

Mamba: Linear-Time Sequence Modeling with Selective State Spaces

Arxiv 2312

Albert Gu, Tri Dao

MambaMixer: 一种高效选择性状态空间模型，它使用跨 token 和通道的双重选择机制，称为选择性 token 和通道混合器，在视觉和时序预测任务上表现出色！代码即将开源👍 单位：康奈尔大学等

Sigma: 第一个成功将 Mamba 应用于多模态语义分割的新工作, 采用 **状态空间模型**, 它使用跨 token 和通道的 Siamese 编码器并创新 Mamba 融合机制, 有效地从不同模态中选择重要 **通道混合器**, 在视觉和时序预测任务信息, 性能表现 SOTA! 代码已开源👍: 康奈尔大学等



Mamba 再下一城! Sigma: 多模态语义分割的孪生 Mamba 网络

双星
上表现

让Mamba涨点！本文在Mamba层和自注意力层之间建立了重要的联系，直接将高效的Mamba层与Transformer层连接起来，还提出Mamba第一个可解释性技术！代码刚开源👍



涨点！Mamba 模型的隐藏注意力

Sigma: 第一个从 Siamese 编码器中提取信息，性能表现 SOTA

和通道的时序预测任务

Mamba 再

定义分割的孪生 Mamba 网络

涨点！Sigma: 多模态语

在不同模态中选择重要信息，采用通道混合和注意力机制，康奈尔大学等

双星
上表现

InsectMamba: 第一个将 Mamba 应用于害虫分类的网络, 集成了 SSM, CNN, 多头自注意力机制和 MLP, 在五个害虫分类数据集上性能表现出色👍

InsectMamba: 基于状态空间模型的害虫分类

! 本文在 Mamba 层和自注意力层之间建立了重要的联系, 直接将高效的 Mamba 层与 Transformer 层连接起来, 还提出 Mamba 第一个可解释性技术! 代码刚开源👍

涨点! Mamba 模型的隐藏注意力

和通道的
时序预测任务

Mamba 再

定义分割的孪生 Mamba 网络

Sigma: 多模态语

不同模态中选择重要, 采用通道混合, 康奈尔大学等

双星
上表现

InsectMamba: 第一个将 Mamba 应用于害虫分类的网络, CNN, 多头自注意力机制和 MLP, 在五个害虫分类数据集表现出色👍
据称, RS3Mamba: 第一个专门为遥感图像语义分割设计的视觉 Mamba, 利用 VSS 块, 并引入协作完成模块 (CCM), 性能表现出色!
代码即将开源👍
单位: 香港中文大学, 武汉科技大学

InsectMamba: 基于状态空间模型的害虫分类
性能表现 SOTA

RS3Mamba: 遥感图像语义分割的视觉状态空间模型

直接将高效的 Mamba 层与 Transformer 层连接
第一个可解释性技术! 代码刚开源👍

涨点! Mamba 模型的隐藏注意力

Mamba 再进阶! Sigma: 多模态语义分割的孪生 Mamba 网络

双星上表现

排序任务

不同模态中选择重要通道混合



Albert Gu @_albertgu · 2023年12月5日

Quadratic attention has been indispensable for information-dense modalities such as language... until now.

Announcing Mamba: a new SSM arch. that has linear-time scaling, ultra long context, and most importantly--outperforms Transformers everywhere we've tried.

With @tri_dao 1/



InsectMamba: 第一个将 Mamba 应用于 CNN, 多头自注意力机制和 MLP, 在色👍

InsectMamba: 基于状态... 害虫分类... 直接将... 第一个可解... 性能表现 S...

Mamba 再... 义分割的孪生 Mamba 网络... 哦! Sign... Mamba 网...

义分割设计的视觉... (CM), 性能表现出色!

排序... 任务

Background

- Foundation models
- Large models pretrained on massive data then adapted for downstream tasks
- Backbone : sequence models

Transformer

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

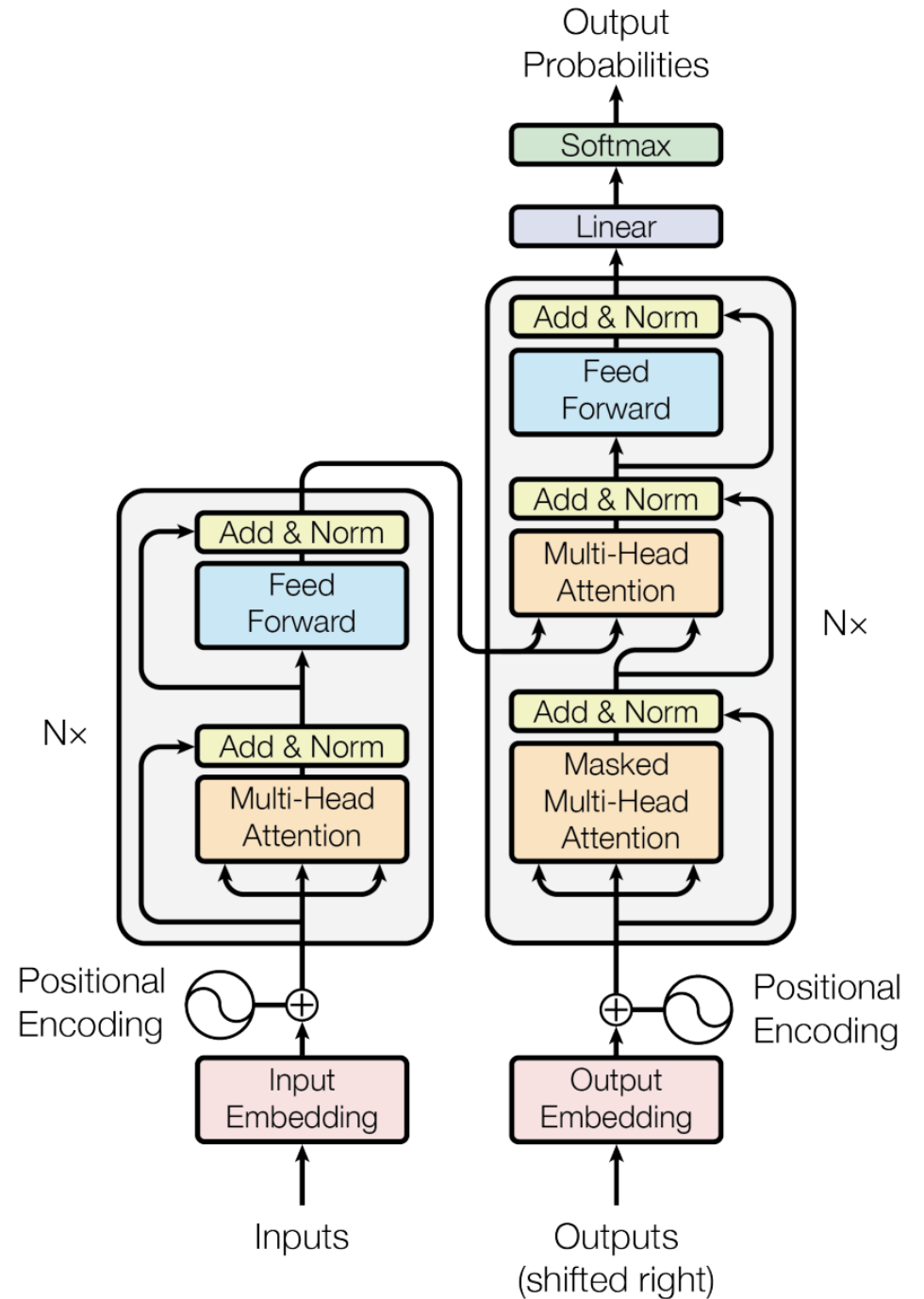
Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

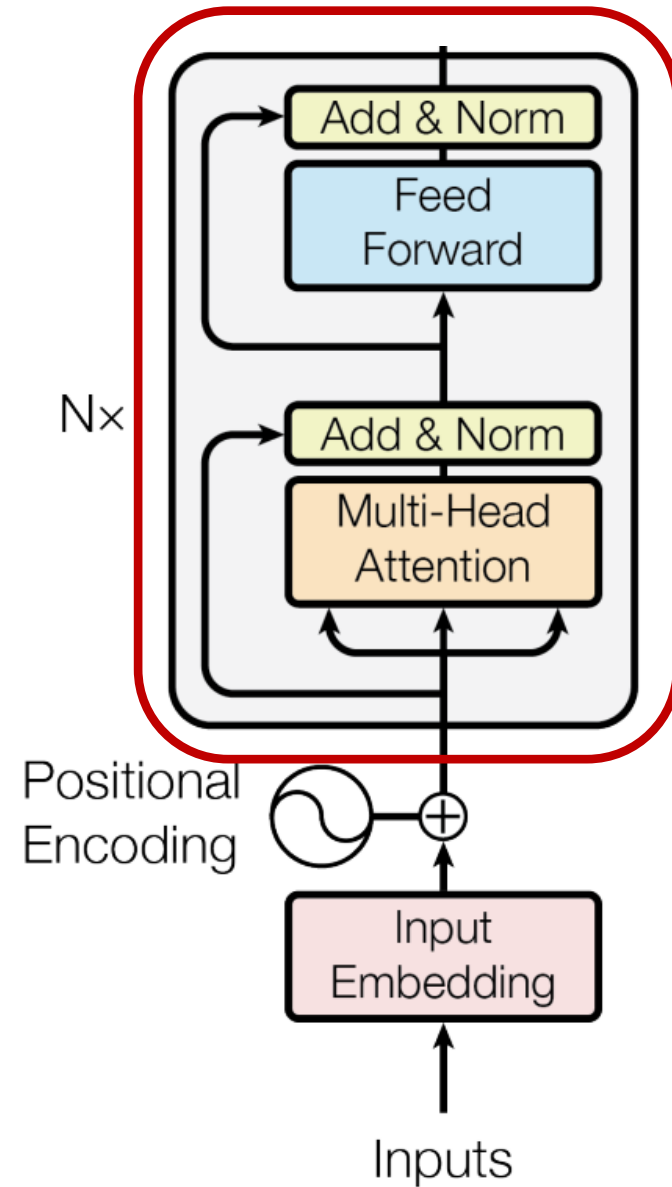
Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* †
illia.polosukhin@gmail.com



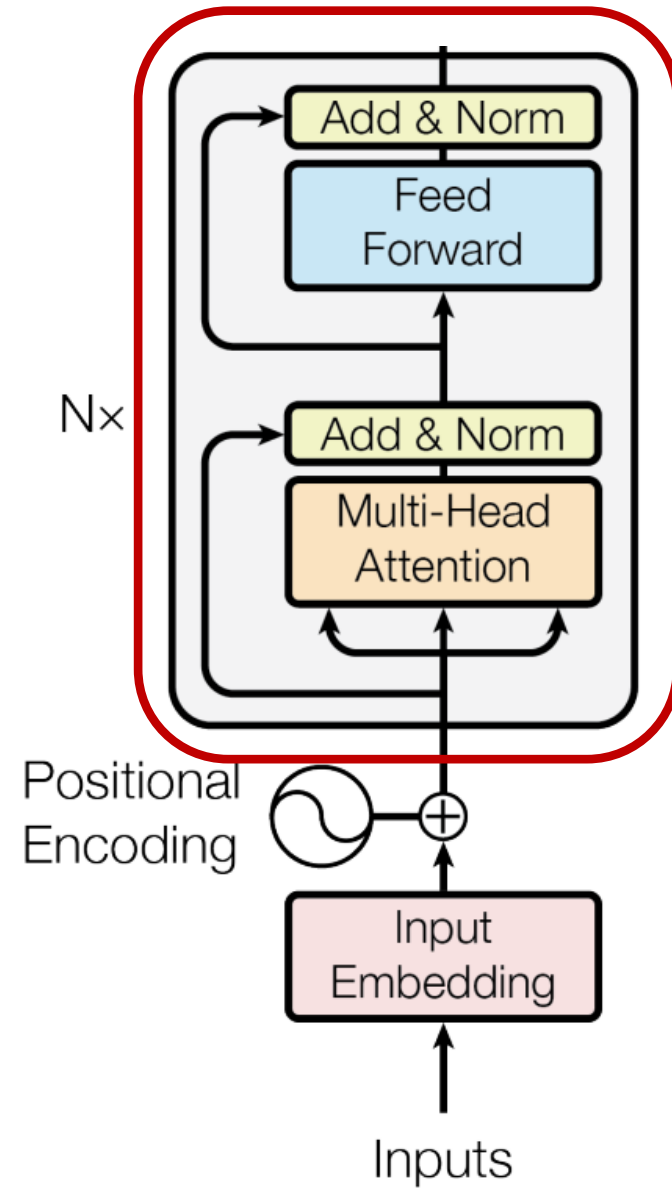
Transformer

- Encoder:
- Multi-Head Self-Attention
- Feed Forward
- Residual Connection & Layer Norm



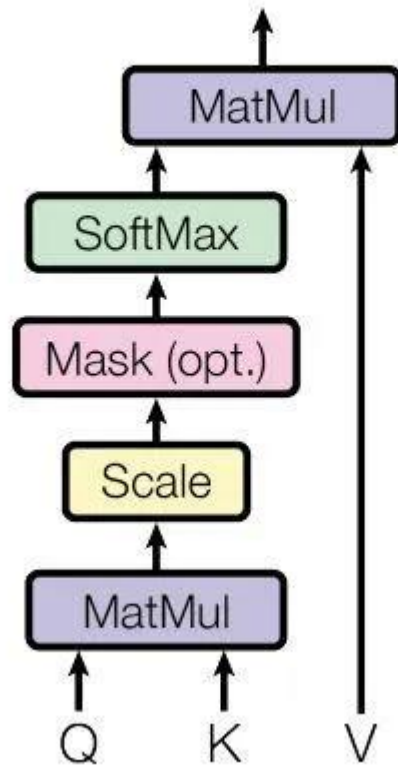
Transformer

- Encoder:
- Multi-Head Self-Attention
- Feed Forward
- Residual Connection & Layer Norm



Self-Attention

Scaled Dot-Product Attention



Input

Thinking

Machines

Embedding

x_1

x_2

Queries

q_1

q_2



W^Q

Keys

k_1

k_2



W^K

Values

v_1

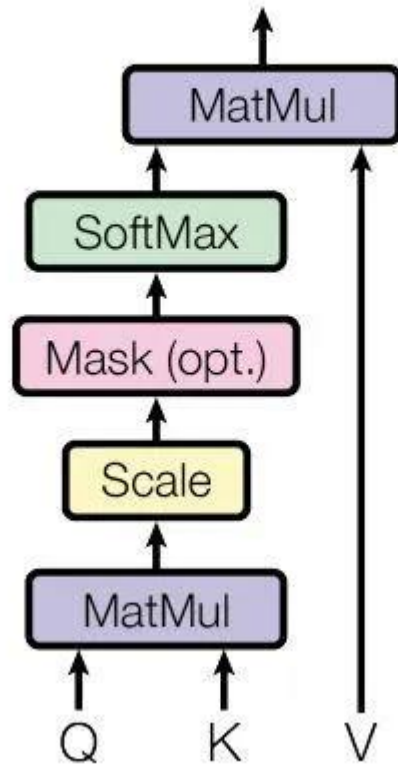
v_2



W^V

Self-Attention

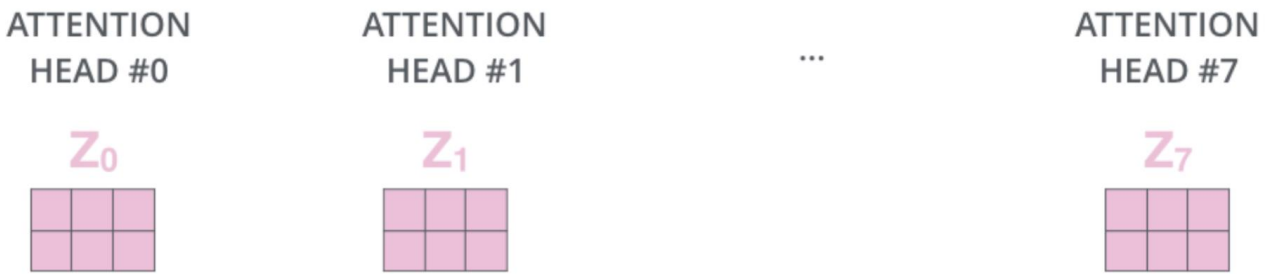
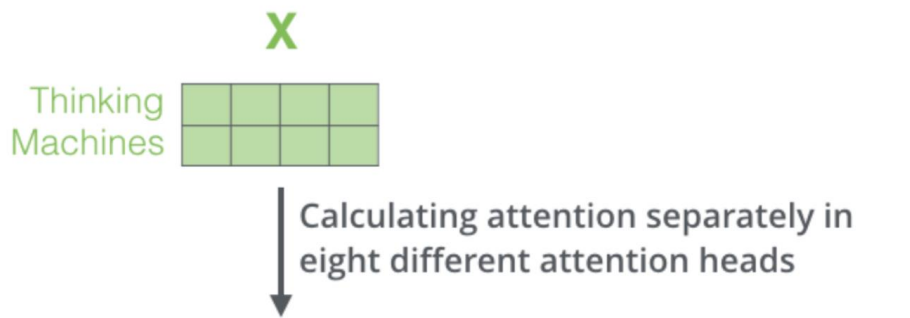
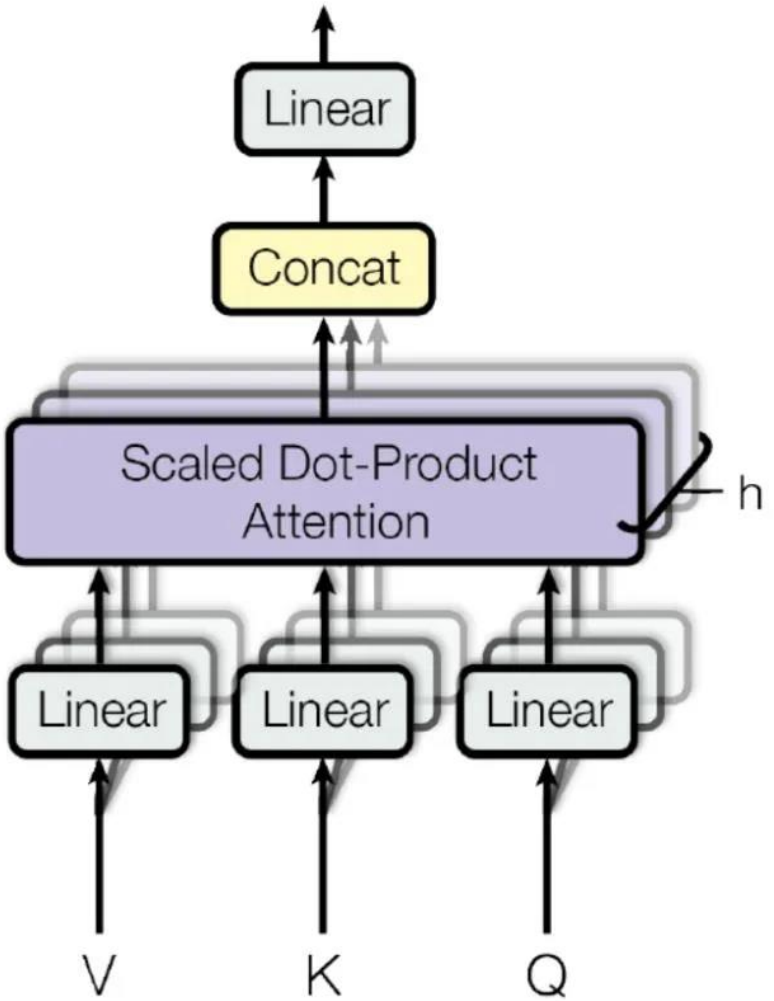
Scaled Dot-Product Attention



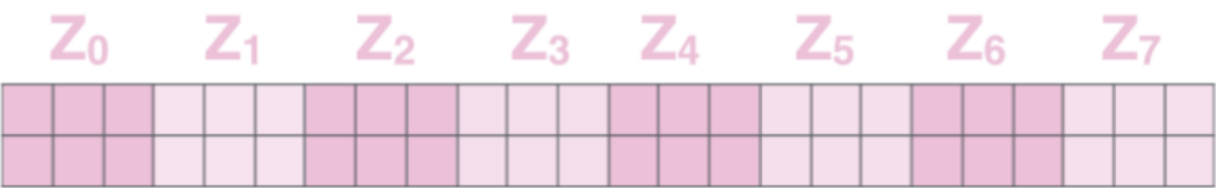
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Multi-Head Self-Attention

Multi-Head Attention

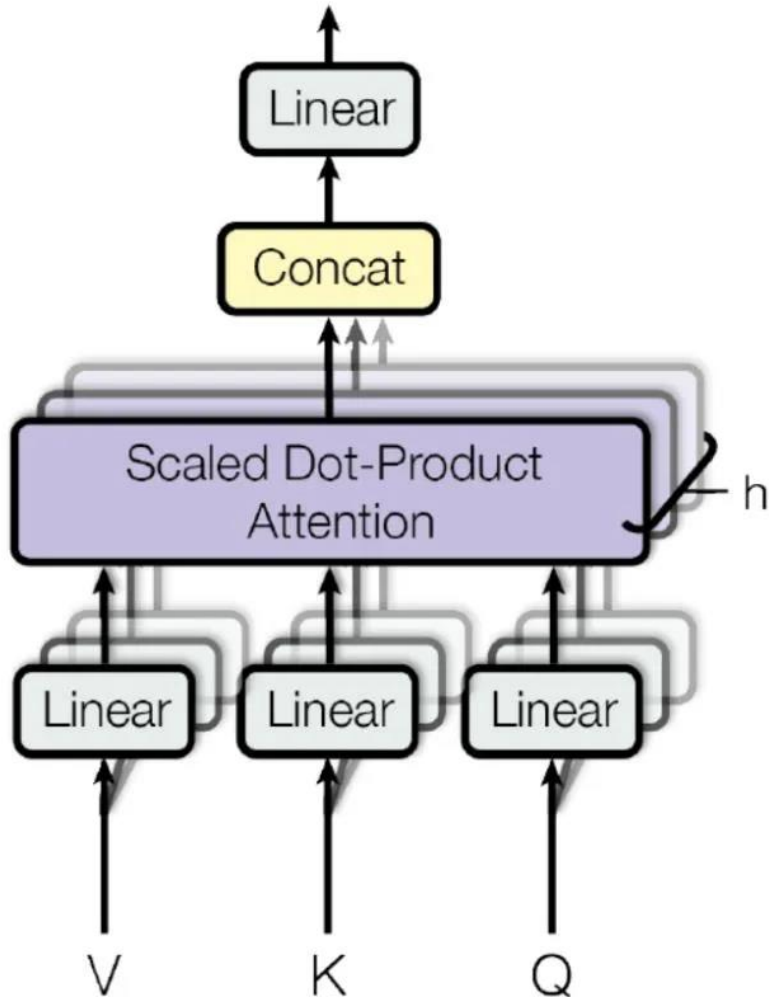


1) Concatenate all the attention heads



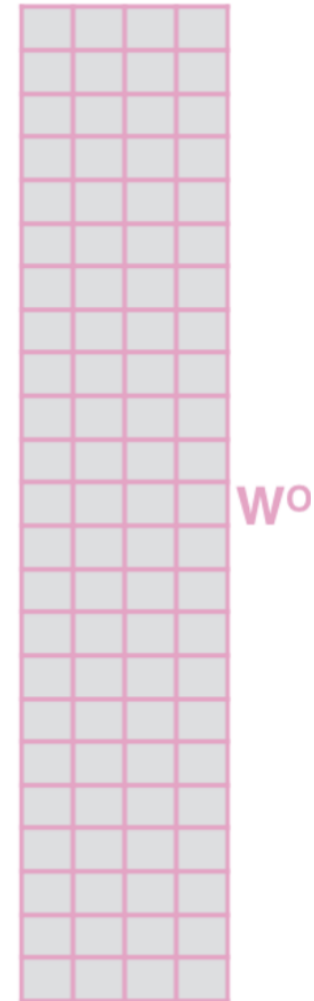
Multi-Head Self-Attention

Multi-Head Attention



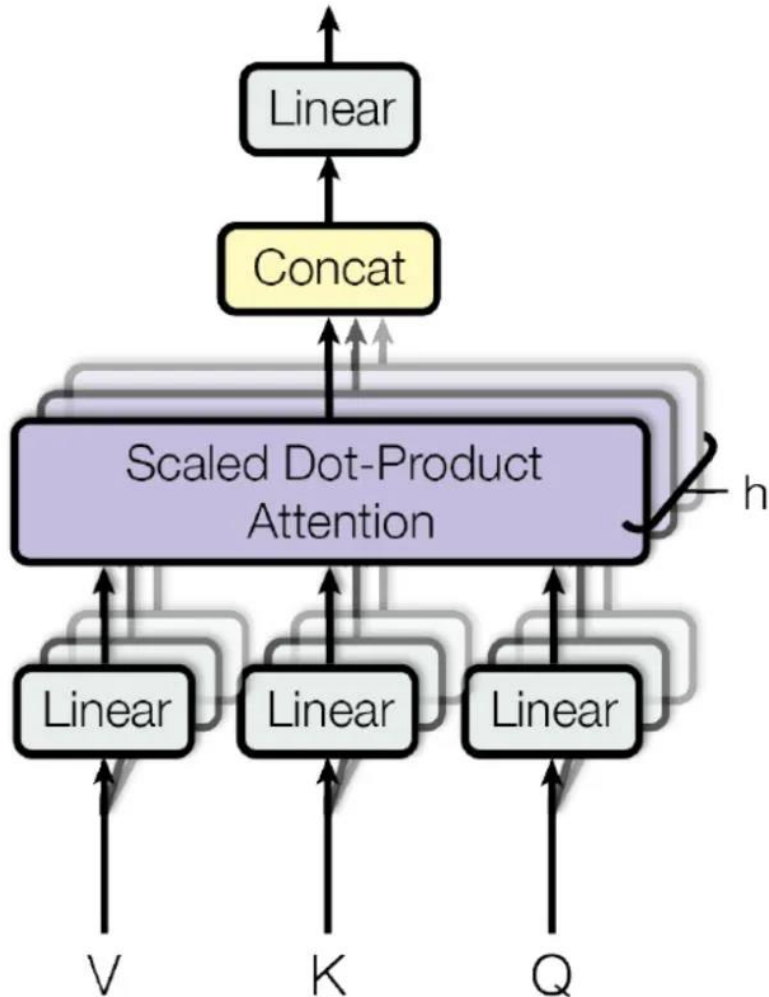
2) Multiply with a weight matrix W^O that was trained jointly with the model

\times



Multi-Head Self-Attention

Multi-Head Attention

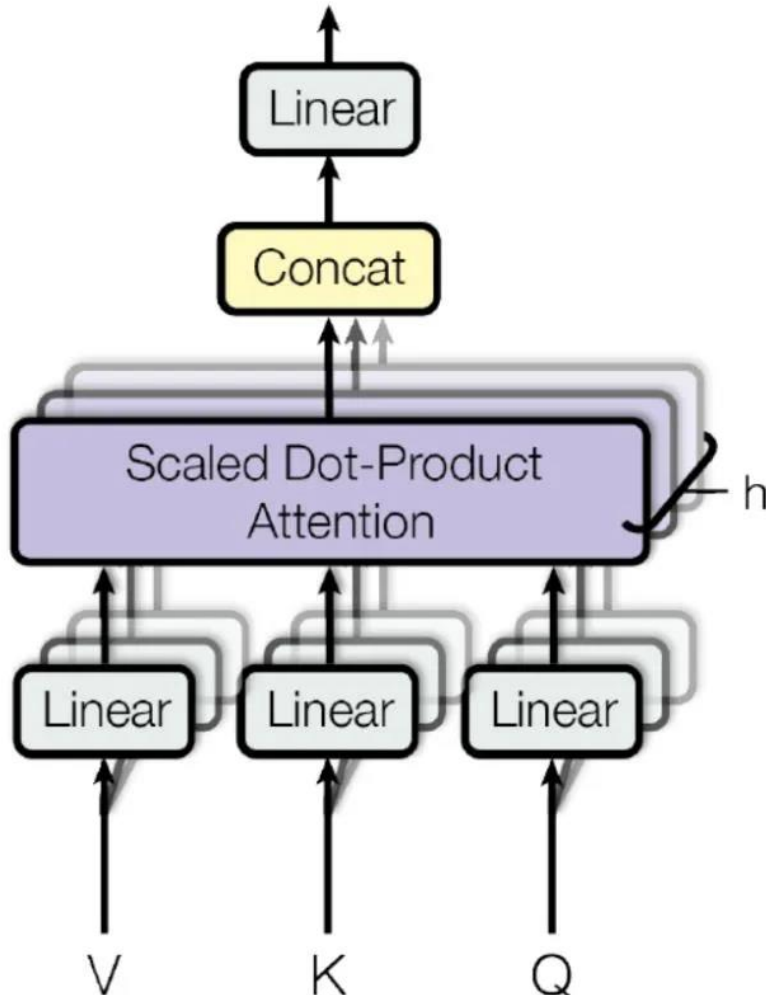


3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



Multi-Head Self-Attention

Multi-Head Attention



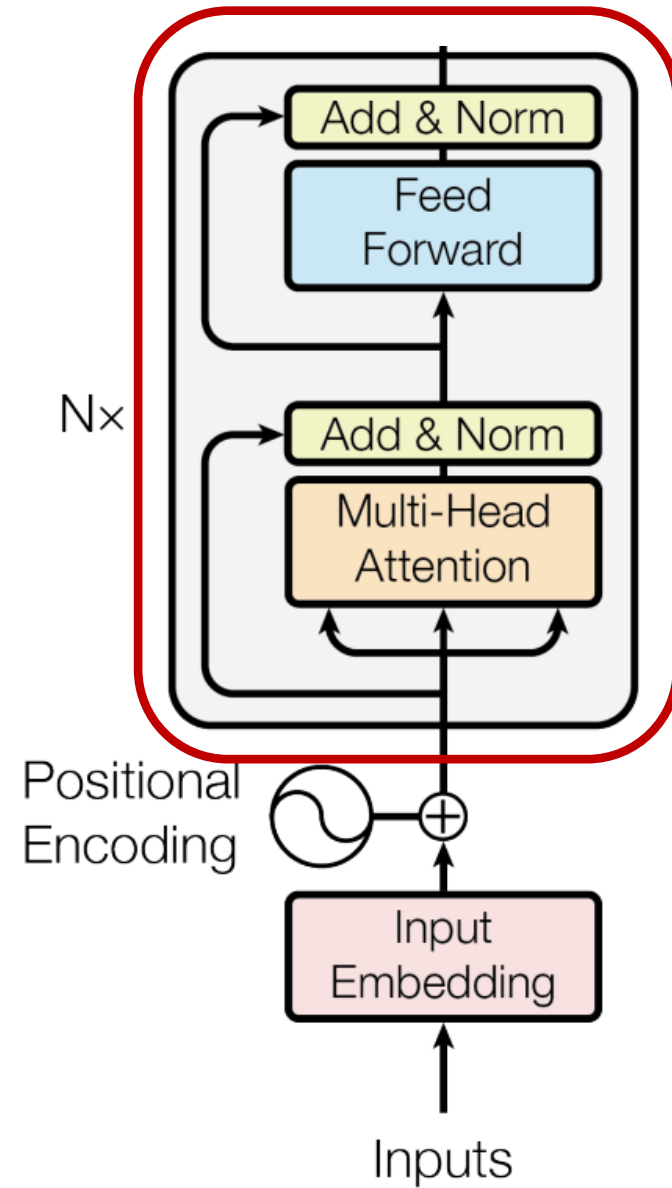
3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



‘Multi-head attention allows the model to jointly attend to information from different representation subspaces at different positions. With a single attention head, averaging inhibits this.’

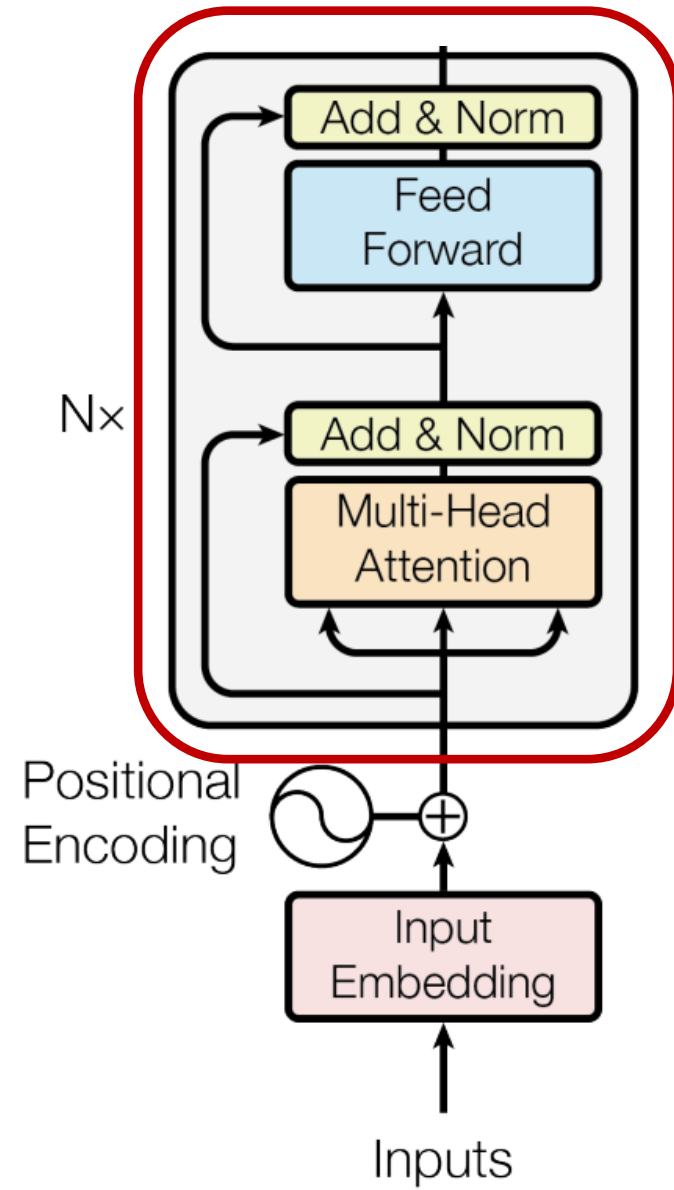
Transformer

- Encoder:
 - Multi-Head Self-Attention
- Feed Forward
 - Residual Connection & Layer Norm



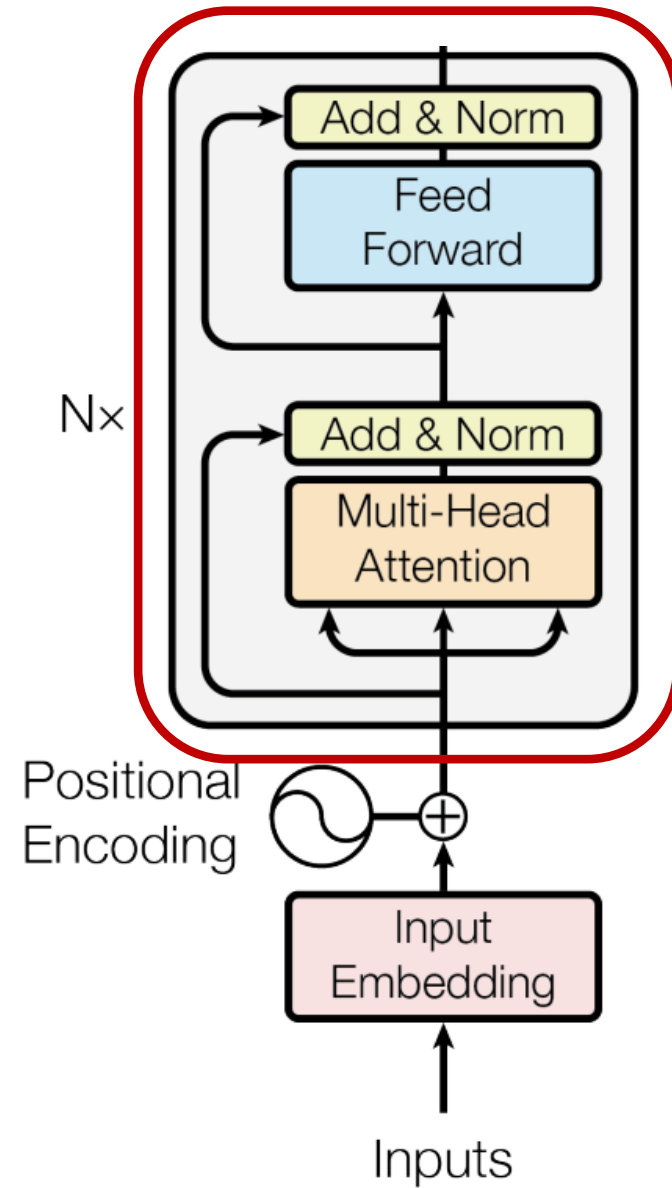
Transformer

$$\max(0, XW_1 + b_1)W_2 + b_2$$



Transformer

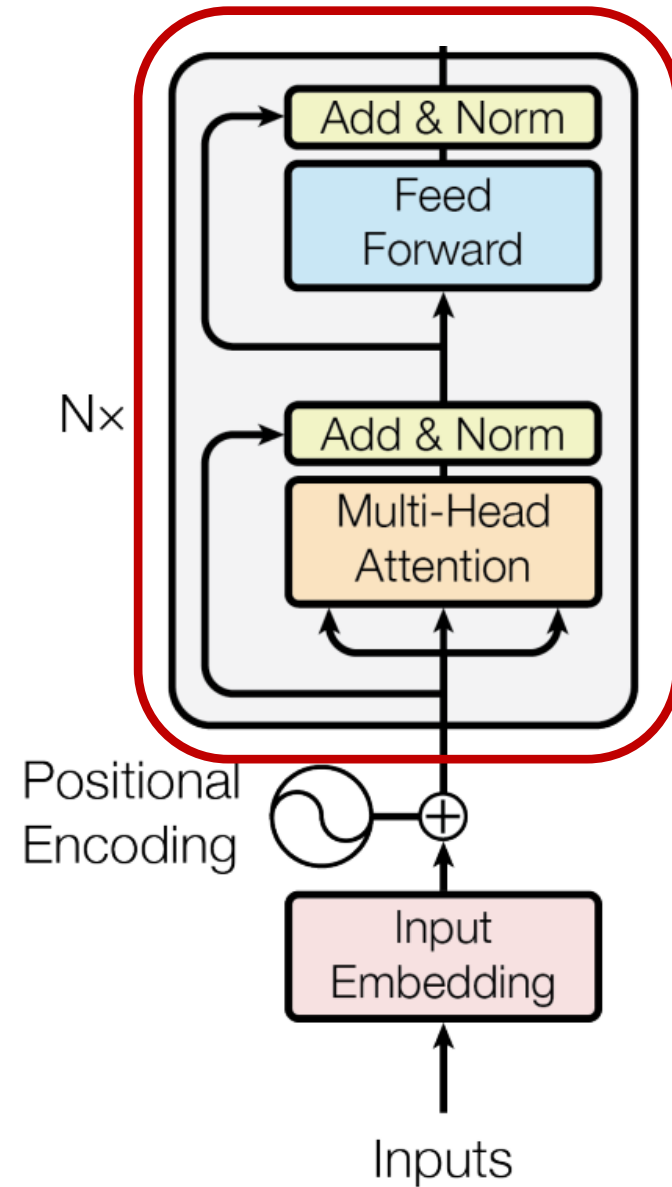
- Encoder:
- Multi-Head Self-Attention
- Feed Forward
- Residual Connection & Layer Norm



Add & Norm

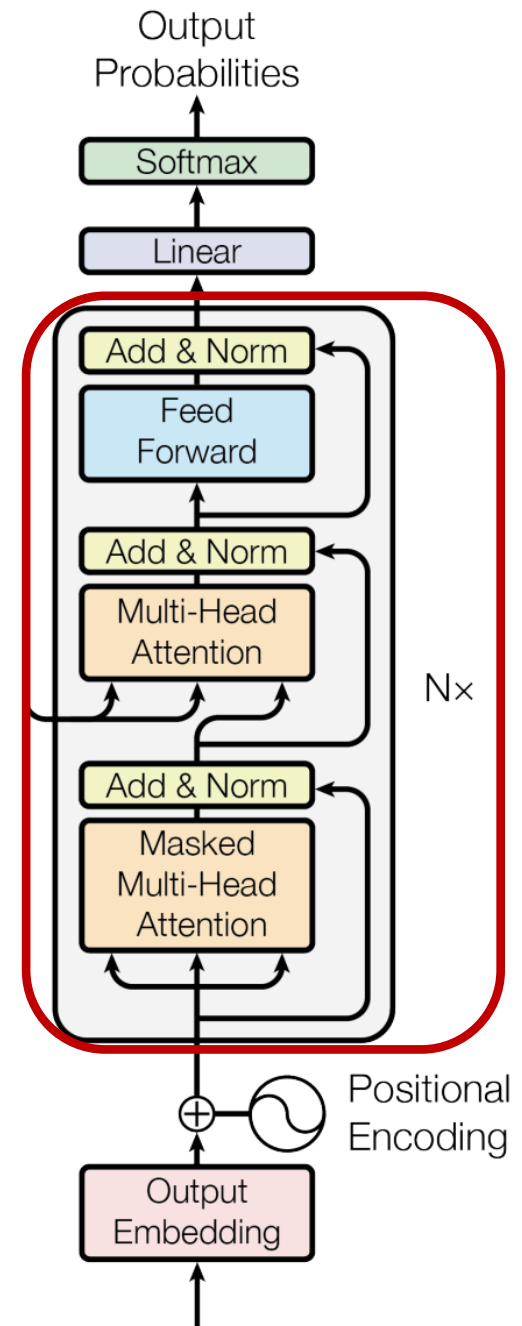
$\text{LayerNorm}(X + \text{MultiHeadAttention}(X))$

$\text{LayerNorm}(X + \text{FeedForward}(X))$



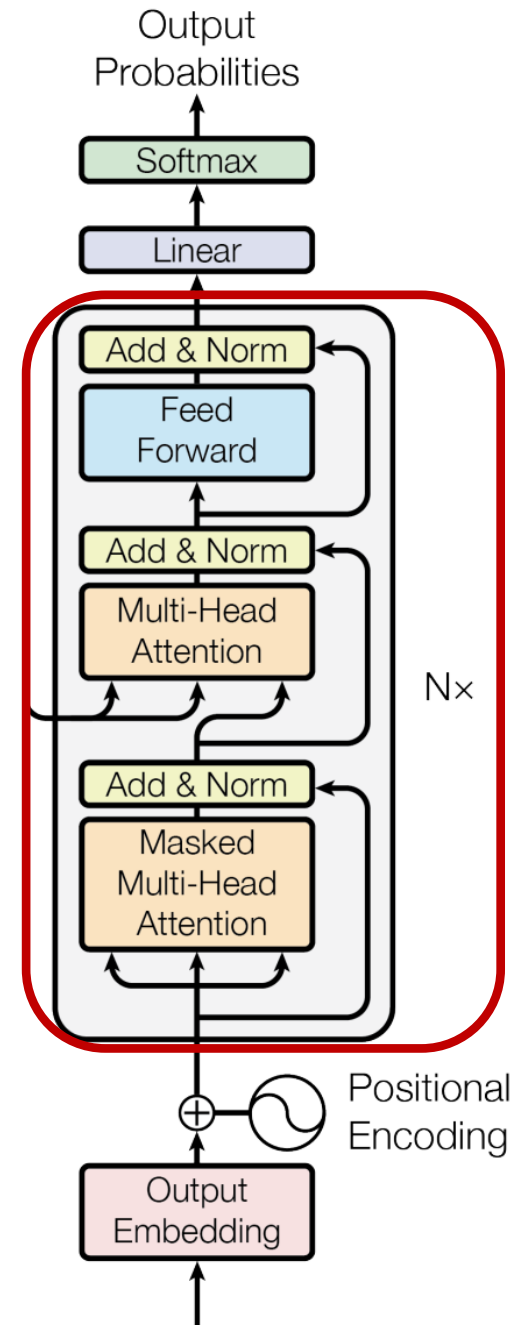
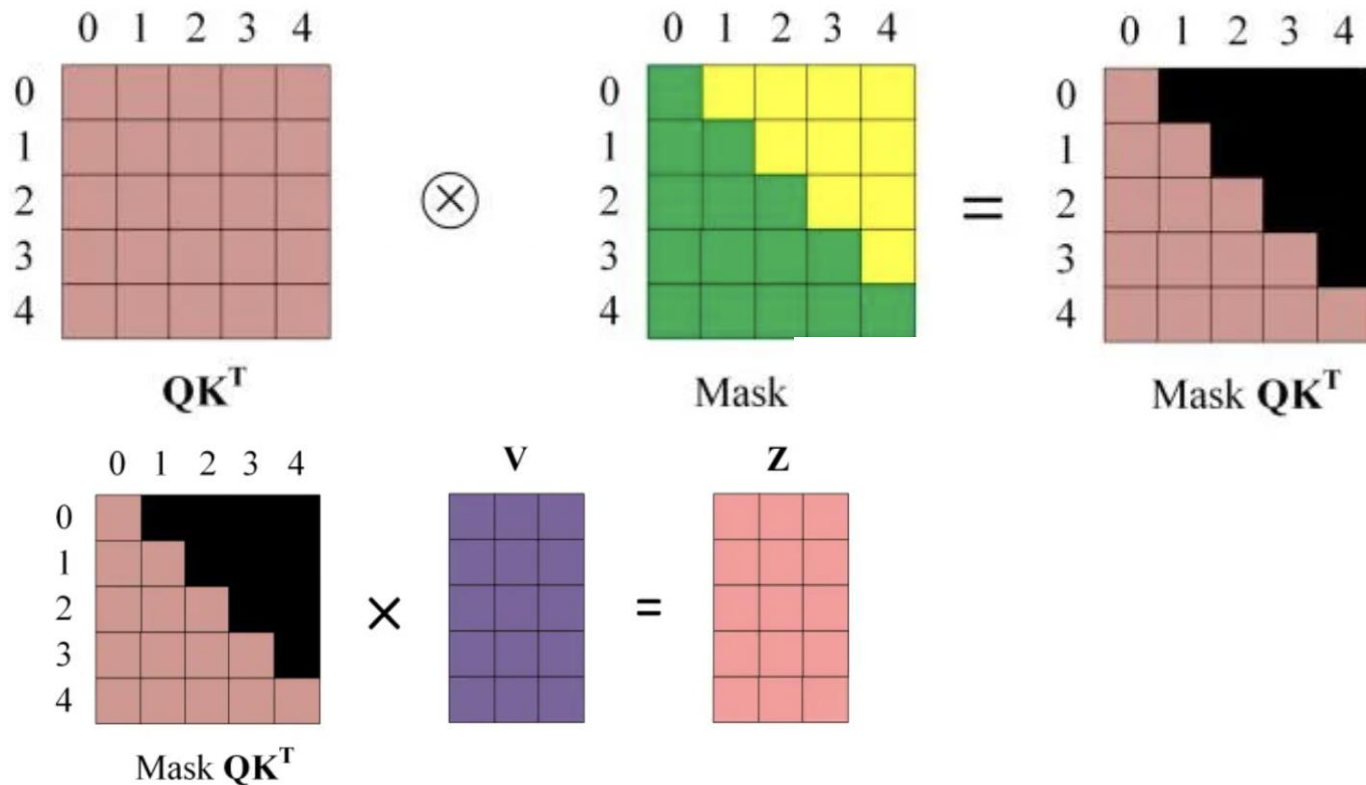
Transformer

- Decoder:
- 2 x Multi-Head Self-Attention
- 1st Masked Multi-Head Self-Attention
- 2nd Q K V
- Softmax



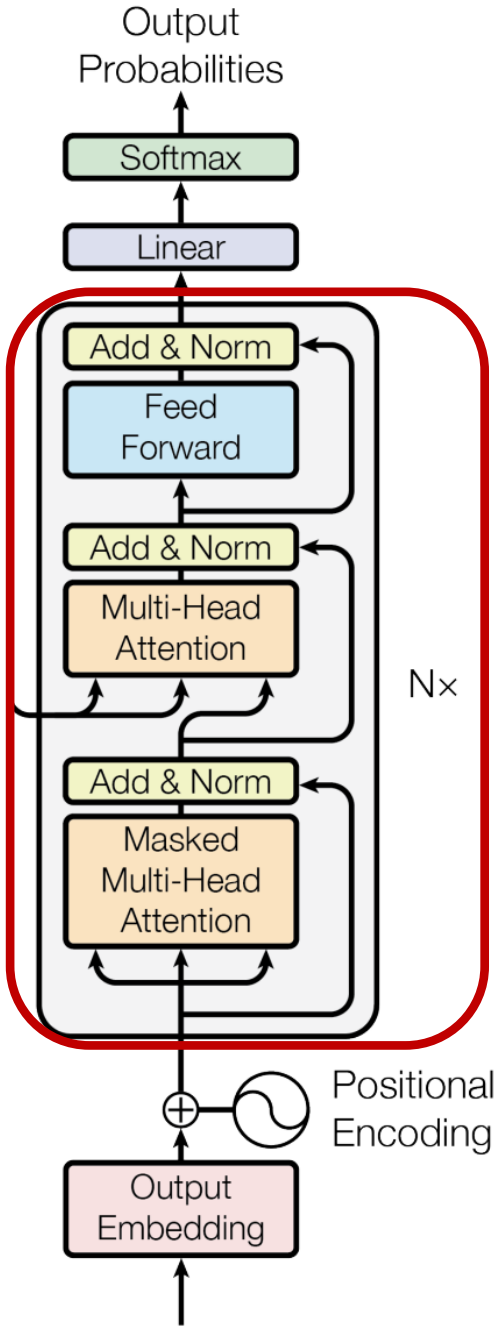
Transformer

- Masked Multi-Head Self-Attention



Transformer

- Softmax

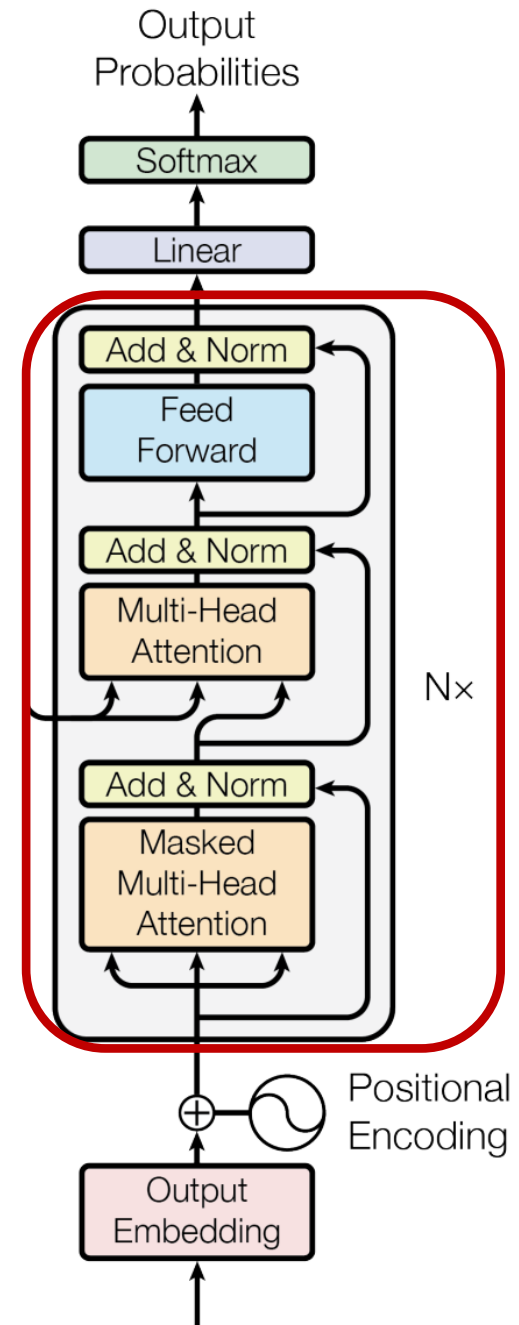


Transformer

- Positional Encoding

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10000^{2k/d}}$$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$



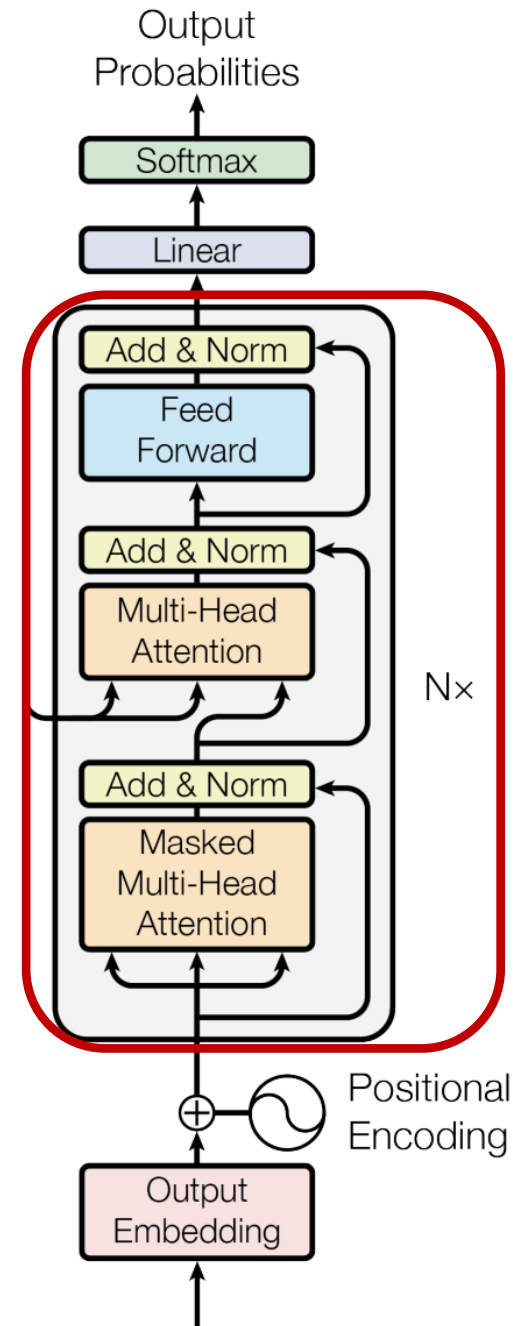
Transformer

- Positional Encoding

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10000^{2k/d}}$$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

$$M \cdot \begin{bmatrix} \sin(\omega_k \cdot t) \\ \cos(\omega_k \cdot t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k \cdot (t + \phi)) \\ \cos(\omega_k \cdot (t + \phi)) \end{bmatrix}$$



Advantages of Transformer

- Enhanced Parallelization Capabilities
- Capturing Long-Distance Dependencies
- Dynamic Weight Allocation

Disadvantages of Transformer

- Computational Efficiency Issues
 - Quadratic Time Complexity
 - High Memory Consumption
- Limited Capability with Long Sequences
 - Limited Effective Resolution Window
 - Extended Training Times
- Overparameterization

SSMs—S4

Efficiently Modeling Long Sequences with Structured State Spaces

Albert Gu, Karan Goel, and Christopher Ré

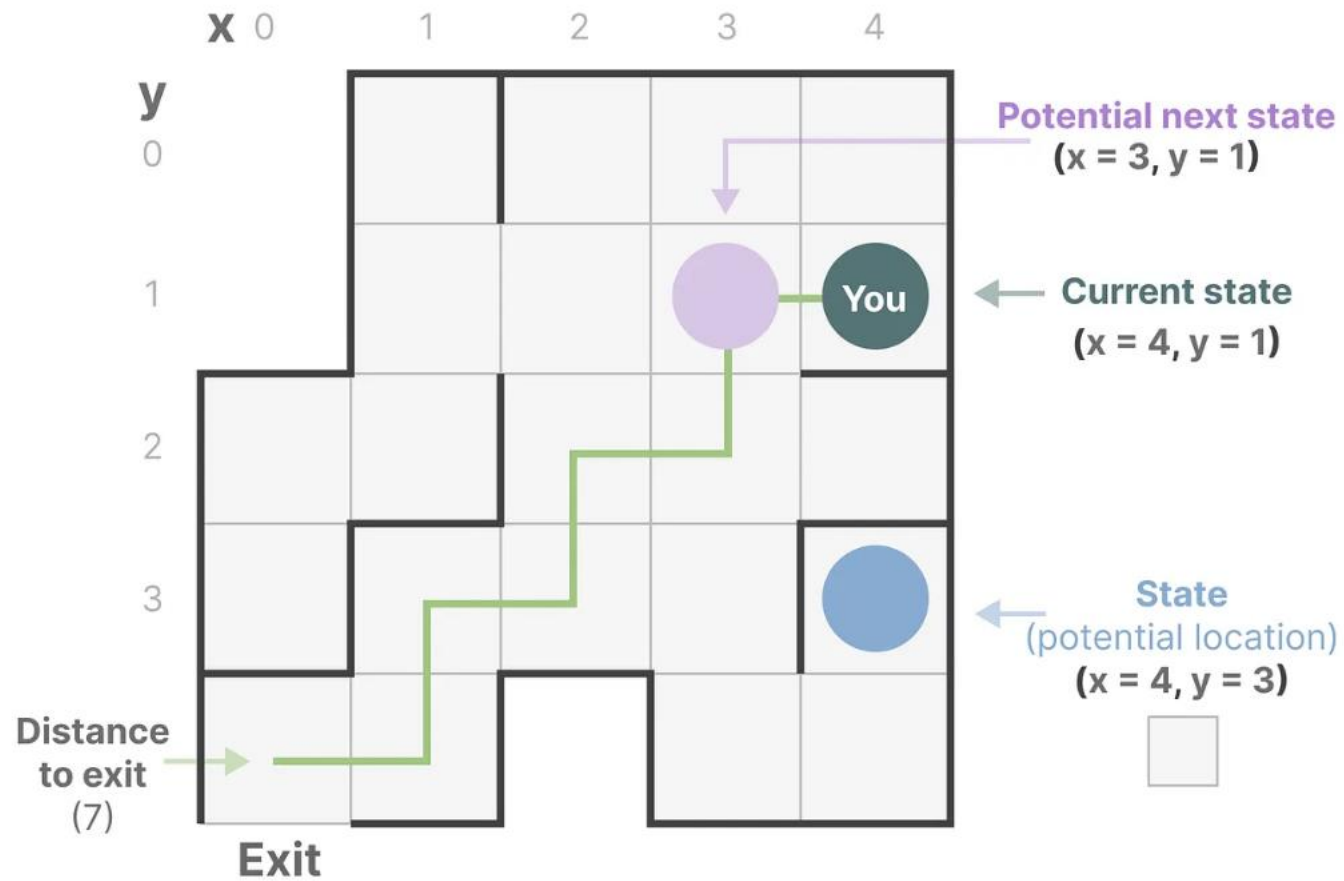
Department of Computer Science, Stanford University

{albertgu,krng}@stanford.edu, chrismre@cs.stanford.edu

Maarten Grootendorst «A Visual Guide to Mamba and State Space Models»

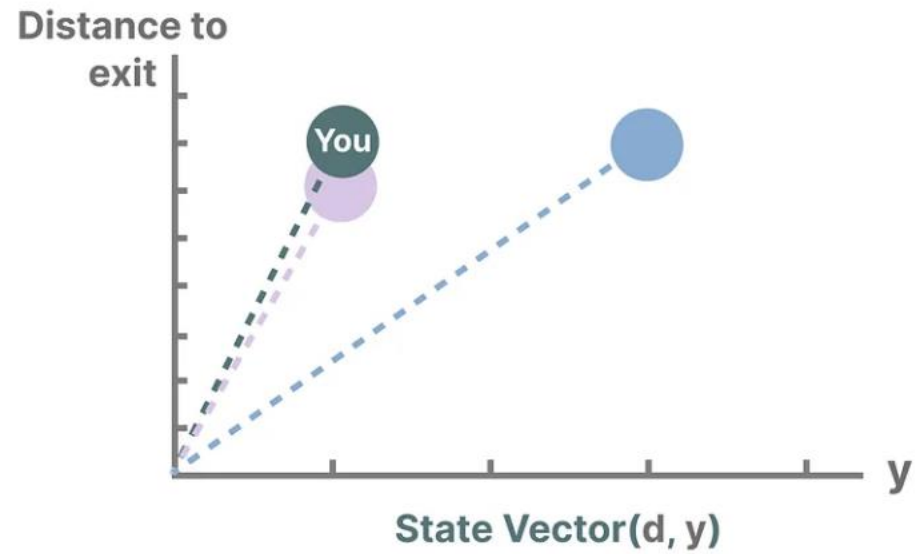
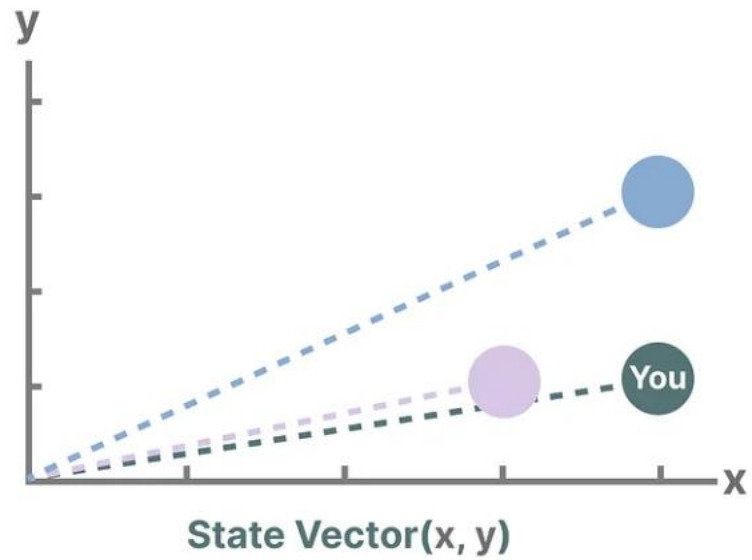
SSMs — — S4

- State Space



SSMs—S4

- State Space



SSMs — — S4

- SSM

Input
(sequence)



$x(t)$

**State Space Model
(SSM)**

Output
(sequence)



$y(t)$

SSMs — — S4

- SSM

Input
(sequence)



$x(t)$

SSM

state equation

$$\mathbf{h}'(t) = \mathbf{A}\mathbf{h}(t) + \mathbf{B}\mathbf{x}(t)$$

output equation

$$\mathbf{y}(t) = \mathbf{C}\mathbf{h}(t) + \mathbf{D}\mathbf{x}(t)$$

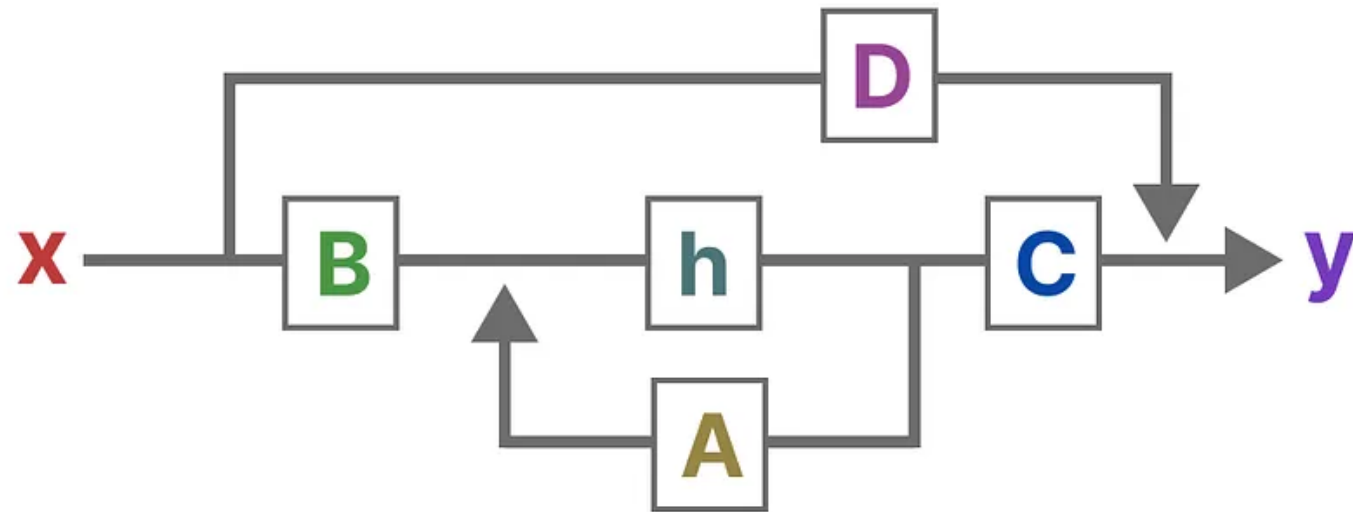
Output
(sequence)



$y(t)$

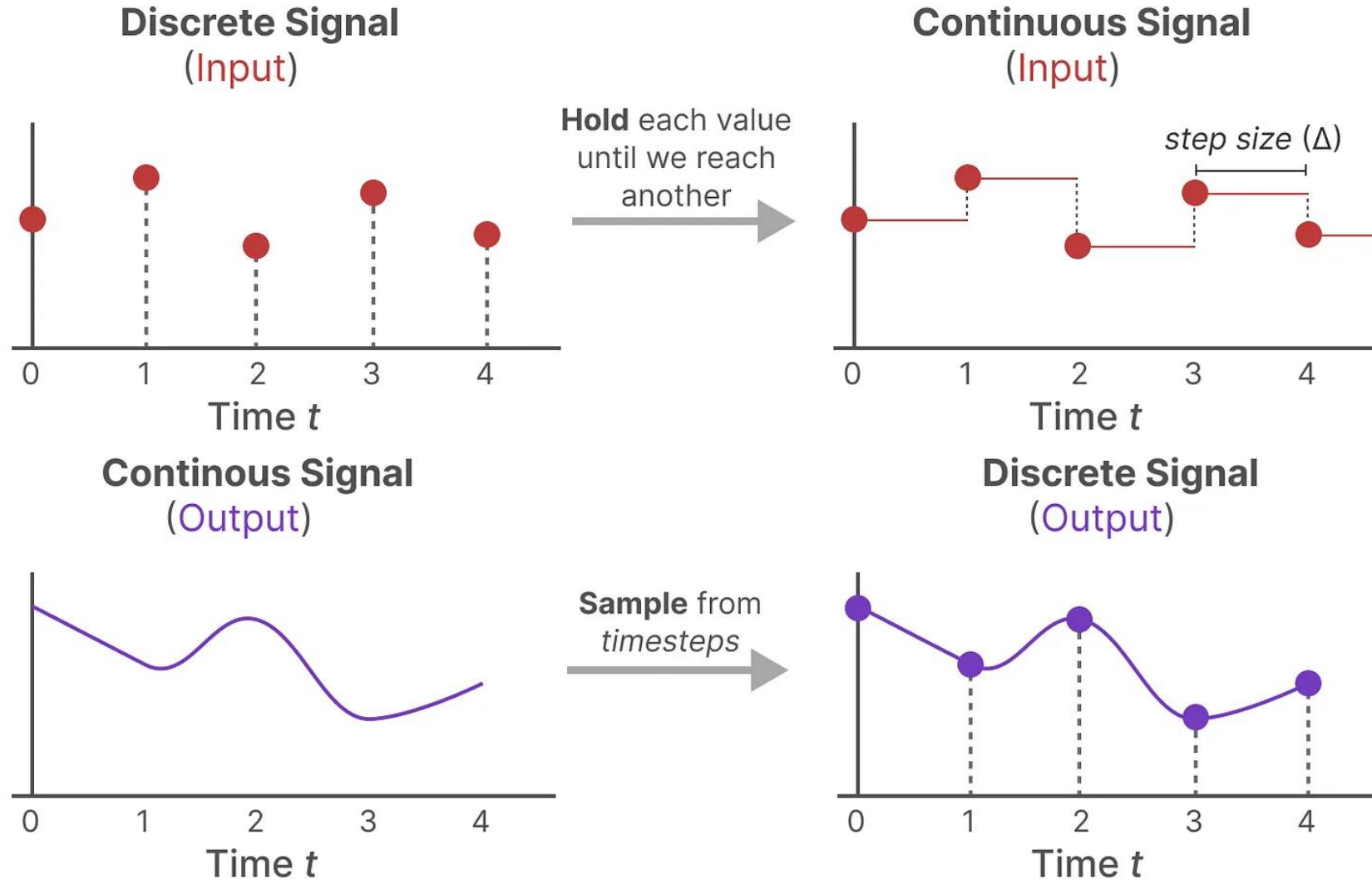
SSMs — S4

- SSM



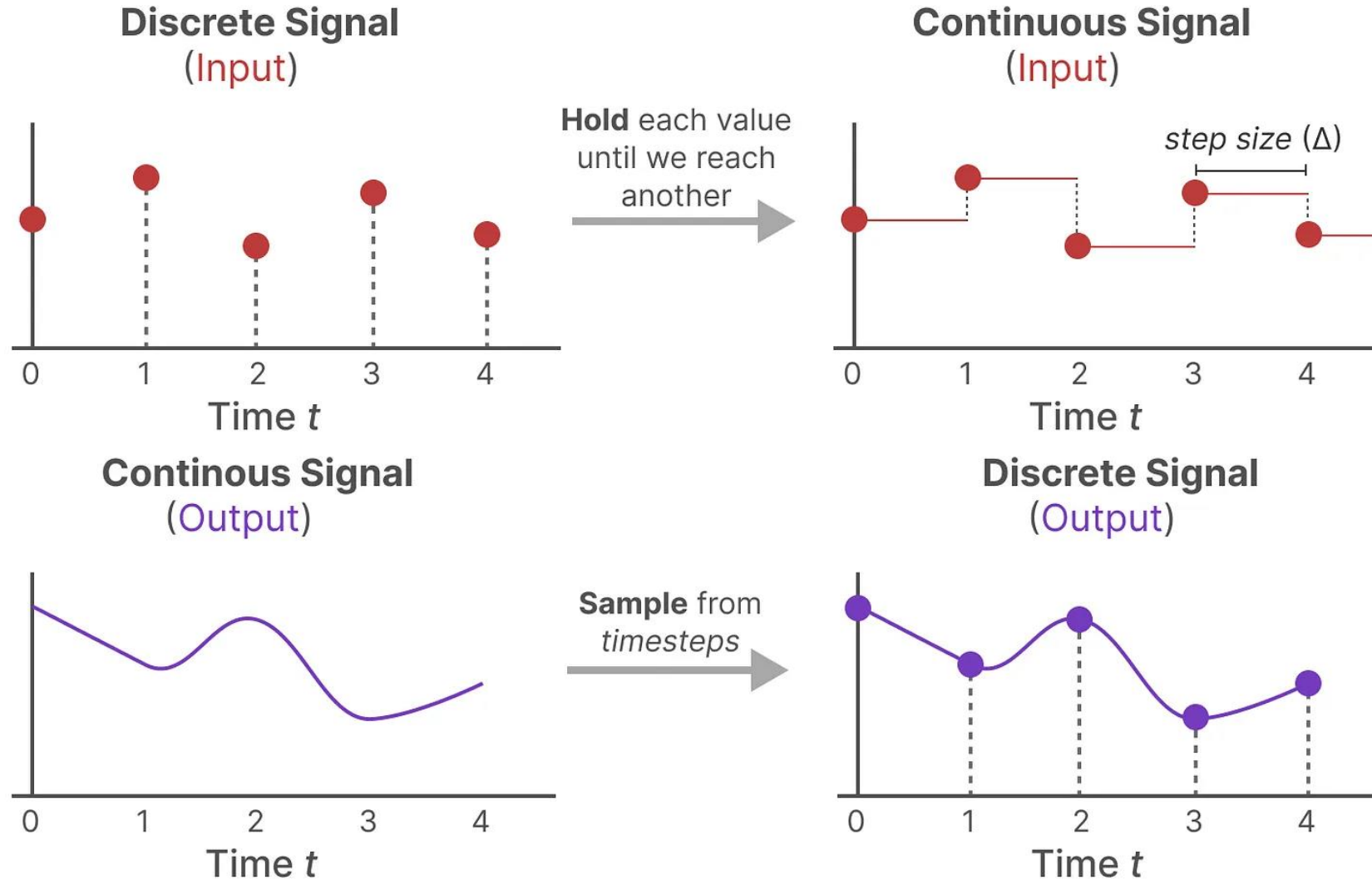
SSMs—S4

- S4



SSMs—S4

- S4



SSMs — S4

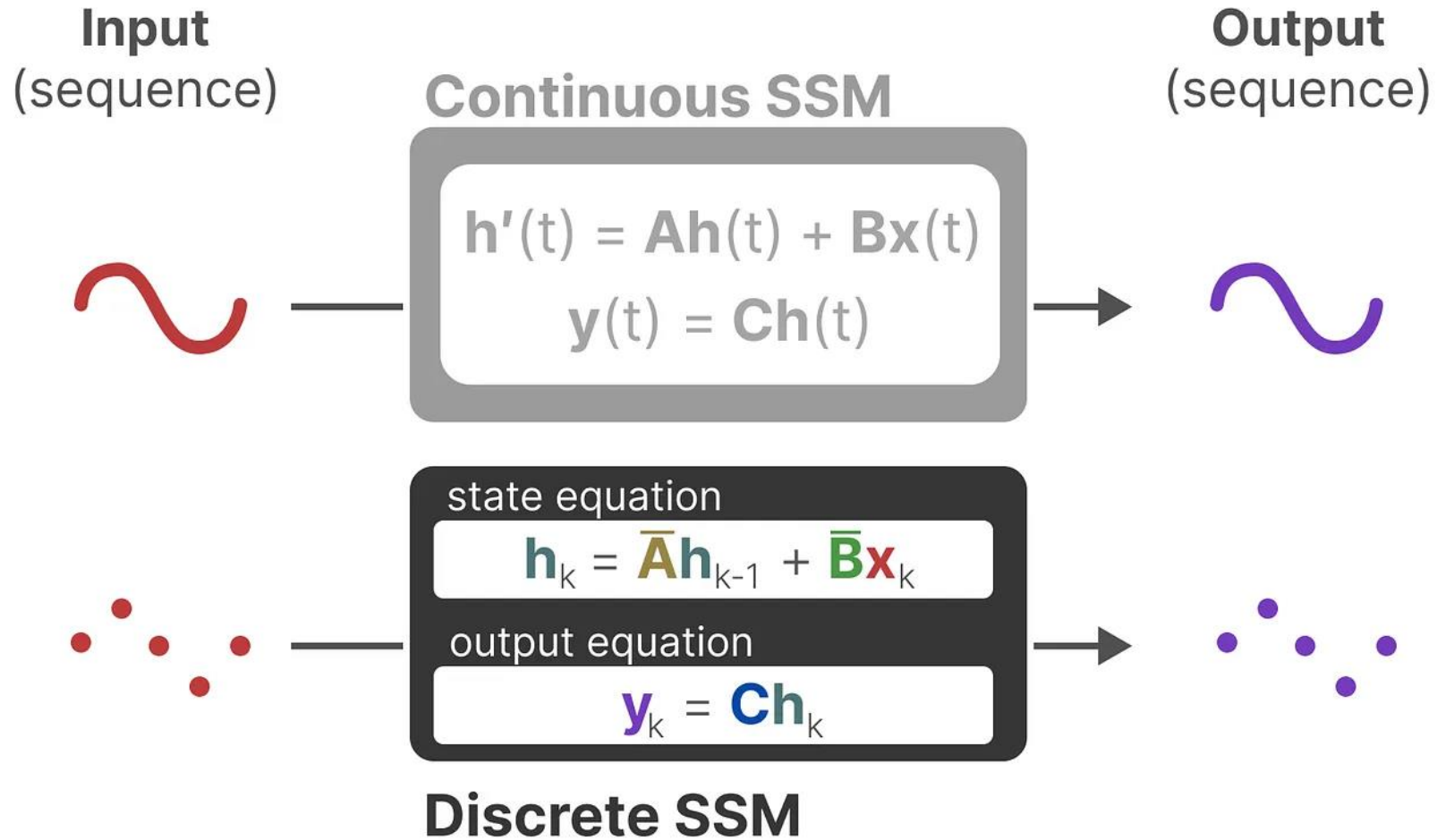
- S4

Discretized matrix **A**

$$\bar{\mathbf{A}} = \exp(\Delta\mathbf{A})$$

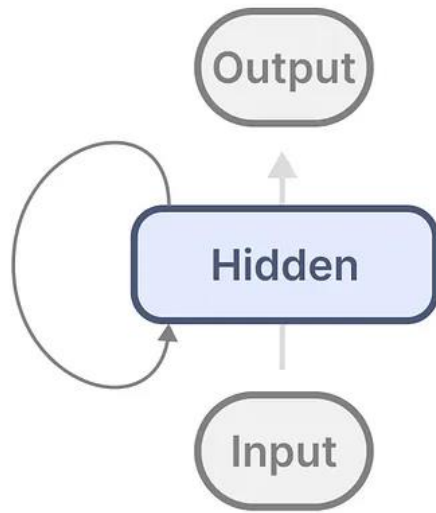
Discretized matrix **B**

$$\bar{\mathbf{B}} = (\Delta\mathbf{A})^{-1} (\exp(\Delta\mathbf{A}) - \mathbf{I}) \cdot \Delta\mathbf{B}$$

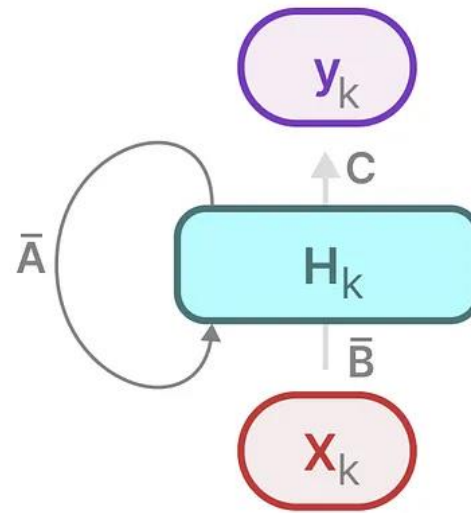


SSMs — — S4

- S4



RNN



SSM
(Recurrent)

SSMs—S4

- S4

Timestep 0

$$\mathbf{h}_0 = \bar{\mathbf{B}}\mathbf{x}_0$$

$$\mathbf{y}_0 = \mathbf{C}\mathbf{h}_0$$

Timestep -1
does not exist so

$\mathbf{A}\mathbf{h}_{-1}$
can be ignored



Timestep 1

$$\mathbf{h}_1 = \bar{\mathbf{A}}\mathbf{h}_0 + \bar{\mathbf{B}}\mathbf{x}_1$$

$$\mathbf{y}_1 = \mathbf{C}\mathbf{h}_1$$

State of
previous timestep

State of
current timestep



Timestep 2

$$\mathbf{h}_2 = \bar{\mathbf{A}}\mathbf{h}_1 + \bar{\mathbf{B}}\mathbf{x}_2$$

$$\mathbf{y}_2 = \mathbf{C}\mathbf{h}_2$$

State of
previous timestep

State of
current timestep

$$y_2 = \mathbf{C}\mathbf{h}_2$$

$$= \mathbf{C} (\bar{\mathbf{A}}\mathbf{h}_1 + \bar{\mathbf{B}}\mathbf{x}_2)$$

$$= \mathbf{C} (\bar{\mathbf{A}} (\bar{\mathbf{A}}\mathbf{h}_0 + \bar{\mathbf{B}}\mathbf{x}_1) + \bar{\mathbf{B}}\mathbf{x}_2)$$

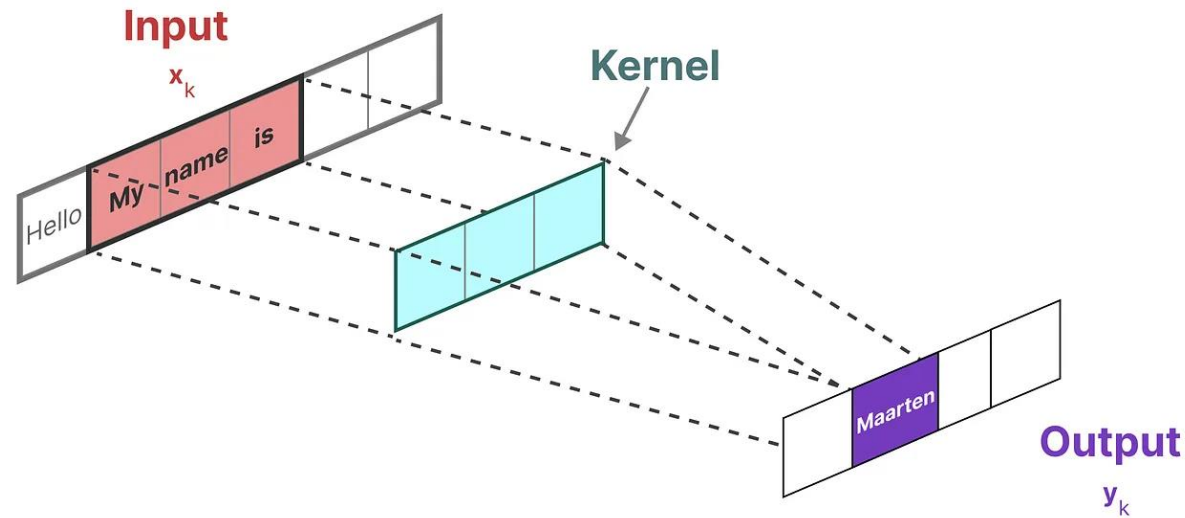
$$= \mathbf{C} (\bar{\mathbf{A}} (\bar{\mathbf{A}} \cdot \bar{\mathbf{B}}\mathbf{x}_0 + \bar{\mathbf{B}}\mathbf{x}_1) + \bar{\mathbf{B}}\mathbf{x}_2)$$

$$= \mathbf{C} (\bar{\mathbf{A}} \cdot \bar{\mathbf{A}} \cdot \bar{\mathbf{B}}\mathbf{x}_0 + \bar{\mathbf{A}} \cdot \bar{\mathbf{B}}\mathbf{x}_1 + \bar{\mathbf{B}}\mathbf{x}_2)$$

$$= \mathbf{C} \cdot \bar{\mathbf{A}}^2 \cdot \bar{\mathbf{B}}\mathbf{x}_0 + \mathbf{C} \cdot \bar{\mathbf{A}} \cdot \bar{\mathbf{B}} \cdot \mathbf{x}_1 + \mathbf{C} \cdot \bar{\mathbf{B}}\mathbf{x}_2$$

SSMs — S4

- S4



$$\text{kernel} \rightarrow \bar{\mathbf{K}} = (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^k\bar{\mathbf{B}}, \dots)$$

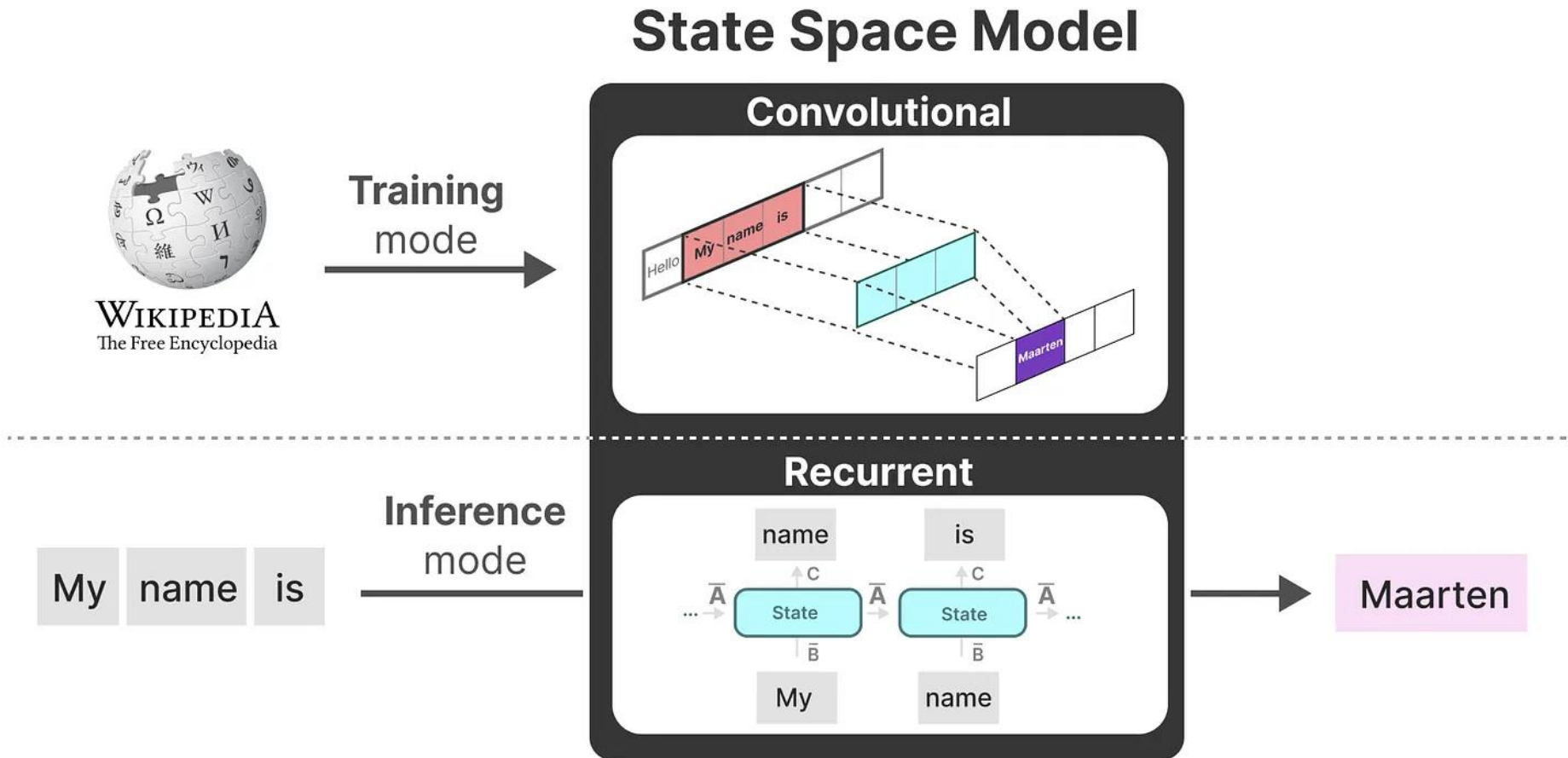
$$\mathbf{y} = \mathbf{x} * \bar{\mathbf{K}}$$

output input kernel

$$\begin{aligned} y_2 &= Ch_2 \\ &= C(\bar{A}h_1 + \bar{B}x_2) \\ &= C(\bar{A}(\bar{A}h_0 + \bar{B}x_1) + \bar{B}x_2) \\ &= C(\bar{A}(\bar{A} \cdot \bar{B}x_0 + \bar{B}x_1) + \bar{B}x_2) \\ &= C(\bar{A} \cdot \bar{A} \cdot \bar{B}x_0 + \bar{A} \cdot \bar{B}x_1 + \bar{B}x_2) \\ &= C \cdot \bar{A}^2 \cdot \bar{B}x_0 + C \cdot \bar{A} \cdot \bar{B} \cdot x_1 + C \cdot \bar{B}x_2 \end{aligned}$$

SSMs — — S4

- S4



SSMs — — S4

HiPPO: Recurrent Memory with Optimal Polynomial Projections

Albert Gu^{*†}, Tri Dao^{*†}, Stefano Ermon[†], Atri Rudra[‡], Christopher Ré[†]

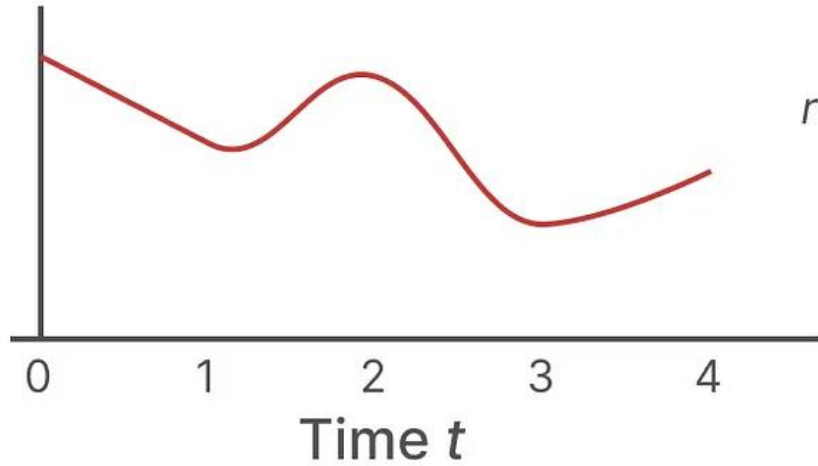
[†] Department of Computer Science, Stanford University

[‡] Department of Computer Science and Engineering, University at Buffalo, SUNY

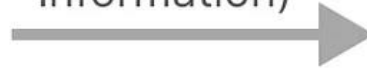
{albertgu, trid}@stanford.edu, ermon@cs.stanford.edu, atri@buffalo.edu, chrismre@cs.stanford.edu

SSMs—S4

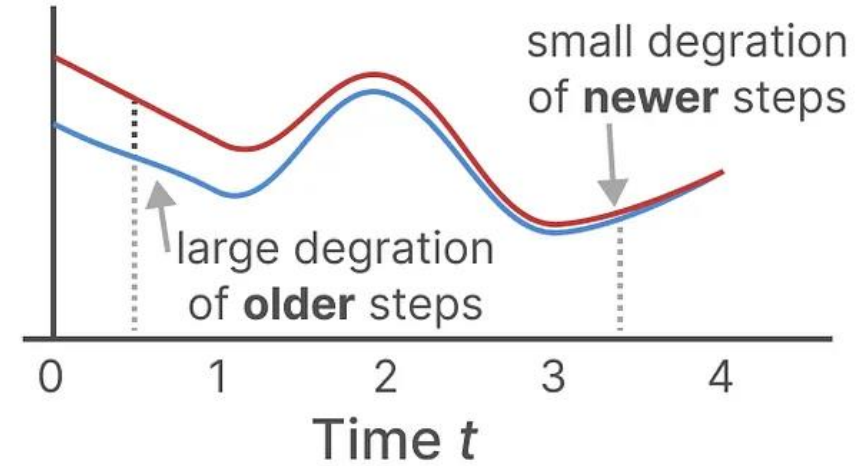
Input Signal



HiPPO
(compress and
reconstruct signal
information)



Reconstructed Signal



SSMs—S4

HiPPO Matrix \mathbf{A}_{nk} $\left\{ \begin{array}{l} (2n + 1)^{1/2} (2k + 1)^{1/2} \leftarrow \text{everything **below** the diagonal} \\ n + 1 \leftarrow \text{the diagonal} \\ 0 \leftarrow \text{everything **above** the diagonal} \end{array} \right.$

HiPPO Matrix

1	0	0	0
1	2	0	0
1	3	3	0
1	3	5	4

n

k

SSMs — — S4

Structured State Spaces for Sequences (S4)

||

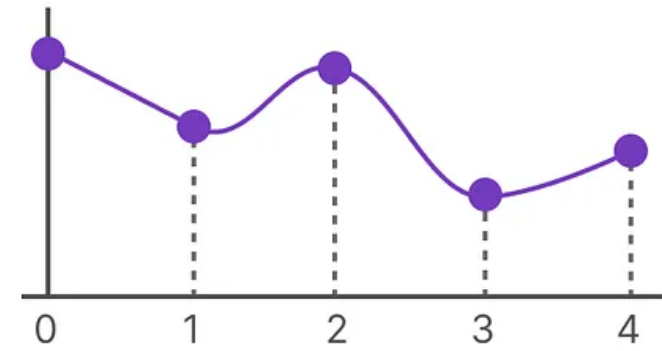
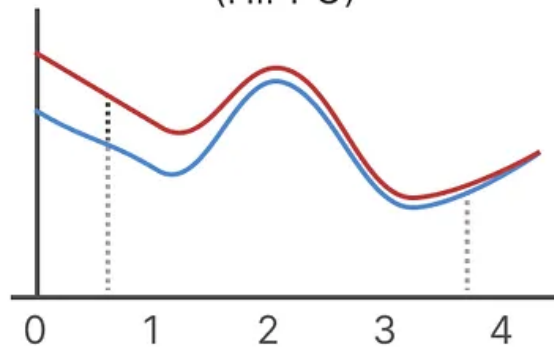
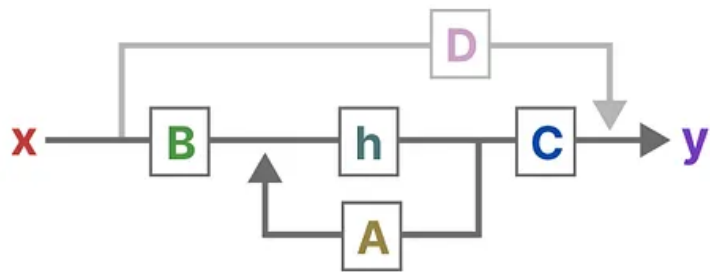
Continuous State Space

+

Long-Range Dependencies (HiPPO)

+

Discrete Representations



Training mode (convolutional)

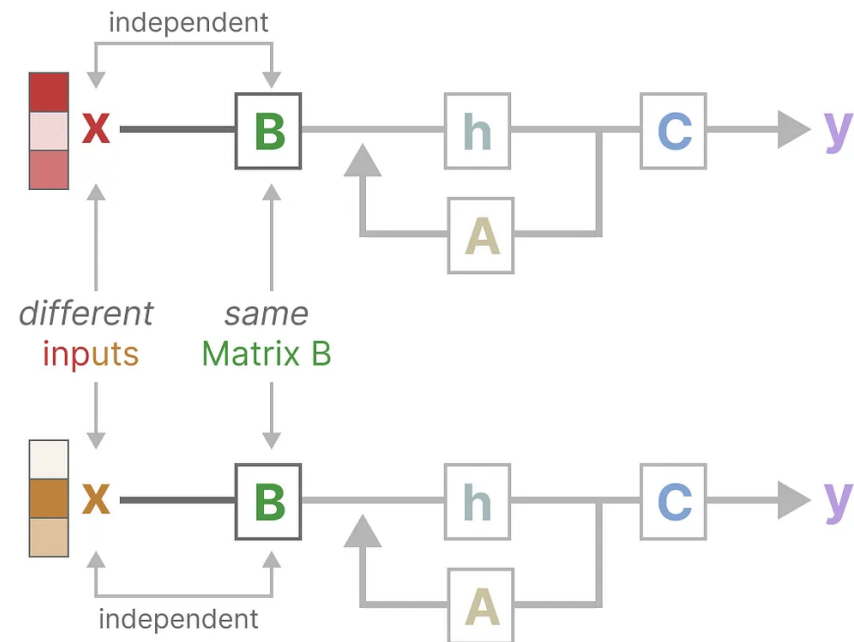
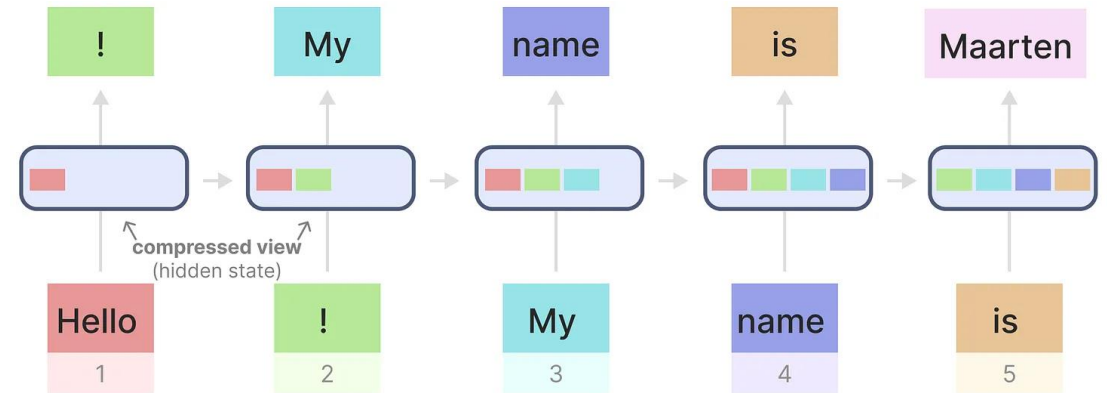
Inference mode (recurrence)

Mamba— —S6

- Selection Mechanism
- Hardware-aware Algorithm
- Simpler SSM Architecture

Disadvantages of Previous Works

- Transformer—long context
- RNN—forget far context
- S4—fixed A, B, C



Mamba—S6

- Selection Mechanism

Algorithm 1 SSM (S4)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

1: $\mathbf{A} : (D, N) \leftarrow \text{Parameter}$

▷ Represents structured $N \times N$ matrix

2: $\mathbf{B} : (D, N) \leftarrow \text{Parameter}$

3: $\mathbf{C} : (D, N) \leftarrow \text{Parameter}$

4: $\Delta : (D) \leftarrow \tau_{\Delta}(\text{Parameter})$

5: $\overline{\mathbf{A}}, \overline{\mathbf{B}} : (D, N) \leftarrow \text{discretize}(\Delta, \mathbf{A}, \mathbf{B})$

6: $y \leftarrow \text{SSM}(\overline{\mathbf{A}}, \overline{\mathbf{B}}, \mathbf{C})(x)$

▷ Time-invariant: recurrence or convolution

7: **return** y

Algorithm 2 SSM + Selection (S6)

Input: $x : (B, L, D)$

Output: $y : (B, L, D)$

1: $\mathbf{A} : (D, N) \leftarrow \text{Parameter}$

▷ Represents structured $N \times N$ matrix

2: $\mathbf{B} : (B, L, N) \leftarrow s_B(x)$

3: $\mathbf{C} : (B, L, N) \leftarrow s_C(x)$

4: $\Delta : (B, L, D) \leftarrow \tau_{\Delta}(\text{Parameter} + s_{\Delta}(x))$

5: $\overline{\mathbf{A}}, \overline{\mathbf{B}} : (B, L, D, N) \leftarrow \text{discretize}(\Delta, \mathbf{A}, \mathbf{B})$

6: $y \leftarrow \text{SSM}(\overline{\mathbf{A}}, \overline{\mathbf{B}}, \mathbf{C})(x)$

▷ Time-varying: recurrence (*scan*) only

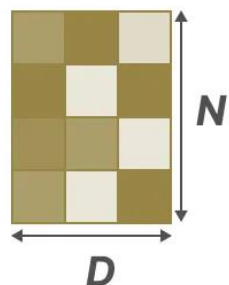
7: **return** y

Mamba — — S6

Matrix A

How the **current state** evolves over time

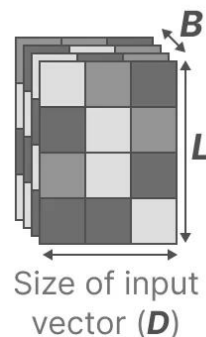
Structured State Space Model (S4)



Step size (Δ)

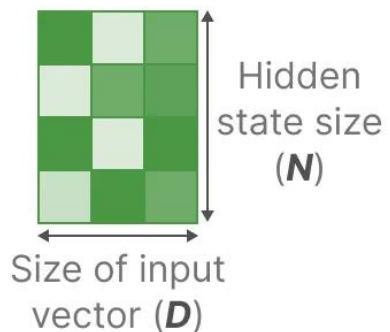
Resolution of the **input** (discretization parameter)

SSM + Selection



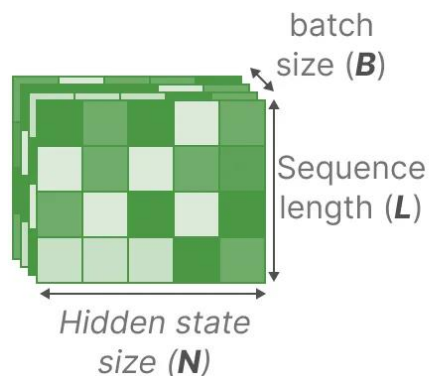
Matrix B

How the **input** influences the state



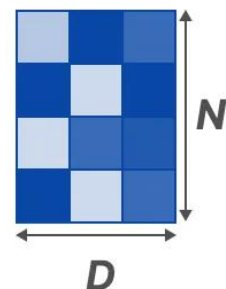
Matrix B

How the **input** influences the state



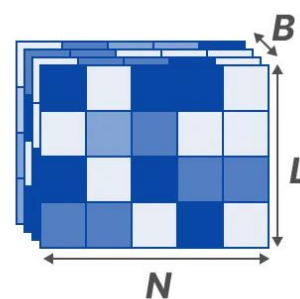
Matrix C

How the **current state** translates to the **output**



Matrix C

How the **current state** translates to the **output**



$$s_B(x) = \text{Linear}_N(x)$$

$$s_C(x) = \text{Linear}_N(x)$$

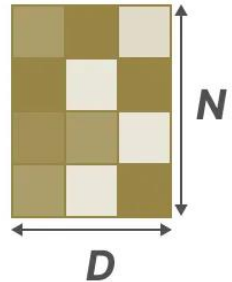
$$s_\Delta(x) = \text{Linear}_D(x)$$

$$\tau_\Delta = \text{softplus}$$

Mamba — — S6

Matrix A

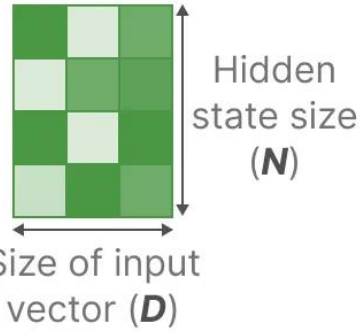
How the **current state** evolves over time



Structured State Space Model (S4)

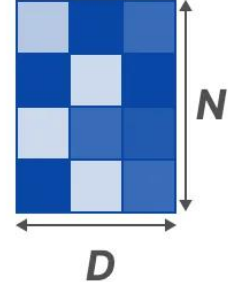
Matrix B

How the **input** influences the state



Matrix C

How the **current state** translates to the **output**



Step size (Δ)

Resolution of the **input** (discretization parameter)

Matrix B

How the **input** influences the state

Matrix C

How the **current state** translates to the **output**

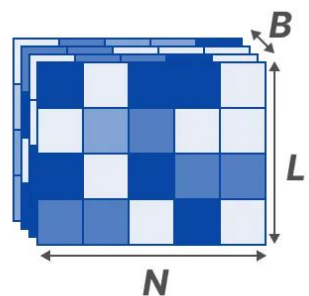
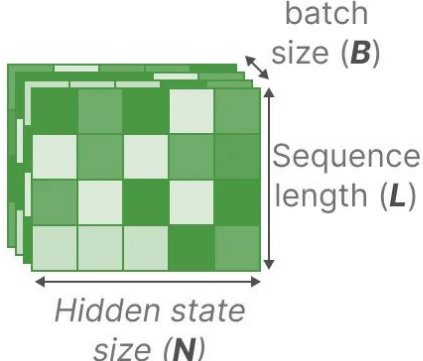
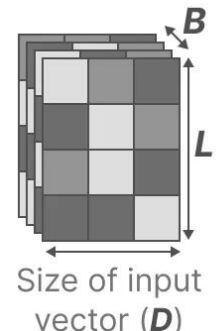
$$s_B(x) = \text{Linear}_N(x)$$

$$s_C(x) = \text{Linear}_N(x)$$

$$s_\Delta(x) = \text{Linear}_D(x)$$

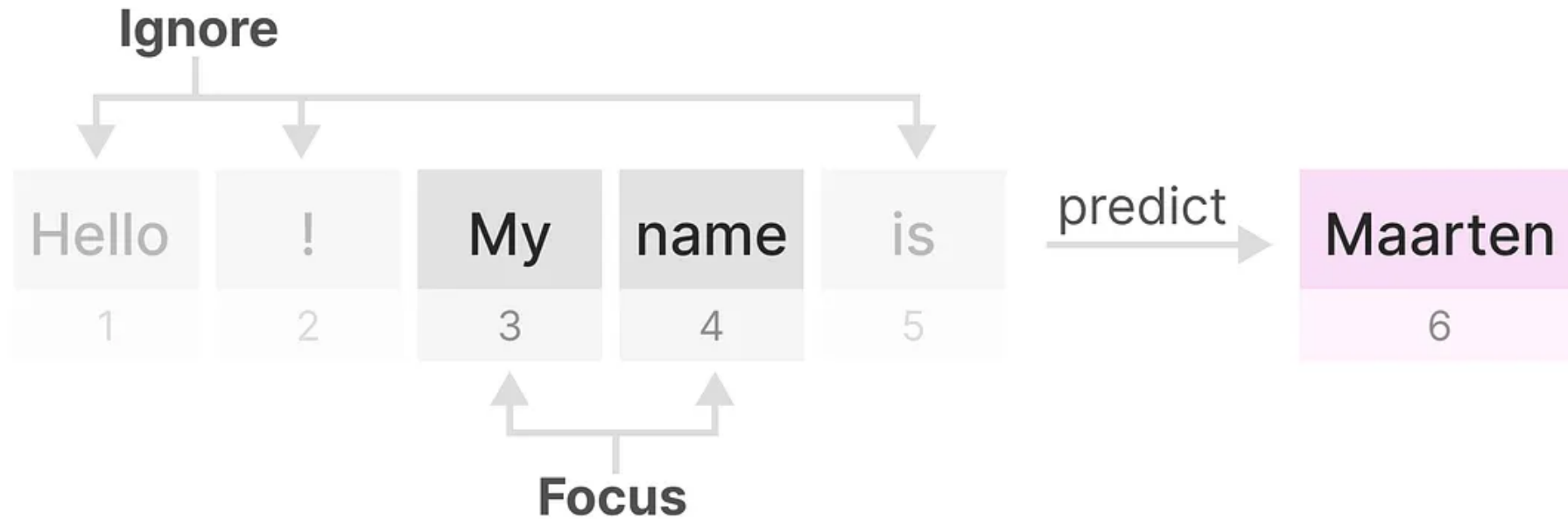
$$\tau_\Delta = \text{softplus}$$

SSM + Selection



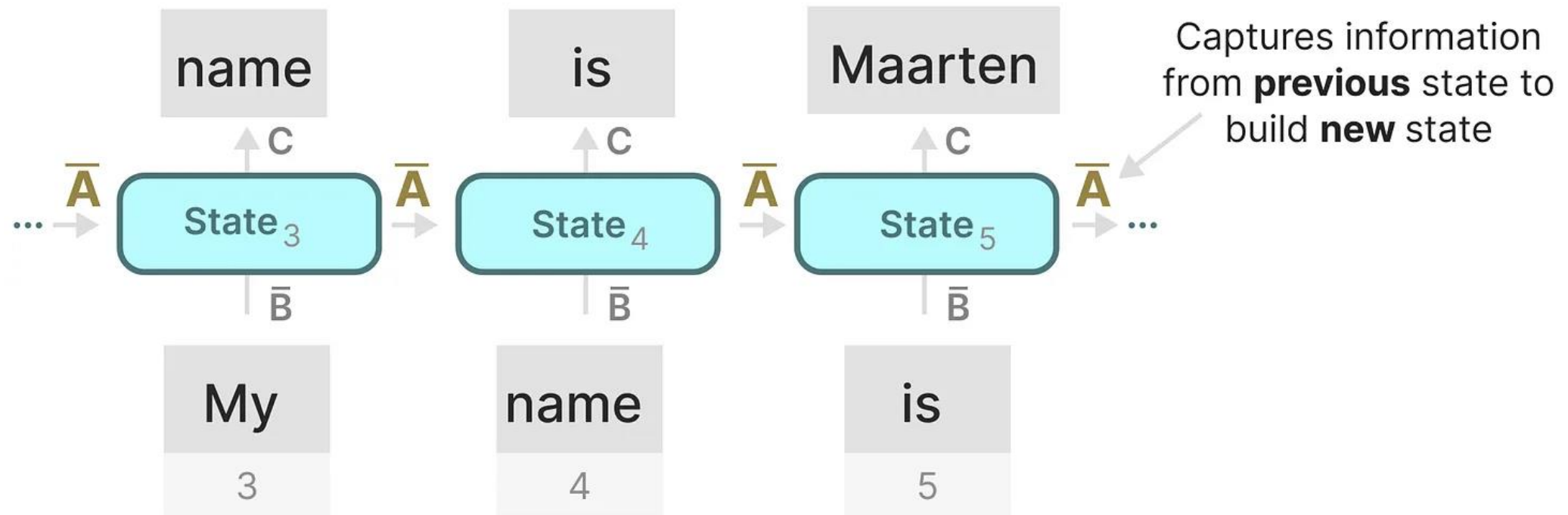
Mamba—S6

- Selection Mechanism



Mamba—S6

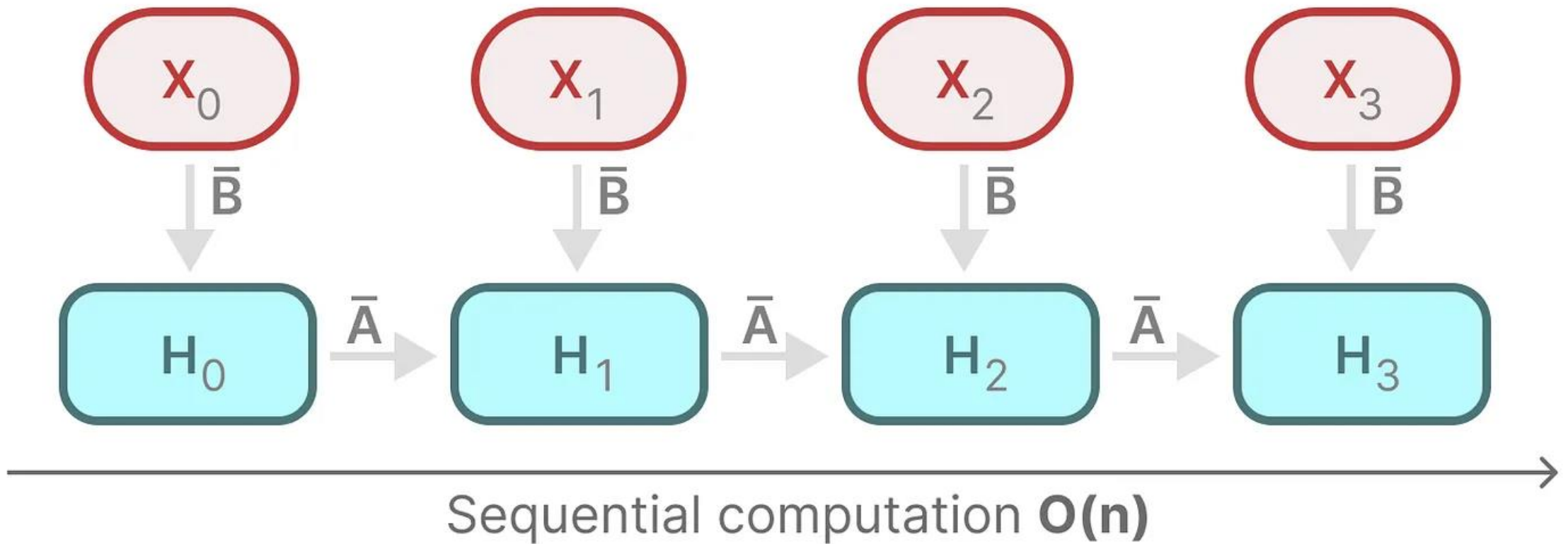
- Selection Mechanism



SSM
(Recurrent + Unfolded)

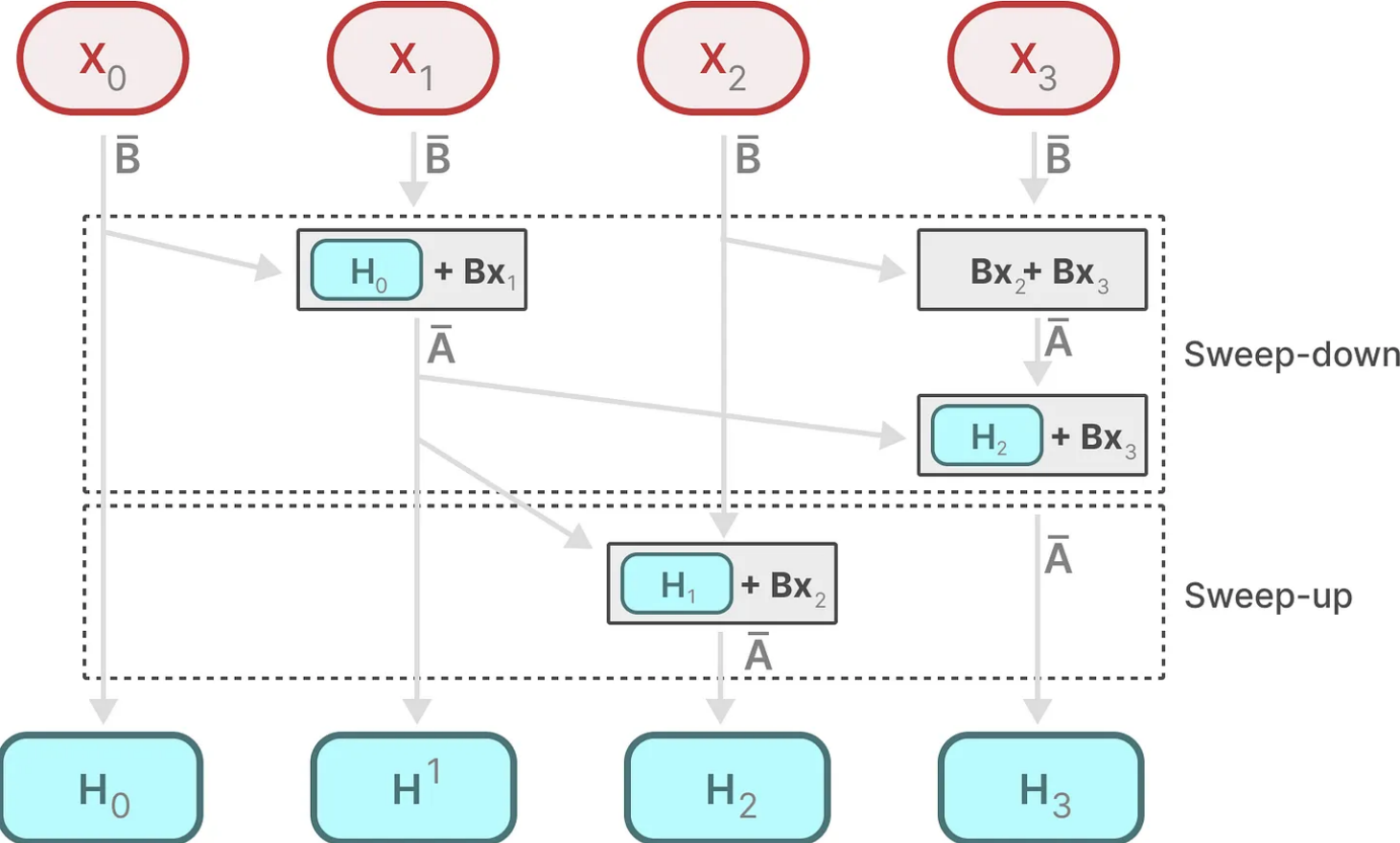
Mamba—S6

- Hardware-aware Algorithm



Mamba——S6

- Hardware-aware Algorithm

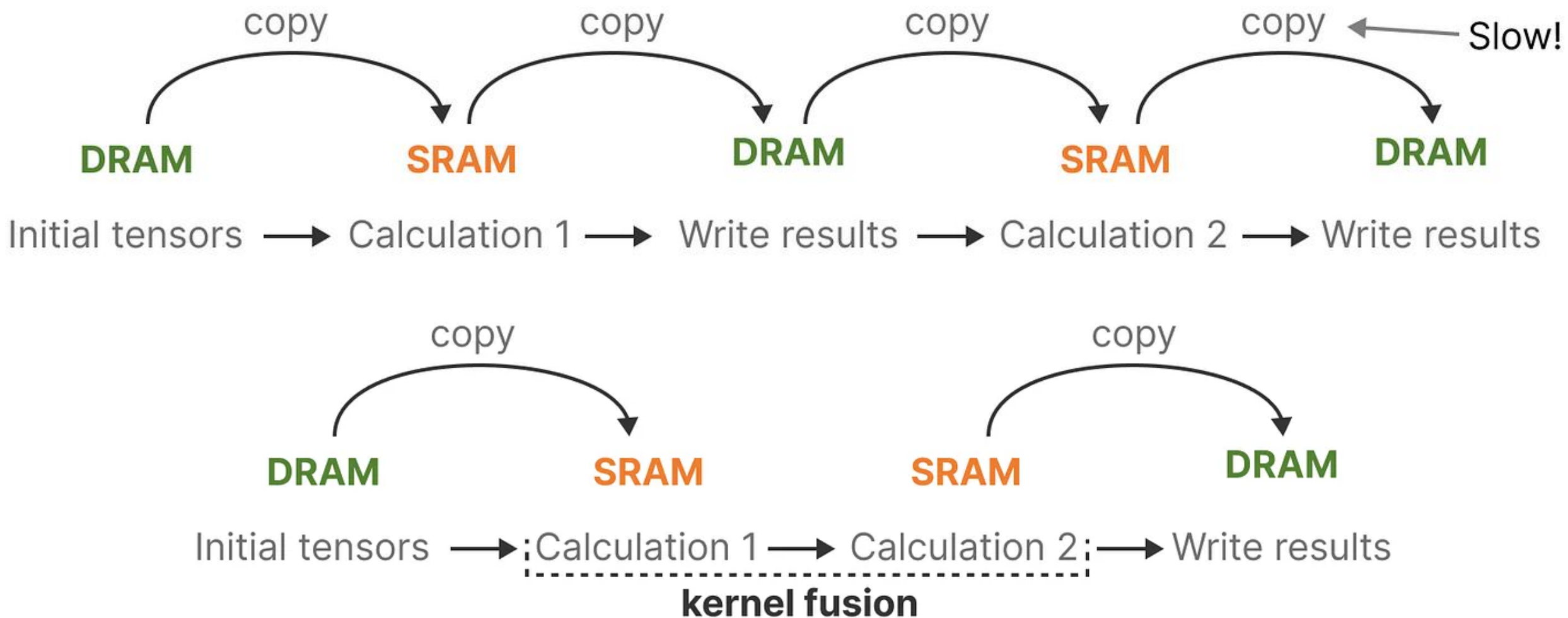


$$\begin{aligned}
 H_1 &= \bar{A} \cdot H_0 + BX_1 \\
 H_2 &= \bar{A} \cdot H_1 + BX_2 \\
 &= \bar{A} \cdot (\bar{A}H_0 + BX_1) + BX_2 \\
 H_3 &= \bar{A} \cdot H_2 + BX_3 \\
 &= \bar{A} \cdot (\bar{A} \cdot H_1 + BX_2) + BX_3 \\
 &= \bar{A} \cdot (\bar{A} \cdot (\bar{A}H_0 + BX_1) + BX_2) + BX_3
 \end{aligned}$$

Parallel computation $O(n/t)$

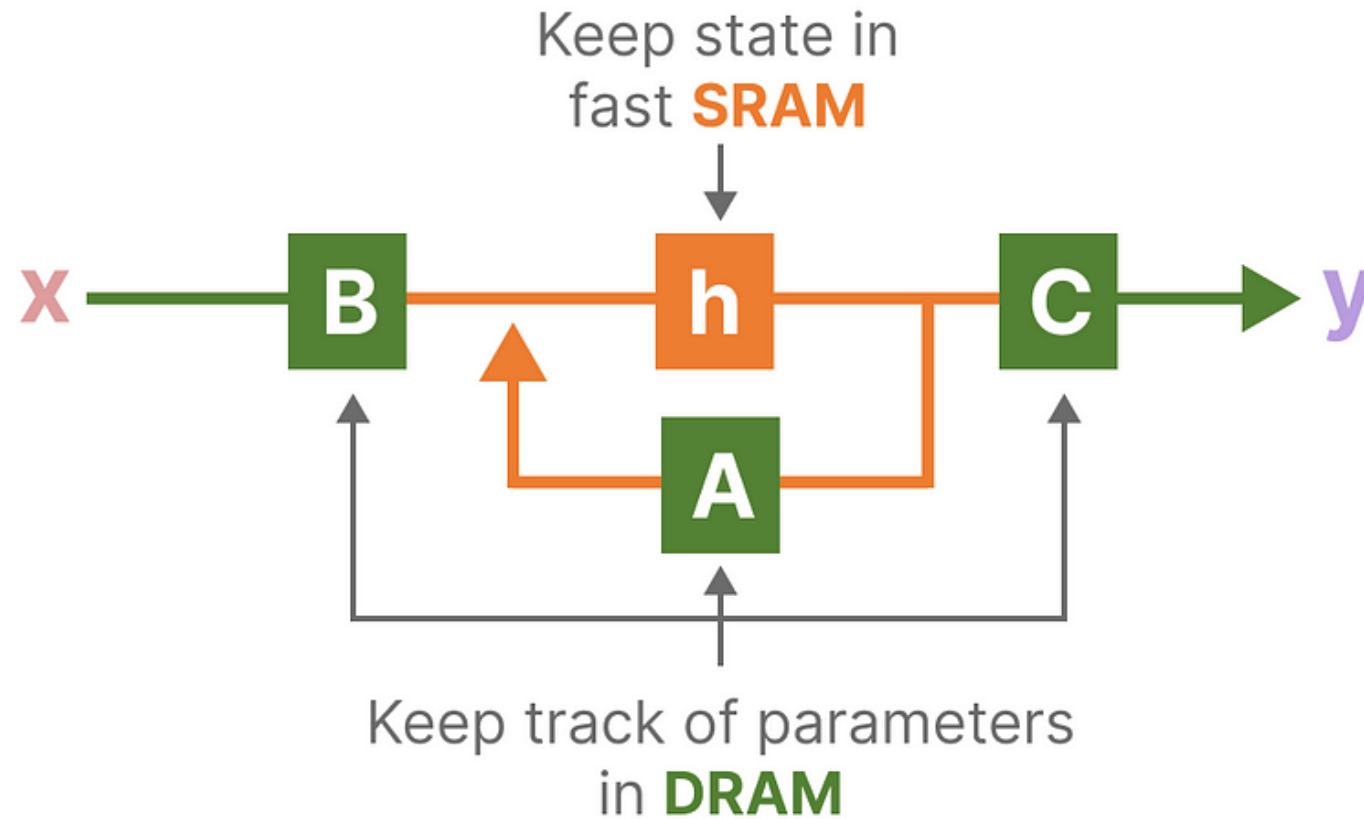
Mamba—S6

- Hardware-aware Algorithm (2)



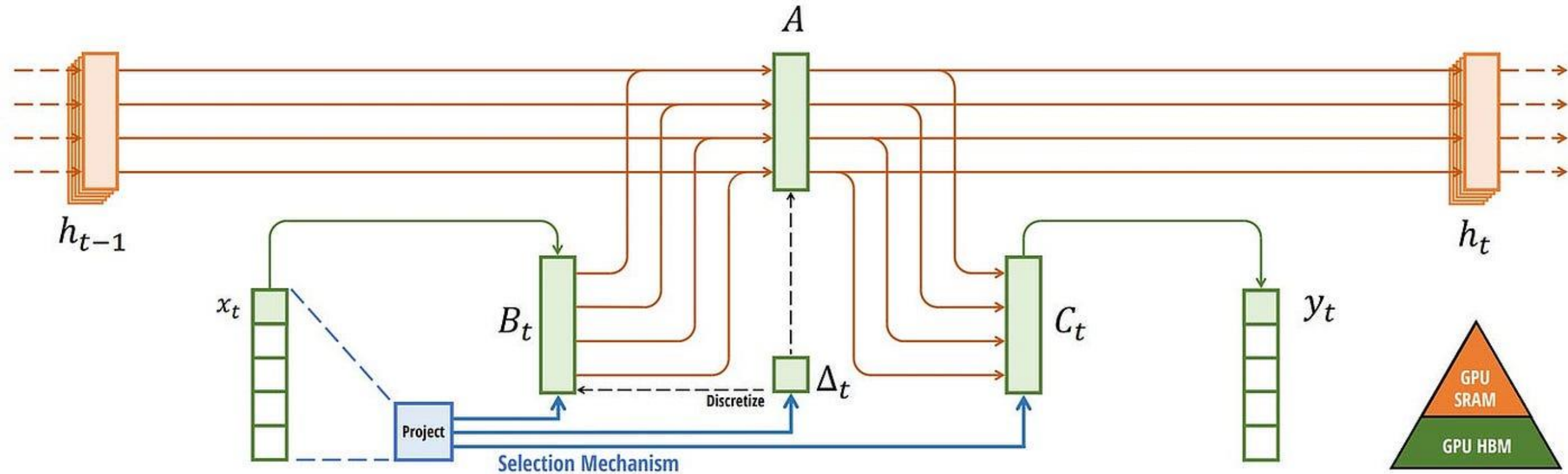
Mamba—S6

- Hardware-aware Algorithm (2)

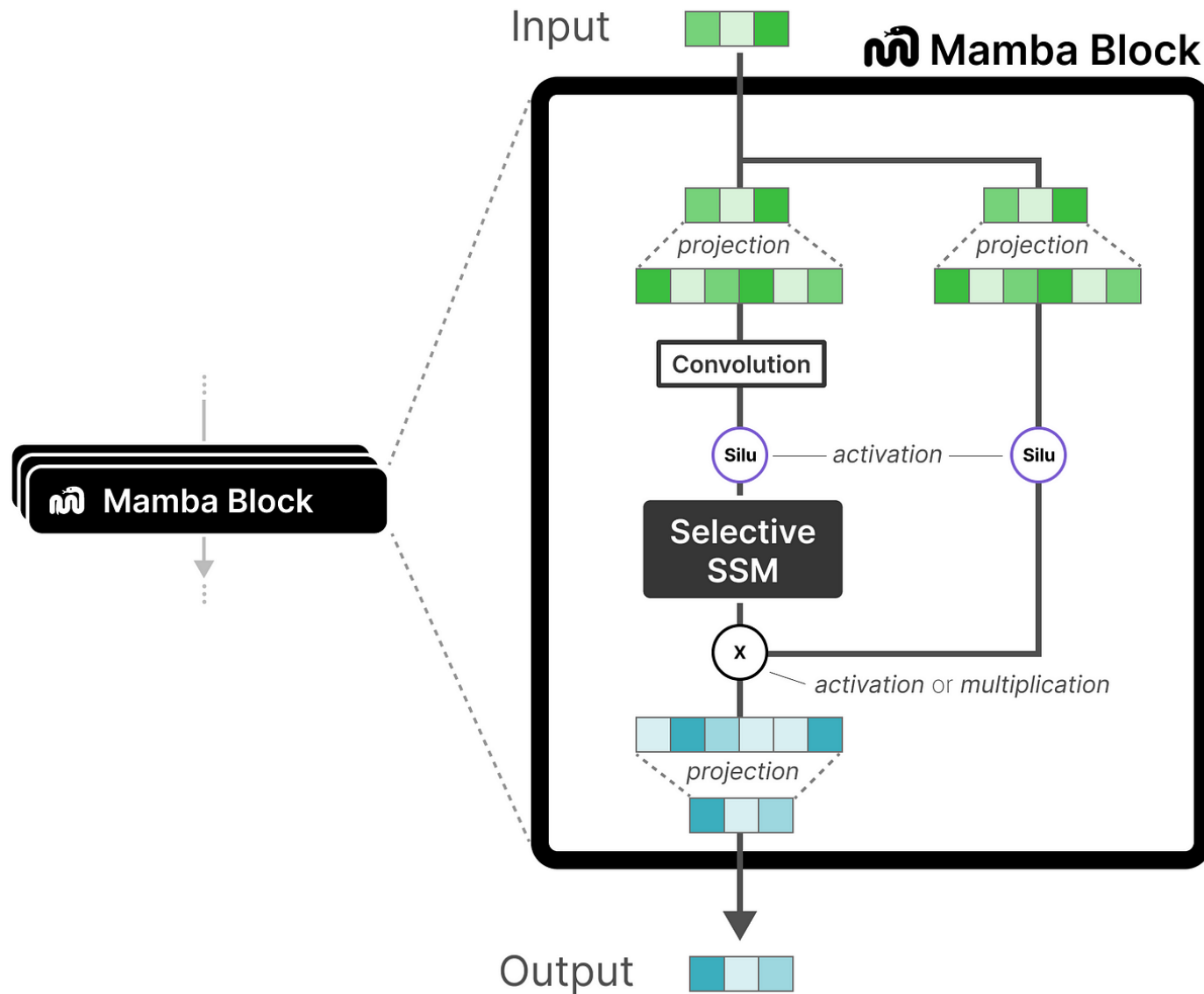


Mamba—S6

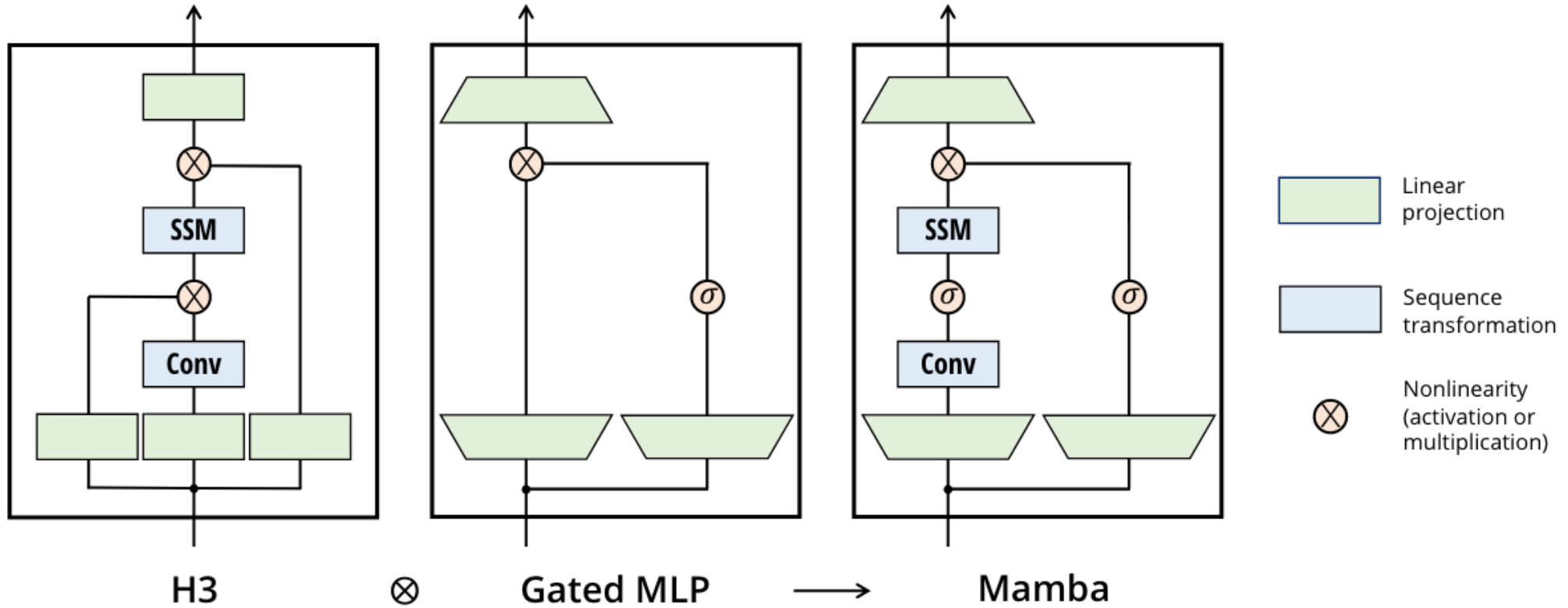
- Hardware-aware Algorithm (2)



Mamba—S6



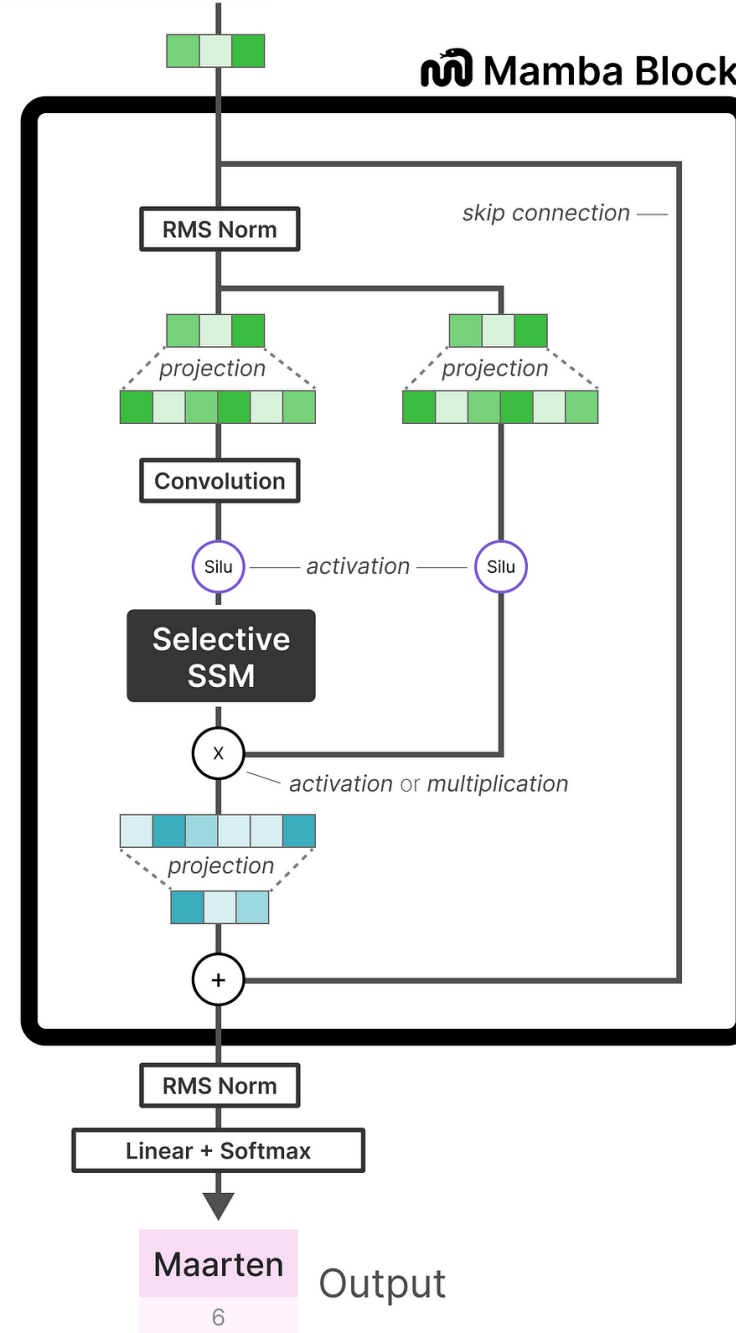
Mamba—S6



Mamba—S6

Hello	!	My	name	is	Input
1	2	3	4	5	

Repeated n times →

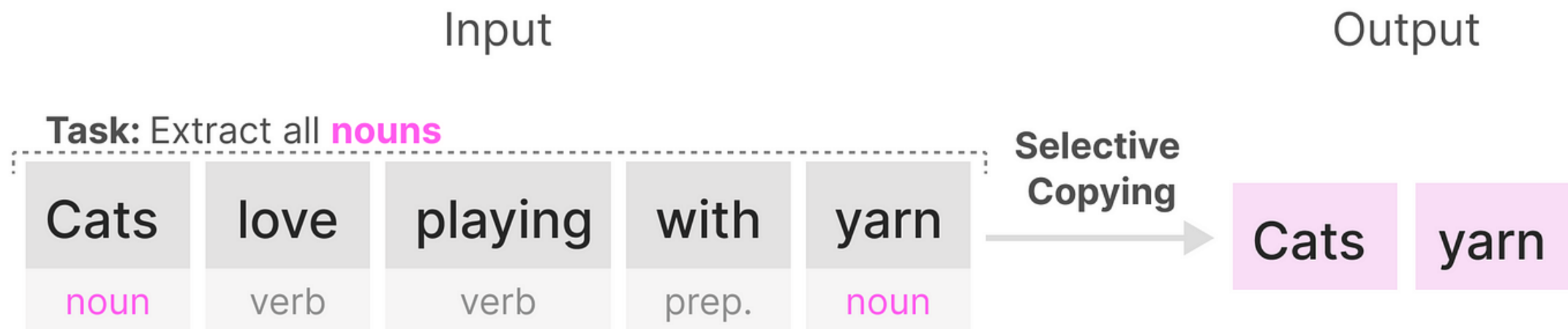


Experiments

Table 3: (**Zero-shot Evaluations.**) Best results for each size in bold. We compare against open source LMs with various tokenizers, trained for up to 300B tokens. Pile refers to the validation split, comparing only against models trained on the same dataset and tokenizer (GPT-NeoX-20B). For each model size, Mamba is best-in-class on every single evaluation result, and generally matches baselines at twice the model size.

Model	Token.	Pile ppl ↓	LAMBADA ppl ↓	LAMBADA acc ↑	HellaSwag acc ↑	PIQA acc ↑	Arc-E acc ↑	Arc-C acc ↑	WinoGrande acc ↑	Average acc ↑
Hybrid H3-130M	GPT2	—	89.48	25.77	31.7	64.2	44.4	24.2	50.6	40.1
Pythia-160M	NeoX	29.64	38.10	33.0	30.2	61.4	43.2	24.1	51.9	40.6
Mamba-130M	NeoX	10.56	16.07	44.3	35.3	64.5	48.0	24.3	51.9	44.7
Hybrid H3-360M	GPT2	—	12.58	48.0	41.5	68.1	51.4	24.7	54.1	48.0
Pythia-410M	NeoX	9.95	10.84	51.4	40.6	66.9	52.1	24.6	53.8	48.2
Mamba-370M	NeoX	8.28	8.14	55.6	46.5	69.5	55.1	28.0	55.3	50.0
Pythia-1B	NeoX	7.82	7.92	56.1	47.2	70.7	57.0	27.1	53.5	51.9
Mamba-790M	NeoX	7.33	6.02	62.7	55.1	72.1	61.2	29.5	56.1	57.1
GPT-Neo 1.3B	GPT2	—	7.50	57.2	48.9	71.1	56.2	25.9	54.9	52.4
Hybrid H3-1.3B	GPT2	—	11.25	49.6	52.6	71.3	59.2	28.1	56.9	53.0
OPT-1.3B	OPT	—	6.64	58.0	53.7	72.4	56.7	29.6	59.5	55.0
Pythia-1.4B	NeoX	7.51	6.08	61.7	52.1	71.0	60.5	28.5	57.2	55.2
RWKV-1.5B	NeoX	7.70	7.04	56.4	52.5	72.4	60.5	29.4	54.6	54.3
Mamba-1.4B	NeoX	6.80	5.04	64.9	59.1	74.2	65.5	32.8	61.5	59.7
GPT-Neo 2.7B	GPT2	—	5.63	62.2	55.8	72.1	61.1	30.2	57.6	56.5
Hybrid H3-2.7B	GPT2	—	7.92	55.7	59.7	73.3	65.6	32.3	61.4	58.0
OPT-2.7B	OPT	—	5.12	63.6	60.6	74.8	60.8	31.3	61.0	58.7
Pythia-2.8B	NeoX	6.73	5.04	64.7	59.3	74.0	64.1	32.9	59.7	59.1
RWKV-3B	NeoX	7.00	5.24	63.9	59.6	73.7	67.8	33.1	59.6	59.6
Mamba-2.8B	NeoX	6.22	4.23	69.2	66.1	75.2	69.7	36.3	63.5	63.3
GPT-J-6B	GPT2	—	4.10	68.3	66.3	75.4	67.0	36.6	64.1	63.0
OPT-6.7B	OPT	—	4.25	67.7	67.2	76.3	65.6	34.9	65.5	62.9
Pythia-6.9B	NeoX	6.51	4.45	67.1	64.0	75.2	67.3	35.5	61.3	61.7
RWKV-7.4B	NeoX	6.31	4.38	67.2	65.5	76.1	67.8	37.5	61.0	62.5

Experiments

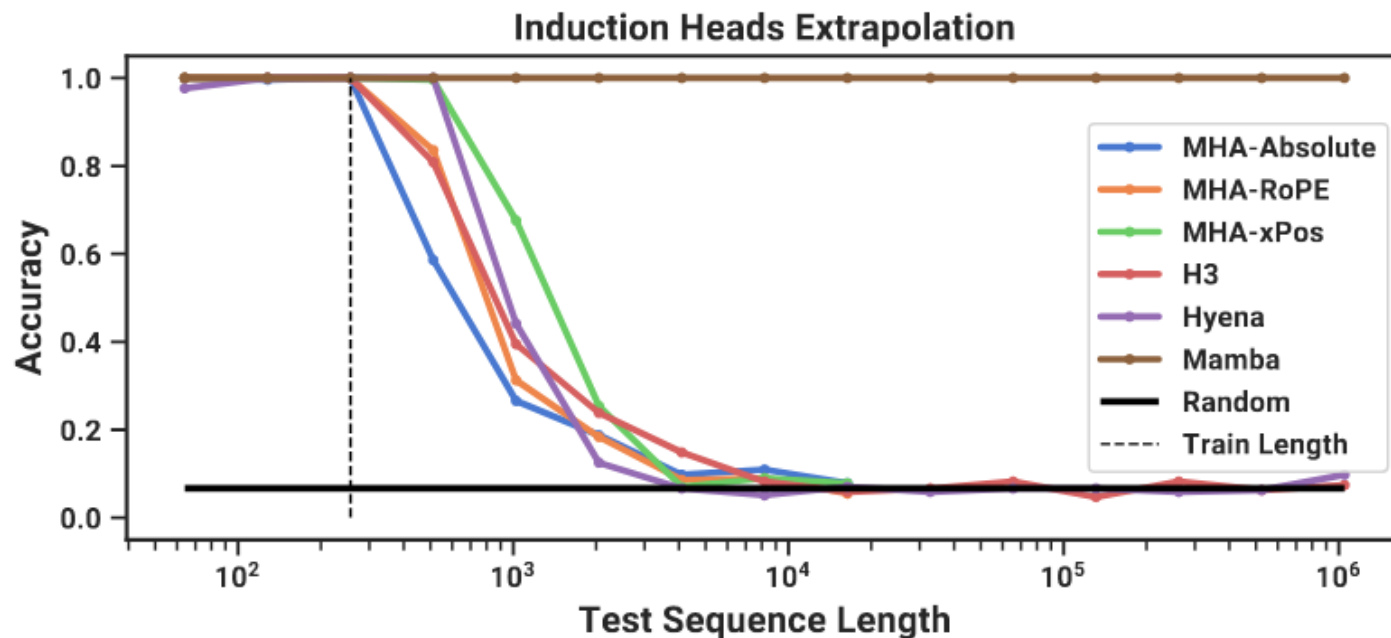
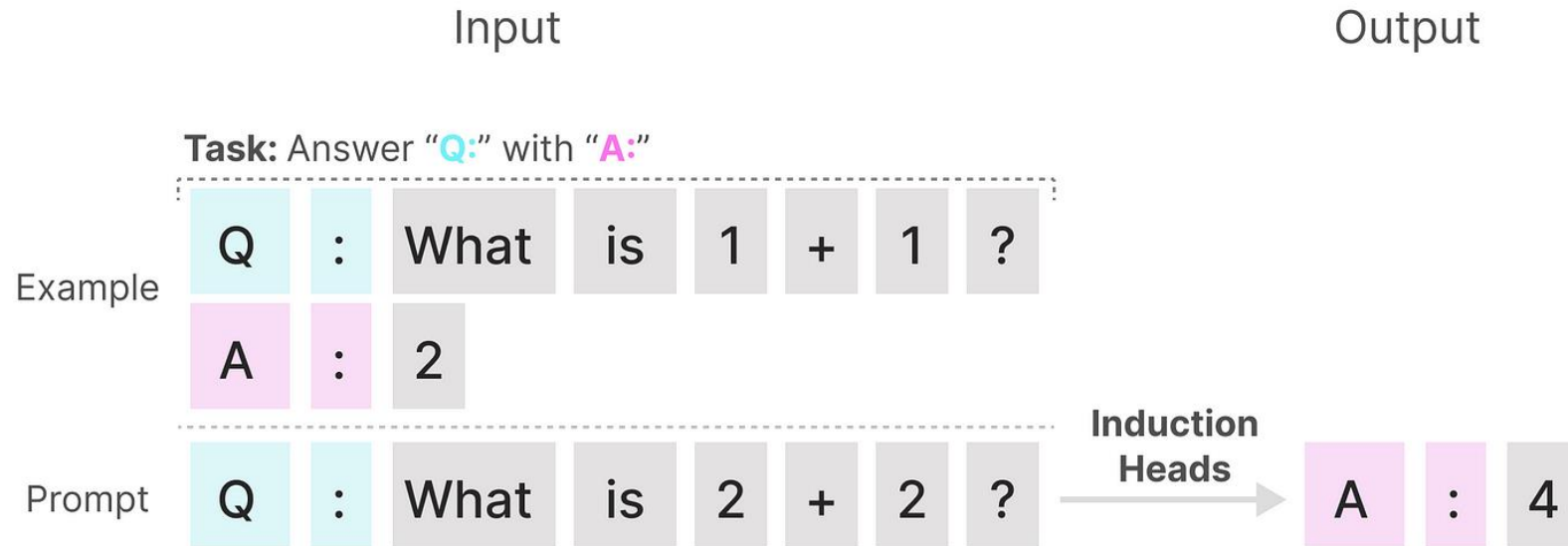


Model	Arch.	Layer	Acc.
S4	No gate	S4	18.3
-	No gate	S6	97.0
H3	H3	S4	57.0
Hyena	H3	Hyena	30.1
-	H3	S6	99.7
-	Mamba	S4	56.4
-	Mamba	Hyena	28.4
Mamba	Mamba	S6	99.8

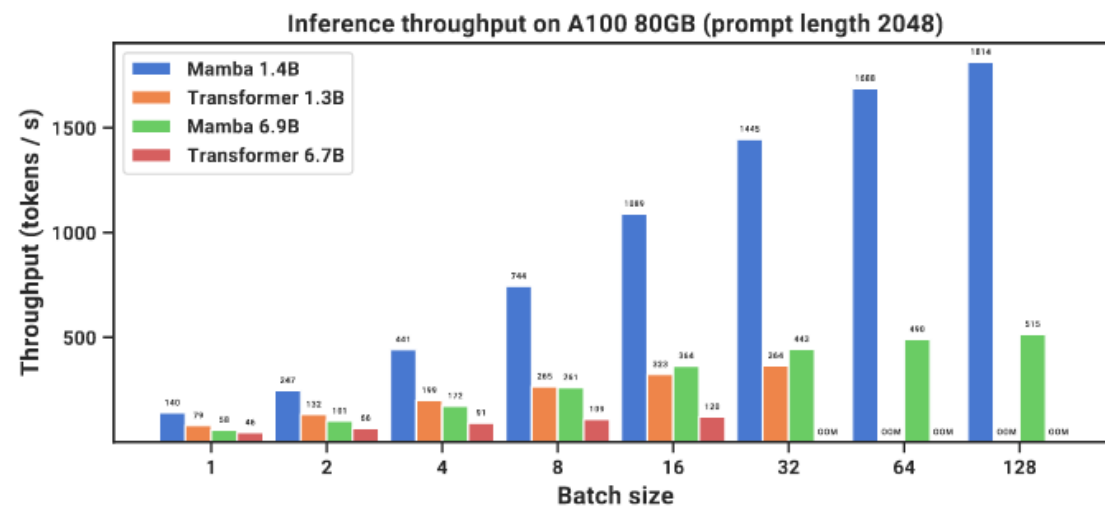
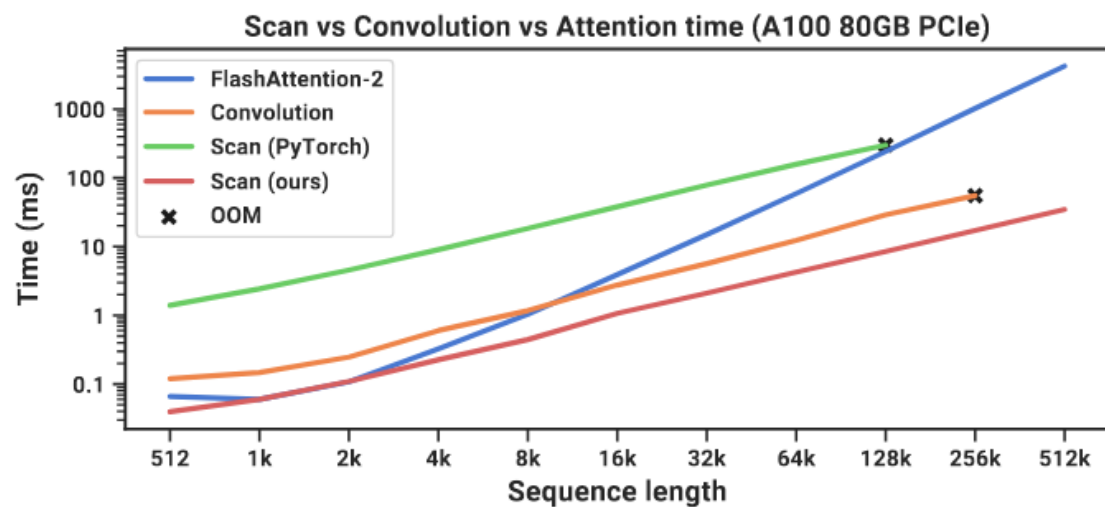
Table 1: (**Selective Copying.**)

Accuracy for combinations of architectures and inner sequence layers.

Experiments



Experiments



Experiments

Table 6: (**Ablations: Architecture and SSM layer.**) The Mamba block performs similarly to H3 while being simpler. In the inner layer, there is little difference among different parameterizations of LTI models, while selective SSMs (S6) provide a large improvement. More specifically, the S4 (real) variant is S4D-Real and the S4 (complex) variant is S4D-Lin.

Model	Arch.	SSM Layer	Perplexity
Hyena	H3	Hyena	10.24
H3	H3	S4 (complex)	10.30
-	H3	S4 (real)	10.34
-	H3	S6	8.95

Model	Arch.	SSM Layer	Perplexity
-	Mamba	Hyena	10.75
-	Mamba	S4 (complex)	10.54
-	Mamba	S4 (real)	10.56
Mamba	Mamba	S6	8.69

Why Rejected by ICLR24

- Need comparison with H3

Model	Lambada	HellaSwag	PIQA	Arc-E	Arc-C	WinoGrande	Avg
Hybrid H3-130M	25.8	31.7	64.2	44.4	24.2	50.6	40.2
Mamba-130M	44.2	35.2	64.5	48.0	24.2	52.3	44.7
Hybrid H3-360M	48.0	41.5	68.1	51.4	24.7	54.1	48.0
Mamba-370M	55.6	46.5	69.5	55.1	28.0	55.3	51.7
Hybrid H3-1.3B	49.6	52.6	71.3	59.2	28.1	56.9	53.0
Mamba-1.4B	65.0	59.1	74.2	65.5	32.8	61.5	59.7
Hybrid H3-2.7B	55.7	59.7	73.3	65.6	32.3	61.4	58.0
Mamba-2.8B	69.2	66.1	75.2	69.7	36.3	63.5	63.3

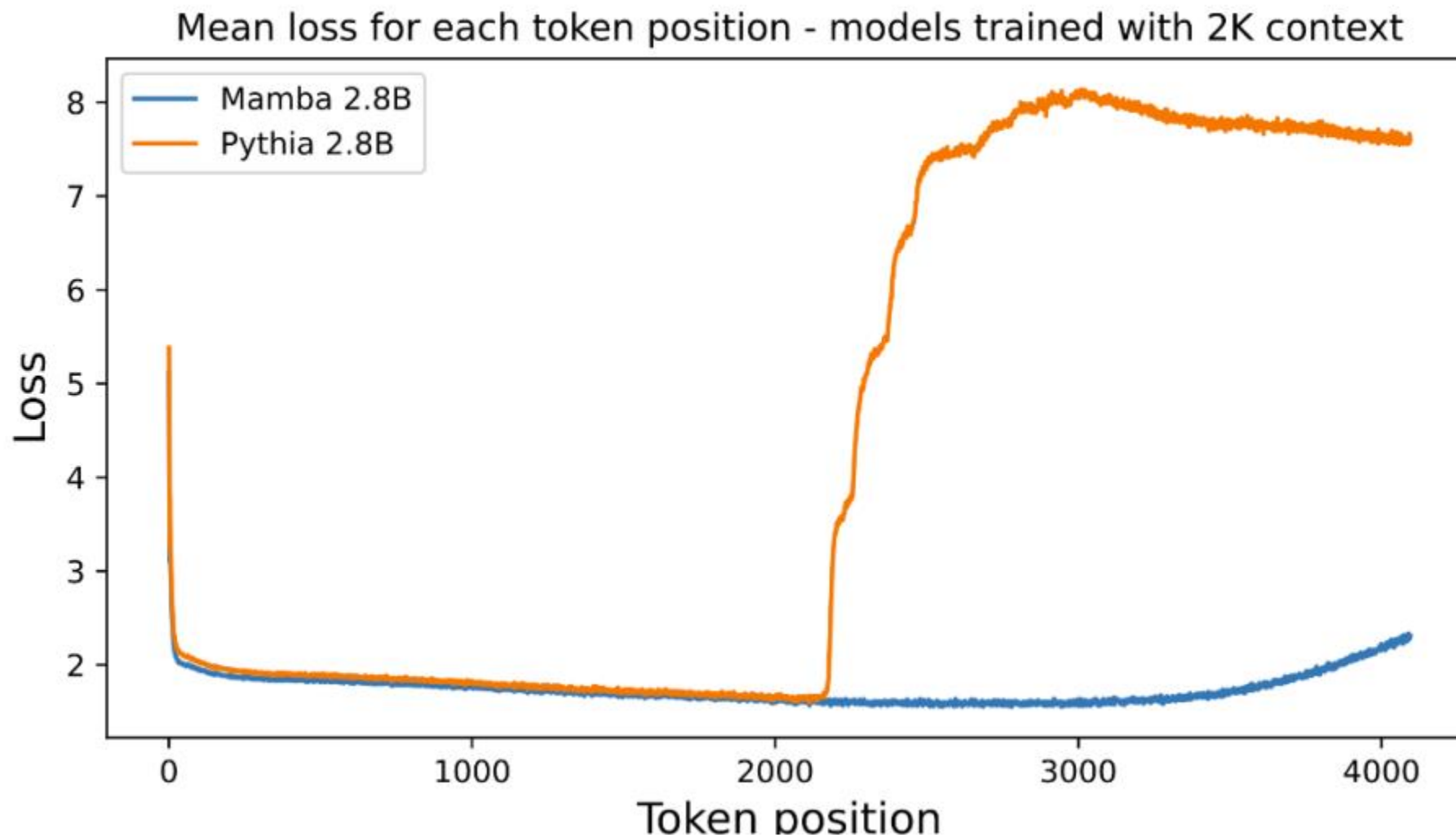
Why Rejected by ICLR24

- Scaling beyond 1.4B, vs Transformer 10B

Model	Lambada	HellaSwag	PIQA	Arc-E	Arc-C	WinoGrande	Avg
GPT-Neo 2.7B	62.2	55.8	72.1	61.1	30.2	57.6	56.5
Hybrid H3-2.7B	55.7	59.7	73.3	65.6	32.3	61.4	58.0
OPT-2.7B	63.6	60.6	74.8	60.8	31.3	61.0	58.7
Pythia-2.8B	64.7	59.3	74.0	64.1	32.9	59.7	59.1
RWKV-3B	63.9	59.6	73.7	67.8	33.1	59.6	59.6
Mamba-2.8B	69.2	66.1	75.2	69.7	36.3	63.5	63.3
OPT-6.7B	67.7	67.2	76.3	65.6	34.9	65.5	62.9
Pythia-6.9B	67.1	64.0	75.2	67.3	35.5	61.3	61.7
RWKV-7.4B	67.2	65.5	76.1	67.8	37.5	61.0	62.5

Why Rejected by ICLR24

- Beyond Training Length



Conclusions

- Selection Mechanism
- Hardware-aware Algorithm
- Simpler SSM Architecture

Thank You !