

# A View of Mamba

杨泽康

中科院计算技术研究所

2024 年 4 月 19 日

## ① Mamba

## ② Vision Mamba

## ③ MIL Mamba

## ④ Multimodal Mamba

## ⑤ Audio Mamba

## ⑥ Diffusion Mamba

## ⑦ 3D Mamba

## ⑧ XAI Mamba

## ① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

⑦ 3D Mamba

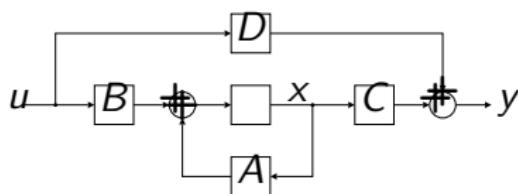
⑧ XAI Mamba

# State Space Model

$\mu(t)$  为输入变量,  $x(t)$  为隐变量,  $y(t)$  为输出变量,  $\dot{x}(t)$  为隐变量的梯度

$$\dot{x}(t) = A(t)x(t) + B(t)\mu(t)$$

$$y(t) = C(t)x(t) + D(t)\mu(t)$$

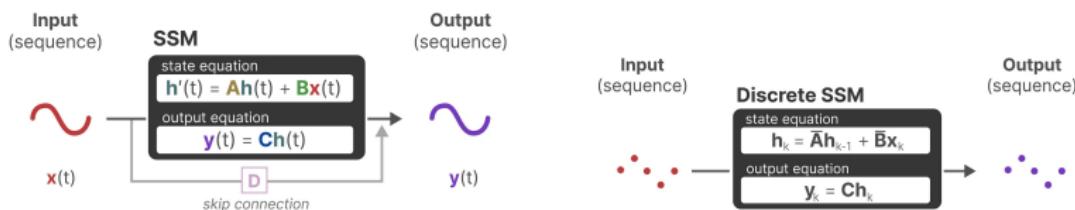


System type	State-space model
Continuous time-invariant	$\dot{x}(t) = Ax(t) + Bu(t)$ $y(t) = Cx(t) + Du(t)$
Continuous time-variant	$\dot{x}(t) = A(t)x(t) + B(t)u(t)$ $y(t) = C(t)x(t) + D(t)u(t)$
Explicit discrete time-invariant	$x(k+1) = Ax(k) + Bu(k)$ $y(k) = Cx(k) + Du(k)$
Explicit discrete time-variant	$x(k+1) = A(k)x(k) + B(k)u(k)$ $y(k) = C(k)x(k) + D(k)u(k)$
Laplace domain of continuous time-invariant	$sX(s) - x(0) = AX(s) + BU(s)$ $Y(s) = CX(s) + DU(s)$
Z-domain of discrete time-invariant	$zX(z) - zx(0) = AX(z) + BU(z)$ $Y(z) = CX(z) + DU(z)$

S4

## Efficiently Modeling Long Sequences with Structured State Spaces

想通过 SSM 来处理长序列问题，采用的是 SSM 的一种连续时间不变形式



- 离散化: 表达式中存在导数, 没有办法直接处理离散的序列

$$\bar{A} = (\mathbf{I} - \Delta/2 \cdot A)^{-1}(\mathbf{I} + \Delta/2 \cdot A)$$

$$\bar{B} = (\mathbf{I} - \Delta/2 \cdot A)^{-1} \Delta B$$

$$\bar{C} = C$$

## S4

- 循环-卷积表示：等效成卷积形式，训练时采用卷积方式，推理时循环形式，加快训练和推理速度

$$\begin{aligned}x_0 &= \overline{B}u_0 & x_1 &= \overline{AB}u_0 + \overline{B}u_1 & x_2 &= \overline{A}^2\overline{B}u_0 + \overline{AB}u_1 + \overline{B}u_2 & \dots \\y_0 &= \overline{CB}u_0 & y_1 &= \overline{CAB}u_0 + \overline{CB}u_1 & y_2 &= \overline{CA}^2\overline{B}u_0 + \overline{CAB}u_1 + \overline{CB}u_2 & \dots\end{aligned}$$

$$\begin{aligned}y_k &= \overline{CA}^k\overline{B}u_0 + \overline{CA}^{k-1}\overline{B}u_1 + \dots + \overline{CAB}u_{k-1} + \overline{CB}u_k \\y &= \overline{\mathbf{K}} * u\end{aligned}$$

$$\overline{\mathbf{K}} \in \mathbb{R}^L = (\overline{CB}, \overline{CAB}, \dots, \overline{CA}^{L-1}\overline{B})$$

卷积核大小固定为序列长度大小：傅里叶变换后进行加法  
序列为因果关系：掩码

## S4

- 基于 HiPPO 处理长序列：解决长距离依赖问题

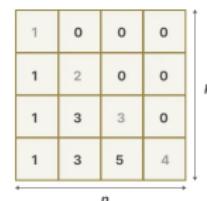
	x	<b>Wxh</b>	Hidden	<b>Whh</b>	<b>Why</b>	y
SNN	1	<b>1xd'</b>	d'	<b>d'xd'</b>	<b>1x1</b>	1
	d	<b>dxd'</b>			<b>1xd</b>	d'
S4	1	<b>1xd'</b>	d'	<b>1xd'</b>	1	
	d		dxd'		d	

相比 RNN 更显式的建模隐藏层状态对过去信息的保存，此时  $A$  表示了对过去信息的记忆方式，采用 HiPPO 方式初始化  $A$  能更好的记忆长距离信息

HiPPO Matrix

HiPPO Matrix  $A_{nk}$

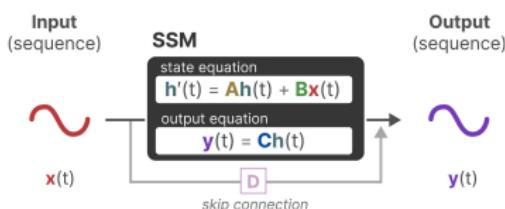
$$\left\{ \begin{array}{l} (2n+1)^{1/2} (2k+1)^{1/2} \text{ everything below the diagonal} \\ n+1 \text{ the diagonal} \\ 0 \text{ everything above the diagonal} \end{array} \right.$$



# Mamba

## 相比 S4

- 将  $\Delta, B, C$  变成输入依赖的，通过将输入通过全连接层预测得到， $\bar{A}$  间接依赖输入（能对输入进行不同程度的记忆，之前  $A, B, C$  共享对不同输入的关注程度相同）
- 每个维度用一个 SSM, A, B, C 不再在不同维度共享



$$\begin{aligned}\bar{A} &= (\mathbf{I} - \Delta/2 \cdot \mathbf{A})^{-1}(\mathbf{I} + \Delta/2 \cdot \mathbf{A}) \\ \bar{B} &= (\mathbf{I} - \Delta/2 \cdot \mathbf{A})^{-1}\Delta\mathbf{B} \\ \bar{C} &= \mathbf{C}\end{aligned}$$

### Algorithm 1 SSM (S4)

```

Input:  $x : (B, L, D)$ 
Output:  $y : (B, L, D)$ 
1:  $A : (D, N) \leftarrow \text{Parameter}$ 
     $\triangleright$  Represents structured  $N \times N$  matrix
2:  $B : (D, N) \leftarrow \text{Parameter}$ 
3:  $C : (D, N) \leftarrow \text{Parameter}$ 
4:  $\Delta : (D) \leftarrow \tau_\Delta(\text{Parameter})$ 
5:  $\bar{A}, \bar{B} : (D, N) \leftarrow \text{discretize}(\Delta, A, B)$ 
6:  $y \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x)$ 
     $\triangleright$  Time-invariant: recurrence or convolution
7: return  $y$ 
```

### Algorithm 2 SSM + Selection (S6)

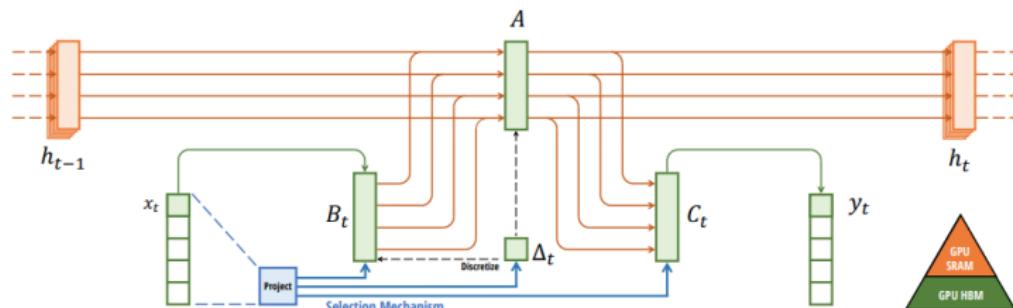
```

Input:  $x : (B, L, D)$ 
Output:  $y : (B, L, D)$ 
1:  $A : (D, N) \leftarrow \text{Parameter}$ 
     $\triangleright$  Represents structured  $N \times N$  matrix
2:  $B : (B, L, N) \leftarrow s_B(x)$ 
3:  $C : (B, L, N) \leftarrow s_C(x)$ 
4:  $\Delta : (B, L, D) \leftarrow \tau_\Delta(\text{Parameter} + s_\Delta(x))$ 
5:  $\bar{A}, \bar{B} : (B, L, D, N) \leftarrow \text{discretize}(\Delta, A, B)$ 
6:  $y \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x)$ 
     $\triangleright$  Time-varying: recurrence (scan) only
7: return  $y$ 
```

Mamba

### 带来的问题

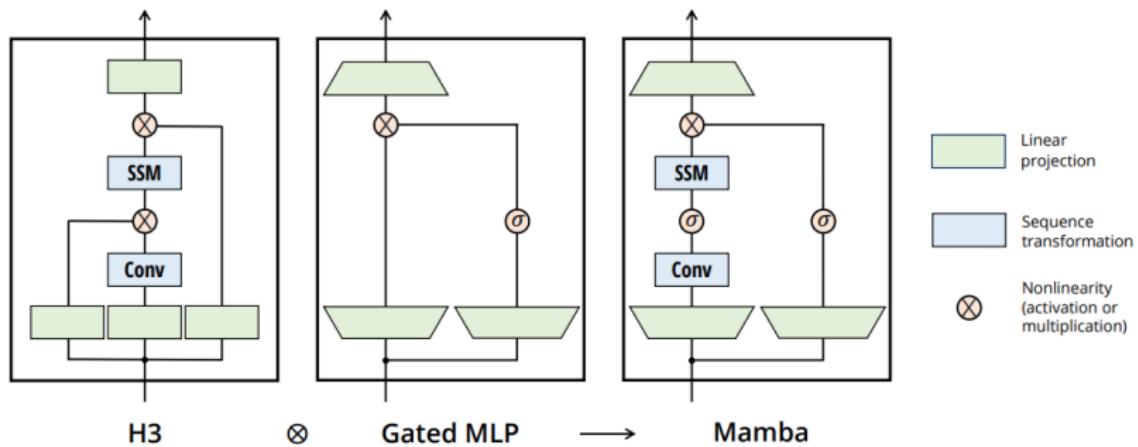
- $\bar{A}, \bar{B}, \bar{C}$  与输入相关，没有办法再提前计算卷积核并行训练减少反复读写显存这个瓶颈：并行扫描 + 核融合 + 重计算



**Figure 1: (Overview.)** Structured SSMs independently map each channel (e.g.  $D = 5$ ) of an input  $x$  to output  $y$  through a higher dimensional latent state  $h$  (e.g.  $N = 4$ ). Prior SSMs avoid materializing this large effective state ( $DN$ , times batch size  $B$  and sequence length  $L$ ) through clever alternate computation paths requiring time-invariance: the  $(\Delta, A, B, C)$  parameters are constant across time. Our selection mechanism adds back input-dependent dynamics, which also requires a careful hardware-aware algorithm to only materialize the expanded states in more efficient levels of the GPU memory hierarchy.

Mamba

与现代神经网络结合，得到 Mamba 架构



Mamba  
oooooooo

Vision Mamba  
●ooo

MIL Mamba  
ooo

Multimodal Mamba  
oooo

Audio Mamba  
oo

Diffusion Mamba  
○○

3D Mamba  
○○

XAI Mamba  
ooo

① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

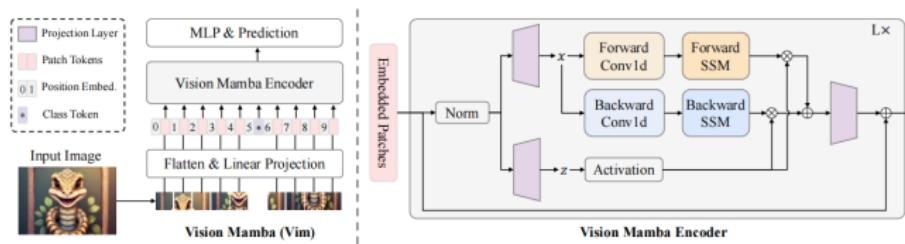
⑦ 3D Mamba

⑧ XAI Mamba

# Vision Mamba

## Vision Mamba: Efficient Visual Representation Learning with Bidirectional State Space Model

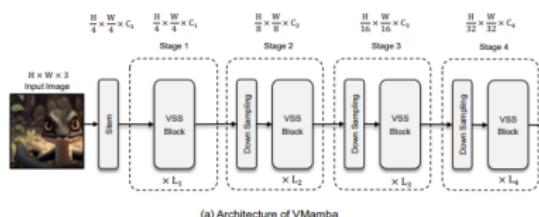
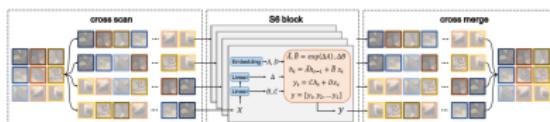
- Vision Transformer 结构
- 增加 backward 扫描分支，和 forward 扫描融合



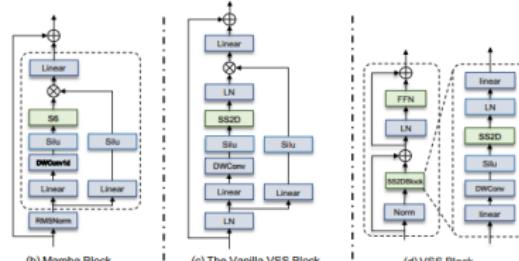
VMamba

VMamba: Visual State Space Model

- ConvNext 形式的 CNN 架构
  - 图像 patch 展开后四个方向进行扫描 (左-> 右, 右-> 左, 上-> 下, 下-> 上)
  - 对比了几种不同的扫描方式



### (a) Architecture of VMambo



(b) Mamba Block

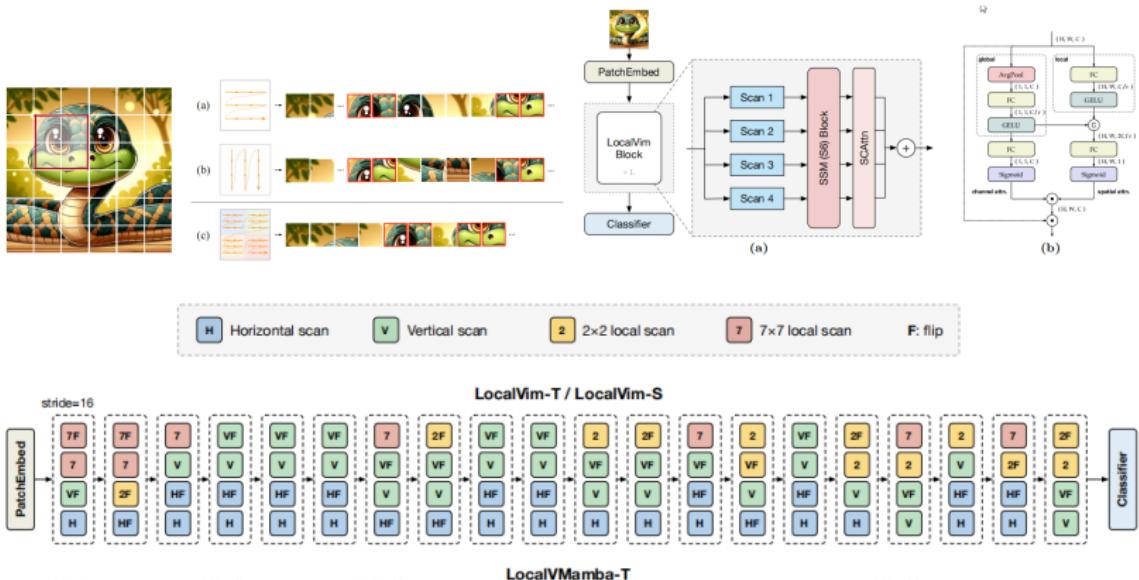
### (c) The Vanilla VSS Block

**(d) VSS Block**

## Local Mamba

LocalMamba: Visual State Space Model with Windowed Selective Scan

- 提出了四种扫描方式，对四种扫描方式进行融合
  - 进行神经架构搜索，搜索四种扫描方式的最佳组合



Mamba  
oooooooo

Vision Mamba  
oooo

MIL Mamba  
●oo

Multimodal Mamba  
oooo

Audio Mamba  
oo

Diffusion Mamba  
oo

3D Mamba  
oo

XAI Mamba  
ooo

① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

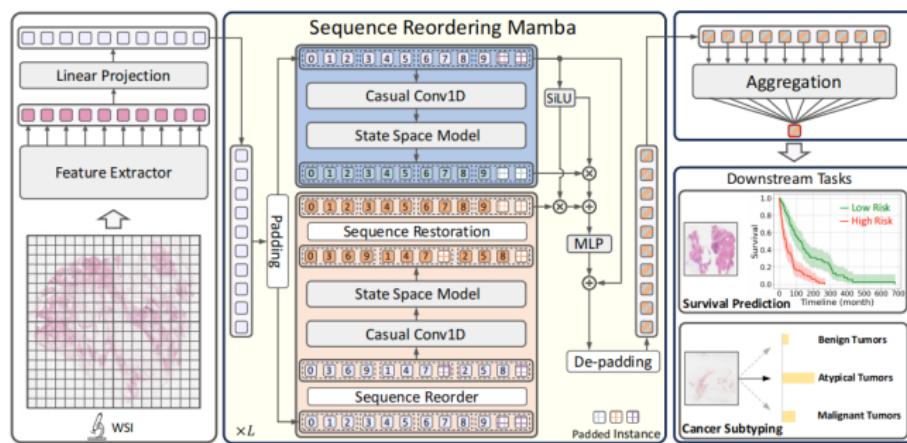
⑦ 3D Mamba

⑧ XAI Mamba

MambaMIL

# MambaMIL: Enhancing Long Sequence Modeling with Sequence Reordering in Computational Pathology

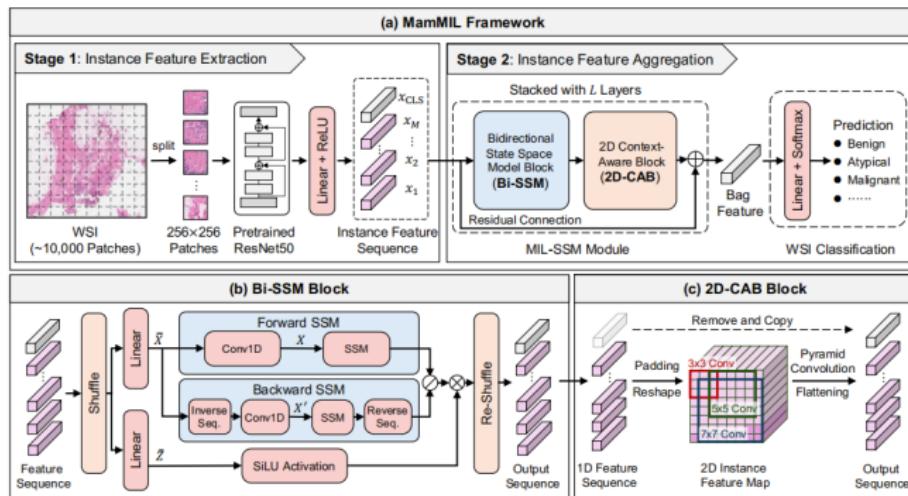
- 添加一个对序列重新排序后通过 Mamba 的分支
  - 重新排序为 (近似) 上-> 下扫描



# MamMIL

## MamMIL: Multiple Instance Learning for Whole Slide Images with State Space Models

- 结构与 Vision Mamba 一样
- 增加一个 shuffle 和 reshuffle 操作



① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

⑦ 3D Mamba

⑧ XAI Mamba

# Fusion Mamba

## FusionMamba: Efficient Image Fusion with State Space Model

- 提出 FSSM 模块，由另一个模态的数据去产生  $B, C, \Delta$

**Algorithm 1** SSM Block

```

Input:  $x : (\text{HW}, \text{C})$ 
Output:  $y : (\text{HW}, \text{C})$ 
1:  $A : (\text{C}, \text{N}) \leftarrow \text{Parameter}_A$ 
   /* A represents C sets of structured  $\text{N} \times \text{N}$  matrices [12] */
2:  $B : (\text{HW}, \text{N}) \leftarrow \text{Linear}_B(x)$ 
3:  $C : (\text{HW}, \text{N}) \leftarrow \text{Linear}_C(x)$ 
4:  $\Delta : (\text{HW}, \text{C}) \leftarrow \log(1 + \exp(\text{Linear}_\Delta(x) + \text{Parameter}_\Delta))$ 
   /* Parameter $_\Delta$  is a bias vector with a size of C */
5:  $\bar{A} : (\text{HW}, \text{C}, \text{N}) \leftarrow \exp(\Delta \otimes A)$ 
6:  $\bar{B} : (\text{HW}, \text{C}, \text{N}) \leftarrow \Delta \otimes B$ 
7:  $y \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x)$ 
   /* SSM represents Eq. 3 implemented by selective scan [11] */
return  $y$ 

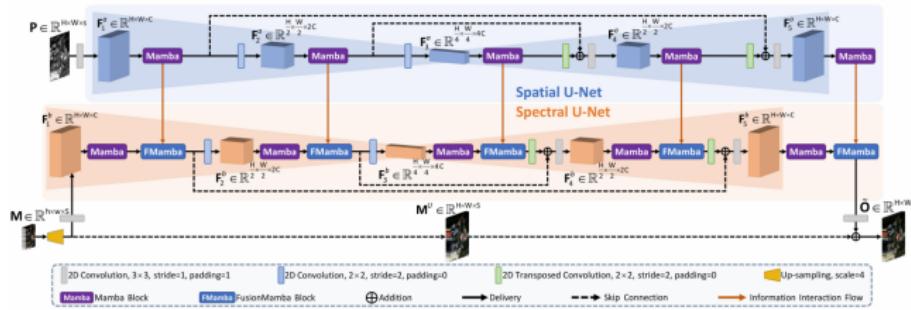
```

**Algorithm 2** FSSM Block

```

Input:  $x^a, x^b : (\text{HW}, \text{C})$ 
Output:  $y^d : (\text{HW}, \text{C})$ 
1:  $A : (\text{C}, \text{N}) \leftarrow \text{Parameter}_A$ 
   /* A represents C sets of structured  $\text{N} \times \text{N}$  matrices [12] */
2:  $B : (\text{HW}, \text{N}) \leftarrow \text{Linear}_B(x^b)$ 
3:  $C : (\text{HW}, \text{N}) \leftarrow \text{Linear}_C(x^b)$ 
4:  $\Delta : (\text{HW}, \text{C}) \leftarrow \log(1 + \exp(\text{Linear}_\Delta(x^b) + \text{Parameter}_\Delta))$ 
   /* Parameter $_\Delta$  is a bias vector with a size of C */
5:  $\bar{A} : (\text{HW}, \text{C}, \text{N}) \leftarrow \exp(\Delta \otimes A)$ 
6:  $\bar{B} : (\text{HW}, \text{C}, \text{N}) \leftarrow \Delta \otimes B$ 
7:  $y^d \leftarrow \text{SSM}(\bar{A}, \bar{B}, C)(x^a)$ 
   /* SSM represents Eq. 3 implemented by selective scan [11] */
return  $y^d$ 

```

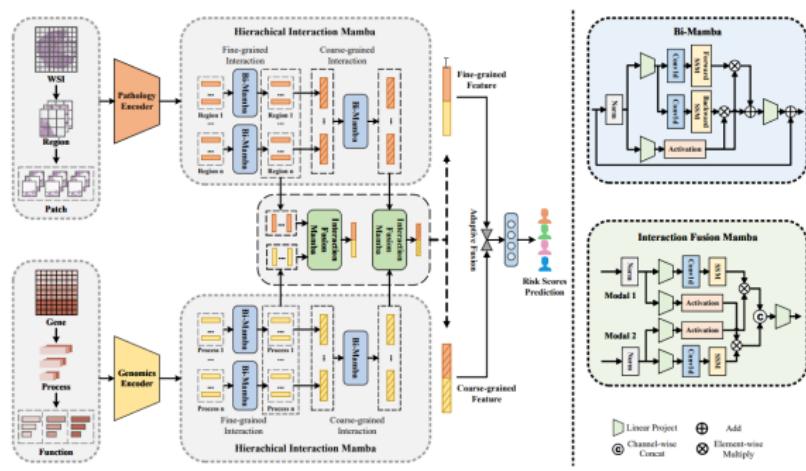


## SurvMamba

### SurvMamba: State Space Model with Multi-grained Multi-modal Interaction for Survival Prediction

- 不同模态交叉 gated
- 拼接融合特征

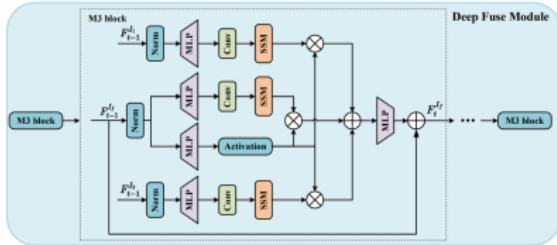
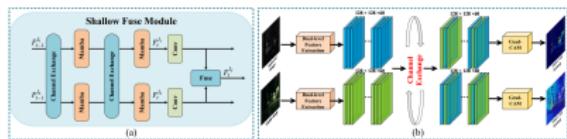
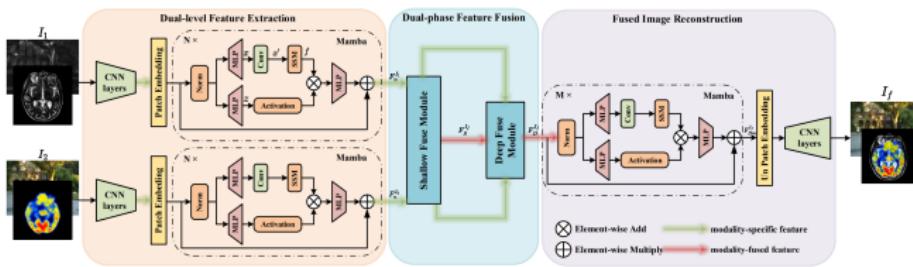
存在问题：这种融合方式需要图像序列的长度和基因序列的长度一致才能融合，而且序列相应位置由对应关系融合才有意义



# MambaDFuse

## MambaDFuse: A Mamba-based Dual-phase Model for Multi-modality Image Fusion

- 浅层融合: 交换通道
- 深层融合: 浅层融合的作 gate



① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

⑦ 3D Mamba

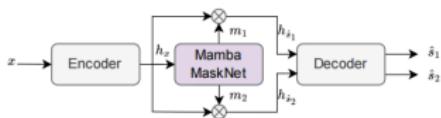
⑧ XAI Mamba

## Frame Title

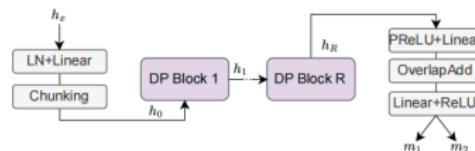
# Dual-path Mamba: Short and Long-term Bidirectional Selective Structured State Space Models for Speech Separation

- 分成多个 chunk, intra-chunk 和 inter-chunk 交替
  - 前向和反向 Mamba 结合

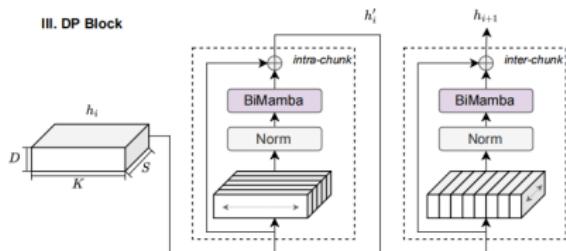
## I. Dual-path Mamba



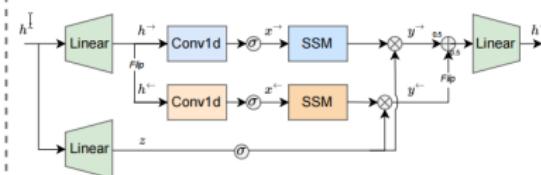
## II. Mamba MaskNet



### III. DP Block



#### **IV. BiMamba**



① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

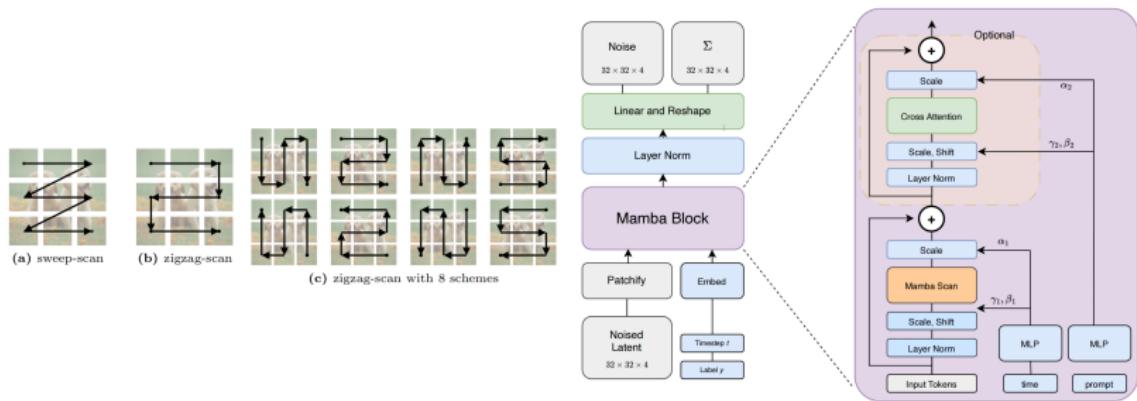
⑦ 3D Mamba

⑧ XAI Mamba

ZigMa

## ZigMa: A DiT-style Zigzag Mamba Diffusion Model

- 基于 DiT 修改，时间信息影响 SSM 模块，prompt 信息对应 CrossAttention
  - 不同的扫描方式



① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

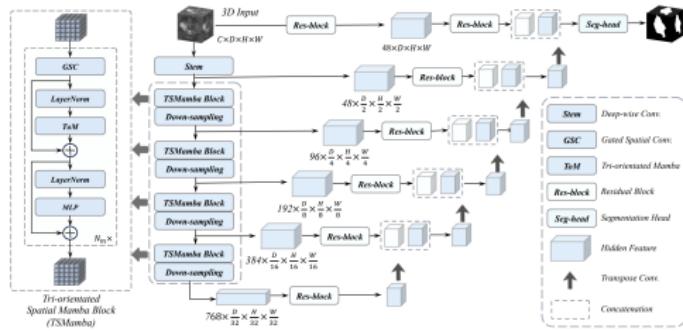
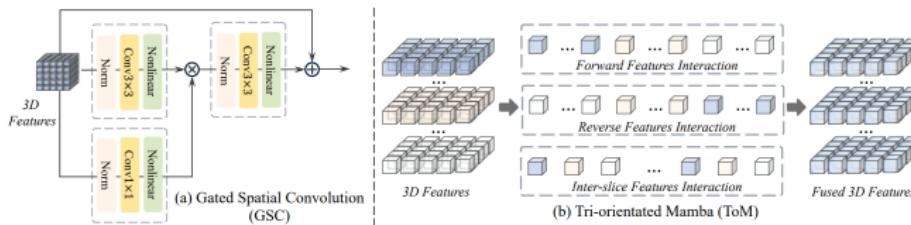
⑦ 3D Mamba

⑧ XAI Mamba

# SegMamba

## SegMamba: Long-range Sequential Modeling Mamba For 3D Medical Image Segmentation

- 将 3D 图像按不同方式展平成 1D 序列



① Mamba

② Vision Mamba

③ MIL Mamba

④ Multimodal Mamba

⑤ Audio Mamba

⑥ Diffusion Mamba

⑦ 3D Mamba

⑧ XAI Mamba

## The Hidden Attention of Mamba Models

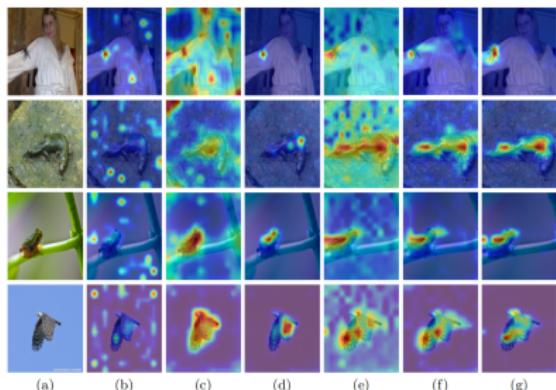
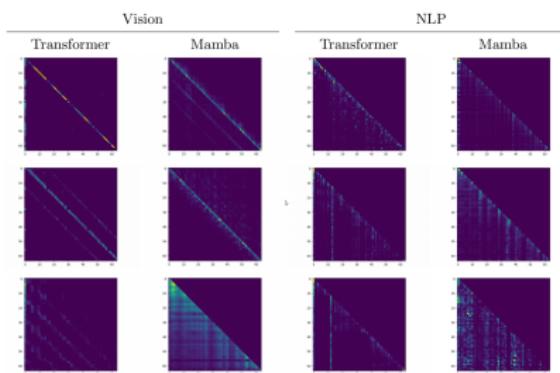
## The Hidden Attention of Mamba Models

$$y = \bar{\alpha}x, \quad \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_L \end{bmatrix} = \begin{bmatrix} C_1\bar{B}_1 & 0 & \cdots & 0 \\ C_2\bar{A}_2\bar{B}_1 & C_2\bar{B}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ C_L\bar{H}_L^L\bar{A}_k\bar{B}_1 & C_L\bar{H}_L^L\bar{A}_k\bar{B}_2 & \cdots & C_L\bar{B}_L \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_L \end{bmatrix}$$

$$\begin{aligned}\vec{Q}_i &:= S_C(\hat{x}_i), \quad \vec{K}_j := \text{ReLU}(S_{\Delta}(\hat{x}_j))S_B(\hat{x}_j), \quad \vec{H}_{i,j} := \exp\left(\sum_{\substack{k=j+1 \\ S_{\Delta}(\hat{x}_k) > 0}}^i S_{\Delta}(\hat{x}_k)\right) A\end{aligned}\tag{18}$$

Eq. 17 can be further simplified to:

$$\tilde{O}_{\perp i} \approx \tilde{O}_i \tilde{H}_{\perp i} \tilde{K}_i \quad (19)$$



**Fig. 5:** Qualitative results for the different explanation methods for the ViT-small and the Mamba-small models. (a) the original image, (b) the aggregated Raw-Attention of ViT-Small, (c) Attention Rollout for ViT-Small, (d) Transformer-Attribution for ViT-Small, (e) the Raw-Attention of Mamba-Small, (f) Attention-Rollout of Mamba-Small and (g) the Mamba-Attribution method for the Mamba-Small model.

Mamba  
oooooooo

Vision Mamba  
oooo

MIL Mamba  
ooo

Multimodal Mamba  
oooo

Audio Mamba  
oo

Diffusion Mamba  
oo

3D Mamba  
oo

XAI Mamba  
ooo●

*Thanks!*