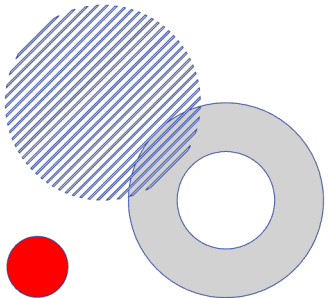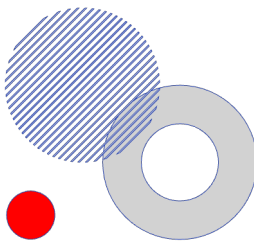# Real-Time Object Detection
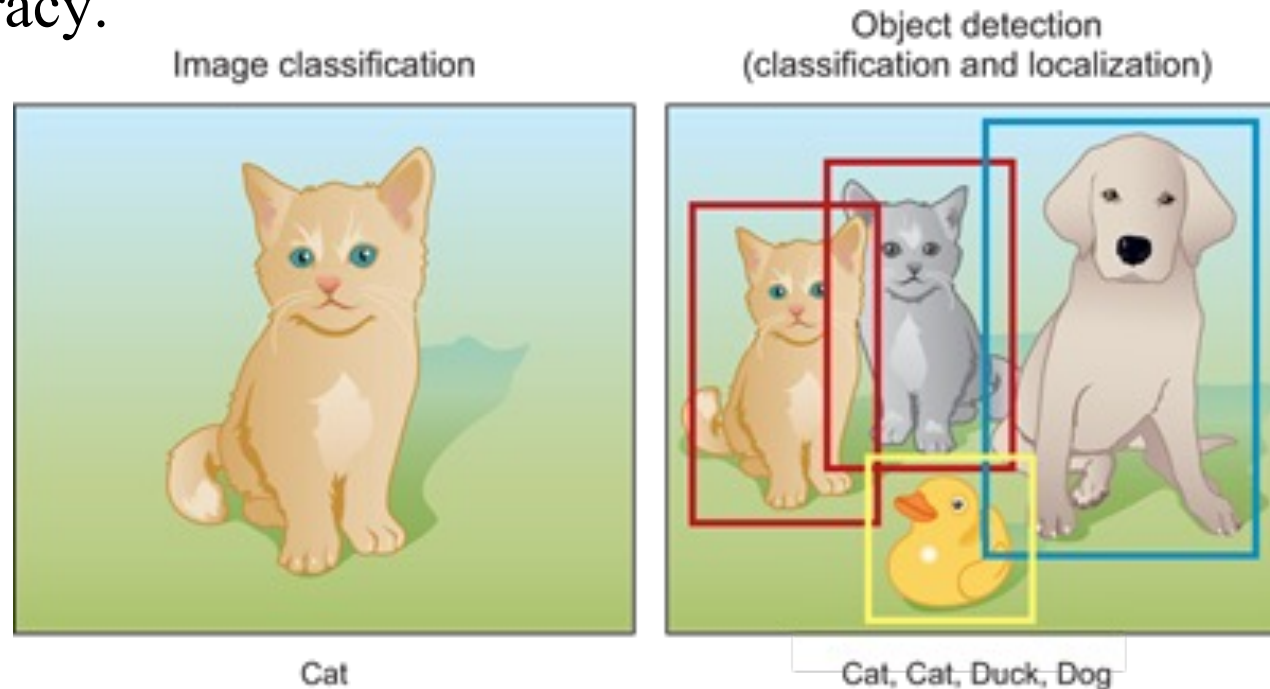
Lecturer: Dr. Thittaporn Ganokratanaa
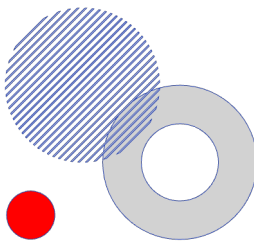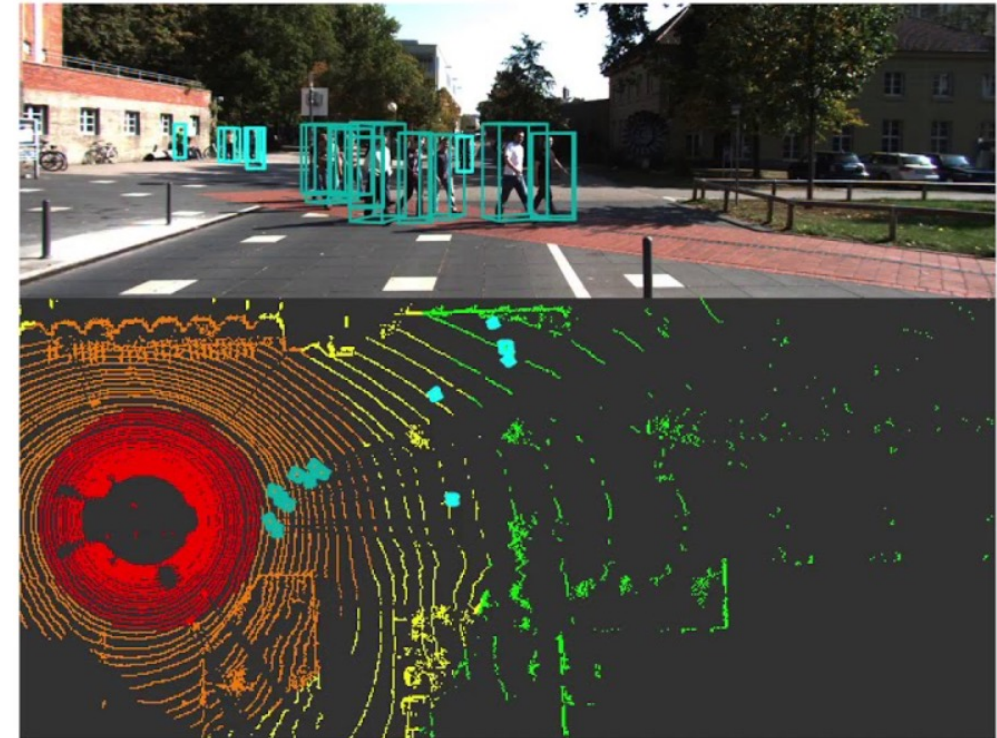
## ❖ **Problem Addressed: Object Detection**

➤ Object detection is the problem of both locating AND classifying objects

➤ Goal of object detection algorithm is to do object detection both fast AND
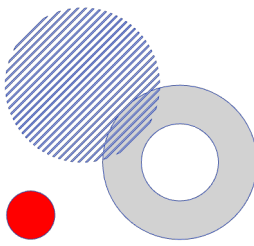   with high accuracy.



Image Source  Object Detection

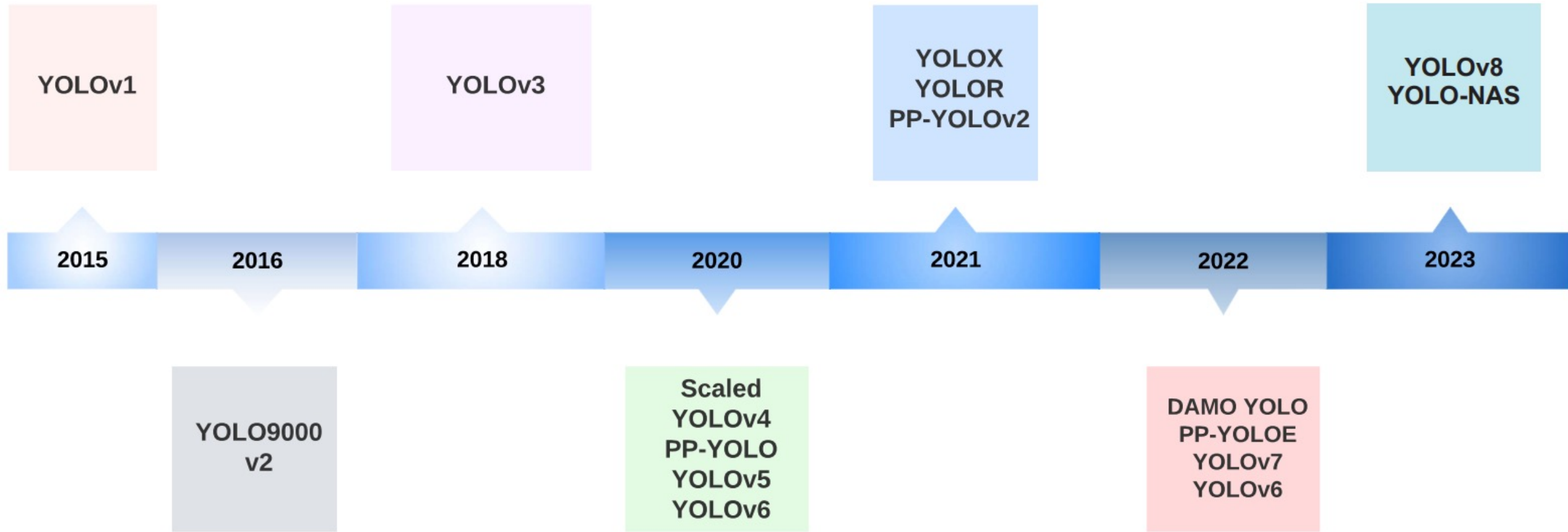Prepared by: Dr. THITTAPORN GANOKRATANAA

## ❖ **Importance of Object Detection**
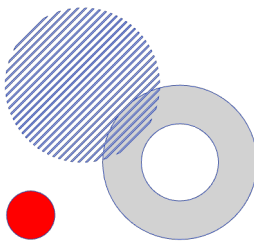
➢ Visual modality is very powerful

➢ Humans are able to detect objects and do perception using just this modality in real-time (not needing radar)

➢ If we want responsive robot systems that work real-time (without specialized sensors), almost real-time vision based object detection can help greatly.

❖ **A timeline of YOLO versions**

Prepared by: Dr. THITTAPORN GANOKRATANAA

## ❖ YOLO Overview

➢ First, image is split into a S×S grid

➢ For each grid square, generate B bounding boxes

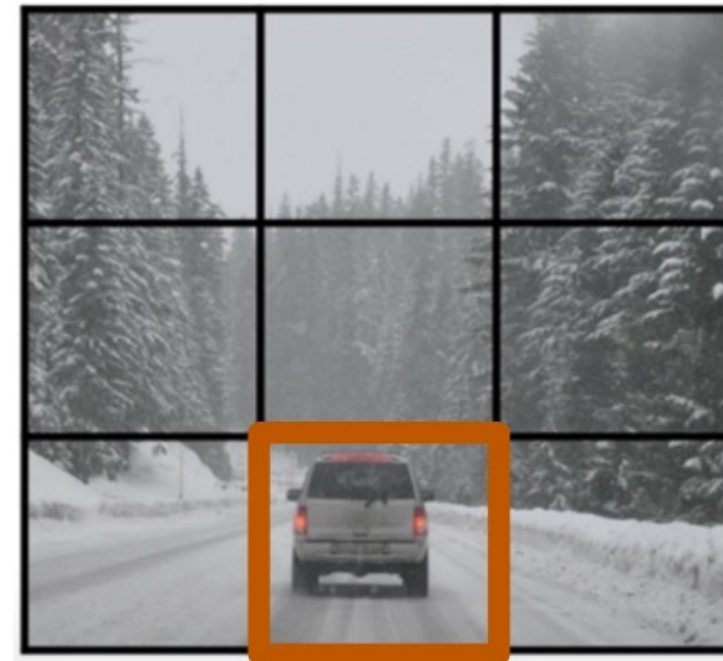➢ For each bounding box, there are 5 predictions:
  x, y, w, h, confidence



S = 3, B = 2

Prepared by: Dr. THITTAPORN GANOKRATANAA

❖ **YOLO Training**

➢ YOLO is a regression algorithm. What is X? What is Y?

➢ X is simple, just an image width (in pixels) * height (in pixels) * RGB values

➢ Y is a tensor of size S * S * (B * 5 + C)

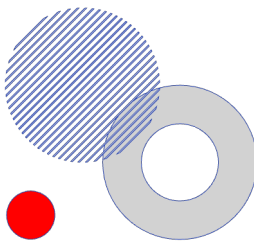➢ B*5 + C term represents the predictions + class predicted distribution for a grid block

For each grid block, we have a vector like this. For this example B is 2 and C is 2

| |
|---|
| $p\_1$ |
| $b\_x\_1$ |
| $b\_y\_1$ |
| $b\_h\_1$ |
| $b\_w\_1$ |
| $p\_2$ |
| $b\_x\_2$ |
| $b\_y\_2$ |
| $b\_h\_2$ |
| $b\_w\_2$ |
| $c\_1$ |
| $c\_2$ |

GT label example:

| |
|---|
| 1 |
| $b\_x\_1$ |
| $b\_y\_1$ |
| $b\_h\_1$ |
| $b\_w\_1$ |
| 0 |
| ? |
| ? |
| ? |
| ? |
| $c\_1 = 1$ |
| $c\_2 = 0$ |

Prepared by: Dr. THITTAPORN GANOKRATANAA

## ❖ YOLO Architecture

➢ Now that we know the input and output, we can discuss the model

➢ We are given 448 by 448 by 3 as our input.

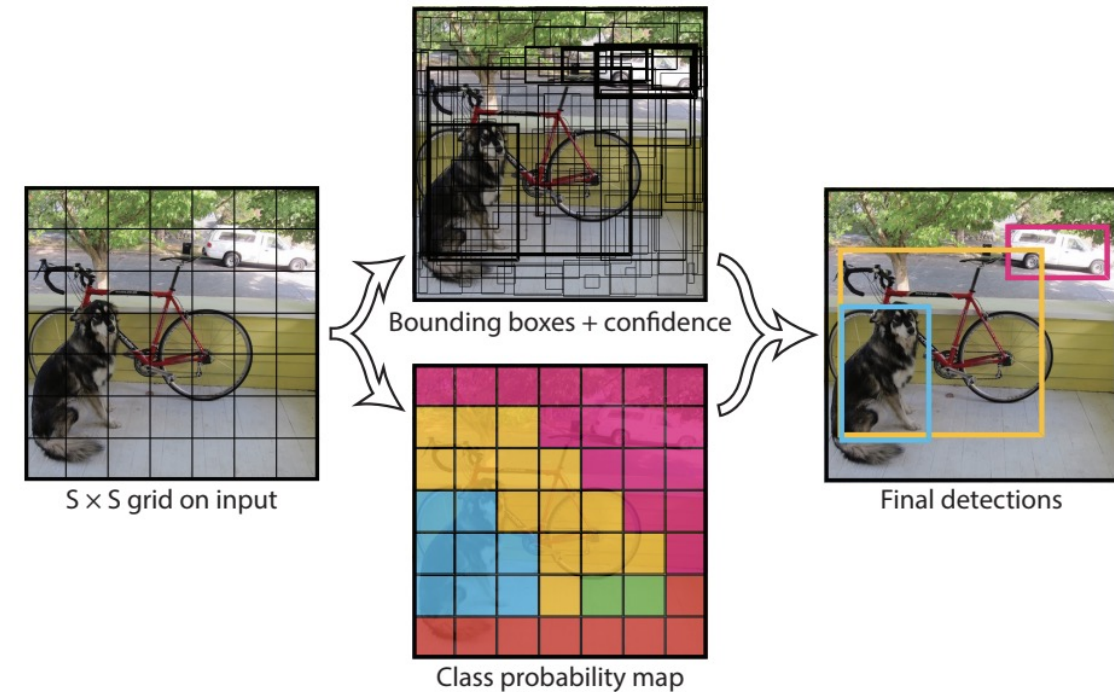➢ Implementation uses 7 convolution layers

➢ Paper parameters: S = 7, B = 2, C = 20

➢ Output is S*S*(5B+C) = 7*7*(5*2+20) = 7*7*30

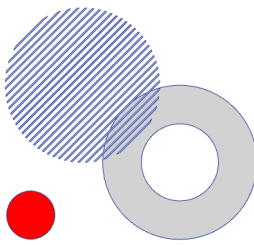Prepared by: Dr. THITTAPORN GANOKRATANAA

❖ **Non-maximal suppression**

➢ We then use the output to make final detections

➢ Use a threshold to filter out bounding boxes with low P(Object)

➢ In order to know the class for the bounding box compute score take argmax over the distribution Pr(Class|Object) for the grid the bounding box's center is in



Bounding boxes + confidence
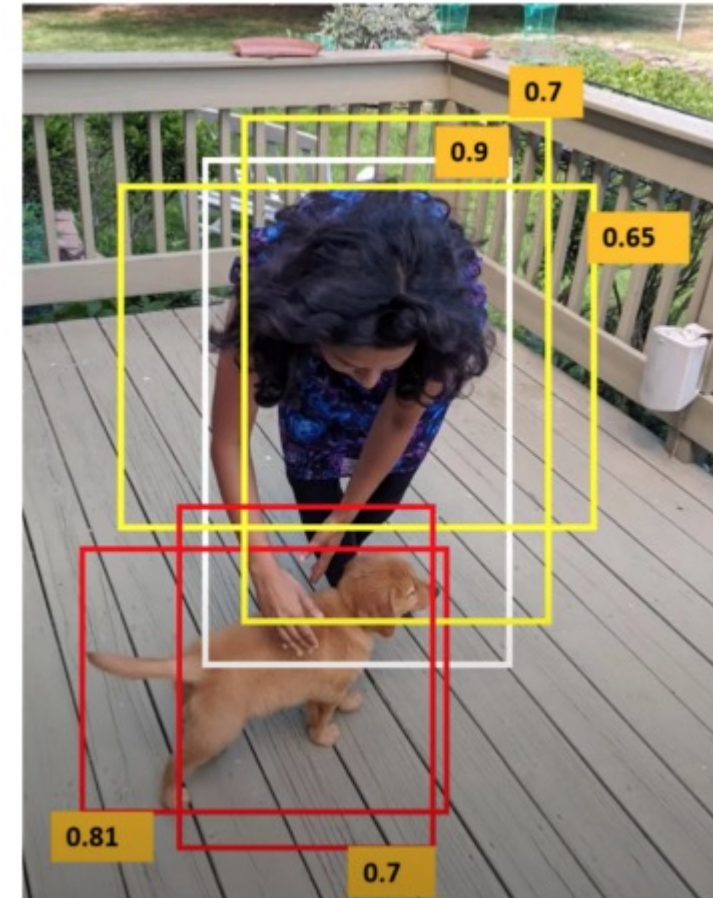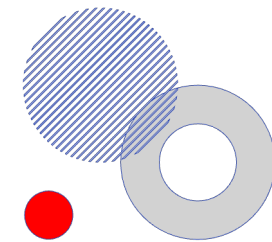
S × S grid on input

Class probability map

Final detections

$$\Pr(\text{Class}_i|\text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

Prepared by: Dr. THITTAPORN GANOKRATANAA

## ❖ YOLO Prediction

➤ Most of the time objects fall in one grid, however it is still possible to get redundant boxes (rare case as object must be close to multiple grid cells for this to happen)

➤ Discard bounding box with high overlap (keeping the bounding box with highest confidence)

➤ Adds 2-3% on final mAP score

Localization loss

Set to 5 to increase the loss of bounding box predictions

GT bbox x-coordinate in the ith cell

GT bbox y-coordinate in the jth cell

Predicted bbox x-coordinate in the ith cell

Predicted bbox y-coordinate in the ith cell

Sum-squared error

$$\lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{obj} \left[ (x_i - \widehat{x_i})^2 + (y_i - \widehat{y_i})^2 \right]$$

For each grid cell

For each grid box

'1' if object appears in the ith cell and the jth box detect it, '0' otherwise

GT bbox height in the ith cell

$$+\lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{obj} \left[ \left(\sqrt{w_i} - \sqrt{\widehat{w_i}}\right)^2 + \left(\sqrt{h_i} - \sqrt{\widehat{h_i}}\right)^2 \right]$$
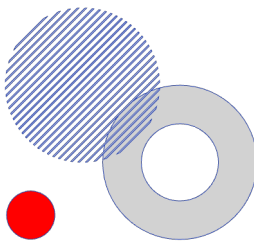
Square root to reduce the range of the values

GT bbox width in the ith cell

Predicted bbox width in the ith cell

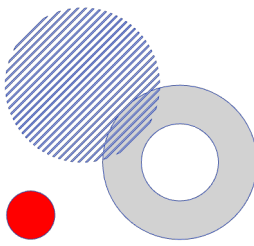Predicted bbox height in the ith cell

❖ **YOLO Objective Function (Cont.)**



**Confidence loss**

$$+\sum_{i=0}^{S^2}\sum_{j=0}^{B}\mathbb{1}_{ij}^{obj}\left[\left(C_i-\widehat{C}_i\right)^2\right]$$

GT confidence score

Predicted confidence score

Confidence error when an object is detected in the ith cell

Set to 0.5 to decrease the loss for empty boxes

$$+\lambda_{noobj}\sum_{i=0}^{S^2}\sum_{j=0}^{B}\mathbb{1}_{ij}^{noobj}\left[\left(C_i-\widehat{C}_i\right)^2\right]$$

'1' if there is no object in the ith cell, '0' otherwise

Confidence error when an object not detected in the ith cell

**Classification loss**

Predicted conditional probability of an object of class c appearing in the ith cell

$$+\sum_{i=0}^{S^2}\mathbb{1}_i^{obj}\sum_{c\in classes}\left[(p_i(c)-\widehat{p}_i(c))^2\right]$$

For each grid cell

For each class

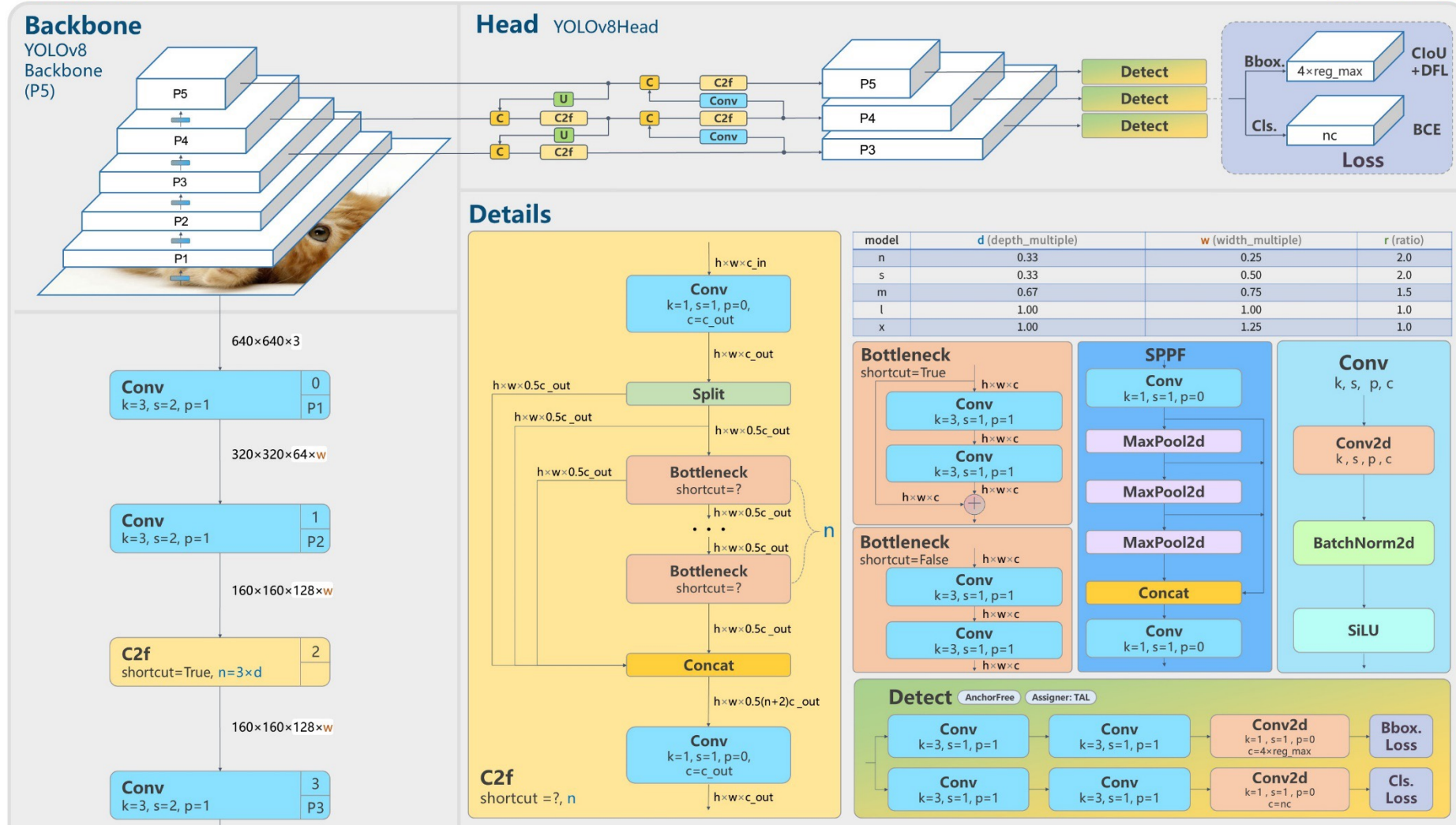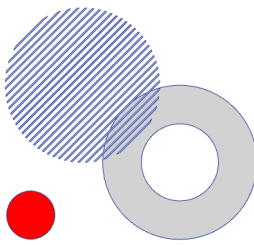GT conditional probability of class c appearing in the ith cell

## ❖ **YOLO V8**

➢ YOLOv8 uses a similar backbone as YOLOv5 with some changes on the CSPLayer, now called the C2f module.

➢ The C2f module (cross-stage partial bottleneck with two convolutions) combines high-level features with contextual information to improve detection accuracy
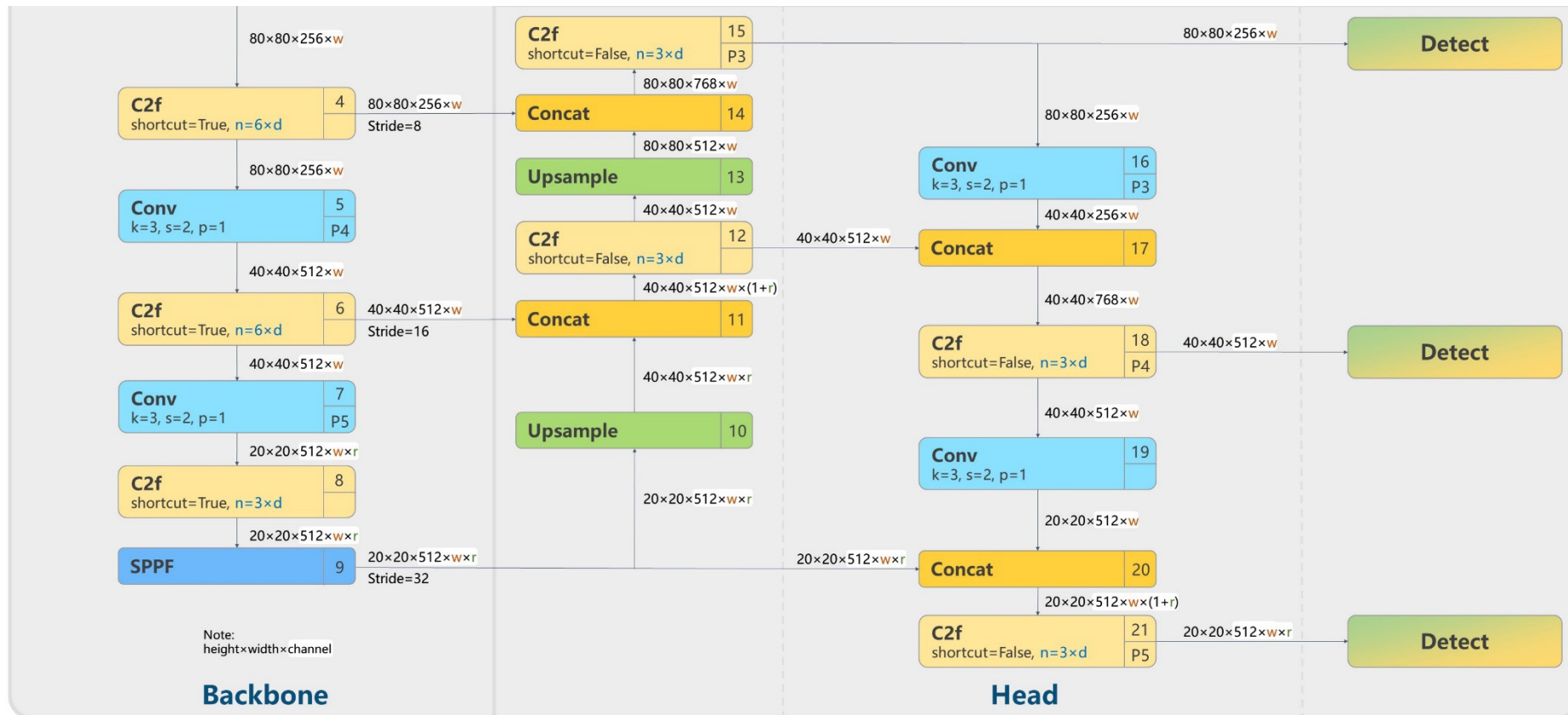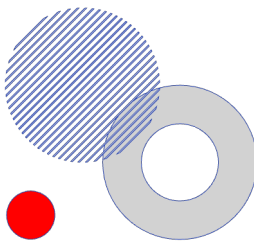
# ❖ YOLO V8 Architecture

Prepared by: Dr. THITTAPORN GANOKRATANAA

Prepared by: Dr. THITTAPORN GANOKRATANAA

❖ **YOLO V8 Experiment**

➢ Using this Google Colab:

https://colab.research.google.com/drive/14x7_B44tBvAe8RzuETDVJ14cYWstnT2D?usp=sharing

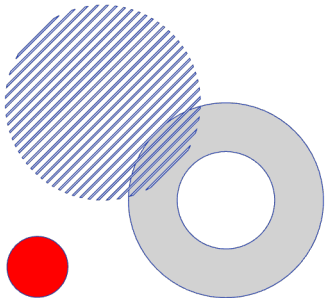Prepared by: Dr. THITTAPORN GANOKRATANAA

# Exercise

Extract this video into frame and label it into four classes (bus, taxi, car, and pedestrian), then generate the model to classify those four classes using yolov8

# Conclusion

- The research focused on utilizing AI technology to augment police efficiency in Thailand.
- We aimed to enhance law enforcement capabilities and bolster public trust in crime prevention measures.
- By employing AI in crime data analysis, leveraging intelligent CCTV technology for crime monitoring, and integrating real-time alerts for suspicious activities to police.

# Q&A