

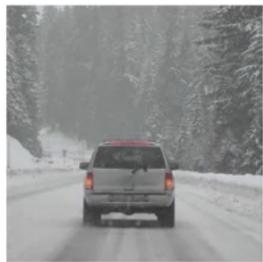
## Detection Algorithms

Alec Dewulf

### Object Localization

## What are localization and detection?

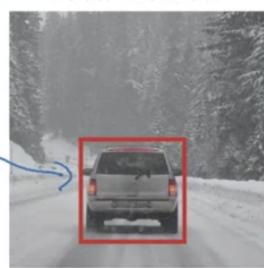
Image classification



"Car"

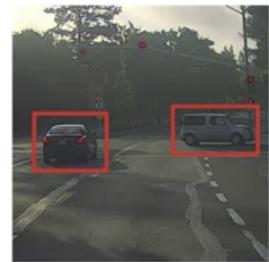
1 object

Classification with localization



"Car"

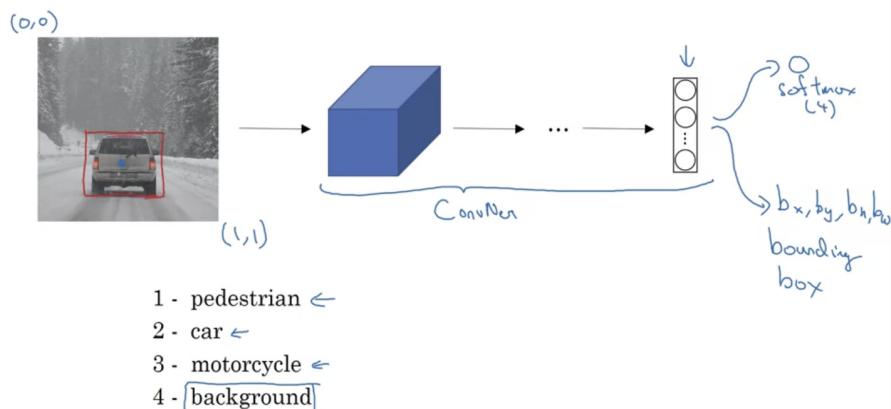
Detection



multiple objects

You can have a network output four more numbers to parameterize a bounding box to surround the detected object.

### Classification with localization



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)



$x =$

is there any object?

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

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

Annotations

$$\ell(\hat{y}, y) =$$

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

Loss function has two cases. If there's no object we only care whether or not there was an object. Not what the predictions were.

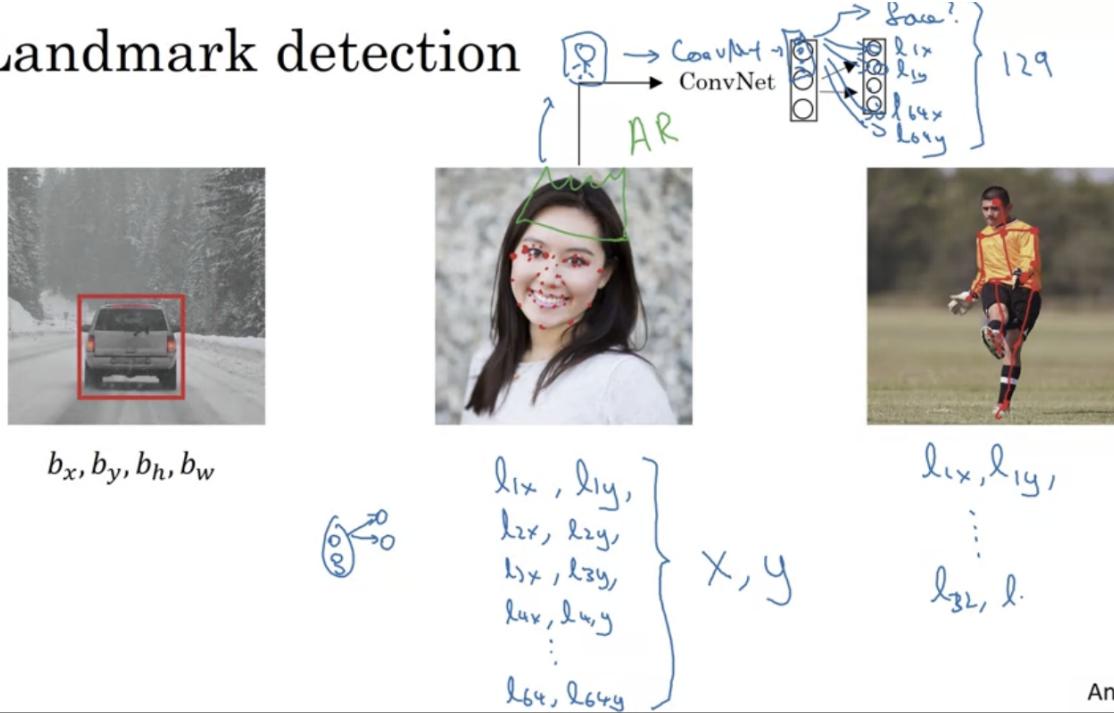
## Landmark Detection

Landmarks are important points in an image that you might want the algorithm to recognize. This can be represented by two numbers: the x and y coordinates.

You can output a vector of landmarks.

You can also define key positions if you are interested in pose estimation. By outputting landmarks in key positions you can output the pose of a person.

## Landmark detection



The landmark labels have to be consistent across images.

## Object Detection

Train a ConvNet with closely cropped images to use with sliding window detection. The area outlined by the sliding window will be input to this model.

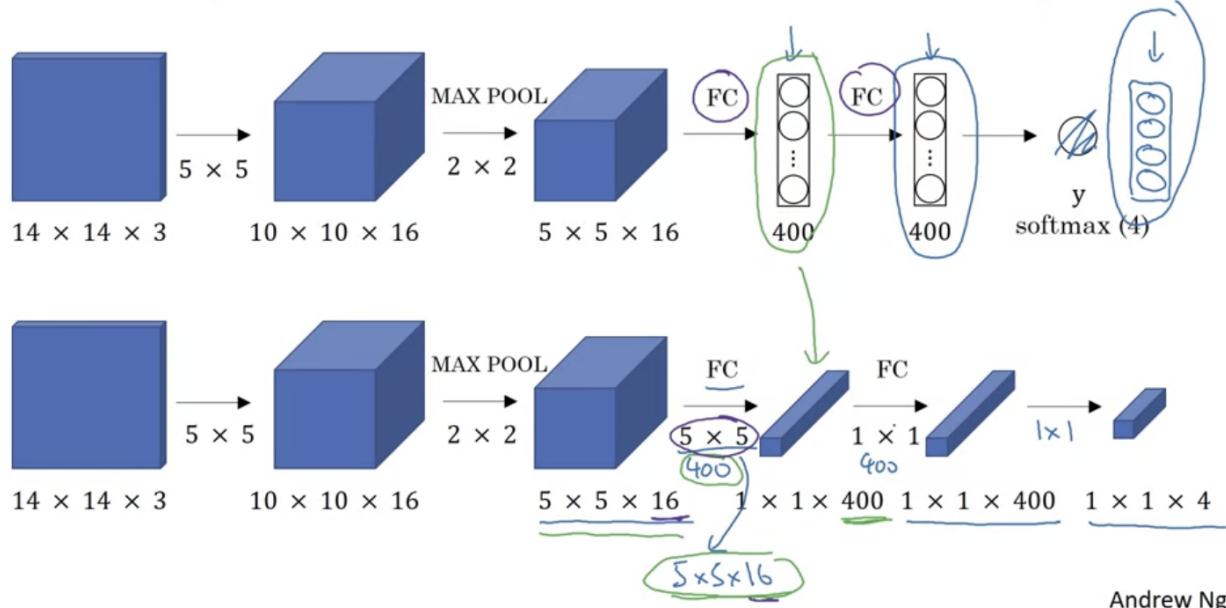
Use the sliding window multiple times with larger windows each time. Using sliding windows has a huge computational cost that is infeasible.

The sliding window's object detector can be implemented convolutionally which is much more efficient.

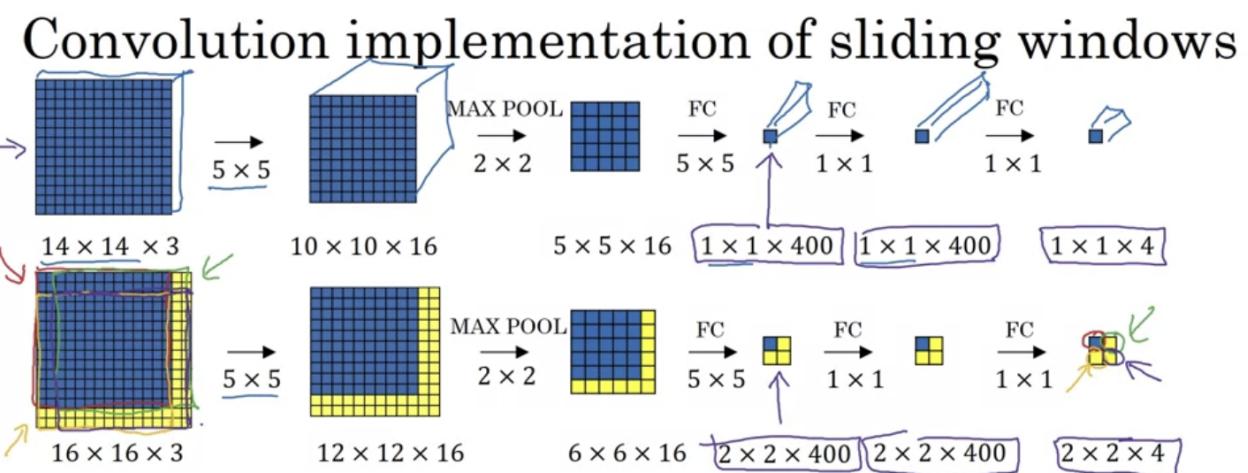
## Convolutional Implementation of Sliding Windows

We can replace a FC with a Conv layer by using a large number of filters to build a  $1 \times 1 \times n_C$  volume.

## Turning FC layer into convolutional layers



A convolutional implementation allows sliding window steps to share computation.



The output volume has channels that give the same information as implementing the full sliding window process.

The max pooling dimension dictates the corresponding stride of the equivalent sliding window.

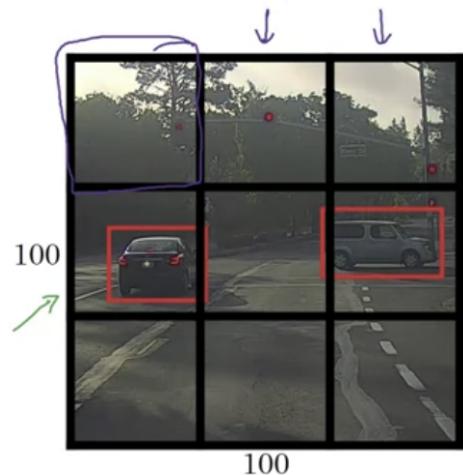
## Bounding Box Predictions

The convolutional implementation of sliding windows doesn't output very accurate bounding boxes.

The YOLO (You Only Look Once) algorithm outputs more accurate bounding boxes. The idea is to divide the image into boxes and apply the image localization to each grid cell.

You have an output for each grid cell as shown in the image below.

## YOLO algorithm



Labels for training  
For each grid cell:

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

The object is assigned to the cell that contains its midpoint.

The outputs are concatenated into an output volume where each channel is an output vector for a specific grid cell.

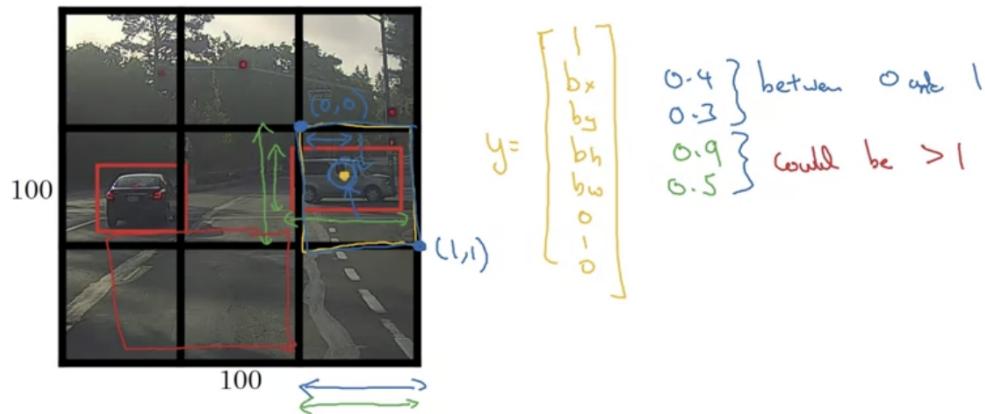
This algorithm works well as long as there is only one object in each grid cell.

The coordinate system allows the algorithm to output a bounding box of non-square dimensions.

This algorithm runs very fast because it's a convolutional implementation.

The convention in the YOLO algorithm is to specify the upper left corner of the grid cell as (0,0). The midpoint coordinates are specified relative to this coordinate and the height and width are specified relative to the overall dimensions of the grid cell.

## Specify the bounding boxes



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

Andrew Ng

## Intersection Over Union

The IoU function computes the intersection over union of the best bounding box and the bounding box outputted by the model. This is the ratio of the size of the union to the size of the remaining area in the output box. This ratio would be one if it was perfect.

## Evaluating object localization

Diagram showing two overlapping bounding boxes on a snowy road scene. One is blue and yellow, the other is red and yellow. Handwritten notes define Intersection over Union (IoU) as the size of the intersection divided by the size of the union, and state that it is "Correct" if  $\text{IoU} \geq 0.5$ .

$$\text{IoU} = \frac{\text{size of intersection}}{\text{size of union}}$$

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

0.6

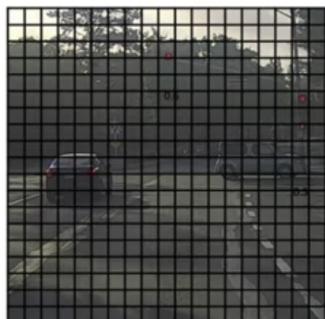
Andrew Ng

## Non-max Suppression

This is a way to make sure your algorithm detects each object only once.

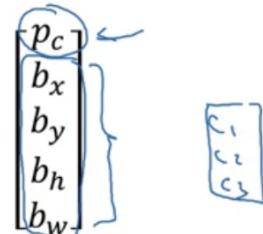
It first looks at the probabilities associated with each of the detections. It takes the largest one and then looks at the bounding boxes with high IoUs. These will be suppressed because they mostly overlap.

## Non-max suppression algorithm



19 × 19

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

Andrew Ng

If you have more than one object, you should carry out non-max suppression on each of them.

## Anchor Boxes

### Anchor box algorithm

Previously:

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

Output  $y$ :

$3 \times 2 \times 8$

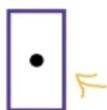
With two anchor boxes:

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.

## Anchor box example



Anchor box 1: Anchor box 2:



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

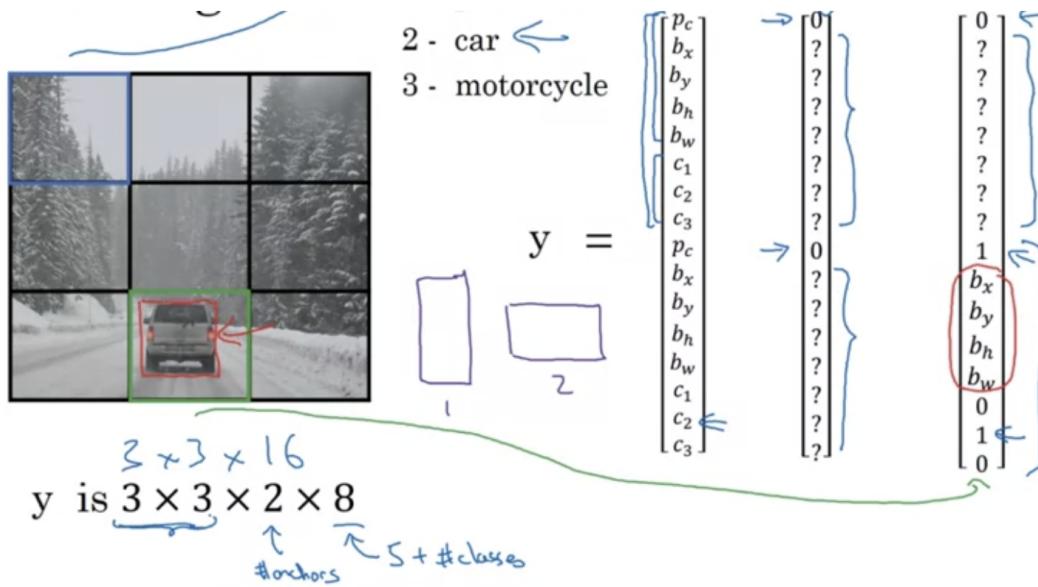
Andrew Ng

This algorithm doesn't handle the case where two objects are in the same grid cell and have a similar anchor box.

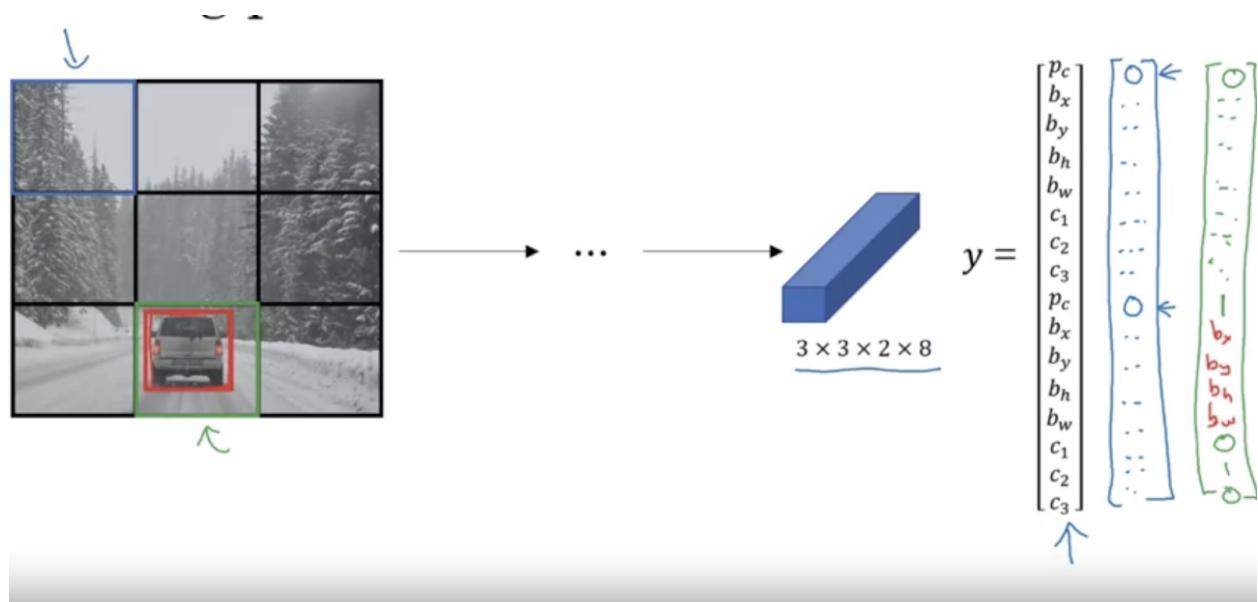
This is a method to deal with the case when two objects appear in the same grid cell (which doesn't happen often).

It also allows your algorithm to specialize. For example, in detecting tall skinny things and wide things. People usually choose the anchor boxes by hand.

## YOLO Algorithm



For each grid cell you attempt to find an object and see which bounding box fits it best. You then output parameters for each bounding box (which is why y is double the size you may expect).



# 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.

## Region Proposals

The sliding window algorithm tends to look at a lot of grid cells that have nothing interesting in them.

The R-CNN algorithm tries to pick a few regions where it makes more sense to run the sliding window classifier.

They run a segmentation algorithm to find “blobs” that might be interesting. Bounding boxes are then placed around the interesting blobs and the ConvNet is run on these boxes.

## 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. ↩

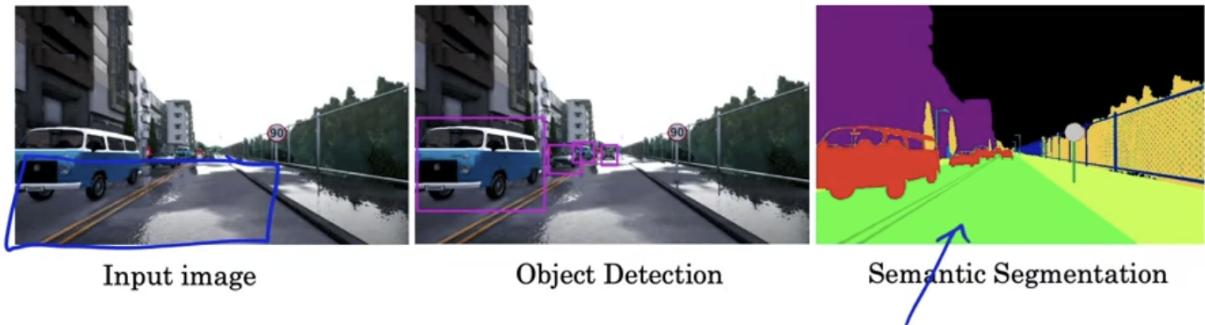
Faster R-CNN: Use convolutional network to propose regions.

Andrew Ng thinks region proposal is a cool idea but the two-step process doesn't seem promising.

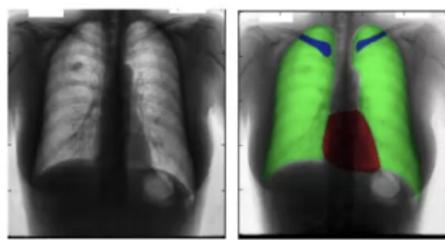
### Semantic Segmentation with U-Net

The goal of semantic segmentation is to draw a careful outline on the detected object.

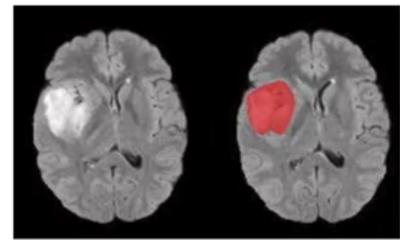
## Object Detection vs. Semantic Segmentation



## Motivation for U-Net



Chest X-Ray

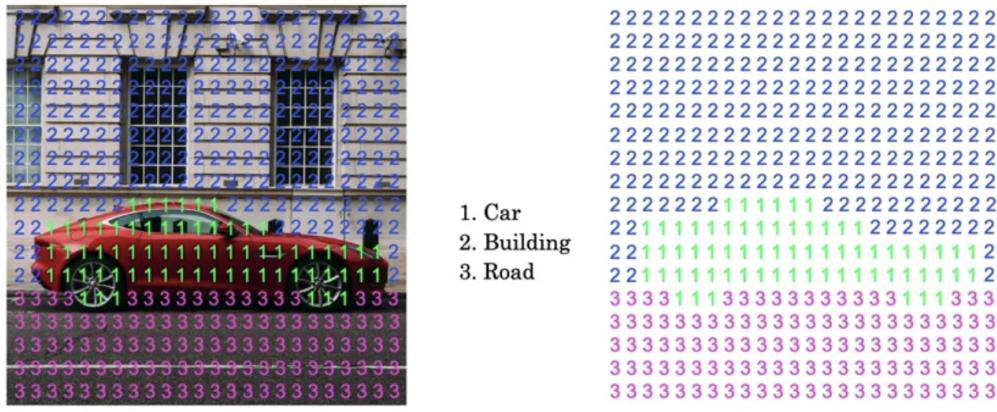


Brain MRI

The job of the U-Net algorithm is to output a value for every pixel depending on its class. You may have two or more classes.

## Per-pixel class labels

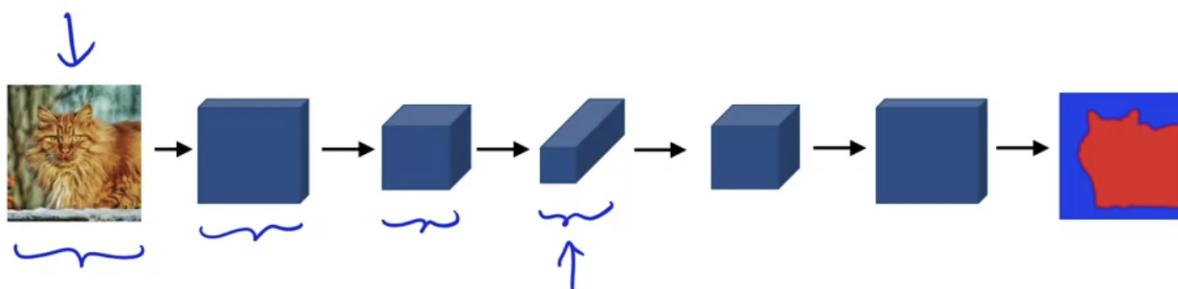
↓



## Segmentation Map

Andrew Ng

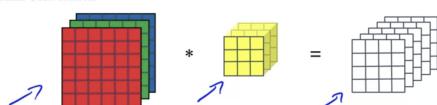
# Deep Learning for Semantic Segmentation



## Transpose Convolutions

## Transpose Convolution

### Normal Convolution



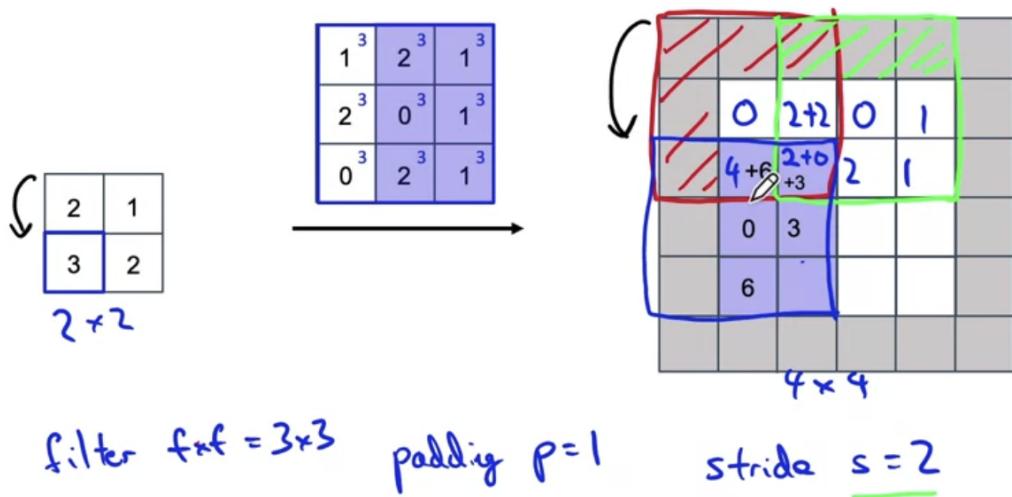
### Transpose Convolution



A transpose convolution increases the size of the matrix in contrast to the normal convolution.

You place the filter on the output rather than the input. You multiply the input by everything in the filter to get a matrix which you paste in the output matrix. Padding removes values from the output matrix (they aren't added outside the output).

## Transpose Convolution



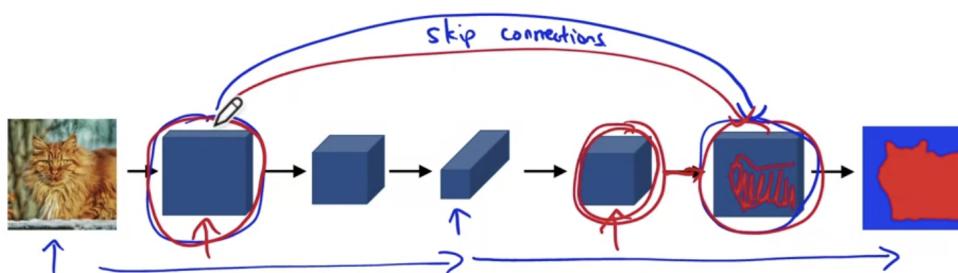
Andrew Ng

You add the elements in the overlapping sections.

## U-Net Architecture Intuition

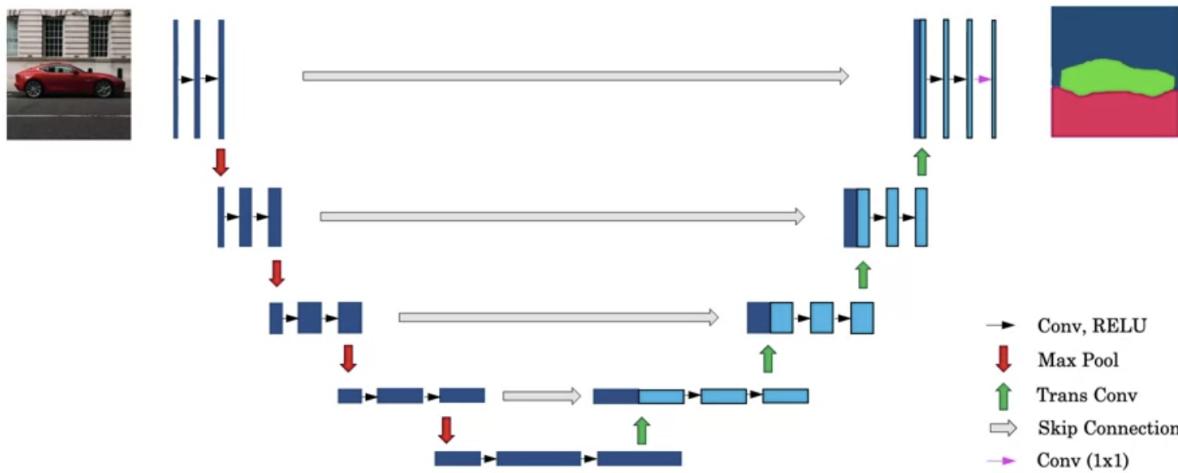
The skip connection in U-Net allows the model to take the high resolution low level information and pass it to the later layer. The last layer then has the high level and low-level features.

## Deep Learning for Semantic Segmentation



## U-Net Architecture

### U-Net



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

Andrew Ng

The skip connection always runs by appending the set of activations to the connected layer deeper in the network.