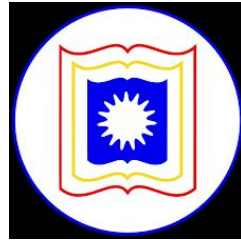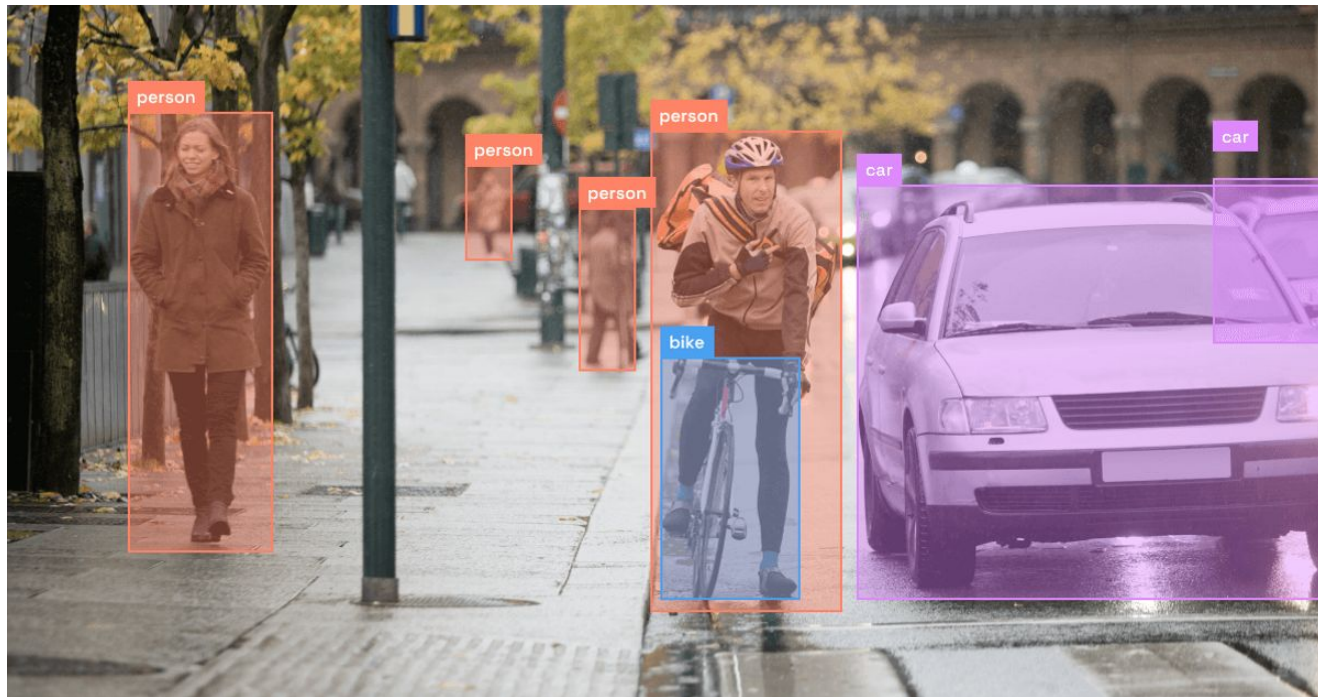# CSE4261: Neural Network and Deep Learning

Lecture: 25.06.2025

Sangeeta Biswas, Ph.D.
Associate Professor,
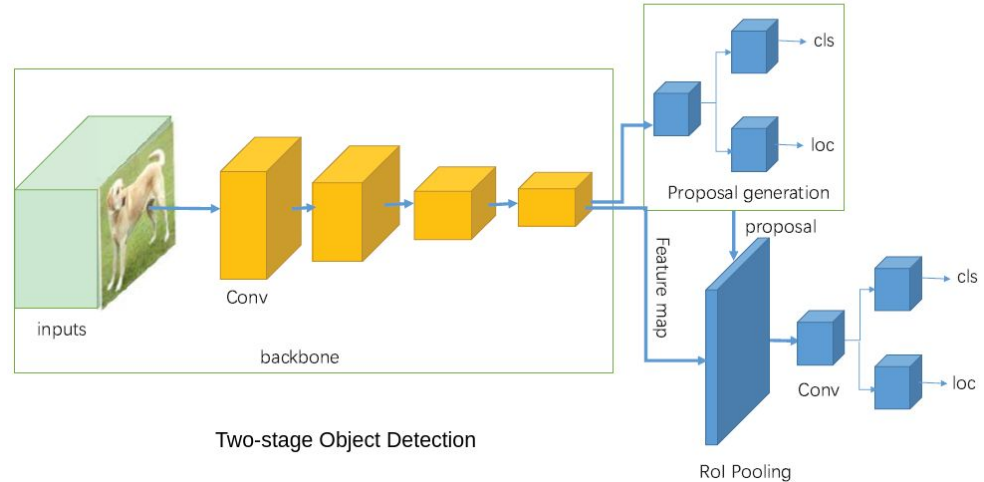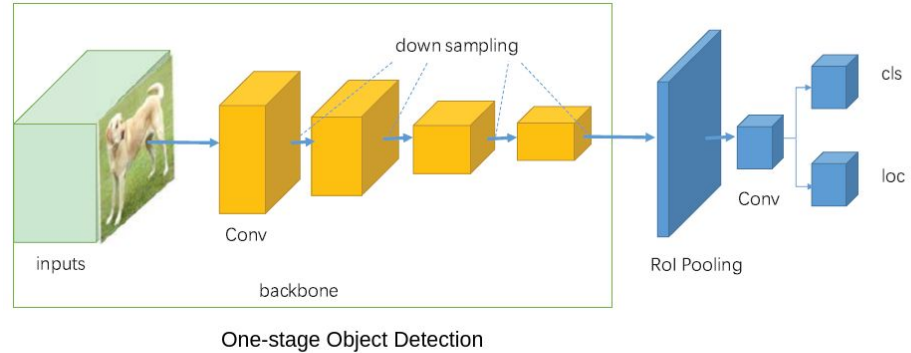University of Rajshahi, Rajshahi-6205, Bangladesh
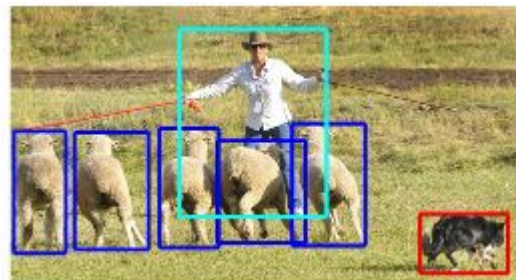
# Object Detection

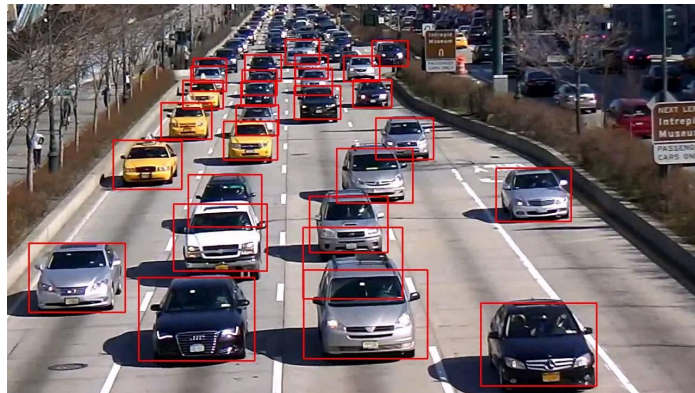# Types of Object Detector

Depending on how many times input image is seen by the object detector system, it is divided into two categories:

1. One Stage Object Detector:
   a. YOLO, SSD
2. Two-Stage Object Detector:
   a. R-CNN, Fast R-CNN, Faster R-CNN



One-stage Object Detection



Two-stage Object Detection

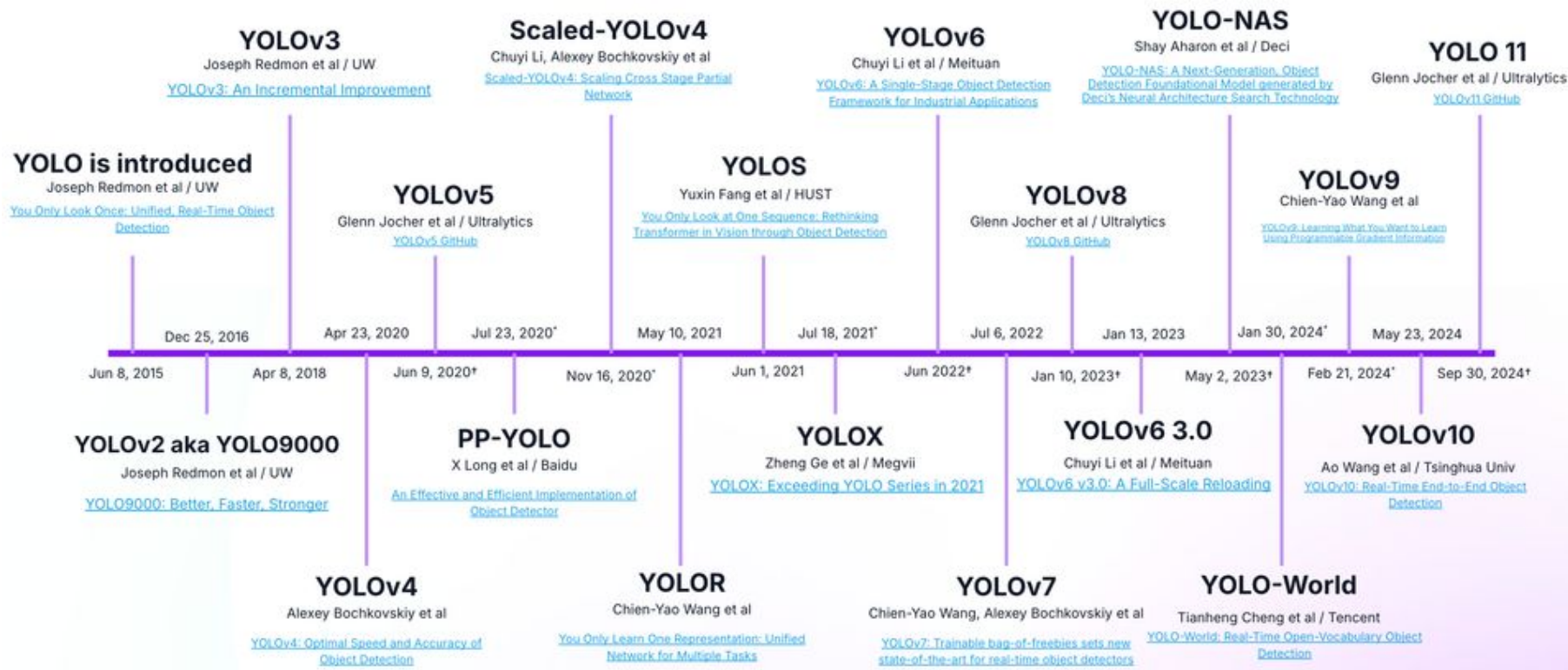# One-Stage vs. Two-Stage Object Detector

- One-stage object detection is better suited for speedy real-time applications



- Two-stage object detection is better for applications where accuracy is more important than speed.
- Two-stage is more computationally expensive than one-stage.

# YOLO

- **Y**ou **O**nly **L**ook **O**nce
- It is a one-stage object detector
- It is the most popular **real-time** object detection system.
- It processes images in one time. Therefore, it is faster than two-stage object detection systems such as Faster R-CNN.
  - YOLOv8 has a GPU latency of 1.3ms, while Faster R-CNN has a GPU latency of 54ms.
  - YOLOv3 can operate in real time with a frames per second (FPS) that is more than eight times faster than Faster R-CNN.
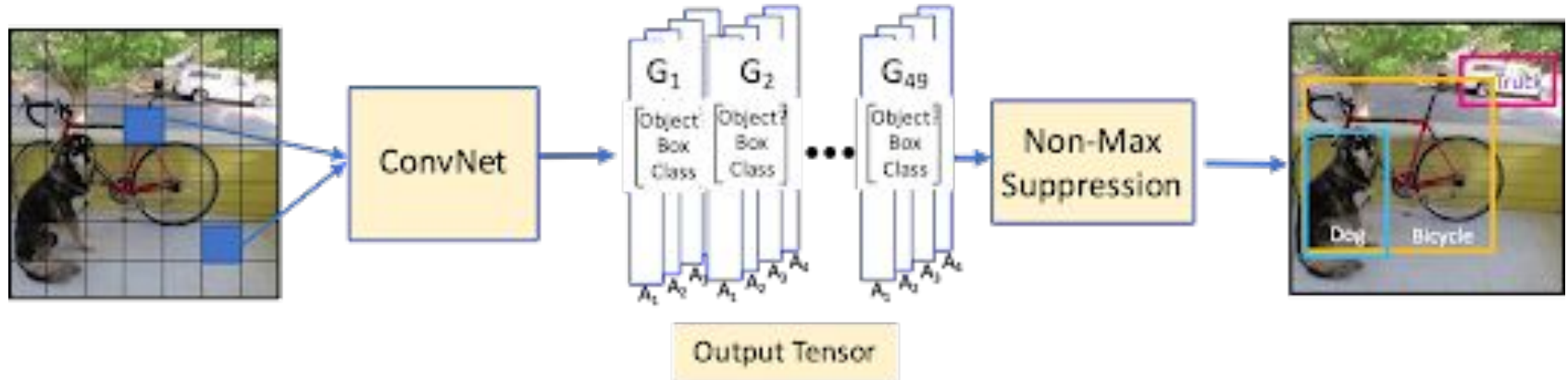
# Different Versions of YOLO



On Feb, 2025, YOLOv12: https://docs.ultralytics.com/models/yolo12/

# How YOLO Works

- Instead of using two separate stages for box prediction and category prediction like R-CNN, YOLO employs a single network to simultaneously predict bounding boxes and object categories in a single stage.

# How YOLO WOrks

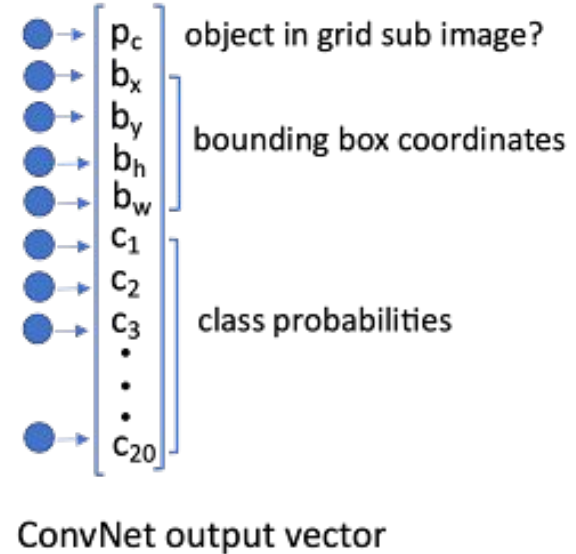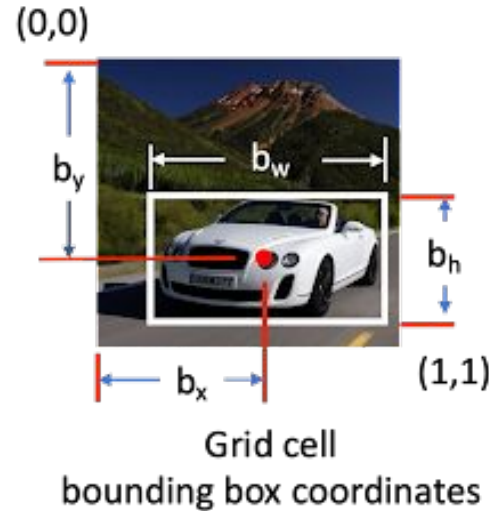YOLO follows the following steps:

- Step-1: Divide the original image into NxN grid cells of equal shape
- Step-2: Determine the bounding boxes corresponding to all the objects in the image.
- Step-3: Keep relevant grid-box candidates based on Intersection over Union (IoU) value
- Step-4: Keep only the boxes with the highest probability detection score using Non-Max Suppression (NMS) technique

# Output Vector of YOLO Model

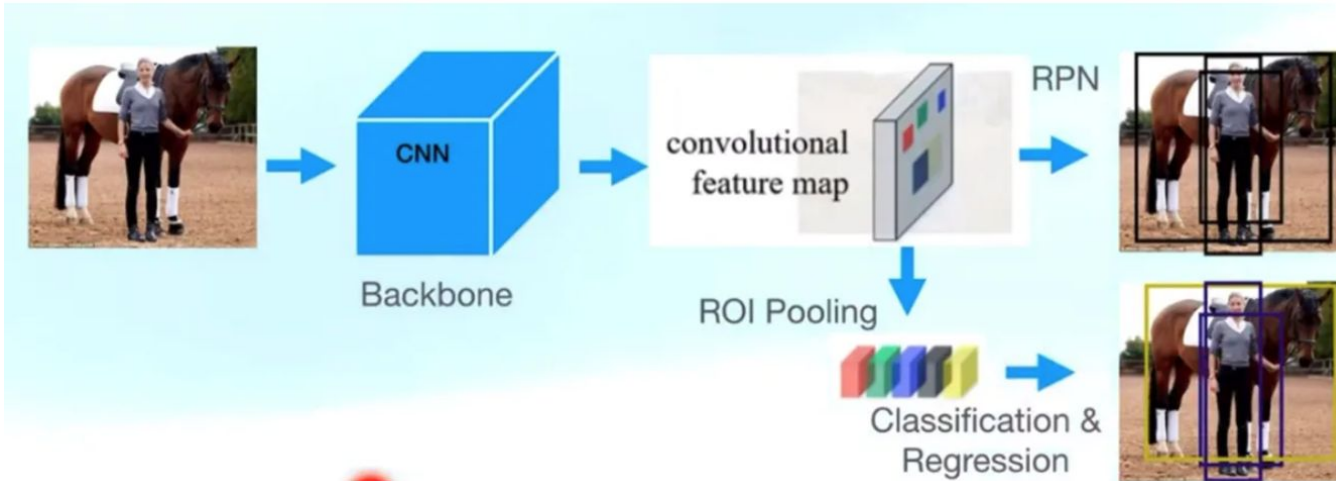size_of_output_vector = (5 + size_of_one_hot_vector) * number of objects per grid

Here,
size_of_output_vect
or= (5+20)*1 = 25



Grid cell
bounding box coordinates

ConvNet output vector

# Two Stage Object Detector

Stage-1: Generate bounding boxes to locate potential objects (i.e., region proposals) in an image

Stage-2: Classify bounding boxes to determine the object's class.

# Cons of R-CNN

- Feature extraction and classification need to be done for 2000 region proposals, which are computationally intensive and time-consuming
- Cannot be implemented for large datasets or real-time scenarios.

If a surveillance camera that needs to quickly identify multiple objects, then R-CNN or other two-stage approaches may struggle to provide fast results due to its sequential nature, potentially causing delays in critical applications.