

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



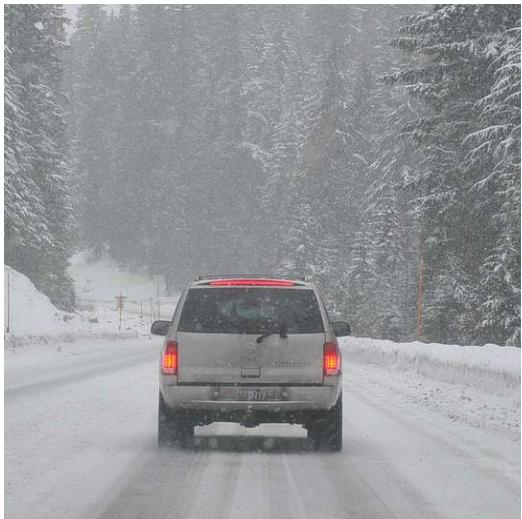
deeplearning.ai

Object Detection

Object localization

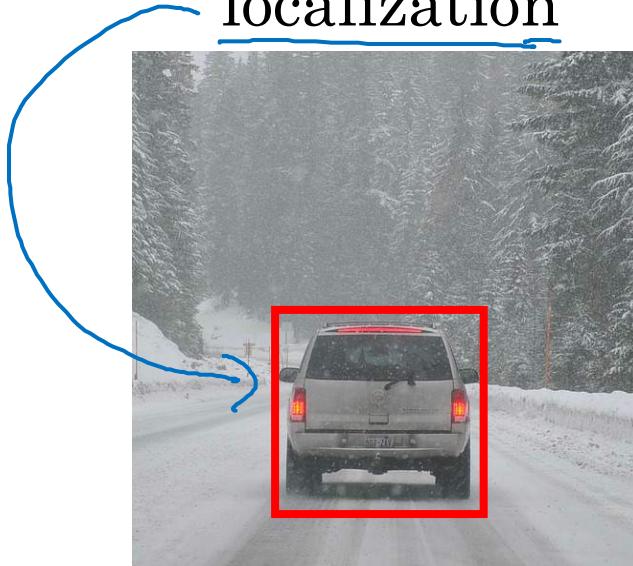
What are localization and detection?

Image classification



"Car"

Classification with
localization



"Car"

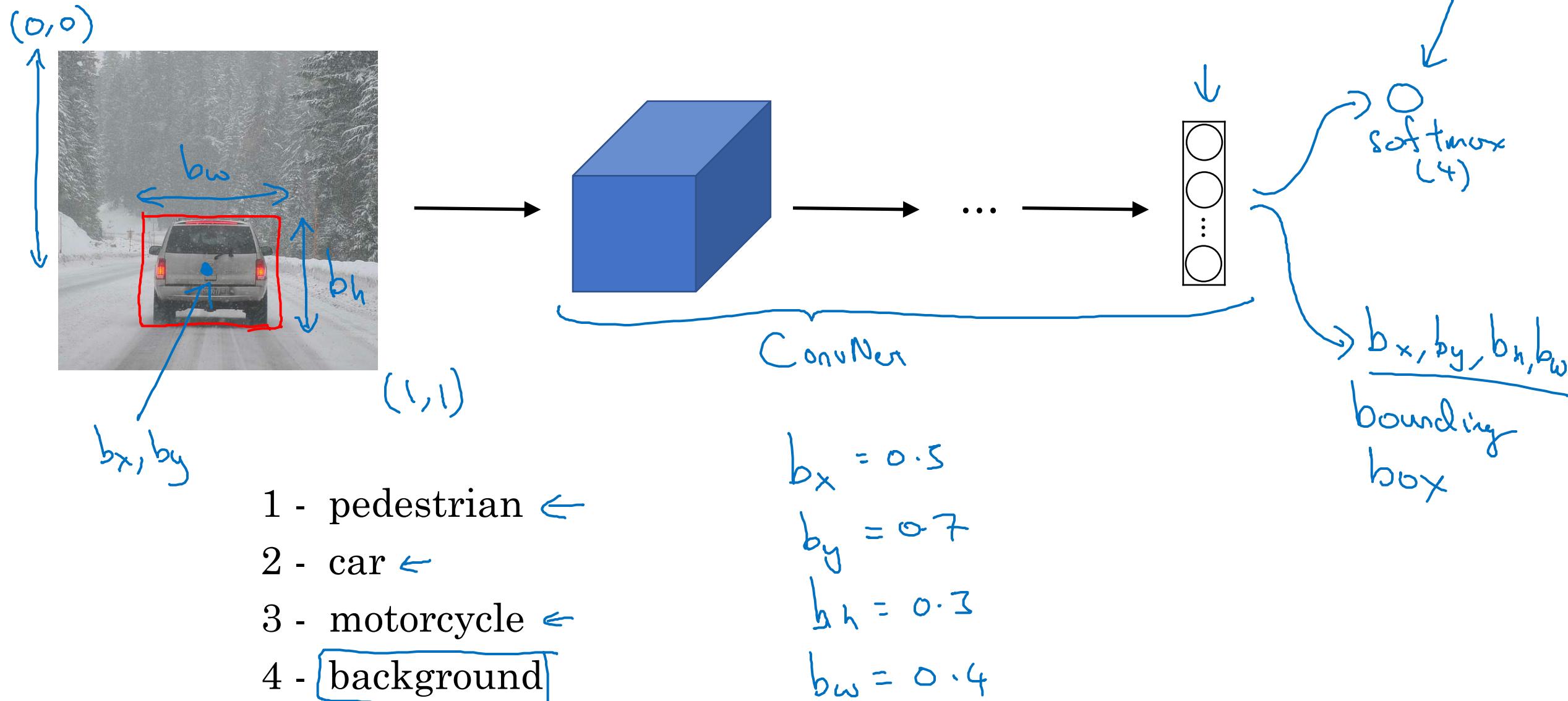
1 object

Detection



multiple
objects

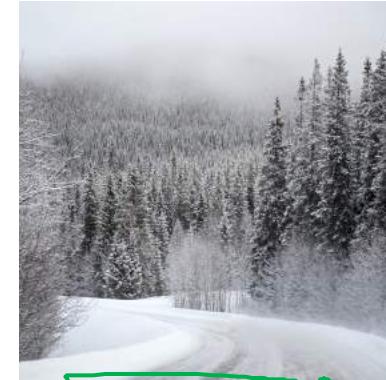
Classification with localization



Defining the target label y

- 1 - pedestrian
- 2 - car 
- 3 - motorcycle
- 4 - background 

Need to output b_x, b_y, b_h, b_w , class label (1-4)

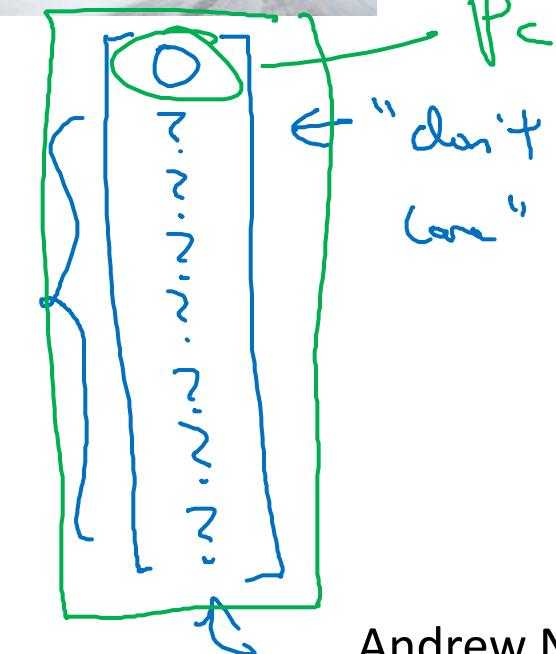
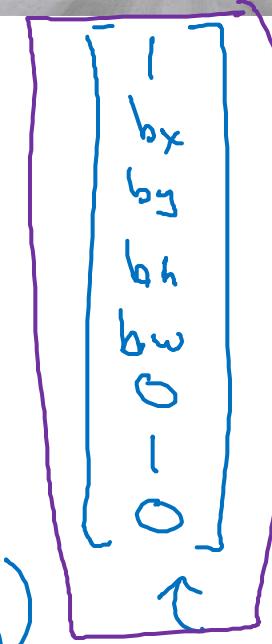


$x =$

$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 \\ + \dots + (\hat{y}_8 - y_8)^2 & \text{if } \underline{y_1 = 1} \\ (\hat{y}_1 - y_1)^2 & \text{if } \underline{y_1 = 0} \end{cases}$$

$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix} \quad \rightarrow \quad \left\{ \begin{array}{l} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{array} \right\} \text{ is there any object?}$$

(x, y)



Andrew Ng

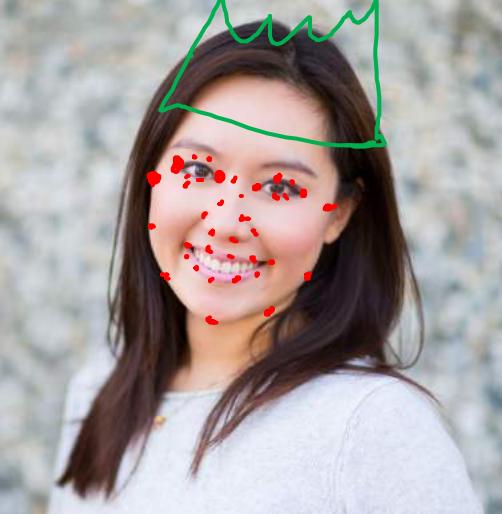
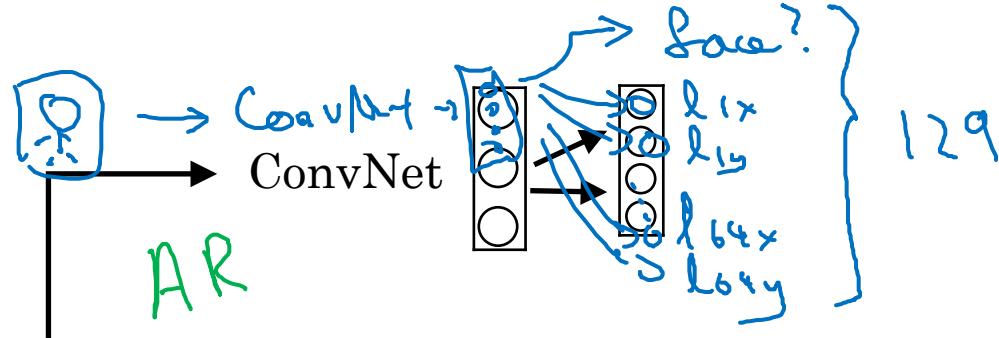


deeplearning.ai

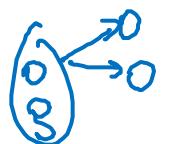
Object Detection

Landmark detection

Landmark detection



b_x, b_y, b_h, b_w



$l_{1x}, l_{1y},$
 $l_{2x}, l_{2y},$
 $l_{3x}, l_{3y},$
 $l_{4x}, l_{4y},$
:
 l_{64x}, l_{64y}

x, y

$l_{1x}, l_{1y},$
:
 l_{32x}, l_{32y}



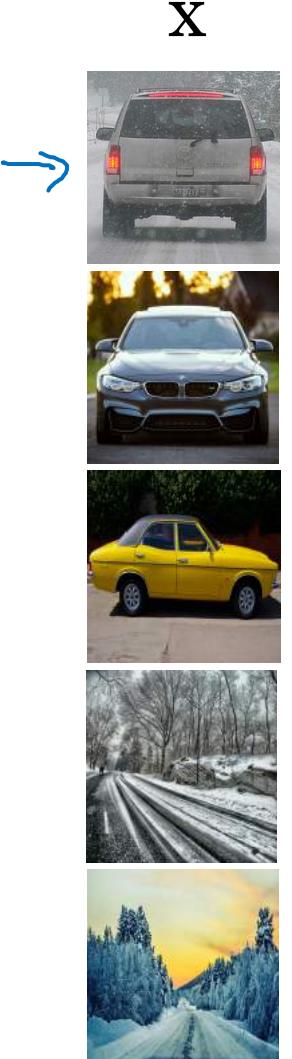
deeplearning.ai

Object Detection

Object detection

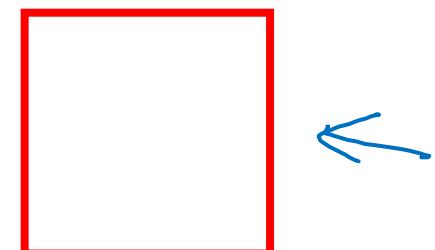
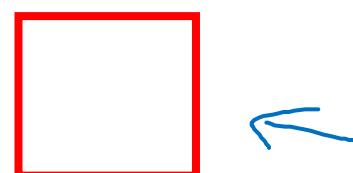
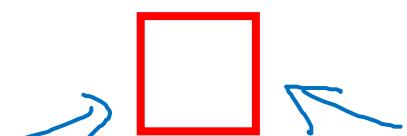
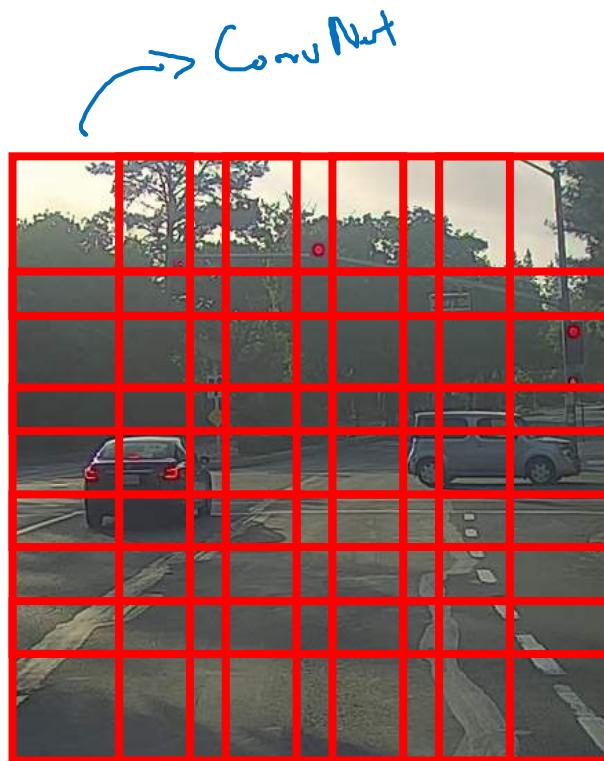
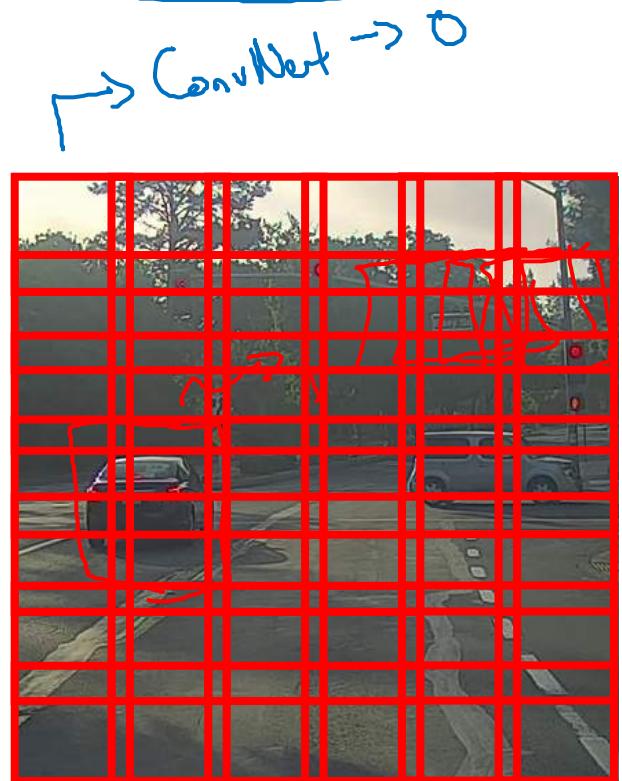
Car detection example

Training set:



$\xrightarrow{\hspace{1cm}}$ ConvNet $\rightarrow y$

Sliding windows detection



Computation cost

Idea: instead of 1-training a CNN on images of cars precisely cropped and 2-
at prediction time operate multiple predictions on many crops (sliding window) of
a test image , opt for this alternative.



deeplearning.ai

↳ Design the CNN with filters and strides of a specific size so that the final layers are not standard FC , but still 2D layers corresponding to various locations on the image.

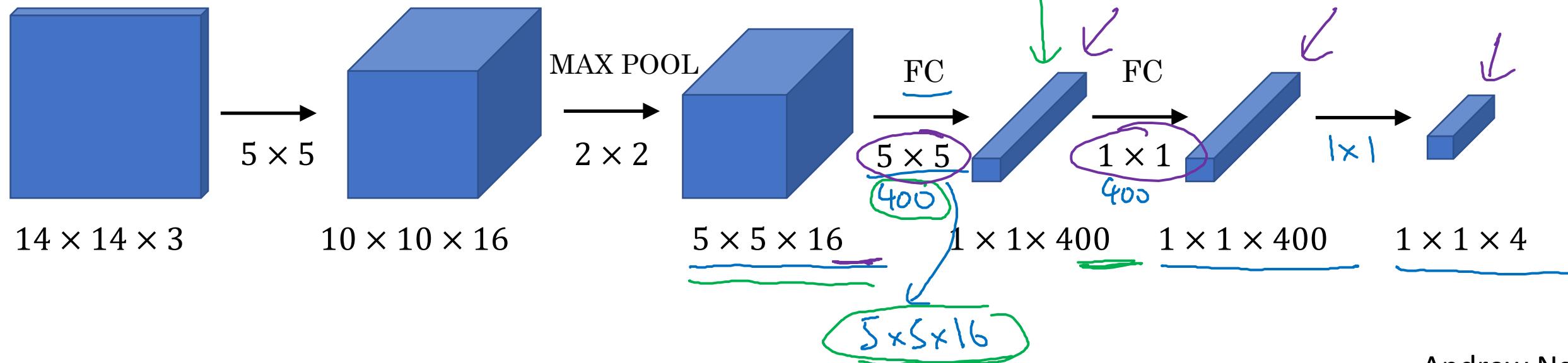
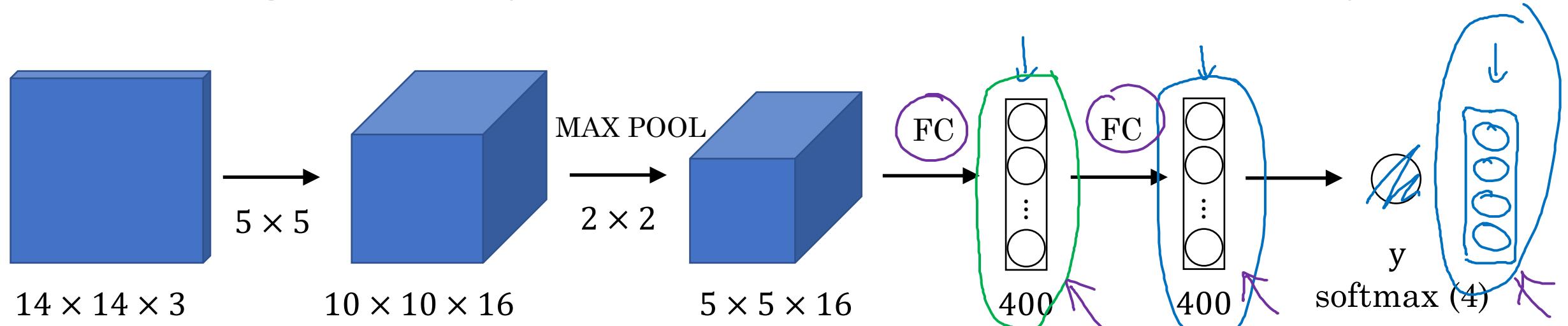
Each pixel on this final layer will correspond to a specific FC neuron due to the properties of the convolution .

↳ idea: conv 5,5 on a 5x5 layer corresponds to an FC operation .

Object Detection

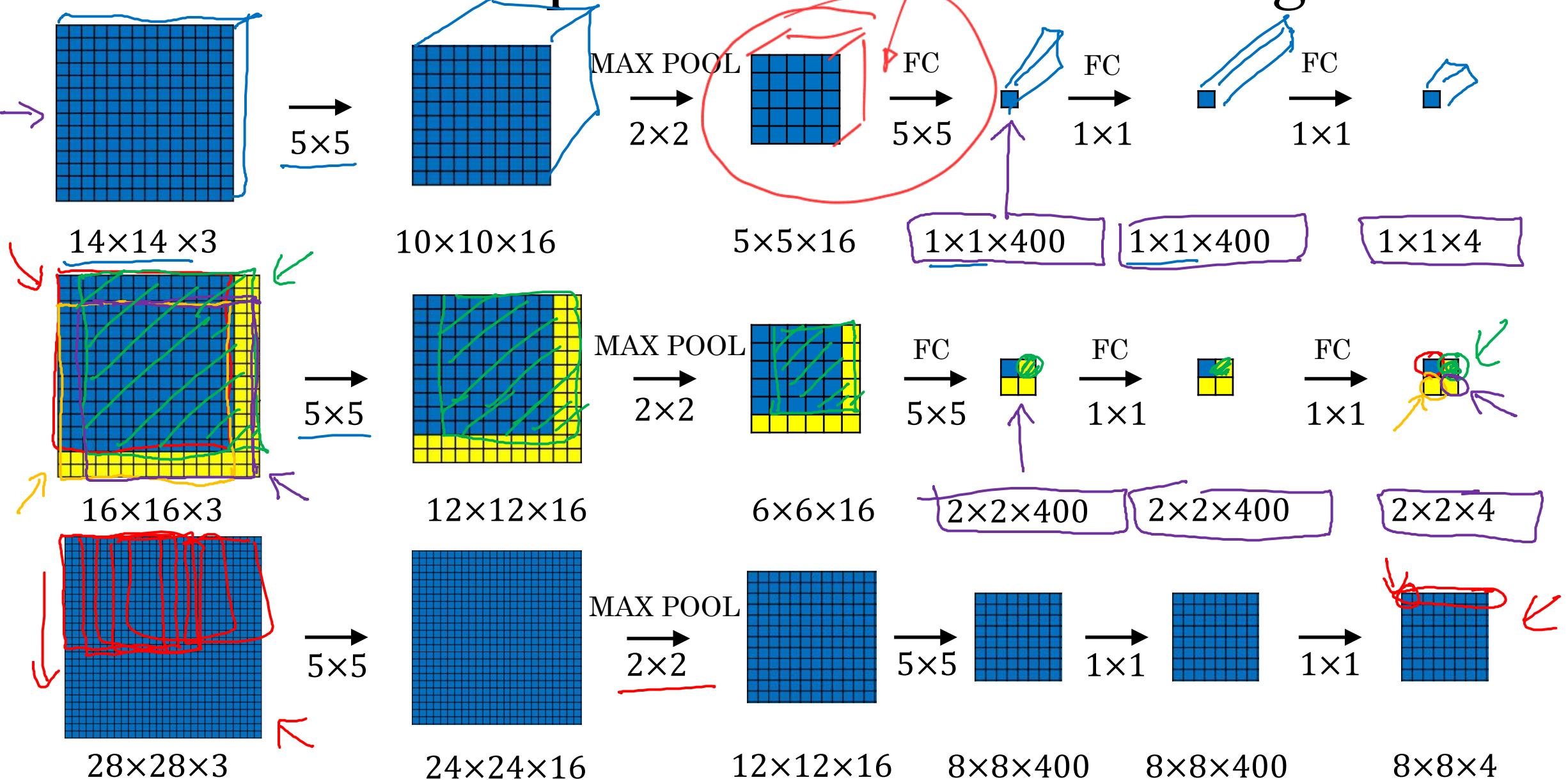
Convolutional implementation of sliding windows

Turning FC layer into convolutional layers

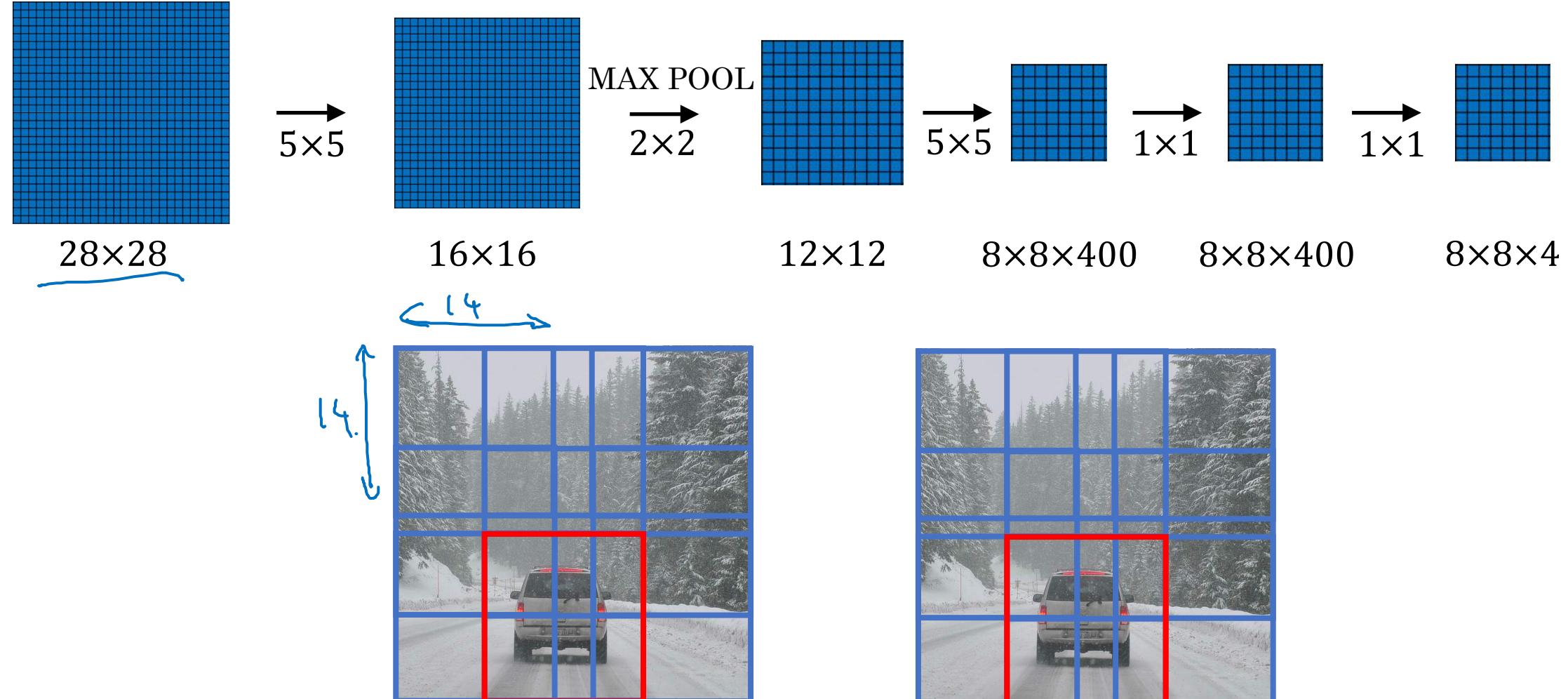


Convolution implementation of sliding windows

Idea: conv 5,5 on a 5x5 layer corresponds to FC!



Convolution implementation of sliding windows



The YOLO architecture allows for a pretty efficient object detection solution. The idea is to integrate the simultaneous object recognition and localization task described before with the convolutional implementation of the sliding window approach. IN PRACTICE, this means designing an architecture



that outputs an

3x3x8 output

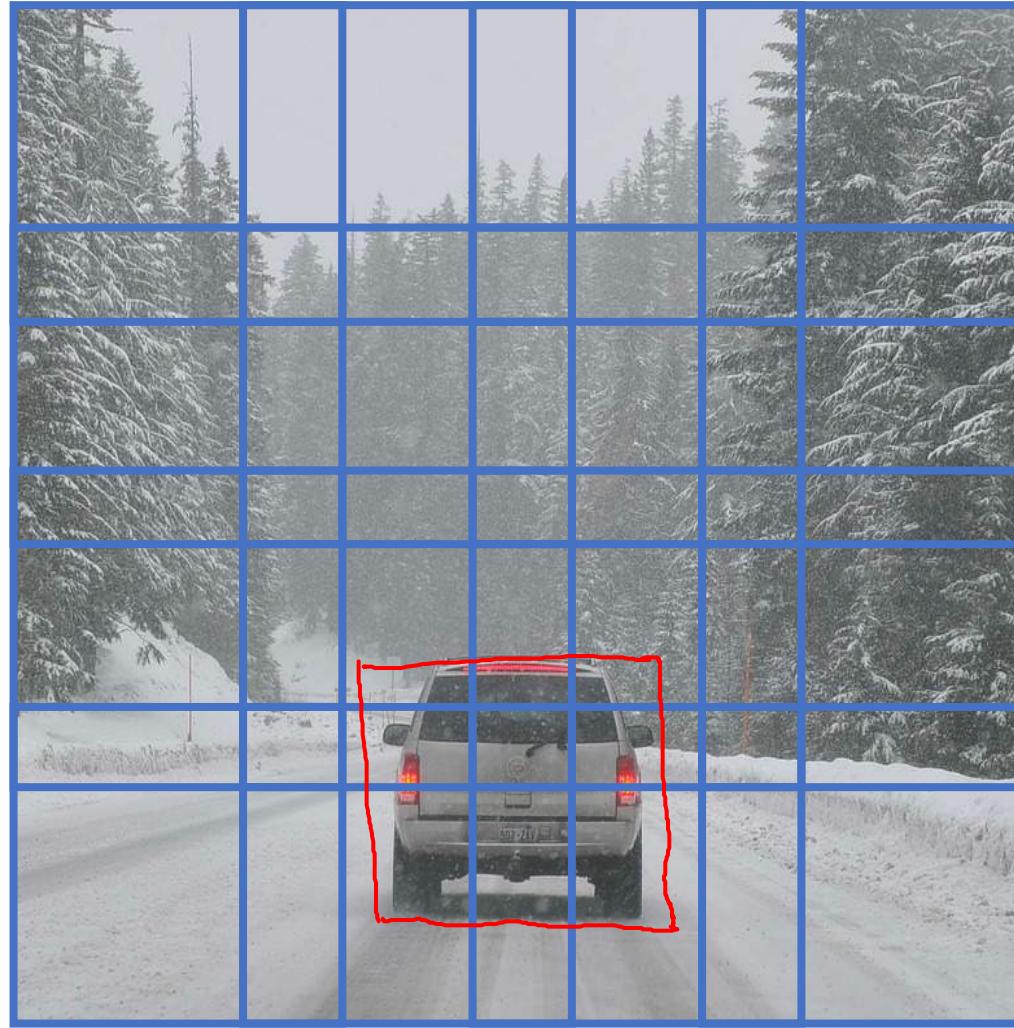
That implicitly carries out simultaneous object recognition and localization in
3x3 distinct non overlapping squares of the image.

Object Detection

Bounding box predictions

The point of dividing the image in these parts is to avoid that one object appears twice in one section.

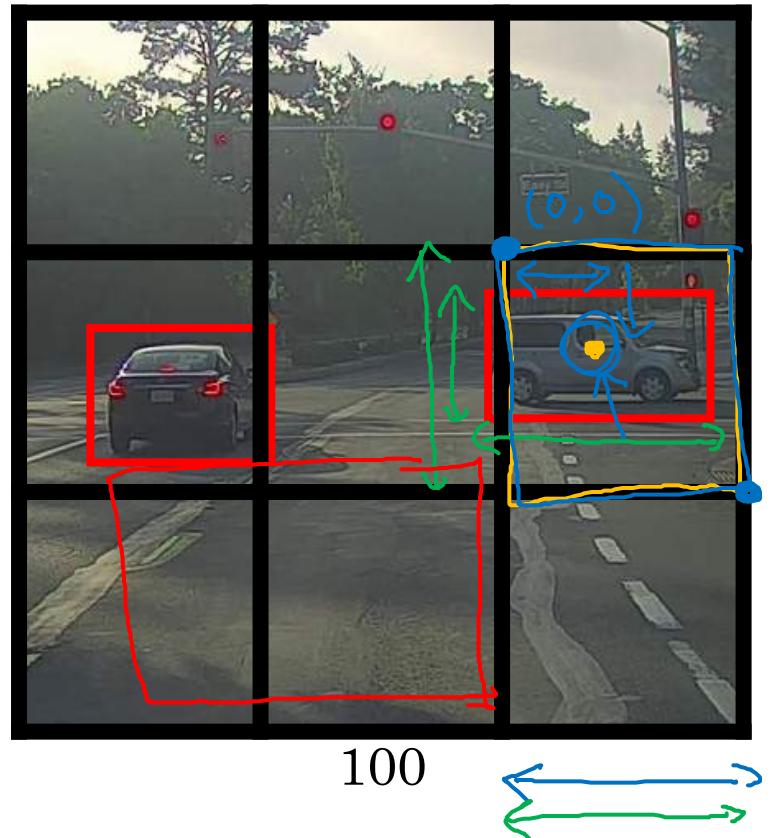
Output accurate bounding boxes



YOLO algorithm



Specify the bounding boxes



$$y = \begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ o \end{bmatrix}$$

0.4 } between 0 and 1
0.3 }
0.5 }
0.9 } Could be > 1

How to EVALUATE PERFORMANCE OF OBJECT DETECTION?

By comparing true and predicted bounding boxes using the intersection over union function.
If $\text{IoU} \geq 0.5$ is ok (heuristically...).

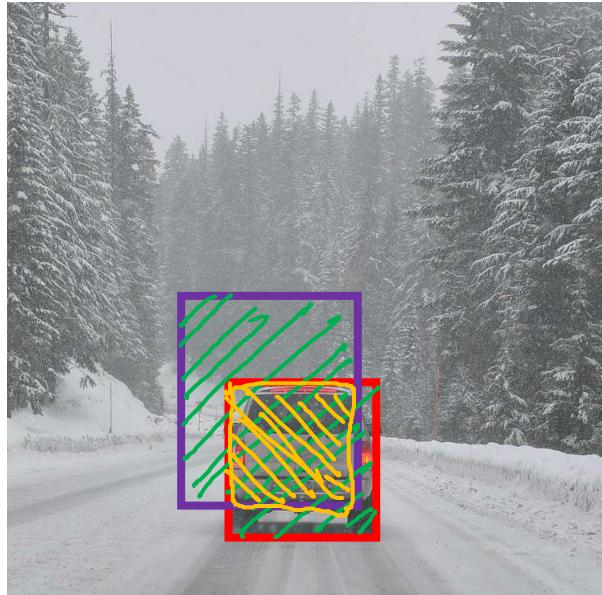


deeplearning.ai

Object Detection

Intersection
over union

Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{Size of intersection}}{\text{Size of union}}$$
A diagram illustrating the calculation of IoU. It shows two overlapping rectangles: one yellow and one green. The overlapping area is shaded with diagonal lines, representing the intersection. The total area covered by either rectangle or both is shaded with horizontal lines, representing the union.

“Correct” if $\text{IoU} \geq 0.5$

0.6

More generally, IoU is a measure of the overlap between two bounding boxes.

HOW TO OBTAIN ONLY ONE BB FOR EACH OBJECT?

By analyzing the image in cells, the object detection algorithm may come up with multiple BBs for each object (one BB created for each grid cell that covers the object).



deeplearning.ai

Object Detection

The solution is a "post-processing" elimination of overlapping BB using the Non-max suppression.

The idea is to take all the BB found in the image for a certain class and then iteratively

Non-max

take the one
with the most
confident
prediction

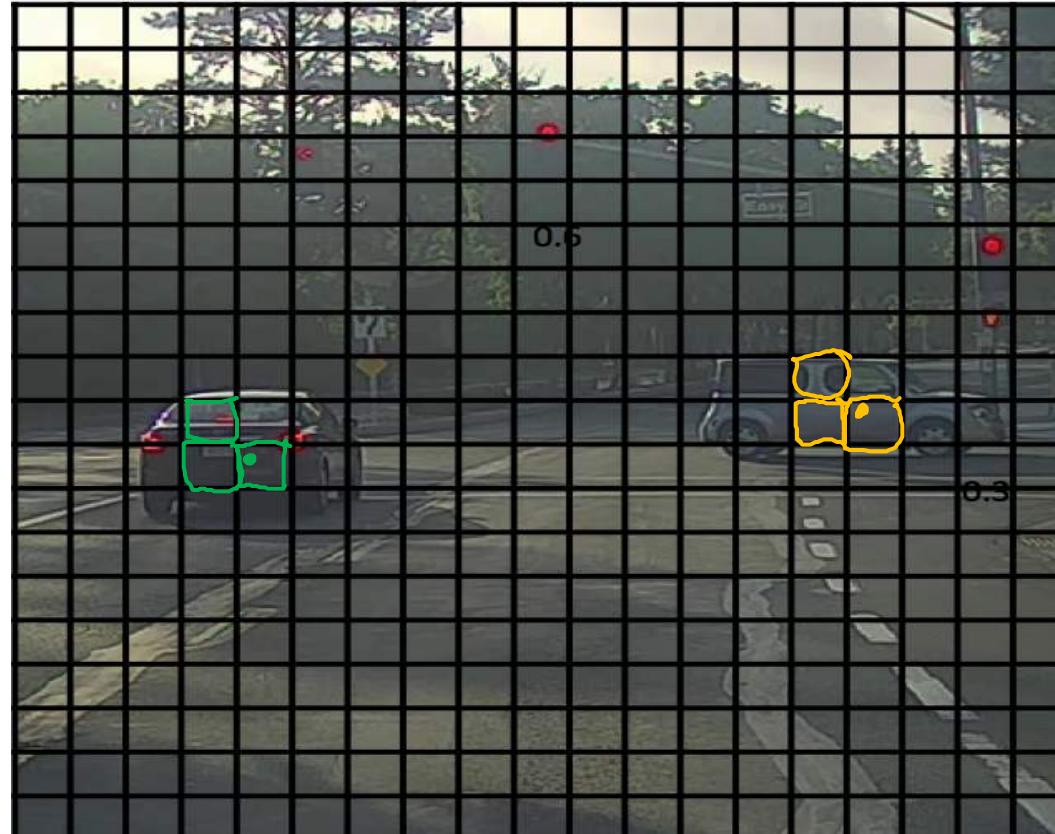
suppression

and eliminate all the BBs overlapping with it; then continue with this process until all the overlaps are solved.

Non-max suppression example

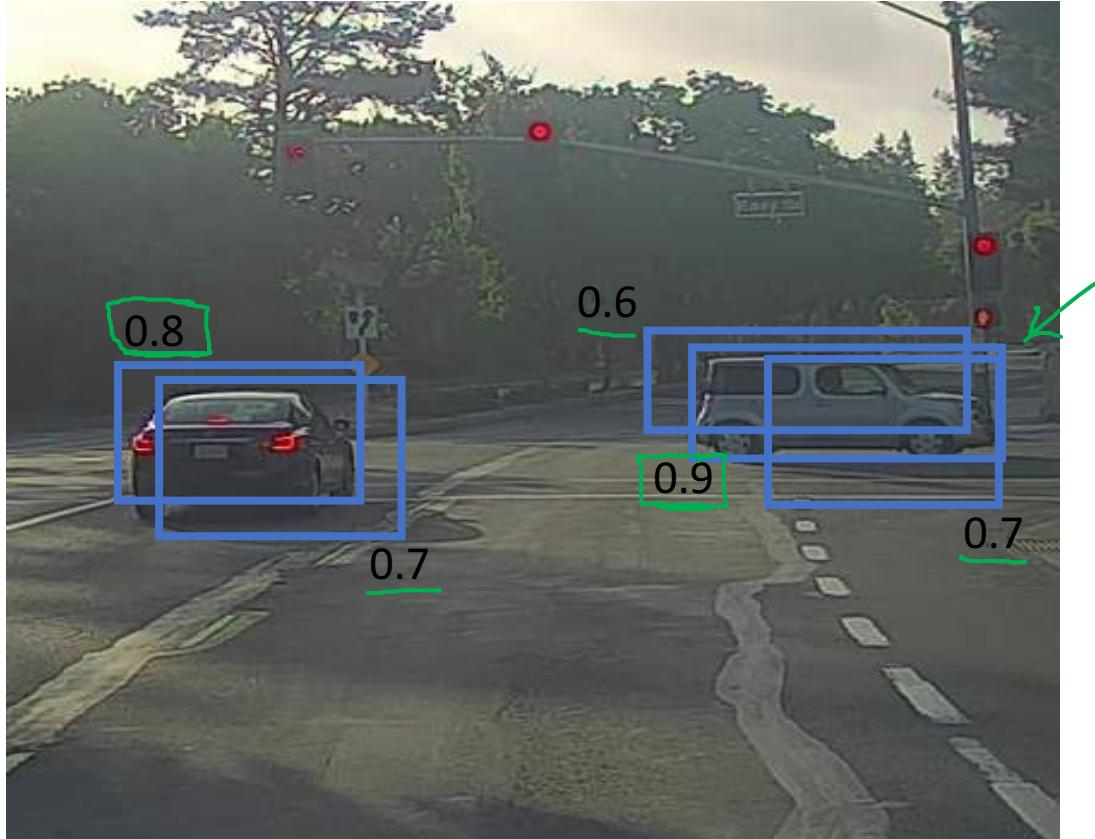


Non-max suppression example

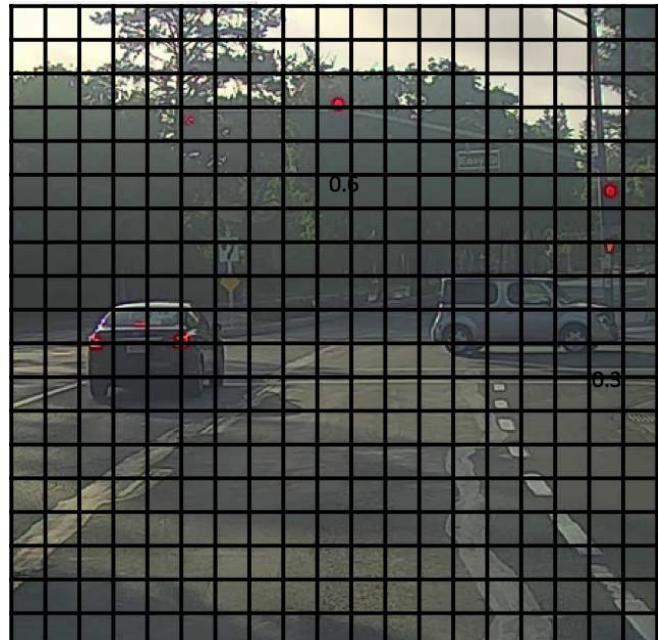


19x19

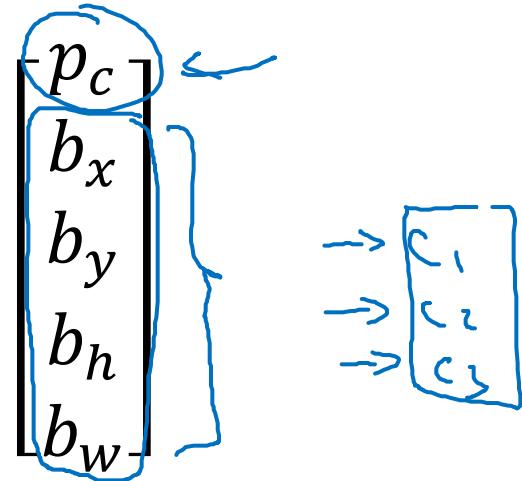
Non-max suppression example



Non-max suppression algorithm



Each output prediction is:



Discard all boxes with $\underline{p_c \leq 0.6}$

→ While there are any remaining boxes:

- Pick the box with the largest $\underline{p_c}$
Output that as a prediction.
- Discard any remaining box with
 $\underline{\text{IoU} \geq 0.5}$ with the box output
in the previous step

How to detect multiple objects in each grid cell?

The idea is to set the network to try identify n (2 in this example) distinct objects in each grid cell of the image. It is assumed that the n objects have very different shapes, and can therefore be distinguished from one another by setting up the network to

Object Detection

look up for bounding boxes of very different shapes
(referred to as anchor boxes) in each grid cell.



deeplearning.ai

Given that the distinction in anchor boxes shapes is encoded in the training labels and learned by the network. That is, the training data is labeled so

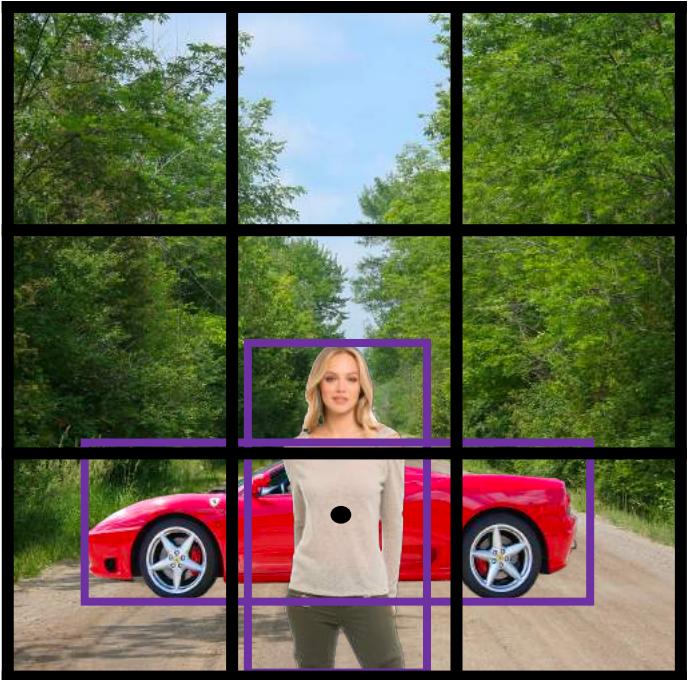
Anchor boxes

? It remains unclear to see how these boxes of different shapes are parameterized within the network.

I given that the distinction in anchor boxes shapes is encoded in the training labels and learned by the network. That is, the training data is labeled so

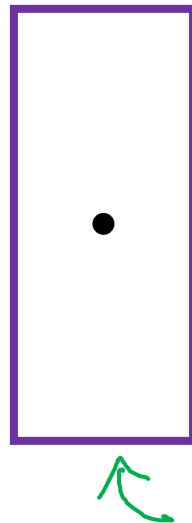
$$y = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \begin{array}{l} \text{--- tall and skinny object reported here (if any)} \\ \text{--- short and wide object reported here (if any)} \end{array}$$

Overlapping objects:

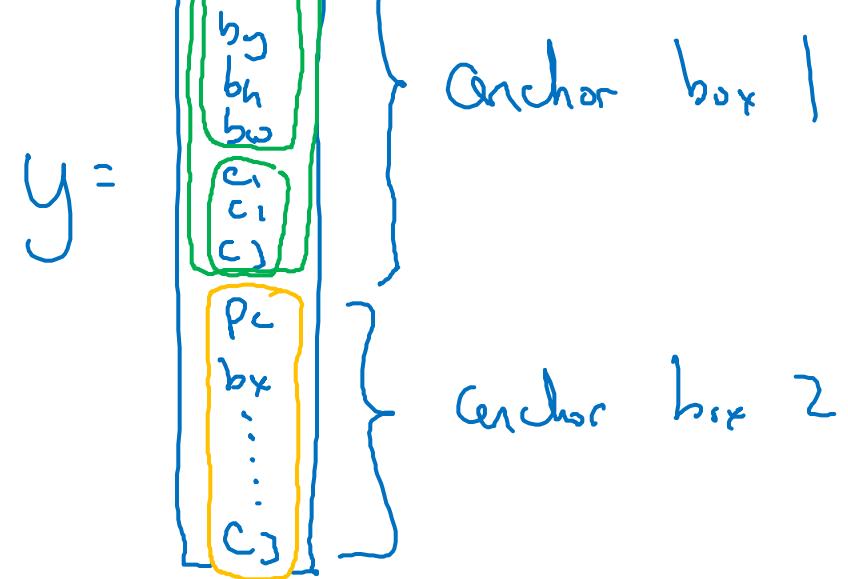
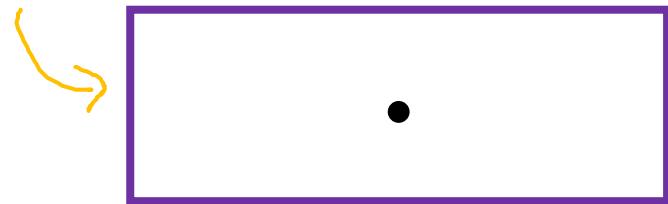


$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Anchor box 1:



Anchor box 2:



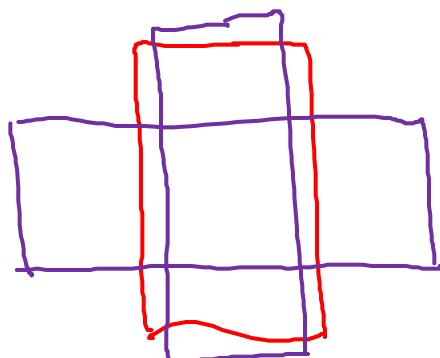
Anchor box algorithm

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint.

Output y:

$3 \times 3 \times 8$



With two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

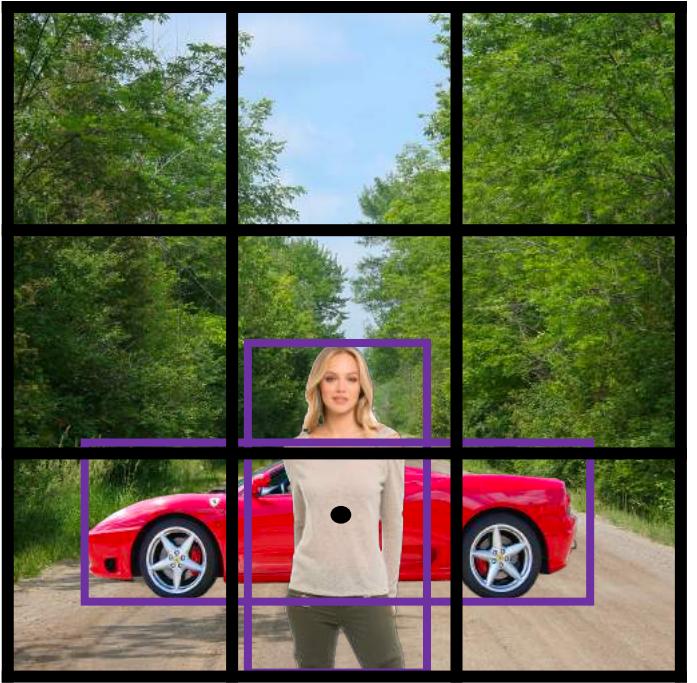
(grid cell, anchor box)

Output y:

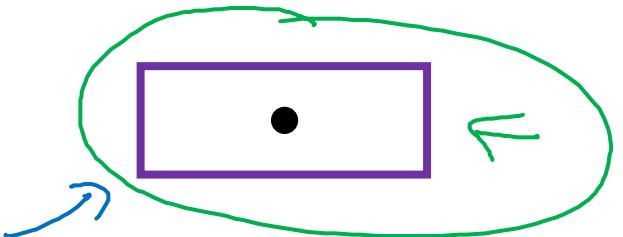
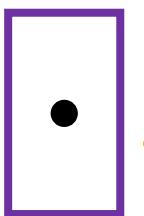
$3 \times 3 \times 16$

$3 \times 3 \times 2 \times 8$

Anchor box example



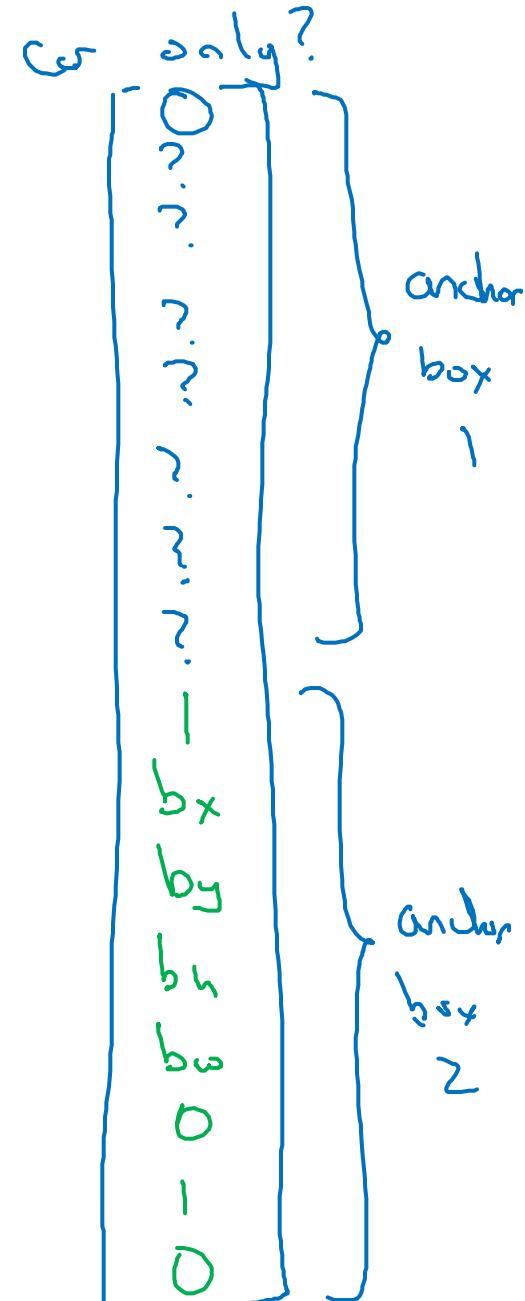
Anchor box 1: Anchor box 2:



$$y =$$

$$\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ - \\ 0 \\ 0 \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$



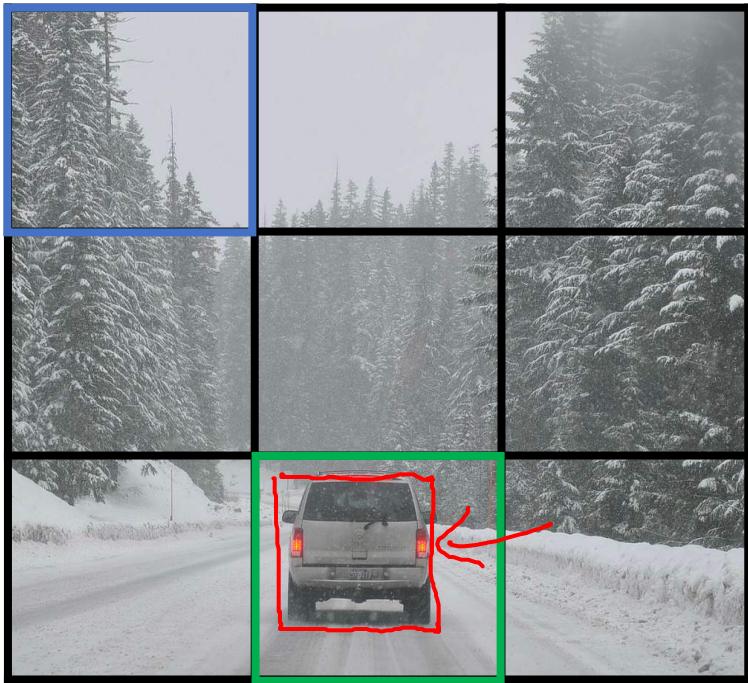


deeplearning.ai

Object Detection

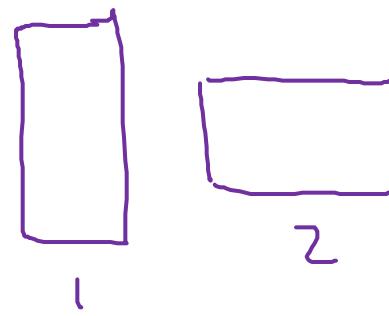
Putting it together:
YOLO algorithm

Training



- 1 - pedestrian
- 2 - car ←
- 3 - motorcycle

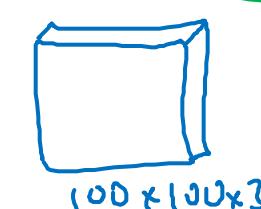
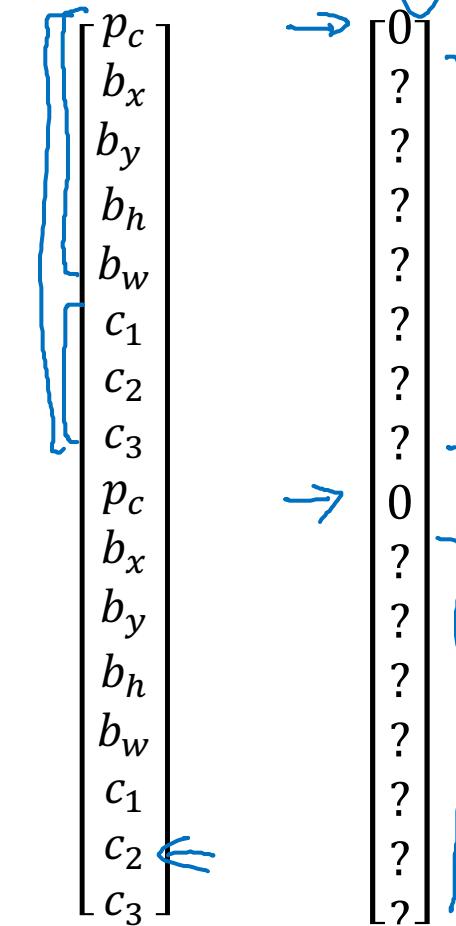
$y =$



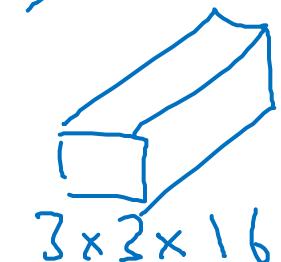
y is $3 \times 3 \times 2 \times 8$

$19 \times 19 \times 16$
 $19 \times 19 \times 40$

#anchors → $5 + \#classes$

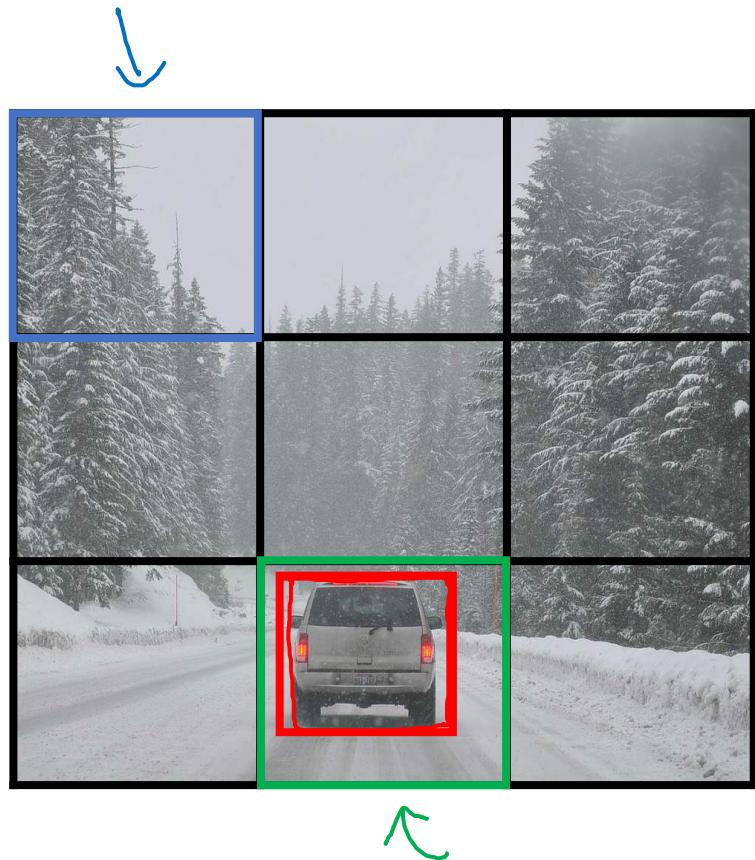


→ ConvNet



Andrew Ng

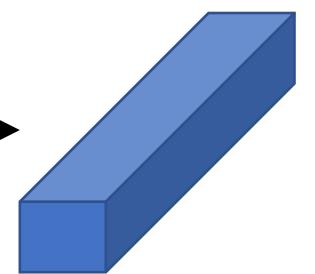
Making predictions



→

...

→



$y =$

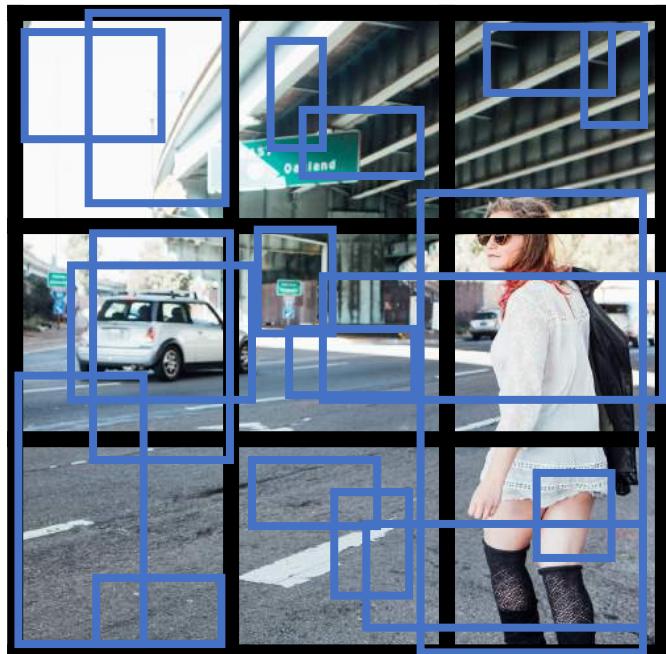
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Diagram illustrating the output vector y for object detection. The vector is composed of 17 elements, grouped into three main categories:

- Object Class Probability (p_c):** Represented by blue circles in the first row.
- Object Bounding Box Coordinates (b_x, b_y, b_h, b_w):** Represented by red numbers in the second row.
- Class-specific Feature Vectors (c_1, c_2, c_3):** Represented by green dashed lines in the third row.

Arrows point from the labels to their corresponding elements in the vector. A blue arrow points to the first p_c element, a green arrow points to the first c_1 element, and another green arrow points to the first b_x element.

Outputting the non-max suppressed outputs



- For each grid call, get 2 predicted bounding boxes.
- Get rid of low probability predictions.
- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.

In order to reduce the number of grids in which to check for objects , it was also proposed to use REGION PROPOSALS. The idea is to preliminary segment the image , then run the object recognition only in the regions where the segmentation found some shapes of interest.



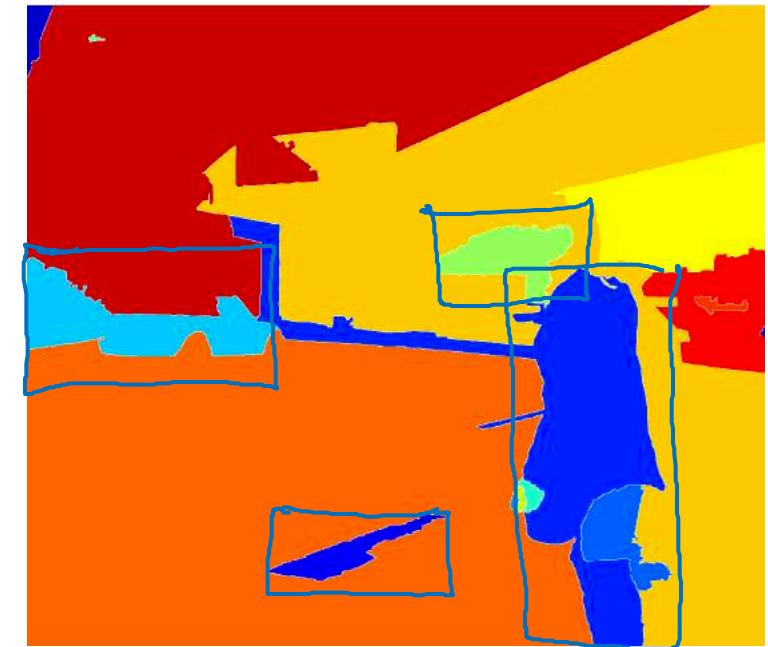
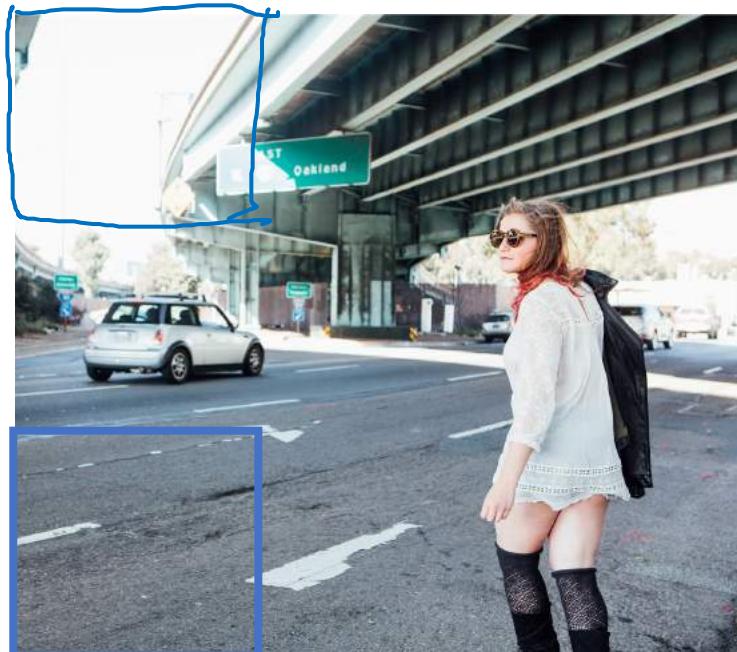
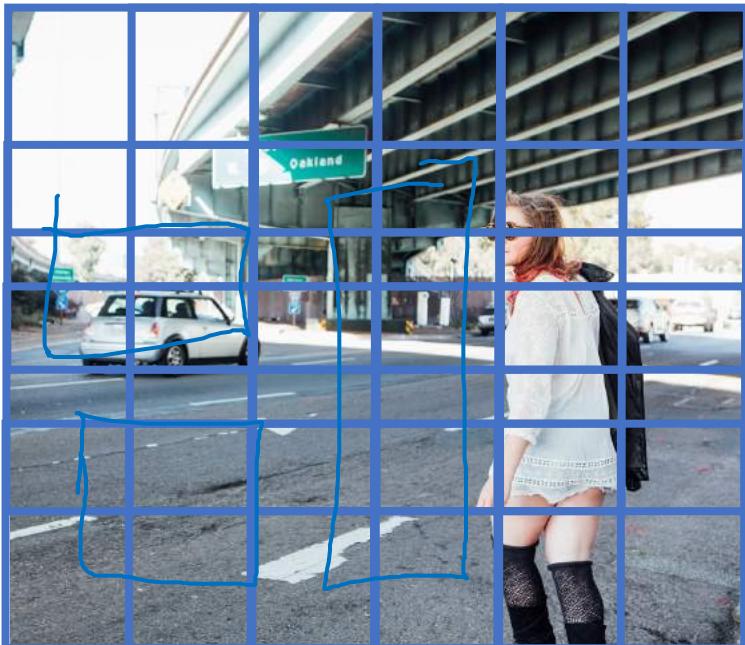
deeplearning.ai

Object Detection

NG prefers "one step/one look" solutions such as YOLO

Region proposals (Optional)

Region proposal: R-CNN



Segmentation algorithm

~2,000

Faster algorithms

→ R-CNN: Propose regions. Classify proposed regions one at a time. Output label + bounding box. ←

Fast R-CNN: Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions. ←

Faster R-CNN: Use convolutional network to propose regions.
↳ actually *STILL SLOWER THAN YOLO*

[Girshik et. al, 2013. Rich feature hierarchies for accurate object detection and semantic segmentation]

[Girshik, 2015. Fast R-CNN]

[Ren et. al, 2016. Faster R-CNN: Towards real-time object detection with region proposal networks]

Andrew Ng



deeplearning.ai

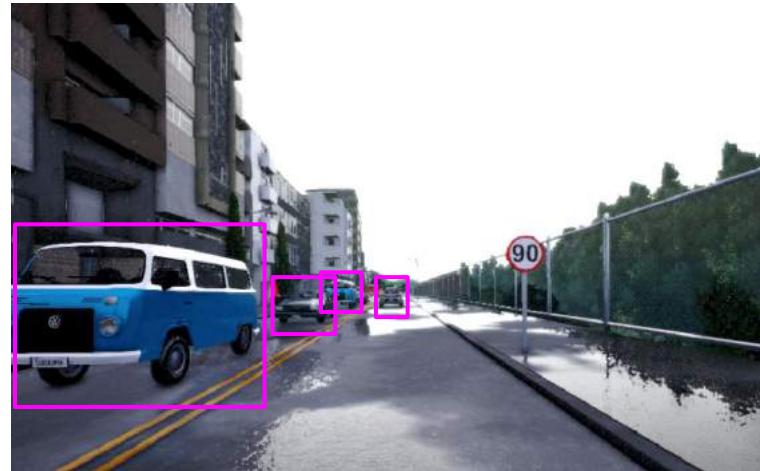
Convolutional Neural Networks

Semantic segmentation with U-Net

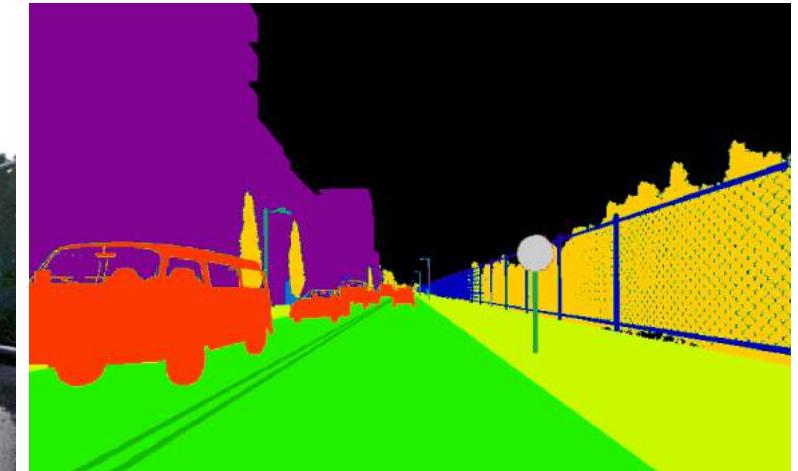
Object Detection vs. Semantic Segmentation



Input image

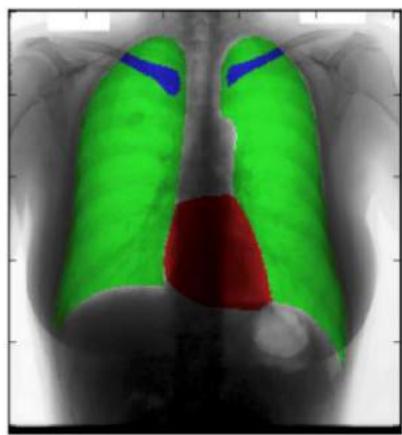


Object Detection

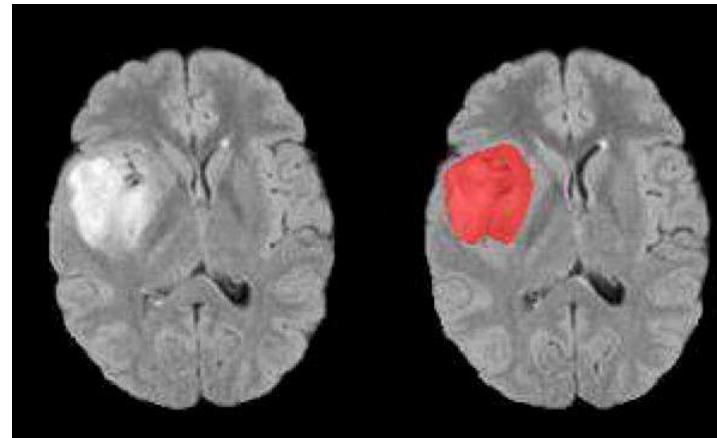


Semantic Segmentation

Motivation for U-Net



Chest X-Ray



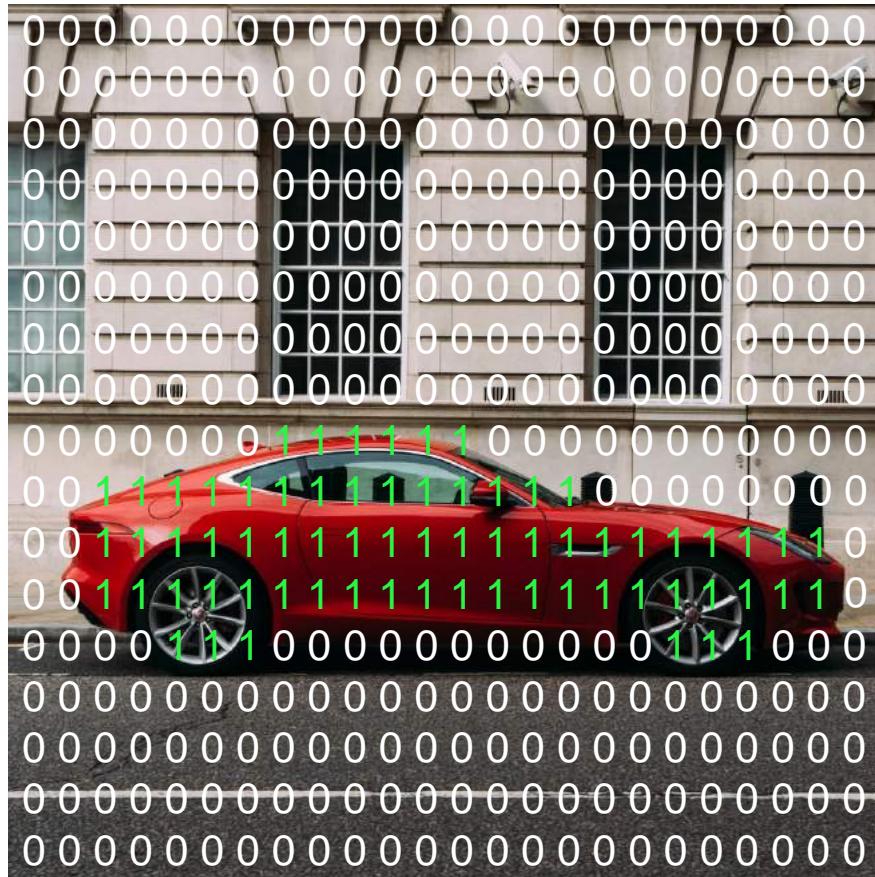
Brain MRI

[Novikov et al., 2017, Fully Convolutional Architectures for Multi-Class Segmentation in Chest Radiographs]

[Dong et al., 2017, Automatic Brain Tumor Detection and Segmentation Using U-Net Based Fully Convolutional Networks]

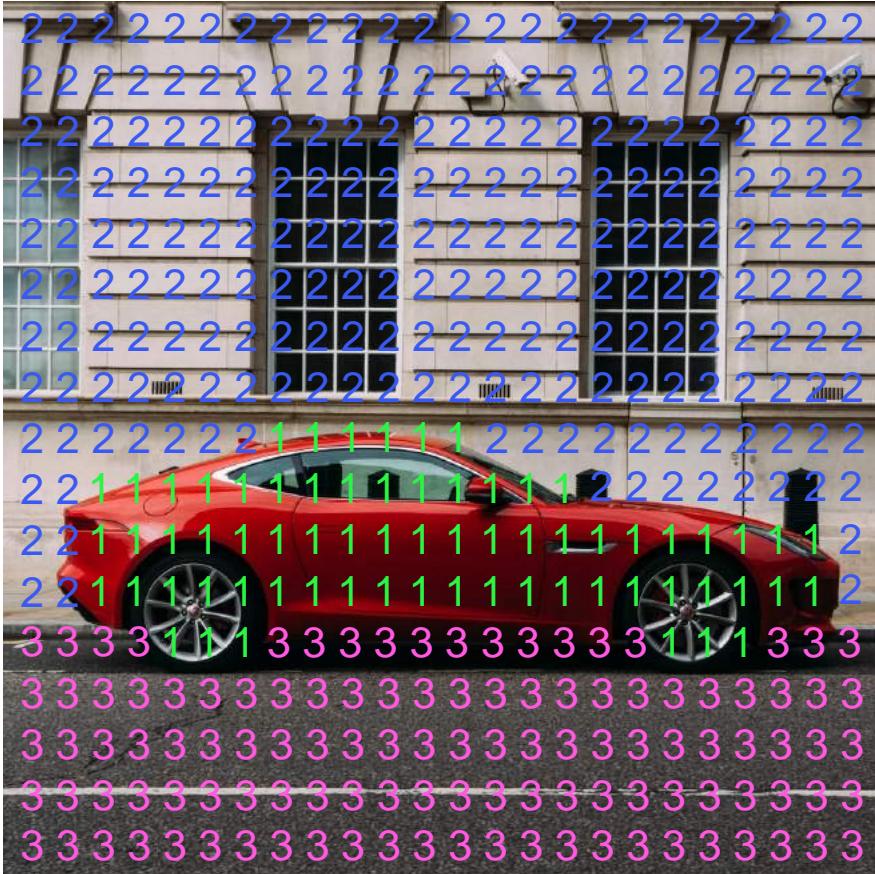
Andrew Ng

Per-pixel class labels



1. Car
0. Not Car

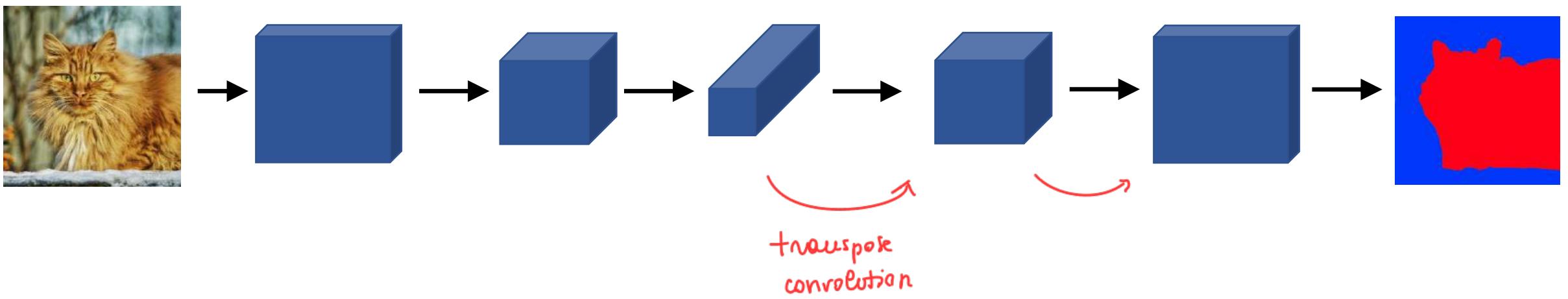
Per-pixel class labels



1. Car
 2. Building
 3. Road

Segmentation Map

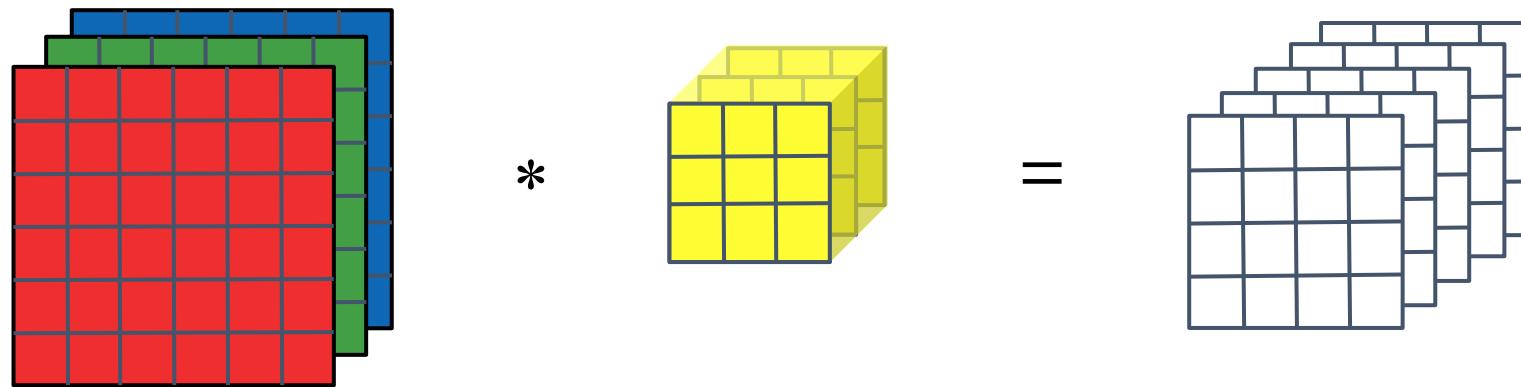
Deep Learning for Semantic Segmentation



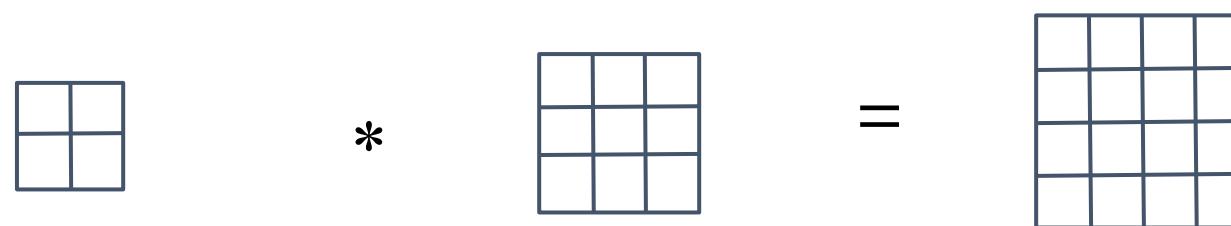
intuition: compress image to extract high-level semantic information, then reexpand that information attempting to combine it with low-level spatial information.

Transpose Convolution

Normal Convolution

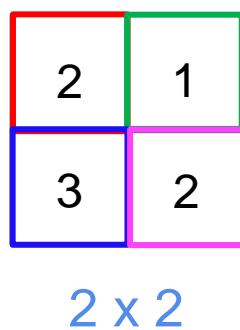


Transpose Convolution



- idea: place the filter on the output.
- Define stride and padding, which will be applied on the output
 - Take the first element in the input, multiply it by the filter, overlap the filter with the output (account for padding!), copy filter values.

Transpose Convolution

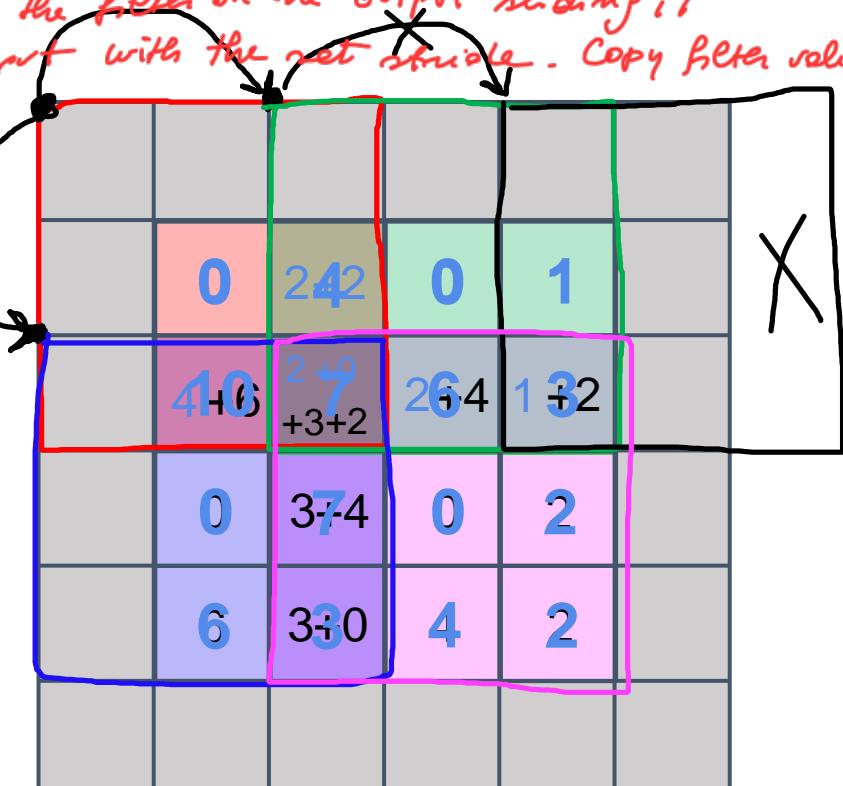


1	2	1
2	0	1
0	2	1

weight filter

• Repeat for second element of the input, place the filter on the output sliding it in the same direction as in the input with the net stride. Copy filter values but add values to pre-existing values in the output.

• Continue.

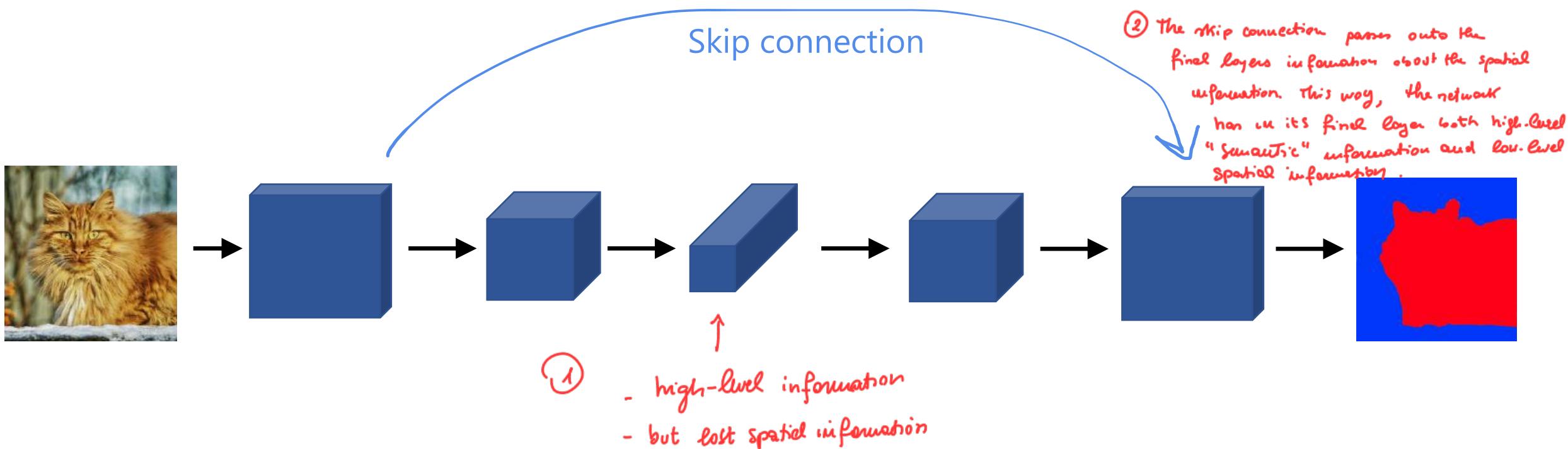


filter $f \times f = 3 \times 3$

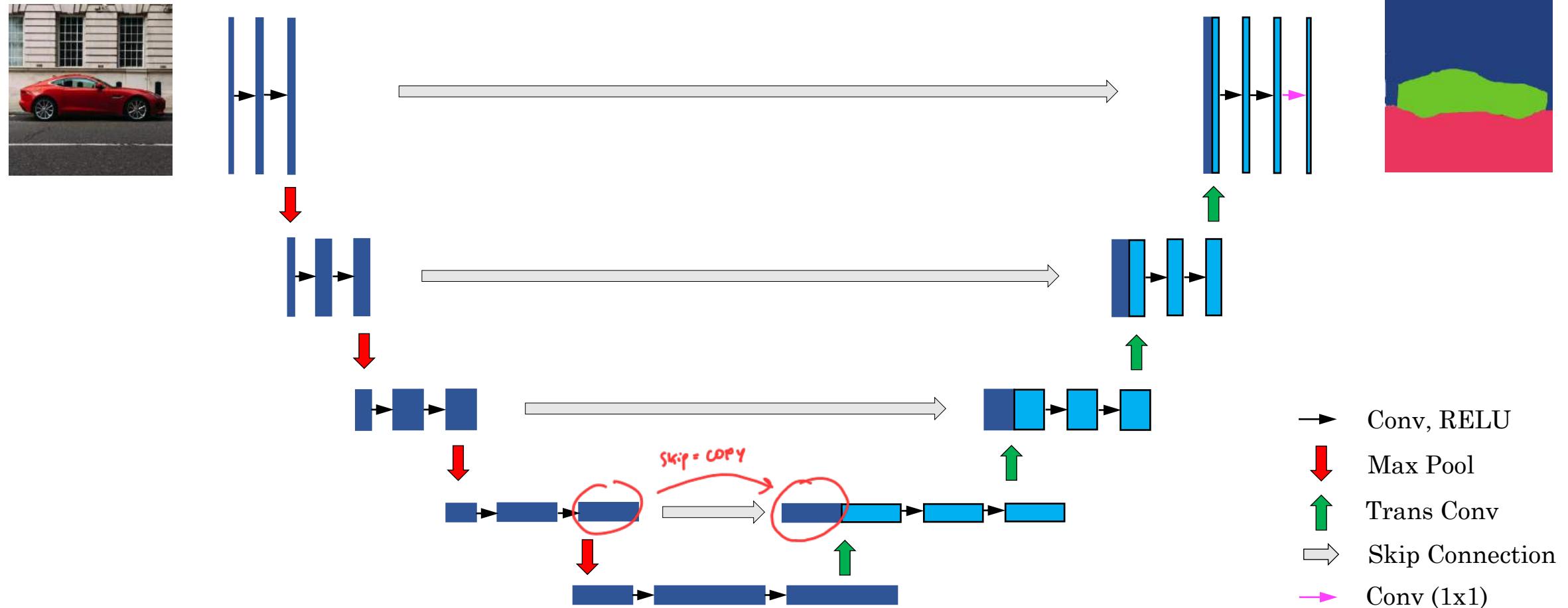
padding $p = 1$

stride $s = 2$

Deep Learning for Semantic Segmentation



U-Net



U-Net

