
Are Transformers universal approximators of sequence-to-sequence functions?

Chulhee Yun*

MIT

chulheey@mit.edu

Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J. Reddi, Sanjiv Kumar

Google Research New York

{bsrinadh, ankitsrawat, sashank, sanjivk}@google.com

Abstract

Despite the widespread adoption of Transformer models for NLP tasks, the expressive power of these models is not well-understood. In this paper, we establish that Transformer models are universal approximators of continuous *permutation equivariant* sequence-to-sequence functions with compact support, which is quite surprising given the amount of shared parameters in these models. Furthermore, using positional encodings, we circumvent the restriction of permutation equivariance, and show that Transformer models can universally approximate *arbitrary* continuous sequence-to-sequence functions on a compact domain. Interestingly, our proof techniques clearly highlight the different roles of the self-attention and the feed-forward layers in Transformers. In particular, we prove that fixed width self-attention layers can compute *contextual mappings* of the input sequences, playing a key role in the universal approximation property of Transformers. Based on this insight from our analysis, we consider other simpler alternatives to self-attention layers and empirically evaluate them.

1 Introduction

Self-attention based Transformer networks [Vaswani et al., 2017] have been at the center of the recent progress on various natural language processing (NLP) tasks, including machine translation [Vaswani et al., 2017], language modeling [Radford et al., 2018, 2019], and question answering [Devlin et al., 2018, Yang et al., 2019, Liu et al., 2019]. All these tasks involve learning models that map an input sequence of tokens to an output sequence of tokens. Transformers make it feasible to train large models to approximate these sequence-to-sequence functions due to their ability to process the input tokens in a parallel way, as opposed to the sequential nature of RNNs and LSTMs.

A Transformer block consists of two kinds of layers: a self-attention layer and a token-wise feed-forward layer, with skip connections present in both layers. The self-attention layer transforms each input token embedding using a weighted combination of the embeddings of all tokens in the input sequence, where weights are generated by pairwise dot-products among the input token embeddings. The token-wise feed-forward layer then independently processes each of these modified input token embeddings without any interaction among them. Notably, Transformers employ parameter reuse across tokens, as both layers use the same parameters to process each token. Moreover, Transformers have to rely solely on the pairwise dot-products to capture interaction between the input tokens.

*Based on work performed at Google Research New York

Given the parameter sharing and limited interactions between tokens, it is natural to wonder: what class of sequence-to-sequence functions can the Transformers represent? What is the role of the two different kinds of layers? Are both layers needed to obtain the representation power of Transformers? In the existing literature, the advantage of Transformers has often been attributed to their capability of computing *contextual* embeddings/mappings of the input, as opposed to fixed word embeddings as in word2vec [Mikolov et al., 2013]. Is it possible to formalize the notion of contextual mappings? Can Transformers compute such mappings? Such questions still remain elusive.

1.1 Our contributions

In this paper, we present results that shed light on the expressive power of Transformer networks. We prove that Transformers are universal approximators of any permutation equivariant sequence-to-sequence function. We also improve this result using positional encodings, and show that Transformers can represent any sequence-to-sequence function; i.e., the restriction of permutation equivariance can be removed by positional encodings.

Interestingly, our proof techniques clearly highlight different roles of the self-attention and the feed-forward layers. Due to limited space, the proofs are outlined in the appendix; see Sections C–G. In particular, in Section E, we provide a mathematical definition of *contextual mappings* and surprisingly show that multi-head self-attention layers can indeed compute contextual mappings of the input sequences, where unique contexts are mapped to unique vectors. This ability to compute contextual mappings coupled with the value mapping ability of the feed-forward layers makes Transformers universal approximators.

Given that self-attention layers can compute contextual mappings, is it possible to have a more efficient architecture to do the same task, consequently, preserving the ability to universally approximate sequence-to-sequence functions? Towards this, we discuss other simpler architectures that can implement contextual mappings to some extent, and experimentally evaluate their performance. In our experiments, we notice that the models that combine these simpler architectures with Transformers have better performance, compared to the standalone Transformers. Due to space limits, we present experimental results in Section H.

1.2 Related works & notation

Analysis of attention-based models. Given the popularity of Transformers, there have been numerous works trying to understand the role of attention layers in natural language processing models. One such line of work focuses on probing the output of attention layers to understand the attention mechanism and internal language representation [Hewitt and Manning, 2019, Clark et al., 2019, Coenen et al., 2019, Vig and Belinkov, 2019]. Although these results give valuable insights, a consistent theoretical analysis corroborating these findings is missing.

Universal approximation theorems. Universal approximation theorems are classical results in neural network theory, dating back many decades [Cybenko, 1989, Hornik, 1991]. These results show that given unbounded width, a shallow neural network can approximate arbitrary continuous function with compact support, up to any accuracy. Other results focusing on depth appeared more recently [Lu et al., 2017, Hanin and Sellke, 2017, Lin and Jegelka, 2018]. In particular, Lu et al. [2017], Hanin and Sellke [2017] consider fully-connected ReLU networks whose input dimension is d , and show that networks with width $d + 1$ and unbounded depth are universal approximators of scalar-valued continuous functions. Lin and Jegelka [2018] show that a residual network with one hidden neuron per residual block is a universal approximator of scalar-valued functions, given unbounded depth. Although Transformer networks do have residual connections, due to their heavy parameter sharing, the existing analyses for residual networks do not extend to Transformers. Sannai et al. [2019] consider universally approximating permutation invariant/equivariant functions using fully-connected ReLU networks.

Turing completeness results on Transformers. Recently, Pérez et al. [2019] have shown that Transformers with infinite precision are Turing complete, which is not the case in finite precision setting [Dehghani et al., 2018]. We note that Turing completeness deals with computation on formal languages (thus discrete objects), while universal approximation focuses on functions on a continuum. In other words, these are two different concepts; and one does not imply another.

Notation. Given a matrix \mathbf{A} , we use $\|\mathbf{A}\|_p$ to denote the entry-wise ℓ^p norm of \mathbf{A} . Let $\sigma[\cdot]$ be the softmax operator, which takes a matrix as input and applies softmax operation to each column of the

matrix, which results in a column stochastic matrix, i.e., a matrix that has non-negative entries with all columns summing to 1. We use $\mathbf{1}_n$ to denote a vector of length n whose entries are all 1. We use d and n to denote the embedding dimension and the sequence length, respectively. We assume throughout that $n \geq 2$, as the Transformers reduce to residual networks when $n = 1$.

2 Transformer networks

A Transformer block is a sequence-to-sequence function mapping $\mathbb{R}^{d \times n}$ to $\mathbb{R}^{d \times n}$. It consists of two layers: a self-attention layer and a token-wise feed-forward layer, with both layers having a skip connection. More concretely, for an input $\mathbf{X} \in \mathbb{R}^{d \times n}$ consisting of d -dimensional embeddings of n tokens, a Transformer block with *multiplicative* or *dot-product* attention [Luong et al., 2015] consists of the following two layers¹:

$$\text{Attn}(\mathbf{X}) = \mathbf{X} + \sum_{i=1}^h \mathbf{W}_O^i \mathbf{W}_V^i \mathbf{X} \cdot \sigma[(\mathbf{W}_K^i \mathbf{X})^T \mathbf{W}_Q^i \mathbf{X}], \quad (1)$$

$$\text{FF}(\mathbf{X}) = \text{Attn}(\mathbf{X}) + \mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1 \cdot \text{Attn}(\mathbf{X}) + \mathbf{b}_1 \mathbf{1}_n^T) + \mathbf{b}_2 \mathbf{1}_n^T, \quad (2)$$

where $\mathbf{W}_O^i \in \mathbb{R}^{d \times m}$, $\mathbf{W}_V^i, \mathbf{W}_K^i, \mathbf{W}_Q^i \in \mathbb{R}^{m \times d}$, $\mathbf{W}_2 \in \mathbb{R}^{d \times r}$, $\mathbf{W}_1 \in \mathbb{R}^{r \times d}$, $\mathbf{b}_2 \in \mathbb{R}^d$, $\mathbf{b}_1 \in \mathbb{R}^r$, and $\text{FF}(\mathbf{X})$ is the output of the Transformer block. The number of heads h and the head size m are two main parameters of the attention layer; and r denotes the hidden layer size of the feed-forward layer.

Here, we would like to point out that our definition of the self-attention layer (1) is an equivalent reformulation of [Vaswani et al., 2017], where they concatenate attention heads and multiply a matrix $\mathbf{W}_O \in \mathbb{R}^{d \times mh}$ to the concatenation. One difference in our setup is the absence of layer normalization, which simplifies our analysis while preserving the basic architecture of the Transformer.

We define the Transformer networks as the composition of Transformer blocks. The family of the sequence-to-sequence functions corresponding to the Transformers can be defined as:

$$\mathcal{T}^{h,m,r} := \{g : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n} \mid g \text{ is a composition of Transformer blocks } t^{h,m,r}\text{'s}\}. \quad (3)$$

where $t^{h,m,r} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ denotes a Transformer block defined by an attention layer with h heads of size m each, and a feed-forward layer with r hidden nodes.

We say that a function $f : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ is *permutation equivariant* if for any permutation matrix \mathbf{P} , we have $f(\mathbf{X}\mathbf{P}) = f(\mathbf{X})\mathbf{P}$; i.e., if we permute the columns of \mathbf{X} , then the columns of $f(\mathbf{X})$ are permuted in the same way. A Transformer block is permutation equivariant, which we formally prove in Section B. This consequently establishes the permutation equivariance of the class $\mathcal{T}^{h,m,r}$.

Claim 1. A Transformer block $t^{h,m,r}$ defines a permutation equivariant map from $\mathbb{R}^{d \times n}$ to $\mathbb{R}^{d \times n}$.

3 Transformers are universal approximators of seq-to-seq functions

In this section, we present our results showing that the Transformer networks are universal approximators of sequence-to-sequence functions. Let us start by defining the target function class \mathcal{F}_{PE} , which consists of all continuous permutation equivariant functions with compact support that map $\mathbb{R}^{d \times n}$ to $\mathbb{R}^{d \times n}$. Here, continuity is defined with respect to any entry-wise ℓ^p norm, $1 \leq p < \infty$. Given two functions $f_1, f_2 : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$, for $1 \leq p < \infty$, we define a distance between them as

$$d_p(f_1, f_2) := \left(\int \|f_1(\mathbf{X}) - f_2(\mathbf{X})\|_p^p d\mathbf{X} \right)^{1/p}.$$

The following result shows that a Transformer network with a constant number of heads h , head size m , and hidden layer of size r can approximate any function in \mathcal{F}_{PE} .

Theorem 2. Let $1 \leq p < \infty$ and $\epsilon > 0$, then for any given $f \in \mathcal{F}_{\text{PE}}$, there exists a Transformer network $g \in \mathcal{T}^{2,1,4}$, such that $d_p(f, g) \leq \epsilon$.

Here, we would like to point out that Theorem 2 appears quite surprising at a first glance, given the parameter sharing across all the tokens in a sequence, e.g., feed-forward layers are applied token-wise and the projection matrices in the self-attention layers are the same across different tokens. Furthermore, attention layers can only capture pairwise interaction between different tokens in the sequence. In Section E.1, we make it clear how Transformer networks are able to circumvent these apparent restrictions imposed by parameter sharing and limited interaction among tokens, and achieve universal approximation of the function class \mathcal{F}_{PE} .

¹In our proof we use bias vectors \mathbf{b}_Q^i for query projections in attention layers. We omit them here for brevity.

3.1 Transformers with trainable positional encodings

In order to endow the Transformer networks with the ability to capture the information about the position of tokens in the input sequence, it is a common practice to add positional encodings $\mathbf{E} \in \mathbb{R}^{d \times n}$ to the input sequence before feeding it to the Transformer network [Vaswani et al., 2017, Devlin et al., 2018]. Consider the functions represented by Transformers with positional encodings:

$$\mathcal{T}_P^{h,m,r} := \{g_P(\mathbf{X}) = g(\mathbf{X} + \mathbf{E}) \mid g \in \mathcal{T}^{h,m,r} \text{ and } \mathbf{E} \in \mathbb{R}^{d \times n}\}. \quad (4)$$

Here we show that if \mathbf{E} is trainable, these positional encodings are sufficient to remove the permutation equivariance restriction of the Transformers. Towards this, we define \mathcal{F}_{CD} to be the set of all continuous functions that map a compact domain in $\mathbb{R}^{d \times n}$ to $\mathbb{R}^{d \times n}$. Note that \mathcal{F}_{CD} does not have the restriction of permutation equivariance as in \mathcal{F}_{PE} , but any $f \in \mathcal{F}_{CD}$ is defined on a compact domain instead of the whole $\mathbb{R}^{d \times n}$. The following result states that, equipped with the trainable positional encodings, Transformers can approximate any sequence-to-sequence function in \mathcal{F}_{CD} .

Theorem 3. *Let $1 \leq p < \infty$ and $\epsilon > 0$, then for any given $f \in \mathcal{F}_{CD}$, there exists a Transformer network $g \in \mathcal{T}_P^{2,1,4}$ such that we have $d_p(f, g) \leq \epsilon$.*

Theorems 2 and 3 provide an interesting characterization of the representation power of fixed-width Transformer networks. Since the function classes $\mathcal{T}^{h,m,r}$ and $\mathcal{T}_P^{h,m,r}$ become richer as we increase the values of (h, m, r) , our results establish that general Transformer networks are also universal approximators of sequence-to-sequence functions. Remarkably, none of the parameters (h, m, r) depend on the input sequence length n or embedding dimension d .

3.2 Tightness of constructions

We refer the reader to Section C and Section G for the proofs of Theorem 2 and Theorem 3, respectively. Even though Theorems 2 and 3 do not specifically mention the required depth for approximation, our proof techniques do characterize it. The rough idea of the proof is to approximate continuous functions by piece-wise constant functions on cubes of width δ , and then construct Transformer networks that approximate the piece-wise constant functions.

We show in Section E that Theorem 2 requires $O(n(1/\delta)^{dn}/n!)$ Transformer blocks for approximation, where δ is the width of the cubes. Each transformer block is of constant width, so it has $O(d)$ parameters; this means that the total number of parameters is $O(dn(1/\delta)^{dn}/n!)$. We note that this exponential dependence cannot be avoided in the worse case. If we assume continuity without any additional smoothness, quantizing the domain to cubes and approximating the function with constants require memorizing $(\text{output dim}) \times (\text{num cubes})/n!$ real numbers, where the factor of $1/n!$ is due to permutation equivariance. Thus, Theorem 2 is optimal in the order of parameters.

If we compare with the residual network result [Lin and Jegelka, 2018], we can consider “flattening” \mathbf{X} into a dn -dimensional vector and fitting the function. The proof technique in [Lin and Jegelka, 2018] requires $O((1/\delta)^{dn})$ layers, where each layer has $O(dn)$ parameters: the total parameter requirement is $O(dn(1/\delta)^{dn})$. This shows that Transformers can approximate permutation equivariant functions in a more efficient way than residual networks.

In Section G, our proof of Theorem 3 shows that we require $O(n(1/\delta)^{dn})$ layers to approximate continuous (not permutation equivariant) sequence-to-sequence functions. As seen from the argument above, this construction is also optimal in the order of parameters.

4 Conclusion

In this paper, we prove that Transformer networks are universal approximators of any continuous and permutation equivariant sequence-to-sequence functions, which shed light on the expressive power of Transformer networks. We also theoretically validate the use of additive positional encodings in Transformers, as they can remove the permutation equivariance restriction and make Transformers universal approximators of arbitrary continuous sequence-to-sequence functions.

In the supplementary material, we present the proofs of our theorems, which reveal that self-attention layers in Transformer networks can compute *contextual mappings*; this is one of the crucial components that make Transformer networks universal. We also discuss and experiment with other simpler layers that can implement weaker forms of contextual mappings.

References

- Nir Ailon and Bernard Chazelle. The fast Johnson–Lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on computing*, 39(1):302–322, 2009.
- François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. What does BERT look at? an analysis of BERT’s attention. *arXiv preprint arXiv:1906.04341*, 2019.
- Andy Coenen, Emily Reif, Ann Yuan, Been Kim, Adam Pearce, Fernanda Viégas, and Martin Wattenberg. Visualizing and measuring the geometry of BERT. *arXiv preprint arXiv:1906.02715*, 2019.
- George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal transformers. *arXiv preprint arXiv:1807.03819*, 2018.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Yunchao Gong, Sanjiv Kumar, Henry A Rowley, and Svetlana Lazebnik. Learning binary codes for high-dimensional data using bilinear projections. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 484–491, 2013.
- Boris Hanin and Mark Sellke. Approximating continuous functions by relu nets of minimal width. *arXiv preprint arXiv:1710.11278*, 2017.
- John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138, 2019.
- Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- Lukasz Kaiser, Aidan N Gomez, and Francois Chollet. Depthwise separable convolutions for neural machine translation. *arXiv preprint arXiv:1706.03059*, 2017.
- Hongzhou Lin and Stefanie Jegelka. ResNet with one-neuron hidden layers is a universal approximator. In *Advances in Neural Information Processing Systems*, pages 6169–6178, 2018.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT pre-training approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Zhou Lu, Hongming Pu, Feicheng Wang, Zhiqiang Hu, and Liwei Wang. The expressive power of neural networks: A view from the width. In *Advances in neural information processing systems*, pages 6231–6239, 2017.
- Minh-Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1412–1421, Lisbon, Portugal, September 2015. Association for Computational Linguistics.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- Jorge Pérez, Javier Marinković, and Pablo Barceló. On the Turing completeness of modern neural network architectures. *arXiv preprint arXiv:1901.03429*, 2019.

- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *Technical Report, OpenAI*, 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *Technical Report, OpenAI*, 2019.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, 2016.
- Akiyoshi Sannai, Yuuki Takai, and Matthieu Cordonnier. Universal approximations of permutation invariant/equivariant functions by deep neural networks. *arXiv preprint arXiv:1903.01939*, 2019.
- Laurent Sifre and Stéphane Mallat. Rigid-motion scattering for image classification. *Ph. D. dissertation*, 2014.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language model. *arXiv preprint arXiv:1906.04284*, 2019.
- Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics, 2018. URL <http://aclweb.org/anthology/N18-1101>.
- Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. *arXiv preprint arXiv:1901.10430*, 2019.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. XLNet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27, 2015.

A Notation

We first summarize the additional notation used in the appendix.

For $a, b \in \mathbb{N}$ where $a \leq b$, let $[a] = \{1, \dots, a\}$ and $[a : b] = \{a, a+1, \dots, b-1, b\}$. For $a, b, c \in \mathbb{R}$ where $b-a > 0$ is an integer multiple of $c > 0$, we write $[a : c : b] := \{a, a+c, a+2c, \dots, b-c, b\}$. Given a matrix \mathbf{A} , let $A_{i,j}$, $\mathbf{A}_{i,:}$, and $\mathbf{A}_{:,j}$ denote its (i, j) -th entry, i -th row, and j -th column, respectively. In a similar way as the softmax operator $\sigma[\cdot]$, we define $\sigma_{\text{H}}[\cdot]$ to be the hardmax operator, which outputs the one-hot representation of the $\arg \max$ entry for each column of the input matrix. If there are k $\arg \max$ entries, then the output is $1/k$ for such entries. We denote the 0-1 indicator function by $\mathbb{1}\{\cdot\}$.

B Proof of Claim 1

Suppose $\mathbf{X}\mathbf{P}$ was given as input, where \mathbf{P} is a permutation matrix. First note that

$$(\mathbf{W}_K^i \mathbf{X}\mathbf{P})^T (\mathbf{W}_Q^i \mathbf{X}\mathbf{P}) = \mathbf{P}^T (\mathbf{W}_K^i \mathbf{X})^T (\mathbf{W}_Q^i \mathbf{X}) \mathbf{P}$$

After the softmax operation, we get

$$\sigma[\mathbf{P}^T (\mathbf{W}_K^i \mathbf{X})^T (\mathbf{W}_Q^i \mathbf{X}) \mathbf{P}] = \mathbf{P}^T \sigma[(\mathbf{W}_K^i \mathbf{X})^T (\mathbf{W}_Q^i \mathbf{X})] \mathbf{P}.$$

Then,

$$\text{Attn}(\mathbf{X}\mathbf{P}) = \mathbf{X}\mathbf{P} + \sum_{i=1}^h \mathbf{W}_O^i (\mathbf{W}_V^i \mathbf{X}\mathbf{P}) \cdot \mathbf{P}^T \sigma[(\mathbf{W}_K^i \mathbf{X})^T (\mathbf{W}_Q^i \mathbf{X})] \mathbf{P} = \text{Attn}(\mathbf{X})\mathbf{P},$$

where we used $\mathbf{P}\mathbf{P}^T = \mathbf{I}$. Permutation equivariance of the token-wise feed-forward layer can be shown similarly:

$$\begin{aligned} \text{FF}(\mathbf{X}\mathbf{P}) &= \text{Attn}(\mathbf{X})\mathbf{P} + \mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1 \cdot \text{Attn}(\mathbf{X})\mathbf{P} + \mathbf{b}_1 \mathbf{1}_n^T \mathbf{P}) + \mathbf{b}_2 \mathbf{1}_n^T \mathbf{P} \\ &= \text{Attn}(\mathbf{X})\mathbf{P} + \mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1 \cdot \text{Attn}(\mathbf{X}) + \mathbf{b}_1 \mathbf{1}_n^T) \mathbf{P} + \mathbf{b}_2 \mathbf{1}_n^T \mathbf{P} = \text{FF}(\mathbf{X})\mathbf{P}, \end{aligned}$$

where $\text{ReLU}(\mathbf{X}\mathbf{P}) = \text{ReLU}(\mathbf{X})\mathbf{P}$ was used. This analysis shows that the function class $\mathcal{T}^{h,m,r}(\cdot)$ is restricted to permutation equivariant functions.

C Proof of Theorem 2

Recall that we want to show that given a function $f \in \mathcal{F}_{\text{PE}}$, we can find a Transformer network $g \in \mathcal{T}^{2,1,4}$ such that $d_p(f, g) \leq \epsilon$. Without loss of generality, we can assume that the compact support of f is contained in $[0, 1]^{d \times n}$. We achieve our desired objective in three key steps:

Step 1. Approximate \mathcal{F}_{PE} with piece-wise constant functions. We first use (a variant of) the classical result that any continuous function can be approximated up to arbitrary accuracy by piece-wise constant functions. For $\delta > 0$, we define the following class of piece-wise constant functions.

$$\overline{\mathcal{F}}_{\text{PE}}(\delta) := \left\{ f : \mathbf{X} \mapsto \sum_{\mathbf{L} \in \mathbb{G}_\delta} \mathbf{A}_{\mathbf{L}} \mathbb{1}\{\mathbf{X} \in \mathbb{S}_{\mathbf{L}}\} \mid f \text{ is permutation equivariant, } \mathbf{A}_{\mathbf{L}} \in \mathbb{R}^{d \times n} \right\},$$

where $\mathbb{G}_\delta := \{0, \delta, \dots, 1 - \delta\}^{d \times n}$ and, for a grid point $\mathbf{L} \in \mathbb{G}_\delta$, $\mathbb{S}_{\mathbf{L}} := \prod_{j=1}^d \prod_{k=1}^n [L_{j,k}, L_{j,k} + \delta) \subset [0, 1]^{d \times n}$ denotes the associated cube of width δ .

The following result states that the underlying function $f \in \mathcal{F}_{\text{PE}}$ can be approximated using the function class $\overline{\mathcal{F}}_{\text{PE}}(\delta)$.

Lemma 4. *For any given $f \in \mathcal{F}_{\text{PE}}$ and $1 \leq p < \infty$, one can find a $\delta^* > 0$ such that $\exists \overline{f} \in \overline{\mathcal{F}}_{\text{PE}}(\delta^*)$ which satisfies $d_p(f, \overline{f}) \leq \epsilon/3$.*

Step 2. Approximate $\overline{\mathcal{F}}_{\text{PE}}(\delta)$ with modified Transformers. We then consider a slightly modified architecture for Transformer networks, where the softmax operator $\sigma[\cdot]$ and $\text{ReLU}(\cdot)$ are replaced by the hardmax operator $\sigma_{\text{H}}[\cdot]$ and an activation function $\phi \in \Phi$, respectively. Here, the set of allowed activations Φ consists of all piece-wise linear functions with at most three pieces, where at least one piece is constant. Let $\overline{\mathcal{T}}^{h,m,r}$ denote the function class corresponding to the sequence-to-sequence functions defined by the modified Transformer networks. The following result establishes that the modified Transformer networks in $\overline{\mathcal{T}}^{2,1,1}$ can closely approximate functions in $\overline{\mathcal{F}}_{\text{PE}}(\delta)$.

Proposition 5. For each $\bar{f} \in \bar{\mathcal{F}}_{PE}(\delta)$ and $1 \leq p < \infty$, $\exists \bar{g} \in \bar{\mathcal{T}}^{2,1,1}$ such that $d_p(\bar{f}, \bar{g}) = O(\delta^{d/p})$.

Step 3. Approximate modified Transformers with (original) Transformers. Finally, we show that $\bar{g} \in \bar{\mathcal{T}}^{2,1,1}$ can be approximated by $\mathcal{T}^{2,1,4}$. Let $g \in \mathcal{T}^{2,1,4}$ be such that $d_p(\bar{g}, g) \leq \epsilon/3$.

The following result relies on the connection between the softmax operator and the hardmax operator; and the fact each activation $\phi \in \Phi$ can be approximated by the sum of four ReLU's.

Lemma 6. For each $\bar{g} \in \bar{\mathcal{T}}^{2,1,1}$ and $1 \leq p < \infty$, $\exists g \in \mathcal{T}^{2,1,4}$ such that $d_p(\bar{g}, g) \leq \epsilon/3$.

Theorem 2 now follows from these three steps, because we have

$$d_p(f, g) \leq d_p(f, \bar{f}) + d_p(\bar{f}, \bar{g}) + d_p(\bar{g}, g) \leq 2\epsilon/3 + O(\delta^{d/p}).$$

Choosing $\delta \leq \delta^*$ small enough ensures that $d_p(f, g) \leq \epsilon$. \square

Next, we present the proof of Lemma 4, Proposition 5, and Lemma 6, in Section D, Section E, and Section F, respectively.

D Proof of Lemma 4

Since $f : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ is a continuous function with compact support, the function is uniformly continuous. Since continuity is defined using entry-wise ℓ_p norm, and entry-wise ℓ_p norm is equivalent to entry-wise ℓ_∞ norm when the number of entries are finite, uniform continuity implies that

$$\forall \epsilon > 0, \exists \delta > 0 \text{ such that } \forall \mathbf{X}, \mathbf{Y}, \|\mathbf{X} - \mathbf{Y}\|_\infty < \delta \implies \|f(\mathbf{X}) - f(\mathbf{Y})\|_p < \epsilon.$$

This means that given any $\epsilon/3 > 0$, we have such a $\delta > 0$. Using this δ , we can create a grid \mathbb{G}_δ and corresponding cubes \mathbb{S}_L , as described in the main text. For any $L \in \mathbb{G}_\delta$, we define $C_L \in \mathbb{S}_L$ to be the center point of the cube \mathbb{S}_L . Then, we can define a piece-wise constant approximation $\bar{f}(\mathbf{X}) = \sum_{L \in \mathbb{G}_\delta} f(C_L) \mathbb{1}\{\mathbf{X} \in \mathbb{S}_L\}$. Note that, for any $\mathbf{X} \in \mathbb{S}_L$, we have $\|\mathbf{X} - C_L\|_\infty < \delta$, so by uniform continuity, we have $\|f(\mathbf{X}) - \bar{f}(\mathbf{X})\|_p = \|f(\mathbf{X}) - f(C_L)\|_p < \epsilon/3$. This proves that $d_p(f, \bar{f}) < \epsilon/3$.

As for permutation equivariance, since f is permutation equivariant, we have $f(C_L P) = f(C_L) P$ for any permutation matrix P . For any $\mathbf{X} \in \mathbb{S}_L$, we have $\mathbf{X} P \in \mathbb{S}_{LP}$, so

$$\bar{f}(\mathbf{X} P) = f(C_{LP}) = f(C_L P) = f(C_L) P = \bar{f}(\mathbf{X}) P.$$

Thus, the approximation \bar{f} is also permutation equivariant. This proves the lemma. \square

E Proof of Proposition 5

E.1 Proof sketch: Demystifying Transformers

As mentioned earlier, the heavy parameter sharing in Transformers makes the goal of universally approximating sequence-to-sequence functions seemingly difficult. Both the self-attention and the feed-forward layer weights inside a Transformer block are fixed across n tokens. In this section, we show that Transformers are able to overcome this architectural constraint, and compute contextual mappings of the entire input sequence just based on the pair-wise interactions. The token-wise feedforward layers then transform these contextual mappings to the desired output sequence.

We highlight these inner workings of Transformers en route to proving Proposition 5. We want to show that given a piece-wise constant function $\bar{f} \in \bar{\mathcal{F}}_{PE}(\delta)$, there exists a modified Transformer network $\bar{g} \in \bar{\mathcal{T}}^{2,1,1}$ that closely approximates \bar{f} . We achieve this goal by establishing the following three claims.

1. Given an input $\mathbf{X} \in \mathbb{R}^{d \times n}$, a series of feed-forward layers in the modified Transformer network can quantize \mathbf{X} to an element L on the extended grid $\mathbb{G}_\delta^+ := \{-\delta^{-nd}, 0, \delta, \dots, 1 - \delta\}^{d \times n}$.
2. Next, a series of self-attention layers in the modified Transformer network can take the input L and implement a *contextual mapping* q such that, for L and L' that are not permutation of each other, all the elements in $q(L)$ and $q(L')$ are distinct.

3. Finally, a series of feed-forward layers in the modified Transformer network can map elements of the contextual embedding $q(\mathbf{L})$ to the desired output value of $\bar{f} \in \bar{\mathcal{F}}_{\text{PE}}$ at the input \mathbf{X} .

Before discussing these three claims in detail, we note that even though a Transformer network stacks self-attention and feed-forward layers in an alternate manner, the skip connections enable these networks to employ a composition of multiple self-attention or feed-forward layers. Furthermore, as alluded earlier, these three steps clearly highlight the different roles that self-attention and feed-forward layers play in realizing the ability to universally approximate sequence-to-sequence functions: 1) self-attention layers compute precise contextual maps; and 2) feed-forward layers then assign the results of these contextual maps to the desired output values.

E.1.1 Quantization by feed-forward layers

Since our objective in Proposition 5 is to approximate the function $\bar{f} \in \bar{\mathcal{F}}_{\text{PE}}(\delta)$, which takes a constant value on the cubes $\mathbb{S}_{\mathbf{L}}$'s, the (modified) Transformer network approximating \bar{f} first quantizes the input \mathbf{X} according to these cubes. In particular, we want each input $\mathbf{X} \in \mathbb{S}_{\mathbf{L}}$ to be mapped to the point \mathbf{L} . The following result shows that a modified Transformer network can indeed implement this quantization map with a composition of multiple feed-forward layers.

Lemma 7. *Consider a scalar quantization map $g_q^{\text{ent}} : \mathbb{R} \rightarrow \{-\delta^{-nd}, 0, \delta, \dots, 1 - \delta\}$:*

$$g_q^{\text{ent}}(t) = \begin{cases} k\delta & \text{if } k\delta \leq t < (k+1)\delta, \quad k = 0, \dots, 1/\delta - 1, \\ -\delta^{-nd} & \text{otherwise.} \end{cases}$$

There exists a function $g_q : \mathbb{R}^{d \times n} \mapsto \mathbb{G}_{\delta}^+$ composed of $\frac{d}{\delta} + d$ token-wise feed-forward layers with $r = 1$ and activations in Φ , which employs the scalar quantization g_q^{ent} to each entry of its input.

As desired, the function g_q maps any $\mathbf{X} \in \mathbb{S}_{\mathbf{L}}$ to \mathbf{L} . Furthermore, if any element of \mathbf{X} is not in $[0, 1]$, the element is mapped to $-\delta^{-nd}$, indicating that \mathbf{X} is outside the compact support of $\bar{f} \in \bar{\mathcal{F}}_{\text{PE}}(\delta)$.

E.1.2 Contextual mapping by self-attention layers

In this subsection, we show that the (modified) Transformer network can compute contextual maps from the output of the quantized map g_q (cf. Section E.1.1) by using a composition of self-attention layers.

Let us consider a setting where we are interested in embedding two sentences: 1) I am happy; and 2) I am Bob. These sentences are fed to a sequence-to-sequence model as

$$\mathbf{X} = [\mathbf{X}_{:,1}, \mathbf{X}_{:,2}, \mathbf{X}_{:,3}] = [\mathbf{v}_I, \mathbf{v}_{\text{am}}, \mathbf{v}_{\text{happy}}] \quad \text{and} \quad \tilde{\mathbf{X}} = [\tilde{\mathbf{X}}_{:,1}, \tilde{\mathbf{X}}_{:,2}, \tilde{\mathbf{X}}_{:,3}] = [\mathbf{v}_I, \mathbf{v}_{\text{am}}, \mathbf{v}_{\text{Bob}}],$$

where $\mathbf{v}_I, \mathbf{v}_{\text{am}}, \mathbf{v}_{\text{happy}}$, and \mathbf{v}_{Bob} denote d -dimensional embedding for the tokens ‘I’, ‘am’, ‘happy’, and ‘Bob’, respectively. Since the word ‘I’ occurs in different contexts in these sentences, the contextual mapping of this word should be different. We formally define this requirement below.

Definition E.1 (Contextual mapping). *Consider a finite set $\mathbb{L} \subset \mathbb{R}^{d \times n}$. A map $q : \mathbb{L} \rightarrow \mathbb{R}^{1 \times n}$ defines a contextual mapping if the map satisfies the following:*

1. For any $\mathbf{L} \in \mathbb{L}$, the n entries in $q(\mathbf{L})$ are all distinct.
2. For any $\mathbf{L}, \mathbf{L}' \in \mathbb{L}$, with $\mathbf{L} \neq \mathbf{L}'$, all entries of $q(\mathbf{L})$ and $q(\mathbf{L}')$ are distinct.

In other words, a contextual mapping maps each token (column) of $\mathbf{L} \in \mathbb{L}$ to a unique value which depends on the entire \mathbf{L} ; as a result, capturing the precise context of \mathbf{L} . This allows the subsequent token-wise function (e.g., defined by the feed-forward layers in case of Transformer networks) to realize the outputs of any arbitrary sequence-to-sequence functions.

At first thought, we can consider getting a contextual mapping by simply averaging all the tokens, because this can capture the one-word difference (e.g., “happy” vs. “Bob”) in two different contexts. However, if there are multiple words that are different, it is not guaranteed that the average will be different. Indeed, requiring unique mappings for all the tokens for any change in any number of tokens, is a steep requirement.

While the self-attention layer does consider *pair-wise* interactions among different input tokens, it is not clear if this weak form of pair-wise interaction with shared projection weights is sufficient to

extract the underlying context. The following result shows that the (modified) Transformer networks can implement a *permutation equivariant* contextual mapping over almost all elements of \mathbb{G}_δ .

Lemma 8. *Consider the following subset of $\mathbb{G}_\delta = \{0, \delta, \dots, 1 - \delta\}^{d \times n}$:*

$$\tilde{\mathbb{G}}_\delta := \{\mathbf{L} \in \mathbb{G}_\delta \mid \mathbf{L}_{:,i} \neq \mathbf{L}_{:,j} \text{ for all } i \neq j\}.$$

Assume that $n \geq 2$ and $\delta^{-1} \geq 2$. Then, there exist a function $g_c : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ composed of $\delta^{-d} + 1$ self-attention layers ($h = 2, m = 1$) that employ the σ_H operator, a vector $\mathbf{u} \in \mathbb{R}^d$, constants $t_l, t_r \in \mathbb{R}$ ($0 < t_l < t_r$), such that $q(\mathbf{L}) := \mathbf{u}^T g_c(\mathbf{L})$ satisfies the following properties:

1. *For any $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$, the entries of $q(\mathbf{L})$ are all distinct.*
2. *For any $\mathbf{L}, \mathbf{L}' \in \tilde{\mathbb{G}}_\delta$ such that \mathbf{L} is not a permutation of \mathbf{L}' , all entries of $q(\mathbf{L}), q(\mathbf{L}')$ are distinct.*
3. *For any $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$, all the entries of $q(\mathbf{L})$ are in $[t_l, t_r]$.*
4. *For any $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$, all the entries of $q(\mathbf{L})$ are outside $[t_l, t_r]$.*

At this point, a few remarks about the result in Lemma 8 are in order. First, since the Transformer networks are bound to implement permutation invariant maps, we require the Property 8.2 to hold for the pair of sequences that cannot be mapped to each other via permutation of columns. Furthermore, the self-attention layers implement the desirable contextual map for only $\tilde{\mathbb{G}}_\delta \subseteq \mathbb{G}_\delta$, where all columns of \mathbf{L} are distinct. Note that for small δ , $\mathbb{G}_\delta \setminus \tilde{\mathbb{G}}_\delta$ constitutes a negligible fraction of \mathbb{G}_δ because $|\mathbb{G}_\delta \setminus \tilde{\mathbb{G}}_\delta| = O(\delta^d |\mathbb{G}_\delta|)$. The function q in Lemma 8 maps the elements of $\mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$ outside $[t_l, t_r]$ —the interval where the outputs of the contextual mapping for $\tilde{\mathbb{G}}_\delta$ reside.

E.1.3 Function value mapping by feed-forward layers

This brings us to the final step, which demonstrates the key utility of the feed-forward layers. After the contextual mapping performed by self-attention layers, each token captures the entire context available in the input sequence. The following result shows that token-wise application of a composition of feed-forward layers can map these tokens to the desired output values as required by an arbitrary function \bar{f} .

Lemma 9. *Let $g_c : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ be the function from Lemma 8. Then, there exists a function $g_v : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{d \times n}$ composed of $O(n(\frac{1}{\delta})^{dn}/n!)$ token-wise feed-forward layers ($r = 1$) with activations in Φ such that g_v is defined by a token-wise function $g_v^{\text{tkn}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ on each column,*

$$g_v(\mathbf{X}) = [g_v^{\text{tkn}}(\mathbf{X}_{:,1}) \quad \dots \quad g_v^{\text{tkn}}(\mathbf{X}_{:,n})],$$

where for all $j \in \{1, \dots, n\}$,

$$g_v^{\text{tkn}}(g_c(\mathbf{L})_{:,j}) = \begin{cases} (\mathbf{A}_L)_{:,j} & \text{if } \mathbf{L} \in \tilde{\mathbb{G}}_\delta, \\ \mathbf{0}_d & \text{if } \mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta. \end{cases}$$

E.1.4 Putting all these together

With the results established in Lemmas 7, 8, and 9, we are now equipped with all the tools necessary to complete the proof of Proposition 5. Let us recall the functions g_q, g_c , and g_v from Lemma 7, 8, and 9, respectively. We now show that the (modified) Transformer network $\bar{g} = g_v \circ g_c \circ g_q$ approximates the underlying peicewise constant function $\bar{f} \in \mathcal{F}_{\text{PE}}$ over all points in its support except for a set of of measure $O(\delta^d)$.

Consider a point $\mathbf{X} \in \mathbb{S}_L \subset [0, 1]^{d \times n}$, where $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$. By Lemma 7, we have that $g_q(\mathbf{X}) = \mathbf{L}$. Thus, it follows from Lemmas 8 and 9 that

$$g_v \circ g_c \circ g_q(\mathbf{X}) = g_v \circ g_c(\mathbf{L}) = [g_v^{\text{tkn}}(g_c(\mathbf{L})_{:,1}) \quad g_v^{\text{tkn}}(g_c(\mathbf{L})_{:,2}) \quad \dots \quad g_v^{\text{tkn}}(g_c(\mathbf{L})_{:,n})] = \mathbf{A}_L.$$

On the other hand, any point $\mathbf{X} \in \bigcup_{\mathbf{L} \in \mathbb{G}_\delta \setminus \tilde{\mathbb{G}}_\delta} \mathbb{S}_L \cup (\mathbb{R}^{d \times n} \setminus [0, 1]^{d \times n})$ is mapped by g_q to $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$; as a result, we get $g_v \circ g_c \circ g_q(\mathbf{X}) = g_v \circ g_c(\mathbf{L}) = \mathbf{0}$.

Therefore, we have $\bar{g}(\mathbf{X}) = g_v \circ g_c \circ g_q(\mathbf{X}) = \mathbf{A}_L = \bar{f}(\mathbf{X})$ for $\mathbf{X} \in \bigcup_{\mathbf{L} \in \tilde{\mathbb{G}}_\delta} \mathbb{S}_L$, and $\mathbf{0}$ everywhere else. Recall that \bar{f} has its compact support in $[0, 1]^d$, thus bounded; i.e., there exists $B \geq 0$ such

that $\|\bar{f}(\mathbf{X})\|_p \leq B$. The modified Transformer network \bar{g} takes the same value as \bar{f} on all points in $[0, 1]^d$ except for a set $\bigcup_{L \in \mathbb{G}_\delta \setminus \bar{\mathbb{G}}_\delta} \mathbb{S}_L$ that has measure $O(\delta^d)$. This implies that $d_p(\bar{f}, \bar{g}) \leq (B^p \delta^d)^{1/p} = O(\delta^{d/p})$.

E.2 Proof of Lemma 7

The proof strategy is simple; using $\frac{1}{\delta} + 1$ token-wise feed-forward layers, we implement the quantization function g_q^{ent} that works on the first row of the input. Then stack another $\frac{1}{\delta} + 1$ layers that quantizes the second row, and so on.

Given input \mathbf{X} , we first start by clipping $\mathbf{X}_{1,:}$ in the set $(-\infty, 0) \cup [1, +\infty)$ and mapping the intervals to $-\delta^{-nd}$. This can be done by the following layer:

$$\mathbf{Z} \mapsto \mathbf{Z} + \mathbf{e}^{(1)} \phi((\mathbf{e}^{(1)})^T \mathbf{Z}), \quad \phi(t) = \begin{cases} -t - \delta^{-nd} & \text{if } t < 0 \text{ or } t \geq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Next, add $1/\delta$ layers of the following form, for $k = 0, \delta, \dots, 1 - \delta$.

$$\mathbf{Z} \mapsto \mathbf{Z} + \mathbf{e}^{(1)} \phi((\mathbf{e}^{(1)})^T \mathbf{Z} - k\delta \mathbf{1}_n^T), \quad \phi(t) = \begin{cases} 0 & t < 0 \text{ or } t \geq \delta \\ -t & 0 \leq t < \delta. \end{cases}$$

Each layer quantizes $\mathbf{X}_{1,:}$ in $[k\delta, k\delta + \delta)$ to $k\delta$, without modifying other intervals.

Note that both ϕ 's used in this construction are piecewise linear functions with three pieces, and at least one of them are constant. Thus, both ϕ 's are in Φ . We can repeat the same thing for the other rows, and at the end we will get a map from $\mathbb{R}^{d \times n}$ to \mathbb{G}_δ^+ .

E.3 Proof of Lemma 8

Selective shift operation. Before starting the proof, we first describe the key component of our proof, which we refer to the *selective shift operation*. Consider the following function, which can be expressed with a multiplicative attention head, with head size $m = 1$ and hardmax σ_H :

$$\psi(\mathbf{Z}; b_Q) = \mathbf{e}^{(1)} \mathbf{u}^T \mathbf{Z} \sigma_H[(\mathbf{u}^T \mathbf{Z})^T (\mathbf{u}^T \mathbf{Z} - b_Q \mathbf{1}_n^T)]$$

where $\mathbf{u} \in \mathbb{R}^d$ is a vector that we will choose later, and $\mathbf{e}^{(1)} = (1, 0, 0, \dots, 0) \in \mathbb{R}^d$ is the standard basis vector.

To see what this function computes, first consider the j -th column of the attention score matrix: $(\mathbf{u}^T \mathbf{Z})^T (\mathbf{u}^T \mathbf{Z}_{:,j} - b_Q)$. Note that, if $\mathbf{u}^T \mathbf{Z}_{:,j} > b_Q$, σ_H will calculate $\arg \max$ of $\mathbf{u}^T \mathbf{Z}$, whereas if $\mathbf{u}^T \mathbf{Z}_{:,j} < b_Q$, it will calculate $\arg \min$. Therefore, the $(1, j)$ -th entry of $\psi(\mathbf{Z}; b_Q) \in \mathbb{R}^{d \times n}$ can be written as

$$\psi(\mathbf{Z}; b_Q)_{1,j} = \mathbf{u}^T \mathbf{Z} \sigma_H[(\mathbf{u}^T \mathbf{Z})^T (\mathbf{u}^T \mathbf{Z}_{:,j} - b_Q)] = \begin{cases} \max_k \mathbf{u}^T \mathbf{Z}_{:,k} & \text{if } \mathbf{u}^T \mathbf{Z}_{:,j} > b_Q, \\ \min_k \mathbf{u}^T \mathbf{Z}_{:,k} & \text{if } \mathbf{u}^T \mathbf{Z}_{:,j} < b_Q, \end{cases}$$

for $j \in [n]$. Note that due to $\mathbf{e}^{(1)}$, all rows of $\psi(\mathbf{Z}; b_Q)$ except the first row are zero. From this observation, one can define a function parametrized by b_Q and b'_Q , where $b_Q < b'_Q$, which consists of two attention heads:

$$\begin{aligned} \Psi(\mathbf{Z}; b_Q, b'_Q) &:= \psi(\mathbf{Z}; b_Q) - \psi(\mathbf{Z}; b'_Q), \\ \Psi(\mathbf{Z}; b_Q, b'_Q)_{1,j} &= \begin{cases} \max_k \mathbf{u}^T \mathbf{Z}_{:,k} - \min_k \mathbf{u}^T \mathbf{Z}_{:,k} & \text{if } b_Q < \mathbf{u}^T \mathbf{Z}_{:,j} < b'_Q, \\ 0 & \text{if } \mathbf{u}^T \mathbf{Z}_{:,j} < b_Q \text{ or } \mathbf{u}^T \mathbf{Z}_{:,j} > b'_Q. \end{cases} \end{aligned}$$

What this means is that, if we define an attention layer of the form $\mathbf{Z} \mapsto \mathbf{Z} + \Psi(\mathbf{Z}; b_Q, b'_Q)$, then any column $\mathbf{Z}_{:,j}$ satisfying $\mathbf{u}^T \mathbf{Z}_{:,j} \in (b_Q, b'_Q)$ is shifted up in its first coordinate $\mathbf{Z}_{1,j}$ by $\max_k \mathbf{u}^T \mathbf{Z}_{:,k} - \min_k \mathbf{u}^T \mathbf{Z}_{:,k}$, **while all the other coordinates stay untouched**. We call this the selective shift operation, because we can choose b_Q and b'_Q to selectively shift certain entries of the input.

Bijjective column id mapping. Recall that the input to this step is from the range of g_q (Lemma 7), which is $\mathbb{G}_\delta^+ = \{-\delta^{-nd}, 0, \delta, \dots, 1 - \delta\}^{d \times n}$. Now consider $\mathbf{L} \in \mathbb{G}_\delta^+$ and $\mathbf{u} = (1, \delta^{-1}, \delta^{-2}, \dots, \delta^{-d+1})$.

For any $j \in [n]$, it is easy to check two following facts:

1. If $\mathbf{L}_{i,j} \neq -\delta^{-nd}$ for all $i \in [d]$, i.e., $\mathbf{L}_{:,j} \in \{0, \delta, \dots, 1 - \delta\}^d$, then $\mathbf{u}^T \mathbf{L}_{:,j} \in [0 : \delta : \delta^{-d+1} - \delta]$, and the map $\mathbf{L}_{:,j} \mapsto \mathbf{u}^T \mathbf{L}_{:,j}$ from $\{0, \delta, \dots, 1 - \delta\}^d$ to $[0 : \delta : \delta^{-d+1} - \delta]$ is a bijection.
2. If there exists $i \in [d]$ such that $\mathbf{L}_{i,j} = -\delta^{-nd}$, then $\mathbf{u}^T \mathbf{L}_{:,j} \leq -\delta^{-nd} + \delta^{-d+1} - 1 < 0$.

Therefore, one can say that $\mathbf{u}^T \mathbf{L}_{:,j}$ gives the ‘‘column id’’ for each possible value of $\mathbf{L}_{:,j} \in \{0, \delta, \dots, 1 - \delta\}^d$.

The rough idea of the construction is to apply the selective shift operation to each column id, by setting \mathbf{u} in the definition of $\Psi(\cdot)$ to be $(1, \delta^{-1}, \delta^{-2}, \dots, \delta^{-d+1})$ and choosing $b_Q = l - \delta/2$ and $b'_Q = l + \delta/2$ for each $l \in [0 : \delta : \delta^{-d+1} - \delta]$. More concretely, we stack $(1/\delta)^d$ attention layers, with attention parts $\delta^{-d} \Psi(\cdot; l - \delta/2, l + \delta/2)$ for each $l \in [0 : \delta : \delta^{-d+1} - \delta]$, in increasing order of l . After that, we add an extra single-head attention layer with attention part $\delta^{-(n+1)d} \psi(\cdot; 0)$.

We now divide possible input values $\mathbf{L} \in \mathbb{G}_\delta^+$ into three disjoint categories, and show how these layers change the input values at the end of all the layers. Recall the hierarchy $\tilde{\mathbb{G}}_\delta \subset \mathbb{G}_\delta \subset \mathbb{G}_\delta^+$. The categories are defined as follows:

1. $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$. All entries are between 0 and $1 - \delta$, and all columns are unique.
2. $\mathbf{L} \in \mathbb{G}_\delta \setminus \tilde{\mathbb{G}}_\delta$. All entries are between 0 and $1 - \delta$, but there are duplicate columns.
3. $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \mathbb{G}_\delta$. The point has at least one entry that equals to $-\delta^{-nd}$.

E.3.1 Category 1

In Category 1, we have $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$. Let $l_j := \mathbf{u}^T \mathbf{L}_{:,j}$. Due to permutation equivariance, we can assume without loss of generality that l_j 's are in increasing order: $l_1 < l_2 < \dots < l_n$. The first $(1/\delta)^d$ layers sweep the set $[0 : \delta : \delta^{-d+1} - \delta]$ and apply selective shift operation on each element in the set. This means that selective shift operation will be applied to l_1 first, then l_2 , and then l_3 , and so on, regardless of the specific values of l_j 's.

First shift operation. In the first selective shift operation, the $(1, 1)$ -th entry of \mathbf{L} ($L_{1,1}$) is shifted by the operation, while the other entries are left untouched. The updated value $\tilde{L}_{1,1}$ is

$$\tilde{L}_{1,1} = L_{1,1} + \delta^{-d}(\max_k \mathbf{u}^T \mathbf{L}_{:,k} - \min_k \mathbf{u}^T \mathbf{L}_{:,k}) = L_{1,1} + \delta^{-d}(l_n - l_1).$$

Therefore, after the operation, the output of the layer is $[\tilde{\mathbf{L}}_{:,1} \quad \mathbf{L}_{:,2} \quad \dots \quad \mathbf{L}_{:,n}]$, and the new value of the first column $\tilde{\mathbf{L}}_{:,1}$ results in

$$\mathbf{u}^T \tilde{\mathbf{L}}_{:,1} = \tilde{L}_{1,1} + \sum_{i=2}^d \delta^{-i+1} L_{i,1} = L_{1,1} + \delta^{-d}(l_n - l_1) + \sum_{i=2}^d \delta^{-i+1} L_{i,1} = l_1 + \delta^{-d}(l_n - l_1).$$

Let us denote the updated ‘‘column id’’ $\mathbf{u}^T \tilde{\mathbf{L}}_{:,1}$ as \tilde{l}_1 . We can show that $l_n < \tilde{l}_1$, because

$$\tilde{l}_1 := l_1 + \delta^{-d}(l_n - l_1) \geq 0 + \delta^{-d} \cdot \delta = \delta^{-d+1} > l_n.$$

Therefore, after updating,

$$\max \mathbf{u}^T [\tilde{\mathbf{L}}_{:,1} \quad \mathbf{L}_{:,2} \quad \dots \quad \mathbf{L}_{:,n}] = \max\{\tilde{l}_1, l_2, \dots, l_n\} = \tilde{l}_1,$$

and the new minimum is l_2 .

Second shift operation. The second selective shift operation is applied to l_2 , by which only one entry $L_{1,2}$ will be shifted. The updated value $\tilde{L}_{1,2}$ is

$$\tilde{L}_{1,2} = L_{1,2} + \delta^{-d}(\tilde{l}_1 - l_2) = L_{1,2} + \delta^{-d}(l_1 - l_2) + \delta^{-2d}(l_n - l_1).$$

After updating, the new inner product of \mathbf{u} and $\tilde{\mathbf{L}}_{:,2}$ results in

$$\tilde{l}_2 := \mathbf{u}^T \tilde{\mathbf{L}}_{:,2} = l_2 + \delta^{-d}(l_1 - l_2) + \delta^{-2d}(l_n - l_1).$$

We can show that $\tilde{l}_1 < \tilde{l}_2$, because

$$\begin{aligned} l_1 + \delta^{-d}(l_n - l_1) &< l_2 + \delta^{-d}(l_1 - l_2) + \delta^{-2d}(l_n - l_1) \\ \Leftrightarrow (\delta^{-d} - 1)(l_2 - l_1) &< \delta^{-d}(\delta^{-d} - 1)(l_n - l_1), \end{aligned}$$

and the last inequality is true because $\delta^{-d} > 1$ and $l_n > l_2$. Since we have $\tilde{l}_1 < \tilde{l}_2$, and the new maximum in $\mathbf{u}^T [\tilde{\mathbf{L}}_{:,1} \quad \tilde{\mathbf{L}}_{:,2} \quad \mathbf{L}_{:,3} \quad \dots \quad \mathbf{L}_{:,n}]$ is now \tilde{l}_2 , and the new minimum is l_3 .

Repeating the process. More generally, we can repeat this process, and show that the j -th shift operation shifts $L_{1,j}$ by $\delta^{-d}(\tilde{l}_{j-1} - l_j)$, and results in the new column id

$$\tilde{l}_j := \mathbf{u}^T \tilde{\mathbf{L}}_{:,j} = l_j + \sum_{k=1}^{j-1} \delta^{-kd}(l_{j-k} - l_{j-k+1}) + \delta^{-jd}(l_n - l_1).$$

In the general case, $\tilde{l}_{j-1} < \tilde{l}_j$ holds $j = [2 : n]$, because

$$\begin{aligned} \tilde{l}_{j-1} &= l_{j-1} + \sum_{k=2}^{j-1} \delta^{-kd+d}(l_{j-k} - l_{j-k+1}) + \delta^{-(j-1)d}(l_n - l_1) \\ &< \tilde{l}_j = l_j + \sum_{k=1}^{j-1} \delta^{-kd}(l_{j-k} - l_{j-k+1}) + \delta^{-jd}(l_n - l_1) \\ \Leftrightarrow \sum_{k=1}^{j-1} \delta^{-kd+d}(\delta^{-d} - 1)(l_{j-k+1} - l_{j-k}) &< \delta^{-(j-1)d}(\delta^{-d} - 1)(l_n - l_1), \end{aligned}$$

and the last inequality holds because

$$\delta^{-(j-1)d}(l_n - l_1) > \delta^{-(j-1)d} \sum_{k=1}^{j-1} (l_{j-k+1} - l_{j-k}) > \sum_{k=1}^{j-1} \delta^{-kd+d}(l_{j-k+1} - l_{j-k}).$$

Therefore, after the j -th selective shift operation, \tilde{l}_j is the new maximum among $\{\tilde{l}_1, \dots, \tilde{l}_j, l_{j+1}, \dots, l_n\}$ and l_{j+1} is the new minimum, which makes us possible to continue the process until the n -th operation.

After n shift operations. As a result, after the whole sweep from 0 to $\delta^{-d+1} - \delta$ by the first $(1/\delta)^d$ layers, a total of n shift operations are applied, and the input \mathbf{L} is mapped to a new point $\tilde{\mathbf{L}}$, where $\mathbf{u}^T \tilde{\mathbf{L}} = [\tilde{l}_1 \quad \tilde{l}_2 \quad \dots \quad \tilde{l}_n]$ and $\tilde{l}_1 < \tilde{l}_2 < \dots < \tilde{l}_n$.

We can now prove the following technical lemma, whose proof is deferred to Appendix E.3.4:

Lemma 10. *After n shift operations, $\tilde{l}_n = \mathbf{u}^T \tilde{\mathbf{L}}_{:,n}$ satisfies the following bounds:*

$$\delta^{-(n-1)d+1}(\delta^{-d} - 1) \leq \tilde{l}_n \leq \delta^{-nd+1}(\delta^{-d} - 1) - \delta(\delta^{-d} - 1)^2.$$

Also, the map from $[l_1 \quad l_2 \quad \dots \quad l_n] \in [0 : \delta : \delta^{-d+1} - \delta]$ (where $l_1 < l_2 < \dots < l_n$) to \tilde{l}_n is one-to-one.

Global shifting by the last layer. As mentioned earlier, after this sweep, there is another attention layer with attention part $\delta^{-(n+1)d}\psi(\cdot; 0)$. Since $0 < \tilde{l}_1 < \dots < \tilde{l}_n$, what it does to $\tilde{\mathbf{L}}$ is that it adds $\delta^{-(n+1)d} \max_k \mathbf{u}^T \tilde{\mathbf{L}}_{:,k} = \delta^{-(n+1)d} \tilde{l}_n$ to each entry in the first row of $\tilde{\mathbf{L}}$. The output of this layer is defined to be the function $g_c(\mathbf{L})$.

Now, in summary, for any $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$, $i \in [d]$, and $j \in [n]$, we have

$$g_c(\mathbf{L})_{i,j} = \begin{cases} L_{1,j} + \sum_{k=1}^{j-1} \delta^{-kd}(l_{j-k} - l_{j-k+1}) + \delta^{-jd}(l_n - l_1) + \delta^{-(n+1)d} \tilde{l}_n & \text{if } i = 1, \\ L_{i,j} & \text{if } i \in [2, d], \end{cases}$$

and for any $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$ and $j \in [n]$,

$$\mathbf{u}^T g_c(\mathbf{L})_{:,j} = \tilde{l}_j + \delta^{-(n+1)d} \tilde{l}_n.$$

Checking Properties 8.1 and 8.2. Given this result so far, it is now left to check if the constructed network is really a permutation equivariant contextual mapping, i.e., if it satisfies Properties 8.1 and 8.2 in Lemma 8.

First, for any $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$, Property 8.1 holds because we already know $\tilde{l}_1 < \tilde{l}_2 < \dots < \tilde{l}_n$, so they are all distinct. As for Property 8.2, note that the upper bound on \tilde{l}_n from Lemma 10 also holds for other \tilde{l}_j 's, so

$$\mathbf{u}^T g_c(\mathbf{L})_{:,j} \in [\delta^{-(n+1)d} \tilde{l}_n, \delta^{-(n+1)d} \tilde{l}_n + \delta^{-(n+1)d+1}],$$

for all $j \in [n]$. Now, from Lemma 10, two $\mathbf{L}, \mathbf{L}' \in \tilde{\mathbb{G}}_\delta$ (that are not permutations of each other) map to different \tilde{l}_n and \tilde{l}'_n , and they differ at least by δ . This means that two intervals $[\delta^{-(n+1)d} \tilde{l}_n, \delta^{-(n+1)d} \tilde{l}_n + \delta^{-(n+1)d+1}]$ and $[\delta^{-(n+1)d} \tilde{l}'_n, \delta^{-(n+1)d} \tilde{l}'_n + \delta^{-(n+1)d+1}]$ are guaranteed to be disjoint, so the entries of $\mathbf{u}^T g_c(\mathbf{L})$ and $\mathbf{u}^T g_c(\mathbf{L}')$ are all distinct. This proves Property 8.2.

Therefore, we finished showing that the map $g_c(\cdot)$ we constructed using $(1/\delta)^d + 1$ attention layers implements a permutation equivariant contextual mapping on $\tilde{\mathbb{G}}_\delta$.

Checking Property 8.3. It is now left to check if the map g_c satisfies the other properties. At this point, we can check Property 8.3. From $\mathbf{u}^T g_c(\mathbf{L})_{:,j} \in [\delta^{-(n+1)d} \tilde{l}_n, \delta^{-(n+1)d} \tilde{l}_n + \delta^{-(n+1)d+1}]$ and Lemma 10, we can show that for any $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$, we have

$$\begin{aligned} \delta^{-2nd+1}(\delta^{-d} - 1) &\leq \mathbf{u}^T g_c(\mathbf{L})_{:,j} < \delta^{-(n+1)d}(\delta^{-nd+1}(\delta^{-d} - 1) - \delta(\delta^{-d} - 1)^2) + \delta^{-(n+1)d+1} \\ &\leq \delta^{-(2n+1)d+1}(\delta^{-d} - 1), \end{aligned}$$

where we used $\delta^{-1} \geq 2$. This proves that all $\mathbf{u}^T g_c(\mathbf{L})_{:,j}$ are between $t_l = \delta^{-2nd+1}(\delta^{-d} - 1)$ and $t_r = \delta^{-(2n+1)d+1}(\delta^{-d} - 1)$. For the remaining input points $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$, we will check that $\mathbf{u}^T g_c(\mathbf{L})_{:,j}$ is outside the interval $[t_l, t_r]$ (Property 8.4).

E.3.2 Category 2

In Category 2, we have $\mathbf{L} \in \mathbb{G}_\delta \setminus \tilde{\mathbb{G}}_\delta$. Here, all entries are between 0 and $1 - \delta$, but there are duplicate columns. Again, let $l_j := \mathbf{u}^T \mathbf{L}_{:,j}$, and assume without loss of generality that $l_1 \leq l_2 \leq \dots \leq l_n$. For the input \mathbf{L} in Category 2, there exist some $j, j' \in [n]$, $j \neq j'$, such that $l_j = l_{j'}$. This means that when the input passes through the attention layer $\delta^{-d}\Psi(\cdot; l_j - \delta/2, l_j + \delta/2)$, the selective shift operation for l_j is applied to both j -th and j' -th columns; the two columns are coupled together. More generally, suppose we have $n' < n$ distinct columns.

If $n' = 1$. In the extreme case of $n' = 1$, we have $\max_j l_j = \min_j l_j$, so the selective shift operation applied at l_j does not shift the entry at all; therefore, at the end of the first $(1/\delta)^d$ attention layers, $\tilde{\mathbf{L}} = \mathbf{L}$.

If $1 < n' \leq n - 1$. When $1 < n' \leq n - 1$, let the n' distinct values of l_j 's be $l'_1, \dots, l'_{n'}$. The shift operation is applied n' times, to $l'_1, \dots, l'_{n'}$, and shifts one or more entries at a time. After the first $(1/\delta)^d$ layers, the output $\tilde{\mathbf{L}}$ has n' distinct $\tilde{l}_j = \mathbf{u}^T \tilde{\mathbf{L}}_{:,j}$, $0 \leq \tilde{l}_1 \leq \tilde{l}_2 \leq \dots \leq \tilde{l}_n$, whose distinct values are the same as the numbers we get when we apply shift operations to a length- n' sequence $[l'_1 \dots l'_{n'}]$. Then, applying the same calculations from Category 1 shows that

$$\tilde{l}_n = \mathbf{u}^T \tilde{\mathbf{L}}_{:,n} = l'_{n'} + \sum_{k=1}^{n'-1} \delta^{-kd} (l'_{n'-k} - l'_{n'-k+1}) + \delta^{-n'd} (l'_{n'} - l'_1),$$

and it follows from the upper bound in Lemma 10 that

$$\tilde{l}_n \leq \delta^{-n'd+1}(\delta^{-d} - 1) - \delta(\delta^{-d} - 1)^2 < \delta^{-(n-1)d+1}(\delta^{-d} - 1).$$

Note that the RHS matches the lower bound in Lemma 10. This implies that the value of \tilde{l}_n calculated from the input $\mathbf{L} \in \mathbb{G}_\delta \setminus \tilde{\mathbb{G}}_\delta$ (Category 2) is always strictly less (by at least δ) than that calculated from $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$ (Category 1).

Checking Property 8.4. After the global shifting by the last layer with attention part $\delta^{-(n+1)d}\psi(\cdot; 0)$, we get the output $g_c(\mathbf{L})$ which satisfies

$$\begin{aligned} \mathbf{u}^T g_c(\mathbf{L})_{:,j} &= \tilde{l}_j + \delta^{-(n+1)d} \tilde{l}_n \leq (\delta^{-(n+1)d} + 1)(\delta^{-(n-1)d+1}(\delta^{-d} - 1) - \delta(\delta^{-d} - 1)^2) \\ &< \delta^{-2nd+1}(\delta^{-d} - 1) =: t_l. \end{aligned}$$

where the RHS is a lower bound on possible values of $\mathbf{u}^T g_c(\mathbf{L})_{:,j}$ for $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$ (Category 1). This means that the entries of $\mathbf{u}^T g_c(\mathbf{L})$ for Category 2 are outside $[t_l, t_r]$, which satisfies Property 8.4.

E.3.3 Category 3

In Category 3, we have $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \mathbb{G}_\delta$; the point \mathbf{L} has at least one entry that equals to $-\delta^{-nd}$. Let $l_j := \mathbf{u}^T \mathbf{L}_{:,j}$, and recall that whenever a column $\mathbf{L}_{:,j}$ has an entry that equals to $-\delta^{-nd}$, we have $l_j = \mathbf{u}^T \mathbf{L}_{:,j} \leq -\delta^{-nd} + \delta^{-d+1} - 1 < 0$. Assume without loss of generality that $l_1 \leq l_2 \leq \dots \leq l_n$.

Recall that the selective shift operation is applied to each element of $[0 : \delta : \delta^{-d+1} - \delta]$, not to negative values. In case of Category 3, we have $\min_k \mathbf{u}^T \mathbf{L}_{:,k} = l_1 < 0$, and l_1 never gets shifted upwards, so it remains as the minimum for the whole time.

If all l_j 's are negative. In case where all l_j 's are negative, selective shift operation never changes the input \mathbf{L} , so we get $\tilde{\mathbf{L}} = \mathbf{L}$. Since we have $\mathbf{u}^T \tilde{\mathbf{L}} < \mathbf{0}_n^T$ (entry-wise), the last layer with attention part $\delta^{-(n+1)d}\psi(\cdot; 0)$ adds $\delta^{-(n+1)d} \min_k \mathbf{u}^T \tilde{\mathbf{L}}_{:,k} < 0$ to each entry in the first row of $\tilde{\mathbf{L}}$, further pushing it to the negative side. Therefore, the final output $g_c(\mathbf{L})$ satisfies $\mathbf{u}^T g_c(\mathbf{L}) < \mathbf{0}_n^T < t_l \mathbf{1}_n^T$.

If not all l_j 's are negative. Now consider the case where at least one l_j is positive. Let i be the index that satisfies $l_{i-1} < 0 \leq l_i$. Then, selective shift operation does not affect l_1, \dots, l_{i-1} , and then it shifts l_i by

$$\delta^{-d}(\max_k \mathbf{u}^T \mathbf{L}_{:,k} - \min_k \mathbf{u}^T \mathbf{L}_{:,k}) = \delta^{-d}(l_n - l_1) \geq \delta^{-d}(0 + \delta^{-nd} - \delta^{-d+1} + 1) \geq \delta^{-(n+1)d+1},$$

where we used $\delta^{-1} \geq 2$ at the last inequality. The next shift operations shift l_{i+1}, \dots, l_n by even larger amount, so at the end of the first $(1/\delta)^d$ layers, we have $\delta^{-(n+1)d+1} \leq \tilde{l}_i \leq \dots \leq \tilde{l}_n$, while $\tilde{l}_j = l_j < 0$ for $j \in [i - 1]$.

Shifts by the last layer. Here, the last layer with attention part $\delta^{-(n+1)d}\psi(\cdot; 0)$ acts differently for negative and positive \tilde{l}_j 's. For negative \tilde{l}_j 's, it adds $\delta^{-(n+1)d} \min_k \tilde{l}_k = \delta^{-(n+1)d} l_1 < 0$ to $\tilde{l}_1, \dots, \tilde{l}_{i-1}$, pushing them further to the negative side. For positive \tilde{l}_j 's, the layer adds $\delta^{-(n+1)d} \max_k \tilde{l}_k = \delta^{-(n+1)d} \tilde{l}_n \geq \delta^{-(2n+2)d+1}$ to $\tilde{l}_i, \dots, \tilde{l}_n$, so that they are all greater than or equal to $\delta^{-(2n+2)d+1}$. Note that $\delta^{-(2n+2)d+1} > t_r$.

Checking Property 8.4. Therefore, in both cases, we can see that the final output $g_c(\mathbf{L})$ satisfies $\mathbf{u}^T g_c(\mathbf{L})_{:,j} \notin [t_l, t_r]$, for all $j \in [n]$. This completes the verification of Property 8.4.

E.3.4 Proof of Lemma 10

Proof of lower and upper bounds on \tilde{l}_n are straightforward:

$$\begin{aligned} \tilde{l}_n &:= l_n + \sum_{k=1}^{n-1} \delta^{-kd} (l_{n-k} - l_{n-k+1}) + \delta^{-nd} (l_n - l_1) \\ &\geq \delta^{-(n-1)d} \sum_{k=1}^{n-1} (l_{n-k} - l_{n-k+1}) + \delta^{-nd} (l_n - l_1) = (\delta^{-nd} - \delta^{-(n-1)d}) (l_n - l_1) \\ &\geq \delta^{-(n-1)d+1} (\delta^{-d} - 1), \\ \tilde{l}_n &\leq l_n + \delta^{-d} (l_1 - l_n) + \delta^{-nd} (l_n - l_1) \leq \delta^{-d+1} - \delta + (\delta^{-nd} - \delta^{-d}) (\delta^{-d+1} - \delta) \\ &= \delta^{-nd+1} (\delta^{-d} - 1) - \delta (\delta^{-2d} - 2\delta^{-d} + 1) = \delta^{-nd+1} (\delta^{-d} - 1) - \delta (\delta^{-d} - 1)^2. \end{aligned}$$

For one-to-one property of the map, consider $[l_1 \ l_2 \ \dots \ l_n]$ and $[l'_1 \ l'_2 \ \dots \ l'_n]$ with increasing entries, which are mapped to \tilde{l}_n and \tilde{l}'_n , respectively. Suppose $\tilde{l}_n = \tilde{l}'_n$. By definition,

$$\begin{aligned} \tilde{l}_n - \tilde{l}'_n &= (l_n - l'_n) + \delta^{-d} (l_{n-1} - l_n - l'_{n-1} + l'_n) + \delta^{-2d} (l_{n-2} - l_{n-1} - l'_{n-2} + l'_{n-1}) + \dots \\ &\quad + \delta^{-(n-1)d} (l_1 - l_2 - l'_1 + l'_2) + \delta^{-nd} (l_n - l_1 - l'_n + l'_1) = 0. \end{aligned}$$

Now assume for contradiction that $l_n \neq l'_n$. Then, we have $-\delta^{-d+1} + \delta \leq l_n - l'_n \leq \delta^{-d+1} - \delta$. However, the remaining terms have ‘‘coarse resolution’’, and they can never cancel $l_n - l'_n$ and make the sum zero, because for example, $\delta^{-d} (l_{n-1} - l_n - l'_{n-1} + l'_n)$ can only have values $0, \delta^{-d+1}, -\delta^{-d+1}, 2\delta^{-d+1}, -2\delta^{-d+1}, \dots$. Thus, $l_n = l'_n$ must hold and the first term must be zero.

Similarly, assume that $l_{n-1} \neq l'_{n-1}$. Then, the second term is in the interval $[-\delta^{-2d+1} + \delta^{-d+1}, \delta^{-2d+1} - \delta^{-d+1}]$. Again, the remaining terms cannot cancel the second term, hence $l_{n-1} = l'_{n-1}$ must hold. We can proceed this way, and show that $l_j = l'_j$ must hold for all $j \in [n]$, hence proving that the map is one-to-one.

E.4 Proof of Lemma 9

Note that $|\mathbb{G}_\delta^+| = (\frac{1}{\delta} + 1)^{dn}$, so the image of $g_c(\mathbb{G}_\delta^+)$ (from Lemma 8) has finite number of distinct real numbers. Let M be the maximum over all these numbers. By construction of g_c , we know that $M > 0$.

To construct a function g_v^{tkn} that satisfies the statement of the lemma, we first implement the second part: $g_v^{\text{tkn}}(g_c(\mathbf{L})_{:,j}) = \mathbf{0}_d$ if $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$. Note from Lemma 8 that, for any $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$, we have $\mathbf{u}^T g_c(\mathbf{L})_{:,j} \in [t_l, t_r]$ for all j , and for any $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$, $\mathbf{u}^T g_c(\mathbf{L})_{:,j} \notin [t_l, t_r]$ for all j . Using this, we add the following feed-forward layer:

$$\mathbf{Z} \mapsto \mathbf{Z} - (M+1)\mathbf{1}_n \phi(\mathbf{u}^T \mathbf{Z}), \quad \phi(t) = \begin{cases} 0 & \text{if } t \in [t_l, t_r] \\ 1 & \text{if } t \notin [t_l, t_r]. \end{cases}$$

Input to this layer is $g_c(\mathbf{L})$. If $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$, then $\phi(\mathbf{u}^T g_c(\mathbf{L})) = \mathbf{0}_n^T$, so the output stays the same as the input. If $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$, then $\phi(\mathbf{u}^T g_c(\mathbf{L})) = \mathbf{1}_n^T$, so all the entries of the input are shifted by $-M-1$, and become strictly negative.

Recall that by definition of $\tilde{\mathbb{G}}_\delta$, all the entries of $g_c(\mathbf{L})$ for $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$ are nonnegative. So the next thing to do is mapping all strictly negative entries to zero. This can be done in a similar way as Lemma 7. For $i \in [d]$, add the following layer:

$$\mathbf{Z} \mapsto \mathbf{Z} + \mathbf{e}^{(i)} \phi((\mathbf{e}^{(i)})^T \mathbf{Z}), \quad \phi(t) = \begin{cases} -t & \text{if } t < 0 \\ 0 & \text{if } t \geq 0. \end{cases}$$

After these d layers, the output for $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$ is a zero matrix, while the output for $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$ is $g_c(\mathbf{L})$.

Now, it is left to map $g_c(\mathbf{L})$ to $\mathbf{A}_\mathbf{L}$, for $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$. Up to permutation equivariance, each different context \mathbf{L} maps to n unique numbers $\mathbf{u}^T g_c(\mathbf{L})$, which are at least δ apart from each other. The idea of value mapping is to map each unique number to the corresponding output column.

More precisely, choose any $\bar{\mathbf{L}} \in \tilde{\mathbb{G}}_\delta$. For each value of $\mathbf{u}^T g_c(\bar{\mathbf{L}})_{:,j}$, $j \in [n]$, we add one feed-forward layer

$$\mathbf{Z} \mapsto \mathbf{Z} + ((\mathbf{A}_{\bar{\mathbf{L}}})_{:,j} - g_c(\bar{\mathbf{L}})_{:,j}) \phi(\mathbf{u}^T \mathbf{Z} - \mathbf{u}^T g_c(\bar{\mathbf{L}})_{:,j} \mathbf{1}_n^T), \quad \phi(t) = \begin{cases} 0 & t < -\delta/2 \text{ or } t \geq \delta/2, \\ 1 & -\delta/2 \leq t < \delta/2. \end{cases}$$

If the input \mathbf{Z} is a zero matrix, which is the case for $\mathbf{L} \in \mathbb{G}_\delta^+ \setminus \tilde{\mathbb{G}}_\delta$, $\mathbf{u}^T \mathbf{Z} = \mathbf{0}_n^T$. Since t_l is much larger than 0, activation is all zero. Thus, zero input matrix remains the same at the output.

If the input \mathbf{Z} is $g_c(\mathbf{L})$, where $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$ is not a permutation of $\bar{\mathbf{L}}$, then

$$\phi(\mathbf{u}^T g_c(\mathbf{L}) - \mathbf{u}^T g_c(\bar{\mathbf{L}})_{:,j} \mathbf{1}_n^T) = \mathbf{0}_n^T,$$

so $g_c(\mathbf{L})$ is left untouched.

If some other \mathbf{L} is a permutation of $\bar{\mathbf{L}}$, and $\mathbf{L}_{:,i} = \bar{\mathbf{L}}_{:,j}$, then

$$\phi(\mathbf{u}^T g_c(\mathbf{L}) - \mathbf{u}^T g_c(\bar{\mathbf{L}})_{:,j} \mathbf{1}_n^T) = (\mathbf{e}^{(i)})^T,$$

so i -th column of $g_c(\mathbf{L})$ will turn to

$$g_c(\mathbf{L})_{:,i} \mapsto g_c(\mathbf{L})_{:,i} + ((\mathbf{A}_{\bar{\mathbf{L}}})_{:,j} - g_c(\bar{\mathbf{L}})_{:,j}) = g_c(\mathbf{L})_{:,i} + ((\mathbf{A}_\mathbf{L})_{:,i} - g_c(\mathbf{L})_{:,i}) = (\mathbf{A}_\mathbf{L})_{:,i},$$

which is the desired output. In conclusion, this layer maps the column $g_c(\bar{\mathbf{L}})_{:,j}$ to $(\mathbf{A}_{\bar{\mathbf{L}}})_{:,j}$, without affecting any other columns.

As seen above, we need one layer per each unique value of $\mathbf{u}^T g_c(\mathbf{L})_{:,j}$ for each $\mathbf{L} \in \tilde{\mathbb{G}}_\delta$. Note that there are $O(n(1/\delta)^{dn}/n!)$ such numbers, so we can use $O(n(1/\delta)^{dn}/n!)$ layers to finish our construction.

F Proof of Lemma 6

Recall that $T^{h,m,r}$ refers to the class of functions representable with composition of Transformer blocks with h heads of size m in self-attention layers and r hidden nodes in feed-forward layers. The same notation holds for the modified Transformers $\bar{T}^{h,m,r}$.

Note that the softmax operator on a matrix \mathbf{A} can be made arbitrarily close to hardmax by scaling up \mathbf{A} . That is,

$$\sigma[\lambda \mathbf{A}] \rightarrow \sigma_{\text{H}}[\mathbf{A}] \quad \text{as } \lambda \rightarrow \infty.$$

This means that by scaling up parameters inside σ , we can approximate σ_{H} arbitrarily closely. Thus, the modified self-attention layers can be approximated with the original self-attention layers of the same number of heads h and head size m .

Also, any arbitrary (possibly discontinuous) piecewise linear function $\phi \in \Phi$ can be approximated arbitrarily closely by four ReLU's. Note that $\phi \in \Phi$ as at most three pieces, and at least one of the pieces is constant. For example, consider the following function $\phi \in \Phi$:

$$\phi(t) = \begin{cases} b_1 & \text{if } t < c_1, \\ a_2 t + b_2 & \text{if } c_1 \leq t < c_2, \\ a_3 t + b_3 & \text{if } c_2 \leq t. \end{cases}$$

This function can be approximated by four ReLU's, as claimed by the lemma:

$$\begin{aligned} \tilde{\phi}(t) &= b_1 + \frac{a_2 c_1 + b_2 - b_1}{\epsilon} \text{ReLU}(t - c_1 + \epsilon) + \left(a_2 - \frac{a_2 c_1 + b_2 - b_1}{\epsilon} \right) \text{ReLU}(t - c_1) \\ &\quad + \left(\frac{a_3 c_2 + b_3 - a_2(c_2 - \epsilon) - b_2}{\epsilon} - a_2 \right) \text{ReLU}(t - c_2 + \epsilon) \\ &\quad + \left(a_3 - \frac{a_3 c_2 + b_3 - a_2(c_2 - \epsilon) - b_2}{\epsilon} \right) \text{ReLU}(t - c_2) \\ &= \begin{cases} b_1 & \text{if } t < c_1 - \epsilon, \\ \frac{a_2 c_1 + b_2 - b_1}{\epsilon} (t - c_1) + a_2 c_1 + b_2 & \text{if } c_1 - \epsilon \leq t < c_1, \\ a_2 t + b_2 & \text{if } c_1 \leq t < c_2 - \epsilon, \\ \frac{a_3 c_2 + b_3 - a_2(c_2 - \epsilon) - b_2}{\epsilon} (t - c_2) + a_3 c_2 + b_3 & \text{if } c_2 - \epsilon \leq t < c_2, \\ a_3 t + b_3 & \text{if } c_2 \leq t. \end{cases} \end{aligned}$$

Also, as we make $\epsilon \rightarrow 0$, we can approximate ϕ as closely as possible using $\tilde{\phi}$. The cases where the second or third piece is constant can be shown similarly. This means that the modified feed-forward layers (whose activation is $\phi \in \Phi$) with single hidden node can be approximated with the original feed-forward layers (ReLU) with four hidden nodes.

Thus, given any $\bar{g} \in \overline{\mathcal{T}}^{2,1,1}$, there exists a function $g \in \mathcal{T}^{2,1,4}$ arbitrarily close to \bar{g} , by appropriately choosing the parameters to be large enough. This finishes the proof. \square

G Proof of Theorem 3

Proof of Theorem 3 can be done in a similar way as Theorem 2. As in the proof of Theorem 2, there are three parts: Lemma 4, Proposition 5, and Lemma 6. The statement and proof of Lemmas 4 and 6 can be done in almost the same way, this time without permutation equivariance.

For the proof of the second part, which corresponds to Proposition 5, we construct the network in a similar way. Recall that we can assume without loss of generality that $\mathbf{X} \in [0, 1]^{d \times n}$. Choose

$$\mathbf{E} = \begin{bmatrix} 0 & 1 & 2 & \cdots & n-1 \\ 0 & 1 & 2 & \cdots & n-1 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 1 & 2 & \cdots & n-1 \end{bmatrix}.$$

Then, the first column of $\mathbf{X} + \mathbf{E}$ is in $[0, 1]^d$, second is in $[1, 2]^d$, and so on; this means that for all rows, the coordinates are monotonically increasing. So we can use the same technique as the proof of Proposition 5 to divide the input values into cubes, quantize them to \mathbf{L} , apply contextual mapping, and then value mapping. We describe each step in the following.

G.1 Quantization by feed-forward layers

In a similar way as Lemma 7, the goal of this step is to quantize the input in $[0, 1]^d \times [1, 2]^d \times \cdots \times [n-1, n]^d$ to its discrete version:

$$[0 : \delta : 1 - \delta]^d \times [1 : \delta : 2 - \delta]^d \times \cdots \times [n-1 : \delta : n - \delta]^d.$$

This can be done by dn/δ feed-forward layers. We add dn/δ layers of the following form, for $k = 0, \delta, \dots, n - \delta$ and $i = 1, \dots, d$:

$$\mathbf{Z} \mapsto \mathbf{Z} + e^{(i)} \phi((e^{(i)})^T \mathbf{Z} - k\delta \mathbf{1}_n^T), \quad \phi(t) = \begin{cases} 0 & t < 0 \text{ or } t \geq \delta \\ -t & 0 \leq t < \delta. \end{cases}$$

After dn/δ layers, any input entry of $\mathbf{X} + \mathbf{E}$ in $[k\delta, k\delta + \delta)$ is quantized to $k\delta$.

G.2 Contextual mapping by attention layers

By Step 1, we quantized any input $\mathbf{X} + \mathbf{E}$ to its quantized version. We call this quantized version \mathbf{L} :

$$\mathbf{L} \in [0 : \delta : 1 - \delta]^d \times [1 : \delta : 2 - \delta]^d \times \cdots \times [n - 1 : \delta : n - \delta]^d.$$

As done in Lemma 8, we define $\mathbf{u} := (1, \delta^{-1}, \dots, \delta^{-d+1})$ and $l_j := \mathbf{u}^T \mathbf{L}_{:,j}$, for all $j \in [n]$. Note that, because $\mathbf{L}_{:,j} \in [j - 1 : \delta : j - \delta]^d$, we have

$$(j - 1)(1 + \delta^{-1} + \cdots + \delta^{-d+1}) \leq l_j \leq (j - 1)(1 + \delta^{-1} + \cdots + \delta^{-d+1}) + \delta^{-d+1} - \delta,$$

and $l_1 < l_2 < \cdots < l_n$. Notice that this corresponds to the Category 1 in the proof of Lemma 8.

For simplicity of notation, let $s_j = (j - 1) \sum_{k=0}^{d-1} \delta^{-k}$. We stack $n(1/\delta)^d$ attention layers, with attention parts $\delta^{-d}\Psi(\cdot; l - \delta/2, l + \delta/2)$ for each $l \in \bigcup_{j=1}^n [s_j : \delta : s_j + \delta^{-d+1} - \delta]$, in increasing order of l .

These $n(1/\delta)^d$ attention layers perform selective shift operations on l_j 's, in increasing order of j . As seen in Appendix E.3.1, shift operations result in $\tilde{l}_1 < \tilde{l}_2 < \cdots < \tilde{l}_n$. Also, the map from \mathbf{L} to \tilde{l}_n is one-to-one, which can be shown in the same way as Appendix E.3.4. Since the range of l_j 's are a bit different, we have a different upper bound on \tilde{l}_n :

$$\begin{aligned} \tilde{l}_n &:= l_n + \sum_{k=1}^{n-1} \delta^{-kd} (l_{n-k} - l_{n-k+1}) + \delta^{-nd} (l_n - l_1) \\ &\leq l_n + \delta^{-d} (l_1 - l_n) + \delta^{-nd} (l_n - l_1) \leq s_n + \delta^{-d+1} - \delta + (\delta^{-nd} - \delta^{-d}) (s_n + \delta^{-d+1} - \delta) \\ &= (\delta^{-nd} - \delta^{-d} + 1) \left((n-1) \frac{\delta^{-d} - 1}{\delta^{-1} - 1} + \delta^{-d+1} - \delta \right) \\ &\leq (\delta^{-nd} - \delta^{-d} + 1) (\delta^{-d} - 1) (n - 1 + \delta) < n\delta^{-(n+1)d}. \end{aligned}$$

Finally, we add an extra single-head attention layer with attention part $n\delta^{-(n+1)d-1}\psi(\cdot; 0)$. We define the output of this layer as $g_c(\mathbf{L})$. In a similar way as Appendix E.3.1, this layer shifts all the layers by $n\delta^{-(n+1)d-1}\tilde{l}_n$, thus making the intervals corresponding to different values of \tilde{l}_n disjoint from each other. This ensures that different contexts \mathbf{L} are mapped to distinct numbers in $\mathbf{u}^T g_c(\mathbf{L})$, thus implementing a contextual mapping.

G.3 Function value mapping by feed-forward layers

Now, it is left to map $g_c(\mathbf{L})$ to the desired output. As seen in the last step, each different context \mathbf{L} maps to n unique numbers $\mathbf{u}^T g_c(\mathbf{L})$, which are at least δ apart from each other. The value mapping step can be done in a similar way as Lemma 9. The construction now requires $O(n(1/\delta)^{dn})$ layers because there is no permutation equivariance.

H Discussion on other architectures and Experiments

As detailed in Section E.1, the ability of the self-attention layers to compute contextual mappings plays a crucial role in the universal approximation property. Interestingly, our analysis shows that replacing the dot-product attention in Transformers with any other component capable of computing contextual mappings should preserve this universal approximation property. This leads naturally to questions about the alternative architectures that realize certain kinds of contextual mappings at different computational and memory costs. We explore and discuss some examples of such alternatives in this section. Our preliminary empirical study demonstrates their practical utility.

H.1 Bi-linear projection

Given token embeddings \mathbf{X} as input, the bi-linear projection layer computes the following update.

$$\text{BProj}(\mathbf{X}) = \mathbf{X} + \mathbf{W}_O \cdot \mathbf{X} \cdot \mathbf{W}_P. \quad (5)$$

The bi-linear projection layer [Gong et al., 2013] is motivated from the ability of random (Gaussian) matrices to map sparse differences to dense vectors [Ailon and Chazelle, 2009]. If there are two

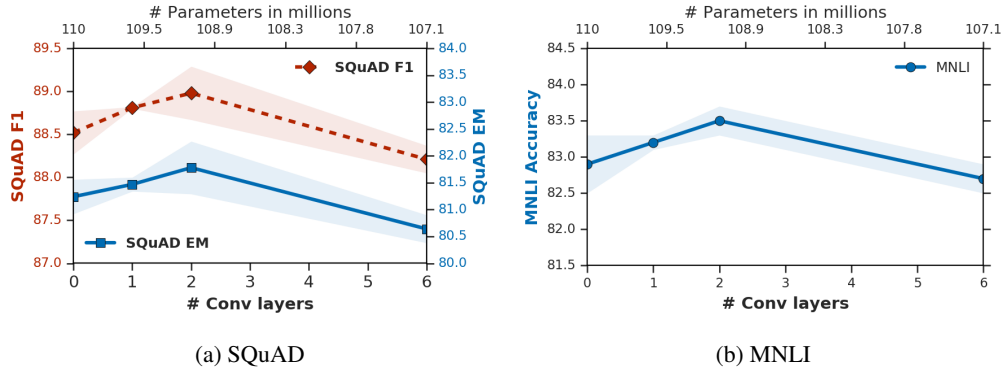


Figure 1: Performance of hybrid models constructed by first taking $\text{BERT}_{\text{BASE}}$, a 12 layer Transformer model, and replacing the self-attention layers with depth-wise separable convolution layers, in a varying number of the Transformer blocks closer to the input. Surprisingly, replacing 1 or 2 self-attention layers with convolutions improves the performance, while replacing more hurts the performance. This suggests both that Transformers have functionality beyond just computing contextual mappings, and having simpler layers to realize contextual mapping can aid Transformers.

input contexts \mathbf{X}_1 and \mathbf{X}_2 that differ in one token, their difference $\mathbf{X}_1 - \mathbf{X}_2$ is sparse; however, after random projection, the difference $(\mathbf{X}_1 - \mathbf{X}_2)\mathbf{W}_P$ will be dense, and the numbers are distinct with high probability, implementing a form “pair-wise contextual mapping,”² although different from the contextual mapping in Definition E.1.

This layer advantageously incurs smaller number of matrix multiplications as compared to the dot-product attention. That said, the number of parameters in this layer depend on the sequence length, making it harder to reuse the model across tasks with different input sequence lengths. Moreover, the weights used to compute the contextual embeddings (\mathbf{W}_P) are independent of the inputs (\mathbf{X}), whereas in self-attention the weights ($\sigma[(\mathbf{W}_K^i \mathbf{X})^T \mathbf{W}_Q^i \mathbf{X}]$) depend on \mathbf{X} . The first drawback can be addressed by replacing the linear projection with a depth-wise separable convolution layer, which is discussed in the next subsection.

H.2 Depth-wise separable convolutions

A depth-wise convolution layer [Sifre and Mallat, 2014, Chollet, 2017, Kaiser et al., 2017] involves convolving each dimension of \mathbf{X} with a corresponding convolution filter of size k :

$$\text{SepConv}(\mathbf{X}) = \mathbf{X} + \mathbf{W}_O(\mathbf{X} * \mathbf{W}_C), \quad (6)$$

where $\mathbf{W}_C \in \mathbb{R}^{d \times k}$ and $(\mathbf{X} * \mathbf{W}_C)_{i,:} := \mathbf{X}_{i,:} * (\mathbf{W}_C)_{i,:}$. Unlike bi-linear projection, this layer can be used across tasks with different input sequence lengths as the number of parameters are independent of the sequence length. While a single layer is unable to compute contextual mappings when the filter size is small, stacking multiple such layers can potentially provide a cheaper way to compute contextual mappings. In fact, based on depth-wise separable convolutions, Wu et al. [2019] proposed a light-weight dynamic convolution architecture that performs competitively with Transformers on machine translation.

H.3 Experiments

We now present our experiments with these other architectures, with the goal of understanding the extent to which computing contextual mappings can capture the performance of Transformers. As discussed earlier, BProj and SepConv do *not* implement contextual mappings (cf. Definition E.1), so we do not expect that either BProj or SepConv based models to have the same performance as the expensive Transformers. These models do not use input dependent weights to compute attention, and hence have weaker representation power. Instead, our goal is to see if we can use these cheaper layers to replace (some of) the expensive self-attention layers.

²This guarantee only holds for a finite set (can be exponential in n) of fixed vectors in \mathbb{R}^n .

We follow the experimental setting from Devlin et al. [2018] to train the Transformers, with the masked language model pre-training followed by a task specific fine-tuning, and work with a 12 layer architecture based on BERT_{BASE}. We present our results on a question answering task (SQuAD) [Rajpurkar et al., 2016] and a sentence entailment task (MNLI) [Williams et al., 2018]. In our first set of experiments we train models that employ BProj and SepConv layers, instead of the self-attention layer in eq.(1). We notice that, as expected, these simpler models have weaker performance than the self-attention layer. See Table 1 in Section H.4 for a comparison of these models on MNLI.

Next, we swap a varying number of the first few self-attention layers in BERT_{BASE} with SepConv, implemented with filter reuse across dimensions [Wu et al., 2019]³. Fig. 1 illustrates the performance of these hybrid models. Interestingly, models with 1 or 2 convolution layers and rest the self-attention layers, perform better than models with only the self-attention layers. Note that, replacing self-attention layer with SepConv also reduces the computational cost and the number of parameters. One explanation we have is that the first few attention layers tend to attend broadly to the whole sequence (as empirically observed in [Clark et al., 2019]), and the cheaper convolution layers can perform this job more efficiently. A detailed evaluation of such hybrid architectures will be interesting future research.

Our experiments also call for a deeper understanding of the exact nature of the embeddings computed by practical attention models. Since Transformers in practice have fixed depth, we believe that they might not be able to exactly implement contextual mappings as we defined in Definition E.1. However, there is some preliminary empirical evidence that Transformers do implement some sort of “contextual mappings.” For example, Fig. 4 of Coenen et al. [2019] presents visualizations of embeddings of a single word in different contexts (sentences). They experimentally notice that Transformers, in addition to computing contextual mappings, also map a word into semantic clusters. Formalizing and evaluating this property of Transformers is an interesting direction for future work. We again note that Wu et al. [2019] have proposed an alternative way to compute such embeddings based on dynamic convolution layers. Evaluating the mappings computed by these models should shed more light on the workings of attention models and inspire efficient and better performing architectures.

H.4 Detailed experiment setup

For our experiments we follow the same setting as in BERT [Devlin et al., 2018]. We first pre-train the models on the masked language modeling task and the next sentence prediction task. We use English Wikipedia corpus and BooksCorpus dataset [Zhu et al., 2015] for this pre-training. We use BERT_{BASE}, a 12 layer Transformer model as the baseline. This model uses an embedding size of 768 and has 12 head self-attention layers and 3072 wide feed forward layers. We train it with the Adam optimizer, with .01 dropout and weight decay. We do pre-training for 250k steps with a batch size of 1024 and a max sequence length of 512. Pre-training takes around 2 days on 16 TPUv3 chips. We take the pre-train models and finetune them on the MNLI and SQuAD datasets separately using the same hyper-parameters as in Devlin et al. [2018]. MNLI is a sentence entailment task in which, given a premise sentence, requires us to classify a hypothesis sentence into neutral, contradiction or entailment classes. We report the classification accuracy on this task. SQuAD is a question answering task, in which given a paragraph and a question, requires us to identify the answer as a span of the words in the paragraph. For this task we report both the F1 score and the Exact Match (EM) percentage. The metrics are reported on the dev sets of these datasets.

For our experiments with the depth-wise separable convolution layers, we follow the implementation in [Wu et al., 2019]. We first use a GLU layer followed by the convolution layer. We use 16 separable convolution filters, of filter length 128, and reuse them, with each filter operating on 48 of the 768 dimensions of the input. This layer also has a skip connection and the output is normalized using layer normalization, similar to the self-attention layer. In our experiments, we replace the self-attention layers of the Transformers, in the lower layers, with this convolution layer. We keep the feed forward layer of the Transformer block the same.

For the experiments performed in this paper, one might consider an alternate explanation that the tasks considered maybe are easy, and do not require any advanced architecture to solve them, and even a simple architecture (bi-linear projection or separable convolution) might solve these tasks. To

³We refer to Section H.4 for a complete description of the setup.

Architecture	Average Attention	BProj	SepConv	Transformer
# params	88.3M	90M	102.5M	110M
Masked LM accuracy (%)	28	59	60	63
MNLI accuracy (%)	66	72.3	73	78.2

Table 1: Performance of bi-linear projection and separable convolution layers on masked LM pre-training task and MNLI. Note that we expect these computationally cheaper models to have lower performance than the expensive Transformers as they do not compute input dependent attention weights and have weaker representation power. Our goal in studying them is to see if they can substitute some of the expensive attention layers for computing the contextual mappings. These models are trained in a large batch setting, with a batch size of 8192 for 60k steps, unlike the other set of experiments reported in Fig. 1. Note that average attention has clearly worse performance, showing that these tasks indeed require an advanced architecture.

rule out this case we consider an even simpler architecture, namely average attention, as a baseline for our experiments.

Average attention. An average attention layer replaces the self-attention layer, and just computes the average of projections of all the other tokens. That is, we replace $\sigma[(\mathbf{W}_K^i \mathbf{X})^T \mathbf{W}_Q^i \mathbf{X}]$ in (1) with a matrix full of $1/n$. The model still has the skip connections and the feed-forward layers like Transformer.