

# VMamba: Visual State Space Model

---

문학준

gloriel621@g.skku.edu

**TNT Vision**

2024/4/30



TRAIN AND TEST

# Contents

---

1. Introduction + Related Work
2. Preliminaries (SSM, S4, Mamba)
3. Cross-Scan Module(CSM)
4. 2D-Selective-Scan (SS2D)
5. Overall Architecture
6. Experiments

# Introduction

---

## 1. ViT의 단점

- a. Self Attention의 Complexity

## 2. VMamba의 contribution

- a. SSM에 기반한 novel network architecture 제시.
- b. Cross-Scan Module (CSM) 제시 -> bridge the gap between 1d array scanning and 2d plain traversing
- c. Classification, Object Detection, Semantic Segmentation에서 좋은 성능을 보임.

# Preliminaries

---

- SSM(State Space Models)

$$h'(t) = Ah(t) + Bx(t) \quad (1a)$$

$$y(t) = Ch(t) \quad (1b)$$

- Discretization of SSM (0|산화) :  $(\Delta, A, B, C) \mapsto (\bar{A}, \bar{B}, C)$

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t \quad (2a) \quad \bar{A} = \exp(\Delta A) \quad \bar{B} = (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta B$$

$$y_t = Ch_t \quad (2b)$$

- Kernel을 Precomputing 하여 값 계산, Convolution 가능.

$$\bar{K} = (C\bar{B}, C\bar{A}\bar{B}, \dots, C\bar{A}^k\bar{B}, \dots) \quad (3a)$$

$$y = x * \bar{K} \quad (3b)$$

# Preliminaries

---

- Linear Time Invariant(선형 시간 불변성)

$$\bar{A} = \exp(\Delta A) \quad \bar{B} = (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta B \quad (\bar{A}, \bar{B}, C \text{는 모두 고정적인 값})$$

(변화량  $\Delta$ 는 항상 일정한 값)

$$h_t = \bar{A}h_{t-1} + \bar{B}x_t \quad (2a)$$

$$y_t = Ch_t \quad (2b)$$

# Preliminaries

- Structured State Space Model (S4)

- a. A의 값을 어떻게 정의하는 것이 좋은가?
  - b. 연산 효율을 어떻게 높일 수 있는가?
- 분해 가능하고, 좋은 성능을 보이는 A를 사용하여 Convolution Kernel  $\bar{K}$  를 구한다.

- 1. S4에서 Random A보다 HiPPO A를 사용했을 때 더 나은 성능을 보임.

$$\text{(HiPPO Matrix)} \quad A_{nk} = - \begin{cases} (2n+1)^{1/2}(2k+1)^{1/2} & \text{if } n > k \\ n+1 & \text{if } n = k \\ 0 & \text{if } n < k \end{cases}$$

- 2. 이 HiPPO A는 Normal Plus Low Rank(NPLR) 로 decompose 할 수 있다.  
(= A를 정규 행렬과 low-rank 행렬의 합의 형태로 표현할 수 있다.)  
V = unitary matrix,  $\Lambda$  = diagonal matrix, P, Q = low-rank matrix

$$A = V\Lambda V^* - PQ^T = V(\Lambda - (V^*P)(V^*Q)^*)V^*$$

# Preliminaries

- Structured State Space Model (S4)
  - a. A의 값을 어떻게 정의하는 것이 좋은가?
  - b. 연산 효율을 어떻게 높일 수 있는가?
- 3.  $\bar{A}, \bar{B}$  를 A, B에 대한 식으로 나타낼 수 있다.  
그리고 iFFT를 통해  $\bar{K}$ 도 구할 수 있다.

$$\bar{A} = A_1 A_0$$

$$\bar{B} = \frac{2}{\Delta} A_1 \Delta B = 2A_1 B.$$

$$\begin{aligned} x_k &= \bar{A}x_{k-1} + \bar{B}u_k \\ &= A_1 A_0 x_{k-1} + 2A_1 B u_k \\ y_k &= Cx_k. \end{aligned}$$

---

## Algorithm 1 S4 CONVOLUTION KERNEL (SKETCH)

---

**Input:** S4 parameters  $\Lambda, P, Q, B, C \in \mathbb{C}^N$  and step size  $\Delta$

**Output:** SSM convolution kernel  $\bar{K} = \mathcal{K}_L(\bar{A}, \bar{B}, \bar{C})$  for  $A = \Lambda - PQ^*$  (equation (5))

- 1:  $\tilde{C} \leftarrow (I - \bar{A}^L)^* \bar{C}$  ▷ Truncate SSM generating function (SSMGF) to length  $L$
  - 2:  $\begin{bmatrix} k_{00}(\omega) & k_{01}(\omega) \\ k_{10}(\omega) & k_{11}(\omega) \end{bmatrix} \leftarrow [\tilde{C} Q]^* \left( \frac{2}{\Delta} \frac{1-\omega}{1+\omega} - \Lambda \right)^{-1} [B P]$  ▷ Black-box Cauchy kernel
  - 3:  $\hat{K}(\omega) \leftarrow \frac{2}{1+\omega} [k_{00}(\omega) - k_{01}(\omega)(1 + k_{11}(\omega))^{-1}k_{10}(\omega)]$  ▷ Woodbury Identity
  - 4:  $\hat{K} = \{\hat{K}(\omega) : \omega = \exp(2\pi i \frac{k}{L})\}$  ▷ Evaluate SSMGF at all roots of unity  $\omega \in \Omega_L$
  - 5:  $\bar{K} \leftarrow \text{iFFT}(\hat{K})$  ▷ Inverse Fourier Transform
-

# Preliminaries

---

- Mamba(S6)
  - a. Linear-Time Invariant -> Time Variant Model 로 변환.
    - LTI 모델(S4)의 한계 : Selective copying, Inductive head task를 잘 수행하지 못함. (Selectivity가 부족)
    - 따라서, 정보를 선택적으로 저장해야 할 필요가 있음.
    - S4에서는 입력  $x$ 에 상관없이,  $A$ ,  $B$ ,  $C$ 는 항상 같은 정적인 특성을 가진다.
    - Mamba는 입력에 따라  $B$ ,  $C$ ,  $\Delta$ 가 달라진다. ( $A$ 는 여전히 고정된 값)
    - $\Delta$ 가 작을 경우, 입력의 비중이 낮아지고, 이전 맥락을 사용하는 데에 초점을 맞춘다.



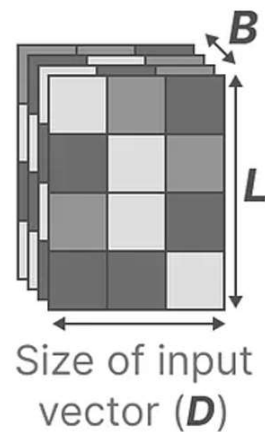
# Preliminaries

- Mamba(S6)

## Step size ( $\Delta$ )

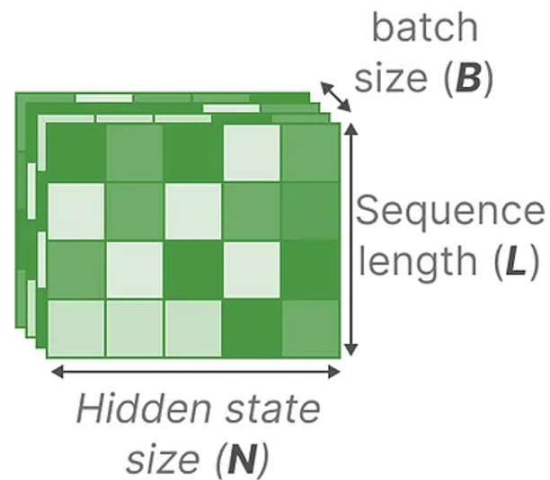
Resolution of the **input**  
(discretization parameter)

SSM +  
Selection



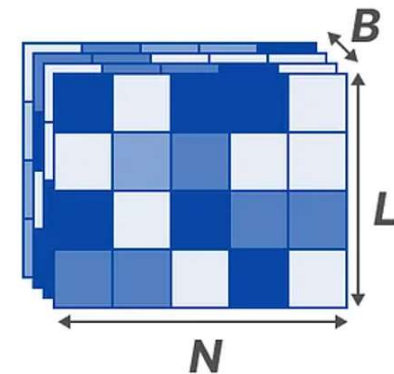
## Matrix **B**

**How** the **input**  
influences the state



## Matrix **C**

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



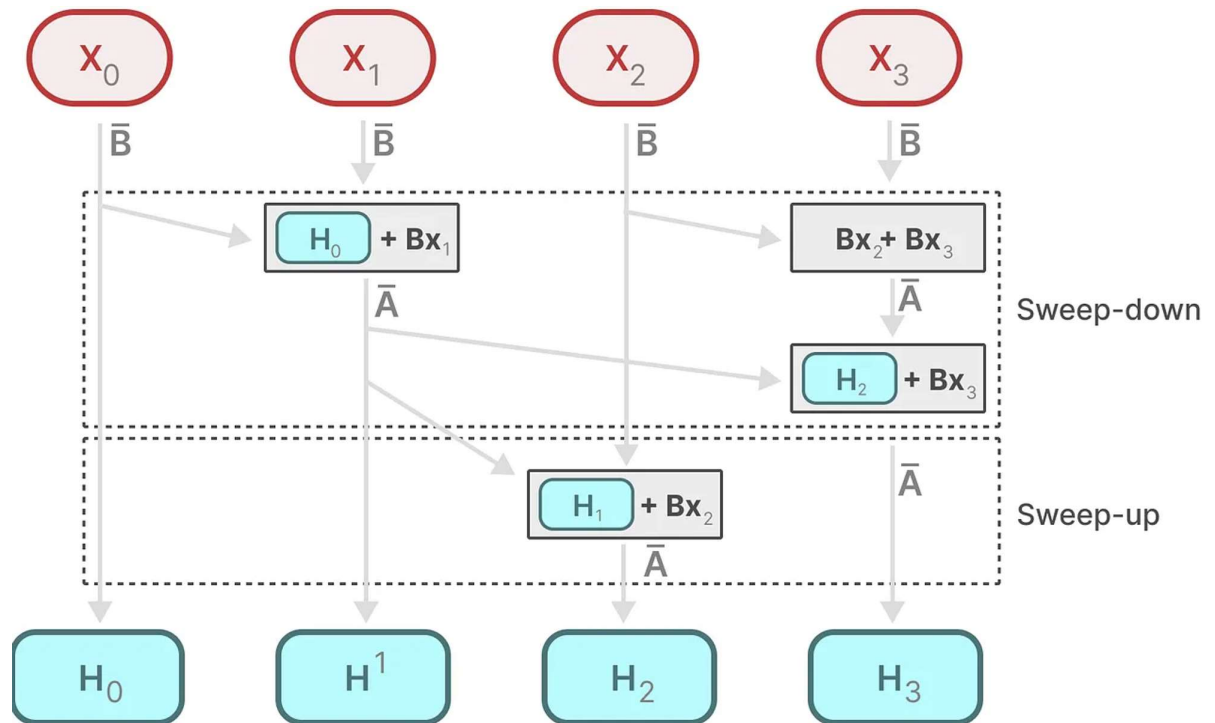
# Preliminaries

---

- Mamba(S6)
  - b. Selective Scan
    - 입력에 따라  $B, C, \Delta$  가 달라지기 때문에, Kernel을 미리 계산해 둘 수 없다. (Convolution 사용 불가)
    - Parallel Associative Scan, Kernel Fusion, Recomputation을 사용하여 해당 단점을 극복.

# Preliminaries

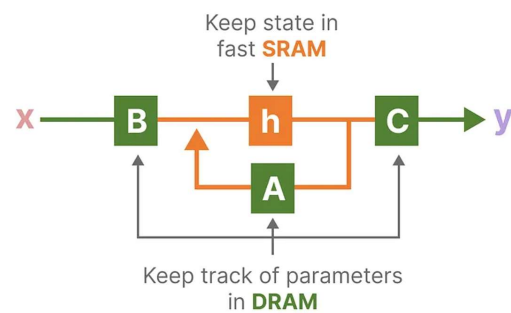
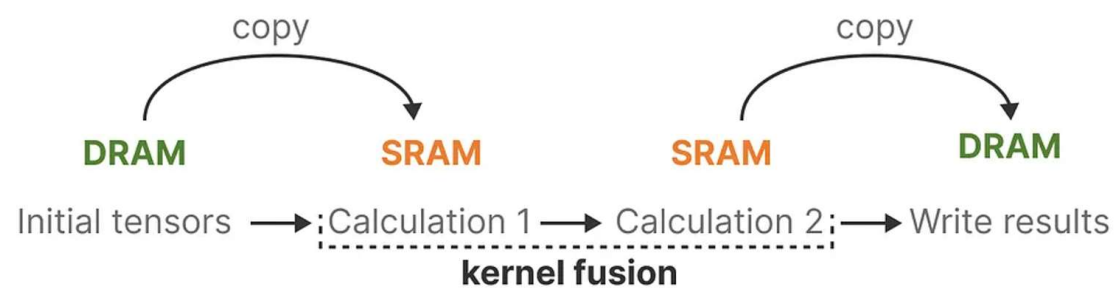
- Mamba(S6)
  - Parallel Associative Scan  
(Convolution 대신 사용)



Parallel computation  $O(n/t)$

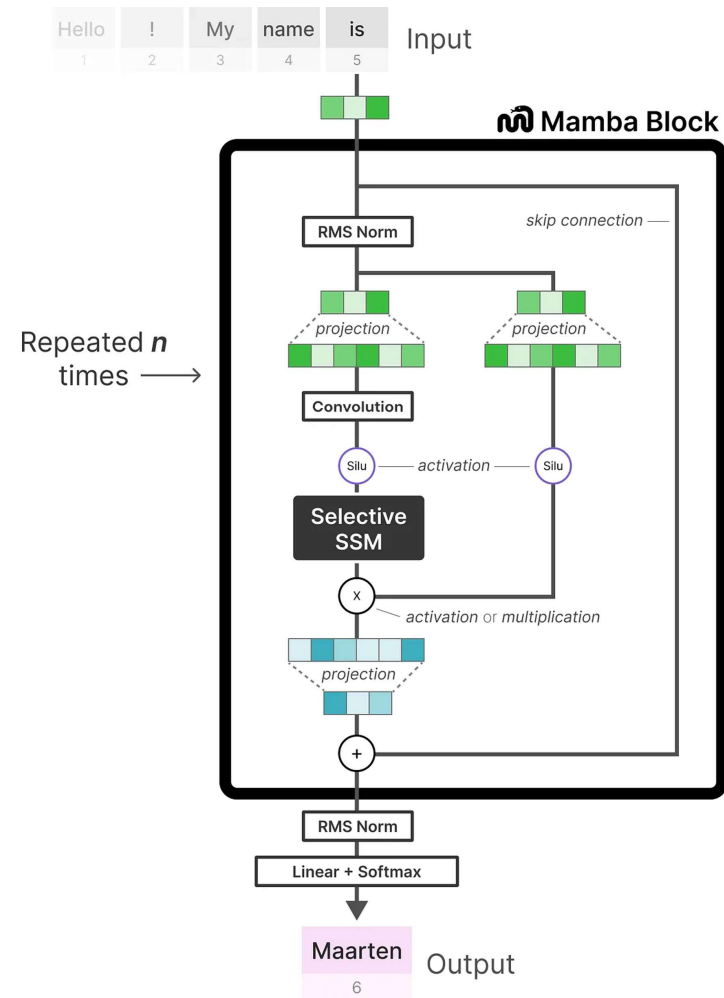
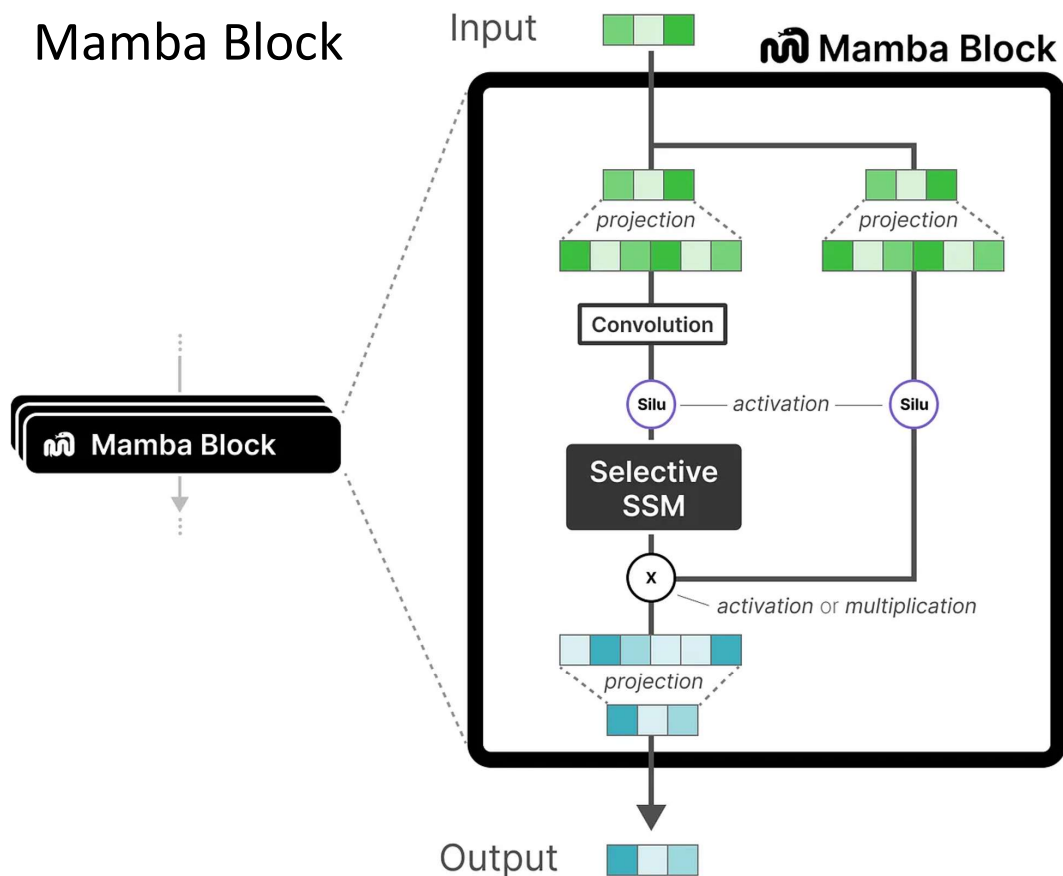
# Preliminaries

- Mamba(S6)
  - Kernel Fusion & Recomputation



# Preliminaries

- Mamba Block

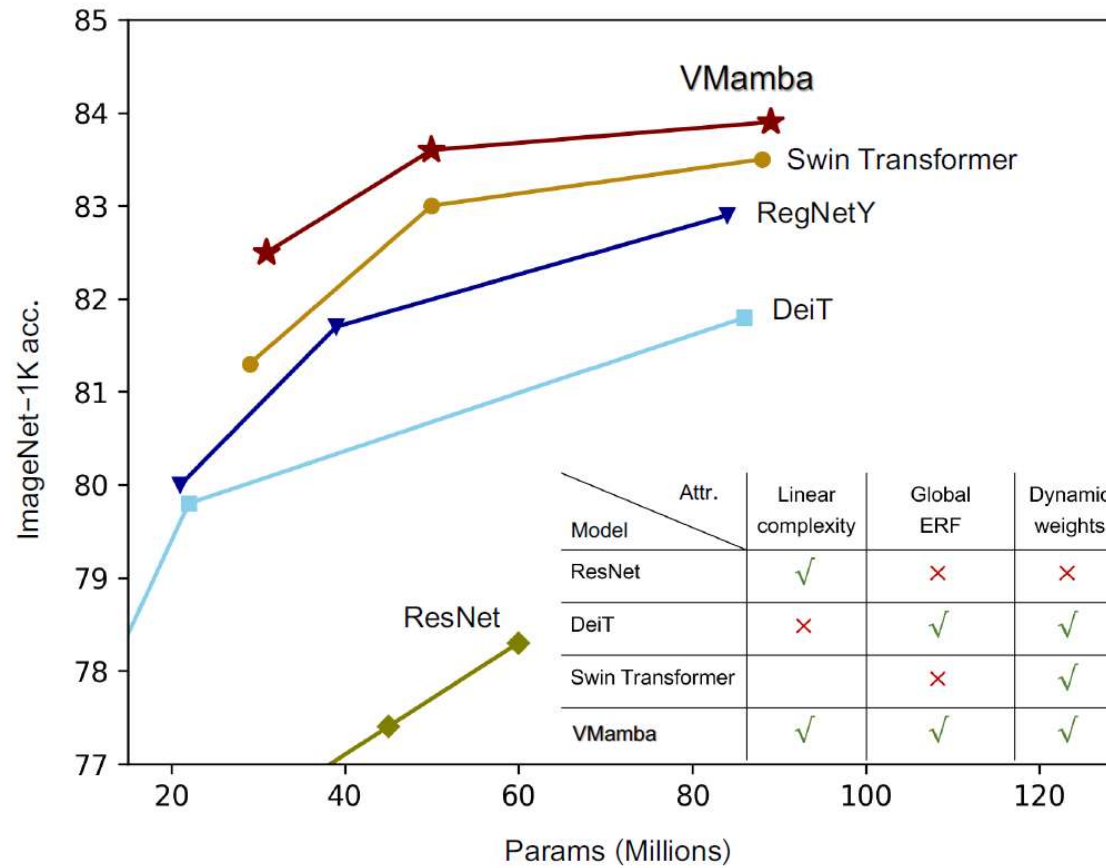


---

# VMamba: Visual State Space Model

# VMamba Performance

- ImageNet 1K



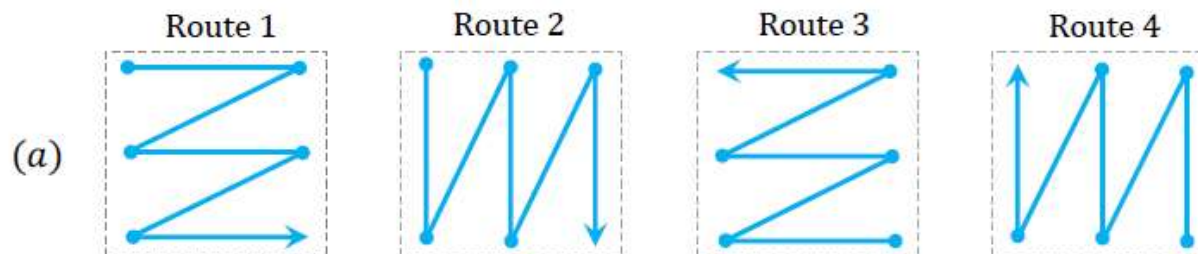
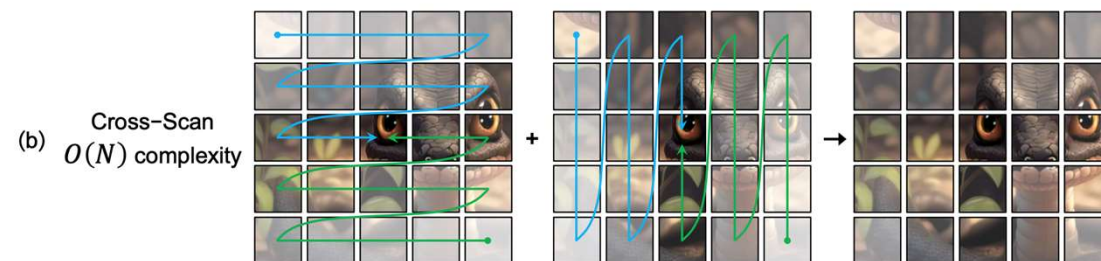
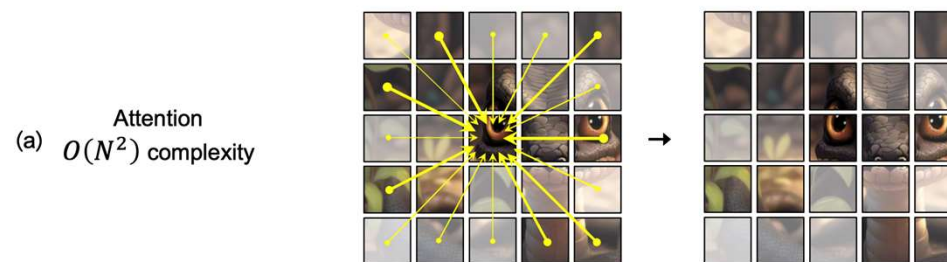
# Cross-Scan Module(CSM)

S4ND: S4를 Vision Task에 적용시킨 논문

- Convolution 사용
- Weight이 input independent.
- Limited capacity.

Attention VS Cross-Scan Module

- 한 patch를 다른 patch들과 모두 attention 함.  $O(N^2)$
- CSM은 기존의 SSM 기반 방식과 달리 4방향의 Cross Scan을 적용하여 처리함.





# 2D-Selective-Scan (SS2D)

1. Cross-Scan
2. Parallel S6 block processing
3. Cross-Merge

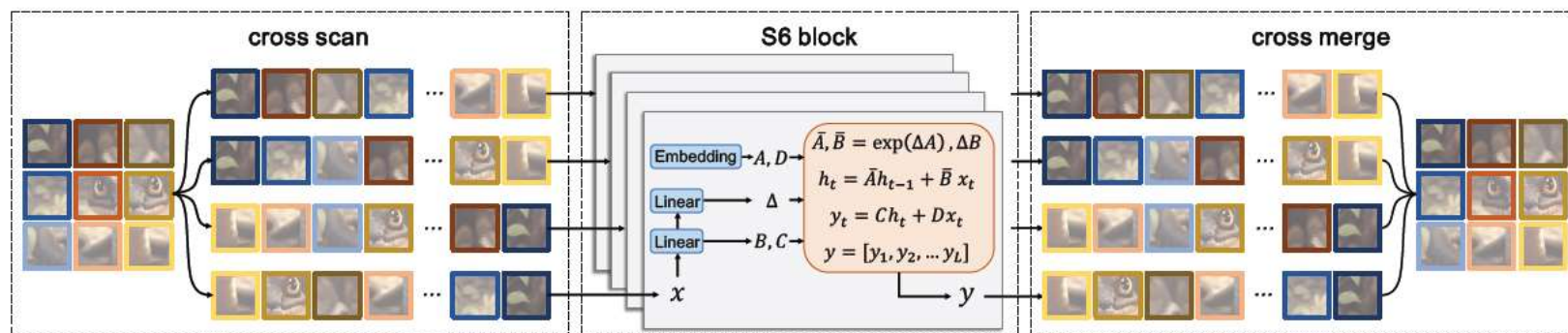


Figure 3: **Illustration of the 2D-Selective-Scan (SS2D) operation.** The input image is first divided into patches and then flattened along four different scanning paths (*Cross-Scan*). The resulting sequences of image patches are then separately processed by distinct S6 blocks. Subsequently, the outputs are merged to construct the 2D feature map as the final output (*Cross-Merge*).

# The Relationship between SS2D and Self-attention

In this section, we clarify the relationship between SS2D and the self-attention mechanism commonly employed in existing backbone models. Subsequently, visualization results are provided to substantiate our explanation.

**Notation Definition.** Let  $T$  denote the length of the sequence with indices from  $a$  to  $b$ , we define the following variables

$$V := [V_1; \dots; V_T] \in \mathbb{R}^{T \times D_v}, \text{ where } V_i := u_{a+i-1} \Delta_{a+i-1} \in \mathbb{R}^{1 \times D_v} \quad (10)$$

$$K := [K_1; \dots; K_T] \in \mathbb{R}^{T \times D_k}, \text{ where } K_i := B_{a+i-1} \in \mathbb{R}^{1 \times D_k} \quad (11)$$

$$Q := [Q_1; \dots; Q_T] \in \mathbb{R}^{T \times D_k}, \text{ where } Q_i := C_{a+i-1} \in \mathbb{R}^{1 \times D_k} \quad (12)$$

$$w := [w_1; \dots; w_T] \in \mathbb{R}^{T \times D_k \times D_v}, \text{ where } w_i := \prod_{j=1}^i e^{A \Delta_{a-1+j}} \in \mathbb{R}^{D_k \times D_v} \quad (13)$$

$$H := [h_a; \dots; h_b] \in \mathbb{R}^{T \times D_k \times D_v}, \text{ where } h_i \in \mathbb{R}^{D_k \times D_v} \quad (14)$$

$$Y := [y_a; \dots; y_b] \in \mathbb{R}^{T \times D_v}, \text{ where } y_i \in \mathbb{R}^{D_v} \quad (15)$$

# The Relationship between SS2D and Self-attention

**Mathematical Derivation.** Based on these notations, the discretized solution to time-varying SSMs (Eq. 3) can be written as

$$h_b = w_T \odot h_a + \sum_{i=1}^T \frac{w_T}{w_i} \odot (K_i^\top V_i) \quad (16)$$

where  $\odot$  denotes the element-wise product between matrices, and the division is also elements-wise. Therefore, the first term of the output of SSM can be computed by

$$y_b = Q_T h_b \quad (17)$$

$$= Q_T (w_T \odot h_a) + Q_T \sum_{i=1}^T \frac{w_T}{w_i} \odot (K_i^\top V_i) \quad (18)$$

Therefore, the  $j$ -th slice along dimension  $D_v$  of  $Y$ , denoted as  $X := Y^{(j)} \in \mathbb{R}^{T \times 1}$ , can be expressed as

$$X = (Q_T \odot w^{(j)}) h_a^{(j)} + \left[ \left( Q_T \odot w^{(j)} \right) \left( \frac{K}{w^{(j)}} \right)^\top \odot M \right] V^{(j)} \quad (19)$$

# Activation Map

---

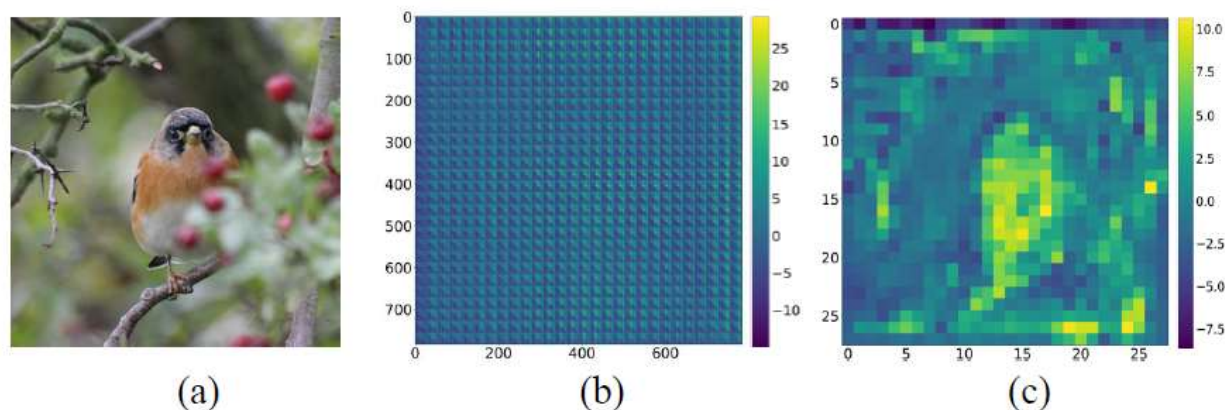
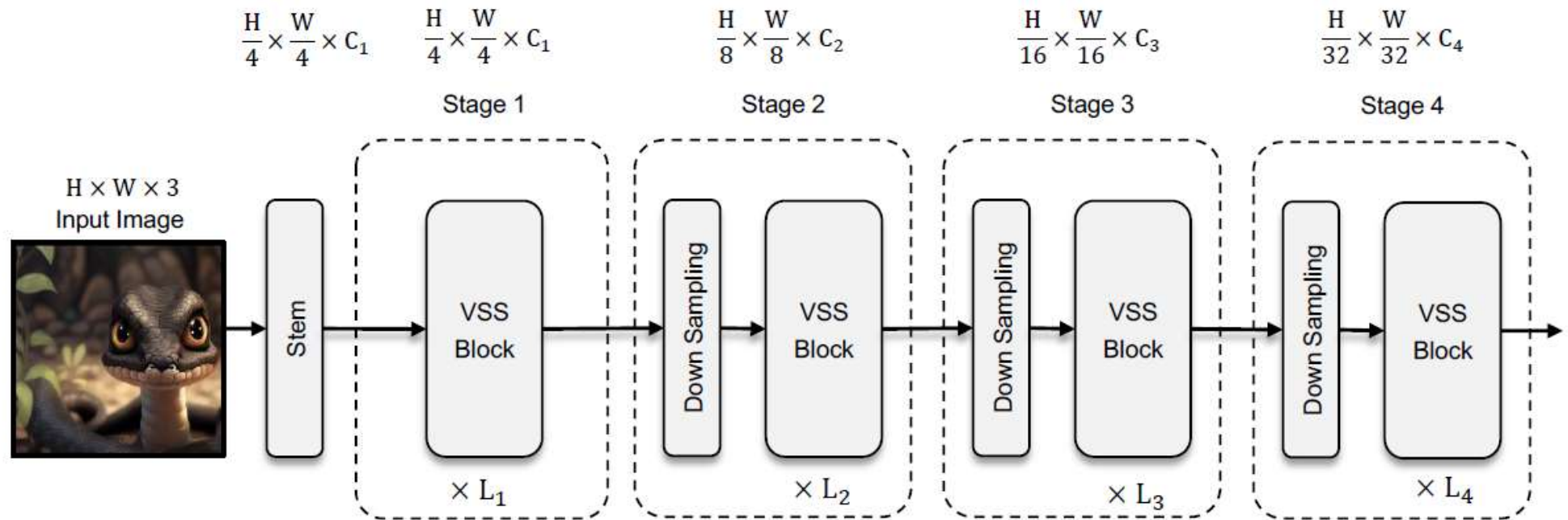


Figure 5: The visualization of the (a) input image, (b) overall activation map, and (c) heat map of CSM.

**Activation Map Visualization.** We have demonstrated that the computational process of selective SSMs encompasses mechanisms resembling self-attention, and this enables us to probe the internal operations of SS2D by visualizing its weight matrices.



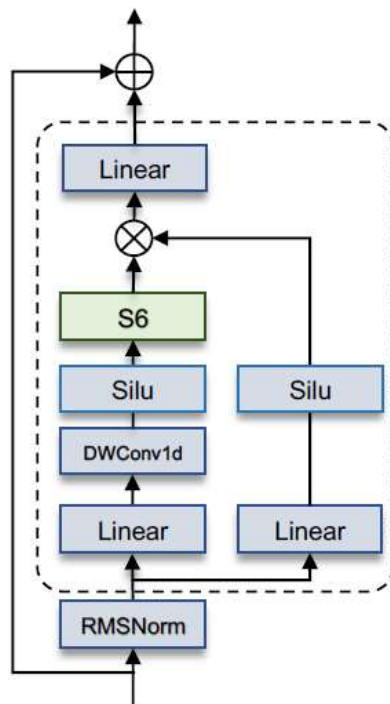
# Overall Architecture



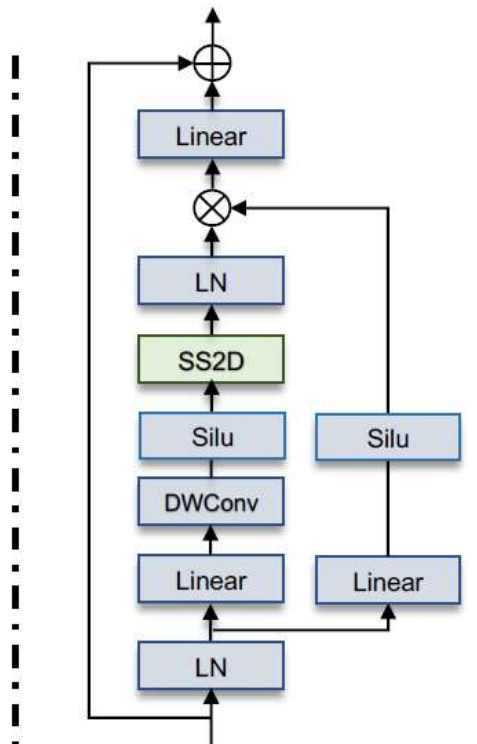
(a) Architecture of VMamba

# VSS Block

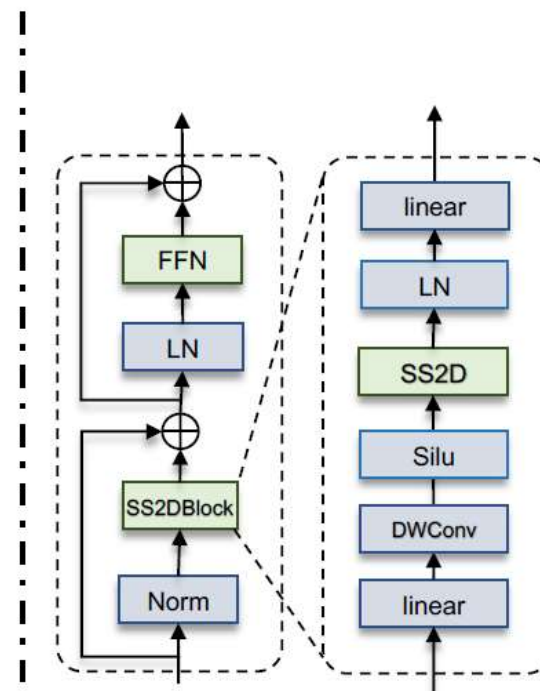
VSS: Visual State Space



(b) Mamba Block



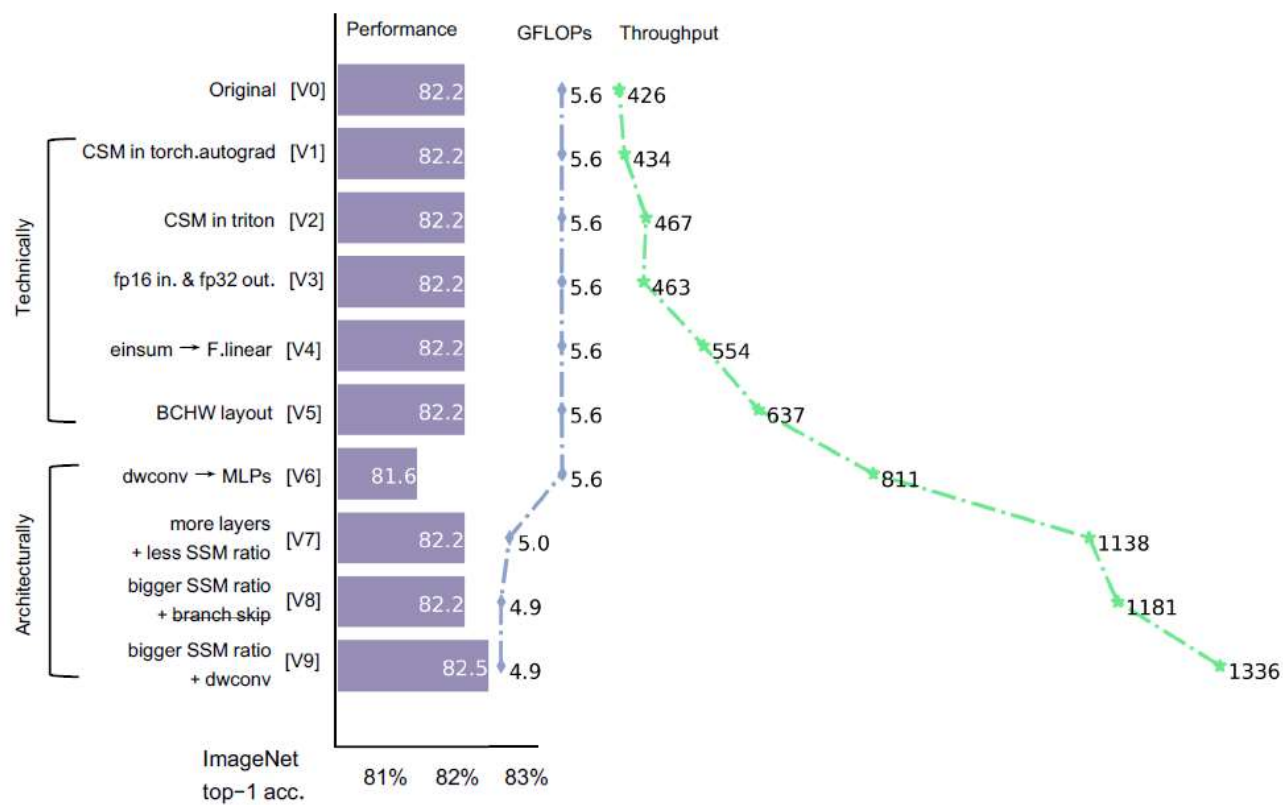
(c) The Vanilla VSS Block



(d) VSS Block

# Accelerating VMamba

(= 많은 시행착오를 통해 조금씩 성능을 높였습니다)



# Experiments

CNN

CNN

Method	Image size	#Param.	FLOPs	Throughput	Train Throughput	ImageNet top-1 acc.
RegNetY-4G [41]	224 <sup>2</sup>	21M	4.0G	–	–	80.0
RegNetY-8G [41]	224 <sup>2</sup>	39M	8.0G	–	–	81.7
RegNetY-16G [41]	224 <sup>2</sup>	84M	16.0G	–	–	82.9
EffNet-B3 [47]	300 <sup>2</sup>	12M	1.8G	–	–	81.6
EffNet-B4 [47]	380 <sup>2</sup>	19M	4.2G	–	–	82.9
EffNet-B5 [47]	456 <sup>2</sup>	30M	9.9G	–	–	83.6
EffNet-B6 [47]	528 <sup>2</sup>	43M	19.0G	–	–	84.0
ViT-B/16 [12]	384 <sup>2</sup>	86M	55.4G	–	–	77.9
ViT-L/16 [12]	384 <sup>2</sup>	307M	190.7G	–	–	76.5
DeiT-S [50]	224 <sup>2</sup>	22M	4.6G	1759	2397	79.8
DeiT-B [50]	224 <sup>2</sup>	86M	17.5G	500	1024	81.8
DeiT-B [50]	384 <sup>2</sup>	86M	55.4G	498	344	83.1

ViT + Knowledge Distillation

ViT 구조 참조하여 ResNet 수정

ConvNeXt-T [33]	224 <sup>2</sup>	29M	4.5G	1189	701	82.1
ConvNeXt-S [33]	224 <sup>2</sup>	50M	8.7G	682	444	83.1
ConvNeXt-B [33]	224 <sup>2</sup>	89M	15.4G	435	334	83.8
HiViT-T [64]	224 <sup>2</sup>	19M	4.6G	1391	1300	82.1
HiViT-S [64]	224 <sup>2</sup>	38M	9.1G	711	697	83.5
HiViT-B [64]	224 <sup>2</sup>	66M	15.9G	456	541	83.8
Swin-T [32]	224 <sup>2</sup>	28M	4.6G	1247	985	81.3
Swin-S [32]	224 <sup>2</sup>	50M	8.7G	719	640	83.0
Swin-B [32]	224 <sup>2</sup>	88M	15.4G	457	494	83.5
S4ND-ConvNeXt-T [40]	224 <sup>2</sup>	30M	-	684	331	82.2
S4ND-ViT-B [40]	224 <sup>2</sup>	89M	-	404	340	80.4
ViM-S [68]	224 <sup>2</sup>	26M	-	811	232 <sup>†</sup>	80.5
VMamba-T	224 <sup>2</sup>	31M	4.9G	1335	464	82.5
VMamba-S	224 <sup>2</sup>	50M	8.7G	874	313	83.6
VMamba-B	224 <sup>2</sup>	89M	15.4G	645	246	83.9

Table 1: Performance comparison on ImageNet-1K. Throughput values are measured with an A100 GPU and an AMD EPYC 7542 CPU, using the toolkit released by [56], following the protocol proposed in [32]. † denotes that the training process only supports float32 datatype.

Swin: ViT + Hierarchical Structure + Shifted Window Self-Attention



# Experiments

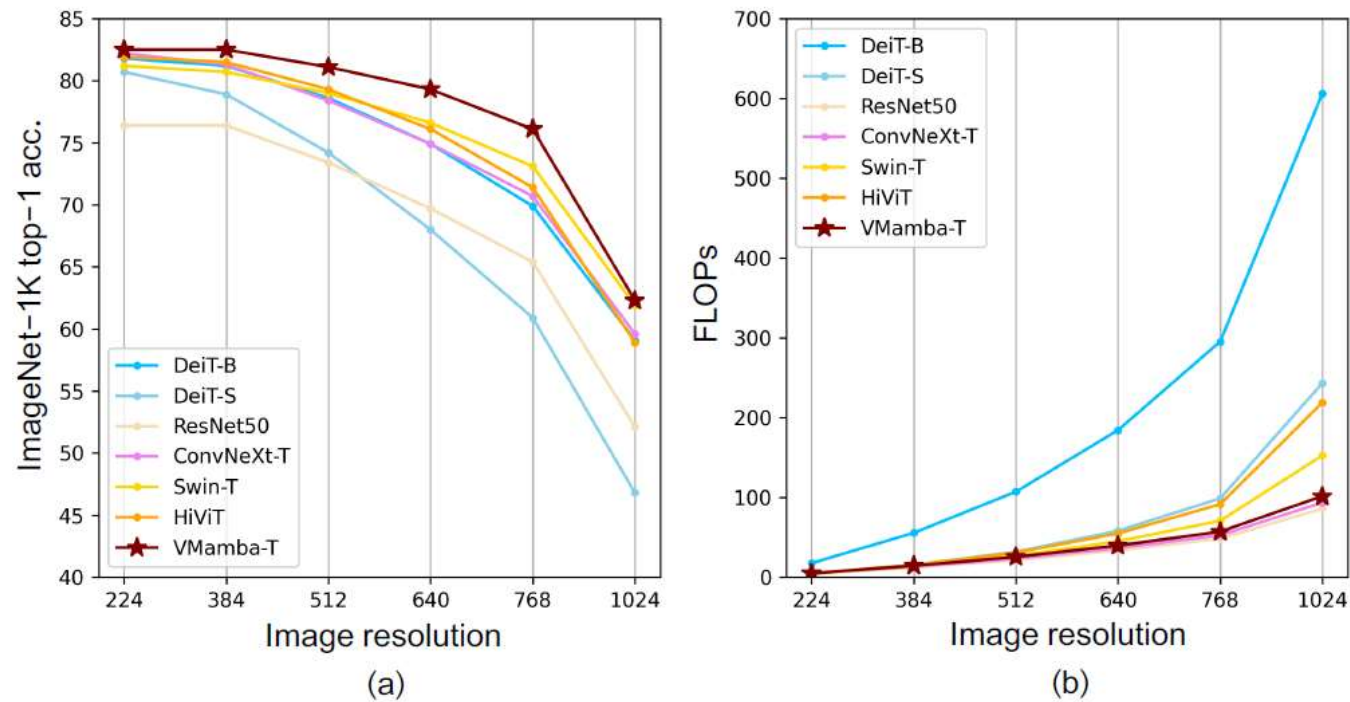


Figure 11: Illustration of the change in (a) classification accuracy and (b) FLOPs with progressively larger test image resolutions.

# Experiments

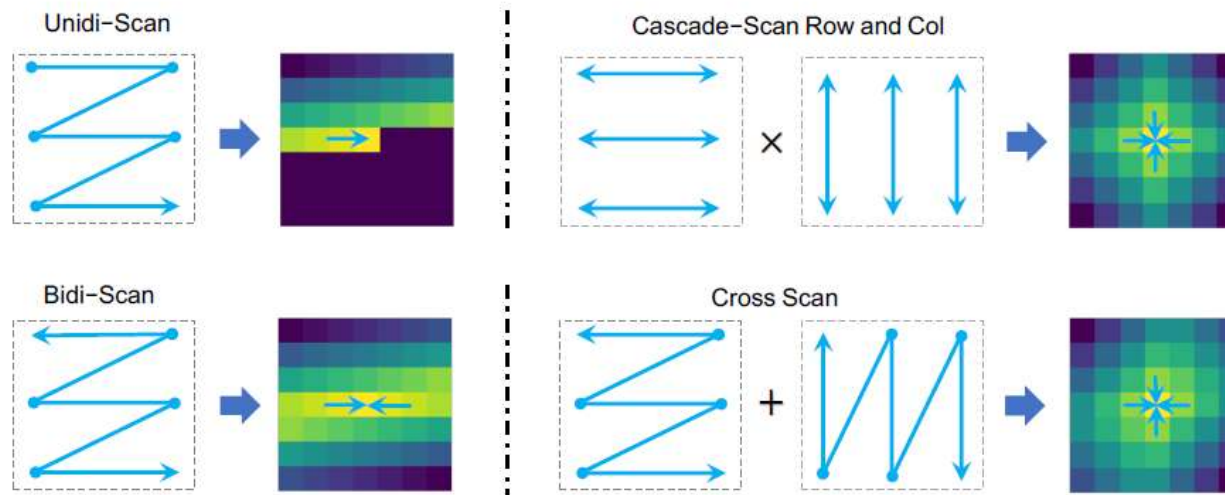


Figure 12: Illustration of different scanning methods for selective scan.

Method	#Param.	FLOPs	Throughput	Train Throughput.	ImageNet top-1 acc.
Unidi-Scan	30.70	4.86	1342	464	82.2
Bidi-Scan	30.70	4.86	1344	465	74.8 <sup>†</sup>
Cascade-Scan	30.70	4.86	817	253	–
CSM	30.70	4.86	1343	464	82.5

Table 5: Performance comparison of different scanning approaches. The proposed CSM achieves

# Inference Results

[2024-04-27 18:26:27 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [60/391]	Time 0.129 (0.472)	Loss 0.7212 (0.5593)	Acc@1 83.594 (88.589)	Acc@5 95.312 (97.797)	Mem 10417MB
[2024-04-27 18:26:28 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [70/391]	Time 0.128 (0.424)	Loss 0.6528 (0.5886)	Acc@1 84.375 (87.797)	Acc@5 99.219 (97.700)	Mem 10417MB
[2024-04-27 18:26:29 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [80/391]	Time 0.129 (0.387)	Loss 0.6416 (0.6034)	Acc@1 85.156 (87.375)	Acc@5 96.094 (97.598)	Mem 10417MB
[2024-04-27 18:26:30 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [90/391]	Time 0.128 (0.359)	Loss 0.8384 (0.6052)	Acc@1 72.656 (87.303)	Acc@5 96.875 (97.588)	Mem 10417MB
[2024-04-27 18:26:32 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [100/391]	Time 0.128 (0.336)	Loss 0.4221 (0.6045)	Acc@1 92.188 (87.106)	Acc@5 97.656 (97.649)	Mem 10417MB
[2024-04-27 18:26:33 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [110/391]	Time 0.128 (0.317)	Loss 1.0830 (0.6134)	Acc@1 60.938 (86.691)	Acc@5 99.219 (97.677)	Mem 10417MB
[2024-04-27 18:26:34 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [120/391]	Time 0.129 (0.302)	Loss 0.5562 (0.6121)	Acc@1 89.062 (86.848)	Acc@5 98.438 (97.708)	Mem 10417MB
[2024-04-27 18:26:36 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [130/391]	Time 0.129 (0.289)	Loss 0.2549 (0.6049)	Acc@1 96.875 (87.083)	Acc@5 100.000 (97.722)	Mem 10417MB
[2024-04-27 18:26:37 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [140/391]	Time 0.129 (0.277)	Loss 1.0137 (0.6041)	Acc@1 77.344 (86.963)	Acc@5 96.875 (97.800)	Mem 10417MB
[2024-04-27 18:26:38 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [150/391]	Time 0.129 (0.267)	Loss 0.7349 (0.6128)	Acc@1 78.125 (86.796)	Acc@5 98.438 (97.724)	Mem 10417MB
[2024-04-27 18:26:39 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [160/391]	Time 0.129 (0.259)	Loss 0.5918 (0.6161)	Acc@1 85.156 (86.753)	Acc@5 97.656 (97.656)	Mem 10417MB
[2024-04-27 18:26:41 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [170/391]	Time 0.129 (0.251)	Loss 1.5020 (0.6349)	Acc@1 66.406 (86.230)	Acc@5 89.844 (97.492)	Mem 10417MB
[2024-04-27 18:26:42 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [180/391]	Time 0.129 (0.244)	Loss 1.6621 (0.6506)	Acc@1 58.594 (85.855)	Acc@5 89.844 (97.302)	Mem 10417MB
[2024-04-27 18:26:43 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [190/391]	Time 0.129 (0.238)	Loss 1.3262 (0.6703)	Acc@1 70.312 (85.320)	Acc@5 92.969 (97.129)	Mem 10417MB
[2024-04-27 18:26:45 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [200/391]	Time 0.129 (0.233)	Loss 0.4934 (0.6863)	Acc@1 88.281 (84.966)	Acc@5 99.219 (96.984)	Mem 10417MB
[2024-04-27 18:26:46 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [210/391]	Time 0.130 (0.228)	Loss 0.8276 (0.6998)	Acc@1 79.688 (84.575)	Acc@5 97.656 (96.901)	Mem 10417MB
[2024-04-27 18:26:47 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [220/391]	Time 0.129 (0.224)	Loss 0.2949 (0.7015)	Acc@1 94.531 (84.523)	Acc@5 98.438 (96.847)	Mem 10417MB
[2024-04-27 18:26:49 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [230/391]	Time 0.129 (0.220)	Loss 1.2549 (0.7075)	Acc@1 71.875 (84.409)	Acc@5 89.844 (96.774)	Mem 10417MB
[2024-04-27 18:26:50 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [240/391]	Time 0.129 (0.216)	Loss 0.7974 (0.7079)	Acc@1 85.938 (84.430)	Acc@5 92.969 (96.713)	Mem 10417MB
[2024-04-27 18:26:51 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [250/391]	Time 0.129 (0.212)	Loss 0.4014 (0.7225)	Acc@1 94.531 (83.945)	Acc@5 99.219 (96.592)	Mem 10417MB
[2024-04-27 18:26:52 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [260/391]	Time 0.129 (0.209)	Loss 0.7847 (0.7338)	Acc@1 81.250 (83.675)	Acc@5 96.875 (96.483)	Mem 10417MB
[2024-04-27 18:26:54 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [270/391]	Time 0.130 (0.206)	Loss 1.4209 (0.7398)	Acc@1 67.188 (83.525)	Acc@5 90.625 (96.411)	Mem 10417MB
[2024-04-27 18:26:55 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [280/391]	Time 0.130 (0.204)	Loss 0.8794 (0.7422)	Acc@1 77.344 (83.502)	Acc@5 96.875 (96.391)	Mem 10417MB
[2024-04-27 18:26:56 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [290/391]	Time 0.129 (0.201)	Loss 1.1279 (0.7476)	Acc@1 53.906 (83.336)	Acc@5 97.656 (96.333)	Mem 10417MB
[2024-04-27 18:26:58 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [300/391]	Time 0.129 (0.199)	Loss 0.6025 (0.7519)	Acc@1 89.062 (83.285)	Acc@5 95.312 (96.257)	Mem 10417MB
[2024-04-27 18:26:59 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [310/391]	Time 0.129 (0.196)	Loss 0.8398 (0.7558)	Acc@1 80.469 (83.224)	Acc@5 94.531 (96.184)	Mem 10417MB
[2024-04-27 18:27:00 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [320/391]	Time 0.130 (0.194)	Loss 0.4866 (0.7615)	Acc@1 91.406 (83.100)	Acc@5 98.438 (96.111)	Mem 10417MB
[2024-04-27 18:27:01 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [330/391]	Time 0.130 (0.192)	Loss 1.1514 (0.7716)	Acc@1 67.969 (82.841)	Acc@5 94.531 (96.009)	Mem 10417MB
[2024-04-27 18:27:03 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [340/391]	Time 0.129 (0.191)	Loss 0.5874 (0.7740)	Acc@1 89.062 (82.790)	Acc@5 98.438 (95.988)	Mem 10417MB
[2024-04-27 18:27:04 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [350/391]	Time 0.129 (0.189)	Loss 0.7593 (0.7753)	Acc@1 82.812 (82.737)	Acc@5 97.656 (95.989)	Mem 10417MB
[2024-04-27 18:27:05 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [360/391]	Time 0.130 (0.187)	Loss 1.1582 (0.7828)	Acc@1 71.875 (82.535)	Acc@5 96.094 (95.947)	Mem 10417MB
[2024-04-27 18:27:07 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [370/391]	Time 0.130 (0.186)	Loss 0.9873 (0.7798)	Acc@1 68.750 (82.560)	Acc@5 98.438 (95.999)	Mem 10417MB
[2024-04-27 18:27:08 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [380/391]	Time 0.129 (0.184)	Loss 0.7271 (0.7823)	Acc@1 82.031 (82.495)	Acc@5 98.438 (95.977)	Mem 10417MB
[2024-04-27 18:27:09 vssm1_tiny_0230]	(main.py 346):	INFO	Test: [390/391]	Time 0.082 (0.183)	Loss 1.3779 (0.7805)	Acc@1 57.500 (82.492)	Acc@5 95.000 (95.996)	Mem 10417MB
[2024-04-27 18:27:10 vssm1_tiny_0230]	(main.py 353):	INFO	* Acc@1 82.492 Acc@5 95.996					
[2024-04-27 18:27:10 vssm1_tiny_0230]	(main.py 197):	INFO	Accuracy of the network ema on the 50000 test images: 82.5%					



TRAIN AND TEST