

Vocabulary In-Context Learning in Transformers: Benefits of Positional Encoding

Anonymous Authors¹

Abstract

Numerous studies have demonstrated that the Transformer architecture possesses the capability for in-context learning (ICL). In scenarios involving function approximation, context can serve as a control parameter for the model, endowing it with the universal approximation property (UAP). In practice, context is represented by tokens from a finite set, referred to as a vocabulary, which is the case considered in this paper, *i.e.*, vocabulary in-context learning (VICL). We demonstrate that VICL in single-layer Transformers, without positional encoding, does not possess the UAP; however, it is possible to achieve the UAP when positional encoding is included. Several sufficient conditions for the positional encoding are provided. Our findings reveal the benefits of positional encoding from an approximation theory perspective in the context of ICL.

1. Introduction

Transformers have emerged as a dominant architecture in deep learning over the past few years. Thanks to their remarkable performance in language tasks, they have become the preferred framework in the natural language processing (NLP) field. A major trend in modern NLP is the development and integration of various black-box models, along with the construction of extensive text datasets. In addition, improving model performance in specific tasks through techniques such as in-context learning (ICL) (Dong et al., 2024; Brown et al., 2020), chain of thought (CoT) (Wei et al., 2022b; Chu et al., 2024), and retrieval-augmented generation (RAG) (Gao et al., 2024) has become a significant research focus. While the practical success of these models and techniques is well-documented, the theoretical understanding of why they perform so well remains incomplete.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

To explore the capabilities of Transformers in handling ICL tasks, it is essential to examine their approximation power. The universal approximation property (UAP) (Cybenko, 1989; Hornik et al., 1989; Hornik, 1991; Leshno et al., 1993) has long been a key topic in the theoretical study of neural networks (NNs), with much of the focus historically on feed-forward neural networks (FNNs). (Yun et al., 2020) was the first to investigate the UAP of Transformers, demonstrating that any sequence-to-sequence function could be approximated by a Transformer network with fixed positional encoding. (Luo et al., 2022) highlighted that a Transformer with relative positional encoding does not possess the UAP. Meanwhile, (Petrov et al., 2024b) explored the role of prompting in Transformers, proving that prompting a pre-trained Transformer can act as a universal functional approximator.

However, one limitation of these studies is that, in practical scenarios, the inputs to language models are derived from a finite set embedded in high-dimensional Euclidean space—commonly referred to as a vocabulary. Whether examining the work on prompts in (Petrov et al., 2024b) or the research on ICL in (Ahn et al., 2024; Cheng et al., 2024), these studies assume inputs from the entire Euclidean space, which differs significantly from the discrete nature of vocabularies used in real-world applications.

1.1. Contributions

Starting with the connection between FNNs and Transformers, we turn to the finite restriction of vocabularies and study the benefits of positional encoding. Leveraging the UAP of FNNs, we explore the approximation properties of Transformers for ICL tasks in two scenarios: one where the inputs are from the entire Euclidean space, and the other where the inputs are from a finite vocabulary.

1. Without the restriction of a finite vocabulary, we establish a connection between FNNs and Transformers in processing ICL tasks, as demonstrated in Lemma 3. Using this lemma, we show that Transformers can function as universal approximators (Lemma 4), where the context serves as control parameters, while the weights and biases of the Transformer remain fixed.

2. When the vocabulary is finite and positional encoding is not used, we prove that single-layer Transformers cannot achieve the UAP for ICL tasks (Theorem 7). However, when the vocabulary is finite and positional encoding is used, it becomes possible for single-layer Transformers to achieve the UAP (Theorem 9). In particular, for Transformers with ReLU activation functions, the conditions on the positional encoding are discussed (Theorem 10).

1.2. Related Works

Universal approximation property. NNs, through multi-layer nonlinear transformations and feature extraction, are capable of learning deep feature representations from raw data. From the early FNNs (Rosenblatt, 1958), to later advancements like recurrent neural networks (RNNs) (Waibel et al., 1989; Hochreiter & Schmidhuber, 1997), convolutional neural networks (CNNs) (Waibel et al., 1989; Lecun et al., 1998), and residual neural networks (ResNets) (He et al., 2016), remarkable progress has been made. As the application of NNs becomes more widespread, efforts have been directed toward understanding the theoretical foundations behind their effectiveness, particularly through the UAP of NNs. Research on the UAP of NNs generally falls into two categories: the first considers networks with any number of neurons in each layer but a fixed number of layers (Cybenko, 1989; Hornik et al., 1989; Hornik, 1991; Leshno et al., 1993), while the second examines networks with an arbitrary number of layers but a finite number of neurons in each layer (Lu et al., 2017; Park et al., 2021; Cai, 2023; Li et al., 2024). Since our study builds on existing results regarding the approximation capabilities of FNNs, we focus on investigating the approximation abilities of single-layer Transformers in modulating context for ICL tasks. Consequently, our work relies more on the findings from the first category of research. The realization of the UAP depends on the architecture of the network itself, providing constructive insights for exploring the connection between FNNs and Transformers, and offering valuable guidance for our study. Recently, (Petrov et al., 2024b) also explored UAP in the context of in-context learning, but without considering vocabulary constraints or positional encodings.

Transformers. The Transformer is a widely used neural network architecture for modeling sequences (Vaswani et al., 2017; Devlin et al., 2019; Yang et al., 2019; Raffel et al., 2020; Zhenzhong et al., 2021; Liu et al., 2020). This non-recurrent architecture relies entirely on the attention mechanism to capture global dependencies between inputs and outputs (Vaswani et al., 2017). The highly effective neural sequence transduction model is typically structured using an encoder-decoder framework (Bahdanau et al., 2014; Sutskever et al., 2014). The encoder maps the input sequence X into a continuous representation S , from

which the decoder generates the output sequence Y . In the Transformer, both the encoder and decoder are composed of stacked self-attention layers and fully connected layers.

For simplicity, we describe the Transformer using a simplified self-attention sequence encoder. Without positional encoding, the Transformer can be viewed as a stack of N blocks, each consisting of a self-attention layer followed by a feed-forward layer with skip connections. In this paper, we focus on the case of a single-layer self-attention sequence encoder.

In-context learning. The Transformer has demonstrated remarkable performance in the field of NLP, and large language models (LLMs) are gaining increasing popularity. ICL has emerged as a new paradigm in NLP, enabling LLMs to make better predictions through prompts provided within the context (Brown et al., 2020; Chowdhery et al., 2023; Touvron et al., 2023; OpenAI et al., 2024; Xun et al., 2017). We chose ICL as the focus of our research primarily due to its wide range of applications and superior performance, which motivated us to explore its underlying theoretical foundations. ICL delivers high performance with high-quality data at a lower cost (Wang et al., 2021b; Khorashadizadeh et al., 2023; Ding et al., 2023). It enhances retrieval-augmented methods by prepending grounding documents to the input (Ram et al., 2023) and can effectively update or refine the model’s knowledge base through well-designed prompts (De Cao et al., 2021).

Positional Encoding. The following explanation clarifies the significance of incorporating positional encoding into the Transformer architecture. RNNs capture sequential order by encoding the changes in hidden states over time. In contrast, for Transformers, the self-attention mechanism is permutation equivariant, meaning that for any model f , any permutation matrix π , and any input x , the following holds: $f(\pi(x)) = \pi(f(x))$.

We aim to explore the impact of positional encoding on the performance of a single-layer Transformer when performing ICL tasks with a finite vocabulary. Therefore, we focus on analyzing existing positional encoding methods. There are two fundamental methods for encoding positional information in a sequence within the Transformer: absolute positional encodings (APEs) *e.g.* (He et al., 2021; Liu et al., 2020; Wang et al., 2021a; Ke et al., 2021), relative positional encodings (RPEs) *e.g.* (Shaw et al., 2018; Dai et al., 2019; Ke et al., 2021) and rotary positional embedding (RoPE) (Su et al., 2024). The commonly used APE is implemented by directly adding the positional encodings to the word embeddings, and we follow this implementation.

UAP of ICL. Regarding the understanding of the mechanism of ICL, various explanations have been proposed, including those based on Bayesian theory (Xie et al., 2022;

Wang et al., 2024) and gradient descent theory (Dai et al., 2023). Fine-tuning the Transformer through ICL alters the presentation of the input rather than the model parameters, which is driven by successful few-shot and zero-shot learning (Wei et al., 2022a; Kojima et al., 2022). This success raises the question of whether we can achieve the UAP through context adjustment.

(Yun et al., 2020) demonstrated that Transformers can serve as universal sequence-to-sequence approximators, while (Alberti et al., 2023) extended the UAP to architectures with non-standard attention mechanisms. These works represent significant efforts in enabling Transformers to achieve sequence-to-sequence approximation; however, their implementations allow the internal parameters of the Transformers to vary, which does not fully reflect the characteristics of ICL. In contrast, (Likhoshesterov et al., 2021) showed that while the parameters of self-attention remain fixed, various sparse matrices can be approximated by altering the inputs. Fixing self-attention parameters aligns more closely with practical scenarios and provides valuable insights for our work. However, this approach has the limitation of excluding the full Transformer architecture. Furthermore, (Deora et al., 2024) illustrated the convergence and generalization of single-layer multi-head self-attention models trained using gradient descent, supporting the feasibility of our research by emphasizing the robust generalization of Transformers. Nevertheless, (Petrov et al., 2024a) indicated that the presence of a prefix does not alter the attention focus within the context, prompting us to explore variations in input context and introduce flexibility in positional encoding.

1.3. Outline

We will introduce the notations and background results in Section 2. Section 3 addresses the case where the vocabulary is finite and positional encoding is not used. Section 4 discusses the benefits of using positional encoding. A summary is provided in Section 5. All proof of lemmas and theorems are provided in Appendix.

2. Background Materials

We consider the approximation problem as follows. For a target continuous function $f : \mathcal{K} \rightarrow \mathbb{R}^{d_y}$ with a compact domain $\mathcal{K} \subset \mathbb{R}^{d_x}$, we aim to adjust the content of the context so that the output of the Transformer network can approximate f . First, we present the concrete forms and notations for the inputs of ICL, FNNs, and Transformers.

2.1. Notations

Input of in-context learning. In the ICL task, the given n demonstrations are denoted as $z^{(i)} = (x^{(i)}, y^{(i)})$ for $i = 1, 2, \dots, n$, where $x^{(i)} \in \mathbb{R}^{d_x}$ and $y^{(i)} \in \mathbb{R}^{d_y}$. Unlike the

setting in (Ahn et al., 2024) and (Cheng et al., 2024) where $y^{(i)}$ was related to $x^{(i)}$ (for example $y^{(i)} = \phi(x^{(i)})$ for some function ϕ), in this paper, we do not assume any correspondence between $x^{(i)}$ and $y^{(i)}$, i.e., $x^{(i)}$ and $y^{(i)}$ are chosen freely. To predict the target at a query vector $x \in \mathbb{R}^{d_x}$ or $z = (x, 0) \in \mathbb{R}^{d_x+d_y}$, we define the following matrix Z as the input:

$$Z = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(n)} & z \end{bmatrix} \\ := \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} & x \\ y^{(1)} & y^{(2)} & \dots & y^{(n)} & 0 \end{bmatrix} \in \mathbb{R}^{(d_x+d_y) \times (n+1)}. \quad (1)$$

Furthermore, let $\mathcal{P} : \mathbb{N}^+ \rightarrow \mathbb{R}^{d_x+d_y}$ represent a positional encoding function, and define $\mathcal{P}^{(i)} := \mathcal{P}(i)$. Denote the demonstrations with positional encoding as $z_{\mathcal{P}}^{(i)} = z^{(i)} + \mathcal{P}^{(i)}$ and $z_{\mathcal{P}} = z + \mathcal{P}^{(n+1)}$. The context with positional encoding can then be represented as:

$$Z_{\mathcal{P}} = \begin{bmatrix} z_{\mathcal{P}}^{(1)} & z_{\mathcal{P}}^{(2)} & \dots & z_{\mathcal{P}}^{(n)} & z_{\mathcal{P}} \end{bmatrix} \\ := \begin{bmatrix} x_{\mathcal{P}}^{(1)} & x_{\mathcal{P}}^{(2)} & \dots & x_{\mathcal{P}}^{(n)} & x_{\mathcal{P}} \\ y_{\mathcal{P}}^{(1)} & y_{\mathcal{P}}^{(2)} & \dots & y_{\mathcal{P}}^{(n)} & y_{\mathcal{P}} \end{bmatrix} \in \mathbb{R}^{(d_x+d_y) \times (n+1)}. \quad (2)$$

Here, the vectors $x_{\mathcal{P}}^{(i)}$ and $y_{\mathcal{P}}^{(i)}$ represent the corresponding components of $z_{\mathcal{P}}^{(i)}$. Additionally, we denote:

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} \end{bmatrix} \in \mathbb{R}^{d_x \times n}, \quad (3)$$

$$X_{\mathcal{P}} = \begin{bmatrix} x_{\mathcal{P}}^{(1)} & x_{\mathcal{P}}^{(2)} & \dots & x_{\mathcal{P}}^{(n)} \end{bmatrix} \in \mathbb{R}^{d_x \times n}, \quad (4)$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(n)} \end{bmatrix} \in \mathbb{R}^{d_y \times n}, \quad (5)$$

$$Y_{\mathcal{P}} = \begin{bmatrix} y_{\mathcal{P}}^{(1)} & y_{\mathcal{P}}^{(2)} & \dots & y_{\mathcal{P}}^{(n)} \end{bmatrix} \in \mathbb{R}^{d_y \times n}. \quad (6)$$

Feed-forward neural networks. One-hidden-layer FNNs have sufficient capacity to approximate continuous functions on any compact domain. In this article, all the FNNs we refer to and use are one-hidden-layer networks. We denote a one-hidden-layer FNN with activation function σ as N^{σ} , and the set of all such networks is denoted as \mathcal{N}^{σ} , i.e.,

$$\mathcal{N}^{\sigma} = \left\{ N^{\sigma} := A \sigma(Wx + b) \mid \right. \\ \left. A \in \mathbb{R}^{d_y \times k}, W \in \mathbb{R}^{k \times d_x}, b \in \mathbb{R}^k, k \in \mathbb{N}^+ \right\} \\ = \left\{ N^{\sigma} := \sum_{i=1}^k a_i \sigma(w_i \cdot x + b_i) \mid \right. \\ \left. (a_i, w_i, b_i) \in \mathbb{R}^{d_y} \times \mathbb{R}^{d_x} \times \mathbb{R}, k \in \mathbb{N}^+ \right\}. \quad (7)$$

For elementwise activations, such as ReLU, the above notation is well-defined. However, if the activation function is not elementwise, especially in the case of softmax activation,

we need to give more details for the notation:

$$\mathcal{N}^{\text{softmax}} = \left\{ \mathcal{N}^{\text{softmax}} = \frac{\sum_{i=1}^k a_i e^{w_i \cdot x + b_i}}{\sum_{i=1}^k e^{w_i \cdot x + b_i}} \mid (a_i, w_i, b_i) \in \mathbb{R}^{d_y} \times \mathbb{R}^{d_x} \times \mathbb{R}, k \in \mathbb{N}^+ \right\}. \quad (8)$$

Transformers. We define the general attention mechanism following (Ahn et al., 2024; Cheng et al., 2024) as:

$$\text{Attn}_{Q,K,V}^\sigma(Z) := V Z M \sigma((QZ)^\top K Z), \quad (9)$$

where V, Q, K are the value, query, and key matrices in $\mathbb{R}^{(d_x+d_y) \times (d_x+d_y)}$, respectively, $M = \text{diag}(I_n, 0)$ is the mask matrix in $\mathbb{R}^{(n+1) \times (n+1)}$, and σ is the activation function. Here the softmax activation of a matrix $G \in \mathbb{R}^{m \times n}$ is defined as:

$$\text{softmax}(G) := \left[\frac{\exp(G_{i,j})}{\sum_{l=1}^m \exp(G_{l,j})} \right]_{i,j}. \quad (10)$$

With this formulation of the general attention mechanism, we can define a single-layer Transformer without positional encoding as:

$$\begin{aligned} \mathcal{T}^\sigma(x; X, Y) := \\ (Z + V Z M \sigma((QZ)^\top K Z))_{d_x+1:d_x+d_y, n+1}, \end{aligned} \quad (11)$$

where $[a : b, c : d]$ denotes the submatrix from the a -th row to the b -th row and from the c -th column to the d -th column. If $a = b$ (or $c = d$), the row (or column) index is reduced to a single number. Similarly to the notation for FNNs, \mathcal{T}^σ denotes the set of all \mathcal{T}^σ with different parameters.

Vocabulary. In the above notations, the parameters are general and unrestricted. When we refer to a ‘‘vocabulary’’, we mean that the parameters are drawn from a finite set. For networks and their corresponding sets, we use the subscript $*$ to indicate the use of a vocabulary \mathcal{V} .

In the context of ICL, we refer to it as vocabulary ICL if all input vectors $z^{(i)}$ come from a finite vocabulary $\mathcal{V} = \mathcal{V}_x \times \mathcal{V}_y \subset \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$. In this case, we use $\mathcal{T}_*^\sigma(x; X, Y)$ to represent the Transformer $\mathcal{T}^\sigma(x; X, Y)$ defined in equation (11), and denote the set of such Transformers as \mathcal{T}_*^σ :

$$\mathcal{T}_*^\sigma = \left\{ \mathcal{T}_*^\sigma(x; X, Y) := \mathcal{T}^\sigma(x; X, Y) \mid z^{(i)} \in \mathcal{V}, i \in \{1, 2, \dots, n\}, n \in \mathbb{N}^+ \right\}. \quad (12)$$

When positional encoding \mathcal{P} is involved, we add the subscript \mathcal{P} , i.e.,

$$\mathcal{T}_{*,\mathcal{P}}^\sigma = \left\{ \mathcal{T}_{*,\mathcal{P}}^\sigma(x; X, Y) := \mathcal{T}^\sigma(x; X_{\mathcal{P}}, Y_{\mathcal{P}}) \mid z^{(i)} \in \mathcal{V}, i \in \{1, 2, \dots, n\}, n \in \mathbb{N}^+ \right\}. \quad (13)$$

Note that the context length n in \mathcal{T}^σ , \mathcal{T}_*^σ , and $\mathcal{T}_{*,\mathcal{P}}^\sigma$ are unbounded.

For FNNs, we denote a network with a finite set of weights as \mathcal{N}_*^σ , and the corresponding set of such networks as \mathcal{N}_*^σ . When the activation function is an elementwise activation:

$$\mathcal{N}_*^\sigma = \left\{ \mathcal{N}_*^\sigma := \sum_{i=1}^k a_i \sigma(w_i \cdot x + b_i) \mid (a_i, w_i, b_i) \in \mathcal{A} \times \mathcal{W} \times \mathcal{B}, k \in \mathbb{N}^+ \right\}. \quad (14)$$

where $\mathcal{A} \subset \mathbb{R}^{d_y}$, $\mathcal{W} \subset \mathbb{R}^{d_x}$, and $\mathcal{B} \subset \mathbb{R}$ are finite sets. When the activation function is softmax:

$$\mathcal{N}^{\text{softmax}} = \left\{ \mathcal{N}^{\text{softmax}} = \frac{\sum_{i=1}^k a_i e^{w_i \cdot x + b_i}}{\sum_{i=1}^k e^{w_i \cdot x + b_i}} \mid (a_i, w_i, b_i) \in \mathcal{A} \times \mathcal{W} \times \mathcal{B}, k \in \mathbb{N}^+ \right\} \quad (15)$$

where \mathcal{A}, \mathcal{W} , and \mathcal{B} are defined as in the previous context. To simplify calculations and expressions, we introduce the following assumptions throughout the remainder of the article similar to the setting in (Cheng et al., 2024).

Assumption 1. The matrices $Q, K, V \in \mathbb{R}^{(d_x+d_y) \times (d_x+d_y)}$ have the following sparse partition:

$$Q = \begin{bmatrix} B & 0 \\ 0 & 0 \end{bmatrix}, \quad K = \begin{bmatrix} C & 0 \\ 0 & 0 \end{bmatrix}, \quad V = \begin{bmatrix} D & E \\ F & U \end{bmatrix}, \quad (16)$$

where $B, C, D \in \mathbb{R}^{d_x \times d_x}$, $E \in \mathbb{R}^{d_x \times d_y}$, $F \in \mathbb{R}^{d_y \times d_x}$ and $U \in \mathbb{R}^{d_y \times d_y}$. Furthermore, the matrices B, C and U are non-singular, and the matrix $F = 0$.

In addition, we assume the elementwise activation σ is non-polynomial, locally bounded, and continuous.

We present all our notations in Table 1 in Appendix A for easy reference.

2.2. Universal Approximation Property

The vanilla form of the universal approximation property for FNNs plays a crucial role in our study. We state it in the following lemma:

Lemma 2 (UAP of FNNs (Leshno et al., 1993)). *Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a non-polynomial, locally bounded, piecewise continuous activation function. For any continuous function $f : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ defined on a compact domain \mathcal{K} , and for any $\varepsilon > 0$, there exist $k \in \mathbb{N}$, $A \in \mathbb{R}^{d_y \times k}$, $b \in \mathbb{R}^k$, and $W \in \mathbb{R}^{k \times d_x}$ such that*

$$\|A\sigma(Wx + b) - f(x)\| < \varepsilon, \quad \forall x \in \mathcal{K}. \quad (17)$$

The theorem presented above is well-known and primarily applies to activation functions operating elementwise. However, it can be readily extended to the case of the softmax

activation function. In fact, this can be achieved using neural networks with exponential activation functions. The specific approach for this generalization is detailed in Appendix B.1.

2.3. Feed-forward neural networks and Transformers

It is important to emphasize the connection between FNNs and Transformers.

Lemma 3. *Let σ be an elementwise activation and T^σ be a single-layer Transformer. For any one-hidden-layer network $N^\sigma : \mathbb{R}^{d_x-1} \rightarrow \mathbb{R}^{d_y} \in \mathcal{N}^\sigma$ with n hidden neurons, there exist matrices $X \in \mathbb{R}^{d_x \times n}$ and $Y \in \mathbb{R}^{d_y \times n}$ such that*

$$T^\sigma(\tilde{x}; X, Y) = N^\sigma(x), \quad \forall x \in \mathbb{R}^{d_x-1}. \quad (18)$$

There is a difference in the input dimensions of T^σ and N^σ , as the latter includes a bias dimension absent in the former. To connect the two inputs, \tilde{x} and x , we use a tilde, where \tilde{x} is formed by augmenting x with an additional one appended to the end.

By employing the structure of query, key, and value matrices in (16), the output forms of the Transformer $T^\sigma(\tilde{x}; X, Y)$ can be simplified as follows:

$$\begin{aligned} & T^\sigma(\tilde{x}; X, Y) \\ &= \left(\begin{bmatrix} X & \tilde{x} \\ Y & 0 \end{bmatrix} + \begin{bmatrix} DX + EY & 0 \\ FX + UY & 0 \end{bmatrix} \right) \\ & \quad \sigma \left(\begin{bmatrix} X^\top B^\top CX & X^\top B^\top C\tilde{x} \\ \tilde{x}^\top B^\top CX & \tilde{x}^\top B^\top C\tilde{x} \end{bmatrix} \right)_{d_x+1:d_x+d_y, n+1} \\ &= (FX + UY)\sigma(X^\top B^\top C\tilde{x}) \\ &= UY\sigma(X^\top B^\top C\tilde{x}). \end{aligned} \quad (19)$$

Comparing this with the output form of FNNs, $N^\sigma(x) = A\sigma(Wx + b)$, it becomes evident that setting $X = (C^\top B)^{-1} [W \quad b]^\top$ and $Y = U^{-1}A$ is sufficient to finish the proof.

It can be observed that the form in equation (19) exhibits the structure of an FNN. Consequently, Lemma 3 implies that single-layer Transformers T^σ with in-context learning and FNNs N^σ are equivalent. However, this equivalence does not hold for the case of softmax activation due to differences in the normalization operations between FNNs and Transformers. Therefore, in the subsequent sections of this article, we employ different analytical methods to address the two types of activation functions.

Moreover, the equivalence in equation (38) suggests that the context in Transformers can act as a control parameter for the model, thereby endowing it with the universal approximation property. This offers a novel perspective on the parameterization of FNNs.

2.4. Universal Approximation Property of In-context Learning

We now present the UAP of Transformers in the context of ICL.

Lemma 4. *Let T^σ be a single-layer Transformer with elementwise or softmax activation, and \mathcal{K} be a compact domain in \mathbb{R}^{d_x-1} . Then for any continuous function $f : \mathcal{K} \rightarrow \mathbb{R}^{d_y}$ and any $\varepsilon > 0$, there exist matrices $X \in \mathbb{R}^{d_x \times n}$ and $Y \in \mathbb{R}^{d_y \times n}$ such that*

$$\|T^\sigma(\tilde{x}; X, Y) - f(x)\| < \varepsilon, \quad \forall x \in \mathcal{K}. \quad (20)$$

For the case of elementwise activation, the result follows directly by combining Lemma 2 and Lemma 3. However, for the softmax activation, the normalization operation requires an additional technique in the proof. The key idea is to consider an FNN with the exponential function as its activation and introduce an additional neuron to account for the normalization effect. Detailed proofs are provided in Appendix B. Similar results have also been reported in recent work (Petrov et al., 2024b), albeit using different techniques.

3. The Non-Universal Approximation Property of \mathcal{N}_*^σ and \mathcal{T}_*^σ

One key aspect of ICL is that the context can act as a control parameter for the model. We now consider the case where the context is restricted to a finite vocabulary. A natural question arises: can a single-layer Transformer with a finite vocabulary, $T_*^\sigma \in \mathcal{T}_*^\sigma$, still achieve the UAP? Given the established connection between FNNs and Transformers, we first analyze $N_*^\sigma \in \mathcal{N}_*^\sigma$ for simplicity.

The answer is that \mathcal{N}_*^σ cannot achieve the UAP because the parameters can only take on a finite number of values. For elementwise activations, the span of \mathcal{N}_*^σ , $\text{span}\{\mathcal{N}_*^\sigma\}$, forms a finite-dimensional function space. According to results from functional analysis, $\text{span}\{\mathcal{N}_*^\sigma\}$ is closed under the function norm (see e.g. Theorem 1.21 of (Rudin, 1991) or Corollary C.4 of (Cannarsa & D’Aprile, 2015)). This implies that the set of functions approximable by $\text{span}\{\mathcal{N}_*^\sigma\}$ is precisely the set of functions within $\text{span}\{\mathcal{N}_*^\sigma\}$. Consequently, any function not in $\text{span}\{\mathcal{N}_*^\sigma\}$ cannot be arbitrarily approximated, meaning that the UAP cannot be achieved.

For softmax networks, the normalization operation introduces further limitations. Even though N_*^{softmax} consists of weighted units drawn from a fixed finite collection of basic units, normalization prevents these networks from being simple linear combinations of one another. While the span of $\mathcal{N}_*^{\text{softmax}}$ might theoretically have infinite dimensionality, its expressive power remains constrained.

To better understand the behavior of functions within

$\mathcal{N}_*^{\text{softmax}}$, we present the following proposition as an introduction.

Proposition 5. *The scalar function $h_k(x) = \sum_{i=1}^k a_i e^{b_i x}$, where $a_i, b_i, x \in \mathbb{R}$ and at least one a_i is nonzero, has at most $k - 1$ zero points.*

The function $h_k(x)$ is commonly referred to as a sum of exponentials. Proposition 5 establishes the maximum number of zero points for this class of functions. The result can be proved using mathematical induction. The cases for $k = 1$ and $k = 2$ are straightforward. Assuming the proposition holds for $k = N$, we proceed with a proof by contradiction for $k = N + 1$. Assume $a_{N+1} \neq 0$ and h_{N+1} has $N + 1$ zero points. We can define a new function g that shares the same zero points as h_{N+1} , given by

$$g(x) = \frac{h_{N+1}(x)}{a_{N+1}e^{b_{N+1}x}} = 1 + \sum_{i=1}^N \frac{a_i}{a_{N+1}} e^{(b_i - b_{N+1})x}. \quad (21)$$

The derivative of g is the sum of N exponentials. By applying the intermediate value theorem, we show that if the number of zero points exceeds N , it leads to a contradiction.

In the proof of Lemma 6, we demonstrated through Proposition 5 that the number of zeros of $\mathcal{N}_*^{\text{softmax}}$ depends solely on a finite set of parameters and constitutes a bounded quantity. Functions can be explicitly constructed whose number of zeros exceeds this bound, thereby preventing their approximation within $\mathcal{N}_*^{\text{softmax}}$.

Now we can summarize the non-universal approximation property of \mathcal{N}_*^σ in the following lemma.

Lemma 6. *The function class \mathcal{N}_*^σ , with elementwise or softmax activation σ , cannot achieve the UAP. Specifically, for any compact domain $\mathcal{K} \subset \mathbb{R}^{d_x}$, there exists a continuous function $f : \mathcal{K} \rightarrow \mathbb{R}^{d_y}$ and $\varepsilon_0 > 0$ such that $\max_{x \in \mathcal{K}} \|f(x) - \mathcal{N}_*^\sigma(\tilde{x})\| \geq \varepsilon_0$ for all $\mathcal{N}_*^\sigma \in \mathcal{N}_*^\sigma$.*

By leveraging the connection between FNNs and Transformers, we establish Theorem 7.

Theorem 7. *The function class \mathcal{T}_*^σ , with elementwise or softmax activation σ , cannot achieve the UAP. Specifically, for any compact domain $\mathcal{K} \subset \mathbb{R}^{d_x-1}$, there exists a continuous function $f : \mathcal{K} \rightarrow \mathbb{R}^{d_y}$ and $\varepsilon_0 > 0$ such that*

$$\max_{x \in \mathcal{K}} \|f(x) - \mathcal{T}_*^\sigma(\tilde{x})\| \geq \varepsilon_0, \quad \forall \mathcal{T}_*^\sigma \in \mathcal{T}_*^\sigma. \quad (22)$$

The result for elementwise activations follows directly from the application of Lemma 3 and Lemma 6. However, the case of the softmax activation is more intricate, as it requires additional techniques to account for the normalization effect. The proof, which utilizes Proposition 5 once again, is presented in the Appendix C.

It is worth noting that Theorem 7 holds even without imposing any constraints on the value, query, and key matrices, V , Q , and K (e.g., the sparse partition described in equation (16)). For further details, refer to Appendix E.

4. Universal Approximation Property of $\mathcal{T}_{*,\mathcal{P}}^\sigma$

After establishing that neither \mathcal{N}_*^σ nor \mathcal{T}_*^σ can achieve the UAP, we aim to leverage a key feature of Transformers: their ability to incorporate absolute positional encodings during token input. This motivates us to investigate whether $\mathcal{T}_{*,\mathcal{P}}^\sigma$ can realize the UAP.

To facilitate our constructive proof, we introduce Lemma 8 as an auxiliary tool to support the main theorem.

Lemma 8 (Kronecker Approximation Theorem (see e.g. (Apostol, 1989))). *Given real n -tuples $\alpha^{(i)} = (\alpha_1^{(i)}, \alpha_2^{(i)}, \dots, \alpha_n^{(i)}) \in \mathbb{R}^n$ for $i = 1, \dots, m$ and $\beta = (\beta_1, \beta_2, \dots, \beta_n) \in \mathbb{R}^n$, the following condition holds: for any $\varepsilon > 0$, there exist $q_i, l_i \in \mathbb{Z}$ such that*

$$\left\| \beta_j - \sum_{i=1}^m q_i \alpha_j^{(i)} + l_j \right\| < \varepsilon, \quad 1 \leq j \leq n, \quad (23)$$

if and only if for any $r_1, \dots, r_n \in \mathbb{Z}, i = 1, \dots, m$ with

$$\sum_{j=1}^n \alpha_j^{(i)} r_j \in \mathbb{Z}, \quad i = 1, \dots, m, \quad (24)$$

the number $\sum_{j=1}^n \beta_j r_j$ is also an integer. In the case of $m = 1$

and $n = 1$, for any $\alpha, \beta, \varepsilon \in \mathbb{R}$ with α irrational and $\varepsilon > 0$, there exist integers l and q with $q > 0$ such that $|\beta - q\alpha + l| < \varepsilon$.

Lemma 8 indicates that if the condition in equation (24) is satisfied only when all r_i are zeros, then the set $\{Mq + l \mid q \in \mathbb{Z}^m, l \in \mathbb{R}^n\}$ is dense in \mathbb{R}^n , where the matrix $M \in \mathbb{R}^{n \times m}$ is assembled with vectors $\alpha^{(i)}$, i.e., $M = [\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(m)}]$. In the case of $m = n = 1$, let $\alpha = \sqrt{2}$. Then, Lemma 8 implies that the set $\{q\sqrt{2} \pm l \mid l \in \mathbb{N}^+, q \in \mathbb{N}^+\}$ is dense in \mathbb{R} . We will build upon this result to prove one of the most significant theorems in this article.

Theorem 9. *Let $\mathcal{T}_{*,\mathcal{P}}^\sigma$ be the class of functions $\mathcal{T}_{*,\mathcal{P}}^\sigma$, where σ is an elementwise activation, the subscript refers the finite vocabulary $\mathcal{V} = \mathcal{V}_x \times \mathcal{V}_y$, $\mathcal{P} = \mathcal{P}_x \times \mathcal{P}_y$ represents the positional encoding map, and denote the set S as:*

$$S := \mathcal{V}_x + \mathcal{P}_x = \left\{ x_i + \mathcal{P}_x^{(j)} \mid x_i \in \mathcal{V}_x, i, j \in \mathbb{N}^+ \right\}.$$

If S is dense in \mathbb{R}^{d_x} , $\{1, -1, \sqrt{2}, 0\}^{d_y} \subset \mathcal{V}_y$ and $\mathcal{P}_y = 0$, then $\mathcal{T}_{,\mathcal{P}}^\sigma$ can achieve the UAP. That is, for any continuous function $f : \mathbb{R}^{d_x-1} \rightarrow \mathbb{R}^{d_y}$ defined on a compact domain \mathcal{K} , and for any $\varepsilon > 0$, there always exist $X \in \mathbb{R}^{d_x \times n}$ and*

$Y \in \mathbb{R}^{d_y \times n}$ from the vocabulary \mathcal{V} (i.e., $x^{(i)} \in \mathcal{V}_x, y^{(i)} \in \mathcal{V}_y$) with some length $n \in \mathbb{N}^+$ such that

$$\|T_{*,\mathcal{P}}^\sigma(\tilde{x}; X, Y) - f(x)\| < \varepsilon, \quad \forall x \in \mathcal{K}. \quad (25)$$

We provide a constructive proof in Appendix D, and here we only demonstrate the proof idea by considering the specific case of $d_y = 1$ and assuming the matrices U , B , C , and D in the Transformer are identity matrices. In this case, the Transformer $T_{*,\mathcal{P}}^\sigma(\tilde{x}; X, Y)$ can be simplified to an FNN, N_*^σ , similar to the calculation in equation (19):

$$N_*^\sigma(x) = Y \sigma(X_{\mathcal{P}}^\top \tilde{x}) = \sum_{j=1}^n y^{(j)} \sigma\left(\left(x^{(j)} + \mathcal{P}_x^{(j)}\right) \cdot \tilde{x}\right). \quad (26)$$

The UAP of FNNs shown in Lemma 2 implies that the target function f can be approximated by an FNN $N^\sigma(x)$ with k hidden neurons:

$$N^\sigma(x) = A \sigma(W \tilde{x}) = \sum_{i=1}^k a_i \sigma(w_i \cdot \tilde{x}). \quad (27)$$

Since we are considering a continuous activation function σ , we can conclude that slightly perturbing the parameters A and W will lead to new FNNs that can still approximate f , provided the perturbations are small enough. This observation motivates us to construct a proof using the property that each $w_i \in \mathbb{R}^{d_x}$ can be approximated by vectors in $S = \mathcal{V}_x + \mathcal{P}_x$, and each $a_i \in \mathbb{R}$ can be approximated by numbers of the form $q_i \sqrt{2} \pm l_i$, with positive integers q_i and l_i . Note that the summation $\sum_{i=1}^k (q_i \sqrt{2} \pm l_i) \sigma(w_i \cdot \tilde{x})$ can be reformulated as $\sum_{i'=1}^{k'} y_{i'} \sigma(w_{i'} \cdot \tilde{x})$ with $k' = \sum_{i=1}^k (q_i + l_i)$, $y_{i'} \in \{\sqrt{2}, \pm 1\}$ and $w_{i'} \in \{w_1, \dots, w_k\}$. For each $w_{i'}$, we can choose a vector $\hat{w}_{i'} := x_{j_{i'}} + \mathcal{P}_x^{(j_{i'})} \in S$ that approximates $w_{i'}$ well, where $j_{i'} \in \mathbb{N}^+$ and $x_{j_{i'}} \in \mathcal{V}_x$. The integers $j_{i'}$ can be chosen to be distinct from each other.

Now, the FNN in (26) can be constructed by using $n = \max\{j_1, j_2, \dots, j_{k'}\}$ neurons, where the j -th neuron is assigned by setting $y^{(j)} = y_{i'} \in \mathcal{V}_y$ and $x^{(j)} = x_{j_{i'}} \in \mathcal{V}_x$ for the case of $j = j_{i'} \in \{j_1, j_2, \dots, j_{k'}\}$, and $y^{(j)} = 0 \in \mathcal{V}_y$ for the case of $j \notin \{j_1, j_2, \dots, j_{k'}\}$. Here, the nonzero value of $y^{(j)}$ highlights useful positions and demonstrations.

In the proof idea above, we take the density of the set S in \mathbb{R}^{d_x} as a fundamental assumption. \mathcal{V}_x contains only finitely many elements, rendering it bounded. For S to be dense in the entire space, \mathcal{P}_x must be unbounded. Next, we relax this requirement, eliminating the need for \mathcal{P}_x to be unbounded, making the conditions more aligned with practical scenarios. Particularly, we consider the specific activation function in the following Theorem 10, where the notations not explicitly mentioned remain consistent with those in Theorem 9.

Theorem 10. *If the set S is dense in $[-1, 1]^{d_x}$, then $\mathcal{T}_{*,\mathcal{P}}^{\text{ReLU}}$ is capable of achieving the UAP. Additionally, if S is only dense in a neighborhood $B(w^*, \delta)$ of a point $w^* \in \mathbb{R}^{d_x}$ with radius $\delta > 0$, then the class of transformers with exponential activation, $\mathcal{T}_{*,\mathcal{P}}^{\text{exp}}$, is capable of achieving the UAP.*

The density condition on S is significantly refined here. This improvement is possible because the proof of Theorem 9 relies directly on the UAP of FNNs, where the weights take values from the entire parameter space. However, for FNNs with specific activations, we can restrict the weights to a small set without losing the UAP.

For ReLU networks, we can use the positive homogeneity property, i.e., $A\text{ReLU}(W\tilde{x}) = \frac{1}{\lambda} A\text{ReLU}(\lambda W\tilde{x})$ for any $\lambda > 0$, to restrict the weight matrix W . In fact, the restriction that all elements of W take values in the interval $[-1, 1]$ does not affect the UAP of ReLU FNNs because the scale of W can be recovered by adjusting the scale of A via choosing a proper λ .

For exponential networks, the condition on S is much weaker than in the ReLU case. This relaxation is nontrivial, and the proof stems from a property of the derivatives of exponential functions. Consider the exponential function $\exp(w \cdot x)$ as a function of $w \in B(w^*, \delta)$, and denote it as $h(w)$,

$$\begin{aligned} h(w) &= \exp(w \cdot x) \\ &= \exp(w_1 x_1 + \dots + w_d x_d), \quad w, x \in \mathbb{R}^d, d = d_x, \end{aligned} \quad (28)$$

where w_i and $x_i \in \mathbb{R}$ are the components of w and x , respectively. Calculating the partial derivatives of $h(w)$, we observe the following relations:

$$\frac{\partial^\alpha h}{\partial w^\alpha} := \frac{\partial^{|\alpha|} h}{\partial w_1^{\alpha_1} \dots \partial w_d^{\alpha_d}} = x_1^{\alpha_1} \dots x_d^{\alpha_d} h(w), \quad (29)$$

where $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}^d$ is the index vector representing the order of partial derivatives, and $|\alpha| := \alpha_1 + \dots + \alpha_d$. This relationship allows us to link exponential FNNs to polynomials since any polynomial $P(x)$ can be represented in the following form:

$$P(x) = \exp(-w^* \cdot x) \left(\sum_{\alpha \in \Lambda} a_\alpha \frac{\partial^{|\alpha|} h}{\partial w^\alpha} \right) \Big|_{w=w^*}, \quad (30)$$

where a_α are the coefficients of the polynomials, Λ is a finite set of indices, and the partial derivatives can be approximated by finite differences, which are FNNs. For example, the first-order partial derivative $\frac{\partial h}{\partial w_1} \Big|_{w=w^*} = x_1 h(w^*)$ can be approximated by the following difference with a small

nonzero number $\lambda \in (0, \delta)$,

$$\begin{aligned} & \frac{h(w^* + \lambda e_1) - h(w^*)}{\lambda} \\ &= \frac{1}{\lambda} \exp((w^* + \lambda e_1) \cdot x) - \frac{1}{\lambda} \exp(w^* \cdot x). \end{aligned} \quad (31)$$

This is an exponential FNN with two neurons. Finally, employing the well-known Stone-Weierstrass theorem, which states that any continuous function f on compact domains can be approximated by polynomials, and combining the above relations between FNNs and polynomials, we can establish the UAP of exponential FNNs with weight constraints.

Remark 11. When discussing density, one of the most immediate examples that comes to mind is the density of rational numbers in \mathbb{R} . How can we effectively enumerate rational numbers? The work by (Calkin & Wilf, 2000) introduces an elegant method for enumerating positive rational numbers, synthesizing ideas from (Stern, 1858) and (Berndt et al., 1990). It demonstrates the computational feasibility of enumeration through an effective algorithm. Thus, we assume that positional encodings can be implemented using computer algorithms, such as iterative functions.

5. Conclusion

In this paper, we establish a connection between FNNs and Transformers through ICL. By leveraging the universal approximation property of FNNs, we demonstrate that the UAP of ICL holds when the context is selected from the entire vector space. When the context is drawn from a finite set, we explore the approximation power of vocabulary-based in-context learning, showing that the UAP is achievable only when appropriate positional encodings are incorporated, underscoring the importance of positional encodings.

In our work, we consider Transformers with input sequences of arbitrary length, implying that the positional encoding \mathcal{P}_x consists of a countably infinite set of elements, independent of the target function. As a result, the set S is also infinitely large and may or may not be dense in \mathbb{R}^d . In Theorem 9, we assume a strong density condition, which is later relaxed in Theorem 10. However, in practical applications, input sequences are finite, typically truncated for computational feasibility. This shift allows our conclusions to be interpreted through an approximation lens, where the objective is to approximate functions within a specified error margin, rather than achieving infinitesimal precision. Additionally, to achieve universal approximation, it is insightful to compare the function approximation capabilities of our approach (outlined in Lemma 4) with the direct use of FNNs, particularly when the Transformer parameters are trainable.

It is important to note that this paper is limited to single-

layer Transformers with absolute positional encodings, and the main results (Theorem 9 and Theorem 10) focus on elementwise activations. Future research should extend these findings to multi-layer Transformers, general positional encodings (such as RPEs and RoPE), and softmax activations. For softmax Transformers, our analysis in Sections 2 and 3 highlighted their connection to Transformers with exponential activations. However, extending this connection to the scenario in Section 4 proves challenging and requires more sophisticated techniques.

Although this paper primarily addresses theoretical issues, we believe our results can offer valuable insights for practitioners. Specifically, in Remark 11, we observe that certain algorithms use function composition to enumerate numbers dense in \mathbb{R} . This idea could inspire the design of positional encodings via compositions of fixed functions, similar to RNN approaches. RNNs capture the sequential nature of information by integrating the importance of word order in sentence meaning. However, to the best of our knowledge, existing research on RNNs has not explored the denseness properties of the sets formed by their hidden state sequences. We hope this unexplored property will inspire experimental research in future studies. Furthermore, our construction for Theorem 9 relies on the sparse partition assumption in equation (16). The practical validity of this assumption remains uncertain, and we leave this question open for future exploration.

In fact, some recent studies on continuous thought chains and continuous states have certain connections to our work—specifically, leveraging positional encoding to enable Transformers to achieve UAP for functions whose domain is a finite set while the range covers the entire Euclidean space. Moreover, a paper proposing an approach for automatically adjusting prompts for function fitting is also related to our theoretical findings. Therefore, with further research, our theory holds practical significance.

References

- Ahn, K., Cheng, X., Daneshmand, H., and Sra, S. Transformers learn to implement preconditioned gradient descent for in-context learning. In *Advances in Neural Information Processing Systems*, 2024.
- Alberti, S., Dern, N., Thesing, L., and Kutyniok, G. Sumformer: Universal approximation for efficient transformers. *Annual Workshop on Topology, Algebra, and Geometry in Machine Learning*, pp. 72–86, 2023.
- Apostol, T. M. *Modular Functions and Dirichlet Series in Number Theory (Graduate Texts in Mathematics, 41)*. Springer, 1989. ISBN 978-0387971278.

- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate, 2014.
- Berndt, B. C., Diamond, H. G., Halberstam, H., and Hildebrand, A. *Analytic Number Theory: Proceedings of a Conference in Honor of Paul T. Bateman*. Birkhäuser, 1990. ISBN 978-1461280347.
- Boole, G. *A Treatise on the Calculus of Finite Differences*. Cambridge University Press, 2009. ISBN 978-0511693014.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020.
- Cai, Y. Achieve the minimum width of neural networks for universal approximation. In *International Conference on Learning Representations*, 2023.
- Calkin, N. J. and Wilf, H. S. Recounting the rationals. *The American Mathematical Monthly*, 107:360–363, 2000.
- Cannarsa, P. and D’Aprile, T. *Introduction to Measure Theory and Functional Analysis*. Springer Cham, 2015. ISBN 978-3319170183.
- Cheng, X., Chen, Y., and Sra, S. Transformers implement functional gradient descent to learn non-linear functions in context. In *International Conference on Machine Learning*, 2024.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., Schuh, P., Shi, K., Tsvyashchenko, S., Maynez, J., Rao, A., Barnes, P., Tay, Y., Shazeer, N., Prabhakaran, V., Reif, E., Du, N., Hutchinson, B., Pope, R., Bradbury, J., Austin, J., Isard, M., Gur-Ari, G., Yin, P., Duke, T., Levskaya, A., Ghemawat, S., Dev, S., Michalewski, H., Garcia, X., Misra, V., Robinson, K., Fedus, L., Zhou, D., Ippolito, D., Luan, D., Lim, H., Zoph, B., Spiridonov, A., Sepassi, R., Dohan, D., Agrawal, S., Omernick, M., Dai, A. M., Pillai, T. S., Pellat, M., Lewkowycz, A., Moreira, E., Child, R., Polozov, O., Lee, K., Zhou, Z., Wang, X., Saeta, B., Diaz, M., Firat, O., Catasta, M., Wei, J., Meier-Hellstern, K., Eck, D., Dean, J., Petrov, S., and Fiedel, N. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24:1–113, 2023.
- Chu, Z., Chen, J., Chen, Q., Yu, W., He, T., Wang, H., Peng, W., Liu, M., Qin, B., and Liu, T. Navigate through enigmatic labyrinth a survey of chain of thought reasoning: Advances, frontiers and future. In *Annual Meeting of the Association for Computational Linguistics*, 2024.
- Cybenko, G. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2:303–314, 1989.
- Dai, D., Sun, Y., Dong, L., Hao, Y., Ma, S., Sui, Z., and Wei, F. Why can GPT learn in-context? language models secretly perform gradient descent as meta-optimizers. In *Annual Meeting of the Association for Computational Linguistics*, 2023.
- Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q., and Salakhutdinov, R. Transformer-xl: Attentive language models beyond a fixed-length context. In *Annual Meeting of the Association for Computational Linguistics*, 2019.
- De Cao, N., Aziz, W., and Titov, I. Editing factual knowledge in language models. In *Empirical Methods in Natural Language Processing*, 2021.
- Deora, P., Ghaderi, R., Taheri, H., and Thrampoulidis, C. On the optimization and generalization of multi-head attention. *Transactions on Machine Learning Research*, 2024.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Annual Meeting of the Association for Computational Linguistics*, 2019.
- Ding, B., Qin, C., Liu, L., Chia, Y. K., Li, B., Joty, S., and Bing, L. Is gpt-3 a good data annotator? In *Annual Meeting of the Association for Computational Linguistics*, 2023.
- Dong, Q., Li, L., Dai, D., Zheng, C., Ma, J., Li, R., Xia, H., Xu, J., Wu, Z., Chang, B., Sun, X., Li, L., and Sui, Z. A survey on in-context learning, 2024.
- Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Wang, M., and Wang, H. Retrieval-augmented generation for large language models: A survey, 2024.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- He, P., Liu, X., Gao, J., and Chen, W. DeBERTa: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2021.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 9:1735–1780, 1997.

- Hornik, K. Approximation capabilities of multilayer feed-forward networks. *Neural Networks*, 4:251–257, 1991.
- Hornik, K., Stinchcombe, M., and White, H. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2:359–366, 1989.
- Ke, G., He, D., and Liu, T.-Y. Rethinking positional encoding in language pre-training. In *International Conference on Learning Representations*, 2021.
- Khorashadizadeh, H., Mihindukulasooriya, N., Tiwari, S., Groppe, J., and Groppe, S. Exploring in-context learning capabilities of foundation models for generating knowledge graphs from text, 2023.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, 2022.
- Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86:2278–2324, 1998.
- Leshno, M., Lin, V. Y., Pinkus, A., and Schocken, S. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6:861–867, 1993.
- Li, L., Duan, Y., Ji, G., and Cai, Y. Minimum width of leaky-relu neural networks for uniform universal approximation. In *arXiv:2305.18460v3*, 2024.
- Likhoshervstov, V., Choromanski, K., and Weller, A. On the expressive power of self-attention matrices, 2021.
- Liu, X., Yu, H.-F., Dhillon, I., and Hsieh, C.-J. Learning to encode position for transformer with continuous dynamical model. In *International Conference on Machine Learning*, 2020.
- Lu, Z., Pu, H., Wang, F., Hu, Z., and Wang, L. The expressive power of neural networks: A view from the width. In *Advances in Neural Information Processing Systems*, 2017.
- Luo, S., Li, S., Zheng, S., Liu, T.-Y., Wang, L., and He, D. Your transformer may not be as powerful as you expect. In *Advances in Neural Information Processing Systems*, 2022.
- OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.-L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H. W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S. P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S. S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N. S., Khan, T., Kilpatrick, L., Kim, J. W., Kim, C., Kim, Y., Kirchner, J. H., Kiros, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C. M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S. M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H. P., Michael, Pokorný, Pokrass, M., Pong, V. H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Sel-sam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F. P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M. B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J. F. C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J. J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., and Zoph, B. Gpt-4 technical report, 2024.

- Park, S., Yun, C., Lee, J., and Shin, J. Minimum width for universal approximation. In *International Conference on Learning Representations*, 2021.
- Petrov, A., Torr, P., and Bibi, A. When do prompting and prefix-tuning work? a theory of capabilities and limitations. In *International Conference on Learning Representations*, 2024a.
- Petrov, A., Torr, P., and Bibi, A. Prompting a pretrained transformer can be a universal approximator. In *International Conference on Machine Learning*, 2024b.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21: 1–67, 2020.
- Ram, O., Levine, Y., Dalmedigos, I., Muhlga, D., Shashua, A., Leyton-Brown, K., and Shoham, Y. In-context retrieval-augmented language models. In *Annual Meeting of the Association for Computational Linguistics*, 2023.
- Rosenblatt, F. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65:386–408, 1958.
- Rudin, W. *Functional Analysis*. McGraw-Hill Science, 1991. ISBN 978-0070542365.
- Shaw, P., Uszkoreit, J., and Vaswani, A. Self-attention with relative position representations. In *Annual Meeting of the Association for Computational Linguistics*, 2018.
- Stern, M. Ueber eine zahlentheoretische funktion. *Journal für die reine und angewandte Mathematik*, 1858:193–220, 1858.
- Su, J., Ahmed, M., Lu, Y., Pan, S., Bo, W., and Liu, Y. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- Sutskever, I., Vinyals, O., and Le, Q. V. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, 2014.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. Llama: Open and efficient foundation language models, 2023.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ukasz Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., and Lang, K. Phoneme recognition using time-delay neural networks. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 37:328–339, 1989.
- Wang, B., Shang, L., Lioma, C., Jiang, X., Yang, H., Liu, Q., and Simonsen, J. G. On position embeddings in bert. In *International Conference on Learning Representations*, 2021a.
- Wang, S., Liu, Y., Xu, Y., Zhu, C., and Zeng, M. Want to reduce labeling cost? gpt-3 can help. In *Empirical Methods in Natural Language Processing*, 2021b.
- Wang, X., Zhu, W., Saxon, M., Steyvers, M., and Wang, W. Y. Large language models are latent variable models: Explaining and finding good demonstrations for in-context learning. In *Advances in Neural Information Processing Systems*, 2024.
- Wei, J., Bosma, M., Zhao, V., Guu, K., Yu, A. W., Lester, B., Du, N., Dai, A. M., and Le, Q. V. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022a.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q. V., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, 2022b.
- Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*, 2022.
- Xun, G., Jia, X., Gopalakrishnan, V., and Zhang, A. A survey on context learning. *IEEE Transactions on Knowledge and Data Engineering*, 29:38–56, 2017.
- Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. R., and Le, Q. V. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in Neural Information Processing Systems*, 2019.
- Yun, C., Bhojanapalli, S., Rawat, A. S., Reddi, S., and Kumar, S. Are transformers universal approximators of sequence-to-sequence functions? In *International Conference on Learning Representations*, 2020.
- Zhenzhong, L., Mingda, C., Sebastian, G., Kevin, G., Piyush, S., and Radu, S. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2021.

A. Table of notations

We present all our notations in Table 1 for easy reference.

Table 1. Table of Notations

Notations	Explanations
d_x, d_y	Dimensions of input and output.
\mathcal{P}	Positional encoding.
X, Y	Context without positional encoding.
$X_{\mathcal{P}}, Y_{\mathcal{P}}$	Context with positional encoding \mathcal{P} .
Z	Input without positional encoding.
$Z_{\mathcal{P}}$	Input with positional encoding \mathcal{P} .
\mathcal{V}	Vocabulary.
$\mathcal{V}_x, \mathcal{V}_y$	Vocabulary of $x^{(i)}$ and $y^{(i)}$.
σ	Activation function.
$\#$	The cardinality of a set.
$\mathcal{N}^{\sigma}, \mathcal{N}^{\sigma}$	One-hidden-layer FNN and its collection.
$\mathcal{T}^{\sigma}, \mathcal{T}^{\sigma}$	Single-layer Transformer and its collection.
$\mathcal{N}_{*}^{\sigma}, \mathcal{N}_{*}^{\sigma}$	One-hidden-layer FNN with a finite set of weights and its collection.
$\mathcal{T}_{*}^{\sigma}, \mathcal{T}_{*}^{\sigma}$	Single-layer Transformer with vocabulary restrictions and its collection.
$\mathcal{T}_{*,\mathcal{P}}^{\sigma}, \mathcal{T}_{*,\mathcal{P}}^{\sigma}$	Single-layer Transformer with positional encoding, vocabulary restrictions, and its collection.
$\ \cdot\ $	The uniform norm of vectors, <i>i.e.</i> , a shorthand for $\ \cdot\ _{\infty}$.
\tilde{x}	Append a one to the end of x , <i>i.e.</i> , $\tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix}$.

B. Proof for Section 2

In this appendix, we provide detailed proofs of Lemma 3 and Lemma 4.

B.1. Proof of The UAP of softmax FNNs

Lemma 12. *For any continuous function $f : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ defined on a compact domain \mathcal{K} and $\varepsilon > 0$, there always exist a softmax FNN $\mathcal{N}^{\text{softmax}}(x) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ satisfying*

$$\|\mathcal{N}^{\text{softmax}}(x) - f(x)\| < \varepsilon. \quad (32)$$

Proof. According to Lemma 2 we can construct a network

$$\begin{aligned} \mathcal{N}^{\text{exp}}(x) &= A \exp(Wx + b) \\ &= A \begin{bmatrix} \exp((Wx + b)_1) \\ \exp((Wx + b)_2) \\ \vdots \\ \exp((Wx + b)_k) \end{bmatrix}, \end{aligned} \quad (33)$$

such that $\|\mathcal{N}^{\text{exp}}(x) - f(x)\| < \varepsilon/2$ for all $x \in \mathcal{K}$ and k represents the width of hidden layer. We now construct a softmax network as follows

$$\mathcal{N}^{\text{softmax}}(x) = A' \text{softmax} \left(\begin{bmatrix} Wx + b' \\ 0 \end{bmatrix} \right), \quad (34)$$

where every element in $b' = b'(\varepsilon)$ is sufficiently small to satisfy $\exp((Wx + b')_i) < \frac{\varepsilon'}{k}$ for all $x \in \mathcal{K}, i = 1, 2, \dots, k$, and

$$A'_{i,j} = \begin{cases} A_{i,j} \exp(b_j - b'_j) & j = 1, \dots, k \\ 0 & j = k + 1 \end{cases}, \text{ where } i = 1, \dots, d_y. \text{ We can compute that}$$

$$\|f(x) - N^{\text{softmax}}(x)\| \leq \|f(x) - N^{\text{exp}}(x)\| + \|N^{\text{exp}}(x) - N^{\text{softmax}}(x)\|. \quad (35)$$

We focus on estimating of the upper bound of the second term, since it is evident that the first term does not exceed ε .

$$\begin{aligned} \|N^{\text{exp}}(x) - N^{\text{softmax}}(x)\| &= \max_{1 \leq i \leq d_y} \left\| \sum_{j=1}^k A_{i,j} \exp((Wx + b)_j) - \frac{\sum_{j=1}^k A'_{i,j} \exp((Wx + b')_j)}{1 + \sum_{j=1}^k \exp((Wx + b')_j)} \right\| \\ &= \max_{1 \leq i \leq d_y} \left\| \sum_{j=1}^k A_{i,j} \exp((Wx + b)_j) - \frac{\sum_{j=1}^k A_{i,j} \exp((Wx + b)_j)}{1 + \sum_{j=1}^k \exp((Wx + b')_j)} \right\| \\ &= \|N^{\text{exp}}(x)\| \left(1 - \frac{1}{1 + \sum_{j=1}^k \exp((Wx + b')_j)} \right) \\ &\leq \|N^{\text{exp}}(x)\| \left(1 - \frac{1}{1 + \varepsilon'} \right) \\ &= \|N^{\text{exp}}(x)\| \varepsilon'. \end{aligned} \quad (36)$$

By setting $\varepsilon' = \frac{\varepsilon}{1 + \max_{x \in \mathcal{X}} \|N^{\text{exp}}(x)\|}$, we ensure it is finite, leading to the conclusion that

$$\|f(x) - N^{\text{softmax}}(x)\| \leq \varepsilon. \quad (37)$$

This finishes the proof. \square

B.2. Proof of Lemma 3

Lemma 3. *Let σ be an elementwise activation and T^σ be a single-layer Transformer. For any one-hidden-layer network $N^\sigma : \mathbb{R}^{d_x-1} \rightarrow \mathbb{R}^{d_y} \in \mathcal{N}^\sigma$ with n hidden neurons, there exist matrices $X \in \mathbb{R}^{d_x \times n}$ and $Y \in \mathbb{R}^{d_y \times n}$ such that*

$$T^\sigma(\tilde{x}; X, Y) = N^\sigma(x), \quad \forall x \in \mathbb{R}^{d_x-1}. \quad (38)$$

Proof. We can directly compute the following

$$\begin{aligned} T^\sigma(\tilde{x}; X, Y) &= (Z + \text{Attn}_{Q,K,V}^\sigma(\tilde{x}; X, Y))_{d_x+1:d_x+d_y, n+1} \\ &= (Z + VZM\sigma(Z^\top Q^\top KZ))_{d_x+1:d_x+d_y, n+1} \\ &= \left(Z + \begin{bmatrix} DX + EY & 0 \\ FX + UY & 0 \end{bmatrix} \begin{bmatrix} \sigma(X^\top B^\top CX) & \sigma(X^\top B^\top C\tilde{x}) \\ \sigma(\tilde{x}^\top B^\top CX) & \sigma(\tilde{x}^\top B^\top C\tilde{x}) \end{bmatrix} \right)_{d_x+1:d_x+d_y, n+1}. \end{aligned} \quad (39)$$

It is obvious that

$$T^\sigma(\tilde{x}; X, Y) = (FX + UY)\sigma(X^\top B^\top C\tilde{x}). \quad (40)$$

Assume $N^\sigma(x) = A\sigma(Wx + b)$ is an arbitrary single-layer FNN, where $W \in \mathbb{R}^{k \times d_x}$, $A \in \mathbb{R}^{d_y \times k}$, $b \in \mathbb{R}^k$, and k represents the width of hidden layer.

Let us set the length of context to k , that is $X \in \mathbb{R}^{d_x \times k}$, $Y \in \mathbb{R}^{d_y \times k}$. Through trivial calculation we can find that if we set

$$X = (C^\top B)^{-1} \begin{bmatrix} W^\top \\ b^\top \end{bmatrix}, \quad Y = U^{-1}(A - FX), \quad (41)$$

then $\mathcal{T}^\sigma(\tilde{x}; X, Y) = \mathcal{N}^\sigma(x)$ holds. \square

It is noteworthy that in the above proof, we did not set the matrix F to zero, which differs from the previous assumption. Here, we want to highlight to the reader that F does not necessarily need to be set to zero; in our assumptions, setting F to zero was merely for computational convenience.

B.3. Proof of Lemma 4

Lemma 4. *Let \mathcal{T}^σ be a single-layer Transformer with elementwise or softmax activation, and \mathcal{K} be a compact domain in \mathbb{R}^{d_x-1} . Then for any continuous function $f : \mathcal{K} \rightarrow \mathbb{R}^{d_y}$ and any $\varepsilon > 0$, there exist matrices $X \in \mathbb{R}^{d_x \times n}$ and $Y \in \mathbb{R}^{d_y \times n}$ such that*

$$\|\mathcal{T}^\sigma(\tilde{x}; X, Y) - f(x)\| < \varepsilon, \quad \forall x \in \mathcal{K}. \quad (42)$$

Proof. For elementwise activation cases, with the help of Lemma 2 and Lemma 3, the conclusion follows trivially.

Then we solve the softmax case. Similarly, for any $\varepsilon > 0$, we can construct a softmax FNN $\mathcal{N}^{\text{softmax}}(x) = A^{\text{softmax}}\left(\begin{bmatrix} Wx + b \\ 0 \end{bmatrix}\right)$, with k hidden neurons, using Lemma 12 such that $\|\mathcal{N}^{\text{softmax}}(x) - f(x)\| < \varepsilon/2$ for all $x \in \mathcal{K}$. What we need to do is to approximate this softmax FNN with a softmax Transformer.

We can directly compute the following

$$\begin{aligned} \mathcal{T}^{\text{softmax}}(\tilde{x}; X, Y) &= \left(Z + \text{Attn}_{Q,K,V}^{\text{softmax}}(\tilde{x}; X, Y) \right)_{d_x+1:d_x+d_y, n+1} \\ &= (Z + V Z M^{\text{softmax}}(Z^\top Q^\top K Z))_{d_x+1:d_x+d_y, n+1} \\ &= \left(Z + \begin{bmatrix} DX + EY & 0 \\ FX + UY & 0 \end{bmatrix} \text{softmax} \left(\begin{bmatrix} X^\top B^\top C X & X^\top B^\top C \tilde{x} \\ \tilde{x}^\top B^\top C X & \tilde{x}^\top B^\top C \tilde{x} \end{bmatrix} \right) \right)_{d_x+1:d_x+d_y, n+1}. \end{aligned} \quad (43)$$

Since $F = 0$, we have:

$$\mathcal{T}^{\text{softmax}}(\tilde{x}; X, Y) = UY \text{softmax} \left(\begin{bmatrix} X^\top B^\top C \tilde{x} \\ \tilde{x}^\top B^\top C \tilde{x} \end{bmatrix} \right)_{1:n}. \quad (44)$$

Then through comparing the output of the softmax Transformer with the exponential FNN, we can find out that there is one more bounded positive term $t(x) = \exp(\tilde{x}^\top B^\top C \tilde{x})$ when processing normalization.

Chose $n = k + 1$, $\varepsilon' = \frac{\varepsilon}{1 + \max_{x \in \mathcal{K}} \|\mathcal{N}^{\text{softmax}}(x)\|}$, and X, Y such that $X^\top B^\top C = \begin{bmatrix} W & b + s\mathbf{1} \\ 0 & s \end{bmatrix} \in \mathbb{R}^{n \times d_x}$ and $UY = \begin{bmatrix} A & 0 \end{bmatrix} \in \mathbb{R}^{d_y \times n}$, where s is big enough, making $\exp(\tilde{x}^\top B^\top C \tilde{x} - s) < \varepsilon'$ for all $x \in \mathcal{K}$. Then $X^\top B^\top C \tilde{x} = \begin{bmatrix} W & b + s\mathbf{1} \\ 0 & s \end{bmatrix} \begin{bmatrix} x \\ 1 \end{bmatrix} = \begin{bmatrix} Wx + b + s\mathbf{1} \\ s \end{bmatrix}$, and we can compute a detailed form of $\mathcal{T}^{\text{softmax}}$ as:

$$\begin{aligned} \mathcal{T}^{\text{softmax}}(\tilde{x}; X, Y) &= \begin{bmatrix} \frac{\sum_{j=1}^k A_{1,j} \exp((Wx+b)_j + s)}{\sum_{j=1}^k \exp((Wx+b)_j + s) + \exp(s) + \exp(\tilde{x}^\top B^\top C \tilde{x})} \\ \frac{\sum_{j=1}^k A_{2,j} \exp((Wx+b)_j + s)}{\sum_{j=1}^k \exp((Wx+b)_j + s) + \exp(s) + \exp(\tilde{x}^\top B^\top C \tilde{x})} \\ \vdots \\ \frac{\sum_{j=1}^k A_{d_y,j} \exp((Wx+b)_j + s)}{\sum_{j=1}^k \exp((Wx+b)_j + s) + \exp(s) + \exp(\tilde{x}^\top B^\top C \tilde{x})} \end{bmatrix} = \begin{bmatrix} \frac{\sum_{j=1}^k A_{1,j} \exp((Wx+b)_j)}{\sum_{j=1}^k \exp((Wx+b)_j) + 1 + \exp(\tilde{x}^\top B^\top C \tilde{x} - s)} \\ \frac{\sum_{j=1}^k A_{2,j} \exp((Wx+b)_j)}{\sum_{j=1}^k \exp((Wx+b)_j) + 1 + \exp(\tilde{x}^\top B^\top C \tilde{x} - s)} \\ \vdots \\ \frac{\sum_{j=1}^k A_{d_y,j} \exp((Wx+b)_j)}{\sum_{j=1}^k \exp((Wx+b)_j) + 1 + \exp(\tilde{x}^\top B^\top C \tilde{x} - s)} \end{bmatrix}. \end{aligned} \quad (45)$$

We focus on estimating the upper bound of the distance between N^{softmax} and T^{softmax} , that is

$$\begin{aligned}
 & \|N^{\text{softmax}}(x) - T^{\text{softmax}}(\tilde{x}; X, Y)\| \\
 &= \max_{1 \leq i \leq d_y} \left\| \frac{\sum_{j=1}^k A_{i,j} \exp((Wx + b)_j)}{\sum_{j=1}^k \exp((Wx + b)_j) + 1} - \frac{\sum_{j=1}^k A_{i,j} \exp((Wx + b)_j)}{\sum_{j=1}^k \exp((Wx + b)_j) + 1 + \exp(\tilde{x}^\top B^\top C \tilde{x} - s)} \right\| \\
 &= \|N^{\text{softmax}}(x)\| \left| 1 - \frac{\sum_{j=1}^k \exp((Wx + b)_j) + 1}{\sum_{j=1}^k \exp((Wx + b)_j) + 1 + \exp(\tilde{x}^\top B^\top C \tilde{x} - s)} \right| \\
 &= \|N^{\text{softmax}}(x)\| \left| \frac{\exp(\tilde{x}^\top B^\top C \tilde{x} - s)}{\sum_{j=1}^k \exp((Wx + b)_j) + 1 + \exp(\tilde{x}^\top B^\top C \tilde{x} - s)} \right| \\
 &\leq \|N^{\text{softmax}}(x)\| |\exp(\tilde{x}^\top B^\top C \tilde{x} - s)| \tag{46} \\
 &\leq \|N^{\text{softmax}}(x)\| \varepsilon' \leq \varepsilon. \tag{47}
 \end{aligned}$$

As a consequence, we have $\|T^{\text{softmax}}(\tilde{x}; X, Y) - f(x)\| \leq \varepsilon$ for all $x \in \mathcal{K}$, which finishes the proof. \square

C. Proof for Section 3

In this Appendix, we provide detailed proofs of the Proposition 5, Lemma 6, and Theorem 7 presented in Section 3.

C.1. Proof of Proposition 5

Proposition 5. *The scalar function $h_k(x) = \sum_{i=1}^k a_i e^{b_i x}$, where $a_i, b_i, x \in \mathbb{R}$ and at least one a_i is nonzero, has at most $k - 1$ zero points.*

Proof. We prove this statement by induction. When $k = 1$ and 2, the statement is easy to prove. For the case $k = N$, suppose that every h_N has at most $N - 1$ zero points.

Now consider $k = N + 1$. Let $h_{N+1}(x) = \sum_{i=1}^{N+1} a_i \exp(b_i x)$. Without loss of generality, assume $a_{N+1} \neq 0$. Thus, we can rewrite $h_{N+1}(x)$ as

$$h_{N+1}(x) = a_{N+1} e^{b_{N+1} x} \left(1 + \sum_{i=1}^N \frac{a_i}{a_{N+1}} e^{(b_i - b_{N+1})x} \right). \tag{48}$$

We proceed by contradiction. Suppose $h_{N+1}(x)$ has more than N zero points. This implies

$$g(x) := 1 + \sum_{i=1}^N \frac{a_i}{a_{N+1}} e^{(b_i - b_{N+1})x}, \tag{49}$$

has more than N zero points.

Then, according to Rolle's Theorem, $g'(x)$ must have more than $N - 1$ zero points. Since $g'(x) = \sum_{i=1}^N \frac{a_i(b_i - b_{N+1})}{a_{N+1}} e^{(b_i - b_{N+1})x}$ must have at least N zero points, this leads to a contradiction.

Thus, $h_{N+1}(x) = \sum_{i=1}^{N+1} a_i e^{b_i x}$ can have at most N zero points. The proof is complete. \square

C.2. Proof of Lemma 6

Lemma 6. *The function class \mathcal{N}_*^σ , with elementwise or softmax activation σ , cannot achieve the UAP. Specifically, for any compact domain $\mathcal{K} \subset \mathbb{R}^{d_x}$, there exists a continuous function $f : \mathcal{K} \rightarrow \mathbb{R}^{d_y}$ and $\varepsilon_0 > 0$ such that*

$$\max_{x \in \mathcal{K}} \|f(x) - \mathcal{N}_*^\sigma(\tilde{x})\| \geq \varepsilon_0, \quad \forall \mathcal{N}_*^\sigma \in \mathcal{N}_*^\sigma. \quad (50)$$

Proof. For any elementwise activation σ , the span of \mathcal{N}_*^σ , denoted as $\text{span}\{\mathcal{N}_*^\sigma\}$, forms a finite-dimensional function space because \mathcal{N}_*^σ is spanned by finitely many basis functions $\{\sigma(w_i \cdot x + b_i) \mid (w_i, b_i) \in \mathcal{W} \times \mathcal{B}\}$. $\text{Span}\{\mathcal{N}_*^\sigma\}$ is closed under the uniform norm supported by Theorem 1.21 from (Rudin, 1991) and Corollary C.4 from (Cannarsa & D’Aprile, 2015). This implies that the set of functions approximable by $\text{span}\{\mathcal{N}_*^\sigma\}$ is precisely the set of functions within $\text{span}\{\mathcal{N}_*^\sigma\}$. Consequently, any function not in $\text{span}\{\mathcal{N}_*^\sigma\}$ cannot be arbitrarily approximated, meaning that the UAP cannot be achieved.

Without loss of generality, for any $\mathcal{N}_*^{\text{softmax}} \in \mathcal{N}_*^{\text{softmax}}$, assume $\mathcal{K} = [0, 1]^{d_x}$ and consider only the first component of x . Thus, we may assume $d_x = 1$. Let us consider the output of an arbitrary j -th dimension, that is

$$(\mathcal{N}_*^{\text{softmax}})^{(j)} = \frac{\sum_{i=1}^k A_{j,i} \exp(w_i x + b_i)}{\sum_{l=1}^k \exp(w_l x + b_l)}, \quad (51)$$

where $(A, w_i, b_i) \in \hat{\mathcal{A}} \times \mathcal{W} \times \mathcal{B}$, $k \in \mathbb{N}^+$. Here, the parameter spaces are defined as $\hat{\mathcal{A}} \subset \mathbb{R}^{d_y \times k}$, $\mathcal{W} \subset \mathbb{R}^{d_x}$, and $\mathcal{B} \subset \mathbb{R}$ are finite sets. Consequently, the set $\{(w_i, b_i) \mid (w_i, b_i) \in \mathcal{W} \times \mathcal{B}\}$ is finite, and we denote the number of elements in this set as \mathfrak{R} . By regrouping identical terms in the numerator, we can rewrite the equation as,

$$(\mathcal{N}_*^{\text{softmax}})^{(j)} = \frac{\sum_{h(i)=1}^{\tilde{k}} \tilde{A}_{j,h(i)} \exp(w_{h(i)} x + b_{h(i)})}{\sum_{l=1}^k \exp(w_l x + b_l)}, \quad (52)$$

where $\tilde{k} \leq \mathfrak{R}$ and $\tilde{k} \leq k$. It is important to note that this transformation applies to any $\mathcal{N}_*^{\text{softmax}} \in \mathcal{N}_*^{\text{softmax}}$, ensuring that the number of summation terms in the numerator remains strictly bounded by \mathfrak{R} . Since the case of $\tilde{A}_{j,h(i)}$ for all i is trivial, we further assume that there is at least one $\tilde{A}_{j,h(i)}, i \in \{1, 2, \dots, \tilde{k}\}$, being nonzero. Then the numerator, $\sum_{h(i)=1}^{\tilde{k}} \tilde{A}_{j,h(i)} \exp(w_{h(i)} x + b_{h(i)})$, can have at most $\tilde{k} - 1$ zero points.

Now, we consider a special function $f(x) = \sin(mx)$, where $\lceil \frac{m}{\pi} \rceil > \mathfrak{R} - 1$, and the period is $T = \frac{2\pi}{m}$. $\lceil x \rceil$ is the smallest integer greater than or equal to x .

Let us focus on two adjacent extreme points x_1, x_2 , where $f(x_1) = 1$ and $f(x_2) = -1$. We proceed by contradiction in our proof. Suppose $\mathcal{N}_*^{\text{softmax}}$ can achieve the UAP. There exists $\mathcal{N}_*^{\text{softmax}} \in \mathcal{N}_*^{\text{softmax}}$ such that $|(\mathcal{N}_*^{\text{softmax}})^{(j)} - f(x)| < \varepsilon$ for all $x \in [0, 1]$.

Taking $\varepsilon = 0.1$, we have:

$$\begin{aligned} |(\mathcal{N}_*^{\text{softmax}}(x_1))^{(j)} - f(x_1)| < 0.1 &\Rightarrow (\mathcal{N}_*^{\text{softmax}}(x_1))^{(j)} > -0.1 + f(x_1) = 0.9, \\ |(\mathcal{N}_*^{\text{softmax}}(x_2))^{(j)} - f(x_2)| < 0.1 &\Rightarrow (\mathcal{N}_*^{\text{softmax}}(x_2))^{(j)} < 0.1 + f(x_2) = -0.9, \end{aligned}$$

By the intermediate value theorem, there exists some $x_0 \in (\min(x_1, x_2), \max(x_1, x_2))$, such that $(\mathcal{N}_*^{\text{softmax}}(x_0))^{(j)} = 0$. Therefore, there is at least one zero of $(\mathcal{N}_*^{\text{softmax}}(x))^{(j)}$ between two adjacent extrema of $f(x)$, and the total number of zero points of $(\mathcal{N}_*^{\text{softmax}}(x))^{(j)}$ in the interval $[0, 1]$ is at least $\lceil \frac{m}{\pi} \rceil$. Thus, the number of zeros of $(\mathcal{N}_*^{\text{softmax}}(x))^{(j)}$ exceeds $\mathfrak{R} - 1$, which contradicts the fact that any $(\mathcal{N}_*^{\text{softmax}}(x))^{(j)}$ has at most $\tilde{k} - 1$ zero points.

If approximation cannot be achieved in one dimension, it is evident that it cannot be achieved in higher dimensions either. Therefore, $\mathcal{N}_*^{\text{softmax}}$ cannot achieve the UAP. \square

C.3. Proof of Theorem 7

Theorem 7. *The function class \mathcal{T}_*^σ , with elementwise or softmax activation σ , cannot achieve the UAP. Specifically, for any compact domain $\mathcal{K} \subset \mathbb{R}^{d_x-1}$, there exists a continuous function $f : \mathcal{K} \rightarrow \mathbb{R}^{d_y}$ and $\varepsilon_0 > 0$ such that*

$$\max_{x \in \mathcal{K}} \|f(x) - T_*^\sigma(\tilde{x})\| \geq \varepsilon_0, \quad \forall T_*^\sigma \in \mathcal{T}_*^\sigma. \quad (53)$$

Proof. For cases of elementwise activations, since T_*^σ has a structure similar to N_*^σ , we find that $\text{span}\{T_*^\sigma\}$ is also a finite-dimensional function space. Hence, the same argument from Lemma 6 can be applied here to complete the proof.

Without loss of generality, for any $T_*^{\text{softmax}} \in \mathcal{T}_*^{\text{softmax}}$, assume $\mathcal{K} = [0, 1]^{d_x}$ and consider the output of an arbitrary j -th dimension and one-dimensional input as an example that is

$$(T_*^{\text{softmax}})^{(j)} = \frac{\sum_{i=1}^k A_{j,i} \exp(w_i x + b_i)}{\sum_{i=1}^k \exp(w_i x + b_i) + \exp(\tilde{x}^\top B^\top C^\top \tilde{x})}. \quad (54)$$

The mathematical validity of this formulation is rigorously established through Equation (45) in Lemma 4. We observe that the form of the numerator remains consistent with Lemma 6, and we follow the same proof as above. Building upon the same proof methodology employed in Lemma 6, we have formally established through Proposition 5 that the number of zeros of T_*^{softmax} is exclusively determined by a finite parameter set and remains a bounded quantity. This theoretical framework enables the explicit construction of functions whose zero counts exceed this established bound. \square

D. Proof for Section 4

In this Appendix, we introduce Lemma 13 to assist in the proof of Theorem 9 and utilize Lemma 14 to provide a detailed proof of Theorem 10.

D.1. Proof of Lemma 13

Lemma 13. *For a network with a fixed width and a continuous activation function, it is possible to apply slight perturbations within an arbitrarily small error margin. For any network $N_1^q(x)$ defined on a compact set $\mathcal{K} \subset \mathbb{R}^{d_x}$, with parameters $A \in \mathbb{R}^{d_y \times k}$, $W \in \mathbb{R}^{k \times d_x}$, $b \in \mathbb{R}^{k \times 1}$, there exists $M > 0$, $M_1 > 0$ ($\|x\| < M$ and $\|a_i\| < M_1, i = 1, \dots, k$), and for any $\varepsilon > 0$, there exists $0 < \delta < \frac{\varepsilon}{2M_1k}$ and a perturbed network $N_2^q(x)$ with parameters $\tilde{A} \in \mathbb{R}^{d_y \times k}$, $\tilde{W} \in \mathbb{R}^{k \times d_x}$, $\tilde{b} \in \mathbb{R}^{k \times 1}$ ($\|\sigma(\tilde{w}_i x + \tilde{b}_i)\| < M_1, i = 1, \dots, k$), such that if $\max\{\|a_i - \tilde{a}_i\|, M\|w_i - \tilde{w}_i\| + \|b - \tilde{b}\| \mid i = 1, \dots, k\} < \delta$, then*

$$\|N_1(x) - N_2(x)\| < \varepsilon, \quad \forall x \in \mathcal{K}, \quad (55)$$

where a_i, \tilde{a}_i are the i -th column vectors of A, \tilde{A} , respectively, w_i, \tilde{w}_i are the i -th row vectors of W, \tilde{W} , and b_i, \tilde{b}_i are the i -th components of b, \tilde{b} , respectively, for any $i = 1, \dots, k$.

Proof. We have $N_1^q(x) = \sum_{i=1}^k a_i \sigma(w_i x + b_i)$, where $a_i \in \mathbb{R}^{d_y}$, $w_i \in \mathbb{R}^{d_x}$, $b_i \in \mathbb{R}$, and $\tilde{N}_2^q(x) = \sum_{i=1}^k \tilde{a}_i \sigma(\tilde{w}_i x + \tilde{b}_i)$, where $\tilde{a}_i \in \mathbb{R}^{d_y}$, $\tilde{w}_i \in \mathbb{R}^{d_x}$, $\tilde{b}_i \in \mathbb{R}$. For any $x \in \mathcal{K}$, $\|x\| < M$. There exists a constant $M_1 > 0$ such that for any $i = 1, \dots, k$, the following inequalities hold: $\|a_i\| < M_1$ and $\|\sigma(\tilde{w}_i x + \tilde{b}_i)\| < M_1$.

Due to the continuity of the activation function, for any $\varepsilon > 0$, there exists $0 < \delta < \frac{\varepsilon}{2M_1k}$, such that if $\|w_i x + b_i - (\tilde{w}_i x + \tilde{b}_i)\| \leq \|w_i - \tilde{w}_i\| \|x\| + \|b_i - \tilde{b}_i\| < M\|w_i - \tilde{w}_i\| + \|b - \tilde{b}\| < \delta, i = 1, \dots, k$, then $\|\sigma(w_i x + b_i) - \sigma(\tilde{w}_i x + \tilde{b}_i)\| < \frac{\varepsilon}{2M_1k}, i = 1, \dots, k$, and $\|a_i - \tilde{a}_i\| < \delta, i = 1, \dots, k$. This conclusion is ensured by the continuity of the activation function.

Combining all these inequalities, we can further derive:

$$\begin{aligned}
 \|N_1^\sigma(x) - N_2^\sigma(x)\| &= \left\| \sum_{i=1}^k a_i \sigma(w_i x + b_i) - \sum_{i=1}^k \tilde{a}_i \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| \\
 &\leq \left\| \sum_{i=1}^k a_i \sigma(w_i x + b_i) - \sum_{i=1}^k a_i \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| + \left\| \sum_{i=1}^k a_i \sigma(\tilde{w}_i x + \tilde{b}_i) - \sum_{i=1}^k \tilde{a}_i \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| \\
 &\leq \max_i \|a_i\| \left\| \sum_{i=1}^k \sigma(w_i x + b_i) - \sum_{i=1}^k \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| + \max_i \left\| \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| \left\| \sum_{i=1}^k a_i - \sum_{i=1}^k \tilde{a}_i \right\| \quad (56) \\
 &\leq \max_i \|a_i\| \sum_{i=1}^k \left\| \sigma(w_i x + b_i) - \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| + \max_i \left\| \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| \sum_{i=1}^k \|a_i - \tilde{a}_i\| \\
 &< M_1 k \frac{\varepsilon}{2M_1 k} + M_1 k \frac{\varepsilon}{2M_1 k} = \varepsilon
 \end{aligned}$$

The proof is complete. \square

D.2. Proof of Theorem 9

Theorem 9. Let $\mathcal{T}_{*,\mathcal{P}}^\sigma$ be the class of functions $T_{*,\mathcal{P}}^\sigma$, where σ is an elementwise activation, the subscript refers the finite vocabulary $\mathcal{V} = \mathcal{V}_x \times \mathcal{V}_y$, $\mathcal{P} = \mathcal{P}_x \times \mathcal{P}_y$ represents the positional encoding map, and denote the set S as:

$$S := \mathcal{V}_x + \mathcal{P}_x = \left\{ x_i + \mathcal{P}_x^{(j)} \mid x_i \in \mathcal{V}_x, i, j \in \mathbb{N}^+ \right\}. \quad (57)$$

If S is dense in \mathbb{R}^{d_x} , $\{1, -1, \sqrt{2}, 0\}^{d_y} \subset \mathcal{V}_y$ and $\mathcal{P}_y = 0$, then $\mathcal{T}_{*,\mathcal{P}}^\sigma$ can achieve the UAP. That is, for any continuous function $f : \mathbb{R}^{d_x-1} \rightarrow \mathbb{R}^{d_y}$ defined on a compact domain \mathcal{K} , and for any $\varepsilon > 0$, there always exist $X \in \mathbb{R}^{d_x \times n}$ and $Y \in \mathbb{R}^{d_y \times n}$ from the vocabulary \mathcal{V} (i.e., $x^{(i)} \in \mathcal{V}_x, y^{(i)} \in \mathcal{V}_y$) with some length $n \in \mathbb{N}^+$ such that

$$\|T_{*,\mathcal{P}}^\sigma(\tilde{x}; X, Y) - f(x)\| < \varepsilon, \quad \forall x \in \mathcal{K}. \quad (58)$$

Proof. Our conclusion holds for all element-wise continuous activation functions in $\mathcal{T}_{*,\mathcal{P}}^\sigma$. We assume that with $d_y = 1$. Similar cases can be inferred by analogy.

We are reformulating the problem.

Using Lemma 3, we have,

$$T_{*,\mathcal{P}}^\sigma(\tilde{x}; X, Y) = UY_{\mathcal{P}} \sigma \left((X + \mathcal{P})^\top B^\top C \tilde{x} \right) = UY_{\mathcal{P}} \sigma \left(X_{\mathcal{P}}^\top B^\top C \tilde{x} \right). \quad (59)$$

Since $\mathcal{P}_y = 0$, it follows that $Y_{\mathcal{P}} = Y$. For any continuous function $f : \mathbb{R}^{d_x-1} \rightarrow \mathbb{R}^{d_y}$ defined on a compact domain \mathcal{K} and for any $\varepsilon > 0$, we aim to show that there exists $T_{*,\mathcal{P}}^\sigma \in \mathcal{T}_{*,\mathcal{P}}^\sigma$ such that:

$$\begin{aligned}
 &\left\| T_{*,\mathcal{P}}^\sigma \left(\begin{bmatrix} x \\ 1 \end{bmatrix}; X, Y \right) - Uf(x) \right\| < \|U\|\varepsilon, \quad \forall x \in \mathcal{K}, \\
 &\Leftrightarrow \|Y \sigma(X_{\mathcal{P}}^\top B^\top C \tilde{x}) - f(x)\| < \varepsilon, \quad \forall x \in \mathcal{K}.
 \end{aligned} \quad (60)$$

Let $N_*^\sigma(x) := Y \sigma(X_{\mathcal{P}}^\top B^\top C \tilde{x}) = \sum_{i=1}^n y^{(i)} \sigma(\tilde{R}_i \tilde{x}) \in \mathcal{N}_*^\sigma$, where $n \in \mathbb{N}^+$, $y^{(i)} \in \mathbb{R}^{d_y}$ and $\tilde{R}_i \in \mathbb{R}^{d_x}$ (the i -th row of $\tilde{R} \in \mathbb{R}^{n \times d_x}$). The proof is divided into four steps:

Step (1): Approximating $f(x)$ Using $N^\sigma(x)$

For any $\varepsilon > 0$, there exists a neural network $N^\sigma(x) = A \sigma(Wx + b) = \sum_{i=1}^k a_i \sigma(w_i x + b_i) \in \mathcal{N}^\sigma$, with parameters $k \in \mathbb{N}^+$, $A \in \mathbb{R}^{d_y \times k}$, $b \in \mathbb{R}^k$, and $W \in \mathbb{R}^{k \times (d_x-1)}$ (where a_i and w_i denote the i -th column of A and the i -th row of W),

$$\|A \sigma(Wx + b) - f(x)\| < \frac{\varepsilon}{3}, \quad \forall x \in \mathcal{K}, \quad (61)$$

which is supported by Lemma 2.

Step (2): Approximating $N^\sigma(x)$ Using $N'(x)$

Using Lemma 8 and Lemma 13, a neural network $N^\sigma(x) = \sum_{i=1}^k a_i \sigma(w_i x + b_i) \in \mathcal{N}^\sigma$ can be perturbed into $N'(x) = \sum_{i=1}^k (q\sqrt{2} \pm l)_i \sigma(\tilde{w}_i x + \tilde{b}_i)$ (with $q_i \in \mathbb{N}^+$ and $l_i \in \mathbb{N}^+, i = 1, \dots, k$), such that for any $\varepsilon > 0$, there exists $0 < \delta < \frac{\varepsilon}{6M_1 k}$ satisfying:

$$\max\{\|a_i - (q\sqrt{2} \pm l)_i\|, M\|w_i - \tilde{w}_i\| + \|b - \tilde{b}\| \mid i = 1, \dots, k\} < \delta, \quad (62)$$

ensuring:

$$\|N^\sigma(x) - N'(x)\| = \left\| \sum_{i=1}^k a_i \sigma(w_i x + b_i) - \sum_{i=1}^k (q\sqrt{2} \pm l)_i \sigma(\tilde{w}_i x + \tilde{b}_i) \right\| < \frac{\varepsilon}{3}, \quad \forall x \in \mathcal{K}. \quad (63)$$

Step (3): Approximating $N'(x)$ Using $N_*^\sigma(x)$

Next, we show that $N_*^\sigma(x) = \sum_{i=1}^n y^{(i)} \sigma(\tilde{R}_i \tilde{x}) \in \mathcal{N}_*^\sigma$ can approximate $N'(x) = \sum_{i=1}^k (q\sqrt{2} \pm l)_i \sigma(\tilde{w}_i \tilde{x})$. As a demonstration, we approximate a single term $(q\sqrt{2} \pm l)_1 \sigma(\tilde{w}_1 \tilde{x})$.

Since S is dense in \mathbb{R}^{d_x} and $B^\top C$ is non-singular, the set $G := \{\tilde{R} \mid \tilde{R} = X_{\mathcal{P}}^\top B^\top C, X_{\mathcal{P}} \subset 2^S\}$ remains dense. Since $y^{(i)} \in \{1, -1, \sqrt{2}, 0\}$, we require $q_1 + l_1$ elements of \tilde{R}_i to approximate \tilde{w}_1 such that

$$\begin{aligned} & \left\| \sum_{j \in K_1} y^{(j)} \sigma(\tilde{R}_j \tilde{x}) - (q\sqrt{2} \pm l)_1 \sigma(\tilde{w}_1 \tilde{x}) \right\| \\ &= \left\| \sqrt{2} \sum_{j \in Q_1} \sigma(\tilde{R}_j \tilde{x}) \pm \sum_{j \in L_1} \sigma(\tilde{R}_j \tilde{x}) - (q\sqrt{2} \pm l)_1 \sigma(\tilde{w}_1 \tilde{x}) \right\| \\ &< \frac{\varepsilon}{3k}, \quad \forall x \in \mathcal{K}. \end{aligned} \quad (64)$$

Here, $\#(K_1) = q_1 + l_1$ and $K_1 = Q_1 \cup L_1$, where Q_1, L_1 are disjoint subsets of positive integer indices satisfying $\#(Q_1) = q_1$ and $\#(L_1) = l_1$. For this construction, we assign $y^{(j)} = \sqrt{2}$ for $j \in Q_1$ and $y^{(j)} = \pm 1$ for $j \in L_1$. For $j \notin \bigcup_{l=1}^k K_l$, we set $y^{(j)} = 0$. We then define $n = \max\{j \mid j \in \bigcup_{l=1}^k K_l\}$.

The multi-term approximation employs parallel construction via disjoint node subsets $K_i = Q_i \cup L_i$, where Q_i (q_i nodes) and L_i (l_i nodes) implement $\sqrt{2}$ and ± 1 coefficients respectively. Each term achieves:

$$\left\| \sum_{j \in K_i} y^{(j)} \sigma(\tilde{R}_j \tilde{x}) - (q\sqrt{2} \pm l)_i \sigma(\tilde{w}_i \tilde{x}) \right\| < \frac{\varepsilon}{3k}. \quad (65)$$

The complete network combines these approximations through:

$$\|N_*^\sigma(x) - N'(x)\| = \left\| \sum_{i=1}^n y^{(i)} \sigma(\tilde{R}_i \tilde{x}) - \sum_{i=1}^k (q\sqrt{2} \pm l)_i \sigma(\tilde{w}_i \tilde{x}) \right\| < \frac{\varepsilon}{3}, \quad \forall x \in \mathcal{K}. \quad (66)$$

Step (4): Combining Results

Combining all results, we have:

$$\begin{aligned} \|Y \sigma(X_{\mathcal{P}}^\top B^\top C \tilde{x}) - f(x)\| &= \|N_*^\sigma(x) - f(x)\| \\ &< \|N_*^\sigma(x) - N'(x)\| + \|N'(x) - N^\sigma(x)\| + \|N^\sigma(x) - f(x)\| \\ &< \varepsilon, \quad \forall x \in \mathcal{K}. \end{aligned} \quad (67)$$

The presented results for scalar outputs ($d_y = 1$) generalize naturally to multidimensional cases through component-wise approximation. For any continuous vector-valued function $f : \mathbb{R}^{d_x-1} \rightarrow \mathbb{R}^{d_y}$ defined on compact domain \mathcal{K} , we achieve

uniform approximation by independently handling each output dimension. The key observation is that the approximation problem decouples across dimensions - each component function f_j can be approximated by a scalar network $N_{*,j}^\sigma(x)$ satisfying:

$$\|N_{*,j}^\sigma(x) - f_j(x)\| < \frac{\varepsilon}{\sqrt{d_y}}, \quad \forall x \in \mathcal{K}. \quad (68)$$

The complete approximator assembles these component networks through vertical concatenation:

$$N_*^\sigma(x) = \begin{pmatrix} N_{*,1}^\sigma(x) \\ \vdots \\ N_{*,d_y}^\sigma(x) \end{pmatrix}, \quad \|N_*^\sigma(x) - f(x)\| < \varepsilon, \quad (69)$$

$$N_{*,i}^\sigma(x) = \sum_{j=1}^n y_j^{(i)} \sigma(\tilde{R}_i \tilde{x}), \quad (70)$$

where $y_j^{(i)}$ is the j -th row of the $y^{(i)}$. We require that the index sets satisfy $\mathcal{K}_i^{(o)} \cap \mathcal{K}_j^{(u)} = \emptyset$ for all $o, u, i, j \in \mathbb{N}^+$, where $\mathcal{K}_i^{(o)}$ denotes the index set constructed for the i -th term approximation in the o -th output dimension. Furthermore, each $y^{(j)}$ must have at most one non-zero element across its dimensions. This ensures we achieve uniform approximation by independently handling each output dimension. The proof is complete. \square

D.3. Example of Theorem 9

We present a concrete example with 2D input ($d_x = 2$) and 2D output ($d_y = 2$) to demonstrate the universal approximation capability of our architecture. Consider approximating a continuous vector-valued function defined on the unit square:

$$f(x_1, x_2) = \begin{bmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{bmatrix}, \quad (x_1, x_2) \in [0, 1]^2. \quad (71)$$

Our goal is to construct a Transformer-like module $T_{*,\mathcal{P}}^\sigma$ such that:

$$\left\| T_{*,\mathcal{P}}^\sigma \left(\begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}; X, Y \right) - f(x_1, x_2) \right\| < \varepsilon. \quad (72)$$

Step 1: Component-wise Approximation

For each output dimension, we independently construct approximating networks:

For the first component, there exists a single-hidden-layer network $N_1^\sigma(x) = A_1 \sigma(W_1 x + b_1)$ satisfying,

$$\sup_{x \in [0,1]^2} \|f_1(x_1, x_2) - N_1^\sigma(x)\| < \frac{\varepsilon}{6\sqrt{2}}. \quad (73)$$

For the second component, there exists $N_2^\sigma(x) = A_2 \sigma(W_2 x + b_2)$ with

$$\sup_{x \in [0,1]^2} \|f_2(x_1, x_2) - N_2^\sigma(x)\| < \frac{\varepsilon}{6\sqrt{2}}. \quad (74)$$

Step 2: Rational Perturbation We consider the rationally perturbed neural networks $N_i^\sigma(x), i = 1, 2$: $N_1' = \sum_{i=1}^{k_1} y_1^{(i)} \sigma(\tilde{w}_i \tilde{x}) = (3\sqrt{2} - 2) \sigma(\tilde{w}_i \tilde{x})$ and $N_2' = \sum_{i=1}^{k_2} y_2^{(i)} \sigma(\tilde{R}_i \tilde{x}) = (2\sqrt{2} + 1) \sigma(\tilde{w}_i \tilde{x})$ satisfying,

$$\sup_{x \in [0,1]^2} \|N_i^\sigma(x) - N_i'(x)\| < \frac{\varepsilon}{6\sqrt{2}}, \quad i = 1, 2. \quad (75)$$

Step 3: Architecture Realization The Transformer-like module is constructed as: $Y = \begin{bmatrix} \sqrt{2} & \sqrt{2} & \sqrt{2} & -1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \sqrt{2} & \sqrt{2} & 1 \end{bmatrix}$ and $\tilde{R} = [\approx \tilde{w}_1 \approx w_1 \approx \tilde{w}_1 \approx \tilde{w}_1 \approx \tilde{w}_1 \approx \tilde{w}_2 \approx \tilde{w}_2 \approx \tilde{w}_2]^\top$.

Such that,

$$N_*^\sigma = \left[\begin{array}{c} \sum_{i=1}^8 y_1^{(i)} \sigma(\tilde{R}_i \tilde{x}) \\ \sum_{i=1}^8 y_2^{(i)} \sigma(\tilde{R}_i \tilde{x}) \end{array} \right] = \left[\begin{array}{c} (3\sqrt{2} - 2)\sigma(\approx \tilde{w}_1 \tilde{x}) \\ (2\sqrt{2} + 1)\sigma(\approx \tilde{w}_2 \tilde{x}) \end{array} \right], \quad (76)$$

Step 4: Error Analysis The total error is controlled via:

$$\begin{aligned} \|f(x) - N_{*,i}^\sigma(x)\| &\leq \sqrt{\sum_{i=1}^2 (\|f_i(x) - N_i^\sigma(x)\| + \|N_i^\sigma(x) - N_i^\sigma(x)\| + \|N_i^\sigma(x) - N_{*,i}^\sigma(x)\|)^2} \\ &\leq \sqrt{\left(\frac{\varepsilon}{2\sqrt{2}}\right)^2 + \left(\frac{\varepsilon}{2\sqrt{2}}\right)^2} = \frac{\varepsilon}{2} < \varepsilon. \end{aligned} \quad (77)$$

Implementation Details The construction ensures dimensional independence through:

Table 2. Node allocation for 2D output example

Node Index	$y^{(i)}$	\tilde{R}_i	Purpose
1-3	$(\sqrt{2}, 0)$	$\approx w_1$	$3\sqrt{2}$ term for $\sigma(\tilde{w}_1 \tilde{x})$
4-5	$(-1, 0)$	$\approx w_1$	-2 term for $\sigma(\tilde{w}_1 \tilde{x})$
6-7	$(0, \sqrt{2})$	$\approx w_2$	$2\sqrt{2}$ term for $\sigma(\tilde{w}_2 \tilde{x})$
8	$(0, 1)$	$\approx w_2$	1 term for $\sigma(\tilde{w}_2 \tilde{x})$

Alternative Construction The Transformer-like module is constructed as: $Y = \begin{bmatrix} \sqrt{2} & \sqrt{2} & \sqrt{2} & \sqrt{2} & \sqrt{2} \\ -1 & -1 & 0 & 1 & 0 \end{bmatrix}$ and $\tilde{R} = [\approx \tilde{w}_1 \approx w_1 \approx \tilde{w}_1 \approx \tilde{w}_2 \approx \tilde{w}_2]^\top$. This alternative construction may potentially reduce the total number of required tokens (denoted as n in our proof). However, the corresponding analysis in higher-dimensional spaces becomes substantially more intricate from a theoretical perspective. Consequently, we elect to employ the disjoint index sets methodology throughout our proof, as it provides both analytical tractability and mathematical rigor.

D.4. Proof of Theorem 10

Lemma 14. For any continuous function $f : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ defined on a compact domain \mathcal{K} and $\varepsilon > 0$, there always exist a exp FNN $N^{\text{exp}}(x) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}, x \mapsto A \exp(Wx + b)$ satisfying

$$\|N^{\text{exp}}(x) - f(x)\| < \varepsilon, \quad \forall x \in \mathcal{K},$$

where $b = 0$ and all row vector of W are restricted in a neighborhood $B(w^*, \delta)$ with any prefixed $w^* \in \mathbb{R}^{d_x}$ and $\delta > 0$.

Proof. According to Stone-Weierstrass theorem we know that, for any continuous function f and $i = 1, \dots, d_y$ and $\varepsilon' > 0$, there exists a polynomial $P_i(x)$ which can approximate $\exp(-w^* \cdot x)(f(x))_i$, i.e.

$$\|P_i(x) - \exp(-w^* \cdot x)(f(x))_i\| < \varepsilon', \quad \forall x \in \mathcal{K}. \quad (78)$$

The inequation above indicates that

$$\|\exp(w^* \cdot x)P_i(x) - (f(x))_i\| < \|\exp(w^* \cdot x)\|\varepsilon' := \varepsilon, \quad \forall x \in \mathcal{K}. \quad (79)$$

Then we construct a single layer FNN with exponential activation function to approximate $\exp(w^* \cdot x)P_i(x)$. Without loss of generality, let us consider the first hidden neuron of a exponential FNN, whose output can be represented as

$$\exp(w_1 x_1 + \dots + w_{d_x} x_{d_x}). \quad (80)$$

And the multiple derivatives of $h(w) := \exp(w \cdot x)$ with respect to w_1, \dots, w_{d_x} is

$$\frac{\partial^{|\alpha|} h}{\partial w^\alpha} = \frac{\partial^\alpha h}{\partial w_1^{\alpha_1} \dots \partial w_{d_x}^{\alpha_{d_x}}}, \quad (81)$$

where $\alpha \in \mathbb{N}^{d_x}$ represents the index and $|\alpha| := \alpha_1 + \dots + \alpha_{d_x}$. Actually, the form of $\frac{\partial^\alpha h}{\partial w_1^{\alpha_1} \dots \partial w_{d_x}^{\alpha_{d_x}}}$ is a polynomial of $|\alpha|$ degree with respect to x_1, \dots, x_k times $h(w)$. Note that $\exp(w^* \cdot x)P_i(x)$ can be written as a finite sum of some multiple derivatives of $h(x)$, that is

$$\exp(w^* \cdot x)P_i(x) = \left(\sum_{\alpha \in \Lambda_i} a_\alpha \frac{\partial^{|\alpha|} h}{\partial w^\alpha} \right) \Big|_{w=w^*}, \quad (82)$$

where $\alpha \in \mathbb{N}^{d_x}$ is the index of multiple derivative and Λ_i is a finite multiple set of indexes. As for multiple derivatives, they can be approximated by finite difference method, and the approach of finite difference method can be done by a one hidden layer. For example,

$$\begin{aligned} x_1 \exp(w^* \cdot x) &= \frac{\partial h}{\partial w_1} \Big|_{w=w^*} \\ &= \frac{h(w^* + \lambda e_1) - h(w^*)}{\lambda} + R_1(\lambda, w^*) \\ &= \frac{1}{\lambda} \exp((w^* + \lambda e_1) \cdot x) - \frac{1}{\lambda} \exp(w^* \cdot x) + R_1(\lambda, w^*), \end{aligned} \quad (83)$$

and

$$\begin{aligned} x_1 x_2 \exp(w^* \cdot x) &= \frac{\partial^2 h}{\partial w_1 \partial w_2} \Big|_{w=w^*} \\ &= \frac{h(w^* + \lambda(e_1 + e_2)) - h(w^* + \lambda e_1) - h(w^* + \lambda e_2) + h(w^*)}{\lambda^2} + R_2(\lambda, w^*) \\ &= \frac{1}{\lambda^2} \exp((w^* + \lambda(e_1 + e_2)) \cdot x) - \frac{1}{\lambda^2} \exp((w^* + \lambda e_1) \cdot x) - \\ &\quad \frac{1}{\lambda^2} \exp((w^* + \lambda e_2) \cdot x) + \frac{1}{\lambda^2} \exp(w^* \cdot x) + R_2(\lambda, w^*), \end{aligned} \quad (84)$$

where $e_1 = (1, 0, 0, \dots, 0)$, $e_2 = (0, 1, 0, \dots, 0)$ are unit vectors and $R_1(\lambda, w^*)$ and $R_2(\lambda, w^*)$ are error terms with respect to λ and w^* . The error term $R_1(\lambda, w^*) = \lambda \frac{\partial^2 h}{\partial w_1^2} \Big|_{w=\xi}$ for some ξ between w^* and $w^* + \lambda e_1$. It is obvious that the partial differential term is bounded in $B(w^*, \delta)$, so the error can be controlled by λ . For $R_2(\lambda, w^*)$ it is similar. Equation (84) holds, as shown in Chapter X of (Boole, 2009).

Since λ is very small and the exponential terms $\exp(w^* \cdot x)$ only involve the parameters w^* , $w^* + \lambda e_1$ and $w^* + \lambda e_2$, which all lie within a small neighborhood of w^* the desired conclusion can be drawn, and this means we can actually restrict that all row vectors of W are restricted in the neighborhood $B(w^*, \delta)$. \square

Theorem 10. *If the set S is dense in $[-1, 1]^{d_x}$, then $\mathcal{T}_{*, \mathcal{P}}^{\text{ReLU}}$ is capable of achieving the UAP. Additionally, if S is only dense in a neighborhood $B(w^*, \delta)$ of a point $w^* \in \mathbb{R}^{d_x}$ with radius $\delta > 0$, then the class of transformers with exponential activation, $\mathcal{T}_{*, \mathcal{P}}^{\text{exp}}$, is capable of achieving the UAP.*

Proof. For the proof of ReLU case, we follow the same reasoning as in the previous one, noting that $\text{ReLU}(ax) = a \text{ReLU}(x)$ holds for any positive a . In the proof of Theorem 9, we construct a $\hat{T}_{*, \mathcal{P}}^{\text{ReLU}}$ to approximate a FNN $A \text{ReLU}(Wx + b)$. Here we can do the similar construction to find another $\hat{T}_{*, \mathcal{P}}^{\text{ReLU}}$ to approximate $tA \text{ReLU}(\frac{W}{t}x + b)$ as the second to the forth steps in Theorem 9, where t is big enough to make the elements in $\frac{W}{t}$ is small enough so $S = \{x_i + \mathcal{P}_x^{(j)} \mid x_i \in \mathcal{V}_x, i, j \in \mathbb{N}^+\}$ is dense in $[-1, 1]^{d_x}$ is sufficient. For the exponential, by using Lemma 14, we can do step the second to the forth steps in Theorem 9 again, which is similar to ReLU case. \square

E. General case for Theorem 7

It is important to note that Theorem 7 remains valid even without imposing specific constraints on the value, query, and key matrices V , Q , and K (e.g., the sparse partition described in equation (16)). Below, we outline the reasoning.

In general, we decompose the matrices as follows:

$$Q^\top K = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix}, V = \begin{bmatrix} D & E \\ F & U \end{bmatrix}, \quad (85)$$

where $M_{11}, D \in \mathbb{R}^{d_x \times d_x}$, $M_{12}, E \in \mathbb{R}^{d_x \times d_y}$, $M_{21}, F \in \mathbb{R}^{d_y \times d_x}$, and $M_{22}, U \in \mathbb{R}^{d_y \times d_y}$, respectively.

The attention mechanism can then be computed as:

$$\begin{aligned} \text{Attn}_{Q,K,V}^\sigma(Z) &= V Z M \sigma(Z^\top Q^\top K Z) \\ &= \begin{bmatrix} D & E \\ F & U \end{bmatrix} \begin{bmatrix} X & x \\ Y & 0 \end{bmatrix} \begin{bmatrix} I_n & 0 \end{bmatrix} \sigma \left(\begin{bmatrix} X^\top & Y^\top \\ x^\top & 0 \end{bmatrix} \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix} \begin{bmatrix} X & x \\ Y & 0 \end{bmatrix} \right) \\ &= \begin{bmatrix} DX + EY & 0 \\ FX + UY & 0 \end{bmatrix} \sigma \left(\begin{bmatrix} M & (X^\top M_{11} + Y^\top M_{21})x \\ x^\top (M_{11}X + M_{12}Y) & x^\top M_{11}x \end{bmatrix} \right), \end{aligned}$$

where M represents the matrix $X^\top M_{11}X + X^\top M_{12}Y + Y^\top M_{21}X + Y^\top M_{22}Y$. As a result, we have:

$$T^\sigma(\tilde{x}; X, Y) = (FX + UY) \sigma((X^\top M_{11} + Y^\top M_{21})\tilde{x}), \quad (86)$$

for the case of elementwise activations, and:

$$T^{\text{softmax}}(\tilde{x}; X, Y) = (FX + UY) \left[\text{softmax} \left(\begin{bmatrix} (X^\top M_{11} + Y^\top M_{21})\tilde{x} \\ \tilde{x}^\top M_{11}\tilde{x} \end{bmatrix} \right) \right]_{1:n}, \quad (87)$$

for the case of softmax activation.

By revisiting the definition of $T_*^\sigma(x; X, Y)$ and comparing T_*^σ and T_*^{softmax} presented here with those in Appendix C, it is clear that the only distinction lies in the specific matrices involved. Consequently, the proof process for Theorem 7 can be directly applied to obtain the same results.