TRANSFORMERS ARE DEEP OPTIMIZERS: PROVABLE IN-CONTEXT LEARNING FOR DEEP MODEL TRAINING

Anonymous authors

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

032

Paper under double-blind review

ABSTRACT

We investigate the transformer's capability for in-context learning (ICL) to simulate the training process of deep models. Our key contribution is providing a positive example of using a transformer to train a deep neural network by gradient descent in an implicit fashion via ICL. Specifically, we provide an explicit construction of a (2N + 4)L-layer transformer capable of simulating L gradient descent steps of an N-layer ReLU network through ICL. We also give the theoretical guarantees for the approximation within any given error and the convergence of the ICL gradient descent. Additionally, we extend our analysis to the more practical setting using Softmax-based transformers. We validate our findings on synthetic datasets for 3layer, 4-layer, and 6-layer neural networks. The results show that ICL performance matches that of direct training.

022 1 INTRODUCTION

We study transformers' ability to simulate the training process of deep models. This analysis is not only practical but also timely. On one hand, transformers and deep models (Brown, 2020; Radford et al., 2019) are so powerful, popular and form a new machine learning paradigm — foundation models. These large-scale machine learning models, trained on vast data, provide a general-purpose foundation for various tasks with minimal supervision (Team et al., 2023; Touvron et al., 2023; Zhang et al., 2022). On the other hand, the high cost of pretraining these models often makes them prohibitive outside certain industrial labs (Jiang et al., 2024; Bi et al., 2024; Achiam et al., 2023). In this work, we aim to advance the "one-for-many" modeling philosophy of foundation model paradigm (Bommasani et al., 2021) by considering the following research problem:

Question 1. Is it possible to train one deep model with the ICL of another foundation model?

The implication of Question 1 is profound: if true, one foundation model could lead to many others without pertaining. In this work, we provide an affirmative example for Question 1. Specifically, we show that transformer models are capable of simulating the training of a deep ReLU-based feed-forward neural network with provable guarantees through In-Context Learning (ICL). Our analysis assumes that we have well-pretrained the transformer using the data generated by the deep network. We require the deep network to maintain consistent hyperparameters (e.g., model width and depth) during the pretraining and testing. However, during the testing, we vary the parameter distribution and input data distribution of the deep network to generate data for the transformer.

041 In ICL, the models learn to solve new tasks during inference by using task-specific examples provided 042 as part of the input prompt, rather than through parameter updates (Wei et al., 2023; Bubeck et al., 043 2023; Achiam et al., 2023; Bai et al., 2023; Min et al., 2022; Garg et al., 2022; Brown, 2020). Unlike 044 standard supervised learning, ICL enables models to adapt to new tasks during inference using only the provided examples. In this work, the new task of our interest is algorithmic approximation via ICL (Bai et al., 2023; Zhang et al., 2023; Wang et al.). Specifically, we aim to use transformer's ICL 046 capability to replace/simulate the standard supervised training algorithms for N-layer networks. To 047 be concrete, we formalize the learning problem of how transformers learn (i) a given function and (ii) 048 a machine learning algorithm (e.g., gradient descent) via ICL, following (Bai et al., 2023). 049

(i) ICL for Function f. Let $f : \mathbb{R}^d \to \mathbb{R}$ be the function of our interest. Suppose we have a dataset $\mathcal{D}_n := \{(x_i, y_i)\}_{i \in [n]}$, where $\{x_i\}_{i \in [n]} \subseteq \mathbb{R}^d$ and $\{y_i\}_{i \in [n]} \subseteq \mathbb{R}$ are the input and output of f, respectively. Let x_{n+1} be the test input. The goal of ICL is to use a transformer, denoted by \mathcal{T} , to predict y_{n+1} based on the test input and the in-context dataset autoregresively: $\hat{y}_{n+1} \sim \mathcal{T}(\mathcal{D}_n, x_{n+1})$. The goal is for the prediction \hat{y}_{n+1} to be close to $y_{n+1} = f(x)$. 054 (ii) ICL for Gradient Descent of a Parametrized Model $f(w, \cdot)$. Bai et al. (2023) generalize (i) 055 to include algorithmic approximations of Gradient Descent (GD) training algorithms and explore 056 how transformers simulate gradient descent during inference without parameter updates. They term 057 the simulated GD algorithm "In-Context Gradient Descent (ICGD)." In essence, ICGD enables 058 transformers to approximate gradient descent on a loss function $L_n(w)$ for a parameterized model $f(w, \cdot)$ based on a dataset \mathcal{D}_n . Traditional gradient descent updates w iteratively as $w_{t+1} = w_t - w_t$ $\eta \nabla L_n(w_t)$. In contrast, ICGD uses a transformer \mathcal{T} to simulate these updates within a forward 060 pass. Given example data \mathcal{D}_n and test input x_{n+1} , the transformer performs gradient steps in an 061 implicit fashion by inferring parameter updates through its internal representations, using input 062 context without explicit weight changes. Please see Section 2 for explicit formulation. 063

In this work, we investigate the case where $f(w, \cdot)$ is a deep feed-forward neural network. We defer the detailed problem setting to Section 2. In comparison to standard ICGD (Bai et al., 2023), ICGD for deep feed-forward networks is not trivial. This is due to two technical challenges:

067 (C1) Analytical feasibility of gradient computation for these thick networks.

068 (C2) Explicit construction capable of approximating ICGD for such layers and their gradients.

A work similar to ours is (Wang et al.). It demonstrates that the transformer implements multiple steps of ICGD on deep neural networks. However, it requires more layers in the transformer model and fails to consider the Softmax-transformer. We provide a more detailed comparison in Appendix B.1.
 To this end, we present the first explicit expression for gradient computation of *N*-layer feed-forward network (Lemma 1). Importantly, its term-by-term tractability provides key insights for the detailed construction of a specific transformer to train this network via ICGD (Theorem 1).

Contributions. We offer a positive early investigation of Question 1. Our contributions are threefold:

- Approximation by ReLU-Transformer. For simplicity, we begin with the ReLU-based transformer. For a broad class of smooth empirical risks, we construct a (2N + 4)L-layer transformer to approximate *L* steps of in-context gradient descent on the *N*-layer feed-forward networks with the same input and output dimensions (Theorem 1). We then extend this to accommodate varying dimensions (Theorem 5). We also provide the theoretical guarantees for the approximation within any given error (Corollary 1.1) and the convergence of the ICL gradient descent (Lemma 14).
- Approximation by Softmax-Transformer. We extend our analysis to the Softmax-transformer to better reflect realistic applications. The key technique is to ensure a qualified approximation error at each point to achieve universal approximation capabilities of the Softmax-based Transformer (Lemma 16). We give a construction of a 4*L*-layer Softmax transformer to approximate *L* steps of gradient descent, and guarantee the approximation and the convergence (Theorem 6).
- Experimental Validation. We validate our theory with ReLU- and Softmax-transformers, specifically, ICGD for the *N*-layer networks (Theorem 1, Theorem 5, and Theorem 6). We assess the ICL capabilities of transformers by training 3-, 4-, and 6-layer networks in Appendix G. The numerical results show that the performance of ICL matches that of training *N*-layer networks.

Organization. We show our main results in Theorem 1 and the informal version in Appendix A.
 Section 2 presents preliminaries. Section 3 presents the problem setup and ICL approximation to
 GD steps of *N*-layer feed-forward network based on ReLU-transformer. The appendix includes the
 related works (Appendix B.1), the detailed proofs of the main text (Appendix D), ICL approximation
 to GD steps of *N*-layer network based on Softmax-transformer (Appendix F), the experimental
 results (Appendix G), and the application to diffusion models (Appendix H).

Notations. We use lower case letters to denote vectors and upper case letters to denote matrices. The index set $\{1, ..., I\}$ is denoted by [I], where $I \in \mathbb{N}^+$. For any matrices $A \in \mathbb{R}^{n \times n}$, let ℓ_p norm of *A* be induced by vector ℓ_p -norm, defined as $||A||_p := \sup\{||Ax||_p : x \in \mathbb{R}^n \text{ with } ||x||_p = 1\}$. For any function *f* and distribution *P*, we denote $L^2(P)$ norm of *f* as $||f||_{L^2(P)} = \mathbb{E}_P^{1/2}[||f||_2^2]$. We use A[i, j] to denote the element in *i*-th row and *j*-th column of matrix *A*. For any matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{m \times n}$, let \odot denotes the Hadamard product: $(A \odot B)[i, j] := A[i, j] \cdot B[i, j]$. For any matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$, let \otimes denote the Kronecker product:

$$A \otimes B := \begin{bmatrix} A[1,1]B & \cdots & A[1,n]B \\ \vdots & \ddots & \vdots \\ A[m,1]B & \cdots & A[m,n]B \end{bmatrix}.$$

¹⁰⁶ 107 2 PRELIMINARIES: ICL AND ICGD

104 105

We present the ideas we built upon: In-Context Gradient Descent (ICGD).

(i) ICL for Function f. Let $f : \mathbb{R}^d \to \mathbb{R}$ be the function of our interest. Suppose we have a dataset $\mathcal{D}_n := \{(x_i, y_i)\}_{i \in [n]}, \text{ where } \{x_i\}_{i \in [n]} \subseteq \mathbb{R}^d \text{ and } \{y_i\}_{i \in [n]} \subseteq \mathbb{R} \text{ are the input and output of } f,$ respectively. Let x_{n+1} be the test input. The goal of ICL is to use a transformer, denoted by \mathcal{T} , to predict y_{n+1} based on the test input and the in-context dataset autoregresively: $\hat{y}_{n+1} \sim \mathcal{T}(\mathcal{D}_n, x_{n+1})$. For convenience in our analysis, we adopt the ICL notation from (Bai et al., 2023). Specifically, we shorthand (\mathcal{D}_n, x_{n+1}) into an input sequence (i.e., prompt) of length n + 1 and represent it as a compact matrix $H \in \mathbb{R}^{D \times (n+1)} := [h_1, \dots, h_{n+1}]$ in the form:

- 115 116
- 117

$$H \coloneqq \begin{bmatrix} x_1 & x_2 & \cdots & x_n & x_{n+1} \\ y_1 & y_2 & \cdots & y_n & 0 \\ q_1 & q_2 & \cdots & q_n & q_{n+1} \end{bmatrix} \in \mathbb{R}^{D \times (n+1)}, \quad q_i \coloneqq \begin{bmatrix} 0_{D-(d+3)} \\ 1 \\ t_i \end{bmatrix} \in \mathbb{R}^{D-(d+1)}.$$
(2.1)

Here, we choose $D := \dim x_i + \dim y_i + \dim q_i = \Theta(d)$. We use q_i to fill in the remain D - (d+1)entries in addition to $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$. The last entry $t_i := \mathbb{1}(i < n+1)$ of q_i is the position indicator to distinguish the n in-context examples and the test data. The **problem of "ICL for** f" is to show the existence of a transformer \mathcal{T} that, when given H, outputs $\mathcal{T}(H) \in \mathbb{R}^{D \times (n+1)}$ of the same shape, and the "(d+1, n+1) entry of $\mathcal{T}(H)$ " provides the prediction \hat{y}_{n+1} . The goal is for the prediction \hat{y}_{n+1} to be close to $y_{n+1} = f(x)$ measured by some proper loss.

(ii) ICL for Gradient Descent of a Parametrized Model $f(w, \cdot)$. In this work, we aim to use ICL to replace/simulate the standard supervised training procedure for *N*-layer neural networks. To achieve this, we introduce the concept of In-Context Gradient Descent (ICGD) for a parameterized model.

Consider a machine learning model $f(w, \cdot) : \mathbb{R}^{D_w} \times \mathbb{R}^d \to \mathbb{R}^d$, parametrized by $w \in \mathbb{R}^{D_w}$. Given a dataset $\mathcal{D}_n := \{(x_i, y_i)\}_{i \in [n]} \stackrel{\text{iid}}{\sim} \mathbb{P}$, a typical learning task is to find parameters w^* such that $f(w^*, \cdot)$ becomes closest to the true data distribution \mathbb{P} . Then, for any test input x_{n+1} , we predict: $\widehat{y}_{n+1} = f(w^*, x_{n+1})$. To find w^* , Bai et al. (2023) configure a transformer to implement gradient descent on $f(w, \cdot)$ through ICL, simulating optimization algorithms during inference without explicit parameter updates. We formalize this **In-Context Gradient Descent (ICGD)** problem: using a pretrained model to simulate gradient descent on $f(w, \cdot)$ w.r.t. the provided context (\mathcal{D}_n, x_{n+1}) .

Problem 1 (In-Context Gradient Descent (ICGD) on Model $f(w, \cdot)$ (Bai et al., 2023)). Let $\epsilon > 0$ and $L \ge 1$. Consider a machine learning model $f(w, x) : \mathbb{R}^{D_w} \times \mathbb{R}^d \to \mathbb{R}^d$ parameterized by $w \in \mathbb{R}^{D_w}$. Given a dataset $\mathcal{D}_n := \{(x_i, y_i)\}_{i \in [n]} \stackrel{\text{iid}}{\sim} \mathbb{P}$ with $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}^d$, define the empirical risk function:

$$\mathcal{L}_n(w) \coloneqq \frac{1}{2n} \sum_{i=1}^n \ell(f(w, x_i), y_i), \quad \text{where } \ell : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R} \text{ is a loss function.}$$
(2.2)

Let $W \subseteq \mathbb{R}^{D_w}$ be a closed domain, and Proj_W denote the projection onto W. The problem of "ICGD on model $f(w, \cdot)$ " is to find a transformer \mathcal{T} with L blocks, each approximating one step of gradient descent using T layers. For any input $H^{(0)} \in \mathbb{R}^{D \times (n+1)}$ in the form of (2.1), the transformer $\mathcal{T}(H^{(0)})$ approximates L steps of gradient descent. Specifically, for $l \in [L]$ and $i \in [n+1]$, the output at layer Tl is: $h_i^{(Tl)} = [x_i; y_i; \overline{w}^{(l)}; \mathbf{0}; 1; t_i]$, where, with $\overline{w}^{(0)} = \mathbf{0}$,

$$\overline{w}^{(l)} = \operatorname{Proj}_{\mathcal{W}} \left(\overline{w}^{(l-1)} - \eta \left(\nabla \mathcal{L}_n(\overline{w}^{(l-1)}) + \epsilon^{(l-1)} \right) \right) \text{ is updated recursively,}$$
(2.3)

and $\|\epsilon^{(l-1)}\|_2 \leq \epsilon$ represents the approximation error at step l-1.

Problem 1 aims to find a transformers \mathcal{T} to perform L steps gradient descent on loss $\mathcal{L}_n(w)$ in an implicit fashion (i.e., no explicit parameter update). More precisely, Bai et al. (2023) configure \mathcal{T} with L identical blocks, each approximating one gradient descent step using T layers. In this work, we investigate the case where $f(w, \cdot)$ is an "N-layer neural network."

Transformer. We defer standard definition of transformer to Appendix C.1 due to the page limit.

156 157 158

140 141 142

143

144

145

146

151

152

153

154

155

3 IN-CONTEXT GRADIENT DESCENT ON N-LAYER NEURAL NETWORKS

We now show that transformers is capable of implementing gradient descent on *N*-layer neural networks through ICL. In Section 3.1, we define the *N*-layer ReLU neural network and state its ICGD problem. In Section 3.2, we derive explicit gradient descent expression for *N*-layer NN. In Section 3.3, we show how transformers execute gradient descent on *N*-layer NN via ICL.

176

177 178

179 180 181

182

183

185 186 187

189

190

191

192

194

195

197 198

204 205

206

207

208

162 3.1 PROBLEM SETUP: ICGD FOR N-LAYER NEURAL NETWORK 163

To begin, we introduce the construction of our N-Layer Neural Network which we aims to implement 164 gradient descent on its empirical loss function. 165

166 **Definition 1** (N-Layer Neural Network). An N-Layer Neural Network comprises N - 1 hidden 167 layers and 1 output layer, all constructed similarly. Let $r : \mathbb{R} \to \mathbb{R}$ be the activation function. For the hidden layers: for any $i \in [n+1], j \in [N-1]$, and $k \in [K]$, the output for the first j layers w.r.t. input $x_i \in \mathbb{R}^d$, denoted by $\operatorname{pred}_h(x_i; j) \in \mathbb{R}^K$, is defined as recursive form: 168 169

$$\operatorname{pred}_{h}(x_{i};1)[k] := r(v_{1_{k}}^{\top}x_{i}), \text{ and } \operatorname{pred}_{h}(x_{i};j)[k] := r(v_{j_{k}}^{\top}\operatorname{pred}_{h}(x_{i};j-1)),$$

171 where $v_{1_k} \in \mathbb{R}^d$ and $v_{j_k} \in \mathbb{R}^K$ for $j \in \{2, ..., N-1\}$ are the k-th parameter vectors in the first layer and the j-th layer, respectively. For the output layer (N-th layer), the output for the first N 172 173 layers (i.e the entire neural network) w.r.t. input $x_i \in \mathbb{R}^d$, denoted by $\operatorname{pred}_o(x_i; w, N) \in \mathbb{R}^d$, is 174 defined for any $k \in [d]$ as follows: 175

$$\operatorname{pred}_{o}(x_{i}; w, N)[k] := r(v_{N_{k}}^{\top} \operatorname{pred}_{h}(x_{i}; N-1)), \qquad (3.1)$$

where $v_{N_k} \in \mathbb{R}^K$ are the k-th parameter vectors in the N-th layer and $w \in \mathbb{R}^{2dK+(N-2)K^2}$ denotes the vector containing all parameters in the neural network,

$$w := \left[v_{1_1}^{\top}, \dots, v_{1_K}^{\top}, \dots, v_{j_k}^{\top}, \dots, v_{N-1_1}^{\top}, \dots, v_{N-1_K}^{\top}, v_{N_1}^{\top}, \dots, v_{N_d}^{\top} \right]^{\top}.$$
 (3.2)

Remark 1 (Prediction Function for *j*-th layer on *i*-th Data: $p_i(j)$). For simplicity, we abbreviate the output from the first j-th layer of the N-layer neural networks NN with input x_i as $p_i(j)$,

$$p_i(j) := \begin{cases} x_i \in \mathbb{R}^d, & \text{for } j = 0, \\ \text{pred}_h(x_i; j) \in \mathbb{R}^K, & \text{for } j \in [N-1], \\ \text{pred}_o(x_i; w, N) \in \mathbb{R}^d, & \text{for } j = N. \end{cases}$$
(3.3)

Additionally, we define $p_i := [p_i(1); \ldots; p_i(N)] \in \mathbb{R}^{(N-1)K+d}$. 188

We formalize the problem of using a transformer to simulate gradient descent algorithms for training the N-layer NN defined in Definition 1, by optimizing loss (2.2). Specifically, we consider the ICGD (Problem 1) with the parameterized model $f(w, \cdot) := \text{pred}_o(\cdot; w, N)$.

Problem 2 (ICGD on N-Layer Neural Networks). Let the N-layer neural networks, activation 193 function r, and prediction function $p_i(j)$ for all layers follow Definition 1 and Remark 1. Assume we under the identical setting as Problem 1, considering model $f(w, \cdot) := \text{pred}_o(\cdot; w, N)$ and specifying \mathcal{W} is a closed domain such that for any $j \in [N-1]$ and $k \in [K]$, 196

$$\mathcal{W} \subset \left\{ w = [v_{j_k}] \in \mathbb{R}^{D_N} : \|v_{j_k}\|_2 \le B_v \right\}.$$
(3.4)

The problem of "ICGD on N-layer neural networks" is to find a TL layers transformer \mathcal{T} , capable of 199 implementing L steps gradient descent as in Problem 1. Specifically, for any $L \in [L], i \in [n+1]$, the *Tl*-th layer outputs $h_i^{(Tl)} = [x_i; y_i; \overline{w}^{(l)}; \mathbf{0}; 1; t_i]$, where:

$$\overline{w}^{(l)} = \operatorname{Proj}_{\mathcal{W}} \left(\overline{w}^{(l-1)} - \eta \left(\nabla \mathcal{L}_n(\overline{w}^{(l-1)}) + \epsilon^{(l-1)} \right) \right), \quad \overline{w}^{(0)} = \mathbf{0},$$
(3.5)

and $\|\epsilon^{(l-1)}\|_2 \leq \epsilon$ is the error term generated from the approximation in the *l*-th step.

Remark 2 (Necessary for bounded domain \mathcal{W}). For using a sum of ReLU to approximate functions like r, which illustrated in the consequent section, we need to avoid gradient exploding. Therefore, we require \mathcal{W} to be a bounded domain, and utilize $\operatorname{Proj}_{\mathcal{W}}$ to project w into bounded domain \mathcal{W} .

3.2 EXPLICIT GRADIENT DESCENT OF N-LAYER NEURAL NETWORK 209

210 Intuitively, Problem 2 asks whether there exists a transformer capable of simulating the gradient 211 descent algorithm on the loss function of an N-layer neural network. We answer Problem 2 by 212 providing an explicit construction for such a transformer \mathcal{T} in Theorem 1. To facilitate our proof, we 213 first introduce the necessary notations for explicit expression of the gradient $\nabla_w \mathcal{L}_n(w)$.

214 **Definition 2** (Abbreviations). Fix $i \in [n+1]$, and consider an N-layer neural network with activation 215 function r and prediction function $p_i(j)$ as defined in Definition 1.

• Let $D_j \in \mathbb{R}$ denote the total number of parameters in the first j layers. By (3.2), we have:

$$D_j = \begin{cases} 0, & j = 0\\ dK, & j = 1\\ (j-1)K^2 + dK, & 2 \le j \le N - 1\\ (N-2)K^2 + 2dK, & j = N. \end{cases}$$

The parameter vector w := [v₁₁^T,..., v_{1K}^T,..., v_{N-11}^T,..., v_{N-1K}^T, v_{N1}^T,..., v_{Nd}^T]^T follows (3.2). Define φ_i := (∂ℓ(p_i(N),y_i) → ∂p_i(N)/∂w)^T ∈ ℝ^{D_N}. For any j ∈ [N], let A_i(j) denote the derivative of ℓ(p_i(N), y_i) with respect to the parameters in the j-th layer:A_i(j) = φ_i[D_{j-1} : D_j], where φ_i[a : b] selects elements from the a-th to b-th position in φ_i.
For a structure parameters with a structure parameters are parameters and parameters are parameters.

• For activation function r(t), let r'(t) be its derivative. Define $r'_i(j) \in \mathbb{R}^K$ as:

$$r'_i(j)[k] := r'(v_{j+1_k}^{\top} p_i(j)).$$

• Define $r'_i := [r'_i(0); ...; r'_i(N-1)]$ and $R_i(j)$ as:

$$R_{i}(j) := \begin{cases} \operatorname{diag}\{r'(v_{j+1,1}^{\top}p_{i}(j)), \dots, r'(v_{j+1,K}^{\top}p_{i}(j))\} \in \mathbb{R}^{K \times K}, & j \in \{0, 1, \dots, N-2\}\\ \operatorname{diag}\{r'(v_{j+1,1}^{\top}p_{i}(j)), \dots, r'(v_{j+1,d}^{\top}p_{i}(j))\} \in \mathbb{R}^{d \times d}, & j = N-1. \end{cases}$$

• For any $j \in [N]$, let V_j denote the parameters in the j-th layer as:

$$V_j := \begin{cases} \begin{bmatrix} v_{1_1}, \dots, v_{1_K} \end{bmatrix}^\top \in \mathbb{R}^{K \times d}, & j = 1 \\ \begin{bmatrix} v_{j_1}, \dots, v_{j_K} \end{bmatrix}^\top \in \mathbb{R}^{K \times K}, & j \in 2, \dots, N-1 \\ \begin{bmatrix} v_{N_1}, \dots, v_{N_d} \end{bmatrix}^\top \in \mathbb{R}^{d \times K}, & j = N. \end{cases}$$

Definition 2 splits the gradient of $\mathcal{L}_n(w)$ into N parts. This makes $\nabla_w \mathcal{L}_n(w)$ more interpretable and tractable, since all parts follows a recursion formula according to chain rule. With above notations, we calculate the gradient descent step (3.5) of N-layer neural network as follows:

Lemma 1 (Decomposition of One Gradient Descent Step). Fix any $B_v, \eta > 0$. Suppose loss function $\mathcal{L}_n(w)$ on n data points $\{(x_i, y_i)\}_{i \in [n]}$ follows (2.2). Suppose closed domain \mathcal{W} and projection function $\operatorname{Proj}_{\mathcal{W}}(w)$ follows (3.4). Let $A_i(j), r'_i(j), R_i(j), V_j$ be as defined in Definition 2. Then the explicit form of gradient $\nabla \mathcal{L}_n(w)$ becomes

 $\nabla \mathcal{L}_n(w) = \frac{1}{2n} \sum_{i=1}^n \begin{bmatrix} A_i(1) \\ \vdots \\ A_i(N) \end{bmatrix},$ (3.6)

where $A_i(j)$ denote the derivative of $\ell(p_i(N), y_i)$ with respect to the parameters in the *j*-th layer,

$$A_{i}(j) = \begin{cases} (R_{i}(N-1) \cdot V_{N} \cdot \ldots \cdot R_{i}(j-1) \cdot \left[\mathbf{I}_{K \times K} \otimes p_{i}(j-1)^{\top}\right])^{\top} \cdot \left(\frac{\partial \ell(p_{i}(N), y_{i})}{\partial p_{i}(N)}\right)^{\top}, & j \neq N\\ (R_{i}(N-1) \cdot \left[\mathbf{I}_{d \times d} \otimes p_{i}(N-1)^{\top}\right])^{\top} \cdot \left(\frac{\partial \ell(p_{i}(N), y_{i})}{\partial p_{i}(N)}\right)^{\top}, & j = N. \end{cases}$$

Proof Sketch. Using the chain rule and product rule, we decompose the gradient as follows: $\nabla_w \mathcal{L}_n(w) = \frac{1}{2n} \sum_{i=1}^N \left[\frac{\partial p_i(N)}{\partial w} \right]^\top \cdot \left[\frac{\partial \ell(p_i(N), y_i)}{\partial p_i(N)} \right]^\top$. Thus, we only need to compute $\frac{\partial p_i(N)}{\partial w}$. By Definition 1 and the chain rule, we prove that $\frac{\partial p_i(N)}{\partial w}$ satisfies the recursive formulation (D.4). Combining these, we derive the explicit form of gradient $\nabla_w \mathcal{L}_n(w)$, and the gradient step follows directly. Please see Appendix D.1 for a detailed proof.

Since it is hard to calculate the elements in $A_i(j)$ in a straightforward mannar, we calculate each parts of it successively. Specifically, we define the intermediate terms $s_i(j)$ and u as follows

Definition 3 (Definition of intermediate terms). Let $A_i(j), r'_i(j), R_i(j), V_j$ be as defined in Definition 2. By Lemma 1, the derivative of $\ell(p_i(N), y_i)$ w.r.t the parameters in the *j*-th layer follows,

$$A_{i}(j) = \begin{cases} (R_{i}(N-1) \cdot V_{N} \cdot \ldots \cdot R_{i}(j-1) \cdot \left[\mathbf{I}_{K \times K} \otimes p_{i}(j-1)^{\top}\right])^{\top} \cdot \left(\frac{\partial \ell(p_{i}(N), y_{i})}{\partial p_{i}(N)}\right)^{\top}, & j \neq N \\ (R_{i}(N-1) \cdot \left[\mathbf{I}_{d \times d} \otimes p_{i}(N-1)^{\top}\right])^{\top} \cdot \left(\frac{\partial \ell(p_{i}(N), y_{i})}{\partial p_{i}(N)}\right)^{\top}, & j = N. \end{cases}$$

271

272

273 274 275

276 277

278

279

281 282

283

284

285

286

287

288

289

295

296

297

298

299

300

301

302

303

304

305

306

307

308

310

312

313

314 315 316

320 321

For any $t, y \in \mathbb{R}^d$, we define vector function $u(t, y) := (\frac{\partial \ell(t, y)}{\partial t})^\top : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}^d$. Moreover, for any $j \in [N], i \in [n+1]$, we define $s_i(j)$ as $s_i(j) := \begin{cases} R_i(j-1)V_{j+1}^\top \dots R_i(N-2)V_N^\top \cdot R_i(N-1) \cdot u(p_i(N), y_i), & \in \mathbb{R}^K, j \neq N, \\ R_i(N-1) \cdot u(p_i(N), y_i), & \in \mathbb{R}^d, j = N. \end{cases}$ Let \odot denotes hadamard product. For any $j \in [N-1], i \in [N+1]$, Definition 3 leads to $s_i(j) = r'_i(j-1) \odot (V_{j+1}^{\top} \cdot s_i(j+1)),$ Moreover, by Definition 3, it holds $A_i(j) = \begin{cases} \left[\mathbf{I}_{K \times K} \otimes p_i(j-1) \right] \cdot s_i(j), & j \neq N, \\ \left[\mathbf{I}_{d \times d} \otimes p_i(N-1) \right] \cdot s_i(N), & j = N. \end{cases}$

TRANSFORMERS APPROXIMATE GRADIENT DESCENT OF N-LAYER NEURAL NETWORKS 3.3 WITH ICL

(3.7)

(3.8)

For using neural networks to approximate (2.2), which contains smooth functions changeable, we need to approximate these smooth functions by simple combination of activation functions. Our key approximation theory is using a sum of ReLUs to approximate any smooth function (Bai et al., 2023).

Definition 4 (Approximability by Sum of ReLUs, Definition 12 of (Bai et al., 2023)). Let $z \in \mathbb{R}^k$. We say that a function $g: \mathbb{R}^k \to \mathbb{R}$ is $(\epsilon_{approx}, R, H, C)$ -approximable by sum of ReLUs if there exist a "(H, C)-sum of ReLUs" function $f_{H,C}(z)$ defined as

$$f_{H,C}(z) = \sum_{h=1}^{H} c_h \sigma(a_h^{\top}[z;1]) \quad \text{with} \quad \sum_{h=1}^{H} |c_h| \le C, \quad \max_{h \in [H]} ||a_h||_1 \le 1, \quad a_h \in \mathbb{R}^{k+1}, \quad c_h \in \mathbb{R},$$

such that $\sup_{z \in [-R,R]^k} |g(z) - f_{H,C}(z)| \le \epsilon_{\text{approx}}$.

Overview of Our Proof Strategy. Lemma 1 and Definition 4 motivate the following strategy: term-by-term approximation for our gradient descent step (3.6).

- **Step 1.** Given (x_i, w) , we use N attention layers to approximate the output of the first j layers with input $x_i, p_i(j) \coloneqq \text{pred}_h(x_i; j) \in \mathbb{R}^k$ (Definition 1) for any $j \in [N]$. Then we use 1 attention layer to approximate chain-rule intermediate terms $r'_i(j-1)[k] := r'(v_{j_k} p_i(j-1))$ (Definition 2) for any $i \in [n], j \in [N]$ and $k \in [K]$: Lemma 2 and Lemma 3.
- **Step 2.** Given (r'_i, p_i, w) , we use an MLP layer to approximate $u(p_i(N), y_i)$ (Definition 3), for $i \in [n]$, and use N element-wise multiplication layers to approximate $s_i(j)$ (Definition 3), for any $j \in [N]$: Lemma 4 and Lemma 5. Moreover, Lemma 6 shows the closeness result for approximating $s_i(j)$, which leads to the final error accumulation in Theorem 1.
- **Step 3.** Given $(p_i, r'_i, g_i s_i(j), w)$, we use an attention layer to approximate $w \eta \nabla \mathcal{L}_n(w)$. Then we use an MLP layer to approximate $\operatorname{Proj}_{\mathcal{W}}(w)$. And implementing L steps gradient descent by a (2N+4)L-layer neural network NN $_{\theta}$ constructed based on Step 1 and 2. Finally, we arrive our main result: Theorem 1. Furthermore, Lemma 14 shows closeness results to the true gradient descent path.

Step 1. We start with approximation for $p_i(j)$. 311

Lemma 2 (Approximate $p_i(j)$). Let upper bounds $B_v, B_x > 0$ such that for any $k \in [K], j \in$ [N] and $i \in [n]$, $||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. For any $j \in [N]$, $i \in [n]$, define

$$B_r^j := \max_{|t| \le B_v B_r^{j-1}} |r(t)|, \quad B_r^0 := B_x, \text{ and } B_r := \max_j B_r^j.$$

Let function r(t) be $(\epsilon_r, R_1, M_1, C_1)$ -approximable for $R_1 = \max\{B_v B_r, 1\}, M_1 \leq \mathcal{O}(C_1^2 \epsilon_r^{-2}),$ 317 where C_1 depends only on R_1 and the C^2 -smoothness of r. Then, for any $\epsilon_r > 0$, there exist N 318 attention layers $\operatorname{Attn}_{\theta_1}, \ldots, \operatorname{Attn}_{\theta_N}$ such that for any input $h_i \in \mathbb{R}^D$ takes from (2.1), they map 319

$$h_i = [x_i; y_i; w; \bar{p}_i(1); \dots; \bar{p}_i(j-1); \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{Attn}_{\theta_j}} \widetilde{h_i} = [x_i; y_i; w; \bar{p}_i(1); \dots; \bar{p}_i(j); \mathbf{0}; 1; t_i],$$

322 where $\bar{p}_i(j)$ is approximation for $p_i(j)$ (Definition 1). In the expressions of h_i and h_i , the dimension 323 of 0 differs. Specifically, the 0 in h_i is larger than in h_i . The dimensional difference between these 0 vectors equals the dimension of $\bar{p}_i(j)$. Suppose function r is L_r -smooth in bounded domain \mathcal{W} , then for any $i \in [n+1], j \in [N], \bar{p}_i(j)$ such that

$$\bar{p}_i(j) = p_i(j) + \epsilon(i,j), \quad \|\epsilon(i,j)\|_2 \le \begin{cases} (\sum_{l=0}^{j-1} K^{l/2} L_r^l B_v^l) \sqrt{K} \epsilon_r, & 1 \le j \le N-1\\ (\sum_{l=0}^{N-1} K^{l/2} L_r^l B_v^l) \sqrt{d} \epsilon_r, & j = N \end{cases}$$
(3.9)

Additionally, for any $j \in [N]$, the norm of parameters B_{θ_j} defined as (C.1) such that $B_{\theta_j} \leq 1 + KC_1$.

Proof Sketch. By Definition 1, we provide term-by-term approximations for $p_i(j)$ as forward propagation. Specifically, we construct Attention layers to implement forward propagation algorithm. Then we establish upper bounds for the errors $\|\bar{p}_i(j) - p_i(j)\|_2$ inductively. Finally, we present the norms (C.1) of the Transformers constructed. Please see Appendix D.2 for a detailed proof.

Notice that the form of error accumulation in Lemma 2 is complicated. For the ease of later presentations, we define the upper bound of coefficient in (3.9) as

$$E_r := \max_{j \in [N]} \frac{\|\epsilon(i,j)\|_2}{\epsilon_r} = \max_{j \in [N]} \{ (\sum_{l=0}^{j-1} K^{l/2} L_r^l B_v^l) \sqrt{K}, (\sum_{l=0}^{N-1} K^{l/2} L_r^l B_v^l) \sqrt{d} \},$$
(3.10)

such that (3.9) becomes

$$\bar{p}_i(j) = p_i(j) + \epsilon(i,j), \quad \|\epsilon(i,j)\|_2 \le E_r \epsilon_r.$$
(3.11)

343 Moreover, we abbreviate $\bar{p}_i := [\bar{p}_i(1); \ldots; \bar{p}_i(N)] \in \mathbb{R}^{(N-1)K+d}$, such that the output of $\operatorname{Attn}_{\theta_1} \circ \cdots \circ \operatorname{Attn}_{\theta_N}$ is

$$h_i = [x_i; y_i; w; \bar{p}_i; \mathbf{0}; 1; t_i].$$
(3.12)

Then, the next lemma approximates $r'_i(j)$ base on $\bar{p}_i(j)$ obtained in Lemma 2.

Lemma 3 (Approximate $r'_i(j)$). Let upper bounds $B_v, B_x > 0$ such that for any $k \in [K], j \in [N]$ and $i \in [n], ||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. For any $j \in [N], i \in [n]$, define

$$B_{r}^{\prime j} := \max_{|t| \le B_{v} B_{r'}^{j-1}} |r'(t)|, \quad B_{r'}^{0} := B_{x}, \quad \text{and} \quad B_{r'} := \max_{j} B_{r'}^{j}$$

Suppose function r'(t) is $(\epsilon_{r'}, R_2, M_2, C_2)$ -approximable for $R_2 = \max\{B_v B_{r'}, 1\}, M_2 \leq \widetilde{\mathcal{O}}(C_2^2 \epsilon_r'^{-2})$, where C_2 depends only on R_2 and the C^2 -smoothness of r'. Then, for any $\epsilon_r > 0$, there exist an attention layer $\operatorname{Attn}_{\theta_{N+1}}$ such that for any input $h_i \in \mathbb{R}^D$ takes from (3.12), it maps

$$h_i = [x_i; y_i; w; \bar{p}_i; \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{Attn}_{\theta_{N+1}}} \widetilde{h_i} = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \mathbf{0}; 1; t_i],$$

where $\bar{r}'_i(j)$ is approximation for $r'_i(j)$ (Definition 2) and $\bar{r}'_i := [\bar{r}'_i(0); \ldots; \bar{r}'_i(N-1)] \in \mathbb{R}^{(N-2)K+d}$. Similar to Lemma 2, in the expressions of h_i and \tilde{h}_i , the dimension of 0 differs. In addition, let E_r be defined in (3.11), for any $i \in [n+1], j \in [N], k \in [K], \bar{r}'_i(j)$ such that

$$\bar{r}'_{i}(j-1)[k] = r'_{i}(j-1)[k] + \epsilon(i,j,k), \quad |\epsilon(i,j,k)| \le \epsilon_{r'} + L_{r'}B_{v}E_{r}\epsilon_{r},$$
(3.13)

where ϵ_r denotes the error generated in approximating r by sum of ReLUs \bar{r} follows (D.5). Additionally, the norm of parameters $B_{\theta_{N+1}}$ defined as (C.1) such that $B_{\theta_{N+1}} \leq 1 + K(N-1)C_2$.

Proof Sketch. By Lemma 2, we obtain $\bar{p}_i(j)$, the approximation for $p_i(j)$ (3.3). Using $\bar{p}_i(j)$, we construct an Attention layer to approximate $r'_i(j)$. We then establish upper bounds for the errors $|\bar{r}'_i(j)[k] - r'_i(j)[k]|$ by applying Cauchy-Schwarz inequality and Lemma 2. Finally we present the norms (C.1) of the Transformers constructed. Please see Appendix D.3 for a detailed proof.

Let $\operatorname{Attn}_{\theta_j}(j \in [N])$ be as defined in Lemma 2, then Lemma 3 implies that for the input takes from Problem 2, the output of $\operatorname{Attn}_{\theta_1} \circ \cdots \circ \operatorname{Attn}_{\theta_{N+1}}$ is

$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \mathbf{0}; 1; t_i].$$
(3.14)

Step 2. Now, we construct an approximation for $u(p_i(N), y_i) = (\frac{\partial \ell(p_i(N), y_i)}{\partial p_i(N)})^\top$.

1375 Lemma 4 (Approximate $u(p_i(N), y_i)$). Let upper bounds $B_v, B_x, > 0$ such that for any $k \in [K], j \in [N]$ and $i \in [n], ||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. For any $k \in [d]$, suppose function **1377** u(t, y)[k] be $(\epsilon_l, R_3, M_3^k, C_3^k)$ -approximable for $R_3 = \max\{B_v B_r, B_y, 1\}, M_3 \leq \widetilde{\mathcal{O}}((C_3^k)^2 \epsilon_l^{-2}),$

where C_3^k depends only on R_3^k and the C^3 -smoothness of u(t, y)[k]. Then, there exists an MLP layer $MLP_{\theta_{N+2}}$ such that for any input sequences $h_i \in \mathbb{R}^D$ takes from (3.14), it maps

$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{MLP}_{\theta_{N+2}}} \tilde{h_i} = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \mathbf{0}; 1; t_i],$$

where $g_i \in \mathbb{R}^d$ is an approximation for $u(p_i(N), y_i)$. For any $k \in [d]$, assume $u(p_i(N), y_i)$ is L_l -Lipschitz continuous. Then the approximation g_i such that,

$$g_i[k] = u(p_i(N), y_i)[k] + \epsilon(i, k), \quad \text{with} \quad |\epsilon(i, k)| \le \epsilon_l + L_l E_r \epsilon_r. \tag{3.15}$$

Additionally, the parameters θ_{N+2} such that $B_{\theta_{N+2}} \leq \max\{R_3 + 1, C_3\}$.

Proof Sketch. By Lemma 2, we obtain $\bar{p}_i(N)$, the approximation for $p_i(N)$ (3.3). Using $\bar{p}_i(N)$, we construct an MLP layer to approximate u. We then establish upper bounds for the errors $|g_i[k] - u[k]|$ and present the norms (C.1) of Transformers constructed. See Appendix D.4 for a detailed proof. \Box

Let $\operatorname{Attn}_{\theta_j}(j \in [N+1])$ be as defined in Lemma 2 and Lemma 3, then for any input sequences $h_i \in \mathbb{R}^D$ takes from (2.1), the output of $\operatorname{Attn}_{\theta_1} \circ \cdots \circ \operatorname{Attn}_{\theta_{N+1}} \circ \operatorname{MLP}_{\theta_{N+2}}$ is

$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \mathbf{0}; 1; t_i].$$
(3.16)

Before introducing our next approximation lemma, we define an element-wise multiplication layer, since attention mechanisms and MLPs are unable to compute self-products (e.g., output xy from input [x; y]). To enable self-multiplication, we introduce a function γ . This function, for any square matrix, preserves the diagonal elements and sets all others to zero.

Definition 5 (Operator Function γ). For any square matrix $A \in \mathbb{R}^{n \times n}$, define $\gamma(A) := \text{diag}(A[1,1], \dots A[n,n]) \in \mathbb{R}^{n \times n}$.

By Definition 5, we introduce the following element-wise multiplication layer, capable of performing self-multiplication operations such as the Hadamard product.

Definition 6 (Element-wise Multiplication Layer). Let γ be defined as Definition 5. An element-wise multiplication layer with m heads is denoted as $Attn_{\theta}(\cdot)$ with parameters $\theta = \{Q_m, K_m, V_m\}_{m \in [M]}$. On any input sequence $H \in \mathbb{R}^{D \times n}$,

$$\text{EWML}_{\theta}(H) = H + \sum_{i=1}^{m} (V_m H) \cdot \gamma((Q_m H)^{\top}(K_m H)).$$
(3.17)

where $Q_m, K_m, V_m \in \mathbb{R}^{D \times D}$ and $\gamma(\cdot)$ is operator function follows Definition 5. In vector form, for for each token $h_i \in \mathbb{R}^D$ in H, it outputs $[\text{EWML}_{\theta}(H)]_i = h_i + \sum_{m=1}^M \gamma(\langle Q_m h_i, K_m h_i \rangle) \cdot V_m h_i$. In addition, we define *L*-layer neural networks $\text{EWML}_{\theta}^L := \text{EWML}_{\theta_1} \circ \cdots \circ \text{EWML}_{\theta_L}$.

Remark 3 (Necessary for Element-Wise Multiplication Layer). As we shall show in subsequent 418 sections, element-wise multiplication layer is capable of implementing multiplication in h_i . Specifi-419 cally, it allows us to multiply some elements in h_i in Lemma 5. By Definition 7, it is impossible for 420 transformer layers to achieve our goal without any other assumptions.

422 Similar to (C.1), we define the norm for L-layer transformer EWML^L_{θ} as:

$$B_{\theta} := \max_{m \in [M], l \in [L]} \{ \|Q_m^l\|_1, \|K_m^l\|_1, \|V_m^l\|_1 \}.$$
(3.18)

Then, given the approximations for $p_i(j)$ and $r'_i(j)$, we use N element-wise multiplication layer (Definition 6) to approximate $s_i(j)$, the chain-rule intermediate terms defined as Definition 3.

428 Lemma 5 (Approximate $\bar{s}_t(j)$). Recall that $s_i(j) = r'_i(j-1) \odot (V_{j+1}^\top \cdot s_i(j+1))$ follows Definition 3. 429 Let the initial input take from (3.16). Then, there exist N element-wise multiplication layers: 430 EWML_{θ_{N+3}},..., EWML_{θ_{2N+2}} such that for input sequences, $j \in [N]$,

431
$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \bar{s}_i(N); \dots; \bar{s}_i(j+1); \mathbf{0}; 1; t_i],$$

they map EWML_{θ_{2N+3-j}} $(h_i) = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \bar{s}_i(N); \dots; \bar{s}_i(j); \mathbf{0}; 1; t_i]$, where the approximation $\bar{s}_i(j)$ is defined as recursive form: for any $i \in [n+1], j \in [N]$,

$$\bar{s}_i(j) := \begin{cases} \bar{r}'_i(j-1) \odot (V_{j+1}^\top \cdot \bar{s}_i(j+1)), & j \in [N-1] \\ \bar{r}'_i(N-1) \odot g_i, & j = N. \end{cases}$$
(3.19)

Additionally, for any $j \in [N]$, $B_{\theta_{N+2+j}}$ defined in (C.1) satisfies $B_{\theta_{N+2+j}} \leq 1$.

Proof Sketch. By Lemma 2 and Lemma 3, we obtain $\bar{p}_i(j)$ and $\bar{r}'_i(j)$, the approximation for $p_i(j)$ (3.3) and $r'_i(j)$ respectively. Using $\bar{p}_i(j)$ and $\bar{r}'_i(j)$, we construct N element-wise multiplication layers to approximate $s_i(j)$. We then present the norms (3.18) of the EWMLs constructed. Please see Appendix D.5 for a detailed proof.

Let $\operatorname{Attn}_{\theta_j}(j \in [N+1])$, $\operatorname{MLP}_{\theta_{N+2}}$ be as defined in Lemma 2, Lemma 3 and Lemma 4 respectively. Define $\bar{s}_i := [\bar{s}_i(N); \ldots; \bar{s}_i(1)] \in \mathbb{R}^{(N-1)K+d}$, then for any input sequences $h_i \in \mathbb{R}^D$ takes from Problem 2, the output of neural network

$$\operatorname{Attn}_{\theta_1} \circ \cdots \circ \operatorname{Attn}_{\theta_{N+1}} \circ \operatorname{MLP}_{\theta_{N+2}} \circ \operatorname{EWML}_{\theta_{N+3}} \circ \cdots \circ \operatorname{EWML}_{\theta_{2N+1}}, \tag{3.20}$$

is

$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \bar{s}_i; \mathbf{0}; 1; t_i].$$
(3.21)

For the sake of simplicity, we consider ReLU Attention layer and MLP layer are both a special kind of transformer. In this way, by Definition 9, (3.20) becomes

$$\Gamma F_{\theta}^{N+2} \circ EWML_{\theta}^{N-1}$$

Next we calculate the error accumulation $|\bar{s}_i(j)[k] - s_i(j)[k]|$ based on Lemma 3 and Lemma 4.

Lemma 6 (Error for $\bar{s}_i(j)$). Suppose the upper bounds $B_v, B_x > 0$ such that for any $k \in [K], j \in [N]$ and $i \in [n]$, $||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. Let $r'_i(j) \in \mathbb{R}^K$ such that $r'_i(j)[k] := r'(v_{j+1_k}^\top p_i(j))$ follows Definition 2. Let $s_i(j) = R_i(j-1)V_{j+1}^\top \dots R_i(N-2)V_N^\top \cdot R_i(N-1)u$ follows Definition 3. Let $\bar{r}'_i(j), g_i, \bar{s}_i(j)$ be the approximations for $r'_i(j), u(p_i(N), y_i), s_i(j)$ follows Lemma 3, Lemma 4 and Lemma 5 respectively. Let $B_{r'}$ be the upper bound of $\bar{r}'_i(j)[k]$ and $r'_i(j)[k]$ as defined in Lemma 3. Let B_l be the upper bound of g_i and $u(p_i(N), y_i)$ as defined in Lemma 4. Then for any $i \in [n+1], j \in [N], k \in [K]$,

$$\bar{s}_i(j)[k] \le B_s$$
, and $|\bar{s}_i(j)[k] - s_i(j)[k]| \le E_s^r \epsilon_r + E_s^{r'} \epsilon_{r'} + E_s^l \epsilon_l$, where,

$$P := \max\{\sqrt{K}, \sqrt{d}\}\tag{3.22}$$

$$B_s := \max_{j \in [N]} \{ (P \cdot B_{r'} B_v)^{N-j} B_{r'} B_l \},$$
(3.23)

$$E_{s}^{r} := \max_{j \in [N]} \{ L_{r'} E_{r} P B_{s} B_{v}^{2} [\sum_{l=0}^{N-j-1} (B_{r'} B_{v} P)^{l}] + (B_{r'} B_{v} P)^{N-j} (B_{l} L_{r'} B_{v} E_{r} + B_{r'} L_{l} E_{r}) \},$$

$$E_{s}^{r'} := \max_{j \in [N]} \{ P B_{s} B_{v} [\sum_{l=0}^{N-j-1} (B_{r'} B_{v} P)^{l}] + (B_{r'} B_{v} P)^{N-j} B_{l} \},$$

$$E_{s}^{l} := \max_{i \in [N]} \{ (B_{r'} B_{v} P)^{N-j} B_{r'} \}.$$
(3.24)

Above, B_s is the upper bound of $\bar{s}_i(j)[k]$ and $E_s^r, E_s^{r'}, E_s^l$ are the coefficients of $\epsilon_r, \epsilon'_r, \epsilon_l$ in the upper bounds of $|\bar{s}_i(j)[k] - s_i(j)[k]|$, respectively.

Proof Sketch. By Lemma 5, we manage to approximate $s_i(j)$ by $\bar{s}_i(j)$. By triangle inequality, 480 we have $|\bar{s}_i(j)[k] - s_i(j)[k]| \le |\bar{r}'_i(n-1)[k] - r'_i(n-1)[k]| \cdot |v_{n+1_k}^\top \bar{s}_i(n+1)| + |r'_i(n-1)[k]| \cdot |v_{n+1_k}^\top \bar{s}_i(n+1)| + |v_i(n-1)[k]| + |$

Lemma 6 offers the explicit form of the error $|\bar{s}_i(j)[k] - s_i(j)[k]|$, which is crucial for calculating the error $\|\nabla_w \bar{\mathcal{L}}_n(w) - \nabla_w \mathcal{L}_n(w)\|_2$ in Theorem 1.

486 **Step 3.** Combining the above, we prove the existence of a neural network, that implements L487 in-context GD steps on our N-layer neural network. And finally we arrive our main result: a neural 488 network \mathcal{T} for Problem 2.

489 **Theorem 1** (In-Context Gradient Descent on N-layer NNs). Fix any $B_v, \eta, \epsilon > 0, L \ge 1$. For 490 any input sequences takes from (2.1), their exist upper bounds B_x, B_y such that for any $i \in [n]$, 491 $\|y_i\|_2 \leq B_y$, $\|x_i\|_2 \leq B_x$. Assume functions r(t), r'(t) and u(t,y)[k] are $L_r, L_{r'}, L_l$ -Lipschitz continuous. Suppose \mathcal{W} is a closed domain such that for any $j \in [N-1]$ and $k \in [K]$, 492

$$\mathcal{W} \subset \left\{ w = [v_{j_k}] \in \mathbb{R}^{D_N} : \|v_{j_k}\|_2 \le B_v \right\}$$

and $\operatorname{Proj}_{\mathcal{W}}$ project w into bounded domain \mathcal{W} . Assume $\operatorname{Proj}_{\mathcal{W}} = \operatorname{MLP}_{\theta}$ for some MLP layer with hidden dimension D_w parameters $\|\theta\| \leq C_w$. If functions r(t), r'(t) and u(t,y)[k] are C^4 smoothness, then for any $\epsilon > 0$, there exists a transformer model NN_{θ} with (2N + 4)L hidden layers consists of L neural network blocks $\mathrm{TF}_{\theta}^{N+2} \circ \mathrm{EWML}_{\theta}^{N} \circ \mathrm{TF}_{\theta}^{2}$, 498

$$\mathrm{NN}_{\theta} := \mathrm{TF}_{\theta}^{N+2} \circ \mathrm{EWML}_{\theta}^{N} \circ \mathrm{TF}_{\theta}^{2} \circ \ldots \circ \mathrm{TF}_{\theta}^{N+2} \circ \mathrm{EWML}_{\theta}^{N} \circ \mathrm{TF}_{\theta}^{2}$$

such that the heads number M^l , parameter dimensions D^l , and the parameter norms B_{θ^l} suffice

$$\max_{e \in [(2N+4)L]} M^{l} \le \widetilde{O}(e^{-2}), \quad \max_{l \in [(2N+4)L]} D^{l} \le O(NK^{2}) + D_{w}, \quad \max_{l \in [(2N+4)L]} B_{\theta^{l}} \le O(\eta) + C_{w} + 1,$$

where $O(\cdot)$ hides the constants that depend on d, K, N, the radius parameters B_x, B_y, B_v and the smoothness of r and ℓ . And this neural network such that for any input sequences $H^{(0)}$, take from (2.1), $NN_{\theta}(H^{(0)})$ implements L steps in-context gradient descent on risk Eqn (2.2): For every $l \in [L]$, the (2N+4)l-th layer outputs $h_i^{((2N+4)l)} = [x_i; y_i; \overline{w}^{(l)}; \mathbf{0}; 1; t_i]$ for every $i \in [n+1]$, and approximation gradients $\overline{w}^{(l)}$ such that

$$\overline{w}^{(l)} = \operatorname{Proj}_{\mathcal{W}}(\overline{w}^{(l-1)} - \eta \nabla \mathcal{L}_n(\overline{w}^{(l-1)}) + \epsilon^{(l-1)}), \quad \overline{w}^{(0)} = \mathbf{0},$$

where $\|\epsilon^{(l-1)}\|_2 < \eta\epsilon$ is an error term. 511

493 494

495

496

497

499 500

501 502

504

505

506

507

508 509 510

522

512 *Proof Sketch.* Let the first 2N + 2 layers of NN_{θ} are Transformers and EWMLs constructed in 513 Lemma 2, Lemma 3, Lemma 4, and Lemma 5. Explicitly, we design the last two layers to implement 514 the gradient descent step (Lemma 1). We then establish the upper bounds for error $\|\nabla_w \mathcal{L}_n(w) - \nabla_w \mathcal{L}_n(w)\|$ $\nabla_w \mathcal{L}_n(w) \|_2$, where $\nabla_w \overline{\mathcal{L}}_n(w)$, derived from the outputs of NN_{θ}, approximates $\nabla_w \mathcal{L}_n(w)$. Next, 515 516 for any $\epsilon > 0$, we select appropriate parameters ϵ_l , ϵ_r and $\epsilon_{r'}$ to ensure that $\|\nabla_w \overline{\mathcal{L}}_n(\overline{w}^{(l-1)}) - \overline{\mathcal{L}}_n(\overline{w}^{(l-1)})\|$ 517 $\nabla_w \mathcal{L}_n(\overline{w}^{(l-1)}) \|_2 \leq \epsilon$ holds for any $l \in [L]$. Please see Appendix D.7 for a detailed proof.

518 As a direct result, the neural networks NN_{θ} constructed earlier is able to approximate the true gradient 519 descent trajectory $\{w_{\text{GD}}^l\}_{l\geq 0}$, defined by $w_{\text{GD}}^0 = \mathbf{0}$ and $w_{\text{GD}}^{l+1} = w_{\text{GD}}^l - \eta \nabla_w \mathcal{L}_n(w_{\text{GD}}^l)$ for any $l\geq 0$. Consequently, Theorem 1 motivates us to investigate the error accumulation under setting 520 521

$$\overline{w}^{(l)} = \operatorname{Proj}_{\mathcal{W}}(\overline{w}^{(l-1)} - \eta \nabla \mathcal{L}_n(\overline{w}^{(l-1)}) + \epsilon^{(l-1)}), \quad \overline{w}^{(0)} = \mathbf{0},$$

523 where $\|\epsilon^{(l-1)}\|_2 \leq \eta \epsilon$ represents error terms. Moreover, Corollary 1.1 shows NN_{θ} constructed in Theorem 1 implements \hat{L} steps ICGD with exponential error accumulation to the true GD paths. 524

525 **Corollary 1.1** (Error for implementing ICGD on N-layer neural network). Fix $L \ge 1$, under the 526 same setting as Theorem 1, (2N + 4)L-layer neural networks NN_{θ} approximates the true gradient 527 descent trajectory $\{w_{\text{GD}}^l\}_{l\geq 0} \in \mathbb{R}^{D_N}$ with the error accumulation $\|\overline{w}^l - w_{\text{GD}}^l\|_2 \leq L_f^{-1}(1 + nL_f)^l \epsilon$, 528 where L_f denotes the Lipschitz constant of $\mathcal{L}_n(w)$ within \mathcal{W} . 529

Proof. Please see Appendix D.8 for a detailed proof. 530

531 4 DISCUSSION AND CONCLUSION

532 We provide an explicit characterization of the ICL capabilities of a transformer model in approximating the gradient descent training process of a N-layer feed-forward neural network. Our results 534 include approximation (Theorem 1) and convergence (Corollary 1.1) guarantees. We further extend 535 our analysis in two ways: (i) from N-layer networks with the same input and output dimensions 536 to scenarios with arbitrary dimensions (Appendix E); (ii) from ReLU-transformers (aligned with (Bai et al., 2023)) to more practical Softmax-transformers for ICGD of N-layer neural network (Appendix F). We support our theory with numerical validations in Appendix G, and apply our results 538 to learn the score function of the diffusion model through ICL in Appendix H. Please see the related 539 works, a detailed comparison with (Wang et al.), broader impact, and limitations in Appendix B.

540 REFERENCES

549

550

551

552

556

578

579

580

581

- 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.
- Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556, 2022.
 - Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*, 2023.
- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding,
 Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek Ilm: Scaling open-source language models with
 longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 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.
- Stanley H Chan. Tutorial on diffusion models for imaging and vision. arXiv preprint arXiv:2403.18103, 2024.
- Mingda Chen, Jingfei Du, Ramakanth Pasunuru, Todor Mihaylov, Srini Iyer, Veselin Stoyanov, and
 Zornitsa Kozareva. Improving in-context few-shot learning via self-supervised training. *arXiv preprint arXiv:2205.01703*, 2022.
- 571
 572
 573
 574
 Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In *International Conference on Machine Learning*, pages 4672–4712. PMLR, 2023.
- 575 Minshuo Chen, Song Mei, Jianqing Fan, and Mengdi Wang. An overview of diffusion models: Applications, guided generation, statistical rates and optimization. *arXiv preprint arXiv:2404.07771*, 2024.
 - Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers. *arXiv* preprint arXiv:2212.10559, 2022.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdi nov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Pre-training to learn in context. *arXiv preprint arXiv:2305.09137*, 2023.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.

- Jerry Yao-Chieh Hu, Wei-Po Wang, Ammar Gilani, Chenyang Li, Zhao Song, and Han Liu. Funda mental limits of prompt tuning transformers: Universality, capacity and efficiency. *arXiv preprint*, 2024a. To appear.
- Jerry Yao-Chieh Hu, Weimin Wu, Yi-Chen Lee, Yu-Chao Huang, Minshuo Chen, and Han Liu. On statistical rates of conditional diffusion transformer: Approximation and estimation. *arXiv preprint*, 2024b. To appear.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.
 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Tokio Kajitsuka and Issei Sato. Are transformers with one layer self-attention using low-rank
 weight matrices universal approximators? In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.
- Shuai Li, Zhao Song, Yu Xia, Tong Yu, and Tianyi Zhou. The closeness of in-context learning and
 weight shifting for softmax regression. *arXiv preprint arXiv:2304.13276*, 2023.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? In
 Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 11048–11064, 2022.
- Abhishek Panigrahi, Sadhika Malladi, Mengzhou Xia, and Sanjeev Arora. Trainable transformer in
 transformer. *arXiv preprint arXiv:2307.01189*, 2023.
- Madhur Panwar, Kabir Ahuja, and Navin Goyal. In-context learning through the bayesian prism.
 arXiv preprint arXiv:2306.04891, 2023.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Itay Safran and Ohad Shamir. Depth-width tradeoffs in approximating natural functions with neural networks. In *International conference on machine learning*, pages 2979–2987. PMLR, 2017.
- Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Victoria Lin, Noah A Smith,
 Luke Zettlemoyer, Scott Yih, and Mike Lewis. In-context pretraining: Language modeling beyond
 document boundaries. *arXiv preprint arXiv:2310.10638*, 2023.
- Seongjin Shin, Sang-Woo Lee, Hwijeen Ahn, Sungdong Kim, HyoungSeok Kim, Boseop Kim, Kyunghyun Cho, Gichang Lee, Woomyoung Park, Jung-Woo Ha, et al. On the effect of pretraining corpora on in-context learning by a large-scale language model. *arXiv preprint arXiv:2204.13509*, 2022.
- Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach
 to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584. PMLR,
 2020a.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
 Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020b.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 647 Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.

648 649 650	Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In <i>International Conference on Machine Learning</i> , pages 35151–35174. PMLR, 2023.
652 653	Zhijie Wang, Bo Jiang, and Shuai Li. In-context learning on function classes unveiled for transformers. In <i>Forty-first International Conference on Machine Learning</i> .
654 655 656	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, et al. Larger language models do in-context learning differently. <i>arXiv</i> preprint arXiv:2303.03846, 2023.
657 658 659	Noam Wies, Yoav Levine, and Amnon Shashua. The learnability of in-context learning. Advances in Neural Information Processing Systems, 36, 2024.
660 661	Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. <i>arXiv preprint arXiv:2111.02080</i> , 2021.
663 664 665	Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. <i>ACM Computing Surveys</i> , 56(4):1–39, 2023.
666 667 668	Kang Min Yoo, Junyeob Kim, Hyuhng Joon Kim, Hyunsoo Cho, Hwiyeol Jo, Sang-Woo Lee, Sang-goo Lee, and Taeuk Kim. Ground-truth labels matter: A deeper look into input-label demonstrations. <i>arXiv preprint arXiv:2205.12685</i> , 2022.
669 670 671	Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. <i>arXiv preprint arXiv:2306.09927</i> , 2023.
672 673 674 675	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022.
676	
677	
678	
679	
680	
681	
682	
003 604	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

Appendix

705			
706	Α	Informal Version of Results	15
707	P	Palated Works Broader Impact and Limitations	15
708	D	B 1 Related Works	15
709		B.2 Broader Impact	16
710		B.3 Limitations	16
711			
712	С	Supplementary Theoretical Backgrounds	18
713		C.1 Transformers	18
714		C.2 ReLU Provably Approximates Smooth <i>k</i> -Variable Functions	18
715	п	Proofe of Main Tayt	20
716	υ	D 1 Proof of Lemma 1	20
717		D 2 Proof of Lemma 2	21
718		D.3 Proof of Lemma 3	24
719		D.4 Proof of Lemma 4	26
720		D.5 Proof of Lemma 5	27
721		D.6 Proof of Lemma 6	28
722		D.7 Proof of Theorem 1	30
723		D.8 Proof of Corollary 1.1	33
724	Б	Extension, Different Input and Output Dimensions	25
725	Ľ	Extension: Different input and Output Dimensions	33
726	F	Extension: Softmax Transformer	37
727	-	F.1 Universal Approximation of Softmax Transformer	37
728		F.2 In-Context Gradient Descent with Softmax Transformer	37
729		F.3 Proof of Theorem 6	38
730		F.4 Proof of Lemma 16	38
731	C		477
732	G	Experimental Details	47
733		G. 1.1. Derformance of Pel II Transformer	47
734		G.1.2 Performance of Softmax Transformer	40
735		G.2 Experiments for Objective 3	48
736		G.3 Experiments for Objective 4	49
737		r	
738	Η	Application: ICL for Diffusion Score Approximation	51
739		H.1 Score Matching Generative Diffusion Models	51
740		H.2 ICL for Score Approximation	51
741			
742			
743			
744			
745			
746			
747			
748			

A INFORMAL VERSION OF RESULTS

In this section, we show the informal version of our main results.

Theorem 2 (ICGD on *N*-layer NNs, informal version of Theorem 1). Assume functions r(t), r'(t) and u(t, y)[k] (Definition 1) are C^4 -smoothness. Let all parameter and data are bounded, then their exist a explicit constructed transformer capable of solving Problem 1 on *N*-layer NNs.

Theorem 3 (ICGD on General Risk Function, informal version of Theorem 6). Assume all parameters and data are bounded. Let $l(w, x_i, y_i)$ be a loss function with *L*-Lipschitz gradient. Let $\mathcal{L}_n(w) = \frac{1}{n} \sum_{i=1}^n \ell(w, x_i, y_i)$ denote the empirical loss function, then there exists a transformer NN_{θ}, such that for any input sequences $H^{(0)}$, take from (2.1), NN_{θ}($H^{(0)}$) solves Problem 1 on $\mathcal{L}_n(w)$.

B RELATED WORKS, BROADER IMPACT AND LIMITATIONS

In this section, we show the related works, broader impact and limitations.

772 B.1 RELATED WORKS

760

761

762 763

764

765 766

767 768

769 770

771

773

774 **In-Context Learning.** Large language models (LLMs) demonstrate the in-context learning (ICL) ability (Brown, 2020), an ability to flexibly adjust their prediction based on additional data given in 775 context. In recent years, a number of studies investigate enhancing ICL capabilities (Chen et al., 2022; 776 Gu et al., 2023; Shi et al., 2023), exploring influencing factors (Shin et al., 2022; Yoo et al., 2022), and 777 interpreting ICL theoretically (Xie et al., 2021; Wies et al., 2024; Panwar et al., 2023; Li et al., 2023; 778 Bai et al., 2023; Dai et al., 2022). The works most relevant to ours are as follows. (Von Oswald et al., 779 2023) showed that linear attention-only Transformers with manually set parameters closely resemble models trained via gradient descent. (Bai et al., 2023) providing a more efficient construction for in-781 context gradient descent and established quantitative error bounds for simulating multi-step gradient 782 descent. However, these results focused on simple ICL algorithms or specific tasks like least squares, 783 ridge regression, and gradient descent on two-layer neural networks. These algorithms are inadequate 784 for practical applications. For example: (i) Approximating the diffusion score function requires 785 neural networks with multiple layers (Chen et al., 2023). (ii) Approximating the indicator function requires at least 3-layer networks (Safran and Shamir, 2017). Therefore, the explicit construction 786 of transformers to implement in-context gradient descent (ICGD) on deep models is necessary to 787 better align with real-world in-context settings. Our work achieves this by analyzing the gradient 788 descent on N-layer neural networks through the use of ICL. We provide a more efficient construction 789 for in-context gradient descent. Furthermore, we extend our analysis to Softmax-transformer in 790 Appendix F to better align with real-world uses. 791

In-Context Gradient Descent on Deep Models (Wang et al.; Panigrahi et al., 2023). A work
 similar to ours is (Wang et al.). It constructs a family of transformers with flexible activation functions
 to implement multiple steps of ICGD on deep neural networks. This work emphasizes the generality
 of activation functions and demonstrates the theoretical feasibility of such constructions. Our work
 adopts a different approach by enhancing the efficiency of transformers and better aligning with
 practical applications. We introduce the following novelties:

- 798• More structured and efficient transformer architecture. While the work (Wang et al.) uses a
 $O(N^2L)$ -layer transformer to approximate L gradient descent steps on N-layer neural networks,
our approach achieves more efficient simulation for ICGD. We approximate specific terms in the
gradient expression to reduce computational costs, requiring only a (2N+4)L-layer transformer for
L gradient descent steps. Our method focuses on selecting and approximating the most impactful
intermediate terms in the explicit gradient descent expression (Lemmas 3 to 5), optimizing layer
complexity to O(NL).
- Less restrictive input and output dimensions for *N*-layer neural networks. The work (Wang et al.) simplifies the output of *N*-layer networks to a scalar. Our work expands this by considering cases where output dimensions exceed one, as detailed in Appendix E. This includes scenarios where input and output dimensions differ.
- More practical transformer model. The work (Wang et al.) discusses activation functions in the attention layer that meet a general decay condition ((Wang et al., Definition 2.3)) without consid-

ering the Softmax activation function. We extend our analysis to include Softmax-transformers.
 Our analysis reflects more realistic applications, as detailed in Appendix F.

More advanced and complicated applications. The work (Wang et al.) discusses the applications to functions, including indicators, linear, and smooth functions. We explore more advanced and complicated scenarios, i.e., the score function in diffusion models discussed in Appendix H. The score function (Chen et al., 2023) falls outside the smooth function class. This enhancement broadens the applicability of our results.

817 Another work similar to ours is (Panigrahi et al., 2023). It proposes a new efficient construction, 818 Transformer in Transformer (TINT), to allow a transformer to simulate and finetune more complex 819 models (e.g., one transformer). The main distinction between ours and (Panigrahi et al., 2023) 820 lies in the different aims: Our approach focuses on using a standard transformer for the simulator 821 (with a minor modification: the "element-wise multiplication layer"), and we provide a theoretical 822 understanding of how a standard transformer can learn the ICGD of an N-layer network using ICL. 823 In contrast, the work (Panigrahi et al., 2023) aims to build even stronger transformers by introducing several structural modifications that enable running gradient descent on auxiliary transformers. While 824 it demonstrates in-context gradient descent for a more advanced model, i.e., one transformer, our 825 work offers the following potential advantages: 826

- Explicit transformer construction. We provide an explicit construction of the transformer, whereas the work (Panigrahi et al., 2023) does not detail the explicit construction of model parameters within their transformer.
- Exact gradient descent. We compute the exact and explicit gradient descent for an *N*-layer network (Lemma 1). Building on this, we employ the transformer's ICL to perform gradient descent on all parameters. However, the work (Panigrahi et al., 2023) stops the gradient computation through attention scores in the self-attention layer and only updates the value parameter in the self-attention module. Additionally, it uses Taylor expansion to approximate the gradient.
- Rigorous error and convergence guarantees. We provide rigorous gradient descent approximation errors (for multiple steps) and convergence guarantees for the ICGD on an *N*-layer network (Corollary 1.1 and Lemma 14). However, the work (Panigrahi et al., 2023) only presents the gradient approximation error for each specific part of the parameters in a single step.
- Attention layer better aligned with practice. Our analysis is based on ReLU-attention (Theorem 1) or Softmax-attention (Theorem 6), whereas the work (Panigrahi et al., 2023) utilizes linear attention. Our choice of attention layer better aligns with practical applications.
 - B.2 BROADER IMPACT

This theoretical work aims to shed light on the foundations of large transformer-based models and is not expected to have negative social impacts.

B.3 LIMITATIONS

842

843 844

845

846 847

848

849

Our work has the following four limitations:

Although we provide a theoretical guarantee for the ICL of the Softmax-Transformer to approximate gradient descent in *N*-layer NN, characterizing the weight matrices construction in Softmax-Transformer remains challenging. This motivates us to rethink transformer universality and explore more accurate proof techniques for ICL in Softmax-Transformer, which we leave for future work.

- The hidden dimension and MLP dimension of the transformer in Theorem 1 are both $\tilde{O}(NK^2) + D_w$, which is very large. The reason for the large dimensions is that if we use ICL to perform ICGD on the *N*-layer network, we need to allow the transformer to realize the *N*-layer network parameters. This means that it is reasonable for the input dimension to be so large. However, it is possible to reduce the hidden dimension and MLP dimension of the transformer through smarter construction. We leave this for future work.
- The generalization capabilities are limited compared with traditional transformers. In our setting, the pretraining task refers to using in-context examples generated by an N-layer network for a given N. Specifically, during pretraining, the distribution of the N-layer network parameters is predetermined (e.g., N(0, I)). The input data distribution of N-layer network for generating the

864	in contact examples is also predatornized (e.g. $N(-2, I)$). The consultation conshibities include
865	the following two aspects: (i) Varying the input data distribution for the N layer network to generate
866	the in-context examples. For example, we change the input data distribution for $N(-2, I)$ to
867	0.9N(-2, I) + 0.1N(2, I) during the testing in Appendix G 1 (ii) Varying the distribution of the
868	N-layer network parameters. For example, we change the distribution from $N(0, I)$ to $N(0, 5, I)$
869	in Appendix G.2. The above points lead to differences between the distributions of in-context
870	examples during pretraining and testing. However, we must generate the in-context examples by
871	the N -layer network with the same hyperparameters, including the network width and depth. We
872	leave the theoretical analysis of broader generalization capabilities for future work.
873 •	In theory, the FLOPs (Hoffmann et al., 2022) required to perform one forward pass of the trans-
874	former are greater than those required for the direct training of an N-layer network. (i) For the
875	forward pass of the transformer, the FLOPs for in-context learning (ICL) are $O(nLN^3K^5/\epsilon^2)$,
876	where ϵ is the approximation error in the sum of ReLU. (ii) For direct training of the N-layer
877	network, the FLOPs without ICL are $O(nLNK^2)$. Therefore, the FLOPs required for ICL ex-
878	ceed those needed for direct training of the N-layer network. However, experimental results in
879	Appendix G demonstrate that the transformer with ICL can achieve the performance of a trained
880	o-layer network using rewer FLOPs in practice (3.3 billion vs. 7.6 billion FLOPs). This finding
881	research
882	research.
883	
884	
885	
886	
887	
888	
880	
800	
801	
802	
803	
894	
805	
806	
807	
808	
800	
000	
900	
002	
902	
903	
005	
905	
900	
908	
000	
910	
911	
912	
013	
914	
015	
916	
917	
411	

918 C SUPPLEMENTARY THEORETICAL BACKGROUNDS

Here we present some ideas we built on.

C.1 TRANSFORMERS

Lastly, we introduce key components for constructing a transformer for ICGD: ReLU-Attention, MLP, and element-wise multiplication layers. We begin with the ReLU-Attention layer.

Definition 7 (ReLU-Attention Layer). For any input sequence $H \in \mathbb{R}^{D \times n}$, an *M*-head ReLUattention layer with parameters $\theta = \{Q_m, K_m, V_m\}_{m \in [M]}$ outputs

$$\operatorname{Attn}_{\theta}(H) \coloneqq H + \frac{1}{n} \sum_{m=1}^{M} (V_m H) \cdot \sigma((Q_m H)^{\top} (K_m H))$$

where $Q_m, K_m, V_m \in \mathbb{R}^{D \times D}$ and $\sigma(\cdot)$ is element-wise ReLU activation function. In vector form, for each token $h_i \in \mathbb{R}^D$ in H, it outputs $[Attn_{\theta}(H)]_i = h_i + \frac{1}{n} \sum_{m=1}^{M} \sum_{s=1}^{n} \sigma(\langle Q_m h_i, K_m h_s \rangle) \cdot V_m h_s.$

Notably, Definition 7 uses normalized ReLU activation σ/n , instead of the standard Softmax. We adopt this for technical convenience following (Bai et al., 2023). Next we define the MLP layer.

Definition 8 (MLP Layer). For any input sequence $H \in \mathbb{R}^{D \times n}$, an d'-hidden dimensions MLP layer with parameters $\theta = (W_1, W_2)$ outputs $\mathrm{MLP}_{\theta}(H) := H + W_2 \sigma(W_1 H)$, where $W_1 \in \mathbb{R}^{d' \times D}$, $W_2 \in \mathbb{R}^{D \times d'}$ and $\sigma(\cdot) : \mathbb{R} \to \mathbb{R}$ is element-wise ReLU activation function. In vector form, for each token $h_i \in \mathbb{R}^D$ in H, it outputs $\mathrm{MLP}_{\theta}(H)_i := h_i + W_2 \sigma(W_1 h_i)$.

Then, we consider a transformer architecture with $L \ge 1$ transformer layers, each consisting of a self-attention layer followed by an MLP layer.

Definition 9 (Transformer). For any input sequence $H \in \mathbb{R}^{D \times n}$, an *L*-layer transformer with parameters $\theta = \{\theta_{Attn}, \theta_{MLP}\}$ outputs

$$\mathrm{TF}^{\mathrm{L}}_{\theta}(H) := \mathrm{MLP}_{\theta_{\mathrm{mlp}}^{(L)}} \circ \mathrm{Attn}_{\theta_{\mathrm{attn}}^{(L)}} \dots \mathrm{MLP}_{\theta_{\mathrm{mlp}}^{(1)}} \circ \mathrm{Attn}_{\theta_{\mathrm{attn}}^{(1)}}(H),$$

950 where $\theta = \{\theta_{Attn}, \theta_{MLP}\}$ consists of Attention layers $\theta_{Attn} = \{(Q_m^l, K_m^l, V_m^l)\}_{l \in [L], m \in [M^l]}$ and 951 MLP layers $\theta_{MLP} = \{(W_1^l, W_2^l)\}_{l \in [L]}$. Above, for any $l \in [L], m \in [M^l], Q_m^l, K_m^l, V_m^l \in \mathbb{R}^{D \times D}$ 952 and $(W_1^l, W_2^l) \in \mathbb{R}^{d' \times D} \times \mathbb{R}^{D \times d'}$ In this section, we consider ReLU Attention layer and MLP layer 953 are both a special kind of 1-layer transformer, which is for technical convenience.

For later proof use, we define the norm for L-layer transformer TF_{θ} as:

$$B_{\theta} := \max_{l \in [L]} \left\{ \max_{m \in [M]} \left\{ \|Q_m^l\|_1, \|K_m^l\|_1 \right\} + \sum_{i=1}^m \|V_m^l\|_1 + \|W_1\|_1 + \|W_2\|_1 \right\}.$$
(C.1)

The choice of operation norm and max/sum operation is for convenience in later proof only, as our result depends only on B_{θ} .

C.2 RELU PROVABLY APPROXIMATES SMOOTH k-VARIABLE FUNCTIONS

Following lemma expresses that the smoothness enables the approximability of sum of ReLU.

Lemma 7 (Approximating Smooth k-Variable Functions, modified from Proposition A.1 of (Bai et al., 2023)). For any ϵ , $C_l > 0$, $R \ge 1$. If function $g : \mathbb{R}^k \to \mathbb{R}$ such that for $s := \lceil (k-1)/2 \rceil + 1$, g is a C^s function on $B^k_{\infty}(R)$, and for all $i \in \{0, 1, \ldots, s\}$,

970
971
$$\sup_{z \in B^k_{\infty}(R)} \|\nabla^i g(z)\|_{\infty} \le L_i, \quad \max_{0 \le i \le s} L_i R^i \le C_l,$$

972 973	then function $C(k)C_l^2\log(1 - C_l)$	g is (ϵ, R, H, C) -a + $C_1/\epsilon)/\epsilon^2$ and $C \leq$	pproximable by $C(k)C_l$ where C	sum of ReLUs (k) is a constant	(Definition 4) with that depends only on k	$H \leq $
974		, ,,		< /	1 2	
975						
976						
977						
978						
979						
980						
981						
982						
983						
984						
985						
900						
907						
900						
905						
991						
992						
993						
994						
995						
996						
997						
998						
999						
1000						
1001						
1002						
1003						
1004						
1005						
1006						
1007						
1008						
1009						
1010						
1011						
1012						
1014						
1015						
1016						
1017						
1018						
1019						
1020						
1021						
1022						
1023						
1024						
1025						

1026 D PROOFS OF MAIN TEXT

1029 D.1 PROOF OF LEMMA 1

Lemma 8 (Lemma 1 Restated: Decomposition of One Gradient Descent Step). Fix any $B_v, \eta > 0$. Suppose loss function $\mathcal{L}_n(w)$ on n data points $\{(x_i, y_i)\}_{i \in [n]}$ follows (2.2). Suppose closed domain \mathcal{W} and projection function $\operatorname{Proj}_{\mathcal{W}}(w)$ follows (3.4). Let $A_i(j), r'_i(j), R_i(j), V_j$ be as defined in Definition 2. Then the explicit form of gradient $\nabla \mathcal{L}_n(w)$ becomes

$$\nabla \mathcal{L}_n(w) = \frac{1}{2n} \sum_{i=1}^n \begin{bmatrix} A_i(1) \\ \vdots \\ A_i(N) \end{bmatrix}$$

where $A_i(j)$ denote the derivative of $\ell(p_i(N), y_i)$ with respect to the parameters in the *j*-th layer,

$$A_{i}(j) = \begin{cases} (R_{i}(N-1) \cdot V_{N} \cdot \ldots \cdot R_{i}(j-1) \cdot \left[\mathbf{I}_{K \times K} \otimes p_{i}(j-1)^{\top}\right])^{\top} \cdot \left(\frac{\partial \ell(p_{i}(N), y_{i})}{\partial p_{i}(N)}\right)^{\top}, & j \neq N \\ (R_{i}(N-1) \cdot \left[\mathbf{I}_{d \times d} \otimes p_{i}(N-1)^{\top}\right])^{\top} \cdot \left(\frac{\partial \ell(p_{i}(N), y_{i})}{\partial p_{i}(N)}\right)^{\top}, & j = N. \end{cases}$$

Proof of Lemma 1. We start with calculating $\nabla_w \mathcal{L}_n(w)$. By chain rule and (2.2),

$$\underbrace{\nabla_{w}\mathcal{L}_{n}(w)}_{\mathbb{R}^{D_{N}\times 1}} = \frac{1}{2n} \sum_{i=1}^{n} \underbrace{\left[\frac{\partial}{\partial w} p_{i}(N)\right]^{\top}}_{\mathbb{R}^{D_{N}\times d}} \underbrace{\left[\frac{\partial}{\partial p_{i}(N)} \ell(p_{i}(N), y_{i})\right]^{\top}}_{\mathbb{R}^{d\times 1}} \qquad (By (2.2) \text{ and chain rule})$$

Thus we only need to calculate $\frac{\partial}{\partial w}p_i(N)$. For a vector x and a function $r : \mathbb{R} \to \mathbb{R}$, we use $\mathbf{r}(x)$ to denote the vector that *i*-th coordinate is $r(x_i)$. Let $R_i(j), V_j$ follows Definition 2, then it holds

$$\underbrace{\frac{\partial p_i(N)}{\partial w}}_{\mathbb{R}^{d \times D_N}} = \underbrace{\frac{\partial \mathbf{r}(\overbrace{V_N}^{\mathbb{R}^{d \times K}}, \overbrace{p_i(N-2)}^{\mathbb{R}^{K}})}{\partial w}}_{\mathbb{R}^{d \times D_N}}$$
(By Definition 1)

$$= \underbrace{\frac{\partial \mathbf{r}(V_N \cdot p_i(N-1))}{\partial V_N \cdot p_i(N-1)}}_{\mathbb{R}^{d \times d}} \cdot \underbrace{\frac{\partial V_N \cdot p_i(N-1)}{\partial w}}_{\mathbb{R}^{d \times D_N}}$$
(By chain rule)

)

$$= \operatorname{diag}\{r'(v_{N_1}^{\top}p_i(N-1)), \dots, r'(v_{N_K}^{\top}p_i(N-1))\} \cdot \frac{\partial V_N \cdot p_i(N-1)}{\partial w}$$
(By Definition 2)

$$= R_i(N-1) \cdot \frac{\partial V_N \cdot p_i(N-1)}{\partial w}.$$
 (D.1)

1072 Notice that for any $k \in [d]$, v_{N_k} is a part of w, thus

1074
1075
1076
$$\frac{\partial v_{N_k}}{\partial w} = \begin{bmatrix} \begin{array}{c} D_{N-1} + (k-1)K & K & D_N - D_{N-1} - kK \\ \hline \mathbf{0} & \mathbf{I} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{d \times D_N}. \quad (D.2)$$

1077 Therefore, letting \otimes denotes Kronecker product, it holds

1079
$$\frac{\partial V_N \cdot p_i(N-1)}{\partial w}$$

$$\begin{bmatrix} v_{N_1}^\top \cdot \frac{\partial p_i(N-1)}{\partial w} + p_i(N-1)^\top \cdot \frac{\partial v_{N_1}}{\partial w} \end{bmatrix}$$

(By chain rule and product rule)

$$\begin{bmatrix} v_{N_d}^\top \cdot \frac{\partial p_i(N-1)}{\partial w} + p_i(N-1)^\top \cdot \frac{\partial v_{N_d}}{\partial w} \end{bmatrix}$$

= $V_N \cdot \frac{\partial p_i(N-1)}{\partial w} + \begin{bmatrix} \mathbf{0}_{D_{N-1}}; \mathbf{I}_{K \times K} \otimes p_i(N-1)^\top \end{bmatrix},$ (D.3)

where the last step follows from the definition of V_N (i.e., Definition 2) and (D.2).

Substituting (D.3) into (D.1), we obtain

$$\frac{\partial p_i(N)}{\partial w} = R_i(N-1) \cdot (V_N \cdot \frac{\partial p_i(N-1)}{\partial w} + \left[\mathbf{0}_{D_{N-1}}; \mathbf{I}_{d \times d} \otimes p_i(N-1)^\top\right]).$$

Similarly, for any $j \in [N]$, we can proof

$$\frac{\partial p_i(j)}{\partial w} = R_i(j-1) \cdot (V_j \cdot \frac{\partial p_i(j-1)}{\partial w} + \left[\mathbf{0}_{D_{j-1}}; \mathbf{I}_{K \times K} \otimes p_i(j-1)^\top; \mathbf{0}_{D_N - D_j} \right]).$$
(D.4)

By the recursion formula (D.4), for any $j \in [N-1]$, we calculate $A_i(j)$ as follows,

$$A_{i}(j) = \left(\left(\frac{\partial \ell(p_{i}(N), y_{i})}{\partial p_{i}(N)} \cdot \frac{\partial p_{i}(N)}{\partial w} \right)^{\mathsf{T}} \right) [D_{j-1} : D_{j}]$$
(By Definition 2)
$$\frac{\partial p_{i}(N)}{\partial p_{i}(N)} = \frac{\partial \ell(p_{i}(N), y_{i})}{\partial w} = 0$$

$$= \left(\frac{\partial p_i(N)}{\partial w}\right)^\top \cdot \left(\frac{\partial \ell(p_i(N), y_i)}{\partial p_i(N)}\right)^\top [D_{j-1} : D_j]$$
(By transpose property)
$$= \left(\frac{\partial p_i(N)}{\partial y_i(N)}\right)^\top [*, D_{j-1} : D_j] \cdot \left(\frac{\partial \ell(p_i(N), y_j)}{\partial y_i(N)}\right)^\top$$

1105
1106
1107
1108

$$= \left(\frac{\partial p_i(N)}{\partial w}\right)^\top [*, D_{j-1} : D_j] \cdot \left(\frac{\partial \ell(p_i(N), N)}{\partial p_i(N)}\right)$$

$$= \left(\frac{\partial p_i(N)}{\partial w}\right)^\top [*, D_{j-1} : D_j] \cdot \left(\frac{\partial \ell(p_i(N), N)}{\partial p_i(N)}\right)$$

$$= (R_i(N-1) \cdot V_N \cdot \ldots \cdot R_i(j-1) \cdot \left[\mathbf{I}_{K \times K} \otimes p_i(j-1)^{\top}\right])^{\top} \cdot \left(\frac{\partial \ell(p_i(N), y_i)}{\partial p_i(N)}\right)^{\top},$$
(By (D.4))

where M[*, a : b] denotes a sub-matrix of M, which includes all the columns but only the rows from the *a*-th row to the *b*-th row of A. Similarly, for j = N, it holds

$$A_i(N) = (R_i(N-1) \cdot \left[\mathbf{I}_{d \times d} \otimes p_i(N-1)^\top \right])^\top \cdot \left(\frac{\partial \ell(p_i(N), y_i)}{\partial p_i(N)} \right)^\top.$$

1117 Thus we completes the proof.

1121 D.2 PROOF OF LEMMA 2

Lemma 9 (Lemma 2 Restated: Approximate $p_i(j)$). Let upper bounds $B_v, B_x > 0$ such that for any $k \in [K], j \in [N]$ and $i \in [n], ||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. For any $j \in [N], i \in [n]$, define

$$B_r^j := \max_{|t| \le B_v B_r^{j-1}} |r(t)|, \quad B_r^0 := B_x, \quad \text{and} \quad B_r := \max_j B_r^j.$$

1129 Let function r(t) be $(\epsilon_r, R_1, M_1, C_1)$ -approximable for $R_1 = \max\{B_v B_r, 1\}, M_1 \leq \widetilde{\mathcal{O}}(C_1^2 \epsilon_r^{-2})$, 1130 where C_1 depends only on R_1 and the C^2 -smoothness of r. Then, for any $\epsilon_r > 0$, there exist N1131 attention layers $\operatorname{Attn}_{\theta_1}, \ldots, \operatorname{Attn}_{\theta_N}$ such that for any input $h_i \in \mathbb{R}^D$ takes from (2.1), they map

1132
1133
$$h_i = [x_i; y_i; w; \bar{p}_i(1); \dots; \bar{p}_i(j-1); \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{Attn}_{\theta_j}} \widetilde{h_i} = [x_i; y_i; w; \bar{p}_i(1); \dots; \bar{p}_i(j); \mathbf{0}; 1; t_i],$$

where $\bar{p}_i(j)$ is approximation for $p_i(j)$ (Definition 1). In the expressions of h_i and \bar{h}_i , the dimension of 0 differs. Specifically, the 0 in h_i is larger than in \tilde{h}_i . The dimensional difference between these 0 vectors equals the dimension of $\bar{p}_i(j)$. Suppose function r is L_r -smooth in bounded domain \mathcal{W} , then for any $i \in [n+1], j \in [N], \bar{p}_i(j)$ such that

1139 1140

1141 1142

1152 1153 1154

1134

$$\bar{p}_i(j) = p_i(j) + \epsilon(i,j), \quad \|\epsilon(i,j)\|_2 \le \begin{cases} (\sum_{l=0}^{j-1} K^{l/2} L_r^l B_v^l) \sqrt{K} \epsilon_r , & 1 \le j \le N-1 \\ (\sum_{l=0}^{N-1} K^{l/2} L_r^l B_v^l) \sqrt{d} \epsilon_r , & j = N \end{cases}$$

Additionally, for any $j \in [N]$, the norm of parameters B_{θ_i} defined as (C.1) such that

$$B_{\theta_i} \leq 1 + KC_1$$

1149 *Proof of Lemma* 2. First we need to give a approximation for activation function r(t). By our 1150 assumption and Definition 4, r(t) is $(\epsilon_r, R_1, M_1, C_1)$ -approximable by sum of ReLUs, there exists: 1151

$$\bar{r}(t) = \sum_{m=1}^{M_1} c_m^1 \sigma(\langle a_m^1, [t; 1] \rangle) \text{ with } \sum_{m=1}^{M_1} \left| c_m^1 \right| \le C_1, \ \|a_m^1\|_1 \le 1, \ \forall m \in [M_1],$$
(D.5)

such that $\sup_{t \in [-R_1,R_1]} |\bar{r}(t) - r(t)| \le \epsilon_r$. Let $\bar{p}_i(0) := p_i(0) = x_i$. Similar to $p_i(j)$ follows Definition 1, we pick $\bar{p}_i(j)$ such that for any $j \in [N]$,

$$\bar{p}_i(j)[k] := \bar{r}(v_{j_k}^\top \bar{p}_i(j-1)).$$
 (D.6)

Fix any $j \in [N]$, suppose the input sequences $h_i = [x_i; y_i; w; \bar{p}_i(1); \dots; \bar{p}_i(j-1); 0; 1; t_i]$. Then for every $m \in [M_1], k \in [K]$ (or $k \in [d]$ if j = N), we define matrices $Q_{m,k}^j, K_{m,k}^j, V_{m,k}^j \in \mathbb{R}^{D \times D}$ such that for all $i \in [n+1]$,

1163 1164 1165

1166 1167

1158 1159

$$Q_{m,k}^{j}h_{i} = \begin{bmatrix} a_{m}^{1}[1] \cdot \bar{p}_{i}(j-1) \\ a_{m}^{1}[2] \\ \mathbf{0} \end{bmatrix}, \quad K_{m,k}^{j}h_{i} = \begin{bmatrix} v_{j_{k}} \\ 1 \\ \mathbf{0} \end{bmatrix}, \quad V_{m,k}^{j}h_{i} = c_{m}^{1}e_{j,k}^{1}, \qquad (D.7)$$

1168 where $e_{j,k}^1$ denotes the position unit vector of element $\bar{p}_i(j)[k]$ because this position only depends on 1169 j, k. Since input $h_i = [x_i; y_i; w; \bar{p}_i(1); \ldots; \bar{p}_i(j-1); \mathbf{0}; 1; t_i]$, those matrices indeed exist. In fact, it 1170 is simple to check that

are suffice to (D.7). $I_K(j)$, $I_K(j,k)$, $c_m^1(j,k)$ represents their positions are related to variables in parentheses. In Addition, by (C.1), notice that they have operator norm bounds

1185
$$\max_{j,m,k} \|Q_{m,k}^j\|_1 \le 1, \quad \max_{j,m,k} \|K_{m,k}^j\|_1 \le 1, \quad \max_j \sum_{k,m} \|V_{m,k}^j\|_1 \le KC_1$$
1186

1187

1184

Consequently, for any $j \in [N]$, $B_{\theta_j} \leq 1 + C_1$.

By our construction follows (D.7), a simple calculation shows that

$$\sum_{\substack{m \in [M_1], k \in [K] \\ m \in [M_1], k \in [M_1]$$

Therefore, by definition of ReLU Attention layer follows Definition 7, the output h_i becomes

Therefore, let the attention layer $\theta_j = \{(Q_{m,k}^j, K_{m,k}^j, V_{m,k}^j)\}_{(k,m)}$, we construct $Attn_{\theta_j}$ such that

$$h_i = [x_i; y_i; w; \bar{p}_i(1); \dots; \bar{p}_i(j-1); \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{Attn}_{\theta_j}} \widetilde{h_i} = [x_i; y_i; w; \bar{p}_i(1); \dots; \bar{p}_i(j); \mathbf{0}; 1; t_i].$$

In addition, by setting $R_1 = \max\{B_v B_r, 1\}$, the lemma then follows directly by induction on j. For the base case j = 1, it holds

$$\begin{aligned} |\bar{p}_i(1)[k] - p_i(1)[k]| &= \left|\bar{r}_i(v_{1_k}^\top x_i)[k] - r(v_{1_k}^\top x_i)\right| & (By \text{ Definition 1}) \\ &\leq \epsilon_r. & (By \text{ definition of } \bar{r} \text{ follows (D.5)}) \end{aligned}$$

(By inductive hypothesis)

 $\leq \left|\bar{p}_i(j)[k] - r(v_{j_k}^\top \bar{p}_i(j-1))\right| + \left|r(v_{j_k}^\top \bar{p}_i(j-1)) - p_i(j)[k]\right|$ (By triangle inequality)

Suppose the claim holds for iterate j - 1 and function r is L_r -smooth in bounded domain W. Then for iterate j,

$$\leq \epsilon_r + L_r \|v_{j_k}^{\top}\|_2 \|\bar{p}_i(j-1) - p_i(j-1)\|_2 \qquad (By (D.5) \text{ and } Cauchy–Schwarz inequality))$$

$$\leq \epsilon_r + \sqrt{K} L_r B_v (\epsilon_r \sum_{l=0}^{j-2} K^{l/2} L_r^l B_v^l) \qquad (By inductive hypothesis)$$

$$\leq \epsilon_r \sum_{l=0}^{j-1} K^{l/2} L_r^l B_v^l,$$

 $\left|\bar{p}_i(j)[k] - p_i(j)[k]\right|$

Thus, it holds

1242
1243
1244
1245

$$\leq \sqrt{K} (\epsilon_r \sum_{l=0}^{j-1} K^{l/2} L_r^l B_v^l).$$

1246 This finish the induction. Then for the output layer j = N, it holds

$$\|\bar{p}_{i}(N) - p_{i}(N)\|_{2} = \sqrt{\sum_{k=1}^{d} |\bar{p}_{i}(N)[k] - p_{i}(N)[k]|^{2}}$$
$$\leq \sqrt{d} (\epsilon_{r} \sum^{N-1} K^{l/2} L_{r}^{l} B_{v}^{l}).$$

l = 0

Thus we complete the proof.

1258 D.3 PROOF OF LEMMA 3

Lemma 10 (Lemma 3 Restated: Approximate $r'_i(j)$). Let upper bounds $B_v, B_x > 0$ such that for any $k \in [K], j \in [N]$ and $i \in [n], ||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. For any $j \in [N], i \in [n]$, define

$$B_r'^j := \max_{|t| \le B_v B_{r'}^{j-1}} |r'(t)|, \quad B_{r'}^0 := B_x, \quad \text{and} \quad B_{r'} := \max_j B_{r'}^j.$$

Suppose function r'(t) is $(\epsilon_{r'}, R_2, M_2, C_2)$ -approximable for $R_2 = \max\{B_v B_{r'}, 1\}, M_2 \leq \widetilde{\mathcal{O}}(C_2^2 \epsilon_r'^{-2})$, where C_2 depends only on R_2 and the C^2 -smoothness of r'. Then, for any $\epsilon_r > 0$, there exist an attention layer $\operatorname{Attn}_{\theta_{N+1}}$ such that for any input $h_i \in \mathbb{R}^D$ takes from (3.12), it maps

$$h_i = [x_i; y_i; w; \bar{p}_i; \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{Attn}_{\theta_{N+1}}} \widetilde{h_i} = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \mathbf{0}; 1; t_i],$$

where $\bar{r}'_i(j)$ is approximation for $r'_i(j)$ (Definition 2) and $\bar{r}'_i := [\bar{r}'_i(0); \ldots; \bar{r}'_i(N-1)] \in \mathbb{R}^{(N-2)K+d}$. Similar to Lemma 2, in the expressions of h_i and \tilde{h}_i , the dimension of 0 differs. In addition, let E_r be defined in (3.11), for any $i \in [n+1], j \in [N], k \in [K], \bar{r}'_i(j)$ such that

$$\overline{r}'_i(j-1)[k] = r'_i(j-1)[k] + \epsilon(i,j,k), \quad |\epsilon(i,j,k)| \le \epsilon_{r'} + L_{r'}B_vE_r\epsilon_r,$$

where ϵ_r denotes the error generated in approximating r by sum of ReLUs \bar{r} follows (D.5). Additionally, the norm of parameters $B_{\theta_{N+1}}$ defined as (C.1) such that $B_{\theta_{N+1}} \leq 1 + K(N-1)C_2$.

Proof of Lemma 3. By Definition 2, recall that for any $j \in [N], i \in [n+1], k \in [K]$,

$$r'_{i}(j)[k] = r'(v_{j+1_{k}}^{\top} p_{i}(j)).$$
(D.9)

Therefore we need to give a approximation for r'. By our assumption and Definition 4, r'(t) is $(\epsilon_{r'}, R_2, M_2, C_2)$ -approximable by sum of relus. In other words, there exists:

$$\bar{r}'(t) = \sum_{m=1}^{M_2} c_m^2 \sigma(\langle a_m^2, [t;1] \rangle) \text{ with } \sum_{m=1}^{M_2} \left| c_m^2 \right| \le C_2, \ \|a_m^2\|_2 \le 1, \ \forall m \in [M_2], \tag{D.10}$$

such that $\sup_{t \in [-R_2, R_2]} |\bar{r}'(t) - r'(t)| \le \epsilon_{r'}$. Similar to (D.9), we pick $\bar{r}'_i(j)$ such that

$$\bar{r}'_i(j)[k] := \bar{r}'(v_{j+1_k}^\top \bar{p}_i(j)).$$
(D.11)

To ensure (D.11), we construct our attention layer as follows: for every $j \in [N], m \in [M_2], k \in [K]$, we define matrices $Q_{j,m,k}^{N+1}, K_{j,m,k}^{N+1}, V_{j,m,k}^{N+1} \in \mathbb{R}^{D \times D}$ such that

$$Q_{j,m,k}^{N+1}h_i = \begin{bmatrix} a_m^2[1] \cdot \bar{p}_i(j-1) \\ a_m^2[2] \\ \mathbf{0} \end{bmatrix}, \quad K_{j,m,k}^{N+1}h_i = \begin{bmatrix} v_{j_k} \\ 1 \\ \mathbf{0} \end{bmatrix}, \quad V_{j,m,k}^{N+1}h_i = c_m^2 e_{j,k}^2 , \qquad (D.12)$$

for all $i \in [n+1]$ and $e_{i,k}^2$ denotes the position unit vector of element $\bar{r}'_i(j)[k]$. Since input $h_i = [x_i; y_i; w; \bar{p}_i; 0; 1; t_i]$, similar to (D.8), those matrices indeed exist. In addition, they have operator norm bounds

$$\max_{j,m,k} \|Q_{j,m,k}^{N+1}\|_1 \le 1, \quad \max_{j,m,k} \|K_{j,m,k}^{N+1}\|_1 \le 1, \quad \sum_{j,m,k} \|V_{j,m,k}^{N+1}\|_1 \le K(N-1)C_2.$$

Consequently, by definition of parameter norm follows (C.1), $B_{\theta_{N+1}} \leq 1 + K(N-1)C_2$.

A simple calculation shows that

$$\begin{split} &\sum_{j \in [N], m \in [M_2], k \in [K]} \sigma(\langle Q_{j,m,k}^{N+1} h_i, K_{j,m,k}^{N+1} h_s \rangle) V_{j,m,k}^{N+1} h_s \\ &= \sum_{j=1}^N \sum_{k=1}^K \sum_{m=1}^{M_2} c_m^2 \sigma(\langle a_m^2, [v_{j_k}^\top \bar{p}_i(j-1); 1] \rangle) e_{j,k}^2 \qquad \text{(By our construction follows (D.12))} \\ &= \sum_{j=1}^N \sum_{k=1}^K (\bar{r}'(v_{j_k}^\top \bar{p}_i(j-1))) e_{j,k}^2 \qquad \text{(By definition of } \bar{r}' \text{ follows (D.5))} \\ &= [\mathbf{0}; \bar{r}'_i(0); \dots; \bar{r}'_i(N-1); \mathbf{0}] \qquad \text{(By definition of } \bar{r}'_i(j) \text{ follows (D.11))} \\ &= [\mathbf{0}; \bar{r}'_i; \mathbf{0}], \qquad \text{(By definition of } \bar{r}'_i) \end{split}$$

Therefore, by definition of ReLU Attention layer follows Definition 7, the output \tilde{h}_i becomes

$$\begin{split} \widetilde{h}_{i} &= [\operatorname{Attn}_{\theta_{N}}(h_{i})] \\ &= h_{i} + \frac{1}{n+1} \sum_{s=1}^{n+1} \sum_{j \in [N-1], m \in [M_{2}], k \in [K]}^{n} \sigma(\langle Q_{j,m,k}^{N} h_{i}, K_{j,m,k}^{N} h_{s} \rangle) V_{j,m,k}^{N} h_{s} \\ &= h_{i} + \frac{1}{n+1} \sum_{s=1}^{n+1} (n+1) [\mathbf{0}; \vec{r}_{i}'; \mathbf{0}] \\ &= [x_{i}; y_{i}; w; \bar{p}_{i}; \mathbf{0}; 1; t_{i}] + [\mathbf{0}; \vec{r}_{i}'; \mathbf{0}] \\ &= [x_{i}; y_{i}; w; \bar{p}_{i}; \vec{r}_{i}'; \mathbf{0}; 1; t_{i}]. \end{split}$$

Next, we calculate the error accumulation in this approximation layer. By our assumption, $R_2 =$ $\max\{B_v B_{r'}, 1\}$. Thus, for any $j \in [N], k \in [K], i \in [n+1]$, it holds

$$v_{j_k}^\top p_i(j-1) \le R_2.$$

As our assumption, we suppose function r' is L_r -smooth in bounded domain W. Combining above, the upper bound of error accumulation $|\bar{r}'_i(j)[k] - r'_i(j)[k]|$ becomes

 $\left|\bar{r}_{i}'(j)[k] - r_{i}'(j)[k]\right|$ $\leq \left| \bar{r}'_{i}(j)[k] - r'(v_{j_{k}}^{\top}\bar{p}_{i}(j-1)) \right| + \left| r'(v_{j_{k}}^{\top}\bar{p}_{i}(j-1)) - r'_{i}(j)[k] \right|$ (By triangle inequality) $\leq \epsilon_{r'} + L_{r'} \|v_{j_k}^{\top}\|_2 \|\bar{p}_i(j-1) - p_i(j-1)\|_2$ (By (D.10) and Cauchy–Schwarz inequality) $\leq \epsilon_{r'} + L_{r'} B_v E_r \epsilon_r.$ (By definition of E_r follows (3.11))

Thus we complete the proof.

D.4 PROOF OF LEMMA 4

Lemma 11 (Lemma 4 Restated: Approximate $\partial_1 \ell(p_i(N), y_i)$). Let upper bounds $B_v, B_x, > 0$ such that for any $k \in [K], j \in [N]$ and $i \in [n], ||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. For any $k \in [d]$, suppose function u(t, y)[k] be $(\epsilon_l, R_3, M_3^k, C_3^k)$ -approximable for $R_3 = \max\{B_v B_r, B_y, 1\}, M_3 \leq 1$ $\widetilde{\mathcal{O}}((C_3^k)^2 \epsilon_l^{-2})$, where C_3^k depends only on R_3^k and the C^3 -smoothness of u(t,y)[k]. Then, there exists an MLP layer MLP_{θ_{N+2}} such that for any input sequences $h_i \in \mathbb{R}^D$ takes from (3.14), it maps

$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{MLP}_{\theta_{N+2}}} \widetilde{h_i} = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \mathbf{0}; 1; t_i],$$

where $g_i \in \mathbb{R}^d$ is an approximation for $u(p_i(N), y_i)$. For any $k \in [d]$, assume $u(p_i(N), y_i)$ is L_l -Lipschitz continuous. Then the approximation g_i such that,

$$g_i[k] = u(p_i(N), y_i)[k] + \epsilon(i, k), \text{ with } |\epsilon(i, k)| \le \epsilon_l + L_l E_r \epsilon_r.$$

Additionally, the parameters θ_{N+2} such that $B_{\theta_{N+2}} \leq \max\{R_3 + 1, C_3\}$.

Proof of Lemma 4. By our assumption and Definition 4, for any $k \in [d]$, function u[k](t,y) is $(\epsilon_l, R_3, M_3^k, C_3^k)$ -approximable by sum of relus, there exists :

$$g_k(t,y) = \sum_{m=1}^{M_3^k} c_m^{3,k} \sigma(\langle a_m^{3,k}, [t;y;1] \rangle) \text{ with } \sum_{m=1}^{M_3^k} \left| c_m^{3,k} \right| \le C_3, \ \|a_m^{3,k}\|_2 \le 1, \ \forall m \in [M_3^k], \quad (D.13)$$

such that $\sup_{(t,y)\in [-R_3,R_3]^2} |g_k(t,y) - u[k](t,y)| \le \epsilon_l$. Then we construct our MLP layer.

Let $M_3 := \sum_{k=1}^d M_3^k$, we pick matrices $W_1^{N+1} \in \mathbb{R}^{M_3 \times D}, W_2^{N+1} \in \mathbb{R}^{D \times M_3}$ such that for any $i \in [n+1], m \in [M_3]$,

$$W_{1}^{N+1}h_{i} = \begin{bmatrix} a_{1}^{3,1}[1] \cdot \bar{p}_{i}(N) + a_{1}^{3,1}[2] \cdot y_{i} + a_{1}^{3,1}[3] - R_{3}(1 - t_{i}) \\ \vdots \\ a_{M_{3}^{1}}^{3,1}[1] \cdot \bar{p}_{i}(N) + a_{M^{1}}^{3,1}[2] \cdot y_{i} + a_{M_{3}^{1}}^{3,1}[3] - R_{3}(1 - t_{i}) \\ \vdots \\ a_{1}^{3,d}[1] \cdot \bar{p}_{i}(N) + a_{1}^{3,d}[2] \cdot y_{i} + a_{1}^{3,d}[3] - R_{3}(1 - t_{i}) \\ \vdots \\ a_{M_{3}^{3,d}}^{3,d}[1] \cdot \bar{p}_{i}(N) + a_{M^{d}}^{3,d}[2] \cdot y_{i} + a_{M_{3}^{3,d}}^{3,d}[3] - R_{3}(1 - t_{i}) \end{bmatrix} \in \mathbb{R}^{M_{3}},$$

$$W_2^{N+1}[j,m] = c_m^{3,k} \cdot 1\{j = D_g^k, M_3^{k-1} < m \le M_3^k\},$$
(D.14)

where D_a^k denotes the position of element $g_i[k]$. Since input $h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \mathbf{0}; 1; t_i]$, similar to (D.8), those matrices indeed exist. Furthermore, by (C.1), they have operator norm bounds

 $||W_1^{N+1}||_1 \le R_3 + 1, \quad ||W_2^{N+1}||_1 \le C_3$

Consequently, $B_{\theta_{N+2}} \le \max\{R_3 + 1, C_3\}.$

By our construction (D.14), a simple calculation shows that

$$W_2^{N+1}\sigma(W_1^{N+1}h_i) = \sum_{k=1}^d \sum_{m=1}^{M_3^k} \sigma(\langle a_m^{3,k}, [\bar{p}_i(N); y_i; 1] \rangle - R_3(1-t_i)) \cdot c_m^{3,k} e_{D_g^k}$$

1404
 0

 1405

$$g_1(\bar{p}_i(N), y_i)$$

 1406
 \vdots

 1407
 $g_d(\bar{p}_i(N), y_i)$

 1408
 0

For $k \in [d]$, we let $g_i[k] = 1\{t_j = 1\} \cdot g_k(\bar{p}_i(N), y_i)e_{D_a^k}$ for $i \in [n+1]$. Hence, $\text{MLP}_{\theta_{N+2}}$ maps

$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; \mathbf{0}; 1; t_i] \xrightarrow{\operatorname{MLP}_{\theta_{N+2}}} \widetilde{h_i} = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \mathbf{0}; 1; t_i],$$

Next, we calculate the error generated in this approximation. By setting $R_3 = \max\{B_v B_r, B_u, 1\}$, for any $i \in [n+1]$, it holds

$$p_i(N) \le R_3, \quad y_i \le R_3$$

Moreover, as our assumption, we suppose function $\partial_1 \ell$ is L_l -smooth in bounded domain \mathcal{W} . Therefore, by the definition of the function g, for each $i \in [n]$, the error becomes

 $|q_i[k] - u(p_i(N), y_i)[k]|$ $\leq |g_i[k] - u(\bar{p}_i(N), y_i)[k]| + |u(\bar{p}_i(N), y_i)[k] - u(p_i(N), y_i)[k]|$ (By triangle inequality) $\leq \epsilon_l + L_l \|\bar{p}_i(N) - p_i(N)\|_2$ (By the definition of g_k follows (D.13) and L_l -smooth assumption) $<\epsilon_l + L_l E_r \epsilon_r$,. (By the definition of E_r follows (3.11))

Combining above, we complete the proof.

PROOF OF LEMMA 5 D.5

Lemma 12 (Lemma 5 Restated: Approximate $\bar{s}_t(j)$). Recall that $s_i(j) = r'_i(j-1) \odot (V_{j+1}^\top \cdot s_i(j+1))$ follows Definition 3. Let the initial input takes from (3.16). Then, there exist N element-wise multiplication layers: EWML_{θ_{N+3}},..., EWML_{θ_{2N+2}} such that for input sequences, $j \in [N]$,

$$h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \bar{s}_i(N); \dots; \bar{s}_i(j+1); \mathbf{0}; 1; t_i],$$

they map $\text{EWML}_{\theta_{2N+3-j}}(h_i) = [x_i; y_i; w; \overline{p}_i; \overline{r}'_i; g_i; \overline{s}_i(N); \dots; \overline{s}_i(j); \mathbf{0}; 1; t_i]$, where the approxi-mation $\bar{s}_i(j)$ is defined as recursive form: for any $i \in [n+1], j \in [N-1]$,

$$\bar{s}_i(j) := \begin{cases} \bar{r}'_i(j-1) \odot (V_{j+1}^\top \cdot \bar{s}_i(j+1)), & j \in [N-1] \\ \bar{r}'_i(N-1) \odot g_i, & j = N. \end{cases}$$

Additionally, for any $j \in [N]$, $B_{\theta_{N+2+j}}$ defined in (C.1) satisfies $B_{\theta_{N+2+j}} \leq 1$.

Proof of Lemma 5. We give the construction of parameters directly. For every $j \in [N-1], k \in [K]$, we define matrices $Q_k^{2N+3-j}, K_k^{2N+3-j}, V_k^{2N+3-j} \in \mathbb{R}^{D \times D}$ such that for all $i \in [n+1]$,

$$\begin{array}{l} \mathbf{1451} \\ \mathbf{1452} \\ \mathbf{1453} \\ \mathbf{1454} \\ \mathbf{1455} \\ \mathbf{1455} \\ \mathbf{1456} \end{array} \qquad Q_k^{2N+3-j} h_i = \begin{bmatrix} v_{j+1_1}[k] \\ \vdots \\ v_{j+1_K}[k] \\ \mathbf{0} \end{bmatrix}, \quad K_k^{2N+3-j} h_i = \begin{bmatrix} \bar{s}_i(j+1) \\ \mathbf{0} \end{bmatrix}, \quad V_k^{2N+3-j} h_i = \bar{r}_i'(j-1)[k] \cdot e_{j,k}^3,$$

$$\begin{array}{l} \mathbf{1456} \\ \mathbf{0} \end{bmatrix}$$

$$\begin{array}{l} \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$

where e_{ik}^3 denotes the position unit vector of element $\bar{s}_i(j)[k]$.

1458 Since input $h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \bar{s}_i(N); \dots; \bar{s}_i(j+1); \mathbf{0}; 1; t_i]$, similar to (D.8), those matrices 1459 indeed exist. Thus, it is straightforward to check that 1460 $\sum_{k \in [K]} \gamma(\langle Q_k^{2N+3-j} h_i, K_k^{2N+3-j} h_i \rangle) V_k^{2N+3-j} h_i$ 1461 1462 1463 $=\sum_{k=1}^{K} (V_{j+1}^{\top}[k,*] \cdot \bar{s}_i(j+1)) \bar{r}'_i(j-1)[k] e_{j,k}^3 \quad (\text{By definition of EWML layer follows Definition 6})$ 1464 1465 1466 $= \begin{bmatrix} \mathbf{0} \\ r_i(j-1)[1]V_{j+1}^{\top}[1,*] \cdot \bar{s}_i(j+1) \\ \vdots \\ \bar{r}'_i(j-1)[k]V_{j+1}^{\top}[K,*] \cdot \bar{s}_i(j+1) \end{bmatrix}$ 1467 1468 (By definition of $e_{i,k}^3$) 1469 1470 1471 1472 $= \begin{bmatrix} \mathbf{0} \\ \bar{r}'_i(j-1) \odot (V_{j+1}^\top \cdot \bar{s}_i(j+1)) \\ \mathbf{0} \end{bmatrix}$ 1473 (By definition of hadamard product) 1474 1475 $= [\mathbf{0}; \bar{s}_i(j); \mathbf{0}].$ (By definition of $\bar{s}_i(j)$ follows (3.19)) 1476 1477 Therefore, by the definition of EWML layer follows Definition 6, the output h_i becomes 1478 1479 $\widetilde{h}_i = [\operatorname{Attn}_{\theta_{2N+3-i}}(h_i)]$ 1480 $=h_i + \sum_{m \in [2], k \in [K]} \sigma(\langle Q_{m,k}^{2N+3-j} h_i, K_{m,k}^{2N+3-j} h_s \rangle) V_{m,k}^{2N+3-j} h_s$ 1481 1482 1483 $= h_i + [\mathbf{0}; s(j); \mathbf{0}]$ 1484 $= [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; \bar{s}_i(N-1); \dots; \bar{s}_i(j+1); \mathbf{0}; 1; t_i] + [\mathbf{0}; \bar{s}_i(j); \mathbf{0}]$ 1485 $= [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; q_i; \bar{s}_i(N-1); \dots; \bar{s}_i(j); \mathbf{0}; 1; t_i].$ 1486 1487 Finally we come back to approximate the initial approximation $\bar{s}_i(N) = \bar{r}'_i(N-1) \odot q_i$. Notice 1488 that g_i and $r'_i(N-1)$ are already in the input $h_i = [x_i; y_i; w; \bar{p}_i; \bar{r}'_i; g_i; 0; 1; t_i]$, thus it is simple to 1489 construct EWML_{N+3}, similar to (D.15), such that it maps, 1490 1491 $[x_i: y_i: w: \bar{p}_i: \bar{r}': q_i: \mathbf{0}: 1: t_i] \xrightarrow{\text{EWML}_{N+3}} [x_i: y_i: w: \mathbf{0}: 1: \bar{p}_i: \bar{r}': q_i: \bar{s}_i(N): \mathbf{0}: 1: t_i].$

1494

1500

1502

Since we don't using the sum of ReLU to approximate any variables, these step don't generate extra error. Besides, by (3.18), matrices have operator norm bounds

$$\max_{j,k} \|Q_k^{N+2+j}\|_1 \le 1, \quad \max_{j,k} \|K_k^{N+2+j}\|_1 \le 1, \quad \max_{j,k} \|V_k^{N+2+j}\|_1 \le 1.$$

1499 Consequently, for any $j \in [N]$, $B_{\theta_{N+2+j}} \leq 1$. Thus we complete the proof.

1501 D.6 PROOF OF LEMMA 6

Lemma 13 (Lemma 6 Restated: Error for $g_i \bar{s}_i(j)$). Suppose the upper bounds $B_v, B_x > 0$ such that for any $k \in [K], j \in [N]$ and $i \in [n], ||v_{j_k}||_2 \leq B_v$, and $||x_i||_2 \leq B_x$. Let $r'_i(j) \in \mathbb{R}^K$ such that $r'_i(j)[k] := r'(v_{j+1_k}^\top p_i(j))$ follows Definition 2. Let $s_i(j) = R_i(j - 1)V_{j+1}^\top \dots R_i(N-2)V_N^\top \cdot R_i(N-1)u$ follows Definition 3. Let $\bar{r}'_i(j), g_i, \bar{s}_i(j)$ be the approximations for $r'_i(j), u(p_i(N), y_i), s_i(j)$ follows Lemma 3, Lemma 4 and Lemma 5 respectively. Let $B_{r'}$ be the upper bound of $\bar{r}'_i(j)[k]$ and $r'_i(j)[k]$ as defined in Lemma 3. Let B_l be the upper bound of $g_i[k]$ and $u(p_i(N), y_i)[k]$ as defined in Lemma 4. Then for any $i \in [n+1], j \in [N], k \in [K]$,

1511
$$\bar{s}_i(j)[k] \le B_s,$$

$$|\bar{s}_i(j)[k] - s_i(j)[k]| \le E_s^r \epsilon_r + E_s^{r'} \epsilon_{r'} + E_s^l \epsilon_l,$$

1515 where

$$\begin{split} P &:= \max\{\sqrt{K}, \sqrt{d}\}\\ B_s &:= \max_{j \in [N]} \{(P \cdot B_{r'} B_v)^{N-j} B_{r'} B_l\},\\ E_s^r &:= \max_{j \in [N]} \{L_{r'} E_r P B_s B_v^2 [\sum_{l=0}^{N-j-1} (B_{r'} B_v P)^l] + (B_{r'} B_v P)^{N-j} (B_l L_{r'} B_v E_r + B_{r'} L_l E_r)\},\\ E_s^{r'} &:= \max_{j \in [N]} \{P B_s B_v [\sum_{l=0}^{N-j-1} (B_{r'} B_v P)^l] + (B_{r'} B_v P)^{N-j} B_l\},\\ E_s^l &:= \max_{j \in [N]} \{(B_{r'} B_v P)^{N-j} B_{r'}\}. \end{split}$$

Above, B_s is the upper bound of $\bar{s}_i(j)[k]$ and $E_s^r, E_s^{r'}, E_s^l$ are the coefficients of $\epsilon_r, \epsilon'_r, \epsilon_l$ in the upper bounds of $|\bar{s}_i(j)[k] - s_i(j)[k]|$, respectively.

1531 Proof of Lemma 6. We use induction to prove the first two statements. To begin with, we illustrate 1532 the recursion formula for $\bar{s}_i(j)$. By (3.19), recall that for any $j \in [N]$,

$$\bar{s}_i(j) := \begin{cases} \bar{r}'_i(j-1) \odot (V_{j+1}^\top \cdot \bar{s}_i(j+1)), & j \in [N-1] \\ \bar{r}'_i(N-1) \odot g_i, & j = N. \end{cases}$$

We consider applying induction to prove the first statement:

$$\bar{s}_i(j)[k] \le (P \cdot B_{r'} B_v)^{N-n} B_{r'} B_l$$

1540 As for the base case, j = N:

$$\bar{s}_i(N)[k] = \bar{r}'_i(N-1)[k] \cdot g_i[k] \le B_{r'}B_l.$$

1544 Therefore, if the statement holds for j = n + 1, by (3.19) and our assumption, it holds

$$\begin{split} \bar{s}_{i}(n)[k] &= \bar{r}'_{i}(n-1)[k] \cdot (v_{j+1_{k}}^{\top} \bar{s}_{i}(n+1)) & (\text{By recursion formula (3.19)}) \\ &\leq \bar{r}'_{i}(n-1)[k] \cdot \|v_{n+1_{k}}\|_{2} \cdot \|\bar{s}_{i}(n+1)\|_{2} & (\text{By Cauchy-schwarz inequality}) \\ &\leq \bar{r}'_{i}(n-1)[k] \cdot \|v_{n+1_{k}}\|_{2} \cdot \max\{\sqrt{K}, \sqrt{d}\} \cdot \max_{k} |\bar{s}_{i}(n+1)[k]| \\ &\leq (B_{r'}B_{v}) \cdot \max\{\sqrt{K}, \sqrt{d}\} \cdot (\max\{\sqrt{K}, \sqrt{d}\} \cdot B_{r'}B_{v})^{N-n-1}B_{r'}B_{l} \\ & (\text{By inductive hypothesis}) \\ &= (P \cdot B_{r'}B_{v})^{N-n}B_{r'}B_{l}. & (\text{By definition of } P \text{ follows (3.22)}) \end{split}$$

Thus, by the principle of induction, the first statement is true for all integers $j \in [N]$. Moreover, by the definition of B_s follows (3.23), we know B_s is the upper bound of $\bar{s}_i(j)[k]$. Next we apply induction to prove the second statement:

$$\begin{aligned} |\bar{s}_i(j)[k] - s_i(j)[k]| &\leq (\epsilon_{r'} + L_{r'}B_vE_r\epsilon_r)PB_vB_s[\sum_{l=0}^{N-n-1} (B_{r'}B_vP)^l] \\ &+ (B_{r'}B_vP)^{N-n}[(B_lL_{r'}B_vE_r + B_{r'}L_lE_r)\epsilon_r + B_l\epsilon_{r'} + B_{r'}\epsilon_l]. \end{aligned}$$

1563 For the base case, j = N:

$$\begin{aligned} & |\bar{s}_i(N)[k] - s_i(N)[k]| \\ &= |\bar{r}'_i(N-1)[k] \cdot g_i[k] - r'_i(N-1)[k] \cdot u(p_i(N), y_i)[k]| \end{aligned}$$
 (By definition (3.19) and (3.7))

1566 $\leq |\bar{r}'_{i}(N-1)[k] - r'_{i}(N-1)[k]| \cdot |g_{i}[k]| + |r'_{i}(N-1)[k]| \cdot |g_{i}[k] - u(p_{i}(N), y_{i})[k]|$ 1567 (By triangle inequality) 1568 $\leq (\epsilon_{r'} + L_{r'}B_v E_r \epsilon_r)B_l + B_{r'}(\epsilon_l + L_l E_r \epsilon_r).$ (By (3.13) and (3.15))1569 $= (B_l L_{r'} B_v E_r + B_{r'} L_l E_r) \epsilon_r + B_l \epsilon_{r'} + B_{r'} \epsilon_l$ 1570 1571 Therefore, if the statement holds for j = n + 1, by (3.19) and our assumption, it holds 1572 1573 $|\bar{s}_i(n)[k] - s_i(n)[k]|$ $|s_i(n)[\kappa] - s_i(n)[n]| = \left| \bar{r}'_i(n-1)[k] \cdot (v_{n+1_k}^\top \bar{s}_i(n+1)) - r'_i(n-1)[k] \cdot (v_{n+1_k}^\top s_i(n+1)) \right|$ (By the recursion formula (3.7) and (3.19)) 1574 1575 1576 $\leq |\bar{r}'_{i}(n-1)[k] - r'_{i}(n-1)[k]| \cdot |v_{n+1}^{\top} \bar{s}_{i}(n+1)|$ 1577 1578 $+ |r'_{i}(n-1)[k]| \cdot |(v_{n+1}^{\top} \bar{s}_{i}(n+1)) - (v_{n+1}^{\top} \bar{s}_{i}(n+1))|$ (By triangle inequality) 1579 $\leq (\epsilon_{r'} + L_{r'} B_v E_r \epsilon_r) P B_v B_s + B_{r'} B_v \|\bar{s}_i(n+1) - s_i(n+1)\|_2$ 1580 (By error accumulation of approximating r' follows (3.13)) 1581 $\leq (\epsilon_{r'} + L_{r'} B_v E_r \epsilon_r) P B_v B_s + B_{r'} B_v P \max_k |\bar{s}_i(n+1)[k] - s_i(n+1)[k]|$ 1582 $\leq (\epsilon_{r'} + L_{r'}B_vE_r\epsilon_r)PB_vB_s + B_{r'}B_vP\Big\{(\epsilon_{r'} + L_{r'}B_vE_r\epsilon_r)PB_vB_s[\sum_{l=0}^{N-n-2}(B_{r'}B_vP)^l]\Big\}$ 1585 $+ (B_{r'}B_vP)^{N-n-1}[(B_lL_{r'}B_vE_r + B_{r'}L_lE_r)\epsilon_r + B_l\epsilon_{r'} + B_{r'}\epsilon_l]\Big\} (By inductive hypothesis)$ 1586 1587 $\leq (\epsilon_{r'} + L_{r'}B_vE_r\epsilon_r)PB_vB_s[\sum_{l=0}^{N-n-1}(B_{r'}B_vP)^l]$ 1589 1590 $+ (B_{r'}B_{v}P)^{N-n} [(B_{l}L_{r'}B_{v}E_{r} + B_{r'}L_{l}E_{r})\epsilon_{r} + B_{l}\epsilon_{r'} + B_{r'}\epsilon_{l}].$ 1591 1592 Thus, by the principle of induction, the second statement is true for all integers $j \in [N-1]$. By the 1593 1594

definition of E_s follows (3.24), it is simple to check that

$$\bar{s}_i(j)[k] - s_i(j)[k]| \le E_s^r \epsilon_r + E_s^{r'} \epsilon_{r'} + E_s^l \epsilon_l.$$

Thus we complete the proof.

1598

1602

1604

1607 1608 1609

1614 1615 1616

1617

161

1595 1596 1597

D.7 **PROOF OF THEOREM 1**

Theorem 4 (Theorem 1 Restated: In-context gradient descent on N-layer NNs). Fix any $B_{\eta}, \eta, \epsilon > 0$ $0, L \ge 1$. For any input sequences takes from (2.1), their exist upper bounds B_x, B_y such that for any $i \in [n]$, $\|y_i\|_2 \leq B_y$, $\|x_i\|_2 \leq B_x$. Assume functions r(t), r'(t) and u(t, y)[k] are $L_r, L_{r'}, L_l$ -Lipschitz continuous. Suppose \mathcal{W} is a closed domain such that for any $j \in [N-1]$ and $k \in [K]$,

$$\mathcal{W} \subset \left\{ w = [v_{j_k}] \in \mathbb{R}^{D_N} : \|v_{j_k}\|_2 \le B_v \right\},\$$

1610 and $\operatorname{Proj}_{\mathcal{W}}$ project w into bounded domain \mathcal{W} . Assume $\operatorname{Proj}_{\mathcal{W}} = \operatorname{MLP}_{\theta}$ for some MLP layer with hidden dimension D_w parameters $\|\theta\| \leq C_w$. If functions r(t), r'(t) and u(t,y)[k] are C^4 -1611 smoothness, then for any $\epsilon > 0$, there exists a transformer model NN_{θ} with (2N + 4)L hidden layers 1612 consists of L neural network blocks $\mathrm{TF}_{\theta}^{N+2} \circ \mathrm{EWML}_{\theta}^{N} \circ \mathrm{TF}_{\theta}^{2}$, 1613

$$NN_{\theta} := TF_{\theta}^{N+2} \circ EWML_{\theta}^{N} \circ TF_{\theta}^{2} \circ \ldots \circ TF_{\theta}^{N+2} \circ EWML_{\theta}^{N} \circ TF_{\theta}^{2},$$

such that the heads number M^l , embedding dimensions D^l , and the parameter norms B_{θ^l} suffice

1618
1619
$$\max_{l \in [(2N+4)L]} M^l \le \widetilde{O}(\epsilon^{-2}), \quad \max_{l \in [(2N+4)L]} D^l \le O(NK^2) + D_w, \quad \max_{l \in [(2N+4)L]} B_{\theta^l} \le O(\eta) + C_w + 1,$$

where $\tilde{O}(\cdot)$ hides the constants that depend on d, K, N, the radius parameters B_x, B_y, B_v and the smoothness of r and ℓ . And this neural network such that for any input sequences $H^{(0)}$, take from (2.1), NN $_{\theta}(H^{(0)})$ implements L steps in-context gradient descent on risk Eqn (2.2): For every $l \in [L]$, the (2N + 4)l-th layer outputs $h_i^{((2N+4)l)} = [x_i; y_i; \overline{w}^{(l)}; \mathbf{0}; 1; t_i]$ for every $i \in [n + 1]$, and approximation gradients $\overline{w}^{(l)}$ such that

$$\overline{w}^{(l)} = \operatorname{Proj}_{\mathcal{W}}(\overline{w}^{(l-1)} - \eta \nabla \mathcal{L}_n(\overline{w}^{(l-1)}) + \epsilon^{(l-1)}), \quad \overline{w}^{(0)} = \mathbf{0},$$

where $\|\epsilon^{(l-1)}\|_2 \leq \eta \epsilon$ is an error term.

Proof of Theorem 1. We consider the first N + 2 transformer layers TF_{θ}^{N+2} are layers in Lemma 2, Lemma 3 and Lemma 4. Then we let the middle N element-wise multiplication layers $EWML_{\theta}^{N}$ be layers in Lemma 5. We only need to check approximability conditions. By Lemma 7 and our assumptions, for any $\epsilon_r, \epsilon_{r'}, \epsilon_l$, it holds

- Function r(t) is $(\epsilon_r, R_1, M_1, C_1)$ -approximable for $R_1 = \max\{B_v B_r, 1\}, M_1 \leq \tilde{\mathcal{O}}(C_1^2 \epsilon_r^{-2}),$ where C_1 depends only on R_1 and the C^2 -smoothness of r(t).
- Function r'(t) is $(\epsilon_{r'}, R_2, M_2, C_2)$ -approximable for $R_2 = \max\{B_v B_{r'}, 1\}, M_2 \leq \tilde{\mathcal{O}}(C_2^2 \epsilon_r'^{-2}),$ where C_2 depends only on R_2 and the C^2 -smoothness of r'(t).
- Function $\partial_1 \ell(t, y)$ is $(\epsilon_l, R_3, M_3, C_3)$ -approximable for $R_3 = \max\{B_v B_r, 1\}, M_3 \leq \tilde{\mathcal{O}}(C_3^2 \epsilon_l^{-2}),$ where C_3 depends only on R_3 and the C^3 -smoothness of u(t, y)[k].

which suffice approximability conditions in Lemma 2, Lemma 3 and Lemma 4.

Now we construct the last two layers to implement $w - \eta \nabla \mathcal{L}_n(w)$ and $\operatorname{Proj}_{\mathcal{W}}(w)$. First we construct a attention layer to approximate $w - \eta \nabla \mathcal{L}_n(w)$. For every $m \in [2], j \in [N], k \in [K]$, we consider matrices $Q_{m,j,k}^{2N+3}, j_{m,j,k}^{2N+3} \in \mathbb{R}^{D \times D}$ such that

$$Q_{1,j,k}^{2N+3}h_{i} = \begin{bmatrix} 1\\ \mathbf{0} \end{bmatrix}, \quad K_{1,j,k}^{2N+3}h_{i} = \begin{bmatrix} \bar{s}_{i}(j)[k] \\ \mathbf{0} \end{bmatrix}, \quad V_{1,j,k}^{2N+3}h_{i} = -\frac{\eta(n+1)}{2n} \begin{bmatrix} \mathbf{0} \\ \bar{p}_{i}(j-1) \\ \mathbf{0} \end{bmatrix}, \quad Q_{2,j,k}^{2N+3}h_{i} = \begin{bmatrix} -1\\ \mathbf{0} \end{bmatrix}, \quad K_{2,j,k}^{2N+3}h_{i} = \begin{bmatrix} \bar{s}_{i}(j)[k] \\ \mathbf{0} \end{bmatrix}, \quad V_{2,j,k}^{2N+3}h_{i} = -\frac{\eta(n+1)}{2n} \begin{bmatrix} \mathbf{0} \\ -\bar{p}_{i}(j-1) \\ \mathbf{0} \end{bmatrix}.$$
(D.16)

Furthermore, we define approximation gradient $\nabla_w \overline{\mathcal{L}}_n(w)$ as follows,

$$\nabla_{w} \bar{\mathcal{L}}_{n}(w) := -\frac{1}{\eta(n+1)} \sum_{t=1}^{n+1} \sum_{m \in [2], j \in [N], k \in [K]} \sigma(\langle Q_{m,j,k}^{2N+3}h_{i}, K_{m,j,k}^{2N+3}h_{t}\rangle) V_{m,j,k}^{2N+3}h_{t}$$

$$= \frac{1}{2n} \sum_{t=1}^{n+1} \sum_{k=1}^{K} \sum_{j=1}^{N} (\sigma(\bar{s}_{t}(j)[k]) - \sigma(-\bar{s}_{t}(j)[k])) \begin{bmatrix} \mathbf{0} \\ \bar{p}_{t}(j-1) \\ \mathbf{0} \end{bmatrix}$$
(By our construction (D.16))

$$= \frac{1}{2n} \sum_{t=1}^{n+1} \sum_{k=1}^{K} \sum_{j=1}^{N} \bar{s}_t(j) [k] \cdot \begin{bmatrix} \mathbf{0} \\ \bar{p}_t(j-1) \\ \mathbf{0} \end{bmatrix} \qquad (\text{By } f(x) = \sigma(x) - \sigma(-x))$$

1672
1673
$$= \frac{1}{2n} \sum_{t=1}^{n+1} \sum_{j=1}^{N} \begin{bmatrix} \mathbf{I}_{K \times K} \otimes \bar{p}_t(j-1) \cdot \bar{s}_t(j) \\ \mathbf{0} \end{bmatrix}$$
(By definition of Kronecker product)

 $h_i = [\operatorname{Attn}_{\theta_{2N+3}}(h_i)]$

definition of ReLU attention layer follows Definition 7, for any $i \in [n+1]$,

 $= [x_i; y_i; w - \eta \nabla_w \overline{\mathcal{L}}_n(w); \overline{p}_i; \overline{r}'_i; g_i; \mathbf{0}; 1; t_i].$

where $\bar{A}_t(j) := \mathbf{I}_{K \times K} \otimes \bar{p}_t(j-1) \cdot \bar{s}_t(j)$ denotes the approximation for $A_t(j)$. Therefore, by the

 $= h_{i} + \frac{1}{n+1} \sum_{i=1}^{n+1} \sum_{m \in [2], j \in [N], k \in [K]}^{n+1} \sigma(\langle Q_{m,j,k}^{2N+3}h_{s}, K_{m,j,k}^{2N+3}h_{i} \rangle) V_{m,j,k}^{2N+3}h_{i}$ $= [x_{i}; y_{i}; w; \bar{p}_{i}; \bar{r}_{i}'; g_{i}; \bar{s}_{i}; \mathbf{0}; 1; t_{i}] - \frac{\eta}{2n} \sum_{t=1}^{n} \begin{bmatrix} \mathbf{0} \\ \bar{A}_{t}(1) \\ \vdots \\ \bar{A}_{t}(N) \\ \mathbf{0} \end{bmatrix}$

 Since we do not use approximation technique like Definition 4, this step do not generate extra error. Besides, by (C.1), matrices have operator norm bounds

$$\max_{j,m,k} \|Q_{j,m,k}^{2N+3}\|_1 \le 1, \quad \max_{j,m,k} \|K_{j,m,k}^{2N+3}\|_1 \le 1, \quad \sum_{j,m,k} \|V_{j,m,k}^{2N+3}\|_1 \le 2\eta NK.$$

Consequently, $B_{\theta_{2N+3}} \leq 1 + 2\eta N K$. Fix any $\epsilon > 0$, then we pick appropriate $\epsilon_r, \epsilon'_r, \epsilon_l$ such that

$$\|\epsilon^{(l-1)}\|_2 = \eta \|\nabla_w \overline{\mathcal{L}}_n(\overline{w}^{(l-1)}) - \nabla_w \mathcal{L}_n(\overline{w}^{(l-1)})\|_2 \le \eta \epsilon.$$

By Definition 3 and Lemma 6, for any $j \in [N-1], i \in [n]$, it holds

$$\leq P[(E_s^r\epsilon_r + E_s^{r'}\epsilon_{r'} + E_s^l\epsilon_l)\sqrt{P}B_r + B_sE_r\epsilon_r], \qquad (By (3.11) \text{ and } Lemma 6)$$

п

where B_s is the upper bound of $\bar{s}_i(j)[k]$ and $E_s^r, E_s^{r'}, E_s^l$ are the coefficients of $\epsilon_r, \epsilon_r', \epsilon_l$ in the upper bounds of $|\bar{s}_i(j)[k] - s_i(j)[k]|$ follow Lemma 6, respectively. We can drive similar results as j = N. Actually, by $P = \max\{\sqrt{K}, \sqrt{d}\}$ follows Lemma 6, above inequality also holds for j = N. Therefore, the error in total such that for any w,

Г О

$$\|\nabla_w \overline{\mathcal{L}}_n(w) - \nabla_w \mathcal{L}_n(w)\|_2$$

$$\begin{array}{ccc} \mathbf{1723} & & & \mathbf{0} \\ \mathbf{1724} & & & \\ \mathbf{1725} & & & \\ \mathbf{1726} & & & \\ \mathbf{1727} & & & \\ \mathbf{1727} & & & \\ \end{array} = \|\frac{1}{2n} \sum_{t=1}^{n} \begin{bmatrix} \mathbf{0} \\ \overline{A}_{t}(1) \\ \vdots \\ \overline{A}_{t}(N) \\ \mathbf{0} \end{bmatrix} - \frac{1}{2n} \sum_{t=1}^{n} \begin{bmatrix} \mathbf{0} \\ A_{t}(1) \\ \vdots \\ A_{t}(N) \\ \mathbf{0} \end{bmatrix} \|_{2}$$

$$\begin{bmatrix} 1720 \\ 1727 \end{bmatrix} \begin{bmatrix} t-1 \\ A_t(N) \\ 0 \end{bmatrix} \begin{bmatrix} t-1 \\ A_t(N) \\ 0 \end{bmatrix}$$

(By definition of $\overline{\mathcal{L}}_n(w)$ and $\mathcal{L}_n(w)$)

(By definition of $\nabla_w \overline{\mathcal{L}}_n(w)$)

1728
1729
$$\leq \frac{1}{2} \max_{1 \leq t \leq n} \{\sum_{i=1}^{N} \|\bar{A}_{t}(j) - A_{t}(j)\|_{2}\}$$

1733 1734

1736

1737

1739 1740

$$\leq \frac{N}{2} P[(E_s^r \epsilon_r + E_s^{r'} \epsilon_{r'} + E_s^l \epsilon_l) \sqrt{P} B_r + B_s E_r \epsilon_r].$$

(By the error accumulation results derived before)

1735 Let $C_l, C_r, C_{r'}$ denotes coefficients in front of $\epsilon_l, \epsilon_r, \epsilon_{r'}$ respectively. Then it holds

$$C_l = NP^{\frac{3}{2}}B_rE_s^l,$$

$$C_r = NP^{\frac{3}{2}}B_rE_s^r + NPB_sE_r,$$

$$C_{r'} = NP^{\frac{3}{2}}B_r E_s^{r'}.$$

Thus, to ensure $\|\nabla_w \overline{\mathcal{L}}_n(w) - \nabla_w \mathcal{L}_n(w)\|_2 \le \epsilon$, we only need to select $\epsilon_l, \epsilon_r, \epsilon'_r$ as

$$\epsilon_l = \frac{2\epsilon}{3C_l}, \quad \epsilon_r = \frac{2\epsilon}{3C_r}, \quad \epsilon'_r = \frac{2\epsilon}{3C_{r'}}$$

1743 1744 1745

Therefore, we only need to pick the last MLP layer MLP_{2N+4} such that it maps

1748
1749
$$[x_i; y_i; w - \eta \nabla_w \mathcal{L}_n(w); \bar{p}_i; \bar{r}'_i; g_i; \bar{s}_i; \mathbf{0}; 1; t_i] \xrightarrow{\mathrm{MLP}_{2N+4}} [x_i; y_i; \mathrm{Proj}_{\mathcal{W}}(w - \eta \nabla_w \mathcal{L}_n(w)); \mathbf{0}; 1; t_i].$$

1750 By our assumption on the map $\operatorname{Proj}_{\mathcal{W}}$, this is easy.

Finally, we analyze how many embedding dimensions of Transformers are needed to implement the above ICGD. Recall that

$$\begin{array}{l} \text{1754} \\ \text{1755} \end{array} \quad x_i, y_i \in \mathbb{R}^d, w \in \mathbb{R}^{2dK + (N-2)K^2}, \bar{p}_i \in \mathbb{R}^{(N-1)K+d}, \bar{r}'_i \in \mathbb{R}^{(N-2)K+d}, g_i \in \mathbb{R}^d, \bar{s}_i \in \mathbb{R}^{(N-1)K+d}. \end{array}$$

1756 Therefore, $\max{\{\Omega(NK^2), D_w\}}$ embedding dimensions of Transformer are required to implement 1757 ICGD on deep models.

Combining the above, we complete the proof.

1760

1761 1762

1763 D.8 PROOF OF COROLLARY 1.1

1764

1769

1770

1772

1774

Corollary 4.1 (Corollary 1.1 Restated: Error for implementing ICGD on *N*-layer neural network). Fix $L \ge 1$, under the same setting as Theorem 1, (2N+4)L-layer neural networks NN_{θ} approximates the true gradient descent trajectory $\{w_{GD}^l\}_{l\ge 0} \in \mathbb{R}^{D_N}$ with the error accumulation

$$\overline{w}^l - w_{\mathrm{GD}}^l \|_2 \le L_f^{-1} (1 + nL_f)^l \epsilon$$

1771 where L_f denotes the Lipschitz constant of $\mathcal{L}_N(w)$ within \mathcal{W} .

ŀ

1773 First we introduce a helper lemma.

Lemma 14 (Error for Approximating GD, Lemma G.1 of (Bai et al., 2023)). Let $\mathcal{W} \subset \mathbb{R}^d$ is a convex bounded domain and $\operatorname{Proj}_{\mathcal{W}}$ projects all vectors into \mathcal{W} . Suppose $f: \mathcal{W} \to R$ and ∇f is L_f -Lipschitz on \mathcal{W} . Fix any $\epsilon > 0$, let sequences $\{\overline{w}^l\}_{l\geq 0} \in \mathbb{R}^d$ and $\{w_{\mathrm{GD}}^l\}_{l\geq 0} \in \mathbb{R}^d$ are given by $\overline{w}^0 = w_{\mathrm{GD}}^0 = \mathbf{0}$, then for all $l \geq 0$,

1779
1780
$$\overline{w}^{l} = \operatorname{Proj}_{\mathcal{W}}(\overline{w}^{l-1} - \eta \nabla \mathcal{L}_{n}(\overline{w}^{l-1}) + \epsilon^{l-1}), \quad \|\epsilon^{l-1}\|_{2} \le \eta \epsilon,$$

1781
$$w_{\rm GD}^l = \operatorname{Proj}_{\mathcal{W}}(w_{\rm GD}^{l-1} - \eta \nabla \mathcal{L}_n(w_{\rm GD}^{l-1}))$$

To show the convergence, we define the gradient mapping at w with step size η as,

$$G_{\mathcal{W},\eta}^f := \frac{w - \operatorname{Proj}_{\mathcal{W}}(w - \eta \nabla \mathcal{L}_n(w))}{n}.$$

1787 Then if $\eta \leq L_f$, for all $L \geq 1$, convergence holds

$$\min_{l \in [L-1]} \|\mathbf{G}_{\mathcal{W},\eta}^{f}(\overline{w}^{l})\|_{2}^{2} \leq \frac{1}{L} \sum_{l=1}^{L-1} \|\mathbf{G}_{\mathcal{W},\eta}^{f}(\overline{w}^{l})\|_{2}^{2} \leq \frac{8(f(\mathbf{0}) - \inf_{w \in \mathcal{W}} f(w))}{\eta L} + 10\epsilon^{2}.$$

1792 Moreover, for any $l \ge 0$, the error accumulation is

$$\|\overline{w}^l - w_{\mathrm{GD}}^l\|_2 \le L_f^{-1} (1 + nL_f)^l \epsilon$$

Lemma 14 shows Theorem 1 leads to exponential error accumulation in the general case. Moreover,
 Lemma 14 also provides convergence of approximating GD. Then we proof Corollary 1.1.

1799 *Proof.* For any small ϵ , by Theorem 1, the neural network NN_{θ} implements each gradient descent step with error bounded by ϵ . Then we simply apply Lemma 14 to complete the proof.

¹⁸³⁶ E EXTENSION: DIFFERENT INPUT AND OUTPUT DIMENSIONS

In this section, we explore the ICGD on N-layer neural networks under the setting where the dimensions of input x_i and label y_i can be different. Specifically, we consider our prompt datasets $\{(x_i, y_i)\}_{i \in [n]}$ where $x_i \in \mathbb{R}^{d_x}$ and $y_i \in \mathbb{R}^{d_y}$. We start with our new N-layer neural network.

Definition 10 (*N*-Layer Neural Network). An *N*-Layer Neural Network comprises N - 1 hidden layers and 1 output layer, all constructed similarly. Let $r : \mathbb{R} \to \mathbb{R}$ be the activation function. For the hidden layers: for any $i \in [n + 1], j \in [N - 1]$, and $k \in [K]$, the output for the first j layers w.r.t. input $x_i \in \mathbb{R}^d$, denoted by $\operatorname{pred}_h(x_i; j) \in \mathbb{R}^K$, is defined as recursive form:

$$\operatorname{pred}_h(x_i;1)[k] := r(v_{1_k}^\top x_i), \quad \text{and} \quad \operatorname{pred}_h(x_i;j)[k] := r(v_{j_k}^\top \operatorname{pred}_h(x_i;j-1)),$$

1849 where $v_{1_k} \in \mathbb{R}^d$ and $v_{j_k} \in \mathbb{R}^K$ for $j \in \{2, ..., N-1\}$ are the k-th parameter vectors in the first 1850 layer and the j-th layer, respectively. For the output layer (N-th layer), the output for the first N 1851 layers (i.e the entire neural network) w.r.t. input $x_i \in \mathbb{R}^{d_x}$, denoted by $\operatorname{pred}_o(x_i; w, N) \in \mathbb{R}^{d_y}$, is 1852 defined for any $k \in [d_y]$ as follows:

$$\operatorname{pred}_o(x_i; w, N)[k] := r(v_{N_k}^{\perp} \operatorname{pred}_h(x_i; N-1))$$

where $v_{N_k} \in \mathbb{R}^K$ are the k-th parameter vectors in the N-th layer and $w \in \mathbb{R}^{(d_x+d_y)K+(N-2)K^2}$ denotes the vector containing all parameters in the neural network,

$$w := \left[v_{1_1}^{\top}, \dots, v_{1_K}^{\top}, \dots, v_{j_k}^{\top}, \dots, v_{N-1_1}^{\top}, \dots, v_{N-1_K}^{\top}, v_{N_1}^{\top}, \dots, v_{N_{d_y}}^{\top} \right]^{\top}.$$

1861 Notice that our new *N*-layer neural network only modify the output layer compared to Definition 1. 1862 Intuitively, this results in minimal change in output, which allows our framework in Section 3.3 1863 to function across varying input/output dimensions. Theoretically, we derive the explicit form of 1864 gradient $\nabla \mathcal{L}_n(w)$.

Lemma 15 (Decomposition of One Gradient Descent Step). Fix any $B_v, \eta > 0$. Suppose the empirical loss function $\mathcal{L}_n(w)$ on n data points $\{(x_i, y_i)\}_{i \in [n]}$ is defined as

$$\mathcal{L}_n(w) \coloneqq \frac{1}{2n} \sum_{i=1}^n \ell(f(w, x_i), y_i), \quad \text{where } \ell : \mathbb{R}^{d_y} \times \mathbb{R}^{d_y} \to \mathbb{R} \text{ is a loss function,}$$

where $f(w, x_i), y_i$ is the output of *N*-layer neural networks (Definition 10) with modified output layer. Suppose closed domain \mathcal{W} and projection function $\operatorname{Proj}_{\mathcal{W}}(w)$ follows (3.4). Let $A_i(j), r'_i(j), R_i(j), V_j$ be as defined in Definition 2 (with modified dimensions), then the explicit form of gradient $\nabla \mathcal{L}_n(w)$ becomes

$$\nabla \mathcal{L}_n(w) = \frac{1}{2n} \sum_{i=1}^n \begin{bmatrix} A_i(1) \\ \vdots \\ A_i(N) \end{bmatrix}$$

where $A_i(j)$ denote the derivative of $\ell(p_i(N), y_i)$ with respect to the parameters in the j-th layer,

$$A_i(j) = \begin{cases} (R_i(N-1) \cdot V_N \cdot \ldots \cdot R_i(j-1) \cdot \left[\mathbf{I}_{K \times K} \otimes p_i(j-1)^\top\right])^\top \cdot \left(\frac{\partial \ell(p_i(N), y_i)}{\partial p_i(N)}\right)^\top, & j \neq N\\ (R_i(N-1) \cdot \left[\mathbf{I}_{d_y \times d_y} \otimes p_i(N-1)^\top\right])^\top \cdot \left(\frac{\partial \ell(p_i(N), y_i)}{\partial p_i(N)}\right)^\top, & j = N. \end{cases}$$

1881

1847 1848

1859 1860

1865

1868

1870 1871

Proof. Simply follow the proof of Lemma 1. We show the different terms compared to Definition 2:

• Let $D_j \in \mathbb{R}$ denote the total number of parameters in the first *j* layers.

$$D_j = \begin{cases} 0, & j = 0\\ d_x K, & j = 1\\ (j-1)K^2 + d_x K, & 2 \le j \le N-1\\ (N-2)K^2 + (d_x + d_y)K, & j = N, \end{cases}$$

• The intermediate term $R_i(N-1)$,

1890

1892

1894 1895

1896 1897

1903 1904 1905

1907

1910

1919

1920

1926 1927 1928

1929 1930 1931

$$R_i(N-1) = \operatorname{diag}\{r'(v_{j+1_1}^\top p_i(j)), \dots, r'(v_{j+1_{d_y}}^\top p_i(j))\} \in \mathbb{R}^{d_y \times d_y}.$$

• The parameters matrices of the first and the last layers:

$$V_j := \begin{cases} \begin{bmatrix} v_{1_1}, \dots, v_{1_K} \end{bmatrix}^\top \in \mathbb{R}^{K \times d_x}, & j = 1 \\ \begin{bmatrix} v_{N_1}, \dots, v_{N_{d_y}} \end{bmatrix}^\top \in \mathbb{R}^{d_y \times K}, & j = N. \end{cases}$$

1909 Thus we complete the proof.

1911 Lemma 15 shows that the explicit form of gradient $\nabla \mathcal{L}_n(w)$ holds the same structure as Lemma 1. 1912 Therefore, it is simple to follow our framework in Section 3.3 to approximate $\nabla \mathcal{L}_n(w)$ term by term. 1913 Finally, we introduce the generalized version of main result Theorem 1.

Theorem 5 (In-Context Gradient Descent on N-layer NNs). Fix any B_v , η , $\epsilon > 0$, $L \ge 1$. For any input sequences takes from (2.1), where $\{(x_i, y_i)\}_{i \in [n]}$ and $x_i \in \mathbb{R}^{d_x}$ and $y_i \in \mathbb{R}^{d_y}$, their exist upper bounds B_x , B_y such that for any $i \in [n]$, $||y_i||_2 \le B_y$, $||x_i||_2 \le B_x$. Assume functions r(t), r'(t)and u(t, y)[k] are L_r , $L_{r'}$, L_l -Lipschitz continuous. Suppose \mathcal{W} is a closed domain such that for any $j \in [N-1]$ and $k \in [K]$,

$$\mathcal{W} \subset \left\{ w = [v_{j_k}] \in \mathbb{R}^{D_N} : \|v_{j_k}\|_2 \le B_v \right\},$$

and $\operatorname{Proj}_{\mathcal{W}}$ project w into bounded domain \mathcal{W} . Assume $\operatorname{Proj}_{\mathcal{W}} = \operatorname{MLP}_{\theta}$ for some MLP layer with hidden dimension D_w parameters $\|\theta\| \leq C_w$. If functions r(t), r'(t) and u(t, y)[k] are C^4 smoothness, then for any $\epsilon > 0$, there exists a transformer model $\operatorname{NN}_{\theta}$ with (2N + 4)L hidden layers consists of L neural network blocks $\operatorname{TF}_{\theta}^{N+2} \circ \operatorname{EWML}_{\theta}^N \circ \operatorname{TF}_{\theta}^2$,

$$\mathrm{NN}_{\theta} := \mathrm{TF}_{\theta}^{N+2} \circ \mathrm{EWML}_{\theta}^{N} \circ \mathrm{TF}_{\theta}^{2} \circ \ldots \circ \mathrm{TF}_{\theta}^{N+2} \circ \mathrm{EWML}_{\theta}^{N} \circ \mathrm{TF}_{\theta}^{2},$$

such that the heads number M^l , parameter dimensions D^l , and the parameter norms B_{θ^l} suffice

$$\max_{l \in [(2N+4)L]} M^{l} \le \widetilde{O}(\epsilon^{-2}), \quad \max_{l \in [(2N+4)L]} D^{l} \le O(K^{2}N) + D_{w}, \quad \max_{l \in [(2N+4)L]} B_{\theta^{l}} \le O(\eta) + C_{w} + 1,$$

where $\widetilde{O}(\cdot)$ hides the constants that depend on d, K, N, the radius parameters B_x, B_y, B_v and the smoothness of r and ℓ . And this neural network such that for any input sequences $H^{(0)}$, take from (2.1), $NN_{\theta}(H^{(0)})$ implements L steps in-context gradient descent on risk $\mathcal{L}_n(w)$ follows Lemma 15: For every $l \in [L]$, the (2N + 4)*l*-th layer outputs $h_i^{((2N+4)l)} = [x_i; y_i; \overline{w}^{(l)}; \mathbf{0}; 1; t_i]$ for every $i \in [n + 1]$, and approximation gradients $\overline{w}^{(l)}$ such that

$$\overline{w}^{(l)} = \operatorname{Proj}_{\mathcal{W}}(\overline{w}^{(l-1)} - \eta \nabla \mathcal{L}_n(\overline{w}^{(l-1)}) + \epsilon^{(l-1)}), \quad \overline{w}^{(0)} = \mathbf{0},$$

1940 where $\|\epsilon^{(l-1)}\|_2 \le \eta \epsilon$ is an error term. 1941

1942

1937 1938 1939

¹⁹⁴⁴ F EXTENSION: SOFTMAX TRANSFORMER

1950

1952

1956 1957

1984

1986

In this part, we demonstrate the existence of pretrained Softmax transformers capable of implementing ICGD on an *N*-layer neural network. First, we introduce our main technique: the universal approximation property of softmax transformers in Appendix F.1. Then, we prove the existence of pretrained softmax transformers that implement ICGD on *N*-layer neural networks in Appendix F.2.

1951 F.1 UNIVERSAL APPROXIMATION OF SOFTMAX TRANSFORMER

Softmax-Attention Layer. We replace modified normalized ReLU activation σ/n in ReLU attention layer (Definition 7) by standard softmax. Thus, for any input sequence $H \in \mathbb{R}^{D \times n}$, a single head attention layer outputs

Attn
$$(H) = H + W^{(O)}(VH)$$
 Softmax $[(KH)^{\top}(QH)],$ (F.1)

where $W^{(O)}, Q, K, V \in \mathbb{R}^{D \times D} \in \mathbb{R}^{d \times d}$ are the weight matrices. Then we introduce the softmax transformer block, which consists of two feed-forward neural network layers and a single-head self-attention layer with the softmax function.

Definition 11 (Transformer Block $\mathcal{T}_{Softmax}$). For any input sequences $H \in \mathbb{R}^{D \times n}$, let $FF(H) := H + W_2 \cdot \operatorname{ReLU}(W_1H + b_1\mathbb{1}_L^T) + b_2\mathbb{1}_L^T$ be the Feed-Forward layer, where d' is hidden dimensions, $W_1 \in \mathbb{R}^{d' \times D}, W_2 \in \mathbb{R}^{D \times d'}, b_1 \in \mathbb{R}^l$, and $b_2 \in \mathbb{R}^d$. We configure a transformer block with Softmax-attention layer as $\mathcal{T}_{Softmax} := \{FF \circ Attn \circ FF : \mathbb{R}^{d \times L} \to \mathbb{R}^{d \times L}\}.$

1966 1967 Universal Approximation of Softmax-Transformer. We show the universal approximation theorem for Transformer blocks (Definition 11). Specifically, Transformer blocks $\mathcal{T}_{Softmax}$ are universal approximators for continuous permutation equivariant functions on bounded domain.

Lemma 16 (Universal Approximation of $\mathcal{T}_{Softmax}$). Let $f(\cdot) := \mathbb{R}^{d \times n} \to \mathbb{R}^{d \times n}$ be any *L*-Lipschitz permutation equivariant function supported on $[0, B_x]^{d \times n}$. We denote the discrete input domain of $[0, B_x]^{d \times n}$ by a grid \mathbb{G}_D with granularity $D \in \mathbb{N}$ defined as $\mathbb{G}_D = \{B_x/D, 2B_x/D, \ldots, B_x\}^{d \times n} \subset \mathbb{R}^{d \times n}$. For any $\kappa > 0$, there exists a transformer network $f_{Softmax} \in \mathcal{T}_{Softmax}$, such that for any $Z \in [0, B_x]^{d \times n}$, it approximate f(Z) as: $||f_{Softmax}(Z) - f(Z)||_2 \leq \frac{L(8\sqrt{dn} + \sqrt{n})B_x}{2D} = \kappa$.

1976
1977Proof Sketch. First, we use a piece-wise constant function to approximate f and derive an upper
bound based on its L-Lipschitz property. Next, we demonstrate how the feed-forward neural network
 $\mathcal{F}_1^{(FF)}$ quantizes the continuous input domain into the discrete domain \mathbb{G}_D through a multiple-step
function, using ReLU functions to create a piece-wise linear approximation. Then, we apply the
self-attention layer $\mathcal{F}_1^{(SA)}$ on $\mathcal{F}_1^{(FF)}$, establishing a bounded output region for $\mathcal{F}_S^{(SA)} \circ \mathcal{F}_1^{(FF)}$.1980
1981Finally, we employ a second feed-forward network $\mathcal{F}_2^{(FF)}$ to predict $f_{\text{Softmax}}(Z)$ and assess the
approximation error relative to the actual output f(Z). See Appendix F.4 for a detailed proof.

F.2 IN-CONTEXT GRADIENT DESCENT WITH SOFTMAX TRANSFORMER

In-Context Gradient Descent with Softmax Transformer. By applying universal approximation theory (Lemma 16), we now illustrate how to use Transformer block $\mathcal{T}_{Softmax}$ (Definition 11) and MLP layers (Definition 8) to implement ICGD on general risk function $\mathcal{L}_n(w)$.

Theorem 6 (In-Context Gradient Descent on General Risk Function). Fix any B_w , η , $\epsilon > 0$, $L \ge 1$. For any input sequences takes from (2.1), their exist upper bounds B_x , B_y such that for any $i \in [n]$, $\|y_i\|_{\max} \le B_y$, $\|x_i\|_{\max} \le B_x$. Suppose W is a closed domain such that $\|w\|_{\max} \le B_w$ and Proj_W project w into bounded domain W. Assume $\operatorname{Proj}_W = \operatorname{MLP}_{\theta}$ for some MLP layer. Define $l(w, x_i, y_i)$ as a loss function with L-Lipschitz gradient. Let $\mathcal{L}_n(w) = \frac{1}{n} \sum_{i=1}^n \ell(w, x_i, y_i)$ denote the empirical loss function, then there exists a transformer $\operatorname{NN}_{\theta}$, such that for any input sequences $H^{(0)}$, take from (2.1), $\operatorname{NN}_{\theta}(H^{(0)})$ implements L steps in-context gradient descent on $\mathcal{L}_n(w)$: For every $l \in [L]$, the 4*l*-th layer outputs $h_i^{(4l)} = [x_i; y_i; \overline{w}^{(l)}; \mathbf{0}; 1; t_i]$ for every $i \in [n + 1]$, and approximation gradients

 $\overline{w}^{(l)}$ such that

$$\overline{w}^{(l)} = \operatorname{Proj}_{\mathcal{W}}(\overline{w}^{(l-1)} - \eta \nabla \mathcal{L}_n(\overline{w}^{(l-1)}) + \epsilon^{(l-1)}), \quad \overline{w}^{(0)} = \mathbf{0},$$

where $\|\epsilon^{(l-1)}\|_2 \leq \eta \epsilon$ is an error term.

Proof Sketch. By our assumption $\operatorname{Proj}_{\mathcal{W}} = \operatorname{MLP}_{\theta}$, we only need to find a transformer to implement gradient descent $w^+ := w - \eta \nabla \mathcal{L}_n(w)$. For ant input takes from (F.2), let function $f : \mathbb{R}^{D \times n} \to$ $\mathbb{R}^{D \times n}$ maps w into $w - \eta \nabla \mathcal{L}_n(w)$ and preserve other elements. By Lemma 7, their exist a transformer block f_{Softmax} capable of approximating f with any desired small error. Therefore, $f_{\text{Softmax}} \circ \text{MLP}$ suffices our requirements. Please see Appendix F.3 for a detailed proof.

F.3 PROOF OF THEOREM 6

Proof of Theorem 6. We only need to construct a 4 layers transformer capable of implementing single step gradient descent. With out loss of generality, we assume $w \in \mathbb{R}^{D_w}$. Recall that the input sequences $H \in \mathbb{R}^{D \times n}$ takes form

$$H \coloneqq \begin{bmatrix} x_1 & x_2 & \cdots & x_n & x_{n+1} \\ y_1 & y_2 & \cdots & y_n & 0 \\ q_1 & q_2 & \cdots & q_n & q_{n+1} \end{bmatrix} \in \mathbb{R}^{D \times (n+1)}, \quad q_i \coloneqq \begin{bmatrix} w \\ 0 \\ 1 \\ t_i \end{bmatrix} \in \mathbb{R}^{D - (d+1)}.$$
(F.2)

Let function $f : \mathbb{R}^{D \times n} \to \mathbb{R}^{D \times n}$ output

$$f(H) = \begin{bmatrix} x_1 & x_2 & \cdots & x_n & x_{n+1} \\ y_1 & y_2 & \cdots & y_n & 0 \\ q_1 & q_2 & \cdots & q_n & q_{n+1} \end{bmatrix}, \quad q_i \coloneqq \begin{bmatrix} w - \eta \nabla \mathcal{L}_n(w) \\ 0 \\ 1 \\ t_i \end{bmatrix} \in \mathbb{R}^{D-(d+1)}.$$

By Lemma 16, for any $\kappa > 0$, there exists a transformer network $f_{
m Softmax} \in \mathcal{T}_{
m Softmax}$, such that for any input $H \in [-B, B]^{d \times L}$, we have $\|f_{\text{Softmax}}(H) - f(H)\|_2 \leq \kappa$. Therefore, by the equivalence of matrix norms, $\|f_{\text{Softmax}}(H) - f(H)\|_{\text{max}} \leq \kappa$ holds without loss of generality. Above $B := \max\{B_x, B_y, B_w, 1\}$ denotes the upper bound for every elements in H. Thus, we obtain \overline{w} from the identical position of w in $f_{\text{Softmax}}(H)$. Suppose we choose $\kappa = \frac{\epsilon}{\sqrt{D_w}}$, then it holds

$$\|\overline{w} - (w - \eta \nabla \mathcal{L}_n(w)\|_2 \le \sqrt{D_w} \|\overline{w} - (w - \eta \nabla \mathcal{L}_n(w)\|_{\max})$$

$$\leq \|f_{\text{Softmax}} - f(H)\|_{\text{max}}$$

Finally, by our assumption, there exists an MLP layer such that for any $i \in [n+1]$, it maps

 $< \epsilon$.

$$[x_i; y_i; w - \eta \nabla \mathcal{L}_n(w); \mathbf{0}; 1; t_i] \xrightarrow{\mathrm{MLP}} [x_i; y_i; \mathrm{Proj}_{\mathcal{W}}(w - \eta \nabla_w \mathcal{L}_n(w)); \mathbf{0}; 1; t_i].$$

Therefore, a four-layer transformer $f_{\text{Softmax}} \circ \text{MLP}$ is capable of implementing one-step gradient descent through ICL. As a direct corollary, there exist a 4L-layer transformer consists of L identical blocks $f_{\text{Softmax}} \circ \text{MLP}$ to approximate L steps gradient descent algorithm. Each block approximates a one-step gradient descent algorithm on general risk function $\mathcal{L}_n(w)$.

F.4 PROOF OF LEMMA 16

In this section, we introduce a helper lemma Lemma 17 to prove Lemma 16. At the beginning, we assume all input sequences are separated by a certain distance.

2052 **Definition 12** (Token-wise Separateness, Definition 1 of (Kajitsuka and Sato, 2024)). Let $N \ge 1$ 2053 and $Z^{(1)}, \ldots, Z^{(N)} \in \mathbb{R}^{d \times n}$ be input sequences. Then, $Z^{(1)}, \ldots, Z^{(N)}$ are called token-wise 2054 $(r_{\min}, r_{\max}, \delta)$ -separated if the following three conditions hold. 2055 • For any $i \in [N]$ and $k \in [n]$, $\left\| Z_{:,k}^{(i)} \right\|_2 > r_{\min}$ holds. 2056 2057 • For any $i \in [N]$ and $k \in [n]$, $\left\| Z_{:,k}^{(i)} \right\|_{2} < r_{\max}$ holds. 2058 • For any $i, j \in [N]$ and $k, l \in [n]$ with $Z_{:,k}^{(i)} \neq Z_{:,l}^{(j)}, \left\| Z_{:,k}^{(i)} - Z_{:,l}^{(j)} \right\|_2 > \delta$ holds. 2059 2060 Note that we refer to $Z^{(1)}, \ldots, Z^{(N)}$ as token-wise (r_{\max}, ϵ) -separated instead if the sequences 2061 satisfy the last two conditions. 2062 Then we introduce the definition of contextual mapping. Intuitively, a contextual mapping can provide 2063 every input sequence with a unique id, which enables us to construct approximation for labels. 2064 2065 Definition 13 (Contextual mapping, Definition 2 of (Kajitsuka and Sato, 2024)). Let input sequences 2066 $Z^{(1)}, \ldots, Z^{(N)} \in \mathbb{R}^{d \times n}$. Then, a map $q: \mathbb{R}^{d \times n} \to \mathbb{R}^{d \times n}$ is called an (r, δ) -contextual mapping if 2067 the following two conditions hold: 2068 • For any $i \in [N]$ and $k \in [n]$, $\left\|q\left(Z^{(i)}\right)_{k}\right\|_{0} < r$ holds. 2069 • For any $i, j \in [N]$ and $k, l \in [n]$, if $Z_{:,k}^{(i)} \neq Z_{:,l}^{(j)}$, then $\left\|q\left(Z^{(i)}\right)_{:,k} - q\left(Z^{(j)}\right)_{:,l}\right\|_2 > \delta$ holds. 2070 2071 In particular, $q(Z^{(i)})$ for $i \in [N]$ is called a context id of $Z^{(i)}$. 2072 2073 Next, we show that a softmax-based 1-layer attention block with low-rank weight matrices is a 2074 contextual mapping for almost all input sequences. 2075

Lemma 17 (Softmax attention is contextual mapping, Theorem 2 of (Kajitsuka and Sato, 2024)). Let $Z^{(1)}, \ldots, Z^{(N)} \in \mathbb{R}^{d \times n}$ be input sequences with no duplicate word token in each sequence, that is,

$$Z_{:,k}^{(i)} \neq Z_{:,l}^{(i)}$$

for any $i \in [N]$ and $k, l \in [n]$. Also assume that $Z^{(1)}, \ldots, Z^{(N)}$ are token-wise $(r_{\min}, r_{\max}, \epsilon)$ separated. Then, there exist weight matrices $W^{(O)} \in \mathbb{R}^{d \times s}$ and $V, K, Q \in \mathbb{R}^{s \times d}$ such that the ranks of V, K and Q are all 1, and 1-layer single head attention with softmax, i.e., $\mathcal{F}_{S}^{(SA)}$ with h = 1 is an (r, δ) -contextual mapping for the input sequences $Z^{(1)}, \ldots, Z^{(N)} \in \mathbb{R}^{d \times n}$ with r and δ defined by

$$\begin{aligned} r &= r_{\max} + \frac{1}{4} \\ \delta &= \frac{2(\log n)^2 \epsilon^2 r_{\min}}{r_{\max}^2 (|\mathcal{V}| + 1)^4 (2\log n + 3)\pi d} \exp\left(-(|\mathcal{V}| + 1)^4 \frac{(2\log n + 3)\pi dr_{\max}^2}{4\epsilon r_{\min}}\right) \end{aligned}$$

2090 Applying Lemma 17, we extends Proposition 1 of (Kajitsuka and Sato, 2024) to our Lemma 16.¹ We provide explicit upper bound of error $||f_{Softmax}(Z) - f(Z)||_2$ and analysis with function f of a 2092 broader supported domain. 2093

Lemma 18 (Lemma 16 Restated: Universal Approximation of $\mathcal{T}_{\text{Softmax}}$). Let $f(\cdot) \coloneqq \mathbb{R}^{d \times n} \to$ $\mathbb{R}^{d \times n}$ be any *L*-Lipschitz permutation equivalent function supported on $[0, B_x]^{d \times n}$. We denote the discrete input domain of $[0, B_x]^{d \times n}$ by a grid \mathbb{G}_D with granularity $D \in \mathbb{N}$ defined as $\mathbb{G}_D = \{B_x/D, 2B_x/D, \ldots, B_x\}^{d \times n} \subset \mathbb{R}^{d \times n}$. For any $\kappa > 0$, there exists a transformer network $f_{\text{Softmax}} \in \mathcal{T}_{\text{Softmax}}$ (Definition 11), such that for any $Z \in [0, B_x]^{d \times n}$, it approximate f(Z) as: $\|f_{\text{Softmax}}(Z) - f(Z)\|_2 \le \frac{L(8\sqrt{dn} + \sqrt{n})B_x}{2D} = \kappa.$

2099 2100 2101

2098

2081

2082

2083

2084 2085 2086

2089

2091

2094

2095 2096 2097

2102 *Proof.* We begin our 3-step proof.

²¹⁰³

²¹⁰⁴ ¹This extension builds on the results of (Hu et al., 2024a), which extend the rank-1 requirement to any rank 2105 for attention weights. Additionally, Hu et al. (2024b) apply similar techniques to analyze the statistical rates of diffusion transformers (DiTs).

Approximation of f **by piece-wise constant function.** Since f is a continuous function on a compact set, f has maximum and minimum values on the domain. By scaling with $\mathcal{F}_1^{(FF)}$ and $\mathcal{F}_2^{(FF)}$, f is assumed to be normalized: for any $Z \in \mathbb{R}^{d \times n} \setminus [0, B_x]^{d \times n}$

$$f(Z) = 0$$

and for any $Z \in [0, B_x]^{d \times n}$

$$-B_y \le f(Z) \le B_y.$$

2116 Let $D \in \mathbb{N}$ be the granularity of a grid \mathbb{G}_D :

$$\mathbb{G}_D = \{\frac{B_x}{D}, \frac{2B_x}{D}, \dots, B_x\}^{d \times n} \subset \mathbb{R}^{d \times n}$$

where each coordinate only take discrete value B_x/D , $2B_x/D$, ..., B_x . Now with a continuous input *Z*, we approximate *f* by using a piece-wise constant function \overline{f} evaluating on the nearest grid point *L* of *Z* in the following way:

$$\bar{f}(Z) = \sum_{L \in \mathbb{G}_D} f(L) \, \mathbf{1}_{Z \in L + [-B_x/D, 0)^{d \times n}}.$$
(F.3)

Additionally if $Z \in L + [-1/D, 0)^{d \times n}$, denote it as Q(Z) = L.

Now we bound the piece-wise constant approximation error $||f - \overline{f}||$ as follows.

2131 Define set $P_D = \{L + [-B_x/D, 0)^{d \times n} | L \in \mathbb{G}_D\}$. It is a set of regions of size $(\frac{B_x}{D})^{d \times n}$, whose 2132 vertexes are the points in \mathbb{G}_D .

For any subset $U \in P_D$, the maximal difference of f and \overline{f} in this region is:

$$\max_{Z \in U} \|f(Z) - \bar{f}(Z)\|_{2} = \max_{Z \in U} \|f(Z) - f(Q(Z))\|_{2}$$

$$\leq \max_{Z, Z' \in U} \|f(Z) - f(Z')\|_{2}$$

$$\leq L \cdot \max_{Z, Z' \in U} \|Z - Z'\|_{2} \qquad (By \ f \ is \ a \ L-Lipschitz \ function)$$

$$= L \cdot \sqrt{dn \cdot (\frac{B_{x}}{D})^{2}} \quad (Z, Z' \ are \ in \ the \ same \ \frac{B_{x}}{D} \text{-wide} \ (d \cdot n) \text{-dimension} \ U.)$$

$$= \frac{L\sqrt{dn}B_{x}}{D}. \qquad (F.4)$$

Quantization of input using $\mathcal{F}_1^{(FF)}$. In the second step, we use $\mathcal{F}_1^{(FF)}$ to quantize the continuous input domain into \mathbb{G}_D . This process is achieved by a multiple-step function, and we use ReLU functions to approximate this multiple-step functions. This ReLU function can be easily implemented by a one-layer feed-forward network.

First for any small $\delta > 0$ and $z \in \mathbb{R}$, we construct a δ -approximated step function using ReLU functions:

$$\frac{\sigma_R \left[\frac{z}{\delta}\right] - \sigma_R \left[\frac{z}{\delta} - B_x\right]}{D} = \begin{cases} 0 & z < 0\\ \frac{z}{\delta D} & 0 \le z < \delta B_x \\ \frac{B_x}{D} & \delta B_x \le z \end{cases}$$
(F.5)

where a one-hidden-layer feed-forward neural network is able to implement this. By shifting (F.5) by B_x , for any $t \in [D-1]$, we have:

$$\frac{\sigma_R \left[\frac{z}{\delta} - \frac{tB_x}{\delta D}\right] - \sigma_R \left[\frac{z}{\delta} - B_x - \frac{tB_x}{\delta D}\right]}{D} = \begin{cases} 0 & z < \frac{tB_x}{D} \\ \frac{z}{\delta D} & \frac{tB_x}{D} \le z < \delta B_x + \frac{tB_x}{D} \end{cases},$$
(F.6)

when δ is small the above function approximates to a step function:

$$\operatorname{quant}_D^{(t)}(z) = \begin{cases} 0 & z \leq \frac{tB_x}{D} \\ \frac{B_x}{D} & \frac{tB_x}{D} \leq z \end{cases}.$$

2173 By adding up (F.6) at every $t \in [D-1]$, we have an approximated multiple-step function

2174
2175
2176
2177
2178
2179
2180
2181
2182
2184
2184
2184

$$\sum_{t=0}^{D-1} \frac{\sigma_R \left[\frac{z}{\delta} - \frac{tB_x}{\delta D}\right] - \sigma_R \left[\frac{z}{\delta} - B_x - \frac{tB_x}{\delta D}\right]}{D} \quad (F.7)$$
(F.7)
(when δ is small.)
(when δ is small.)
(F.7)
(when δ is small.)
(F.8)

$$2184 = \begin{cases} 2 & 2 \\ 2185 & 2 \\ 2185 & 2185 \\$$

 $\begin{array}{c} 2105\\ 2186 \end{array} \qquad \qquad \left(\begin{array}{c} B_x \quad B_x - \frac{B_x}{D} \le z \end{array} \right)$

2188 Note that the error of approximation at z here estimated as:

$$\left|\sum_{t=0}^{D-1} \frac{\sigma_R \left[\frac{z}{\delta} - \frac{tB_x}{\delta D}\right] - \sigma_R \left[\frac{z}{\delta} - B_x - \frac{tB_x}{\delta D}\right]}{D} - \text{quant}_D(z)\right| \le \frac{B_x}{D},\tag{F.9}$$

2193 and for matrix $Z \in \mathbb{R}^{d \times n}$:

$$\begin{split} \|\sum_{t=0}^{D-1} \frac{\sigma_R \left[\frac{Z}{\delta} - \frac{Q(Z)}{\delta D} \right] - \sigma_R \left[\frac{Z}{\delta} - B_x E - \frac{Q(Z)}{\delta D} \right]}{D} - \operatorname{quant}_D(Z) \|_2 \\ &\leq \sqrt{d \times n \times \left(\frac{B_x}{D} \right)^2} \\ &= \frac{B_x \sqrt{dn}}{D}. \end{split}$$

Subtract the last step function from (F.7) we get the desired result:

$$\sum_{t=0}^{D-1} \frac{\sigma_R \left[\frac{z}{\delta} - \frac{tB_x}{\delta D} \right] - \sigma_R \left[\frac{z}{\delta} - B_x - \frac{tB_x}{\delta D} \right]}{D} - \left(\sigma_R \left[\frac{z}{\delta} - \frac{B_x}{\delta} \right] - \sigma_R \left[\frac{z}{\delta} - 1 - \frac{B_x}{\delta} \right] \right).$$
(F.10)

This equation approximate the quantization of input domain $[0, B_x]$ into $\{B_x/D, \ldots, B_x\}$ and making $\mathbb{R} \setminus [0, B_x]$ to 0. In addition to the quantization of input domain $[0, B_x]$, we add a penalty term for input out of $[0, B_x]$ in the following way:

2216 2217

2218

 $-B_x \sigma_R \left[\frac{(z - B_x)}{\delta} \right] + B_x \sigma_R \left[\frac{(z - B_x)}{\delta} - 1 \right] - B_x \sigma_R \left[\frac{-z}{\delta} \right] + B_x \sigma_R \left[\frac{-z}{\delta} - 1 \right]$ $\approx \text{penalty}(z) = \begin{cases} -B_x & z \le 0\\ 0 & 0 < z \le B_x \\ -B_r & B_r < z \end{cases}$

2219 2220

2222 2223

2224

2244 2245

Both (F.10) and (F.11) can be realized by the one-layer feed-forward neural network. Also, it is straightforward to show that generate both of them to input $Z \in \mathbb{R}^{d \times n}$.

2225 2226 Combining both components together, the first feed-forward neural network layer $\mathcal{F}_1^{(FF)}$ approxi-2227 mates the following function $\overline{\mathcal{F}}_1^{(FF)}(Z)$:

$$\mathcal{F}_1^{(FF)} \approx \overline{\mathcal{F}}_1^{(FF)}(Z) = \operatorname{quant}_D^{d \times n}(Z) + \sum_{t=1}^d \sum_{k=1}^n \operatorname{penalty}(Z_{t,k}).$$
(F.12)

(F.11)

Note how we generalize penalty(·) to multi-dimensional occasions in the above equation. Whenever an input sequence Z has one entry $Z_{t,k}$ out of $[0, B_x]^{d \times n}$, we penalize the whole input sequence by adding a $-B_x$ to all entries. This makes all entries of this quantization lower bounded by $-dnB_x$.

(F.12) quantizes inputs in $[0, B_x]^{d \times n}$ with granularity D, while every element of the output is nonpositive for inputs outside $[0, B_x]^{d \times n}$. In particular, the norm of the output is upper-bounded when every entry in Z is out of $[0, B_x]$, this adds $-dnB_x$ penalties to all entries:

$$\max_{Z \in \mathbb{R}^{d \times n}} \left\| \mathcal{F}_1^{(FF)}(Z)_{:,k} \right\|_2 = \sqrt{d \cdot (-dnB_x)^2} \qquad (\text{One column is } d-\text{dimension.})$$
$$\leq dn \cdot \sqrt{d}B_x, \qquad (F.13)$$

for any $k \in [n]$.

Estimating the influence of self-attention $\mathcal{F}^{(SA)}$. Define $\widetilde{\mathbb{G}}_D \subset \mathbb{G}_D$ as:

$$\widehat{\mathbb{G}}_D = \{ L \in \mathbb{G}_D \mid \forall k, l \in [n], L_{:,k} \neq L_{:,l} \} .$$
(F.14)

It is a set of all the input sequences that don't have have identical tokens after quantization.

Within this set, the elements are at least $\frac{B_x}{D}$ separated by the quantization. Thus Lemma 17 allows us to construct a self-attention $\mathcal{F}(SA)$ to be a contextual mapping for such input sequences.

Since when D is sufficiently large, originally different tokens will still be different after quantization. In this context, we omit $\mathbb{G}_D/\widetilde{\mathbb{G}}_D$ for simplicity.

From the proof of Lemma 17 in (Kajitsuka and Sato, 2024), we follow their way to construct self-attention and have following equation:

$$\left\|\mathcal{F}_{S}^{(SA)}(Z)_{:,k} - Z_{:,k}\right\|_{2} < \frac{1}{4\sqrt{dD}} \max_{k' \in [n]} \|Z_{:,k'}\|_{2},\tag{F.15}$$

for any $k \in [n]$ and $Z \in \mathbb{R}^{d \times n}$.

2264 Combining this upper-bound with (F.13) we have

2265

2259 2260

2266

$$<rac{1}{4\sqrt{d}D} imes dn\sqrt{d}B_x$$
 (By (F.13))

$$=\frac{dnB_x}{4D}.$$
(F.16)

We show that if we take large enough D, every element of the output for $Z \in \mathbb{R}^{d \times n} \setminus [0, B_x]^{d \times n}$ is upper-bounded by

 $\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)} \left(Z \right)_{t,k} < \frac{B_{x}}{4D} \quad (\forall t \in [d], \ k \in [n]).$ (F.17)

2281 To show (F.17) holds, we consider the opposite occasion that there exists a $\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t_{0},k_{0}} \geq B_{x}/4D$. Then we divide the case into two sub cases:

1. The whole $\mathcal{F}_{1}^{(FF)}(Z)$ receives no less than 2 penalties. In this occasion, since every entry consists of two counterparts in (F.12): the quantization part quant $_{D}^{d \times n}(Z) \in [0, B_{x}]$ and aggregated with a penalty part $\sum_{t=1}^{d} \sum_{k=1}^{n} \text{penalty}(Z_{t,k}) \leq -2B_{x}$, for every entry we have $\mathcal{F}^{(FF)}(Z)_{t,k} \leq -B_{x}$. This yields that:

$$\begin{aligned} \|\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{:,k_{0}} - \mathcal{F}^{(FF)}(Z)_{:,k_{0}} \|_{2} &\geq \|\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t_{0},k_{0}} - \mathcal{F}^{(FF)}(Z)_{t_{0},k_{0}} \|_{2} \\ &\geq |\frac{B_{x}}{4D} - (-B_{x})| \\ &\geq \frac{dn}{4D} B_{x}, \end{aligned}$$
 (for a large enough D)

thus we derive a contradiction towards (F.16) from the assumption, proving it to be incorrect.

2. The whole $\mathcal{F}_1^{(FF)}(Z)$ receives only one penalty. In this case all entries in Z is penalized by $-B_x$ and satisfies:

$$\mathcal{F}_{1}^{(FF)}(Z)_{t,k} \in [-B_{x}, 0]^{d \times n}.$$
(F.18)

By (F.15), this further denotes:

$$\left\|\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{:,k} - \mathcal{F}_{1}^{(FF)}(Z)_{:,k}\right\|_{2} < \frac{1}{4\sqrt{d}D} \max_{k' \in [n]} \|\mathcal{F}_{1}^{(FF)}(Z)_{:,k'}\|_{2} \qquad (\text{By (F.15)})$$

$$\leq \frac{1}{4\sqrt{d}D}\sqrt{d\times B_x^2} \qquad (By (F.18))$$

$$=\frac{B_x}{4D}.$$
(F.19)

Yet by our assumption, there exists such an entry $\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}^{(FF)}(Z)_{t_{0},k_{0}} \geq B_{x}/4D$, which since $\mathcal{F}_{1}^{(FF)}(Z)_{t_{0},k_{0}} \leq 0$, yields:

$$\begin{split} \left\| \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{:,k_{0}} - \mathcal{F}_{1}^{(FF)}(Z)_{:,k_{0}} \right\|_{2} &\geq \left\| \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t_{0},k_{0}} - \mathcal{F}_{1}^{(FF)}(Z)_{t_{0},k_{0}} \right\|_{2} \\ &\geq \left| \frac{B_{x}}{4D} - 0 \right| \\ &= \frac{B_{x}}{4D} \end{split}$$

The final conclusion contradict the former result, suggesting the prerequisite to be fallacious.

Joining the incorrectness of the two sub-cases of the opposite occasion, we confirm the upper bound when input Z is outside $[0, B_x]^{d \times n}$ in (F.17).

For the input Z inside $[0, B_x]^{d \times n}$, we now show it is lower-bounded by

$$\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)} \left(Z \right)_{t,k} > \frac{3B_{x}}{4D} \quad (\forall t \in [d], \ k \in [n]).$$
(F.20)

By our construction, every entry Z in $[0, B_x]^{d \times n}$ satisfies:

$$\mathcal{F}_1^{(FF)}(Z)_{t,k} \in [\frac{B_x}{D}, B_x]. \tag{F.21}$$

By (F.15):

$$\begin{aligned} \left\| \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{:,k} - \mathcal{F}_{1}^{(FF)}(Z)_{:,k} \right\|_{2} \\ < \frac{1}{4\sqrt{dD}} \max_{k' \in [n]} \left\| \mathcal{F}_{1}^{(FF)}(Z)_{:k'} \right\|_{2} \qquad (By (F.15)) \\ \leq \frac{1}{4\sqrt{dD}} \sqrt{d \times B_{x}^{2}} \qquad (d-\text{dimension vector with each entry has maximum value } B_{x}.) \\ = \frac{B_{x}}{4D}. \end{aligned}$$
(F.22)

This yields:

$$\left|\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t,k} - \mathcal{F}_{1}^{(FF)}(Z)_{t,k}\right| \leq \left\|\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{:,k} - \mathcal{F}_{1}^{(FF)}(Z)_{:,k}\right\|_{2} < \frac{B_{x}}{4D}.$$
(F.23)

2352 Finally, we have:

$$\begin{aligned} \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t,k} \\ &> \mathcal{F}_{1}^{(FF)}(Z)_{t,k} - \left\| \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t,k} - \mathcal{F}_{1}^{(FF)}(Z)_{t,k} \right\|_{2} \\ &> \frac{B_{x}}{D} - \left\| \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t,k} - \mathcal{F}_{1}^{(FF)}(Z)_{t,k} \right\|_{2} \\ &> \frac{B_{x}}{D} - \frac{B_{x}}{4D} \\ &= \frac{3B_{x}}{4D}. \end{aligned}$$
(By (F.23)

Hence we finally finish the proof for the upper bound of $\mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)_{t,k}$ for Z outside $[0, B_x]$ in (F.17) and lower bound for Z inside $[0, B_x]$ in (F.20).

Approximation error Now, we can conclude our work by constructing the final feed-forward network $\mathcal{F}_2^{(FF)}$. It receives the output of the self-attention layer and maps the ones in $\widetilde{\mathbb{G}}_D \subset$ $(3B_x/4D,\infty)^{d\times n}$ to the corresponding value of the target function, and the rest in $(-\infty, B_x/4D)^{d\times n}$ to 0.

In order to adapt to the L_2 norm, we use a continuous and Lipschitz function to map the input Z to its targeted corresponding output f(Q(Z)).

According to piece-wise linear approximation, function $\mathcal{F}_2^{(FF)}$ exists such that for any input $L \in G_D$, it maps it to corresponding f(L), and for an arbitrary input Z, its output suffices:

$$\mathcal{F}_{2}^{(FF)}(Z) \in [\min_{\|L-Z\|_{\max} \le \frac{B_{x}}{2D}} f(L), \max_{\|L-Z\|_{\max} \le \frac{B_{x}}{2D}} f(L)].$$
(F.24)

Next we estimate the difference between $\mathcal{F}_2^{(FF)} \circ \mathcal{F}_S^{(SA)} \circ \mathcal{F}_1^{(FF)}$ and $\mathcal{F}_2^{(FF)} \circ \mathcal{F}_S^{(SA)} \circ \overline{\mathcal{F}}_1^{(FF)}$.

The difference is caused by the difference between $\overline{\mathcal{F}}_1^{(FF)}$ and $\mathcal{F}_1^{(FF)}$. By (F.9), this difference is bounded by $\frac{1}{D}$ in every dimension, for any input $Z \in \mathbb{R}^{d \times n}$:

$$\|\overline{\mathcal{F}}_1^{(FF)}(Z) - \mathcal{F}_1^{(FF)}(Z)\|_2 < \frac{\sqrt{dn}B_x}{D}$$

By (F.19):

$$\begin{aligned} \|\mathcal{F}_{S}^{(SA)} \circ \overline{\mathcal{F}}_{1}^{(FF)}(Z) - \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)\|_{2} \\ &\leq \|\mathcal{F}_{S}^{(SA)} \circ \overline{\mathcal{F}}_{1}^{(FF)}(Z) - \overline{\mathcal{F}}_{1}^{(FF)}(Z)\|_{2} + \|\overline{\mathcal{F}}_{1}^{(FF)}(Z) - \mathcal{F}_{1}^{(FF)}(Z)\|_{2} \\ &\leq \|\mathcal{F}_{S}^{(SA)} \circ \overline{\mathcal{F}}_{1}^{(FF)}(Z) - \overline{\mathcal{F}}_{1}^{(FF)}(Z)\|_{2} + \|\overline{\mathcal{F}}_{1}^{(FF)}(Z) - \mathcal{F}_{1}^{(FF)}(Z)\|_{2} \\ &+ \|\mathcal{F}_{1}^{(FF)}(Z) - \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}(Z)\|_{2} \qquad \text{(By triangle inequality)} \\ &\leq \frac{\sqrt{dn}B_{x}}{D} + 2 \cdot \frac{\sqrt{n}B_{x}}{4D}. \qquad \text{(By } \|A\|_{2} \leq \|A\|_{F} \text{ and } (F.19) \end{aligned}$$

In the section on quantization of the input, we used piece-wise linear functions (F.7) to approximate piece-wise-constant functions (F.8), this creates a deviation for the inputs on the boundaries of the constant regions. Consider Z as one of these inputs whose value deviated from $\mathcal{F}_2^{(FF)} \circ \mathcal{F}_S^{(SA)} \circ$ $\overline{\mathcal{F}}_1^{(FF)}(Q(Z))$. Let $f(L_1)$ denote the value given to $\mathcal{F}_2^{(FF)} \circ \mathcal{F}_S^{(SA)} \circ \mathcal{F}_1^{(FF)}(Z)$. Because the deviation take the output to a grid at most $\sqrt{dn}B_x/D + \sqrt{n}B_x/2D$ away from its original grid, under the quantization of the output, $f(L_1)$ at most deviate from its original output $\mathcal{F}_2^{(FF)} \circ \mathcal{F}_S^{(SA)} \circ$ $\overline{\mathcal{F}}_1^{(FF)}(Z)$ by the distance of $\sqrt{dn}B_x/D + \sqrt{n}B_x/2D$ aggregated with 2 times of the maximal distance within a grid. They sum up to be:

Lastly, by condition we neglect the $\mathbb{G}_D \setminus \widetilde{\mathbb{G}}_D$ part. This yields:

$$\mathcal{F}_2^{(FF)} \circ \mathcal{F}_S^{(SA)} \circ \overline{\mathcal{F}}_1^{(FF)} = \overline{f}.$$

Thus, adding up the errors yields:

$$\begin{split} \|f - \mathcal{F}_{2}^{(FF)} \circ \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}\|_{2} \\ \leq \|f - \overline{f}\|_{2} + \|\overline{f} - \mathcal{F}_{2}^{(FF)} \circ \mathcal{F}_{S}^{(SA)} \circ \mathcal{F}_{1}^{(FF)}\|_{2} \qquad \text{(By triangle inequality)} \\ = L \frac{6\sqrt{dn}B_{x} + \sqrt{n}B_{x}}{2D} + L \frac{\sqrt{dn}B_{x}}{D} \qquad \text{(By (F.4))} \\ = \frac{L(8\sqrt{dn} + \sqrt{n})B_{x}}{2D}. \end{split}$$

2430	This secondates the need
2431	This completes the proof.
2432	
2433	
2434	
2435	
2436	
2437	
2438	
2439	
2440	
2441	
2442	
2443	
2444	
2445	
2446	
2447	
2448	
2449	
2450	
2451	
2452	
2453	
2454	
2455	
2456	
2457	
2458	
2459	
2460	
2461	
2462	
2463	
2464	
2465	
2466	
2467	
2468	
2469	
2470	
2471	
2472	
2473	
2474	
2475	
2476	
2477	
2478	
2479	
2480	
2481	
2482	
2403	

2484 G EXPERIMENTAL DETAILS

2489

2490 2491

2492

2493

2494

2495

2496

2497

2498

2499

2500

2501

2502

2524

2527 2528

2536

2537

In this section, we conduct experiments to verify the capability of ICL to learn deep feed-forward
neural networks. We conduct the experiments based on 3-layer NN, 4-layer NN and 6-layer NN using
both ReLU-Transformer and Softmax-Transformer based on the GPT-2 backbone.

Experimental Objectives. Our objectives include the following three parts:

- **Objective 1.** Validating the performance of ICL matches that of training *N*-layer networks, i.e., the results in Theorem 1, Theorem 5, and Theorem 6.
- **Objective 2.** Validating the ICL performance in scenarios where the testing distribution diverges from the pretraining one or where prompt lengths exceed those used in pretraining.
 - **Objective 3.** Validating the ICL performance in scenarios where the distribution of parameters in the *N*-layer network diverges from that of the pretraining phase.
 - **Objective 4.** Validating that a deeper transformer achieves better ICL performance, supporting the idea that scaling up the transformer enables it to perform more ICGD steps.

Computational Resource. We conduct all experiments using 1 NVIDIA A100 GPU with 80GB of memory. Our code is based on the PyTorch implementation of the in-context learning for the transformer (Garg et al., 2022) at https://github.com/dtsip/in-context-learning.

2503 G.1 EXPERIMENTS FOR OBJECTIVES 1 AND 2 2504

In this section, we conduct experiments to validate Objectives 1 and 2. We sample the input of feed-forward network $x \in \mathbb{R}^d$ from the Gaussian mixture distribution: $w_1N(-2, I_d) + w_2N(2, I_d)$, where $w_1, w_2 \in \mathbb{R}$. We consider three kinds of network $f : \mathbb{R}^d \to \mathbb{R}$, (i) 3-layer NN, (ii) 4-layer NN, and (iii) 6-layer NN. We generate the true output by y = f(x). In our setting, we use d = 20.

Model Architecture. The sole difference between ReLU-Transformer and Softmax-Transformer is the activation function in the attention layer. Both models comprise 12 transformer blocks, each with 8 attention heads, and share the same hidden and MLP dimensions of 256.

Transformer Pretraining. We pretrain the ReLU-Transformer and Softmax-Transformer based on the GPT-2 backbone. In our setting, we sample the pertaining data from $N(-2, I_d)$, i.e., $w_1 = 1$ and $w_2 = 0$. Following the pre-training method in (Garg et al., 2022), we use the batch size as 64. To construct each sample in a batch, we use the following steps (take the generation for the *i*-th sample as an example):

- 1. Initialize the parameters in f_i with a standard Gaussian distribution, i.e., N(0, I).
- 2. Generate *n* queries $\{x_{i,j}\}_{j=1}^{n}$ (i.e., input of f_i) from the Gaussian mixture model $\omega_1 N(-2, I_d) + \omega_2 N(2, I_d)$. Here we take n = 51.
- 2521 3. For each query $x_{i,j}$, use $y_{i,j} = f_i(x_{i,j})$ to calculate the true output.

²⁵²² This generates a training sample for the transformer model with inputs

$$[x_{i,1}, y_{i,1}, \cdots, x_{i,50}, y_{i,50}, x_{i,51}],$$

and training target

$$o_i = [y_{i,1}, \cdots, y_{i,50}, y_{i,51}].$$

We use the MSE loss between prediction and true value of o_i . The pretraining process iterates for 500k steps.

Testing Method. We generate samples similar to the pretraining process. The batch size is 64, and the number of batch is 100, i.e., we have 6400 samples totally. For each sample, we extend the value n from 51 to 76 to learn the performance of in-context learning when the prompt length is longer than we used in pretraining. The input to the model becomes

 $[x_{i,1}, y_{i,1}, \cdots, x_{i,75}, y_{i,75}, x_{i,76}].$

We assess performance using the mean R-squared value for all 6400 samples.

Baseline. We use the 3-layer, 4-layer, and 6-layer feed-forward neural networks with 200 hidden dimensions as baselines by training them with in-context examples. Specially, given a testing sample (take the *i*-th sample as an example), which includes prompts $\{x_{i,j}, y_{i,j}\}_{j=1}^{k-1}$ and a test query $x_{i,k}$. We use $\{x_{i,j}, y_{i,j}\}_{j=1}^{k-1}$ to train the network with MSE loss for 100 epochs. We select the highest R-squared value from each epoch as the testing measure and calculate the average across all 6400 samples.

G.1.1 PERFORMANCE OF RELU TRANSFORMER.

We use three different Gaussian mixture distribution $\omega_1 N(-2, I_d) + \omega_2 N(2, I_d)$ for the testing data: (i) $\omega_1 = 1, \omega_2 = 0$, (ii) $\omega_1 = 0.9, \omega_2 = 0.1$, (iii) $\omega_1 = 0.7, \omega_2 = 0.3$, (iv) $\omega_1 = 0.5, \omega_2 = 0.5$. Here the distribution in the first setting matches the distribution in pretraining. We show the results in Figure 1.



Figure 1: **Performance of ICL in ReLU-Transformer:** ICL learns 3-layer, 4-layer, and 6-layer NN and achieves R-squared values comparable to those from training with prompt samples. The results also show the ICL performance declines as the testing distribution diverges from the pretraining one.

G.1.2 PERFORMANCE OF SOFTMAX TRANSFORMER.

We use three different Gaussian mixture distribution $\omega_1 N(-2, I_d) + \omega_2 N(2, I_d)$ for the testing data: (i) $\omega_1 = 1, \omega_2 = 0$, (ii) $\omega_1 = 0.9, \omega_2 = 0.1$, (iii) $\omega_1 = 0.7, \omega_2 = 0.3$, (iv) $\omega_1 = 0.5, \omega_2 = 0.5$. Here the distribution in the first setting matches the distribution in pretraining. We show the results in Figure 2.



Figure 2: Performance of ICL in Softmax-Transformer: ICL learns 3-layer, 4-layer, and 6-layer
NN and achieves R-squared values comparable to those from training with prompt samples. The
results also show the ICL performance declines as the testing distribution diverges from the pretraining
one. Note that performance decreases when the prompt length exceeds the pretraining length (i.e.,
s0), a well-known issue (Dai et al., 2019; Anil et al., 2022). We believe this is due to the absolute
positional encodings in GPT-2, as noted in (Zhang et al., 2023)

The results in Appendix G.1.1 and Appendix G.1.2 show that the performance of ICL in the transformer matches that of training N-layer networks, regardless of whether the prompt lengths are within or exceed those used in pretraining. Furthermore, the ICL performance declines as the testing distribution diverges from the pretraining one.

2588 2589

2590

2545

2546 2547

2548

2549

2550

2563

2564 2565

2566

2567

- G.2 EXPERIMENTS FOR OBJECTIVE 3
- In this section, we conduct experiments to validate Objective 3. For these experiments, we use testing data that is identical to the training data, which follows a distribution of $N(-2, I_d)$. We vary the

2592 distribution of parameters in the N-layer network. During the training process, we set the distribution 2593 as N(0, I). In the testing process, we examine different distributions, including N(0, I), N(-0.5, I), 2594 and N(0.5, I). All other model hyperparameters and experimental details remain consistent with 2595 those described in Appendix G.1. We evaluate the ICL performance of both the ReLU-Transformer 2596 and the Softmax-Transformer for 4-layer networks, as shown in Figure 3 and Figure 4. The results demonstrate that the ICL performance in the transformer matches that of training N-layer networks, 2597 regardless of whether the parameter distribution in the N-layer network diverges from that of the 2598 pretraining phase. 2599



Figure 3: Performance of ICL Across Various *N*-layer Network Parameter Distributions for the **ReLU-Transformer:** ICL learns 4-layer NN and achieves R-squared values comparable to those from training with prompt samples, even when the parameter distribution in the *N*-layer network during testing diverges from that in the pretraining phase (N(0, I)).



Figure 4: Performance of ICL Across Various N-layer Network Parameter Distributions for the Softmax-Transformer: ICL learns 4-layer NN and achieves R-squared values comparable to those from training with prompt samples, even when the parameter distribution in the N-layer network during testing diverges from that in the pretraining phase (N(0, I)).

2625 2626 G.3 EXPERIMENTS FOR OBJECTIVE 4

2627 In this section, we conduct experiments to validate Objective 4. For these experiments, we use testing 2628 data identical to the pertaining data from $N(-2, I_d)$. We vary the number of layers in the transformer 2629 architecture, testing configurations with 4, 6, 8 and 10 layers. All other model hyperparameters 2630 and experimental details remain consistent with those described in Appendix G.1. We evaluate the 2631 ICL performance of both the ReLU-Transformer and the Softmax-Transformer with 15, 30, and 45 in-context examples, as shown in Figure 5. The results show that a deeper transformer achieves better 2632 ICL performance, supporting the idea that scaling up the transformer enables it to perform more 2633 ICGD steps. 2634

2635 2636

2609

2610

- 2637 2638
- 2030
- 2039 2640
- 2641
- 2642
- 2643
- 2644
- 2645



Figure 5: Performance of ICL Across Varying Transformer Depths: We use the number of in-context examples as 15, 30, or 45 for both the ReLU-Transformer and the Softmax-Transformer. The results show that a deeper transformer achieves better ICL performance, supporting the idea that scaling up the transformer enables it to perform more ICGD steps.

2700 H APPLICATION: ICL FOR DIFFUSION SCORE APPROXIMATION

In this part, we give an important application of our work, i.e., learn the score function of diffusion models by the in-context learning of transformer models. We give the preliminaries about score matching generative diffusion models in Appendix H.1. Then, we give the analysis for ICL to approximate the diffusion score function in Appendix H.2.

2707 H.1 SCORE MATCHING GENERATIVE DIFFUSION MODELS

Diffusion Model. Let $x_0 \in \mathbb{R}^d$ be initial data following target data distribution $x_0 \sim P_0$. In essence, a diffusion generative model consists of two stochastic process in \mathbb{R}^d :

• A forward process gradually add noise to the initial data (e.g., images): $x_0 \rightarrow x_1 \rightarrow \cdots \rightarrow x_T$.

• A backward process gradually remove noise from pure noise: $y_T \to y_{T-1} \to \cdots \to y_0$.

Importantly, the backward process is the reversed forward process, i.e., $y_t \approx x_{T-t}$ for $i \in 0, \dots, T^2$. 2714 This allows the backward process to reconstruct the initial data from noise, and hence generative. To 2715 achieve this time-reversal, a diffusion model learns the reverse process by ensuring the backward 2716 conditional distributions mirror the forward ones. The most prevalent technique for aligning these 2717 conditional dynamics is through "score matching" — a strategy training a model to match score 2718 function, i.e., the gradients of the log marginal density of the forward process (Song et al., 2020b;a; 2719 Vincent, 2011). To be precise, let $P_t, p_t(\cdot)$ denote the distribution function and destiny function of 2720 x_t . The score function is given by $\nabla \log p_t(\cdot)$. In this work, we focus on leveraging the in-context 2721 learning (ICL) capability of transformers to emulate the score-matching training process.

2722 Score Matching Loss. We introduce the basic setting of score-matching as follows³. To estimate the score function, we use the following loss to train a score network $s_W(\cdot, t)$ with parameters W:

$$\min_{W} \int_{T_0}^T \gamma(t) \mathbb{E}_{x_t \sim P_t} \left[\|s_W(x_t, t) - \nabla \log p_t(x_t)\|_2^2 \right] \mathrm{d}t, \quad \text{where } \gamma(t) \text{ is a weight function, (H.1)}$$

and T_0 is a small value for stabilizing training and preventing the score function from diverging. In practice, as $\nabla \log p_t(\cdot)$ is unknown, we minimize the following equivalent loss (Vincent, 2011).

$$\min_{W} \int_{T_0}^{T} \gamma(t) \mathbb{E}_{x_0 \sim P_0} \left[\mathbb{E}_{x_t \mid x_0} \left[\|s_W(x_t, t) - \nabla \log p(x_t \mid x_0)\|_2^2 \right] \right] \mathrm{d}t, \tag{H.2}$$

where $p(x_t|x_0)$ is distribution of x_t conditioned on x_0 .

2736 H.2 ICL FOR SCORE APPROXIMATION

2738 We first give the problem setup about the ICL for score approximation as the following:

Problem 3 (In-Context Learning (ICL) for Score Function $\nabla \log p_t(\cdot)$). Consider the score function $\nabla \log p_t(\cdot)$ for any $t \ge 0$. Given a dataset $\mathcal{D}_n \coloneqq \{(x_i, y_i)\}_{i \in [n]}$, where $\{x_i\}_{i \in [n]} \subseteq \mathbb{R}^d$ and $y_i = \nabla \log p_{t_i}(x_i) \subseteq \mathbb{R}^d$ $(t_i \ge 0)$, and a test input x_{n+1} , the goal of "ICL for Score Function" is to find a transformer \mathcal{T} to predict y_{n+1} based on x_{n+1} and the in-context dataset \mathcal{D}_n . In essence, the desired transformer \mathcal{T} serves as the trained score network $s_W(\cdot, t)$.

To solve Problem 3, we follow two steps: (i) Approximate the diffusion score function $\nabla \log p_t(\cdot)$ with a multi-layer feed-forward network with ReLU activation functions under the given training dataset \mathcal{D}_n . (ii) Approximate the gradient descent used to train this network by the in-context learning of the Transformer until convergence, using the same training set \mathcal{D}_n as the prompts of ICL.

For the first step, we follow the score approximation results based on a multi-layer feed-forward network with ReLU activation in (Chen et al., 2023), stated as next lemma.

2751 2752 2753

2706

2708

2711

2712

2713

2725 2726 2727

2731 2732 2733

2735

2737

2739

2740

2741

2742

2743

 $^{{}^{2} \}stackrel{a}{\approx}$ denotes distributional equivalence.

³Please also see Appendix B.1 and (Chen et al., 2024; Chan, 2024; Yang et al., 2023) for overviews.

2754 2755 2756	Lemma 19 (Score Approximation by Feed-Forward Networks, Theorem 1 of (Chen et al., 2023)). Given an approximation error $\epsilon > 0$, for any initial data distribution P_0 , there exist a multi-layer feed-forward network with ReI II activation $f(w, x, t) : \mathbb{R}^{D_w} \times \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^d$. Then for any
2757	$t \in [T_0, T]$, we have $\ f(w, \cdot, t) - \nabla \log p_t(\cdot)\ _{L^2(P_t)} \leq \mathcal{O}(\epsilon)$.
2759 2760 2761	With the approximation result, we reduce the Problem 3 to Problem 2, where the loss function is (H.1). Following Theorem 1, we show that the in-context learning of transformer models can approximate the score function of diffusion model.
2762	
2763	
2764	
2765	
2766	
2767	
2768	
2769	
2770	
2771	
2772	
2773	
2774	
2775	
2776	
2777	
2778	
2779	
2780	
2781	
2782	
2783	
2784	
2785	
2786	
2787	
2788	
2789	
2790	
2791	
2792	
2793	
2794	
2795	
2797	
2798	
2799	
2800	
2801	
2802	
2803	
2804	
2805	
2806	
2807	