

# A PRIMAL-DUAL FRAMEWORK FOR TRANSFORMERS AND NEURAL NETWORKS

**Tan M. Nguyen\***

Department of Mathematics  
University of California, Los Angeles  
tanmnguyen89@ucla.edu

**Tam Nguyen\***

Department of ECE  
Rice University  
nguyenminhtam9520@gmail.com

**Nhat Ho**

Department of Statistics & Data Sciences  
University of Texas at Austin  
minhnhat@utexas.edu

**Andrea L. Bertozzi**

Department of Mathematics  
University of California, Los Angeles  
bertozzi@math.ucla.edu

**Richard G. Baraniuk\*\***

Department of ECE  
Rice University  
richb@rice.edu

**Stanley J. Osher\*\***

Department of Mathematics  
University of California, Los Angeles  
sjo@math.ucla.edu

## ABSTRACT

Self-attention is key to the remarkable success of transformers in sequence modeling tasks including many applications in natural language processing and computer vision. Like neural network layers, these attention mechanisms are often developed by heuristics and experience. To provide a principled framework for constructing attention layers in transformers, we show that the self-attention corresponds to the support vector expansion derived from a support vector regression problem, whose primal formulation has the form of a neural network layer. Using our framework, we derive popular attention layers used in practice and propose two new attentions: 1) the Batch Normalized Attention (Attention-BN) derived from the batch normalization layer and 2) the Attention with Scaled Head (Attention-SH) derived from using less training data to fit the SVR model. We empirically demonstrate the advantages of the Attention-BN and Attention-SH in reducing head redundancy, increasing the model’s accuracy, and improving the model’s efficiency in a variety of practical applications including image and time-series classification.

## 1 INTRODUCTION

Transformer models (Vaswani et al., 2017) have achieved impressive success with state-of-the-art performance in a myriad of sequence processing tasks, including those in computer vision (Dosovitskiy et al., 2021; Liu et al., 2021; Touvron et al., 2020; Ramesh et al., 2021; Radford et al., 2021; Arnab et al., 2021; Liu et al., 2022; Zhao et al., 2021; Guo et al., 2021), natural language processing (Devlin et al., 2018; Al-Rfou et al., 2019; Dai et al., 2019; Child et al., 2019; Raffel et al., 2020; Baevski & Auli, 2019; Brown et al., 2020; Dehghani et al., 2018), reinforcement learning (Chen et al., 2021; Janner et al., 2021), and other important applications (Rives et al., 2021; Jumper et al., 2021; Zhang et al., 2019; Gulati et al., 2020; Wang & Sun, 2022). Transformers can also effectively transfer knowledge from pre-trained models to new tasks with limited supervision (Radford et al., 2018; 2019; Devlin et al., 2018; Yang et al., 2019; Liu et al., 2019). The driving force behind the success of transformers is the self-attention mechanism (Cho et al., 2014; Parikh et al., 2016; Lin et al., 2017), which computes a weighted average of feature representations of the tokens in the sequence with the weights proportional to similarity scores between pairs of representations. The weights calculated by the self-attention determine the relative importance between tokens and thus capture the contextual representations of the sequence (Bahdanau et al., 2014; Vaswani et al., 2017; Kim et al., 2017). It has

\* Co-first authors. \*\* Co-last authors. Please correspond to: tanmnguyen89@ucla.edu

been argued that the flexibility in capturing diverse syntactic and semantic relationships is critical for the success of transformers (Tenney et al., 2019; Vig & Belinkov, 2019; Clark et al., 2019).

### 1.1 BACKGROUND: SELF-ATTENTION

For a given input sequence  $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D_x}$  of  $N$  feature vectors, self-attention transforms  $\mathbf{X}$  into the output sequence  $\mathbf{H}$  in the following two steps:

**Step 1.** The input sequence  $\mathbf{X}$  is projected into the query matrix  $\mathbf{Q}$ , the key matrix  $\mathbf{K}$ , and the value matrix  $\mathbf{V}$  via three linear transformations

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q^\top; \mathbf{K} = \mathbf{X}\mathbf{W}_K^\top; \mathbf{V} = \mathbf{X}\mathbf{W}_V^\top,$$

where  $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{D \times D_x}$ , and  $\mathbf{W}_V \in \mathbb{R}^{D_v \times D_x}$  are the weight matrices. We denote  $\mathbf{Q} := [\mathbf{q}_1, \dots, \mathbf{q}_N]^\top, \mathbf{K} := [\mathbf{k}_1, \dots, \mathbf{k}_N]^\top$ , and  $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_N]^\top$ , where the vectors  $\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i$  for  $i = 1, \dots, N$  are the query, key, and value vectors, respectively.

**Step 2.** The output sequence  $\mathbf{H} := [\mathbf{h}_1, \dots, \mathbf{h}_N]^\top$  is then computed as follows

$$\mathbf{H} = \text{softmax}\left(\mathbf{Q}\mathbf{K}^\top / \sqrt{D}\right) \mathbf{V} := \mathbf{A}\mathbf{V}, \quad (1)$$

where the softmax function is applied to each row of the matrix  $\mathbf{Q}\mathbf{K}^\top / \sqrt{D}$ . The matrix  $\mathbf{A} := \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right) \in \mathbb{R}^{N \times N}$  and its component  $a_{ij}$  for  $i, j = 1, \dots, N$  are called the attention matrix and attention scores, respectively. For each query vector  $\mathbf{q}_i$  for  $i = 1, \dots, N$ , an equivalent form of Eqn. (1) to compute the output vector  $\mathbf{h}_i$  is given by

$$\mathbf{h}_i = \sum_{j=1}^N \text{softmax}\left(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{D}\right) \mathbf{v}_j. \quad (2)$$

The self-attention computed by Eqn. (1) and (2) is called the scaled dot-product or softmax attention. In our paper, we call a transformer that uses this attention the softmax transformer. The structure that the attention matrix  $\mathbf{A}$  learns from training determines the ability of the self-attention to capture contextual representation for each token. Additionally, a residual connection can be added to the output of the self-attention layer,  $\mathbf{h}_i = \mathbf{x}_i + \sum_{j=1}^N \text{softmax}\left(\mathbf{q}_i^\top \mathbf{k}_j / \sqrt{D}\right) \mathbf{v}_j$ .

**Multi-head Attention (MHA).** In MHA, multiple heads are concatenated to compute the final output. This MHA mechanism allows transformers to capture more diverse attention patterns and increase the capacity of the model. Let  $H$  be the number of heads and  $\mathbf{W}_O^{\text{multi}} = [\mathbf{W}_O^1, \dots, \mathbf{W}_O^H] \in \mathbb{R}^{D_v \times H D_v}$  be the projection matrix for the output where  $\mathbf{W}_O^1, \dots, \mathbf{W}_O^H \in \mathbb{R}^{D_v \times D_v}$ . The MHA is defined as

$$\text{MultiHead}(\{\mathbf{H}\}_{s=1}^H) = \text{Concat}(\mathbf{H}^1, \dots, \mathbf{H}^H) \mathbf{W}_O^{\text{multi}\top} = \sum_{s=1}^H \mathbf{H}^s \mathbf{W}_O^{s\top} = \sum_{s=1}^H \mathbf{A}^s \mathbf{V}^s \mathbf{W}_O^{s\top}. \quad (3)$$

Despite their remarkable success, most attention layers are developed based on heuristic approaches, and a coherent principled framework for synthesizing attention layers has remained elusive.

### 1.2 CONTRIBUTION

We derive the self-attention as the support vector expansion of a given support vector regression (SVR) problem. The primal representation of the regression function has the form of a neural network layer. Thus, we establish a primal-dual connection between an attention layer in transformers and a neural network layer in deep neural networks. Our framework suggests a principled approach to developing an attention mechanism: *Starting from a neural network layer and a support vector regression problem, we derive the dual as a support vector expansion to attain the corresponding attention layer.* We then employ this principled approach to invent two novel classes of attentions: the Batch Normalized Attention (Attention-BN) derived from the batch normalization layer in deep neural networks and the Attention with Scaled Heads (Attention-SH) resulting from solving the support vector regression model with less amount of training data. Our contribution is three-fold.

1. We derive self-attention as a support vector expansion that solves a SVR problem, thus providing a principled primal-dual framework to study and develop self-attentions.
2. We re-derive popular attentions, such as the linear attention (Katharopoulos et al., 2020), the sparse attention (Child et al., 2019), and the multi-head attention (Vaswani et al., 2017), from our proposed framework.

3. We develop two new attention mechanism: the Batch Normalized Attention (Attention-BN) and the Attention with Scaled Heads (Attention-SH) using our proposed framework.

We empirically demonstrate that 1) the Attention-BN significantly outperforms the baseline softmax and linear attention and 2) the Attention-SH performs better while being more efficient than the same baselines on a variety of practical tasks including image and time-series classification.

## 2 PRIMAL-DUAL INTERPRETATION OF SELF-ATTENTION

We first provide a primal-dual interpretation of self-attention as a support vector regression problem in Section 2.1. Based on that primal-dual framework, we derive popular attention mechanisms as the support vector expansion in Section 2.2. Finally, we introduce two new attention mechanisms in Section 2.3, the Attention-BN and Attention-SH.

### 2.1 ATTENTION AS A SUPPORT VECTOR REGRESSION MODEL

In this section, we derive self-attention from a support vector regression problem. Suppose we are given a training data  $\{(\mathbf{k}_1, \mathbf{y}_1), \dots, (\mathbf{k}_N, \mathbf{y}_N)\} \subset \mathcal{K} \times \mathcal{Y}$ , where  $\mathcal{K} = \mathbb{R}^D$  and  $\mathcal{Y} = \mathbb{R}^{D_v}$ . Here,  $\mathbf{k}_1, \dots, \mathbf{k}_N$  are attention keys in self-attention, and  $\mathbf{y}_1, \dots, \mathbf{y}_N$  are the training targets. We consider the function  $f$ , taking the form

$$\mathbf{y} = f(\mathbf{x}) := \mathbf{W} \frac{\Phi(\mathbf{x})}{h(\mathbf{x})} + \mathbf{b}, \quad (4)$$

where  $\mathbf{x} \in \mathcal{K} = \mathbb{R}^D$ ,  $\Phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_{D_\phi}(\mathbf{x})] \in \mathbb{R}^{D_\phi}$ ,  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{D_v}]^\top \in \mathbb{R}^{D_v \times D_\phi}$ ,  $\mathbf{b} \in \mathbb{R}^{D_v}$ , and  $h(\mathbf{x})$  is a vector-scalar function. We fit the function  $f$  to the training data  $\{(\mathbf{k}_1, \mathbf{y}_1), \dots, (\mathbf{k}_N, \mathbf{y}_N)\}$  with an  $\mathbb{L}_2$  regularization on  $\mathbf{W}$ , i.e., a ridge regression, by solving the following convex optimization problem:

$$\begin{aligned} \underset{\mathbf{W}}{\text{minimize}} \quad & \frac{1}{2} \|\mathbf{W}\|_F^2 + C \sum_{j=1}^N \sum_{d=1}^{D_v} \left( \xi_j(d) + \tilde{\xi}_j(d) \right) = \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d\|^2 + C \sum_{j=1}^N \sum_{d=1}^{D_v} \left( \xi_j(d) + \tilde{\xi}_j(d) \right) \\ \text{subject to} \quad & \begin{cases} \mathbf{y}_j(d) - \mathbf{w}_d^\top \Phi(\mathbf{k}_j)/h(\mathbf{k}_j) - \mathbf{b}(d) \leq \epsilon + \xi_j(d) \\ \mathbf{w}_d^\top \Phi(\mathbf{k}_j)/h(\mathbf{k}_j) + \mathbf{b}(d) - \mathbf{y}_j(d) \leq \epsilon + \tilde{\xi}_j(d) \\ \xi_j(d), \tilde{\xi}_j(d) \geq 0 \end{cases}, \quad j = 1, \dots, N, \quad d = 1, \dots, D_v. \end{aligned} \quad (5)$$

The Eqn. 5 implies that there exists a function  $f$  that can approximate all pairs  $(\mathbf{k}_j, \mathbf{y}_j)$  with  $\epsilon$  precision. The additional slack variables  $\xi_j, \tilde{\xi}_j$  relax this assumption and allows some of the training set data points to have the training error greater than  $\epsilon$  just as in the soft-margin SVM (Cortes & Vapnik, 1995; Schölkopf et al., 2002).  $C > 0$  is a constant determining the trade between the complexity penalizer  $\sum_{d=1}^{D_v} \|\mathbf{w}_d\|^2$ , i.e., the flatness of  $f$ , and the amount up to which deviations larger than  $\epsilon$  are tolerated.

In order to derive the self-attention from the support vector regression defined by the optimization problem 5, the key idea to construct the Lagrangian from Eqn. 5 and find the representation of the  $\mathbf{w}_d, d = 1, \dots, D_v$ , in terms of the dual variables. We define the Lagrangian function as follows:

$$\begin{aligned} \mathcal{L} := & \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d\|^2 + C \sum_{j=1}^N \sum_{d=1}^{D_v} \left( \xi_j(d) + \tilde{\xi}_j(d) \right) - \sum_{j=1}^N \sum_{d=1}^{D_v} \left( \eta_j(d) \xi_j(d) + \tilde{\eta}_j(d) \tilde{\xi}_j(d) \right) \\ & - \sum_{j=1}^N \sum_{d=1}^{D_v} \alpha_j(d) \left( \epsilon + \xi_j(d) - \mathbf{y}_j(d) + \mathbf{w}_d^\top \frac{\Phi(\mathbf{k}_j)}{h(\mathbf{k}_j)} + \mathbf{b}(d) \right) \\ & - \sum_{j=1}^N \sum_{d=1}^{D_v} \tilde{\alpha}_j(d) \left( \epsilon + \tilde{\xi}_j(d) + \mathbf{y}_j(d) - \mathbf{w}_d^\top \frac{\Phi(\mathbf{k}_j)}{h(\mathbf{k}_j)} - \mathbf{b}(d) \right), \end{aligned} \quad (6)$$

where  $\eta_j, \tilde{\eta}_j, \alpha_j$  and  $\tilde{\alpha}_j$  are Lagrange multipliers. These dual variables have to satisfy positivity constraints, i.e.,  $\eta_j(d), \tilde{\eta}_j(d), \alpha_j(d), \tilde{\alpha}_j(d) \geq 0, \forall j = 1, \dots, N, \forall d = 1, \dots, D_v$ . It follows from the saddle point condition that the partial derivatives of the Lagrangian function  $\mathcal{L}$  with respect to the primal variables  $(\mathbf{w}_d, \mathbf{b}(d), \{\xi_j(d), \tilde{\xi}_j(d)\}_{j=1}^N), d = 1, \dots, D_v$ , have to vanish for optimality,

namely, we have:

$$\partial_{\mathbf{b}(d)} \mathcal{L} = \sum_{j=1}^N (\tilde{\alpha}_j(d) - \alpha_j(d)) = 0 \Rightarrow \sum_{j=1}^N (\alpha_j(d) - \tilde{\alpha}_j(d)) = 0, \quad (7)$$

$$\partial_{\mathbf{w}_d} \mathcal{L} = \mathbf{w}_d - \sum_{j=1}^N (\alpha_j(d) - \tilde{\alpha}_j(d)) \frac{\Phi(\mathbf{k}_j)}{h(\mathbf{k}_j)} = 0 \Rightarrow \mathbf{w}_d = \sum_{j=1}^N (\alpha_j(d) - \tilde{\alpha}_j(d)) \frac{\Phi(\mathbf{k}_j)}{h(\mathbf{k}_j)}, \quad (8)$$

$$\partial_{\xi_j(d)} \mathcal{L} = C - \alpha_j(d) - \eta_j(d) = 0, \quad \partial_{\tilde{\xi}_j(d)} \mathcal{L} = C - \tilde{\alpha}_j(d) - \tilde{\eta}_j(d) = 0. \quad (9)$$

Let  $\mathbf{v}_j = [\frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(\mathbf{k}_j)}, \dots, \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(\mathbf{k}_j)}]^\top$ ,  $j = 1, \dots, N$ , and substitute Eqn. 8 into Eqn. 4, we obtain the following support vector expansion of the linear basis function  $f$ :

$$\begin{aligned} f(\mathbf{x}) &= \left[ \sum_{j=1}^N \frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(\mathbf{k}_j)} \frac{\Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j)}{h(\mathbf{x})}, \dots, \sum_{j=1}^N \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(\mathbf{k}_j)} \frac{\Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j)}{h(\mathbf{x})} \right]^\top + \mathbf{b}, \\ &= \sum_{j=1}^N \frac{\Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j)}{h(\mathbf{x})} \mathbf{v}_j + \mathbf{b}. \end{aligned} \quad (10)$$

**Remark 1** Notice that from Eqn. 9 and the conditions  $\eta_j(d), \tilde{\eta}_j(d), \alpha_j(d), \tilde{\alpha}_j(d) \geq 0$ , we can prove that  $\alpha_j(d), \tilde{\alpha}_j(d) \in [0, C]$ . Furthermore, we can show that  $\alpha_j(d) * \tilde{\alpha}_j(d) = 0$  (Smola & Schölkopf, 2004; Schölkopf et al., 2002). As a result,  $\mathbf{v}_j(d) \in [-\frac{C}{h(\mathbf{k}_j)}, \frac{C}{h(\mathbf{k}_j)}]$ ,  $d = 1, \dots, D_v$ .

**Deriving Softmax Attention.** Choosing the appropriate  $h(\mathbf{x})$  and  $\Phi(\mathbf{x})$  allows us to derive the popular softmax attention given in Eqn. 1 and 2. In particular, if we choose  $h(\mathbf{x}) := \sum_{j=1}^N \Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j)$ , Eqn. 10 becomes

$$f(\mathbf{x}) = \sum_{j=1}^N \frac{\Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j)}{\sum_{j'=1}^N \Phi(\mathbf{x})^\top \Phi(\mathbf{k}_{j'})} \mathbf{v}_j + \mathbf{b} = \frac{\sum_{j=1}^N \Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j) \mathbf{v}_j}{\sum_{j'=1}^N \Phi(\mathbf{x})^\top \Phi(\mathbf{k}_{j'})} + \mathbf{b}. \quad (11)$$

We then select  $\Phi(\mathbf{x}) = (a_{l_0}^{(0)}, a_1^{(1)}, \dots, a_{l_1}^{(1)}, \dots, a_1^{(t)}, \dots, a_{l_t}^{(t)}, \dots)$  where  $l_t = \binom{D+t-1}{t}$  and

$$a_l^{(t)} = \frac{(x_1/\sqrt[4]{D})^{n_1} \dots (x_D/\sqrt[4]{D})^{n_D}}{\sqrt{n_1! \dots n_D!}} \mid n_1 + \dots + n_D = t, 1 \leq l \leq l_t. \quad (12)$$

Since

$$\exp(\mathbf{x}^\top \mathbf{y}) = \sum_{t=0}^{\infty} \frac{(\mathbf{x}^\top \mathbf{y})^t}{t!} = \sum_{t=0}^{\infty} \sum_{n_1 + \dots + n_D = t} \left( \frac{x_1^{n_1} \dots x_D^{n_D}}{\sqrt{n_1! \dots n_D!}} \right) \left( \frac{y_1^{n_1} \dots y_D^{n_D}}{\sqrt{n_1! \dots n_D!}} \right), \quad (13)$$

then Eqn. 27 becomes

$$f(\mathbf{x}) = \sum_{j=1}^N \frac{\exp(\mathbf{x}^\top \mathbf{k}_j / \sqrt{D})}{\sum_{j'=1}^N \exp(\mathbf{x}^\top \mathbf{k}_{j'} / \sqrt{D})} \mathbf{v}_j + \mathbf{b} = \sum_{j=1}^N \text{softmax}(\mathbf{x}^\top \mathbf{k}_j / \sqrt{D}) \mathbf{v}_j + \mathbf{b}. \quad (14)$$

Let  $\mathbf{x} = \mathbf{q}_i$ ,  $\mathbf{b} = 0$  and relax the boundness constraint of  $\mathbf{v}_j$  in Remark 1. Eqn. 30 becomes Eqn. 2 of the softmax attention (Vaswani et al., 2017). We summarize our results in the following theorem.

**Theorem 1 (Softmax Attention as a Support Vector Expansion)** Given the function  $f$  defined in Eqn. 4 with  $h(\mathbf{x}) := \sum_{j=1}^N \Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j)$  and the support vector regression problem defined in Eqn. 5, we set  $\mathbf{b} = 0$ , choose  $\Phi(\mathbf{x})$  as in Eqn. 28, and relax the boundness constraint of the variables  $\mathbf{v}_j = [\frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(\mathbf{k}_j)}, \dots, \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(\mathbf{k}_j)}]^\top$ , where  $\alpha_j$  and  $\tilde{\alpha}_j$  are dual variables of Eqn. 5,  $j = 1, \dots, N$ . Then, the support vector expansion of  $f$  derived from Eqn. 5 has the form of a softmax attention

$$f(\mathbf{x}) = \sum_{j=1}^N \text{softmax}(\mathbf{x}^\top \mathbf{k}_j / \sqrt{D}) \mathbf{v}_j. \quad (15)$$

**Remark 2** Since  $\mathbf{b}$  is set to 0, the centering constraint of  $\alpha_j$  and  $\tilde{\alpha}_j$  in Eqn. 7 can be ignored.

**Remark 3** Theorem 1 and its derivation can be easily extended to capture the full form of the softmax attention with the residual connection, the query matrix projection  $\mathbf{W}_Q$ , the key matrix projection  $\mathbf{W}_K$ , and the value matrix projection  $\mathbf{W}_V$ . We include this result in Appendix F.

**Remark 4** The primal representation of the function  $f$  as in Eqn. 4 has the form of a neural network layer where  $\mathbf{W}$  is the weight,  $\mathbf{b}$  is the bias term,  $\Phi(\mathbf{x})$  is the input, and  $h(\mathbf{x})$  is the normalization term. Thus, an attention layer and a neural network layer are primal-dual of each other.

**A principled approach to developing an attention mechanism.** The observation in Remark 4 suggests a principled way to construct an attention layer: Starting from a neural network layer and a support vector regression problem, we derive the dual as a support vector expansion to attain the corresponding attention layer. Using this approach, we derive popular attention mechanisms in Section 2.2 and propose our new attention mechanisms in Section 2.3.

## 2.2 DERIVING POPULAR ATTENTION MECHANISMS AS THE SUPPORT VECTOR EXPANSION

In this section, we derive popular attentions such as the linear attention (Katharopoulos et al., 2020), the sparse attention (Child et al., 2019), and the multi-head attention (Vaswani et al., 2017).

### 2.2.1 LINEAR ATTENTION

The Eqn. 27, which is obtained when choosing  $h(\mathbf{x}) := \sum_j^N \Phi(\mathbf{x})^T \Phi(\mathbf{k}_j)$ , already matches the formula of the linear attention. Here, we can let  $\mathbf{b} = 0$  as above and select the function  $\Phi$  that results in a positive similarity function, e.g.  $\Phi(\mathbf{x}) = \text{elu}(\mathbf{x}) + 1$ , as in (Katharopoulos et al., 2020).

### 2.2.2 SPARSE ATTENTION

The sparse attention (Child et al., 2019) can be derived by fitting the function  $f$  in Eqn. 4 using a different subset  $\{(\mathbf{k}_{m_x(1)}, \mathbf{y}_{m_x(1)}), \dots, (\mathbf{k}_{m_x(M)}, \mathbf{y}_{m_x(M)})\}$  of training data  $\{(\mathbf{k}_1, \mathbf{y}_1), \dots, (\mathbf{k}_N, \mathbf{y}_N)\}$  for each input data  $\mathbf{x}$ , where  $\mathcal{M}_x = \{m_x(1), \dots, m_x(M)\} \subset \{1, \dots, N\}$ . The support vector expansion of  $f$  is then given by

$$f(\mathbf{x}) = \sum_{j=1}^N \mathbf{1}_{\mathcal{M}_x}(j) \frac{\Phi(\mathbf{x})^T \Phi(\mathbf{k}_j)}{h(\mathbf{x})} \mathbf{v}_j + \mathbf{b} \quad (16)$$

where  $\mathbf{1}_{\mathcal{M}_x}(j) = [j \in \mathcal{M}_x] := \begin{cases} 1 & \text{if } j \in \mathcal{M}_x \\ 0 & \text{otherwise} \end{cases}$ . Note that the subsets  $\mathcal{M}_x$  are different for different  $\mathbf{x}$ . When letting  $\mathbf{x} = \mathbf{q}_i$  where  $\mathbf{q}_i, i = 1, \dots, N$ , are the query vectors and choosing  $\Phi, h, \mathbf{b}$  as in Section 2.1, we can obtain the sparse attention in (Child et al., 2019) where the binary matrix  $\mathbf{M} = (\mathbf{1}_{\mathcal{M}_{q_i}}(j))_{i,j=1}^N$  becomes the sparse masking matrix.

### 2.2.3 MULTI-HEAD ATTENTION (MHA)

The MHA can be derived by solving multiple support vector regression problems and then linearly combining their outputs. In particular, given  $H$  training datasets  $\{(\mathbf{k}_1^1, \mathbf{y}_1^1), \dots, (\mathbf{k}_N^1, \mathbf{y}_N^1)\}, \dots, \{(\mathbf{k}_1^H, \mathbf{y}_1^H), \dots, (\mathbf{k}_N^H, \mathbf{y}_N^H)\} \subset \mathcal{K} \times \mathcal{Y}$ , where  $\mathcal{K} = \mathbb{R}^D$  and  $\mathcal{Y} = \mathbb{R}^{D_v}$ . We define the function  $f$  applied on the input vector  $\mathbf{x} = [\mathbf{x}^1, \dots, \mathbf{x}^H]$  as follows

$$\mathbf{y} = f(\mathbf{x}) := \sum_{s=1}^H \mathbf{W}_O^s \mathbf{y}^s = \sum_{s=1}^H \mathbf{W}_O^s f^s(\mathbf{x}^s) = \sum_{s=1}^H \mathbf{W}_O^s \left( \mathbf{W}^s \frac{\Phi^s(\mathbf{x}^s)}{h^s(\mathbf{x}^s)} + \mathbf{b}^s \right), \quad (17)$$

where each function  $f^s(\mathbf{x}^s) = \mathbf{W}^s \frac{\Phi^s(\mathbf{x}^s)}{h^s(\mathbf{x}^s)} + \mathbf{b}^s$  is fitted to the training dataset  $\{(\mathbf{k}_1^s, \mathbf{y}_1^s), \dots, (\mathbf{k}_N^s, \mathbf{y}_N^s)\}$ . Following the same derivation and choosing  $\{\Phi^s, h^s, \mathbf{b}^s\}_{s=1}^H$  as in Section 2.1, we can rewrite  $f(\mathbf{x})$  in terms of the support vector expansions of the individual functions  $f^s(\mathbf{x}^s)$ , which are the individual softmax attentions

$$f(\mathbf{x}) = \sum_{s=1}^H \mathbf{W}_O^s \left( \sum_{j=1}^N \frac{\Phi^s(\mathbf{x}^s)^T \Phi^s(\mathbf{k}_j^s)}{h^s(\mathbf{x}^s)} \mathbf{v}_j^s + \mathbf{b}^s \right) = \sum_{s=1}^H \mathbf{W}_O^s \left( \sum_{j=1}^N \text{softmax}(\mathbf{x}^{s\top} \mathbf{k}_j^s / \sqrt{D}) \mathbf{v}_j^s \right). \quad (18)$$

Comparing Eqn. 18 and Eqn. 3, we see that Eqn. 18 computes the MHA when choosing  $\mathbf{x}^s = \mathbf{q}_i^s$  where  $\mathbf{q}_i^s, i = 1, \dots, N$ , are the query vectors at the  $s^{\text{th}}$  head.

## 2.3 DERIVING NEW ATTENTION MECHANISMS: BATCH NORMALIZED ATTENTION AND MULTIREOLUTION HEAD ATTENTION

In this section, we employ our primal-dual framework to develop new attention mechanisms. In particular, we derive: 1) the Batch Normalized Attention from employing the batch normalization (Ioffe & Szegedy, 2015); and 2) the Attention with Scaled Heads from using different amounts of training data. By 1) and 2), we demonstrate that *new attentions can be invented by modifying the primal neural network layer and the support vector regression problem in our framework, respectively.*

### 2.3.1 BATCH NORMALIZED ATTENTION

We incorporate the batch normalization into the primal form of the function  $f$  in Eqn. 4. Given a training data  $\{(\mathbf{k}_1, \mathbf{y}_1), \dots, (\mathbf{k}_N, \mathbf{y}_N)\} \subset \mathcal{K} \times \mathcal{Y}$ , where  $\mathcal{K} = \mathbb{R}^D$  and  $\mathcal{Y} = \mathbb{R}^{D_v}$  as in Section 2.1, the resultant  $f$  is defined as follows

$$f(\mathbf{x}) := \mathbf{W} \frac{\Phi((\mathbf{x} - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})}{h((\mathbf{x} - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})} + \mathbf{b}, \quad (19)$$

where

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{j=1}^N \mathbf{k}_j, \quad \mathbf{s}^{-1} = \left[ \frac{1}{\sqrt{\sigma_1^2 + \epsilon}}, \dots, \frac{1}{\sqrt{\sigma_D^2 + \epsilon}} \right]^\top, \quad \sigma_d^2 = \frac{1}{N} \sum_{j=1}^N (\mathbf{k}_j(d) - \boldsymbol{\mu}(d))^2. \quad (20)$$

Here,  $d = 1, \dots, D$ , and the mean subtraction and division by the standard deviation is performed element-wise along the feature dimension of  $\mathbf{x}$ . Following the same derivation as in Section 2.1, we derive the following support vector expansion of  $f$

$$f(\mathbf{x}) = \sum_{j=1}^N \frac{\Phi((\mathbf{x} - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})^\top \Phi((\mathbf{k}_j - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})}{h((\mathbf{x} - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})} \mathbf{v}_j + \mathbf{b}. \quad (21)$$

Here,  $\mathbf{v}_j = \left[ \frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h((\mathbf{k}_j - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})}, \dots, \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h((\mathbf{k}_j - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})} \right]^\top$ , where  $\alpha_j$  and  $\tilde{\alpha}_j$  are the dual variables,  $j = 1, \dots, N$ . Same as in Section 2.1, in Eqn. 21, we choose  $\Phi$  as in Eqn. 28,  $h(\mathbf{x}) := \sum_j \Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j)$ , and  $\mathbf{b} = 0$  to obtain the Batch Normalized Attention, which is defined as follows.

**Definition 1 (Batch Normalized Attention)** Given a set of the key and value vectors  $\{\mathbf{k}_j, \mathbf{v}_j\}_{j=1}^N$ , for each query vector  $\mathbf{q}_i$ ,  $i = 1, \dots, N$ , the Batch Normalized Attention (Attention-BN) computes the corresponding output vector  $\mathbf{h}_i$  of the query  $\mathbf{q}_i$  by the following attention formula:

$$\mathbf{h}_i = \sum_{j=1}^N \text{softmax} \left( ((\mathbf{q}_i - \boldsymbol{\mu}) \odot \mathbf{s}^{-1})^\top ((\mathbf{k}_j - \boldsymbol{\mu}) \odot \mathbf{s}^{-1}) / \sqrt{D} \right) \mathbf{v}_j, \quad (22)$$

where

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{j=1}^N \mathbf{k}_j, \quad \mathbf{s}^{-1} = \left[ \frac{1}{\sqrt{\sigma_1^2 + \epsilon}}, \dots, \frac{1}{\sqrt{\sigma_D^2 + \epsilon}} \right]^\top, \quad \sigma_d^2 = \frac{1}{N} \sum_{j=1}^N (\mathbf{k}_j(d) - \boldsymbol{\mu}(d))^2. \quad (23)$$

**The Effect of Normalization.** Expanding the dot product in the Attention-BN (see Appendix E), Eqn. 22 becomes

$$\mathbf{h}_i = \sum_{j=1}^N \text{softmax} \left( \frac{\sum_{d=1}^D \mathbf{q}_i(d) \mathbf{k}_j(d) - \frac{1}{N} \sum_{j'=1}^N \mathbf{k}_{j'}(d) \mathbf{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \mathbf{v}_j. \quad (24)$$

Eqn. 24 implies that in the Attention-BN, the similarity between the query  $\mathbf{q}_i$  and the key  $\mathbf{k}_j$  is adjusted by the similarity between the key  $\mathbf{k}_j$  and all the keys  $\mathbf{k}_{j'}$ ,  $j' = 1, \dots, N$ . In particular, if the key  $\mathbf{k}_j$  is too similar to other keys, the query  $\mathbf{q}_i$  will attend to it less and vice versa.

### 2.3.2 ATTENTION WITH SCALED HEADS

The *Attention with Scaled Heads*, named Attention-SH, is derived based on the derivation of the MHA in Section 2.2.3. The key idea underlying the Attention-SH is to train multiple support vector regression problems using different amounts of training data. In particular, the Attention-SH follows Eqn. 17 in Section 2.2.3 and defines the same regression function  $f$  as the MHA. However, the Attention-SH fits the function  $f^s$ ,  $s = 1, \dots, H$ , in Eqn. 17 with training sets  $\{(\mathbf{k}_1^1, \mathbf{y}_1^1), \dots, (\mathbf{k}_{N_1}^1, \mathbf{y}_{N_1}^1)\}, \dots, \{(\mathbf{k}_1^H, \mathbf{y}_1^H), \dots, (\mathbf{k}_{N_H}^H, \mathbf{y}_{N_H}^H)\} \subset \mathcal{K} \times \mathcal{Y}$  of different sizes  $N_1, \dots, N_H$ , where  $\mathcal{K} = \mathbb{R}^{\bar{D}}$  and  $\mathcal{Y} = \mathbb{R}^{D_v}$ . The resultant support vector expansion yields the formula of the Attention-SH as in the following definition.

**Definition 2 (Attention with Scaled Heads)** Given  $H$  sets of the key and value vectors  $\{\mathbf{k}_j^1, \mathbf{v}_j^1\}_{j=1}^{N_1}, \dots, \{\mathbf{k}_j^H, \mathbf{v}_j^H\}_{j=1}^{N_H}$ , for each set of  $H$  query vectors  $\mathbf{q}_i^1, \dots, \mathbf{q}_i^H$ ,  $i = 1, \dots, N$ , the Attention with Scaled Heads (Attention-SH) computes the corresponding output vector  $\mathbf{h}_i$  of the queries  $\mathbf{q}_i^1, \dots, \mathbf{q}_i^H$  by the following attention formula:

$$\mathbf{h}_i = \sum_{s=1}^H \mathbf{W}_O^s \left( \sum_{j=1}^{N_s} \text{softmax} \left( \mathbf{q}_i^s{}^\top \mathbf{k}_j^s / \sqrt{D} \right) \mathbf{v}_j^s \right). \quad (25)$$



Table 1: Test Accuracy (%) of the Attention-BN/SH/BN+SH vs. the baseline softmax attention on a subset of the UEA Time Series Classification Archive benchmark (Bagnall et al., 2018). Our proposed attentions significantly outperform the baseline. We also include the reported results from (Zerveas et al., 2021) and (Wu et al., 2022) (in parentheses) in addition to our reproduced results.

Dataset/Model	Baseline Softmax	Attention-BN	Attention-SH	Attention-BN+SH
ETHANOLCONCENTRATION	32.08 $\pm$ 1.24 (33.70)	33.33 $\pm$ 0.44	33.59 $\pm$ 0.58	<b>34.35 <math>\pm</math> 0.53</b>
FACEDETECTION	68.70 $\pm$ 0.61 (68.10)	68.62 $\pm$ 0.26	<b>68.83 <math>\pm</math> 0.16</b>	68.67 $\pm$ 0.13
HANDWRITING	32.08 $\pm$ 0.88 (30.50)	33.17 $\pm$ 0.20	33.29 $\pm$ 0.42	<b>33.45 <math>\pm</math> 0.61</b>
HEARTBEAT	75.77 $\pm$ 1.01 (77.60)	76.10 $\pm$ 0.98	76.25 $\pm$ 1.02	<b>76.26 <math>\pm</math> 1.07</b>
JAPANESEVOWELS	99.46 $\pm$ 0.27 (99.40)	<b>99.55 <math>\pm</math> 0.31</b>	99.46 $\pm$ 0.27	<b>99.55 <math>\pm</math> 0.31</b>
PEMS-SF	82.66 $\pm$ 0.51 (82.10)	<b>84.77 <math>\pm</math> 0.33</b>	83.04 $\pm$ 0.86	83.81 $\pm$ 0.62
SELFREGULATIONSCP1	91.46 $\pm$ 0.35 (92.50)	91.58 $\pm$ 0.39	91.70 $\pm$ 0.39	<b>92.04 <math>\pm</math> 0.36</b>
SELFREGULATIONSCP2	54.72 $\pm$ 0.74 (53.90)	56.11 $\pm$ 0.71	55.93 $\pm$ 0.85	<b>57.04 <math>\pm</math> 0.82</b>
SPOKENARABICDIGITS	99.33 $\pm$ 0.02 (99.30)	99.23 $\pm$ 0.09	99.34 $\pm$ 0.11	<b>99.42 <math>\pm</math> 0.28</b>
UWAVEGESTURELIBRARY	84.45 $\pm$ 0.72 (85.60)	86.46 $\pm$ 0.81	86.77 $\pm$ 0.78	<b>87.60 <math>\pm</math> 0.67</b>
AVERAGE ACCURACY	72.07 $\pm$ 0.47(72.27)	72.89 $\pm$ 0.09	72.82 $\pm$ 0.12	<b>73.22 <math>\pm</math> 0.33</b>

**Remark 5** For a given input sequence  $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D_x}$  of  $N$  feature vectors in self-attention, in order to generate the sets of  $\{\mathbf{k}_j^s, \mathbf{v}_j^s\}_{j=1}^{N_s}$  at the scale  $s^{th}$ , we can downsample the input  $\mathbf{X}$  before projecting into the key matrix  $\mathbf{K}$  and the value matrix  $\mathbf{V}$ . There are multiples approaches to downsampling  $\mathbf{X}$ , such as using the average-pooling, max-pooling, 1-D convolution, or  $K$ -means clustering. In this paper, we employ the average-pooling to downsample  $\mathbf{X}$ .

**Linear Attention with Batch Normalization and Scaled Heads.** The Attention-BN/SH can be extended to use with the linear attention. In particular, in the Linear Attention-BN/SH, we replace the softmax kernel in Eqn. 22 and Eqn. 25 by the linear kernel, respectively.

### 3 EXPERIMENTAL RESULTS

In this section, we empirically demonstrate the advantages of our Attention-BN, Attention-SH, and their combination (Attention-BN+SH) over the baseline softmax attention on the UEA time-series classification benchmark (Bagnall et al., 2018), the Long Range Arena benchmark (Tay et al., 2021), and the image classification task on the Imagenet dataset (Deng et al., 2009; Russakovsky et al., 2015). We aim to show that: (i) Attention-BN significantly outperforms the softmax baseline across tasks; (ii) Attention-SH achieves better or comparable accuracy while saving computation and memory compared to the baseline; (iii) Attention-BN+SH, which combines both Attention-BN and Attention-SH, results in the best model performance in term of accuracy and efficiency; (iv) all our proposed models help reduce the redundancy in multi-head attention and benefit learning of the long-term dependency in long input sequences; (v) Attention-BN and Attention-SH can be applied on other attention mechanisms beyond the softmax attention. When combined with the linear attention (Katharopoulos et al., 2020), the resultant Linear Attention-BN and Linear Attention-SH yield similar advantages mentioned in (i), (ii), (iii) and (iv) over the baseline linear attention.

In our experiments, we compare the proposed models with the baseline softmax and linear attentions of the same configuration. For the Attention-BN and Attention-BN+SH, we observe that recentering queries and keys alone is sufficient for improving the model performance. In addition, weighting  $\mu$  with a constant  $\beta$ , as in Eqn. 26 in the Appendix, enables the Attention-BN/BN+SH to adjust the effect of normalization to the attention score and help increase the accuracy. Our results are averaged over 5 runs. Details on datasets, models, and training are provided in Appendix A.

**UEA Time Series Classification.** Table 1 compares the accuracy of the Attention-BN and Attention-SH with the baseline softmax attention on 10 tasks in the UEA Time Series Classification benchmark (Bagnall et al., 2018). Both Attention-BN and Attention-SH significantly outperform the softmax baseline on most tasks and on average among all tasks. When combining two models, the resulting Attention-BN+SH yields the best accuracy with more than 1% overall improvement over the softmax baseline. Notably, the Attention-SH and Attention-BN+SH are much more efficient than the baseline since they need much fewer keys and values in computing the attention output. The efficiency advantage of the Attention-SH/BN+SH is discussed and analyzed in detail in Section 4.

**Long Range Arena (LRA) benchmark.** In this experiment, we verify the advantage of our methods over the softmax baseline on tasks that involve very long sequences (e.g., the sequence length can be up to 4K) in the LRA benchmark (Tay et al., 2021). Those tasks require the model to capture long-range dependency in the input sequence. The summarized results in Table 2 indicate significant improvements of Attention-BN/SH/BN+SH over the baseline softmax attention. Same as in the UEA

Table 2: Test Accuracy (%) of the Attention-BN/SH/BN+SH vs. the baseline softmax attention on the LRA benchmark (Tay et al., 2021). Our models significantly outperform the softmax baseline.

Dataset/Model	Baseline Softmax	Attention-BN	Attention-SH	Attention-BN+SH
LISTOPS	36.76 $\pm$ 0.42	37.32 $\pm$ 0.06	37.08 $\pm$ 0.48	<b>37.33 <math>\pm</math> 0.53</b>
TEXT	64.90 $\pm$ 0.07	65.07 $\pm$ 0.08	<b>65.19 <math>\pm</math> 0.19</b>	65.03 $\pm$ 0.35
RETRIEVAL	79.68 $\pm$ 0.52	81.05 $\pm$ 0.08	80.74 $\pm$ 0.32	<b>81.20 <math>\pm</math> 0.11</b>
IMAGE	39.23 $\pm$ 1.35	<b>39.75 <math>\pm</math> 1.21</b>	38.87 $\pm$ 1.24	39.18 $\pm$ 1.27
PATHFINDER	72.72 $\pm$ 0.75	73.24 $\pm$ 0.67	<b>74.00 <math>\pm</math> 0.29</b>	73.98 $\pm$ 0.55
AVERAGE ACCURACY	58.66 $\pm$ 0.26	59.28 $\pm$ 0.25	59.18 $\pm$ 0.22	<b>59.35 <math>\pm</math> 0.29</b>

Table 3: Top-1 and top-5 accuracy (%) of the Attention-BN/SH/SH+BN Deit vs. the baseline softmax attention Deit on the ImageNet benchmark. The Attention-BN Deit outperforms the baseline in terms of accuracy. The Attention-SH/BN+SH Deit achieve comparable accuracy with the baseline while being more efficient.

Metric/Model	Baseline Softmax Deit	Attention-BN Deit	Attention-SH Deit	Attention-BN+SH Deit
Top-1 Acc (%)	72.23 $\pm$ 0.23	<b>72.79 <math>\pm</math> 0.21</b>	72.08 $\pm$ 0.23	72.25 $\pm$ 0.22
Top-5 Acc (%)	91.13 $\pm$ 0.15	<b>91.43 <math>\pm</math> 0.12</b>	91.05 $\pm$ 0.14	91.14 $\pm$ 0.12

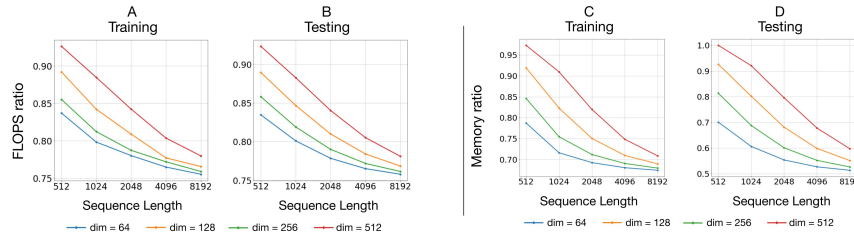


Figure 1: (Left) FLOPS ratios and (Right) memory usage ratios between the Attention-BN+SH and the softmax baseline trained on retrieval task for different model dimensions and sequence lengths. The reduction in computation and memory when using our models improves with sequence length. When scaling up the model, our methods remain significantly more beneficial than the baseline.

Time Series experiment, on this LRA benchmark, Attention-BN and Attention-SH both outperform the softmax attention on most five tasks. Moreover, Attention-BN+SH, which combines these two attention mechanisms, results in the most accurate models on average across tasks. Specifically, for the retrieval task, the most challenging task with the largest sequence length in the LRA benchmark, Attention-BN+SH achieve a remarkable improvement of more than 1.5% over the baseline.

**Image Classification on Imagenet.** We corroborate the advantage of our proposed attention over the baseline softmax attention when scaled up for the large-scale ImageNet image classification task. We summarize the results in Table 3. The Deit model (Touvron et al., 2021) equipped with the Attention-BN yields better performance than the softmax baseline. Meanwhile, Attention-SH/BN+SH Deit perform on par with the baseline while being more efficient. These results, together with other results above justify the benefits of our proposed methods across various tasks and data modalities, proving the effectiveness of our primal-dual approach to develop new attentions.

## 4 EMPIRICAL ANALYSIS

**Efficiency Analysis.** The Attention-BN+SH not only improves the accuracy of the model remarkably but also help reduce the computational and memory cost significantly. Fig. 1 presents the efficiency benefits of our Attention-BN+SH trained on the retrieval task when the model dimension  $D$  and sequence lengths  $N$  grow. The efficiency advantage of our model increase as  $N$  increase. In addition, the scaled-up models (with large  $D$ ) remains significantly more efficient than the baseline. When the model dimension is 64 and sequence length is 4096, which is the standard configuration of the task, the model’s FLOPS, in both training and inference, reduce almost 25%, whereas the reductions for memory usage in training and testing are 31.9% and 47.3%, respectively. Notably, this efficient model also outperforms the baseline with more than 1.5% improvement in accuracy. These results prove the benefits of applying the Attention-BN+SH for long-sequence tasks and large-scale models.

**New Attentions Helps Reduce Head Redundancy.** We compute the average  $\mathcal{L}_2$  distances between heads to analyze the attention diversity. Given our trained models for the retrieval task, the layer-average mean and standard deviation of distances between heads are reported in Table 4. All our introduced attentions attain greater  $\mathcal{L}_2$  distances compared to the baseline, reducing the risk of learning redundant heads. In particular, Attention-SH has the highest head difference, indicating the model’s attention patterns are most spread out between heads.



Table 4: Layer-average mean and standard deviation of  $\mathcal{L}_2$  distances between heads of Attention-BN/SH/BN+SH versus dot-product attention transformer trained on the retrieval task. Our attentions attain greater  $\mathcal{L}_2$  distances between heads than the baseline, suggesting that they capture more diverse attention patterns.

MetricModel	Baseline Softmax	Attention-BN	Attention-SH	Attention-BN+SH
Mean	2.01	2.45	<b>3.81</b>	3.19
Std	0.39	0.66	0.75	<b>1.01</b>

**Combining Attention-BN and Attention-SH with Other Attentions.** Our methods can be extended to combine with other attention mechanisms. We study the Linear Attention-BN/SH/BN+SH, that combine the Attention-BN/SH/BN+SH with the linear attention (Katharopoulos et al., 2020) as explained at the end of Section 2.3. We summarize our results in Table 5 in Appendix B.1.

## 5 RELATED WORK

**Interpretation of Attention Mechanism.** Recent works have focused on understanding the attention mechanism in transformers from different perspectives. (Tsai et al., 2019) considers attention as a weighted moving average over the inputs via a smoothing kernel. (Nguyen et al., 2022) draws a connection between self-attention and nonparametric kernel regression. With this understanding, the work explores better regression estimators, e.g. the generalized Fourier nonparametric regression estimator, to improve transformers. In addition, (Cao, 2021) then shows that the linear transformer (Katharopoulos et al., 2020) corresponds to a Petrov-Galerkin projection (Reddy, 2004) and proves that the softmax normalization in the softmax attention is sufficient but not necessary. Other works that employ ordinary/partial differential equations to provide an interpretation for attention include (Lu et al., 2019; Sander et al., 2022). From a probabilistic perspective, (Tang & Matteson, 2021; Gabbur et al., 2021; Zhang & Feng, 2021) propose Gaussian mixture model frameworks to study the self-attention in transformers. Using graph-structured learning and message passing in graphical models is another attempt at understanding the attention mechanism Wang et al. (2018); Shaw et al. (2018); Kreuzer et al. (2021). Optimization perspectives of attention mechanisms are recently explored. (Sander et al., 2022) connects transformers with an optimization process across iterations by specifically constructing the core energy function. (Sahiner et al., 2022) derive finite-dimensional convex equivalence of attentions that can be solved for global optimality. Different from these approaches, our primal-dual framework focuses on deriving attention as the dual expansion of a primal neural network layer via solving a support vector regression problem. This framework allows us to not only explain many different types of attention mechanisms but also create new ones.

**Efficient Transformers.** Recently, efficient transformers have been studied (Roy et al., 2021). Among them are sparse transformers which incorporate sparse structures into the attention matrix (Parmar et al., 2018; Liu et al., 2018; Qiu et al., 2019; Child et al., 2019; Beltagy et al., 2020). Another class of efficient transformers are models that aim to have better coverage by integrating different access patterns (Child et al., 2019; Ho et al., 2019), which can also be learned from the data (Kitaev et al., 2020; Roy et al., 2021; Tay et al., 2020). An emerging body of work is proposed to distill and prune the model, including (Sanh et al., 2019; Sun et al., 2019; Voita et al., 2019; Sajjad et al., 2020). In other works, a side memory module is utilized in order to access multiple tokens simultaneously (Lee et al., 2019; Sukhbaatar et al., 2019; Asai & Choi, 2020; Beltagy et al., 2020). Low-rank and kernelization methods have been proposed to improve the efficiency of self-attention calculation (Tsai et al., 2019; Wang et al., 2020; Katharopoulos et al., 2020; Choromanski et al., 2021; Shen et al., 2021; Peng et al., 2021). Our Attention-SH/BN+SH is orthogonal to these methods.

## 6 CONCLUDING REMARKS

In this paper, we derive self-attention as a support vector expansion that solves a support vector regression (SVR) problem and provide a principled primal-dual framework to analyze and synthesize attention mechanisms. We then use our framework to invent two new attention mechanisms, the Batch Normalized Attention (Attention-BN) and the Attention with Scaled Heads (Attention-SH), that improve the accuracy and the efficiency of the baseline softmax attention. In our work, we approximate and learn the dual variables  $\alpha_j$  and  $\tilde{\alpha}_j$  using the value vector  $v_j$ ,  $j = 1, \dots, N$  in self-attention. It is natural to include more inductive biases and structures of those dual variables from solving the dual optimization problem of SVR into the value vectors  $v_j$ . Furthermore, extending our framework to explain the attention modules that compute the attention using neural network layers or convolutional neural network layers applying on the input feature, such as the Convolutional Block Attention Module (Woo et al., 2018), is an interesting research direction. We leave these exciting research ideas as future work.

## ACKNOWLEDGEMENTS

This material is based on research sponsored by the NSF under Grant# 2030859 to the Computing Research Association for the CIFellows Project (CIF2020-UCLA-38). SJO acknowledges support from the ONR N00014-20-1-2093 and N00014-20-1-2787 and the NSF DMS 2208272 and 1952339. RGB acknowledges support from the NSF grants CCF-1911094, IIS-1838177, and IIS-1730574; ONR grants N00014-18-12571, N00014-20-1-2534, and MURI N00014-20-1-2787; AFOSR grant FA9550-22-1-0060; and a Vannevar Bush Faculty Fellowship, ONR grant N00014-18-1-2047. ALB acknowledges support from the NSF grants DMS-2152717 and DMS-1952339. NH acknowledges support from the NSF IFML 2019844 and the NSF AI Institute for Foundations of Machine Learning.

## REFERENCES

- Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. Character-level language modeling with deeper self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3159–3166, 2019.
- Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 6816–6826, 2021. doi: 10.1109/ICCV48922.2021.00676.
- Akari Asai and Eunsol Choi. Challenges in information seeking qa: Unanswerable questions and paragraph retrieval. *arXiv preprint arXiv:2010.11915*, 2020.
- Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ByxZX20qFQ>.
- Anthony Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom, Paul Southam, and Eamonn Keogh. The uea multivariate time series classification archive, 2018. *arXiv preprint arXiv:1811.00075*, 2018.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Shuhao Cao. Choose a transformer: Fourier or galerkin. *Advances in Neural Information Processing Systems*, 34, 2021.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Arvind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1179. URL <https://www.aclweb.org/anthology/D14-1179>.
- Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin Belanger, Lucy J Colwell, and Adrian Weller. Rethinking attention with performers. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=Ua6zuk0WRH>.

- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look at? an analysis of BERT’s attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 276–286, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4828. URL <https://www.aclweb.org/anthology/W19-4828>.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal transformers. *arXiv preprint arXiv:1807.03819*, 2018.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- 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.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Prasad Gabbur, Manjot Bilkhu, and Javier Movellan. Probabilistic attention for interactive segmentation. *Advances in Neural Information Processing Systems*, 34, 2021.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*, 2020.
- Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021.
- Jonathan Ho, Nal Kalchbrenner, Dirk Weissenborn, and Tim Salimans. Axial attention in multidimensional transformers. *arXiv preprint arXiv:1912.12180*, 2019.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. PMLR, 2015.
- Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. *Advances in neural information processing systems*, 34:1273–1286, 2021.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 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.
- Yoon Kim, Carl Denton, Luong Hoang, and Alexander M Rush. Structured attention networks. *arXiv preprint arXiv:1702.00887*, 2017.

- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020.
- Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, and Prudencio Tossou. Rethinking graph transformers with spectral attention. *Advances in Neural Information Processing Systems*, 34, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosioerek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International Conference on Machine Learning*, pp. 3744–3753. PMLR, 2019.
- Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. *CoRR*, abs/1703.03130, 2017. URL <http://arxiv.org/abs/1703.03130>.
- Drew Linsley, Junkyung Kim, Vijay Veerabadran, Charles Windolf, and Thomas Serre. Learning long-range spatial dependencies with horizontal gated recurrent units. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/ec8956637a99787bd197eacd77acce5e-Paper.pdf>.
- Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. *arXiv preprint arXiv:1801.10198*, 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 pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022, 2021.
- Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- Yiping Lu, Zhuohan Li, Di He, Zhiqing Sun, Bin Dong, Tao Qin, Liwei Wang, and Tie-Yan Liu. Understanding and improving transformer from a multi-particle dynamic system point of view. *arXiv preprint arXiv:1906.02762*, 2019.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/P11-1015>.
- Nikita Nangia and Samuel Bowman. ListOps: A diagnostic dataset for latent tree learning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pp. 92–99, New Orleans, Louisiana, USA, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-4013. URL <https://www.aclweb.org/anthology/N18-4013>.
- Tam Minh Nguyen, Tan Minh Nguyen, Dung DD Le, Duy Khuong Nguyen, Viet-Anh Tran, Richard Baraniuk, Nhat Ho, and Stanley Osher. Improving transformers with probabilistic attention keys. In *International Conference on Machine Learning*, pp. 16595–16621. PMLR, 2022.
- Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2249–2255, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1244. URL <https://www.aclweb.org/anthology/D16-1244>.

- Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4055–4064. PMLR, 10–15 Jul 2018. URL <http://proceedings.mlr.press/v80/parmar18a.html>.
- Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah Smith, and Lingpeng Kong. Random feature attention. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=QtTKTdVrFBB>.
- Jiezhong Qiu, Hao Ma, Omer Levy, Scott Wen-tau Yih, Sinong Wang, and Jie Tang. Blockwise self-attention for long document understanding. *arXiv preprint arXiv:1911.02972*, 2019.
- Dragomir R Radev, Pradeep Muthukrishnan, Vahed Qazvinian, and Amjad Abu-Jbara. The acl anthology network corpus. *Language Resources and Evaluation*, 47(4):919–944, 2013.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *OpenAI report*, 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pp. 8821–8831. PMLR, 2021.
- JN Reddy. *An introduction to the finite element method*, volume 1221. McGraw-Hill New York, 2004.
- Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15), 2021.
- Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9:53–68, 2021. URL <https://www.aclweb.org/anthology/2021.tacl-1.4>.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- Arda Sahiner, Tolga Ergen, Batu Ozturkler, John Pauly, Morteza Mardani, and Mert Pilanci. Unraveling attention via convex duality: Analysis and interpretations of vision transformers. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 19050–19088. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/sahiner22a.html>.
- Hassan Sajjad, Fahim Dalvi, Nadir Durrani, and Preslav Nakov. Poor man’s bert: Smaller and faster transformer models. *arXiv e-prints*, pp. arXiv–2004, 2020.



- Michael E Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Sinkformers: Transformers with doubly stochastic attention. In *International Conference on Artificial Intelligence and Statistics*, pp. 3515–3530. PMLR, 2022.
- 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.
- Bernhard Schölkopf, Alexander J Smola, Francis Bach, et al. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 464–468, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2074. URL <https://aclanthology.org/N18-2074>.
- Zhuoran Shen, Mingyuan Zhang, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Efficient attention: Attention with linear complexities. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3531–3539, 2021.
- Alex J Smola and Bernhard Schölkopf. A tutorial on support vector regression. *Statistics and computing*, 14(3):199–222, 2004.
- Sainbayar Sukhbaatar, Edouard Grave, Guillaume Lample, Herve Jegou, and Armand Joulin. Augmenting self-attention with persistent memory. *arXiv preprint arXiv:1907.01470*, 2019.
- Siqi Sun, Yu Cheng, Zhe Gan, and Jingjing Liu. Patient knowledge distillation for bert model compression. *arXiv preprint arXiv:1908.09355*, 2019.
- Chang Wei Tan, C. Bergmeir, François Petitjean, and Geoffrey I. Webb. Monash university, uea, ucr time series regression archive. *ArXiv*, abs/2006.10996, 2020.
- Binh Tang and David S. Matteson. Probabilistic transformer for time series analysis. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=HfpNVDg3ExA>.
- Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse Sinkhorn attention. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 9438–9447. PMLR, 13–18 Jul 2020. URL <http://proceedings.mlr.press/v119/tay20a.html>.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for efficient transformers. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=qVyeW-grC2k>.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4593–4601, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1452. URL <https://www.aclweb.org/anthology/P19-1452>.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers and distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers distillation through attention, 2021.
- 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. *arXiv preprint arXiv:1908.11775*, 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.
- Jesse Vig and Yonatan Belinkov. Analyzing the structure of attention in a transformer language model. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 63–76, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-4808. URL <https://www.aclweb.org/anthology/W19-4808>.
- 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. *arXiv preprint arXiv:1905.09418*, 2019.
- Sinong Wang, Belinda Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7794–7803, 2018.
- Zifeng Wang and Jimeng Sun. TransTab: Learning Transferable Tabular Transformers Across Tables. In *Advances in Neural Information Processing Systems (NeurIPS 2022)*, 2022.
- Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018.
- Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22–31, 2021.
- Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Flowformer: Linearizing transformers with conservation flows. In *International Conference on Machine Learning*, 2022.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237*, 2019.
- George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff. A transformer-based framework for multivariate time series representation learning. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 2114–2124, 2021.
- Shaolei Zhang and Yang Feng. Modeling concentrated cross-attention for neural machine translation with Gaussian mixture model. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 1401–1411, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-emnlp.121. URL <https://aclanthology.org/2021.findings-emnlp.121>.
- Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. Deep learning based recommender system: A survey and new perspectives. *ACM Computing Surveys (CSUR)*, 52(1):1–38, 2019.
- Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16259–16268, 2021.
- Chen Zhu, Wei Ping, Chaowei Xiao, Mohammad Shoeybi, Tom Goldstein, Anima Anandkumar, and Bryan Catanzaro. Long-short transformer: Efficient transformers for language and vision. In *NeurIPS*, 2021.

## Supplement to “A Primal-Dual Framework for Transformers and Neural Networks”

### A ADDITIONAL DETAILS ON THE EXPERIMENTS

This section provides datasets, models, and training details for experiments in Section 3. As mentioned in Section 3, for Attention-BN models, recentering queries and keys alone is sufficient for accuracy improvement, and we weight the mean  $\mu$  in Eqn 22 with a constant  $\beta$ . Hence Eqn 22 is simplified to:

$$h_i = \sum_{j=1}^N \text{softmax} \left( (q_i - \beta\mu)^\top (k_j - \beta\mu) / \sqrt{D} \right) v_j. \quad (26)$$

In our experiments, we consider the constant  $\beta$  in Attention-BN/BN+SH and the different downsampling scales in Attention-SH/SH+BN as hyper-parameters to finetune. All of our experiments are conducted on a server with 4 NVIDIA A100 GPUs.

#### A.1 UEA TIME SERIES CLASSIFICATION

**Datasets and metrics** The benchmark (Bagnall et al., 2018) consists of 30 datasets. Following (Wu et al., 2022), we choose 10 datasets, which vary in input sequence lengths, the number of classes, and dimensionality, to evaluate our models on temporal sequences. We report the test accuracy as evaluation for the benchmark.

**Models and baselines** The experiment setups and configurations for the softmax/linear baseline and our models are the same as in (Wu et al., 2022) <sup>1</sup> (for the PEMS-SF, SelfRegulationSCP2, UWaveGestureLibrary datasets) and (Zerveas et al., 2021) <sup>2</sup> (for the other tasks). In all models, the number of heads is 8, whereas the model dimension and number of transformer layers are varied. For Attention-SH/SH+BN, we downsample keys and values by the factor of 2, after every two successive heads.

#### A.2 LONG RANGE ARENA BENCHMARK

**Datasets and metrics** We adopt the tasks: Listops (Nangia & Bowman, 2018), byte-level IMDb reviews text classification (Maas et al., 2011), byte-level document retrieval (Radev et al., 2013), CIFAR-10 image classification (Krizhevsky et al., 2009) and the Pathfinder challenge (Linsley et al., 2018) in the LRA benchmark for our experiments. They consist of long sequences of length  $2K$ ,  $4K$ ,  $4K$ ,  $1K$ , and  $1K$  respectively. The evaluation protocol and metric are the same as in (Tay et al., 2021).

**Models and baselines** All our models and softmax/linear baselines follow the same architecture and configuration as in (Zhu et al., 2021) <sup>3</sup>. Each model consists of two layers and 64 embedding dimensions. While one head at each layer remains intact, the keys and values of the other heads are halved in our Attention-SH/SH+BN experiments.

#### A.3 IMAGE CLASSIFICATION ON IMAGENET

**Datasets and metrics** The ImageNet dataset (Deng et al., 2009; Russakovsky et al., 2015) consists of  $1.28M$  training images and  $50K$  validation images. The task is to classify 1000 categories. Top-1 and top-5 accuracies are reported.

**Models and baselines** Our baseline is DeiT-tiny model (Touvron et al., 2021) with 12 transformer layers, 4 attention heads per layer, and the model dimension of 192. For model setting and setting and configuration, we follow (Touvron et al., 2021) <sup>4</sup>. The downsampling scales in Attention-SH/BN+SH models are  $[1, 1, 2, 4]$  for 4 heads at each layer, respectively.

<sup>1</sup>Implementation available at <https://github.com/thuml/Flowformer>.

<sup>2</sup>Implementation available at <https://github.com/gzerveas/mvts.transformer>.

<sup>3</sup>Implementation available at <https://github.com/NVIDIA/transformer-ls>.

<sup>4</sup>Implementation available at <https://github.com/facebookresearch/deit>.

Table 5: Test Accuracy (%) of the Linear Attention-BN/SH/BN+SH vs. the baseline Linear Attention (Katharopoulos et al., 2020) on the UEA Time Series Classification Archive benchmark (Bagnall et al., 2018). Our proposed attentions outperform the baseline.

Dataset/Model	Baseline Linear	Linear Attention-BN	Linear Attention-SH	Linear Attention-BN+SH
ETHANOLCONCENTRATION	33.84 $\pm$ 0.66	<b>34.98 <math>\pm</math> 0.74</b>	34.76 $\pm$ 0.69	34.35 $\pm$ 0.70
FACEDETECTION	69.17 $\pm$ 0.32	69.22 $\pm$ 0.17	<b>69.38 <math>\pm</math> 0.17</b>	69.12 $\pm$ 0.19
HANDWRITING	32.87 $\pm$ 0.27	32.86 $\pm$ 0.49	32.82 $\pm$ 0.12	<b>32.98 <math>\pm</math> 0.36</b>
HEARTBEAT	75.61 $\pm$ 0.73	75.78 $\pm$ 0.71	74.96 $\pm$ 0.62	<b>75.94 <math>\pm</math> 0.68</b>
JAPANESEVOWELS	99.37 $\pm$ 0.16	<b>99.60 <math>\pm</math> 0.19</b>	99.28 $\pm$ 0.41	99.33 $\pm$ 0.19
PEMS-SF	83.43 $\pm$ 0.88	85.74 $\pm$ 0.67	<b>86.51 <math>\pm</math> 0.88</b>	84.97 $\pm$ 0.76
SELFREGULATIONSCP1	90.90 $\pm$ 0.40	91.81 $\pm$ 0.69	90.76 $\pm$ 0.59	<b>91.92 <math>\pm</math> 0.60</b>
SELFREGULATIONSCP2	55.18 $\pm$ 0.89	<b>56.11 <math>\pm</math> 0.94</b>	54.44 $\pm$ 0.88	55.74 $\pm$ 0.92
SPOKENARABICDIGITS	<b>99.07 <math>\pm</math> 0.10</b>	99.01 $\pm$ 0.07	99.03 $\pm$ 0.18	98.91 $\pm$ 0.17
UWAVEGESTURELIBRARY	85.63 $\pm$ 0.81	<b>86.04 <math>\pm</math> 0.86</b>	84.89 $\pm$ 1.00	85.78 $\pm$ 0.75
AVERAGE ACCURACY	72.51 $\pm$ 0.34	<b>73.12 <math>\pm</math> 0.20</b>	72.68 $\pm$ 0.40	72.90 $\pm$ 0.23

Table 6: Top-1 and top-5 accuracy (%) of the Attention-Conv2D Deit vs. the baseline Deit with the softmax attention on the ImageNet image classification task. The Attention-Conv2D Deit significantly outperforms the baseline in both top-1 and top-5 accuracy.

Metric/Model	Baseline Softmax Deit	Attention-Conv2D Deit
Top-1 Acc (%)	72.23 $\pm$ 0.23	<b>73.18 <math>\pm</math> 0.24</b>
Top-5 Acc (%)	91.13 $\pm$ 0.15	<b>91.52 <math>\pm</math> 0.13</b>

Table 7: Test Accuracy (%) of the Attention-Conv1D vs. the baseline softmax attention on 5 tasks of the LRA benchmark (Tay et al., 2021). Our models outperform the softmax baseline.

Dataset/Model	Baseline Softmax	Attention-Conv1D
LISTOPS	36.76 $\pm$ 0.42	<b>37.20 <math>\pm</math> 0.05</b>
TEXT	64.90 $\pm$ 0.07	<b>64.92 <math>\pm</math> 0.44</b>
RETRIEVAL	79.68 $\pm$ 0.52	<b>80.75 <math>\pm</math> 0.15</b>
IMAGE	<b>39.23 <math>\pm</math> 1.35</b>	39.18 $\pm$ 0.59
PATHFINDER	72.72 $\pm$ 0.75	<b>73.01 <math>\pm</math> 0.24</b>
AVERAGE ACCURACY	58.66 $\pm$ 0.26	<b>59.01 <math>\pm</math> 0.20</b>

## B ADDITIONAL EXPERIMENTAL RESULTS

### B.1 UEA TIME SERIES CLASSIFICATION USING THE LINEAR ATTENTION-BN/SH/BN+SH

Table 5 summarizes the comparison between the Linear Attention-BN/SH/BN+SH and the baseline Linear Attention on the UEA Time Series Classification task. The Linear Attention-BN/SH/BN+SH achieve better accuracy than the Linear Attention baseline while being more efficient.

### B.2 CONVOLUTION ATTENTION

Table 6 demonstrates the advantage of Attention-Conv2D (Def. 3, Section G) over softmax Deit on the ImageNet image classification task. Furthermore, as shown in Table 7, the Attention-Conv1D (Def. 4, Section G) outperforms the baseline softmax attention on 5 tasks of the LRA benchmark (Tay et al., 2021).

### B.3 ADDITIONAL EXPERIMENTS ON THE UEA TIMESERIES CLASSIFICATION BENCHMARK AND THE UCR TIME SERIES REGRESSION ARCHIVE

In this section, we further demonstrate the advantage of our Attention-BN/SH/BN+SH on additional 15 tasks in the UEA Time Series Classification benchmark and on 6 tasks in the UCR Time Series Regression benchmark. The results in Table 8 and 9 show that our Attention-BN and Attention-SH+BN outperform the baseline softmax transformers significantly on both of these benchmarks, while the attention-SH has comparable performance with the baseline but being more efficient.

### B.4 UEA TIME SERIES CLASSIFICATION USING THE SPARSE ATTENTION-BN/SH/BN+SH

Table 10 summarizes the comparison between the Sparse Attention-BN/SH/BN+SH and the Sparse Attention baseline on a subset of the UEA Time Series Classification benchmark. Our models when combined with Sparse Attention achieve significantly better accuracy than the Sparse Attention baseline while the Sparse Attention-SH/BN+SH are more efficient (See Fig. 3 and Fig. 4 in Appendix C).

Table 8: Root mean square error (RMSE) of the Attention-BN/SH/BN+SH vs. the baseline softmax attention on 6 UCR Time Series Regression tasks (Tan et al., 2020). Smaller RMSE indicates better performance.

Dataset/Model	Baseline Softmax	Attention-BN	Attention-SH	Attention-BN+SH
APPLIANCESENERGY	3.44 $\pm$ 0.06	3.38 $\pm$ 0.34	3.39 $\pm$ 0.02	<b>3.37 <math>\pm</math> 0.23</b>
BENZENECONCENTRATION	0.91 $\pm$ 0.03	<b>0.89 <math>\pm</math> 0.17</b>	1.00 $\pm$ 0.09	0.90 $\pm$ 0.08
BEIJINGPM10	92.31 $\pm$ 1.06	<b>92.00 <math>\pm</math> 0.89</b>	92.82 $\pm$ 0.92	92.40 $\pm$ 0.85
BEIJINGPM25	59.73 $\pm$ 1.21	59.55 $\pm$ 0.92	59.66 $\pm$ 0.88	<b>59.24 <math>\pm</math> 1.22</b>
LIVEFUELMOISTURE	43.08 $\pm$ 0.17	<b>43.01 <math>\pm</math> 0.50</b>	43.65 $\pm$ 0.09	43.79 $\pm$ 0.49
IEEPPPG	32.12 $\pm$ 1.25	<b>30.69 <math>\pm</math> 0.64</b>	31.38 $\pm$ 1.02	30.73 $\pm$ 1.20
AVERAGE RMSE	38.60 $\pm$ 0.67	<b>38.25 <math>\pm</math> 0.30</b>	38.65 $\pm$ 0.27	38.40 $\pm$ 0.51

Table 9: Accuracy (%) of the Attention-BN/SH/BN+SH vs. the baseline softmax attention on other 15 UEA Time Series classification tasks (Bagnall et al., 2018).

Dataset/Model	Baseline Softmax	Attention-BN	Attention-SH	Attention-BN+SH
ARTICULARYWORDRECOGNITION	97.44 $\pm$ 0.42	98.22 $\pm$ 0.87	97.22 $\pm$ 0.95	<b>98.44 <math>\pm</math> 0.41</b>
BASICMOTIONS	98.75 $\pm$ 1.25	99.38 $\pm$ 1.08	99.37 $\pm$ 1.06	<b>99.78 <math>\pm</math> 0.51</b>
EPILEPSY	<b>93.71 <math>\pm</math> 1.23</b>	92.27 $\pm$ 0.74	89.13 $\pm$ 1.07	92.02 $\pm$ 1.02
ERING	95.18 $\pm$ 0.52	95.18 $\pm$ 0.37	94.72 $\pm$ 0.66	<b>95.46 <math>\pm</math> 0.40</b>
FINGERMOVEMENTS	59.67 $\pm$ 0.47	63.00 $\pm$ 0.41	61.33 $\pm$ 0.70	<b>63.66 <math>\pm</math> 0.64</b>
LIBRAS	85.00 $\pm$ 0.45	<b>85.37 <math>\pm</math> 0.69</b>	83.88 $\pm$ 0.45	85.00 $\pm$ 0.78
NATOPS	95.00 $\pm$ 0.45	95.37 $\pm$ 0.26	<b>96.29 <math>\pm</math> 0.69</b>	95.74 $\pm$ 0.94
RACKETSPORTS	87.28 $\pm$ 0.82	87.93 $\pm$ 0.31	88.16 $\pm$ 0.54	<b>89.03 <math>\pm</math> 0.64</b>
ATRIALFIBRILLATION	33.33 $\pm$ 2.71	41.67 $\pm$ 2.88	35.00 $\pm$ 2.89	<b>41.68 <math>\pm</math> 2.80</b>
CRICKET	94.90 $\pm$ 0.65	95.37 $\pm$ 0.65	93.98 $\pm$ 0.65	<b>96.29 <math>\pm</math> 0.65</b>
STANDWALKJUMP	50.00 $\pm$ 2.33	<b>55.55 <math>\pm</math> 2.14</b>	50.01 $\pm$ 2.34	55.00 $\pm$ 2.08
HANDMOVEMENTDIRECTION	63.96 $\pm$ 2.30	64.41 $\pm$ 2.76	61.71 $\pm$ 2.64	<b>66.66 <math>\pm</math> 2.54</b>
LSST	58.54 $\pm$ 0.54	57.05 $\pm$ 0.26	<b>60.34 <math>\pm</math> 0.73</b>	59.91 $\pm$ 0.34
DUCKDUCKGEESE	64.50 $\pm$ 1.96	65.00 $\pm$ 1.73	64.50 $\pm$ 1.95	<b>65.50 <math>\pm</math> 1.66</b>
MOTORIMAGERY	58.66 $\pm$ 1.25	60.67 $\pm$ 1.69	59.00 $\pm$ 1.41	<b>62.00 <math>\pm</math> 0.81</b>
AVERAGE ACCURACY	75.73 $\pm$ 0.51	77.10 $\pm$ 0.22	75.61 $\pm$ 0.18	<b>77.74 <math>\pm</math> 0.24</b>

Table 10: Test Accuracy (%) of the Sparse Attention-BN/SH/BN+SH vs. the baseline Sparse Attention (Child et al., 2019) on a subset of the UEA Time Series Classification Archive benchmark (Bagnall et al., 2018). Our proposed attentions outperform the baseline.

Dataset/Model	Baseline Sparse	Sparse Attention-BN	Sparse Attention-SH	Sparse Attention-BN+SH
ETHANOLCONCENTRATION	33.33 $\pm$ 1.23	33.33 $\pm$ 0.78	32.50 $\pm$ 0.57	<b>33.46 <math>\pm</math> 0.71</b>
FACEDETECTION	68.58 $\pm$ 0.95	68.65 $\pm$ 0.44	<b>68.67 <math>\pm</math> 0.78</b>	68.44 $\pm$ 0.51
HANDWRITING	31.08 $\pm$ 0.38	31.79 $\pm$ 0.44	32.75 $\pm$ 0.39	<b>33.37 <math>\pm</math> 0.61</b>
HEARTBEAT	74.95 $\pm$ 0.81	75.98 $\pm$ 0.72	74.96 $\pm$ 0.80	<b>76.09 <math>\pm</math> 0.75</b>
JAPANESEVOWELS	99.45 $\pm$ 0.10	<b>99.54 <math>\pm</math> 0.12</b>	99.18 $\pm$ 0.14	99.36 $\pm$ 0.34
PEMS-SF	82.08 $\pm$ 0.63	83.81 $\pm$ 0.47	82.66 $\pm$ 0.63	<b>84.01 <math>\pm</math> 0.89</b>
SELFREGULATIONSCP1	91.24 $\pm$ 0.85	91.69 $\pm$ 0.42	91.47 $\pm$ 0.84	<b>91.70 <math>\pm</math> 0.16</b>
SELFREGULATIONSCP2	55.18 $\pm$ 0.69	<b>58.52 <math>\pm</math> 0.71</b>	55.92 $\pm$ 0.94	56.67 $\pm$ 0.68
SPOKENARABICDIGITS	99.04 $\pm$ 0.06	99.10 $\pm$ 0.15	99.06 $\pm$ 0.13	<b>99.15 <math>\pm</math> 0.09</b>
UWAVEGESTURELIBRARY	84.90 $\pm$ 0.39	85.73 $\pm$ 0.38	85.31 $\pm$ 0.88	<b>86.56 <math>\pm</math> 0.25</b>
AVERAGE ACCURACY	71.98 $\pm$ 0.38	72.81 $\pm$ 0.15	72.25 $\pm$ 0.24	<b>72.88 <math>\pm</math> 0.35</b>

Table 11: Test Accuracy (%) of the Attention-BN/BN+SH with  $\beta$  is learnable or set as a hyperparameter on the retrieval task (Tay et al., 2021).

Model	Retrieval
Attention-BN (learn $\beta$ )	80.77 $\pm$ 0.23
Attention-BN+SH (learn $\beta$ )	<b>81.31 <math>\pm</math> 0.25</b>
Attention-BN ( $\beta$ as a hyperparameter)	81.05 $\pm$ 0.08
Attention-BN+SH ( $\beta$ as a hyperparameter)	81.20 $\pm$ 0.11

## B.5 ATTENTION-BN/BN+SH WITH LEARNABLE $\beta$

We experiment with our Attention-BN/BN+SH with learnable  $\beta$  on the retrieval task. Table 11 shows that learning  $\beta$  does not improve much over setting  $\beta$  to be a hyperparameter.

## C ADDITIONAL RESULTS ON EFFICIENCY ANALYSIS

This section provides more efficiency analysis on our models.

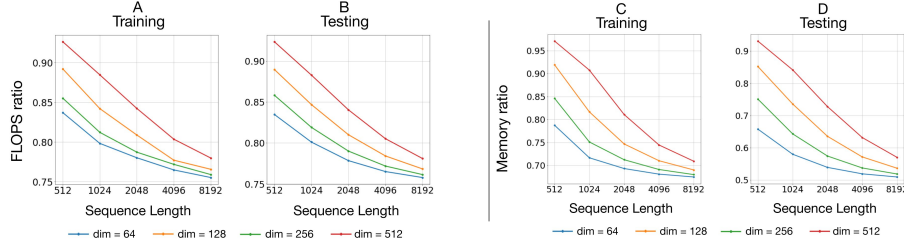


Figure 2: (Left) FLOPS ratios and (Right) memory usage ratios between the Attention-SH and the softmax attention baseline trained on the LRA retrieval task for different model dimensions and sequence lengths.

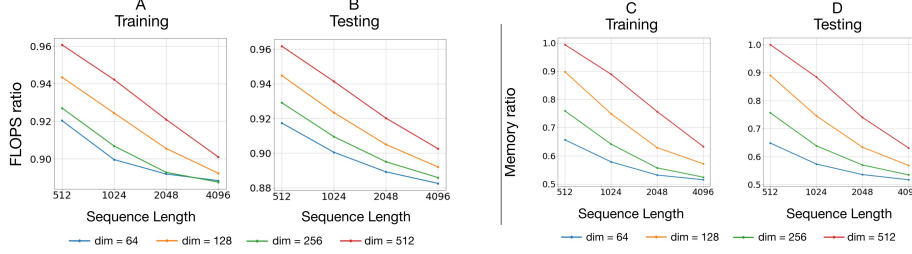


Figure 3: (Left) FLOPS ratios and (Right) memory usage ratios between the Sparse Attention-BN+SH and the Sparse Attention baseline trained on the LRA retrieval task for different model dimensions and sequence lengths. When using our models, the reduction in computation and memory improves with sequence length. When scaling up the model with greater model dimension, our methods remain significantly more efficient than the baseline.

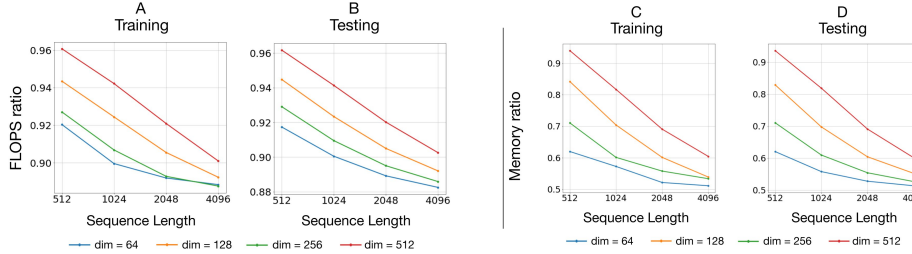


Figure 4: (Left) FLOPS ratios and (Right) memory usage ratios between the Sparse Attention-SH and the Sparse Attention baseline trained on the LRA retrieval task for different model dimensions and sequence lengths. When using our models, the reduction in computation and memory improves with sequence length. When scaling up the model with greater model dimension, our methods remain significantly more efficient than the baseline.

**Attention-SH.** Fig.2 shows the efficiency benefits of our Attention-SH when trained on the retrieval task. Same as in the case of Attention-SH+BN, the efficiency benefits of our Attention-SH over the baseline Softmax attention grows when  $N$  and  $D$  increase.

**Sparse Attention-SH/BN+SH.** Fig.3 and Fig.4 show that the efficiency advantages of our Sparse Attention-BN+SH and Sparse Attention-SH, respectively, increase as the model dimension  $D$  and sequence length  $N$  grow. All models are trained on the LRA retrieval task. In addition to the efficiency advantage, the Sparse Attention-BN+SH also significantly outperforms the Sparse Attention baseline in terms of accuracy in this task (79.86% vs. 78.20%) while the Sparse Attention-SH achieves a comparable result to the baseline. More accuracy advantages of the Sparse Attention-BN/SH/BN+SH over the Sparse Attention baseline are given in Table 10.

## D DERIVING SOFTMAX ATTENTION.

Choosing the appropriate  $h(x)$  and  $\Phi(x)$  allows us to derive the popular softmax attention given in Eqn. 1 and 2. In particular, if we choose  $h(x) := \sum_{j=1}^N \Phi(x)^T \Phi(k_j)$ , Eqn. 10 becomes

$$f(x) = \sum_{j=1}^N \frac{\Phi(x)^T \Phi(k_j)}{\sum_{j'=1}^N \Phi(x)^T \Phi(k_{j'})} v_j + b = \frac{\sum_{j=1}^N \Phi(x)^T \Phi(k_j) v_j}{\sum_{j'=1}^N \Phi(x)^T \Phi(k_{j'})} + b. \quad (27)$$

We then select  $\Phi(\mathbf{x}) = (a_{l_0}^{(0)}, a_1^{(1)}, \dots, a_{l_1}^{(1)}, \dots, a_1^{(t)}, \dots, a_{l_t}^{(t)}, \dots)$  where  $l_t = \binom{D+t-1}{t}$  and

$$a_l^{(t)} = \frac{(x_1/\sqrt[4]{D})^{n_1} \dots (x_D/\sqrt[4]{D})^{n_D}}{\sqrt{n_1! \dots n_D!}} \mid n_1 + \dots + n_D = t, 1 \leq l \leq l_t. \quad (28)$$

Since

$$\exp(\mathbf{x}^T \mathbf{y}) = \sum_{t=0}^{\infty} \frac{(\mathbf{x}^T \mathbf{y})^t}{t!} = \sum_{t=0}^{\infty} \sum_{n_1 + \dots + n_D = t} \left( \frac{x_1^{n_1} \dots x_D^{n_D}}{\sqrt{n_1! \dots n_D!}} \right) \left( \frac{y_1^{n_1} \dots y_D^{n_D}}{\sqrt{n_1! \dots n_D!}} \right), \quad (29)$$

then Eqn. 27 becomes

$$\begin{aligned} f(\mathbf{x}) &= \sum_{j=1}^N \frac{\sum_{t=0}^{\infty} \sum_{n_1 + \dots + n_D = t} \left( \frac{(x_1/\sqrt[4]{D})^{n_1} \dots (x_D/\sqrt[4]{D})^{n_D}}{\sqrt{n_1! \dots n_D!}} \right) \left( \frac{(k_{j1}/\sqrt[4]{D})^{n_1} \dots (k_{jD}/\sqrt[4]{D})^{n_D}}{\sqrt{n_1! \dots n_D!}} \right)}{\sum_{j'=1}^N \sum_{t=0}^{\infty} \sum_{n_1 + \dots + n_D = t} \left( \frac{(x_1/\sqrt[4]{D})^{n_1} \dots (x_D/\sqrt[4]{D})^{n_D}}{\sqrt{n_1! \dots n_D!}} \right) \left( \frac{(k_{j'1}/\sqrt[4]{D})^{n_1} \dots (k_{j'D}/\sqrt[4]{D})^{n_D}}{\sqrt{n_1! \dots n_D!}} \right)} \mathbf{v}_j + \mathbf{b} \\ &= \sum_{j=1}^N \frac{\exp\left(\left(\frac{\mathbf{x}}{\sqrt[4]{D}}\right)^\top \frac{\mathbf{k}_j}{\sqrt[4]{D}}\right)}{\sum_{j'=1}^N \exp\left(\left(\frac{\mathbf{x}}{\sqrt[4]{D}}\right)^\top \frac{\mathbf{k}_{j'}}{\sqrt[4]{D}}\right)} \mathbf{v}_j + \mathbf{b} = \sum_{j=1}^N \frac{\exp(\mathbf{x}^\top \mathbf{k}_j / \sqrt{D})}{\sum_{j'=1}^N \exp(\mathbf{x}^\top \mathbf{k}_{j'} / \sqrt{D})} \mathbf{v}_j + \mathbf{b}. \end{aligned} \quad (30)$$

Let  $\mathbf{x} = \mathbf{q}_i$ ,  $\mathbf{b} = 0$  and relax the boundness constraint of  $\mathbf{v}_j$  in Remark 1. Eqn. 30 becomes Eqn. 2 of the softmax attention (Vaswani et al., 2017).

## E BATCH NORMALIZED ATTENTION: DERIVATION OF EQN. 24

$$\begin{aligned} h_i &= \sum_{j=1}^N \text{softmax} \left( \sum_{d=1}^D \frac{(\mathbf{q}_i(d) - \boldsymbol{\mu}(d))(\mathbf{k}_j(d) - \boldsymbol{\mu}(d))}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \mathbf{v}_j \\ &= \sum_{j=1}^N \text{softmax} \left( \sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_j(d) - \mathbf{q}_i(d)\boldsymbol{\mu}(d) - \boldsymbol{\mu}(d)\mathbf{k}_j(d) + \boldsymbol{\mu}(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \mathbf{v}_j \\ &= \sum_{j=1}^N \frac{\exp\left(\sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_j(d) - \mathbf{q}_i(d)\boldsymbol{\mu}(d) - \boldsymbol{\mu}(d)\mathbf{k}_j(d) + \boldsymbol{\mu}(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right)}{\sum_{j'=1}^N \exp\left(\sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_{j'}(d) - \mathbf{q}_i(d)\boldsymbol{\mu}(d) - \boldsymbol{\mu}(d)\mathbf{k}_{j'}(d) + \boldsymbol{\mu}(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right)} \mathbf{v}_j \\ &= \sum_{j=1}^N \frac{\exp\left(\sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_j(d) - \boldsymbol{\mu}(d)\mathbf{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right) \exp\left(\sum_{d=1}^D \frac{\boldsymbol{\mu}(d)\boldsymbol{\mu}(d) - \mathbf{q}_i(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right)}{\sum_{j'=1}^N \exp\left(\sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_{j'}(d) - \boldsymbol{\mu}(d)\mathbf{k}_{j'}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right) \exp\left(\sum_{d=1}^D \frac{\boldsymbol{\mu}(d)\boldsymbol{\mu}(d) - \mathbf{q}_i(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right)} \mathbf{v}_j \\ &= \sum_{j=1}^N \frac{\exp\left(\sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_j(d) - \boldsymbol{\mu}(d)\mathbf{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right)}{\sum_{j'=1}^N \exp\left(\sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_{j'}(d) - \boldsymbol{\mu}(d)\mathbf{k}_{j'}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)}\right)} \mathbf{v}_j \\ &= \sum_{j=1}^N \text{softmax} \left( \sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_j(d) - \boldsymbol{\mu}(d)\mathbf{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \mathbf{v}_j \\ &= \sum_{j=1}^N \text{softmax} \left( \sum_{d=1}^D \frac{\mathbf{q}_i(d)\mathbf{k}_j(d) - \frac{1}{N} \sum_{j'=1}^N \mathbf{k}_{j'}(d)\mathbf{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \mathbf{v}_j. \end{aligned} \quad (31)$$

## F ATTENTION WITH THE RESIDUAL CONNECTION AND MATRIX PROJECTIONS

In this supplement, we first discuss attention with the residual connection and matrix projections in Appendix F.

Suppose we are given a training data  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\} \subset \mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{X} = \mathbb{R}^{D_x}$  and  $\mathcal{Y} = \mathbb{R}^{D_y}$ . Here,  $\mathbf{x}_1, \dots, \mathbf{x}_N$  are the training inputs, and  $\mathbf{y}_1, \dots, \mathbf{y}_N$  are the training targets. In



order to derive the attention with the residual connection and query, key, and value matrix projections, we define the function  $f$  as follows

$$\mathbf{y} = f(\mathbf{x}) := \mathbf{W} \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x})}{h(\mathbf{x})} + \mathbf{x} + \mathbf{b}, \quad (32)$$

where  $\mathbf{x} \in \mathcal{X} = \mathbb{R}^{D_x}$ ,  $\mathbf{W}^{\text{proj}} = [\mathbf{w}_1^{\text{proj}}, \dots, \mathbf{w}_D^{\text{proj}}]^\top \in \mathbb{R}^{D \times D_x}$ ,  $\Phi(\cdot) = [\phi_1(\cdot), \dots, \phi_{D_\phi}(\cdot)] : \mathbb{R}^D \rightarrow \mathbb{R}^{D_\phi}$ ,  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{D_v}]^\top \in \mathbb{R}^{D_v \times D_\phi}$ ,  $\mathbf{b} \in \mathbb{R}^{D_v}$ , and  $h(\mathbf{x})$  is a vector-scalar function. We fit the function  $f$  to the training data  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$  with an  $\mathbb{L}_2$  regularization on  $\mathbf{W}$  and  $\mathbf{W}^{\text{proj}}$  by solving the following convex optimization problem:

$$\begin{aligned} & \underset{\substack{\mathbf{W}, \mathbf{W}^{\text{proj}} \\ \boldsymbol{\xi}_j, \tilde{\boldsymbol{\xi}}_j, j=1, \dots, N}}{\text{minimize}} & \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d\|^2 + \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d^{\text{proj}}\|^2 + C \sum_{j=1}^N \sum_{d=1}^{D_v} (\boldsymbol{\xi}_j(d) + \tilde{\boldsymbol{\xi}}_j(d)) \\ & \text{subject to} & \begin{cases} \mathbf{y}_j(d) - \mathbf{w}_d^\top \Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j) / h(\mathbf{x}_j) - \mathbf{x}_j - \mathbf{b}(d) \leq \epsilon + \boldsymbol{\xi}_j(d) \\ \mathbf{w}_d^\top \Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j) / h(\mathbf{x}_j) + \mathbf{x}_j + \mathbf{b}(d) - \mathbf{y}_j(d) \leq \epsilon + \tilde{\boldsymbol{\xi}}_j(d) \\ \boldsymbol{\xi}_j(d), \tilde{\boldsymbol{\xi}}_j(d) \geq 0 \end{cases}, \quad j = 1, \dots, N, \quad d = 1, \dots, D_v. \end{aligned} \quad (33)$$

The Lagrangian of the optimization problem 33 is given by

$$\begin{aligned} \mathcal{L}_1 := & \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d\|^2 + \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d^{\text{proj}}\|^2 + C \sum_{j=1}^N \sum_{d=1}^{D_v} (\boldsymbol{\xi}_j(d) + \tilde{\boldsymbol{\xi}}_j(d)) - \sum_{j=1}^N \sum_{d=1}^{D_v} (\boldsymbol{\eta}_j(d) \boldsymbol{\xi}_j(d) + \tilde{\boldsymbol{\eta}}_j(d) \tilde{\boldsymbol{\xi}}_j(d)) \\ & - \sum_{j=1}^N \sum_{d=1}^{D_v} \alpha_j(d) \left( \epsilon + \boldsymbol{\xi}_j(d) - \mathbf{y}_j(d) + \mathbf{w}_d^\top \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x}_j)} + \mathbf{x}_j + \mathbf{b}(d) \right) \\ & - \sum_{j=1}^N \sum_{d=1}^{D_v} \tilde{\alpha}_j(d) \left( \epsilon + \tilde{\boldsymbol{\xi}}_j(d) + \mathbf{y}_j(d) - \mathbf{w}_d^\top \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x}_j)} - \mathbf{x}_j - \mathbf{b}(d) \right), \end{aligned} \quad (34)$$

Similar to the derivation in Section 2.1, the partial derivatives of  $\mathcal{L}_1$  with respect to the primal variable  $\mathbf{w}_d$ ,  $d = 1, \dots, D_v$ , have to vanish for optimality, which leads to

$$\partial_{\mathbf{w}_d} \mathcal{L}_1 = \mathbf{w}_d - \sum_{j=1}^N (\alpha_j(d) - \tilde{\alpha}_j(d)) \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x}_j)} = 0 \Rightarrow \mathbf{w}_d = \sum_{j=1}^N (\alpha_j(d) - \tilde{\alpha}_j(d)) \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x}_j)}. \quad (35)$$

Note that here we only find the form of the optimal solution for  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{D_v}]^\top$ . The optimal value of  $\mathbf{W}^{\text{proj}}$  can then be found by optimization algorithm such as the (stochastic) gradient descent when training the transformer.

Let  $\mathbf{v}_j = [\frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(\mathbf{x}_j)}, \dots, \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(\mathbf{x}_j)}]^\top$ ,  $j = 1, \dots, N$ , we obtain the following support vector expansion of the function  $f$ :

$$\begin{aligned} f(\mathbf{x}) &= \left[ \sum_{j=1}^N (\alpha_j(1) - \tilde{\alpha}_j(1)) \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x}_j)}, \dots, \sum_{j=1}^N (\alpha_j(D_v) - \tilde{\alpha}_j(D_v)) \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x}_j)} \right]^\top \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x})}{h(\mathbf{x})} + \mathbf{x} + \mathbf{b}, \\ &= \left[ \sum_{j=1}^N \frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(\mathbf{x}_j)} \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x})^\top \Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x})}, \dots, \sum_{j=1}^N \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(\mathbf{x}_j)} \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x})^\top \Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x})} \right]^\top + \mathbf{x} + \mathbf{b}, \\ &= \underbrace{\sum_{j=1}^N \frac{\Phi(\mathbf{W}^{\text{proj}} \mathbf{x})^\top \Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)}{h(\mathbf{x})} \mathbf{v}_j}_{\text{Residual connection}} + \mathbf{x} + \mathbf{b}. \end{aligned} \quad (36)$$

Here, the support vector expansion of  $f$  already includes a residual connection. The softmax attention can then be derived by selecting  $h(\mathbf{x}) := \sum_{j=1}^N \Phi(\mathbf{W}^{\text{proj}} \mathbf{x})^\top \Phi(\mathbf{W}^{\text{proj}} \mathbf{x}_j)$  and choosing  $\Phi$  as in Eqn. 28 in Section 2.1. Note that in Eqn. 36,  $\{\mathbf{x}_j\}_{j=1}^N$  and  $\mathbf{x}$  are the training samples and test sample, respectively. In order to derive the key, query, and value matrix projections in attention, we can then relax Eqn. 36 by letting  $\mathbf{W}^{\text{proj}} \mathbf{x}_j = \mathbf{W}_K \mathbf{x}_j$ ,  $\mathbf{W}^{\text{proj}} \mathbf{x} = \mathbf{W}_Q \mathbf{x}$ ,  $\mathbf{v}_j = \mathbf{W}_V \mathbf{x}_j$  and choosing the test sample  $\mathbf{x}$  among the training samples  $\{\mathbf{x}_j\}_{j=1}^N$ .

**Remark 6** Here, for self-attention, we choose the test sample  $\mathbf{x}$  among the training samples  $\{\mathbf{x}_j\}_{j=1}^N$  to compute the attention score of a token to other tokens in the same sequence. For cross-attention where a token in a sequence attends to tokens in another sequence, this constraint can be removed.

## G 2D-CONVOLUTION ATTENTION

In this section, we discuss attention with 2D-convolution. Suppose we are given a training data  $\{(\mathbf{x}_1^{train}, \mathbf{y}_1^{train}), \dots, (\mathbf{x}_{N_H \times N_W}^{train}, \mathbf{y}_{N_H \times N_W}^{train})\} \subset \mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{X} = \mathbb{R}^{D_x}$  and  $\mathcal{Y} = \mathbb{R}^{D_v}$ . Here,  $\mathbf{x}_1^{train}, \dots, \mathbf{x}_{N_H \times N_W}^{train}$  are the training inputs, and  $\mathbf{y}_1^{train}, \dots, \mathbf{y}_{N_H \times N_W}^{train}$  are the training targets. Let  $\mathbf{X}^{train} \in \mathbb{R}^{N_H \times N_W \times D_x}$  be the 3D-tensor of training inputs, where  $\mathbf{X}^{train}(h, w, d) = \mathbf{x}_{N_W \times (h-1) + w}^{train}(d)$ . Given a new set of inputs  $\{\mathbf{x}_1, \dots, \mathbf{x}_{N_H \times N_W}\} \subset \mathcal{X}$  and the corresponding 3D-tensor  $\mathbf{X} \in \mathbb{R}^{N_H \times N_W \times D_x}$  of these inputs, where  $\mathbf{X}(h, w, d) = \mathbf{x}_{N_W \times (h-1) + w}(d)$ . We consider the function  $f$  applying on the 3D-tensor  $\mathbf{X}$  and taking the following form

$$f(\mathbf{x}_i) = \mathbf{W} \frac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}, s))(i))}{h(\mathbf{x}_i)}, \quad i = 1, \dots, N_H \times N_W \quad (37)$$

where Conv2D is the depth-wise 2D-convolution (Howard et al., 2017), with the kernel size  $s \times s$  and identical kernel channels, applied on the input tensor  $\mathbf{X}$ . Here, the last dimension of  $\mathbf{X}$ , i.e.,  $D_x$ , is the depth. Also,  $\Phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_{D_\phi}(\mathbf{x})] \in \mathbb{R}^{D_\phi}$ ,  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_{D_v}]^\top \in \mathbb{R}^{D_v \times D_\phi}$ ,  $\mathbf{b} \in \mathbb{R}^{D_v}$ , and  $h$  is a vector-scalar function. We fit the function  $f$  to the training data  $\{(\mathbf{x}_1^{train}, \mathbf{y}_1^{train}), \dots, (\mathbf{x}_{N_H \times N_W}^{train}, \mathbf{y}_{N_H \times N_W}^{train})\}$  with an  $\mathbb{L}_2$  regularization on  $\mathbf{W}$ , i.e., a ridge regression, by solving the following convex optimization problem:

$$\begin{aligned} & \underset{\mathbf{W}}{\text{minimize}} && \frac{1}{2} \|\mathbf{W}\|_F^2 + C \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} (\xi_j(d) + \tilde{\xi}_j(d)) = \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d\|^2 + C \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} (\xi_j(d) + \tilde{\xi}_j(d)) \\ & \text{subject to} && \begin{cases} \mathbf{y}_j^{train}(d) - \mathbf{w}_d^\top \frac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\mathbf{x}_j^{train})} - \mathbf{b}(d) \leq \epsilon + \xi_j(d) \\ \mathbf{w}_d^\top \frac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\mathbf{x}_j^{train})} + \mathbf{b}(d) - \mathbf{y}_j^{train}(d) \leq \epsilon + \tilde{\xi}_j(d) \\ \xi_j(d), \tilde{\xi}_j(d) \geq 0, \quad j = 1, \dots, N_H \times N_W, \quad d = 1, \dots, D_v. \end{cases} \end{aligned} \quad (38)$$

The Lagrangian of the optimization problem 38 is given by

$$\begin{aligned} \mathcal{L} := & \frac{1}{2} \sum_{d=1}^{D_v} \|\mathbf{w}_d\|^2 + C \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} (\xi_j(d) + \tilde{\xi}_j(d)) - \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} (\eta_j(d) \xi_j(d) + \tilde{\eta}_j(d) \tilde{\xi}_j(d)) \\ & - \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} \alpha_j(d) \left( \epsilon + \xi_j(d) - \mathbf{y}_j^{train}(d) + \mathbf{w}_d^\top \frac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\mathbf{x}_j^{train})} + \mathbf{b}(d) \right) \\ & - \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} \tilde{\alpha}_j(d) \left( \epsilon + \tilde{\xi}_j(d) + \mathbf{y}_j^{train}(d) - \mathbf{w}_d^\top \frac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\mathbf{x}_j^{train})} - \mathbf{b}(d) \right), \end{aligned} \quad (39)$$

Similar to the derivation in Section 2.1 in the main text, the partial derivatives of  $\mathcal{L}$  with respect to the primal variable  $\mathbf{w}_d$ ,  $d = 1, \dots, D_v$ , have to vanish for optimality, which leads to

$$\partial_{\mathbf{w}_d} \mathcal{L} = \mathbf{w}_d - \sum_{j=1}^{N_H \times N_W} (\alpha_j(d) - \tilde{\alpha}_j(d)) \frac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\mathbf{x}_j^{train})} = 0 \quad (40)$$

$$\Rightarrow \mathbf{w}_d = \sum_{j=1}^{N_H \times N_W} (\alpha_j(d) - \tilde{\alpha}_j(d)) \frac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\mathbf{x}_j^{train})}. \quad (41)$$

Let  $\mathbf{v}_j = [\frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(\mathbf{x}_j^{train})}, \dots, \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(\mathbf{x}_j^{train})}]^\top$ ,  $j = 1, \dots, N_H \times N_W$ , and substitute Eqn. 41 into Eqn. 38, we obtain the following support vector expansion of the linear basis function  $f$ :

$$\begin{aligned} f(\mathbf{x}_i) &= \left[ \sum_{j=1}^{N_H \times N_W} \frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(\mathbf{x}_j^{train})} \frac{\mathbf{A}_{ij}}{h(\mathbf{x}_i)}, \dots, \sum_{j=1}^{N_H \times N_W} \frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(\mathbf{x}_j^{train})} \frac{\mathbf{A}_{ij}}{h(\mathbf{x}_i)} \right]^\top + \mathbf{b}, \\ &= \sum_{j=1}^{N_H \times N_W} \frac{\mathbf{A}_{ij}}{h(\mathbf{x}_i)} \mathbf{v}_j + \mathbf{b}, \end{aligned} \quad (42)$$

where  $\mathbf{A}_{ij} := \Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}, s))(i))^\top \Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))$ .

Same as in Section 2.1, we set  $\mathbf{b}_s = 0$ . To derive the softmax normalization in attention, we choose  $h(\mathbf{x}_i) := \sum_{j=1}^N \mathbf{A}_{ij}$  and select  $\Phi$  as in Eqn. 28. Let the training inputs  $\{\mathbf{x}_1^{train}, \dots, \mathbf{x}_{N_H \times N_W}^{train}\} \subset \mathcal{X}$  be the attention keys  $\{\mathbf{k}_1, \dots, \mathbf{k}_{N_H \times N_W}\} \subset \mathcal{K}$ , where  $\mathcal{K} = \mathbb{R}^D$ , in self-attention. Also, let the new inputs  $\{\mathbf{x}_1, \dots, \mathbf{x}_{N_H \times N_W}\} \subset \mathcal{X}$  be the attention queries  $\{\mathbf{q}_1, \dots, \mathbf{q}_{N_H \times N_W}\} \subset \mathcal{K}$  in self-attention. We define the 2D-Convolution Attention (Attention-Conv2D) as follows:

**Definition 3 (2D-Convolution Attention)** Given a set of the key and value vectors  $\{\mathbf{k}_j, \mathbf{v}_j\}_{j=1}^{N_H \times N_W}$ , and a set of the query vectors  $\{\mathbf{q}_i\}_{i=1}^{N_H \times N_W}$ . Denote the key tensor and the query tensor by  $\mathbf{K} \in \mathbb{R}^{N_H \times N_W \times D}$  and  $\mathbf{Q} \in \mathbb{R}^{N_H \times N_W \times D}$ , respectively, where  $\mathbf{K}(h, w, d) = \mathbf{k}_{N_W \times (h-1) + w}(d)$  and  $\mathbf{Q}(h, w, d) = \mathbf{q}_{N_W \times (h-1) + w}(d)$ . The 2D-Convolution Attention (Attention-Conv2D) computes the corresponding output vector  $\mathbf{h}_i$  of the query  $\mathbf{q}_i$  by the following attention formula:

$$\mathbf{h}_i = \sum_{j=1}^N \text{softmax} \left( \text{Flatten}(\text{Conv2D}(\mathbf{Q}, s))(i)^\top \text{Flatten}(\text{Conv2D}(\mathbf{K}, s))(j) / \sqrt{D} \right) \mathbf{v}_j, \quad (43)$$

where the  $\text{Conv2D}(\cdot, s)$  is the depth-wise 2D-convolution (Howard et al., 2017) with the kernel size  $s \times s$  and identical kernel channels.

**Remark 7 (Convolutional Projection for Attention in the Convolutional vision Transformer)**

The convolutional projections used in the Convolutional vision Transformer (CvT) (Wu et al., 2021) can be derived from Eqn. 42 by letting the training input tensor  $\mathbf{X}^{train}$  to be the 2D input matrix of size  $N \times D_x$  of the self-attention layer (see Section 1.1 in the main text) reshaped into a 3D tensor of size  $N_H \times N_W \times D_x$  where  $N = N_H \times N_W$ . Here, to avoid confusion, we denote the input of the self-attention layer by  $\mathbf{X}^{input}$  and its reshaped version by  $\text{Reshape2D}(\mathbf{X}^{input})$ . We then replace the depth-wise 2D-convolution by the depth-wise separable 2D-convolution in (Wu et al., 2021) and remove the constraint that the kernels have identical channels. In order to derive the convolutional projections for the keys, queries, and values in CvT, for  $i, j = 1, \dots, N$ , we let

$$\begin{aligned} \mathbf{k}_j &= \text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j) = \Phi(\text{Flatten}(\text{Conv2D}(\text{Reshape2D}(\mathbf{X}^{input}), s, \mathbf{W}_K))(j)), \\ \mathbf{q}_i &= \text{Flatten}(\text{Conv2D}(\mathbf{X}, s))(i) = \Phi(\text{Flatten}(\text{Conv2D}(\text{Reshape2D}(\mathbf{X}^{input}), s, \mathbf{W}_Q))(i)), \\ \mathbf{v}_j &= \text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j) = \Phi(\text{Flatten}(\text{Conv2D}(\text{Reshape2D}(\mathbf{X}^{input}), s, \mathbf{W}_V))(j)). \end{aligned} \quad (44)$$

Here, we specify the kernel/filter  $\mathbf{W}_K$ ,  $\mathbf{W}_Q$ , and  $\mathbf{W}_V$  to emphasize that the convolutional projections in CvT uses different kernels to compute keys, queries, and values in self-attention. Eqn. 44 matches the convolutional projects in CvT. By choosing  $h$  and  $\Phi$  similar to above, we can derive the convolutional attention in CvT.

## H 1D-CONVOLUTION ATTENTION

Following the derivation for the Attention-Conv2D in Appendix G above, we can derive the 1D-Convolution Attention (Attention-Conv1D) in a similar way by letting  $\mathbf{X}^{train} \in \mathbb{R}^{N \times D_x}$  and  $\mathbf{X} \in \mathbb{R}^{N \times D_x}$  be 2D-matrices of training inputs and new inputs, respectively, and by replacing Conv2D by Conv1D, which is the depth-wise 1D-convolution, with the kernel size  $s \times 1$  and identical kernel channels, applied on the input tensor  $\mathbf{X}$ . Here, the last dimension of  $\mathbf{X}$ , i.e.,  $D_x$ , is the depth. We define the 1D-Convolution Attention (Attention-Conv1D) as follows:

**Definition 4 (1D-Convolution Attention)** Given a set of the key and value vectors  $\{\mathbf{k}_j, \mathbf{v}_j\}_{j=1}^N$ , and a set of the query vectors  $\{\mathbf{q}_i\}_{i=1}^N$ . Denote the key matrix and the query matrix by  $\mathbf{K} :=$

$[\mathbf{k}_1, \dots, \mathbf{k}_N]^\top \in \mathbb{R}^{N \times D}$  and  $\mathbf{Q} := [\mathbf{q}_1, \dots, \mathbf{q}_N]^\top \in \mathbb{R}^{N \times D}$ , respectively. The 1D-Convolution Attention (Attention-Conv1D) computes the corresponding output vector  $\mathbf{h}_i$  of the query  $\mathbf{q}_i$  by the following attention formula:

$$\mathbf{h}_i = \sum_{j=1}^N \text{softmax} \left( \text{Conv1D}(\mathbf{Q}, s)(i)^\top \text{Conv1D}(\mathbf{K}, s)(j) / \sqrt{D} \right) \mathbf{v}_j, \quad (45)$$

where the  $\text{Conv1D}(\cdot, s)$  is the depth-wise 1D-convolution with the kernel size  $s \times 1$  and identical kernel channels.

## I ATTENTION WITH BATCH NORMALIZATION AND SCALED HEADS

The Attention-BN+SH combines both the Attention-BN and Attention-SH. The Attention-BN+SH fits the function  $f^s$ ,  $s = 1, \dots, H$ , in Eqn. 17 with training sets  $\{(\mathbf{k}_1^1, \mathbf{y}_1^1), \dots, (\mathbf{k}_{N_1}^1, \mathbf{y}_{N_1}^1)\}, \dots, \{(\mathbf{k}_1^H, \mathbf{y}_1^H), \dots, (\mathbf{k}_{N_H}^H, \mathbf{y}_{N_H}^H)\} \subset \mathcal{K} \times \mathcal{Y}$  of different sizes  $N_1, \dots, N_H$ , where  $\mathcal{K} = \mathbb{R}^D$  and  $\mathcal{Y} = \mathbb{R}^{D_v}$ . The function  $f^s$  is defined as:

$$f^s(\mathbf{x}) := \mathbf{W}^s \frac{\Phi((\mathbf{x} - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})}{h^s((\mathbf{x} - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})} + \mathbf{b}^s, \quad (46)$$

where

$$\boldsymbol{\mu}^s = \frac{1}{N_s} \sum_{j=1}^{N_s} \mathbf{k}_j^s, \quad \mathbf{s}^{s-1} = \left[ \frac{1}{\sqrt{\sigma_1^{s^2} + \epsilon}}, \dots, \frac{1}{\sqrt{\sigma_D^{s^2} + \epsilon}} \right]^\top, \quad \sigma_d^{s^2} = \frac{1}{N_s} \sum_{j=1}^{N_s} (\mathbf{k}_j^s(d) - \boldsymbol{\mu}^s(d))^2. \quad (47)$$

Following the same derivation as in Section 2.1, we derive the following support vector expansion of  $f^s$

$$f^s(\mathbf{x}) = \sum_{j=1}^{N_s} \frac{\Phi((\mathbf{x} - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})^\top \Phi((\mathbf{k}_j^s - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})}{h^s((\mathbf{x} - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})} \mathbf{v}_j^s + \mathbf{b}^s. \quad (48)$$

Here,  $\mathbf{v}_j^s = \left[ \frac{\alpha_j^s(1) - \tilde{\alpha}_j^s(1)}{h^s((\mathbf{k}_j^s - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})}, \dots, \frac{\alpha_j^s(D_v) - \tilde{\alpha}_j^s(D_v)}{h^s((\mathbf{k}_j^s - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})} \right]^\top$ , where  $\alpha_j^s$  and  $\tilde{\alpha}_j^s$  are the dual variables,  $j = 1, \dots, N$ . Same as in Section 2.1, in Eqn. 48, we choose  $\Phi$  as in Eqn. 28,  $h^s(\mathbf{x}) := \sum_{j=1}^{N_s} \Phi(\mathbf{x})^\top \Phi(\mathbf{k}_j^s)$ , and  $\mathbf{b}^s = 0$  to obtain the Batch Normalized Attention with Scaled Heads (Attention-BN+SH), which is defined as follows:

**Definition 5 (Batch Normalized Attention with Scaled Heads)** Given  $H$  sets of the key and value vectors  $\{\mathbf{k}_j^1, \mathbf{v}_j^1\}_{j=1}^{N_1}, \dots, \{\mathbf{k}_j^H, \mathbf{v}_j^H\}_{j=1}^{N_H}$ , for each set of  $H$  query vectors  $\mathbf{q}_1^1, \dots, \mathbf{q}_i^H$ ,  $i = 1, \dots, N$ , the Batch Normalized Attention with Scaled Heads (Attention-BN+SH) computes the corresponding output vector  $\mathbf{h}_i$  of the queries  $\mathbf{q}_1^1, \dots, \mathbf{q}_i^H$  by the following attention formula:

$$\mathbf{h}_i = \sum_{s=1}^H \mathbf{W}_O^s \left( \sum_{j=1}^{N_s} \text{softmax} \left( ((\mathbf{q}_i^s - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1})^\top ((\mathbf{k}_j^s - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s-1}) / \sqrt{D} \right) \mathbf{v}_j^s \right), \quad (49)$$

where

$$\boldsymbol{\mu}^s = \frac{1}{N_s} \sum_{j=1}^{N_s} \mathbf{k}_j^s, \quad \mathbf{s}^{s-1} = \left[ \frac{1}{\sqrt{\sigma_1^{s^2} + \epsilon}}, \dots, \frac{1}{\sqrt{\sigma_D^{s^2} + \epsilon}} \right]^\top, \quad \sigma_d^{s^2} = \frac{1}{N_s} \sum_{j=1}^{N_s} (\mathbf{k}_j^s(d) - \boldsymbol{\mu}^s(d))^2. \quad (50)$$

Following the same Remark 5 in Section 2.3.2, given input sequence  $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D_x}$  of  $N$  feature vectors in self-attention, in order to generate the sets of  $\{\mathbf{k}_j^s, \mathbf{v}_j^s\}_{j=1}^{N_s}$  at the scale  $s^{th}$ , we can downsample the input  $\mathbf{X}$  before projecting into the key matrix  $\mathbf{K}$  and the value matrix  $\mathbf{V}$ . In this paper, we use the average-pooling to downsample  $\mathbf{X}$ .

As in the same case of Attention-BN, for Attention-BN+SH, recentering queries and keys alone are sufficient for accuracy improvement, and we weight the mean  $\boldsymbol{\mu}$  in Eqn 49 with a constant  $\beta$ . Hence Eqn. 49 is simplified to:

$$\mathbf{h}_i = \sum_{s=1}^H \mathbf{W}_O^s \left( \sum_{j=1}^{N_s} \text{softmax} \left( (\mathbf{q}_i^s - \beta \boldsymbol{\mu}^s)^\top (\mathbf{k}_j^s - \beta \boldsymbol{\mu}^s) / \sqrt{D} \right) \mathbf{v}_j^s \right). \quad (51)$$

Table 12: The values of  $\beta$  for Linear Attention-BN/BN+SH and Sparse Attention-BN/BN+SH trained on the selected 10 UEA tasks.

Dataset/Model	Linear Attention-BN	Linear Attention-BN+SH	Sparse Attention-BN	Sparse Attention-BN+SH
ETHANOLCONCENTRATION	0.15	0.95	0.8	0.2
FACEDETECTION	0.6	0.6	0.6	0.6
HANDWRITING	0.25	0.3	0.3	0.3
HEARTBEAT	0.6	0.15	0.4	0.5
JAPANESEVOWELS	0.6	0.6	0.6	0.6
PEMS-SF	0.35	0.65	0.5	0.6
SELFREGULATIONSCP1	0.35	0.25	0.1	0.9
SELFREGULATIONSCP2	0.75	0.15	0.5	0.3
SPOKENARABICDIGITS	0.6	0.6	0.6	0.6
UWAVEGESTURELIBRARY	0.65	0.55	0.9	0.3

Table 13: The values of  $\beta$  for Attention-BN/BN+SH trained on 25 UEA Time Series classification tasks (Bagnall et al., 2018) and 6 UEA Time Series Regression tasks.

Dataset/Model	Attention-BN	Attention-BN+SH
ETHANOLCONCENTRATION	0.25	0.15
FACEDETECTION	0.6	0.6
HANDWRITING	0.65	0.25
HEARTBEAT	0.55	0.85
JAPANESEVOWELS	0.6	0.6
PEMS-SF	0.25	0.35
SELFREGULATIONSCP1	0.5	0.85
SELFREGULATIONSCP2	1.2	0.9
SPOKENARABICDIGITS	0.65	0.6
UWAVEGESTURELIBRARY	0.1	0.2
ARTICULARYWORDRECOGNITION	0.2	0.6
BASICMOTIONS	0.1	0.1
EPILEPSY	0.3	0.2
ERING	0.1	0.9
FINGERMOVEMENTS	1.0	0.3
LIBRAS	0.3	0.7
NATOPS	1.0	0.4
RACKETSPORTS	1.0	0.4
ATRIALFIBRILLATION	0.9	0.6
CRICKET	0.2	1.0
STANDWALKJUMP	1.0	0.6
HANDMOVEMENTDIRECTION	0.5	0.5
LSST	0.3	0.2
DUCKDUCKGEESE	0.6	0.3
MOTORIMAGERY	0.2	0.3
APPLIANCESENERGY	0.4	0.2
BENZENECONCENTRATION	0.2	0.1
BEIJINGPM10	0.1	0.5
BEIJINGPM25	0.1	0.2
LIVEFUELMOISTURE	0.1	0.3
IEEEPPG	0.4	0.2

Table 14: The values of  $\beta$  of Attention-BN/BN+SH trained on the 5 tasks of the LRA benchmark (Tay et al., 2021).

Dataset/Model	Attention-BN	Attention-BN+SH
LISTOPS	0.5	0.2
TEXT	0.5	0.8
RETRIEVAL	1.0	1.0
IMAGE	0.2	0.2
PATHFINDER	0.2	0.4

## J HYPERPARAMETERS

In this section, we provide the hyper-parameters for our best models.

### J.1 UEA TIME SERIES CLASSIFICATION AND REGRESSION

For these two benchmarks, use the set of downsampling factors  $\mathbf{s} = [1, 1, 2, 2, 4, 4, 8, 8]$  for Attention-SH/BN+SH and Linear/Sparse Attention-SH/BN+SH models trained on the UEA benchmark. Table 13 and Table 12 provide the values of  $\beta$  used for our best Attention-BN/BN+SH and Linear Attention-BN/BN+SH, Sparse Attention-BN/BN+SH models trained on subsets of the two benchmarks.

## J.2 LONG RANGE ARENA BENCHMARK

For all 5 tasks of the LRA benchmark, we set the downsampling factors  $\mathbf{s}$  of Attention-SH/BN+SH, Linear/Sparse Attention-SH/BN+SH is  $[1, 2]$  and kernel size of Attention-Conv1D models is 5. In addition, Table 14 provides the values  $\beta$  of Attention-BN/BN+SH models trained on the benchmark.

## J.3 IMAGENET CLASSIFICATION

This task's  $\beta$  of Attention-BN/BN+SH is 1. Attention-SH/BN+SH has the downsampling factor of  $[1, 1, 2, 4]$ , and the kernel size of Attention-Conv2D is  $(2, 2)$ .