



From Mamba to Structured Masked Attention: Past, Recent, and Present

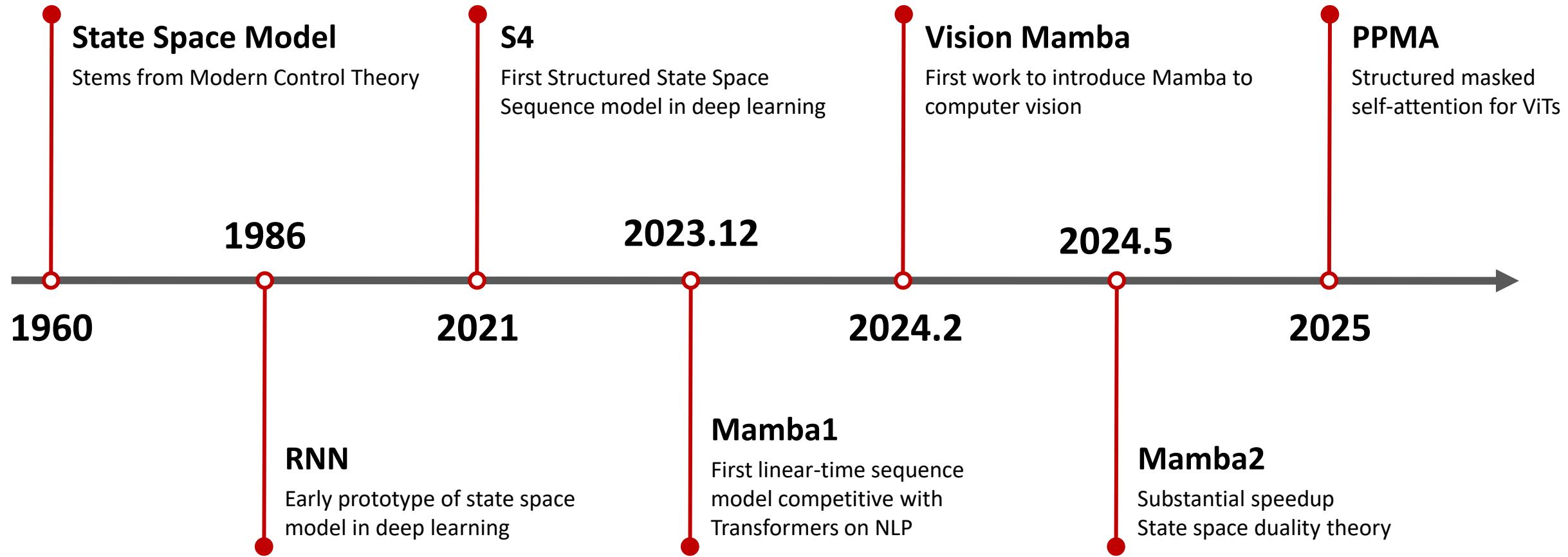
Zhongchen Zhao

2025/09/22

Past

Recent

Present





01 State Space Model

02 Mamba1

03 Vision Mamba

04 Mamba2

05 PPMA

1. State Space Model: 现代控制理论 (1960s)

4

■ 状态空间方程 (源于现代控制理论)

- 考虑一个质量-弹簧-阻尼系统的受力分析，一个质量为m的物体块质点

- 系统输入：外力 $x(t)$
- 系统输出：质量块的位置 $y(t)$
- 微分方程 (牛顿第二定律: $F = ma$)

$$m \times y''(t) = x(t) + mg - k \times y(t) - b \times y'(t)$$

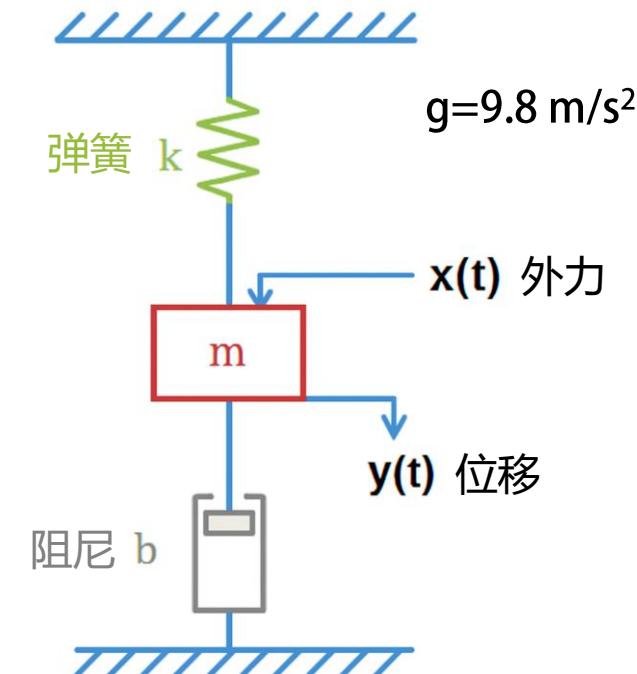
- 状态空间方程：使用一阶微分方程组就可以表示任意阶线性系统

状态向量 $\begin{bmatrix} y'(t) \\ y''(t) \end{bmatrix}$

$$\begin{bmatrix} y'(t) \\ y''(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{k}{m} & -\frac{b}{m} \end{bmatrix} \begin{bmatrix} y(t) \\ y'(t) \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{1}{m} \end{bmatrix} (x(t) + mg)$$

输出变量 $y(t) = [1 \quad 0] \begin{bmatrix} y(t) \\ y'(t) \end{bmatrix}$

输入变量 $x(t) + mg$

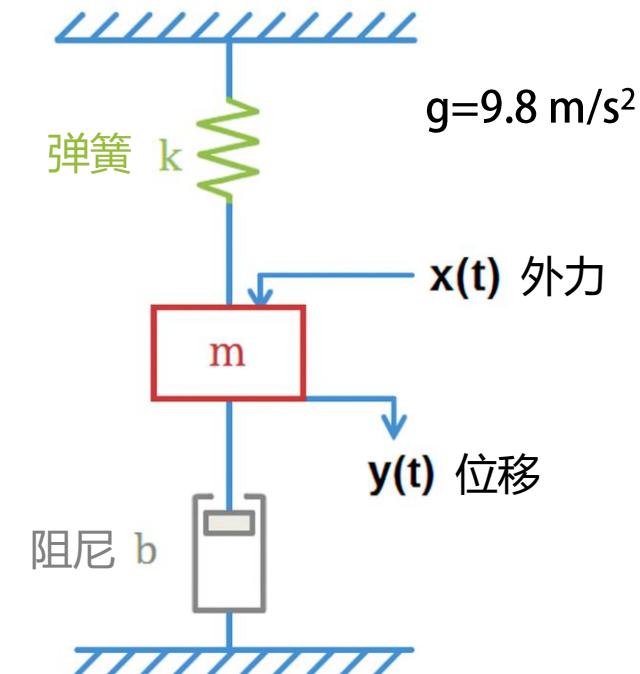


1. State Space Model: 现代控制理论 (1960s)

■ 状态空间方程 (源于现代控制理论)

- 考虑一个质量-弹簧-阻尼系统的受力分析，一个质量为m的物体块质点
 - 状态空间方程：使用一阶微分方程组就可以表示任意阶线性系统

$$\begin{aligned} \text{参数 } A & \quad \text{参数 } B \\ \begin{bmatrix} y'(t) \\ y''(t) \end{bmatrix} &= \begin{bmatrix} 0 & 1 \\ -\frac{k}{m} & -\frac{b}{m} \end{bmatrix} \begin{bmatrix} y(t) \\ y'(t) \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{1}{m} \end{bmatrix} (x(t) + mg) \\ \text{状态向量} & \quad \text{输入变量} \\ \text{隐变量} & \quad \text{输出变量} \quad \text{参数 } C \\ \text{记 } h(t) = \begin{bmatrix} y(t) \\ y'(t) \end{bmatrix}, \quad \begin{cases} h'(t) = A \times h(t) + B \times x(t) \\ y(t) = C \times h(t) \end{cases} \end{aligned}$$



1. State Space Model: 现代控制理论 (1960s)

■ 状态空间方程 (源于现代控制理论)

- 用向量和矩阵形式描述**动态系统**的方法，能够统一表示单输入单输出 (**SISO**) 和多输入多输出 (**MIMO**) 系统，适用于**连续或离散、线性或非线性、时变或时不变**系统
- 对于**连续时间、线性时不变**系统，状态空间方程为：

$$\begin{cases} \dot{\mathbf{h}}(t) = \mathbf{A} \times \mathbf{h}(t) + \mathbf{B} \times \mathbf{x}(t) \\ \mathbf{y}(t) = \mathbf{C} \times \mathbf{h}(t) \end{cases} \quad \gg \quad \text{Structured State Space Sequence Model (S4)}$$

- 对于**连续时间、线性时变**系统，状态空间方程为：

$$\begin{cases} \dot{\mathbf{h}}(t) = \boxed{\mathbf{A}(t)} \times \mathbf{h}(t) + \boxed{\mathbf{B}(t)} \times \mathbf{x}(t) \\ \mathbf{y}(t) = \boxed{\mathbf{C}(t)} \times \mathbf{h}(t) \end{cases} \quad \gg \quad \text{Selective State Space Model (Mamba)}$$

1. State Space Model: 现代控制理论 (1960s)

7

■ 状态空间方程 (源于现代控制理论)

- 用向量和矩阵形式描述动态系统的方法，能够统一表示单输入单输出 (**SISO**) 和多输入多输出 (**MIMO**) 系统，适用于连续或离散、线性或非线性、时变或时不变系统

系统类型	状态空间方程	系统类型	状态空间方程
连续线性时不变	$\begin{cases} \dot{\mathbf{h}}(t) = \mathbf{A} \times \mathbf{h}(t) + \mathbf{B} \times x(t) \\ y(t) = \mathbf{C} \times \mathbf{h}(t) \end{cases}$	离散线性时不变 S4	$\begin{cases} \mathbf{h}[k+1] = \bar{\mathbf{A}} \times \mathbf{h}[k] + \bar{\mathbf{B}} \times x[k] \\ y[k] = \bar{\mathbf{C}} \times \mathbf{h}[k] \end{cases}$
连续线性时变	$\begin{cases} \dot{\mathbf{h}}(t) = \mathbf{A}(t) \times \mathbf{h}(t) + \mathbf{B}(t) \times x(t) \\ y(t) = \mathbf{C}(t) \times \mathbf{h}(t) \end{cases}$	离散线性时变 Mamba	$\begin{cases} \mathbf{h}[k+1] = \bar{\mathbf{A}}[k] \times \mathbf{h}[k] + \bar{\mathbf{B}}[k] \times x[k] \\ y[k] = \bar{\mathbf{C}}[k] \times \mathbf{h}[k] \end{cases}$
连续非线性时不变	$\begin{cases} \dot{\mathbf{h}}(t) = f(\mathbf{h}(t), x(t)) \\ y(t) = g(\mathbf{h}(t), x(t)) \end{cases}$	离散非线性时不变 RNN	$\begin{cases} \mathbf{h}[k+1] = f(\mathbf{h}[k], x[k]) \\ y[k] = g(\mathbf{h}[k], x[k]) \end{cases}$
连续非线性时变	$\begin{cases} \dot{\mathbf{h}}(t) = f(\mathbf{h}(t), x(t), t) \\ y(t) = g(\mathbf{h}(t), x(t), t) \end{cases}$	离散非线性时变	$\begin{cases} \mathbf{h}[k+1] = f(\mathbf{h}[k], x[k], k) \\ y[k] = g(\mathbf{h}[k], x[k], k) \end{cases}$

1. State Space Model: 现代控制理论 (1960s)

■ 连续状态空间方程的离散化

- e^x 的多种近似方法:

- a) $\lim_{x \rightarrow 0} e^x \approx 1 + x \quad \checkmark$

- b) $\lim_{x \rightarrow 0} e^x \approx \frac{1+x/2}{1-x/2}$

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

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

$$\Rightarrow \mathbf{h}'(t)e^{-At} - \mathbf{A}\mathbf{h}(t)e^{-At} = \mathbf{B}x(t) e^{-At}$$

$$\Rightarrow \int \mathbf{h}(t)e^{-At} dt = \int \mathbf{B}x(t) e^{-At} dt$$

Δ : 时间采样间隔

设 $\mathbf{g}(t) = \mathbf{h}(t)e^{-At} \Rightarrow \mathbf{g}'(t) = e^{-At}\mathbf{B}x(t)$

$$\therefore \mathbf{h}(t + \Delta)e^{-At-\Delta A} - \mathbf{h}(t)e^{-At} = \mathbf{g}(t + \Delta) - \mathbf{g}(t)$$

$$= \int_t^{t+\Delta} \mathbf{g}'(t) dt$$

$$= \int_t^{t+\Delta} e^{-At}\mathbf{B}x(t) dt$$

$$\approx \int_t^{t+\Delta} e^{-At} \mathbf{B} dt \cdot \mathbf{x}(t) \quad (if \Delta \rightarrow 0)$$

$$= \mathbf{A}^{-1}e^{-At}(\mathbf{I} - e^{-\Delta A})\mathbf{B} \cdot \mathbf{x}(t)$$

$$\Rightarrow \mathbf{h}(t + \Delta) = e^{\Delta A}\mathbf{h}(t) + \mathbf{A}^{-1}(e^{\Delta A} - \mathbf{I})\mathbf{B}x(t)$$

$$\Rightarrow \mathbf{h}_t = \bar{\mathbf{A}} \mathbf{h}_{t-1} + \bar{\mathbf{B}} \mathbf{x}_t$$

where $\bar{\mathbf{A}} = e^{\Delta A}$, $\bar{\mathbf{B}} = \mathbf{A}^{-1}(e^{\Delta A} - \mathbf{I})\mathbf{B} \approx \Delta \mathbf{B}$

1. State Space Model: Structured State Space Model

9

■ Structured State Space Model (S4)

- S4将系统建模为线性时不变模型，旨在弥补 Transformer/RNN 在长序列建模中的不足

- 连续时间的状态空间方程

$$\begin{cases} \dot{\mathbf{h}}(t) = \mathbf{A} \times \mathbf{h}(t) + \mathbf{B} \times \mathbf{x}(t) \\ \mathbf{y}(t) = \mathbf{C} \times \mathbf{h}(t) + \mathbf{D} \times \mathbf{x}(t) \end{cases}$$

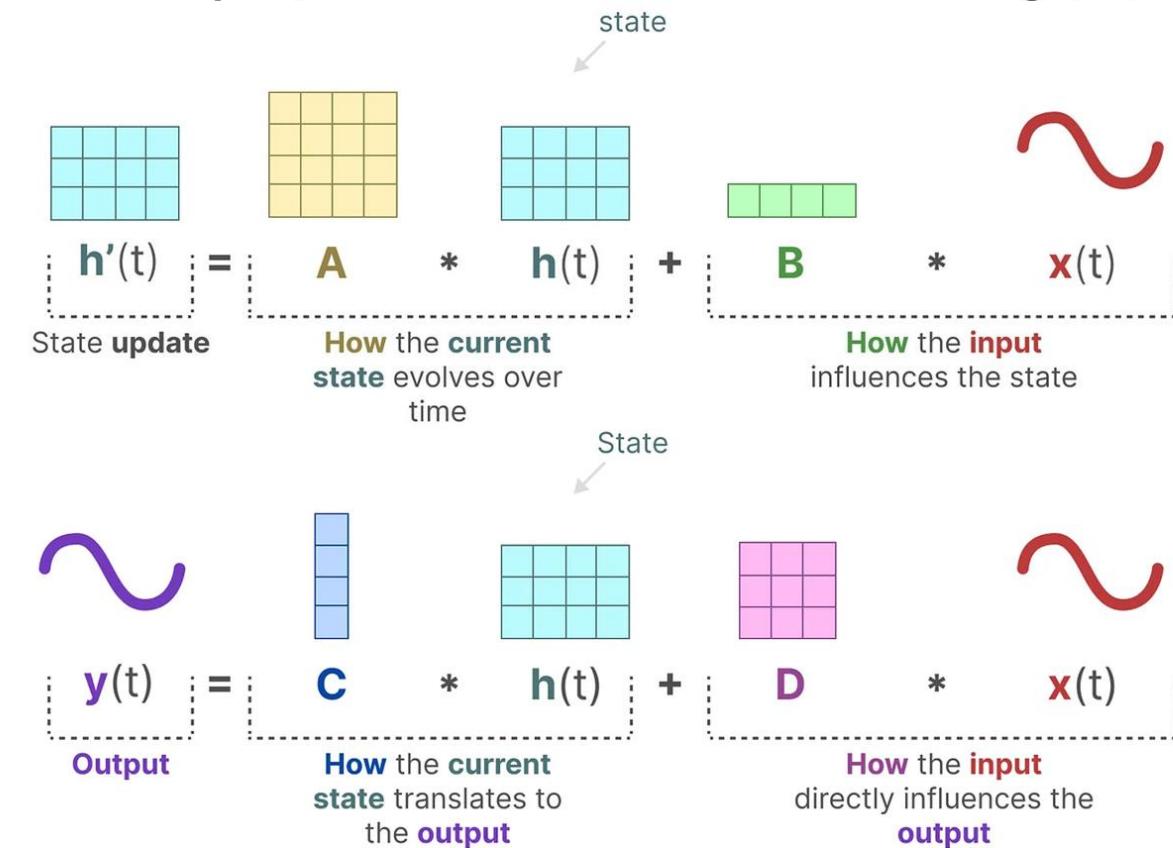
- 离散时间的状态空间方程

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times \mathbf{x}_t \\ \mathbf{y}_t = \mathbf{C} \times \mathbf{h}_t + \mathbf{D} \times \mathbf{x}_t \end{cases}$$

- 变量 $\mathbf{x}_t, \mathbf{y}_t \in \mathbb{R}$, $\mathbf{h}_t \in \mathbb{R}^D$

- 参数 $\bar{\mathbf{A}} = e^{\Delta \mathbf{A}} \in \mathbb{R}^{D \times D}$, $\bar{\mathbf{B}} = \Delta \mathbf{B} \in \mathbb{R}^{D \times 1}$

- 参数 $\mathbf{C} \in \mathbb{R}^{1 \times D}$, $\mathbf{D} \in \mathbb{R}$



1. State Space Model: Structured State Space Model

10

■ Structured State Space Model (S4)

- S4 将系统建模为线性时不变模型，旨在弥补 Transformer/RNN 在长序列建模中的不足

- 连续时间的状态空间方程

$$\begin{cases} \dot{\mathbf{h}}(t) = \mathbf{A} \times \mathbf{h}(t) + \mathbf{B} \times x(t) \\ y(t) = \mathbf{C} \times \mathbf{h}(t) + D \times x(t) \end{cases}$$

- 离散时间的状态空间方程

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times x_t \\ y_t = \mathbf{C} \times \mathbf{h}_t + D \times x_t \end{cases}$$

- 变量 $x_t, y_t \in \mathbb{R}$, $\mathbf{h}_t \in \mathbb{R}^D$

- 参数 $\bar{\mathbf{A}} = e^{\Delta A} \in \mathbb{R}^{D \times D}$, $\bar{\mathbf{B}} = \Delta \mathbf{B} \in \mathbb{R}^{D \times 1}$

- 参数 $\mathbf{C} \in \mathbb{R}^{1 \times D}$, $D \in \mathbb{R}$

- 参数 $\bar{\mathbf{A}} \in \mathbb{R}^{D \times D}$ 是一个 HiPPO 矩阵，定义为：

$$\bar{A}_{i,j} = \begin{cases} \sqrt{(2i+1)(2j+1)} & \text{if } i > j \text{ and } i+j \text{ is odd} \\ 0 & \text{otherwise} \end{cases}$$

- 例如， $D = 5$

$$\bar{\mathbf{A}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ \sqrt{3} & 0 & 0 & 0 & 0 \\ 0 & \sqrt{15} & 0 & 0 & 0 \\ \sqrt{7} & 0 & \sqrt{35} & 0 & 0 \\ 0 & \sqrt{27} & 0 & \sqrt{63} & 0 \end{bmatrix}$$

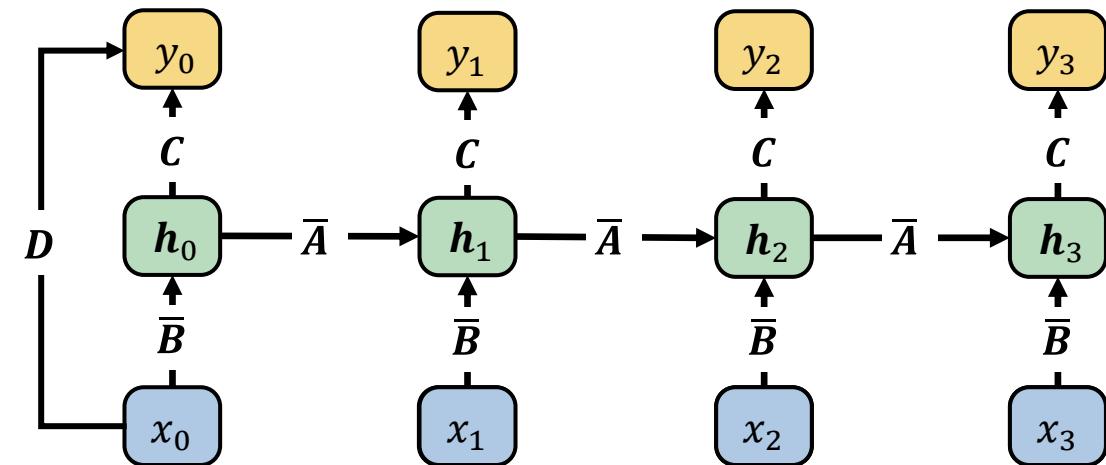
1. Structured State Space Model: Parallelizable Training

11

■ Structured State Space Model (S4)

- 状态空间模型的循环递推形式（串行计算，速度慢）

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times \mathbf{x}_t \\ \mathbf{y}_t = \mathbf{C} \times \mathbf{h}_t + \mathbf{D} \times \mathbf{x}_t \end{cases}$$



1. Structured State Space Model: Parallelizable Training

12

■ Structured State Space Model (S4)

- 状态空间模型的循环递推形式 (串行计算, 速度慢)

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times \mathbf{x}_t \\ y_t = \mathbf{C} \times \mathbf{h}_t + \cancel{\mathbf{D} \times \mathbf{x}_t} \end{cases}$$

- 状态空间模型的一维卷积形式 (并行计算, 速度快)

$$y_0 = \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_0$$

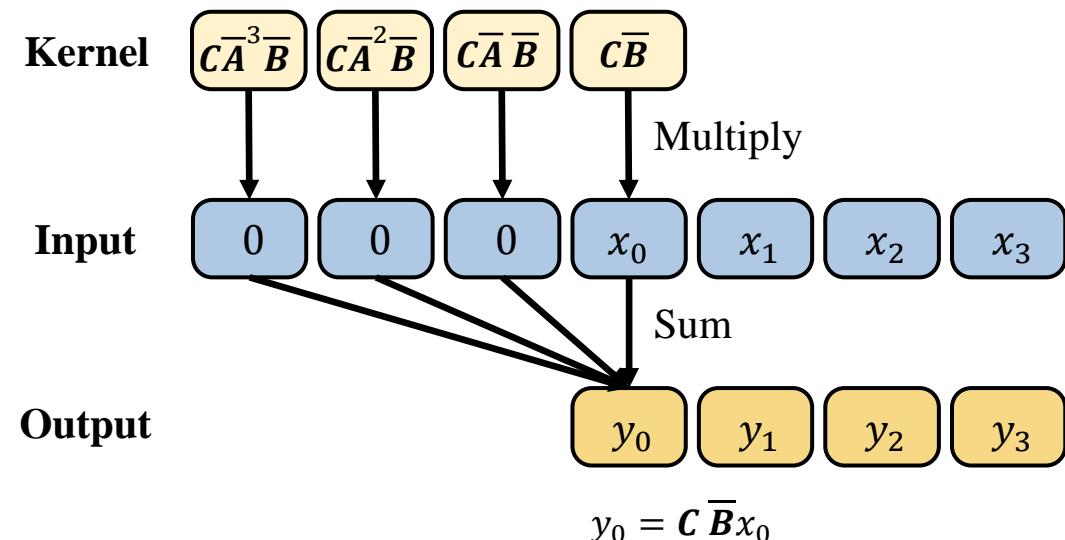
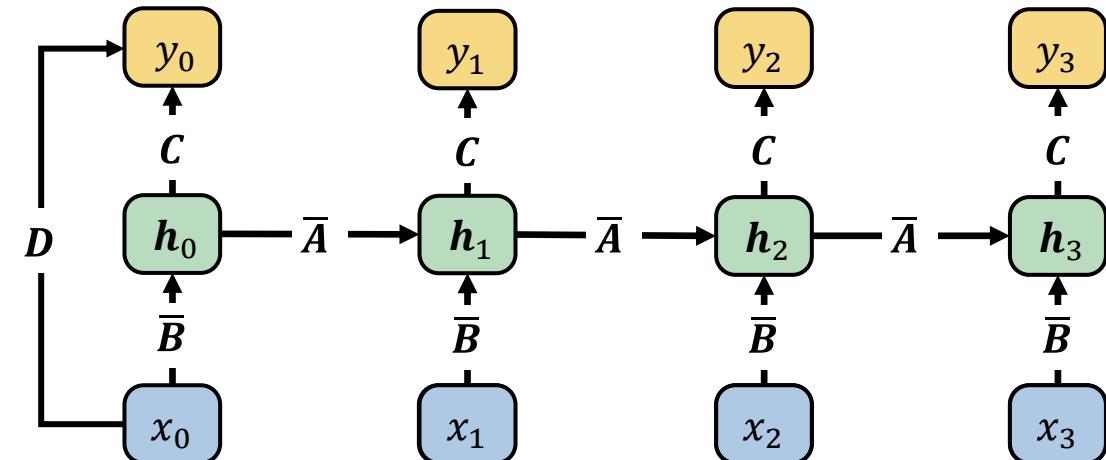
$$y_1 = \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_1$$

$$y_2 = \mathbf{C} \bar{\mathbf{A}}^2 \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_1 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_2$$

...

$$y_t = \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}}^{t-1} \bar{\mathbf{B}} \mathbf{x}_1 + \cdots + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_{t-1} + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_t$$

$$y_t = \bar{\mathbf{K}} * \vec{\mathbf{x}} = \sum_{n=0}^t \bar{\mathbf{K}}_n \mathbf{x}_{t-n}, \quad \bar{\mathbf{K}} = [\mathbf{C} \bar{\mathbf{B}}, \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}}, \dots, \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}}]$$



1. Structured State Space Model: Parallelizable Training

13

■ Structured State Space Model (S4)

- 状态空间模型的循环递推形式 (串行计算, 速度慢)

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times \mathbf{x}_t \\ y_t = \mathbf{C} \times \mathbf{h}_t + \cancel{\mathbf{D} \times \mathbf{x}_t} \end{cases}$$

- 状态空间模型的一维卷积形式 (并行计算, 速度快)

$$y_0 = \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_0$$

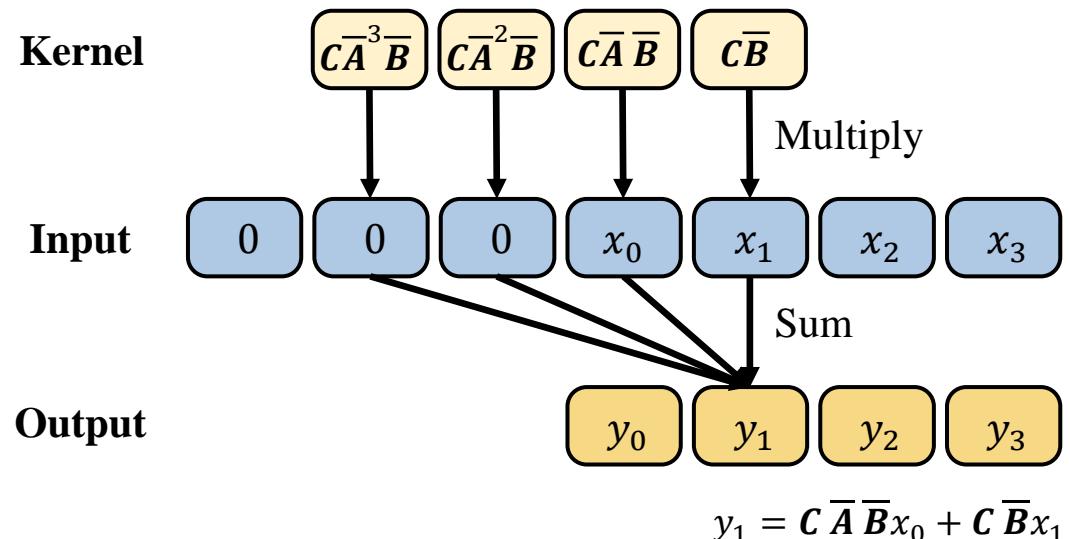
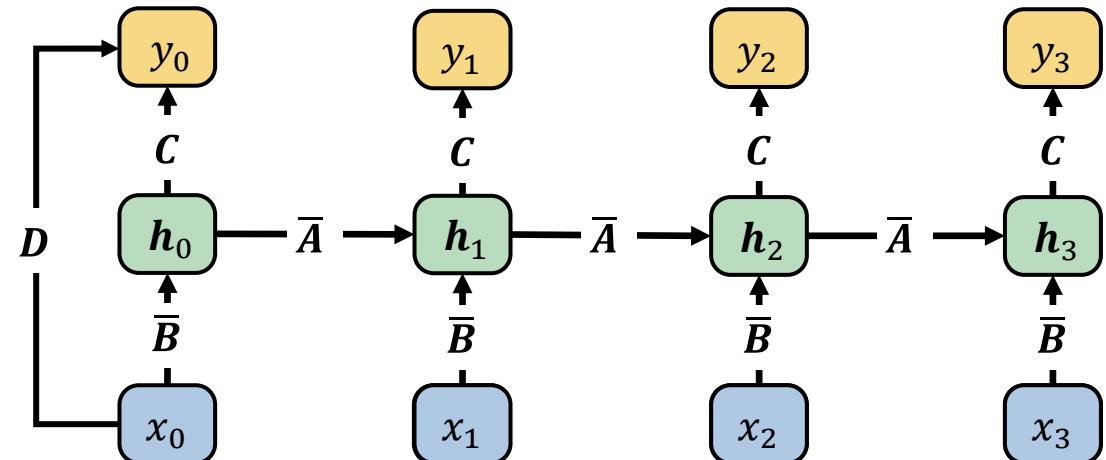
$$y_1 = \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_1$$

$$y_2 = \mathbf{C} \bar{\mathbf{A}}^2 \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_1 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_2$$

...

$$y_t = \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}}^{t-1} \bar{\mathbf{B}} \mathbf{x}_1 + \cdots + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_{t-1} + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_t$$

$$y_t = \bar{\mathbf{K}} * \vec{\mathbf{x}} = \sum_{n=0}^t \bar{\mathbf{K}}_n \mathbf{x}_{t-n}, \quad \bar{\mathbf{K}} = [\mathbf{C} \bar{\mathbf{B}}, \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}}, \dots, \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}}]$$



1. Structured State Space Model: Parallelizable Training

14

■ Structured State Space Model (S4)

- 状态空间模型的循环递推形式 (串行计算, 速度慢)

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times \mathbf{x}_t \\ y_t = \mathbf{C} \times \mathbf{h}_t + \cancel{\mathbf{D} \times \mathbf{x}_t} \end{cases}$$

- 状态空间模型的一维卷积形式 (并行计算, 速度快)

$$y_0 = \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_0$$

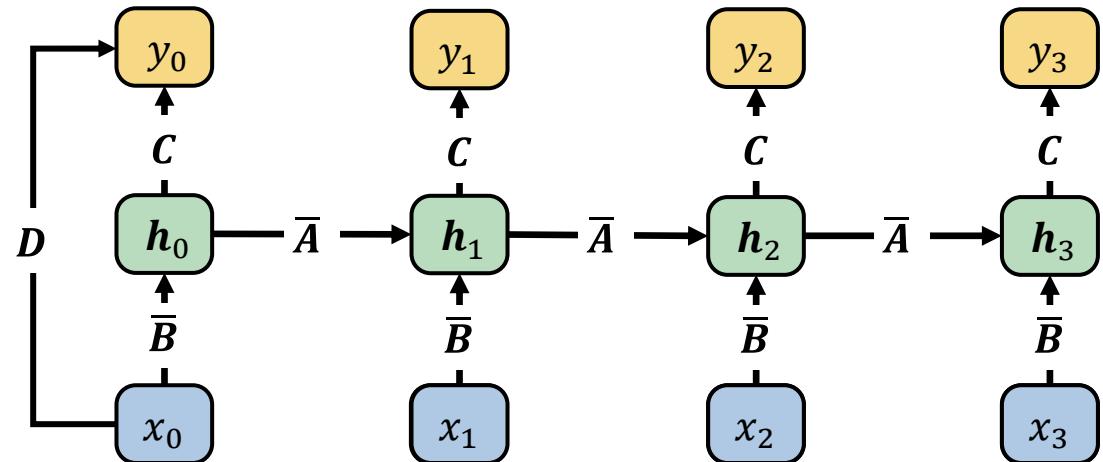
$$y_1 = \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_1$$

$$y_2 = \mathbf{C} \bar{\mathbf{A}}^2 \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_1 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_2$$

...

$$y_t = \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}}^{t-1} \bar{\mathbf{B}} \mathbf{x}_1 + \cdots + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_{t-1} + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_t$$

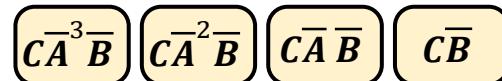
$$y_t = \bar{\mathbf{K}} * \vec{\mathbf{x}} = \sum_{n=0}^t \bar{\mathbf{K}}_n \mathbf{x}_{t-n}, \quad \bar{\mathbf{K}} = [\mathbf{C} \bar{\mathbf{B}}, \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}}, \dots, \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}}]$$



Kernel

Input

Output



Multiply

Sum



$$y_2 = \mathbf{C} \bar{\mathbf{A}}^2 \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_1 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_2$$

1. Structured State Space Model: Parallelizable Training

15

■ Structured State Space Model (S4)

- 状态空间模型的循环递推形式 (串行计算, 速度慢)

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times \mathbf{x}_t \\ y_t = \mathbf{C} \times \mathbf{h}_t + \cancel{\mathbf{D} \times \mathbf{x}_t} \end{cases}$$

- 状态空间模型的一维卷积形式 (并行计算, 速度快)

$$y_0 = \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_0$$

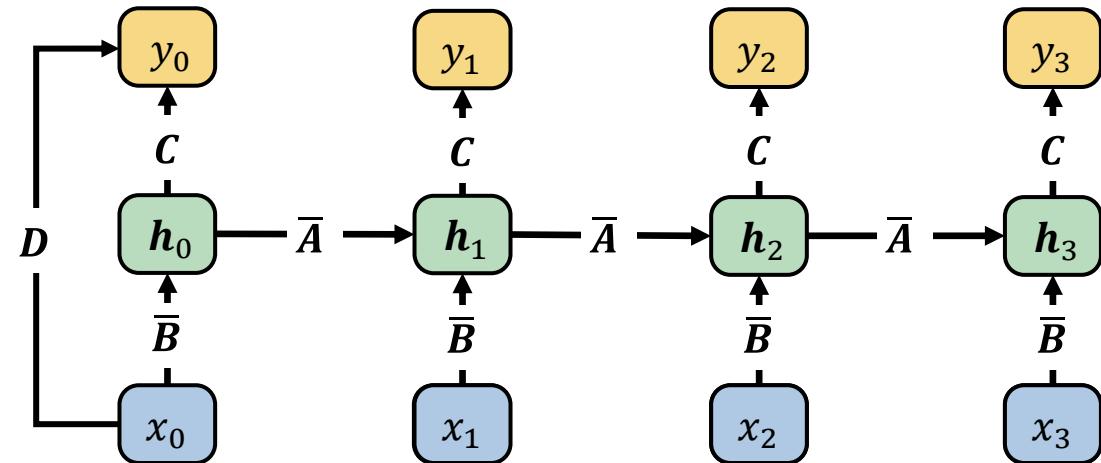
$$y_1 = \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_1$$

$$y_2 = \mathbf{C} \bar{\mathbf{A}}^2 \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_1 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_2$$

...

$$y_t = \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}}^{t-1} \bar{\mathbf{B}} \mathbf{x}_1 + \cdots + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_{t-1} + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_t$$

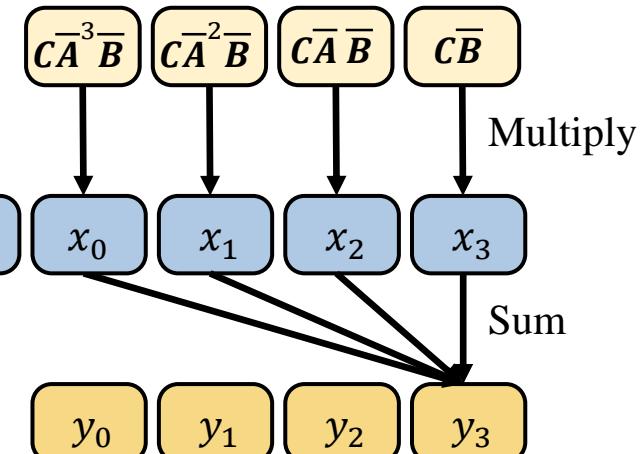
$$y_t = \bar{\mathbf{K}} * \vec{\mathbf{x}} = \sum_{n=0}^t \bar{\mathbf{K}}_n \mathbf{x}_{t-n}, \quad \bar{\mathbf{K}} = [\mathbf{C} \bar{\mathbf{B}}, \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}}, \dots, \mathbf{C} \bar{\mathbf{A}}^t \bar{\mathbf{B}}]$$



Kernel

Input

Output



$$y_3 = \mathbf{C} \bar{\mathbf{A}}^3 \bar{\mathbf{B}} \mathbf{x}_0 + \mathbf{C} \bar{\mathbf{A}}^2 \bar{\mathbf{B}} \mathbf{x}_1 + \mathbf{C} \bar{\mathbf{A}} \bar{\mathbf{B}} \mathbf{x}_2 + \mathbf{C} \bar{\mathbf{B}} \mathbf{x}_3$$



01 State Space Model

02 Mamba1

03 Vision Mamba

04 Mamba2

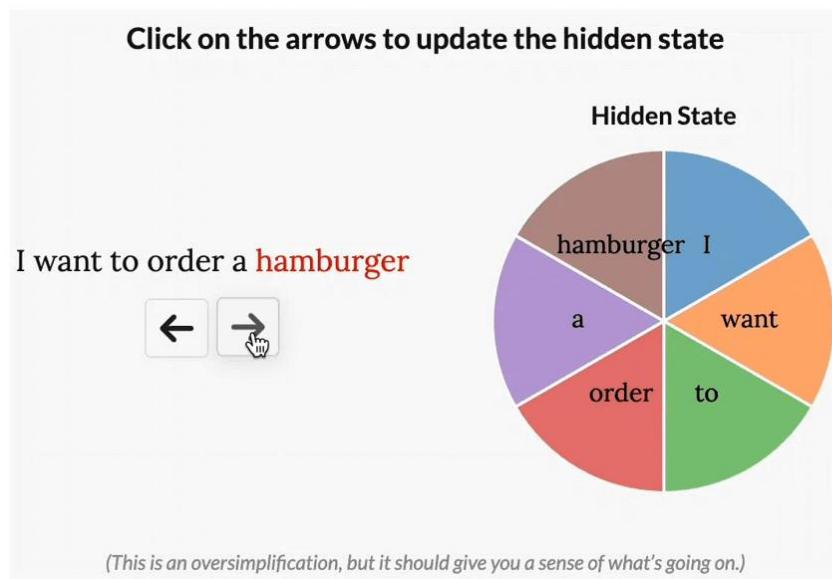
05 PPMA

2. Mamba1: Selectivity

17

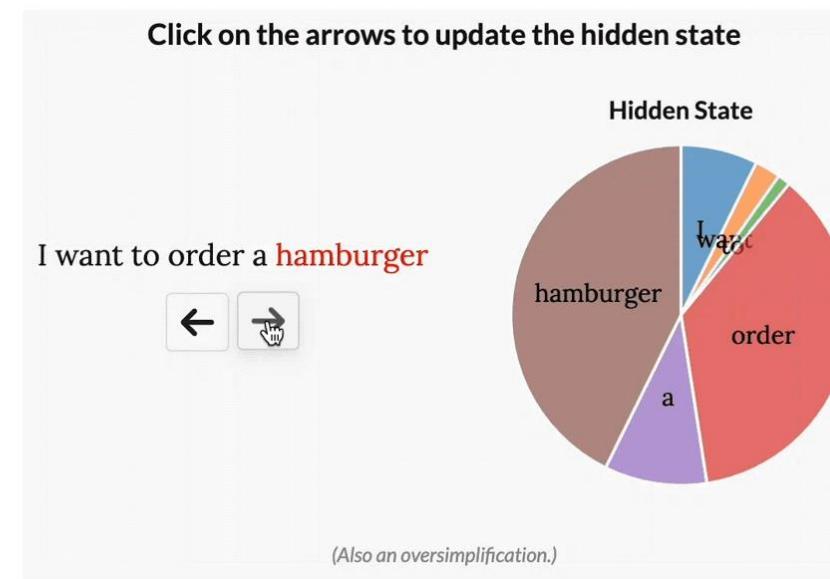
■ Structured State Space Model (S4)

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}} \times \mathbf{h}_{t-1} + \bar{\mathbf{B}} \times \mathbf{x}_t \\ \mathbf{y}_t = \mathbf{C} \times \mathbf{h}_t \end{cases}$$



■ Selective State Space Models (Mamba)

$$\begin{cases} \mathbf{h}_t = \boxed{\bar{\mathbf{A}}_t} \times \mathbf{h}_{t-1} + \boxed{\bar{\mathbf{B}}_t} \times \mathbf{x}_t \\ \mathbf{y}_t = \boxed{\mathbf{C}_t} \times \mathbf{h}_t \end{cases}$$



Selectivity: 动态权重，对不同的输入赋以不同的权重

2. Mamba1: Selectivity

18

■ Selective State Space Models (Mamba)

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}}_t \times \mathbf{h}_{t-1} + \bar{\mathbf{B}}_t \times x_t \\ y_t = \mathbf{C}_t \times \mathbf{h}_t \end{cases}$$

网络的改变

$$\begin{aligned} B &= nn.\text{Parameter}() \rightarrow B_t = nn.\text{Linear}(x_t) \\ C &= nn.\text{Parameter}() \rightarrow C_t = nn.\text{Linear}(x_t) \\ \Delta &= nn.\text{Parameter}() \rightarrow \Delta_t = nn.\text{Linear}(x_t) \end{aligned}$$

Algorithm 1 SSM (S4)

Input: $x : (\mathcal{B}, \mathcal{L}, \mathcal{D})$

Output: $y : (\mathcal{B}, \mathcal{L}, \mathcal{D})$

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

▷ Represents structured $N \times N$ matrix

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

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

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

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

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

▷ Time-invariant: recurrence or convolution

7: **return** y

Algorithm 2 SSM + Selection (S6)

Input: $x : (\mathcal{B}, \mathcal{L}, \mathcal{D})$

Output: $y : (\mathcal{B}, \mathcal{L}, \mathcal{D})$

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

▷ Represents structured $N \times N$ matrix

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

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

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

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

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

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

7: **return** y

2. Mamba1: Selectivity

■ Selective State Space Models (Mamba)

- Input tokens $x = [x_0 \quad x_1 \quad \cdots \quad x_{N-1}] \in \mathbb{R}^{N \times C}$, output tokens $y \in \mathbb{R}^{N \times C}$ and hidden states $h \in \mathbb{R}^{N \times D \times D}$

$$\begin{cases} h_t = \bar{A}_t \times h_{t-1} + \bar{B}_t \times x_t \\ y_t = C_t \times h_t \end{cases}$$

- Learnable input-dependent parameters

$$A = -\exp(nn.\text{Parameter}(torch.diag(D))) \in \mathbb{R}^{D \times D} (< 0)$$

$$B_t = nn.\text{Linear}(x_t) = x_t W_B \in \mathbb{R}^{1 \times D}$$

$$C_t = nn.\text{Linear}(x_t) = x_t W_C \in \mathbb{R}^{1 \times D}$$

$$\Delta_t = \text{Softplus}(\text{MLP}_\Delta(x_t)) \in \mathbb{R}^1 (> 0)$$

$$\bar{A}_t = e^{\Delta_t A} = \exp(\Delta_t A) \in \mathbb{R}^{D \times D} \text{ (0~1)}$$

$$\bar{B}_t = \Delta_t B_t \in \mathbb{R}^{1 \times D}$$

降低多个矩阵乘法 $\bar{A}_t \cdots \bar{A}_1$ 的计算量

$$h_t = \begin{matrix} A & \times & h_{t-1} + B & \times & x_t \end{matrix}$$

(a) General (Unstructured) SSMs

$$h_t = \begin{matrix} A_t & \times & h_{t-1} + B_t & \times & x_t \end{matrix}$$

(b) Mamba1 (Structured SSMs)

2. Mamba1: Parallelizable Training

20

■ 深度学习模型的计算效率

- 训练推理速度 Throughput (images/s) $\propto \frac{\text{并行度} \times \text{GPU算力}}{\text{计算复杂度}}$

训练推理时间/速度	基本概念	相关因素
计算复杂度 / 浮点运算次数 (FLOPs)	一次前向传播所需的总计算次数 (浮点数乘法加法的总次数)	模型的结构、输入数据的尺寸
并行度	可并发计算的单元数	模型的结构、进程/线程数、硬件架构 (CUDA 核心数)

■ 存储空间

存储空间	基本概念	相关因素
显存/内存占用	模型在运行过程中占用的内存资源	模型参数、前向激活值、反向传播的梯度等
模型存储占用	模型储存所占用的内存资源	模型参数量、参数的数据精度类型 (FP32/FP16)

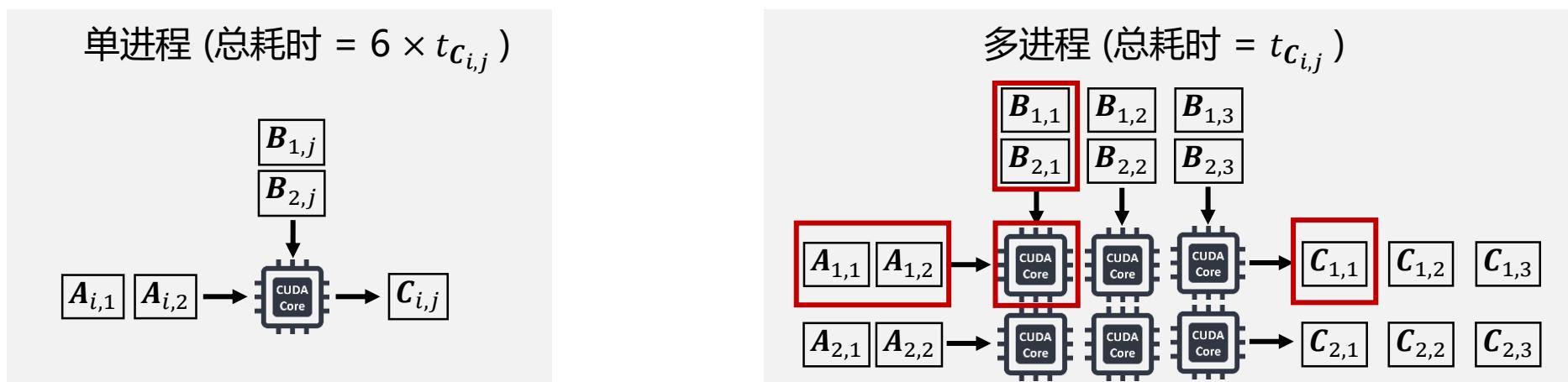
2. Mamba1: Parallelizable Training

■ 深度学习模型的计算效率

- 训练推理速度 Throughput (images/s): 前向计算时间 $\propto \frac{\text{计算复杂度}}{\text{并行度} \times \text{GPU算力}}$

训练推理时间/速度	基本概念	相关因素
计算复杂度 / 浮点运算次数 (FLOPs)	一次前向传播所需的总计算次数 (浮点数乘法加法的总次数)	模型的结构、输入数据的尺寸
并行度	可并发计算的单元数	模型的结构、进程/线程数、硬件架构 (CUDA 核心数)

- 例如，矩阵乘法可利用的并行计算极大加速: $C = A \times B, A \in \mathbb{R}^{2 \times 2}, B \in \mathbb{R}^{2 \times 3}$ (FLOPs=12)



2. Mamba1: Parallelizable Training

- Selective State Space Models (Mamba)

- Mamba1 的循环递推形式 (串行计算, 速度慢)

$$\begin{cases} \mathbf{h}_t = \bar{\mathbf{A}}_t \times \mathbf{h}_{t-1} + \bar{\mathbf{B}}_t \times \mathbf{x}_t \\ \mathbf{y}_t = \mathbf{C}_t \times \mathbf{h}_t \end{cases}$$

- Mamba1 的累计和形式 (并行计算, 速度快)

$$\mathbf{y}_0 = \mathbf{C}_0 \bar{\mathbf{B}}_0 \mathbf{x}_0$$

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

$$\mathbf{y}_2 = \mathbf{C}_2 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1 \bar{\mathbf{B}}_0 \mathbf{x}_0 + \mathbf{C}_2 \bar{\mathbf{A}}_2 \bar{\mathbf{B}}_1 \mathbf{x}_1 + \mathbf{C}_2 \bar{\mathbf{B}}_2 \mathbf{x}_2$$

...

$$\mathbf{y}_t = \mathbf{C}_t (\bar{\mathbf{A}}_t \cdots \bar{\mathbf{A}}_1) \bar{\mathbf{B}}_0 \mathbf{x}_0 + \mathbf{C}_t (\bar{\mathbf{A}}_t \cdots \bar{\mathbf{A}}_2) \bar{\mathbf{B}}_1 \mathbf{x}_1 + \cdots + \mathbf{C}_t \bar{\mathbf{A}}_t \bar{\mathbf{B}}_{t-1} \mathbf{x}_{t-1} + \mathbf{C}_t \bar{\mathbf{B}}_t \mathbf{x}_t$$

$$y_0 = \mathcal{C}_0 \bar{\mathbf{B}}_0 x_0$$

$$y_1 = \mathcal{C}_1 \bar{\mathbf{A}}_1 \bar{\mathbf{B}}_0 x_0 + \mathcal{C}_1 \bar{\mathbf{B}}_1 x_1$$

$$y_2 = \mathcal{C}_2 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1 \bar{\mathbf{B}}_0 x_0 + \mathcal{C}_2 \bar{\mathbf{A}}_2 \bar{\mathbf{B}}_1 x_1 + \mathcal{C}_2 \bar{\mathbf{B}}_2 x_2$$

...

$$y_t = \mathcal{C}_t (\bar{\mathbf{A}}_t \cdots \bar{\mathbf{A}}_1) \bar{\mathbf{B}}_0 x_0 + \mathcal{C}_t (\bar{\mathbf{A}}_{t-1} \cdots \bar{\mathbf{A}}_1) \bar{\mathbf{B}}_1 x_1 + \cdots + \mathcal{C}_t \bar{\mathbf{A}}_{t-1} \bar{\mathbf{B}}_{t-1} x_{t-1} + \mathcal{C}_t \bar{\mathbf{B}}_t x_t$$

$$\text{令 } \bar{x}_t = \mathbf{B}_t x_t, \quad y_t = \mathcal{C}_t \bar{y}_t$$

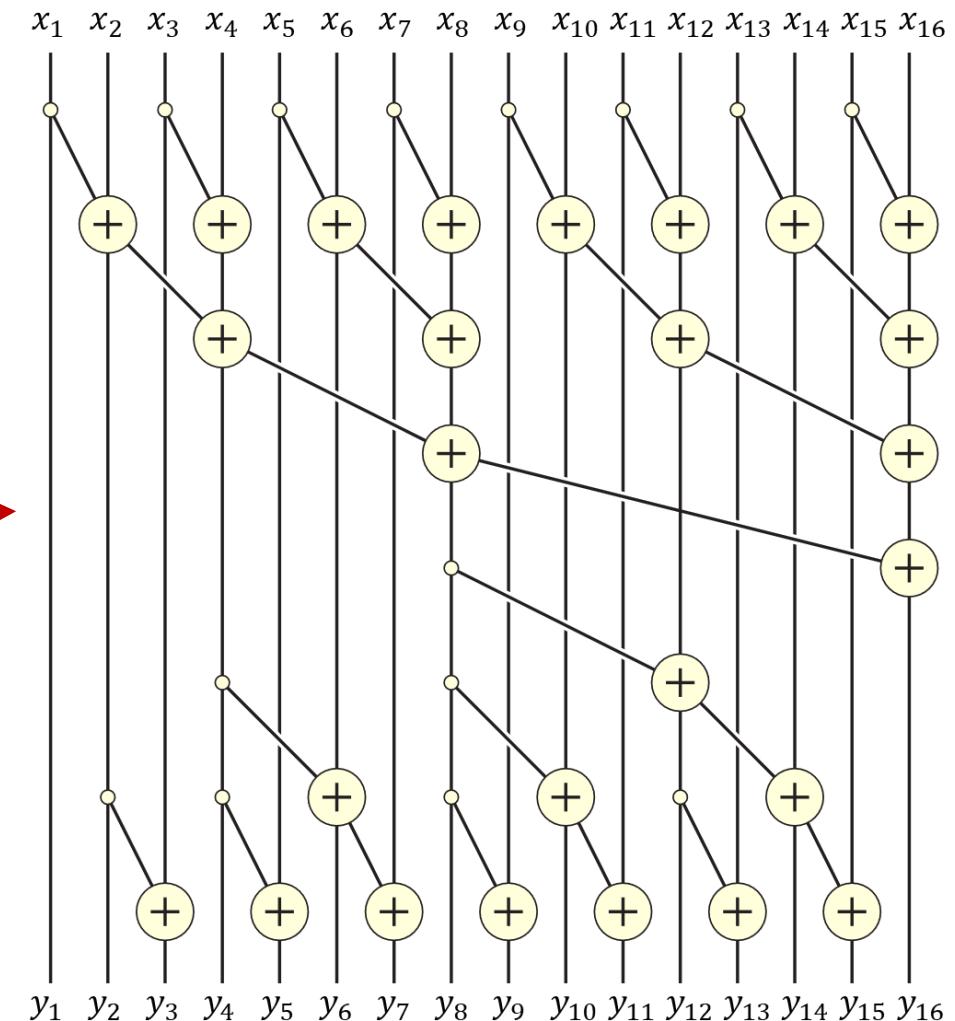
$$\bar{y}_0 = \bar{x}_0$$

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

$$\bar{y}_2 = \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1 \bar{x}_0 + \bar{\mathbf{A}}_2 \bar{x}_1 + \bar{x}_2$$

$$\bar{y}_3 = \bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{\mathbf{A}}_1 \bar{x}_0 + \bar{\mathbf{A}}_3 \bar{\mathbf{A}}_2 \bar{x}_1 + \bar{\mathbf{A}}_3 \bar{x}_2 + \bar{x}_3$$

...



$$y_0 = x_0$$

$$y_1 = A_1 x_0 + x_1$$

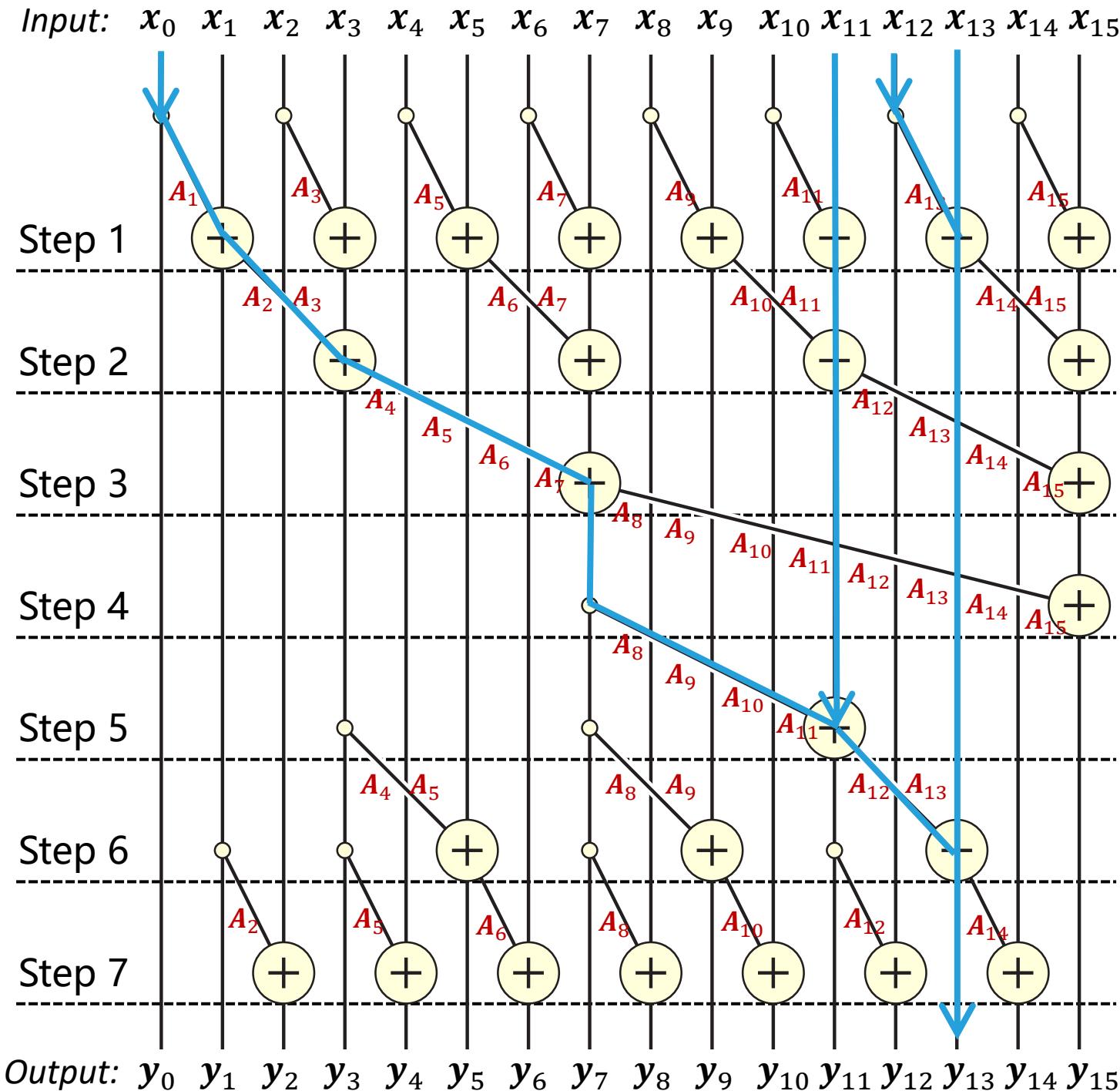
$$y_2 = A_2 A_1 x_0 + A_2 x_1 + x_2$$

$$y_3 = A_3 A_2 A_1 x_0 + A_3 A_2 x_1 + A_3 x_2 + x_3$$

...

■ 累积和 (Cumulative Sum)

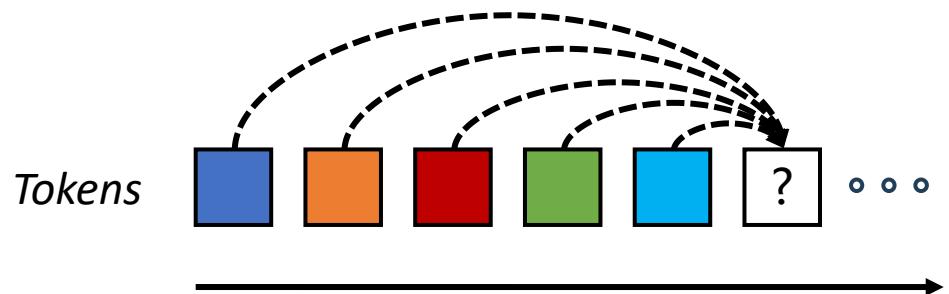
- 串行计算 (从左向右求和)
 - 时间复杂度: $\mathcal{O}(N)$
- 并行计算 (从上向下求和)
 - 时间复杂度: $\mathcal{O}(\log N)$



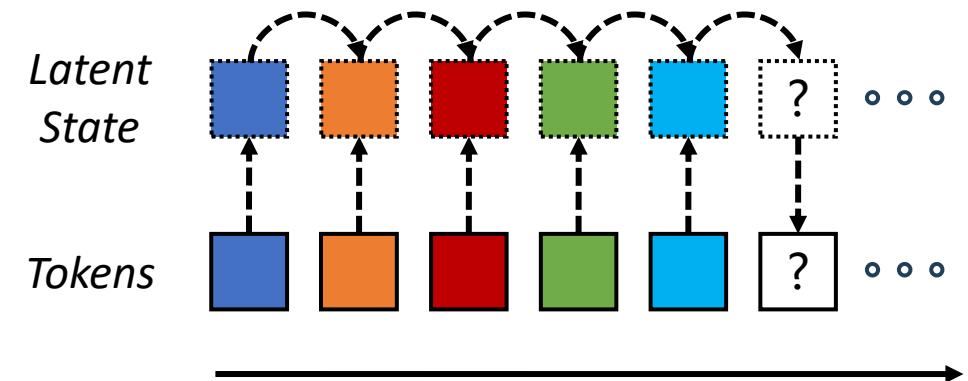
2. Mamba1

25

■ Selective State Space Models (Mamba)



(a) Transformer (complexity: $\mathcal{O}(N^2)$)



(b) State Space Model (complexity: $\mathcal{O}(N)$)

Recurrent(RNN)

- ✓ Global receptive fields
Stateful inference
- ✗ Inefficient training
Vanishing gradients

Convolution (CNN)

- ✓ Easy optimization
Parallelizable training
- ✗ Bounded receptive fields

Transformer

- ✓ Global receptive fields
Parallelizable training
- ✗ Quadratic complexity

SSM (Mamba)

- ✓ Global receptive fields
Dynamic weights
Linear computational
- ✓ Parallelizable training

2. Mamba1: SOTA performance on NLP tasks

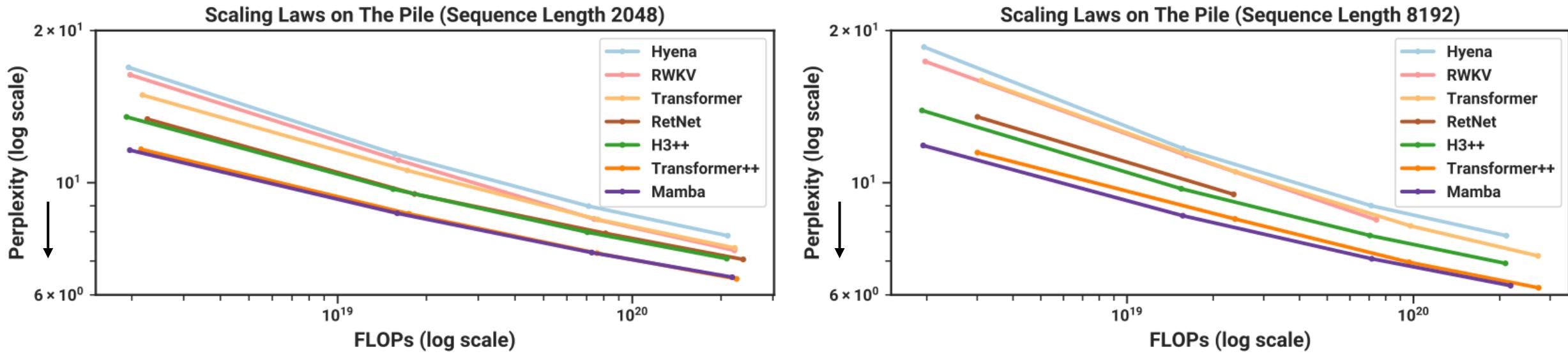


Figure 4: (**Scaling Laws.**) Models of size $\approx 125M$ to $\approx 1.3B$ parameters, trained on the Pile. Mamba scales better than all other attention-free models and is the first to match the performance of a very strong “Transformer++” recipe that has now become standard, particularly as the sequence length grows.

2. Mamba1: SOTA performance on NLP tasks

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



01 State Space Model

02 Mamba1

03 Vision Mamba

04 Mamba2

05 PPMA

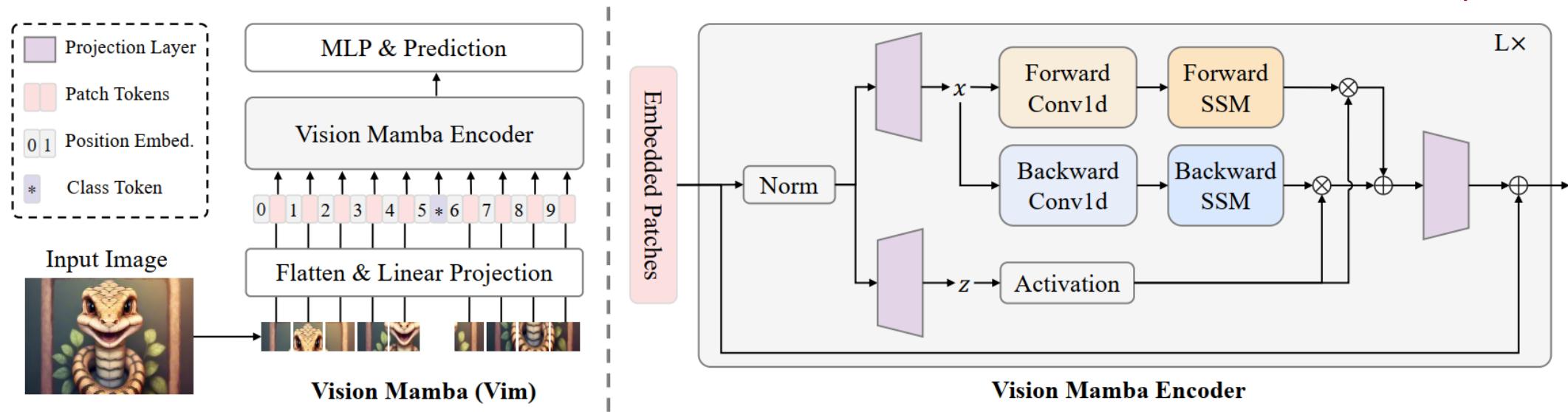
3. Vision Mamba

29

■ ViM: 如何将Mamba应用于视觉任务

- Input tokens: $x \in \mathbb{R}^{N \times C}$; output tokens: $y \in \mathbb{R}^{N \times C}$
- 类似于 ViTs, 将图像裁剪为 patches, 展开成像文本一样的一维序列 输入Mamba block
- 单向扫描 → 双向扫描

First work to introduce
Mamba to computer vision



3. Vision Mamba

30

■ VMamba

- Cross scanning: four scanning paths in different directions

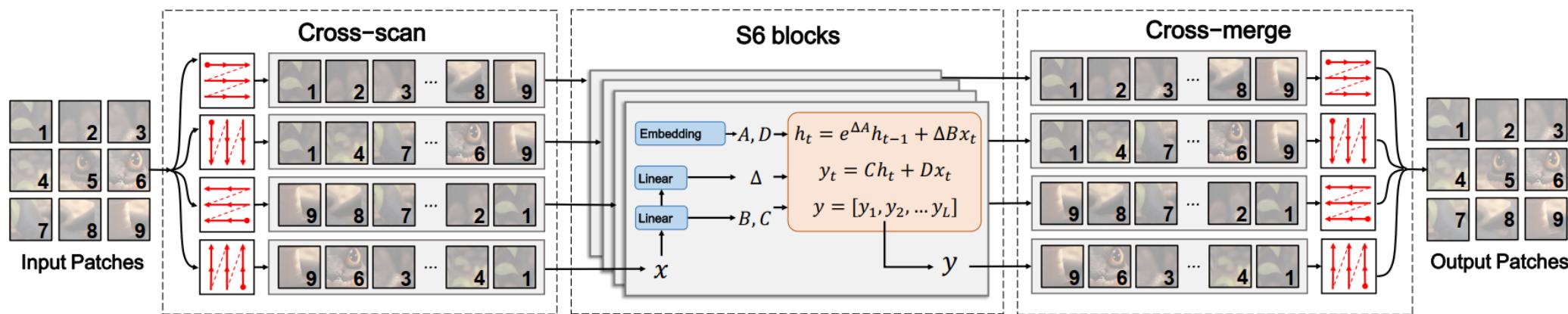
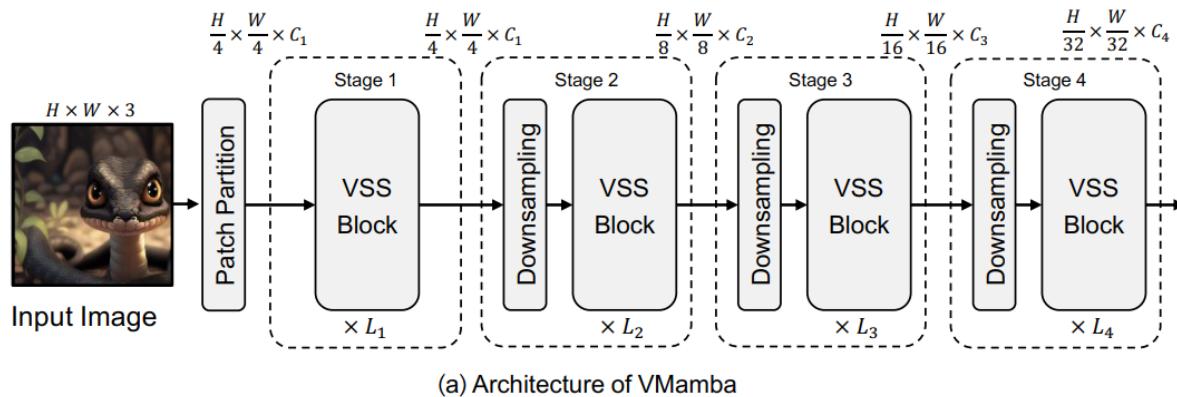


Figure 2: Illustration of 2D-Selective-Scan (SS2D). Input patches are traversed along four different scanning paths (*Cross-Scan*), and each sequence is independently processed by distinct S6 blocks. Subsequently, the results are merged to construct a 2D feature map as the final output (*Cross-Merge*).

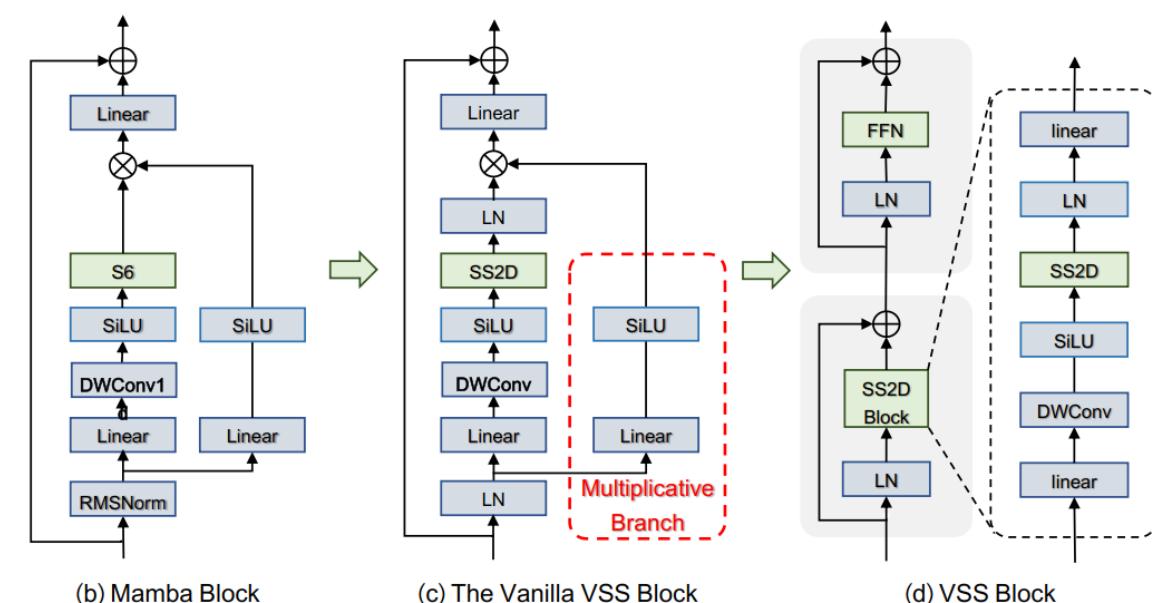
3. Vision Mamba

31

■ VMamba

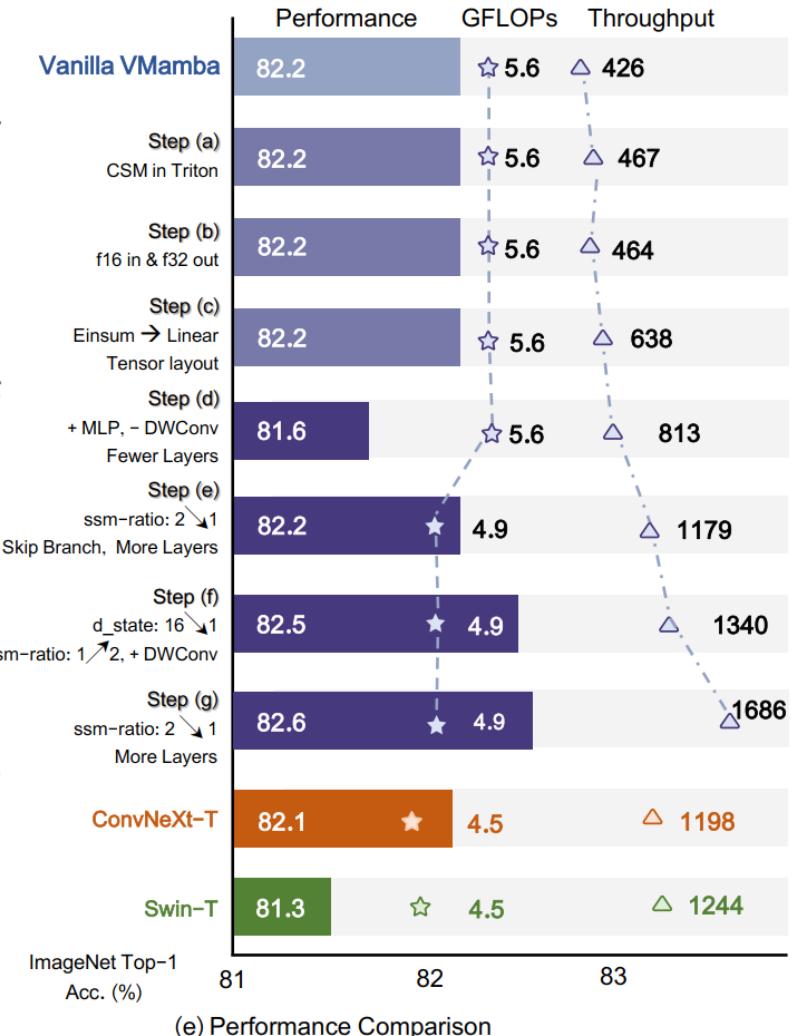


(a) Architecture of VMamba



Architectural

Implementational



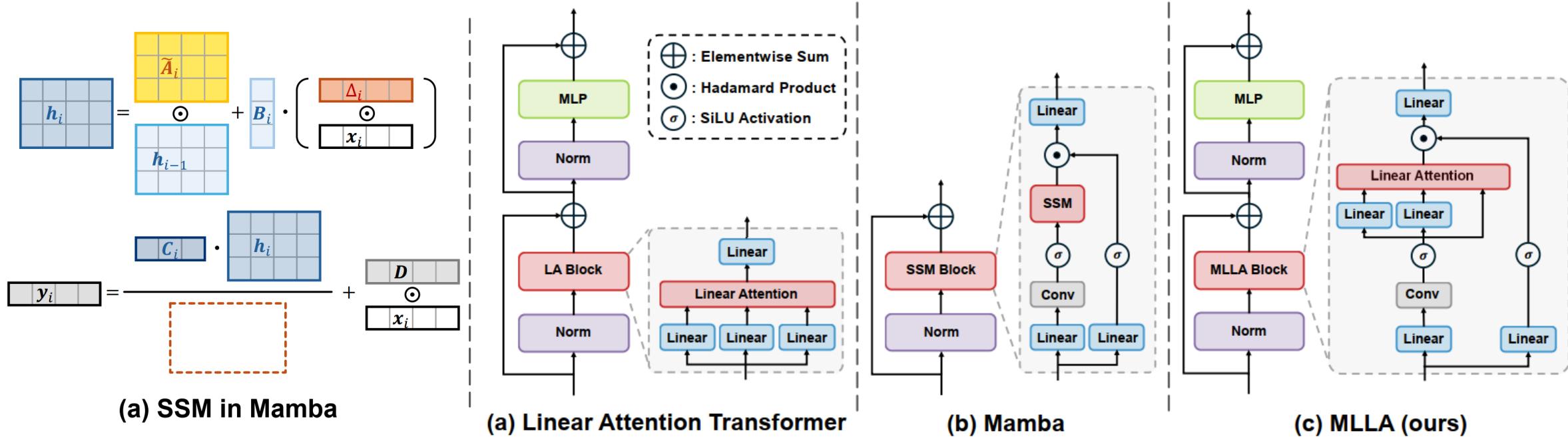
(e) Performance Comparison

3. Vision Mamba

32

■ MLLA: Mamba-Like Linear Attention

- Reveal that Mamba model shares surprising similarities with linear attention Transformer



3. Vision Mamba

33

■ Mamba Vision

- First Mamba-Transformer hybrid architecture for computer vision applications

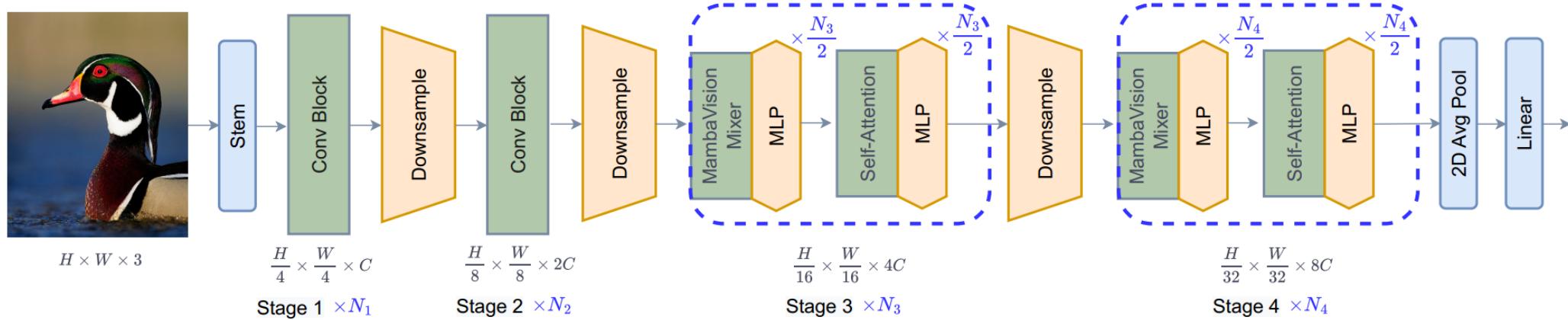


Figure 2 – The architecture of hierarchical MambaVision models. The first two stages use residual convolutional blocks for fast feature extraction. Stage 3 and 4 employ both MambaVision and Transformer blocks. Specifically, given N layers, we use $\frac{N}{2}$ MambaVision and MLP blocks which are followed by additional $\frac{N}{2}$ Transformer and MLP blocks. The Transformer blocks in final layers allow for recovering lost global context and capture long-range spatial dependencies.

3. Vision Mamba

■ MambaIR

- First to work to adapt state space models for low-level image restoration

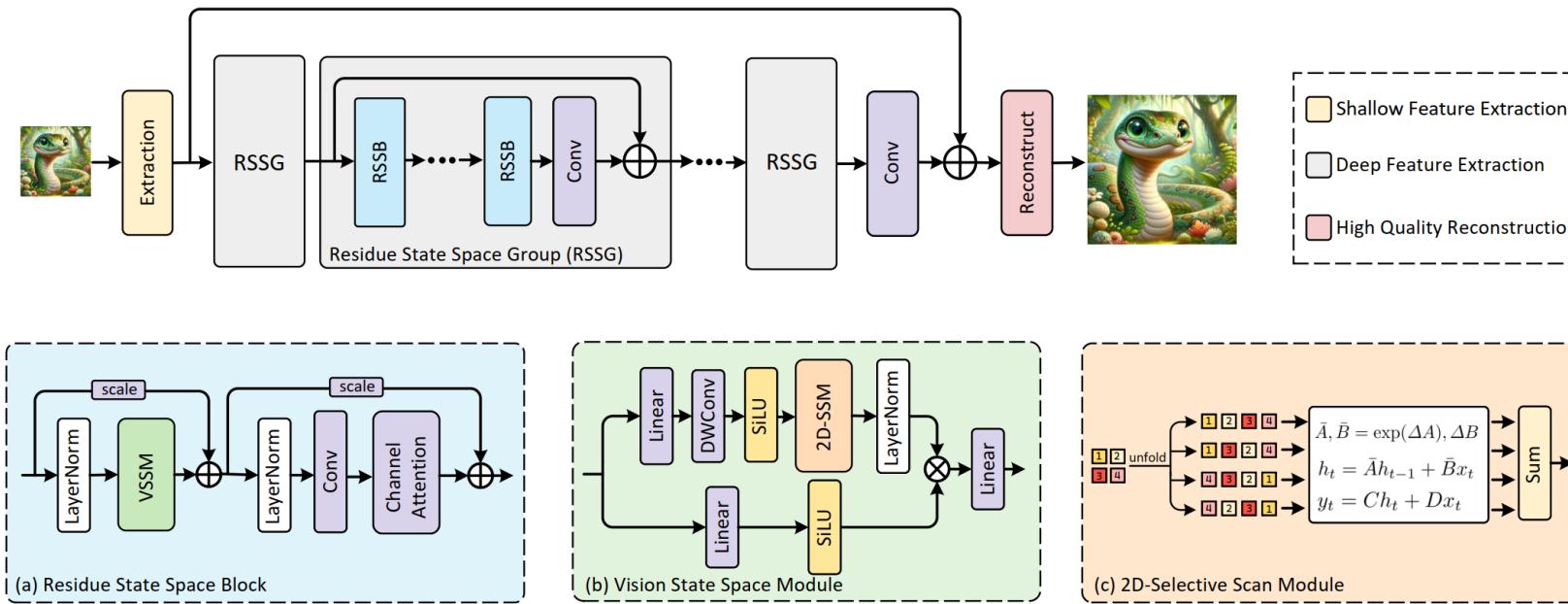
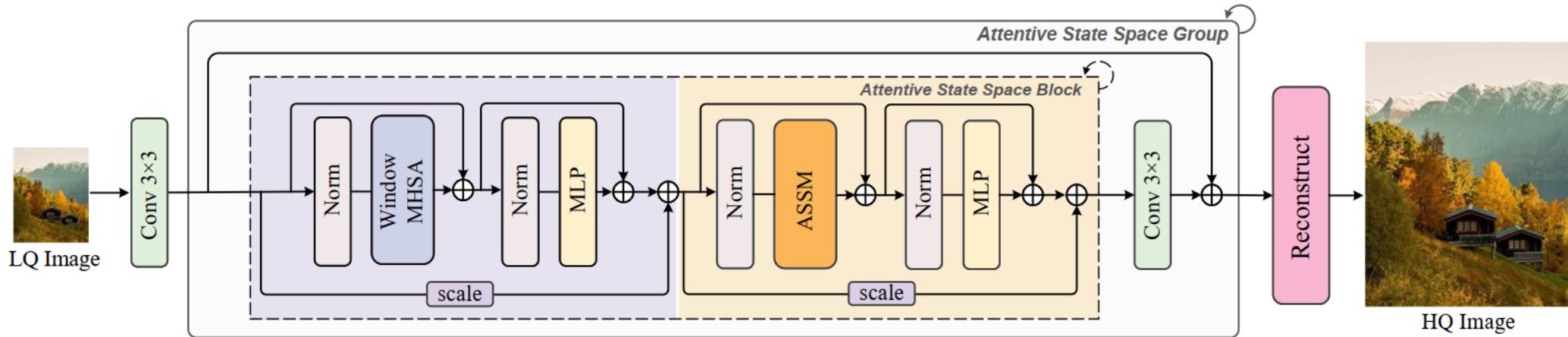


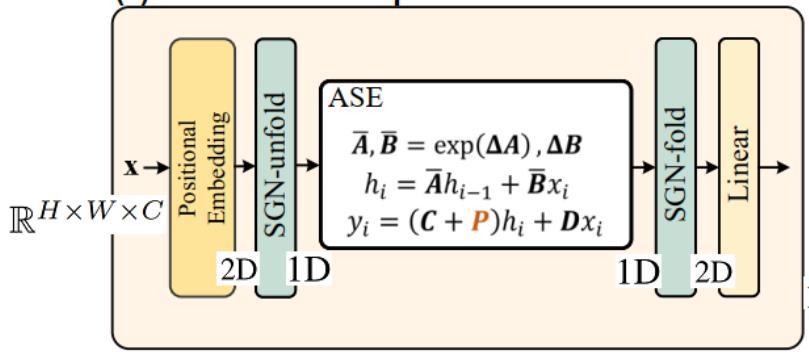
Fig. 2: The overall network architecture of our MambaIR, as well as the (a) Residual State-Space Block (RSSB), the (b) Vision State-Space Module (VSSM), and the (c) 2D Selective Scan Module (2D-SSM).

3. Vision Mamba

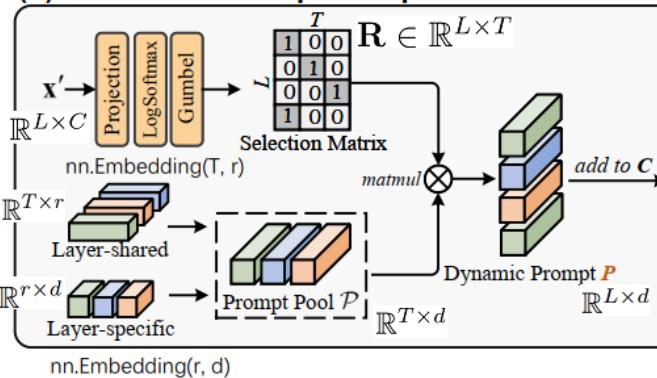
■ MambaIR V2



(a) Attentive State Space Module



(b) Attentive State-space Equation



(c) Semantic Guided Neighboring

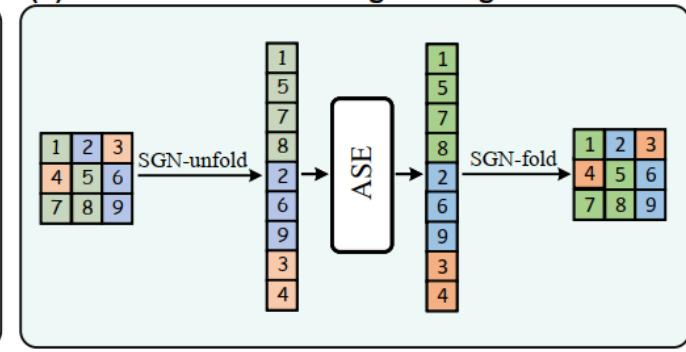


Figure 3. The overall architecture of our proposed MambaIRv2, as well as the (a) Attentive State Space Module (ASSM), (b) Attentive State-space Equation (ASE), and (c) Semantic Guided Neighboring (SGN).

Table 3. Quantitative comparison on ***lightweight image super-resolution*** with state-of-the-art methods. The best and the second best results are in **red** and **blue**.

Method	scale	#param	MACs	Set5		Set14		BSDS100		Urban100		Manga109	
				PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
CARN [1]	2×	1,592K	222.8G	37.76	0.9590	33.52	0.9166	32.09	0.8978	31.92	0.9256	38.36	0.9765
LatticeNet [38]	2×	756K	169.5G	38.13	0.9610	33.78	0.9193	32.25	0.9005	32.43	0.9302	-	-
SwinIR-light [33]	2×	910K	244.2G	38.14	0.9611	33.86	0.9206	32.31	0.9012	32.76	0.9340	39.12	0.9783
MambaIR-light [23]	2×	905K	334.2G	38.13	0.9610	33.95	0.9208	32.31	0.9013	32.85	0.9349	39.20	0.9782
ELAN [63]	2×	621K	203.1G	38.17	0.9611	33.94	0.9207	32.30	0.9012	32.76	0.9340	39.11	0.9782
SRFormer-light [70]	2×	853K	236.3G	38.23	0.9613	33.94	0.9209	32.36	0.9019	32.91	0.9353	39.28	0.9785
MambaIRv2-light	2×	774K	286.3G	38.26	0.9615	34.09	0.9221	32.36	0.9019	33.26	0.9378	39.35	0.9785
CARN [1]	3×	1,592K	118.8G	34.29	0.9255	30.29	0.8407	29.06	0.8034	28.06	0.8493	33.50	0.9440
LatticeNet [38]	3×	765K	76.3G	34.53	0.9281	30.39	0.8424	29.15	0.8059	28.33	0.8538	-	-
SwinIR-light [33]	3×	918K	111.2G	34.62	0.9289	30.54	0.8463	29.20	0.8082	28.66	0.8624	33.98	0.9478
MambaIR-light	3×	913K	148.5G	34.63	0.9288	30.54	0.8459	29.23	0.8084	28.70	0.8631	34.12	0.9479
ELAN [63]	3×	629K	90.1G	34.61	0.9288	30.55	0.8463	29.21	0.8081	28.69	0.8624	34.00	0.9478
SRformer-light [70]	3×	861K	105.4G	34.67	0.9296	30.57	0.8469	29.26	0.8099	28.81	0.8655	34.19	0.9489
MambaIRv2-light	3×	781K	126.7G	34.71	0.9298	30.68	0.8483	29.26	0.8098	29.01	0.8689	34.41	0.9497
CARN [1]	4×	1,592K	90.9G	32.13	0.8937	28.60	0.7806	27.58	0.7349	26.07	0.7837	30.47	0.9084
LatticeNet [38]	4×	777K	43.6G	32.30	0.8962	28.68	0.7830	27.62	0.7367	26.25	0.7873	-	-
SwinIR-light [33]	4×	930K	63.6G	32.44	0.8976	28.77	0.7858	27.69	0.7406	26.47	0.7980	30.92	0.9151
MambaIR-light [23]	4×	924K	84.6G	32.42	0.8977	28.74	0.7847	27.68	0.7400	26.52	0.7983	30.94	0.9135
ELAN [63]	4×	640K	54.1G	32.43	0.8975	28.78	0.7858	27.69	0.7406	26.54	0.7982	30.92	0.9150
SRformer-light [70]	4×	873K	62.8G	32.51	0.8988	28.82	0.7872	27.73	0.7422	26.67	0.8032	31.17	0.9165
MambaIRv2-light	4×	790K	75.6G	32.51	0.8992	28.84	0.7878	27.75	0.7426	26.82	0.8079	31.24	0.9182

3. Vision Mamba

37

■ Video Mamba

- First to work to adapt state space models for video understanding

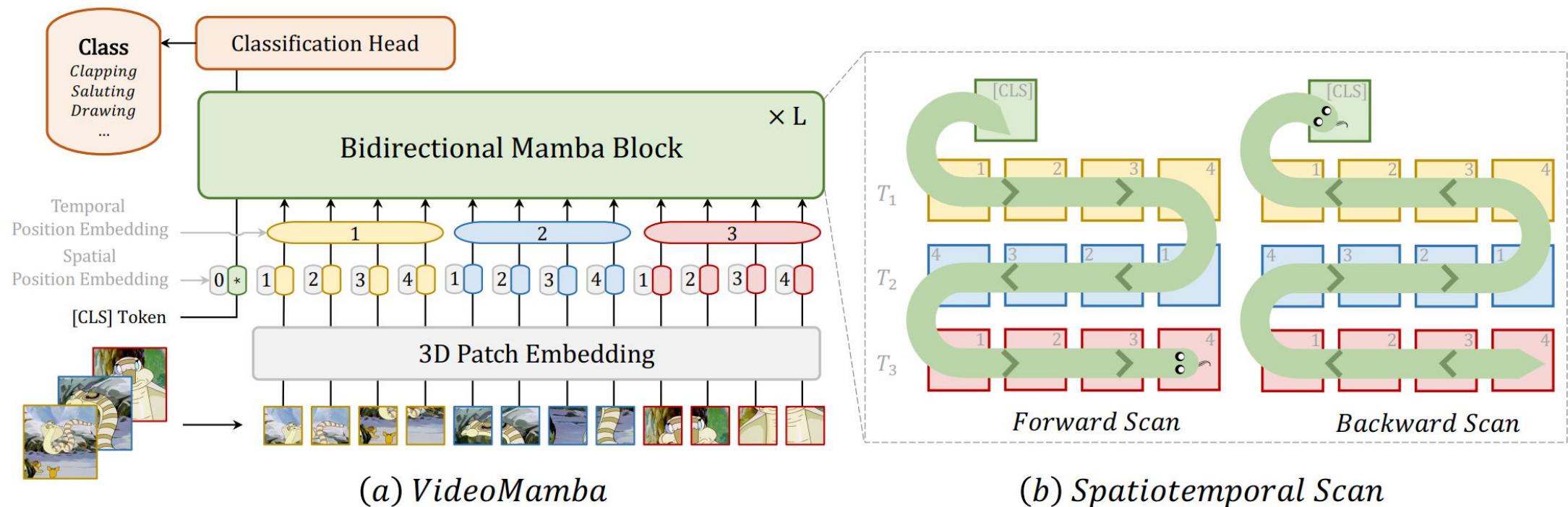


Fig. 3: Framework of VideoMamba. We strictly follow the architecture of vanilla ViT [15], and adapt the bidirectional mamba block [91] for 3D video sequences.

Mamba Models in Computer Vision



Generation and Restoration



Video Processing

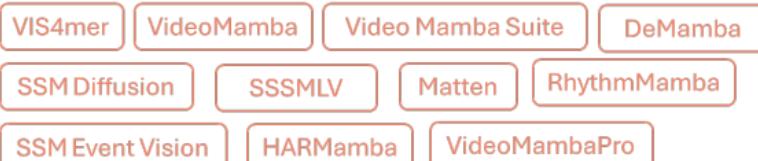
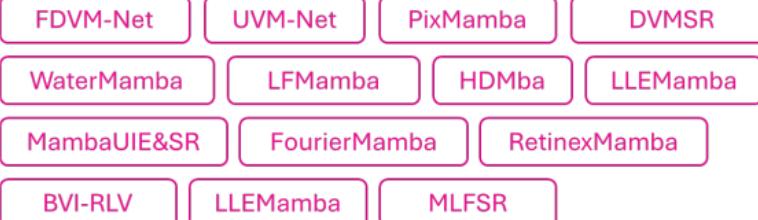


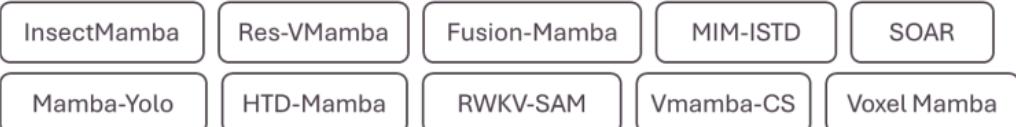
Image Enhancement



Multimodal



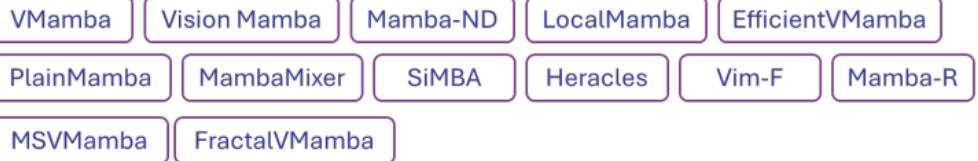
Classification ,Detection and Segmentation



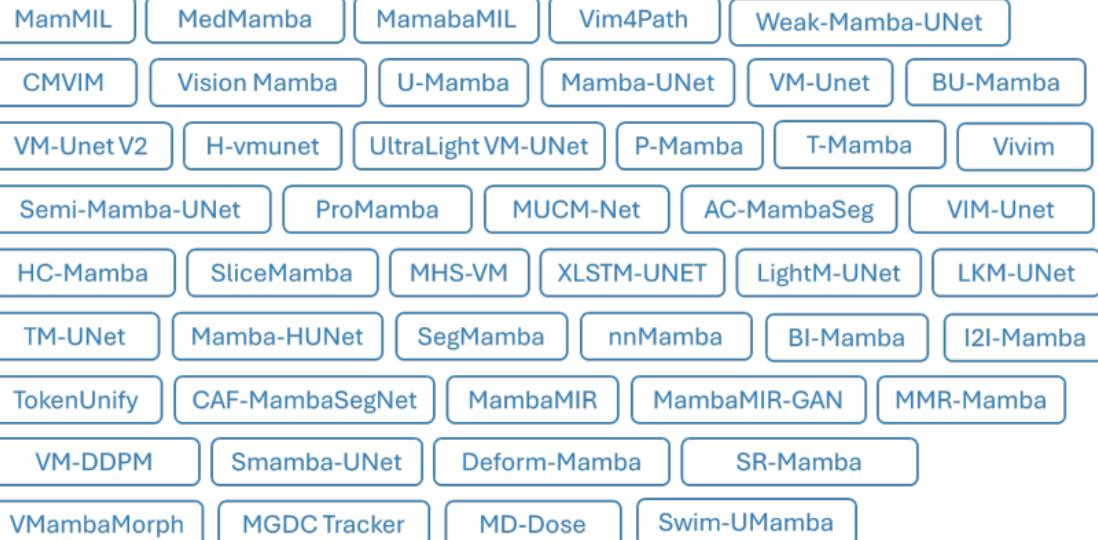
Point Cloud Analysis



General Purpose



Medical Image Analysis

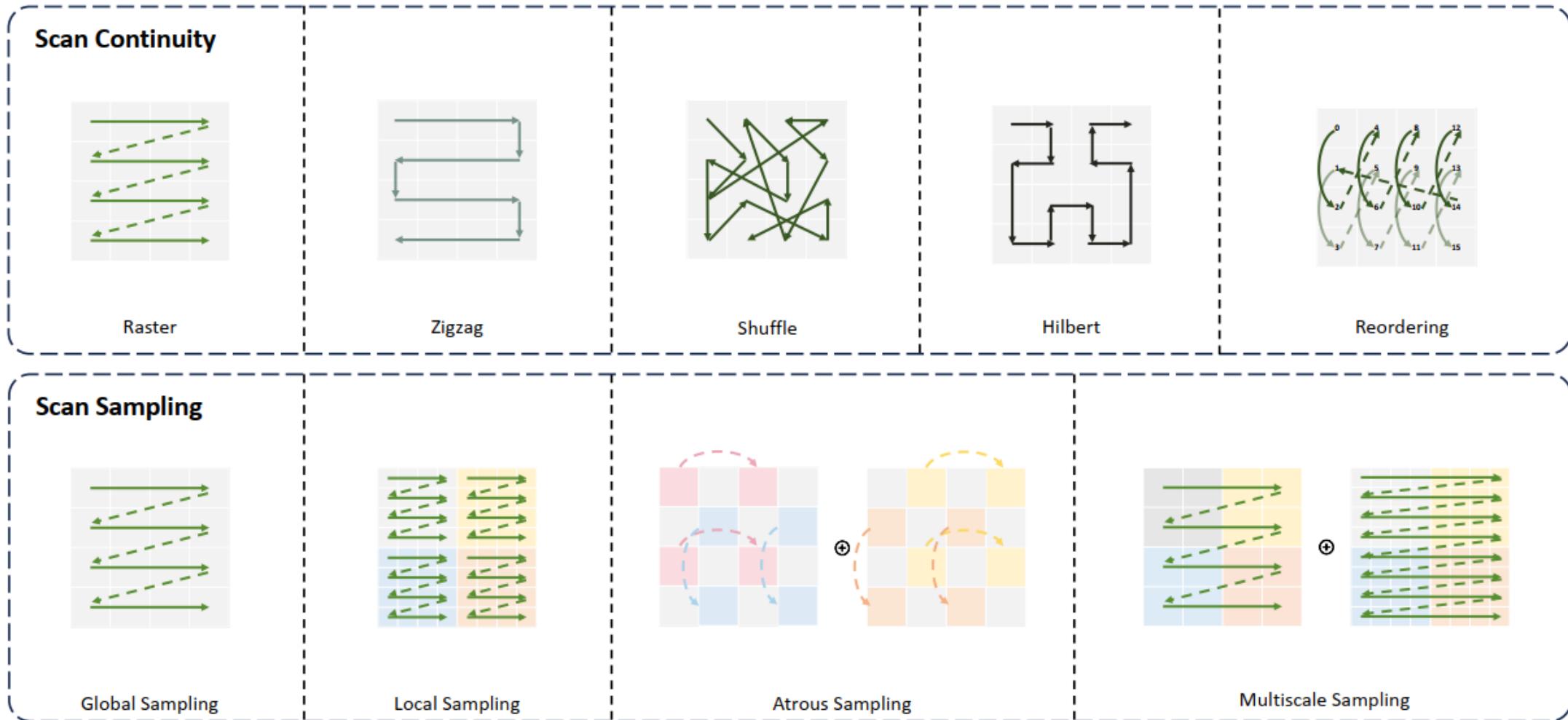


Remote Sensing



3. Vision Mamba

- Key issue: design a proper scanning strategy for 2D images





01 State Space Model

02 Mamba1

03 Vision Mamba

04 Mamba2

05 PPMA

4. Mamba2

41

Insight

Model

Algorithm

Transformers are SSMs: Generalized Models and Efficient Algorithms
Through Structured State Space Duality

Theory

Tri Dao^{*1} and Albert Gu^{*2}

¹Department of Computer Science, Princeton University

²Machine Learning Department, Carnegie Mellon University

tri@tridao.me, agu@cs.cmu.edu

1. Part I - The Model
2. [Part II - The Theory](#)
3. [Part III - The Algorithm](#)
4. [Part IV - The Systems](#)

- 理论层面 (Structured State Space Duality): 统一 Mamba (SSMs) 与 Attention variants (Transformers)
- 算法层面 (Efficient Algorithms): 使用矩阵乘法大幅加速Mamba的训练

[1] Gu A, Dao T. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality [C]. ICML, 2024/05/31. (cited: 575) (<https://tridao.me/blog/2024/mamba2-part1-model/>)

4. Mamba2: Problems in Mamba1

42

■ Problem 1 (Understanding)

- **Mamba范式的理解欠缺**: SSMs 与多种**序列模型范式**有着联系，包括**continuous, convolutional, and recurrent** sequence models

Question 1: What are the **conceptual connections** between **state space models** and **attention**? Can we combine them?

Yes, **SSD Theory**

■ Problem 2 (Efficiency)

- **计算速度慢**: Mamba 的计算效率和硬件适配度仍比注意力机制低 (Small FLOPs but **low training speed**)
- **硬件不适配**: Accelerators such as **GPUs** and **TPUs** are highly specialized for **matrix multiplications**

Question 2: Can we **speed up the training** of Mamba models by recasting them as **matrix multiplications**?

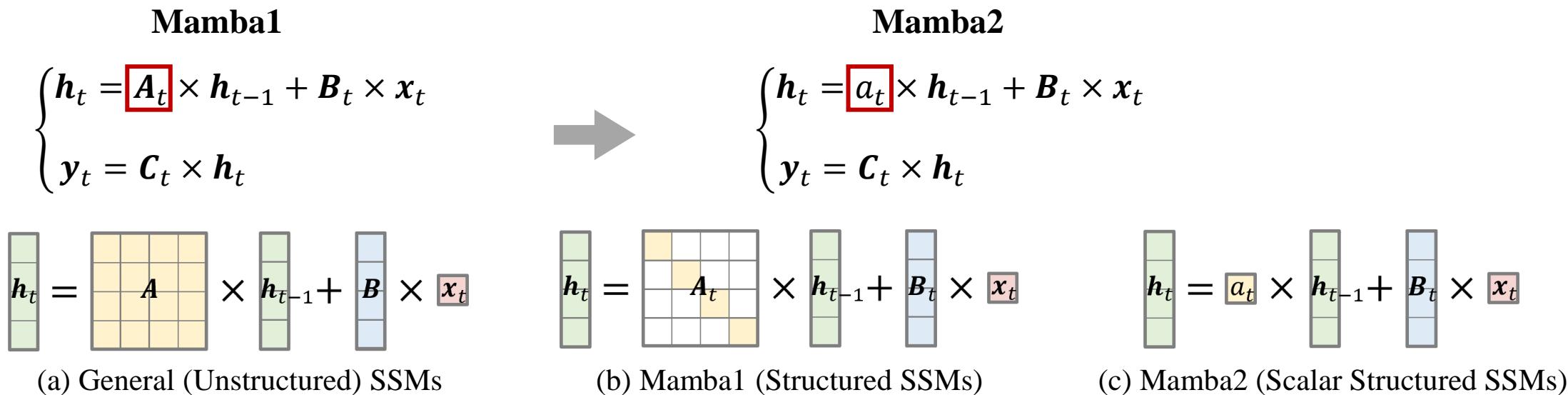
Yes, speed up 2-8x

4.1 Mamba2: The Model

43

■ Mamba2 模型上的改进

- **Scalar Structured SSMs:** 将mamba1中的对角矩阵 A_t 进一步的简化为了一个标量 a_t



where $\mathbf{x} = [\mathbf{x}_0 \quad \mathbf{x}_1 \quad \cdots \quad \mathbf{x}_{N-1}] \in \mathbb{R}^{N \times C}$, $\mathbf{y} \in \mathbb{R}^{N \times C}$, $\mathbf{h}_t \in \mathbb{R}^{D \times C}$; $a_t \in \mathbb{R}$ (**0~1**), $\mathbf{B}_t \in \mathbb{R}^{D \times 1}$, $\mathbf{C}_t \in \mathbb{R}^{1 \times D}$

- **Multihead SSMs:** 将Mamba1中的单头改为多头，将token x_t 沿着通道维度分为多头，多头之间相互独立（类似于多头自注意力）

4.1 Mamba2: The Model

44

■ Mamba2 模型上的改进

- **Scalar Structured SSMs:** 将mamba1中的对角矩阵 A_t 进一步的简化为了一个标量 a_t
- **Multihead SSMs:** 将Mamba1中的单头改为多头，将token x_t 沿着通道维度分为多头，多头之间相互独立（类似于多头自注意力）
- 做这些改进的**原始动机**: 提高计算效率 (推导出对偶的注意力形式以便使用矩阵乘法)
- 改进后的**优势**: Mamba-2 allows much larger state dimensions (from D=16 in Mamba-1 to D=64 or 256 or even higher in Mamba-2)

$$h_t = \begin{matrix} A \\ \hline \end{matrix} \times h_{t-1} + B \times x_t$$

(a) General (Unstructured) SSMs

$$h_t = \begin{matrix} & & A_t & & \\ \hline & & \text{---} & & \\ & & \text{---} & & \\ \hline & & A_t & & \\ \hline & & \text{---} & & \\ & & \text{---} & & \end{matrix} \times h_{t-1} + B_t \times x_t$$

(b) Mamba1 (Structured SSMs)

$$h_t = a_t \times h_{t-1} + B_t \times x_t$$

(c) Mamba2 (Scalar Structured SSMs)

4.1 Mamba2: The Model

- Mamba2 模型上的改进
 - **Scalar Structured SSMs**: 将mamba1中的对角矩阵 A_t 进一步的简化为了一个标量 a_t
 - **Multihead SSMs**: 将Mamba1中的单头改为多头，将token x_t 沿着通道维度分为多头，多头之间相互独立（类似于多头自注意力）
 - 做这些改进的**原始动机**: 提高计算效率（推导出对偶的注意力形式以便使用矩阵乘法）
 - 改进后的**优势**: Mamba-2 allows much larger state dimensions (from D=16 in Mamba-1 to D=64 to N=256 or even higher in Mamba-2)
 - **训练速度**: 2-8× faster during training
 - **模型性能**: Mamba-2 is not strictly better than Mamba-1 (it should generally be on par or better across the board)

4.1 Mamba2: The Model

Table 1: (**Zero-shot Evaluations.**) Best results for each size in bold, second best unlined. 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-2 outperforms Mamba, and generally matches Pythia at twice the model size. Full results in Table 10.

MODEL	TOKEN.	PILE PPL ↓	LAMBADA PPL ↓	LAMBADA ACC ↑	HELLASWAG ACC ↑	PIQA ACC ↑	ARC-E ACC ↑	ARC-C ACC ↑	WINOGRANDE ACC ↑	OPENBOOKQA ACC ↑	AVERAGE ACC ↑
Pythia-1B	NeoX	7.82	7.92	56.1	47.2	70.7	57.0	27.1	53.5	31.4	49.0
Mamba-790M	NeoX	<u>7.33</u>	<u>6.02</u>	62.7	55.1	72.1	61.2	29.5	56.1	<u>34.2</u>	<u>53.0</u>
Mamba-2-780M	NeoX	7.26	5.86	<u>61.7</u>	<u>54.9</u>	<u>72.0</u>	<u>61.0</u>	<u>28.5</u>	<u>60.2</u>	<u>36.2</u>	<u>53.5</u>
Hybrid H3-1.3B	GPT2	—	11.25	49.6	52.6	71.3	59.2	28.1	56.9	34.4	50.3
Pythia-1.4B	NeoX	7.51	6.08	61.7	52.1	71.0	60.5	28.5	57.2	30.8	51.7
RWKV4-1.5B	NeoX	7.70	7.04	56.4	52.5	72.4	60.5	29.4	54.6	34.0	51.4
Mamba-1.4B	NeoX	<u>6.80</u>	<u>5.04</u>	<u>65.0</u>	<u>59.1</u>	74.2	65.5	<u>32.8</u>	61.5	<u>36.4</u>	<u>56.4</u>
Mamba-2-1.3B	NeoX	6.66	5.02	65.7	59.9	<u>73.2</u>	<u>64.3</u>	<u>33.3</u>	<u>60.9</u>	<u>37.8</u>	<u>56.4</u>
Hybrid H3-2.7B	GPT2	—	7.92	55.7	59.7	73.3	65.6	32.3	61.4	33.6	54.5
Pythia-2.8B	NeoX	6.73	5.04	64.7	59.3	74.0	64.1	32.9	59.7	35.2	55.7
RWKV4-3B	NeoX	7.00	5.24	63.9	59.6	73.7	67.8	33.1	59.6	37.0	56.4
Mamba-2.8B	NeoX	<u>6.22</u>	<u>4.23</u>	<u>69.2</u>	<u>66.1</u>	<u>75.2</u>	69.7	<u>36.3</u>	<u>63.5</u>	39.6	<u>59.9</u>
Mamba-2-2.7B	NeoX	6.09	4.10	69.7	66.6	76.4	<u>69.6</u>	<u>36.4</u>	<u>64.0</u>	<u>38.8</u>	<u>60.2</u>

4.2 Mamba2: The Theory

- Mamba2 状态空间方程的两种**对偶形式**

- **循环递推形式** (Recurrent Form: 串行计算, 速度慢)

$$\begin{cases} \mathbf{h}_t = a_t \times \mathbf{h}_{t-1} + \mathbf{B}_t \times \mathbf{x}_t \\ \mathbf{y}_t = \mathbf{C}_t \times \mathbf{h}_t \end{cases}$$

- **注意力形式** (Attention Form: 矩阵乘法, 支持并行计算, 速度快)

$$\mathbf{y}_0 = \mathbf{C}_0 \mathbf{B}_0 \mathbf{x}_0$$

$$\mathbf{y}_1 = \mathbf{C}_1 a_1 \mathbf{B}_0 \mathbf{x}_0 + \mathbf{C}_1 \mathbf{B}_1 \mathbf{x}_1$$

$$\mathbf{y}_2 = \mathbf{C}_2 a_2 a_1 \mathbf{B}_0 \mathbf{x}_0 + \mathbf{C}_2 a_2 \mathbf{B}_1 \mathbf{x}_1 + \mathbf{C}_2 \mathbf{B}_2 \mathbf{x}_2$$

...

$$\mathbf{y}_t = \mathbf{C}_t (a_t \cdots a_1) \mathbf{B}_0 \mathbf{x}_0 + \mathbf{C}_t (a_t \cdots a_2) \mathbf{B}_1 \mathbf{x}_1 + \cdots + \mathbf{C}_t a_t \mathbf{B}_{t-1} \mathbf{x}_{t-1} + \mathbf{C}_t \mathbf{B}_t \mathbf{x}_t$$

- 循环递推形式 (Recurrent Form: 串行计算, 速度慢)

$$\begin{cases} \mathbf{h}_t = a_t \times \mathbf{h}_{t-1} + \mathbf{B}_t \times \mathbf{x}_t \\ \mathbf{y}_t = \mathbf{C}_t \times \mathbf{h}_t \end{cases}$$

- 注意力形式 (Attention Form: 矩阵乘法, 支持并行计算, 速度快)

$$\mathbf{y}_0 = \mathbf{C}_0 \mathbf{B}_0 \mathbf{x}_0$$

$$\mathbf{y}_1 = \mathbf{C}_1 a_1 \mathbf{B}_0 \mathbf{x}_0 + \mathbf{C}_1 \mathbf{B}_1 \mathbf{x}_1$$

$$\mathbf{y}_2 = \mathbf{C}_2 a_2 a_1 \mathbf{B}_0 \mathbf{x}_0 + \mathbf{C}_2 a_2 \mathbf{B}_1 \mathbf{x}_1 + \mathbf{C}_2 \mathbf{B}_2 \mathbf{x}_2$$

...

$$\mathbf{y}_t = \mathbf{C}_t (a_t \cdots a_1) \mathbf{B}_0 \mathbf{x}_0 + \mathbf{C}_t (a_t \cdots a_2) \mathbf{B}_1 \mathbf{x}_1 + \cdots + \mathbf{C}_t a_t \mathbf{B}_{t-1} \mathbf{x}_{t-1} + \mathbf{C}_t \mathbf{B}_t \mathbf{x}_t$$

$$\Rightarrow \begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{C}_0 \mathbf{B}_0^\top & 0 & \cdots & 0 & 0 \\ \mathbf{C}_1 a_1 \mathbf{B}_0^\top & \mathbf{C}_1 A_{1:1} \mathbf{B}_1^\top & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2} (a_{N-2} \cdots a_1) \mathbf{B}_0^\top & \mathbf{C}_{N-2} (a_{N-2} \cdots a_2) \mathbf{B}_1^\top & \cdots & \mathbf{C}_{N-2} \mathbf{B}_{N-2}^\top & 0 \\ \mathbf{C}_{N-1} (a_{N-1} \cdots a_1) \mathbf{B}_0^\top & \mathbf{C}_{N-1} (a_{N-1} \cdots a_2) \mathbf{B}_1^\top & \cdots & \mathbf{C}_{N-1} a_{N-1} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1} \mathbf{B}_{N-1}^\top \end{bmatrix}}_{\mathbb{R}^{N \times N}} \underbrace{\begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{N-2} \\ \mathbf{x}_{N-1} \end{bmatrix}}_{\mathbb{R}^{N \times C}}$$

■ 注意力形式 (Attention Form: 矩阵乘法, 支持并行计算, 速度快)

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \begin{bmatrix} \mathbf{C}_0 \mathbf{B}_0^\top & 0 & \dots & 0 & 0 \\ \mathbf{C}_1 a_1 \mathbf{B}_0^\top & \mathbf{C}_1 A_{1:1} \mathbf{B}_1^\top & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2}(a_{N-2} \cdots a_1) \mathbf{B}_0^\top & \mathbf{C}_{N-2}(a_{N-2} \cdots a_2) \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-2} \mathbf{B}_{N-2}^\top & 0 \\ \mathbf{C}_{N-1}(a_{N-1} \cdots a_1) \mathbf{B}_0^\top & \mathbf{C}_{N-1}(a_{N-1} \cdots a_2) \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-1} a_{N-1} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1} \mathbf{B}_{N-1}^\top \end{bmatrix} \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{N-2} \\ \mathbf{x}_{N-1} \end{bmatrix}$$

■ 令 $A_{i:j} = \begin{cases} a_{i+1} \times \cdots \times a_j & \text{if } i < j \\ 1 & \text{if } i = j \\ a_{j+1} \times \cdots \times a_i & \text{if } i > j \end{cases}$

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \begin{bmatrix} \mathbf{C}_0 A_{0:0} \mathbf{B}_0^\top & 0 & \dots & 0 & 0 \\ \mathbf{C}_1 A_{1:0} \mathbf{B}_0^\top & \mathbf{C}_1 A_{1:1} \mathbf{B}_1^\top & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2} A_{N-2:0} \mathbf{B}_0^\top & \mathbf{C}_{N-2} A_{N-2:1} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-2} A_{N-2:N-2} \mathbf{B}_{N-2}^\top & 0 \\ \mathbf{C}_{N-1} A_{N-1:0} \mathbf{B}_0^\top & \mathbf{C}_{N-1} A_{N-1:1} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-1} A_{N-1:N-2} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1} A_{N-1:N-1} \mathbf{B}_{N-1}^\top \end{bmatrix} \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{N-2} \\ \mathbf{x}_{N-1} \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{C}_0 \mathbf{B}_0^\top & \mathbf{C}_0 \mathbf{B}_1^\top & \dots & \mathbf{C}_0 \mathbf{B}_{N-2}^\top & \mathbf{C}_0 \mathbf{B}_{N-1}^\top \\ \mathbf{C}_1 \mathbf{B}_0^\top & \mathbf{C}_1 \mathbf{B}_1^\top & \dots & \mathbf{C}_1 \mathbf{B}_{N-2}^\top & \mathbf{C}_1 \mathbf{B}_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2} \mathbf{B}_0^\top & \mathbf{C}_{N-2} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-2} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-2} \mathbf{B}_{N-1}^\top \\ \mathbf{C}_{N-1} \mathbf{B}_0^\top & \mathbf{C}_{N-1} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-1} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1} \mathbf{B}_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix} \right) \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{N-2} \\ \mathbf{x}_{N-1} \end{bmatrix}$$

■ The Paradigm of Self-attention

- For input tokens $x = [x_0 \quad x_1 \quad \cdots \quad x_{N-1}] \in \mathbb{R}^{N \times C}$, output tokens $y \in \mathbb{R}^{N \times C}$ is formulated as:

$$Q = \text{nn.Linear}(x) = xW_Q \in \mathbb{R}^{N \times D},$$

$$K = \text{nn.Linear}(x) = xW_K \in \mathbb{R}^{N \times D},$$

$$V = \text{nn.Linear}(x) = xW_V \in \mathbb{R}^{N \times C},$$

$$Y = \text{softmax}(QK^\top) \times V \in \mathbb{R}^{N \times C}$$

Complexity: $\mathcal{O}(N^2)$

■ The Paradigm of Vanilla Linear Attention

$$Y = (QK^\top) \times V = Q \times (K^\top V)$$

Complexity: $\mathcal{O}(N)$

■ The Paradigm of Causal Linear Attention^[1]

$$Y = ((QK^\top) \odot L) \times V = Q \star \text{cumsum}(K \star V) \quad \text{Complexity: } \mathcal{O}(N)$$

$$\mathcal{Z} = K \star V = \text{einsum}[ND, NC \rightarrow NDC](K, V)$$

$$\mathcal{H} = \text{cumsum}(K \star V) = \text{cumsum}(\mathcal{Z}) = [\mathcal{Z}_0 \quad \mathcal{Z}_0 + \mathcal{Z}_1 \quad \cdots \quad \sum_{i=0}^{N-1} \mathcal{Z}_i]$$

$$Y = Q \star \mathcal{H} = \text{einsum}[ND, NDC \rightarrow NC](Q, \mathcal{H})$$

$$L = \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 \\ 1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & \cdots & 1 & 0 \\ 1 & 1 & \cdots & 1 & 1 \end{bmatrix}$$

[1] Katharopoulos A, Vyas A, Pappas N, et al. Transformers are rnns: Fast autoregressive transformers with linear attention[C]. 2020 ICML (cited: 2181)

■ The Paradigm of Self-attention

- For input tokens $x = [x_0 \quad x_1 \quad \cdots \quad x_{N-1}] \in \mathbb{R}^{N \times C}$, output tokens $y \in \mathbb{R}^{N \times C}$ is formulated as:

$$Q = \text{nn. Linear}(x) = xW_Q \in \mathbb{R}^{N \times D},$$

$$K = \text{nn. Linear}(x) = xW_K \in \mathbb{R}^{N \times D},$$

$$V = \text{nn. Linear}(x) = xW_V \in \mathbb{R}^{N \times C},$$

$$Y = \text{softmax}(QK^\top) \times V \in \mathbb{R}^{N \times C}$$

Complexity: $\mathcal{O}(N^2)$

■ The Paradigm of Vanilla Linear Attention

$$Y = (QK^\top) \times V = Q \times (K^\top V)$$

Complexity: $\mathcal{O}(N)$

■ The Paradigm of Causal Linear Attention^[1]

$$Y = ((QK^\top) \odot L) \times V = Q \star \text{cumsum}(K \star V) \quad \text{Complexity: } \mathcal{O}(N)$$

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{N-2} \\ y_{N-1} \end{bmatrix} = \left(\begin{bmatrix} Q_0 K_0^\top & Q_0 K_1^\top & \dots & Q_0 K_{N-2}^\top & Q_0 K_{N-1}^\top \\ Q_1 K_0^\top & Q_1 K_1^\top & \dots & Q_1 K_{N-2}^\top & Q_1 K_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ Q_{N-2} K_0^\top & Q_{N-2} K_1^\top & \dots & Q_{N-2} K_{N-2}^\top & Q_{N-2} K_{N-1}^\top \\ Q_{N-1} K_0^\top & Q_{N-1} K_1^\top & \dots & Q_{N-1} K_{N-2}^\top & Q_{N-1} K_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix} \right) \begin{bmatrix} V_0 \\ V_1 \\ \vdots \\ V_{N-2} \\ V_{N-1} \end{bmatrix}$$

■ Causal Linear Attention

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{Q}_0 \mathbf{K}_0^\top & \mathbf{Q}_0 \mathbf{K}_1^\top & \dots & \mathbf{Q}_0 \mathbf{K}_{N-2}^\top & \mathbf{Q}_0 \mathbf{K}_{N-1}^\top \\ \mathbf{Q}_1 \mathbf{K}_0^\top & \mathbf{Q}_1 \mathbf{K}_1^\top & \dots & \mathbf{Q}_1 \mathbf{K}_{N-2}^\top & \mathbf{Q}_1 \mathbf{K}_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{Q}_{N-2} \mathbf{K}_0^\top & \mathbf{Q}_{N-2} \mathbf{K}_1^\top & \dots & \mathbf{Q}_{N-2} \mathbf{K}_{N-2}^\top & \mathbf{Q}_{N-2} \mathbf{K}_{N-1}^\top \\ \mathbf{Q}_{N-1} \mathbf{K}_0^\top & \mathbf{Q}_{N-1} \mathbf{K}_1^\top & \dots & \mathbf{Q}_{N-1} \mathbf{K}_{N-2}^\top & \mathbf{Q}_{N-1} \mathbf{K}_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix} \right) \begin{bmatrix} \mathbf{V}_0 \\ \mathbf{V}_1 \\ \vdots \\ \mathbf{V}_{N-2} \\ \mathbf{V}_{N-1} \end{bmatrix}$$

$$\mathbf{Q} = \text{nn. Linear}(\mathbf{x}) = \mathbf{x} \mathbf{W}_Q \in \mathbb{R}^{N \times D},$$

$$\mathbf{K} = \text{nn. Linear}(\mathbf{x}) = \mathbf{x} \mathbf{W}_K \in \mathbb{R}^{N \times D},$$

$$\mathbf{V} = \text{nn. Linear}(\mathbf{x}) = \mathbf{x} \mathbf{W}_V \in \mathbb{R}^{N \times C},$$

■ Mamba2 的注意力形式

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{C}_0 \mathbf{B}_0^\top & \mathbf{C}_0 \mathbf{B}_1^\top & \dots & \mathbf{C}_0 \mathbf{B}_{N-2}^\top & \mathbf{C}_0 \mathbf{B}_{N-1}^\top \\ \mathbf{C}_1 \mathbf{B}_0^\top & \mathbf{C}_1 \mathbf{B}_1^\top & \dots & \mathbf{C}_1 \mathbf{B}_{N-2}^\top & \mathbf{C}_1 \mathbf{B}_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2} \mathbf{B}_0^\top & \mathbf{C}_{N-2} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-2} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-2} \mathbf{B}_{N-1}^\top \\ \mathbf{C}_{N-1} \mathbf{B}_0^\top & \mathbf{C}_{N-1} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-1} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1} \mathbf{B}_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix} \right) \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{N-2} \\ \mathbf{x}_{N-1} \end{bmatrix}$$

$$\mathbf{B}_t = \Delta_t \times \text{nn. Linear}(\mathbf{x}_t) = \Delta_t \times \mathbf{x}_t \mathbf{W}_B \in \mathbb{R}^{1 \times D}$$

$$\mathbf{C}_t = \text{nn. Linear}(\mathbf{x}_t) = \mathbf{x}_t \mathbf{W}_C \in \mathbb{R}^{1 \times D}$$

$$a = -\exp(\text{nn. Paramter}(torch.diag(1))) \in \mathbb{R}^1 (< 0)$$

$$\Delta_t = \text{Softplus}(\text{MLP}_\Delta(\mathbf{x}_t)) \in \mathbb{R}^1 (> 0)$$

$$a_t = e^{\Delta_t a} = \exp(\Delta_t a) \in \mathbb{R}^1 (\mathbf{0} \sim \mathbf{1})$$

■ Causal Linear Attention

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{Q}_0 \mathbf{K}_0^\top & \mathbf{Q}_0 \mathbf{K}_1^\top & \dots & \mathbf{Q}_0 \mathbf{K}_{N-2}^\top & \mathbf{Q}_0 \mathbf{K}_{N-1}^\top \\ \mathbf{Q}_1 \mathbf{K}_0^\top & \mathbf{Q}_1 \mathbf{K}_1^\top & \dots & \mathbf{Q}_1 \mathbf{K}_{N-2}^\top & \mathbf{Q}_1 \mathbf{K}_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{Q}_{N-2} \mathbf{K}_0^\top & \mathbf{Q}_{N-2} \mathbf{K}_1^\top & \dots & \mathbf{Q}_{N-2} \mathbf{K}_{N-2}^\top & \mathbf{Q}_{N-2} \mathbf{K}_{N-1}^\top \\ \mathbf{Q}_{N-1} \mathbf{K}_0^\top & \mathbf{Q}_{N-1} \mathbf{K}_1^\top & \dots & \mathbf{Q}_{N-1} \mathbf{K}_{N-2}^\top & \mathbf{Q}_{N-1} \mathbf{K}_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix} \right) \begin{bmatrix} \mathbf{V}_0 \\ \mathbf{V}_1 \\ \vdots \\ \mathbf{V}_{N-2} \\ \mathbf{V}_{N-1} \end{bmatrix}$$

$$\mathbf{Q} = \text{nn. Linear}(\mathbf{x}) = \mathbf{x} \mathbf{W}_Q \in \mathbb{R}^{N \times D},$$

$$\mathbf{K} = \text{nn. Linear}(\mathbf{x}) = \mathbf{x} \mathbf{W}_K \in \mathbb{R}^{N \times D},$$

$$\mathbf{V} = \text{nn. Linear}(\mathbf{x}) = \mathbf{x} \mathbf{W}_V \in \mathbb{R}^{N \times C},$$

■ Mamba2 的注意力形式

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \left(\begin{bmatrix} \mathbf{C}_0 \mathbf{B}_0^\top & \mathbf{C}_0 \mathbf{B}_1^\top & \dots & \mathbf{C}_0 \mathbf{B}_{N-2}^\top & \mathbf{C}_0 \mathbf{B}_{N-1}^\top \\ \mathbf{C}_1 \mathbf{B}_0^\top & \mathbf{C}_1 \mathbf{B}_1^\top & \dots & \mathbf{C}_1 \mathbf{B}_{N-2}^\top & \mathbf{C}_1 \mathbf{B}_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2} \mathbf{B}_0^\top & \mathbf{C}_{N-2} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-2} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-2} \mathbf{B}_{N-1}^\top \\ \mathbf{C}_{N-1} \mathbf{B}_0^\top & \mathbf{C}_{N-1} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-1} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1} \mathbf{B}_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix} \right) \begin{bmatrix} \Delta_0 \mathbf{x}_0 \\ \Delta_1 \mathbf{x}_1 \\ \vdots \\ \Delta_{N-2} \mathbf{x}_{N-2} \\ \Delta_{N-1} \mathbf{x}_{N-1} \end{bmatrix}$$

$$\mathbf{B}_t = \text{nn. Linear}(\mathbf{x}_t) = \mathbf{x}_t \mathbf{W}_B \in \mathbb{R}^{1 \times D}$$

$$\mathbf{C}_t = \text{nn. Linear}(\mathbf{x}_t) = \mathbf{x}_t \mathbf{W}_C \in \mathbb{R}^{1 \times D}$$

$$a = -\exp(\text{nn. Paramter}(\text{torch. diag}(1))) \in \mathbb{R}^1 (< 0)$$

$$\Delta_t = \text{Softplus}(\text{MLP}_\Delta(\mathbf{x}_t)) \in \mathbb{R}^1 (> 0)$$

$$a_t = e^{\Delta_t a} = \exp(\Delta_t a) \in \mathbb{R}^1 (\mathbf{0} \sim \mathbf{1})$$

■ Causal Linear Attention

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{N-2} \\ y_{N-1} \end{bmatrix} = \left(\begin{bmatrix} Q_0 K_0^\top & Q_0 K_1^\top & \dots & Q_0 K_{N-2}^\top & Q_0 K_{N-1}^\top \\ Q_1 K_0^\top & Q_1 K_1^\top & \dots & Q_1 K_{N-2}^\top & Q_1 K_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ Q_{N-2} K_0^\top & Q_{N-2} K_1^\top & \dots & Q_{N-2} K_{N-2}^\top & Q_{N-2} K_{N-1}^\top \\ Q_{N-1} K_0^\top & Q_{N-1} K_1^\top & \dots & Q_{N-1} K_{N-2}^\top & Q_{N-1} K_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix} \right) \begin{bmatrix} V_0 \\ V_1 \\ \vdots \\ V_{N-2} \\ V_{N-1} \end{bmatrix}$$

$\boxed{Q} = \text{nn.Linear}(x) = xW_Q \in \mathbb{R}^{N \times D},$

$\boxed{K} = \text{nn.Linear}(x) = xW_K \in \mathbb{R}^{N \times D},$

$\boxed{V} = \text{nn.Linear}(x) = xW_V \in \mathbb{R}^{N \times C},$

■ Mamba2 的注意力形式

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{N-2} \\ y_{N-1} \end{bmatrix} = \left(\begin{bmatrix} C_0 B_0^\top & C_0 B_1^\top & \dots & C_0 B_{N-2}^\top & C_0 B_{N-1}^\top \\ C_1 B_0^\top & C_1 B_1^\top & \dots & C_1 B_{N-2}^\top & C_1 B_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ C_{N-2} B_0^\top & C_{N-2} B_1^\top & \dots & C_{N-2} B_{N-2}^\top & C_{N-2} B_{N-1}^\top \\ C_{N-1} B_0^\top & C_{N-1} B_1^\top & \dots & C_{N-1} B_{N-2}^\top & C_{N-1} B_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix} \right) \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{N-2} \\ x_{N-1} \end{bmatrix}$$

$\boxed{B_t} = \text{nn.Linear}(x_t) = x_t W_B \in \mathbb{R}^{1 \times D}$

$\boxed{C_t} = \text{nn.Linear}(x_t) = x_t W_C \in \mathbb{R}^{1 \times D}$

$a = -\exp(\text{nn.Paramter}(\text{torch.diag}(1))) \in \mathbb{R}^1 (< 0)$

$\Delta_t = \text{Softplus}(\text{MLP}_\Delta(x_t)) \in \mathbb{R}^1 (> 0)$

$a_t = e^{\Delta_t a} = \exp(\Delta_t a) \in \mathbb{R}^1 (\mathbf{0} \sim \mathbf{1})$

4.2 Mamba2: The Theory

■ Causal Linear Attention

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{N-2} \\ y_{N-1} \end{bmatrix} = \left(\begin{bmatrix} Q_0 K_0^\top & Q_0 K_1^\top & \dots & Q_0 K_{N-2}^\top & Q_0 K_{N-1}^\top \\ Q_1 K_0^\top & Q_1 K_1^\top & \dots & Q_1 K_{N-2}^\top & Q_1 K_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ Q_{N-2} K_0^\top & Q_{N-2} K_1^\top & \dots & Q_{N-2} K_{N-2}^\top & Q_{N-2} K_{N-1}^\top \\ Q_{N-1} K_0^\top & Q_{N-1} K_1^\top & \dots & Q_{N-1} K_{N-2}^\top & Q_{N-1} K_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & \dots & 1 & 0 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix} \right) \begin{bmatrix} V_0 \\ V_1 \\ \vdots \\ V_{N-2} \\ V_{N-1} \end{bmatrix}$$

■ Mamba2 的注意力形式

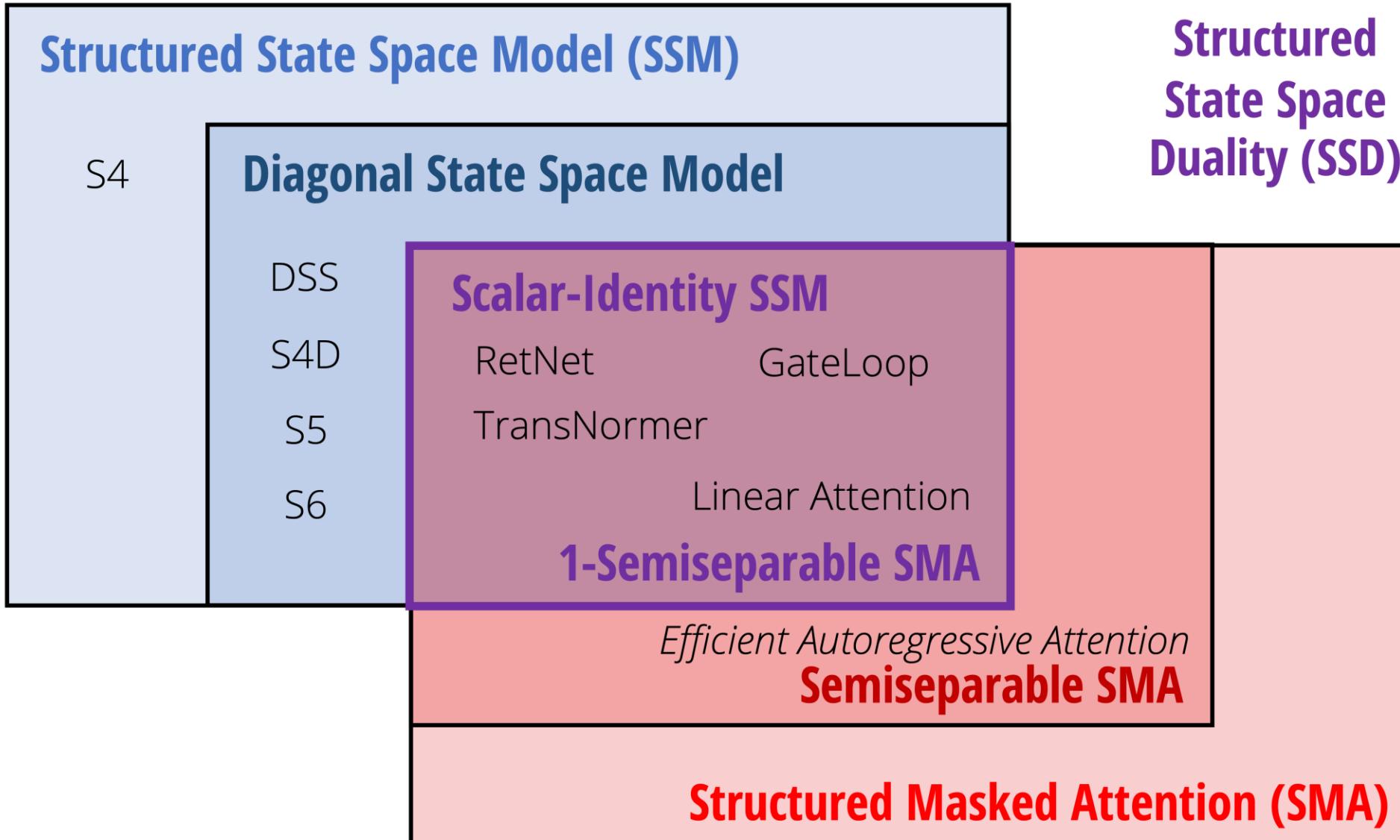
$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{N-2} \\ y_{N-1} \end{bmatrix} = \left(\begin{bmatrix} C_0 B_0^\top & C_0 B_1^\top & \dots & C_0 B_{N-2}^\top & C_0 B_{N-1}^\top \\ C_1 B_0^\top & C_1 B_1^\top & \dots & C_1 B_{N-2}^\top & C_1 B_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ C_{N-2} B_0^\top & C_{N-2} B_1^\top & \dots & C_{N-2} B_{N-2}^\top & C_{N-2} B_{N-1}^\top \\ C_{N-1} B_0^\top & C_{N-1} B_1^\top & \dots & C_{N-1} B_{N-2}^\top & C_{N-1} B_{N-1}^\top \end{bmatrix} \odot \begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix} \right) \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{N-2} \\ x_{N-1} \end{bmatrix}$$

结论：Mamba和Linear Attention本质上是统一的



Theory: Structured State Space Duality (SSD)

Theory: Structured State Space Duality (SSD)



Theory: Structured State Space Duality (SSD)

Structured State Space Model		Structured Masked Attention	
C	(contraction matrix)	Q	(queries)
B	(expansion matrix)	K	(keys)
X	(input sequence)	V	(values)
$A_{j:i}$	(state matrix)	L_{ji}	(mask)
N	(state expansion dim.)	N	(kernel feature dim.)
H $= L \cdot XB$	(hidden states (8b) linear mode)	SMA linear dual (15)	
SSM quadratic dual (16)		G $= Q \cdot K^T$	(Gram matrix (13a)) (quadratic mode)

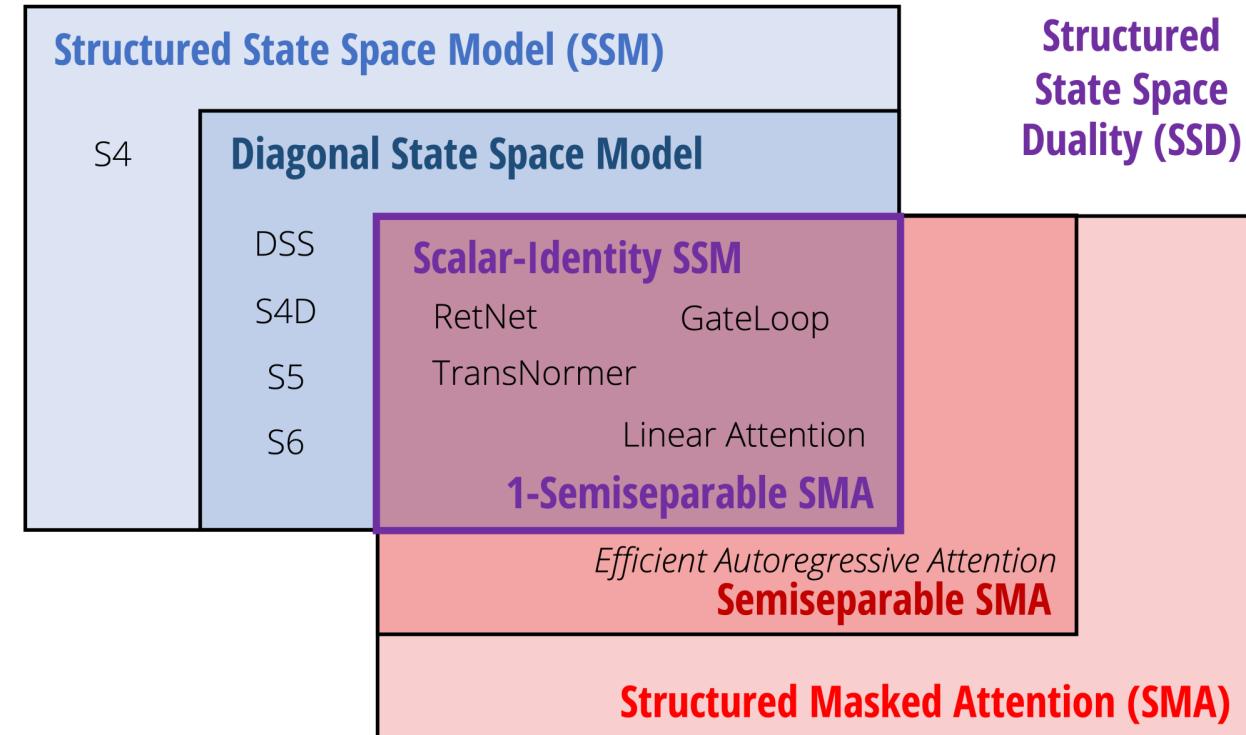
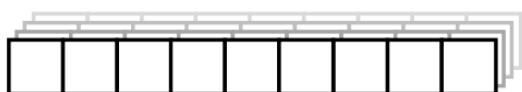


Figure 4: (**Structured State Space Duality**.) State space duality describes the close relationship between state space models and masked attention. (Left) General SSMs and SMA both possess linear and quadratic forms, with direct analogs in notation. (Right) SSMs and SMA intersect at a large class of *state space dual models* (SSD) which capture many sequence models as special cases.

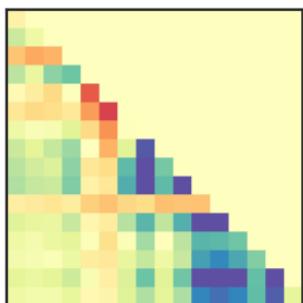
■ Mamba2: Structured Masked Attention

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \underbrace{\left(\begin{array}{cccc} \mathbf{C}_0 \mathbf{B}_0^\top & \mathbf{C}_0 \mathbf{B}_1^\top & \dots & \mathbf{C}_0 \mathbf{B}_{N-2}^\top & \mathbf{C}_0 \mathbf{B}_{N-1}^\top \\ \mathbf{C}_1 \mathbf{B}_0^\top & \mathbf{C}_1 \mathbf{B}_1^\top & \dots & \mathbf{C}_1 \mathbf{B}_{N-2}^\top & \mathbf{C}_1 \mathbf{B}_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2} \mathbf{B}_0^\top & \mathbf{C}_{N-2} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-2} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-2} \mathbf{B}_{N-1}^\top \\ \mathbf{C}_{N-1} \mathbf{B}_0^\top & \mathbf{C}_{N-1} \mathbf{B}_1^\top & \dots & \mathbf{C}_{N-1} \mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1} \mathbf{B}_{N-1}^\top \end{array} \right)}_{\text{Linear Attention}} \odot \underbrace{\begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix}}_{\text{Structured Mask } \mathbf{L}} \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{N-2} \\ \mathbf{x}_{N-1} \end{bmatrix}$$

Outputs Y



Matrix multiplication ↑



Sequence Transformation Matrix M

$$\begin{bmatrix} C_0^\top A_{0:0} B_0 & & & & & & \\ C_1^\top A_{1:0} B_0 & C_1^\top A_{1:1} B_1 & & & & & \\ C_2^\top A_{2:0} B_0 & C_2^\top A_{2:1} B_1 & C_2^\top A_{2:2} B_2 & & & & \\ C_3^\top A_{3:0} B_0 & C_3^\top A_{3:1} B_1 & C_3^\top A_{3:2} B_2 & C_3^\top A_{3:3} B_3 & & & \\ C_4^\top A_{4:0} B_0 & C_4^\top A_{4:1} B_1 & C_4^\top A_{4:2} B_2 & C_4^\top A_{4:3} B_3 & C_4^\top A_{4:4} B_4 & & \\ C_5^\top A_{5:0} B_0 & C_5^\top A_{5:1} B_1 & C_5^\top A_{5:2} B_2 & C_5^\top A_{5:3} B_3 & C_5^\top A_{5:4} B_4 & C_5^\top A_{5:5} B_5 & \\ C_6^\top A_{6:0} B_0 & C_6^\top A_{6:1} B_1 & C_6^\top A_{6:2} B_2 & C_6^\top A_{6:3} B_3 & C_6^\top A_{6:4} B_4 & C_6^\top A_{6:5} B_5 & C_6^\top A_{6:6} B_6 \\ C_7^\top A_{7:0} B_0 & C_7^\top A_{7:1} B_1 & C_7^\top A_{7:2} B_2 & C_7^\top A_{7:3} B_3 & C_7^\top A_{7:4} B_4 & C_7^\top A_{7:5} B_5 & C_7^\top A_{7:6} B_6 & C_7^\top A_{7:7} B_7 \end{bmatrix}$$

Head dim.

Sequence dim. T

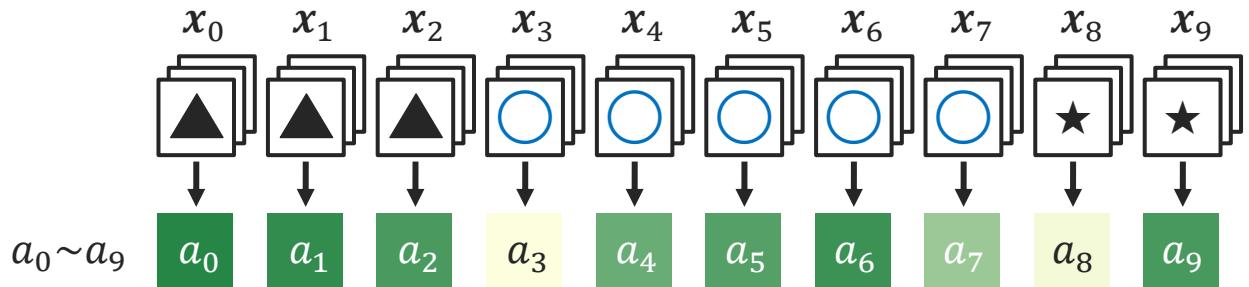


State Space Models are Semiseparable Matrix Transformations

Inputs X

■ Structured Mask L

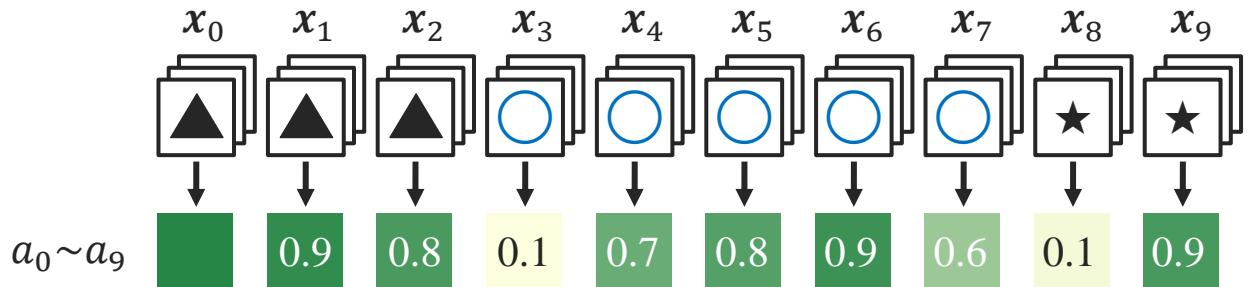
- 1-semiseparable structured matrix
- 衰减因子 a_t 介于 **(0~1)**
- $L_{i,j} = a_{i+1} \times \dots \times a_j$ (即 x_i 到 y_j 途径所有衰减因子的连乘)
- 可以理解为一种 input-dependent relative positional encodings



$$L = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ a_1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{N-2} \cdots a_1 & a_{N-2} \cdots a_2 & \dots & 1 & 0 \\ a_{N-1} \cdots a_1 & a_{N-1} \cdots a_2 & \dots & a_{N-1} & 1 \end{bmatrix}$$

■ Structured Mask L

- 1-semiseparable structured matrix
- 衰减因子 a_t 介于 **(0~1)**
- $L_{i,j} = a_{i+1} \times \dots \times a_j$ (即 x_i 到 y_j 途径所有衰减因子的连乘)
- 可以理解为一种 input-dependent relative positional encodings

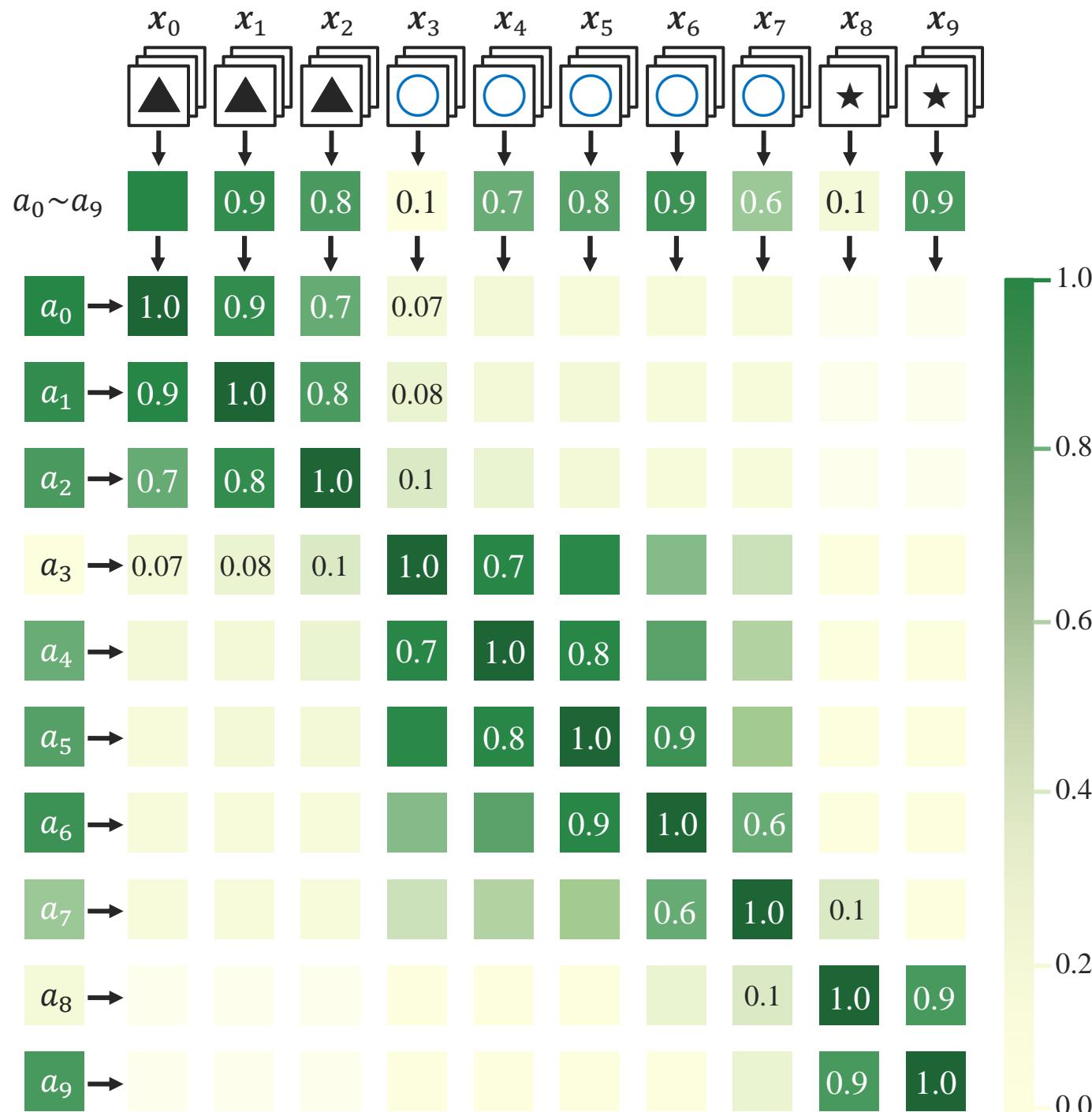


$$L = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ a_1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{N-2} \cdots a_1 & a_{N-2} \cdots a_2 & \dots & 1 & 0 \\ a_{N-1} \cdots a_1 & a_{N-1} \cdots a_2 & \dots & a_{N-1} & 1 \end{bmatrix}$$

■ Structured Mask L

- 1-semiseparable structured matrix
- 衰减因子 a_t 介于 **(0~1)**
- $L_{i,j} = a_{i+1} \times \dots \times a_j$ (即 x_i 到 y_j 途径所有衰减因子的连乘)
- 可以理解为一种 input-dependent relative positional encodings

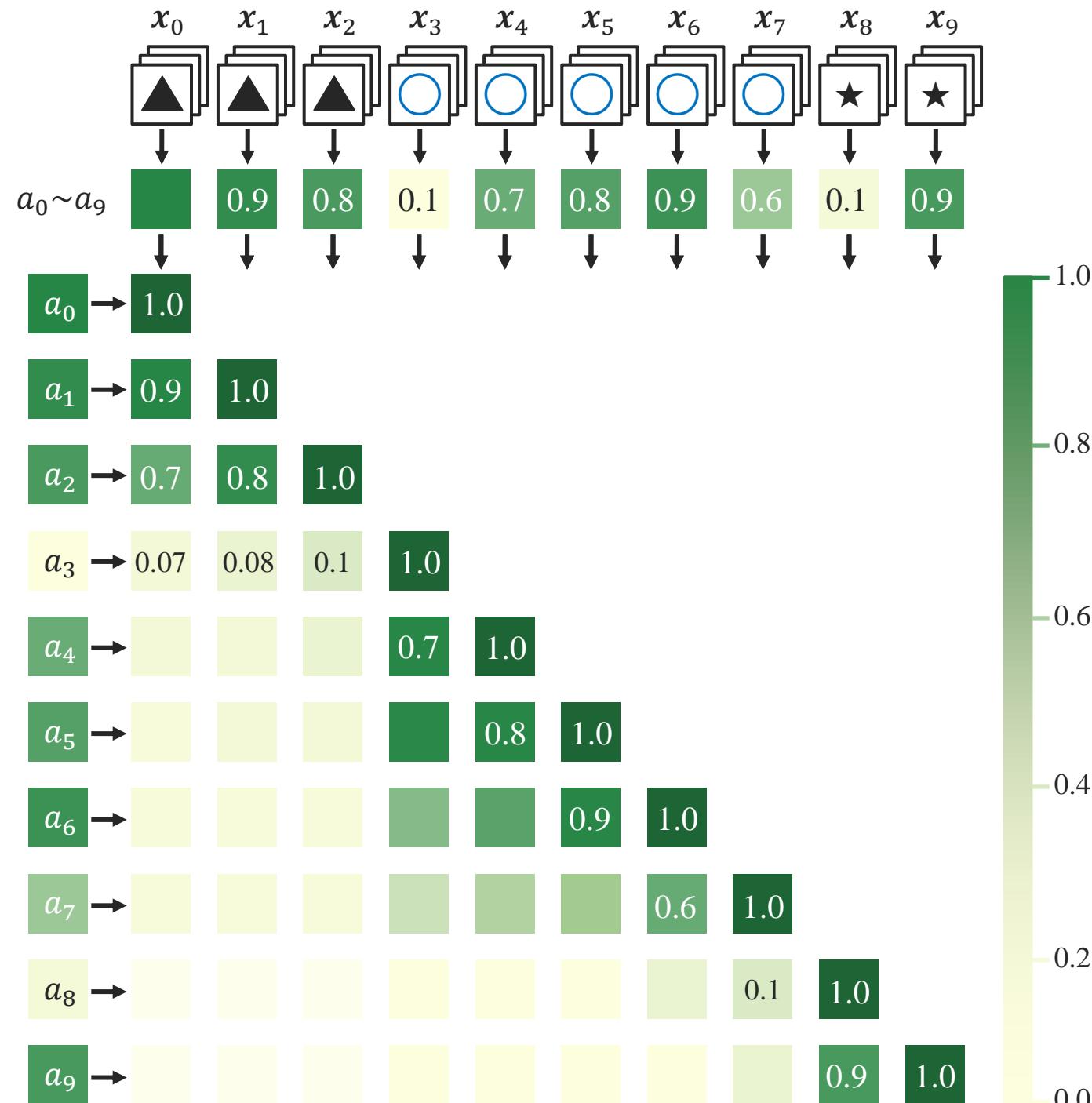
$$L = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ a_1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{N-2} \cdots a_1 & a_{N-2} \cdots a_2 & \dots & 1 & 0 \\ a_{N-1} \cdots a_1 & a_{N-1} \cdots a_2 & \dots & a_{N-1} & 1 \end{bmatrix}$$



■ Structured Mask L

- 1-semiseparable structured matrix
 - 衰减因子 a_t 介于 **(0~1)**
 - $L_{i,j} = a_{i+1} \times \cdots \times a_j$ (即 x_i 到 y_j 途径所有衰减因子的连乘)
 - 可以理解为一种 input-dependent
relative positional encodings

$$L = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ a_1 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{N-2} \cdots a_1 & a_{N-2} \cdots a_2 & \dots & 1 & 0 \\ a_{N-1} \cdots a_1 & a_{N-1} \cdots a_2 & \dots & a_{N-1} & 1 \end{bmatrix}$$



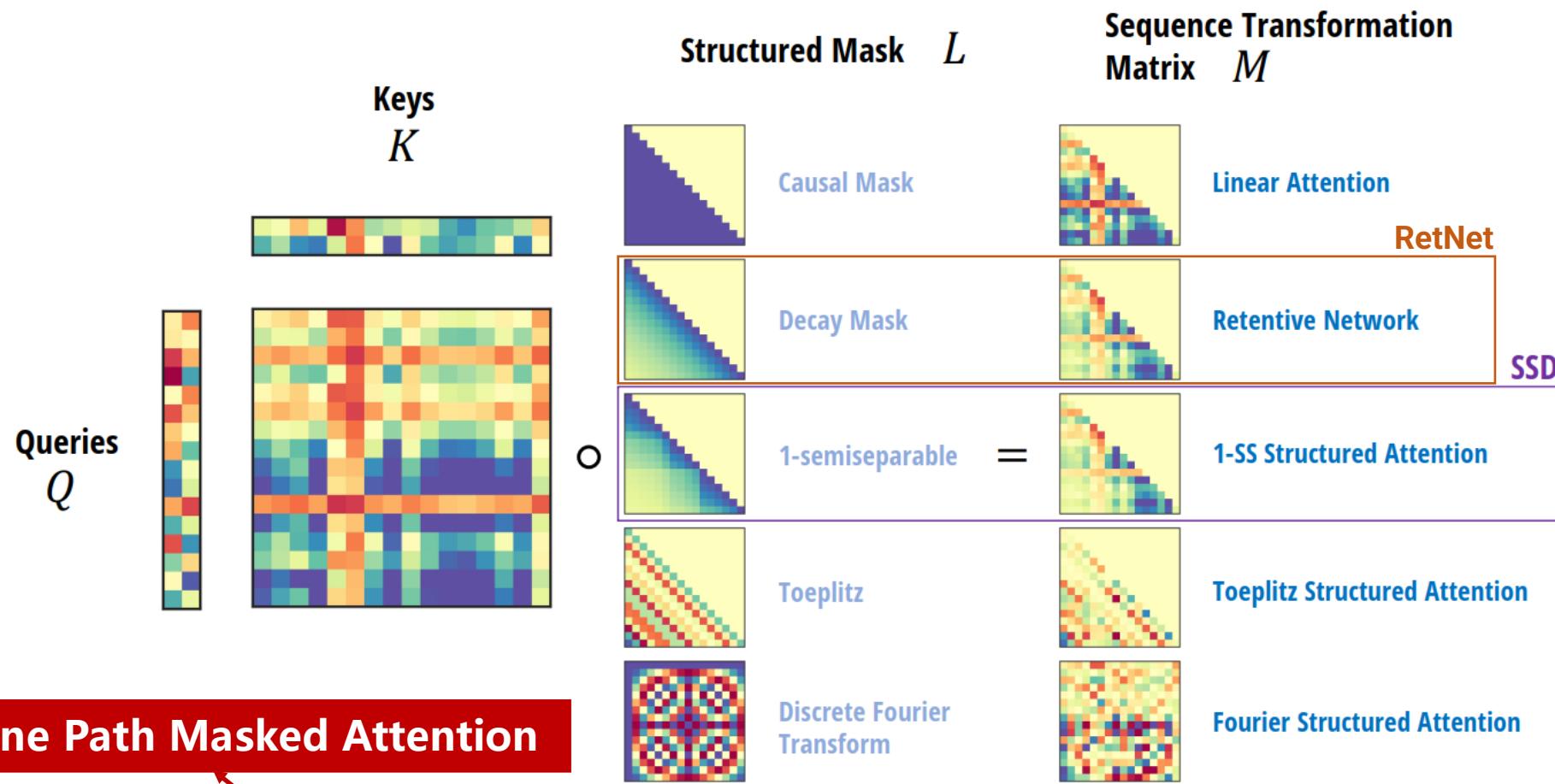


Figure 3: (**Structured Masked Attention**.) SMA constructs a masked attention matrix $M = QK^T \circ L$ for any structured matrix L , which defines a matrix sequence transformation $Y = MV$. All instances of SMA have a dual subquadratic form induced by a different contraction ordering, combined with the efficient structured matrix multiplication by L . Previous examples include Linear Attention (Katharopoulos et al. 2020) and RetNet (Y. Sun et al. 2023). Beyond SSD (1-semiseparable SMA), the focus of this paper, many other potential instantiations of structured attention are possible.

- The decay mask could be generalized to a Toeplitz matrix $L_{ij} = \alpha_{i-j}$ for some learnable (or input-dependent) set of parameters $\alpha \in \mathbb{R}^T$. This can be interpreted as a form of **relative positional encoding**, reminiscent of other methods such as AliBi (Press, N. Smith, and Lewis 2022) but multiplicative instead of additive.

4.3 Mamba2: The Algorithm

64

- Mamba2的两种对偶形式

- 循环递推形式 (串行计算, 速度慢) *Complexity: $\mathcal{O}(N)$*

$$\begin{cases} \mathbf{h}_t = a_t \times \mathbf{h}_{t-1} + \mathbf{B}_t \times \mathbf{x}_t \\ \mathbf{y}_t = \mathbf{C}_t \times \mathbf{h}_t \end{cases}$$

- Structured Masked Attention形式 (并行计算, 矩阵乘法, 速度快) *Complexity: $\mathcal{O}(N^2)$* ×

$$\mathbf{Y} = \mathbf{M}\mathbf{X} = (\mathbf{C}\mathbf{B}^\top \odot \mathbf{L})\mathbf{X}$$

$$\begin{bmatrix} \mathbf{y}_0 \\ \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_{N-2} \\ \mathbf{y}_{N-1} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{C}_0\mathbf{B}_0^\top & \mathbf{C}_0\mathbf{B}_1^\top & \dots & \mathbf{C}_0\mathbf{B}_{N-2}^\top & \mathbf{C}_0\mathbf{B}_{N-1}^\top \\ \mathbf{C}_1\mathbf{B}_0^\top & \mathbf{C}_1\mathbf{B}_1^\top & \dots & \mathbf{C}_1\mathbf{B}_{N-2}^\top & \mathbf{C}_1\mathbf{B}_{N-1}^\top \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{C}_{N-2}\mathbf{B}_0^\top & \mathbf{C}_{N-2}\mathbf{B}_1^\top & \dots & \mathbf{C}_{N-2}\mathbf{B}_{N-2}^\top & \mathbf{C}_{N-2}\mathbf{B}_{N-1}^\top \\ \mathbf{C}_{N-1}\mathbf{B}_0^\top & \mathbf{C}_{N-1}\mathbf{B}_1^\top & \dots & \mathbf{C}_{N-1}\mathbf{B}_{N-2}^\top & \mathbf{C}_{N-1}\mathbf{B}_{N-1}^\top \end{bmatrix}}_{\mathbb{R}^{N \times C}} \odot \underbrace{\begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix}}_{\mathbb{R}^{N \times N}} \underbrace{\begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{N-2} \\ \mathbf{x}_{N-1} \end{bmatrix}}_{\mathbb{R}^{N \times C}}$$

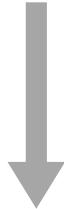
$\downarrow \mathcal{O}(N^2)$ $\downarrow \mathcal{O}(N^2)$

Theorem 1. For any matrices $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{N \times D}$, $\mathbf{L} \in \mathbb{R}^{N \times N}$, $\mathbf{V} \in \mathbb{R}^{N \times C}$, the following equation holds:

$$\mathbf{Z} = \mathbf{K} \star \mathbf{V} = \text{einsum}[ND, NC \rightarrow NDC](\mathbf{K}, \mathbf{V}) \quad \mathcal{O}(NDC)$$

$$\mathbf{Y} = (\mathbf{Q}\mathbf{K}^\top \odot \mathbf{L})\mathbf{V} \iff \mathbf{H} = \text{einsum}[MN, NDC \rightarrow MDC](\mathbf{L}, \mathbf{Z}), \quad (M = N) \quad \mathcal{O}(N^2DC)$$

$$\mathbf{Y} = \mathbf{Q} \star \mathbf{H} = \text{einsum}[MD, MDC \rightarrow MC](\mathbf{Q}, \mathbf{H}) \quad \mathcal{O}(NDC)$$



$$\begin{aligned}
 \text{Proof: } \mathbf{Y}_{m,c} &= \sum_{n=1}^N (\mathbf{L}_{m,n} \sum_{d=1}^D \mathbf{Q}_{m,d} \mathbf{K}_{n,d}) \mathbf{V}_{n,c} & \mathbf{Z}_{n,d,c} &= \mathbf{K}_{n,d} \mathbf{V}_{n,c} \\
 && \mathbf{H}_{m,d,c} &= \sum_{n=1}^N \mathbf{L}_{m,n} \mathbf{Z}_{n,d,c} = \sum_{n=1}^N \mathbf{L}_{m,n} \mathbf{K}_{n,d} \mathbf{V}_{n,c} \\
 \mathbf{Y}_{m,c} &= \sum_{d=1}^D \mathbf{Q}_{m,d} \mathbf{H}_{m,d,c} = \sum_{d=1}^D \mathbf{Q}_{m,d} \left(\sum_{n=1}^N \mathbf{L}_{m,n} \mathbf{K}_{n,d} \mathbf{V}_{n,c} \right) \\
 &= \sum_{d=1}^D \sum_{n=1}^N \mathbf{Q}_{m,d} \mathbf{L}_{m,n} \mathbf{K}_{n,d} \mathbf{V}_{n,c} \\
 &= \sum_{n=1}^N \left(\sum_{d=1}^D \mathbf{Q}_{m,d} \mathbf{L}_{m,n} \mathbf{K}_{n,d} \right) \mathbf{V}_{n,c}
 \end{aligned}$$

Corollary 1. The complexity of " $(\mathbf{Q}\mathbf{K}^\top \odot \mathbf{L})\mathbf{V}$ " with respect to N depends on "einsum[MN, NDC → MDC](\mathbf{L}, \mathbf{Z})"

Theorem 1. For any matrices $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{N \times D}$, $\mathbf{L} \in \mathbb{R}^{N \times N}$, $\mathbf{V} \in \mathbb{R}^{N \times C}$, the following equation holds:

$$\mathbf{Z} = \mathbf{K} \star \mathbf{V} = \text{einsum}[ND, NC \rightarrow NDC](\mathbf{K}, \mathbf{V}) \quad \mathcal{O}(NDC)$$

$$\mathbf{Y} = (\mathbf{Q}\mathbf{K}^\top \odot \mathbf{L})\mathbf{V} \iff \mathbf{H} = \text{einsum}[MN, NDC \rightarrow MDC](\mathbf{L}, \mathbf{Z}), \quad (M = N) \quad \mathcal{O}(N^2DC)$$

$$\mathbf{Y} = \mathbf{Q} \star \mathbf{H} = \text{einsum}[MD, MDC \rightarrow MC](\mathbf{Q}, \mathbf{H}) \quad \mathcal{O}(NDC)$$

Corollary 1. The complexity of " $(\mathbf{Q}\mathbf{K}^\top \odot \mathbf{L})\mathbf{V}$ " with respect to N depends on " $\text{einsum}[MN, NDC \rightarrow MDC](\mathbf{L}, \mathbf{Z})$ "

- Structured Masked Attention形式

$$\mathbf{Z} = \mathbf{B} \star \mathbf{X} = \text{einsum}[ND, NC \rightarrow NDC](\mathbf{B}, \mathbf{X})$$

$$\mathbf{Y} = (\mathbf{C}\mathbf{B}^\top \odot \mathbf{L})\mathbf{X} \iff \mathbf{H} = \text{einsum}[MN, NDC \rightarrow MDC](\mathbf{L}, \mathbf{Z}), \quad (M = N)$$

$$\mathbf{Y} = \mathbf{C} \star \mathbf{H} = \text{einsum}[MD, MDC \rightarrow MC](\mathbf{C}, \mathbf{H})$$

$$\mathbf{H} = \text{reshape}[\mathbb{R}^{N \times D \times C} \rightarrow \mathbb{R}^{N \times DC}](\mathbf{H}), \quad \mathbf{Z} = \text{reshape}[\mathbb{R}^{N \times D \times C} \rightarrow \mathbb{R}^{N \times DC}](\mathbf{Z})$$

$$\mathbf{H} = \text{einsum}[MN, NDC \rightarrow MDC](\mathbf{L}, \mathbf{Z}) \Rightarrow \mathbf{H} = \mathbf{L}\mathbf{Z} \iff \begin{bmatrix} \mathbf{H}_0 \\ \mathbf{H}_1 \\ \vdots \\ \mathbf{H}_{N-2} \\ \mathbf{H}_{N-1} \end{bmatrix} = \begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix} \begin{bmatrix} \mathbf{Z}_0 \\ \mathbf{Z}_1 \\ \vdots \\ \mathbf{Z}_{N-2} \\ \mathbf{Z}_{N-1} \end{bmatrix}$$

Theorem 2. For any 1-semiseparable structured matrix $\mathbf{L} \in \mathbb{R}^{N \times N}$ and input tokens $\mathbf{X} \in \mathbb{R}^{N \times C}$, the calculation of $\mathbf{Y} = \mathbf{L}\mathbf{X}$ can be achieved with a complexity of $\mathcal{O}(N)$ by applying the **chunkwise algorithm**.

$$\underbrace{\begin{bmatrix} \mathbf{Y}_0 \\ \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_{N-2} \\ \mathbf{Y}_{N-1} \end{bmatrix}}_{\mathbb{R}^{N \times C}} = \underbrace{\begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix}}_{\mathbb{R}^{N \times N}} \underbrace{\begin{bmatrix} \mathbf{X}_0 \\ \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_{N-2} \\ \mathbf{X}_{N-1} \end{bmatrix}}_{\mathbb{R}^{N \times C}}$$

- Structured Masked Attention形式 Complexity: $\mathcal{O}(N^2) \rightarrow \mathcal{O}(N)$

$$\underbrace{\begin{bmatrix} \mathbf{H}_0 \\ \mathbf{H}_1 \\ \vdots \\ \mathbf{H}_{N-2} \\ \mathbf{H}_{N-1} \end{bmatrix}}_{I\text{-semiseparable structured matrix}} = \underbrace{\begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix}}_{I\text{-semiseparable structured matrix}} \underbrace{\begin{bmatrix} \mathbf{Z}_0 \\ \mathbf{Z}_1 \\ \vdots \\ \mathbf{Z}_{N-2} \\ \mathbf{Z}_{N-1} \end{bmatrix}}_{\mathbb{R}^{N \times C}}$$

Chunkwise algorithm For example, let sequence length $N = 3136$, chunk size $K = 56$

$$\begin{bmatrix} Y_0 \\ Y_1 \\ \vdots \\ Y_{N-2} \\ Y_{N-1} \end{bmatrix} = \underbrace{\begin{bmatrix} A_{0:0} & 0 & \dots & 0 & 0 \\ A_{1:0} & A_{1:1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N-2:0} & A_{N-2:1} & \dots & A_{N-2:N-2} & 0 \\ A_{N-1:0} & A_{N-1:1} & \dots & A_{N-1:N-2} & A_{N-1:N-1} \end{bmatrix}}_{I\text{-semiseparable structured matrix}} \begin{bmatrix} X_0 \\ X_1 \\ \vdots \\ X_{N-2} \\ X_{N-1} \end{bmatrix}$$

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{54} \\ y_{55} \end{bmatrix} = \begin{bmatrix} 1 \\ A_{1:0} & 1 \\ \vdots & \vdots & \ddots \\ A_{54:0} & A_{54:1} & \cdots & 1 \\ A_{55:0} & A_{55:1} & \cdots & A_{55:54} & 1 \\ A_{56:0} & A_{56:1} & \cdots & A_{56:54} & A_{56:55} \\ A_{57:0} & A_{57:1} & \cdots & A_{57:54} & A_{57:55} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{110:0} & A_{110:1} & \cdots & A_{110:54} & A_{110:55} \\ A_{111:0} & A_{111:1} & \cdots & A_{111:54} & A_{111:55} \end{bmatrix} \begin{bmatrix} 1 \\ A_{57:56} & 1 \\ \vdots & \vdots & \ddots \\ A_{110:56} & A_{110:57} & \cdots & 1 \\ A_{111:56} & A_{111:57} & \cdots & A_{111:110} & 1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{54} \\ x_{55} \end{bmatrix}$$

$$= \begin{bmatrix} y_{3024} \\ y_{3025} \\ \vdots \\ y_{3078} \\ y_{3079} \\ y_{3080} \\ y_{3081} \\ \vdots \\ y_{3134} \\ y_{3135} \end{bmatrix} = \begin{bmatrix} A_{3024:0} & A_{3024:1} & \cdots & A_{3024:54} & A_{3024:55} \\ A_{3025:0} & A_{3025:1} & \cdots & A_{3025:54} & A_{3025:55} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{3078:0} & A_{3078:1} & \cdots & A_{3078:54} & A_{3078:55} \\ A_{3079:0} & A_{3079:1} & \cdots & A_{3079:54} & A_{3079:55} \\ A_{3080:0} & A_{3080:1} & \cdots & A_{3080:54} & A_{3080:55} \\ A_{3081:0} & A_{3081:1} & \cdots & A_{3081:54} & A_{3081:55} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{3134:0} & A_{3134:1} & \cdots & A_{3134:54} & A_{3134:55} \\ A_{3135:0} & A_{3135:1} & \cdots & A_{3135:54} & A_{3135:55} \end{bmatrix} \begin{bmatrix} A_{3024:56} & A_{3024:57} & \cdots & A_{3024:110} & A_{3024:111} \\ A_{3025:56} & A_{3025:57} & \cdots & A_{3025:110} & A_{3025:111} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{3078:56} & A_{3078:57} & \cdots & A_{3078:110} & A_{3078:111} \\ A_{3079:56} & A_{3079:57} & \cdots & A_{3079:110} & A_{3079:111} \\ A_{3080:56} & A_{3080:57} & \cdots & A_{3080:110} & A_{3080:111} \\ A_{3081:56} & A_{3081:57} & \cdots & A_{3081:110} & A_{3081:111} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{3134:56} & A_{3134:57} & \cdots & A_{3134:110} & A_{3134:111} \\ A_{3135:56} & A_{3135:57} & \cdots & A_{3135:110} & A_{3135:111} \end{bmatrix} \begin{bmatrix} 1 \\ A_{3025:3024} & 1 \\ \vdots & \vdots & \ddots \\ A_{3078:3024} & A_{3078:3025} & \cdots & 1 \\ A_{3079:3024} & A_{3079:3025} & \cdots & A_{3079:3078} & 1 \\ A_{3080:3024} & A_{3080:3025} & \cdots & A_{3080:3078} & A_{3080:3079} \\ A_{3081:3024} & A_{3081:3025} & \cdots & A_{3081:3078} & A_{3081:3079} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{3134:3024} & A_{3134:3025} & \cdots & A_{3134:3078} & A_{3134:3079} \\ A_{3135:3024} & A_{3135:3025} & \cdots & A_{3135:3078} & A_{3135:3079} \end{bmatrix} \begin{bmatrix} x_{3024} \\ x_{3025} \\ \vdots \\ x_{3078} \\ x_{3079} \\ x_{3080} \\ x_{3081} \\ \vdots \\ x_{3134} \\ x_{3135} \end{bmatrix}$$

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{54} \\ y_{55} \end{bmatrix} \left[\begin{array}{c} 1 \\ A_{1:0} \quad 1 \\ \vdots \quad \vdots \quad \ddots \\ A_{54:0} \quad A_{54:1} \quad \cdots \quad 1 \\ A_{55:0} \quad A_{55:1} \quad \cdots \quad A_{55:54} \quad 1 \end{array} \right] \begin{bmatrix} y_{56} \\ y_{57} \\ \vdots \\ y_{110} \\ y_{111} \\ \vdots \end{bmatrix} = \left[\begin{array}{c} 1 \\ A_{56:0} \quad A_{56:1} \quad \cdots \quad A_{56:54} \quad A_{56:55} \\ A_{57:0} \quad A_{57:1} \quad \cdots \quad A_{57:54} \quad A_{57:55} \\ \vdots \quad \vdots \quad \ddots \quad \vdots \quad \vdots \\ A_{110:0} \quad A_{110:1} \quad \cdots \quad A_{110:54} \quad A_{110:55} \\ A_{111:0} \quad A_{111:1} \quad \cdots \quad A_{111:54} \quad A_{111:55} \end{array} \right] \begin{bmatrix} 1 \\ A_{57:56} \quad 1 \\ \vdots \quad \vdots \quad \ddots \\ A_{110:56} \quad A_{110:57} \quad \cdots \quad 1 \\ A_{111:56} \quad A_{111:57} \quad \cdots \quad A_{111:110} \quad 1 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{54} \\ x_{55} \\ \vdots \\ x_{56} \\ x_{57} \\ \vdots \\ x_{110} \\ x_{111} \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} y_{3024} \\ y_{3025} \\ \vdots \\ y_{3078} \\ y_{3079} \\ y_{3080} \\ y_{3081} \\ \vdots \\ y_{3134} \\ y_{3135} \end{bmatrix} = \left[\begin{array}{c} 1 \\ A_{3024:0} \quad A_{3024:1} \quad \cdots \quad A_{3024:54} \quad A_{3024:55} \\ A_{3025:0} \quad A_{3025:1} \quad \cdots \quad A_{3025:54} \quad A_{3025:55} \\ \vdots \quad \vdots \quad \ddots \quad \vdots \quad \vdots \\ A_{3078:0} \quad A_{3078:1} \quad \cdots \quad A_{3078:54} \quad A_{3078:55} \\ A_{3079:0} \quad A_{3079:1} \quad \cdots \quad A_{3079:54} \quad A_{3079:55} \\ A_{3080:0} \quad A_{3080:1} \quad \cdots \quad A_{3080:54} \quad A_{3080:55} \\ A_{3081:0} \quad A_{3081:1} \quad \cdots \quad A_{3081:54} \quad A_{3081:55} \\ \vdots \quad \vdots \quad \ddots \quad \vdots \quad \vdots \\ A_{3134:0} \quad A_{3134:1} \quad \cdots \quad A_{3134:54} \quad A_{3134:55} \\ A_{3135:0} \quad A_{3135:1} \quad \cdots \quad A_{3135:54} \quad A_{3135:55} \end{array} \right] \begin{bmatrix} 1 \\ A_{3024:56} \quad A_{3024:57} \quad \cdots \quad A_{3024:110} \quad A_{3024:111} \\ A_{3025:56} \quad A_{3025:57} \quad \cdots \quad A_{3025:110} \quad A_{3025:111} \\ \vdots \quad \vdots \quad \ddots \quad \vdots \quad \vdots \\ A_{3078:56} \quad A_{3078:57} \quad \cdots \quad A_{3078:110} \quad A_{3078:111} \\ A_{3079:56} \quad A_{3079:57} \quad \cdots \quad A_{3079:110} \quad A_{3079:111} \\ A_{3080:56} \quad A_{3080:57} \quad \cdots \quad A_{3080:110} \quad A_{3080:111} \\ A_{3081:56} \quad A_{3081:57} \quad \cdots \quad A_{3081:110} \quad A_{3081:111} \\ \vdots \quad \vdots \quad \ddots \quad \vdots \quad \vdots \\ A_{3134:56} \quad A_{3134:57} \quad \cdots \quad A_{3134:110} \quad A_{3134:111} \\ A_{3135:56} \quad A_{3135:57} \quad \cdots \quad A_{3135:110} \quad A_{3135:111} \end{bmatrix} \cdots \left[\begin{array}{c} 1 \\ A_{3025:3024} \quad 1 \\ \vdots \quad \vdots \quad \ddots \\ A_{3078:3024} \quad A_{3078:3025} \quad \cdots \quad 1 \\ A_{3079:3024} \quad A_{3079:3025} \quad \cdots \quad A_{3079:3078} \quad 1 \\ A_{3080:3024} \quad A_{3080:3025} \quad \cdots \quad A_{3080:3078} \quad A_{3080:3079} \\ A_{3081:3024} \quad A_{3081:3025} \quad \cdots \quad A_{3081:3078} \quad A_{3081:3079} \\ \vdots \quad \vdots \quad \ddots \quad \vdots \quad \vdots \\ A_{3134:3024} \quad A_{3134:3025} \quad \cdots \quad A_{3134:3078} \quad A_{3134:3079} \\ A_{3135:3024} \quad A_{3135:3025} \quad \cdots \quad A_{3135:3078} \quad A_{3135:3079} \end{array} \right] \begin{bmatrix} x_{3024} \\ x_{3025} \\ \vdots \\ x_{3078} \\ x_{3079} \\ x_{3080} \\ x_{3081} \\ \vdots \\ x_{3134} \\ x_{3135} \end{bmatrix}$$

$$\begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_{54} \\ y_{55} \end{bmatrix} \left[\begin{array}{c} 1 \\ A_{1:0} \quad 1 \\ \vdots \quad \vdots \quad \ddots \\ A_{54:0} \quad A_{54:1} \quad \cdots \quad 1 \\ A_{55:0} \quad A_{55:1} \quad \cdots \quad A_{55:54} \quad 1 \end{array} \right] \begin{bmatrix} y_{56} \\ y_{57} \\ \vdots \\ y_{110} \\ y_{111} \\ \vdots \end{bmatrix} = \left[\begin{array}{c} 1 \\ A_{57:56} \\ \vdots \\ A_{110:56} \\ A_{111:56} \end{array} \right] \left[\begin{array}{c} 1 \\ A_{57:56} \quad 1 \\ \vdots \quad \vdots \quad \ddots \\ A_{110:56} \quad A_{110:57} \quad \cdots \quad 1 \\ A_{111:56} \quad A_{111:57} \quad \cdots \quad A_{111:110} \quad 1 \end{array} \right] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{54} \\ x_{55} \\ \vdots \\ x_{56} \\ x_{57} \\ \vdots \\ x_{110} \\ x_{111} \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} y_{3024} \\ y_{3025} \\ \vdots \\ y_{3078} \\ y_{3079} \\ y_{3080} \\ y_{3081} \\ \vdots \\ y_{3134} \\ y_{3135} \end{bmatrix} = \left[\begin{array}{c} 1 \\ A_{3025:3024} \\ \vdots \\ A_{3078:3024} \\ A_{3079:3024} \end{array} \right] A_{3024:56}[A_{56:0} \quad A_{56:1} \quad \cdots \quad A_{56:54} \quad A_{56:55}] \left[\begin{array}{c} 1 \\ A_{3025:3024} \\ \vdots \\ A_{3078:3024} \\ A_{3079:3024} \end{array} \right] \cdots \left[\begin{array}{c} 1 \\ A_{3025:3024} \quad 1 \\ \vdots \quad \vdots \quad \ddots \\ A_{3078:3024} \quad A_{3078:3025} \quad \cdots \quad 1 \\ A_{3079:3024} \quad A_{3079:3025} \quad \cdots \quad A_{3079:3078} \quad 1 \end{array} \right] \begin{bmatrix} x_{3024} \\ x_{3025} \\ \vdots \\ x_{3078} \\ x_{3079} \\ x_{3080} \\ x_{3081} \\ \vdots \\ x_{3134} \\ x_{3135} \end{bmatrix}$$

Y_{diag}

$$Y^{diag} = \begin{bmatrix} 1 & & & & \\ A_{1:0} & 1 & & & \\ \vdots & \vdots & \ddots & & \\ A_{54:0} & A_{54:1} & \cdots & 1 & \\ A_{55:0} & A_{55:1} & \cdots & A_{55:54} & 1 \end{bmatrix} \begin{bmatrix} 1 & & & & \\ A_{57:56} & 1 & & & \\ \vdots & \vdots & \ddots & & \\ A_{110:56} & A_{110:57} & \cdots & 1 & \\ A_{111:56} & A_{111:57} & \cdots & A_{111:110} & 1 \end{bmatrix} \begin{bmatrix} \ddots & & & & \\ & 1 & & & \\ & A_{3025:3024} & 1 & & \\ & \vdots & \vdots & \ddots & \\ & A_{3078:3024} & A_{3078:3025} & \cdots & 1 \\ & A_{3079:3024} & A_{3079:3025} & \cdots & A_{3079:3078} & 1 \end{bmatrix} \begin{bmatrix} x_0 & & & & \\ x_1 & & & & \\ \vdots & & & & \\ x_{54} & & & & \\ x_{55} & & & & \\ x_{56} & & & & \\ x_{57} & & & & \\ \vdots & & & & \\ x_{110} & & & & \\ x_{111} & & & & \\ \vdots & & & & \\ x_{3024} & & & & \\ x_{3025} & & & & \\ \vdots & & & & \\ x_{3078} & & & & \\ x_{3079} & & & & \\ x_{3080} & & & & \\ x_{3081} & & & & \\ \vdots & & & & \\ x_{3134} & & & & \\ x_{3135} & & & & \end{bmatrix}$$

$$Y^{tril} = \begin{bmatrix} 0 & & & & \\ 0 & 0 & & & \\ \vdots & \vdots & \ddots & & \\ 0 & 0 & \cdots & 0 & \\ 0 & 0 & \cdots & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ A_{57:56} \\ \vdots \\ A_{110:56} \\ A_{111:56} \end{bmatrix} A_{56:56} [A_{56:0} \quad A_{56:1} \quad \cdots \quad A_{56:54} \quad A_{56:55}] \begin{bmatrix} 0 & & & & \\ 0 & 0 & & & \\ \vdots & \vdots & \ddots & & \\ 0 & 0 & \cdots & 0 & \\ 0 & 0 & \cdots & 0 & 0 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{54} \\ x_{55} \\ x_{56} \\ x_{57} \\ \vdots \\ x_{110} \\ x_{111} \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ A_{3025:3024} \\ \vdots \\ A_{3078:3024} \\ A_{3079:3024} \end{bmatrix} A_{3024:56}[A_{56:0} \ A_{56:1} \ \cdots \ A_{56:54} \ A_{56:55}] \begin{bmatrix} 1 \\ A_{3025:3024} \\ \vdots \\ A_{3078:3024} \\ A_{3079:3024} \end{bmatrix} A_{3024:112}[A_{112:56} \ A_{112:57} \ \cdots \ A_{112:110} \ A_{112:111}] \ \cdots \begin{bmatrix} 0 \\ 0 \ 0 \\ \vdots \ \vdots \ \ddots \\ 0 \ 0 \ \cdots \ 0 \\ 0 \ 0 \ \cdots \ 0 \ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{bmatrix} A_{3080:56}[A_{56:0} \ A_{56:1} \ \cdots \ A_{56:54} \ A_{56:55}] \begin{bmatrix} 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{bmatrix} A_{3080:112}[A_{112:56} \ A_{112:57} \ \cdots \ A_{112:110} \ A_{112:111}] \ \cdots \begin{bmatrix} 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{bmatrix} A_{3080:3080}[A_{3080:3024} \ A_{3080:3025} \ \cdots \ A_{3080:3078} \ A_{3080:3079}] \begin{bmatrix} 0 \\ 0 \ 0 \\ \vdots \ \vdots \ \ddots \\ 0 \ 0 \ \cdots \ 0 \\ 0 \ 0 \ \cdots \ 0 \ 0 \end{bmatrix}$$

$$\begin{aligned}
& \left[\begin{array}{c} 0 \\ 0 \quad 0 \\ \vdots \quad \vdots \quad \ddots \\ 0 \quad 0 \quad \cdots \quad 0 \\ 0 \quad 0 \quad \cdots \quad 0 \quad 0 \end{array} \right] \\
& \left[\begin{array}{c} 1 \\ A_{57:56} \\ \vdots \\ A_{110:56} \\ A_{111:56} \end{array} \right] A_{56:56}[A_{56:0} \quad A_{56:1} \quad \cdots \quad A_{56:54} \quad A_{56:55}] \quad \left[\begin{array}{c} 0 \\ 0 \quad 0 \\ \vdots \quad \vdots \quad \ddots \\ 0 \quad 0 \quad \cdots \quad 0 \\ 0 \quad 0 \quad \cdots \quad 0 \quad 0 \end{array} \right] \\
Y^{tril} = & \vdots \quad \vdots \quad \ddots \\
& \left[\begin{array}{c} 1 \\ A_{3025:3024} \\ \vdots \\ A_{3078:3024} \\ A_{3079:3024} \end{array} \right] A_{3024:56}[A_{56:0} \quad A_{56:1} \quad \cdots \quad A_{56:54} \quad A_{56:55}] \quad \left[\begin{array}{c} 1 \\ A_{3025:3024} \\ \vdots \\ A_{3078:3024} \\ A_{3079:3024} \end{array} \right] A_{3024:112}[A_{112:56} \quad A_{112:57} \quad \cdots \quad A_{112:110} \quad A_{112:111}] \quad \cdots \quad \left[\begin{array}{c} 0 \\ 0 \quad 0 \\ \vdots \quad \vdots \quad \ddots \\ 0 \quad 0 \quad \cdots \quad 0 \\ 0 \quad 0 \quad \cdots \quad 0 \quad 0 \end{array} \right] \\
& \left[\begin{array}{c} 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{array} \right] A_{3080:56}[A_{56:0} \quad A_{56:1} \quad \cdots \quad A_{56:54} \quad A_{56:55}] \quad \left[\begin{array}{c} 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{array} \right] A_{3080:112}[A_{112:56} \quad A_{112:57} \quad \cdots \quad A_{112:110} \quad A_{112:111}] \quad \cdots \quad \left[\begin{array}{c} 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{array} \right] A_{3080:3080}[A_{3080:3024} \quad A_{3080:3025} \quad \cdots \quad A_{3080:3078} \quad A_{3080:3079}] \quad \left[\begin{array}{c} 0 \\ 0 \quad 0 \\ \vdots \quad \vdots \quad \ddots \\ 0 \quad 0 \quad \cdots \quad 0 \\ 0 \quad 0 \quad \cdots \quad 0 \quad 0 \end{array} \right] \\
& \left[\begin{array}{c} x_0 \\ x_1 \\ \vdots \\ x_{54} \\ x_{55} \end{array} \right] \quad \left[\begin{array}{c} x_{56} \\ x_{57} \\ \vdots \\ x_{110} \\ x_{111} \end{array} \right] \quad \left[\begin{array}{c} x_{3024} \\ x_{3025} \\ \vdots \\ x_{3078} \\ x_{3079} \end{array} \right] \quad \left[\begin{array}{c} x_{3080} \\ x_{3081} \\ \vdots \\ x_{3134} \\ x_{3135} \end{array} \right]
\end{aligned}$$

$$\begin{aligned}
& \left[\begin{array}{c} [A_{56:0} \quad A_{56:1} \quad \cdots \quad A_{56:54} \quad A_{56:55}] \quad \left[\begin{array}{c} x_0 \\ x_1 \\ \vdots \\ x_{54} \\ x_{55} \end{array} \right] \\ [A_{112:56} \quad A_{112:57} \quad \cdots \quad A_{112:110} \quad A_{112:111}] \quad \left[\begin{array}{c} x_{56} \\ x_{57} \\ \vdots \\ x_{110} \\ x_{111} \end{array} \right] \\ \vdots \\ [A_{3080:3024} \quad A_{3080:3025} \quad \cdots \quad A_{3080:3078} \quad A_{3080:3079}] \quad \left[\begin{array}{c} x_{3024} \\ x_{3025} \\ \vdots \\ x_{3078} \\ x_{3079} \end{array} \right] \\ [0 \quad 0 \quad \cdots \quad 0 \quad 0] \quad \left[\begin{array}{c} x_{3080} \\ x_{3081} \\ \vdots \\ x_{3134} \\ x_{3135} \end{array} \right] \end{array} \right], \\
& Y^{tril} = \left[\begin{array}{c} 0 \\ 1 \\ A_{57:56} \\ \vdots \\ A_{110:56} \\ A_{111:56} \end{array} \right] A_{56:56} \quad \left[\begin{array}{c} 0 \\ 1 \\ A_{3025:3024} \\ \vdots \\ A_{3078:3024} \\ A_{3079:3024} \end{array} \right] A_{3024:112} \quad \cdots \quad \left[\begin{array}{c} 0 \\ 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{array} \right] A_{3080:112} \quad \cdots \quad \left[\begin{array}{c} 0 \\ 1 \\ A_{3081:3080} \\ \vdots \\ A_{3134:3080} \\ A_{3135:3080} \end{array} \right] A_{3080:3080} \\
& \left[\begin{array}{c} s_0 \\ s_1 \\ \vdots \\ s_{54} \\ s_{55} \end{array} \right] = \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right] \quad \left[\begin{array}{c} C_{56} \\ C_{57} \\ \vdots \\ C_{110} \\ C_{111} \end{array} \right] A_{57:56} \quad \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right] \quad \left[\begin{array}{c} C_{56} \\ C_{57} \\ \vdots \\ C_{110} \\ C_{111} \end{array} \right] A_{110:56} \quad \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right] \quad \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right] \\
& \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right] = \left[\begin{array}{c} s_0 \\ s_1 \\ \vdots \\ s_{54} \\ s_{55} \end{array} \right] \quad \left[\begin{array}{c} C_{2024} \\ C_{2025} \\ \vdots \\ C_{3078} \\ C_{3079} \end{array} \right] A_{3025:3024} \quad \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right] \quad \left[\begin{array}{c} C_{2024} \\ C_{2025} \\ \vdots \\ C_{3078} \\ C_{3079} \end{array} \right] A_{3080:3079} \quad \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right] \quad \left[\begin{array}{c} s_0 \\ s_1 \\ \vdots \\ s_{54} \\ s_{55} \end{array} \right] \\
& \text{Loop}
\end{aligned}$$

Complexity: $\mathcal{O}(\frac{N}{K} \times K)$

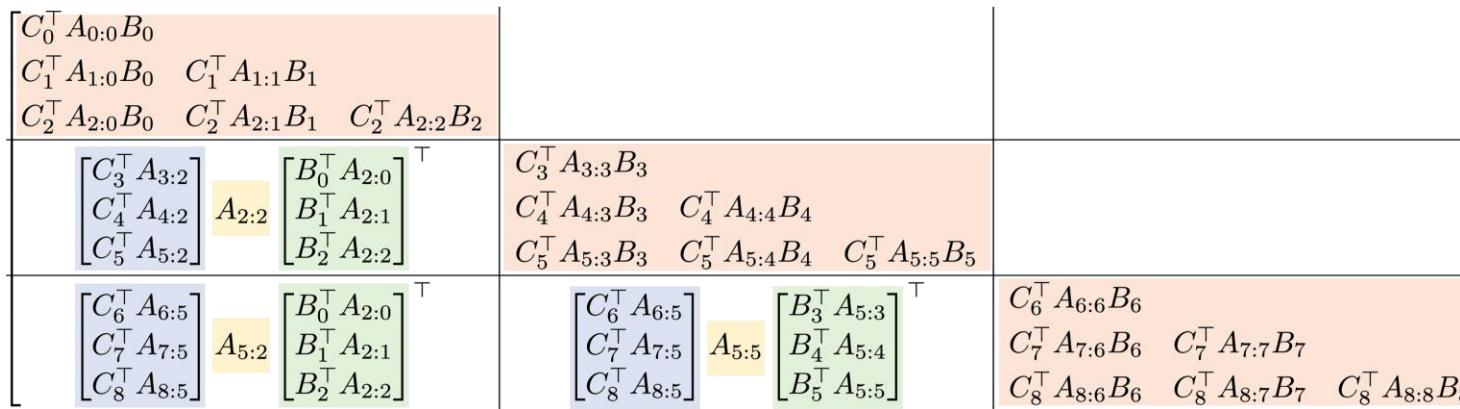
Complexity: $\mathcal{O}(\frac{N}{K} + N)$

4.3 Mamba2: The Algorithm

72

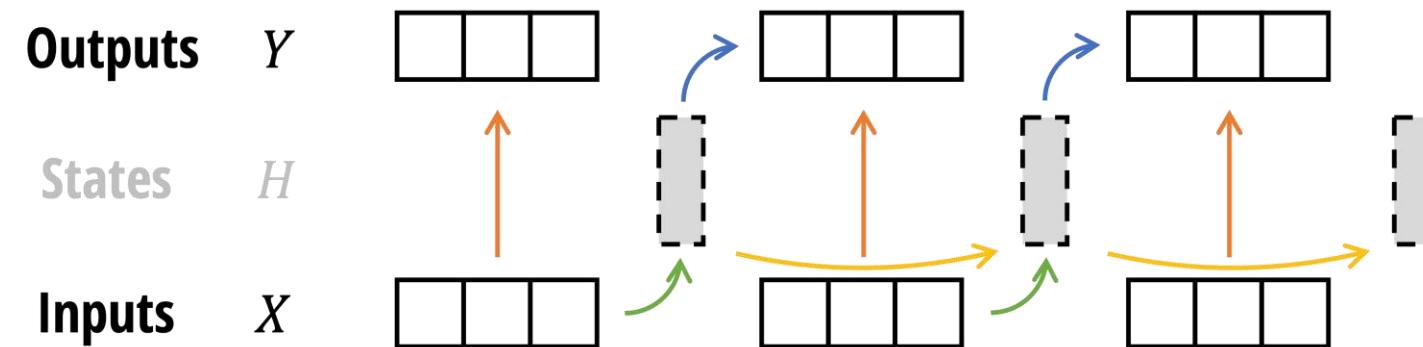
Mamba2: Structured Masked Attention

$$Y = (\mathbf{C}B^\top \odot \mathbf{L})\mathbf{X} = \mathbf{C} \star (\mathbf{L} \times (\mathbf{B} \star \mathbf{X})) \quad \text{Complexity: } \mathcal{O}(N^2) \rightarrow \mathcal{O}(N)$$



Semiseparable Matrix M

Block Decomposition



- Diagonal Block: Input → Output
- Low-Rank Block: Input → State
- Low-Rank Block: State → State
- Low-Rank Block: State → Output

4.4 Mamba2: The Systems

73

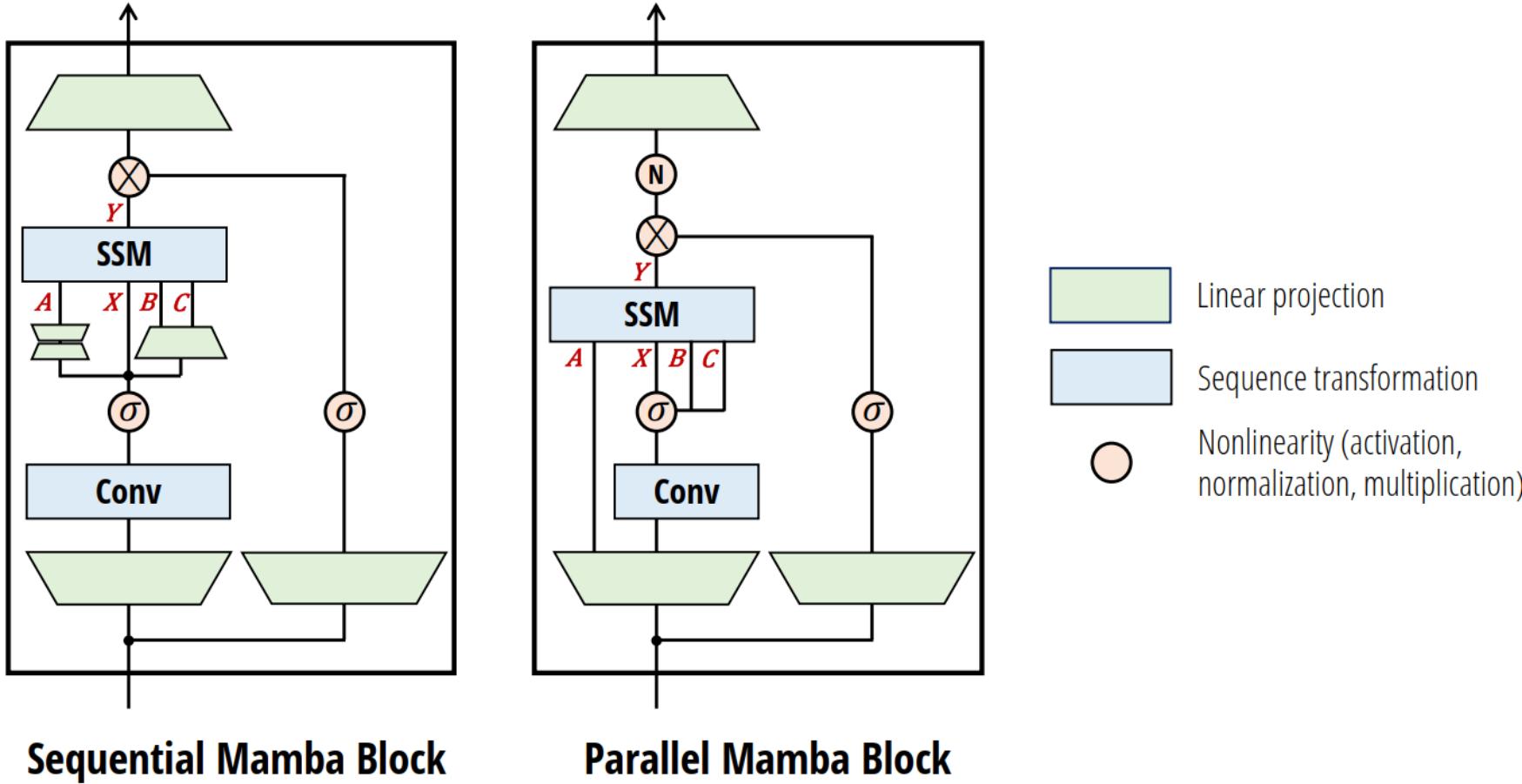


Figure 6: (**Mamba-2 Architecture**.) The Mamba-2 block simplifies the Mamba block by removing sequential linear projections; the SSM parameters A, B, C are produced at the beginning of the block instead of as a function of the SSM input X . An additional normalization layer is added as in NormFormer (Shleifer, Weston, and Ott 2021), improving stability. The B and C projections only have a single head shared across the X heads, analogous to multi-value attention (MVA).

4.4 Mamba2: The Systems

74

	Self-Attention	Mamba (SSM)	Mamba2 (SSD)
State size	N	1	1
Training FLOPs	N^2	N	N
Inference FLOPs	N	1	1
(Naive) memory implementation	N^2	N	N
Matrix multiplication	CUDA ✓	CUDA ✗	Triton ✓

	PyTorch	CUDA kernel	Triton kernel
性能表现	🟡 中等, 通用优化	🟢🟢 最高, 可精细优化	🟢 高, 结构友好自动并行优化
易用性	🟢 非常高, 调用即用	🔴 最低, 需熟悉 CUDA 编程	🟡 需要写 kernel
灵活性	🟡 中等	🟢🟢 最高, 完全掌控硬件细节	🟢 高, 支持结构感知优化
编译时间	🟢 无编译, 直接运行	🔴 静态编译, 耗时且复杂	🟡 动态 JIT 编译
适合场景	通用矩阵/张量操作	极限性能场景、重复运行模型	Transformer, SSM等结构优化

4.4 Mamba2: The Systems

- Experiment results
 - Mamba-2 (SSD) is significantly **faster** and scales much better than Mamba-1(Scan)

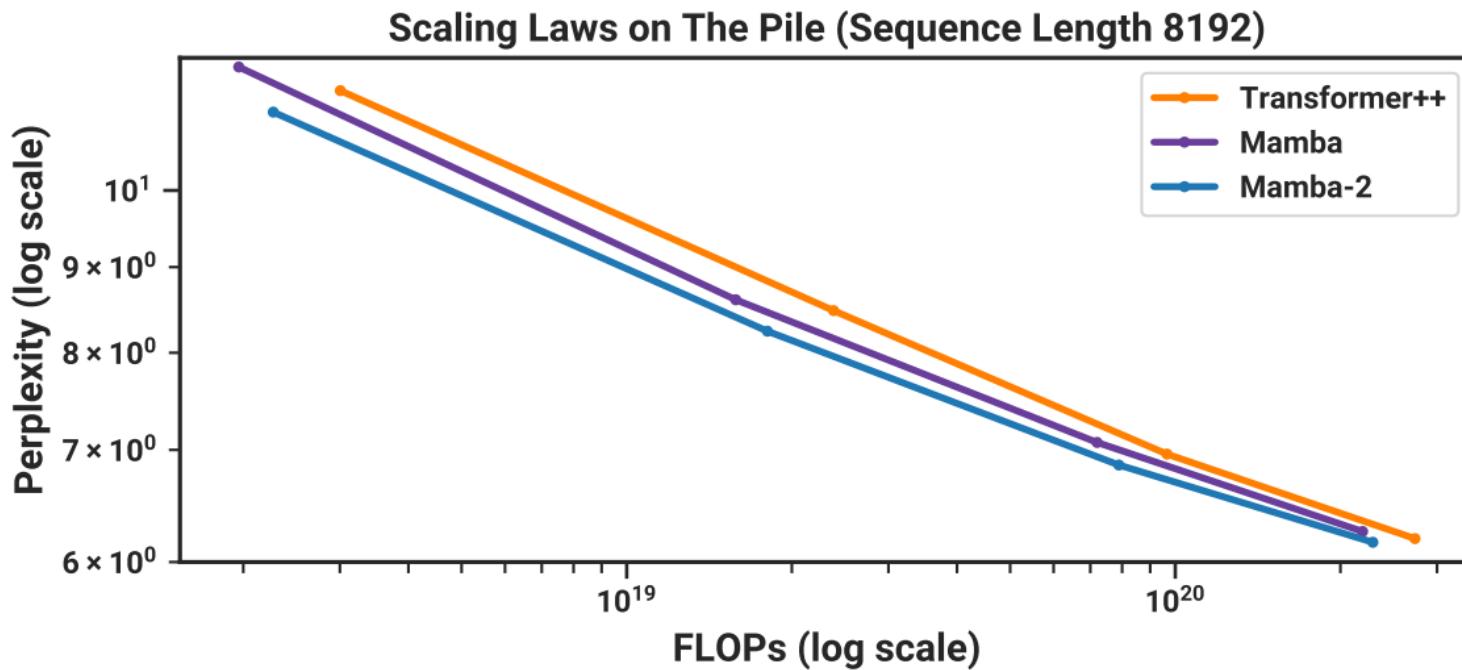


Figure 9: (**Scaling Laws.**) Models of size $\approx 125M$ to $\approx 1.3B$ parameters, trained on the Pile. Mamba-2 matches or exceeds the performance of Mamba as well as a strong “Transformer++” recipe. Compared to our Transformer baseline, Mamba-2 is Pareto dominant on performance (perplexity), theoretical FLOPs, and actual wall-clock time.

4.4 Mamba2: The Systems

76

- Experiment results

- Mamba-2 (SSD) is significantly **faster** and scales much better than Mamba-1(Scan)

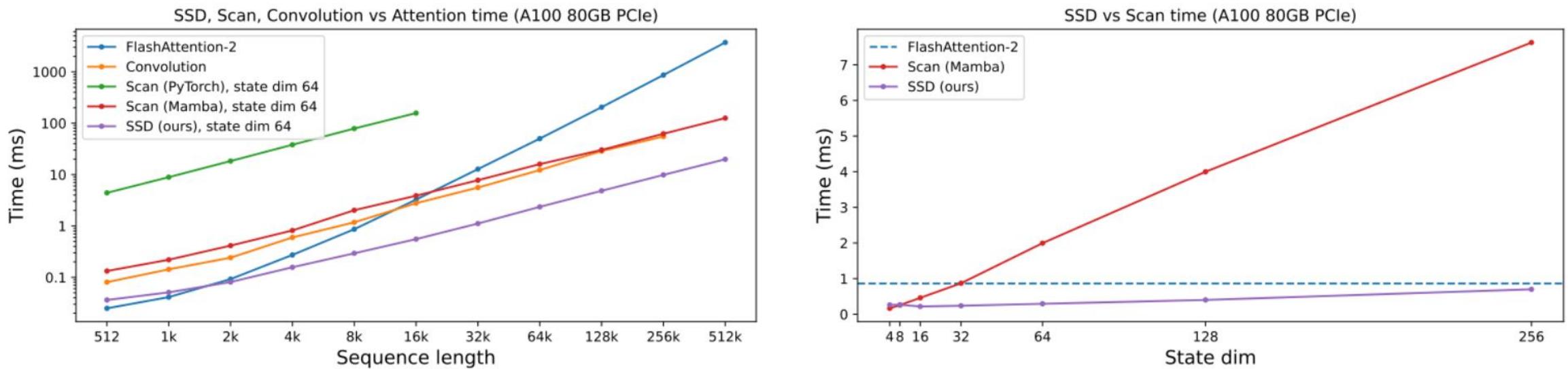


Figure 10: (**Efficiency Benchmarks.**) (Left) Our SSD is 2 – 8× faster than a Mamba fused scan for large state expansion ($N = 64$) and faster than FlashAttention-2 for sequence length 2k and above. (Right) Sequence length 4K: Increasing state expansion slows down the Mamba optimized scan implementation linearly. SSD can handle much larger state expansion factors without much slowdown.



01 State Space Model

02 Mamba1

03 Vision Mamba

04 Mamba2

05 PPMA

5.1 PPMA: Motivation (Current Issue)

Question: Have SOTA Mamba-based models outperform SOTA ViTs in computer vision domain, especially on high-level vision tasks?

Not yet, experimental results speak louder than words

Table 5: Image classification performance on the ImageNet-1K.

Model	Arch.	#Param. (M)	FLOPs (G)	Top-1 (%)
Spatial-Mamba-T [46]	SSM	29	4.5	82.1
		19	4.2	82.9
		30	4.9	82.6
		30	4.8	83.4
	Trans.	27	4.5	83.5
		25	4.2	83.5
		29	4.5	82.1
		28	4.3	83.2
BiFormer-S [56]	Trans.	26	4.5	83.8
RMT-S [10]		27	4.5	84.0
ConvNeXt-S [32]		50	8.7	83.1
EffNet-B5 [41]	CNN	30	9.9	83.6
VMamba-S [30]	SSM	50	8.7	83.6
GrootVL-S [48]		51	8.5	84.2
MLLA-S [15]		43	7.3	84.4
Spatial-Mamba-S [46]		43	7.1	84.6
Swin-S [31]		50	8.7	83.0
NAT-S [16]	Trans.	51	7.8	83.7
BiFormer-B [56]		57	9.8	84.3
iFormer-B [37]		48	9.4	84.6
RMT-B [10]		54	9.7	84.9

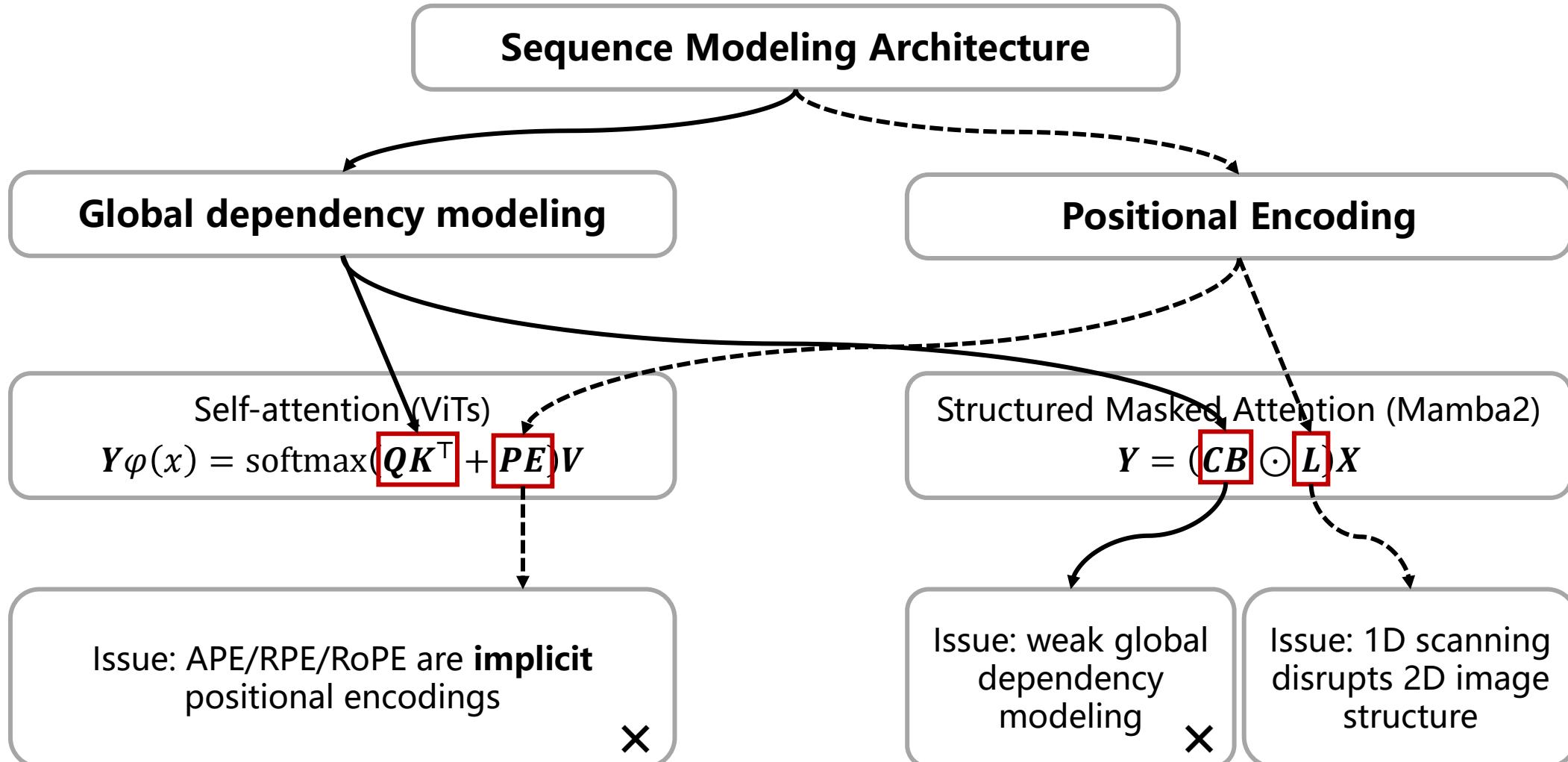
Table 6: Object detection and instance segmentation performance on COCO.

Backbone	Arch.	#Param. (M)	FLOPs (G)	AP ^b	AP ^m		
Spatial-Mamba-T [46]	SSM	ResNet-50 [19]	CNN	44	260	38.2	34.7
		ConvNeXt-T [32]	CNN	48	262	44.2	40.1
		MLLA-T [15]	CNN	44	255	46.8	42.1
		GrootVL-T [48]	CNN	49	265	47.0	42.7
	Trans.	VMamba-T [30]	CNN	50	271	47.3	42.7
		Spatial-Mamba-T [46]	CNN	46	216	47.6	42.9
		Swin-T [31]	Trans.	48	267	43.7	39.8
		CSWin-T [8]	Trans.	42	279	46.7	42.2
Spatial-Mamba-S [46]	SSM	BiFormer-S [56]	Trans.	—	—	47.8	43.2
		RMT-S [10]	Trans.	46	262	48.8	43.6
	Trans.	ResNet-101 [19]	CNN	63	336	40.4	36.4
		ConvNeXt-S [32]	CNN	70	348	45.4	41.8
		GrootVL-S [48]	CNN	70	341	48.6	43.6
		VMamba-S [30]	CNN	70	349	48.7	43.7
	Trans.	Spatial-Mamba-S [46]	CNN	63	315	49.2	44.0
		MLLA-S [15]	CNN	63	319	49.2	44.2
		Swin-S [31]	Trans.	69	359	45.7	41.1
		CSWin-S [8]	Trans.	54	342	47.9	43.2
Spatial-Mamba-B [46]	Trans.	BiFormer-B [56]	Trans.	—	—	48.6	43.7
		RMT-B [10]	Trans.	73	373	50.7	45.1

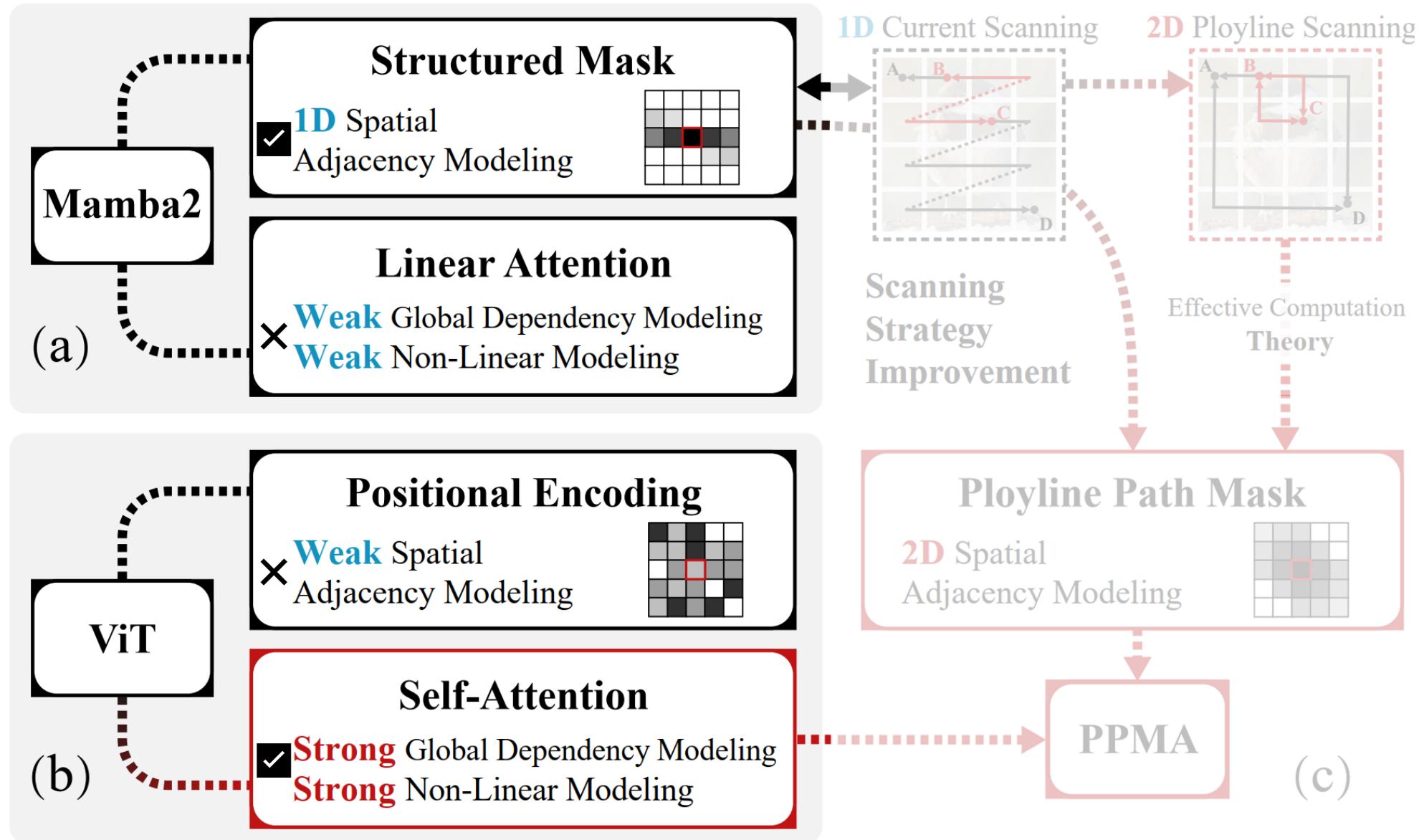
Table 7: Semantic segmentation performance on ADE20K.

Backbone	Arch.	#Param. (M)	FLOPs (G)	mIoU(%) SS MS			
Spatial-Mamba-S [46]	SSM	ResNet-50 [19]	CNN	67	953	42.1	42.8
		ConvNeXt-T [32]	CNN	60	939	46.0	46.7
		VMamba-T [30]	CNN	62	949	48.0	48.8
		2DMamba-T [52]	CNN	62	950	48.6	49.3
	Trans.	GrootVL-T [48]	CNN	60	941	48.5	49.4
		Spatial-Mamba-S [46]	CNN	57	936	48.6	49.4
		Swin-T [31]	Trans.	60	945	44.4	45.8
		NAT-T [16]	Trans.	58	934	47.1	48.4
Spatial-Mamba-B [46]	Trans.	BiFormer-S [56]	Trans.	—	—	49.8	50.8
		RMT-S [10]	Trans.	56	937	49.8	49.7
	Trans.	ResNet-101 [19]	CNN	85	1030	42.9	44.0
		ConvNeXt-S [32]	CNN	82	1027	48.7	49.6
		VMamba-S [30]	CNN	82	1028	50.6	51.2
	SSM	Spatial-Mamba-S [46]	CNN	73	992	50.6	51.4
		GrootVL-S [48]	CNN	82	1019	50.7	51.7
		Swin-S [31]	Trans.	81	1039	47.6	49.5
Spatial-Mamba-B [46]	Trans.	NAT-S [16]	Trans.	82	1010	48.0	49.5
		BiFormer-B [56]	Trans.	—	—	51.0	51.7
		RMT-B [10]	Trans.	83	1051	52.0	52.1

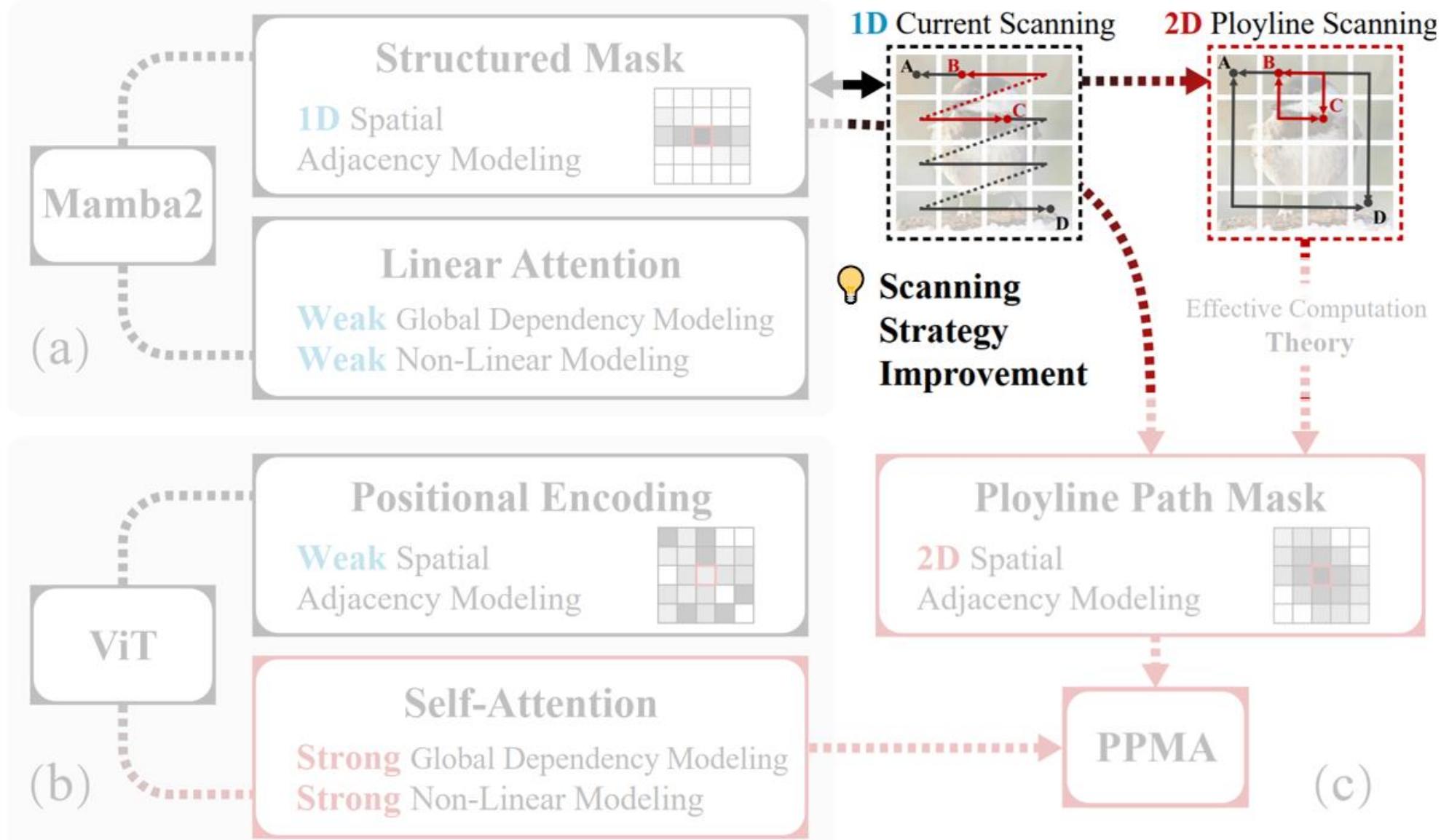
5.1 PPMA: Motivation (Current Issue)



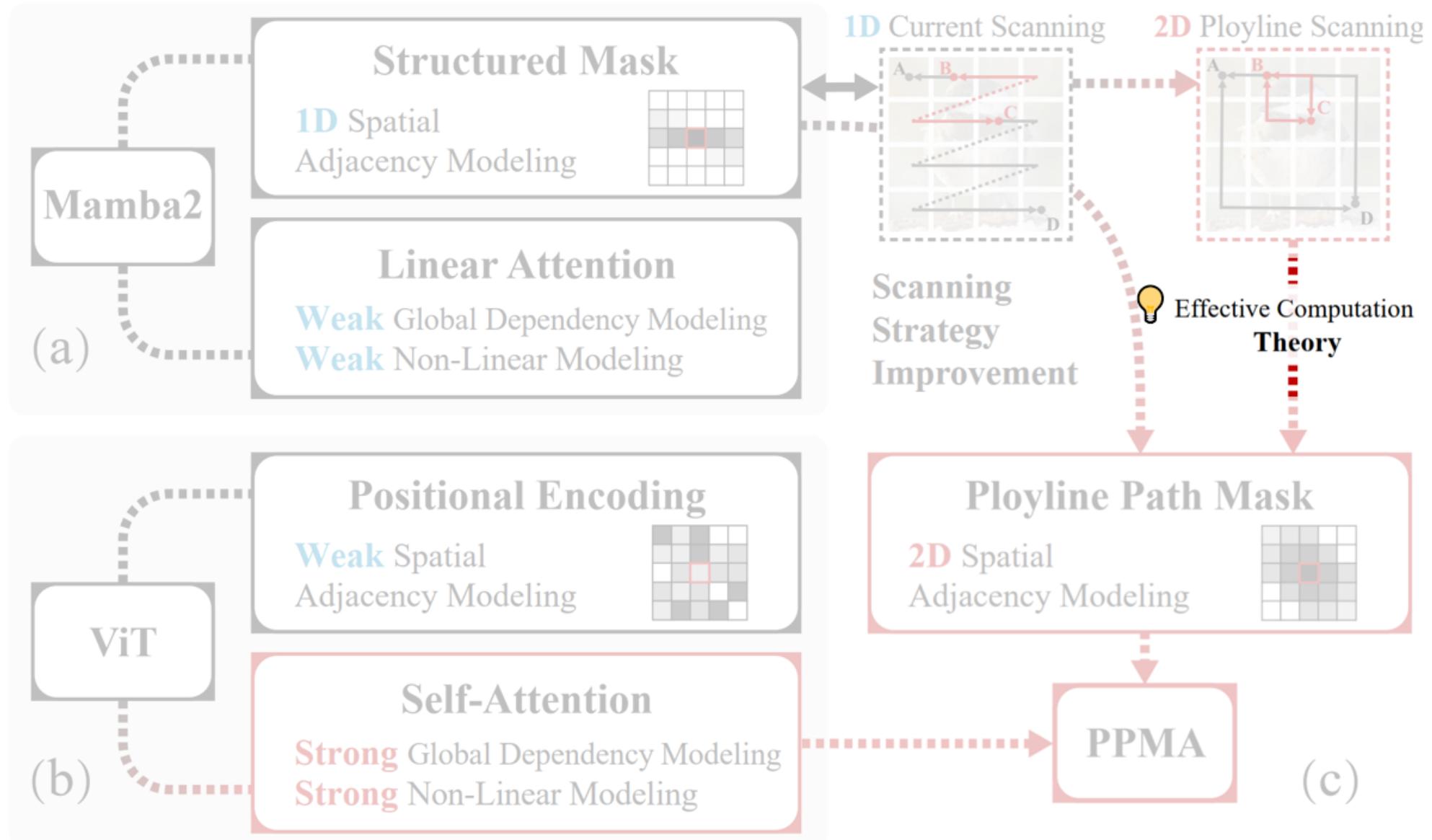
We adapt Mamba2's 1D structured mask to **2D polyline path mask** and integrate it into the self-attention mechanism of ViTs as an **explicit** positional encoding.



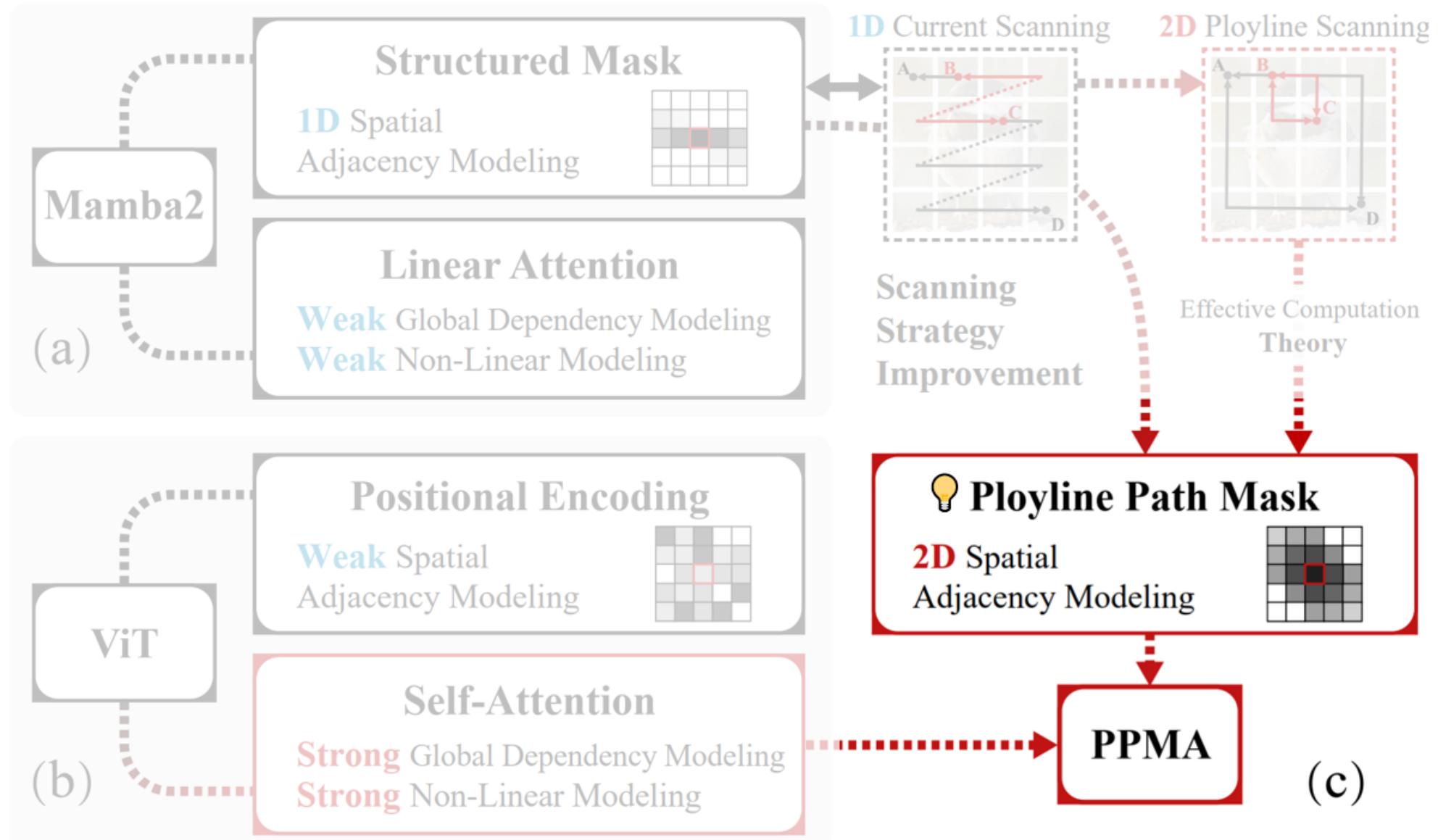
We adapt Mamba2's 1D structured mask to **2D polyline path mask** and integrate it into the self-attention mechanism of ViTs as an **explicit** positional encoding.



We adapt Mamba2's 1D structured mask to **2D polyline path mask** and integrate it into the self-attention mechanism of ViTs as an **explicit** positional encoding.



We adapt Mamba2's 1D structured mask to **2D polyline path mask** and integrate it into the self-attention mechanism of ViTs as an **explicit** positional encoding.

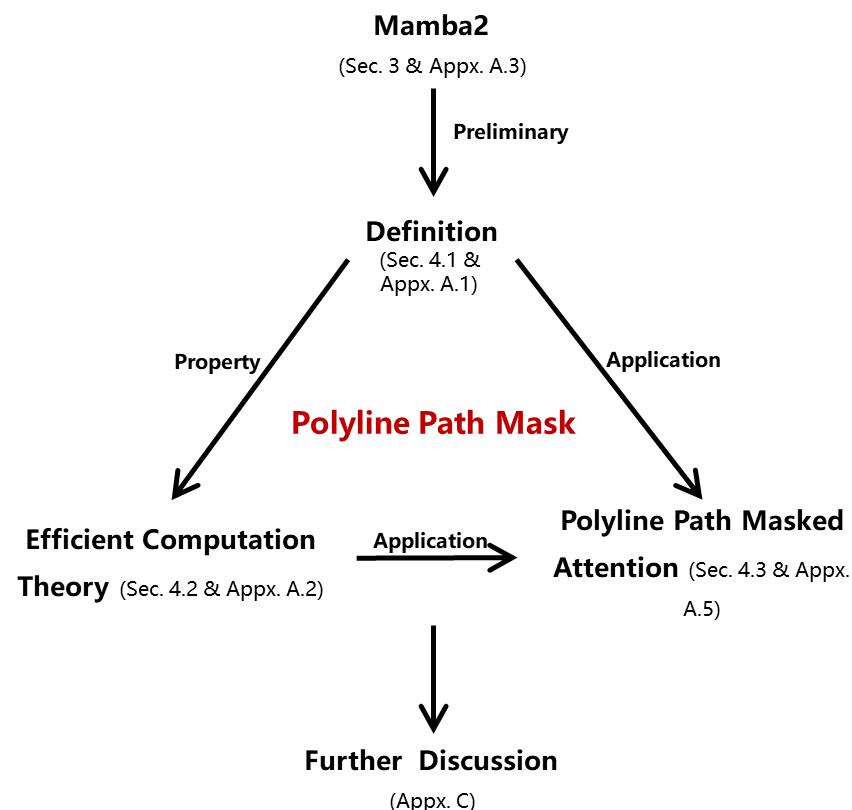


5.1 PPMA: Insight & Contribution

84

- Our insight: Mamba2 core mechanism is **structured mask**
 - **Explicit positional encoding** through the recursive propagation mechanism
 - **Semantic continuity awareness** in sequences through the **selective** mechanism

- Our contribution
 - A novel **2D polyline path structured mask**
 - An **efficient algorithm** for the calculation of the polyline path mask
 - **Polyline Path Masked (Sparse) Attention**
 - **SOTA performance** on image classification, object detection, and segmentation tasks compared to SSM-based models and ViTs



5.2 PPMA: Method (Definition of Polyline Path Mask)

85

- 2D polyline path scanning

- Current issue: fail to preserve the distance relationship of 2D tokens
- Our solution: scan the 2D tokens along multiple paths

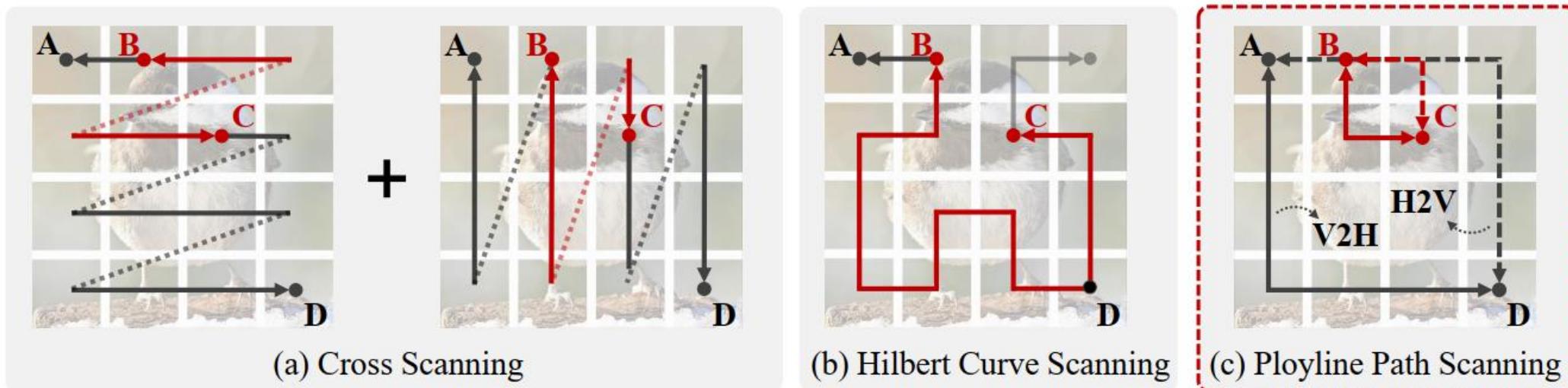


Figure 2: Compared to existing scanning strategies (a) and (b), which flatten 2D tokens into a 1D sequence, our polyline path scanning (c) better preserves the adjacency of 2D tokens.

5.2 PPMA: Method (Definition of Polyline Path Mask)

86

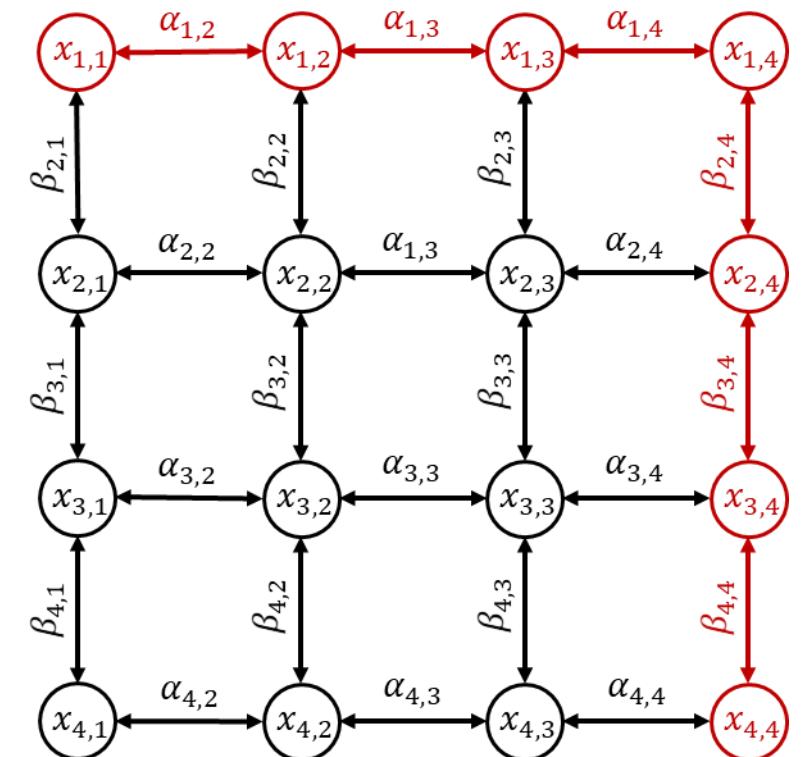
- 2D polyline path mask
 - Learn the **horizontal** factor and **vertical** decay factors

$$\begin{cases} \alpha_{i,j} = \exp(-\text{Softplus}(\text{MLP}_\alpha(x_{i,j}))) \in \mathbb{R}^1 (0 \sim 1) \\ \beta_{i,j} = \exp(-\text{Softplus}(\text{MLP}_\beta(x_{i,j}))) \in \mathbb{R}^1 (0 \sim 1) \end{cases}$$

- Calculate the decay weight of each polyline path, i.e., for path from $x_{k,l}$ to $x_{i,j}$:

$$\mathcal{L}_{i,j,k,l} = \alpha_{i,j:l} \beta_{i:k,l}$$

$$\alpha_{i,j:l} = \begin{cases} \alpha_{i,j+1} \times \cdots \times \alpha_{i,l} & \text{if } j < l \\ 1 & \text{if } j = l \\ \alpha_{i,l+1} \times \cdots \times \alpha_{i,j} & \text{if } j > l \end{cases}, \quad \beta_{i:k,l} = \begin{cases} \beta_{i+1,l} \times \cdots \times \beta_{k,l} & \text{if } i < k \\ 1 & \text{if } i = k \\ \beta_{k+1,l} \times \cdots \times \beta_{i,l} & \text{if } i > k \end{cases}$$



$$\mathcal{L}_{1,1,4,4} = \alpha_{1,2} \alpha_{1,3} \alpha_{1,4} \beta_{2,4} \beta_{3,4} \beta_{4,4}$$

5.2 PPMA: Method (Definition of Polyline Path Mask)

87

- 2D polyline path mask

- Learn the horizontal factor and vertical decay factors
- Calculate the decay weight of each polyline path, i.e., for path from $x_{k,l}$ to $x_{i,j}$:

$$\mathcal{L}_{i,j,k,l} = \alpha_{i,j:l} \beta_{i:k,l}$$

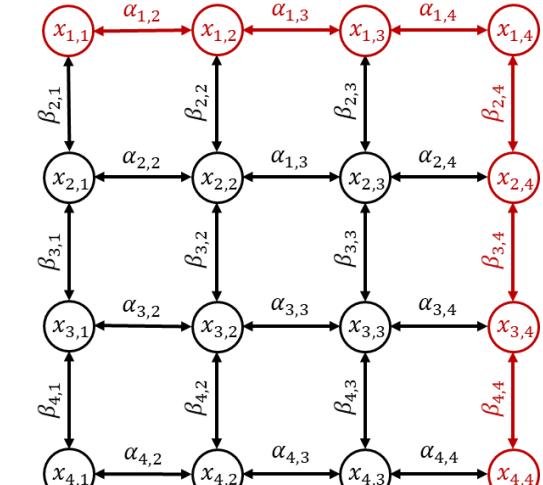
- Combine bidirectional vertical-then-horizontal path and horizontal-then-vertical path:

$$\mathcal{L}^{2D} = \mathcal{L} + \tilde{\mathcal{L}}, \quad \tilde{\mathcal{L}}_{i,j,k,l} = \alpha_{k,j:l} \beta_{i:k,j} = \mathcal{L}_{k,l,i,j}$$

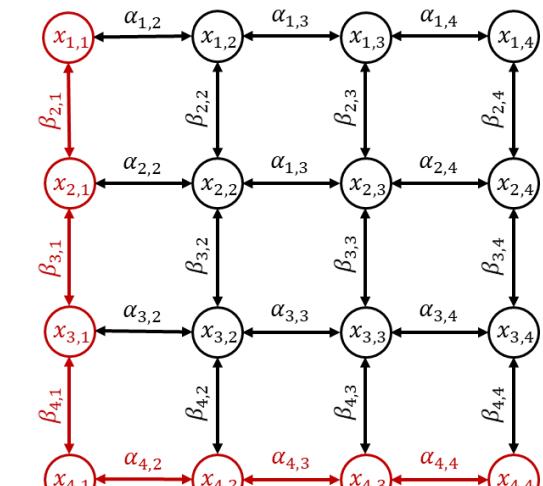
- Unfold 4D tensors $\mathcal{L}^{2D} \in \mathbb{R}^{H \times W \times H \times W}$ to 2D matrix $\mathbf{L}^{2D} \in \mathbb{R}^{HW \times HW}$:

$$\mathbf{L}^{2D} = \text{unfold}(\mathcal{L}^{2D}), \quad \mathbf{L}_{(i-1) \times W + j, (k-1) \times W + l,}^{2D} = \mathcal{L}_{i,j,k,l}^{2D}$$

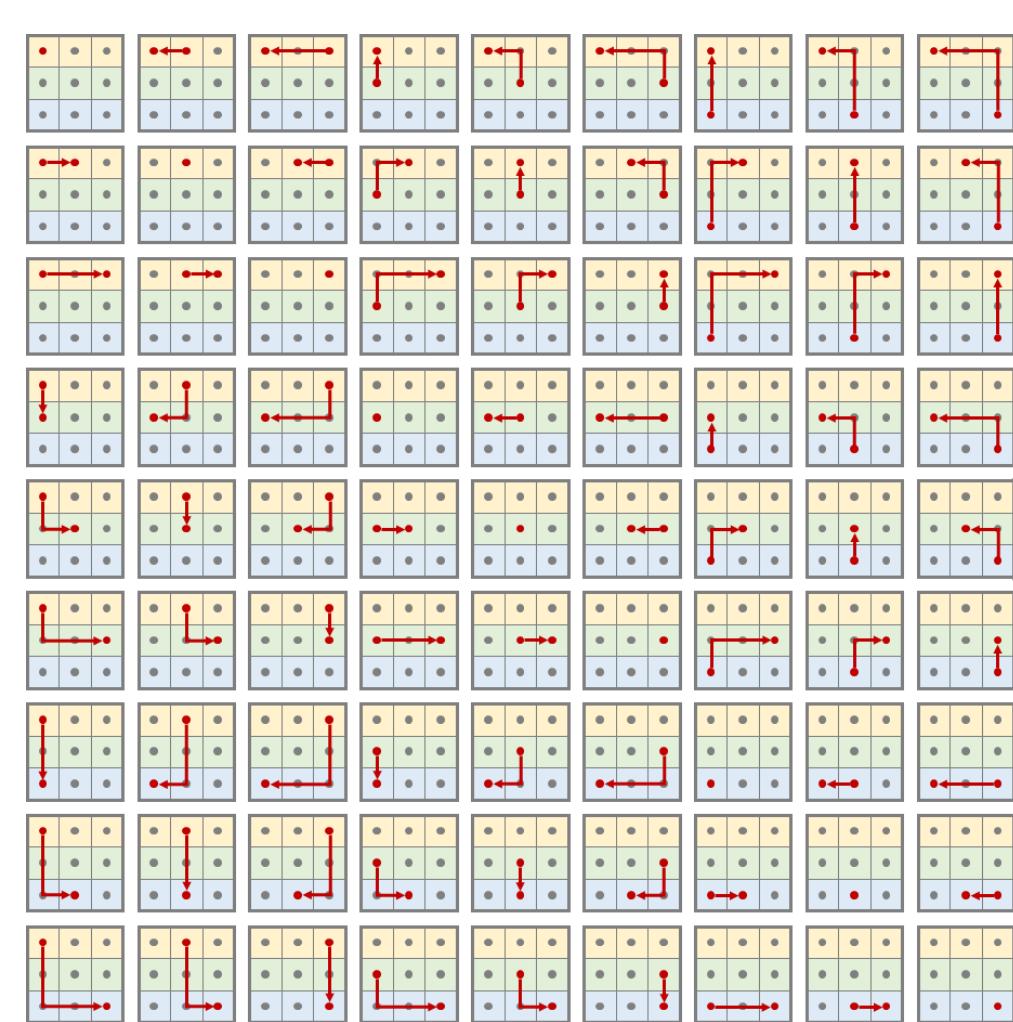
$$\mathbf{L} = \text{unfold}(\mathcal{L})$$



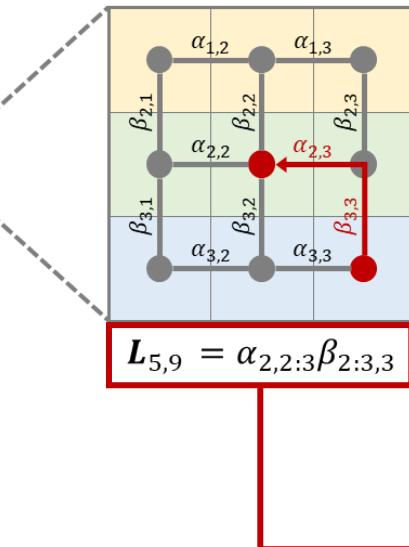
(a) V2H polyline path $\mathcal{L}_{1,1,4,4}$



(b) H2V polyline path $\tilde{\mathcal{L}}_{1,1,4,4}$



(c) An illustration of the V2H polyline path mask on a 3×3 grid (with a total of 9 tokens).



$$L_{5,9} = \alpha_{2,2;3} \beta_{2,3,3}$$

$$L = \begin{bmatrix} \alpha_{1,1:1}\beta_{1:1,1} & \alpha_{1,1:2}\beta_{1:1,2} & \alpha_{1,1:3}\beta_{1:1,3} & \alpha_{1,1:1}\beta_{1:2,1} & \alpha_{1,1:2}\beta_{1:2,2} & \alpha_{1,1:3}\beta_{1:2,3} & \alpha_{1,1:1}\beta_{1:3,1} & \alpha_{1,1:2}\beta_{1:3,2} & \alpha_{1,1:3}\beta_{1:3,3} \\ \alpha_{1,2:1}\beta_{1:1,1} & \alpha_{1,2:2}\beta_{1:1,2} & \alpha_{1,2:3}\beta_{1:1,3} & \alpha_{1,2:1}\beta_{1:2,1} & \alpha_{1,2:2}\beta_{1:2,2} & \alpha_{1,2:3}\beta_{1:2,3} & \alpha_{1,2:1}\beta_{1:3,1} & \alpha_{1,2:2}\beta_{1:3,2} & \alpha_{1,2:3}\beta_{1:3,3} \\ \alpha_{1,3:1}\beta_{1:1,1} & \alpha_{1,3:2}\beta_{1:1,2} & \alpha_{1,3:3}\beta_{1:1,3} & \alpha_{1,3:1}\beta_{1:2,1} & \alpha_{1,3:2}\beta_{1:2,2} & \alpha_{1,3:3}\beta_{1:2,3} & \alpha_{1,3:1}\beta_{1:3,1} & \alpha_{1,3:2}\beta_{1:3,2} & \alpha_{1,3:3}\beta_{1:3,3} \\ \alpha_{2,1:1}\beta_{2:1,1} & \alpha_{2,1:2}\beta_{2:1,2} & \alpha_{2,1:3}\beta_{2:1,3} & \alpha_{2,1:1}\beta_{2:2,1} & \alpha_{2,1:2}\beta_{2:2,2} & \alpha_{2,1:3}\beta_{2:2,3} & \alpha_{2,1:1}\beta_{2:3,1} & \alpha_{2,1:2}\beta_{2:3,2} & \alpha_{2,1:3}\beta_{2:3,3} \\ \alpha_{2,2:1}\beta_{2:1,1} & \alpha_{2,2:2}\beta_{2:1,2} & \alpha_{2,2:3}\beta_{2:1,3} & \alpha_{2,2:1}\beta_{2:2,1} & \alpha_{2,2:2}\beta_{2:2,2} & \alpha_{2,2:3}\beta_{2:2,3} & \alpha_{2,2:1}\beta_{2:3,1} & \alpha_{2,2:2}\beta_{2:3,2} & \alpha_{2,2:3}\beta_{2:3,3} \\ \alpha_{2,3:1}\beta_{2:1,1} & \alpha_{2,3:2}\beta_{2:1,2} & \alpha_{2,3:3}\beta_{2:1,3} & \alpha_{2,3:1}\beta_{2:2,1} & \alpha_{2,3:2}\beta_{2:2,2} & \alpha_{2,3:3}\beta_{2:2,3} & \alpha_{2,3:1}\beta_{2:3,1} & \alpha_{2,3:2}\beta_{2:3,2} & \alpha_{2,3:3}\beta_{2:3,3} \\ \alpha_{3,1:1}\beta_{3:1,1} & \alpha_{3,1:2}\beta_{3:1,2} & \alpha_{3,1:3}\beta_{3:1,3} & \alpha_{3,1:1}\beta_{3:2,1} & \alpha_{3,1:2}\beta_{3:2,2} & \alpha_{3,1:3}\beta_{3:2,3} & \alpha_{3,1:1}\beta_{3:3,1} & \alpha_{3,1:2}\beta_{3:3,2} & \alpha_{3,1:3}\beta_{3:3,3} \\ \alpha_{3,2:1}\beta_{3:1,1} & \alpha_{3,2:2}\beta_{3:1,2} & \alpha_{3,2:3}\beta_{3:1,3} & \alpha_{3,2:1}\beta_{3:2,1} & \alpha_{3,2:2}\beta_{3:2,2} & \alpha_{3,2:3}\beta_{3:2,3} & \alpha_{3,2:1}\beta_{3:3,1} & \alpha_{3,2:2}\beta_{3:3,2} & \alpha_{3,2:3}\beta_{3:3,3} \\ \alpha_{3,3:1}\beta_{3:1,1} & \alpha_{3,3:2}\beta_{3:1,2} & \alpha_{3,3:3}\beta_{3:1,3} & \alpha_{3,3:1}\beta_{3:2,1} & \alpha_{3,3:2}\beta_{3:2,2} & \alpha_{3,3:3}\beta_{3:2,3} & \alpha_{3,3:1}\beta_{3:3,1} & \alpha_{3,3:2}\beta_{3:3,2} & \alpha_{3,3:3}\beta_{3:3,3} \end{bmatrix}$$

5.2 PPMA: Method (Efficient Computation Theory)

- Efficient Computation of Polyline Path Mask

- **Naive Computation:** the polyline mask $L \in \mathbb{R}^{N \times N}$ is large in size and each element $\mathcal{L}_{i,j,k,l} = \alpha_{i,j:l}\beta_{i:k,l}$ require numerous multiplications, resulting in a total complexity of $\mathcal{O}(N^{\frac{5}{2}})$

$$\mathcal{L}_{i,j,k,l} = \alpha_{i,j:l}\beta_{i:k,l}, \quad \text{where } i, k = 0, 1, \dots, H - 1; \quad k, l = 0, 1, \dots, W - 1$$

$$L = \begin{bmatrix} \alpha_{1,1:1}\beta_{1:1,1} & \alpha_{1,1:2}\beta_{1:1,2} & \alpha_{1,1:3}\beta_{1:1,3} & \alpha_{1,1:1}\beta_{1:2,1} & \alpha_{1,1:2}\beta_{1:2,2} & \alpha_{1,1:3}\beta_{1:2,3} & \alpha_{1,1:1}\beta_{1:3,1} & \alpha_{1,1:2}\beta_{1:3,2} & \alpha_{1,1:3}\beta_{1:3,3} \\ \alpha_{1,2:1}\beta_{1:1,1} & \alpha_{1,2:2}\beta_{1:1,2} & \alpha_{1,2:3}\beta_{1:1,3} & \alpha_{1,2:1}\beta_{1:2,1} & \alpha_{1,2:2}\beta_{1:2,2} & \alpha_{1,2:3}\beta_{1:2,3} & \alpha_{1,2:1}\beta_{1:3,1} & \alpha_{1,2:2}\beta_{1:3,2} & \alpha_{1,2:3}\beta_{1:3,3} \\ \alpha_{1,3:1}\beta_{1:1,1} & \alpha_{1,3:2}\beta_{1:1,2} & \alpha_{1,3:3}\beta_{1:1,3} & \alpha_{1,3:1}\beta_{1:2,1} & \alpha_{1,3:2}\beta_{1:2,2} & \alpha_{1,3:3}\beta_{1:2,3} & \alpha_{1,3:1}\beta_{1:3,1} & \alpha_{1,3:2}\beta_{1:3,2} & \alpha_{1,3:3}\beta_{1:3,3} \\ \alpha_{2,1:1}\beta_{2:1,1} & \alpha_{2,1:2}\beta_{2:1,2} & \alpha_{2,1:3}\beta_{2:1,3} & \alpha_{2,1:1}\beta_{2:2,1} & \alpha_{2,1:2}\beta_{2:2,2} & \alpha_{2,1:3}\beta_{2:2,3} & \alpha_{2,1:1}\beta_{2:3,1} & \alpha_{2,1:2}\beta_{2:3,2} & \alpha_{2,1:3}\beta_{2:3,3} \\ \alpha_{2,2:1}\beta_{2:1,1} & \alpha_{2,2:2}\beta_{2:1,2} & \alpha_{2,2:3}\beta_{2:1,3} & \alpha_{2,2:1}\beta_{2:2,1} & \alpha_{2,2:2}\beta_{2:2,2} & \alpha_{2,2:3}\beta_{2:2,3} & \alpha_{2,2:1}\beta_{2:3,1} & \alpha_{2,2:2}\beta_{2:3,2} & \alpha_{2,2:3}\beta_{2:3,3} \\ \alpha_{2,3:1}\beta_{2:1,1} & \alpha_{2,3:2}\beta_{2:1,2} & \alpha_{2,3:3}\beta_{2:1,3} & \alpha_{2,3:1}\beta_{2:2,1} & \alpha_{2,3:2}\beta_{2:2,2} & \alpha_{2,3:3}\beta_{2:2,3} & \alpha_{2,3:1}\beta_{2:3,1} & \alpha_{2,3:2}\beta_{2:3,2} & \alpha_{2,3:3}\beta_{2:3,3} \\ \alpha_{3,1:1}\beta_{3:1,1} & \alpha_{3,1:2}\beta_{3:1,2} & \alpha_{3,1:3}\beta_{3:1,3} & \alpha_{3,1:1}\beta_{3:2,1} & \alpha_{3,1:2}\beta_{3:2,2} & \alpha_{3,1:3}\beta_{3:2,3} & \alpha_{3,1:1}\beta_{3:3,1} & \alpha_{3,1:2}\beta_{3:3,2} & \alpha_{3,1:3}\beta_{3:3,3} \\ \alpha_{3,2:1}\beta_{3:1,1} & \alpha_{3,2:2}\beta_{3:1,2} & \alpha_{3,2:3}\beta_{3:1,3} & \alpha_{3,2:1}\beta_{3:2,1} & \alpha_{3,2:2}\beta_{3:2,2} & \alpha_{3,2:3}\beta_{3:2,3} & \alpha_{3,2:1}\beta_{3:3,1} & \alpha_{3,2:2}\beta_{3:3,2} & \alpha_{3,2:3}\beta_{3:3,3} \\ \alpha_{3,3:1}\beta_{3:1,1} & \alpha_{3,3:2}\beta_{3:1,2} & \alpha_{3,3:3}\beta_{3:1,3} & \alpha_{3,3:1}\beta_{3:2,1} & \alpha_{3,3:2}\beta_{3:2,2} & \alpha_{3,3:3}\beta_{3:2,3} & \alpha_{3,3:1}\beta_{3:3,1} & \alpha_{3,3:2}\beta_{3:3,2} & \alpha_{3,3:3}\beta_{3:3,3} \end{bmatrix}$$

$\mathbb{R}^{N \times N}$

5.2 PPMA: Method (Efficient Computation Theory)

90

- Efficient Computation of Polyline Path Mask

- Naive Computation:** the polyline mask $L \in \mathbb{R}^{N \times N}$ is large in size and each element $L_{i,j,k,l} = \alpha_{i,j;l}\beta_{i:k,l}$ require numerous multiplications, resulting in a total complexity of $\mathcal{O}(N^{\frac{5}{2}})$
- Efficient Computation:** L can be decomposed as the multiplication of two **sparse matrices**: $L = L^H \times L^V = \hat{L}^H \odot \hat{L}^V$

Theorem 1 (Matrix Decomposition). *For any matrix $M \in \mathbb{R}^{HW \times HW}$ and $\mathcal{M} = \text{fold}(M)$, if for $\forall i, j, k, l, \exists A^i \in \mathbb{R}^{W \times W}$ and $B^l \in \mathbb{R}^{H \times H}$, s.t., $\mathcal{M}_{i,j,k,l} = [A^i]_{j,l} \times [B^l]_{i,k}$, then M can be decomposed as:*

$$M = M^A \times M^B = \hat{M}^A \odot \hat{M}^B, \quad (6)$$

where $M^A, M^B, \hat{M}^A, \hat{M}^B \in \mathbb{R}^{HW \times HW}$, which satisfy

$$M^A = \text{unfold}(\mathcal{M}^A), M^B = \text{unfold}(\mathcal{M}^B), \text{s.t., } \mathcal{M}_{i,:,:k,:}^A = \begin{cases} A^i & k=i \\ 0 & k \neq i \end{cases}, \mathcal{M}_{:,:,j,:l}^B = \begin{cases} B^l & j=l \\ 0 & j \neq l \end{cases}, \quad (7)$$

$$\hat{M}^A = \text{unfold}(\hat{\mathcal{M}}^A), \hat{M}^B = \text{unfold}(\hat{\mathcal{M}}^B), \text{s.t., } \hat{\mathcal{M}}_{i,:,:k,:}^A = A^i, \hat{\mathcal{M}}_{:,:,j,:l}^B = B^l. \quad (8)$$

Corollary 1 (Mask Complexity). *The complexity of directly computing polyline path mask L is $\mathcal{O}(N^{\frac{5}{2}})$, which can be reduced to $\mathcal{O}(N^2)$ by applying Theorem 1, where $N = H \times W$.*

Theorem 1 (Matrix Decomposition). *For any matrix $M \in \mathbb{R}^{HW \times HW}$ and $\mathcal{M} = \text{fold}(M)$, if for $\forall i, j, k, l, \exists A^i \in \mathbb{R}^{W \times W}$ and $B^l \in \mathbb{R}^{H \times H}$, s.t., $\mathcal{M}_{i,j,k,l} = [A^i]_{j,l} \times [B^l]_{i,k}$, then M can be decomposed as:*

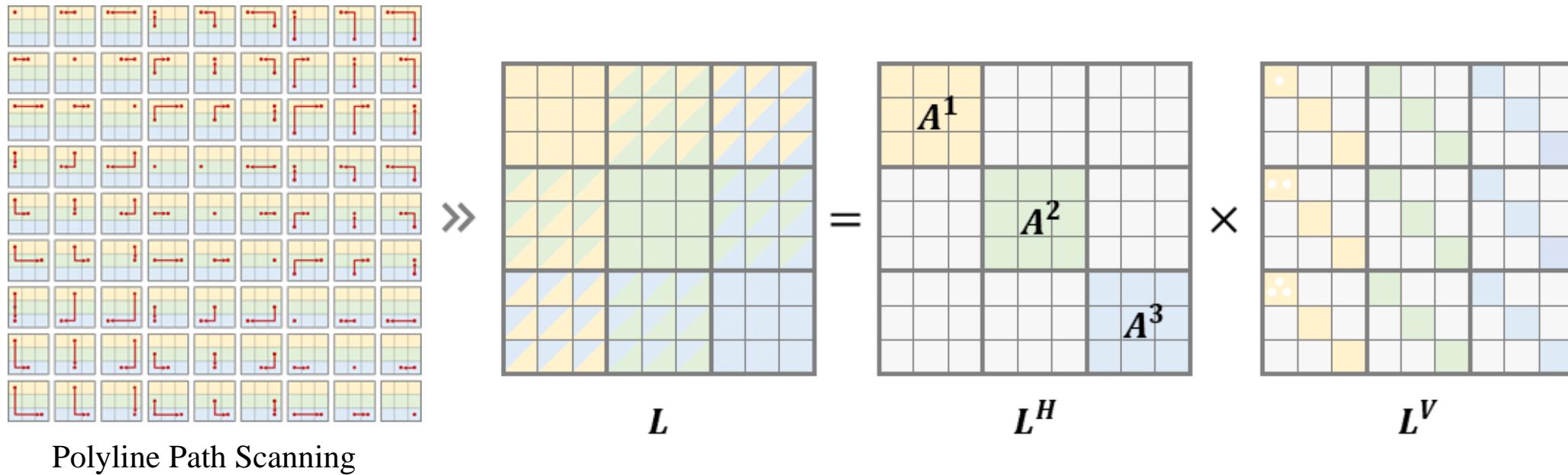
$$M = M^A \times M^B = \hat{M}^A \odot \hat{M}^B, \quad (6)$$

where $M^A, M^B, \hat{M}^A, \hat{M}^B \in \mathbb{R}^{HW \times HW}$, which satisfy

$$M^A = \text{unfold}(\mathcal{M}^A), M^B = \text{unfold}(\mathcal{M}^B), \text{s.t., } \mathcal{M}_{i,:,:k,:}^A = \begin{cases} A^i & k=i \\ 0 & k \neq i \end{cases}, \mathcal{M}_{:,:,j,l}^B = \begin{cases} B^l & j=l \\ 0 & j \neq l \end{cases}, \quad (7)$$

$$\hat{M}^A = \text{unfold}(\hat{\mathcal{M}}^A), \hat{M}^B = \text{unfold}(\hat{\mathcal{M}}^B), \text{s.t., } \hat{\mathcal{M}}_{i,:,:k,:}^A = A^i, \hat{\mathcal{M}}_{:,:,j,l}^B = B^l. \quad (8)$$

The matrix L satisfies the conditions in Theorem 1 with $[A^i]_{j,l} = \alpha_{i,j:l}$ and $[B^l]_{i,k} = \beta_{i:k,l}$.



(a) An illustration of the Polyline Path Mask Decomposition

$$L = \begin{bmatrix} \alpha_{1,1:1}\beta_{1:1,1} & \alpha_{1,1:2}\beta_{1:1,2} & \alpha_{1,1:3}\beta_{1:1,3} & \alpha_{1,1:1}\beta_{1:2,1} & \alpha_{1,1:2}\beta_{1:2,2} & \alpha_{1,1:3}\beta_{1:2,3} & \alpha_{1,1:1}\beta_{1:3,1} & \alpha_{1,1:2}\beta_{1:3,2} & \alpha_{1,1:3}\beta_{1:3,3} \\ \alpha_{1,2:1}\beta_{1:1,1} & \alpha_{1,2:2}\beta_{1:1,2} & \alpha_{1,2:3}\beta_{1:1,3} & \alpha_{1,2:1}\beta_{1:2,1} & \alpha_{1,2:2}\beta_{1:2,2} & \alpha_{1,2:3}\beta_{1:2,3} & \alpha_{1,2:1}\beta_{1:3,1} & \alpha_{1,2:2}\beta_{1:3,2} & \alpha_{1,2:3}\beta_{1:3,3} \\ \alpha_{1,3:1}\beta_{1:1,1} & \alpha_{1,3:2}\beta_{1:1,2} & \alpha_{1,3:3}\beta_{1:1,3} & \alpha_{1,3:1}\beta_{1:2,1} & \alpha_{1,3:2}\beta_{1:2,2} & \alpha_{1,3:3}\beta_{1:2,3} & \alpha_{1,3:1}\beta_{1:3,1} & \alpha_{1,3:2}\beta_{1:3,2} & \alpha_{1,3:3}\beta_{1:3,3} \\ \alpha_{2,1:1}\beta_{2:1,1} & \alpha_{2,1:2}\beta_{2:1,2} & \alpha_{2,1:3}\beta_{2:1,3} & \alpha_{2,1:1}\beta_{2:2,1} & \alpha_{2,1:2}\beta_{2:2,2} & \alpha_{2,1:3}\beta_{2:2,3} & \alpha_{2,1:1}\beta_{2:3,1} & \alpha_{2,1:2}\beta_{2:3,2} & \alpha_{2,1:3}\beta_{2:3,3} \\ \alpha_{2,2:1}\beta_{2:1,1} & \alpha_{2,2:2}\beta_{2:1,2} & \alpha_{2,2:3}\beta_{2:1,3} & \alpha_{2,2:1}\beta_{2:2,1} & \alpha_{2,2:2}\beta_{2:2,2} & \alpha_{2,2:3}\beta_{2:2,3} & \alpha_{2,2:1}\beta_{2:3,1} & \alpha_{2,2:2}\beta_{2:3,2} & \alpha_{2,2:3}\beta_{2:3,3} \\ \alpha_{2,3:1}\beta_{2:1,1} & \alpha_{2,3:2}\beta_{2:1,2} & \alpha_{2,3:3}\beta_{2:1,3} & \alpha_{2,3:1}\beta_{2:2,1} & \alpha_{2,3:2}\beta_{2:2,2} & \alpha_{2,3:3}\beta_{2:2,3} & \alpha_{2,3:1}\beta_{2:3,1} & \alpha_{2,3:2}\beta_{2:3,2} & \alpha_{2,3:3}\beta_{2:3,3} \\ \alpha_{3,1:1}\beta_{3:1,1} & \alpha_{3,1:2}\beta_{3:1,2} & \alpha_{3,1:3}\beta_{3:1,3} & \alpha_{3,1:1}\beta_{3:2,1} & \alpha_{3,1:2}\beta_{3:2,2} & \alpha_{3,1:3}\beta_{3:2,3} & \alpha_{3,1:1}\beta_{3:3,1} & \alpha_{3,1:2}\beta_{3:3,2} & \alpha_{3,1:3}\beta_{3:3,3} \\ \alpha_{3,2:1}\beta_{3:1,1} & \alpha_{3,2:2}\beta_{3:1,2} & \alpha_{3,2:3}\beta_{3:1,3} & \alpha_{3,2:1}\beta_{3:2,1} & \alpha_{3,2:2}\beta_{3:2,2} & \alpha_{3,2:3}\beta_{3:2,3} & \alpha_{3,2:1}\beta_{3:3,1} & \alpha_{3,2:2}\beta_{3:3,2} & \alpha_{3,2:3}\beta_{3:3,3} \\ \alpha_{3,3:1}\beta_{3:1,1} & \alpha_{3,3:2}\beta_{3:1,2} & \alpha_{3,3:3}\beta_{3:1,3} & \alpha_{3,3:1}\beta_{3:2,1} & \alpha_{3,3:2}\beta_{3:2,2} & \alpha_{3,3:3}\beta_{3:2,3} & \alpha_{3,3:1}\beta_{3:3,1} & \alpha_{3,3:2}\beta_{3:3,2} & \alpha_{3,3:3}\beta_{3:3,3} \end{bmatrix}$$

Complexity: $\mathcal{O}(N^{\frac{5}{2}})$

$$\begin{array}{ccccccccc} \alpha_{1,1:1} & \alpha_{1,1:2} & \alpha_{1,1:3} & 0 & 0 & 0 & 0 & 0 & 0 \\ \alpha_{1,2:1} & \alpha_{1,2:2} & \alpha_{1,2:3} & 0 & 0 & 0 & 0 & 0 & 0 \\ \alpha_{1,3:1} & \alpha_{1,3:2} & \alpha_{1,3:3} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha_{2,1:1} & \alpha_{2,1:2} & \alpha_{2,1:3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha_{2,2:1} & \alpha_{2,2:2} & \alpha_{2,2:3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha_{2,3:1} & \alpha_{2,3:2} & \alpha_{2,3:3} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \alpha_{3,1:1} & \alpha_{3,1:2} & \alpha_{3,1:3} \\ 0 & 0 & 0 & \alpha_{3,2:1} & \alpha_{3,2:2} & \alpha_{3,2:3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha_{3,3:1} & \alpha_{3,3:2} & \alpha_{3,3:3} & 0 & 0 & 0 \end{array} \times \begin{array}{ccccccccc} \beta_{1:1,1} & 0 & 0 & \beta_{1:2,1} & 0 & 0 & \beta_{1:3,1} & 0 & 0 \\ 0 & \beta_{1:1,2} & 0 & 0 & \beta_{1:2,2} & 0 & 0 & \beta_{1:3,2} & 0 \\ 0 & 0 & \beta_{1:1,3} & 0 & 0 & \beta_{1:2,3} & 0 & 0 & \beta_{1:3,3} \\ \beta_{2:1,1} & 0 & 0 & \beta_{2:2,1} & 0 & 0 & \beta_{2:3,1} & 0 & 0 \\ 0 & \beta_{2:1,2} & 0 & 0 & \beta_{2:2,2} & 0 & 0 & \beta_{2:3,2} & 0 \\ 0 & 0 & \beta_{2:1,3} & 0 & 0 & \beta_{2:2,3} & 0 & 0 & \beta_{2:3,3} \\ \beta_{3:1,1} & 0 & 0 & \beta_{3:2,1} & 0 & 0 & \beta_{3:3,1} & 0 & 0 \\ 0 & \beta_{3:1,2} & 0 & 0 & \beta_{3:2,2} & 0 & 0 & \beta_{3:3,2} & 0 \\ 0 & 0 & \beta_{3:1,3} & 0 & 0 & \beta_{3:2,3} & 0 & 0 & \beta_{3:3,3} \end{array}$$

$$L^H: \mathcal{O}(N^{\frac{3}{2}})$$

$$L^V: \mathcal{O}(N^{\frac{3}{2}})$$

Complexity: $\mathcal{O}(N^2)$

5.2 PPMA: Method (Efficient Computation Theory)

93

- Efficient Computation of Polyline Path Mask Matrix Multiplication

- Naive Computation:** for a rank-N matrix $\mathbf{L} \in \mathbb{R}^{N \times N}$ and a vector $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{L}\mathbf{x}$ requires a complexity of $\mathcal{O}(N^2)$
- Efficient Computation:** for decomposed polyline path mask, $\mathbf{L}\mathbf{x} = \mathbf{L}^H \times \mathbf{L}^V \times \mathbf{x}$ requires a complexity of $\mathcal{O}(N)$

Theorem 2 (Efficient Matrix Multiplication). *For matrices $\mathbf{M}^A, \mathbf{M}^B$ defined in Theorem 1 $\forall \mathbf{x} \in \mathbb{R}^{HW}$, the following equation holds:*

$$\mathbf{y} = \mathbf{M}^A \times \mathbf{M}^B \times \mathbf{x} \Leftrightarrow \mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}, \mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}, \quad (9)$$

where $\mathbf{y} \in \mathbb{R}^{HW}$, $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$, $\mathbf{Y} = \text{unvec}(\mathbf{y}) \in \mathbb{R}^{H \times W}$, $\mathbf{Z} \in \mathbb{R}^{H \times W}$, and the operator $\text{vec}(\cdot)$ vectorizes a matrix by stacking its columns and $\text{unvec}(\cdot)$ is its inverse operator.

Algorithm 1: Efficient Masked Attention Computation.

Input: decay factors α, β of \mathbf{L} , vector $\mathbf{x} \in \mathbb{R}^{HW}$;

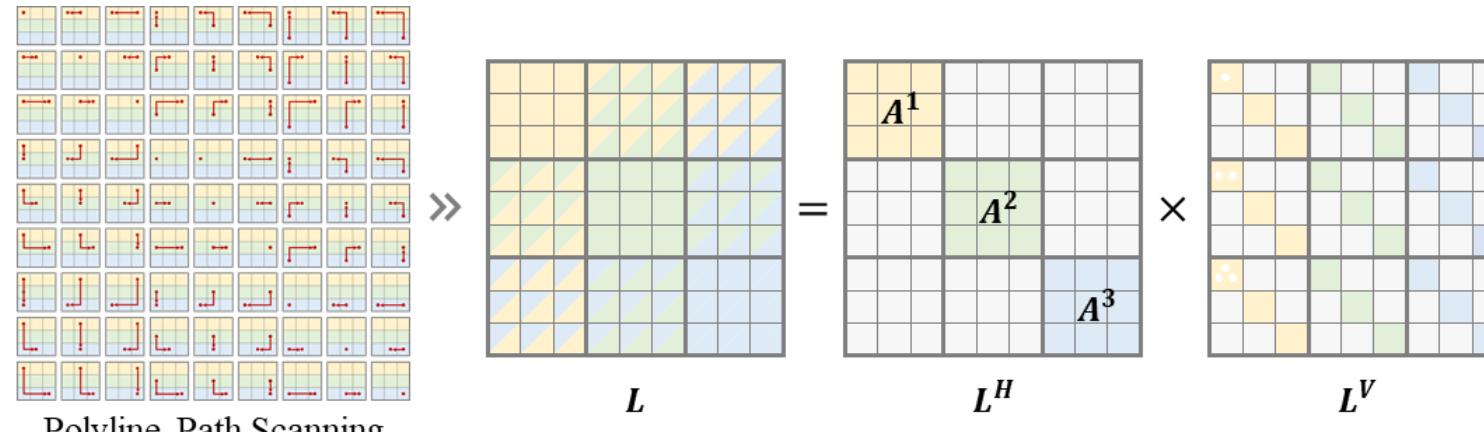
- 1: Compute $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$;
- 2: Compute $\mathbf{B}^l \in \mathbb{R}^{H \times H}$, where for $l = 1:W$, $[\mathbf{B}^l]_{i,k} = \beta_{i:k,l}$;
- 3: Compute $\mathbf{Z} \in \mathbb{R}^{H \times W}$, where $\mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}$;
- 4: Compute $\mathbf{A}^i \in \mathbb{R}^{W \times W}$, where for $i = 1:H$, $[\mathbf{A}^i]_{j,l} = \alpha_{i,j:l}$;
- 5: Compute $\mathbf{Y} \in \mathbb{R}^{H \times W}$, where $\mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}$;

Output: $\mathbf{y} = \text{vec}(\mathbf{Y})$;

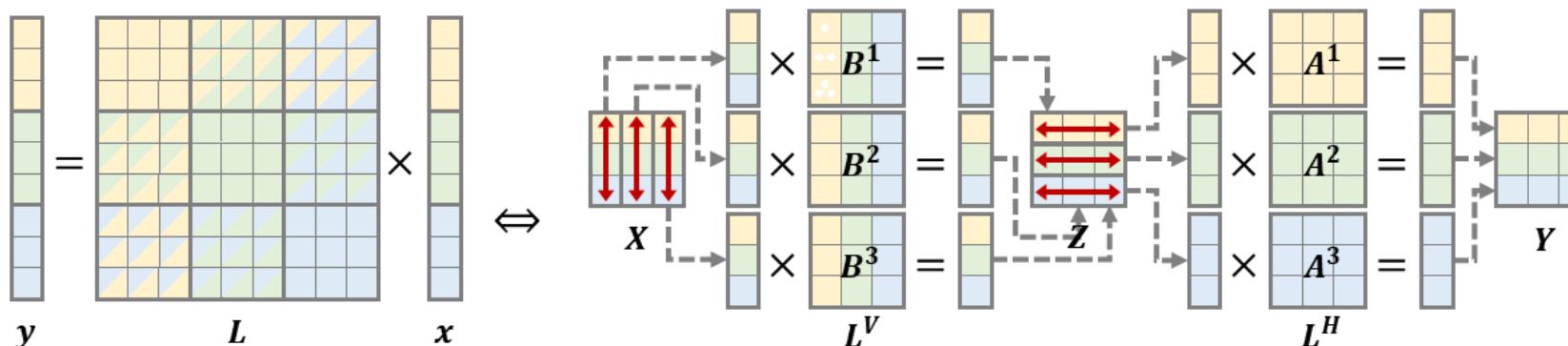
Theorem 2 (Efficient Matrix Multiplication). *For matrices M^A, M^B defined in Theorem 1 $\forall x \in \mathbb{R}^{HW}$, the following equation holds:*

$$\mathbf{y} = M^A \times M^B \times x \Leftrightarrow \mathbf{Z}_{:,l} = B^l \times \mathbf{X}_{:,l}, \mathbf{Y}_{i,:} = A^i \times \mathbf{Z}_{i,:}, \quad (9)$$

where $\mathbf{y} \in \mathbb{R}^{HW}$, $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$, $\mathbf{Y} = \text{unvec}(\mathbf{y}) \in \mathbb{R}^{H \times W}$, $\mathbf{Z} \in \mathbb{R}^{H \times W}$, and the operator $\text{vec}(\cdot)$ vectorizes a matrix by stacking its columns and $\text{unvec}(\cdot)$ is its inverse operator.



(a) Polyline Path Mask Decomposition



(b) Polyline Path Mask Multiplication

Theorem 2 (Efficient Matrix Multiplication). *For matrices $\mathbf{M}^A, \mathbf{M}^B$ defined in Theorem 1 $\forall \mathbf{x} \in \mathbb{R}^{HW}$, the following equation holds:*

$$\mathbf{y} = \mathbf{M}^A \times \mathbf{M}^B \times \mathbf{x} \Leftrightarrow \mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}, \mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}, \quad (9)$$

where $\mathbf{y} \in \mathbb{R}^{HW}$, $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$, $\mathbf{Y} = \text{unvec}(\mathbf{y}) \in \mathbb{R}^{H \times W}$, $\mathbf{Z} \in \mathbb{R}^{H \times W}$, and the operator $\text{vec}(\cdot)$ vectorizes a matrix by stacking its columns and $\text{unvec}(\cdot)$ is its inverse operator.

Algorithm 1: Efficient Masked Attention Computation.

Input: decay factors α, β of \mathbf{L} , vector $\mathbf{x} \in \mathbb{R}^{HW}$;

- 1: Compute $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$;
- 2: Compute $\mathbf{B}^l \in \mathbb{R}^{H \times H}$, where for $l=1:W$, $[\mathbf{B}^l]_{i,k} = \beta_{i:k,l}$;
- 3: Compute $\mathbf{Z} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}}$;
- 4: Compute $\mathbf{A}^i \in \mathbb{R}^{W \times W}$, where for $i=1:H$, $[\mathbf{A}^i]_{j,l} = \alpha_{i,j:l}$;
- 5: Compute $\mathbf{Y} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}}$;

Complexity: $\mathcal{O}(H^2W) \rightarrow \mathcal{O}(HW)$

Complexity: $\mathcal{O}(HW^2) \rightarrow \mathcal{O}(HW)$

Output: $\mathbf{y} = \text{vec}(\mathbf{Y})$;

$$L = \begin{bmatrix} \alpha_{1,1:1}\beta_{1:1,1} & \alpha_{1,1:2}\beta_{1:1,2} & \alpha_{1,1:3}\beta_{1:1,3} & \alpha_{1,1:1}\beta_{1:2,1} & \alpha_{1,1:2}\beta_{1:2,2} & \alpha_{1,1:3}\beta_{1:2,3} & \alpha_{1,1:1}\beta_{1:3,1} & \alpha_{1,1:2}\beta_{1:3,2} & \alpha_{1,1:3}\beta_{1:3,3} \\ \alpha_{1,2:1}\beta_{1:1,1} & \alpha_{1,2:2}\beta_{1:1,2} & \alpha_{1,2:3}\beta_{1:1,3} & \alpha_{1,2:1}\beta_{1:2,1} & \alpha_{1,2:2}\beta_{1:2,2} & \alpha_{1,2:3}\beta_{1:2,3} & \alpha_{1,2:1}\beta_{1:3,1} & \alpha_{1,2:2}\beta_{1:3,2} & \alpha_{1,2:3}\beta_{1:3,3} \\ \alpha_{1,3:1}\beta_{1:1,1} & \alpha_{1,3:2}\beta_{1:1,2} & \alpha_{1,3:3}\beta_{1:1,3} & \alpha_{1,3:1}\beta_{1:2,1} & \alpha_{1,3:2}\beta_{1:2,2} & \alpha_{1,3:3}\beta_{1:2,3} & \alpha_{1,3:1}\beta_{1:3,1} & \alpha_{1,3:2}\beta_{1:3,2} & \alpha_{1,3:3}\beta_{1:3,3} \\ \alpha_{2,1:1}\beta_{2:1,1} & \alpha_{2,1:2}\beta_{2:1,2} & \alpha_{2,1:3}\beta_{2:1,3} & \alpha_{2,1:1}\beta_{2:2,1} & \alpha_{2,1:2}\beta_{2:2,2} & \alpha_{2,1:3}\beta_{2:2,3} & \alpha_{2,1:1}\beta_{2:3,1} & \alpha_{2,1:2}\beta_{2:3,2} & \alpha_{2,1:3}\beta_{2:3,3} \\ \alpha_{2,2:1}\beta_{2:1,1} & \alpha_{2,2:2}\beta_{2:1,2} & \alpha_{2,2:3}\beta_{2:1,3} & \alpha_{2,2:1}\beta_{2:2,1} & \alpha_{2,2:2}\beta_{2:2,2} & \alpha_{2,2:3}\beta_{2:2,3} & \alpha_{2,2:1}\beta_{2:3,1} & \alpha_{2,2:2}\beta_{2:3,2} & \alpha_{2,2:3}\beta_{2:3,3} \\ \alpha_{2,3:1}\beta_{2:1,1} & \alpha_{2,3:2}\beta_{2:1,2} & \alpha_{2,3:3}\beta_{2:1,3} & \alpha_{2,3:1}\beta_{2:2,1} & \alpha_{2,3:2}\beta_{2:2,2} & \alpha_{2,3:3}\beta_{2:2,3} & \alpha_{2,3:1}\beta_{2:3,1} & \alpha_{2,3:2}\beta_{2:3,2} & \alpha_{2,3:3}\beta_{2:3,3} \\ \alpha_{3,1:1}\beta_{3:1,1} & \alpha_{3,1:2}\beta_{3:1,2} & \alpha_{3,1:3}\beta_{3:1,3} & \alpha_{3,1:1}\beta_{3:2,1} & \alpha_{3,1:2}\beta_{3:2,2} & \alpha_{3,1:3}\beta_{3:2,3} & \alpha_{3,1:1}\beta_{3:3,1} & \alpha_{3,1:2}\beta_{3:3,2} & \alpha_{3,1:3}\beta_{3:3,3} \\ \alpha_{3,2:1}\beta_{3:1,1} & \alpha_{3,2:2}\beta_{3:1,2} & \alpha_{3,2:3}\beta_{3:1,3} & \alpha_{3,2:1}\beta_{3:2,1} & \alpha_{3,2:2}\beta_{3:2,2} & \alpha_{3,2:3}\beta_{3:2,3} & \alpha_{3,2:1}\beta_{3:3,1} & \alpha_{3,2:2}\beta_{3:3,2} & \alpha_{3,2:3}\beta_{3:3,3} \\ \alpha_{3,3:1}\beta_{3:1,1} & \alpha_{3,3:2}\beta_{3:1,2} & \alpha_{3,3:3}\beta_{3:1,3} & \alpha_{3,3:1}\beta_{3:2,1} & \alpha_{3,3:2}\beta_{3:2,2} & \alpha_{3,3:3}\beta_{3:2,3} & \alpha_{3,3:1}\beta_{3:3,1} & \alpha_{3,3:2}\beta_{3:3,2} & \alpha_{3,3:3}\beta_{3:3,3} \end{bmatrix}$$

Theorem 2 (Efficient Matrix Multiplication). *For matrices $\mathbf{M}^A, \mathbf{M}^B$ defined in Theorem 1 $\forall \mathbf{x} \in \mathbb{R}^{HW}$, the following equation holds:*

$$\mathbf{y} = \mathbf{M}^A \times \mathbf{M}^B \times \mathbf{x} \Leftrightarrow \mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}, \mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}, \quad (9)$$

where $\mathbf{y} \in \mathbb{R}^{HW}$, $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$, $\mathbf{Y} = \text{unvec}(\mathbf{y}) \in \mathbb{R}^{H \times W}$, $\mathbf{Z} \in \mathbb{R}^{H \times W}$, and the operator $\text{vec}(\cdot)$ vectorizes a matrix by stacking its columns and $\text{unvec}(\cdot)$ is its inverse operator.

Algorithm 1: Efficient Masked Attention Computation.

Input: decay factors α, β of \mathbf{L} , vector $\mathbf{x} \in \mathbb{R}^{HW}$;

- 1: Compute $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$;
- 2: Compute $\mathbf{B}^l \in \mathbb{R}^{H \times H}$, where for $l=1:W$, $[\mathbf{B}^l]_{i,k} = \beta_{i:k,l}$;
- 3: Compute $\mathbf{Z} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}}$;
- 4: Compute $\mathbf{A}^i \in \mathbb{R}^{W \times W}$, where for $i=1:H$, $[\mathbf{A}^i]_{j,l} = \alpha_{i,j:l}$;
- 5: Compute $\mathbf{Y} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}}$;

Complexity: $\mathcal{O}(H^2W) \rightarrow \mathcal{O}(HW)$

Complexity: $\mathcal{O}(HW^2) \rightarrow \mathcal{O}(HW)$

Output: $\mathbf{y} = \text{vec}(\mathbf{Y})$;

$$L = \begin{bmatrix} \alpha_{1,1:1} & \alpha_{1,1:2} & \alpha_{1,1:3} & 0 & 0 & 0 & 0 & 0 & 0 \\ \alpha_{1,2:1} & \alpha_{1,2:2} & \alpha_{1,2:3} & 0 & 0 & 0 & 0 & 0 & 0 \\ \alpha_{1,3:1} & \alpha_{1,3:2} & \alpha_{1,3:3} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha_{2,1:1} & \alpha_{2,1:2} & \alpha_{2,1:3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha_{2,2:1} & \alpha_{2,2:2} & \alpha_{2,2:3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \alpha_{2,3:1} & \alpha_{2,3:2} & \alpha_{2,3:3} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \alpha_{3,1:1} & \alpha_{3,1:2} & \alpha_{3,1:3} \\ 0 & 0 & 0 & 0 & 0 & 0 & \alpha_{3,2:1} & \alpha_{3,2:2} & \alpha_{3,2:3} \\ 0 & 0 & 0 & 0 & 0 & 0 & \alpha_{3,3:1} & \alpha_{3,3:2} & \alpha_{3,3:3} \end{bmatrix} \times \begin{bmatrix} \beta_{1:1,1} & 0 & 0 & \beta_{1:2,1} & 0 & 0 & \beta_{1:3,1} & 0 & 0 \\ 0 & \beta_{1:1,2} & 0 & 0 & \beta_{1:2,2} & 0 & 0 & \beta_{1:3,2} & 0 \\ 0 & 0 & \beta_{1:1,3} & 0 & 0 & \beta_{1:2,3} & 0 & 0 & \beta_{1:3,3} \\ \beta_{2:1,1} & 0 & 0 & \beta_{2:2,1} & 0 & 0 & \beta_{2:3,1} & 0 & 0 \\ 0 & \beta_{2:1,2} & 0 & 0 & \beta_{2:2,2} & 0 & 0 & \beta_{2:3,2} & 0 \\ 0 & 0 & \beta_{2:1,3} & 0 & 0 & \beta_{2:2,3} & 0 & 0 & \beta_{2:3,3} \\ \beta_{3:1,1} & 0 & 0 & \beta_{3:2,1} & 0 & 0 & \beta_{3:3,1} & 0 & 0 \\ 0 & \beta_{3:1,2} & 0 & 0 & \beta_{3:2,2} & 0 & 0 & \beta_{3:3,2} & 0 \\ 0 & 0 & \beta_{3:1,3} & 0 & 0 & \beta_{3:2,3} & 0 & 0 & \beta_{3:3,3} \end{bmatrix}$$

Theorem 2 (Efficient Matrix Multiplication). *For matrices $\mathbf{M}^A, \mathbf{M}^B$ defined in Theorem 1 $\forall x \in \mathbb{R}^{HW}$, the following equation holds:*

$$\mathbf{y} = \mathbf{M}^A \times \mathbf{M}^B \times x \Leftrightarrow \mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}, \mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}, \quad (9)$$

where $\mathbf{y} \in \mathbb{R}^{HW}$, $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$, $\mathbf{Y} = \text{unvec}(\mathbf{y}) \in \mathbb{R}^{H \times W}$, $\mathbf{Z} \in \mathbb{R}^{H \times W}$, and the operator $\text{vec}(\cdot)$ vectorizes a matrix by stacking its columns and $\text{unvec}(\cdot)$ is its inverse operator.

Algorithm 1: Efficient Masked Attention Computation.

Input: decay factors α, β of \mathbf{L} , vector $\mathbf{x} \in \mathbb{R}^{HW}$;

- 1: Compute $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$;
- 2: Compute $\mathbf{B}^l \in \mathbb{R}^{H \times H}$, where for $l=1:W$, $[\mathbf{B}^l]_{i,k} = \beta_{i:k,l}$;
- 3: Compute $\mathbf{Z} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}}$;
- 4: Compute $\mathbf{A}^i \in \mathbb{R}^{W \times W}$, where for $i=1:H$, $[\mathbf{A}^i]_{j,l} = \alpha_{i,j:l}$;
- 5: Compute $\mathbf{Y} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}}$;

Output: $\mathbf{y} = \text{vec}(\mathbf{Y})$;

Chunkwise algorithm

Complexity: $\mathcal{O}(H^2W) \rightarrow \mathcal{O}(HW)$

Complexity: $\mathcal{O}(HW^2) \rightarrow \mathcal{O}(HW)$

$$L = \begin{bmatrix} \alpha_{1,1:1} & \alpha_{1,1:2} & \alpha_{1,1:3} \\ \alpha_{1,2:1} & \alpha_{1,2:2} & \alpha_{1,2:3} \\ \alpha_{1,3:1} & \alpha_{1,3:2} & \alpha_{1,3:3} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{matrix} \mathbf{A}^1 \\ \mathbf{A}^2 \\ \mathbf{A}^3 \end{matrix} \times \begin{bmatrix} \beta_{1:1,1} & 0 & 0 & \beta_{1:2,1} & 0 & 0 & \beta_{1:3,1} & 0 & \mathbf{B}^1 \\ 0 & \beta_{1:1,2} & 0 & 0 & \beta_{1:2,2} & 0 & 0 & \beta_{1:3,2} & 0 \\ 0 & 0 & \beta_{1:1,3} & 0 & 0 & \beta_{1:2,3} & 0 & 0 & \beta_{1:3,3} \\ \beta_{2:1,1} & 0 & 0 & \beta_{2:2,1} & 0 & 0 & \beta_{2:3,1} & 0 & 0 \\ 0 & \beta_{2:1,2} & 0 & 0 & \beta_{2:2,2} & 0 & 0 & \beta_{2:3,2} & 0 \\ 0 & 0 & \beta_{2:1,3} & 0 & 0 & \beta_{2:2,3} & 0 & 0 & \beta_{2:3,3} \\ \beta_{3:1,1} & 0 & 0 & \beta_{3:2,1} & 0 & 0 & \beta_{3:3,1} & 0 & 0 \\ 0 & \beta_{3:1,2} & 0 & 0 & \beta_{3:2,2} & 0 & 0 & \beta_{3:3,2} & 0 \\ 0 & 0 & \beta_{3:1,3} & 0 & 0 & \beta_{3:2,3} & 0 & 0 & \beta_{3:3,3} \end{bmatrix}$$

Theorem 2 (Efficient Matrix Multiplication). *For matrices $\mathbf{M}^A, \mathbf{M}^B$ defined in Theorem 1 $\forall x \in \mathbb{R}^{HW}$, the following equation holds:*

$$\mathbf{y} = \mathbf{M}^A \times \mathbf{M}^B \times x \Leftrightarrow \mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}, \mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}, \quad (9)$$

where $\mathbf{y} \in \mathbb{R}^{HW}$, $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$, $\mathbf{Y} = \text{unvec}(\mathbf{y}) \in \mathbb{R}^{H \times W}$, $\mathbf{Z} \in \mathbb{R}^{H \times W}$, and the operator $\text{vec}(\cdot)$ vectorizes a matrix by stacking its columns and $\text{unvec}(\cdot)$ is its inverse operator.

Algorithm 1: Efficient Masked Attention Computation.

Input: decay factors α, β of \mathbf{L} , vector $\mathbf{x} \in \mathbb{R}^{HW}$;

- 1: Compute $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$;
- 2: Compute $\mathbf{B}^l \in \mathbb{R}^{H \times H}$, where for $l=1:W$, $[\mathbf{B}^l]_{i,k} = \beta_{i:k,l}$;
- 3: Compute $\mathbf{Z} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}}$;
- 4: Compute $\mathbf{A}^i \in \mathbb{R}^{W \times W}$, where for $i=1:H$, $[\mathbf{A}^i]_{j,l} = \alpha_{i,j:l}$;
- 5: Compute $\mathbf{Y} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}}$;

Output: $\mathbf{y} = \text{vec}(\mathbf{Y})$;

Chunkwise algorithm

Complexity: $\mathcal{O}(H^2W) \rightarrow \mathcal{O}(HW)$

Complexity: $\mathcal{O}(HW^2) \rightarrow \mathcal{O}(HW)$

$$L = \begin{bmatrix} \alpha_{1,1:1} & \alpha_{1,1:2} & \alpha_{1,1:3} \\ \alpha_{1,2:1} & \alpha_{1,2:2} & \alpha_{1,2:3} \\ \alpha_{1,3:1} & \alpha_{1,3:2} & \alpha_{1,3:3} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} \alpha_{2,1:1} & \alpha_{2,1:2} & \alpha_{2,1:3} \\ \alpha_{2,2:1} & \alpha_{2,2:2} & \alpha_{2,2:3} \\ \alpha_{2,3:1} & \alpha_{2,3:2} & \alpha_{2,3:3} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} \alpha_{3,1:1} & \alpha_{3,1:2} & \alpha_{3,1:3} \\ \alpha_{3,2:1} & \alpha_{3,2:2} & \alpha_{3,2:3} \\ \alpha_{3,3:1} & \alpha_{3,3:2} & \alpha_{3,3:3} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} \beta_{1:1,1} & 0 & 0 & \beta_{1:2,1} & 0 & 0 & \beta_{1:3,1} & 0 & \mathbf{B}^2 \\ 0 & \boxed{\beta_{1:1,2}} & 0 & 0 & \boxed{\beta_{1:2,2}} & 0 & 0 & \boxed{\beta_{1:3,2}} & 0 \\ 0 & 0 & \beta_{1:1,3} & 0 & 0 & \beta_{1:2,3} & 0 & 0 & \beta_{1:3,3} \\ \beta_{2:1,1} & 0 & 0 & \beta_{2:2,1} & 0 & 0 & \beta_{2:3,1} & 0 & 0 \\ 0 & \boxed{\beta_{2:1,2}} & 0 & 0 & \boxed{\beta_{2:2,2}} & 0 & 0 & \boxed{\beta_{2:3,2}} & 0 \\ 0 & 0 & \beta_{2:1,3} & 0 & 0 & \beta_{2:2,3} & 0 & 0 & \beta_{2:3,3} \\ \beta_{3:1,1} & 0 & 0 & \beta_{3:2,1} & 0 & 0 & \beta_{3:3,1} & 0 & 0 \\ 0 & \boxed{\beta_{3:1,2}} & 0 & 0 & \boxed{\beta_{3:2,2}} & 0 & 0 & \boxed{\beta_{3:3,2}} & 0 \\ 0 & 0 & \beta_{3:1,3} & 0 & 0 & \beta_{3:2,3} & 0 & 0 & \beta_{3:3,3} \end{bmatrix}$$

Theorem 2 (Efficient Matrix Multiplication). *For matrices $\mathbf{M}^A, \mathbf{M}^B$ defined in Theorem 1 $\forall x \in \mathbb{R}^{HW}$, the following equation holds:*

$$\mathbf{y} = \mathbf{M}^A \times \mathbf{M}^B \times x \Leftrightarrow \mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}, \mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}, \quad (9)$$

where $\mathbf{y} \in \mathbb{R}^{HW}$, $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$, $\mathbf{Y} = \text{unvec}(\mathbf{y}) \in \mathbb{R}^{H \times W}$, $\mathbf{Z} \in \mathbb{R}^{H \times W}$, and the operator $\text{vec}(\cdot)$ vectorizes a matrix by stacking its columns and $\text{unvec}(\cdot)$ is its inverse operator.

Algorithm 1: Efficient Masked Attention Computation.

Input: decay factors α, β of \mathbf{L} , vector $\mathbf{x} \in \mathbb{R}^{HW}$;

- 1: Compute $\mathbf{X} = \text{unvec}(\mathbf{x}) \in \mathbb{R}^{H \times W}$;
- 2: Compute $\mathbf{B}^l \in \mathbb{R}^{H \times H}$, where for $l=1:W$, $[\mathbf{B}^l]_{i,k} = \beta_{i:k,l}$;
- 3: Compute $\mathbf{Z} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}}$;
- 4: Compute $\mathbf{A}^i \in \mathbb{R}^{W \times W}$, where for $i=1:H$, $[\mathbf{A}^i]_{j,l} = \alpha_{i,j:l}$;
- 5: Compute $\mathbf{Y} \in \mathbb{R}^{H \times W}$, where $\boxed{\mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}}$;

Output: $\mathbf{y} = \text{vec}(\mathbf{Y})$;

Chunkwise algorithm

Complexity: $\mathcal{O}(H^2W) \rightarrow \mathcal{O}(HW)$

Complexity: $\mathcal{O}(HW^2) \rightarrow \mathcal{O}(HW)$

$$L = \begin{bmatrix} \alpha_{1,1:1} & \alpha_{1,1:2} & \alpha_{1,1:3} \\ \alpha_{1,2:1} & \alpha_{1,2:2} & \alpha_{1,2:3} \\ \alpha_{1,3:1} & \alpha_{1,3:2} & \alpha_{1,3:3} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{matrix} \mathbf{A}^1 \\ \mathbf{A}^2 \\ \mathbf{A}^3 \end{matrix} \times \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{matrix} \times \begin{bmatrix} \beta_{1:1,1} & 0 & 0 & \beta_{1:2,1} & 0 & 0 & \beta_{1:3,1} & 0 & \mathbf{B}^3 \\ 0 & \beta_{1:1,2} & 0 & 0 & \beta_{1:2,2} & 0 & 0 & \beta_{1:3,2} & 0 \\ 0 & 0 & \boxed{\beta_{1:1,3}} & 0 & 0 & \boxed{\beta_{1:2,3}} & 0 & 0 & \boxed{\beta_{1:3,3}} \\ \beta_{2:1,1} & 0 & 0 & \beta_{2:2,1} & 0 & 0 & \beta_{2:3,1} & 0 & 0 \\ 0 & \beta_{2:1,2} & 0 & 0 & \beta_{2:2,2} & 0 & 0 & \beta_{2:3,2} & 0 \\ 0 & 0 & \boxed{\beta_{2:1,3}} & 0 & 0 & \boxed{\beta_{2:2,3}} & 0 & 0 & \boxed{\beta_{2:3,3}} \\ \beta_{3:1,1} & 0 & 0 & \beta_{3:2,1} & 0 & 0 & \beta_{3:3,1} & 0 & 0 \\ 0 & \beta_{3:1,2} & 0 & 0 & \beta_{3:2,2} & 0 & 0 & \beta_{3:3,2} & 0 \\ 0 & 0 & \boxed{\beta_{3:1,3}} & 0 & 0 & \boxed{\beta_{3:2,3}} & 0 & 0 & \boxed{\beta_{3:3,3}} \end{bmatrix}$$

5.2 PPMA: Method (Efficient Computation Theory)

100

- Efficient Computation of Polyline Path Mask Matrix Multiplication

- Naive Computation:** for a rank- N matrix $L \in \mathbb{R}^{N \times N}$ and a vector $x \in \mathbb{R}^N$, Lx requires a complexity of $\mathcal{O}(N^2)$
- Efficient Computation:** for decomposed polyline path mask, $Lx = L^H \times L^V \times x$ requires a complexity of $\mathcal{O}(N)$

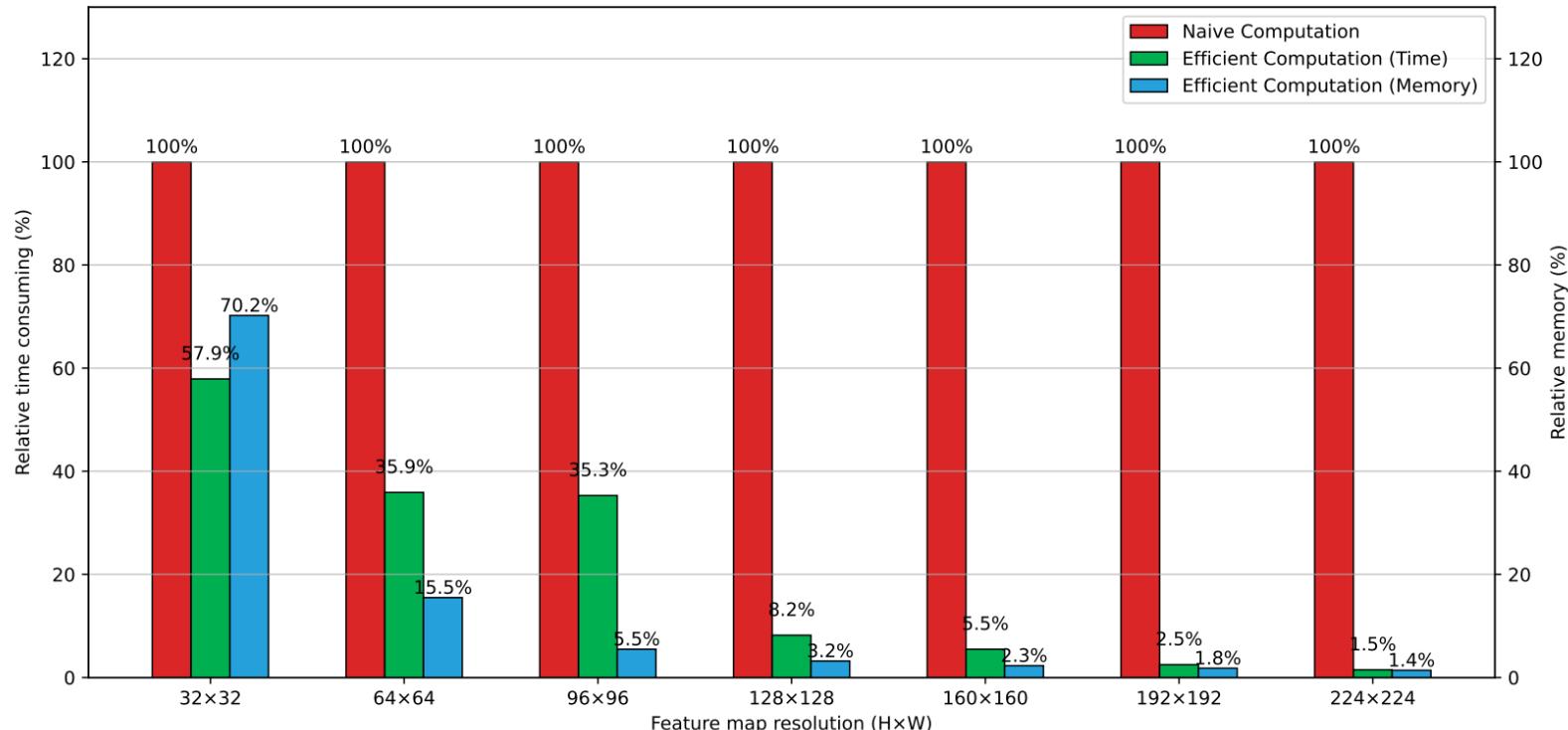


Figure 10: The comparison of the relative time consuming and memory usage between the naive computation and efficient computation (Algorithm 1) of Lx .

Remarks. Intuitively, as illustrated in Fig. 4, Algorithm 1 shows that the 2D polyline path scanning on 2D tokens (i.e., Lx) can be decomposed as the 1D vertical scanning along each column of X (i.e., $Z_{:,l} = B^l \times X_{:,l}$) followed by the 1D horizontal scanning along each row of Z (i.e., $Y_{i,:} = A^i \times Z_{i,:}$). This equivalence offers an intuitive understanding of the physical meaning of the decomposed polyline path mask $L = L^H L^V$ and enables its natural extension to 3D or higher-dimensional tokens, as detailed in Appendix C.2.

Algorithm 1: Efficient Masked Attention Computation.

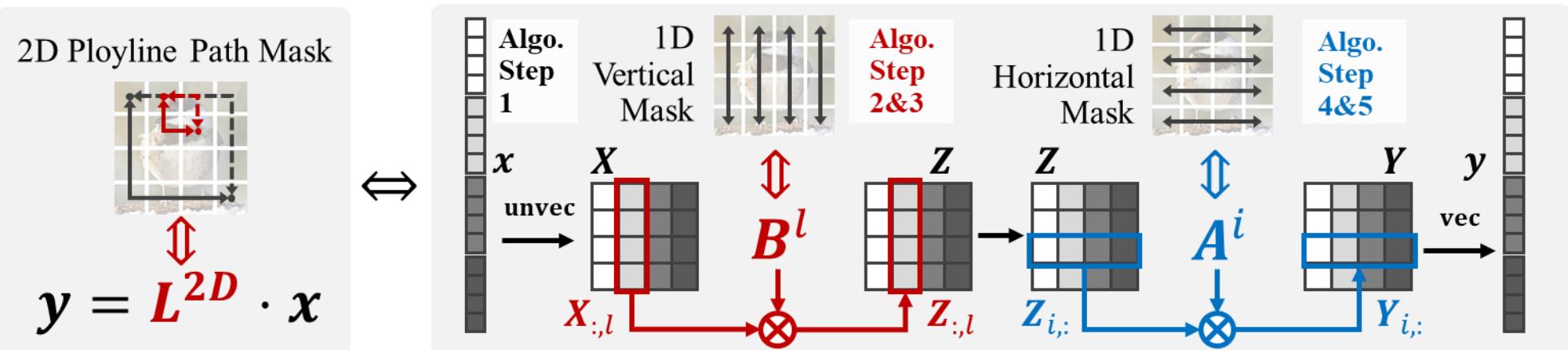
Input: decay factors α, β of L , vector $x \in \mathbb{R}^{HW}$;

- 1: Compute $X = \text{unvec}(x) \in \mathbb{R}^{H \times W}$;
- 2: Compute $B^l \in \mathbb{R}^{H \times H}$, where for $l=1:W$, $[B^l]_{i,k} = \beta_{i:k,l}$;
- 3: Compute $Z \in \mathbb{R}^{H \times W}$, where $Z_{:,l} = B^l \times X_{:,l}$;
- 4: Compute $A^i \in \mathbb{R}^{W \times W}$, where for $i=1:H$, $[A^i]_{j,l} = \alpha_{i,j:l}$;
- 5: Compute $Y \in \mathbb{R}^{H \times W}$, where $Y_{i,:} = A^i \times Z_{i,:}$;

Output: $y = \text{vec}(Y)$;

Complexity: $\mathcal{O}(H^2W) \rightarrow \mathcal{O}(HW)$

Complexity: $\mathcal{O}(HW^2) \rightarrow \mathcal{O}(HW)$



Remarks. Intuitively, as illustrated in Fig. 4, Algorithm 1 shows that the 2D polyline path scanning on 2D tokens (i.e., Lx) can be decomposed as the 1D vertical scanning along each column of \mathbf{X} (i.e., $\mathbf{Z}_{:,l} = \mathbf{B}^l \times \mathbf{X}_{:,l}$) followed by the 1D horizontal scanning along each row of \mathbf{Z} (i.e., $\mathbf{Y}_{i,:} = \mathbf{A}^i \times \mathbf{Z}_{i,:}$). This equivalence offers an intuitive understanding of the physical meaning of the decomposed polyline path mask $\mathbf{L} = \mathbf{L}^H \mathbf{L}^V$ and enables its natural extension to 3D or higher-dimensional tokens, as detailed in Appendix C.2.

C.2 3D Extension of Polyline Path Mask

Based on the decomposability, we naturally extend the 2D polyline path mask to 3D applications. As illustrated in Fig. 14, the 3D polyline path mask \mathbf{L}^{3D} can be decomposed as the multiplication of three 1D structured masks, $\mathbf{L}^H \times \mathbf{L}^V \times \mathbf{L}^D$, representing the horizontal, vertical, and depth scanning masks, respectively. Specifically, for each token pair $(\mathbf{x}_{i,j,k}, \mathbf{x}_{l,m,n})$ in the 3D grid, the 3D polyline path mask is defined as:

$$\mathcal{L}_{(i,j,k),(l,m,n)}^{3D} = \alpha_{i,j,k:n} \beta_{i,j:m,n} \gamma_{i:l,m,n}, \quad (40)$$

where \mathcal{L}^{3D} is the tensor form of matrix \mathbf{L}^{3D} , α , β , and γ are the decay factors along the horizontal, vertical, and depth axes, respectively. Compared to the cross-scanning strategy [30], the 3D polyline path scanning strategy better preserves the adjacency relationships of 3D tokens.

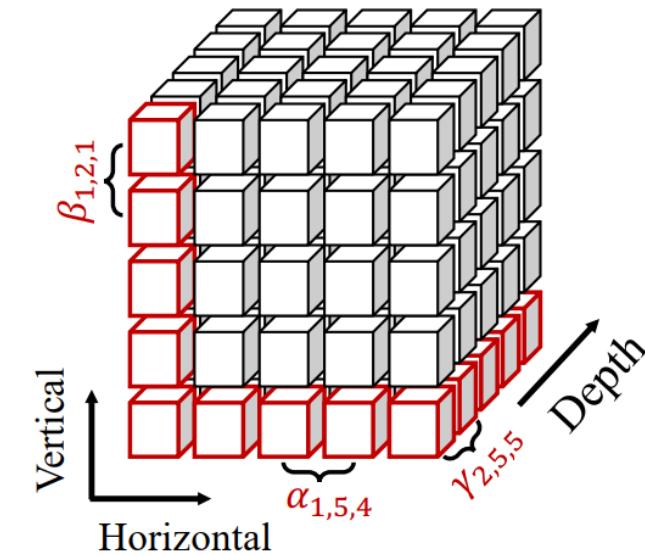


Figure 14: Illustration of the 3D extension of polyline path mask.

5.2 PPMA: Method (Polyline Path Masked Attention)

103

- Polyline path mask can be integrated into various attention in a plug-and-play manner

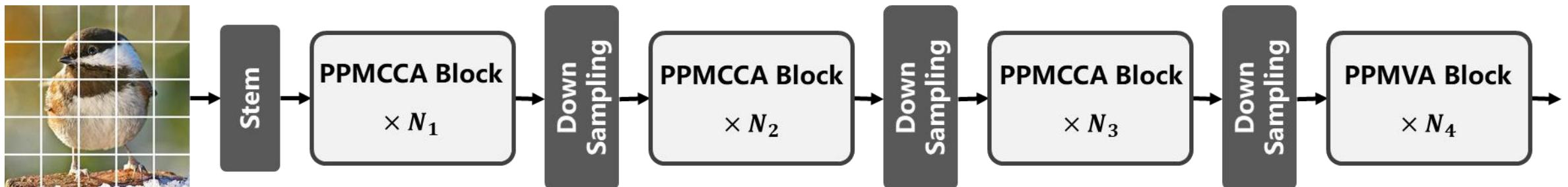
Basic Paradigm: $\text{PPMA}(X) = (\text{Attn}(\mathbf{Q}, \mathbf{K}) \odot \mathbf{L}^{2D})\mathbf{V} = (\text{Attn}(\mathbf{Q}, \mathbf{K}) \odot \mathbf{L})\mathbf{V} + (\text{Attn}(\mathbf{Q}, \mathbf{K}) \odot \tilde{\mathbf{L}})\mathbf{V}$

1) Masked **Vanilla self-attention**: $\text{PPMVA}(X) = (\text{softmax}(\mathbf{Q}\mathbf{K}^T) \odot \mathbf{L}^{2D})\mathbf{V}$ Complexity: $\mathcal{O}(N^2)$

2) Masked **linear attention**: $\text{PPMLA}(X) = (\mathbf{Q}\mathbf{K}^T \odot \mathbf{L}^{2D})\mathbf{V} = \mathbf{Q} \star (\mathbf{L}^{2D} \times (\mathbf{K} \star \mathbf{V}))$ Complexity: $\mathcal{O}(N)$

3) Masked **criss-cross attention**: $\text{PPMCCA}(X) = ((\mathbf{S}^H \times \mathbf{S}^V) \odot \mathbf{L}^{2D})\mathbf{V}$ Complexity: $\mathcal{O}(N^{\frac{3}{2}})$

4) Masked **decomposed attention** : $\text{PPMDA}(X) = ((\mathbf{S}_1 \times \mathbf{S}_2) \odot \mathbf{L}^{2D})\mathbf{V} = \mathbf{S}_1 \star (\mathbf{L}^{2D} \times (\mathbf{S}_2 \star \mathbf{V}))$



5.2 PPMA: Method (Polyline Path Masked Attention)

104

- Decomposed Criss-Cross Attention^[1]

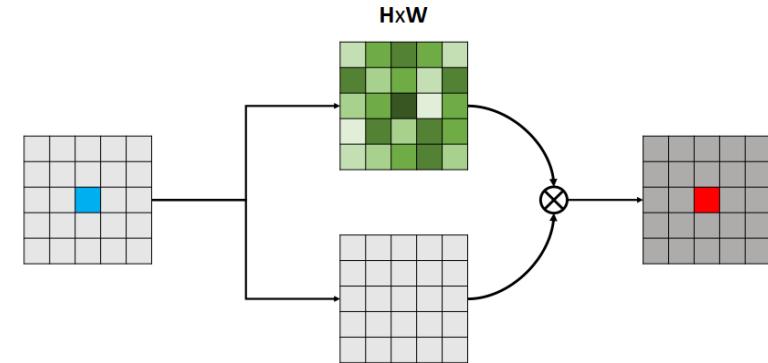
$$\mathbf{S}^H = \text{softmax}(\mathbf{Q}_H \mathbf{K}_H^T) \odot \mathbf{L}_H$$

$$\mathbf{S}^V = \text{softmax}(\mathbf{Q}_V \mathbf{K}_V^T) \odot \mathbf{L}_V$$

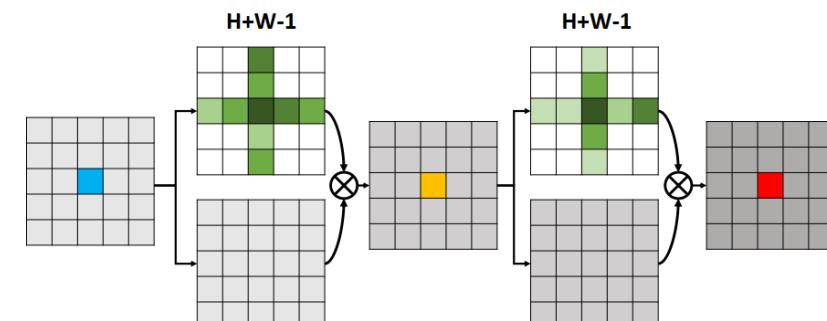
$$\mathbf{Y} = 0.5 \times (\mathbf{S}^V (\mathbf{S}^H \mathbf{V})^\top)^\top + 0.5 \times (\mathbf{S}^H (\mathbf{S}^V \mathbf{V}^\top)^\top)$$

3) Masked criss-cross attention: $\text{PPMCCA}(X) = ((\mathbf{S}^H \times \mathbf{S}^V) \odot \mathbf{L}^{2D}) \mathbf{V}$

$$\begin{aligned} ((\mathbf{S}^H \times \mathbf{S}^V) \odot \mathbf{L}) \mathbf{V} &= ((\mathbf{S}^H \times \mathbf{S}^V) \odot (\mathbf{L}^H \times \mathbf{L}^V)) \mathbf{V} \\ &= ((\widehat{\mathbf{S}}^H \odot \widehat{\mathbf{S}}^V) \odot (\widehat{\mathbf{L}}^H \odot \widehat{\mathbf{L}}^V)) \mathbf{V} \\ &= ((\widehat{\mathbf{S}}^H \odot \widehat{\mathbf{L}}^H) \odot (\widehat{\mathbf{S}}^V \odot \widehat{\mathbf{L}}^V)) \mathbf{V} \\ &= ((\mathbf{S}^H \odot \mathbf{L}^H) \times (\mathbf{S}^V \odot \mathbf{L}^V)) \mathbf{V} \\ &= (\mathbf{S}^H \odot \mathbf{L}^H) \times ((\mathbf{S}^V \odot \mathbf{L}^V) \times \mathbf{V}) \end{aligned}$$



(a) Vanilla Attention Block



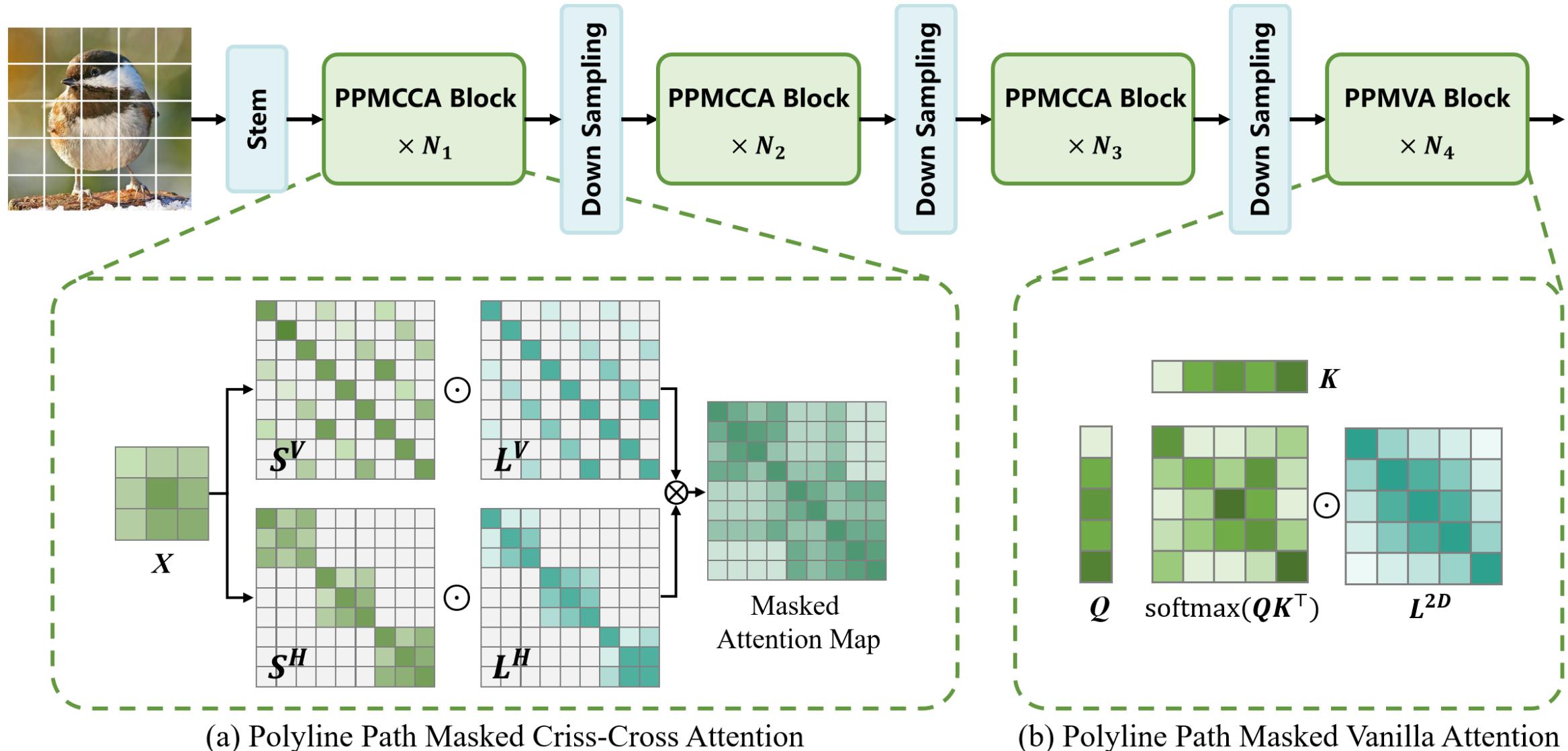
(b) Criss-Cross Attention block



[1] Huang Z, Wang X, Huang L, et al. Ccnet: Criss-cross attention for semantic segmentation[C], 2019ICCV & 2020TPAMI. (cited:3646)

5.2 PPMA: Method (Polyline Path Masked Attention)

105



5.3 PPMA: Experiments (Image Classification)

Table 1: Image classification performance on the ImageNet-1K validation set.

Model	Arch.	#Param. (M)	FLOPs (G)	Top-1 (%)	Model	Arch.	#Param. (M)	FLOPs (G)	Top-1 (%)
RegNetY-1.6G [34]	SSM	11	1.6	78.0	NAT-T [16]	Trans.	28	4.3	83.2
EffNet-B3 [41]		12	1.8	81.6	BiFormer-S [56]		26	4.5	83.8
Vim-T [57]		7	1.5	76.1	RMT-S [10]		27	4.5	84.0
MSVMamba-M [36]		12	1.5	79.8	PPMA-S		27	4.9	84.2
BiFormer-T [56]		13	2.2	81.4	RegNetY-8G [34]	CNN	39	8.0	81.7
NAT-M [16]		20	2.7	81.8	ConvNeXt-S [32]		50	8.7	83.1
SMT-T [29]		12	2.4	82.2	EffNet-B5 [41]		30	9.9	83.6
RMT-T [10]		14	2.5	82.4	VMamba-S [30]		50	8.7	83.6
PPMA-T		14	2.7	82.6	2DMamba-S [52]		50	8.8	83.8
RegNetY-4G [34]	SSM	21	4.0	80.0	GrootVL-S [48]		51	8.5	84.2
ConvNeXt-T [32]		29	4.5	82.1	MLLA-S [15]		43	7.3	84.4
EffNet-B4 [41]		19	4.2	82.9	Spatial-Mamba-S [46]		43	7.1	84.6
VMamba-T [30]		30	4.9	82.6	Swin-S [31]		50	8.7	83.0
2DMamba-T [52]		31	4.9	82.8	NAT-S [16]		51	7.8	83.7
GrootVL-T [48]		30	4.8	83.4	CSWin-B [8]		78	15.0	84.2
Spatial-Mamba-T [46]		27	4.5	83.5	MambaVision-B [17]		98	15.0	84.2
MLLA-T [15]		25	4.2	83.5	BiFormer-B [56]		57	9.8	84.3
Swin-T [31]		29	4.5	82.1	iFormer-B [37]		48	9.4	84.6
CSWin-T [8]		23	4.3	82.7	RMT-B [10]		54	9.7	84.9
MambaVision-T2 [17]	Trans.	35	5.1	82.7	PPMA-B		54	10.6	85.0

5.3 PPMA: Experiments (Object Detection)

Table 2: Object detection and instance segmentation performance with Mask R-CNN [18] detector on COCO val2017. FLOPs are calculated with input resolution of 1280×800 .

Backbone	#Param. (M)	FLOPs (G)	AP ^b	AP ^b ₅₀	AP ^b ₇₅	AP ^m	AP ^m ₅₀	AP ^m ₇₅
Vim-T [57]	–	–	45.7	63.9	49.6	39.2	60.9	41.7
MSVMamba-M [36]	32	201	43.8	65.8	47.7	39.9	62.9	42.9
	30	231	44.2	66.7	48.4	40.4	63.4	43.4
	33	218	46.7	68.6	51.6	42.1	65.3	45.2
	33	218	47.1	68.7	51.7	42.4	65.9	45.7
	33	218	47.1	68.7	51.7	42.4	65.9	45.7
ResNet-50 [19]	44	260	38.2	58.8	41.4	34.7	55.7	37.2
ConvNeXt-T [32]	48	262	44.2	66.6	48.3	40.1	63.3	42.8
MLLA-T [15]	44	255	46.8	69.5	51.5	42.1	66.4	45.0
GrootVL-T [48]	49	265	47.0	69.4	51.5	42.7	66.4	46.0
VMamba-T [30]	50	271	47.3	69.3	52.0	42.7	66.4	45.9
Spatial-Mamba-T [46]	46	216	47.6	69.6	52.3	42.9	66.5	46.2
Swin-T [31]	48	267	43.7	66.6	47.7	39.8	63.3	42.7
CSWin-T [8]	42	279	46.7	68.6	51.3	42.2	65.6	45.4
BiFormer-S [56]	–	–	47.8	69.8	52.3	43.2	66.8	46.5
RMT-S [10]	46	262	48.8	70.8	53.4	43.6	67.4	47.3
PPMA-S	46	263	49.2	70.7	54.0	43.8	67.4	47.1
ResNet-101 [19]	63	336	40.4	61.1	44.2	36.4	57.7	38.8
ConvNeXt-S [32]	70	348	45.4	67.9	50.0	41.8	65.2	45.1
GrootVL-S [48]	70	341	48.6	70.3	53.5	43.6	67.5	47.1
VMamba-S [30]	70	349	48.7	70.0	53.4	43.7	67.3	47.0
Spatial-Mamba-S [46]	63	315	49.2	70.8	54.2	44.0	67.9	47.5
MLLA-S [15]	63	319	49.2	71.5	53.9	44.2	68.5	47.2
Swin-S [31]	69	359	45.7	67.9	50.4	41.1	64.9	44.2
CSWin-S [8]	54	342	47.9	70.1	52.6	43.2	67.1	46.2
BiFormer-B [56]	–	–	48.6	70.5	53.8	43.7	67.6	47.1
RMT-B [10]	73	373	50.7	72.0	55.7	45.1	69.2	49.0
PPMA-B	73	374	51.1	72.5	55.9	45.5	69.7	49.1

5.3 PPMA: Experiments (Semantic Segmentation)

Table 3: Semantic segmentation performance with UPerNet [47] segmentor on ADE20K val set. ‘SS’ and ‘MS’ represent single-scale and multi-scale testing, respectively.

Backbone	#Param. (M)	FLOPs (G)	mIoU(%)		Backbone	#Param. (M)	FLOPs (G)	mIoU(%)	
			SS	MS				SS	MS
LocalVim-T [21]	36	181	43.4	44.4	BiFormer-S [56]	–	–	49.8	50.8
MSVMamba-M [36]	42	875	45.1	45.4	RMT-S [10]	56	937	49.8	49.7
NAT-M [16]	50	900	45.1	46.4	PPMA-S	56	984	51.1	52.0
RMT-T [10]	43	977	48.0	48.8	ResNet-101 [19]	85	1030	42.9	44.0
PPMA-T	43	983	48.7	49.1	ConvNeXt-S [32]	82	1027	48.7	49.6
ResNet-50 [19]	67	953	42.1	42.8	VMamba-S [30]	82	1028	50.6	51.2
ConvNeXt-T [32]	60	939	46.0	46.7	Spatial-Mamba-S [46]	73	992	50.6	51.4
VMamba-T [30]	62	949	48.0	48.8	GrootVL-S [48]	82	1019	50.7	51.7
2DMamba-T [52]	62	950	48.6	49.3	Swin-S [31]	81	1039	47.6	49.5
GrootVL-T [48]	60	941	48.5	49.4	NAT-S [16]	82	1010	48.0	49.5
Spatial-Mamba-S [46]	57	936	48.6	49.4	MambaVision-S [17]	84	1135	48.2	–
Swin-T [31]	60	945	44.4	45.8	CSWin-S [8]	65	1027	50.4	51.5
MambaVision-T [17]	55	945	46.6	–	BiFormer-B [56]	–	–	51.0	51.7
NAT-T [16]	58	934	47.1	48.4	RMT-B [10]	83	1051	52.0	52.1
CSWin-S [8]	60	959	49.3	50.7	PPMA-B	83	1137	52.3	53.0

5.3 PPMA: Experiments (Ablation Study)

109

- Ablation study on the polyline path mask design

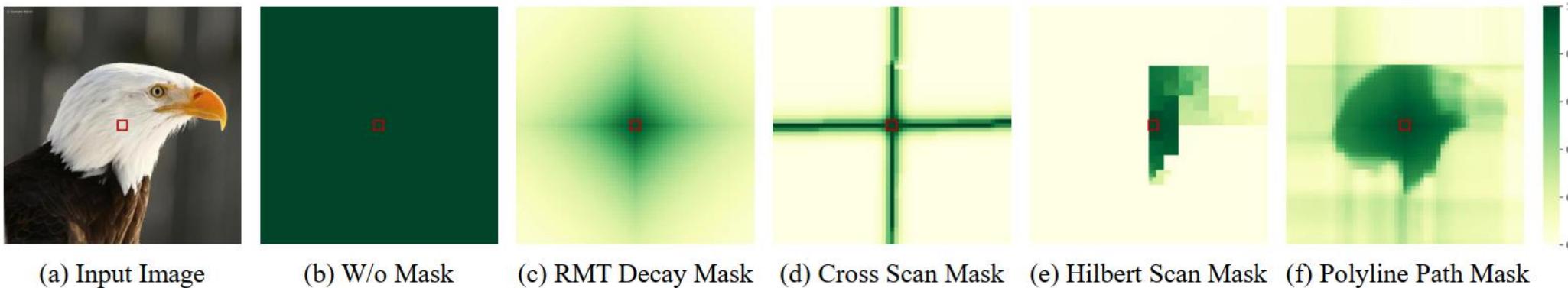


Figure 6: Illustration of various structured masks.

Table 9: Ablation study of structured mask designs in PPMA-T on ImageNet-1K and ADE20K.

Structured Mask	#Param. (M)	FLOPs (G)	Top-1 (%)	mIoU SS (%)
Baseline (w/o mask)	14.33	2.65	82.28	47.78
+ RMT Decay Mask	14.33	2.65	82.35	48.01
+ Cross Scan Mask	14.34	2.71	82.44	48.14
+ V2H Polyline Path Mask	14.34	2.71	82.44	48.57
+ 2D Polyline Path Mask	14.34	2.71	82.60	48.73
Shared Decay factors ($\alpha_{i,j} = \beta_{i,j}$)	14.33	2.71	82.37	48.27
Different Decay factors ($\alpha_{i,j} \neq \beta_{i,j}$)	14.34	2.71	82.60	48.73

5.3 PPMA: Experiments (Visualization)

110

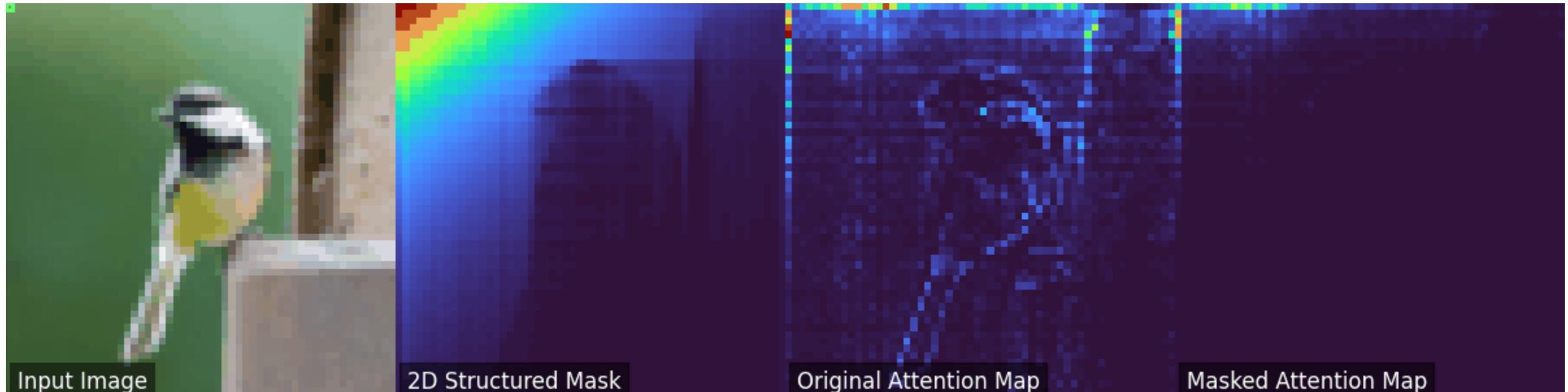
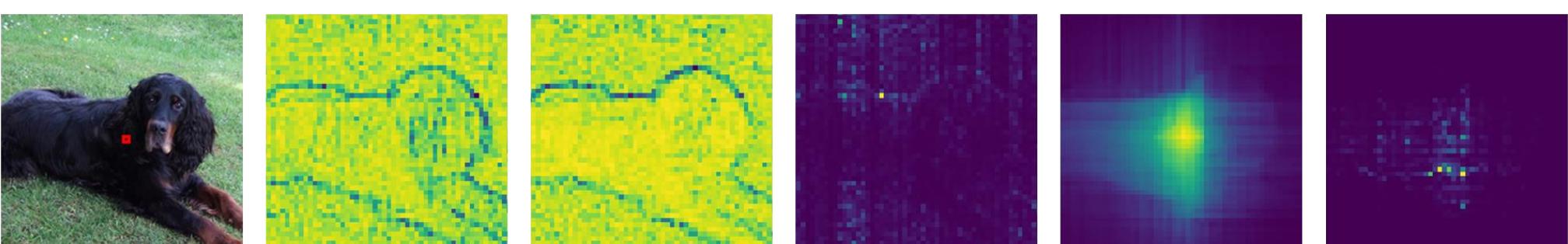
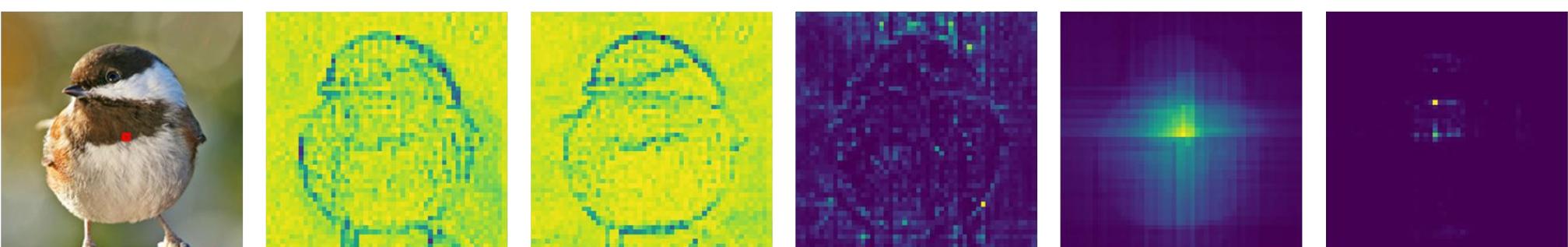
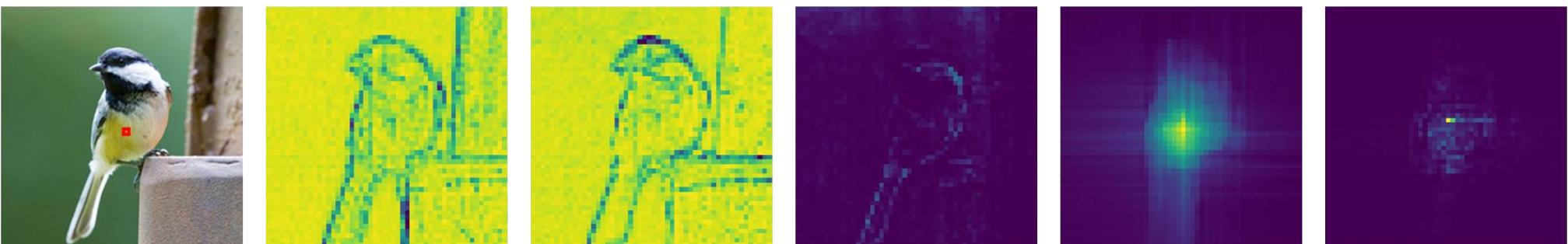
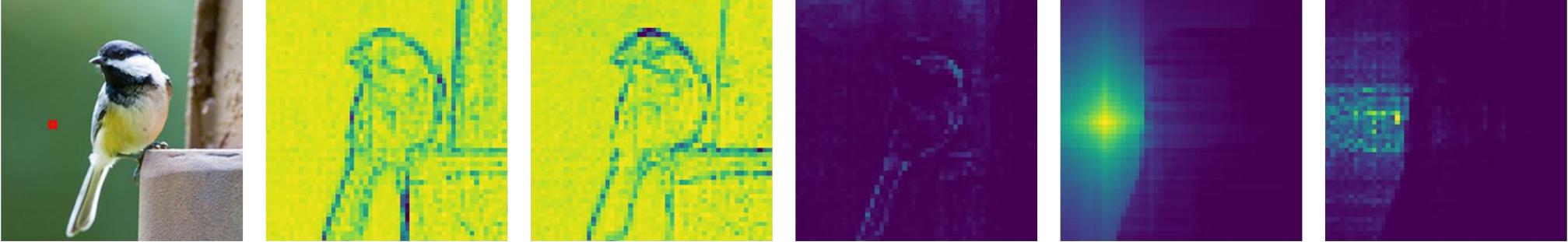


Figure 8: Visualizations of the Polyline Path Masked Attention





(a) Input Image

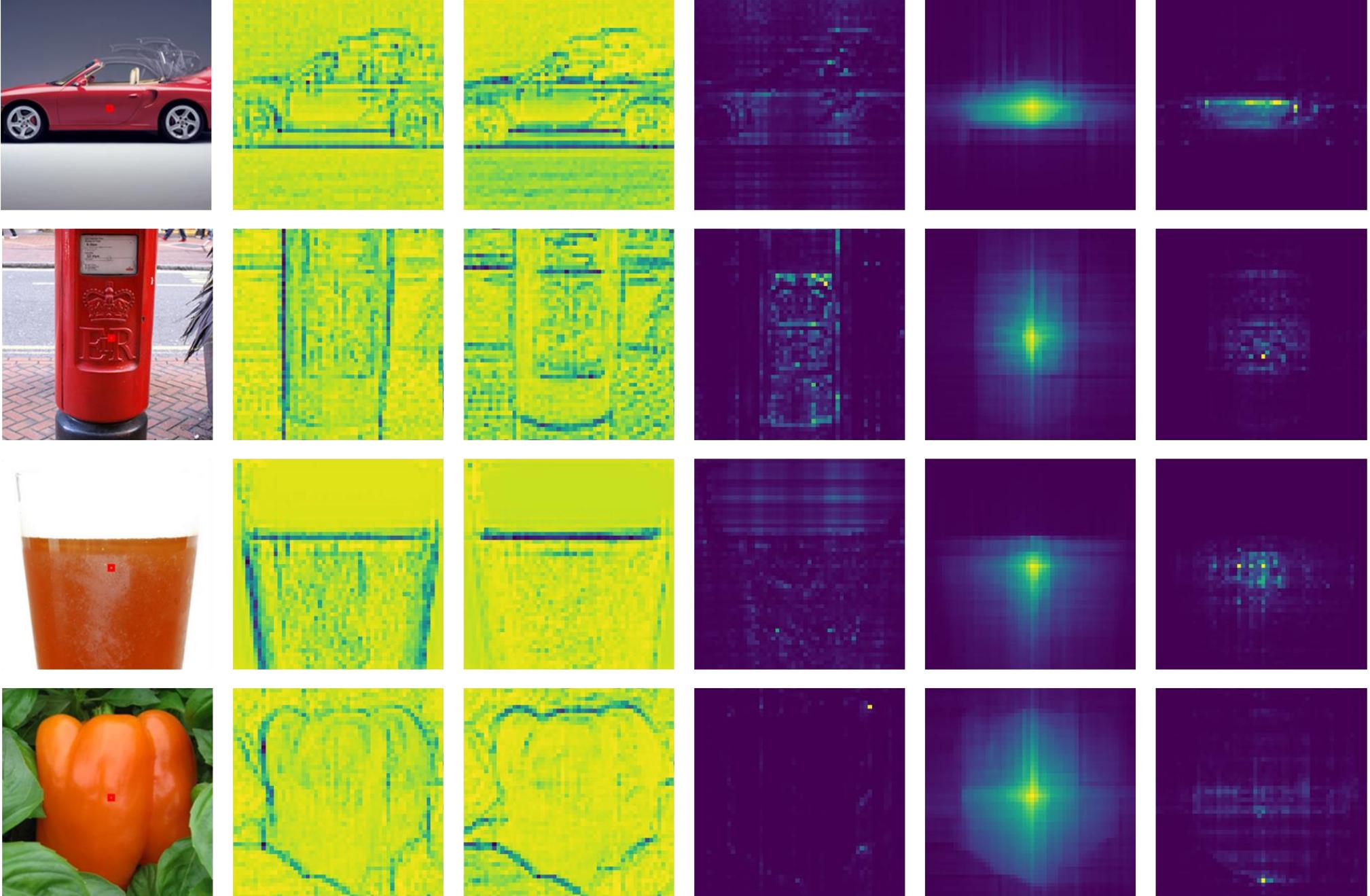
(b) Horizontal Decay Factor α

(c) Vertical Decay Factor β

(d) Original Attention Map $\text{softmax}(\mathbf{Q}\mathbf{K}^\top)$

(e) Polyline Path Mask \mathbf{L}^{2D}

(f) Polyline Path Masked Attention Map



(a) Input Image

(b) Horizontal
Decay Factor α

(c) Vertical
Decay Factor β

(d) Original Attention
Map $\text{softmax}(QK^\top)$

(e) Polyline Path
Mask L^{2D}

(f) Polyline Path
Masked Attention Map

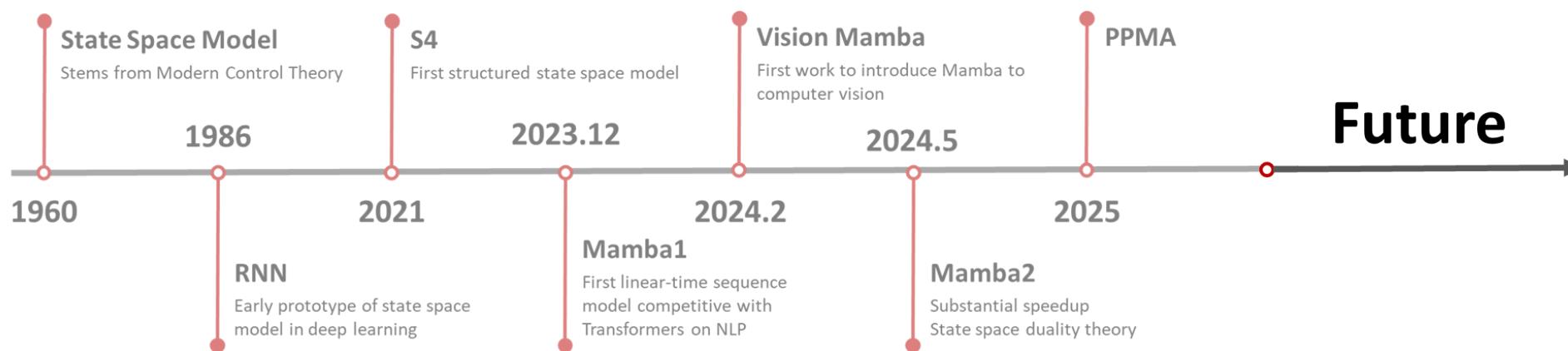
What is next for Mamba?

113

- Stable and scalable linear-time foundational model remains a worthwhile subject
- Hybrid architecture maybe the future

9.2.3 Hybrid Models: Combining SSD Layer with MLP and Attention

Recent and concurrent work (Dao, D. Y. Fu, et al. 2023; De et al. 2024; Glorioso et al. 2024; Lieber et al. 2024) suggests that a hybrid architecture with both SSM layers and attention layers could improve the model quality over that of a Transformer, or a pure SSM (e.g., Mamba) model, especially for in-context learning. We explore the different ways that SSD layers can be combined with attention and MLP to understand the benefits of each. Empirically we find that having around 10% of the total number of layers being attention performs best. Combining SSD layers, attention layers, and MLP also works better than either pure Transformer++ or Mamba-2.





西安交通大学
XI'AN JIAOTONG UNIVERSITY



Thanks!

