

TUFORMER: DATA-DRIVEN DESIGN OF EXPRESSIVE TRANSFORMERS BY TUCKER TENSOR REPRESENTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformers are neural network architectures that achieve remarkable performance in many areas. However, the core component of Transformers, multi-head self-attention (MHSA), is mainly derived from heuristics, and the interactions across its components are not well understood. To address the problem, we first introduce a mathematically rigorous and yet intuitive *tensor diagram* representation of MHSA. Guided by tensor diagram representations, we formulate a design space where we can analyze the expressive power of the network structure, providing new directions and possibilities for enhanced performance. We then propose a novel design, namely *Tucker Transformer* (Tuformers), inspired by a variant of Tucker representation with a guaranteed higher expressive power than MHSA. Unlike vanilla Transformer models, where the number of heads is a pre-defined fixed constant, Tuformer’s structure is data-driven, and the number of heads is trainable. Training of Tuformers could be made very efficient as it allows initialization from existing pre-trained Transformer models. We test Tuformers on various tasks across multiple domains and show competitive results under a wide range of model sizes.

1 INTRODUCTION

Transformer models are first introduced by Vaswani et al. (2017) in the context of sequence modelling. They have demonstrated impressive results on a wide variety of tasks in many fields, such as language model pre-training (Sanh et al., 2019), speech recognition (Moritz et al., 2020), image classification (Dosovitskiy et al., 2020), and generation (Parmar et al., 2018). The core component of Transformer models is the *multi-head self-attention* (MHSA) which is extended from the standard attention mechanism (Bahdanau et al., 2014). Each attention head in MHSA has a global receptive field, i.e., each token’s representation is updated by attending to all other tokens, and H attention heads are computed in parallel and concatenated together.

The current MHSA design is mainly derived from empirical studies or heuristics, leaving some unresolved challenges. **(1) Lack of solid theoretical understanding.** The theory behind MHSA is only starting to catch up with practice. The role of the components and their interactions in MHSA are not well understood rigorously from a theoretical perspective, which may lead to inefficiencies in the design. For example, Michel et al. (2019) show that most attention heads can be removed in the testing phase without much performance compromise, while Cordonnier et al. (2020) find empirical evidence of redundancy in key/query projections and propose a re-parameterization scheme. These works focus more on practical solutions, leading to questions about whether theoretical patterns exist in these designs. **(2) The number of heads is not trainable.** Intuitively, the heads in multi-heads are expected to capture different context information through each head. However, the number of heads is fixed in training. Thus, although we could tune the hyper-parameter, an exhaustive search would be time-consuming or not practical for large-scale tasks. **(3) Hard to analyze the expressive power.** Analyzing the expressive power of a neural network, i.e., proving that some architectures are more expressive than others, is a non-trivial and challenging task. Cohen et al. (2016) showed the relationship between Hierarchical Tucker (HT) Decompositions (Hackbusch & Kühn, 2009) and Convolutional Neural Networks (CNNs) (LeCun et al., 1995) and compared the expressive power between CNNs and shallow neural networks. Khrulkov et al. (2018) extended Cohen et al. (2016) and analyzed the expressive power of Recurrent Neural Networks (RNNs) (Mikolov et al., 2010) using

Tensor-Train (TT) (Oseledets, 2011) decompositions. However, no such work exists for interpreting MHSA or guiding the structural design.

In response to the above challenges, we first **analyze the expressive power** of MHSA from a tensor representation perspective using the intuitive graphical tool, the *tensor diagram*. Current prevalent descriptions of MHSA use flow charts to convey high-level intuitions, which could cause ambiguities. Therefore, it is inevitable to pair those flow charts with mathematical formulas to understand the mechanism precisely. However, these two separated descriptions create difficulties for interpretations and inspections of the operations implemented in MHSA. To address this issue, we propose a graphical representation of MHSA, which is both semantically intuitive and mathematically rigorous. Specifically, we modify and extend the vanilla *tensor diagrams* (Penrose, 1971), which conveniently allow for rigorous graphical representation of multi-linear operations between multi-dimensional arrays (i.e., higher-order tensors) to represent nonlinear operations.

We then **formulate a design space** for MHSA. Using tensor diagram representation, we project the current design of MHSA into its tensor form, which renders a holistic view of the weights in MHSA for better interpretation of the information flow and exchange among the components of MHSA. More importantly, we formulate a design space that allows improving the expressive power of the network by finding better tensor representations of the weights through tensor representation theory.

We also **propose a novel data-driven structure** that is guaranteed to have higher expressive power. Inspired from the tensor diagram representation of MHSA, we obtain a novel design, namely *Tucker-Head Self-Attention (THSA)*, which is a re-parameterization of the weight matrices via a variant of Tucker Decomposition or Tucker Tensor Representation (Rabanser et al., 2017). Transformers with THSA, named *Tucker Transformers (Tuformers)*, have several advantages compared against vanilla Transformers: (1) A guaranteed higher expressive power. We prove that MHSA is a special case of THSA. (2) The number of heads is trainable. The concept of the number of heads in THSA generalizes to the stable rank of the core matrix, allowing data-driven implicit training. (3) Tuformers allow initialization from pre-trained Transformers such as BERT (Devlin et al., 2019) and its variants.

We experiment Tuformers with several tasks across multiple domains, from language modeling, machine translation to image generation under a wide range of model sizes. We demonstrate competitive results not only on Tuformers but also in cases where Tuformers are initialized with pre-trained Transformers for other downstream tasks and when combined with Linear Transformer (Katharopoulos et al., 2020) on the image generation task.

Summary of Contributions:

- (1) We propose a mathematically rigorous and semantically intuitive *tensor diagram* representation of the multi-head self-attention, introducing a new tool to the ML community for future studies on interpretation and improvements of Transformers.
- (2) We formulate a design space for attention unit design and a framework for analyzing their expressive power using tensor representation theory under the intuitive guidance of tensor diagrams.
- (3) We propose a novel design of the MHSA, *Tucker-Head Self-Attention (THSA)*, resulting in *Tuformers*, which has a theoretical guaranteed improvement in expressive power.
- (4) We provide a constructive re-parameterization scheme from vanilla Transformers to *Tuformers* so that *Tuformers* can be easily used in fine-tuning tasks with pre-trained models and efficient Transformers such as Linear Transformer (Katharopoulos et al., 2020) and Performer (Choromanski et al., 2020) with the state-of-the-art computation and memory complexities.

2 TENSOR DIAGRAM REPRESENTATION

Notations. We use lower case letters (e.g., v) to denote vectors, upper case letters (e.g., M) to denote matrices, and curly letters (e.g., \mathcal{T}) to denote general tensors. For a tensor $\mathcal{T} \in \mathbb{R}^{I_1 \times \dots \times I_M}$, we refer to the number of indices as *order*, each individual index as *mode*, and the length of one mode as *dimension*. For instance, \mathcal{T} is an M^{th} order tensor that has dimension I_m at its m^{th} mode. We reserve superscripts to distinguish similar arrays (e.g., $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V$ are query/key/value weight matrices), and subscripts to index the elements in an array (e.g., \mathbf{W}_{ij} is the $(i, j)^{\text{th}}$ element of \mathbf{W}). We use colon $:$ to slice an array (e.g., $\mathbf{W}_{i,:}$ denotes the i^{th} row of \mathbf{W}).

We propose to use *tensor diagrams* (Penrose, 1971), a commonly used rigorous/precise and intuitive graphical representation for multi-linear operations among higher-order arrays (i.e., tensors), to

represent multi-head self-attention (MHSA). Since MHSA consists of multilinear operations and a nonlinear softmax function, we will introduce tensor diagrams and our novel design extending tensor diagrams to denote the composition of nonlinear and multilinear operations in MHSA.

2.1 TENSOR DIAGRAM BASICS

Arrays denoted as nodes with legs. An array is represented as a node with leg(s) in a tensor diagram as shown in Figure 1. We denote the *order* (the number of dimensions) of the array by the number of legs extending from the node. Each labeled leg represents one *mode* of a tensor. Every mode of the tensor needs to be uniquely labeled. We usually use the dimension of the mode as the label (i.e., an associated positive integer written on top of each leg). The legs do not need to be straight lines, and their orientations do not matter. Matrices $M \in \mathbb{R}^{A \times B}$ and $M^T \in \mathbb{R}^{B \times A}$ can be represented via the same tensor diagram as long as the M node has two legs (to denote that it is a matrix, wherever the legs extend to), labeled as A and B (to denote its size).

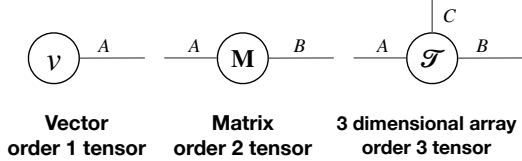


Figure 1: Arrays in tensor diagram. A vector is denoted as a node with 1 leg, a matrix as a node with 2 legs and an N -dimensional array as a node with N legs.

Operations in tensor diagrams. There are three types of operations in the calculation of MHSA: *contraction*, *softmax*, and *batch multiplication*, as shown in Figure 2 and explained in its caption.

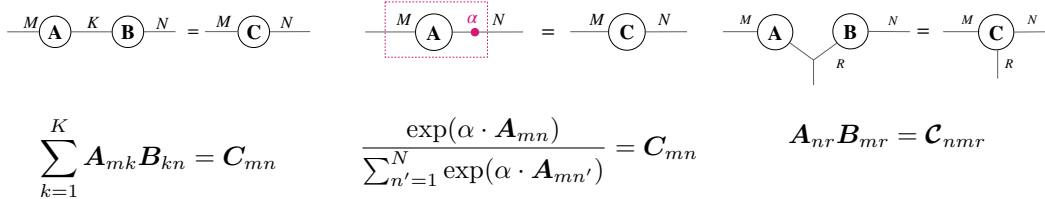


Figure 2: Tensor diagrams for atomic operations. (1) (left) **Contraction** is an operation that generalizes the matrix multiplication. It sums element-wise products on a mode in object A and a corresponding mode (with the same dimension) in object B (i.e., along a leg in node A and a corresponding leg in node B). In the tensor diagram, multiplying two matrices (or higher-order tensors with more than 2 legs) corresponds to “gluing” their corresponding legs (on a certain mode). (2) (middle) **Softmax** is an element-wise exponential function normalized along a certain mode. We propose to denote the α -scaled softmax function $\text{softmax}(\alpha A)$ on A as a dotted box with a labeled filled ball (to distinguish itself from tensor objects, i.e., nodes which are blank circles) attached to one leg. (3) (right) **Batch multiplication** is an element-wise product along the connected legs.

Evaluation of tensor diagrams: (1) Evaluation order. We can evaluate a tensor diagram in any pair-wise order except for the nodes in the softmax box. Since the softmax function is nonlinear, we must first evaluate the nodes in the softmax box with arbitrary order before those outside the box. (2) **Reading a tensor diagram.** We can easily identify the output shape by the dangling edges. For instant, in Section 2.1, there are three dangling legs M, N, R . Thus the output is a 3rd order tensor with dimensions M, N, R . Note that the softmax function does not change the shape of the tensor.

Advantages of tensor diagrams: (1) Tensor diagram is orientation invariant, meaning that we can represent X and X^T using the same tensor diagram. Thus in multi-head, we obtain a universal graphical representation regardless of whether we represent an embedding of each token in the sequence as rows or columns of the input embedding matrix X . (2) Multi-linear operations are concisely represented and interpreted. (3) The representation is both precise in math and intuitive, and the information flow is clear, making analysis of network structure more accessible. In addition, with the labels of the legs, we can read the model complexity explicitly.

We include a comprehensive introduction to tensor diagram representation in Appendix A.

2.2 TENSOR DIAGRAM REPRESENTATION OF MULTI-HEAD SELF-ATTENTION

The core of Transformer models is a *multi-head self-attention* (MHSA) module which allows the model to jointly attend to information from different representation sub-spaces at different posi-

tions (Vaswani et al., 2017). To distinguish inputs for query, key, value matrices, we use $\mathbf{X}^Q \in \mathbb{R}^{N \times F}$, $\mathbf{X}^K, \mathbf{X}^V \in \mathbb{R}^{M \times F}$ respectively. The MHSA module outputs a matrix $\mathbf{M} \in \mathbb{R}^{N \times F}$ as

$$\mathbf{Q}_{[h]} = \mathbf{X}^Q \mathbf{W}_{[h]}^Q; \mathbf{K}_{[h]} = \mathbf{X}^K \mathbf{W}_{[h]}^K; \mathbf{V}_{[h]} = \mathbf{X}^V \mathbf{W}_{[h]}^V, \quad (1a)$$

$$\mathbf{head}_{[h]} = \text{softmax} \left(\mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top / \sqrt{D} \right) \mathbf{V}_{[h]}, \quad (1b)$$

$$\mathbf{M} = [\mathbf{head}_{[1]}, \mathbf{head}_{[2]}, \dots, \mathbf{head}_{[H]}] \mathbf{W}^O. \quad (1c)$$

In the above equations, we see how MHSA computes its output in three *steps*:

(1a) *Linear transformation on input embedding to get latent features.* For each h , the query/key/value matrices $\mathbf{Q}_{[h]}, \mathbf{K}_{[h]}, \mathbf{V}_{[h]}$ are linear transformations, parameterized by $\mathbf{W}_{[h]}^Q, \mathbf{W}_{[h]}^K$, and $\mathbf{W}_{[h]}^V \in \mathbb{R}^{F \times D}$ respectively, of the input embedding.

(1b) *Scaled dot-product attention.* Each $\mathbf{head}_{[h]} \in \mathbb{R}^{N \times D}$ computes a scaled dot-product between a latent feature matrix $\mathbf{Q}_{[h]} \in \mathbb{R}^{N \times D}$ (i.e., a *query* matrix) and another latent feature matrix $\mathbf{K}_{[h]} \in \mathbb{R}^{M \times D}$ (i.e., a *key* matrix) along D , and applies a softmax function to obtain the weights to be multiplied with a third latent feature matrix $\mathbf{V}_{[h]} \in \mathbb{R}^{M \times D}$ (i.e., a *value* matrix).

(1c) *Concatenation and contraction for MHSA.* After concatenation of the H heads, a linear transformation of the concatenated heads, parameterized by $\mathbf{W}^O \in \mathbb{R}^{HD \times F}$, is implemented to get the outcome \mathbf{M} of the MHSA as a $\mathbb{R}^{N \times F}$ object.

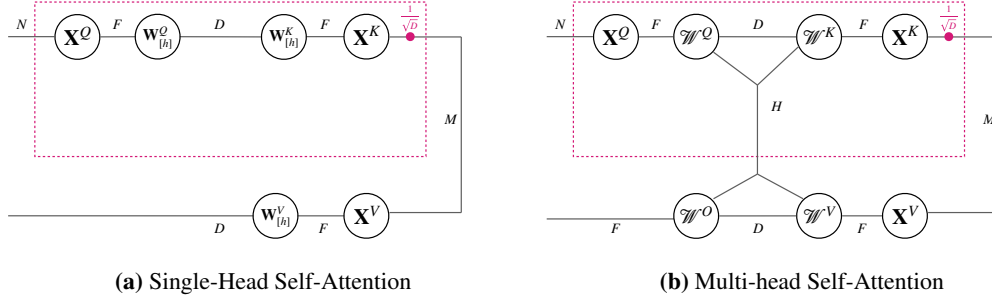


Figure 3: Figure 3a shows the tensor diagram of a single attention head. Figure 3b shows MHSA. See Appendix B for more details and the proof.

Now we will see how to represent MHSA in the tensor diagram. We first explain the representation of a single head and then extend it to multiple heads.

Single-head self-attention in tensor diagram. Figure 3b shows how to extend SHSA to MHSA as in *Step (1c)*. The *concatenation* of H matrices $\{\mathbf{head}_{[h]}\}_{h=1}^H$ is interpreted as constructing a three-dimensional array with the third mode being the number of heads, i.e., three-legged nodes in tensor diagram notation. This way, the weight matrices $\{\mathbf{W}_{[h]}^Q\}_{h=1}^H$ are stacked to be an order 3 tensor \mathcal{W}^Q . It is the same with other weight matrices. Since in MHSA, the attention calculation among key, query, and value matrices are done separately within each head, batch multiplication on the mode H is a proper way to denote these “head-wise” operations. Finally, we connect the resultant tensor with \mathcal{W}^O to obtain MHSA as in Figure 3b. See Appendix B for more details and the proof.

Multi-head self-attention in tensor diagram. Figure 3b shows how to extend SHSA to MHSA as in *Step (1c)*. The *concatenation* of H matrices $\{\mathbf{head}_{[h]}\}_{h=1}^H$ can be viewed as constructing a three dimensional array with the third mode being the number of heads, i.e., three-legged nodes in tensor diagram notation. This way, the weight matrices $\{\mathbf{W}_{[h]}^Q\}_{h=1}^H$ are stacked to be an order 3 tensor \mathcal{W}^Q . It is the same with other weight matrices. Since in MHSA, the attention calculation among key, query and value matrices are done separately within each head, batch multiplication along mode H is a proper way to denote these “head-wise” operations. Finally, we connect the resultant tensor with \mathcal{W}^O to obtain MHSA as in Figure 3b. See Appendix B for more details and the proof.

3 TUFOMERS: DATA-DRIVEN DESIGN WITH HIGHER EXPRESSIVE POWER

We have seen the tensor diagram of MHSA in Section 2.2. A natural question to ask then is whether there exists a paradigm that outperforms MHSA. We introduce our improved design of the self-

attention module, *Tucker-head self-attention (THSA)*, inspired by a Tucker tensor form. We will demonstrate the tensor diagram of THSA and its mathematical formulation. We will then show how MHTSA motivates THSA and prove that the expressive power of THSA is higher than that of MHTSA.

3.1 A NOVEL DESIGN: TUCKER-HEAD SELF-ATTENTION (THSA)

We have seen the tensor diagram of MHTSA in Section 2.2. A natural question to ask then is whether there exists a paradigm that outperforms MHTSA. In response, we introduce our improved design of the self-attention module, *Tucker-head self-attention (THSA)*, inspired by a Tucker tensor form. First, we demonstrate our THSA design in the form of the tensor diagram and mathematical equations. We then prove that our THSA has higher expressive power than MHTSA.

With the tensor diagram of MHTSA, it is also clear that if we take a holistic consideration of the weight tensors $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V, \mathbf{W}^O \in \mathbb{R}^{H \times D \times F}$ and consider them as components of an order 4 tensor of size $F \times F \times F \times F$ (the “super node” represented as the shaded ellipse in Figure 4a), then we can use the tensor representation theory to “re-parameterize” the 4th order tensor.

Inspired by a variant of Tucker decomposition (Kolda & Bader, 2009) in tensor representation theory, we propose to use a variant of Tucker decomposition to represent the $F \times F \times F \times F$ “super node” (tensor), which replaces these two types of contractions with a *Tucker Contraction*, i.e., a global contraction among all weights and the core matrix. Note that the core matrix is trainable.

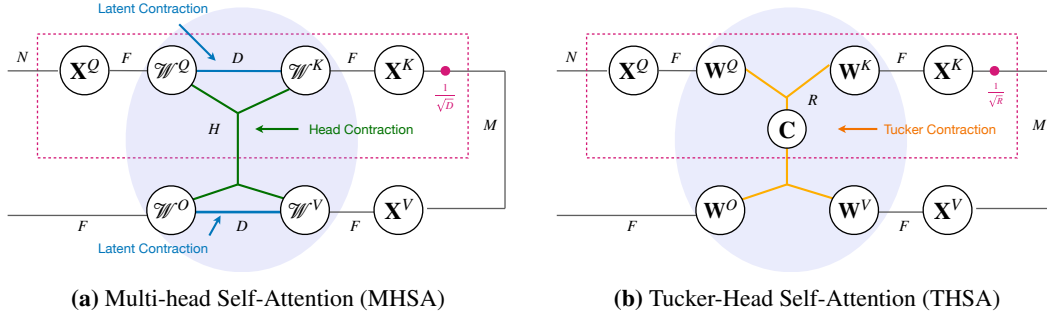


Figure 4: From Multi-Head Self-Attention to Tucker-Head Self-Attention.

We propose a Tucker-head self-attention (THSA), which generalizes the inner product among weight matrices in MHTSA to a simultaneous contraction among weight matrices and a core matrix \mathbf{C} , as shown in Figure 4b. A Tucker-head self-attention (THSA) module has the same input and output domains as an MHTSA, i.e., it takes the input embedding $\mathbf{X}^Q \in \mathbb{R}^{N \times F}$ and $\mathbf{X}^K, \mathbf{X}^V \in \mathbb{R}^{M \times F}$ as inputs and returns a matrix $\mathbf{T} \in \mathbb{R}^{N \times F}$ as output. A THSA, with five weight matrices $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{F \times R}$, $\mathbf{W}^O \in \mathbb{R}^{R \times F}$, and $\mathbf{C} \in \mathbb{R}^{R \times R}$, is mathematically described as

$$\mathbf{Q} = \mathbf{X}^Q \mathbf{W}^Q; \mathbf{K} = \mathbf{X}^K \mathbf{W}^K; \mathbf{V} = \mathbf{X}^V \mathbf{W}^V, \quad (2a)$$

$$\mathbf{head}_r = \text{softmax} \left(\sum_{s=1}^R \mathbf{C}_{rs} \mathbf{Q}_s \mathbf{K}_s^\top / \sqrt{R} \right) \mathbf{V}_r, \quad (2b)$$

$$\mathbf{T} = [\mathbf{head}_1, \mathbf{head}_2, \dots, \mathbf{head}_R] \mathbf{W}^O. \quad (2c)$$

In Equation (2a), we compute the query, key, value matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{F \times R}$ from the corresponding input and weight matrices. In Equation (2b), we calculate the head matrix using $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ as well as the core matrix \mathbf{C} . Here, we use $\mathbf{Q}_r, \mathbf{K}_r, \mathbf{V}_r \in \mathbb{R}^F$, $\mathbf{head}_r \in \mathbb{R}^R$ to denote the r^{th} column of the corresponding matrix, and \mathbf{C}_{rs} to denote the $(r, s)^{\text{th}}$ element of \mathbf{C} . In MHTSA, the softmax function is scaled, which normalizes its input by \sqrt{D} , i.e., the square root of the latent dimension. The scaling ensures that the distribution of the attention matrix approximates a standard Gaussian distribution. To achieve a similar effect, in THSA, we change the constant to \sqrt{R} accordingly. In Equation (2c), we transform the concatenation of head matrices to THSA through the output weight matrix \mathbf{W}^O .

The tensor diagram representation of THSA is given in Figure 4b. It has an additional core matrix \mathbf{C} compared to MHTSA in Figure 4a. As we will show in Proposition 1, *THSA generalizes MHTSA*;

MHSA is a special case in which the core matrix is fixed and not trainable. In contrast, the core matrix in THSA is trainable.

Proposition 1 (THSA generalizes MHSA). *Given $R = DH$ and other hyperparameters being the same, THSA module includes MHSA as a special case. Specifically, if $C = \mathbf{I}_H \otimes (\mathbf{1}_D \mathbf{1}_D^\top)$, i.e., a Kronecker product of an all-ones matrix $\mathbf{1}_D \mathbf{1}_D^\top \in \mathbb{R}^{D \times D}$ and an identity matrix $\mathbf{I}_H \in \mathbb{R}^{H \times H}$, the THSA reduces to an MHSA with H heads and latent dimension D .*

Proof. We prove Proposition 1 in Appendix D.1. \square

Remarks. (1) We call a Transformer with THSA module as *Tucker Transformers (Tuformers)*. **(2)** *Tuformers can be initialized with Pre-trained Transformers.* Although different in the structure, Tuformers can always be initialized with pre-trained Transformer models through re-parameterization. MHSA is updated by weight matrices $\{\mathbf{W}_{[h]}^K\}_{h=1}^H$, $\{\mathbf{W}_{[h]}^Q\}_{h=1}^H$, $\{\mathbf{W}_{[h]}^V\}_{h=1}^H$ and \mathbf{W}^O while in THSA we have \mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V , \mathbf{W}^O and C . To initialize THSA with MHSA, we follow the scheme proposed in Proposition 1. C is the kronecker product between a $D \times D$ all-one matrix and an $H \times H$ identity matrix, where D is the latent dimension and H is the number of heads in MHSA. \mathbf{W}^O remains the same as that of MHSA. As for \mathbf{W}^Q , \mathbf{W}^K and $\mathbf{W}^V \in \mathbb{R}^{F \times R}$, we concatenate $\{\mathbf{W}_{[h]}^K\}_{h=1}^H$, $\{\mathbf{W}_{[h]}^Q\}_{h=1}^H$, $\{\mathbf{W}_{[h]}^V\}_{h=1}^H \in \mathbb{R}^{F \times D}$ along mode D .

The notion of heads in THSA. We define the number of heads in THSA as the stable rank of the core matrix C , which is $\sum_i \sigma_i^2 / \max_i \sigma_i$, where σ_i is the singular value of C . Since we replace head-contractions and latent-contractions with tucker-contractions in THSA, it is not immediate clear how to migrate the notions of heads in MHSA to THSA. By Proposition 1, we know that when C is the Kronecker product of the all-one matrix and the identical matrix, THSA is equivalent to MHSA. The stable rank of C of MHSA is the number of heads H . Thus we generalize this correspondence to THSA, allowing an implicit trainable structure. In MHSA, in contrast with THSA, this number does not change during training since C is fixed in MHSA.

3.2 EXPRESSIVE POWER OF TUFORMERS

We first formally describe how to compare the expressive power between two models.

Definition 2 (Expressive Power). *Suppose we have two function classes \mathbb{F}, \mathbb{G} with the same source domain \mathcal{X} and target domain \mathcal{Y} , i.e., each function $f \in \mathbb{F}$ (or $g \in \mathbb{G}$) is a mapping from \mathcal{X} to \mathcal{Y} . We say \mathbb{G} is more expressive than \mathbb{F} if $\mathbb{G} \supseteq \mathbb{F}$: for any $f \in \mathbb{F}$, there exists $g \in \mathbb{G}$ such that $g(x) = f(x), \forall x \in \mathcal{X}$. Furthermore, we say \mathbb{G} is strictly more expressive than \mathbb{F} if $\mathbb{G} \supset \mathbb{F}$: besides $\mathbb{G} \supseteq \mathbb{F}$, there exists $g \in \mathbb{G}$ such that given any f there exists $x \in \mathcal{X}$ and $g(x) \neq f(x)$.*

According to the definition, expressive power is a partial order set. Therefore, expressive power of different models is usually not comparable as two sets might have overlap and non-overlap regions. However, when one is a superset or subset of the other, we can unambiguously compare their expressive power. In the following theorem, we formally show that a Tucker-head self-attention (THSA) module has a higher expressive than multi-head self-attention (MHSA) if the rank R in THSA is equal to the product of the number of heads H and the latent dimension D in MHSA.

Theorem 3 (THSA is more expressive than MHSA). *Given $R = DH$ and other hyper-parameters being the same, a Tucker-head self-attention (THSA) module with rank R is more expressive than a multi-head self-attention (MHSA) module with H heads and latent dimension D .*

Proof. We prove Theorem 3 in Appendix D.2. \square

4 EXPERIMENTS

In this section, we show that THSA (used in Tuformers) universally outperforms MHSA (used in Transformers) under a wide range of model scales in diverse tasks, including natural language processing, automatic speech recognition and image generation.

We also demonstrate that THSA is simple to use in practice. First, THSA is a plug-and-play module that can be used in existing Transformer architectures. Second, THSA can be initialized from pre-trained MHSA modules (to avoid expensive training from scratch). For instance, fine-tuning

Tuformers initialized from pre-trained BERT and its variants further improves their performance. Finally, THSA can be combined with other kernel-based efficient Transformers to obtain linear computational and memory complexities regarding the sequence length.

Datasets, tasks and evaluation metrics. We evaluate Tuformers on 7 datasets for 5 tasks: word-level Language Modeling (LM) on Penn Treebank (PTB) (Marcus et al., 1993), Neural Machine Translation (NMT) on WMT16 ‘English-German’ (Sennrich et al., 2016), Automatic Speech Recognition (ASR) on LibriSpeech (Panayotov et al., 2015), Natural Language Inference (NLI) on MNLI (Williams et al., 2018) and QNLI (Wang et al., 2018) and Image Generation on CIFAR10 (Krizhevsky, 2009) and MNIST (Deng, 2012) datasets.

Baselines. We compare the performance of MHSA module with THSA module in the following 5 backbone architectures: (1) a vanilla Transformer model (Vaswani et al., 2017), (2) an TransformerXL model (Dai et al., 2019), (3) a BERT-large model (Devlin et al., 2019), (4) a RoBERTa-large model (Liu et al., 2019), and (5) an ALBERT model (Lan et al., 2019). In all baselines, we adopt all default settings used in the original paper (hyperparameters specified in Appendix E), and only replace/re-parameterize the MHSA to THSA.

Performance Results of Tuformers.

(1) *THSA consistently outperforms MHSA on diverse datasets for a variety of tasks under different model scales.* In LM, NMT and ASR tasks, we compare THSA against the MHSA models under different model scales. MHSA is parameterized by four weight matrices W^Q , W^K , W^V and W^O whereas THSA adds an additional core matrix C of the same size. Adopting the size of the four original weight matrices in MHSA to THSA incurs an 25% increase in the number of parameters. For a fair comparison, we first compare MHSA against THSA with the same number of parameters (thus with smaller query/key/value/output weight matrices). To get a full picture, we also compare MHSA against THSA with the same sized four original weight matrices. As shown in Figure 5, under the same number of parameters, our THSA (blue curve) universally outperforms MHSA (red curve) under all model scales. By increasing the number of parameters of the attention units by 25%¹ (i.e., adding the core matrix C while keeping the size of the query/key/value/output weight matrices the same), our THSA (purple curve) further improves the performance universally under all model scales.

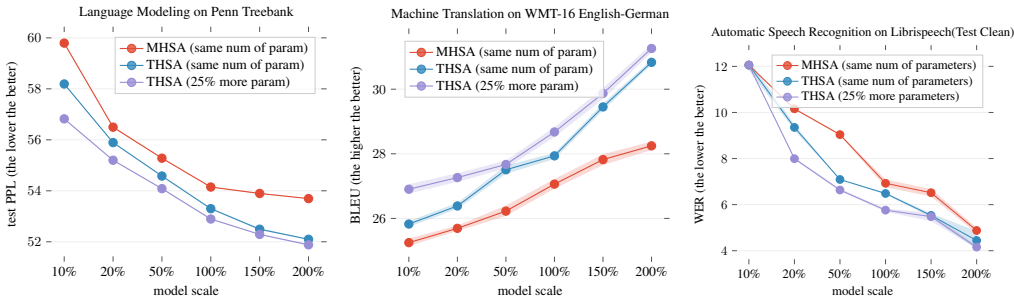


Figure 5: Performance comparison of THSA and MHSA in (left) language modeling on Penn Treebank dataset on TransformerXL, (middle) neural machine translation on WMT16 English-German dataset and (right) automatic speech recognition on Librispeech on Transformers. Under the same number of parameters, our THSA universally outperforms MHSA under all model sizes. By increasing the number of parameters of the attention units parameter by 25% (i.e., adding the core matrix C while keeping the size of the query/key/value weight matrices the same), our THSA further improves the performance universally under all model scales. The error bars displayed in shades are generated using 10 random runs.

(2) *Tuformers initialized with pre-trained strong models further improve the performance.* In NLI task, we initialize Tuformers with pre-trained BERT-large, RoBERTa-large and ALBERT models, and perform the same fine-tuning tasks. As shown in Table 1 (Left), THSA is compatible with existing pre-trained Transformer models and there is no need to train Tuformers from scratch: fine-tuning THSA initialized from pre-trained BERT and its variants further improves their performance.

(3) *Combining THSA with kernels to obtain linear computation and memory complexities in sequence length.* Although an effective design with guaranteed higher expressive power, Tuformers come with

¹As will be discussed later, the overall model size is not drastically increased.

an extra memory overhead when calculating values within the softmax. In Tuformers, if done naively, the memory overhead is MNR where N, M denote the sequence length and $R = DH$ is the rank of the core tensor, incurring a D times larger memory overhead than vanilla Transformers in calculating softmax. Combining Tuformers with some state-of-the-art kernel-based efficient Transformers such as Linear Transformers using ELU (Katharopoulos et al., 2020) and polynomial (Tsai et al., 2019) kernels, we reduce the computation and memory complexities of Tuformers to be linear in the sequence length. Kernel-based Transformers apply kernel methods to MHSA to remove the nonlinearity constrain in the softmax function, obtaining linear computation and memory complexities. For example, in Linear Transformer ($N = M$), the space complexity is $N \times D \times H$. Applying the same idea to THSA leads to the same space complexity. We experiment on image generation tasks on CIFAR10 (Krizhevsky, 2009) and MNIST (Deng, 2012) datasets as in Table 1(Right) and show that Tuformers obtain strong performance and linear space complexity when combined with kernels.

Table 1: (Left) THSA initialized with pre-trained models. Comparison of accuracies obtained for language inference task on MNLI and QNLI datasets, the higher the better. Results show that Tuformers can be conveniently initialized with pre-trained models for downstream tasks, and fine-tuning obtains impressive performance. **(Right) THSA +kernels.** Comparison of bits per dimension (BPD) for image generation task on MNIST and CIFAR-10 datasets, the lower the better. Results show that Tuformers can be extended to image generation tasks and improve the performance of other efficient designs with linear computational and memory complexities.

Models	MNLI	QNLI	Models	MNIST	CIFAR 10
BERT (Devlin et al., 2019)	84.1%	92.1%	ELU kernel (Katharopoulos et al., 2020)	0.72	3.51
BERT + THSA	84.9%	92.8%	ELU kernel + THSA	0.65	3.42
RoBERTa (Liu et al., 2019)	89.4%	94.3%	Polynomial kernel (Tsai et al., 2019)	0.64	3.47
RoBERTa + THSA	90.2%	94.7%	Polynomial kernel THSA	0.57	3.37
ALBERT (Lan et al., 2019)	89.1%	97.9%			
ALBERT + THSA	89.5%	98.1%			

Parameter efficiency of Tuformers. THSA introduces an additional trainable core matrix compared with MHSA. However, THSA does not increase the complexity in big O notation compared with MHSA. Moreover, in Transformer models, attention modules are not the only source of parameters and a 25% increase in MHSA often does not increase the total number of parameters significantly. For instance, if we replace MHSA with THSA by keeping the size of the query/key/value/output weight matrices to be the same, a Transformer encoder with $1.3M$ parameters will be transferred to a Tuformer with $1.6M$ parameters. Finally, we empirically verify that the core matrix design still maintains parameter efficiency. In Figure 6, we benchmark the training time of Tuformers in multiple tasks and show that THSA with 25% more number of parameters (blue curve) is slightly slower to train compared with MHSA (red curve) but significantly improves the performance (as shown in Figure 5). In addition, as shown previously in Figure 5(middle), using 10% parameters of a vanilla Transformer, a Tuformer obtains a 25.8 BLEU score whereas the original vanilla Transformer has a BLEU score of 27.1 (95% of the performance), demonstrating competitive parameter efficiency.

Ablation study: check-pointing to alleviate memory overhead without kernels. Although combining THSA with kernel-based efficient Transformers relieves the memory overhead as shown above, we investigate other solutions to alleviate the D times larger memory overhead compared with vanilla Transformers in calculating softmax when not using kernels. Specifically, we use check-pointing to eliminate the memory overhead in THSA. Since no intermediate results are saved in the forward pass using check-pointing, we recalculate the intermediate results of the softmax box in the back-propagation, which introduces some computational overhead. Figure 6 (purple curve) shows the per iteration training time overhead induced by check-pointing. From Figure 6 we see that applying check-pointing to Tuformers leads to only a slight increase in training time while relieving the memory overhead. We can reduce or even avoid the repetitive calculations by developing a proper backward module, although such development is beyond the scope of this paper and deferred for future works.

5 RELATED WORKS

Analysis of MHSA. Given the popularity of Transformers in a wide range of domains, a line of works focuses on understanding and improving MHSA. Voita et al. (2019) proposed a practical

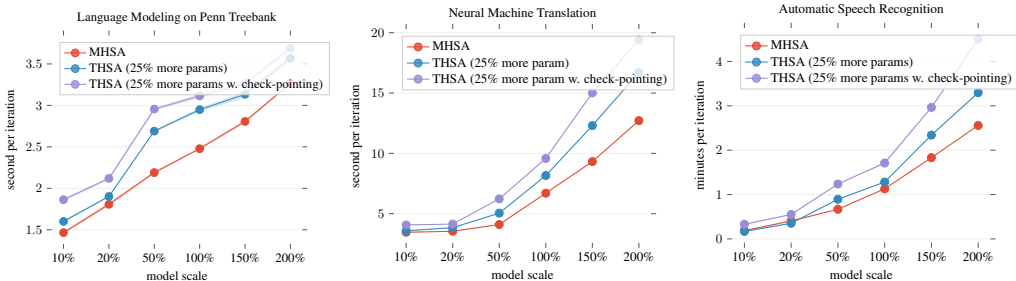


Figure 6: Training run-time comparison between MHSA, THSA with same sized query/key/value/output matrices as MHSA with an additional core matrix (thus 25% more number of parameters in the attention units), and THSA 25% more parameters with check-pointing. **(Left)** In language modeling on PTB dataset, the batch size is 256. **(Middle)** In neural machine translation on WMT16 English-German dataset, batch size is 512. **(Right)** In Automatic Speech Recognition task, the batch size is 2.

scheme to prune less informative heads. [Michel et al. \(2019\)](#) also showed that most attention heads are removable in testing without much performance compromise. [Cordonnier et al. \(2020\)](#) found empirical evidence of redundancy in key/query projections and proposed a re-parameterization scheme. These works bring valuable insights into MHSA, but a consistent theoretical analysis for these findings is missing. [Ma et al. \(2019\)](#) proposed Tensorized Transformers using Block-Term Tensor Decomposition ([De Lathauwer, 2008](#)), but it focused on model compression and improvement on parameter efficiency. Our work differs from these works in that we focus on the theoretical understanding of MHSA in terms of expressive power, and we provide theoretical guarantees for a novel design, Tuformers, with higher expressive power and data-driven trainable heads.

The expressive power of neural networks. Now we present the most related research analyzing Transformer models using expressive power. [Cohen et al. \(2016\)](#) rigorously defined the expressive power in convolutional neural networks from a tensor perspective. [Khrulkov et al. \(2018\)](#) followed [Cohen et al. \(2016\)](#)’s work and extended the theory to recurrent neural networks. [Cordonnier et al. \(2019\)](#) used expressive power in comparing the relationship between the multi-head self-attention and convolutional layers. [Yun et al. \(2019\)](#) established that Transformer models are universal approximators of continuous permutation equivariant sequence-to-sequence functions with compact support. Our work focuses on MHSA. Guided by tensor diagram representation, we compare the expressive power of different tensor structures among the weight matrices in MHSA.

Efficient Transformers. MHSA has a global receptive field, i.e., each token’s representation is updated by attending to all other tokens, therefore incurring a quadratic memory and computation complexities concerning the sequence length. An extensive line of works focusing on relieving such the dependencies on the sequence length, such as Performer ([Choromanski et al., 2020](#)), Reformer ([Kitaev et al., 2020](#)), Linformer ([Wang et al., 2020](#)) and Linear Transformer ([Katharopoulos et al., 2020](#)). Our proposed model, Tuformers, can be incorporated into some kernel-based efficient Transformers, which applies kernel-based methods to MHSA to remove the nonlinearity constraint brought by the softmax function, achieving the same state-of-the-art linear computation and memory complexities in the sequence length.

6 CONCLUSION

This paper introduces a mathematically rigorous yet intuitive tensor diagram representation of MHSA, formulates a design space where we can analyze the expressive power of MHSA and its variants and proposes Tuformers, a novel model design with a guaranteed higher expressive power. Furthermore, Tuformers have a data-driven structure where heads can be trained implicitly and initialized with pre-trained Transformer models. The introduction of the tensor diagram representations and the theory presented in this paper provide new tools and open new directions for future research searching for expressive and efficient architectures.

ETHICS STATEMENT

Transformers are powerful neural network structures that achieve remarkable results in many domains. Understanding the expressive power of such neural networks is desirable but difficult. In this paper, we introduce tensor diagram representation of MHSA (the core component of Transformers) and formulate a design space where we can analyze the expressive power from a tensor perspective. Guided by the tensor diagram representation, we propose a plug-and-play THSA module to replace MHSA with THSA, obtaining Tuformers. Tuformers are theoretically guaranteed to have high expressive power than Transformers.

Our work benefits the Machine Learning community in following points:

- (1) We introduce tensor diagram representation to the ML community as an analytical tool for understanding expressive power of neural networks. Equipping Transformer models with such a precise and yet intuitive tool, we are able to close the gap between the currently used simple but intuitive flowcharts and the precise but complex code/mathematical formulation.
- (2) By tensorizing MHSA, we provide new possibilities for network analysis from a tensor representation perspective. We formulate a design space for MHSA, which can foster future studies for architecture design.
- (3) We explore the theory behind MHSA and propose Tuformers with THSA. The theoretical framework we provide can lead to future explorations and developments of expressive and yet efficient models.
- (4) Our proposed Tuformers are energy and resource efficient. (a) Tuformers can be initialized with pre-trained strong models such as BERT, RoBERTa and ALBERT models, drastically reducing the training cost for new downstream tasks. (b) THSA is a plug-and-play module that can be paired with ELU or polynomial kernels to obtain linear computation and memory efficiency in sequence length.

REPRODUCIBILITY STATEMENT

For theoretical results, full details are given in Appendix D.

For experimental results, firstly, all the datasets used in this paper are public and contain no sensitive information. Then a detailed description of the datasets, tasks, the hyper-parameters and how they are chosen, and train/test data splitting are in Appendix E. The main code can be found in the supplementary materials.

REFERENCES

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. In *International Conference on Learning Representations*, 2020.
- Nadav Cohen, Or Sharir, and Amnon Shashua. On the expressive power of deep learning: A tensor analysis. In *Conference on learning theory*, pp. 698–728. PMLR, 2016.
- Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-attention and convolutional layers. In *International Conference on Learning Representations*, 2019.
- Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. Multi-head attention: Collaborate instead of concatenate. *arXiv preprint arXiv:2006.16362*, 2020.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2978–2988, 2019.

- Lieven De Lathauwer. Decompositions of a higher-order tensor in block terms—part ii: Definitions and uniqueness. *SIAM Journal on Matrix Analysis and Applications*, 30(3):1033–1066, 2008.
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- W. Hackbusch and S. Kühn. A new scheme for the tensor representation. *J. Fourier Anal. Appl.*, 15(5):706–722, 2009. doi: 10.1007/s00041-009-9094-9.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pp. 5156–5165. PMLR, 2020.
- Valentin Khrulkov, Alexander Novikov, and Ivan Oseledets. Expressive power of recurrent neural networks. In *International Conference on Learning Representations*, 2018.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020.
- Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3): 455–500, 2009.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995.
- Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 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 pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Jerry Ma and Denis Yarats. On the adequacy of untuned warmup for adaptive optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- Xindian Ma, Peng Zhang, Shuai Zhang, Nan Duan, Yuexian Hou, Ming Zhou, and Dawei Song. A tensorized transformer for language modeling. *Advances in Neural Information Processing Systems*, 32:2232–2242, 2019.
- Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of english: the penn treebank. *Computational Linguistics*, 19(2):313–330, 1993.
- Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in Neural Information Processing Systems*, 32:14014–14024, 2019.
- Tomas Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur. Recurrent neural network based language model. In *INTERSPEECH*, 2010.

- Niko Moritz, Takaaki Hori, and Jonathan Le. Streaming automatic speech recognition with the transformer model. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6074–6078. IEEE, 2020.
- Ivan V Oseledets. Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5): 2295–2317, 2011.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pp. 48–53, 2019.
- Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 5206–5210. IEEE, 2015.
- Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International Conference on Machine Learning*, pp. 4055–4064. PMLR, 2018.
- Roger Penrose. Applications of negative dimensional tensors. *Combinatorial mathematics and its applications*, 1:221–244, 1971.
- Stephan Rabanser, Oleksandr Shchur, and Stephan Günnemann. Introduction to tensor decompositions and their applications in machine learning. *arXiv preprint arXiv:1711.10781*, 2017.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Edinburgh neural machine translation systems for wmt 16. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pp. 371–376, 2016.
- Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhutdinov. Transformer dissection: An unified understanding for transformer’s attention via the lens of kernel. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2019.
- 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*, pp. 5998–6008, 2017.
- Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5797–5808, 2019.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- 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)*, pp. 1112–1122. Association for Computational Linguistics, 2018. URL <http://aclweb.org/anthology/N18-1101>.
- Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? In *International Conference on Learning Representations*, 2019.

Appendix of Tuformer: Data-driven Design of Expressive Transformers by Tucker Tensor Representation

A SUPPLEMENTARY MATERIAL FOR TENSOR DIAGRAMS

In this section, we provide a comprehensive introduction to the tensor diagram representation.

A.1 COMPONENTS OF TENSOR DIAGRAM: ARRAYS AS NODES WITH LEG(S)

Arrays denoted as Nodes with Legs. An array is represented as a node with leg(s) in a tensor diagram as shown in Figure 1. The *order* (the number of dimensions) of the array is denoted by the number of legs extending from the node. The legs do not need to be straight lines, and the orientation does not matter. For example, matrices A and A^T are equivalent in this notation. Each labeled leg represents one *mode* of a tensor. Every mode of the tensor needs to be uniquely labeled. We usually use the dimension of the mode as the label (i.e., an associated positive integer written on top of each leg). If multiple modes of the same tensor have the same length, we use subscripts to differentiate them.

A.2 OPERATIONS IN TENSOR DIAGRAM

$$\begin{aligned} & \overset{M}{\text{A}} \overset{K}{\text{B}} \overset{N}{\text{C}} = \overset{M}{\text{C}} \overset{N}{\text{C}} \quad \overset{M}{\text{A}} \overset{\alpha}{\text{N}} = \overset{M}{\text{C}} \overset{N}{\text{C}} \quad \overset{M}{\text{A}} \overset{R}{\text{B}} \overset{N}{\text{C}} = \overset{M}{\text{C}} \overset{N}{\text{C}} \\ & \sum_{k=1}^K A_{mk} B_{kn} = C_{mn} \quad \frac{\exp(\alpha \cdot A_{mn})}{\sum_{n'=1}^N \exp(\alpha \cdot A_{mn'})} = C_{mn} \quad A_{nr} B_{mr} = C_{nmr} \end{aligned}$$

Figure 7: (1) Contraction as shown in Appendix A.2 is an operation that generalizes the matrix multiplication. It does summation on element-wise products along a mode in object 1 and a corresponding mode (with the same dimension) in object 2 (i.e., in tensor diagram language, along a leg in node 1 and a corresponding leg in node 2). **In tensor diagram, multiplying two matrices (or higher-order tensors with more than 2 legs) corresponds to “gluing” their corresponding legs (along a certain mode).** (2) **Softmax** as shown in Appendix A.2 is an element-wise exponential function normalized along a certain mode. Tensor diagram has a convention to denote the contraction, but not the softmax. We propose to denote the α -scaled softmax function $\text{softmax}(\alpha A)$ on A as a **dotted box with a labeled filled ball** (to distinguish itself from tensor objects, i.e., nodes which are blank circles) attached to one leg. (3) **Batch Multiplication**, as shown in appendix A.2, is an elementwise product along the connected legs.

There are three types of operations: *contraction*, *softmax* and *batch multiplication*. *Contraction* is an operation that does summation on element-wise products along a mode in object 1 and a corresponding mode (with the same dimension) in object 2 (i.e., in tensor diagram language, along a leg in node 1 and a corresponding leg in node 2), whereas *softmax* is an element-wise exponential function normalized along a certain mode. Tensor diagram has a convention to denote the contraction, but not the softmax. We will first introduce contraction in tensor diagram notation and then propose our design of softmax in tensor diagram language.

Contractions denoted as edges connecting the node legs. In tensor diagram, multiplying two matrices (or higher-order tensors with more than 2 legs) corresponds to “gluing” their corresponding legs (along a certain mode) as shown in Appendix A.2, and it is called *tensor contraction*. Since the representation of arrays in the tensor diagram is orientation invariant, contractions are also orientation invariant. This is especially useful in representing the self-attention unit. In some cases, the data object are represented as row vectors while in some other scenarios they are denoted as column vectors. If we use mathematical formula, we need to take care of the ordering of the matrix multiplication because $A^T B$, AB^T , BA^T , $B^T A$ and more are all different. In contrast, tensor diagram provides a universal graphical representation as long as the corresponding legs are connected correctly.

Nonlinear activation denoted as dotted box with a labelled filled ball attached. We propose to denote the α -scaled softmax function $\text{softmax}(\alpha\mathbf{A})$ on \mathbf{A} as a dotted box with a labeled filled ball (to distinguish itself from tensor objects, i.e., nodes which are blank circles) attached to one leg. The label right above the filled ball denote the scaling parameter α , and the leg that the ball is attached to indicate the mode where normalization is implemented along in the softmax operation. As shown in Appendix A.2, let \mathbf{A} denote the resultant of some operations (for example $\mathbf{Q}_{[h]}\mathbf{K}_{[h]}^\top$), the $\frac{1}{\sqrt{D}}$ -scaled softmax on \mathbf{A} which normalize along leg N is depicted. Softmax does not change the shape of the input. That is why we introduce a dotted box, which does not change the shape either, to denote it.

Batch Multiplication Given \mathbf{A} and \mathbf{B} , the batch multiplication is an element-wise product along the connected legs. The resultant tensor \mathbf{C} is obtained by merging connected nodes, maintaining the dangling legs M , N and R . Since in MHSA, the attention calculation among key, query and value matrices are done separately within each head, batch multiplication is a proper way to denote these “head-wise” operations.

A.3 ADVANTAGES OF TENSOR DIAGRAM

Tensor diagram representations enjoy a few advantages compared with using flow charts together with math formulas.

(1) The notation is orientation invariant. The input embedding matrix \mathbf{X} can be represented as either $\mathbb{R}^{N \times F}$ or $\mathbb{R}^{F \times N}$ with rows or columns being the embedding of each token in the sequence respectively. Although we take the former convention in our mathematical notations in Equation B.1, there are papers that use the latter, creating discrepancies between notations. However, in tensor diagram, since the node legs are orientation invariant, we obtain a universal graphical representation of multi-head self-attention, irrespective of how the input embedding matrix is represented, as shown in Figure 8c.

(2) The representation is both precise in math and intuitive, making analysis easier. The tensor diagram itself is mathematically rigorous.

With the labels of the legs, the model complexity (number of parameters) is explicitly displayed, requiring no supplementary information such as the size of the parameters \mathbf{W}^Q , \mathbf{W}^K and \mathbf{W}^V as in mathematical formulas. In addition, the resultant from an operation or a sequence of operations (i.e., a connected sub-graph of the entire tensor diagram graph) can be treated as “merging” any connected nodes in the sub-graph. This is particularly convenient to obtain the size of the resultant: any the dangling leg becomes a leg of the resultant.

(3) Multi-linear operations are concisely represented and interpreted. Using tensor diagram, we can concisely represent the multi-head self-attention rigorously in one diagram as shown in Figure 8c. More important, the multi-linear interaction between the latent features $\mathbf{Q}_{[h]}$, $\mathbf{K}_{[h]}$ and $\mathbf{V}_{[h]}$, which essentially attribute to the multi-linear interaction between the weight tensors \mathbf{W}^Q , \mathbf{W}^K and \mathbf{W}^B , is clearly illustrated in the tensor diagram. This illustration, which is not achieved by any existing representations, is crucial for our principled understanding of the multi-head self-attention.

(4) The information flow is clear. For example, in Figure 8a there is no direct information pass from $\mathbf{W}_{[h]}^K$, $\mathbf{W}_{[h]}^Q$ to $\mathbf{W}_{[h]}^V$ while in Figure 8c, such exchange is done through edge D s and H s. It is such observation (there are multiple paths that connect the weight parameter matrices) that leads to our idea of comparing expressive power between different structures.

B SUPPLEMENTARY MATERIAL FOR MULTI-HEAD SELF-ATTENTION

In the section, we prove the equivalence between the tensor diagram and the equations for *multi-head self-attention* (MHSA). To establish such an equivalence, we will need to show that both tensor diagram and the matrix equations lead to the same result in element-wise notation.

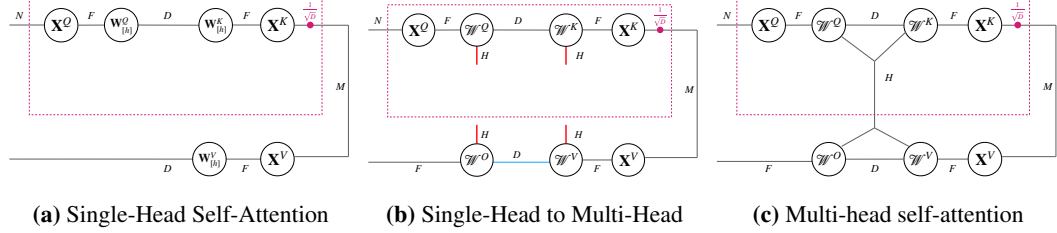


Figure 8: Tensor diagrams from single-head self-attention to multi-head self-attention. Figure 8a is the tensor diagram representation of a single-head self-attention. Figure 8c is the tensor-diagram representation of a multi-head self-attention.

For convenience of reference, we recap both representations in Equation (B.1) and Figure 8.

$$\mathbf{Q}_{[h]} = \mathbf{X}^Q \mathbf{W}_{[h]}^Q; \mathbf{K}_{[h]} = \mathbf{X}^K \mathbf{W}_{[h]}^K; \mathbf{V}_{[h]} = \mathbf{X}^V \mathbf{W}_{[h]}^V. \quad (\text{B.1a})$$

$$\text{head}_{[h]} = \text{softmax} \left(\mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top / \sqrt{D} \right) \mathbf{V}_{[h]}, \quad (\text{B.1b})$$

$$\mathbf{M} = [\text{head}_{[1]}, \text{head}_{[2]}, \dots, \text{head}_{[H]}] \mathbf{W}^O, \quad (\text{B.1c})$$

Note that the query weight tensor \mathbf{W}^Q comes from a concatenation of H query weight matrices $\{\mathbf{W}_{[h]}^Q\}_{h=1}^H$ such that $\mathbf{W}_{h,:,:}^Q = \mathbf{W}_{[h]}^Q$. Similarly, $\mathbf{W}^K, \mathbf{W}^V$ are concatenated from $\{\mathbf{W}_{[h]}^K\}_{h=1}^H, \{\mathbf{W}_{[h]}^V\}_{h=1}^H$ such that $\mathbf{W}_{h,:,:}^K = \mathbf{W}_{[h]}^K, \mathbf{W}_{h,:,:}^V = \mathbf{W}_{[h]}^V$. Furthermore, we denote $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ as the reshaped tensors from $\{\mathbf{Q}_{[h]}\}_{h=1}^H, \{\mathbf{K}_{[h]}\}_{h=1}^H, \{\mathbf{V}_{[h]}\}_{h=1}^H$. Finally, \mathbf{W}^O is reshaped from \mathbf{W}^O .

To differentiate between equations and tensor diagrams in the proof, we add an over-line to every object in the tensor diagram throughout this section.

B.1 SINGLE-HEAD SELF-ATTENTION (SHSA)

We will prove that the tensor diagram in Figure 8a is equivalent to Equations (B.1a) and (B.1b).

Proof. For simplicity, we omit the subscript $[h]$ in this subsection.

We first show that a contraction of the nodes within the softmax box is equivalent to $\mathbf{Q}\mathbf{K}^\top$. Let $\mathbf{A} = \mathbf{Q}\mathbf{K}^\top = \mathbf{X}^Q \mathbf{W}^Q (\mathbf{W}^K \mathbf{X}^K)^\top$, where $\mathbf{X}^Q \in \mathbb{R}^{N \times F}$, $\mathbf{X}^K \in \mathbb{R}^{M \times F}$, $\mathbf{W}^Q, \mathbf{W}^K \in \mathbb{R}^{F \times D}$. We use subscripts in the element-wise notation to differentiate edges when multiple edges with same dimension appear in the same equation.

$$\mathbf{Q}_{nd} = \sum_{f_1=1}^F \mathbf{X}_{nf_1}^Q \mathbf{W}_{f_1d}^Q, \quad (\text{B.2a})$$

$$\mathbf{K}_{md} = \sum_{f_2=1}^F \mathbf{X}_{mf_2}^K \mathbf{W}_{f_2d}^K, \quad (\text{B.2b})$$

$$\mathbf{A}_{nm} = \sum_{d=1}^D \mathbf{Q}_{nd} \mathbf{K}_{dm} = \sum_{d=1}^D \sum_{f_1=1}^F \sum_{f_2=1}^F \mathbf{X}_{nf_1}^Q \mathbf{W}_{f_1d}^Q \mathbf{W}_{f_2d}^K \mathbf{X}_{mf_2}^K. \quad (\text{B.2c})$$

Denote be the result of the subset of nodes in the softmax box as \mathbf{A} .

$$\overline{\mathbf{A}}_{nm} = \sum_{d=1}^D \sum_{f_1=1}^F \sum_{f_2=1}^F \overline{\mathbf{X}}_{nf_1}^Q \overline{\mathbf{W}}_{f_1d}^Q \overline{\mathbf{W}}_{f_2d}^K \overline{\mathbf{X}}_{mf_2}^K. \quad (\text{B.3})$$

Comparing Equations (B.2c) and (B.3), we have $\mathbf{A} = \overline{\mathbf{A}}$. Let $\mathbf{B} = \text{softmax}(\mathbf{A})$ and $\overline{\mathbf{B}} = \text{softmax}(\overline{\mathbf{A}})$, we further have $\mathbf{B} = \text{softmax}(\mathbf{A}) = \text{softmax}(\overline{\mathbf{A}}) = \overline{\mathbf{B}}$.

Let $\overline{\mathbf{M}}$ be the result of the tensor diagram and \mathbf{M} be the matrix representation of the single-head attention. We have

$$\mathbf{M}_{nd} = \sum_{m=1}^N \sum_{f=1}^F \mathbf{B}_{nm} \mathbf{X}_{mf}^V \mathbf{W}_{fd}^V. \quad (\text{B.4})$$

$$\overline{\mathbf{M}}_{nd} = \sum_{m=1}^N \sum_{f=1}^F \overline{\mathbf{B}}_{nm} \overline{\mathbf{X}}_{mf}^V \overline{\mathbf{W}}_{fd}^V. \quad (\text{B.5})$$

Comparing Equations (B.4) and (B.5), we have \mathbf{C} is equivalent to $\overline{\mathbf{C}}$, which completes the proof. \square

B.2 MULTI-HEAD SELF-ATTENTION (MHSA)

We now prove that tensor diagram in Figure 8c is equivalent to Equation (1c).

Proof. Let H be the number of heads. Let $\overline{\mathbf{A}} \in \mathbb{R}^{H \times N \times M}$ be the result of a contraction of the nodes within the softmax box. From the last section, we know that $\overline{\mathbf{A}}_{h, :, :} = \mathbf{A}_{[h]} = \mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top$. Applying the softmax function $\overline{\mathbf{B}} = \text{softmax}(\overline{\mathbf{A}})$, we have $\overline{\mathbf{B}}_{h, :, :} = \text{softmax}(\mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top)$.

Let $\overline{\mathbf{T}}$ be result of the tensor diagram and \mathbf{T} be that of Equation (B.1c).

$$\mathbf{T}_{nf_4} = \sum_{m, h, d_1, d_2} \text{softmax}(\mathbf{Q}_{[h], nd_1} \mathbf{Q}_{[h], md_1}) \mathbf{V}_{[h], md_2} \mathbf{W}_{[h], d_2 f_4}^O. \quad (\text{B.6})$$

$$\overline{\mathbf{T}}_{n1f_4} = \sum_{m, h, d_1, d_2} \text{softmax}(\overline{\mathbf{Q}}_{hnd_1} \overline{\mathbf{K}}_{hmd_1}) \overline{\mathbf{V}}_{hmd_2} \overline{\mathbf{W}}_{hmd_2}^O. \quad (\text{B.7})$$

Comparing $\overline{\mathbf{T}}$ and \mathbf{T} , we complete the proof. \square

C CLOSER INVESTIGATION OF MHSA AND ITS COMPARISON TO THSA

Rethinking MHSA. There are two types contractions in MHSA as shown in Figure 4a: *Head-Contraction*, a global contraction among all weights \mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V and \mathbf{W}^O along mode H , and *Latent-Contraction*, a local contraction between \mathbf{W}^Q and \mathbf{W}^K as well as between \mathbf{W}^V and \mathbf{W}^O both along mode D . If we remove the Head-Contraction and maintain the Latent-Contraction only, multi-head is reduced to single-head with an additional contraction with the weight matrix \mathbf{W}^O . To understand the role of the two types of contractions, we propose ablation studies by removing one of them and study the changes in the structure and the expressive power. It turns out that in both situation, *MHSA falls into a special case of a variant of Tucker Tensor Representation* with the core tensor taking different values. Note that the core tensor is of order 2, which means it is a matrix.

Repeated Proposition 1 with two special cases. Given $R = DH$ and other hyperparameters being the same, THSA module includes MHSA as a special case. Specifically,

- (1) If $\mathbf{C} = \mathbf{1}_D \mathbf{1}_D^\top$, i.e., an $D \times D$ all-ones matrix, the THSA module reduces to an MHSA with a single head and latent dimension D (i.e., a single-head self-attention with latent dimension D).
- (2) If $\mathbf{C} = \mathbf{I}_H$, i.e., an $H \times H$ identity matrix, the THSA reduces to a MHSA with H heads and latent dimension 1 (i.e., a heads-only self-attention with H heads).
- (3) If $\mathbf{C} = \mathbf{I}_H \otimes (\mathbf{1}_D \mathbf{1}_D^\top)$, i.e., a Kronecker product of an all-ones matrix $\mathbf{1}_D \mathbf{1}_D^\top \in \mathbb{R}^{D \times D}$ and an identity matrix $\mathbf{I}_H \in \mathbb{R}^{H \times H}$, the THSA reduces to an MHSA with H heads and latent dimension D .

The proof of the above propositions is given in Appendix D.

Definition 4 (Stable Rank). Let $\mathbf{A} \in \mathbb{R}^{N \times M}$, the stable rank of \mathbf{A} of is defined as $\sum_i \sigma_i^2 / \max_i \sigma_i$ where σ_i is the singular value of \mathbf{A} .

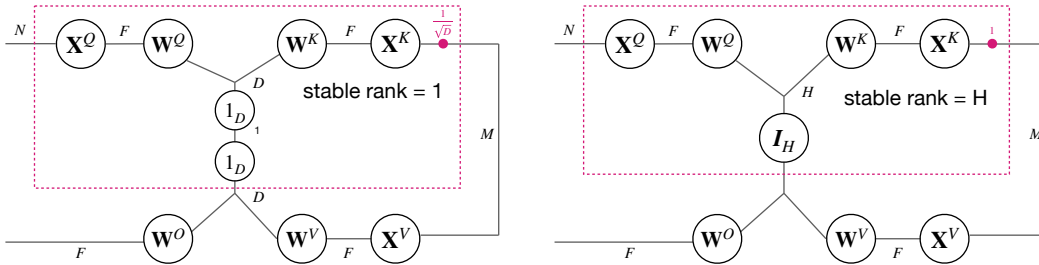


Figure 9: (Left) THSA includes a single-head self-attention. In this case, the stable rank of the core tensor C is 1. It is interesting to observe that this is a single-head self-attention and this structure can be obtained by removing the head contraction from Figure 4a. (Right) THSA includes a MHSA with H heads and latent dimension = 1. In this case, the stable rank equals to the rank of the identical matrix, which is H , corresponding to the number of heads of MHSA. Also, the structure on the right can be obtained by removing the latent contraction edges D from Figure 4a.

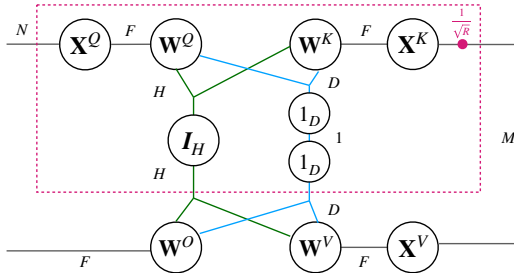


Figure 10: This figure shows how to initialize THSA with MHSA. MHSA is a special case of THSA when the core tensor C takes the Kronecker product of a $D \times D$ all-one matrix and a $H \times H$ identical matrix. Note that in this case, the number of heads H in MHSA is equivalent to the stable rank of C .

The notion of heads in THSA. The concept of the number of heads in THSA generalizes to the stable rank of the core matrix, allowing a data-driven implicit training. The number of heads H in MHSA corresponds to the *stable rank* of the core matrix C , defined as $\sum_i \sigma_i^2 / \max_i \sigma_i$ where σ_i is the singular value. It is a useful surrogate for the rank because it is largely unaffected by tiny singular values.

As shown in Appendix C when C is constructed by the outer product of two vectors, the stable rank of C is 1, while in this case, the whole structure is equivalent to a single-head self-attention. When C is an identical matrix as shown in appendix C, the stable rank of C equals to the rank of the identical matrix, which is H , corresponding to the number of heads of MHSA. Also, the structure on the right can be obtained by removing the latent contraction edges D from Figure 4a. When C is the Kronecker product of the all-one matrix and the identical matrix, which corresponds to MHSA, the stable rank of C is H , and this number does not change in the training because that C is a fixed matrix in MHSA.

D SUPPLEMENTARY MATERIAL FOR TUFORMER AND PROOFS

In this section, we will prove Proposition 1 and Theorem 3 in Section 3. In both theorems, we assume the rank in *Tucker-head self-attention* (THSA) equals to the product of the number of heads H and the latent dimension D in *multi-head self-attention* (MHSA), i.e., $R = DH$.

For convenience of reference, we recap the mathematical expressions and tensor diagrams for both multi-head self-attention (MHSA) and Tucker-head self-attention (THSA) in the following.

Multi-head Self-Attention (MHSA). An MHSA module has its learnable parameters as H sets of (query, key, value) weight matrices $\{\mathbf{W}_{[h]}^Q, \mathbf{W}_{[h]}^K, \mathbf{W}_{[h]}^V\}_{h=1}^H$ (with $\mathbf{W}_{[h]}^Q, \mathbf{W}_{[h]}^K, \mathbf{W}_{[h]}^V \in \mathbb{R}^{F \times D}$ for each h) and an output weight matrix $\mathbf{W}^O \in \mathbb{R}^{HD \times F}$. The module maps three input matrices $\mathbf{X}^Q \in \mathbb{R}^{N \times F}$, $\mathbf{X}^K, \mathbf{X}^V \in \mathbb{R}^{M \times F}$ into an output matrix $\mathbf{M} \in \mathbb{R}^{N \times F}$.

$$\mathbf{Q}_{[h]} = \mathbf{X}^Q \mathbf{W}_{[h]}^Q; \mathbf{K}_{[h]} = \mathbf{X}^K \mathbf{W}_{[h]}^K; \mathbf{V}_{[h]} = \mathbf{X}^V \mathbf{W}_{[h]}^V, \quad (\text{D.1a})$$

$$\mathbf{head}_{[h]} = \text{softmax}\left(\mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top\right) \mathbf{V}_{[h]}, \quad (\text{D.1b})$$

$$\mathbf{M} = [\mathbf{head}_{[1]}, \mathbf{head}_{[2]}, \dots, \mathbf{head}_{[H]}] \mathbf{W}^O. \quad (\text{D.1c})$$

In Equation (D.1a), $\mathbf{Q}_{[h]} \in \mathbb{R}^{N \times D}$, $\mathbf{K}_{[h]}, \mathbf{V}_{[h]} \in \mathbb{R}^{M \times D}$ are query, key, value matrices respectively. In Equation (D.1b), we omit the scaling factor $1/\sqrt{D}$ for the softmax function (the scalar can be merged into the learnable parameters), and $\mathbf{head}_{[h]} \in \mathbb{R}^{N \times D}$ is the resulted matrix for head h .

Tucker-Head Self-Attention. A THSA module is parameterized by five matrices: a query weight matrix $\overline{\mathbf{W}}^Q \in \mathbb{R}^{F \times R}$, a key weight matrix $\overline{\mathbf{W}}^K \in \mathbb{R}^{F \times R}$, a value weight matrix $\overline{\mathbf{W}}^V \in \mathbb{R}^{F \times R}$, an output weight matrix $\overline{\mathbf{W}}^O \in \mathbb{R}^{R \times F}$, and an additional core matrix $\mathbf{C} \in \mathbb{R}^{R \times R}$. The module has the same input and output domains as in MHSA, i.e., it takes three matrices $\mathbf{X}^Q \in \mathbb{R}^{N \times F}$, $\mathbf{X}^K, \mathbf{X}^V \in \mathbb{R}^{M \times F}$ as inputs and returns a matrix $\mathbf{T} \in \mathbb{R}^{N \times F}$ as output.

$$\overline{\mathbf{Q}} = \mathbf{X}^Q \overline{\mathbf{W}}^Q; \overline{\mathbf{K}} = \mathbf{X}^K \overline{\mathbf{W}}^K; \overline{\mathbf{V}} = \mathbf{X}^V \overline{\mathbf{W}}^V, \quad (\text{D.2a})$$

$$\overline{\mathbf{head}}_r = \text{softmax}\left(\sum_{s=1}^R \mathbf{C}_{rs} \overline{\mathbf{Q}}_s \overline{\mathbf{K}}_s^\top\right) \overline{\mathbf{V}}_r, \quad (\text{D.2b})$$

$$\mathbf{T} = \overline{\mathbf{head}} \overline{\mathbf{W}}^O = [\overline{\mathbf{head}}_1, \overline{\mathbf{head}}_2, \dots, \overline{\mathbf{head}}_R]. \quad (\text{D.2c})$$

To distinguish THSA from MHSA, we use over-scored symbols whenever needed. In Equation (D.2a), $\overline{\mathbf{Q}} \in \mathbb{R}^{N \times R}$, $\overline{\mathbf{K}}, \overline{\mathbf{V}} \in \mathbb{R}^{M \times R}$ are query, key, value matrices respectively. In Equation (D.2b), we again omit the scaling factor $1/\sqrt{R}$ in the softmax function. Equation (D.2b) leads to a head matrix $\overline{\mathbf{head}} \in \mathbb{R}^{F \times R}$, where $\overline{\mathbf{head}}_r \in \mathbb{R}^F$ denotes the r^{th} column of the matrix $\overline{\mathbf{head}}$.

D.1 PROOF OF PROPOSITION 1

Proof of Proposition 1. We constructively prove that THSA reduces to three cases of MHSA if the query/key/value weight matrices in THSA concatenate the corresponding matrices in MHSA:

$$\overline{\mathbf{W}}^Q = [\mathbf{W}_{[1]}^Q, \mathbf{W}_{[2]}^Q, \dots, \mathbf{W}_{[H]}^Q], \quad (\text{D.3a})$$

$$\overline{\mathbf{W}}^K = [\mathbf{W}_{[1]}^K, \mathbf{W}_{[2]}^K, \dots, \mathbf{W}_{[H]}^K], \quad (\text{D.3b})$$

$$\overline{\mathbf{W}}^V = [\mathbf{W}_{[1]}^V, \mathbf{W}_{[2]}^V, \dots, \mathbf{W}_{[H]}^V]. \quad (\text{D.3c})$$

In addition, we set $\overline{\mathbf{W}}^O = \mathbf{W}^O$. We have $\overline{\mathbf{W}}^Q, \overline{\mathbf{W}}^K, \overline{\mathbf{W}}^V \in \mathbb{R}^{F \times R}$, and $\overline{\mathbf{W}}^O \in \mathbb{R}^{R \times F}$.

(1) For $\mathbf{C} = \mathbf{1}_R \mathbf{1}_R^\top$, i.e. an $R \times R$ all-ones matrix, we aim to prove that the THSA module reduces to a single-head self-attention with $H = 1$ and $D = R$. Since there is only one head in the module, we have $\overline{\mathbf{W}}^Q = \mathbf{W}_{[1]}^Q$, $\overline{\mathbf{W}}^K = \mathbf{W}_{[1]}^K$, and $\overline{\mathbf{W}}^V = \mathbf{W}_{[1]}^V$. As an immediate result,

$$\overline{\mathbf{Q}} = \mathbf{X}^Q \overline{\mathbf{W}}^Q = \mathbf{X}^Q \mathbf{W}_{[1]}^Q = \mathbf{Q}_{[1]}, \quad (\text{D.4a})$$

$$\overline{\mathbf{K}} = \mathbf{X}^K \overline{\mathbf{W}}^K = \mathbf{X}^K \mathbf{W}_{[1]}^K = \mathbf{K}_{[1]}, \quad (\text{D.4b})$$

$$\overline{\mathbf{V}} = \mathbf{X}^V \overline{\mathbf{W}}^V = \mathbf{X}^V \mathbf{W}_{[1]}^V = \mathbf{V}_{[1]}. \quad (\text{D.4c})$$

Since \mathbf{C} is an all-ones matrix, i.e., $\mathbf{C}_{rs} = 1, \forall r, s$, we rewrite Equation (D.2b) as:

$$\overline{\mathbf{head}}_r = \text{softmax}\left(\sum_{s=1}^R \overline{\mathbf{Q}}_s \overline{\mathbf{K}}_s^\top\right) \overline{\mathbf{V}}_r = \text{softmax}\left(\overline{\mathbf{Q}} \overline{\mathbf{K}}^\top\right) \overline{\mathbf{V}}_r. \quad (\text{D.5})$$

Notice that the equation holds for each column r , we further write all R equations jointly as:

$$\overline{\mathbf{head}} = \text{softmax}\left(\overline{\mathbf{Q}} \overline{\mathbf{K}}^\top\right) \overline{\mathbf{V}} = \text{softmax}\left(\mathbf{Q}_{[1]} \mathbf{K}_{[1]}^\top\right) \mathbf{V}_{[1]} = \mathbf{head}_{[1]}. \quad (\text{D.6})$$

Finally, we express the output matrix as:

$$\mathbf{T} = \overline{\mathbf{head}} \overline{\mathbf{W}}^O = \mathbf{head}_{[1]} \mathbf{W}^O = \mathbf{M}, \quad (\text{D.7})$$

which concludes the reduction for the special case (1).

(2) For $\mathbf{C} = \mathbf{I}_R$, an $R \times R$ identity matrix, we aim to prove that the THSA module reduces to a heads-only self-attention with $H = R$ and $D = 1$. Since the latent dimension $D = 1$, i.e., each $\mathbf{W}_{[h]}^Q$ (or $\mathbf{W}_{[h]}^K, \mathbf{W}_{[h]}^V$) is a vector, we have $\overline{\mathbf{W}}_r^Q = \mathbf{W}_r^Q$, $\overline{\mathbf{W}}_r^K = \mathbf{W}_r^K$, and $\overline{\mathbf{W}}_r^V = \mathbf{W}_r^V$. Therefore,

$$\overline{\mathbf{Q}}_r = \mathbf{X}^Q \overline{\mathbf{W}}_r^Q = \mathbf{X}^Q \mathbf{W}_r^Q = \mathbf{Q}_{[r]}, \quad (\text{D.8a})$$

$$\overline{\mathbf{K}}_r = \mathbf{X}^K \overline{\mathbf{W}}_r^K = \mathbf{X}^K \mathbf{W}_r^K = \mathbf{K}_{[r]}, \quad (\text{D.8b})$$

$$\overline{\mathbf{V}}_r = \mathbf{X}^V \overline{\mathbf{W}}_r^V = \mathbf{X}^V \mathbf{W}_r^V = \mathbf{V}_{[r]}. \quad (\text{D.8c})$$

Since \mathbf{C} is an identity matrix, i.e., $\mathbf{C}_{rr} = 1$ and $\mathbf{C}_{rs} = 0, \forall r \neq s$, we rewrite Equation (D.2b) as:

$$\overline{\mathbf{head}}_r = \text{softmax}\left(\overline{\mathbf{Q}}_r \overline{\mathbf{K}}_r^\top\right) \overline{\mathbf{V}}_r = \text{softmax}\left(\mathbf{Q}_{[r]} \mathbf{K}_{[r]}^\top\right) \mathbf{V}_{[r]} = \mathbf{head}_{[r]}. \quad (\text{D.9})$$

Finally, we express the output matrix as:

$$\mathbf{T} = \overline{\mathbf{head}} \overline{\mathbf{W}}^O = [\mathbf{head}_{[1]}, \mathbf{head}_{[2]}, \dots; \mathbf{head}_{[R]}] \mathbf{W}^O = \mathbf{M}, \quad (\text{D.10})$$

which concludes the reduction for the special case (2).

(3) For $\overline{\mathbf{C}} = \mathbf{I}_H \otimes (\mathbf{1}_D \mathbf{1}_D^\top)$, i.e., a Kronecker product between an identity matrix $\mathbf{I}_H \in \mathbb{R}^{H \times H}$ and an all-ones matrix $\mathbf{1}_D \mathbf{1}_D^\top \in \mathbb{R}^{D \times D}$, we aim to prove that the THSA module reduces to a multi-head self-attention (MHSA) module with H heads and latent dimension D . In this general case, we have

$$\overline{\mathbf{Q}} = \mathbf{X}^Q \overline{\mathbf{W}}^Q = \mathbf{X}^Q \left[\mathbf{W}_{[1]}^Q, \mathbf{W}_{[2]}^Q, \dots, \mathbf{W}_{[H]}^Q \right] = [\mathbf{Q}_{[1]}, \mathbf{Q}_{[2]}, \dots, \mathbf{Q}_{[H]}] \quad (\text{D.11})$$

$$\overline{\mathbf{K}} = \mathbf{X}^K \overline{\mathbf{W}}^K = \mathbf{X}^K \left[\mathbf{W}_{[1]}^K, \mathbf{W}_{[2]}^K, \dots, \mathbf{W}_{[H]}^K \right] = [\mathbf{K}_{[1]}, \mathbf{K}_{[2]}, \dots, \mathbf{K}_{[H]}] \quad (\text{D.12})$$

$$\overline{\mathbf{V}} = \mathbf{X}^V \overline{\mathbf{W}}^V = \mathbf{X}^V \left[\mathbf{W}_{[1]}^V, \mathbf{W}_{[2]}^V, \dots, \mathbf{W}_{[H]}^V \right] = [\mathbf{V}_{[1]}, \mathbf{V}_{[2]}, \dots, \mathbf{V}_{[H]}] \quad (\text{D.13})$$

Notice that $\overline{\mathbf{C}}$ is a block diagonal matrix $\text{blkdiag}(\mathbf{1}_D \mathbf{1}_D^\top, \dots, \mathbf{1}_D \mathbf{1}_D^\top)$, where each block is an all-ones matrix. Therefore we can rewrite Equation (D.2b) for $r = (h-1)D + d$ as:

$$\overline{\mathbf{head}}_{(h-1)D+d} = \text{softmax}\left(\sum_{s=1}^R \mathbb{1}_{(h-1)D < s \leq hD} \left(\overline{\mathbf{Q}}_s \overline{\mathbf{K}}_s^\top\right)\right) \overline{\mathbf{V}}_r \quad (\text{D.14})$$

$$= \text{softmax}\left(\sum_{d=1}^D \mathbf{Q}_{[h],d} \mathbf{K}_{[h],d}^\top\right) \mathbf{V}_{[h],d} \quad (\text{D.15})$$

$$= \text{softmax}\left(\mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top\right) \mathbf{V}_{[h],d}. \quad (\text{D.16})$$

Since Equation (D.16) holds for each column d , we write all D equations jointly as:

$$[\overline{\mathbf{head}}_{(h-1)D+d}, \dots, \overline{\mathbf{head}}_{hD}] = \text{softmax}\left(\mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top\right) \mathbf{V}_{[h]} = \mathbf{head}_{[h]}. \quad (\text{D.17})$$

Again, since Equation (D.17) holds for each block h , we combine all H equations as:

$$\overline{\mathbf{head}} = [\mathbf{head}_{[1]}, \mathbf{head}_{[2]}, \dots, \mathbf{head}_{[H]}]. \quad (\text{D.18})$$

Finally, we express the output matrix as:

$$\mathbf{T} = \overline{\mathbf{head}} \overline{\mathbf{W}}^O = [\mathbf{head}_{[1]}, \mathbf{head}_{[2]}, \dots, \mathbf{head}_{[H]}] \mathbf{W}^O = \mathbf{M}, \quad (\text{D.19})$$

which concludes the reduction for the case (3). \square

D.2 PROOF OF THEOREM 3

Proof of Theorem 3. We have shown that THSA reduces to MHSA when its Tucker core \mathbf{C} takes specific forms. This shows that THSA is more expressive than MHSA — any mapping by an MHSA module can also be realized by a THSA module.

To prove that THSA is *strictly* more expressive than MHSA, we need find a mapping by a THSA module that can not be realized by an MHSA module. It suffices to show that the matrix \mathbf{M} is rank-deficient, while the matrix \mathbf{T} can be full-rank. We assume $MH < \min(N, F)$.

For convenience, we divide the matrix \mathbf{W}^O for MHSA into H blocks.

$$\mathbf{W}^O = \left[\mathbf{W}_{[1]}^{O\top}, \mathbf{W}_{[2]}^{O\top}, \dots, \mathbf{W}_{[H]}^{O\top} \right]^\top, \quad (\text{D.20})$$

where $\mathbf{W}_{[h]}^O \in \mathbb{R}^{D \times F}$ is the matrix for the h^{th} head. We now rewrite the output matrix \mathbf{M} as:

$$\mathbf{M} = \sum_{h=1}^H \underbrace{\text{softmax} \left(\mathbf{Q}_{[h]} \mathbf{K}_{[h]}^\top \right)}_{\mathbb{R}^{N \times M}} \underbrace{\left(\mathbf{V}_{[h]} \mathbf{W}_{[h]}^O \right)}_{\mathbb{R}^{M \times F}}. \quad (\text{D.21})$$

Since the matrix inside summation is a product of two matrices of size $\mathbb{R}^{N \times M}$ and $\mathbb{R}^{M \times F}$, its rank is at most M . Use the property $\text{rank}(\mathbf{A} + \mathbf{B}) \leq \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B})$, we have $\text{rank}(\mathbf{M}) \leq MH < \min(N, F)$. Similarly, we divide the matrix $\overline{\mathbf{W}}^O$ for THSA into R rows and rewrite the output matrix \mathbf{T} as:

$$\mathbf{T} = \sum_{r=1}^R \underbrace{\text{softmax} \left(\sum_{s=1}^R \mathbf{C}_{rs} \overline{\mathbf{Q}}_{:,r} \overline{\mathbf{K}}_{:,r}^\top \right)}_{\mathbb{R}^{N \times M}} \underbrace{\left(\overline{\mathbf{V}}_{:,r} \overline{\mathbf{W}}_{r,:}^O \right)}_{\mathbb{R}^{M \times F}}. \quad (\text{D.22})$$

Using the same argument, we have $\text{rank}(\mathbf{T}) \leq MR$. Since $H < R$ (in fact, $R = DH$), we can always find an example such that $\text{rank}(\mathbf{T}) > \text{rank}(\mathbf{M})$. This completes the proof. \square

E SUPPLEMENTARY MATERIAL FOR EXPERIMENTS

All the experiments are run on computing nodes with 4 NVIDIA GeForce GPUs.

(1) Language Modeling on Penn Treebank. Language modeling is the task of computing the probability of a sentence or a sequence of words. The model performance is measured by *per-word Perplexity (PPL)* which is the lower the better.

Penn Treebank (Marcus et al., 1993) is under the LDC User Agreement for Non-Members. We adopt the default settings as in the vanilla Transformer models (Vaswani et al., 2017). Specifically, we have N (length of the sequence) to be 2048, F (length of the initial embedding) to be 512, 8 number of heads, 6 layers encoder-decoder structure, D (latent feature dimension) to be 64. We use Adam optimizer and have 4000 warm-up steps. We also use label smoothing(0.1). The model is trained for 50 epochs.

(2) Neural Machine Translation on WMT16. The goal of NMT is to generate a corresponding sequence in one language given a sequence in another. The performance is measured by *BLEU* scores. (the higher, the better).

WMT16 English-German dataset (Sennrich et al., 2016) is under the MIT license. In this experiment, we use a Transformer model with 8 layers and we have $F = 512$ and $R = 1024$. The learning rate is $1e^{-5}$ and 4000 warmups. We also set label-smoothing to be 0.1. To speed up the experiment, we use the mix-precision trick with an optimization level being O1.

(3) Image Generation on MNIST and CIFAR-10. The image generation task is to predict an image pixel by pixel. The performance is evaluated by *bits per dimension (BPD)*. The image generation task is chosen to demonstrate, firstly, that Tuformers work beyond the language domain, and secondly, Tuformers can improve other efficient designs.

CIFAR10 (Krizhevsky, 2009) dataset is under the MIT license (MIT), MNIST (LeCun et al., 2010) is under the Creative Commons Attribution-Share Alike 3.0 license. The reason why we choose to evaluate Tuformer in image generation tasks on these two datasets is that we want to show the compatibility of our model to other efficient designs. Katharopoulos et al. (2020) introduces a linear transformer that has a state-of-the-art complexity. As a result, we evaluate our design on the same experiments as they do in their paper. We adopt basically the same settings except for GPU resources. Concretely, we use a 8 layer transformer model with 8 heads. The embedding size F is 256 and D is 32. N is set to be 1024. We also use “RAdam” (Ma & Yarats, 2021) optimizer for a state-of-the-art result. The model is trained for 250 epochs for MNIST dataset and 50 epochs for CIFAR10 dataset. The batch size is 10 in MNIST training and 4 in CIFAR training. We additionally use the mix-precision trick with an optimization level set to be O1 to speed up the training.

(4) Automatic Speech Recognition on LibriSpeech Dataset ASR consists of transcribing audio speech segments into text. The performance is measured by *Word Error Rate (WER)*. (the lower, the better). We adopt the default settings as in Fairseq (Ott et al., 2019), except that we replace MHSA with THSA.

(5) Natural Language Inference on MNLI and QNLI. NLI is a task of determining whether the given “hypothesis” and “premise” logically follow (entailment) or unfollow (contradiction) or are undetermined (neutral) to each other. The results are evaluated by accuracy. Using MNLI dataset, the model tries to predict whether sentence A entails or contradicts B while in QNLI dataset, the model is trained to answer whether sentence B contains answers to the question in sentence A. We use a batch size of 32 for 10 epochs, with a learning rate $1e^{-5}$.