



# Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction without Convolutions

Wenhai Wang<sup>1</sup>, Enze Xie<sup>2</sup>, Xiang Li<sup>3</sup>, Deng-Ping Fan<sup>4</sup>,  
Kaitao Song<sup>3</sup>, Ding Liang<sup>5</sup>, Tong Lu<sup>1</sup>, Ping Luo<sup>2</sup>, Ling Shao<sup>4</sup>

<sup>1</sup>Nanjing University    <sup>2</sup>The University of Hong Kong

<sup>3</sup>Nanjing University of Science and Technology    <sup>4</sup>IIAI    <sup>5</sup>SenseTime Research

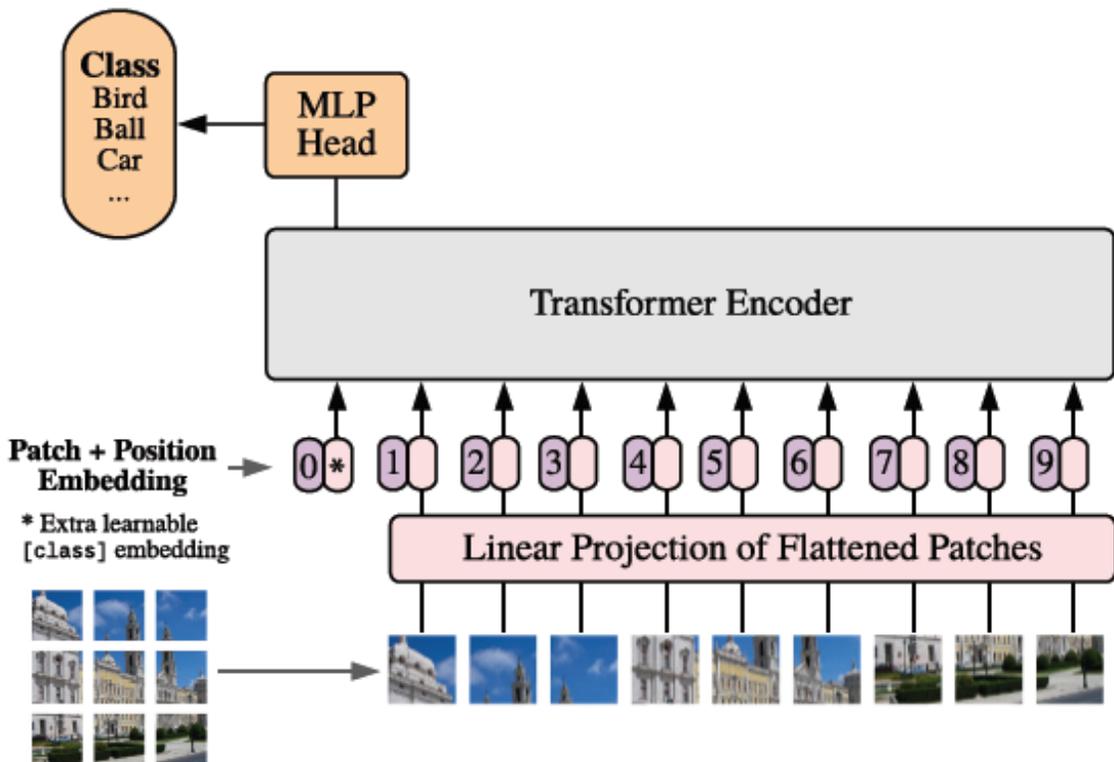
wangwenhai362@smail.nju.edu.cn



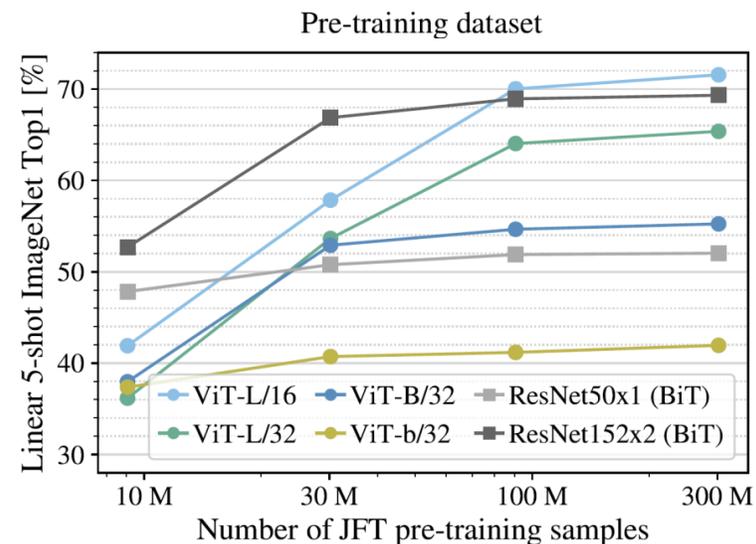
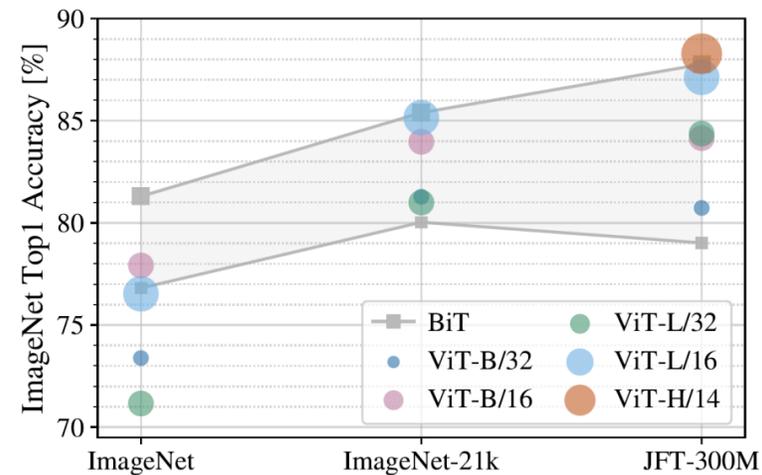
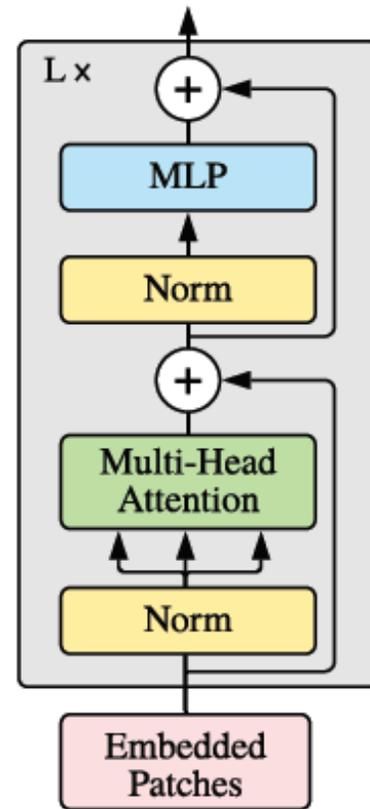
# Related Work



## Vision Transformer (ViT)



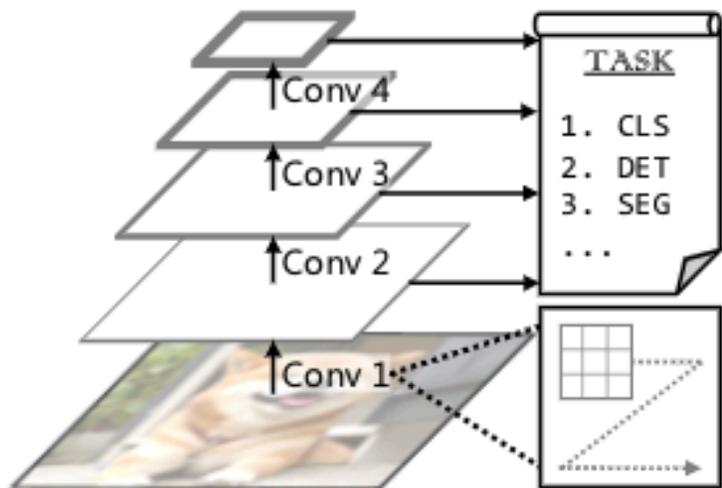
### Transformer Encoder



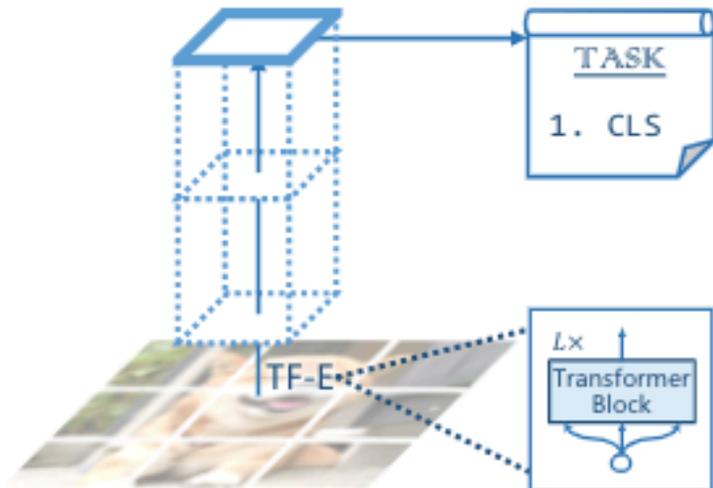
Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." in ICLR, 2020.



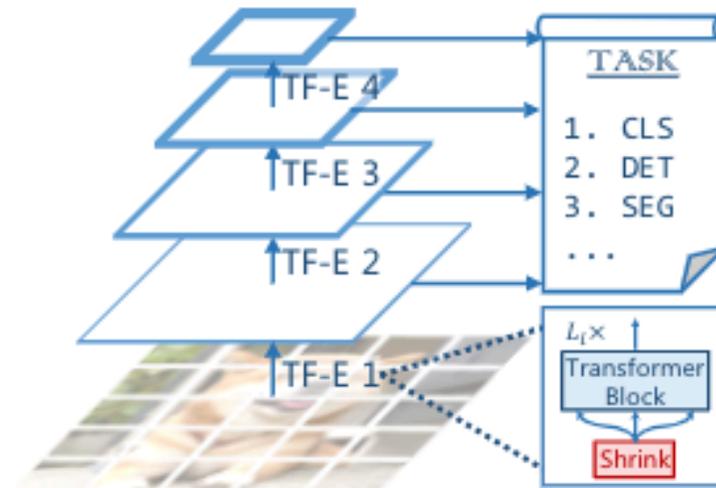
# Pyramid Vision Transformer (PVT)



(a) CNNs: VGG [41], ResNet [15], etc.



(b) Vision Transformer [10]



(c) Pyramid Vision Transformer (ours)

## CNN's Limitations

- Local Receptive Field
- Fixed Weights

## ViT's Limitations

- Columnar Structure
- Low-Resolution Output
- Unsuitable for Det/Seg

## PVT (ours)

- A Transformer backbone as versatile as CNN

Wang, Wenhai, et al. "Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction without Convolutions." arXiv preprint arXiv:2102.12122 (2021).



# Overall Architecture



## Key Points

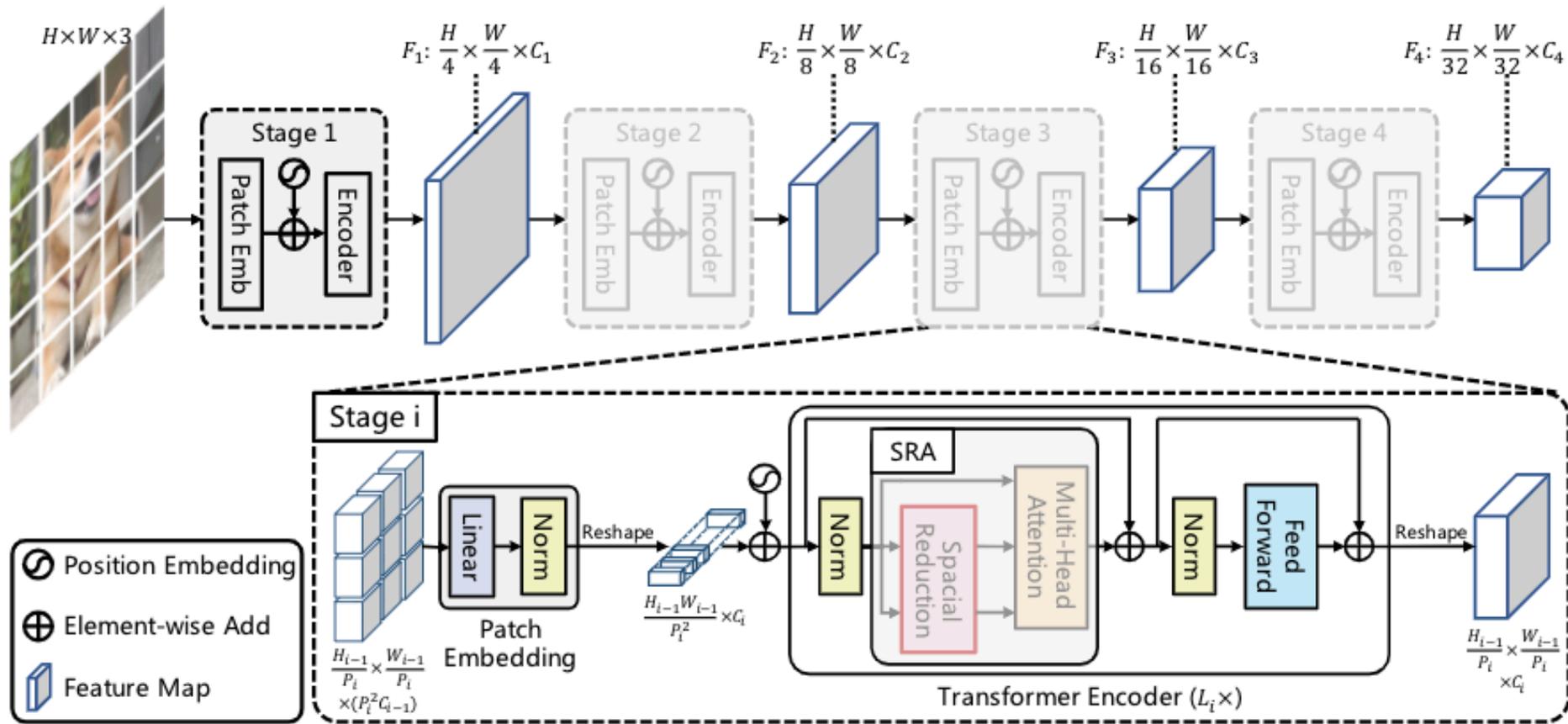
- Four Stages

- Each Stage:

(1) Patch Emb.

(2) Transformer Enc.

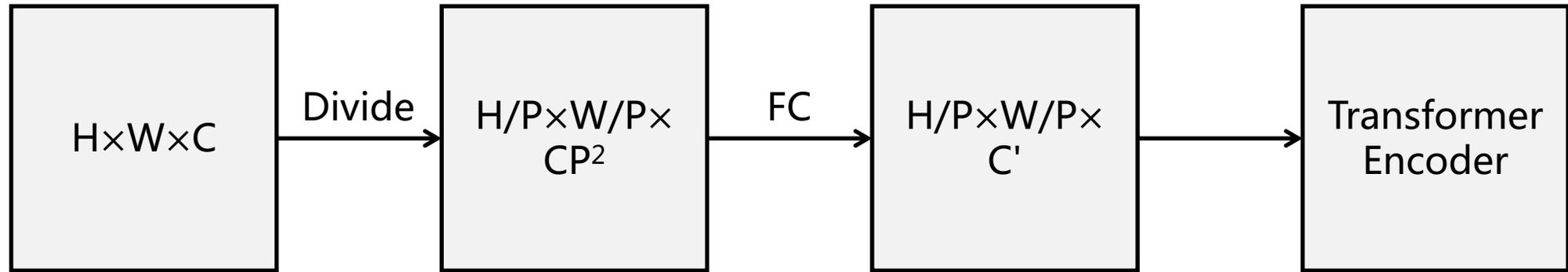
- Spatial-Reduction Attention (SRA) for high-resolution input





# How PVT obtains the feature pyramid?

- Adjusting the **patch size** ( $P_i$ ) in Stage  $i$



The process of the patch embedding in Stage  $i$



# Spatial-Reduction Attention



- SRA

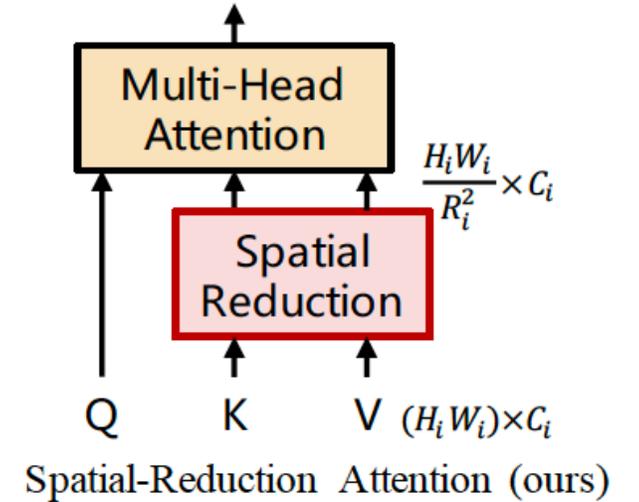
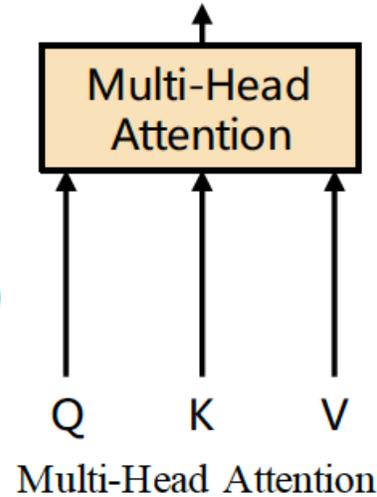
$$\text{SRA}(Q, K, V) = \text{Concat}(\text{head}_0, \dots, \text{head}_{N_i})W^O$$

$$\text{head}_j = \text{Attention}(QW_j^Q, \text{SR}(K)W_j^K, \text{SR}(V)W_j^V)$$

$$\text{SR}(\mathbf{x}) = \text{Norm}(\text{Reshape}(\mathbf{x}, R_i)W^S)$$

$$\text{Attention}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \text{Softmax}\left(\frac{\mathbf{qk}^T}{\sqrt{d_{\text{head}}}}\right)\mathbf{v}$$

Compared to original multi-head attention, the complexity of SRA is  $R_i^2$  times lower!





# Detailed settings

- $P_i$ : the patch size of the stage  $i$ ;
- $C_i$ : the channel number of the output of the stage  $i$ ;
- $L_i$ : the number of encoder layers in the stage  $i$ ;
- $R_i$ : the reduction ratio of the SRA in the stage  $i$ ;
- $N_i$ : the head number of the SRA in the stage  $i$ ;
- $E_i$ : the expansion ratio of the feed-forward layer [51] in the stage  $i$ ;



	Output Size	Layer Name	PVT-Tiny	PVT-Small	PVT-Medium	PVT-Large
Stage 1	$\frac{H}{4} \times \frac{W}{4}$	Patch Embedding	$P_1 = 4; C_1 = 64$			
		Transformer Encoder	$\begin{bmatrix} R_1 = 8 \\ N_1 = 1 \\ E_1 = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} R_1 = 8 \\ N_1 = 1 \\ E_1 = 8 \end{bmatrix} \times 3$	$\begin{bmatrix} R_1 = 8 \\ N_1 = 1 \\ E_1 = 8 \end{bmatrix} \times 3$	$\begin{bmatrix} R_1 = 8 \\ N_1 = 1 \\ E_1 = 8 \end{bmatrix} \times 3$
Stage 2	$\frac{H}{8} \times \frac{W}{8}$	Patch Embedding	$P_2 = 2; C_2 = 128$			
		Transformer Encoder	$\begin{bmatrix} R_2 = 4 \\ N_2 = 2 \\ E_2 = 8 \end{bmatrix} \times 2$	$\begin{bmatrix} R_2 = 4 \\ N_2 = 2 \\ E_2 = 8 \end{bmatrix} \times 3$	$\begin{bmatrix} R_2 = 4 \\ N_2 = 2 \\ E_2 = 8 \end{bmatrix} \times 3$	$\begin{bmatrix} R_2 = 4 \\ N_2 = 2 \\ E_2 = 8 \end{bmatrix} \times 8$
Stage 3	$\frac{H}{16} \times \frac{W}{16}$	Patch Embedding	$P_3 = 2; C_3 = 320$			
		Transformer Encoder	$\begin{bmatrix} R_3 = 2 \\ N_3 = 5 \\ E_3 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} R_3 = 2 \\ N_3 = 5 \\ E_3 = 4 \end{bmatrix} \times 6$	$\begin{bmatrix} R_3 = 2 \\ N_3 = 5 \\ E_3 = 4 \end{bmatrix} \times 18$	$\begin{bmatrix} R_3 = 2 \\ N_3 = 5 \\ E_3 = 4 \end{bmatrix} \times 27$
Stage 4	$\frac{H}{32} \times \frac{W}{32}$	Patch Embedding	$P_4 = 2; C_4 = 512$			
		Transformer Encoder	$\begin{bmatrix} R_4 = 1 \\ N_4 = 8 \\ E_4 = 4 \end{bmatrix} \times 2$	$\begin{bmatrix} R_4 = 1 \\ N_4 = 8 \\ E_4 = 4 \end{bmatrix} \times 3$	$\begin{bmatrix} R_4 = 1 \\ N_4 = 8 \\ E_4 = 4 \end{bmatrix} \times 3$	$\begin{bmatrix} R_4 = 1 \\ N_4 = 8 \\ E_4 = 4 \end{bmatrix} \times 3$

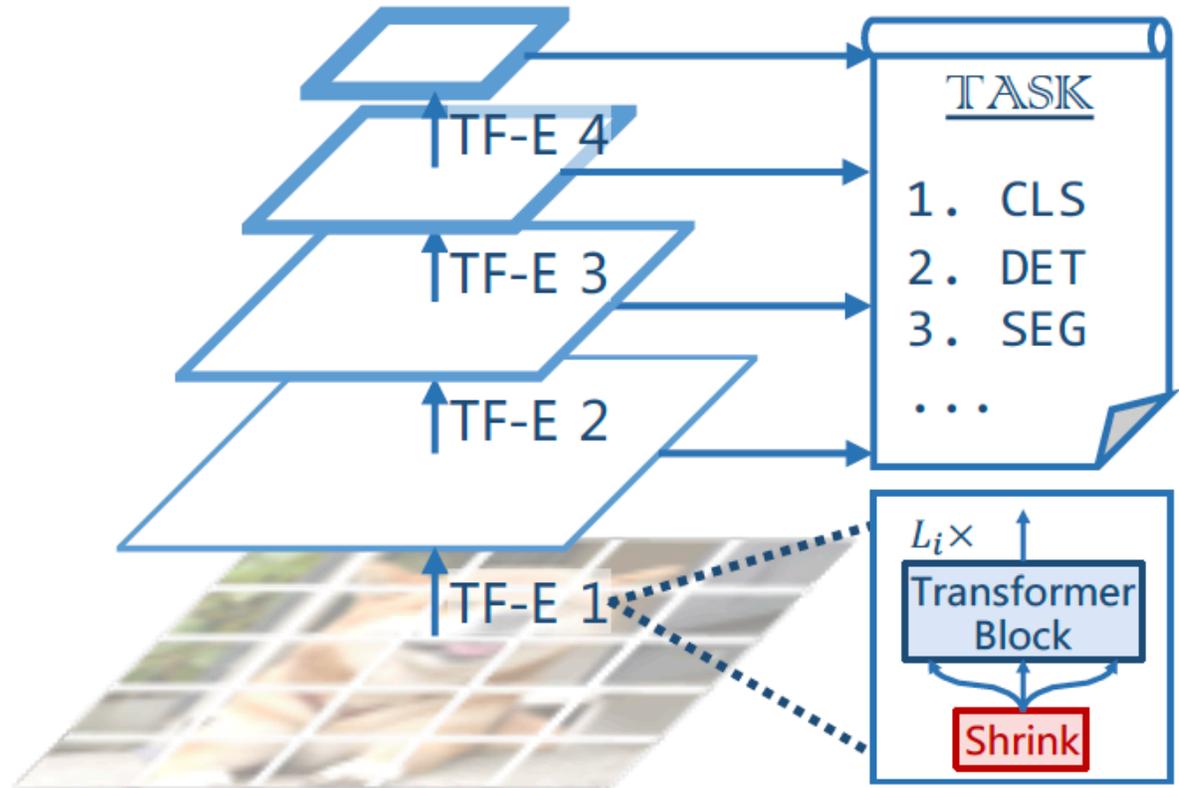
Table A1: **Detailed settings of Pyramid Vision Transformer (PVT) series.** The design follows the two rules of ResNet [5]. (1) With the growth of network depth, the hidden dimension gradually increases, and the output resolution progressively shrinks; (2) The major computation resource is concentrated in Stage 3.



# Advantages



- Multi-Scale/High-Resolution Output
- As versatile as CNN, can be apply to detection/segmentation
- Making pure Transformer detection/segmentation possible, for example
  - (1) PVT + DETR
  - (2) PVT + Trans2Seg





# Performance



- PVT-S vs. R50

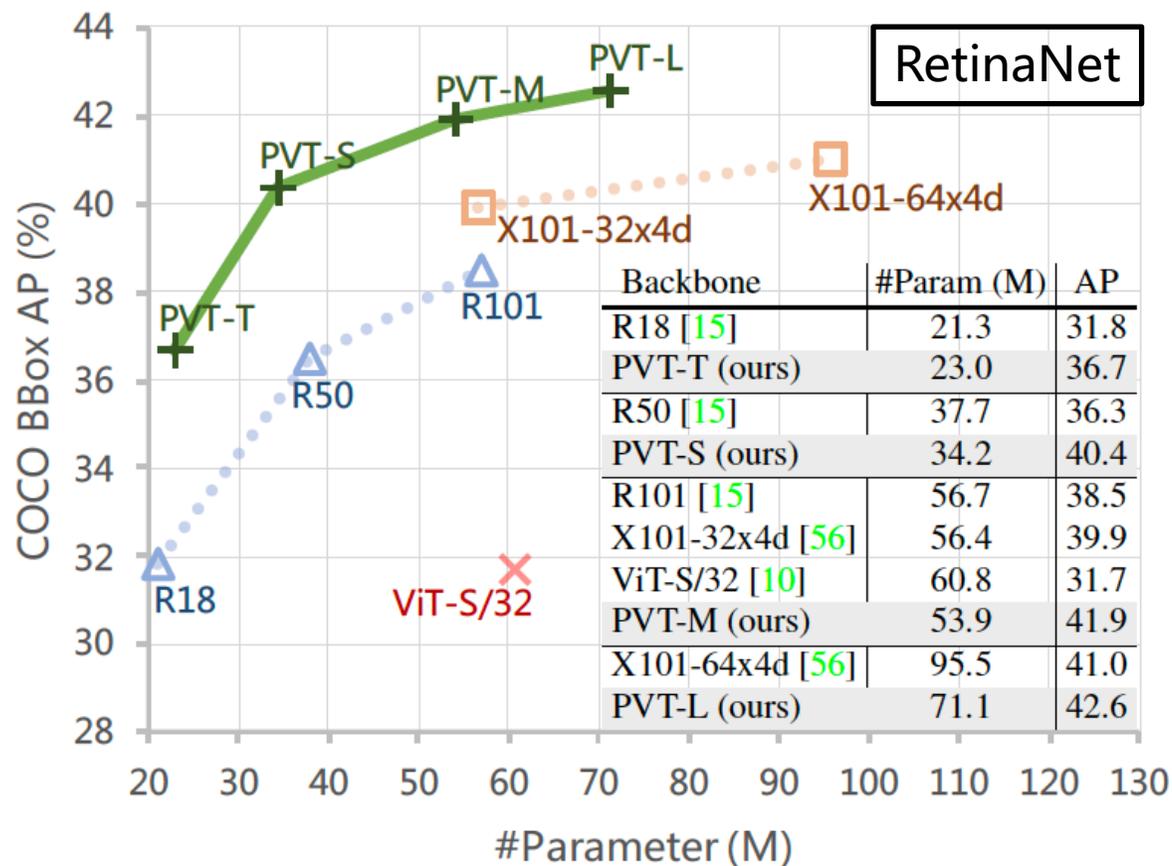
AP: 40.4 vs. 36.3 (+4.1)

#Param: 34.2 vs. 37.7

- PVT-L vs. X101-64x4d

AP: 42.6 vs. 41.0 (+1.6)

#Param: 71.1 vs. 95.5 (20% fewer)





# Deeper vs. Wider & Pretrained



## - Deeper vs. Wider

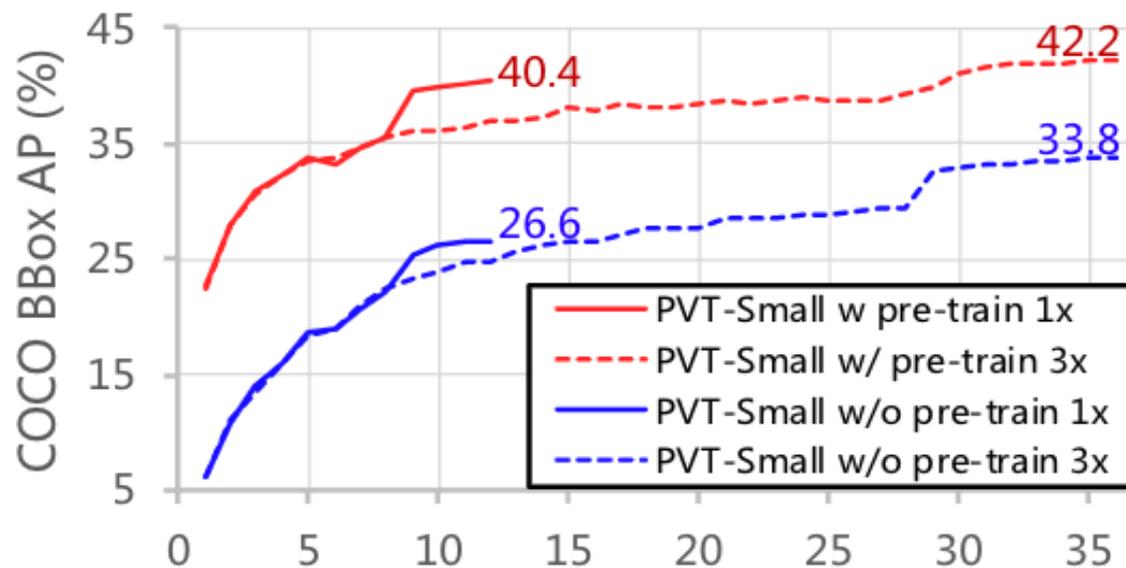
Going deeper is better than going wider.

## - Pretrained

Pretrained weights can help PVT converge faster and better.

Method	#Param (M)	Top-1	RetinaNet 1x		
			AP	AP <sub>50</sub>	AP <sub>75</sub>
Wider PVT-Small	46.8	19.3	40.8	61.8	43.3
Deeper PVT-Small	44.2	18.8	41.9	63.1	44.3

Table A3: **Deeper vs. Wider.** “Top-1” denotes the top-1 error on the ImageNet validation set. “AP” denotes the bounding box AP on COCO val2017. The deeper model obtains better performance than the wider model under comparable parameter number.





# Computation Cost

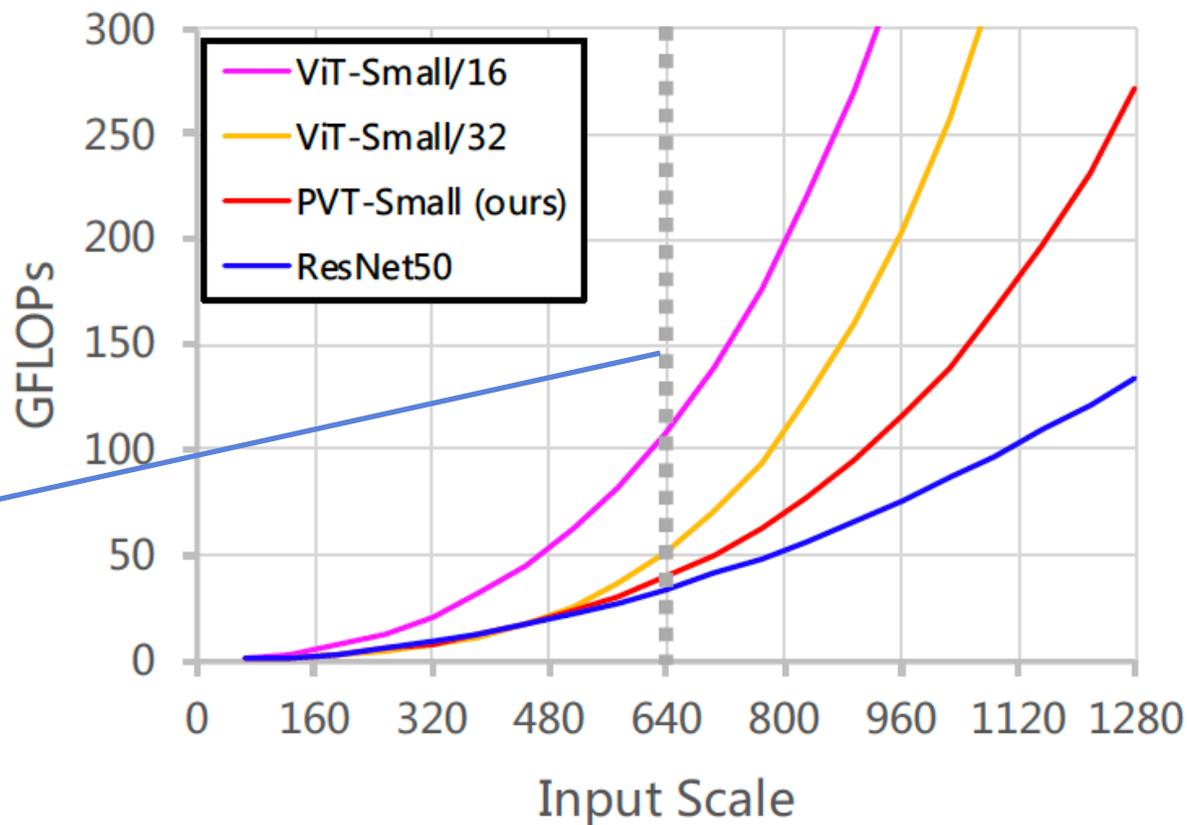


## - FLOPs Growth Rate

ViT-S/16 > ViT-S/32 > PVT-S > R50

PVT is more suitable for medium-resolution.

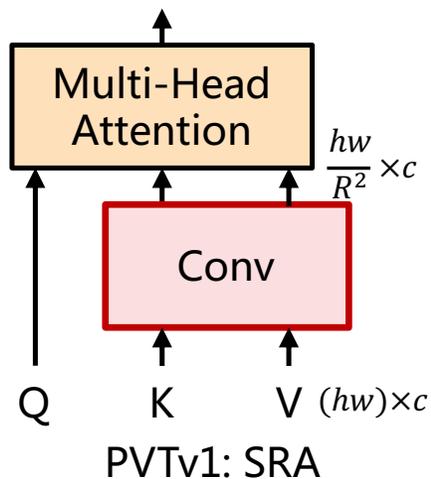
PVTv2 solve this problem by Linear SRA!!!



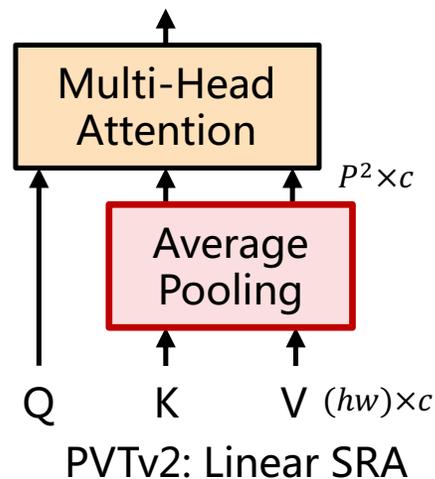


# Improvements

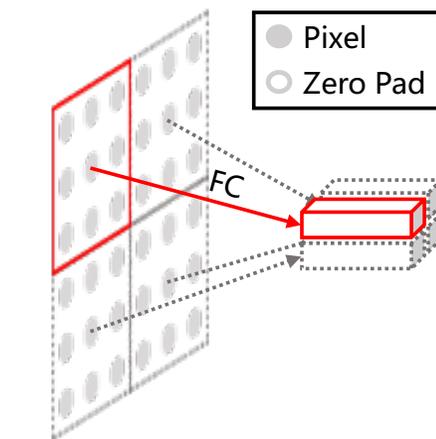
- Overlapping Patch Embedding
- Convolutional Feed-Forward
- Linear SRA



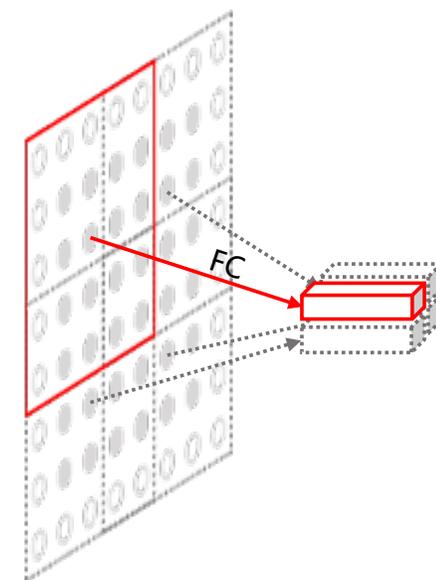
PVTv1: SRA



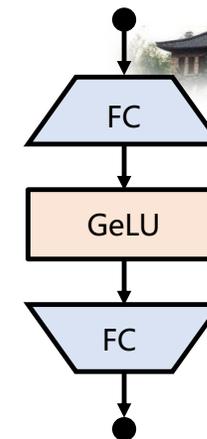
PVTv2: Linear SRA



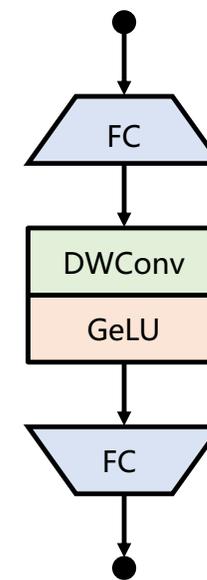
Before: Original Patch Embedding



Improved: Overlapping Patch Embedding



Before: Original Feed-Forward



Improved: Convolutional Feed-Forward





# Performance



## Classification on ImageNet

Method	#Param (M)	GFLOPS	Top-1 Acc (%)
ResNeXt101-64x4d [33]	83.5	15.6	79.6
RegNetY-16G [24]	84.0	16.0	82.9
ViT-Base/16 [7]	86.6	17.6	81.8
DeiT-Base/16 [29]	86.6	17.6	81.8
Swin-B [21]	88.0	15.4	83.3
Twins-SVT-L [4]	99.2	14.8	83.3
PVTv2-B5 (ours)	<b>82.0</b>	<b>11.8</b>	<b>83.8</b>

## Detection on COCO

Backbone	Method	AP <sup>b</sup>	AP <sub>50</sub> <sup>b</sup>	AP <sub>75</sub> <sup>b</sup>	#P (M)	GFLOPS
ResNet50 [13]	Cascade Mask R-CNN [1]	46.3	64.3	50.5	82	739
Swin-T [21]		50.5	69.3	54.9	86	745
PVTv2-B2-Li (ours)		50.9	69.5	55.2	<b>80</b>	<b>725</b>
PVTv2-B2 (ours)		<b>51.1</b>	<b>69.8</b>	<b>55.3</b>	83	788
ResNet50 [13]	ATSS [37]	43.5	61.9	47.0	32	205
Swin-T [21]		47.2	66.5	51.3	36	215
PVTv2-B2-Li (ours)		48.9	68.1	53.4	<b>30</b>	<b>194</b>
PVTv2-B2 (ours)		<b>49.9</b>	<b>69.1</b>	<b>54.1</b>	33	258
ResNet50 [13]	GFL [17]	44.5	63.0	48.3	32	208
Swin-T [21]		47.6	66.8	51.7	36	215
PVTv2-B2-Li (ours)		49.2	68.2	53.7	<b>30</b>	<b>197</b>
PVTv2-B2 (ours)		<b>50.2</b>	<b>69.4</b>	<b>54.7</b>	33	261
ResNet50 [13]	Sparse R-CNN [26]	44.5	63.4	48.2	106	166
Swin-T [21]		47.9	67.3	52.3	110	172
PVTv2-B2-Li (ours)		48.9	68.3	53.4	<b>104</b>	<b>151</b>
PVTv2-B2 (ours)		<b>50.1</b>	<b>69.5</b>	<b>54.9</b>	107	215



# Future Direction



- Efficient Attention Layer

Deformable Attention, Linear SRA, ...

- Position Embedding for 2D/3D Images

CPVT, Local ViT, ...

- Pure Transformer Vision Models

Segformer, YOLOS, ...

- Transformer + NAS/Pruning/Distillation/Quantification

Visual Transformer Pruning, Patch Slimming, ...

- Multimodal Transformer (*e.g.*, CV+NLP)

CLIP, Kaleido-BERT, ...



# Thanks

Code: <https://github.com/whai362/PVT>