)) - \mathcal{L}_{\mathbb{H}_2}(\theta(0))| &\leq 2\sigma_{\text{init}} \left(\left| \frac{\partial \mathcal{L}_{\mathbb{H}_2}}{\partial w} \right| + \left| \frac{\partial \mathcal{L}_{\mathbb{H}_2}}{\partial h} \right| \right) \\ &\leq 2\sigma_{\text{init}} \left(\mathcal{O}\left(\frac{1}{(1 + \alpha^*)L}\right) + o\left(\frac{1}{(1 + \alpha^*)L}\right) \right) \\ &\leq 0.01\mathcal{L}_{\mathbb{H}_2}(\theta(0)). \end{aligned}$$

Thus,

$$\mathcal{L}_{\mathbb{H}_2}(\theta(T_1)) \geq 0.99\mathcal{L}_{\mathbb{H}_2}(\theta(0)).$$

Proof of Phase II (plateau) + Phase III (emergence).

First, (27) and (29) ensures that w grows exponentially before $t < T_1^w$:

$$\sigma_{\text{init}} \exp\left(\frac{w^{*2}}{(1 + \alpha^*)^2(L - 2)}t\right) \leq w \leq \sigma_{\text{init}} \exp\left(\frac{4w^{*2}(1 - 4w^{*4})^{\frac{3}{2}}}{(1 + \alpha^*)(L - 2)}t\right).$$

Thus, we have:

$$w(t) = \sigma_{\text{init}} \exp\left(\Theta\left(\frac{w^{*2}t}{(1 + \alpha^*)^2L}\right)\right), \quad t < \Theta\left(\frac{(1 + \alpha^*)^2L}{w^{*2}} \log\left(\frac{1}{\sigma_{\text{init}}}\right)\right).$$

Now we define the observation time $T_o := T_1^h = \Theta(L)$. Notably,

$$h(T_o) = h^*, \quad w(T_o) < 0.01w^*.$$

The exponential growth of w further implies:

$$T_{0.01}^w := \{t > 0 : w(t) > 0.01w^*\} = \Theta\left(\frac{(1 + \alpha^*)^2L}{w^{*2}} \log\left(\frac{1}{\sigma_{\text{init}}}\right)\right).$$

Regarding the dynamics of h , by (26), we have $|h(t) - h(T_o)| < 0.02|h(T_o)|$, $\forall t \geq T_o$.

Now we incorporate these facts ($0 < w(T_o) < 0.01w^*$, $0 < w(T_{0.01}^w) \leq 0.01w^*$, $|h(T_{0.01}^w) - h(T_o)| < 0.02|h(T_o)|$, $h(T_o) = h^*$) into the loss (22). By the Lipschitz continuity of $\mathcal{L}_{\mathbb{H}_2}$, it is straightforward that

$$\mathcal{L}_{\mathbb{H}_2}(\theta(T_{0.01}^w)) \geq 0.99\mathcal{L}(\theta(T_o)).$$

Thus, we have established the lower bound for T_{II} :

$$T_{\text{II}} := \inf \{t > T_o : \mathcal{L}_{\mathbb{H}_2}(\theta(t)) \leq 0.99 \cdot \mathcal{L}_{\mathbb{H}_2}(\theta(T_o))\}$$

$$\geq T_{0.01}^w = \Omega \left(\frac{(1 + \alpha^*)^2 L}{w^{*2}} \log \left(\frac{1}{\sigma_{\text{init}}} \right) \right).$$

Combining the loss (22) and our parameter estimates (30), we obtain:

$$\mathcal{L}_{\text{IH}_2}(\theta(t)) = \mathcal{O} \left(\exp \left(-\Omega \left(\frac{w^{*2} t}{L(1 + \alpha^*)^2} \right) \right) \right), \quad t > T_2^h = \Theta \left(\frac{(1 + \alpha^*)^2 L}{w^{*2}} \log \left(\frac{1}{\sigma_{\text{init}}} \right) \right).$$

This implies the upper bound for T_{III} :

$$\begin{aligned} T_{\text{III}} &:= \inf \{ t > T_o : \mathcal{L}_{\text{IH}_2}(\theta(t)) \leq 0.01 \cdot \mathcal{L}_{\text{IH}_2}(\theta(T_o)) \} \\ &= T_{1/2}^w + \mathcal{O} \left((\alpha^* + 1)^2 L \log(1/\sigma_{\text{init}}) / w^{*2} \right) = \mathcal{O} \left((\alpha^* + 1)^2 L \log(1/\sigma_{\text{init}}) / w^{*2} \right). \end{aligned}$$

Combining the fact $T_{\text{II}} < T_{\text{III}}$, the lower bound for T_{II} , and the upper bound for T_{III} , we obtain the two-sided bounds for both T_{II} and T_{III} :

$$T_{\text{II}}, T_{\text{III}} = \Theta \left((\alpha^* + 1)^2 L \log(1/\sigma_{\text{init}}) / w^{*2} \right).$$

Proof of Phase IV (convergence).

By combining the loss (19), (22), and our parameter estimates (21), (30), it follows that:

$$\mathcal{L}_{\text{G}_4}(\theta(t)) = \mathcal{O} \left(\frac{1}{t} \right), \quad \mathcal{L}_{\text{IH}_2}(\theta(t)) = \mathcal{O} \left(\exp \left(-\Omega \left(\frac{w^{*2} t}{L(1 + \alpha^*)^2} \right) \right) \right), \quad t > T_{\text{III}}.$$

D USEFUL INEQUALITIES

Lemma D.1 (Corollary A.7 in Edelman et al. (2022)). *For any $\theta, \theta' \in \mathbb{R}^d$, we have*

$$\|\text{softmax}(\theta) - \text{softmax}(\theta')\|_1 \leq 2\|\theta - \theta'\|_\infty$$

Lemma D.2 (lemma E.1 in Wang and E (2024)). *For any $T \in \mathbb{N}_+$, $q, m \in \mathbb{N}_+$, there exists and absolute constant $C > 0$ and a $\phi_m^{\text{exp}}(t) = \sum_{k=1}^m \alpha_k e^{-\beta_k t}$ such that*

$$\|\mathbb{I}(\cdot = T) - \phi_m^{\text{exp}}(\cdot)\|_{\ell_1(\mathbb{N})} \leq \frac{C e^{q+0.01(q+1)T}}{m^q}.$$

where $\beta_k > 0$ holds for any $k \in [m]$.

Lemma D.3. $\mathbb{E}_{X,Y,Z} \exp(aXY)Z^2 = (1 - a^2)^{-1/2}$, $a < 1$.

Proof of Lemma D.3.

$$\begin{aligned} & \int \exp(aXY)Z^2 \left(\frac{1}{2\pi} \right)^{-3/2} \exp\left(-\frac{1}{2}X^2 - \frac{1}{2}Y^2 - \frac{1}{2}Z^2\right) dXdYdZ \\ &= \int \frac{1}{2\pi} \exp\left(-\frac{1}{2}(X - aY)^2 - \frac{1}{2}Y^2 + \frac{1}{2}a^2Y^2\right) d(X - aY)dY \\ &= \int \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}W^2\right) dW \quad (W = (1 - a^2)^{1/2}Y) \\ &= (1 - a^2)^{-1/2} \end{aligned}$$

□

Lemma D.4. Let $M(p) := \frac{1-e^{-p(L-2)}}{1-e^{-p}}$, then it holds that

$$\|\text{softmax}((-p(L-1-s))_{s=1}^{L-1})\|_2^2 = \frac{M(2p)}{M(p)^2}.$$

Definition D.5 (weakly majorizes). A vector $\mathbf{x} \in \mathbb{R}^n$ is said to *weakly majorize* another vector $\mathbf{y} \in \mathbb{R}^n$, denoted by $\mathbf{x} \prec_w \mathbf{y}$, if the following conditions hold:

1. $\sum_{i=1}^k x_{[i]} \leq \sum_{i=1}^k y_{[i]}$ for all $k = 1, 2, \dots, n-1$,
2. $\sum_{i=1}^n x_{[i]} = \sum_{i=1}^n y_{[i]}$,

where $x_{[i]}$ and $y_{[i]}$ are the components of \mathbf{x} and \mathbf{y} , respectively, arranged in decreasing order.

Lemma D.6 (Weighted Karamata Inequality). Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be a convex function, and let $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\mathbf{y} = (y_1, y_2, \dots, y_n)$ be two vectors in \mathbb{R}^n . If \mathbf{x} weakly majorizes \mathbf{y} (i.e., $\mathbf{x} \prec_w \mathbf{y}$), and w_1, w_2, \dots, w_n are non-negative weights such that

$$\sum_{i=1}^n w_i = 1,$$

then the following inequality holds:

$$\sum_{i=1}^n w_i f(x_i) \leq \sum_{i=1}^n w_i f(y_i).$$

E EXPERIMENTS

E.1 EXPERIMENTAL DETAILS FOR FIGURE 2

In line with our theoretical setting, we examine a simplified two-layer transformer, as described in (10). Specifically, the first layer only contains RPE (3) and the second layer consists of two heads: one uses only RPE and the other employs only dot-product structure. The target function is specified by (8) with $\alpha^* = 1$, $w^* = 0.49$, $\sigma_{\text{init}} = 0.01$, $L = 40$, and the distribution of each token is Gaussian, i.e., $x_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ for $i \in [L]$. Training is conducted by minimizing the squared loss (11) using online SGD with learning rate 0.1 and batch size $B = 1,000$. Following our theoretical analysis, the two layers are trained sequentially:

- Training Stage I: only the first layer is trained for 100,000 iterations;
- Training Stage II: Subsequently, only the second layer undergoes training for another 100,000 iterations.

The dynamical behavior of the Training Stage II is visualized in Figure 2.

E.2 ADDITIONAL EXPERIMENTS SUPPORTING OPTIMIZATION DYNAMICS

1. Standard transformers on real-world natural language dataset.

Setup. We train a two-layer two-head **standard transformer** with RPE (3) (without any simplification) on the **wikitext-2** dataset, a natural language dataset (Merity et al., 2016). The transformer has an embedding dimension $D = 128$ and FFN width $W = 512$. For this dataset, the input dimension is $d = 33278$. We use a context length $L = 200$ and batch size $B = 32$. The parameters are initialized with the scale 0.01. The model is trained for 1,500 epochs on 1 H100, using cross-entropy loss and SGD with learning rate 0.1, and the initialization scale is 0.01. It is important to note that **both layers are trained simultaneously**. The results are presented in Figure 3.

2. Discrete token distribution in toy setting.

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

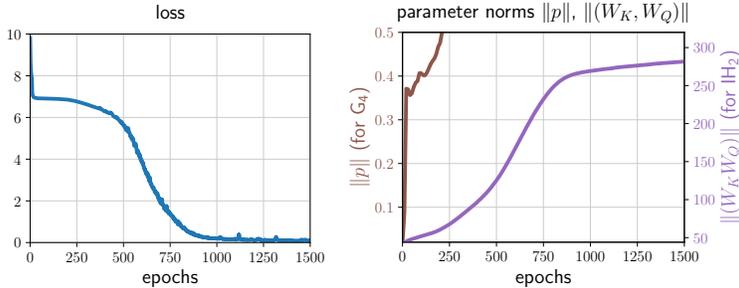


Figure 3: The loss and parameters for the experiment training a two-layer two-head **standard transformer** (without any simplification) on the **wikitext-2** dataset (Merity et al., 2016). Here, $\|p\|$ and $\|(W_K, W_Q)\|$ denote the Frobenius norms of all positional encoding parameters and all W_K, W_Q parameters across layers and heads, respectively. The results show that: the loss exhibits a clear plateau; position encoding p 's are learned first; and the dot-product structure W_K, W_Q are learned slowly at the beginning, resembling an exponential increase; additionally, as W_K, W_Q are learned, the loss escapes that plateau. These findings closely resemble the behavior observed in our toy model (Figure 2). This experiment provides further support for our theoretical insights regarding the **time-scale separation** between the learning of positional encoding and the dot-product structure.

Setup. We modified the Gaussian input distribution used in the setup for Figure 2 to a boolean input distribution, where each input token, where each input token $x_i \stackrel{iid}{\sim} \text{Unif}(\{\pm 1\})$ for $i \in [L]$. All other experimental setups remain the same as in the setup for Figure 2. The training dynamics of Stage (ii) are presented in Figure 4. We can see clearly that the dynamical behavior of the learning process is nearly the same as the one observed for Gaussian inputs in Figure 2.

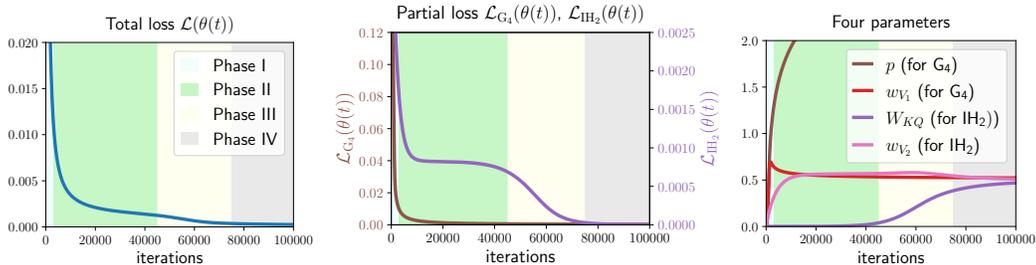


Figure 4: Visualization of the total loss, partial loss, and the parameter dynamics, for the experiment on **discrete token distribution** (Boolean, $X \sim \text{Unif}(\{\pm 1\}^L)$) in our toy setting with $\alpha^* = 1, w^* = 0.49, \sigma_{\text{init}} = 0.01, L = 40$. The figure clearly shows that transformer learns the 4-gram component first and then, starts to learn the induction head mechanism. Notably, the entire dynamics exhibit four phases. These results are **extremely similar** to that observed with Gaussian inputs, as shown in Figure 2.

3. Adam in high-dimensional toy setting.

Setup. We modified the setup for Figure 2 to employ a high-dimensional model ($D = 100$). Specifically, the target is $w^* = 0.49I_D/D$, the dot-product parameters are $W_K, W_Q \in \mathbb{R}^D$, initialized such that $\|W_K\|_F, \|W_Q\|_F = \sigma_{\text{init}}$. Additionally, for the Adam optimizer, we use learning rate $5e-4$. All other experimental setups remain the same as in the setup for Figure 2.

The training dynamics are depicted in Figure 5, where, for comparison, results using GD are also presented. In both scenarios, the learning process begins with the 4-gram pattern, followed by a gradual learning phase of the induction head mechanism. Notably, within the given number of iterations, GD remains stuck in the plateau, whereas Adam successfully escapes that plateau.

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

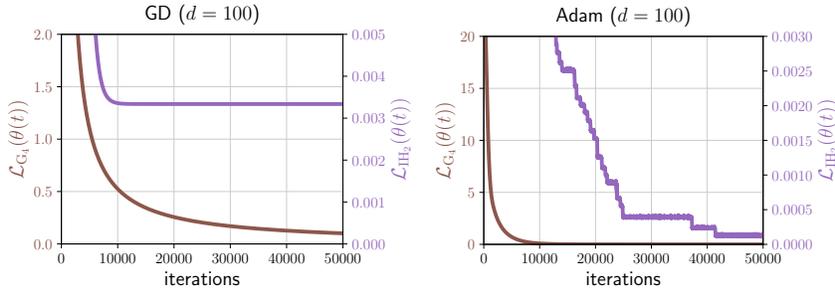


Figure 5: Partial loss for the experiment comparing **GD v.s. Adam optimizer** in high-dimensional settings ($D = 100$). In this setting, a larger D increases the difficulty of the transition from the lazy regime (learning 4-gram) to the rich regime (learning induction head). The results indicate that: (1) GD learns the 4-gram component first but becomes stuck in a plateau when learning induction head; (2) Adam, while eventually transitioning from the lazy regime (learning 4-gram) to the rich regime (learning induction head), experiences a **challenging** transition characterized by **multiple plateaus** during learning induction heads. This finding closely resembles the dynamics for GD.

E.3 EXPERIMENTS SUPPORTING APPROXIMATION RESULTS

1. Supporting the **necessity** of the required H and D in Theorem 4.3.

Setup. We train two-layer transformers (without FFN layers) with varying H and D to learn the generalized induction head (6) with $n = 4$. The input sequence $X = (x_1, \dots, x_L)$ is boolean, with $x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\{\pm 1\})$ and $L = 10$. Each model is trained for 200,000 iterations using squared loss and (online) Adam optimizer with learning rate $5e-4$ and batch size $B = 100$. Both layers are trained simultaneously. The results for the models with $D = H = 8$ and $D = H = 2$ are presented in Figure 3.

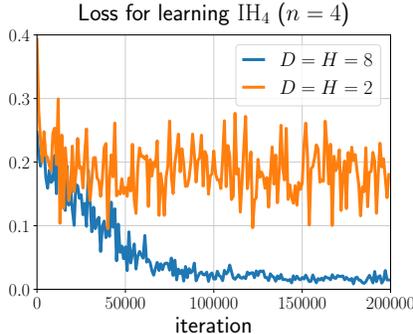


Figure 6: Results supporting the **necessity** of the required number of heads H and embedding dimension D in Theorem 4.3. We train two-layer transformers with varying H and D to learn the target in Eq. (6) with $n = 4$. The results indicate that the transformer with $H = D = 8 (> n)$ successfully expresses this task, while the transformer with $H = D = 2 (< n)$ fails. These results confirm that the sufficient conditions provided in Theorem 4.3 ($H \gtrsim n$ and $D \geq nd$, where $d = 1$ in our setting) are also nearly necessary.

2. Supporting our **construction** in Theorem 4.3.

Setup. We linear probing experiments (Alain and Bengio, 2016) on the transformers with $H = D = 8$ trained in the above experiment (Figure 6). For each checkpoint model TF, we denote its output in the first layer on the input sequence X as $\text{TF}^{(1)}(X)$. The probing loss is measured by $\text{dist}(X_{:-n+1:}; \text{TF}^{(1)}(X)) = \min_{P \in \mathbb{R}^{D \times n}} : \sum_{s=n}^L \|X_{s-n+1:s} - \text{TF}_s^{(1)}(X)P\|$, where $n = 4, L = 10$,

and $X = (x_1, \dots, x_L)$ is generated by $x_i \stackrel{i.i.d.}{\sim} \text{Unif}(\{\pm 1\})$ with testing batch 1000. The results are shown in Figure 7.

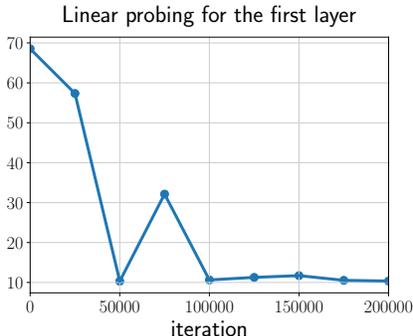


Figure 7: Probing results supporting our **construction** in Theorem 4.3. First, we train a two-layer transformer with head $H = 8$ and embedding dimension $D = 8$ to learn Eq. (6) with $n = 4$, and the checkpoints are stored during training. For each checkpoint model TF, we denote its output in the *first layer* on the input sequence X as $\text{TF}^{(1)}(X)$. To validate whether it encodes the semantic information $X_{s-n+2:s}$ near each x_s , as predicted by our construction, we conduct a standard linear probing experiment (Alain and Bengio, 2016). Specifically, we measured $\text{dist}(X_{-n+1:}; \text{TF}^{(1)}(X)) = \min_{P \in \mathbb{R}^{D \times n}} : \sum_{s=n}^L \|X_{s-n+1:s} - \text{TF}_s^{(1)}(X)P\|$. As the results shown, the probing loss decreases significantly during training, confirming our key construction in Theorem 4.3: **the first layer is responsible for extracting local semantic information $X_{s-n+2:s}$ near each x_s** , enabling the second layer to generate the final output.

F DETAILED COMPARISON WITH RELATED WORKS

In this section, we discuss the relationship between our work and two closely related studies: Bietti et al. (2024) and Edelman et al. (2024).

Comparison with Bietti et al. (2024).

- **Approximation analysis:**

- Bietti et al. (2024) focus primarily on the implementation of the vanilla induction head. In contrast, our study extends this analysis by investigating not only how two-layer transformers achieve vanilla induction heads (Eq. (4)) but also how they implement generalized induction heads, i.e., in-context n-grams (Eqs. (6) and (7)).
- Furthermore, our work provides explicit approximation rate results, offering insights into the distinct roles of multiple heads, positional encoding, dot-product structure, and FFNs in implementing these induction heads.

- **Optimization analysis:**

- *Study objective:* While Bietti et al. (2024) examines the transition from 2-gram to induction head, our work focuses on the transition from 4-gram to induction head.
- *study methods:* Bietti et al. (2024) conducts extensive experiments supported by partial theoretical properties but does not fully characterize the training dynamics theoretically. In contrast, our study provides a precise theoretical analysis of the entire training process in a toy model, uncovering the sharp transition from 4-gram to induction head.
- *Main insights:* Bietti et al. (2024) emphasizes the the role of weight matrices as associative memories and the impact of data distributional properties. Our analysis, on the other hand, identifies two primary drivers of the transition: (1) the time-scale separation due to low-

2322 and high-order parameter dependencies in self-attention; (2) the speed differences caused
2323 by the relative proportions of the two components in the mixed target.
2324

2325 **Comparison with Edelman et al. (2024).** The primary connection between Edelman et al. (2024) and
2326 our work lies in the optimization analysis. Specifically, Edelman et al. (2024) focuses on the transition
2327 from uni-gram to bi-gram mechanisms in Markov Chain data. In contrast, our study investigates
2328 the transition from 4-gram to in-context 2-gram mechanisms (induction head). Additionally, we
2329 theoretically identify two primary drivers of the transition: (1) the time-scale separation due to low-
2330 and high-order parameter dependencies in self-attention; (2) the speed differences caused by the
2331 relative proportions of the two components in the mixed target.

2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375