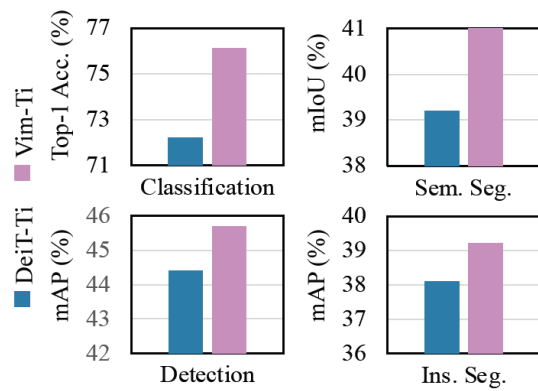


# Vision Mamba: Efficient Visual Representation Learning with Bidirectional State Space Model

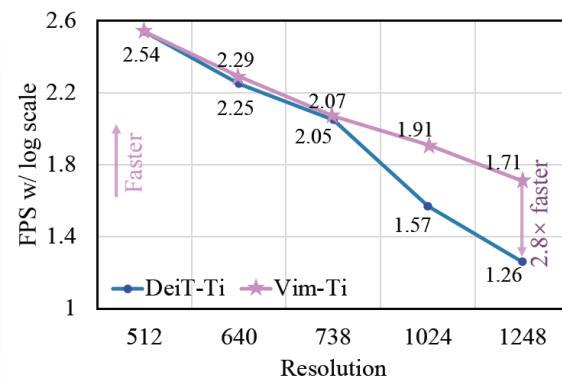
Authors: Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang,  
Wenyu Liu, Xinggang Wang

# Motivation

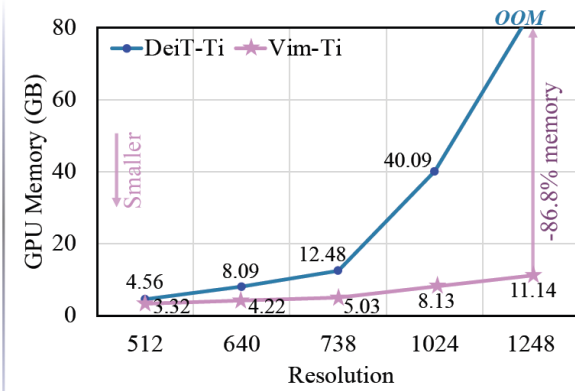
- Introducing Vision Mamba (Vim) with **Bidirectional SSM**
- Improving existing SOTA Transformer based models (DeiT) for high resolution in terms of:
  - **Memory efficiency**
  - **Performance in Vision Tasks**



(a) Accuracy Comparison



(b) Speed Comparison



(c) GPU Memory Comparison

# State Space Models (SSM)

They are inspired in basic 1-D continuous differential models for sequences

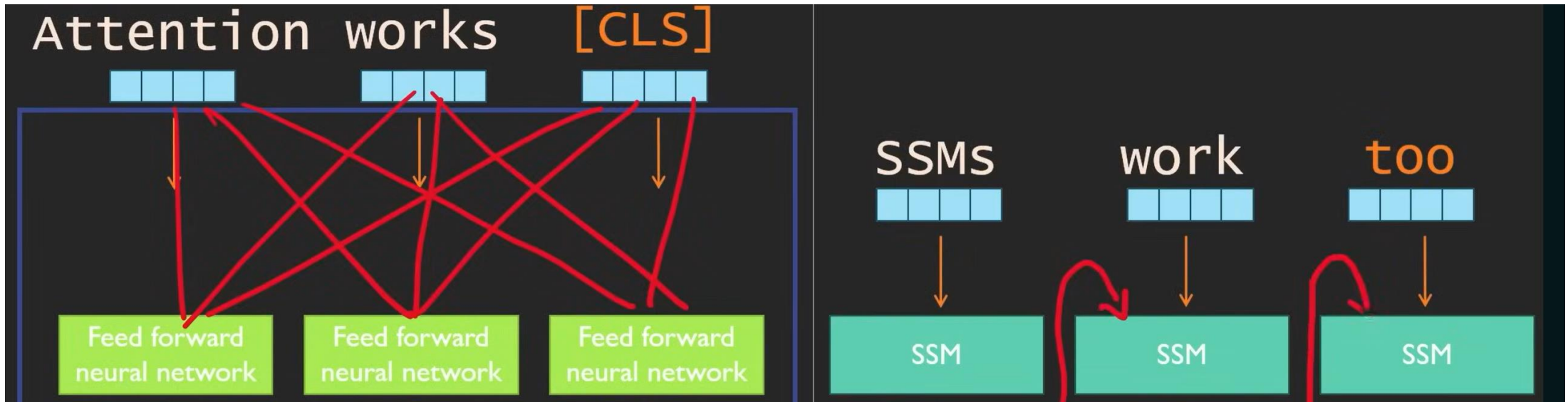
$$\begin{array}{ccc} \begin{array}{l} h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t), \\ y(t) = \mathbf{C}h(t). \end{array} & \begin{array}{c} \xrightarrow{\text{Discretization}} \\ \overline{\mathbf{A}} = \exp(\Delta \mathbf{A}), \\ \overline{\mathbf{B}} = (\Delta \mathbf{A})^{-1}(\exp(\Delta \mathbf{A}) - \mathbf{I}) \cdot \Delta \mathbf{B}. \end{array} & \begin{array}{l} h_t = \overline{\mathbf{A}}h_{t-1} + \overline{\mathbf{B}}x_t, \\ y_t = \mathbf{C}h_t. \end{array} \end{array}$$

**Learnable parameters:** Step size ( $\Delta$ )

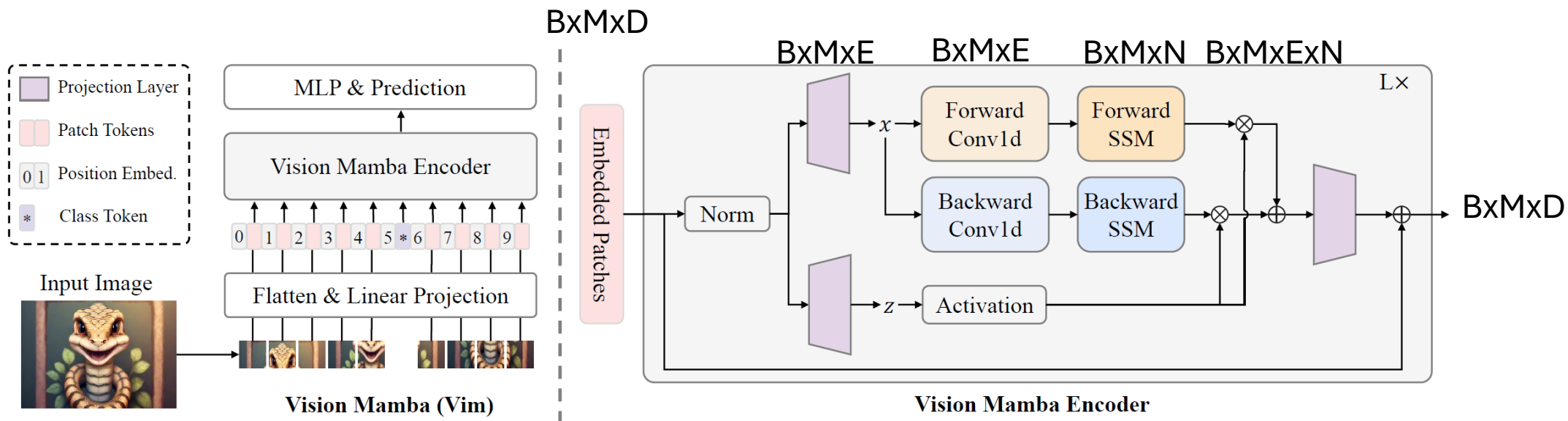
**Method:** Convolution (Efficient in GPU)

$$\begin{array}{l} \overline{\mathbf{K}} = (\mathbf{C}\overline{\mathbf{B}}, \mathbf{C}\overline{\mathbf{A}}\overline{\mathbf{B}}, \dots, \mathbf{C}\overline{\mathbf{A}}^{M-1}\overline{\mathbf{B}}), \\ \mathbf{y} = \mathbf{x} * \overline{\mathbf{K}}, \end{array}$$

# SSMs vs. Transformers (Efficiency)



# Vision Mamba



$$\mathbf{T}_0 = [\mathbf{t}_{cls}; \mathbf{t}_p^1 \mathbf{W}; \mathbf{t}_p^2 \mathbf{W}; \cdots; \mathbf{t}_p^J \mathbf{W}] + \mathbf{E}_{pos},$$

$$\mathbf{T}_l = \mathbf{Vim}(\mathbf{T}_{l-1}) + \mathbf{T}_{l-1},$$

$$\mathbf{f} = \mathbf{Norm}(\mathbf{T}_L^0),$$

$$\hat{p} = \mathbf{MLP}(\mathbf{f}),$$

$L$  : Number of vim blocks

$D$ : Hidden state dimension

$E$ : Expanded state dimension

$N$ : SSM dimensión

	Tiny	Small
$L$ : Number of vim blocks	24	24
$D$ : Hidden state dimension	192	384
$E$ : Expanded state dimension	384	768
$N$ : SSM dimensión	16	16

# Experiments: Classification

Method	image size	#param.	ImageNet top-1 acc.
<b>Convnets</b>			
ResNet-18	224 <sup>2</sup>	12M	69.8
ResNet-50	224 <sup>2</sup>	25M	76.2
ResNet-101	224 <sup>2</sup>	45M	77.4
ResNet-152	224 <sup>2</sup>	60M	78.3
ResNeXt50-32×4d	224 <sup>2</sup>	25M	77.6
RegNetY-4GF	224 <sup>2</sup>	21M	80.0
<b>Transformers</b>			
ViT-B/16	384 <sup>2</sup>	86M	77.9
ViT-L/16	384 <sup>2</sup>	307M	76.5
DeiT-Ti	224 <sup>2</sup>	6M	72.2
DeiT-S	224 <sup>2</sup>	22M	79.8
DeiT-B	224 <sup>2</sup>	86M	81.8
<b>SSMs</b>			
S4ND-ViT-B	224 <sup>2</sup>	89M	80.4
Vim-Ti	224 <sup>2</sup>	7M	76.1
Vim-Ti <sup>†</sup>	224 <sup>2</sup>	7M	78.3 +2.2
Vim-S	224 <sup>2</sup>	26M	80.5
Vim-S <sup>†</sup>	224 <sup>2</sup>	26M	81.6 +1.1

Table 1. Comparison with different backbones on ImageNet-1K validation set. <sup>†</sup> represents the model is fine-tuned with our long sequence setting.

## ImageNet-1K Dataset:

- **1.28M training** images
- **50K validation** images
- 1,000 categories

**Long Sequence Fine-Tuning:** Double of patches than DeiT with the same size (stride 8, 16x16).

## Results:

# Experiments: Classification

Method	image size	#param.	ImageNet top-1 acc.
<b>Convnets</b>			
ResNet-18	224 <sup>2</sup>	12M	69.8
ResNet-50	224 <sup>2</sup>	25M	76.2
ResNet-101	224 <sup>2</sup>	45M	77.4
ResNet-152	224 <sup>2</sup>	60M	78.3
ResNeXt50-32×4d	224 <sup>2</sup>	25M	77.6
RegNetY-4GF	224 <sup>2</sup>	21M	80.0
<b>Transformers</b>			
ViT-B/16	384 <sup>2</sup>	86M	77.9
ViT-L/16	384 <sup>2</sup>	307M	76.5
DeiT-Ti	224 <sup>2</sup>	6M	72.2
DeiT-S	224 <sup>2</sup>	22M	79.8
DeiT-B	224 <sup>2</sup>	86M	81.8
<b>SSMs</b>			
S4ND-ViT-B	224 <sup>2</sup>	89M	80.4
Vim-Ti	224 <sup>2</sup>	7M	76.1
Vim-Ti <sup>†</sup>	224 <sup>2</sup>	7M	78.3 +2.2
Vim-S	224 <sup>2</sup>	26M	80.5
Vim-S <sup>†</sup>	224 <sup>2</sup>	26M	81.6 +1.1

**ImageNet-1K Dataset:**

- **1.28M training** images
- **50K validation** images
- 1,000 categories

**Long Sequence Fine-Tuning:** Double of patches than DeiT with the same size (stride 8, 16x16).

**Results:**

- **3.9 points** higher for **Vim-Tiny** over **DeiT-Tiny**

+3.9

Table 1. Comparison with different backbones on ImageNet-1K validation set. <sup>†</sup> represents the model is fine-tuned with our long sequence setting.

# Experiments: Classification

Method	image size	#param.	ImageNet top-1 acc.
<b>Convnets</b>			
ResNet-18	224 <sup>2</sup>	12M	69.8
ResNet-50	224 <sup>2</sup>	25M	76.2
ResNet-101	224 <sup>2</sup>	45M	77.4
ResNet-152	224 <sup>2</sup>	60M	78.3
ResNeXt50-32×4d	224 <sup>2</sup>	25M	77.6
RegNetY-4GF	224 <sup>2</sup>	21M	80.0
<b>Transformers</b>			
ViT-B/16	384 <sup>2</sup>	86M	77.9
ViT-L/16	384 <sup>2</sup>	307M	76.5
DeiT-Ti	224 <sup>2</sup>	6M	72.2
<b>DeiT-S</b>	224 <sup>2</sup>	22M	<b>79.8</b>
DeiT-B	224 <sup>2</sup>	86M	81.8
<b>SSMs</b>			
S4ND-ViT-B	224 <sup>2</sup>	89M	80.4
Vim-Ti	224 <sup>2</sup>	7M	76.1
Vim-Ti <sup>†</sup>	224 <sup>2</sup>	7M	78.3 +2.2
<b>Vim-S</b>	224 <sup>2</sup>	26M	<b>80.5</b>
Vim-S <sup>†</sup>	224 <sup>2</sup>	26M	81.6 +1.1

**ImageNet-1K Dataset:**

- **1.28M training** images
- **50K validation** images
- 1,000 categories

**Long Sequence Fine-Tuning:** Double of patches than DeiT with the same size (stride 8, 16x16).

**Results:**

- **3.9 points** higher for **Vim-Tiny** over **DeiT-Tiny**
- **0.7 points** higher for **Vim-Small** over **DeiT-Small**

+0.7

Table 1. Comparison with different backbones on ImageNet-1K validation set. <sup>†</sup> represents the model is fine-tuned with our long sequence setting.



# Experiments: Classification

Method	image size	#param.	ImageNet top-1 acc.
<b>Convnets</b>			
ResNet-18	224 <sup>2</sup>	12M	69.8
ResNet-50	224 <sup>2</sup>	25M	76.2
ResNet-101	224 <sup>2</sup>	45M	77.4
ResNet-152	224 <sup>2</sup>	60M	78.3
ResNeXt50-32×4d	224 <sup>2</sup>	25M	77.6
RegNetY-4GF	224 <sup>2</sup>	21M	80.0
<b>Transformers</b>			
ViT-B/16	384 <sup>2</sup>	86M	77.9
ViT-L/16	384 <sup>2</sup>	307M	76.5
DeiT-Ti	224 <sup>2</sup>	6M	72.2
DeiT-S	224 <sup>2</sup>	22M	79.8
DeiT-B	224 <sup>2</sup>	86M	81.8
<b>SSMs</b>			
S4ND-ViT-B	224 <sup>2</sup>	89M	80.4
Vim-Ti	224 <sup>2</sup>	7M	76.1
Vim-Ti <sup>†</sup>	224 <sup>2</sup>	7M	78.3 +2.2
Vim-S	224 <sup>2</sup>	26M	80.5
Vim-S <sup>†</sup>	224 <sup>2</sup>	26M	81.6 +1.1

Table 1. Comparison with different backbones on ImageNet-1K validation set. <sup>†</sup> represents the model is fine-tuned with our long sequence setting.

## ImageNet-1K Dataset:

- **1.28M training** images
- **50K validation** images
- 1,000 categories

**Long Sequence Fine-Tuning:** Double of patches than DeiT with the same size (stride 8, 16x16).

## Results:

- **3.9 points** higher for **Vim-Tiny** over **DeiT-Tiny**
- **0.7 points** higher for **Vim-Small** over **DeiT-Small**
- **Vim-S** achieves results **similar to DeiT-B** with LSFT

# Experiments: Classification

Method	image size	#param.	ImageNet top-1 acc.
<b>Convnets</b>			
ResNet-18	224 <sup>2</sup>	12M	69.8
ResNet-50	224 <sup>2</sup>	25M	76.2
ResNet-101	224 <sup>2</sup>	45M	77.4
ResNet-152	224 <sup>2</sup>	60M	78.3
ResNeXt50-32×4d	224 <sup>2</sup>	25M	77.6
RegNetY-4GF	224 <sup>2</sup>	21M	80.0
<b>Transformers</b>			
ViT-B/16	384 <sup>2</sup>	86M	77.9
ViT-L/16	384 <sup>2</sup>	307M	76.5
DeiT-Ti	224 <sup>2</sup>	6M	72.2
DeiT-S	224 <sup>2</sup>	22M	79.8
DeiT-B	224 <sup>2</sup>	86M	81.8
<b>SSMs</b>			
S4ND-ViT-B	224 <sup>2</sup>	89M	80.4
Vim-Ti	224 <sup>2</sup>	7M	76.1
Vim-Ti <sup>†</sup>	224 <sup>2</sup>	7M	78.3 +2.2
Vim-S	224 <sup>2</sup>	26M	80.5
Vim-S <sup>†</sup>	224 <sup>2</sup>	26M	81.6 +1.1

Table 1. Comparison with different backbones on ImageNet-1K validation set. <sup>†</sup> represents the model is fine-tuned with our long sequence setting.

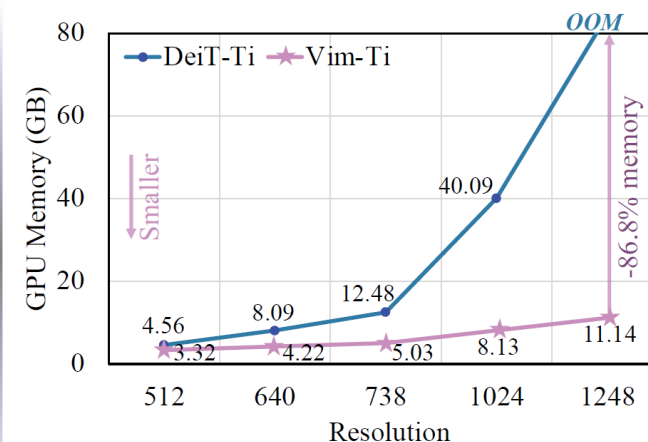
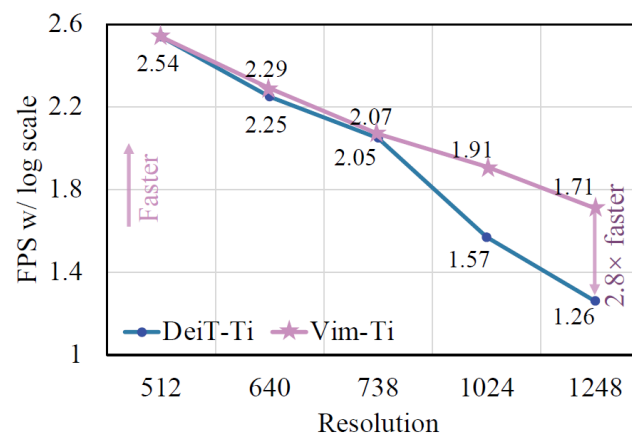
**ImageNet-1K Dataset:**

- **1.28M training** images
- **50K validation** images
- 1,000 categories

**Long Sequence Fine-Tuning:** Double of patches than DeiT with the same size (stride 8, 16x16).

**Results:**

- **3.9 points** higher for **Vim-Tiny** over **DeiT-Tiny**
- **0.7 points** higher for **Vim-Small** over **DeiT-Small**
- **Vim-S** achieves results **similar to DeiT-B** with LSFT
- **1248×1248: Vim is 2.8× faster** than DeiT and **saves 86.8% GPU memory**



# Experiments: Semantic Segmentation

Method	Backbone	image size	#param.	<i>val</i> mIoU
DeepLab v3+	ResNet-101	$512^2$	63M	44.1
UperNet	ResNet-50	$512^2$	67M	41.2
UperNet	ResNet-101	$512^2$	86M	44.9
UperNet	DeiT-Ti	$512^2$	11M	39.2
UperNet	DeiT-S	$512^2$	43M	44.0
UperNet	Vim-Ti	$512^2$	13M	41.0
UperNet	Vim-S	$512^2$	46M	44.9

Table 2. Results of semantic segmentation on the ADE20K *val* set.

## ADE20K Dataset:

- **20K training** images
- **2K validation** images
- 150 categories
- **UperNet** framework

## Results:

# Experiments: Semantic Segmentation

Method	Backbone	image size	#param.	<i>val</i> mIoU
DeepLab v3+	ResNet-101	512 <sup>2</sup>	63M	44.1
UperNet	ResNet-50	512 <sup>2</sup>	67M	41.2
UperNet	ResNet-101	512 <sup>2</sup>	86M	44.9
UperNet	DeiT-Ti	512 <sup>2</sup>	11M	39.2
UperNet	DeiT-S	512 <sup>2</sup>	43M	44.0
UperNet	Vim-Ti	512 <sup>2</sup>	13M	41.0
UperNet	Vim-S	512 <sup>2</sup>	46M	44.9

**ADE20K Dataset:**

- **20K training** images
- **2K validation** images
- 150 categories
- **UperNet** framework

**Results:**

- **1.8 mIoU** higher for **Vim-Ti** over **DeiT-Ti**

Table 2. Results of semantic segmentation on the ADE20K *val* set.

# Experiments: Semantic Segmentation

Method	Backbone	image size	#param.	<i>val</i> mIoU
DeepLab v3+	ResNet-101	512 <sup>2</sup>	63M	44.1
UperNet	ResNet-50	512 <sup>2</sup>	67M	41.2
UperNet	ResNet-101	512 <sup>2</sup>	86M	44.9
UperNet	DeiT-Ti	512 <sup>2</sup>	11M	39.2
UperNet	<b>DeiT-S</b>	512 <sup>2</sup>	43M	<b>44.0</b>
UperNet	Vim-Ti	512 <sup>2</sup>	13M	41.0
UperNet	<b>Vim-S</b>	512 <sup>2</sup>	46M	<b>44.9</b>

**ADE20K Dataset:**

- **20K training** images
- **2K validation** images
- 150 categories
- **UperNet** framework

**Results:**

- **1.8 mIoU** higher for **Vim-Ti** over **DeiT-Ti**
- **0.9 mIoU** higher for **Vim-S** over **DeiT-S**

Table 2. Results of semantic segmentation on the ADE20K *val* set.

# Experiments: Semantic Segmentation

Method	Backbone	image size	#param.	val mIoU
DeepLab v3+	ResNet-101	512 <sup>2</sup>	63M	44.1
UperNet	ResNet-50	512 <sup>2</sup>	67M	41.2
UperNet	ResNet-101	512 <sup>2</sup>	86M	44.9
UperNet	DeiT-Ti	512 <sup>2</sup>	11M	39.2
UperNet	DeiT-S	512 <sup>2</sup>	43M	44.0
UperNet	Vim-Ti	512 <sup>2</sup>	13M	41.0
UperNet	Vim-S	512 <sup>2</sup>	46M	44.9

Table 2. Results of semantic segmentation on the ADE20K *val* set.

## ADE20K Dataset:

- **20K training** images
- **2K validation** images
- 150 categories
- **UperNet** framework

## Results:

- **1.8 mIoU** higher for **Vim-Ti** over **DeiT-Ti**
- **0.9 mIoU** higher for **Vim-S** over **DeiT-S**
- **Vim-S similar to ResNet-101** but **2x fewer parameters**

# Experiments: Object Detection and Instance Segmentation

Backbone	$AP^{\text{box}}$	$AP_{50}^{\text{box}}$	$AP_{75}^{\text{box}}$	$AP_s^{\text{box}}$	$AP_m^{\text{box}}$	$AP_l^{\text{box}}$
DeiT-Ti	44.4	63.0	47.8	26.1	47.4	61.8
Vim-Ti	45.7	63.9	49.6	26.1	49.0	63.2
Backbone	$AP^{\text{mask}}$	$AP_{50}^{\text{mask}}$	$AP_{75}^{\text{mask}}$	$AP_s^{\text{mask}}$	$AP_m^{\text{mask}}$	$AP_l^{\text{mask}}$
DeiT-Ti	38.1	59.9	40.5	18.1	40.5	58.4
Vim-Ti	39.2	60.9	41.7	18.2	41.8	60.2

Table 3. Results of object detection and instance segmentation on the COCO *val* set using Cascade Mask R-CNN [4] framework.

## COCO 2017 Dataset:

- **118K training** images
- **5K validation** images
- Cascade Mask R-CNN base framework

## Results:

- Vim-Ti **surpasses** DeiT-Ti for **medium-size** and **big** objects
- **Better long-range context learning**
- Not necessary window attention

# Experiments: Ablation Study for Design

Bidirectional strategy	ImageNet top-1 acc.	ADE20K mIoU
None	73.2	32.3
Bidirectional Layer	70.9	33.6
Bidirectional SSM	72.8	33.2
Bidirectional SSM + Conv1d	73.9	35.9

Table 4. Ablation study on the bidirectional design. To ensure a fair comparison, we do not use the class token for each experiment. The default setting for Vim is marked in blue .

Classification strategy	ImageNet top-1 acc.
Mean pool	73.9
Max pool	73.4
Head class token	75.2
Double class token	74.3
Middle class token	76.1

Table 5. Ablation study on the classification design. The default setting for Vim is marked in blue .

- Unidirectionality makes Mamba Block (None) fail in dense classification tasks (i.e. segmentation)
  - Bidirectional Block improves segmentation (+1.3 mIoU).
  - **Further enhancement with Bidirectional SSM + Conv1D.**
- 
- Concatenating class token to the visual sequence and performing classification on it outperforms pooling strategy .
  - The best design is by adding **class token at the middle of the visual sequence** and then perform classification on the final middle class token.



# Conclusions

- Possible alternative to Transformer based backbones
- **Computational complexity** linear on sequence length as shown for text
- **Modeling power** similar to DeiT and superior for higher resolution images thanks to efficient long sequences management

# Possible Improvements and Future Work

- Different Datasets and Frameworks
- Ablation study on Hyperparameters
- Self-Supervised Learning
- Comparison of improvements for SOTA systems based on Transformers