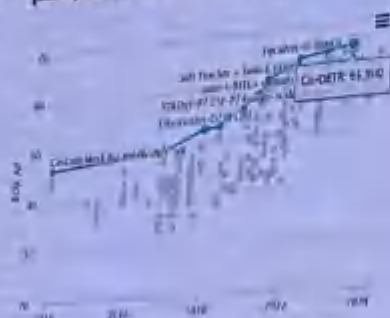


Transformer “一统江湖”

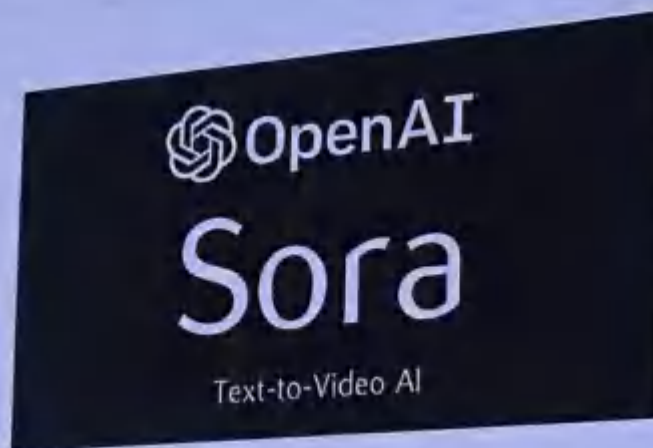
- 在各种理解、生成任务上性能领先



图像理解

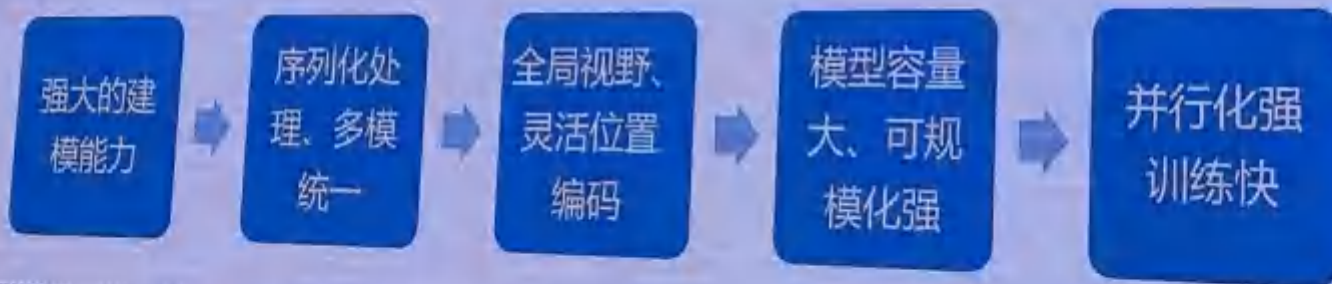


文本生成



视频生成

- why?



VALSE 2024 APR: 面向大模型的新型高效率网络架构

Transformer的局限性

- 计算和存储复杂度为 $O(n^2)$ ，随着序列长度 n 呈二次方增长
- 在多模态大模型、视频生成、具身基础模型等场景下都存在长序列的问题



黄仁勋和Transformer八子

"I think the world needs something better than the transformer," said Gomez. "I think all of us here hope it gets succeeded by something that will carry us to a new plateau of performance."

NLP: RWKV

- 核心思路：在Attention Free Transformer (AFT) 的成对位置偏置矩阵中引入时序衰退机制，能够在训练时和Transformer一样并行，推理的时候可以等价于RNN。
- 关键方法：引入门控机制调制对历史信息的接收，引入Token shift 技术加强对相邻局部特征的融合。

$$AFT(W, K, V)_t = \frac{\sum_{i=1}^t e^{w_{i,t} + k_i} \odot v_i}{\sum_{i=1}^t e^{w_{i,t} + k_i}}$$

↓ 引入时序衰退机制

$$w_{i,t} = -(t-i)w$$

$$wkv_t = \frac{\sum_{i=1}^{t-1} e^{-(t-1-i)w + k_i} \odot v_i + e^{u+k_t} \odot v_t}{\sum_{i=1}^{t-1} e^{-(t-1-i)w + k_i} + e^{u+k_t}}$$

↓ 引入时序门控机制

$$RWKV: o_t = W_o \cdot (\sigma(r_t) \odot wkv_t)$$

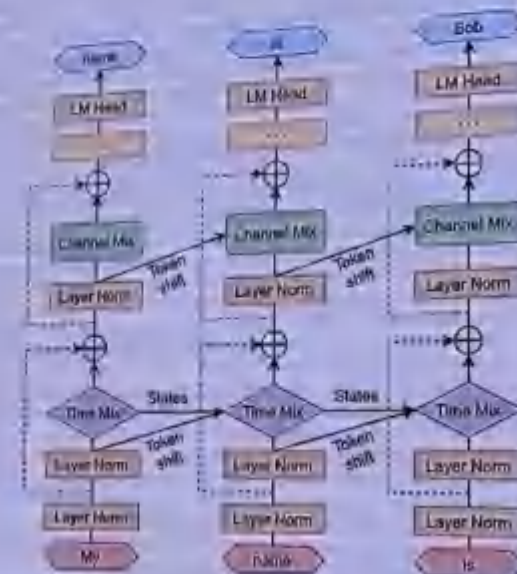
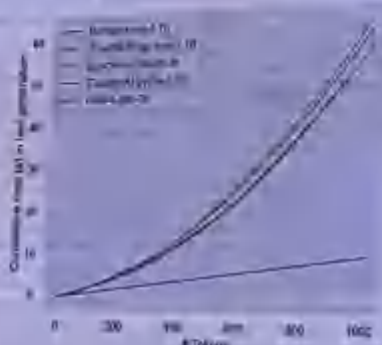
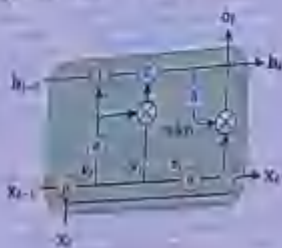


Figure 3: RWKV architecture for language modeling

Peng et al. RWKV: Reinventing RNNs for the Transformer Era, In arXiv, 2023
 VALSE 2024 APR: 面向大模型的新型高效率网络架构

NLP: RetNet

- 引入Retention机制, 同时具有并行和递归的表达机制, 使其能够在训练时像Transformer一样并行, 推理时和RNN一样高效
- 提出时序衰减技术和多头门控机制, 提升模型的表达性

递归形式, 用于推理

$$s_n = A s_{n-1} + K_n^T v_n, \quad A \in \mathbb{R}^{d \times d}, K_n \in \mathbb{R}^{1 \times d}$$

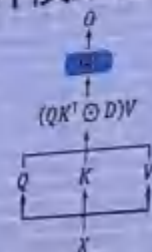
$$o_n = Q_n s_n = \sum_{m=1}^n Q_n A^{n-m} K_m^T v_m, \quad Q_n \in \mathbb{R}^{1 \times d}$$

并行形式, 用于训练

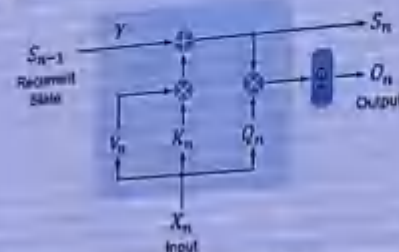
$$Q = (XW_Q) \odot \Theta, \quad K = (XW_K) \odot \bar{\Theta}, \quad V = XW_V$$

$$\Theta_n = e^{in\theta}, \quad D_{nm} = \begin{cases} \gamma^{n-m}, & n \geq m \\ 0, & n < m \end{cases}$$

$$\text{Retention}(X) = (QK^T \odot D)V$$



(a) Parallel representation.

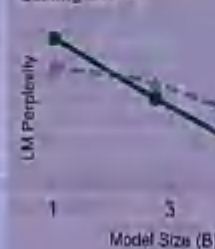


(b) Recurrent representation.

Inference Cost



Scaling Curve



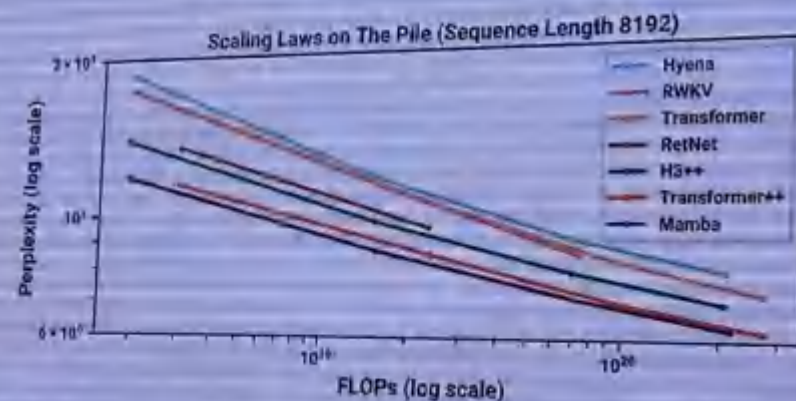
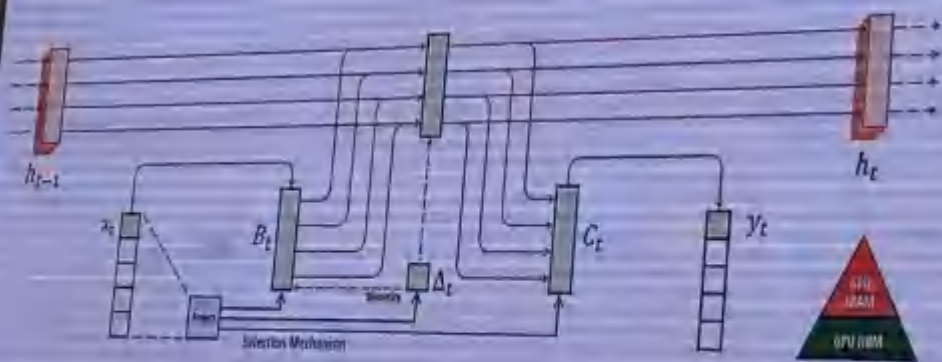
Transformer RetNet

Sun et al. Retentive Network: A Successor to Transformer for Large Language Models, In arXiv, 2023

VALSE 2024 APR: 面向大模型的新型高效率网络架构

NLP: Mamba

- 使用状态空间模型完成次二次方的序列建模任务
- 提出了 Mamba 模块：输入相关建模，次二次方复杂度，硬件感知算法
- 在 NLP 任务中取得了比 Transformer++ 更好的 efficiency 和 effectiveness



Gu et al. Mamba: Linear-time sequence modeling with selective state spaces, In arXiv, 2023

VALUE 2024 APR: 面向大模型的新型高效率网络架构

Transformer vs 新型高效网络架构

Architecture	Inference		Parallel	Training	
	Time	Memory		Time	Memory
LSTM/LMU	$O(1)$	$O(1)$	✗	$O(N)$	$O(N)$
Transformer	$O(N)$	$O(N)^a$	✓	$O(N^2)$	$O(N)^b$
Linear Transformer	$O(1)$	$O(1)$	✓	$O(N)$	$O(N)$
H3/S4	$O(1)$	$O(1)$	✓	$O(N \log N)$	$O(N)$
Hyena	$O(N)$	$O(N)$	✓	$O(N \log N)$	$O(N)$
RWKV/Mamba/RetNet	$O(1)$	$O(1)$	✓	$O(N)$	$O(N)$

^a $O(1)$ without KV cache; ^b With Flash Attention

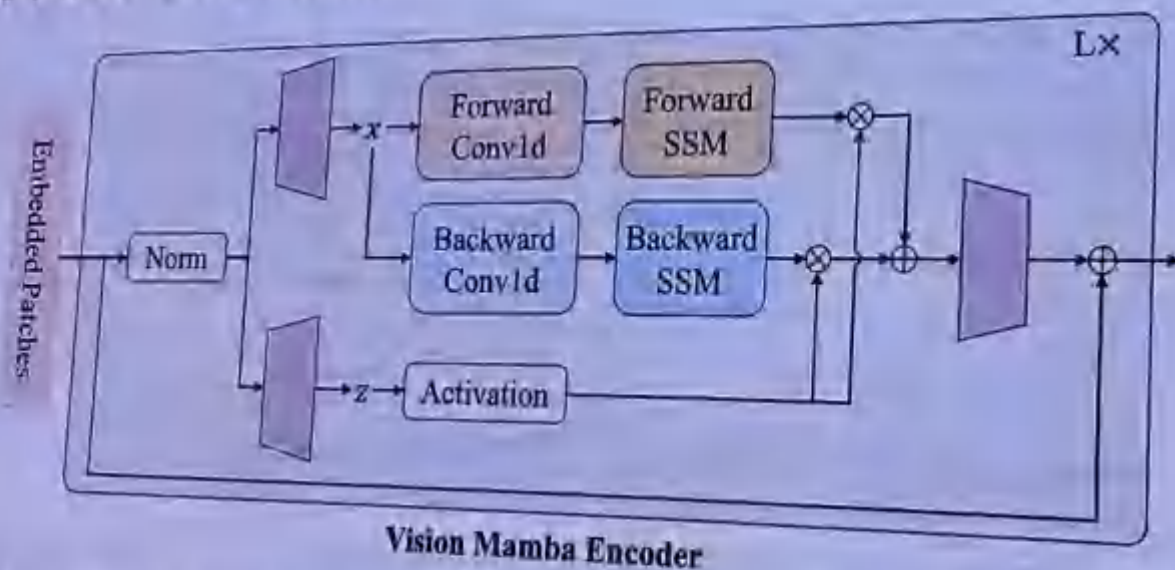
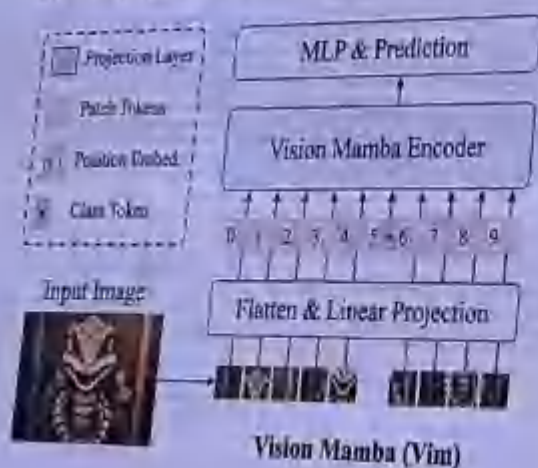
RWKV-5 & -6 中关于Transformer和其它新型高效网络架构之间的推理、训练复杂度对比

Peng, Bo, et al. "Eagle and Finch: RWKV with Matrix-Valued States and Dynamic Recurrence." arXiv:2404.05892 (2024)

VALSE 2024 APR: 面向大模型的新型高效率网络架构

计算机视觉: Vision Mamba

- 通过双向的结构设计使得 Mamba 能够拥有全局上下文感知
- 首次提出具有双向 Mamba 模块的通用视觉基础模型
- 在分类、检测、分割等视觉任务上获得了超越 ViT、DeiT 的结果, 并且再 1248 分辨率的大图上减少了 86% 的显存占用, 提速 2.8 倍

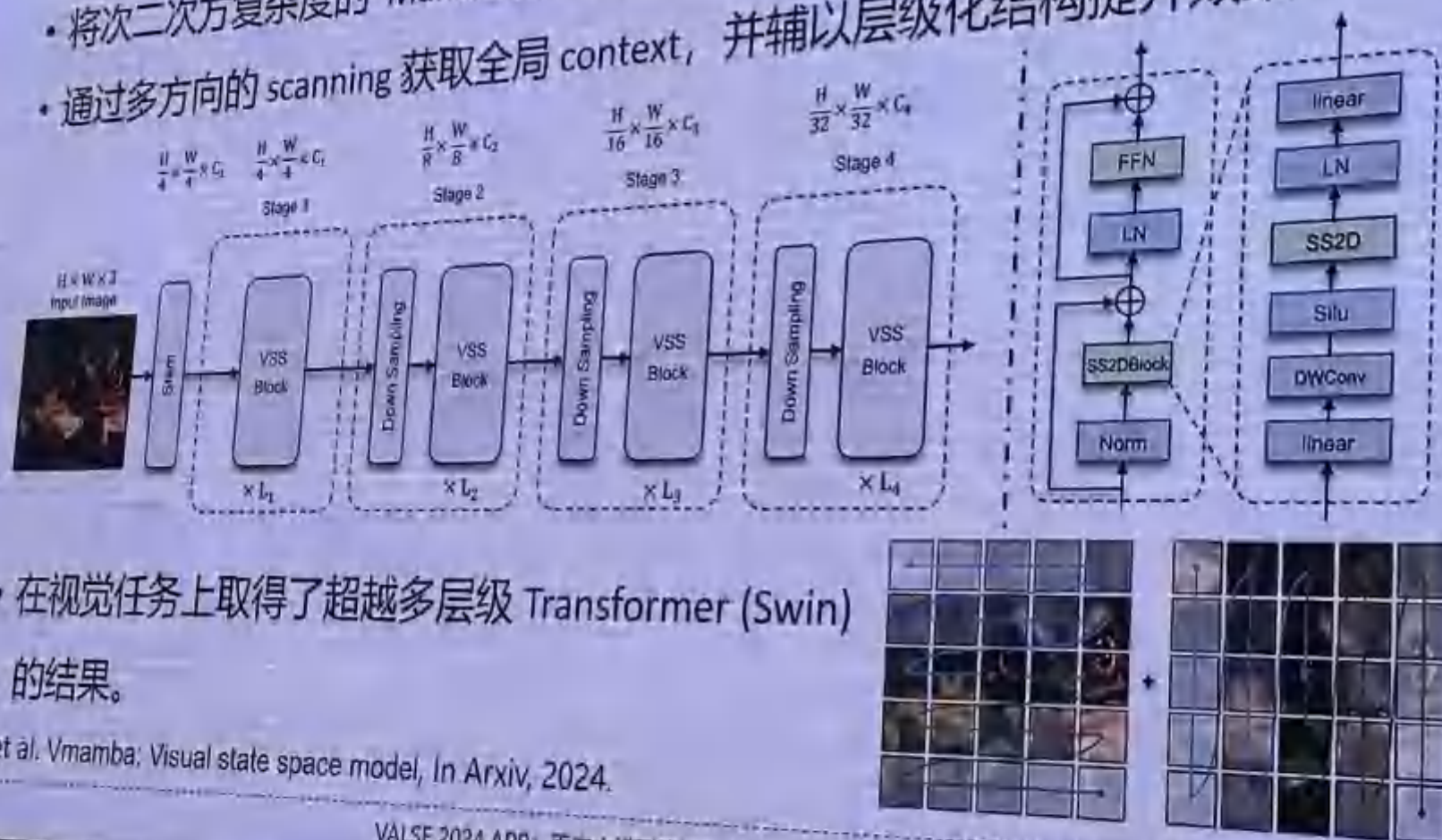


Zhu et al. Vision mamba: Efficient visual representation learning with bidirectional state space model, In ICML, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

计算机视觉: VMamba

- 将次二次方复杂度的 Mamba 方法引入视觉任务
- 通过多方向的 scanning 获取全局 context, 并辅以层级化结构提升效果



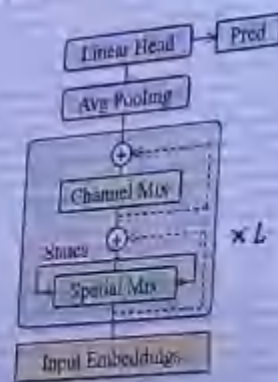
- 在视觉任务上取得了超越多层级 Transformer (Swin) 的结果。

Liu et al. Vmamba: Visual state space model, In Arxiv, 2024.

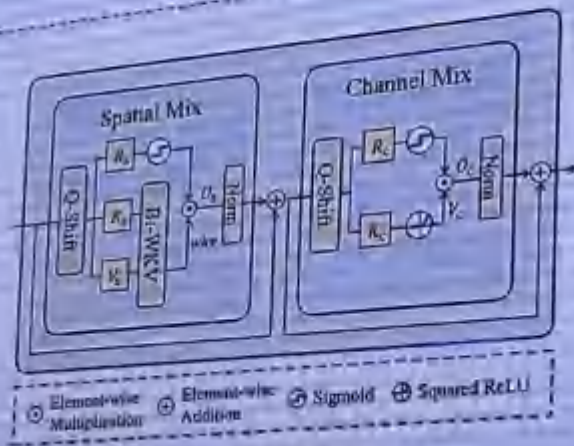
VALSE 2024 APR: 面向大模型的新型高效率网络架构

计算机视觉: Vision RWKV

- 将NLP里的线性算子RWKV引入到视觉领域，并证明了其具有良好的可规模化能力

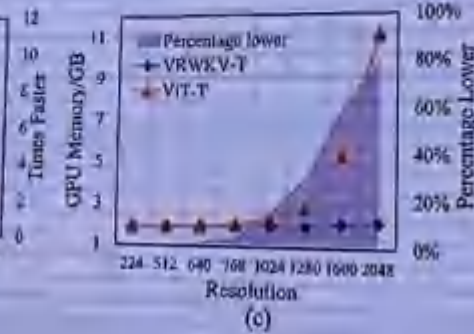
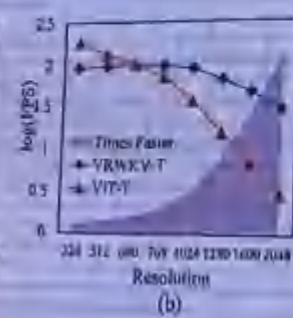
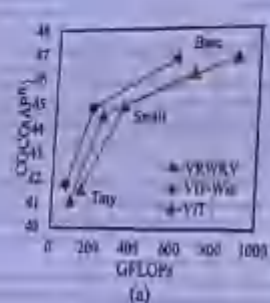


(a) Vision-RWKV Architecture



(b) Vision-RWKV Encoder Layer

- 引入双向扫描机制，和上下左右四个方向的Token shift机制，极大地增强了RWKV处理2D视觉信号的能力

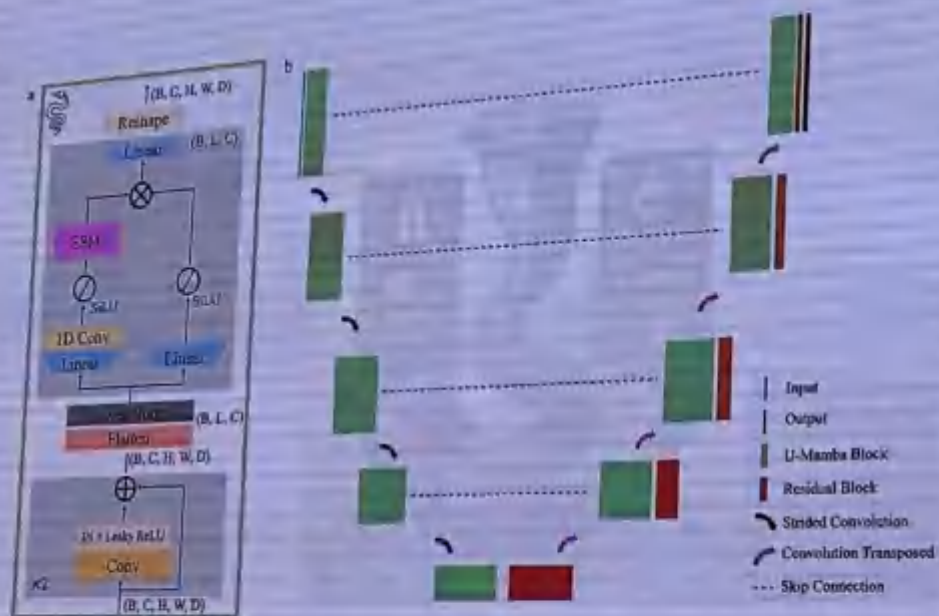


Duan et al. Vision-RWKV: Efficient and Scalable Visual Perception with RWKV-Like Architectures, In arXiv, 2024

VALUE 2024 APR: 面向大模型的新型高效率网络架构

医学图像: U-Mamba

- 首次将 Mamba 引入医学图像分割领域
- 通过 Mamba Block 和卷积结合的方式融合细粒度特征和长距离依赖
- 在 2D 和 3D 的多个医学图像分割任务上取得了更好的结果



Ma et al. U-mamba: Enhancing long-range dependency for biomedical image segmentation, In Arxiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

医学图像

- VM-Unet 和 SegMamba 将 Mamba 和 Unet 结合以处理不同的医学视觉任务
- 以上混合架构方法均在不同基准上取得了更优的结果



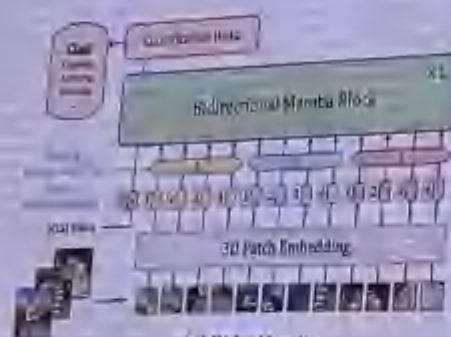
Ruan et al. VM-UNet: Vision Mamba UNet for Medical Image Segmentation, In Arxiv, 2024

Xing et al. SegMamba: Long-range Sequential Modeling Mamba For 3D Medical Image Segmentation, In Arxiv, 2024

VALSE 2024 APR. 面向大模型的新型高效率网络架构

视频理解

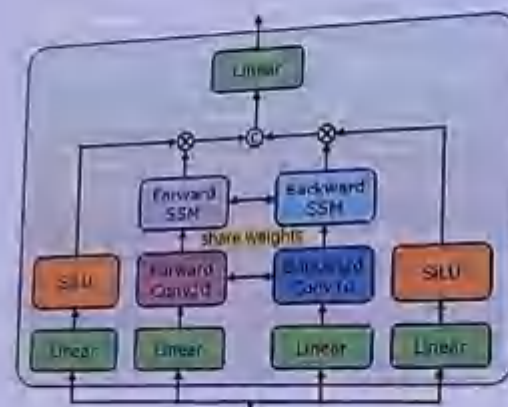
- VideoMamba / Video Mamba Suite 均使用了 Vim 构建视频理解网络
- VideoMamba 使用了多方向的 scanning 策略, Video Mamba Suite 使用了共享空间状态模块权重的设计以更好地利用方向偏置



(a) VideoMamba



(b) Spatiotemporal Scan



- 以上两方法均取得了包括视频理解、跨模态交互等多项任务的最好结果

Chen et al. Video Mamba Suite: State Space Model as a Versatile Alternative for Video Understanding, In Arxiv, 2024

Li et al. VideoMamba: State Space Model for Efficient Video Understanding, In Arxiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

点云理解: PointMamba

- 将NLP里的Mamba架构引入到3D点云处理
- 提出点云序列重排序方法, 使得Mamba中的SSM模块在处理点云序列时具有更为合理的扫描顺序

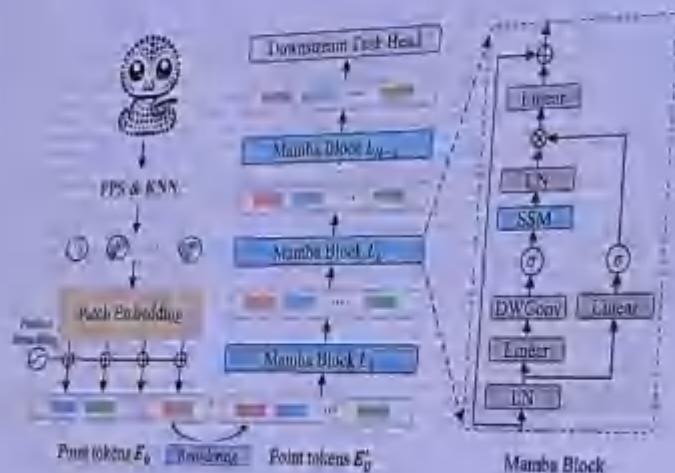


Table 3. Object classification on ScanObjectNN [30]. We evaluate PointMamba on three variants, with PB-T50-RS being the most challenging. Accuracy (%) is reported. * denotes reproduced results, † indicates that using simple rotational augmentation [9] for training.

Method	Backbone	Param. (M) ↓	FLOPs (G) ↓	OBJ-BG †	OBJ-ONLY †	PB-T50-RS †
PointNet [40]	-	3.5	0.5	73.3	79.2	68.0
PointNet++ [31]	-	1.3	1.7	82.3	84.3	77.9
PointCNN [26]	-	0.6	0.9	86.1	85.5	78.5
DGCNN [36]	-	1.8	2.4	82.8	86.2	76.1
PRNet [6]	-	-	-	-	-	82.8
MVTN [30]	-	11.2	43.7	-	-	87.7
PointNeX [43]	-	1.4	1.6	-	-	85.4
PointMLP [14]	-	13.2	31.4	-	-	84.3
RepStair-U [46]	-	1.5	0.8	-	-	87.5
ADS [33]	-	-	-	-	-	-
Training from scratch						
Transformer [64]	Transformer-based	22.1	4.8	79.86	80.55	77.24
PointMAE* [37]	Transformer-based	22.1	4.8	86.75	86.92	80.78
PointMamba (ours)	Mamba-based	12.3	3.6	88.30	87.78	82.48
Training from pre-training						
Point-BERT [64]	Transformer-based	22.1	4.8	87.43	88.12	83.07
PointMAE [37]	Transformer-based	22.1	4.8	90.02	88.29	85.18
PointMAE† [37]	Transformer-based	22.1	4.8	92.77	91.22	89.04
PointMamba (vars)	Mamba-based	12.3	3.6	90.71	88.47	84.87
PointMamba† (vars)	Mamba-based	12.3	3.6	93.29	91.91	88.17

Liang et al. PointMamba: A Simple State Space Model for Point Cloud Analysis, In arXiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

底层视觉: MambaIR

- 将NLP中的Mamba引入到图像复原任务
- 设计了残差空间状态模块, 通过额外引入卷积和channel attention来缓解Mamba在处理low-level视觉中出现的局部像素遗忘和通道冗余问题

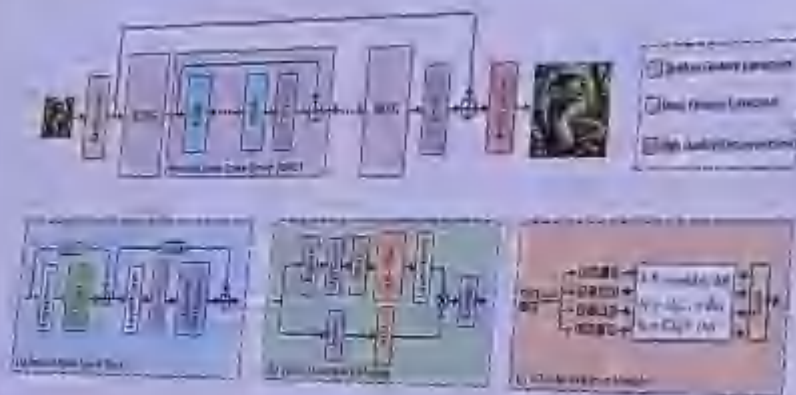


Table 4: Quantitative comparison on classic image super-resolution with state-of-the-art methods. The best and the second best results are in red and blue.

Model	Scale	Set5		Set11		BSDS100		Urban100		Manga109	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
EDSR [27]	x2	40.13	0.9607	38.82	0.9194	41.32	0.9613	38.93	0.9351	38.10	0.9773
RCAN [7]	x2	39.27	0.9612	36.12	0.9214	41.41	0.9617	38.24	0.9364	39.44	0.9780
SAN [3]	x2	39.51	0.9620	36.07	0.9213	41.42	0.9620	38.19	0.9370	39.32	0.9792
HAN [34]	x2	39.27	0.9614	36.10	0.9217	41.41	0.9623	38.35	0.9385	39.46	0.9785
IGRN [78]	x2	39.24	0.9619	36.07	0.9217	41.40	0.9624	38.15	0.9386	39.37	0.9789
CHSNet [40]	x2	39.29	0.9619	36.13	0.9228	41.43	0.9627	38.42	0.9394	39.88	0.9789
NLSA [47]	x2	39.34	0.9618	36.05	0.9211	41.43	0.9630	38.44	0.9391	39.82	0.9793
ELAN [79]	x2	39.36	0.9630	36.10	0.9228	41.45	0.9630	38.44	0.9391	39.82	0.9793
UT [4]	x2	39.37	-	36.43	-	41.45	-	38.76	-	-	-
SwiSR [38]	x2	39.42	0.9633	36.46	0.9280	41.53	0.9641	38.81	0.9427	39.92	0.9797
SRFormer [39]	x2	39.51	0.9637	36.44	0.9283	41.57	0.9648	38.93	0.9440	40.07	0.9802
MambaIR	x2	39.57	0.9637	36.47	0.9289	41.58	0.9648	38.95	0.9446	40.26	0.9806
MambaIR+	x2	39.63	0.9639	36.48	0.9289	41.60	0.9648	38.97	0.9448	40.28	0.9806
EDSR [27]	x4	41.85	0.9930	39.89	0.9462	43.25	0.9930	40.85	0.9653	34.17	0.9476
RCAN [7]	x4	41.78	0.9928	39.85	0.9462	43.22	0.9931	40.89	0.9659	34.04	0.9480
SAN [3]	x4	41.75	0.9930	39.83	0.9465	43.23	0.9932	40.93	0.9671	34.30	0.9484
HAN [34]	x4	41.76	0.9934	39.87	0.9469	43.24	0.9933	40.94	0.9675	34.45	0.9490
IGRN [78]	x4	41.72	0.9933	39.86	0.9464	43.21	0.9935	40.92	0.9666	34.39	0.9486
CHSNet [40]	x4	41.73	0.9930	39.88	0.9467	43.23	0.9935	40.93	0.9672	34.43	0.9492
NLSA [47]	x4	41.85	0.9935	39.79	0.9465	43.34	0.9937	40.95	0.9678	34.57	0.9500
ELAN [79]	x4	41.88	0.9933	39.80	0.9468	43.35	0.9934	40.97	0.9675	34.79	0.9517
UT [4]	x4	41.82	-	40.83	-	43.26	-	40.48	-	-	-
SwiSR [38]	x4	41.87	0.9935	39.95	0.9504	43.45	0.9935	40.95	0.9678	35.12	0.9507
SRFormer [39]	x4	41.92	0.9938	39.94	0.9508	43.43	0.9936	40.94	0.9683	35.20	0.9503
MambaIR	x4	41.98	0.9938	39.99	0.9516	43.51	0.9937	40.93	0.9681	35.43	0.9510
MambaIR+	x4	42.12	0.9939	40.15	0.9541	43.63	0.9938	40.98	0.9688	35.55	0.9549
EDSR [27]	x8	42.48	0.9968	40.80	0.9789	43.71	0.9938	41.64	0.9813	31.62	0.9148
RCAN [7]	x8	42.62	0.9968	40.87	0.9789	43.77	0.9938	41.63	0.9807	31.52	0.9173
SAN [3]	x8	42.64	0.9970	40.92	0.9796	43.75	0.9936	41.79	0.9809	31.18	0.9169
HAN [34]	x8	42.68	0.9970	40.93	0.9790	43.80	0.9942	41.85	0.9814	31.43	0.9177
IGRN [78]	x8	42.57	0.9964	40.89	0.9789	43.77	0.9934	41.84	0.9800	31.38	0.9182
CHSNet [40]	x8	42.68	0.9974	40.95	0.9788	43.80	0.9940	41.82	0.9810	31.43	0.9184
NLSA [47]	x8	42.59	0.9969	40.87	0.9784	43.78	0.9944	41.84	0.9810	31.57	0.9184
ELAN [79]	x8	42.75	0.9972	40.98	0.9784	43.83	0.9945	41.83	0.9817	31.68	0.9226
UT [4]	x8	42.64	-	40.81	-	43.82	-	41.79	-	-	-
SwiSR [38]	x8	42.92	0.9974	40.89	0.9790	43.92	0.9949	41.85	0.9814	31.68	0.9226
SRFormer [39]	x8	42.93	0.9974	40.95	0.9793	43.94	0.9952	41.86	0.9811	31.71	0.9271
MambaIR	x8	43.03	0.9976	41.03	0.9801	43.98	0.9953	41.87	0.9817	31.82	0.9272
MambaIR+	x8	43.12	0.9978	41.13	0.9811	44.01	0.9953	41.88	0.9818	31.85	0.9281

Guo et al. MambaIR: A Simple Baseline for Image Restoration with State-Space Model, In arXiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

遥感图像: RS-Mamba

- 将Mamba应用到遥感图像分析领域, 用于处理非常高分辨率 (VHR) 遥感图像的密集预测任务
- 设计了全方向扫描机制, 在水平、垂直、对角和反对角方向上进行选择性扫描, 增强了在多个方向上对上下文信息的全局建模能力

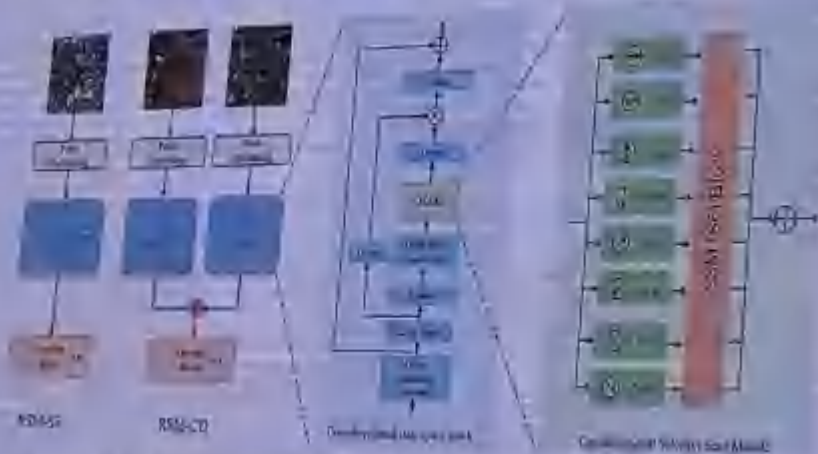


Fig. 1. Overview of the Overall Structure of RSM-SS and RSM-CD. RSM-SS and RSM-CD can globally model the context of images in multiple directions with dense connectivity using the reconfigurable selective scan.

TABLE III
ACCURACY COMPARISON ON THE MASSACHUSETTS ROAD DATASET. THE BEST VALUES ARE HIGHLIGHTED IN BOLD.

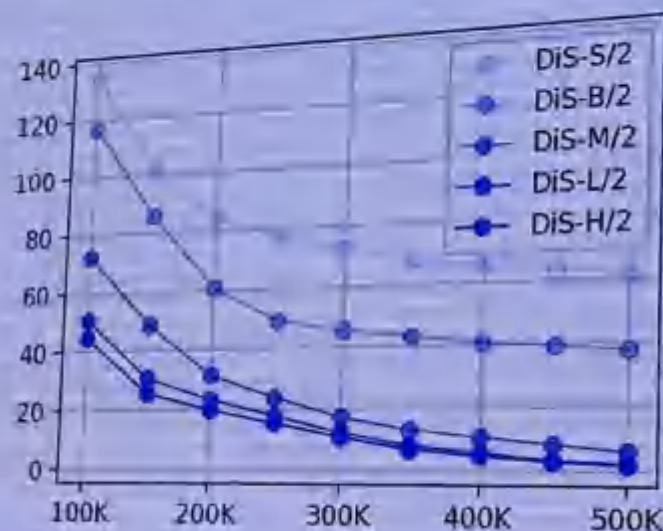
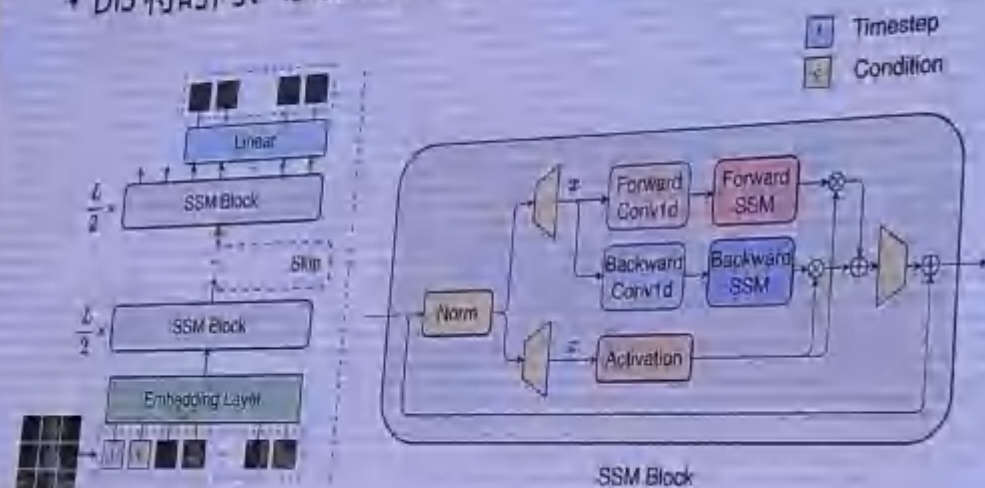
Methods	P (%)	R (%)	F1 (%)	IoU (%)
SegNet [44]	76.09	78.23	77.15	62.79
U-Net [17]	77.53	77.82	77.67	63.50
ResU-Net [23]	78.77	77.45	78.10	64.07
D-LinkNet [51]	78.34	77.91	78.12	64.10
HRNetv2 [46]	79.01	78.20	78.60	64.75
DeepLabv3+ [48]	75.14	72.56	73.83	58.51
SiTNet [52]	85.36	74.13	79.35	65.77
RoadFormer [8]	80.54	78.90	79.71	66.27
BDTNet [54]	82.99	76.37	79.54	66.03
RSM-SS	86.52	75.24	80.49	67.35

Chen et al. RSMamba: Remote Sensing Image Classification with State Space Model, In arXiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

扩散模型: DiS

- DiS 首次将 Vim 模块引入 Diffusion 生成领域
- DiS 将时间、条件和噪声图像块视作 token, 并通过跳跃连接贯通浅层和深层



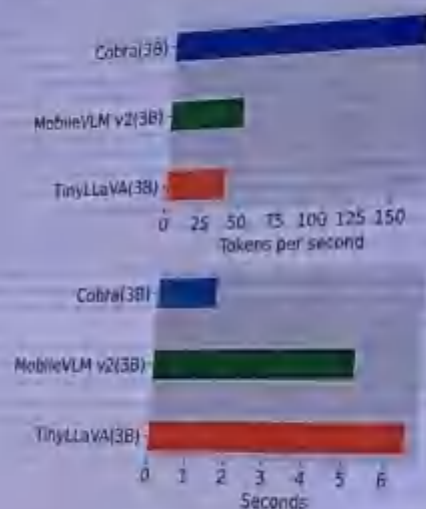
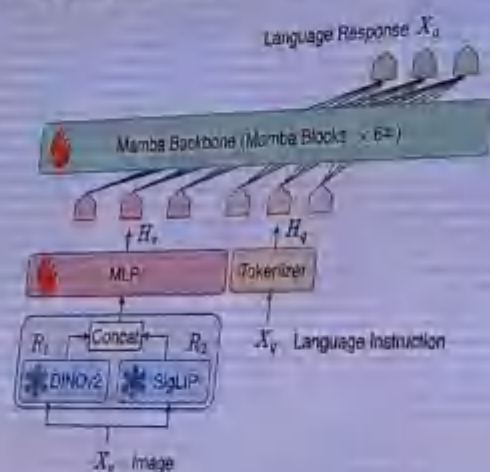
- DiS 展现了较强的可扩展性并在多个图片生成基准上展现了卓越的结果

Fei et al. Scalable Diffusion Models with State Space Backbone, In Arxiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

多模态模型: Cobra

- Cobra 首次将 Mamba 引入多模态大语言模型领域
- Cobra 通过线性层将视觉基础模型的输出映射为 Mamba 可接受的向量



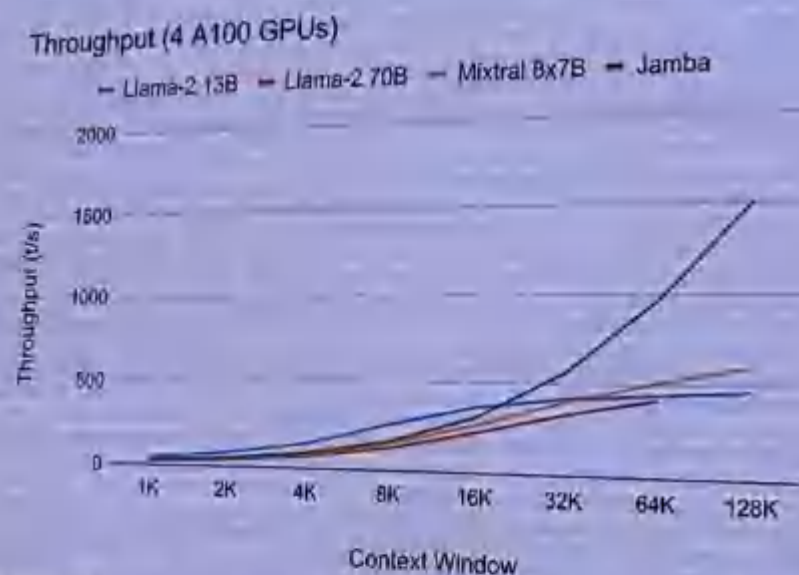
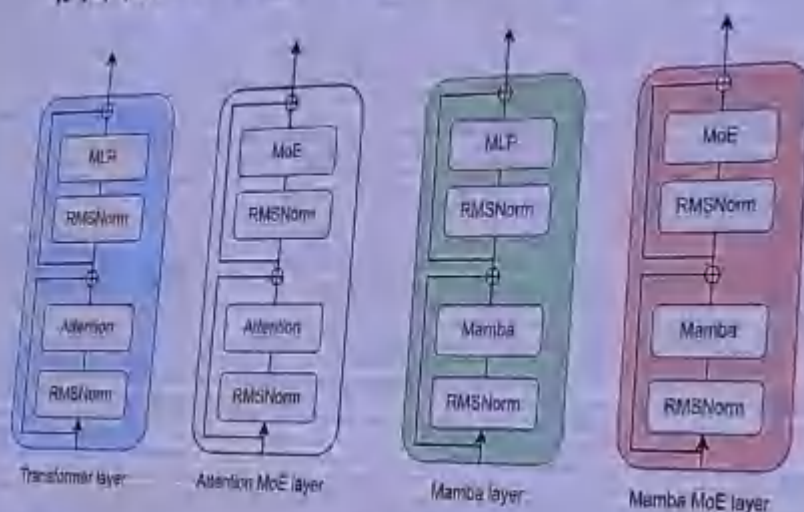
- Cobra 在多个基准取得了更好的结果, 以及更高的效率

Zhao et al. Cobra: Extending Mamba to Multi-Modal Large Language Model for Efficient Inference, In Arxiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

混合架构: Jamba

- Jamba 希望同时结合 Transformer (高效用) 和 Mamba (高效率) 的优点
- 通过交错排列Transformer和Mamba层的块来实现, 同时在某些层中加入MoE以增加模型容量, 并控制被激活的参数数量
- 仅采用1个80G现存的GPU, 实现了大规模数据集上的训练



Lieber et al. Jamba: A Hybrid Transformer-Mamba Language Model, In Arxiv, 2024

VALSE 2024 APR: 面向大模型的新型高效率网络架构

总结和展望

- 长序列的表征学习在Sora视频生成、多模态大模型领域非常关键，探索类似Mamba的低复杂度序列建模网络是一个非常重要的问题。
- Mamba等模型尚存在的挑战：（1）缺乏超大规模预训练上的验证；（2）缺乏灵活的多模态自监督训练方法；（3）串行模型的并行训练与硬件实现息息相关，更高效率的并行训练方法需要进一步挖掘。
- 值得进一步探索的方向：（1）混合架构，有机的结合transformer、mamba等方法的优点；（2）算子优化，新型架构和异构硬件的协同优化；（3）视频生成：充分发挥新型架构在长序列建模方面的低复杂度优势；（4）多模态建模：构建类似于fuyu模型那种原生的多模态模型。