

# Lecture 08 – Image Detection and Segmentation

Prof. André Gustavo Hochuli

[gustavo.hochuli@pucpr.br](mailto:gustavo.hochuli@pucpr.br)

[aghochuli@ppgia.pucpr.br](mailto:aghochuli@ppgia.pucpr.br)

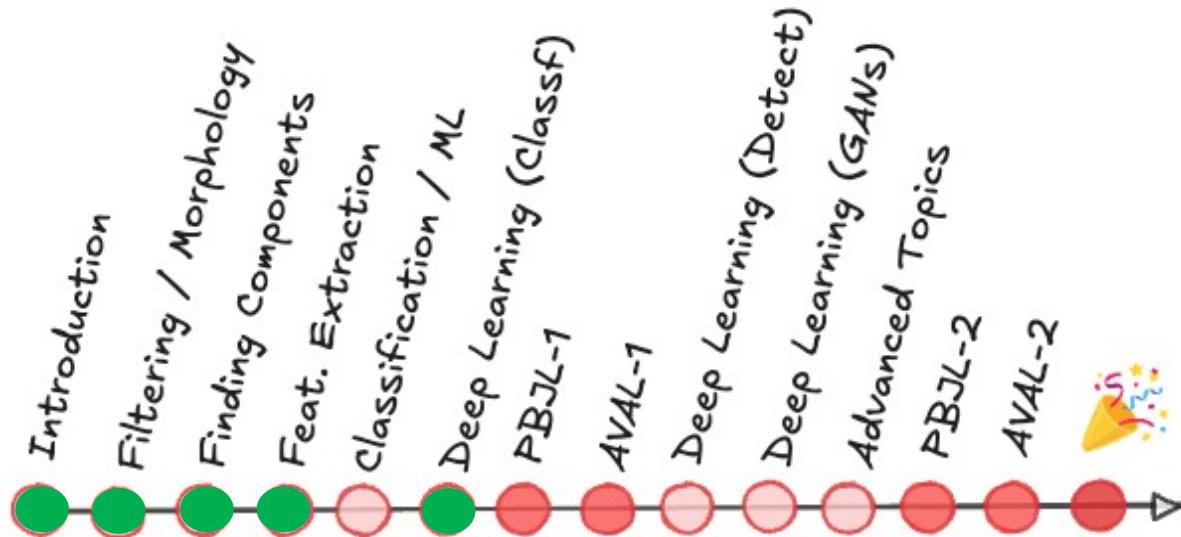
# Topics

- Review of Lecture 10 – CNN Applications and Tricks

- Classification vs Segmentation

- Classification
- Object Detection
- Segmentation

- Practice

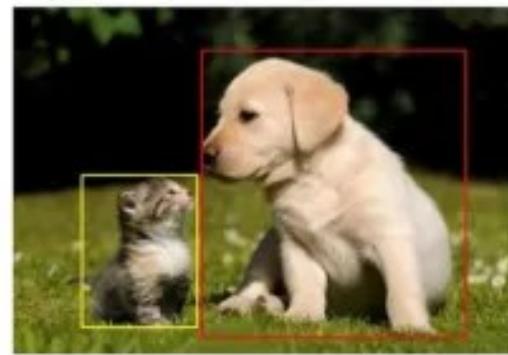


# Classification vs Segmentation

Is this a dog?



What is there in image  
and where?



Which pixels belong to  
which object?

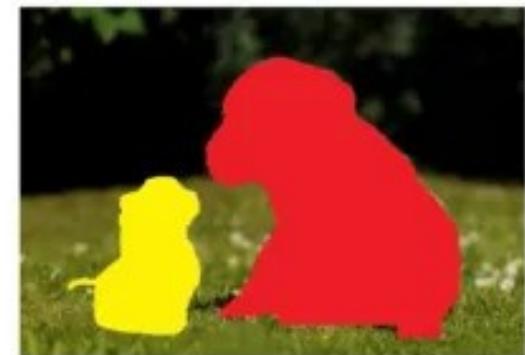
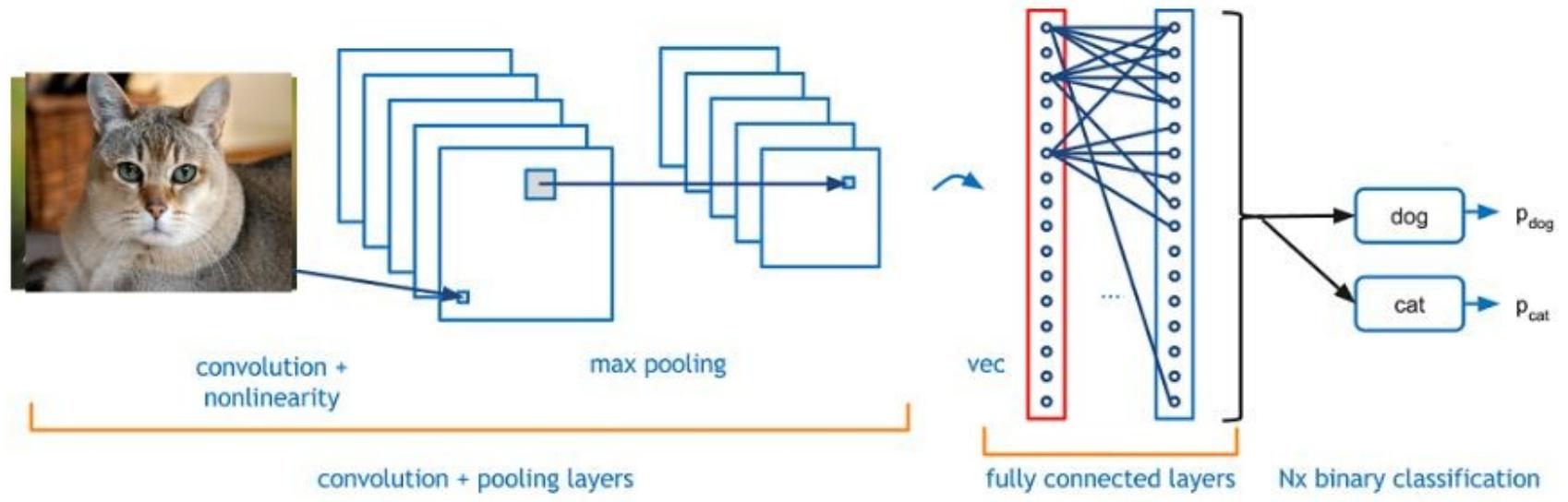
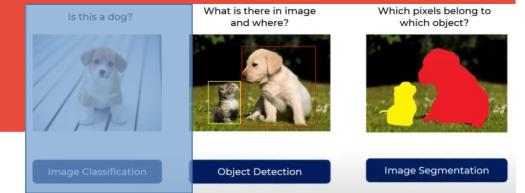


Image Classification

Object Detection

Image Segmentation

# Classification



# Object Detection

Is this a dog?



Image Classification

What is there in image  
and where?



Object Detection

Which pixels belong to  
which object?

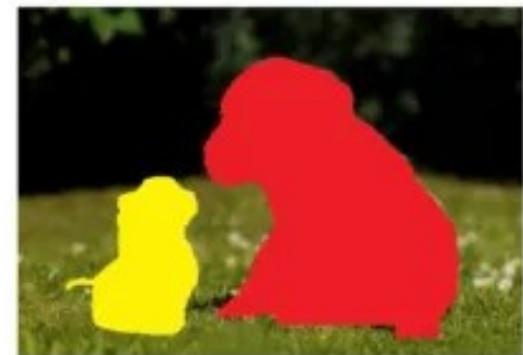
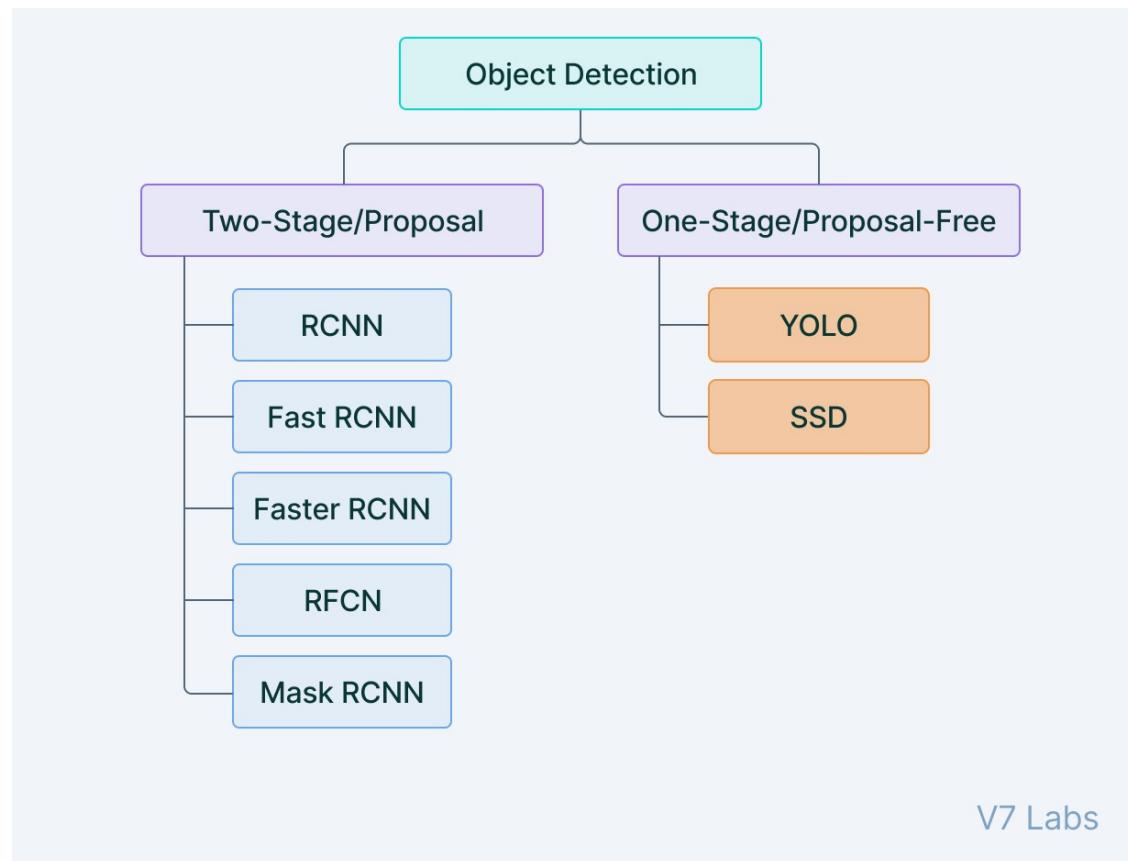


Image Segmentation

# Object Detection

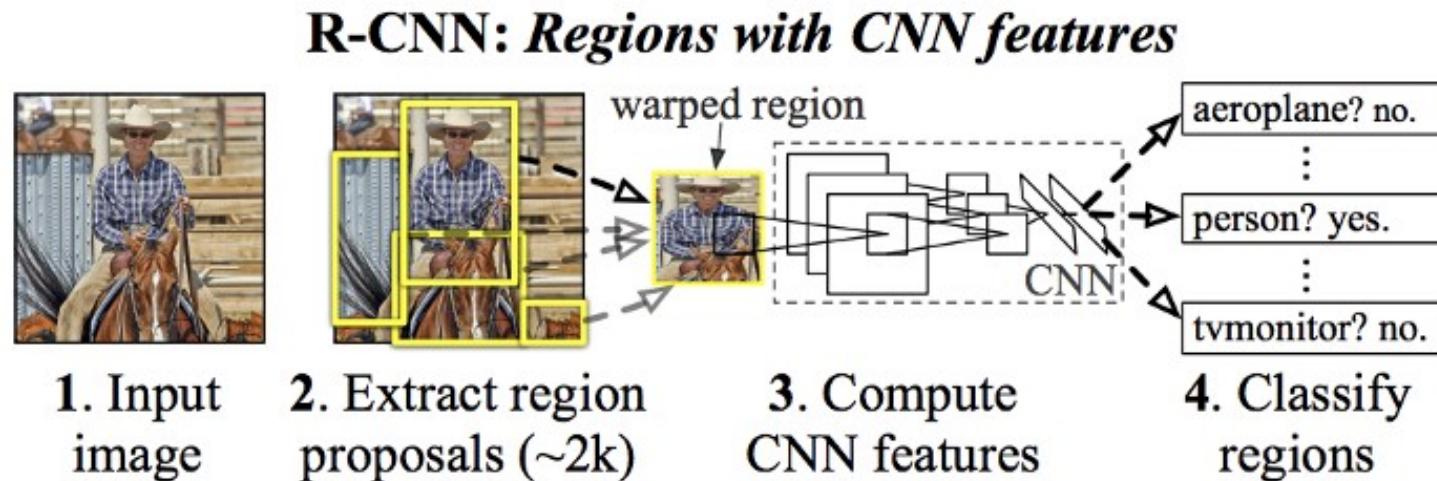


V7 Labs

# Object Detection - RCNN



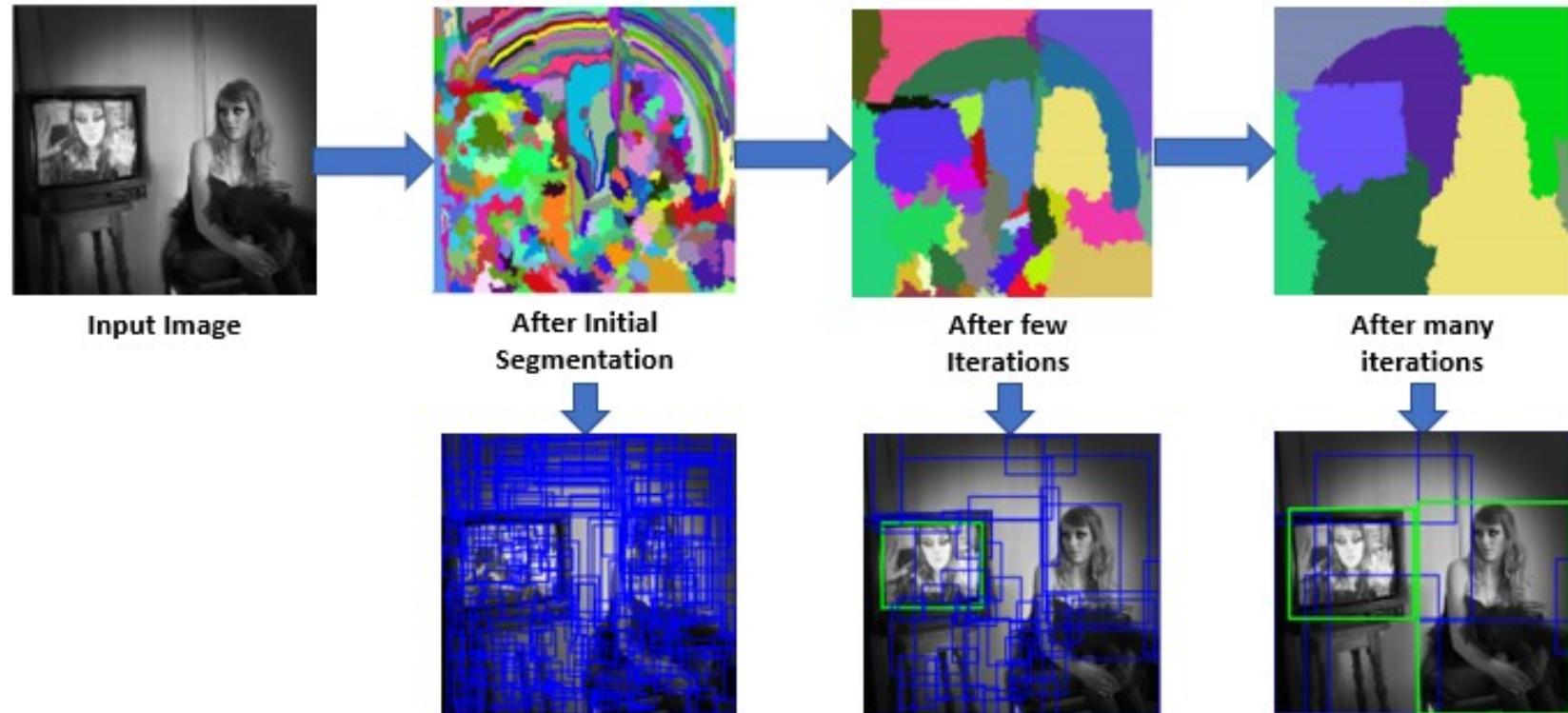
- Region Based Convolutional Neural Network (2014) - Ross Girshick
- Selective Search Algorithm (Region Proposal)
- CNN (Classification)



# Object Detection - RCNN

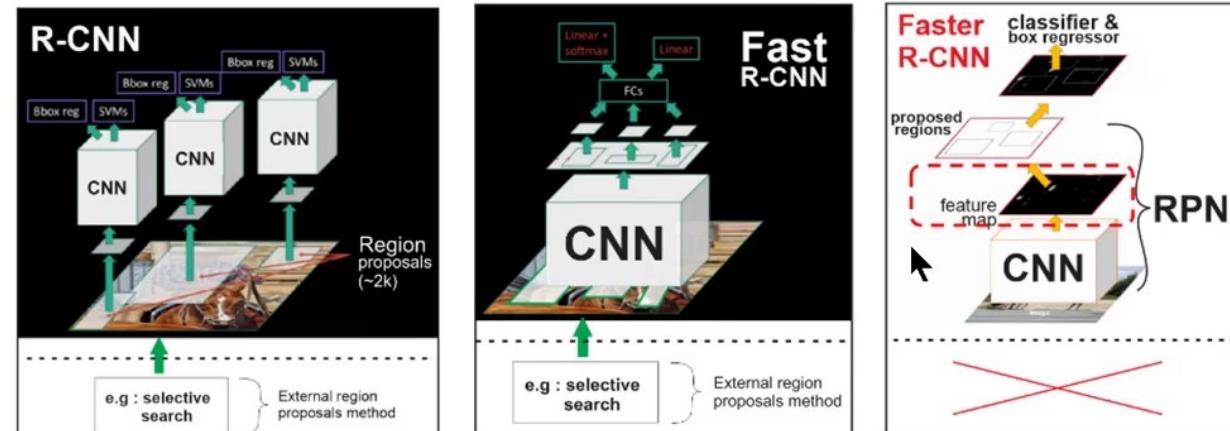


- Selective Search Algorithm (Region Proposal)



# Object Detection - RCNN

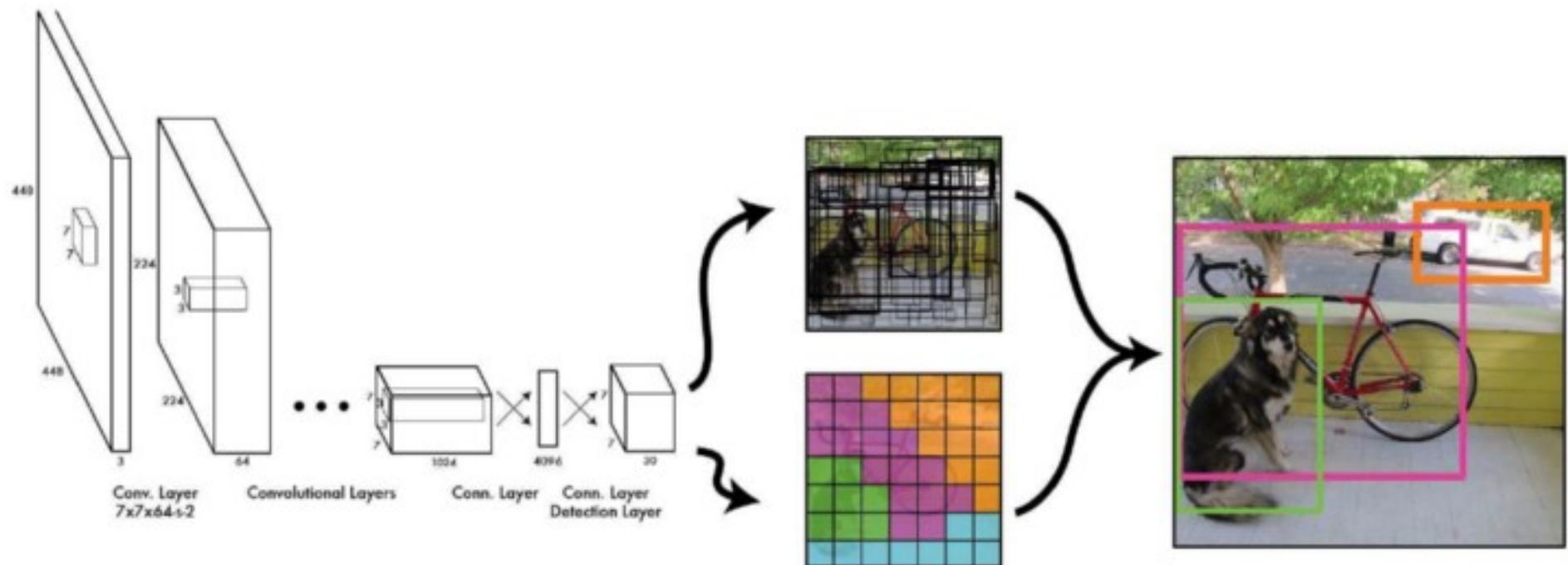
- R-CNN: Selective Search->CNN
- Fast: End-to-end (Sel. Search->ROI Pooling→FC)
- Faster: Region Proposal Network (RPN)



	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image	50 seconds	2 seconds	0.2 seconds
Speed-up	1x	25x	250x
mAP (VOC 2007)	66.0%	66.9%	66.9%

# Object Detection - Yolo

- You Look Once (YoLo - 2015 - now)
  - Joseph Redmon / Ross Girshick
- Fast End-to-End Architecture



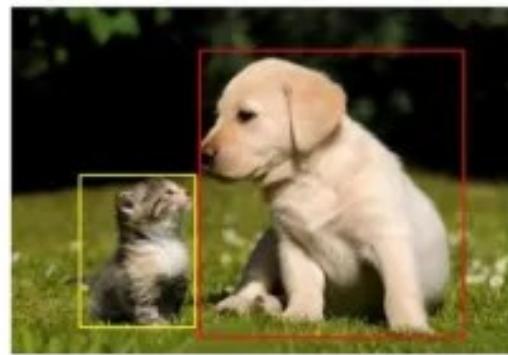
# Segmentation

Is this a dog?



Image Classification

What is there in image  
and where?



Object Detection

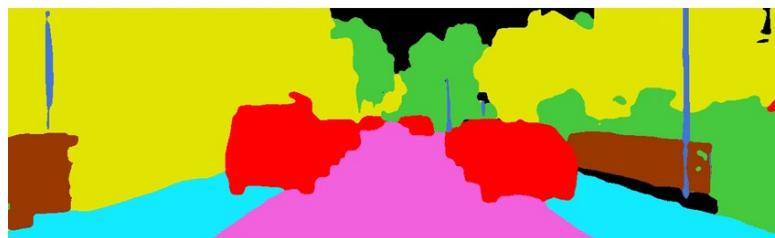
Which pixels belong to  
which object?



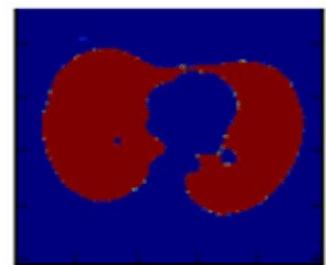
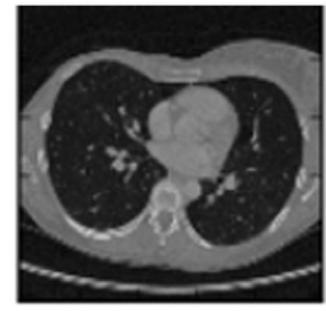
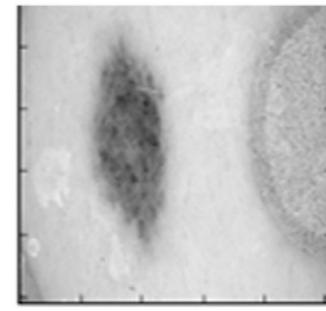
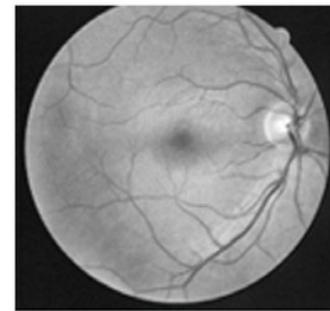
Image Segmentation

# Segmentation

- Classification at pixel level

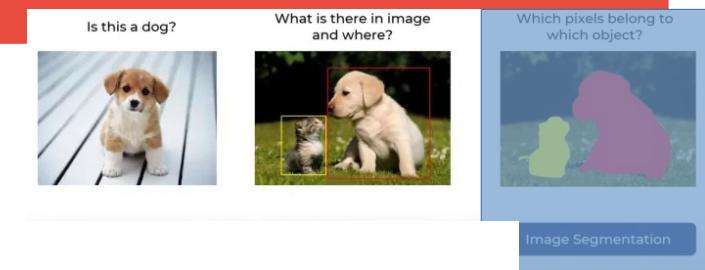


Road	Sidewalk	Building	Fence
Pole	Vegetation	Vehicle	Unlabel

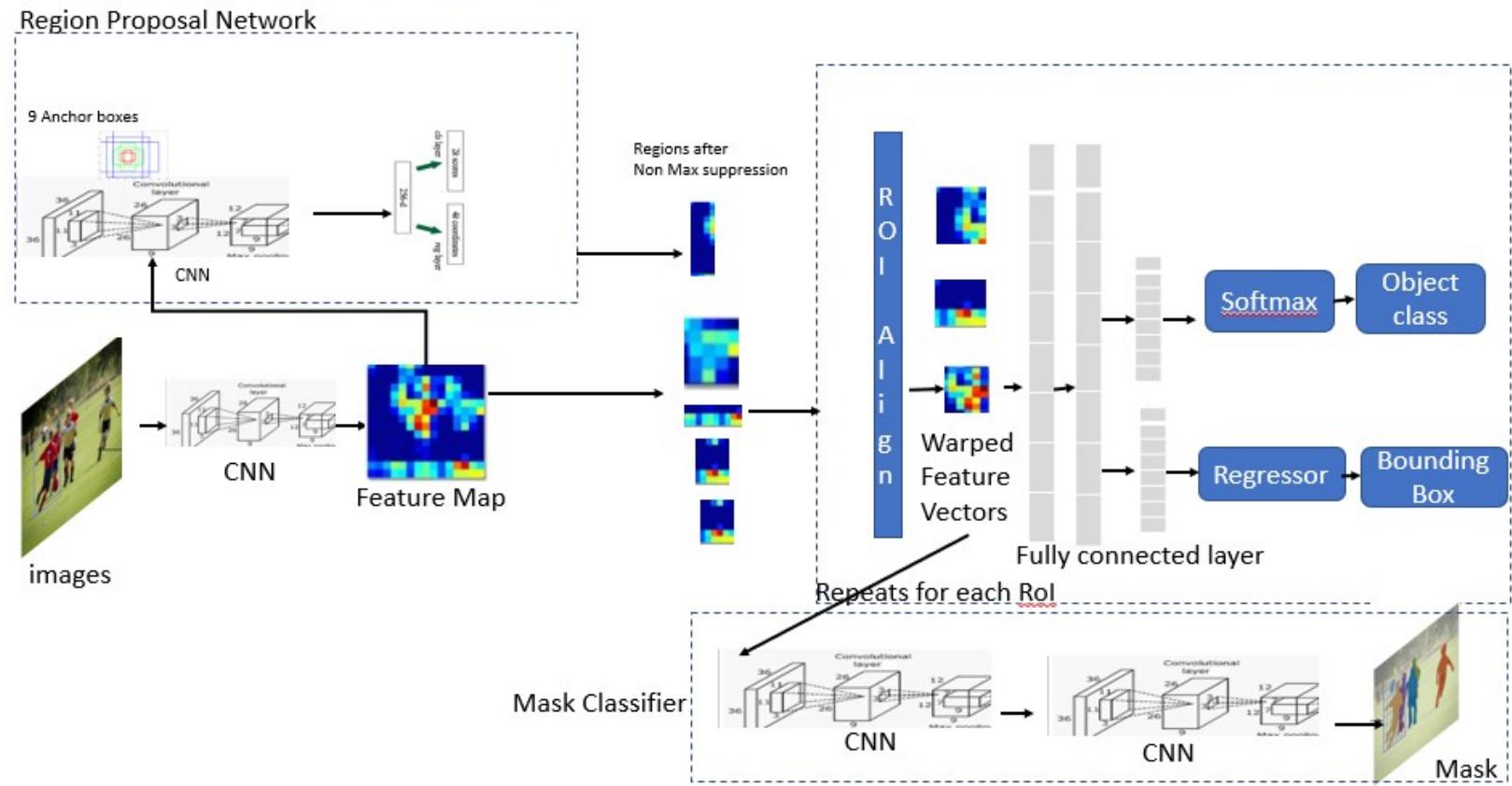


# Segmentation – Mask RCNN

- Faster R-CNN with Binary Mask (2017)

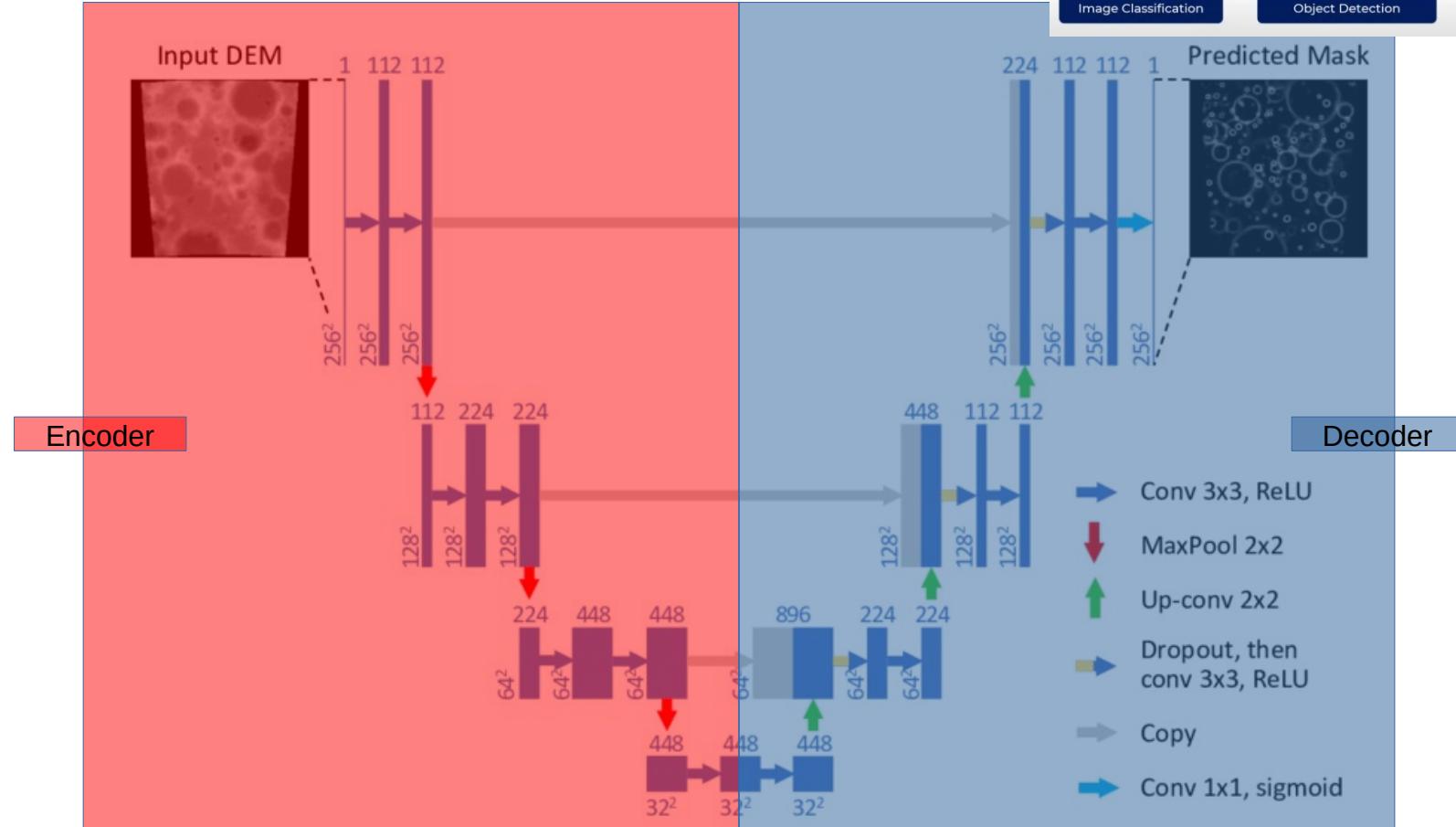


## Mask RCNN



# Segmentation - UNET

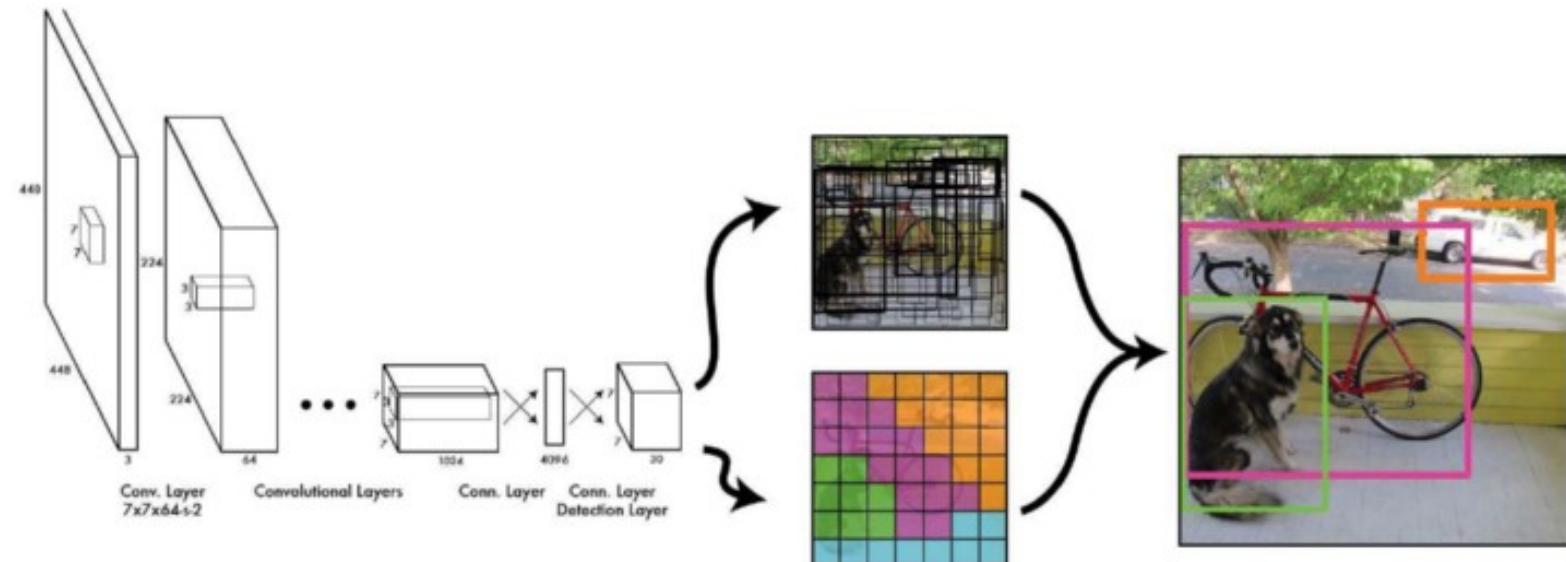
- U-Net (Encoder and Decoder)



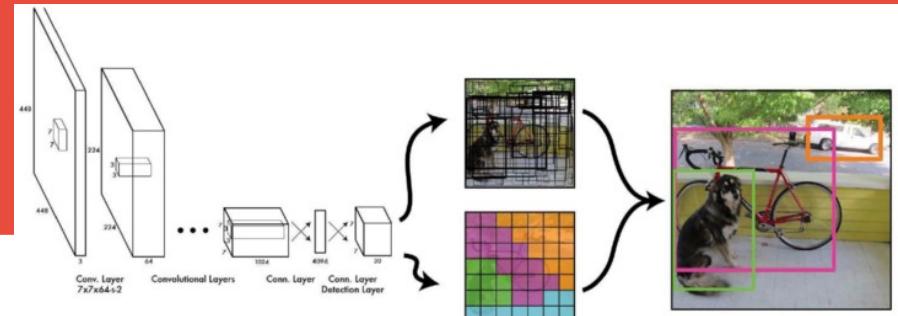
**YoLo – You Look Once**

# Object Detection - Yolo

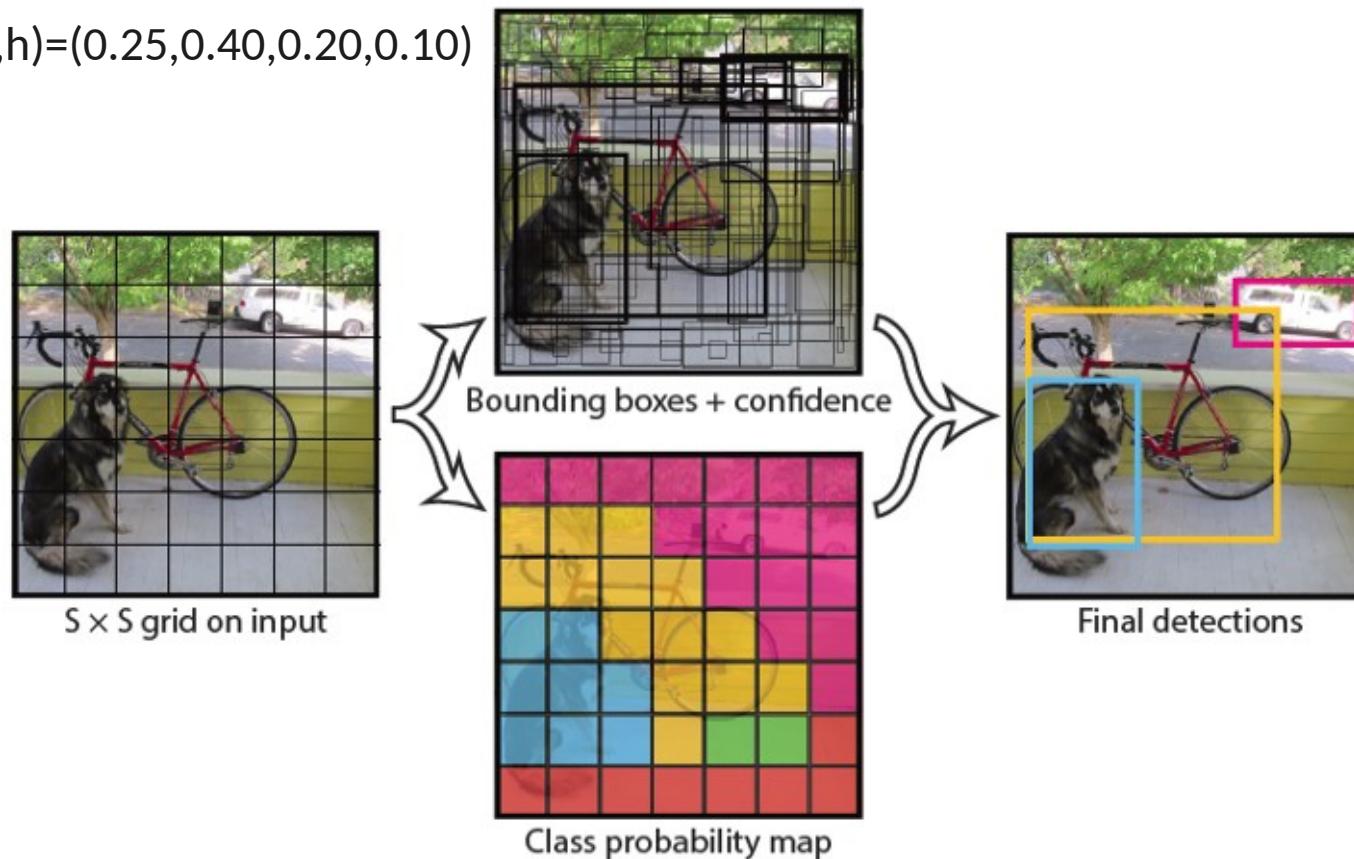
- You Look Once (YoLo - 2015 - now)
  - The input image passes through convolutional backbone (e.g., Darknet)
  - The output is a feature map of lower spatial resolution (e.g for instance, 80x80, 20x20 )
  - The split is applied on the latent (feature map), not the raw input. Each feature cell corresponds to a specific spatial region (receptive field) of the input image.



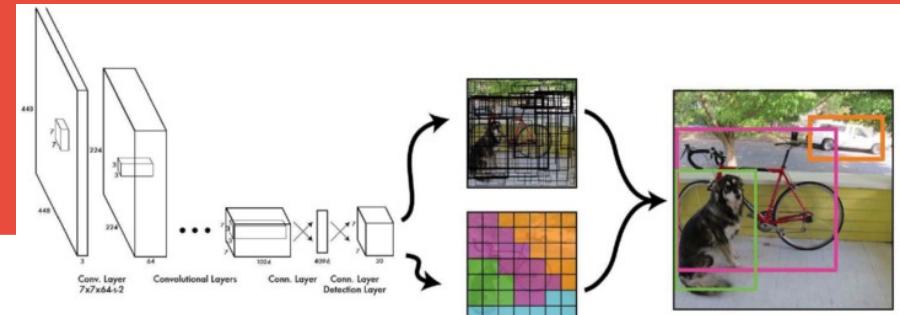
# Object Detection - Yolo



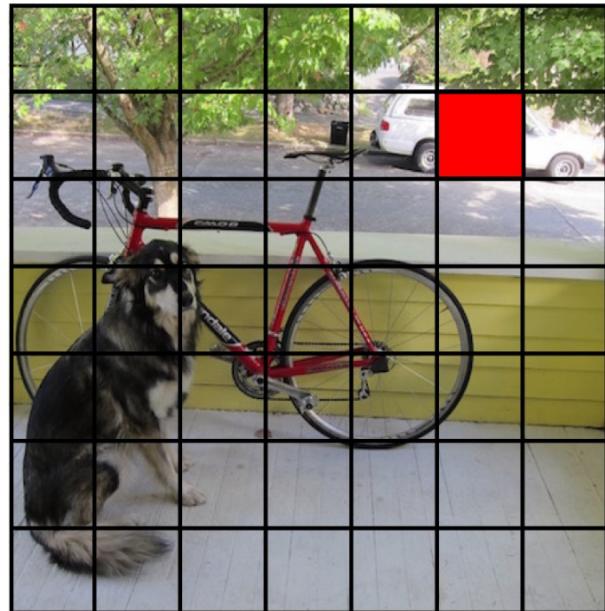
- Regression is the key!
  - Bounding boxes are treated as continuous variables in normalized image coordinates.
  - $(x, y, w, h) = (0.25, 0.40, 0.20, 0.10)$



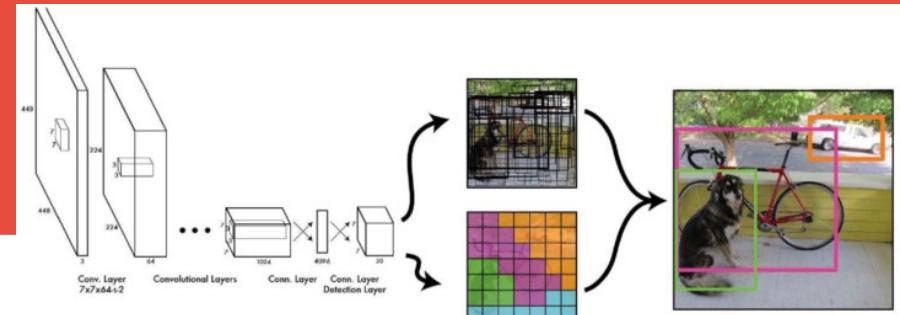
# Object Detection - Yolo



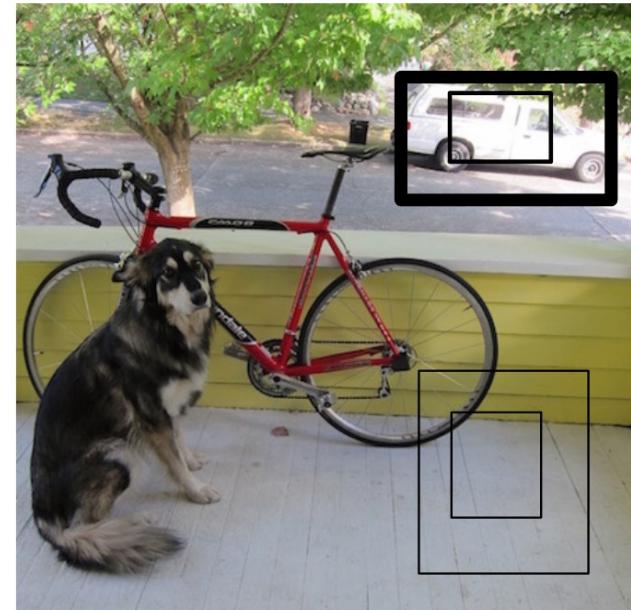
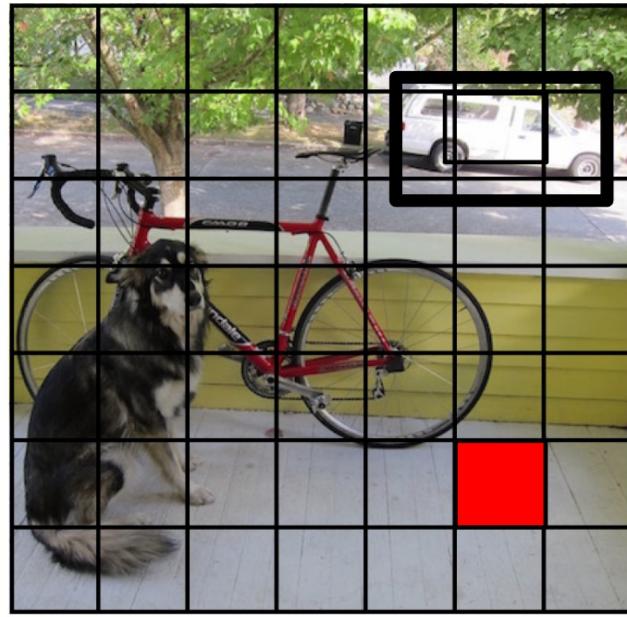
- Each cell predicts bounding boxes and confidences
  - $P(\text{object})$ : [0,1] quantifies the confidence that any object occupies this box (not background)



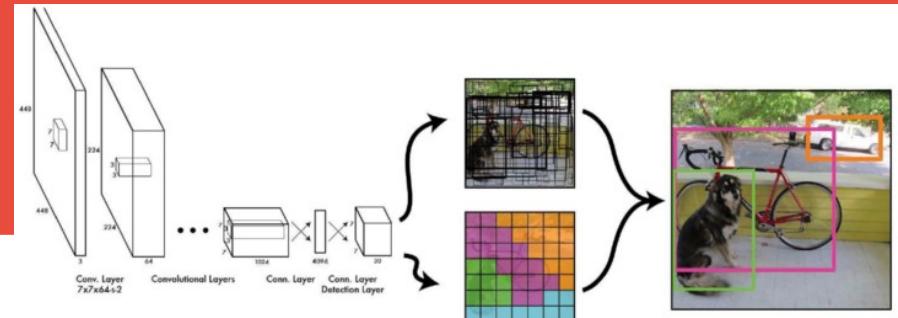
# Object Detection - Yolo



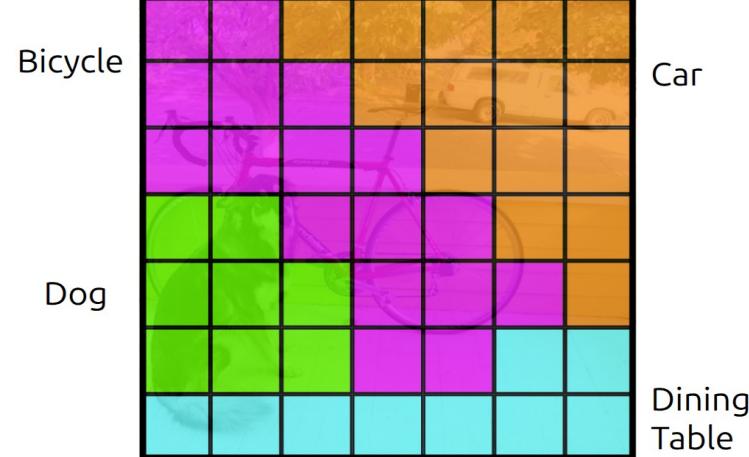
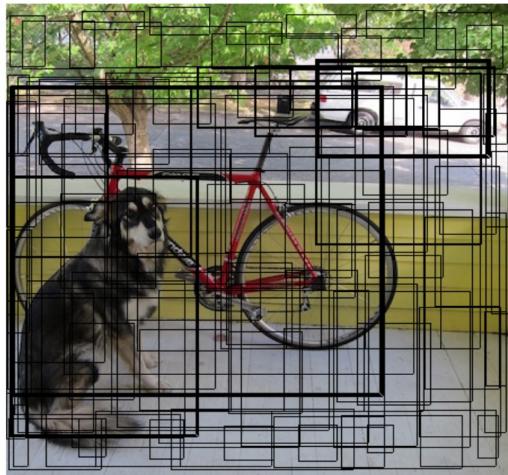
- Each cell predicts bounding boxes and confidences
  - $P(\text{object})$ : [0,1] quantifies the confidence that any object occupies this box (not background)



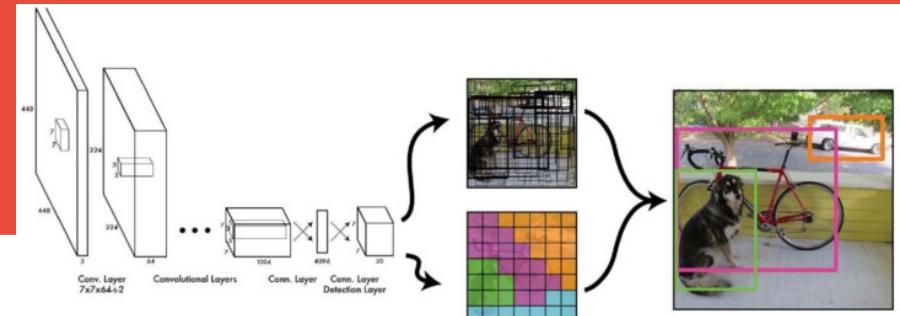
# Object Detection - Yolo



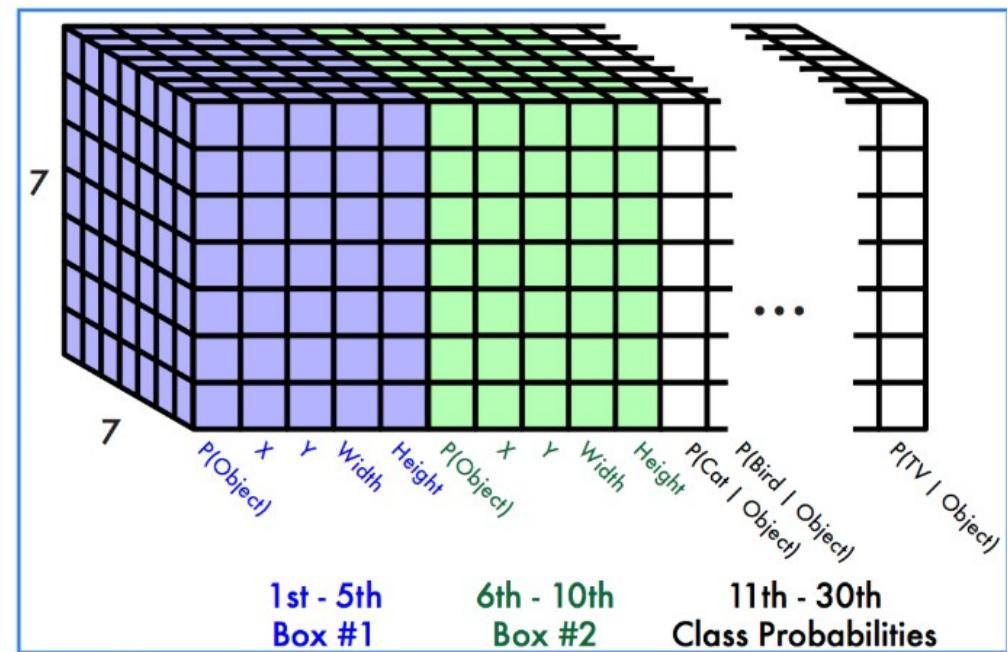
- Class Prob =  $P(\text{class}) \Rightarrow P(\text{car}) = 0.8$
- Conditionated Prob: e.g  $P(\text{class} \mid \text{object}) \Rightarrow P(\text{car}) = 0.9$
- Confidence:
  - $P(\text{Object}) * P(\text{Car} \mid \text{Object}) = 0.72$



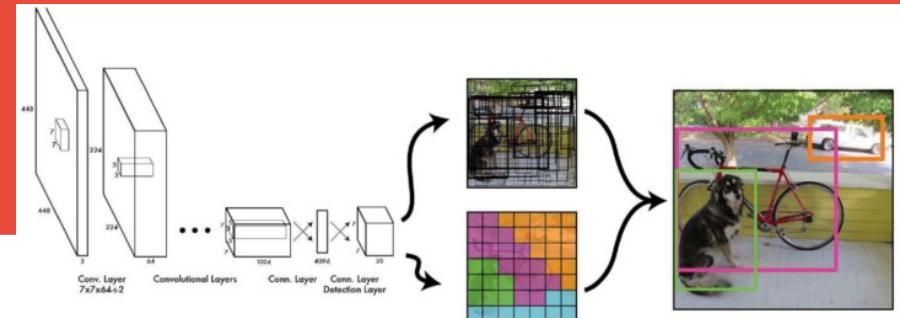
# Object Detection - Yolo



- Each cell predicts:
  - Bounding Boxes:
    - 4 Coordinates (X,Y,W,H)
    - 1 Confidence
  - I.E PASCAL VOC
    - 7x7 Grid
    - 2 Bounding Box / Cell
    - 20 Classes
    - $7 * 7 * (2 * 5 + 20) = 7x7^*30$  tensors per cell => 1470 predictions per image



# Object Detection - Yolo

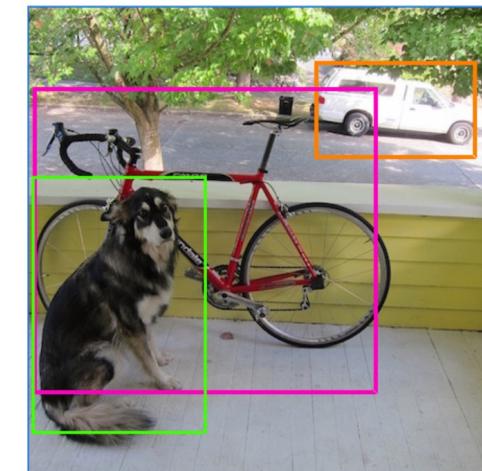


- Non-Maximum Supression (NMS)
  - Sort all boxes by confidence score ( $P(\text{object}) \times P(\text{class})$ ).
  - Pick the box with the highest score → keep it as the best detection.
    - Compute IoU between this box and all others
    - Remove all boxes with IoU above the suppression threshold (e.g., 0.5)
    - Repeat until all boxes are processed

$$\text{IoU} = \frac{\text{INTERSECTION}}{\text{UNION}}$$

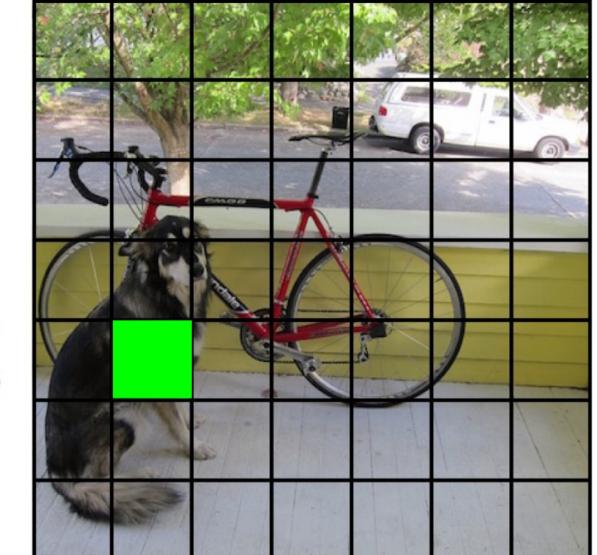
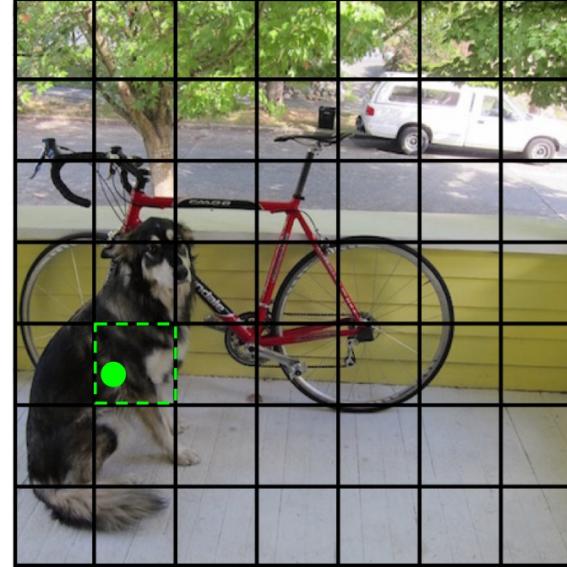
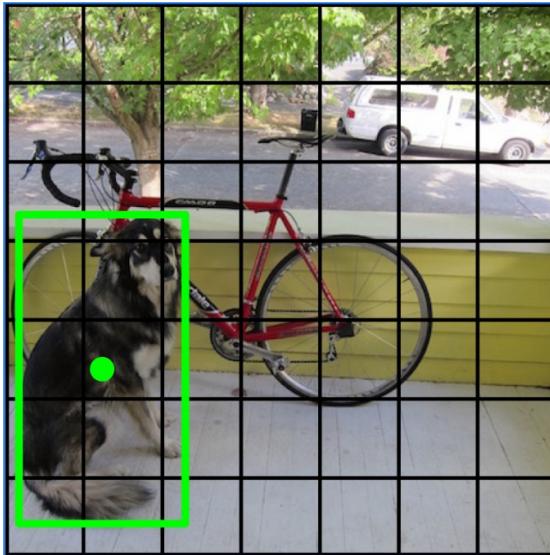
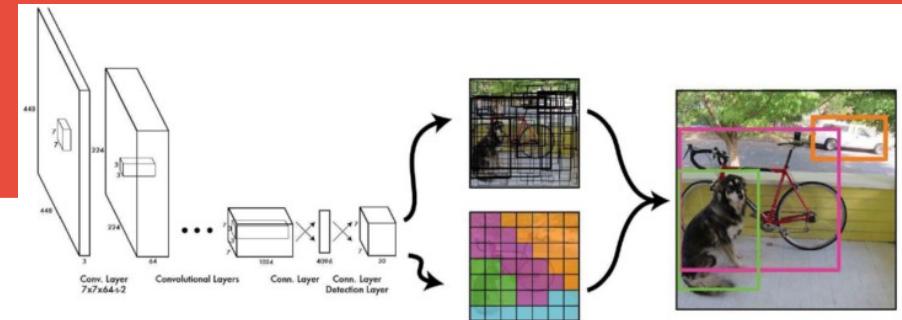
Diagram illustrating Intersection over Union (IoU) calculations for overlapping rectangles:

- Two non-overlapping rectangles:  $\text{IoU} = 0.0$
- One rectangle inside another:  $\text{IoU} = 0.08$
- Two overlapping rectangles:  $\text{IoU} = 0.18$
- Two rectangles that overlap significantly:  $\text{IoU} = 0.43$
- Two rectangles that completely cover each other:  $\text{IoU} = 1.0$



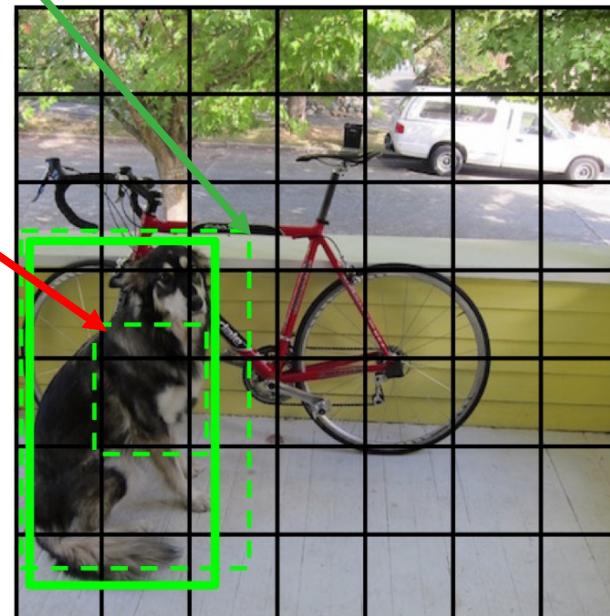
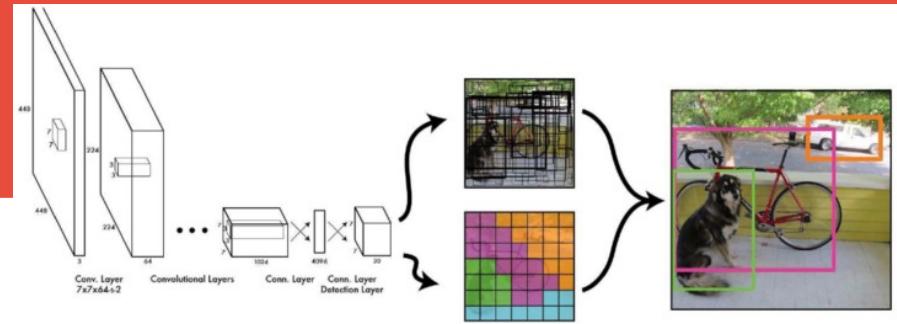
# Object Detection - Yolo

- Training
  - Match example to the right cell (ground-truth)



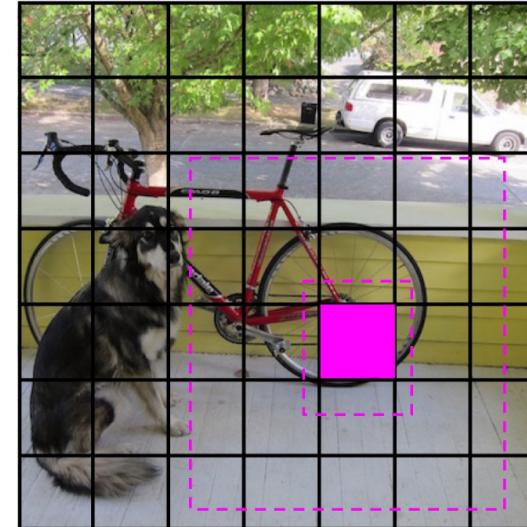
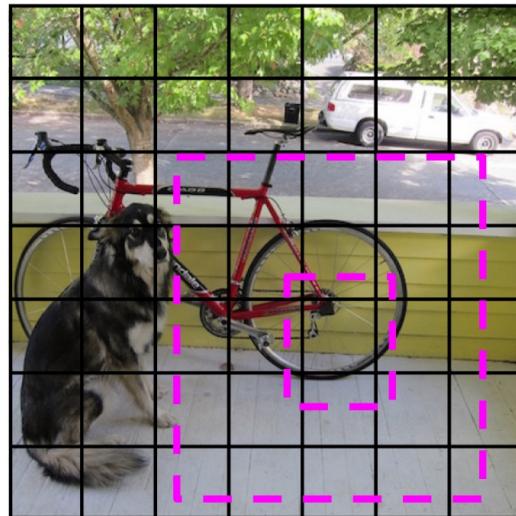
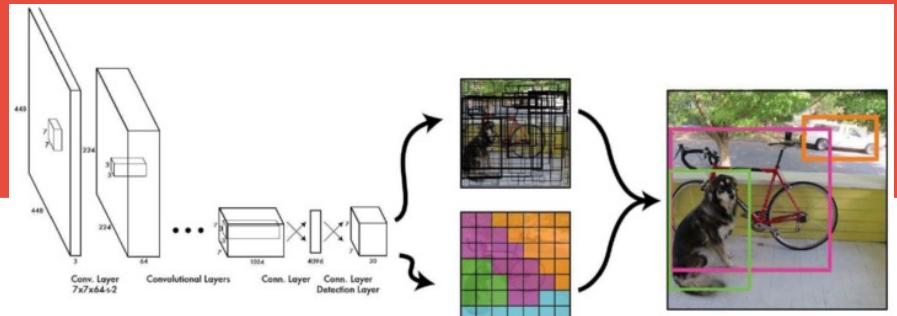
# Object Detection - Yolo

- Training
  - Predict Bounding-Boxes
    - Selects the best fit and increases its confidence
    - Penalizes all other predictions

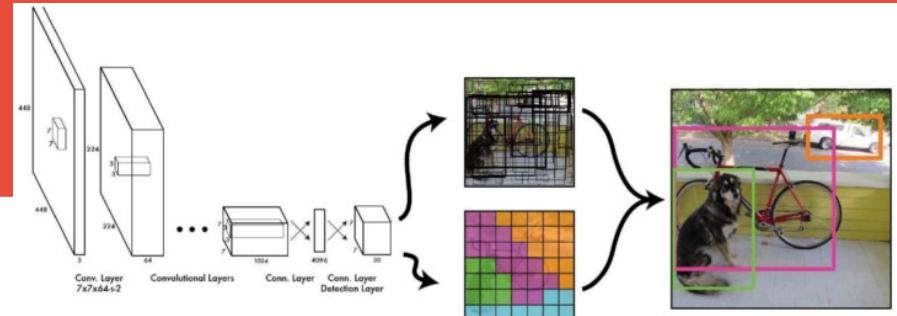


# Object Detection - Yolo

- Training
  - Penalizes when the prediction does not match any class (i.e., background).



# Object Detection - Yolo



loss function:

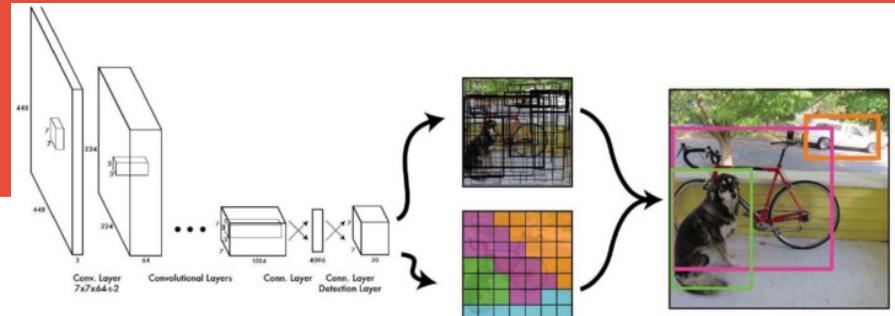
$$\begin{aligned}
 & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
 & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
 & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \\
 & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \\
 & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad (3)
 \end{aligned}$$

model. We use sum-squared error because it is easy to optimize, however it does not perfectly align with our goal of maximizing average precision. It weights localization error equally with classification error which may not be ideal. Also, in every image many grid cells do not contain any object. This pushes the “confidence” scores of those cells towards zero, often overpowering the gradient from cells that do contain objects. This can lead to model instability, causing training to diverge early on.

To remedy this, we increase the loss from bounding box coordinate predictions and decrease the loss from confidence predictions for boxes that don't contain objects. We use two parameters,  $\lambda_{\text{coord}}$  and  $\lambda_{\text{noobj}}$  to accomplish this. We set  $\lambda_{\text{coord}} = 5$  and  $\lambda_{\text{noobj}} = .5$ .

$$\lambda_{\text{coord}} = 5, \lambda_{\text{noobj}} = 0.5$$

# Object Detection - Yolo



loss function:

$$\begin{aligned}
 & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
 & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \\
 & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \\
 & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \\
 & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad (3)
 \end{aligned}$$

$\mathbb{1}_{ij}^{\text{obj}}$

The  $j$ th bbox predictor in **cell  $i$**  is “responsible” for that prediction

$\mathbb{1}_{ij}^{\text{noobj}}$

$\mathbb{1}_i^{\text{obj}}$

If object appears in **cell  $i$**

Note that the loss function only penalizes classification error if an object is present in that grid cell (hence the conditional class probability discussed earlier). It also only penalizes bounding box coordinate error if that predictor is “responsible” for the ground truth box (i.e. has the highest IOU of any predictor in that grid cell).

# Object Detection - Yolo

- Datasets
- PASCAL VOC 2007

&

VOC 2012

20 classes:

- Person: person
- Animal: bird, cat, cow, dog, horse, sheep
- Vehicle: aeroplane, bicycle, boat, bus, car, motorbike, train
- Indoor: bottle, chair, dining table, potted plant, sofa, tv/monitor

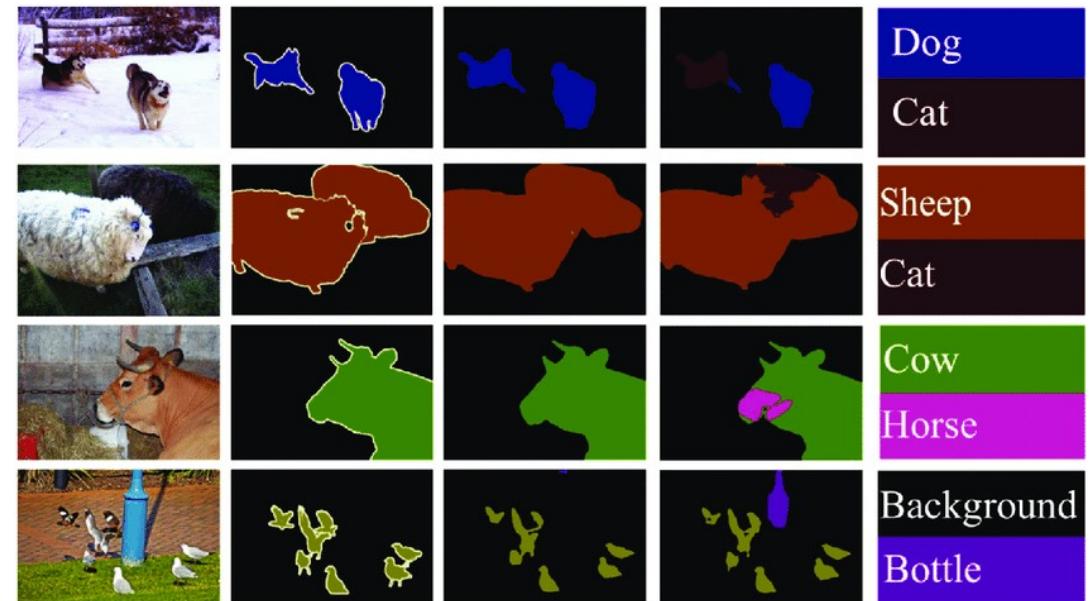
2007

Train/validation/test: 9,963 images containing 24,640 annotated objects.

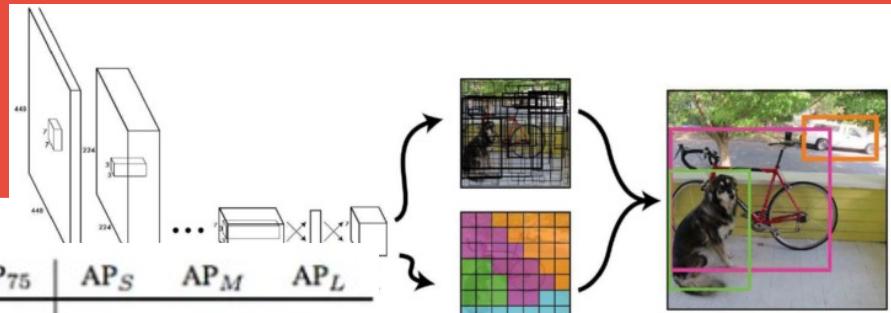


2012

20 classes. The train/val data has 11,530 images containing 27,450 ROI annotated objects and 6,929 segmentations.



# Object Detection - Yolo

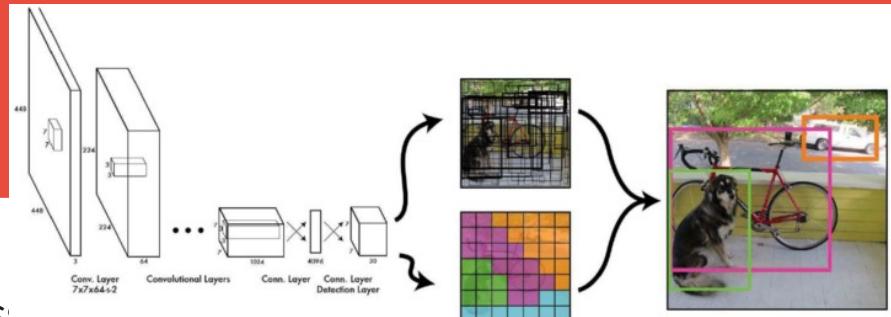


	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<i>Two-stage methods</i>							
Faster R-CNN+++ [3]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [6]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [4]	Inception-ResNet-v2 [19]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [18]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	<b>52.1</b>
<i>One-stage methods</i>							
YOLOv2 [13]	DarkNet-19 [13]	21.6	44.0	19.2	5.0	22.4	35.5
SSD513 [9, 2]	ResNet-101-SSD	31.2	50.4	33.3	10.2	34.5	49.8
DSSD513 [2]	ResNet-101-DSSD	33.2	53.3	35.2	13.0	35.4	51.1
RetinaNet [7]	ResNet-101-FPN	39.1	59.1	42.3	21.8	42.7	50.2
RetinaNet [7]	ResNeXt-101-FPN	<b>40.8</b>	<b>61.1</b>	<b>44.1</b>	<b>24.1</b>	<b>44.2</b>	51.2
YOLOv3 608 × 608	Darknet-53	33.0	57.9	34.4	18.3	35.4	41.9

	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img
Faster R-CNN	73.2	7 FPS	140 ms/img
YOLO	<del>63.4</del> 69.0	45 FPS	22 ms/img

# Object Detection - Yolo

- mAP measures a detector's average precision across:



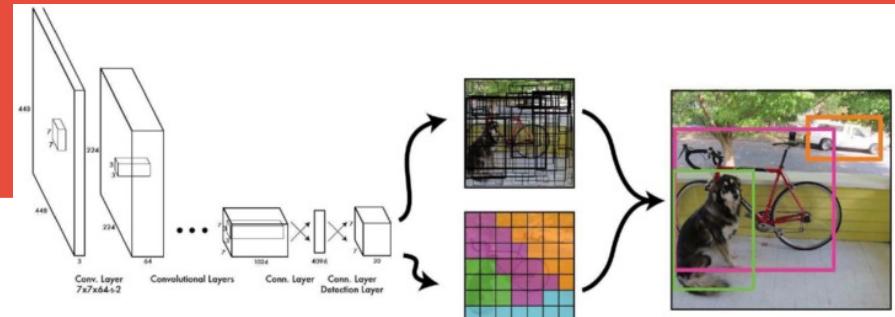
	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<i>Two-stage methods</i>							
Faster R-CNN+++ [3]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [6]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [4]	Inception-ResNet-v2 [19]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [18]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	<b>52.1</b>
<i>One-stage methods</i>							
YOLOv2 [13]	DarkNet-19 [13]	21.6	44.0	19.2	5.0	22.4	35.5
SSD513 [9, 2]	ResNet-101-SSD	31.2	50.4	33.3	10.2	34.5	49.8
DSSD513 [2]	ResNet-101-DSSD	33.2	53.3	35.2	13.0	35.4	51.1
RetinaNet [7]	ResNet-101-FPN	39.1	59.1	42.3	21.8	42.7	50.2
RetinaNet [7]	ResNeXt-101-FPN	<b>40.8</b>	<b>61.1</b>	<b>44.1</b>	<b>24.1</b>	<b>44.2</b>	51.2
YOLOv3 608 × 608	Darknet-53	33.0	57.9	34.4	18.3	35.4	41.9

$$mAP = \frac{1}{N} \sum_{i=1}^N AP_i$$

	Pascal 2007 mAP	Speed
DPM v5	33.7	.07 FPS   14 s/img
R-CNN	66.0	.05 FPS   20 s/img
Fast R-CNN	70.0	.5 FPS   2 s/img
Faster R-CNN	73.2	7 FPS   140 ms/img
YOLO	<del>63.4</del> 69.0	45 FPS   22 ms/img

# Object Detection - Yolo

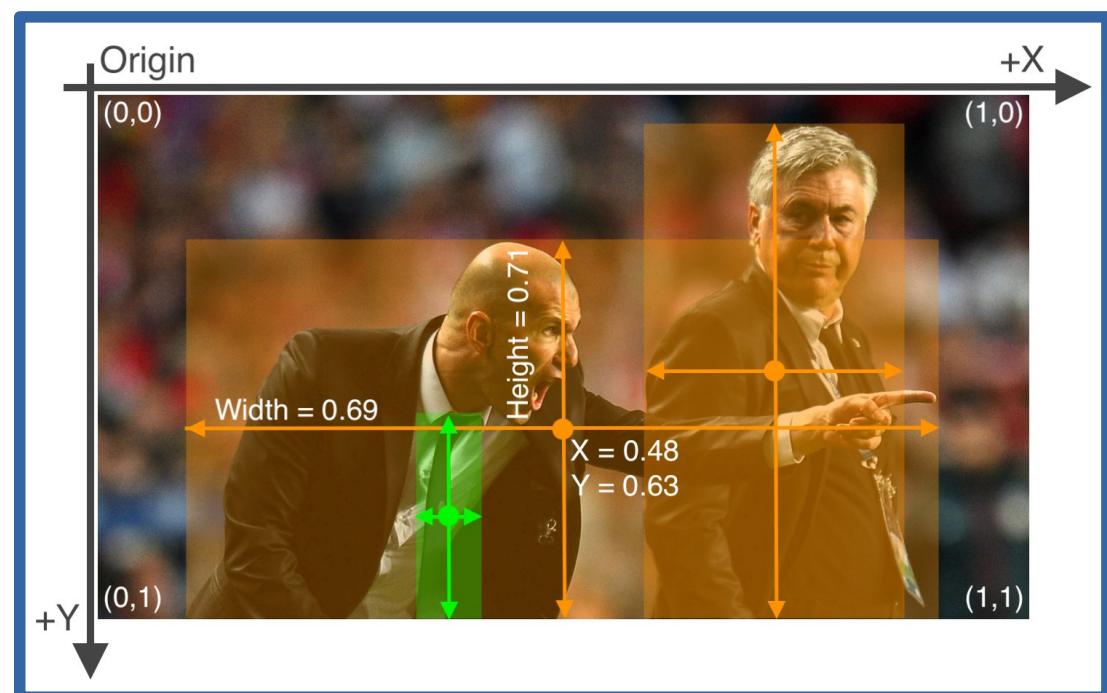
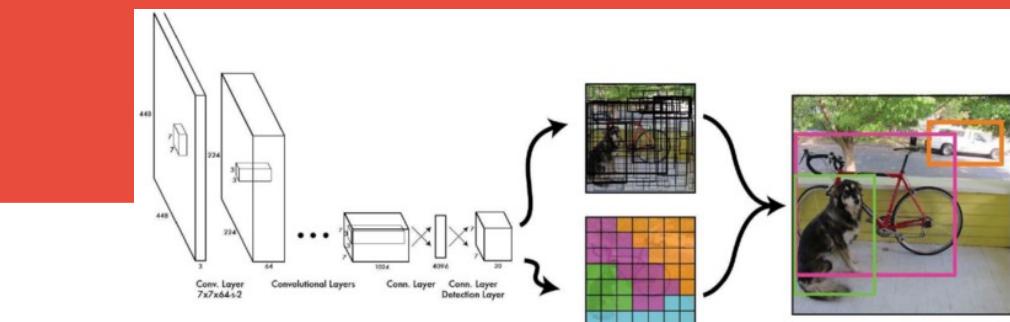
- Let's Code!
- This exercise will utilize Ultralytics ([www.ultralytics.com](http://www.ultralytics.com)) as the framework.
  - Single Image
  - Frame by Frame (Video / Camera)
- Check out the GitHub repository, specifically the `yolo-ultralytics` folder.



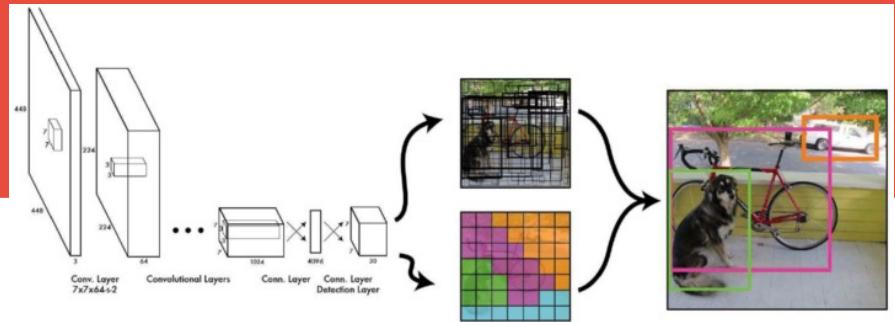
# YoLo - Training and Annotation

# Yolo Training

- Each image must include:
  - Bounding box coordinates (4 points) and the corresponding class label
  - Bounding boxes should be relative (normalized), not absolute pixel coordinates
- Available annotation tools:
  - LabelImg (simple and lightweight)
  - LabelMe (browser-based)
  - CVAT (advanced, collaborative)
  - Roboflow Annotate

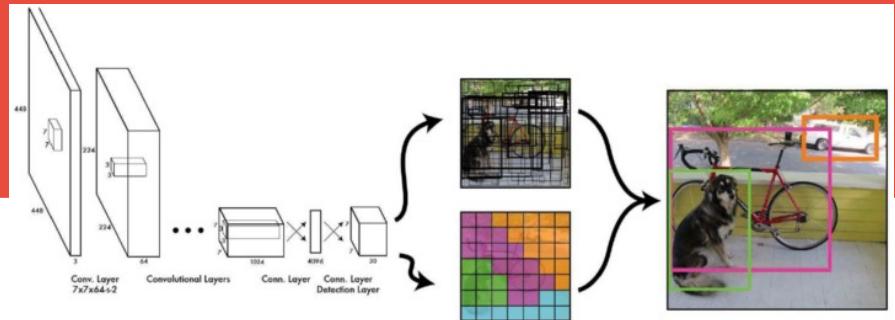


# Yolo Training



5 0.360417 0.459375 0.229167 0.165625  
6 0.575000 0.543750 0.216667 0.162500

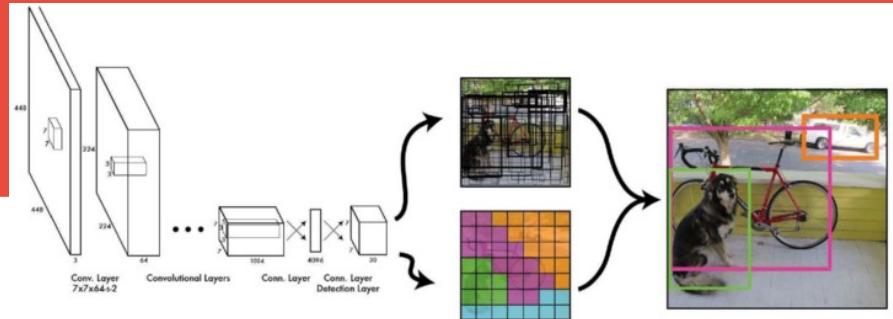
# Yolo Training



5 0.360417 0.459375 0.229167 0.165625  
6 0.575000 0.543750 0.216667 0.162500

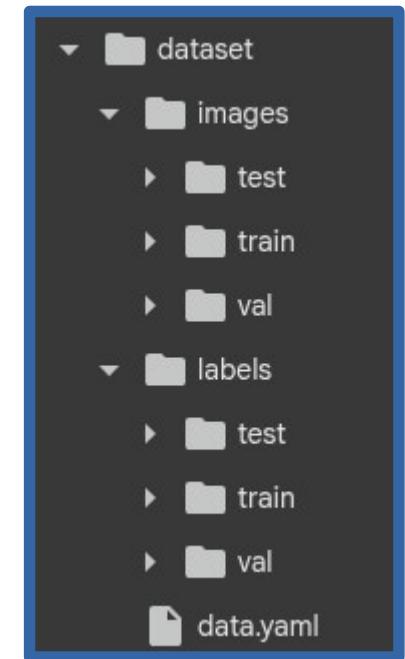
```
class_names = [  
    '1C', #0  
    '2C', #1  
    '5C', #2  
    '10C', #3  
    '20C', #4  
    '50C', #5  
    '1Eu', #6  
    '2Eu' #7  
]
```

# Yolo Training



`data.yaml` ×

```
1 path: /content/dataset
2 train: images/train
3 val: images/val
4 test: images/test
5 nc: 8
6 names:
7 - 1C
8 - 2C
9 - 5C
10 - 10C
11 - 20C
12 - 50C
13 - 1Eu
14 - 2Eu
15
```



Let's Code!

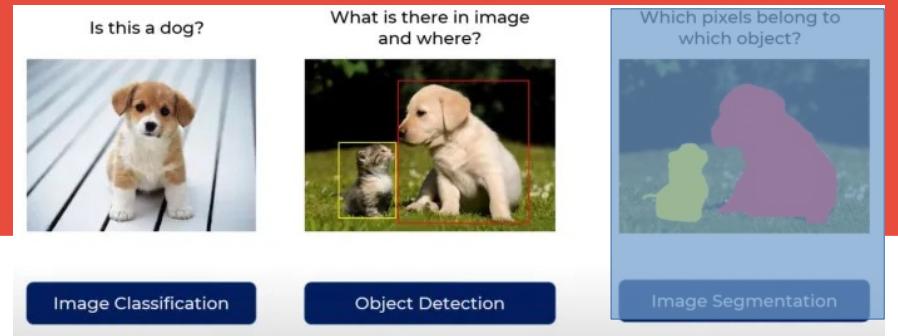
# **U-Net**

# Segmentation Task



- Process of assigning a label to each pixel in an image to partition it into meaningful regions.
  - Precise object localization and delineation
  - Precise object localization and delineation
- State-of-The-Art Architectures
  - FCN (Fully Conv. Networks)(\*)
  - U-NET(\*)
  - DeepLab
  - Mask R-CNN

# Segmentation Task



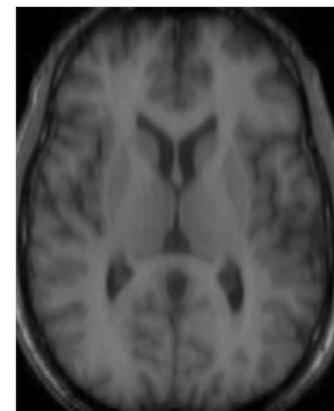
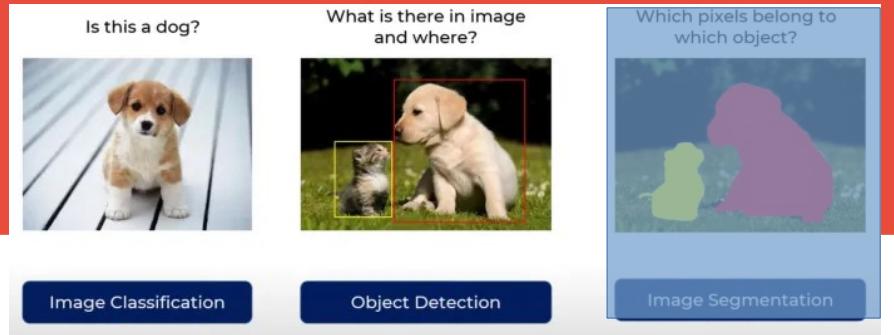
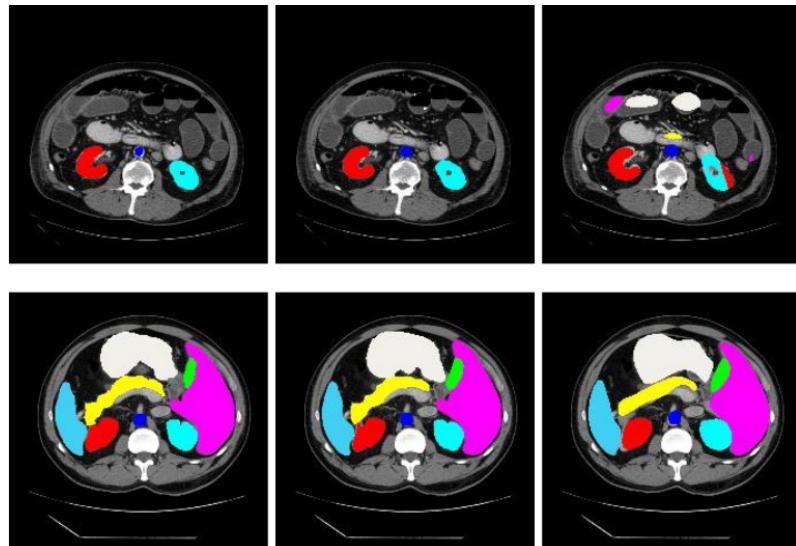
- Types of Segmentation

- Semantic Segmentation: Classify each pixel into predefined classes (e.g., road, car, sky).
- Instance Segmentation: Detect and delineate each object instance separately.
- Panoptic Segmentation: Combines semantic and instance segmentation in a unified view.

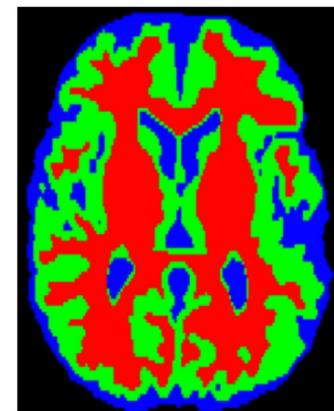


# Segmentation Task

- Applications:
  - Medical imaging: Tumor or organ segmentation



(a) Axial slice

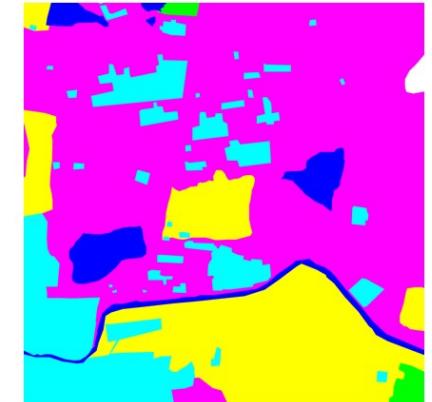


(b) Tissue segmentation

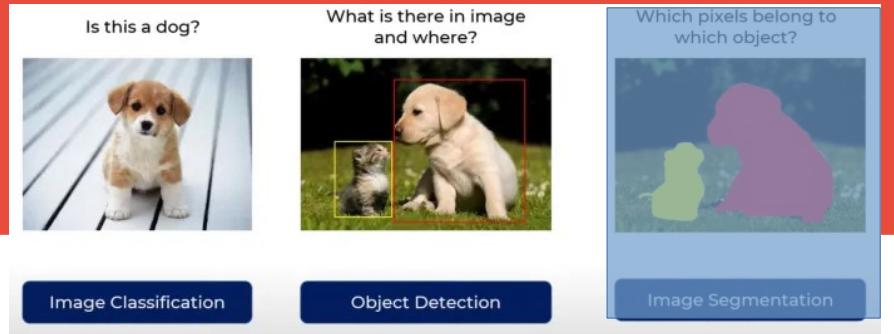
# Segmentation Task



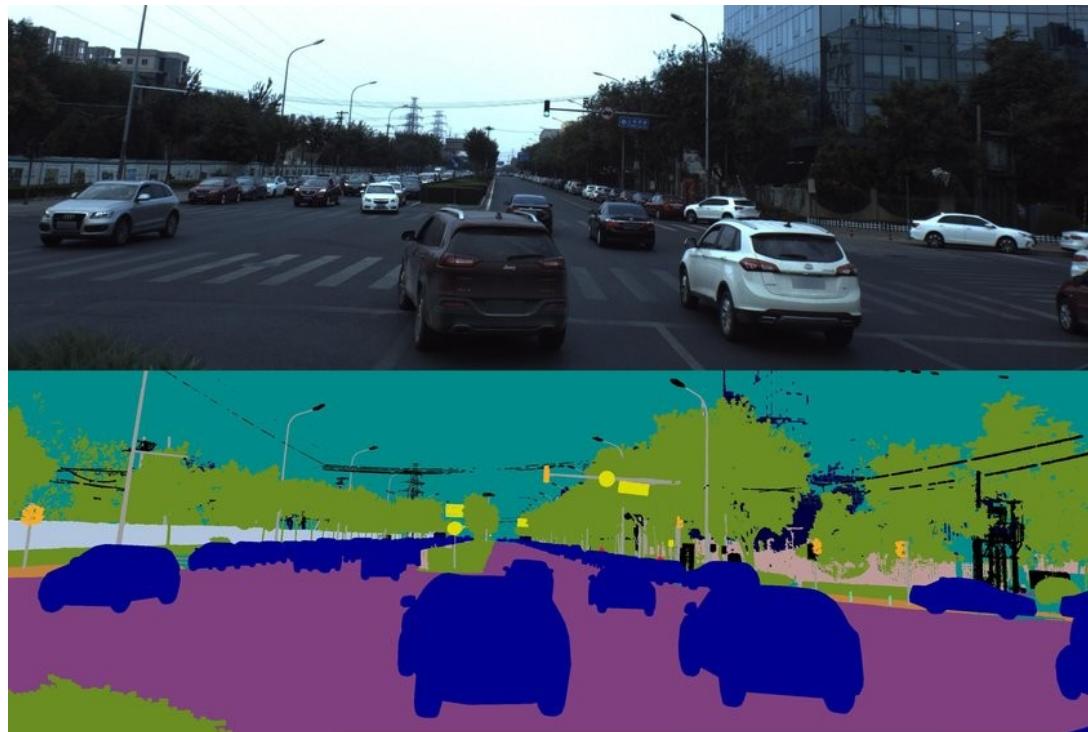
- Applications:
  - Remote sensing: Land cover mapping



# Segmentation Task



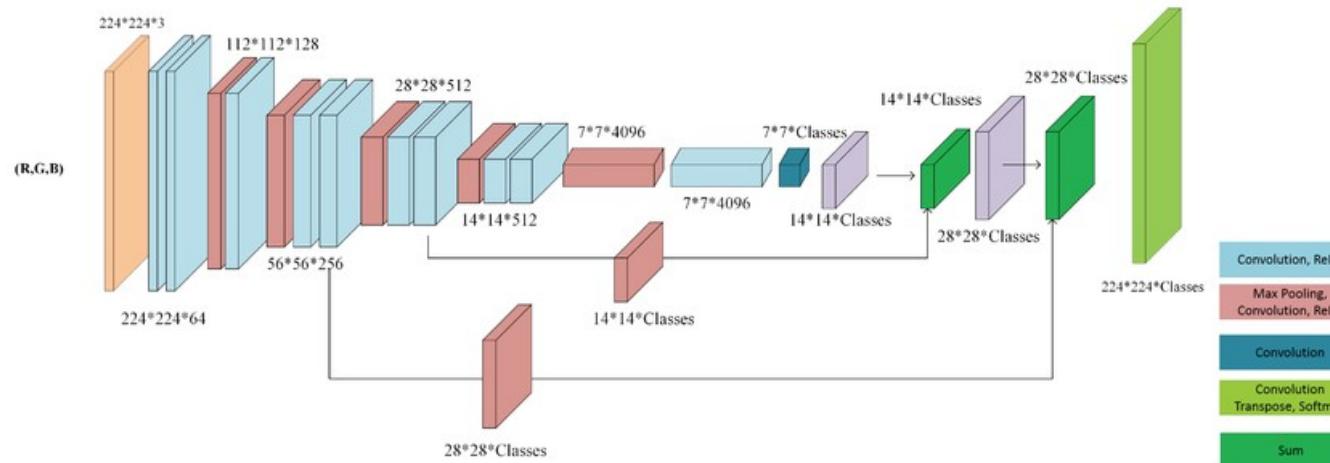
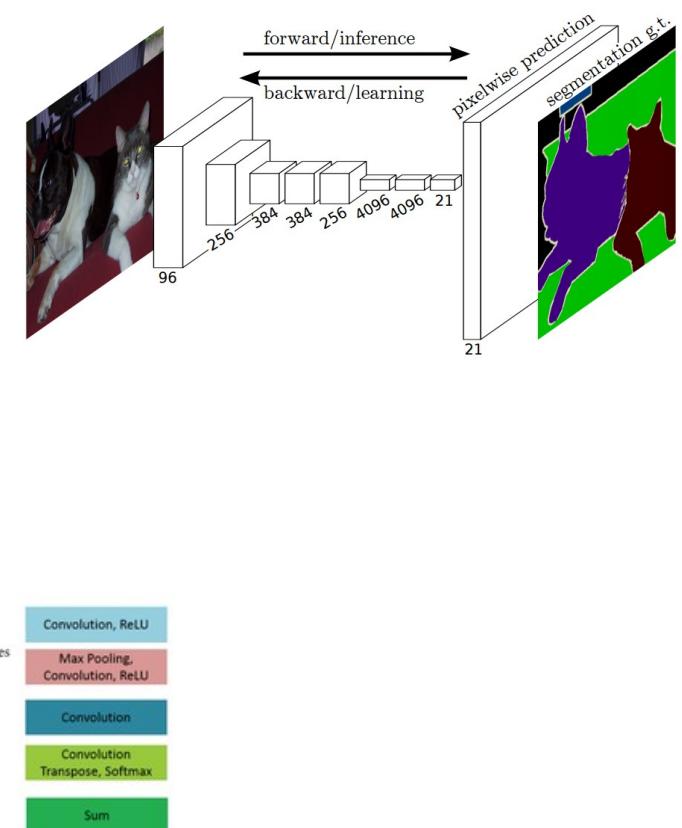
- Applications:
  - Autonomous vehicles: Road scene understanding.



# Fully Convolutional Network (FCN)

- Applications:

- Replace fully connected layers in CNNs with convolutional layers to produce dense pixel-level predictions.
- Encoder derived from classical classification networks (e.g., VGG, AlexNet).
- Decoder built using upsampling (deconvolution) layers to recover spatial resolution.
- Skip connections from shallow layers to refine details



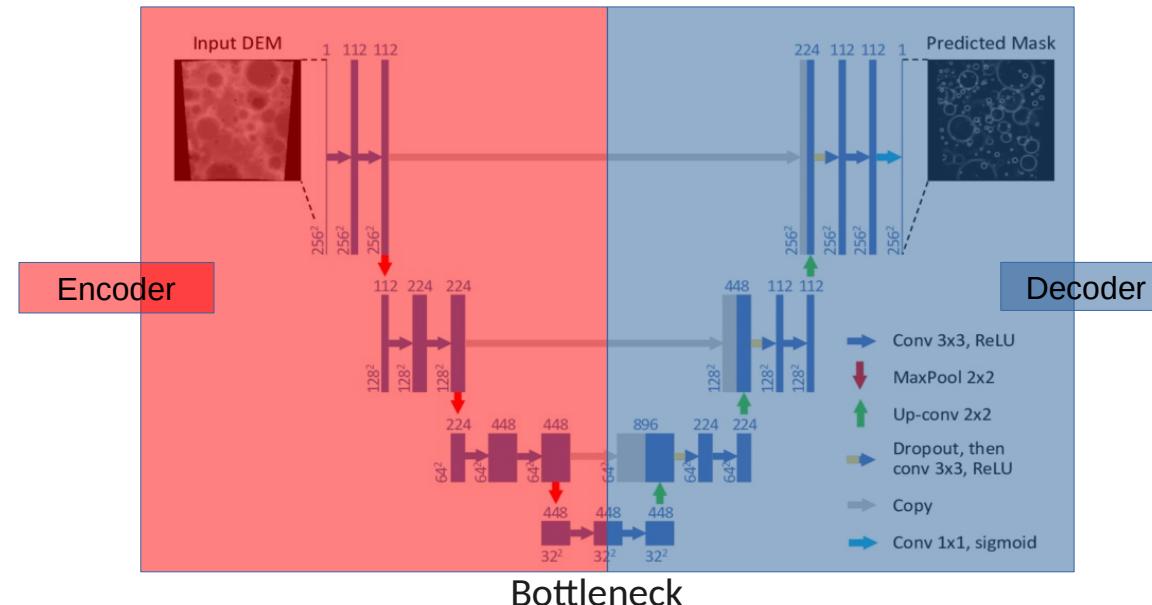
# UNET



- U-Net (Encoder and Decoder)

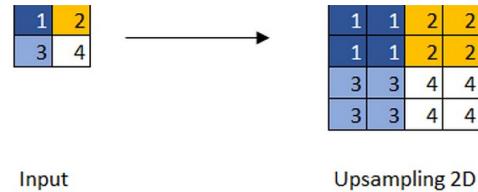
- Symmetric encoder-decoder with dense skip links

- Encoder (contracting path) captures context through convolution + pooling.
    - Decoder (expanding path) upsamples using transposed convolutions.

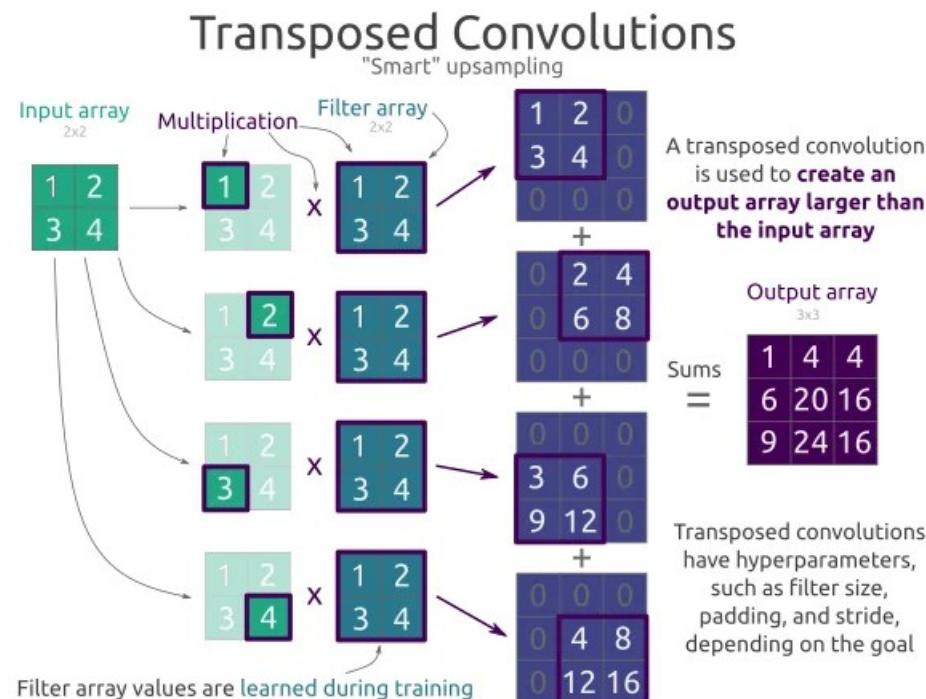


# Learnable Upsampling

- Simple Upsampling: Fixed Interpolation

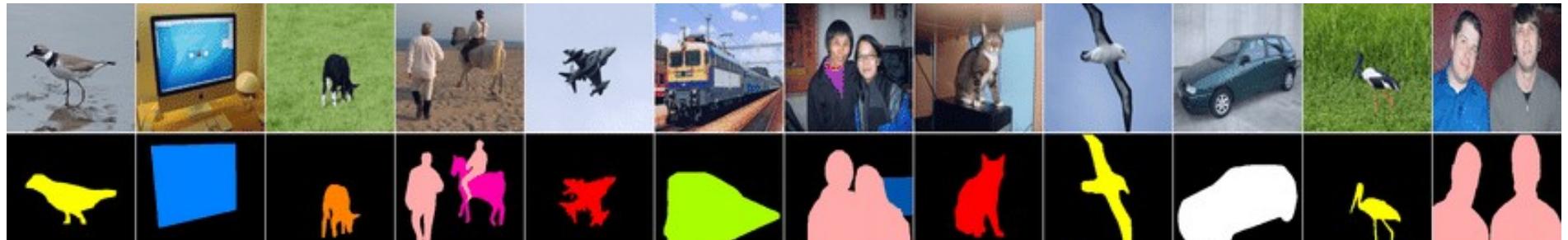
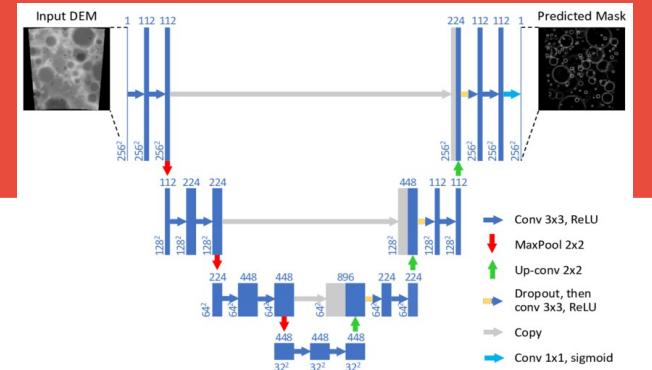


- Learnable Upsampling: Transpose convolution learns the inverse mapping of a convolution to produce larger feature maps.



# UNET

- Training UNET
  - Paired Input → Mask (Segmentation)
  - High Annotation Cost



# UNET

- Training UNET

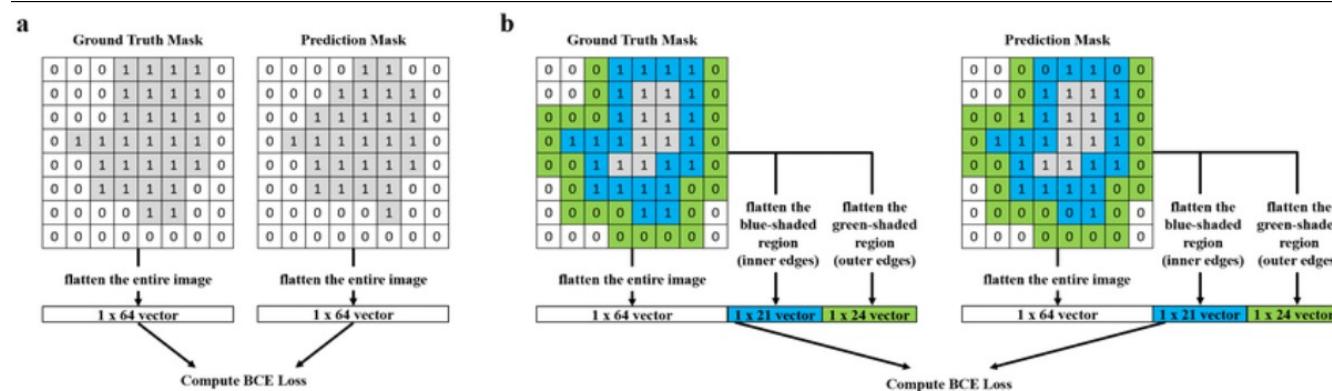
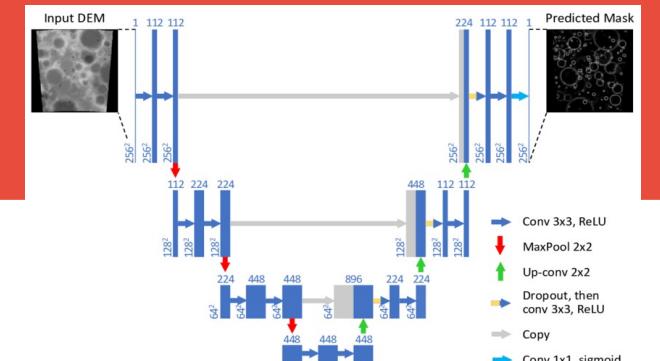
- Losses

- Binary Cross-Entropy (BCE)

- Minimize pixel-level differences

- Treats each pixel independently, penalizing wrong foreground/background classification.

- Sensitive to class imbalance (e.g., small objects vs. large background)



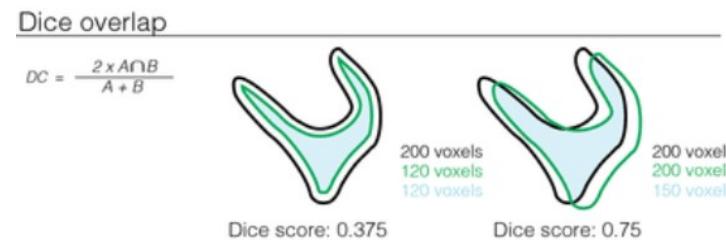
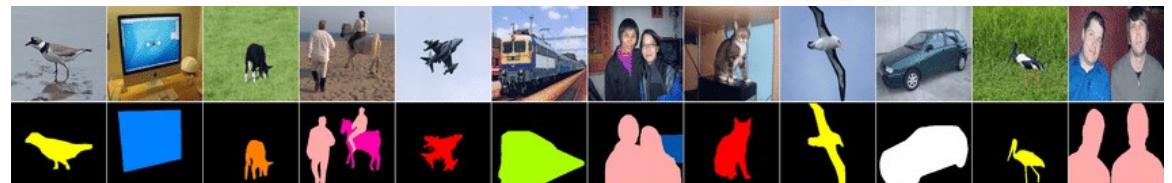
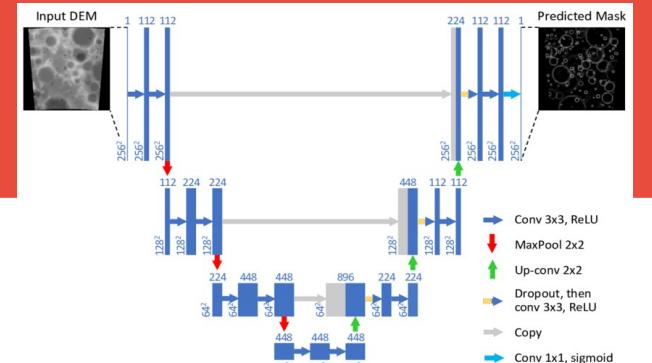
# UNET

- Training UNET

- Losses

- Dice Score

- Maximize the overlap between predicted and ground truth masks.
- Emphasizes correct segmentation of regions, not individual pixels.
- Robust when foreground regions are small (Class Imbalance)



# UNET

- Let's Code!!

