

# Intro to Object Detection

John Sigman

# Image Classification

Question: "What is this an image of?"

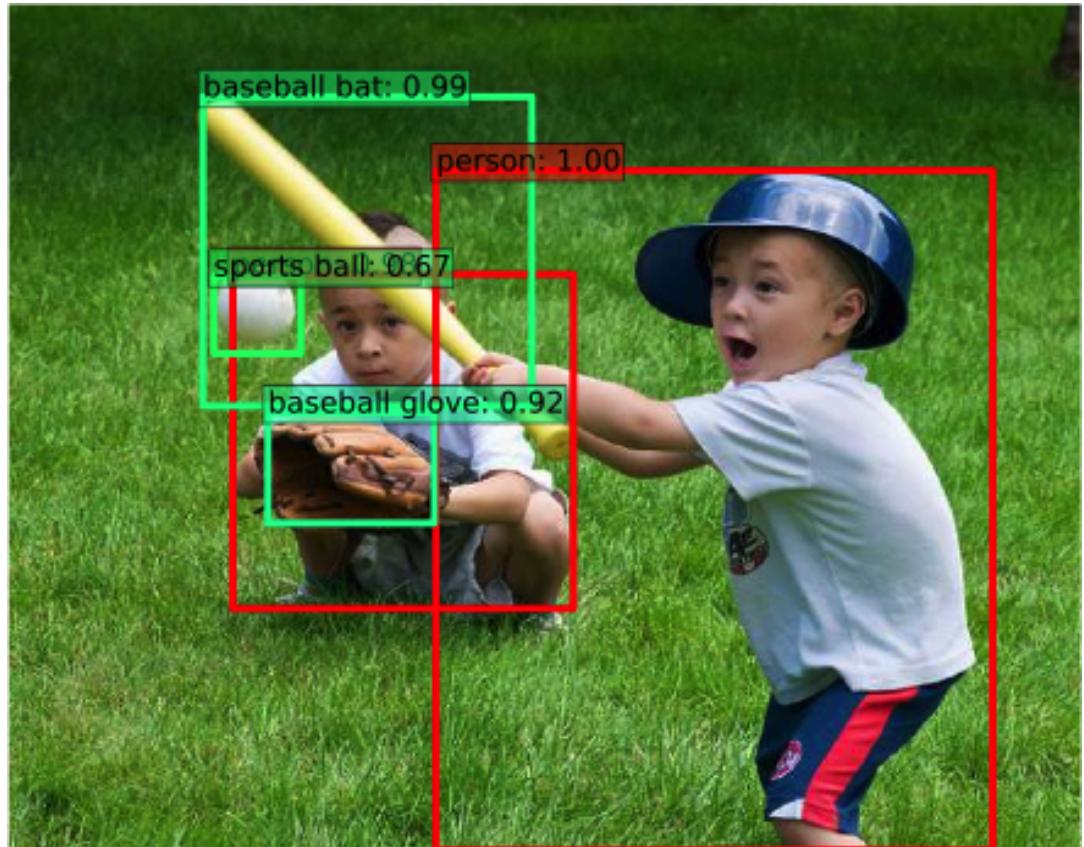
Answer: "95% probability this is a ballplayer"



# Object Detection

Question: “What are all the objects in this object and where are they?”

Answer:



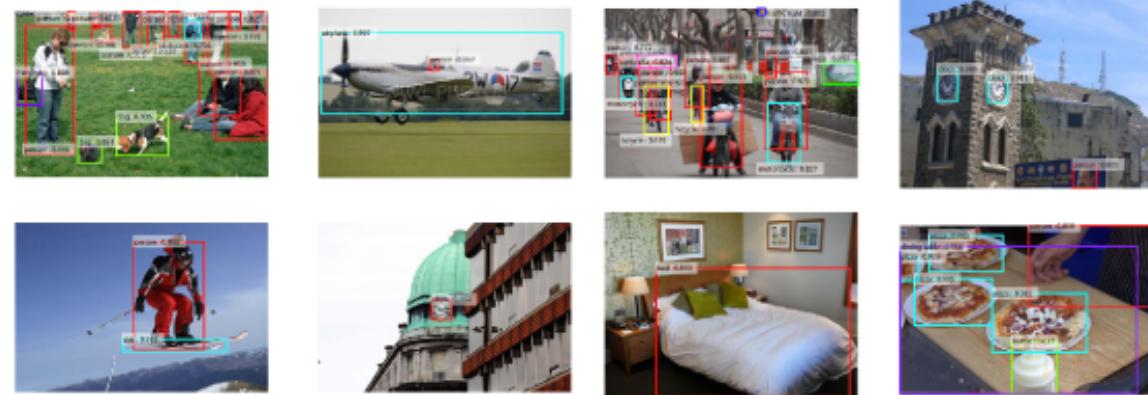
# Classification versus detection

Very related, but different problems.

Many aspects will be shared between the two problems



mite	container ship	motor scooter	leopard
mite	container ship	motor scooter	leopard
black widow	lifeboat	go-kart	jaguar
cockroach	amphibian	moped	cheetah
tick	fireboat	bumper car	snow leopard
starfish	drilling platform	golfcart	Egyptian cat

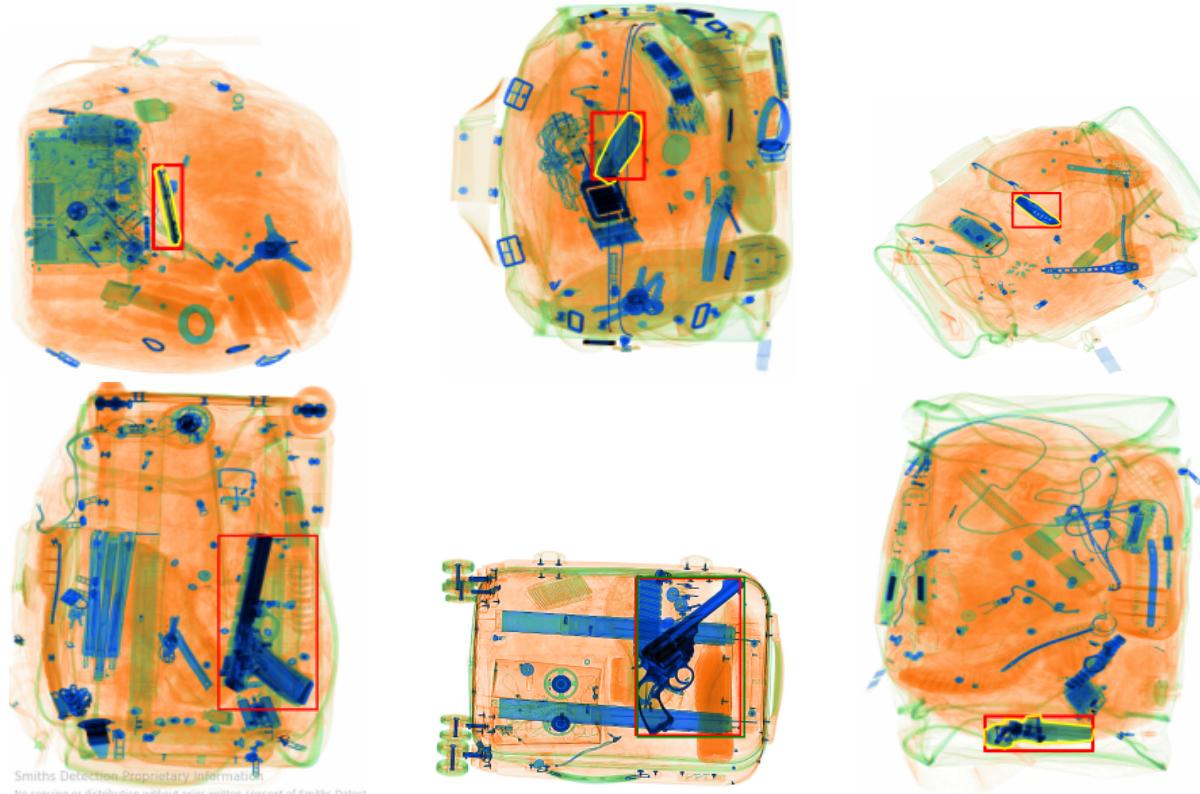


# Self-Driving Cars



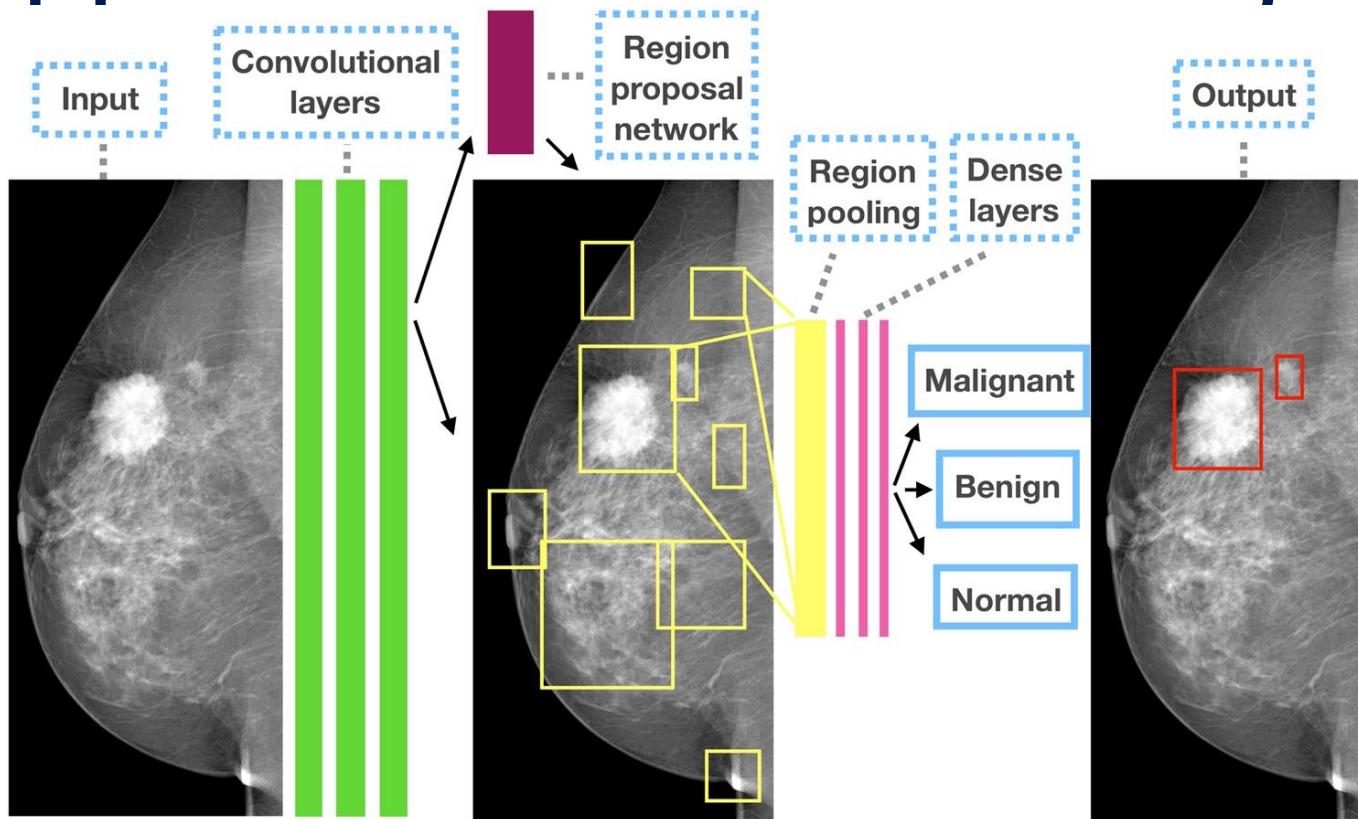
Feng, Rosenbaum,  
and Dietmayer,  
“Towards Safe  
Autonomous  
Driving: Capture  
Uncertainty in the  
Deep Neural  
Network For Lidar  
-In-Vehicle  
Detection,” UNIVERSITY

# Application: TSA



Liang et al., "Automatic threat recognition of prohibited items at aviation checkpoints with x-ray imaging: a deep learning approach".

# Application: Medical Imaging



Ribli, Dezső, et al. "Detecting and classifying lesions in mammograms with deep learning." *Scientific reports* 8.1 (2018): 4165.

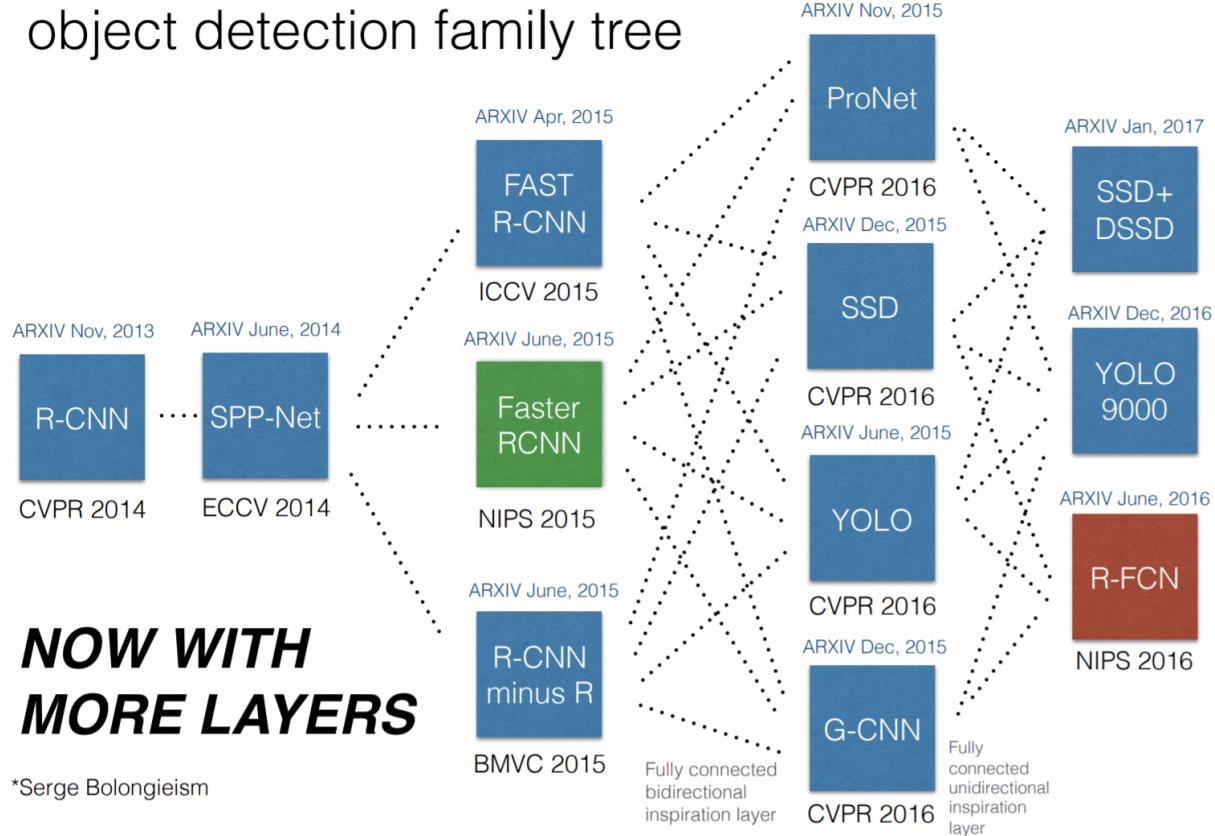
# Pieces of a Deep Algorithm

Object detection is very related to image classification. We use similar types of algorithms and related data types.

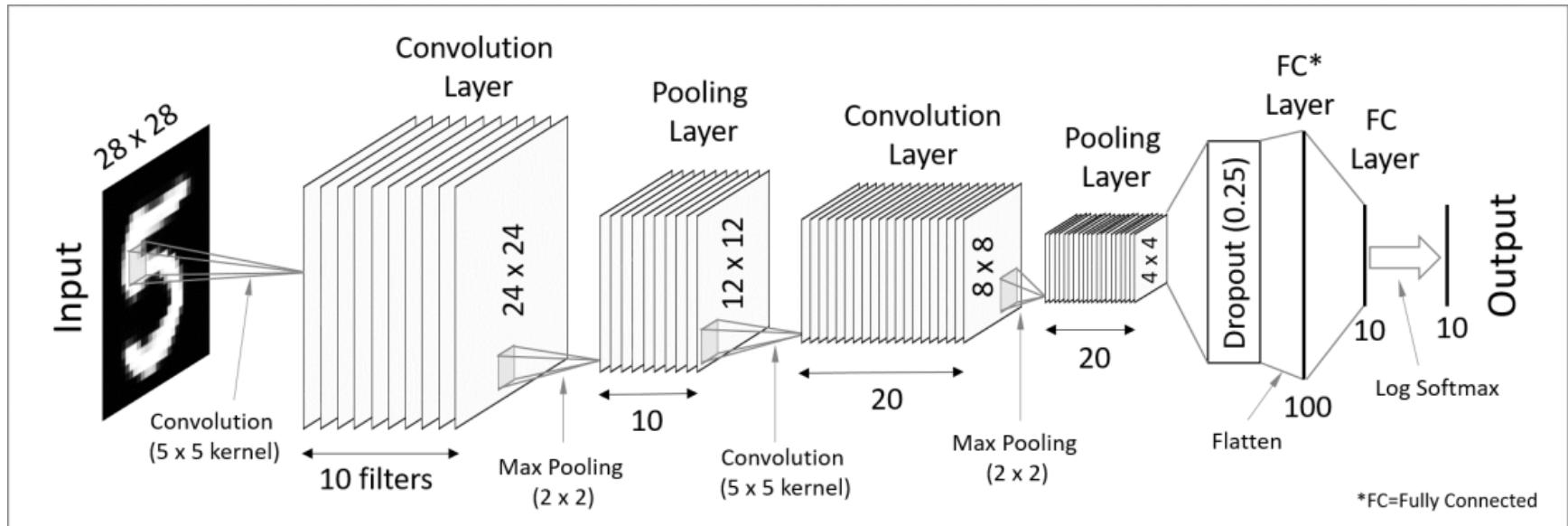
	<b>Image Classification</b>	<b>Object Detection</b>
<b>Model</b>	CNN Feature Extractor + Softmax Classifier	SSD, Faster R-CNN, e.g.
<b>Data</b>	Imagenet, CIFAR-10, e.g.	COCO, PASCAL VOC, e.g.
<b>Scalar loss function (Objective)</b>	CE loss + Regularization	SSD Loss, Faster R-CNN Loss, (both very complicated)
<b>Optimization Algorithm and Hyperparameters</b>	Stochastic Gradient Descent (SGD), Adam, Momentum, etc.	

Rapidly changing  
area of research

## Some members of the *postdeepluvian\** object detection family tree



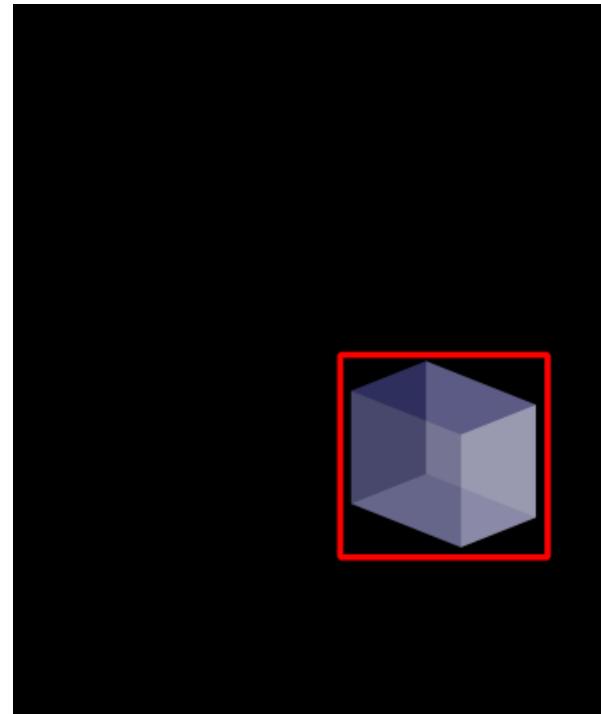
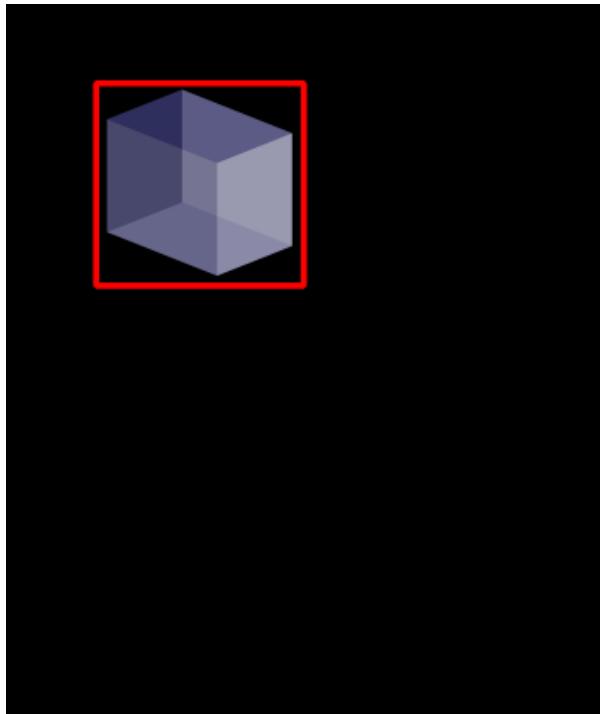
# Review: Convolutional Neural Network (MNIST)



# What do we want in object detection?

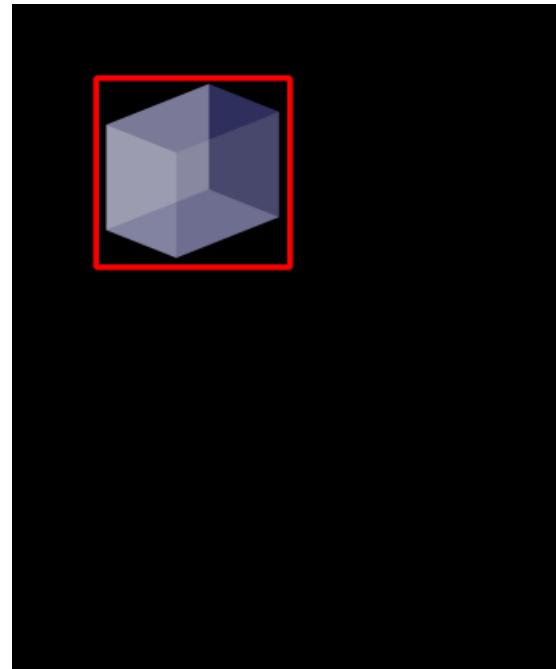
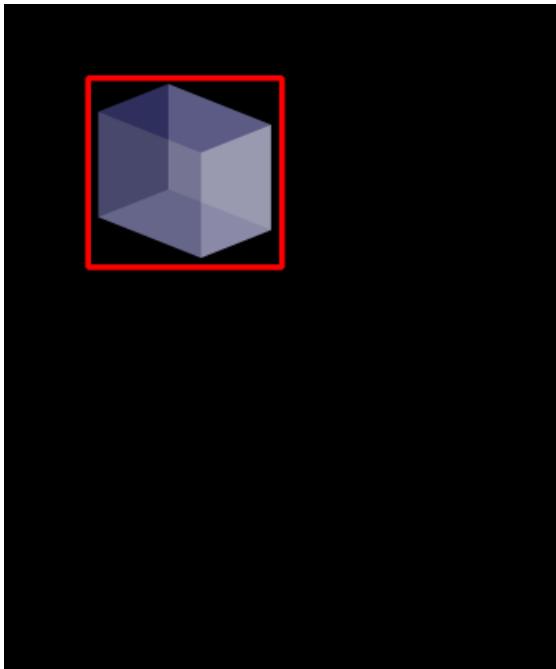
- We have some important properties that are useful in object detection
  - Translation Invariance
  - Scale Invariance
  - Rotation Invariance
- Some are built into the algorithm, and some come from the structure of the dataset

# Translation Invariance



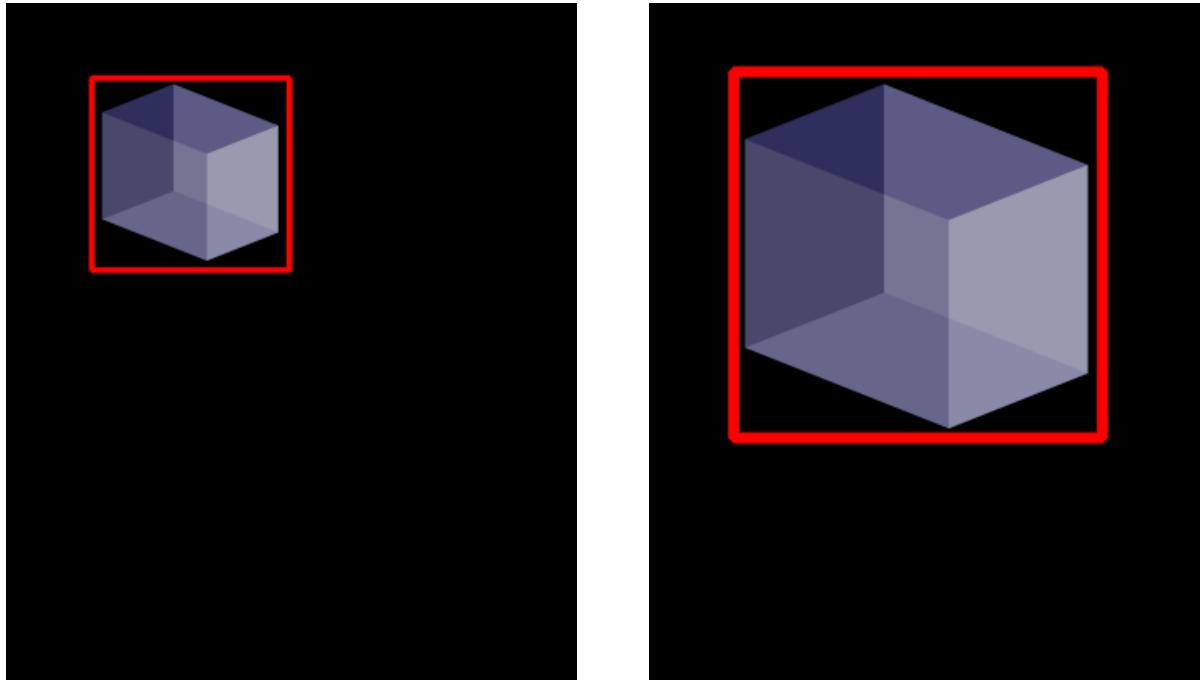
This is addressed with convolutional features...

# Rotation Invariance



Will be handled by data augmentation and other techniques in the dataset.

# Scale Invariance



Will be handled by data diversity and rescaling feature maps.

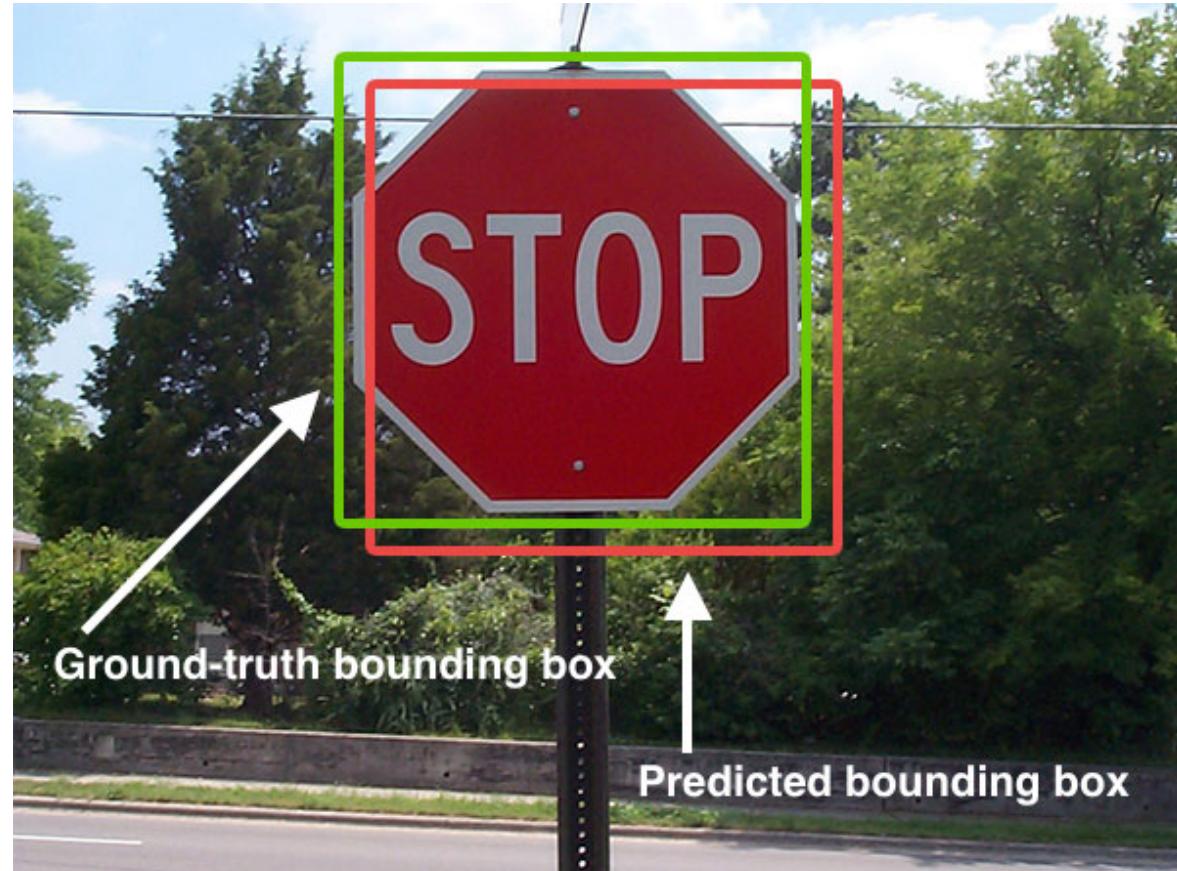
## We need a metric to evaluate our model

Classification – did we choose the correct class?

Object detection – is our object correct?

How do we determine if we got the correct location?

We want the bounding box to be good enough.



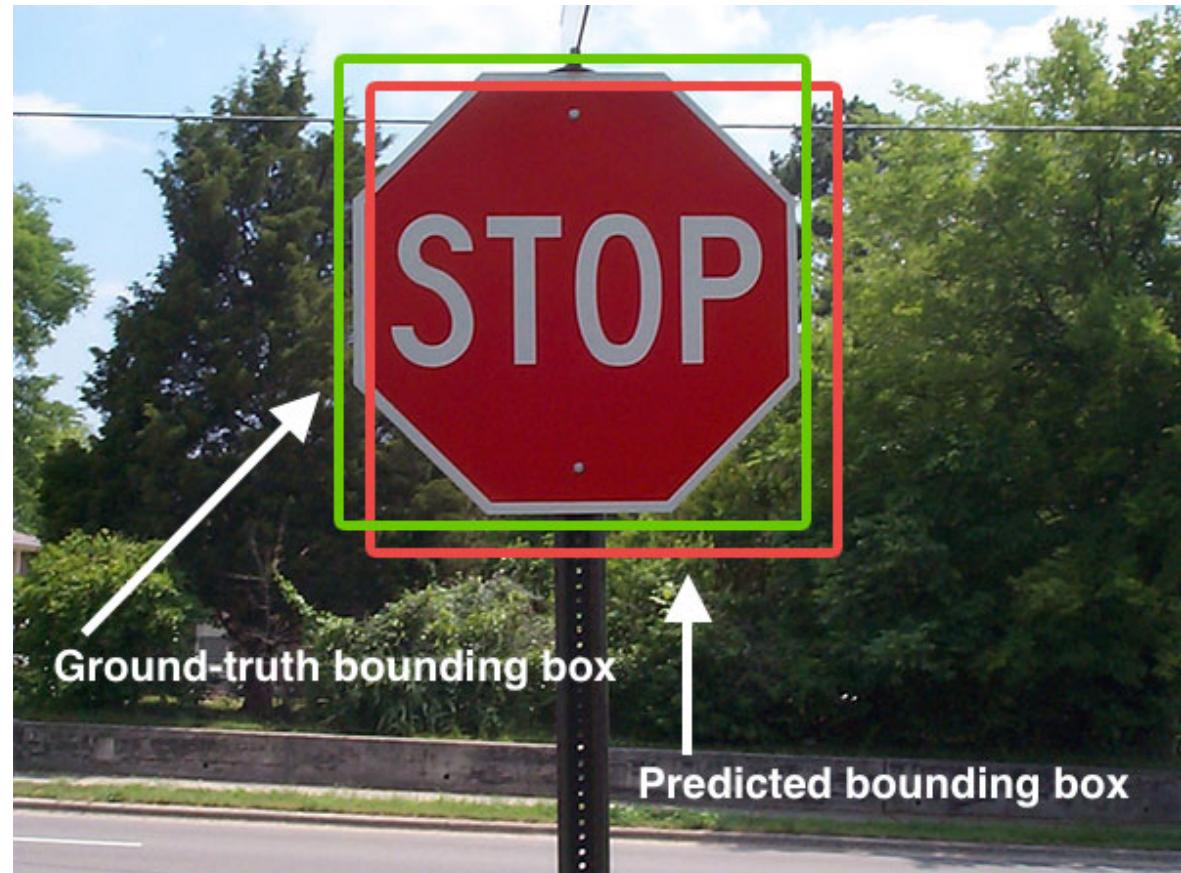
## Intersection-over-Union

The intersection-over-union metric is one way of defining this.

Intersection-the common area covered by the ground-truth bounding box and the predicted bounding box

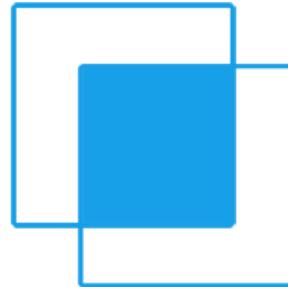
Union-the total area covered by either the ground-truth bounding box and predicted bounding box.

IoU=Intersection/Union



# Intersection-over-Union

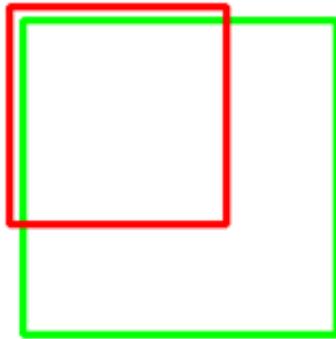
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



4

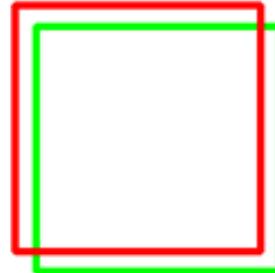
# Intersection-over-Union

IoU: 0.4034



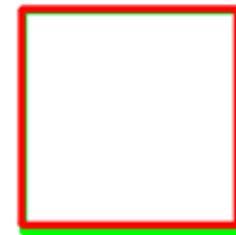
Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

5

# Two Types of Detectors

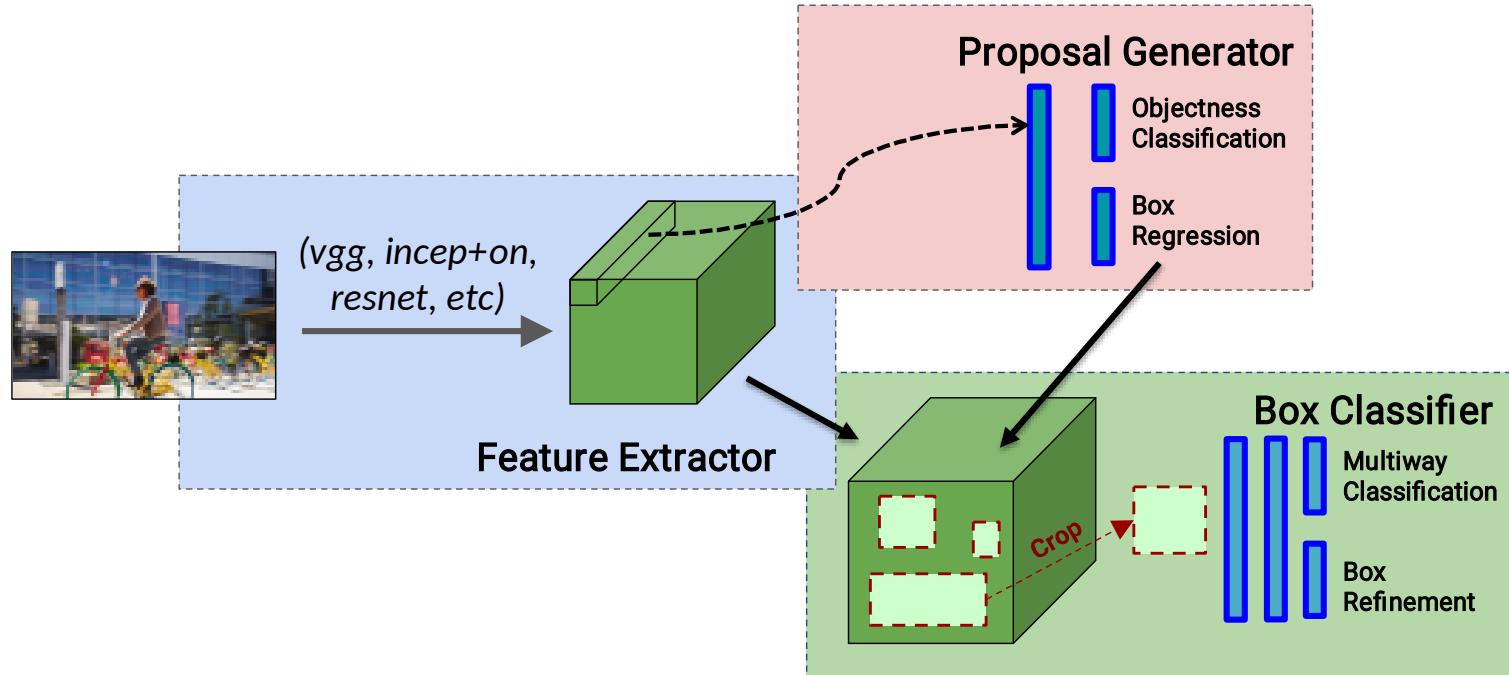
	One-Stage	Two-Stage
Examples	Single-Shot Multibox Detector (SSD), YOLO	Faster R-CNN
Speed	Faster (59 fps <sup>6</sup> )	Slower (5 fps <sup>7</sup> )
Accuracy	Less accurate (≈.24 mAP on COCO)	More accurate (≈.37 mAP on COCO)

<sup>6</sup>Liu et al., “Ssd: Single shot multibox detector”.

<sup>7</sup>Ren et al., “Faster r-cnn: Towards real-time object detection with region proposal networks” .

# **HOW DO CONVOLUTIONAL FEATURES LEAD TO OBJECT DETECTION?**

# Faster RCNN

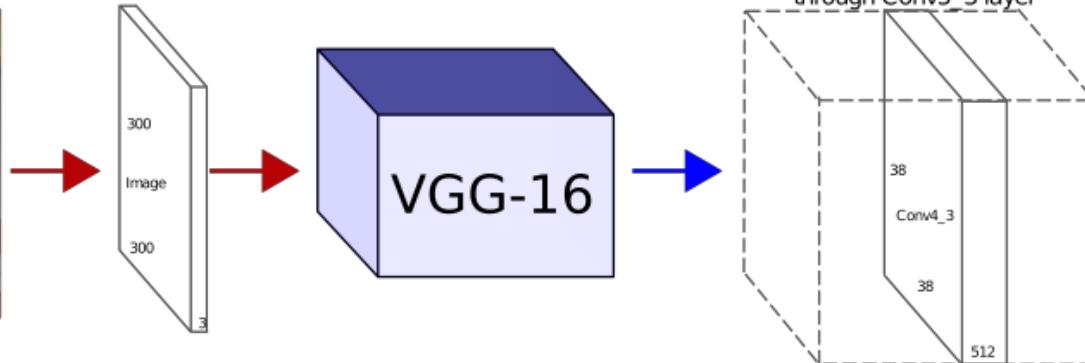


Huang et al., "Speed/accuracy trade-offs for modern convolutional object detectors"

Ren et al

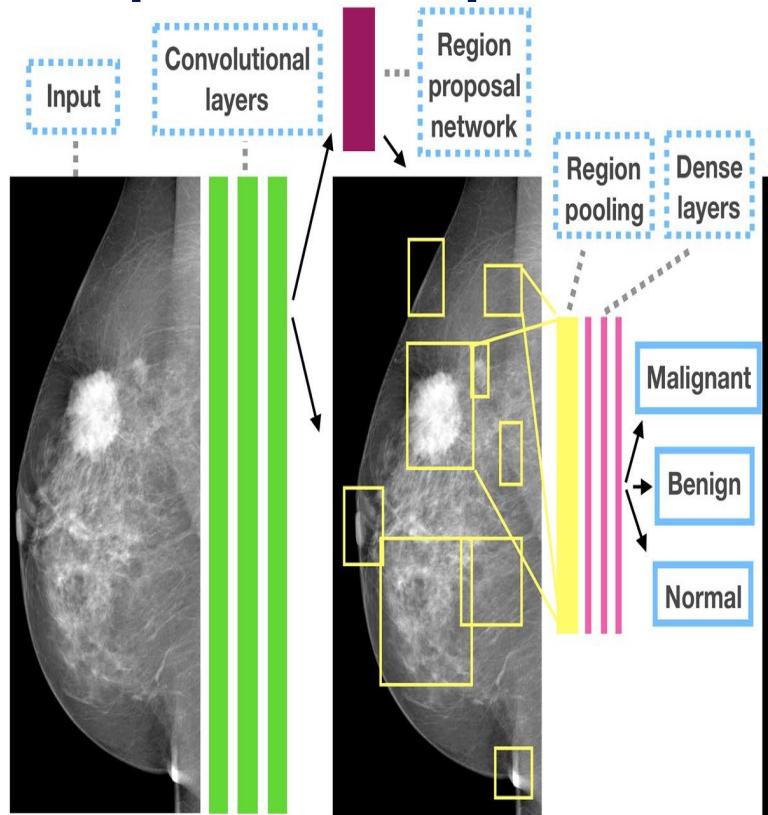
# Stage 1, Step 1: Feature Extraction

Input Image  
(479,479,3)



Feature Maps  
(38,38,512)

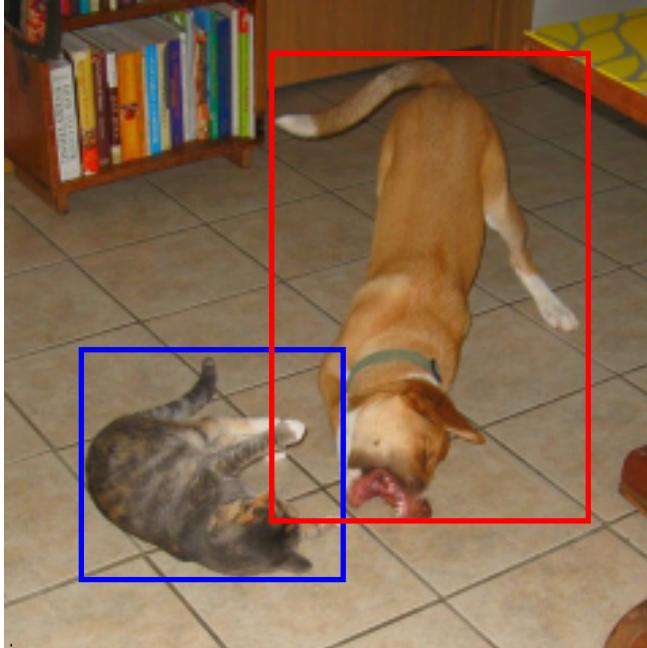
# Stage 1, Step 2: Proposal Generation



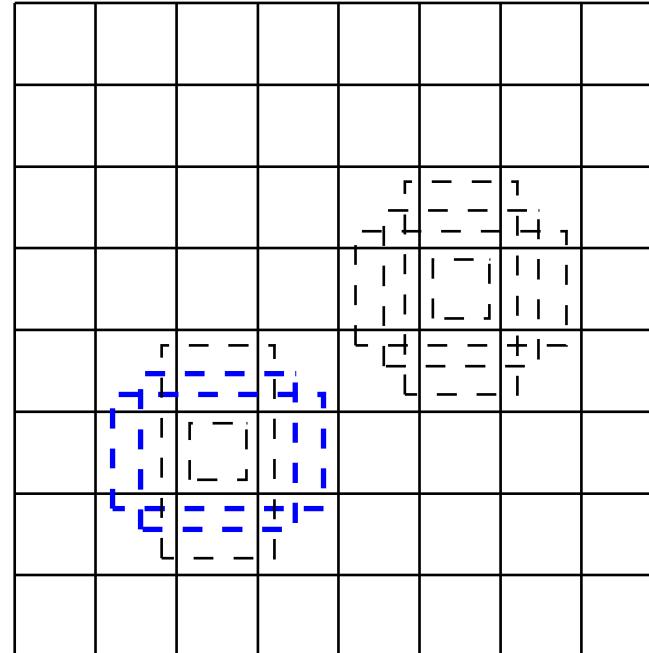
Ribli, Dezső, et al. "Detecting and classifying lesions in mammograms with deep learning." *Scientific reports* 8.1 (2018): 4165.

# Anchor Boxes

Image

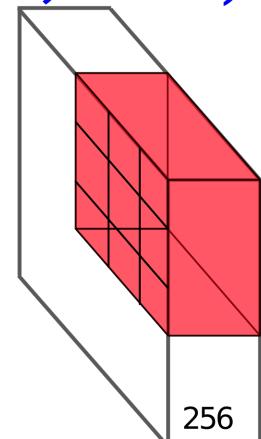


Features

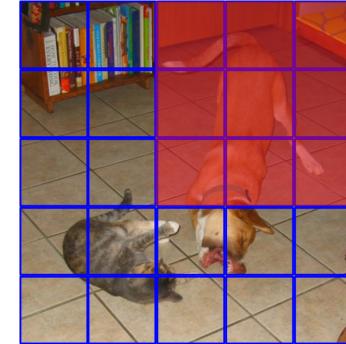


Consider only one section  
of an image at a time

Feature Maps  
(5,5,256)

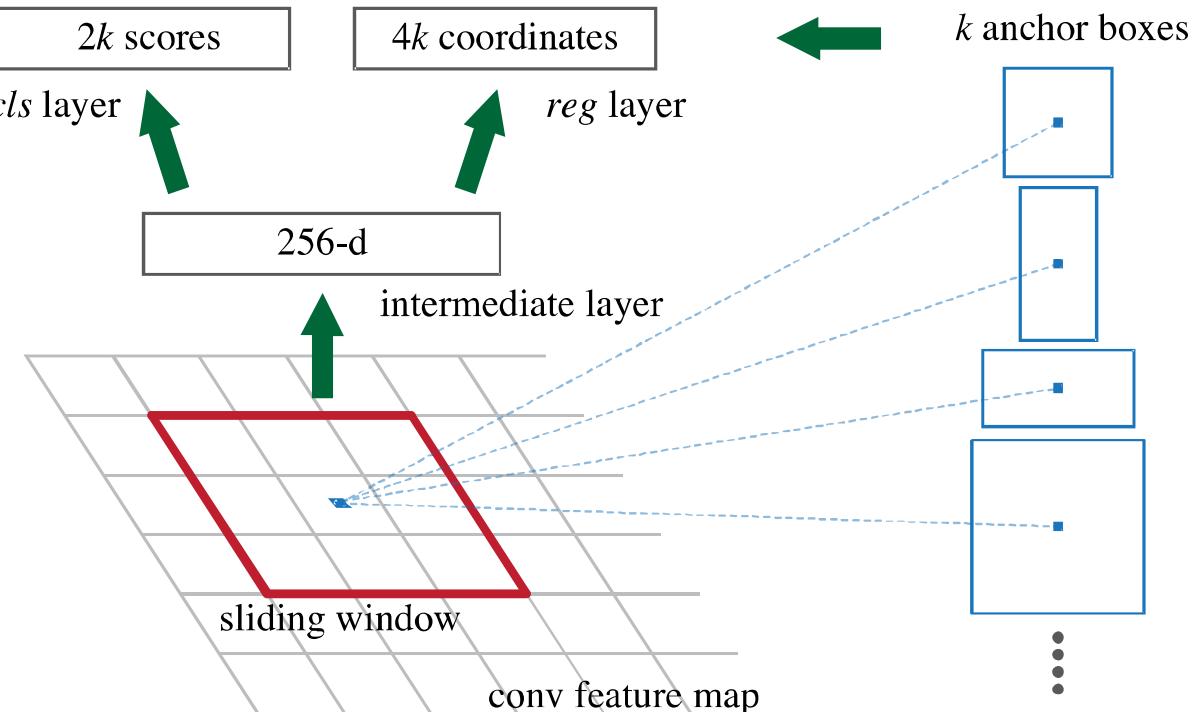


Conv: 1x1x128  
Conv: 3x3x256-s2



Sliding Convolutional  
Kernel (3,3,256)

# Region Proposal Network



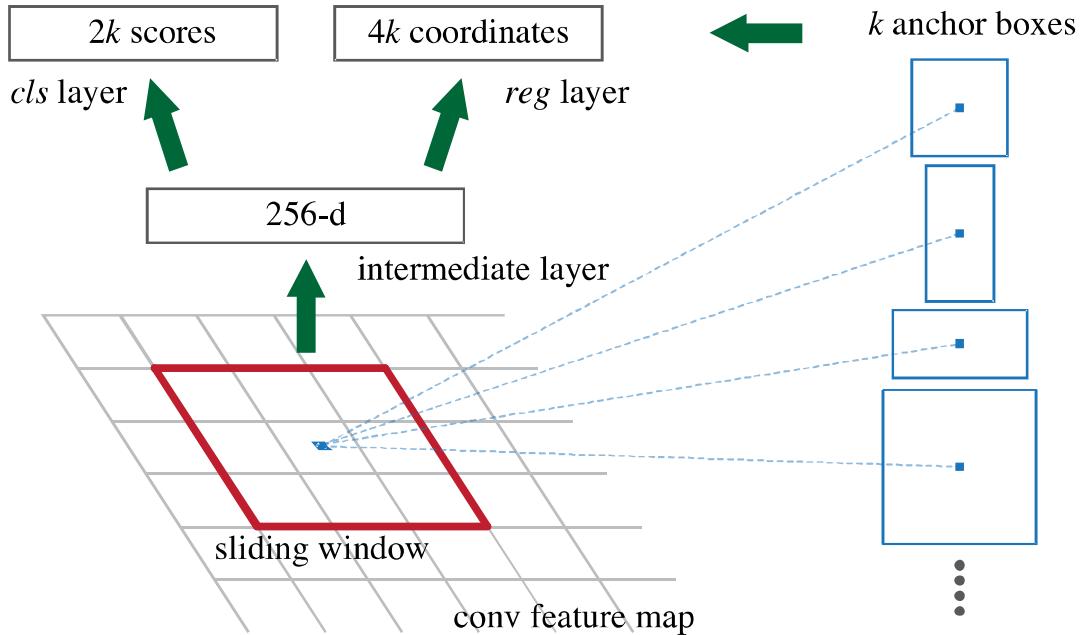
Ren et al., "Faster r-cnn: Towards real-time object detection with region proposal networks"

# Region proposal network

Using a sliding window on our convolutional filters, we first produce an intermediate layer of features.

We predict  $4k$  coordinates per bound box.

What does this mean?



Ren et al., "Faster r-cnn: Towards real-time object detection with region proposal networks"

## How do we predict bounding boxes?

Start with a “base” bounding box.

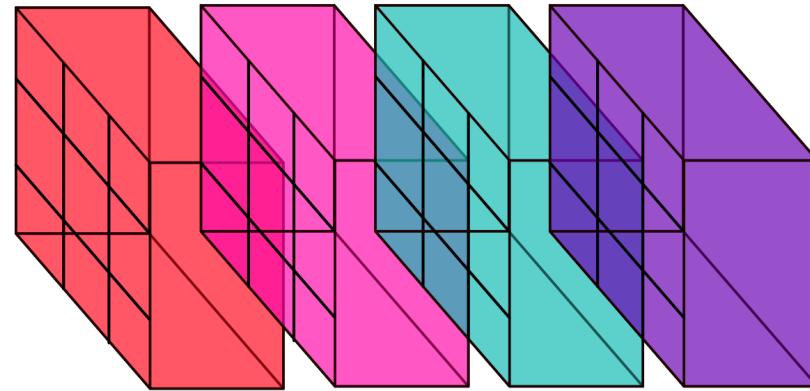
Predict four terms:

Center ( $x, y$ ) of the bounding box

Width of the bounding box  
(typically a multiple, 1x is the default box)

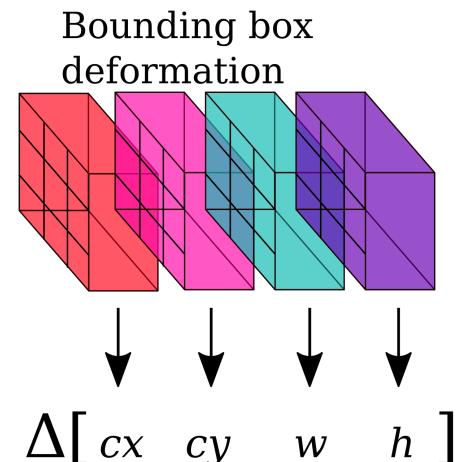
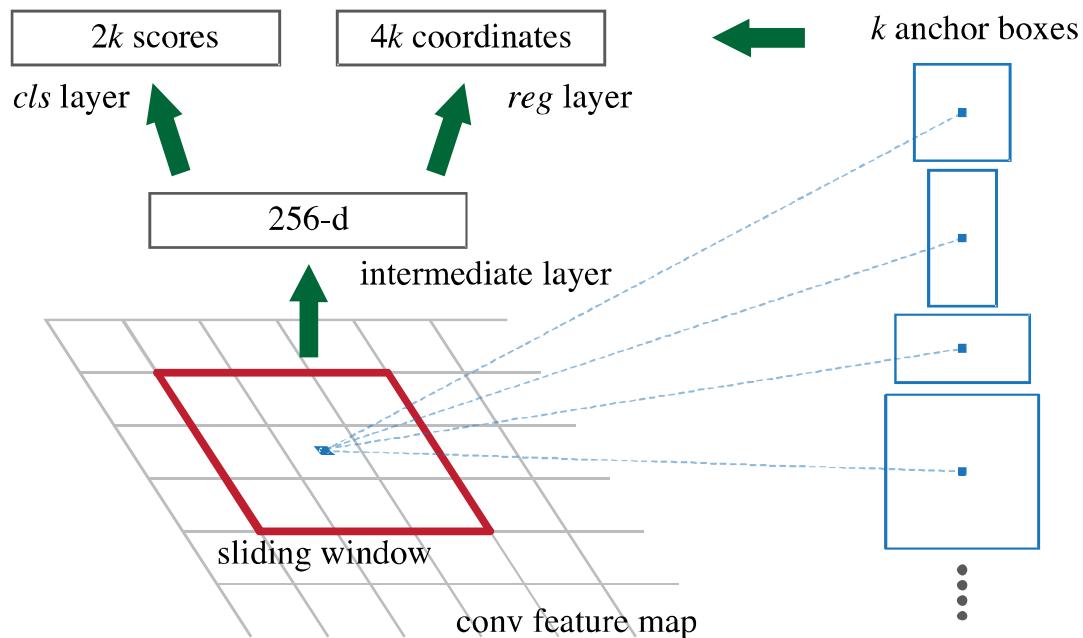
Height of the bounding box  
(typically a multiple, 1x is the default box)

## Bounding box deformation



$$\Delta [ \begin{matrix} cx & cy & w & h \end{matrix} ]$$

# How do we predict bounding boxes?



# Mathematically, what do we predict?

Let  $a$  denote the “anchor” box and  $a^*$  denote ground truth boxes.

Then we can define our predictions as:

Then we can define our predictions as:

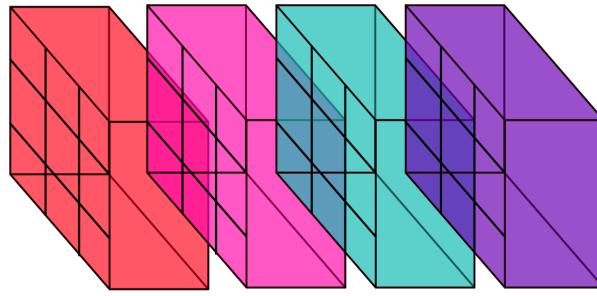
$$tx = \frac{cx - x_a}{w_a}$$

Can calculate same quantities with respect to ground truth (necessary for creating loss functions).

$$tw = \log\left(\frac{w}{w_a}\right)$$

Can calculate same quantities with respect to ground truth (necessary for creating loss functions).

## Bounding box deformation



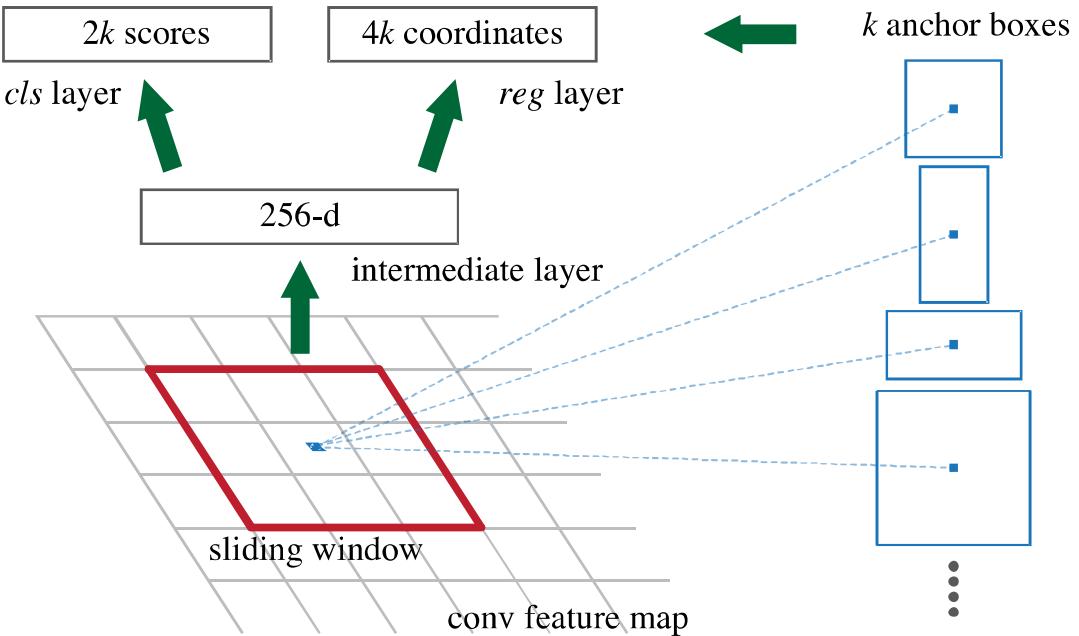
$$\Delta [ cx \quad cy \quad w \quad h ]$$

# What do our scores do?

Our scores simply return probabilities of “box” or “no box.”

Box means that this is a real bounding box (i.e. we’re confident that we have a real box)

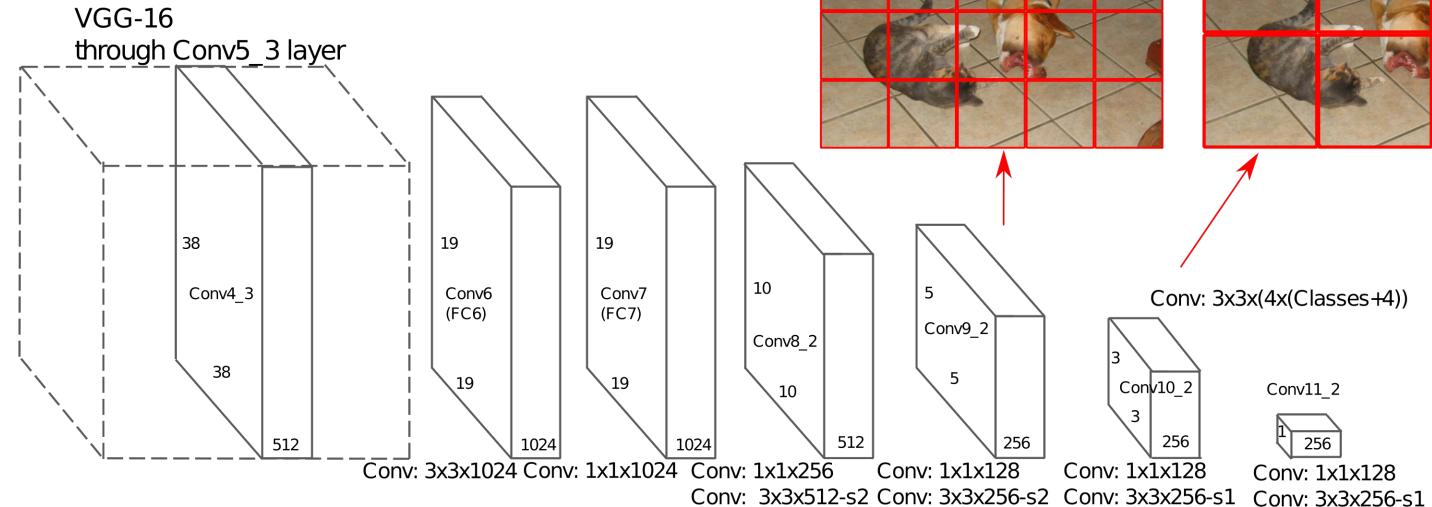
Want to ignore everything else. There are *lots* of incorrect proposals...



Ren et al., “Faster r-cnn: Towards real-time object detection with region proposal networks”

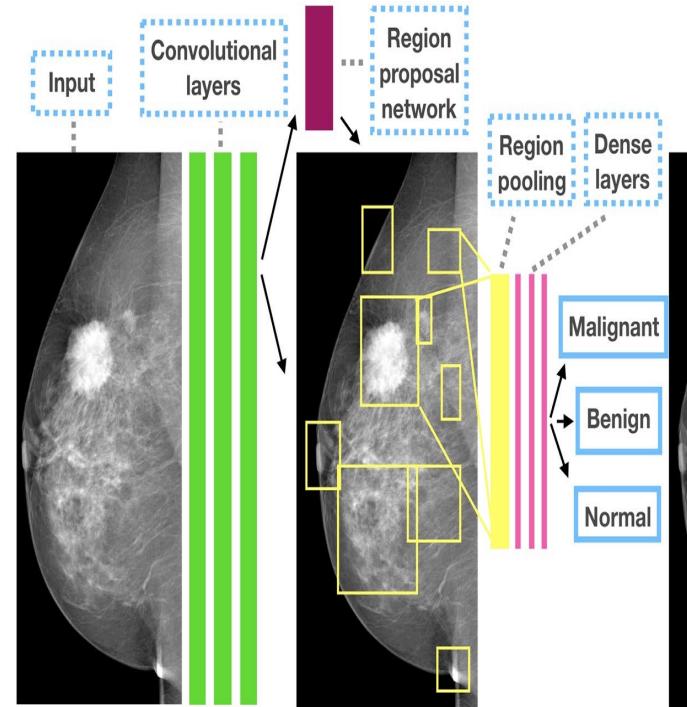
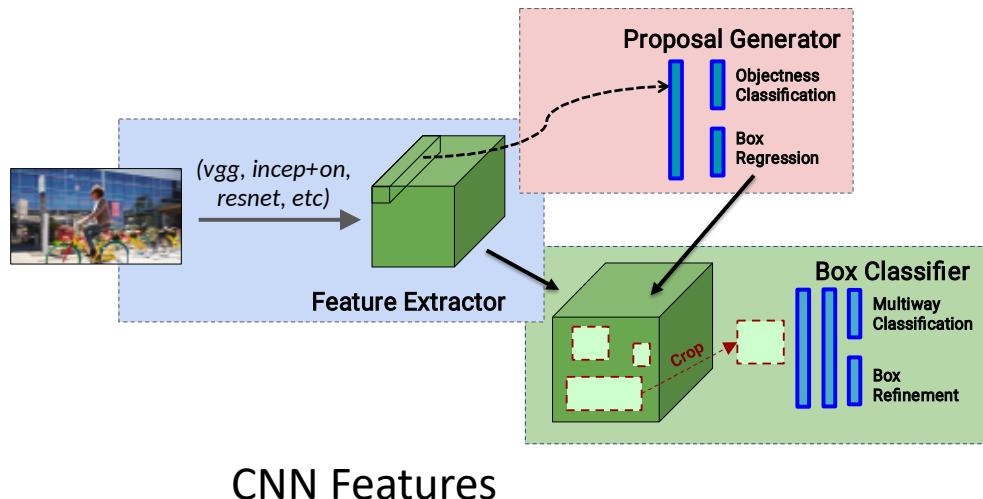
# Can Consider Multiple Scales

Feature Maps  
(30,30,1024)



# Stage 2: Classification

Second stage object classifier

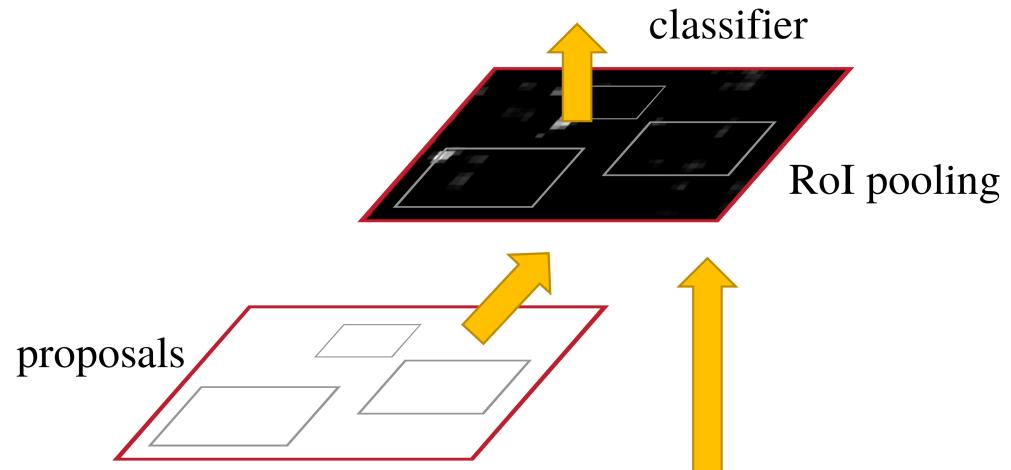


Ren et al., "Faster r-cnn: Towards real-time object detection with region proposal networks"

# How to we determine what's in a box?

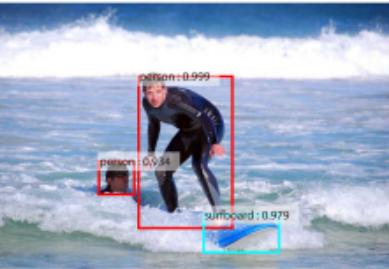
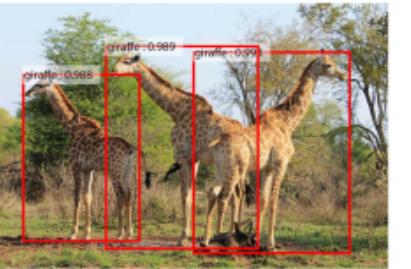
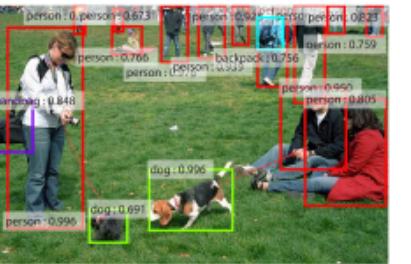
“Region-of-Interest (ROI) pooling” is a critical step to reduce complexity of later bounding boxes. It forces all boxes to have the same feature dimensionality.

Vector of features for the box from the CNN goes through a classifier (typically an MLP) to predict one of the classes or **background**.



Ren et al., “Faster r-cnn: Towards real-time object detection with region proposal networks”

# Sample Detections



# How can this be learned?

- Object detection adds *a lot* of complexity compared to a CNN
- Many practical details become important for effective training
- I will go over key details; read the papers if you go to implement this yourself

# Positive and Negative Boxes

- For evaluation, we consider the following a correct box:
  - $\geq 0.5$  IOU: Correct
  - $< 0.5$  IOU: Incorrect
- For training:
  - Positive:  $0.7 \leq$  IOU:
  - Ignore:  $0.3 \leq$  IOU  $< 0.7$
  - Negative: IOU  $\leq 0.3$

# Training Loss for the RPN

- Define
  - =  $p_i$ : probability anchor  $i$  is an object
  - =  $p_i^*$ : indicator object is positive (i.e. 0/1)
  - =  $t_i$ : vector containing coordinates for bounding box  $i$
  - =  $t_i^*$ : coordinates of ground truth bounding box  $i$
- The total loss for the model is given as:

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*).$$

# Training Loss for the RPN

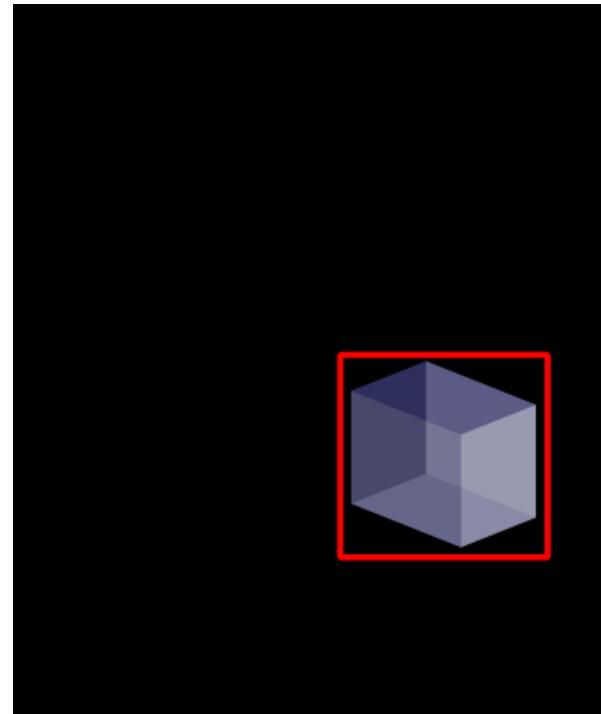
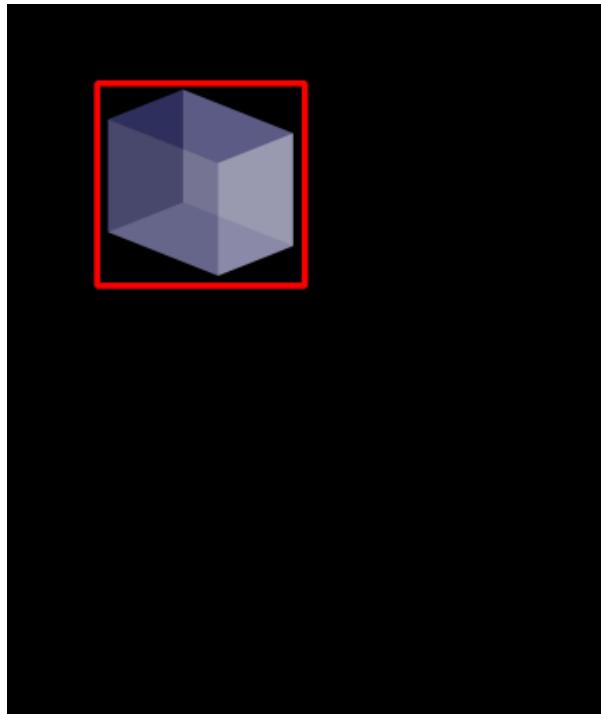
$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*).$$

- The first term is a cross-entropy loss on the object predictor (i.e. object/not object)
- Second term is on the box dimensions, which is an  $\text{I}1$  loss
- Given the boxes, it's just a multi-class classification (often treated with a separate loss given just the boxes).

# Let's go back to our goals:

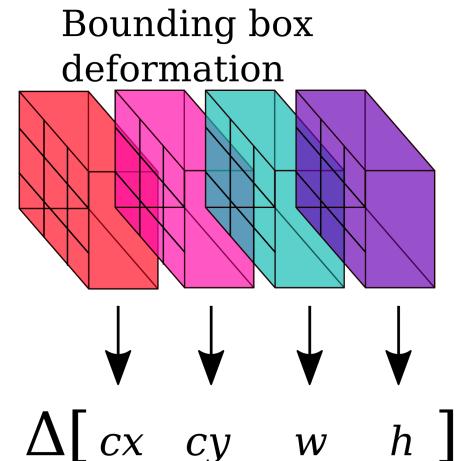
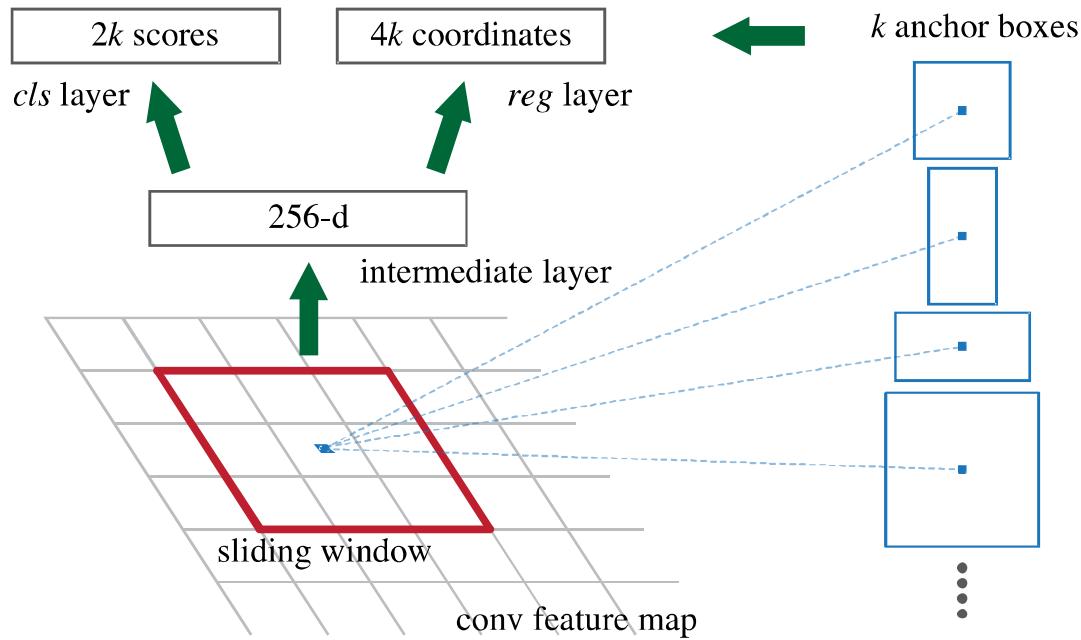
- We have some important properties that are useful in object detection
  - Translation Invariance
  - Scale Invariance
  - Rotation Invariance
- Some are built into the algorithm, and some come from the structure of the dataset

# Translation Invariance

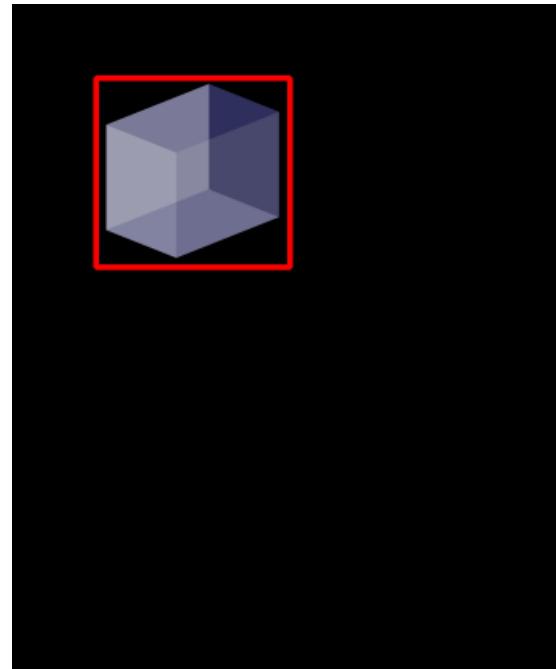
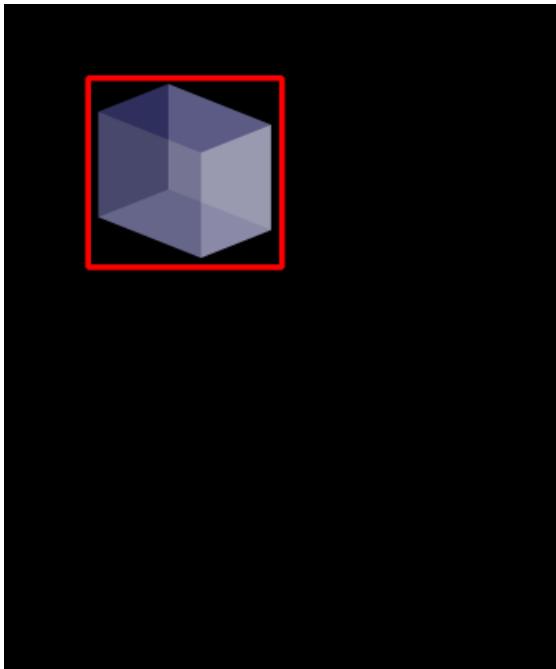


This is addressed with convolutional features...

# Translation invariance comes from the sliding windows (same as conv.)

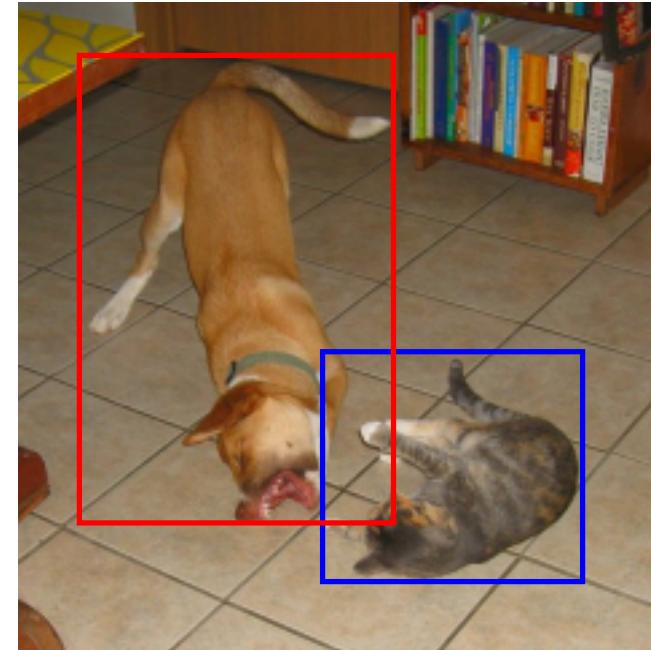
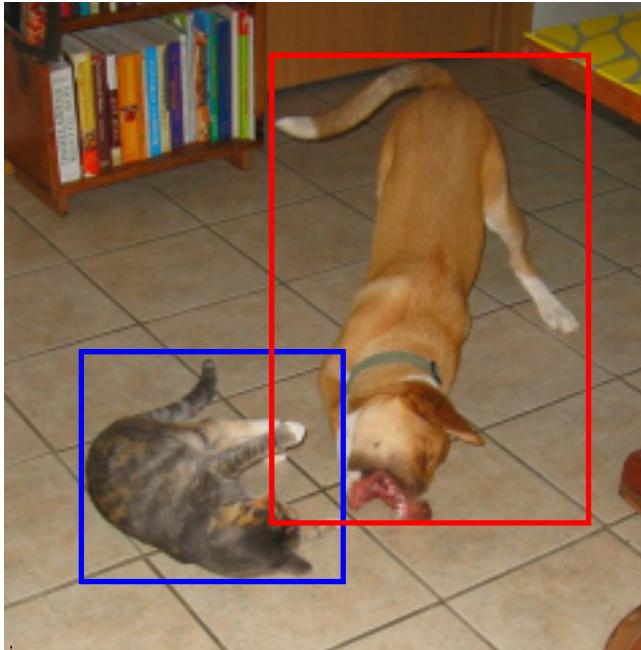


# Rotation Invariance



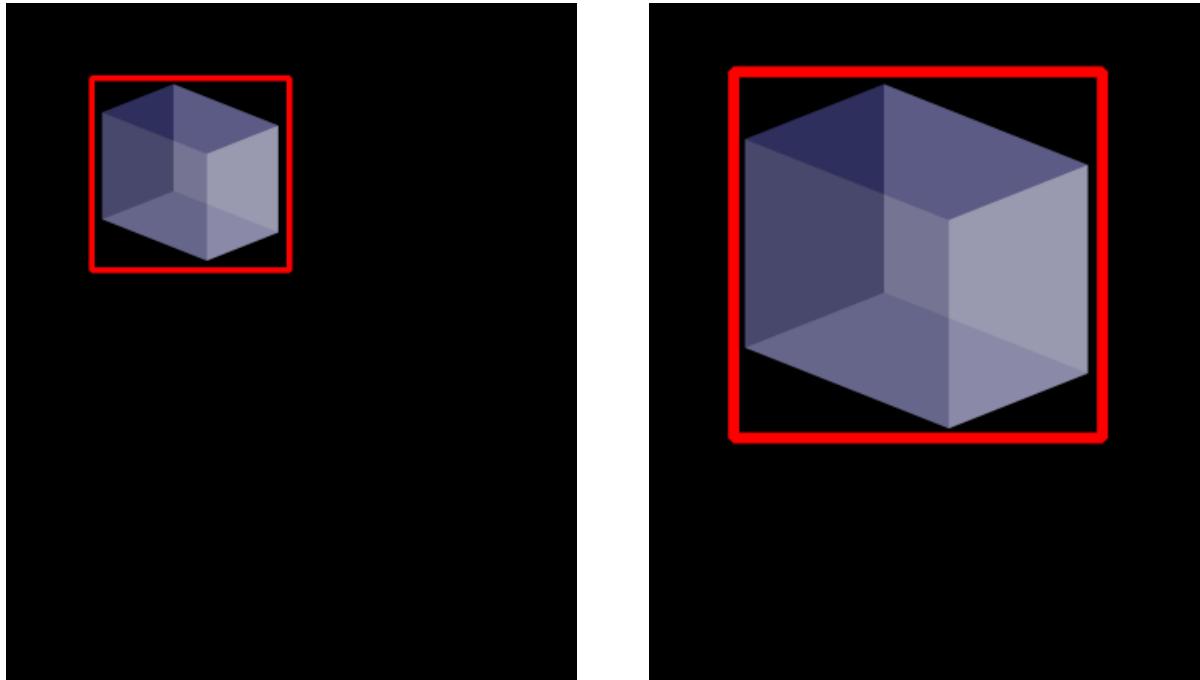
Will be handled by data augmentation and other techniques in the dataset.

# Data Augmentation



Feed in mirrors, rotations, shifts, etc.

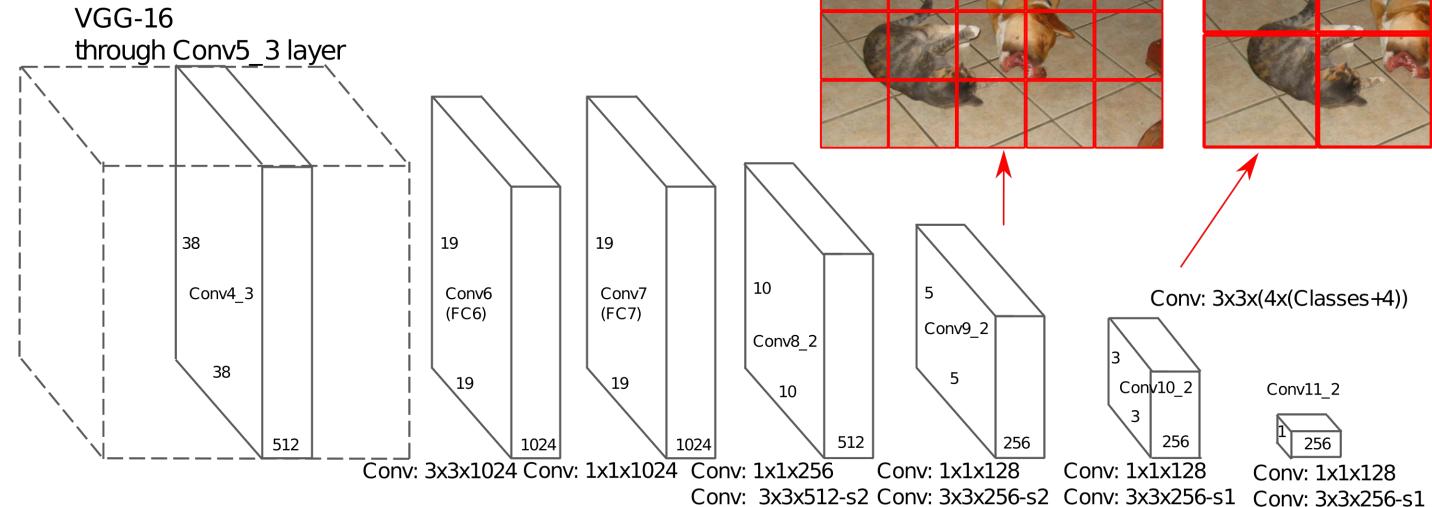
# Scale Invariance



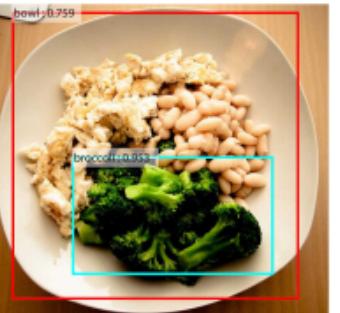
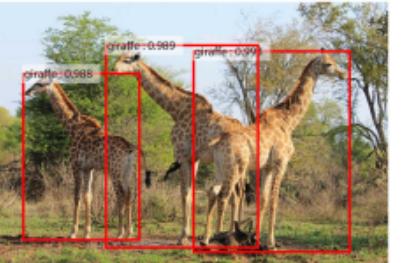
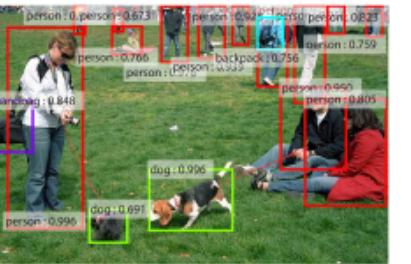
Will be handled by data diversity and rescaling feature maps.

# Pyramid or Multiscale

Feature Maps  
(30,30,1024)



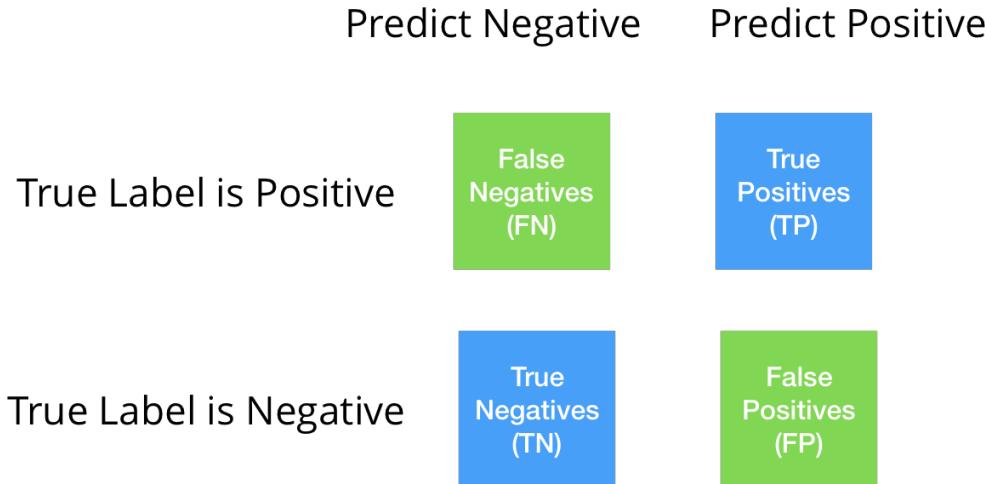
# Sample Detections



# How do we measure performance

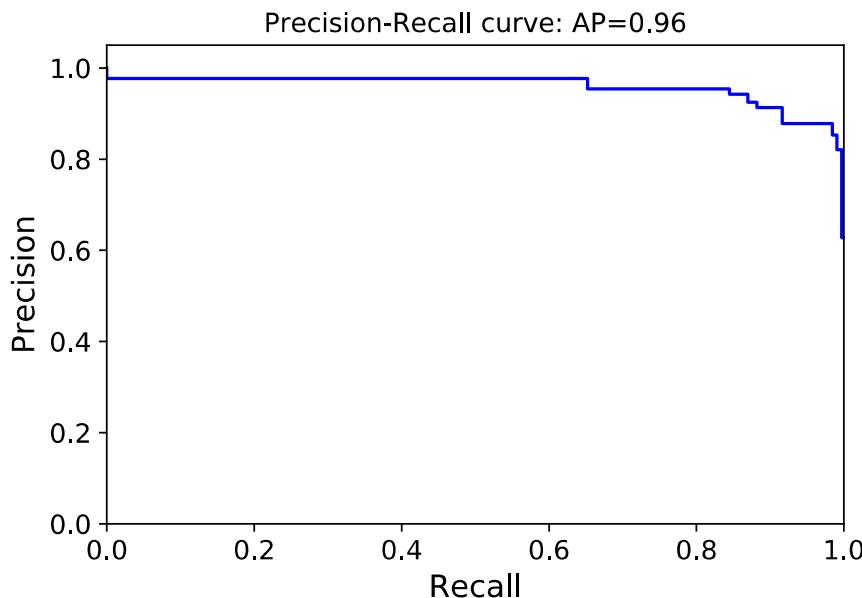
- We use mean average precision (mAP) over all of our classes
- Must use a precision-recall metric  
*because there are so many ways to get things wrong*
  - Classes are super imbalanced

# Types of Predictions



*Precision* is given by  $(TP)/(TP+FP)$ .  
*Recall* is given by  $(TP)/(TP+FN)$ .

# Precision-Recall

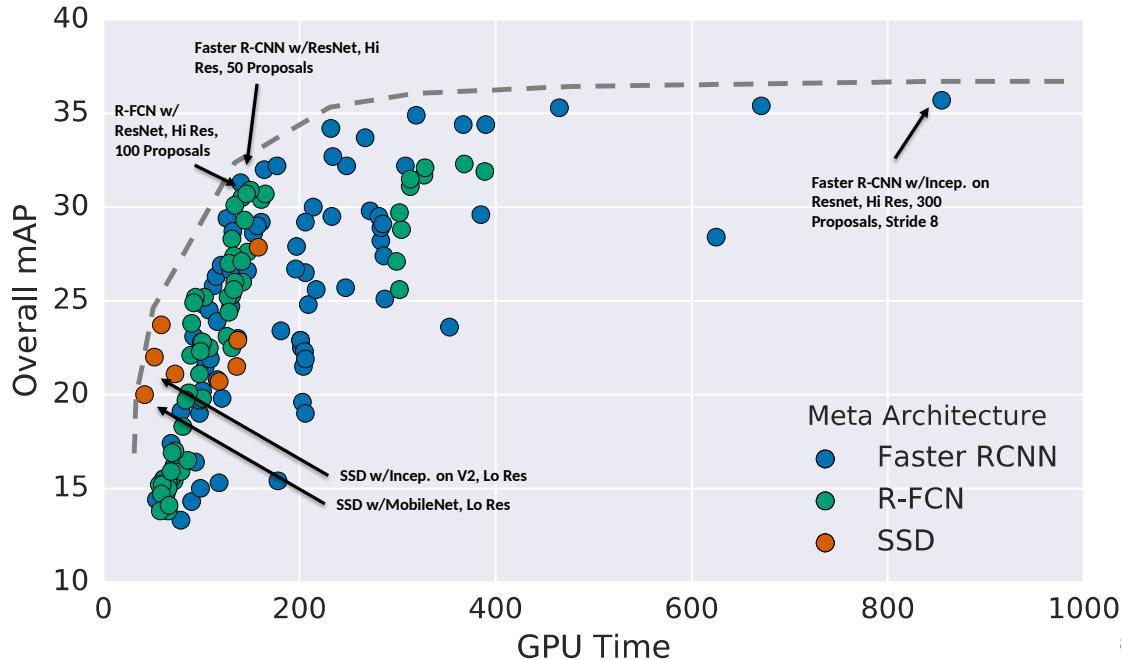


*Average Precision* is the integral of this curve.

# Mean Average Precision

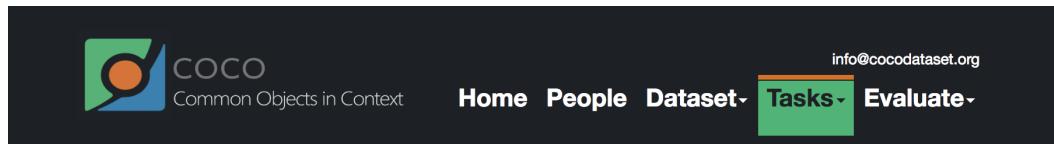
- When we have many classes or individualized responses, it can be difficult to evaluate. Accuracy is dependent on decision thresholds, which means that it can be gamed a bit and may be hard to compare across algorithms.
- Average Precision (AP) isn't perfect, but works well for a single class.
- When we have multiple classes, we can take the "Mean Average Precision," which is the average precision over multiple classes:  
 $mAP$  (or  $MAP$ ) = mean (AP for each class).
- Other common metrics include:  
 $F1=2 \cdot (\text{Precision} \cdot \text{Recall}) / (\text{Precision} + \text{Recall})$
- $MRR=$ mean reciprocal rank (useful when only one answer is appropriate out of many)

# Performance Tradeoffs



Huang et al., "Speed/accuracy trade-offs for modern convolutional object detectors"

# Where is object detection going?



## COCO 2018 Object Detection Task



### 1. Overview

The COCO Object Detection Task is designed to push the state of the art in object detection forward. COCO features two object detection tasks: using either bounding box output or object segmentation output (the latter is also known as instance segmentation). For full details of this task please see the [detection evaluation](#) page. Note: **only the detection task with object segmentation output will be featured at the COCO 2018 challenge** (more details follow below).

<http://cocodataset.org/#detection-2018>

# Resources for training your own object detector

## Tensorflow Object Detection API

Creating accurate machine learning models capable of localizing and identifying multiple objects in a single image remains a core challenge in computer vision. The TensorFlow Object Detection API is an open source framework built on top of TensorFlow that makes it easy to construct, train and deploy object detection models. At Google we've certainly found this codebase to be useful for our computer vision needs, and we hope that you will as well.



Contributions to the codebase are welcome and we would love to hear back from you if you find this API useful. Finally if you use the Tensorflow Object Detection API for a research publication, please consider citing:

# Conclusions

- Object detection leads to a similar task as object classification, but with different goals
- Lots of technical details, but feasible in a lot of different challenges
- *Many* real-world applications