# 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 $\boldsymbol{\mu}$ in Eqn 22 with a constant $\beta$. Hence Eqn 22 is simplified to:

$$\boldsymbol{h}_i = \sum_{j=1}^{N} \text{softmax}\left((\boldsymbol{q}_i - \beta\boldsymbol{\mu})^\top (\boldsymbol{k}_j - \beta\boldsymbol{\mu})/\sqrt{D}\right) \boldsymbol{v}_j. \tag{26}$$

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) [1] (for the PEMS-SF, SelfRegulationSCP2, UWaveGestureLibrary datasets) and (Zerveas et al., 2021) [2] (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)[3]. 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)[4]. The downsampling scales in Attention-SH/BN+SH models are $[1, 1, 2, 4]$ for 4 heads at each layer, respectively.

---

[1] Implementation available at https://github.com/thuml/Flowformer.
[2] Implementation available at https://github.com/gzerveas/mvts_transformer.
[3] Implementation available at https://github.com/NVIDIA/transformer-ls.
[4] 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 | $33.84 \pm 0.66$ | $\mathbf{34.98 \pm 0.74}$ | $34.76 \pm 0.69$ | $34.35 \pm 0.70$ |
| FaceDetection | $69.17 \pm 0.32$ | $69.22 \pm 0.17$ | $\mathbf{69.38 \pm 0.17}$ | $69.12 \pm 0.19$ |
| HandWriting | $32.87 \pm 0.27$ | $32.86 \pm 0.49$ | $32.82 \pm 0.12$ | $\mathbf{32.98 \pm 0.36}$ |
| HeartBeat | $75.61 \pm 0.73$ | $75.78 \pm 0.71$ | $74.96 \pm 0.62$ | $\mathbf{75.94 \pm 0.68}$ |
| JapaneseVowels | $99.37 \pm 0.16$ | $\mathbf{99.60 \pm 0.19}$ | $99.28 \pm 0.41$ | $99.33 \pm 0.19$ |
| PEMS-SF | $83.43 \pm 0.88$ | $85.74 \pm 0.67$ | $\mathbf{86.51 \pm 0.88}$ | $84.97 \pm 0.76$ |
| SelfRegulationSCP1 | $90.90 \pm 0.40$ | $91.81 \pm 0.69$ | $90.76 \pm 0.59$ | $\mathbf{91.92 \pm 0.60}$ |
| SelfRegulationSCP2 | $55.18 \pm 0.89$ | $\mathbf{56.11 \pm 0.94}$ | $54.44 \pm 0.88$ | $55.74 \pm 0.92$ |
| SpokenArabicDigits | $\mathbf{99.07 \pm 0.10}$ | $99.01 \pm 0.07$ | $99.03 \pm 0.18$ | $98.91 \pm 0.17$ |
| UWaveGestureLibrary | $85.63 \pm 0.81$ | $\mathbf{86.04 \pm 0.86}$ | $84.89 \pm 1.00$ | $85.78 \pm 0.75$ |
| Average Accuracy | $72.51 \pm 0.34$ | $\mathbf{73.12 \pm 0.20}$ | $72.68 \pm 0.40$ | $72.90 \pm 0.23$ |

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 (%) | $72.23 \pm 0.23$ | $\mathbf{73.18 \pm 0.24}$ |
| Top-5 Acc (%) | $91.13 \pm 0.15$ | $\mathbf{91.52 \pm 0.13}$ |

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 | $36.76 \pm 0.42$ | $\mathbf{37.20 \pm 0.05}$ |
| Text | $64.90 \pm 0.07$ | $\mathbf{64.92 \pm 0.44}$ |
| Retrieval | $79.68 \pm 0.52$ | $\mathbf{80.75 \pm 0.15}$ |
| Image | $\mathbf{39.23 \pm 1.35}$ | $39.18 \pm 0.59$ |
| Pathfinder | $72.72 \pm 0.75$ | $\mathbf{73.01 \pm 0.24}$ |
| Average Accuracy | $58.66 \pm 0.26$ | $\mathbf{59.01 \pm 0.20}$ |

# 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. Furthemore, 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 effiicient.

## 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 | $3.44 \pm 0.06$ | $3.38 \pm 0.34$ | $3.39 \pm 0.02$ | $\mathbf{3.37 \pm 0.23}$ |
| BENZENECONCENTRATION | $0.91 \pm 0.03$ | $\mathbf{0.89 \pm 0.17}$ | $1.00 \pm 0.09$ | $0.90 \pm 0.08$ |
| BEIJINGPM10 | $92.31 \pm 1.06$ | $\mathbf{92.00 \pm 0.89}$ | $92.82 \pm 0.92$ | $92.40 \pm 0.85$ |
| BEIJINGPM25 | $59.73 \pm 1.21$ | $59.55 \pm 0.92$ | $59.66 \pm 0.88$ | $\mathbf{59.24 \pm 1.22}$ |
| LIVEFUELMOISTURE | $43.08 \pm 0.17$ | $\mathbf{43.01 \pm 0.50}$ | $43.65 \pm 0.09$ | $43.79 \pm 0.49$ |
| IEEEPPG | $32.12 \pm 1.25$ | $\mathbf{30.69 \pm 0.64}$ | $31.38 \pm 1.02$ | $30.73 \pm 1.20$ |
| AVERAGE RMSE | $38.60 \pm 0.67$ | $\mathbf{38.25 \pm 0.30}$ | $38.65 \pm 0.27$ | $38.40 \pm 0.51$ |

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

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

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 $\beta$) | $\mathbf{81.31 \pm 0.25}$ |
| 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

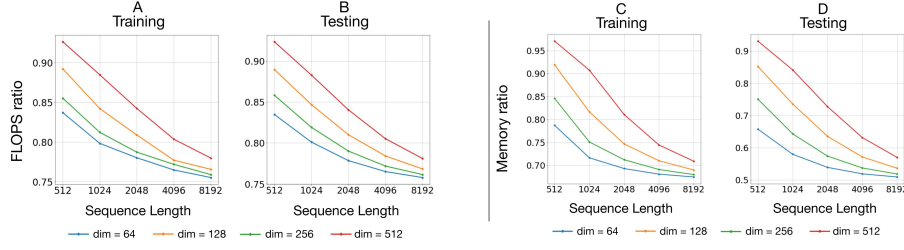This section provides more efficiency analysis on our models.

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.
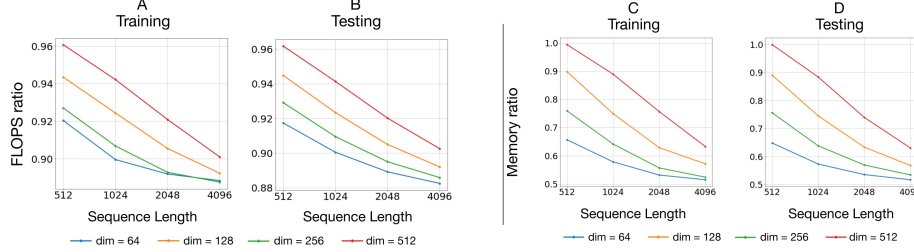


Figure 3: (Left) FLOPS ratios and (Right) memory usage ratios between the Sparse Attention-BN+SH and the Sparse Attention baseline trained on the LRA retrieval task for different model dimensions and sequence lengths. When using our models, the reduction in computation and memory improves with sequence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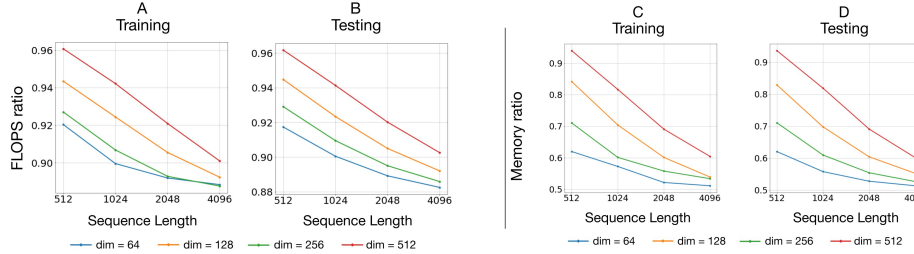


Figure 4: (Left) FLOPS ratios and (Right) memory usage ratios between the Sparse Attention-SH and the Sparse Attention baseline trained on the LRA retrieval task for different model dimensions and sequence lengths. When using our models, the reduction in computation and memory improves with sequence 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 $h(\boldsymbol{x})$ and $\Phi(\boldsymbol{x})$ allows us to derive the popular softmax attention given in Eqn. 1 and 2. In particular, if we choose $h(\boldsymbol{x}) := \sum_j^N \Phi(\boldsymbol{x})^T \Phi(\boldsymbol{k}_j)$, Eqn. 10 becomes

$$f(\boldsymbol{x}) = \sum_{j=1}^N \frac{\Phi(\boldsymbol{x})^\top \Phi(\boldsymbol{k}_j)}{\sum_{j'}^N \Phi(\boldsymbol{x})^T \Phi(\boldsymbol{k}_{j'})} \boldsymbol{v}_j + \boldsymbol{b} = \frac{\sum_{j=1}^N \Phi(\boldsymbol{x})^\top \Phi(\boldsymbol{k}_j) \boldsymbol{v}_j}{\sum_{j'}^N \Phi(\boldsymbol{x})^T \Phi(\boldsymbol{k}_{j'})} + \boldsymbol{b}. \qquad (27)$$

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

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

Since

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

then Eqn. 27 becomes

$$f(\boldsymbol{x}) = \sum_{j=1}^{N} \frac{\sum_{t=0}^{\infty} \sum_{n_1 + \cdots + n_D = t} \left( \frac{\left(\frac{x_1}{\sqrt[4]{D}}\right)^{n_1} \cdots \left(\frac{x_D}{\sqrt[4]{D}}\right)^{n_D}}{\sqrt{n_1! \ldots n_D!}} \right) \left( \frac{\left(\frac{k_{j1}}{\sqrt[4]{D}}\right)^{n_1} \cdots \left(\frac{k_{jD}}{\sqrt[4]{D}}\right)^{n_D}}{\sqrt{n_1! \ldots n_D!}} \right)}{\sum_{j'=1}^{N} \sum_{t=0}^{\infty} \sum_{n_1 + \cdots + n_D = t} \left( \frac{\left(\frac{x_1}{\sqrt[4]{D}}\right)^{n_1} \cdots \left(\frac{x_D}{\sqrt[4]{D}}\right)^{n_D}}{\sqrt{n_1! \ldots n_D!}} \right) \left( \frac{\left(\frac{k_{j'1}}{\sqrt[4]{D}}\right)^{n_1} \cdots \left(\frac{k_{j'D}}{\sqrt[4]{D}}\right)^{n_D}}{\sqrt{n_1! \ldots n_D!}} \right)} \boldsymbol{v}_j + \boldsymbol{b}$$

$$= \sum_{j=1}^{N} \frac{\exp\left( \left(\frac{\boldsymbol{x}}{\sqrt[4]{D}}\right)^{\top} \frac{\boldsymbol{k}_j}{\sqrt[4]{D}} \right)}{\sum_{j'=1}^{N} \exp\left( \left(\frac{\boldsymbol{x}}{\sqrt[4]{D}}\right)^{\top} \frac{\boldsymbol{k}_{j'}}{\sqrt[4]{D}} \right)} \boldsymbol{v}_j + \boldsymbol{b} = \sum_{j=1}^{N} \frac{\exp\left( \boldsymbol{x}^{\top} \boldsymbol{k}_j / \sqrt{D} \right)}{\sum_{j'=1}^{N} \exp\left( \boldsymbol{x}^{\top} \boldsymbol{k}_{j'} / \sqrt{D} \right)} \boldsymbol{v}_j + \boldsymbol{b}. \tag{30}$$

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

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

$$\boldsymbol{h}_i = \sum_{j=1}^{N} \mathrm{softmax}\left( \sum_{d=1}^{D} \frac{(\boldsymbol{q}_i(d) - \boldsymbol{\mu}(d))(\boldsymbol{k}_j(d) - \boldsymbol{\mu}(d))}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \boldsymbol{v}_j$$

$$= \sum_{j=1}^{N} \mathrm{softmax}\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_j(d) - \boldsymbol{q}_i(d)\boldsymbol{\mu}(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_j(d) + \boldsymbol{\mu}(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \boldsymbol{v}_j$$

$$= \sum_{j=1}^{N} \frac{\exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_j(d) - \boldsymbol{q}_i(d)\boldsymbol{\mu}(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_j(d) + \boldsymbol{\mu}(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right)}{\sum_{j'=1}^{N} \exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_{j'}(d) - \boldsymbol{q}_i(d)\boldsymbol{\mu}(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_{j'}(d) + \boldsymbol{\mu}(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right)} \boldsymbol{v}_j$$

$$= \sum_{j=1}^{N} \frac{\exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_j(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{\mu}(d)\boldsymbol{\mu}(d) - \boldsymbol{q}_i(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right)}{\sum_{j'=1}^{N} \exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_{j'}(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_{j'}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{\mu}(d)\boldsymbol{\mu}(d) - \boldsymbol{q}_i(d)\boldsymbol{\mu}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right)} \boldsymbol{v}_j$$

$$= \sum_{j=1}^{N} \frac{\exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_j(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right)}{\sum_{j'=1}^{N} \exp\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_{j'}(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_{j'}(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right)} \boldsymbol{v}_j$$

$$= \sum_{j=1}^{N} \mathrm{softmax}\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_j(d) - \boldsymbol{\mu}(d)\boldsymbol{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \boldsymbol{v}_j$$

$$= \sum_{j=1}^{N} \mathrm{softmax}\left( \sum_{d=1}^{D} \frac{\boldsymbol{q}_i(d)\boldsymbol{k}_j(d) - \frac{1}{N}\sum_{j'=1}^{N} \boldsymbol{k}_{j'}(d)\boldsymbol{k}_j(d)}{\sqrt{D}(\sigma_d^2 + \epsilon)} \right) \boldsymbol{v}_j. \tag{31}$$

## 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 $\{(\boldsymbol{x}_1, \boldsymbol{y}_1), \ldots, (\boldsymbol{x}_N, \boldsymbol{y}_N)\} \subset \mathcal{X} \times \mathcal{Y}$, where $\mathcal{X} = \mathbb{R}^{D_x}$ and $\mathcal{Y} = \mathbb{R}^{D_v}$. Here, $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$ are the training inputsd, and $\boldsymbol{y}_1, \ldots, \boldsymbol{y}_N$ 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

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

where $x \in \mathcal{X} = \mathbb{R}^{D_x}$, $W^{\text{proj}} = [w_1^{\text{proj}}, \ldots, w_D^{\text{proj}}]^\top \in \mathbb{R}^{D \times D_x}$, $\Phi(\cdot) = [\phi_1(\cdot), \ldots, \phi_{D_\phi}(\cdot)] : \mathbb{R}^D \to \mathbb{R}^{D_\phi}$, $W = [w_1, \ldots, w_{D_v}]^\top \in \mathbb{R}^{D_v \times D_\phi}$, $b \in \mathbb{R}^{D_v}$, and $h(x)$ is a vector-scalar function. We fit the function $f$ to the training data $\{(x_1, y_1), \ldots, (x_N, y_N)\}$ with an $\mathbb{L}_2$ regularization on $W$ and $W^{proj}$ by solving the following convex optimization problem:

$$\underset{\substack{W, W^{\text{proj}} \\ \xi_j, \tilde{\xi}_j, j=1,\ldots,N}}{\text{minimize}} \quad \frac{1}{2}\sum_{d=1}^{D_v}\|w_d\|^2 + \frac{1}{2}\sum_{d=1}^{D_v}\|w_d^{\text{proj}}\|^2 + C\sum_{j=1}^{N}\sum_{d=1}^{D_v}\left(\xi_j(d) + \tilde{\xi}_j(d)\right)$$

$$\text{subject to} \quad \begin{cases} y_j(d) - w_d^\top\Phi(W^{\text{proj}}x_j)/h(x_j) - x_j - b(d) \le \epsilon + \xi_j(d) \\ w_d^\top\Phi(W^{\text{proj}}x_j)/h(x_j) + x_j + b(d) - y_j(d) \le \epsilon + \tilde{\xi}_j(d) \\ \xi_j(d), \tilde{\xi}_j(d) \ge 0 \end{cases}, \; j=1,\ldots,N, \; d=1,\ldots,D_v.$$

$$\tag{33}$$

The Lagrangian of the optimization problem 33 is given by

$$\mathcal{L}_1 := \frac{1}{2}\sum_{d=1}^{D_v}\|w_d\|^2 + \frac{1}{2}\sum_{d=1}^{D_v}\|w_d^{\text{proj}}\|^2 + C\sum_{j=1}^{N}\sum_{d=1}^{D_v}\left(\xi_j(d) + \tilde{\xi}_j(d)\right) - \sum_{j=1}^{N}\sum_{d=1}^{D_v}\left(\eta_j(d)\xi_j(d) + \tilde{\eta}_j(d)\tilde{\xi}_j(d)\right)$$

$$- \sum_{j=1}^{N}\sum_{d=1}^{D_v}\alpha_j(d)\left(\epsilon + \xi_j(d) - y_j(d) + w_d^\top\frac{\Phi(W^{\text{proj}}x_j)}{h(x_j)} + x_j + b(d)\right)$$

$$- \sum_{j=1}^{N}\sum_{d=1}^{D_v}\tilde{\alpha}_j(d)\left(\epsilon + \tilde{\xi}_j(d) + y_j(d) - w_d^\top\frac{\Phi(W^{\text{proj}}x_j)}{h(x_j)} - x_j - b(d)\right),$$

$$\tag{34}$$

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

$$\partial_{w_d}\mathcal{L}_1 = w_d - \sum_{j=1}^{N}(\alpha_j(d) - \tilde{\alpha}_j(d))\frac{\Phi(W^{\text{proj}}x_j)}{h(x_j)} = 0 \Rightarrow w_d = \sum_{j=1}^{N}(\alpha_j(d) - \tilde{\alpha}_j(d))\frac{\Phi(W^{\text{proj}}x_j)}{h(x_j)}.$$

$$\tag{35}$$

Note that here we only find the form of the optimal solution for $W = [w_1, \ldots, w_{D_v}]^\top$. The optimal value of $W^{\text{proj}}$ can then be found by optimization algorithm such as the (stochastic) gradient descent when training the transformer.

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

$$f(x) = \left[\sum_{j=1}^{N}(\alpha_j(1) - \tilde{\alpha}_j(1))\frac{\Phi(W^{\text{proj}}x_j)}{h(x_j)}, \ldots, \sum_{j=1}^{N}(\alpha_j(D_v) - \tilde{\alpha}_j(D_v))\frac{\Phi(W^{\text{proj}}x_j)}{h(x_j)}\right]^\top \frac{\Phi(W^{\text{proj}}x)}{h(x)} + x + b,$$

$$= \left[\sum_{j=1}^{N}\frac{\alpha_j(1) - \tilde{\alpha}_j(1)}{h(x_j)}\frac{\Phi(W^{\text{proj}}x)^\top\Phi(W^{\text{proj}}x_j)}{h(x)}, \ldots, \sum_{j=1}^{N}\frac{\alpha_j(D_v) - \tilde{\alpha}_j(D_v)}{h(x_j)}\frac{\Phi(W^{\text{proj}}x)^\top\Phi(W^{\text{proj}}x_j)}{h(x)}\right]^\top + x + b,$$

$$= \underbrace{\sum_{j=1}^{N}\frac{\Phi(W^{\text{proj}}x)^\top\Phi(W^{\text{proj}}x_j)}{h(x)}v_j + x}_{\text{Residual connection}} + b. \tag{36}$$

Here, the support vector expansion of $f$ already includes a residual connection. The softmax attention can then be derived by selecting $h(x) := \sum_{j}^{N}\Phi(W^{\text{proj}}x)^T\Phi(W^{\text{proj}}x_j)$ and choosing $\Phi$ as in Eqn. 28 in Section 2.1. Note that in Eqn. 36, $\{x_j\}_{j=1}^{N}$ and $x$ 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 $W^{\text{proj}}x_j = W_Kx_j$, $W^{\text{proj}}x = W_Qx$, $v_j = W_Vx_j$ and choosing the test sample $x$ among the training samples $\{x_j\}_{j=1}^{N}$.

**Remark 6** *Here, for self-attention, we choose the test sample $\boldsymbol{x}$ among the training samples $\{\boldsymbol{x}_j\}_{j=1}^{N}$ 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 $\{(\boldsymbol{x}_1^{train}, \boldsymbol{y}_1^{train}), \dots, (\boldsymbol{x}_{N_H \times N_W}^{train}, \boldsymbol{y}_{N_H \times N_W}^{train})\} \subset \mathcal{X} \times \mathcal{Y}$, where $\mathcal{X} = \mathbb{R}^{D_x}$ and $\mathcal{Y} = \mathbb{R}^{D_v}$. Here, $\boldsymbol{x}_1^{train}, \dots, \boldsymbol{x}_{N_H \times N_W}^{train}$ are the training inputs, and $\boldsymbol{y}_1^{train}, \dots, \boldsymbol{y}_{N_H \times N_W}^{train}$ 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 $\mathbf{X}^{train}(h, w, d) = \boldsymbol{x}_{N_W \times (h-1)+w}^{train}(d)$. Given a new set of inputs $\{\boldsymbol{x}_1, \dots, \boldsymbol{x}_{N_H \times N_W}\} \subset \mathcal{X}$ and the corresponding 3D-tensor $\mathbf{X} \in \mathbb{R}^{N_H \times N_W \times D_x}$ of these inputs, where $\mathbf{X}(h, w, d) = \boldsymbol{x}_{N_W \times (h-1)+w}(d)$. We consider the function $f$ applying on the 3D-tensor $\mathbf{X}$ and taking the following form

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

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

$$\underset{\substack{\mathbf{W} \\ \boldsymbol{\xi}_j, \tilde{\boldsymbol{\xi}}_j, j=1,\dots,N_H \times N_W}}{\text{minimize}} \quad \frac{1}{2}\|\mathbf{W}\|_{\text{F}}^2 + C \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} \left(\boldsymbol{\xi}_j(d) + \tilde{\boldsymbol{\xi}}_j(d)\right) = \frac{1}{2} \sum_{d=1}^{D_v} \|\boldsymbol{w}_d\|^2 + C \sum_{j=1}^{N_H \times N_W} \sum_{d=1}^{D_v} \left(\boldsymbol{\xi}_j(d) + \tilde{\boldsymbol{\xi}}_j(d)\right)$$

subject to

$$\begin{cases} \boldsymbol{y}_j^{train}(d) - \boldsymbol{w}_d^\top \dfrac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\boldsymbol{x}_j^{train})} - \boldsymbol{b}(d) \leq \epsilon + \boldsymbol{\xi}_j(d) \\ \boldsymbol{w}_d^\top \dfrac{\Phi(\text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j))}{h(\boldsymbol{x}_j^{train})} + \boldsymbol{b}(d) - \boldsymbol{y}_j^{train}(d) \leq \epsilon + \tilde{\boldsymbol{\xi}}_j(d) \\ \boldsymbol{\xi}_j(d), \tilde{\boldsymbol{\xi}}_j(d) \geq 0, \; j = 1, \dots, N_H \times N_W, \; d = 1, \dots, D_v. \end{cases}$$
$$\tag{38}$$

The Lagrangian of the optimization problem 38 is given by

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

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

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

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

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

$$f(\boldsymbol{x}_i) = \left[ \sum_{j=1}^{N_H \times N_W} \frac{\boldsymbol{\alpha}_j(1) - \tilde{\boldsymbol{\alpha}}_j(1)}{h(\boldsymbol{x}_j^{train})} \frac{\mathbf{A}_{ij}}{h(\boldsymbol{x}_i)}, \ldots, \sum_{j=1}^{N_H \times N_W} \frac{\boldsymbol{\alpha}_j(D_v) - \tilde{\boldsymbol{\alpha}}_j(D_v)}{h(\boldsymbol{x}_j^{train})} \frac{\mathbf{A}_{ij}}{h(\boldsymbol{x}_i)} \right]^\top + \boldsymbol{b},$$

$$= \sum_{j=1}^{N_H \times N_W} \frac{\mathbf{A}_{ij}}{h(\boldsymbol{x}_i)} \boldsymbol{v}_j + \boldsymbol{b}, \tag{42}$$

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 $\boldsymbol{b}_s = 0$. To derive the softmax normalization in attention, we choose $h(\boldsymbol{x}_i) := \sum_{j=1}^{N} \mathbf{A}_{ij}$ and select $\Phi$ as in Eqn. 28. Let the training inputs $\{\boldsymbol{x}_1^{train}, \ldots, \boldsymbol{x}_{N_H \times N_W}^{train}\} \subset \mathcal{X}$ be the attention keys $\{\boldsymbol{k}_1, \ldots, \boldsymbol{k}_{N_H \times N_W}\} \subset \mathcal{K}$, where $\mathcal{K} = \mathbb{R}^D$, in self-attention. Also, let the new inputs $\{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_{N_H \times N_W}\} \subset \mathcal{X}$ be the attention queries $\{\boldsymbol{q}_1, \ldots, \boldsymbol{q}_{N_H \times N_W}\} \subset \mathcal{K}$ 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 $\{\boldsymbol{k}_j, \boldsymbol{v}_j\}_{j=1}^{N_H \times N_W}$, and a set of the query vectors $\{\boldsymbol{q}_i\}_{i=1}^{N_H \times N_W}$. 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 $\mathbf{K}(h, w, d) = \boldsymbol{k}_{N_W \times (h-1)+w}(d)$ and $\mathbf{Q}(h, w, d) = \boldsymbol{q}_{N_W \times (h-1)+w}(d)$. The 2D-Convolution Attention (Attention-Conv2D) computes the corresponding output vector $\boldsymbol{h}_i$ of the query $\boldsymbol{q}_i$ by the following attention formula:*

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

*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 $i, j = 1, \ldots, N$, we let*

$$\boldsymbol{k}_j = \text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j)) = \Phi(\text{Flatten}(\text{Conv2D}(\text{Reshape2D}(\mathbf{X}^{input}), s, \mathbf{W}_K))(j),$$

$$\boldsymbol{q}_i = \text{Flatten}(\text{Conv2D}(\mathbf{X}, s))(i)) = \Phi(\text{Flatten}(\text{Conv2D}(\text{Reshape2D}(\mathbf{X}^{input}), s, \mathbf{W}_Q))(i),$$

$$\boldsymbol{v}_j = \text{Flatten}(\text{Conv2D}(\mathbf{X}^{train}, s))(j)) = \Phi(\text{Flatten}(\text{Conv2D}(\text{Reshape2D}(\mathbf{X}^{input}), s, \mathbf{W}_V))(j).$$

$$\tag{44}$$

*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 $\{\boldsymbol{k}_j, \boldsymbol{v}_j\}_{j=1}^{N}$, and a set of the query vectors $\{\boldsymbol{q}_i\}_{i=1}^{N}$. Denote the key matrix and the query matrix by $\mathbf{K} :=$*

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

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

where the $Conv1D(\cdot, s)$ 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 $f^s$, $s = 1, \ldots, H$, in Eqn. 17 with training sets $\{(\boldsymbol{k}_1^1, \boldsymbol{y}_1^1), \ldots, (\boldsymbol{k}_{N_1}^1, \boldsymbol{y}_{N_1}^1)\}, \ldots, \{(\boldsymbol{k}_1^H, \boldsymbol{y}_1^H), \ldots, (\boldsymbol{k}_{N_H}^H, \boldsymbol{y}_{N_H}^H)\} \subset \mathcal{K} \times \mathcal{Y}$ of different sizes $N_1, \ldots, N_H$, where $\mathcal{K} = \mathbb{R}^D$ and $\mathcal{Y} = \mathbb{R}^{D_v}$. The function $f^s$ is defined as:

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

where

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

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

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

Here, $\boldsymbol{v}_j^s = \left[ \frac{\boldsymbol{\alpha}_j^s(1) - \tilde{\boldsymbol{\alpha}}_j^s(1)}{h^s((\boldsymbol{k}_j^s - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s^{-1}})}, \ldots, \frac{\boldsymbol{\alpha}_j^s(D_v) - \tilde{\boldsymbol{\alpha}}_j^s(D_v)}{h^s((\boldsymbol{k}_j^s - \boldsymbol{\mu}^s) \odot \mathbf{s}^{s^{-1}})} \right]^\top$, where $\boldsymbol{\alpha}_j^s$ and $\tilde{\boldsymbol{\alpha}}_j^s$ are the dual variables, $j = 1, \ldots, N$. Same as in Section 2.1, in Eqn. 48, we choose $\Phi$ as in Eqn. 28, $h^s(\boldsymbol{x}) := \sum_j^{N_s} \Phi(\boldsymbol{x})^T \Phi(\boldsymbol{k}_j^s)$, and $\boldsymbol{b}^s = 0$ 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 $\{\boldsymbol{k}_j^1, \boldsymbol{v}_j^1\}_{j=1}^{N_1}, \ldots, \{\boldsymbol{k}_j^H, \boldsymbol{v}_j^H\}_{j=1}^{N_H}$, for each set of $H$ query vectors $\boldsymbol{q}_i^1, \ldots, \boldsymbol{q}_i^H$, $i = 1, \ldots, N$, the Batch Normalized Attention with Scaled Heads (Attention-BN+SH) computes the corresponding output vector $\boldsymbol{h}_i$ of the queries $\boldsymbol{q}_i^1, \ldots, \boldsymbol{q}_i^H$ by the following attention formula:*

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

where

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

Following the same Remark 5 in Section 2.3.2, given input sequence $\mathbf{X} := [\boldsymbol{x}_1, \cdots, \boldsymbol{x}_N]^\top \in \mathbb{R}^{N \times D_x}$ of $N$ feature vectors in self-attention, in order to generate the sets of $\{\boldsymbol{k}_j^s, \boldsymbol{v}_j^s\}_{j=1}^{N_s}$ 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:

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

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 **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 $(2, 2)$.