

# Approximation Bounds for Transformer Networks with Application to Regression

Yuling Jiao, Yanming Lai, Defeng Sun, Yang Wang, Bokai Yan

## Abstract

We explore the approximation capabilities of Transformer networks for Hölder and Sobolev functions, and apply these results to address nonparametric regression estimation with dependent observations. First, we establish novel upper bounds for standard Transformer networks approximating sequence-to-sequence mappings whose component functions are Hölder continuous with smoothness index  $\gamma \in (0, 1]$ . To achieve an approximation error  $\varepsilon$  under the  $L^p$ -norm for  $p \in [1, \infty]$ , it suffices to use a fixed-depth Transformer network whose total number of parameters scales as  $\varepsilon^{-d_* n / \gamma}$ . This result not only extends existing findings to include the case  $p = \infty$ , but also matches the best known upper bounds on number of parameters previously obtained for fixed-depth FNNs and RNNs. Similar bounds are also derived for Sobolev functions. Second, we derive explicit convergence rates for the nonparametric regression problem under various  $\beta$ -mixing data assumptions, which allow the dependence between observations to weaken over time. Our bounds on the sample complexity impose no constraints on weight magnitudes. Lastly, we propose a novel proof strategy to establish approximation bounds, inspired by the Kolmogorov-Arnold representation theorem. We show that if the self-attention layer in a Transformer can perform column averaging, the network can approximate sequence-to-sequence Hölder functions, offering new insights into the interpretability of self-attention mechanisms.

## 1 Introduction

Transformers [60] have become the cornerstone of modern deep learning, driving breakthroughs across multiple domains, including natural language processing [10], large language models [44], computer vision [15], and generative models [46]. Although their high performance has led to widespread use in practice, significant theoretical efforts are still underway to explain exactly what contributes to their success.

An important aspect of Transformers is their expressive capacity, which refers to their ability to effectively approximate target functions. As early as the 1980s, researchers established the universal approximation property for neural networks (e.g., [8, 25]), demonstrating that feed-forward neural networks (FNNs) can approximate any continuous function to any precision. With the rise of deep neural networks in recent years, many works have focused on the approximation theory of neural networks. For example, [31, 39, 66] studied approximation rates of deep ReLU FNNs for smooth functions, while [64, 65] and [32] examined, respectively, shallow ReLU FNNs and ReLU recurrent neural networks (RNNs). However, analyses based on Transformer architectures have rarely been observed. For a representative example, [69] showed the universal

---

Yuling Jiao is with the School of Artificial Intelligence and the School of Mathematics and Statistics, Wuhan University, Wuhan, China (email: yulingjiaomath@whu.edu.cn).

Yanming Lai, Yang Wang and Bokai Yan are with the Department of Mathematics, Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong, China (email: ylaiam@connect.ust.hk, yangwang@ust.hk, byanac@connect.ust.hk).

Defeng Sun is with the Department of Applied Mathematics and the Research Center for Intelligent Operations Research, The Hong Kong Polytechnic University, Hung Hom, Hong Kong, China (email: defeng.sun@polyu.edu.hk).

approximation of Transformers, which can approximate sequence-to-sequence continuous functions under the  $L^p$ -norm with  $p \in [1, \infty)$ . [33] revealed that even a Transformer with a single attention layer is a universal approximator. [58] investigated a special class of Transformers with infinite dimensional inputs. [63] established the approximation rates of looped Transformers, which reuse the same Transformer layer iteratively, by defining the modulus of continuity for sequence-to-sequence functions. [28] derived approximation rate estimates for continuous Transformers by defining novel complexity measures for nonlinear sequence relationships. Most recently, [59] showed that Transformers can approximate column-symmetric polynomials. Another line of research has explored replacing the softmax function in the attention mechanism with more tractable alternatives [21, 22, 30]. Despite these advances, the approximation rates for general functions and the performance under the  $L^\infty$ -norm using the standard Transformer architecture still remain unclear.

Another significant concern pertains to Transformer performance in sequence modeling, specifically regarding how Transformers capture relationships within sequential data. Recent theoretical studies have provided diverse insights into this issue. For example, [16] analyzed the inductive biases inherent in self-attention mechanisms, demonstrating that sample complexity scales only logarithmically in the sequence length. [7] introduced architectural adjustments to Transformers, enabling exact recognition of formal languages through enhanced long-range dependency modeling. [5] offered an interpretation of Transformer weight matrices as associative memory systems, distinguishing between long-term parametric storage and short-term contextual memory. [48] proved that a pretrained Transformer, when appropriately prompted or prefix-tuned, can approximate any sequence-to-sequence function. [27] showed that soft prompt tuning yields a universal approximator for Lipschitz sequence-to-sequence mappings. Furthermore, [62] conducted a rigorous analysis emphasizing how architectural components such as depth, attention heads, and feed-forward layers influence performance in tasks necessitating extensive, sparse memories. Recent developments in nonparametric regression estimation based on neural networks have mainly relied on assumptions of independent and identically distributed (i.i.d.) observations drawn from an unknown distribution [19, 31, 37, 43, 50, 57, 64]. However, sequential tasks typically exhibit temporal dependence, making the i.i.d. assumption too restrictive. Some recent studies have relaxed this assumption by considering that observations are drawn from a stationary mixing distribution, where the dependence between observations weakens over time [20, 29, 32, 49]. Despite these advancements, a gap remains regarding the capability of Transformers, explicitly designed to handle sequential data and temporal dependencies, in regression tasks involving dependent observations.

In this paper, we investigate the approximation of Hölder and Sobolev functions using Transformer networks and study nonparametric regression estimation under dependent observations. Our main contributions are as follows:

- We derive novel upper bounds on the approximation of standard Transformer architectures for Hölder and Sobolev functions. Specifically, to approximate a sequence-to-sequence mapping, where each component function is Hölder continuous with smoothness index  $\gamma \in (0, 1]$ , to approximation error  $\varepsilon$  under the  $L^p$ -norm for  $p \in [1, \infty]$ , it suffices to use a Transformer network whose total number of parameters scales as  $\varepsilon^{-d_x n / \gamma}$ . Our result establishes explicit approximation rates and extends existing findings to include the case  $p = \infty$ . Moreover, the number of parameters matches the best known upper bounds previously established for fixed-depth FNNs and RNNs. Similar results are also derived for Sobolev functions.
- We present a comprehensive error analysis for the nonparametric regression problem with weakly dependent data. We achieve rates of  $m^{-\frac{\gamma}{\gamma + d_x n}}$ ,  $m^{-\frac{r\gamma}{(r+2)\gamma + (r+1)d_x n}}$  and  $m^{-\frac{\gamma}{\gamma + d_x n}}$  up to logarithmic factors corresponding respectively to geometrically  $\beta$ -mixing, algebraically  $\beta$ -mixing, and i.i.d. data assumptions, where  $m$  denotes the sample size and the parameter  $r$  controls the strength of dependence. We also establish upper bounds on the sample complex-

ity of Transformer networks, notably without imposing constraints on the weight size of the network.

- We propose a novel proof strategy for establishing approximation bounds of Transformer networks, inspired by the Kolmogorov-Arnold representation theorem. By observing that the self-attention layers merely compute column-wise averages in our analysis, we demonstrate that the softmax function can be generalized to broader alternatives. This viewpoint provides new insights into the interpretability of self-attention mechanisms.

## 1.1 Organization

The rest of the paper is organized as follows. In Section 2, we define the Transformer architecture, describe the setup of the nonparametric regression problem, and list our main results. In Section 3, we present discussions and related works. All proofs are provided in Section 4.

## 2 Summary of Results

*Notation.* We use bold lowercase letters to represent vectors and bold uppercase letters to represent matrices. For any vector  $\mathbf{v} \in \mathbb{R}^d$ , we denote by  $v_i$  the  $i$ -th element of  $\mathbf{v}$ . For any matrix  $\mathbf{A} \in \mathbb{R}^{d \times n}$ , we denote its  $i$ -th row by  $\mathbf{A}_{i,:}$ , its  $j$ -th column by  $\mathbf{A}_{:,j}$  and the element at its  $i$ -th row and  $j$ -th column by  $A_{i,j}$ . We denote the all-zero and all-one vectors of length  $n$  by  $\mathbf{0}_n$  and  $\mathbf{1}_n$ , respectively. The identity matrix of size  $n$  is denoted by  $\mathbf{I}_n$ . The zero matrix of size  $m \times n$  is denoted by  $\mathbf{O}_{m,n}$ . When the dimensions are clear from the context, we omit the subscripts for brevity. For  $m \in \mathbb{N}$ , we write  $[m] := 1, \dots, m$ . We use  $\mathbb{N}_0$  to denote the set of nonnegative integers and  $\mathbb{N}_0^d = \{(\alpha_1, \alpha_2, \dots, \alpha_d) : \alpha_k \in \mathbb{N}_0, \forall k \in [d]\}$  to denote the set of  $d$ -dimensional multi-index. For a multi-index  $\boldsymbol{\alpha} \in \mathbb{N}_0^d$ , we denote  $\|\boldsymbol{\alpha}\|_{\ell_1} = \alpha_1 + \alpha_2 + \dots + \alpha_d$ . For a finite set  $\mathbb{G}$ , we use  $|\mathbb{G}|$  to denote its cardinality. For two sequences  $\{a_n\}$  and  $\{b_n\}$ , we use the notation  $a_n \lesssim b_n$  and  $a_n \gtrsim b_n$  to indicate  $a_n \leq c_1 b_n$  and  $a_n \geq c_2 b_n$ , respectively, for some constants  $c_1, c_2 > 0$  that are independent of  $n$ . Furthermore,  $a_n \asymp b_n$  means that both  $a_n \lesssim b_n$  and  $a_n \gtrsim b_n$  hold. In our analysis, we use  $\sigma_S$  to represent the column-wise softmax function. Specifically, for a matrix  $\mathbf{A} \in \mathbb{R}^{d \times n}$ ,  $\sigma_S[\mathbf{A}] \in \mathbb{R}^{d \times n}$  is computed as  $\sigma_S[\mathbf{A}]_{i,j} := \exp(A_{i,j}) / \sum_{k=1}^d \exp(A_{k,j})$ . The ReLU activation function is denoted by  $\sigma_R[x] := \max\{x, 0\}$ . In contrast to  $\sigma_S$ ,  $\sigma_R$  operates element-wise, regardless of whether the input is a vector or a matrix. Let  $\Omega \subseteq \mathbb{R}^{d \times n}$  be a bounded domain. For  $1 \leq p < \infty$ , the  $L^p$ -norm of a real-valued function  $f : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}$  is defined as  $\|f\|_{L^p(\Omega)} := (\int_{\Omega} |f(\mathbf{X})|^p d\mathbf{X})^{1/p}$ , and for  $p = \infty$ , it is given by  $\|f\|_{L^\infty(\Omega)} := \text{ess sup}_{\mathbf{X} \in \Omega} |f(\mathbf{X})|$ . For a matrix-valued function  $\mathbf{F} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{m \times n}$ , the  $L^p$ -norm is defined as  $\|\mathbf{F}\|_{L^p(\Omega)} := (\int_{\Omega} \|\mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X})^{1/p}$  for  $1 \leq p < \infty$ , and for  $p = \infty$ ,  $\|\mathbf{F}\|_{L^\infty(\Omega)} := \text{ess sup}_{\mathbf{X} \in \Omega} \|\mathbf{F}(\mathbf{X})\|_F$ .

### 2.1 Approximation Rates for Hölder and Sobolev Functions

We begin by introducing the architecture of Transformers, following the notations in [36] and [34]. A Transformer network is a sequence-to-sequence function  $\mathbb{R}^{d_x \times n} \rightarrow \mathbb{R}^{d_y \times n}$ , comprising three main components: the self-attention layer, the (token-wise) feed-forward layer, and the embedding layer.

**Embedding and projection layer:** For embedding dimension  $D \in \mathbb{N}$ , the embedding and projection layers connect the input, hidden, and output spaces. The embedding layer  $\mathcal{E}_{in} : \mathbb{R}^{d_x \times n} \rightarrow \mathbb{R}^{D \times n}$  is defined as

$$\mathcal{E}_{in}(\mathbf{X}) := \mathbf{E}_{in}\mathbf{X} + \mathbf{P} \in \mathbb{R}^{D \times n},$$

where  $\mathbf{E}_{in} \in \mathbb{R}^{D \times d_x}$  is a learnable weight matrix, and  $\mathbf{P} \in \mathbb{R}^{D \times n}$  is a trainable positional encoding matrix. Since self-attention and feed-forward layers are permutation equivariant,  $\mathbf{P}$  is

introduced to provide positional information and break this equivariance. The projection layer  $\mathcal{E}_{out} : \mathbb{R}^{D \times n} \rightarrow \mathbb{R}^{d_x \times n}$  is defined as

$$\mathcal{E}_{out}(\mathbf{Y}) := \mathbf{E}_{out} \mathbf{Y} \in \mathbb{R}^{d_y \times n},$$

where  $\mathbf{E}_{out} \in \mathbb{R}^{d_y \times D}$  maps the high-dimensional hidden representation onto the output space.

**Self-attention layer:** Given a sequence  $\mathbf{Z} \in \mathbb{R}^{D \times n}$ , composed of  $n$  tokens, each with an embedding dimension  $D$ , the  $l$ -th self-attention layer  $\mathcal{F}_l^{(SA)} : \mathbb{R}^{D \times n} \rightarrow \mathbb{R}^{D \times n}$  is defined as

$$\mathcal{F}_l^{(SA)}(\mathbf{Z}) := \mathbf{Z} + \sum_{h=1}^H \mathbf{W}_{h,l}^{(O)} \left( \mathbf{W}_{h,l}^{(V)} \mathbf{Z} \right) \sigma_S \left[ \left( \mathbf{W}_{h,l}^{(K)} \mathbf{Z} \right)^\top \left( \mathbf{W}_{h,l}^{(Q)} \mathbf{Z} \right) \right] \in \mathbb{R}^{D \times n},$$

where  $\mathbf{W}_{h,l}^{(V)}$ ,  $\mathbf{W}_{h,l}^{(K)}$ ,  $\mathbf{W}_{h,l}^{(Q)}$   $\in \mathbb{R}^{S \times D}$  and  $\mathbf{W}_{h,l}^{(O)} \in \mathbb{R}^{D \times S}$  are the value, key, query, and projection matrices for head  $h \in [H]$  with head size  $S$ , respectively.

**Feed-forward layer:** The output  $\mathbf{Z} \in \mathbb{R}^{D \times n}$  of the self-attention layer is then passed to the feed-forward layer, given by

$$\mathcal{F}_l^{(FF)}(\mathbf{Z}) := \mathbf{Z} + \mathbf{W}_l^{(2)} \sigma_R \left[ \mathbf{W}_l^{(1)} \mathbf{Z} + \mathbf{b}_l^{(1)} \mathbf{1}_n^\top \right] + \mathbf{b}_l^{(2)} \mathbf{1}_n^\top \in \mathbb{R}^{D \times n},$$

where  $\mathbf{W}_l^{(1)} \in \mathbb{R}^{W \times D}$  and  $\mathbf{W}_l^{(2)} \in \mathbb{R}^{D \times W}$  are weight matrices with hidden dimension  $W$ , and  $\mathbf{b}_l^{(1)} \in \mathbb{R}^W$ ,  $\mathbf{b}_l^{(2)} \in \mathbb{R}^D$  are bias terms.

The class of Transformer networks is then defined as

$$\mathcal{T}_{d_x, d_y}(D, H, S, W, L) := \left\{ \mathcal{E}_{out} \circ \mathcal{F}_L^{(FF)} \circ \mathcal{F}_L^{(SA)} \circ \dots \circ \mathcal{F}_1^{(FF)} \circ \mathcal{F}_1^{(SA)} \circ \mathcal{E}_{in} \right\},$$

where  $D$  is the embedding dimension,  $H$  is the number of attention heads,  $S$  is the head size,  $W$  is the hidden dimension in the feed-forward layer, and  $L$  is the number of Transformer layers, each consisting of a self-attention and a feed-forward sublayer. When the dimensions are clear from the context, we use the simplified notation  $\mathcal{T}(D, H, S, W, L)$  for convenience. We list some basic properties of the Transformer class in Proposition 6. Let

$$N = Dd_x + Dn + d_y D + L(4HSD + 2WD + W + D) \lesssim (HS + W)DL \quad (1)$$

be the total number of training parameters in the Transformer network.

The purpose of this paper is to study the approximation of Hölder and Sobolev functions by Transformer networks. We recall the definitions of Hölder and Sobolev functions with bounded norm as follows.

**Definition 1** (Hölder functions). *Let  $\Omega$  be a bounded domain in  $\mathbb{R}^{d_x \times n}$  and  $\gamma \in (0, 1]$ . Given  $K_{\mathcal{H}} > 0$ , we denote the Hölder class  $\mathcal{H}^\gamma(\Omega, K_{\mathcal{H}})$  as*

$$\mathcal{H}^\gamma(\Omega, K_{\mathcal{H}}) = \left\{ f : \Omega \rightarrow \mathbb{R} : \|f\|_{L^\infty(\Omega)} + \sup_{\mathbf{X}, \mathbf{Y} \in \Omega, \mathbf{X} \neq \mathbf{Y}} \frac{|f(\mathbf{X}) - f(\mathbf{Y})|}{\|\mathbf{X} - \mathbf{Y}\|_F^\gamma} \leq K_{\mathcal{H}} \right\}.$$

**Definition 2** (Sobolev functions). *Let  $\Omega$  be a bounded domain in  $\mathbb{R}^{d_x \times n}$ . For  $p \in [1, \infty)$  and  $K_{\mathcal{W}} > 0$ , we denote the Sobolev class  $\mathcal{W}^{1,p}(\Omega, K_{\mathcal{W}})$  as*

$$\mathcal{W}^{1,p}(\Omega, K_{\mathcal{W}}) = \left\{ f : \Omega \rightarrow \mathbb{R} : \left( \sum_{\|\alpha\|_{\ell^1} \leq 1} \int_{\Omega} |D^\alpha f|^p d\mathbf{X} \right)^{1/p} \leq K_{\mathcal{W}} \right\},$$

and for  $p = \infty$ , the Sobolev class  $\mathcal{W}^{1,\infty}(\Omega, K_{\mathcal{W}})$  is defined as

$$\mathcal{W}^{1,\infty}(\Omega, K_{\mathcal{W}}) = \left\{ f : \Omega \rightarrow \mathbb{R} : \sum_{\|\alpha\|_{\ell^1} \leq 1} \text{ess sup}_{\Omega} |D^\alpha f| \leq K_{\mathcal{W}} \right\},$$

where  $\alpha \in \mathbb{N}_0^{d_x \times n}$  is a multi-index and  $D^\alpha$  is the weak derivative of order  $\alpha$ .

Hölder and Sobolev functions are central objects in approximation theory due to their close connection with polynomial and spline approximations [11]. A variety of embedding and interpolation results relate these spaces. For example, the Sobolev embedding  $\mathcal{W}^{1,p}(\Omega, K_{\mathcal{W}}) \hookrightarrow \mathcal{H}^{1-\frac{d}{p}}(\Omega, K_{\mathcal{H}})$  holds with the Hölder exponent  $\gamma = 1 - \frac{d}{p}$  and an appropriate constant  $K_{\mathcal{H}}$  depending on  $K_{\mathcal{W}}$  and the geometry of  $\Omega$ . In the limiting case where  $p = \infty$ , the Sobolev space  $\mathcal{W}^{1,\infty}(\Omega, K_{\mathcal{W}})$  consists of functions with essentially bounded first-order weak derivatives, which directly implies that these functions are Lipschitz continuous. In other words,  $\mathcal{W}^{1,\infty}(\Omega, K_{\mathcal{W}}) \hookrightarrow \mathcal{H}^1(\Omega, K_{\mathcal{H}})$ . For applications in machine learning, it is thus important to understand how efficiently Transformer networks can approximate functions in both the Hölder and Sobolev spaces.

We now present our main results on the approximation capabilities of Transformer networks for Hölder and Sobolev functions. We defer the proofs to Section 4.1.

**Theorem 1.** *Given  $\gamma \in (0, 1]$  and  $K_{\mathcal{H}} > 0$ , assume that the target function  $\mathbf{F} : [0, 1]^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  satisfies  $F_{i,j} \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$  for each  $i \in [d_x], j \in [n]$ . For any  $\varepsilon \in (0, 1)$  and  $p \in [1, \infty]$ , there exists a Transformer network*

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(D = C_1, H = C_2, S = C_3, W = C_4 \cdot \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil, L = C_5)$$

such that

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0, 1]^{d_x \times n})} \leq 4(d_x n)^2 K_{\mathcal{H}} \varepsilon,$$

where

1.  $C_1 = d_x, C_2 = 1, C_3 = 1, C_4 = 5n$  and  $C_5 = 2$  if  $p \in [1, \infty)$ ;
2.  $C_1 = 5d_x 3^{d_x n}, C_2 = 3^{d_x n}, C_3 = 1, C_4 = 5n 3^{d_x n}$  and  $C_5 = 2 + 2d_x n$  if  $p = \infty$ .

**Theorem 2.** *Given  $p \in [1, \infty)$  and  $K_{\mathcal{W}} > 0$ , assume that the target function  $\mathbf{F} : [0, 1]^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  satisfies  $F_{i,j} \in \mathcal{W}^{1,p}([0, 1]^{d_x \times n}, K_{\mathcal{W}})$  for each  $i \in [d_x], j \in [n]$ . For any  $\varepsilon \in (0, 1)$ , there exists a Transformer network*

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(D = d_x, H = 1, S = 1, W = 5n \cdot \lceil \varepsilon^{-d_x n} \rceil, L = 2)$$

such that

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0, 1]^{d_x \times n})} \leq 4C(d_x n)^2 K_{\mathcal{W}} \varepsilon,$$

where  $C$  is a constant depending only on  $d_x n$ .

We make several remarks regarding our results. First, ever since [60] proposed the Transformer architecture, there have been various theoretical analyses on its expressive capacity. A series of works established the universal approximation for sequence-to-sequence continuous functions under the  $L^p$ -norm for  $1 \leq p < \infty$ . More precisely, [69] showed that

$$\sup_{\mathbf{F}: F_{i,j} \in \mathcal{C}(\Omega)} \inf_{\mathcal{N} \in \mathcal{T}(D, H, S, W, L)} \|\mathcal{N} - \mathbf{F}\|_{L^p(\Omega)} \rightarrow 0$$

for fixed and sufficiently large  $D, H, S, W$  and as  $L \rightarrow \infty$ , where  $\mathcal{C}(\Omega)$  denotes the space of continuous functions on  $\Omega$ , and subsequently [33] showed that

$$\sup_{\mathbf{F}: F_{i,j} \in \mathcal{C}(\Omega)} \inf_{\mathcal{N} \in \mathcal{T}(D, H, S, W, L)} \|\mathcal{N} - \mathbf{F}\|_{L^p(\Omega)} \rightarrow 0$$

for fixed and sufficiently large  $D, H, S, L$  and as  $W \rightarrow \infty$ . We show that for  $1 \leq p \leq \infty$ ,

$$\sup_{\mathbf{F}: F_{i,j} \in \mathcal{H}^\gamma(\Omega, K_{\mathcal{H}})} \inf_{\mathcal{N} \in \mathcal{T}(D, H, S, W, L)} \|\mathcal{N} - \mathbf{F}\|_{L^p(\Omega)} \lesssim K_{\mathcal{H}} W^{-\gamma/(d_x n)}$$

and

$$\sup_{\mathbf{F}: F_{i,j} \in \mathcal{W}^{1,p}(\Omega, K_{\mathcal{W}})} \inf_{\mathcal{N} \in \mathcal{T}(D, H, S, W, L)} \|\mathcal{N} - \mathbf{F}\|_{L^p(\Omega)} \lesssim K_{\mathcal{W}} W^{-1/(d_x n)},$$

for fixed  $D, H, S, L$  and adjustable  $W$ . Our results not only provide the approximation rates for general Hölder and Sobolev functions, but also extend to the case  $p = \infty$ , which previous works were unable to address. These improvements are largely attributed to the use of the horizontal shift technique, which was originally introduced in [39] and further developed in [53, 71]. While their technique was developed for ReLU FNNs, we find that the ideas can be applied to the Transformer architecture. We summarize the related results in Table 1.

Second, We emphasize that our approximation results are established in a sequence-to-sequence sense; that is, every entry of the Transformer network  $\mathcal{N}$  simultaneously approximates the corresponding entry of the matrix-valued target function  $\mathbf{F}$ . It is not hard to extend the target function  $\mathbf{F} : [0, 1]^{d_x \times n} \rightarrow \mathbb{R}^{d_y \times n}$  to general  $d_y \in \mathbb{N}$  in Theorem 1 and 2.

Our results further demonstrate that Transformer networks possess stronger expressive capabilities than RNNs in approximating sequence-to-sequence functions. As discussed in [24, 32], RNNs are inherently limited to approximating past-dependent sequence-to-sequence functions because, at each time step, only the current and past tokens are utilized, leaving future tokens unprocessed. In contrast, Transformers have the advantage of accessing the entire input sequence. In other words, even the first output token of a Transformer depends on the entire input sequence, whereas in an RNN the first output token depends only on the first input token, the second on the first two, and so forth, owing to the sequential nature of RNNs. This distinction underpins our assertion that Transformer architectures outperform RNNs in terms of expressive power.

Third, we observe that to achieve an approximation error of  $\varepsilon$ , the total number of training parameters required scales as  $\varepsilon^{-d_x n / \gamma}$  for Hölder functions and as  $\varepsilon^{-d_x n}$  for Sobolev functions, matching the best known upper bounds previously established for fixed-depth FNNs and RNNs with input dimension  $d_x n$  [31, 32, 39, 53, 54, 66].

## 2.2 Nonparametric Regression

We then study the regression problem, which seeks to estimate an unknown target regression function from finite observations. We consider the following  $n$ -step prediction model

$$Y = f^*(X_1, X_2, \dots, X_n) + \varepsilon, \quad (2)$$

where  $Y \in \mathbb{R}$  is a response,  $f^*(\mathbf{x}_1, \dots, \mathbf{x}_n) = \mathbb{E}[Y | X_1 = \mathbf{x}_1, \dots, X_n = \mathbf{x}_n] : [0, 1]^{d_x \times n} \rightarrow \mathbb{R}$  is an unknown regression function and  $\varepsilon$  is a sub-Gaussian noise, independent of  $X_i, i = 1, \dots, n$ , with  $\mathbb{E}[\varepsilon] = 0$  and

$$\mathbb{E}[\exp(s\varepsilon)] \leq \exp\left(\frac{\sigma^2 s^2}{2}\right) \text{ for any } s \in \mathbb{R}.$$

Our purpose is to estimate the unknown target regression function  $f^*$  given observations  $\mathcal{D}_m = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$  which may not be i.i.d.

As observed in real sequence modeling applications, the sequential observations often exhibit temporal dependence, rendering the usual i.i.d. assumption inapplicable. This motivates us to consider dependent data. A frequently used alternative is to assume that observations are drawn from a stationary mixing distribution, where the dependence between observations weakens over time. This scenario has become standard and has been discussed extensively in previous studies [1, 38, 40–42, 49, 52, 55, 68]. We now introduce the relevant definitions.

**Definition 3** (Stationarity). *A sequence of random variables  $\{\mathbf{x}_t\}_{t=-\infty}^{\infty}$  is said to be stationary if for any  $t$  and non-negative integers  $m$  and  $k$ , the random vectors  $(\mathbf{x}_t, \dots, \mathbf{x}_{t+m})$  and  $(\mathbf{x}_{t+k}, \dots, \mathbf{x}_{t+m+k})$  have the same distribution.*

Table 1: Comparison of approximation rates.

Reference	Type	Target Function <sup>1</sup>	Metric <sup>2</sup>	Activations in Self-Attention Layers <sup>3</sup>	Width <sup>4</sup>	Depth
[69]	Universality	$\mathcal{C}^0$	$L^p$	$\sigma_S$ with bias		
[33]	Universality	$\mathcal{C}^0$	$L^p$	$\sigma_S$		
[18]	Universality	$\mathcal{C}^0$	$L^\infty$	$\sigma_H$		
[30]	Rate	$\mathcal{H}^\gamma$ $\mathcal{C}^m$	$L^\infty$	$X \odot \sigma_H(X)$	$\mathcal{O}(\varepsilon^{-d_x n/\gamma})$ $\mathcal{O}(\varepsilon^{-d_x n/m})$	$\mathcal{O}(\log \frac{1}{\varepsilon})$ $\mathcal{O}(\log \frac{1}{\varepsilon})$
[22]	Rate	$\mathcal{H}^\gamma$	$L^\infty$	$\sigma_R$	$\mathcal{O}(\varepsilon^{-d_x n/\gamma})$	$\mathcal{O}(\log \frac{1}{\varepsilon})$
<b>Ours</b> (Theorem 1)	Rate	$\mathcal{H}^\gamma$	$L^\infty$	$\sigma_S$	$\mathcal{O}(\varepsilon^{-d_x n/\gamma})$	$\mathcal{O}(1)$
<b>Ours</b> (Theorem 2)	Rate	$\mathcal{W}^{1,p}$	$L^p$	$\sigma_S$	$\mathcal{O}(\varepsilon^{-d_x n})$	$\mathcal{O}(1)$

<sup>1</sup>The space  $\mathcal{C}^m$  consists of all functions whose first  $m$  derivatives exist and are continuous, and  $\mathcal{C}^0$  denotes the space of continuous functions. <sup>2</sup> $p \in [1, \infty)$ . <sup>3</sup>Different Transformer architectures are obtained by replacing the softmax function in the self-attention layer with various activation functions. The symbol  $\odot$  denotes the Hadamard product. <sup>4</sup>Following [36], the width of a Transformer network is defined as  $\max\{D, HS, W\}$ . We omit constants independent of  $\varepsilon \in (0, 1)$ .

**Definition 4** ( $\beta$ -mixing). Let  $\{\mathbf{x}_t\}_{t=-\infty}^\infty$  be a stationary sequence of random variables. For any  $i, j \in \mathbb{Z} \cup \{-\infty, +\infty\}$ , let  $\sigma_i^j$  denote the  $\sigma$ -algebra generated by the random variables  $\mathbf{x}_k, i \leq k \leq j$ . Then, for any positive integer  $k$ , the  $\beta$ -mixing coefficient of the stochastic process  $\{\mathbf{x}_t\}_{t=-\infty}^\infty$  is defined as

$$\beta(k) = \sup_n \mathbb{E}_{B \in \sigma_{-n}^\infty} \left[ \sup_{A \in \sigma_{n+k}^\infty} |\mathbb{P}(A | B) - \mathbb{P}(A)| \right].$$

$\{\mathbf{x}_t\}_{t=-\infty}^\infty$  is said to be  $\beta$ -mixing if  $\beta(k) \rightarrow 0$  as  $k \rightarrow \infty$ . It is said to be algebraically  $\beta$ -mixing if there exist real numbers  $\beta_0 > 0$  and  $r > 0$  such that  $\beta(k) \leq \beta_0/k^r$  for all  $k$ , and geometrically  $\beta$ -mixing if there exist real numbers  $\beta_0, \beta_1 > 0$  and  $r > 0$  such that  $\beta(k) \leq \beta_0 \exp(-\beta_1 k^r)$  for all  $k$ .

In this work, we assume that the sequence of random variables  $\mathcal{X} = \{\mathbf{x}_t\}_{t=1}^m$  is drawn from a stationary  $\beta$ -mixing process. By Definition 3, the time index  $t$  does not affect the distribution of  $\mathbf{x}_t$  in a stationary sequence. Moreover, any  $n$  consecutive observations,  $(\mathbf{x}_{t-n+1}, \dots, \mathbf{x}_t)$ , share the same joint distribution, which we denote by  $\Pi$ . We assume that  $\Pi$  is supported on  $[0, 1]^{d_x \times n}$  and is absolutely continuous with respect to the Lebesgue measure, with its probability density function uniformly bounded above by a finite constant on  $[0, 1]^{d_x \times n}$ .

Definition 4 states that a sequence of random variables is mixing if the influence of past events on future events diminishes as the temporal gap increases. This definition provides a standard measure of the dependence among the random variables  $\{\mathbf{x}_t\}$  within a stationary sequence. We note that in certain special cases, such as Markov chains, the mixing coefficients admit upper bounds that can be estimated from data [26]. If  $\{\mathbf{x}_t\}_{t=1}^m$  are i.i.d. random variables, then by definition,  $\beta(k) = 0$  for all  $k \geq 1$ .

A fundamental method for estimating  $f^*$  is to minimize the mean squared error or the  $L^2$  risk, i.e., to solve

$$\arg \min_f \mathcal{R}(f) := \mathbb{E}_{(X_1, \dots, X_n) \sim \Pi, Y} [(f(X_1, \dots, X_n) - Y)^2].$$

Under the assumption that  $\mathbb{E}[\varepsilon|X_1, \dots, X_n] = 0$ , the underlying regression function  $f^*$  is the optimal solution, that is, the global minimizer of  $\mathcal{R}(f)$ . However, in practical applications the joint distribution of  $((X_1, \dots, X_n), Y)$  is typically unknown, and only a random sample  $\mathcal{D}_m = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$  with sample size  $m$  is available. Given that each evaluation of the Transformer requires a sequence of length  $n$ , we consider a sliding window training approach. Specifically, we first group the observations into overlapping sequences, each of length  $n$ , to construct new sequences

$$\{((\mathbf{x}_1, \dots, \mathbf{x}_n), y_n), ((\mathbf{x}_2, \dots, \mathbf{x}_{n+1}), y_{n+1}), \dots, ((\mathbf{x}_{m-n+1}, \dots, \mathbf{x}_m), y_m)\},$$

and then estimate the unknown target function  $f^*$  using the empirical risk minimizer

$$\hat{f}_m \in \arg \min_{f \in \mathcal{F}} \mathcal{R}_m(f) := \frac{1}{m-n+1} \sum_{t=n}^m (f(\mathbf{x}_{t-n+1}, \dots, \mathbf{x}_t) - y_t)^2, \quad (3)$$

where we choose the hypothesis class

$$\mathcal{F} = \mathcal{F}(D_m, H_m, S_m, W_m, L_m) = \{\langle \mathcal{N}(X), \mathbf{E} \rangle : \mathcal{N} \in \mathcal{T}_{d_x, d_x}(D_m, H_m, S_m, W_m, L_m)\}.$$

Here,  $\langle \cdot, \cdot \rangle$  denotes the matrix inner product, and  $\mathbf{E} \in \mathbb{R}^{d_x \times n}$  is an arbitrary weight matrix. The performance of the estimator is evaluated by the excess risk, defined as the difference between the  $L^2$  risks of  $\hat{f}_m$  and  $f^*$ , given by

$$\mathcal{R}(\hat{f}_m) - \mathcal{R}(f^*) = \mathbb{E}_{(X_1, \dots, X_n)} [(\hat{f}_m(X_1, \dots, X_n) - f^*(X_1, \dots, X_n))^2].$$

To control the sample complexity, we require that the hypothesis class is uniformly bounded. We define the truncation operator  $\mathcal{C}_B$  with level  $B > 0$  for a real-valued function  $f$  as

$$\mathcal{C}_B f(x) := \begin{cases} f(x) & \text{if } |f(x)| \leq B, \\ \text{sgn}(f(x)) \cdot B & \text{if } |f(x)| > B. \end{cases}$$

For a class of real-valued functions  $\mathcal{F}$ , we use the notation  $\mathcal{C}_B \mathcal{F} := \{\mathcal{C}_B f : f \in \mathcal{F}\}$ . Note that the truncation can be implemented by a feed-forward layer that applies the operation  $\sigma_R[x] - \sigma_R[-x] - \sigma_R[x-B] + \sigma_R[-x-B]$  element-wise. Our next theorem provides a convergence rate for estimating the target function  $f^*$  using the truncated empirical risk minimizer  $\mathcal{C}_{B_m} \hat{f}_m$ . We defer the proof to Section 4.2.

**Theorem 3.** *Under model (2), assume that the regression function  $f^* \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$  for some  $\gamma \in (0, 1]$  and  $K_{\mathcal{H}} > 0$ , and that the probability measure of the covariate  $\Pi$  is supported on  $[0, 1]^{d_x \times n}$  and is absolutely continuous with respect to the Lebesgue measure, with its density function uniformly bounded by a finite constant. Let  $\hat{f}_m$  be the empirical risk minimizer defined in (3) over a random sample  $\mathcal{D}_m = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ . Then, the following excess risk bounds hold:*

- If  $\{\mathbf{x}_i\}_{i=1}^m$  is a geometrically  $\beta$ -mixing sequence, i.e.,  $\beta(k) \leq \beta_0 \exp(-\beta_1 k^r)$  for some  $r, \beta_0, \beta_1 > 0$ , then by choosing  $B_m \asymp \log m$  and the hypothesis class

$$\mathcal{F}(D_m \lesssim 1, H_m \lesssim 1, S_m \lesssim 1, W_m \lesssim m^{\frac{d_x n}{2\gamma + 2d_x n}}, L_m \lesssim 1),$$

we have

$$\mathbb{E}_{\mathcal{D}_m} [\mathcal{R}(\mathcal{C}_{B_m} \hat{f}_m) - \mathcal{R}(f^*)] \lesssim m^{-\frac{\gamma}{\gamma + d_x n}} (\log m)^{3+1/r}.$$

- If  $\{\mathbf{x}_i\}_{i=1}^m$  is an algebraically  $\beta$ -mixing sequence, i.e.,  $\beta(k) \leq \beta_0/k^r$  for some  $r, \beta_0 > 0$ , then by choosing  $B_m \asymp \log m$  and the hypothesis class

$$\mathcal{F}(D_m \lesssim 1, H_m \lesssim 1, S_m \lesssim 1, W_m \lesssim m^{\frac{r d_x n}{2(\gamma+2)\gamma + 2(r+1)d_x n}}, L_m \lesssim 1),$$

we have

$$\mathbb{E}_{\mathcal{D}_m} [\mathcal{R}(\mathcal{C}_{B_m} \hat{f}_m) - \mathcal{R}(f^*)] \lesssim m^{-\frac{r\gamma}{(r+2)\gamma + (r+1)d_x n}} (\log m)^3.$$

Table 2: Comparison of convergence rates.

Reference	Hypothesis Class	Dependence Assumption	Convergence Rate <sup>1</sup>
[20]	FNN	geometrically $\beta$ -mixing	$\tilde{\mathcal{O}}(m^{-\frac{\gamma}{2\gamma+2d+2}})$
[49]	FNN	geometrically $\beta$ -mixing	$\tilde{\mathcal{O}}(m^{-\frac{2\gamma}{2\gamma+d}})$
[32]	RNN	geometrically $\beta$ -mixing	$\tilde{\mathcal{O}}(m^{-\frac{2\gamma}{2\gamma+d_x n}})$
		algebraically $\beta$ -mixing i.i.d.	$\tilde{\mathcal{O}}(m^{-\frac{2r\gamma}{(2r+4)\gamma+(r+1)d_x n}})$ $\tilde{\mathcal{O}}(m^{-\frac{2\gamma}{2\gamma+d_x n}})$
<b>Ours</b> (Theorem 3)	Transformer	geometrically $\beta$ -mixing	$\tilde{\mathcal{O}}(m^{-\frac{\gamma}{\gamma+d_x n}})$
		algebraically $\beta$ -mixing	$\tilde{\mathcal{O}}(m^{-\frac{r\gamma}{(r+2)\gamma+(r+1)d_x n}})$
		i.i.d.	$\tilde{\mathcal{O}}(m^{-\frac{\gamma}{\gamma+d_x n}})$

<sup>1</sup>We omit constants independent of  $m$  and logarithmic factors in  $m$ .  $d$  and  $d_x n$  denote the input dimensions for vector and sequence inputs, respectively.

- If  $\{\mathbf{x}_i\}_{i=1}^m$  is a sequence of i.i.d. random variables, then by choosing  $B_m \asymp \log m$  and the hypothesis class

$$\mathcal{F}(D_m \lesssim 1, H_m \lesssim 1, S_m \lesssim 1, W_m \lesssim m^{\frac{d_x n}{2\gamma+2d_x n}}, L_m \lesssim 1),$$

we have

$$\mathbb{E}_{\mathcal{D}_m}[\mathcal{R}(\mathcal{C}_{B_m} \hat{f}_m) - \mathcal{R}(f^*)] \lesssim m^{-\frac{\gamma}{\gamma+d_x n}} (\log m)^3.$$

It is well known that the optimal convergence rate in nonparametric regression with squared loss for i.i.d. data is  $m^{-2\gamma/(d_x n+2\gamma)}$  [14, 56], and that the same rate remains optimal for certain  $\beta$ -mixing sequences [61, 67]. Therefore, the rates in Theorem 3 are suboptimal. We attribute this suboptimality to the loose upper bound on the VC-dimension (see Lemma 13). We consider the i.i.d. case for an illustration. Classical empirical process techniques yield a decomposition of the excess risk into an approximation error and a statistical error, and by trading off these two errors one obtains the optimal convergence rate, as in [31]. For the approximation error, to approximate a Hölder function with smoothness index  $\gamma$  up to accuracy  $\varepsilon$ , it suffices to use a ReLU FNN with a total number of adjustable parameters  $N \lesssim \varepsilon^{-d/\gamma}$  (up to logarithmic factors), where  $d$  denotes the input dimension and in our setting  $d = d_x n$ . Meanwhile, the VC-dimension, which governs the statistical error, grows linearly with  $N$  (assuming fixed depth) due to the piecewise linear nature of ReLU FNNs. In fact, for FNNs with piecewise polynomial activations (such as  $\text{ReLU}^k$ , see [3, 13]), or for self-attention layers with piecewise polynomial activation functions (e.g., by replacing the softmax  $\sigma_S[\mathbf{Z}]$  with the hardmax  $\sigma_H[\mathbf{Z}]$  [34] or  $\mathbf{Z} \odot \sigma_H[\mathbf{Z}]$  [21, 30]), the VC-dimension grows linearly in the total number of parameters  $N$ . However, for function classes involving exponential operations, such as those defined by sigmoid networks or radial basis function networks, the best known upper bounds on the VC-dimension grow quadratically in  $N$  [2, 35]. We use the same method to establish an upper bound on the VC-dimension of Transformer networks, and hence it also exhibits quadratic growth in  $N$ . As noted in [4], there is a gap between the best known upper and lower bounds for function classes that involve exponential operations, and it remains open whether these bounds are optimal. [35] conjectured that the upper bounds could be improved. To prove Theorem 3, we decompose the excess risk into the approximation error and the statistical error. By Theorem 1 and (1), to achieve an approximation error of at most  $\varepsilon$ , it suffices to use a Transformer network with total parameters  $N \lesssim \varepsilon^{-d_x n/\gamma}$ . However, since the VC-dimension scales as  $N^2$ , it grows faster

than in the ReLU FNN case, leading to the suboptimal rate after trade-off. We leave possible improvements to this gap as an open problem for future study.

We observe that, ignoring logarithmic factors, the convergence rates for the geometrically  $\beta$ -mixing and i.i.d. cases are identical. In addition, for the algebraically  $\beta$ -mixing case the convergence rate is given by  $m^{-\frac{r\gamma}{(r+2)\gamma+(r+1)d_x n}}$ , which improves as the mixing parameter  $r$  increases. When  $r$  is sufficiently large, this rate approaches  $m^{-\frac{\gamma}{\gamma+d_x n}}$ , matching that of the geometrically  $\beta$ -mixing and i.i.d. cases. This observation is consistent with the findings in [32]. We remark that a completely analogous theorem holds for estimating a Sobolev target function  $f^* \in \mathcal{W}^{1,p}([0,1]^{d_x \times n}, K_{\mathcal{W}})$  for some  $p \geq 2$  and  $K_{\mathcal{W}} > 0$ . We summarize the related results in Table 2.

### 2.3 Approximation by Generalized Transformer Networks

We begin by introducing generalized Transformer networks, which extend the original definitions of Transformer layers to enable more flexible functional representations.

**Generalized feed-forward layer:** We define the generalized feed-forward layer as

$$\mathcal{F}_l^{(GFF)}(\mathbf{Z}) := \mathbf{Z} + \mathbf{W}_l^{(2)} \sigma_R \left[ \mathbf{W}_l^{(1)} \mathbf{Z} + \mathbf{B}_l^{(1)} \right] + \mathbf{B}_l^{(2)} \in \mathbb{R}^{D \times n},$$

where  $\mathbf{W}_l^{(1)} \in \mathbb{R}^{W \times D}$  and  $\mathbf{W}_l^{(2)} \in \mathbb{R}^{D \times W}$  are weight matrices, and  $\mathbf{B}_l^{(1)} \in \mathbb{R}^{W \times n}$ ,  $\mathbf{B}_l^{(2)} \in \mathbb{R}^{D \times n}$  are bias matrices. In contrast to the standard feed-forward layer, we allow different bias terms for each column, thereby generalizing the original formulation.

**Generalized self-attention layer:** We define the generalized self-attention layer as

$$\mathcal{F}_l^{(GSA)}(\mathbf{Z}) := \mathbf{Z} + \sum_{h=1}^H \mathbf{W}_{h,l}^{(O)} \sigma_G[\mathbf{Z}] \in \mathbb{R}^{D \times n},$$

where  $\mathbf{W}_{h,l}^{(O)} \in \mathbb{R}^{D \times D}$  is a weight matrix with rank at most  $S$  for all  $h$  and  $l$ , and  $\sigma_G[\mathbf{Z}]$  is a general function (may vary across different  $h$  and  $l$ ) with the only requirement that, for a particular parameter choice, it computes the column average of  $\mathbf{Z}$ , namely,

$$\sigma_G[\mathbf{Z}] = \left( \frac{1}{n} \sum_{j=1}^n \mathbf{Z}_{:,j}, \dots, \frac{1}{n} \sum_{j=1}^n \mathbf{Z}_{:,j} \right).$$

Clearly, both softmax-based self-attention [60]

$$\sigma_G[\mathbf{Z}] = \mathbf{Z} \cdot \sigma_S \left[ \left( \mathbf{W}^{(K)} \mathbf{Z} \right)^\top \left( \mathbf{W}^{(Q)} \mathbf{Z} \right) \right]$$

and (averaging) hardmax-based self-attention [47]

$$\sigma_G[\mathbf{Z}] = \mathbf{Z} \cdot \sigma_H \left[ \left( \mathbf{W}^{(K)} \mathbf{Z} \right)^\top \left( \mathbf{W}^{(Q)} \mathbf{Z} \right) \right]$$

satisfy the above definition, since they compute the column average of  $\mathbf{Z}$  when  $\mathbf{W}^{(K)} = \mathbf{W}^{(Q)} = \mathbf{O}$ .

The class of generalized Transformer networks is then defined as

$$\mathcal{GT}_{d_x, d_y}(D, H, S, W, L) := \left\{ \mathcal{E}_{out} \circ \mathcal{F}_L^{(GFF)} \circ \mathcal{F}_L^{(GSA)} \circ \dots \circ \mathcal{F}_1^{(GFF)} \circ \mathcal{F}_1^{(GSA)} \circ \mathcal{E}_{in} \right\}.$$

A central challenge in understanding the expressivity of Transformers lies in explaining why the self-attention layer effectively captures complex token-wise interactions. Indeed, the self-attention mechanism is the only component within Transformer architectures explicitly designed

to integrate token-level information, thus playing a central role in modeling dependencies across sequences. However, the highly nonlinear softmax function commonly used in self-attention presents significant analytical difficulties. Previous research has approached this problem from several perspectives: some studies first analyzed the simpler hardmax function, viewing softmax attention as a smoother approximation whose limiting behavior converges to hardmax [58, 70]; some restricted the analysis to finite discrete sequences and showed by construction that softmax function acts as a contextual mapping, sending distinct input sequences to distinct output sequences [33, 36, 69]; yet another approach replaced softmax with more analytically tractable functions, such as piecewise polynomial activation functions [21, 22, 30]. In contrast, we provide a fundamentally different proof strategy inspired by the Kolmogorov-Arnold representation theorem (see Proposition 14). Our key observation is that representing an arbitrary  $d_x n$ -dimensional function exactly requires only one inner and one outer function. The inner function in this construction completely separates each entry of the input sequence, effectively simplifying complex token interactions into structured summations. Consequently, we show that the essential role of the self-attention layer can be simplified to performing column-wise summations, providing a more direct and general theoretical justification for the expressive power of Transformer architectures.

The following theorem provides explicit approximation bounds for Hölder continuous functions using generalized Transformer networks. We defer the proof to Section 4.3.

**Theorem 4.** *Given  $\gamma \in (0, 1]$  and  $K_{\mathcal{H}} > 0$ , assume that the target function  $\mathbf{F} : [0, 1]^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  satisfies  $F_{i,j} \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$  for each  $i \in [d_x], j \in [n]$ . For any  $\varepsilon \in (0, 1)$  and  $p \in [1, \infty)$ , there exists a generalized Transformer network*

$$\mathcal{N} \in \mathcal{GT}_{d_x, d_x}(D = 4d_x n, H = 1, S = d_x, W = 3d_x n \cdot \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil, L = 6 \lceil \frac{1}{\gamma} \log_2 \frac{1}{\varepsilon} \rceil)$$

such that

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0, 1]^{d_x \times n})} \leq 4(d_x n)^3 K_{\mathcal{H}} \varepsilon.$$

### 3 Discussions and Related Works

**Nonparametric regression using neural networks.** The convergence rates of neural network regression estimators have been extensively analyzed in the literature. Minimax optimal rates have been established across various neural network architectures, including under-parameterized sparse deep FNNs [50, 57], under-parameterized fully connected deep FNNs [31], over-parameterized shallow FNNs [64, 65], and RNNs [32]. In contrast, convergence rates for Transformer-based estimators have rarely been observed. [58] investigated the approximation and estimation capabilities of Transformers as sequence-to-sequence functions operating on infinite-dimensional inputs, where variable-length sliding window attention was considered. Additionally, modifications to the Transformer architecture have been explored, such as replacing the standard softmax function  $\sigma_S[\mathbf{Z}]$  in the self-attention layer by  $\mathbf{Z} \odot \sigma_H[\mathbf{Z}]$  [21] and by  $\sigma_R[\mathbf{Z}]$  [22]. Although these studies provide insightful constructions and analyses, their avoidance of the standard softmax function does not fully explain the successes observed ever since the introduction of the Transformer mechanism [60]. Our result (Theorem 3) directly addresses this gap by analyzing convergence rates for the standard Transformer architecture explicitly using the original softmax attention mechanism. It has also been shown that neural networks are able to circumvent the curse of dimensionality under certain conditions, for example, when the intrinsic dimension of the regression function is low [6, 22, 31, 43], or the regression function has certain hierarchical structures [37, 50]. Exploring the conditions under which Transformers similarly mitigate the curse of dimensionality within our framework is an important direction for future research.

**Assumptions on the smoothness of target function.** In this work, we require the target function to be either Hölder continuous with smoothness index  $\gamma \in (0, 1]$  or a Sobolev function with bounded first-order weak derivative. It remains an interesting problem whether our methods are adaptive to higher regularity. In the proof of Theorems 1 and 2, we approximate the target function by piecewise constant functions defined on a uniform partition of  $[0, 1]^{d_x \times n}$  into  $K^{d_x n}$  cells. The approximation orders  $K^{-\gamma}$  in (4) and  $K^{-1}$  in (10) cannot be improved in general. We refer to [11, Section 6.2] for a detailed discussion of saturation and inverse theorems, where certain smoothness properties of a function are deduced from the order of its approximation by multivariate piecewise polynomials. For example, consider a real-valued function  $f$  defined on a bounded domain  $\Omega \subseteq \mathbb{R}^d$ , and let  $\Delta$  be a partition of  $\Omega$  into a finite number of subsets. Suppose  $f$  is approximated by a piecewise constant function

$$s(\mathbf{x}) = \sum_{\omega \in \Delta} c_\omega \mathbb{1}_\omega(\mathbf{x}),$$

and assume that  $\inf_s \|f - s\|_{L^\infty(\Omega)} = o(\text{diam}(\Delta))$  as  $\text{diam}(\Delta) := \max_{\omega \in \Delta} \text{diam}(\omega) \rightarrow 0$  for all partitions  $\Delta$ . We can then easily show that  $f$  is a constant function. Indeed, for any  $\mathbf{x}, \mathbf{y} \in \Omega$ , there exists a partition  $\Delta$  such that  $\mathbf{x}$  and  $\mathbf{y}$  belong to the same cell  $\omega$ , where  $\text{diam}(\omega) = \text{diam}(\Delta) \leq 2\|\mathbf{x} - \mathbf{y}\|_2$ . Let  $\bar{s}$  be a best piecewise constant approximation to  $f$  under the  $L^\infty$ -norm (the existence of which is ensured by a compactness argument). Then  $|f(\mathbf{x}) - f(\mathbf{y})| \leq |f(\mathbf{x}) - \bar{s}(\mathbf{x})| + |\bar{s}(\mathbf{x}) - \bar{s}(\mathbf{y})| + |\bar{s}(\mathbf{y}) - f(\mathbf{y})| \leq 2 \inf_s \|f - s\|_{L^\infty(\Omega)}$  since  $\bar{s}$  is constant on  $\omega$ . Hence, by assumption,  $|f(\mathbf{x}) - f(\mathbf{y})| = o(\text{diam}(\Delta)) = o(\|\mathbf{x} - \mathbf{y}\|_2)$  as  $\mathbf{y} \rightarrow \mathbf{x}$ , which implies that  $f$  has zero derivative at every point in  $\Omega$ , i.e.,  $f$  is constant. For the uniform partition used in our proof, a more subtle analysis is required; see [12, Chapter 12.2].

In the proof of Theorem 4, we use the Kolmogorov-Arnold representation  $f(x_1, \dots, x_d) = g(3 \sum_{p=1}^d 3^{-p} \phi(x_p))$  for any  $d$ -variate function  $f$  [51]. Although this representation allows the transfer of smoothness properties of  $f$  to the function  $g$  for Hölder continuous  $f$  with smoothness index  $\gamma \in (0, 1]$ , it remains unclear whether this representation can be generalized to higher order smoothness or anisotropic smoothness.

## 4 Proofs

Before proceeding, we clarify some simplifications used in the proof.

1. In a self-attention layer, if we set  $\mathbf{W}^{(O)} = \mathbf{O}$ , the layer behaves as an identity mapping due to the presence of the skip connection. Similarly, in a feed-forward layer, setting  $\mathbf{W}^{(2)} = \mathbf{O}$  causes the layer to degenerate into an identity mapping. Therefore, as long as identity mappings are appropriately introduced, the composition of multiple self-attention layers or feed-forward layers remains consistent with our definition of a Transformer.
2. Since a feed-forward layer applies the same operation to each column of the input matrix, we do not distinguish between matrix and vector inputs when the context is clear, with a slight abuse of notation. For example, given a feed-forward layer  $\mathcal{F}^{(FF)} : \mathbb{R}^{D \times n} \rightarrow \mathbb{R}^{D \times n}$  defined as

$$\mathcal{F}^{(FF)}(\mathbf{H}) = \mathbf{H} + \mathbf{W}^{(2)} \sigma_R \left[ \mathbf{W}^{(1)} \mathbf{H} + \mathbf{b}^{(1)} \mathbf{1}_n^\top \right] + \mathbf{b}^{(2)} \mathbf{1}_n^\top,$$

we also define  $\mathcal{F}^{(FF)} : \mathbb{R}^D \rightarrow \mathbb{R}^D$  as

$$\mathcal{F}^{(FF)}(\mathbf{H}_{:,i}) = \mathbf{H}_{:,i} + \mathbf{W}^{(2)} \sigma_R \left[ \mathbf{W}^{(1)} \mathbf{H}_{:,i} + \mathbf{b}^{(1)} \right] + \mathbf{b}^{(2)}$$

for each  $i$ , so that  $\mathcal{F}^{(FF)}(\mathbf{H}) = (\mathcal{F}^{(FF)}(\mathbf{H}_{:,1}), \dots, \mathcal{F}^{(FF)}(\mathbf{H}_{:,n}))$ .

If all self-attention layers degenerate into identity mappings, the resulting Transformer reduces to a token-wise ResNet [23]. We will demonstrate that any token-wise FNN can be represented by a token-wise ResNet, thereby naturally extending existing results on FNNs to Transformers. Specifically, an FNN  $\mathcal{N} : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$  is a function that can be parameterized in the form

$$\begin{aligned}\mathcal{N}_0(\mathbf{x}) &= \mathbf{x}, \\ \mathcal{N}_{l+1}(\mathbf{x}) &= \sigma_R[\mathbf{A}_l \mathcal{N}_l(\mathbf{x}) + \mathbf{b}_l], \quad l = 0, \dots, L-1, \\ \mathcal{N}(\mathbf{x}) &= \mathbf{A}_L \mathcal{N}_L(\mathbf{x}) + \mathbf{b}_L,\end{aligned}$$

where  $\mathbf{A}_l \in \mathbb{R}^{W_{l+1} \times W_l}$ ,  $\mathbf{b}_l \in \mathbb{R}^{W_{l+1}}$  with  $W_0 = d_x$ ,  $W_{L+1} = d_y$  and  $W_l = W$  for  $l = 1, \dots, L$ . The parameters  $W$  and  $L$  are referred to as the width and depth of the neural network, respectively. We denote by  $\mathcal{FNN}_{d_x, d_y}(W, L)$  the set of functions that can be parameterized in this form with width  $W$  and depth  $L$ .

**Lemma 5.** *Let  $d_x, d_y$  be positive integers. For any  $\mathcal{N} \in \mathcal{FNN}_{d_x, d_y}(W, L)$  with width  $W \geq \max\{d_x, d_y\}$  and depth  $L \geq 2$ , there exist an embedding map  $\mathcal{E}_{in} : \mathbf{X} \in \mathbb{R}^{d_x \times n} \mapsto \begin{pmatrix} \mathbf{X} \\ \mathbf{O} \end{pmatrix} \in \mathbb{R}^{W \times n}$ , a projection map  $\mathcal{E}_{out} : \begin{pmatrix} \mathbf{Y} \\ \mathbf{O} \end{pmatrix} \in \mathbb{R}^{W \times n} \mapsto \mathbf{Y} \in \mathbb{R}^{d_y \times n}$ , and  $L$  feed-forward layers with width at most  $3W$ , such that for any  $\mathbf{X} \in \mathbb{R}^{d_x \times n}$ ,*

$$\mathcal{E}_{out} \circ \mathcal{F}_L^{(FF)} \circ \dots \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in}(\mathbf{X}) = (\mathcal{N}(\mathbf{X}_{:,1}), \dots, \mathcal{N}(\mathbf{X}_{:,n})) \in \mathbb{R}^{d_y \times n}.$$

*Proof.* The idea is to use the identity  $\sigma_R[x] - \sigma_R[-x] = x$  to eliminate the skip connection. For any  $\mathbf{x} \in \mathbb{R}^{d_x}$ , direct computation yields

$$\begin{aligned}\mathcal{F}_1^{(FF)} \begin{pmatrix} \mathbf{x} \\ \mathbf{0} \end{pmatrix} &= \begin{pmatrix} \mathbf{x} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{I}_{d_x} & \mathbf{O} & -\mathbf{I}_{d_x} & \mathbf{I}_{d_x} \\ \mathbf{O} & \mathbf{I}_{W-d_x} & \mathbf{O} & \mathbf{O} \end{pmatrix} \sigma_R \left[ \begin{pmatrix} \mathbf{A}_0 & \mathbf{O} \\ \mathbf{I}_{d_x} & \mathbf{O} \\ -\mathbf{I}_{d_x} & \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{b}_0 \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix} \right] \\ &= \sigma_R[\mathbf{A}_0 \mathbf{x} + \mathbf{b}_0] + \begin{pmatrix} \mathbf{x} + -\sigma_R[\mathbf{x}] + \sigma_R[-\mathbf{x}] \\ \mathbf{0} \end{pmatrix} \\ &= \mathcal{N}_1(\mathbf{x}),\end{aligned}$$

where we have used the identity  $\sigma_R[\mathbf{x}] - \sigma_R[-\mathbf{x}] = \mathbf{x}$ .

Now, assuming that  $\mathcal{F}_l^{(FF)} \circ \dots \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in}(\mathbf{x}) = \mathcal{N}_l(\mathbf{x})$ , we define the  $(l+1)$ -th feed-forward layer  $\mathcal{F}_{l+1}^{(FF)}$  as

$$\begin{aligned}\mathcal{F}_{l+1}^{(FF)}(\mathcal{N}_l(\mathbf{x})) &= \mathcal{N}_l(\mathbf{x}) + (\mathbf{I}_W, -\mathbf{I}_W, \mathbf{I}_W) \sigma_R \left[ \begin{pmatrix} \mathbf{A}_l \\ \mathbf{I}_W \\ -\mathbf{I}_W \end{pmatrix} \mathcal{N}_l(\mathbf{x}) + \begin{pmatrix} \mathbf{b}_l \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix} \right] \\ &= \sigma_R[\mathbf{A}_l \mathcal{N}_l(\mathbf{x}) + \mathbf{b}_l] = \mathcal{N}_{l+1}(\mathbf{x}).\end{aligned}$$

By induction, it follows that

$$\mathcal{F}_{l+1}^{(FF)} \circ \mathcal{F}_l^{(FF)} \circ \dots \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in}(\mathbf{x}) = \mathcal{N}_{l+1}(\mathbf{x}).$$

By the principle of induction, we establish that  $\mathcal{F}_{L-1}^{(FF)} \circ \dots \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in}(\mathbf{x}) = \mathcal{N}_{L-1}(\mathbf{x})$ . For the last feed-forward layer, we calculate that

$$\begin{aligned}\mathcal{F}_L^{(FF)}(\mathcal{N}_{L-1}(\mathbf{x})) &= \mathcal{N}_{L-1}(\mathbf{x}) + \begin{pmatrix} \mathbf{A}_L & -\mathbf{I}_{d_y} & \mathbf{O} & \mathbf{I}_{d_y} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & -\mathbf{I}_{W-d_y} & \mathbf{O} & \mathbf{I}_{W-d_y} \end{pmatrix} \sigma_R \left[ \begin{pmatrix} \mathbf{A}_{L-1} \\ \mathbf{I}_W \\ -\mathbf{I}_W \end{pmatrix} \mathcal{N}_{L-1}(\mathbf{x}) + \begin{pmatrix} \mathbf{b}_{L-1} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix} \right] + \begin{pmatrix} \mathbf{b}_L \\ \mathbf{0} \end{pmatrix}\end{aligned}$$

$$\begin{aligned}
&= \begin{pmatrix} \mathbf{A}_L \sigma_R[\mathbf{A}_{L-1} \mathcal{N}_{L-1}(\mathbf{x}) + \mathbf{b}_{L-1}] + \mathbf{b}_L \\ \mathbf{0} \end{pmatrix} \\
&= \begin{pmatrix} \mathcal{N}(\mathbf{x}) \\ \mathbf{0} \end{pmatrix}.
\end{aligned}$$

Thus, we obtain

$$\mathcal{E}_{out} \circ \mathcal{F}_L^{(FF)} \circ \dots \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in}(\mathbf{x}) = \mathcal{N}(\mathbf{x}).$$

Since each feed-forward layer in our construction has width at most  $3W$ , the proof is complete by considering  $\mathbf{x} = \mathbf{X}_{:,i}$  for  $i \in [n]$ .  $\square$

The following proposition gives basic properties of Transformer networks that enable the recursive construction of complex architectures.

**Proposition 6.** *Let  $\mathcal{N}_i \in \mathcal{T}_{d_i, k_i}(D_i, H_i, S_i, W_i, L_i)$  for  $i = 1, 2$ .*

1. *If  $d_1 = d_2$ ,  $k_1 = k_2$ , and  $D_1 \leq D_2, H_1 \leq H_2, S_1 \leq S_2, W_1 \leq W_2, L_1 \leq L_2$ , then*

$$\mathcal{T}_{d_1, k_1}(D_1, H_1, S_1, W_1, L_1) \subseteq \mathcal{T}_{d_2, k_2}(D_2, H_2, S_2, W_2, L_2).$$

2. *(Concatenation) If define  $\mathcal{N} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathcal{N}_1(\mathbf{X}) \\ \mathcal{N}_2(\mathbf{Y}) \end{pmatrix}$ , then*

$$\mathcal{N} \in \mathcal{T}_{d_1+d_2, k_1+k_2}(D_1 + D_2, H_1 + H_2, \max\{S_1, S_2\}, W_1 + W_2, \max\{L_1, L_2\}).$$

3. *(Summation) If  $d_1 = d_2$  and  $k_1 = k_2$ , then*

$$\mathcal{N}_1 + \mathcal{N}_2 \in \mathcal{T}_{d_1, k_1}(D_1 + D_2, H_1 + H_2, \max\{S_1, S_2\}, W_1 + W_2, \max\{L_1, L_2\}).$$

*Proof.* We provide the proof for (2), as the arguments for (1) and (3) follow analogously.

Let

$$\mathcal{N}_i = \mathcal{E}_{i,out} \circ \mathcal{F}_{i,L_i}^{(FF)} \circ \mathcal{F}_{i,L_i}^{(SA)} \circ \dots \circ \mathcal{F}_{i,1}^{(FF)} \circ \mathcal{F}_{i,1}^{(SA)} \circ \mathcal{E}_{i,in} \in \mathcal{T}_{d_i, k_i}(D_i, H_i, S_i, W_i, L_i), \quad i = 1, 2.$$

Without loss of generality, assume that  $L_1 = L_2 = L$ , since the identity mapping can be viewed as a special self-attention layer or a feed-forward layer. We define the following components:

1. Input Embedding:

$$\mathcal{E}_{in} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{E}_{1,in} & \\ & \mathbf{E}_{2,in} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} + \begin{pmatrix} \mathbf{P}_1 \\ \mathbf{P}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{E}_{1,in} \mathbf{X} + \mathbf{P}_1 \\ \mathbf{E}_{2,in} \mathbf{Y} + \mathbf{P}_2 \end{pmatrix} = \begin{pmatrix} \mathcal{E}_{1,in}(\mathbf{X}) \\ \mathcal{E}_{2,in}(\mathbf{Y}) \end{pmatrix}$$

2. Feed-forward Layer:

$$\begin{aligned}
&\mathcal{F}_l^{(FF)} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} + \begin{pmatrix} \mathbf{W}_{1,l}^{(2)} & \\ & \mathbf{W}_{2,l}^{(2)} \end{pmatrix} \sigma_R \left[ \begin{pmatrix} \mathbf{W}_{1,l}^{(1)} & \\ & \mathbf{W}_{2,l}^{(1)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} + \begin{pmatrix} \mathbf{b}_{1,l}^{(1)} \\ \mathbf{b}_{2,l}^{(1)} \end{pmatrix} \mathbf{1}_n^\top \right] + \begin{pmatrix} \mathbf{b}_{1,l}^{(2)} \\ \mathbf{b}_{2,l}^{(2)} \end{pmatrix} \mathbf{1}_n^\top \\
&= \begin{pmatrix} \mathbf{X} + \mathbf{W}_{1,l}^{(2)} \sigma_R[\mathbf{W}_{1,l}^{(1)} \mathbf{X} + \mathbf{b}_{1,l}^{(1)} \mathbf{1}_n^\top] + \mathbf{b}_{1,l}^{(2)} \mathbf{1}_n^\top \\ \mathbf{Y} + \mathbf{W}_{2,l}^{(2)} \sigma_R[\mathbf{W}_{2,l}^{(1)} \mathbf{Y} + \mathbf{b}_{2,l}^{(1)} \mathbf{1}_n^\top] + \mathbf{b}_{2,l}^{(2)} \mathbf{1}_n^\top \end{pmatrix} \\
&= \begin{pmatrix} \mathcal{F}_{1,l}^{(FF)}(\mathbf{X}) \\ \mathcal{F}_{2,l}^{(FF)}(\mathbf{Y}) \end{pmatrix}
\end{aligned}$$

3. Self-Attention Layer:

$$\begin{aligned}
& \mathcal{F}_l^{(SA)} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \\
&= \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} + \sum_{h=1}^{H_1} \begin{pmatrix} \mathbf{W}_{1,h,l}^{(O)} \\ \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{W}_{1,h,l}^{(V)} \\ \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sigma_S \left[ \left( \begin{pmatrix} \mathbf{W}_{1,h,l}^{(K)} \\ \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \right)^\top \begin{pmatrix} \mathbf{W}_{1,h,l}^{(Q)} \\ \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \right] \\
&+ \sum_{h=1}^{H_2} \begin{pmatrix} \mathbf{O} \\ \mathbf{W}_{2,h,l}^{(O)} \end{pmatrix} \begin{pmatrix} \mathbf{O}, \mathbf{W}_{2,h,l}^{(V)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sigma_S \left[ \left( \begin{pmatrix} \mathbf{O}, \mathbf{W}_{2,h,l}^{(K)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \right)^\top \begin{pmatrix} \mathbf{O}, \mathbf{W}_{2,h,l}^{(Q)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \right] \\
&= \begin{pmatrix} \mathbf{X} + \sum_{h=1}^{H_1} \mathbf{W}_{1,h,l}^{(O)} \begin{pmatrix} \mathbf{W}_{1,h,l}^{(V)} \\ \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sigma_S \left[ \begin{pmatrix} \mathbf{W}_{1,h,l}^{(K)} \\ \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \right]^\top \begin{pmatrix} \mathbf{W}_{1,h,l}^{(Q)} \\ \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \\ \mathbf{Y} + \sum_{h=1}^{H_2} \begin{pmatrix} \mathbf{O} \\ \mathbf{W}_{2,h,l}^{(O)} \end{pmatrix} \begin{pmatrix} \mathbf{O}, \mathbf{W}_{2,h,l}^{(V)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sigma_S \left[ \begin{pmatrix} \mathbf{O}, \mathbf{W}_{2,h,l}^{(K)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \right]^\top \begin{pmatrix} \mathbf{O}, \mathbf{W}_{2,h,l}^{(Q)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \end{pmatrix} \\
&= \begin{pmatrix} \mathcal{F}_{1,l}^{(SA)}(\mathbf{X}) \\ \mathcal{F}_{2,l}^{(SA)}(\mathbf{Y}) \end{pmatrix}
\end{aligned}$$

4. Output Projection:

$$\mathcal{E}_{out} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{E}_{1,out} \\ \mathbf{E}_{2,out} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{E}_{1,out} \mathbf{X} \\ \mathbf{E}_{2,out} \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathcal{E}_{1,out}(\mathbf{X}) \\ \mathcal{E}_{2,out}(\mathbf{Y}) \end{pmatrix}$$

By direct verification, we obtain

$$\mathcal{N} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} := \mathcal{E}_{out} \circ \mathcal{F}_L^{(FF)} \circ \mathcal{F}_L^{(SA)} \circ \dots \circ \mathcal{F}_1^{(FF)} \circ \mathcal{F}_1^{(SA)} \circ \mathcal{E}_{in} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathcal{N}_1(\mathbf{X}) \\ \mathcal{N}_2(\mathbf{Y}) \end{pmatrix}.$$

Furthermore, it follows that  $\mathcal{N} \in \mathcal{T}_{d_1+d_2, k_1+k_2}(D_1 + D_2, H_1 + H_2, \max\{S_1, S_2\}, W_1 + W_2, L)$ , thus completing the proof.  $\square$

#### 4.1 Proof of Theorems 1 and 2

Given  $K \in \mathbb{N}$  and  $\delta \in (0, \frac{1}{K})$ , we define a trifling region  $\Omega([0, 1]^{D \times n}, K, \delta) \subseteq [0, 1]^{D \times n}$  as

$$\Omega([0, 1]^{D \times n}, K, \delta) := \left\{ \mathbf{X} \in [0, 1]^{D \times n} : \exists X_{i,j} \in \cup_{t=1}^{K-1} \left( \frac{t}{K}, \frac{t}{K} + \delta \right) \right\}.$$

The introduction of the trifling region serves to identify the "bad" areas where mismatches occur when approximating a discontinuous multi-step function using continuous piecewise linear functions, which can be implemented by a feed-forward layer. Since the trifling region has arbitrarily small Lebesgue measure, we focus on function approximation in the "good" region, namely, the complement domain  $[0, 1]^{D \times n} \setminus \Omega([0, 1]^{D \times n}, K, \delta)$ .

**Proposition 7.** *Given  $\gamma \in (0, 1]$  and  $K_{\mathcal{H}} > 0$ , assume that  $\mathbf{F} : [0, 1]^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  satisfies  $F_{i,j} \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$  for each  $i \in [d_x], j \in [n]$ . For any  $K \in \mathbb{N}$  and  $\delta \in (0, \frac{1}{K})$ , there exists a Transformer network  $\mathcal{N} \in \mathcal{T}_{d_x, d_x}(d_x, 1, 1, 5nK^{d_x n}, 2)$  such that*

1.  $|\mathcal{N}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq K_{\mathcal{H}}(d_x n)^{\gamma/2} K^{-\gamma}$  for any  $i \in [d_x], j \in [n]$  and  $\mathbf{X} \in [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta)$ ,
2.  $\|\mathcal{N}(\mathbf{X})\|_F \leq \sqrt{d_x n} K_{\mathcal{H}}$  for any  $\mathbf{X} \in \mathbb{R}^{d_x \times n}$ .

*Proof.* We basically follow the proof of [33, Proposition 1].

**Step 1:** We begin by uniformly partitioning the domain  $[0, 1]^{d_x \times n}$  into  $K^{d_x n}$  subregions and constructing a piecewise constant function  $\overline{\mathbf{F}}$  that approximates the target function  $\mathbf{F}$ , with an approximation error scales as  $K^{-\gamma}$ . Specifically, let  $K \in \mathbb{N}$  denote the granularity of the grid

$$\mathbb{G}_K = \left\{ \frac{1}{K}, \frac{2}{K}, \dots, 1 \right\}^{d_x \times n}.$$

We define each subregion as

$$\omega_{\mathbf{G}} := \prod_{i \in [d_x], j \in [n]} \begin{cases} [G_{i,j} - \frac{1}{K}, G_{i,j}], & \text{if } G_{i,j} = \frac{1}{K} \\ (G_{i,j} - \frac{1}{K}, G_{i,j}], & \text{otherwise} \end{cases}$$

associated with  $\mathbf{G} \in \mathbb{G}_K$ . Clearly, these subregions form a partition of the domain  $[0, 1]^{d_x \times n} = \bigcup_{\mathbf{G} \in \mathbb{G}_K} \omega_{\mathbf{G}}$ . Given a target function  $\mathbf{F}$  with  $F_{i,j} \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$ , we define a piecewise constant approximation of  $\mathbf{F}$  as

$$\bar{\mathbf{F}}(\mathbf{X}) = \sum_{\mathbf{G} \in \mathbb{G}_K} \mathbf{F}(\mathbf{G}) \mathbb{1}_{\omega_{\mathbf{G}}}(\mathbf{X}),$$

where  $\mathbb{1}_{\omega}$  denotes the indicator function of set  $\omega$ . That is, within each subregion  $\omega_{\mathbf{G}}$ , we approximate  $\mathbf{F}$  using its value at the grid point  $\mathbf{G}$ . By the regularity of  $\mathbf{F}$ , we have the error estimate

$$\begin{aligned} |F_{i,j}(\mathbf{X}) - \bar{F}_{i,j}(\mathbf{X})| &= \left| \sum_{\mathbf{G} \in \mathbb{G}_K} (F_{i,j}(\mathbf{X}) - \bar{F}_{i,j}(\mathbf{X})) \mathbb{1}_{\omega_{\mathbf{G}}}(\mathbf{X}) \right| \\ &= \left| \sum_{\mathbf{G} \in \mathbb{G}_K} (F_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{G})) \mathbb{1}_{\omega_{\mathbf{G}}}(\mathbf{X}) \right| \\ &\leq \sum_{\mathbf{G} \in \mathbb{G}_K} |(F_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{G}))| \mathbb{1}_{\omega_{\mathbf{G}}}(\mathbf{X}) \\ &\leq \sum_{\mathbf{G} \in \mathbb{G}_K} K_{\mathcal{H}} \|\mathbf{X} - \mathbf{G}\|_F^\gamma \mathbb{1}_{\omega_{\mathbf{G}}}(\mathbf{X}) \\ &\leq K_{\mathcal{H}} (d_x n)^{\gamma/2} K^{-\gamma} \sum_{\mathbf{G} \in \mathbb{G}_K} \mathbb{1}_{\omega_{\mathbf{G}}}(\mathbf{X}) \\ &= K_{\mathcal{H}} (d_x n)^{\gamma/2} K^{-\gamma}, \end{aligned} \tag{4}$$

for any  $i \in [d_x]$ ,  $j \in [n]$  and  $\mathbf{X} \in [0, 1]^{d_x \times n}$ .

**Step 2:** Given a positional encoding matrix  $\mathbf{P}$  and a spatial discretization function  $\bar{\mathcal{F}}_1^{(FF)}$  satisfying

$$\bar{\mathcal{F}}_1^{(FF)}(\mathbf{X} + \mathbf{P}) = \mathbf{G} + \mathbf{P}, \quad \text{for all } \mathbf{X} \in \omega_{\mathbf{G}},$$

our objective is to construct a feed-forward layer  $\mathcal{F}_1^{(FF)}$ , with width at most  $2nd_x(K+1)$ , that accurately represents  $\bar{\mathcal{F}}_1^{(FF)}$  outside the trifling region  $\Omega([0, 1]^{d_x \times n}, K, \delta)$ , that is,

$$\mathcal{F}_1^{(FF)}(\mathbf{X} + \mathbf{P}) = \bar{\mathcal{F}}_1^{(FF)}(\mathbf{X} + \mathbf{P}), \quad \text{for any } \mathbf{X} \in [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta).$$

To achieve this goal, we first approximate a univariate multiple-step function using a piecewise linear function, and then extend this function to matrix elements by stacking.

We define the embedding layer as

$$\mathcal{E}_{in}(\mathbf{X}) = \mathbf{X} + \mathbf{P} \in \mathbb{R}^{d_x \times n},$$

where the positional encoding matrix  $\mathbf{P}$  is chosen as

$$\mathbf{P} = \begin{pmatrix} 0 & 2 & \cdots & 2(n-1) \\ \vdots & \vdots & & \vdots \\ 0 & 2 & \cdots & 2(n-1) \end{pmatrix}.$$

Since  $\mathbf{X} \in [0, 1]^{d_x \times n}$ , the positional encoding ensures that the columns of  $\mathbf{X} + \mathbf{P}$  are mapped to distinct intervals, that is,  $[\mathbf{X} + \mathbf{P}]_{i,j} \in [2j - 2, 2j - 1]$  for each  $j \in [n]$ . Now, consider a multiple-step function  $\text{step}_K(z)$  defined on  $[0, 1]$  as

$$\text{step}_K(z) = \begin{cases} \frac{1}{K}, & 0 \leq z \leq \frac{1}{K} \\ \frac{2}{K}, & \frac{1}{K} < z \leq \frac{2}{K} \\ \frac{3}{K}, & \frac{2}{K} < z \leq \frac{3}{K} \\ \vdots & \vdots \\ 1, & 1 - \frac{1}{K} < z \leq 1 \end{cases}.$$

Given  $\delta \in (0, \frac{1}{K})$ , by translations, scalings and summations of the  $\delta$ -approximated step function

$$\sigma_R[z/\delta] - \sigma_R[z/\delta - 1] = \begin{cases} 0 & z \leq 0 \\ z/\delta & 0 < z < \delta \\ 1 & \delta \leq z \end{cases},$$

we define

$$\begin{aligned} f(z) &= \frac{1}{K} + \sum_{j=1}^n \sum_{t=1}^{K-1} \frac{1}{K} \left( \sigma_R \left[ \frac{z - 2(j-1)}{\delta} - \frac{t}{\delta K} \right] - \sigma_R \left[ \frac{z - 2(j-1)}{\delta} - 1 - \frac{t}{\delta K} \right] \right) \\ &\quad + \sum_{j=1}^{n-1} \left( 1 + \frac{1}{K} \right) (\sigma_R[z - (2j-1)] - \sigma_R[z - 2j]). \end{aligned}$$

It is straightforward to verify that  $f(z + (2j-2)) = \text{step}_K(z) + (2j-2)$  for all  $z \in [0, 1] \setminus \Omega([0, 1], K, \delta)$  and  $j \in [n]$ . Moreover, the function  $f$  can be represented by a shallow ReLU network with  $2nK - 2$  units in the hidden layer.

We then concatenate multiple  $f$  in parallel to construct a feed-forward layer  $\mathcal{F}_1^{(FF)} : \mathbb{R}^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  with width at most  $2nd_x(K+1)$ , satisfying

$$\begin{aligned} \mathcal{F}_1^{(FF)}(\mathbf{X} + \mathbf{P}) &= \begin{pmatrix} f(X_{1,1} + P_{1,1}) & \cdots & f(X_{1,n} + P_{1,n}) \\ \vdots & \ddots & \vdots \\ f(X_{d_x,1} + P_{d_x,1}) & \cdots & f(X_{d_x,n} + P_{d_x,n}) \end{pmatrix} \\ &\approx \begin{pmatrix} \text{step}_K(X_{1,1}) + P_{1,1} & \cdots & \text{step}_K(X_{1,n}) + P_{1,n} \\ \vdots & \ddots & \vdots \\ \text{step}_K(X_{d_x,1}) + P_{d_x,1} & \cdots & \text{step}_K(X_{d_x,n}) + P_{d_x,n} \end{pmatrix} \\ &= \begin{pmatrix} \text{step}_K(X_{1,1}) & \cdots & \text{step}_K(X_{1,n}) \\ \vdots & \ddots & \vdots \\ \text{step}_K(X_{d_x,1}) & \cdots & \text{step}_K(X_{d_x,n}) \end{pmatrix} + \mathbf{P}. \end{aligned}$$

Noting that

$$\Omega([0, 1]^{d_x \times n}, K, \delta) = \bigcup_{i \in [d_x], j \in [n]} \{\mathbf{X} : X_{i,j} \in \Omega([0, 1], K, \delta)\},$$

we conclude that  $\mathcal{F}_1^{(FF)}(\mathbf{X} + \mathbf{P}) = \overline{\mathcal{F}}_1^{(FF)}(\mathbf{X} + \mathbf{P})$  for any  $\mathbf{X} \in [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta)$ .

**Step 3:** Since  $\{\mathbf{G} + \mathbf{P} : \mathbf{G} \in \mathbb{G}_K\}$  can be regarded as sequences, each of which has no duplicate token due to positional encoding, it follows from [33, Theorem 2] that there exists a self-attention layer  $\mathcal{F}^{(SA)} : \mathbb{R}^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  with  $H = 1$  and  $s = 1$  that serves as a contextual mapping for such input sequences (see [33, 69] for further discussion). In essence, a contextual

mapping is a bijection between sequences that satisfies  $\mathcal{F}^{(SA)}(\mathbf{G}^{(i)} + \mathbf{P})_{:,k} \neq \mathcal{F}^{(SA)}(\mathbf{G}^{(j)} + \mathbf{P})_{:,l}$  if  $\mathbf{G}^{(i)} \neq \mathbf{G}^{(j)} \in \mathbb{G}_K$  or  $k \neq l \in [n]$ . The remaining is to associate each output token with its corresponding function value using a feed-forward layer, which reduces to a memorization task. Lemma 11 gives a construction of such a feed-forward layer, denoted as  $\mathcal{F}_2^{(FF)}$ , with width at most  $5nK^{d_x n}$  (set  $r = n \cdot |\mathbb{G}_K| \leq nK^{d_x n}$  therein), such that

$$\mathcal{F}_2^{(FF)}(\mathcal{F}^{(SA)}(\mathbf{G} + \mathbf{P})) = \mathbf{F}(\mathbf{G}) \quad \text{for all } \mathbf{G} \in \mathbb{G}_K,$$

and  $\|\mathcal{F}_2^{(FF)}(\mathbf{Z})\|_F \leq \sqrt{d_x n} K_{\mathcal{H}}$ .

Let  $\mathcal{E}_{out}$  be the identity mapping. It holds that  $\mathcal{E}_{out} \circ \mathcal{F}_2^{(FF)} \circ \mathcal{F}^{(SA)} \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in} \in \mathcal{T}_{d_x, d_x}(d_x, 1, 1, 5nK^{d_x n}, 2)$ . Note that for any  $\mathbf{X} \in \omega_{\mathbf{G}} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta)$ ,

$$\begin{aligned} & \mathcal{E}_{out} \circ \mathcal{F}_2^{(FF)} \circ \mathcal{F}^{(SA)} \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in}(\mathbf{X}) \\ &= \mathcal{F}_2^{(FF)} \circ \mathcal{F}^{(SA)} \circ \mathcal{F}_1^{(FF)}(\mathbf{X} + \mathbf{P}) \\ &= \mathcal{F}_2^{(FF)} \circ \mathcal{F}^{(SA)} \circ \overline{\mathcal{F}}_1^{(FF)}(\mathbf{X} + \mathbf{P}) \\ &= \mathcal{F}_2^{(FF)} \circ \mathcal{F}^{(SA)}(\mathbf{G} + \mathbf{P}) \\ &= \mathbf{F}(\mathbf{G}) = \overline{\mathbf{F}}(\mathbf{X}). \end{aligned}$$

Thus, for any  $\mathbf{X} \in [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta) = \bigcup_{\mathbf{G} \in \mathbb{G}_K} \omega_{\mathbf{G}} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta)$ , we have

$$\mathcal{E}_{out} \circ \mathcal{F}_2^{(FF)} \circ \mathcal{F}^{(SA)} \circ \mathcal{F}_1^{(FF)} \circ \mathcal{E}_{in}(\mathbf{X}) = \overline{\mathbf{F}}(\mathbf{X}),$$

which completes the proof by noting (4). □

**Proposition 8.** *Given  $\gamma \in (0, 1]$  and  $K_{\mathcal{H}} > 0$ , assume that  $\mathbf{F} : [0, 1]^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  with each entry  $F_{i,j} \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$ . For any  $\varepsilon > 0$ ,  $K \in \mathbb{N}$  and  $\delta \in (0, \frac{1}{3K}]$ , if  $\tilde{\mathcal{N}} \in \mathcal{T}_{d_x, d_x}(D, H, S, W, L)$  is a Transformer network that satisfies*

$$|\tilde{\mathcal{N}}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq \varepsilon$$

for any  $i \in [d_x]$ ,  $j \in [n]$  and  $\mathbf{X} \in [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta)$ , then there exists a new Transformer network

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(3^{d_x n} \max\{D, 5d_x\}, 3^{d_x n} H, S, 3^{d_x n} \max\{W, 14d_x\}, L + 2d_x n),$$

such that

$$|\mathcal{N}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq \varepsilon + d_x n K_{\mathcal{H}} \delta^\gamma$$

for any  $i \in [d_x]$ ,  $j \in [n]$  and  $\mathbf{X} \in [0, 1]^{d_x \times n}$ .

*Proof.* We basically follow the proof of [39, Theorem 2.1].

**Step 1:** We prove that, given  $i \in [d_x]$ ,  $j \in [n]$ ,  $F_{i,j} \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$ , and a general function  $G_{i,j} : \mathbb{R}^{d_x \times n} \rightarrow \mathbb{R}$  satisfying

$$|G_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq \varepsilon \quad \text{for any } \mathbf{X} \in [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta), \quad (5)$$

then

$$|\Phi_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq \varepsilon + d_x n K_{\mathcal{H}} \delta^\gamma \quad \text{for any } \mathbf{X} \in [0, 1]^{d_x \times n}, \quad (6)$$

where  $\Phi_{i,j} := \Phi_{i,j}^{(d_x n)}$  is defined by induction through

$$\Phi_{i,j}^{(k)}(\mathbf{X}) := \text{mid} \left( \Phi_{i,j}^{(k-1)}(\mathbf{X} - \delta \mathbf{E}^{(k)}), \Phi_{i,j}^{(k-1)}(\mathbf{X}), \Phi_{i,j}^{(k-1)}(\mathbf{X} + \delta \mathbf{E}^{(k)}) \right) \quad (7)$$

for  $k = 1, 2, \dots, d_x n$ ,  $\Phi_{i,j}^{(0)} = G_{i,j}$ ,  $\text{mid}(\cdot, \cdot, \cdot)$  is a function returning the middle value of three inputs, and  $\mathbf{E}^{(u+(v-1)d_x)}$  denotes the matrix with 1 at the  $(u, v)$ -th position and 0 elsewhere for  $u \in [d_x], v \in [n]$ . In other words, if  $G_{i,j}$  provides a uniform approximation outside the trifling region, then the carefully constructed  $\Phi_{i,j}$  extends this uniform approximation to the entire domain, with only a slight increase in the approximation error.

Note that  $\{\mathbf{E}^{(k)}\}_{k=1}^{d_x n}$  defined above is a re-indexing of the standard basis in  $\mathbb{R}^{d_x \times n}$ . We re-index the elements of  $\mathbf{X} = (X_{u,v})$  in the same manner. Let  $X^{(u+(v-1)d_x)} = X_{u,v}$  for  $u \in [d_x], v \in [n]$ . Using this notation, define

$$\Omega_k := \left\{ \mathbf{X} : X^{(i)} \in \begin{cases} [0, 1], & \text{if } i \leq k \\ [0, 1] \setminus \Omega([0, 1], K, \delta), & \text{if } i > k \end{cases} \right\}.$$

Clearly,  $\Omega_0 = [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta)$  and  $\Omega_{d_x n} = [0, 1]^{d_x \times n}$ .

We will prove by induction that for each  $k \in \{0, 1, \dots, d_x n\}$ ,

$$|\Phi_{i,j}^{(k)}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq \varepsilon + k \cdot K_{\mathcal{H}} \delta^\gamma \quad \text{for any } \mathbf{X} \in \Omega_k. \quad (8)$$

As the final step of the induction, we derive

$$\begin{aligned} |\Phi_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| &= |\Phi_{i,j}^{(d_x n)}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \\ &\leq \varepsilon + d_x n K_{\mathcal{H}} \delta^\gamma \quad \text{for any } \mathbf{X} \in \Omega_{d_x n} = [0, 1]^{d_x \times n}, \end{aligned}$$

which completes the proof of (6).

In the base case, it follows from (5) that

$$\begin{aligned} |\Phi_{i,j}^{(0)}(\mathbf{X}) - F_{i,j}(\mathbf{X})| &= |G_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \\ &\leq \varepsilon \quad \text{for any } \mathbf{X} \in \Omega_0 = [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta). \end{aligned}$$

Now, assume that for some  $k \in \{1, 2, \dots, d_x n\}$ ,

$$|\Phi_{i,j}^{(k-1)}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq \varepsilon + (k-1) K_{\mathcal{H}} \delta^\gamma \quad \text{for any } \mathbf{X} \in \Omega_{k-1}.$$

For fixed  $X^{(1)}, \dots, X^{(k-1)} \in [0, 1]$  and  $X^{(k+1)}, \dots, X^{(d_x n)} \in [0, 1] \setminus \Omega([0, 1], K, \delta)$ , define

$$\phi(t) = \Phi_{i,j}^{(k-1)}(X^{(1)}, \dots, X^{(k-1)}, t, X^{(k+1)}, \dots, X^{(d_x n)})$$

and

$$f(t) = F_{i,j}(X^{(1)}, \dots, X^{(k-1)}, t, X^{(k+1)}, \dots, X^{(d_x n)}).$$

The induction hypothesis gives

$$|\phi(t) - f(t)| \leq \varepsilon + (k-1) \cdot K_{\mathcal{H}} \delta^\gamma \quad \text{for any } t \in [0, 1] \setminus \Omega([0, 1], K, \delta).$$

Since  $F_{i,j} \in \mathcal{H}^\gamma([0, 1]^{d_x \times n}, K_{\mathcal{H}})$  implies  $f \in \mathcal{H}^\gamma([0, 1], K_{\mathcal{H}})$ , applying Lemma 10 to the univariate functions  $\phi(t)$  and  $f(t)$  yields

$$|\tilde{\phi}(t) - f(t)| \leq \varepsilon + (k-1) \cdot K_{\mathcal{H}} \delta^\gamma + K_{\mathcal{H}} \delta^\gamma = \varepsilon + k \cdot K_{\mathcal{H}} \delta^\gamma \quad \text{for any } t \in [0, 1],$$

where

$$\begin{aligned}
\tilde{\phi}(t) &= \text{mid}(\phi(t - \delta), \phi(t), \phi(t + \delta)) \\
&= \text{mid}(\Phi_{i,j}^{(k-1)}(X^{(1)}, \dots, X^{(k-1)}, t - \delta, X^{(k+1)}, \dots, X^{(d_x n)}), \\
&\quad \Phi_{i,j}^{(k-1)}(X^{(1)}, \dots, X^{(k-1)}, t, X^{(k+1)}, \dots, X^{(d_x n)}), \\
&\quad \Phi_{i,j}^{(k-1)}(X^{(1)}, \dots, X^{(k-1)}, t + \delta, X^{(k+1)}, \dots, X^{(d_x n)})) \\
&= \Phi_{i,j}^{(k)}(X^{(1)}, \dots, X^{(k-1)}, t, X^{(k+1)}, \dots, X^{(d_x n)})
\end{aligned}$$

by definition of  $\Phi_{i,j}^{(k)}$ . Since  $X^{(1)}, \dots, X^{(k-1)} \in [0, 1]$ ,  $X^{(k)} = t \in [0, 1]$  and  $X^{(k+1)}, \dots, X^{(d_x n)} \in [0, 1] \setminus \Omega([0, 1], K, \delta)$  are arbitrary, we obtain

$$|\Phi_{i,j}^{(k)}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq \varepsilon + k \cdot K_{\mathcal{H}} \delta^\gamma \quad \text{for any } \mathbf{X} \in \Omega_k.$$

This completes the induction.

**Step 2:** Recall that  $\Phi = (\Phi_{i,j})_{i \in [d_x], j \in [n]}$  is defined by (7). We now prove that if  $\mathbf{G} \in \mathcal{T}_{d_x, d_x}(D, H, S, W, L)$ , then

$$\Phi \in \mathcal{T}_{d_x, d_x}(3^{d_x n} \max\{D, 5d_x\}, 3^{d_x n} H, S, 3^{d_x n} \max\{W, 14d_x\}, L + 2d_x n).$$

The observation here is that, to compute  $\Phi = \Phi^{(d_x n)}$ , we first evaluate

$$\Phi^{(d_x n - 1)}(\cdot + c_{d_x n} \delta \mathbf{E}^{(d_x n)}) \quad \text{for each } c_{d_x n} \in \{-1, 0, 1\}.$$

Each such evaluation, in turn, requires computing

$$\Phi^{(d_x n - 2)}(\cdot + c_{d_x n - 1} \delta \mathbf{E}^{(d_x n - 1)} + c_{d_x n} \delta \mathbf{E}^{(d_x n)}) \quad \text{for each } c_{d_x n - 1} \in \{-1, 0, 1\}.$$

Continuing this process recursively, determining  $\Phi$  ultimately requires evaluating

$$\Phi^{(0)}(\cdot + \sum_{l=1}^{d_x n} c_l \delta \mathbf{E}^{(l)}) \quad \text{for every } (c_1, \dots, c_{d_x n}) \in \{-1, 0, 1\}^{d_x n}.$$

Conversely, since  $\Phi^{(0)} = \mathbf{G}$  by definition, assuming that we have access to all functions  $\mathbf{G}(\cdot + \sum_{l=1}^{d_x n} c_l \delta \mathbf{E}^{(l)})$ , we can iteratively apply the mid function to  $\Phi^{(k)}$  to recover  $\Phi^{(k+1)}$ , following the same construction as in (7), and ultimately get  $\Phi$ . On the other hand, each function  $\mathbf{G}(\cdot + \sum_{k=1}^{d_x n} c_k \delta \mathbf{E}^{(k)})$  is a Transformer network thanks to positional encoding, and the mid function can be implemented by feed-forward layers and vectorized operations, thereby completing the construction. The details are given below.

Fixing  $k \in \{0, 1, \dots, d_x n - 1\}$ , we reindex the functions

$$\left\{ \Phi^{(k)}(\cdot + \sum_{l=k+1}^{d_x n} c_l \delta \mathbf{E}^{(l)}) : (c_{k+1}, \dots, c_{d_x n}) \in \{-1, 0, 1\}^{d_x n - k} \right\}$$

as  $\{\Phi_l^{(k)}\}_{l=1}^{3^{d_x n - k}}$  (set  $\Phi_1^{(d_x n)} = \Phi^{(d_x n)}$  for notational convenience), such that

$$\left[ \Phi_l^{(k+1)} \right]_{i,j} = \text{mid} \left( \left[ \Phi_{3l-2}^{(k)} \right]_{i,j}, \left[ \Phi_{3l-1}^{(k)} \right]_{i,j}, \left[ \Phi_{3l}^{(k)} \right]_{i,j} \right)$$

for all  $i \in [d_x], j \in [n]$ , which aligns with (7). Since  $\text{mid}(\cdot, \cdot, \cdot) \in \mathcal{FNN}_{3,1}(14, 2)$  by Lemma 9, there exists an FNN  $\tilde{\mathcal{N}} \in \mathcal{FNN}_{3d_x, d_x}(14d_x, 2)$ , such that for all  $j \in [n]$ ,

$$\tilde{\mathcal{N}} \left( \begin{bmatrix} \left[ \Phi_{3l-2}^{(k)} \right]_{:,j} \\ \left[ \Phi_{3l-1}^{(k)} \right]_{:,j} \\ \left[ \Phi_{3l}^{(k)} \right]_{:,j} \end{bmatrix} \right) = \left[ \Phi_l^{(k+1)} \right]_{:,j}.$$

We then concatenate  $\tilde{\mathcal{N}}$  in parallel to construct a new FNN

$$\tilde{\mathcal{N}}^{(k)} \in \mathcal{FNN}_{3d_x \cdot 3^{d_x n - k - 1}, d_x \cdot 3^{d_x n - k - 1}}(14d_x \cdot 3^{d_x n - k - 1}, 2)$$

such that

$$\tilde{\mathcal{N}}^{(k)} \begin{pmatrix} \left[ \Phi_1^{(k)} \right]_{:,j} \\ \left[ \Phi_2^{(k)} \right]_{:,j} \\ \left[ \Phi_3^{(k)} \right]_{:,j} \\ \vdots \\ \left[ \Phi_{3^{d_x n - k - 2}}^{(k)} \right]_{:,j} \\ \left[ \Phi_{3^{d_x n - k - 1}}^{(k)} \right]_{:,j} \\ \left[ \Phi_{3^{d_x n - k}}^{(k)} \right]_{:,j} \end{pmatrix} = \begin{pmatrix} \left[ \Phi_1^{(k+1)} \right]_{:,j} \\ \vdots \\ \left[ \Phi_{3^{d_x n - k - 1}}^{(k+1)} \right]_{:,j} \end{pmatrix}.$$

By recursively composing  $\tilde{\mathcal{N}}^{(k)}$  for each  $k \in \{0, 1, \dots, d_x n - 1\}$ , we obtain

$$\tilde{\mathcal{N}}^{(d_x n - 1)} \circ \tilde{\mathcal{N}}^{(d_x n - 2)} \circ \dots \circ \tilde{\mathcal{N}}^{(0)} \in \mathcal{FNN}_{d_x 3^{d_x n}, d_x}(14d_x 3^{d_x n - 1}, 2d_x n),$$

which, by construction, satisfies

$$\begin{aligned} & \tilde{\mathcal{N}}^{(d_x n - 1)} \circ \tilde{\mathcal{N}}^{(d_x n - 2)} \circ \dots \circ \tilde{\mathcal{N}}^{(0)} \begin{pmatrix} \left[ \Phi_1^{(0)} \right]_{:,j} \\ \left[ \Phi_2^{(0)} \right]_{:,j} \\ \left[ \Phi_3^{(0)} \right]_{:,j} \\ \vdots \\ \left[ \Phi_{3^{d_x n - 2}}^{(0)} \right]_{:,j} \\ \left[ \Phi_{3^{d_x n - 1}}^{(0)} \right]_{:,j} \\ \left[ \Phi_{3^{d_x n}}^{(0)} \right]_{:,j} \end{pmatrix} \\ &= \tilde{\mathcal{N}}^{(d_x n - 1)} \circ \tilde{\mathcal{N}}^{(d_x n - 2)} \circ \dots \circ \tilde{\mathcal{N}}^{(1)} \begin{pmatrix} \left[ \Phi_1^{(1)} \right]_{:,j} \\ \vdots \\ \left[ \Phi_{3^{d_x n - 1}}^{(1)} \right]_{:,j} \end{pmatrix} \\ &\vdots \\ &= \left[ \Phi^{(d_x n)} \right]_{:,j}, \end{aligned}$$

for each  $j \in [n]$ . Furthermore, Lemma 5 guarantees that any token-wise FNN can be expressed in terms of feed-forward layers. We have (set  $W = 14d_x 3^{d_x n - 1}$  and  $L = 2d_x n$  therein) an embedding map  $\mathcal{E}_{in} : \mathbb{R}^{d_x 3^{d_x n} \times n} \rightarrow \mathbb{R}^{14d_x 3^{d_x n - 1} \times n}$ , a projection map  $\mathcal{E}_{out} : \mathbb{R}^{14d_x 3^{d_x n - 1} \times n} \rightarrow \mathbb{R}^{d_x \times n}$ , and  $2d_x n$  feed-forward layers  $\mathcal{F}_{L+2d_x n}^{(FF)}, \dots, \mathcal{F}_{L+1}^{(FF)}$ , each with width at most  $3 \cdot 14d_x 3^{d_x n - 1} = 14d_x 3^{d_x n}$ , such that

$$\mathcal{E}_{out} \circ \mathcal{F}_{L+2d_x n}^{(FF)} \circ \dots \circ \mathcal{F}_{L+1}^{(FF)} \circ \mathcal{E}_{in} \begin{pmatrix} \Phi_1^{(0)} \\ \Phi_2^{(0)} \\ \vdots \\ \Phi_{3^{d_x n}}^{(0)} \end{pmatrix} = \Phi^{(d_x n)}. \quad (9)$$

Due to the positional encoding, each function

$$\mathbf{G}(\cdot + \sum_{l=1}^{d_x n} c_k \delta \mathbf{E}^{(l)}) \in \mathcal{T}_{d_x, d_x}(D, H, S, W, L).$$

Recall that  $\mathbf{G} = \Phi^{(0)}$  and  $\{\Phi_l^{(0)}\}_{l=1}^{3^{d_x n}}$  is a reordering of  $\{\Phi^{(0)}(\cdot + \sum_{l=1}^{d_x n} c_l \delta \mathbf{E}^{(l)})\}$ . By concatenation of Transformers (see Proposition 6), there exists a Transformer network

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x, 3^{d_x n}}(3^{d_x n} D, 3^{d_x n} H, S, 3^{d_x n} W, L)$$

such that

$$\mathcal{N}(\mathbf{X}) = \begin{pmatrix} \Phi_1^{(0)} \\ \Phi_2^{(0)} \\ \vdots \\ \Phi_{3^{d_x n}}^{(0)} \end{pmatrix}.$$

Together with (9) and  $\Phi^{(d_x n)} = \Phi$ , we have

$$\begin{aligned} & \mathcal{E}_{out} \circ \mathcal{F}_{L+2d_x n}^{(FF)} \circ \dots \circ \mathcal{F}_{L+1}^{(FF)} \circ \mathcal{E}_{in} \circ \mathcal{N}(\mathbf{X}) \\ &= \mathcal{E}_{out} \circ \mathcal{F}_{L+2d_x n}^{(FF)} \circ \dots \circ \mathcal{F}_{L+1}^{(FF)} \circ \mathcal{E}_{in} \begin{pmatrix} \Phi_1^{(0)} \\ \Phi_2^{(0)} \\ \vdots \\ \Phi_{3^{d_x n}}^{(0)} \end{pmatrix} \\ &= \Phi^{(d_x n)} = \Phi(\mathbf{X}), \end{aligned}$$

thereby

$$\begin{aligned} \Phi &\in \mathcal{T}_{d_x, d_x}(\max\{3^{d_x n} D, 14d_x 3^{d_x n-1}\}, 3^{d_x n} H, S, \max\{3^{d_x n} W, 14d_x 3^{d_x n}\}, L + 2d_x n) \\ &\subseteq \mathcal{T}_{d_x, d_x}(3^{d_x n} \max\{D, 5d_x\}, 3^{d_x n} H, S, 3^{d_x n} \max\{W, 14d_x\}, L + 2d_x n), \end{aligned}$$

which completes the proof.  $\square$

**Lemma 9** (Lemma 3.1 of [39]). *The middle value function  $\text{mid}(x_1, x_2, x_3) \in \mathcal{FNN}_{3,1}(14, 2)$ .*

**Lemma 10** (Lemma 3.3 of [39]). *Given any  $\varepsilon > 0$ ,  $K \in \mathbb{N}$ , and  $\delta \in (0, \frac{1}{3K}]$ , assume that  $f \in \mathcal{H}^\gamma([0, 1], K_{\mathcal{H}})$  and  $g : \mathbb{R} \rightarrow \mathbb{R}$  is a general function with*

$$|g(x) - f(x)| \leq \varepsilon, \text{ for any } x \in [0, 1] \setminus \Omega([0, 1], K, \delta).$$

Then

$$|\phi(x) - f(x)| \leq \varepsilon + K_{\mathcal{H}} \delta^\gamma \text{ for any } x \in [0, 1],$$

where

$$\phi(x) := \text{mid}(g(x - \delta), g(x), g(x + \delta)) \text{ for any } x \in \mathbb{R}.$$

*Proof of Theorem 1. Case 1:*  $p \in [1, \infty)$ . Let

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(d_x, 1, 1, 5nK^{d_x n}, 2)$$

be as in Proposition 7. By Proposition 7 and noting that the Lebesgue measure of  $\Omega([0, 1]^{d_x \times n}, K, \delta)$  is at most  $d_x n K \delta$ , we have

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0, 1]^{d_x \times n})}^p$$

$$\begin{aligned}
&= \int_{[0,1]^{d_x \times n}} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} \\
&= \int_{\Omega([0,1]^{d_x \times n}, K, \delta)} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} + \int_{[0,1]^{d_x \times n} \setminus \Omega([0,1]^{d_x \times n}, K, \delta)} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} \\
&\leq \int_{\Omega([0,1]^{d_x \times n}, K, \delta)} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} \\
&\quad + \int_{[0,1]^{d_x \times n} \setminus \Omega([0,1]^{d_x \times n}, K, \delta)} (d_x n)^{\max\{0, \frac{p}{2} - 1\}} \sum_{i=1}^{d_x} \sum_{j=1}^n |\mathcal{N}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})|^p d\mathbf{X} \\
&\leq (2\sqrt{d_x n} K_{\mathcal{H}})^p \cdot d_x n K \delta + (d_x n)^{1 + \max\{0, \frac{p}{2} - 1\}} (K_{\mathcal{H}}(d_x n)^{\gamma/2} K^{-\gamma})^p \\
&\leq 2^p (d_x n)^{2p} K_{\mathcal{H}}^p ((K\delta)^{\frac{1}{p}} + K^{-\gamma})^p,
\end{aligned}$$

using for the last inequality that  $\gamma \in (0, 1]$ ,  $\max\{a, b\} \leq a + b$  for any  $a, b \geq 0$ , and  $a^p + b^p \leq (a + b)^p$  for all  $p \geq 1$  and  $a, b \geq 0$ . Hence,

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})} \leq 2(d_x n)^2 K_{\mathcal{H}} ((K\delta)^{\frac{1}{p}} + K^{-\gamma}).$$

Choosing  $\delta \leq K^{-p\gamma-1}$  and  $K \geq \varepsilon^{-1/\gamma}$  so that  $K^{d_x n} = \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil$ , we conclude

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})} \leq 4(d_x n)^2 K_{\mathcal{H}} \varepsilon$$

and

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(d_x, 1, 1, 5n \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil, 2).$$

**Case 2:**  $p = \infty$ . By Proposition 7, there exists a Transformer network

$$\tilde{\mathcal{N}} \in \mathcal{T}_{d_x, d_x}(d_x, 1, 1, 5n K^{d_x n}, 2)$$

such that

$$|\tilde{\mathcal{N}}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq (d_x n)^{\gamma/2} K_{\mathcal{H}} K^{-\gamma}$$

for any  $i \in [d_x]$ ,  $j \in [n]$  and  $\mathbf{X} \in [0, 1]^{d_x \times n} \setminus \Omega([0, 1]^{d_x \times n}, K, \delta)$ . By Proposition 8 (assume that  $5n K^{d_x n} \geq 14d_x$ ), there exists a new Transformer network

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(5d_x \mathfrak{Z}^{d_x n}, \mathfrak{Z}^{d_x n}, 1, 5n \mathfrak{Z}^{d_x n} K^{d_x n}, 2 + 2d_x n),$$

such that

$$|\mathcal{N}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \leq (d_x n)^{\gamma/2} K_{\mathcal{H}} K^{-\gamma} + d_x n K_{\mathcal{H}} \delta^\gamma$$

for any  $i \in [d_x]$ ,  $j \in [n]$  and  $\mathbf{X} \in [0, 1]^{d_x \times n}$ . This implies

$$\begin{aligned}
\|\mathcal{N} - \mathbf{F}\|_{L^\infty([0,1]^{d_x \times n})} &= \sup_{\mathbf{X} \in [0,1]^{d_x \times n}} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F \\
&\leq \sup_{\mathbf{X} \in [0,1]^{d_x \times n}} \sum_{i=1}^{d_x} \sum_{j=1}^n |\mathcal{N}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \\
&\leq \sum_{i=1}^{d_x} \sum_{j=1}^n \sup_{\mathbf{X} \in [0,1]^{d_x \times n}} |\mathcal{N}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})| \\
&\leq (d_x n)^{1 + \gamma/2} K_{\mathcal{H}} K^{-\gamma} + (d_x n)^2 K_{\mathcal{H}} \delta^\gamma.
\end{aligned}$$

Choosing  $\delta \in (0, \frac{1}{3K}]$  sufficiently small and  $K \geq \varepsilon^{-1/\gamma}$  so that  $K^{d_x n} = \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil$ , we conclude

$$\|\mathcal{N} - \mathbf{F}\|_{L^\infty([0,1]^{d_x \times n})} \leq 4(d_x n)^2 K_{\mathcal{H}} \varepsilon$$

and

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(5d_x 3^{d_x n}, 3^{d_x n}, 1, 5n 3^{d_x n} \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil, 2 + 2d_x n).$$

This completes the proof.  $\square$

*Proof of Theorem 2.* Let  $K$ ,  $\mathbb{G}_K$  and  $\omega_{\mathbf{G}}$  be as defined in the proof of Proposition 7. We approximate the target function  $\mathbf{F}$  using a piecewise constant function, where the value in each cell is given by the average of  $\mathbf{F}$  over that cell. Define

$$\overline{\mathbf{F}}(\mathbf{X}) = \sum_{\mathbf{G} \in \mathbb{G}_K} \mathbf{F}_{\mathbf{G}} \mathbb{1}_{\omega_{\mathbf{G}}}(\mathbf{X}),$$

where

$$[\mathbf{F}_{\mathbf{G}}]_{i,j} = K^{d_x n} \int_{\omega_{\mathbf{G}}} F_{i,j}(\mathbf{X}) d\mathbf{X}, \quad i \in [d_y], j \in [n].$$

Since each cell  $\omega_{\mathbf{G}}$  is a bounded convex domain, Poincaré inequality gives, for any  $p \in [1, \infty]$ ,

$$\|[\mathbf{F}_{\mathbf{G}}]_{i,j} - F_{i,j}\|_{L^p(\omega_{\mathbf{G}})} \leq C \|\nabla F_{i,j}\|_{L^p(\omega_{\mathbf{G}})} K^{-1},$$

where  $C$  is a constant depending only on  $d_x n$ , and  $\|\nabla F\|_{L^p(\omega)}$  denotes the  $L^p$ -norm of the Frobenius norm of  $\nabla F$  (see [9,17]). Summing over all grid cells and using that  $F_{i,j} \in \mathcal{W}^{1,p}([0,1]^{d_x \times n}, K_{\mathcal{W}})$  implies  $\|\nabla F_{i,j}\|_{L^p([0,1]^{d_x \times n})} \leq (d_x n)^{\max\{0, \frac{1}{2} - \frac{1}{p}\}} K_{\mathcal{W}}$ , we obtain

$$\begin{aligned} \|F_{i,j} - \overline{F}_{i,j}\|_{L^p([0,1]^{d_x \times n})} &= \begin{cases} \left( \sum_{\mathbf{G} \in \mathbb{G}_K} \|F_{i,j} - [\mathbf{F}_{\mathbf{G}}]_{i,j}\|_{L^p(\omega_{\mathbf{G}})}^p \right)^{1/p} & \text{if } p < \infty \\ \sup_{\mathbf{G} \in \mathbb{G}_K} \|F_{i,j} - [\mathbf{F}_{\mathbf{G}}]_{i,j}\|_{L^\infty(\omega_{\mathbf{G}})} & \text{if } p = \infty \end{cases} \\ &\leq C \|\nabla F_{i,j}\|_{L^p([0,1]^{d_x \times n})} K^{-1} \\ &\leq C (d_x n)^{\max\{0, \frac{1}{2} - \frac{1}{p}\}} K_{\mathcal{W}} K^{-1}, \end{aligned} \quad (10)$$

for any  $p \in [1, \infty]$ .

From **Step 2** and **Step 3** of Proposition 7, there exists a Transformer network

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(d_x, 1, 1, 5n K^{d_x n}, 2)$$

such that

$$\mathcal{N}(\mathbf{X}) = \overline{\mathbf{F}}(\mathbf{X}) \quad \text{for any } \mathbf{X} \in [0,1]^{d_x \times n} \setminus \Omega([0,1]^{d_x \times n}, K, \delta)$$

and

$$\|\mathcal{N}(\mathbf{X})\|_F \leq \sqrt{d_x n} K_{\mathcal{W}} \quad \text{for any } \mathbf{X} \in \mathbb{R}^{d_x \times n}.$$

Since the Lebesgue measure of  $\Omega([0,1]^{d_x \times n}, K, \delta)$  is at most  $d_x n K \delta$ , for  $p \in [1, \infty)$ , we have

$$\begin{aligned} &\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})}^p \\ &= \int_{[0,1]^{d_x \times n}} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} \end{aligned}$$

$$\begin{aligned}
&= \int_{\Omega([0,1]^{d_x \times n}, K, \delta)} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} + \int_{[0,1]^{d_x \times n} \setminus \Omega([0,1]^{d_x \times n}, K, \delta)} \|\overline{\mathbf{F}}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} \\
&\leq \int_{\Omega([0,1]^{d_x \times n}, K, \delta)} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} \\
&\quad + \int_{[0,1]^{d_x \times n} \setminus \Omega([0,1]^{d_x \times n}, K, \delta)} (d_x n)^{\max\{0, \frac{p}{2} - 1\}} \sum_{i=1}^{d_x} \sum_{j=1}^n |\overline{F}_{i,j}(\mathbf{X}) - F_{i,j}(\mathbf{X})|^p d\mathbf{X} \\
&\leq (2\sqrt{d_x n} K_{\mathcal{W}})^p \cdot d_x n K \delta + (d_x n)^{\max\{0, \frac{p}{2} - 1\}} \sum_{i=1}^{d_x} \sum_{j=1}^n \left( C (d_x n)^{\max\{0, \frac{1}{2} - \frac{1}{p}\}} K_{\mathcal{W}} K^{-1} \right)^p \\
&\leq (2C)^p (d_x n)^{2p} K_{\mathcal{W}}^p \left( (K\delta)^{\frac{1}{p}} + K^{-1} \right)^p,
\end{aligned}$$

using for the last inequality that  $p \geq 1$ ,  $\max\{a, b\} \leq a + b$  for any  $a, b \geq 0$ , and  $a^p + b^p \leq (a + b)^p$  for all  $p \geq 1$  and  $a, b \geq 0$ . Hence,

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})} \leq 2C (d_x n)^2 K_{\mathcal{W}} \left( (K\delta)^{\frac{1}{p}} + K^{-1} \right).$$

Choosing  $\delta \leq K^{-p-1}$  and  $K \geq \varepsilon^{-1}$  so that  $K^{d_x n} = \lceil \varepsilon^{-d_x n} \rceil$ , we conclude

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})} \leq 4C (d_x n)^2 K_{\mathcal{W}} \varepsilon$$

and

$$\mathcal{N} \in \mathcal{T}_{d_x, d_x}(d_x, 1, 1, 5n \lceil \varepsilon^{-d_x n} \rceil, 2).$$

This completes the proof.  $\square$

**Lemma 11.** *Let  $d_x, d_y, r \in \mathbb{N}$  with  $d_x \geq d_y$ . Let  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^r$  be a set of input-output pairs such that  $\mathbf{x}_i \in \mathbb{R}^{d_x}$ ,  $\mathbf{y}_i \in \mathbb{R}^{d_y}$ ,  $i \in [r]$  and  $\mathbf{x}_i \neq \mathbf{x}_j$  if  $i \neq j$ . Then, there exists a feed-forward layer  $\mathcal{F}^{(FF)} : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_x}$  with width at most  $3r + 2d_x$  such that*

$$\mathcal{F}^{(FF)}(\mathbf{x}_i) = \begin{pmatrix} \mathbf{y}_i \\ \mathbf{0} \end{pmatrix} \quad \text{for all } i \in [r],$$

and  $\|\mathcal{F}^{(FF)}(\mathbf{z})\| \leq \max_i \|\mathbf{y}_i\|$  for any  $\mathbf{z} \in \mathbb{R}^{d_x}$ .

*Proof.* Let  $R > 0$  be determined later. Since  $\mathbf{x}_i, i \in [r]$  are pairwise distinct, we can find  $\mathbf{v} \in \mathbb{R}^{d_x}$  such that  $\mathbf{v}^\top \mathbf{x}_i, i \in [r]$  are distinct. The existence of  $\mathbf{v}$  can be found in [45, Lemma 13]. We define

$$\mathbf{A}_i^{(1)} = R \mathbf{1}_3 \mathbf{v}^\top, \quad \mathbf{b}_i^{(1)} = \begin{pmatrix} -R \mathbf{v}^\top \mathbf{x}_i - 1 \\ -R \mathbf{v}^\top \mathbf{x}_i \\ -R \mathbf{v}^\top \mathbf{x}_i + 1 \end{pmatrix}, \quad \mathbf{A}_i^{(2)} = \begin{pmatrix} \mathbf{y}_i \\ \mathbf{0} \end{pmatrix} (1, -2, 1), \quad \mathbf{b}_i^{(2)} = \mathbf{0}.$$

Then, by direct calculation, we obtain

$$\begin{aligned}
&\mathbf{A}_i^{(2)} \sigma_R[\mathbf{A}_i^{(1)} \mathbf{x} + \mathbf{b}_i^{(1)}] + \mathbf{b}_i^{(2)} \\
&= \begin{pmatrix} \mathbf{y}_i \\ \mathbf{0} \end{pmatrix} (\sigma_R[R \mathbf{v}^\top (\mathbf{x} - \mathbf{x}_i) - 1] - 2\sigma_R[R \mathbf{v}^\top (\mathbf{x} - \mathbf{x}_i)] + \sigma_R[R \mathbf{v}^\top (\mathbf{x} - \mathbf{x}_i) + 1]) \\
&= \begin{pmatrix} \mathbf{y}_i \\ \mathbf{0} \end{pmatrix} I_i(\mathbf{x}),
\end{aligned}$$

where  $I_i(\mathbf{x})$  is the hat function with  $I_i(\mathbf{x}_i) = 1$  and  $I_i(\mathbf{x}) = 0$  if  $|\mathbf{v}^\top(\mathbf{x} - \mathbf{x}_i)| \geq 1/R$ . To ensure the supports of  $I_i(\mathbf{x})$  for all  $i \in [r]$  are disjoint, we choose  $R > 2/\min_{i \neq j} |\mathbf{v}^\top(\mathbf{x}_i - \mathbf{x}_j)|$ . Define

$$\mathbf{A}^{(1)} = \begin{pmatrix} \mathbf{A}_1^{(1)} \\ \vdots \\ \mathbf{A}_r^{(1)} \\ \mathbf{I}_{d_x} \\ -\mathbf{I}_{d_x} \end{pmatrix}, \quad \mathbf{b}^{(1)} = \begin{pmatrix} \mathbf{b}_1^{(1)} \\ \vdots \\ \mathbf{b}_r^{(1)} \\ \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{A}^{(2)} = \left( \mathbf{A}_1^{(2)}, \dots, \mathbf{A}_r^{(2)}, -\mathbf{I}_{d_x}, \mathbf{I}_{d_x} \right), \quad \mathbf{b}^{(2)} = \mathbf{0},$$

and let

$$\begin{aligned} \mathcal{F}^{(FF)}(\mathbf{x}) &= \mathbf{x} + \mathbf{A}^{(2)} \sigma_R[\mathbf{A}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}] + \mathbf{b}^{(2)} \\ &= \sum_{i=1}^r \begin{pmatrix} \mathbf{y}_i \\ \mathbf{0} \end{pmatrix} I_i(\mathbf{x}). \end{aligned}$$

We complete the proof by verifying that

$$\mathcal{F}^{(FF)}(\mathbf{x}_k) = \sum_{i=1}^r \begin{pmatrix} \mathbf{y}_i \\ \mathbf{0} \end{pmatrix} I_i(\mathbf{x}_k) = \begin{pmatrix} \mathbf{y}_k \\ \mathbf{0} \end{pmatrix}$$

and

$$\begin{aligned} \|\mathcal{F}^{(FF)}(\mathbf{x})\| &= \left\| \sum_{i=1}^r \begin{pmatrix} \mathbf{y}_i \\ \mathbf{0} \end{pmatrix} I_i(\mathbf{x}) \right\| \\ &\leq \max_i \|\mathbf{y}_i\| \left\| \sum_{i=1}^r I_i(\mathbf{x}) \right\| \\ &\leq \max_i \|\mathbf{y}_i\|. \end{aligned}$$

□

## 4.2 Proof of Theorem 3

We introduce sample complexities, which measure the richness of the function class in different aspects, and use them to bound the generalization error.

**Definition 5** (VC-dimension). *Let  $\mathcal{H}$  be a class of real-valued functions defined on  $\Omega$ . The VC-dimension of  $\mathcal{H}$ , denoted by  $\text{VCDim}(\mathcal{H})$ , is the largest integer  $N$  for which there exist points  $x_1, \dots, x_N \in \Omega$  such that*

$$|\{\text{sgn}(h(x_1)), \dots, \text{sgn}(h(x_N)) : h \in \mathcal{H}\}| = 2^N.$$

**Definition 6** (Pseudo-dimension). *Let  $\mathcal{H}$  be a class of real-valued functions defined on  $\Omega$ . The pseudo-dimension of  $\mathcal{H}$ , denoted by  $\text{Pdim}(\mathcal{H})$ , is the largest integer  $N$  for which there exist points  $x_1, \dots, x_N \in \Omega$  and constants  $c_1, \dots, c_N \in \mathbb{R}$  such that*

$$|\{\text{sgn}(h(x_1) - c_1), \dots, \text{sgn}(h(x_N) - c_N) : h \in \mathcal{H}\}| = 2^N.$$

**Definition 7** (Covering number). *Let  $\rho$  be a pseudo-metric on  $\mathcal{M}$  and  $S \subseteq \mathcal{M}$ . For any  $\delta > 0$ , a set  $A \subseteq \mathcal{M}$  is called a  $\delta$ -covering of  $S$  if for any  $x \in S$  there exists  $y \in A$  such that  $\rho(x, y) \leq \delta$ . The  $\delta$ -covering number of  $S$ , denoted by  $\mathcal{N}(\delta, S, \rho)$ , is the minimum cardinality of any  $\delta$ -covering of  $S$ .*

**Theorem 12** (Theorem 8.14 of [2]). *Let  $h$  be a function from  $\mathbb{R}^d \times \mathbb{R}^n$  to  $\{0, 1\}$ , determining the class*

$$\mathcal{H} = \{x \mapsto h(a, x) : a \in \mathbb{R}^d\}.$$

*Suppose that  $h$  can be computed by an algorithm that takes as input the pair  $(a, x) \in \mathbb{R}^d \times \mathbb{R}^n$  and returns  $h(a, x)$  after no more than  $t$  of the following operations:*

- *the exponential function  $\alpha \mapsto e^\alpha$  on real numbers,*
- *the arithmetic operations  $+$ ,  $-$ ,  $\times$ , and  $/$  on real numbers,*
- *jumps conditioned on  $>$ ,  $\geq$ ,  $<$ ,  $\leq$ ,  $=$ , and  $\neq$  comparisons of real numbers, and*
- *output 0 or 1.*

*Then  $\text{VCdim}(\mathcal{H}) \leq t^2 d(d + 19 \log_2(9d))$ . Furthermore, if the  $t$  steps include no more than  $q$  in which the exponential function is evaluated, then*

$$\text{VCdim}(\mathcal{H}) \leq (d(q + 1))^2 + 11d(q + 1)(t + \log_2(9d(q + 1))).$$

Theorem 12 gives bounds on the VC-dimension of a function class in terms of the number of arithmetic operations required to compute the functions. This result immediately implies a bound on the VC-dimension (or pseudo-dimension) for Transformer networks. By applying standard techniques in learning theory, one can further derive upper bounds for the covering number. The following lemma summarizes these bounds.

**Lemma 13.** *Recall that  $\mathcal{F} = \{f = \langle \mathcal{N}, \mathbf{E} \rangle : \mathcal{N} \in \mathcal{T}_{d_x, d_x}(D, H, S, W, L)\}$ . Then the following bounds hold:*

- $\text{VCdim}(\mathcal{F}) \lesssim (HS + W)^2 D^2 H^2 L^4,$
- $\text{Pdim}(\mathcal{F}) \lesssim (HS + W)^2 D^2 H^2 L^4,$
- $\sup_{\mathcal{X}} \log \mathcal{N}(\delta, \mathcal{C}_K \mathcal{F}, d_{\mathcal{X}, \infty}) \lesssim (HS + W)^2 D^2 H^2 L^4 \log \frac{mK}{\delta},$  where  $\mathcal{X} = \{\mathbf{X}_i\}_{i=1}^m$  and

$$d_{\mathcal{X}, \infty}(f, g) = \max_{i \in [m]} |f(\mathbf{X}_i) - g(\mathbf{X}_i)|.$$

*We hide constants that depend on  $d_x$  and  $n$ .*

*Proof.* Recall that  $\mathcal{F} = \{f = \langle \mathcal{N}, \mathbf{E} \rangle : \mathcal{N} \in \mathcal{T}_{d_x, d_x}(D, H, S, W, L)\}$ . By carefully counting the computational steps required to evaluate any  $f \in \mathcal{F}$ , we deduce that

- the total number of parameters is bounded by  $d \lesssim (HS + W)DL,$
- the total number of computational operations is bounded by  $t \lesssim L(HDSn + HS n^2 + WDn),$
- the number of evaluations of the exponential function is bounded by  $q \lesssim LHn^2.$

Theorem 12 immediately implies that

$$\begin{aligned} \text{VCdim}(\mathcal{F}) &\leq (d(q + 1))^2 + 11d(q + 1)(t + \log_2(9d(q + 1))) \\ &\lesssim (HS + W)^2 D^2 H^2 L^4, \end{aligned}$$

where we use  $\log(x) \leq x$  for  $x \geq 1$  and suppress constants that depend on  $d_x$  and  $n$ .

For the pseudo-dimension, note that by definition  $\text{VCdim}(\{f(x) - r : f \in \mathcal{F}, r \in \mathbb{R}\}) = \text{Pdim}(\mathcal{F})$ . Using the same reasoning as above, we have

$$\text{Pdim}(\mathcal{F}) \lesssim (HS + W)^2 D^2 H^2 L^4,$$

again by Theorem 12.

Finally, by Theorem 12.2 of [2], we have

$$\begin{aligned}\log \mathcal{N}(\delta, \mathcal{C}_K \mathcal{F}, d_{\mathcal{X}, \infty}) &\leq \text{Pdim}(\mathcal{C}_K \mathcal{F}) \log \frac{emK}{\delta} \\ &\leq \text{Pdim}(\mathcal{F}) \log \frac{emK}{\delta} \\ &\lesssim (HS + W)^2 D^2 H^2 L^4 \log \frac{mK}{\delta}.\end{aligned}$$

Taking the supremum over all possible sample sets  $\mathcal{X}$  completes the proof.  $\square$

*Proof of Theorem 3.* Let  $\mathcal{X}$  and  $d_{\mathcal{X}, \infty}$  be defined as in Lemma 13. Similar to the proof of [32, Theorem 5], given a random sample  $\mathcal{D}_m = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , the excess risk can be decomposed as

$$\mathbb{E}_{\mathcal{D}_m}[\mathcal{R}(\mathcal{C}_{B_m} \hat{f}_m) - \mathcal{R}(f^*)] \lesssim \mathcal{E}_{app} + \mathcal{E}_{gen} + \mathcal{E}_{den},$$

where

$$\begin{aligned}\mathcal{E}_{app} &:= \inf_{f \in \mathcal{F}} \mathbb{E}[(f - f^*)^2], \\ \mathcal{E}_{gen} &:= \frac{B_m^2 k_m}{m} \sup_{|\mathcal{X}|=m} \log \mathcal{N}(m^{-1}, \mathcal{C}_{B_m} \mathcal{F}, d_{\mathcal{X}, \infty}), \\ \mathcal{E}_{den} &:= \frac{B_m^2 m}{k_m} \beta(k_m).\end{aligned}$$

Here,  $k_m \in \mathbb{N}$  is a parameter to be chosen. It can be seen that the excess risk is bounded by the sum of the approximation error  $\mathcal{E}_{app}$ , the generalization error  $\mathcal{E}_{gen}$ , and the dependence error  $\mathcal{E}_{den}$ . Note that as  $k_m$  increases,  $\mathcal{E}_{den}$  decreases due to the monotonic decrease of the  $\beta$ -mixing coefficient  $\beta(k_m)$ , whereas  $\mathcal{E}_{gen}$  increases. Besides, if we select a larger hypothesis class  $\mathcal{F}$ , then  $\mathcal{E}_{app}$  decreases but  $\mathcal{E}_{gen}$  increases because the covering number grows with the size of the hypothesis class. Therefore, to obtain a better convergence rate, we must carefully trade off these three errors by choosing an appropriate hypothesis class  $\mathcal{F}$  and tuning the parameter  $k_m$ .

Since by assumption the density of  $\Pi$  is upper bounded, Theorem 1 implies that

$$\mathcal{E}_{app} \leq \inf_{f \in \mathcal{F}} \|f - f^*\|_{L^2([0,1]^{d_x \times n})}^2 \lesssim \varepsilon^2,$$

where the hypothesis class

$$\mathcal{F} = \mathcal{F}(D_m \lesssim 1, H_m \lesssim 1, S_m \lesssim 1, W_m \lesssim \varepsilon^{-\frac{d_x n}{\gamma}}, L_m \lesssim 1).$$

Then by Lemma 13,

$$\begin{aligned}\mathcal{E}_{gen} &\lesssim \frac{B_m^2 k_m}{m} (HS + W)^2 D^2 H^2 L^4 \log(m^2 B_m) \\ &\lesssim \frac{(\log m)^3 k_m}{m} \varepsilon^{-\frac{2d_x n}{\gamma}},\end{aligned}$$

where we take  $B_m \asymp \log m$ .

We now consider three cases for the sequence  $\{\mathbf{x}_i\}_{i=1}^m$ .

**Case 1:** if  $\{\mathbf{x}_i\}_{i=1}^m$  is geometrically  $\beta$ -mixing, i.e.,  $\beta(k_m) \leq \beta_0 \exp(-\beta_1 k_m^r)$  for some  $r, \beta_0, \beta_1 > 0$ , we set  $k_m \asymp (\log m)^{1/r}$  so that  $\beta(k_m) \lesssim 1/m^{100}$ . Then,

$$\mathbb{E}_{\mathcal{D}_m}[\mathcal{R}(\mathcal{C}_{B_m} \hat{f}_m) - \mathcal{R}(f^*)] \lesssim \varepsilon^2 + \frac{(\log m)^{3+1/r}}{m} \varepsilon^{-\frac{2d_x n}{\gamma}}$$

$$\lesssim m^{-\frac{\gamma}{\gamma+d_x n}} (\log m)^{3+1/r},$$

where  $\varepsilon$  is chosen as  $\varepsilon \asymp m^{-\frac{\gamma}{2\gamma+2d_x n}}$ .

**Case 2:** if  $\{\mathbf{x}_i\}_{i=1}^m$  is algebraically  $\beta$ -mixing, that is,  $\beta(k_m) \leq \beta_0/k_m^r$  for some  $r, \beta_0 > 0$ , then

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_m}[\mathcal{R}(\mathcal{C}_{B_m} \hat{f}_m) - \mathcal{R}(f^*)] &\lesssim \varepsilon^2 + \frac{(\log m)^3 k_m \varepsilon^{-\frac{2d_x n}{\gamma}}}{m} + \frac{(\log m)^2 m}{k_m^{r+1}} \\ &\lesssim m^{-\frac{r\gamma}{(r+2)\gamma+(r+1)d_x n}} (\log m)^3, \end{aligned}$$

where we use the AM-GM inequality and choose  $\varepsilon \asymp m^{-\frac{r\gamma}{2(r+2)\gamma+2(r+1)d_x n}}$  and  $k_m \asymp m^{\frac{2\gamma+d_x n}{(r+2)\gamma+(r+1)d_x n}}$ .

**Case 3:** if  $\{\mathbf{x}_i\}_{i=1}^m$  is a sequence of i.i.d. random variables, then  $\beta(k_m) = 0$  for all  $k_m \geq 1$ . This implies

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_m}[\mathcal{R}(\mathcal{C}_{B_m} \hat{f}_m) - \mathcal{R}(f^*)] &\lesssim \varepsilon^2 + \frac{(\log m)^3}{m} \varepsilon^{-\frac{2d_x n}{\gamma}} \\ &\lesssim m^{-\frac{\gamma}{\gamma+d_x n}} (\log m)^3, \end{aligned}$$

where we choose  $\varepsilon \asymp m^{-\frac{\gamma}{2\gamma+2d_x n}}$ . So we complete the proof.  $\square$

### 4.3 Proof of Theorem 4

The original Kolmogorov-Arnold representation theorem states that for any continuous function  $f : [0, 1]^d \rightarrow \mathbb{R}$ , there exist univariate continuous functions  $g_q, \psi_{p,q}$  such that

$$f(x_1, \dots, x_d) = \sum_{q=0}^{2d} g_q \left( \sum_{p=1}^d \psi_{p,q}(x_p) \right).$$

[51] derived modifications of this representation that transfer smoothness properties of the represented function to the outer function.

**Proposition 14** (Theorem 2 of [51]). *For any fixed dimension  $d \geq 2$ , there exists a monotone function  $\phi : [0, 1] \rightarrow \mathcal{C}$  (the Cantor set) such that for any function  $f \in \mathcal{H}^\gamma([0, 1]^d, K_{\mathcal{H}})$  with some  $\gamma \in (0, 1]$ , we can find a function  $g \in \mathcal{H}^{\frac{\gamma \log 2}{d \log 3}}(\mathcal{C}, 2\sqrt{d}K_{\mathcal{H}})$  such that*

$$f(x_1, \dots, x_d) = g \left( 3 \sum_{p=1}^d 3^{-p} \phi(x_p) \right). \quad (11)$$

Moreover, for any  $x \in [0, 1]$  with its binary representation  $x = [0.a_1^x a_2^x \dots]_2$ , the function  $\phi$  is given explicitly by

$$\phi(x) = \sum_{j=1}^{\infty} \frac{2a_j^x}{3^{1+d(j-1)}} = [0.(2a_1^x) \underbrace{0 \dots 0}_{(d-1)\text{-times}} (2a_2^x) \underbrace{0 \dots 0}_{(d-1)\text{-times}}]_3,$$

where  $[\cdot]_B$  denotes the  $B$ -adic expansion of a real number.

We note that a given real number can have multiple  $B$ -adic representations (for example,  $[1]_{10} = [0.999 \dots]_{10}$ ), which may make  $\phi$  not well-defined. To eliminate this ambiguity, we adopt the convention of using a unique  $B$ -adic representation for all real numbers. Observe that the argument of  $g$  in (11) satisfies

$$3 \sum_{p=1}^d 3^{-p} \phi(x_p) = [0.(2a_1^{x_1})(2a_1^{x_2}) \dots (2a_1^{x_d})(2a_2^{x_1}) \dots]_3. \quad (12)$$

By construction, the Cantor set consists precisely of those numbers in  $[0, 1]$  whose ternary expansion contains only the digits 0 and 2. This shows that the mapping  $3 \sum_{p=1}^d 3^{-p} \phi(x_p)$  indeed defines a bijection between  $[0, 1]^d$  and the Cantor set  $\mathcal{C}$ . Additionally, an approximation of  $\phi$  with a truncation parameter  $K$  is defined by

$$\phi_K(x) := \sum_{j=1}^K \frac{2a_j^x}{3^{1+d(j-1)}}, \quad (13)$$

which will be used in our construction.

*Proof of Theorem 4.* By Proposition 14 and the fact that  $\phi_K$  approximates  $\phi$ , there exists a function  $\mathbf{G} : \mathbb{R}^{d_x \times n} \rightarrow \mathbb{R}^{d_x \times n}$  with each entry  $G_{u,v} \in \mathcal{H}_{\frac{\gamma \log 2}{d_x n \log 3}}(\mathcal{C}, 2\sqrt{d_x n} K_{\mathcal{H}})$  such that

$$\mathbf{F}(\mathbf{X}) = \mathbf{G} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi(X_{p,q}) \right) \approx \mathbf{G} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi_K(X_{p,q}) \right). \quad (14)$$

We will construct a generalized Transformer network that approximates the latter mapping. In the proof below, for simplicity, we omit the placeholder zeros used for alignment.

**Step 1:** We first show that there exist  $2K+2$  generalized feed-forward layers  $\mathcal{F}_1^{(GFF)}, \dots, \mathcal{F}_{2K+2}^{(GFF)}$  such that

$$\mathcal{F}_{2K+2}^{(GFF)} \circ \dots \circ \mathcal{F}_1^{(GFF)} : \mathbf{X} \mapsto \mathbf{Z}_1,$$

where

$$\mathbf{Z}_1 = \begin{pmatrix} 3 \sum_{p=1}^{d_x} a_{p,1} \tilde{\phi}_K(X_{p,1}) & 3 \sum_{p=1}^{d_x} a_{p,2} \tilde{\phi}_K(X_{p,2}) & \cdots & 3 \sum_{p=1}^{d_x} a_{p,n} \tilde{\phi}_K(X_{p,n}) \\ \vdots & \vdots & & \vdots \\ 3 \sum_{p=1}^{d_x} a_{p,1} \tilde{\phi}_K(X_{p,1}) & 3 \sum_{p=1}^{d_x} a_{p,2} \tilde{\phi}_K(X_{p,2}) & \cdots & 3 \sum_{p=1}^{d_x} a_{p,n} \tilde{\phi}_K(X_{p,n}) \end{pmatrix}.$$

Here,  $a_{p,q} = \frac{1}{3^{(q-1)d_x+p}}$  and  $\tilde{\phi}_K$  is a function that satisfies

$$\tilde{\phi}_K(x) = \begin{cases} 0, & \text{if } x < 0, \\ \phi_K(x), & \text{if } x \in \Omega_K \subseteq [0, 1], \\ 1, & \text{if } x > 1, \end{cases}$$

where  $\phi_K$  is defined in (13) and  $\Omega_K \subseteq [0, 1]$  has Lebesgue measure at least  $1 - 2^{-K\gamma p}$ . [51, Theorem 3] guarantees the existence of an FNN  $\tilde{\phi}_K \in \mathcal{FNN}_{1,1}(4, 2K)$  with the above properties.

By parallel computation, we can construct an FNN  $\tilde{\mathcal{N}}_1 \in \mathcal{FNN}_{1,1}(4n, 2K)$  such that

$$\begin{aligned} \tilde{\mathcal{N}}_1(x) &= \left( 1, 3^{-d_x}, 3^{-2d_x}, \dots, 3^{-(n-1)d_x} \right) \begin{pmatrix} \tilde{\phi}_K(x) \\ \tilde{\phi}_K(x-2) \\ \tilde{\phi}_K(x-4) \\ \vdots \\ \tilde{\phi}_K(x-2(n-1)) \end{pmatrix} \\ &= \sum_{q=1}^n 3^{-(q-1)d_x} \tilde{\phi}_K(x-2(q-1)). \end{aligned}$$

If  $x \in [2(j-1), 2j-1]$  for some  $j \in [n]$ , then

$$\tilde{\mathcal{N}}_1(x) = \sum_{q=1}^{j-1} 3^{-(q-1)d_x} \tilde{\phi}_K(x-2(q-1)) + 3^{-(j-1)d_x} \tilde{\phi}_K(x-2(j-1))$$

$$\begin{aligned}
& + \sum_{q=j+1}^n 3^{-(q-1)d_x} \tilde{\phi}_K(x - 2(q-1)) \\
& = \sum_{q=1}^{j-1} 3^{-(q-1)d_x} + 3^{-(j-1)d_x} \tilde{\phi}_K(x - 2(j-1)) \\
& = 3^{-(j-1)d_x} \tilde{\phi}_K(x - 2(j-1)) + b_j,
\end{aligned}$$

where we define  $b_1 = 0$  and  $b_j = \sum_{q=1}^{j-1} 3^{-(q-1)d_x}$  for  $j \geq 2$ .

Now consider  $\mathbf{x} \in \mathbb{R}^{d_x}$ . Fixing  $j \in [n]$ , if  $x_i \in [2(j-1), 2j-1]$  for all  $i \in [d_x]$ , we can construct an FNN  $\tilde{\mathcal{N}}_2 \in \mathcal{FNN}_{d_x, d_x}(4d_x n, 2K)$  such that

$$\begin{aligned}
\tilde{\mathcal{N}}_2(\mathbf{x}) & = \mathbf{1}_{d_x} \left( 1, 3^{-1}, \dots, 3^{1-d_x} \right) \begin{pmatrix} \tilde{\mathcal{N}}_1(x_1) \\ \tilde{\mathcal{N}}_1(x_2) \\ \vdots \\ \tilde{\mathcal{N}}_1(x_{d_x}) \end{pmatrix} \\
& = \mathbf{1}_{d_x} \left( 1, 3^{-1}, \dots, 3^{1-d_x} \right) \begin{pmatrix} 3^{-(j-1)d_x} \tilde{\phi}_K(x_1 - 2(j-1)) + b_j \\ 3^{-(j-1)d_x} \tilde{\phi}_K(x_2 - 2(j-1)) + b_j \\ \vdots \\ 3^{-(j-1)d_x} \tilde{\phi}_K(x_{d_x} - 2(j-1)) + b_j \end{pmatrix} \\
& = \left( \sum_{p=1}^{d_x} 3^{1-p-(j-1)d_x} \tilde{\phi}_K(x_p - 2(j-1)) \right) \mathbf{1}_{d_x} + \left( b_j \sum_{p=1}^{d_x} 3^{1-p} \right) \mathbf{1}_{d_x} \\
& = \left( 3 \sum_{p=1}^{d_x} a_{p,j} \tilde{\phi}_K(x_p - 2(j-1)) \right) \mathbf{1}_{d_x} + c_j \mathbf{1}_{d_x},
\end{aligned}$$

where we define  $c_j = b_j \sum_{p=1}^{d_x} 3^{1-p}$ . Using that  $X_{p,j} \in [0, 1]$  implies  $X_{p,j} + 2(j-1) \in [2(j-1), 2j-1]$ , set  $\mathbf{x} = \mathbf{X}_{:,j} + 2(j-1)\mathbf{1}_{d_x}$  in the above equation to obtain

$$\tilde{\mathcal{N}}_2(\mathbf{X}_{:,j} + 2(j-1)\mathbf{1}_{d_x}) = \left( 3 \sum_{p=1}^{d_x} a_{p,j} \tilde{\phi}_K(X_{p,j}) \right) \mathbf{1}_{d_x} + c_j \mathbf{1}_{d_x}.$$

By Lemma 5, there exist  $2K$  feed-forward layers  $\mathcal{F}_2^{(FF)}, \dots, \mathcal{F}_{2K+1}^{(FF)}$ , each with width at most  $3 \cdot 4d_x n = 12d_x n$ , such that

$$\mathcal{F}_{2K+1}^{(FF)} \circ \dots \circ \mathcal{F}_2^{(FF)}(\mathbf{X}_{:,1}, \dots, \mathbf{X}_{:,n} + 2(n-1)\mathbf{1}_{d_x}) = \left( \tilde{\mathcal{N}}_2(\mathbf{X}_{:,1}), \dots, \tilde{\mathcal{N}}_2(\mathbf{X}_{:,n} + 2(n-1)\mathbf{1}_{d_x}) \right),$$

where we omit placeholder zeros for simplicity. Finally, to add and then remove the bias terms, we use two generalized feed-forward layers. We define

$$\mathcal{F}_1^{(GFF)}(\mathbf{X}_{:,1}, \dots, \mathbf{X}_{:,n}) := (\mathbf{X}_{:,1}, \dots, \mathbf{X}_{:,n} + 2(n-1)\mathbf{1}_{d_x})$$

and

$$\mathcal{F}_{2K+2}^{(GFF)}(\mathbf{Z}_{:,1}, \dots, \mathbf{Z}_{:,n}) := (\mathbf{Z}_{:,1} - c_1 \mathbf{1}_{d_x}, \dots, \mathbf{Z}_{:,n} - c_n \mathbf{1}_{d_x}).$$

It is straightforward to verify that

$$\begin{aligned}
& \mathcal{F}_{2K+2}^{(GFF)} \circ \mathcal{F}_{2K+1}^{(FF)} \circ \dots \circ \mathcal{F}_2^{(FF)} \circ \mathcal{F}_1^{(GFF)}(\mathbf{X}) \\
& = \mathcal{F}_{2K+2}^{(GFF)} \circ \mathcal{F}_{2K+1}^{(FF)} \circ \dots \circ \mathcal{F}_2^{(FF)}(\mathbf{X}_{:,1}, \dots, \mathbf{X}_{:,n} + 2(n-1)\mathbf{1}_{d_x})
\end{aligned}$$

$$\begin{aligned}
&= \mathcal{F}_{2K+2}^{(GFF)} \left( \tilde{\mathcal{N}}_2(\mathbf{X}_{:,1}), \dots, \tilde{\mathcal{N}}_2(\mathbf{X}_{:,n} + 2(n-1)\mathbf{1}_{d_x}) \right) \\
&= \mathcal{F}_{2K+2}^{(GFF)} \left( \left( 3 \sum_{p=1}^{d_x} a_{p,1} \tilde{\phi}_K(X_{p,1}) \right) \mathbf{1}_{d_x} + c_1 \mathbf{1}_{d_x}, \dots, \left( 3 \sum_{p=1}^{d_x} a_{p,n} \tilde{\phi}_K(X_{p,n}) \right) \mathbf{1}_{d_x} + c_n \mathbf{1}_{d_x} \right) \\
&= \left( \left( 3 \sum_{p=1}^{d_x} a_{p,1} \tilde{\phi}_K(X_{p,1}) \right) \mathbf{1}_{d_x}, \dots, \left( 3 \sum_{p=1}^{d_x} a_{p,n} \tilde{\phi}_K(X_{p,n}) \right) \mathbf{1}_{d_x} \right) \\
&= \mathbf{Z}_1.
\end{aligned}$$

**Step 2:** We show that there exist a generalized self-attention layer  $\mathcal{F}^{(GSA)}$  and a generalized feed-forward layer  $\mathcal{F}_{2K+3}^{(GFF)}$  such that

$$\mathcal{F}_{2K+3}^{(GFF)} \circ \mathcal{F}^{(GSA)} : \begin{pmatrix} \mathbf{Z}_1 \\ \mathbf{O} \end{pmatrix} \mapsto \begin{pmatrix} \mathbf{Z}_2 \\ \mathbf{O} \end{pmatrix},$$

where

$$\mathbf{Z}_2 = \begin{pmatrix} 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) & 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2 & \cdots & 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2(n-1) \\ \vdots & \vdots & & \vdots \\ 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) & 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2 & \cdots & 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2(n-1) \end{pmatrix}.$$

Note that  $\mathbf{Z}_2$  is obtained by summing the columns of  $\mathbf{Z}_1$  and then adding different bias terms to each column.

We now prove the existence of such layers by first considering a standard self-attention layer. In fact, we only require the softmax function to compute the column average, so it can be replaced by a generalized self-attention layer. We define a self-attention layer by choosing the parameters as follows:

$$H = 1, \quad S = d_x, \quad \mathbf{W}^{(O)} = n \begin{pmatrix} \mathbf{O}_{d_x} \\ \mathbf{I}_{d_x} \end{pmatrix}, \quad \mathbf{W}^{(V)} = (\mathbf{I}_{d_x}, \mathbf{O}_{d_x}), \quad \mathbf{W}^{(K)} = \mathbf{O}, \quad \mathbf{W}^{(Q)} = \mathbf{O}.$$

Then, by direct calculation based on the definition, we have

$$\mathcal{F}^{(SA)} \begin{pmatrix} \mathbf{Z}_1 \\ \mathbf{O} \end{pmatrix} = \begin{pmatrix} \left( 3 \sum_{p=1}^{d_x} a_{p,1} \tilde{\phi}_K(X_{p,1}) \right) \mathbf{1}_{d_x} & \cdots & \left( 3 \sum_{p=1}^{d_x} a_{p,n} \tilde{\phi}_K(X_{p,n}) \right) \mathbf{1}_{d_x} \\ \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) \right) \mathbf{1}_{d_x} & \cdots & \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) \right) \mathbf{1}_{d_x} \end{pmatrix}.$$

Next, we define a generalized feed-forward layer with the following parameters:

$$\mathbf{W}^{(1)} = \begin{pmatrix} \mathbf{I}_{d_x} & \mathbf{O}_{d_x} \\ -\mathbf{I}_{d_x} & \mathbf{O}_{d_x} \\ \mathbf{O}_{d_x} & \mathbf{I}_{d_x} \\ \mathbf{O}_{d_x} & -\mathbf{I}_{d_x} \end{pmatrix}, \quad \mathbf{B}^{(1)} = \mathbf{O},$$

$$\mathbf{W}^{(2)} = \begin{pmatrix} -\mathbf{I}_{d_x} & \mathbf{I}_{d_x} & \mathbf{I}_{d_x} & -\mathbf{I}_{d_x} \\ \mathbf{O}_{d_x} & \mathbf{O}_{d_x} & -\mathbf{I}_{d_x} & \mathbf{I}_{d_x} \end{pmatrix}, \quad \mathbf{B}^{(2)} = \begin{pmatrix} \mathbf{0}_{d_x} & 2\mathbf{1}_{d_x} & \cdots & 2(n-1)\mathbf{1}_{d_x} \\ \mathbf{0}_{d_x} & \mathbf{0}_{d_x} & \cdots & \mathbf{0}_{d_x} \end{pmatrix}.$$

It can then be verified that

$$\mathcal{F}_{2K+3}^{(GFF)} \circ \mathcal{F}^{(SA)} \begin{pmatrix} \mathbf{Z}_1 \\ \mathbf{O} \end{pmatrix} = \begin{pmatrix} \mathbf{Z}_2 \\ \mathbf{O} \end{pmatrix},$$

where we have used the identity  $x = \sigma_R[x] - \sigma_R[-x]$ .

**Step 3:** We construct a generalized feed-forward layer  $\mathcal{F}_{2K+4}^{(GFF)}$  interpolating the outer function  $\mathbf{G}$  in (14) at the  $2^{d_x n K} + 1$  interpolation points

$$3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi_K(X_{p,q}) \in \left\{ \sum_{j=1}^{d_x n K} 2t_j 3^{-j} : (t_1, \dots, t_{d_x n K}) \in \{0, 1\}^{d_x n K} \right\} \cup \{1\}.$$

Denote these points by  $0 =: s_0 < s_1 < \dots < s_{2^{d_x n K}-1} < s_{2^{d_x n K}} := 1$  and fix  $u \in [d_x]$ . For any  $x \in \mathbb{R}$ , we define a scalar function

$$\begin{aligned} \tilde{G}_u(x) &:= G_{u,1}(s_0) + \sum_{j=1}^{2^{d_x n K}} \frac{G_{u,1}(s_j) - G_{u,1}(s_{j-1})}{s_j - s_{j-1}} (\sigma_R[x - s_{j-1}] - \sigma_R[x - s_j]) \\ &+ \frac{G_{u,2}(s_0) - G_{u,1}(s_{2^{d_x n K}})}{s_0 + 2 - s_{2^{d_x n K}}} (\sigma_R[x - s_{2^{d_x n K}}] - \sigma_R[x - (s_0 + 2)]) \\ &+ \sum_{j=1}^{2^{d_x n K}} \frac{G_{u,2}(s_j) - G_{u,2}(s_{j-1})}{s_j - s_{j-1}} (\sigma_R[x - (s_{j-1} + 2)] - \sigma_R[x - (s_j + 2)]) + \dots \\ &+ \frac{G_{u,n}(s_0) - G_{u,n-1}(s_{2^{d_x n K}})}{s_0 + 2 - s_{2^{d_x n K}}} (\sigma_R[x - (s_{2^{d_x n K}} + 2(n-2))] - \sigma_R[x - (s_0 + 2(n-1))]) \\ &+ \sum_{j=1}^{2^{d_x n K}} \frac{G_{u,n}(s_j) - G_{u,n}(s_{j-1})}{s_j - s_{j-1}} (\sigma_R[x - (s_{j-1} + 2(n-1))] - \sigma_R[x - (s_j + 2(n-1))]) \\ &= G_{u,1}(s_0) \\ &+ \sum_{v=1}^n \sum_{j=1}^{2^{d_x n K}} \frac{G_{u,v}(s_j) - G_{u,v}(s_{j-1})}{s_j - s_{j-1}} (\sigma_R[x - (s_{j-1} + 2(v-1))] - \sigma_R[x - (s_j + 2(v-1))]) \\ &+ \sum_{v=2}^n \frac{G_{u,v}(s_0) - G_{u,v-1}(s_{2^{d_x n K}})}{s_0 + 2 - s_{2^{d_x n K}}} (\sigma_R[x - (s_{2^{d_x n K}} + 2(v-2))] - \sigma_R[x - (s_0 + 2(v-1))]). \end{aligned}$$

In other words,  $\tilde{G}_u$  is defined as the piecewise linear interpolation of the points

$$\left\{ (s_j + 2(v-1), G_{u,v}(s_j)) : j = 0, 1, \dots, 2^{d_x n K}, v = 1, \dots, n \right\},$$

with the function being constant outside the interval  $[0, 2n-1]$ . We observe that

- $\tilde{G}_u(s_j + 2(v-1)) = G_{u,v}(s_j)$  for every  $j \in \{0\} \cup [2^{d_x n K}]$ ,  $u \in [d_x]$ ,  $v \in [n]$ ,
- $\|\tilde{G}_u\|_{L^\infty(\mathbb{R})} \leq \max_{v \in [n]} \|G_{u,v}\|_{L^\infty(\mathcal{C})}$ ,
- $\tilde{G}_u \in \mathcal{FNN}_{1,1}(n(2^{d_x n K} + 1), 1)$ .

By stacking the functions  $\tilde{G}_u$  for  $u \in [d_x]$  vertically, we obtain a feed-forward layer, with width at most  $d_x n(2^{d_x n K} + 1) + 2d_x$ , such that

$$\mathcal{F}_{2K+4}^{(GFF)}(\mathbf{Z}) = \begin{pmatrix} \tilde{G}_1(Z_{1,1}) & \tilde{G}_1(Z_{1,2}) & \dots & \tilde{G}_1(Z_{1,n}) \\ \tilde{G}_2(Z_{2,1}) & \tilde{G}_2(Z_{2,2}) & \dots & \tilde{G}_2(Z_{2,n}) \\ \vdots & \vdots & \dots & \vdots \\ \tilde{G}_{d_x}(Z_{d_x,1}) & \tilde{G}_{d_x}(Z_{d_x,2}) & \dots & \tilde{G}_{d_x}(Z_{d_x,n}) \end{pmatrix}.$$

Together with **Step 1** and **Step 2**, we define the overall generalized Transformer network as

$$\begin{aligned} \mathcal{N} &:= \mathcal{E}_{out} \circ \mathcal{F}_{2K+4}^{(GFF)} \circ \mathcal{F}_{2K+3}^{(GFF)} \circ \mathcal{F}^{(GSA)} \circ \mathcal{F}_{2K+2}^{(GFF)} \circ \dots \circ \mathcal{F}_1^{(GFF)} \circ \mathcal{E}_{in} \\ &\in \mathcal{GT}_{d_x, d_x}(D = 4d_x n, H = 1, S = d_x, W = d_x n(2^{d_x n K} + 1) + 2d_x, L = 2K + 4), \end{aligned} \tag{15}$$

where  $\mathcal{E}_{in}$  and  $\mathcal{E}_{out}$  are appropriately chosen to add and remove zeros to match the hidden dimension  $D$ . In particular, we have

$$\begin{aligned}
& \mathcal{N}(\mathbf{X}) \\
&= \mathcal{E}_{out} \circ \mathcal{F}_{2K+4}^{(GFF)} \circ \mathcal{F}_{2K+3}^{(GFF)} \circ \mathcal{F}^{(GSA)} \circ \mathcal{F}_{2K+2}^{(GFF)} \circ \dots \circ \mathcal{F}_1^{(GFF)} \circ \mathcal{E}_{in}(\mathbf{X}) \\
&= \mathcal{E}_{out} \circ \mathcal{F}_{2K+4}^{(GFF)} \circ \mathcal{F}_{2K+3}^{(GFF)} \circ \mathcal{F}^{(GSA)} \left( \begin{array}{c} \mathbf{Z}_1 \\ \mathbf{O} \end{array} \right) \\
&= \mathcal{E}_{out} \circ \mathcal{F}_{2K+4}^{(GFF)} \left( \begin{array}{c} \mathbf{Z}_2 \\ \mathbf{O} \end{array} \right) \\
&= \begin{pmatrix} \tilde{G}_1 \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) \right) & \tilde{G}_1 \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2 \right) & \dots & \tilde{G}_1 \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2(n-1) \right) \\ \tilde{G}_2 \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) \right) & \tilde{G}_2 \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2 \right) & \dots & \tilde{G}_2 \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2(n-1) \right) \\ \vdots & \vdots & & \vdots \\ \tilde{G}_{d_x} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) \right) & \tilde{G}_{d_x} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2 \right) & \dots & \tilde{G}_{d_x} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2(n-1) \right) \end{pmatrix}.
\end{aligned}$$

**Step 4:** We now conduct an error analysis. We have

$$\begin{aligned}
& \|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})}^p \\
&= \int_{[0,1]^{d_x \times n}} \|\mathcal{N}(\mathbf{X}) - \mathbf{F}(\mathbf{X})\|_F^p d\mathbf{X} \\
&\leq \int_{[0,1]^{d_x \times n}} (d_x n)^{\max\{0, \frac{p}{2}-1\}} \sum_{u=1}^{d_x} \sum_{v=1}^n |\mathcal{N}_{u,v}(\mathbf{X}) - F_{u,v}(\mathbf{X})|^p d\mathbf{X} \\
&= (d_x n)^{\max\{0, \frac{p}{2}-1\}} \sum_{u=1}^{d_x} \sum_{v=1}^n \left( \int_{\mathbf{X}: \forall X_{i,j} \in \Omega_K} + \int_{\mathbf{X}: \exists X_{i,j} \notin \Omega_K} \right) |\mathcal{N}_{u,v}(\mathbf{X}) - F_{u,v}(\mathbf{X})|^p d\mathbf{X} \\
&=: (d_x n)^{\max\{0, \frac{p}{2}-1\}} \sum_{u=1}^{d_x} \sum_{v=1}^n (\text{I} + \text{II}).
\end{aligned}$$

To estimate I, using that  $\tilde{\phi}_K(X_{p,q}) = \phi_K(X_{p,q})$  when  $X_{p,q} \in \Omega_K$ ,  $\tilde{G}_u$  interpolates  $G_{u,v}$  by construction, and  $G_{u,v} \in \mathcal{H}_{\frac{\gamma \log 2}{d_x n \log 3}}(\mathcal{C}, 2\sqrt{d_x n} K_{\mathcal{H}})$ , we have

$$\begin{aligned}
& |\mathcal{N}_{u,v}(\mathbf{X}) - F_{u,v}(\mathbf{X})| \\
&= \left| \tilde{G}_u \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \tilde{\phi}_K(X_{p,q}) + 2(v-1) \right) - G_{u,v} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi(X_{p,q}) \right) \right| \\
&= \left| \tilde{G}_u \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi_K(X_{p,q}) + 2(v-1) \right) - G_{u,v} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi(X_{p,q}) \right) \right| \\
&= \left| G_{u,v} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi_K(X_{p,q}) \right) - G_{u,v} \left( 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi(X_{p,q}) \right) \right| \\
&\leq 2(d_x n)^{\frac{1}{2}} K_{\mathcal{H}} \left| 3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} (\phi_K(X_{p,q}) - \phi(X_{p,q})) \right|^{\frac{\gamma \log 2}{d_x n \log 3}} \\
&\leq 2(d_x n)^{\frac{1}{2}} K_{\mathcal{H}} \left| 2 \sum_{q=d_x n K+1}^{\infty} 3^{-q} \right|^{\frac{\gamma \log 2}{d_x n \log 3}}
\end{aligned}$$

$$\begin{aligned}
&\leq 2(d_x n)^{\frac{1}{2}} K_{\mathcal{H}} 3^{-\frac{K\gamma \log 2}{\log 3}} \\
&= 2(d_x n)^{\frac{1}{2}} K_{\mathcal{H}} 2^{-\gamma K},
\end{aligned}$$

where the second inequality follows from the fact that, as indicated in (12),  $3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi(X_{p,q})$  and  $3 \sum_{p=1}^{d_x} \sum_{q=1}^n a_{p,q} \phi_K(X_{p,q})$  are both in the Cantor set  $\mathcal{C}$  and have the same first  $d_x n K$  ternary digits. Thus,

$$\begin{aligned}
\text{I} &= \int_{\mathbf{X}: \forall X_{i,j} \in \Omega_K} |\mathcal{N}_{u,v}(\mathbf{X}) - F_{u,v}(\mathbf{X})|^p d\mathbf{X} \\
&\leq 2^p (d_x n)^{\frac{p}{2}} K_{\mathcal{H}}^p 2^{-p\gamma K}.
\end{aligned}$$

To estimate II, noting that both  $\mathcal{N}_{u,v}$  and  $F_{u,v}$  are bounded, and that  $\Omega_K$  has Lebesgue measure at least  $1 - 2^{-p\gamma K}$ , we obtain

$$\begin{aligned}
\text{II} &= \int_{\mathbf{X}: \exists X_{i,j} \notin \Omega_K} |\mathcal{N}_{u,v}(\mathbf{X}) - F_{u,v}(\mathbf{X})|^p d\mathbf{X} \\
&\leq \int_{\mathbf{X}: \exists X_{i,j} \notin \Omega_K} \left( \|\mathcal{N}_{u,v}\|_{L^\infty([0,1]^{d_x \times n})} + \|F_{u,v}\|_{L^\infty([0,1]^{d_x \times n})} \right)^p d\mathbf{X} \\
&\leq \left( \|\mathcal{N}_{u,v}\|_{L^\infty([0,1]^{d_x \times n})} + \|F_{u,v}\|_{L^\infty([0,1]^{d_x \times n})} \right)^p (1 - (1 - 2^{-p\gamma K})^{d_x n}) \\
&\leq 2^{2p} (d_x n)^{\frac{p}{2}+1} K_{\mathcal{H}}^p 2^{-p\gamma K},
\end{aligned}$$

where we apply Bernoulli's inequality in the last inequality.

We combine the bounds for I and II to obtain

$$\begin{aligned}
\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})}^p &\leq (d_x n)^{\max\{0, \frac{p}{2}-1\}} \sum_{u=1}^{d_x} \sum_{v=1}^n \left( 2^p (d_x n)^{\frac{p}{2}} K_{\mathcal{H}}^p 2^{-p\gamma K} + 2^{2p} (d_x n)^{\frac{p}{2}+1} K_{\mathcal{H}}^p 2^{-p\gamma K} \right) \\
&\leq (d_x n)^{3p} 2^{2p} K_{\mathcal{H}}^p 2^{-p\gamma K},
\end{aligned}$$

where we have used  $p \geq 1$  and  $\max\{a, b\} \leq a + b$  for all  $a, b \geq 0$ , which implies

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})} \leq 4(d_x n)^3 K_{\mathcal{H}} 2^{-\gamma K}.$$

Choose  $K \geq \frac{1}{\gamma} \log_2 \frac{1}{\varepsilon}$  so that  $2^{d_x n K} = \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil$ . Then we have

$$\|\mathcal{N} - \mathbf{F}\|_{L^p([0,1]^{d_x \times n})} \leq 4(d_x n)^3 K_{\mathcal{H}} \varepsilon,$$

and by (15),

$$\mathcal{N} \in \mathcal{GT}_{d_x, d_x}(D = 4d_x n, H = 1, S = d_x, W \leq 3d_x n \lceil \varepsilon^{-\frac{d_x n}{\gamma}} \rceil, L \leq 6 \lceil \frac{1}{\gamma} \log_2 \frac{1}{\varepsilon} \rceil).$$

This completes the proof.  $\square$

## References

- [1] Alekh Agarwal and John C Duchi. The generalization ability of online algorithms for dependent data. *IEEE Transactions on Information Theory*, 59(1):573–587, 2012.
- [2] Martin Anthony and Peter L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999.

- [3] Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20(63):1–17, 2019.
- [4] Peter L Bartlett and Wolfgang Maass. Vapnik-chervonenkis dimension of neural nets. *The handbook of brain theory and neural networks*, pages 1188–1192, 2003.
- [5] Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 2023.
- [6] Minshuo Chen, Haoming Jiang, Wenjing Liao, and Tuo Zhao. Nonparametric regression on low-dimensional manifolds using deep relu networks: Function approximation and statistical recovery. *Information and Inference: A Journal of the IMA*, 11(4):1203–1253, 2022.
- [7] David Chiang and Peter Cholak. Overcoming a theoretical limitation of self-attention. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2022.
- [8] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- [9] Shai Dekel and Dany Leviatan. The bramble–hilbert lemma for convex domains. *SIAM journal on mathematical analysis*, 35(5):1203–1212, 2004.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 4171–4186, 2019.
- [11] Ronald A DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998.
- [12] Ronald A DeVore and George G Lorentz. *Constructive Approximation*, volume 303. Springer Science & Business Media, 1993.
- [13] Zhao Ding, Chenguang Duan, Yuling Jiao, and Jerry Zhijian Yang. Semi-supervised deep sobolev regression: Estimation and variable selection by requ neural network. *IEEE Transactions on Information Theory*, 2025.
- [14] David L Donoho and Iain M Johnstone. Minimax estimation via wavelet shrinkage. *The Annals of Statistics*, 26(3):879–921, 1998.
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [16] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*. PMLR, 2022.
- [17] Lawrence C Evans. *Partial Differential Equations*, volume 19. American Mathematical Soc., 2010.
- [18] Zhiying Fang, Yidong Ouyang, Ding-Xuan Zhou, and Guang Cheng. Attention enables zero approximation error. *arXiv preprint arXiv:2202.12166*, 2022.

- [19] Max H Farrell, Tengyuan Liang, and Sanjog Misra. Deep neural networks for estimation and inference. *Econometrica*, 89(1):181–213, 2021.
- [20] Xingdong Feng, Yuling Jiao, Lican Kang, Baqun Zhang, and Fan Zhou. Over-parameterized deep nonparametric regression for dependent data with its applications to reinforcement learning. *Journal of Machine Learning Research*, 24(383):1–40, 2023.
- [21] Iryna Gurevych, Michael Kohler, and Gözde Gül Şahin. On the rate of convergence of a classifier based on a transformer encoder. *IEEE Transactions on Information Theory*, 68(12):8139–8155, 2022.
- [22] Alexander Havrilla and Wenjing Liao. Understanding scaling laws with statistical and approximation theory for transformer neural networks on intrinsically low-dimensional data. *Advances in Neural Information Processing Systems*, 2024.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [24] Chang hoon Song, Geonho Hwang, Jun ho Lee, and Myungjoo Kang. Minimal width for universal property of deep rnn. *Journal of Machine Learning Research*, 24(121):1–41, 2023.
- [25] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- [26] Daniel J Hsu, Aryeh Kontorovich, and Csaba Szepesvári. Mixing time estimation in reversible markov chains from a single sample path. *Advances in neural information processing systems*, 2015.
- [27] Jerry Yao-Chieh Hu, Wei-Po Wang, Ammar Gilani, Chenyang Li, Zhao Song, and Han Liu. Fundamental limits of prompt tuning transformers: Universality, capacity and efficiency. In *International Conference on Learning Representations*, 2025.
- [28] Haotian Jiang and Qianxiao Li. Approximation rate of the transformer architecture for sequence modeling. *Advances in Neural Information Processing Systems*, 2024.
- [29] Yuling Jiao, Lican Kang, Jin Liu, Xiliang Lu, and Jerry Zhijian Yang. Deep approximate policy iteration. *Annals of Statistics*, 2025.
- [30] Yuling Jiao, Yanming Lai, Yang Wang, and Bokai Yan. Convergence analysis of flow matching in latent space with transformers. *arXiv preprint arXiv:2404.02538*, 2024.
- [31] Yuling Jiao, Guohao Shen, Yuanyuan Lin, and Jian Huang. Deep nonparametric regression on approximate manifolds: Nonasymptotic error bounds with polynomial prefactors. *The Annals of Statistics*, 51(2):691–716, 2023.
- [32] Yuling Jiao, Yang Wang, and Bokai Yan. Approximation bounds for recurrent neural networks with application to regression. *arXiv preprint arXiv:2409.05577*, 2024.
- [33] Tokio Kajitsuka and Issei Sato. Are transformers with one layer self-attention using low-rank weight matrices universal approximators? In *International Conference on Learning Representations*, 2024.
- [34] Tokio Kajitsuka and Issei Sato. On the optimal memorization capacity of transformers. In *International Conference on Learning Representations*, 2025.

- [35] Marek Karpinski and Angus Macintyre. Polynomial bounds for vc dimension of sigmoidal and general pfaffian neural networks. *Journal of Computer and System Sciences*, 54(1):169–176, 1997.
- [36] Junghwan Kim, Michelle Kim, and Barzan Mozafari. Provable memorization capacity of transformers. In *International Conference on Learning Representations*, 2023.
- [37] Michael Kohler and Sophie Langer. On the rate of convergence of fully connected deep neural network regression estimates. *The Annals of Statistics*, 49(4):2231–2249, 2021.
- [38] Vitaly Kuznetsov and Mehryar Mohri. Generalization bounds for non-stationary mixing processes. *Machine Learning*, 106(1):93–117, 2017.
- [39] Jianfeng Lu, Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation for smooth functions. *SIAM Journal on Mathematical Analysis*, 53(5):5465–5506, 2021.
- [40] Ron Meir. Nonparametric time series prediction through adaptive model selection. *Machine Learning*, 39:5–34, 2000.
- [41] Mehryar Mohri and Afshin Rostamizadeh. Rademacher complexity bounds for non-iid processes. In *Advances in Neural Information Processing Systems*, 2008.
- [42] Mehryar Mohri and Afshin Rostamizadeh. Stability bounds for stationary  $\varphi$ -mixing and  $\beta$ -mixing processes. *Journal of Machine Learning Research*, 11(2), 2010.
- [43] Ryumei Nakada and Masaaki Imaizumi. Adaptive approximation and generalization of deep neural network with intrinsic dimensionality. *Journal of Machine Learning Research*, 21(174):1–38, 2020.
- [44] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [45] Sejun Park, Jaeho Lee, Chulhee Yun, and Jinwoo Shin. Provable memorization via deep neural networks using sub-linear parameters. In *Conference on learning theory*. PMLR, 2021.
- [46] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- [47] Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is turing-complete. *Journal of Machine Learning Research*, 22(75):1–35, 2021.
- [48] Aleksandar Petrov, Philip Torr, and Adel Bibi. Prompting a pretrained transformer can be a universal approximator. In *International Conference on Machine Learning*. PMLR, 2024.
- [49] Yinuo Ren, Yiping Lu, Lexing Ying, and Grant M Rotskoff. Statistical spatially inhomogeneous diffusion inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024.
- [50] Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation function. *The Annals of Statistics*, 48(4):1875, 2020.
- [51] Johannes Schmidt-Hieber. The kolmogorov–arnold representation theorem revisited. *Neural Networks*, 137:119–126, 2021.
- [52] Cosma Shalizi and Aryeh Kontorovich. Predictive pac learning and process decompositions. In *Advances in Neural Information Processing Systems*, 2013.

- [53] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation characterized by number of neurons. *Communications in Computational Physics*, 28(5):1768–1811, 2020.
- [54] Jonathan W Siegel. Optimal approximation rates for deep relu neural networks on sobolev and besov spaces. *Journal of Machine Learning Research*, 24(357):1–52, 2023.
- [55] Ingo Steinwart and Andreas Christmann. Fast learning from non-iid observations. In *Advances in Neural Information Processing Systems*, 2009.
- [56] Charles J Stone. Optimal global rates of convergence for nonparametric regression. *The Annals of Statistics*, 10(4):1040–1053, 1982.
- [57] Taiji Suzuki. Adaptivity of deep reLU network for learning in besov and mixed smooth besov spaces: optimal rate and curse of dimensionality. In *International Conference on Learning Representations*, 2019.
- [58] Shokichi Takakura and Taiji Suzuki. Approximation and estimation ability of transformers for sequence-to-sequence functions with infinite dimensional input. In *International Conference on Machine Learning*. PMLR, 2023.
- [59] Naoki Takeshita and Masaaki Imaizumi. Approximation of permutation invariant polynomials by transformers: Efficient construction in column-size. *arXiv preprint arXiv:2502.11467*, 2025.
- [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- [61] Gabrielle Viennet. Inequalities for absolutely regular sequences: application to density estimation. *Probability Theory and Related Fields*, 107:467–492, 1997.
- [62] Mingze Wang and Weinan E. Understanding the expressive power and mechanisms of transformer for sequence modeling. In *Advances in Neural Information Processing Systems*, 2024.
- [63] Kevin Xu and Issei Sato. On expressive power of looped transformers: Theoretical analysis and enhancement via timestep encoding. *arXiv preprint arXiv:2410.01405*, 2024.
- [64] Yunfei Yang and Ding-Xuan Zhou. Nonparametric regression using over-parameterized shallow relu neural networks. *Journal of Machine Learning Research*, 25(165):1–35, 2024.
- [65] Yunfei Yang and Ding-Xuan Zhou. Optimal rates of approximation by shallow relu k neural networks and applications to nonparametric regression. *Constructive Approximation*, pages 1–32, 2024.
- [66] Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114, 2017.
- [67] Bin Yu. Density estimation in the  $l^\infty$  norm for dependent data with applications to the gibbs sampler. *The Annals of Statistics*, 21(2):711–735, 1993.
- [68] Bin Yu. Rates of convergence for empirical processes of stationary mixing sequences. *The Annals of Probability*, 22(1):94–116, 1994.
- [69] Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? In *International Conference on Learning Representations*, 2019.

- [70] Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar.  $O(n)$  connections are expressive enough: Universal approximability of sparse transformers. *Advances in Neural Information Processing Systems*, 2020.
- [71] Shijun Zhang, Zuowei Shen, and Haizhao Yang. Deep network approximation: Achieving arbitrary accuracy with fixed number of neurons. *Journal of Machine Learning Research*, 23(276):1–60, 2022.