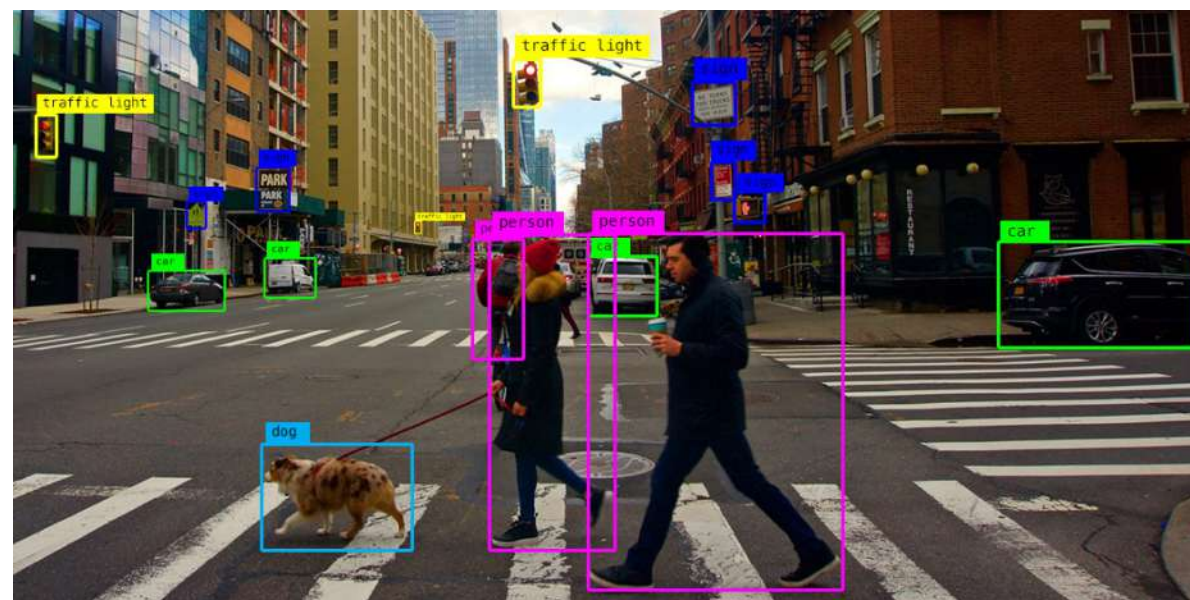


AI VIET NAM  
@aivietnam.edu.vn

# Object Detection

*(Part 2: CNN-based Algorithms)*



[Code & Data](#)

Vinh Dinh Nguyen - PhD in Computer Science

# Outline



➤ **CNN Limitations**

➤ **Region Based Convolutional Neural Networks**

➤ **Spatial Pyramid Pooling**

➤ **Fast R-CNN**

➤ **Faster RCNN**

➤ **YOLOv1-v2**

# Object Detection Milestones

Milestones: Traditional Detectors

*Viola Jones Detectors, SVM + HOG & DPM*

Milestones: CNN based Two-stage Detectors

*RCNN, SPPNet, Faster RCNN, Faster RCNN,..*

Milestones: CNN based One-stage Detectors

*YOLO, SSD, RetinaNet, CornerNet, Center Net,..*

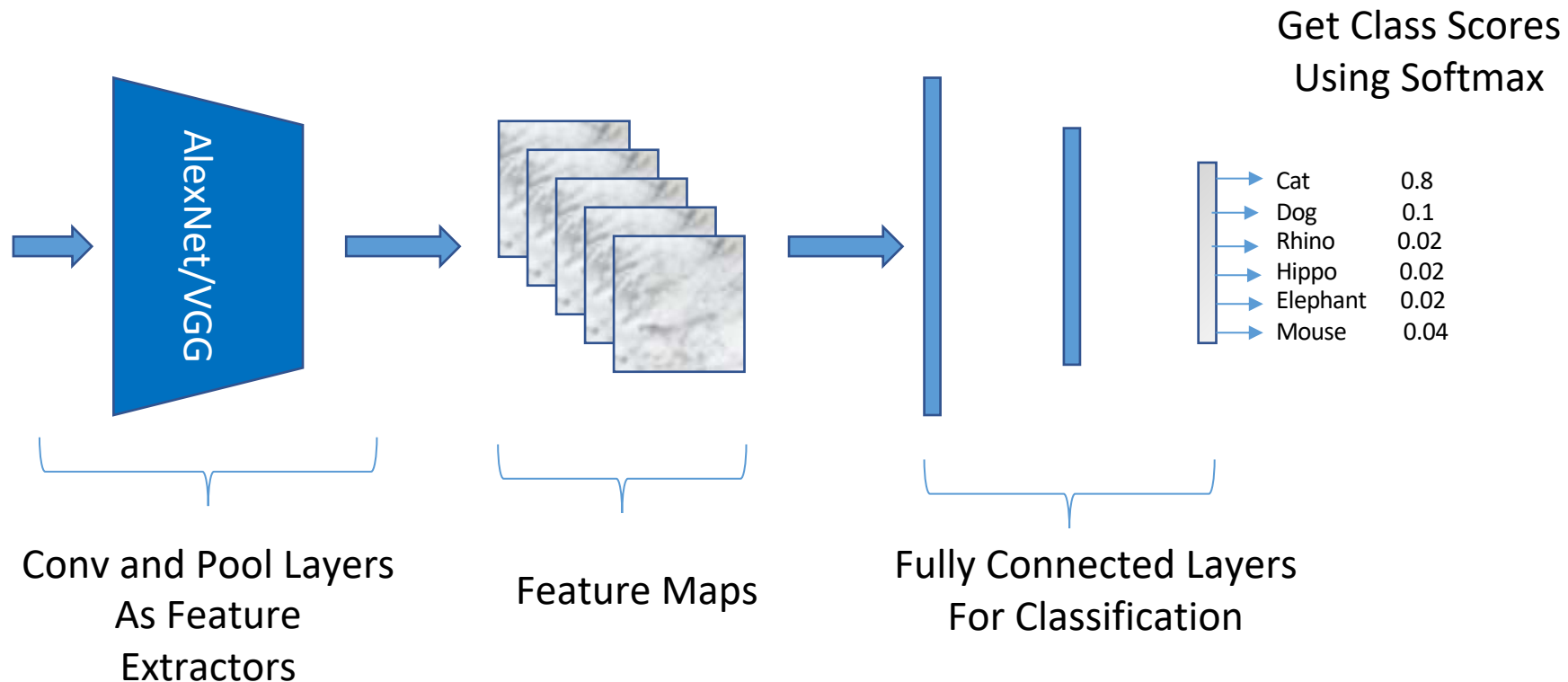
Milestones: Transformer for OD

*DETR, D-DETR, DINO,..*

# Classification Pipeline

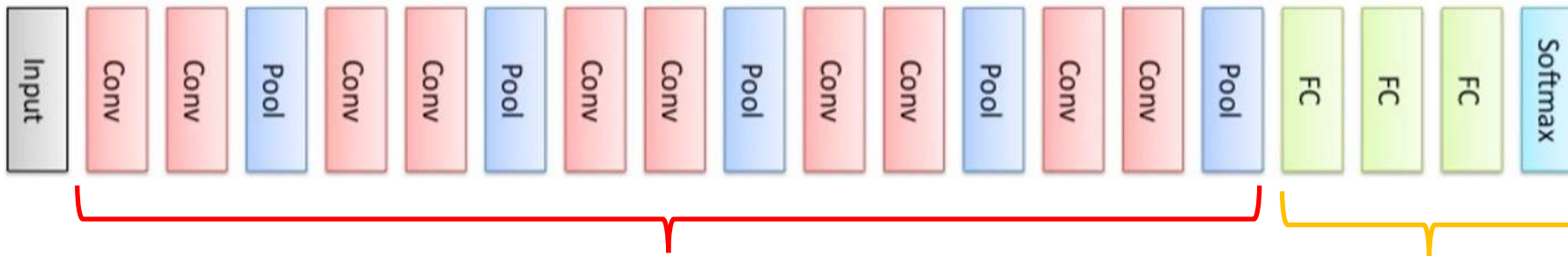


CAT



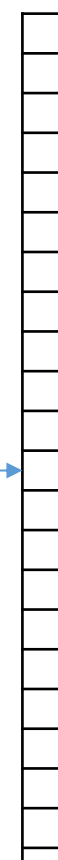
# Classification Pipeline

## VGGNet



Conv and Pool Layers  
As Feature Extractors

Fully Connected Layers  
For Classification

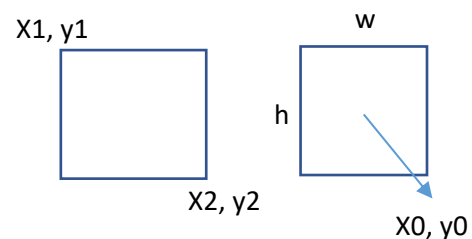
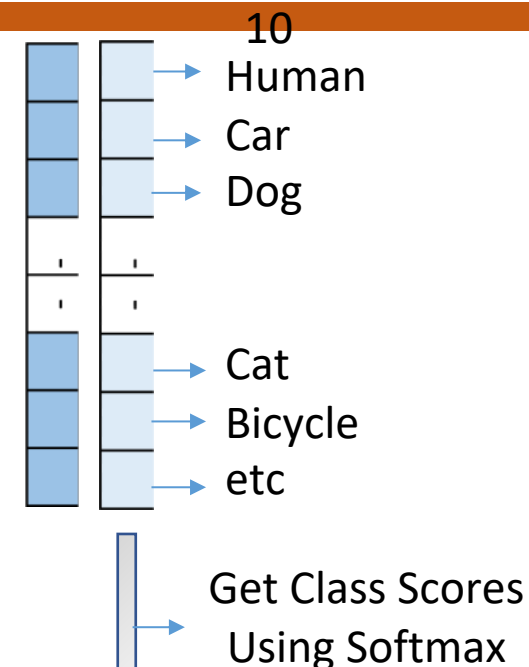
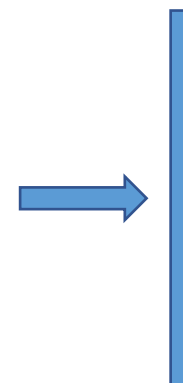
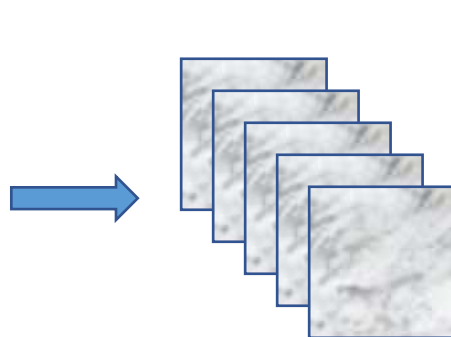
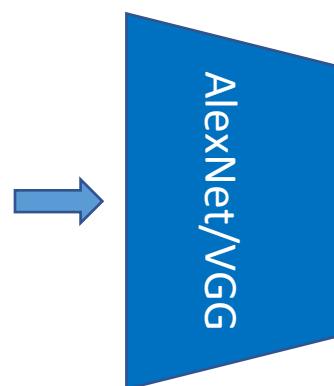


0.8  
0.1  
0.02  
0.03  
0.02  
0.03

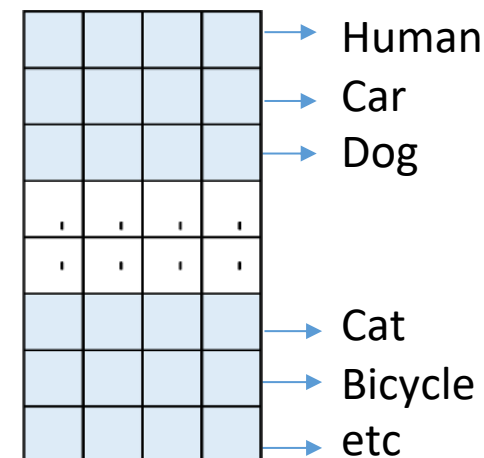


# Ideas for *Localization* using ConvNets

Case #1 – Only one object per image



Get Bounding boxes  
Using L2 loss  
( $x_1, y_1, x_2, y_2$ )



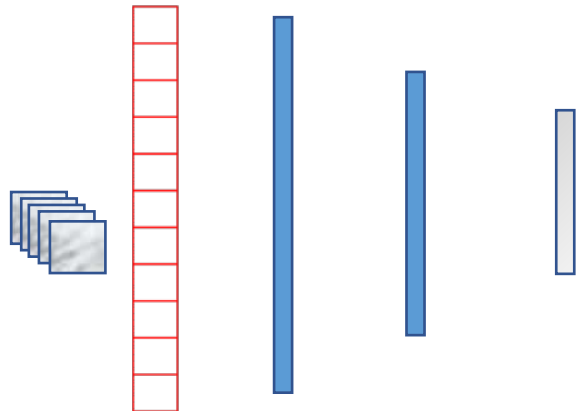
# Bounding Box Regression Training

$(x1,y1) = (200, 250)$   
 $(x2,y2) = (600, 400)$

800



	x1	y1	x2	y2	L2 Loss				
Expected	200	250	600	400					
Prediction	0	0	800	600	$(200-0)^2$	$(250-0)^2$	$(600-800)^2$	$(400-600)^2$	182500
	100	150	700	450	$(200-100)^2$	$(250-150)^2$	$(600-700)^2$	$(400-450)^2$	32500
	210	245	590	405	$(200-210)^2$	$(250-245)^2$	$(600-590)^2$	$(400-405)^2$	250
	200	250	600	400	$(200-200)^2$	$(250-250)^2$	$(600-600)^2$	$(400-400)^2$	0



# About Bounding Boxes



300

500

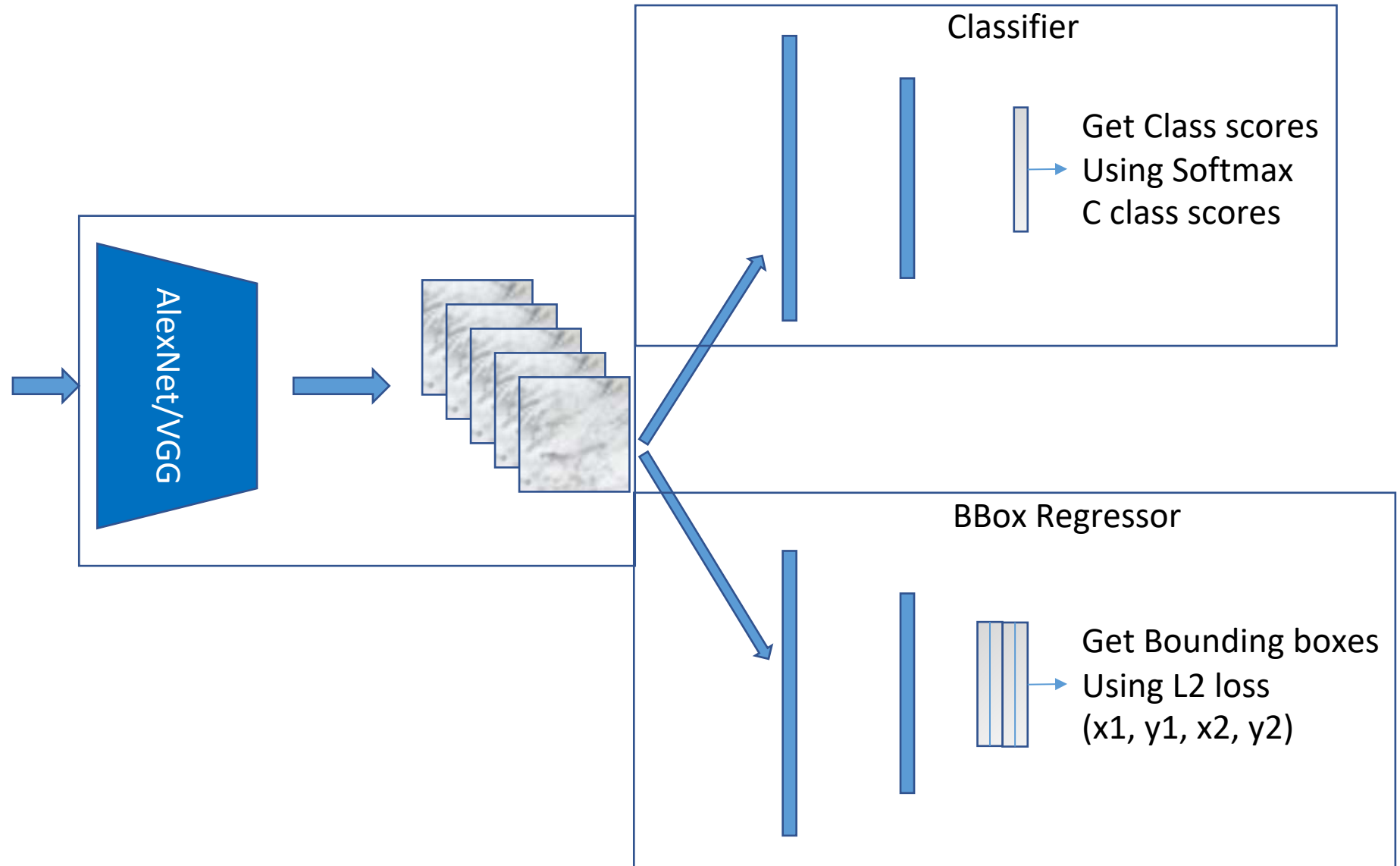


300

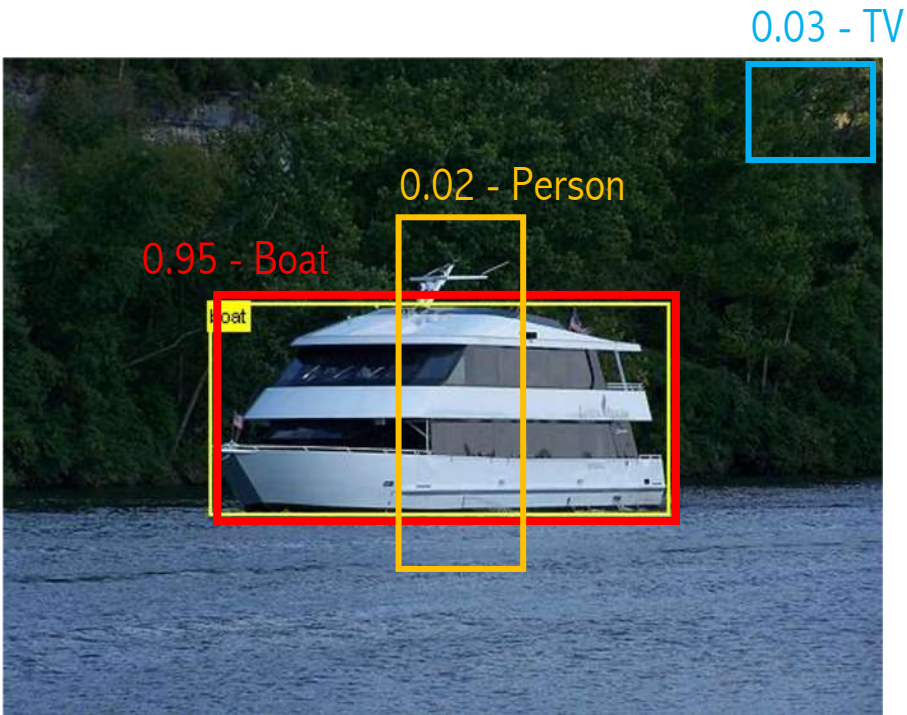
600



# Ideas for *Localization* using ConvNets



# Combining Results

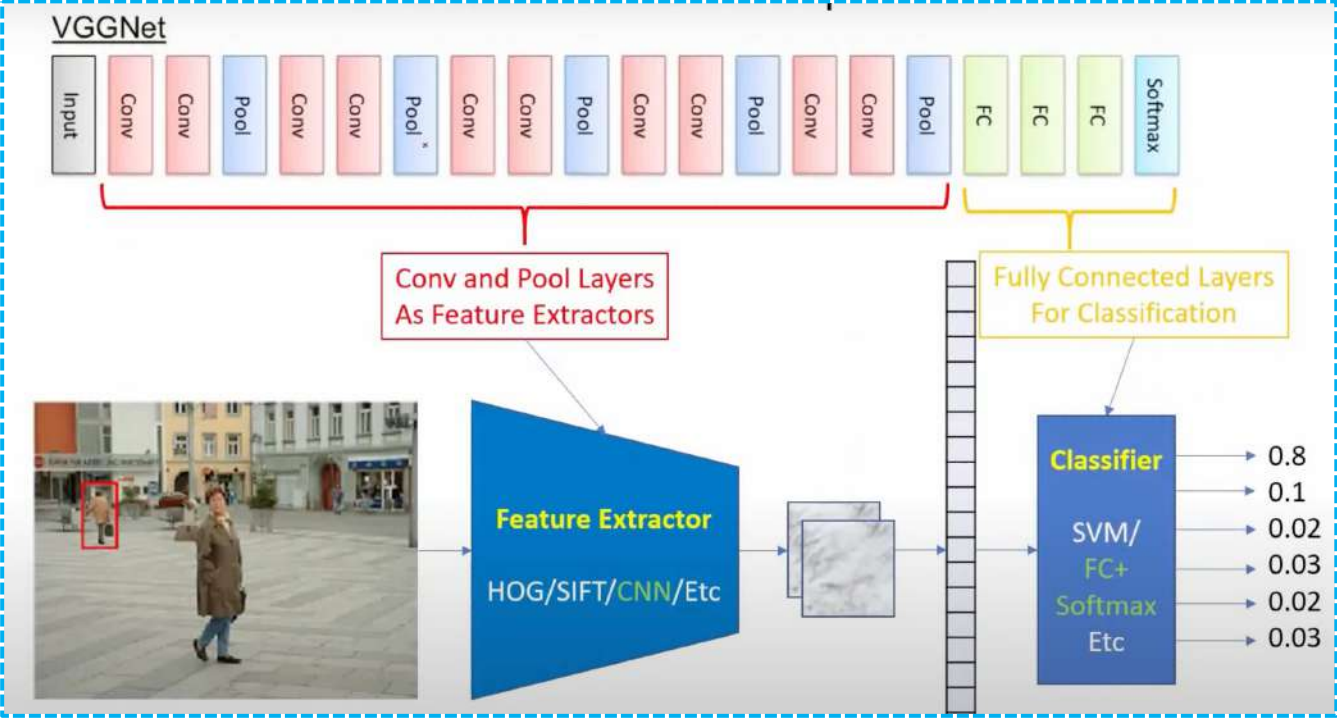
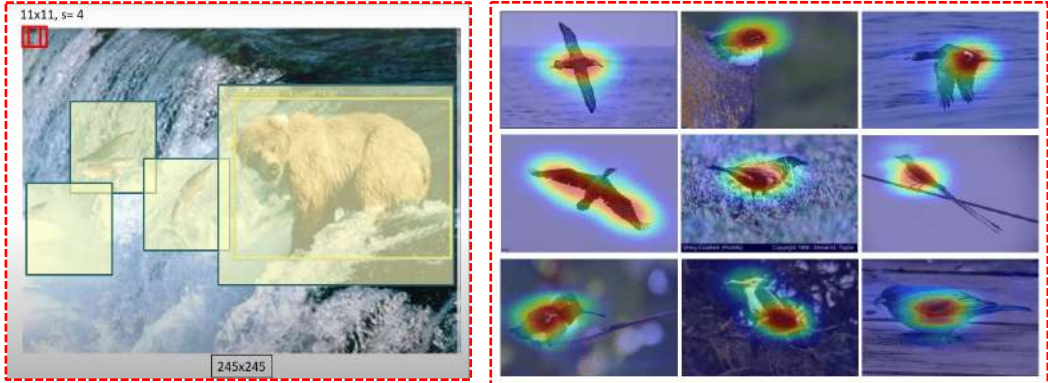
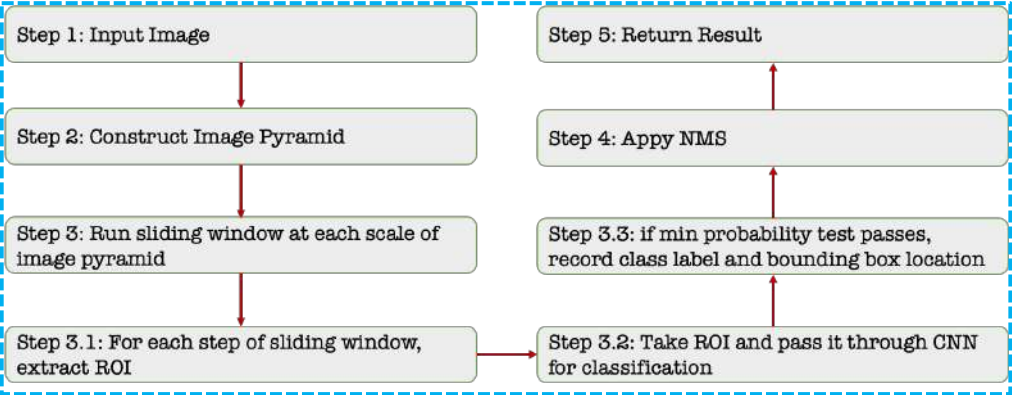


		FC Softmax			
Person					
Boat					
TV					

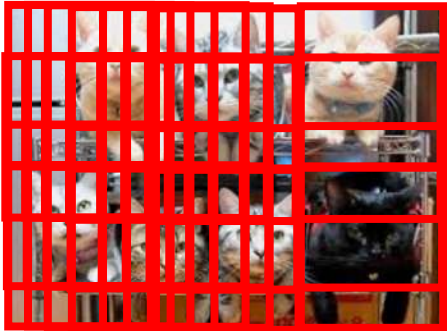
Class	Conf	Bbox coordinates			
Person	0.02	380	200	430	400
Boat	0.95	210	245	590	405
TV	0.03	700	10	790	100

		x1	y1	x2	y2
Person					
Boat					
TV					

# Convolutional Neural Network's Limitations



We need to focus on object only, not entire image





# Region Proposal Methods

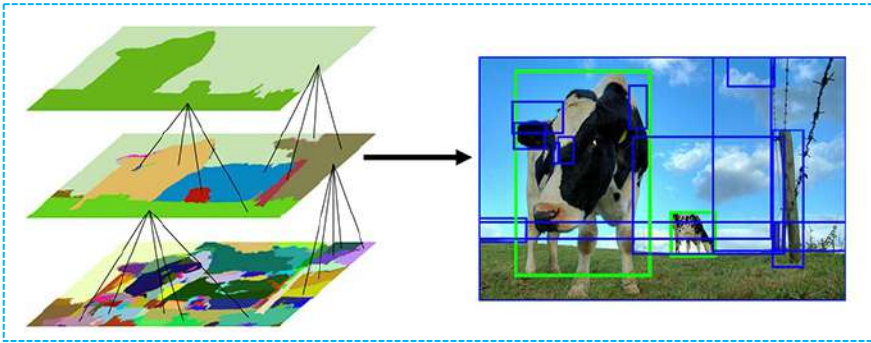
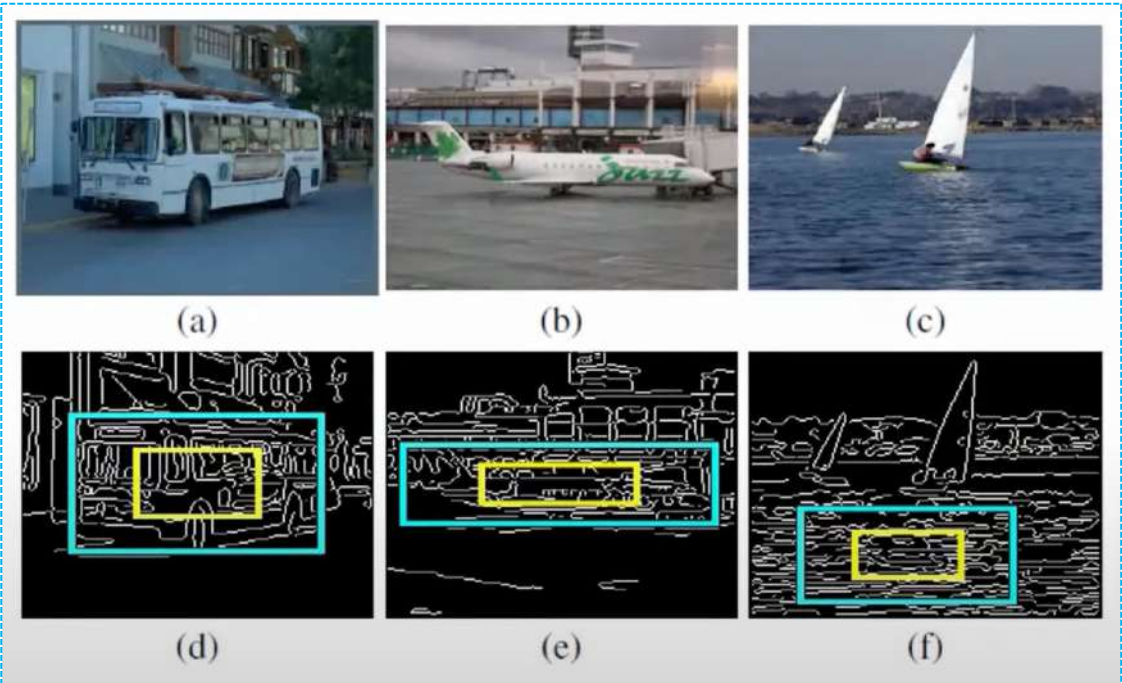


Image Segmentation



Method		Time	Repeatability	Recall	Detection	mAP
Objectness[1]	O	3	.	*	.	25.0/25.4
CPMC[4]	C	250	-	**	*	29.9/30.7
Endres2010[9]	E	100	-	**	**	31.2/31.6
Sel.Search[30]	SS	10	**	***	**	31.7/32.3
Rahtu2011[24]	R1	3	.	.	*	29.6/30.4
Rand.Prim[22]	RP	1	*	*	*	30.5/30.9
Bing[6]	B	0.2	***	*	.	21.8/22.4
MCG[3]	M	30	*	***	**	32.4/32.7
Ranta.2014[25]	R4	10	**	.	*	30.7/31.3
EdgeBoxes[33]	EB	0.3	**	***	**	31.8/32.2
Uniform	U	0	.	.	.	16.6/16.9
Gaussian	G	0	.	.	*	27.3/28.0
SlidingWindow	SW	0	***	.	.	20.7/21.5
Superpixels	SP	1	*	.	.	11.2/11.3

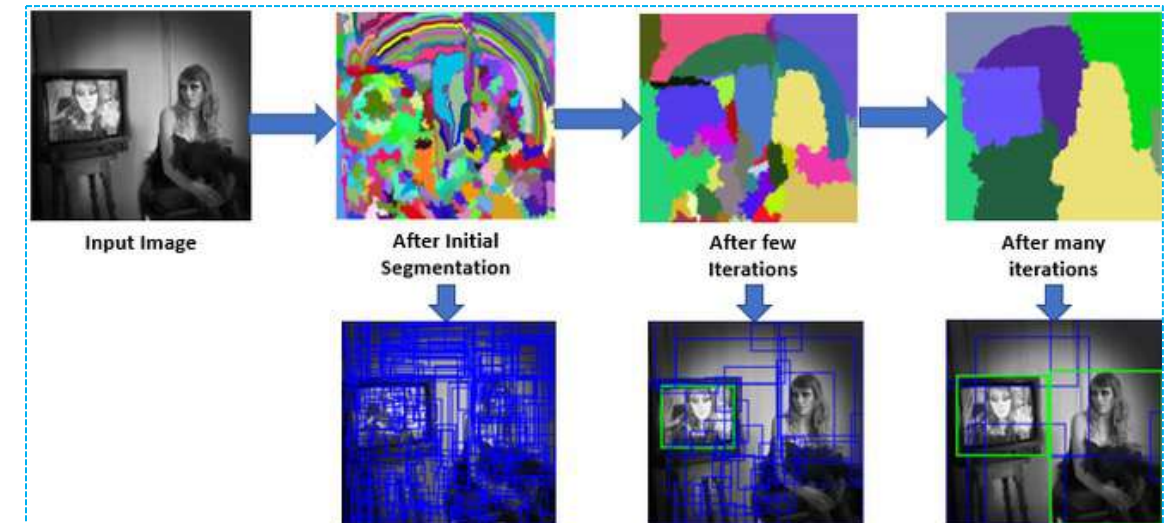


# Edge Boxes & Selective Search



Zitnick, C.L., Dollár, P. (2014). Edge Boxes: Locating Object Proposals from Edges.

How to Apply Edge Boxes / Selective Search for CNN?



# Region Proposal Algorithms

Region proposal algorithms identify prospective objects in an image using segmentation. In segmentation, we group adjacent regions which are similar to each other based on some criteria such as color, texture etc.

RGB Image



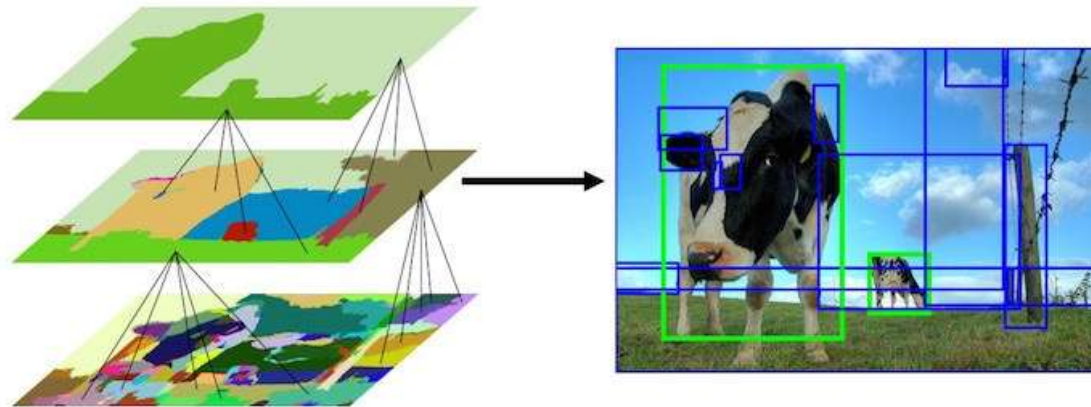
Segmentation Result

Can we use segmented parts in this image as region proposals? The answer is no and there are two reasons why we cannot do that

1. Most of the actual objects in the original image contain 2 or more segmented parts
2. Region proposals for occluded objects such as the plate covered by the cup or the cup filled with coffee cannot be generated using this method

# Selective Search

Selective search uses oversegments from Felzenszwalb and Huttenlocher's method as an initial seed. An oversegmented image looks like this.



Selective Search uses 4 similarity measures based on color, texture, size and shape compatibility.

1. Add all bounding boxes corresponding to segmented parts to the list of regional proposals
2. Group adjacent segments based on similarity
3. Go to step 1

25×3 = 75-dimensional color descriptor

$$s_{color}(r_i, r_j) = \sum_{k=1}^n \min(c_i^k, c_j^k)$$

$$s_{size}(r_i, r_j) = 1 - \frac{size(r_i) + size(r_j)}{size(im)}$$

10×8×3 = 240-dimensional feature descriptor

$$s_{texture}(r_i, r_j) = \sum_{k=1}^n \min(t_i^k, t_j^k)$$

$$s_{fill}(r_i, r_j) = 1 - \frac{size(BB_{ij}) - size(r_i) - size(r_j)}{size(im)}$$

Final Similarity

$$s(r_i, r_j) = a_1 s_{color}(r_i, r_j) + a_2 s_{texture}(r_i, r_j) + a_3 s_{size}(r_i, r_j) + a_4 s_{fill}(r_i, r_j)$$



# Selective Search

Selective search uses oversegments from Felzenszwalb and Huttenlocher's method as an initial seed. An oversegmented image looks like this.



Dogs: top 250 region proposals



Breakfast Table: top 200 region proposals

1. Add all bounding boxes corresponding to segmented parts to the list of regional proposals
2. Group adjacent segments based on similarity
3. Go to step 1

<https://learnopencv.com/selective-search-for-object-detection-cpp-python/>



# Selective Search

```
def generate_region_by_selective_search(img_file):  
  
    img_ss = cv2.imread(img_file)  
  
    ss = cv2.ximgproc.segmentation.createSelectiveSearchSegmentation()  
    ss.setBaseImage(img_ss)  
    ss.switchToSelectiveSearchFast()  
    # "... extract around 2000 region proposals (we use selective search's "fast mode" in all experiments)."  
    rects = ss.process()  
  
    print('Found', len(rects), 'boxes...')  
    for i, rect in enumerate(rects):  
        if i > 2000:  
            break  
        x, y, w, h = rect  
        cv2.rectangle(img_ss, (x, y), (x+w, y+h), (100, 255, 100), 1)  
  
    cv2.imshow(img_ss)  
    # close image show window
```



# Outline



➤ **CNN Limitations**

➤ **Region Based Convolutional Neural Networks**

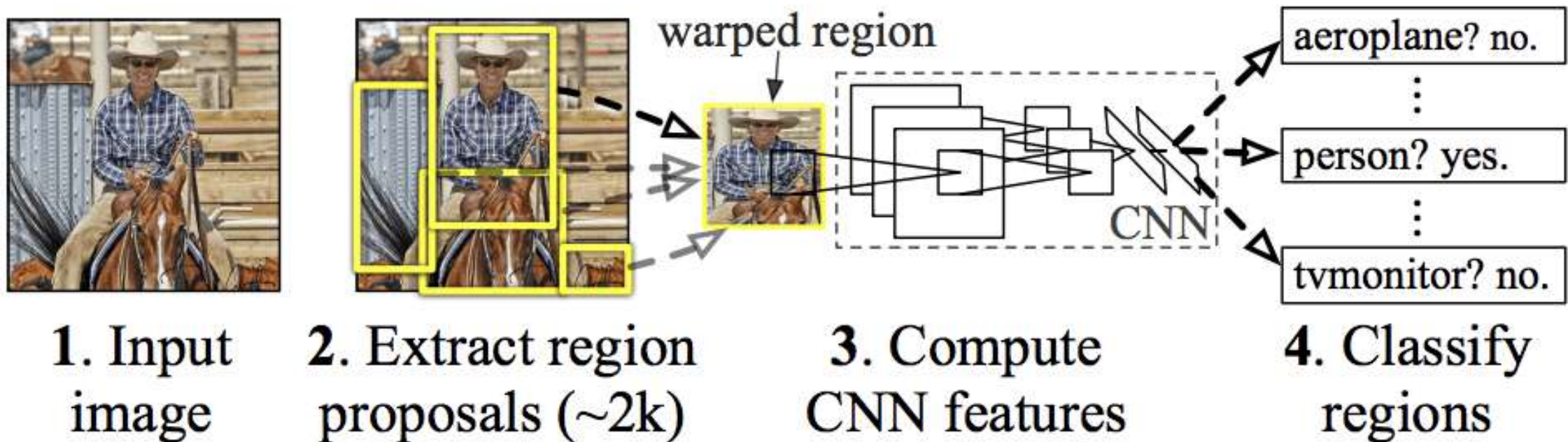
➤ **Spatial Pyramid Pooling**

➤ **Fast R-CNN**

➤ **Faster RCNN**

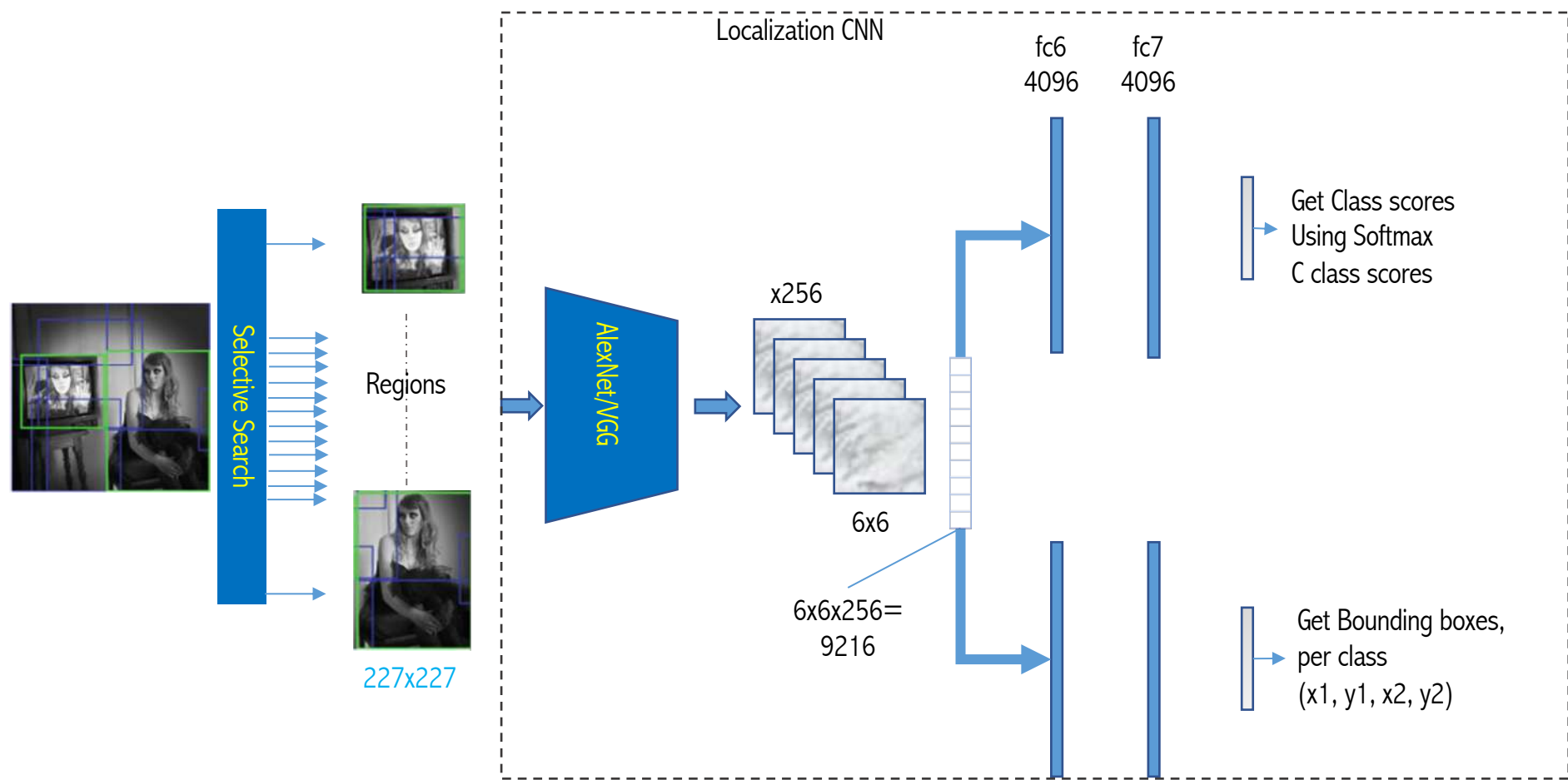
➤ **YOLOv1-v2**

## R-CNN: *Regions with CNN features*



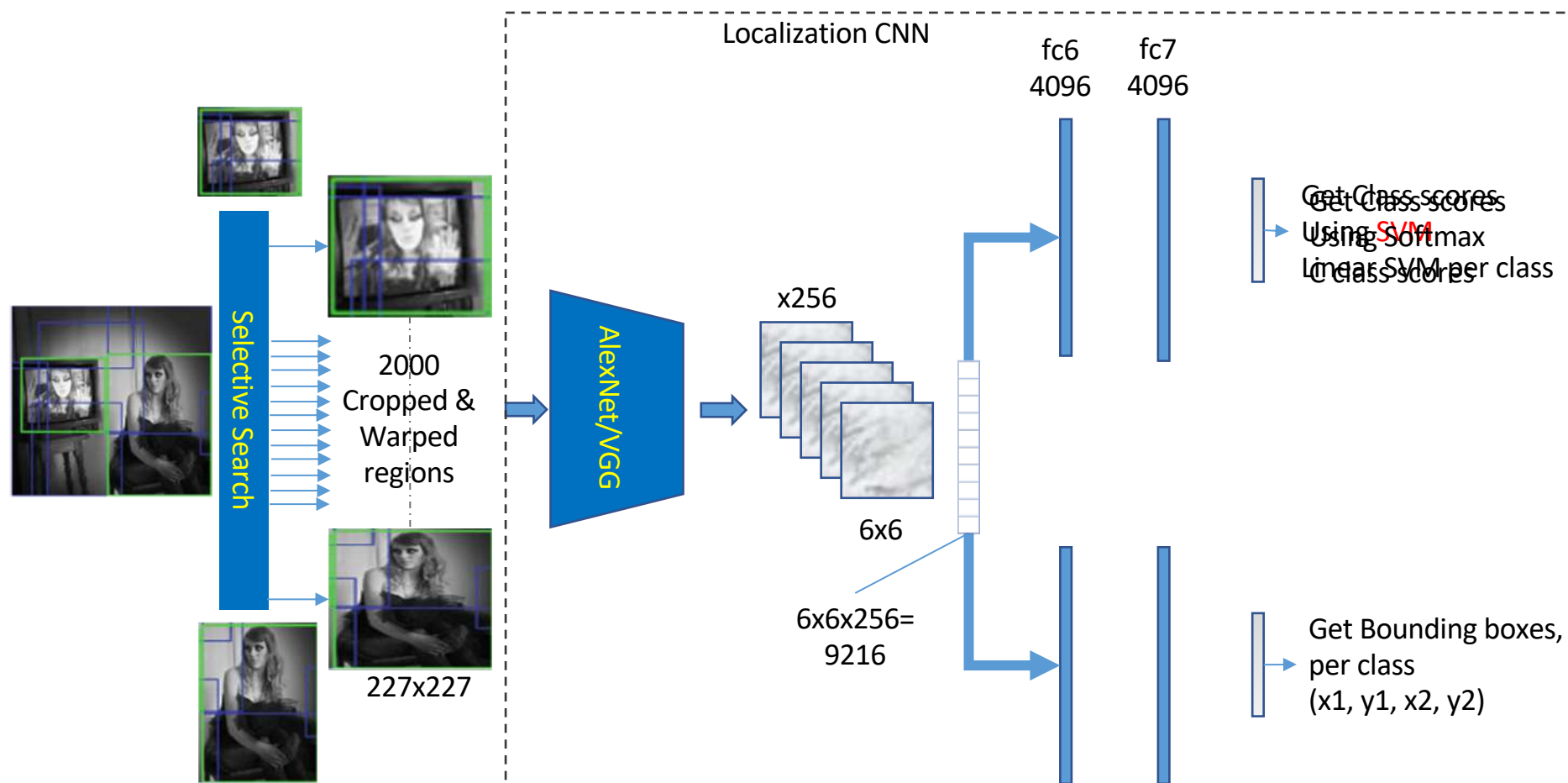
*Ross Girshick, et al. from UC Berkeley titled "Rich feature hierarchies for accurate object detection and semantic segmentation."*

# RCNN - Region proposals with CNNs

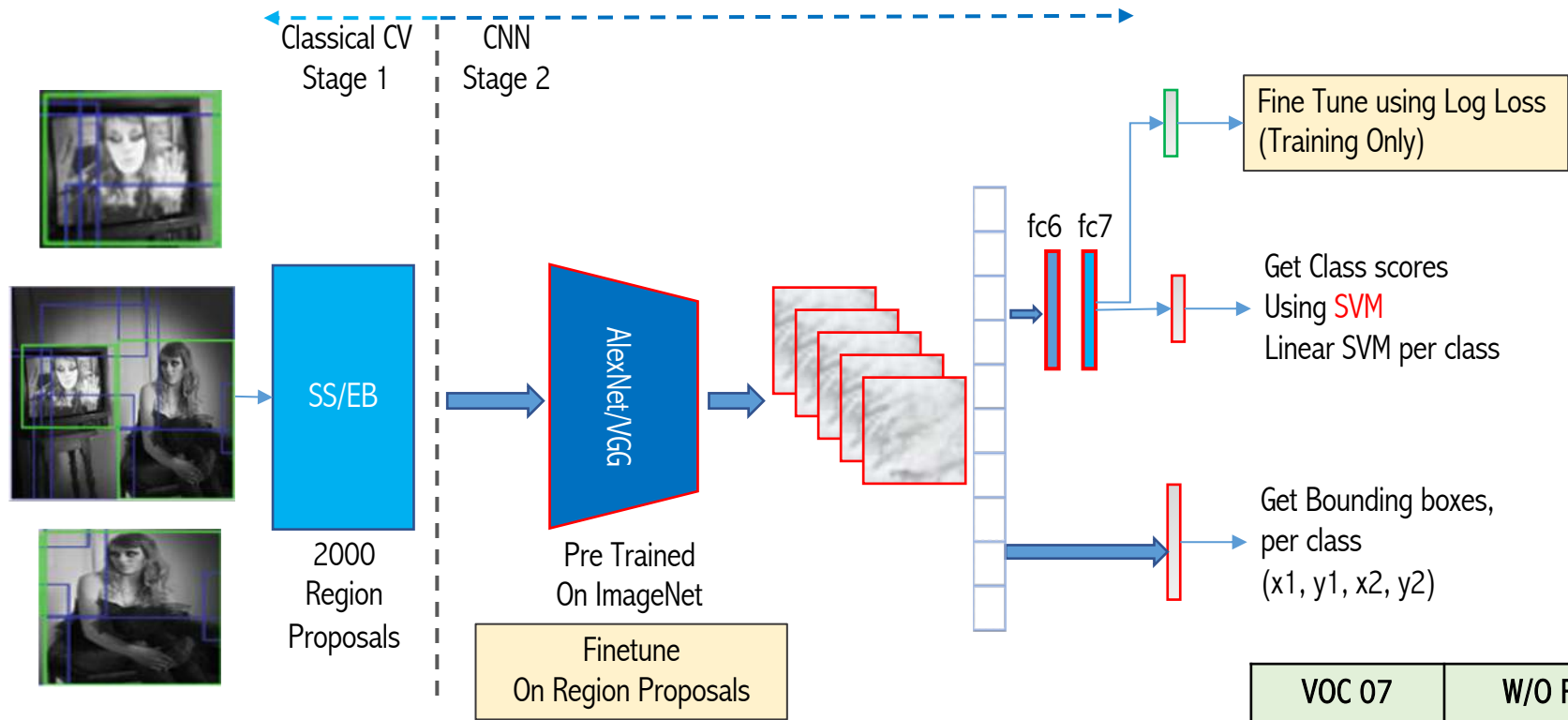




# RCNN - Region proposals with CNNs



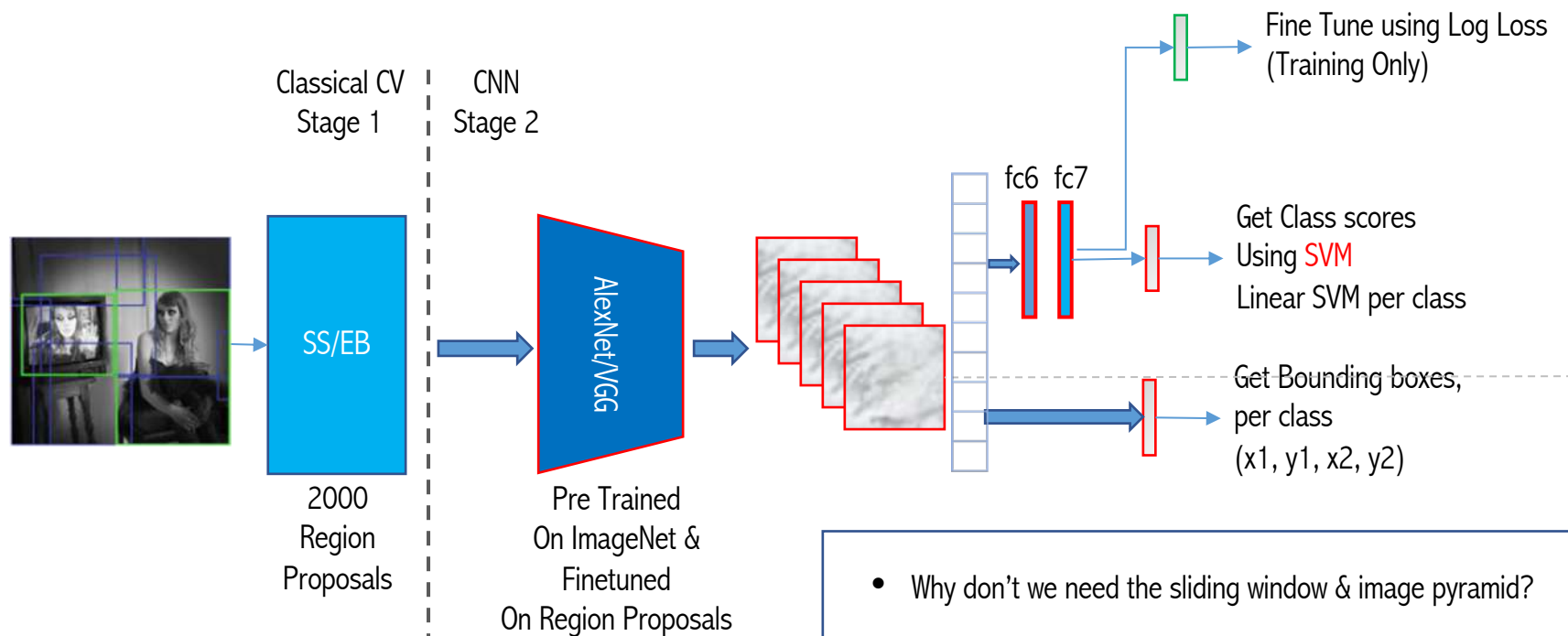
# RCNN - Region proposals with CNNs



- Alex Net
- Before finetune: 44%
  - After finetune: 54%
  - Adding bounding box regression: 58%
    - VGG: 66%

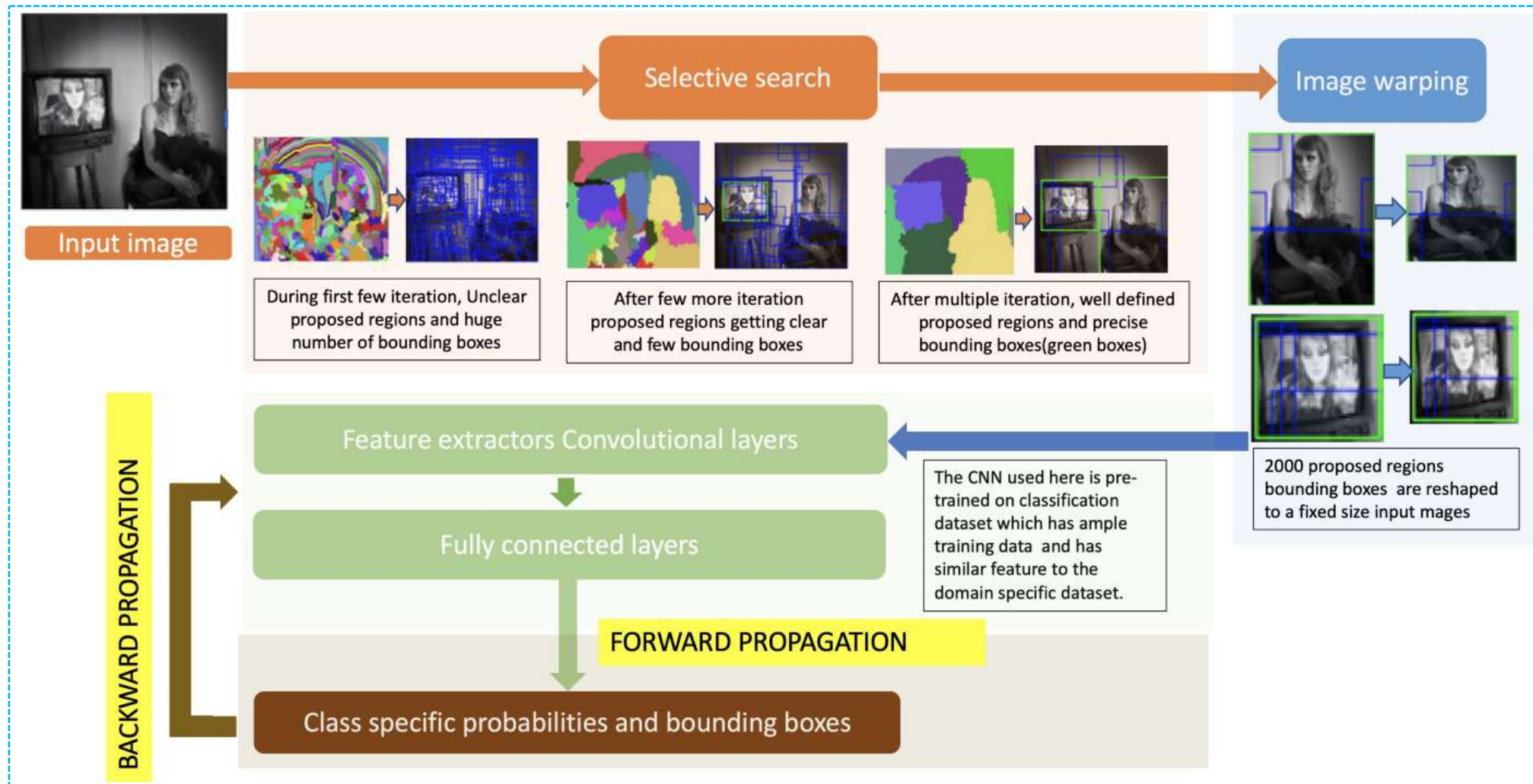
VOC 07	W/O FT	W FT
pool5	44.2	47.3
fc6	46.2	53.1
fc7	44.7	54.2

# RCNN - Region proposals with CNNs



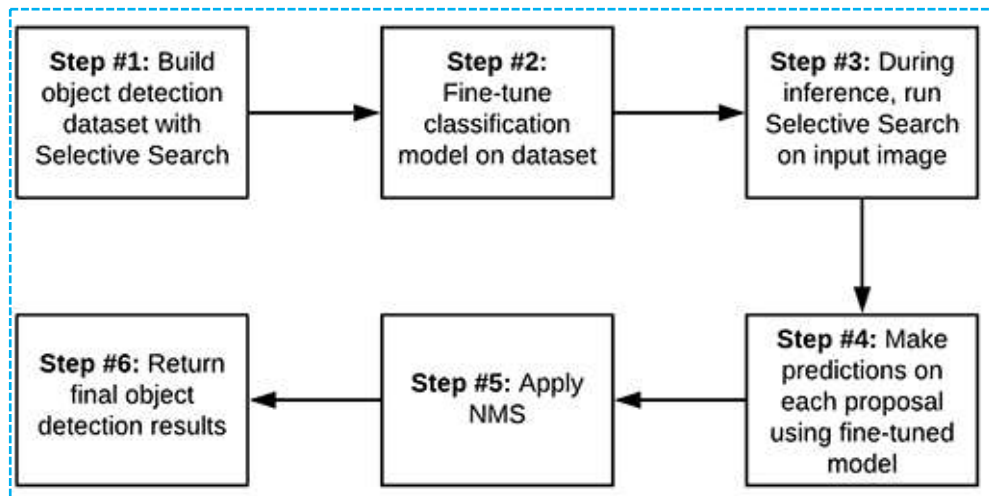
- Why don't we need the sliding window & image pyramid?
- Didn't we end up with too many inputs to the localization network?

# RCNN - Region proposals with CNNs



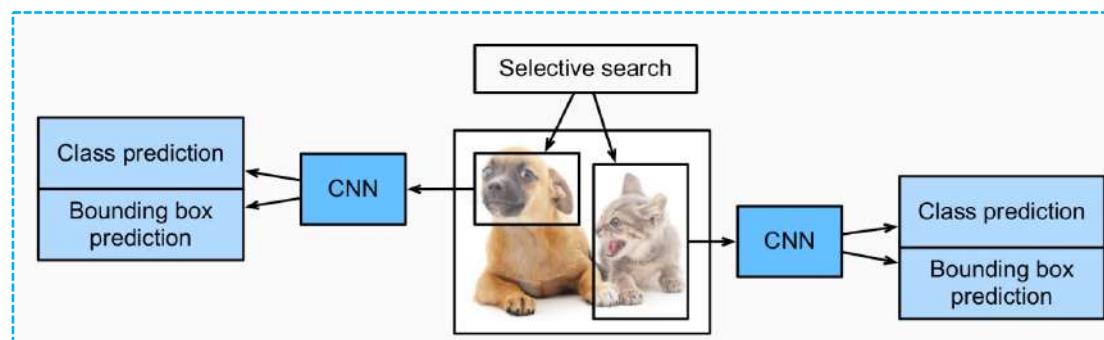
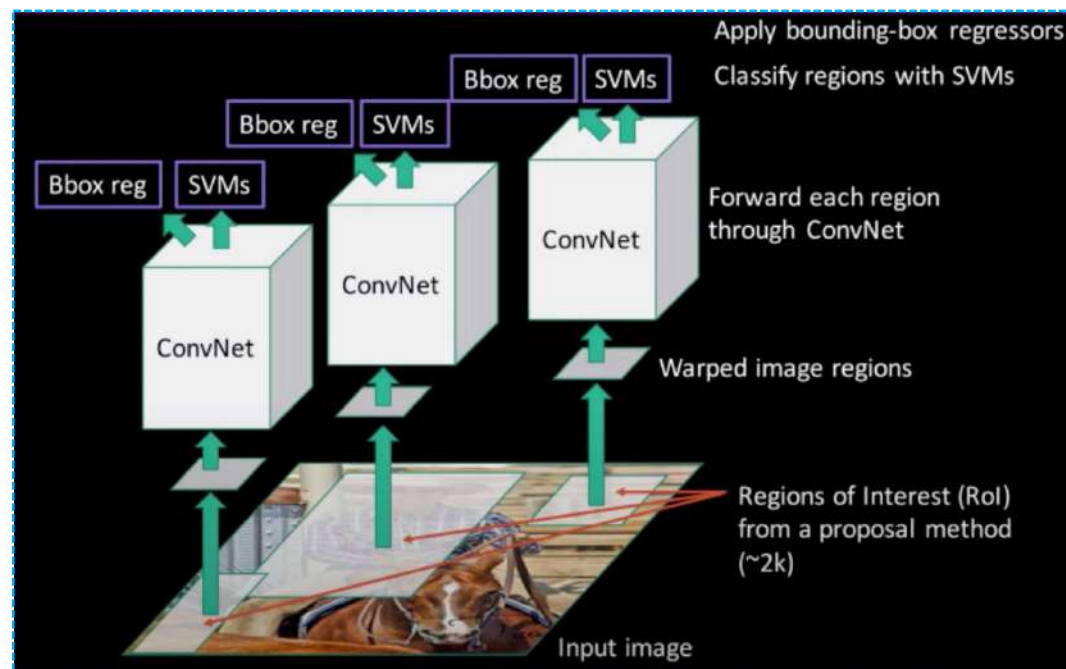


# RCNN - Region proposals with CNNs



1. It consumes a huge amount of time, storage, and computation power.
2. It has a complex multi-stage training pipeline (3 stage — Log loss, SVM, and BBox Regressor's L2 loss).

Why not consider running the CNN just once per image and then find a way to share that computation across the  $\sim 2000$  proposals?



# Outline



➤ **CNN Limitations**

➤ **Region Based Convolutional Neural Networks**

➤ **Spatial Pyramid Pooling**

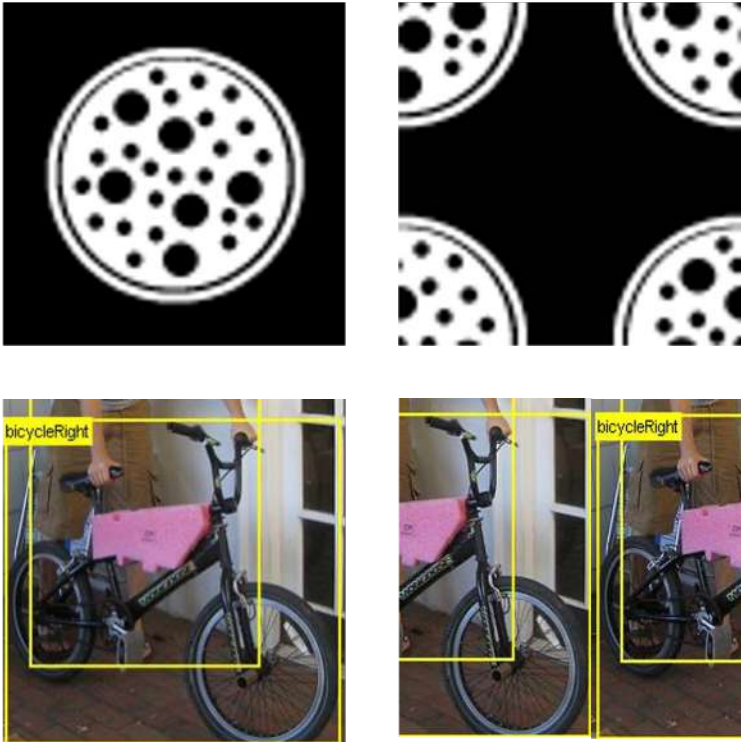
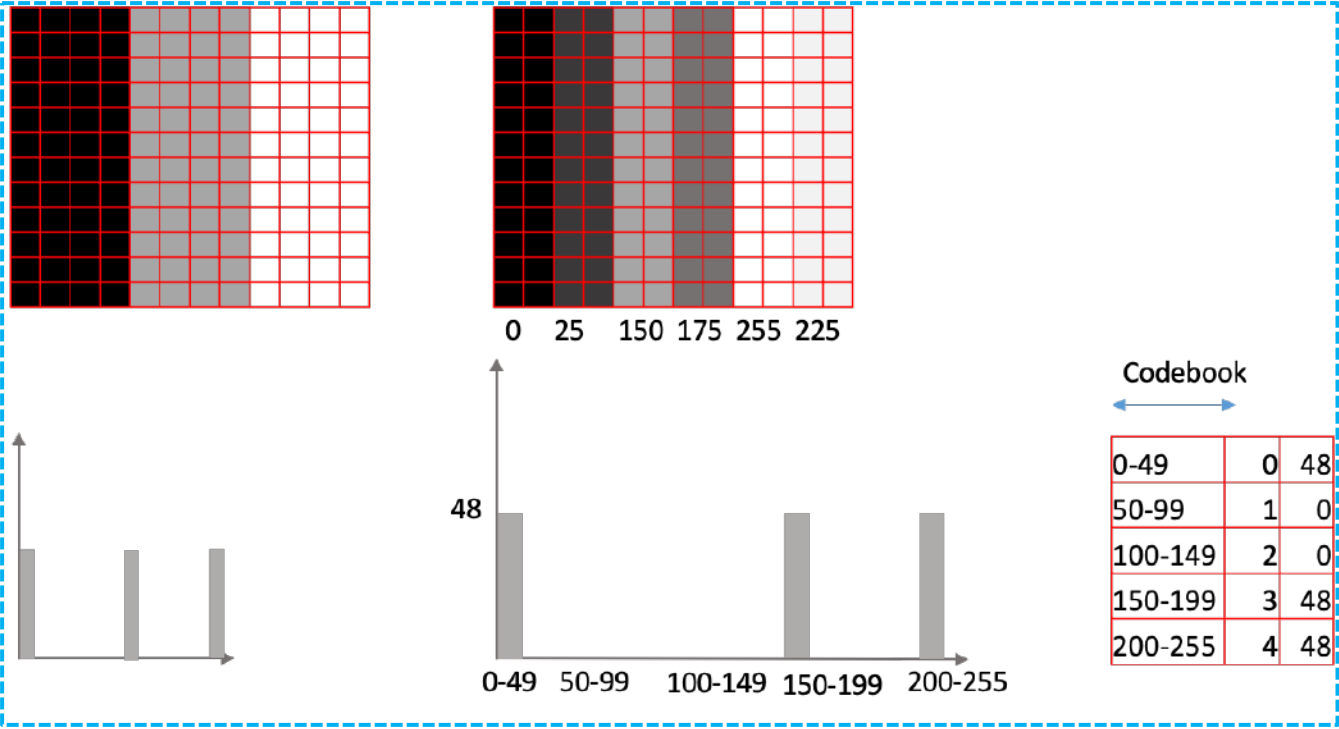
➤ **Fast R-CNN**

➤ **Faster RCNN**

➤ **YOLOv1-v2**

# Spatial Pyramid Matching

Histograms of Images - Bins



LIMITATION

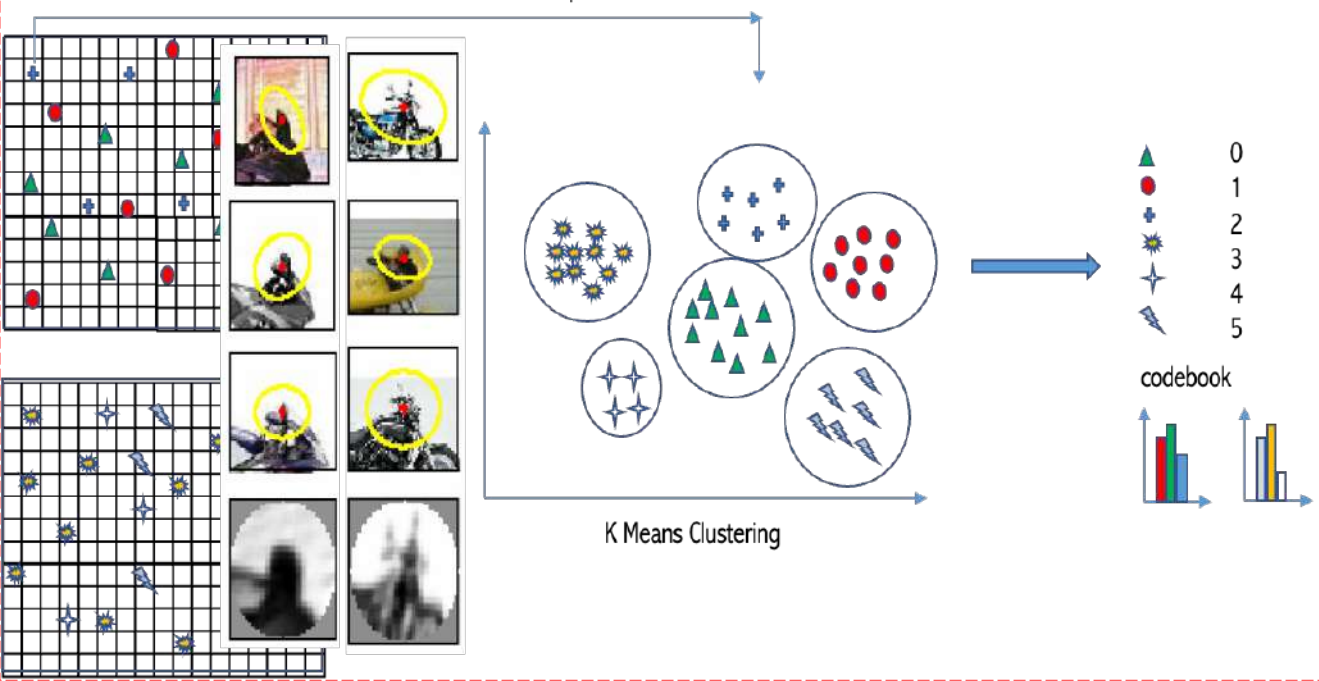




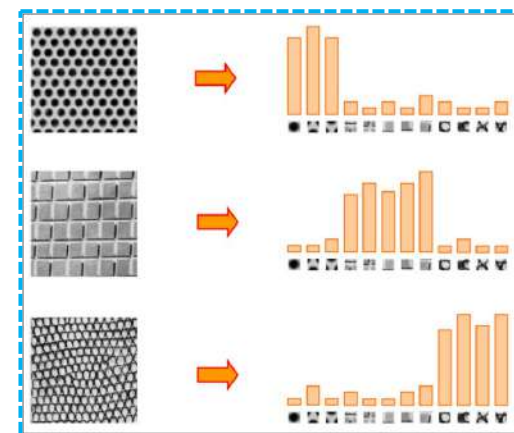
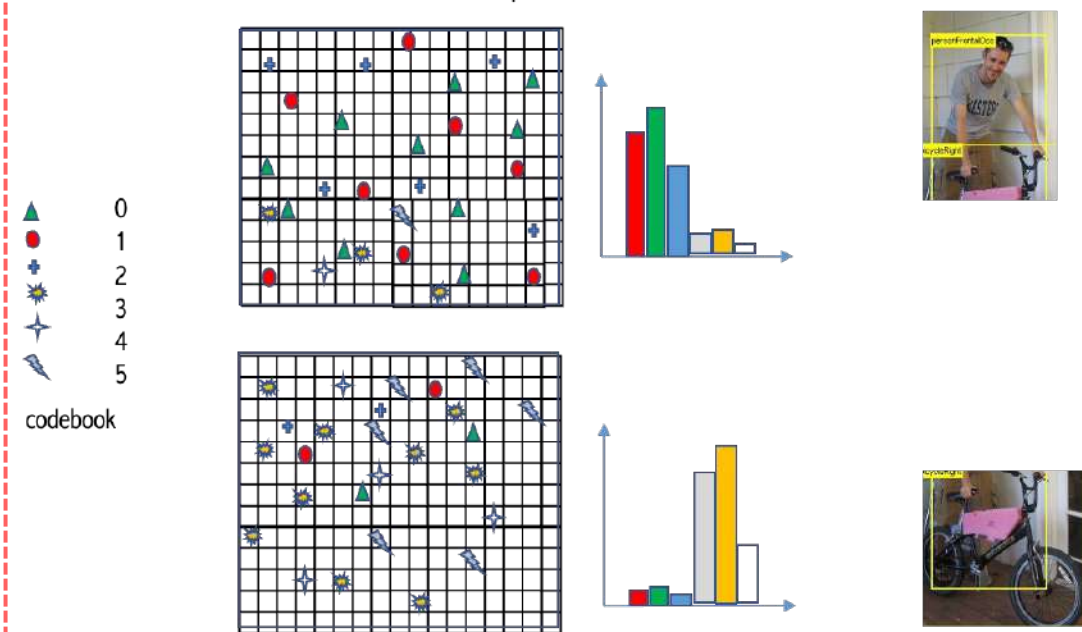
# Spatial Pyramid Matching

## Bag of Visual Words – K Means Clustering

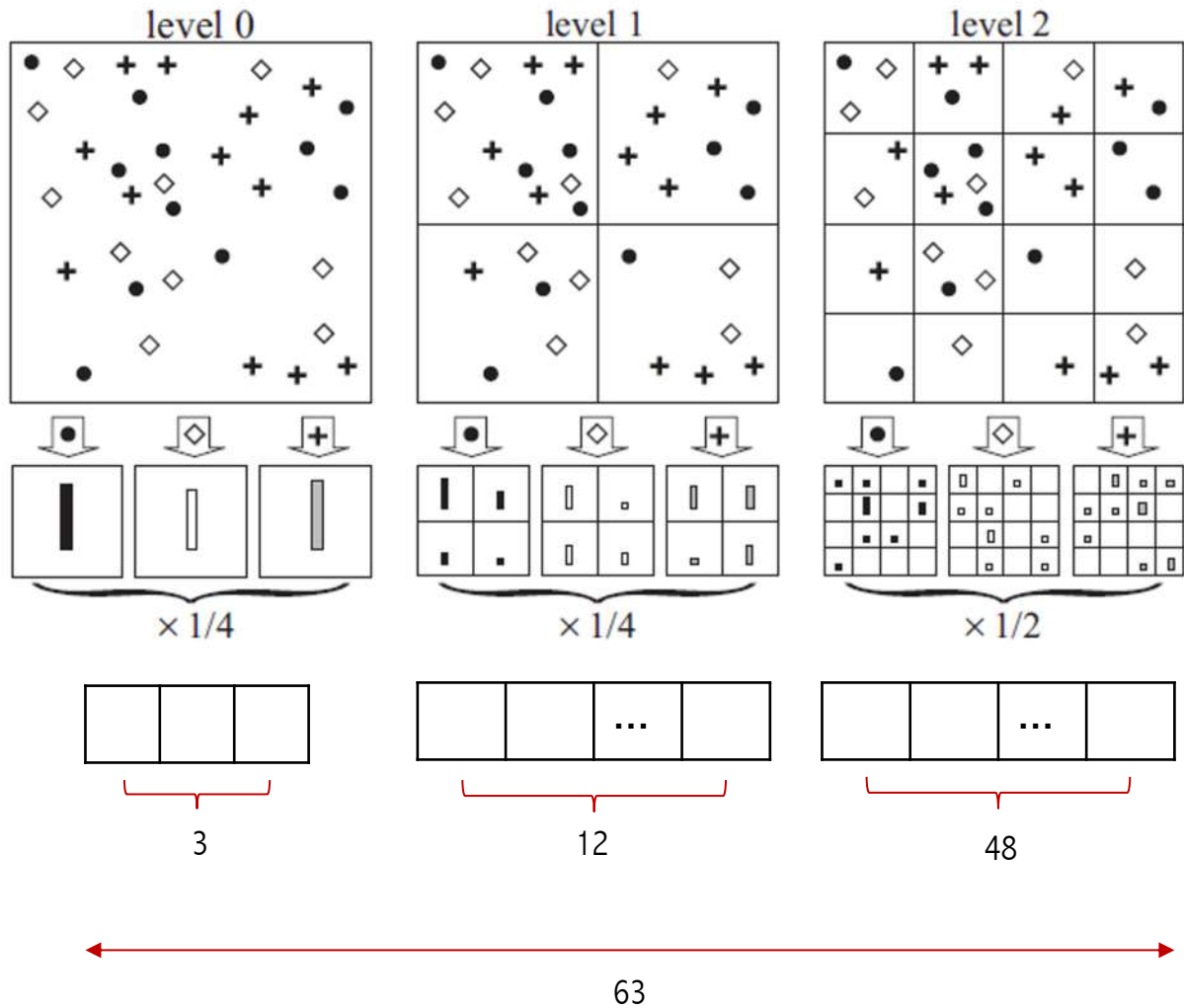
Generate HOG/SIFT Feature Descriptors



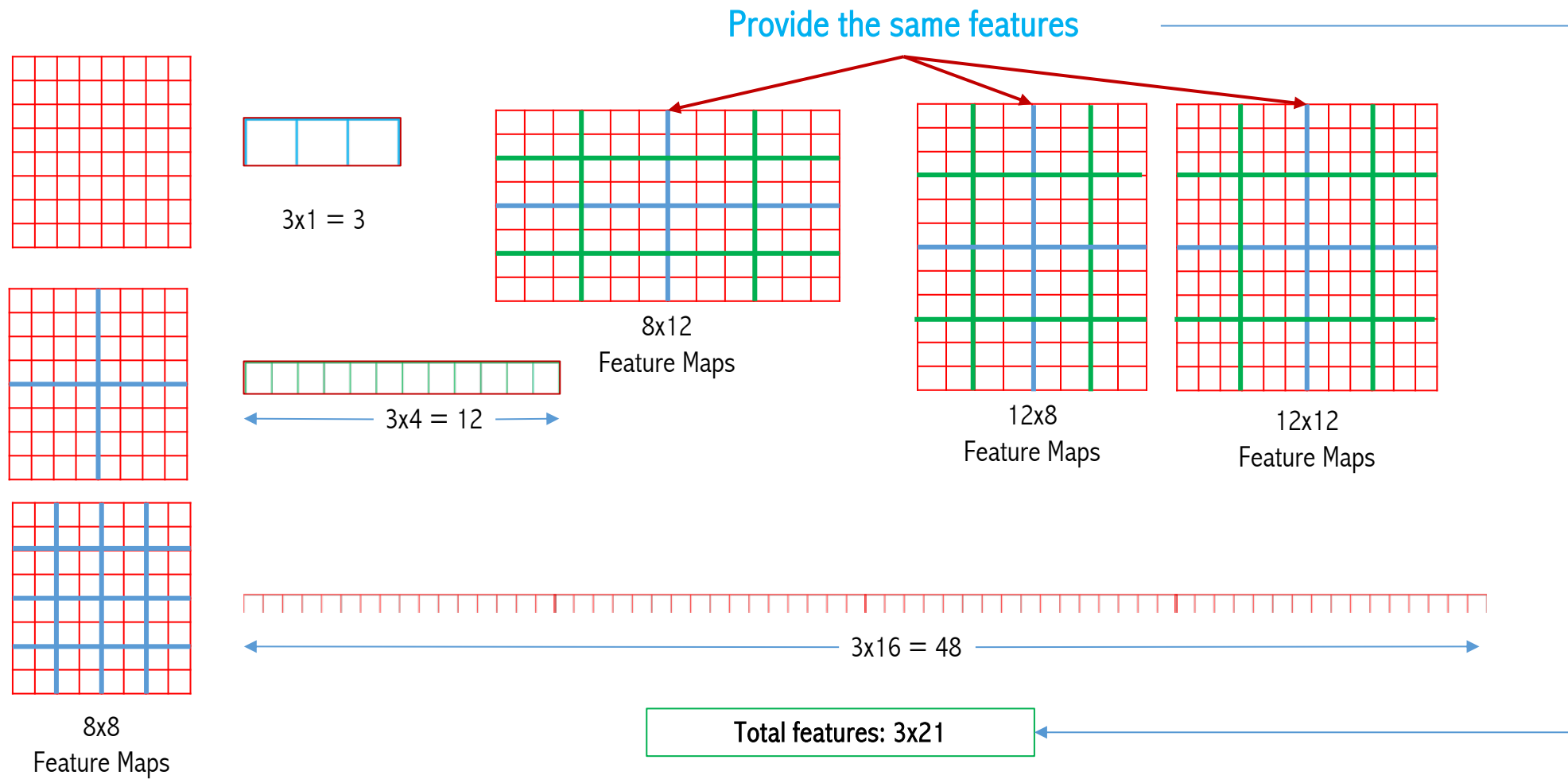
Generate HOG/SIFT Feature Descriptors



# Spatial Pyramid Matching

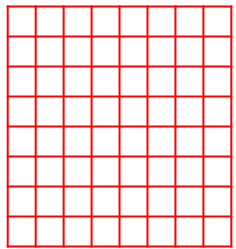


# Spatial Pyramid Matching

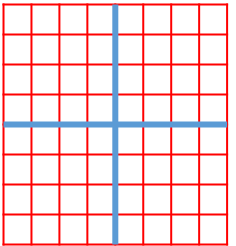


Any Size and Aspect Ratio

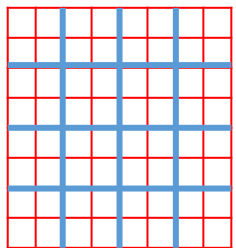
# Spatial Pyramid Pooling



1x1



4x1



16x1

8x8  
Feature Maps

21x1

- Identifying features
- K-means clustering
- Codebooks
- Histograms

*Just Max-Pool*

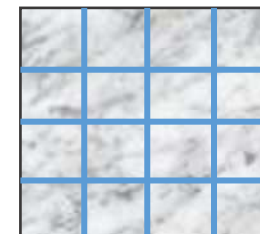
x256



1x256



4x256

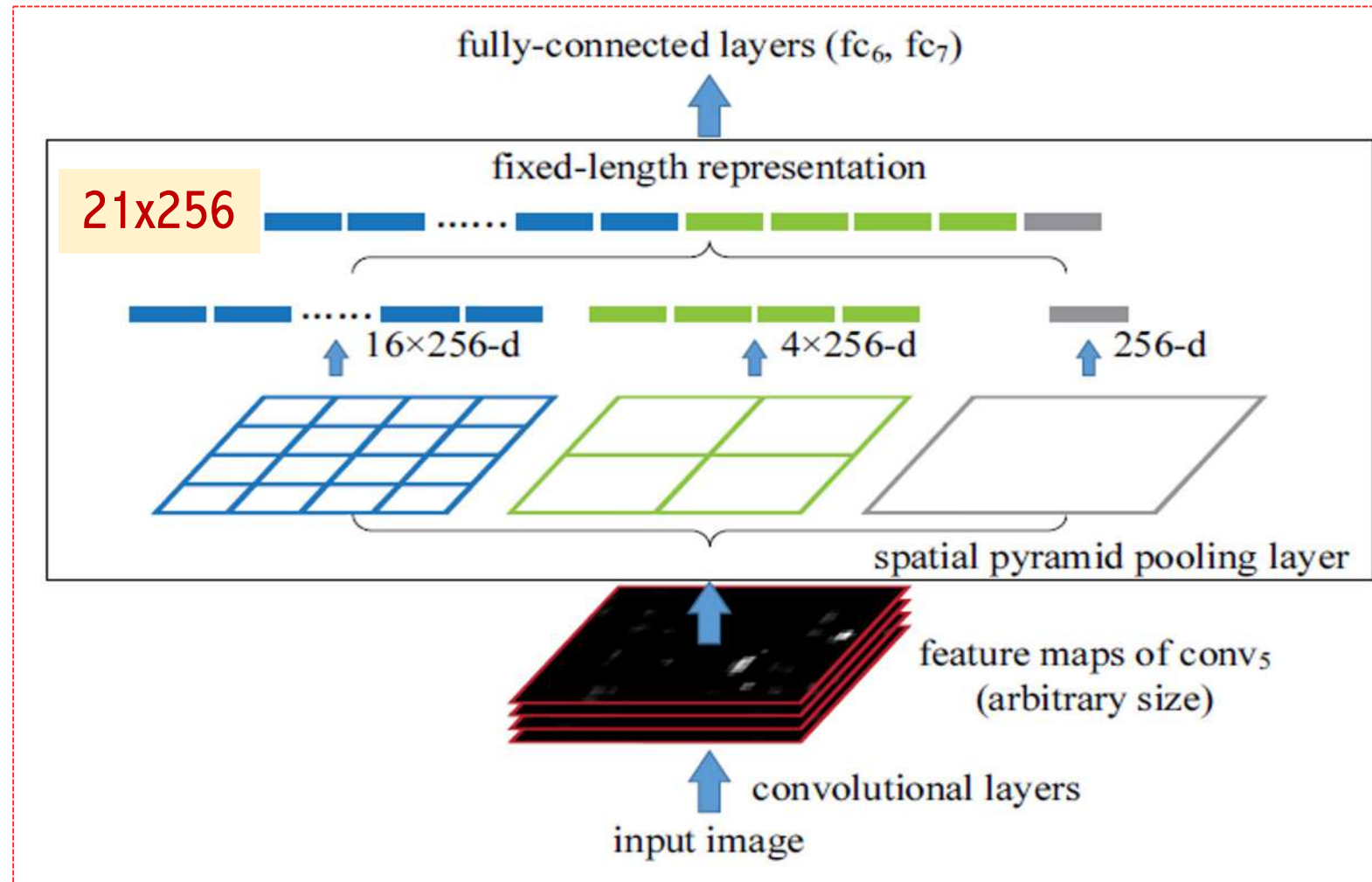


16x256

21x256

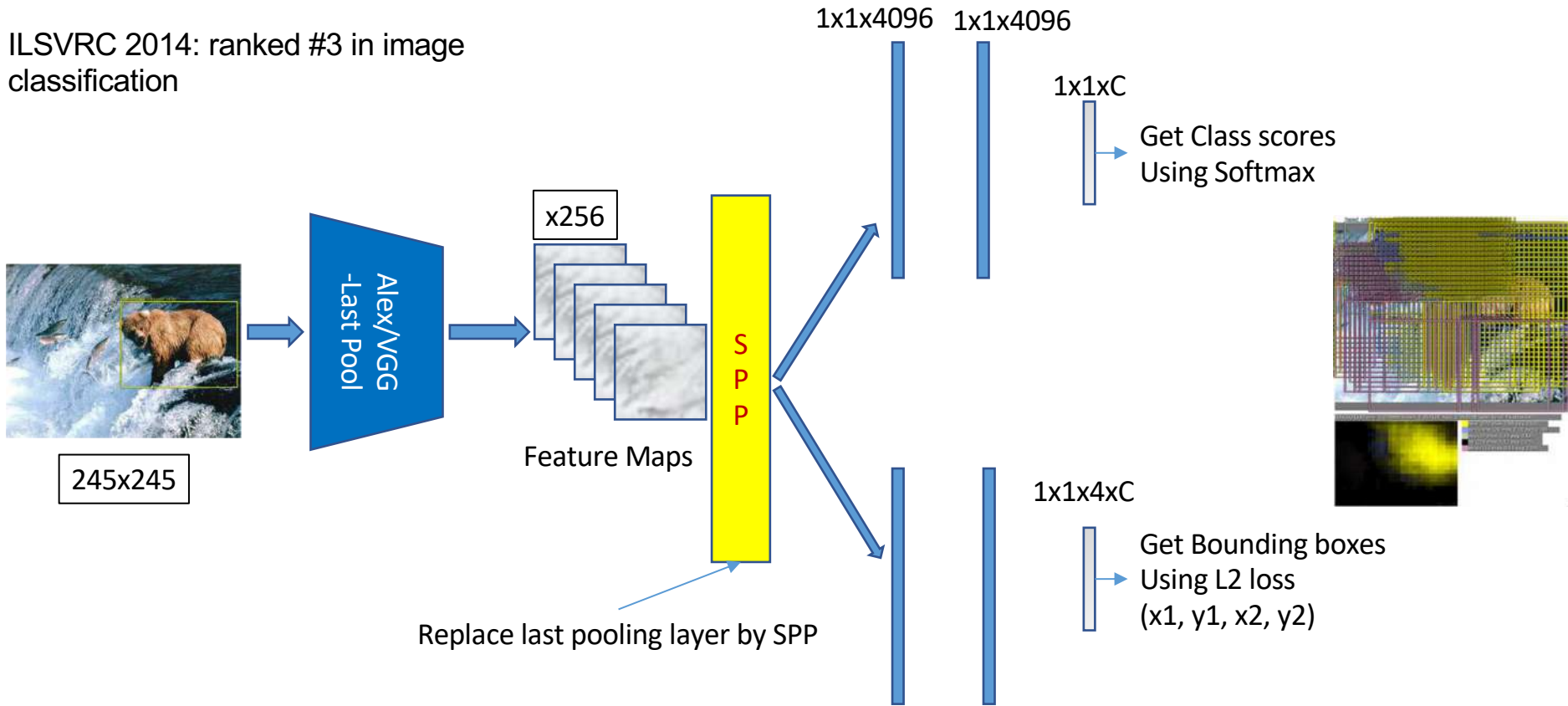


# Spatial Pyramid Pooling

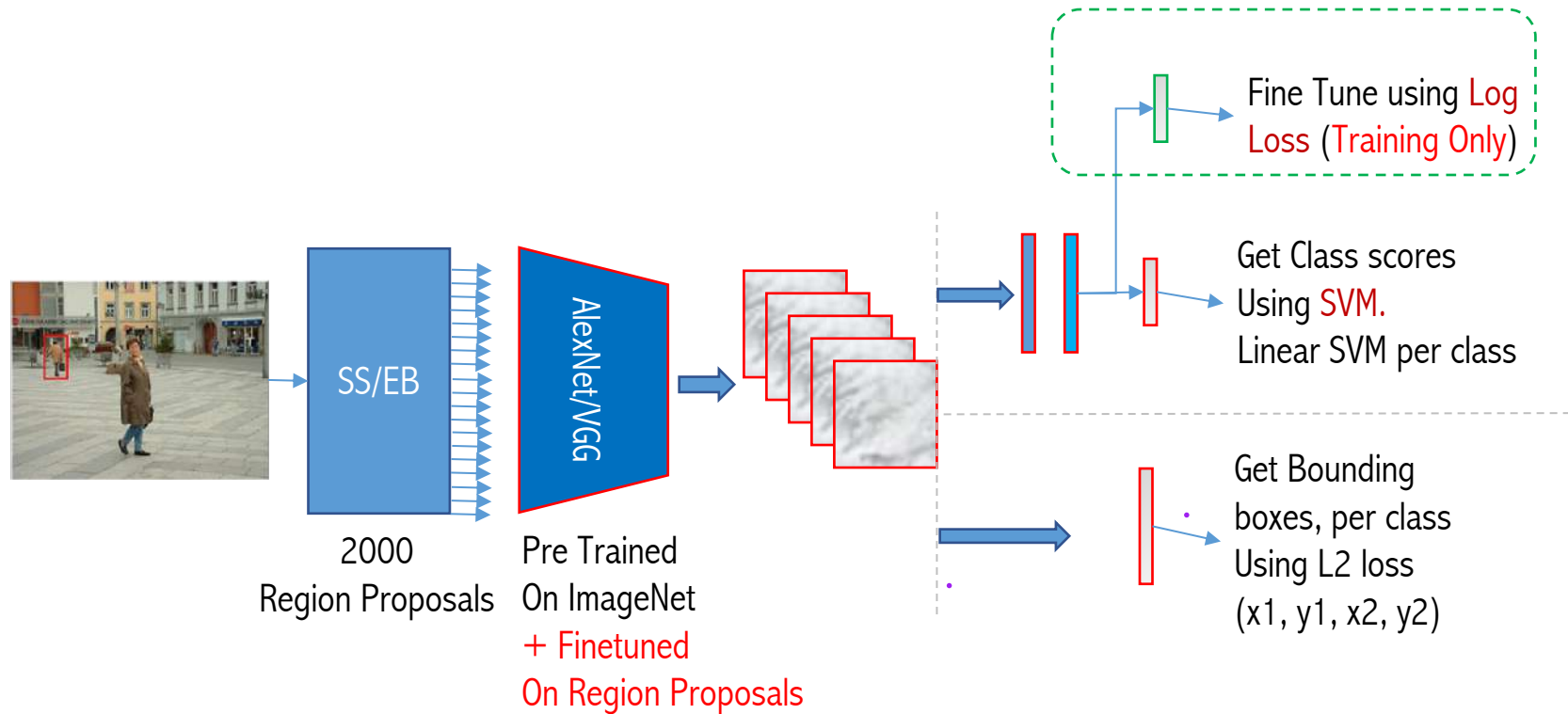


# SPPNet = SPP + Overfeat for Classification

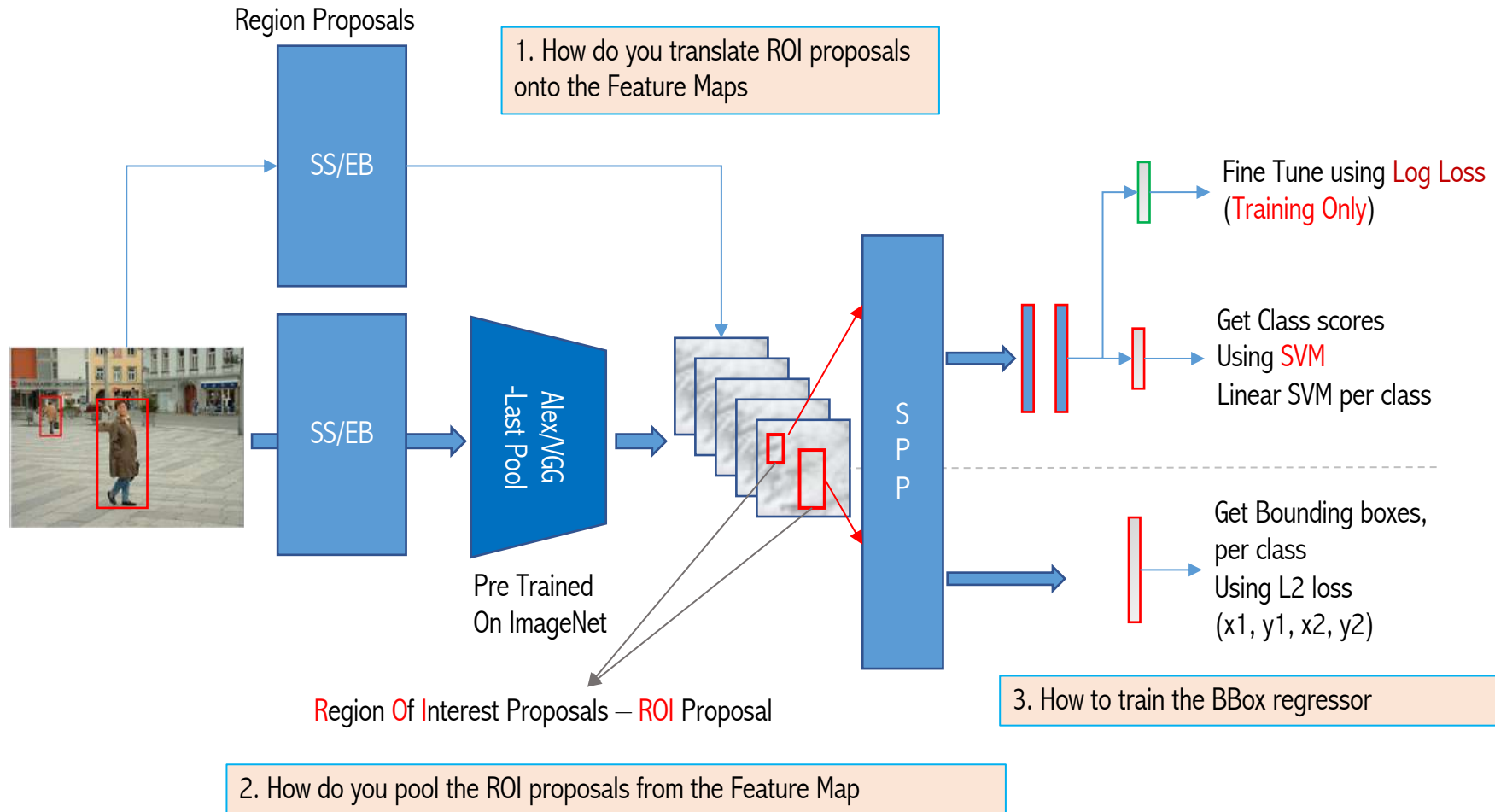
ILSVRC 2014: ranked #3 in image classification



# RCNN – Two stage-based methods



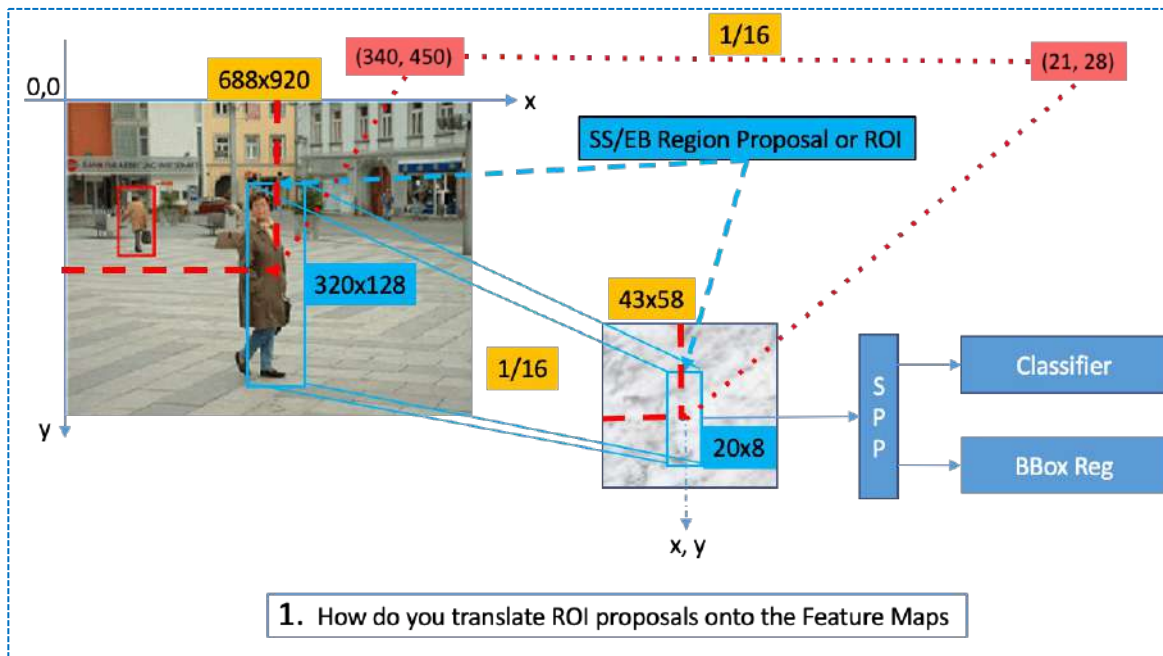
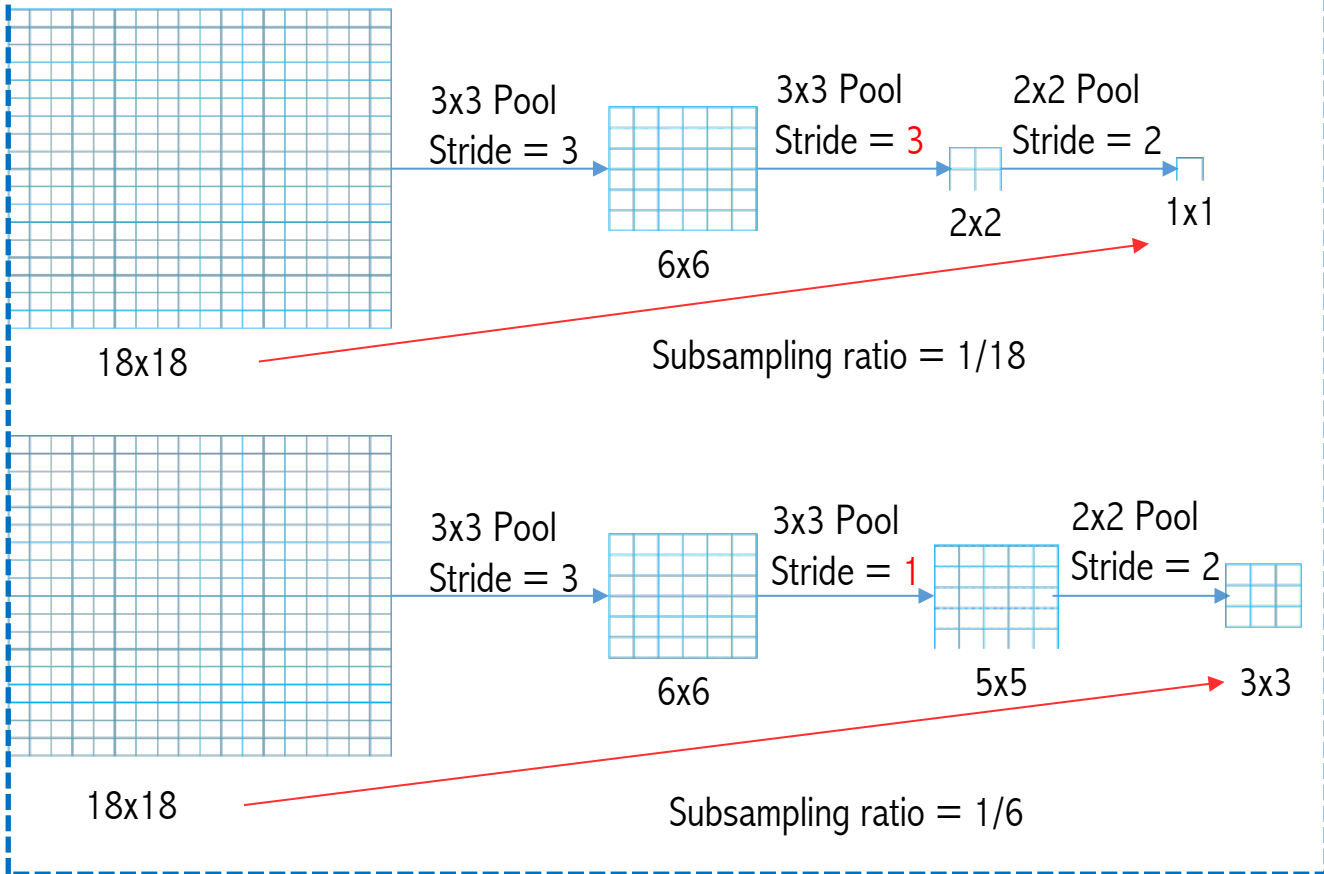
# SPP – Two stage-based methods



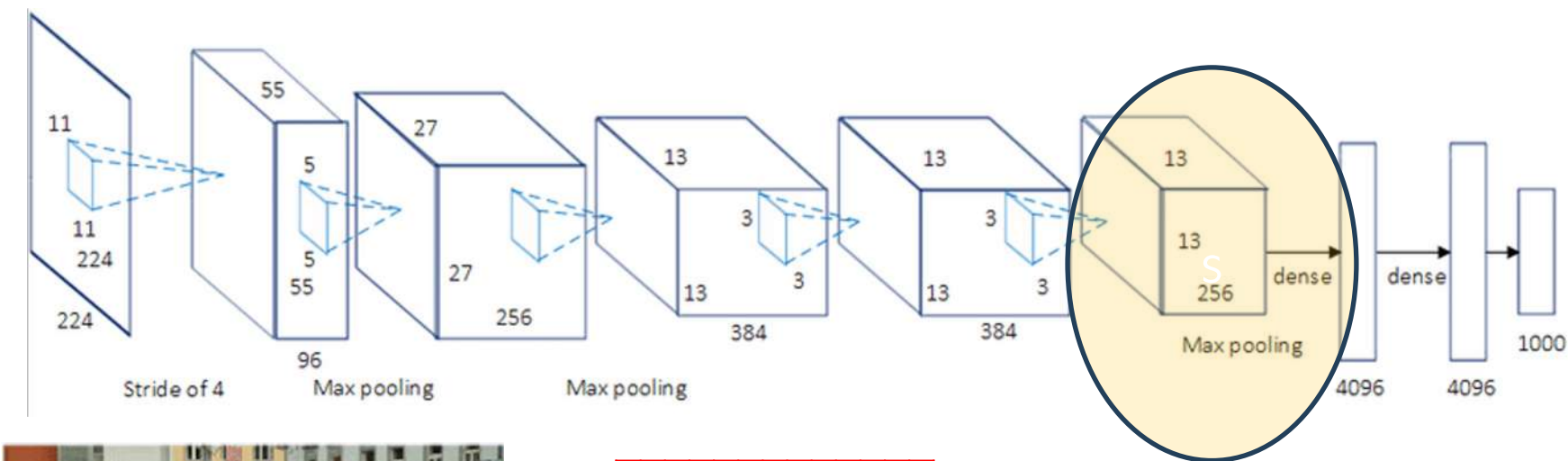


# Subsampling Ratio

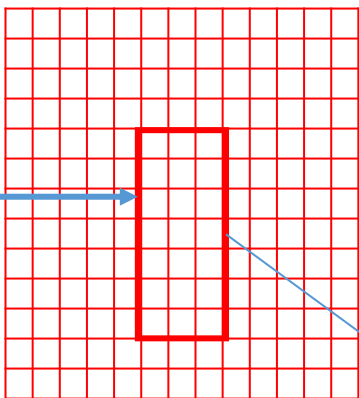
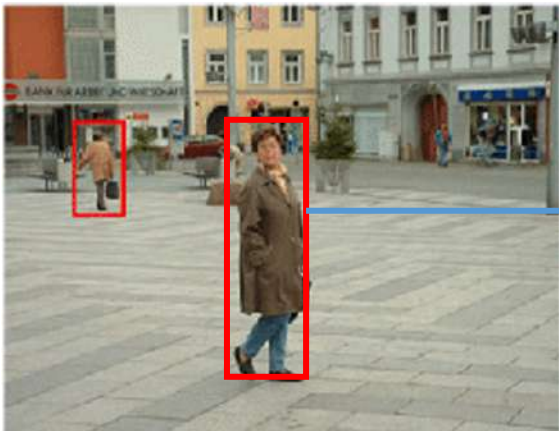
1. How do you translate ROI proposals onto the Feature Maps



# AlexNet Subsampling Ratio



2. How do you pool the ROI proposals from the Feature Map



13x13  
Feature Maps

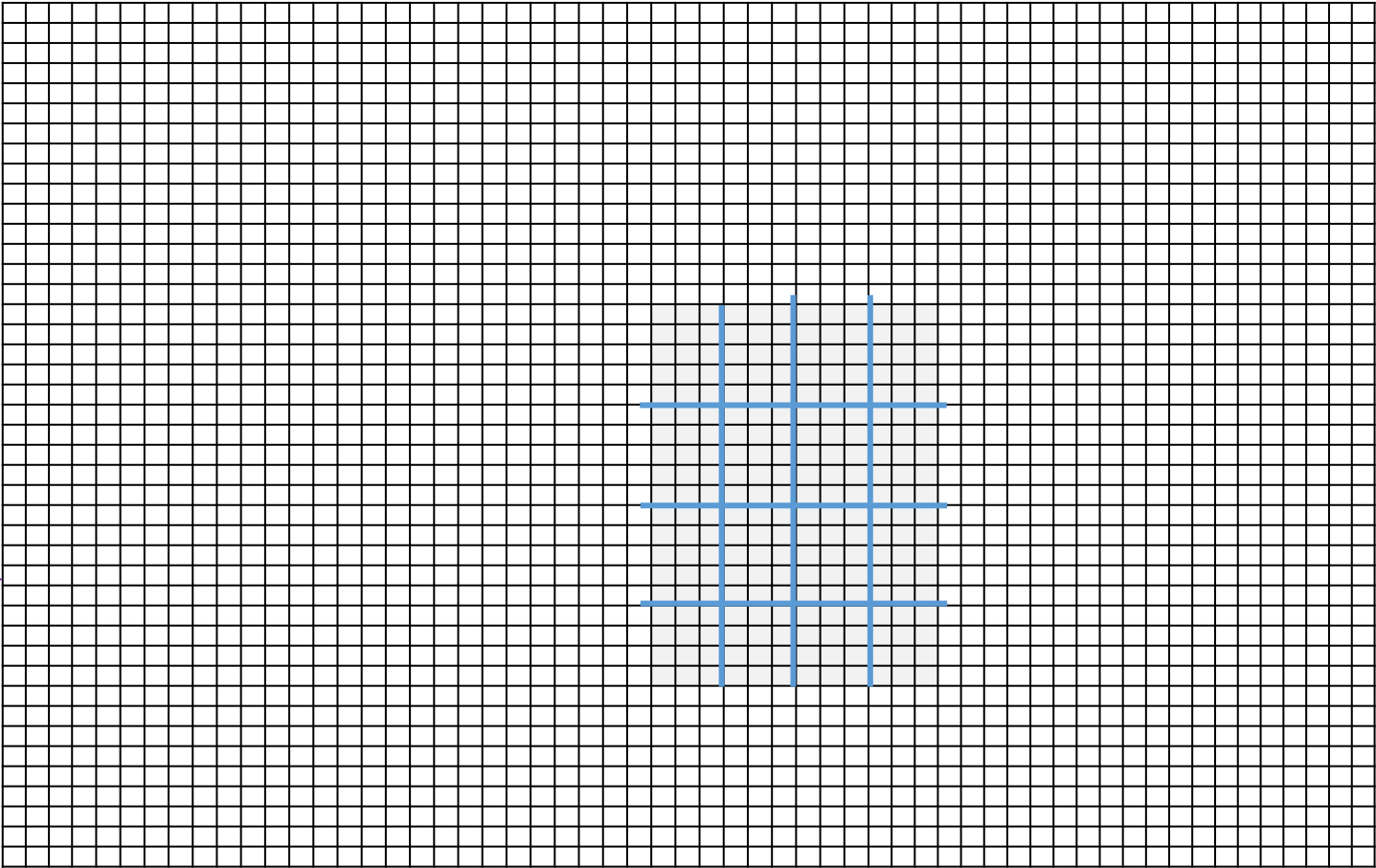
Spatial Pyramid Pooling  
Level 1, Level 2, Level 3

6x6 grid. It is impossible

How to solve

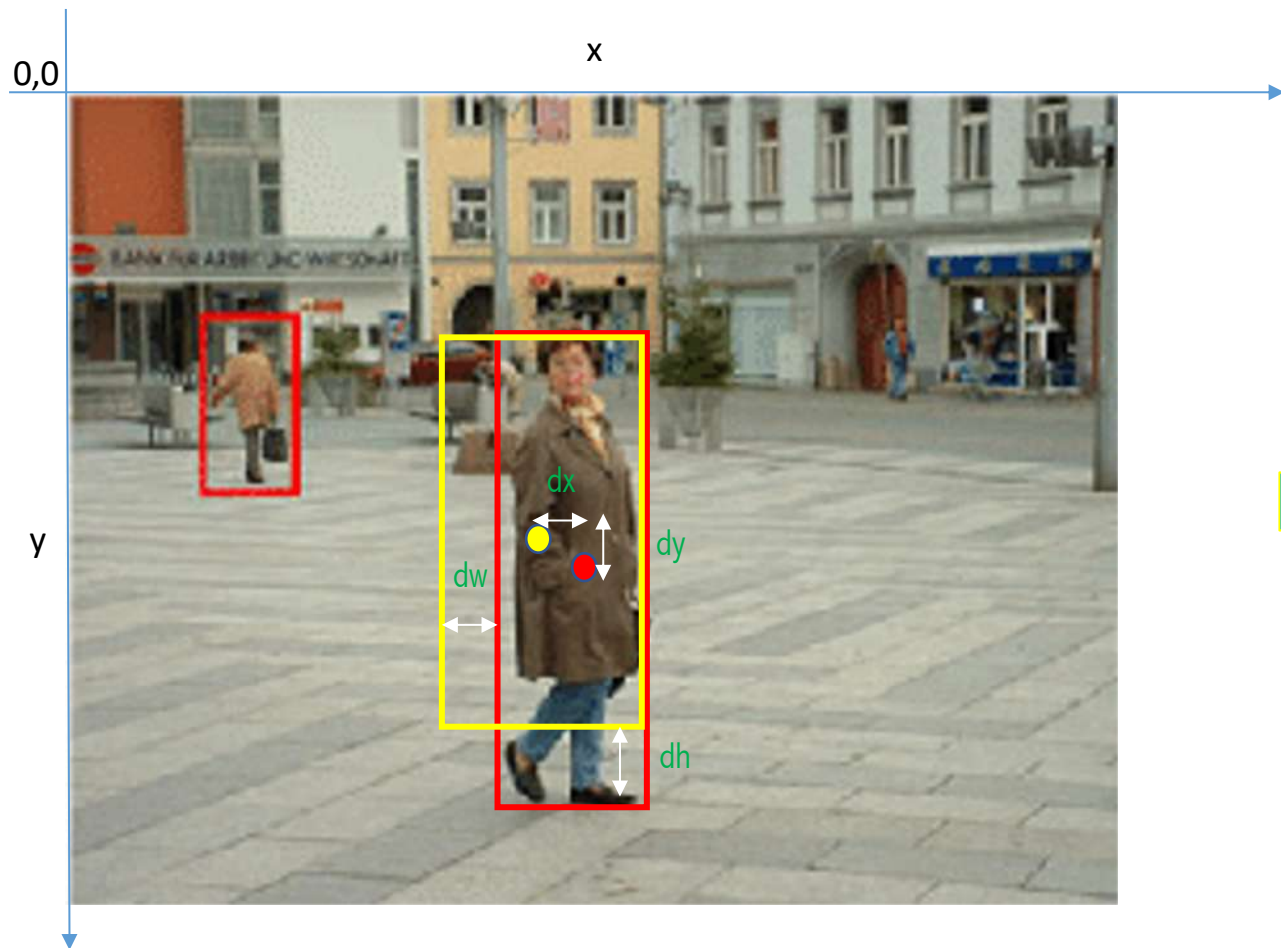
# SPP on Region Proposals

Aspect ratio 3:4

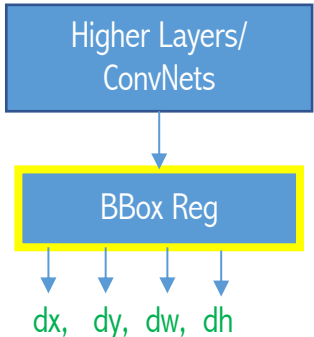


Three level in Practice - {6x6, 3x3, 2x2, 1x1}

# BBox Regression Training



ROI Centre/W/H  
 $x, y, w, h$



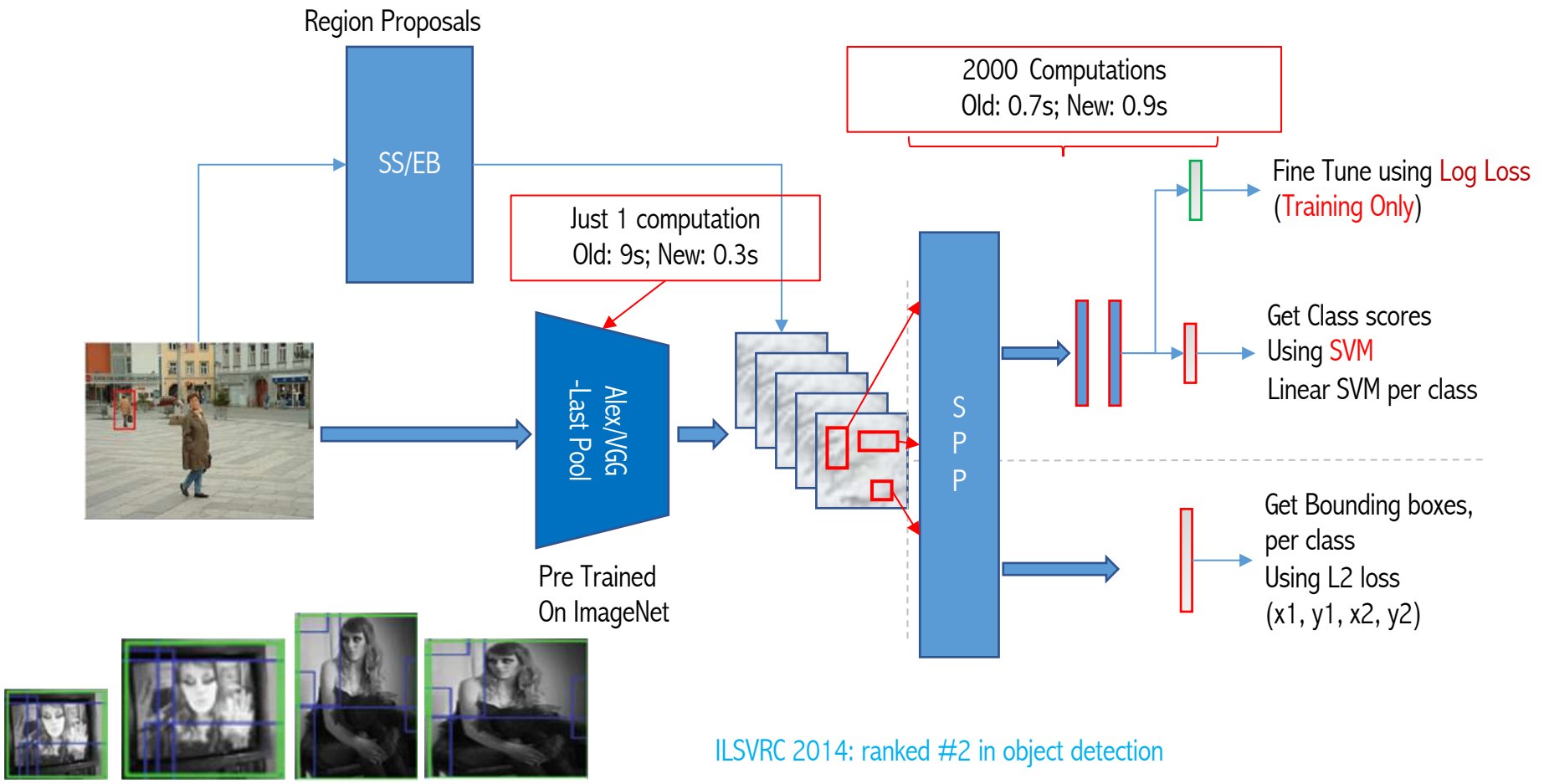
Ground Truth  
 $x_g, y_g, w_g, h_g$

$$(x + dx - x_g)^2 = 0$$

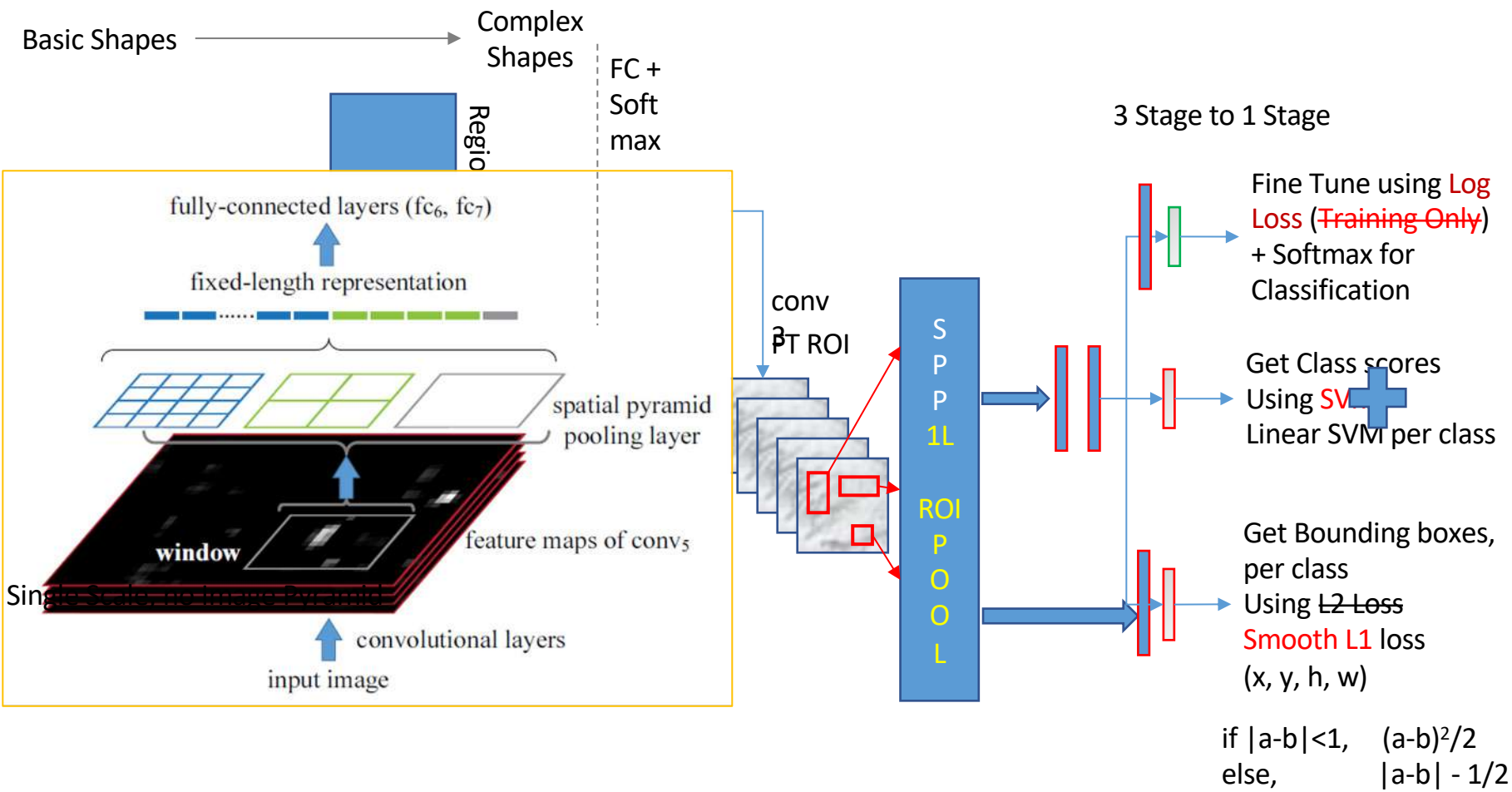
3. How to train the BBox regressor



# SPP – 2 Stage Network - Inference



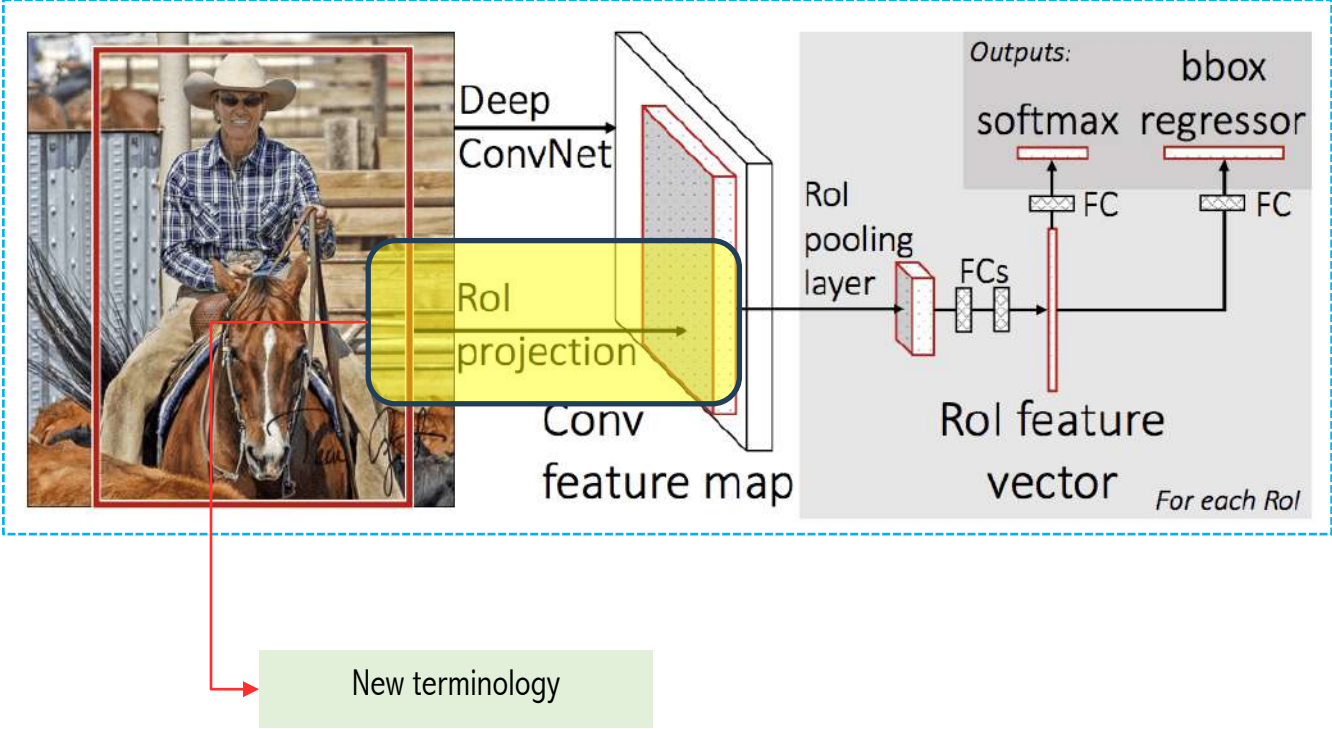
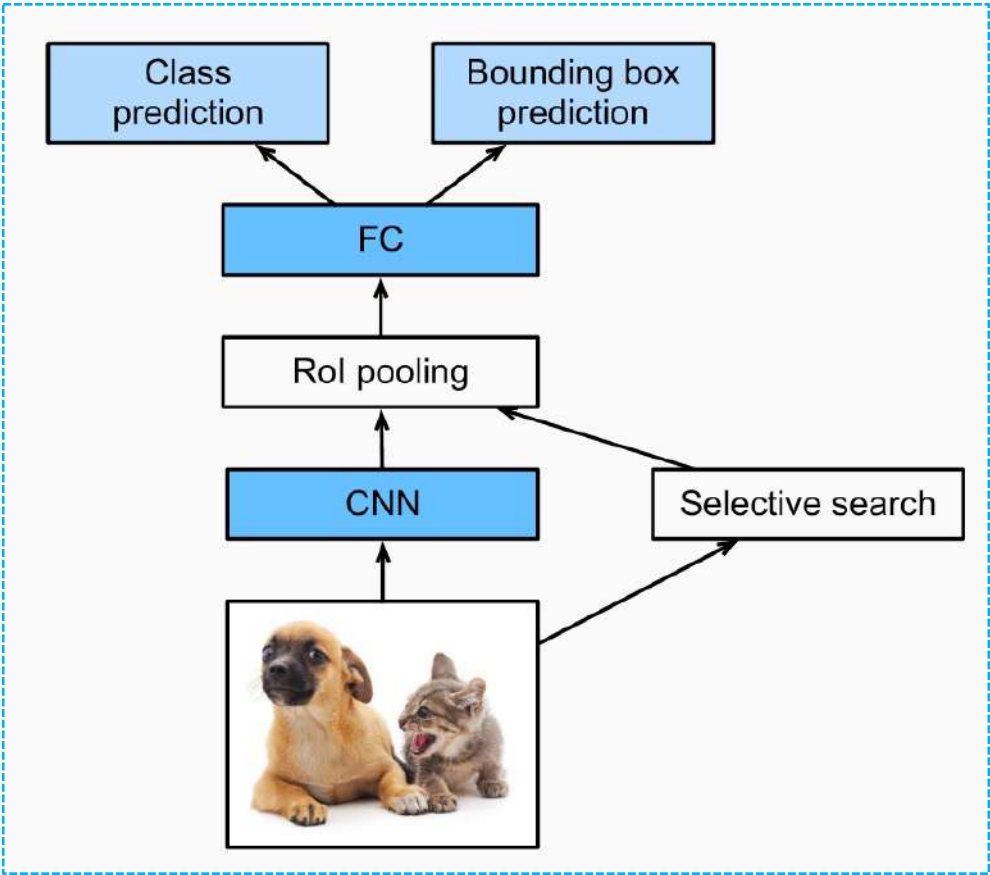
# Speed-up: SPP vs RCNN=>Fast RCNN



	SPP (1-sc) (ZF-5)	SPP (5-sc) (ZF-5)	R-CNN (Alex-5)
pool <sub>5</sub>	43.0	<u>44.9</u>	44.2
fc <sub>6</sub>	42.5	44.8	<u>46.2</u>
ftfc <sub>6</sub>	52.3	<u>53.7</u>	53.1
ftfc <sub>7</sub>	54.5	<u>55.2</u>	54.2
ftfc <sub>7</sub> bb	58.0	<b>59.2</b>	58.5
conv time (GPU)	0.053s	0.293s	8.96s
fc time (GPU)	0.089s	0.089s	0.07s
total time (GPU)	0.142s	0.382s	9.03s
speedup (vs. RCNN)	64×	24×	-

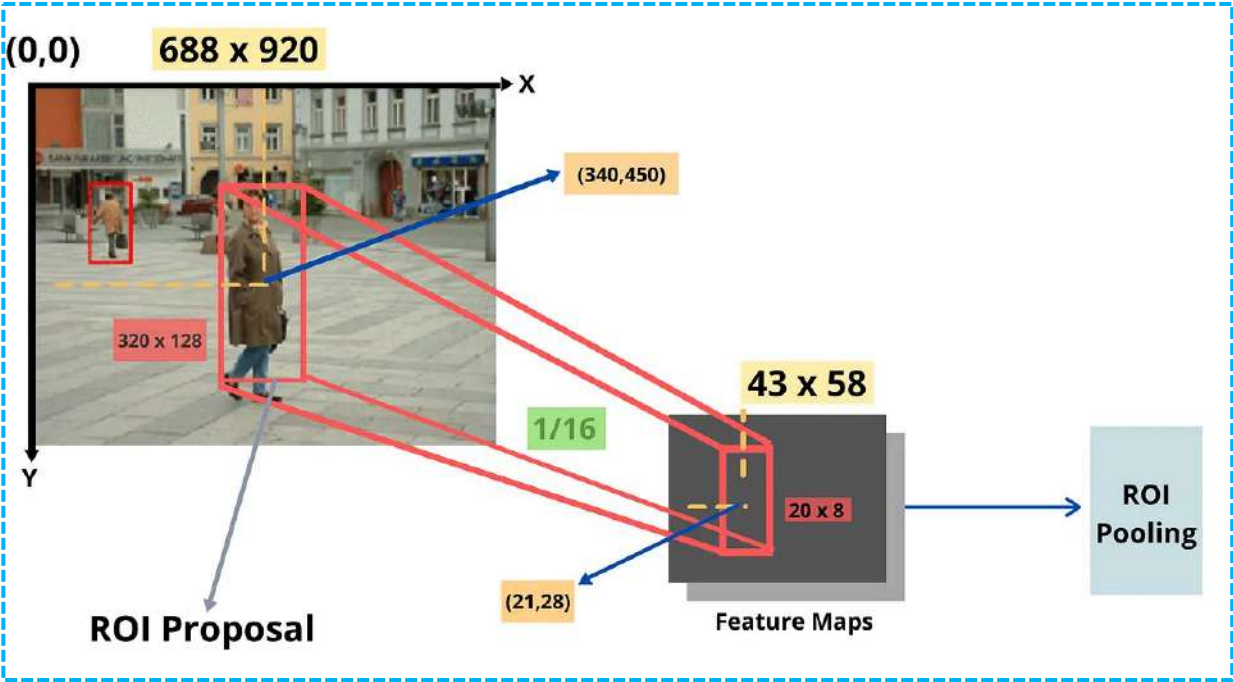
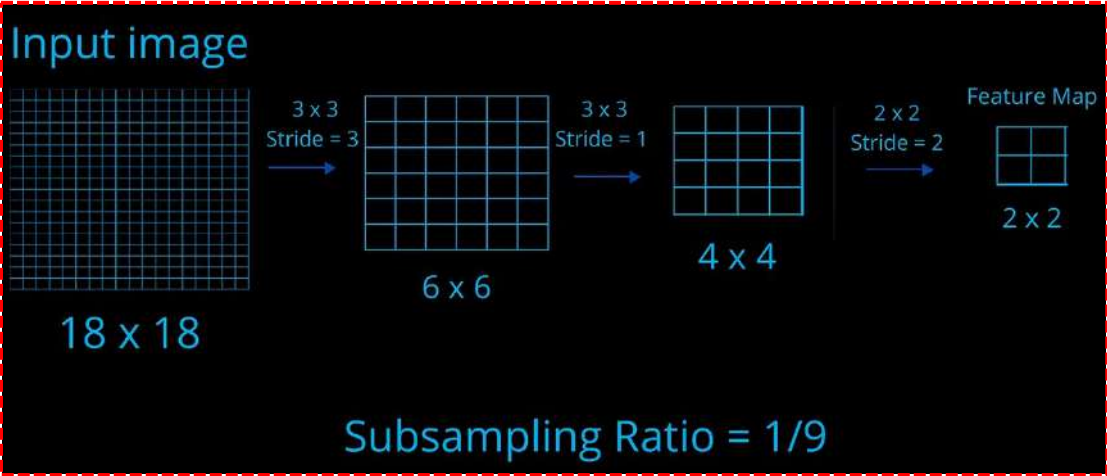
Table 9: Detection results (mAP) on Pascal VOC 2007. "ft" and "bb" denote fine-tuning and bounding box regression.

# Fast R-CNN



Fast R-CNN, which was developed a year later after R-CNN, solves these issues very efficiently and is about 146 times faster than the R-CNN during the test time.

# Sub-Sampling Ratio & Roi Projection



The idea of ROI projection is that we get the coordinates of the bounding box from the ROI proposal and we need to project them onto the feature maps by projecting the ROI proposal with respect to the subsampling ratio.

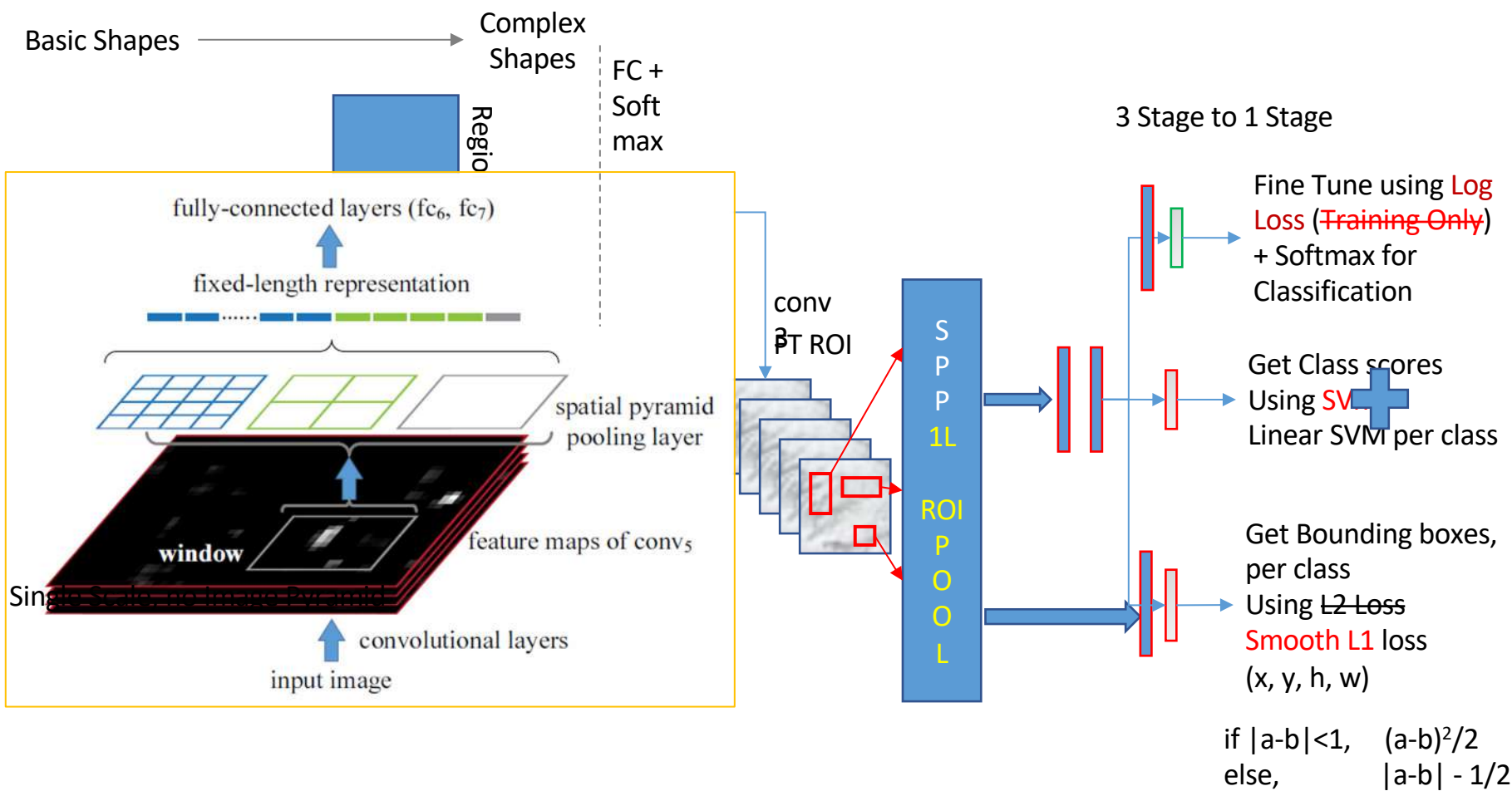


# Outline

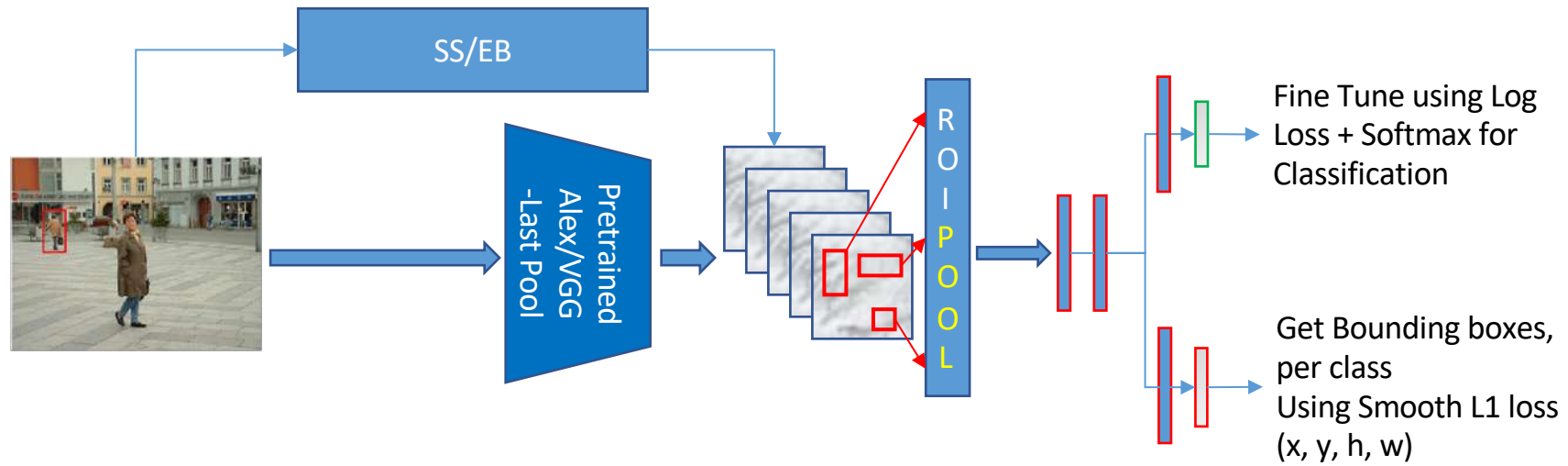
- **CNN Limitations**
- **Region Based Convolutional Neural Networks**
- **Spatial Pyramid Pooling**
- **Fast R-CNN**
- **Faster RCNN**
- **YOLOv1-v2**



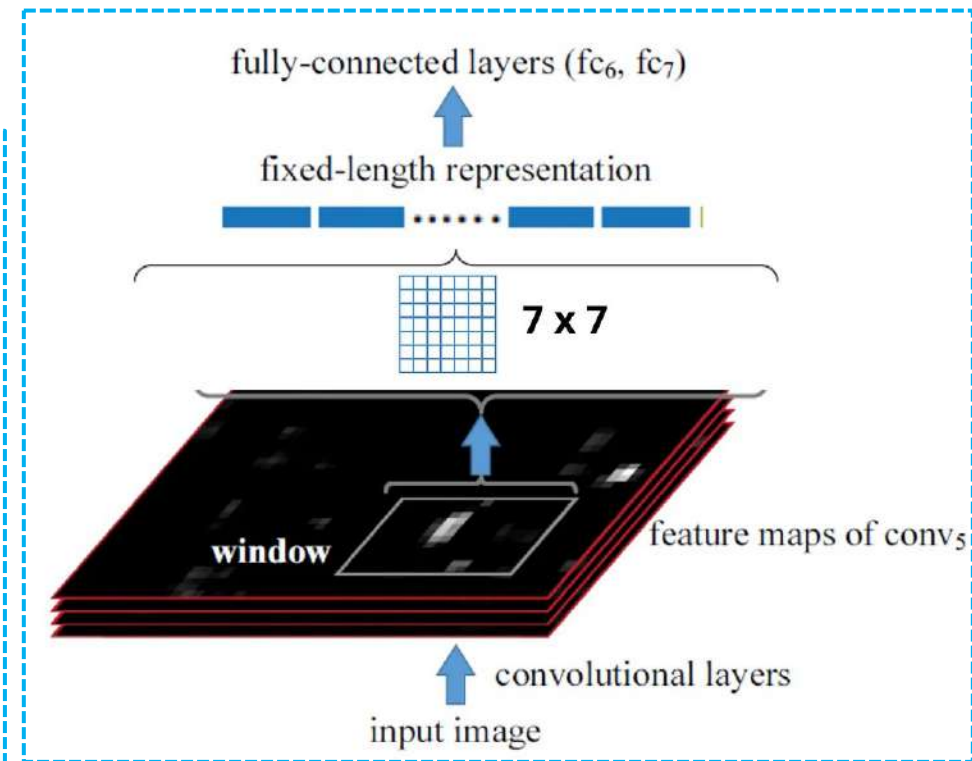
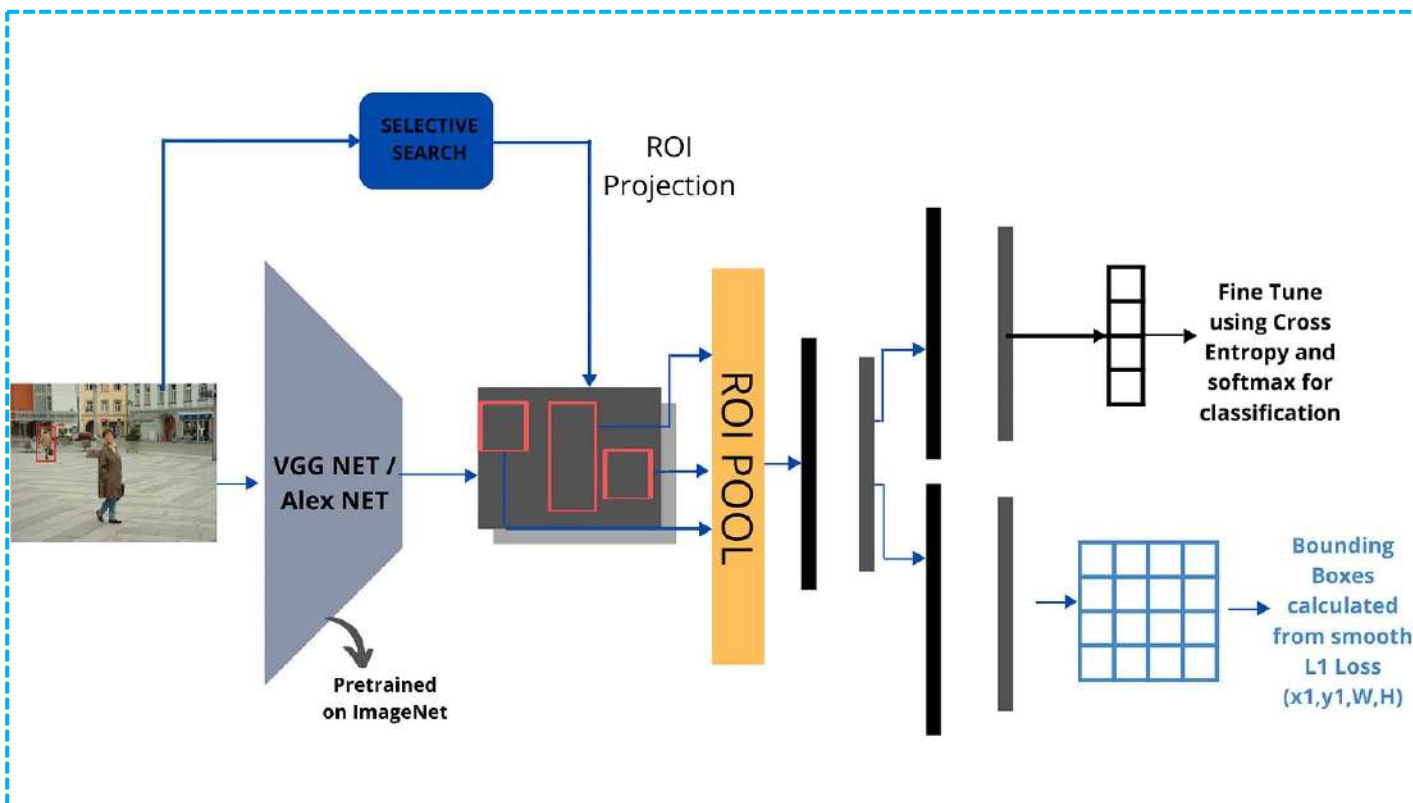
# RCNN -> SPPNet -> Fast RCNN



# RCNN -> SPPNet -> Fast RCNN



# Fast R-CNN



So why don't we reuse those same CNN results for region proposals instead of running a separate selective search algorithm?

	RCNN	Fast RCNN
Training Time	84 hours	9.5 hours
(Speedup)	1x	8.8x
Test time per image	47 seconds	0.32 seconds
(Speedup)	1x	146x
Test time per image with Selective Search	50 seconds	2 seconds
(Speedup)	1x	25x

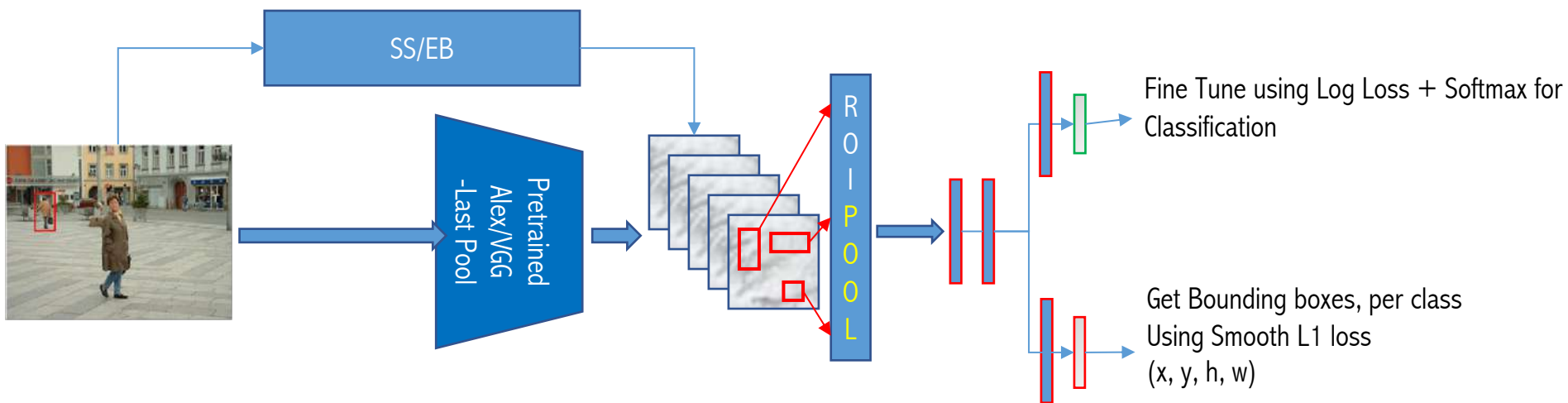
# Outline

- **CNN Limitations**
- **Region Based Convolutional Neural Networks**
- **Spatial Pyramid Pooling**
- **Fast R-CNN**
- **Faster RCNN**
- **YOLOv1-v2**





# Criteria for replacing SS



< 2000 Region Proposals

As fast as SS or better

As Accurate as SS or better

Should be able to propose Overlapping ROIs with different Aspect Ratios and Scale

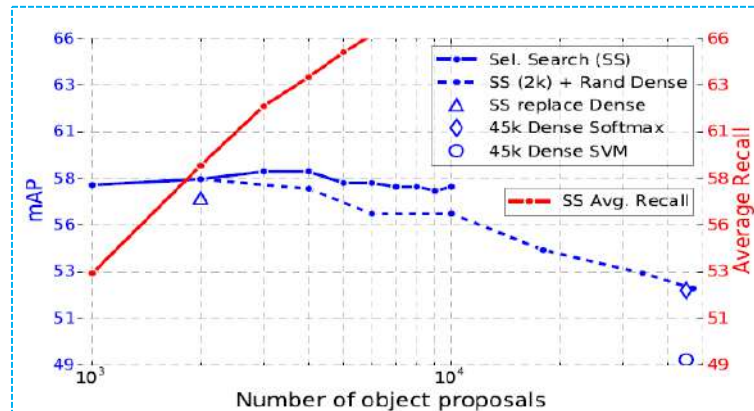
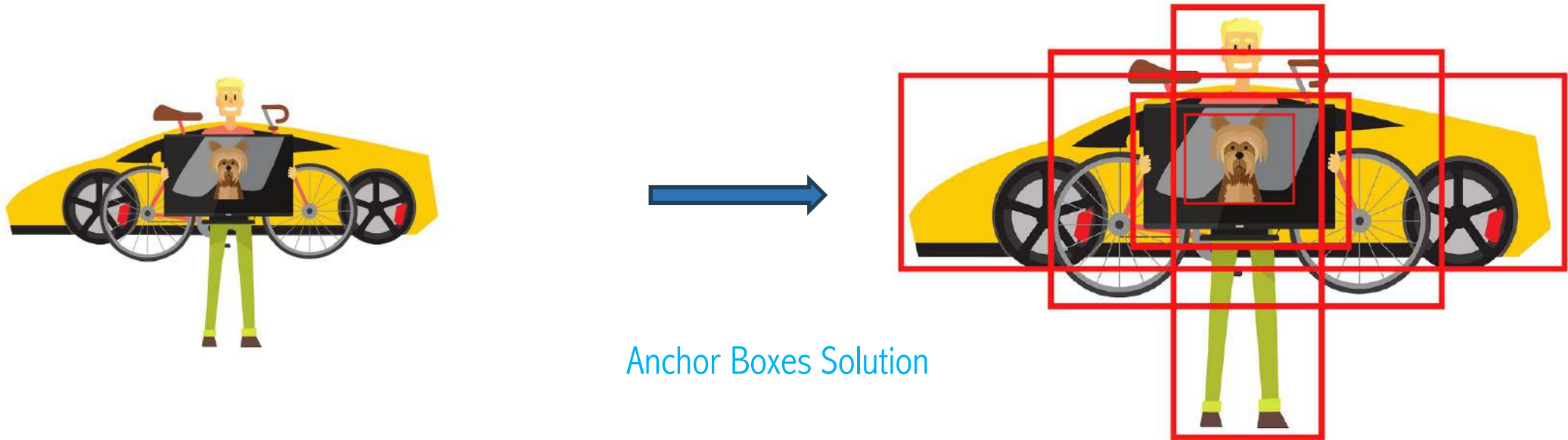


Figure 3. VOC07 test mAP and AR for various proposal schemes.

Role of Region Proposals

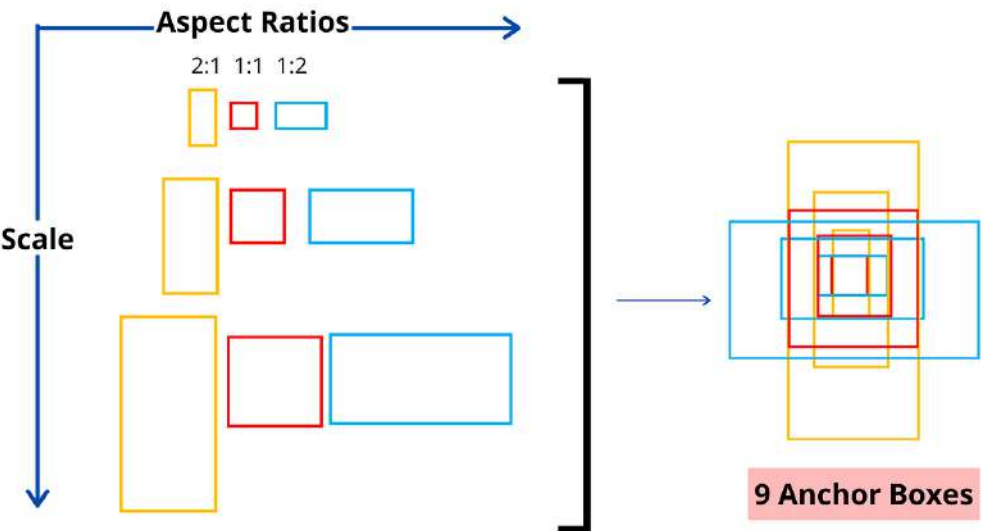
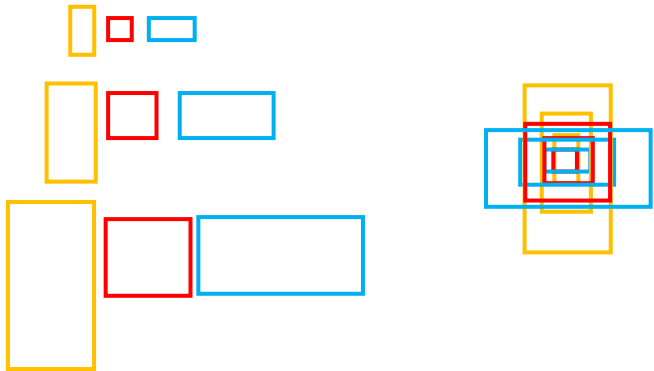
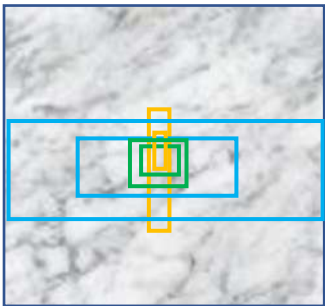
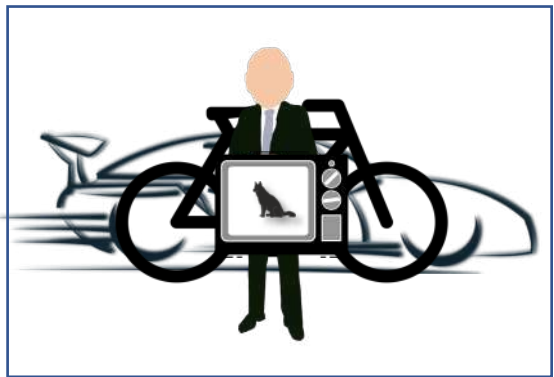
# Overlapping ROIs



Anchor Boxes Solution

From the image, we see a lot of objects overlapping each other. We see a car, a bicycle, a person holding a television, and a dog inside this television. The selective search could solve this problem but we end up with a huge number of ROIs. We need to think of an idea that efficiently solves this.

# Anchor Boxes

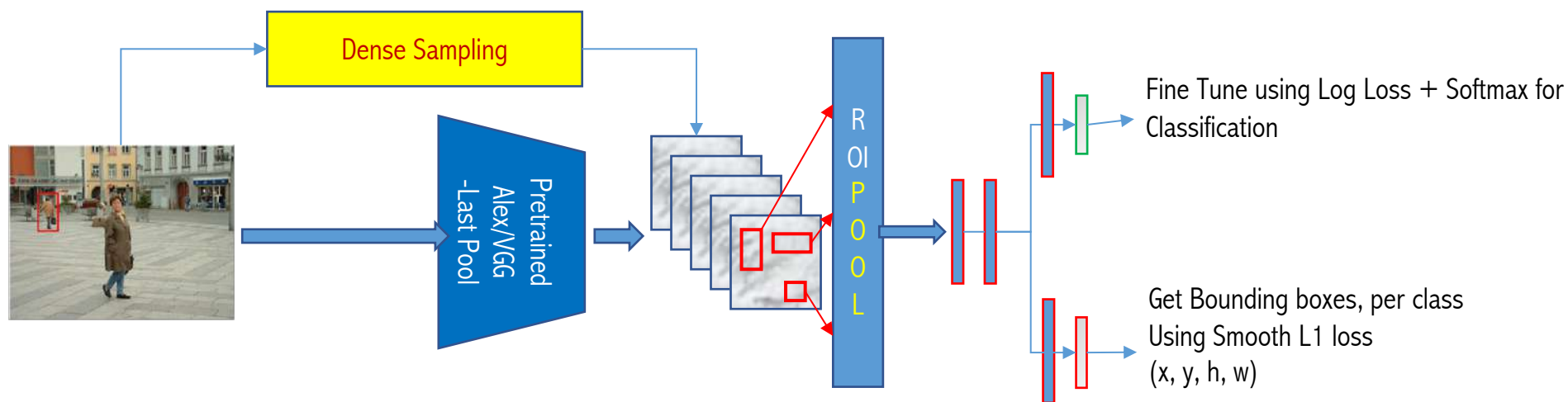


Any object in the image can be detected using boxes of 3 different scales and 3 different aspect ratios.

This could be a technique that can be used to solve our purpose of replacing the region proposal.

# Solutions for Replace Region Proposals

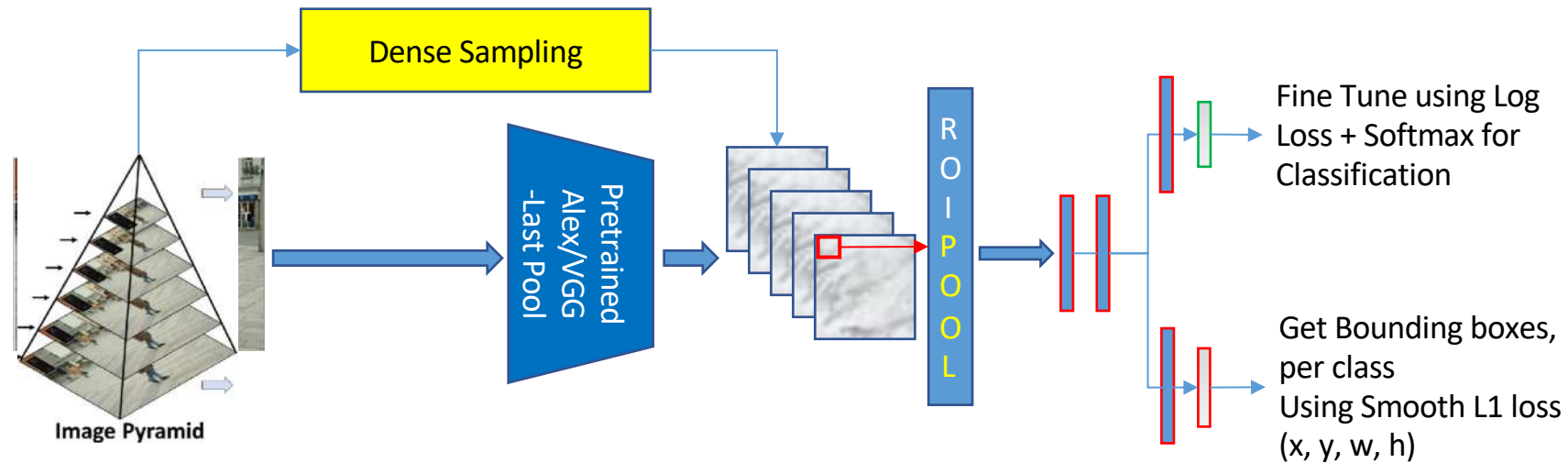
Removing Selective Search and applying a sliding window on top of the Feature Maps. But with this, we end detecting mostly objects of a single scale.



Fast RCNN + Sliding Window

# Solutions for Replace Region Proposals

To take care of multiple scales, we have to use Image Pyramids at the input. But using images of 5 different scales (by which almost every object can be detected) makes the network 4 times slower.

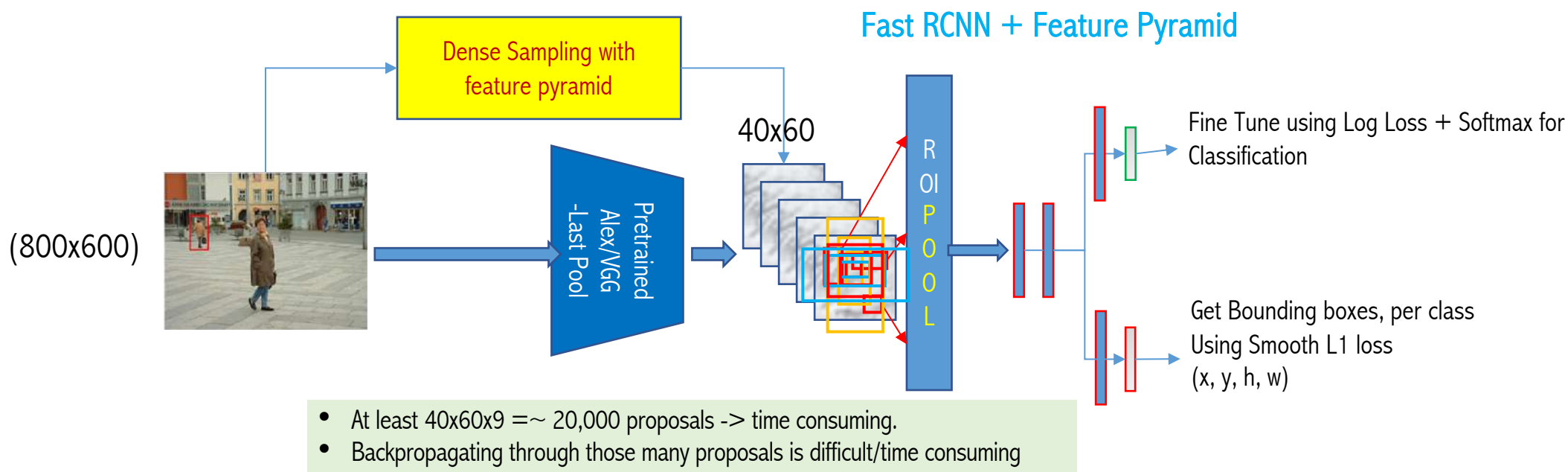
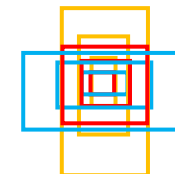


Fast RCNN + Sliding Window + Image Pyramid

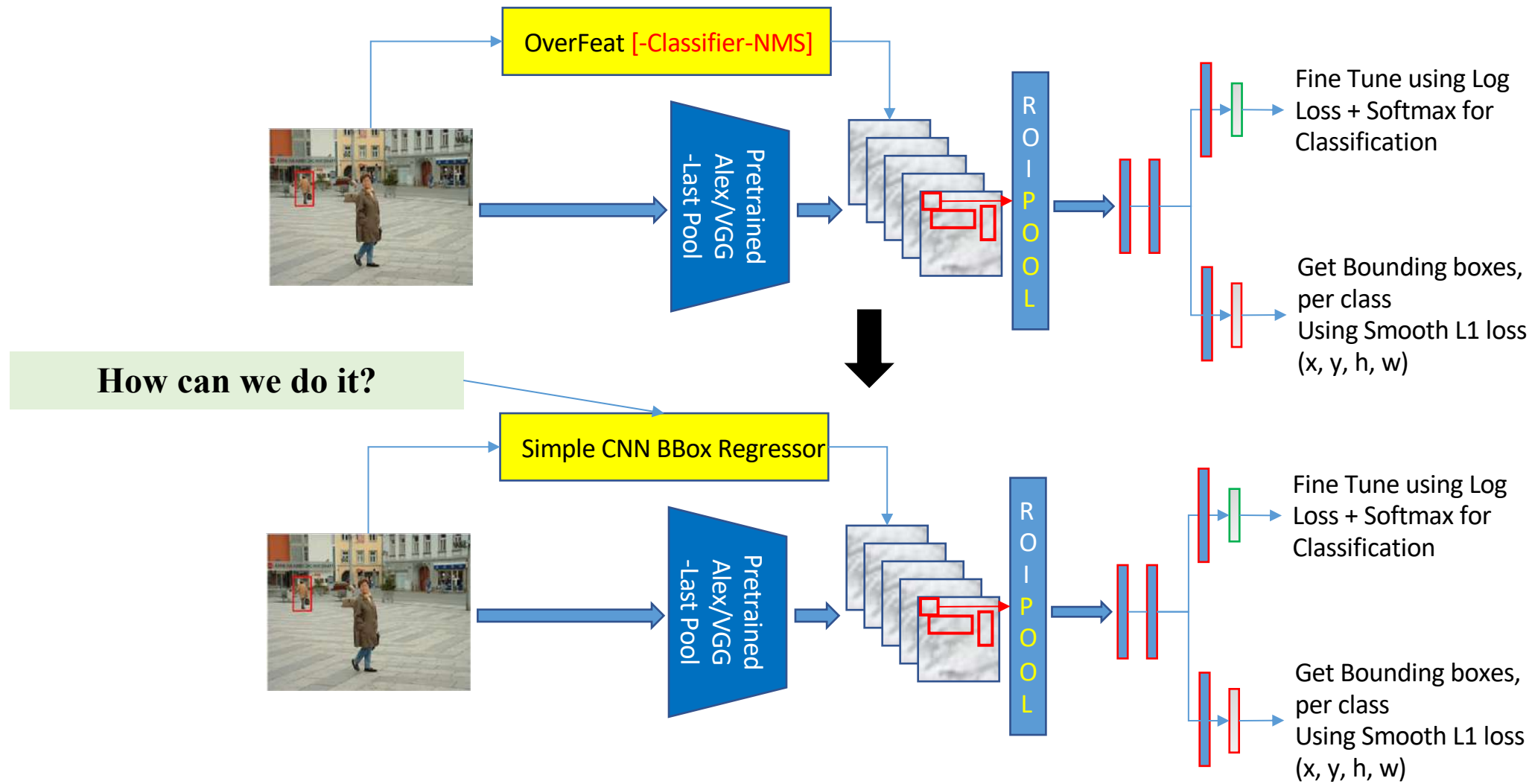


# Solutions for Replace Region Proposals

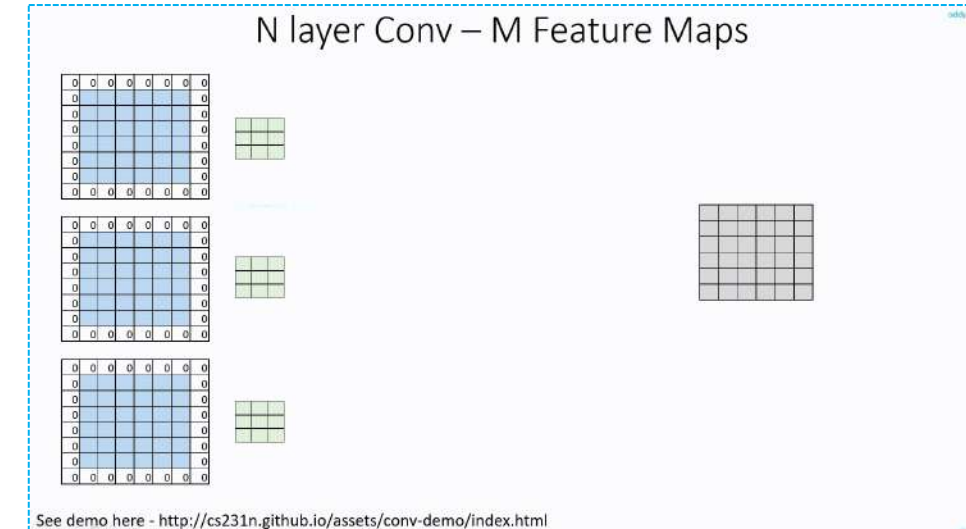
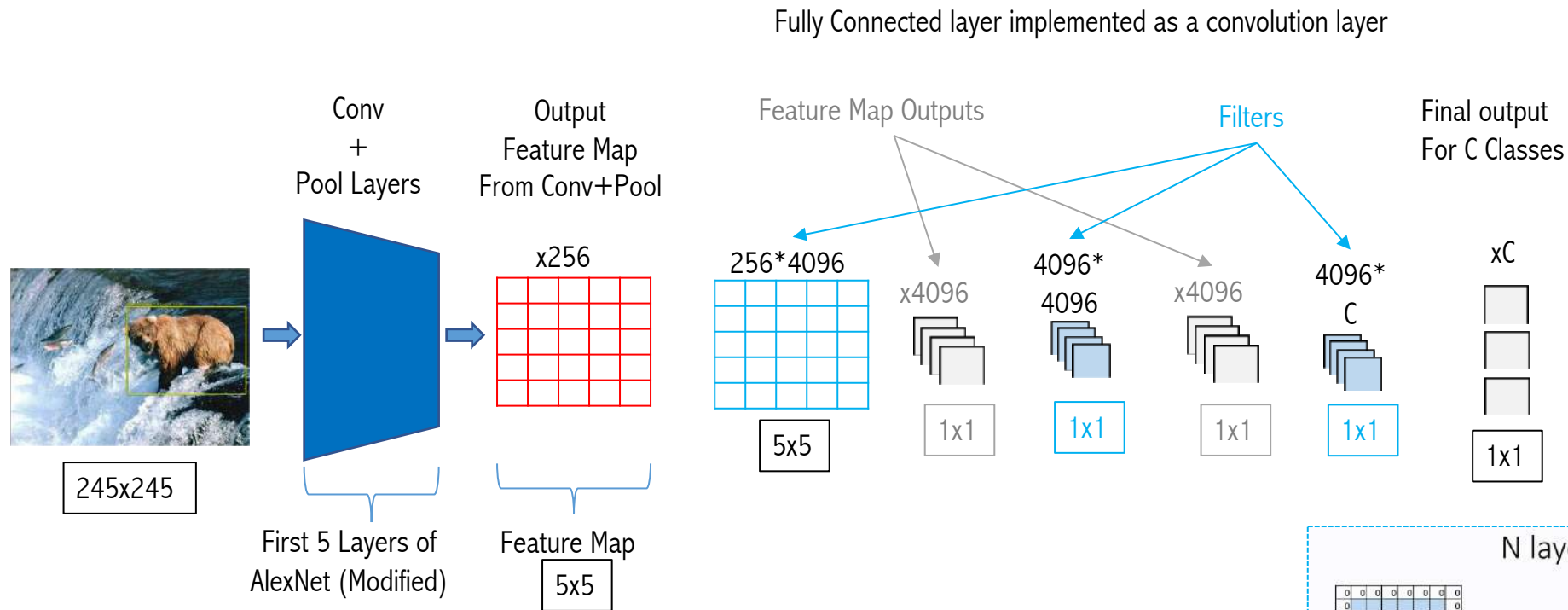
Another option is to use sliding windows of different sizes (9, as shown above) on the Feature Map. This concept is called the Feature Pyramid. This involves the use of sliding windows of 9 different sizes on top of the feature maps.



# Fast RCNN + Neural Network

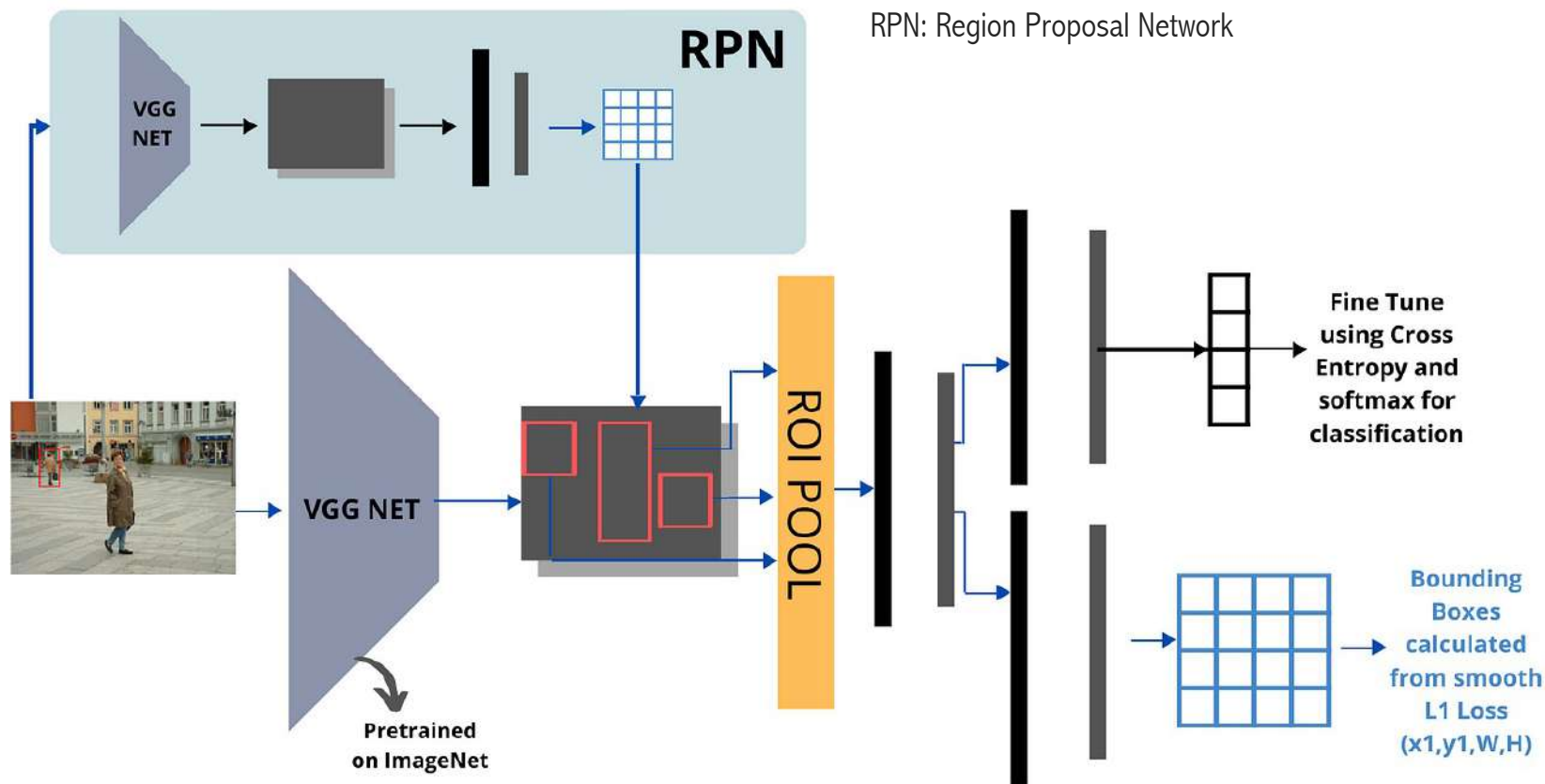


# Overfeat Classification



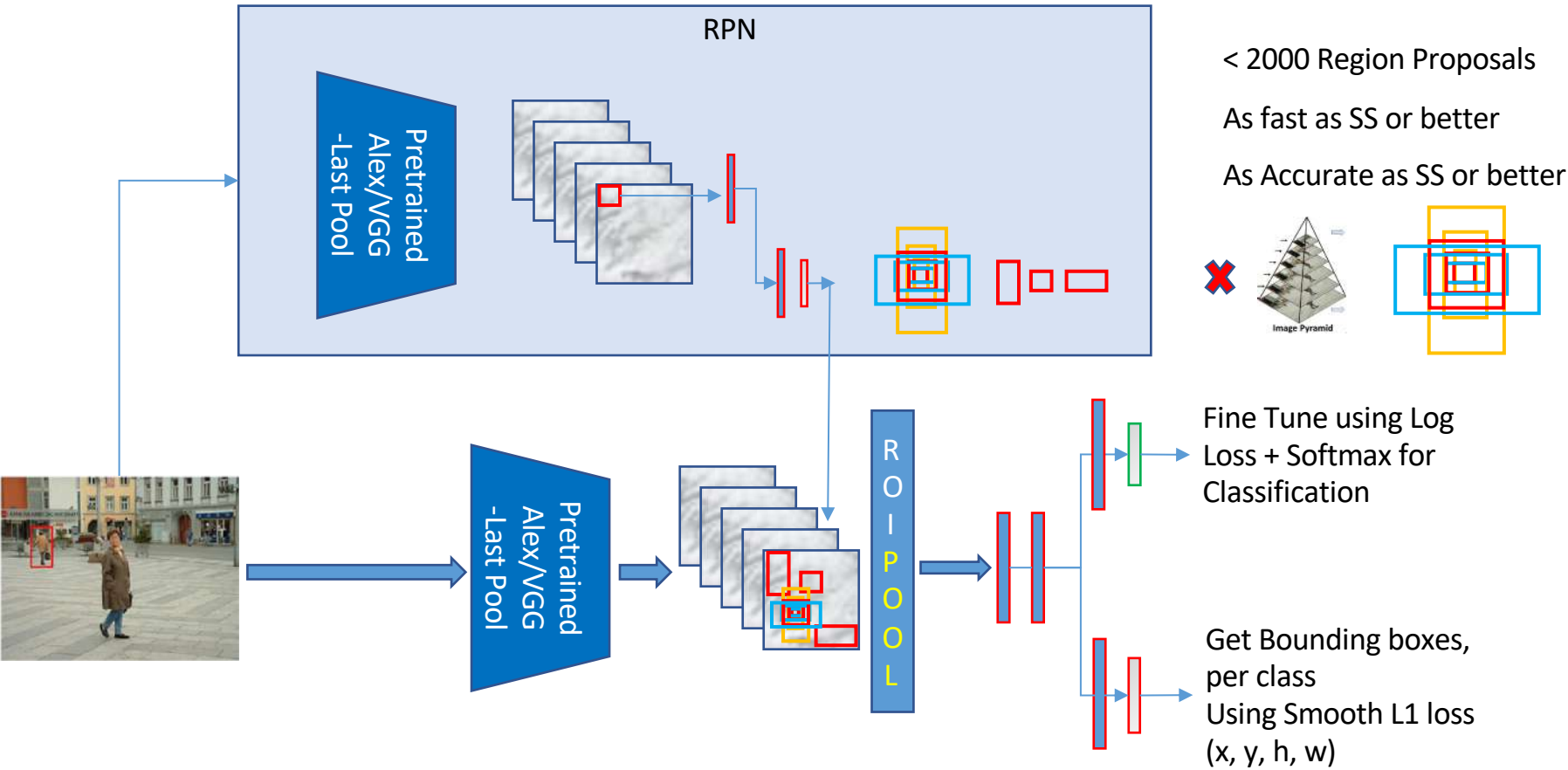
# Solutions for Replace Region Proposals

Consider using a simple CNN BBox regressor in place of Selective Search to get the approximate region proposals of the image which could further be fed to the underlying Fast R-CNN architecture. This is the core idea behind Faster R-CNN.



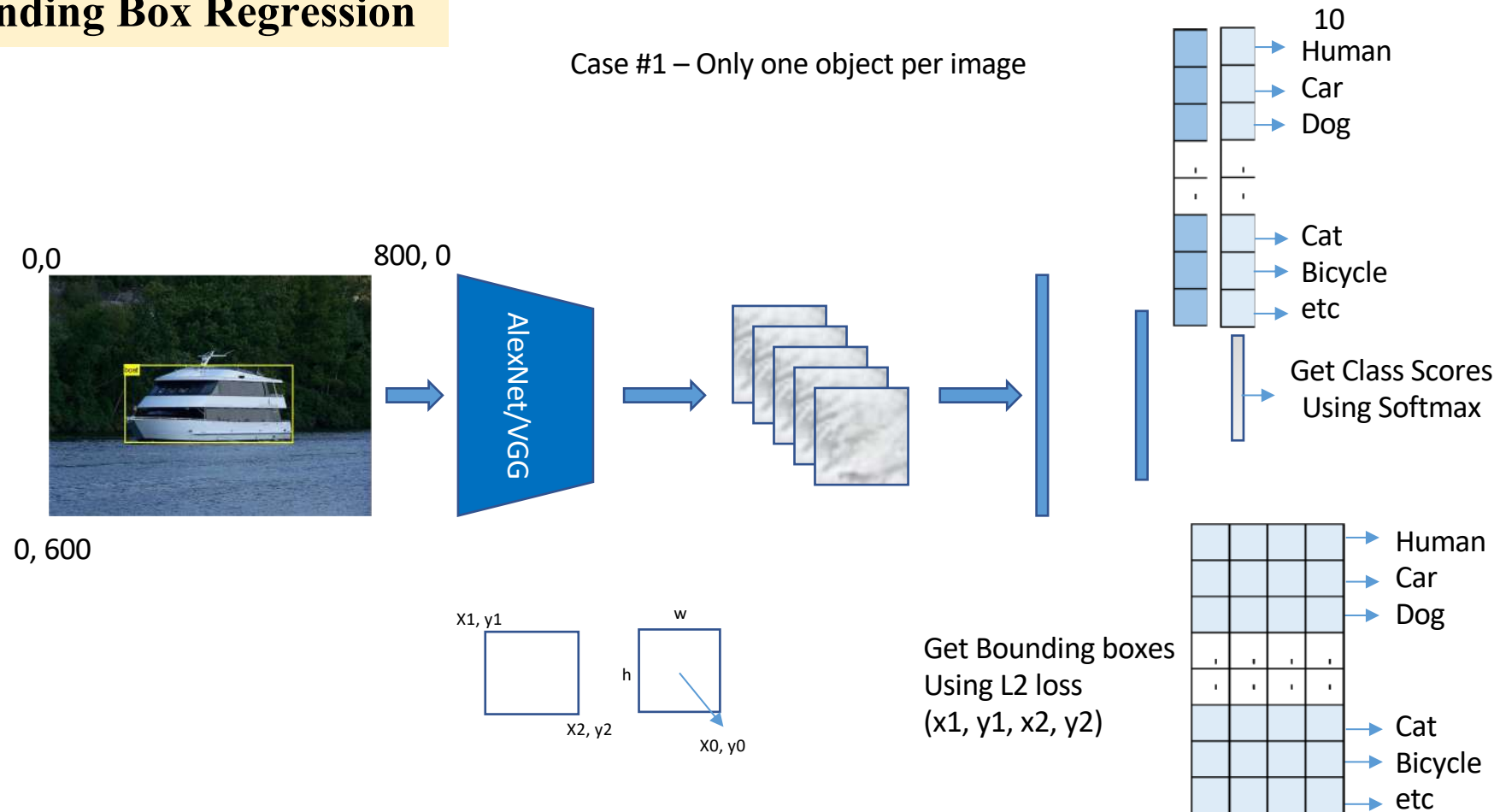
# Fast RCNN + RPN

How can we do it?



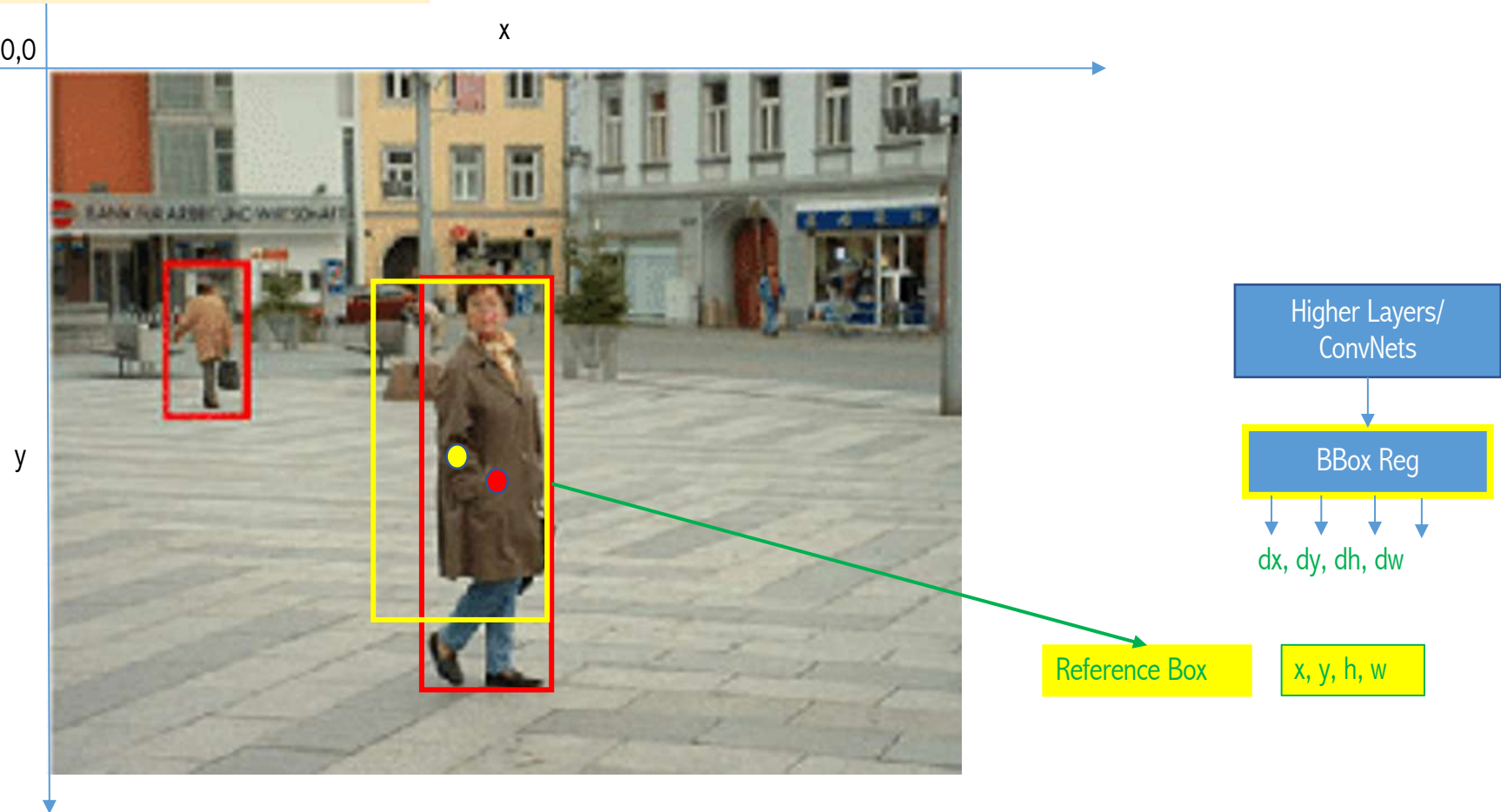


## Absolute Bounding Box Regression



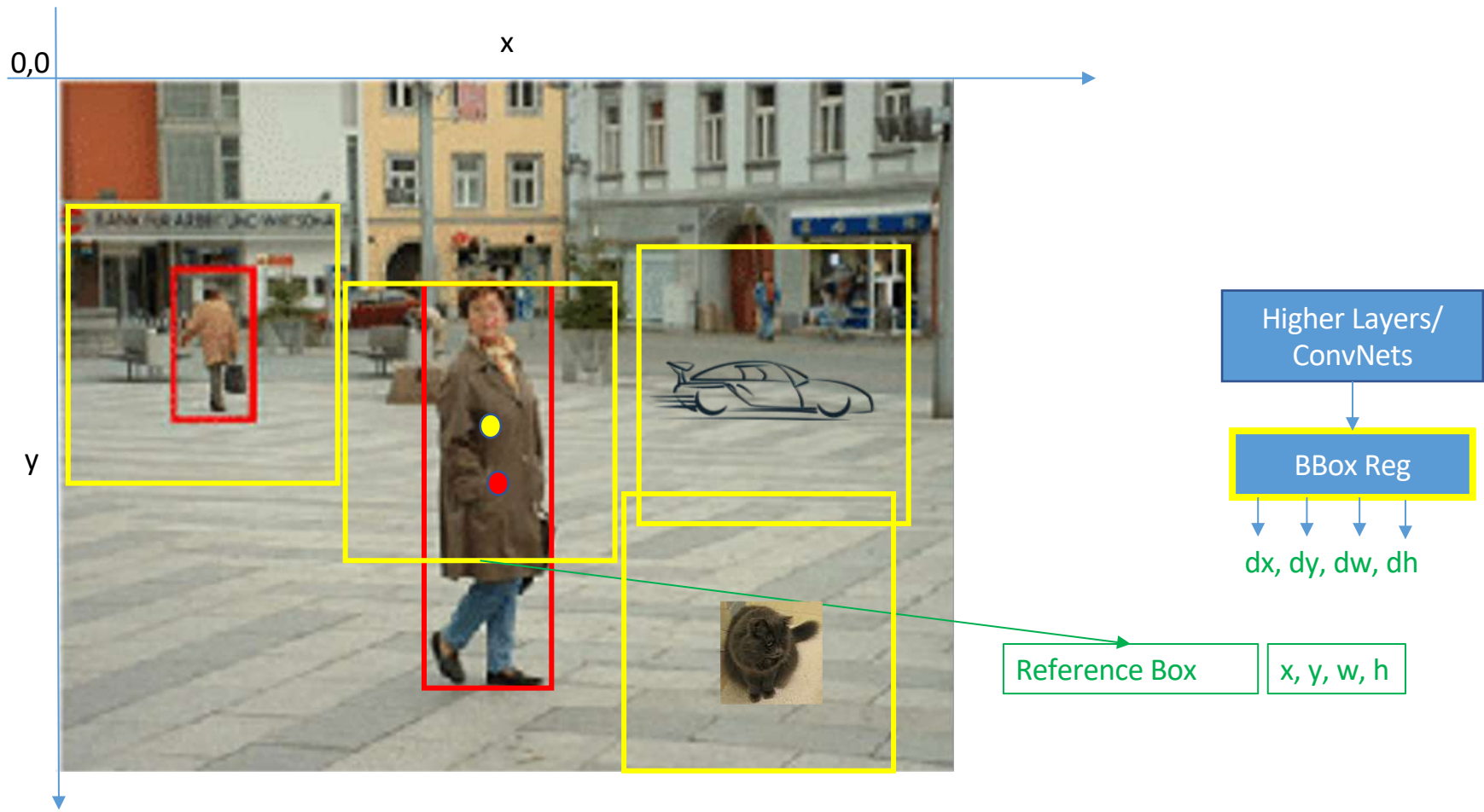
# BBox Regression - Relative

## Relative Bounding Box Regression

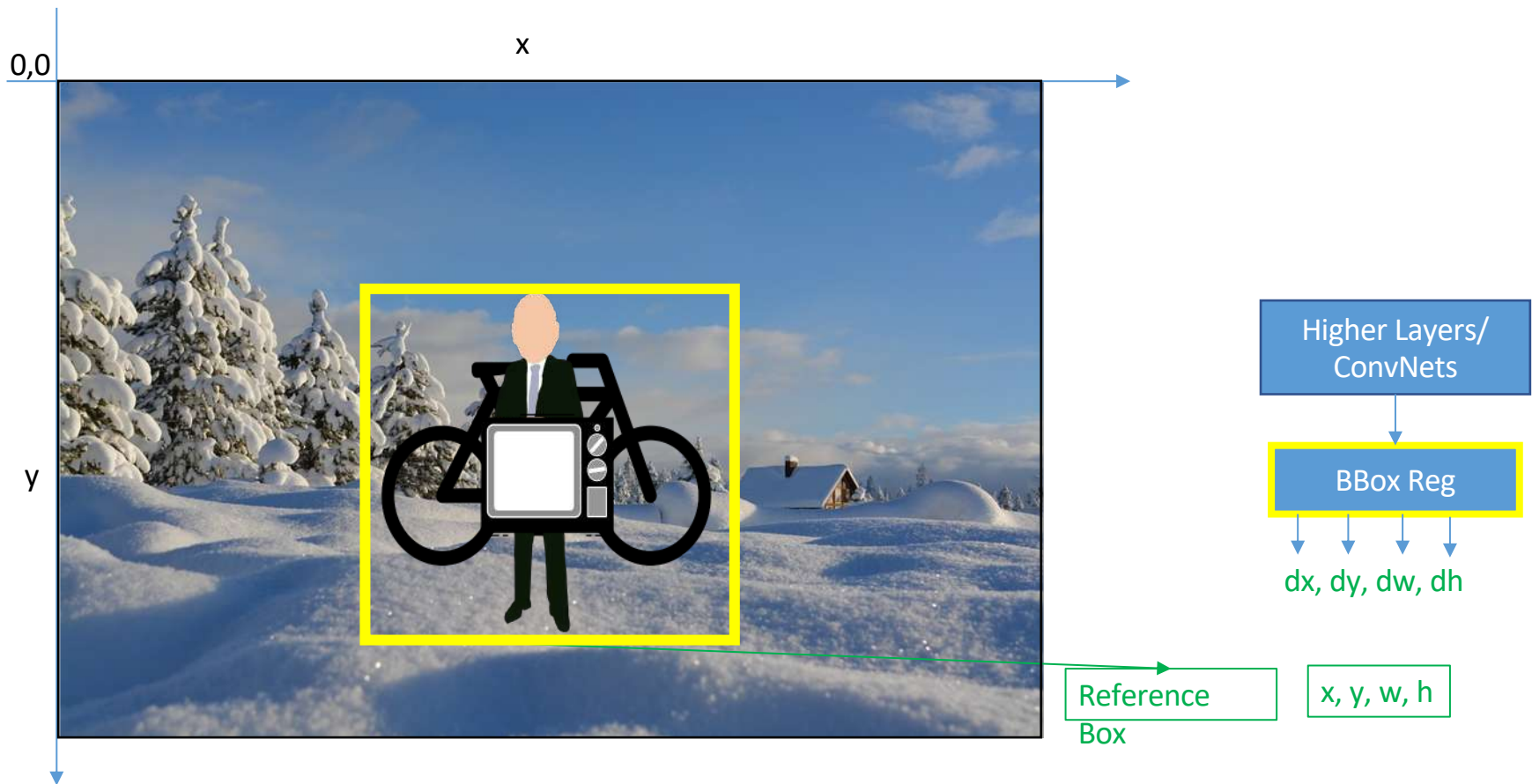


ROI reference				Bbox deltas				Predicted				Expected			
x	y	h	w	dx	dy	dh	dw	x	y	h	w	x	y	h	w
160	240	150	150	18	-22	-30	-125	178	218	120	25	180	220	120	30

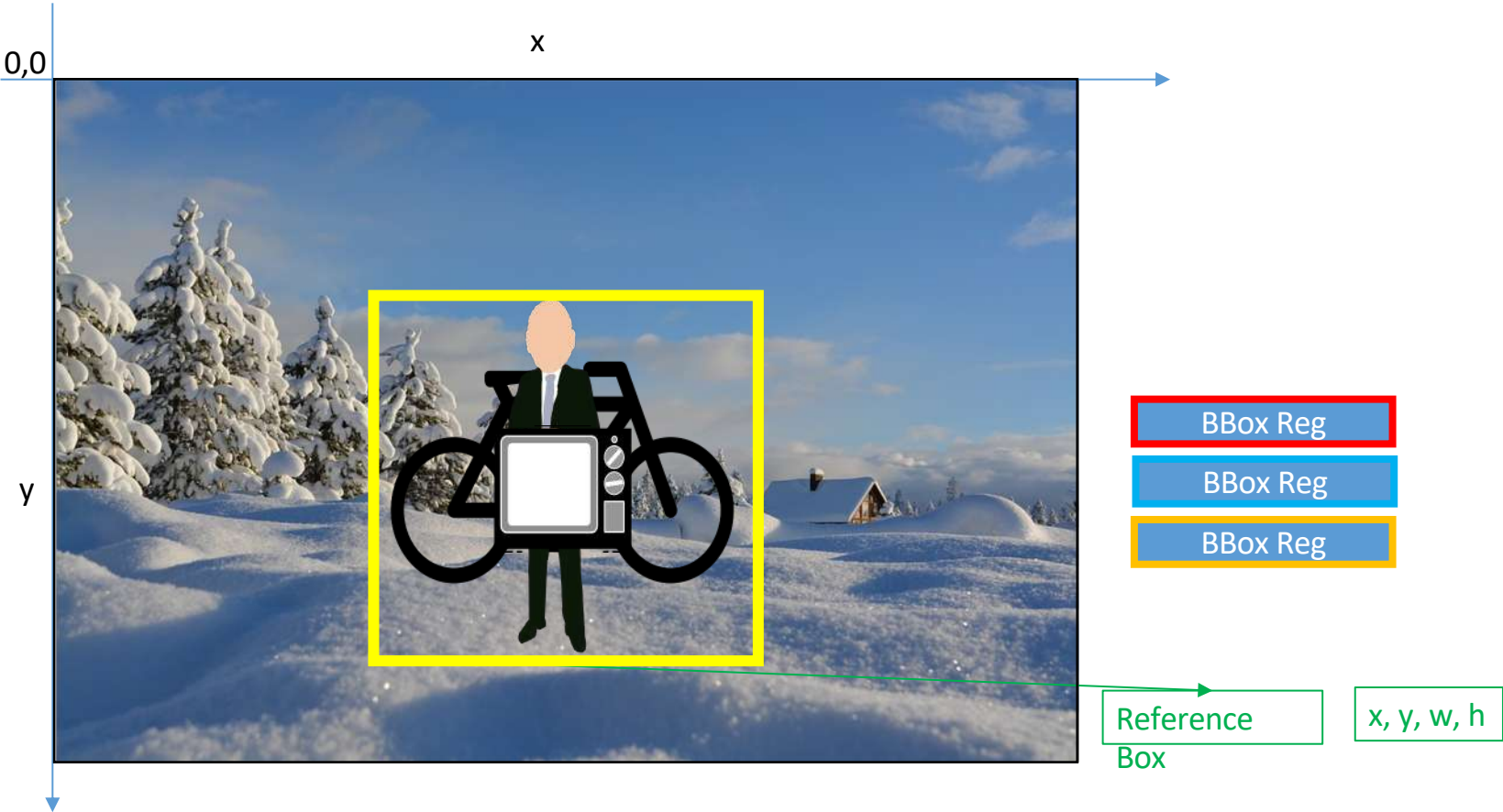
# Sliding Window as Reference Box



# Sliding Window as Reference Box

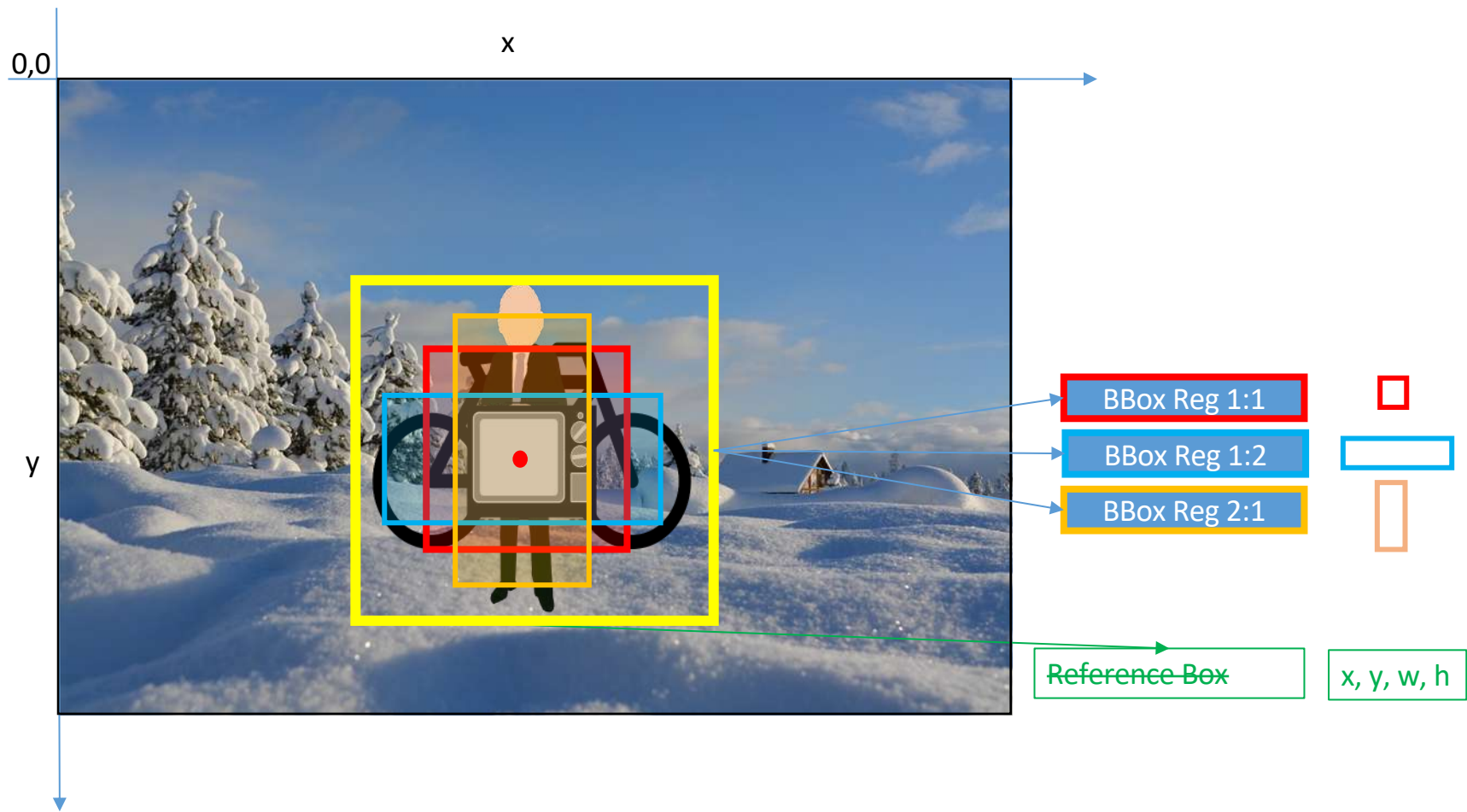


# Sliding Window as Reference Box

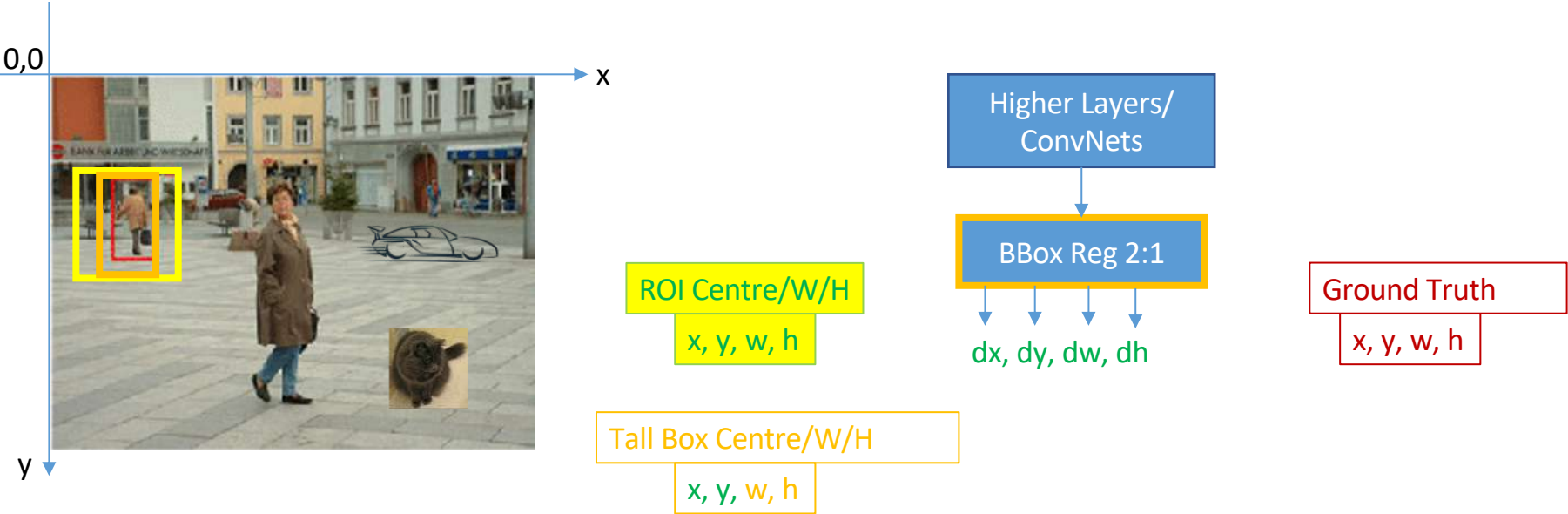




# Sliding Window as Reference Box

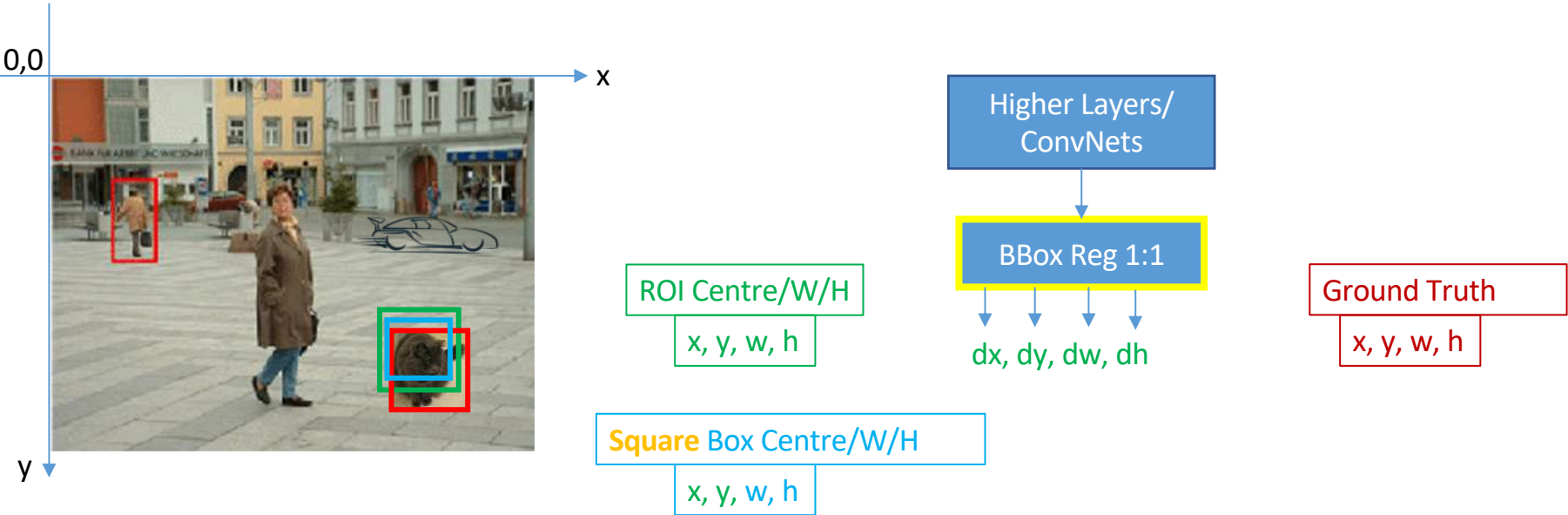


# Square ROI to Rectangular Proposal

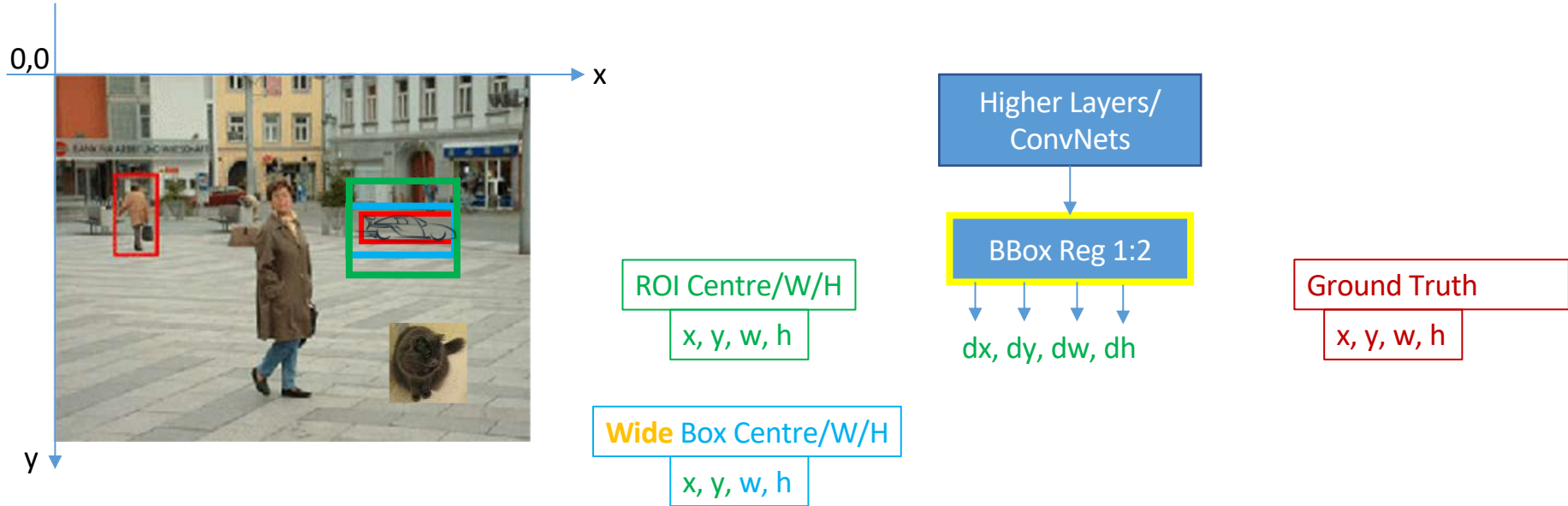


ROI reference				Bbox deltas				Predicted				Expected			
x	y	h	w	dx	dy	dh	dw	x	y	h	w	x	y	h	w
160	240	130	65	18	-22	-10	-32	178	218	120	33	180	220	120	30

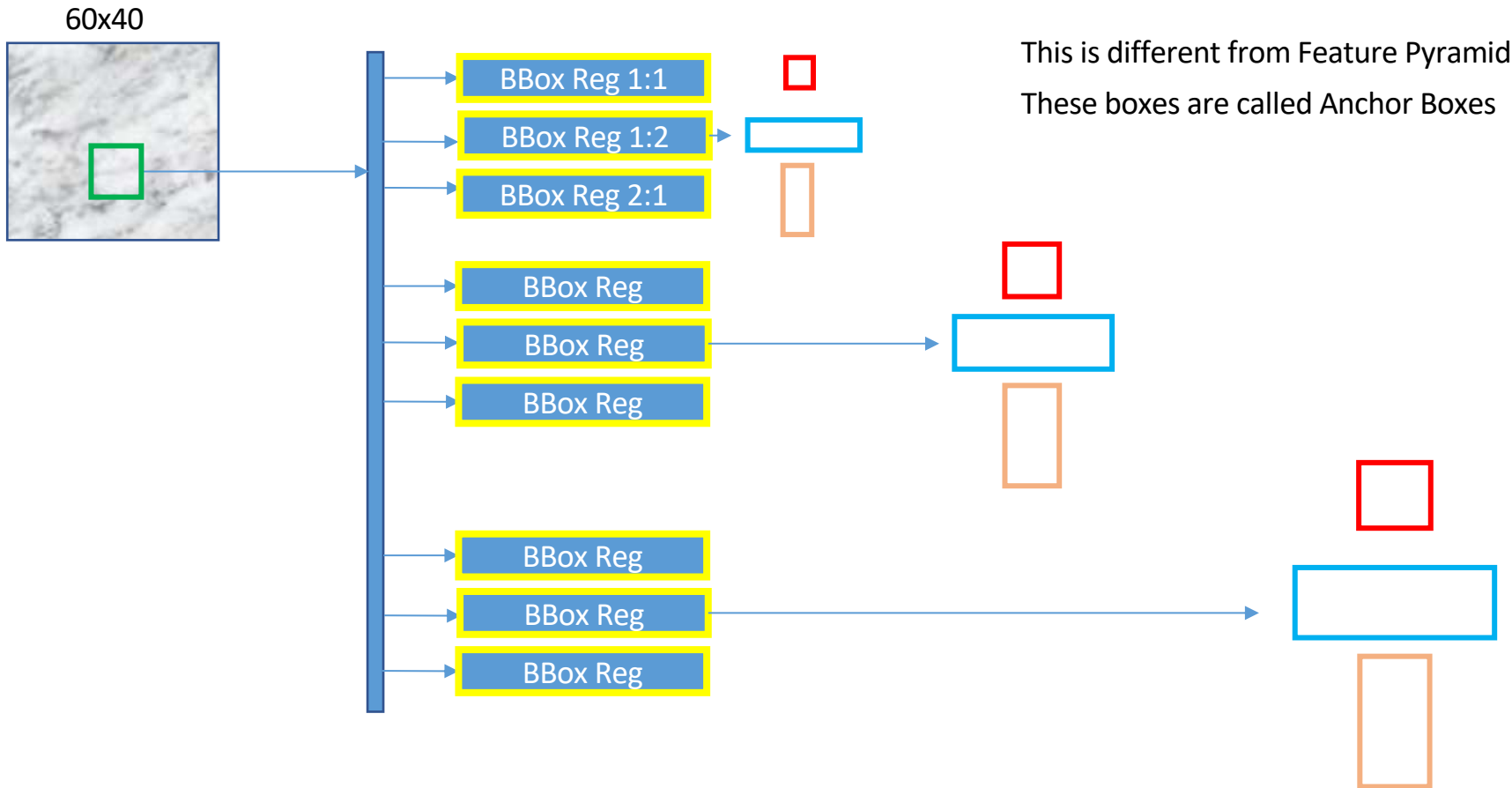
# Square ROI to Rectangular Proposal



# Square ROI to Rectangular Proposal

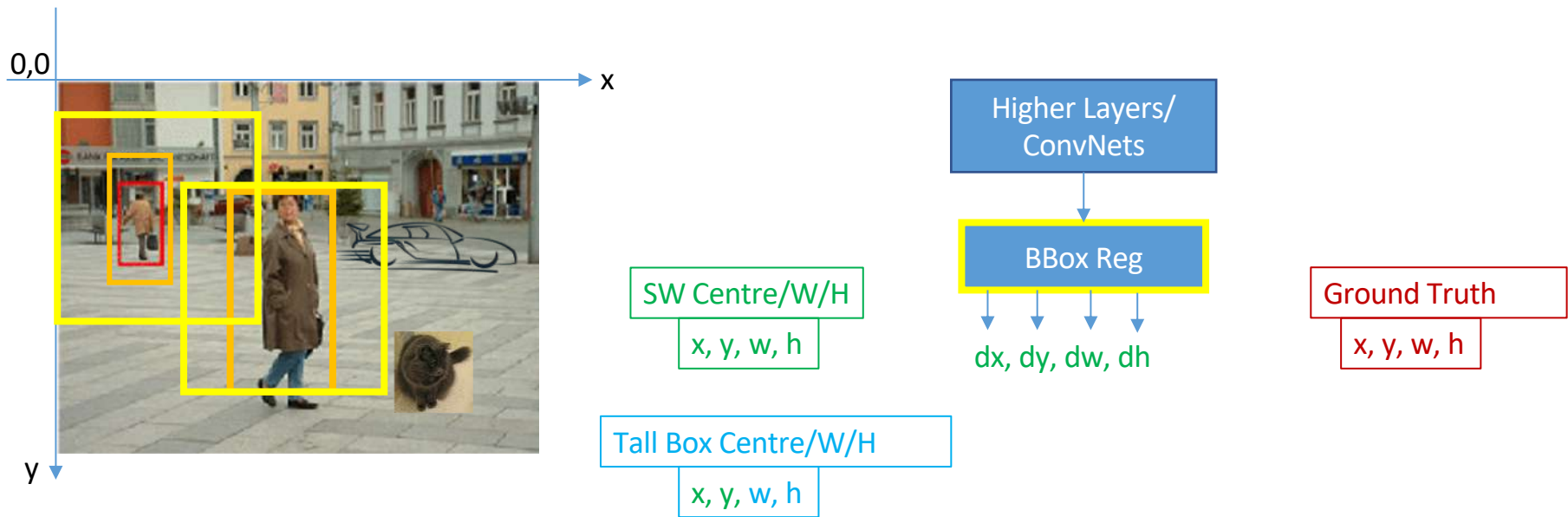


# Multiple BBox Reg using Reference Boxes

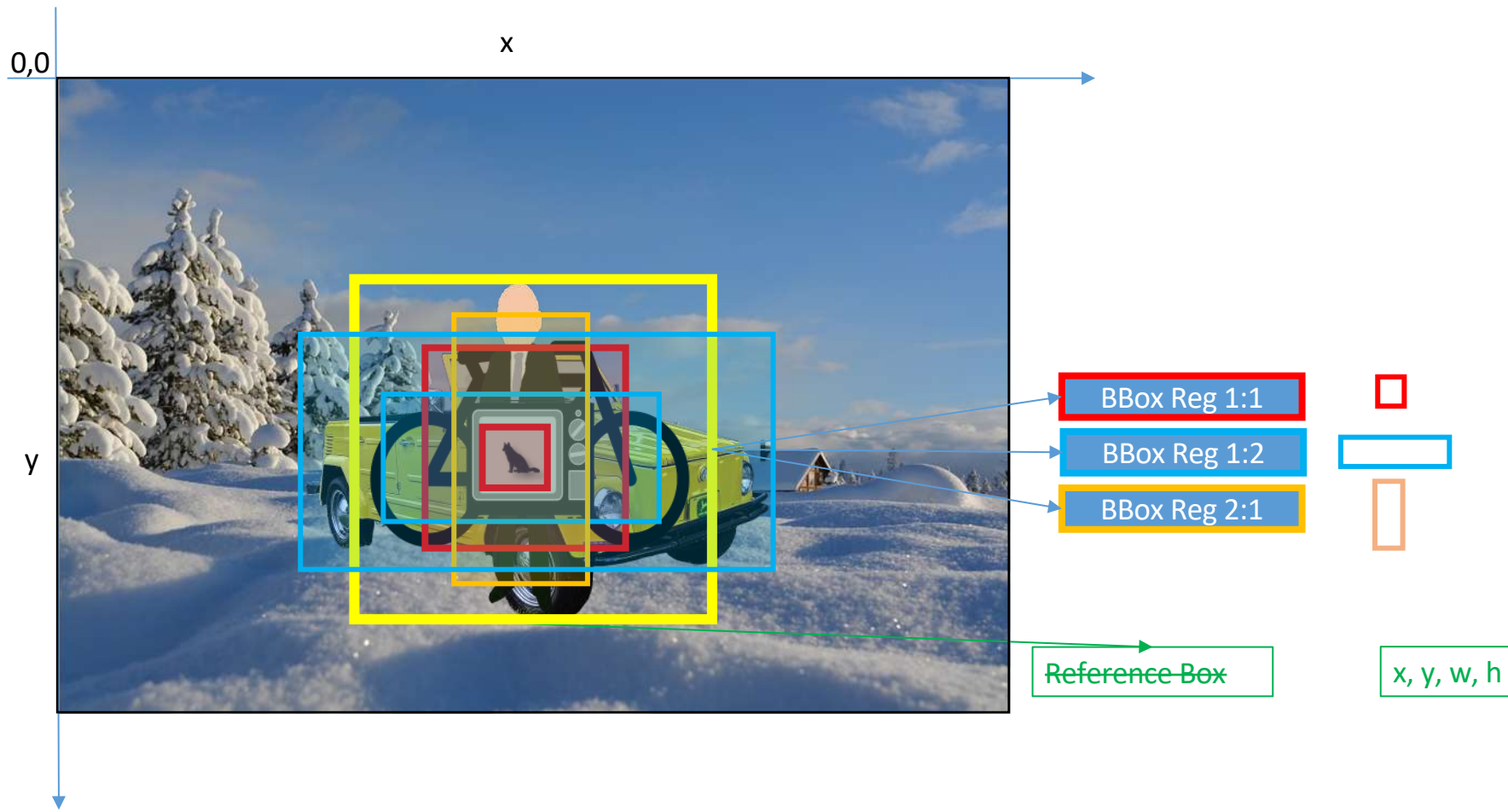




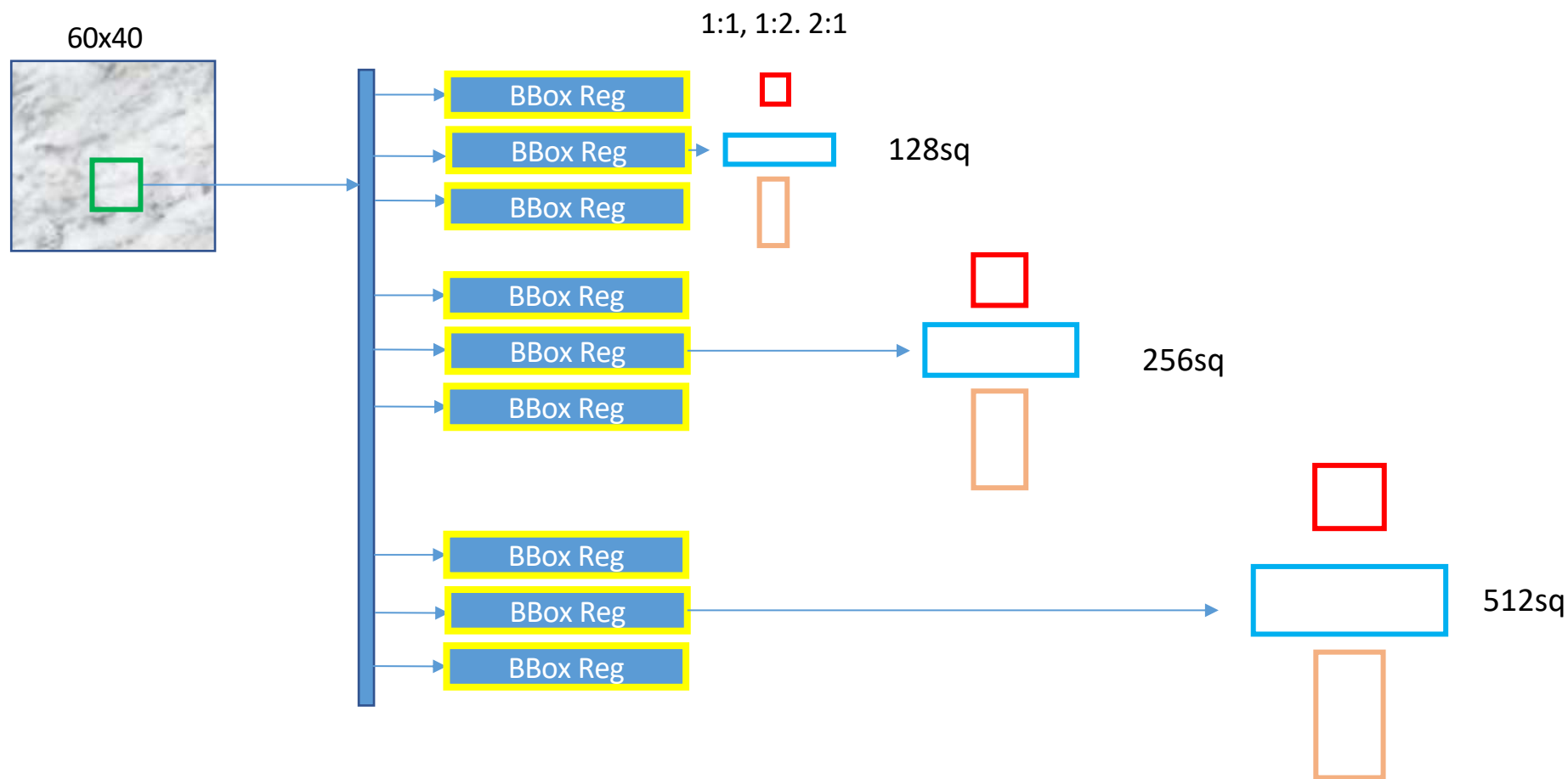
# Bigger Objects?



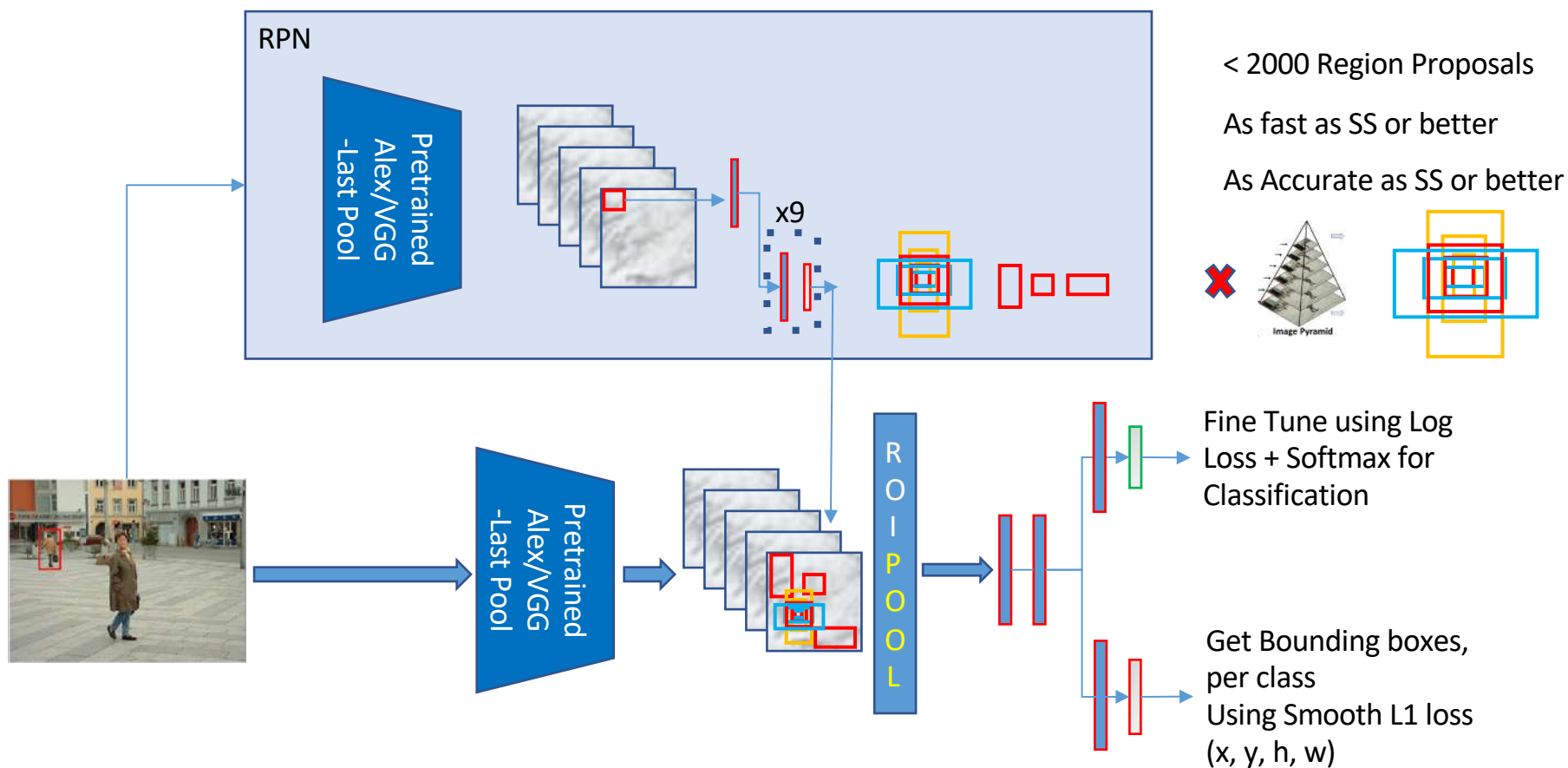
# Bigger Objects?



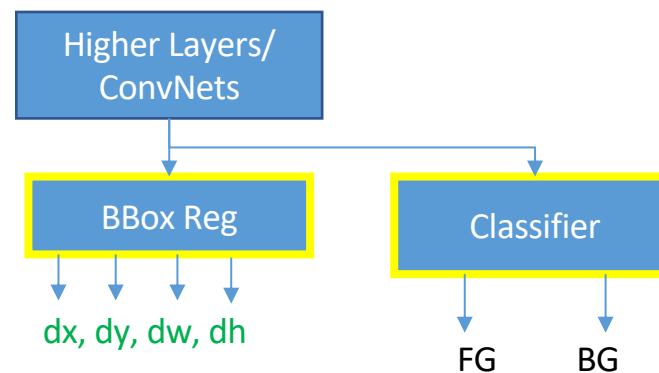
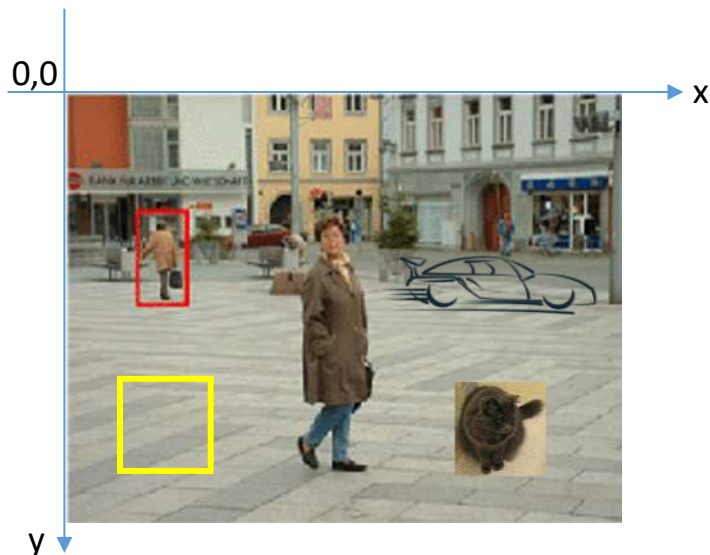
# Multiple BBox Reg using Anchor Boxes



# Fast RCNN + RPN

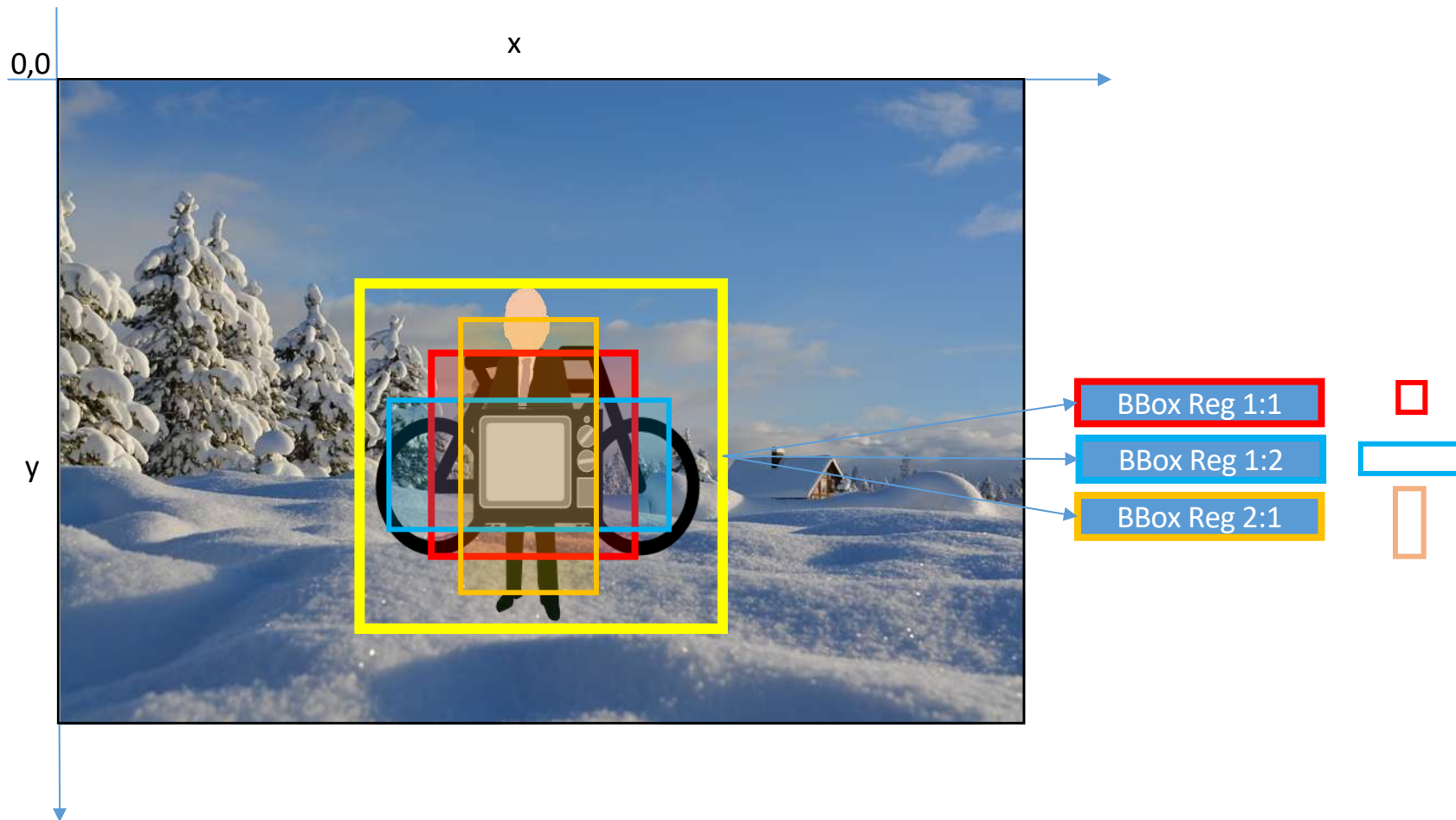


# How to reduce number of proposals?

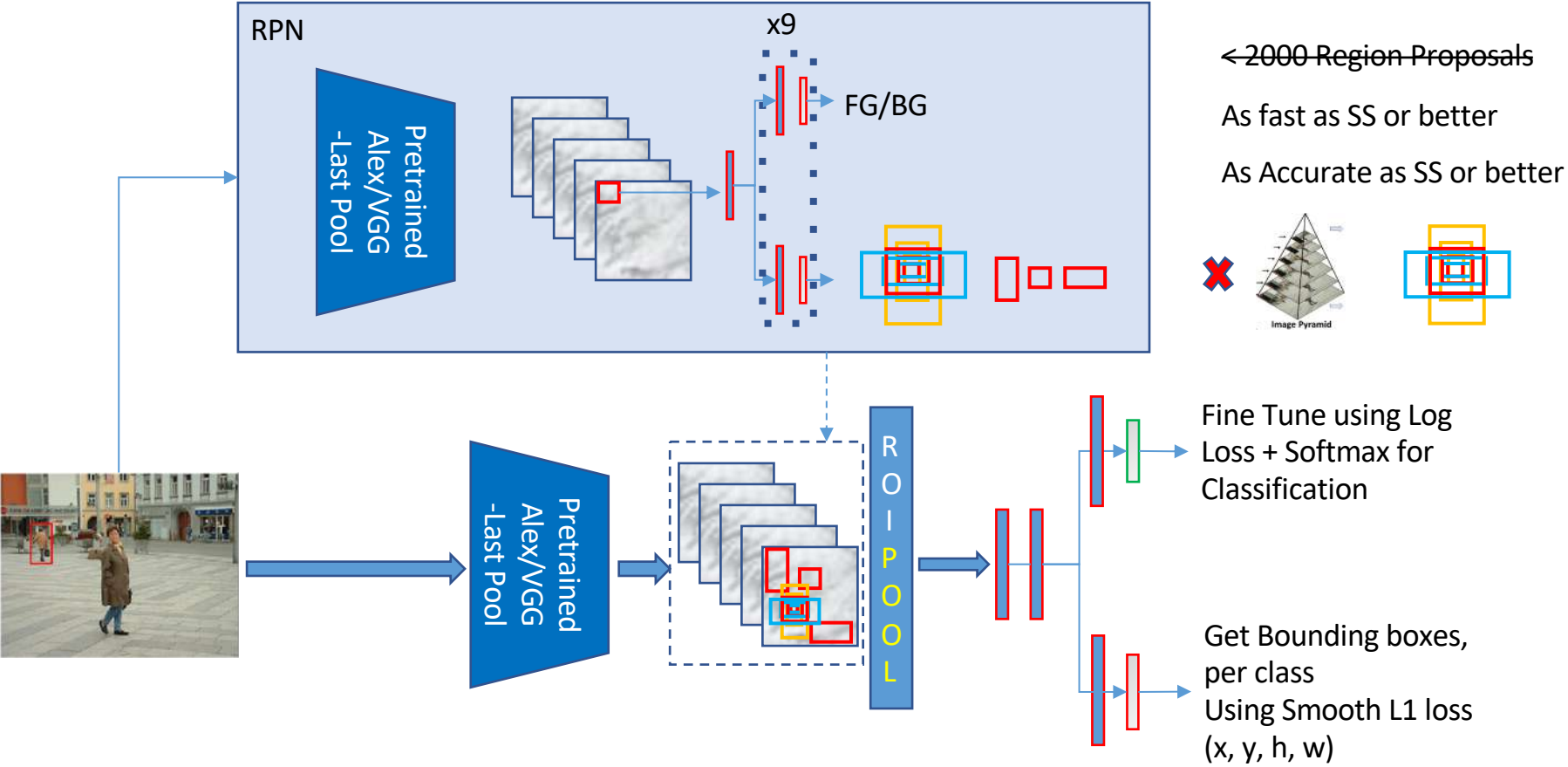




# Faster RCNN – Training Anchor Boxes - Labelling



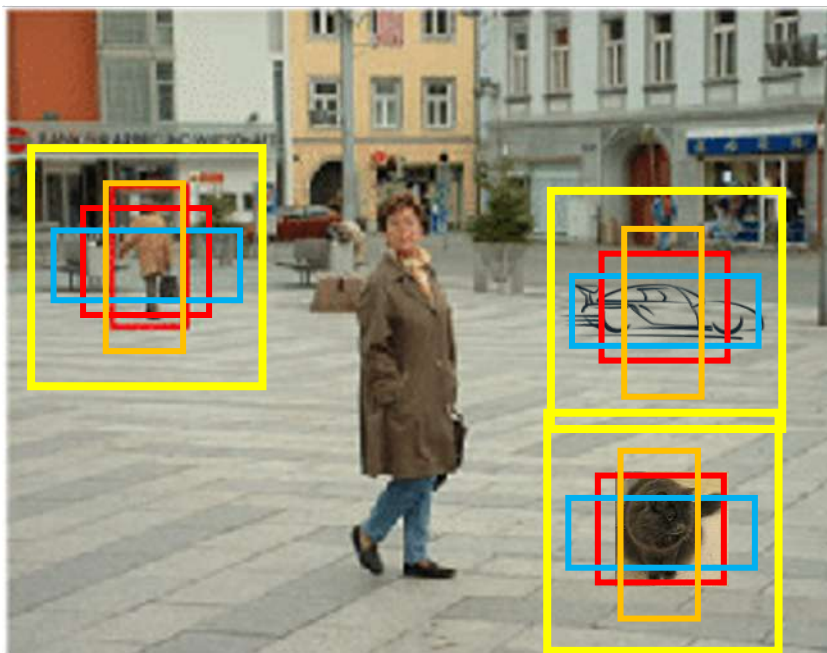
# Fast RCNN + RPN



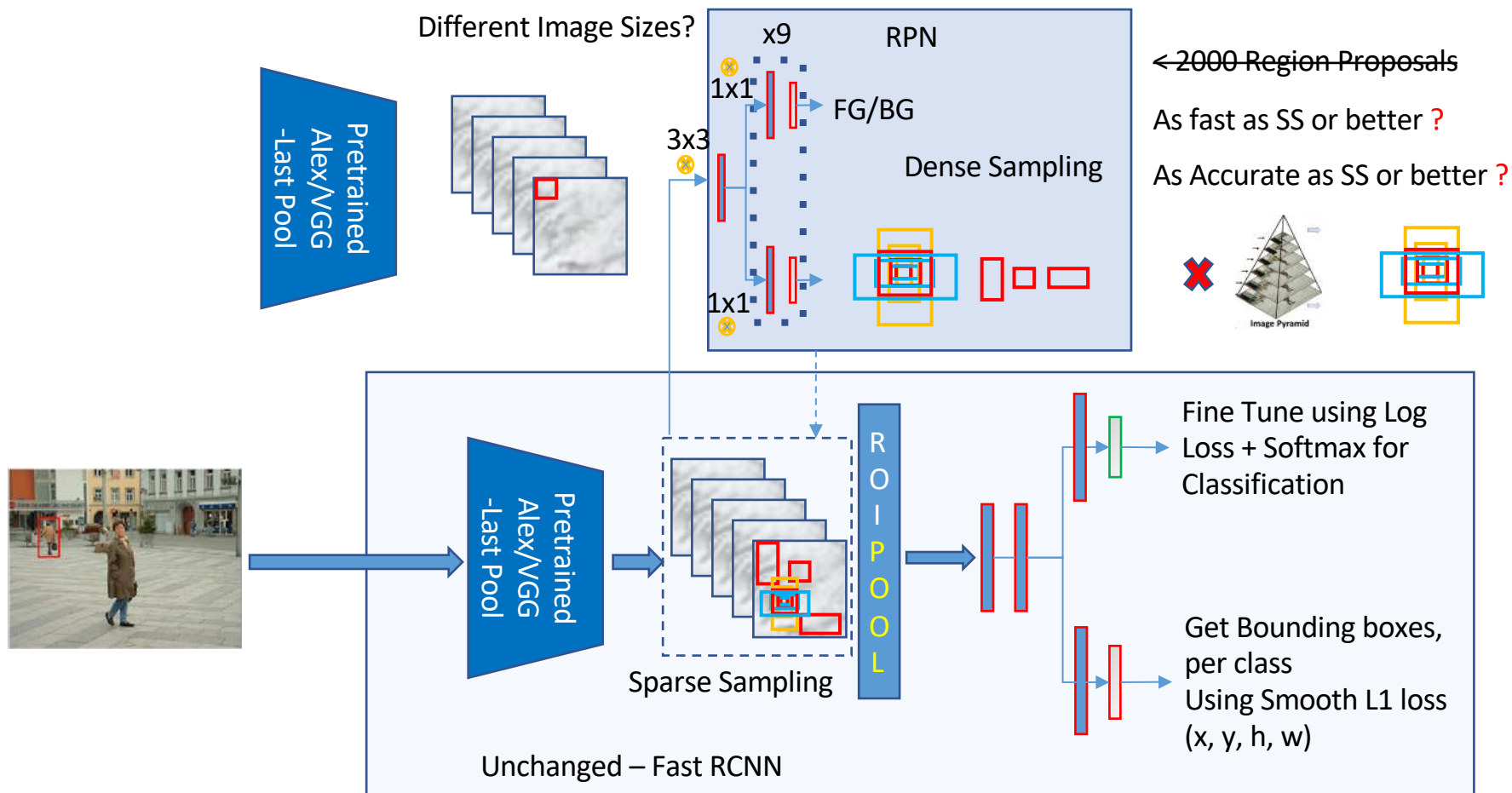
BBox Reg 1:1

BBox Reg 1:2

BBox Reg 2:1



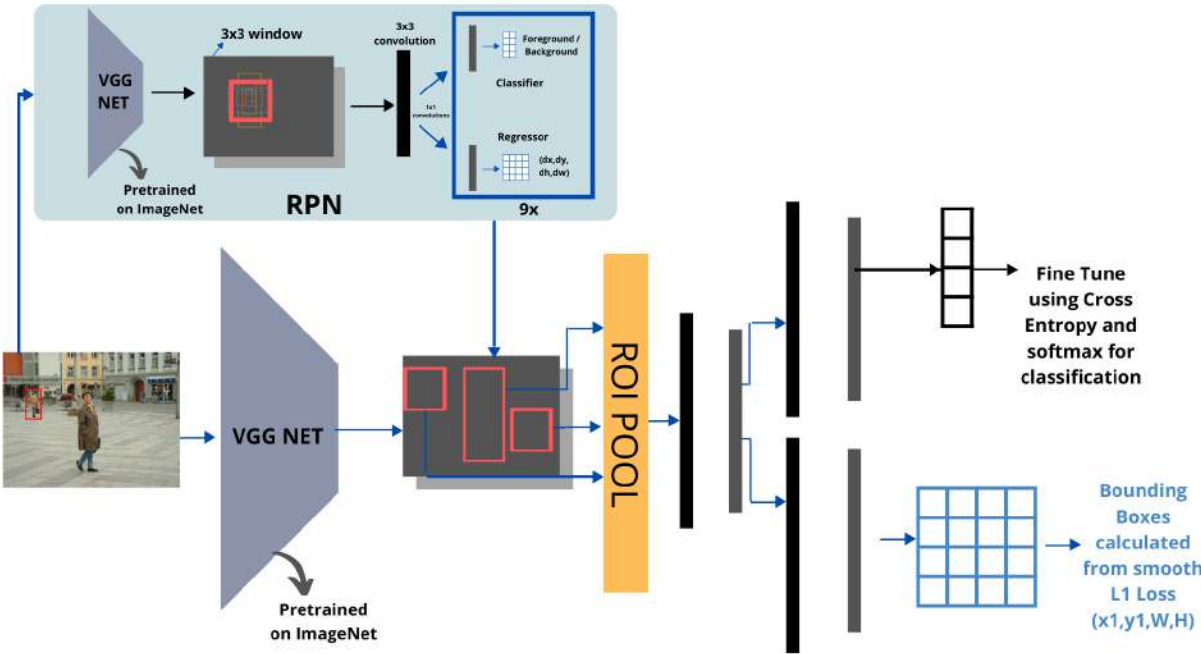
# Fast RCNN + RPN = Faster RCNN



# Faster & Accurate than SS?

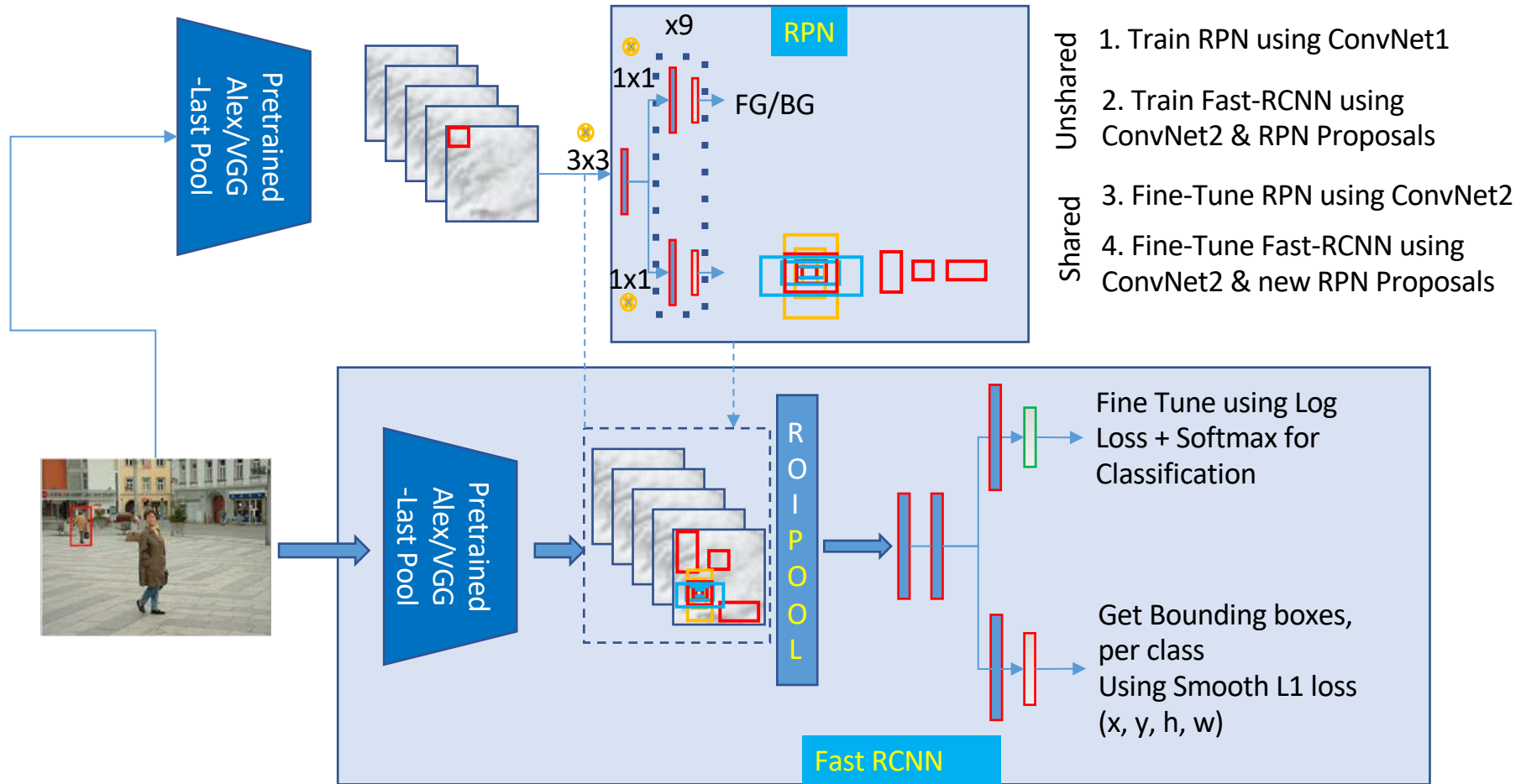
←----- Time in ms ----->

model	system	conv	proposal	region-wise	total	rate	mAP
VGG	<b>SS</b> + Fast R-CNN	146	1510	174	1830	0.5 fps	66.9
VGG	RPN + Fast R-CNN	141	10	47	198	5 fps	69.9





# Faster RCNN - Training





# Object Detection Milestones

Milestones: Traditional Detectors

*Viola Jones Detectors, SVM + HOG & DPM*

Milestones: CNN based Two-stage Detectors

*RCNN, SPPNet, Faster RCNN, Faster RCNN,..*

Milestones: CNN based One-stage Detectors

*YOLO, SSD, RetinaNet, CornerNet, Center Net,..*

Milestones: Transformer for OD

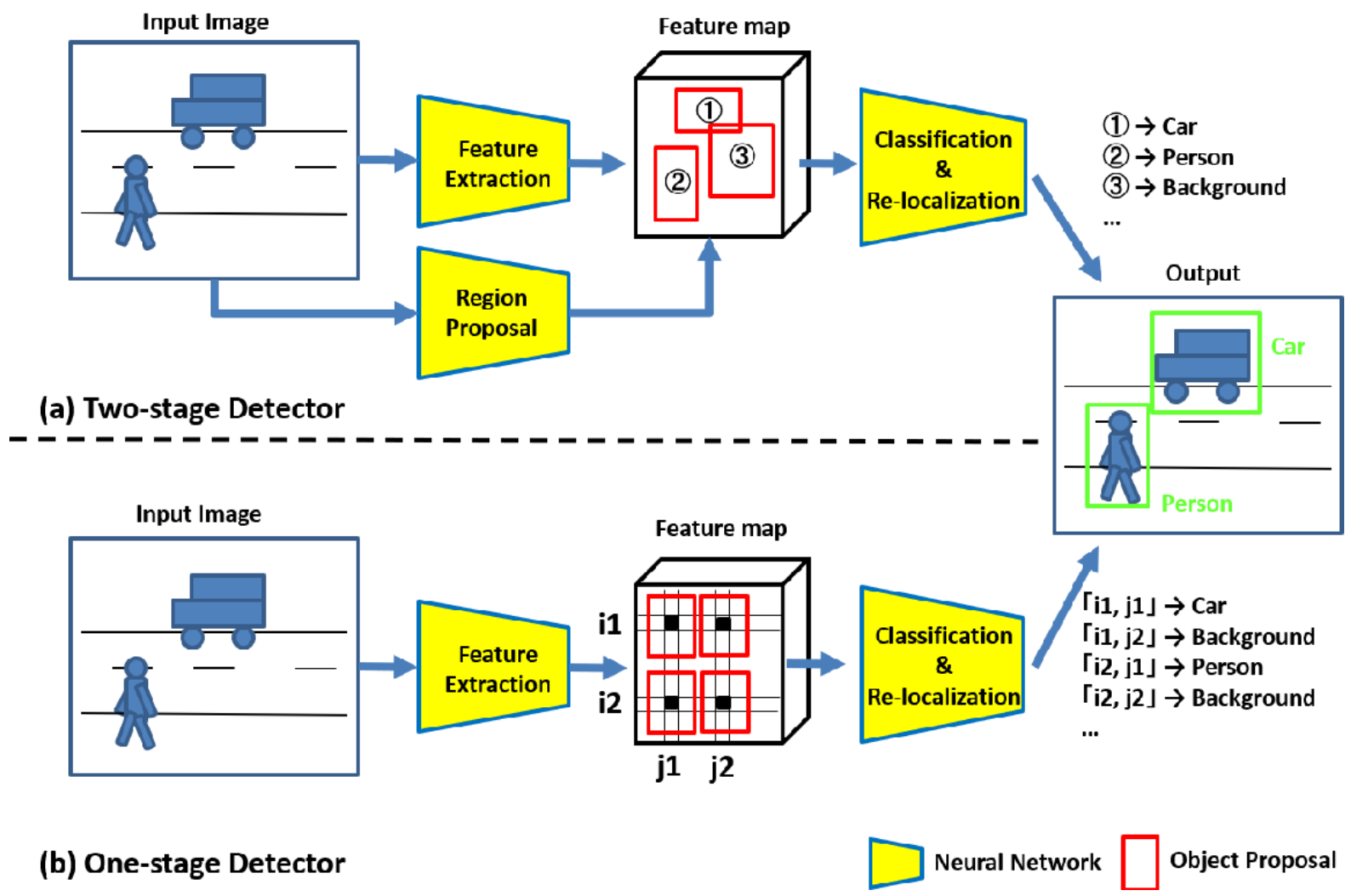
*DETR, D-DETR, DINO,..*

# Outline

- **CNN Limitations**
- **Region Based Convolutional Neural Networks**
- **Spatial Pyramid Pooling**
- **Fast R-CNN**
- **Faster RCNN**
- **YOLOv1-v2**

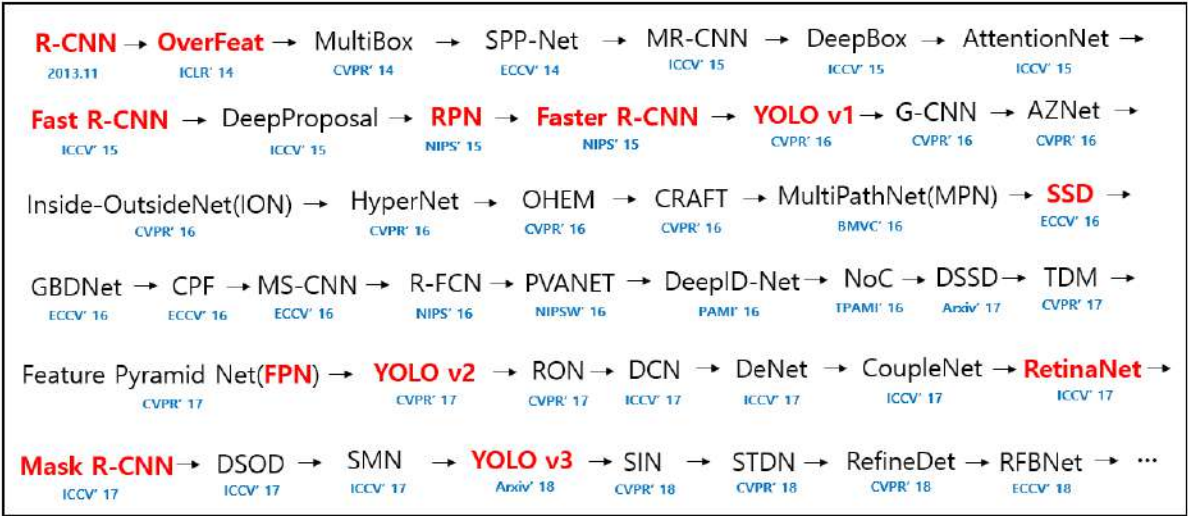
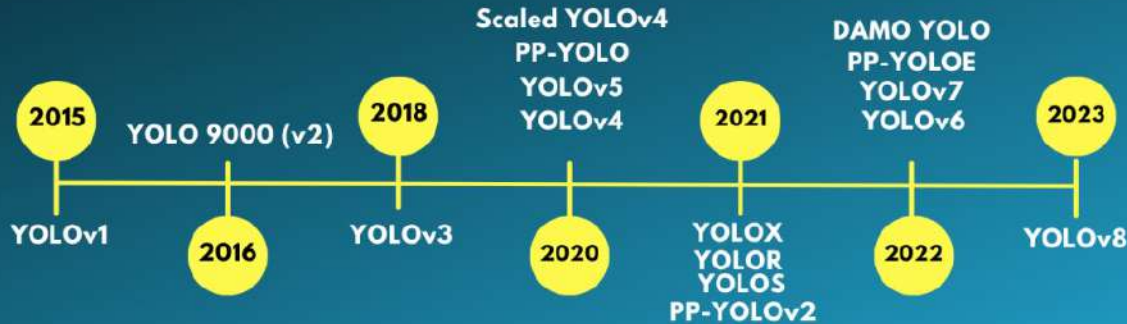


# One-stage vs. Two-stage Detectors



# Object Detection Milestones

## YOLO Object Detection Models Timeline



Z. Zou, K. Chen, Z. Shi, Y. Guo and J. Ye, "Object Detection in 20 Years: A Survey," in Proceedings of the IEEE, vol. 111, no. 3, pp. 257-276, March 2023, doi: 10.1109/JPROC.2023.3238524.



# Why Study YOLO?

YOLO is an abbreviation for the term 'You Only Look Once'. This is an algorithm that detects and recognizes various objects in a picture (in real-time). Object detection in YOLO is done as a regression problem and provides the class probabilities of the detected images.

YOLO algorithm employs convolutional neural networks (CNN) to detect objects in real-time. As the name suggests, the algorithm requires only a single forward propagation through a neural network to detect objects.



**Speed:** This algorithm improves the speed of detection because it can predict objects in real-time.



**High accuracy:** YOLO is a predictive technique that provides accurate results with minimal background errors.

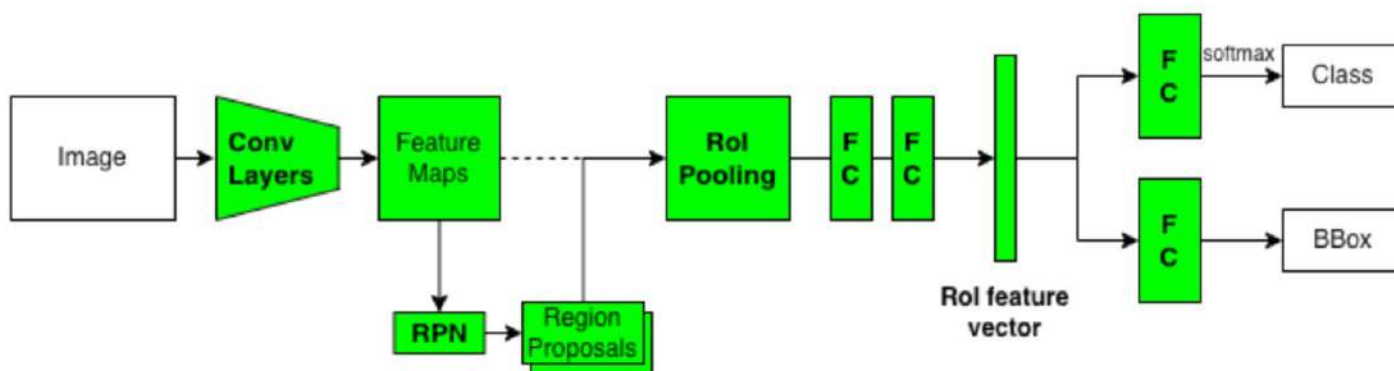


**Learning capabilities:** The algorithm has excellent learning capabilities that enable it to learn the representations of objects and apply them in object detection.

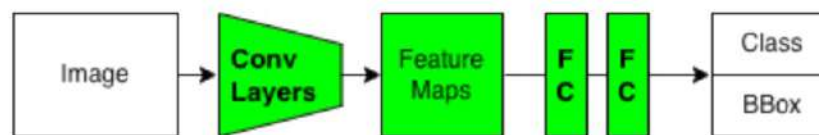
For small datasets and limited computational power, YOLOv8 might be a better choice as it has been optimized for speed and accuracy. However, if you have more significant datasets and more complex object detection tasks, DETR could be a better fit due to its ability to handle object detection without pre-defined anchor boxes

# YOLO: Motivations

In 2015, Joseph Redmon (University of Washington) developed YOLO. One of his co-authors, Ross Girshick (Microsoft Research), published a paper for Faster R-CNN around the same time. They probably shared common ideas in computer vision research as there are some similarities between YOLO and Faster R-CNN. For example, both models apply convolutional layers on input images to generate feature maps. However, Faster R-CNN uses a two-stage object detection pipeline, while YOLO has no separate region proposal step and is much faster than Faster R-CNN.



*Faster R-CNN pipeline*



*YOLO pipeline*

# YOLO: Motivations

In 2015, Joseph Redmon (University of Washington) developed YOLO. One of his co-authors, Ross Girshick (Microsoft Research), published a paper for Faster R-CNN around the same time. They probably shared common ideas in computer vision research as there are some similarities between YOLO and Faster R-CNN. For example, both models apply convolutional layers on input images to generate feature maps. However, Faster R-CNN uses a two-stage object detection pipeline, while YOLO has no separate region proposal step and is much faster than Faster R-CNN.

YOLO has many versions (variants). Joseph Redmon developed the first three versions of YOLO: YOLOv1, v2, and v3. Then, he quit.

After YOLOv3, different groups of people developed their versions of YOLO:

- **YOLOv4** by Alexey Bochkovskiy, et al.
- **YOLOv5** by [Ultralytics](#)
- **YOLOv6** by [Meituan](#)
- **YOLOX** by Zheng Ge et al.
- **YOLOv7** by Chien-Yao Wang et al. (The same people from YOLOv4)



# YOLO: Motivations



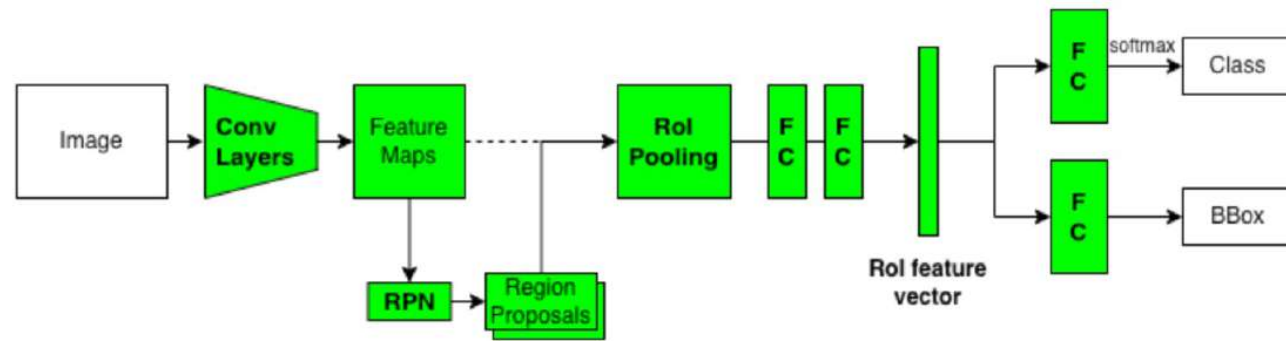
Residual blocks



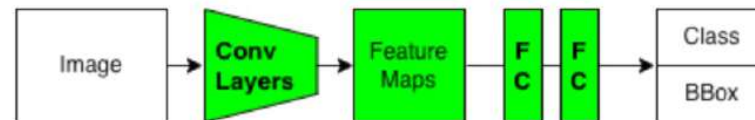
Bounding box regression



Intersection Over Union (IOU)

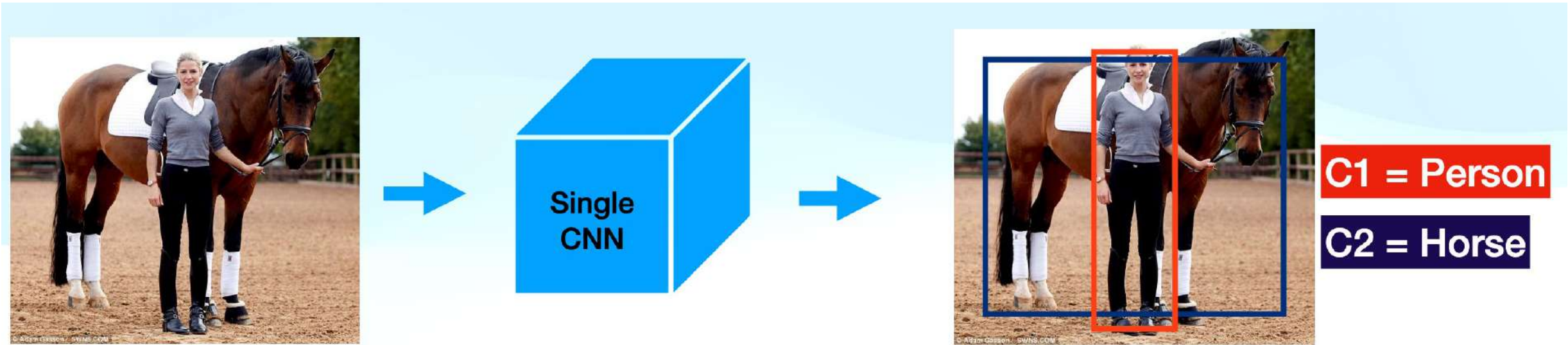


*Faster R-CNN pipeline*



*YOLO pipeline*

# YOLO: Motivations



Single stage regression problem

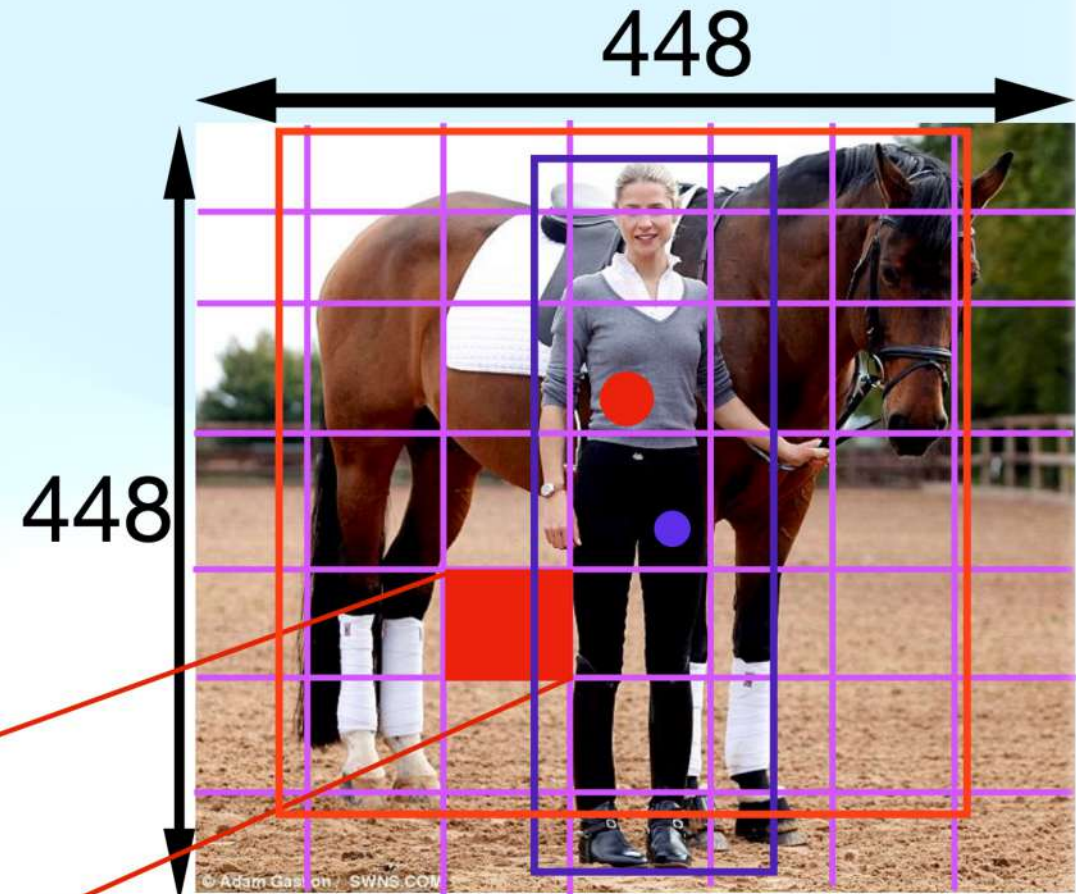


## Steps in YOLO

- Take input image
- Resize to 448x448
- Divide into SxS Grid cells
- S=7 in paper
- Each cell is responsible for predicting one object
- Which cells are responsible for person and horse?



64x64 Cell

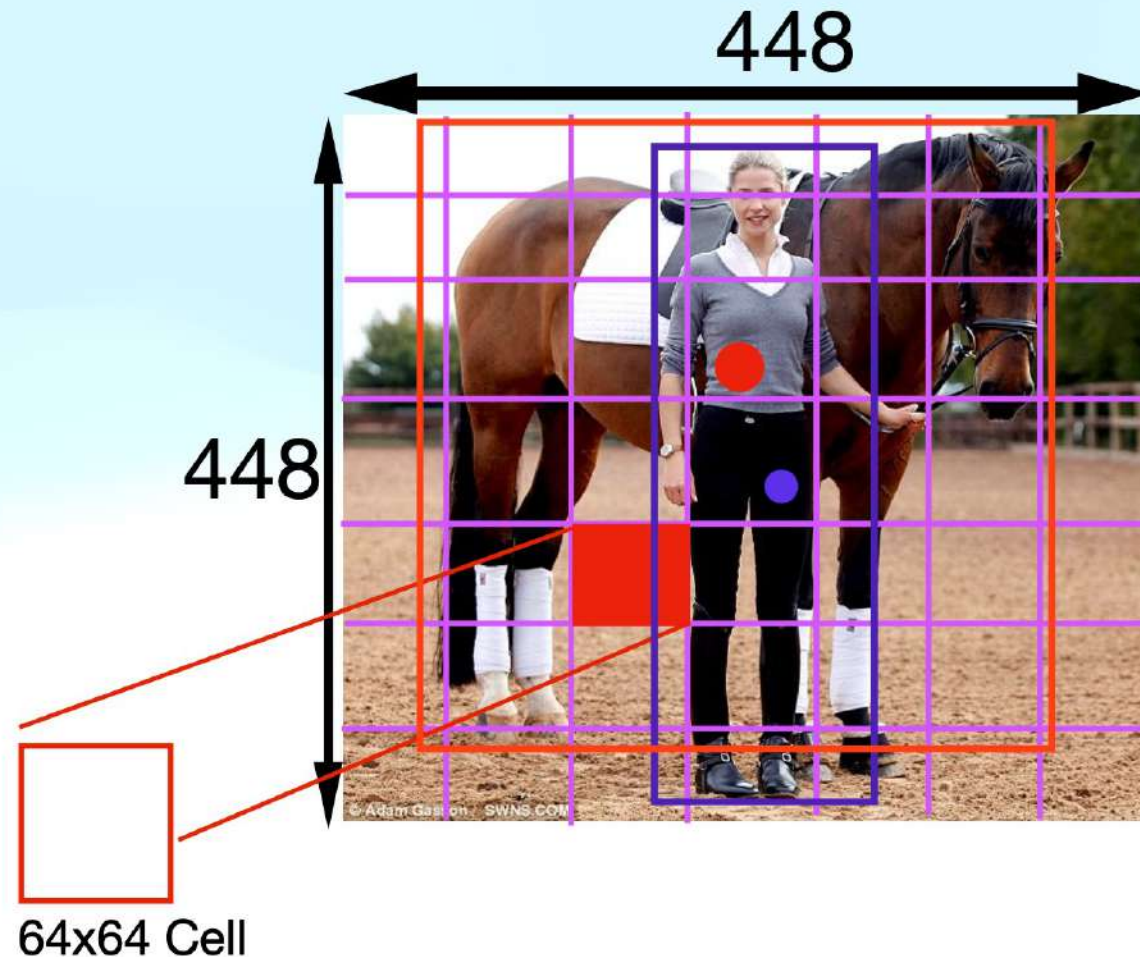




# YOLO: Motivations

## Steps in YOLO

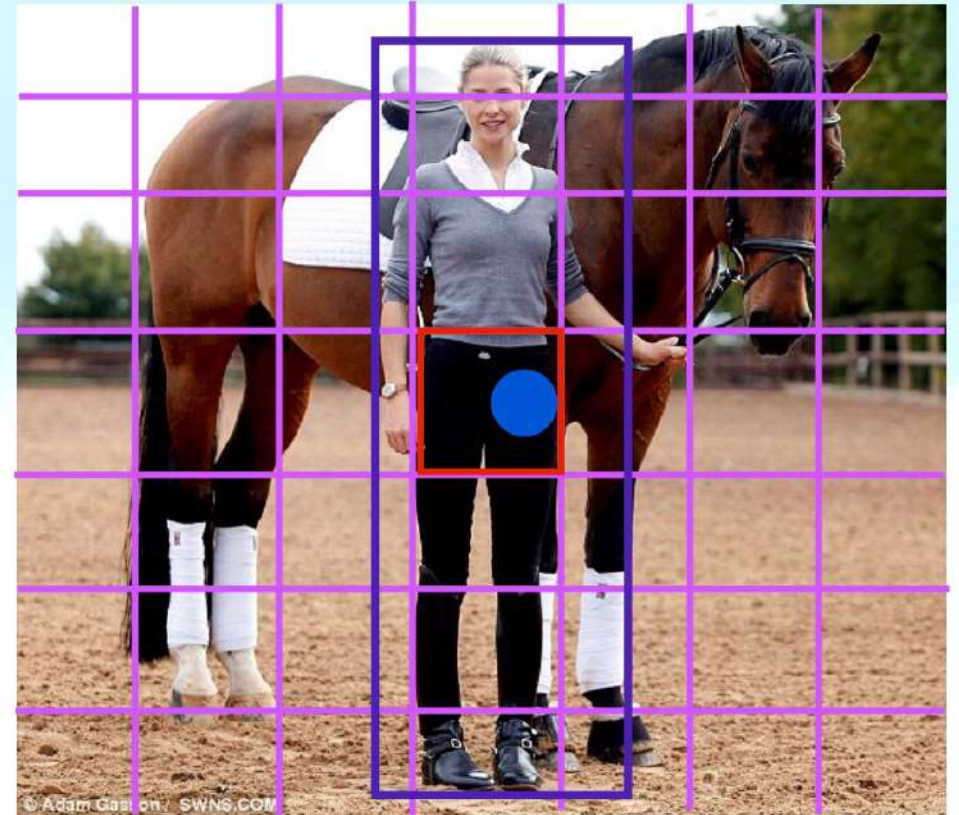
- Take input image
- Resize to 448x448
- Divide into SxS Grid cells
- S=7 in paper
- Each cell is responsible for predicting one object
- Which cell is responsible?
- Where Center of object falls into



# Bounding boxes

- Box -  $x, y, w, h$
- How are these encoded?
- Relative to Grid cell that the object Center falls into.

(200, 311, 142, 250)





# Bounding Boxes

- Center point (x,y): Relative to anchor that (x,y) falls into.

$$\Delta x = (x - x_a) / 64$$

$$\Delta y = (y - y_a) / 64$$

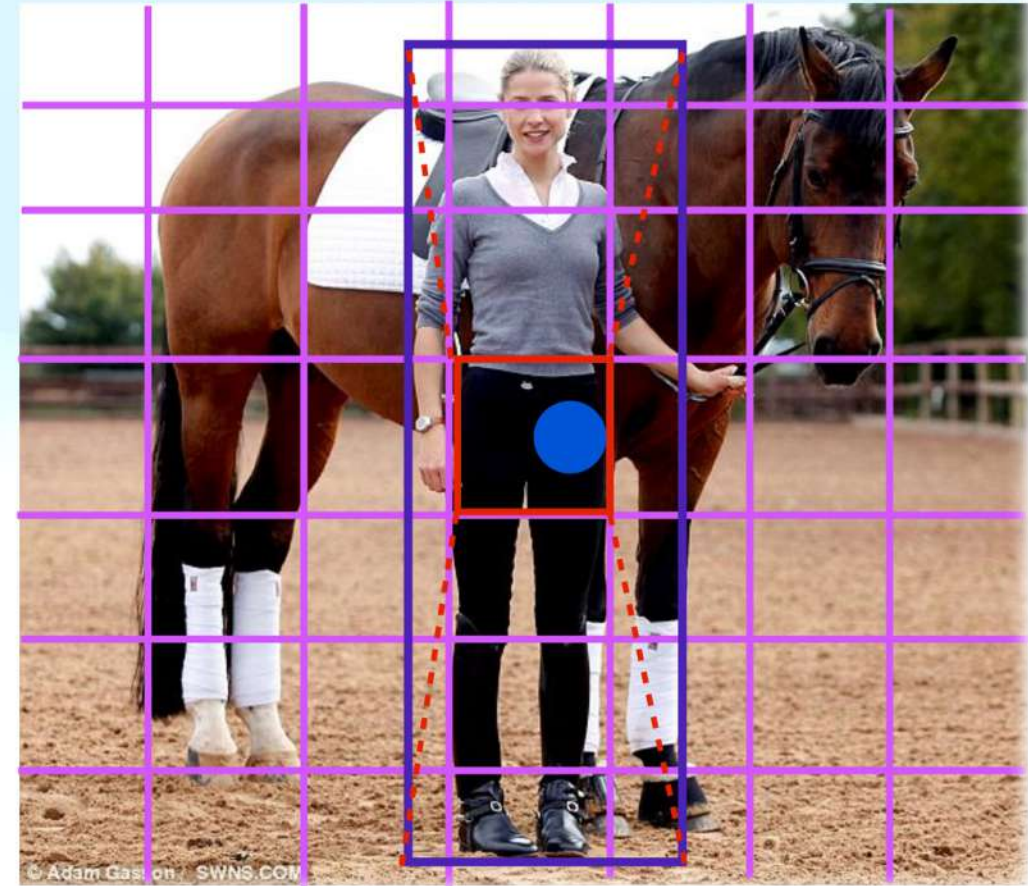
$(x_a, y_a)$ : the coordinate of left-top point

- Width/height (w,h): relative to the whole image

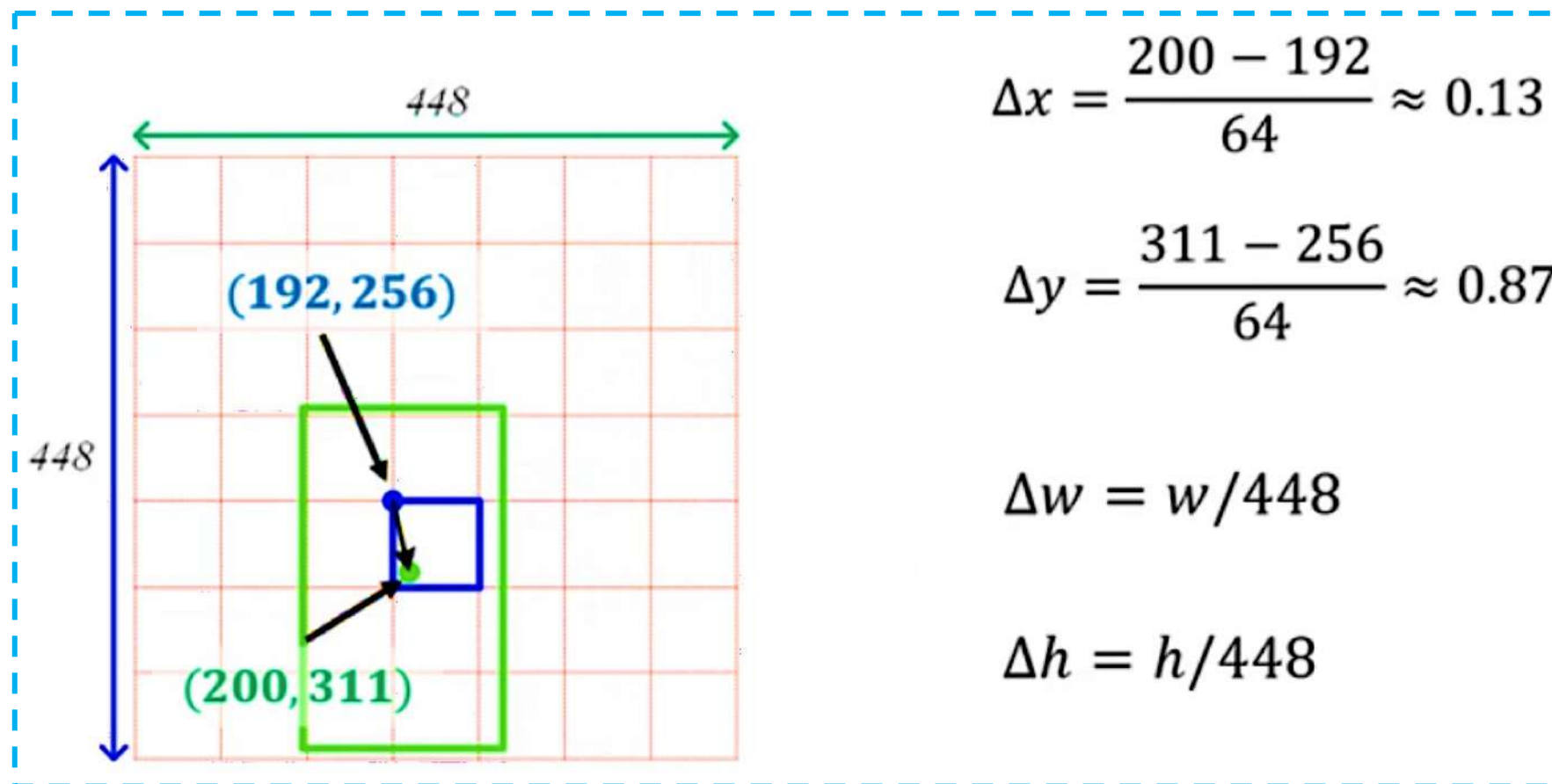
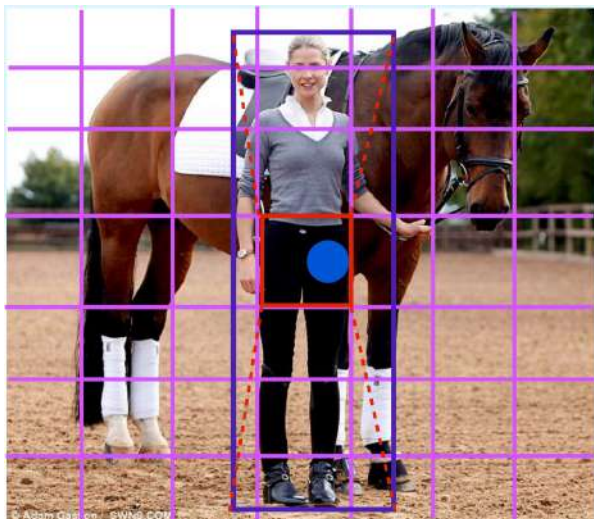
$$\Delta w = w / 448$$

$$\Delta h = h / 448$$

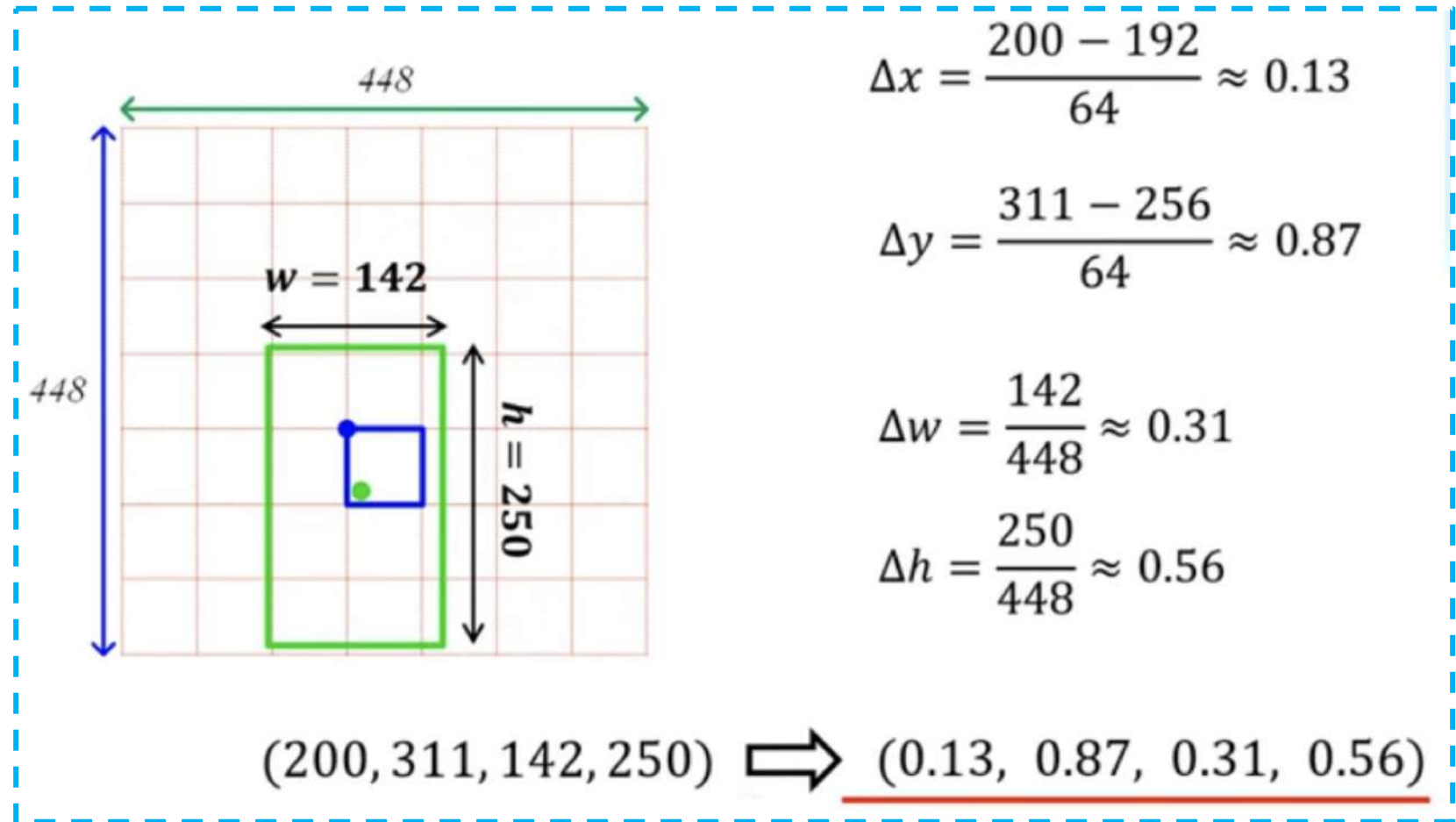
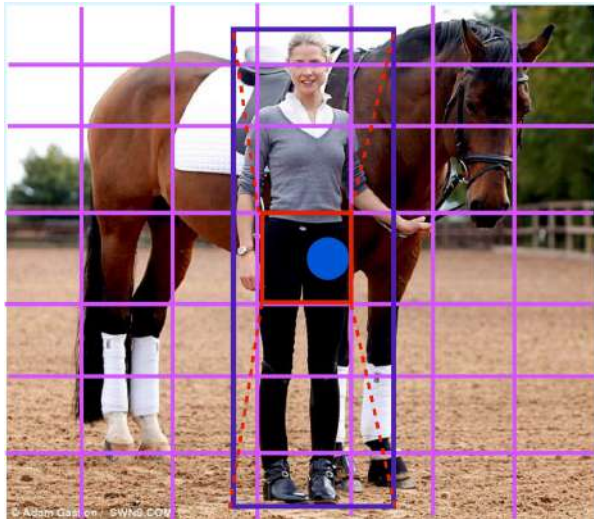
(200, 311, 142, 250)



# Bounding Boxes



# Example Calculation: GT/Target Values





# YOLO: Label Encoding

- For every Grid cell (Anchor box), we need to create targets/labels
- No Object - All zeros
- Object - Relative values w.r.t grid
- Classes - one-hot encoding
- Ex:  $(x,y,w,h,c) = (0.9,0.7,0.1,0.1,1.0)$ 
  - Classes =  $(1.0, 0,0,...0)$  - 20 values

	$(\Delta\hat{x}, \Delta\hat{y}, \Delta\hat{w}, \Delta\hat{h}, \hat{c})$	$(\hat{p}_1, \hat{p}_2, \dots, \hat{p}_{20})$
$A_1$	$(0 \ 0 \ 0 \ 0 \ 0)$	$(0 \ 0 \ \dots \ 0)$
$A_2$	$(0 \ 0 \ 0 \ 0 \ 0)$	$(0 \ 0 \ \dots \ 0)$
	$\dots \dots \dots$	
$A_{11}$	$(0.9 \ 0.7 \ 0.1 \ 0.1 \ 1.0)$	$(0 \ \dots \ 1.0 \ \dots)$
	$\dots \dots \dots$	
$A_{32}$	$(0.1 \ 0.8 \ 0.3 \ 0.5 \ 1.0)$	$(0 \ \dots \ 1.0 \ \dots)$
	$\dots \dots \dots$	

$\hat{p}_{14} = \text{person}$



# YOLO: Motivations



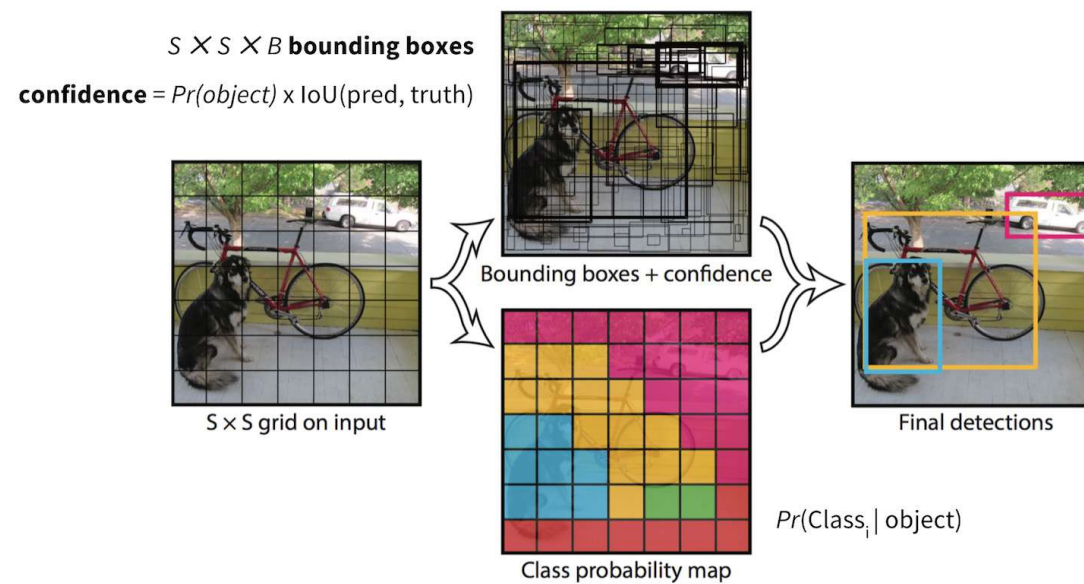
Residual blocks



Bounding box regression



Intersection Over Union (IOU)



# YOLO: Motivations

Image Classification



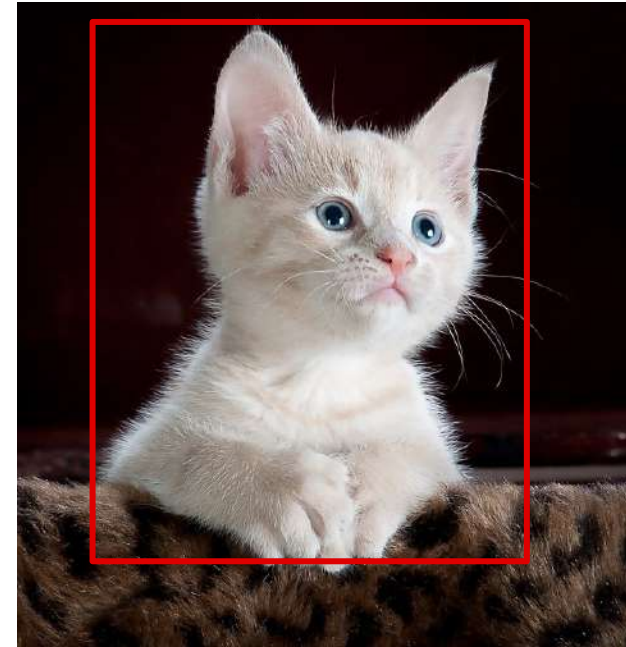
Is this a cat or a dog?

Neural Network Output:

Cat = 1

Dog = 0

Object Localization



Where exactly is the cat in the image?

Neural Network Output:

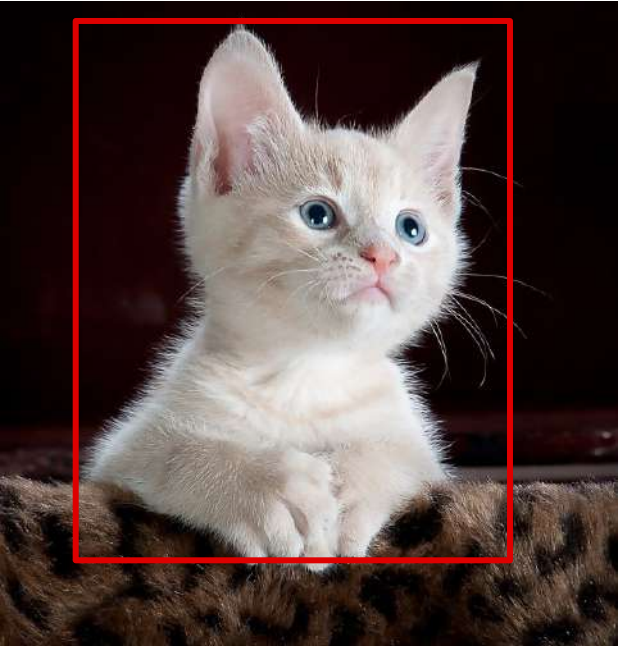
Cat = 1

Dog = 0

Bounding box

# YOLO: Motivations

## Object Localization



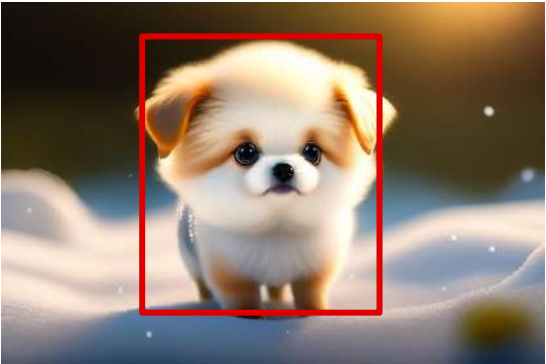
Where exactly is the cat in the image?

$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix} \begin{bmatrix} 1 \\ 50 \\ 70 \\ 60 \\ 70 \\ 1 \\ 0 \end{bmatrix}$

$C_1$ : Cat class  
 $C_2$ : Dog class

Neural Network Output:  
Cat = 1  
Dog = 0

Bounding box



$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix} \begin{bmatrix} 1 \\ 30 \\ 28 \\ 28 \\ 82 \\ 0 \\ 1 \end{bmatrix}$

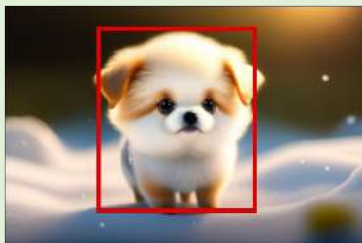
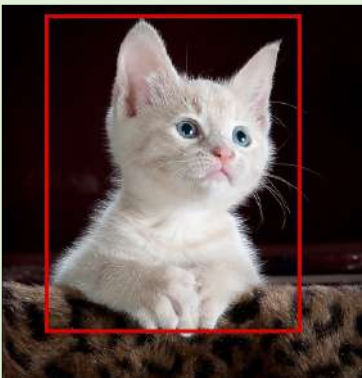


$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

# YOLO: Motivations

X: Train

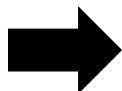
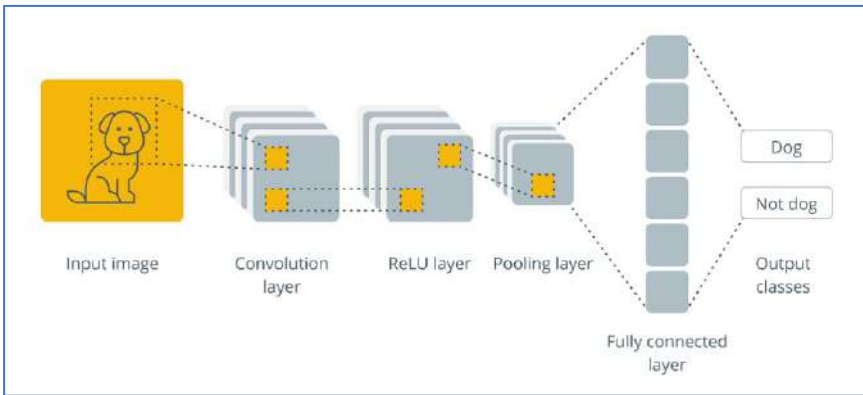
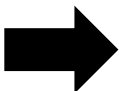
Y: Train



$$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix} \begin{bmatrix} 1 \\ 50 \\ 70 \\ 60 \\ 70 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix} \begin{bmatrix} 1 \\ 30 \\ 28 \\ 28 \\ 82 \\ 0 \\ 1 \end{bmatrix}$$

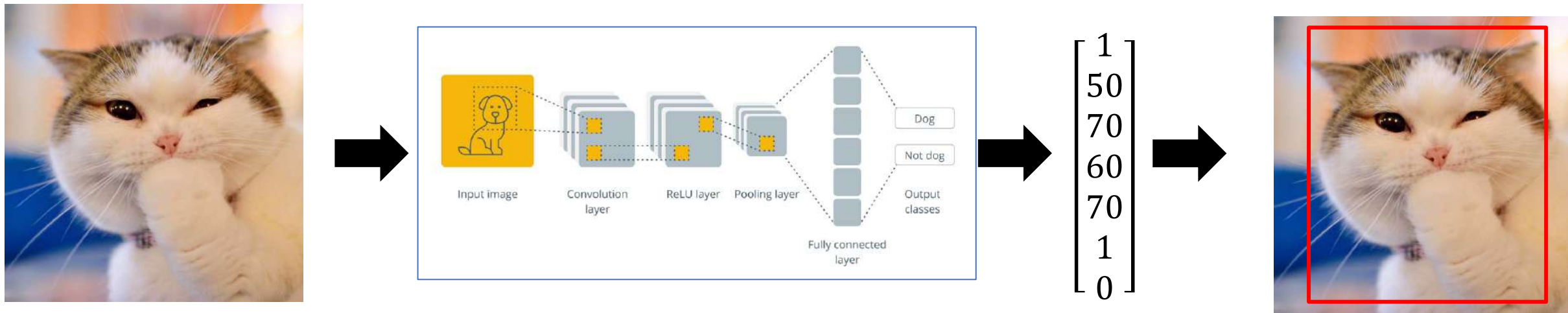
$$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



$$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix}$$



# YOLO: Motivations



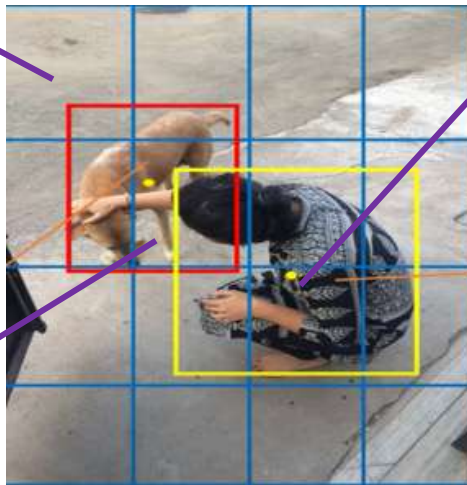
This works OK only for single object. What about multiple objects in an image





# YOLO: Motivations

$$\begin{bmatrix} P_c \\ B_x \\ B_y \\ B_w \\ B_h \\ C_1 \\ C_2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



$$\begin{bmatrix} 1 \\ 0.32 \\ 0.02 \\ 3.0 \\ 2.0 \\ 0 \\ 1 \end{bmatrix}$$

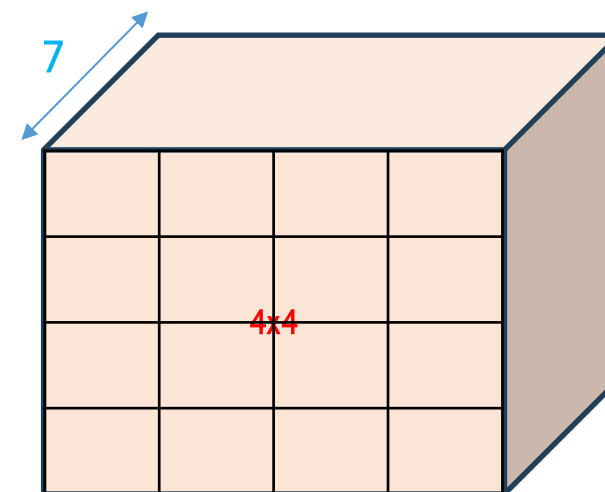
$$\begin{bmatrix} 1 \\ 0.05 \\ 0.3 \\ 2.0 \\ 1.3 \\ 1 \\ 0 \end{bmatrix}$$

(0,0)



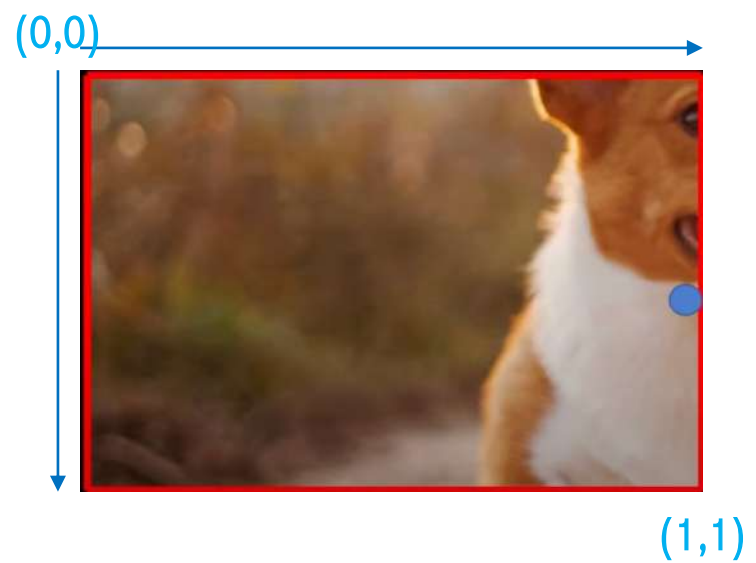
(1,1)

$C_1$ : Dog class  
 $C_2$ : Human class

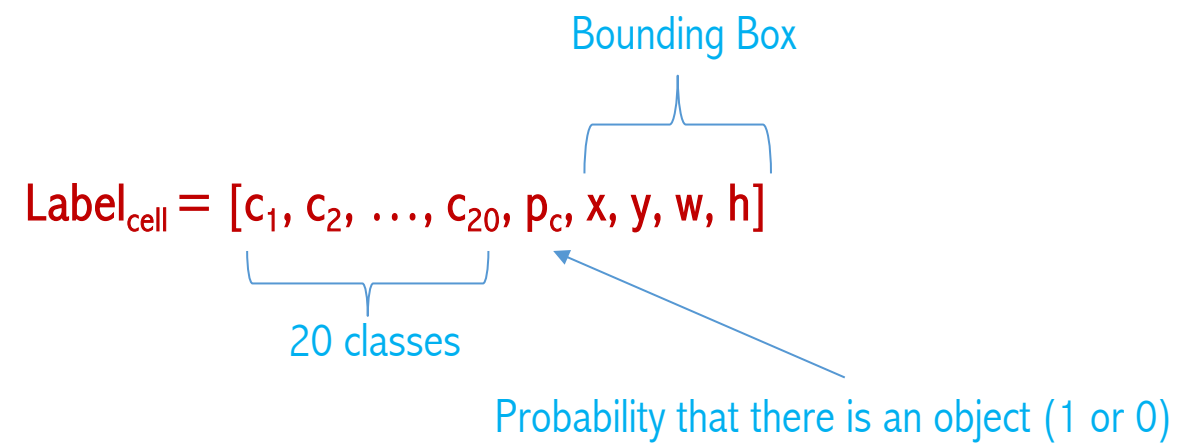


# YOLO: Motivations

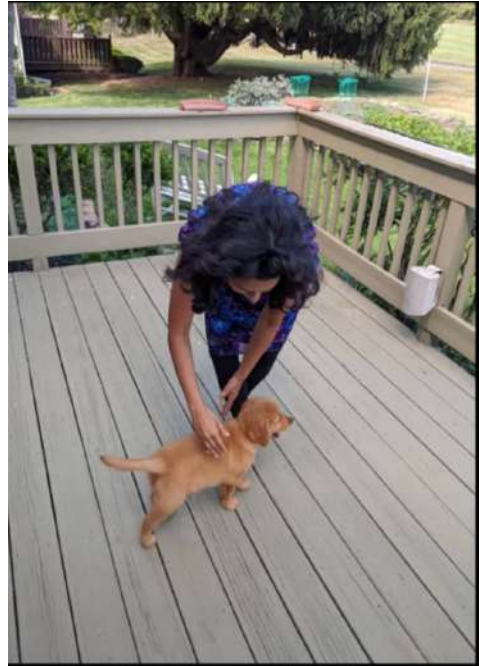
How the labels look like?



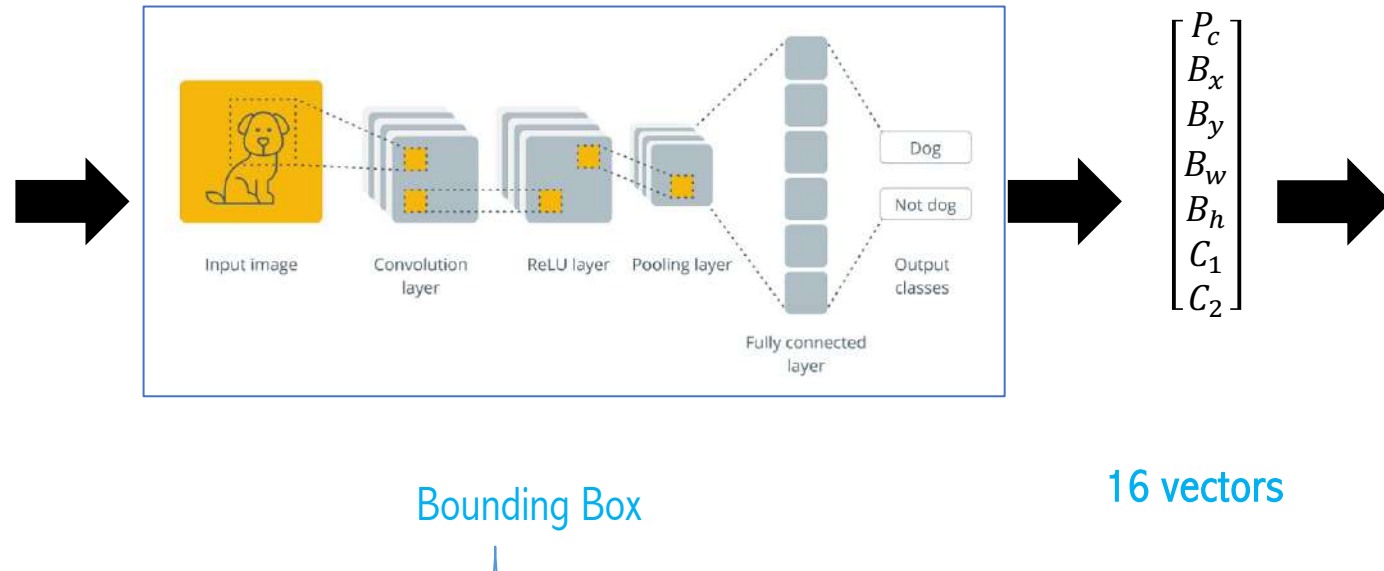
Each output and label will be relative to the cell!  
Each bounding box for each cell will have:  
 $[x,y,w,h] = [0.95, 0.55, 0.5, 1.0]$



# YOLO: Motivations



How the prediction look like?

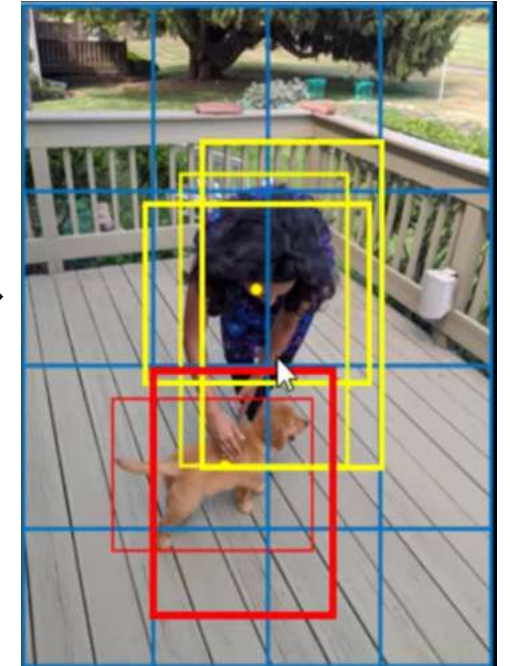


Bounding Box

$$\text{Prediction}_{\text{cell}} = [c_1, c_2, \dots, c_{20}, p_{c1}, x, y, w, h]$$

20 classes

Probability that there is an object (1 or 0)



Multiple bounding boxes

# YOLO: Motivations

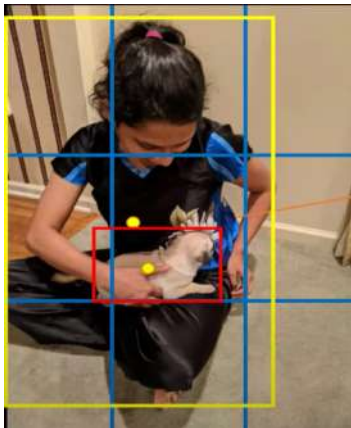
## Non Maximum Suppression



Applying  
NMS



LearnOpenCV.com



What if one grid cell has center of two objects

### Algorithm 1 Non-Max Suppression

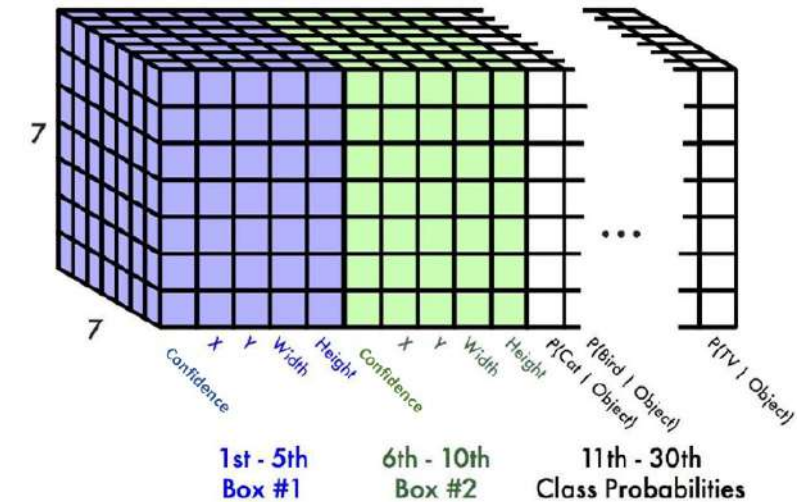
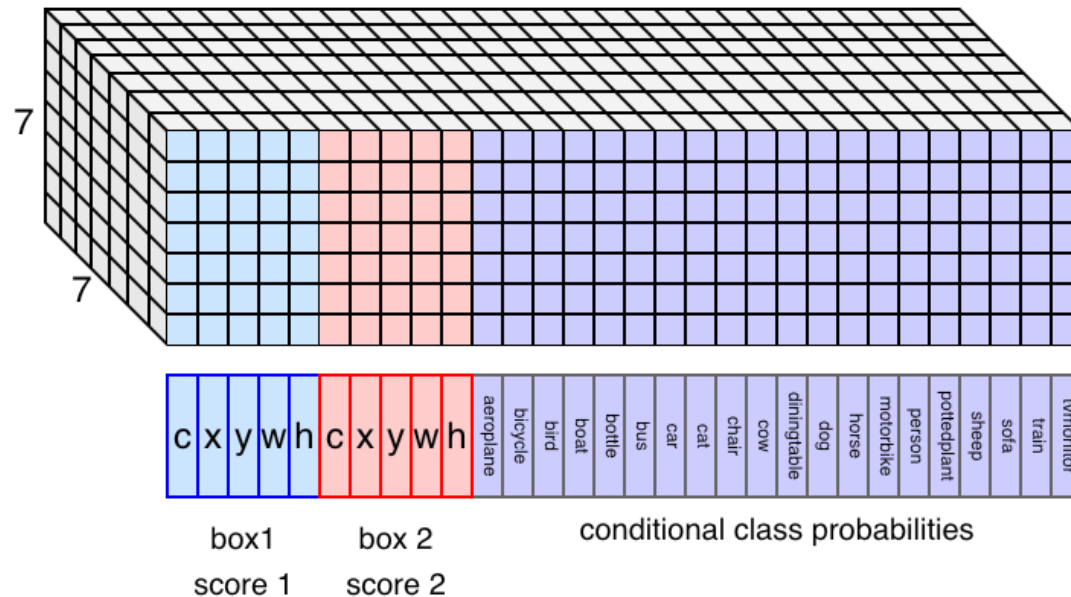
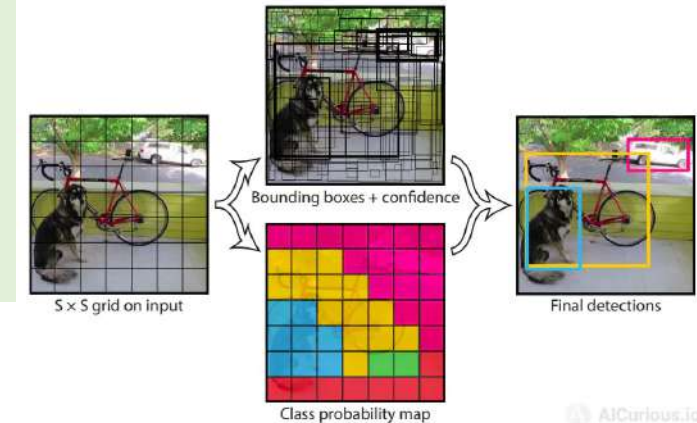
```

1: procedure NMS( $B, c$ )
2:    $B_{nms} \leftarrow \emptyset$  Initialize empty set
3:   for  $b_i \in B$  do  $\Rightarrow$  Iterate over all the boxes
4:      $discard \leftarrow \text{False}$  Take boolean variable and set it as false. This variable indicates whether  $b(i)$  should be kept or discarded
5:     for  $b_j \in B$  do Start another loop to compare with  $b(i)$ 
6:       if  $\text{same}(b_i, b_j) > \lambda_{nms}$  then If both boxes having same IOU
7:         if  $\text{score}(c, b_j) > \text{score}(c, b_i)$  then
8:            $discard \leftarrow \text{True}$  Compare the scores. If score of  $b(i)$  is less than that of  $b(j)$ ,  $b(i)$  should be discarded, so set the flag to True.
9:         if not  $discard$  then Once  $b(i)$  is compared with all other boxes and still the discarded flag is False, then  $b(i)$  should be considered. So add it to the final list.
10:           $B_{nms} \leftarrow B_{nms} \cup b_i$  add it to the final list.
11:   return  $B_{nms}$  Do the same procedure for remaining boxes and return the final list
  
```



# YOLO-v1 Architecture

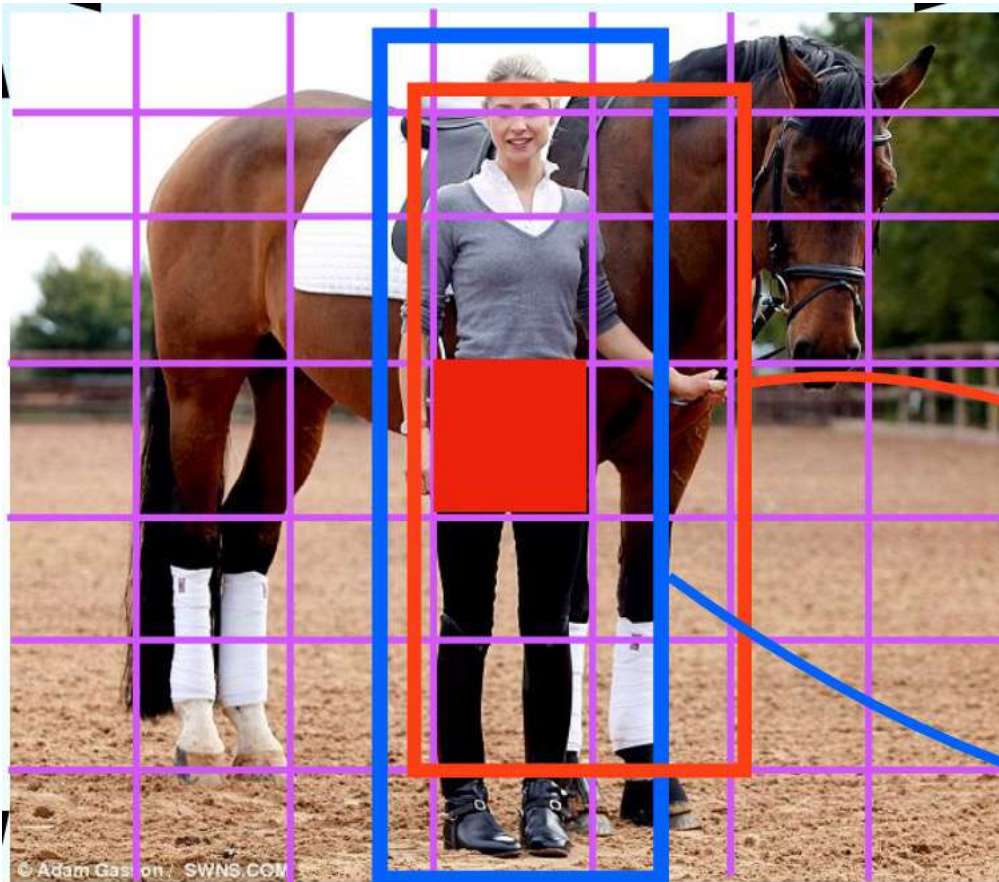
- ❑ The input image is divided into an  $S \times S$  grid ( $S=7$ )
- ❑ Each grid cell predicts  $B$  bounding boxes ( $B=2$ ) and confidence scores for those boxes
- ❑ Each bounding box consists of 5 predictions:  $x, y, w, h$ , and confidence
- ❑ Each grid cell also predicts conditional class probabilities,  $P(\text{Class } i | \text{Object})$ . (Total number of classes=20)



The output size :  $7 \times 7 \times (2 \times 5 + 20) = 1470$



# Output Parsing



Consider the box with highest confidence score per each grid

$(\Delta x_1, \Delta y_1, \Delta w_1, \Delta h_1), \mathbf{c1}$   $(\Delta x_2, \Delta y_2, \Delta w_2, \Delta h_2), \mathbf{c2}$

$$\left\{ \begin{array}{l} x_1 = \Delta x_1 \times 64 + x_a \\ y_1 = \Delta y_1 \times 64 + y_a \\ w_1 = \Delta w_1 \times 448 \\ h_1 = \Delta h_1 \times 448 \end{array} \right. \quad \left\{ \begin{array}{l} x_2 = \Delta x_2 \times 64 + x_a \\ y_2 = \Delta y_2 \times 64 + y_a \\ w_2 = \Delta w_2 \times 448 \\ h_2 = \Delta h_2 \times 448 \end{array} \right.$$

**object class**  $l = \operatorname{argmax}(p_1, p_2, \dots, p_{20})$

**class confidence**  $\hat{c}_1 = c_1 \times p$   $\hat{c}_2 = c_2 \times p$

$p = \max(p_1, p_2, \dots, p_{20})$

# YOLOv1: Summary

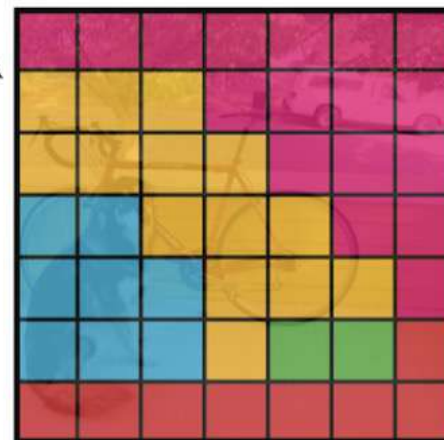
$S \times S \times B$  bounding boxes  
 $\text{confidence} = Pr(\text{object}) \times \text{IoU}(\text{pred}, \text{truth})$



$S \times S$  grid on input

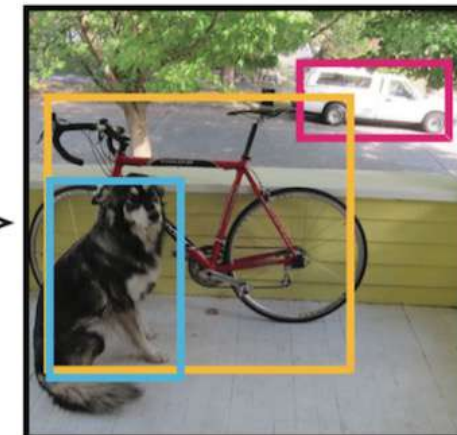


Bounding boxes + confidence



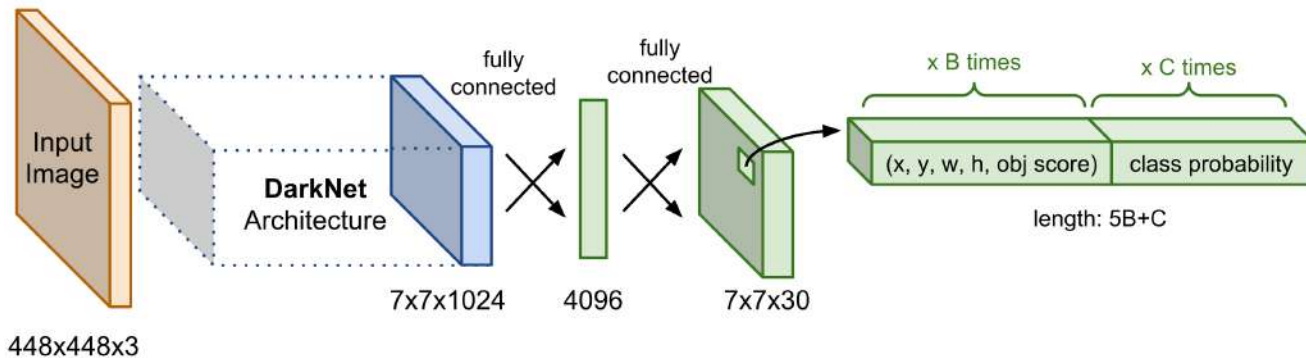
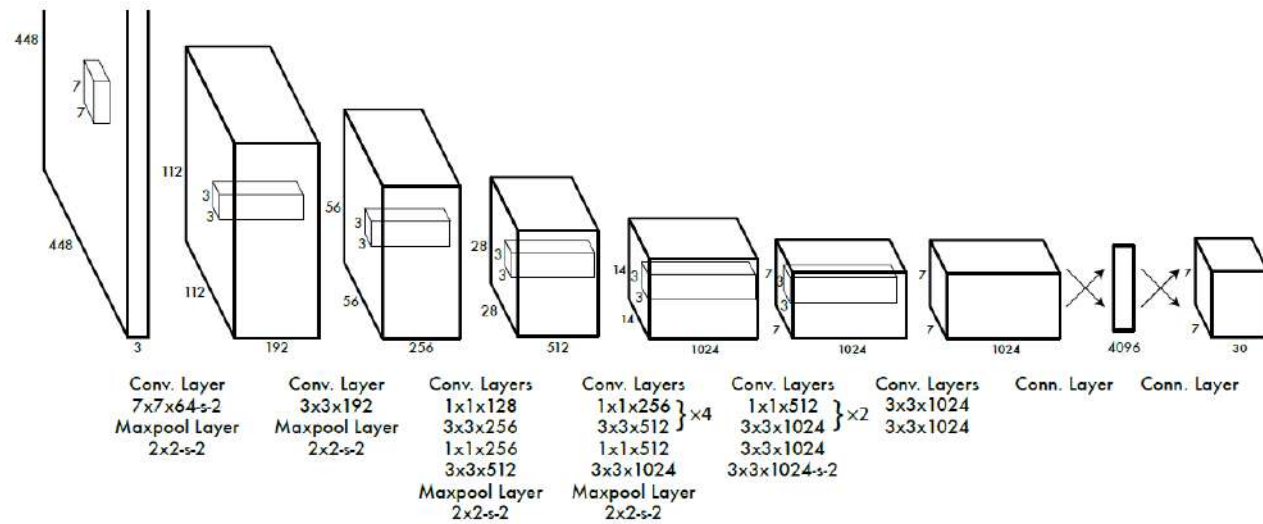
Class probability map

$Pr(\text{Class}_i | \text{object})$



Final detections

# YOLO-v1 Architecture



The YOLO pipeline is simple.

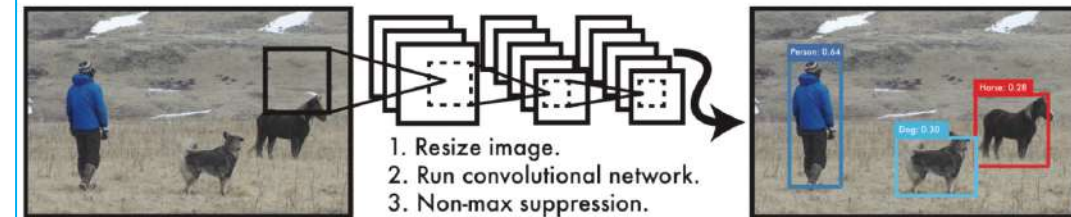


Figure 1 of the [paper](#)

1. It resizes input images to 448 x 448.
2. It runs a single convolutional network on the input images.
3. It thresholds the resulting detections by the model's confidence.

The YOLO is a network was “*inspired by*” [GoogleNet](#). It has 24 [convolutional layers](#) working for feature extractors and 2 dense layers for doing the predictions. This architecture works upon is called Darknet. There is a fast version of YOLO called “Tiny-YOLO” which only has 9 convolution layers



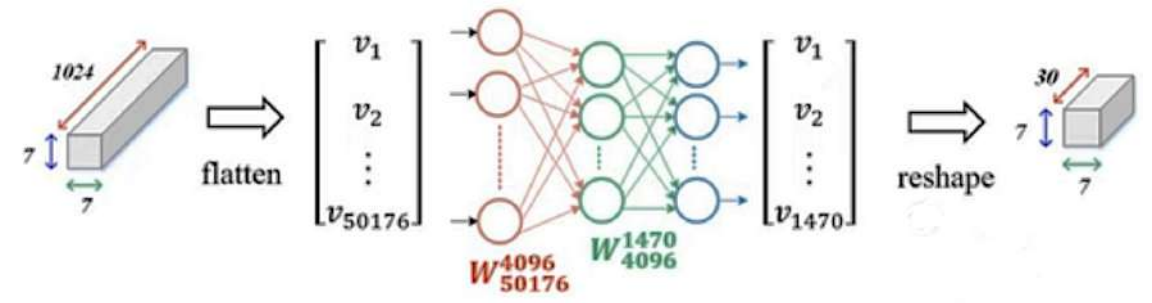
# YOLO-v1 Architecture

Type	Size	Filters	Stride	Output	
Conv.	7 x 7 x 3	64	2	224 x 224 x 64	6
max pool	2 x 2			112 x 112 x 64	
Conv.	3 x 3 x 64	192	1	112 x 112 x 192	
max pool	2 x 2			56 x 56 x 192	
Conv.	1 x 1 x 192	128	1	56 x 56 x 128	
Conv.	3 x 3 x 128	256	1	56 x 56 x 256	
Conv.	1 x 1 x 256	256	1	56 x 56 x 256	4 × 2 = 8
Conv.	3 x 3 x 256	512	1	56 x 56 x 512	
max pool	2 x 2			28 x 28 x 512	2
Conv.	1 x 1 x 512	256	1	28 x 28 x 256	
Conv.	3 x 3 x 256	512	1	28 x 28 x 512	2 × 2 = 4
Conv.	1 x 1 x 512	512	1	28 x 28 x 512	
Conv.	3 x 3 x 512	1024	1	28 x 28 x 1024	4
max pool	2 x 2			14 x 14 x 1024	
Conv.	1 x 1 x 1024	512	1	14 x 14 x 512	
Conv.	3 x 3 x 512	1024	1	14 x 14 x 1024	
Conv.	3 x 3 x 1024	1024	1	14 x 14 x 1024	4
Conv.	3 x 3 x 1024	1024	2	7 x 7 x 1024	
Conv.	3 x 3 x 1024	1024	1	7 x 7 x 1024	
Conv.	3 x 3 x 1024	1024	1	7 x 7 x 1024	

**6 + 8 + 2 + 4 + 4 = 24**

Flatten the last conv map 7x7x1024 to 50176 feature vector

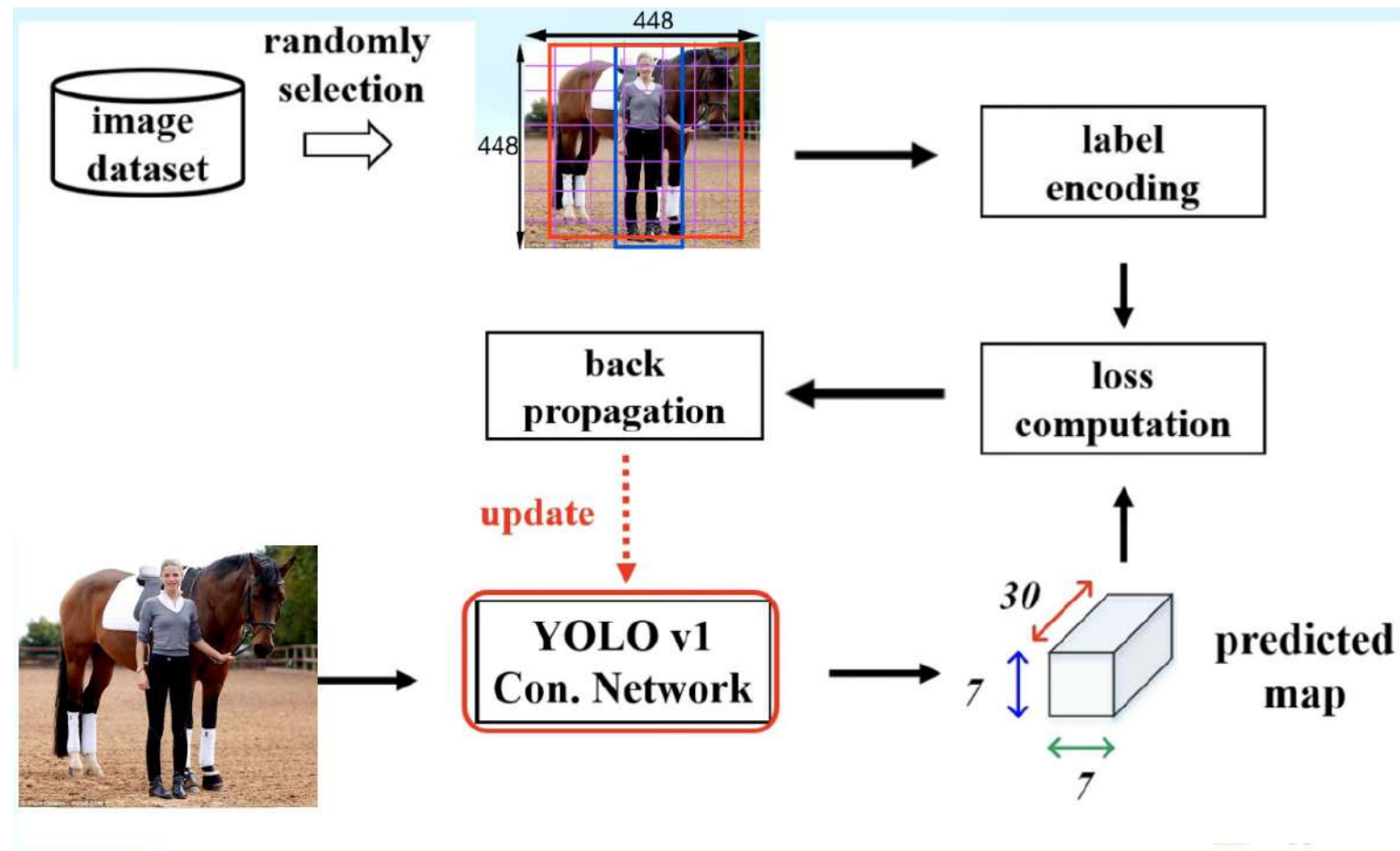
- Pass through 2 fully connected layers
- Output - 1470 feature vector
- Reshape 1470 vector to 7x7x30 feature map



# Training Process

Dataset: Pascal VOC - 20 classes

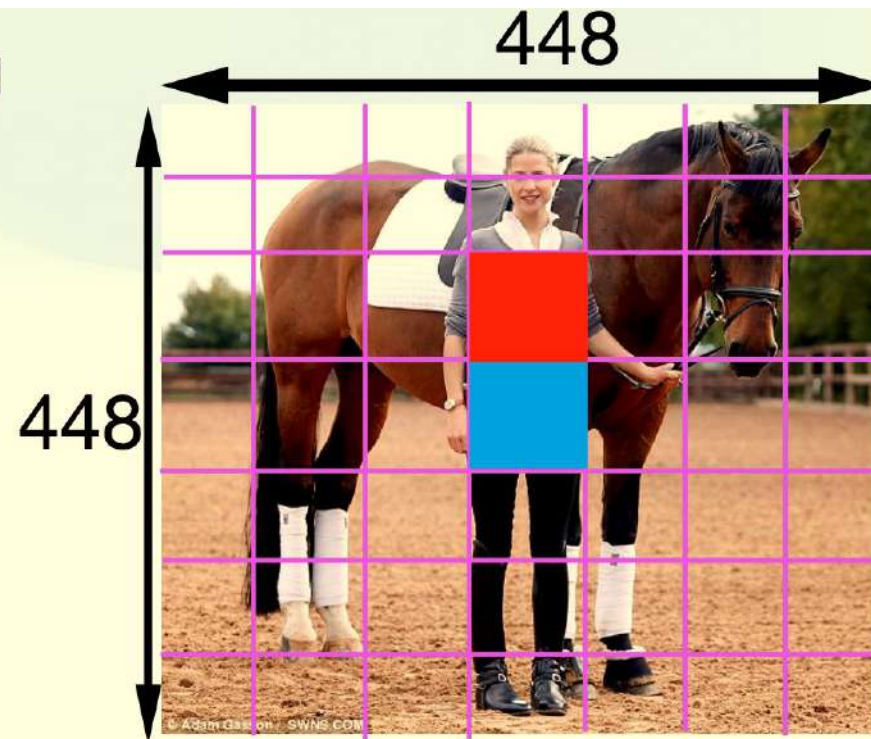
- Network pretrained on Imagenet at 224x224
- Actual training on 448x448 on VOC dataset





# Loss Function

- Loss L is the sum of losses over all grid cells SxS.
- Put more importance on grid cells that contain objects
- Decrease the importance of grid cells having no objects
- Ex: 2 object cells, 47 no-object cells



$$L = \sum_{i=1}^{S^2} L_i$$

$$L = \sum_{i=1}^{S^2} 1_i^{obj} \times L_{i,obj} + \lambda_{no\_obj} \sum_{i=1}^{S^2} 1_i^{no\_obj} \times L_{i,no\_obj}$$

$1_i^{obj} = 1$  if  $i^{th}$  grid is **object anchor**

$1_i^{no\_obj} = 1$  if  $i^{th}$  grid is **no-object anchor**

0.5

# Loss for object cells

Loss = Confidence loss + classification loss + Box Regression loss

- Put more weightage on box parameters

$$L_{i,obj} = \lambda_{coord} \times L_{i,obj}^{box} + L_{i,obj}^{conf} + L_{i,obj}^{cls}$$

$\lambda_{coord} = 5$

$$L_{i,obj}^{box} = (\Delta x_i^* - \Delta \hat{x}_i)^2 + (\Delta y_i^* - \Delta \hat{y}_i)^2 \\ + (\sqrt{\Delta w_i^*} - \sqrt{\Delta \hat{w}_i})^2 + (\sqrt{\Delta h_i^*} - \sqrt{\Delta \hat{h}_i})^2$$

- $(\Delta \hat{x}_i, \Delta \hat{y}_i, \Delta \hat{w}_i, \Delta \hat{h}_i)$ : ground-truth box
- $(\Delta x_i^*, \Delta y_i^*, \Delta w_i^*, \Delta h_i^*)$ : **responsible** predicted box that has the largest IoU with ground-truth box

# Loss for object cells

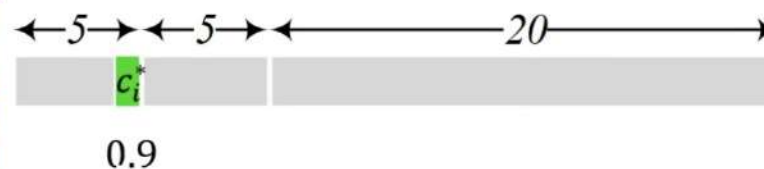
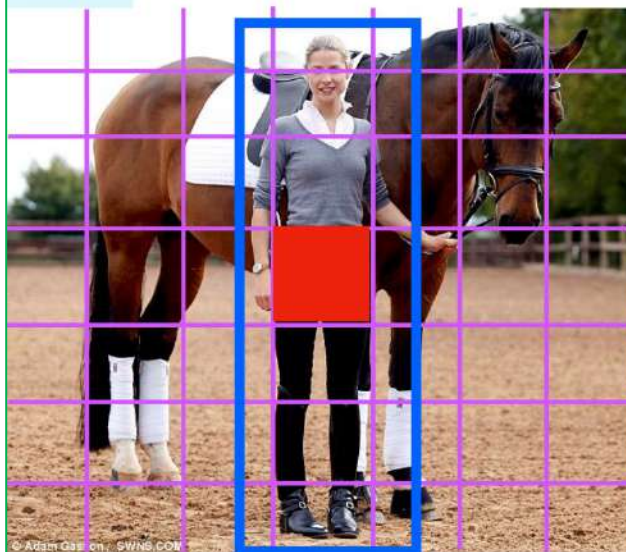
Loss = Confidence loss + classification loss + Box Regression loss

- Put more weightage on box parameters

$$L_{i,obj} = \lambda_{coord} \times L_{i,obj}^{box} + \underline{L_{i,obj}^{conf}} + L_{i,obj}^{cls}$$

= 5

$$L_{i,obj}^{conf} = (c_i^* - \hat{c}_i)^2 = (0.9 - 1.0)^2$$



# Loss for object cells

Loss = Confidence loss + classification loss + Box Regression loss

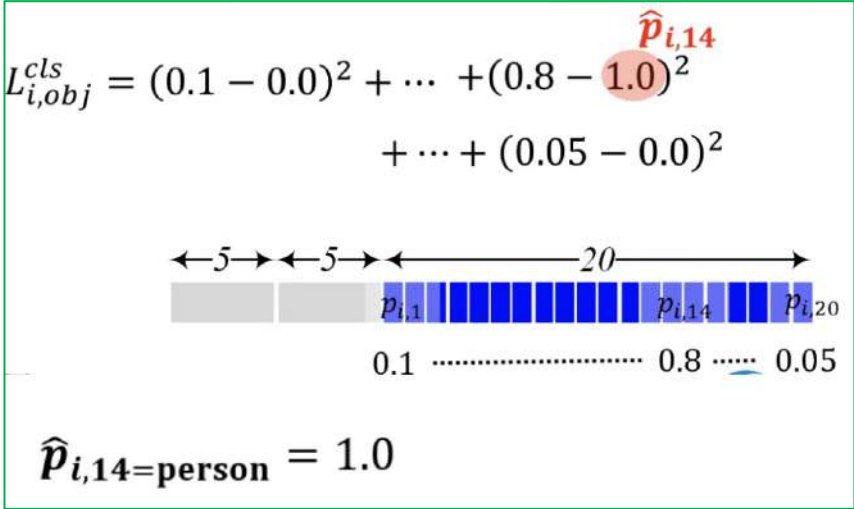
- Put more weightage on box parameters

$$L_{i,obj} = \lambda_{coord} \times L_{i,obj}^{box} + L_{i,obj}^{conf} + \underline{L_{i,obj}^{cls}}$$

$\lambda_{coord} = 5$

$$L_{i,obj}^{cls} = \sum_{c=1}^{20} (p_{i,c} - \hat{p}_{i,c})^2$$

$$L_{i,obj}^{cls} = (p_{i,1} - \hat{p}_{i,1})^2 + \dots + (p_{i,14} - \hat{p}_{i,14})^2 + \dots + (p_{i,20} - \hat{p}_{i,20})^2$$

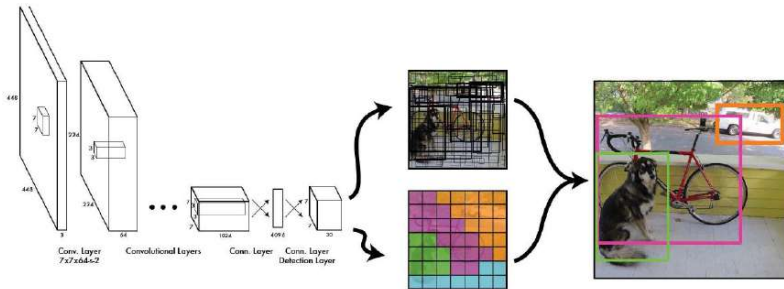


# YOLO-v1 Architecture

Loss Function

Regression  
loss

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$



Confidence  
loss

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\ + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2$$

Classification  
loss

$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

$$\mathcal{L} = \mathcal{L}_{\text{loc}} + \mathcal{L}_{\text{obj}} + \mathcal{L}_{\text{cls}}$$

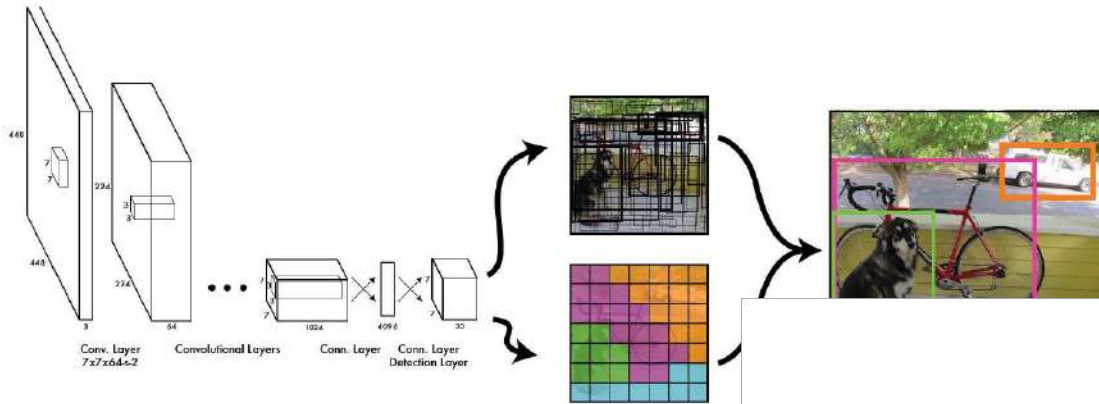
$\mathcal{L}_{\text{loc}}$  - loss function cho dự đoán vị trí bounding box so với ground truth.

$\mathcal{L}_{\text{obj}}$  - loss function cho dự đoán trong cell có object hay không.

$\mathcal{L}_{\text{cls}}$  - loss function cho dự đoán phân phối xác suất cho từng class.



# YOLO-v1 Architecture



$$\text{box confidence score} = Pr(\text{object}) \cdot IOU_{pred}^{truth}$$

$$\text{conditional class probabilities} = Pr(\text{class}_i | \text{object})$$

$$\text{class confidence score} = Pr(\text{class}_i) \cdot IOU_{pred}^{truth}$$

## Loss Function

1 when there is object, 0 when there is no object (in the box j of the cell i)

**Localization loss**

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

Bounding box location (x, y) when there is object

$$+ \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right]$$

Bounding box size (w, h) when there is object

**Confidence loss**

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

Confidence when there is object

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

Confidence when there is no object (in the box j of the cell i)

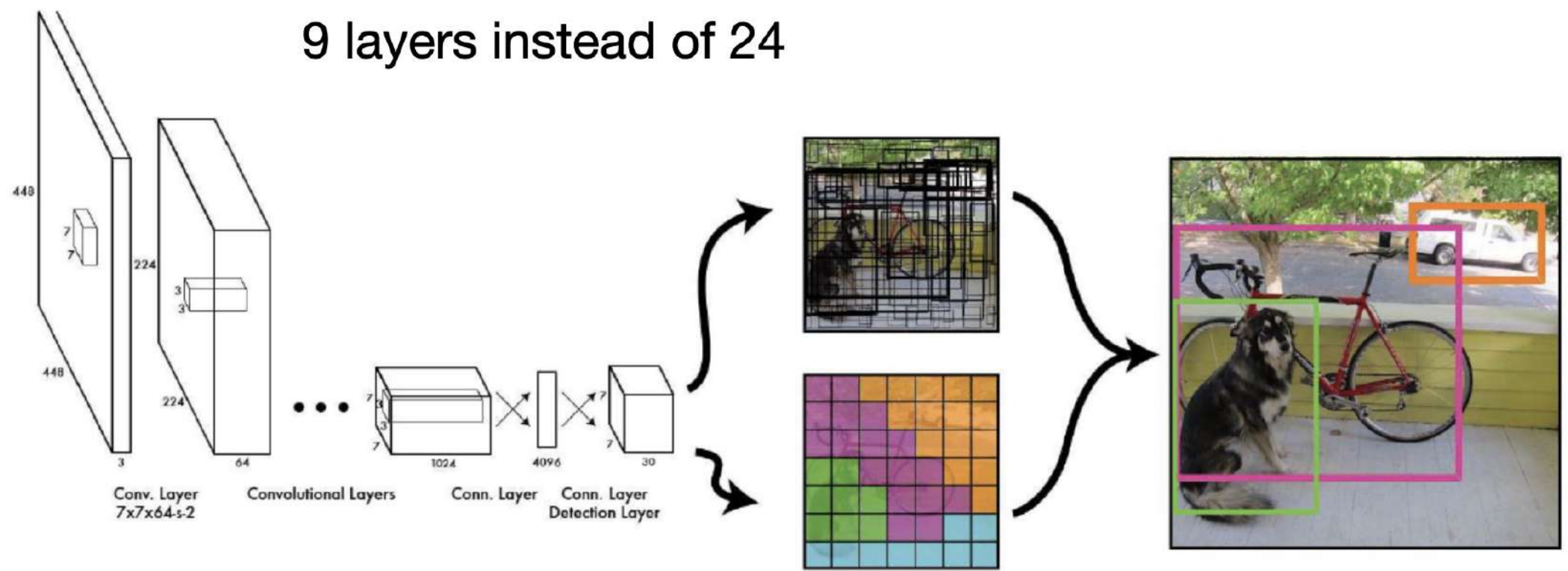
**Classification loss**

$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{obj} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

Class probabilities when there is object (in the cell i)

1 when there is no object, 0 when there is object

# Fast YOLOv1



# YOLO-v1 Limitations



Strong Spatial Constraints: Each cell only predicts two bounding boxes and can only have two class  
YOLO struggles with small objects that appear in groups, like flocks of birds



Relative Coarse Features:

YOLO has many downsampling layers

It struggles to generalize to objects in new or unusual aspect ratios or configurations



Small-Object Localization



# YOLO-v2 Motivations



YOLO v1 was faster than Faster R-CNN, but it was less accurate.



YOLO v1's weakness was the bounding box accuracy. It didn't predict object locations and sizes well, particularly bad at spotting small objects.



SSD, another single-stage detector, broke the record by being better (more accurate) than Faster R-CNN and even faster than YOLO v1.



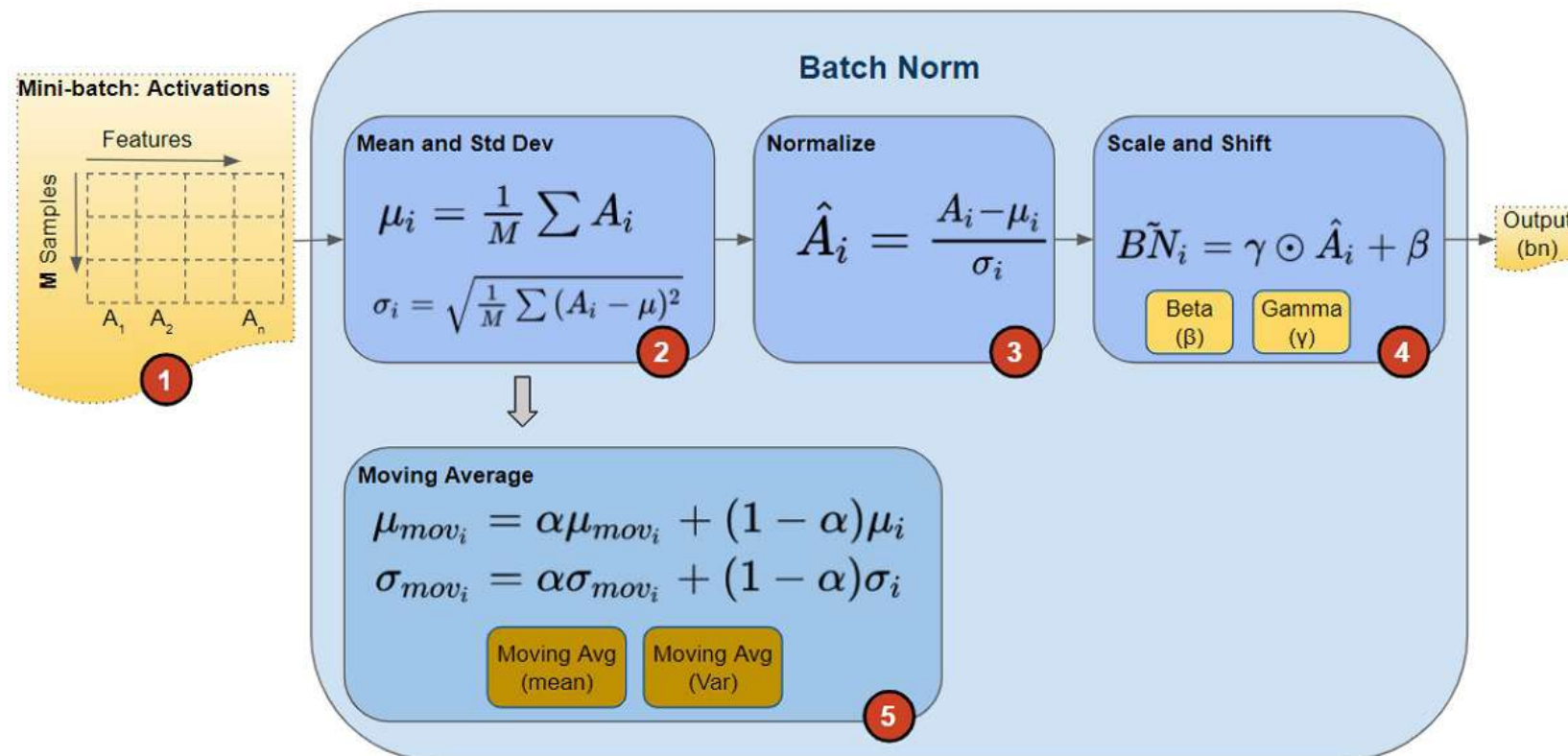
Authors wanted to make their object detector to recognize a wide variety of objects. Pascal VOC object detection dataset contains only 20 classes. They wanted their model to recognize much more classes of objects

# Changes in YOLO-v2

## Batch Normalization

In YOLO v2, they added Batch Normalization to all convolutional layers

- It improved mAP by 2%
- It helped method model to avoid overfitting

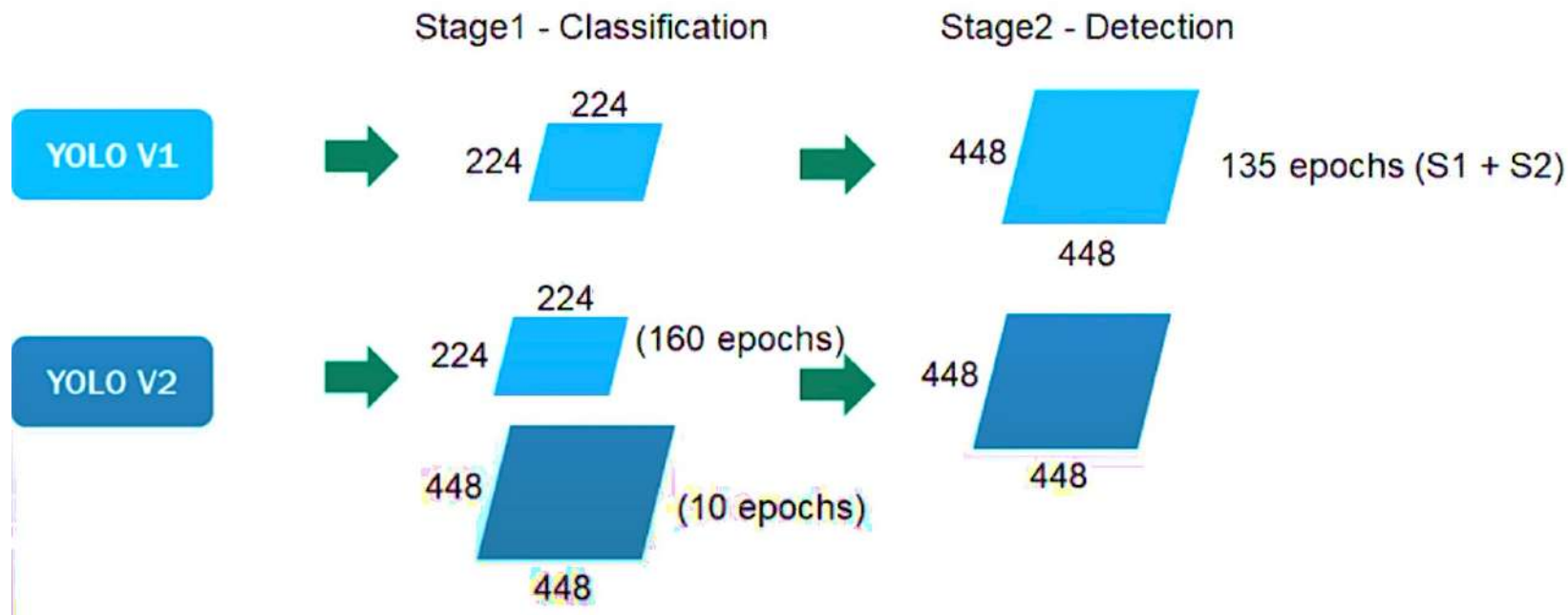




# Changes in YOLO-v2

## High-Resolution Classifier

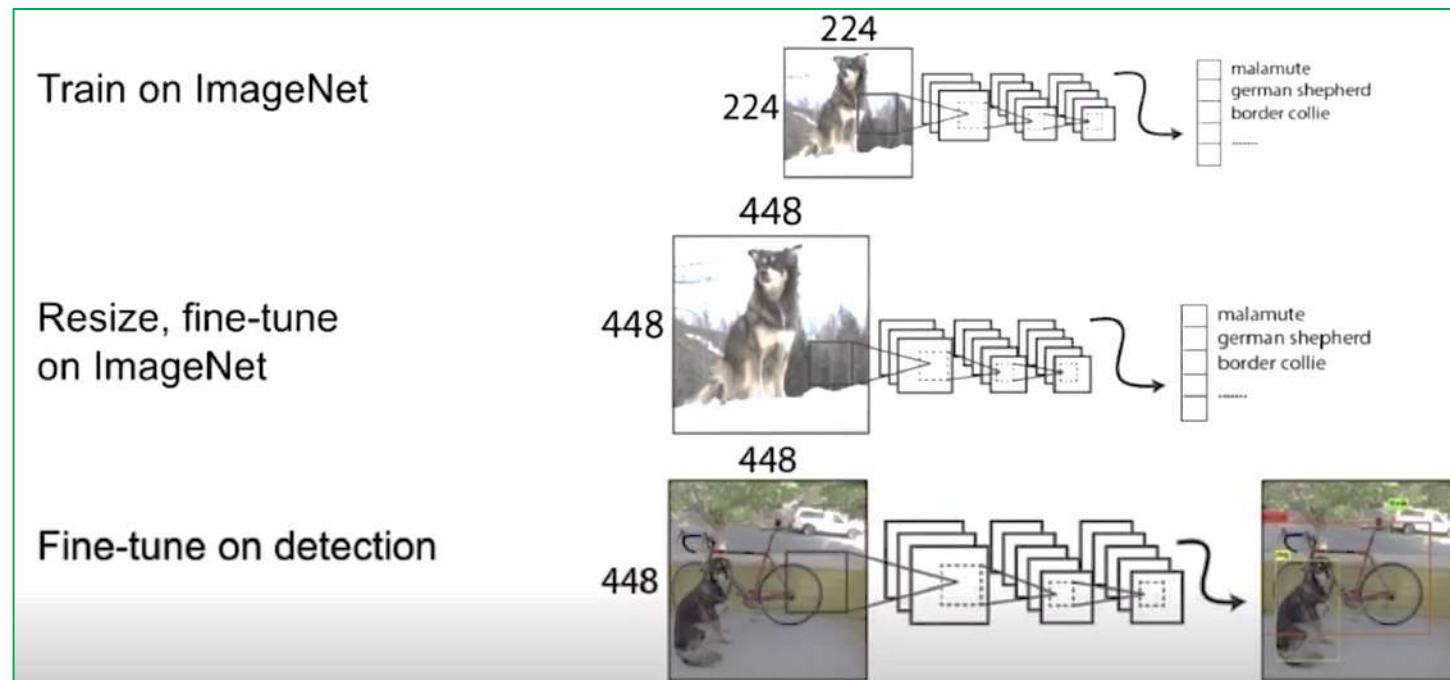
- First, train the classifier on images of size 224 x 224
- Fine-tune the classifier on images of size 448 x 448 for 10 epochs on ImageNet
- Improve mAP by almost 4%



# Changes in YOLO-v2

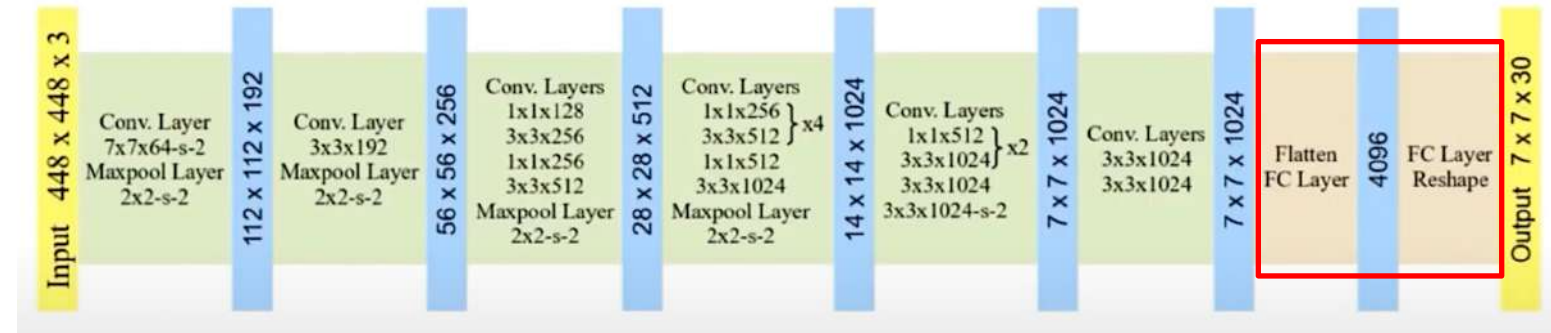
## High-Resolution Classifier

- First, train the classifier on images of size 224 x 224
- Fine-tune the classifier on images of size 448 x 448 for 10 epochs on ImageNet
- Improve mAP by almost 4%

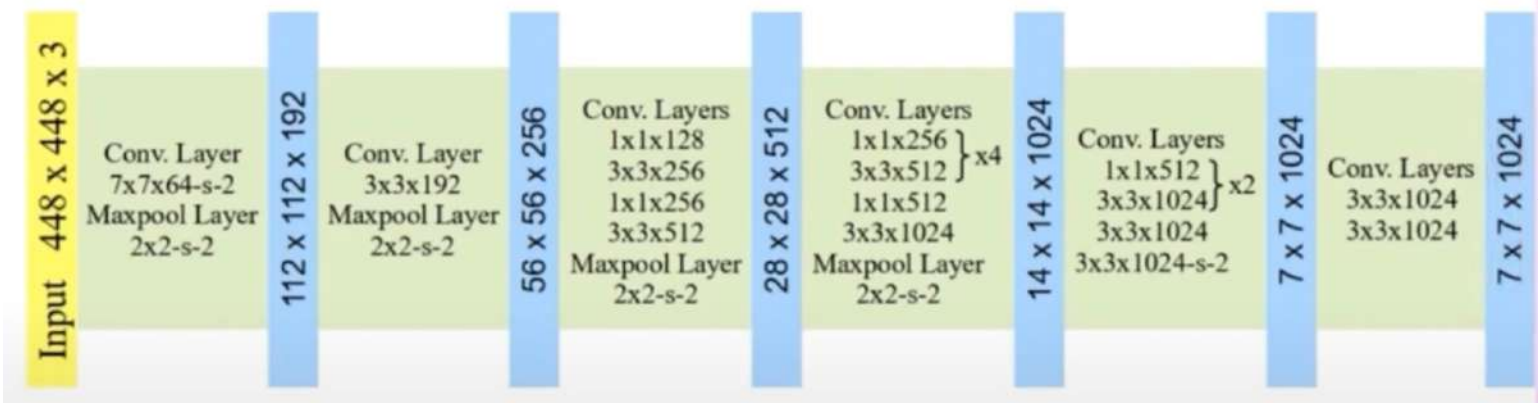


# Changes in YOLO-v2

## YOLOv1

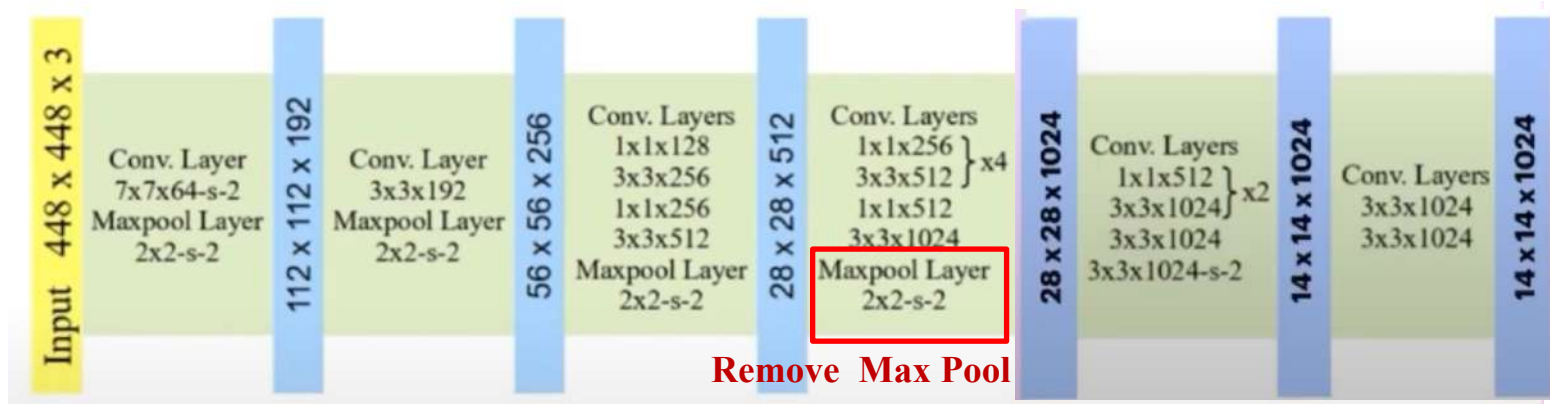


Improve YOLOv1: high resolution feature map



Remove FC

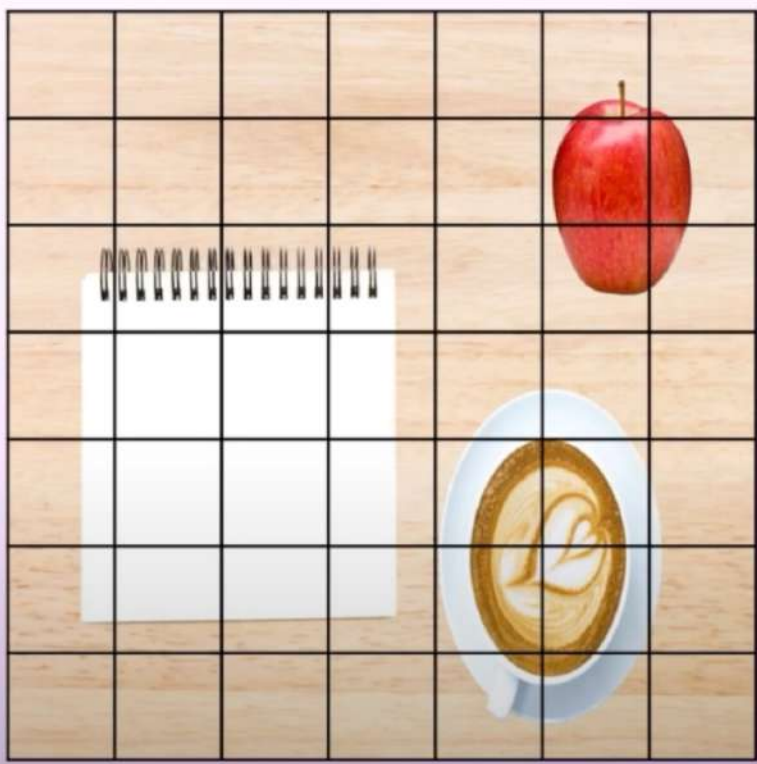
Improve YOLOv1: high resolution feature map



Remove Max Pool

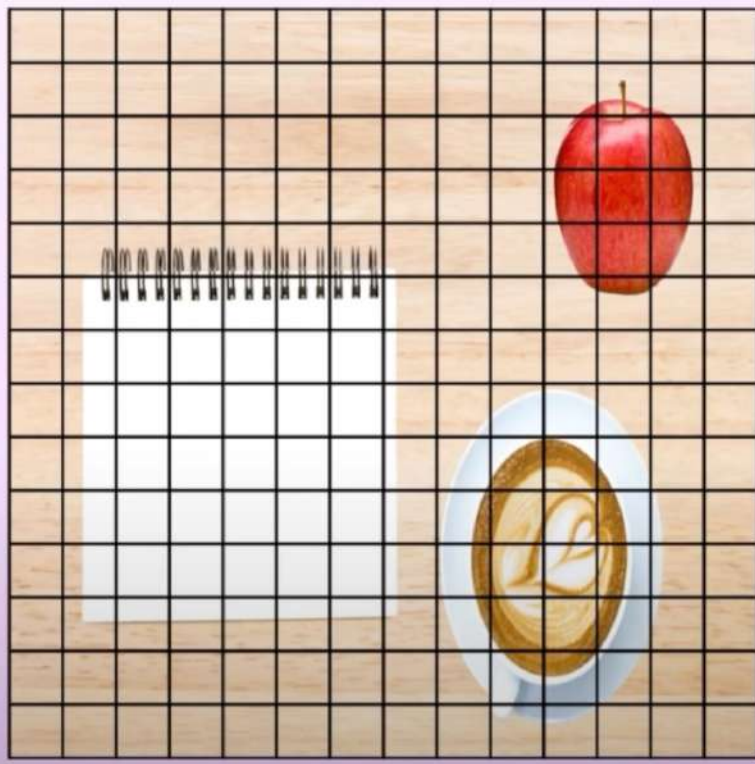
# Changes in YOLO-v2

**YOLO V1**



**7x7**

**YOLO V2**



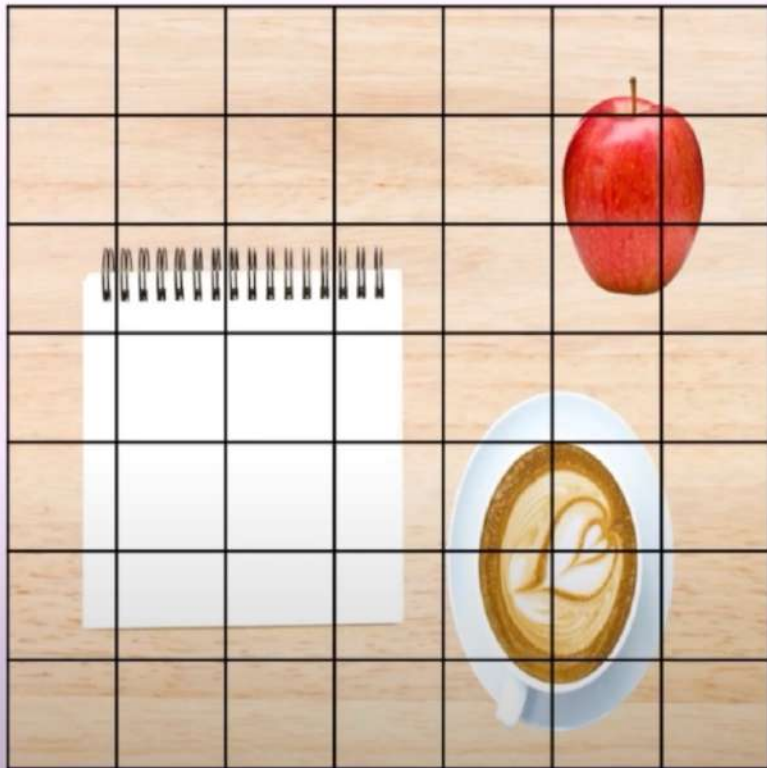
**14x14**

**Problem: If you have even number of grid cells then there is no single center location**



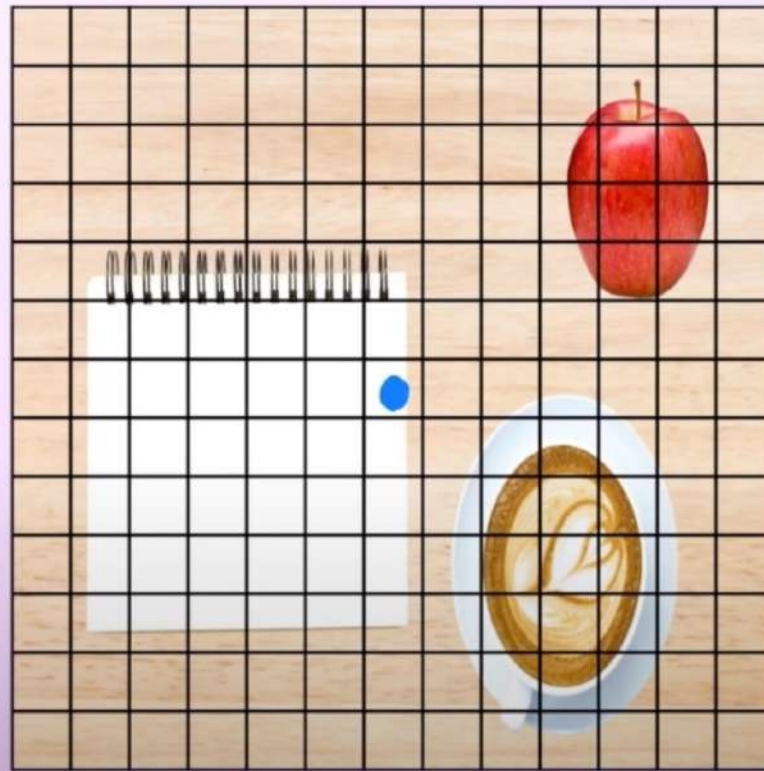
# Changes in YOLO-v2

**YOLO V1**



7x7

**YOLO V2**

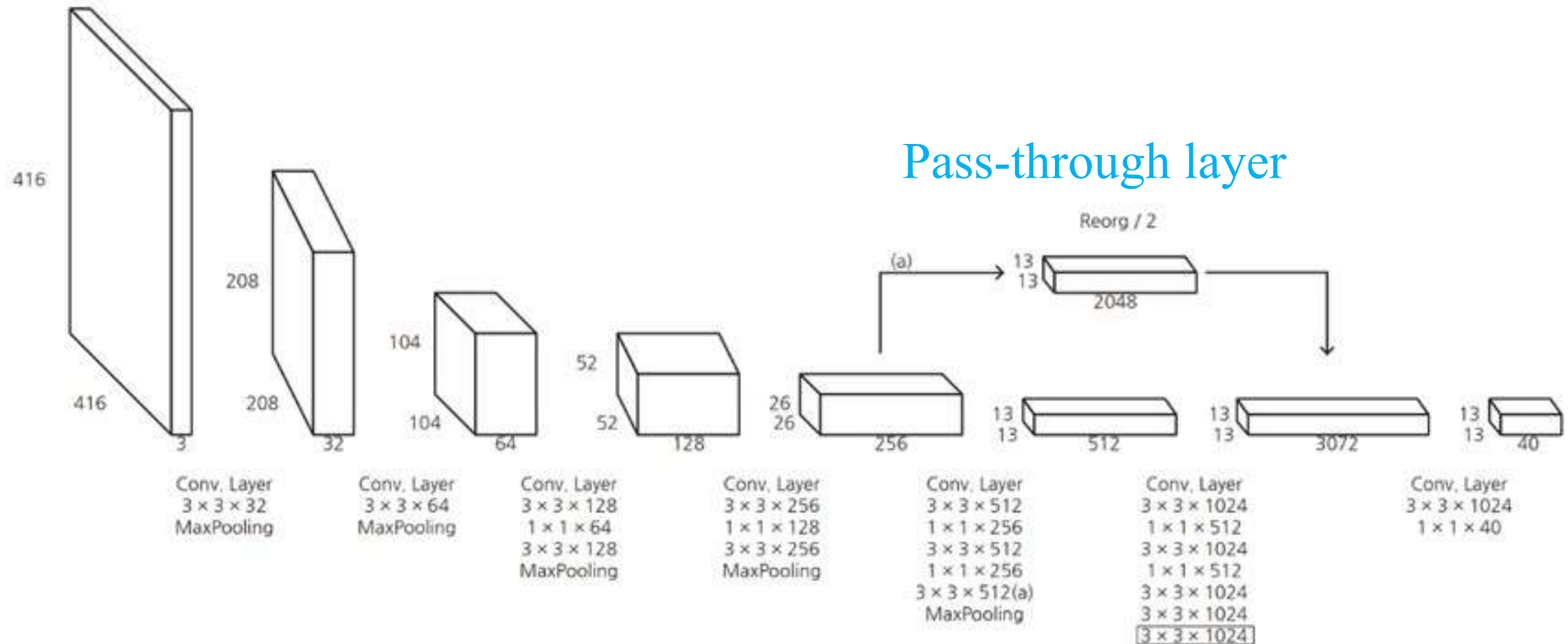


13x13

**Problem: If you have even number of grid cells then there is no single center location**

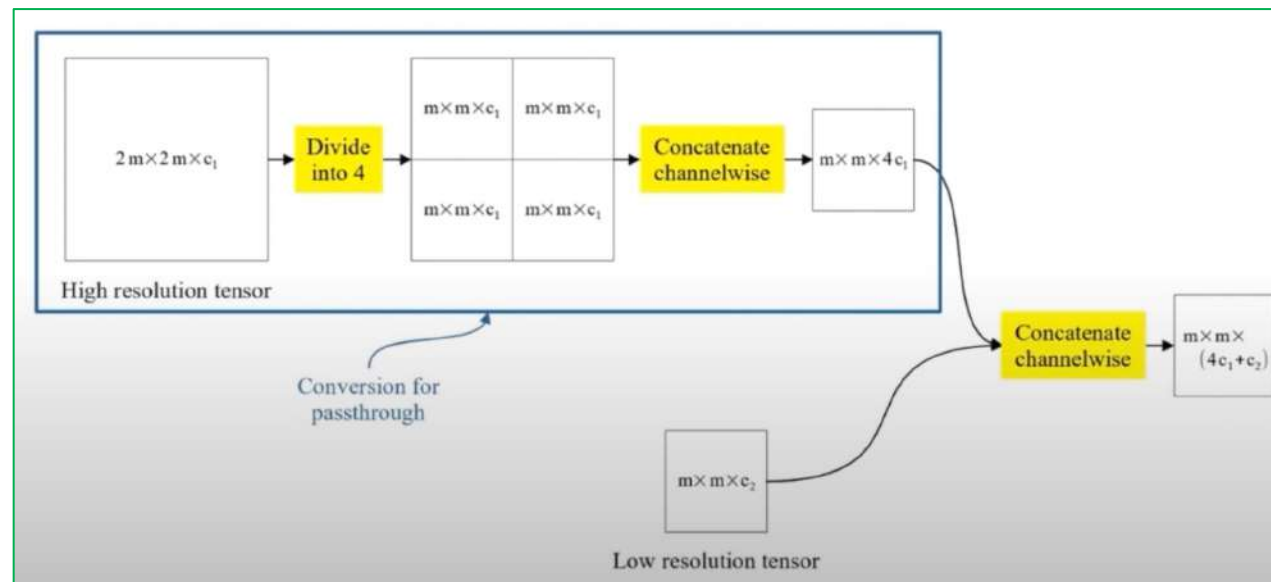
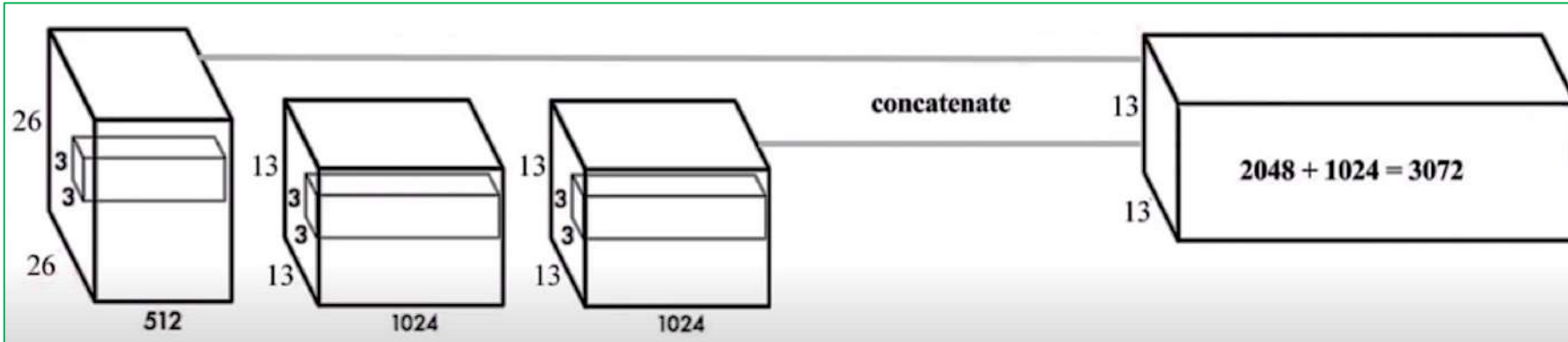


# Changes in YOLO-v2



# Changes in YOLO-v2

## Pass-through layer



# Changes in YOLO-v2

## Convolutional with Anchor Boxes

YOLO v1 suffered from a low recall rate

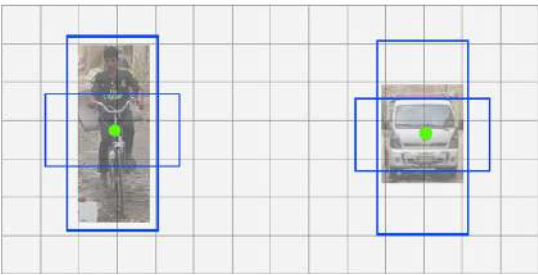
YOLOv2 removes all fully connected layers and uses anchor boxes to predict bounding boxes

YOLO v1 only predicted two bounding boxes per grid cell, which means a total of 98 (= 7 x 7 x 2) bounding boxes per image, much lower than Faster R-CNN.

		Predicted	
		0	1
Actual	0	TN	FP
	1	FN	TP

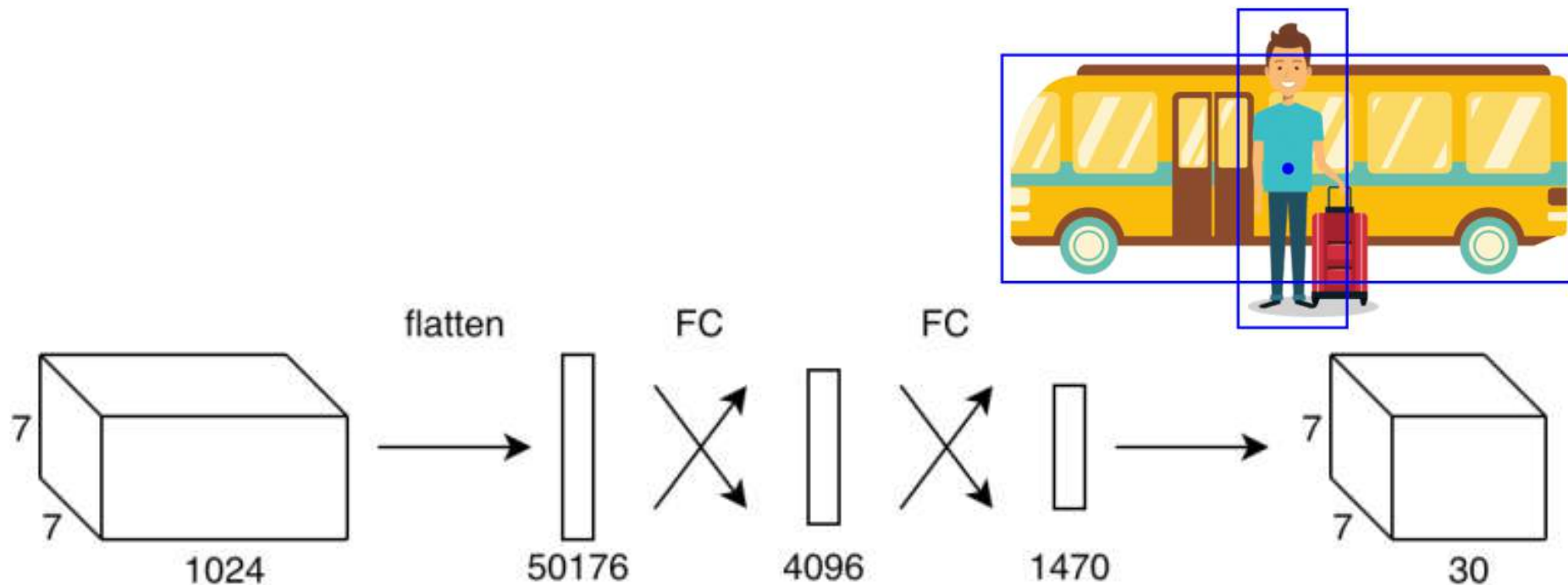
$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$



Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	74.3	46	1	8732	300 × 300
SSD512	76.8	19	1	24564	512 × 512
SSD300	74.3	59	8	8732	300 × 300
SSD512	76.8	22	8	24564	512 × 512

# YOLO-v1 vs. YOLO-v2



Two bounding boxes had to share the same class probabilities. As such, increasing the number of bounding boxes would not benefit much. On the contrary, Faster R-CNN and SSD predicted class probabilities for each bounding box, making it easier to predict multiple classes sharing a similar center location.



# YOLO-v1 vs. YOLO-v2



YOLOv1 was an anchor-free model that predicted the coordinates of B-boxes directly using fully connected layers in each grid cell.



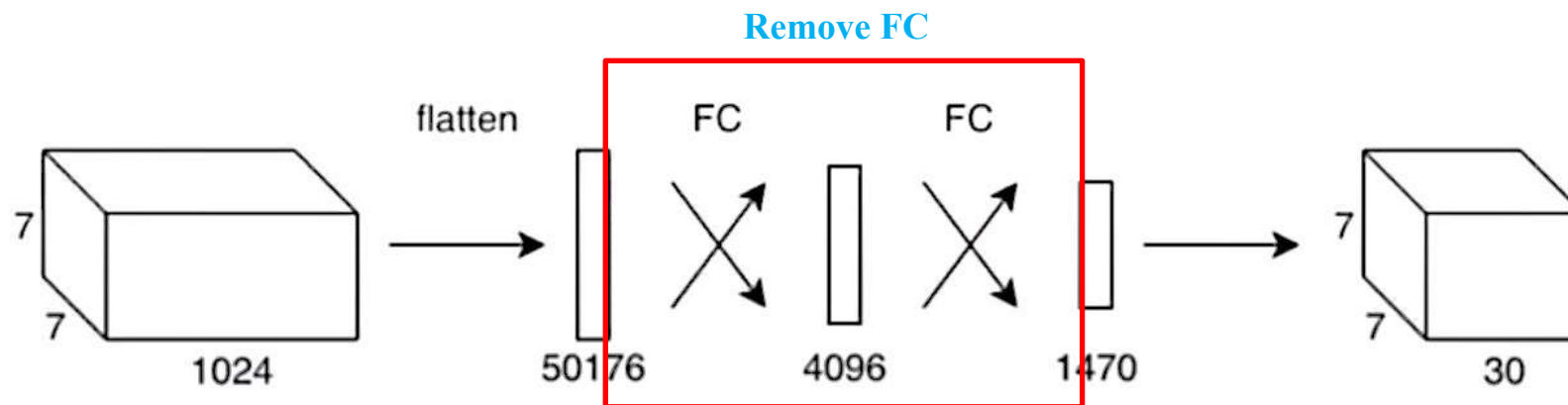
Inspired by Faster-RCNN that predicts B-boxes using hand-picked priors known as anchor boxes, YOLOv2 also works on the same principle.



Unlike YOLOv1, wherein each grid cell, the model predicted one set of class probabilities per grid cell, ignoring the number of boxes  $B$ , YOLOv2 predicted class and objectness for every anchor box.



# Problem with Bounding Boxes

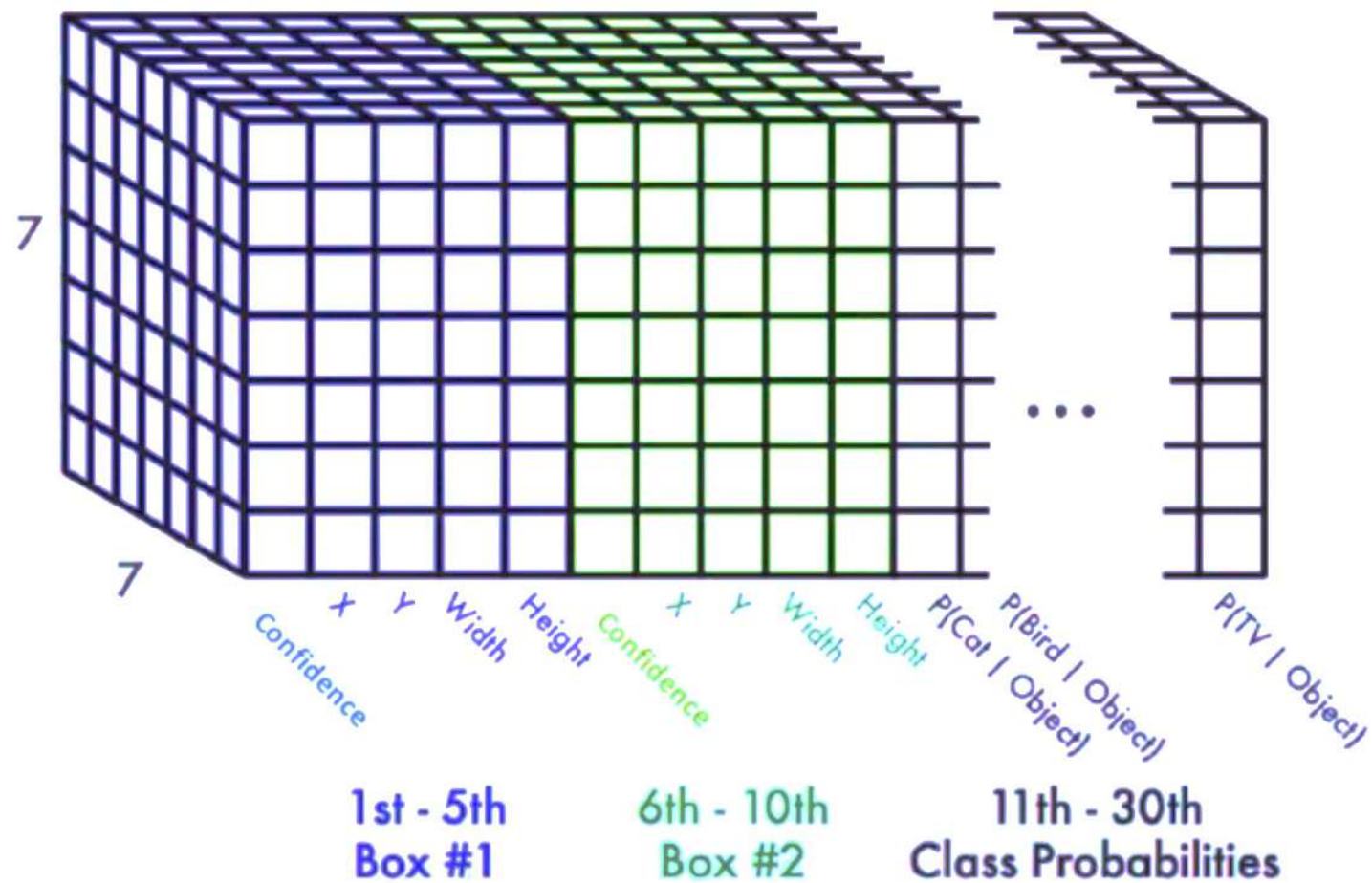


- With 13x13 grid cells -  $13 \times 13 \times 2 = 338$
- Should we increase the grids to 50x50?

Problem with fully connected layer  $\Rightarrow$  YOLOv2 is fully convolutional network

# Problem with Bounding Boxes

- 1 class per grid cell
- Limits the number of objects detected
- Solution: class prediction per box



# Reason for Poor Localization

- Boxes are learnt relative to grid cell
- Objects can be of different shapes



Anchor box is a solution

## Dimension Clusters

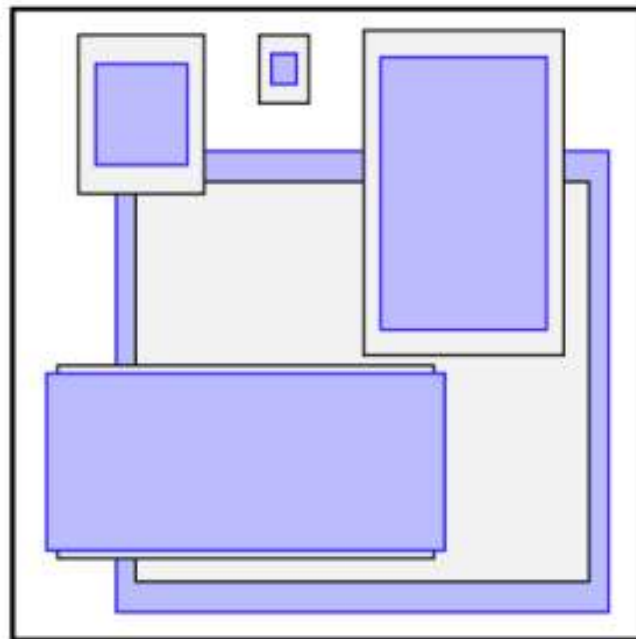
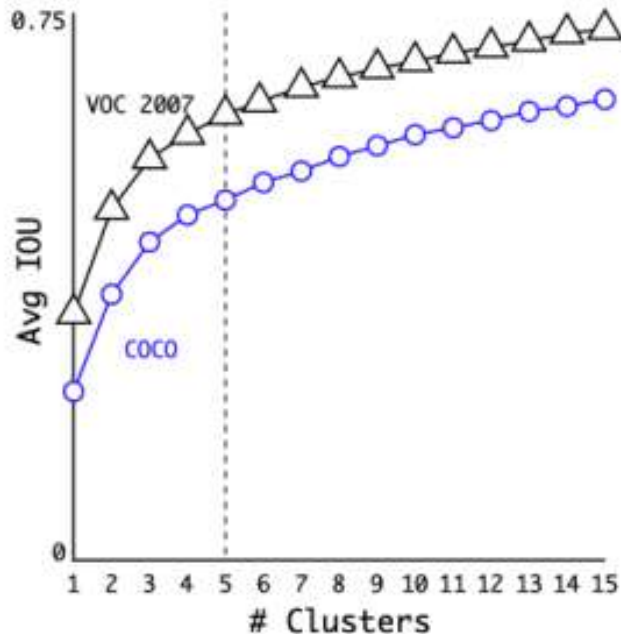
Dimension Clusters Unlike Faster-RCNN, which used hand-picked anchor boxes, YOLOv2 used a smart technique to find anchor boxes for the PASCAL VOC and MS COCO datasets.

Redmon and Farhadi thought that instead of using hand-picked anchor boxes, we pick better priors that reflect the data more closely. It would be a great starting point for the network, and it would become much easier for the network to predict the detections and optimize faster.

# Changes in YOLO-v2

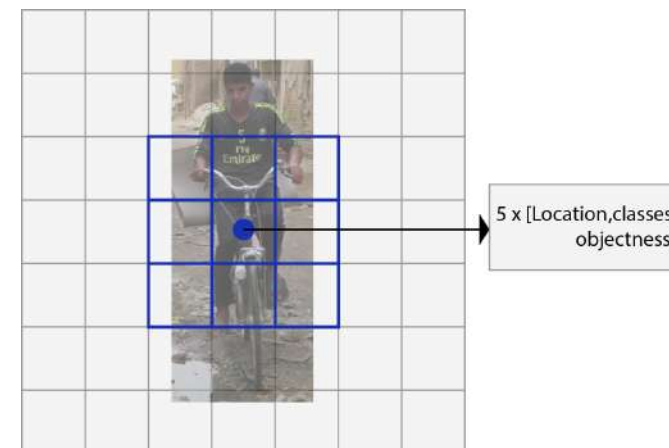
## Dimension Clusters

Using k-means clustering on the training set bounding boxes to find good anchor boxes or priors.



a) They picked the distance function as follows:  $d(\text{box}, \text{centroid}) = 1 - \text{IOU}(\text{box}, \text{centroid})$ .

b) They ran K-Means with a various value of  $k$  and found out that  $k=5$  gives a good tradeoff between between model complexity and high recall.

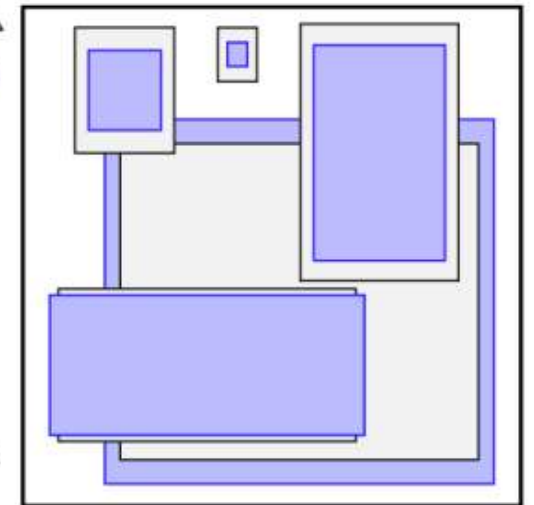
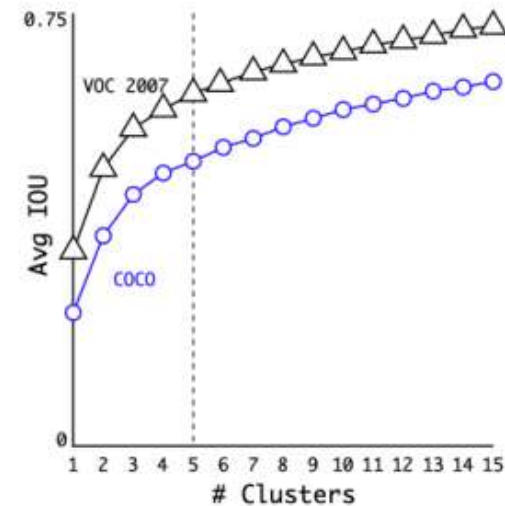
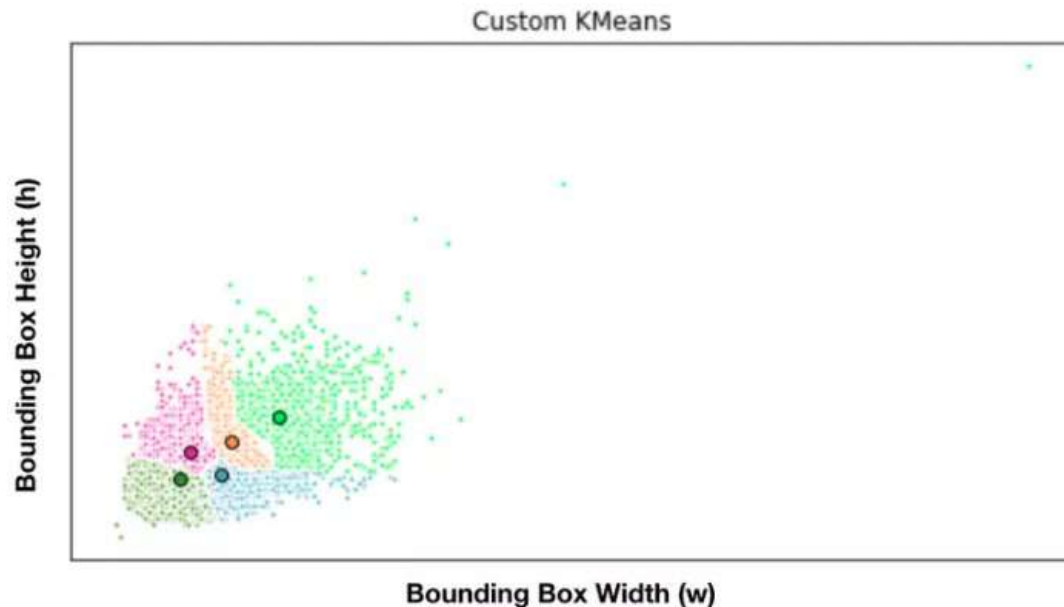




# Changes in YOLO-v2

## Dimension Clusters

Using k-means clustering on the training set bounding boxes to find good anchor boxes or priors.



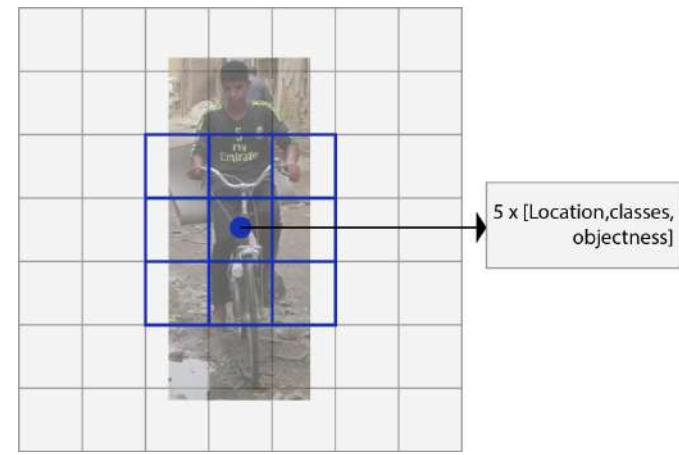
# Changes in YOLO-v2

## Dimension Clusters

Using k-means clustering on the training set bounding boxes to find good anchor boxes or priors.

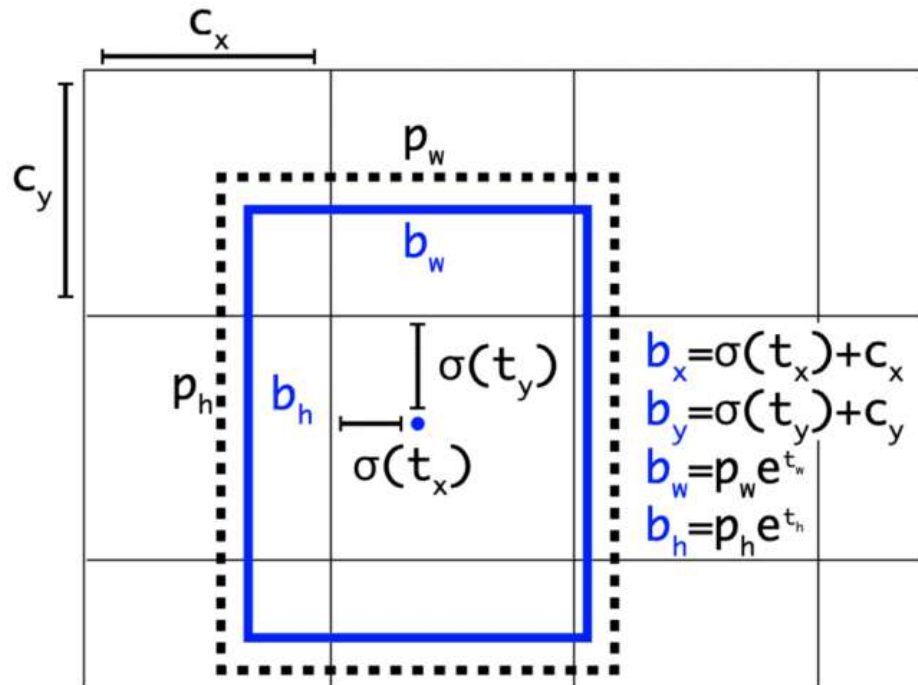
Box Generation	Number of Anchors	Average IOU
Cluster Sum-Squared Distance	5	58.7
Cluster IOU	5	61.0
Anchor Boxes	9	60.9
Cluster IOU	9	67.2

- a) They picked the distance function as follows:  $d(\text{box}, \text{centroid}) = 1 - \text{IOU}(\text{box}, \text{centroid})$ .
- b) They ran K-Means with a various value of  $k$  and found out that  $k=5$  gives a good tradeoff between between model complexity and high recall.



# Changes in YOLO-v2

## Direct Location Prediction vs. Offset location predictions



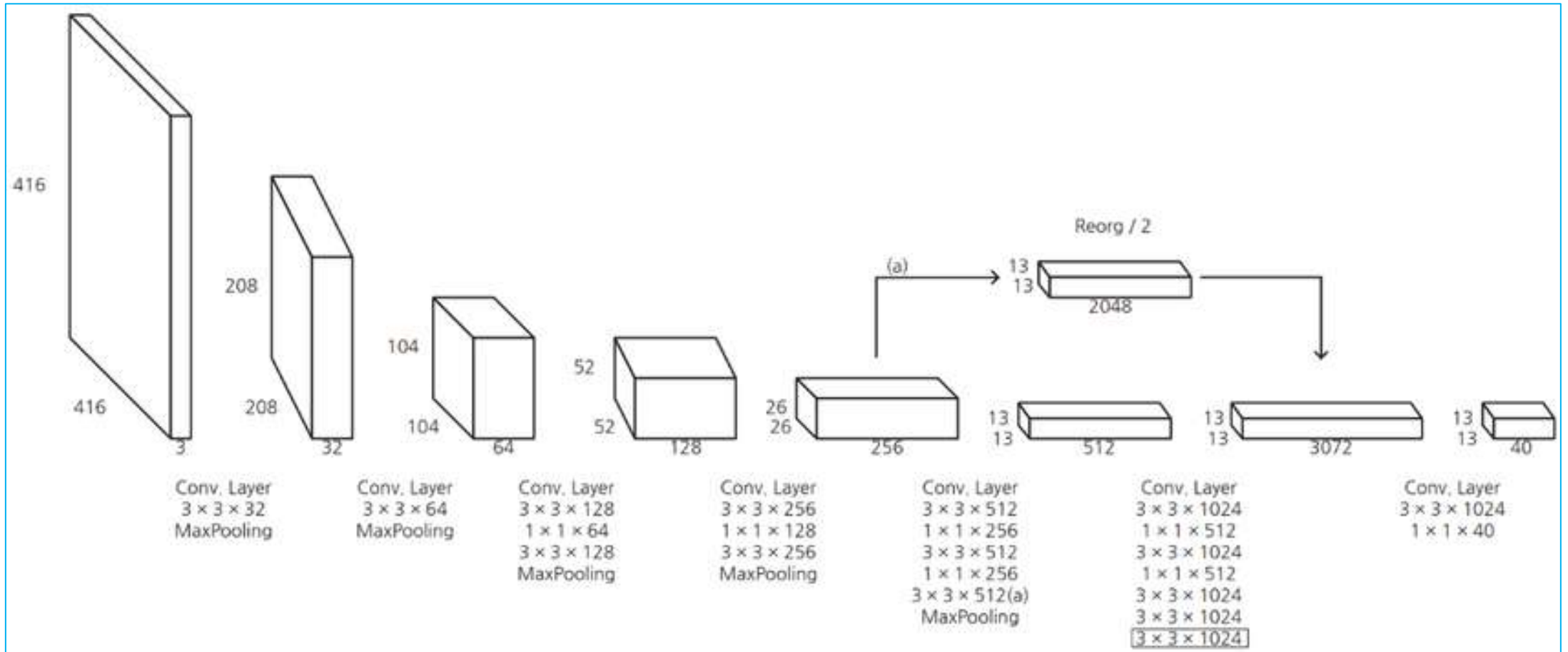
Instead of predicting the direct coordinates , they predict offsets to these bounding boxes during the training.

$$\begin{aligned} b_x &= \sigma(t_x) + c_x \\ b_y &= \sigma(t_y) + c_y \\ b_w &= p_w e^{t_w} \\ b_h &= p_h e^{t_h} \end{aligned}$$

# Changes in YOLO-v2

Add fine-grained feature

*The idea is similar to the skip connections in ResNet*





# Changes in YOLO-v2

## Multiple-Scale Training

Detection Frameworks	Training Data	mAP	FPS
Fast R-CNN	2007+2012	70.0	0.5
Faster R-CNN with VGG-16 backbone	2007+2012	73.2	7.0
Faster R-CNN with ResNet backbone	2007+2012	76.4	5.0
YOLOv1	2007+2012	63.4	45.0
SSD300	2007+2012	74.3	46.0
SSD500	2007+2012	76.8	19.0
YOLOv2 with input size 288 x 288	2007+2012	69.0	91.0
YOLOv2 with input size 352 x 352	2007+2012	73.7	81.0
YOLOv2 with input size 416 x 416	2007+2012	76.8	67.0
YOLOv2 with input size 480 x 480	2007+2012	77.8	59.0
YOLOv2 with input size 544 x 544	2007+2012	78.6	40.0

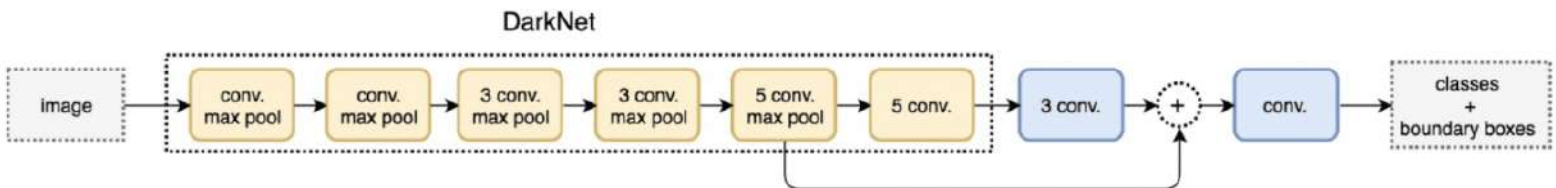


# Changes in YOLO-v2

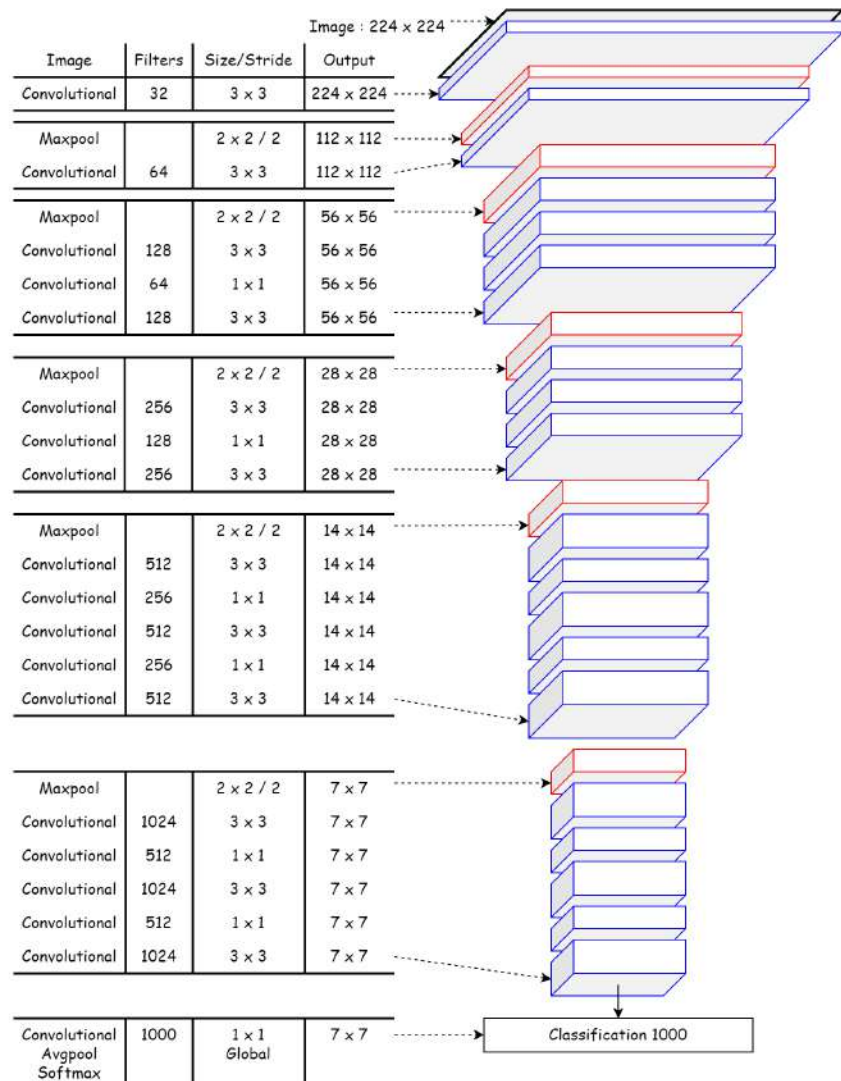
## Light-weight backbone

Darknet-19 A fully convolutional model with 19 convolutional layers and five max-pooling layers was designed.

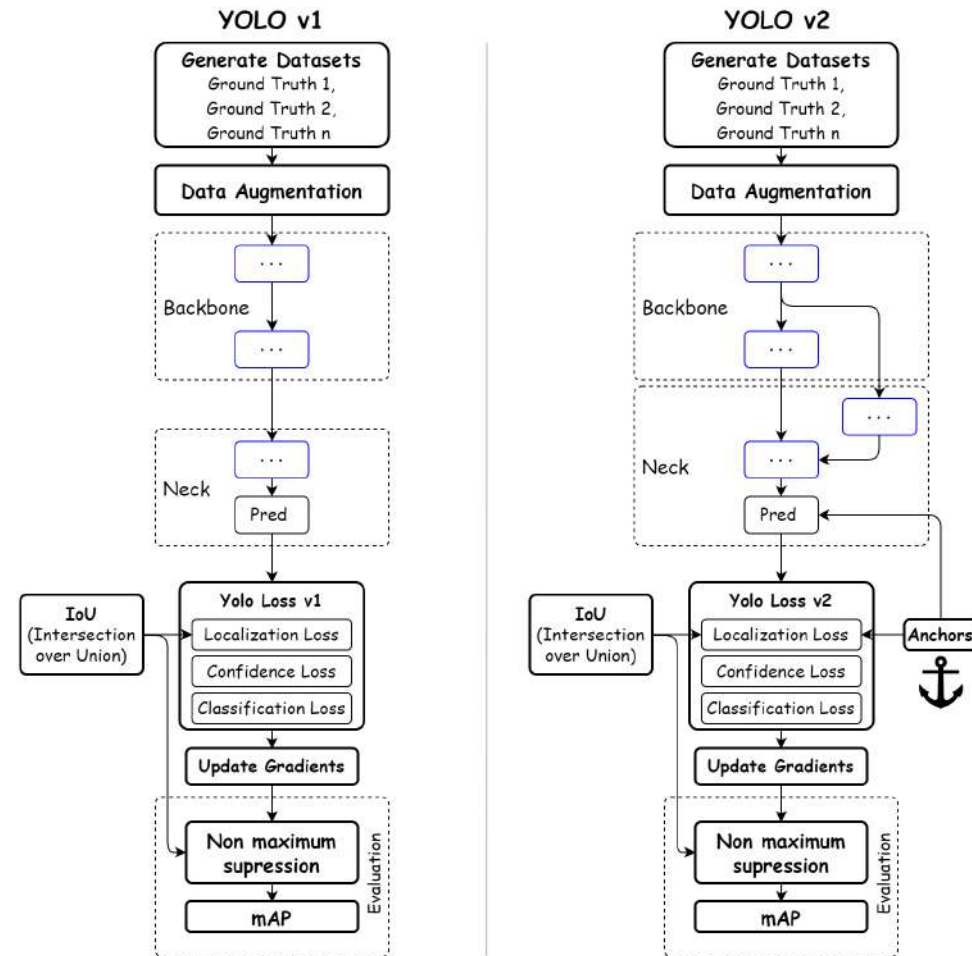
Type	Filters	Size/Stride	Output
Convolutional	32	$3 \times 3$	$224 \times 224$
Maxpool		$2 \times 2/2$	$112 \times 112$
Convolutional	64	$3 \times 3$	$112 \times 112$
Maxpool		$2 \times 2/2$	$56 \times 56$
Convolutional	128	$3 \times 3$	$56 \times 56$
Convolutional	64	$1 \times 1$	$56 \times 56$
Convolutional	128	$3 \times 3$	$56 \times 56$
Maxpool		$2 \times 2/2$	$28 \times 28$
Convolutional	256	$3 \times 3$	$28 \times 28$
Convolutional	128	$1 \times 1$	$28 \times 28$
Convolutional	256	$3 \times 3$	$28 \times 28$
Maxpool		$2 \times 2/2$	$14 \times 14$
Convolutional	512	$3 \times 3$	$14 \times 14$
Convolutional	256	$1 \times 1$	$14 \times 14$
Convolutional	512	$3 \times 3$	$14 \times 14$
Convolutional	256	$1 \times 1$	$14 \times 14$
Convolutional	512	$3 \times 3$	$14 \times 14$
Maxpool		$2 \times 2/2$	$7 \times 7$
Convolutional	1024	$3 \times 3$	$7 \times 7$
Convolutional	512	$1 \times 1$	$7 \times 7$
Convolutional	1024	$3 \times 3$	$7 \times 7$
Convolutional	512	$1 \times 1$	$7 \times 7$
Convolutional	1024	$3 \times 3$	$7 \times 7$
Convolutional	1000	$1 \times 1$	$7 \times 7$
Avgpool		Global	1000
Softmax			



# YOLO-v1 vs. YOLO-v2

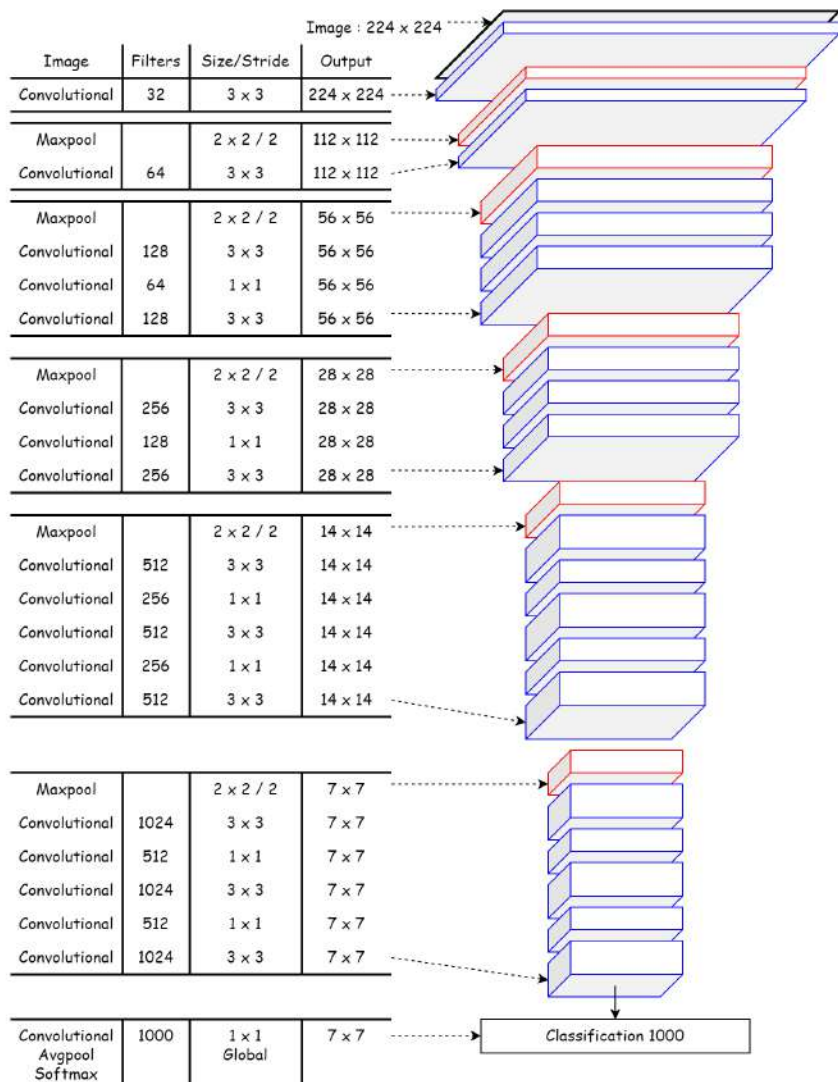


## Workflow Comparison



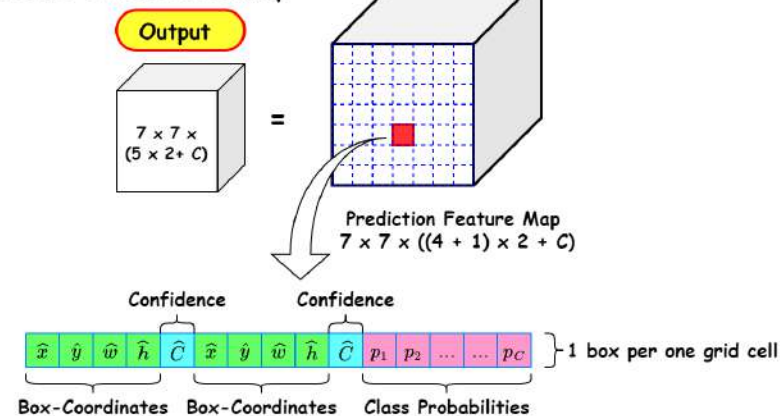
Evaluation will be explained in detail after the Yolo V3 Training.

# YOLO-v1 vs. YOLO-v2

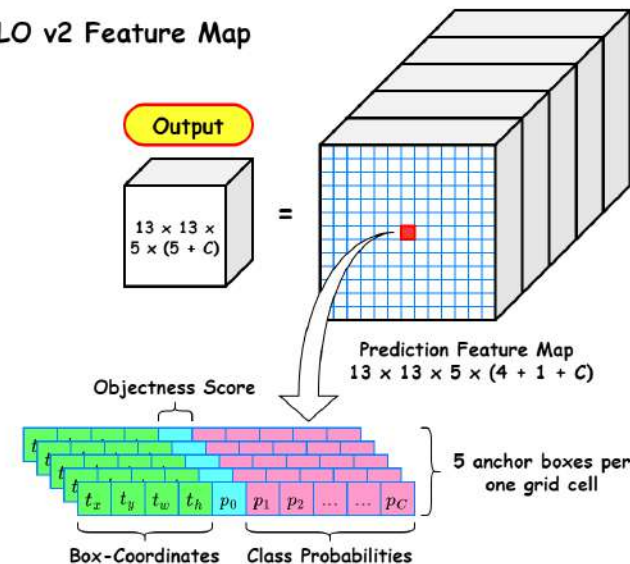


Head

YOLO v1 Feature Map



YOLO v2 Feature Map



# YOLO-v1 vs. YOLO-v2

YOLO v1 Loss

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

YOLO v2 Loss

$$\begin{aligned} \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} & \left[ (x_i - b_{xi})^2 + (y_i - b_{yi})^2 \right] \\ & + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} \left[ \left( \sqrt{w_i} - \sqrt{b_{wi}} \right)^2 + \left( \sqrt{h_i} - \sqrt{b_{hi}} \right)^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{obj} (C_i - \sigma(t_{oi}))^2 \\ & + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{noobj} (C_i - \sigma(t_{oi}))^2 \\ & + \sum_{i=0}^{S^2} \mathbb{1}_i^{obj} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

$$b_x = \sigma(t_x) + c_x$$

$$b_y = \sigma(t_y) + c_y$$

$$b_w = p_w e^{t_w}$$

$$b_h = p_h e^{t_h}$$

$$\text{Pr}(\text{object}) \star \text{IOU}(b, \text{object}) = \sigma(t_o)$$

where

$t_x, t_y, t_w, t_h$

$c_x, c_y$

$p_w, p_h$

$c_x, c_y, p_w, p_h$

$b_x, b_y, b_w, b_h$

$\sigma(t_o)$

are predictions made by YOLO.

is the top left corner of the grid cell of the anchor.

are the width and height of the anchor.

are normalized by the image width and height.

are the predicted boundary box.

is the box confidence score.



## Group DETR: Fast DETR Training with Group-Wise One-to-Many Assignment

Qiang Chen<sup>1\*</sup>, Xiaokang Chen<sup>2\*</sup>, Jian Wang<sup>1</sup>, Shan Zhang<sup>3</sup>  
Kun Yao<sup>1</sup>, Haocheng Feng<sup>1</sup>, Junyu Han<sup>1</sup>, Errui Ding<sup>1</sup>, Gang Zeng<sup>2</sup>, Jingdong Wang<sup>1†</sup>  
<sup>1</sup>Baidu VIS  
<sup>2</sup>Key Lab. of Machine Perception (MoE), School of IST, Peking University  
<sup>3</sup>Australian National University  
{chenqiang13, wangjian33}@baidu.com  
{fenghaocheng, hanjunyu, dingerrui, wangjingdong}@baidu.com  
{pkucxk, gang.zeng}@pku.edu.cn, shan.zhang@anu.edu.au

## Abstract

Detection transformer (DETR) relies on one-to-one assignment, assigning one ground-truth object to one prediction, for end-to-end detection without NMS post-processing. It is known that one-to-many assignment, assigning one ground-truth object to multiple predictions, succeeds in detection methods such as Faster R-CNN and FCOS. While the naive one-to-many assignment does not work for DETR, and it remains challenging to apply one-to-many assignment for DETR training. In this paper, we introduce Group DETR, a simple yet efficient DETR training approach that introduces a group-wise way for one-to-many assignment. This approach involves using multiple groups of object queries, conducting one-to-one assignment within each group, and performing decoder self-attention separately. It resembles data augmentation with automatically-learned object query augmentation. It is also equivalent to simultaneously training parameter-sharing networks of the same architecture, introducing more supervision and thus improving DETR training. The inference process is the same as DETR trained normally and only needs one group of queries without any architecture modification. Group DETR is versatile and is applicable to various DETR variants. The experiments show that Group DETR significantly speeds up the training convergence and improves the performance of various DETR-based models. Code will be available at <https://github.com/Atten4Vis/GroupDETR>.

## Group DETR v2: Strong Object Detector with Encoder-Decoder Pretraining

Qiang Chen<sup>1\*</sup>, Jian Wang<sup>1\*</sup>, Chuchu Han<sup>1\*</sup>, Shan Zhang<sup>2</sup>, Zexian Li<sup>3</sup>, Xiaokang Chen<sup>4</sup>, Jiahui Chen<sup>5</sup>  
Xiaodi Wang<sup>1</sup>, Shuming Han<sup>1</sup>, Gang Zhang<sup>1</sup>, Haocheng Feng<sup>1</sup>, Kun Yao<sup>1</sup>, Junyu Han<sup>1</sup>, Errui Ding<sup>1</sup>  
Jingdong Wang<sup>1†</sup>

<sup>1</sup>Baidu VIS <sup>2</sup>Australian National University <sup>3</sup>Beihang University <sup>4</sup>Peking University

Table 1. Our method establishes a new SoTA on the COCO test-dev leaderboard.

Method	#Params	Encoder Pretraining Data	Detector Pretraining Data	w/ Mask	mAP
Swin-L (HTC++) [16]	284M	IN-22K (14M)	n/a	✓	58.7
DyHead (Swin-L) [6]	213M	IN-22K (14M)	n/a	✓	60.6
Soft-Teacher (Swin-L) [25]	284M	IN-22K (14M)	COCO-unlabeled + C365	✓	61.3
GLIP (DyHead) [11]	≥284M	IN-22K (14M)	FourODs + GoldG + Cap24M	×	61.5
Florence (CoSwin-H) [29]	≥637M	FLD-900M (900M)	FLD-9M	×	62.4
GLIPv2 (CoSwin-H) [29]	≥637M	FLD-900M (900M)	FLD-9M	×	62.4
SwinV2-G (HTC++) [15]	3.0B	IN-22K + ext-70M (84M)	C365	✓	63.1
DINO (Swin-L) [28]	218M	IN-22K (14M)	C365	✓	63.3
BEiT-3 (ViTDet) [23]	1.9B	IN-22K + Image-Text (35M) + Text (160GB)	C365	✓	63.7
FD-SwinV2-G (HTC++) [23]	3.0B	IN-22K + IN-1K + ext-70M (85M)	C365	✓	64.2
FocalNet-H (DINO) [36]	746M	IN-22K (14M)	C365	×	64.3
Group DETR v2 (Our method)	629M	IN-1K (14M)	C365	×	64.5

All the results are achieved with test-time augmentation. In the table, we follow the notations for various datasets used in DINO [28] and FocalNet [36]. 'w/ Mask' means using mask annotations when finetuning the detectors on COCO [13].

## Abstract

We present a strong object detector with encoder-decoder pretraining and finetuning. Our method, called Group DETR v2, is built upon a vision transformer encoder ViT-Huge [8], a DETR variant DINO [28], and an efficient DETR training method Group DETR [3]. The training process consists of self-supervised pretraining and finetuning a ViT-Huge encoder on ImageNet-1K, pretraining the detector on Object365, and finally finetuning it on COCO. Group DETR v2 achieves 64.5 mAP on COCO test-dev, and establishes a new SoTA on the COCO leaderboard<sup>1</sup>.

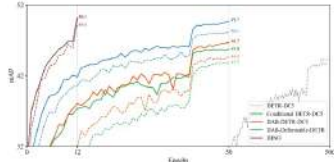


Figure 1. Group DETR accelerates the training process for DETR variants. The training convergence curves are obtained on COCO val2017 [34] with ResNet-50 [22]. Dashed and bold curves correspond to the baseline models and the Group DETR counterparts. Best viewed in color.

crafted components, such as non-maximum suppression (NMS) [23] and anchor generation [44, 53, 43]. The architecture consists of a CNN [22] and transformer encoder [53], and a transformer decoder that consists of self-attention, cross-attention and FFNs, followed by class and box prediction FFNs. During training, one-to-one assignment, where one ground-truth object is assigned to one single prediction, is applied for learning to only promote the predictions assigned to ground-truth objects, and demote the duplicate predictions.

This work explores the solutions to accelerate the DETR training process. Previous solutions contain two main lines. The one line is to modify cross-attention so that informa-

## Detection Transformer with Stable Matching

Shilong Liu<sup>1,2,4†</sup>, Tianhe Ren<sup>2\*</sup>, Jiayu Chen<sup>3\*</sup>,  
Zhaoyang Zeng<sup>2</sup>, Hao Zhang<sup>2,4</sup>, Feng Li<sup>2,4</sup>, Hongyang Li<sup>2,5</sup>,  
Jun Huang<sup>3</sup>, Hang Su<sup>1</sup>, Jun Zhu<sup>1†</sup>, Lei Zhang<sup>2†</sup>

<sup>1</sup>Dept. of Comp. Sci. and Tech., BNRist Center, State Key Lab for Intell. Tech. & Sys.,  
Institute for AI, Tsinghua-Bosch Joint Center for ML, Tsinghua University

<sup>2</sup>International Digital Economy Academy (IDEA) <sup>3</sup>Platform of AI (PAI), Alibaba Group

<sup>4</sup>The Hong Kong University of Science and Technology

<sup>5</sup>South China University of Technology

liusl20@mails.tsinghua.edu.cn {rentianhe, zengzhaoyang}@idea.edu.cn {yunji.cjy, huangjian.hj}@alibaba-inc.com  
{zhhangcx, flay}@connect.ust.hk {fhwang, yeunglei}@mail.scut.edu.cn {suhangss, dcszj}@mail.tsinghua.edu.cn leizhang@idea.edu.cn

## Abstract

This paper is concerned with the matching stability problem across different decoder layers in Detection Transformers (DETR). We point out that the unstable matching in DETR is caused by a multi-optimization path problem, which is highlighted by the one-to-one matching design in DETR. To address this problem, we show that the most important design is to use and only use positional metrics (like IOU) to supervise classification scores of positive examples. Under the principle, we propose two simple yet effective modifications by integrating positional metrics to DETR's classification loss and matching cost, named position-supervised loss and position-modulated cost. We verify our methods on several DETR variants. Our methods show consistent improvements over baselines. By integrating our methods with DINO, we achieve 50.4 and 51.5 AP on the COCO detection benchmark using ResNet-50 backbones under 1× (12 epochs) and 2× (24 epochs) training settings, achieving a new record under the same setting. We achieve 63.8 AP on COCO detection test-dev with a Swin-Large backbone. Our code will be made available at <https://github.com/IDEA-Research/Stable-DINO>.

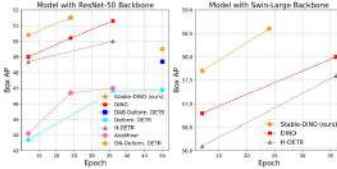


Figure 1: Comparison of our methods (named Stable-DINO in figures) and baselines. We compare models with ResNet-50 backbones in the left figure and models with Swin-Transformer Large backbones in the right figure. All models use a maximum 1/8 resolution feature map from a backbone, except AdaMixer uses a maximum 1/4 resolution feature map.

decades with the development of deep learning, especially the convolutional neural network (CNN) [36, 14, 16, 7].

Detection Transformer (DETR) [3] proposed a novel Transformer-based object detector, which attracted a lot of interest in the research community. It gets rid of the

## SAP-DETR: Bridging the Gap between Salient Points and Queries-Based Transformer Detector for Fast Model Convergence

Yang Liu<sup>1,3\*</sup>, Yao Zhang<sup>1,3\*</sup>, Yixin Wang<sup>2\*</sup>, Yang Zhang<sup>4</sup>, Jiang Tian<sup>4</sup>

Zhongchao Shi<sup>4</sup>, Jianping Fan<sup>4</sup>, Zhiqiang He<sup>1,3,5†</sup>

<sup>1</sup>Institute of Computing Technology (ICT), Chinese Academy of Sciences <sup>2</sup>Stanford University

<sup>3</sup>University of Chinese Academy of Sciences <sup>4</sup>AI Lab, Lenovo Research <sup>5</sup>Lenovo Ltd.

{liuyang20c, zhangyao21s}@mails.uccs.ac.cn yixinwang@stanford.edu  
{zhangyang20, tianjiang1, shizc2, jfan1, hezq}@lenovo.com

## Abstract

Recently, the dominant DETR-based approaches apply central-concept spatial prior to accelerating Transformer detector convergence. These methods gradually refine the reference points to the center of target objects and imbue object queries with the updated central reference information for spatially conditional attention. However, centralizing reference points may severely deteriorate queries' saliency and confuse the detectors due to the indiscriminative spatial prior. To bridge the gap between the reference points of salient queries and Transformer detectors, we propose Salient Point-based DETR (SAP-DETR) by treating object detection as a transformation from salient points to instance objects. Concretely, we explicitly initialize a query-specific reference point for each object query, gradually aggregate them into an instance object, and then predict the distance from each side of the bounding box to these points. By rapidly attending to query-specific reference regions and the conditional box edges, SAP-DETR can effectively bridge the gap between the salient point and the query-based Transformer detector with a significant convergence speed. Experimentally, SAP-DETR achieves 1.4× convergence speed with competitive performance and stably promotes the SoTA approaches by +1.0 AP. Based on ResNet-DC-101, SAP-DETR achieves 46.0 AP. This work will be released as a preprint.

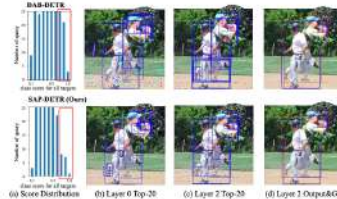


Figure 1. Comparison of SAP-DETR and DAB-DETR under 36 training epochs. (a) Statistics of the query count in different classification score intervals. (b) and (c) Distribution of reference points and the visualization of the query with top-20 classification score (blue proposal bounding boxes and red reference points) in different decoder layers. (d) Visualization of bounding boxes for positive queries (blue) and ground truth (red) during training process.

tors [6, 11, 14, 18, 20, 22] based on Convolutional Neural Networks (CNNs), have received widespread attention and made significant progress. Recently, Carion *et al.* [2] proposed a new end-to-end paradigm for object detection based on the Transformer [24], called DETection Transformer (DETR),

## SoTA

## DETRs with Collaborative Hybrid Assignments Training

Zhuofan Zong Guanglu Song Yu Liu\*

SenseTime Research

{zongzhuofan, liuyuisanai}@gmail.com

songguanglu@sensetime.com

## Abstract

In this paper, we provide the observation that too few queries assigned as positive samples in DETR with one-to-one set matching leads to sparse supervision on the encoder's output which considerably hurt the discriminative feature learning of the encoder and vice versa for attention learning in the decoder. To alleviate this, we present a novel collaborative hybrid assignments training scheme, namely Co-DETR, to learn more efficient and effective DETR-based detectors from versatile label assignment manners. This new training scheme can easily enhance the encoder's learning ability in end-to-end detectors by training the multiple parallel auxiliary heads supervised by one-to-many label assignments such as ATSS and Faster RCNN. In addition, we conduct extra customized positive queries by extracting the positive coordinates from these auxiliary heads to improve the training efficiency of positive samples in the decoder. In inference, these auxiliary heads are discarded and thus our method introduces no additional parameters and computational cost to the original detector while requiring no hand-crafted non-maximum suppression (NMS). We conduct extensive experiments to evaluate the effectiveness of the proposed approach on DETR variants, including DAB-DETR, Deformable-DETR, and DINO-Deformable-DETR. The state-of-the-art DINO-Deformable-DETR with Swin-L can be improved from 58.5% to 59.5% AP on COCO val. Surprisingly, incorporated with ViT-L backbone, we achieve 66.0% AP on COCO test-dev and 67.9% AP on LVIS val, outperforming previous methods by clear margins with much fewer model sizes. Codes are available at <https://github.com/Sense-X/Co-DETR>.

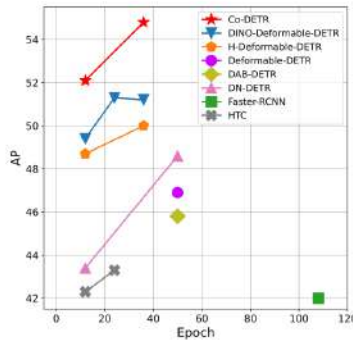


Figure 1. Performance of models with ResNet-50 on COCO val. Co-DETR outperforms other counterparts by a large margin.

a series of variants [31, 37, 44] such as ATSS [41], RetinaNet [21], FCOS [32], and PAA [17] lead to the significant breakthrough of object detection task. One-to-many label assignment is the core scheme of them, where each ground-truth box is assigned to multiple coordinates in the detector's output as the supervised target cooperated with proposals [11, 27], anchors [21] or window centers [32]. Despite their promising performance, these detectors heavily rely on many hand-designed components like a non-maximum suppression procedure or anchor generation [1]. To conduct a more flexible end-to-end detector, DETection Transformer (DETR) [1] is proposed to view the object detection as a set prediction problem and introduce the one-to-one set matching scheme based on a transformer encoder-decoder



