

# EXPLORING NON-LINEARITY IN ATTENTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The representational ability of Transformer architectures arises from two sources of non-linearity: position-wise non-linearity via feed-forward layers and contextual non-linearity through self-attention. In this work, we revisit this distinction and pose two key questions: Can self-attention itself realize position-wise non-linearity? And is contextual non-linearity truly necessary? First, we prove that by appending a fixed bias vector into the input, stacked self-attention layers can approximate deep feed-forward networks—showing that attention alone is sufficient to implement position-wise non-linearity. Second, we prove that contextual non-linearity, i.e., input-dependent attention patterns, is not indispensable: fixed or even randomly chosen patterns, when combined with a feed-forward layer, can still produce context-sensitive representations of the same token in different contexts. As an application, we prove that a two-layer attention-only Transformer can accurately predict masked tokens in masked language modeling. Both theoretical analysis and empirical studies on pre-trained models and synthetic data support our theory.

## 1 INTRODUCTION

Non-linearity is a fundamental source of expressive power in modern deep learning architectures. At the core of the Transformer model (Vaswani et al., 2017) is the interactive composition of feed-forward layers and self-attention layers, which together give rise to two distinct forms of non-linearity. Contextual non-linearity emerges from the self-attention mechanism: through the combined effects of the dot product and the Softmax function, each self-attention layer induces a non-linear mapping from the input sequence to the attention weights. This non-linearity does not operate independently on each token, but instead arises from interactions among tokens within the same sequence. Such a mechanism enables the model to capture long-range dependencies, a capability widely regarded as one of the primary reasons why Transformers surpass recurrent and convolutional neural networks on a variety of sequence modeling tasks (Sherstinsky, 2020; O’shea & Nash, 2015).

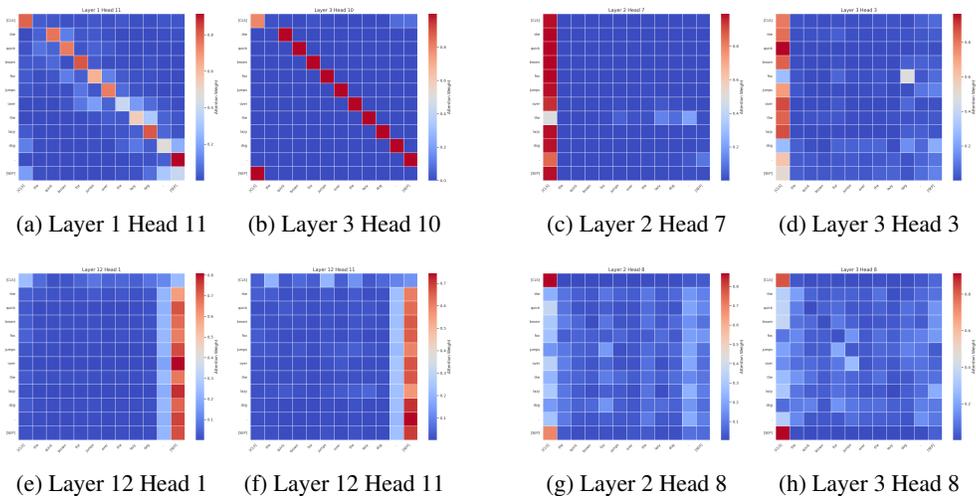
Self-attention, however, is not the sole contributor to the expressivity of Transformer models. Each self-attention layer is followed by a position-wise feed-forward layer, which accounts for roughly two-thirds of the model’s parameters. While the self-attention module computes a weighted linear aggregation of linearly transformed token representations, with weights determined by pairwise token similarity, the feed-forward network applies a non-linear transformation independently to each token embedding. This position-wise non-linearity complements the contextual non-linearity of attention, and together they form the core expressive machinery of the Transformer. In this work, we rethink self-attention mechanism through the lens of non-linear capacity. In particular, we aim to address the following two central questions:

**Is position-wise non-linearity exclusively introduced by feed-forward layers, or do certain attention heads also implicitly contribute to it?** Self-attention is designed to compute the interaction between tokens in different positions while feed-forward neural networks can be viewed as modeling interactions between tokens and a set of fixed tokens that encode global, input-independent information. Based on this similarity, a natural question is whether self-attention layers alone can behave like feed-forward neural networks, even deep ones. Prior works (Sukhbaatar et al., 2019; Huben & Morris, 2023) have investigated this problem from an empirical and theoretical perspective, respectively. As shown in Figure 2, we identify some attention heads in real-world language models that perform similar behavior to that of a feed-forward layer, meaning some attention heads implicitly implement position-wise non-linearity.

054 **Is contextual non-linearity indispensable, or can fixed and even random attention patterns still**  
 055 **yield effective performance in certain tasks.** Unlike feed-forward neural networks, which apply  
 056 exactly the same operations to each token, self-attention layers can incorporate contextual information.  
 057 In the standard attention mechanism, the attention pattern is instance-dependent and vary across  
 058 inputs. However, empirical results (Raganato et al., 2020; Tay et al., 2021; Kovaleva et al., 2019)  
 059 and also Figure 1 have shown that the attention patterns learned by Transformers are often limited in  
 060 diversity and can be replaced by instance-agnostic patterns without a decline in performance. These  
 061 empirical results motivates us to theoretically study the effectiveness of input-independent attention  
 062 mechanism.

063 Our contributions to answer the above two questions are summarized as follows:

- 064 • In Theorem 4.1, we extend the results in (Huben & Morris, 2023) to a more general setting  
 065 by proving that self-attention can approximate deep feed-forward neural networks. In  
 066 Appendix D, we provide empirical verification for Theorem 4.1. In addition, we identify that  
 067 in Figure 2, there are some attention heads in **Bert** (Devlin et al., 2019) mostly computing  
 068 the interaction between data tokens and special tokens [CLS] or [SEP], which perform  
 069 similar behavior to that of a feed-forward neural network, meaning that some attention heads  
 070 in real-world models implicitly provide position-wise non-linearity by letting the data tokens  
 071 interacting with some special tokens that are global and independent of input sequences.
- 072 • In masked language modeling, Transformers are required to predict the original tokens from  
 073 the token [MASK] based on the contexts it appears. As shown in Figure 3, the representation  
 074 of [MASK] tokens are influenced by different contexts. To verify the effectiveness of  
 075 instance-agnostic attention patterns, we theoretically prove that both a fixed attention pattern  
 076 and an arbitrary attention pattern equipped with a feed-forward layer can distinguish the same  
 077 tokens in different contexts. Specifically, we show that in Proposition 5.2 and 5.3, the same  
 078 token [MASK] in different sequences can be mapped to different values by self-attention  
 079 with input-independent attention patterns.
- 080 • As an application, we consider masked language modeling and prove that a two-layer  
 081 attention-only Transformer is sufficient to make accurate predictions for masked tokens,  
 082 where one layer implements position-wise non-linearity and the other aggregates contextual  
 083 information (see Theorem 6.1).
- 084 • Observations on pre-trained language models (see Figure 1, 2, 3) and experiments on  
 085 synthetic datasets (see Appendix D) support our theory.



103 **Figure 1: Repeated Attention Patterns.** We identify several attention patterns repeatedly occurring  
 104 in different attention heads in **Bert** (Devlin et al., 2019) (12 layers with 12 heads in each layer). (a)  
 105 and (b) shows that each token attends to its next one. In (c) and (d), each token mostly interacts with  
 106 the first token, that is, [CLS], while in (e) and (f), each token mainly interacts with the last token  
 107 [SEP]. Similar observation can be found in (Kovaleva et al., 2019; Clark et al., 2019).

## 2 RELATED WORK

**Understanding Attention Mechanism.** Tsai et al. (2019) showed that the attention mechanism can be reformulated as a kernel smoother over inputs, where the kernel scores correspond to token similarity. This perspective provides a more principled understanding of the individual components of Transformer’s attention mechanism, while also enabling competitive performance with deduced computation cost. Building on this view, a number of kernel-based variants have been developed, such as Performer (Choromanski et al., 2020), Skyformer (Chen et al., 2021), Linear Attention (Katharopoulos et al., 2020) and Sumformer (Alberti et al., 2023). Another line of research focuses on introducing sparsity into attention pattern to improve efficiency. Longformer, Beltagy et al. (2020), enables scalable processing of long documents without truncating the input. Zaheer et al. (2020) proposed BiGBIRD, a sparse attention mechanism that is linear in the number of tokens. Theoretically, they proved that BiGBIRD is a universal approximator and Turing complete. Similarly, Zhang et al. (2021) introduced a Softmax-free model, which avoids the categorical distribution and all negative attention scores are pruned out because of the use of ReLU activation function. Tay et al. (2021) proposed SYNTHESIZER, a model that learns synthetic attention weights directly, bypassing the computation of token-token interactions. Beyond the modification of attention mechanism, interpretability studies have investigated the representational properties of attention. Reif et al. (2019) explored the syntactic information encoded in attention matrices in BERT, while (Kovaleva et al., 2019) found that a limited set of recurring attention patterns emerge across different heads, and that manually disabling certain heads can even lead to performance improvement. In addition, Zhao et al. (2021) introduced a method to measure the degree of non-linearity of different components of Transformers. Concurrently, Dong et al. (2025) showed that Transformers with random attention patterns have a stable forward pass, which converges to a stochastic differential equation.

**Understanding Feed-Forward Layers.** The role of feed-forward layers have been closely examined. Geva et al. (2020) showed that feed-forward layers operate as Key-value memories, where each key correlates to textual patterns in the training example. Sukhbaatar et al. (2019) proposed a Transformer variant solely consisting of attention layers, in which the attention mechanism is augmented with persistent memory vectors that are able to substitute for feed-forward layers. Extending this, Huben & Morris (2023) theoretically proved that one attention layer can approximate one feed-forward layer with an extra bias vector appended to the input and special mask strategy. Xu et al. (2024) argued that the primary role of feed-forward layers is to provide non-linearity and they presented an improved FFN module, which is able to enrich non-linear signals while effectively reducing the hidden dimensionality. Kobayashi et al. (2023) showed that FFNs modify the input contextualization to emphasize specific types of linguistic compositions. Ikeda et al. (2025) investigated the layerwise importance of feed-forward layers in Transformer-based language models during pretraining, finding that concentrating FFNs in 70% of the consecutive middle layers consistently outperforms standard configurations for multiple downstream tasks. From another perspective, Bozic et al. (2023) analyzed the effectiveness of using standard shallow feed-forward neural networks to mimic the attention behavior. Zhang et al. (2020) empirically observed that replacing the upper self-attention layers in the encoder with feed-forward layers causes no performance degradation, and sometimes even yields minor gains. Finally, Liu et al. (2021) proposed a simple architecture, based on MLPs with gating, verifying that it can match Transformer’s performance across a variety of tasks.

## 3 PRELIMINARIES

### 3.1 TRANSFORMER

The Transformer architecture is a mapping that takes a sequence  $\mathbf{X} \in \mathbb{R}^{D \times N}$  composed of  $n$  tokens each with an embedding size of  $d$  as an input and outputs another sequence. In a Transformer neural network, there are two primary layers: a self-attention layer and a feed-forward layer. The attention layer mixes information across different positions in the sequence through dot products of token embeddings. Notably, Transformers utilize the Softmax function to transform each dot product into a probability or a weight, and the output of the attention layer becomes a linear combination of token embeddings according to these weights. Subsequently, the token-wise feed-forward layer operates independently on separate tokens without inter-token interaction.

A Transformer neural network is composed of alternating feed-forward layers and self-attention layers. An self-attention layer  $\mathcal{F}_{SA} : \mathbb{R}^{D \times N} \rightarrow \mathbb{R}^{D \times N}$  with  $H$  heads is defined as

$$\mathcal{F}_{SA}(\mathbf{X}) := \mathbf{X} + \sum_{i=1}^H \mathbf{W}_O^{(i)} \mathbf{W}_V^{(i)} \mathbf{X} \sigma_S \left[ \left( \mathbf{W}_K^{(i)} \mathbf{X} \right)^\top \left( \mathbf{W}_Q^{(i)} \mathbf{X} \right) + \mathbf{R}^{(i)} \right],$$

where  $\mathbf{W}_K^{(i)}, \mathbf{W}_Q^{(i)}, \mathbf{W}_V^{(i)} \in \mathbb{R}^{S \times D}$  and  $\mathbf{W}_O^{(i)} \in \mathbb{R}^{D \times S}$  are the weight matrices,  $\mathbf{R}^{(i)} \in \mathbb{R}^{N \times N}$  is the mask matrix, which satisfies  $\mathbf{R}_{i,j}^{(i)} = -\infty$  or 0,  $S$  is the head size, and  $\sigma_S : \mathbb{R}^D \rightarrow \mathbb{R}^D$  is the Softmax function with output  $[\sigma_S(\mathbf{x})]_i = e^{x_i} / \sum_{i=1}^d e^{x_i}$ ,  $i \in [D]$ , for any  $\mathbf{x} \in \mathbb{R}^D$ . Here,  $\sigma_S$  operates on each column of the input matrix. For mathematical simplicity, we ignore the layer normalization in attention layers.

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

$$\mathcal{F}_{FF}(\mathbf{X}) := \mathbf{X} + \mathbf{W}_2 \sigma_R \left[ \mathbf{W}_1 \mathbf{X} + \mathbf{b}_1 \mathbf{1}_N^\top \right] + \mathbf{b}_2 \mathbf{1}_N^\top \in \mathbb{R}^{D \times N},$$

where  $\mathbf{W}_1 \in \mathbb{R}^{W \times D}$  and  $\mathbf{W}_2 \in \mathbb{R}^{W \times D}$  are weight matrices with hidden dimension  $W$ , and  $\mathbf{b}_1 \in \mathbb{R}^W$ ,  $\mathbf{b}_2 \in \mathbb{R}^D$  are bias terms.  $\sigma_R$  denotes the element-wise ReLU activation function, i.e.,  $\sigma_R(x) := \max\{0, x\}$ .

The class of Transformer neural networks is then defined as

$$\mathcal{T}(D, H, S, W, L) := \left\{ \mathcal{E}_{out} \circ \mathcal{F}_{FF}^{(L)} \circ \mathcal{F}_{SA}^{(L)} \circ \dots \circ \mathcal{F}_{FF}^{(1)} \circ \mathcal{F}_{SA}^{(1)} \circ \mathcal{E}_{in} \right\},$$

where  $D$  is the embedding size,  $H$  is the number of heads,  $S$  is the head size,  $W$  is the hidden dimension in the feed-forward layers, and  $L$  is the number of Transformer layers, each consisting of a self-attention layer and a feed-forward layer.  $\mathcal{E}_{in}$  and  $\mathcal{E}_{out}$  represent the embedding layer and the decoding layer, respectively, which are two linear affine function. They are designed to change the input and output dimension of Transformers. In this work, we also consider attention-only Transformers, which discard all the feed-forward layers. The class of attention-only Transformer neural networks can be defined as

$$\mathcal{T}(D, H, S, L) := \left\{ \mathcal{E}_{out} \circ \mathcal{F}_{SA}^{(L)} \circ \dots \circ \mathcal{F}_{SA}^{(1)} \circ \mathcal{E}_{in} \right\}.$$

### 3.2 FEED-FORWARD NEURAL NETWORKS

We denote  $\mathcal{NN}_\sigma(N, L, \mathbb{R}^d \rightarrow \mathbb{R}^{d'})$  as the set of vector-valued functions  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  that can be represented by a feed-forward neural network (FFN) with width  $\leq N \in \mathbb{N}^+$ , depth  $\leq L \in \mathbb{N}^+$ , and activation function  $\sigma$ . The width of a FFN refers to the maximum number of neurons in the hidden layers and the depth corresponds to the number of hidden layers. For instance, suppose  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  is a vector-valued function realized by a feed-forward neural network with  $\sigma$  as the elementwise activation function. Then  $\phi$  can be expressed as

$$\phi = \mathcal{L}_L \circ \sigma \circ \mathcal{L}_{L-1} \circ \dots \circ \sigma \circ \mathcal{L}_1 \circ \sigma \circ \mathcal{L}_0,$$

where each  $\mathcal{L}_\ell$  is an affine linear map given by  $\mathcal{L}_\ell(\mathbf{x}) := \mathbf{W}_\ell \mathbf{x} + \mathbf{b}_\ell$  for  $\ell = 0, 1, \dots, L$ . Here,  $\mathbf{W}_\ell \in \mathbb{R}^{d_{\ell+1} \times d_\ell}$  and  $\mathbf{b}_\ell \in \mathbb{R}^{d_{\ell+1}}$  are the weight matrix and bias term, respectively, with  $d_0 = d$ ,  $d_1, \dots, d_L \in \mathbb{N}^+$ , and  $d_{L+1} = d'$ . Clearly,  $\phi \in \mathcal{NN}_\sigma(N, L, \mathbb{R}^d \rightarrow \mathbb{R}^{d'})$  where  $N = \max\{d_1, \dots, d_L\}$ .

In terms of the choice of the activation function, we consider  $\sigma_R$  (ReLU) and  $\sigma_L$  (SiLU) in this work, with definition given in the following

$$\sigma_R(x) := \max\{0, x\}, \quad \sigma_L(x) := \frac{x}{1 + e^{-x}}.$$

It is natural to extend the input of a FFN from vectors to matrices. We redefine each  $\mathcal{L}_\ell$  as  $\mathcal{L}_\ell(\mathbf{X}) := \mathbf{W}_\ell \mathbf{X} + \mathbf{b}_\ell \mathbf{1}_N^\top$ , where  $\mathbf{W}_\ell \in \mathbb{R}^{d_{\ell+1} \times d_\ell}$ ,  $\mathbf{b}_\ell \in \mathbb{R}^{d_{\ell+1}}$  and  $\mathbf{1}_N \in \mathbb{R}^{N \times 1}$  represents the all-1 vector. Clearly, a FFN with matrix input imposes the same operation on each column of the input. In the sequel, we do not distinguish between feed-forward neural networks with vector input and

feed-forward neural networks with matrix input. Similarly, a residual feed-forward neural network  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is defined as follow

$$\phi = \mathcal{L}_L \circ \mathcal{L}_{L-1} \circ \dots \circ \mathcal{L}_2 \circ \mathcal{L}_1,$$

where each  $\mathcal{L}_\ell$  given by  $\mathcal{L}_\ell(\mathbf{x}) := \mathbf{x} + \mathbf{W}_\ell^{(2)} \sigma(\mathbf{W}_\ell^{(1)} \mathbf{x} + \mathbf{b}_\ell^{(1)}) + \mathbf{b}_\ell^{(2)}$  for  $\ell = 1, \dots, L$ . Here,  $\mathbf{W}_\ell^{(1)} \in \mathbb{R}^{W \times d}$ ,  $\mathbf{W}_\ell^{(2)} \in \mathbb{R}^{d \times W}$ , and  $\mathbf{b}_\ell^{(1)} \in \mathbb{R}^W$ ,  $\mathbf{b}_\ell^{(2)} \in \mathbb{R}^d$  are the weight matrix and bias term, respectively. Let  $\mathcal{NN}_\sigma^{Res}(W, L, \mathbb{R}^d \rightarrow \mathbb{R}^d)$  denote the class of residual feed-forward neural networks with  $\sigma$  being the activation function.

### 3.3 FORMULATION OF MASKED LANGUAGE MODELING

Masked language modeling (MLM) is one of the most widely used pretraining objectives for Transformer-based language models. The key idea is to corrupt an input sequence by replacing some of its tokens with a special placeholder token, denoted by [MASK], and then train the model to recover the original tokens given the corrupted sequence. Moreover, sequences are also augmented with special tokens that serve structural and representational purposes. The token [CLS] is prepended to every sequence and the token [SEP] is appended to mark the end of a sequence, or to separate two sequences in next sentence prediction task. Here, we formally state the assumptions regarding the data.

**Assumption 3.1.** Let  $\mathcal{V} \subset \mathbb{R}^D$  be a finite vocabulary and [MASK] denote the mask token with [MASK]  $\notin \mathcal{V}$ . Each token  $\mathbf{x} \in \mathcal{V}$  corresponds to a token ID  $y \in [C]$ , where  $C = |\mathcal{V}|$ . Let  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)}) \subset \mathbb{R}^{D \times N} \times [C]$  be  $n$  input-output pairs with  $y^{(i)}$  denoting the token ID of the masked token. For any  $m_i \in \{1, 2, \dots, N\}$ , the  $m_i$ -th token in  $\mathbf{X}^{(i)}$  is masked by the [MASK] token, that is,  $\mathbf{X}_{:,m_i}^{(i)} = [\text{MASK}]$ . We assume that the following conditions are satisfied:

1. There exists  $r_1, r_2 > 0$  such that for any  $i \in [n]$  and  $j \in [N]$ ,  $r_1 \leq \|\mathbf{X}_{:,j}^{(i)}\| \leq r_2$ .
2. There exists  $\delta > 0$  such that for any  $i, j \in [n]$  and  $k, l \in [N]$ , either  $\mathbf{X}_{:,k}^{(i)} = \mathbf{X}_{:,l}^{(j)}$  or  $\|\mathbf{X}_{:,k}^{(i)} - \mathbf{X}_{:,l}^{(j)}\| \geq \delta$  holds.
3. For any  $i, j \in [n]$ ,  $\mathbf{X}^{(i)} \neq \mathbf{X}^{(j)}$  up to permutations.
4. There are no duplicated tokens in each  $\mathbf{X}^{(i)}$  for  $i \in [n]$ .

We provide a discussion on Assumption 3.1 in Appendix C, which shows that Assumption 3.1 is mild and holds in a wide range of practical settings. With Assumption 3.1 in hand, we define the masked language modeling task as follows.

**Definition 3.1** (Masked Language Modeling). Let  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)}) \subset \mathbb{R}^{D \times N} \times [C]$  be  $n$  input-output pairs, satisfying Assumption 3.1. The goal of masked language modeling is to find a Transformer  $\mathcal{F}$  such that  $\mathcal{F}(\mathbf{X}^{(i)})_{:,m_i} \approx y^{(i)}$  for any  $i \in [n]$ , which can be formulated as: Given any  $\varepsilon > 0$ , there exists a Transformer  $\mathcal{F} \in \mathcal{T}(D, H, S, W, L)$  for some  $D, H, S, W, L \in \mathbb{Z}^+$  such that

$$\left| \mathcal{F}(\mathbf{X}^{(i)})_{:,m_i} - y^{(i)} \right| < \varepsilon \quad \text{for any } i \in [n].$$

Note that Definition 3.1 describes a simplified variant of the standard masked language modeling task. In practice, approximately 15% of tokens in the sequence are selected for prediction, 80% of which are replaced by [MASK], 10% by a random token, and 10% remain unchanged. In contrast, our setting considers only the case where exactly one token is masked. Furthermore, in real-world models, the output of the Transformer is passed through a linear layer, and the model is trained by minimizing the cross-entropy loss. When the number of data points  $n$  and the vocabulary size  $C$  are large, this linear layer becomes limited to predict accurately the original tokens. Instead, we consider the setup by letting the output correspond directly to the token ID of the masked token. Finally, our formulation takes a constructive perspective, in the sense that we explicitly design a Transformer that satisfies the desired properties.

## 4 POSITION-WISE NON-LINEARITY

In vanilla Transformers, the primary source of non-linearity comes from the feed-forward layers, which apply column-wise non-linear transformations to each token independently. By contrast, self-attention layers are not explicitly designed to provide such position-wise non-linearity. Each token is updated as a weighted linear combination of linearly transformed token representations, and the only non-linear components in the self-attention mechanism are the dot-product similarity and the subsequent Softmax normalization, which serve to compute attention weights rather than to transform token representations directly.

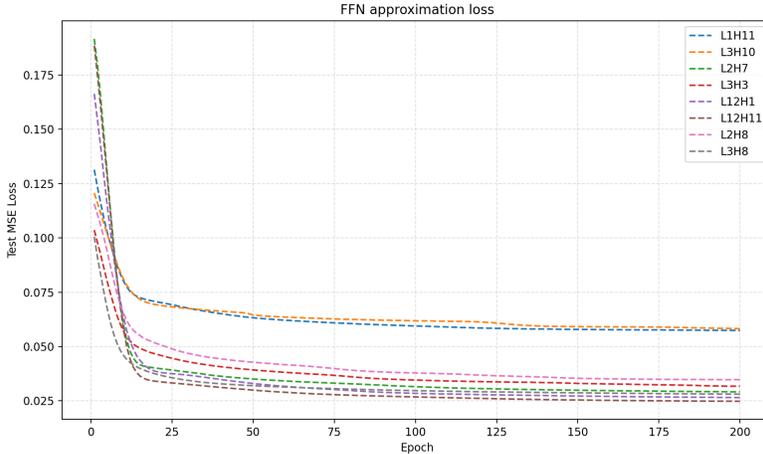


Figure 2: We train a feed-forward layer with hidden dimension 16 to approximate the attention heads mentioned in Figure 1 using a knowledge distillation method. The result shows that the approximation error of the Head 1 in Layer 1 and Head 10 in Layer 3 is higher than others, suggesting that their behavior deviates from that of a feed-forward layer. In contrast, the remaining heads primarily capture the interaction between data tokens and special tokens either [CLS] or [SEP], exhibiting behavior more similar to a feed-forward layer.

In this section, we show that by introducing a fixed, global token into the input sequence (similar to [CLS] or [SEP] tokens in Bert (Devlin et al., 2019)) and a well-designed mask matrix, which prevents the unnecessary interaction among tokens, the Softmax operation in self-attention can act analogously to an activation function and self-attention layers can approximate deep feed-forward neural networks. Under this construction, self-attention layers gain the ability to implement position-wise non-linearity, thereby broadening their functional role beyond contextual interactions. The following Theorem 4.1 summarizes our results.

**Theorem 4.1** (Approximation of feed-forward neural networks). *For any  $\varepsilon > 0$ ,  $M > 0$  and any FFN  $f \in \mathcal{NN}_{\sigma_R}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'})$ , there exists an attention-only Transformer  $\mathcal{F} \in \mathcal{T}(\max\{W, D\} + 2, \max\{W, D\} + 1, 1, L)$  such that*

$$\|\mathcal{F}(\mathbf{X}) - f(\mathbf{X})\|_{\infty} < \varepsilon \quad \text{for any } \mathbf{X} \in [-M, M]^{D \times N}.$$

The embedding layer in  $\mathcal{F}$  is defined by

$$\mathcal{E}_{in}(\mathbf{X}) := \begin{pmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{1}_{1 \times N} & \mathbf{0} \\ \mathbf{0}_{(W-D)^+ \times N} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(\max\{W, D\} + 2) \times (N + 1)}.$$

The proof of Theorem 4.1 is postponed to Appendix E. Note that after the embedding layer, an extra vector  $(\mathbf{0} \ \cdots \ \mathbf{0} \ 1)^{\top}$  is appended to the sequence, which acts as a global bias. With a special design of the mask matrix, we ensure that each data token only interacts with itself and this bias vector and we let the value vector (i.e.  $\mathbf{W}_V \mathbf{x}$ ) of this bias become  $\mathbf{0}$ , meaning that this bias vector contributes no information to each data token, but only bias the attention weights such that it can be regarded as a activation function. Let’s take the case  $N = 2$  as an example: suppose that

$\mathbf{X} = (\mathbf{x}_1 \quad \mathbf{x}_2) \in \mathbb{R}^{D \times 2}$ , the output of a self-attention layer without skip connection given input  $\mathbf{X}$  is  $(a_1 \mathbf{W}_V \mathbf{x}_1 + a_2 \mathbf{W}_V \mathbf{x}_2 \quad \mathbf{0})$ , where  $a_1, a_2$  are attention weights and we only let  $\mathbf{x}_1$  interact with  $\mathbf{x}_1$  and  $\mathbf{x}_2$  by using a mask matrix. If  $\mathbf{W}_V \mathbf{x}_2 = \mathbf{0}$ , we have  $a_1 \mathbf{W}_V \mathbf{x}_1 + a_2 \mathbf{W}_V \mathbf{x}_2 = a_1 \mathbf{W}_V \mathbf{x}_1$ , meaning that  $\mathbf{x}_2$  contributes nothing to  $\mathbf{x}_1$  but only bias its attention weight. Note that  $a_1 \mathbf{W}_V \mathbf{x}_1$  is a non-linear Transformation of  $\mathbf{x}_1$  due to the Softmax function contained in  $a_1$ .

**Remark 4.1.** In Figure 1, we observe that some attention heads in **Bert** (12 layers with 12 heads in each layer) have a special attention pattern, in which the data tokens mostly interacts with [CLS] or [SEP]. Moreover, in Figure 10, we observe that the norm of the value vector of [CLS] or [SEP] is small in certain layers. As a result, in Figure 2, we train a FFN with hidden dimension 16 to mimic these attention heads, finding that their behavior is closer to that of a FFN. Since some works (Melas-Kyriazi, 2021; Bozic et al., 2023; Zhang et al., 2020) proposed to replace some self-attention heads by feed-forward layers, our result may serve as a possible selection standard to determine which heads can be substituted.

**Remark 4.2.** We observe that Theorem 4.1 is similar to prompt tuning (Lester et al., 2021) in several aspects. In the task of prompt tuning, the parameters of the model are kept frozen and the parameters in prompts are trainable. By contrast, in Theorem 4.1, the parameters of the model are trainable (depend on the feed-forward neural network to be approximated), while the appended vector remains fixed. Another correspondence lies in the training stage: prompt tuning is applied after pre-training the model while Theorem 4.1 is realized during pre-training. The results of (Nakada et al., 2025) can therefore be regarded as the prompt tuning counterpart of ours. One additional insight is that one could relax our construction by making the appended vector into trainable or by extending it to a longer sequence of vectors, similar to the case in prompt tuning.

The following corollary demonstrates that the appended vector can be selected arbitrarily, subject only to a mild constraint.

**Corollary 4.1.** *For any  $\varepsilon > 0$ ,  $M > 0$  and any FFN  $\mathbf{f} \in \mathcal{NN}_{\sigma_R}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'})$ , there exists an attention-only Transformer  $\mathcal{F} \in \mathcal{T}(\max\{W, D\} + 2, \max\{W, D\} + 1, 1, L)$  such that*

$$\|\mathcal{F}(\mathbf{X}) - \mathbf{f}(\mathbf{X})\|_\infty < \varepsilon \quad \text{for any } \mathbf{X} \in [-M, M]^{D \times N}.$$

The embedding layer in  $\mathcal{F}$  is defined by

$$\mathcal{E}_{in}(\mathbf{X}) := \begin{pmatrix} \mathbf{X} \\ \mathbf{1}_{1 \times N} \\ \mathbf{0}_{(W-D)^+ \times N} \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{(\max\{W, D\} + 2) \times (N + 1)},$$

where  $\mathbf{v} \in \mathbb{R}^{\max\{W, D\} + 2}$  is an arbitrary vector with the last element  $v_{-1} \neq 0$ .

## 5 CONTEXTUAL NON-LINEARITY

Due to the combined effects of Softmax function and dot product, the self-attention mechanism induces a non-linear mapping from token representations to attention weights. Importantly, the non-linearity is not applied to tokens in isolation but rather arises from interaction among tokens within the same sequence. Since the attention weights are instance-dependent, each input will result in a different attention pattern. We refer to this ability to generate dynamic, input-specific attention patterns as **contextual non-linearity**. In contrast, we define the contextual linearity as the case where the attention pattern is fixed and does not vary with the input. In this setting, the way tokens attend to each other is purely linear, as the attention weights remain constant across all input sequences. This motivates the following definition of a self-attention layer with a constant attention pattern.

**Definition 5.1** (Self-attention with a Fixed Attention Pattern). Given any matrix  $\mathbf{A} \in \mathbb{R}^{N \times N}$  with each element  $A_{i,j} > 0$ . An self-attention layer  $\mathcal{F}_{SA}^{\mathbf{A}}$  with attention pattern  $\mathbf{A}$  and  $H$  heads is defined as

$$\mathcal{F}_{SA}^{\mathbf{A}}(\mathbf{X}) := \mathbf{X} + \sum_{i=1}^H \mathbf{W}_O \mathbf{W}_V \mathbf{X} \mathbf{A}.$$

In the remainder of this work, we allow Transformers to include both standard self-attention layers and fixed-pattern self-attention layers. We will use the superscript to distinguish between the two.

In masked language modeling, there are several tokens masked by a special token denoted by [MASK], which is universal across all input sequences. Transformers are trained to predict the original token before being masked based on the contexts in which it appears. This is a task requiring Transformers to distinguish the same token in different contexts because the tokens before being masked can be different. This is a concept first proposed by (Yun et al., 2019; Kim et al., 2022; Kajitsuka & Sato, 2023). In Figure 3, we trace the cosine similarity between the [MASK] tokens in different contexts across layers. It is shown that the representation of [MASK] tokens becomes more and more distinguishable as the layer goes deeper. To see the ability to differentiate tokens of

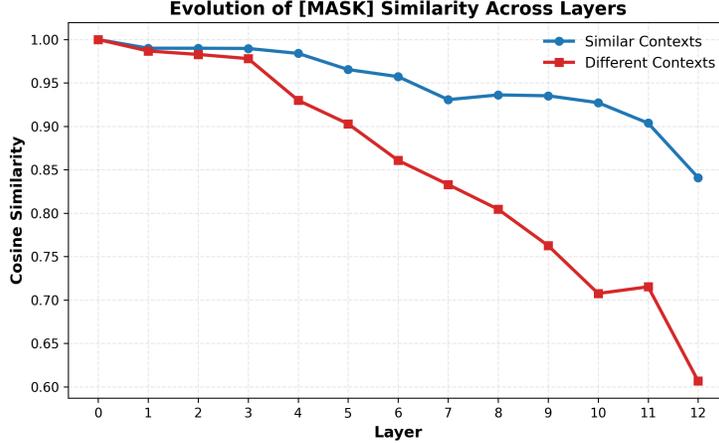


Figure 3: **Cosine similarity between [MASK] tokens in different contexts.** We trace the cosine similarity between [MASK] tokens in Bert (12 layers with 12 heads in each layer). In blue line, we consider two similar contexts, while in red line, we consider two totally different contexts. All [MASK] tokens are placed in the first position of the sequence and the sequence length is chosen to be the same to eliminate the influence brought by position. The result shows that the difference between the [MASK] representation becomes large as the depth goes deeper, especially when the contexts are significantly different.

contextual non-linearity, the following proposition proves that one standard self-attention layer with one head can distinguish the [MASK] in different contexts.

**Proposition 5.1** (Standard Self-attention). *For any input-label pairs  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)})$  satisfying Assumption 3.1, there exists a self-attention layer  $\mathcal{F}_{SA} \in \mathcal{T}(D, 1, 1, 1)$  such that*

$$\mathcal{F}_{SA}(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}_{SA}(\mathbf{X}^{(j)})_{:,m_j} \quad \text{for any } i \neq j \in [n].$$

The proof technique basically follows (Kajitsuka & Sato, 2023), in which we resort to the separation ability of Boltzmann operator and its close relationship with Softmax function. Then, it is natural to ask that whether a self-attention layer with a fixed attention pattern remains this property. The following proposition provides a deterministic answer.

**Proposition 5.2** (Fixed Attention Pattern). *For any input-label pairs  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)})$  satisfying Assumption 3.1, there exists a self-attention layer  $\mathcal{F}_{SA}^A \in \mathcal{T}(D, 1, 1, 1)$  with a fixed attention pattern  $A$  such that*

$$\mathcal{F}_{SA}^A(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}_{SA}^A(\mathbf{X}^{(j)})_{:,m_j} \quad \text{for any } i \neq j \in [n].$$

Note that the attention pattern  $A$  in Proposition 5.2 is trainable, meaning that it depends on the training dataset  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)})$ . So, the total number of parameters becomes  $O(N^2)$  instead of  $O(DN)$  in Proposition 5.1. In the following, we consider any randomly chosen attention pattern, which is independent of the training dataset. Our result shows that any attention pattern is as powerful as trainable one with the help of a feed-forward layer.

**Proposition 5.3** (Random Attention Pattern). *For any input-label pairs  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)})$  satisfying Assumption 3.1, and any attention pattern  $A \in \mathbb{R}^{N \times N}$  with  $A_{i,j} > 0$  for any  $i, j \in [N]$ , there exists a Transformer  $\mathcal{F} = \mathcal{E}_{out} \circ \mathcal{F}_{SA}^A \circ \mathcal{F}_{FF} \circ \mathcal{E}_{in} \in \mathcal{T}(\max\{3(n-1)N, D\}, 1, 1, 3(n-1)N, 1)$  such that*

$$\mathcal{F}(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}(\mathbf{X}^{(j)})_{:,m_j} \quad \text{for any } i \neq j \in [n].$$

All the proofs of this section are postponed to Appendix F. In the proof of Proposition 5.3, the feed-forward layer maps each token into a standard basis in a higher space, whose dimension depends on the total number of different tokens in the vocabulary. That is why the hidden dimension scales with the number of tokens  $(n - 1)N$ . According to the results in (Vardi et al., 2021; Kajitsuka & Sato, 2024), if we use a deep feed-forward neural network with constant width, the depth only grows as  $O(\sqrt{nN})$ , which is much more parameter-efficient. As shown in Figure 3, the difference between the representation of the same [MASK] token in distinct contexts becomes larger as the layer goes deeper. However, in Proposition 5.3, we only use one self-attention layer and one feed-forward layer to distinguish [MASK] tokens.

## 6 MASKED LANGUAGE MODELING

As an application, we consider masked language modeling with attention-only Transformers. In Definition 3.1, we need to predict the masked tokens based on its contexts. Since the [MASK] token is universal across corrupted sequences, we first use one self-attention layer to map each [MASK] token to distinct vectors, and then use another self-attention layer to implement point fitting. Our results are summarized in the following theorem.

**Theorem 6.1.** *For any  $\varepsilon > 0$  and  $n$  input-label pairs  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)}) \subset \mathbb{R}^{D \times N} \times [C]$ , which satisfy the Assumption 3.1, there exists an attention-only Transformer  $\mathcal{F} \in \mathcal{T}(\max\{3n, D\} + 2, \max\{3n, D\} + 1, 1, 2)$  such that*

$$|\mathcal{F}(\mathbf{X}^{(i)})_{:,m_i} - \mathbf{y}^{(i)}| < \varepsilon \quad \text{for any } i = 1, \dots, n.$$

The proof of Theorem 6.1 is in Appendix G. Theorem 6.1 shows that a two-layer attention-only Transformer can accurately predict the masked tokens to any precision. The first Transformer layer is constructed from Proposition 5.1, which distinguishes the same [MASK] token in different contexts by mapping them to different values. The second self-attention layer is augmented with a fixed vector attended to the input and approximately implement a FFN from Lemma H.4, which implements "key-value memory". The input represents the key and it is mapped to the corresponding value. To be specific, for  $n$  key-value pairs,  $\{\mathbf{x}^{(1)} \rightarrow y^{(1)}, \dots, \mathbf{x}^{(n)} \rightarrow y^{(n)}\}$  with  $\mathbf{x}^{(i)} \in \mathbb{R}^D, y^{(i)} \in \mathbb{R}$  and  $\mathbf{x}^{(i)} \neq \mathbf{x}^{(j)}$  for any  $i \neq j \in [n]$ , there exists a FFN  $\mathbf{f} \in \mathcal{NN}_{\sigma_R}(3n, 1, \mathbb{R}^D \rightarrow \mathbb{R})$  such that

$$\mathbf{f}(\mathbf{x}^{(i)}) = y^{(i)} \quad \text{for any } i \in [n].$$

Note that it is natural to restate Theorem 6.1 based on Proposition 5.2 and 5.3, where we only need to modify the first self-attention layer. As a result, all attention patterns in this Transformer are independent of input sequences.

## 7 CONCLUSION

In this paper, we study the position-wise and contextual non-linearity of attention mechanism. Firstly, we prove that by augmenting the input with a fixed bias vector, a stack of self-attention layers are able to approximate deep feed-forward neural networks. This demonstrates that attention layers alone are capable of implementing position-wise non-linearity. We also identify some attention heads in real-world pre-trained language models, which perform similar behavior to that of a FFN. Furthermore, we prove that the instance-dependent interaction patterns of self-attention are not essential for distinguishing the same token in different contexts, calling into question the necessity of contextual non-linearity in certain settings. By integrating both position-wise and contextual non-linearity, we prove that a two-layer attention-only Transformer suffices to make accurate predictions in masked language modeling. Finally, our empirical findings and validation support our theory.

**Limitation:** Sparsity and low-rank condition have not been considered in the construction of attention patterns in Proposition 5.2 and 5.3. Moreover, studying the gradient dynamic during masked language modeling would offer a complementary perspective for our constructive results and yield further insights. Since our work only focuses the pre-training process, fine-tuning or prompt tuning and the performance on downstream tasks should also be studied.

## REFERENCES

- 486  
487  
488 Silas Alberti, Niclas Dern, Laura Thesing, and Gitta Kutyniok. Sumformer: Universal approximation  
489 for efficient transformers. In *Topological, Algebraic and Geometric Learning Workshops 2023*, pp.  
490 72–86. PMLR, 2023.
- 491 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.  
492 *arXiv preprint arXiv:2004.05150*, 2020.
- 493 Vukasin Bozic, Danilo Dordevic, Daniele Coppola, Joseph Thommes, and Sidak Pal Singh.  
494 Rethinking attention: Exploring shallow feed-forward neural networks as an alternative to attention  
495 layers in transformers. *arXiv preprint arXiv:2311.10642*, 2023.
- 496  
497 Xingwu Chen and Difan Zou. What can transformer learn with varying depth? case studies on  
498 sequence learning tasks. *arXiv preprint arXiv:2404.01601*, 2024.
- 499 Yifan Chen, Qi Zeng, Heng Ji, and Yun Yang. Skyformer: Remodel self-attention with gaussian kernel  
500 and nystrom method. *Advances in Neural Information Processing Systems*, 34:2122–2135,  
501 2021.
- 502  
503 Jingpu Cheng, Qianxiao Li, Ting Lin, and Zuowei Shen. A unified framework on the universal  
504 approximation of transformer-type architectures. *arXiv preprint arXiv:2506.23551*, 2025.
- 505 Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas  
506 Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention  
507 with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- 508  
509 Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. What does bert look at?  
510 an analysis of bert’s attention. *arXiv preprint arXiv:1906.04341*, 2019.
- 511 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
512 bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of*  
513 *the North American chapter of the association for computational linguistics: human language*  
514 *technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
- 515  
516 Yihe Dong, Lorenzo Noci, Mikhail Khodak, and Mufan Li. Attention retrieves, mlp memorizes:  
517 Disentangling trainable components in the transformer. *arXiv preprint arXiv:2506.01115*, 2025.
- 518  
519 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are  
520 key-value memories. *arXiv preprint arXiv:2012.14913*, 2020.
- 521  
522 Iryna Gurevych, Michael Kohler, and Gözde Gül Şahin. On the rate of convergence of a classifier  
523 based on a transformer encoder. *IEEE Transactions on Information Theory*, 68(12):8139–8155,  
524 2022.
- 525  
526 Alexander Havrilla and Wenjing Liao. Understanding scaling laws with statistical and approximation  
527 theory for transformer neural networks on intrinsically low-dimensional data. *Advances in Neural*  
528 *Information Processing Systems*, 37:42162–42210, 2024.
- 529  
530 Jerry Yao-Chieh Hu, Wei-Po Wang, Ammar Gilani, Chenyang Li, Zhao Song, and Han Liu.  
531 Fundamental limits of prompt tuning transformers: Universality, capacity and efficiency. *arXiv*  
532 *preprint arXiv:2411.16525*, 2024.
- 533  
534 Jerry Yao-Chieh Hu, Hude Liu, Hong-Yu Chen, Weimin Wu, and Han Liu. Universal approximation  
535 with softmax attention, 2025. URL <https://arxiv.org/abs/2504.15956>.
- 536  
537 Robert Huben and Valerie Morris. Attention-only transformers and implementing mlps with attention  
538 heads. *arXiv preprint arXiv:2309.08593*, 2023.
- 539  
540 Wataru Ikeda, Kazuki Yano, Ryosuke Takahashi, Jaesung Lee, Keigo Shibata, and Jun Suzuki.  
541 Layerwise importance analysis of feed-forward networks in transformer-based language models.  
542 *arXiv preprint arXiv:2508.17734*, 2025.
- 543  
544 Haotian Jiang and Qianxiao Li. Approximation rate of the transformer architecture for sequence  
545 modeling. *Advances in Neural Information Processing Systems*, 37:68926–68955, 2024.

- 540 Yuling Jiao, Yanming Lai, Defeng Sun, Yang Wang, and Bokai Yan. Approximation bounds for  
541 transformer networks with application to regression. *arXiv preprint arXiv:2504.12175*, 2025a.  
542
- 543 Yuling Jiao, Yanming Lai, Yang Wang, and Bokai Yan. Transformers can overcome the curse  
544 of dimensionality: A theoretical study from an approximation perspective. *arXiv preprint*  
545 *arXiv:2504.13558*, 2025b.
- 546 Tokio Kajitsuka and Issei Sato. Are transformers with one layer self-attention using low-rank weight  
547 matrices universal approximators? *arXiv preprint arXiv:2307.14023*, 2023.  
548
- 549 Tokio Kajitsuka and Issei Sato. Optimal memorization capacity of transformers. *arXiv preprint*  
550 *arXiv:2409.17677*, 2024.
- 551 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns:  
552 Fast autoregressive transformers with linear attention. In *International conference on machine*  
553 *learning*, pp. 5156–5165. PMLR, 2020.  
554
- 555 Guolin Ke, Di He, and Tie-Yan Liu. Rethinking positional encoding in language pre-training. *arXiv*  
556 *preprint arXiv:2006.15595*, 2020.
- 557 Junghwan Kim, Michelle Kim, and Barzan Mozafari. Provable memorization capacity of transformers.  
558 In *The Eleventh International Conference on Learning Representations*, 2022.  
559
- 560 Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. Analyzing feed-forward blocks  
561 in transformers through the lens of attention maps. *arXiv preprint arXiv:2302.00456*, 2023.
- 562 Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. Revealing the dark secrets of  
563 bert. *arXiv preprint arXiv:1908.08593*, 2019.  
564
- 565 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt  
566 tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- 567 Hanxiao Liu, Zihang Dai, David So, and Quoc V Le. Pay attention to mlps. *Advances in neural*  
568 *information processing systems*, 34:9204–9215, 2021.  
569
- 570 Hude Liu, Jerry Yao-Chieh Hu, Zhao Song, and Han Liu. Attention mechanism, max-affine partition,  
571 and universal approximation. *arXiv preprint arXiv:2504.19901*, 2025.
- 572 Sadeqh Mahdavi, Renjie Liao, and Christos Thrampoulidis. Memorization capacity of multi-head  
573 attention in transformers. *arXiv preprint arXiv:2306.02010*, 2023.  
574
- 575 Luke Melas-Kyriazi. Do you even need attention? a stack of feed-forward layers does surprisingly  
576 well on imagenet. *arXiv preprint arXiv:2105.02723*, 2021.
- 577 Ryumei Nakada, Wenlong Ji, Tianxi Cai, James Zou, and Linjun Zhang. A theoretical framework for  
578 prompt engineering: Approximating smooth functions with transformer prompts. *arXiv preprint*  
579 *arXiv:2503.20561*, 2025.  
580
- 581 Keiron O’shea and Ryan Nash. An introduction to convolutional neural networks. *arXiv preprint*  
582 *arXiv:1511.08458*, 2015.
- 583 Sejun Park, Jaeho Lee, Chulhee Yun, and Jinwoo Shin. Provable memorization via deep neural  
584 networks using sub-linear parameters. In *Conference on learning theory*, pp. 3627–3661. PMLR,  
585 2021.
- 586 Aleksandar Petrov, Philip HS Torr, and Adel Bibi. Prompting a pretrained transformer can be a  
587 universal approximator. *arXiv preprint arXiv:2402.14753*, 2024.  
588
- 589 Alessandro Raganato, Yves Scherrer, and Jörg Tiedemann. Fixed encoder self-attention patterns in  
590 transformer-based machine translation. *arXiv preprint arXiv:2002.10260*, 2020.  
591
- 592 Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B Viegas, Andy Coenen, Adam Pearce, and  
593 Ben Kim. Visualizing and measuring the geometry of bert. *Advances in neural information*  
*processing systems*, 32, 2019.

- 594 Alex Sherstinsky. Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm)  
595 network. *Physica D: Nonlinear Phenomena*, 404:132306, 2020.  
596
- 597 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced  
598 transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 599 Sainbayar Sukhbaatar, Edouard Grave, Guillaume Lample, Herve Jegou, and Armand Joulin.  
600 Augmenting self-attention with persistent memory. *arXiv preprint arXiv:1907.01470*, 2019.  
601
- 602 Shokichi Takakura and Taiji Suzuki. Approximation and estimation ability of transformers for  
603 sequence-to-sequence functions with infinite dimensional input. In *International Conference on*  
604 *Machine Learning*, pp. 33416–33447. PMLR, 2023.
- 605 Naoki Takeshita and Masaaki Imaizumi. Approximation of permutation invariant polynomials by  
606 transformers: Efficient construction in column-size. *arXiv preprint arXiv:2502.11467*, 2025.  
607
- 608 Yi Tay, Dara Bahri, Donald Metzler, Da-Cheng Juan, Zhe Zhao, and Che Zheng. Synthesizer:  
609 Rethinking self-attention for transformer models. In *International conference on machine learning*,  
610 pp. 10183–10192. PMLR, 2021.
- 611 Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan  
612 Salakhutdinov. Transformer dissection: a unified understanding of transformer’s attention via the  
613 lens of kernel. *arXiv preprint arXiv:1908.11775*, 2019.  
614
- 615 Gal Vardi, Gilad Yehudai, and Ohad Shamir. On the optimal memorization power of relu neural  
616 networks. *arXiv preprint arXiv:2110.03187*, 2021.
- 617 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
618 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*  
619 *systems*, 30, 2017.  
620
- 621 Yihan Wang, Jatin Chaudhan, Wei Wang, and Cho-Jui Hsieh. Universality and limitations of prompt  
622 tuning. *Advances in Neural Information Processing Systems*, 36:75623–75643, 2023.
- 623 Yixing Xu, Chao Li, Dong Li, Xiao Sheng, Fan Jiang, Lu Tian, Ashish Sirasao, and Emad Barsoum.  
624 Enhancing vision transformer: Amplifying non-linearity in feedforward network module. In  
625 *Forty-first International Conference on Machine Learning*, 2024. URL [https://openreview.](https://openreview.net/forum?id=Nv0q2jdw00)  
626 [net/forum?id=Nv0q2jdw00](https://openreview.net/forum?id=Nv0q2jdw00).
- 627
- 628 Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are  
629 transformers universal approximators of sequence-to-sequence functions? In *International*  
630 *Conference on Learning Representations*, 2019.
- 631 Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and  
632 Sanjiv Kumar.  $O(n)$  connections are expressive enough: Universal approximability of sparse  
633 transformers. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.),  
634 *Advances in Neural Information Processing Systems*, volume 33, pp. 13783–13794. Curran  
635 Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/](https://proceedings.neurips.cc/paper_files/paper/2020/file/9ed27554c893b5bad850a422c3538c15-Paper.pdf)  
636 [paper/2020/file/9ed27554c893b5bad850a422c3538c15-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/9ed27554c893b5bad850a422c3538c15-Paper.pdf).
- 637
- 638 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago  
639 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for  
640 longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- 641 Biao Zhang, Ivan Titov, and Rico Sennrich. Sparse attention with linear units. *arXiv preprint*  
642 *arXiv:2104.07012*, 2021.
- 643
- 644 Shijun Zhang, Jianfeng Lu, and Hongkai Zhao. Deep network approximation: Beyond relu to  
645 diverse activation functions. *Journal of Machine Learning Research*, 25(35):1–39, 2024. URL  
646 <http://jmlr.org/papers/v25/23-0912.html>.
- 647
- Shucong Zhang, Erfan Loweimi, Peter Bell, and Steve Renals. When can self-attention be replaced  
by feed forward layers? *arXiv preprint arXiv:2005.13895*, 2020.

648 Sumu Zhao, Damián Pascual, Gino Brunner, and Roger Wattenhofer. Of non-linearity and  
649 commutativity in bert. In *2021 International Joint Conference on Neural Networks (IJCNN)*,  
650 pp. 1–8. IEEE, 2021.  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

# Appendix

<b>A</b>	<b>Notation Table</b>	<b>15</b>
<b>B</b>	<b>Additional Related Work</b>	<b>16</b>
<b>C</b>	<b>A Discussion on Assumption 3.1</b>	<b>17</b>
<b>D</b>	<b>Experiments</b>	<b>18</b>
	D.1 Proof-of-Concept Experiments on Theorem 4.1 . . . . .	18
	D.2 Proof-of-Concept Experiments on Theorem 6.1 . . . . .	18
<b>E</b>	<b>Proof of Section 4</b>	<b>23</b>
	E.1 Proof of Theorem 4.1 . . . . .	23
	E.2 Proof of Lemma E.1 . . . . .	26
	E.3 proof of Lemma E.2 . . . . .	26
	E.4 Proof of Lemma E.3 . . . . .	28
	E.5 Proof of Lemma E.4 . . . . .	29
	E.6 Proof of Corollary 4.1 . . . . .	30
<b>F</b>	<b>Proof of Section 5</b>	<b>32</b>
	F.1 Proof of Proposition 5.1 . . . . .	32
	F.2 Proof of Proposition 5.2 . . . . .	33
	F.3 Proof of Proposition 5.3 . . . . .	34
<b>G</b>	<b>Proof of Section 6</b>	<b>36</b>
	G.1 Proof of Theorem 6.1 . . . . .	36
<b>H</b>	<b>Supporting Lemmas</b>	<b>37</b>
<b>I</b>	<b>The Use of Large Language Models</b>	<b>38</b>

## A NOTATION TABLE

	<b>Functions</b>
756	
757	
758	<b>Functions</b>
759	$f : \mathbb{A} \rightarrow \mathbb{B}$ The function $f$ with domain $\mathbb{A}$ and range $\mathbb{B}$
760	$f \circ g$ Composition of the functions $f$ and $g$
761	
762	$f(\mathbf{x}; \boldsymbol{\theta})$ A function of $\mathbf{x}$ parametrized by $\boldsymbol{\theta}$ . (Sometimes we write
763	$f(\mathbf{x})$ and omit the argument $\boldsymbol{\theta}$ to lighten notation)
764	$\sigma_R(x)$ ReLU function, $\max\{x, 0\}$
765	$\sigma_L(x)$ SiLU function, $\frac{x}{1 + e^{-x}}$
766	
767	$\sigma_S(\mathbf{x})$ Softmax function, $\sigma_S(\mathbf{x})_i = \frac{\exp(\mathbf{x}_i)}{\sum_{i=1}^d \exp(\mathbf{x}_i)}$
768	
769	$\ \mathbf{x}\ _p$ $L^p$ norm of $\mathbf{x}$
770	
771	$\ \mathbf{x}\ $ $L^2$ norm of $\mathbf{x}$
772	
773	$\ \mathbf{x}\ _\infty$ $\infty$ norm of $\mathbf{x}$
774	$x^+$ Positive part of $x$ , i.e., $\max(0, x)$
775	$\mathcal{F}_{SA}$ Standard self-attention layer
776	$\mathcal{F}_{FF}$ Feed-forward layer
777	
778	$\mathcal{F}_{SA}^A$ Self-attention layer with fixed attention pattern $A$
779	$\mathcal{T}(D, H, S, W, L)$ Transformer neural network class with embedding size $D$ ,
780	number of heads $H$ , head size $S$ , hidden dimension $W$ ,
781	number of layers $L$
782	$\mathcal{T}(D, H, S, L)$ Attention-only Transformer neural network class with
783	embedding size $D$ , number of heads $H$ , head size $S$ , number
784	of layers $L$
785	$\mathcal{NN}_\sigma(W, L, \mathbb{R}^d \rightarrow \mathbb{R}^{d'})$ Feed-forward neural network class with width $W$ , depth
786	$L$ , input dimension $d$ and output dimension $d'$ , activation
787	function $\sigma$
788	
789	$\mathcal{NN}_\sigma^{Res}(W, L, \mathbb{R}^d \rightarrow \mathbb{R}^d)$ Residual feed-forward neural network class with hidden
790	dimension $W$ , number of layers $L$ , input dimension and
791	output dimension $d$ , activation function $\sigma$
792	<b>Numbers and Arrays</b>
793	
794	$a$ A scalar (integer or real)
795	$\mathbf{a}$ A vector
796	$\mathbf{A}$ A matrix
797	
798	$I_n$ Identity matrix with $n$ rows and $n$ columns
799	$I$ Identity matrix with dimensionality implied by context
800	
801	$\mathbf{e}_i$ Standard basis vector $[0, \dots, 0, 1, 0, \dots, 0]$ with a 1 at
802	position $i$
803	$\mathbf{1}_{n \times m}$ All-one matrix with dimensionality $n \times m$
804	
805	
806	
807	
808	
809	

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863**Sets**

$\mathbb{R}$	The set of real numbers
$\mathbb{R}^D$	The set of $D$ -dimensional real vectors
$\mathbb{R}_{>0}^D$	The set of $D$ -dimensional positive real vectors
$\{0, 1\}$	The set containing 0 and 1
$\{0, 1, \dots, n\}$	The set of all integers between 0 and $n$
$[n]$	The set of all integers between 1 and $n$ , that is, $[n] = \{1, \dots, n\}$

**Indexing**

$\mathbf{a}_i$	Element $i$ of vector $\mathbf{a}$ , with indexing starting at 1
$\mathbf{a}_{-1}$	The last element of vector $\mathbf{a}$
$\mathbf{A}_{i,j}$	Element $i, j$ of matrix $\mathbf{A}$
$\mathbf{A}_{i,:}$	Row $i$ of matrix $\mathbf{A}$
$\mathbf{A}_{:,i}$	Column $i$ of matrix $\mathbf{A}$

**Asymptotics**

$f(n) = O(g(n))$	$f$ grows at most as fast as $g$ for sufficiently large $n$
$f(n) = \tilde{O}(g(n))$	$f$ grows at most as fast as $g$ for sufficiently large $n$ , up to logarithmic factors
$f(n) = \Omega(g(n))$	$f$ grows at least as fast as $g$ for sufficiently large $n$
$f \lesssim g$	There exists a positive constant $c$ such that $f \leq cg$ holds

**B ADDITIONAL RELATED WORK**

**Theoretical Understanding of Transformers.** One the fundamental study on the expressive capacity of Transformer architecture is to explore its approximation ability, that is, investigating whether Transformers can approximate functions that belong to a given function class. The most seminal work by (Yun et al., 2019) provided the first universal approximation theorem for Transformers, showing that any continuous sequence-to-sequence functions defined on a compact domain can be approximated by Transformer to any precision. They also extended the results to sparse Transformers in (Yun et al., 2020). Gurevych et al. (2022) gave a constructive method to show that Transformers can approximate piecewise polynomials. Jiang & Li (2024) built their results of approximating continuous functions by shallow Transformers based on the Kolmogorov Representation Theorem. Takakura & Suzuki (2023) provided both approximation and estimation error with  $\gamma$ -smooth function class under the assumption that the input is infinite dimensional. Similarly, Havrilla & Liao (2024) leveraged manifold hypothesis, assuming that the input data has a low-dimensional structure and established approximation results for  $\beta$ -Hölder continuous functions. Kajitsuka & Sato (2023) showed that Transformers with one single-head self-attention layer can be a universal approximator, where they dug deeply into the relationship between the softmax function in the self-attention layer and the Boltzmann operator. Takeshita & Imaizumi (2025) proved that Transformers can efficiently approximate column-symmetric polynomials with respect to the number of parameters. Concurrently, Jiao et al. (2025a) established the approximation results of Transformers for Hölder class and Sobolev class under  $L^p$ -norm, where  $p \in [0, +\infty]$ . Besides, their another work (Jiao et al., 2025b) proved that Transformers are able to overcome the curse of dimensionality based on Kolmogorov-Arnold Representation Theorem. Hu et al. (2025) tried to avoid the dependence on large ReLU feed-forward layers by proving that attention layers alone can approximate a generalized version of ReLU function and hence subsumes any known approximators based on ReLU feed-forward neural networks. Similarly, Liu et al. (2025) proved that a single self-attention layer, preceded by sum-

of-linear transformations, is capable of approximating any continuous functions on a compact domain under  $L^\infty$ -norm, highlighting the inherent expressive power of attention mechanism alone. Cheng et al. (2025) investigated the universal approximation property of Transformer-type architectures, providing a unified theoretical framework that incorporates various architecture variants.

As for Transformers with prompts, Wang et al. (2023) showed that prompt tuning Transformers can be a universal approximator to Lipschitz sequence-to-sequence functions and Hu et al. (2024) further extended their results to Transformers with only one self-attention layer. Petrov et al. (2024) proved that prompt tuning Transformers is able to approximate sequence-to-sequence functions defined on the hypersphere. Besides, in Nakada et al. (2025), they proved that a fixed-size Transformer with well-designed prompts can exactly compute a certain class of ReLU feed-forward neural networks.

Another line of theoretical study on the expressive ability of Transformers is to see whether Transformers can achieve zero loss on finite input-output pairs. Kim et al. (2022) is the first work to study the memorization capacity of Transformers. They proved that Transformers with  $\tilde{O}(d + n + \sqrt{nN})$  parameters are able to memorize length  $n$  and  $d$ -dimensional input tokens. Mahdavi et al. (2023) showed that a multi-head self-attention mechanism with  $H$  heads and  $O(Hd^2)$  parameters is capable of memorizing  $O(Hn)$  data samples. Kajitsuka & Sato (2023) proved that a Transformer with only one single-head self-attention layer have data memorization ability. Chen & Zou (2024) built the results of Transformers with ReLU activation function under the assumption that each data label is distinct. Kajitsuka & Sato (2024) established their results with an optimal number of parameters. It was shown that Transformers with  $\tilde{O}(\sqrt{N})$  parameters in the sequence-to-sequence prediction setting and  $\tilde{O}(\sqrt{nN})$  parameters in the sequence-to-sequence prediction task.

## C A DISCUSSION ON ASSUMPTION 3.1

In this section we provide a discussion on Assumption 3.1, which shows that it is mild and can be easily satisfied in real-world settings.

**Assumption C.1** (Restatement of Assumption 3.1). Let  $\mathcal{V} \subset \mathbb{R}^D$  be a finite vocabulary and [MASK] denote the mask token with [MASK]  $\notin \mathcal{V}$ . Each token  $\mathbf{x} \in \mathcal{V}$  corresponds to a token ID  $y \in [C]$ , where  $C = |\mathcal{V}|$ . Let  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)}) \subset \mathbb{R}^{D \times N} \times [C]$  be  $n$  input-output pairs with  $y^{(i)}$  denoting the token ID of the masked token. For any  $m_i \in \{1, 2, \dots, N\}$ , the  $m_i$ -th token in  $\mathbf{X}^{(i)}$  is masked by the [MASK] token, that is,  $\mathbf{X}_{:,m_i}^{(i)} = [\text{MASK}]$ . We assume that the following conditions are satisfied:

1. There exists  $r_1, r_2 > 0$  such that for any  $i \in [n]$  and  $j \in [N]$ ,  $r_1 \leq \|\mathbf{X}_{:,j}^{(i)}\| \leq r_2$ .
2. There exists  $\delta > 0$  such that for any  $i, j \in [n]$  and  $k, l \in [N]$ , either  $\mathbf{X}_{:,k}^{(i)} = \mathbf{X}_{:,l}^{(j)}$  or  $\|\mathbf{X}_{:,k}^{(i)} - \mathbf{X}_{:,l}^{(j)}\| \geq \delta$  holds.
3. For any  $i, j \in [n]$ ,  $\mathbf{X}^{(i)} \neq \mathbf{X}^{(j)}$  up to permutations.
4. There are no duplicated tokens in each  $\mathbf{X}^{(i)}$  for  $i \in [n]$ .

In this following, we provide a step-by-step interpretation of Assumption 3.1, showing that our assumption on data is reasonable and easily satisfied in real-world scenarios:

1. This  $r$  naturally exists in real-world datasets, since every data point is stored with finite precision. Besides, due to the widely utilized normalization technique during training (i.e., layer normalization, batch normalization) or during preprocessing, every data point is well-bounded.
2. Since we consider a discrete dataset, which contains finite data points, so this  $\delta$  inherently exists. This assumption does not impose any extra limitation on the data, but only provides convenience for theoretical analysis.
3. Since there exists one token in each  $\mathbf{X}^{(i)}$  being masked by [MASK], it is possible that there exists two data points  $\mathbf{X}^{(i)}$  and  $\mathbf{X}^{(j)}$ , which are not equal to each other before being masked but equal to each other up to permutation after it. Moreover, It is widely known that

Transformer architecture is permutation equivariant without special positional encoding, that is, given any permutation matrix  $\mathbf{P}$ ,  $f(\mathbf{X}\mathbf{P}) = f(\mathbf{X})\mathbf{P}$  holds for any Transformer  $f$ . As a result,  $f(\mathbf{X}^{(i)})$  is also a permutation of  $f(\mathbf{X}^{(j)})$ , meaning that  $f(\mathbf{X}^{(i)})_{:,m_i} \equiv f(\mathbf{X}^{(j)})_{:,m_j}$ , which contradicts with fact that  $y^{(i)}$  is possible to be different from  $y^{(j)}$ . In natural languages, permutation equivalence of words is not common. Different permutation usually leads to different meanings, which motivates us to assume that no permuted data points are in our setting.

4. In practical scenarios, it is possible that two exactly same tokens appear in one sequence. However, the permutation equivariant limitation of Transformers makes the two tokens undistinguishable. Many works (Su et al., 2024; Ke et al., 2020) have been working on designing effective positional encoding to break this limit. In this work, since we do not focus on positional encoding, it is reasonable that we consider there are no duplicated tokens in each  $\mathbf{X}^{(i)}$ .

## D EXPERIMENTS

This section provides experimental results to back up our theory. We validate that (i) Self-attention augmented with a fixed vector can approximate a FFN effectively (Theorem 4.1), (ii) a two layer attention-only Transformer and its variants can perform well in masked language modeling (Theorem 6.1, Proposition 5.2, Proposition 5.3). We conduct all experiments using one NVIDIA A100 GPU. Our code is based on standard PyTorch modules.

### D.1 PROOF-OF-CONCEPT EXPERIMENTS ON THEOREM 4.1

**Objective: Verifying self-attention augmented with a fixed vector approximates a FFN.** We investigate accuracy of self-attention with a fixed vector appended in the input approximating a FFN, comparing with standard self-attention.

**Data Generation and Model Architecture.** We randomly generate  $\mathbf{X} \in \mathbb{R}^{D \times N}$  drawn from a normal distribution, where  $\mathbf{X} \sim N(0, 1)$ . The number of training data is 75 and test data is 25. We initialize the self-attention layer with  $D = 16$ ,  $N = 8$ ,  $S = 1$ ,  $H = 75$ ,  $L = 3$  where  $D$  is the embedding size,  $N$  is the sequence length,  $S$  is the head size,  $H$  is the number of heads and  $L$  is the number of self-attention layers. Furthermore, we also randomly generate a FFN  $\mathbf{f} \in \mathcal{NN}_{\sigma_R}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^D)$  as the target, with  $W = H = 75$ ,  $L = 3$ . The batch size is chosen to be 1 and optimizer is Adam. Following the original Theorem 4.1, we exclude dropout and layer normalization in our experiments following our theory.

**Results.** We let  $\mathbf{S}$  denote the standard self-attention and  $\mathbf{F}$  represents the self-attention with a fixed vector. As shown in Figure 4, evaluated on  $\log$  test MSE error, we can see that  $\mathbf{F}$  achieves smaller error than  $\mathbf{S}$ , which proves our theory.

**Sensitivity of Approximation to the Number of Heads.** We study the relationship between the approximation error and the number of heads. As shown in Figure 5, increasing the number of heads, the approximation error of  $\mathbf{F}$  decreases stably and reach equilibrium, while the error curve of  $\mathbf{S}$  fluctuates as the number of heads growing.

**Sensitivity of Approximation to the Number of Layers.** We consider to study the approximation ability between the number of self-attention layers. In Theorem 4.1, the required number of self-attention layers equal to the depth of the FFN to be approximated. It is natural to ask that can fewer layers achieve the similar results or can more layers achieve better performance? Fixing a FFN with depth 5, we use a stacked self-attention with layer  $\{1, 2, \dots, 10\}$  to approximate it. As shown in Figure 6, the minimum of  $\mathbf{F}$  is achieved when the number of layers is 5, which proves our theory.

### D.2 PROOF-OF-CONCEPT EXPERIMENTS ON THEOREM 6.1

**Objective: Verifying the Effectiveness in Masked Language Modeling.** We investigate whether four different models that are based on Theorem 6.1, Proposition 5.2 and 5.3, can handle the masked language modeling task.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

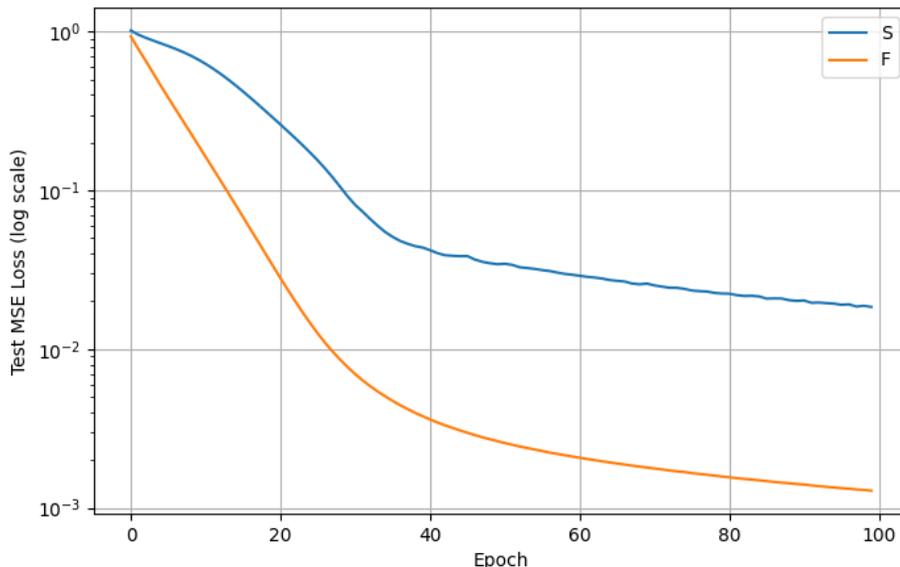


Figure 4: **Test Loss.** We report the test error of both **F** and **S** versus the number of training epochs. We use synthetic data of 75 training data points and 25 testing data points, with sequence length being 8 and embedding dimension 16. We set the batch size to be 16. The optimizer used is Adam with learning rate 0.001. The result shows that after training for 100 epochs, **F** can effectively approximate the target FFN while **S** can not.

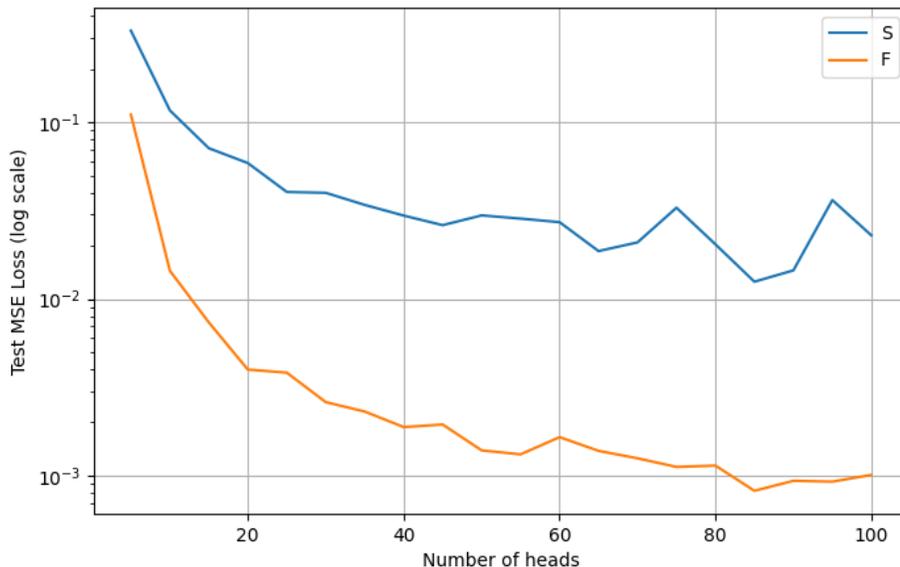


Figure 5: **Test Loss of different number of heads.** We report the test error of both **F** and **S** with number of heads being chosen from  $\{5, 10, \dots, 100\}$ . We use synthetic data of 75 training data points and 25 testing data points, with sequence length being 8 and embedding dimension 16. We set the batch size to be 16. The optimizer used is Adam with learning rate 0.001. The result shows that after training for 100 epochs, increasing the number of heads, **F** can effectively approximate the target FFN, while the error of **S** remains far above zero.

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079

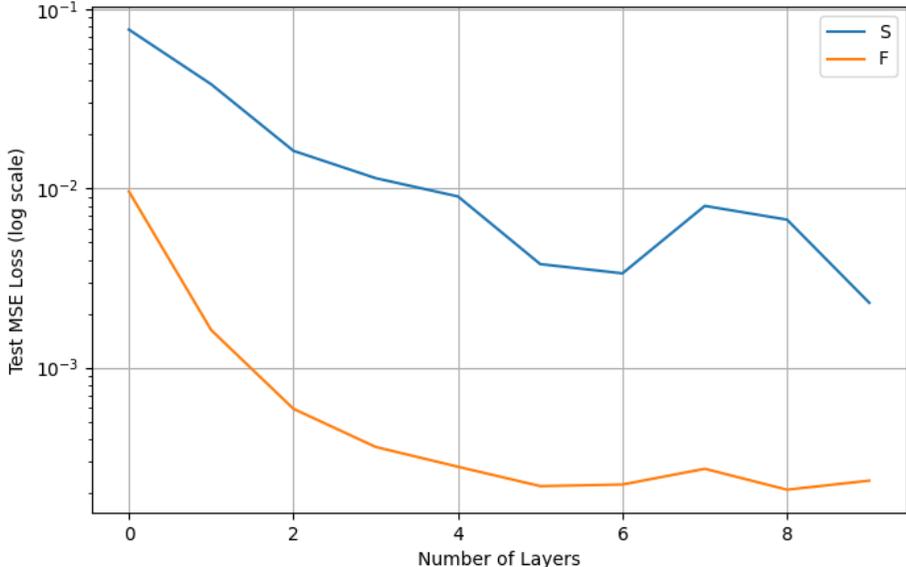


Figure 6: **Test Loss of different number of layers.** We report the test error of both **F** and **S** approximating a FFN with depth 5 versus the number of self-attention layers  $\{1, 2, \dots, 10\}$ . We use synthetic data of 75 training data points and 25 testing data points, with sequence length being 8 and embedding dimension 16. We set the batch size to be 16. The optimizer used is Adam with learning rate 0.001. The result shows that after training for 100 epochs, the minimum error of **F** is achieved at the 5 layer.

**Data Generation and Model Architecture.** We let the vocabulary size  $C$  be 2048, number of training data points is 1024, sequence length 8, embedding size 16. we ensure that there are no repeated tokens in a single sequence even without positional encoding. We consider four models:

1. **S**: A standard Transformer, consisting of one self-attention layer and one feed-forward layer (Baseline).
2. **F**: **S** + replace the feed-forward layer by a self-attention layer augmented by a fixed vector (Theorem 6.1).
3. **F+R**: **F** + replace the self-attention layer with a randomly chosen, non-trainable attention pattern + add an extra layer of self-attention augmented by a fixed vector (Proposition 5.3).
4. **F+T**: **F** + replace the self-attention layer with a trainable and input-independent attention pattern (Proposition 5.2).

Note that We let the number of heads equals to the number of training data points just as the setting in Theorem 6.1, the head size is chosen to be 1. Similarly, we design a mask that only allows the interaction between data tokens and the fixed vector in the second self-attention layer. The batch is 128 and the optimizer is Adam. We only focus on the training loss curve.

**Results.** As shown in Figure 7, we can see that these four models are able to predict accurately in masked language modeling task, by achieving nearly zero loss. Interestingly, **F+T** and **F+R** own a faster convergence rate than that of **S** and **F**, meaning that a randomly chosen or a trainable, universal attention pattern is already powerful.

**Ablation on the extra self-attention layer in F+R.** Note that in Proposition 5.3, we show that a random attention pattern with a feed-forward layer can do well in token distinction task. In the model **F+R**, we use a self-attention with a fixed vector to replace the feed-forward layer. Now, we consider getting rid of this extra self-attention layer, and see how the performance will change. We denote the

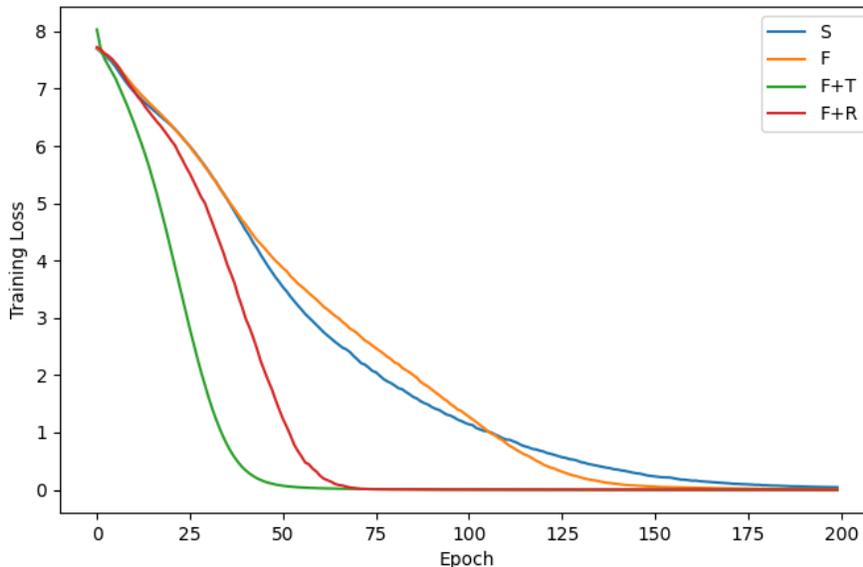


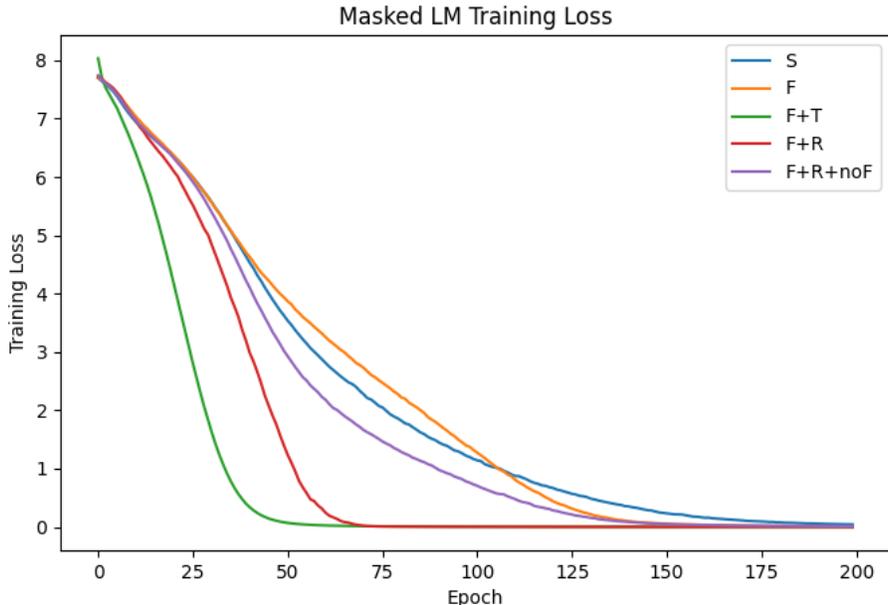
Figure 7: **Training Loss.** We report the training error of four models: **S**, **F**, **F+T**, **F+R** in masked language modeling. We use synthetic data of 1024 training data points, with sequence length being 8 and embedding dimension 16. We set the batch size to be 128. The optimizer used is Adam with learning rate 0.001. The result shows that after training for 200 epochs, all four models achieve nearly zero loss while **F+T** and **F+R** converge much faster.

model by **F+R+noF**. As shown in Figure 8, **F+R+noF** can also reach nearly zero loss only with a slower convergence rate. We leave its theoretical foundation to future reasearch.

**The Linear dependence of Number of heads on the number of training data points.** In Theorem 6.1, we prove that the required number of heads grows linearly as the number of training data points increases. We set the number of the training data points to be  $\{512, 1024, \dots, 4608\}$ , and the number of heads to be  $\{32, 64, \dots, 608\}$ , and see what is the minimum number of heads needed to achieve a training loss less than 0.05 on these training data points. As shown in Figure 9, there is a clear linear relationship between the number of training data points and the number of heads needed to achieve a small loss, which validates our theory.

**Small norm of value vectors.** Note that in the proof of Theorem 4.1, we need the value vector (i.e.,  $\mathbf{W}_V \mathbf{x}$ ) of the bias vector appended to the input equals  $\mathbf{0}$ . We observe that this phenomenon happens in real-world pre-trained language models. Although some attention heads intend to assign large attention weights to special tokens, such as [CLS] or [SEP], the value vector of them are relatively small. We consider tracing the norm of the vector value through layers and results are shown in Figure 10.

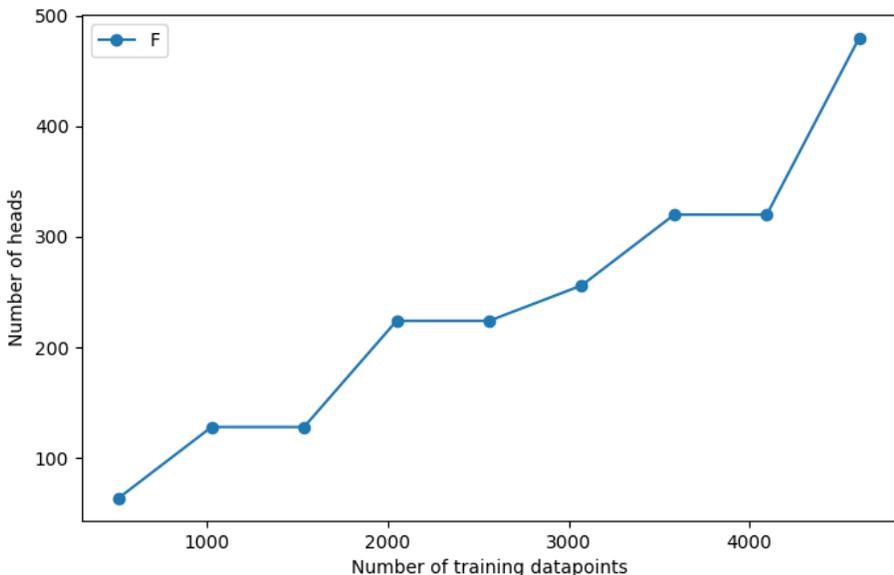
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155



1156  
1157  
1158  
1159  
1160

Figure 8: **Training Loss.** We report the training error of five models: **S**, **F**, **F+T**, **F+R** and **F+R+noF** in masked language modeling. We use synthetic data of 1024 training data points, with sequence length being 8 and embedding dimension 16. We set the batch size to be 128. The optimizer used is Adam with learning rate 0.001. The result shows that after training for 200 epochs, **F+R+noF** has a slower convergence rate than that of **F+R** but still can achieve nearly zero training loss.

1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182



1183  
1184  
1185  
1186  
1187

Figure 9: **Number of heads needed to achieve zero loss under different number of training data points.** We report the number of heads required for producing accurate prediction in masked language modeling. We consider model **F**. We use synthetic data of  $\{512, 1024, \dots, 4608\}$  training data points, with sequence length being 8 and embedding dimension 16. We set the batch size to be 128. The optimizer used is Adam with learning rate 0.001. The result shows that after training for 200 epochs, there is an approximate linear trend in the number of heads.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

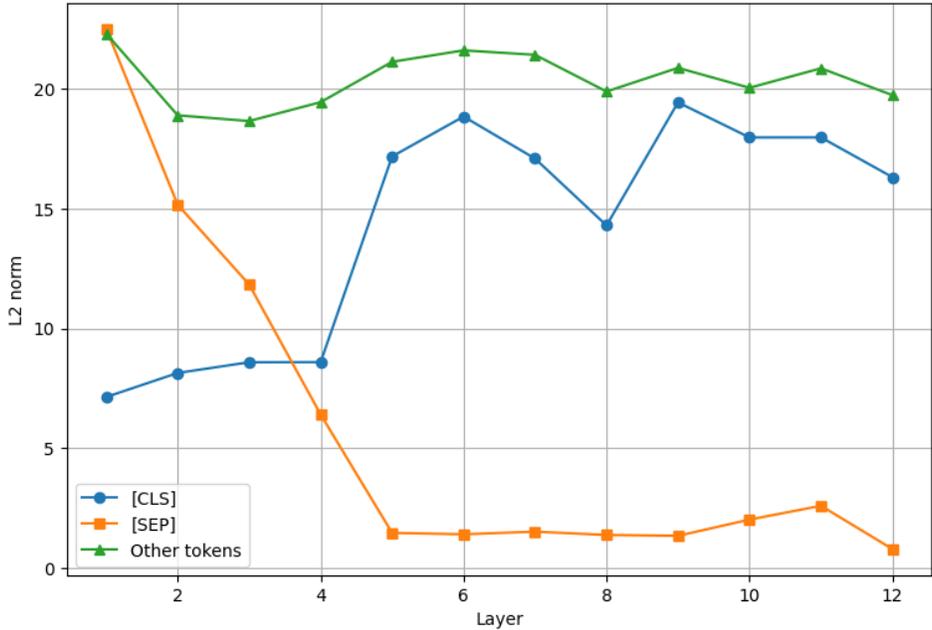


Figure 10: **Norm of the value vectors of different tokens.** We report the norm of the value vector of different tokens across layers. We consider **Bert** (12 layers and 12 heads in each layer). The norm is computed averagely over 100 randomly chosen data points in SST-2 dataset. The result shows the norm of value vector of [CLS] is small before the fourth layer and the norm of value vector of [SEP] becomes nearly zero after layer 5.

## E PROOF OF SECTION 4

### E.1 PROOF OF THEOREM 4.1

**Theorem E.1** (Restatement of Theorem 4.1). *For any  $\varepsilon > 0$ ,  $M > 0$  and any FFN  $f \in \mathcal{NN}_{\sigma_R}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'})$ , there exists an attention-only Transformer  $\mathcal{F} \in \mathcal{T}(\max\{W, D\} + 2, \max\{W, D\} + 1, 1, L)$  such that*

$$\|\mathcal{F}(\mathbf{X}) - f(\mathbf{X})\|_{\infty} < \varepsilon \quad \text{for any } \mathbf{X} \in [-M, M]^{D \times N}.$$

The embedding layer in  $\mathcal{F}$  is defined by

$$\mathcal{E}_{in}(\mathbf{X}) := \begin{pmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{1}_{1 \times N} & \mathbf{0} \\ \mathbf{0}_{(W-D)+ \times N} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(\max\{W, D\}+2) \times (N+1)}.$$

*Proof of Theorem 4.1. Step 1:* In this step, We show that each  $\sigma_R$  or  $\sigma_L$  activated FFN with bias term in each layer can be transformed into a no-bias FFN by slightly increasing the width, which is summarized in the following Lemma.

**Lemma E.1.** *Let  $\sigma = \sigma_R$  or  $\sigma_L$ . For any  $f \in \mathcal{NN}_{\sigma}(W, L, \mathbb{R}^d \rightarrow \mathbb{R}^{d'})$ , there exists  $g \in \mathcal{NN}_{\sigma}(W + 1, L, \mathbb{R}^d \rightarrow \mathbb{R}^{d'})$  with no bias terms in each layer such that*

$$f(\mathbf{X}) = g(\widehat{\mathbf{X}}),$$

where  $\widehat{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \mathbf{1}_{1 \times N} \end{pmatrix}$ .

1242 The proof of Lemma E.1 is postponed to Appendix E.2.  
1243

1244 **Step 2:** We build the connection between FFNs activated by  $\sigma_L$  and FFNs activated by  $\sigma_R$ . The  
1245 following Lemma demonstrates that for any  $\sigma_R$ -activated FFN, there exists a  $\sigma_L$ -activated FFN can  
1246 approximate it to any precision.

1247 **Lemma E.2.** *for any  $\varepsilon > 0$ ,  $M > 0$ , and FFN without bias  $\mathbf{f} \in \mathcal{NN}_{\sigma_R}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'})$ , there  
1248 exists a FFN without bias  $\mathbf{g} \in \mathcal{NN}_{\sigma_L}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'})$  such that*

$$1249 \quad \|\mathbf{f}(\mathbf{X}) - \mathbf{g}(\mathbf{X})\|_{\sup([-M, M]^{D \times N})} < \varepsilon.$$

1251 The proof of Lemma E.2 is placed in Appendix E.3.  
1252

1253 **Step 3:** In this step, we aim to build the connection between residual FFNs and non-residual FFNs.  
1254 Specifically, the Lemma below shows that each non-residual FFN can be represented by a residual  
1255 FFN when the activation function is  $\sigma_L$ .

1256 **Lemma E.3.** *For any FFN without bias terms  $\mathbf{f} \in \mathcal{NN}_{\sigma_L}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'})$  with  $D' \leq W$ , there  
1257 exists a residual FFN  $\mathbf{g} \in \mathcal{NN}_{\sigma_L}^{Res}(\max\{W, D\}, L, \mathbb{R}^{\max\{W, D\}} \rightarrow \mathbb{R}^{\max\{W, D\}})$  such that for any  
1258  $\mathbf{X} \in \mathbb{R}^{D \times N}$ , the following holds*

$$1260 \quad \mathbf{g}(\widehat{\mathbf{X}}) = \begin{pmatrix} \mathbf{f}(\mathbf{X}) \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{\max\{W, D\} \times N},$$

$$1263 \quad \text{where } \widehat{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \mathbf{0}_{(W-D) \times N} \end{pmatrix} \in \mathbb{R}^{\max\{W, D\} \times N}.$$

1267 The proof of Lemma E.3 is postponed to Appendix E.4.  
1268

1269 **Step 4:** In this part, we build a close connection between self-attention and FFNs activated by  $\sigma_L$ .  
1270 The following Lemma shows that we can use a single self-attention layer to implement one layer  
1271 FFN activated by  $\sigma_L$ .

1272 **Lemma E.4.** *for any  $\mathbf{W}_1 \in \mathbb{R}^{W \times D}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{D \times W}$  and  $\mathbf{X} \in \mathbb{R}^{D \times N}$ , there exists a self-attention  
1273 layer  $\mathcal{F}_{SA} \in \mathcal{T}(D+1, W, 1, 1)$  such that*

$$1274 \quad \mathcal{F}_{SA}(\widehat{\mathbf{X}}) = \begin{pmatrix} \mathbf{W}_2 \sigma_L(\mathbf{W}_1 \mathbf{X}) + \mathbf{X} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)},$$

$$1277 \quad \text{where } \widehat{\mathbf{X}} = \begin{pmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)}.$$

1281 The proof of Lemma E.4 is places in Appendix E.5.  
1282

1283 **Step 5:** In this part, we show that for any residual feed-forward neural network without bias terms  
1284  $\mathbf{f} \in \mathcal{NN}_{\sigma_L}^{Res}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^D)$ , there exists an attention-only Transformer  $\mathcal{F} \in \mathcal{T}(D+1, W, 1, L)$   
1285 such that

$$1286 \quad \mathcal{F}(\widehat{\mathbf{X}}) = \begin{pmatrix} \mathbf{f}(\mathbf{X}) & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)},$$

$$1289 \quad \text{where } \widehat{\mathbf{X}} = \begin{pmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)}.$$

1292 According to the definition of residual neural networks,  $\mathbf{f}$  can be written as  
1293

$$1294 \quad \mathbf{f} = \mathcal{L}_L \circ \dots \circ \mathcal{L}_1,$$

1295 where  $\mathcal{L}_\ell(\mathbf{X}) = \mathbf{X} + \mathbf{W}_\ell^{(2)} \sigma_L(\mathbf{W}_\ell^{(1)} \mathbf{X})$  with  $\mathbf{W}_\ell^{(2)} \in \mathbb{R}^{D \times W}$ ,  $\mathbf{W}_\ell^{(1)} \in \mathbb{R}^{W \times D}$ .

1296 Through the analysis in **step 1**, there exists a self-attention layer  $\mathcal{F}_{SA}^{(\ell)} \in \mathcal{T}(D+1, W, 1, 1)$  such that

$$1297$$

$$1298 \quad \mathcal{F}_{SA}^{(\ell)}(\widehat{\mathbf{X}}) = \begin{pmatrix} \mathbf{W}_\ell^{(2)} \sigma_L(\mathbf{W}_\ell^{(1)} \mathbf{X}) + \mathbf{X} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)},$$

$$1299$$

$$1300$$

1301 Let  $\mathcal{F} := \mathcal{F}_{SA}^{(L)} \circ \dots \circ \mathcal{F}_{SA}^{(1)}$ . Since the 1 in the bottom right corner is kept in the whole residual flow,

1302 we can verify that

$$1303 \quad \mathcal{F}(\widehat{\mathbf{X}}) = \mathcal{F}_{SA}^{(L)} \circ \dots \circ \mathcal{F}_{SA}^{(1)}(\widehat{\mathbf{X}})$$

$$1304$$

$$1305 \quad = \mathcal{F}_{SA}^{(L)} \circ \dots \circ \mathcal{F}_{SA}^{(2)} \begin{pmatrix} \mathbf{W}_1^{(2)} \sigma_L(\mathbf{W}_1^{(1)} \mathbf{X}) + \mathbf{X} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix}$$

$$1306$$

$$1307$$

$$1308 \quad \vdots$$

$$1309$$

$$1310 \quad = \begin{pmatrix} \mathbf{f}(\mathbf{X}) & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)}.$$

$$1311$$

$$1312$$

1313 The proof of **step 5** is completed by noting that  $\mathcal{F} \in \mathcal{T}(D+1, W, 1, L)$ .

1314 **Step 6:** Putting everything together.

1315 For any  $\mathbf{f} \in \mathcal{NN}_{\sigma_R}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'})$ , there exists a no-bias  $\mathbf{f}_1 \in \mathcal{NN}_{\sigma_R}(W+1, L, \mathbb{R}^{D+1} \rightarrow$

1317  $\mathbb{R}^{D'})$  such that

$$1318 \quad \mathbf{f}_1(\mathbf{X}_1) = \mathbf{f}(\mathbf{X}),$$

1319

1320 where  $\mathbf{X}_1 = \begin{pmatrix} \mathbf{X} \\ \mathbf{1}_{1 \times N} \end{pmatrix} \in \mathbb{R}^{(D+1) \times N}$ . Moreover, according to Lemma E.2, for any  $\varepsilon > 0$ , there

1322 exists  $\mathbf{f}_2 \in \mathcal{NN}_{\sigma_L}(W+1, L, \mathbb{R}^{D+1} \rightarrow \mathbb{R}^{D'})$  such that

$$1323 \quad \|\mathbf{f}_2(\mathbf{X}_1) - \mathbf{f}_1(\mathbf{X}_1)\|_{\sup([-M, M]^{D' \times N})} < \varepsilon.$$

1324 By applying Lemma E.3, there exists a residual FFN

$$1325 \quad \mathbf{f}_3 \in \mathcal{NN}_{\sigma_L}(\max\{W+1, D+1\}, L, \mathbb{R}^{\max\{W+1, D+1\}} \rightarrow \mathbb{R}^{\max\{W+1, D+1\}})$$

1326 such that

$$1327$$

$$1328 \quad \mathbf{f}_3(\mathbf{X}_2) = \begin{pmatrix} \mathbf{f}_2(\mathbf{X}_1) \\ \mathbf{0}_{(W-D)+ \times N} \end{pmatrix},$$

$$1329$$

$$1330$$

1331 where  $\mathbf{X}_2 = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{0}_{(W-D)+ \times N} \end{pmatrix} \in \mathbb{R}^{(\max\{W, D\}+1) \times N}$ . Finally, Lemma 4.1 and **Step 5** show that there

1333 exists an attention-only Transformer  $\mathcal{F}_{SA} \in \mathcal{T}(\max\{W, D\}+2, \max\{W, D\}+1, 1, L)$  such that

$$1334 \quad \mathcal{F}_{SA}(\mathbf{X}_3) = \begin{pmatrix} \mathbf{f}_3(\mathbf{X}_2) & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(\max\{W, D\}+2) \times (N+1)},$$

$$1335$$

$$1336$$

$$1337$$

$$1338$$

1339 where  $\mathbf{X}_3 = \begin{pmatrix} \mathbf{X}_2 & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(\max\{W, D\}+2) \times (N+1)}$ . The proof is completed by defining  $\mathcal{E}_{in}$  and

1342  $\mathcal{E}_{out}$  as

$$1343 \quad \mathcal{E}_{in}(\mathbf{X}) = \begin{pmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{1}_{1 \times N} & \mathbf{0} \\ \mathbf{0}_{(W-D)+ \times N} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(\max\{W, D\}+2) \times (N+1)} \quad \text{for any } \mathbf{X} \in \mathbb{R}^{D \times N},$$

$$1344$$

$$1345$$

$$1346$$

$$1347$$

$$1348$$

$$1349 \quad \mathcal{E}_{out}(\mathbf{X}) = \mathbf{X}_{1:D, 1:N}.$$

□

## E.2 PROOF OF LEMMA E.1

*Proof of Lemma E.1.* We first consider the case when  $\sigma = \sigma_R$ . According to the definition of FFNs,  $\mathbf{f}$  can be written as

$$\mathbf{f} = \mathcal{L}_L \circ \sigma_R \cdots \circ \sigma_R \circ \mathcal{L}_0,$$

where each  $\mathcal{L}_\ell$  is given by  $\mathcal{L}_\ell(\mathbf{x}) := \mathbf{W}_\ell \mathbf{x} + \mathbf{b}_\ell$  for  $\ell = 0, 1, \dots, L$  with  $\mathbf{W}_\ell \in \mathbb{R}^{D_{\ell+1} \times D_\ell}$ ,  $\mathbf{b}_\ell \in \mathbb{R}^{D_{\ell+1}}$ , and  $D_0 = D, D_1, \dots, D_L \in \mathbb{N}^+$ , and  $D_{L+1} = D'$ .

Let  $\widehat{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \mathbf{1}_{1 \times N} \end{pmatrix}$ , where  $\mathbf{1}_{1 \times N}$  is the all-1 vector of size  $1 \times N$ . We define

$$\widehat{\mathbf{W}}_\ell = \begin{pmatrix} \mathbf{W}_\ell & \mathbf{b}_\ell \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D_{\ell+1}+1) \times (D_\ell+1)}, \quad \text{for any } \ell = 0, \dots, L-1,$$

$$\widehat{\mathbf{W}}_L = \begin{pmatrix} \mathbf{W}_L & \mathbf{b}_L \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{D' \times (D_L+1)}.$$

Let  $\mathbf{g} = \widehat{\mathcal{L}}_L \circ \sigma_R \circ \cdots \circ \sigma_R \circ \widehat{\mathcal{L}}_0$ . Through direct verification, we have

$$\mathbf{g}(\widehat{\mathbf{X}}) = \mathbf{f}(\mathbf{X}).$$

Then, we consider when  $\sigma = \sigma_L$ . Let  $s$  be the solution of the equation  $1 = \frac{x}{1+e^{-x}}$  and let

$\widehat{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \mathbf{1}_{1 \times N} \end{pmatrix}$ . Similarly, we define

$$\widehat{\mathbf{W}}_\ell = \begin{pmatrix} \mathbf{W}_\ell & \mathbf{b}_\ell \\ \mathbf{0} & s \end{pmatrix} \in \mathbb{R}^{(D_{\ell+1}+1) \times (D_\ell+1)}, \quad \text{for any } \ell = 0, \dots, L-1,$$

$$\widehat{\mathbf{W}}_L = \begin{pmatrix} \mathbf{W}_L & \mathbf{b}_L \\ \mathbf{0} & s \end{pmatrix} \in \mathbb{R}^{D' \times (D_L+1)}.$$

Let  $\mathbf{g} = \widehat{\mathcal{L}}_L \circ \sigma_R \circ \cdots \circ \sigma_R \circ \widehat{\mathcal{L}}_0$ . Through direct verification, we have

$$\mathbf{g}(\widehat{\mathbf{X}}) = \mathbf{f}(\mathbf{X}),$$

which completes the proof.  $\square$

## E.3 PROOF OF LEMMA E.2

This proof is an extension of Zhang et al. (2024), in which they consider the bias term in each layer.

*Proof of Lemma E.2.* It is straightforward to verify that

$$\frac{\sigma_L(\eta \cdot x)}{x} \rightarrow \sigma_R(x) \quad \text{as } \eta \rightarrow 0^+ \quad \text{for any } x \in [-M, M].$$

Note that  $\frac{\sigma_L(\eta \cdot x)}{x}$  can be implemented by a 1-layer and 1-width FFN activated by  $\sigma_L$ . Assume that  $\mathbf{f}$  can be represented in the following form

$$\mathbf{f} = \mathcal{L}_L \circ \sigma_R \cdots \circ \sigma_R \circ \mathcal{L}_0,$$

where each  $\mathcal{L}_\ell$  is given by  $\mathcal{L}_\ell(\mathbf{x}) := \mathbf{W}_\ell \mathbf{x}$  for  $\ell = 0, 1, \dots, L$  with  $\mathbf{W}_\ell \in \mathbb{R}^{D_{\ell+1} \times D_\ell}$ , and  $D_0 = D, D_1, \dots, D_L \in \mathbb{N}^+$ , and  $D_{L+1} = D', \max\{D_1, \dots, D_L\} \leq W$ . Let  $\sigma_{L,\eta} = \frac{\sigma_L(\eta \cdot x)}{x}$ , we define

$$\phi_\eta(\mathbf{x}) := \mathcal{L}_L \circ \sigma_{L,\eta} \circ \cdots \circ \sigma_{L,\eta} \circ \mathcal{L}_0 \quad \text{for any } \mathbf{x} \in \mathbb{R}^D.$$

It is easy to verify that

$$\phi_\eta \in \mathcal{NN}_{\sigma_L}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'}),$$

and  $\phi_\eta$  does not have bias terms. Later, we prove there exists  $\eta = \eta_0$  such that

$$\|\phi_\eta - \mathbf{f}\|_{\sup([-M, M]^D)} \rightarrow 0 \quad \text{as } \eta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [-M, M]^D.$$

for  $\ell = 1, \dots, L + 1$ , we define

$$\mathbf{h}_\ell(\mathbf{x}) := \mathcal{L}_{\ell-1} \circ \sigma_R \circ \mathcal{L}_{\ell-2} \circ \dots \circ \sigma_R \circ \mathcal{L}_1 \circ \sigma_R \circ \mathcal{L}_0(\mathbf{x}),$$

and

$$\mathbf{h}_{\ell, \eta}(\mathbf{x}) := \mathcal{L}_{\ell-1} \circ \sigma_{L, \eta} \circ \mathcal{L}_{\ell-2} \circ \dots \circ \sigma_{L, \eta} \circ \mathcal{L}_1 \circ \sigma_{L, \eta} \circ \mathcal{L}_0(\mathbf{x}).$$

It is clear that  $\mathbf{h}_\ell$  and  $\mathbf{h}_{\ell, \eta}$  are mappings from  $\mathbb{R}^D$  to  $\mathbb{R}^{D_\ell}$  for  $\ell = 1, \dots, L + 1$ .

For  $\ell = 1, \dots, L + 1$ , we prove by induction that

$$\|\mathbf{h}_{\ell, \eta} - \mathbf{h}_\ell\|_{\sup([-M, M]^D)} \rightarrow 0 \quad \text{as } \eta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [-M, M]^D. \quad (\text{E.1})$$

First, we consider the case  $\ell = 1$ . Clearly,

$$\mathbf{h}_{1, \eta} = \mathcal{L}_0 = \mathbf{h}_1(\mathbf{x}),$$

which means that Equation (E.4) holds for  $\ell = 1$ .

Next, supposing Equation E.4 holds for  $\ell = i \in \{1, \dots, L\}$ , we aim to prove that is also holds for  $\ell = i + 1$ . Determine  $R > 0$  via

$$R = \sup \{ \|\mathbf{h}_j(\mathbf{x})\|_{\ell_\infty} + 1 : \mathbf{x} \in [-M, M]^D, \quad j = 1, 2, \dots, L + 1 \},$$

where the continuity of  $\sigma_R$  guarantees the above supremum is finite. By the induction hypothesis, we have

$$\|\mathbf{h}_{i, \eta} - \mathbf{h}_i\|_{\sup([-M, M]^D)} \rightarrow 0 \quad \text{as } \eta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [-M, M]^D. \quad (\text{E.2})$$

Since for any  $\mathbf{x} \in [-M, M]^D$ , we have  $\|\mathbf{h}_i(\mathbf{x})\|_\infty \leq M$  and

$$\|\mathbf{h}_{i, \eta}(\mathbf{x})\|_\infty \leq \|\mathbf{h}_i(\mathbf{x})\|_\infty + 1 \leq M \quad \text{for small } \eta > 0.$$

Recall that  $\sigma_{L, \eta}(t) \rightarrow \sigma_R(t)$  as  $\eta \rightarrow 0^+$  for any  $t \in [-R, R]$ . Then, we have

$$\|\sigma_{L, \eta} \circ \mathbf{h}_{\ell, \eta}(\mathbf{x}) - \sigma_R \circ \mathbf{h}_{\ell, h}(\mathbf{x})\|_{\sup([-M, M]^D)} \rightarrow 0 \quad \text{as } \eta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [-M, M]^D. \quad (\text{E.3})$$

Due to the continuity of  $\sigma_R$ , we deduce

$$\|\sigma_R \circ \mathbf{h}_{i, \eta}(\mathbf{x}) - \sigma_R \circ \mathbf{h}_i(\mathbf{x})\|_{\sup([-M, M]^D)} \rightarrow 0 \quad \text{as } \eta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [-M, M]^D. \quad (\text{E.4})$$

Therefore, for any  $\mathbf{x} \in [-M, M]^D$ , as  $\eta \rightarrow 0^+$ , we have

$$\begin{aligned} & \sigma_{L, \eta} \circ \mathbf{h}_{i, \eta}(\mathbf{x}) - \sigma_R \circ \mathbf{h}_i(\mathbf{x}) \\ &= \sigma_{L, h} \circ \mathbf{h}_{i, \eta}(\mathbf{x}) - \sigma_R \circ \mathbf{h}_{i, \eta}(\mathbf{x}) + \sigma_R \circ \mathbf{h}_{i, \eta}(\mathbf{x}) - \sigma_R \circ \mathbf{h}_i(\mathbf{x}) \rightarrow 0 \end{aligned}$$

implying

$$\|\mathbf{h}_{i+1, \eta}(\mathbf{x}) - \mathbf{h}_{i+1}(\mathbf{x})\|_{\sup([-M, M]^D)} = \|\mathcal{L}_i \circ \sigma_{L, \eta} \circ \mathbf{h}_{i, \eta} - \mathcal{L}_i \circ \sigma_R \circ \mathbf{h}_i\|_{\sup([-M, M]^D)} \rightarrow 0,$$

which means that Equation E.4 holds for  $\ell = i + 1$ . So we completes the inductive step.

By the principle of induction, as  $\eta \rightarrow 0^+$  and for any  $\mathbf{x} \in [-M, M]^D$  we have

$$\|\phi_\eta - \mathbf{f}\|_{\sup([-M, M]^D)} = \|\mathbf{h}_{L+1, \eta} - \mathbf{h}_{L+1}\|_{\sup([-M, M]^D)} \rightarrow 0.$$

Then, for any  $\varepsilon > 0$ , there exists a small  $\eta_0 > 0$  such that

$$\|\phi_{\eta_0} - \mathbf{f}\|_{\sup([-M, M]^D)} < \varepsilon.$$

Let  $\mathbf{g} = \phi_{\eta_0}$  and the proof is finished by pointing out that

$$\mathbf{g} \in \mathcal{NN}_{\sigma_L}(W, L, \mathbb{R}^D \rightarrow \mathbb{R}^{D'}).$$

□

## E.4 PROOF OF LEMMA E.3

The following proof basically follows (Jiao et al., 2025a).

*Proof of Lemma E.3.* According to the definition of FFN,  $\mathbf{f}$  has the following form

$$\mathbf{f} = \mathcal{L}_L \circ \dots \circ \mathcal{L}_0,$$

where each  $\mathcal{L}_\ell$  is given by  $\mathcal{L}_\ell(\mathbf{x}) := \mathbf{W}_\ell \mathbf{x}$  for  $\ell = 0, 1, \dots, L$  with  $\mathbf{W}_\ell \in \mathbb{R}^{D_{\ell+1} \times D_\ell}$ ,  $\mathbf{b}_\ell \in \mathbb{R}^{D_{\ell+1}}$ , and  $D_0 = D$ ,  $D_1, \dots, D_L \in \mathbb{N}^+$ , and  $D_{L+1} = D'$ . Without loss of generality, we assume that  $D_1, \dots, D_L = W$ , which can be achieved by simply zero-padding these weight matrices.

In the following, We consider two cases: **(1):**  $D \leq W$ , **(2):**  $D > W$ .

**Case 1:**  $D \leq W$ . Define  $\widehat{\mathcal{L}}_1$  as

$$\widehat{\mathcal{L}}_1(\widehat{\mathbf{X}}) = \widehat{\mathbf{X}} + \widehat{\mathbf{W}}_1^{(2)} \sigma_L \left( \widehat{\mathbf{W}}_1^{(1)} \widehat{\mathbf{X}} \right),$$

where

$$\widehat{\mathbf{W}}_1^{(1)} = \begin{pmatrix} \mathbf{W}_0 & \mathbf{0}_{W \times (W-D)} \\ \mathbf{I}_D & \mathbf{0}_{D \times (W-D)} \\ -\mathbf{I}_D & \mathbf{0}_{D \times (W-D)} \end{pmatrix} \in \mathbb{R}^{(W+2D) \times W},$$

$$\widehat{\mathbf{W}}_1^{(2)} = \begin{pmatrix} \mathbf{I}_D & \mathbf{0} & -\mathbf{I}_D & \mathbf{I}_D \\ \mathbf{0} & \mathbf{I}_{W-D} & \mathbf{0} & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{W \times (W+2D)}.$$

Direct computation yields

$$\begin{aligned} \widehat{\mathcal{L}}_1(\widehat{\mathbf{X}}) &= \begin{pmatrix} \mathbf{X} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{I}_D & \mathbf{0} & -\mathbf{I}_D & \mathbf{I}_D \\ \mathbf{0} & \mathbf{I}_{W-D} & \mathbf{0} & \mathbf{0} \end{pmatrix} \sigma_L \left[ \begin{pmatrix} \mathbf{W}_0 & \mathbf{0}_{W \times (W-D)} \\ \mathbf{I}_D & \mathbf{0}_{D \times (W-D)} \\ -\mathbf{I}_D & \mathbf{0}_{D \times (W-D)} \end{pmatrix} \begin{pmatrix} \mathbf{X} \\ \mathbf{0} \end{pmatrix} \right] \\ &= \sigma_L(\mathbf{W}_0 \mathbf{X}) + \begin{pmatrix} \mathbf{X} - \sigma_L(\mathbf{X}) + \sigma_L(\mathbf{X}) \\ \mathbf{0}_{(W-D) \times N} \end{pmatrix} \\ &= \sigma_L[\mathcal{L}_0(\mathbf{X})] \end{aligned}$$

The last equality comes from the fact that

$$\sigma_L(\mathbf{X}) - \sigma_L(-\mathbf{X}) = \mathbf{X}.$$

For  $\ell = 2, \dots, L-1$ , we define  $\widehat{\mathcal{L}}_\ell(\mathbf{Z}) := \mathbf{Z} + \widehat{\mathbf{W}}_\ell^{(2)} \sigma_L \left( \widehat{\mathbf{W}}_\ell^{(1)} \mathbf{Z} \right)$  for any  $\mathbf{Z} \in \mathbb{R}^{W \times N}$ , where

$$\widehat{\mathbf{W}}_\ell^{(1)} = \begin{pmatrix} \mathbf{W}_{\ell-1} \\ \mathbf{I}_W \\ -\mathbf{I}_W \end{pmatrix} \in \mathbb{R}^{3W \times W},$$

$$\widehat{\mathbf{W}}_\ell^{(2)} = \begin{pmatrix} \mathbf{I}_W & -\mathbf{I}_W & \mathbf{I}_W \end{pmatrix} \in \mathbb{R}^{W \times (3W)}.$$

It is direct to verify that

$$\widehat{\mathcal{L}}_{L-1} \circ \dots \circ \widehat{\mathcal{L}}_2 (\sigma_L(\mathcal{L}_0(\mathbf{X}))) = \sigma_L \circ \mathcal{L}_{L-2} \circ \dots \circ \sigma_L \circ \mathcal{L}_0(\mathbf{X}).$$

For last layer, we define

$$\widehat{\mathbf{W}}_L^{(1)} = \begin{pmatrix} \mathbf{W}_{L-1} \\ \mathbf{I}_W \\ -\mathbf{I}_W \end{pmatrix} \in \mathbb{R}^{3W \times W},$$

$$\widehat{\mathbf{W}}_L^{(2)} = \begin{pmatrix} \mathbf{W}_L & -\mathbf{I}_{D'} & \mathbf{0} & \mathbf{I}_{D'} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & -\mathbf{I}_{W-D'} & \mathbf{0} & \mathbf{I}_{W-D'} \end{pmatrix} \in \mathbb{R}^{W \times 3W}.$$

Let  $\mathbf{Z} = \sigma_L \circ \mathcal{L}_{L-2} \circ \dots \circ \sigma_L \circ \mathcal{L}_0(\mathbf{X})$ . Through direct computation, we have

$$\begin{aligned} \widehat{\mathcal{L}}_L(\mathbf{Z}) &= \mathbf{Z} + \begin{pmatrix} \mathbf{W}_L & -\mathbf{I}_{D'} & \mathbf{0} & \mathbf{I}_{D'} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & -\mathbf{I}_{W-D'} & \mathbf{0} & \mathbf{I}_{W-D'} \end{pmatrix} \sigma_L \left( \begin{pmatrix} \mathbf{W}_{L-1} \\ \mathbf{I}_W \\ -\mathbf{I}_W \end{pmatrix} \mathbf{Z} \right) \\ &= \begin{pmatrix} \mathbf{W}_L \sigma_L(\mathbf{W}_{L-1} \mathbf{Z}) \\ \mathbf{0}_{(W-D') \times N} \end{pmatrix} + \mathbf{Z} - \sigma_L(-\mathbf{Z}) + \sigma_L(\mathbf{Z}) \\ &= \begin{pmatrix} \mathbf{f}(\mathbf{X}) \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{W \times N}. \end{aligned}$$

which finishes the proof of **Case 1**.

**Case 2:**  $D > W$ . It is clear that we can zero-pad the weight matrices such that

$$\begin{aligned} \mathbf{W}_0 &\in \mathbb{R}^{D \times D}, \\ \mathbf{W}_\ell &\in \mathbb{R}^{D \times D}, \quad \text{for } \ell = 1, \dots, L-1, \\ \mathbf{W}_L &\in \mathbb{R}^{D' \times D}. \end{aligned}$$

which does not effect the operation of  $\mathbf{f}$ . Then, **Case 2** is reduced to **Case 1**, which completes the proof.  $\square$

## E.5 PROOF OF LEMMA E.4

The following proof is built upon the techniques in (Huben & Morris, 2023).

*Proof of Lemma E.4.* Since  $\mathbf{W}_2 \mathbf{W}_1 \in \mathbb{R}^{D \times D}$ , we have  $\mathbf{W}_2 \mathbf{W}_1 = \sum_{i=1}^{D'} \mathbf{a}_i \mathbf{b}_i^\top$ , where  $\mathbf{a}_i \in \mathbb{R}^{D \times 1}$  is the  $i$ -th column of  $\mathbf{W}_2$  and  $\mathbf{b}_i^\top \in \mathbb{R}^{1 \times D}$  is the  $i$ -th row of  $\mathbf{W}_1$ . We define

$$\begin{aligned} \mathbf{W}_K^{(i)} &= [0, \dots, 0, -1] \in \mathbb{R}^{1 \times (D+1)}, \\ \mathbf{W}_Q^{(i)} &= [\mathbf{b}_i^\top, 0, \dots, 0] \in \mathbb{R}^{1 \times (D+1)}, \\ \mathbf{W}_V^{(i)} &= [\mathbf{b}_i^\top, 0, \dots, 0] \in \mathbb{R}^{1 \times (D+1)}, \\ \mathbf{W}_O^{(i)} &= [\mathbf{a}_i, 0, \dots, 0]^\top \in \mathbb{R}^{(D+1) \times 1}. \end{aligned}$$

Through direct computation, the following holds

$$\left( \mathbf{W}_K^{(i)} \widehat{\mathbf{X}} \right)^\top \left( \mathbf{W}_Q^{(i)} \widehat{\mathbf{X}} \right) = \begin{pmatrix} 0 & \dots & 0 & 0 \\ \vdots & & \vdots & \vdots \\ 0 & \dots & 0 & 0 \\ -\mathbf{b}_i^\top \mathbf{X} & & 0 & 0 \end{pmatrix} \in \mathbb{R}^{(N+1) \times (N+1)}.$$

We define the positional encoding matrix as

$$\mathbf{R}^{(i)} = \begin{pmatrix} 0 & -\infty & \dots & -\infty & -\infty \\ -\infty & \ddots & & \vdots & \vdots \\ \vdots & & \ddots & -\infty & -\infty \\ -\infty & \dots & -\infty & 0 & -\infty \\ 0 & \dots & 0 & 0 & 0 \end{pmatrix} \in \mathbb{R}^{(N+1) \times (N+1)},$$

where in the first  $N + 1$  columns, the interaction between  $i$ -th and  $(N + 1)$ -th column,  $i$ -th and  $i$ -th token is allowed. In the last column, the interaction between  $(N + 1)$ -th token and  $(N + 1)$ -th token is allowed. Thus, we have

$$\begin{aligned} & \sigma_S \left[ \left( \mathbf{W}_K^{(i)} \widehat{\mathbf{X}} \right)^\top \left( \mathbf{W}_Q^{(i)} \widehat{\mathbf{X}} \right) \right] \\ &= \begin{pmatrix} \frac{1}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})} & 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \frac{1}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})} & 0 & \cdots & \cdots & 0 \\ \vdots & \vdots & \ddots & & & \vdots \\ 0 & 0 & & \ddots & \frac{1}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})} & 0 \\ \frac{\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})} & \frac{\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})} & & & \frac{\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})} & 1 \end{pmatrix}. \end{aligned}$$

Taking  $\mathbf{W}_V$  into consideration, we have

$$\begin{aligned} & \mathbf{W}_V^{(i)} \widehat{\mathbf{X}} \sigma_S \left[ \left( \mathbf{W}_K^{(i)} \widehat{\mathbf{X}} \right)^\top \left( \mathbf{W}_Q^{(i)} \widehat{\mathbf{X}} \right) \right] \\ &= \begin{pmatrix} \frac{\mathbf{b}_i^\top \mathbf{X}_{:,1}}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})} & \frac{\mathbf{b}_i^\top \mathbf{X}_{:,2}}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})} & \cdots & \frac{\mathbf{b}_i^\top \mathbf{X}_{:,N}}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})} & 0 \end{pmatrix} \\ &= \left( \sigma_L(\mathbf{b}_i^\top \mathbf{X}_{:,1}), \sigma_L(\mathbf{b}_i^\top \mathbf{X}_{:,2}), \dots, \sigma_L(\mathbf{b}_i^\top \mathbf{X}_{:,N}), 0 \right) \in \mathbb{R}^{1 \times (N+1)}. \end{aligned}$$

Finally,  $\mathbf{W}_O^{(i)}$  recovers the matrix to the original size

$$\begin{aligned} & \mathbf{W}_O^{(i)} \mathbf{W}_V^{(i)} \widehat{\mathbf{X}} \sigma_S \left[ \left( \mathbf{W}_K^{(i)} \widehat{\mathbf{X}} \right)^\top \left( \mathbf{W}_Q^{(i)} \widehat{\mathbf{X}} \right) \right] \\ &= \begin{pmatrix} \mathbf{a}_i \sigma_L(\mathbf{b}_i^\top \mathbf{X}), \dots, \mathbf{a}_i \sigma_L(\mathbf{b}_i^\top \mathbf{X}) & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)}. \end{aligned}$$

Since  $\mathbf{W}_2 \mathbf{W}_1 = \sum_{i=1}^{D'} \mathbf{a}_i \mathbf{b}_i^\top$  and  $\sigma_L$  is element-wise, it is straightforward to have by incorporating the skip connection

$$\begin{aligned} & \widehat{\mathbf{X}} + \sum_{i=1}^{D'} \mathbf{W}_O^{(i)} \mathbf{W}_V^{(i)} \widehat{\mathbf{X}} \sigma_S \left[ \left( \mathbf{W}_K^{(i)} \widehat{\mathbf{X}} \right)^\top \left( \mathbf{W}_Q^{(i)} \widehat{\mathbf{X}} \right) \right] \\ &= \begin{pmatrix} \mathbf{W}_2 \sigma_L(\mathbf{W}_1 \mathbf{X}) + \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)}. \end{aligned}$$

This part is finished by constructing  $\mathcal{F}_{SA}$  with  $\{\mathbf{W}_O^{(i)}, \mathbf{W}_V^{(i)}, \mathbf{W}_K^{(i)}, \mathbf{W}_Q^{(i)}\}_{i \in [D']}$ .  $\square$

## E.6 PROOF OF COROLLARY 4.1

We only need to modify Lemma E.4 in the proof of Theorem 4.1 and keep other steps unchanged.

*Proof of Corollary 4.1.* For any vector  $\mathbf{v} \in \mathbb{R}^{\max\{W,D\}+2}$  with the last element  $v_{-1} \neq 0$ , We redefine the following matrices in the proof of Lemma E.4

$$\begin{aligned} \mathbf{W}_K^{(i)} &= [0, \dots, 0, -1/v_{-1}] \in \mathbb{R}^{1 \times (D+1)}, \\ \mathbf{W}_Q^{(i)} &= [\mathbf{b}_i^\top, 0, \dots, 0] \in \mathbb{R}^{1 \times (D+1)}, \\ \mathbf{W}_V^{(i)} &= [\mathbf{b}_i^\top, 0, \dots, 0] \in \mathbb{R}^{1 \times (D+1)}, \\ \mathbf{W}_O^{(i)} &= [\mathbf{a}_i, 0 \dots, 0]^\top \in \mathbb{R}^{(D+1) \times 1}. \end{aligned}$$

Through direct computation, the following holds

$$\left(\mathbf{W}_K^{(i)} \widehat{\mathbf{X}}\right)^\top \left(\mathbf{W}_Q^{(i)} \widehat{\mathbf{X}}\right) = \begin{pmatrix} 0 & \cdots & 0 & 0 \\ \vdots & & \vdots & \vdots \\ 0 & \cdots & 0 & 0 \\ -\mathbf{b}_i^\top \mathbf{X} & & [-\mathbf{b}_i^\top, 0] \mathbf{v} \end{pmatrix} \in \mathbb{R}^{(N+1) \times (N+1)}.$$

Similarly, we define the mask matrix as

$$\mathbf{R}^{(i)} = \begin{pmatrix} 0 & -\infty & \cdots & -\infty & -\infty \\ -\infty & \ddots & & \vdots & \vdots \\ \vdots & & \ddots & -\infty & -\infty \\ -\infty & \cdots & -\infty & 0 & -\infty \\ 0 & \cdots & 0 & 0 & -\infty \end{pmatrix} \in \mathbb{R}^{(N+1) \times (N+1)},$$

where in the first  $N + 1$  columns, the interaction between  $i$ -th and  $(N + 1)$ -th column,  $i$ -th and  $i$ -th token is allowed. In the last column, the interaction between  $(N + 1)$ -th token and  $(N + 1)$ -th token is allowed. Thus, we have

$$\begin{aligned} & \sigma_S \left[ \left(\mathbf{W}_K^{(i)} \widehat{\mathbf{X}}\right)^\top \left(\mathbf{W}_Q^{(i)} \widehat{\mathbf{X}}\right) \right] \\ &= \begin{pmatrix} \frac{1}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})} & 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \frac{1}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})} & 0 & \cdots & \cdots & 0 \\ \vdots & \vdots & \ddots & & \cdots & \vdots \\ 0 & 0 & & \ddots & \frac{1}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})} & 0 \\ \frac{\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})} & \frac{\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})} & & & \frac{\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})} & 0 \end{pmatrix}. \end{aligned}$$

Taking  $\mathbf{W}_V^{(i)}$  into consideration, we have

$$\begin{aligned} & \mathbf{W}_V^{(i)} \widehat{\mathbf{X}} \sigma_S \left[ \left(\mathbf{W}_K^{(i)} \widehat{\mathbf{X}}\right)^\top \left(\mathbf{W}_Q^{(i)} \widehat{\mathbf{X}}\right) \right] \\ &= \left( \frac{\mathbf{b}_i^\top \mathbf{X}_{:,1}}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,1})} \quad \frac{\mathbf{b}_i^\top \mathbf{X}_{:,2}}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,2})} \quad \cdots \quad \frac{\mathbf{b}_i^\top \mathbf{X}_{:,N}}{1+\exp(-\mathbf{b}_i^\top \mathbf{X}_{:,N})} \quad 0 \right) \\ &= \left( \sigma_L(\mathbf{b}_i^\top \mathbf{X}_{:,1}), \sigma_L(\mathbf{b}_i^\top \mathbf{X}_{:,2}), \cdots, \sigma_L(\mathbf{b}_i^\top \mathbf{X}_{:,N}), 0 \right) \in \mathbb{R}^{1 \times (N+1)}. \end{aligned}$$

Finally,  $\mathbf{W}_O^{(i)}$  recovers the matrix to the original size

$$\begin{aligned} & \mathbf{W}_O^{(i)} \mathbf{W}_V^{(i)} \widehat{\mathbf{X}} \sigma_S \left[ \left(\mathbf{W}_K^{(i)} \widehat{\mathbf{X}}\right)^\top \left(\mathbf{W}_Q^{(i)} \widehat{\mathbf{X}}\right) \right] \\ &= \begin{pmatrix} \mathbf{a}_i \sigma_L(\mathbf{b}_i^\top \mathbf{X}), \cdots, \mathbf{a}_i \sigma_L(\mathbf{b}_i^\top \mathbf{X}) & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)}. \end{aligned}$$

Since  $\mathbf{W}_2 \mathbf{W}_1 = \sum_{i=1}^{D'} \mathbf{a}_i \mathbf{b}_i^\top$  and  $\sigma_L$  is element-wise, it is straightforward to have by incorporating the skip connection

$$\begin{aligned} & \widehat{\mathbf{X}} + \sum_{i=1}^{D'} \mathbf{W}_O^{(i)} \mathbf{W}_V^{(i)} \widehat{\mathbf{X}} \sigma_S \left[ \left(\mathbf{W}_K^{(i)} \widehat{\mathbf{X}}\right)^\top \left(\mathbf{W}_Q^{(i)} \widehat{\mathbf{X}}\right) \right] \\ &= \begin{pmatrix} \mathbf{W}_2 \sigma_L(\mathbf{W}_1 \mathbf{X}) + \mathbf{X} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(D+1) \times (N+1)}. \end{aligned}$$

The proof is finished by constructing  $\mathcal{F}_{SA}$  with  $\{\mathbf{W}_O^{(i)}, \mathbf{W}_V^{(i)}, \mathbf{W}_K^{(i)}, \mathbf{W}_Q^{(i)}\}_{i \in [D']}$ .  $\square$

1674 F PROOF OF SECTION 5

1675 F.1 PROOF OF PROPOSITION 5.1

1676 **Proposition F.1** (Restatement of Proposition 5.1). *For any input-label pairs*  
 1677  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)})$  *satisfying Assumption 3.1, there exists a self-attention layer*  
 1678  $\mathcal{F}_{SA} \in \mathcal{T}(D, 1, 1, 1)$  *such that*

$$1681 \mathcal{F}_{SA}(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}_{SA}(\mathbf{X}^{(j)})_{:,m_j} \text{ for any } i \neq j \in [n].$$

1682 *Proof of Proposition 5.1.* Let  $\mathcal{X} := \{\mathbf{X}_{:,j}^{(i)} \mid \text{for any } i \in [n], j \in [N]\}$ . According to Lemma H.2, for  
 1683 any  $\delta > 0$ , there exists  $\mathbf{W}_K, \mathbf{W}_Q \in \mathbb{R}^{1 \times D}$  such that

$$1684 \left| (\mathbf{W}_K \mathbf{x}_a)^\top (\mathbf{W}_Q \mathbf{x}_c) - (\mathbf{W}_K \mathbf{x}_b)^\top (\mathbf{W}_Q \mathbf{x}_c) \right| > \delta \quad (\text{F.1})$$

1685 for any  $\mathbf{x}_a, \mathbf{x}_b, \mathbf{x}_c \in \mathcal{X}$  with  $\mathbf{x}_a \neq \mathbf{x}_b$ .  $\mathbf{W}_K, \mathbf{W}_Q$  have the following form

$$1686 \mathbf{W}_K = \mathbf{u}\mathbf{v}^\top, \mathbf{W}_Q = \mathbf{u}'\mathbf{v}^\top$$

1687 where  $\mathbf{u}, \mathbf{u}' \in \mathbb{R}$  and  $\mathbf{v} \in \mathbb{R}^{D \times 1}$ .

1688 Define  $\mathbf{a}^{(i)}$  and  $\mathbf{a}^{(j)}$  by

$$1689 \mathbf{a}^{(i)} = \left( \mathbf{W}_K \mathbf{X}^{(i)} \right)^\top \left( \mathbf{W}^{(Q)} \mathbf{X}_{:,m_i}^{(i)} \right),$$

$$1690 \mathbf{a}^{(j)} = \left( \mathbf{W}_K \mathbf{X}^{(j)} \right)^\top \left( \mathbf{W}^{(Q)} \mathbf{X}_{:,m_j}^{(j)} \right).$$

1691 Then, equation F.1 shows that

$$1692 \|\mathbf{a}^{(i)} - \mathbf{a}^{(j)}\| > \delta.$$

1693 Since there are no duplicated tokens in  $\mathbf{X}^{(i)}$  for any  $i \in [N]$ , it follows from Lemma H.3 that

$$1694 \left( \mathbf{a}^{(i)} \right)^\top \sigma_S \left[ \mathbf{a}^{(i)} \right] \neq \left( \mathbf{a}^{(j)} \right)^\top \sigma_S \left[ \mathbf{a}^{(j)} \right]$$

1695 by letting  $\delta > 2 \log D + 3$ . This implies that there exists  $\delta' > 0$  such that

$$1696 \delta' < \left| \left( \mathbf{a}^{(i)} \right)^\top \sigma_S \left[ \mathbf{a}^{(i)} \right] - \left( \mathbf{a}^{(j)} \right)^\top \sigma_S \left[ \mathbf{a}^{(j)} \right] \right|$$

$$1697 = \left| \left( \mathbf{X}_{:,m_i}^{(i)} \right)^\top \left( \mathbf{W}_Q \right)^\top \mathbf{W}_K \left( \mathbf{X}^{(i)} \sigma_S \left[ \mathbf{a}^{(i)} \right] - \mathbf{X}^{(j)} \sigma_S \left[ \mathbf{a}^{(j)} \right] \right) \right|$$

$$1698 = \left| \left( \mathbf{X}_{:,m_i}^{(i)} \right) \mathbf{v} \mathbf{u}' \mathbf{u} \mathbf{v}^\top \left( \mathbf{X}^{(i)} \sigma_S \left[ \mathbf{a}^{(i)} \right] - \mathbf{X}^{(j)} \sigma_S \left[ \mathbf{a}^{(j)} \right] \right) \right|$$

$$1699 = \left| \mathbf{v}^\top \mathbf{X}_{:,m_i}^{(i)} \right| \cdot |\mathbf{u}\mathbf{u}'| \cdot \left| \left( \mathbf{v}^\top \mathbf{X}^{(i)} \sigma_S \left[ \mathbf{a}^{(i)} \right] - \mathbf{v}^\top \mathbf{X}^{(j)} \sigma_S \left[ \mathbf{a}^{(j)} \right] \right) \right|,$$

1700 which means that

$$1701 \mathbf{v}^\top \mathbf{X}^{(i)} \sigma_S \left[ \mathbf{a}^{(i)} \right] \neq \mathbf{v}^\top \mathbf{X}^{(j)} \sigma_S \left[ \mathbf{a}^{(j)} \right].$$

1702 We construct  $\mathcal{F}$  by

$$1703 \mathbf{W}_K = \mathbf{u}\mathbf{v}^\top,$$

$$1704 \mathbf{W}_Q = \mathbf{u}'\mathbf{v}^\top,$$

$$1705 \mathbf{W}_V = \mathbf{v}^\top,$$

$$1706 \mathbf{W}_O = \mathbf{u}''.$$

where  $\mathbf{u}, \mathbf{u}', \mathbf{v}$  are defined in Lemma H.2,  $\mathbf{u}'' \in \mathbb{R}^{D \times 1}$  is an arbitrary nonzero vector satisfying  $\mathbf{u}'' \mathbf{v}^\top \neq \mathbf{0}$ .

$$\begin{aligned} & \left\| \mathcal{F}_{SA}(\mathbf{X}^{(i)})_{:,m_i} - \mathcal{F}_{SA}(\mathbf{X}^{(j)})_{:,m_j} \right\| \\ &= \left\| \mathbf{W}_O \left( \mathbf{W}_V \mathbf{X}^{(i)} \right) \sigma_S \left[ \mathbf{a}^{(i)} \right] - \mathbf{W}_O \left( \mathbf{W}_V \mathbf{X}^{(j)} \right) \sigma_S \left[ \mathbf{a}^{(j)} \right] \right\| \\ &= \|\mathbf{u}''\| \cdot \left| \left( \mathbf{v}^\top \mathbf{X}^{(i)} \sigma_S \left[ \mathbf{a}^{(i)} \right] \right) - \left( \mathbf{v}^\top \mathbf{X}^{(j)} \sigma_S \left[ \mathbf{a}^{(j)} \right] \right) \right| \\ &\neq \mathbf{0}. \end{aligned}$$

The proof is completed by pointing out that

$$\mathcal{F}_{SA} \in \mathcal{T}(D, 1, 1, 1).$$

□

## F.2 PROOF OF PROPOSITION 5.2

**Proposition F.2** (Restatement of Proposition 5.2). *For any input-label pairs  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)})$  satisfying Assumption 3.1, there exists a self-attention layer  $\mathcal{F}_{SA}^A \in \mathcal{T}(D, 1, 1, 1)$  with a fixed attention pattern  $\mathbf{A}$  such that*

$$\mathcal{F}_{SA}^A(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}_{SA}^A(\mathbf{X}^{(j)})_{:,m_j} \quad \text{for any } i \neq j \in [n].$$

*Proof of Proposition 5.2.* Without loss of generality, we assume that  $m_1, \dots, m_n$  are arranged from low to high. For each  $i = 1, \dots, n$  we define

$$\mathcal{P}_i = \{j : j \in \mathbb{Z}_{>0}, j < i, m_j = m_i\}.$$

For any  $\ell = 1, \dots, n$ , we prove by induction that there exists a set of vectors  $\{\mathbf{a}_{m_1}, \dots, \mathbf{a}_{m_\ell}\} \subset \mathbb{R}^N$  such that

$$\mathbf{X}^{(i)} \mathbf{a}_{m_i} \neq \mathbf{X}^{(j)} \mathbf{a}_{m_j} \quad \text{for any } i \neq j \in [\ell]. \quad (\text{F.2})$$

It is clear that when  $\ell = 1$ , Equation.F.2 is naturally satisfied. Next, supposing Equation.F.2 holds for  $\ell = m \geq 2$ , we aim to prove that this is also holds for  $\ell = m + 1$ .

We consider the following set of vectors

$$\mathcal{B}_{m+1} := \bigcup_{i \in \mathcal{P}_{m+1}} \ker \left( \mathbf{X}^{(i)} - \mathbf{X}^{(m+1)} \right).$$

Each  $\ker \left( \mathbf{X}^{(i)} - \mathbf{X}^{(m+1)} \right)$  is defined by

$$\ker \left( \mathbf{X}^{(i)} - \mathbf{X}^{(m+1)} \right) := \left\{ \mathbf{x} \in \mathbb{R}^N : \left( \mathbf{X}^{(i)} - \mathbf{X}^{(m+1)} \right) \mathbf{x} = \mathbf{0} \right\}.$$

According to Assumption 3.1,  $\mathbf{X}^{(i)} \neq \mathbf{X}^{(j)}$  for any  $i \neq j \in [n]$ , meaning that

$$\ker \left( \mathbf{X}^{(i)} - \mathbf{X}^{(m+1)} \right) \subsetneq \mathbb{R}^N \quad \text{for any } i \in \mathcal{P}_{m+1}.$$

Furthermore, for each  $i = 1, \dots, m - |\mathcal{P}_{m+1}|$ , define

$$\mathcal{Q}_i := \left\{ \mathbf{a} \in \mathbb{R}^N : \mathbf{X}^{(i)} \mathbf{a}_{m_i} = \mathbf{X}^{(m+1)} \mathbf{a} \right\}.$$

For any  $i = 1, \dots, m - |\mathcal{P}_{m+1}|$  and any  $\mathbf{a} \in \mathcal{Q}_i$ , we know that

$$\mathbf{X}^{(i)} \mathbf{a}_{m_i} = \mathbf{X}^{(m+1)} \mathbf{a}.$$

Since  $\mathbf{X}^{(m+1)} \neq \mathbf{0}$ , we know that  $\mathcal{Q}_i \subsetneq \mathbb{R}^N$ .

It is obvious that  $\bigcup_{i=1, \dots, m-|\mathcal{P}_{m+1}|} \mathcal{Q}_i \cup \mathcal{B}_{m+1}$  is a finite union of proper subspaces in  $\mathbb{R}^N$ , which has empty interior. Since  $\mathbb{R}_{>0}^N$  has non empty interior, we know that there exists  $\mathbf{a}_{m+1} \in \mathbb{R}^N$  such that

$$\mathbf{a}_{m+1} \in \mathbb{R}_{>0}^N, \\ \mathbf{a}_{m+1} \notin \bigcup_{i=1, \dots, m-|\mathcal{P}_{m+1}|} \mathcal{Q}_i \cup \mathcal{B}_{m+1},$$

meaning that

$$\mathbf{X}^{(i)} \mathbf{a}_{m_i} \neq \mathbf{X}^{(m+1)} \mathbf{a}_{m+1} \quad \text{for any } i \in \mathcal{P}_{m+1}, \\ \mathbf{X}^{(i)} \mathbf{a}_{m_i} \neq \mathbf{X}^{(m+1)} \mathbf{a}_{m+1} \quad \text{for any } i = 1, \dots, m - |\mathcal{P}_{m+1}|.$$

Since we assume that Equation.F.2 holds for  $\ell = m$ , we have

$$\mathbf{X}^{(i)} \mathbf{a}_{:,m_i} \neq \mathbf{X}^{(j)} \mathbf{a}_{:,m_j} \quad \text{for any } i \neq j \in [\ell],$$

which completes the induction step. Up to now, we have proved that we can find proper  $\{\mathbf{a}_{m_i}\}_{i \in [n]}$  such that

$$\mathbf{X}^{(i)} \mathbf{a}_{m_i} \neq \mathbf{X}^{(j)} \mathbf{a}_{m_j} \quad \text{for any } i \neq j \in [n].$$

By applying Lemma H.1 to  $\mathbf{X}^{(1)} \mathbf{a}_{m_1}, \dots, \mathbf{X}^{(n)} \mathbf{a}_{m_n}$ , there exists a vector  $\mathbf{v} \in \mathbb{R}^D$  such that  $\mathbf{v}^\top \mathbf{X}^{(i)} \mathbf{a}_{m_i}$  are pair-wise distinct. Let  $\mathbf{W}_V = \mathbf{v}^\top$  and  $\mathbf{W}_O$  be an arbitrary non-zero vector in  $\mathbb{R}^D$ . We define

$$\mathcal{F}_{SA}^A := \mathbf{X} + \mathbf{W}_O \mathbf{W}_V \mathbf{X} \mathbf{A},$$

where  $\mathbf{A}_{:,m_i} = \mathbf{a}_{m_i}$  and other columns are arbitrarily chosen from  $\mathbb{R}_{>0}^N$ . It is straightforward to verify that

$$|\mathcal{F}_{SA}^A(\mathbf{X}^{(i)})_{:,m_i} - \mathcal{F}_{SA}^A(\mathbf{X}^{(j)})_{:,m_j}| \\ = |\mathbf{X}_{:,m_i} + \mathbf{W}_O \mathbf{W}_V \mathbf{X} \mathbf{A}_{:,m_i} - \mathbf{X}_{:,m_j} - \mathbf{W}_O \mathbf{W}_V \mathbf{X} \mathbf{A}_{:,m_j}| \\ = \mathbf{W}_O \mathbf{v}^\top (\mathbf{X}^{(i)} \mathbf{a}_{m_i} - \mathbf{X}^{(j)} \mathbf{a}_{m_j}) \\ \neq \mathbf{0},$$

which completes the proof.  $\square$

### F.3 PROOF OF PROPOSITION 5.3

**Proposition F.3** (Restatement of Proposition 5.3). *For any input-label pairs  $(\mathbf{X}^{(1)}, y^{(1)}), \dots, (\mathbf{X}^{(n)}, y^{(n)})$  satisfying Assumption 3.1, and any attention pattern  $\mathbf{A} \in \mathbb{R}^{N \times N}$  with  $A_{i,j} > 0$  for any  $i, j \in [N]$ , there exists a Transformer  $\mathcal{F} = \mathcal{E}_{out} \circ \mathcal{F}_{SA}^A \circ \mathcal{F}_{FF} \circ \mathcal{E}_{in} \in \mathcal{T}(\max\{3(n-1)N, D\}, 1, 1, 3(n-1)N, 1)$  such that*

$$\mathcal{F}(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}(\mathbf{X}^{(j)})_{:,m_j} \quad \text{for any } i \neq j \in [n].$$

*Proof of Proposition 5.3.* Let  $S = \{\mathbf{X}_{:,j}^{(i)} \mid i \in [n], j \in [N]\}$ , which contains all tokens in  $\mathbf{X}^{(i)}$ . It is clear that  $|S| \leq nN$ . Let  $g : S \rightarrow [|S|] = \{1, 2, \dots, |S|\}$  be an arbitrary bijective function. For each  $\mathbf{X}^{(i)}$  with  $i \in [n]$ , we define

$$\mathbf{m}^{(i)} = \sum_{j=1}^N \mathbf{A}_{j,m_i} e_{g(\mathbf{X}_{:,j}^{(i)})} \in \mathbb{N}^{|S|},$$

where  $e_{g(\mathbf{X}_{:,j}^{(i)})}$  is a one-hot vector with 1 in the  $g(\mathbf{X}_{:,j}^{(i)})$ -position. Since there are no repeated tokens in  $\mathbf{X}^{(i)}$  and for any  $i \neq j \in [n]$   $\mathbf{X}^{(i)}$  and  $\mathbf{X}^{(j)}$  are not permutation equivalent, meaning that

$$\|\mathbf{m}^{(i)} - \mathbf{m}^{(j)}\|^2 > 0,$$

for any  $i \neq j \in [n]$ . By applying Lemma H.1 to  $\mathbf{m}^{(1)}, \dots, \mathbf{m}^{(n)}$ , there exists a vector  $\mathbf{v} \in \mathbb{R}^{|S|}$  such that

$$\frac{1}{n^2} \sqrt{\frac{8}{\pi|S|}} \|\mathbf{m}^{(i)} - \mathbf{m}^{(j)}\| \leq \left| \mathbf{v}^\top (\mathbf{m}^{(i)} - \mathbf{m}^{(j)}) \right| \leq \|\mathbf{m}^{(i)} - \mathbf{m}^{(j)}\|$$

holds for any  $i, j \in [n]$ . Let  $\mathbf{h}$  be the function  $\mathbf{h} : S \rightarrow \mathbb{R}, \mathbf{x} \mapsto \mathbf{v}_g(\mathbf{x})$ . Notice that

$$\sum_{j=1}^N \mathbf{A}_{j,m_i} \mathbf{h}(\mathbf{X}_{:,j}^{(i)}) = \sum_{j=1}^N \mathbf{A}_{j,m_i} \mathbf{v}_{g(\mathbf{X}_{:,j}^{(i)})} = \mathbf{v}^\top \mathbf{m}^{(i)}$$

holds for any  $i \in [n]$ . Then, for any  $k \neq l \in [n]$ , we have

$$\begin{aligned} & \left| \sum_{j=1}^N \mathbf{A}_{j,m_k} \mathbf{h}(\mathbf{X}_{:,j}^{(k)}) - \sum_{j=1}^N \mathbf{A}_{j,m_l} \mathbf{h}(\mathbf{X}_{:,j}^{(l)}) \right| \\ &= \left| \mathbf{v}^\top (\mathbf{m}^{(k)} - \mathbf{m}^{(l)}) \right| \\ &\geq \frac{1}{n^2} \sqrt{\frac{8}{\pi|S|}} \|\mathbf{m}^{(k)} - \mathbf{m}^{(l)}\| > 0. \end{aligned}$$

In the following, we use one feed-forward layer to implement  $\mathbf{f}(\cdot)$ . By applying Lemma H.4 to the set  $\{(\mathbf{X}_{:,j}^{(i)})\}_{i \in [n], j \in [N]}$ , there exists a  $\sigma_R$ -activated FFN  $\mathbf{f} \in \mathcal{NN}_{\sigma_R}(3(n-1)N, 1, \mathbb{R}^D \rightarrow \mathbb{R})$  such that

$$\mathbf{f}(\mathbf{X}_{:,j}^{(i)}) = \mathbf{h}(\mathbf{X}_{:,j}^{(i)}).$$

Moreover, according to Lemma E.3, there exists a residual FFN  $\hat{\mathbf{f}} \in \mathcal{NN}_{\sigma_R}^{Res}(\max\{3(n-1)N, D\}, 1, \mathbb{R}^{\max\{3(n-1)N, D\}} \rightarrow \mathbb{R}^{\max\{3(n-1)N, D\}})$  such that

$$\hat{\mathbf{f}}(\widehat{\mathbf{X}}) = \begin{pmatrix} \mathbf{f}(\mathbf{X}) \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{h}(\mathbf{X}) \\ \mathbf{0} \end{pmatrix},$$

where  $\widehat{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \mathbf{0}_{(3(n-1)N-D) \times N} \end{pmatrix}$ . Define  $\mathcal{F}_{SA}$  by

$$\mathcal{F}_{SA}(\mathbf{X}) := \mathbf{X} + \mathbf{W}_O \mathbf{W}_V \mathbf{X} \mathbf{A},$$

where

$$\mathbf{W}_O = \begin{pmatrix} 0 & 1 & 0 & \dots & 0 \end{pmatrix}^\top \in \mathbb{R}^{\max\{3(n-1)N, D\} \times 1},$$

$$\mathbf{W}_V = \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \end{pmatrix} \in \mathbb{R}^{1 \times \max\{3(n-1)N, D\}}.$$

It is straightforward to verify that

$$\mathcal{F}_{SA}(\hat{\mathbf{f}}(\widehat{\mathbf{X}})) = \begin{pmatrix} \mathbf{h}(\mathbf{X}) \\ \mathbf{h}(\mathbf{X}) \mathbf{A} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{h}(\mathbf{X}_{:,1}) & \dots & \mathbf{h}(\mathbf{X}_{:,N}) \\ \sum_{j=1}^N \mathbf{A}_{j,1} \mathbf{h}(\mathbf{X}_{:,j}) & \dots & \sum_{j=1}^N \mathbf{A}_{j,N} \mathbf{h}(\mathbf{X}_{:,j}) \\ \mathbf{0} & \dots & \mathbf{0} \end{pmatrix}.$$

We define

$$\mathcal{E}_{in}(\mathbf{X}) = \begin{pmatrix} \mathbf{X} \\ \mathbf{0}_{(3(n-1)N-D) \times N} \end{pmatrix},$$

and

$$\mathcal{E}_{out}(\mathbf{X}) = \mathbf{X}_{1:2,1:N}.$$

In the last, we let  $\mathcal{F} := \mathcal{E}_{out} \circ \mathcal{F}_{SA} \circ \mathcal{F}_{FF} \circ \mathcal{E}_{in}$ . Through our construction above, it is clear that

$$\mathcal{F} \in \mathcal{T}(\max\{3(n-1)N, D\}, 1, 1, 3(n-1)N, 1).$$

The proof is completed by direct verify that

$$\mathcal{F}(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}(\mathbf{X}^{(j)})_{:,m_j}$$

for any  $i \neq j \in [n]$ .  $\square$

## 1890 G PROOF OF SECTION 6

### 1891 G.1 PROOF OF THEOREM 6.1

1892 *Proof of Theorem 6.1. Step 1:* Identify [MASK] token in different contexts. Note that the [MASK]  
 1893 token in all data points are the same, the following Lemma implies that we can find a single self-  
 1894 attention layer to map each [MASK] to different values according to the contexts. By applying  
 1895 Lemma 5.1x to  $(\mathbf{X}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{X}^{(n)}, \mathbf{y}^{(n)})$ , we know that there exists a self-attention layer  
 1896  $\mathcal{F}_{SA} \in \mathcal{T}(D, 1, 1, 1)$  which has the following form

$$1897 \mathcal{F}_{SA}(\mathbf{X}) = \mathbf{X} + \mathbf{W}_O \mathbf{W}_V \mathbf{X} \sigma_S \left[ (\mathbf{W}_K \mathbf{X})^\top (\mathbf{W}_Q \mathbf{X}) \right]$$

1898 such that

$$1899 \mathcal{F}_{SA}(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}_{SA}(\mathbf{X}^{(j)})_{:,m_j} \quad \text{for any } i \neq j \in [n].$$

1900 Define the embedding layer  $\mathcal{E}_{in}$  as

$$1901 \mathcal{E}_{in}(\mathbf{X}) := \begin{pmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{1}_{1 \times N} & \mathbf{0} \\ \mathbf{0}_{(3n-D)^+ \times N} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbb{R}^{(\max\{3n, D\} + 2) \times (N+1)}.$$

1902 Then, we define the following matrices to adapt to the embedded input

$$1903 \mathbf{W}_K^{(0)} := \begin{pmatrix} \mathbf{W}_K & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{1 \times (\max\{3n, D\})},$$

$$1904 \mathbf{W}_Q^{(0)} := \begin{pmatrix} \mathbf{W}_Q & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{1 \times (\max\{3n, D\})},$$

$$1905 \mathbf{W}_V^{(0)} := \begin{pmatrix} \mathbf{W}_V & \mathbf{0} \end{pmatrix} \in \mathbb{R}^{1 \times (\max\{3n, D\})},$$

$$1906 \mathbf{W}_O^{(0)} := \begin{pmatrix} \mathbf{W}_O \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{(\max\{3n, D\}) \times 1}.$$

1907 Construct self-attention layer  $\mathcal{F}_{SA}^{(0)}$  by  $\mathbf{W}_K^{(0)}, \mathbf{W}_Q^{(0)}, \mathbf{W}_V^{(0)}, \mathbf{W}_O^{(0)}$  and the mask matrix  $\mathbf{P}_0$ , that is

$$1908 \mathcal{F}_{SA}^{(0)}(\mathcal{E}_{in}(\mathbf{X})) := \mathcal{E}_{in}(\mathbf{X}) + \mathbf{W}_O^{(0)} \mathbf{W}_V^{(0)} \mathcal{E}_{in}(\mathbf{X}) \sigma_S \left[ \left( \mathbf{W}_K^{(0)} \mathcal{E}_{in}(\mathbf{X}) \right)^\top \left( \mathbf{W}_Q^{(0)} \mathcal{E}_{in}(\mathbf{X}) \right) + \mathbf{P}_0 \right]$$

$$1909 = \mathcal{E}_{in}(\mathbf{X}) + \begin{pmatrix} \mathbf{W}_O \mathbf{W}_V \mathbf{X} \sigma_S \left[ (\mathbf{W}_K \mathbf{X})^\top (\mathbf{W}_Q \mathbf{X}) \right] & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$

$$1910 = \begin{pmatrix} \mathcal{F}_{SA}(\mathbf{X}) & \mathbf{0} \\ \mathbf{1}_{1 \times N} & \mathbf{0} \\ \mathbf{0}_{(3n-D)^+ \times N} & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix}.$$

1911 **Step 2:** Point fitting by FFNs.

1912 Recall that the effect of  $\mathcal{F}_{SA}$  is to distinguish the [MASK] token in each input, that is

$$1913 \mathcal{F}_{SA}(\mathbf{X}^{(i)})_{:,m_i} \neq \mathcal{F}_{SA}(\mathbf{X}^{(j)})_{:,m_j}.$$

1914 By applying the Lemma H.4 to the set  $\{\mathcal{F}_{SA}(\mathbf{X}^{(i)})_{:,m_i}\}_{i \in [n]}$ , we know that there exists a  $\sigma_R$ -  
 1915 activated FFN  $f \in \mathcal{NN}(3n, 1, \mathbb{R}^D \rightarrow \mathbb{R})$  such that

$$1916 f(\mathcal{F}_{SA}(\mathbf{X}^{(i)})_{:,m_i}) = y_i \quad \text{for any } i \in [n].$$

Then, by applying Proposition 4.1 to  $f$ , there exists an attention-only Transformer  $\mathcal{F}_1 \in \mathcal{T}(\max\{3n, D\} + 2, \max\{3n, D\} + 1, 1, 1)$  and  $M > 0$  with the following form

$$\mathcal{F}_1(\mathbf{X}) = \mathcal{E}_{out} \circ \mathcal{F}_{SA}^{(1)} \circ \mathcal{E}_{in}(\mathbf{X})$$

satisfying

$$\|\mathcal{F}_1(\mathbf{X}) - f(\mathbf{X})\|_\infty < \varepsilon \quad \text{for any } \mathbf{X} \in [-M, M]^{D \times N}.$$

Define  $\mathcal{F}$  as

$$\mathcal{F}(\mathbf{X}) := \mathcal{E}_{out} \circ \mathcal{F}_{SA}^{(1)} \circ \mathcal{F}_{SA}^{(0)} \circ \mathcal{E}_{in}(\mathbf{X}).$$

The proof is completed by verifying that for any  $\varepsilon > 0$

$$\left| \mathcal{F}(\mathbf{X}^{(i)})_{:,m_i} - y_i \right| < \varepsilon \quad \text{for any } i \in [n],$$

and

$$\mathcal{F} \in \mathcal{T}(\max\{3n, D\} + 2, \max\{3n, D\} + 1, 1, 2).$$

□

## H SUPPORTING LEMMAS

**Lemma H.1** (Park et al. (2021)). *Let  $d \in \mathbb{N}$ . Then, for any finite subst  $\mathcal{X} \subset \mathbb{R}^d$ , there exists a unit vector  $\mathbf{v} \in \mathbb{R}^d$  such that*

$$\frac{1}{|\mathcal{X}|^2} \sqrt{\frac{8}{\pi d}} \|\mathbf{x} - \mathbf{x}'\| \leq \|\mathbf{v}'(\mathbf{x} - \mathbf{x}')\| \leq \|\mathbf{x} - \mathbf{x}'\| \leq \|\mathbf{x} - \mathbf{x}'\|.$$

**Lemma H.2** (Kajitsuka & Sato (2023)). *Given a finite subset  $\mathcal{X} \subset \mathbb{R}^d$ . Then, for any  $\delta > 0$ , there exists matrices  $\mathbf{W}_K, \mathbf{W}_Q \in \mathbb{R}^{1 \times d}$  such that*

$$\left| (\mathbf{W}_K \mathbf{x}_a)^\top (\mathbf{W}_Q \mathbf{x}_c) - (\mathbf{W}_K \mathbf{x}_b)^\top (\mathbf{W}_Q \mathbf{x}_c) \right| > \delta$$

for any  $\mathbf{x}_a, \mathbf{x}_b, \mathbf{x}_c \in \mathcal{X}$  with  $\mathbf{x}_a \neq \mathbf{x}_b$ .

*Proof of Lemma H.2.* By applying Lemma H.1 to  $\mathcal{X} \cup \{0\}$ , we know that there exists a unit vector  $\mathbf{v} \in \mathbb{R}^d$  such that for any  $\mathbf{x}_a, \mathbf{x}_b \in \mathcal{X} \cup \{0\}$  such that  $\mathbf{x}_a \neq \mathbf{x}_b$ , we have

$$\frac{1}{(|\mathcal{V} + 1|)^2} \sqrt{\frac{8}{\pi d}} \|\mathbf{x}_a - \mathbf{x}_b\| \leq |\mathbf{v}^\top (\mathbf{x}_a - \mathbf{x}_b)| \leq \|\mathbf{x}_a - \mathbf{x}_b\|$$

In particular, if we let one of  $\mathbf{x}_a$  and  $\mathbf{x}_b$  be 0, we have

$$\frac{1}{(|\mathcal{V} + 1|)^2} \sqrt{\frac{8}{\pi d}} \|\mathbf{x}_c\| \leq |\mathbf{v}^\top \mathbf{x}_c| \leq \|\mathbf{x}_c\| \quad \text{for any } \mathbf{x}_c \in \mathcal{X}.$$

Let  $\varepsilon = \min\{\|\mathbf{x}_a - \mathbf{x}_b\| \mid \mathbf{x}_a \neq \mathbf{x}_b \in \mathcal{X}\}$  and  $r = \min\{\|\mathbf{x}\| \mid \mathbf{x} \in \mathcal{X}\}$ . Assume that  $r \neq 0$ . Then, we can pick up arbitrary vectors  $\mathbf{u}, \mathbf{u}' \in \mathbb{R}$  with

$$|\mathbf{u}\mathbf{u}'| = (|\mathcal{X}| + 1)^4 \frac{\pi d}{8} \frac{\delta}{\varepsilon r}.$$

Let  $\mathbf{W}_K = \mathbf{u}\mathbf{v}^\top$  and  $\mathbf{W}_Q = \mathbf{u}'\mathbf{v}^\top$ , we have

$$\begin{aligned} & \left| (\mathbf{W}_K \mathbf{x}_a)^\top (\mathbf{W}_Q \mathbf{x}_c) - (\mathbf{W}_K \mathbf{x}_b)^\top (\mathbf{W}_Q \mathbf{x}_c) \right| \\ &= \left| (\mathbf{x}_a - \mathbf{x}_b)^\top (\mathbf{W}_K)^\top (\mathbf{W}_Q \mathbf{x}_c) \right| \\ &= \left| (\mathbf{x}_a - \mathbf{x}_b)^\top \mathbf{v} \right| \cdot |\mathbf{u}\mathbf{u}'| \cdot |\mathbf{v}^\top \mathbf{x}_c| \\ &\geq \delta, \end{aligned}$$

which completes the proof. □

**Lemma H.3** (Kajitsuka & Sato (2023)). Let  $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(n)} \in \mathbb{R}^d$  be pair-wise distinct vectors satisfying

$$\|\mathbf{a}^{(i)} - \mathbf{a}^{(j)}\| > 2 \log n + 3 \quad \text{for any } i \neq j \in [n].$$

Then, the following holds

$$\left(\mathbf{a}^{(i)}\right)^\top \sigma_S \left[\mathbf{a}^{(i)}\right] \neq \left(\mathbf{a}^{(j)}\right)^\top \sigma_S \left[\mathbf{a}^{(j)}\right]$$

for any  $i \neq j \in [n]$ .

**Lemma H.4** (Point fitting by ReLU FFNs, modified from (Jiao et al., 2025a)). Let  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  be a set of input-output pairs such that  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$  and  $\mathbf{x}_i \neq \mathbf{x}_j$  if  $i \neq j$ . Then, there exists a  $\sigma_R$ -activated feed-forward network  $\mathbf{f} \in \mathcal{NN}(3n, 1, \mathbb{R}^d \rightarrow \mathbb{R})$  such that

$$\mathbf{f}(\mathbf{x}_i) = y_i \quad \text{for any } i \in [n].$$

*Proof of Lemma H.4.* Let  $K > 0$  be determined later. Since  $\mathbf{x}_i$  are pair-wise distinct. According to Lemma H.1, there exists a vector  $\mathbf{v} \in \mathbb{R}^d$  such that  $\mathbf{v}^\top \mathbf{x}_i, i \in [n]$  are also pair-wise distinct. We define

$$\mathbf{W}_1^{(i)} = K \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \mathbf{v}^\top, \quad \mathbf{b}_1^{(i)} = \begin{pmatrix} -K\mathbf{v}^\top \mathbf{x}_i - 1 \\ -K\mathbf{v}^\top \mathbf{x}_i \\ -K\mathbf{v}^\top \mathbf{x}_i + 1 \end{pmatrix}, \quad \mathbf{W}_2^{(i)} = y_i \begin{pmatrix} 1 & -2 & 1 \end{pmatrix}, \quad \mathbf{b}_2^{(i)} = \mathbf{0},$$

where  $\mathbf{W}_1^{(i)} \in \mathbb{R}^{3 \times d}$ ,  $\mathbf{b}_1^{(i)} \in \mathbb{R}^3$ ,  $\mathbf{W}_2^{(i)} \in \mathbb{R}^{1 \times 3}$ ,  $\mathbf{b}_2^{(i)} \in \mathbb{R}$ . It is straightforward to verify that the following holds for any  $\mathbf{x} \in \mathbb{R}^d$

$$\begin{aligned} & \mathbf{W}_2^{(i)} \sigma_R \left( \mathbf{W}_1^{(i)} \mathbf{x} + \mathbf{b}_1^{(i)} \right) + \mathbf{b}_2^{(i)} \\ &= y_i \left( \sigma_R \left( K\mathbf{v}^\top (\mathbf{x} - \mathbf{x}_i) - 1 \right) - 2\sigma_R \left( K\mathbf{v}^\top (\mathbf{x} - \mathbf{x}_i) \right) + \sigma_R \left( K\mathbf{v}^\top (\mathbf{x} - \mathbf{x}_i) + 1 \right) \right) \\ &= y_i \mathbf{I}_i(\mathbf{x}), \end{aligned}$$

where  $\mathbf{I}_i(\mathbf{x})$  satisfies  $\mathbf{I}_i(\mathbf{x}_i) = 1$  and  $\mathbf{I}_i(\mathbf{x}_j) = 0$  if  $|\mathbf{v}^\top (\mathbf{x} - \mathbf{x}_i)| \geq \frac{1}{K}$ . So, we choose

$$K > \frac{2}{\min_{i \neq j} |\mathbf{v}^\top (\mathbf{x}_i - \mathbf{x}_j)|}.$$

Define

$$\mathbf{W}_1 = \begin{pmatrix} \mathbf{W}_1^{(1)} \\ \vdots \\ \mathbf{W}_1^{(n)} \end{pmatrix}, \quad \mathbf{b}_1 = \begin{pmatrix} \mathbf{b}_1^{(1)} \\ \vdots \\ \mathbf{b}_1^{(n)} \end{pmatrix}, \quad \mathbf{W}_2 = \begin{pmatrix} \mathbf{W}_2^{(1)} & \dots & \mathbf{W}_2^{(n)} \end{pmatrix}, \quad \mathbf{b}_2 = \mathbf{0}$$

where  $\mathbf{W}_1 \in \mathbb{R}^{3n \times d}$ ,  $\mathbf{b}_1 \in \mathbb{R}^{3n}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{1 \times 3n}$ ,  $\mathbf{b}_2 \in \mathbb{R}$ . Let

$$\begin{aligned} \mathbf{f}(\mathbf{x}) &= \mathbf{W}_2 \sigma_R \left( \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 \right) + \mathbf{b}_2 \\ &= \sum_{i=1}^n y_i \mathbf{I}_i(\mathbf{x}). \end{aligned}$$

It is clear that  $\mathbf{f} \in \mathcal{NN}(3n, 1, \mathbb{R}^d \rightarrow \mathbb{R})$ . The proof is completed by pointing out that

$$\mathbf{f}(\mathbf{x}_i) = y_i \quad \text{for any } i \in [n].$$

□

## I THE USE OF LARGE LANGUAGE MODELS

The LLMs are used in the following aspects:

- Improving the clarity of writing and grammar.
- Drafting code snippets for preliminary experiments.
- Summarizing related literature for internal discussion.

We carefully reviewed and revised any text or code suggested by LLMs, and remain fully responsible for the scientific accuracy, originality, and integrity of this work.