

**CHULA ENGINEERING**

Innovation toward Sustainability | ACT NOW



# — OBJECT DETECTION :

## AI VEHICLE TRACKER



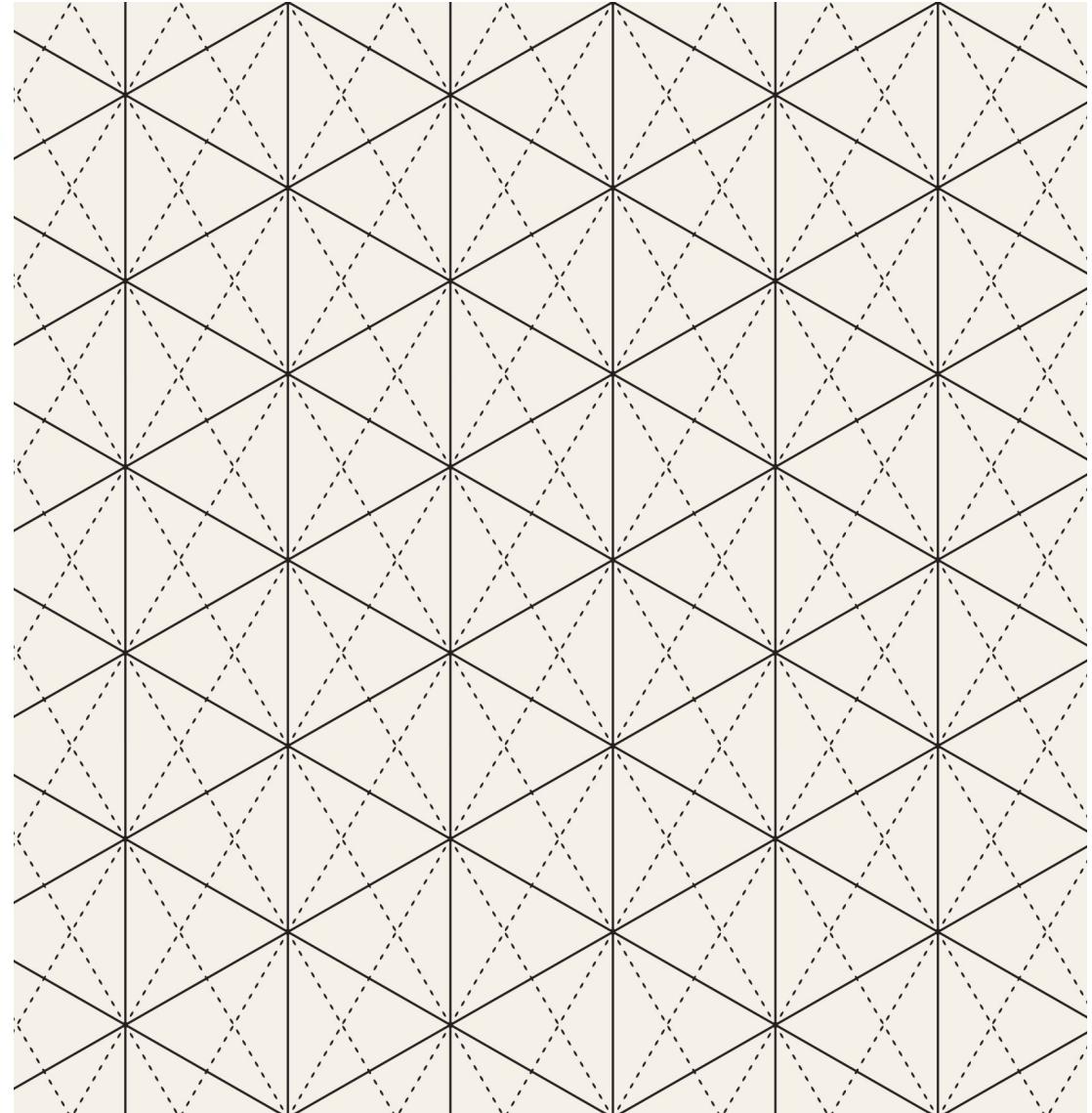
Asst. Prof. Punnarai Siricharoen, PhD

Department of Computer Engineering

Faculty of Engineering

Chulalongkorn University

[Punnarai.s@chula.ac.th](mailto:Punnarai.s@chula.ac.th)



# Outline

---

- Object Detection : Evolution
- Deep Learning Models
- Object Detection Models : YOLO
- Tracking Models : ByteTrack
- DEMOs

“Cat”



Classification



Object Detection



Semantic Segmentation



Instance Segmentation



Image Captioning



Image Generation



Classification

## Features?

**Colour** is powerful - Color Features: Color Histogram, Mean RGB, HSV, CIELAB, YCbCr color model

**Texture Features:** Local Binary Pattern, Co-occurrence matrix, Tamura's textural features, Wavelet Transform

**Statistical Features:** Mean, Standard Deviation and Variance

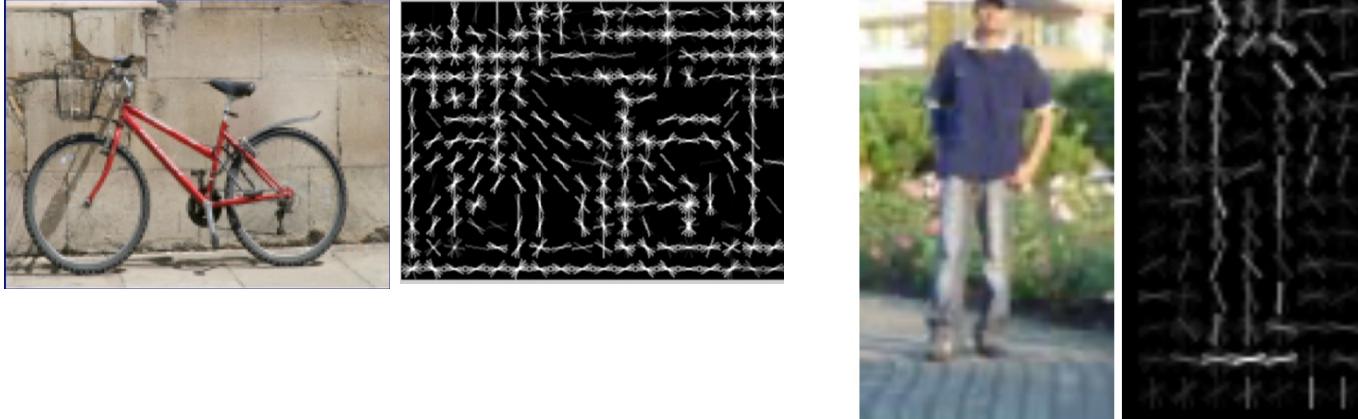
**Structural Features:** Edge, corner, shape (Region properties)

Edge detection (human-perception) - First / second derivatives – sobel, prewitt, canny

# Features – Histogram of Oriented Gradient (HOG)



Overview of HOG feature extraction for person / non-person classification (Dalal and Triggs, 2005)



Dalal and Triggs, 2005, Histograms of Oriented Gradients for Human Detection (CVPR'05)

# Features – SCALE INVARIANT FEATURE TRANSFORM (SIFT)

233x189



729

keypoints after  
gradient threshold

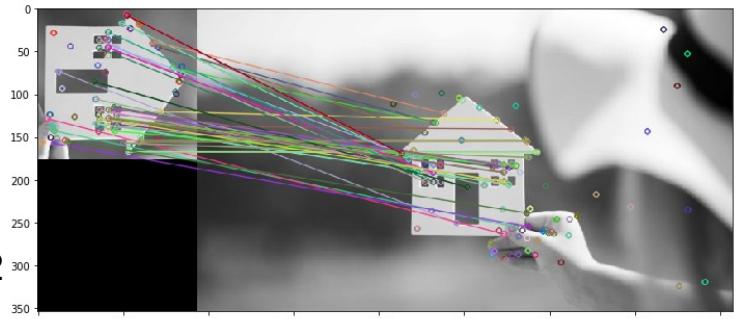


832

initial keypoints

536

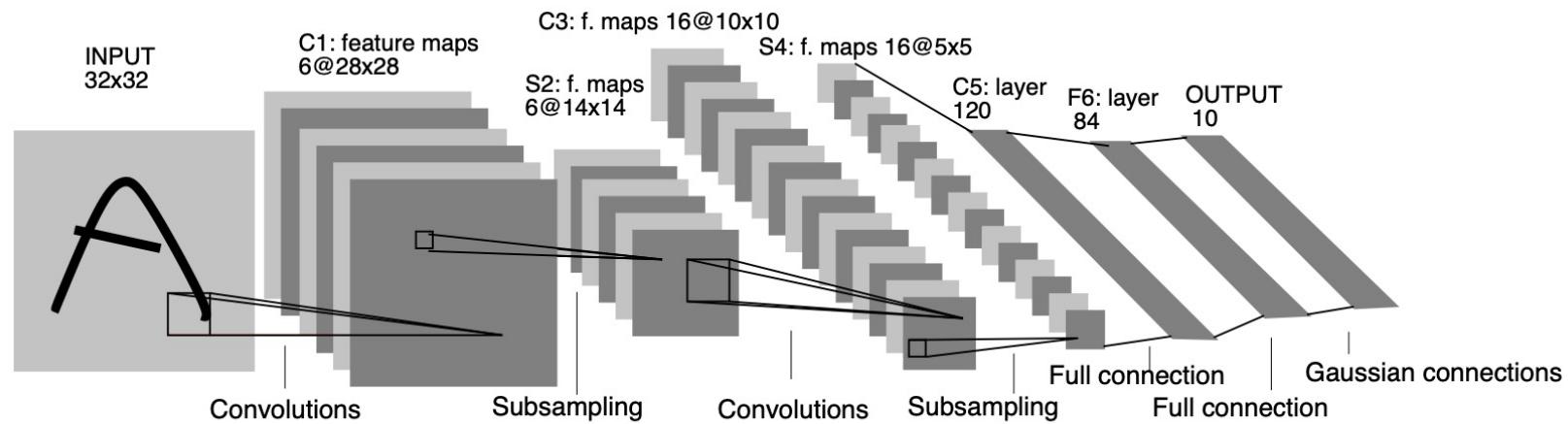
keypoints after  
ratio threshold



David Lowe, IJCV 2004

# Deep Learning Model

---



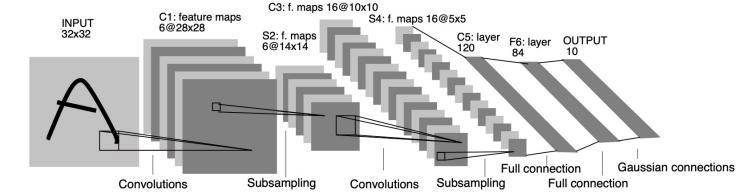
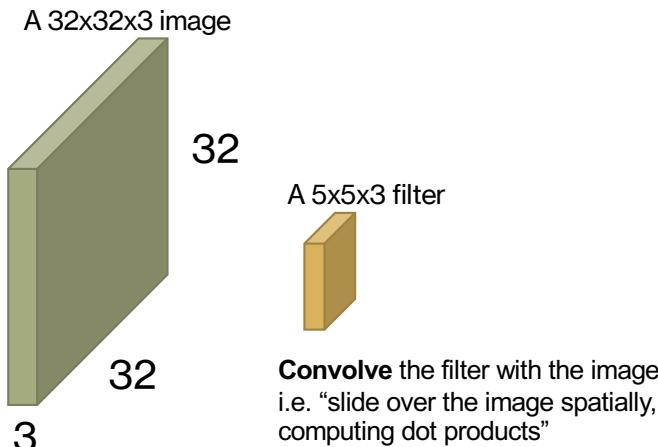
**Fig. 1.** Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

LeCun, Yann, Léon Bottou, Yoshua Bengio and Patrick Haffner. "Gradient-based learning applied to document recognition." Proc. IEEE 86 (1998): 2278-2324.

# Convolutional Neural Networks

## - Convolution Layer

- To extract spatial features from data, such as edges, textures, or patterns, by applying filters across the input.



**Fig. 1.** Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

LeCun, Yann, Léon Bottou, Yoshua Bengio and Patrick Haffner. “Gradient-based learning applied to document recognition.”  
Proc. IEEE 86 (1998): 2278-2324.

## - Filtering (Convolution)

Image $f(x,y)$				
0	0	0	0	0
0	0	0	0	0
0	0	1	0	0
0	0	0	0	0
0	0	0	0	0

kernel  $h(x,y)$

1	2	3
4	5	6
7	8	9

output  $g(x,y)$

0	0	0	0	0
0	9	8	7	0
0	6	5	4	0
0	3	2	1	0
0	0	0	0	0

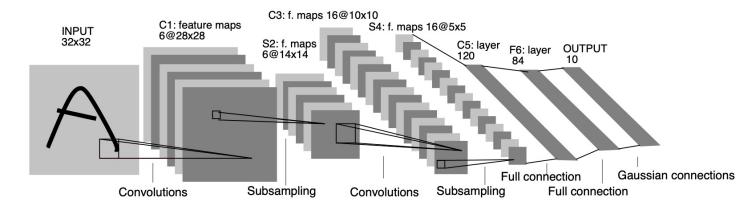
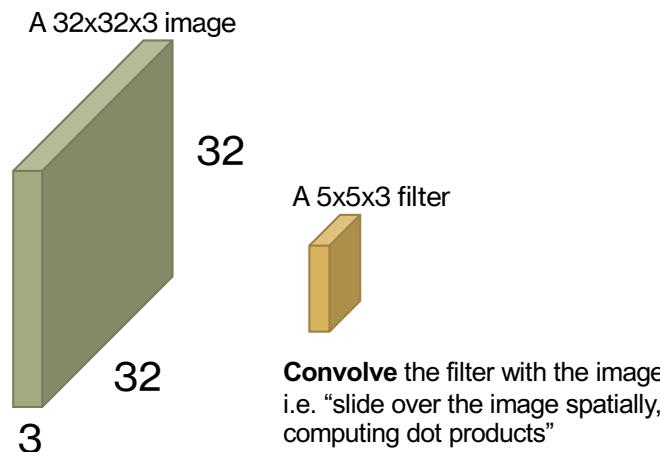
1	2	3
4	5	6
7	8	9
0	0	0
0	0	0
0	0	0
0	0	0

Sliding the kernel over the image and find sum of products between image and kernel.

# Convolutional Neural Networks

## - Convolution Layer

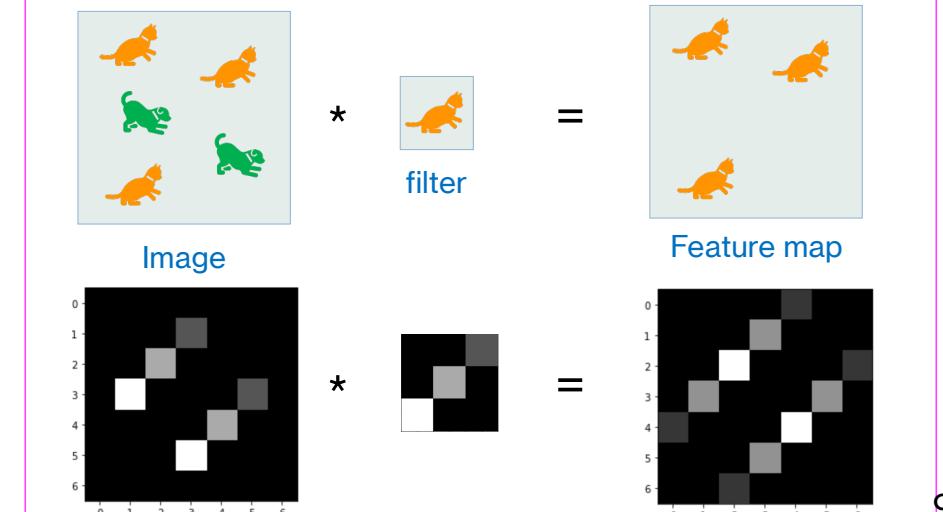
- To extract spatial features from data, such as edges, textures, or patterns, by applying filters across the input.



**Fig. 1.** Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

LeCun, Yann, Léon Bottou, Yoshua Bengio and Patrick Haffner. "Gradient-based learning applied to document recognition." Proc. IEEE 86 (1998): 2278-2324.

### - Filtering (Convolution)



# Vision Transformer

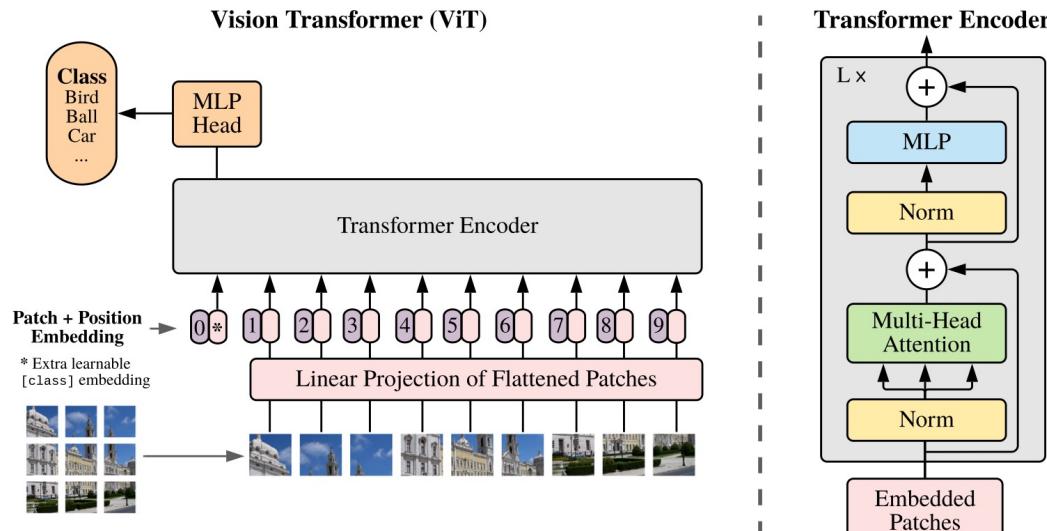
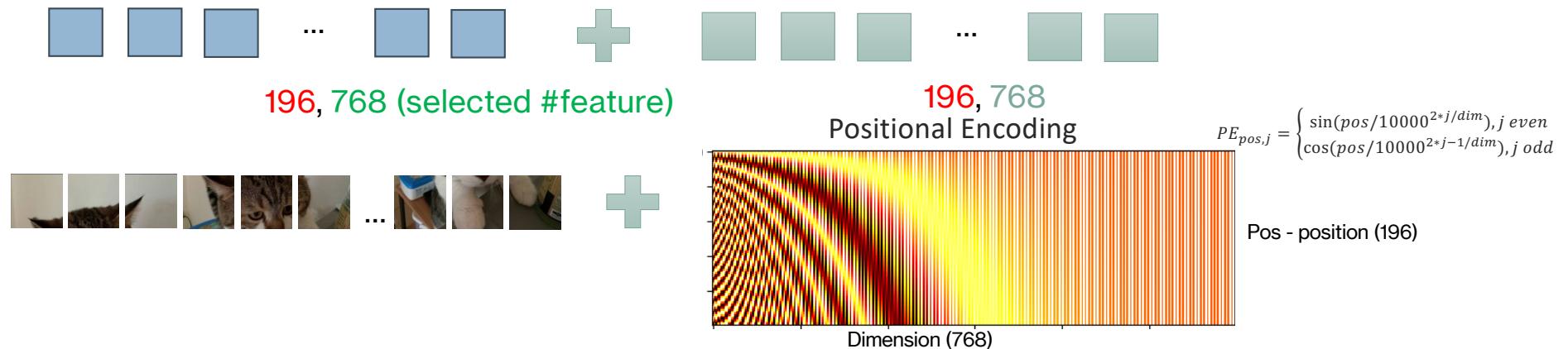


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

- For vision transformer, an input image is split into a set of image patches, called **visual tokens**.
- Position embeddings are added to the patch embeddings to retain positional information.

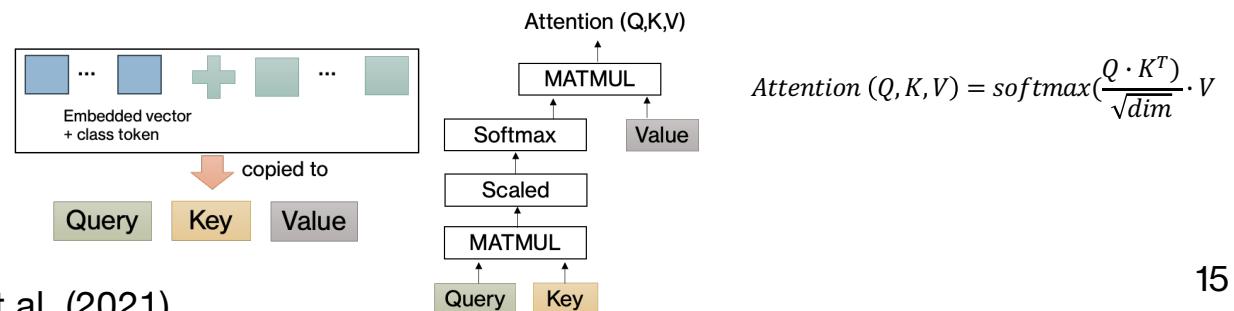
# Evolution of Deep learning: Vision Transformer (2021)

- Positional embedding – model knows where each patched would be placed.



- Transformer Encoder = Layer Norm + Multi-Head Self Attention

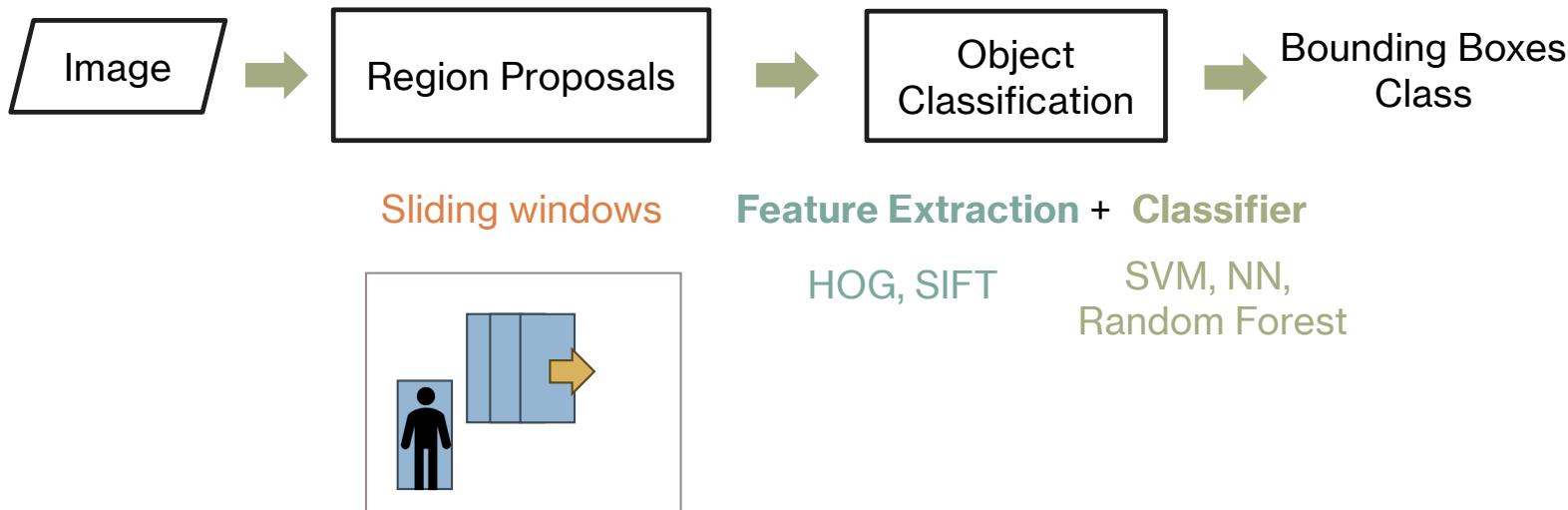
- Multi-Head Self-Attention



Attention is all you need, Vaswani A., et al. (2021)

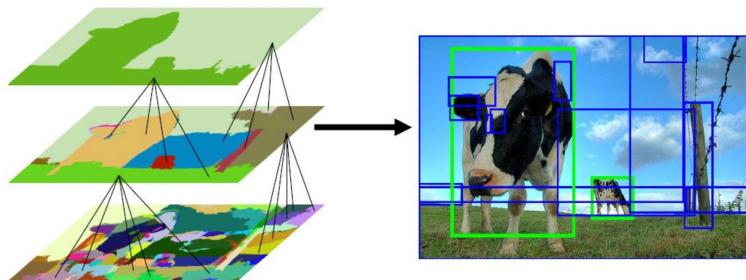
# — Object Detection Models

- Two-stage Detectors 包含 Faster-RCNN



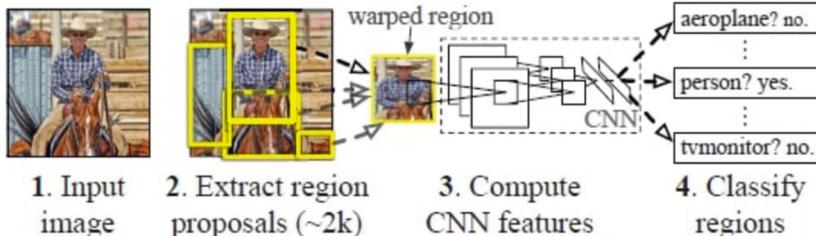
# Object Detection Models

- Two-stage Detectors یعنی Faster-RCNN



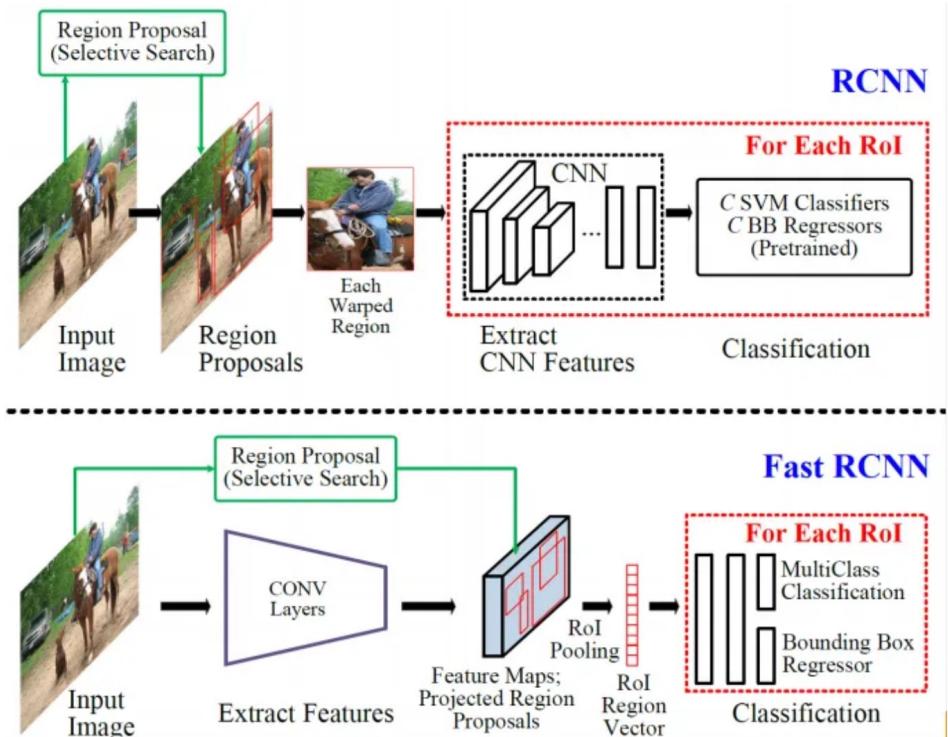
Selective Search for Object Recognition

J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, A. W. M. Smeulders  
In International Journal of Computer Vision 2013.



Architecture of RCNN

Rich feature hierarchies for accurate object detection and semantic segmentation  
Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik



RCNN vs Fast-RCNN

(source: Deep Learning for Generic Object Detection: A Survey)

# Object Detection Models

- Two-stage Detectors  $\hookrightarrow$  Faster-RCNN

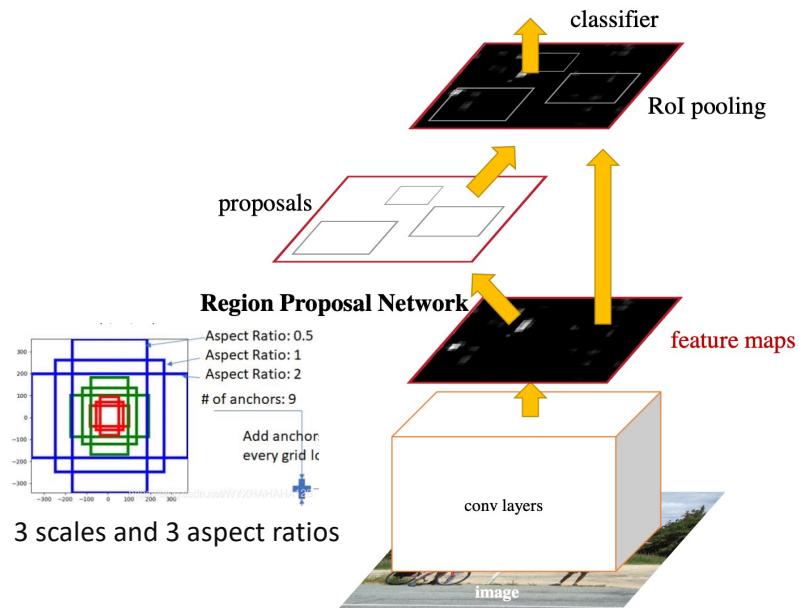


Figure 2: Faster R-CNN is a single, unified network

Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks  
Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun

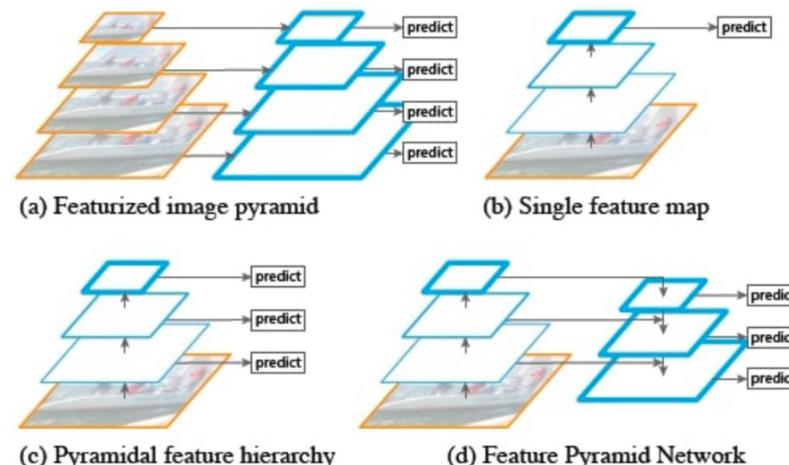


Fig. 5. Feature Pyramid Network [28]. (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

Feature Pyramid Networks for Object Detection  
Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, Serge Belongie

# Object Detection Models

- Two-stage Detectors เช่น Faster-RCNN



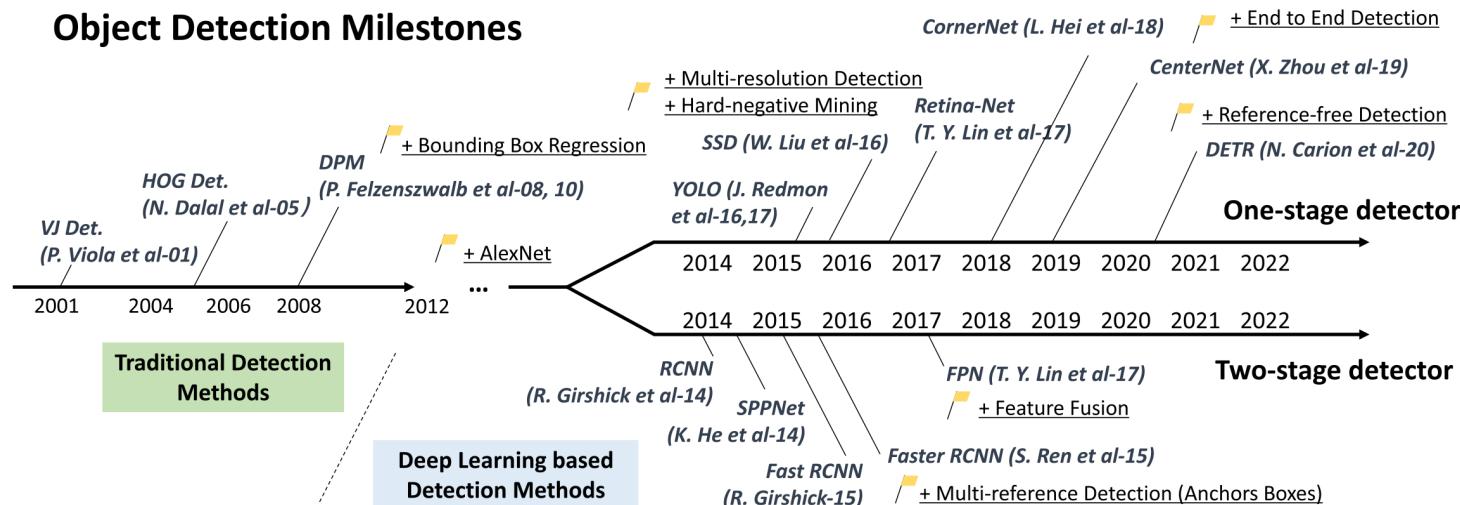
- แม่นยำสูง แต่ค่อนข้างช้า
- One-stage Detectors เช่น YOLO, DETR



- เร็วมาก เหมาะกับการใช้งานแบบเรียลไทม์ เช่นในมือถือหรือกล้อง

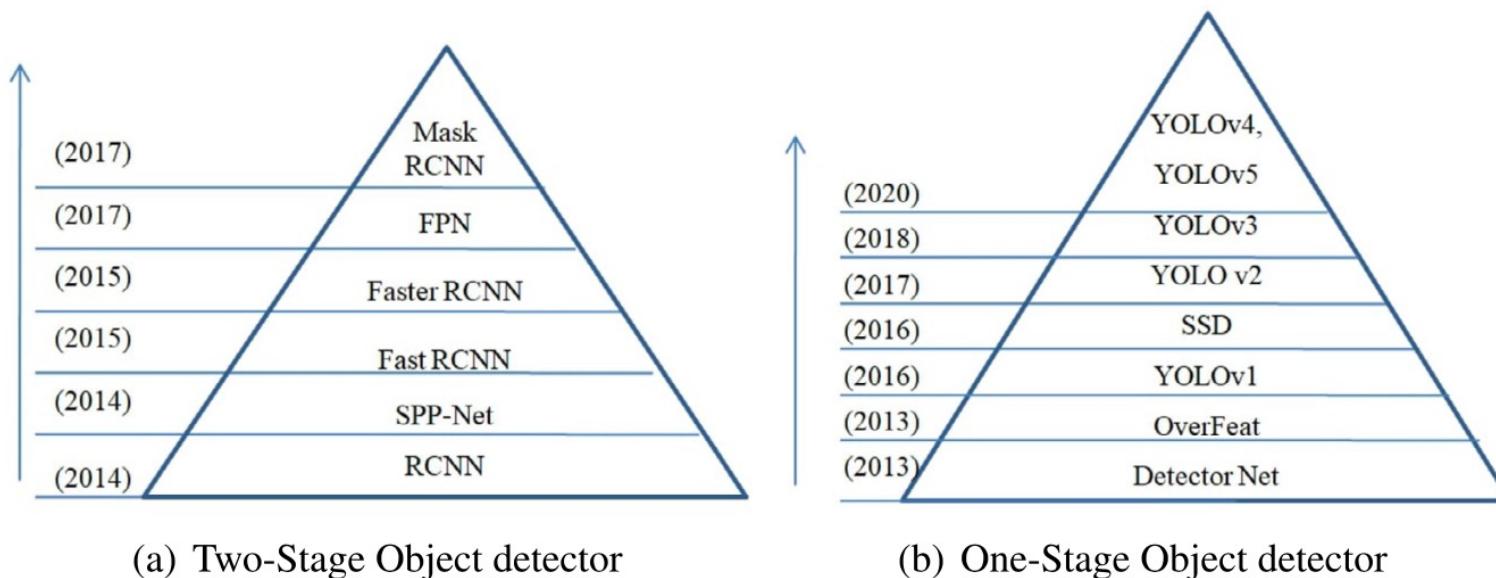
# Object Detection

Zou et al.: Object Detection in 20 Years: A Survey



**Fig. 2.** Road map of object detection. Milestone detectors in this figure: VJ Det. [10], [11], HOG Det. [12], DPM [13], [14], [15], RCNN [16], SPPNet [17], Fast RCNN [18], Faster RCNN [19], YOLO [20], [21], [22], SSD [23], FPN [24], Retina-Net [25], CornerNet [26], CenterNet [27], and DETR [28].

# Object Detection



Kaur & Singh (2023), A comprehensive review of object detection with deep learning, Digital Signal Processing

# — COCO Datasets

- <https://cocodataset.org>

## What is COCO?



COCO is a large-scale object detection, segmentation, and captioning dataset. COCO has several features:

- ✓ Object segmentation
- ✓ Recognition in context
- ✓ Superpixel stuff segmentation
- ✓ 330K images (>200K labeled)
- ✓ 1.5 million object instances
- ✓ 80 object categories
- ✓ 91 stuff categories
- ✓ 5 captions per image
- ✓ 250,000 people with keypoints

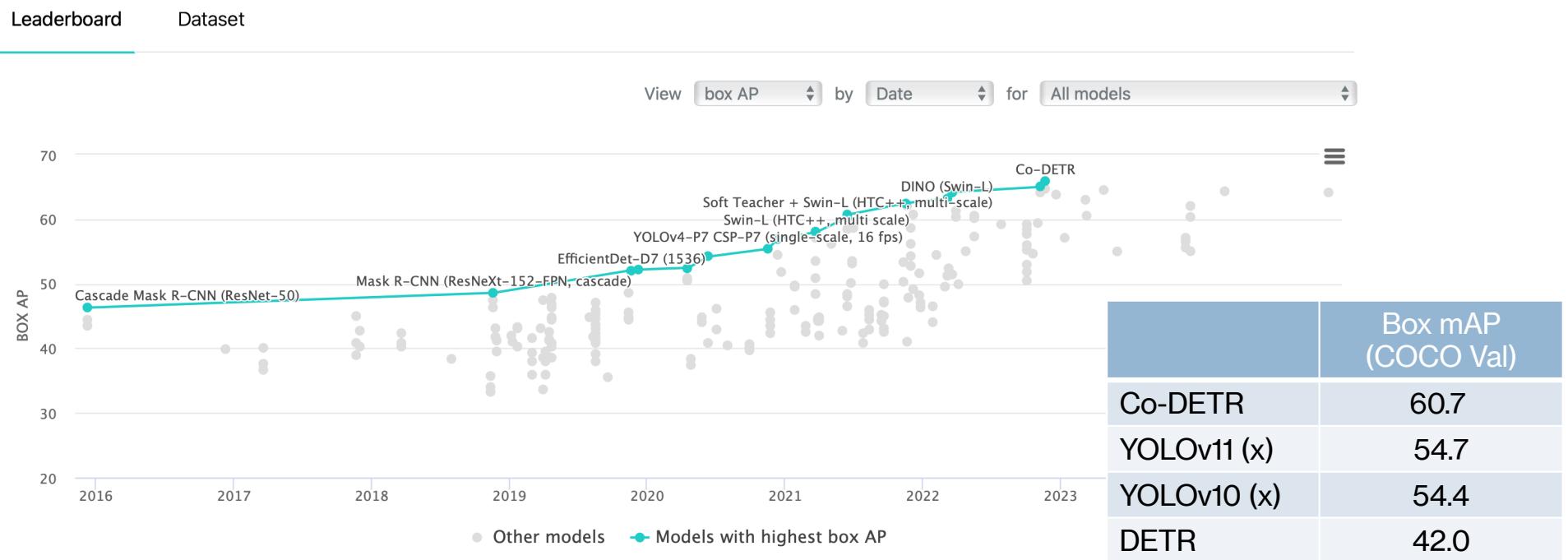
## COCO Explorer

COCO 2017 train/val browser (123,287 images, 886,284 instances). Crowd labels not shown.



# DETR / CO-DETR

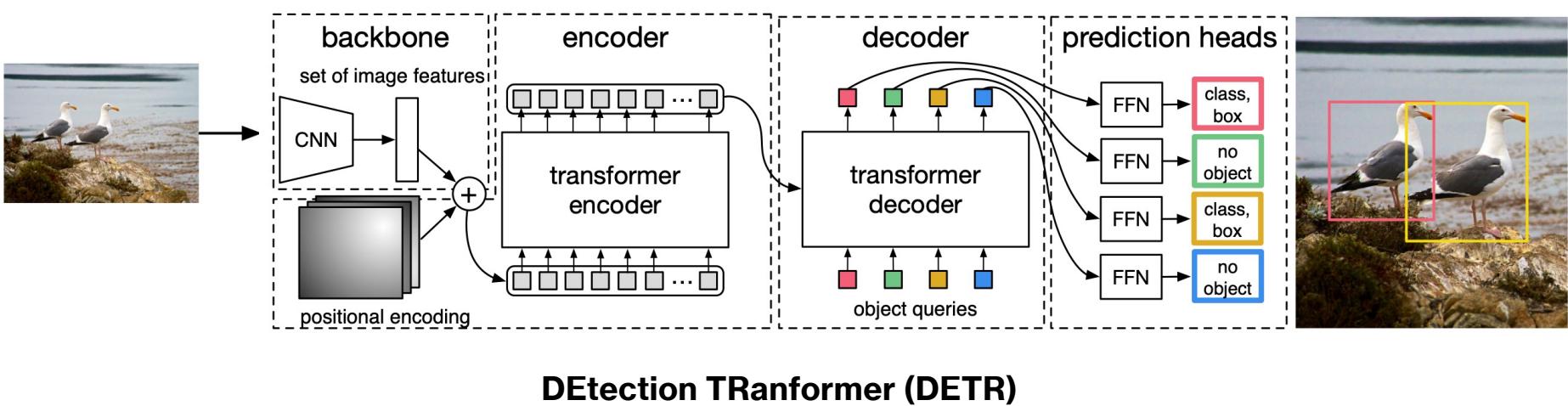
## Object Detection on COCO minival



<https://paperswithcode.com>

23

# Object Detection : DEtection TTransformer (DETR)



N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-End Object Detection with Transformers,” May 2020, [Online]. Available: <http://arxiv.org/abs/2005.12872>

# DETR / CO-DETR

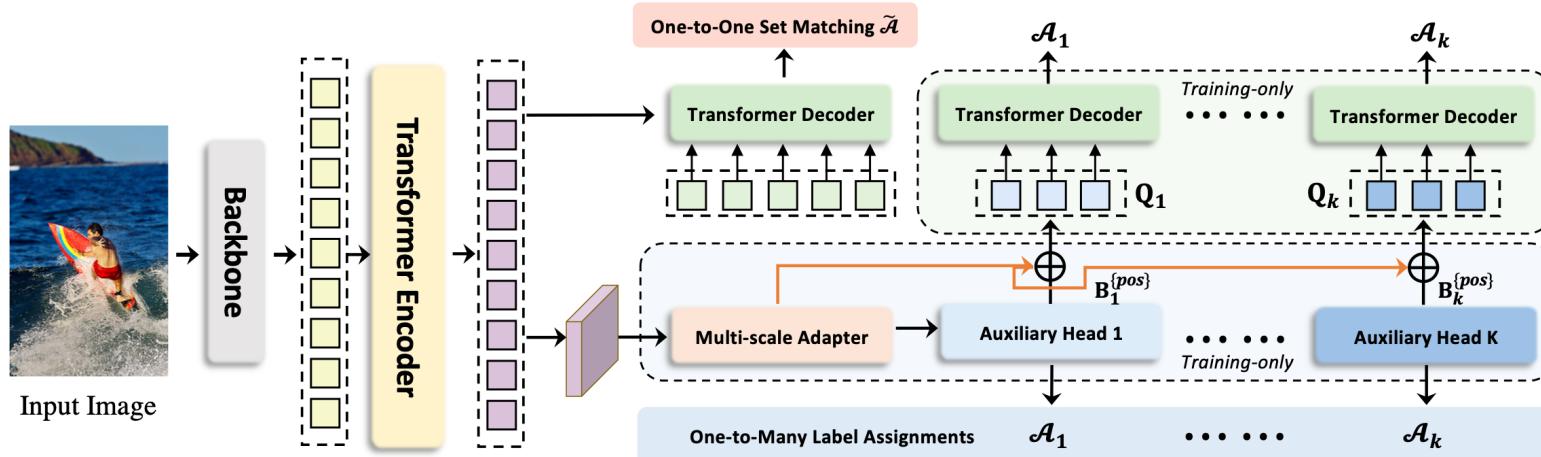
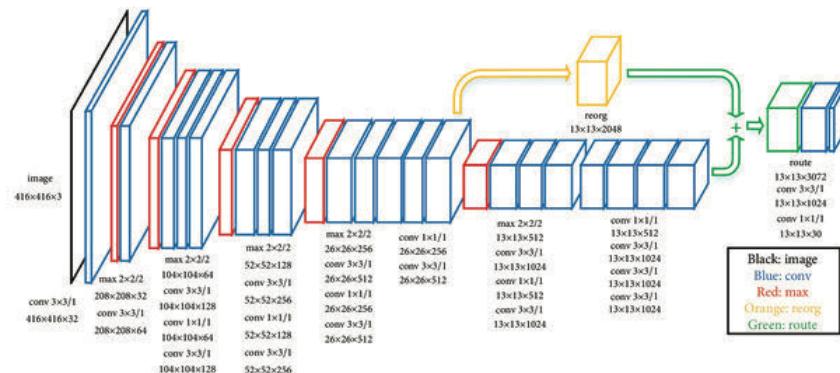


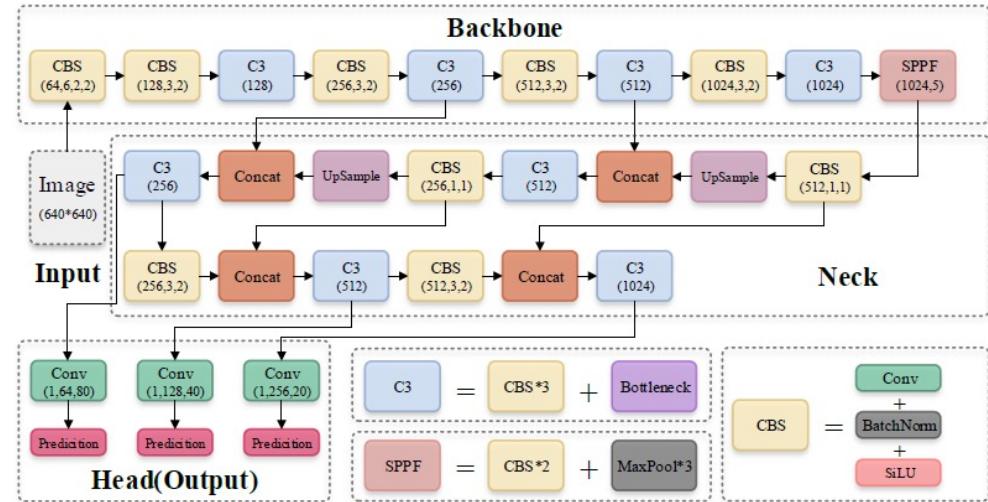
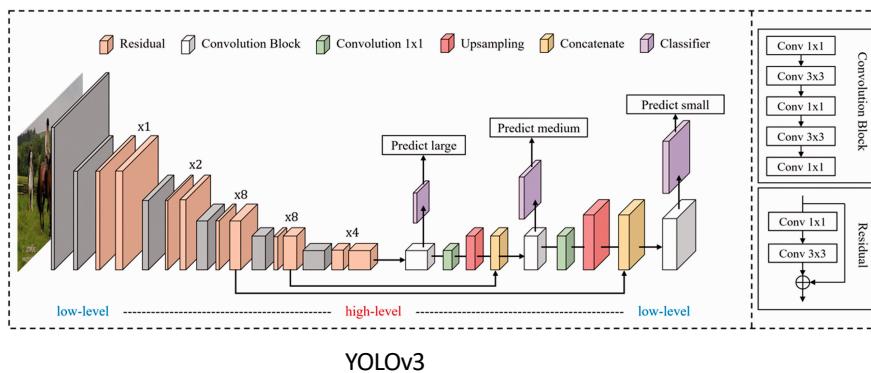
Figure 4. **Framework of our Collaborative Hybrid Assignment Training.** The auxiliary branches are discarded during evaluation.

DETRs with Collaborative Hybrid Assignments Training  
Zhuofan Zong Guanglu Song Yu Liu\* SenseTime Research

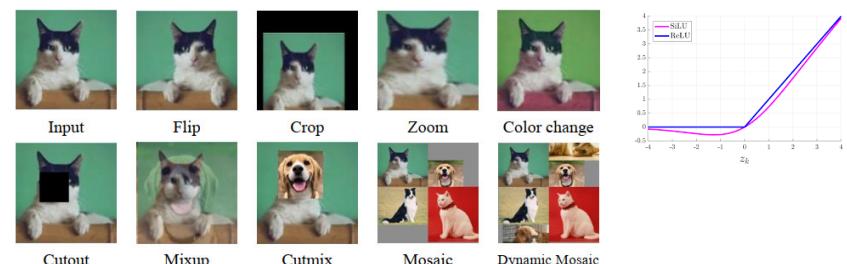
# You Only Look Once (YOLO)



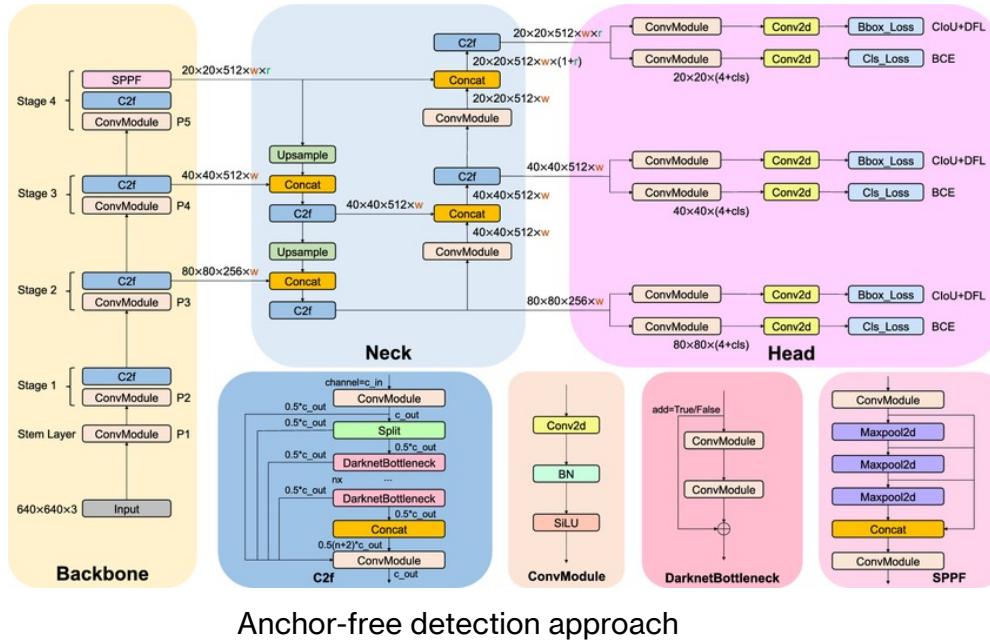
An Efficient Pedestrian Detection Method Based on YOLOv2 - Scientific Figure on ResearchGate. Available from: [https://www.researchgate.net/figure/The-network-architecture-of-YOLOv2\\_fig2\\_330029907](https://www.researchgate.net/figure/The-network-architecture-of-YOLOv2_fig2_330029907) [accessed 8 Nov 2024]



<https://sh-tsang.medium.com/brief-review-yolov5-for-object-detection-84cc6c6a0e3a>



# You Only Look Once (YOLO)

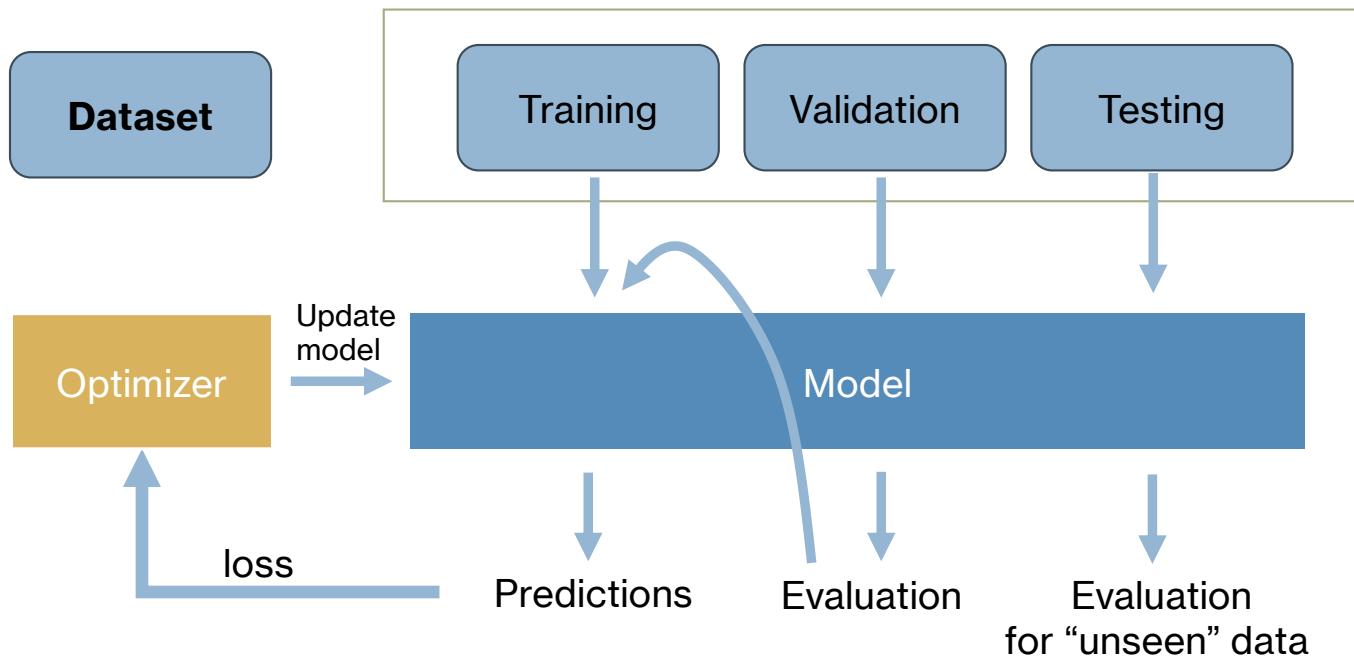


Implementation of Slicing Aided Hyper Inference (SAHI) in YOLOv8 to Counting Oil Palm Trees Using High-Resolution Aerial Imagery Data - Scientific Figure on ResearchGate. Available from: [https://www.researchgate.net/figure/YOLOv8-model-architecture-20\\_fig2\\_382742087](https://www.researchgate.net/figure/YOLOv8-model-architecture-20_fig2_382742087)

- YOLOv8-11
  - Detection, Instance Segmentation, Keypoint Detection
- YOLOv11
  - Better Handling Small Objects
- YOLOv12
  - Attention-Centric Architecture (maintain speed)

# Training Your Model

---



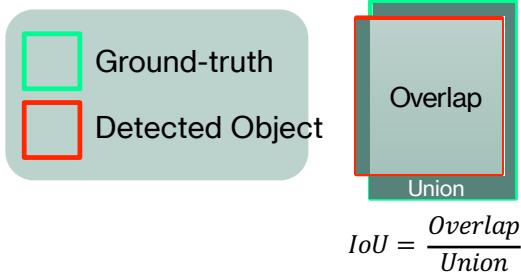
# Loss functions

---

	Two-stage	One-stage
Stage-wise loss separation	✓	
Focus on high precision	✓	focus on speed
Has proposal + refinement	✓	
Detailed loss components	✓	✓ (single pass)

- Two-stage model loss
  - > Region proposal network Loss (Classification Loss + Regression Loss)
  - > Roi Head Loss (Classification Loss + Regression Loss)
- One-stage model loss for a single pass
  - > Classification Loss + Objectness Loss + Localization Loss
  - > YOLOv10-11 – No Objectness Loss, Classification Loss + Localization Loss (Bbox + DFL)

# Evaluation for Object Detection



**mAP** – Mean Average Precision

**mAP50** – The prediction is correct, when  $\text{IoU} \geq 0.5$ .

- 1) Rank the confidence score of the detected objects in the dataset
- 2) Calculate Precision / Recall for each of the object

Rank	IoU	Correct?	Precision	Recall
1	0.7	True	1.0	0.25
2	0.6	True	1.0	0.50 ↑
3	0.4	False	0.67 ↓	0.5
4	0.4	False	0.5 ↓	0.5
5	0.5	True	0.6 ↑	0.75↑
6	0.6	True	0.67 ↑	1.0↑

Rank 1:

$$\text{precision} = \text{TP}/(\text{TP}+\text{FP}) = 1/(1+0) \\ \text{Recall} = \text{TP}/(\text{ALL GTs}) = 1/4 = 0.25$$

Rank 3:

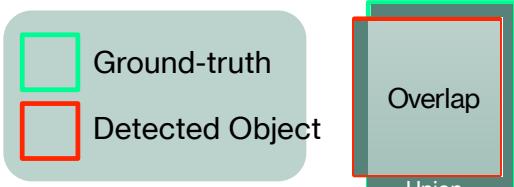
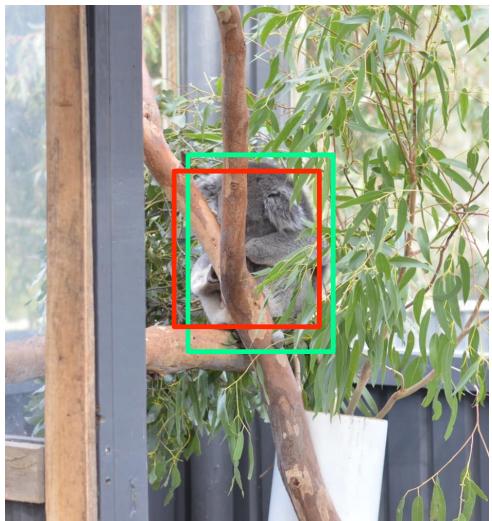
$$\text{Precision} = 2/3 = 0.67 \\ \text{Recall} = 2/4 = 0.5$$

Rank 5:

$$\text{Precision} = 3/5 = 0.6 \\ \text{Recall} = 3/4 = 0.75$$

# Evaluation for Object Detection

---

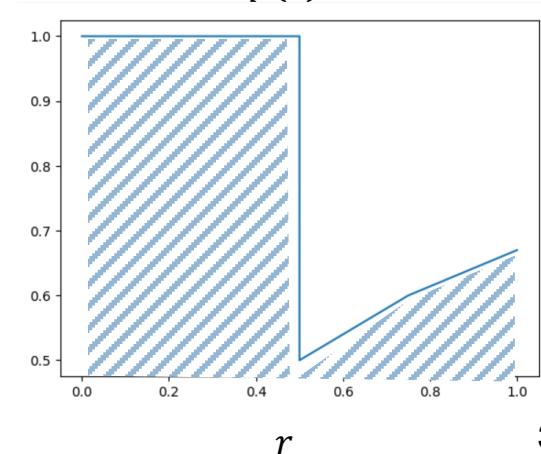


$$IoU = \frac{\text{Overlap}}{\text{Union}}$$

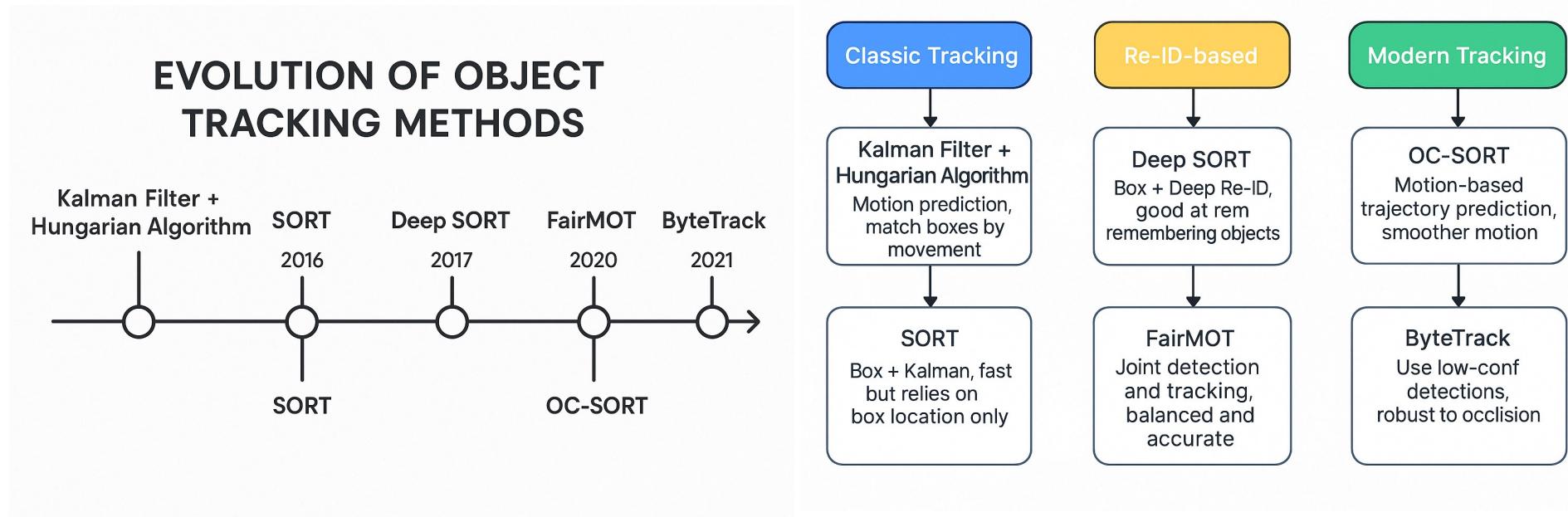
**mAP** – Mean Average Precision

**mAP50** – The prediction is correct, when  $\text{IoU} \geq 0.5$ .

- 1) Rank the confidence score of the detected objects.
- 2) Calculate Precision / Recall for each of the object
- 3) Plot PR curve  $p(r)$
- 4) AP for each class  $= \int_0^1 p(r)dr$ 
  - 11-point interpolation
  - mAP - average over classes
  - mAP50:90 – average over 10 IoUs



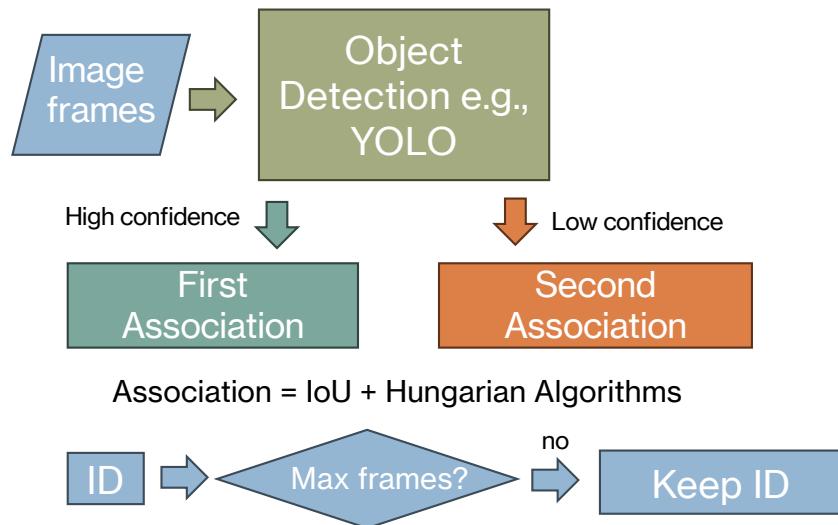
# Object Tracking



MOTA (Multiple Object Tracking Accuracy)

$$\text{MOTA} = 1 - \frac{\text{FP} + \text{NP} + \text{IDSW}}{\text{GT}}$$

# Object Tracking : ByteTrack



- track\_high\_thresh: Minimum score for high-confidence detection (first matching)
- track\_low\_thresh: Minimum score for low-confidence detection (second matching)
- match\_thresh: Minimum IoU to associate detection with a track
- track\_buffer: Max frames to keep unmatched tracks before deletion
- min\_box\_area: Ignore boxes smaller than this area
- mot20: Use settings optimized for dense scenes
- fuse\_score: Final score = objectness  $\times$  class confidence



(a) detection boxes

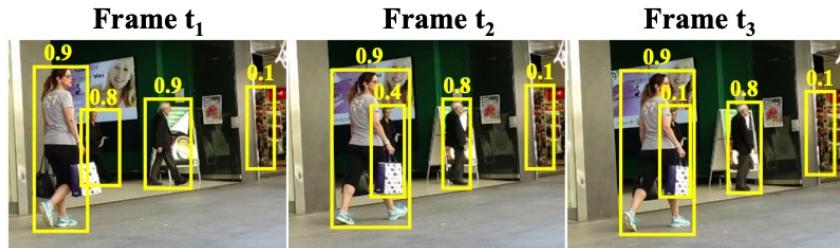


(b) tracklets by associating high score detection boxes



(c) tracklets by associating every detection box

# Object Tracking : ByteTrack



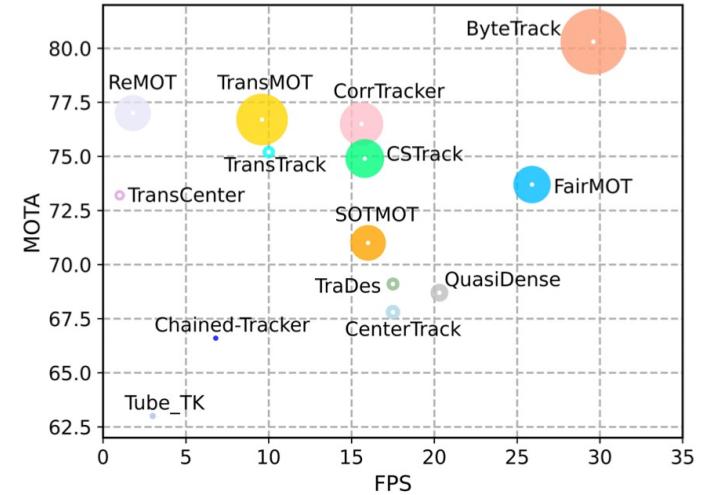
(a) detection boxes



(b) tracklets by associating high score detection boxes



(c) tracklets by associating every detection box



MOTA (Multiple Object Tracking Accuracy)

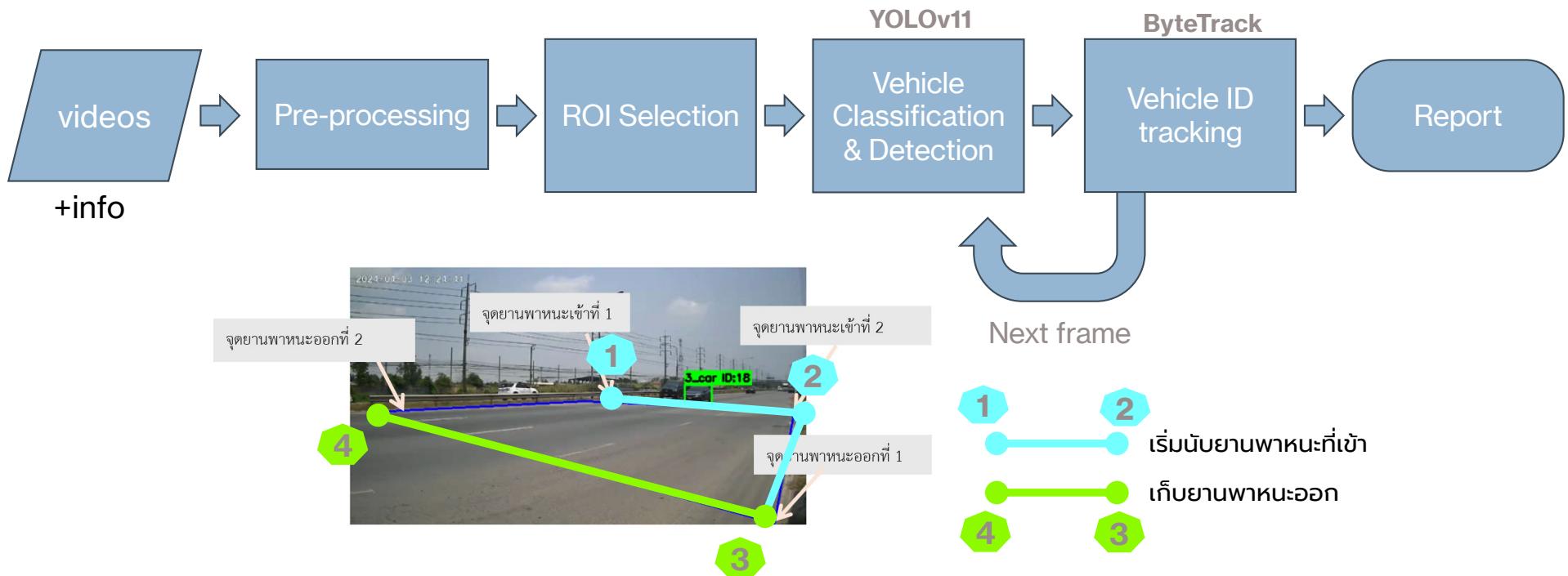
$$\text{MOTA} = 1 - \frac{\text{FP} + \text{NP} + \text{IDSW}}{\text{GT}}$$

# DEMOs

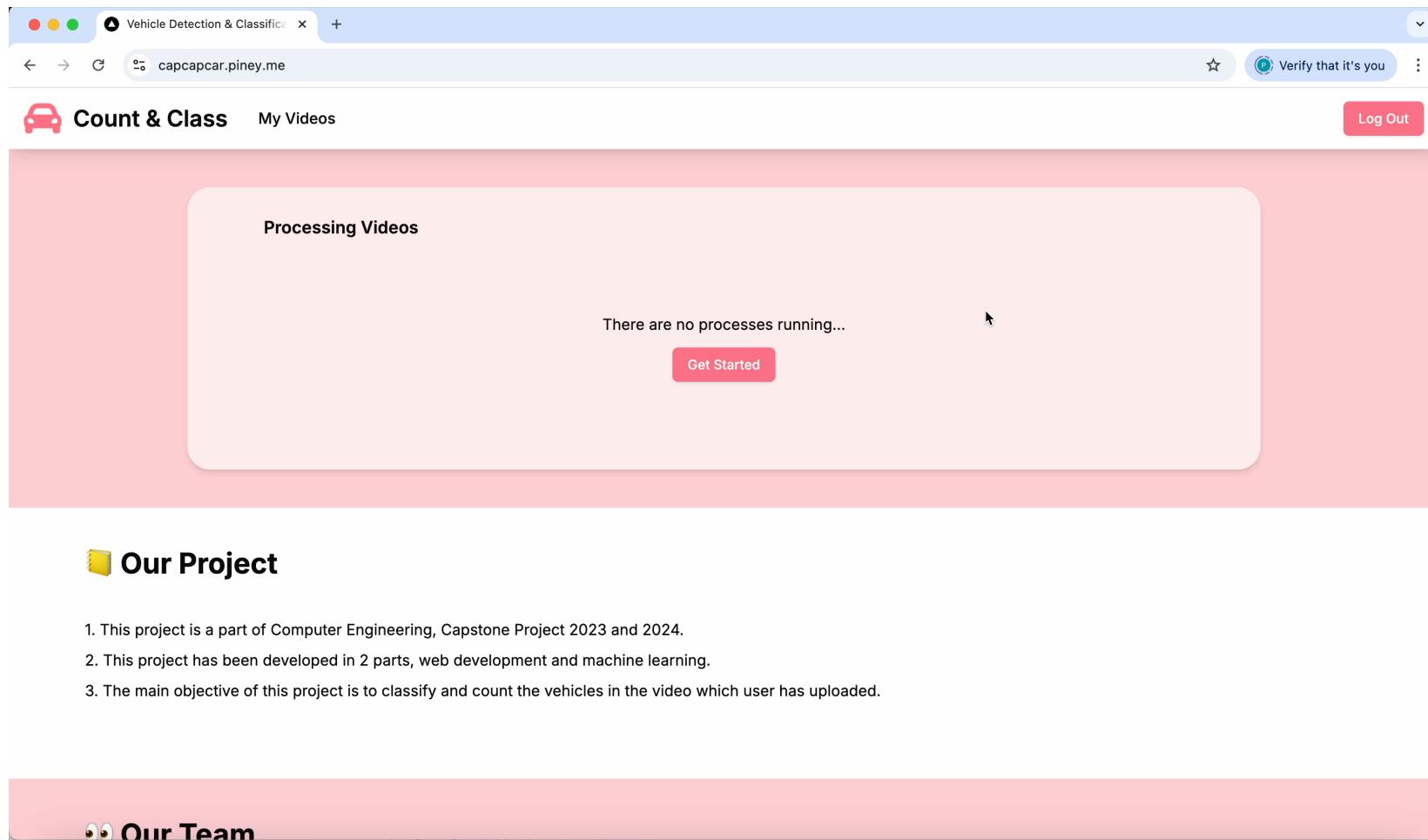
---

- **Demo 1 : YOLO - Cats**
  - <https://colab.research.google.com/drive/11YoWKLkZUNHmJlIAPfVajyJ6o6R6znXf?usp=sharing>
- **Demo 2 : YOLO - ByteTrack**
  - <https://colab.research.google.com/drive/1o5fnniqV1FVKkGvLD2u0zkWK8EHcfNUq?usp=sharing>
- **Demo 3 : YOLO - ByteTrack - Vehicle AI Tracker**
  - <https://colab.research.google.com/drive/1ndn1gPqgyG4TToluezKBDiU0tYZGkf3?usp=sharing>
- **Demo 4 : Vehicle AI Tracker : CapCapCar Year 2**
  - <https://capcapcar.piney.me>
  - Video ตัวอย่าง เลือกสั้น ๆ : <https://youtu.be/Y1jTEyb3wil?si=7u5cbXyLqOjRA9BQ>

# VEHICLE AI TRACKER



# DEMO : VEHICLE AI TRACKER



The screenshot shows a web browser window titled "Vehicle Detection & Classification". The URL is "capcapcar.piney.me". The main content area has a pink header with the text "Count & Class" and "My Videos". A red button in the top right corner says "Log Out". Below the header, a large white box contains the text "Processing Videos" at the top, followed by "There are no processes running..." in the center, and a red "Get Started" button at the bottom. A cursor arrow is visible near the "Get Started" button.

**Our Project**

1. This project is a part of Computer Engineering, Capstone Project 2023 and 2024.
2. This project has been developed in 2 parts, web development and machine learning.
3. The main objective of this project is to classify and count the vehicles in the video which user has uploaded.

**Our Team**