

TOWARDS BETTER MULTI-HEAD ATTENTION VIA CHANNEL-WISE SAMPLE PERMUTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformer plays a central role in many fundamental deep learning models, e.g., the ViT in computer vision and the BERT and GPT in natural language processing, whose effectiveness is mainly attributed to its multi-head attention (MHA) mechanism. In this study, we propose a simple and novel channel-wise sample permutation (CSP) operator, achieving a new structured MHA with fewer parameters and lower complexity. Given an input matrix, CSP circularly shifts the samples of different channels with various steps and then sorts grouped samples of each channel. This operator is equivalent to implicitly implementing cross-channel attention maps as permutation matrices, which achieves linear complexity and suppresses the risk of rank collapse when representing data. We replace the MHA of some representative models with CSP and test the CSP-based models in several discriminative tasks, including image classification and long sequence analysis. Experiments show that the CSP-based models achieve comparable or better performance with fewer parameters and lower computational costs than the classic Transformer and its state-of-the-art variants. The code is available at <https://anonymous.open.science/r/CSP-BA52>.

1 INTRODUCTION

Transformer (Vaswani et al., 2017) has been widely adopted in the deep learning domain. Recent large language models like GPT (Brown et al., 2020; Radford et al.) and LLaMA (Touvron et al., 2023a;b) series are built based on the Transformer and its variants, which demonstrate their remarkable abilities in natural language processing. In the field of computer vision, Vision Transformers (ViTs) (Dosovitskiy et al., 2021), such as EfficientViT (Cai et al., 2023; Liu et al., 2023) and SHViT (Yun & Ro, 2024), exhibit exceptional performance and consistently push their limits. In addition, the Transformer-based models have been designed for the complex structured data in various applications, including the Informer (Zhou et al., 2021) for time series broadcasting, the Transformer Hawkes process (Zuo et al., 2020) for continuous-time event sequence prediction, the Graphormer (Ying et al., 2021) for molecular representation, the Mesh Transformer (Lin et al., 2021) for 3D mesh representation, the Set-Transformer (Lee et al., 2019) and Point-Transformer (Zhao et al., 2021) for point cloud modeling, and so on. Although some new alternatives like Mamba (Gu & Dao, 2023) and RWKV (Peng et al., 2023) have been proposed and shown their competitiveness in some aspects, Transformer still maintains a dominant position when developing deep learning models because of its strong performance and outstanding universality.

The effectiveness of Transformer is mainly attributed to its multi-head attention (MHA) mechanism (Vaswani et al., 2017). However, MHA’s quadratic complexity concerning sequence length leads to a heavy, even unaffordable, computational overhead when modeling long sequences. To improve the efficiency of MHA, many variants of Transformer introduce sparse or low-rank structures into attention maps (Child et al., 2019; Kitaev et al., 2020; Wang et al., 2020; Ma et al., 2021; Wang et al., 2024) and apply algorithms friendly to GPU acceleration (Dao et al., 2022; Dao, 2024). At the same time, many attempts have been made to explore the mathematical reasons for the power of MHA, e.g., analyzing the representation power and rank collapse risk of MHA (Dong et al., 2021; Ying et al., 2021) and revisiting attention maps through the lens of kernel theory (Tsai et al., 2019; Qin et al., 2022) and optimal transport (Tay et al., 2020; Sander et al., 2022). Currently, the above two research directions seem “parallel” in most situations: The acceleration methods of MHA are

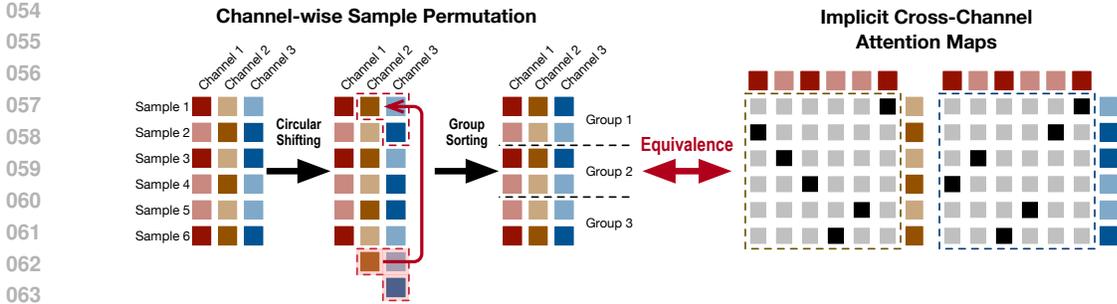


Figure 1: An illustration of the proposed channel-wise sample permutation operator and the equivalent implicit cross-channel attention maps.

often empirical, but the theoretical work mainly analyzes the classic MHA, making it seldom support the rationality of the accelerated MHAs or contribute to developing a new MHA.

In this study, we propose a novel **Channel-wise Sample Permutation (CSP)** operator, which leads to a new multi-head attention mechanism that is solid in theory and efficient in practice. As illustrated in Figure 1, given an input matrix, CSP first shifts the samples of different channels circularly with various steps and then sorts grouped samples of each channel. This operator is equivalent to implicitly implementing cross-channel attention maps as permutation matrices, which introduce inter- and intra-group interactions for the samples across different channels. CSP is much simpler than the classic MHA and its existing variants. It has no learnable parameters and can achieve linear computational complexity regarding sequence length.

The proposed CSP operator is motivated by the recent development of MHA. In particular, the work in (Child et al., 2019; Beltagy et al., 2020; Kitaev et al., 2020; Sander et al., 2022) empirically demonstrate the rationality of pursuing attention maps with sparse doubly stochastic structures, which is further verified by an analytic experiment in this study. CSP achieves permutation-based implicit attention maps that satisfy these structural properties, and thus, it has a good chance of providing a better MHA mechanism. Moreover, such attention maps have all-one spectrums because of their permutation nature. Based on the theoretical analysis framework provided in (Dong et al., 2021), we prove that replacing MHA with CSP can suppress the risk of rank collapse when representing data. In addition, we provide insightful understandings of the CSP operator by explaining its circular shifting and group sorting steps from the perspectives of optimal transport-based attention layer (Sander et al., 2022) and channel-wise mixer (Yu et al., 2022; Lian et al., 2022), respectively.

To demonstrate the usefulness of CSP, we replace the MHA of some state-of-the-art models with CSP and compare the CSP-based models with the original MHA-based ones in representative discriminative tasks, including long sequence analysis and image classification. For each model, replacing its MHA with CSP significantly reduces the number of parameters and the computational cost while maintaining or even improving model performance.

2 PRELIMINARIES AND RELATED WORK

Typically, given an input $\mathbf{X} \in \mathbb{R}^{N \times C}$, where N indicates the length of a sequence or the size of a sample set and C is the number of channels (feature dimensions), an attention head (Vaswani et al., 2017) first obtains the value, query, and key matrices by linear maps, i.e., $\mathbf{V} = \mathbf{X}\mathbf{W}_V \in \mathbb{R}^{N \times D}$, $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q \in \mathbb{R}^{N \times D}$, and $\mathbf{K} = \mathbf{X}\mathbf{W}_K \in \mathbb{R}^{N \times D}$, and then projects \mathbf{V} as follows:

$$\text{Att}(\mathbf{V}; \mathbf{Q}, \mathbf{K}) := \mathbf{P}(\mathbf{Q}, \mathbf{K})\mathbf{V}, \text{ where } \mathbf{P}(\mathbf{Q}, \mathbf{K}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right). \quad (1)$$

Here, we denote \mathbf{V} as the input matrix of the head and $\mathbf{P}(\mathbf{Q}, \mathbf{K}) \in \mathbb{R}^{N \times N}$ as the attention map parametrized by \mathbf{Q} and \mathbf{K} , respectively. The multi-head attention mechanism applies a group of linear maps, i.e., $\theta = \{\mathbf{W}_{V,m}, \mathbf{W}_{Q,m}, \mathbf{W}_{K,m} \in \mathbb{R}^{C \times D}\}_{m=1}^M$, to construct M attention heads and concatenates their outputs, i.e.,

$$\text{MHA}_\theta(\mathbf{X}) := \parallel_{m=1}^M \text{Att}(\mathbf{V}_m; \mathbf{Q}_m, \mathbf{K}_m) \in \mathbb{R}^{N \times MD}, \quad (2)$$

where $V_m = \mathbf{X}\mathbf{W}_{V,m}$, $Q_m = \mathbf{X}\mathbf{W}_{Q,m}$, and $K_m = \mathbf{X}\mathbf{W}_{K,m}$ for $m = 1, \dots, M$, and “ \parallel ” denotes the concatenation operation. In practice, we set $MD = C$ for applying skip connections in the Transformer architecture, i.e., $\text{MHA}_\theta(\mathbf{X}) + \mathbf{X}$.

The attention map in (1) has quadratic computational complexity concerning the sequence length N because of its “query-key-value” (abbreviated, QKV) architecture. Considering the high complexity per attention head, the MHA has to restrict the number of attention heads to achieve a trade-off between model capacity and computational efficiency, which may limit its representation power.

Many efforts have been made to improve the classic MHA. SparseTrans (Child et al., 2019) and Longformer (Beltagy et al., 2020) compute local attention maps based on the subsequences extracted by sliding windows, which leads to sparse global attention maps. To use shorter subsequences while retaining more information, S³ Attention (Wang et al., 2024) integrates global and local information by leveraging Fourier Transformation and a convolutional kernel. Some other models sparsify the key and query matrices directly by locality-sensitive hashing (LSH) (Kitaev et al., 2020) or ReLU (Qin et al., 2022). Besides pursuing sparse attention maps, Performer (Choromanski et al., 2021) and Linformer (Wang et al., 2020) apply low-rank attention maps. Recently, FlashAttention and its variants (Dao et al., 2022; Dao, 2024) further accelerate the computation of attention maps for long sequences by sophisticated I/O design, parallelism, and work partitioning. In addition to simplifying the computation of the attention maps, some work provides new understandings of the attention mechanism. The work in (Tsai et al., 2019; Choromanski et al., 2021; Qin et al., 2022) implements attention maps as various kernel matrices. The work in (Sander et al., 2022) implements doubly stochastic attention maps by the Sinkhorn-Knopp algorithm (Sinkhorn & Knopp, 1967) and explains the computation of each attention map as a discretized Wasserstein gradient flow.

Currently, the above accelerated or structured MHAs often lead to the performance degradation, while the theoretical understandings of MHA seldom help improve its computational efficiency in practice. Our work attempts to bridge the gap, proposing a theoretically solid multi-head attention mechanism with low complexity and competitive performance.

3 PROPOSED METHOD

3.1 MOTIVATION: PURSUING SPARSE DOUBLY STOCHASTIC ATTENTION MAPS

As shown in Section 2, many models apply various strategies to construct *sparse* attention maps, e.g., the locality-sensitive hashing (LSH) in (Kitaev et al., 2020), the subsequence sampling in (Child et al., 2019; Beltagy et al., 2020), and the sparse activation in (Qin et al., 2022). These models achieve encouraging performance and higher efficiency than the vanilla Transformer, demonstrating sparse attention maps’ rationality. Besides making attention maps sparse, the work in (Sander et al., 2022) shows that in various discriminative tasks, the attention maps tend to be *doubly stochastic* automatically (i.e., $\mathbf{P} \in \Pi_N$, where $\Pi_N = \{\mathbf{A} \geq \mathbf{0} \mid \mathbf{A}\mathbf{1}_N = \mathbf{1}_N, \mathbf{A}^\top \mathbf{1}_N = \mathbf{1}_N\}$) during training,¹ and the Transformer applying doubly stochastic attention (called Sinkformer) outperforms the vanilla Transformer in image and text classification.

The above recent models show that sparse attention maps help improve the models’ computational efficiency (thus making increasing attention heads feasible), and doubly stochastic attention maps help improve the models’ discriminative power. These phenomena imply that **designing sparse doubly stochastic attention maps may lead to a better MHA mechanism and further boost model performance**. To verify this claim, we conduct an analytic experiment, replacing the attention maps in a standard ViT (Dosovitskiy et al., 2021) with simple permutation matrices (the doubly stochastic matrices with the strongest sparsity) and evaluating the model performance on the CIFAR-10 dataset (Krizhevsky, 2009). In particular, the ViT used in this experiment consists of six Transformer layers. Each Transformer has eight attention heads (i.e., $M = 8$), and each head sets $N = 64$, $C = 512$, and $D = 64$. For each layer, we replace the attention map of the m -th head with the following permutation matrix:

$$\mathbf{S}_{C(m-1)/D} = \begin{bmatrix} \mathbf{0} & \mathbf{I}_{C(m-1)/D} \\ \mathbf{I}_{N-C(m-1)/D} & \mathbf{0} \end{bmatrix}, \text{ for } m = 1, \dots, M, \quad (3)$$

¹Please refer to Section 3 in (Sander et al., 2022) for more details.

Table 1: A comparison for various MHAs and their classification accuracy (%) on CIFAR-10.

MHA	#Heads per layer	Parameters per layer	Top-1 Acc.	Top-5 Acc.
$\prod_{m=1}^M \mathbf{P}(\mathbf{Q}_m, \mathbf{K}_m) \mathbf{V}_m$	8 (= M)	$\{\mathbf{W}_{Q,m}, \mathbf{W}_{K,m}, \mathbf{W}_{V,m}\}_{m=1}^M$	81.90	98.85
$\prod_{m=1}^M \mathbf{S}_{C(m-1)/D} \mathbf{V}_m$	8 (= M)	$\{\mathbf{W}_{V,m}\}_{m=1}^M$	80.70	98.97
$\prod_{c=1}^C \mathbf{S}_{(c-1) \bmod N} \mathbf{v}_c$	64 (= N)	$\{\mathbf{W}_{V,m}\}_{m=1}^M$	83.84	99.27

where \mathbf{I}_N indicates an identity matrix with a size $N \times N$. Obviously, the permutation matrix $\mathbf{S}_{C(m-1)/D}$ corresponds to a circular shifting operator — $\mathbf{S}_{C(m-1)/D} \mathbf{V}_m$ means shifting the rows of \mathbf{V}_m circularly with $C(m-1)/D$ steps. Furthermore, for each layer, we can concatenate $\{\mathbf{V}_m\}_{m=1}^M$ to get $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_C] \in \mathbb{R}^{N \times C}$ and circularly shift the channels of this matrix by applying $\mathbf{S}_{(c-1) \bmod N} \mathbf{v}_c$ for $c = 1, \dots, C$, where “mod” is the modulo operation. In this case, the number of attention heads, equal to the number of distinguishable permutation matrices, becomes N . As shown in Table 1, even if the sparse doubly stochastic attention maps we designed are extremely simple and have no parameters, applying them with a sufficient number can still result in competitive, even better performance. This experimental result motivates us to construct sufficiently many sparse doubly stochastic attention maps with low complexity, leading to the proposed channel-wise sample permutation operator.

3.2 CHANNEL-WISE SAMPLE PERMUTATION FOR IMPLICIT CROSS-CHANNEL ATTENTION

As shown in Figure 1, given an input matrix $\mathbf{X} \in \mathbb{R}^{N \times C}$, the CSP operator first projects \mathbf{X} to a value matrix with the same size, i.e., $\mathbf{V} = \mathbf{X}\mathbf{W} = [\mathbf{v}_1, \dots, \mathbf{v}_C] \in \mathbb{R}^{N \times C}$, where $\mathbf{W} \in \mathbb{R}^{C \times C}$ and \mathbf{v}_c denotes the N samples in the c -th channel. Given \mathbf{V} , the CSP operator shifts the samples of different channels circularly with various steps and then sorts grouped samples of each channel, i.e.,

$$\text{CSP}_{\mathbf{W}}(\mathbf{X}) := \prod_{c=1}^C \text{GSort}_K(\mathbf{S}_{J_c} \mathbf{v}_c) = \prod_{c=1}^C \mathbf{P}_c \mathbf{v}_c, \quad \text{where } \text{GSort}_K(\mathbf{v}) = \prod_{k=1}^K \text{Sort}(\mathbf{v}^{(k)}). \quad (4)$$

Here, \mathbf{S}_{J_c} is the circular shifting operator defined in (3). $\text{GSort}_K(\mathbf{v})$ denotes grouping the elements of a vector \mathbf{v} into K parts, i.e., $\mathbf{v} = [\mathbf{v}^{(1)}; \dots; \mathbf{v}^{(K)}]$, and sorting each part accordingly. When implementing the CSP operator, we take the first channel \mathbf{v}_1 as the reference in this study: The circular shifting of each \mathbf{v}_c is with respect to \mathbf{v}_1 , and the group sorting permutes the elements of $(\mathbf{S}_{J_c} \mathbf{v}_c)^{(k)}$ according to the element-wise order of $\mathbf{v}_1^{(k)}$, for $c = 2, \dots, C$ and $k = 1, \dots, K$.

The CSP operator is equivalent to implicitly implementing sparse doubly stochastic attention maps as permutation matrices, which builds interactions for the samples across different channels. As shown in (4), we denote each attention map as \mathbf{P}_c . For \mathbf{v}_1 , $\mathbf{P}_1 = \mathbf{I}_N$. For the remaining \mathbf{v}_c , \mathbf{P}_c can be decomposed into the following two parts:

$$\mathbf{P}_c = \mathbf{T}_c \mathbf{S}_{J_c} = \text{BlkDiag}(\{\mathbf{T}_c^{(k)}\}_{k=1}^K) \mathbf{S}_{J_c}, \quad \text{for } c = 2, \dots, C, \quad (5)$$

where \mathbf{T}_c is a block-diagonal permutation matrix determined by the group sorting operation. The k -th block $\mathbf{T}_c^{(k)}$ is a permutation matrix determined by the sorting within the k -th group, which introduces intra-group sample interactions across different channels. The circular shifting operation introduces inter-group sampler interactions across different channels, and the ranges of the interactions are determined by the predefined shifting steps. As a result, for arbitrary two \mathbf{v}_c and $\mathbf{v}_{c'}$, $\mathbf{P}_c \mathbf{v}_c$ and $\mathbf{P}_{c'} \mathbf{v}_{c'}$ captures their interactions determined by $\mathbf{P}_c^\top \mathbf{P}_{c'}$.

3.2.1 ADVANTAGES OVER EXISTING MHAS

High computational efficiency: Replacing MHA with CSP leads to a new variant of Transformer. Table 2 compares the proposed model with the existing MHA-based models. We can find that the computational complexity of CSP can be $\mathcal{O}(N \log \frac{N}{K})$ when applying QuickSort (Hoare, 1962) to implement the group sorting operation, which is much lower than the computational complexity of the existing MHAs. When the group size is 2, we can achieve group sorting by the simple “min-max” operation (Anil et al., 2019; Tanielian & Biau, 2021), and the computational complexity further reduces to $\mathcal{O}(N)$. In addition, as shown in (4), except for the projection matrix \mathbf{W} corresponding to the value matrix, CSP does not require additional projection matrices to construct the query and key matrices. In other words, its parameters are only one-third of the classic MHA.

Table 2: A comparison for existing MHA mechanisms and CSP.

Model	Attention($\mathbf{V}; \mathbf{Q}, \mathbf{K}$)	Complexity	Attention Structure
Transformer	$\text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(CN^2)$	Row-normalized
SparseTrans	$\text{Local2D-Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(CN^{1.5})$	Sparse+Row-normalized
Longformer	$\text{Local1D-Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(CNE)$	Sparse+Row-normalized
Reformer	$\text{LSH-Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}}\right)\mathbf{V}$	$\mathcal{O}(CN \log N)$	Sparse+Row-normalized
CosFormer	$(\mathbf{Q}_{\cos}\mathbf{K}_{\cos}^\top + \mathbf{Q}_{\sin}\mathbf{K}_{\sin}^\top)\mathbf{V}$	$\mathcal{O}(\min\{CE_{QK}, NE_Q\})$	Sparse
MEGA	$f\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{D}} + \mathbf{B}\right)\mathbf{V}$	$\mathcal{O}(CN^2) \sim \mathcal{O}(CNr)$	(Optional) Sparse+Row-normalized
Performer	$\phi_r(\mathbf{Q})\phi_r(\mathbf{K})^\top\mathbf{V}$	$\mathcal{O}(CNr)$	Low-rank
Linformer	$\text{Softmax}\left(\frac{\mathbf{Q}\psi_r(\mathbf{K})^\top}{\sqrt{D}}\right)\psi_r(\mathbf{V})$	$\mathcal{O}(CNr)$	Low-rank+Row-normalized
Proposed	$\ _{c=1}^C \mathbf{P}_c \mathbf{v}_c$	$\mathcal{O}(CN \log \frac{N}{K}) \sim \mathcal{O}(CN)$	Sparse+Doubly stochastic

¹ “Local1D” considers subsequences with length E when computing attention maps. “Local2D” considers the row-wise and column-wise local data for a sequence zigzagging in the 2D space.

² $\phi_r: \mathbb{R}^D \mapsto \mathbb{R}^r$, and $\phi_r(\mathbf{Q}), \phi_r(\mathbf{K}) \in \mathbb{R}^{N \times r}$; $\psi_r: \mathbb{R}^N \mapsto \mathbb{R}^r$, and $\psi_r(\mathbf{K}), \psi_r(\mathbf{V}) \in \mathbb{R}^{r \times D}$.

³ $\mathbf{K}_{\cos} = \text{diag}(\{\cos \frac{\pi i}{2M}\}_{i=1}^N) \text{ReLU}(\mathbf{K})$, $\mathbf{K}_{\sin} = \text{diag}(\{\sin \frac{\pi i}{2M}\}_{i=1}^N) \text{ReLU}(\mathbf{K})$. So are \mathbf{Q}_{\cos} and \mathbf{Q}_{\sin} . E_{QK} and E_Q are the numbers of nonzero elements in $\mathbf{Q}_{\cos}\mathbf{K}_{\cos}^\top$ and \mathbf{Q}_{\cos} , respectively.

⁴ For MEGA, $\mathbf{B} \in \mathbb{R}^{N \times N}$ is a bias matrix. f denotes the Softmax function in NLP tasks and a Laplace function in computer vision tasks. Its complexity becomes $\mathcal{O}(CNr)$ when applying a chunk mechanism to derive sparse attention maps.

A low risk of rank collapse: Besides significantly improving computational efficiency, CSP can suppress an ordinary risk of the classic MHA, rank collapse. In particular, we define the rank-1 estimation residual of a matrix \mathbf{X} associated with an arbitrary matrix norm as

$$\epsilon(\mathbf{X}) = \mathbf{X} - \mathbf{1}\hat{\mathbf{x}}^\top, \text{ where } \hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{X} - \mathbf{1}\mathbf{x}^\top\|. \quad (6)$$

In addition, for a matrix $\mathbf{X} = [x_{nc}] \in \mathbb{R}^{N \times C}$, we can define its $(1, \infty)$ -norm as $\|\mathbf{X}\|_{1, \infty} = \sqrt{\|\mathbf{X}\|_1 \|\mathbf{X}\|_\infty}$, where $\|\mathbf{X}\|_1 = \max_c \sum_{n=1}^N |x_{nc}|$ and $\|\mathbf{X}\|_\infty = \max_n \sum_{c=1}^C |x_{nc}|$, respectively. It has been known that $\|\epsilon(\mathbf{X})\|_{1, \infty}$ measures the rank collapse of \mathbf{X} effectively, i.e., $\|\epsilon(\mathbf{X})\|_{1, \infty} \rightarrow 0$ means that \mathbf{X} collapses to a rank-1 matrix. The work in (Dong et al., 2021) shows that if we construct a Transformer by stacking MHA layers without skip connections, its output matrix will lose its rank doubly exponentially with depth, i.e., $\|(\text{MHA}_L \circ \dots \circ \text{MHA}_1(\mathbf{X}))\|_{1, \infty} = \mathcal{O}(\|\epsilon(\mathbf{X})\|_{1, \infty}^{3L})$, where L is the number of the MHA layers.

Applying CSP can suppress this risk, which is supported by the following theorem.

Theorem 1 Suppose that we construct a layer- L network as $(f \circ \text{CSP})^L = (f_{\lambda_L} \circ \text{CSP}_{\mathbf{W}^{(L)}}) \circ \dots \circ (f_{\lambda_1} \circ \text{CSP}_{\mathbf{W}^{(1)}})$. For $\ell = 1, \dots, L$, $\text{CSP}_{\mathbf{W}^{(\ell)}}$ is a C -channel CSP operator, and $f_{\lambda_\ell}: \mathbb{R}^C \mapsto \mathbb{R}^C$ is a λ_ℓ -Lipschitz function. Denote $\beta = \max_\ell \|\mathbf{W}^{(\ell)}\|_1$ and $\lambda = \max_\ell \lambda_\ell$. Then, we have

$$\|(\text{CSP})^L(\mathbf{X})\|_{1, \infty} \leq C^{\frac{L}{2}} (\lambda\beta)^L \|\epsilon(\mathbf{X})\|_{1, \infty}, \forall \mathbf{X} \in \mathbb{R}^{N \times C}. \quad (7)$$

Theorem 1 indicates a linear convergence rate of the residual. It means that the model applying CSP avoids the rapid decay of the matrix rank. A detailed proof is shown in Appendix A.

3.2.2 IMPLEMENTATION DETAILS

Circular shifting: The shifting step is crucial for the circular shifting operation, which determines the range of sample interaction. When the sequence length is comparable to the number of channels in each layer, i.e., $N \approx C$, we can simply set the shifting step size $J_c = c \lceil \frac{N}{C} \rceil$ for $c = 1, \dots, C$, so that the circular shifting operation can generate sufficient distinguishable attention heads with respect to the sequence length. However, for long sequences, i.e., $N \gg C$, we need to set the shifting steps of different channels with high dynamics, making C attention heads build diverse

interactions in a long sequence. In this study, given a Transformer-based model with L layers, we consider all L value matrices in these layers jointly, and set LC different shifting steps based on power law, as illustrated in Figure 2. In particular, we denote J as the base shifting step. For $c = 1, \dots, LC$, we shift the c -th channel circularly with $J^{c-1} - 1$ steps. In addition, for the last channel, we require $J^{LC-1} - 1 \approx N - 1$. Therefore, we can set $J = \lfloor N^{1/(LC-1)} \rfloor$.

Group sorting: Instead of merely applying the circular shifting operation (as in Section 3.1), we introduce the group sorting operation to CSP, which helps increase the number of attention heads. Given an input matrix with size $N \times C$, the circular shifting operation constructs $\min\{N, C\}$ different attention maps, which results in repeated attention maps when $C > N$. For the channels applying the same circular shifting steps, the group sorting operation can make their attention maps different from each other as long as the orders of their samples are inconsistent. As a result, the group sorting helps CSP increase the number of attention heads from $\min\{C, N\}$ to C .

A special case of CSP: Note that when setting $K = 1$, the group sorting becomes the classic complete sorting, leading to a special case of CSP. Given C channels, the complete sorting can directly generate at most C distinguishable permutation matrices/attention heads. In addition, because of using the complete sorting, the circular shifting step of CSP becomes redundant. In the following experiments, empirically, implementing CSP as the complete sorting often works well when modeling long sequences while the CSP combining circular shifting with group sorting helps represent visual objects.

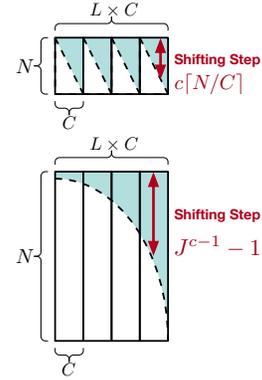


Figure 2: The shifting strategies when $N \approx C$ and $N \gg C$.

3.3 FUNCTIONALITY AND RATIONALITY ANALYSIS OF CSP

Circular shifting works as a channel-wise mixer: The circular shifting of CSP is similar to the channel-wise mixers used in visual representation models. In particular, the convolution neural networks like ShuffleNet (Zhang et al., 2018) and its variant (Ma et al., 2018) apply grouped convolution operation to reduce computational costs and increase inter-group interactions by shuffling the channels across different groups. This shuffling strategy inspires many lightweight channel-wise mixers, e.g., the hierarchical rearrangement in Hira-MLP (Guo et al., 2022), the spatial-shift module in S^2 -MLP (Yu et al., 2022), and the axial-shift module in AS-MLP (Lian et al., 2022). For example, given a visual feature tensor with a size $H \times W \times C$ (i.e., 2D images with C channels), the axial-shift module applies horizontal and vertical shifts with zero padding to the 2D images of different channels. The spatial-shift module first divides the input tensor into four parts by grouping its channels and then shifts the four sub-tensors along four different directions. Both these two modules apply small shifting steps to achieve local shifting. The circular shifting of CSP corresponds to applying these shifting modules to 1D sequences. To capture the short-range and long-range interactions between the samples of different channels simultaneously, we apply various shifting steps to different channels and replace zero padding with circular padding.

Group sorting works as an optimal transport-based MHA: Without causing any ambiguity, we denote $(S_{J_c} v_c)^{(k)}$ as $v_c^{(k)}$ for simplification. It is easy to prove that the $T_c^{(k)}$ in (5) is the optimal transport (OT) between $v_1^{(k)}$ and $v_c^{(k)}$, which can be derived by $\min_{T \in \Pi_{N/K}} \langle -v_1^{(k)} v_c^{(k)\top}, T \rangle$.² From this viewpoint, CSP achieves a new OT-based MHA mechanism. In addition, when approximating $T_c^{(k)}$ as an entropic optimal transport, we can connect CSP to the doubly stochastic attention mechanism used in Sinkformer (Sander et al., 2022). In particular, each attention head in Sinkformer derives a doubly stochastic attention map, denoted as $T_{T,\tau}$, by the Sinkhorn-Knopp algorithm (Sinkhorn & Knopp, 1967), i.e.,

$$T_{t,\tau} = \text{Sinkhorn}_t \left(\exp \left(\frac{QK^\top}{\tau\sqrt{D}} \right) \right), \text{ and } T_{\infty,\tau} = \arg \min_{T \in \Pi_N} \langle -QK^\top, T \rangle + \tau\sqrt{D}H(T). \quad (8)$$

Here, $\text{Sinkhorn}_t(\mathbf{A})$ means normalizing the rows and columns of a nonnegative matrix \mathbf{A} alternatively by t times, i.e., $\mathbf{A}^{(0)} = \mathbf{A}$, and $\mathbf{A}^{(i)} = N_c \circ N_r(\mathbf{A}^{(i-1)})$ for $i = 1, \dots, t$, where N_c and N_r

²See Appendix B for a detailed derivation.

denote column-wise and row-wise normalization, respectively. As shown in (8), the attention map corresponds to the optimal solution of an entropic optimal transport problem (Cuturi, 2013) when $t \rightarrow \infty$, where $\langle \cdot, \cdot \rangle$ denotes the inner product operation, $H(\mathbf{T}) = \langle \mathbf{T}, \log \mathbf{T} \rangle$ denotes the entropy of \mathbf{T} , and its significance is controlled by $\tau > 0$.

We can connect Sinkformer to CSP by modifying the attention mechanism in (8) as follows. Given a value matrix $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_C] \in \mathbb{R}^{N \times C}$, we replace the \mathbf{Q} and \mathbf{K} in (8) with \mathbf{v}_1 and \mathbf{v}_c , respectively, for $c = 1, \dots, C$, and divide each vector into K groups. We can achieve a C -head K -group doubly stochastic attention mechanism by applying the Sinkhorn-Knopp algorithm to $\mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}$, for $c = 1, \dots, C$ and $k = 1, \dots, K$, i.e.,

$$\mathbf{T}_{c,t,\tau} = \text{BlkDiag}\left(\left\{\text{Sinkhorn}_t\left(\exp\left(\frac{1}{\tau}\mathbf{v}_1^{(k)}\mathbf{v}_c^{(k)\top}\right)\right)\right\}_{k=1}^K\right) = \text{BlkDiag}(\{\mathbf{T}_{c,t,\tau}^{(k)}\}_{k=1}^K), \quad (9)$$

where $\mathbf{T}_{c,t,\tau}$ is a doubly stochastic attention map with a block-diagonal structure, and $\mathbf{T}_{c,t,\tau}^{(k)}$ denotes a local attention map, which corresponds to computing the entropic optimal transport between $\mathbf{v}_1^{(k)}$ and $\mathbf{v}_c^{(k)}$ when $t \rightarrow \infty$, i.e., $\mathbf{T}_{c,\infty,\tau}^{(k)} = \arg \min_{\mathbf{T} \in \Pi_{N/K}} \langle -\mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}, \mathbf{T} \rangle + \tau H(\mathbf{T})$.

The connection between the attention map in (9) and CSP is captured by the following theorem.

Theorem 2 *If $\min_{\mathbf{T} \in \Pi_{N/K}} \langle -\mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}, \mathbf{T} \rangle$ admits a unique optimal solution, for $c = 1, \dots, C$ and $k = 1, \dots, K$, then for the $\mathbf{T}_{c,t,\tau}$ in (9), $\lim_{\tau \rightarrow 0} \mathbf{T}_{c,\infty,\tau} \mathbf{S}_{J_c}$ converges to the \mathbf{P}_c in (5) weakly.*

Theorem 2 can be derived directly based on the weak convergence of entropic optimal transport (Theorem 5.10 in (Nutz, 2022)). This theorem indicates that a Sinkformer can implement CSP approximately if it *i*) sets the query and key matrices as the channels of the value matrix and *ii*) applies the Sinkhorn-Knopp algorithm to the grouped samples.

4 EXPERIMENTS

To demonstrate the effectiveness and efficiency of CSP, we conduct comprehensive comparative and analytic experiments in two representative discriminative tasks, image classification and long sequence analysis. The implementation details can be found in Appendix C.

4.1 IMAGE CLASSIFICATION

We conduct comparative experiments and ablation studies on three image datasets, including CIFAR-10, CIFAR-100 (Krizhevsky, 2009), and ImageNet-1k (Deng et al., 2009). For each dataset, we treat the classic ViT as the baseline and replace its MHA layers with *i*) **circular shifting**, *ii*) **group sorting**, and *iii*) the proposed **CSP** operator, respectively. Here, the circular shifting and the group sorting are two simplified CSP operators that help analyze the contributions of different CSP modules. Table 3 shows these models' size and classification accuracy. Applying CSP and its simplified variants can reduce the model size significantly without the query and key matrices. The circular shifting operator achieves competitive performance in all three datasets. In addition, although the standalone group sorting operator results in performance degradation, combining it with the circular shifting operator, i.e., the proposed CSP, can achieve the best performance. These observations are consistent with the experimental results achieved by mixer-MLP models (Yu et al., 2022; Lian et al., 2022): *i*) Simple channel-wise interactions can replace the dense and smoothed attention maps and lead to promising model performance, and *ii*) the shifting operator is crucial for CSP in computer vision tasks because it fully leverages the local similarity nature of the image. Moreover, when we increase the number of channels per layer and make the model size comparable to the original ViT, we can further boost the performance of the CSP-based models and achieve the best performance.

In Figure 3, we illustrate the singular spectrums of the output matrices achieved by different methods on ImageNet-1k. The spectrums achieved by the circular shifting and CSP operators decay much more slowly than the spectrum achieved by MHA. This observed phenomenon serves as a strong validation of the theoretical result in Theorem 1, providing further evidence that the representation model using permutation-based attention maps indeed carries a lower risk of rank collapse compared to the classic MHA-based model.

Table 3: The comparison for various models on the number of parameters ($\times 10^6$) and classification accuracy (%). The best result on each dataset is **bold**, and the second best result is underlined.

Model	Attention	CIFAR-10			CIFAR-100			ImageNet-1k		
		#Param.	Top-1	Top-5	#Param.	Top-1	Top-5	#Param.	Top-1	Top-5
ViT	MHA	9.52	81.90	98.85	9.65	53.30	79.97	22.05	76.53	92.81
	Circular Shifting	6.38	83.84	99.27	6.50	58.38	84.26	18.50	75.64	92.42
	Group Sorting	6.38	79.41	99.03	6.50	51.47	79.67	18.50	64.77	85.28
	CSP (Proposed)	6.38	84.81	99.35	6.50	59.16	84.76	18.50	76.66	93.05
		9.52	85.02	99.37	9.65	59.23	85.09	22.05	77.14	93.23

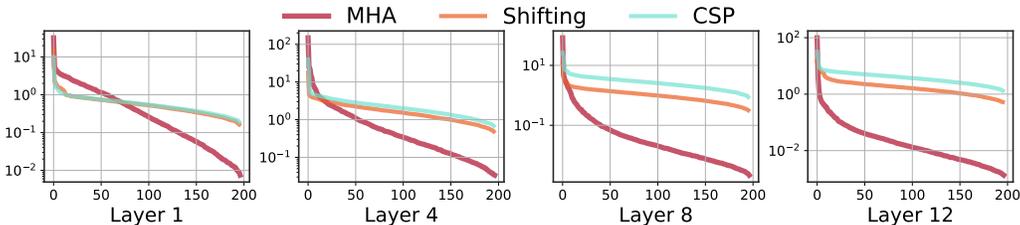


Figure 3: The singular spectrums of the output matrices achieved on ImageNet-1k.

4.2 LONG RANGE ARENA BENCHMARK

Long Range Arena (LRA) is a benchmark designed to evaluate models for long sequence analysis (Tay et al., 2021b), which consists of six discriminative tasks, including ListOps (Nangia & Bowman, 2018), byte-level text classification (Maas et al., 2011), byte-level document retrieval (Radev et al., 2013), and three sequentialized image classification tasks, i.e., CIFAR-10 (Krizhevsky, 2009), PathFind, and Path-X (Linsley et al., 2018).³ Each image is formulated as a long sequence of pixels in the three image classification tasks. We first replace the MHA of the classic Transformer with CSP and compare it with other variants of Transformer. As shown in Figure 4 and the first part of Table 4, the Transformer using CSP outperforms other models on both performance and computational efficiency. It achieves the highest average score and the fastest training speed among all the models, and its memory cost is comparable to the most efficient variant of Transformer. For long sequence modeling, we simply implement CSP as the complete sorting operator in this experiment, which can capture the long-range interactions between the samples with the highest flexibility.

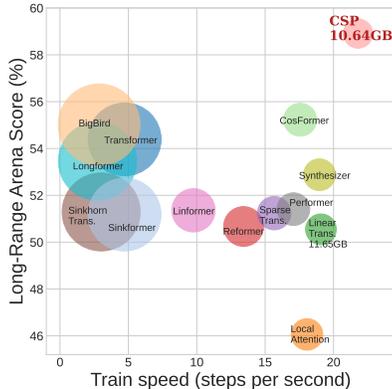


Figure 4: The performance and efficiency of various models on the LRA benchmark. The disk area indicates the memory cost of each method.

Besides improving the classic Transformer, we further plug CSP into the state-of-the-art attention-based model, MEGA (Ma et al., 2022), and analyze its impacts on the model performance. As evidenced in the second part of Table 4, the MEGA with dense attention maps currently outperforms all other methods, including those based on the state space model (SSM), such as S5 (Smith et al., 2022) and SPADE (Zuo et al., 2024), on the LRA benchmark. When MEGA applies chunked attention maps, its performance degrades slightly but its computational complexity can reduce from $\mathcal{O}(CN^2)$ to $\mathcal{O}(CNr)$, where r is the chunk size. When replacing the attention mechanism of MEGA with the proposed CSP operator, the complexity of

³Given a set of gray-level images, each of which plots two points and several curves, PathFind aims to recognize whether there exists a path connecting the points in each image. Path-X is a more challenging version of Pathfind because of applying high-resolution images.

Table 4: Results (%) of various methods on the LRA benchmark. The first group contains the classic Transformer and its variants, and the second group contains the state-of-the-art methods on LRA. The best result on each dataset is **bold**, and the second best result is underlined.

Type	Model	ListOps	Text	Retrieval	Image	PathFind	Path-X	Avg.	
MHA	Transformer (Vaswani et al., 2017)	36.37	64.27	57.46	42.44	71.40	FAIL	54.39	
	LocalAttention (Tay et al., 2021b)	15.82	52.98	53.39	41.46	66.63	FAIL	46.06	
	LinearTrans (Katharopoulos et al., 2020)	16.13	65.90	53.09	42.34	75.30	FAIL	50.55	
	Reformer (Kitaev et al., 2020)	37.27	56.10	53.40	38.07	68.50	FAIL	50.67	
	Sinkformer (Sander et al., 2022)	30.70	64.03	55.45	41.08	64.65	FAIL	51.18	
	SparseTrans (Child et al., 2019)	17.07	63.58	59.59	<u>44.24</u>	71.71	FAIL	51.24	
	SinkhornTrans (Tay et al., 2020)	33.67	61.20	53.83	41.23	67.45	FAIL	51.29	
	Linformer (Wang et al., 2020)	35.70	53.94	52.27	38.56	76.34	FAIL	51.36	
	Performer (Choromanski et al., 2021)	18.01	<u>65.40</u>	53.82	42.77	<u>77.05</u>	FAIL	51.41	
	Synthesizer (Tay et al., 2021a)	36.99	61.68	54.67	41.61	69.45	FAIL	52.88	
	Longformer (Beltagy et al., 2020)	35.63	62.85	56.89	42.22	69.71	FAIL	53.46	
	BigBird (Zaheer et al., 2020)	36.05	64.02	59.29	40.83	74.87	FAIL	55.01	
	Cosformer (Qin et al., 2022)	37.90	63.41	61.36	43.17	70.33	FAIL	55.23	
	Transformer using CSP (Proposed)	<u>37.65</u>	64.60	62.23	48.02	82.04	FAIL	58.91	
Type	Model	Complexity	ListOps	Text	Retrieval	Image	PathFind	Path-X	Avg.
CNN	CCNN (Romero et al., 2022)	$\mathcal{O}(CN^2)$	43.60	84.08	FAIL	88.90	91.51	FAIL	68.02
SSM	ETSMLP (Chu & Lin, 2024)	$\mathcal{O}(CN)$	<u>62.55</u>	88.49	86.72	75.34	91.66	93.78	83.09
	S4 (Gu et al., 2022)		58.35	76.02	87.09	87.26	86.05	88.10	80.48
	S5 (Smith et al., 2022)		62.15	89.31	91.40	88.00	95.33	98.58	87.46
	SPADE (Zuo et al., 2024)		59.70	87.55	90.13	<u>89.11</u>	96.42	94.22	86.19
MHA	MEGA (Ma et al., 2022)	$\mathcal{O}(CN^2)$	63.14	90.43	<u>91.25</u>	90.44	96.01	97.98	88.21
	MEGA-chunk (Ma et al., 2022)	$\mathcal{O}(CNr)$	58.76	90.19	<u>90.97</u>	85.80	94.41	93.81	85.66
	MEGA using CSP (Proposed)	$\mathcal{O}(CN)$	61.85	<u>90.27</u>	90.09	87.42	93.74	91.98	85.89

the CSP-based MEGA becomes linear and thus comparable to that of the chunked MEGA and the SSM-based models. At the same time, the CSP-based MEGA is better than the chunked MEGA in the overall performance. These results serve as compelling evidence, demonstrating the practical rationality of CSP.

5 CONCLUSION & FUTURE WORK

We have proposed a novel channel-wise sample permutation operator, leading to a simple but effective surrogate of existing multi-head attention mechanisms. In theory, we demonstrate that the proposed CSP operator overcomes the rank collapse problem of the classic MHA because of implementing sparse doubly stochastic attention maps as permutation matrices. In addition, we explain the operator from the perspective of channel-wise mixer and optimal transport-based attention. For representative MHA-based models, replacing their MHA layers with the CSP operator helps improve their performance in discriminative tasks and reduce their computational cost at the same time. In summary, our work provides a promising solution to developing a better multi-head attention mechanism, demonstrating the usefulness of discrete algorithms like shifting and sorting in model design.

Limitations and Future Work. Currently, the design of CSP is motivated by pursuing sparse doubly stochastic attention maps, which restricts its application to discriminative tasks — the attention maps of Transformer decoder in generative tasks are lower-triangular, so that imposing the doubly stochastic constraint on the attention maps results in trivial identity matrices. For the Transformers in generative tasks (Radford et al.; Touvron et al., 2023a;b), how to achieve effective and efficient attention maps by simple algorithms is still an open problem, which is left as our future work. In addition, we implement our method based on Pytorch at the current stage. To maximize the computational efficiency of our method, we plan to refactor its underlying code and optimize its I/O, parallelism, and partitioning strategies as FlashAttention (Dao et al., 2022; Dao, 2024) did.

REFERENCES

- 486
487
488 Cem Anil, James Lucas, and Roger Grosse. Sorting out lipschitz function approximation. In *Inter-*
489 *national Conference on Machine Learning*, pp. 291–301. PMLR, 2019.
- 490 Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
491 *arXiv preprint arXiv:2004.05150*, 2020.
- 492 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
493 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, et al. Jax:
494 composable transformations of python+ numpy programs. 2018.
- 495 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
496 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
497 few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- 498 Han Cai, Junyan Li, Muyan Hu, Chuang Gan, and Song Han. Efficientvit: Lightweight multi-scale
499 attention for high-resolution dense prediction. In *Proceedings of the IEEE/CVF International*
500 *Conference on Computer Vision*, pp. 17302–17313, 2023.
- 501 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse
502 transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 503 Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea
504 Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser,
505 David Benjamin Belanger, Lucy J Colwell, and Adrian Weller. Rethinking attention with per-
506 formers. In *International Conference on Learning Representations*, 2021.
- 507 Jiqun Chu and Zuoquan Lin. Incorporating exponential smoothing into mlp: a simple but effective
508 sequence model. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp.
509 326–337, 2024.
- 510 Marco Cuturi. Sinkhorn distances: lightspeed computation of optimal transport. In *Proceedings*
511 *of the 26th International Conference on Neural Information Processing Systems-Volume 2*, pp.
512 2292–2300, 2013.
- 513 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The*
514 *Twelfth International Conference on Learning Representations*, 2024.
- 515 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-
516 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,
517 35:16344–16359, 2022.
- 518 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
519 hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*,
520 pp. 248–255. IEEE, 2009.
- 521 Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure
522 attention loses rank doubly exponentially with depth. In *International Conference on Machine*
523 *Learning*, pp. 2793–2803. PMLR, 2021.
- 524 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
525 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-
526 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at
527 scale. In *International Conference on Learning Representations*, 2021.
- 528 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
529 *preprint arXiv:2312.00752*, 2023.
- 530 Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured
531 state spaces. In *International Conference on Learning Representations*, 2022.
- 532 Jianyuan Guo, Yehui Tang, Kai Han, Xinghao Chen, Han Wu, Chao Xu, Chang Xu, and Yunhe
533 Wang. Hire-mlp: Vision mlp via hierarchical rearrangement. In *Proceedings of the IEEE/CVF*
534 *conference on computer vision and pattern recognition*, pp. 826–836, 2022.

- 540 Charles AR Hoare. Quicksort. *The computer journal*, 5(1):10–16, 1962.
- 541
- 542 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are
543 rnns: Fast autoregressive transformers with linear attention. In *International Conference on Ma-*
544 *chine Learning*, pp. 5156–5165. PMLR, 2020.
- 545 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In
546 *International Conference on Learning Representations*, 2020.
- 547
- 548 Alex Krizhevsky. Learning multiple layers of features from tiny images. *Master’s thesis, University*
549 *of Toronto*, 2009.
- 550
- 551 Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set
552 transformer: A framework for attention-based permutation-invariant neural networks. In *Interna-*
553 *tional Conference on Machine Learning*, pp. 3744–3753. PMLR, 2019.
- 554 Dongze Lian, Zehao Yu, Xing Sun, and Shenghua Gao. As-mlp: An axial shifted mlp architecture
555 for vision. In *International Conference on Learning Representations*, 2022.
- 556
- 557 Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. In *Proceedings of The IEEE/CVF*
558 *International Conference on Computer Vision*, pp. 12939–12948, 2021.
- 559 Drew Linsley, Junkyung Kim, Vijay Veerabadrán, Charles Windolf, and Thomas Serre. Learning
560 long-range spatial dependencies with horizontal gated recurrent units. *Advances in Neural Infor-*
561 *mation Processing Systems*, 31, 2018.
- 562
- 563 Xinyu Liu, Houwen Peng, Ningxin Zheng, Yuqing Yang, Han Hu, and Yixuan Yuan. Efficientvit:
564 Memory efficient vision transformer with cascaded group attention. In *Proceedings of the*
565 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14420–14430, 2023.
- 566 Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for
567 efficient cnn architecture design. In *Proceedings of the European conference on computer vision*
568 *(ECCV)*, pp. 116–131, 2018.
- 569
- 570 Xuezhe Ma, Xiang Kong, Sinong Wang, Chunting Zhou, Jonathan May, Hao Ma, and Luke Zettle-
571 moyer. Luna: Linear unified nested attention. *Advances in Neural Information Processing Sys-*
572 *tems*, 34:2441–2453, 2021.
- 573 Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan
574 May, and Luke Zettlemoyer. Mega: moving average equipped gated attention. In *International*
575 *Conference on Learning Representations*, 2022.
- 576
- 577 Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts.
578 Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the*
579 *association for computational linguistics: Human language technologies*, pp. 142–150, 2011.
- 580 Nikita Nangia and Samuel Bowman. Listops: A diagnostic dataset for latent tree learning. In
581 *Proceedings of the 2018 Conference of the North American Chapter of the Association for Com-*
582 *putational Linguistics: Student Research Workshop*, pp. 92–99, 2018.
- 583
- 584 Marcel Nutz. *Introduction to entropic optimal transport*. Columbia University, 2022.
- 585
- 586 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman,
587 Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, et al. Rwkv: Reinventing rnns for
588 the transformer era. In *Findings of the Association for Computational Linguistics: EMNLP 2023*,
pp. 14048–14077, 2023.
- 589
- 590 Zhen Qin, Weixuan Sun, Hui Deng, Dongxu Li, Yunshen Wei, Baohong Lv, Junjie Yan, Lingpeng
591 Kong, and Yiran Zhong. cosformer: Rethinking softmax in attention. In *International Conference*
592 *on Learning Representations*, 2022.
- 593
- 593 Dragomir R Radev, Pradeep Muthukrishnan, Vahed Qazvinian, and Amjad Abu-Jbara. The acl
anthology network corpus. *Language Resources and Evaluation*, 47:919–944, 2013.

- 594 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language under-
595 standing by generative pre-training.
596
- 597 David W Romero, David M Knigge, Albert Gu, Erik J Bekkers, Efstratios Gavves, Jakub M Tom-
598 czak, and Mark Hoogendoorn. Towards a general purpose cnn for long range dependencies in n
599 d. 2022.
- 600 Michael E Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Sinkformers: Transform-
601 ers with doubly stochastic attention. In *International Conference on Artificial Intelligence and*
602 *Statistics*, pp. 3515–3530. PMLR, 2022.
- 603
- 604 Richard Sinkhorn and Paul Knopp. Concerning nonnegative matrices and doubly stochastic matri-
605 ces. *Pacific Journal of Mathematics*, 21(2):343–348, 1967.
- 606
- 607 Jimmy TH Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for
608 sequence modeling. In *International Conference on Learning Representations*, 2022.
- 609
- 610 Ugo Tanielian and Gerard Biau. Approximating lipschitz continuous functions with groupsort neu-
611 ral networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 442–450.
PMLR, 2021.
- 612
- 613 Yi Tay, Dara Bahri, Liu Yang, Donald Metzler, and Da-Cheng Juan. Sparse sinkhorn attention. In
614 *International Conference on Machine Learning*, pp. 9438–9447. PMLR, 2020.
- 615
- 616 Yi Tay, Dara Bahri, Donald Metzler, Da-Cheng Juan, Zhe Zhao, and Che Zheng. Synthesizer: Re-
617 thinking self-attention for transformer models. In *International Conference on Machine Learning*,
pp. 10183–10192. PMLR, 2021a.
- 618
- 619 Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao,
620 Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient
621 transformers. In *International Conference on Learning Representations*, 2021b.
- 622
- 623 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
624 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 625
- 626 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
627 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
628 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 629
- 630 Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan
631 Salakhutdinov. Transformer dissection: An unified understanding for transformer’s attention via
632 the lens of kernel. In *Proceedings of the 2019 Conference on Empirical Methods in Natural*
Language Processing, pp. 4344–4353, 2019.
- 633
- 634 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
635 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
tion processing systems, 30, 2017.
- 636
- 637 Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention
638 with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- 639
- 640 Xue Wang, Tian Zhou, Jianqing Zhu, Jialin Liu, Kun Yuan, Tao Yao, Wotao Yin, Rong Jin, and
641 HanQin Cai. S 3 attention: Improving long sequence attention with smoothed skeleton sketching.
IEEE Journal of Selected Topics in Signal Processing, 2024.
- 642
- 643 Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and
644 Tie-Yan Liu. Do transformers really perform badly for graph representation? *Advances in Neural*
Information Processing Systems, 34:28877–28888, 2021.
- 645
- 646 Tan Yu, Xu Li, Yunfeng Cai, Mingming Sun, and Ping Li. S2-mlp: Spatial-shift mlp architecture for
647 vision. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*,
pp. 297–306, 2022.

648 Seokju Yun and Youngmin Ro. Shvit: Single-head vision transformer with memory efficient macro
649 design. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
650 pp. 5756–5767, 2024.

651
652 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago
653 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for
654 longer sequences. *Advances in Neural Information Processing Systems*, 33:17283–17297, 2020.

655 Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient
656 convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on*
657 *computer vision and pattern recognition*, pp. 6848–6856, 2018.

658
659 Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In
660 *Proceedings of The IEEE/CVF International Conference on Computer Vision*, pp. 16259–16268,
661 2021.

662 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
663 Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings of*
664 *the AAAI Conference on Artificial Intelligence*, 35(12):11106–11115, 2021.

665 Simiao Zuo, Haoming Jiang, Zichong Li, Tuo Zhao, and Hongyuan Zha. Transformer hawkes
666 process. In *International Conference on Machine Learning*, pp. 11692–11702. PMLR, 2020.

667
668 Simiao Zuo, Xiaodong Liu, Jian Jiao, Denis X Charles, Eren Manavoglu, Tuo Zhao, and Jianfeng
669 Gao. Efficient hybrid long sequence modeling with state space augmented transformers. In *First*
670 *Conference on Language Modeling*, 2024.

671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 A THE PROOF OF THEOREM 1

703
704 **A Single Channel:** Given a CSP's output matrix, we can derive a residue for each channel as

$$705 \epsilon(\mathbf{P}_c \mathbf{X} \mathbf{w}_c) = \mathbf{P}_c \mathbf{X} \mathbf{w}_c - \mathbf{1} \hat{a}_c, \text{ for } c = 1, \dots, C. \quad (10)$$

706 where $\hat{\mathbf{a}} = [\hat{a}_c]$ and

$$707 \hat{a}_c = \arg \min_a \|\mathbf{P}_c \mathbf{X} \mathbf{w}_c - \mathbf{1} a\| = \arg \min_a \|\mathbf{X} \mathbf{w}_c - \mathbf{1} a\|. \quad (11)$$

708 Then, we have

$$709 \|\epsilon(\mathbf{P}_c \mathbf{X} \mathbf{w}_c)\| = \|\mathbf{P}_c \mathbf{X} \mathbf{w}_c - \mathbf{1} \hat{a}_c\| = \|\mathbf{X} \mathbf{w}_c - \mathbf{1} \hat{a}_c\| \leq \|(\mathbf{X} - \mathbf{1} \mathbf{x}^\top) \mathbf{w}_c\| \leq \|\mathbf{w}_c\| \|\epsilon(\mathbf{X})\|, \quad (12)$$

710 where the second equation is based on the permutation invariance of the matrix norm, the first
711 inequation is based on (11), and the second inequation is based on the sub-multiplicativity (or called
712 consistency) of the matrix norm.

713 **A Single CSP:** Considering all C heads and specifying the matrix norm to be 1-norm and ∞ -norm,
714 respectively, we have

$$715 \begin{aligned} 716 \|\epsilon(\text{CSP}_{\mathbf{W}}(\mathbf{X}))\|_1 &= \|(\|\|_{c=1}^C \mathbf{P}_c \mathbf{X} \mathbf{w}_c) - \mathbf{1} \hat{\mathbf{a}}^\top\|_1 \\ 717 &= \max_c \|\mathbf{P}_c \mathbf{X} \mathbf{w}_c - \mathbf{1} \hat{a}_c\|_1 \\ 718 &\leq (\max_c \|\mathbf{w}_c\|_1) \|\epsilon(\mathbf{X})\|_1 \\ 719 &= \|\mathbf{W}\|_1 \|\epsilon(\mathbf{X})\|_1. \end{aligned} \quad (13)$$

$$720 \begin{aligned} 721 \|\epsilon(\text{CSP}_{\mathbf{W}}(\mathbf{X}))\|_\infty &= \|(\|\|_{c=1}^C \mathbf{P}_c \mathbf{X} \mathbf{w}_c) - \mathbf{1} \hat{\mathbf{a}}^\top\|_\infty \\ 722 &\leq \sum_{c=1}^C \|\mathbf{P}_c \mathbf{X} \mathbf{w}_c - \mathbf{1} \hat{a}_c\|_\infty \\ 723 &\leq \sum_{c=1}^C \|\mathbf{w}_c\|_\infty \|\epsilon(\mathbf{X})\|_\infty \\ 724 &\leq \sum_{c=1}^C \sum_{n=1}^N |w_{nc}| \|\epsilon(\mathbf{X})\|_\infty \\ 725 &\leq C (\max_c \|\mathbf{w}_c\|_1) \|\epsilon(\mathbf{X})\|_\infty \\ 726 &= C \|\mathbf{W}\|_1 \|\epsilon(\mathbf{X})\|_\infty. \end{aligned} \quad (14)$$

727 Combining the above two inequations, we have

$$728 \|\epsilon(\text{CSP}_{\mathbf{W}}(\mathbf{X}))\|_{1,\infty} \leq \sqrt{C} \|\mathbf{W}\|_1 \|\epsilon(\mathbf{X})\|_{1,\infty}. \quad (15)$$

729 **A Single CSP followed by a Lipschitz function.** Given a Lipschitz function $f_\lambda : \mathbb{R}^C \mapsto \mathbb{R}^C$, we
730 apply it to each row of a CSP's output matrix. For the residual of $f_\lambda \circ \text{CSP}_{\mathbf{W}}(\mathbf{X})$, we have

$$731 \begin{aligned} 732 \|\epsilon(f_\lambda \circ \text{CSP}_{\mathbf{W}}(\mathbf{X}))\| &= \|f_\lambda \circ \text{CSP}_{\mathbf{W}}(\mathbf{X}) - \mathbf{1} \hat{\mathbf{y}}^\top\| \\ 733 &\leq \|f_\lambda \circ \text{CSP}_{\mathbf{W}}(\mathbf{X}) - \mathbf{1} f_\lambda^\top(\hat{\mathbf{a}})\| \\ 734 &\leq \lambda \|\text{CSP}_{\mathbf{W}}(\mathbf{X}) - \mathbf{1} \hat{\mathbf{a}}^\top\| \\ 735 &= \lambda \|\epsilon(\text{CSP}_{\mathbf{W}}(\mathbf{X}))\|, \end{aligned} \quad (16)$$

736 where $\hat{\mathbf{y}} = \arg \min_{\mathbf{y}} \|f_\lambda \circ \text{CSP}_{\mathbf{W}}(\mathbf{X}) - \mathbf{1} \mathbf{y}^\top\|$ and $\hat{\mathbf{a}}$ is the vector associated to $\epsilon(\text{CSP}_{\mathbf{W}}(\mathbf{X}))$.

737 **Stacking L CSP Operators:** We can recursively leverage the above results and derive the following
738 inequation:

$$739 \begin{aligned} 740 \|\epsilon((f \circ \text{CSP})^L(\mathbf{X}))\|_{1,\infty} &\leq \lambda_L \sqrt{C} \|\mathbf{W}^{(L)}\|_1 \|\epsilon((f \circ \text{CSP})^{L-1}(\mathbf{X}))\|_{1,\infty} \\ 741 &\leq C^{\frac{L}{2}} \left(\prod_{\ell=1}^L \lambda_\ell \|\mathbf{W}^{(\ell)}\|_1 \right) \|\epsilon(\mathbf{X})\|_{1,\infty} \\ 742 &\leq C^{\frac{L}{2}} (\lambda \beta)^L \|\epsilon(\mathbf{X})\|_{1,\infty}, \end{aligned} \quad (17)$$

743 where $\beta = \max_\ell \|\mathbf{W}^{(\ell)}\|_1$ and $\lambda = \max_\ell \lambda_\ell$. When the f 's are the identity map, we have
744 $\|\epsilon(\text{CSP}^L(\mathbf{X}))\|_{1,\infty} \leq C^{\frac{L}{2}} \beta^L \|\epsilon(\mathbf{X})\|_{1,\infty}$.

B THE OT-BASED EXPLANATION OF CROSS-CHANNEL SORTING

For convenience, denote the group size N/K by G . For $\mathbf{v}_1^{(k)}, \mathbf{v}_c^{(k)} \in \mathbb{R}^G$ The optimal transport distance between $\mathbf{v}_1^{(k)}$ and $\mathbf{v}_c^{(k)}$ can be defined as the following linear programming problem:

$$W(\mathbf{v}_1^{(k)}, \mathbf{v}_c^{(k)}) := \min_{\mathbf{T} \in \Pi_G} \langle \mathbf{D}, \mathbf{T} \rangle, \quad (18)$$

where $\mathbf{D} = (\mathbf{v}_1^{(k)} \odot \mathbf{v}_1^{(k)}) \mathbf{1}_G^\top + \mathbf{1}_G (\mathbf{v}_c^{(k)} \odot \mathbf{v}_c^{(k)})^\top - 2\mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}$ is the squared Euclidean distance matrix, and \odot denotes the Hadamard product. Denote the optimal solution of (18) as \mathbf{T}^* . Because $\mathbf{T} \in \Pi_G$, we have

$$\begin{aligned} \mathbf{T}^* &= \arg \min_{\mathbf{T} \in \Pi_G} \langle \mathbf{D}, \mathbf{T} \rangle \\ &= \arg \min_{\mathbf{T} \in \Pi_G} \langle (\mathbf{v}_1^{(k)} \odot \mathbf{v}_1^{(k)}) \mathbf{1}_G^\top + \mathbf{1}_G (\mathbf{v}_c^{(k)} \odot \mathbf{v}_c^{(k)})^\top - 2\mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}, \mathbf{T} \rangle \\ &= \arg \min_{\mathbf{T} \in \Pi_G} \langle \mathbf{v}_1^{(k)} \odot \mathbf{v}_1^{(k)}, \mathbf{T} \mathbf{1}_G \rangle + \langle \mathbf{v}_c^{(k)} \odot \mathbf{v}_c^{(k)}, \mathbf{T} \mathbf{1}_G \rangle - 2\langle \mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}, \mathbf{T} \rangle \\ &= \arg \min_{\mathbf{T} \in \Pi_G} \underbrace{\langle \mathbf{v}_1^{(k)} \odot \mathbf{v}_1^{(k)}, \mathbf{1}_{G \times G} \rangle + \langle \mathbf{v}_c^{(k)} \odot \mathbf{v}_c^{(k)}, \mathbf{1}_{G \times G} \rangle}_{\text{A Constant } C_0} - 2\langle \mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}, \mathbf{T} \rangle \\ &\Leftrightarrow \arg \min_{\mathbf{T} \in \Pi_G} \langle -\mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}, \mathbf{T} \rangle. \end{aligned} \quad (19)$$

In addition, because $\mathbf{v}_1^{(k)}$ and $\mathbf{v}_c^{(k)}$ are 1D vectors, their OT distance can be computed by aligning the elements of $\mathbf{v}_c^{(k)}$ to align to those of $\mathbf{v}_1^{(k)}$, which corresponds to the sorting operation, i.e.,

$$W(\mathbf{v}_1^{(k)}, \mathbf{v}_c^{(k)}) = \|\mathbf{v}_1^{(k)} - \mathbf{T}_c^{(k)} \mathbf{v}_c^{(k)}\|_2^2 = \underbrace{\langle \mathbf{v}_1^{(k)}, \mathbf{v}_1^{(k)} \rangle + \langle \mathbf{v}_c^{(k)}, \mathbf{v}_c^{(k)} \rangle}_{=C_0} - 2\langle \mathbf{v}_1^{(k)} \mathbf{v}_c^{(k)\top}, \mathbf{T}_c^{(k)} \rangle, \quad (20)$$

where $\mathbf{T}_c^{(k)}$ is the permutation matrix. Therefore, as long as $W(\mathbf{v}_1^{(k)}, \mathbf{v}_c^{(k)})$ has a unique optimal transport, $\mathbf{T}_c^{(k)} = \mathbf{T}^*$.

C IMPLEMENTATION DETAILS

C.1 IMAGE CLASSIFICATION

The detailed hyperparameter setups are presented in Table 5. Both training and testing are conducted on 8 NVIDIA GeForce RTX 4080 SUPER GPUs.

Table 5: The hyperparameters of ViT using CSP on image classification tasks.

Dataset	CIFAR-10	CIFAR-100	ImageNet-1k
#Groups K	32	128	98
Shifting step	Linear	Linear	Linear
Batch Size	64	64	256
Epochs	100	100	300
Learning Rate	1E-04	1E-04	5E-04
LR scheduler	cosine	cosine	cosine
Optimizer	Adam	Adam	AdamW
Dropout Rate	0.1	0.1	0.1
Hidden Dims	512	512	386
Num. Layers	6	6	12
Pooling Type	mean	mean	mean
#Param.	6.46M	6.50M	18.50M

C.2 LONG RANGE ARENA BENCHMARK

We strictly follow the LRA benchmark (Tay et al., 2021b)’s default data processing and experimental design. The detailed hyperparameter setups are presented in Table 6 and Table 7. For Image task,

Table 6: The hyperparameters of Transformer using CSP on LRA.

Dataset	ListOps	Text	Retrieval	Image	PathFind
#Groups K	1	1	1	1	1
Shifting step	—	—	—	—	—
Batch Size	32	32	8	256	256
Train steps	5000	20000	5000	35156	125000
Learning Rate	5E-02	5E-02	5E-02	8E-03	1E-03
LR scheduler	sqrt	sqrt	sqrt	cosine	cosine
Optimizer	Adam	Adam	Adam	Adam	Adam
Weight Decay	1E-01	1E-01	1E-01	0	0
Hidden Dims	512	256	128	128	64
Num. Layers	4	4	4	4	6
Pooling Type	cls	cls	cls	cls	cls

Table 7: The hyperparameters of MEGA using CSP on LRA.

Dataset	ListOps	Text	Retrieval	Image	PathFind	Path-X
#Groups K	1	1	1	—	512	8192
Shifting step	—	—	—	Linear	Linear	Exp
Batch Size	64	25	8	50	64	60
Epochs	60	50	40	200	200	100
Learning Rate	1E-03	4E-03	6E-03	1E-02	3E-02	1E-02
LR scheduler	linear	linear	linear	linear	linear	linear
Optimizer	Adam	Adam	Adam	Adam	Adam	Adam
Weight Decay	1E-02	1E-02	4E-02	2E-02	1E-02	1E-02
Dropout Rate	0.1	0.1	0.1	0.0	0.1	0.5
Hidden Dims	160	256	256	1024	256	128
Num. Layers	6	4	6	8	6	4
Pooling Type	mean	mean	mean	mean	mean	mean

we only apply the circular shifting operation. Both training and testing are conducted on 8 NVIDIA RTX A6000 GPUs.

In Figure 4, we compare Transformer using CSP with other baselines based on JAX Bradbury et al. (2018). These models are trained on 4 NVIDIA GeForce RTX 3090 GPUs. The detailed settings are as follows: The length of the sequence is 3K. The x-axis corresponds to the number of training steps per second. The y-axis corresponds to the average score (%) on the LRA benchmark. The peak memory usage of each model is represented as the area of the corresponding circle. For a better comparison, the values (GB) of the top-2 models are shown.