
SSD: Single Shot MultiBox Detector

*Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed,
Cheng-Yang Fu, Alexander C. Berg*
[\[arXiv\]](#)[\[demo\]](#)[\[code\]](#) (Mar 2016)



Image Processing Group
Signal Theory and Communications Department
Universitat Politècnica de Catalunya. BARCELONATECH

Slides by Míriam Bellver
Computer Vision Reading Group, UPC
28th October, 2016

Outline

- ▷ Introduction
- ▷ Related Work
- ▷ The Single-Shot Detector
- ▷ Conclusions



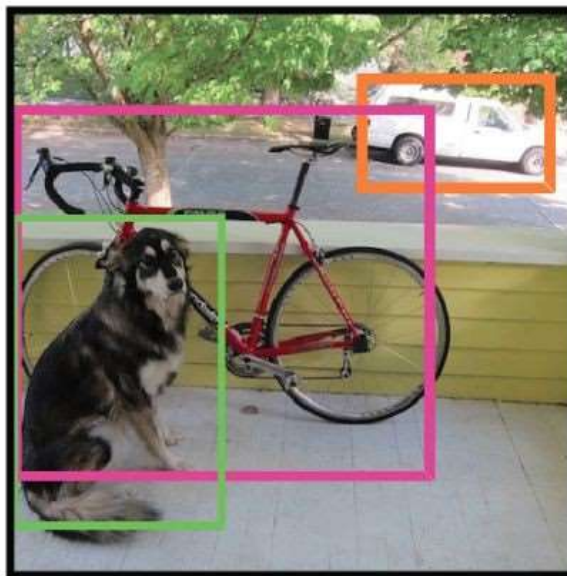
1.

Introduction

SSD: Single Shot MultiBox Detector

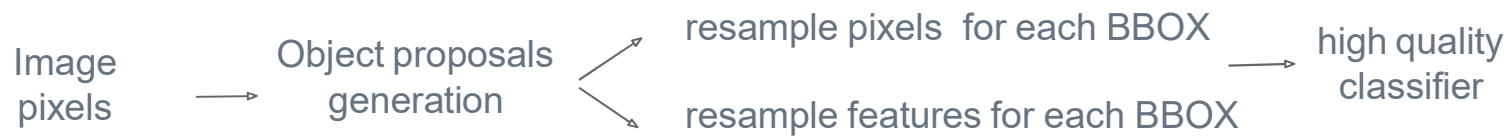
Introduction

Object detection



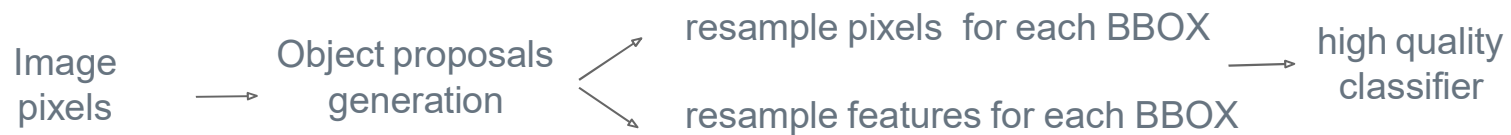
Introduction

Current object detection systems



Introduction

Current object detection systems



Computationally too intensive and too slow for real-time applications

Faster R-CNN 7 FPS



Introduction

Current object detection systems



Computationally too intensive and too slow for real-time applications

Faster R-CNN 7 FPS



Introduction

SSD: First deep network based object detector that does **not resample pixels or features** for bounding box hypotheses and is **as accurate as approaches that do**.

Improvement in **speed vs accuracy trade-off**



Introduction

Method	<i>mAP</i>	FPS	# Boxes
Faster R-CNN [2](VGG16)	73.2	7	300
Faster R-CNN [2](ZF)	62.1	17	300
YOLO [5]	63.4	45	98
Fast YOLO [5]	52.7	155	98
SSD300	72.1	58	7308
SSD500	75.1	23	20097

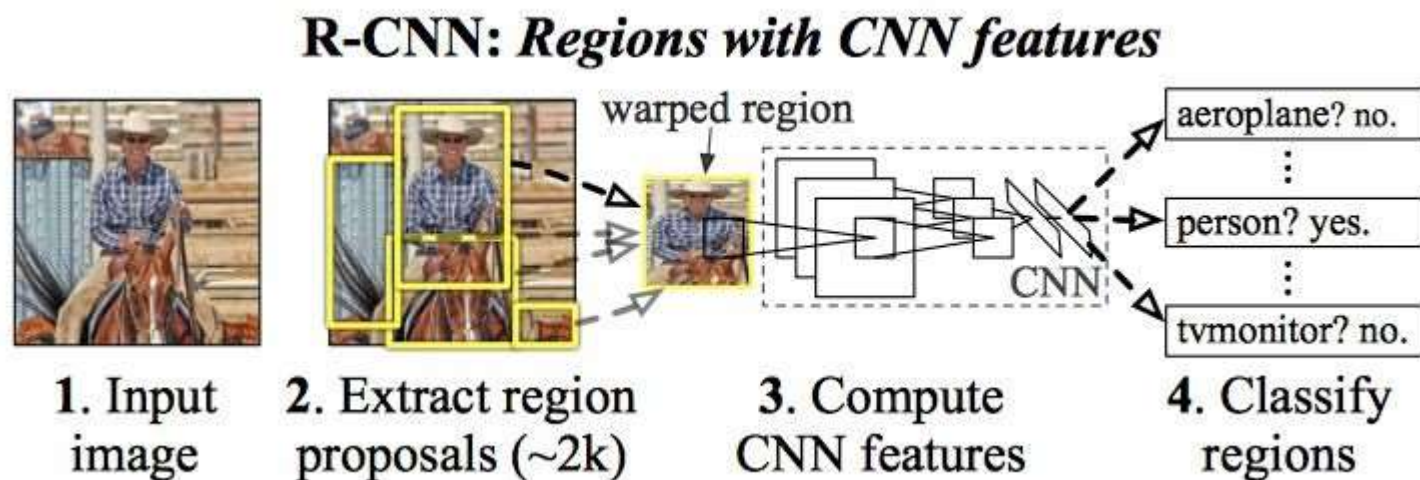
2.

Related Work

SSD: Single Shot MultiBox Detector

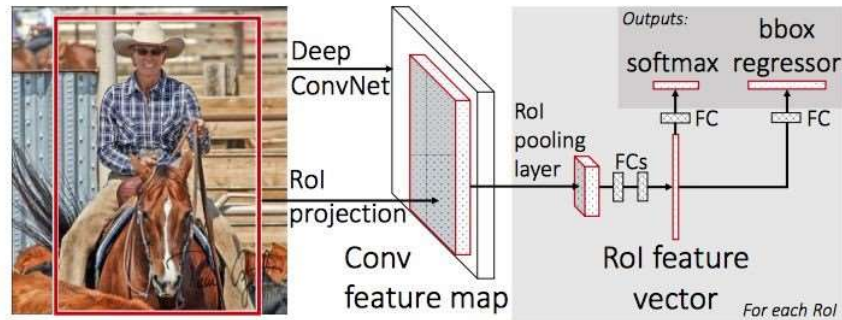
Object detection with CNN's

▷ R-CNN



Leveraging the object proposals bottleneck

▷ Fast R-CNN



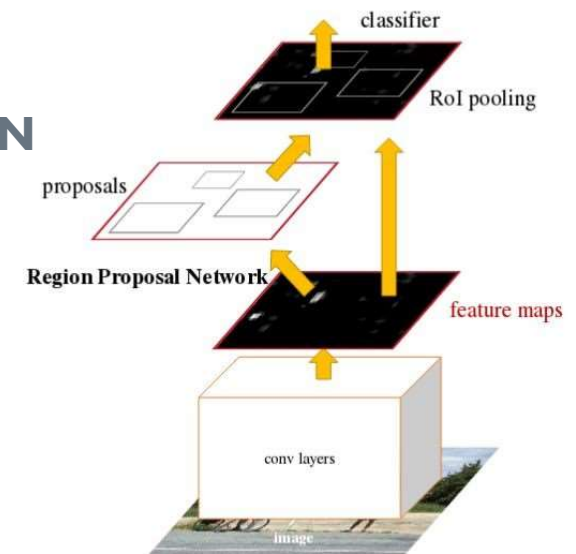
Improving quality of proposals using CNNs

Low-level features object proposals



Proposals generated directly from a DNN

Faster R-CNN



Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster R-CNN: Towards real-time object detection with region proposal networks

Szegedy, C., Reed, S., Erhan, D., & Anguelov, D. (2014). Scalable, high-quality object detection.

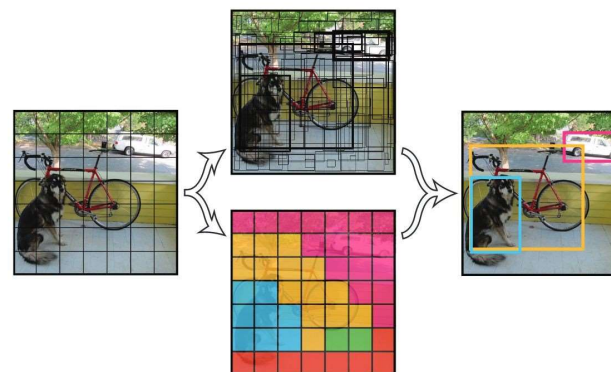
Single-shot detectors

Instead of having two networks

Region Proposals Network + Classifier Network

In Single-shot architectures, bounding boxes and confidences for multiple categories are predicted directly with a single network

e.g.: Overfeat, YOLO



Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2015). You only look once: Unified, real-time object detection
Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., & LeCun, Y. (2013). Overfeat: Integrated recognition, localization and detection using convolutional networks.

Single-shot detectors

Main differences of SSD over YOLO and Overfeat:

Small **conv.filters** to predict object **categories and offsets in BBs locations**,
using :

- separate predictors for different aspect ratios
- different feature maps to perform detection on multiple scales



Single-shot detectors

Contributions:

- ▷ A single-shot detector for multiple categories that is faster than state of the art single shot detectors (YOLO) and as accurate as Faster R-CNN
- ▷ Predicts category scores and boxes offset for a fixed set of default BBs.
- ▷ Predictions of different scales from feature maps of different scales, and separate predictions by aspect ratio
- ▷ End-to-end training and high accuracy, **improving speed vs accuracy trade-off**

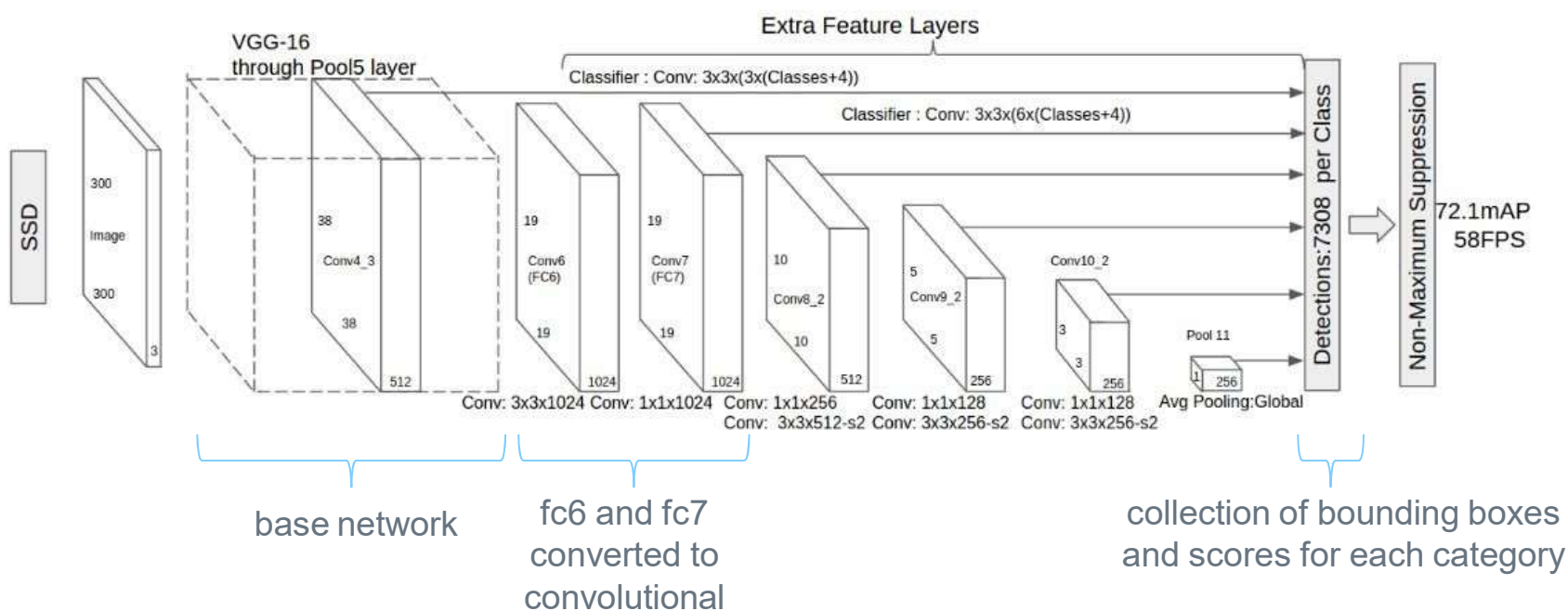


3.1

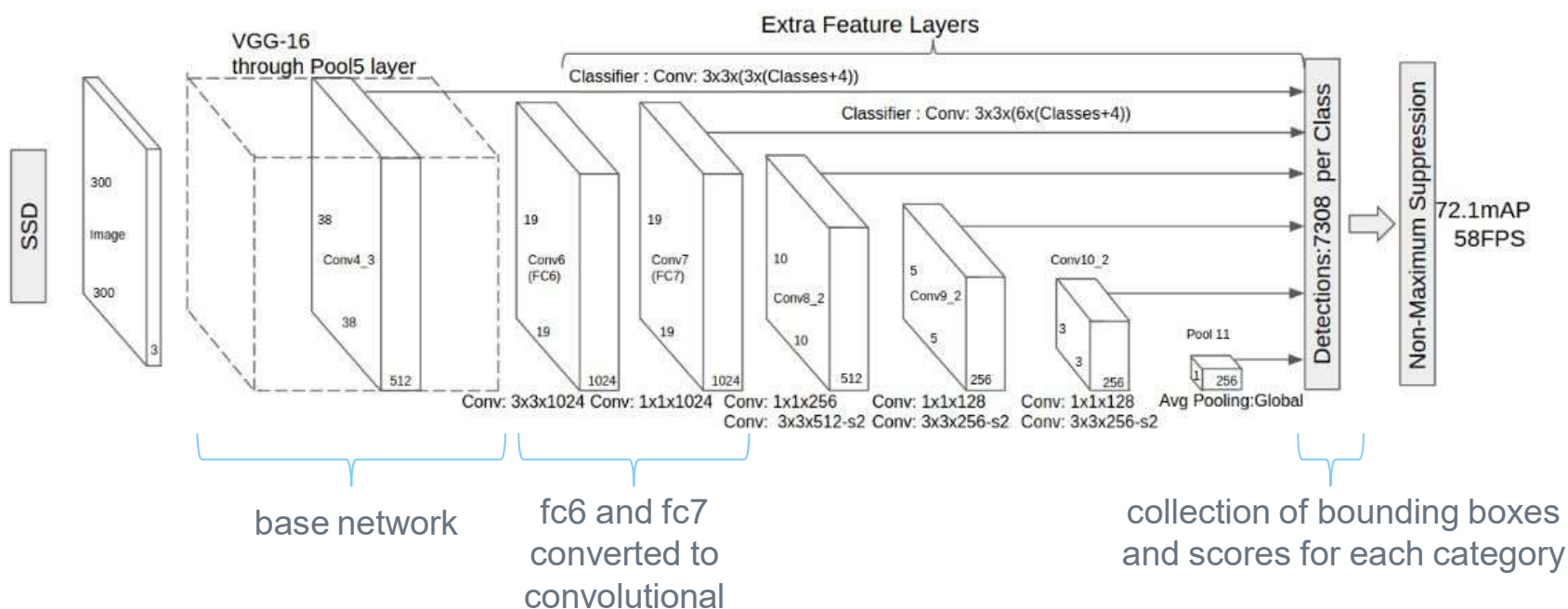
The Single Shot Detector (SSD)

Model

The Single Shot Detector (SSD)



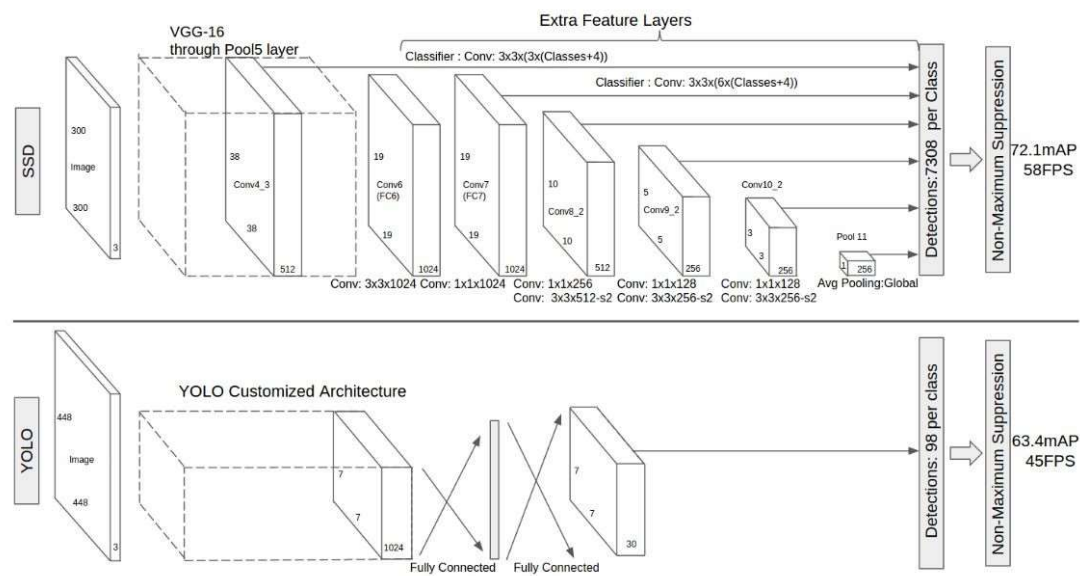
The Single Shot Detector (SSD)



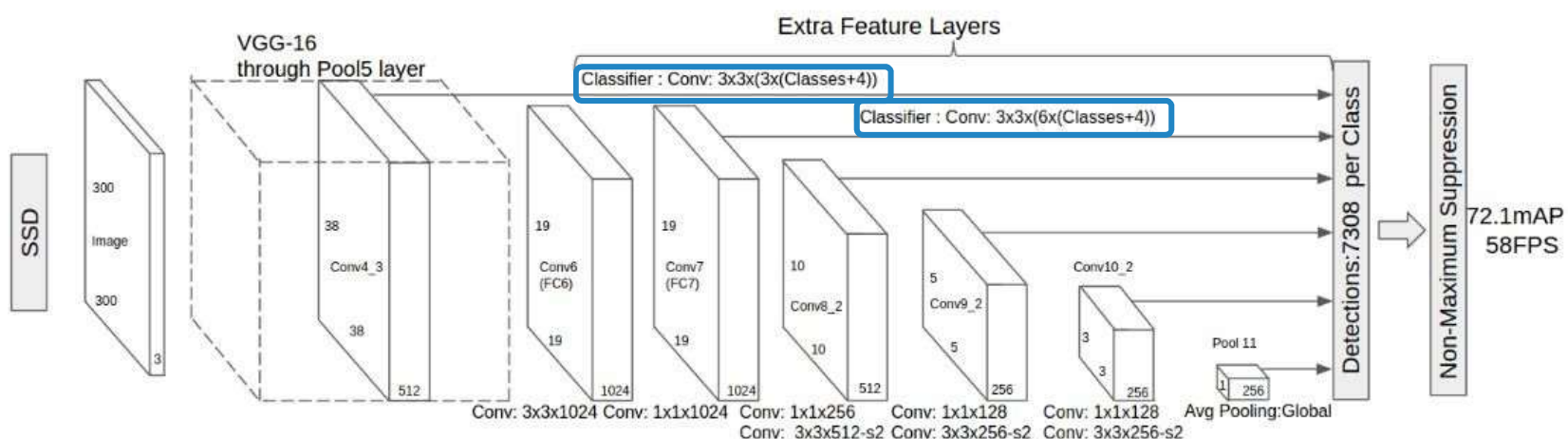
Multi-scale feature maps for detection: observe how conv feature maps decrease in size and allow predictions at multiple scales

The Single Shot Detector (SSD)

Comparison to YOLO



The Single Shot Detector (SSD)



Convolutional predictors for detection: We apply on top of each conv feature map a set of filters that predict detections for different aspect ratios and class categories

The Single Shot Detector (SSD)

What is a detection ?



Described by **four parameters** (center bounding box x and y, width and height)

Class category

For all categories we need for a detection a total of #classes + 4 values



The Single Shot Detector (SSD)

Detector for SSD:

Each detector will output a single value, so we need
(classes + 4) detectors for a detection



The Single Shot Detector (SSD)

Detector for SSD:

Each detector will output a single value, so we need
(classes + 4) detectors for a detection

BUT there are different types of detections!

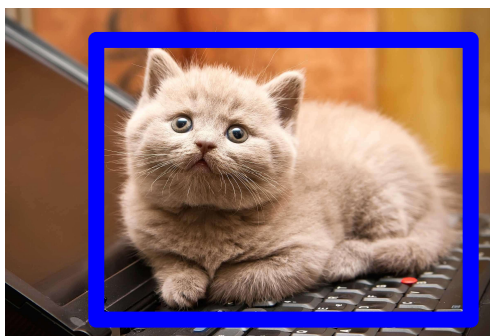


The Single Shot Detector (SSD)

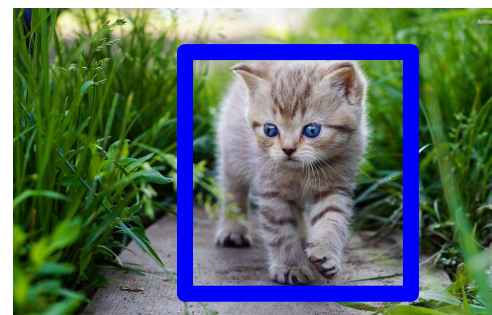
Different “classes” of detections



aspect ratio 2:1
for cats



aspect ratio 1:2
for cats



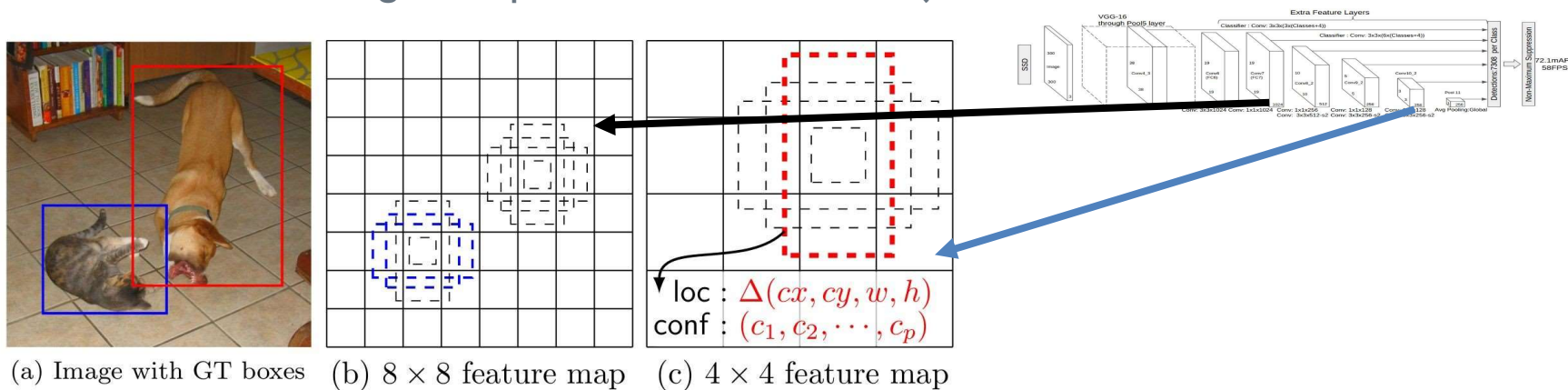
aspect ratio 1:1
for cats



The Single Shot Detector (SSD)

Choosing scales and aspect ratios for default boxes:

- ▷ Feature maps from different layers are used to handle scale variance
- ▷ Specific feature map locations learn to be responsive to specific areas of the image and particular scales of objects



(a) Image with GT boxes

(b) 8×8 feature map


(c) 4×4 feature map

The Single Shot Detector (SSD)

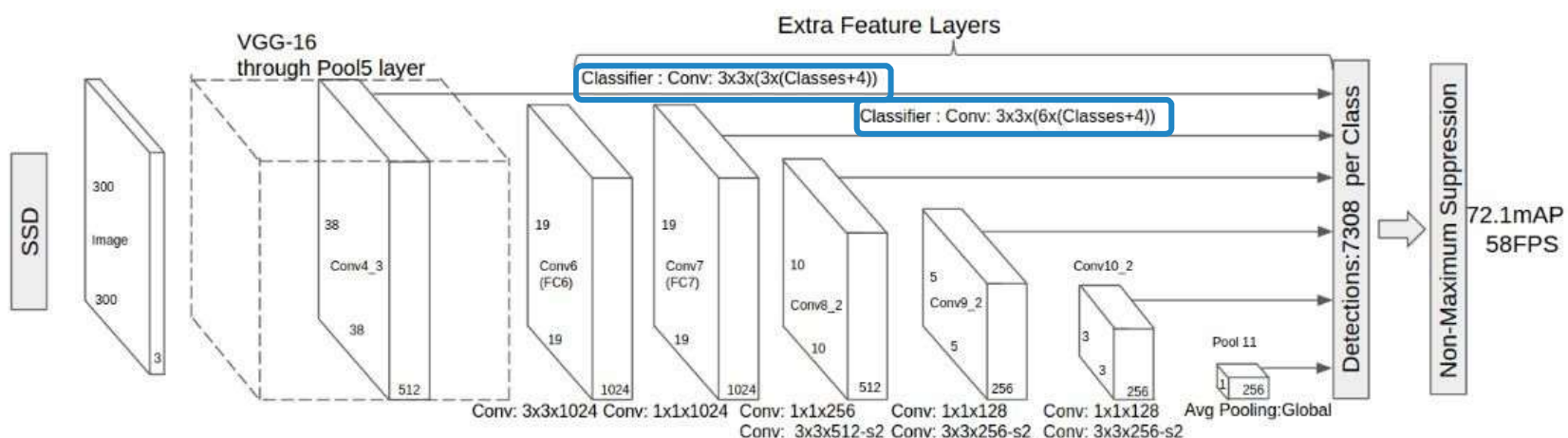
Detector for SSD:

Each detector will output a single value, so we need
(classes + 4) detectors for a detection

as we have **#default boxes**, we need
(classes + 4) x #default boxes detectors



The Single Shot Detector (SSD)



Convolutional predictors for detection: We apply on top of each conv feature map a set of filters that predict detections for different aspect ratios and class categories

The Single Shot Detector (SSD)

For each feature layer of $m \times n$ with p channels we apply kernels of $3 \times 3 \times p$ to produce either a **score for a category**, or a **shape offset** relative to a default bounding box coordinates

So, for each conv layer considered, there are

(classes + 4) x default boxes x m x n
outputs



3.2

The Single Shot Detector (SSD)

Training

The Single Shot Detector (SSD)

SSD requires that ground-truth data is **assigned** to specific outputs in the fixed set of detector outputs



The Single Shot Detector (SSD)

Matching strategy:

For each ground truth box we have to select from **all the default boxes** the ones that best fit in terms of **location, aspect ratio and scale**.

- ▷ We select the default box with **best jaccard overlap**.
- ▷ Default boxes with a jaccard overlap **higher than 0.5** are also selected



The Single Shot Detector (SSD)

Training objective:

Similar to MultiBox but handles multiple categories.

$$L(x, c, l, g) = \frac{1}{N} (L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

confidence loss
softmax loss

localization loss
Smooth L1 loss

is 1 by
cross-validation

N: number of default matched BBs
x: is 1 if the default box is matched to a
determined ground truth box, and 0
otherwise
l: predicted bb parameters
g: ground truth bb parameters
c: class

The Single Shot Detector (SSD)

Hard negative mining:

Significant imbalance between **positive** and **negative** training examples

- ▷ Use negative samples with **higher confidence score**
- ▷ Then the ratio of positive-negative samples is **3:1**



The Single Shot Detector (SSD)

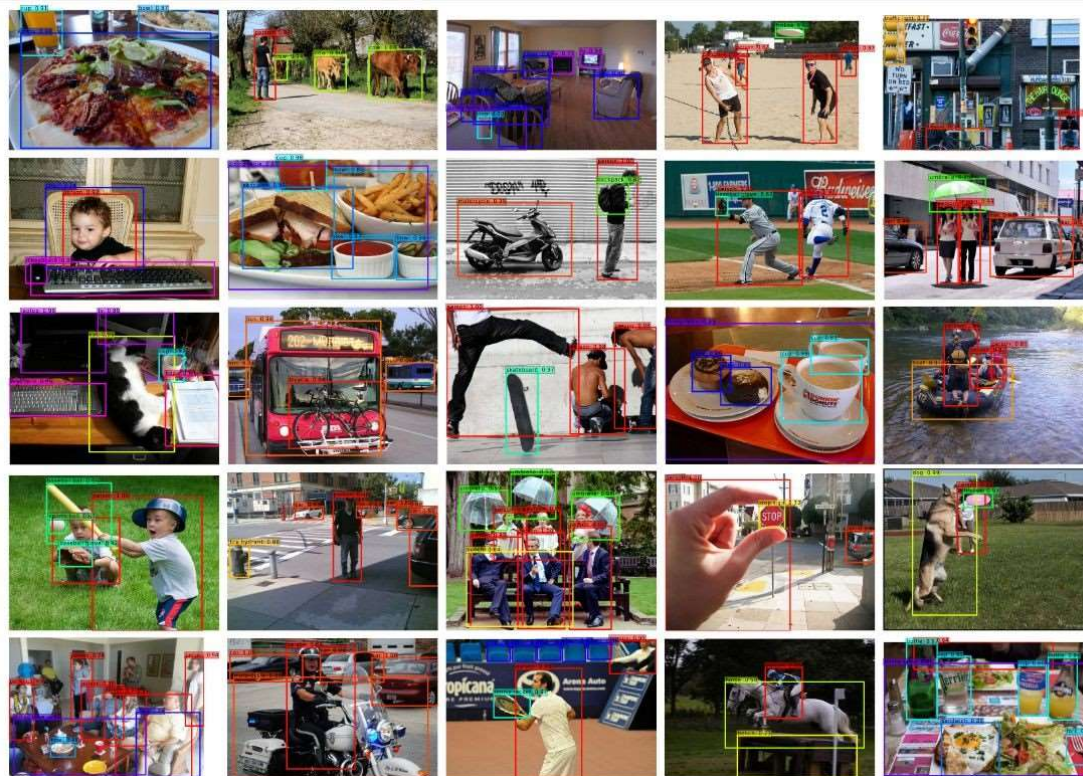
Data augmentation:

Each training sample is randomly sampled by one of the following options:

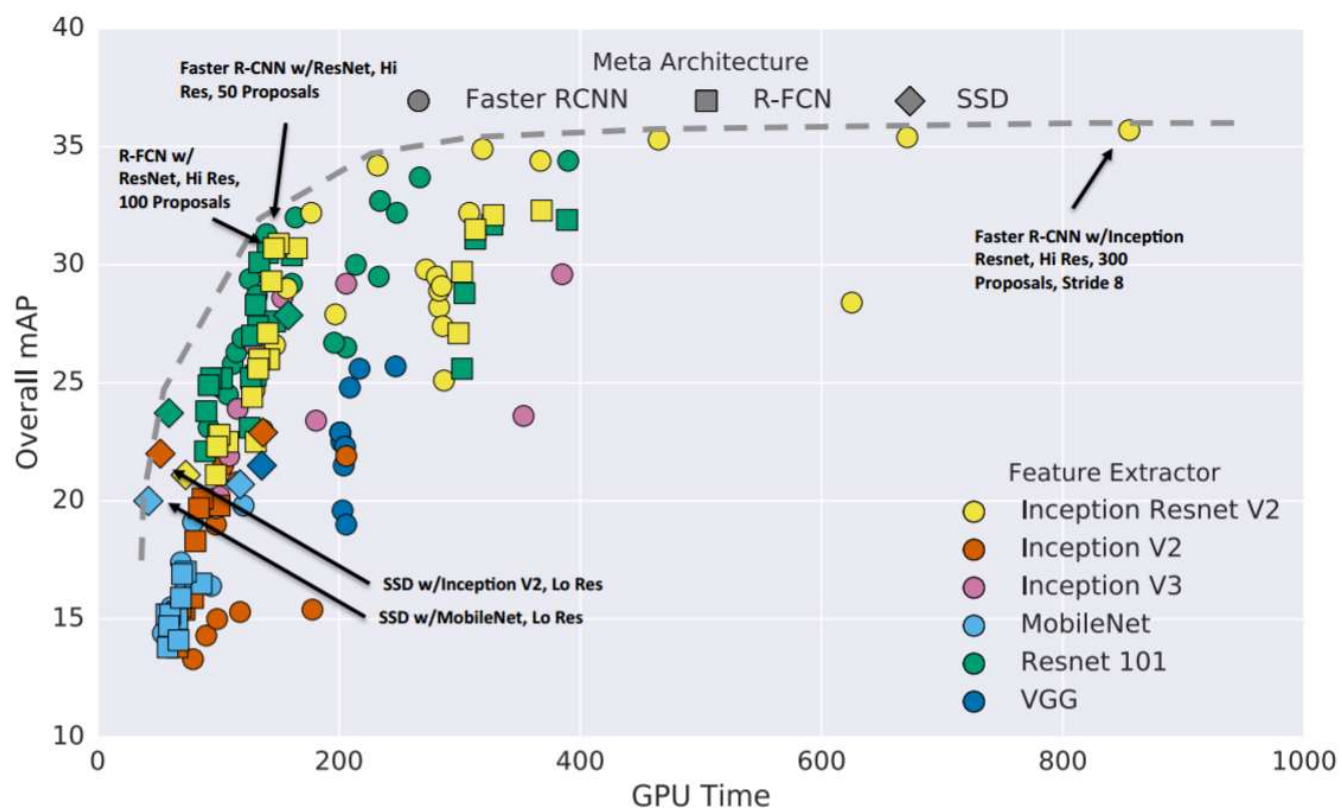
- ▷ Use the original image
- ▷ Sample a path with a minimum jaccard overlap with objects



Visualizations



Google test on TF



5.

Conclusions

SSD: Single Shot MultiBox Detector

Conclusions

- ▷ **Single-shot object** detector for multiple categories
- ▷ One key feature is to use **multiple convolutional maps** to deal with different scales
- ▷ **More default bounding boxes**, the better results obtained
- ▷ Comparable accuracy to state-of-the-art object detectors, but **much faster**



Thank you for your
attention! Questions?