



大模型系列—大模型算法

打破 Transformer



ZOMI

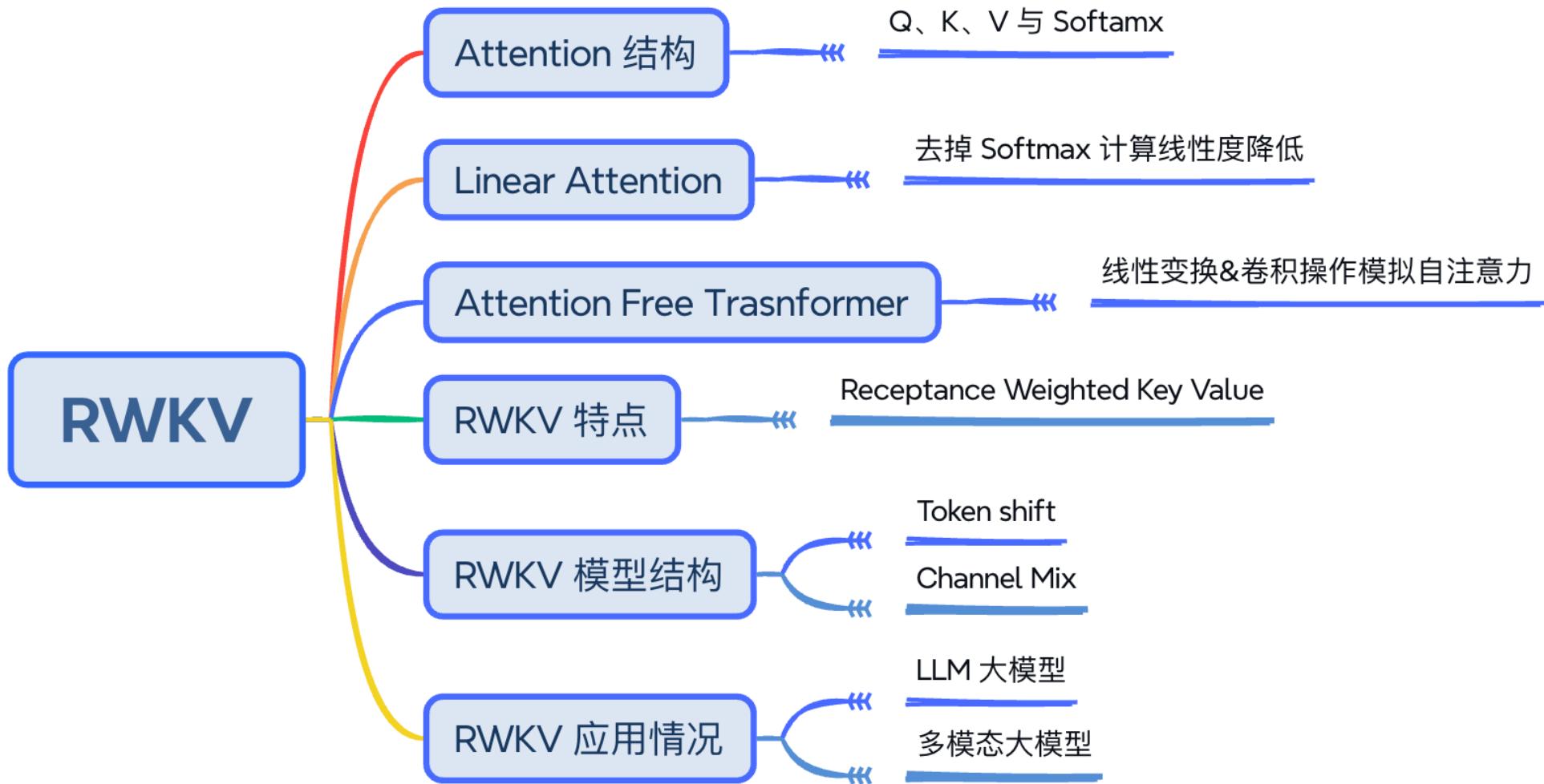
关于本内容

1. Transformer 结构回顾与挑战
2. RWKV 新结构
3. Mamba 新结构
4. 未来打破 Transformer 路在何方

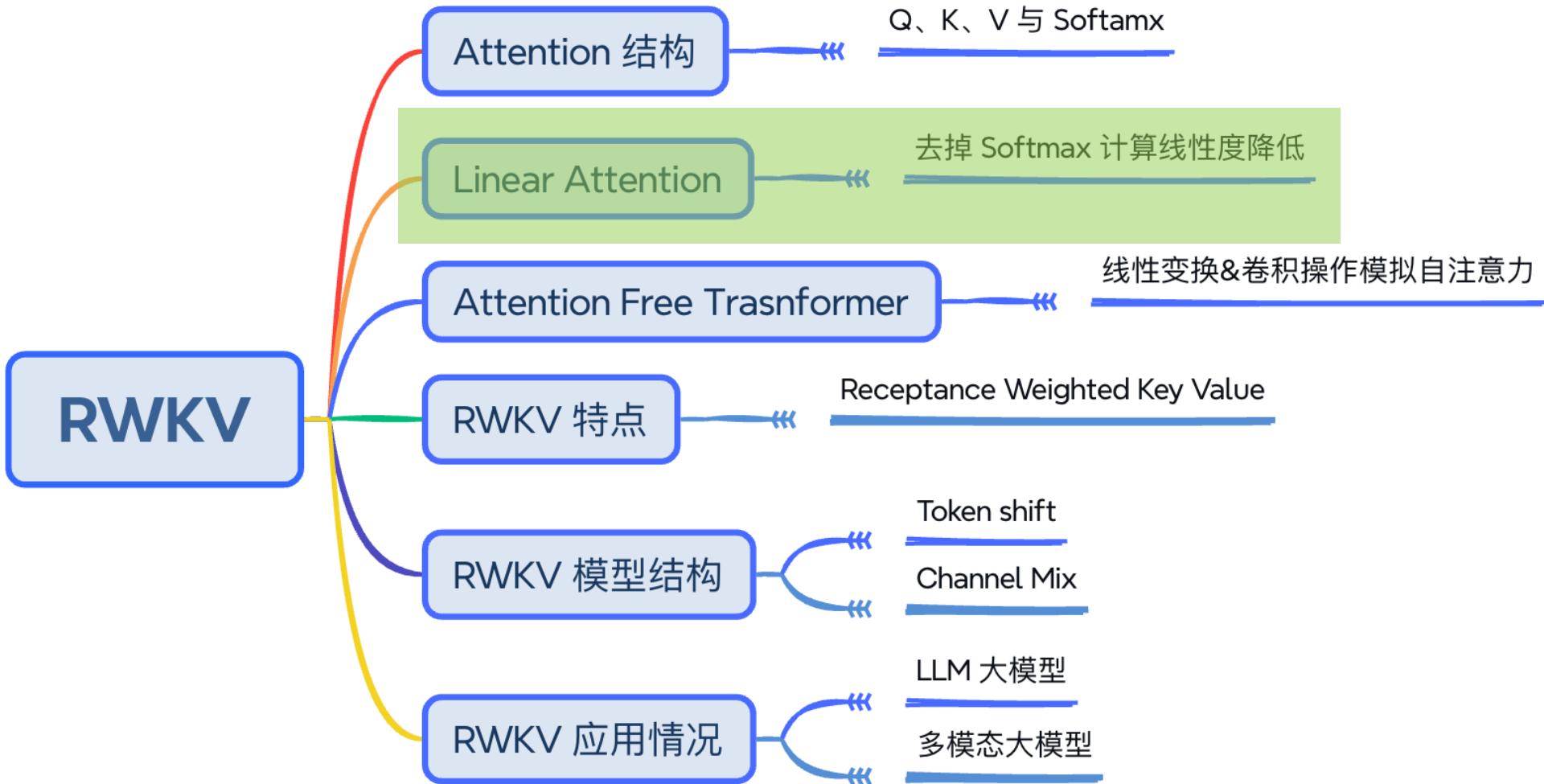
01

RWKV

RWKV



RWKV

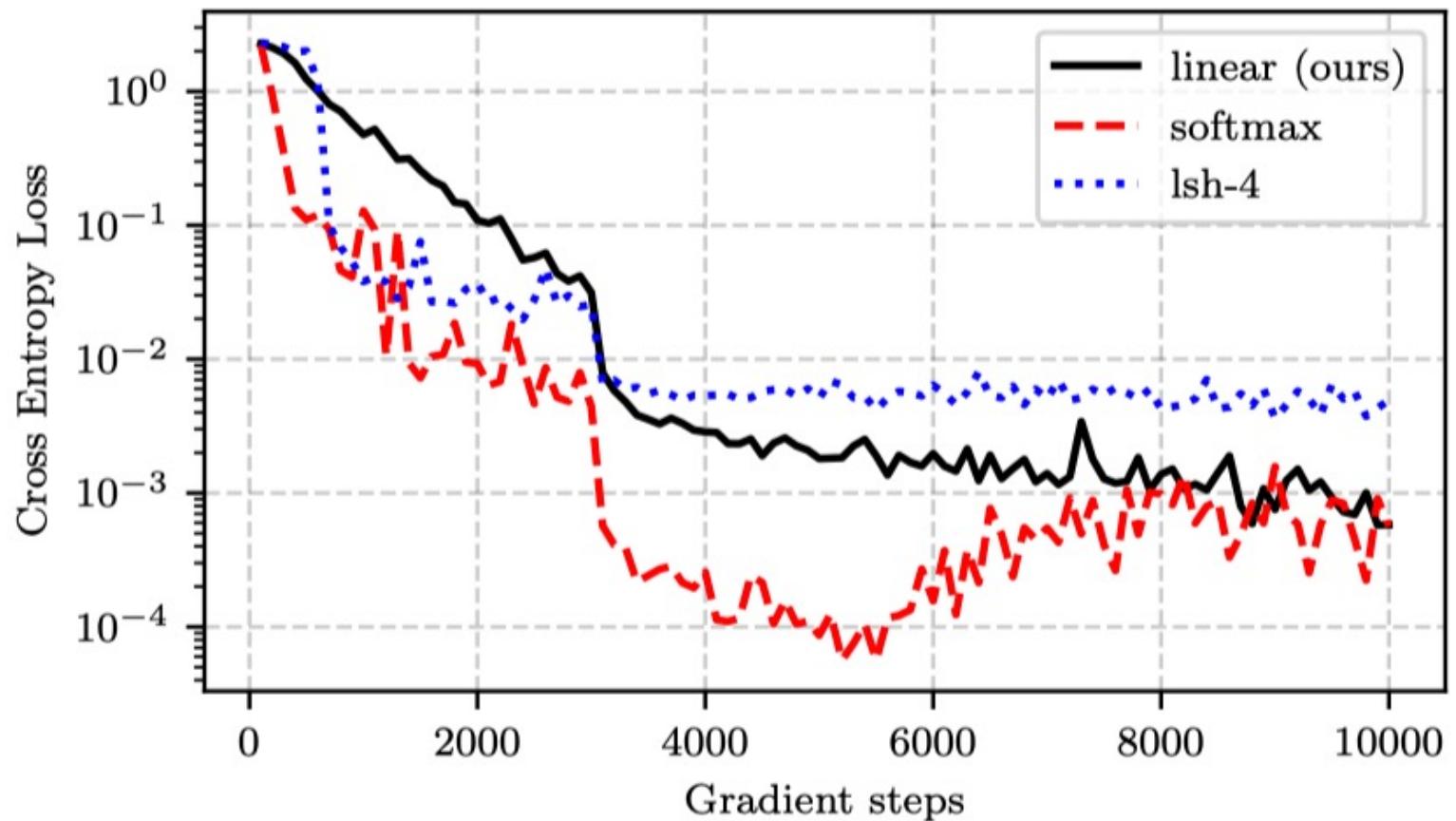


Linear Attention

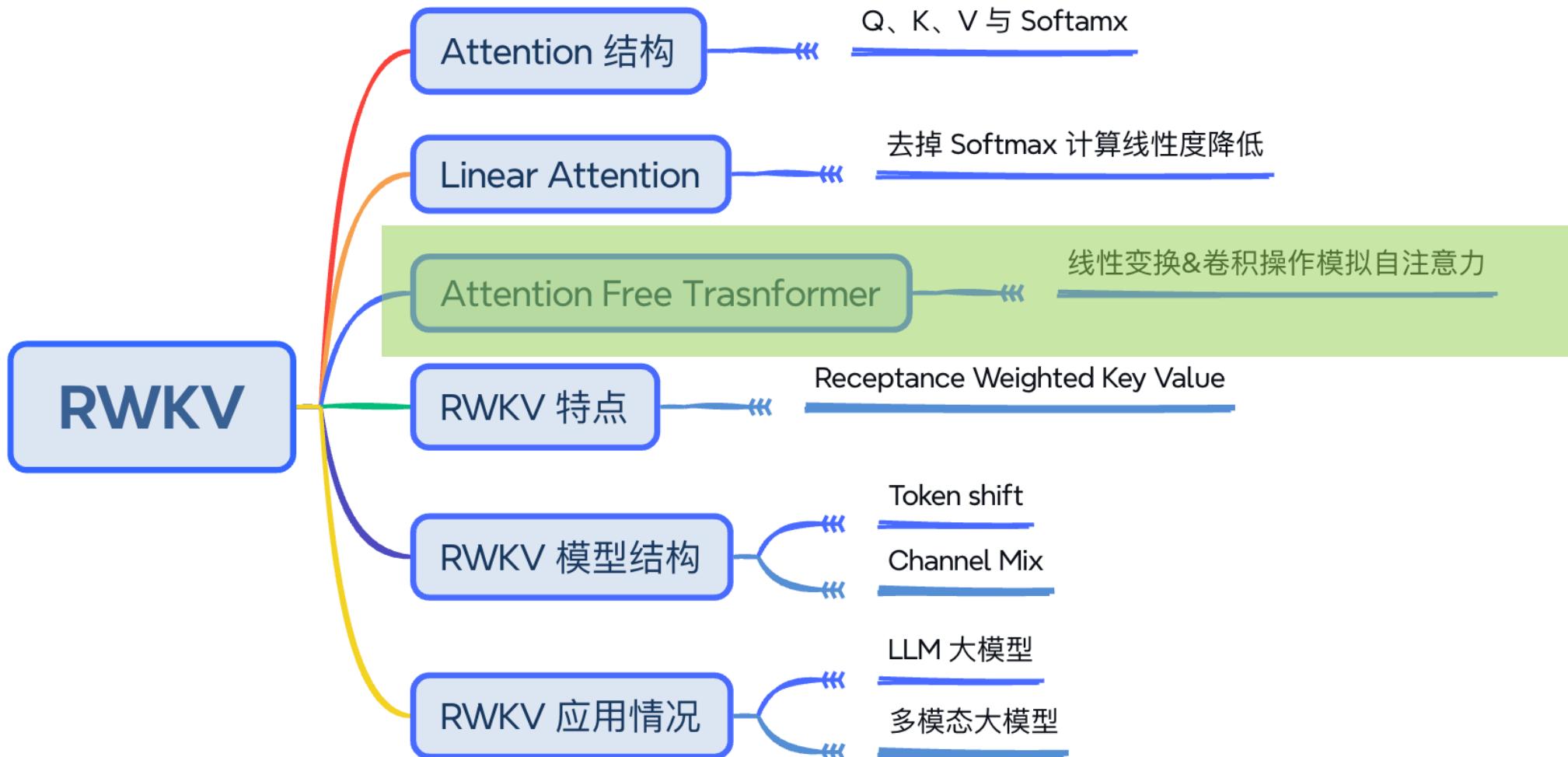
- Transformers are rnns: Fast autoregressive transformers with linear attention.

Algorithm 1 Linear transformers with causal masking

```
function forward( $\phi(Q)$ ,  $\phi(K)$ ,  $V$ ):  
     $V' \leftarrow 0$ ,  $S \leftarrow 0$   
    for  $i = 1, \dots, N$  do  
         $S \leftarrow S + \phi(K_i) V_i^T$  equation 10  
         $\bar{V}_i \leftarrow \phi(Q_i) S$   
    end  
    return  $\bar{V}$   
end  
  
function backward( $\phi(Q)$ ,  $\phi(K)$ ,  $V$ ,  $G$ ):  
    /*  $G$  is the gradient of the loss  
       with respect to the output of  
       forward */  
     $S \leftarrow 0$ ,  $\nabla_{\phi(Q)} \mathcal{L} \leftarrow 0$   
    for  $i = 1, \dots, N$  do  
         $S \leftarrow S + \phi(K_i) V_i^T$  equation 13  
         $\nabla_{\phi(Q_i)} \mathcal{L} \leftarrow G_i S^T$   
    end  
     $S \leftarrow 0$ ,  $\nabla_{\phi(K)} \mathcal{L} \leftarrow 0$ ,  $\nabla_V \mathcal{L} \leftarrow 0$   
    for  $i = N, \dots, 1$  do  
         $S \leftarrow S + \phi(Q_i) G_i^T$  equation 15  
         $\nabla_{V_i} \mathcal{L} \leftarrow S^T \phi(K_i)$  equation 14  
         $\nabla_{\phi(K_i)} \mathcal{L} \leftarrow S V_i$   
    end  
    return  $\nabla_{\phi(Q)} \mathcal{L}$ ,  $\nabla_{\phi(K)} \mathcal{L}$ ,  $\nabla_V \mathcal{L}$   
end
```



RWKV



AFT

- Paper: An attention free transformer.

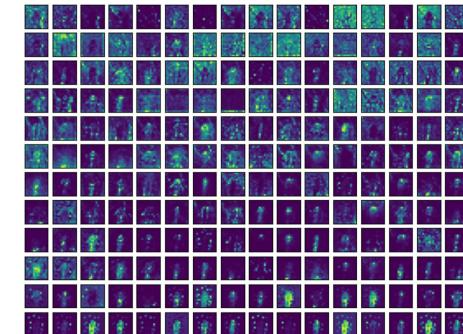
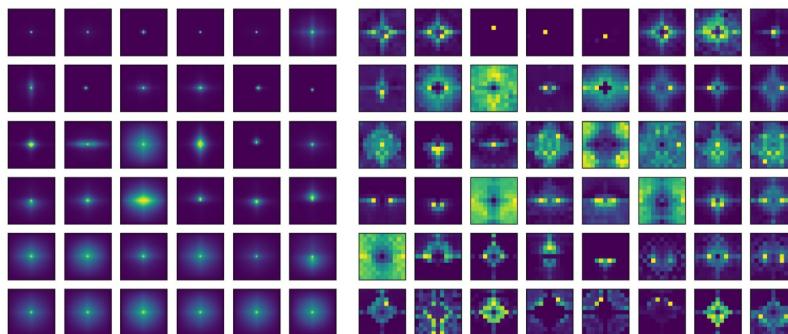
$$\sigma_q(Q_t) \odot \frac{\sum_{t'=1}^T \left[\exp\left(\begin{array}{|c|c|} \hline K & w_t \\ \hline \end{array}\right) + \begin{array}{|c|c|} \hline V & \\ \hline \end{array} \right] \odot \begin{array}{|c|c|} \hline & \\ \hline \end{array}}{\sum_{t'=1}^T \exp\left(\begin{array}{|c|c|} \hline K & w_t \\ \hline \end{array}\right)} = Y_t$$

Figure 2: An illustration of AFT defined in Equation 2, with $T = 3, d = 2$.

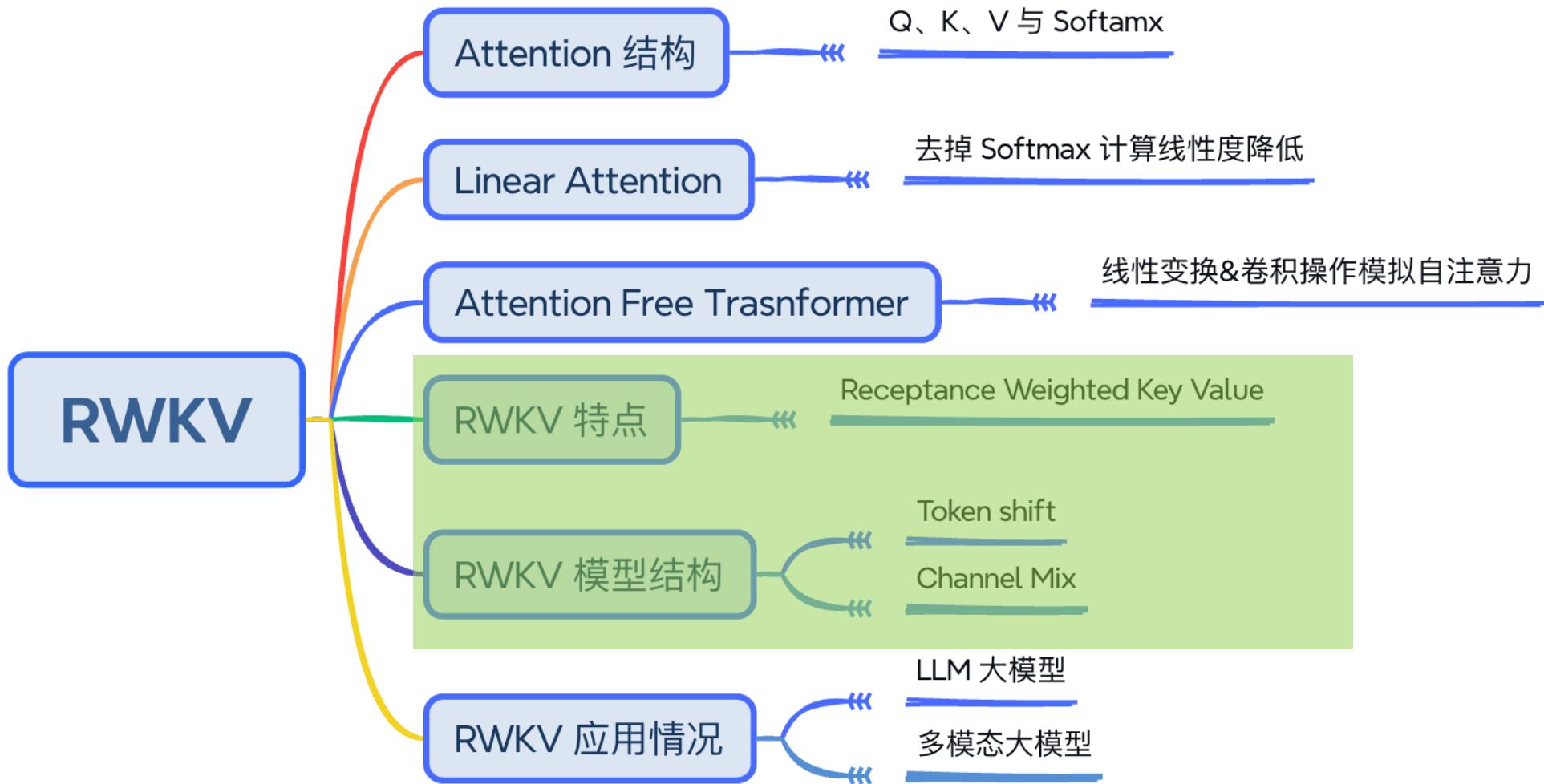
AFT

- Paper: An attention free transformer.

Model	Time	Space
Transformer	$O(T^2d)$	$O(T^2 + Td)$
Reformer	$O(T \log Td)$	$O(T \log T + Td)$
Linear Transformer	$O(Td^2)$	$O(Td + d^2)$
Performer	$O(Td^2 \log d)$	$O(Td \log d + d^2 \log d)$
AFT-simple	$O(\mathbf{Td})$	$O(\mathbf{Td})$
AFT-full	$O(T^2d)$	$O(\mathbf{Td})$
AFT-local (AFT-conv)	$O(Tsd), s < T$	$O(\mathbf{Td})$

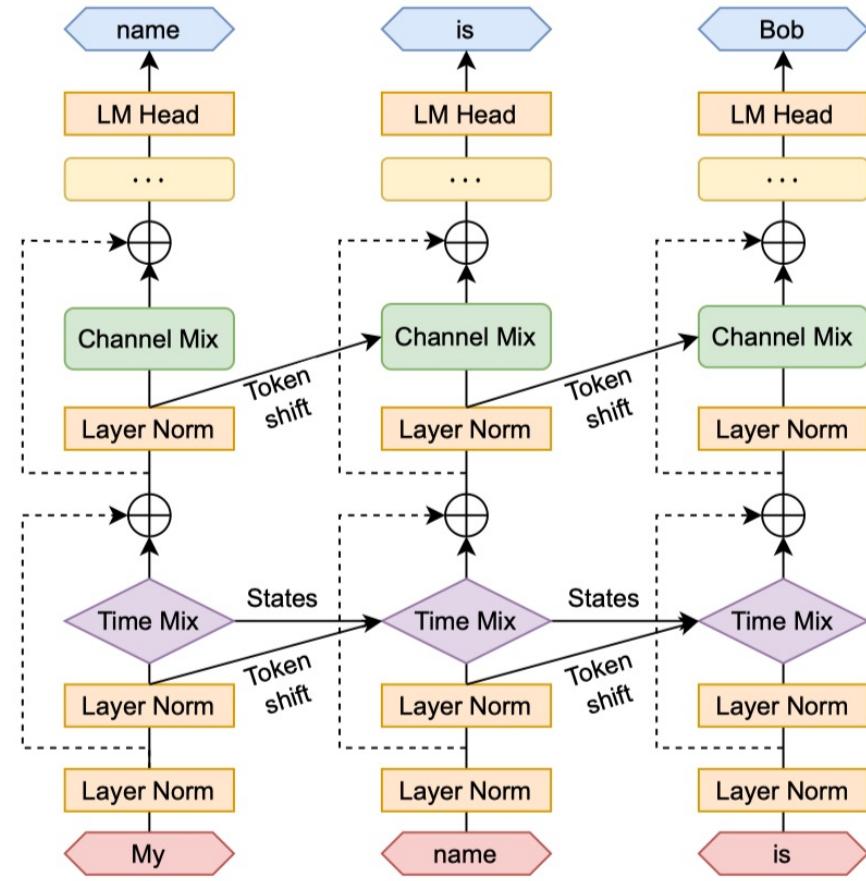
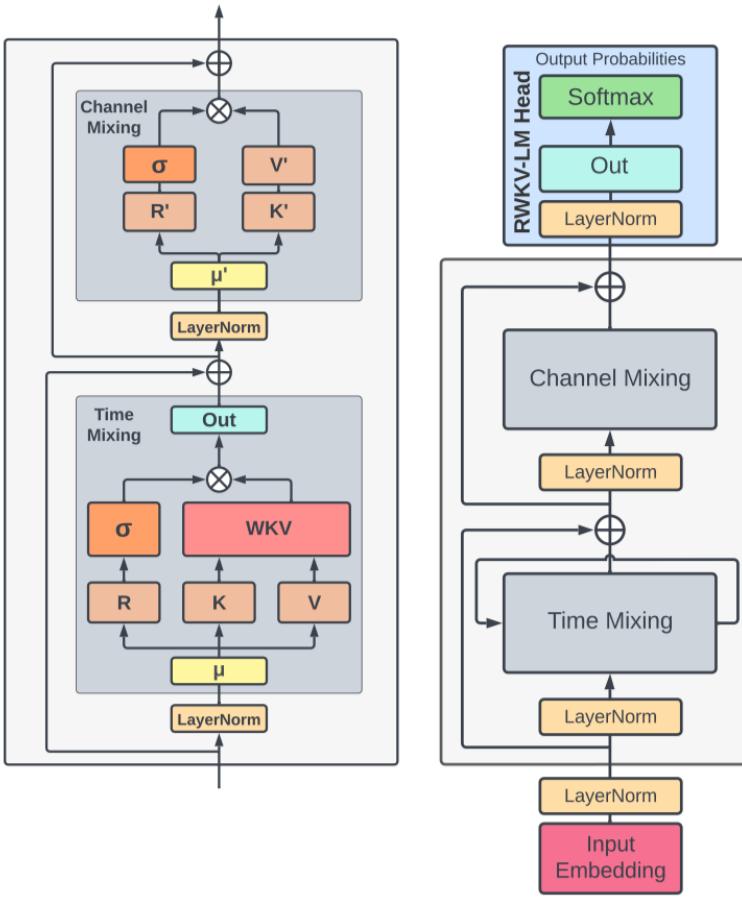


RWKV



RWKV

- RWKV: Reinventing RNNs for the Transformer Era

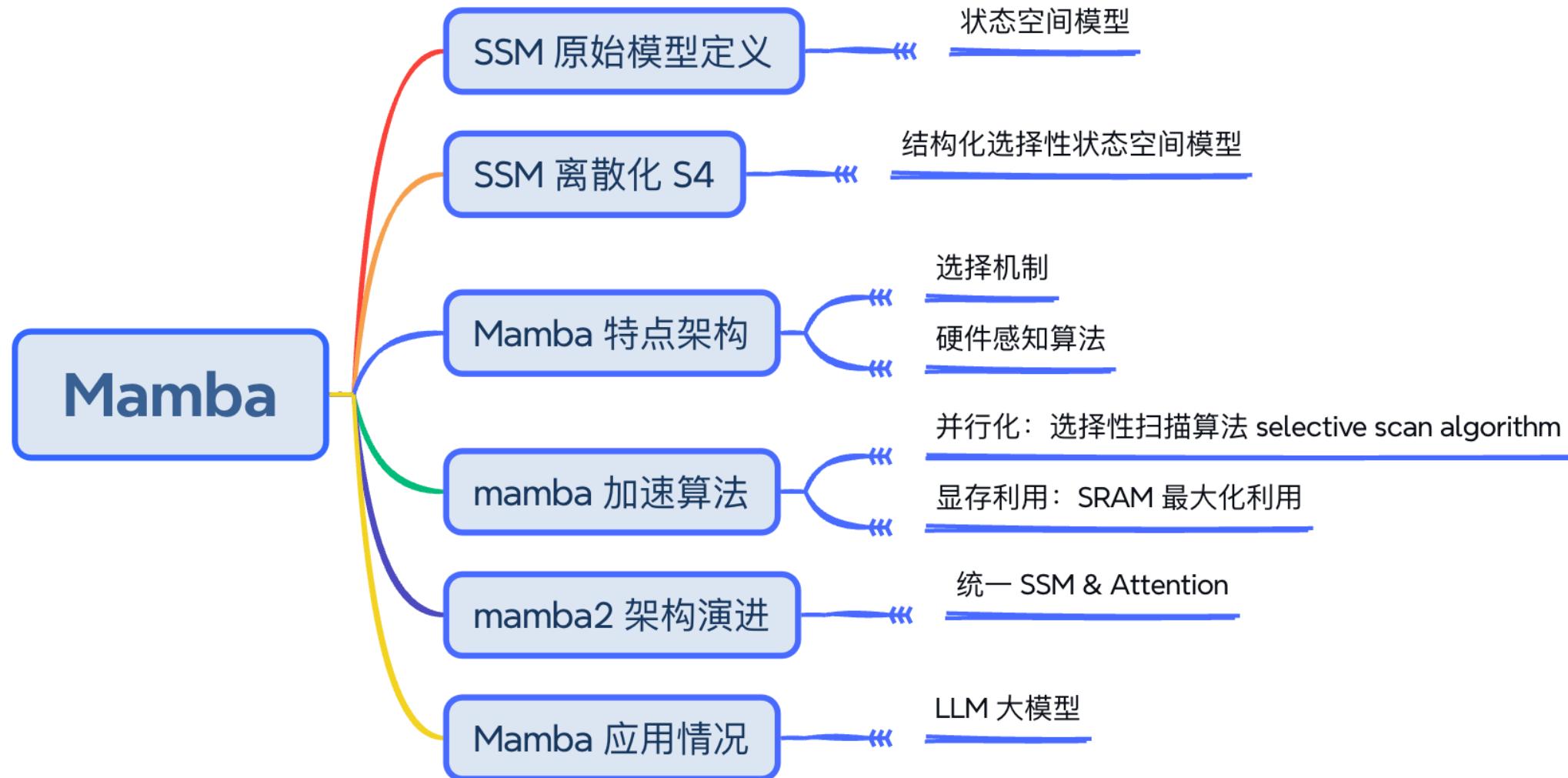


Reference 参考&引用

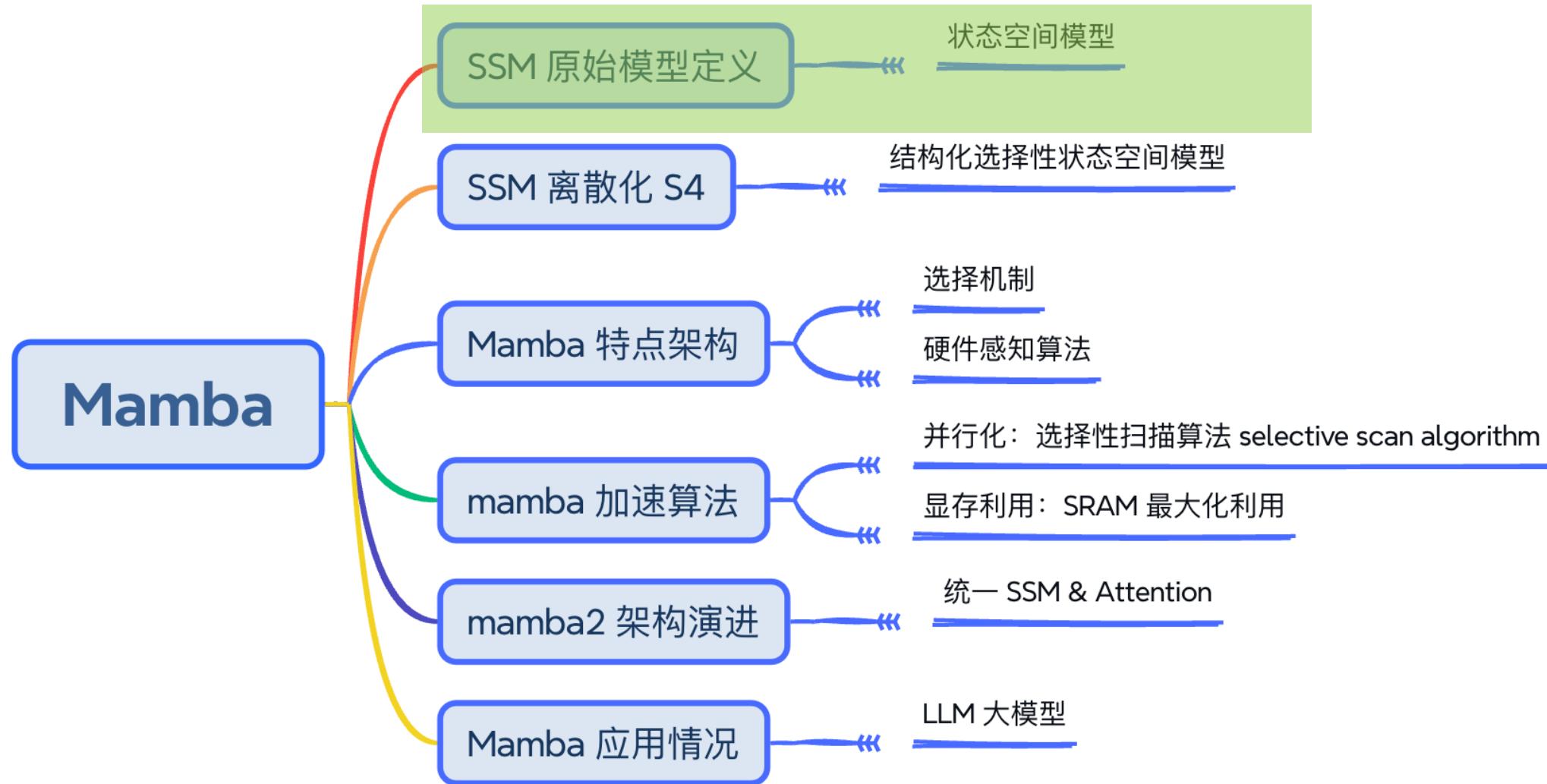
1. <https://arxiv.org/abs/2006.16236>
2. <https://www.rwkv.com/>
3. <https://arxiv.org/abs/2105.14103>
4. <https://github.com/BlinkDL/RWKV-LM>
5. <https://wiki.rwkv.com/>

02 Mamba

Mamba

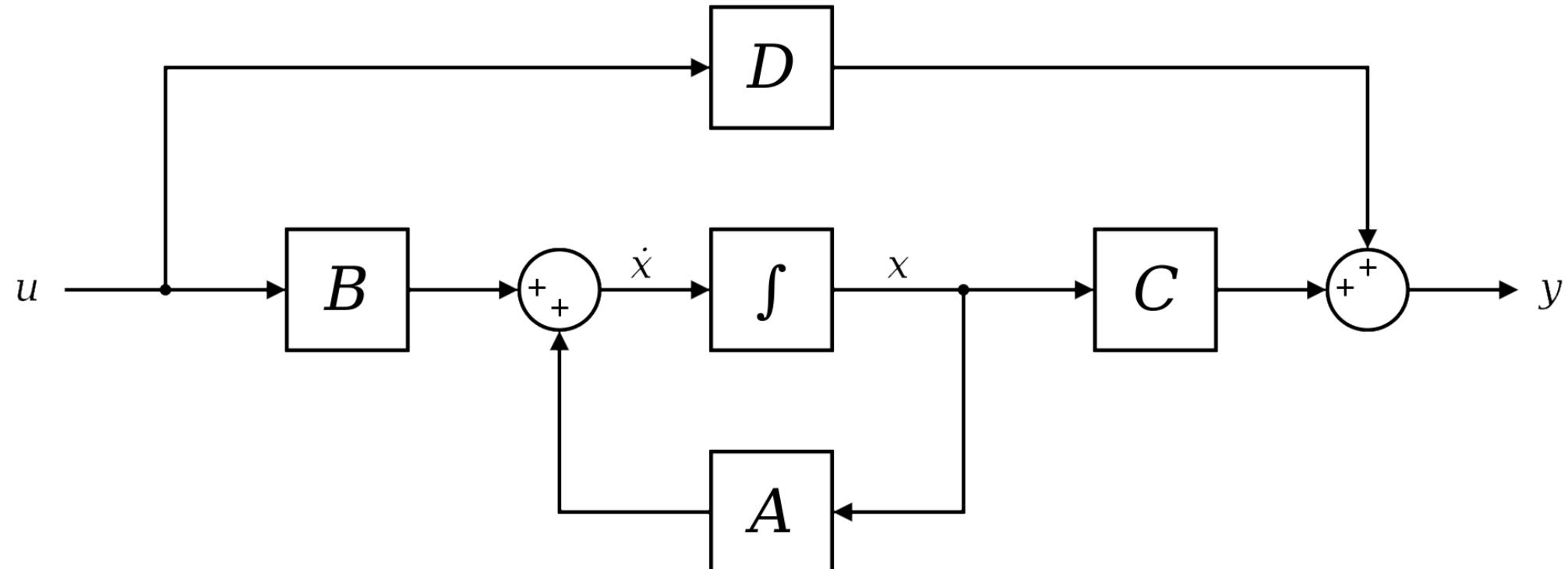


Mamba

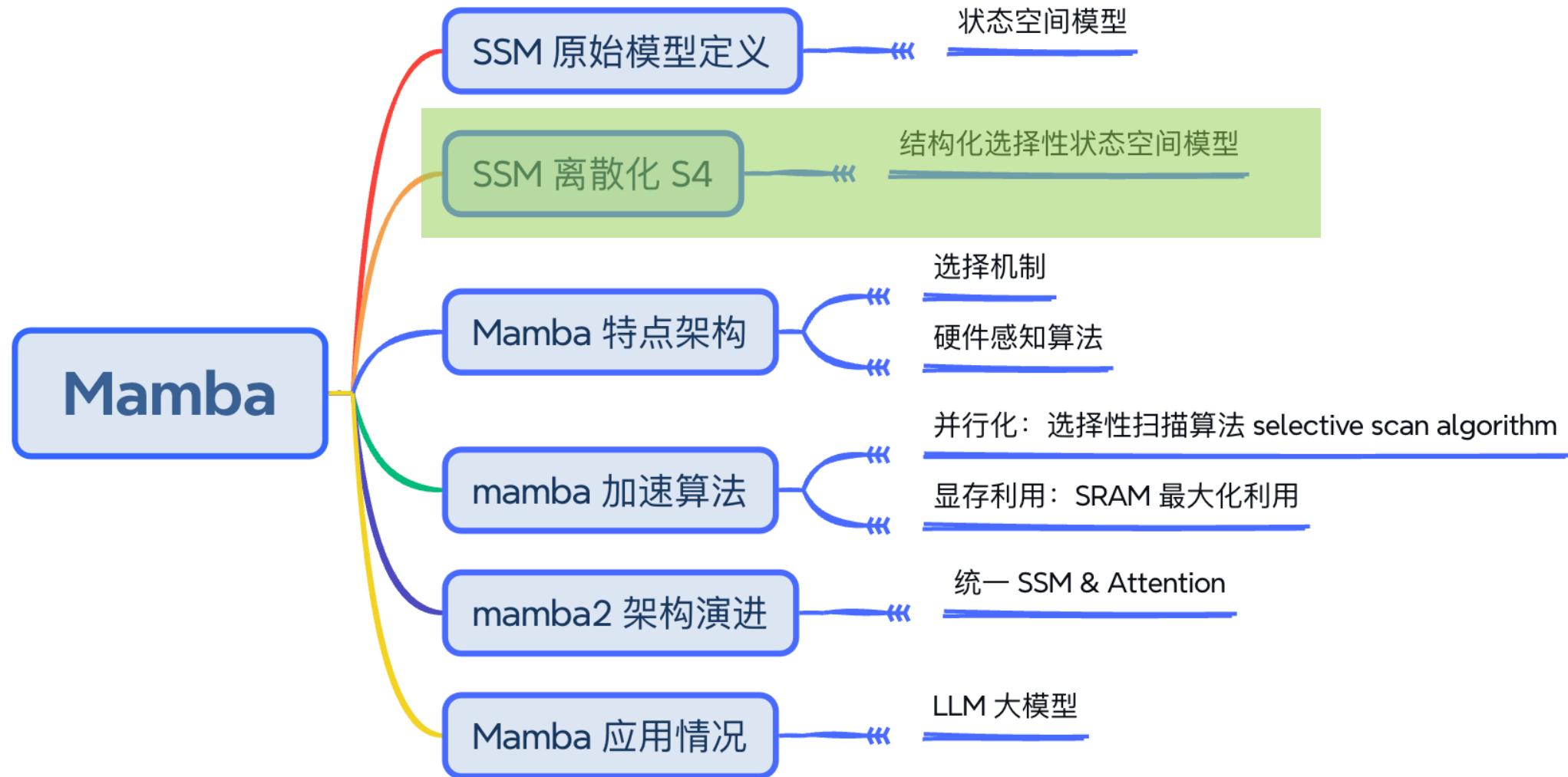


SSM

- Intelligent Control Systems: An Introduction with Examples
- State-space representation → State Space Model

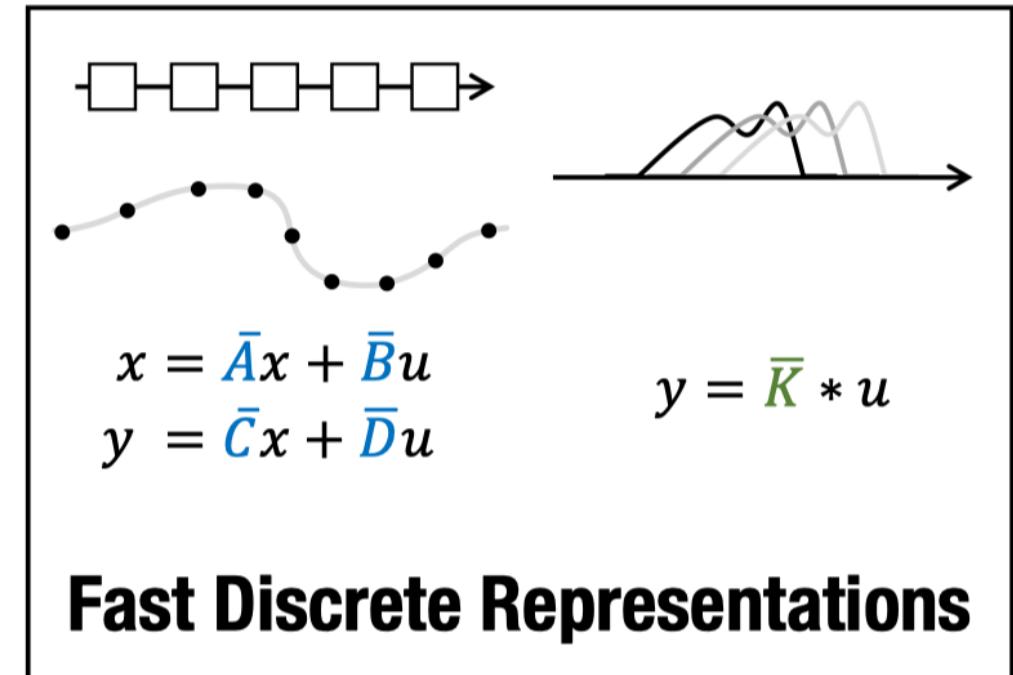
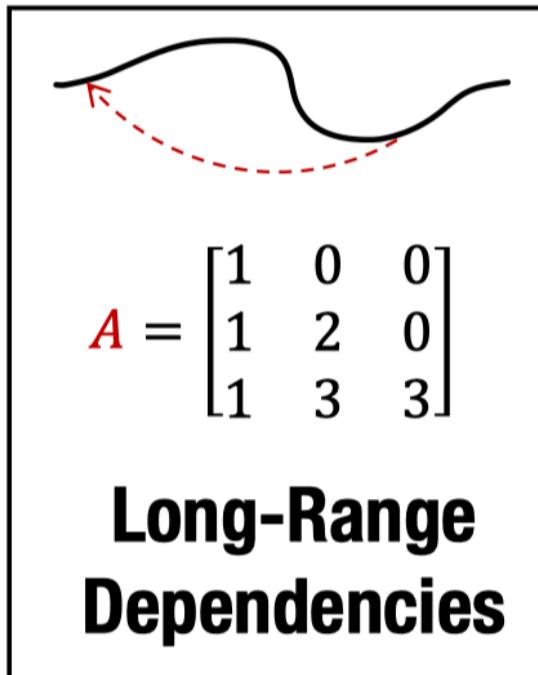
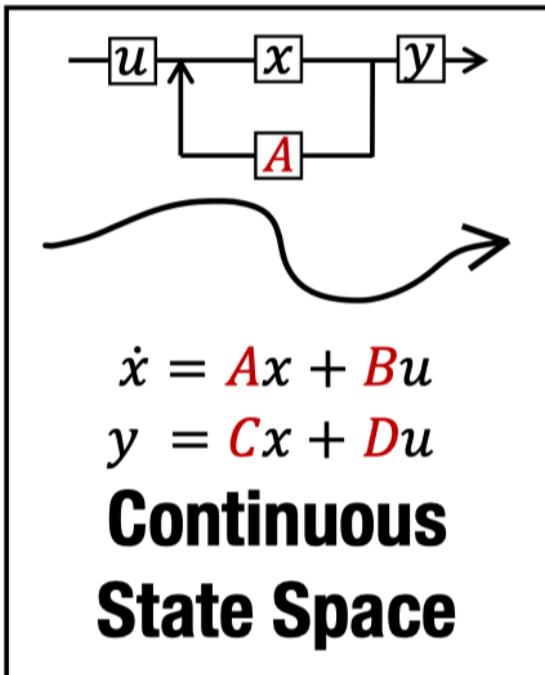


Mamba

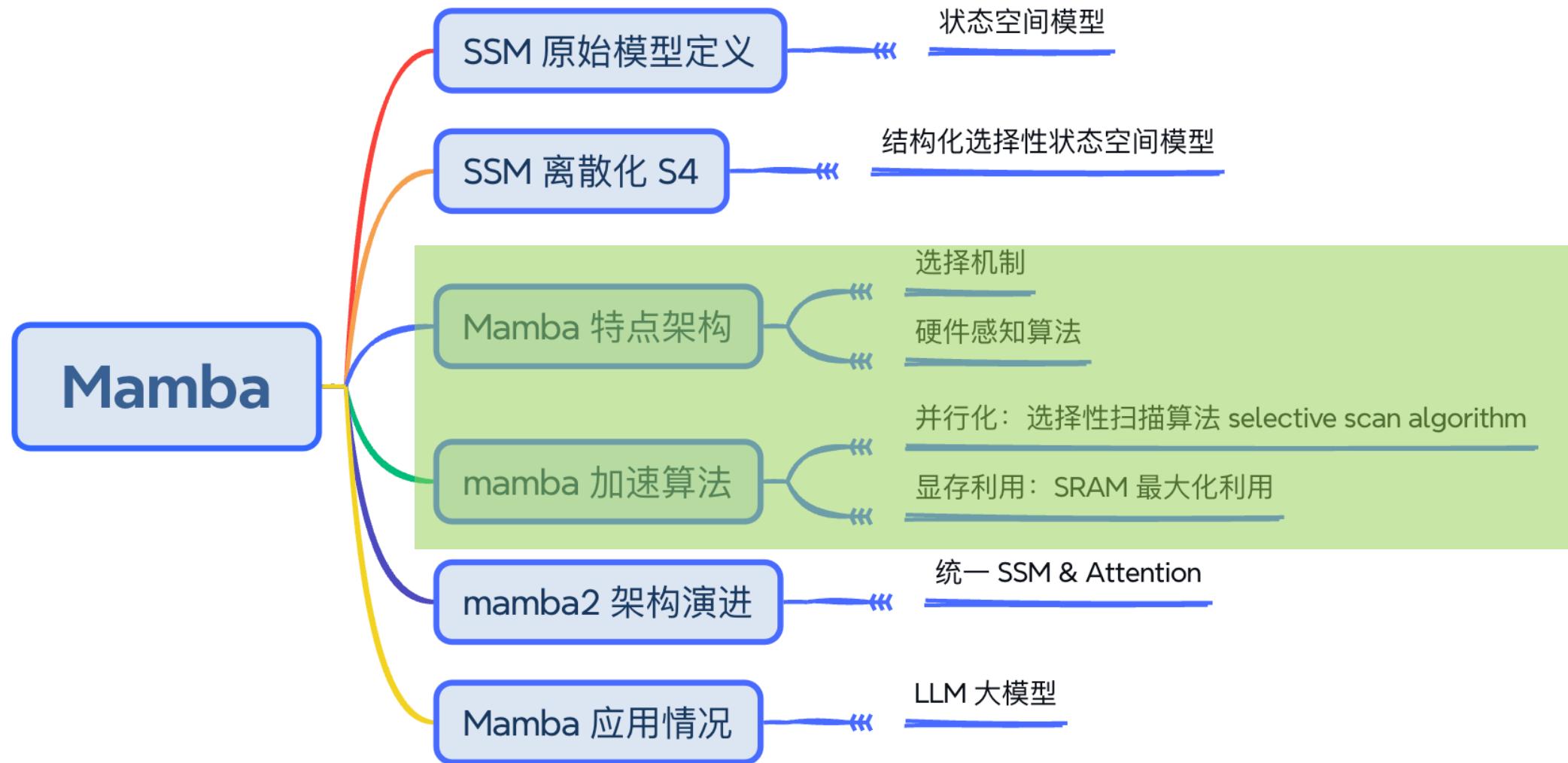


S4

- Paper: Efficiently Modeling Long Sequences with Structured State Spaces
- aka. Structured State Space sequence model (S4)



Mamba



mamba

Mamba: Linear-Time Sequence Modeling with Selective State Spaces

Albert Gu^{*}¹ and Tri Dao^{*}²

¹Machine Learning Department, Carnegie Mellon University

²Department of Computer Science, Princeton University

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

FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness

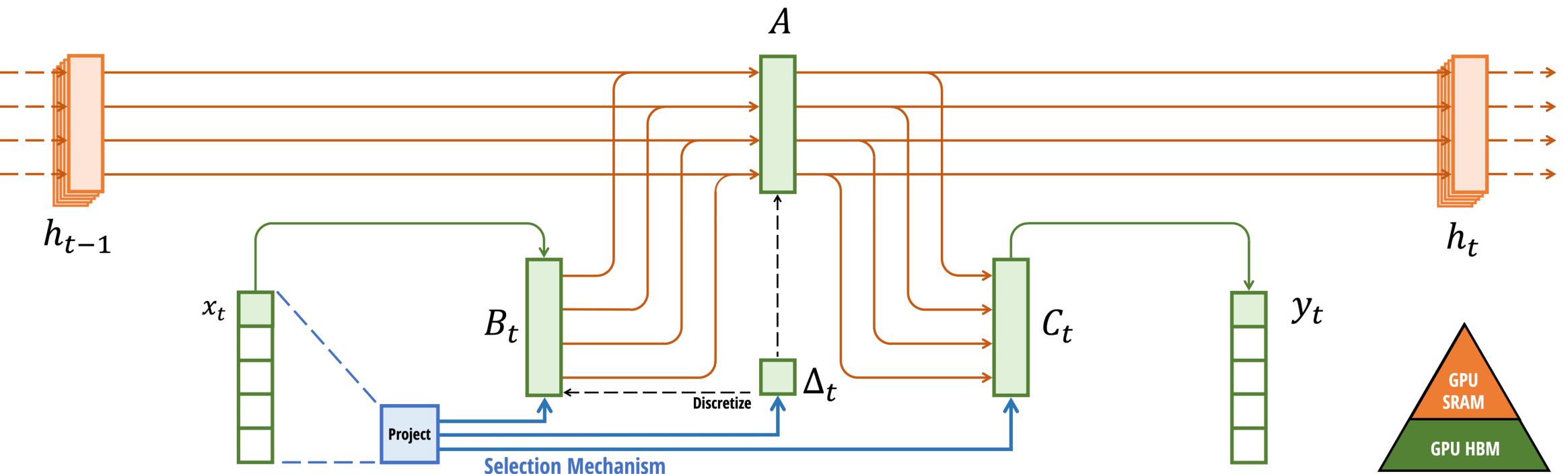
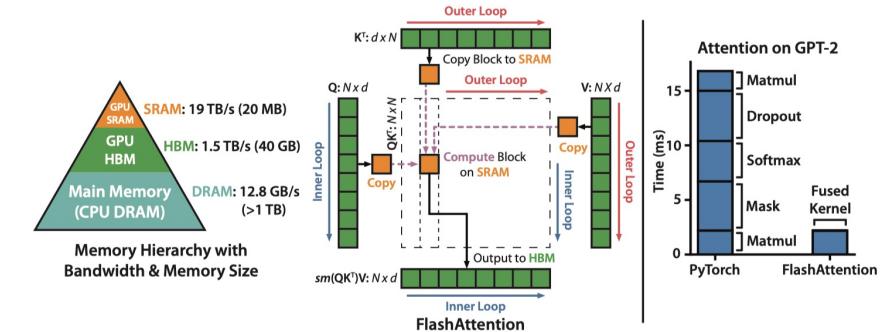
Tri Dao[†], Daniel Y. Fu[†], Stefano Ermon[†], Atri Rudra[‡], and Christopher Ré[†]

[†]Department of Computer Science, Stanford University

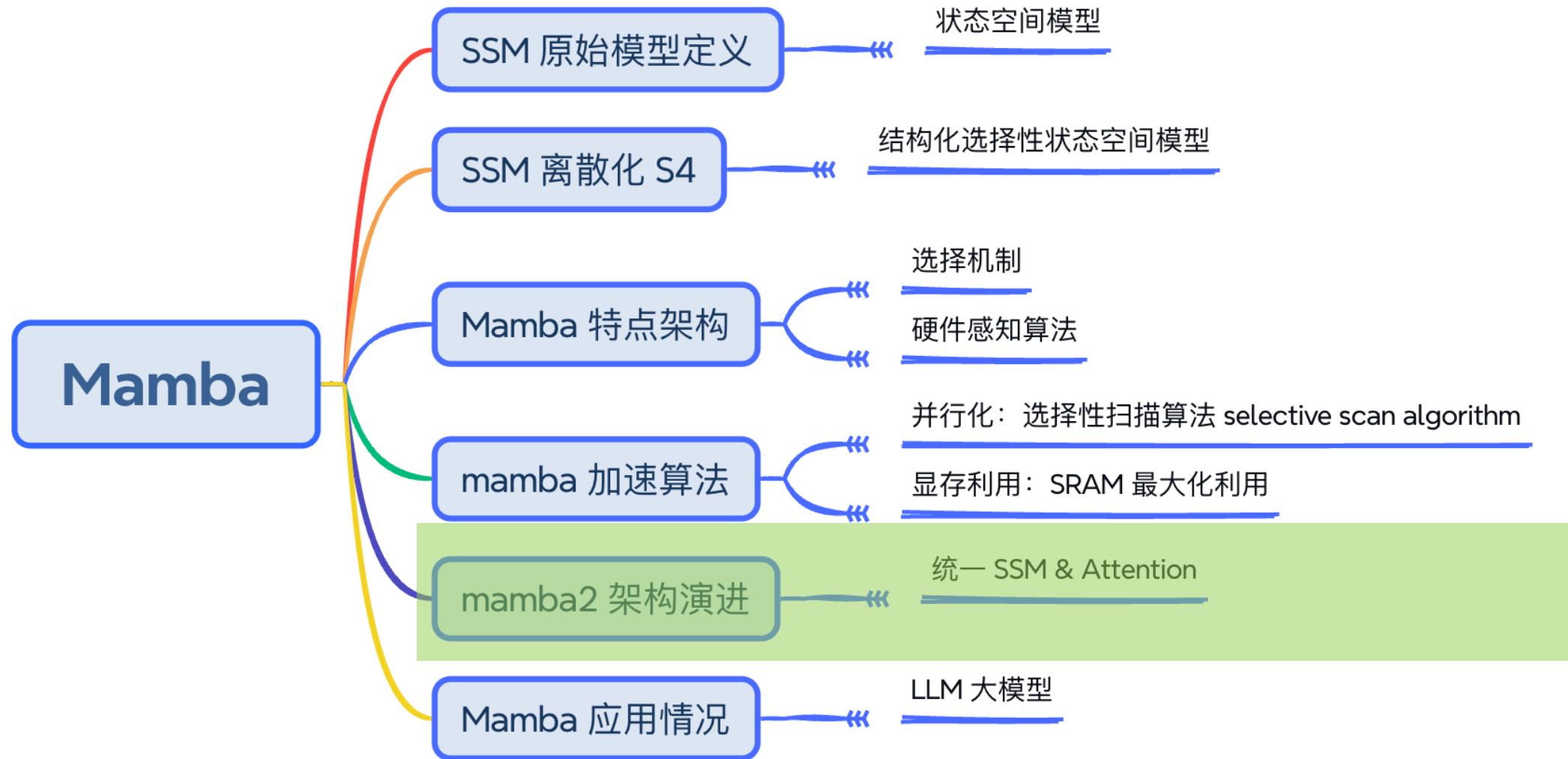
[‡]Department of Computer Science and Engineering, University at Buffalo, SUNY

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

Selective State Space Model with Hardware-aware State Expansion



Mamba



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

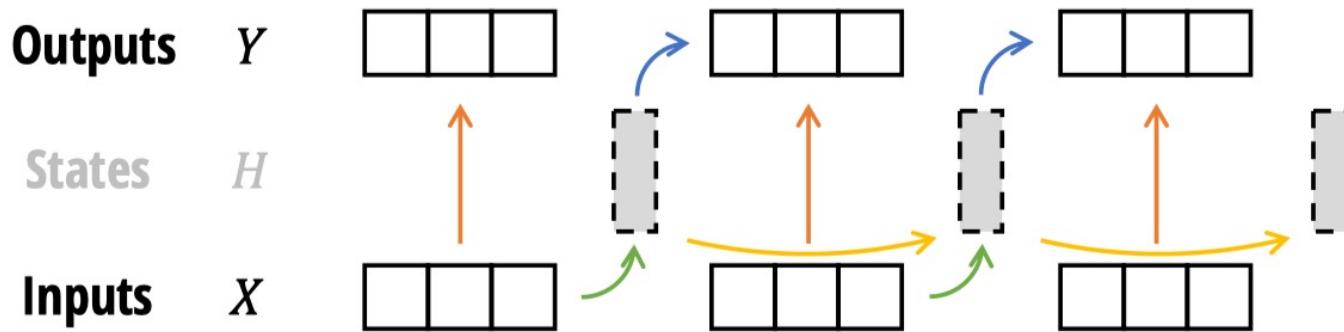
Tri Dao^{*¹} and Albert Gu^{*²}

¹Department of Computer Science, Princeton University

²Machine Learning Department, Carnegie Mellon University
tri@tridao.me, agu@cs.cmu.edu

mamba2

$$\begin{bmatrix}
 C_0^\top A_{0:0}B_0 & & & \\
 C_1^\top A_{1:0}B_0 & C_1^\top A_{1:1}B_1 & & \\
 C_2^\top A_{2:0}B_0 & C_2^\top A_{2:1}B_1 & C_2^\top A_{2:2}B_2 & \\
 \hline
 \left[\begin{array}{c} C_3^\top A_{3:2} \\ C_4^\top A_{4:2} \\ C_5^\top A_{5:2} \end{array} \right] A_{2:2} & \left[\begin{array}{c} B_0^\top A_{2:0} \\ B_1^\top A_{2:1} \\ B_2^\top A_{2:2} \end{array} \right]^\top & C_3^\top A_{3:3}B_3 & \\
 & & C_4^\top A_{4:3}B_3 & C_4^\top A_{4:4}B_4 \\
 & & C_5^\top A_{5:3}B_3 & C_5^\top A_{5:4}B_4 & C_5^\top A_{5:5}B_5 \\
 \hline
 \left[\begin{array}{c} C_6^\top A_{6:5} \\ C_7^\top A_{7:5} \\ C_8^\top A_{8:5} \end{array} \right] A_{5:2} & \left[\begin{array}{c} B_0^\top A_{2:0} \\ B_1^\top A_{2:1} \\ B_2^\top A_{2:2} \end{array} \right]^\top & \left[\begin{array}{c} C_6^\top A_{6:5} \\ C_7^\top A_{7:5} \\ C_8^\top A_{8:5} \end{array} \right] A_{5:5} & \left[\begin{array}{c} B_3^\top A_{5:3} \\ B_4^\top A_{5:4} \\ B_5^\top A_{5:5} \end{array} \right]^\top & C_6^\top A_{6:6}B_6 \\
 & & & C_7^\top A_{7:6}B_6 & C_7^\top A_{7:7}B_7 \\
 & & & C_8^\top A_{8:6}B_6 & C_8^\top A_{8:7}B_7 & C_8^\top A_{8:8}B_8
 \end{bmatrix}$$



Semiseparable Matrix M

Block Decomposition

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

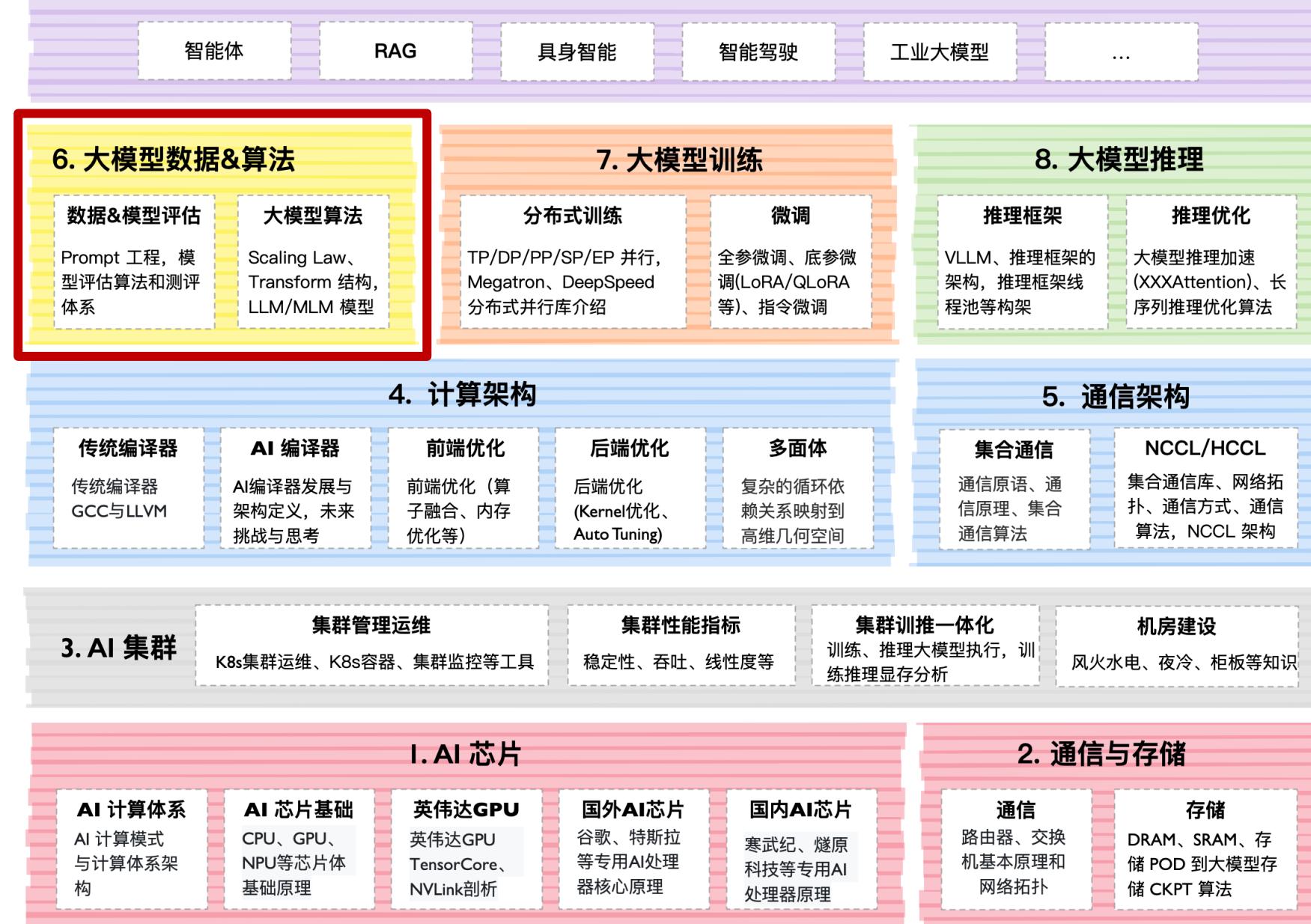
Reference 参考&引用

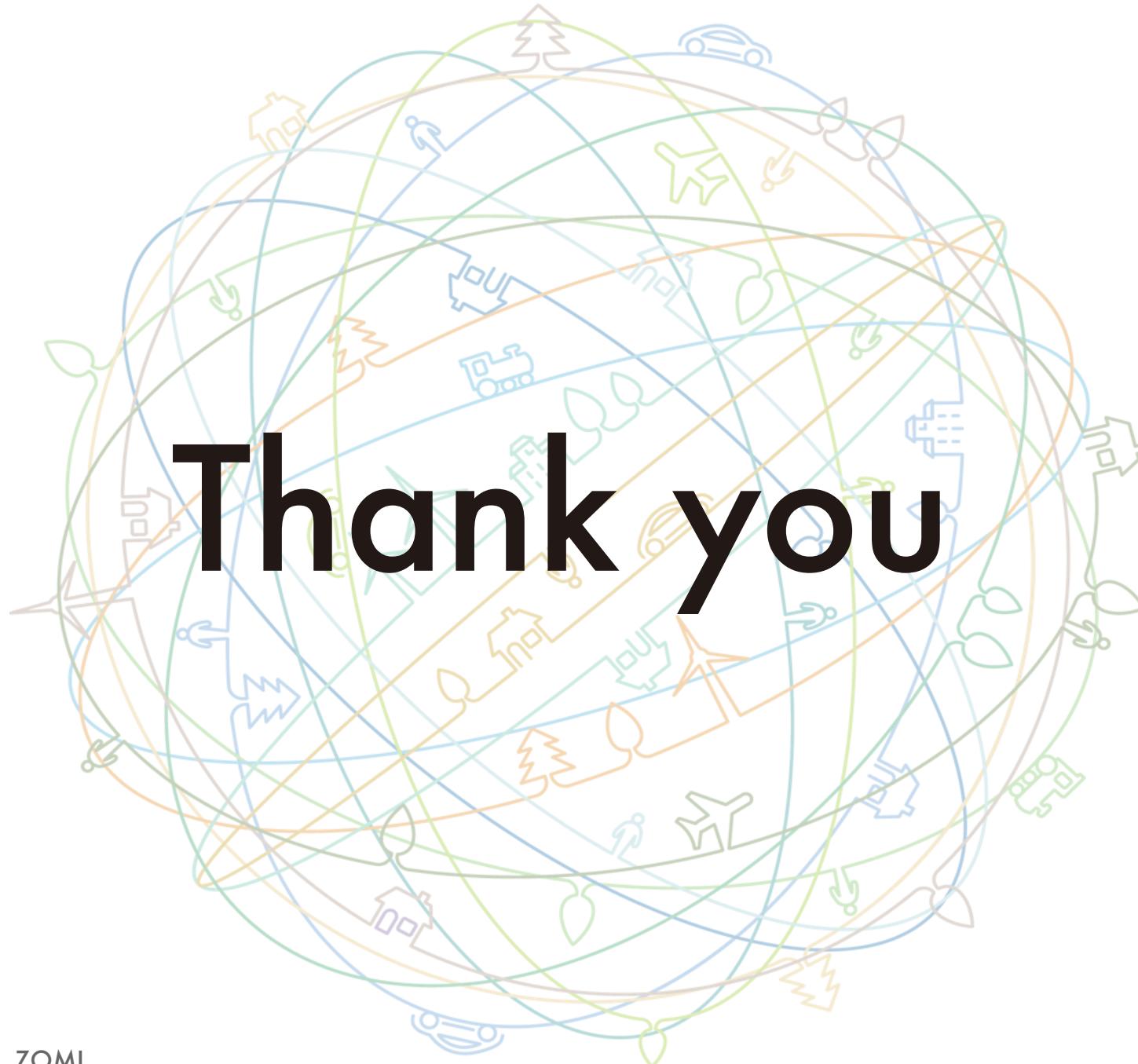
1. <https://srush.github.io/annotated-s4/>
2. <https://huggingface.co/blog/lbourdois/get-on-the-ssm-train>
3. https://en.wikipedia.org/wiki/State-space_representation
4. <https://arxiv.org/abs/2405.21060>
5. <https://github.com/state-spaces/mamba?tab=readme-ov-file>



03

小结与思考





把AI系统带入每个开发者、每个家庭、
每个组织，构建万物互联的智能世界

Bring AI System to every person, home and
organization for a fully connected,
intelligent world.

Copyright © 2023 XXX Technologies Co., Ltd.
All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. XXX may change the information at any time without notice.



Course chenzomi12.github.io

GitHub github.com/chenzomi12/DeepLearningSystem