# Instance Level Detection and Beyond

The University of Hong Kong
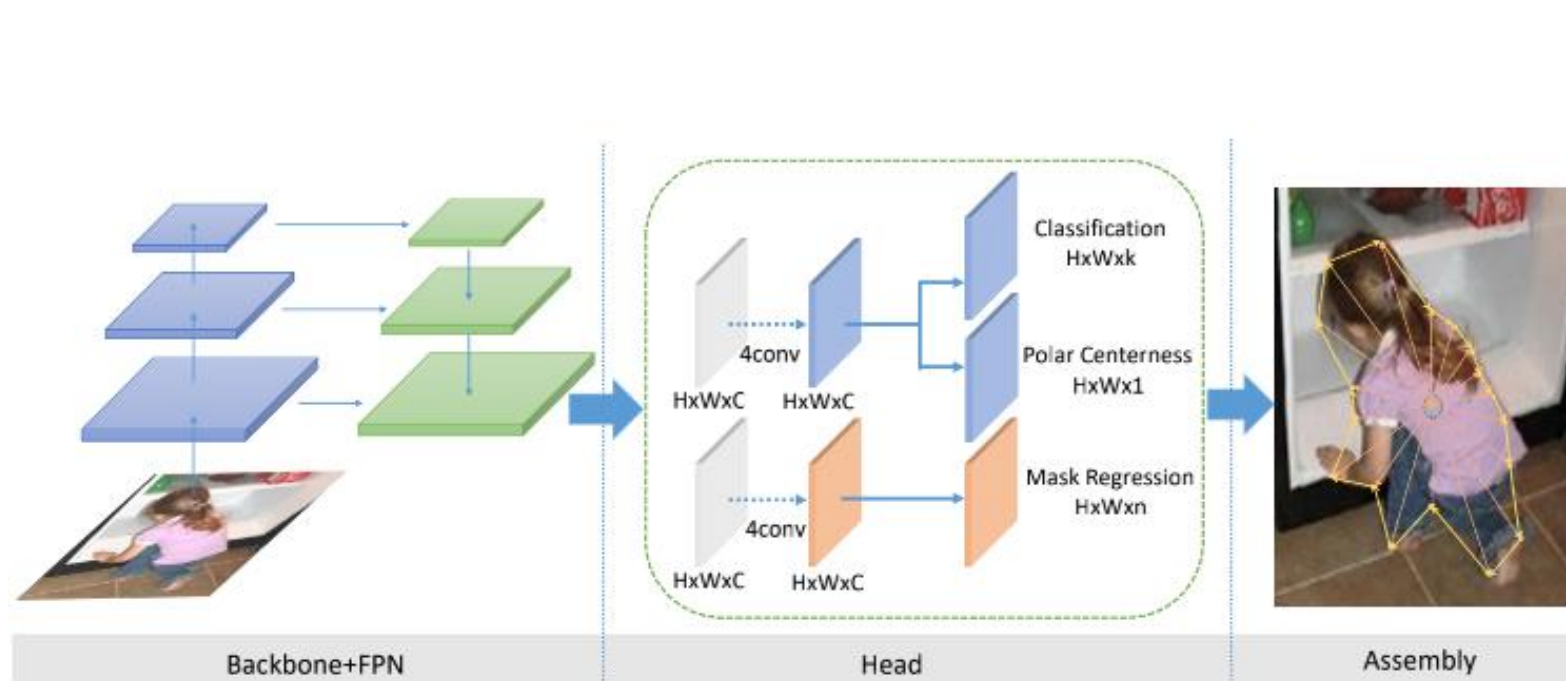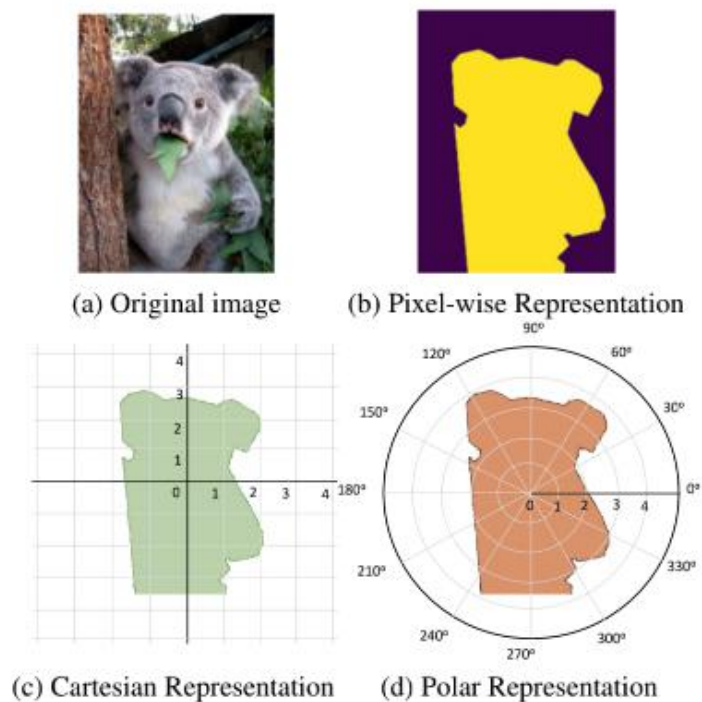Xie Enze
Supervisor: Dr. Luo Ping

# Introduction

- My research is computer vision, especially detection and segmentation with self/semi/weak supervision.

- My work PolarMask was selected as CVPR 2020 Top-10 Influential Papers and with 730+ github star.

- I co-developed OpenSelfSup (1.2k github star), which is a popular self-supervised learning toolbox.

- Publication: 7 first/co-first papers (TPAMI, CVPR, ICCV, ECCV). Citation: 443

- Google scholar: https://scholar.google.com/citations?user=42MVVPgAAAAJ&hl=zh-CN

# Outline

- 1. PolarMask: Single Shot Instance Segmentation with Polar Representation (CPVR20 Oral & TPAMI21, Top-10 influence paper)

- 2. DetCo: Unsupervised Contrastive Learning for Object Detection(a new pretrain friendly for object detection, work in Huawei)

- Pyramid Vision Transformer: A Versatile Backbone for Dense Prediction without Convolutions (say goodbye to CNNs)

# PolarMask

Summary： We introduced a Polar Representation to reformulate the instance segmentation problem.



(a) Original image

(b) Pixel-wise Representation

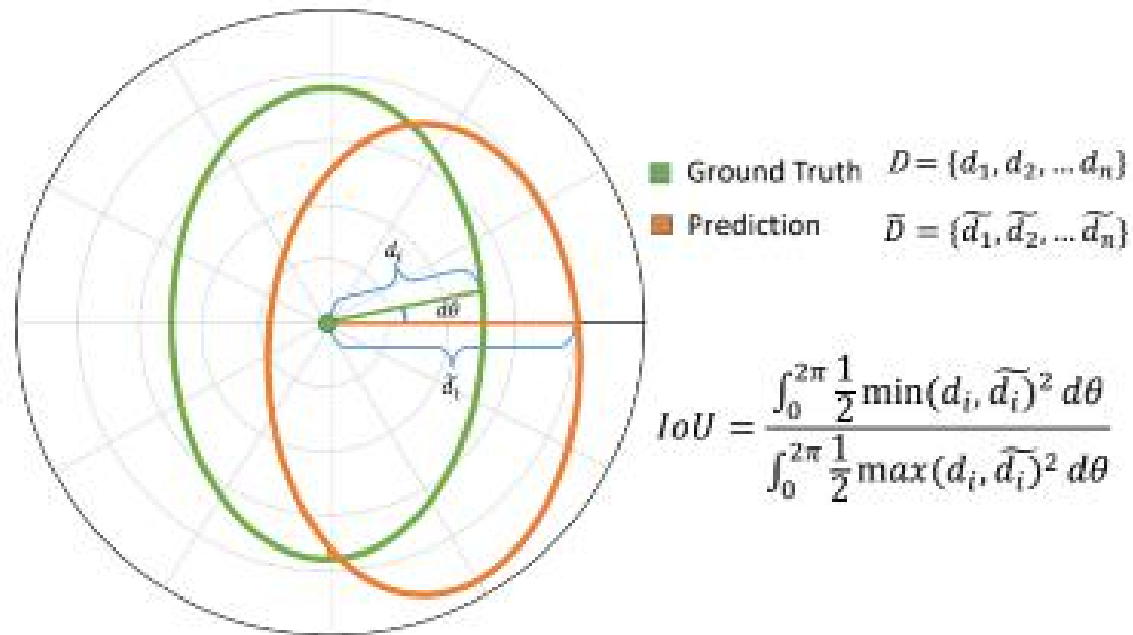(c) Cartesian Representation

(d) Polar Representation

**Figure 2** – The overall pipeline of PolarMask. The left part contains the backbone and feature pyramid to extract features of different levels. The middle part is the two heads for classification and polar mask regression. $H, W, C$ are the height, width, channels of feature maps, respectively, and $k$ is the number of categories (e.g., $k = 80$ on the COCO dataset), $n$ is the number of rays (e.g., $n = 36$)

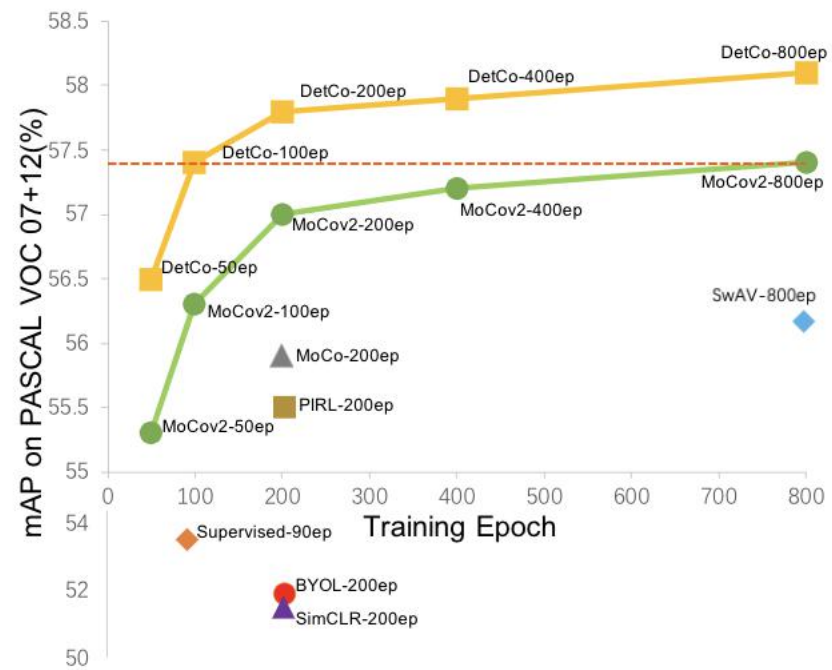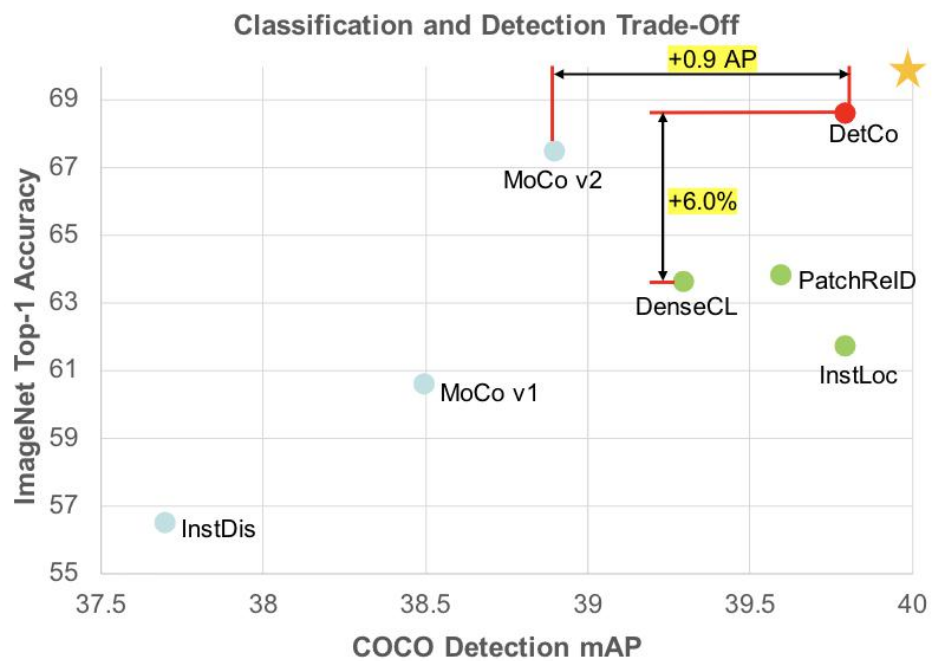Backbone+FPN  Head  Assembly

Classification HxWxk

Polar Centerness HxWx1

Mask Regression HxWxn

4conv HxWxC HxWxC

4conv HxWxC HxWxC

# PolarMask



Ground Truth $D = \{d_1, d_2, \ldots d_n\}$

Prediction $\bar{D} = \{\widetilde{d_1}, \widetilde{d_2}, \ldots \widetilde{d_n}\}$

$$IoU = \frac{\int_0^{2\pi} \frac{1}{2} \min(d_i, \widetilde{d_i})^2 \, d\theta}{\int_0^{2\pi} \frac{1}{2} \max(d_i, \widetilde{d_i})^2 \, d\theta}$$

**Figure 5 – Mask IoU in Polar Representation.** Mask IoU (interaction area over union area) in the polar coordinate can be calculated by integrating the differential IoU area in terms of differential angles.

## Mask IoU Loss

Use calculus to design loss

# DetCo



**Classification and Detection Trade-Off**



Figure 3. Comparisons of mAP on PASCAL VOC 07+12 object...

# DetCo

Summary：We successfully use large-scale unlabeled data to help downstream object detection tasks.
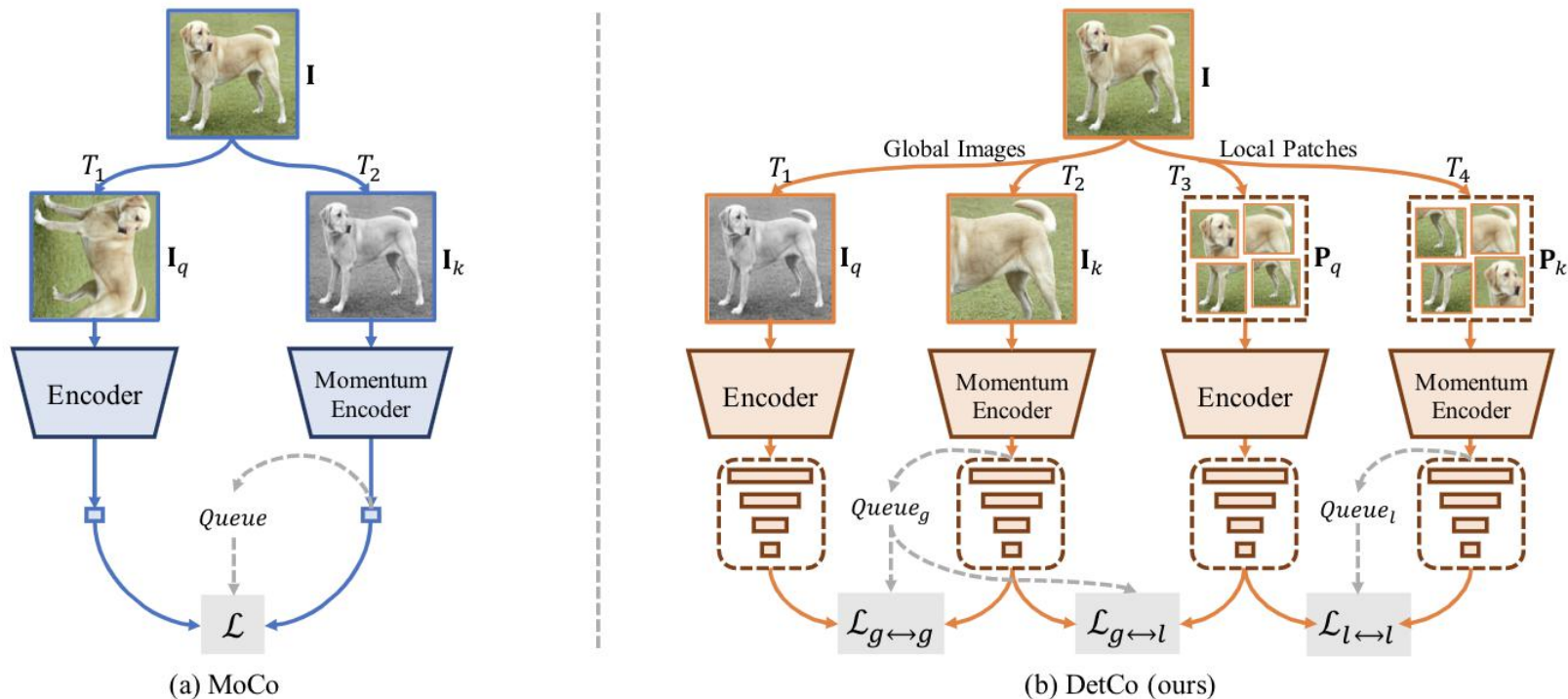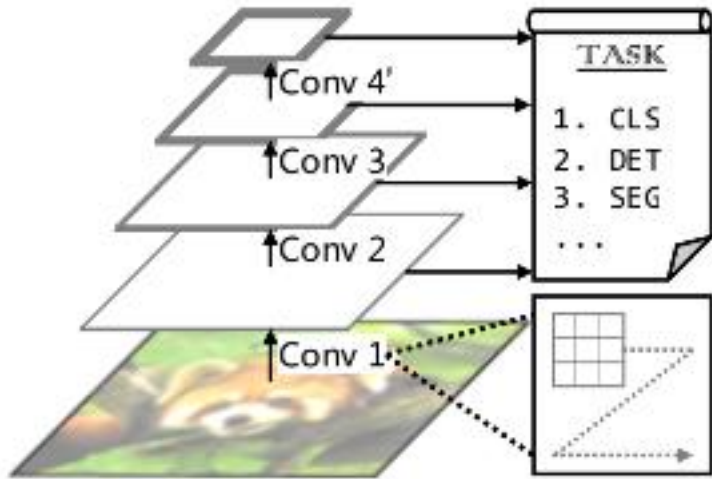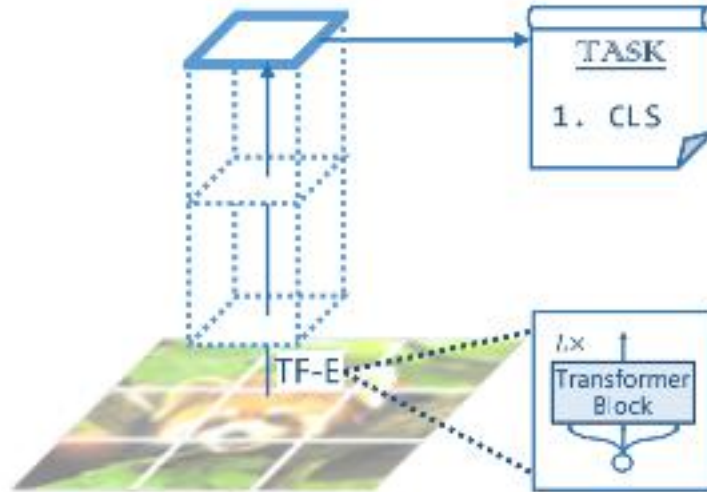


Figure 2. **The overall pipeline of DetCo compared with MoCo [19].** (a) is MoCo's framework, which only considers the single high-level feature and learning contrast from a global perspective. (b) is our DetCo, which learns representation with multi-level supervision and adds two additional local patch sets for input, building contrastive loss cross the global and local views. Note that "$T$" means image transforms. "$Queue_{g/l}$" means different memory banks [40] for global/local features.

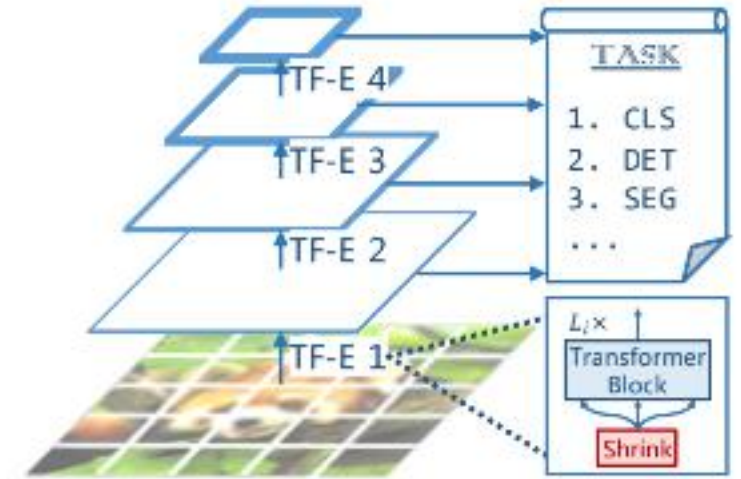# PVT: next generation versatile vision backbone

Summary: A pure Transformer backbone for dense prediction, *e.g.* object detection and semantic segmentation.



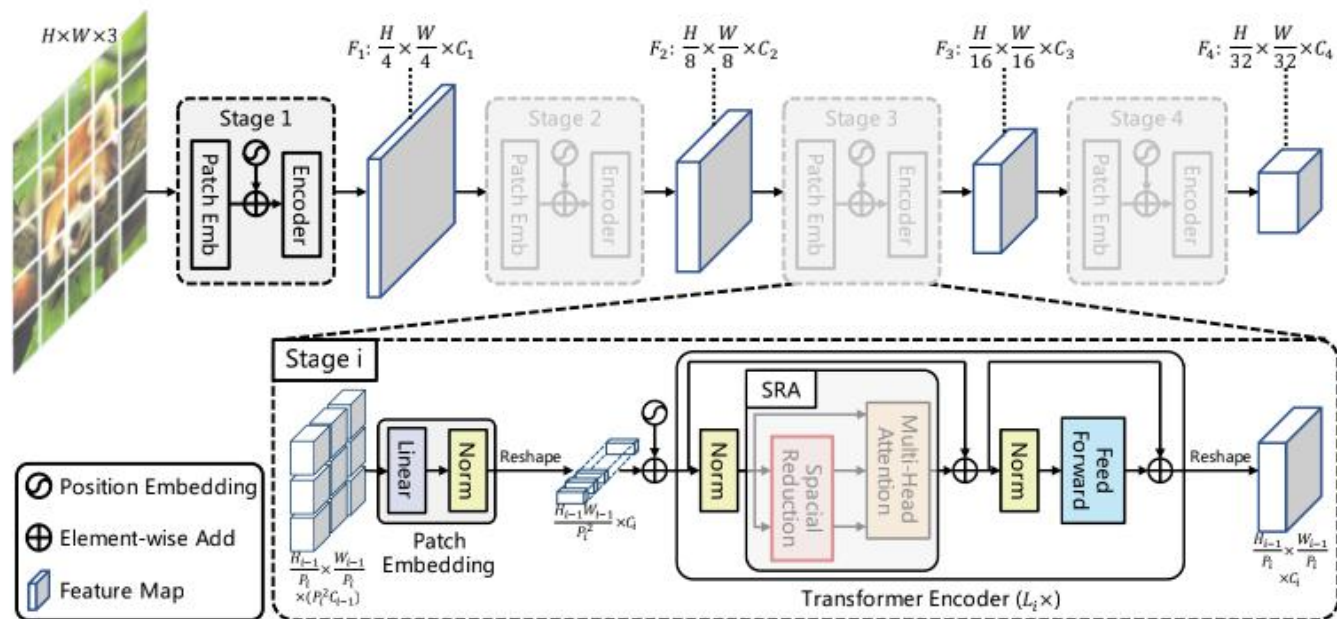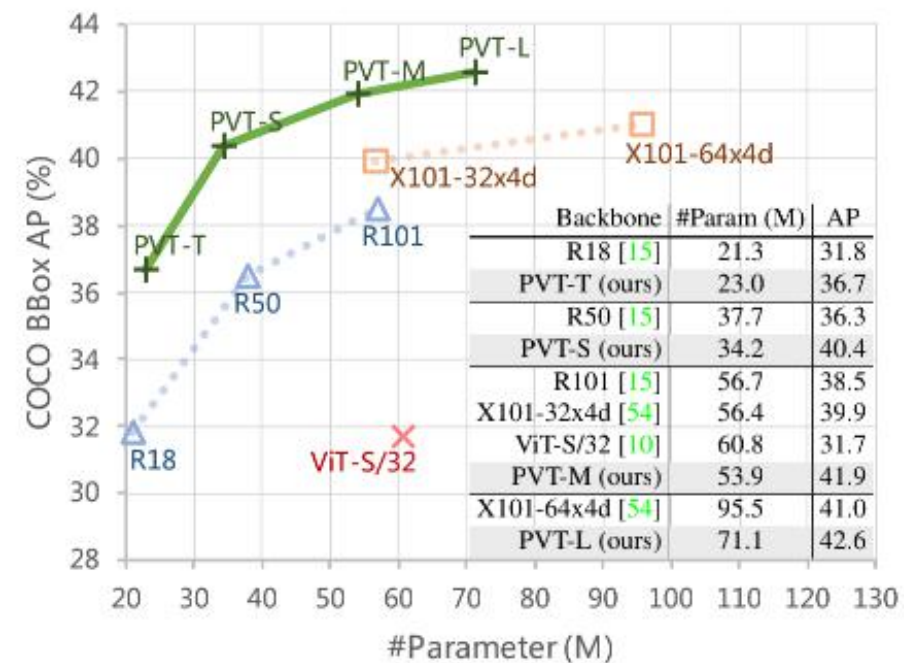(a) CNNs: VGG [41], ResNet [15], *etc.*

(b) Vision Transformer [10]

(c) Pyramid Vision Transformer (ours)

# PVT: next generation versatile vision backbone

# PVT: next generation versatile vision backbone

- 1. Combine DETR and Trans2Seg, we first build a **pure Transformer** Detection and Segmentation pipeline.

- 2. We show that vision tasks can work better without CNNs.(say goodbye to CNNs)

- 3. improve PVT, making it faster and stronger.