

---

# Composite Attention: A Framework for Combining Sequence Mixing Primitives

---

**Harry Jake Cunningham\***  
Centre for Artificial Intelligence  
University College London

**Marc Peter Deisenroth**  
Centre for Artificial Intelligence  
University College London

## Abstract

Hybrid attention architectures have shown promising success in both equipping self attention with inductive bias for long-sequence modelling and reducing the computational burden of transformers without sacrificing quality. This paper introduces Composite Attention, a theoretical framework for analyzing the combination of sequence mixing primitives in modern deep learning architectures. Utilizing the definition of sequence mixers as structured linear maps, we formalize the composition of sequence mixing primitives as either *sequential* or *recurrent* composition.

## 1 Introduction

The design space of sequence models has recently exploded due to the successful introduction of long convolutions [17, 3, 8, 26], state-space models [12, 10, 11, 27], linear attention mechanisms [9, 4, 28, 15], and gating [24, 13]. Hybrid attention models [21, 22, 2] seek to make use of these advancement by combining different sequence mixing primitives to improve efficiency and performance. In this work, we examine the combination of sequence mixing primitives within hybrid attention architectures and propose a framework for enhancing the understanding of composite sequence models through analysis of their matrix structures. Specifically, we aim to: 1) Formalize the composition of sequence mixing primitives, 2) Identify gating as a method for diversifying single-head attention, and 3) Recognize convolution-equipped attention mechanisms as a means of encoding local context.

## 2 Background

**Sequence Mixers** Modern deep learning sequence mixers rely on the aggregation of information as the weighted sum of tokens. Following [14], we can define any sequence mixer which performs aggregation as a linear map  $\mathbf{y} = \mathbf{A}\mathbf{x}$ , where  $\mathbf{A} \in \mathbb{R}^{L \times L}$  is a matrix of *attention* weights and  $\mathbf{x} \in \mathbb{R}^{L \times L}$  is an input sequence of length  $L$  and dimension  $d$ .

**Definition 2.1** (Token Mixing as Linear Maps.). *Let  $\mathbf{x} \in \mathbb{R}^{L \times D}$  be a sequence of tokens, of length  $L$  and embedding dimension  $D$ . Let  $\mathcal{M} \subseteq \mathbb{R}^{L \times L}$  be a space of possible attention matrices. Let  $H$  be the number of heads and  $P$  the head dimension such that  $HP = D$ . For each head  $h \in H$ , we define a generating function  $f_{\mathcal{M}}^h : \mathbb{R}^{L \times D} \times \Theta \rightarrow \mathcal{M}$ , parameterized by  $\Theta$ , that maps the input to a space of possible attention matrices. Denoting the attention weight matrix  $\mathbf{A}^h = f_{\mathcal{M}}^h(\mathbf{x}, \Theta) \in \mathbb{R}^{L \times L}$ , we define sequence mixing via aggregation as,*

$$\mathbf{y}^h = \mathbf{A}^h \mathbf{x}^h \tag{1}$$

where  $\mathbf{x}^h \in \mathbb{R}^{L \times P}$  is the input for a given head  $h$ ,  $\mathbf{y}^h \in \mathbb{R}^{L \times P}$  is the output for a given head  $h$ , and  $\mathbf{y} = \text{Concat}[\mathbf{y}^0, \dots, \mathbf{y}^{h-1}] \in \mathbb{R}^{L \times D}$  is the same shape as the input.

---

\*Correspondence to [jake.cunningham.21@ucl.ac.uk](mailto:jake.cunningham.21@ucl.ac.uk)

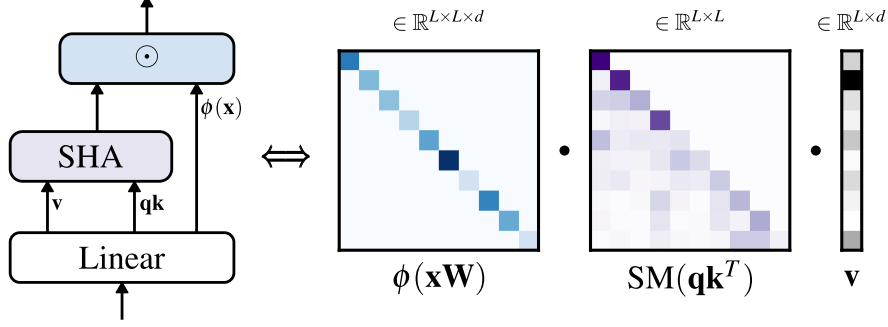


Figure 1: **Gated Attention Unit as Sequential Composition of Sequence Mixers.** (Left) A gated attention unit which consists of single-head attention (SHA) followed by an element-wise product  $\odot$ , where  $\phi$  is an element-wise activation function. (Right) The equivalent compositional form, denoting a *sequential* input dependency structure and their corresponding generating functions.

For notational ease, we define  $f_{\mathcal{M}}(\cdot) = \text{Concat}[f_{\mathcal{M}}^0(\cdot), \dots, f_{\mathcal{M}}^{h-1}(\cdot)]$  as the generating function that returns an attention matrix  $\mathbf{A} \in \mathbb{R}^{H \times L \times L}$  for all heads  $H$ .

**Hybrid Architectures** There is a growing interest in designing *compositional* or hybrid architectures that integrate different sequence mixing primitives to obtain a best-of-all-world solution. Gating has proven a powerful method of enhancing a models input dependency [19, 13, 1, 7, 22]. Depthwise separable convolutions placed before linear attention mechanisms have been shown to reduce the quality gap to SoftMax attention [21, 2, 5, 9, 7, 20, 22]. Hybrid architectures which alternate between SoftMax attention layers and simpler sequence mixers have proven to be more computationally efficient and as performant [17, 2, 23, 25, 18].

### 3 Framework for Composition of Sequence Mixing Primitives

We introduce *Composite Attention*, a framework for analysing the combination of sequence mixing primitives. We can combine input-dependent transformations in 2 different ways: via 1) the *sequential* composition of linear maps where each map depends on the same input  $\mathbf{x}$  and 2) the *recurrent* composition of linear maps where each map depends on the output of the previous transformation.

**Definition 3.1** (Sequential Composition of Sequence Mixers). *Let  $\mathbf{x} \in \mathbb{R}^{L \times d}$  be a sequence of tokens, of length  $L$  and embedding dimension  $d$ . Let  $\{\mathcal{M}^{(i)}\}_{i=0}^{N-1}$  be a set of classes of attention matrices, where  $\mathcal{M}^{(i)} \subseteq \mathbb{R}^{L \times L}$ . Let  $\{f_{\mathcal{M}}^{(i)}\}_{i=0}^{N-1}$  be a set of generating functions, where each  $f_{\mathcal{M}}^{(i)}: \mathbf{X} \times \Phi^{(i)} \rightarrow \mathcal{M}^{(i)}$  is a function that maps from the input to the space of structured matrices  $\mathcal{M}^{(i)}$ . We define the sequential composition of Attention matrices as*

$$\mathbf{y} = g\left(f_{\mathcal{M}}^{(N-1)}(\mathbf{x}), \dots, f_{\mathcal{M}}^{(1)}(\mathbf{x}), f_{\mathcal{M}}^{(0)}(\mathbf{x})\right) \cdot \mathbf{x} \quad (2)$$

$$= \mathbf{A}_S \cdot \mathbf{x} \quad (3)$$

where  $g: \mathcal{M}^{(0)} \times \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(N-1)} \rightarrow \mathbb{R}^{L \times L}$  is a function that takes as input a set of  $N$  structured matrices and returns  $\mathbf{A}_S \in \mathbb{R}^{L \times L}$ , denoted the compositional attention matrix corresponding to the composition of attention matrices, which each depend on the input  $\mathbf{x}$ , by the function  $g(\cdot)$ .

**Remark.** Typically we assume that  $g(\cdot)$  is a composition of binary operations such as matrix additions or matrix multiplications. For the case of matrix multiplications the sequential composition can be constructed as

$$\mathbf{y} = \left(f_{\mathcal{M}}^{(N-1)}(\mathbf{x}) \cdot \dots \cdot f_{\mathcal{M}}^{(1)}(\mathbf{x}) \cdot f_{\mathcal{M}}^{(0)}(\mathbf{x})\right) \cdot \mathbf{x} \quad (4)$$

$$= \mathbf{A}^{(N-1)} \cdot \dots \cdot \mathbf{A}^{(1)} \cdot \mathbf{A}^{(0)} \cdot \mathbf{x} \quad (5)$$

where  $\mathbf{A}_S = \mathbf{A}^{(N-1)} \cdot \dots \cdot \mathbf{A}^{(1)} \cdot \mathbf{A}^{(0)}$  is the compositional attention matrix. We stress though that our definition is general and encompasses potentially non-linear functions of the input matrices.

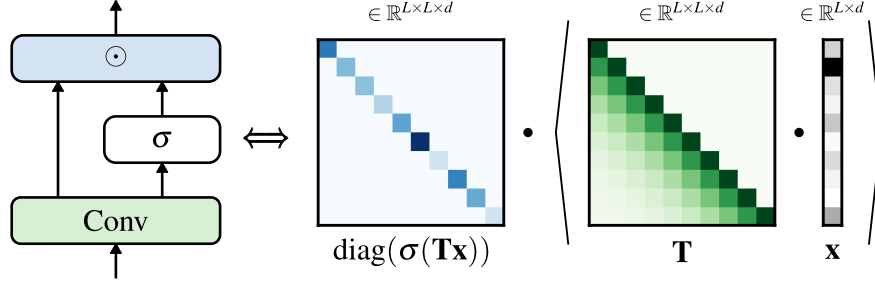


Figure 2: **Silu Activation as Recurrent Composition of Sequence Mixers.** (Left) A global depthwise separable convolution followed by a silu activation function, where  $\sigma$  is the logistic sigmoid applied element-wise. (Right) The equivalent compositional form, denoting a *recurrent* input dependency structure and their corresponding generating functions.

Hierarchical gating, as used in gated attention units (GAU) (See Figure 1), is an example of sequential composition whereby each sequence mixer depends upon the same input.

**Definition 3.2** (Recurrent Composition of Sequence Mixers). *Let  $\mathbf{x} \in \mathbb{R}^{L \times d}$  be a sequence of tokens, of length  $L$  and embedding dimension  $d$ . Let  $\{\mathcal{M}^{(i)}\}_{i=0}^{N-1}$  be a set of classes of attention matrices, where  $\mathcal{M}^{(i)} \subseteq \mathbb{R}^{L \times L}$ . Let  $\{f_{\mathcal{M}}^{(i)}\}_{i=0}^{N-1}$  be a set of generating functions, where each  $f_{\mathcal{M}}^{(i)} : \mathbf{X} \times \Phi^{(i)} \rightarrow \mathcal{M}^{(i)}$  is a function that maps from the input to the space of structure matrices  $\mathcal{M}^{(i)}$ . We define the sequential composition of Attention matrices as*

$$\mathbf{y}_n = f_{\mathcal{M}}^{(i)}(\mathbf{y}_{n-1}) \cdot \mathbf{y}_{n-1} \quad (6)$$

$$\mathbf{y}_0 = \mathbf{x} \quad (7)$$

where the output  $\mathbf{y}_n$  is sequentially generated.

Recurrent composition is typically found in branches of a gated architecture such as a short convolution followed by a silu activation function [6] as used in Mamba [9] (see Figure 2).

## 4 Composite Attention

We will use our framework to analyze combinations of sequence mixing primitives, focusing on the properties of their attention matrices, generating functions, and matrix products, especially looking at gating and convolution-equipped attention.

### 4.1 Exploration of Composite Sequence Mixer Framework

We begin with a Gated Attention Unit (GAU), defined as,

$$\mathbf{u} = \phi_u(\mathbf{x}\mathbf{W}_u), \quad \mathbf{v} = \phi_v(\mathbf{x}\mathbf{W}_v), \quad \mathbf{g} = \phi_g(\mathbf{x}\mathbf{W}_g) \quad (8)$$

$$\mathbf{q} = \kappa_q \odot \mathbf{u} + \mu_q, \quad \mathbf{k} = \kappa_k \odot \mathbf{u} + \mu_k \quad (9)$$

$$\mathbf{a} = \text{SoftMax} \left( \frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d}} \right) \mathbf{v} \quad (10)$$

$$\mathbf{y} = \mathbf{a} \odot \mathbf{g} \quad (11)$$

This structure can be described similarly to hierarchical gating as the *sequential composition* of attention matrices, where both *self-attention* and *gating* correspond to input-dependent sequence mixers whose generating function depends on the input sequence  $\mathbf{x}$  as,

$$\mathbf{y} = \text{diag}(\phi_g(\mathbf{x}\mathbf{W}_g)) \cdot \text{SoftMax} \left( \frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d}} \right) \cdot \phi_u(\mathbf{x}\mathbf{W}_v) \quad (12)$$

$$= f_G(\mathbf{x}) \cdot f_A(\mathbf{x}) \cdot \phi_u(\mathbf{x}\mathbf{W}_v) \quad (13)$$

$$= \mathbf{A}_G \cdot \mathbf{A}_A \cdot \mathbf{v} \quad (14)$$

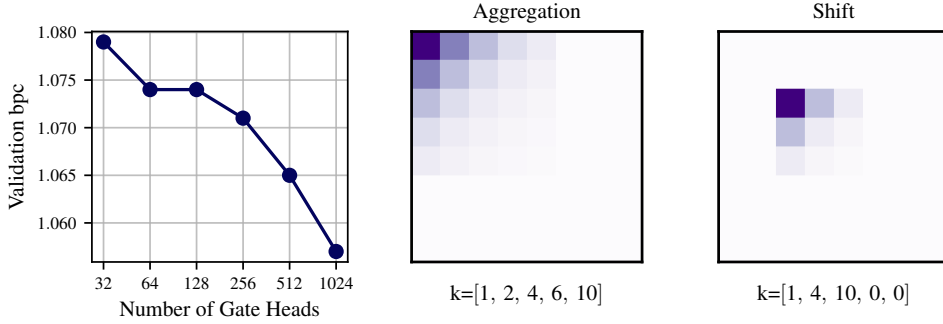


Figure 3: **Left.** Gating Heads Ablation on enwik8. Performance of *GAU* improves with increasing head dimension on enwik8 language modelling task. **Center.** Outer product of unit impulse convolved with an *aggregation* filter. Shows how convolution equipped attention mechanisms spread token information across the attention matrix, encoding context. **Right.** Outer product of unit impulse convolved with a *shift* filter. Demonstrates ability of local convolutions to perform local token shifts, acting as a form of memory.

**Gating Diversifies Single-Head Attention.** As a sequence mixer gating can be represented as an  $L \times L$  diagonal matrix, with generating function  $f_G(\mathbf{x}) = \text{diag}(\phi(\mathbf{x}\mathbf{W})^T) \in \mathbb{R}^{L \times L \times h_g}$  where  $h_g$  corresponds to the number of gating heads, which is typically set to the hidden dimension  $d$ . We *hypothesize* that the effectiveness of gating is derived from its use of a very large number of heads which when composed with single-head attention (SHA) or few-head attention diversifies the attention layer into a multi-headed linear map over the sequence dimension

$$\mathbf{A}_G = f_G(\mathbf{x}) = \text{diag}(\phi(\mathbf{x}\mathbf{W})^T) \in \mathbb{R}^{L \times L \times h_g} \quad (15)$$

$$\mathbf{A}_{SHA} = f_{SHA}(\mathbf{x}) = \text{SoftMax} \left( \frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d}} \right) \in \mathbb{R}^{L \times L \times 1} \quad (16)$$

$$\mathbf{A}_{GAU} = \mathbf{A}_G \cdot \text{Repeat}(\mathbf{A}_{SHA}, h_g) \in \mathbb{R}^{L \times L \times h_g} \quad (17)$$

We emphasize the importance of using many heads when gating by conducting an Ablation study using the MEGA transformer [21]. Indeed, we find that increasing the number of gating heads improves performance on the enwik8 dataset (see Figure 3)

**Convolutions Encode Context.** Placing convolutional layers prior to computing the  $\mathbf{q}, \mathbf{k}$  projections has been proposed as an alternative to positional encodings [16, 21] as well as a means of introducing locality which is required to capture long-range dependencies by means of a hierarchical combination of local dependencies [4, 7, 3]. We examine the effect of convolution equipped attention by *composing together* the convolution operator with the generating function for self-attention.

**Definition 4.1** (Convolution Equipped Attention). *Let  $\mathbf{x} \in \mathbb{R}^{L \times d}$  be an input of length  $L$  and dimension  $d$ ,  $\mathbf{k}_q = [\alpha_0, \dots, \alpha_{L-1}] \in \mathbb{R}^{L \times d}$  and  $\mathbf{k}_k = [\beta_0, \dots, \beta_{L-1}] \in \mathbb{R}^{L \times d}$  the query and key convolution kernels of length  $L$  and  $\mathbf{W}_q, \mathbf{W}_k \in \mathbb{R}^{d \times d}$  dense linear projections over the channel dimension  $d$ . Let the query and keys be defined as  $\mathbf{q} = (\mathbf{k}_q * \mathbf{x})\mathbf{W}_q$  and  $\mathbf{k} = (\mathbf{k}_k * \mathbf{x})\mathbf{W}_k$  respectively. Hence,*

$$\text{SoftMax} \left( \frac{\mathbf{q}_i^T \mathbf{k}_j}{\sqrt{d}} \right) = \text{SoftMax} \left( \frac{1}{\sqrt{d}} \sum_{n=0}^i \sum_{m=0}^j (\alpha_{i-n} \odot \mathbf{x}_n)^T \mathbf{W}_q^T \mathbf{W}_k (\beta_{j-m} \odot \mathbf{x}_m) \right) \quad (18)$$

*corresponds to the convolution equipped attention matrix.*

Here we can see how convolutions are able to spread the influence of tokens across the attention matrix through the convolution weights  $\alpha$  and  $\beta$ . We note that convolutions can do this in 2 ways: 1) by aggregating local tokens providing positional information and 2) by performing local token shifts (e.g.  $[0, 1, 1] * [a, b, c] = [0, a, b]$ ) acting as a form of memory [7, 12]. We illustrate both of these effects in Figure 3 by convolving a unit impulse with convolution kernels corresponding to *aggregations* and *shifts* and plotting the outer product  $\mathbf{q}\mathbf{k}^T$ .

Further, we can see how the dot product between tokens in standard SoftMax attention has been replaced by a sum of dot products that evaluate the similarity between tokens that appear before the tokens of interest. Hence, convolution equipped attention amplifies attention weights for tokens whose contexts are similar.

We define *convolution equipped GAU* by the following composition,

$$\mathbf{y} = f_G(f_{C_1}(\mathbf{x})) \cdot f_A(f_{C_0}(\mathbf{x})) \cdot \mathbf{v} \quad (19)$$

$$= f_G(\mathbf{T}_1 \cdot \mathbf{x}) \cdot f_A(\mathbf{T}_0 \cdot \mathbf{x}) \cdot \mathbf{v} \quad (20)$$

where  $\mathbf{T}_0$  and  $\mathbf{T}_1$  correspond to Toeplitz matrices representing the non-input dependent convolutions, which are followed by a linear projection and silu activation defined within the generating function  $f_A$ , corresponding to a *recurrent* composition.

## 5 Conclusion

We have developed Composite Attention, a framework for combining sequence mixing primitives using both sequential and recurrent compositions of matrix structures. Through our framework, we analyze the impact of gating and convolutions in a gated attention unit. We find that gating enhances single-head attention by incorporating multiple heads, while convolutions encode context by performing local aggregation and shifting. Looking ahead, we aim to utilize our framework to create new sub-quadratic architectures, with a specific focus on improving linear attention mechanisms using cost-effective sequence mixers.

## References

- [1] S. Arora, S. Eyuboglu, A. Timalina, I. Johnson, M. Poli, J. Zou, A. Rudra, and C. Ré. Zoology: Measuring and improving recall in efficient language models. *arXiv preprint arXiv:2312.04927*, 2023.
- [2] S. Arora, S. Eyuboglu, M. Zhang, A. Timalina, S. Alberti, D. Zinsley, J. Zou, A. Rudra, and C. Ré. Simple linear attention language models balance the recall-throughput tradeoff. *arXiv preprint arXiv:2402.18668*, 2024.
- [3] H. J. Cunningham, G. Giannone, M. Zhang, and M. P. Deisenroth. Reparameterized multi-resolution convolutions for long sequence modelling. *arXiv preprint arXiv:2408.09453*, 2024.
- [4] T. Dao and A. Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- [5] S. De, S. L. Smith, A. Fernando, A. Botev, G. Cristian-Muraru, A. Gu, R. Haroun, L. Berrada, Y. Chen, S. Srinivasan, et al. Griffin: Mixing gated linear recurrences with local attention for efficient language models. *arXiv preprint arXiv:2402.19427*, 2024.
- [6] S. Elfving, E. Uchibe, and K. Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 107:3–11, 2018.
- [7] D. Y. Fu, T. Dao, K. K. Saab, A. W. Thomas, A. Rudra, and C. Ré. Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052*, 2022.
- [8] D. Y. Fu, E. L. Epstein, E. Nguyen, A. W. Thomas, M. Zhang, T. Dao, A. Rudra, and C. Ré. Simple hardware-efficient long convolutions for sequence modeling. In *International Conference on Machine Learning*, pages 10373–10391. PMLR, 2023.
- [9] A. Gu and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [10] A. Gu, K. Goel, and C. Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- [11] A. Gu, K. Goel, A. Gupta, and C. Ré. On the parameterization and initialization of diagonal state space models. *Advances in Neural Information Processing Systems*, 35:35971–35983, 2022.
- [12] A. Gu, I. Johnson, A. Timalina, A. Rudra, and C. Ré. How to train your hippo: State space models with generalized orthogonal basis projections. *arXiv preprint arXiv:2206.12037*, 2022.

- [13] W. Hua, Z. Dai, H. Liu, and Q. Le. Transformer quality in linear time. In *International conference on machine learning*, pages 9099–9117. PMLR, 2022.
- [14] S. Hwang, A. Lahoti, T. Dao, and A. Gu. Hydra: Bidirectional state space models through generalized matrix mixers. *arXiv preprint arXiv:2407.09941*, 2024.
- [15] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pages 5156–5165. PMLR, 2020.
- [16] M. Li, X. Zhang, Y. Huang, and S. Oymak. On the power of convolution augmented transformer. *arXiv preprint arXiv:2407.05591*, 2024.
- [17] Y. Li, T. Cai, Y. Zhang, D. Chen, and D. Dey. What makes convolutional models great on long sequence modeling? *arXiv preprint arXiv:2210.09298*, 2022.
- [18] O. Lieber, B. Lenz, H. Bata, G. Cohen, J. Osin, I. Dalmedigos, E. Safahi, S. Meirum, Y. Belinkov, S. Shalev-Shwartz, et al. Jamba: A hybrid transformer-mamba language model. *arXiv preprint arXiv:2403.19887*, 2024.
- [19] H. Liu, Z. Dai, D. So, and Q. V. Le. Pay attention to mlps. *Advances in neural information processing systems*, 34:9204–9215, 2021.
- [20] Z. Liu, S. Li, L. Wang, Z. Wang, Y. Liu, and S. Z. Li. Short-long convolutions help hardware-efficient linear attention to focus on long sequences. *arXiv preprint arXiv:2406.08128*, 2024.
- [21] X. Ma, C. Zhou, X. Kong, J. He, L. Gui, G. Neubig, J. May, and L. Zettlemoyer. Mega: moving average equipped gated attention. *arXiv preprint arXiv:2209.10655*, 2022.
- [22] X. Ma, X. Yang, W. Xiong, B. Chen, L. Yu, H. Zhang, J. May, L. Zettlemoyer, O. Levy, and C. Zhou. Megalodon: Efficient llm pretraining and inference with unlimited context length. *arXiv preprint arXiv:2404.08801*, 2024.
- [23] J. Park, J. Park, Z. Xiong, N. Lee, J. Cho, S. Oymak, K. Lee, and D. Papailiopoulos. Can mamba learn how to learn? a comparative study on in-context learning tasks. *arXiv preprint arXiv:2402.04248*, 2024.
- [24] M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, S. Baccus, Y. Bengio, S. Ermon, and C. Ré. Hyena hierarchy: Towards larger convolutional language models. In *International Conference on Machine Learning*, pages 28043–28078. PMLR, 2023.
- [25] L. Ren, Y. Liu, Y. Lu, Y. Shen, C. Liang, and W. Chen. Samba: Simple hybrid state space models for efficient unlimited context language modeling. *arXiv preprint arXiv:2406.07522*, 2024.
- [26] D. W. Romero, A. Kuzina, E. J. Bekkers, J. M. Tomczak, and M. Hoogendoorn. Ckconv: Continuous kernel convolution for sequential data. *arXiv preprint arXiv:2102.02611*, 2021.
- [27] J. T. Smith, A. Warrington, and S. W. Linderman. Simplified state space layers for sequence modeling. *arXiv preprint arXiv:2208.04933*, 2022.
- [28] S. Yang, B. Wang, Y. Shen, R. Panda, and Y. Kim. Gated linear attention transformers with hardware-efficient training. *arXiv preprint arXiv:2312.06635*, 2023.

## A Appendix

### A.1 Convolutions Encode Context

Let  $\mathbf{x} \in \mathbb{R}^{L \times d}$  be an input. We define the queries and key vectors as,

$$\mathbf{q} = (\mathbf{k}_q * \mathbf{x})\mathbf{W}_q = \mathbf{T}_q \mathbf{x} \mathbf{W}_q \quad (21)$$

$$\mathbf{k} = (\mathbf{k}_k * \mathbf{x})\mathbf{W}_k = \mathbf{T}_k \mathbf{x} \mathbf{W}_k \quad (22)$$

where  $\mathbf{T}_q, \mathbf{T}_k \in \mathbb{R}^{L \times L \times d}$  are Toeplitz matrices which represent a depthwise separable convolution as a matrix-vector product over the sequence length dimension  $L$  and  $\mathbf{W}_q, \mathbf{W}_k \in \mathbb{R}^{d \times d}$  are dense linear projections over the channel dimension  $d$ .

Splitting the convolution matrix and input into  $d$  heads we compute the depthwise separable convolution for a single head  $i$  as

$$\mathbf{T}_q^i \mathbf{x}^i = \begin{bmatrix} \alpha_0^i & 0 & 0 & 0 & 0 & 0 \\ \alpha_1^i & \alpha_0^i & 0 & 0 & 0 & 0 \\ \vdots & \alpha_1^i & \alpha_0^i & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \alpha_{L-2}^i & \ddots & \ddots & \alpha_1 & \alpha_0^i & 0 \\ \alpha_{L-1}^i & \alpha_{L-2}^i & \dots & \dots & \alpha_1^i & \alpha_0^i \end{bmatrix} \begin{bmatrix} x_0^i \\ x_1^i \\ \vdots \\ \vdots \\ \vdots \\ x_{L-1}^i \end{bmatrix} = \begin{bmatrix} \alpha_0^i x_0^i \\ \alpha_1^i x_0^i + \alpha_0^i x_1^i \\ \vdots \\ \vdots \\ \vdots \\ \sum_{l=0}^{L-1} \alpha_{L-1-l}^i x_l^i \end{bmatrix} \in \mathbb{R}^L \quad (23)$$

and hence for all heads,

$$\mathbf{T}_q \mathbf{x} = \begin{bmatrix} \boldsymbol{\alpha}_0 \odot \mathbf{x}_0 \\ \boldsymbol{\alpha}_1 \odot \mathbf{x}_0 + \boldsymbol{\alpha}_0 \odot \mathbf{x}_1 \\ \vdots \\ \vdots \\ \vdots \\ \sum_{l=0}^{L-1} \boldsymbol{\alpha}_{L-1-l} \odot \mathbf{x}_l \end{bmatrix} \in \mathbb{R}^{L \times d} \quad (24)$$

such that the queries and keys are represented as,

$$\mathbf{q} = \begin{bmatrix} \mathbf{W}_q(\boldsymbol{\alpha}_0 \odot \mathbf{x}_0) \\ \mathbf{W}_q(\boldsymbol{\alpha}_1 \odot \mathbf{x}_0 + \boldsymbol{\alpha}_0 \odot \mathbf{x}_1) \\ \vdots \\ \vdots \\ \vdots \\ \sum_{l=0}^{L-1} \mathbf{W}_q(\boldsymbol{\alpha}_{L-1-l} \odot \mathbf{x}_l) \end{bmatrix} \in \mathbb{R}^{L \times d}, \quad \mathbf{k} = \begin{bmatrix} \mathbf{W}_k(\boldsymbol{\beta}_0 \odot \mathbf{x}_0) \\ \mathbf{W}_k(\boldsymbol{\beta}_1 \odot \mathbf{x}_0 + \boldsymbol{\beta}_0 \odot \mathbf{x}_1) \\ \vdots \\ \vdots \\ \vdots \\ \sum_{l=0}^{L-1} \mathbf{W}_k(\boldsymbol{\beta}_{L-1-l} \odot \mathbf{x}_l) \end{bmatrix} \in \mathbb{R}^{L \times d} \quad (25)$$

Considering the case of single-head attention, the outer product  $\mathbf{q}\mathbf{k}^T$  is defined as,

$$\mathbf{q}_i^T \mathbf{k}_j = \sum_{n=0}^i \sum_{m=0}^j (\boldsymbol{\alpha}_{i-n} \odot \mathbf{x}_n)^T \mathbf{W}_q^T \mathbf{W}_k (\boldsymbol{\beta}_{j-m} \odot \mathbf{x}_m) \quad (26)$$