

~~Agenda:~~

- Anchor boxes (concept)
- R-CNN, Fast R-CNN , faster RCNN
(region -CNN) → 2 stage
- YOLO - 1 stage → sota algo
v1 ... v7

Chrome File Edit View History Bookmarks Profiles Tab Window Help

colab.research.google.com/drive/1MTvGvZPDhYaHhf3fkCyfgyFpaHqqAQO1#scrollTo=F3e5Jus83oh8

+ Code + Text Connect |

• And in general, you might use more anchor boxes, maybe five or even more.

Example: Overlapping Objects

1 2 3

YOLO

Anchor Box 1

human

Anchor Box 2

car

B_h

B_x, B_y

B_w

$\leftarrow B_w \rightarrow$

Output (Y) :

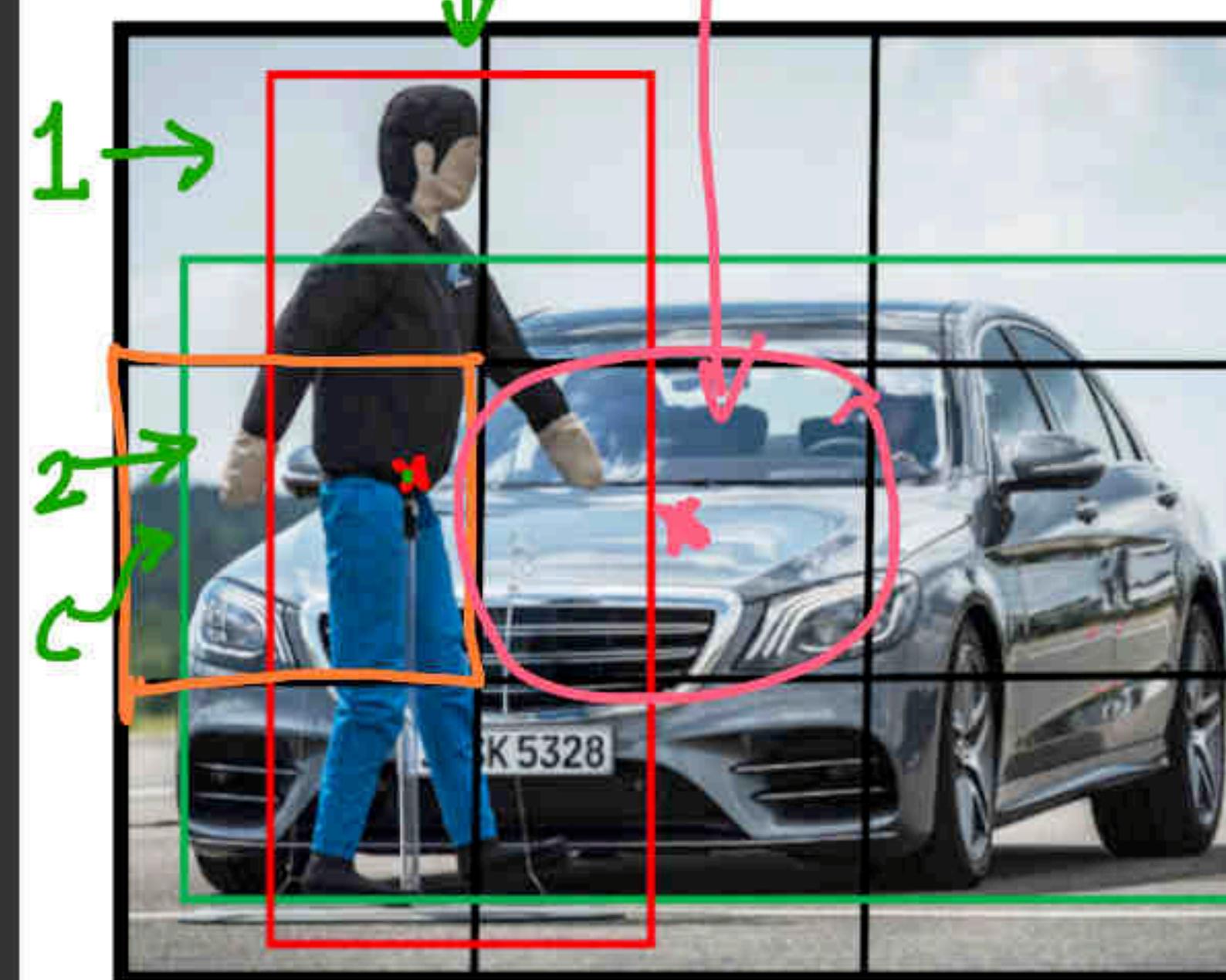
| | | | | | | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| P_c | B_x | B_y | B_h | B_w | C_1 | C_2 | C_3 | P_c | B_x | B_y | B_h | B_w | C_1 | C_2 | C_3 |
| 1 | B_x | B_y | B_h | B_w | 1 | 0 | 0 | 1 | D | D | D | D | 0 | 1 | 0 |

2 / 2

+ Code + Text

Connect

Example: 1 Overlapping Objects



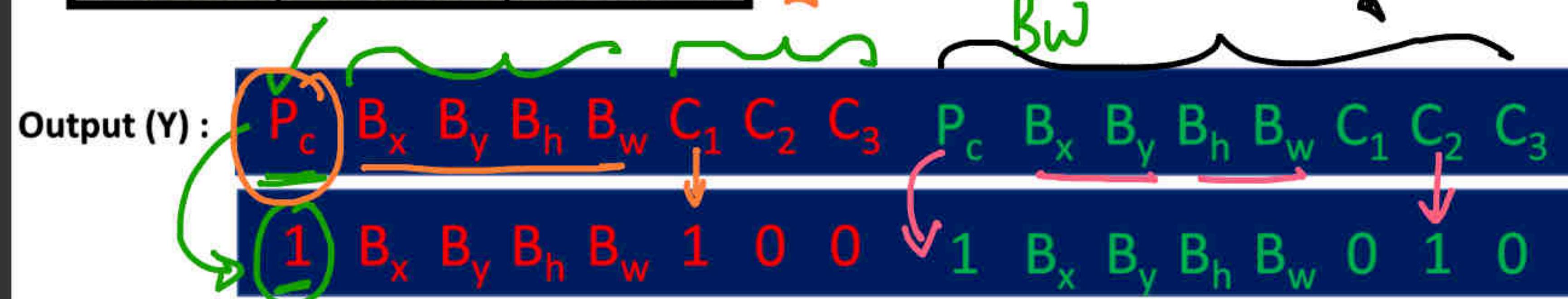
GRID

Anchor Box 1

 B_h B_x, B_y B_w

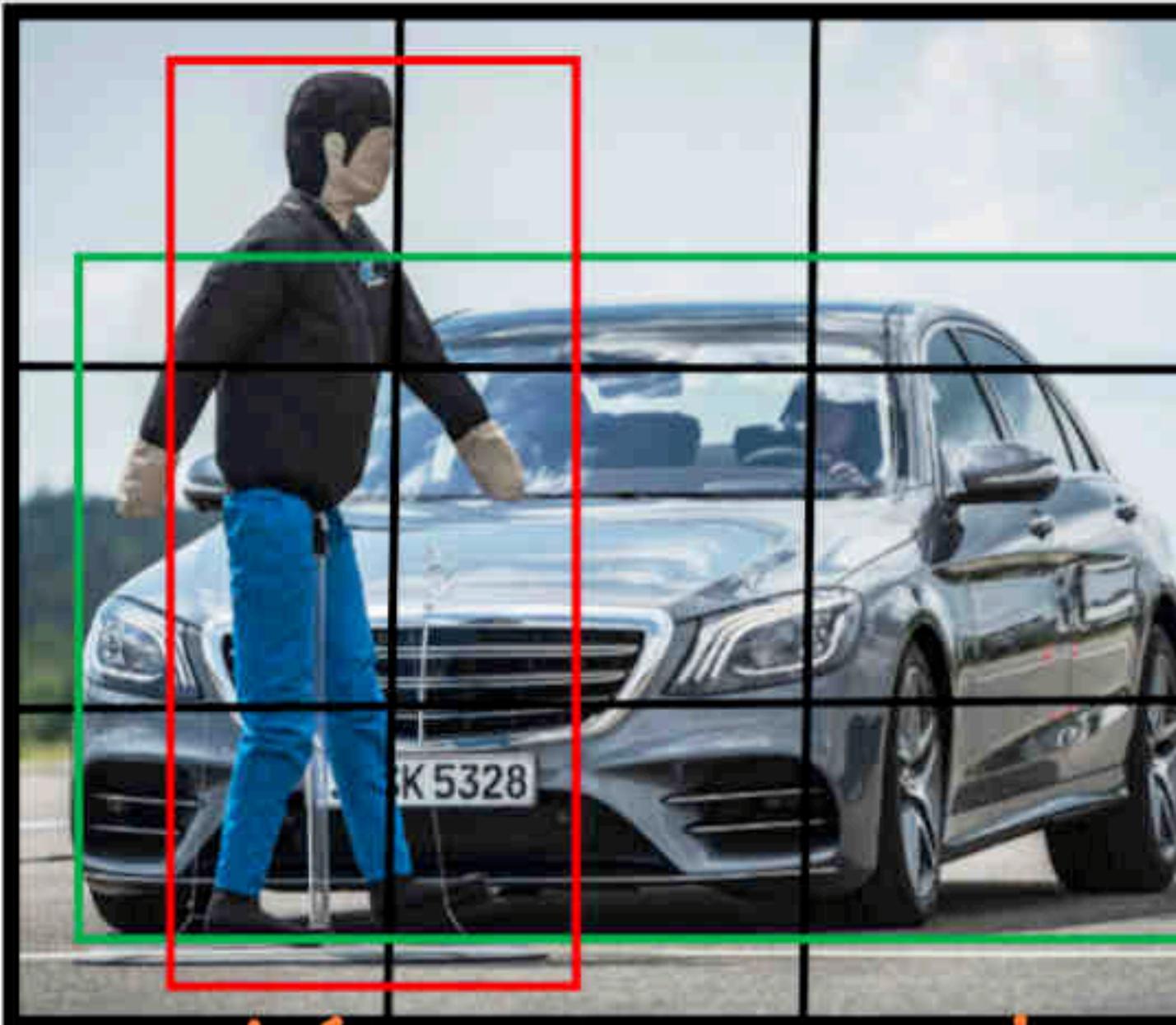
Anchor Box 2

a way to represent
 - 4 params of BBOX
 - prob of presence of
 any object
 - label of the object
 in the BBOX

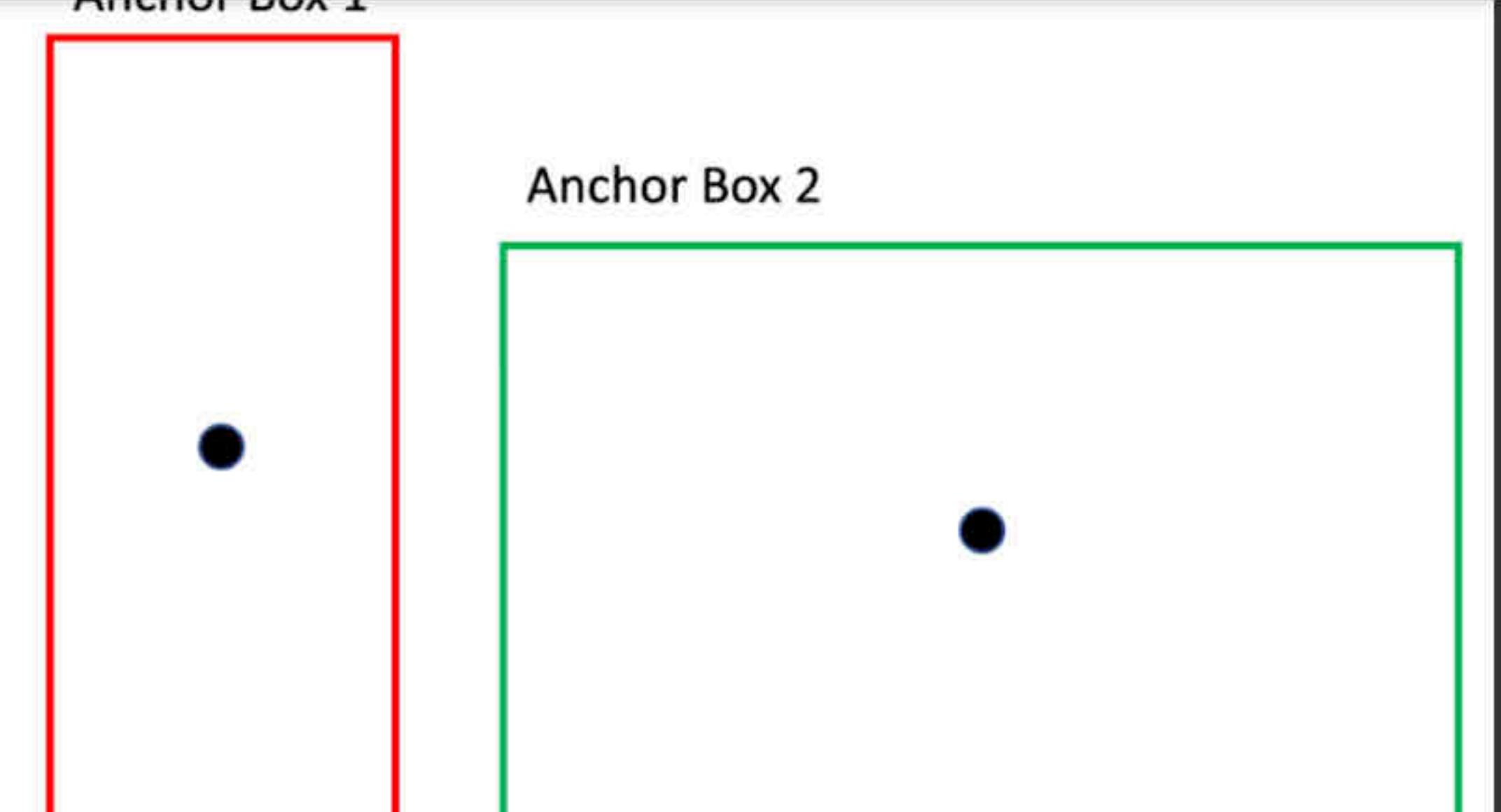


+ Code + Text Connect |

ANCHOR BOX 1



Anchor Box 2



log loss reg loss CE

Output (Y) :

| P_c | B_x | B_y | B_h | B_w | C_1 | C_2 | C_3 | P_c | B_x | B_y | B_h | B_w | C_1 | C_2 | C_3 |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1 | B_x | B_y | B_h | B_w | 1 | 0 | 0 | 1 | B_x | B_y | B_h | B_w | 0 | 1 | 0 |

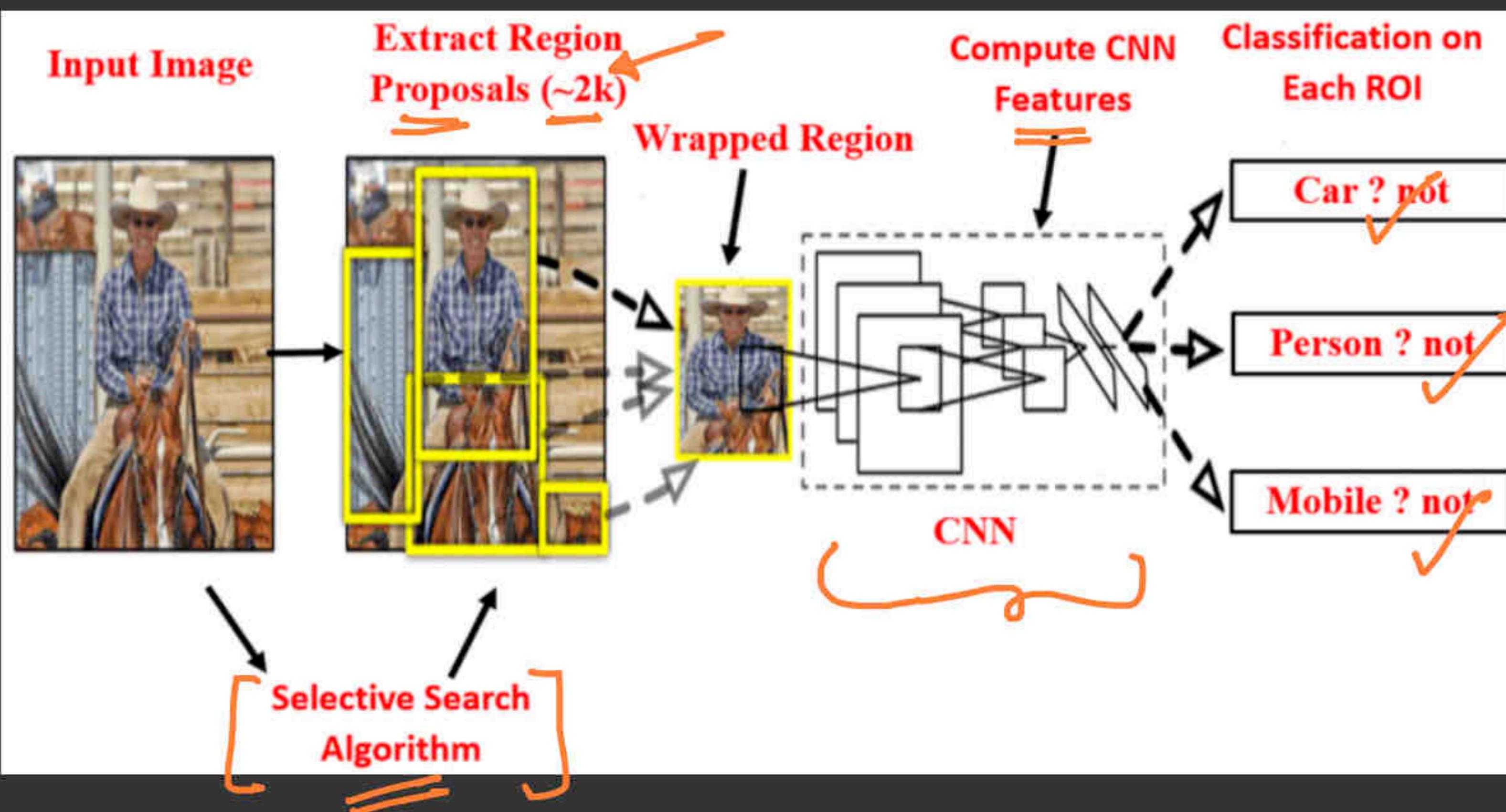
The output table shows two rows of values. The first row corresponds to the red anchor box (Person), and the second row corresponds to the green anchor box (Car). Handwritten annotations in red and orange highlight the first four columns (P_c , B_x , B_y , B_h) under the heading "log loss reg loss CE".

- In the above image Mid point of both Person and car lies in the same Grid Cell.
- We have already Pre-defined two bounding t

+ Code + Text

9:22

2. Then at second stage, the Regions proposed by first stage are classified into different classes.



Nov 2013

2 stages
→ external algo

regions

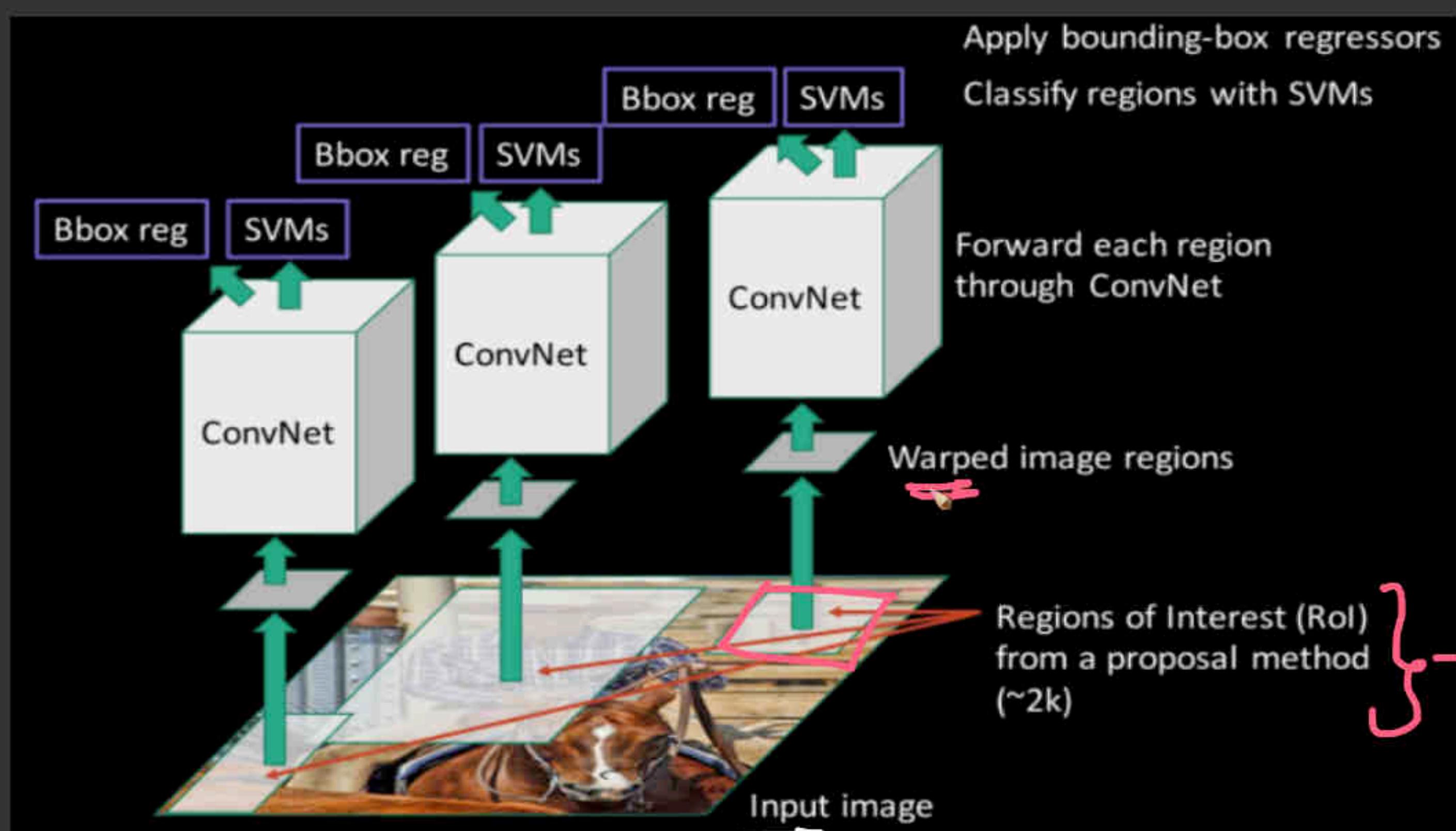
→ + region CNN → feature map

RCNN: Rich feature hierarchies for accurate object detection

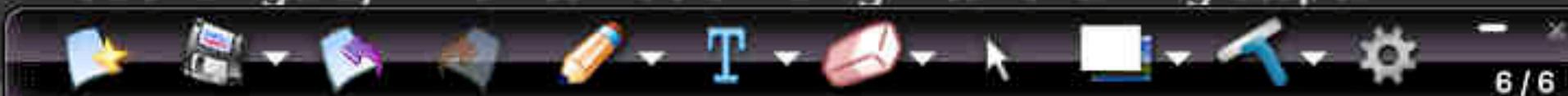
<https://arxiv.org/abs/1311.2524>

+ Code

RCNN: Rich feature hierarchies for accurate object detection

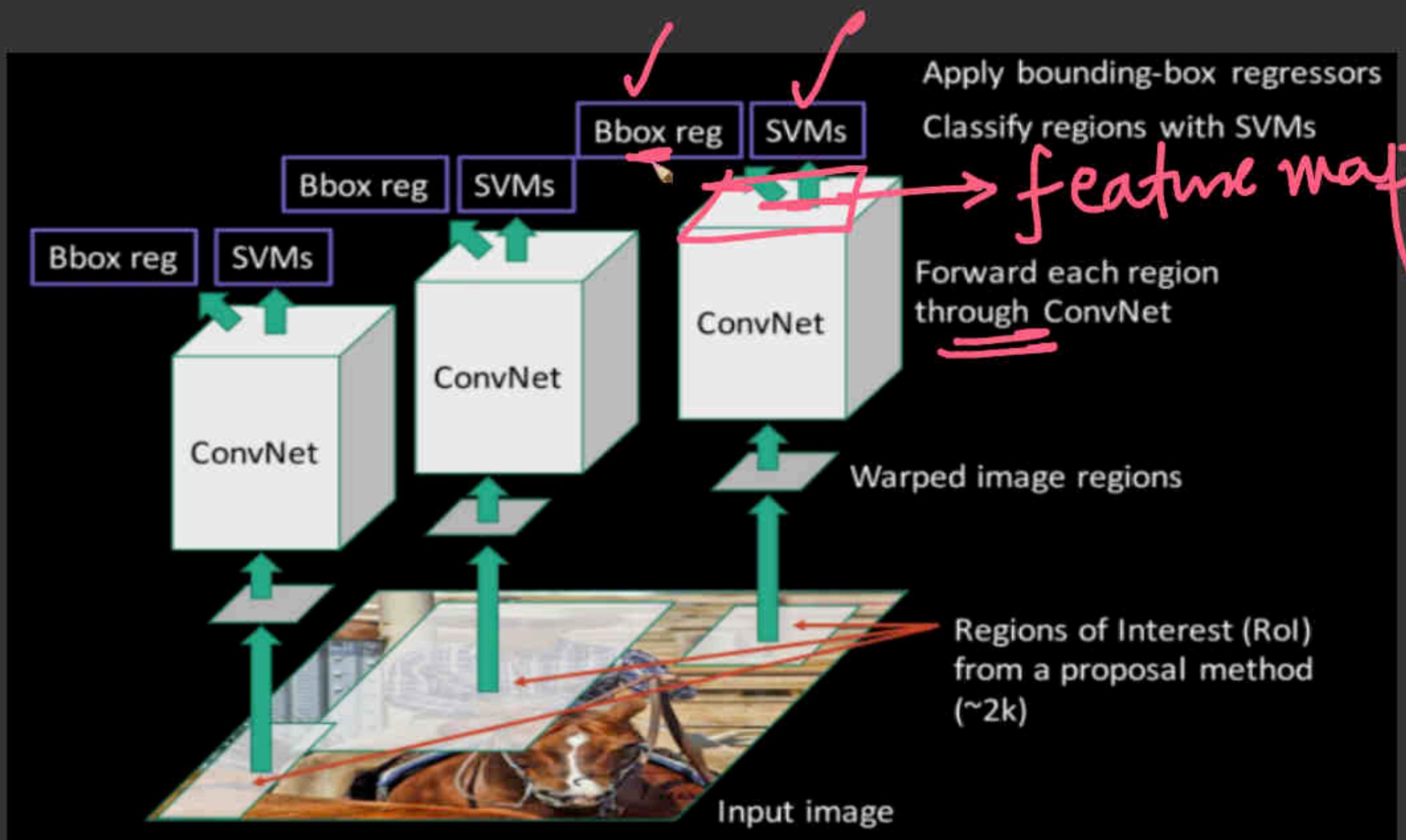
<https://arxiv.org/abs/1311.2524>

An overview of the R-CNN algorithm can be found in above figure, which can be dividing into following steps:



RCNN: Rich feature hierarchies for accurate object detection

<https://arxiv.org/abs/1311.2524>



An overview of the R-CNN algorithm can be found



+ Code + Text

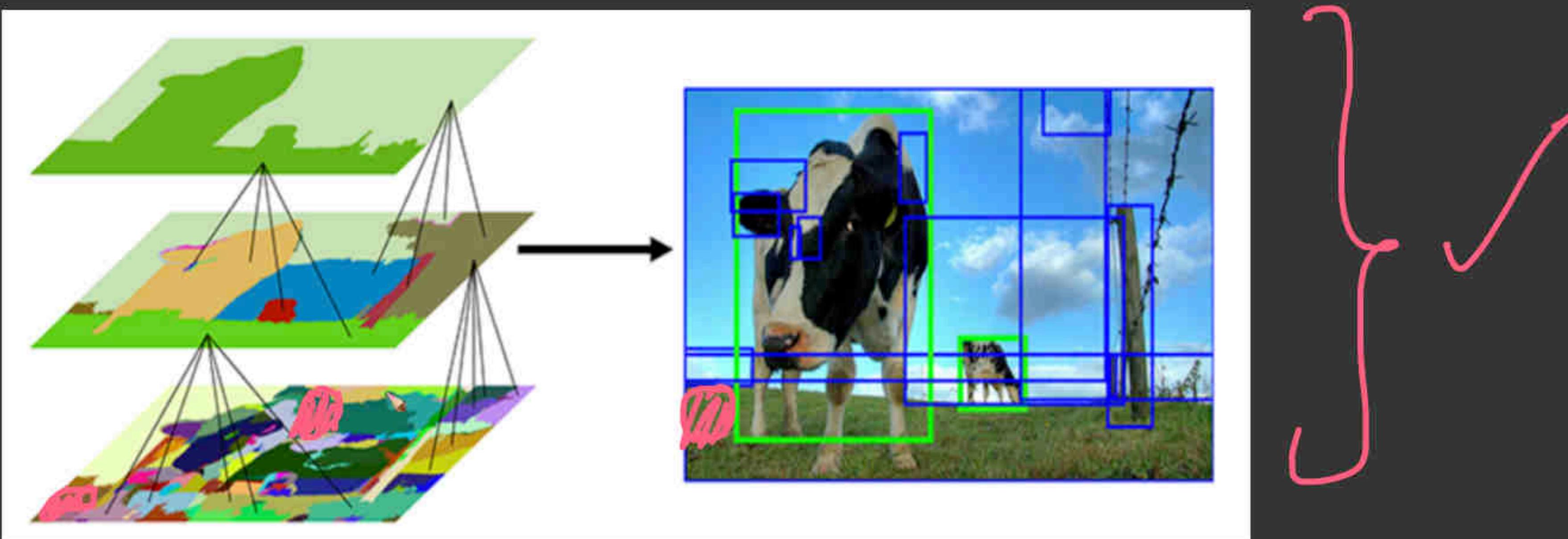
Connect |

2. Group adjacent segments based on similarity

3. Go to step 1 until a predefined number of bounding boxes are reached(typically 2000)

Advantage: Selective Search is far more computationally efficient than exhaustively computing image pyramids and sliding windows

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



Is it possible to obtain end-to-end deep learning-based object detection?

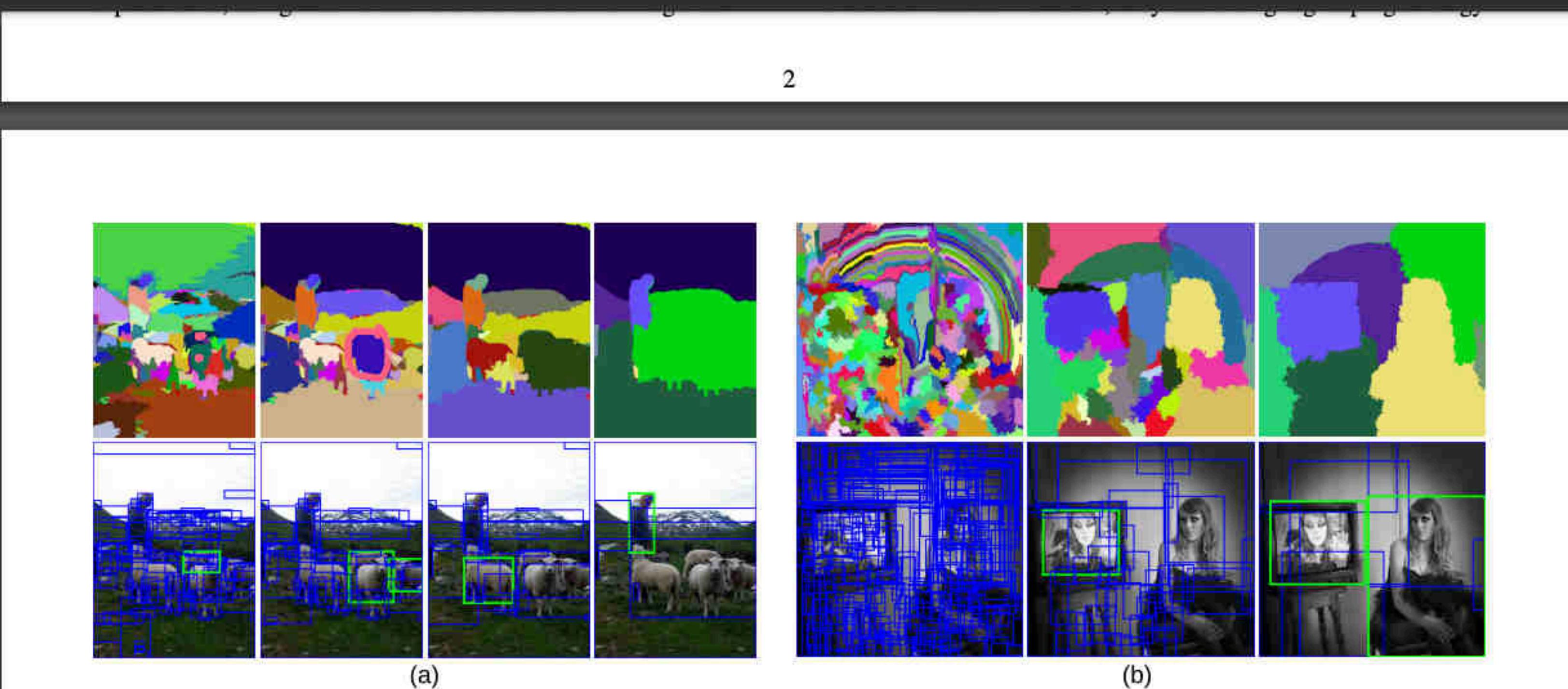


Figure 2: Two examples of our selective search showing the necessity of different scales. On the left we find many objects at different scales. On the right we necessarily find the objects at different scales as the girl is contained by the tv.

whose power of discovering parts or objects is left unevaluated. In this work, we use multiple complementary strategies to deal with as many image conditions as possible. We include the locations generated using [3] in our evaluation.

2.3 Other Sampling

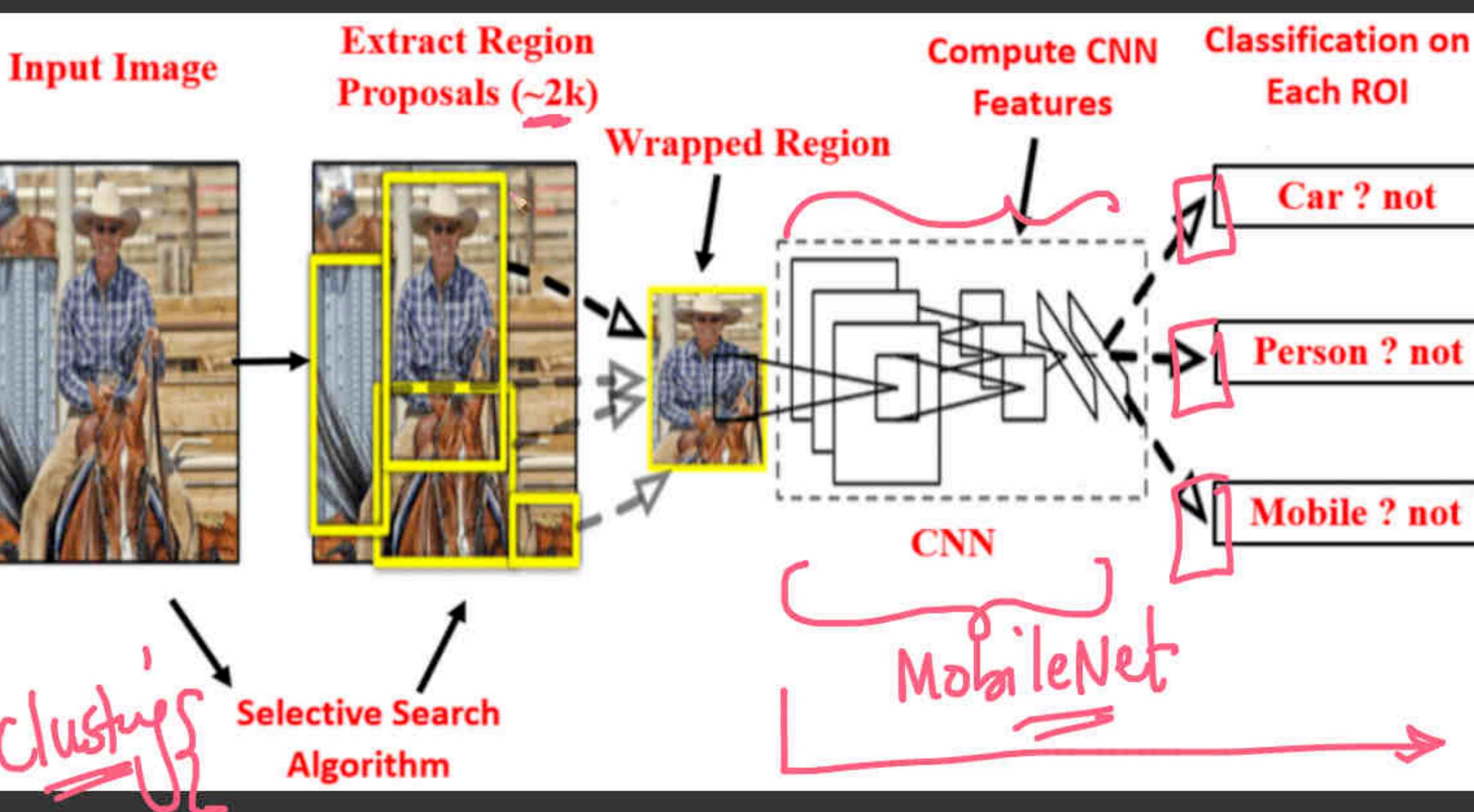
3 Selective Search

In this section we detail our selective search algorithm for object recognition and present a variety of diversification strategies to deal with as many image conditions as possible. A selective search algorithm is subject to the following design considerations:

Capture All Scales. Objects can occur at any scale within the image. Clear boundaries between objects and their backgrounds must be found. All objects of interest must be found, even if they have unclear boundaries or are partially occluded. This is most naturally achieved by using an hierarchical

+ Code + Text

Connect

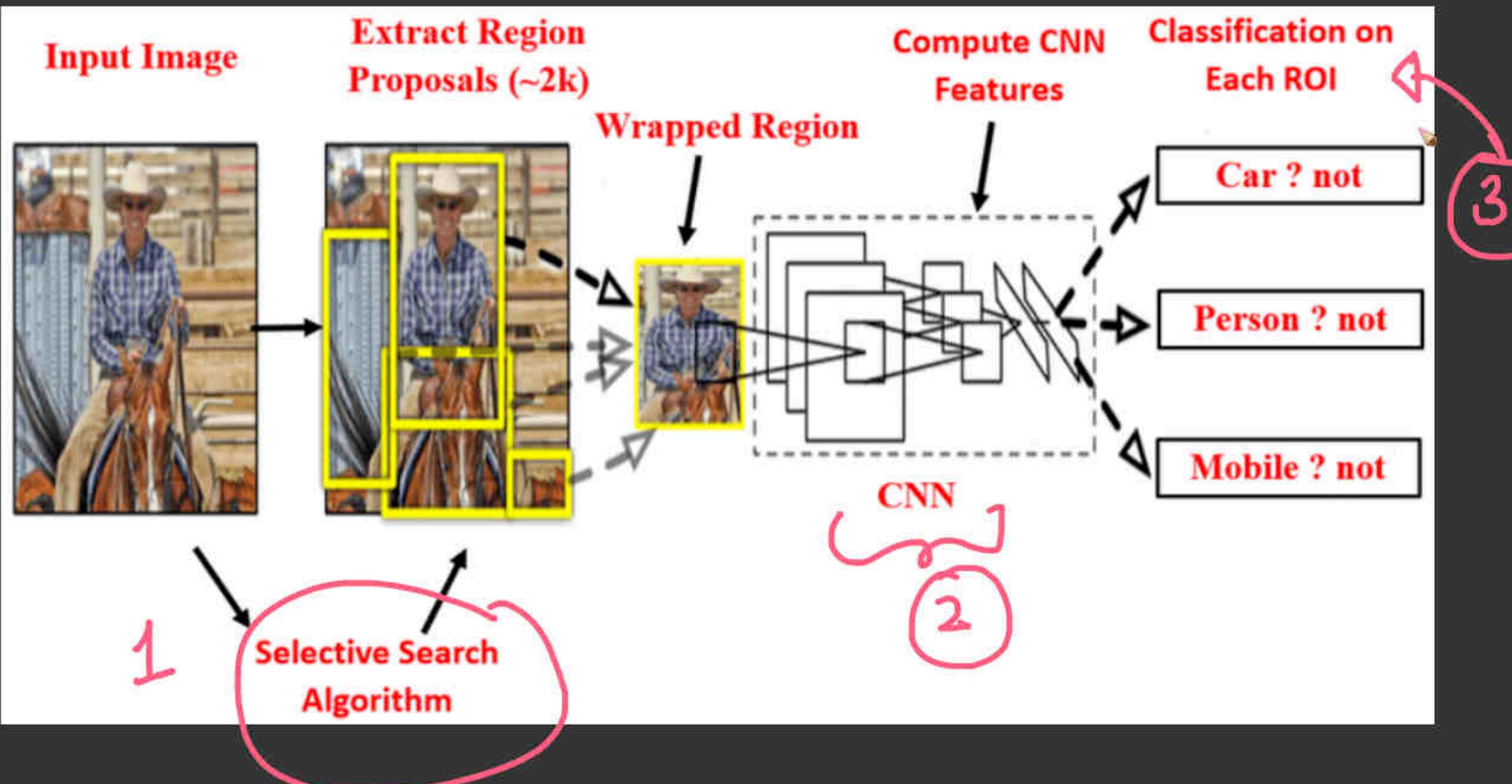


RCNN: Rich feature hierarchies for accurate object detection

<https://arxiv.org/abs/1311.2524>

+ Code + Text

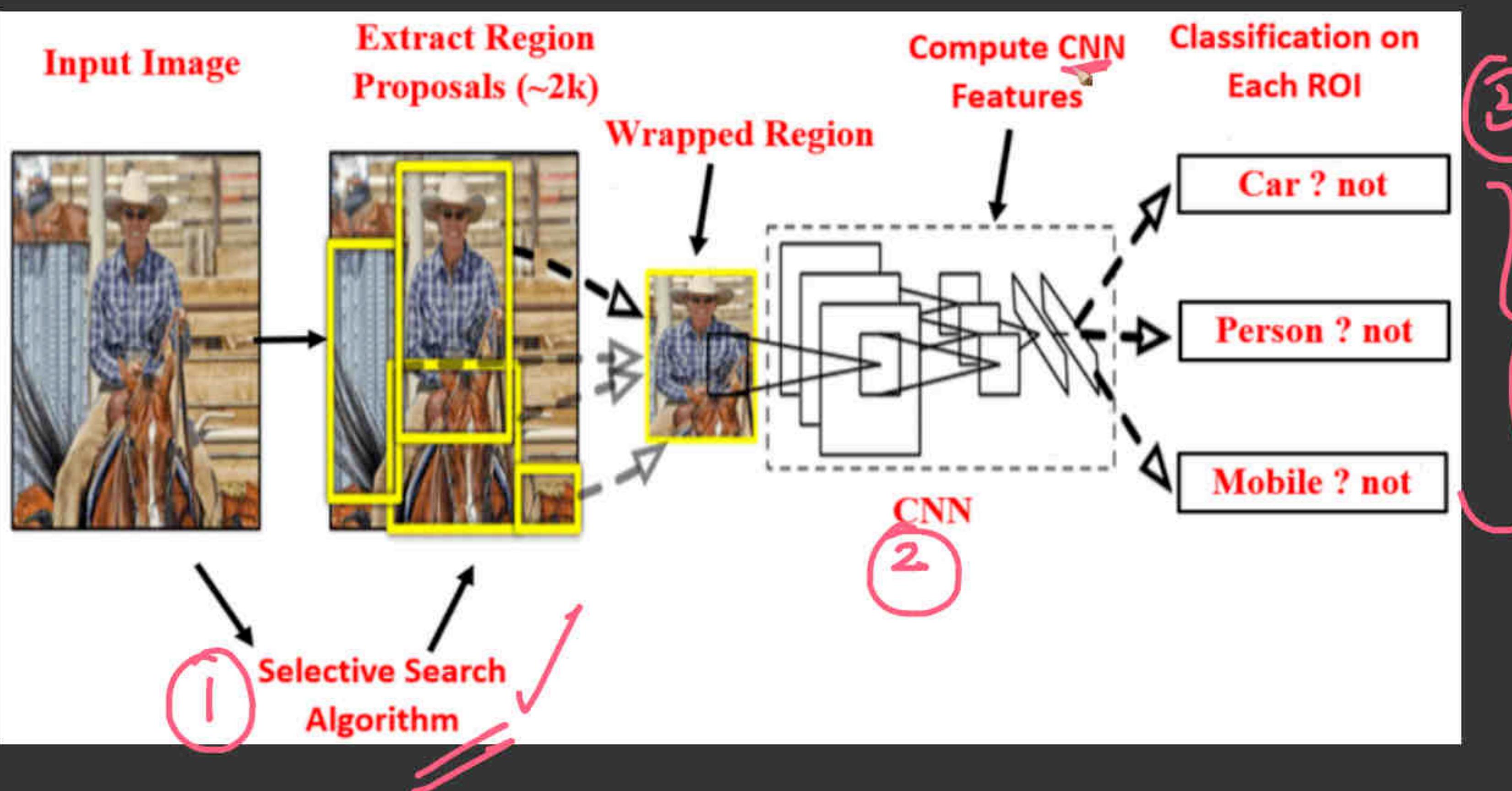
Connect |



▼ RCNN: Rich feature hierarchies for accurate object detection

<https://arxiv.org/abs/1311.2524>

2. Then at second stage, the Regions proposed by first stage are classified into different classes.



RCNN: Rich feature hierarchies for accurate object detection

<https://arxiv.org/abs/1311.2524>



+ Code + Text

Connect



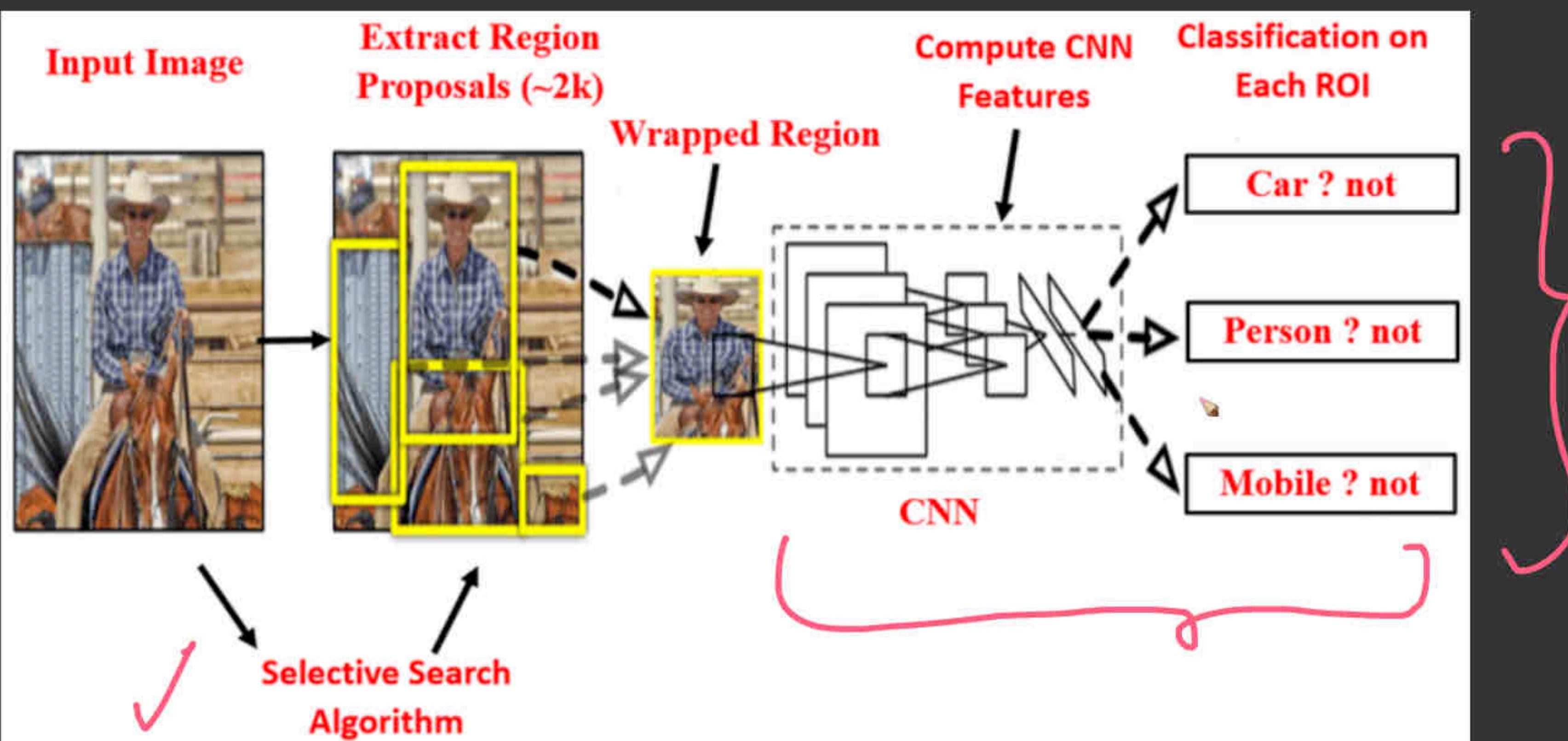
- In the **Fast R-CNN architecture**, we still use **Selective Search** to obtain our proposals, but now we apply **ROI Pooling** by extracting a fixed-size window from the feature map and using these features to obtain the final class label and bounding box.

Let's see the Fast R-CNN Architecture in detail

+ Code

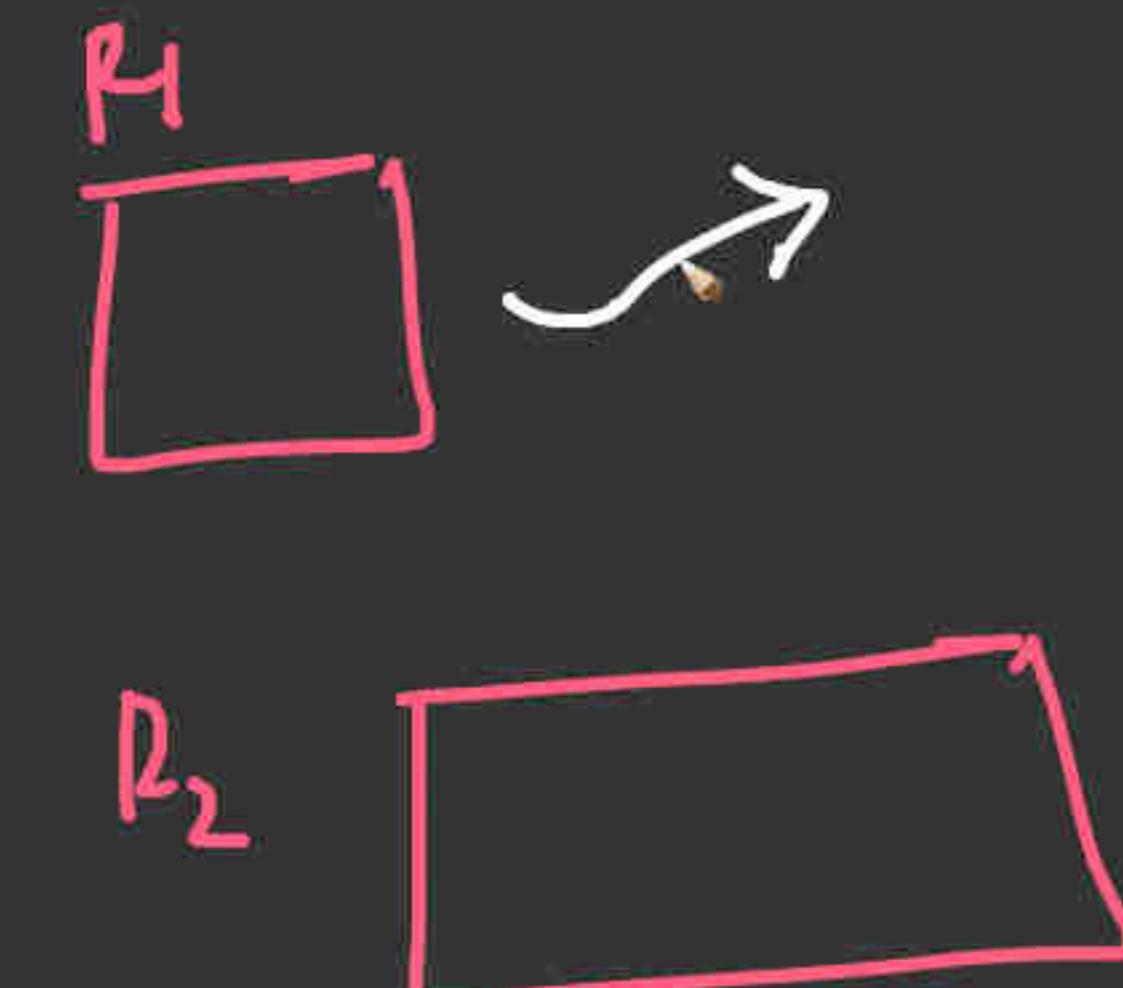
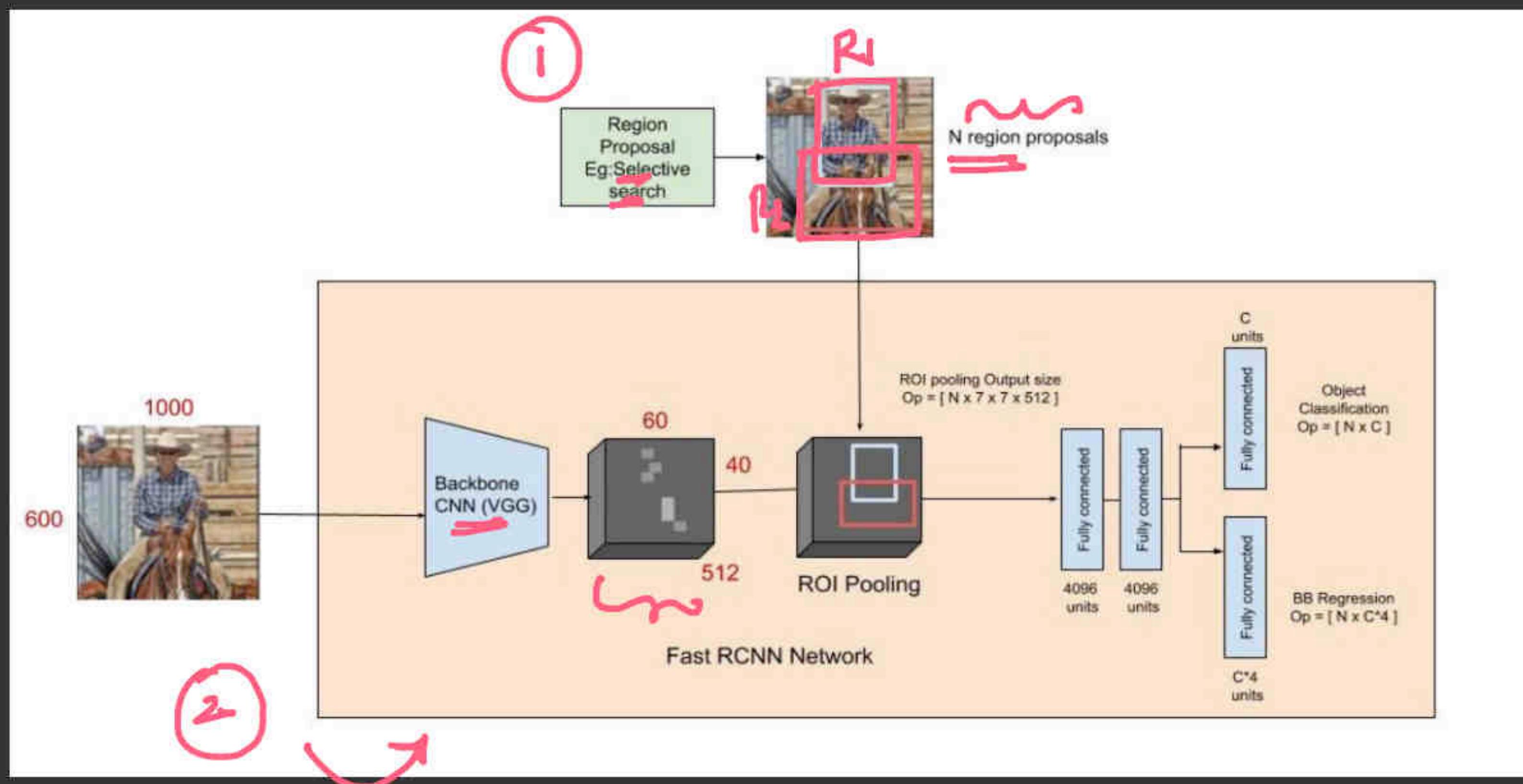
1. At first stage, Propose some Regions from the image where ther is a chance of Object present.

2. Then at second stage, the Regions proposed by first stage are classified into different classes.



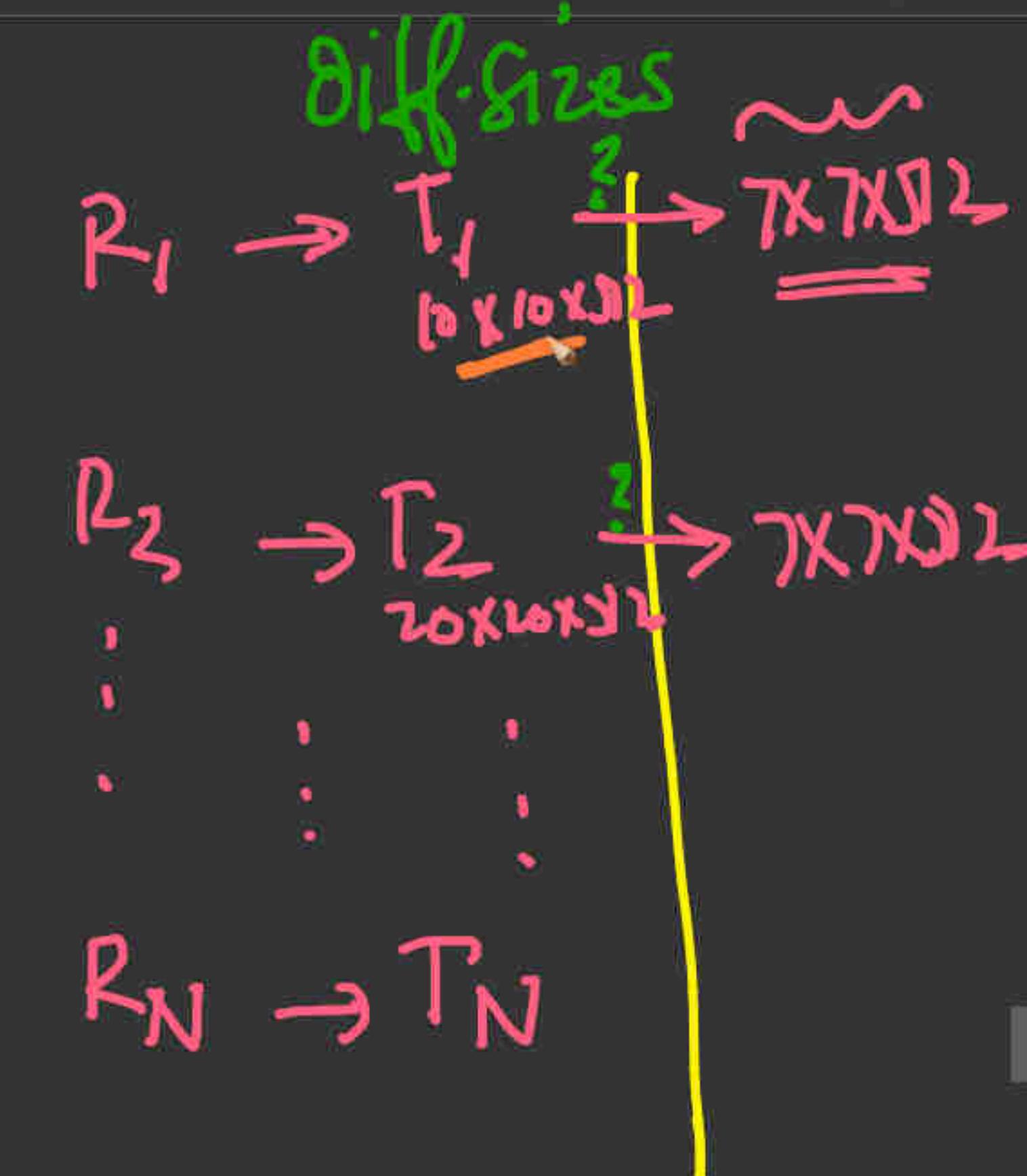
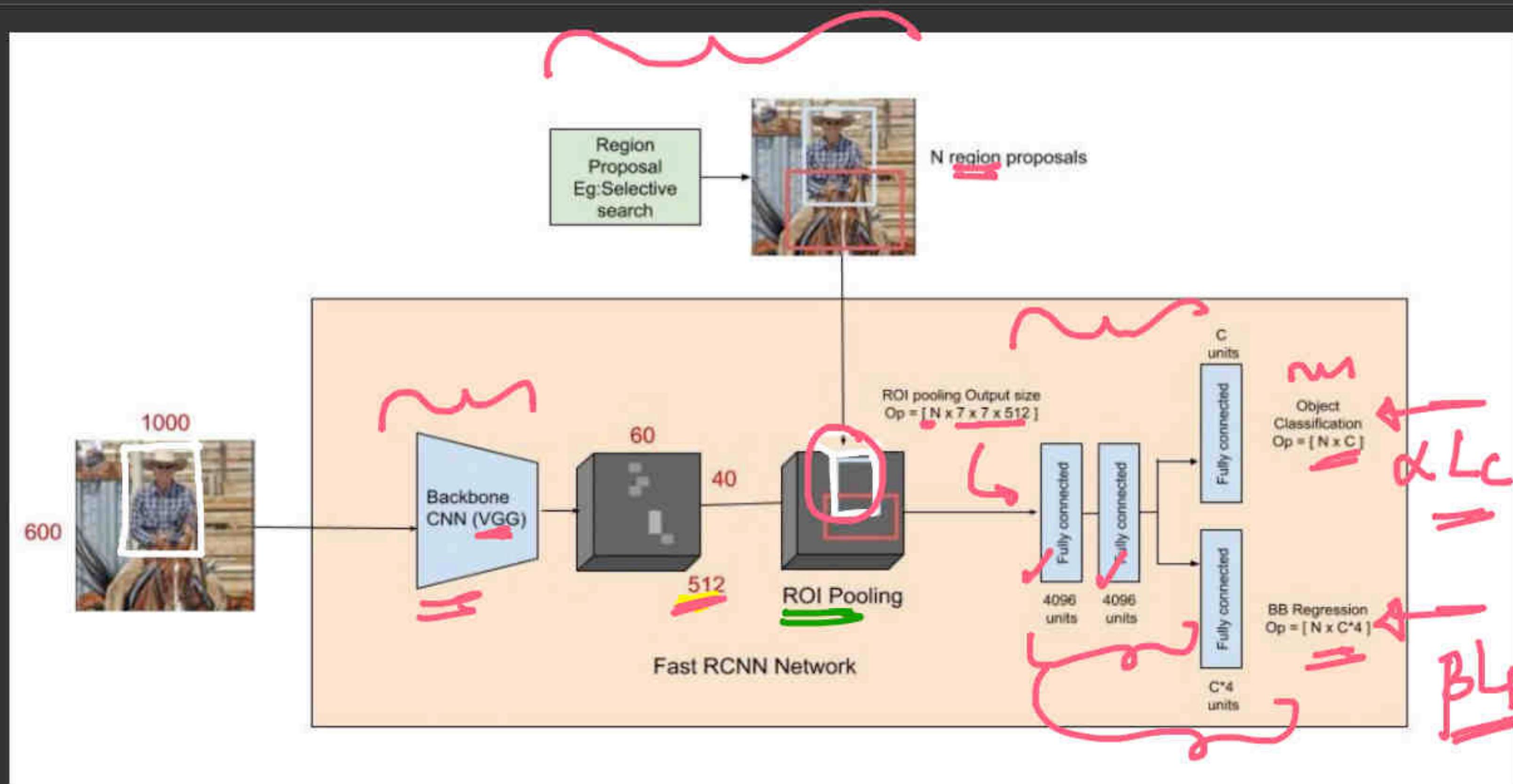
RCNN: Rich feature hierarchies for accurate object detection

+ Code + Text



The primary benefit here is that the network is now, effectively, end-to-end trainable:

1. We input an image and associated ground-truth bounding boxes.
2. Extract the feature map
3. Apply ROI pooling and obtain the ROI feature vector.
4. And finally use two sets of fully-connected layers



The primary benefit here is that the network is now, effectively, end-to-end trainable:

1. we input an image and associated ground-truth bounding boxes.
2. Extract the feature map
3. Apply ROI pooling and obtain the ROI feature vector.
4. And finally use two sets of fully-connected layers

feature maps

+ Code + Text

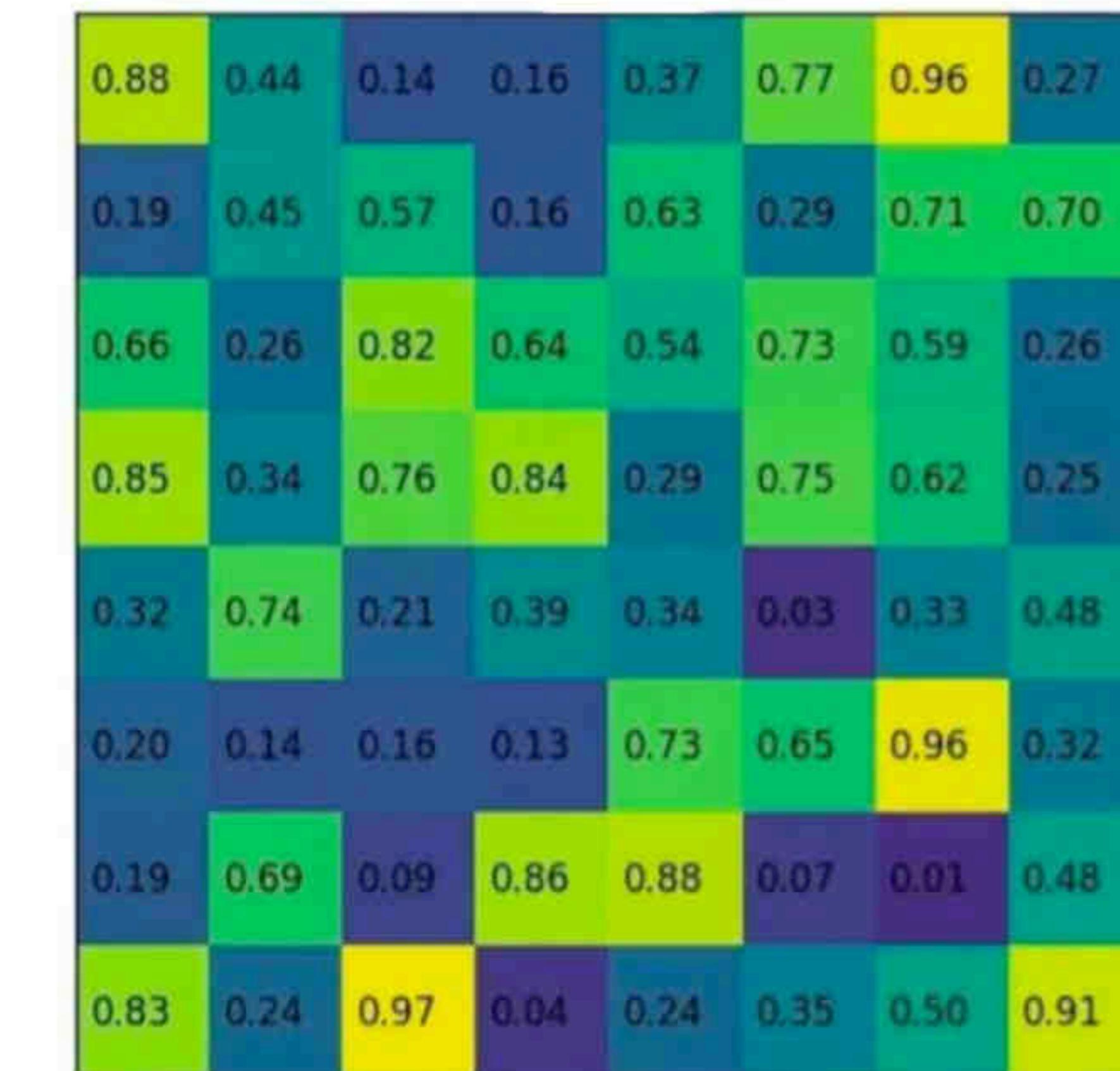
Connect



Input Image



Input feature map

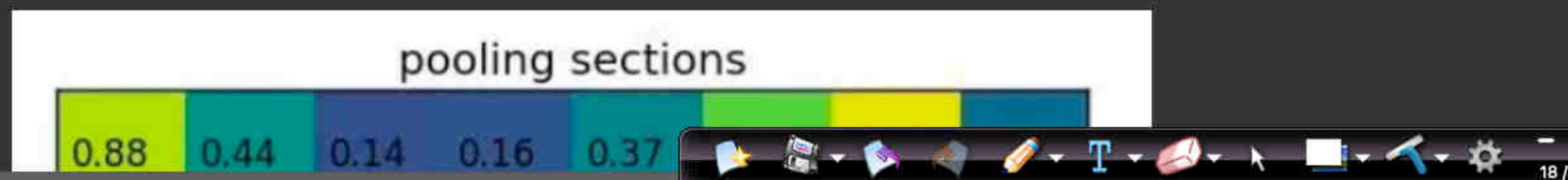


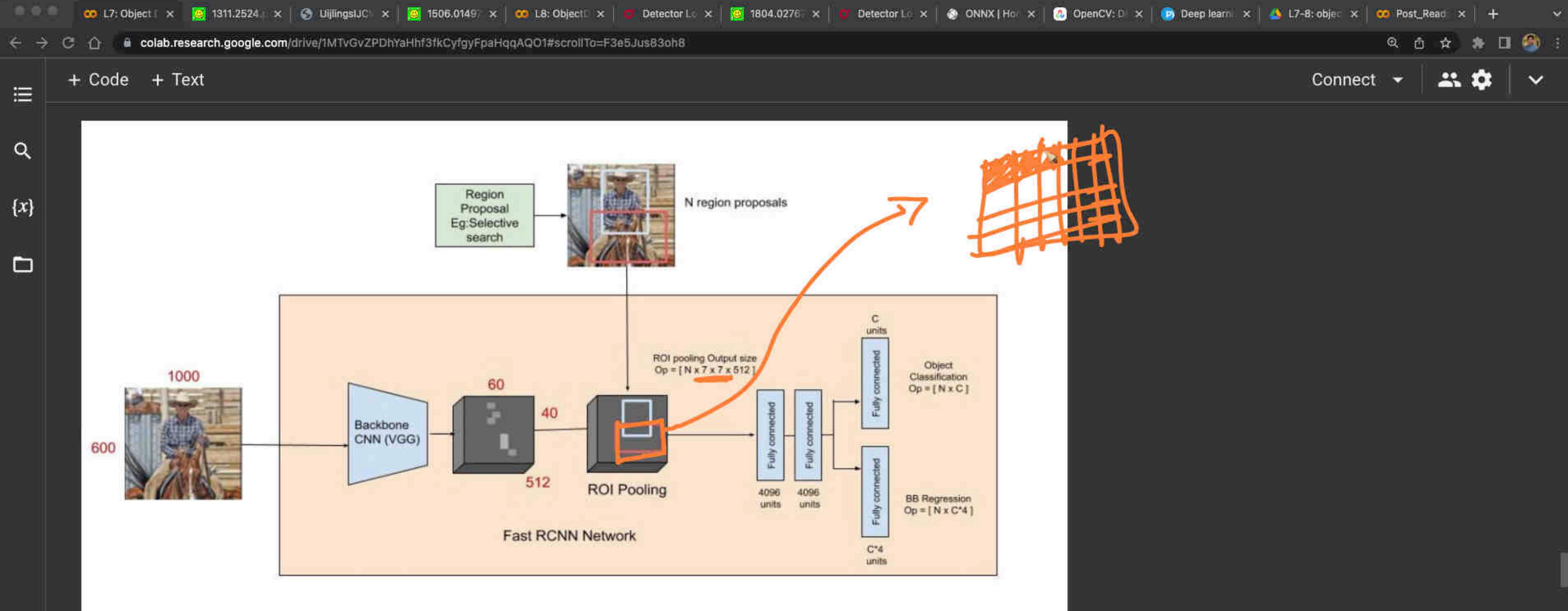
To Include: Corresponding GUN Image





Normally, there'd be multiple feature maps and multiple proposals for each of them, but we're keeping things simple for the example. By dividing it into (2×2) sections (because the output size is 2×2) we get:



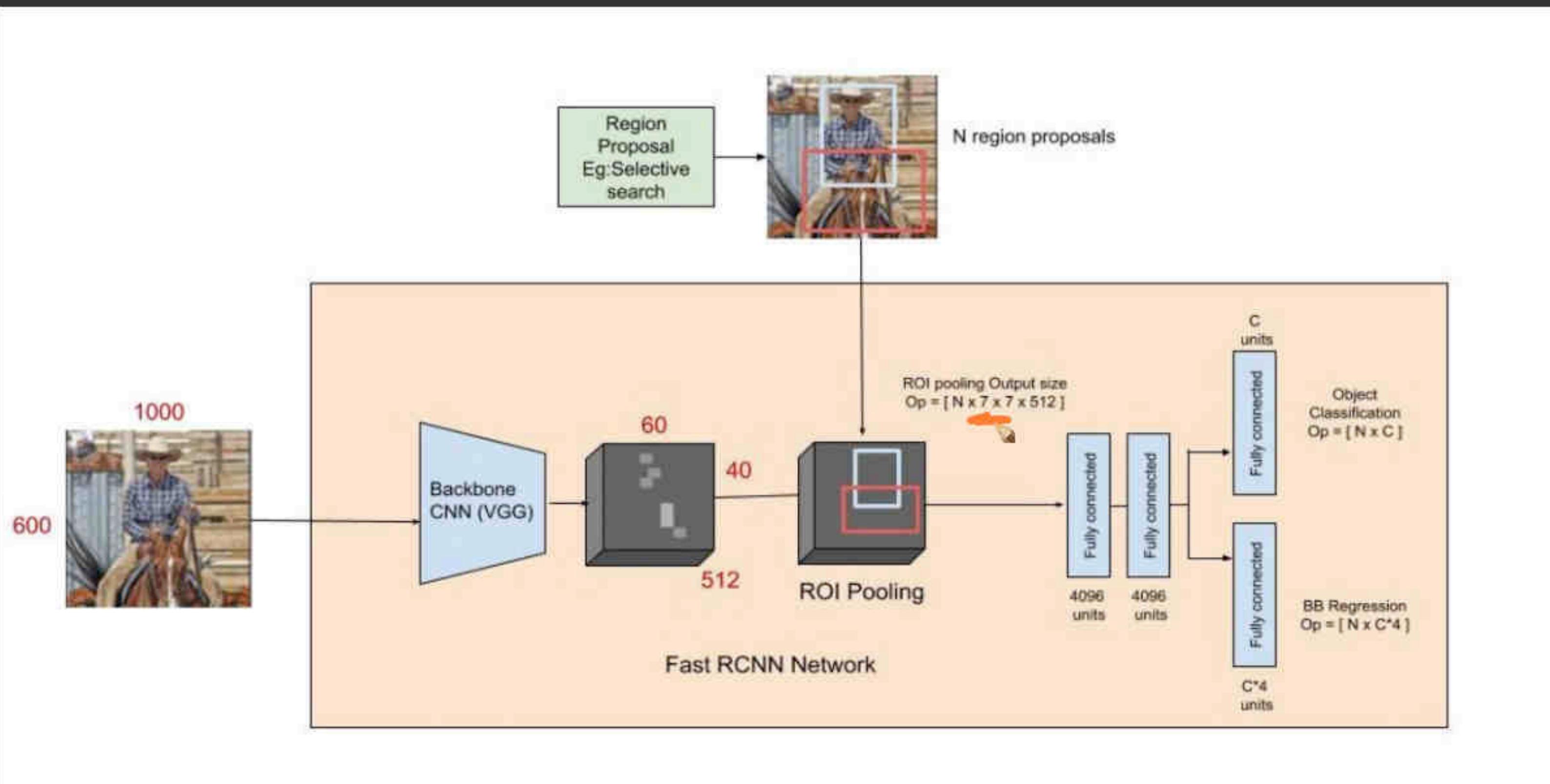


The primary benefit here is that the network is now, effectively, end-to-end trainable:

1. We input an image and associated ground-truth bounding boxes.
2. Extract the feature map
3. Apply ROI pooling and obtain the ROI feature vector.
4. And finally use two sets of fully-connected layers

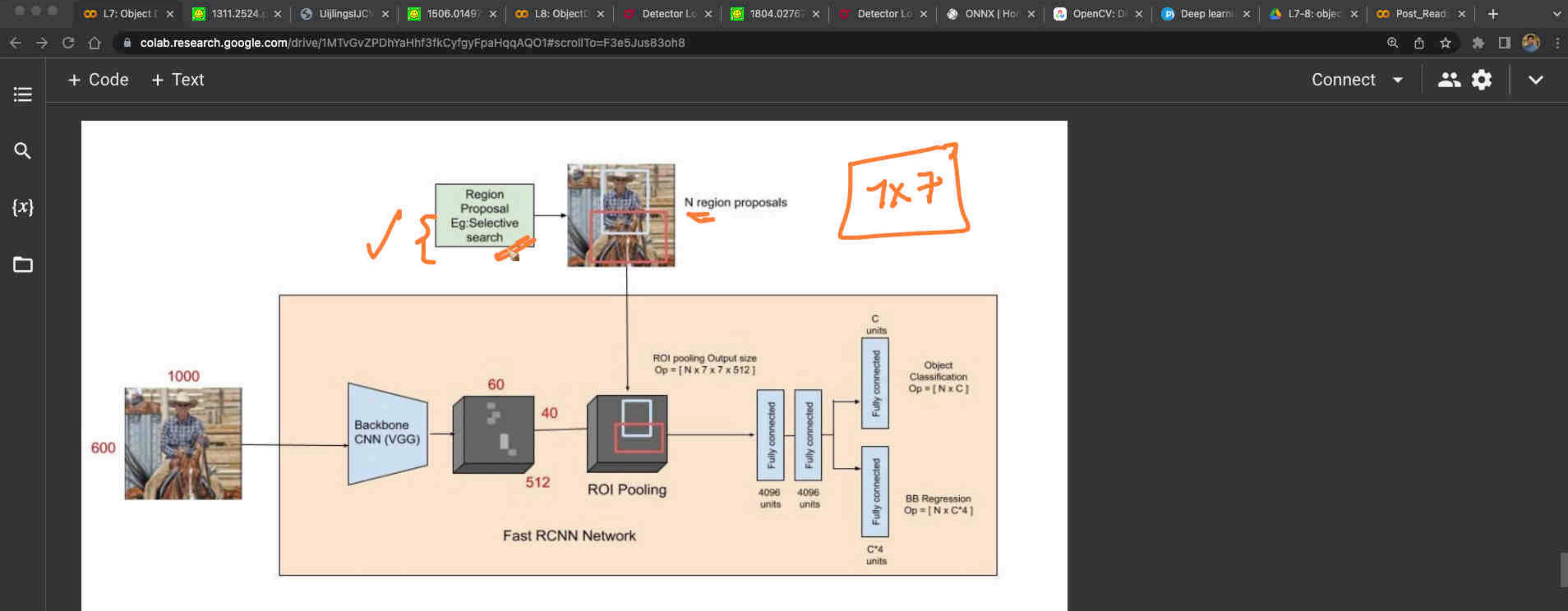
+ Code + Text

Connect



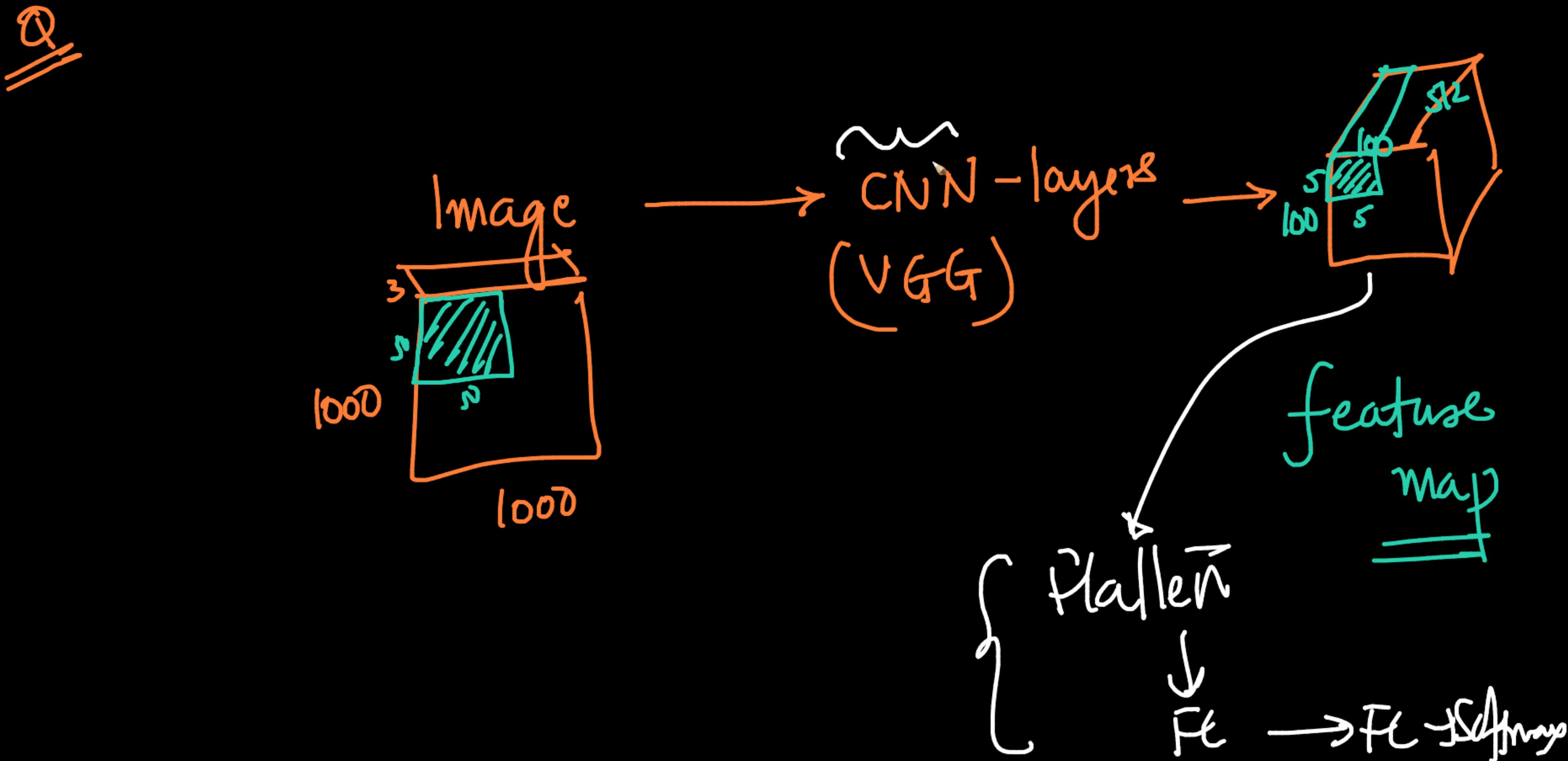
The primary benefit here is that the network is now, effectively, end-to-end trainable:

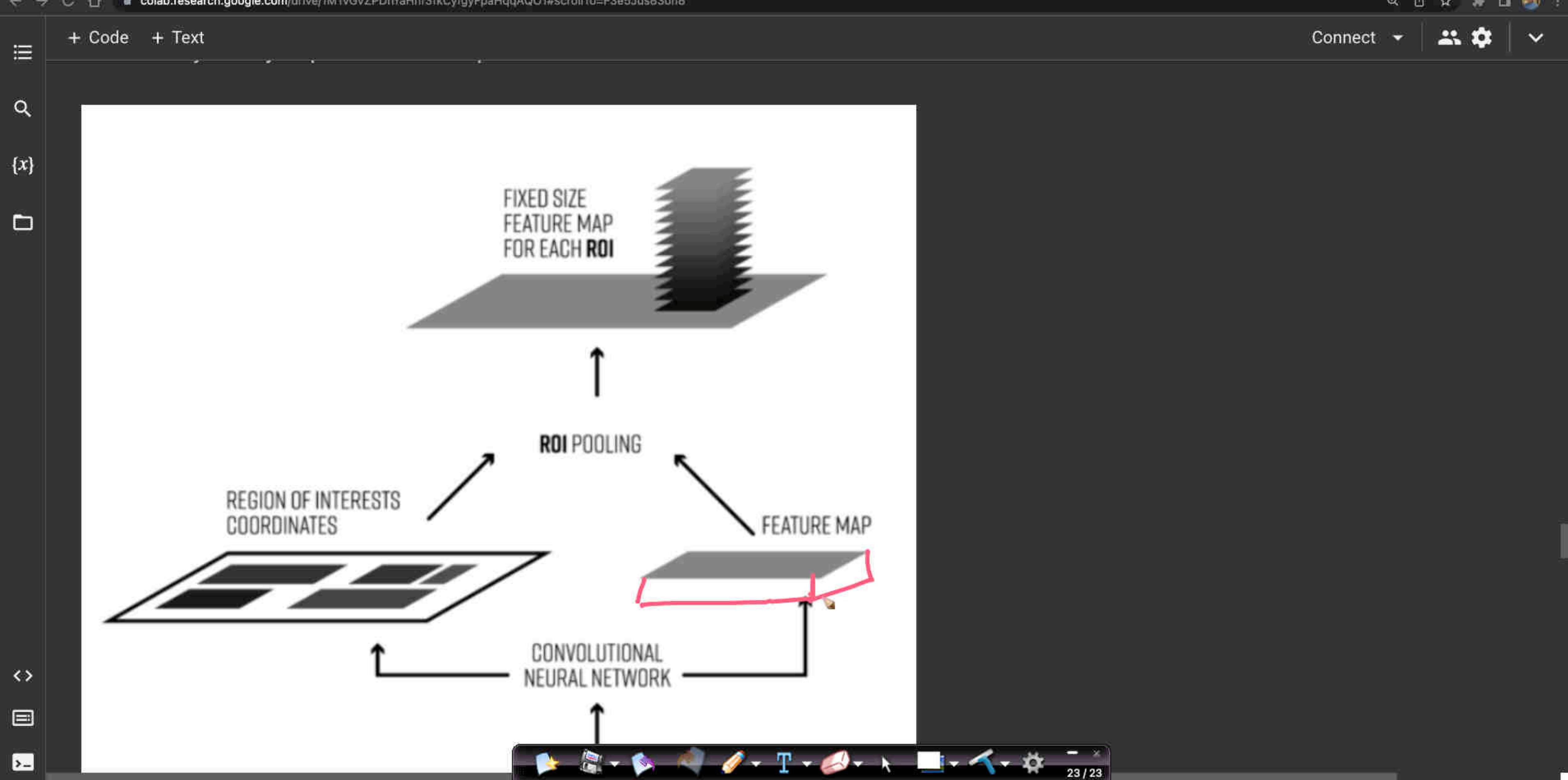
1. We input an image and associated ground-truth bounding boxes.
2. Extract the feature map
3. Apply ROI pooling and obtain the ROI feature vector.
4. And finally use two sets of fully-connected layers



The primary benefit here is that the network is now, effectively, end-to-end trainable:

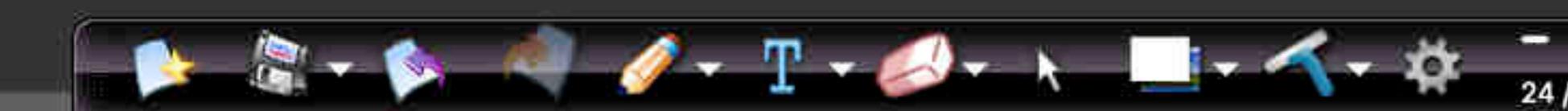
1. We input an image and associated ground-truth bounding boxes.
2. Extract the feature map
3. Apply ROI pooling and obtain the ROI feature vector.
4. And finally use two sets of fully-connected layers







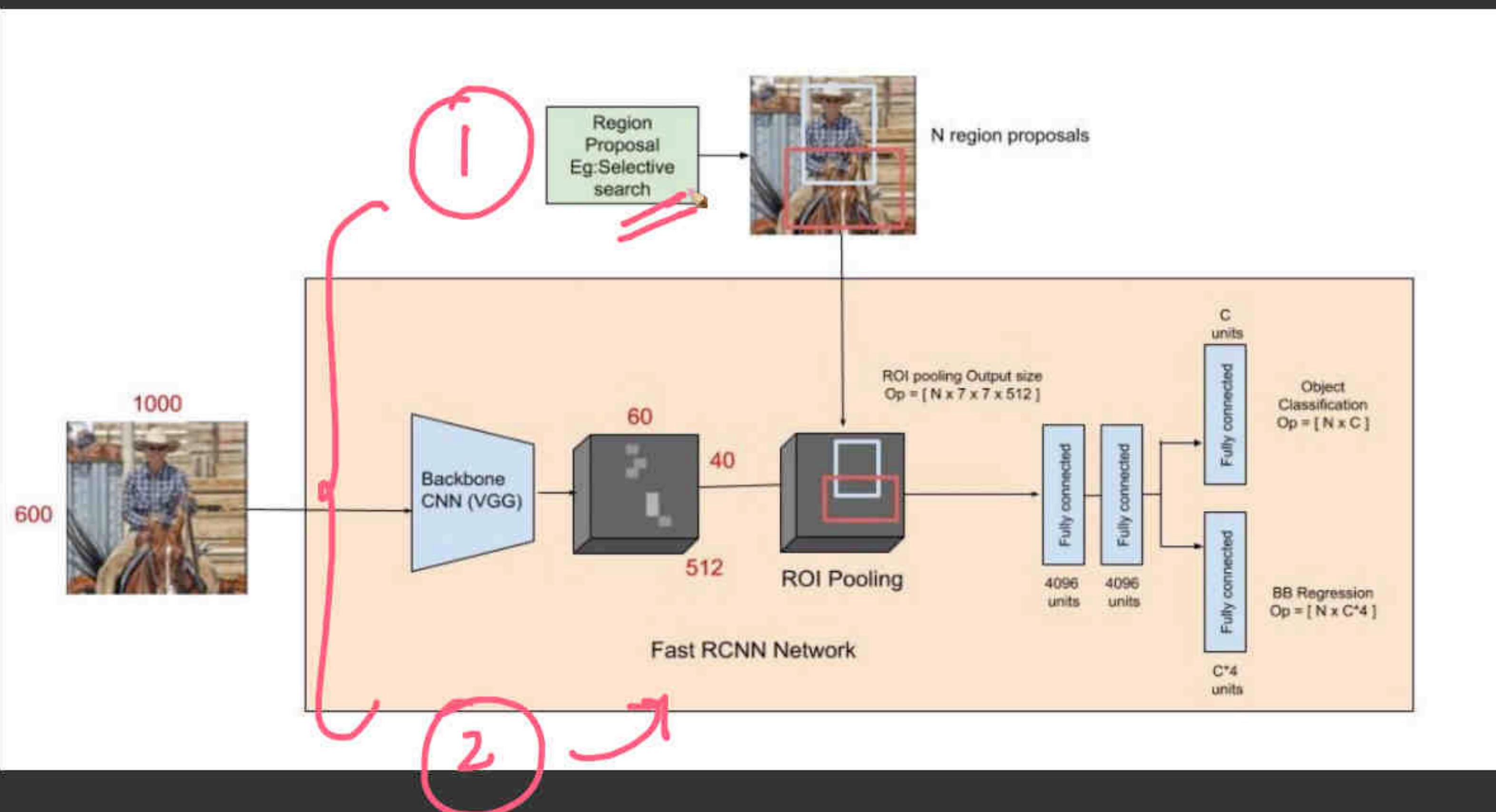
To Include: Corresponding GUN Image



+ Code + Text

Connect |

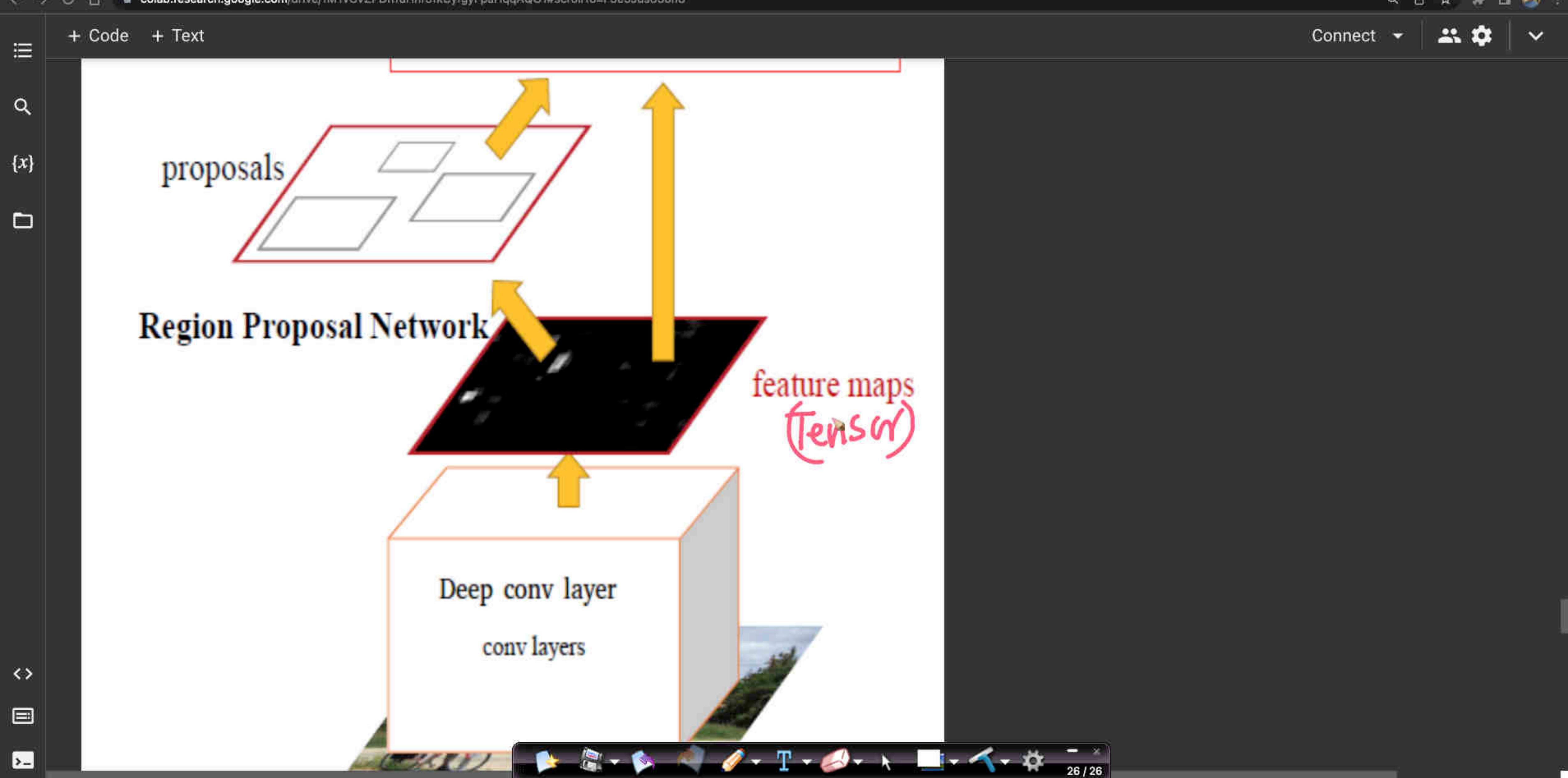
- We'll discuss ROI Pooling in more detail later, but for the time being, understand that **ROI Pooling operates over the feature map extracted from the CNN and extracts a fixed-size window from it.**

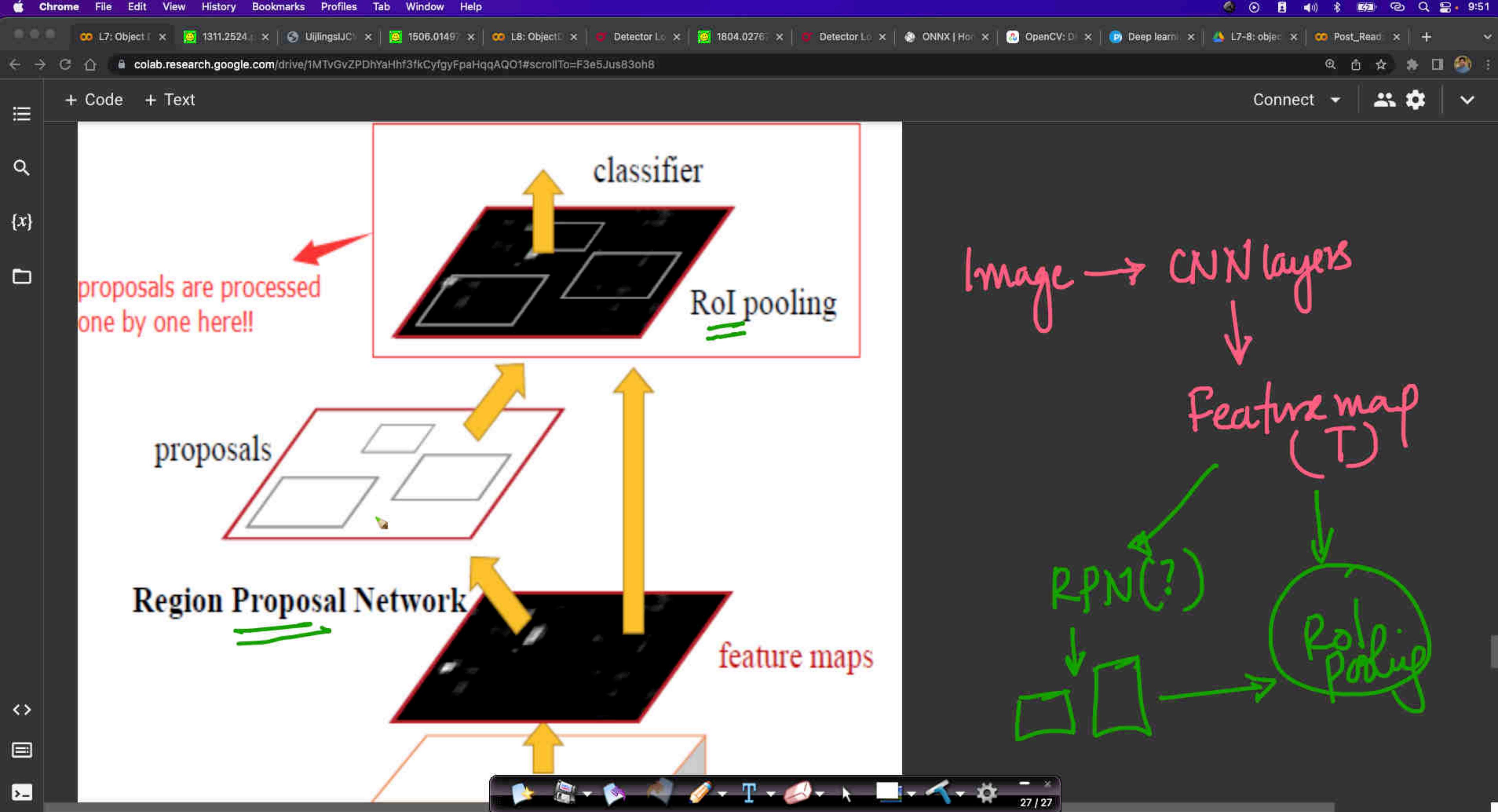


The primary benefit here is that the network is now, effectively, end-to-end trainable:

1. We input an image and associated ground truth

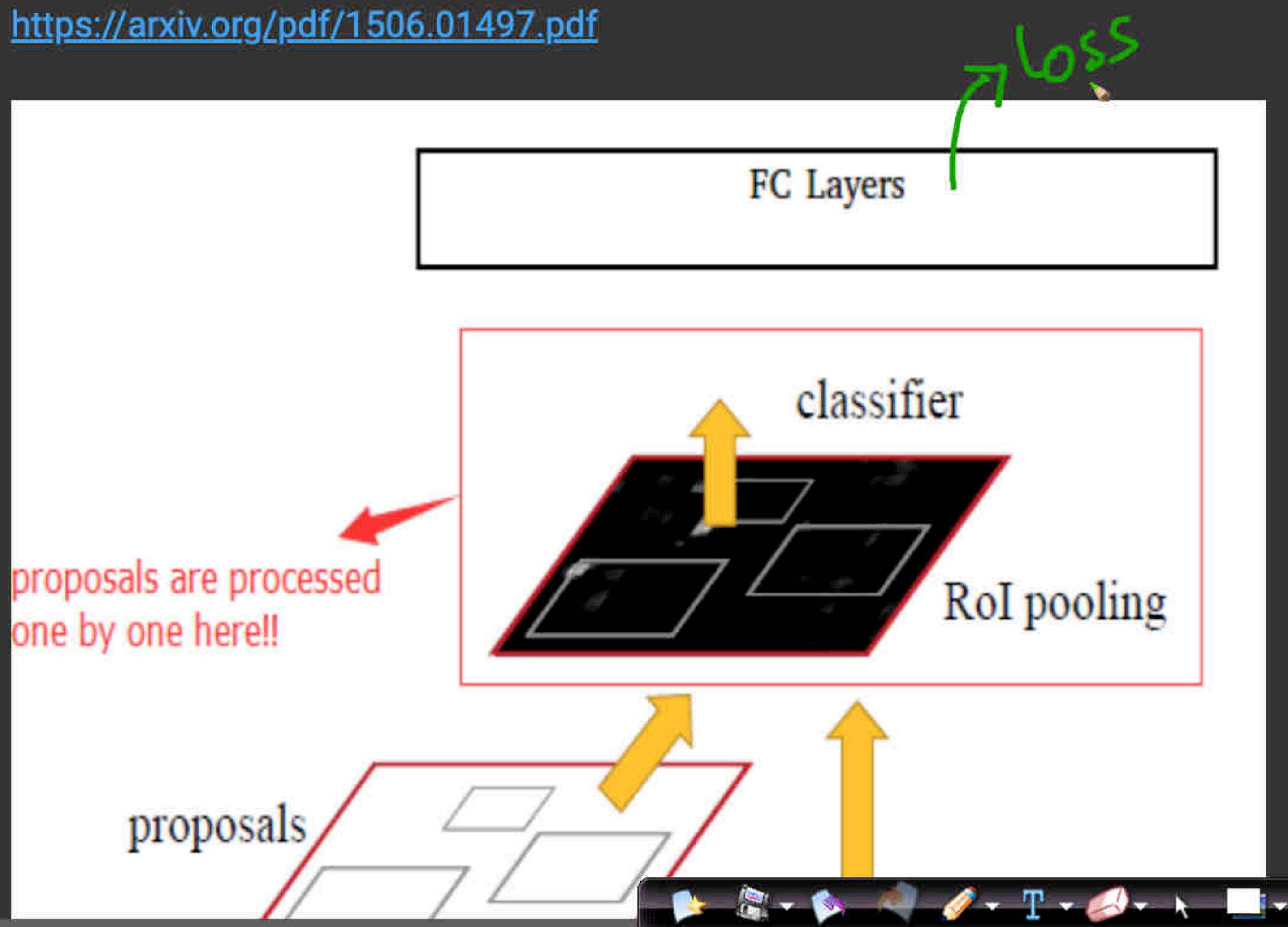


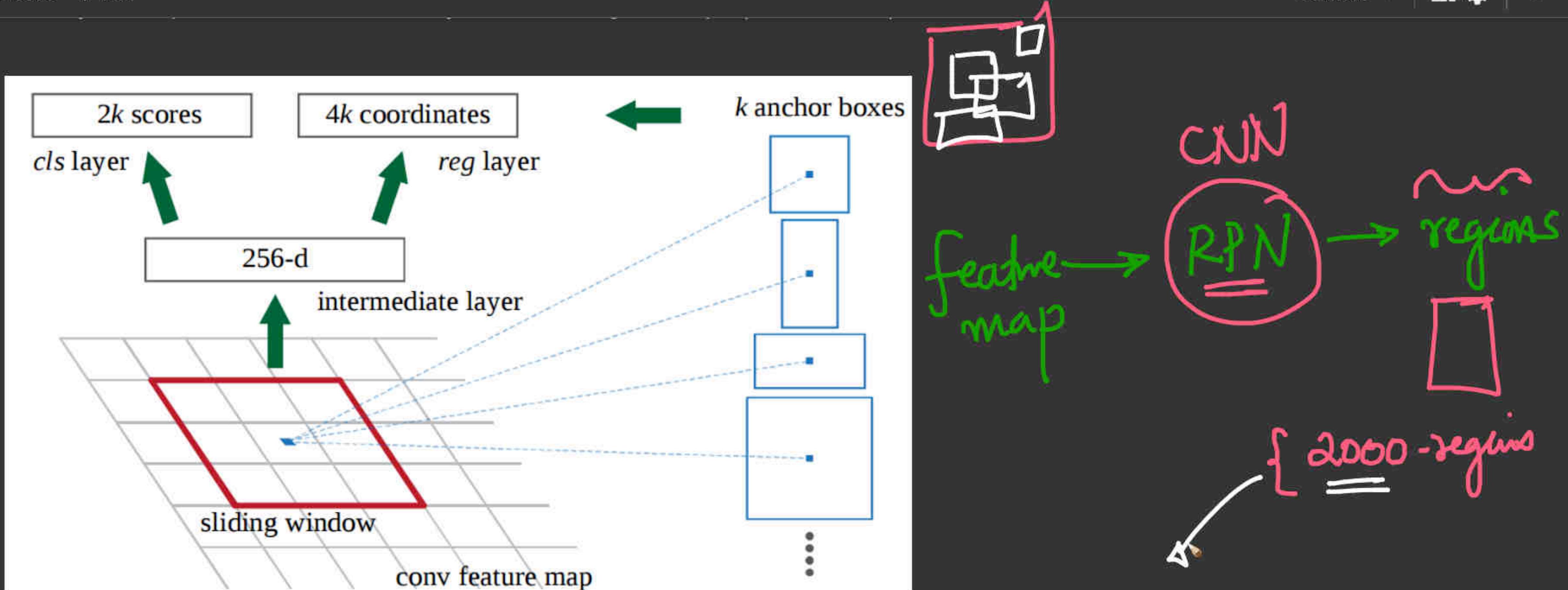




Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks

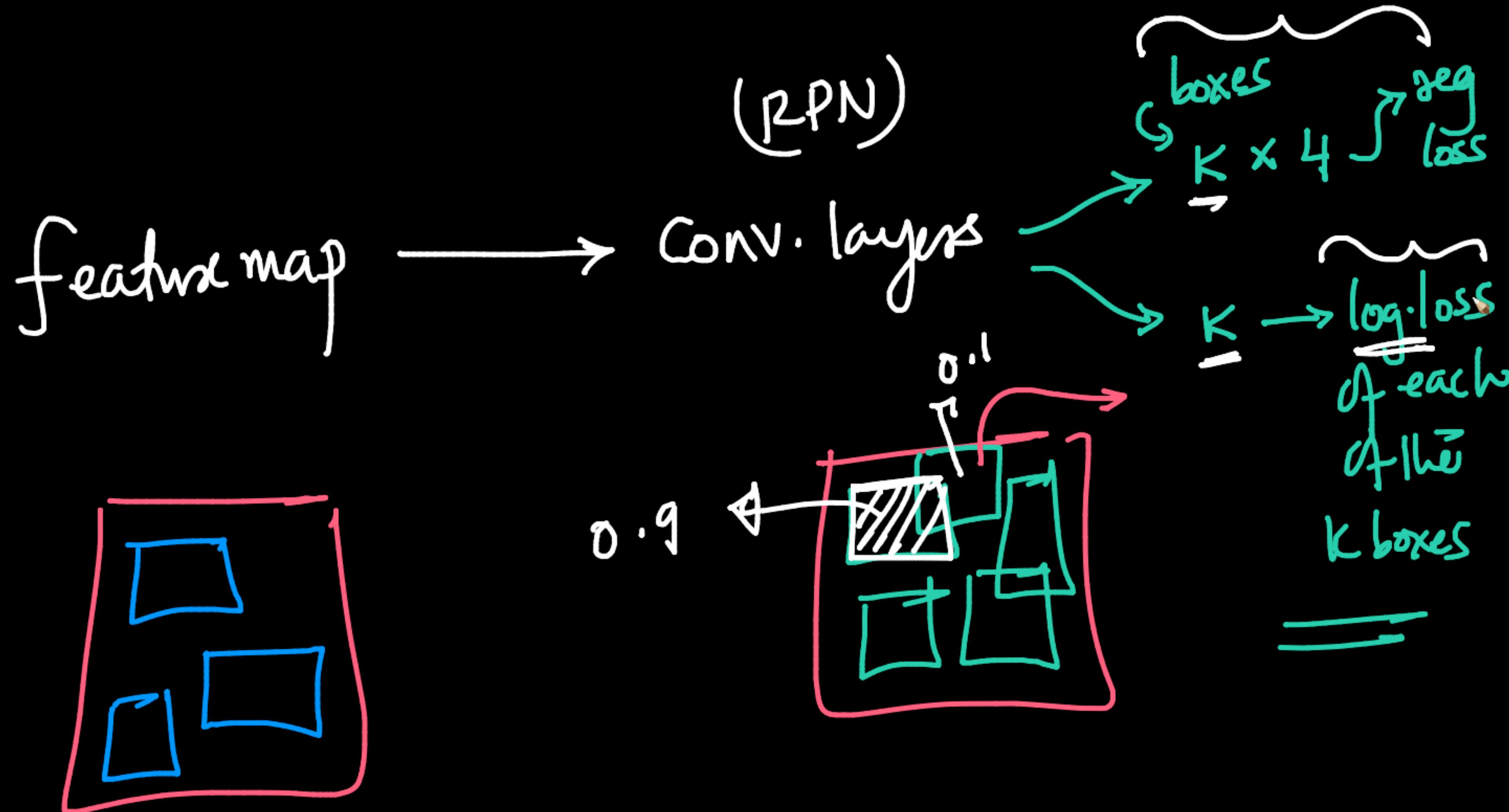
<https://arxiv.org/pdf/1506.01497.pdf>





<> Loss function of Regional Proposal Network is the sum of classification (cls) and regression (reg) loss.

<< Check Equation (1) in the research paper: <https://arxiv.org/pdf/1506.01497.pdf> >>



second condition may find no positive sample. We assign a negative label to a non-positive anchor if its IoU ratio is lower than 0.3 for all ground-truth boxes. Anchors that are neither positive nor negative do not contribute to the training objective.

With these definitions, we minimize an objective function following the multi-task loss in Fast R-CNN [2]. Our loss function for an image is defined as:

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*) \quad (1)$$

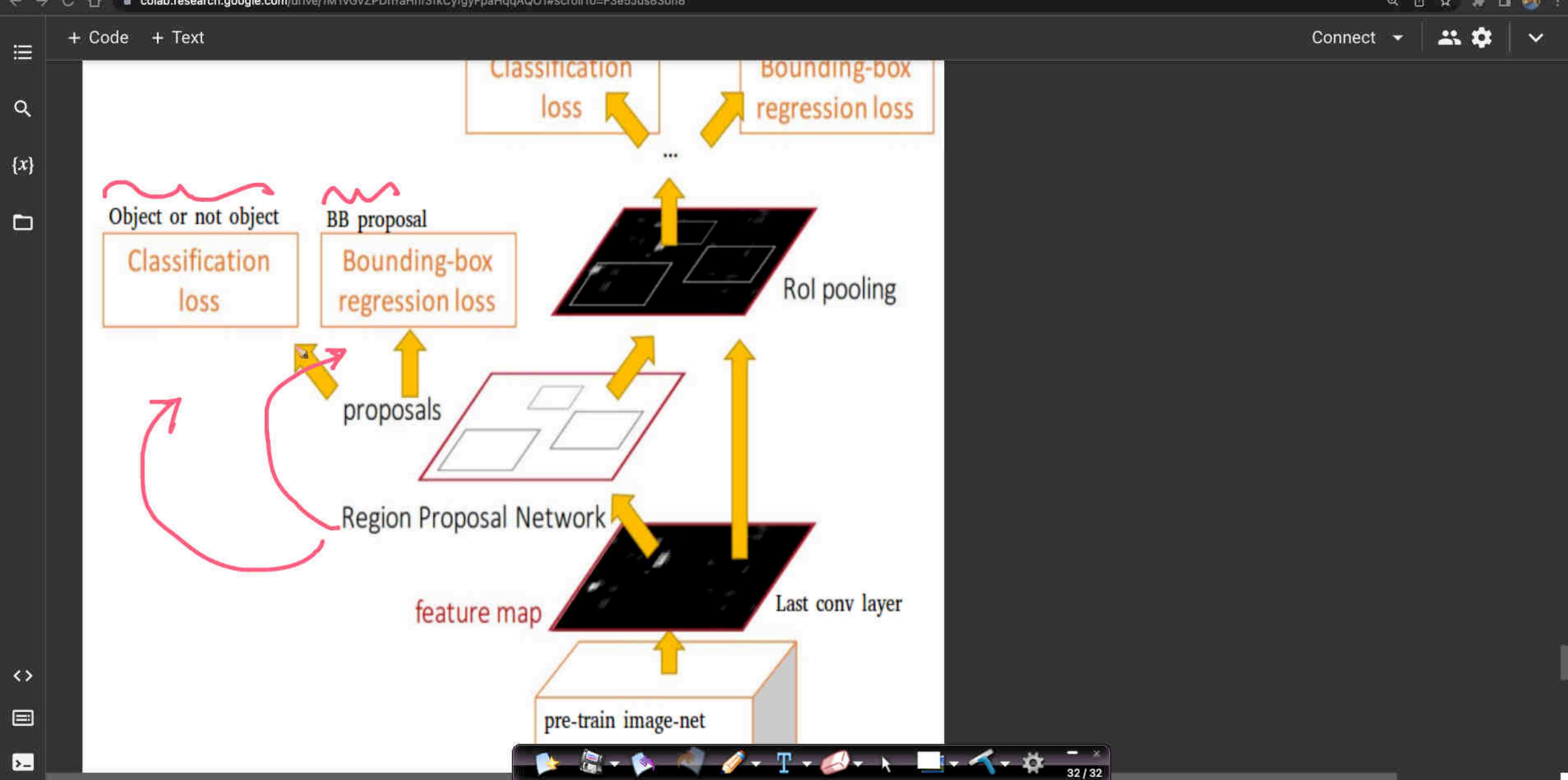
Here, i is the index of an anchor in a mini-batch and p_i is the predicted probability of anchor i being an object. The ground-truth label p_i^* is 1 if the anchor is positive, and is 0 if the anchor is negative. t_i is a vector representing the 4 parameterized coordinates of the predicted bounding box, and t_i^* is that of the ground-truth box associated with a positive anchor. The classification loss L_{cls} is $\log(p_i^*) - p_i$ (object vs. not object). For the regression loss, we use

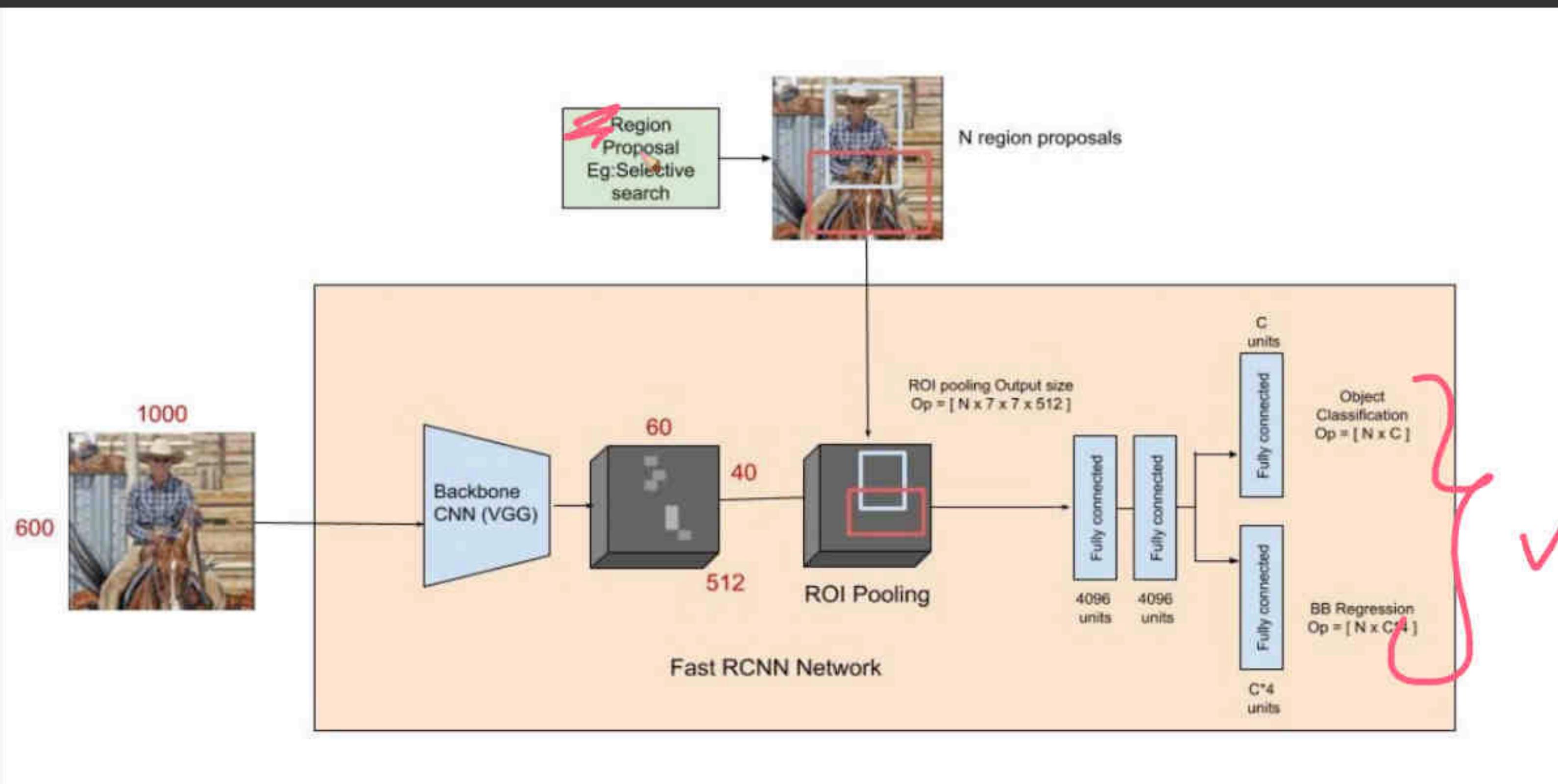
[2], bounding-box regression is performed on features pooled from *arbitrarily* sized RoIs, and the regression weights are *shared* by all region sizes. In our formulation, the features used for regression are of the *same* spatial size (3×3) on the feature maps. To account for varying sizes, a set of k bounding-box regressors are learned. Each regressor is responsible for one scale and one aspect ratio, and the k regressors do *not* share weights. As such, it is still possible to predict boxes of various sizes even though the features are of a fixed size/scale, thanks to the design of anchors.

3.1.3 Training RPNs

The RPN can be trained end-to-end by back-propagation and stochastic gradient descent (SGD) [35]. We follow the “image-centric” sampling strategy from [2] to train this network. Each mini-batch arises from a single image that contains many positive and negative example anchors. It is possible to optimize for the loss functions of all anchors, but this will bias towards negative samples as they are dominate.

For example 256 anchors in an image to compute the loss function of a mini-batch, where

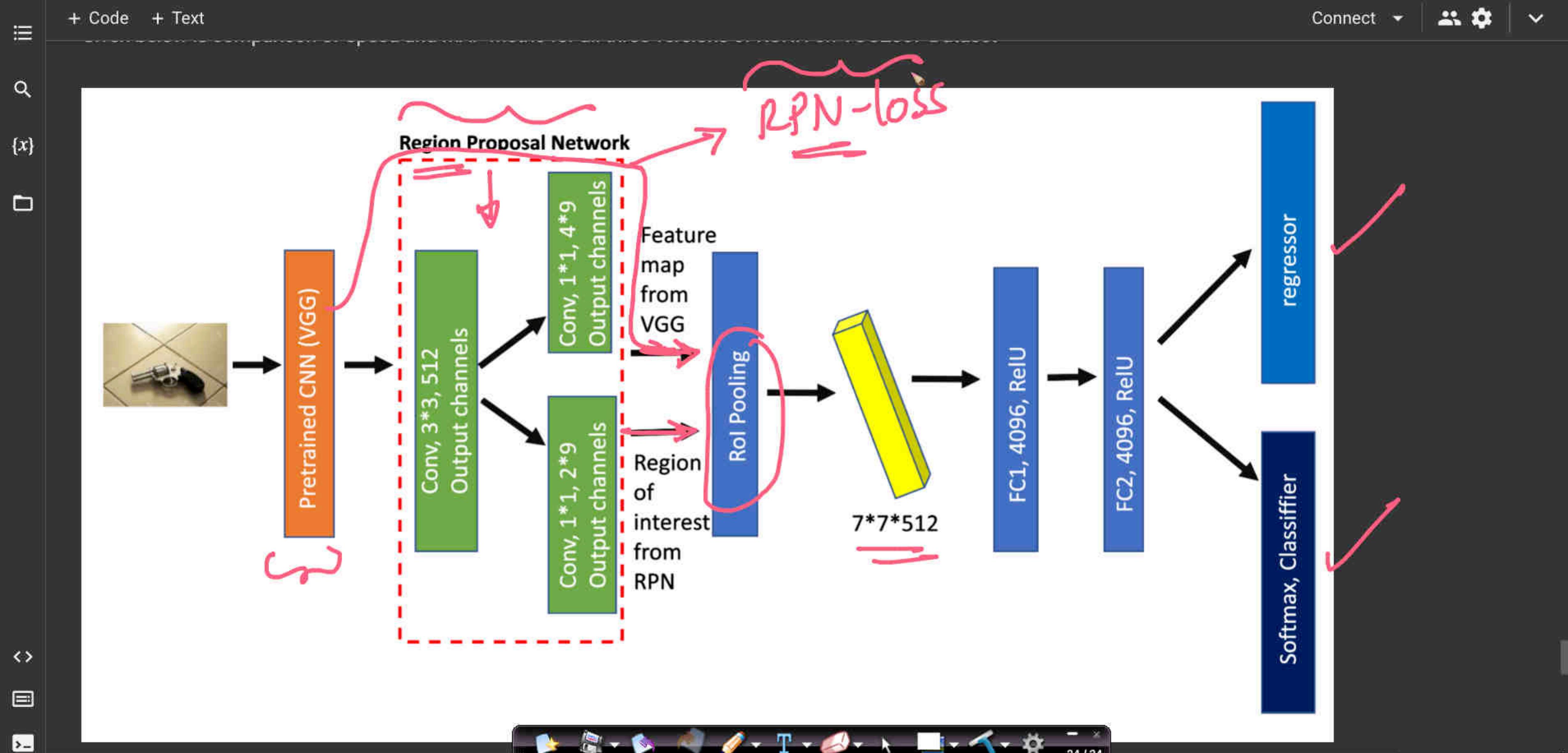


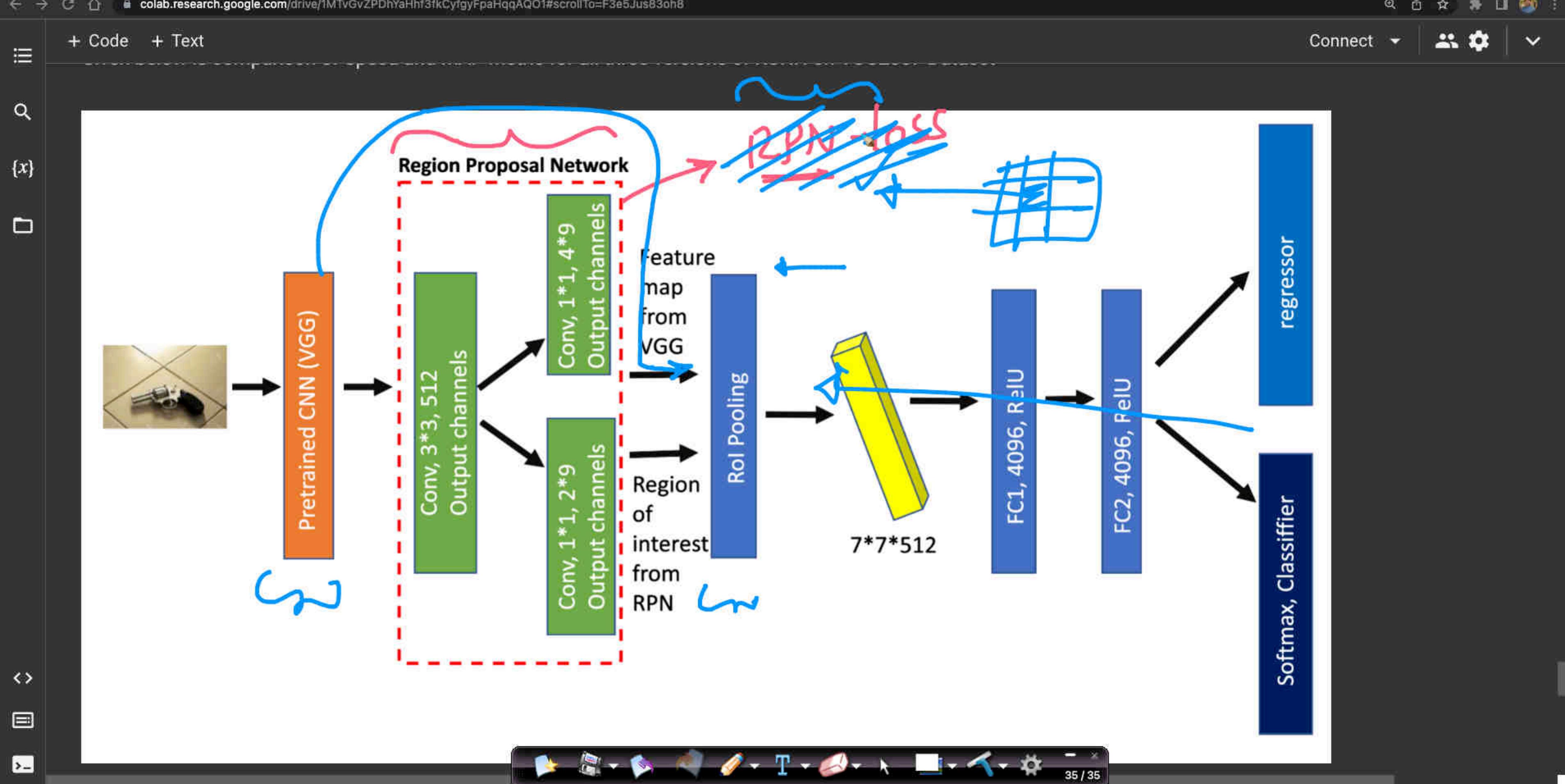


Faster RCNN

The primary benefit here is that the network is now, effectively, end-to-end trainable:

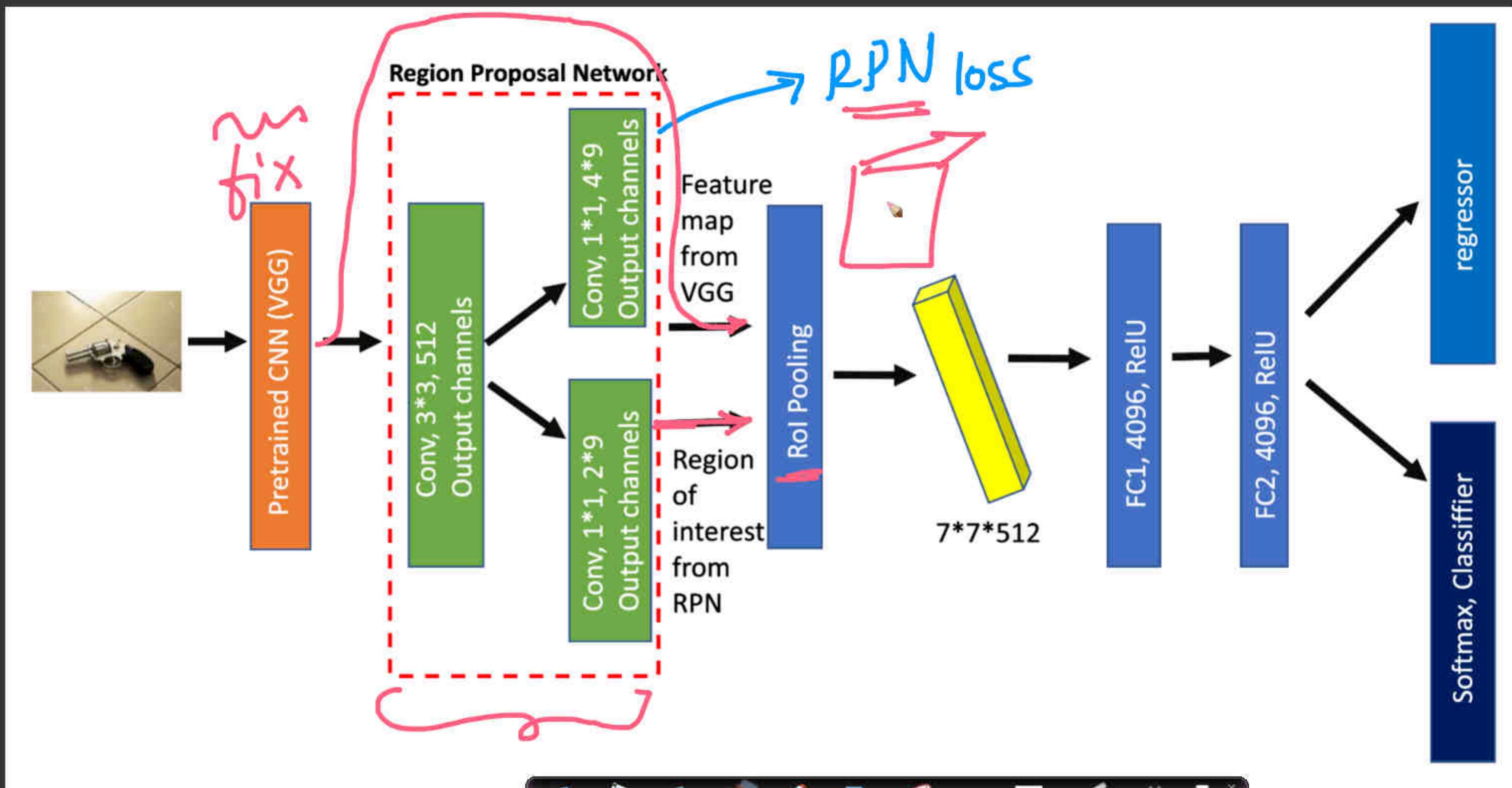
1. We input an image and associated ground-truth bounding boxes.
2. Extract the feature map
3. Apply ROI pooling and obtain the ROI feature vector



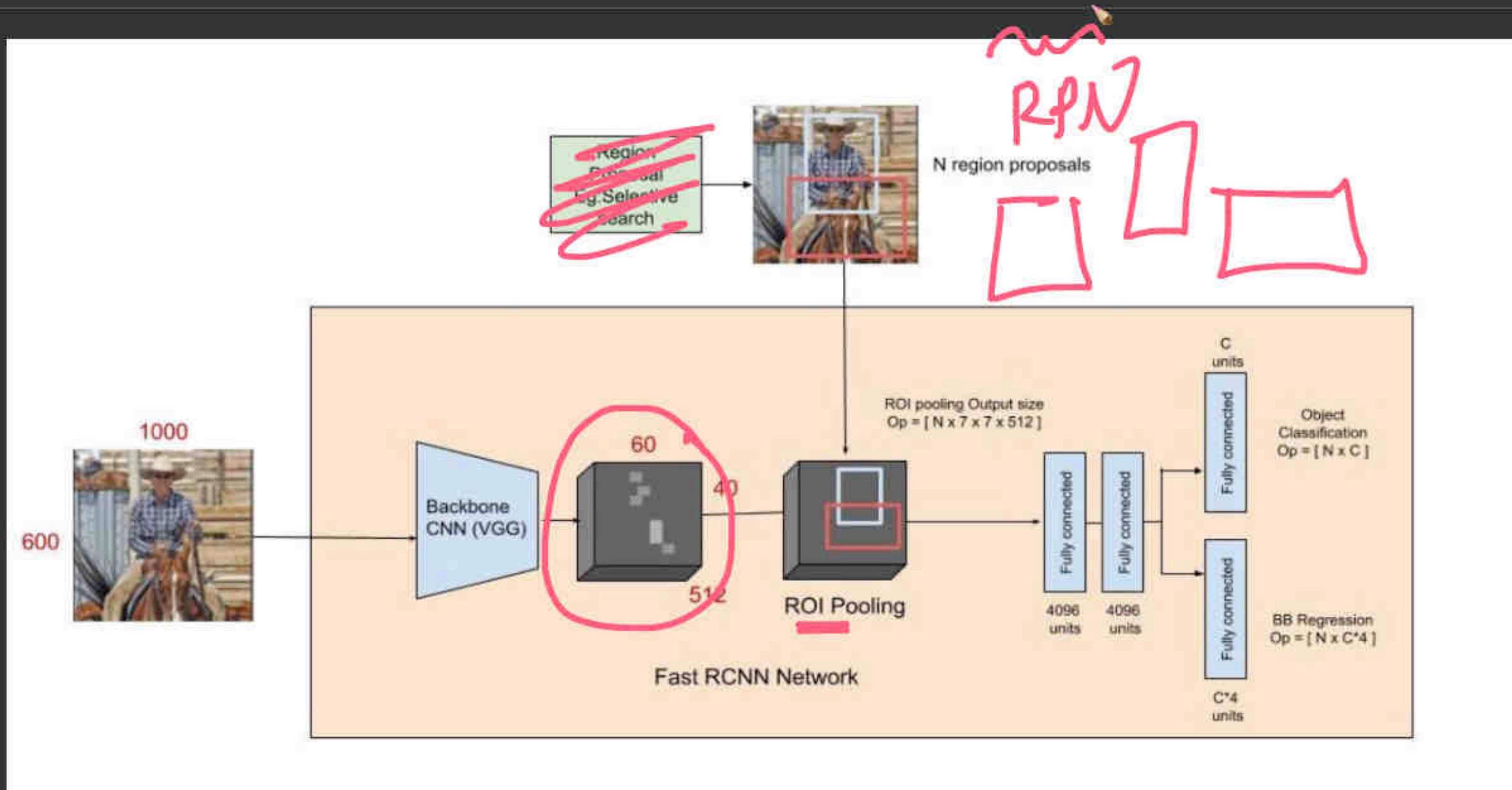


+ Code + Text

Connect



+ Code + Text



The primary benefit here is that the network is now, effectively, end-to-end trainable:

1. we input an image and associated ground-truth bounding boxes.
2. Extract the feature map
3. Apply ROI pooling and obtain the ROI feature vector
4. And finally use two sets of fully-connected layers

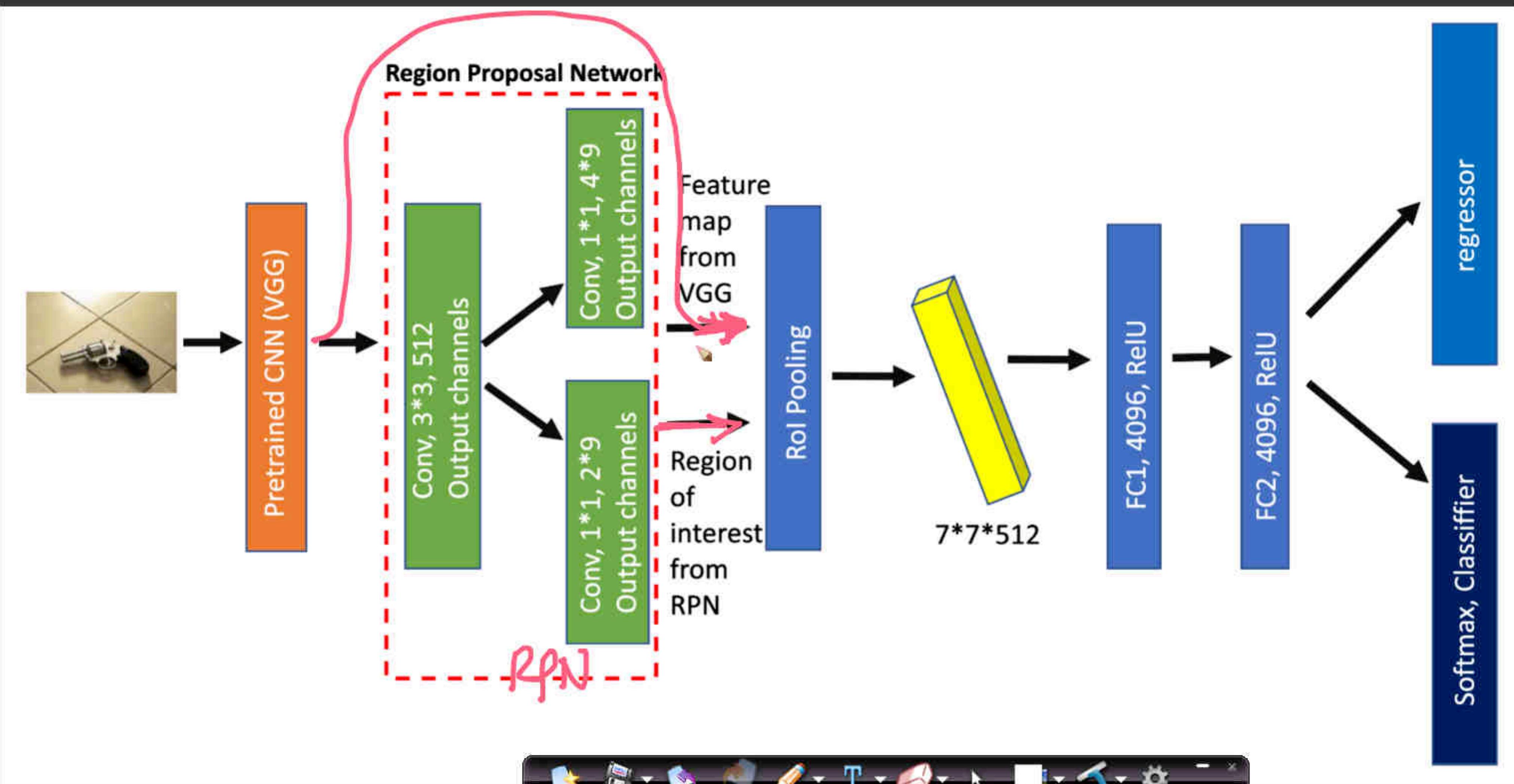
iii

+ Code + Text

Connect



Given below is comparison of Speed and mAP metric for all three versions of RCNN on VOC2007 Dataset



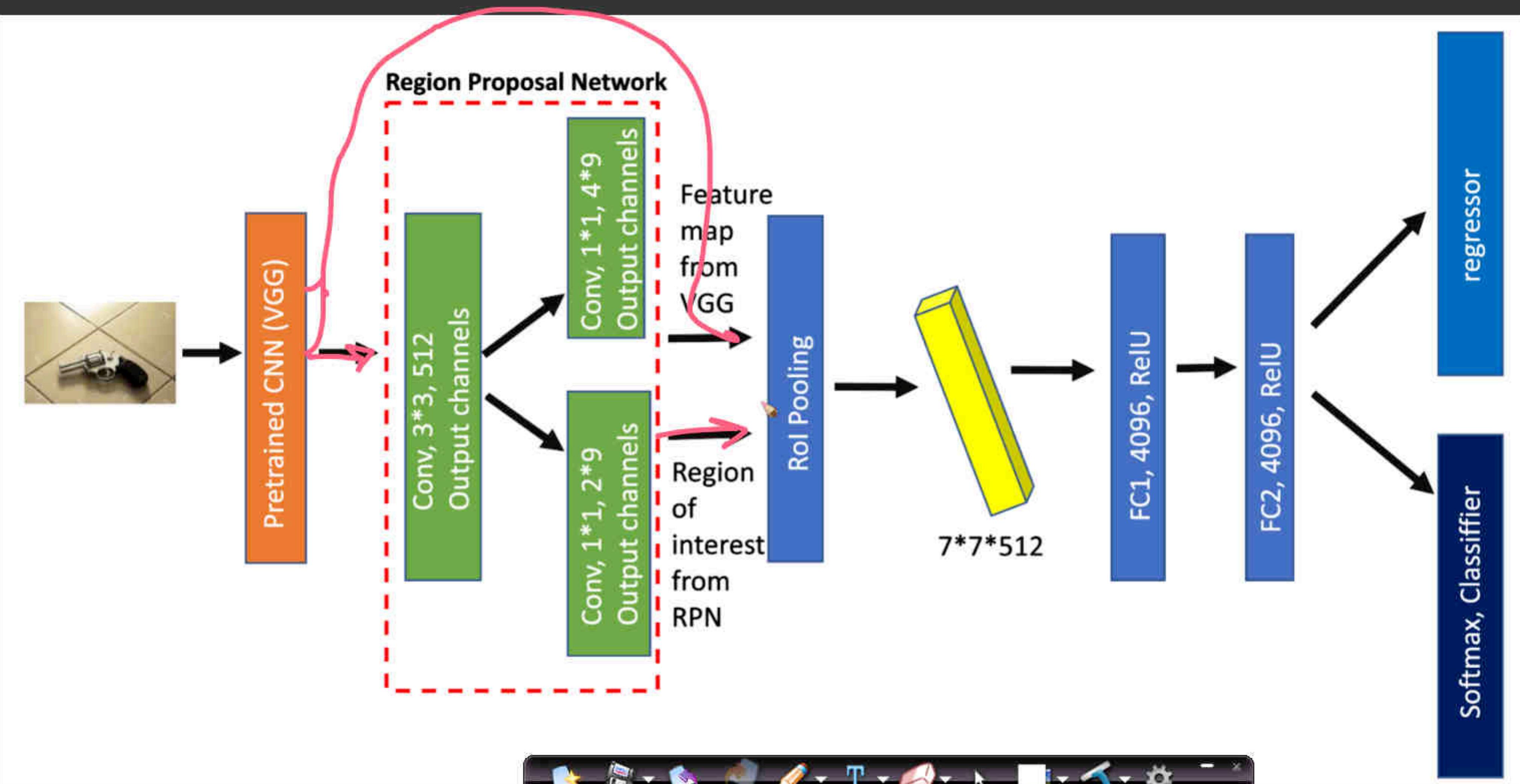
iii

+ Code + Text

Connect

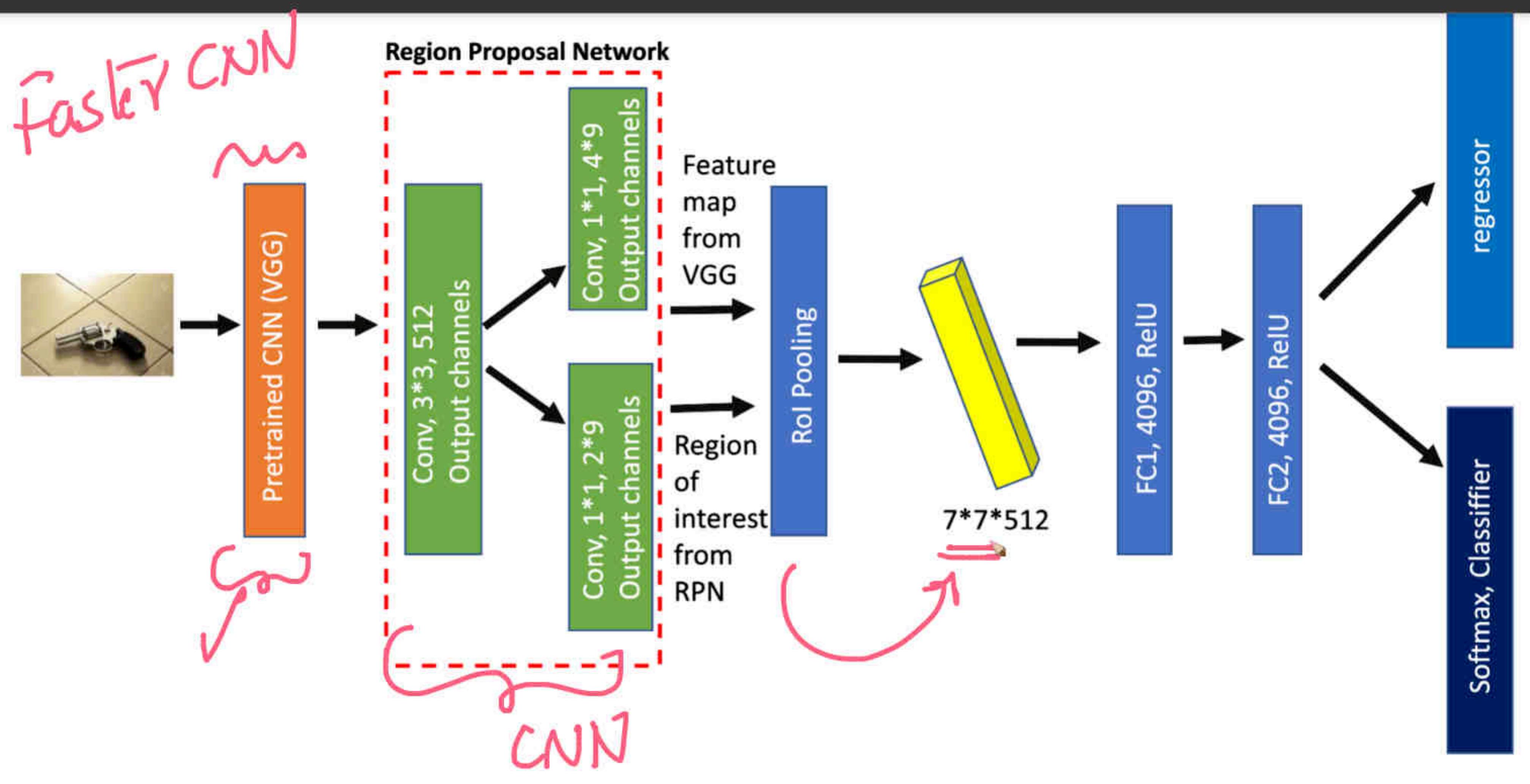


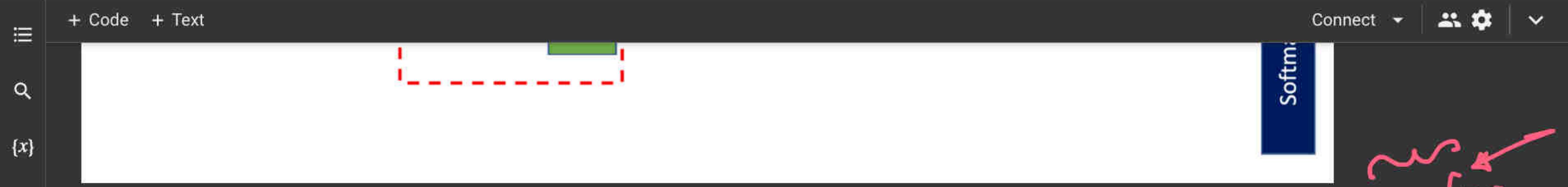
Given below is comparison of Speed and mAP metric for all three versions of RCNN on VOC2007 Dataset



+ Code + Text

Connect





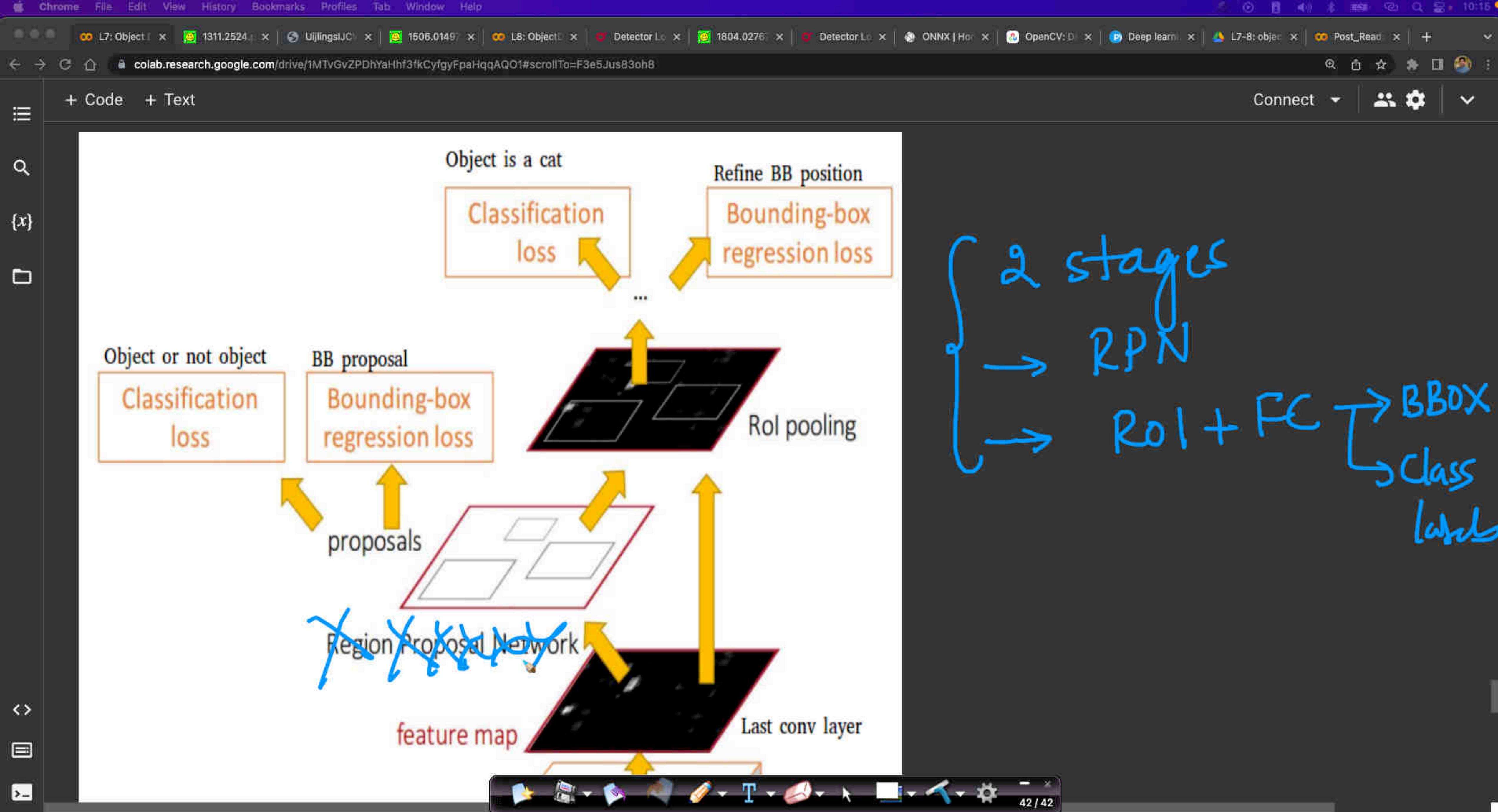
| External region proposals method (Selective Search Algorithm) | | ✓ | ✓ | ✗ |
|--|------------|------------|--------------|------------------|
| | R-CNN | Fast R-CNN | Faster R-CNN | |
| Test time per image | 50 seconds | 2 seconds | 0.2 seconds | { surveillance |
| Speed-up | 1x | 25x | 250x | = 5 |
| mAP (VOC 2007) | 66.0% | 66.9% | 66.9% | self-driving car |

▼ [Below part is Optional]

Code Implementation of Backbone for FasterRCNN(VGG:

$\frac{1}{30} \text{ Sec}$

$= 0.033 \text{ Sec}$



Chrome File Edit View History Bookmarks Profiles Tab Window Help

L7: Object... x 1311.2524 x DijlingsJC x 1506.0149 x L8: Object... x Detector L... x 1804.0276 x DetectorL... x ONNX|Ho... x OpenCV:D... x Deep learn... x L7-8: objec... x Post_Reader x +

colab.research.google.com/drive/1MTvGvZPDhYaHhf3fkCyfgyFpaHqqAQO1#scrollTo=F3e5Jus83oh8

+ Code + Text Connect |

```
x = Conv2D(512, (3, 3), activation='relu', padding='same', name='block5_conv3')(x)
# x = MaxPooling2D((2, 2), strides=(2, 2), name='block5_pool')(x)

# We are not using fully connected layers (3 fc layers) as we need feature maps as output from this network.

{x}
return x
```

Code Implementation OF RPN:

```
[ ] """
#RPN layer
def rpn_layer(base_layers, num_anchors):

    #cnn_used for creating feature maps: vgg, num_anchors: 9
    x = Conv2D(512, (3, 3), padding='same', activation='relu')(base_layers)

    #classification layer: num_anchors (9) channels for 0, 1 sigmoid activation output
    x_class = Conv2D(num_anchors, (1, 1), activation='sigmoid')(x)

    #regression layer: num_anchors*4 (36) channels for computing the regression of bboxes
    x_regr = Conv2D(num_anchors * 4, (1, 1), activation='linear')(x)

    return [x_class, x_regr, base_layers] #classification of object(0 or 1),compute bounding boxes, base layers vgg
```

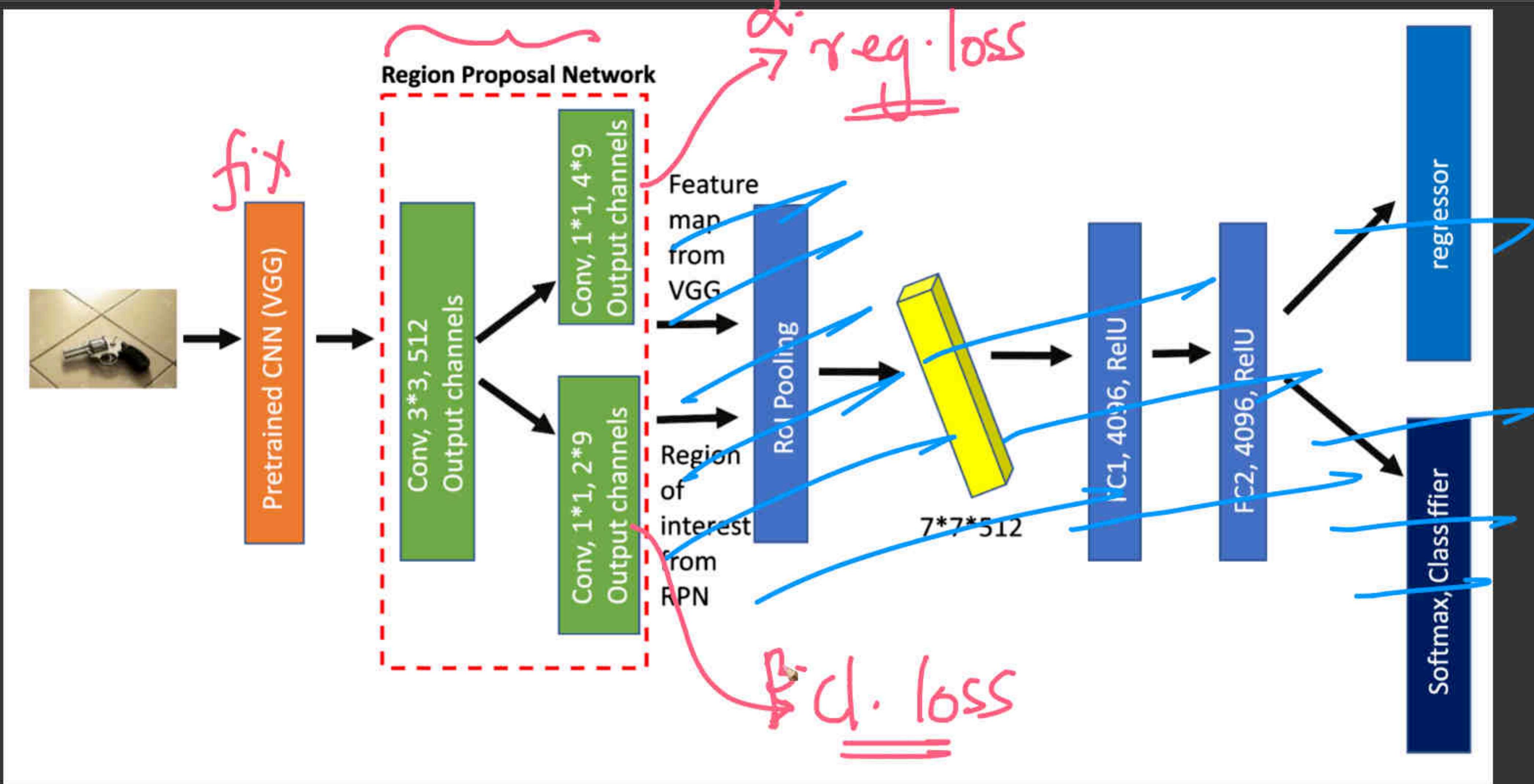
Code Implementation of ROI Pooling Layer:

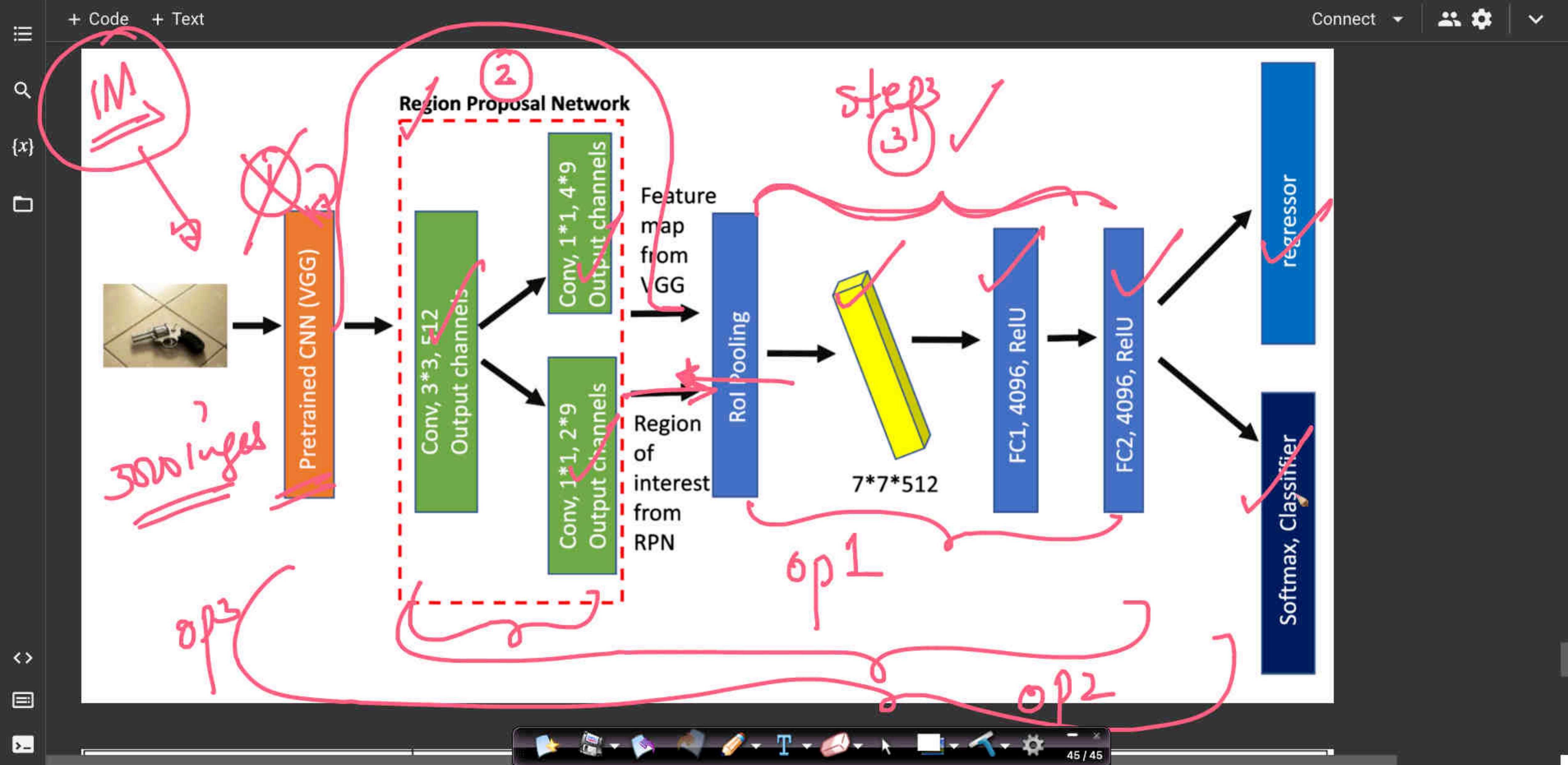
```
[ ] class RoiPoolingConv(Layer):
    """ROI pooling layer for 2D inputs
```



+ Code + Text

Connect







+ Code + Text Cannot save changes

Connect



(a)

Region Proposal Network

- These models don't have the region proposal stage **Region Proposal Network**.

There are majorly two single Stage Detection algorithms:

1. **SSD: Single Shot Detector**

✓ 2. **Yolo: You Only Look Once** (much more popular and faster)

Let's Dive Deeper into Each one of them:

► [Optional Post-READ] Introduction to SSD:

• Original paper: <https://arxiv.org/pdf/1512.02325.pdf>

• The **core idea behind the SSD network** is to have a **CNN** that takes in an image as input and produce detections at different scales, shapes, and locations.

Let's understand this in detail:

↳ 14 cells hidden

Introduction to YOLO algorithm





+ Code + Text Cannot save changes

Connect



(a)

Region Proposal Network

- These models don't have the region proposal stage **Region Proposal Network**.

There are majorly two single Stage Detection algorithms:

1. **SSD**: Single Shot Detector
2. **Yolo**: You Only Look Once (much more popular and faster)

Let's Dive Deeper into Each one of them:

▶ [Optional Post-READ] Introduction to SSD:

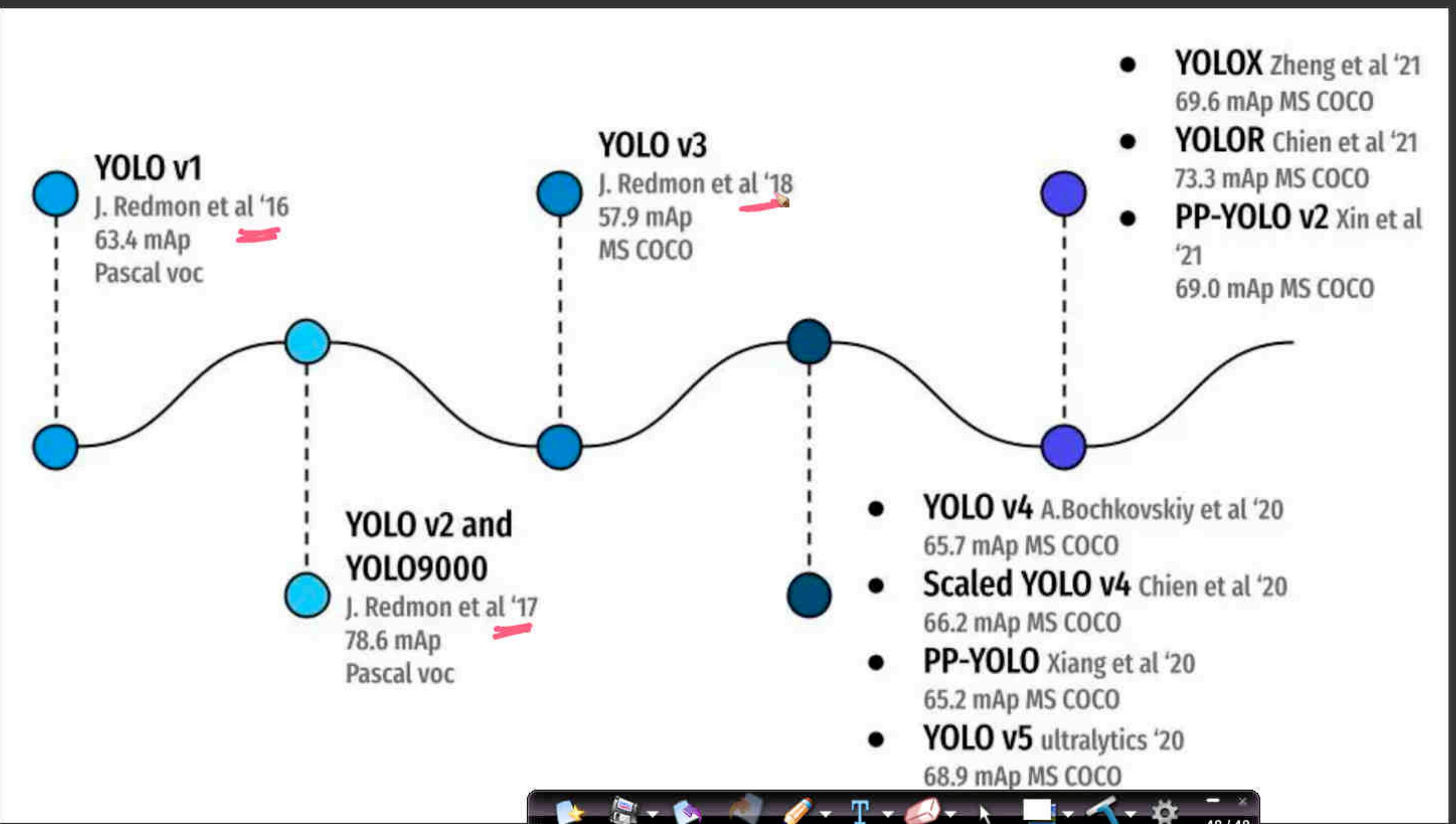
- Original paper: <https://arxiv.org/pdf/1512.02325.pdf>
- The **core idea behind the SSD network** is to have a **CNN** that takes in an image as input and produce detections at different scales, shapes, and locations.

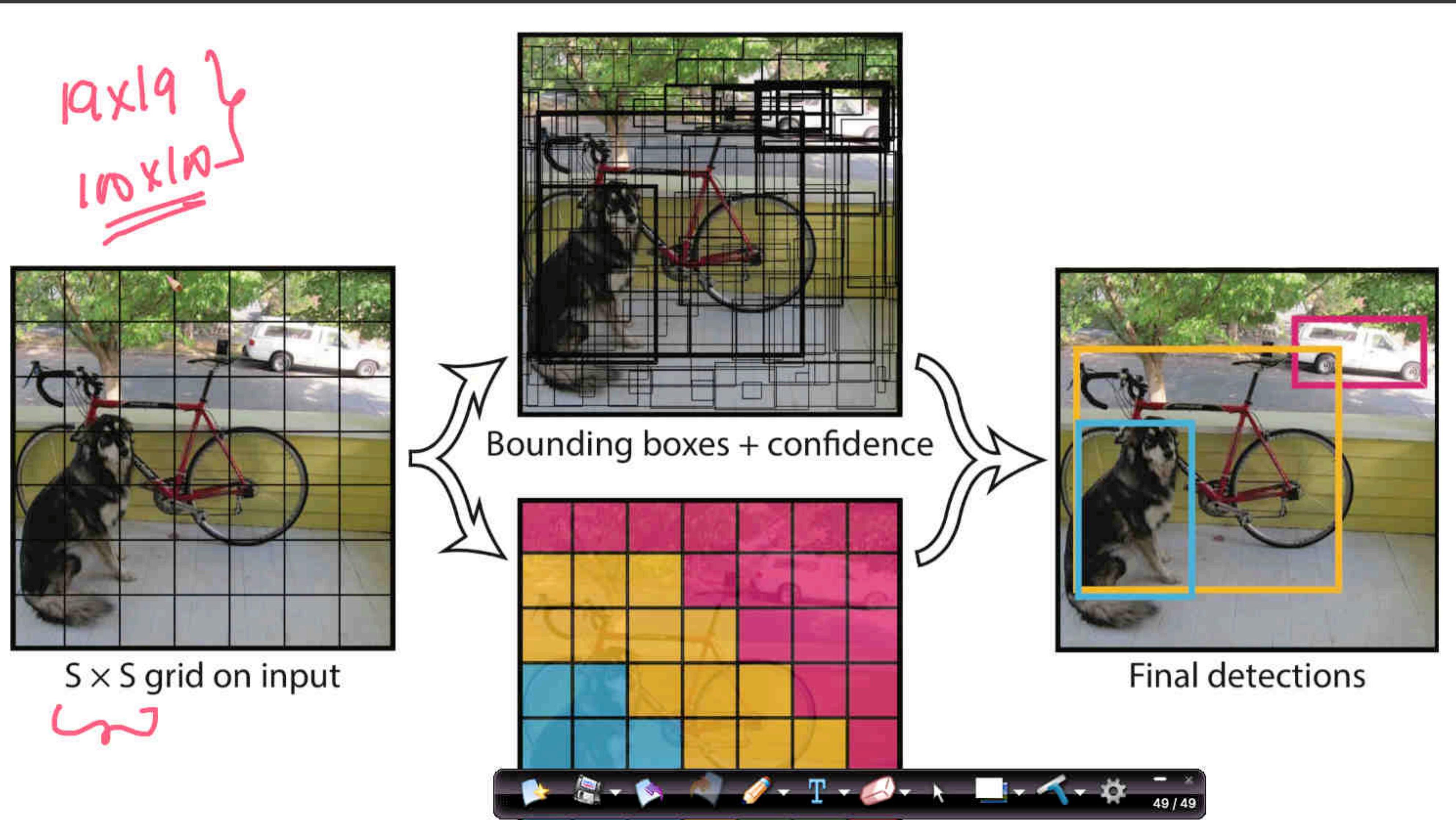
Let's understand this in detail:

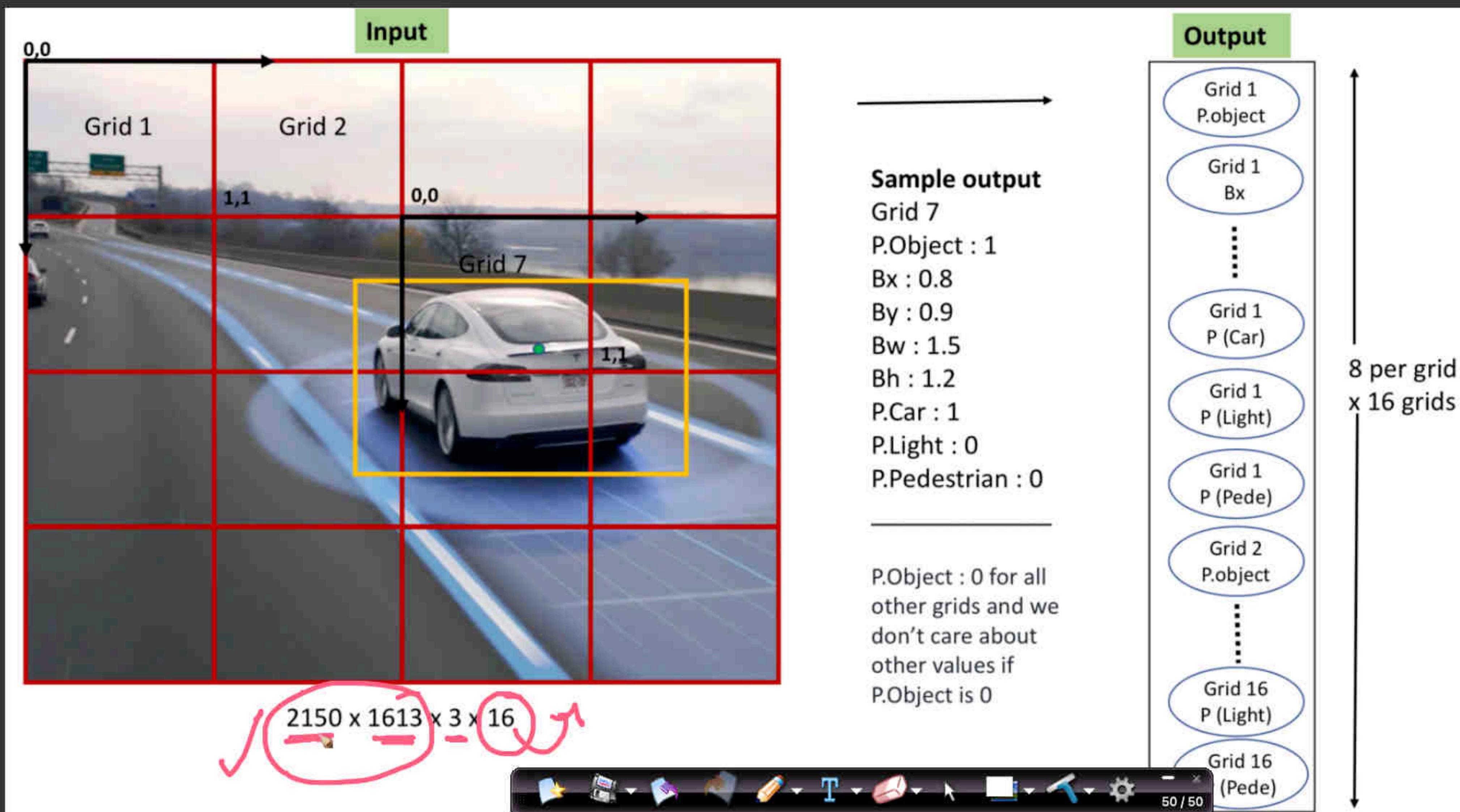
↳ 14 cells hidden

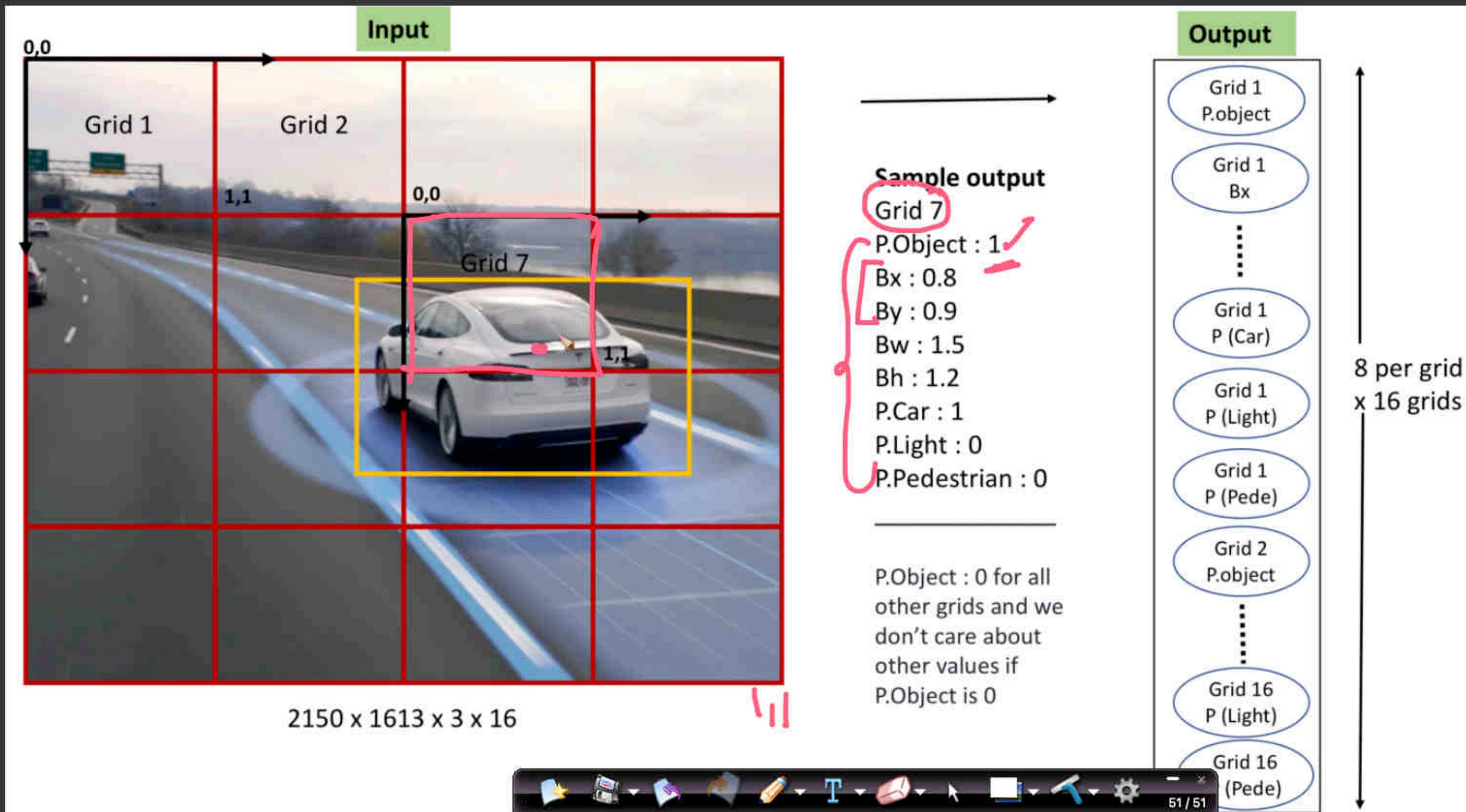
Introduction to YOLO algorithm

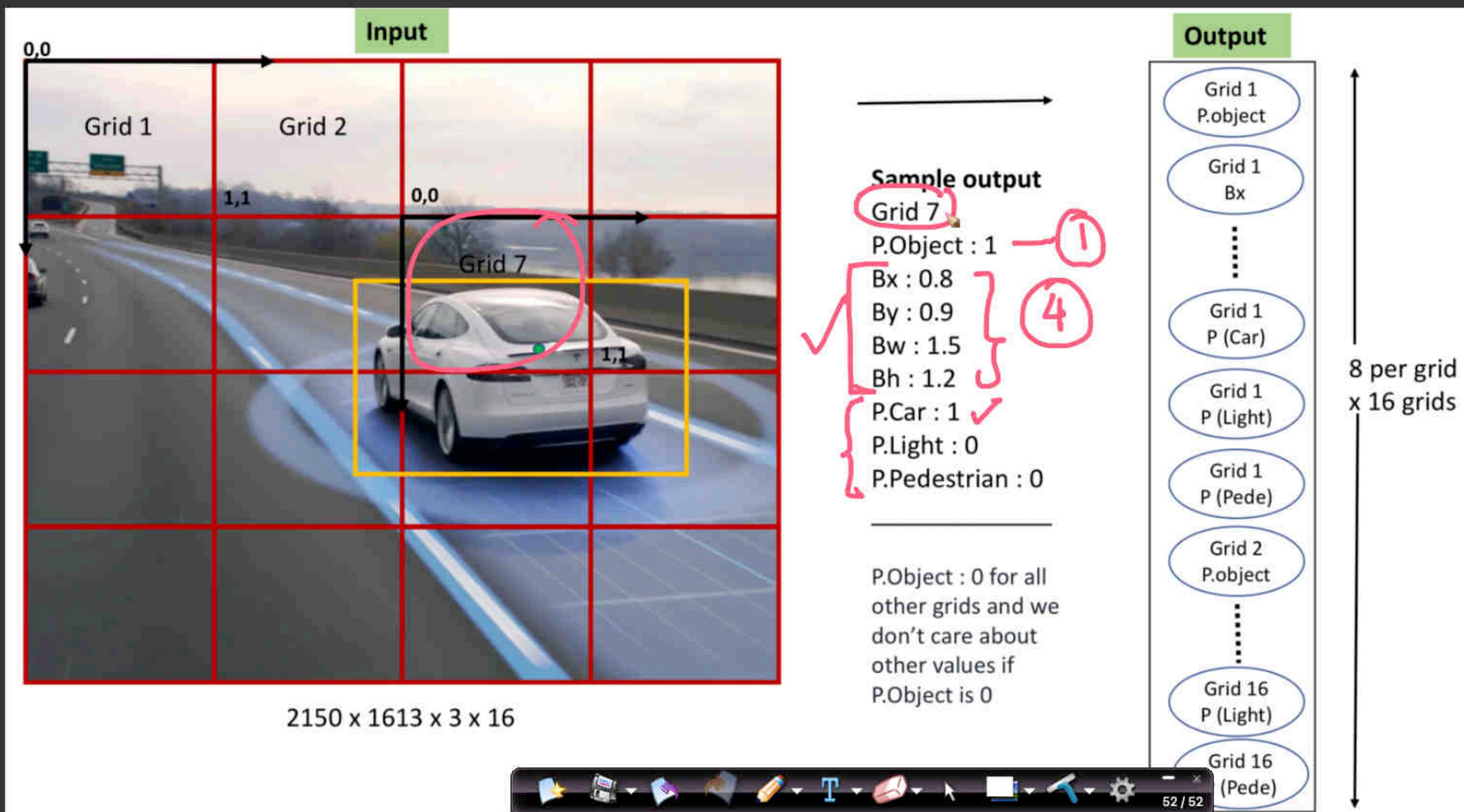


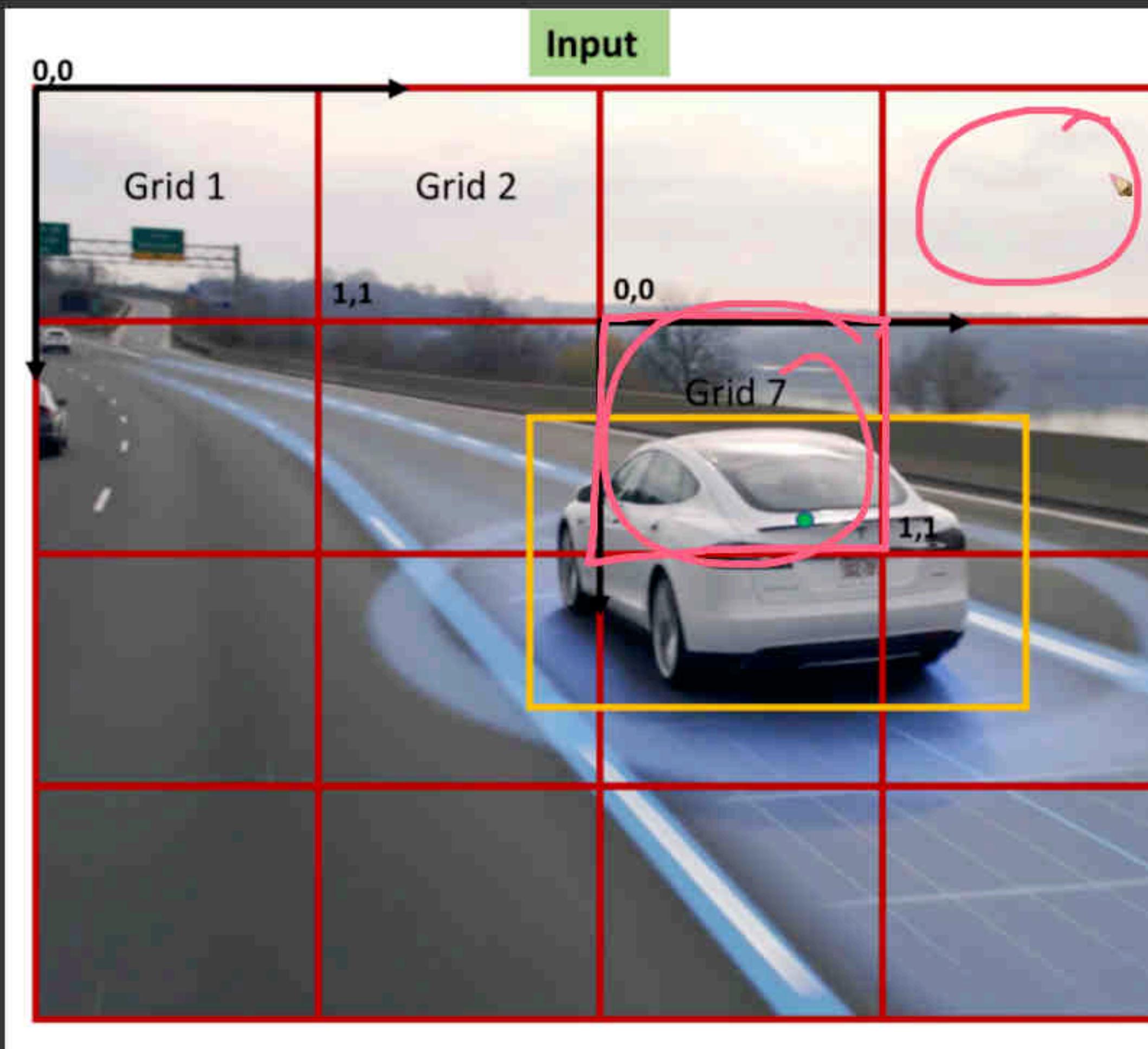




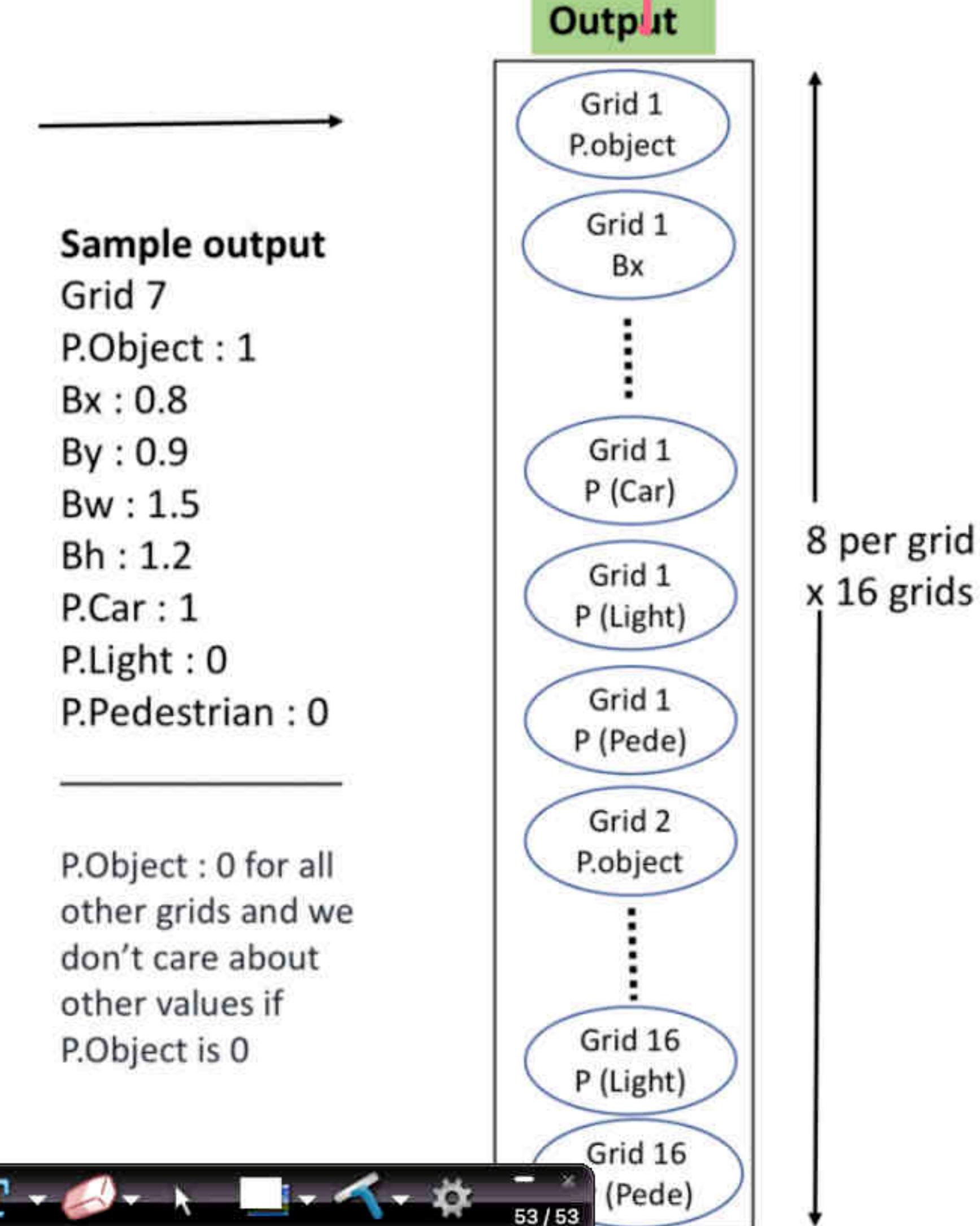


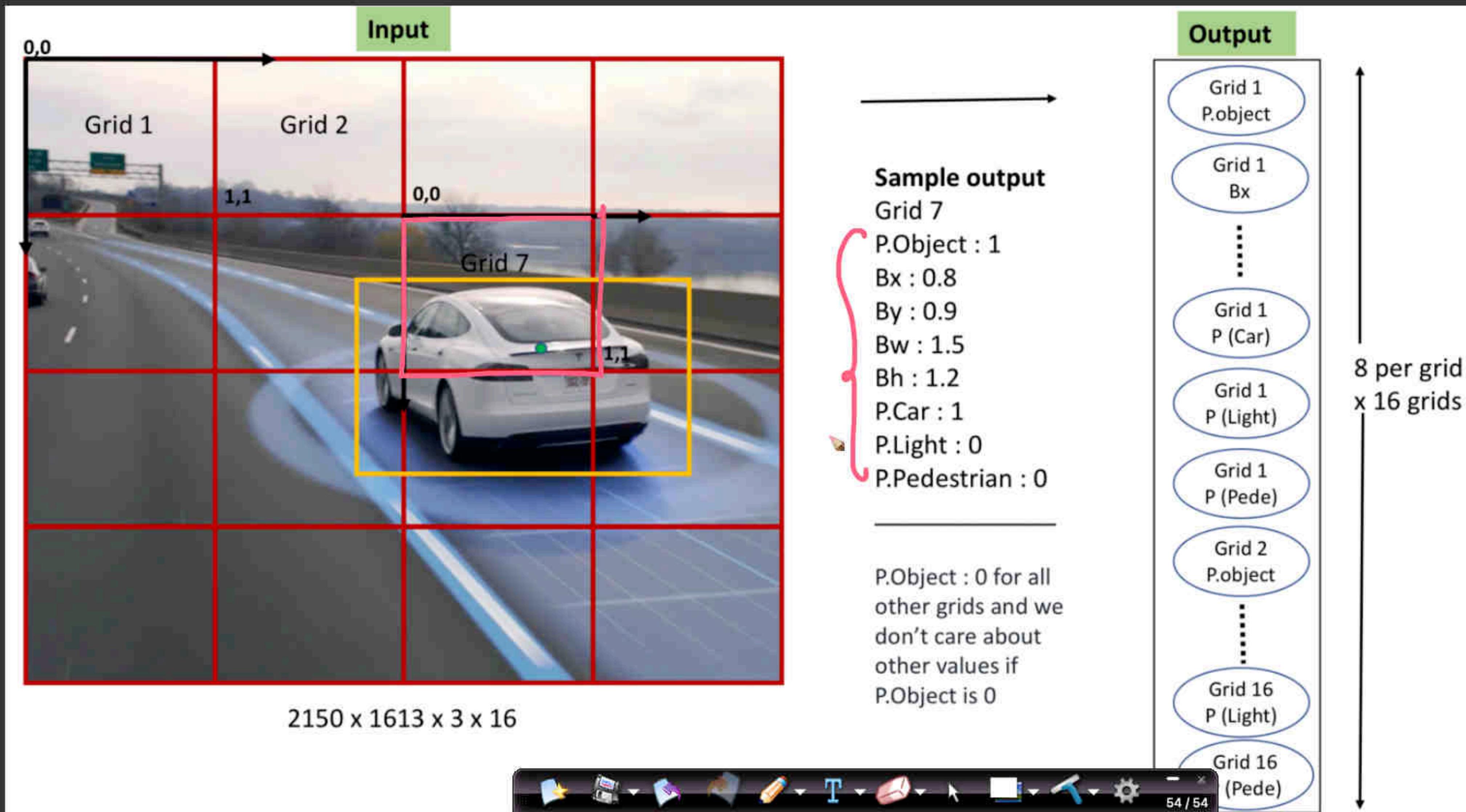


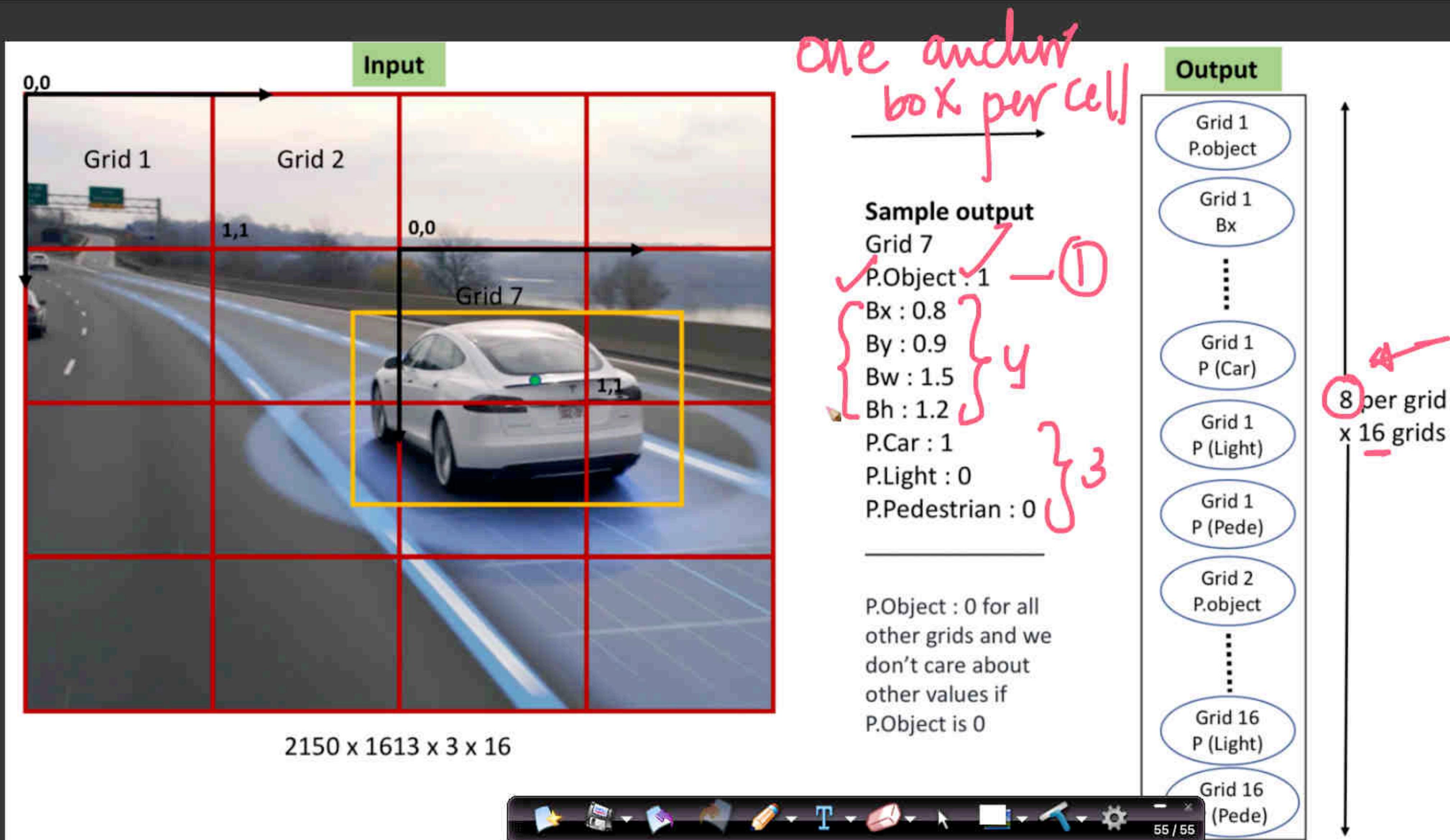


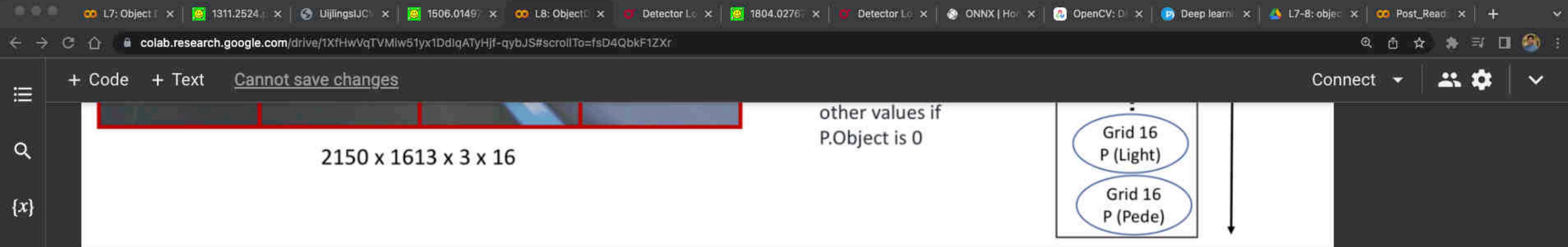


2 anchor boxes per cell









Yolo V3 Architecture details:

As mentioned in the [original paper](#), YOLOv3 has 53 convolutional layers called Darknet-53 is shown in the bellow image. The figure is mainly composed of Convolutional and Residual structures. It should be noted that the last three layers Avgpool, Connected, and softmax layer, are used for classification training on the Imagenet dataset. When using the Darknet-53 layer to extract features from the picture, these three layers are not used:

| Type | Filters | Size | Output |
|------------------|---------|------------------|------------------|
| Convolutional | 32 | 3×3 | 256×256 |
| Convolutional | 64 | $3 \times 3 / 2$ | 128×128 |
| Convolutional | 32 | 1×1 | |
| 1x Convolutional | 64 | 3×3 | |
| Residual | | | 128×128 |
| Convolutional | 128 | $3 \times 3 / 2$ | 64×64 |
| Convolutional | 64 | 1×1 | |
| 2x Convolutional | 128 | 3×3 | |
| Residual | | | 64×64 |
| Convolutional | 256 | $3 \times 3 / 2$ | $32 \times$ |
| Convolutional | 128 | 1×1 | |

Back-bone
Image → Conv. layers

fastest calculation speed. Compared with ResNet-101, the speed of the Darknet-53 network is 1.5 times that of the former; although ResNet-152 and its performance are similar, it takes more than two times.

2. it has highest measurement floating-point operation per second, which means that the network structure can better use the GPU, thereby making it more efficient and faster.

| Backbone | Top-1 | Top-5 | Bn Ops | BFLOP/s | FPS |
|------------------|-------|-------|--------|---------|-----|
| ✓Darknet-19 [15] | 74.1 | 91.8 | 7.29 | 1246 | 171 |
| ResNet-101[5] | 77.1 | 93.7 | 19.7 | 1039 | 53 |
| ResNet-152 [5] | 77.6 | 93.8 | 29.4 | 1090 | 37 |
| ✓Darknet-53 | 77.2 | 93.8 | 18.7 | 1457 | 78 |

YOLO in easy steps:

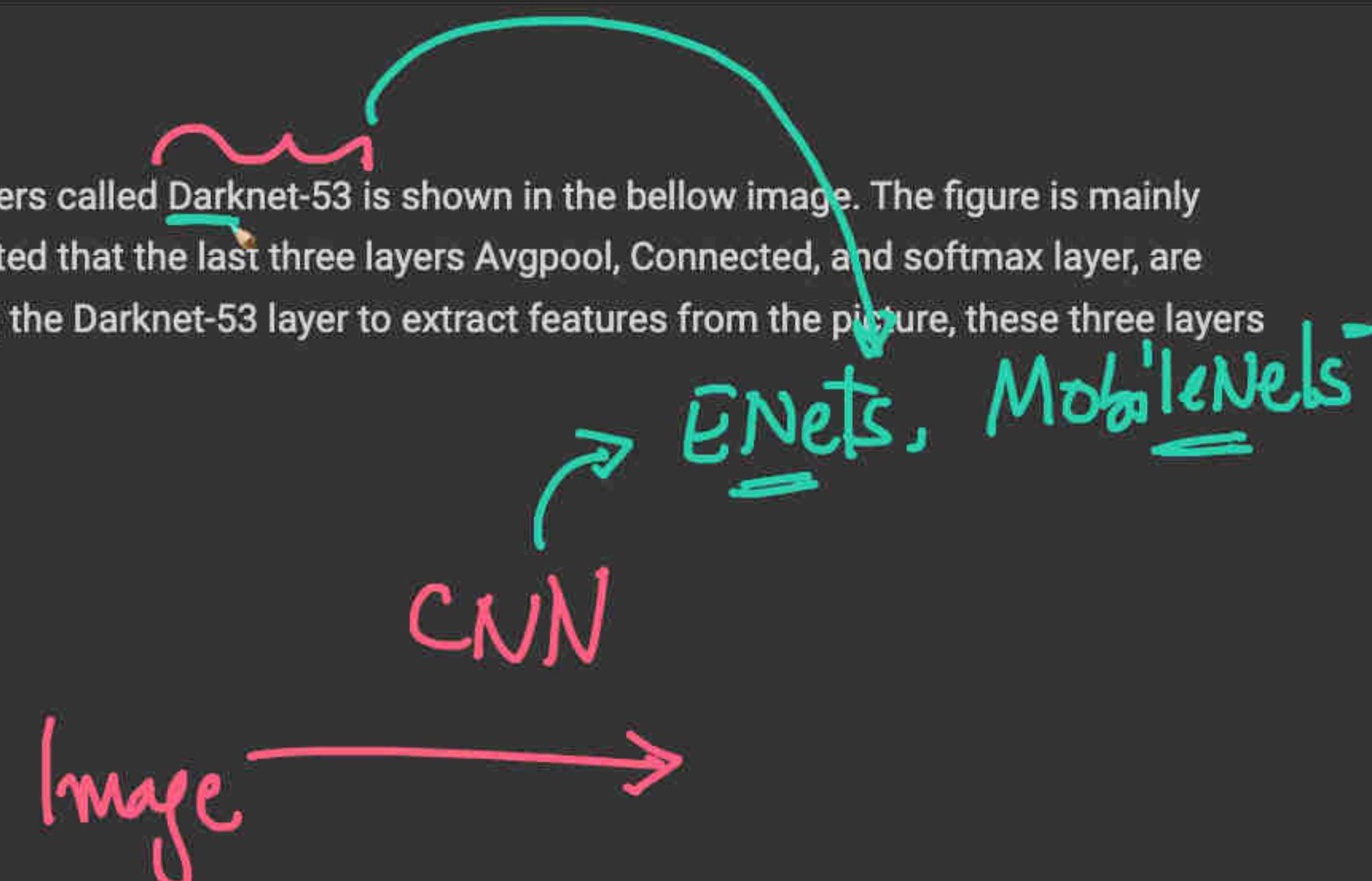
1. Divide the image into multiple grids. For illustration, I have drawn 4×4 grids in above figure.
 2. Label the training data as shown in the above figure.

- If C is number of unique objects in our data, $S \times S$ is number of grids into which we split our image, then our output vector will be of length $S \times S \times (C+5)$.
 - For e.g. in above case, our target vector will be of length $10 \times 10 \times (3+5) = 800$.

Yolo V3 Architecture details:

As mentioned in the [original paper](#), YOLOv3 has 53 convolutional layers called Darknet-53 is shown in the bellow image. The figure is mainly composed of Convolutional and Residual structures. It should be noted that the last three layers Avgpool, Connected, and softmax layer, are used for classification training on the Imagenet dataset. When using the Darknet-53 layer to extract features from the picture, these three layers are not used:

| Type | Filters | Size | Output |
|---------------|---------|------------------|------------------|
| Convolutional | 32 | 3×3 | 256×256 |
| Convolutional | 64 | $3 \times 3 / 2$ | 128×128 |
| 1x | 32 | 1×1 | |
| Convolutional | 64 | 3×3 | |
| Residual | | | 128×128 |
| Convolutional | 128 | $3 \times 3 / 2$ | 64×64 |
| Convolutional | 64 | 1×1 | |
| 2x | 128 | 3×3 | |
| Residual | | | 64×64 |
| Convolutional | 256 | $3 \times 3 / 2$ | 32×32 |
| Convolutional | 128 | 1×1 | |
| 8x | 256 | 3×3 | |
| Residual | | | 32×32 |
| Convolutional | 512 | $3 \times 3 / 2$ | 16×16 |
| Convolutional | 256 | 1×1 | |
| 8x | 512 | 3×3 | |
| Residual | | | 16×16 |



Darknet-53

77.2

93.8

18.7

1457

78

YOLO in easy steps:

1. Divide the image into multiple grids. For illustration, I have drawn 4×4 grids in above figure.

2. Label the training data as shown in the above figure.

- If C is number of unique objects in our data, $S \times S$ is number of grids into which we split our image, then our output vector will be of length $S \times S \times (C+5)$.
- For e.g. in above case, our target vector is $4 \times 4 \times (3+5)$ as we divided our images into 4×4 grids and are training for 3 unique objects: Car, Light and Pedestrian.

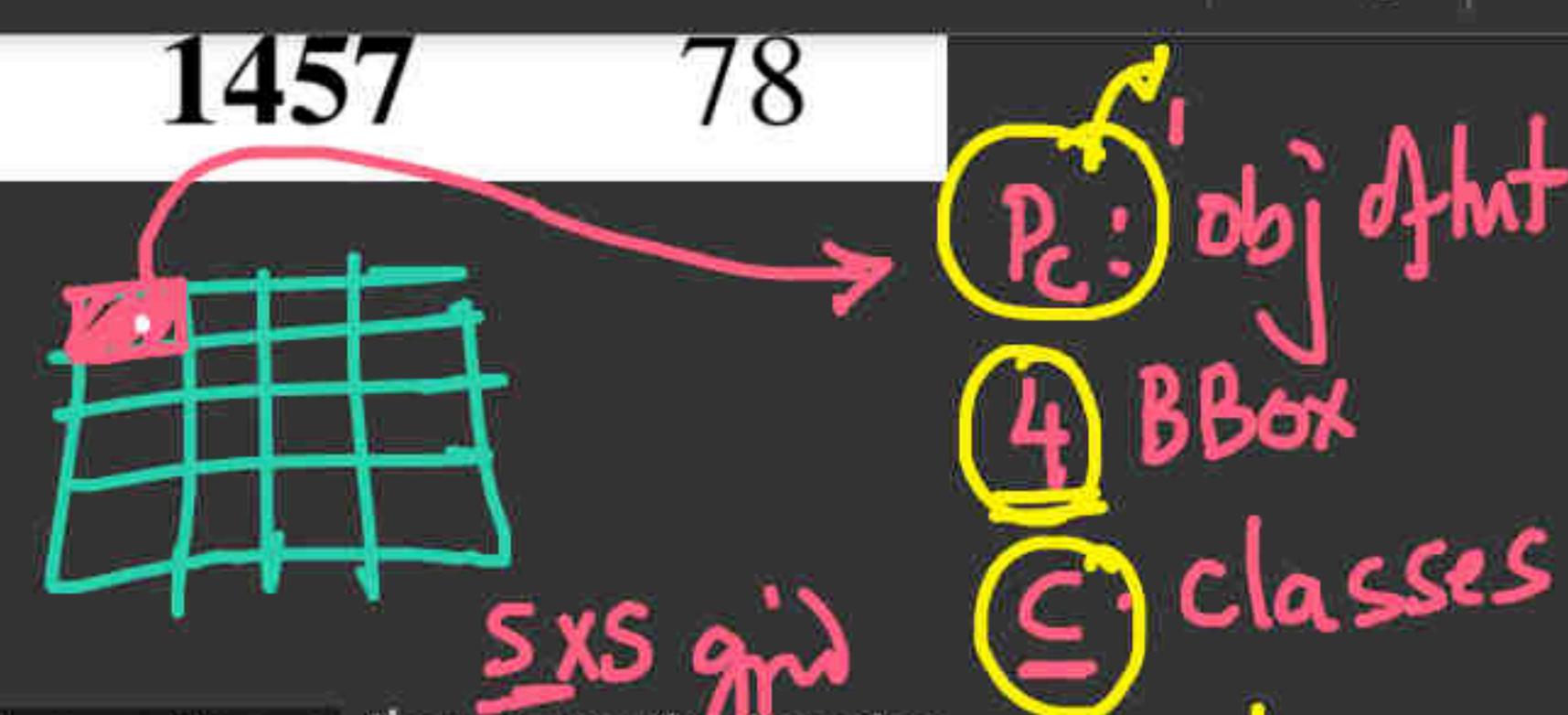
3. Make one deep convolutional neural net with loss function as error between output activations and label vector.

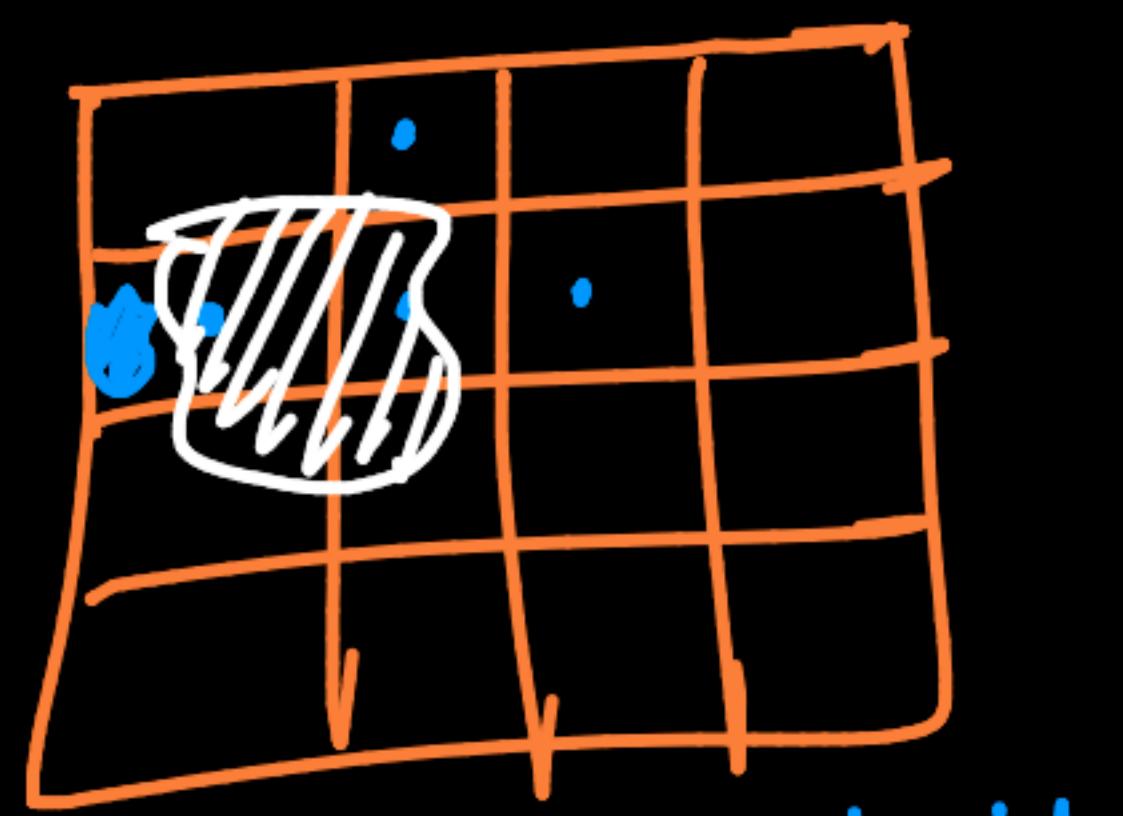
- Basically, the model predicts the output of all the grids in just one forward pass of input image through ConvNet.

4. Keep in mind that the label for object being present in a grid cell (P.Object) is determined by the presence of object's centroid in that grid.

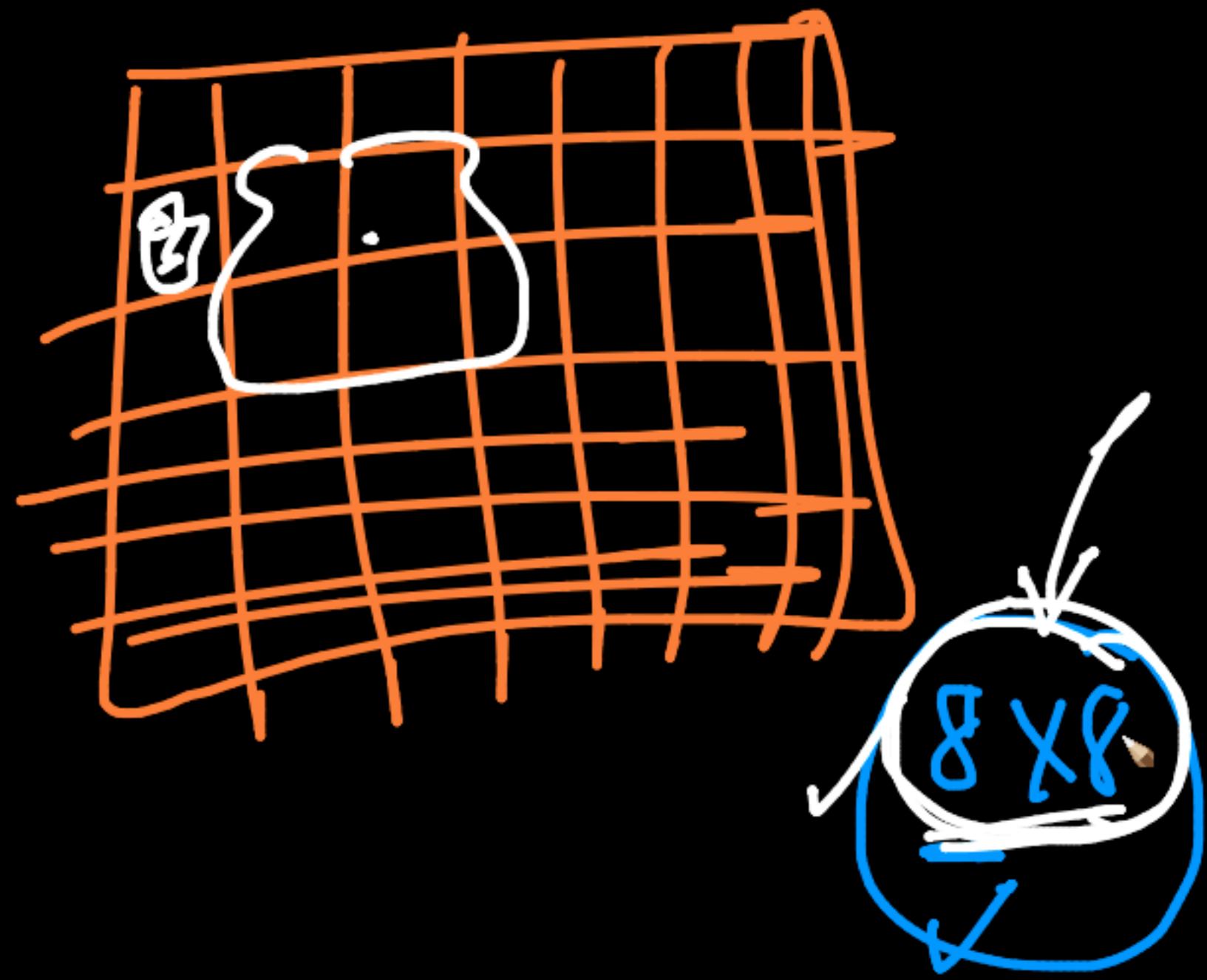
- This is important to not allow one object to be counted multiple times in different grids.

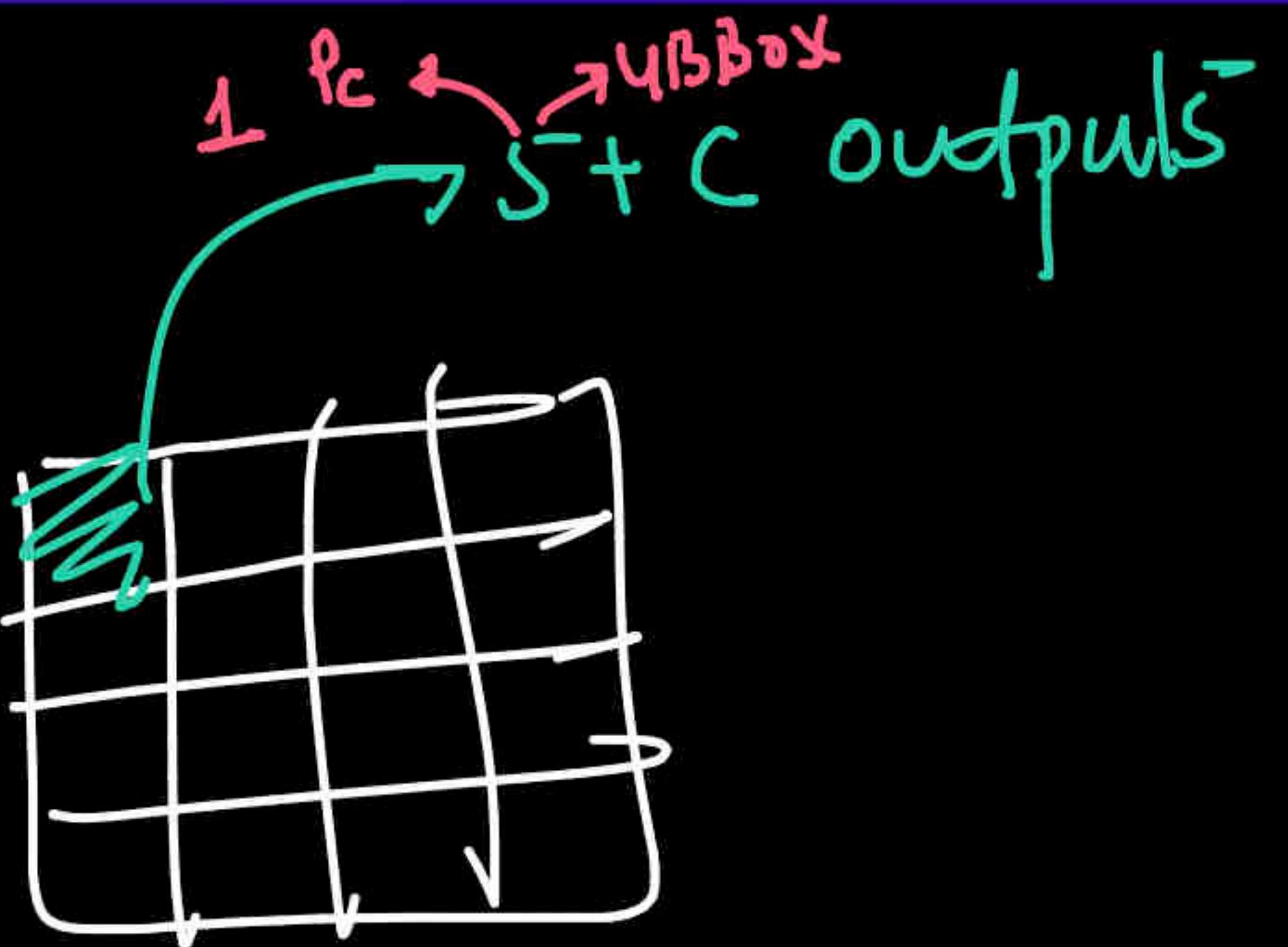
The main advantage Of YOLO is, it is very fast as the model predicts the output of all the grids in just one forward pass of input image through ConvNet.





Case 1 $\underline{u \times u}$





$S \times S$

$S \times S \times (5 + C)$

YOLO in easy steps:

1. Divide the image into multiple grids. For illustration, I have drawn 4x4 grids in above figure.
2. Label the training data as shown in the above figure.
 - If **C** is number of unique objects in our data, $s \times s$ is number of grids into which we split our image, then our output vector will be of length $s \times s \times (C + 5)$.
 - For e.g. in above case, our target vector is $4 \times 4 \times (3 + 5)$ as we divided our images into 4×4 grids and are training for 3 unique objects: Car, Light and Pedestrian.
3. Make one deep convolutional neural net with loss function as error between output activations and label vector.
 - Basically, the model predicts the output of all the grids in just one forward pass of input image through ConvNet.
4. Keep in mind that the label for object being present in a grid cell (**P.Object**) is determined by the presence of object's centroid in that grid.
 - This is important to not allow one object to be counted multiple times in different grids.



The main advantage Of YOLO is, it is very fast as the model predicts the output of all the grids in just one forward pass of input image through ConvNet.

YOLO Loss:





+ Code + Text Cannot save changes

YOLO in easy steps:

1. Divide the image into multiple grids. For illustration, I have drawn 4x4 grids in above figure.

2. Label the training data as shown in the above figure.

- If C is number of unique objects in our data, $S \times S$ is number of grids into which we split our image, then our output vector will be of length $S \times S \times (C+5)$.
- For e.g. in above case, our target vector is $4 \times 4 \times (3+5)$ as we divided our images into 4×4 grids and are training for 3 unique objects: Car, Light and Pedestrian.

3. Make one deep convolutional neural net with loss function as error between output activations and label vector.

- Basically, the model predicts the output of all the grids in just one forward pass of input image through ConvNet.

4. Keep in mind that the label for object being present in a grid cell (P_{Object}) is determined by the presence of object's centroid in that grid.

- This is important to not allow one object to be counted multiple times in different grids.

The main advantage Of YOLO is, it is very fast as the model predicts the output of all the grids in just one forward pass of input image through ConvNet.

YOLO Loss:



+ Code + Text Cannot save changes

YOLO in easy steps:

1. Divide the image into multiple grids. For illustration, I have drawn 4x4 grids in above figure.
2. Label the training data as shown in the above figure.
 - If C is number of unique objects in our data, $S \times S$ is number of grids into which we split our image, then our output vector will be of length $S \times S \times (C + 5)$.
 - For e.g. in above case, our target vector is $4 \times 4 \times (3 + 5)$ as we divided our images into 4×4 grids and are training for 3 unique objects: Car, Light and Pedestrian.
3. Make one deep convolutional neural net with loss function as error between output activations and label vector.
 - Basically, the model predicts the output of all the grids in just one forward pass of input image through ConvNet.
4. Keep in mind that the label for object being present in a grid cell (P.Object) is determined by the presence of object's centroid in that grid.
 - This is important to not allow one object to be counted multiple times in different grids.

The main advantage Of YOLO is, it is very fast as the model predicts the output of all the grids in just one forward pass of input image through ConvNet.

YOLO Loss:



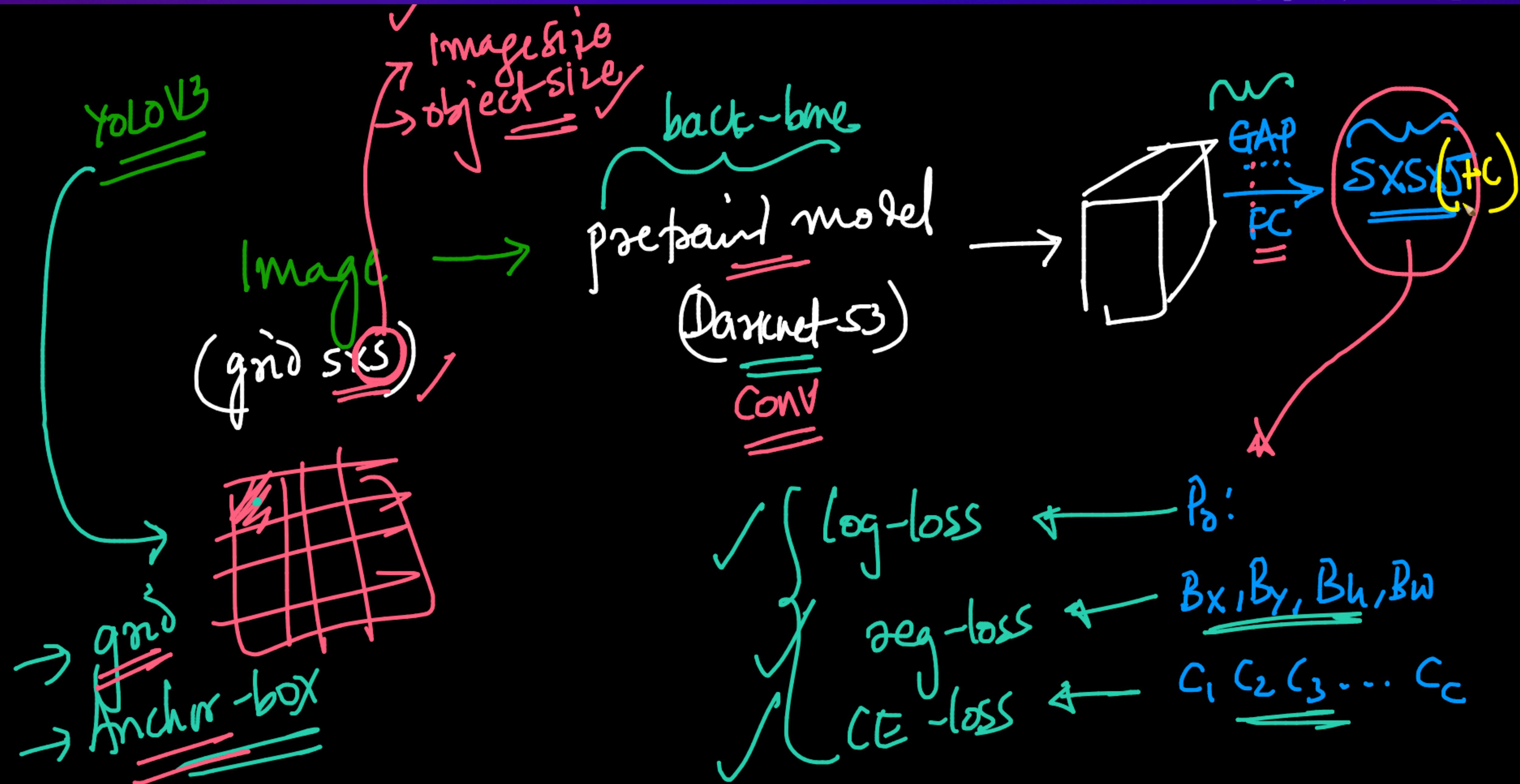
YOLO in easy steps:

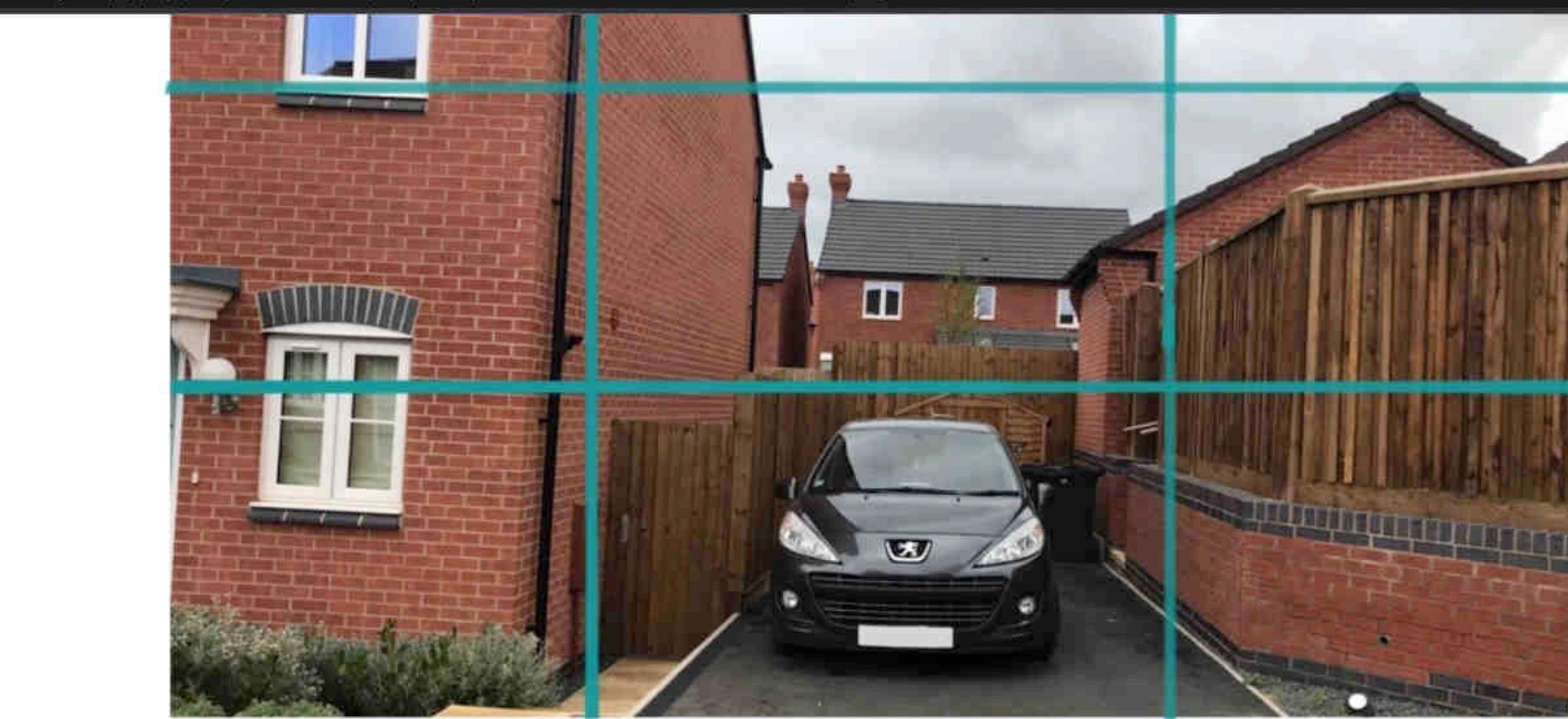
1. Divide the image into multiple grids. For illustration, I have drawn 4x4 grids in above figure.
2. Label the training data as shown in the above figure.
 - If **C** is number of unique objects in our data, $s \times s$ is number of grids into which we split our image, then our output vector will be of length $s \times s \times (C+5)$.
 - For e.g. in above case, our target vector is $4 \times 4 \times (3+5)$ as we divided our images into 4×4 grids and are training for 3 unique objects: Car, Light and Pedestrian
3. Make one deep convolutional neural net with loss function as error between output activations and label vector.
 - Basically, the model predicts the output of all the grids in just one forward pass of input image through ConvNet.
4. Keep in mind that the label for object being present in a grid cell (**P.Object**) is determined by the presence of object's centroid in that grid.
 - This is important to not allow one object to be counted multiple times in different grids.

The main advantage Of YOLO is, it is very fast as the model predicts the output of all the grids in just one forward pass of input image through ConvNet.

YOLO Loss:






 $\lambda_{noobj} = 0, 5$

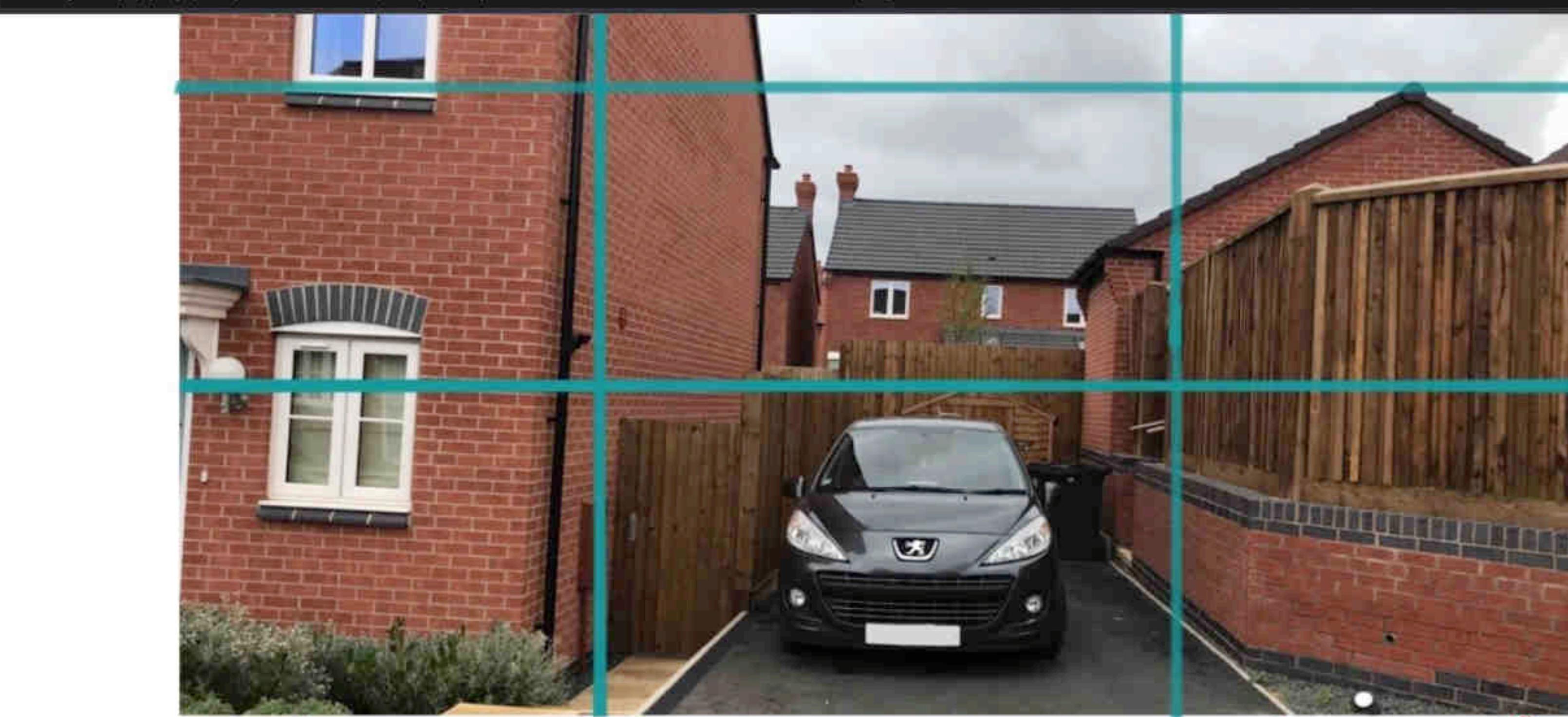
j

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

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

$$\sum_{i=0}^{S^2} 1_{ij}^{noobj} \sum_{C \in classes} \left[(p_i(C) - \hat{p}_i(C))^2 \right] \longrightarrow \text{Classification}$$

$1_{ij}^{noobj} \rightarrow 1$ if box j and cell i no match.
0 otherwise


 $\lambda_{noobj} = 0, 5$

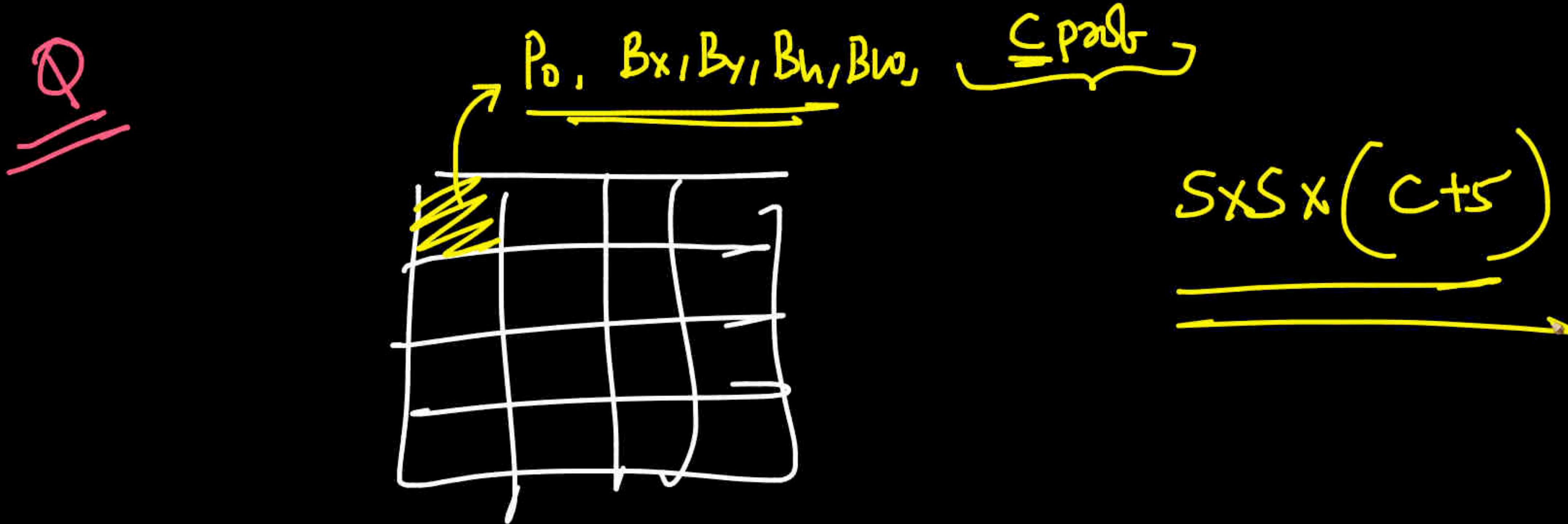
j

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

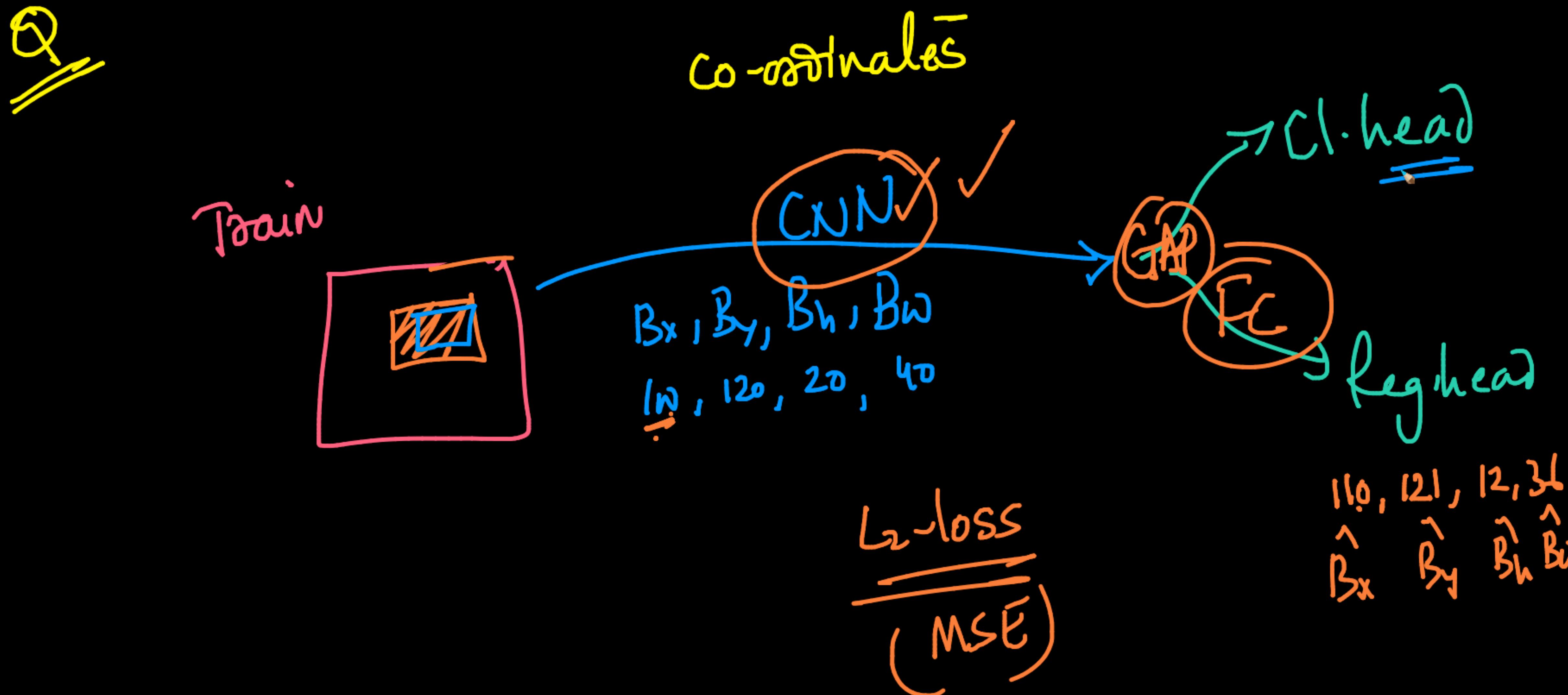
$$\sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{obj} \left[(C_i - \hat{C}_i)^2 \right] + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B 1_{ij}^{noobj} \left[(C_i - \hat{C}_i)^2 \right] \longrightarrow \text{Confidence}$$

$$\sum_{i=0}^{S^2} 1_{ij}^{noobj} \sum_{C \in \text{classes}} \left[(p_i(C) - \hat{p}_i(C))^2 \right] \longrightarrow \text{Classification}$$

 P_{obj}
 $1_{ij}^{noobj} \rightarrow \begin{cases} 1 & \text{if box } j \text{ and cell } i \text{ no match.} \\ 0 & \text{otherwise} \end{cases}$



one anchor box per cell



+ Code + Text Cannot save changes Connect |  

Problem Statement:

You are working as a Machine learning Engineer in Tesla as part of the Autonomous driving team focused on building prototype ML model to identify and detect location of:

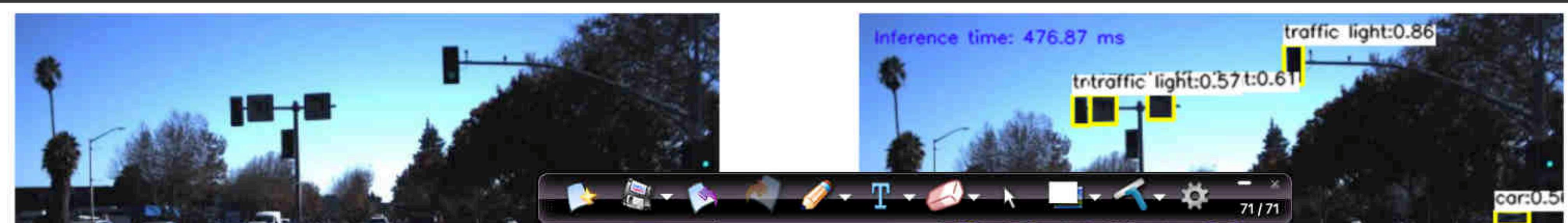
- Traffic Lights, Stop signs
- Other Vehicles such as Car, Bicycles, Truck etc
- Pedestrian(People)
- Animals

on the road in realtime using the Camera installed in the deck of the owners car.

Real time constraints:

- Our Object Detection Algorithm needs to be Fast enough to **run inference in RealTime** on commodity hardware(CPU only) and **highly accurate to ensure safety**.

Let's see how we can achieve this Balance using new family of Object Detection Algorithms: **Single Stage Detectors**.



Inference time: 476.87 ms
traffic light:0.86
traffic light:0.57 t:0.61
car:0.51

71/71