# 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
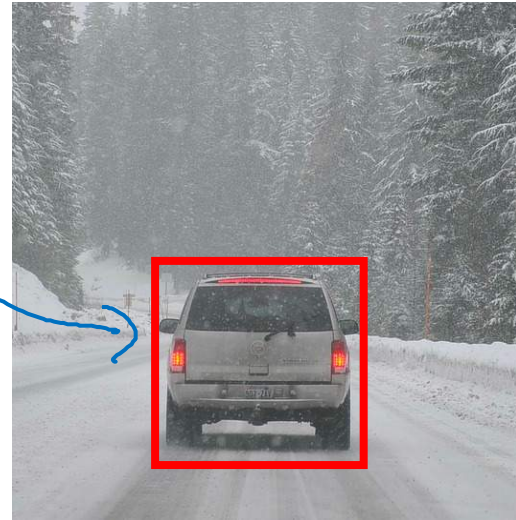_____

Object
localization

# What are localization and detection?

Image classification
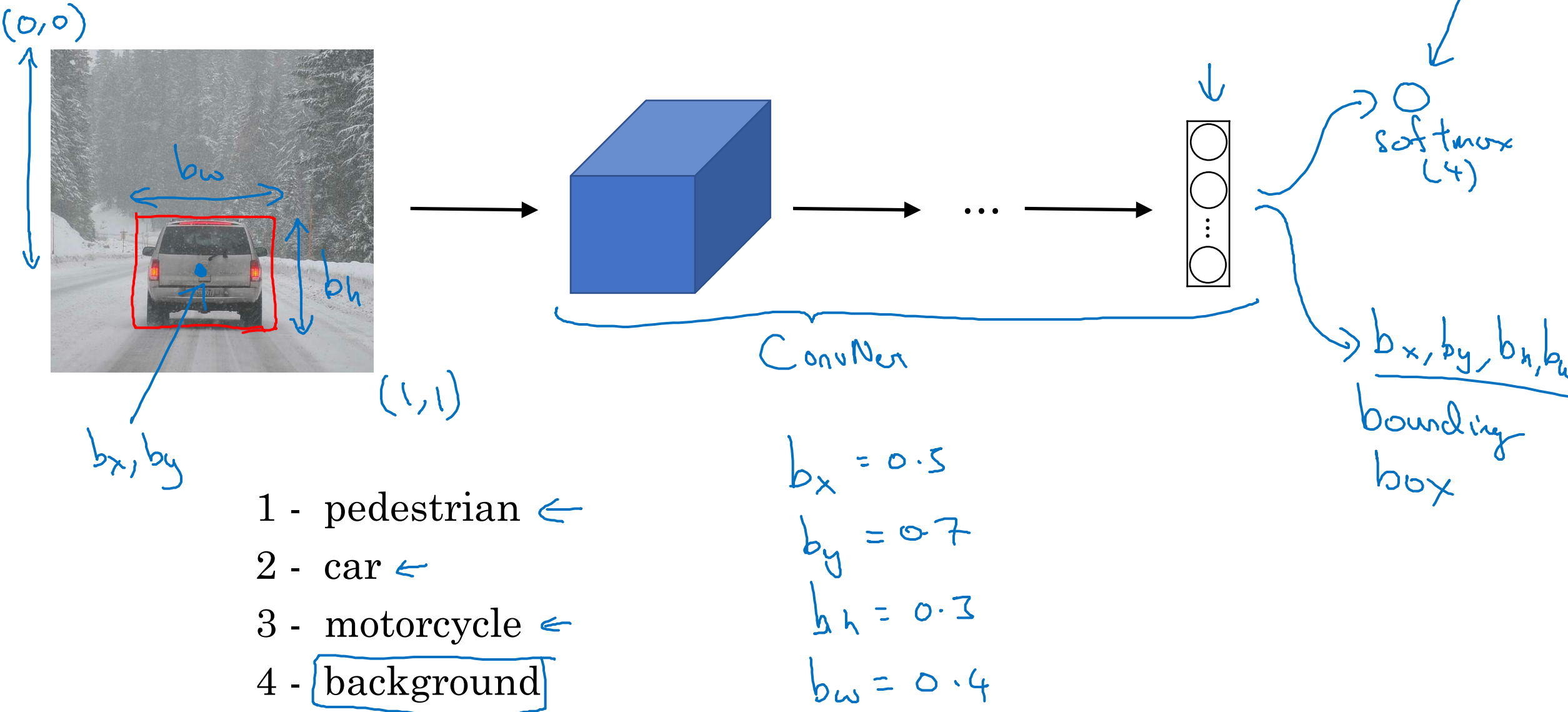


"Car"

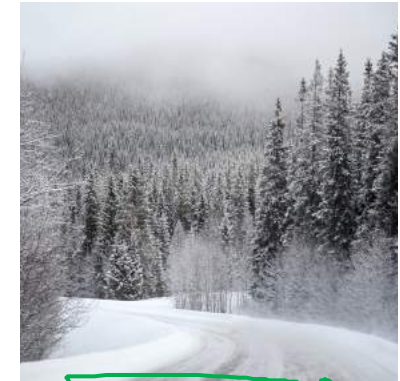Classification with localization
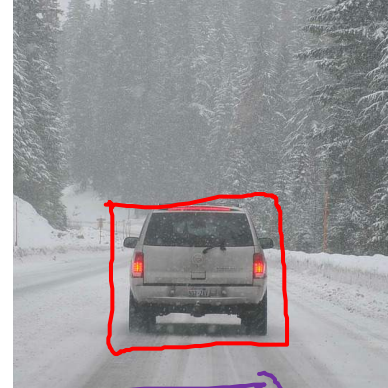


"Car"

Detection



1 object

multiple objects

Andrew Ng

# Classification with localization



(0,0)

$b_w$

$b_h$

$b_x, b_y$

(1,1)

1 - pedestrian

2 - car

3 - motorcycle

4 - background

ConvNet

softmax (4)

$b_x, b_y, b_h, b_w$

bounding box

$b_x = 0.5$

$b_y = 0.7$

$b_h = 0.3$

$b_w = 0.4$

Andrew Ng
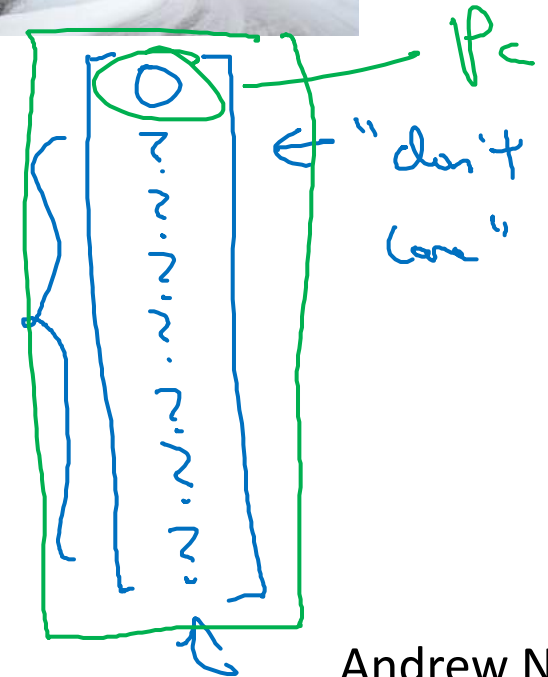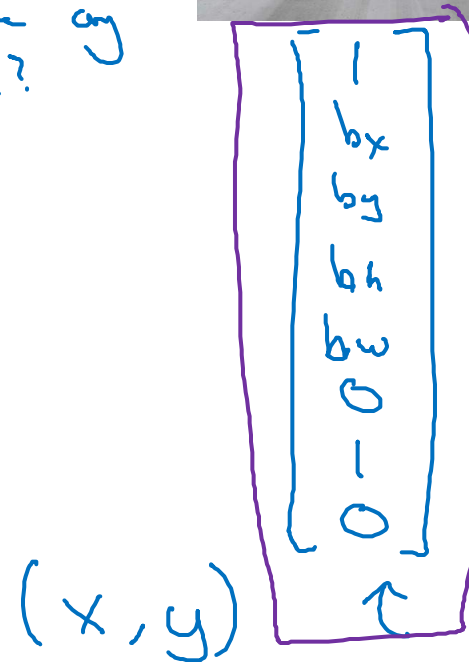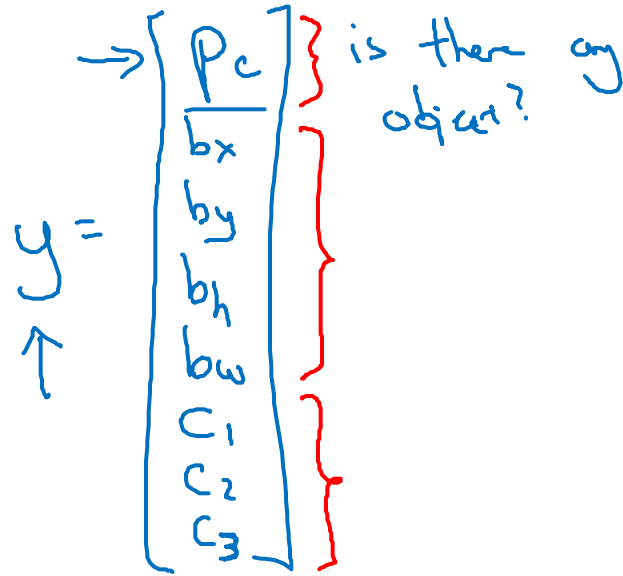
# Defining the target label y

1 - pedestrian

2 - car ←

3 - motorcycle

4 - background ←

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

$x =$



$\mathcal{L}(\hat{y}, y) =$

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

$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$  is there any object?

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

$(x, y)$

$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$  ← "don't care"

$P_c$

Andrew Ng

# Object Detection
---

# Landmark detection

deeplearning.ai

# Landmark detection

ConvNet

$$face?$$
$$l_{1x}$$
$$l_{1y}$$
$$l_{64x}$$
$$l_{64y}$$

$$129$$

AR



$$b_x, b_y, b_h, b_w$$

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

$$X, Y$$

$$l_{1x}, l_{1y},$$
$$\vdots$$
$$l_{32x}, l_{32y}$$

deeplearning.ai

# Object Detection

Object detection

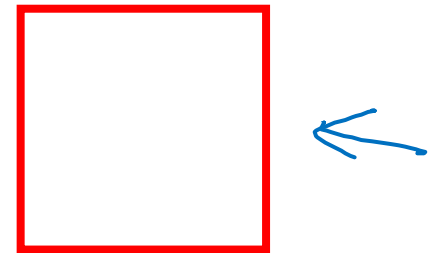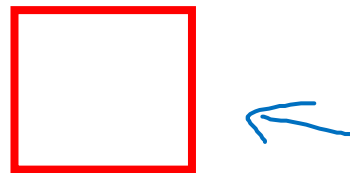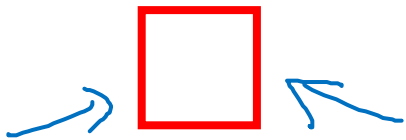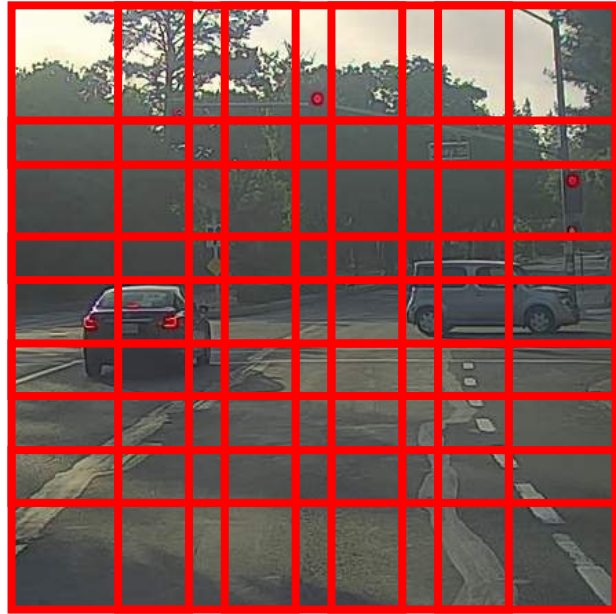# Car detection example

Training set:



x         y

1

1

1

0

0

ConvNet → y

Andrew Ng

# Sliding windows detection



$\rightarrow$ ConvNet $\rightarrow$ 0

$\rightarrow$ ConvNet

Computation cost

Andrew Ng

Object Detection

Convolutional implementation of sliding windows

deeplearning.ai

# Turning FC layer into convolutional layers



$14 \times 14 \times 3$     $5 \times 5$     $10 \times 10 \times 16$     MAX POOL     $2 \times 2$     $5 \times 5 \times 16$     FC     400     FC     400     y softmax (4)

$14 \times 14 \times 3$     $5 \times 5$     $10 \times 10 \times 16$     MAX POOL     $2 \times 2$     $5 \times 5 \times 16$     FC     $5 \times 5$     400     $1 \times 1 \times 400$     FC     $1 \times 1$     400     $1 \times 1 \times 400$     $1 \times 1$     $1 \times 1 \times 4$

$5 \times 5 \times 16$

# Convolution implementation of sliding windows



14×14 ×3 → 5×5 → 10×10×16 → MAX POOL 2×2 → 5×5×16 → FC 5×5 → 1×1×400 → FC 1×1 → 1×1×400 → FC 1×1 → 1×1×4

16×16×3 → 5×5 → 12×12×16 → MAX POOL 2×2 → 6×6×16 → FC 5×5 → 2×2×400 → FC 1×1 → 2×2×400 → FC 1×1 → 2×2×4

28×28×3 → 5×5 → 24×24×16 → MAX POOL 2×2 → 12×12×16 → 5×5 → 8×8×400 → 1×1 → 8×8×400 → 1×1 → 8×8×4

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

Andrew Ng

# Convolution implementation of sliding windows



28×28 → 5×5 → 16×16 → MAX POOL 2×2 → 12×12 → 5×5 → 8×8×400 → 1×1 → 8×8×400 → 1×1 → 8×8×4

deeplearning.ai

# Object Detection

## Bounding box predictions
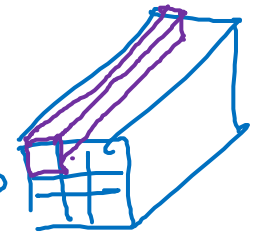
# Output accurate bounding boxes

# YOLO algorithm



100

100

Labels for training
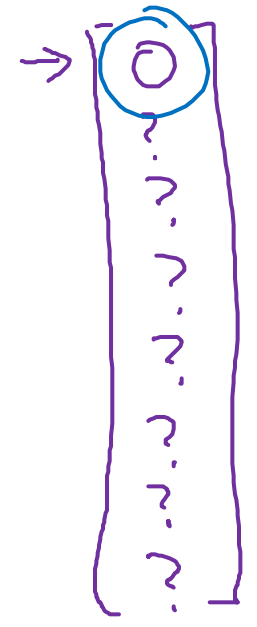
For each grid cell:

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

Target output:

$3 \times 3 \times 8$

$x \rightarrow \boxed{\phantom{x}}$  $100 \times 100 \times 3$  $\rightarrow$  CONV $\rightarrow$ MAXPOOL $\rightarrow$ ....  $\rightarrow$  $3 \times 3 \times 8$  $y$

$\rightarrow 19 \times 19 \times 8$

$3 \cdot 8$

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

# Specify the bounding boxes



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

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

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

Object Detection

Intersection over union

deeplearning.ai

# Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{Size of } \boxed{\phantom{x}}}{\text{Size of } \boxed{\phantom{x}}}$$

"Correct" if IoU ≥ 0.5

0.6

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



19×19

# Non-max suppression example



$P_c$

# Non-max suppression algorithm



19×19

Each output prediction is:

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

$$\begin{matrix} c_1 \\ c_2 \\ c_3 \end{matrix}$$

Discard all boxes with $p_c \le 0.6$

While there are any remaining boxes:

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

- Discard any remaining box with
  IoU $\ge 0.5$ with the box output
  in the previous step

Andrew Ng

Object Detection

Anchor boxes

deeplearning.ai

# Overlapping objects:



Anchor box 1:

Anchor box 2:

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

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \begin{matrix} \\ \\ \text{Anchor box 1} \\ \\ \end{matrix}$$

Anchor box 1

$$\begin{bmatrix} p_c \\ b_x \\ \vdots \\ c_3 \end{bmatrix}$$ Anchor box 2

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

Andrew Ng

# 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.
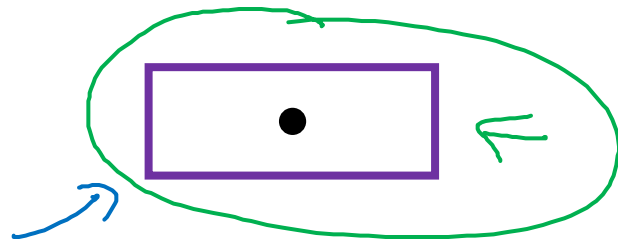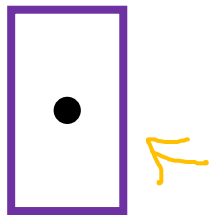
(grid cell, anchor box)

Output y:

$3 \times 3 \times 16$

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

Andrew Ng

# Anchor box example



Anchor box 1:     Anchor box 2:

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

1
$b_x$
$b_y$
$b_h$
$b_w$
1
0
0
1
$b_x$
$b_y$
$b_h$
$b_w$
0
1
0

$c_1$ only?

0
?
?
?
?
?
?
?
1
$b_x$
$b_y$
$b_h$
$b_w$
0
1
0

anchor box 1

anchor box 2

# Object Detection

## Putting it together: YOLO algorithm

deeplearning.ai

# Training

1 - pedestrian

2 - car

3 - motorcycle

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

1

2

$3 \times 3 \times 16$

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

$19 \times 19 \times 16$

$19 \times 19 \times 40$

#anchors

$5 +$ #classes

$100 \times 100 \times 3$

ConvNet

$3 \times 3 \times 16$

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

Andrew Ng

# Making predictions



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

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

Andrew Ng

# Outputting the non-max supressed outputs



- For each grid call, get 2 predicted bounding boxes.

- Get rid of low probability predictions.

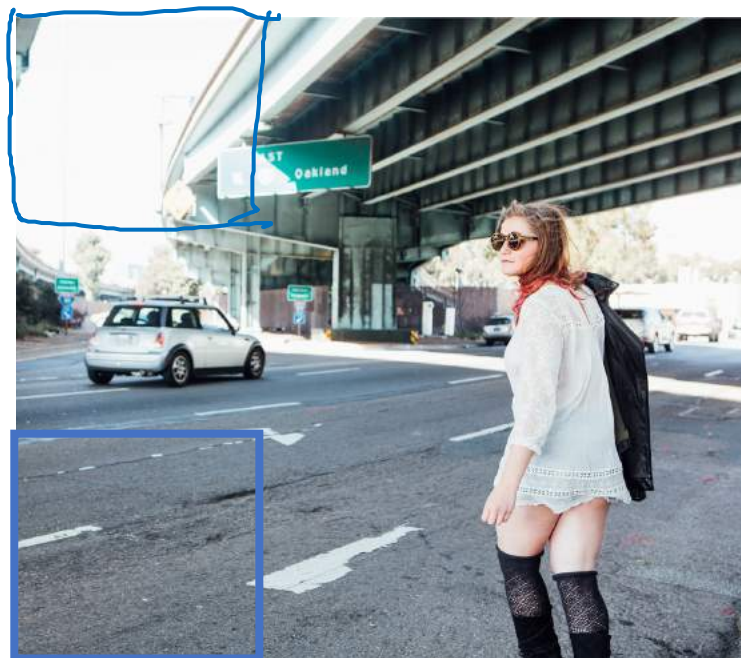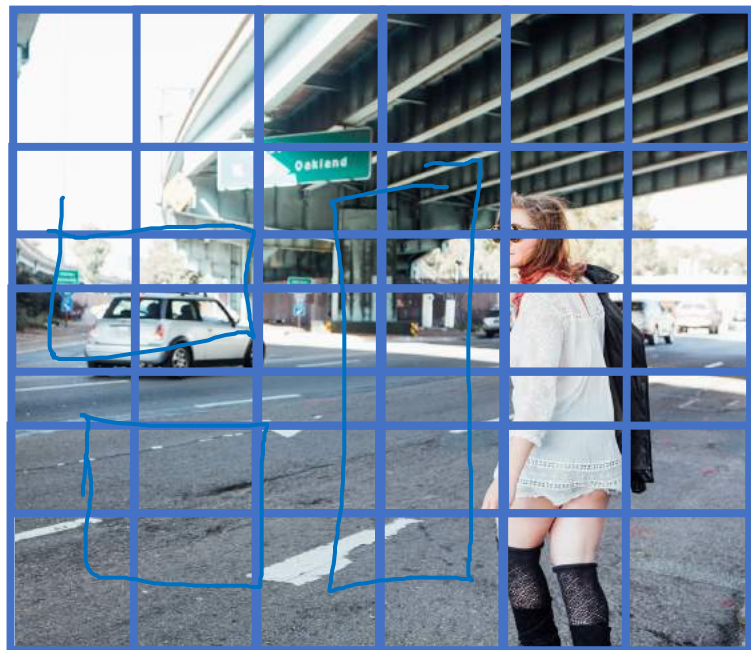- For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions.
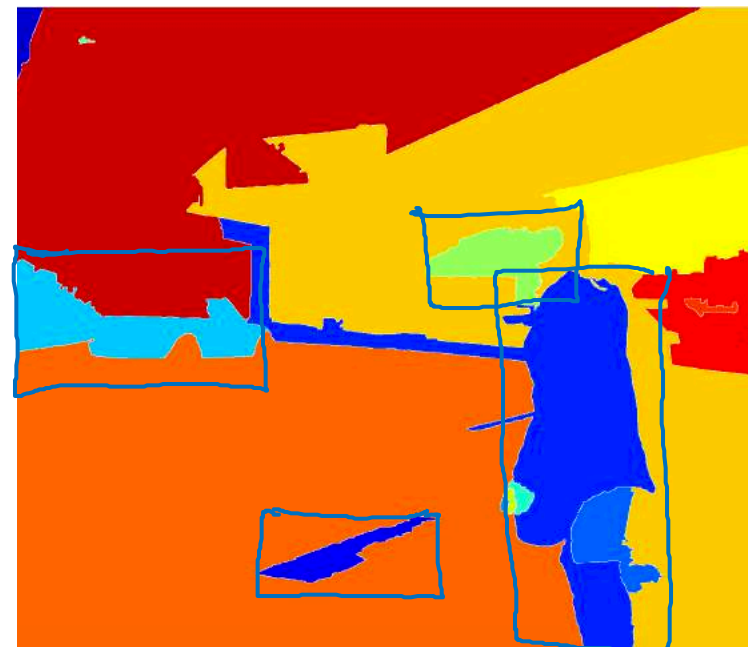
# Object Detection

Region proposals
(Optional)

deeplearning.ai

# Region proposal: R-CNN



Segmentation algorithm

~2,000

[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

# Object Detection vs. Semantic Segmentation



Input image

Object Detection

Semantic Segmentation

Andrew Ng

# 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 ]
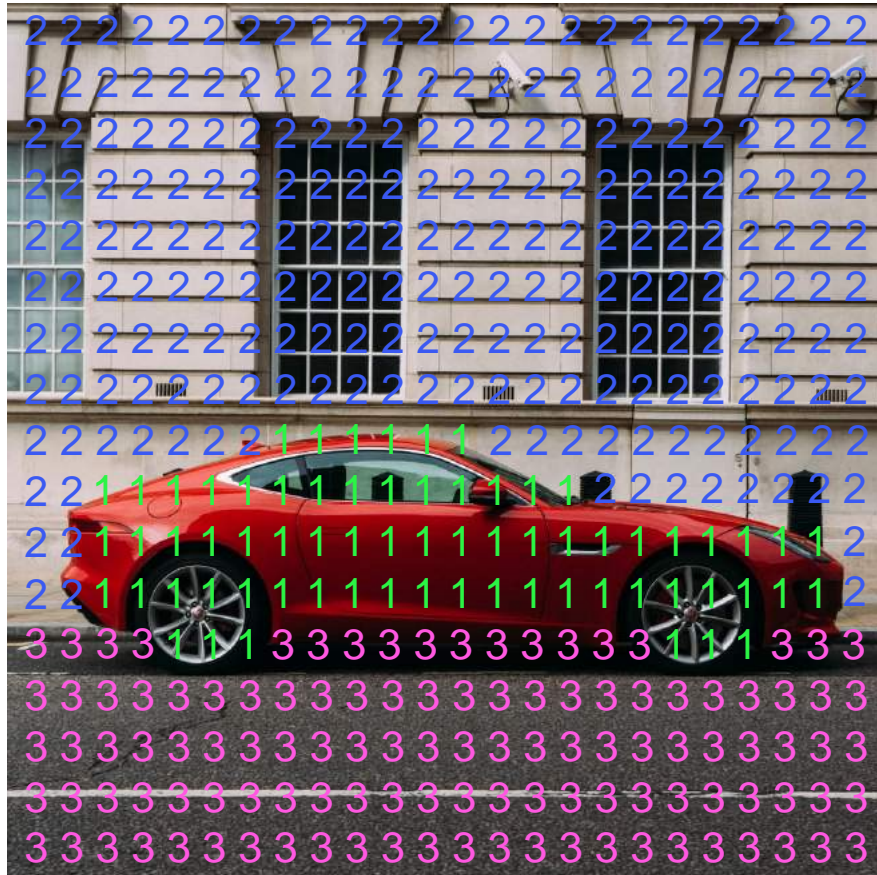
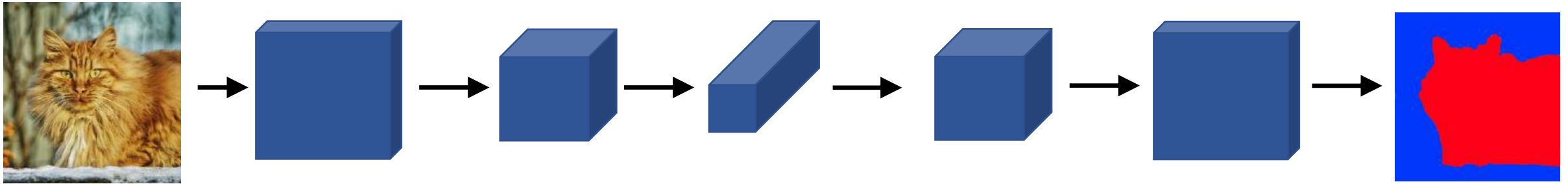Andrew Ng

# Per-pixel class labels



1. Car
0. Not Car

Andrew Ng

# Per-pixel class labels



1. Car
2. Building
3. Road

Segmentation Map

Andrew Ng

# Deep Learning for Semantic Segmentation



Andrew Ng

# Transpose Convolution

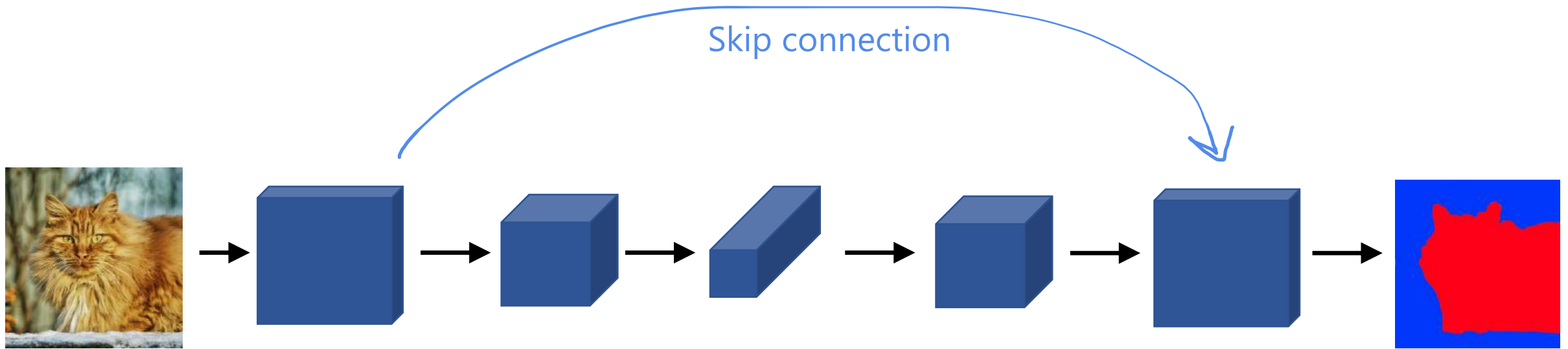Normal Convolution



Transpose Convolution

# Transpose Convolution



2 x 2

weight filter
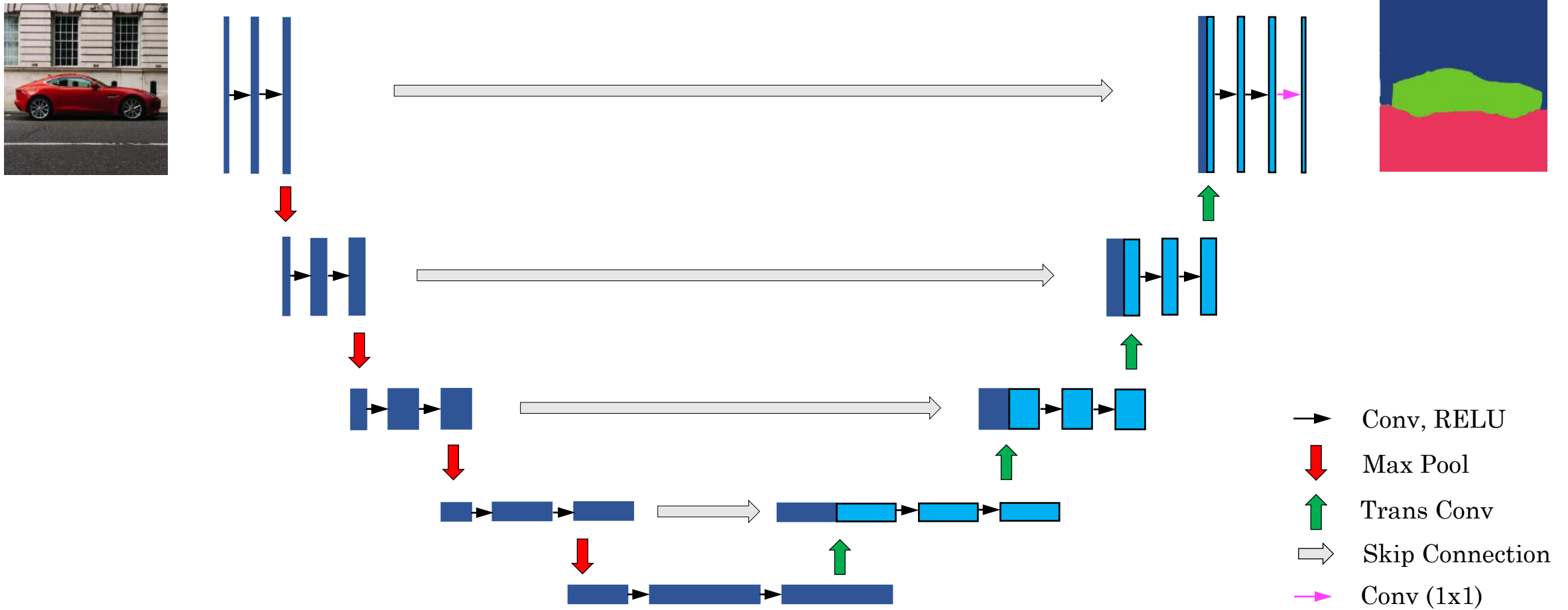
4 x 4

filter f x f = 3 x 3          padding p = 1          stride s =  2

Andrew Ng

# Deep Learning for Semantic Segmentation



Skip connection

Andrew Ng
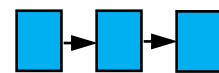
# U-Net



Conv, RELU

Max Pool

Trans Conv

Skip Connection

Conv (1x1)

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

Andrew Ng

# U-Net



h x w x 3

h x w x # classes

→ Conv, RELU

↓ Max Pool

↑ Trans Conv

⇒ Skip Connection

→ Conv (1x1)

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

Andrew Ng