# What are localization and detection?

Image classification

Classification with localization

Detection

"Car"

"Car"

1:25 / 11:53

C4W3L01 Object Localization

**What is the difference between these problems?**

**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)

$x =$

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

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

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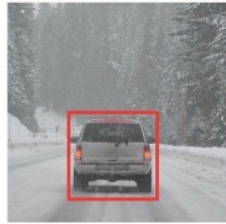
Andrew Ng

**What would the loss function be if you used cross-entropy loss for classification outputs and MSE for regression outputs?**

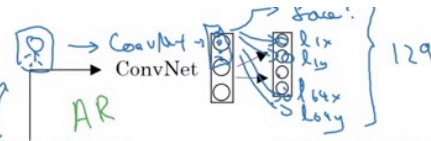**What is the theoretical underpinning of this loss?**

Landmark detection

$b_x, b_y, b_h, b_w$

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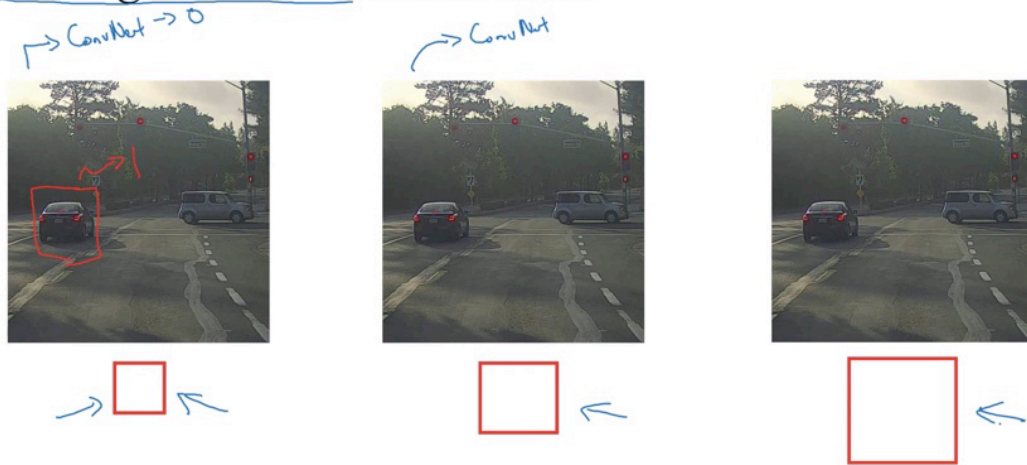$l_{64}, l_{64y}$
$\Big\} \ X, Y$

$l_{1x}, l_{1y},$
$\vdots$
$l_{32x} \cdot l_{22y}$

Andrew Ng

**What are some concerns about doing landmark detection as described?**

# Sliding windows detection

→ ConvNet → 0

→ ConvNet
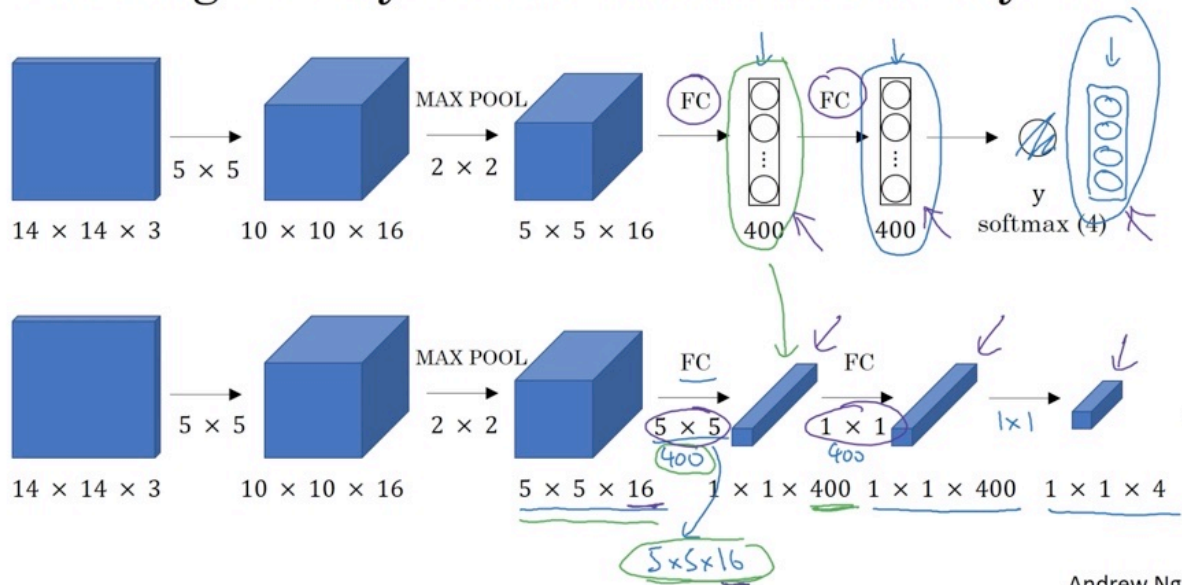
→ ConvNet

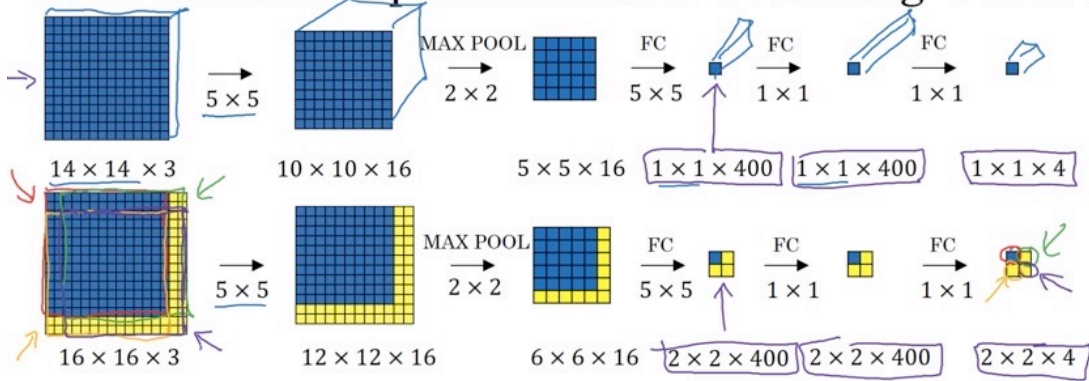**How do you pass larger image fragments into the same classifier?**

# Turning FC layer into convolutional layers



By changing these fully connected layers into convolutional layers, what does that now allow us to do?

# 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

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

Andrew Ng

# Overlapping objects:
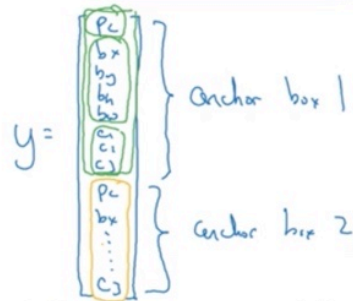


Anchor box 1:

Anchor box 2:

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$
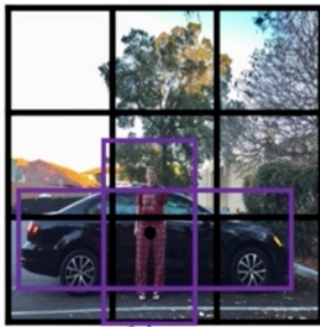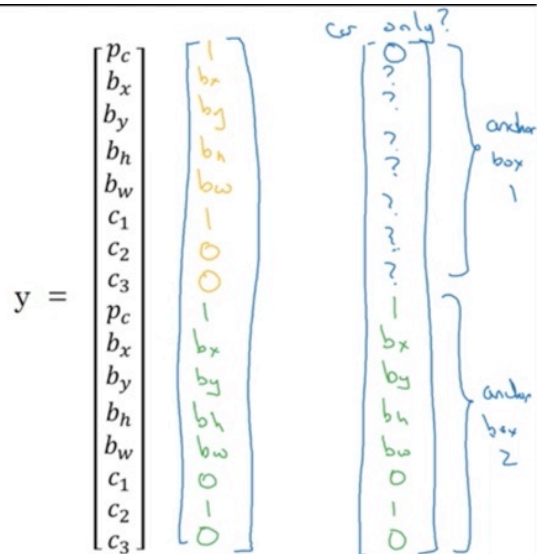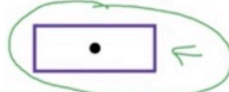
$y =$ (Anchor box 1 / Anchor box 2)

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

Andrew Ng

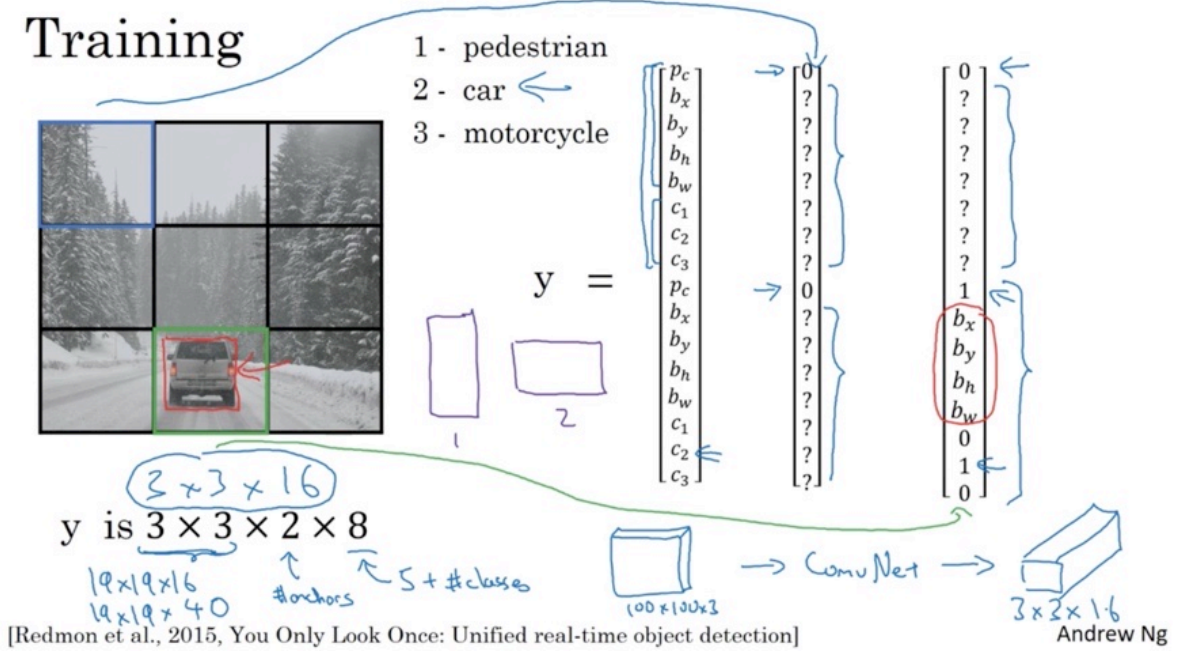# Anchor box example



Anchor box 1:    Anchor box 2:

Car only?

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

anchor box 1

anchor box 2

Andrew Ng

## Training

1 - pedestrian
2 - car
3 - motorcycle

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

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

… →

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

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

$b_x$
$b_y$
$b_h$
$b_w$

Andrew Ng