

ON THE SCALING THEORY OF MULTI-LAYER TRANSFORMERS

Anonymous authors

Paper under double-blind review

ABSTRACT

The scaling law, a cornerstone of Large Language Model (LLM) development, predicts improvements in model performance with increasing computational resources. Yet, while empirically validated, its theoretical underpinnings remain poorly understood. This work formalizes the learning dynamics of transformer-based language models as an ordinary differential equation (ODE) system, then approximates this process to kernel behaviors. Departing from prior toy-model analyses, we rigorously analyze one-pass stochastic gradient descent (SGD) training for multi-layer transformers on sequence-to-sequence data with arbitrary data distribution, closely mirroring real-world conditions. Our analysis characterizes the convergence of generalization error to the irreducible risk as computational resources scale with data. We derive an excess risk of $\Theta(C^{-1/8})$ for computational cost C . The theory reveals a phase transition: under specific conditions, the generalization risk’s upper bound drops sharply to $\exp(-C^{1/4})$ before reverting to its original decay rate. This transition delineates three scaling regimes—*classical*, *over-parameterization*, and *data-limited*—which we analyze for their impact on scaling efficiency and the emergence of grokking.

1 INTRODUCTION

The emergence of transformer-based Large Language Models (LLMs) such as ChatGPT (ChatGPT, 2022), GPT-4 (Achiam et al., 2023; Bubeck et al., 2023), DeepSeek (Bi et al., 2024; Liu et al., 2024b;a; Grattafiori et al., 2024), and LLaMA (Touvron et al., 2023a;b; Grattafiori et al., 2024) has driven transformative advancements across multiple domains (Brown et al., 2020). Tasks like code generalization (Li et al., 2023; Guo et al., 2024), conversational systems (Maaz et al., 2023; Xu et al., 2023; Zheng et al., 2024), and mathematical reasoning (Hendrycks et al., 2020; Yu et al., 2023a; Yao et al., 2023), once considered exclusive to human expertise, are now routinely mastered by these AI systems.

This remarkable progress is fundamentally tied to computational scaling. Empirical evidence reveals a consistent pattern: as the compute budget for optimally training and deploying language model increases, their demonstrated intelligence scales correspondingly (Kaplan et al., 2020; Hoffmann et al., 2022). This phenomenon has been systematically categorized into the following principle:

Scaling Law (Kaplan et al., 2020; Hoffmann et al., 2022): Model capabilities improve predictably with increased computational investment during training, achieved through three-dimensional scaling: model parameter count, training duration, and dataset size.

While the scaling law has been extensively validated through empirical research, a critical gap persists in its theoretical foundation. Current literature lacks a mathematically rigorous framework to explain why these principles exhibit such predictable improvements in model performance, or to systematically justify their reliability in guiding the development of massive-scale transformer architectures (Lin et al., 2024). Existing theoretical studies on scaling laws have predominantly focused on traditional statistical models and simplified neural architectures. Key investigations include analyses of shallow attention networks (Lyu et al., 2025), linear regression (Lin et al., 2024; Daliri et al., 2024), kernel regression (Chen et al., 2025), and data selection methodologies (Jain et al., 2018), among others. However, these works lack generalizability to modern deep learning systems. Notably, the scaling theory underpinning transformer-based LLMs (Vaswani et al., 2017) — the dominant paradigm in contemporary AI — remains largely unexplored.

Our work addresses this gap by establishing theoretical foundations for computational scaling effectiveness in training multi-layer transformer-based language models on sequence-to-sequence data. Specifically, we derive a quantifiable relationship between computational expenditure and model capabilities, formulating an optimized error bound for training objectives that explicitly depends on allocated computational resources. In particular, the goal of this work is to address the following critical question:

What are the theoretical limits for computational resource allocation to ensure the convergence of the generalization error bound during transformer-based language model training?

We present the first comprehensive analysis of the training dynamics of transformer-based language models using the one-pass Stochastic Gradient Descent (SGD) algorithm, along with a convergence guarantee for arbitrary training error. This topic has been largely unexplored due to the inherent complexity of attention mechanisms, the multi-layered structure of transformers, and the extensive matrix computations required for sequence-to-sequence data processing. To address these challenges, we adopt a kernel-based analytical framework to investigate the scaling behavior. By analyzing this scenario, we establish the lower and upper bounds for the expected risk on the whole data distribution. Overall, we make the following contributions:

- We first simplify the complicated matrix computation to the parallel vector computation utilizing the decoder-only property of a generative transformer-based language model. Therefore, we formalize its explicit learning dynamics, which we assume the learning behavior of LLMs is constrained to the kernel regimes due to the over-parameterization. This is referred to the *Lazy Learning* (Jacot et al., 2018; Du et al., 2019), where the model merely memorizes all data points during optimization. (See Section 4)
- We further explore several benefits that emerge under the trend of model scaling, e.g., converging kernel perturbation, which means that once kernel behavior exists at initialization, it will remain stable during training. Moreover, we demonstrate how the training convergence rate exponentially improves with linearly increasing model depth. Finally, we showcase the guaranteed training convergence and approximation bound related to vital scaling factors (model size, dataset size, training time, compute cost). (See Section 5)
- The main contribution of this work is a three-stage upper bound on the in-distribution expected risk, delineating the scaling process into the *classical stage*, the *over-parameterization stage*, and the *data-limited stage*. Our theoretical analysis shows that the generalization error initially decreases at a high rate with reducing scaling profit (i.e., inefficient), then enters a period of slower error reduction but growing scaling profit (i.e., highly efficient), which aligns with the empirically observed grokking phenomenon (Power et al., 2022). Finally, the convergence rate diminishes again in the *data-limited stage*. We provide theoretical guarantees for this risk landscape and present supporting experimental results. (See Section 6 and Section 7)

2 RELATED WORK

Scaling Law. Several recent works have empirically explored scaling laws in deep neural networks (Kaplan et al., 2020; Hoffmann et al., 2022; Rosenfeld et al., 2021; Hestness et al., 2017; Rosenfeld et al., 2019). The study of neural scaling laws can also be traced back to earlier foundational works (Caponnetto & De Vito, 2007; Steinwart et al., 2009; Ahmad & Tesauero, 1988). From a theoretical perspective, various solvable models have been developed using random feature models (Bahri et al., 2024; Atanasov et al., 2021; 2024; Bordelon et al., 2024; Paquette et al., 2024) to analyze neural scaling laws under specific limits. Additionally, theoretical analyses on linear models (Wei et al., 2022a; Bordelon et al., 2020; Seleznova & Kutyniok, 2022; Bordelon & Pehlevan, 2021; Lin et al., 2024; Lyu et al., 2025) have significantly advanced our understanding of scaling laws. In contrast to these studies, our work focuses on the sequence-to-sequence stochastic training of multi-layer transformer-based language models, a topic that has not been widely discussed in prior research.

Neural Tangent Kernel and Learning Theory. The Neural Tangent Kernel (NTK), introduced by Jacot et al. (2018), has become a foundational framework for understanding the gradient flow of neural networks during training. It reveals that, in the infinite-width limit, neural networks are equivalent to Gaussian processes at initialization. This equivalence has been extensively studied in numerous works (Li & Liang, 2018; Du et al., 2019; Song & Yang, 2019; Allen-Zhu et al., 2019; Wei et al.,

2019; Bietti & Mairal, 2019; Lee et al., 2020; Chizat & Bach, 2020; Shi et al., 2021; Zhou et al., 2021; Seleznova & Kutyniok, 2022; Gao et al., 2023; Li et al., 2024; Shi et al., 2024b), which highlight the robust performance and learning capabilities of over-parameterized neural networks. The NTK framework has gained popularity for its ability to elucidate the emerging abilities of large-scale neural networks. Notable advancements include the introduction of the Convolutional NTK (CNTK) by Arora et al. (2019), the Recurrent NTK by Alemohammad et al. (2020), and the concept of infinite attention via NNGP and NTK for attention networks by Hron et al. (2020). Furthermore, Malladi et al. (2023) examined the training dynamics of fine-tuning LLMs using NTK, demonstrating its efficiency in optimizing these systems. These contributions underscore the growing importance of NTK in the theoretical analysis of modern neural networks.

Science of Transformer-based Language Models. The complex architecture and stochastic optimization processes of transformer-based language models pose significant challenges for theoretical analysis. However, developing theoretical guarantees for LLMs is essential for advancing their design and performance. Recent research has addressed various aspects of LLMs, including improving efficiency (Alman & Song, 2023; 2024a;b; Han et al., 2024; Kacham et al., 2023; Addanki et al., 2023; Deng et al., 2024; Shi et al., 2024a), optimizing training processes (Deng et al., 2023; Li et al., 2024), analyzing "white-box" transformers (Yu et al., 2023b;c; Ferrando et al., 2024; Pai et al., 2024), and investigating the emergent abilities of LLMs (Brown et al., 2020; Wei et al., 2022b; Allen-Zhu & Li, 2023a;b;c; 2024). By bridging theoretical understanding with practical advancements, these studies provide valuable insights for the development of the next generation of AI systems.

3 PRELIMINARY

This section provides the preliminary for our analysis, where we introduce some basic notations in Section 3.1, and present the problem setup in Section 3.2. We encourage the reader to refer to Appendix B for the formal technical preliminary.

3.1 BASIC NOTATIONS

Let $[d] = \{1, 2, \dots, d\}$. For $u \in \mathbb{R}^d$, the ℓ_p -norm is $\|u\|_p := (\sum_{k=1}^d |u_k|^p)^{1/p}$. The Frobenius norm of $U \in \mathbb{R}^{d_1 \times d_2}$ is $\|U\|_F := (\sum_{(k_1, k_2) \in [d_1] \times [d_2]} U_{k_1, k_2}^2)^{1/2}$. For a matrix $A \in \mathbb{R}^{d \times d}$, $\lambda_{\min}(A)$ denotes its smallest eigenvalue. The indicator function $\mathbb{I}\{E_1, \dots, E_n\}$ equals 1 if all events E_1, \dots, E_n occur, and 0 otherwise. The mapping $\text{mat} : \mathbb{R}^{d^2} \rightarrow \mathbb{R}^{d \times d}$ reshapes a vector a into a matrix such that $\text{mat}_{k_1, k_2}(a) = a_{(k_1-1)d+k_2}$ for $(k_1, k_2) \in [d] \times [d]$. Conversely, $\text{vec} : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{nd}$ flattens a matrix A into a vector with $\text{vec}_k(A) = A_{\lfloor k/d \rfloor, k - \lfloor k/d \rfloor \cdot d}$ for $k \in [nd]$. For a function $f : X \rightarrow \mathbb{R}^{d_1 \times d_2}$, $f_{k_1, k_2}(x)$ denotes the (k_1, k_2) -entry of $f(x)$. $e_k \in \mathbb{R}^d$ is the standard basis vector with a 1 in the k -th entry and 0 elsewhere. Finally, we denote $a \wedge b := \max\{a, b\}$ for $a, b \in \mathbb{R}$.

3.2 SETUPS

Data Distribution. We consider a sequence-to-sequence regression task with an input space $\mathcal{X} \subseteq \mathbb{R}^{L \times d}$ as the space of encoded input sequence¹, and $\mathcal{Y} \subseteq [C_{\text{lower}}, C_{\text{upper}}]^{L \times d}$ for constant $C_{\text{lower}}, C_{\text{upper}} \in \mathbb{R}$ is the space of encoded target output. We denote F^* as the optimal measurable function mapping $\mathcal{X} \rightarrow \mathcal{Y}$ with minimum risk. Given a model class \mathcal{F} , then for the distribution $\mathcal{D} = \{(X, F^*(X) + \Xi), \Xi \in \mathbb{R}^{L \times d} \text{ is some random noise}\} \subset \mathcal{X} \times \mathcal{Y}$ and model function $F \in \mathcal{F}$, the expected risk (we consider as the generalization error bound) and excess risk of F is defined as:

$$\text{Expected Risk: } \mathcal{R}(F) := \mathbb{E}_{(X, Y) \sim \mathcal{D}} [\|F(X) - Y\|_F^2], \quad \text{Excess Risk: } \Delta \mathcal{R}(F) := \mathcal{R}(F) - \mathcal{R}(F^*).$$

Besides, we have an accessible dataset $\mathbb{D} = \{(X_i, Y_i)\}_{i=1}^n \subset \mathcal{D}$ where each data point independently and uniformly sampled from \mathcal{D} . The random noise $\Xi \in \mathbb{R}^{L \times d}$ is centered by $\mathbf{0}_{L \times d}$. For any input matrix $X \sim \mathcal{X}$, $\|X_\ell\|_2 = \Theta(1)$ holds for each token vector $\ell \in [L]$ due to the utilization of RMS normalization (Zhang & Sennrich, 2019) after the embedding layer.

Model Function. The standard transformer architecture introduced in Vaswani et al. (2017) integrates multiple self-attention layers with token-wise feed-forward layers. The fundamental

¹We choose the max length of sequence L considerably large, for sequences with a length less than L , we use padding to fill them.

architecture, decoder-only transformers (Radford et al., 2019), processing a sequence of L tokens, each represented by a d -dimensional embedding vector, which are compactly arranged into a matrix $X \in \mathbb{R}^{L \times d}$. An N -layer transformer model is formally defined as:

$$F(X, \theta) := \varepsilon \cdot F_{(N)}(F_{(N-1)}(\cdots F_{(2)}(F_{(1)}(X + E, \theta), \theta) \cdots), \theta), \quad (1)$$

where $E \in \mathbb{R}^{L \times d}$ is the positional embedding matrix² and θ is the set of all trainable parameters. $\varepsilon > 0$ is the grokking coefficient, which we show the relationship between its value and the grokking phenomenon Power et al. (2022); Nanda et al. (2023); Liu et al. (2022) in Section 6.2. Each $F_{(\nu)}$ (for $\nu \in [N]$) represents the ν -th transformer block and is given by:

$$F_{(\nu)}(X, \theta) := \frac{\omega}{\sqrt{m}} \text{ReLU} \left(\text{Softmax} \left(\kappa \cdot XU_{(\nu)}X^\top + M \right) XW_{(\nu)} \right) A_{(\nu)} + X,$$

where $U_{(\nu)} \in \mathbb{R}^{d \times d}$, $W_{(\nu)}, A_{(\nu)}^\top \in \mathbb{R}^{d \times m}$ are model parameters. κ is the scaling factor of attention, ω is the scaling coefficient of output. $M \in \mathbb{R}^{L \times L}$ is the causal attention mask. We especially use $w_{(\nu),r}, a_{(\nu),r} \in \mathbb{R}^d$ to denote the r -th column of $W_{(\nu)}$ and the r -th row of $A_{(\nu)}$, respectively.

Initialization and Training. For every layer $\nu \in [N]$, the corresponding parameters of $F_{(\nu)}$ is denoted as $U_{(\nu)}, W_{(\nu)}, A_{(\nu)}$. Each entry of $U_{(\nu)}, W_{(\nu)}$ is initialized from the standard Gaussian distribution $\mathcal{N}(0, 1)$, we then denote them as $U_{(\nu)}(0), W_{(\nu)}(0)$. Besides, each entry of $A_{(\nu)}$ is initialized from a uniform distribution $\text{Uniform}\{-1, +1\}$ and is frozen during training. The flattened vector of the whole trainable parameters is denoted as $\theta(0) \in \mathbb{R}^M$ where $M = N(md + d^2)$ is the number of trainable parameters.

Given the training dataset $\mathbb{D} = \{(X_i, Y_i)\}_{i=1}^n \subset \mathcal{D}$, we define the overall training objective as:

$$\mathcal{L}(t, \mathbb{D}) := \mathbb{E}_{(X, Y) \sim \mathbb{D}} [\|F(X, \theta(t)) - Y\|_F^2].$$

Therefore, we consider a combination of the one-pass stochastic gradient descent (SGD) algorithm and *gradient flow* to update. At t -step optimization, we sample a unbiased subset $\mathbb{B}(t) \subseteq \mathbb{D}$ satisfying $\mathbb{E}[\mathcal{L}(t, \mathbb{B}(t))] = \mathcal{L}(t, \mathbb{D})$ and $\int_0^T \mathbb{B}(t) dt = \mathbb{D}$ for training time $T > 0$. Then the ordinary differential equation (ODE) of $U_{(\nu)}(t), W_{(\nu)}(t)$ and their update rule are given by:

$$\begin{aligned} \frac{d}{dt} U_{(\nu)}(t) &= -\frac{d}{dU_{(\nu)}(t)} \mathcal{L}(t, \mathbb{B}(t)), & \frac{d}{dt} W_{(\nu)}(t) &= -\frac{d}{dW_{(\nu)}(t)} \mathcal{L}(t, \mathbb{B}(t)), \\ U_{(\nu)}(t + \tau) &= U_{(\nu)}(t) + \int_t^{t+\tau} \frac{d}{ds} U_{(\nu)}(s) ds, & W_{(\nu)}(t + \tau) &= W_{(\nu)}(t) + \int_t^{t+\tau} \frac{d}{ds} W_{(\nu)}(s) ds. \end{aligned} \quad (2)$$

Hence, we denote the training algorithm that depends on training time and dataset size as $\mathcal{A}_{T,n}(\theta(0), \mathbb{D}) := \{\theta(0) + \int_0^T -\frac{d}{d\theta} \mathcal{L}(s, \mathbb{B}(s)) ds, \mathbb{B}(s) \subseteq \mathbb{D}, s \in [0, T]\}$.

4 LEARNING DYNAMICS OF SCALING TRANSFORMERS

In this section, we provide the explicit model learning dynamics formulations layer by layer. We first introduce several key simplifications in Section 4.1 to facilitate the preliminary analysis. Thus, Section 4.2 formulates an explicit ODE of the learning dynamics of scaling multi-layer transformer-based language models.

4.1 SIMPLIFICATIONS

Utilizing the attention network’s decoder-only property, it is obvious that the matrix computation can be parallelized to vector computation. We use $X_{i, \leq \ell} \in \mathbb{R}^{\ell \times d}$ to denote the first- ℓ tokens ($\forall \ell \in [L]$) of matrix X_i in \mathbb{D} . Hence, we compact the outputs of the model function and targets on the whole dataset to two matrices $F(t), Y \in \mathbb{R}^{n \times L \times d}$, where $F_{(i-1)L+\ell}(t) = F_\ell(X_i, \theta(t)) \in \mathbb{R}^d$ and $Y_{(i-1)L+\ell} = Y_{i,\ell}$ for each (X_i, Y_i) in training dataset \mathbb{D} and $\ell \in [L]$. Therefore, we derive $\mathcal{L}(t, \mathbb{D}) = \frac{1}{n} \|F(t) - Y\|_F^2$ (See Lemma C.3). Besides, we list following helpful functions ($(i, \ell) \in [n] \times [L]$):

²We choose $E = \mathbf{0}_{L \times d}$ (NoPE, No Positional Embedding, (Kazemnejad et al., 2023)) or ignore it as a fixed matrix (this can be regarded as a part of the training dataset) in the range of this paper.

- (Hidden State) $\Lambda_{(\nu),i}(t) := F_{(\nu)}(\Lambda_{(\nu-1),i}(t), \theta(t)) \in \mathbb{R}^{L \times d}$ for $\nu \in [N]$, $\Lambda_{(0),i}(t) = X_i + E$.
- (Attention Scores) $\sigma_{(\nu),(i-1)L+\ell}(X) = \text{Softmax}_{\ell}(\Lambda_{(\nu),i}(t)U_{(\nu)}(t)\Lambda_{(\nu),i}(t)^{\top} + M) \in \mathbb{R}^L$.
- (Attention Output) $o_{(\nu),(i-1)L+\ell}(t) := \Lambda_{(\nu-1),i}(t)^{\top} \cdot \sigma_{(\nu),(i-1)L+\ell}(t) \in \mathbb{R}^d$.
- (ℓ -th Token of Hidden State)

$$\mu_{(\nu),(i-1)L+\ell}(t) := \frac{\omega}{\sqrt{m}} \sum_{r=1}^m a_{(\nu),r} \cdot \phi(\langle o_{(\nu),(i-1)L+\ell}(t), w_{(\nu),r}(t) \rangle) \in \mathbb{R}^d,$$

where $\phi(x) := \max\{0, x\}, \forall x \in \mathbb{R}$. $\mu_{(0),(i-1)L+\ell}(t) = X_{i,\ell} + E_{\ell}$.

- (Model Output) $F_{(i-1)L+\ell}(t) = \varepsilon \cdot \sum_{\nu=0}^N \mu_{(\nu),(i-1)L+\ell}(t) \in \mathbb{R}^d$.

4.2 KEY DERIVATION FOR LEARNING DYNAMICS

The primary challenge in understanding the learning dynamics of finite-deep transformers is the complex analysis of gradient flow, which differs from the study of shallow or infinite-deep neural networks. We overcome the complexity by cleverly utilizing the multivariable chain rules.

First, we define the kernel matrix at ν -th layer as $H_{(\nu)} \in \mathbb{R}^{nL \times nL}$ and its (i, j) -th entry ($\forall (p, q) \in [nL] \times [nL]$) is defined as:

$$H_{(\nu),p,q}(t) := \underbrace{\langle \beta_{(\nu),p}(t), \beta_{(\nu),q}(t) \rangle}_{\text{kernel w.r.t. } W_{(\nu)}(t)} + \underbrace{\langle \gamma_{(\nu),p}(t), \gamma_{(\nu),q}(t) \rangle}_{\text{kernel w.r.t. } U_{(\nu)}(t)},$$

Here, we let:

$$\begin{aligned} \beta_{(\nu),p}(t) &:= \frac{\omega}{\sqrt{m}} \underbrace{o_{(\nu),p}(t)}_{d \times 1} \otimes \underbrace{\mathbf{1}_{W_{(\nu)}(t)^{\top} o_{(\nu),p}(t) > 0}}_{m \times 1} \in \mathbb{R}^{md}, \\ \gamma_{(\nu),p}(t) &:= \frac{\omega \cdot \kappa}{\sqrt{m}} \underbrace{(\Lambda_{(\nu-1),i,\ell,*}(t) \otimes \Lambda_{(\nu-1),i}(t))^{\top}}_{d^2 \times L} \underbrace{(\text{diag}(\sigma_{(\nu),p}(t)) - \sigma_{(\nu),p}(t)\sigma_{(\nu),p}(t)^{\top})}_{L \times L} \\ &\quad \underbrace{\Lambda_{(\nu-1),i}(t)}_{L \times d} \sum_{r \in [m]} \underbrace{w_{(\nu),r}(t) \mathbb{I}\{o_{(\nu),p}(t)^{\top} w_{(\nu),r}(t) > 0\}}_{d \times 1} \in \mathbb{R}^{d^2}, \end{aligned}$$

where \otimes is the Kronecker product and $i = \lfloor p/L \rfloor$, $\ell = p \bmod L$. The indicator vector $\mathbf{1}_{W_{(\nu)}(t)^{\top} o_{(\nu),p}(t) > 0} \in \{0, 1\}^m$ where its r -th entry is $\mathbb{I}\{(W_{(\nu)}(t)^{\top} o_{(\nu),p}(t))_r > 0\}$ for $r \in [m]$. The layer-wise training dynamics are thereby shown as the following lemma:

Lemma 4.1 (Learning dynamics, informal version of Lemma C.6). *The learning dynamics of the multi-layer transformer Eq. (1) is given by:*

$$\mathbb{E}\left[\frac{d}{dt}\mathcal{L}(t, \mathbb{D})\right] = - \sum_{\nu \in [N]} \underbrace{\text{vec}\left(\frac{d}{d\mu_{(\nu)}(t)}\mathcal{L}(t, \mathbb{D})\right)^{\top}}_{1 \times nLd} \cdot \underbrace{(H_{(\nu)}(t) \otimes I_d)}_{nLd \times nLd} \cdot \underbrace{\text{vec}\left(\frac{d}{d\mu_{(\nu)}(t)}\mathcal{L}(t, \mathbb{D})\right)}_{nLd \times 1}$$

where $\mu_{(\nu)}(t)$ is a $nL \times d$ matrix, $\mu_{(\nu),p}(t)$ is the $(p \bmod L)$ -th row of ν -th layer output regarding to input matrix $X_{\lfloor p/L \rfloor}$ for any $p \in [nL]$ and $\nu \in [N]$.

Proof sketch of Lemma 4.1. Although the derivation of the gradient is complicated, the technique used for the proof is just the chain rule. We provide the complete proof in Appendix C. \square

Lemma 4.1 sums contributions from all layers $\nu \in [N]$, providing a powerful decomposition. This granular view suggests that each layer independently contributes to minimizing the loss based on its own kernel and local gradient. Different layers may learn at varying speeds or contribute differently to the overall task, depending on their respective kernel structures.

5 TRAINING CONVERGENCE AND APPROXIMATION GUARANTEES

In this section, we showcase the training convergence with an arbitrary error by limiting the kernel matrix $H'_{(\nu)}(0)$ at initialization to the Neural Tangent Kernel (NTK) (Jacot et al., 2018; Du et al., 2019) assumption. We state the formal assumption and the basic inductive proof in Section 5.1, and Section 5.2 demonstrates the results of the training convergence.

5.1 ASSUMPTIONS AND INDUCTIONS

Following the setting of Du et al. (2019) (see Assumption 3.1), we align this mild assumption in our analysis as shown:

Assumption 5.1. Defining $H'_{(\nu)}(0) \in \mathbb{R}^{nL \times nL}$ where its (i, j) -th entry $(\forall (i, j) \in [nL] \times [nL])$ is given by $H'_{(\nu),i,j}(0) := \langle \beta_{(\nu),i}(t), \beta_{(\nu),j}(t) \rangle$. For all $\nu \in [N]$, we assume $\frac{1}{\omega} H'_{(\nu)}(0)$ is positive definite (PD), formally, $\lambda_{(\nu)} := \lambda_{\min}(\frac{1}{\omega} H'_{(\nu)}(0)) > 0$.

Converging Kernel Perturbation and Kernel-Based Lazy Learning. Assuming $\lambda_{\min}(H'_{(\nu)}(0)) > 0$, we extend this PD property to the case during the training and the case of $H_{(\nu)}(t)$ that requires bounded summational update of weight at an arbitrary positive time $t > 0$. This connects to the *Lazy Learning* regime (Jacot et al., 2018; Du et al., 2019), where the updates of weights are limited in some high-dimensional ball with radius R . We first give the *Good Properties* requirements for provable arbitrary convergence below:

Definition 5.2 (Good Properties and Good Model Class). We fix the dataset size n , then we say model $F(X, \theta)$ has *Good Properties* once it satisfies (constant $C > 0$):

$$1. \omega = o\left(\frac{1}{NL^2 d^{2.5} B^3}\right); \quad 2. \kappa = \frac{1}{\sqrt{m}}; \quad 3. m = \Omega\left(\frac{n^3 L^5 \exp(Cd)}{\varepsilon^6 \omega^6 \lambda^6 \delta^3 N^2}\right).$$

Good Model Class is $\mathcal{F}_{M,T,N}(\mathbb{D}) := \{F(\cdot, \theta(T)), \theta(0) \sim \mathcal{N}(0, I_M), \theta(T) \in \mathcal{A}_{T,N}(\theta(0), \mathbb{D})\}$ for any $\mathbb{D} \subset \mathcal{D}$ and $|\mathbb{D}| = N$.

Lemma 5.3 (Informal version of Lemma E.1). Assuming Assumption 5.1 and Definition 5.2 hold, denote the failure probability $\delta \in (0, 0.1)$, then the kernel perturbation bound is: $\Pr[\lambda_{\min}(H_{(\nu)}(t)) < \lambda/2] < \delta$. Therefore, bounding loss dynamics is given by ($C > 0$ is some constant): $\Pr\left[\mathbb{E}\left[\frac{d}{dt} \mathcal{L}(t, \mathbb{D})\right] > -C \cdot \omega \lambda N \cdot \mathcal{L}(t, \mathbb{D})\right] < \delta$.

Proof sketch of Lemma 5.3. The stability of the PD property is ensured by $\lambda_{\min}(H') \geq \lambda_{\min}(H) - \|H - H'\|_F$ (Fact B.11) since the considerably small perturbation $\|H - H'\|_F$. The upper bound on perturbation converges due to increasing feed-forward layer width m and decreasing weight perturbation radius R , leading to the *Lazy Learning* regime. The complete proof of this lemma is stated in Appendix E. \square

5.2 RESULTS: CONVERGENCE OF KERNEL REGIMES

Factors of Scaling Law. To begin with, we list several crucial factors of the scaling law below. The number of training time T is ensured by the dataset size and the minimum batch size $|\mathbb{B}|$ due to the one-pass SGD. In addition, the model size refers to the number of trainable parameters in our model, which could be trivially calculated, and the definition of the total compute holds as our compute analysis in Lemma F.1.

Definition 5.4. We define:

- *Model Size:* $M := O(N(md + d^2)) = O(Nmd)$.
- *Training Time:* $T := \frac{N}{|\mathbb{B}|}$.
- *Dataset Size:* $N := n$.
- *Total Compute:* $C := O(MN)$.

Training Convergence and Approximation. The convergence guarantee of implementing Eq (2) as the continuous-time optimization is stated below:

Theorem 5.5 (Informal version of Theorem F.2). *Let all scaling law factors be defined as Definition 5.4 and Assumption 5.1 and Definition 5.2 hold. Denote the failure probability $\delta \in (0, 0.1)$, $\alpha_{\text{approx}} = \text{poly}(\exp(d), L, n, \frac{1}{\delta})$. For the Good Model Class $\mathcal{F}_{M,T,N}(\mathbb{D})$, with a probability at least $1 - \delta$, we have: $\mathcal{L}(t, \mathbb{D}) = \mathbb{E}_{(X,Y) \sim \mathcal{D}}[\|F(X) - Y\|_F^2] \leq \exp(-\frac{\epsilon}{\alpha_{\text{approx}}} \text{MT}), \forall F \in \mathcal{F}_{M,T,N}(\mathbb{D})$.*

Proof Sketch of Theorem 5.5. Firstly, we confirm the connection between hidden-state gradient norm and the training objective in Part 15 of Lemma D.1, the model convergence therefore exponentially benefits from the model size (neural depth and width) and the training time. Besides, the variance produced by the stochastic algorithm is provably reduced with a considerably large m . See Appendix F.2 for the detailed proofs. \square

Corollary 5.6 (Informal version of Corollary F.3). *Assuming we have arbitrary dataset size $N \in (0, +\infty)$. Let all scaling law factors be defined as Definition 5.4 and Assumption 5.1 and Definition 5.2 hold. Denote the failure probability $\delta \in (0, 0.1)$. We define $\alpha = O(Ld\sqrt{\log(1/\delta)})$. For the Good Model Class $\mathcal{F}_{M,T,N}(\mathbb{D}')$ with some $\mathbb{D}' \subset \mathcal{X} \times \mathcal{Y}$ and any function $F' : \mathcal{X} \rightarrow \mathcal{Y}$, arbitrary error $\epsilon > 0$ and compute cost C , with a probability at least $1 - \delta$, we have:*

$$\inf_{F \in \mathcal{F}_{M,T,N}(\mathbb{D}'), \mathbb{D}' \subset \mathcal{X} \times \mathcal{Y}} \mathbb{E}_{(X,Y) \sim \mathcal{D}}[\|F(X) - F'(X)\|_F^2] \leq \epsilon,$$

where $C = O(MN) = \Omega(256\alpha^8\epsilon^{-8} \wedge \epsilon^{-8}\omega^4 \log(2\epsilon^{-1})^4)$, $M = \Omega(N^3)$, $T = N$.

Proof Sketch of Corollary 5.6. We let $\mathbb{D}' := \{(X_i, F'(X_i)), X_i \sim \mathcal{X}\}_{i=1}^N$, then $\mathcal{F}_{M,T,N}(\mathbb{D}') = \{F(\cdot, \theta(T)), \theta(0) \sim \mathcal{N}(0, I_M), \theta(T) \in \mathcal{A}_{T,N}(\theta(0), \mathbb{D}')\}$, where we define $\mathcal{A}_{T,N}(\theta(0), \mathbb{D}') := \{\theta(0) + \int_0^T -\frac{d}{ds}\mathcal{L}(s, \mathbb{B}(s))ds, \mathbb{B}(s) \subseteq \mathbb{D}'\}$. Thus, we are able to obtain the same bound in Theorem 5.5 of the model to approximate the optimal mapping F^* . We combine the Hoeffding inequality with the convergence rate to obtain the results trivially. The formal proof is in Appendix F.3. \square

6 SCALING LAW

This section formally analyzes the generalization error bound of transformer-based language models. In particular, Section 6.1 showcases the general scaling law, which describes how the generalization error bound converges with the growing total computational cost. Moreover, since Theorem 6.1 demonstrates the three-stage upper bound on the in-distribution expected risk, in Section 6.2, we analyze how each variable affects the generalization in different scaling phases.

6.1 GENERAL PRETRAINING SCALING LAW

Excess Risk with Optimal Dataset Size. We state the upper bound on excess risk $\Delta\mathcal{R}(F)$ below:

Theorem 6.1 (Informal version of Theorem G.4). *Let all pre-conditions hold as Corollary 5.6. Let arbitrary grokking coefficient $\epsilon \in (0, 1)$, we provide the main criteria for determining scaling phase,*

$$\epsilon \geq \mathcal{T}(C) := \sqrt{\frac{\log(C^{\frac{1}{8}}/\alpha)}{C^{\frac{1}{4}}\omega}}. \quad (3)$$

Hence, with a probability at least $1 - \delta$, there exists:

$$\inf_{\mathcal{F}_{M,T,N}(\mathbb{D})} \sup_{\mathbb{D} \in \mathcal{D}} \Delta\mathcal{R}(F) \leq R(C) := \begin{cases} O(\frac{\alpha^2}{\omega C^{\frac{1}{8}}}) + O(\frac{\epsilon \cdot d \cdot \log(1/\sqrt{\epsilon})}{C^{\frac{1}{4}}}) + \epsilon^{\frac{3}{2}}, & \text{Eq. (3) holds} \\ \exp(-\epsilon^2\omega C^{\frac{1}{4}}) + O(\frac{\epsilon \cdot d \cdot \log(1/\sqrt{\epsilon})}{C^{\frac{1}{4}}}) + \epsilon^{\frac{3}{2}}, & \text{otherwise} \end{cases},$$

note that $M = \Omega(N^3)$, $T = N$, $C = O(MN) = \Omega(N^4)$.

Proof Sketch of Theorem 6.1. This proof utilizes the results in Lemma 11 of Schmidt-Hieber (2017) for the generalization error and Corollary 5.6 for the approximation error. The full proof is shown in Appendix G.2. \square

Lower Bound. Besides, we also give the lower bound on the excess risk as follows:

Theorem 6.2 (Informal version of Theorem G.5). *Let all pre-conditions hold as Theorem 6.1. For any choice of ϵ , we let $C = \Omega(\epsilon^{-12})$ to offset the negative effect of the grokking coefficient ϵ , thus,*

$$\Pr \left[\inf_{\mathcal{F}_{M,T,N}(\mathbb{D})} \sup_{\mathbb{D} \in \mathcal{D}} \Delta\mathcal{R}(F) \asymp O\left(\frac{\alpha}{C^{\frac{1}{8}}}\right) \right] \geq 1 - \delta.$$

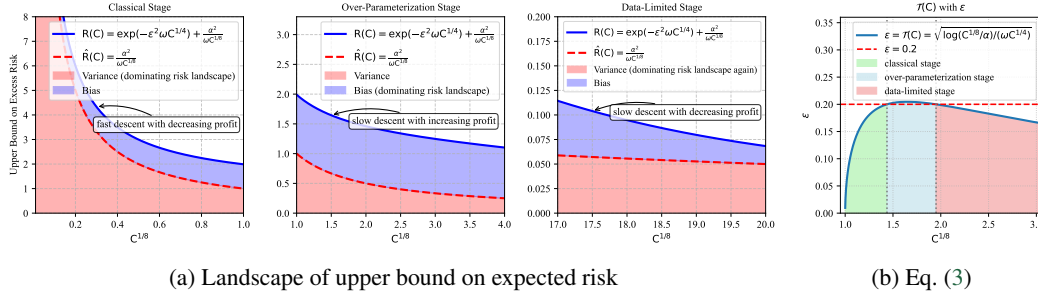


Figure 1: (a) Visualization of the upper bound in Theorem 6.1, comparing generalization error, convergence speed, and scaling profit (defined in Section 6.2) across three scaling phases: *Classical stage* (left): Generalization error decreases rapidly, but scaling profit diminishes. *Over-parameterization stage* (center): Scaling profit increases, reflecting improved generalization efficiency per unit compute; the bias term dominates the error bound. *Data-limited stage* (right): The variance term dominates, leading to slow error reduction and declining scaling profit. (b): Visualization of Eq. (3), illustrating how compute cost C and the grokking coefficient ε influence the three-stage generalization behavior.

Proof Sketch of Theorem 6.2. Since we have no additional assumption on the data distribution, we could only obtain the lower bound using some trivial concentration inequalities (please refer to Appendix G.3 for the formal derivation). \square

6.2 ANALYSIS ON THREE-STAGE BOUND

Discussion of Fix ε situation (Three-Stage Bound). We define three key quantities for our analysis: (i) the generalization error bound $R(C)$ from Theorem 6.1, (ii) the generalization speed $\frac{d}{dC} R(C)$, and (iii) the scaling profit $\frac{d}{dC} (1/R(C))$, which captures generalization efficiency during scaling. For fixed $\varepsilon > 0$, the function $\mathcal{T}(C)$ increases on $[\alpha^8, C^*]$, where C^* satisfies $\frac{d\mathcal{T}(C)}{dC} = 0$, then decreases monotonically to zero. The equation $\varepsilon = \mathcal{T}(C)$ thus has two solutions (Figure 1.b), partitioning the scaling process into three distinct stages: classical, over-parameterization, and data-limited regimes.

In the initial *classical stage*, the term $O(\alpha^2/(\omega C^{1/8}))$ dominates the upper bound. Insufficient neurons prevent perfect data fitting, causing the bound to initially converge rapidly before slowing due to increasing model complexity. Subsequently, in the *over-parameterization stage*, the major term in the bound of Theorem 6.1 is $\exp(-\varepsilon^2 \omega C^{1/4})$. Here, the model becomes over-parameterized to achieve a stable, low-variance solution, where the risk decreases exponentially with compute cost (model size and training duration) and bias dominates. Finally, in the *data-limited stage*, variance again governs the scaling rate owing to limited data samples—the model can easily memorize training data but requires more samples to generalize across the distribution. For more formal analysis, we encourage readers to refer to Appendix A.

Transition from classical stage to over-parameterization stage: Grokking. As shown in the middle of Figure 1.a, the training immediately accelerates to generalize, which perfectly matches the grokking phenomenon (Power et al., 2022; Nanda et al., 2023; Liu et al., 2022). Results in prior works (Kumar, 2024; Chizat et al., 2019; Lyu et al., 2023; Rubin et al., 2023) report that large initialization provably induces the grokking phenomenon, while in this work, it’s equivalent to a large output-scaling coefficient ε . Therefore, we follow (b) in Figure 1 to provide two critical insights below, and we further evaluate its effectiveness in Section 7.

- Insight 1: Lower ε (or smaller initialization) induces later beginning of grokking.
- Insight 2: Lower ε (or smaller initialization) induces slower generalization speed (or smaller convergence rate of generalization) at the grokking stage.

7 NUMERICAL EVALUATIONS

In this section, we conduct experiments to evaluate our scaling theory and the analysis in Section 6.

Setups. We fine-tune the pretrained ViT (Vision Transformer, Dosovitskiy et al. (2020)) on classical image classification task, including Cifar-10, 100 (Krizhevsky et al., 2009), and MNIST (LeCun, 1998) datasets, to compare their loss variety with AdamW optimizer (Kingma & Ba, 2014).

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

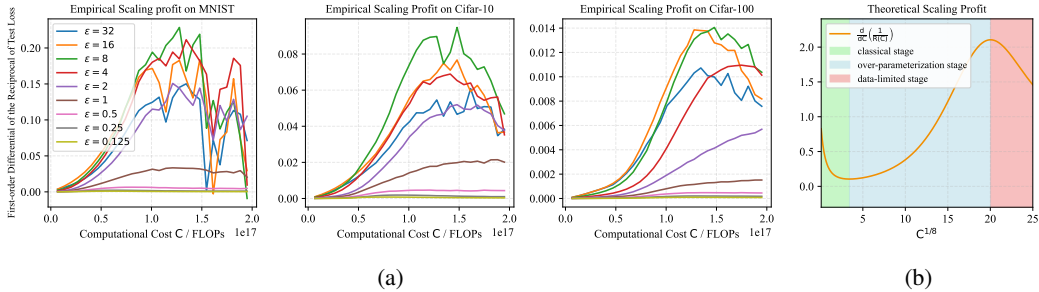


Figure 2: (a): First-order differential of the reciprocal of test error with different ϵ and C on MNIST, Cifar-10, 100 datasets. (b): Visualization of the theoretical scaling profit (see Appendix A).

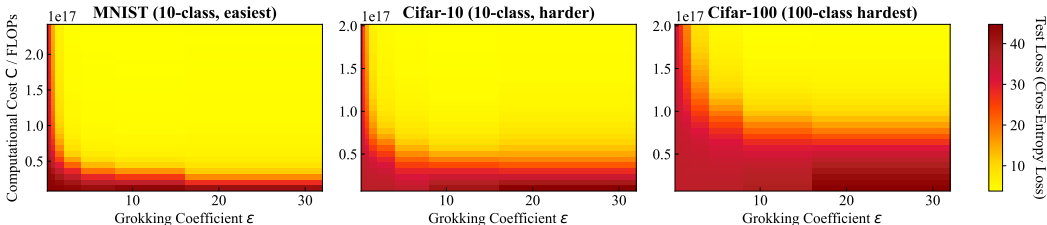


Figure 3: Results of the experiment in Section 7: Heatmaps of test loss on MNIST, Cifar-10, 100 datasets with different computational cost C and grokking coefficient ϵ .

For each dataset, we sample only 1/6 of the entire training set, taking 30 epochs of training to record the training and test cross-entropy loss curves. The value of ϵ is iteratively chosen from $\{32, 16, 8, 4, 2, 1, 0.5, 0.25, 0.125\}$. The computational cost C (computed by the total FLOPs) only depends on the training duration T , which is the number of total training steps in our experiments. We repeat the whole experiment 10 times for stable and fair results.

Results of Three-Phase Scaling Law Validation. We compute the empirical scaling profit of ViT on each dataset as the first-order derivative of the reciprocal test error with computational cost C (Figure 2.1.a). The empirical curves align closely with the theoretical scaling profit within the range of $[5, 25]$ (Figure 2.b), matching and approximating the ideal shape, which further provides strong validation for our three-phase scaling law theory.

Results of Grokking Validation. As shown in Figure 3, as ϵ increases, the transition region (marked by red-yellow boundaries indicative of grokking) occurs at smaller computational resource C , confirming Insight 1 in Section 6.2. Specifically, for smaller ϵ (left side of each panel), these transitions emerge later or are absent during scaling. Moreover, the computational cost range associated with this transition narrows as ϵ grows, implying that the grokking speed improves more sharply with larger ϵ and more gradually with smaller ϵ , validating Insight 2. We also observe that larger ϵ corresponds to higher initial risk (evidenced by darker red regions in the bottom-right) and, as task difficulty increases across the heatmaps, the onset of grokking occurs progressively later. We leave the explanation for these interesting artifacts for future direction.

8 CONCLUSION

This work presents a comprehensive theoretical framework for rigorously analyzing the scaling law phenomenon in LLMs, specifically addressing the empirically observed power-law relationship between model performance and computational resources from a learning theory perspective. By formalizing the dynamics of training sequence-to-sequence multi-layer transformer architectures, we establish a foundational guarantee: under allocation of compute, the generalization error of these models converges asymptotically as computational budgets increase with rate $\Theta(C^{-1/8})$ (Theorem 6.2), where C is the computational cost. Furthermore, our analysis proposes a three-phase pretraining scaling theory, where the scaling process is divided by a sudden pattern transition and transition-back with improved upper bound $\exp(-C^{1/4})$. We match this regime with the grokking phenomenon, giving two insights and conduct experiments to validate its correctness.

ETHICS STATEMENT

This is a purely theoretical paper that studies the scaling theory of training a constructed multi-layer transformer-based language model on sequence-to-sequence data distribution. To the best of our knowledge, this work provides the first analysis of the scaling law of training real-world LLMs from a learning theory perspective. As this work is theoretical and focuses on the capability of deep learning models, we don't foresee direct negative societal impacts and ethics concerns. We follow the ICLR Code of Ethics and affirm that all aspects of this research comply with the principles of fairness, transparency, and integrity.

In particular, we provide our clarification of the use of Large Language Models (LLMs) in Appendix H.

REPRODUCIBILITY STATEMENT

Theoretical Reproducibility. For our theoretical part, the setting is clearly stated in Section 3.2, and we also provide a more technical version, including necessary facts and lemmas, in Appendix B. The only assumption of this paper is Assumption 5.1, and two vital conditions is Definition 5.2 and Definition 5.4.

For all Lemmas/Theorems/Corollaries in the main paper, we have:

- Proof sketch of Lemma 4.1 is provided under itself. The formal version and full proof is in Lemma C.6.
- Proof sketch of Lemma 5.3 is provided under itself. The formal version and full proof is in Lemma E.1.
- Proof sketch of Theorem 5.5 is provided under itself. The formal version and full proof is in Theorem F.2.
- Proof sketch of Corollary 5.6 is provided under itself. The formal version and full proof is in Corollary F.3.
- Proof sketch of Theorem 6.1 is provided under itself. The formal version and full proof is in Theorem G.4.
- Proof sketch of Theorem 6.2 is provided under itself. The formal version and full proof is in Theorem G.5.

Experimental Reproducibility. We provide the code for the experiments in Section 7 and a markdown file "README.md" for explaining how to reproduce our experimental results.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Raghav Addanki, Chenyang Li, Zhao Song, and Chiwun Yang. One pass streaming algorithm for super long token attention approximation in sublinear space. *arXiv preprint arXiv:2311.14652*, 2023.
- Subutai Ahmad and Gerald Tesauro. Scaling and generalization in neural networks: a case study. *Advances in neural information processing systems*, 1, 1988.
- Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk. The recurrent neural tangent kernel. *arXiv preprint arXiv:2006.10246*, 2020.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar. *arXiv preprint arXiv:2305.13673*, 2023a.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.2, knowledge manipulation. *arXiv preprint arXiv:2309.14402*, 2023b.

- 540 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and
541 extraction. *arXiv preprint arXiv:2309.14316*, 2023c.
- 542
- 543 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling
544 laws. *arXiv preprint arXiv:2404.05405*, 2024.
- 545
- 546 Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-
547 parameterization. In *International conference on machine learning*, pp. 242–252. PMLR, 2019.
- 548
- 549 Josh Alman and Zhao Song. Fast attention requires bounded entries. *Advances in Neural Information
550 Processing Systems*, 36, 2023.
- 551
- 552 Josh Alman and Zhao Song. The fine-grained complexity of gradient computation for training large
553 language models. *arXiv preprint arXiv:2402.04497*, 2024a.
- 554
- 555 Josh Alman and Zhao Song. How to capture higher-order correlations? generalizing matrix soft-
556 max attention to kronecker computation. In *The Twelfth International Conference on Learning
557 Representations*, 2024b.
- 558
- 559 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On
560 exact computation with an infinitely wide neural net. *Advances in neural information processing
561 systems*, 32, 2019.
- 562
- 563 Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The
564 silent alignment effect. *arXiv preprint arXiv:2111.00034*, 2021.
- 565
- 566 Alexander Atanasov, Jacob A Zavatone-Veth, and Cengiz Pehlevan. Scaling and renormalization in
567 high-dimensional regression. *arXiv preprint arXiv:2405.00592*, 2024.
- 568
- 569 Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural
570 scaling laws. *Proceedings of the National Academy of Sciences*, 121(27):e2311878121, 2024.
- 571
- 572 Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding,
573 Kai Dong, Qiusi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with
574 longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- 575
- 576 Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. *Advances in Neural
577 Information Processing Systems*, 32, 2019.
- 578
- 579 Blake Bordelon and Cengiz Pehlevan. Learning curves for sgd on structured features. *arXiv preprint
580 arXiv:2106.02713*, 2021.
- 581
- 582 Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in
583 kernel regression and wide neural networks. In *International Conference on Machine Learning*,
584 pp. 1024–1034. PMLR, 2020.
- 585
- 586 Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. A dynamical model of neural scaling
587 laws. *arXiv preprint arXiv:2402.01092*, 2024.
- 588
- 589 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
590 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
591 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 592
- 593 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,
Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:
Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm.
Foundations of Computational Mathematics, 7:331–368, 2007.
- ChatGPT. Optimizing language models for dialogue. *OpenAI Blog*, November 2022. URL <https://openai.com/blog/chatgpt/>.

- 594 Yifang Chen, Xuyang Guo, Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Scaling
595 law phenomena across regression paradigms: Multiple and kernel approaches. *arXiv preprint*
596 *arXiv:2503.01314*, 2025.
- 597 Lенаic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks
598 trained with the logistic loss. In *Conference on learning theory*, pp. 1305–1338. PMLR, 2020.
- 600 Lенаic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming.
601 *Advances in neural information processing systems*, 32, 2019.
- 602 Majid Daliri, Zhao Song, and Chiwun Yang. Unlocking the theory behind scaling 1-bit neural
603 networks. *arXiv preprint arXiv:2411.01663*, 2024.
- 605 Yichuan Deng, Zhihang Li, and Zhao Song. Attention scheme inspired softmax regression. *arXiv*
606 *preprint arXiv:2304.10411*, 2023.
- 607 Yichuan Deng, Zhao Song, and Chiwun Yang. Attention is naturally sparse with gaussian distributed
608 input. *arXiv preprint arXiv:2404.02690*, 2024.
- 610 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
611 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
612 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
613 *arXiv:2010.11929*, 2020.
- 614 Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
615 over-parameterized neural networks. In *ICLR*. *arXiv preprint arXiv:1810.02054*, 2019.
- 617 Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. A primer on the inner
618 workings of transformer-based language models. *arXiv preprint arXiv:2405.00208*, 2024.
- 619 Yeqi Gao, Sridhar Mahadevan, and Zhao Song. An over-parameterized exponential regression. *arXiv*
620 *preprint arXiv:2303.16504*, 2023.
- 622 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
623 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of
624 models. *arXiv preprint arXiv:2407.21783*, 2024.
- 625 Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi,
626 Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the
627 rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- 628 Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David Woodruff, and Amir Zandieh.
629 Hyperattention: Long-context attention in near-linear time. In *The Twelfth International Confer-*
630 *ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Eh0Od2BJIM)
631 [Eh0Od2BJIM](https://openreview.net/forum?id=Eh0Od2BJIM).
- 633 Satoshi Hayakawa. Note on the generalization bounds of the empirical risk minimizer, 2019. URL
634 https://ibis.t.u-tokyo.ac.jp/suzuki/paper/note_ERM_genererror.pdf.
- 635 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
636 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*
637 *arXiv:2009.03300*, 2020.
- 638 Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad,
639 Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable,
640 empirically. *arXiv preprint arXiv:1712.00409*, 2017.
- 642 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
643 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
644 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 645 Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and ntk
646 for deep attention networks. In *International Conference on Machine Learning*, pp. 4376–4386.
647 PMLR, 2020.

- 648 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
649 generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
650
- 651 Prateek Jain, Sham M Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. Parallelizing
652 stochastic gradient descent for least squares regression: mini-batching, averaging, and model
653 misspecification. *Journal of machine learning research*, 18(223):1–42, 2018.
- 654 Praneeth Kacham, Vahab Mirrokni, and Peilin Zhong. Polysketchformer: Fast transformers via
655 sketches for polynomial kernels. *arXiv preprint arXiv:2310.01655*, 2023.
656
- 657 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
658 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
659 *arXiv preprint arXiv:2001.08361*, 2020.
- 660 Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva Reddy.
661 The impact of positional encoding on length generalization in transformers. *Advances in Neural
662 Information Processing Systems*, 36:24892–24928, 2023.
663
- 664 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
665 arXiv:1412.6980*, 2014.
- 666 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Tech
667 report*, 2009.
668
- 669 Tanishq Kumar. Grokking as the transition from lazy to rich training dynamics. *The Twelfth
670 International Conference on Learning Representations*, 2024.
- 671 Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
672
- 673 Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and
674 Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances in
675 Neural Information Processing Systems*, 33:15156–15172, 2020.
- 676 Chenyang Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Exploring the frontiers of softmax: Provable
677 optimization, applications in diffusion model, and beyond. *arXiv preprint arXiv:2405.03251*,
678 2024.
679
- 680 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou,
681 Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with
682 you! *arXiv preprint arXiv:2305.06161*, 2023.
683
- 684 Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient
685 descent on structured data. *Advances in neural information processing systems*, 31, 2018.
- 686 Licong Lin, Jingfeng Wu, Sham M Kakade, Peter L Bartlett, and Jason D Lee. Scaling laws in linear
687 regression: Compute, parameters, and data. *arXiv preprint arXiv:2406.08466*, 2024.
688
- 689 Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong
690 Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-
691 experts language model. *arXiv preprint arXiv:2405.04434*, 2024a.
- 692 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
693 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint
694 arXiv:2412.19437*, 2024b.
- 695 Ziming Liu, Eric J Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic data.
696 *arXiv preprint arXiv:2210.01117*, 2022.
697
- 698 Bochen Lyu, Di Wang, and Zhanxing Zhu. A solvable attention for neural scaling laws. In *The
699 Thirteenth International Conference on Learning Representations*, 2025.
- 700 Kaifeng Lyu, Jikai Jin, Zhiyuan Li, Simon S Du, Jason D Lee, and Wei Hu. Dichotomy of early and
701 late phase implicit biases can provably induce grokking. *arXiv preprint arXiv:2311.18817*, 2023.

- 702 Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt:
703 Towards detailed video understanding via large vision and language models. *arXiv preprint*
704 *arXiv:2306.05424*, 2023.
- 705
706 Sadhika Malladi, Alexander Wettig, Dingli Yu, Danqi Chen, and Sanjeev Arora. A kernel-based
707 view of language model fine-tuning. In *International Conference on Machine Learning*, pp.
708 23610–23641. PMLR, 2023.
- 709
710 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures
711 for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- 712
713 Druv Pai, Sam Buchanan, Ziyang Wu, Yaodong Yu, and Yi Ma. Masked completion via structured
714 diffusion with white-box transformers. In *The Twelfth International Conference on Learning*
Representations, 2024.
- 715
716 Elliot Paquette, Courtney Paquette, Lechao Xiao, and Jeffrey Pennington. 4+ 3 phases of compute-
717 optimal neural scaling laws. *arXiv preprint arXiv:2405.15074*, 2024.
- 718
719 Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Gen-
720 eralization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*,
721 2022.
- 722
723 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
724 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 725
726 Jonathan S Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive prediction
727 of the generalization error across scales. *arXiv preprint arXiv:1909.12673*, 2019.
- 728
729 Jonathan S Rosenfeld, Jonathan Frankle, Michael Carbin, and Nir Shavit. On the predictability of
730 pruning across scales. In *International Conference on Machine Learning*, pp. 9075–9083. PMLR,
731 2021.
- 732
733 Noa Rubin, Inbar Seroussi, and Zohar Ringel. Grokking as a first order phase transition in two layer
734 networks. *arXiv preprint arXiv:2310.03789*, 2023.
- 735
736 Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation
737 function. *arXiv preprint arXiv:1708.06633*, 2017.
- 738
739 Mariia Seleznova and Gitta Kutyniok. Neural tangent kernel beyond the infinite-width limit: Effects
740 of depth and initialization. In *International Conference on Machine Learning*, pp. 19522–19560.
741 PMLR, 2022.
- 742
743 Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in neural
744 networks: Emergence from inputs and advantage over fixed features. In *International Conference*
on Learning Representations, 2021.
- 745
746 Zhenmei Shi, Yifei Ming, Xuan-Phi Nguyen, Yingyu Liang, and Shafiq Joty. Discovering the gems
747 in early layers: Accelerating long-context llms with 1000x input token reduction. *arXiv preprint*
arXiv:2409.17422, 2024a.
- 748
749 Zhenmei Shi, Junyi Wei, and Yingyu Liang. Provable guarantees for neural networks via gradient
750 feature learning. *Advances in Neural Information Processing Systems*, 36, 2024b.
- 751
752 Zhao Song and Xin Yang. Quadratic suffices for over-parametrization via matrix chernoff bound.
753 *arXiv preprint arXiv:1906.03593*, 2019.
- 754
755 Ingo Steinwart, Don R Hush, Clint Scovel, et al. Optimal rates for regularized least squares regression.
In *COLT*, pp. 79–93, 2009.
- 756
757 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
758 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
759 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

- 756 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
757 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
758 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 759 Aad W Van Der Vaart and Jon A Wellner. Weak convergence. In *Weak convergence and empirical*
760 *processes: with applications to statistics*, pp. 16–28. Springer, 1996.
- 762 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
763 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
764 *systems*, 30, 2017.
- 765 Alexander Wei, Wei Hu, and Jacob Steinhardt. More than a toy: Random matrix models predict how
766 real-world neural representations generalize. In *International conference on machine learning*, pp.
767 23549–23588. PMLR, 2022a.
- 769 Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and
770 optimization of neural nets vs their induced kernel. *Advances in Neural Information Processing*
771 *Systems*, 32, 2019.
- 772 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
773 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.
774 *arXiv preprint arXiv:2206.07682*, 2022b.
- 776 Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with
777 parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*, 2023.
- 778 Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of conver-
779 gence. *Annals of Statistics*, pp. 1564–1599, 1999.
- 780 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
781 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*
782 *Information Processing Systems*, 36, 2023.
- 784 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo
785 Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for
786 large language models. *arXiv preprint arXiv:2309.12284*, 2023a.
- 788 Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin
789 Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *Advances in Neural*
790 *Information Processing Systems*, 36:9422–9457, 2023b.
- 791 Yaodong Yu, Tianzhe Chu, Shengbang Tong, Ziyang Wu, Druv Pai, Sam Buchanan, and Yi Ma. Emer-
792 gence of segmentation with minimalistic white-box transformers. *arXiv preprint arXiv:2308.16271*,
793 2023c.
- 794 Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural*
795 *Information Processing Systems*, 32, 2019.
- 797 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
798 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
799 chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.
- 800 Mo Zhou, Rong Ge, and Chi Jin. A local convergence theory for mildly over-parameterized two-layer
801 neural network. In *Conference on Learning Theory*, pp. 4577–4632. PMLR, 2021.
- 802
803
804
805
806
807
808
809

Appendix

810		
811		
812		
813	CONTENTS	
814		
815	1 Introduction	1
816		
817	2 Related Work	2
818		
819		
820	3 Preliminary	3
821	3.1 Basic Notations	3
822	3.2 Setups	3
823		
824		
825	4 Learning Dynamics of Scaling Transformers	4
826	4.1 Simplifications	4
827	4.2 Key Derivation for Learning Dynamics	5
828		
829		
830	5 Training Convergence and Approximation Guarantees	6
831	5.1 Assumptions and Inductions	6
832	5.2 Results: Convergence of Kernel Regimes	6
833		
834		
835	6 Scaling Law	7
836	6.1 General Pretraining Scaling Law	7
837	6.2 Analysis on Three-Stage Bound	8
838		
839		
840	7 Numerical Evaluations	8
841		
842		
843	8 Conclusion	9
844		
845		
846	A Formal Analysis of Three-Stage Generalization	18
847		
848	B Technical Preliminary	20
849	B.1 Notations	20
850	B.2 Original Model Definitions	20
851	B.3 Basic Facts and Lemmas	21
852		
853		
854	C Gradient Computation and Learning Dynamics	22
855	C.1 Simplified Model Definitions	22
856	C.2 Gradient Computation	22
857	C.3 Learning Dynamics	23
858		
859		
860		
861	D Toolkit: Helpful Boundaries	25
862		
863	E Kernel Perturbation	30

864	F Convergence and Approximation Bound	34
865		
866	F.1 Complexity Analysis	34
867	F.2 Training Convergence	34
868	F.3 Approximation	35
869		
870		
871	G Scaling Law	36
872	G.1 Generalization Bound of Empirical Risk Minimizer	36
873	G.2 General Pretraining Scaling Law	37
874	G.3 Upper Bound	38
875		
876		
877	H LLMs Usage Disclosure	40
878		
879		
880		
881		
882		
883		
884		
885		
886		
887		
888		
889		
890		
891		
892		
893		
894		
895		
896		
897		
898		
899		
900		
901		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		

A FORMAL ANALYSIS OF THREE-STAGE GENERALIZATION

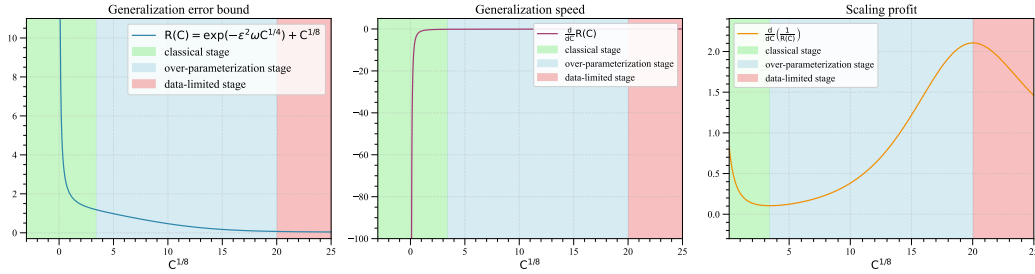


Figure 4: Visualization of Generalization error bound, Generalization speed, and Scaling profit. We use three different colors to distinguish three scaling stages.

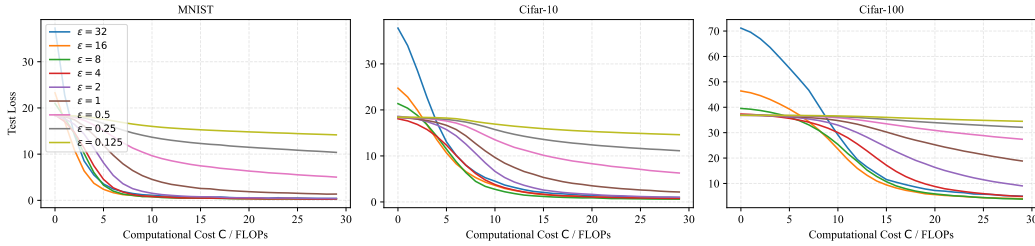


Figure 5: Results of the experiment in Section 7: Test loss curves on MNIST, Cifar-10, 100 datasets with different computational cost C and grokking coefficient ϵ .

We list the key concepts for our scaling law analysis as follows:

- The generalization error bound, denoted by $R(C)$, is a function of the computational cost C . Theorem 6.1 proves that $R(C)$ converges as C increases. This bound is visualized in the left panel of Figure 4.
- The generalization speed, defined as the derivative $R'(C) = \frac{d}{dC} R(C)$, measures the rate of convergence of the generalization error bound. The sign of $R'(C)$ indicates whether the bound is decreasing (converging) or increasing, and its magnitude represents the convergence rate. This is visualized in the middle panel of Figure 4.
- The scaling profit is defined as the derivative of the inverse of the generalization error bound, i.e., $\frac{d}{dC} \left(\frac{1}{R(C)} \right)$. A higher value of this metric indicates a greater benefit (profit) from increasing the computational cost, and conversely, a lower value suggests diminishing returns. This is shown in the right panel of Figure 4.

For the left and middle images in Figure 4, generalization error and speed decrease monotonically to zero as computational cost is scaled up. The right image displays a more complex scaling profit curve, where the computational cost range is divided into three phases by the curve’s maximum and minimum points.

The upper bound during the *classical stage* is principally determined by the term $O(\alpha^2 / (\omega C^{1/8}))$. The scarcity of neurons initially facilitates a phase of rapid convergence by constraining the model’s capacity to fit the data perfectly. This phase is transient, however, as the convergence rate is soon curtailed by the escalating complexity of the model.

The *over-parameterization stage* is marked by a transition in the scaling law. Here, the dominant term in the bound of Theorem 6.1 simplifies to $\exp(-\epsilon^2 \omega C^{1/4})$, indicating that the risk decreases exponentially with compute (model size and training duration). This stage corresponds to a stable, low-variance solution where bias dominates. We identify this phase with the grokking phenomenon

972 (Power et al., 2022). Our analysis shows that grokking is expected for any $\varepsilon < \max_{C>0} \mathcal{T}(C)$, a
973 relationship illustrated in Figure 1.b.

974
975 The scaling behavior enters a *data-limited stage* when the number of samples becomes the primary
976 constraint. In this regime, variance is the dominant factor because models can easily memorize the
977 small training set but fail to generalize. The scaling bound for this phase is therefore $O(\alpha^2/(\omega C^{1/8}))$.

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

B TECHNICAL PRELIMINARY

B.1 NOTATIONS

In this paper, we use integer d to denote the dimension of networks. We use L to denote the input length in language models. $\nabla_x f(x)$ and $\frac{df(x)}{dx}$ are both means to take the derivative of $f(x)$ with x . Let a vector $z \in \mathbb{R}^n$. We denote the ℓ_2 norm as $\|z\|_2 := (\sum_{i=1}^n z_i^2)^{1/2}$, the ℓ_1 norm as $\|z\|_1 := \sum_{i=1}^n |z_i|$, $\|z\|_0$ as the number of non-zero entries in z , $\|z\|_\infty$ as $\max_{i \in [n]} |z_i|$. We use z^\top to denote the transpose of a z . We use $\langle \cdot, \cdot \rangle$ to denote the inner product. Let $A \in \mathbb{R}^{n \times d}$, we use $\text{vec}(A)$ to denote a length nd vector. We denote the Frobenius norm as $\|A\|_F := (\sum_{i \in [n], j \in [d]} A_{i,j}^2)^{1/2}$. For any positive integer n , we use $[n]$ to denote set $\{1, 2, \dots, n\}$. We use $\mathbb{E}[\cdot]$ to denote the expectation. We use $\Pr[\cdot]$ to denote the probability. We use ϵ to denote the error. We define $\lambda_{\min}(\cdot)$ as a function that outputs the minimum eigenvalues of the input matrix, e.g. matrix $A \in \mathbb{R}^{n \times n}$ has eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$, $\lambda_{\min}(A) = \min\{\lambda_1, \lambda_2, \dots, \lambda_n\}$. For a vector $a \in \mathbb{R}^n$, we use a_i to denote its i -th entry for $i \in [n]$. For a matrix $A \in \mathbb{R}^{n \times d}$, vector $A_i \in \mathbb{R}^d$ is the i -th row for $i \in [n]$ and vector $A_{*,j} \in \mathbb{R}^n$ is the j -th column for $j \in [d]$. For a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$, we use f_i to denote the i -th entry of its output. For a function $f: \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$, we use $f_i \in \mathbb{R}^d$ to denote the i -th row of its output for $i \in [n]$ and us $f_{*,j} \in \mathbb{R}^n$ to denote the j -th column of its output for $j \in [d]$. We use $\mathbb{I}\{E_1, E_2, \dots, E_n\}$ to denote the indicator for event set $\{E_1, E_2, \dots, E_n\}$, only when E_1, E_2, \dots, E_n are all true, $\mathbb{I}\{E_1, E_2, \dots, E_n\} = 1$; otherwise, it equals to 0. For a vector $a \in \mathbb{R}^{d^2}$, function mat reshapes a to a $d \times d$ matrix, where its (i, j) -th entry is $\text{mat}_{i,j}(a) = a_{(i-1)d+j}$ for $(i, j) \in [d] \times [d]$. For a matrix $A \in \mathbb{R}^{n \times d}$, function vec flattens A to a nd -dimensional vector, where its i -th entry is $\text{vec}_i(A) = A_{\lfloor i/d \rfloor, i - \lfloor i/d \rfloor \cdot d}$ for $i \in [nd]$.

B.2 ORIGINAL MODEL DEFINITIONS

Definition B.1 (Weights and Initialization). *The weights of the model are denoted as $\theta(t)$ where $t \geq 0$ is time. It contains the weights of each layer $\theta(t) = \{\theta_{(\nu)}(t)\}_{\nu=1}^N$, and $\theta_{(\nu)} = \{U_{(\nu)}(t), W_{(\nu)}(t), A_{(\nu)}\}$. For each matrix:*

- Each entry of $U_{(\nu)}(0) \in \mathbb{R}^{d \times d}$ is initialized from the standard Gaussian distribution, formally, $U_{(\nu),k_1,k_2}(0) \sim \mathcal{N}(0, 1)$ for any $k_1, k_2 \in [d]$.
- Each entry of $W_{(\nu)}(0) \in \mathbb{R}^{d \times m}$ is initialized from the standard Gaussian distribution, formally, $w_{(\nu),r,k}(0) \sim \mathcal{N}(0, 1)$ for any $r \in [m], k \in [d]$, where $w_{(\nu),r}(0) \in \mathbb{R}^d$ is the r -th column of $W_{(\nu)}(0)$.
- Each entry of $A_{(\nu)}(0) \in \mathbb{R}^{m \times d}$ is initialized from a ± 1 uniform distribution, formally, $a_{(\nu),r,k}(0) \sim \text{Uniform}\{-1, +1\}$ for any $r \in [m], k \in [d]$, where $a_{(\nu),r}(0) \in \mathbb{R}^d$ is the r -th row of $A_{(\nu)}(0)$.

Definition B.2 (Data Distribution). *We denote $F^* : \mathcal{X} \rightarrow [C_1, C_2]^{L \times d}$ as the target function. The $\mathcal{X} \in \mathbb{R}^{L \times d}$ is a combination of L d -dimensional balls where $\|X_\ell\|_2^2 = \Theta(1), \forall X \in \mathcal{X}, \ell \in [L]$. $C_1, C_2 \in \mathbb{R}$ are two fixed constants. The distribution: $\mathcal{D} = \{(X, F^*(X) + \Xi), \Xi \in \mathbb{R}^{L \times d} \text{ is some random noise}\} \subset \mathcal{X} \times \mathcal{Y}$. The random noise $\Xi \in \mathbb{R}^{L \times d}$ is centered by $\mathbf{0}_{L \times d}$.*

Definition B.3 (Original Dataset). *The original dataset is $\mathbb{D} = \{(X_i, Y_i)\}_{i=1}^n \subset \mathcal{D}$, where $X_i, Y_i = F^*(X_i) + \Xi_i \in \mathbb{R}^{L \times d}$. For each data point $(X_i, Y_i), \forall i \in [n]$, it holds that:*

- $\|X_{i,\ell}\|_2 = \Theta(1)$ for $\ell \in [L]$.

Definition B.4 (Model Functions). *Given an input matrix $X \in \mathbb{R}^{L \times d}$, the model function is given by ($\epsilon > 0$ is the grokking coefficient):*

$$F(X, \theta(t)) := \epsilon \cdot F_{(N)}(F_{(N-1)}(\dots F_{(2)}(F_{(1)}(X + E, \theta(t)), \theta(t)) \dots), \theta(t)).$$

We list the original definition of each function as follows:

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

•

$$F_{(\nu)}(X, \theta(t)) := \frac{\omega}{\sqrt{m}} \text{ReLU}(\text{Softmax}(\kappa \cdot XU_{(\nu)}(t)X^\top + M) XW_{(\nu)}(t)) A_{(\nu)}$$

• $\text{Softmax}(A) := \text{diag}(\exp(A)\mathbf{1}_L)^{-1} \cdot \exp(A) \in \mathbb{R}^{L \times L}$ for $A \in \mathbb{R}^{L \times L}$.• $\text{ReLU}_{\ell,k}(X) := \max\{X_{\ell,k}, 0\}$ for any $\ell \in [L], k \in [d], X \in \mathbb{R}^{L \times d}$.• $M_{\ell_1, \ell_2} := \begin{cases} 0, & \ell_1 \geq \ell_2 \\ -\infty. & \ell_1 < \ell_2 \end{cases}, \forall \ell_1, \ell_2 \in [L]$.**Definition B.5.** We denote the special notations:• We denote $w_{(\nu),r}(t) \in \mathbb{R}^d$ as the r -th column of $W_{(\nu)} \in \mathbb{R}^{d \times m}$ for $r \in [m], \nu \in [N]$.• We denote $a_{(\nu),r} \in \mathbb{R}^d$ as the r -th row of $A_{(\nu)} \in \mathbb{R}^{m \times d}$ for $r \in [m], \nu \in [N]$.**Definition B.6.** We define the training objective:

$$\mathcal{L}(t, \mathbb{D}) := \mathbb{E}_{(X,Y) \sim \mathbb{D}} [\|F(X, \theta(t)) - Y\|_F^2]$$

B.3 BASIC FACTS AND LEMMAS

Fact B.7. For a variable $x \sim \mathcal{N}(0, \sigma^2)$, then with probability at least $1 - \delta$, we have:

$$|x| \leq C\sigma\sqrt{\log(1/\delta)}$$

Fact B.8. For an 1-Lipschitz function $f(\cdot)$, we have:

$$|f(x) - f(y)| \leq |x - y|, \forall x, y \in \mathbb{R}^d$$

Fact B.9. For a Gaussian variable $x \sim \mathcal{N}(0, \sigma^2 \cdot I_d)$ where $\sigma \in \mathbb{R}$, then for any $t > 0$, we have:

$$\Pr[x \leq t] \leq \frac{2t}{\sqrt{2\pi}\sigma}$$

Fact B.10. For a Gaussian vector $w \sim \mathcal{N}(0, \sigma^2 \cdot I_d)$ where $\sigma \in \mathbb{R}$, and a fixed vector $x \in \mathbb{R}^d$, we have:

$$w^\top x \sim \mathcal{N}(0, \sigma^2 \|x\|_2^2)$$

Fact B.11. For two matrices $H, \tilde{H} \in \mathbb{R}^{n \times n}$, we have:

$$\lambda_{\min}(\tilde{H}) \geq \lambda_{\min}(H) - \|\tilde{H} - H\|_F$$

Fact B.12. The Lipschitz constant of the softmax function is bounded by $O(1)$, such that:

$$\|\langle \exp(x), \mathbf{1}_L \rangle^{-1} \exp(x) - \langle \exp(y), \mathbf{1}_L \rangle^{-1} \exp(y)\|_2 \leq O(1) \cdot \|x - y\|_2, \forall x, y \in \mathbb{R}^L.$$

Fact B.13. For a matrix $H \in \mathbb{R}^{n \times n}$, there is $\lambda_{\min}(H \otimes I_d) = \lambda_{\min}(H)$.

In addition, we state four vital lemmas of concentration inequalities for simplifying analysis:

Lemma B.14 (Hoeffding bound). Let X_1, \dots, X_n denote n independent bounded variables in $[a_i, b_i]$ for $a_i, b_i \in \mathbb{R}$. Let $X := \sum_{i=1}^n X_i$, then we have

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

Lemma B.15 (Markov's inequality). If X is a non-negative random variable and $a > 0$, then the probability that X is at least a is at most the expectation of X divided by a :

$$\Pr[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

Lemma B.16 (Chernoff bound). Let $X = \sum_{i=1}^n X_i$, where $X_i = 1$ with probability p_i and $X_i = 0$ with probability $1 - p_i$, and all X_i are independent. Let $\mu = \mathbb{E}[X] = \sum_{i=1}^n p_i$. Then• $\Pr[X \geq (1 + \delta)\mu] \leq \exp(-\delta^2\mu/3), \forall \delta > 0;$ • $\Pr[X \leq (1 - \delta)\mu] \leq \exp(-\delta^2\mu/1), \forall 0 < \delta < 1.$

C GRADIENT COMPUTATION AND LEARNING DYNAMICS

C.1 SIMPLIFIED MODEL DEFINITIONS

Definition C.1 (Rearranged Dataset). *Given the origin dataset $\mathbb{D} = \{X_i, Y_i\}_{i=1}^n \subset \mathbb{R}^{L \times d} \times \mathbb{R}^{L \times d}$. The rearranged dataset is $\mathbb{D}_{\text{rearrange}} = \{(X_{i, \leq \ell} + E_{\leq \ell}, Y_{i, \ell})\}_{(i, \ell)=(1, 1)}^{(n, L)}$, where $X_{i, \leq \ell} + E_{\leq \ell} \in \mathbb{R}^{\ell \times d}$ and $Y_{i, \ell} \in \mathbb{R}^d$.*

Definition C.2 (Simplified Model Function). *Given the origin dataset $\mathbb{D} = \{X_i, Y_i\}_{i=1}^n \subset \mathbb{R}^{L \times d} \times \mathbb{R}^{L \times d}$, we define the compact form of the model function as:*

$$F(t) \in \mathbb{R}^{nL \times d}, \text{ where its } i\text{-th row is given by } F_p(t) := F_\ell(X_i, \theta(t)) \in \mathbb{R}^d, p \in [nL],$$

$$Y \in \mathbb{R}^{nL \times d}, \text{ where its } i\text{-th row is given by } Y_p := Y_{i, \ell} \in \mathbb{R}^d, p \in [nL].$$

Here, $i = \lfloor p/L \rfloor$ and $\ell = p \bmod L$.

We list the notation-simplified definitions of all functions as follows:

- (Hidden State) $\Lambda_{(\nu), i}(t) := F_{(\nu)}(\Lambda_{(\nu-1), i}(t), \theta(t)) \in \mathbb{R}^{L \times d}$ for $\nu \in [N]$, $\Lambda_{(0), i}(t) = X_i + E$.
- (Attention Scores) $\sigma_{(\nu), (i-1)L+\ell}(X) = \text{Softmax}_\ell(\Lambda_{(\nu), i}(t)U_{(\nu)}(t)\Lambda_{(\nu), i}(t)^\top + M) \in \mathbb{R}^L$.
- (Attention Output) $o_{(\nu), (i-1)L+\ell}(t) := \Lambda_{(\nu-1), i}(t)^\top \cdot \sigma_{(\nu), (i-1)L+\ell}(t) \in \mathbb{R}^d$.
- (ℓ -th Token of Hidden State) $\mu_{(\nu), (i-1)L+\ell}(t) := \frac{\omega}{\sqrt{m}} \sum_{r=1}^m a_{(\nu), r} \cdot \phi(\langle o_{(\nu), (i-1)L+\ell}(t), w_{(\nu), r}(t) \rangle) \in \mathbb{R}^d$, where $\phi(x) := \max\{0, x\}, \forall x \in \mathbb{R}$.
 $\mu_{(0), (i-1)L+\ell}(t) = X_{i, \ell} + E_\ell$.
- (Model Output) $F_{(i-1)L+\ell}(t) = \varepsilon \cdot \sum_{\nu=0}^N \mu_{(\nu), (i-1)L+\ell}(t) \in \mathbb{R}^d$.

Lemma C.3. *We have:*

$$\mathcal{L}(t, \mathbb{D}) = \frac{1}{n} \|F(t) - Y\|_F^2$$

Proof. Since the decoder-only property of the model function, we have:

$$F_{(i-1)L+\ell}(t) = F_\ell(X_i, \theta(t)), \forall (X_i, Y_i) \in \mathbb{D}, i \in [n], \ell \in [L].$$

We then separate each token vector with:

$$\begin{aligned} \mathcal{L}(t, \mathbb{D}) &= \mathbb{E}_{(X, Y) \sim \mathbb{D}} [\|F(X, \theta(t)) - Y\|_F^2] \\ &= \frac{1}{n} \sum_{(X, Y) \in \mathbb{D}} \|F(X, \theta(t)) - Y\|_F^2 \\ &= \frac{1}{n} \sum_{(X, Y) \in \mathbb{D}} \sum_{\ell=1}^L \|F_\ell(X, \theta(t)) - Y_\ell\|_2^2 \\ &= \frac{1}{n} \|F(t) - Y\|_F^2, \end{aligned}$$

where the first three steps follow from simple algebra, and the last step follows from the definitions of F and Y . \square

C.2 GRADIENT COMPUTATION

Lemma C.4. *For $\nu \in [N]$, we have:*

- **Part 1.** *We have:*

$$\frac{d\mathcal{L}(t, \mathbb{D})}{d \text{vec}(U_{(\nu)}(t))} = \frac{\omega \cdot \kappa}{\sqrt{m}} \sum_{p=1}^{nL} (\Lambda_{(\nu-1), i, \ell, *}(t) \otimes \Lambda_{(\nu-1), i}(t))^\top (\text{diag}(\sigma_{(\nu), p}(t)) - \sigma_{(\nu), p}(t)\sigma_{(\nu), p}(t)^\top)$$

$$\Lambda_{(\nu-1),i}(t) \sum_{r \in [m]} \left\langle \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu),p}(t)}, a_{(\nu),r} \right\rangle \cdot w_{(\nu),r}(t) \mathbb{I}\{o_{(\nu),p}(t)^\top w_{(\nu),r}(t) > 0\},$$

where $i = \lfloor p/L \rfloor$ and $\ell = p \bmod L$.

- **Part 2.** For any $r \in [m]$, we have:

$$\frac{d\mathcal{L}(t, \mathbb{D})}{dw_{(\nu),r}(t)} = \frac{\omega}{\sqrt{m}} \sum_{p=1}^{nL} \left\langle \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu),p}(t)}, a_{(\nu),r} \right\rangle \cdot o_{(\nu),p}(t) \cdot \mathbb{I}\{w_{(\nu),r}(t)^\top o_{(\nu),p}(t) > 0\}.$$

- **Part 3.** For $\nu \in [N]$ and $p \in [nL]$, we have:

$$\frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu),p}(t)} = \frac{2\varepsilon}{n} \cdot (I_d + \text{diag}(G_{(\nu),p}(t))) (F_p(t) - Y_p),$$

where

$$\begin{aligned} G_{(\nu),p}(t) &= \frac{\omega}{\sqrt{m}} \left(\sigma_{(\nu),p,\ell}(t) \cdot I_d \right. \\ &\quad \left. + \kappa U_{(\nu)}(t) \Lambda_{(\nu-1),i}(t)^\top (\text{diag}(\mathbf{1}_n - \frac{1}{2}e_\ell)) \cdot (\text{diag}(\sigma_{(\nu),p}(t)) - \sigma_{(\nu),p}(t)\sigma_{(\nu),p}(t)^\top) \Lambda_{(\nu-1),i}(t) \right) \\ &\quad \cdot \sum_{r \in [m]} \left\langle \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu+1),p}(t)}, a_{(\nu),r} \right\rangle \cdot w_{(\nu),r}(t) \mathbb{I}\{o_{(\nu),p}(t)^\top w_{(\nu),r}(t) > 0\}, \end{aligned}$$

and $G_{(N),p}(t) = \mathbf{0}_d$.

Here, we denote $i = \lfloor p/L \rfloor$ and $\ell = p \bmod L$. e_ℓ is a n -dimensional one-hot vector with the ℓ -th entry equal 1.

Proof. We omit the lemma proof since the gradient computation trivially follows from basic algebra and multivariable chain rules. \square

C.3 LEARNING DYNAMICS

We give the definition of NTK:

Definition C.5. We define the kernel matrix at ν -th layer as $H_{(\nu)} \in \mathbb{R}^{nL \times nL}$ and its (i, j) -th entry ($\forall (p, q) \in [nL] \times [nL]$) is defined as:

$$H_{(\nu),p,q}(t) := \underbrace{\langle \beta_{(\nu),p}(t), \beta_{(\nu),q}(t) \rangle}_{\text{kernel w.r.t. } W_{(\nu)}(t)} + \underbrace{\langle \gamma_{(\nu),p}(t), \gamma_{(\nu),q}(t) \rangle}_{\text{kernel w.r.t. } U_{(\nu)}(t)},$$

Here, we let:

$$\begin{aligned} \beta_{(\nu),p}(t) &:= \frac{\omega}{\sqrt{m}} \underbrace{o_{(\nu),p}(t)}_{d \times 1} \otimes \underbrace{\mathbf{1}_{W_{(\nu)}(t)^\top o_{(\nu),p}(t) > 0}}_{m \times 1} \in \mathbb{R}^{md}, \\ \gamma_{(\nu),p}(t) &:= \frac{\omega \cdot \kappa}{\sqrt{m}} \underbrace{(\Lambda_{(\nu-1),i,\ell,*}(t) \otimes \Lambda_{(\nu-1),i}(t))^\top}_{d^2 \times L} \underbrace{(\text{diag}(\sigma_{(\nu),p}(t)) - \sigma_{(\nu),p}(t)\sigma_{(\nu),p}(t)^\top)}_{L \times L} \\ &\quad \underbrace{\Lambda_{(\nu-1),i}(t)}_{L \times d} \sum_{r \in [m]} \underbrace{w_{(\nu),r}(t) \mathbb{I}\{o_{(\nu),p}(t)^\top w_{(\nu),r}(t) > 0\}}_{d \times 1} \in \mathbb{R}^{d^2}, \end{aligned}$$

where \otimes is the Kronecker product and $i = \lfloor p/L \rfloor$, $\ell = p \bmod L$. The indicator vector $\mathbf{1}_{W_{(\nu)}(t)^\top o_{(\nu),p}(t) > 0} \in \{0, 1\}^m$ where its r -th entry is $\mathbb{I}\{(W_{(\nu)}(t)^\top o_{(\nu),p}(t))_r > 0\}$ for $r \in [m]$.

We restate Lemma 4.1 below as its formal version:

Lemma C.6 (Formal version of Lemma 4.1). *The learning dynamics of the multi-layer transformer Eq. (1) is given by:*

$$\mathbb{E}\left[\frac{d}{dt}\mathcal{L}(t, \mathbb{D})\right] = - \sum_{\nu \in [N]} \underbrace{\text{vec}\left(\frac{d}{d\mu_{(\nu)}(t)}\mathcal{L}(t, \mathbb{D})\right)}_{1 \times nLd} \cdot \underbrace{(H_{(\nu)}(t) \otimes I_d)}_{nLd \times nLd} \cdot \underbrace{\text{vec}\left(\frac{d}{d\mu_{(\nu)}(t)}\mathcal{L}(t, \mathbb{D})\right)}_{nLd \times 1}$$

where $\mu_{(\nu)}(t)$ is a $nL \times d$ matrix, $\mu_{(\nu),p}(t)$ is the $(p \bmod L)$ -th row of ν -th layer output regarding to input matrix $X_{\lfloor p/L \rfloor}$ for any $p \in [nL]$ and $\nu \in [N]$.

Proof. We have:

$$\begin{aligned} \mathbb{E}\left[\frac{d}{dt}\mathcal{L}(t, \mathbb{D})\right] &= \mathbb{E}\left[\sum_{\nu=1}^N \left(\frac{d\mathcal{L}(t, \mathbb{D})}{d \text{vec}(W_{(\nu)}(t))} \frac{d \text{vec}(W_{(\nu)}(t))}{dt} + \frac{d\mathcal{L}(t, \mathbb{D})}{d \text{vec}(U_{(\nu)}(t))} \frac{d \text{vec}(U_{(\nu)}(t))}{dt}\right)\right] \\ &= - \sum_{\nu=1}^N \left(\frac{d\mathcal{L}(t, \mathbb{D})}{d \text{vec}(W_{(\nu)}(t))} \frac{d\mathcal{L}(t, \mathbb{D})}{d \text{vec}(W_{(\nu)}(t))} + \frac{d\mathcal{L}(t, \mathbb{D})}{d \text{vec}(U_{(\nu)}(t))} \frac{d\mathcal{L}(t, \mathbb{D})}{d \text{vec}(U_{(\nu)}(t))}\right) \\ &= - \sum_{\nu=1}^N \sum_{p=1}^{nL} \sum_{q=1}^{nL} \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu),p}(t)} (\langle \beta_{(\nu),p}(t), \beta_{(\nu),q}(t) \rangle + \langle \gamma_{(\nu),p}(t), \gamma_{(\nu),q}(t) \rangle) \cdot I_d \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu),q}(t)} \\ &= - \sum_{\nu \in [N]} \text{vec}\left(\frac{d}{d\mu_{(\nu)}(t)}\mathcal{L}(t, \mathbb{D})\right)^\top \cdot (H_{(\nu)}(t) \otimes I_d) \cdot \text{vec}\left(\frac{d}{d\mu_{(\nu)}(t)}\mathcal{L}(t, \mathbb{D})\right), \end{aligned}$$

where the first step follows from the chain rule, the second step follows from the gradient flow (Eq. (2)), the third step follows from Part 1, Part 2 and Part 3 of Lemma C.4 and the definitions of $\beta_{(\nu),p}(t)$ and $\gamma_{(\nu),p}(t)$, the last step follows from the definition of $H_{(\nu)}(t)$. \square

D TOOLKIT: HELPFUL BOUNDARIES

We give the lemma about all help bounds in the range of this paper as follows:

Lemma D.1. Denote failure probability $\delta \in (0, 0.1)$. Define $B := \max\{O(\sqrt{\log(Lmd/\delta)}), 1\}$. Assuming there exists a constant $R \in (0, 1)$ satisfying $\|w_{(\nu),r}(t) - w_{(\nu),r}(0)\|_2 \leq R$ and $\|U_{(\nu)}(t) - U_{(\nu)}(0)\|_F \leq R$ for $r \in [m]$ and $\nu \in [N]$ and $R \in (0, 1)$.

We mark the indices as: $r \in [m]$, $p \in [nL]$, $\ell_1, \ell_2 \in [L]$, and $\nu \in [N]$. We denote $i = \lfloor p/L \rfloor$, $\ell = p \bmod L$, $\ell' \in [\ell]$.

If Definition 5.2 holds, then with a probability at least $1 - \delta$, we have:

- Basic Bounds.

- **Part 1.** $\|w_{(\nu),r}(0)\|_2 \leq O(\sqrt{dB})$.
- **Part 2.** $\|U_{(\nu)}(0)\|_F \leq O(dB)$.
- **Part 3.** $\|w_{(\nu),r}(t)\|_2 \leq O(\sqrt{dB})$.
- **Part 4.** $\|U_{(\nu)}(t)\|_F \leq O(dB)$.
- **Part 5.** $\|\Lambda_{(\nu),i,\ell_1}(t)\|_2 = \Theta(1)$.
- **Part 6.** $\|\Lambda_{(\nu),i}(t)U_{(\nu)}(t)\Lambda_{(\nu),i}(t)^\top\|_\infty \leq O(dB)$.
- **Part 7.** $\exp_{\ell_1,\ell_2}(\Lambda_{(\nu),i}(t)U_{(\nu)}(t)\Lambda_{(\nu),i}(t)^\top) \in [\exp(-O(dB)), \exp(O(dB))]$.
- **Part 8.** For $\ell' \in [\ell]$, $\sigma_{(\nu),p,\ell'}(t) \in [\exp(-O(dB))/L, 1]$.

- Perturbation Bounds.

- **Part 9.** $\|w_{(\nu),r}(t) - w_{(\nu),r}(0)\|_2 \leq R$.
- **Part 10.** $\|U_{(\nu)}(t) - U_{(\nu)}(0)\|_F \leq R$.
- **Part 11.** $\|\Lambda_{(\nu),i,\ell}(t) - \Lambda_{(\nu),i,\ell}(0)\|_2 \leq O(R)$.
- **Part 12.** $\|\sigma_{(\nu),p}(t) - \sigma_{(\nu),p}(0)\|_2 \leq O(\sqrt{LR})$.
- **Part 13.** $\|o_{(\nu),p}(t) - o_{(\nu),p}(0)\|_2 \leq O(\sqrt{LR})$.

- Gradient and Function Norms.

- **Part 14.** $\mathcal{L}(t, \mathbb{D}) \leq O(Ld)$.
- **Part 15.** $\|\frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu)}(t)}\|_F^2 \asymp \varepsilon^2 \cdot \mathcal{L}(t, \mathbb{D})$.
- **Part 16.** $\|\gamma_{(\nu),p}(t)\|_2 \leq o(1/\sqrt{m})$.

Proof. **Proof of Part 1.** This proof follows from the initialization of $W_{(\nu)}(0)$ and the union bound of the tail bound of the Gaussian distribution (Fact B.7).

Proof of Part 2. This proof follows from initialization of $U_{(\nu)}(0)$ and the union bound of the tail bound of the Gaussian distribution (Fact B.7).

Proof of Part 3. This proof combines $\|w_{(\nu),r}(t) - w_{(\nu),r}(0)\|_2 \leq R$, triangle inequality and $R \leq B$.

Proof of Part 4. This proof combines $\|U_{(\nu)}(t) - U_{(\nu)}(0)\|_F \leq R$, triangle inequality and $R \leq B$.

Proof of Part 5. We have:

$$\|\Lambda_{(0),i,\ell}(t)\|_2 = \Theta(1).$$

Therefore, we have:

$$\|o_{(1),p}(t)\|_2 = \|\Lambda_{(0),i}(t)^\top \sigma_{(\nu),p}(t)\|_2 = \Theta(1),$$

where this step follows from the property that the sum of any softmax vector is 1.

Therefore, we can show that:

$$\|\mu_{(\nu),p}(t)\|_2 = \left\| \frac{\omega}{\sqrt{m}} \sum_{r=1}^m a_{(1),r} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle) \right\|_2$$

$$\begin{aligned}
&\leq \frac{\omega\sqrt{d}}{\sqrt{m}} \left\| \sum_{r=1}^m a_{(1),r} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle) \right\|_{\infty} \\
&\leq \frac{\omega\sqrt{d}}{\sqrt{m}} \max_{k \in [d]} \left| \sum_{r=1}^m a_{(1),r,k} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle) \right|
\end{aligned}$$

where the first step follows from the definition of $\mu_{(\nu),p}(t)$, the second step follows from simple algebras, the third step follows from the definition of infinite norm.

We apply Hoeffding bound to each variable $a_{(1),r,k} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle)$, we have:

$$\begin{aligned}
|a_{(1),r,k} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle)| &\leq O(\sqrt{dB}), \\
\mathbb{E}[a_{(1),r,k} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle)] &= 0.
\end{aligned}$$

With a probability at least $1 - \delta$, we have:

$$\begin{aligned}
\|\mu_{(\nu),p}(t)\|_2 &\leq \frac{\omega\sqrt{d}}{\sqrt{m}} \max_{k \in [d]} \left| \sum_{r=1}^m a_{(1),r,k} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle) \right| \\
&\leq \frac{\omega\sqrt{d}}{\sqrt{m}} \cdot O(\sqrt{mdB}) \sqrt{\log(m/\delta)} \leq O(\omega dB^2).
\end{aligned}$$

Hence,

$$\begin{aligned}
\|\Lambda_{(1),i,\ell}(t)\|_2 &= \left\| \frac{\omega}{\sqrt{m}} \sum_{r=1}^m a_{(1),r} \cdot \phi(\langle o_{(1),p}(t), w_{(1),r}(t) \rangle) + \Lambda_{(0),i,\ell}(t) \right\|_2 \\
&= \Theta(1) \pm O(\omega dB^2) \\
&= \Theta(1) \pm o\left(\frac{1}{N}\right)
\end{aligned}$$

the last step follows from choosing $\omega = o\left(\frac{1}{NdB^2}\right)$.

By induction, we can get:

$$\|\Lambda_{(\nu),i,\ell}(t)\|_2 = \Theta(1),$$

and

$$\|o_{(\nu),p}(t)\|_2 = \Theta(1). \quad (4)$$

Proof of Part 6. We have:

$$\begin{aligned}
\|\Lambda_{(\nu),i}(t) U_{(\nu)}(t) \Lambda_{(\nu),i}(t)^\top\|_{\infty} &= \max_{(\ell_1, \ell_2) \in [L] \times [L]} \left| \Lambda_{(\nu),i,\ell_1}(t)^\top U_{(\nu)}(t) \Lambda_{(\nu),i,\ell_2}(t) \right| \\
&\leq O(dB)
\end{aligned}$$

where the first step follows from simple algebra, the second step follows from Part 4 and Part 5 of this lemma and Cauchy-Schwarz inequality.

Proof of Part 7. This proof follows from Part 3 of this lemma and simple algebra.

Proof of Part 8. Following Part 7 of this lemma, we can show that:

$$\sum_{\ell' \in [\ell]} \exp(\Lambda_{(\nu),i,\ell}(t)^\top U_{(\nu)}(t) \Lambda_{(\nu),i,\ell'}(t)) \leq \ell \exp(O(dB)) \leq L \exp(O(dB)).$$

Then we can show that:

$$\begin{aligned}
\sigma_{(\nu),p,\ell'}(t) &\geq \exp(\Lambda_{(\nu),i,\ell}(t)^\top U_{(\nu)}(t) \Lambda_{(\nu),i,\ell'}(t)) / \sum_{\ell' \in [\ell]} \exp(\Lambda_{(\nu),i,\ell}(t)^\top U_{(\nu)}(t) \Lambda_{(\nu),i,\ell'}(t)) \\
&\geq \frac{\exp(-O(dB))}{L \exp(O(dB))}
\end{aligned}$$

$$\geq \frac{\exp(-O(dB))}{L}.$$

Combining the previous results, we obtain the result of this part.

Proof of Part 9. Directly from the lemma condition $\|w_{(\nu),r}(t) - w_{(\nu),r}(0)\|_F \leq R$.

Proof of Part 10. Directly from the lemma condition $\|U_{(\nu)}(t) - U_{(\nu)}(0)\|_F \leq R$.

Proof of Part 11, 12 and 13. When $\nu = 1$, we have:

$$\begin{aligned} & \|\Lambda_{(1),i,\ell}(t) - \Lambda_{(1),i,\ell}(0)\|_F \\ &= \sqrt{\sum_{k \in [d]} \left(\frac{\omega}{\sqrt{m}} \cdot \sum_{r=1}^m a_{(1),r,k} (\phi(\langle w_{(1),r}(t), o_{(1),p}(t) \rangle) - \phi(\langle w_{(1),r}(0), o_{(1),p}(0) \rangle)) \right)^2}, \end{aligned}$$

besides, $\Lambda_{(0),i,\ell}(t) = \Lambda_{(0),i,\ell}(0)$.

Next, we have:

$$\begin{aligned} \|\sigma_{(1),p}(t) - \sigma_{(1),p}(0)\|_2 &\leq O(1) \cdot \|\Lambda_{(0),i}(0)(U_{(1)}(t) - U_{(1)}(0))^\top \Lambda_{(0),i,\ell}(0)\|_2 \\ &\leq O(\sqrt{LR}) \end{aligned}$$

where the first step follows from the definition of $\sigma_{(1),i}(t)$ and Fact B.12, the second step follows from $\|U_{(1)}(t) - U_{(1)}(0)\|_F \leq R$, $\|\Lambda_{(0),i}(0)\|_F = O(\sqrt{L})$, $\|\Lambda_{(0),i,\ell}(0)\|_2 = O(1)$ and Cauchy-Schwarz inequality.

Thus, we can show that

$$\begin{aligned} \|o_{(1),i}(t) - o_{(1),i}(0)\|_2 &= \|\Lambda_{(0),i}(0)^\top \sigma_{(1),p}(t) - \Lambda_{(0),i}(0)^\top \sigma_{(1),p}(0)\|_2 \\ &= \sqrt{\sum_{\ell'=1}^{\ell} (\sigma_{(1),p,\ell'}(t) - \sigma_{(1),p,\ell'}(0))^2 \|\Lambda_{(0),i,\ell'}(0)\|_2^2} \\ &\leq \sqrt{\sum_{\ell'=1}^{\ell} (\sigma_{(1),p,\ell'}(t) - \sigma_{(1),p,\ell'}(0))^2 O(1)} \\ &= O(1) \cdot \|\sigma_{(1),p}(t) - \sigma_{(1),p}(0)\|_2 \\ &\leq O(\sqrt{LR}), \end{aligned}$$

where these steps follow from some basic algebra, Fact B.12 and Part 10 of this Lemma.

For a certain index $k \in [d]$, we apply Hoeffding's inequality (Lemma B.14) to each random variable $a_{(1),r} \cdot \phi(\langle w_{(1),r}(t), o_{(1),i}(t) \rangle)$, besides, we have:

$$\begin{aligned} & |a_{(1),r,k} \cdot (\phi(\langle w_{(1),r}(t), o_{(1),i}(t) \rangle) - \phi(\langle w_{(1),r}(0), o_{(1),i}(0) \rangle))| \\ &\leq \max\{\|w_{(1),r}(0)\| \cdot O(\sqrt{LR}), \|o_{(1),i}(t)\|_2 \cdot O(R)\} \\ &= O(\sqrt{Ld}BR), \\ &\mathbb{E}[a_{(1),r,k} \cdot (\phi(\langle w_{(1),r}(t), o_{(1),i}(t) \rangle) - \phi(\langle w_{(1),r}(0), o_{(1),i}(0) \rangle))] = 0. \end{aligned}$$

Then with a probability at least $1 - \delta$, we have:

$$\frac{\omega}{\sqrt{m}} \sum_{r=1}^m a_{(1),r,k} \cdot \phi(\langle w_{(1),r}(t), o_{(1),i}(t) \rangle) \leq O\left(\frac{\omega}{\sqrt{m}} BR\right) \sqrt{Lmd \log(1/\delta)} \leq O(R/(Bd\sqrt{Ld}))$$

where the last step follows from choosing $\omega = o(\frac{1}{\sqrt{Ld^2 B^3}})$.

We therefore get $\|\Lambda_{(1),i,\ell}(t) - \Lambda_{(1),i,\ell}(0)\|_2 \leq O(R)$ for all $\ell \in [\ell]$ and $\|\Lambda_{(1),i,\ell}(t) - \Lambda_{(1),i,\ell}(0)\|_\infty \leq O(R/\sqrt{d})$.

When $\nu = 2$, we have:

$$\|\sigma_{(2),p}(t) - \sigma_{(2),p}(0)\|_2$$

$$\begin{aligned}
&\leq O(1) \cdot \|\Lambda_{(1),i}(t)U_{(2)}(t)^\top \Lambda_{(1),i,\ell}(t) - \Lambda_{(1),i}(0)U_{(2)}(0)^\top \Lambda_{(1),i,\ell}(0)\|_2 \\
&= O(1) \cdot \|(\Lambda_{(1),i}(t) \otimes \Lambda_{(1),i,\ell}(t)) \text{vec}(U_{(2)}(t)) - (\Lambda_{(1),i}(0) \otimes \Lambda_{(1),i,\ell}(0)) \text{vec}(U_{(2)}(0))\|_2 \\
&\leq O(1) \cdot \max\{\|\Lambda_{(1),i}(t) \otimes \Lambda_{(1),i,\ell}(t)\|_F \cdot O(R), \\
&\quad \|U_{(2)}(0)\|_F \cdot \|\Lambda_{(1),i}(t) \otimes \Lambda_{(1),i,\ell}(t) - \Lambda_{(1),i}(0) \otimes \Lambda_{(1),i,\ell}(0)\|_F\} \\
&\leq O(\sqrt{LR})
\end{aligned}$$

where the first step follows from simple algebra and Fact B.12, the second step follows from a basic tensor trick, the third step follows from simple algebra, triangle inequality, Part 4 of this lemma and $\|\Lambda_{(1),i,\ell}(t) - \Lambda_{(1),i,\ell}(0)\|_\infty \leq O(R/(Bd\sqrt{Ld}))$.

Hence, we have:

$$\begin{aligned}
\|o_{(2),i}(t) - o_{(2),i}(0)\|_2 &= \|\Lambda_{(1),i}(t)^\top \sigma_{(2),p}(t) - \Lambda_{(1),i}(0)^\top \sigma_{(2),p}(0)\|_2 \\
&\leq \max\{\|\Lambda_{(1),i}(t)\| \cdot O(RLdB), O(1) \cdot O(\sqrt{LR})\} \\
&\leq O(\sqrt{LR}),
\end{aligned}$$

where these steps follow from some basic algebra, Fact B.12 and Part 10 of this Lemma.

Here, we similarly apply Hoeffding inequality (Lemma B.14), and get: with a probability at least $1 - \delta$, we have:

$$\|\Lambda_{(2),i,\ell}(t) - \Lambda_{(2),i,\ell}(0)\|_2 \leq O(R).$$

By induction, we obtain the results of Part 11, 12, and 13.

Proof of Part 14. We have:

$$\begin{aligned}
\mathcal{L}(t, \mathbb{D}) &= \frac{1}{n} \|\mathbb{F}(t) - \mathbb{Y}\|_F^2 \\
&\leq L \max_{p \in [nL]} \|\mathbb{F}_p(t) - \mathbb{Y}_p\|_2^2 \\
&\leq L \max_{p \in [nL]} (\|\mathbb{F}_p(t)\|_2 + \|\mathbb{Y}_p\|_2)^2 \\
&\leq L \left(O(1) + O(\sqrt{d}) \right)^2 \\
&\leq O(Ld),
\end{aligned}$$

where the first step follows from Lemma C.3, the second step follows from simple algebras, the third step follows from the Cauchy-Schwartz inequality, the last two steps follow from Definition B.3 and simple algebras.

Proof of Part 15. It is easy to prove:

$$\left\| \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(N)}(t)} \right\|_F \leq O(\sqrt{d}).$$

Besides, we can obtain:

$$\begin{aligned}
&\kappa U_{(N)}(t) \Lambda_{(N-1),i}(t)^\top \left(\text{diag}(\mathbf{1}_n - \frac{1}{2}e_\ell) \cdot (\text{diag}(\sigma_{(N),p}(t)) - \sigma_{(N),p}(t)\sigma_{(N),p}(t)^\top) \right) \Lambda_{(N-1),i}(t) \\
&\leq O(\kappa L^2 dB) \\
&\leq O(L^2 dB).
\end{aligned}$$

Then, by Hoeffding inequality (Lemma B.14), we have:

$$\|G_{(N-1),p}(t)\|_2^2 \leq \omega \cdot O(L^2 d^{2.5} B^3).$$

Choosing $\omega = o(\frac{1}{L^2 d^{2.5} B^3})$, we have: $\|G_{(N-1),p}(t)\|_2^2 \leq o(1)$.

By induction, we can show that: $\left\| \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(N)}(t)} \right\|_F \leq O(\sqrt{d})$.

1512 Similarly, we have:

1513

1514

1515

1516

$$\left\| \frac{d\mathcal{L}(t, \mathbb{D})}{d\mu_{(\nu)}(t)} \right\|_F^2 \asymp \varepsilon^2 \cdot \mathcal{L}(t, \mathbb{D}).$$

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

1540

1541

1542

1543

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565

Proof of Part 16. This proof is trivially similar with bounding $\|G_{(\nu),p}(t)\|_2^2 \leq o(1)$, then we choose $\kappa = 1/\sqrt{m}$ to meet the result.

This completes the proof of all parts of the lemma. \square

1566 E KERNEL PERTURBATION

1567 Here is a formal version of confirmation of Lemma 5.3:

1568 **Lemma E.1** (Formal version of Lemma 5.3). *Assuming Assumption 5.1 and Definition 5.2 hold,*
 1569 *denote the failure probability $\delta \in (0, 0.1)$, then the kernel perturbation bound is:*

$$1570 \Pr [\lambda_{\min}(H_{(\nu)}(t)) < \lambda/2] < \delta.$$

1571 Therefore, bounding loss dynamics is given by ($C > 0$ is some constant):

$$1572 \Pr \left[\mathbb{E} \left[\frac{d}{dt} \mathcal{L}(t, \mathbb{D}) \right] > -C \cdot \omega \lambda N \cdot \mathcal{L}(t, \mathbb{D}) \right] < \delta.$$

1573 *Proof. Proof of Part 1.* First, we assume the condition $\|w_{(\nu),r}(t) - w_{(\nu),r}(0)\|_2 \leq R$ and
 1574 $\|U_{(\nu)}(t) - U_{(\nu)}(0)\|_F \leq R$ as Lemma D.1, which will be confirmed as a provable property in the
 1575 further analysis.

1576 For any $p, q \in [nL]$, we have:

$$1577 \begin{aligned} & |H'_{(\nu),p,q}(t) - H'_{(\nu),p,q}(0)| \\ &= \frac{1}{m} |o_{(\nu),p}(t)^\top o_{(\nu),p}(t) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(t)^\top w_{(\nu),r}(t) > 0, o_{(\nu),q}(t)^\top w_{(\nu),r}(t) > 0\} \\ &\quad - o_{(\nu),p}(0)^\top o_{(\nu),p}(0) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(0)^\top w_{(\nu),r}(0) > 0, o_{(\nu),q}(0)^\top w_{(\nu),r}(0) > 0\}| \\ &\leq \frac{1}{m} (Q_{(\nu),p,q,1} + Q_{(\nu),p,q,2} + Q_{(\nu),p,q,3}), \end{aligned}$$

1578 where the first step follows from the definition of $H'_{(\nu),p,q}(t)$, and the second step follows from the
 1579 triangle inequality and defining:

$$1580 \begin{aligned} Q_{(\nu),p,q,1} &:= |o_{(\nu),p}(t)^\top o_{(\nu),p}(t) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(t)^\top w_{(\nu),r}(t) > 0, o_{(\nu),q}(t)^\top w_{(\nu),r}(t) > 0\} \\ &\quad - o_{(\nu),p}(0)^\top o_{(\nu),p}(0) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(0)^\top w_{(\nu),r}(0) > 0, o_{(\nu),q}(0)^\top w_{(\nu),r}(0) > 0\}|, \\ Q_{(\nu),p,q,2} &:= |o_{(\nu),p}(0)^\top o_{(\nu),p}(t) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(t)^\top w_{(\nu),r}(t) > 0, o_{(\nu),q}(t)^\top w_{(\nu),r}(t) > 0\} \\ &\quad - o_{(\nu),p}(0)^\top o_{(\nu),p}(0) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(0)^\top w_{(\nu),r}(0) > 0, o_{(\nu),q}(0)^\top w_{(\nu),r}(0) > 0\}|, \\ Q_{(\nu),p,q,3} &:= |o_{(\nu),p}(0)^\top o_{(\nu),p}(0) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(t)^\top w_{(\nu),r}(t) > 0, o_{(\nu),q}(t)^\top w_{(\nu),r}(t) > 0\} \\ &\quad - o_{(\nu),p}(0)^\top o_{(\nu),p}(0) \sum_{r \in [m]} \mathbb{I}\{o_{(\nu),p}(0)^\top w_{(\nu),r}(0) > 0, o_{(\nu),q}(0)^\top w_{(\nu),r}(0) > 0\}|. \end{aligned}$$

1581 We assume a constant R that satisfies $\|w_{(\nu),r}(t) - w_{(\nu),r}(0)\|_2 \leq R$ and $\|U_{(\nu)}(t) - U_{(\nu)}(0)\|_F \leq R$.

1582 **Bounding $Q_{(\nu),p,q,1}$.** We have:

$$1583 \begin{aligned} Q_{(\nu),p,q,1} &\leq \left| \left(o_{(\nu),p}(t) - o_{(\nu),p}(0) \right)^\top o_{(\nu),q}(t) \right. \\ &\quad \cdot \left. \sum_{r=1}^m \mathbb{I}\{ \langle o_{(\nu),p}(t), w_{(\nu),r}(t) \rangle > 0, \langle o_{(\nu),q}(t), w_{(\nu),r}(t) \rangle > 0 \} \right| \\ &\leq m \left| \left(o_{(\nu),p}(t) - o_{(\nu),p}(0) \right)^\top o_{(\nu),q}(t) \right| \\ &\leq m \|o_{(\nu),p}(t) - o_{(\nu),p}(0)\|_2 \cdot \|o_{(\nu),q}(t)\| \end{aligned}$$

$$\leq m \cdot O(\sqrt{LR}),$$

where the first step follows from the definition of $Q_{(\nu),p,q,1}$ and Cauchy-Schwarz inequality, the second step follows from $\mathbb{I}\{\langle o_{(\nu),p}(t), w_{(\nu),r}(t) \rangle > 0, \langle o_{(\nu),q}(t), w_{(\nu),r}(t) \rangle > 0\} \leq 1$, the third step follows from Cauchy-Schwarz inequality, the last step follows from Part 6 and Part 13 of Lemma D.1 and $\|o_{(\nu),q}(t)\| \leq O(1)$.

Bounding $Q_{(\nu),p,q,2}$. We omit this proof since it is similar to the proof of bounding $Q_{(\nu),p,q,1}$, we have:

$$Q_{(\nu),p,q,2} \leq m \cdot O(\sqrt{LR}).$$

Bounding $Q_{(\nu),p,q,3}$. We define the following event:

$$\begin{aligned} E_{(\nu),p,r} &:= \{\exists w \in \mathbb{R}^d : \|w - w_{(\nu),r}(0)\|_2 \leq R, \\ &\quad \mathbb{I}\{\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle > 0\} \neq \mathbb{I}\{\langle o_{(\nu),p}(t), w_{(\nu),r}(t) \rangle > 0\}\}. \end{aligned}$$

It is easy to hold that, once:

$$\begin{aligned} |\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle| &\geq O(\sqrt{LR}) \cdot O(\sqrt{dB}) + O(1) \cdot R \\ \iff |\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle| &\geq O(\sqrt{LdB}R), \end{aligned}$$

the event $E_{(\nu),p,r}$ is false, since we combining Part 1, Part 9 and Part 13 of Lemma D.1, $\|o_{(\nu),q}(t)\| \leq O(1)$ and some simple algebra.

Following Fact B.10, we have:

$$\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle \sim \mathcal{N}(0, \|o_{(\nu),p}(0)\|_2^2).$$

We have:

$$\|o_{(\nu),p}(0)\|_2^2 = \sum_{\ell'=1}^{\ell} \sigma_{(\nu),p,\ell'}(t)^2 \|\Lambda_{(\nu),p,\ell'}(t)\|_2^2 = O(1) \cdot \sum_{\ell'=1}^{L_i} \sigma_{(\nu),p,\ell'}(t)^2 \geq \ell \exp(-O(dB)) \geq \exp(-O(dB)),$$

where $\ell = p \bmod L$.

Then, the anti-concentration of $\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle$ shows that:

$$\begin{aligned} \Pr[E_{(\nu),p,r}] &\leq \frac{1}{\Theta(1) \cdot \|o_{(\nu),p}(0)\|_2} O(\sqrt{LdB}R) \\ &\leq O(\sqrt{LR}) \cdot \exp(O(dB)), \end{aligned}$$

where the first step follows from Fact B.9, the second step follows from simple algebra and $\|o_{(\nu),p}(0)\|_2^2 \geq \exp(-O(dB))$.

We have:

$$\begin{aligned} \mathbb{E}[Q_{(\nu),p,q,3}] &= \mathbb{E} \left[\left| \langle o_{(\nu),p}(0)^\top, o_{(\nu),q}(0) \rangle \right| \right. \\ &\quad \cdot \sum_{r=1}^m \left| \mathbb{I}\{\langle o_{(\nu),p}(t), w_{(\nu),r}(t) \rangle > 0, \langle o_{(\nu),q}(t), w_{(\nu),r}(t) \rangle > 0\} \right. \\ &\quad \left. \left. - \mathbb{I}\{\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle > 0, \langle o_{(\nu),q}(0), w_{(\nu),r}(0) \rangle > 0\} \right| \right] \\ &\leq \mathbb{E} \left[O(1) \cdot \sum_{r=1}^m \left| \mathbb{I}\{\langle o_{(\nu),p}(t), w_{(\nu),r}(t) \rangle > 0, \langle o_{(\nu),q}(t), w_{(\nu),r}(t) \rangle > 0\} \right. \right. \\ &\quad \left. \left. - \mathbb{I}\{\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle > 0, \langle o_{(\nu),q}(0), w_{(\nu),r}(0) \rangle > 0\} \right| \right] \\ &\leq O(1) \cdot \sum_{r=1}^m \mathbb{E} \left[\left| \mathbb{I}\{\langle o_{(\nu),p}(t), w_{(\nu),r}(t) \rangle > 0, \langle o_{(\nu),q}(t), w_{(\nu),r}(t) \rangle > 0\} \right| \right] \end{aligned}$$

$$\begin{aligned}
& - \mathbb{I}\{\langle o_{(\nu),p}(0), w_{(\nu),r}(0) \rangle > 0, \langle o_{(\nu),q}(0), w_{(\nu),r}(0) \rangle > 0\} \Big] \\
& \leq O(1) \cdot \sum_{r=1}^m \mathbb{E} \left[\mathbb{I}\{\mathbb{E}_{(\nu),p,r} \cup \mathbb{E}_{(\nu),q,r}\} \right] \\
& \leq O(1) \cdot \sum_{r=1}^m \exp(O(dB)) \cdot RL^{\nu/2} \\
& \leq mO(\sqrt{LR}) \cdot \exp(O(dB)),
\end{aligned}$$

where the first step follows from the definition of $Q_{(\nu),p,q,3}$, the second step follows from $\|o_{(\nu),p}(0)\|_2 \leq O(1)$, the third and fourth step follow from the rules of expectation, the last two steps follow from $\Pr[\mathbb{E}_{(\nu),p,r}] \leq O(\sqrt{LR}) \cdot \exp(O(dB))$ and simple algebra.

Hence, using Markov's inequality (Lemma B.15), with a probability at least $1 - \delta$, we have:

$$Q_{(\nu),p,q,3} \leq m \cdot O(\sqrt{LR}) \cdot \exp(O(dB))/\delta.$$

Combine the upper bounds on three terms, we have:

$$|H'_{(\nu),p,q}(t) - H'_{(\nu),p,q}(0)| \leq O(\sqrt{LR}) \cdot \exp(O(dB))/\delta.$$

We have:

$$\begin{aligned}
\|H'_{(\nu)}(t) - H'_{(\nu)}(0)\|_F &= \sqrt{\sum_{p=1}^{nL} \sum_{q=1}^{nL} (H'_{(\nu),p,q}(t) - H'_{(\nu),p,q}(0))^2} \\
&\leq O(nL^{1.5}R) \cdot \exp(O(dB))/\delta.
\end{aligned}$$

Following Fact B.11 and choose

$$R \leq \frac{\omega\lambda\delta}{nL^{1.5} \exp(O(dB))}, \quad (5)$$

we have:

$$\lambda_{\min}(H'_{(\nu)}(t)) \geq \frac{3}{4}\omega\lambda.$$

It is trivial to have (Part 16 of Lemma D.1):

$$\|\gamma_{(\nu),p}(t)\|_2 \leq o\left(\frac{1}{\sqrt{m}}\right).$$

Therefore, we have:

$$\begin{aligned}
|H_{(\nu),p,q}(t) - H'_{(\nu),p,q}(t)| &= |\langle \gamma_{(\nu),p}(t), \gamma_{(\nu),q}(t) \rangle| \\
&\leq o(1/m).
\end{aligned}$$

Furthermore, we have:

$$\|H_{(\nu)}(t) - H'_{(\nu)}(t)\|_F \leq o\left(\frac{n}{m}\right) \leq \omega\lambda/4,$$

where the last step follows from $m = \Omega(n/(\omega\lambda))$.

Similarly, following Fact B.11 and combining the previous result, we have:

$$\lambda_{\min}(H_{(\nu)}(t)) \geq \omega\lambda/2.$$

Following Fact B.13, we have:

$$\lambda_{\min}(H_{(\nu)}(t) \otimes I_d) \geq \omega\lambda/2.$$

1728 Combining Lemma C.6, Part 15 of Lemma D.1 and Fact B.11, we can show that:

$$1729 \mathbb{E}\left[\frac{d}{dt}\mathcal{L}(t, \mathbb{D})\right] \leq -C \cdot \omega \lambda N \cdot \varepsilon^2 \cdot \mathcal{L}(t, \mathbb{D}).$$

1732 Solving this ODE, we have:

$$1733 \mathbb{E}[\mathcal{L}(T, \mathbb{D})] \leq \exp(-C \cdot \varepsilon^2 \omega \lambda N T) \cdot \mathcal{L}(0, \mathbb{D}). \quad (6)$$

1736 **Bounding Gradient Norm.** Following Part 3 of Lemma D.1, we have:

$$1737 \begin{aligned} & \max_{T \geq 0} \max_{\nu \in [N]} \max_{r \in [m]} \|w_{(\nu), r}(T) - w_{(\nu), r}(0)\|_2 \\ & \leq \max_{T \geq 0} \int_0^T \max_{\nu \in [N]} \max_{r \in [m]} \left\| \frac{d\mathcal{L}(s, \mathbb{D})}{dw_{(\nu), r}(s)} \right\|_2 ds \\ & \leq \max_{T \geq 0} \int_0^T \exp(-C \cdot \varepsilon^2 \omega \lambda N s) \cdot \mathcal{L}(0, \mathbb{D}) \cdot \frac{\omega}{\sqrt{m}} \cdot O(\sqrt{dB}) ds \\ & \leq \max_{T \geq 0} \int_0^T \exp(-C \cdot \varepsilon^2 \omega \lambda N s) \cdot \frac{1}{\sqrt{m}} ds \\ & = \frac{1}{\sqrt{m}} \max_{T \geq 0} \left[-\frac{1}{C \varepsilon^2 \omega \lambda N} \exp(-C \cdot \varepsilon^2 \omega \lambda N s) \right]_{s=0}^{s=T} \\ & \leq O\left(\frac{1}{\sqrt{m} \varepsilon^2 \omega \lambda N}\right) \\ & \leq R \end{aligned} \quad (7)$$

1752 where the first step follows from Eq. (2) and Cauchy-Schwartz inequality, the second step follows
1753 from Eq (6), the third step follows from Part 14 of Lemma D.1 and the choice of ω , the last two steps
1754 follow from simple algebras and choosing

$$1755 \begin{aligned} m &= \Omega\left(\frac{n^2 L^3 \exp(O(dB))}{\varepsilon^4 \omega^4 \lambda^4 \delta^2 N^2}\right) \iff m/\text{polylog}(N, m) = \Omega\left(\frac{n^2 L^3 \exp(Cd)}{\varepsilon^4 \omega^4 \lambda^4 \delta^2 N^2}\right) \\ &\iff m/\text{polylog}(m) = \Omega\left(\frac{n^2 L^3 \exp(Cd)}{\varepsilon^4 \omega^4 \lambda^4 \delta^2 N^{\frac{4}{3}}}\right) \\ &\iff m^{\frac{2}{3}} = \Omega\left(\frac{n^2 L^3 \exp(Cd)}{\varepsilon^4 \omega^4 \lambda^4 \delta^2 N^{\frac{4}{3}}}\right) \\ &\iff m = \Omega\left(\frac{n^3 L^5 \exp(Cd)}{\varepsilon^6 \omega^6 \lambda^6 \delta^3 N^2}\right) \end{aligned}$$

1765 where these steps follow from simple algebras.

1766 Similarly, we have:

$$1767 \max_{T \geq 0} \max_{\nu \in [N]} \|U_{(\nu)}(T) - U_{(\nu)}(0)\|_F \leq O\left(\frac{1}{\sqrt{m} \varepsilon^2 \omega \lambda N}\right) \leq R$$

1771 \square

1782 F CONVERGENCE AND APPROXIMATION BOUND

1784 F.1 COMPLEXITY ANALYSIS

1786 **Lemma F.1.** *We have*

- 1787 • **Part 1.** *The time complexity (forward/backward) of the transformer on a single data point is*
1788 $O(NLmd) = O(ML) \leq O(M)$.
1789
- 1790 • **Part 2.** *The number of neurons in transformer is* $M = O(N(md + d^2)) \leq O(Nmd)$.
1791

1792 *Proof. Proof of Part 1.* We first note that the naive time complexity of matrix multiplication is
1793 $d_1 d_2 d_3$ for matrices $A \in \mathbb{R}^{d_1 \times d_2}$ and $B \in \mathbb{R}^{d_2 \times d_3}$. Formally, $\text{Multiply}(A, B) = O(d_1 d_2 d_3)$.
1794

1795 The following analysis follows from $v \in [N]$, $p \in [nL]$ and $i = \lfloor p/L \rfloor$.

1796 The complexity of $\text{Multiply}(\Lambda_{(\nu-1),i}(t), U_{(\nu)}(t))$ is $O(Ld^2)$.

1797 The complexity of $\text{Multiply}(\Lambda_{(\nu-1),i}(t)U_{(\nu)}(t), \Lambda_{(\nu-1),i}(t)^\top)$ is $O(L^2d)$.
1798

1799 The complexity of naive softmax is $O(L^2)$.

1800 The complexity of $\text{Multiply}(\Lambda_{(\nu-1),i}(t)^\top, \sigma_{(\nu),p}(t))$ is $O(Ld)$. Since the parallelization, total com-
1801 plexity should be $O(L^2d)$.
1802

1803 The complexity of $\text{Multiply}(W_{(\nu)}(t)^\top, o_{(\nu),p}(t))$ is $O(Lmd)$.

1804 Since $A_{(\nu)} \in \{-1, +1\}^{m \times d}$, there exists a algorithm using addition operation to implement
1805 $\text{Multiply}(A_{(\nu)}^\top, W_{(\nu)}(t)^\top o_{(\nu),p}(t)^\top)$ with complexity $O(Lmd)$.
1806

1807 Summing all complexity, we have the total complexity:

$$1808 O(Ld^2 + L^2d + L^2 + L^2d + Lmd + Lmd) \leq O(Lmd),$$

1809 where this step follows from m is the major term.

1810 **Proof of Part 2.** This is trivial. □

1813 F.2 TRAINING CONVERGENCE

1815 **Theorem F.2** (Formal version of Theorem 5.5). *Assuming Assumption 5.1, Definition 5.4 and*
1816 *Definition 5.2 hold, denote the failure probability $\delta \in (0, 0.1)$, then with a probability at least $1 - \delta$,*
1817 *we have:*

- 1818 • $(n, |\mathbb{B}|)$ -**Fixed Convergence.** *Denote $\alpha = \text{poly}(n, L, \exp(d), \frac{1}{\delta})$, we have:*

$$1820 \mathcal{L}(\mathbb{T}, \mathbb{D}) \leq \exp(-\frac{\varepsilon^2}{\alpha} MT).$$

- 1823 • $(n, |\mathbb{B}|)$ -**Dependent Convergence.** *Denote $N = \Theta(C^{\frac{1}{4}})$, $M = \Theta(C^{\frac{3}{4}})$, $T = \Theta(N)$, we have:*

$$1825 \mathcal{L}(\mathbb{N}, \mathbb{D}) \leq \exp(-\varepsilon^2 \omega C^{\frac{1}{4}}),$$

1826 *Proof. Proof of $(n, |\mathbb{B}|)$ -Fixed Convergence.* Following Lemma E.1, it is easy to have:

$$1828 \left| \frac{d}{dt} \mathcal{L}(t, \mathbb{D}) \right| \leq O\left(\frac{N}{\sqrt{m}}\right) \cdot \mathcal{L}(t, \mathbb{D}) \leq o\left(\frac{\sqrt{|\mathbb{B}|}}{\sqrt{nB}} \varepsilon^2 \omega \lambda N\right) \cdot \mathcal{L}(t, \mathbb{D}).$$

1831 where the second step follows from $m = \Omega(n/(\varepsilon^2 \omega \lambda)^2)$ and $|\mathbb{B}| = \min_{t \geq 0} |\mathbb{B}(t)| \geq 1$ is the
1832 minimum batch size.

1833 Then, following Hoeffding inequality (Lemma B.14), with a probability at least $1 - \delta$, we have:

$$1834 \left| \frac{d}{dt} \mathcal{L}(t, \mathbb{D}) - \mathbb{E}\left[\frac{d}{dt} \mathcal{L}(t, \mathbb{D})\right] \right|$$

$$\leq o\left(\frac{\sqrt{|\mathbb{B}|}}{\sqrt{nB}}\varepsilon^2\omega\lambda N\right) \cdot \sqrt{\frac{n}{|\mathbb{B}|}\log(1/\delta)} \cdot \mathcal{L}(t, \mathbb{D}) \leq o(\varepsilon^2\omega\lambda N) \cdot \mathcal{L}(t, \mathbb{D}).$$

Combining Lemma E.1, we have:

$$\frac{d}{dt}\mathcal{L}(t, \mathbb{D}) \leq -C\varepsilon^2\omega\lambda N \cdot \mathcal{L}(t, \mathbb{D}).$$

Solving this ODE, we get:

$$\mathcal{L}(T, \mathbb{D}) \leq \exp\left(-\frac{\varepsilon^2}{\alpha}MT\right),$$

where this step follows from $\omega = o\left(\frac{1}{\text{poly}(L, d, \frac{1}{\delta})}\right)$ and $m = \Omega(n, L, \exp(d), \frac{1}{\delta})$.

$(n, |\mathbb{B}|)$ -Dependent Convergence. Similarly, we have:

$$\mathcal{L}(T, \mathbb{D}) \leq \exp(-\varepsilon^2\omega N \cdot \frac{n}{|\mathbb{B}|}) \leq \exp(-\varepsilon^2\omega N) = \exp(-\varepsilon^2\omega C^{\frac{1}{4}}),$$

this inequality holds since m is $\Omega(n^3)$ and $\omega = o\left(\frac{1}{NLd^{2.5}B^3}\right)$ in Definition 5.2, $|\mathbb{B}| = \min_{t \geq 0} |\mathbb{B}(t)| \geq 1$, we choose $|\mathbb{B}| = 1$ ($T = N$) and $C = O(MN) = \Omega(N^4)$. \square

F.3 APPROXIMATION

Corollary F.3 (Formal version of Corollary 5.6). *Assuming we have arbitrary dataset size $N \in (0, +\infty)$. Let all scaling law factors be defined as Definition 5.4 and Assumption 5.1 and Definition 5.2 hold. Denote the failure probability $\delta \in (0, 0.1)$. We define $\alpha = O(Ld\sqrt{\log(1/\delta)})$. For the Good Model Class $\mathcal{F}_{M, T, N}(\mathbb{D}')$ with some $\mathbb{D}' \subset \mathcal{X} \times \mathcal{Y}$ and any function $F' : \mathcal{X} \rightarrow \mathcal{Y}$, arbitrary error $\epsilon > 0$ and compute cost C , with a probability at least $1 - \delta$, we have:*

$$\inf_{F \in \mathcal{F}_{M, T, N}(\mathbb{D}'), \mathbb{D}' \subset \mathcal{X} \times \mathcal{Y}} \mathbb{E}_{(X, Y) \sim \mathcal{D}} [\|F(X) - F'(X)\|_F^2] \leq \epsilon,$$

where $C = O(MN) = \Omega(256\alpha^8\epsilon^{-8} \wedge \varepsilon^{-8}\omega^4 \log(2\epsilon^{-1})^4)$, $M = \Omega(N^3)$, $T = N$.

Proof. First, we have (Part 14 of Lemma D.1):

$$\mathcal{L}(t, \mathbb{D}) \leq O(Ld).$$

We let $\mathbb{D}' := \{(X_i, F'(X_i)), X_i \sim \mathcal{X}\}_{i=1}^N$, then $\mathcal{F}_{M, T, N}(\mathbb{D}') = \{F(\cdot, \theta(T)), \theta(0) \sim \mathcal{N}(0, I_M), \theta(T) \in \mathcal{A}_{T, N}(\theta(0), \mathbb{D}')\}$, where we define $\mathcal{A}_{T, N}(\theta(0), \mathbb{D}') := \{\theta(0) + \int_0^T -\frac{d}{ds}\mathcal{L}(s, \mathbb{B}(s))ds, \mathbb{B}(s) \subseteq \mathbb{D}'\}$. Following Theorem F.2, we have:

$$\mathcal{L}(N, \mathbb{D}) \leq \exp(-\varepsilon^2\omega C^{\frac{1}{4}}).$$

We denote $\mathcal{R}_{\mathcal{F}_{M, T, N}(\mathbb{D})}(F) := \inf_{F \in \mathcal{F}_{M, T, N}(\mathbb{D})} \mathbb{E}_{(X, Y) \sim \mathcal{D}} [\|F(X) - F'(X)\|_F^2]$

Then we apply Hoeffding inequality (Lemma B.14) to each data point, we have:

$$|\mathcal{R}_{\mathcal{F}_{M, T, N}(\mathbb{D})}(F) - \mathcal{L}(N, \mathbb{D})| \leq O\left(\frac{Ld}{\sqrt{N}}\sqrt{\log(1/\delta)}\right). \quad (8)$$

Therefore, we have:

$$\begin{aligned} \mathcal{R}_{\mathcal{F}_{M, T, N}(\mathbb{D})}(F) &\leq O\left(\frac{Ld}{\sqrt{N}}\sqrt{\log(1/\delta)}\right) + \mathcal{L}(N, \mathbb{D}) \\ &\leq \frac{\alpha}{\sqrt{N}} + \exp(-\varepsilon^2\omega C^{\frac{1}{4}}), \end{aligned}$$

where $\alpha = O(Ld\sqrt{\log(1/\delta)})$.

Then we choose $C = \Omega(256\alpha^8\epsilon^{-8} \wedge \varepsilon^{-8}\omega^4 \log(2\epsilon^{-1})^4)$ and following Theorem F.2, then $M = \Omega(N^3)$, $T = N$.

\square

1890 G SCALING LAW

1891 G.1 GENERALIZATION BOUND OF EMPIRICAL RISK MINIMIZER

1892 **Definition G.1** (Van Der Vaart & Wellner (1996) and Yang & Barron (1999)). For a metric space
1893 (\mathcal{S}, d) and $\varepsilon > 0$, a finite set $\mathcal{S}' \subset \overline{\mathcal{S}}$ is called ε -covering if for any $x \in \mathcal{S}$ there exists $y \in \mathcal{S}'$ such
1894 that $d(x, y) \leq \varepsilon$, and the logarithm of the minimum cardinality of ε -covering is called covering
1895 ε -entropy and denoted by $V_{(\mathcal{S}, d)}(\varepsilon)$. Here, $\overline{\mathcal{S}}$ is the completion of \mathcal{S} with respect to the metric d .

1896 **Lemma G.2** (Lemma 11 of Schmidt-Hieber (2017), modified by Theorem 2.1 in Hayakawa (2019)
1897 and in the setting of this paper). We denote $\mathcal{F}_{M, T, N} = \cup_{\mathbb{D} \subset \mathcal{D}} \mathcal{F}_{M, T, N}(\mathbb{D})$. Then for any $\mathbb{D} \subset \mathcal{D}$, let
1898 F be the empirical risk minimizer taking values in $\mathcal{F}_{M, T, N}$. Suppose every element $F' \in \mathcal{F}_{M, T, N}$
1899 satisfies $d(F') := \max_{\ell \in [L]} \|F'_\ell\|_{L^\infty(\mathcal{X})} \leq B_F$ for some fixed $B_F > 0$. Then, for an arbitrary $\varepsilon > 0$,
1900 if $V_{(\mathcal{F}, d)}(\varepsilon) \geq 1$, then

$$1901 \|F - F^*\|_{L^F(\mathcal{X})}^2 \leq 4 \inf_{F' \in \mathcal{F}_{M, T, N}} \|F' - F^*\|_{L^F(\mathcal{X})}^2 + L \cdot O\left(\frac{B_F^2 V_{(\mathcal{F}_{M, T, N}, d)}(\varepsilon)}{n} + B_F \varepsilon\right).$$

1902 **Lemma G.3.** $\delta \in (0, 0.1)$. With a probability at least $1 - \delta$, we have:

- 1903 • **Part 1.** $\max_{\ell \in [L]} \|F_\ell\|_{L^\infty(\mathcal{X})} \leq O(\sqrt{\varepsilon}) =: B_F$.
- 1904 • **Part 2.** For any constant $\varepsilon > 0$, we can show that $1 \leq V_{(\mathcal{F}, d)}(\varepsilon) \leq O(d \cdot \log(1/\sqrt{\varepsilon}))$.
- 1905 • **Part 3.** For any $F' \in \mathcal{F}_{M', T', N'}$, we have:

$$1906 \sup_{X \in \mathcal{X}} \inf_{F' \in \mathcal{F}} \|F(X) - F'(X)\|_F^2 \leq O\left(\frac{1}{C^{\frac{3}{8}} \omega}\right)$$

1907 *Proof. Proof of Part 1.* Following Part 5 of Lemma D.1 and Eq. (1), we have:

$$1908 \max_{\ell \in [L], X \in \mathcal{X}} \|F_\ell(X)\|_2^2 \leq O(\varepsilon).$$

1909 We have:

$$1910 \|F_\ell\|_{L^\infty(\mathcal{X})} \leq \left(\max_{\ell \in [L]} \int_{\mathcal{X}} \|F_\ell(X)\|_\infty dX \right)^{\frac{1}{2}} \\ 1911 \leq \left(\frac{1}{L} \text{Volume}(\mathcal{X}) \cdot O(\varepsilon) \right)^{\frac{1}{2}} \\ 1912 \leq \left(\frac{\Theta(1)^{\frac{d}{2}}}{\frac{d!}{2}} \cdot O(\varepsilon) \right)^{\frac{1}{2}} \leq O(\sqrt{\varepsilon}),$$

1913 where the first step follows from the L^∞ norm, the second step follows from
1914 $\max_{\ell \in [L], X \in \mathcal{X}} \|F_\ell(X)\|_2^2 \leq O(\varepsilon)$, the third step follows from \mathcal{X} is a d -dimensional ball with
1915 radius $\Theta(1)$, the last step follows from simple algebras.

1916 **Proof of Part 2.** Following Part 1 of this Lemma, we have:

$$1917 \sup_{F, F' \in \mathcal{F}_{M, T, N}} \max_{\ell \in [L]} \|F_\ell - F'_\ell\|_{L^\infty(\mathcal{X})} \leq \sup_{F, F' \in \mathcal{F}_{M, T, N}} \max_{\ell \in [L]} (\|F_\ell\|_{L^\infty(\mathcal{X})} + \|F'_\ell\|_{L^\infty(\mathcal{X})}) \leq 2B_F,$$

1918 where the first step follows from simple algebras. Then we have the logarithm of the covering number:

$$1919 V_{(\mathcal{F}_{M, T, N}, d)}(\varepsilon) = \log\left(\left(\frac{2B_F}{\varepsilon}\right)^d\right) \leq O(d \cdot \log(1/\sqrt{\varepsilon})).$$

1920 **Proof of Part 3.** We choose $C = C' = \Omega(256\alpha^8 \varepsilon^{-8} \wedge \varepsilon^{-8} \omega^4 \log(2\varepsilon^{-1})^4)$, $M = M'$, $T = T'$ and
1921 $N = N'$ as Corollary F.3, let $N' = \Theta(1)$, $M' = N'^3$, then we have:

$$1922 m' = \Theta(C'^{\frac{3}{4}}) = \Omega(64\alpha^6 \varepsilon^{-6} \wedge \varepsilon^{-6} \omega^3 \log(2\varepsilon^{-1})^3).$$

Next,

$$\sup_{X \in \mathcal{X}} \inf_{F \in \mathcal{F}_{M,T,N}} \|F(X) - F'(X)\|_F^2 \leq O(\varepsilon^2 R),$$

where this step follows from Part 11 of Lemma D.1 and Eq (1).

Relax R . We combine Eq. (5) and Eq (7) to relax the bound on R , we get:

$$R = \frac{1}{\sqrt{m\varepsilon^2\omega\lambda N}} \leq \frac{\omega\lambda\delta}{nL^{1.5} \exp(O(dB))}.$$

Hence, we show:

$$\sup_{X \in \mathcal{X}} \inf_{F \in \mathcal{F}_{M,T,N}} \|F(X) - F'(X)\|_F^2 \leq O\left(\frac{1}{C^{\frac{3}{8}}\omega}\right),$$

where this step follows that λ is some fixed constant (Assumption 5.1) and $N = \Theta(1)$ and $M = C^{\frac{3}{4}}$ (Corollary F.3). \square

G.2 GENERAL PRETRAINING SCALING LAW

Theorem G.4 (Formal version of Theorem 6.1). *Let all pre-conditions hold as Corollary F.3. For any $\varepsilon \in (0, 1)$, with a probability at least $1 - \delta$, there exists:*

$$\inf_{\mathcal{F}_{M,T,N}(\mathbb{D})} \sup_{\mathbb{D} \in \mathcal{D}} \Delta \mathcal{R}(F) \leq \begin{cases} O\left(\frac{\alpha^2}{\omega C^{\frac{1}{8}}}\right) + O\left(\frac{\varepsilon \cdot d \cdot \log(1/\sqrt{\varepsilon})}{C^{\frac{1}{4}}}\right) + \varepsilon^{\frac{3}{2}}, & \varepsilon \geq \sqrt{\frac{\log(C^{\frac{1}{8}}/\alpha)}{C^{\frac{1}{4}}\omega}}, \\ \exp(-\varepsilon^2\omega C^{\frac{1}{4}}) + O\left(\frac{\varepsilon \cdot d \cdot \log(1/\sqrt{\varepsilon})}{C^{\frac{1}{4}}}\right) + \varepsilon^{\frac{3}{2}}, & \text{otherwise} \end{cases},$$

note that $M = \Omega(N^3)$, $T = N$, $C = O(MN) = \Omega(N^4)$.

Proof. We have:

$$\begin{aligned} \Delta \mathcal{R}(F) &= \mathcal{R}(F) - \mathcal{R}(F^*) \\ &= \mathbb{E}_{(X,Y) \sim \mathcal{D}} [\|F(X) - Y\|_F^2 - \|F^*(X) - Y\|_F^2] \\ &= \mathbb{E}_{(X,Y) \sim \mathcal{D}} [\|F(X) - F^*(X) - \Xi\|_F^2 - \|\Xi\|_F^2] \\ &= \mathbb{E}_{(X,Y) \sim \mathcal{D}} [\|F(X) - F^*(X)\|_F^2 - \text{vec}(F(X) - F^*(X))^\top \text{vec}(\Xi) + \|\Xi\|_F^2 - \|\Xi\|_F^2] \\ &= \mathbb{E}_{(X,Y) \sim \mathcal{D}} [\|F(X) - F^*(X)\|_F^2] \\ &= \frac{1}{\text{Volume}(\mathcal{X})} \|F - F^*\|_{L^F(\mathcal{X})}^2 \\ &\leq \Theta\left(\frac{1}{L}\right) \cdot \|F - F^*\|_{L^F(\mathcal{X})}^2 \end{aligned}$$

where the first step follows from the definitions of \mathcal{R} and $\Delta \mathcal{R}$, the second step follows from $Y = F^*(X) + \Xi$, the third step follows from simple algebras, the fourth step follows from $\mathbb{E}[\Xi] = \mathbf{0}_{L \times d}$, the fifth step follows from the definition of L^F norm, the last step follows from simple algebras that \mathcal{X} is a space summing L balls.

Bounding Approximation Error. We have:

$$\begin{aligned} &\inf_{F \in \mathcal{F}_{M,T,N}} \|F - F^*\|_{L^F(\mathcal{X})} \\ &\leq \inf_{F' \in \mathcal{F}_{M,T,N}(\mathbb{D}'), \mathbb{D}' \subset \mathcal{X} \times \mathcal{Y}} \|F'(X) - F^*(X)\|_{L^F(\mathcal{X})} + \inf_{F \in \mathcal{F}_{M,T,N}} \|F(X) - F'(X)\|_{L^F(\mathcal{X})} \\ &\leq \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\inf_{F' \in \mathcal{F}_{M,T,N}(\mathbb{D}'), \mathbb{D}' \subset \mathcal{X} \times \mathcal{Y}} \|F'(X) - F^*(X)\|_F + \inf_{F \in \mathcal{F}_{M,T,N}} \|F(X) - F'(X)\|_F \right] \\ &\leq \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\inf_{F' \in \mathcal{F}_{M,T,N}(\mathbb{D}'), \mathbb{D}' \subset \mathcal{X} \times \mathcal{Y}} \|F'(X) - F^*(X)\|_F \right] + O\left(\frac{1}{C^{\frac{3}{16}}\sqrt{\omega}}\right) \end{aligned}$$

$$\leq \sqrt{\epsilon} + O\left(\frac{1}{C^{\frac{3}{16}}\sqrt{\omega}}\right),$$

where the first step follows from Cauchy-Schwartz inequality, the second step follows that minima is smaller than the average, the third step follows from the Part 3 of Lemma G.3, the last step follows from Corollary F.3.

Condition 1. When $\epsilon \geq \sqrt{\frac{\log(C^{\frac{1}{8}}/\alpha)}{C^{\frac{1}{4}}\omega}}$, we have:

$$\epsilon = O\left(\frac{\alpha}{C^{\frac{1}{8}}}\right).$$

Then,

$$\inf_{F \in \mathcal{F}_{M,T,N}} \|F - F^*\|_{L^F(\mathcal{X})} \leq O\left(\frac{\alpha}{C^{\frac{1}{16}}}\right) + O\left(\frac{1}{C^{\frac{3}{16}}\sqrt{\omega}}\right) \leq O\left(\frac{\alpha}{C^{\frac{1}{16}}\sqrt{\omega}}\right).$$

Condition 2. When $\epsilon \leq \sqrt{\frac{\log(C^{\frac{1}{8}}/\alpha)}{C^{\frac{1}{4}}\omega}}$, we have:

$$\epsilon = \exp(-\epsilon^2\omega C^{\frac{1}{4}}),$$

where this step follows from $\epsilon = o(N^{-2})$ and $N = C^{\frac{1}{4}}$.

Then,

$$\inf_{F \in \mathcal{F}_{M,T,N}} \|F - F^*\|_{L^F(\mathcal{X})} \leq \exp(-\epsilon^2\omega C^{\frac{1}{4}}) + O\left(\frac{1}{C^{\frac{3}{16}}\sqrt{\omega}}\right) \leq \exp(-\epsilon^2\omega C^{\frac{1}{4}}),$$

where the second inequality follows from $\exp(-\epsilon^2\omega C^{\frac{1}{4}}) \leq O\left(\frac{\alpha}{C^{\frac{1}{8}}}\right)$ when $\epsilon \leq \sqrt{\frac{\log(C^{\frac{1}{8}}/\alpha)}{C^{\frac{1}{4}}\omega}}$.

Bounding Generalization Error. We have:

$$B_F \epsilon \leq \epsilon^{\frac{3}{2}},$$

where this step follows from Part 1 of Lemma G.3.

Next,

$$\frac{B_F^2 V_{(\mathcal{F}_{M,T,N}, d)}(\epsilon)}{N} \leq O\left(\frac{\epsilon \cdot d \cdot \log(1/\sqrt{\epsilon})}{N}\right) \leq O\left(\frac{\epsilon \cdot d \cdot \log(1/\sqrt{\epsilon})}{C^{\frac{1}{4}}}\right),$$

where the first step follows from Part 1 and 2 of Lemma G.3.

Result. We combine all results above and Theorem G.2, obtaining:

$$\inf_{\mathcal{F}_{M,T,N}(\mathbb{D})} \sup_{\mathbb{D} \in \mathcal{D}} \Delta \mathcal{R}(F) \leq \begin{cases} O\left(\frac{\alpha^2}{\omega C^{\frac{1}{8}}}\right) + O\left(\frac{\epsilon \cdot d \cdot \log(1/\sqrt{\epsilon})}{C^{\frac{1}{4}}}\right) + \epsilon^{\frac{3}{2}}, & \epsilon \geq \sqrt{\frac{\log(C^{\frac{1}{8}}/\alpha)}{C^{\frac{1}{4}}\omega}} \\ \exp(-\epsilon^2\omega C^{\frac{1}{4}}) + O\left(\frac{\epsilon \cdot d \cdot \log(1/\sqrt{\epsilon})}{C^{\frac{1}{4}}}\right) + \epsilon^{\frac{3}{2}}, & \text{otherwise} \end{cases}.$$

□

G.3 UPPER BOUND

Theorem G.5 (Formal version of Theorem 6.2). *Let all pre-conditions hold as Theorem 6.1. For any choice of ϵ , we let $C = \Omega(\epsilon^{-12})$ to offset the negative effect of the grokking coefficient ϵ , thus,*

$$\Pr \left[\inf_{\mathcal{F}_{M,T,N}(\mathbb{D})} \sup_{\mathbb{D} \in \mathcal{D}} \Delta \mathcal{R}(F) \asymp O\left(\frac{\alpha}{C^{\frac{1}{8}}}\right) \right] \geq 1 - \delta.$$

2052 *Proof.* We follow Eq. (8), for all $\mathbb{D} \in \mathcal{D}$, we can show that:

$$\begin{aligned} 2053 & \\ 2054 & \inf_{F \in \mathcal{F}_{M,T,N}(\mathbb{D})} \mathcal{R}(F) \leq \frac{\alpha}{\sqrt{N}} - \exp(-\varepsilon^2 \omega \mathbf{C}^{\frac{1}{4}}) \\ 2055 & \\ 2056 & \leq \frac{\alpha}{\sqrt{N}} - \exp(-\omega \mathbf{C}^{\frac{1}{12}}) \\ 2057 & \\ 2058 & \leq O\left(\frac{\alpha}{\sqrt{N}}\right) \leq O\left(\frac{\alpha}{\mathbf{C}^{\frac{1}{8}}}\right), \\ 2059 & \end{aligned}$$

2060 where these steps follow from simple algebras and $N = o(\mathbf{C}^{-4})$.

2061 At the same time, the upper bound is trivial.

2062 \square

2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

2106 H LLMs USAGE DISCLOSURE
2107

2108 LLMs were used only to polish language, such as grammar and wording. These models did not
2109 contribute to idea creation or writing, and the authors take full responsibility for this paper's content.
2110

2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159