



# Object Detection

---

Object  
localization

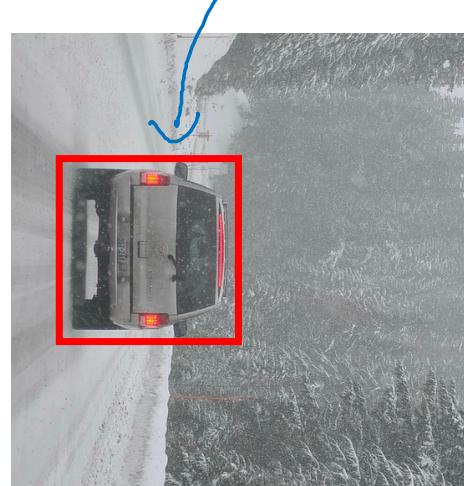
deeplearning.ai

# What are localization and detection?

Image classification



Classification with  
localization



Detection

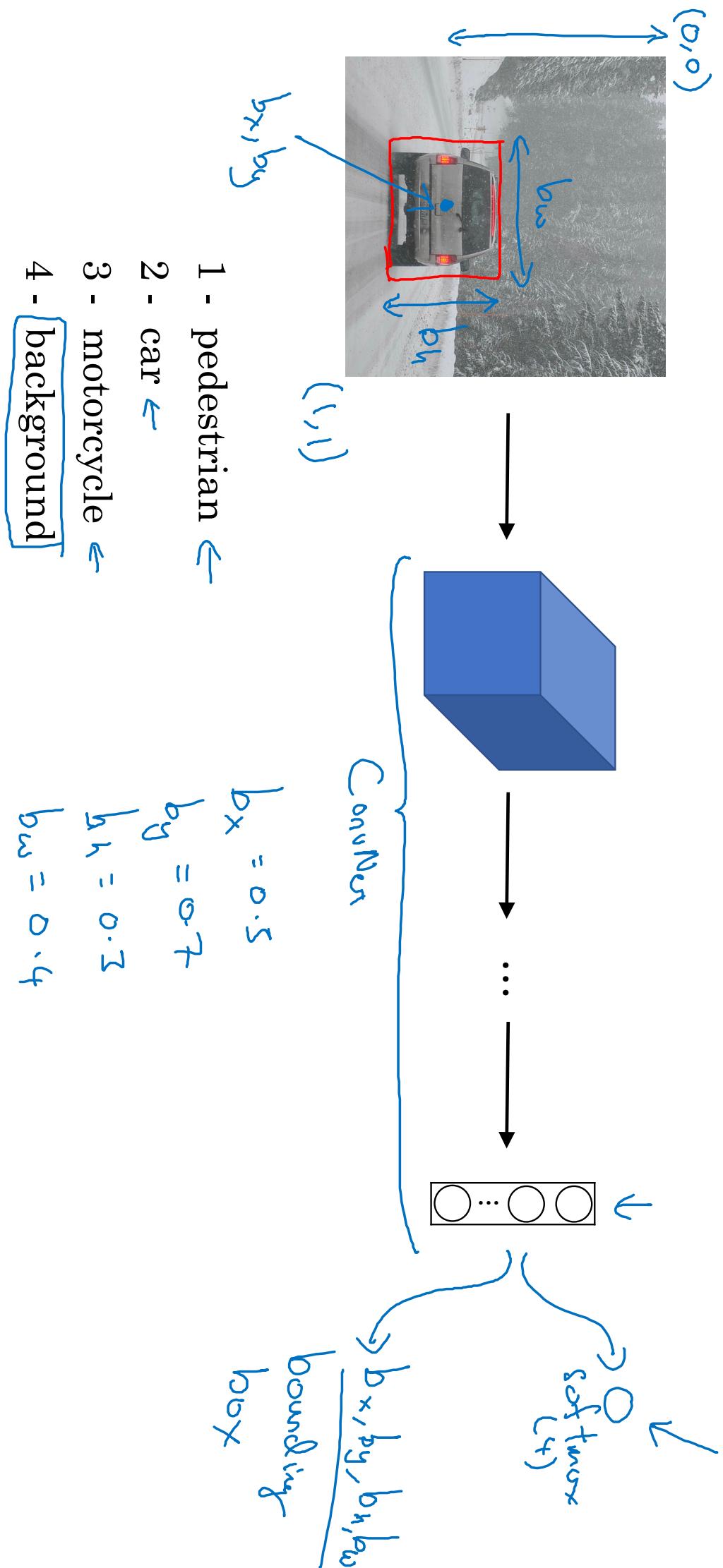


“Car”  
1 object

“Car”

multiple  
objects

# Classification with localization



# Defining the target label $y$

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

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

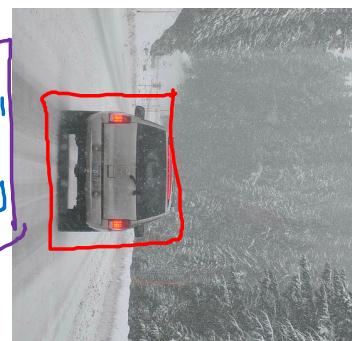
$$l(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 \\ + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

$\rightarrow y =$

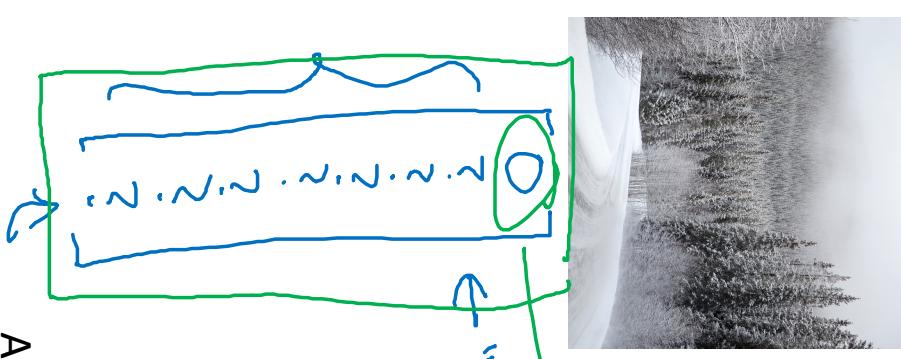
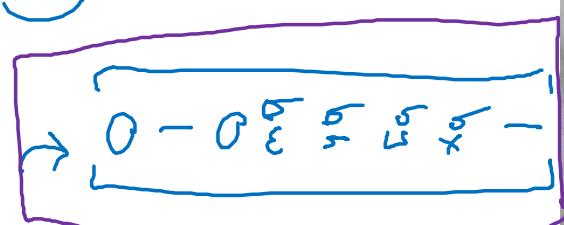
$$\left[ \begin{array}{c} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{array} \right]$$

is there  
object?

$x =$



$(x, y)$



"don't  
care"



# Object Detection

---

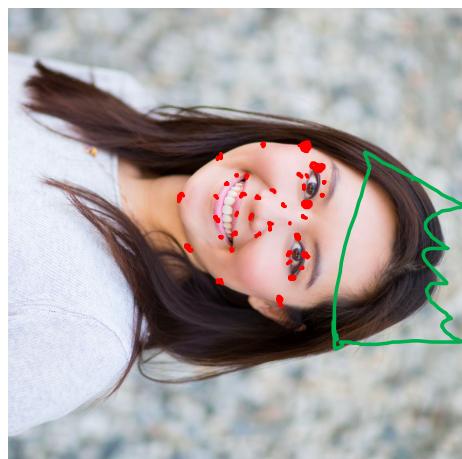
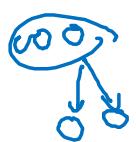
Landmark  
detection

deeplearning.ai

# Landmark detection



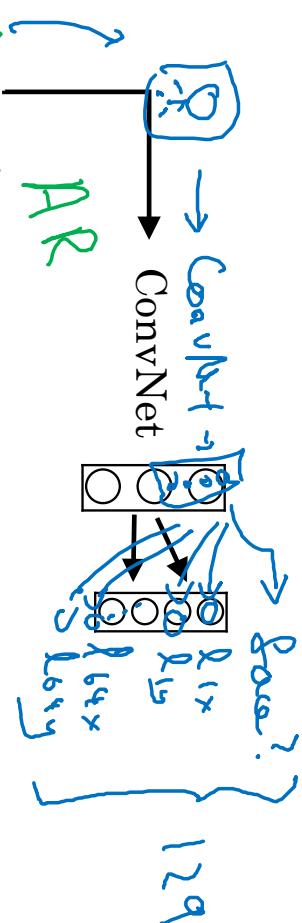
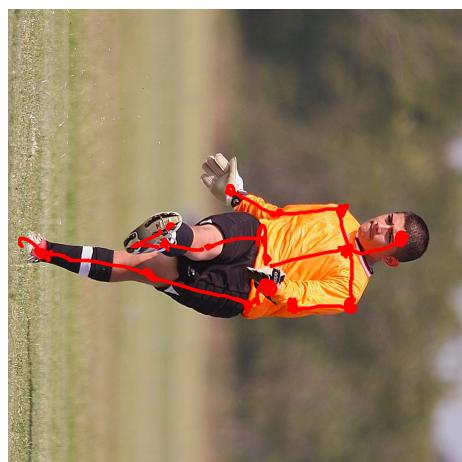
$b_x, b_y, b_h, b_w$



$\{x_1, y_1\}$ ,  
 $\{x_2, y_2\}$ ,  
 $\{x_3, y_3\}$ ,  
 $\{x_4, y_4\}$ ,  
 $\vdots$   
 $\{x_{64}, y_{64}\}$

X, Y

$\{x_1, y_1\}$ ,  
 $\vdots$   
 $\{x_{32}, y_{32}\}$





deeplearning.ai

---

# Object Detection

---

Object  
detection

# Car detection example

Training set:

X  
Y



0 0 1 1 1

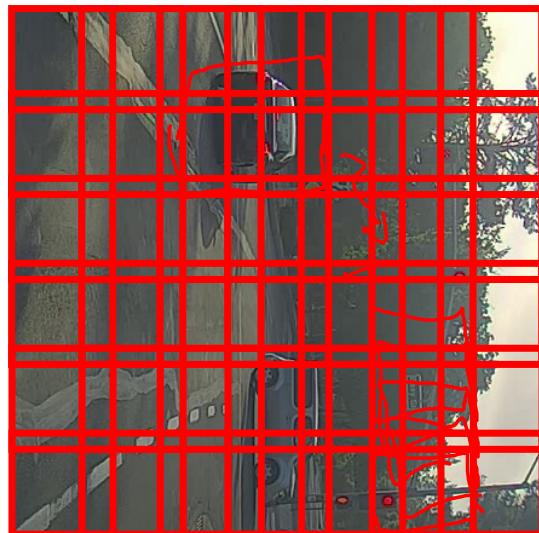


→ ConvNet → y

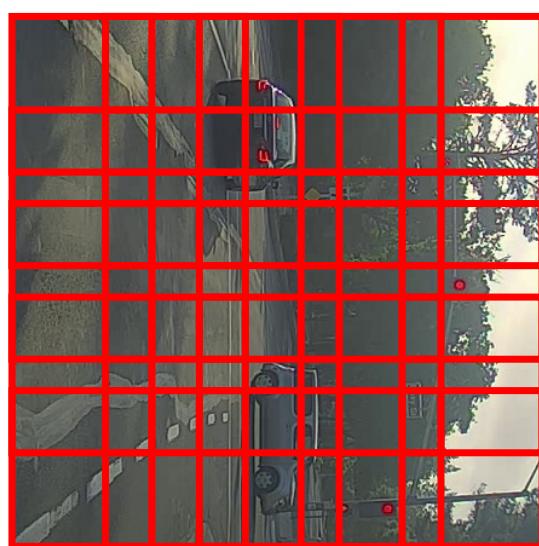


# Sliding windows detection

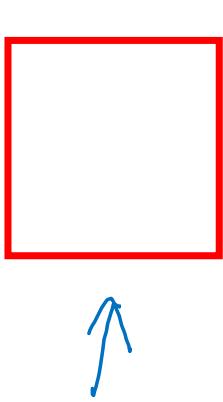
→ ConvNet → 0



→ ConvNet



Computational cost





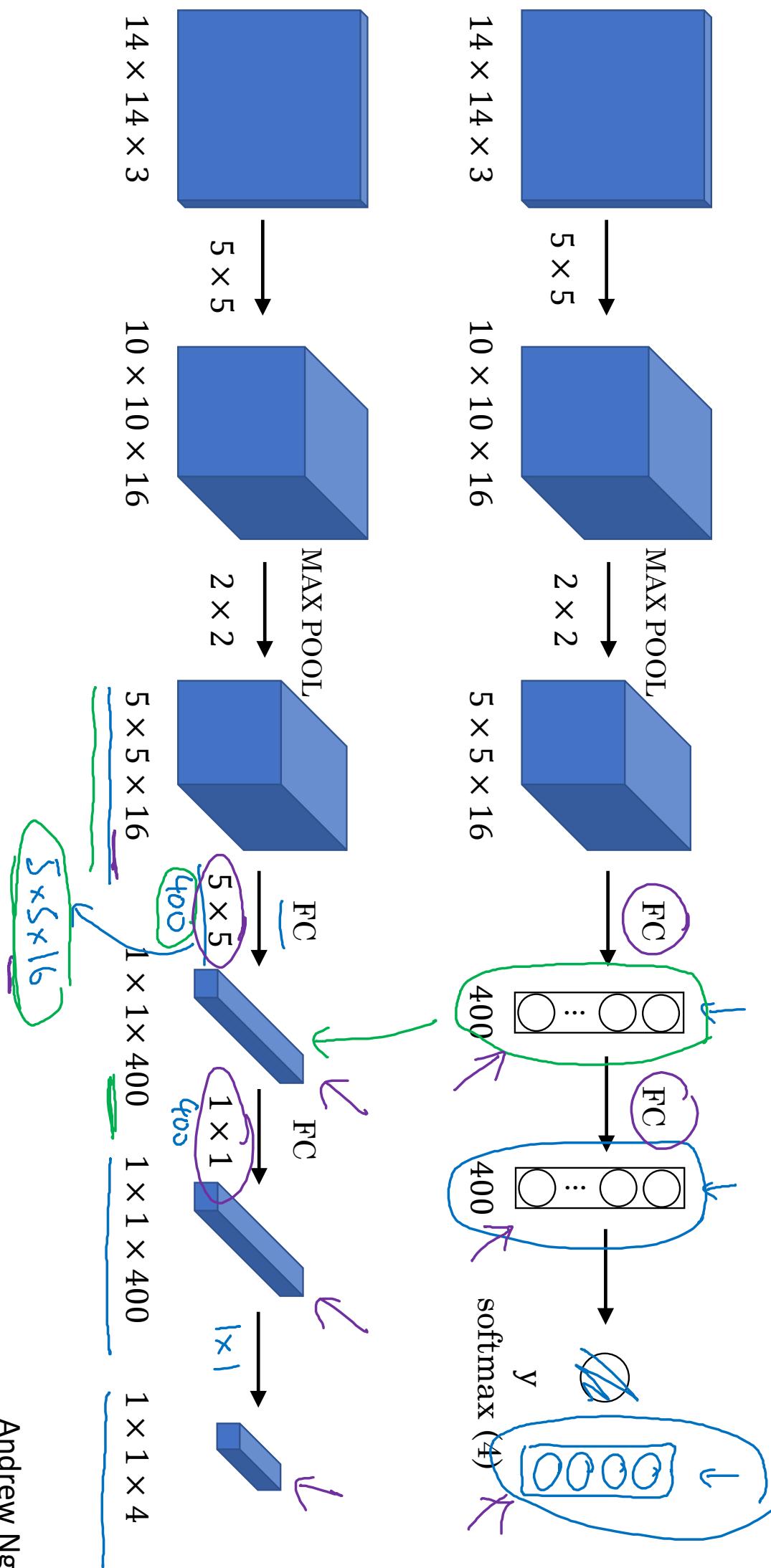
# Object Detection

---

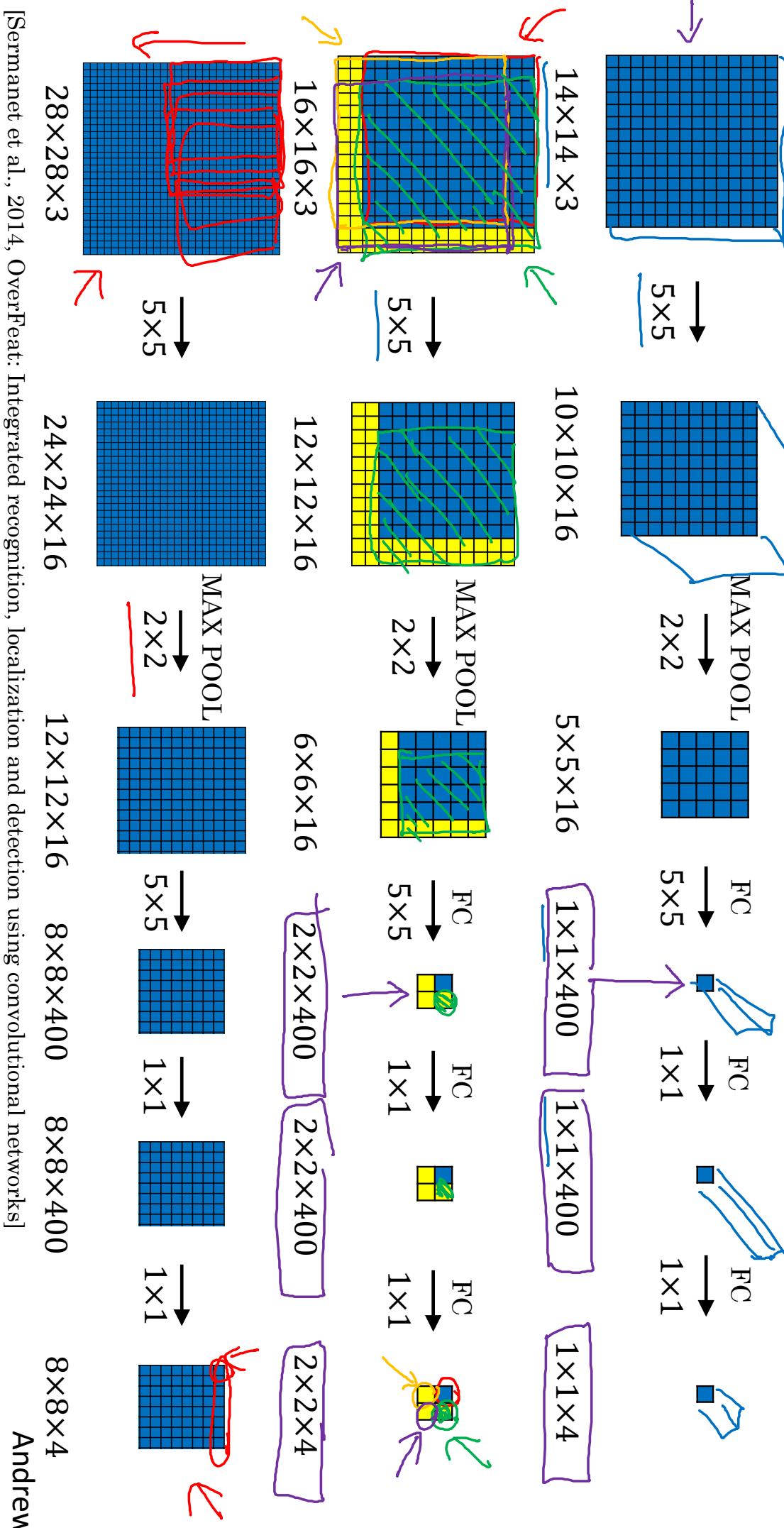
Convolutional  
implementation of  
sliding windows

deeplearning.ai

# Turning FC layer into convolutional layers



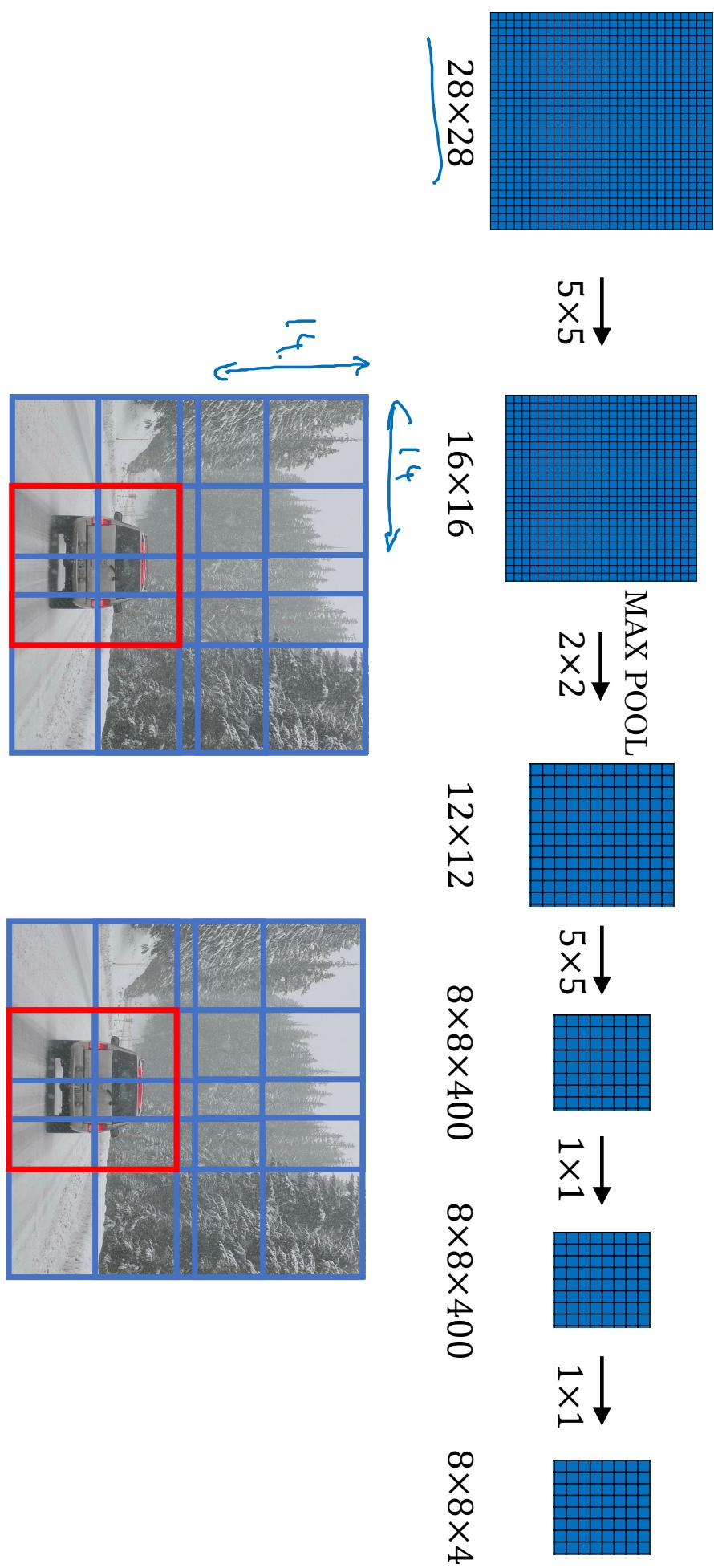
# Convolution implementation of sliding windows



[Sermanet et al., 2014, OverFeat: Integrated recognition, localization and detection using convolutional networks]

Andrew Ng

# Convolution implementation of sliding windows





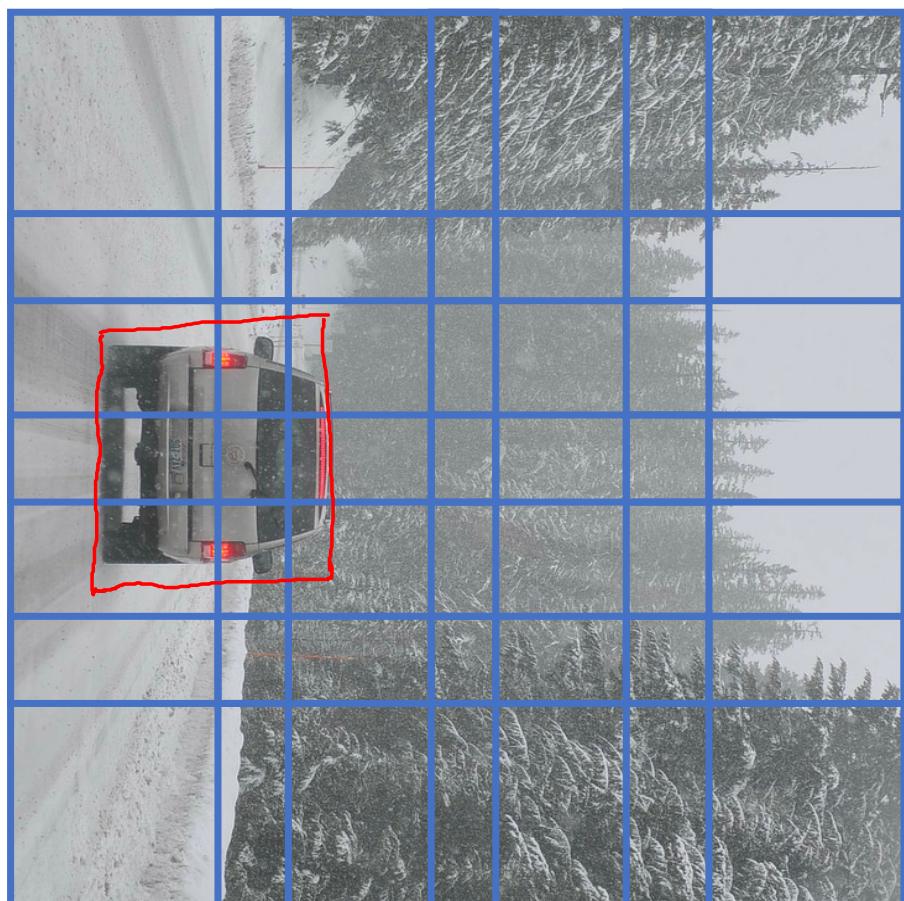
# Object Detection

---

Bounding box  
predictions

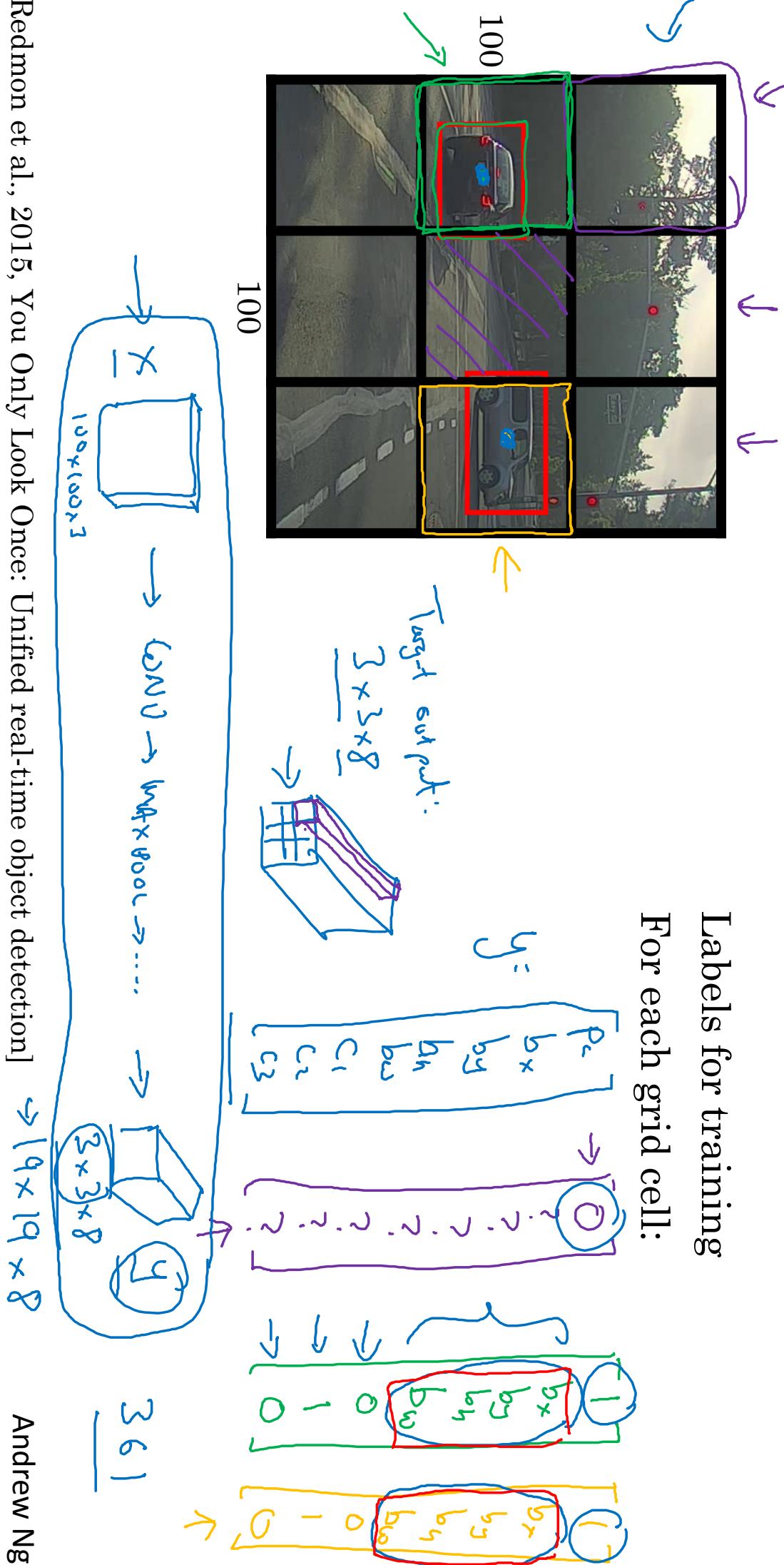
deeplearning.ai

# Output accurate bounding boxes

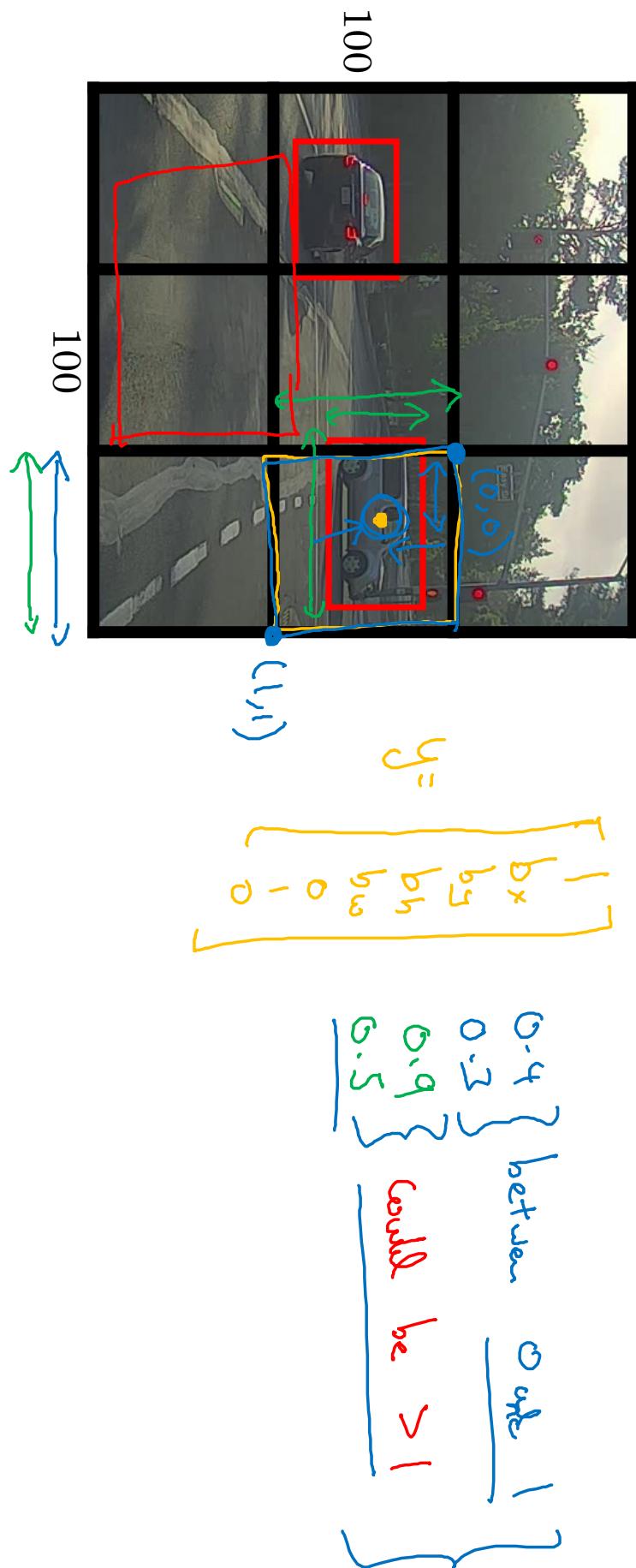


# YOLO algorithm

Labels for training  
For each grid cell:



Specify the bounding boxes



[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]



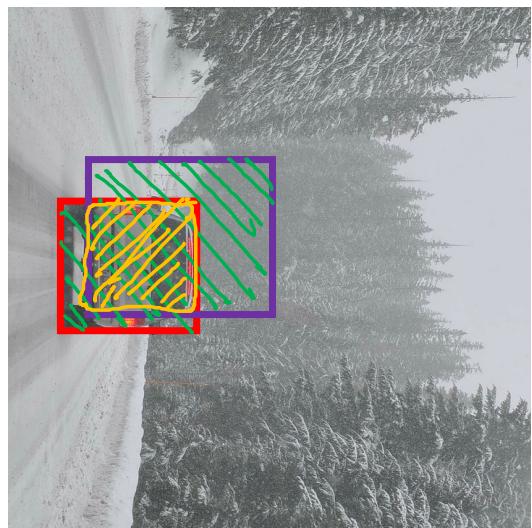
# Object Detection

---

Intersection  
over union

deeplearning.ai

# Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{Size of intersection}}{\text{Size of union}}$$

"Correct" if  $\text{IoU} \geq 0.5$

0.6  $\rightarrow$

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



# Object Detection

---

Non-max  
suppression

deeplearning.ai

# Non-max suppression example

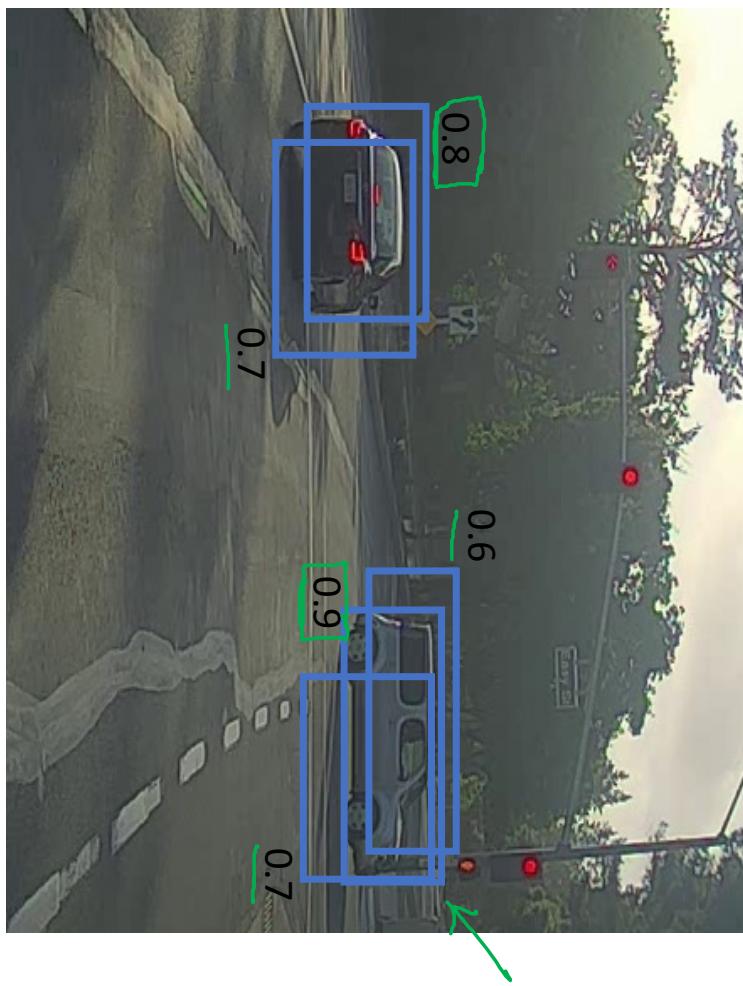


# Non-max suppression example

$19 \times 19$

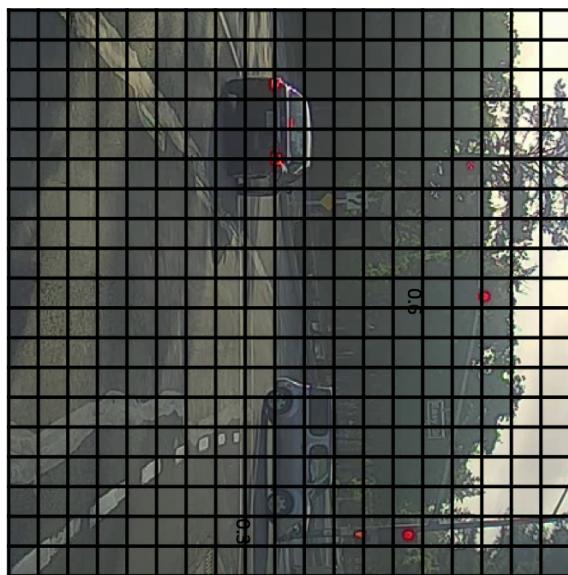


# Non-max suppression example

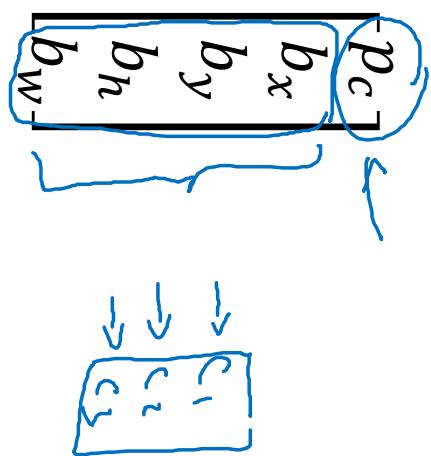


$P_c$

# Non-max suppression algorithm



Each output prediction is:



Discard all boxes with  $p_c \leq 0.6$

→ While there are any remaining boxes:

- Pick the box with the largest  $p_c$
- Output that as a prediction.

- Discard any remaining box with  $\text{IoU} \geq 0.5$  with the box output in the previous step



deeplearning.ai

# Object Detection

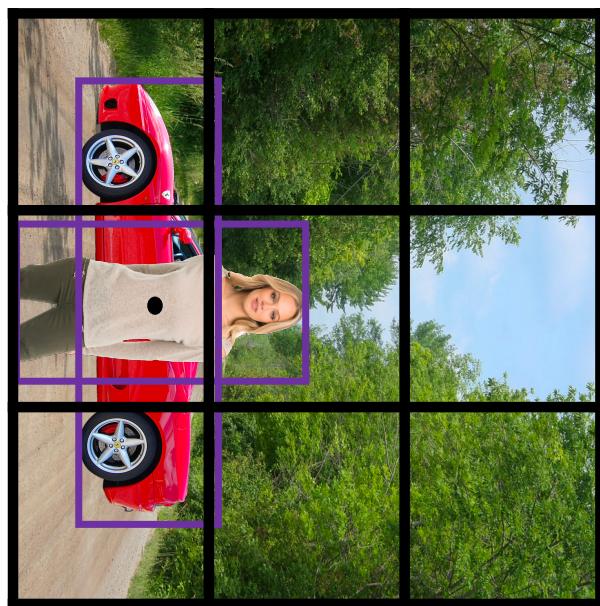
---

## Anchor boxes

---

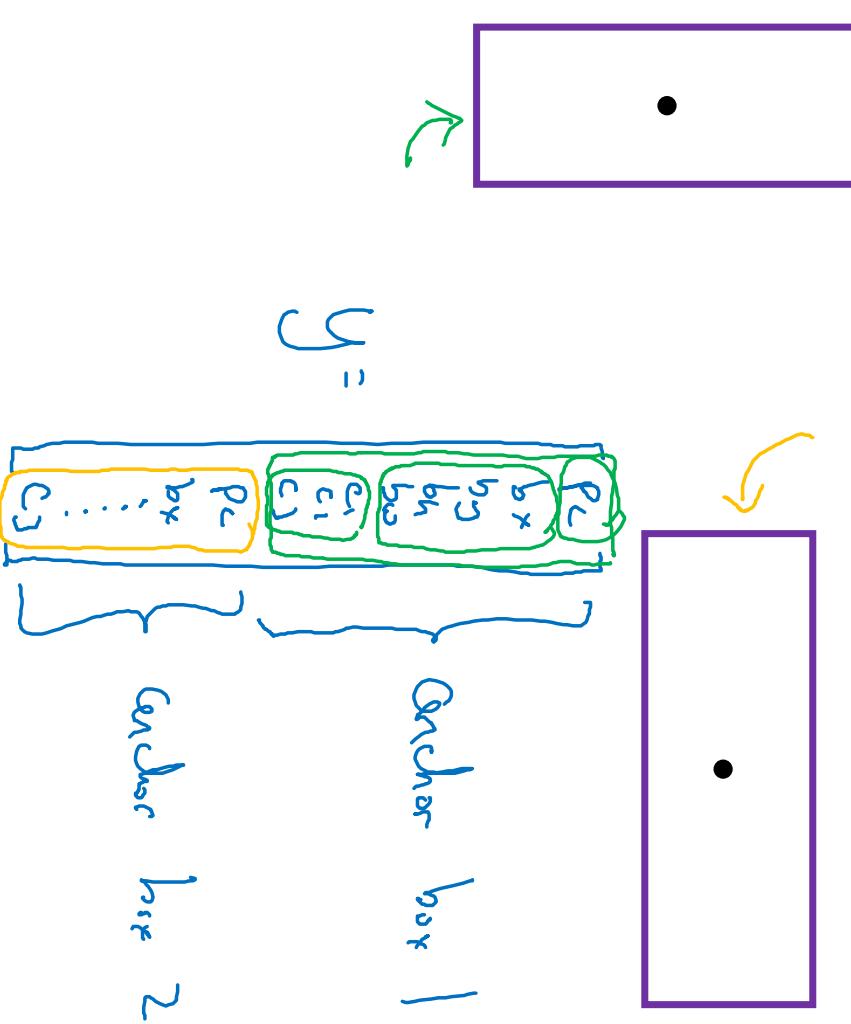
# Overlapping objects:

Anchor box 1:  
Anchor box 2:



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

Annotations: A blue brace groups the last three elements ( $c_1, c_2, c_3$ ). A blue curly brace groups the first four elements ( $p_c, b_x, b_y, b_h, b_w$ ). A green curly brace groups the last two elements ( $b_x, b_y$ ).

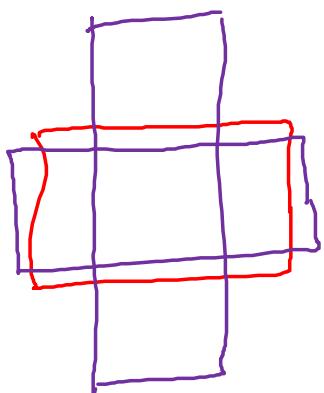


[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]

# Anchor box algorithm

Previously:

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



Output y:  
 $\underline{3 \times 3 \times 8}$

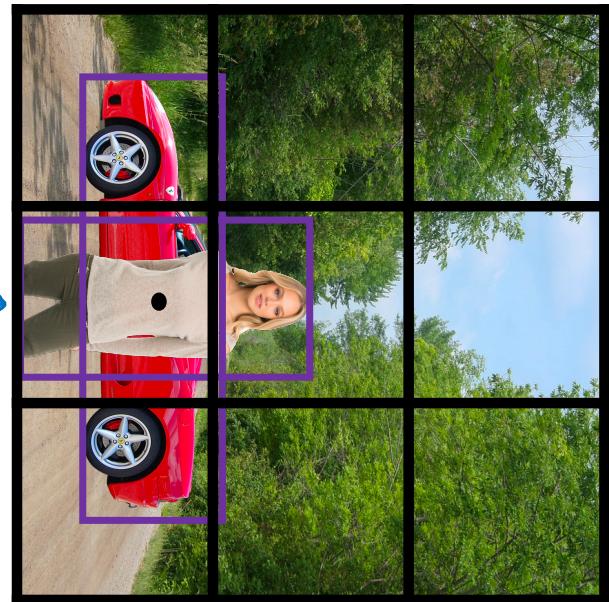
With two anchor boxes:

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:  
 $\underline{3 \times 3 \times \frac{16}{2} \times 8}$

# Anchor box example



Anchor box 1: Anchor box 2:

$$y =$$

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

$$\left[ \begin{array}{ccccccccc} 0 & 0 & b_w & b_h & b_x & b_y & b_w & b_h & b_x & b_y \\ \vdots & \vdots \end{array} \right]$$

$$\left[ \begin{array}{ccccccccc} 0 & 0 & b_w & b_h & b_x & b_y & \dots & \dots & \dots & \dots \\ \vdots & \vdots \end{array} \right]$$

Can't only?

Anchor box 1

Anchor box 2



# Object Detection

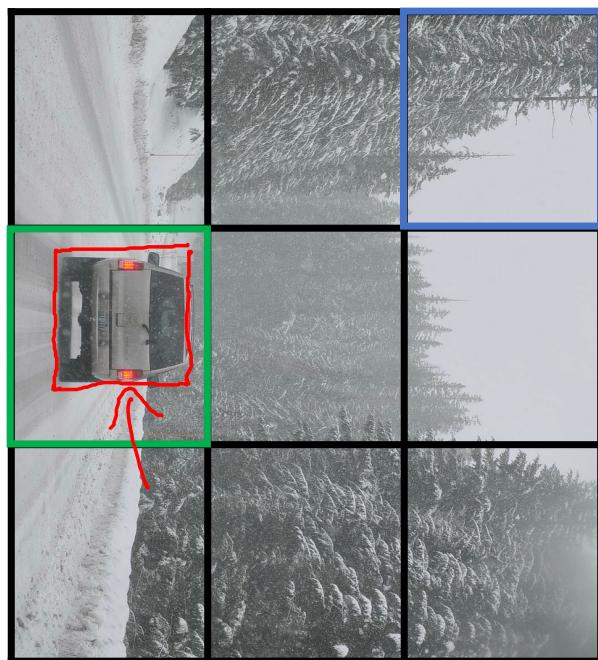
---

Putting it together:  
YOLO algorithm

# Training

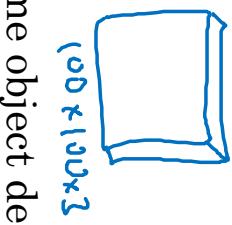
1 - pedestrian  
2 - car ←  
3 - motorcycle

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$



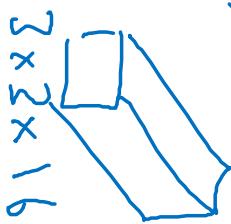
$y$  is  $3 \times 3 \times 2 \times 8$   
 $2 \times 3 \times 16$   
 $16 \times 16 \times 16$   
 $16 \times 16 \times 40$   
 $\# \text{ anchors}$

[Redmon et al., 2015, You Only Look Once: Unified real-time object detection]



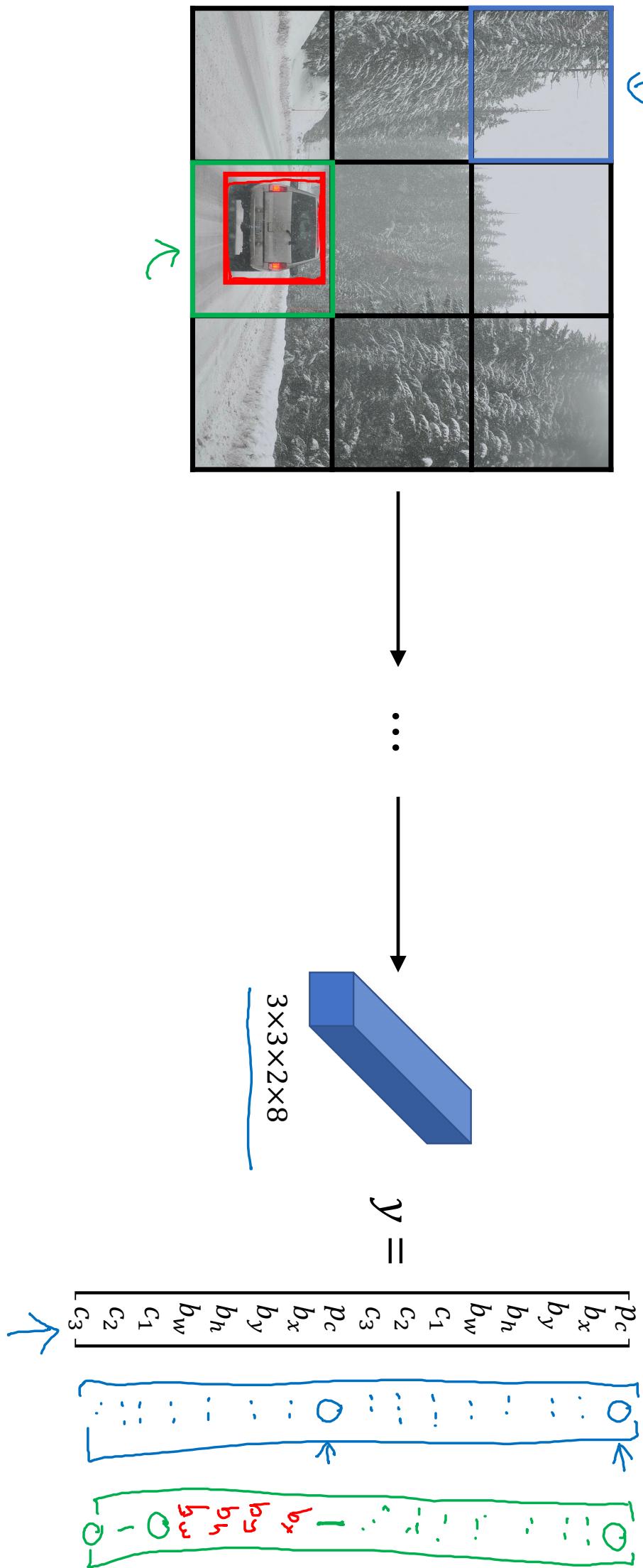
$100 \times 100 \times 3$

→ ConvNet →

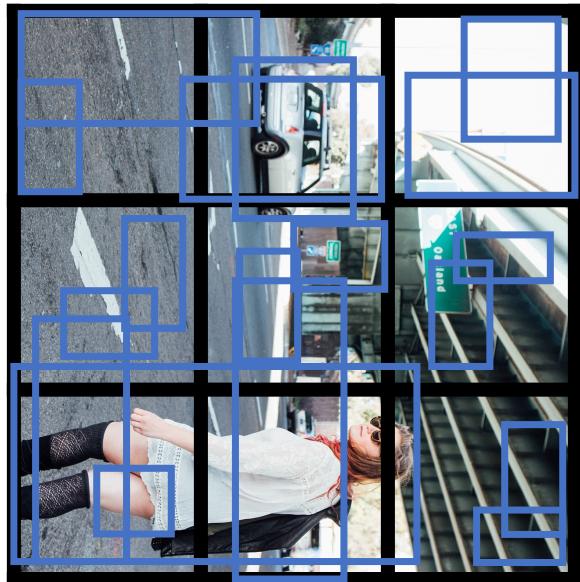


$3 \times 3 \times 16$   
Andrew Ng

# Making predictions



# Outputting the non-max suppressed outputs



- For each grid cell, 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.

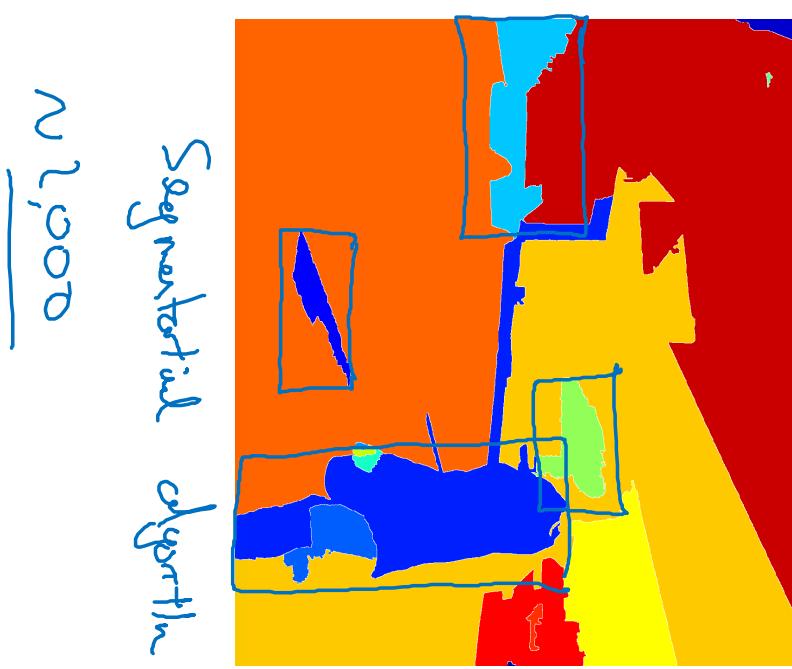
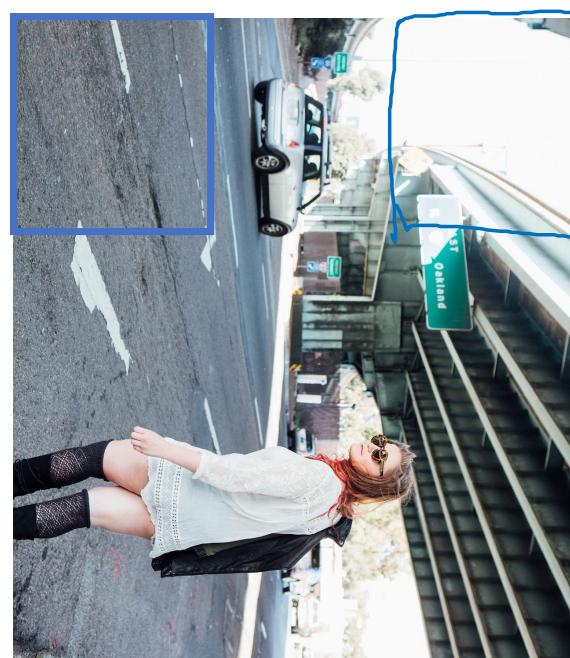
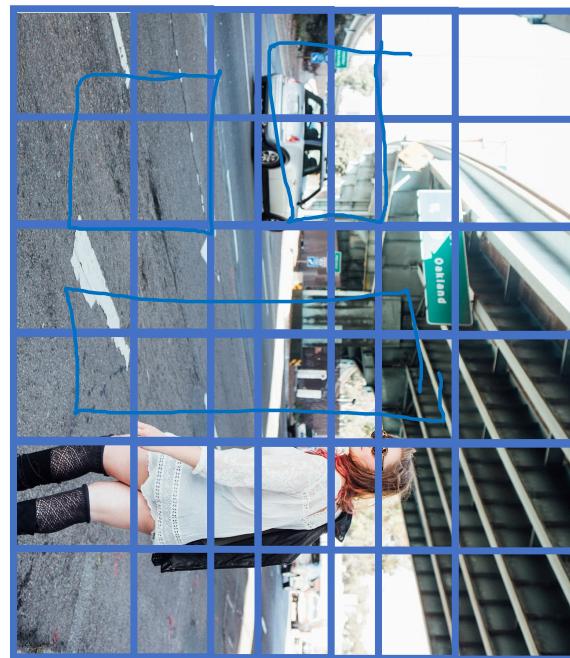


# Object Detection

Region proposals  
(optional)

deeplearning.ai

# Region proposal: R-CNN



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

# 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.

[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

# Convolutional Neural Networks

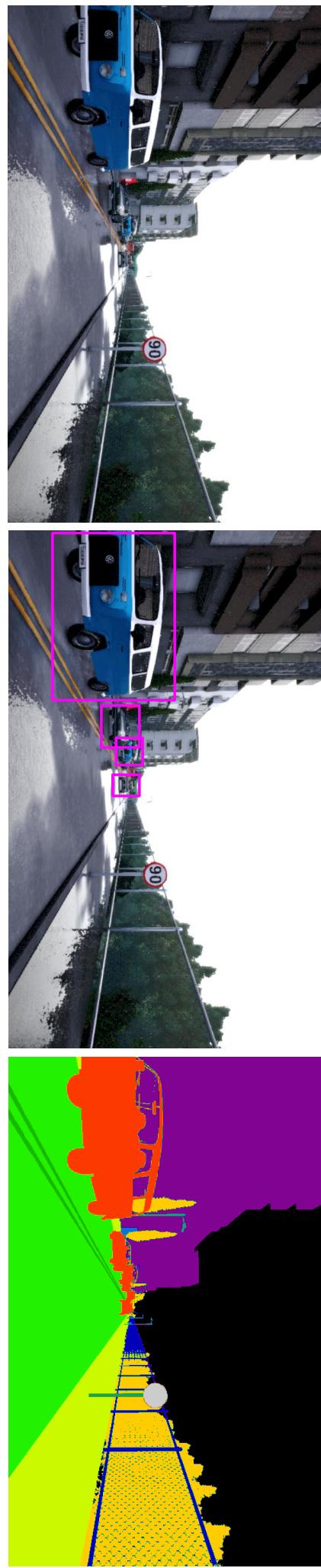
---

Semantic segmentation  
with U-Net

deeplearning.ai



# Object Detection vs. Semantic Segmentation

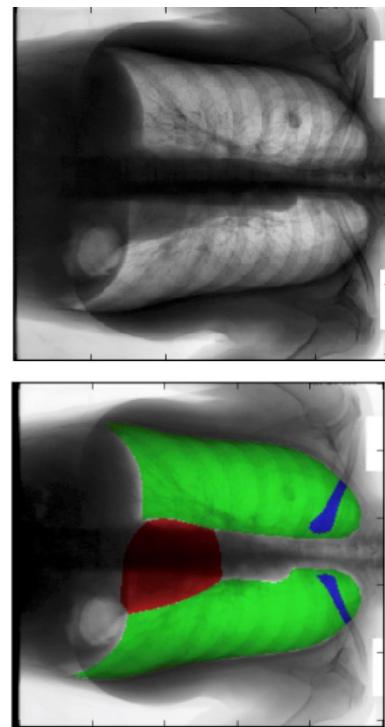


Input image

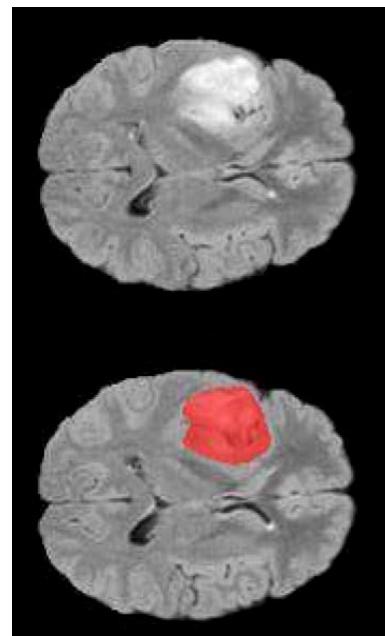
Object Detection

Semantic Segmentation

# Motivation for U-Net



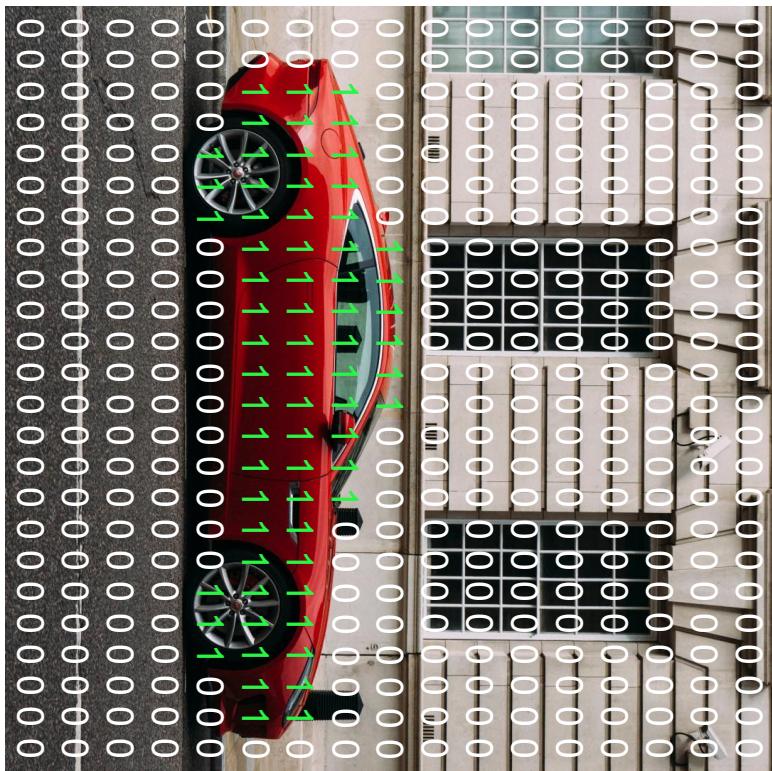
Chest X-Ray



Brain MRI

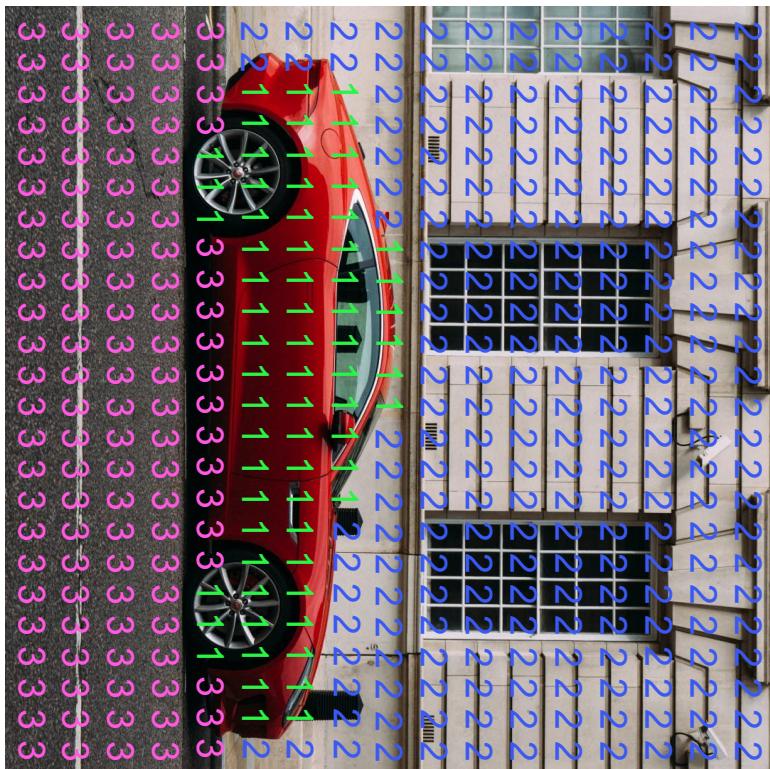
[Novikov et al., 2017, Fully Convolutional Architectures for Multi-Class Segmentation in Chest Radiographs]  
[Dong et al., 2017, Automatic Brain Tumor Detection and Segmentation Using U-Net Based Fully Convolutional Networks ]

# Per-pixel class labels



1. Car  
0. Not Car

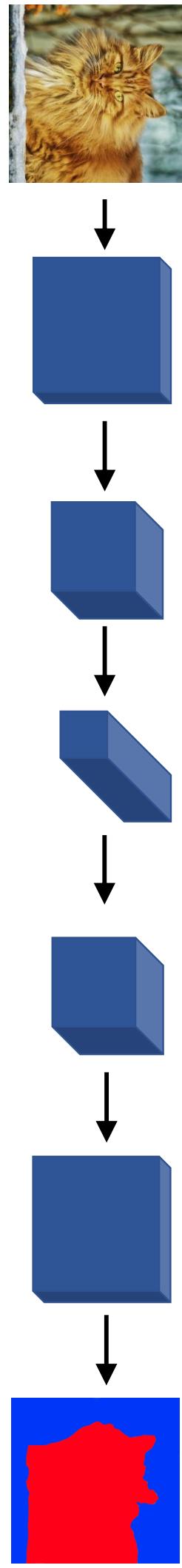
## Per-pixel class labels



1. Car
2. Building
3. Road

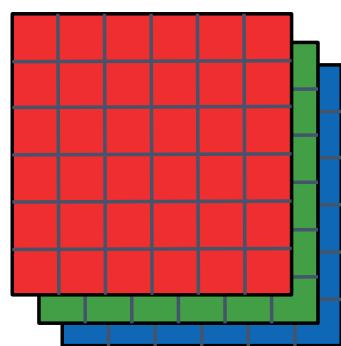
## Segmentation Map

# Deep Learning for Semantic Segmentation



# Transpose Convolution

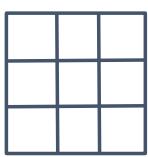
Normal Convolution



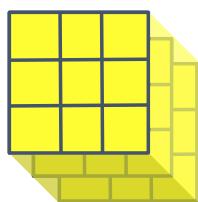
Transpose Convolution



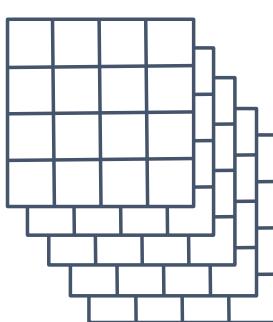
\*



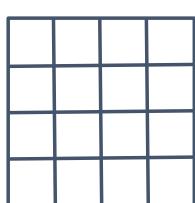
=



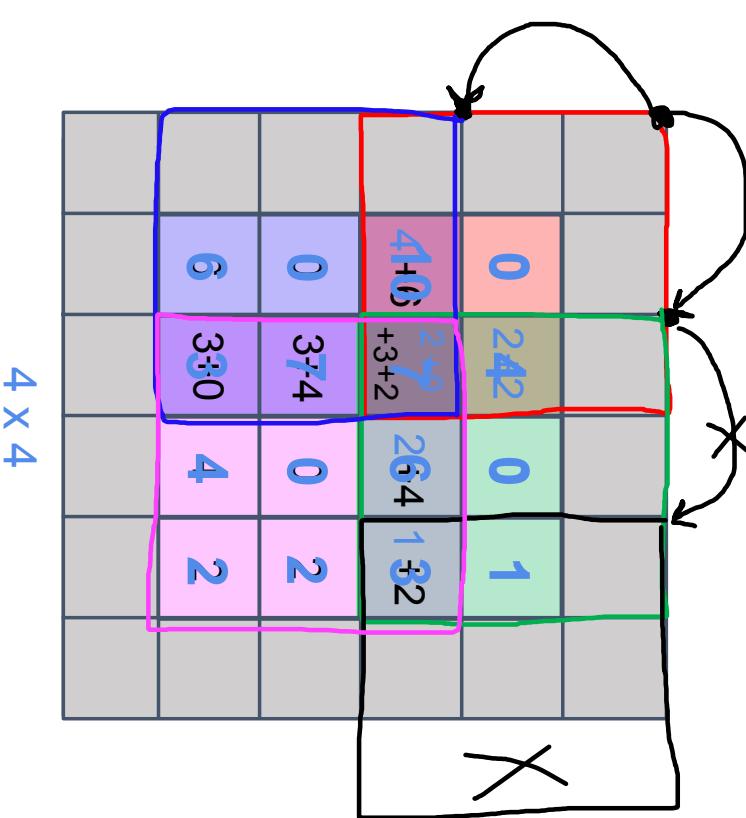
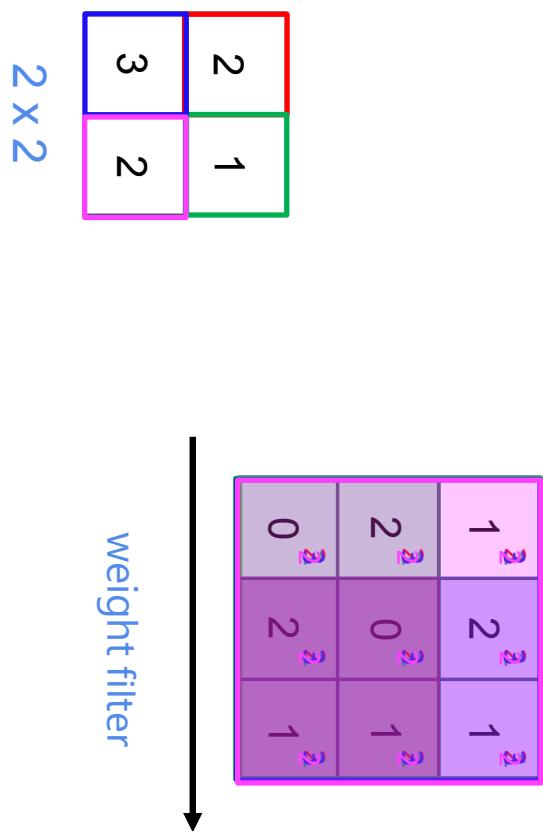
\*



=



# Transpose Convolution

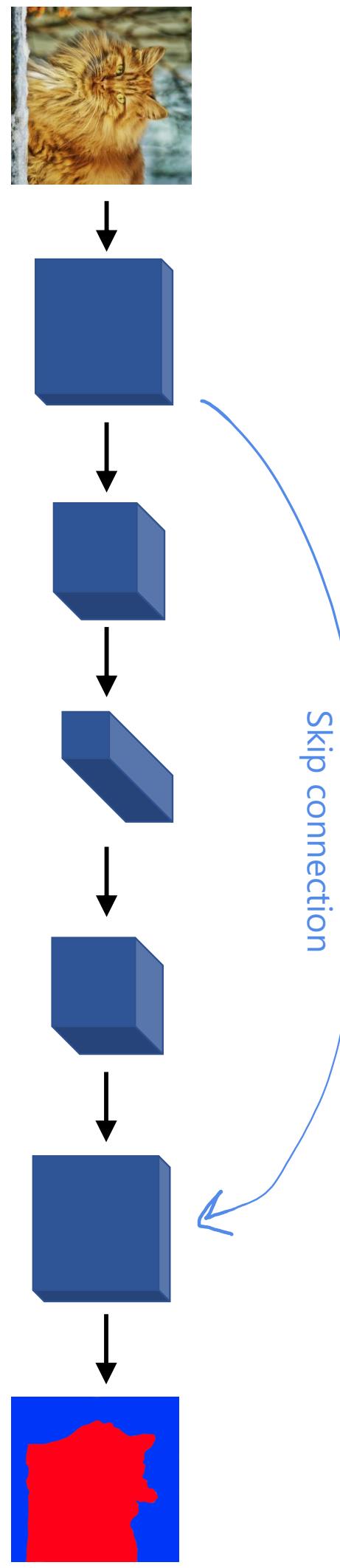


filter  $f \times f = 3 \times 3$

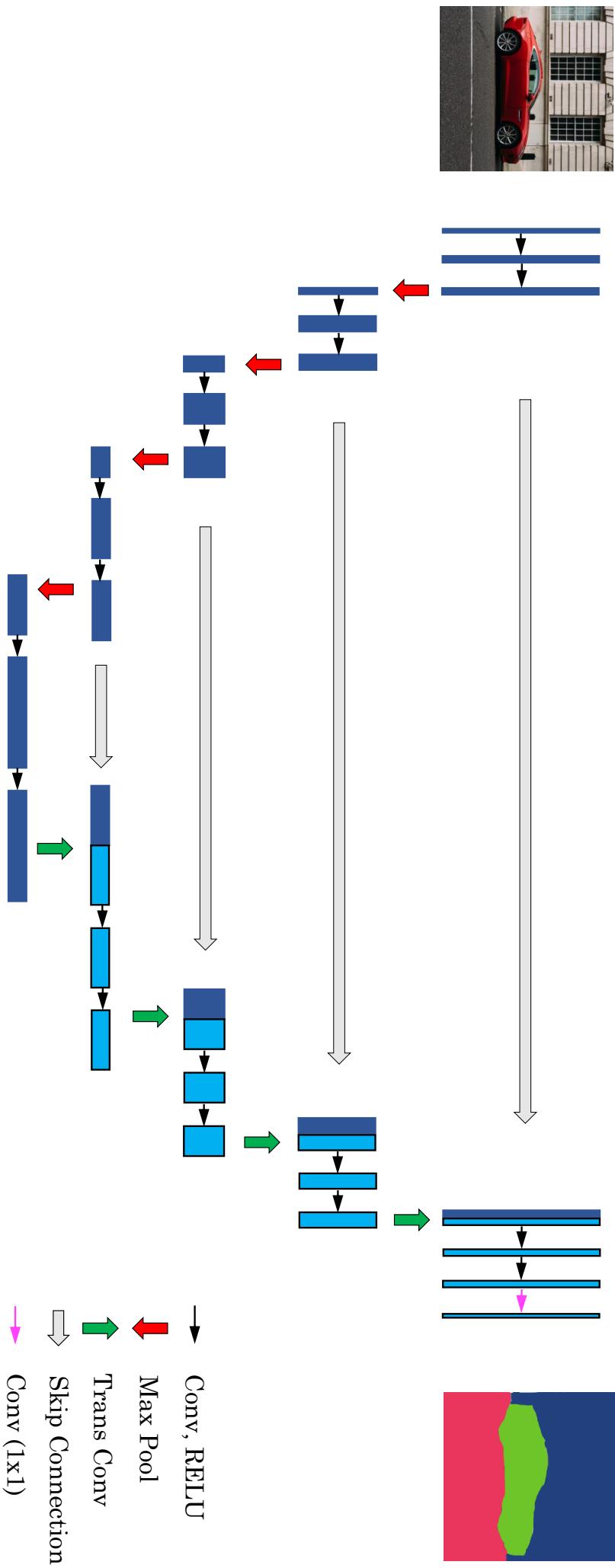
padding  $p = 1$

stride  $s = 2$

# Deep Learning for Semantic Segmentation



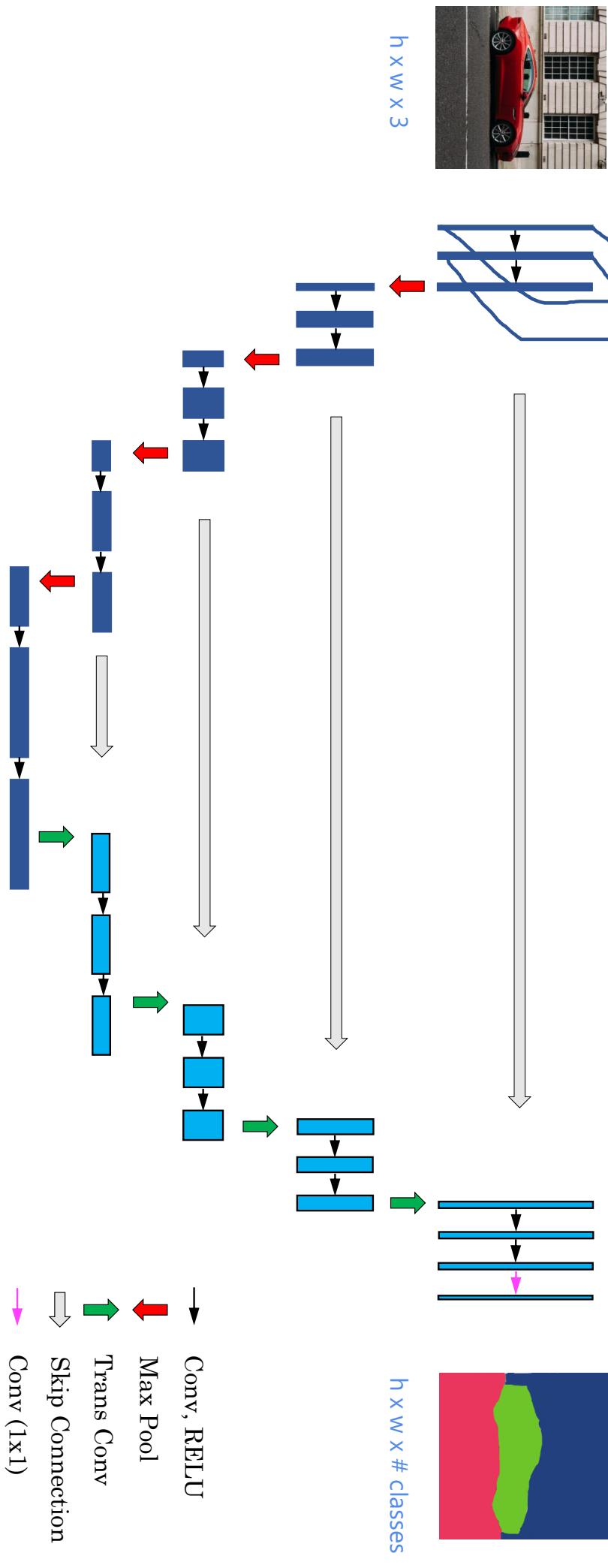
# U-Net



[Ronneberger et al., 2015, U-Net: Convolutional Networks for Biomedical Image Segmentation]

Andrew Ng

# U-Net



[Ronneberger et al., 2015, U-Net: Convolutional Networks for Biomedical Image Segmentation]

Andrew Ng