



EVROPSKÁ UNIE
Evropské strukturální a investiční fondy
Operační program Výzkum, vývoj a vzdělávání

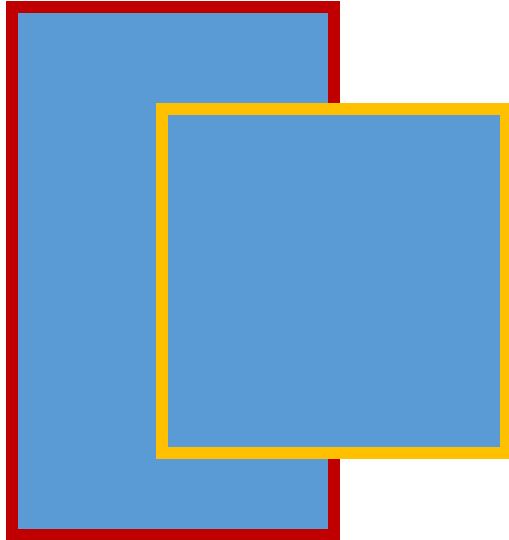


Image Analysis II

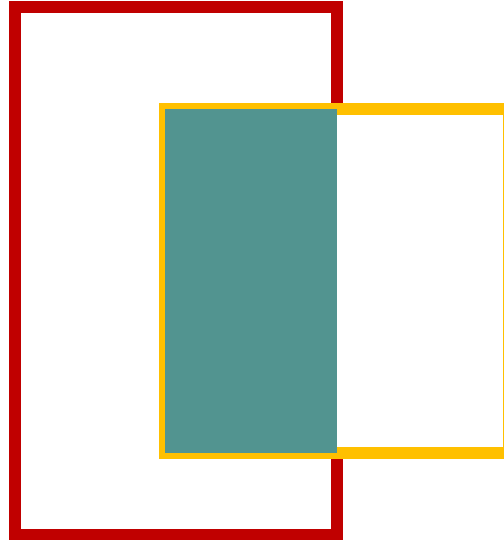
Radovan Fusek

Intersection, Union

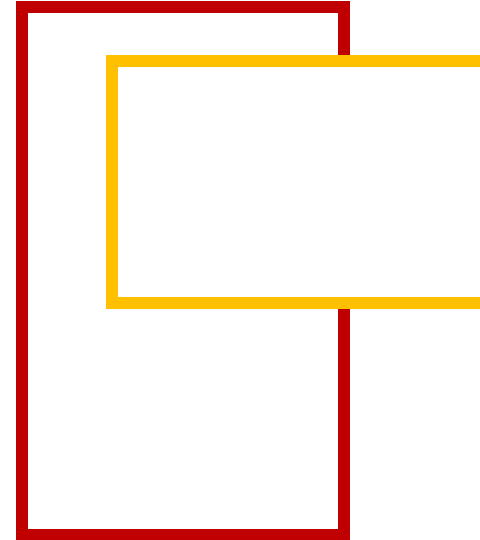
Union



Intersection



IoU = ?



IoU = ?

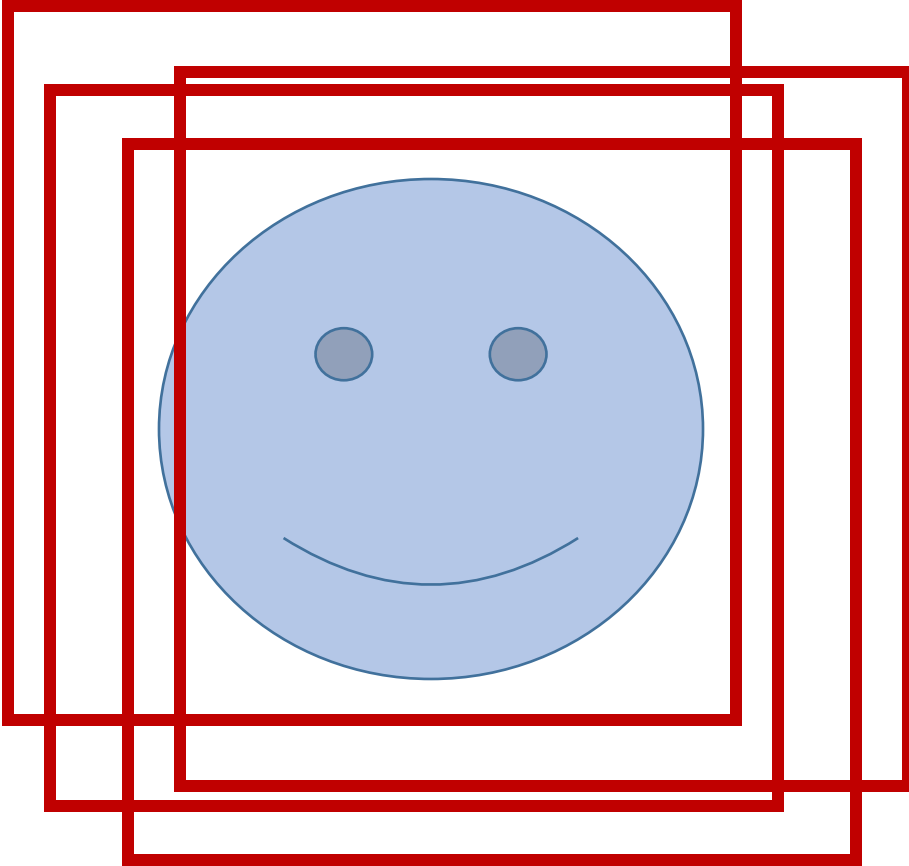


$$\text{IoU} = \text{Area of Overlap} / \text{Area of Union}$$

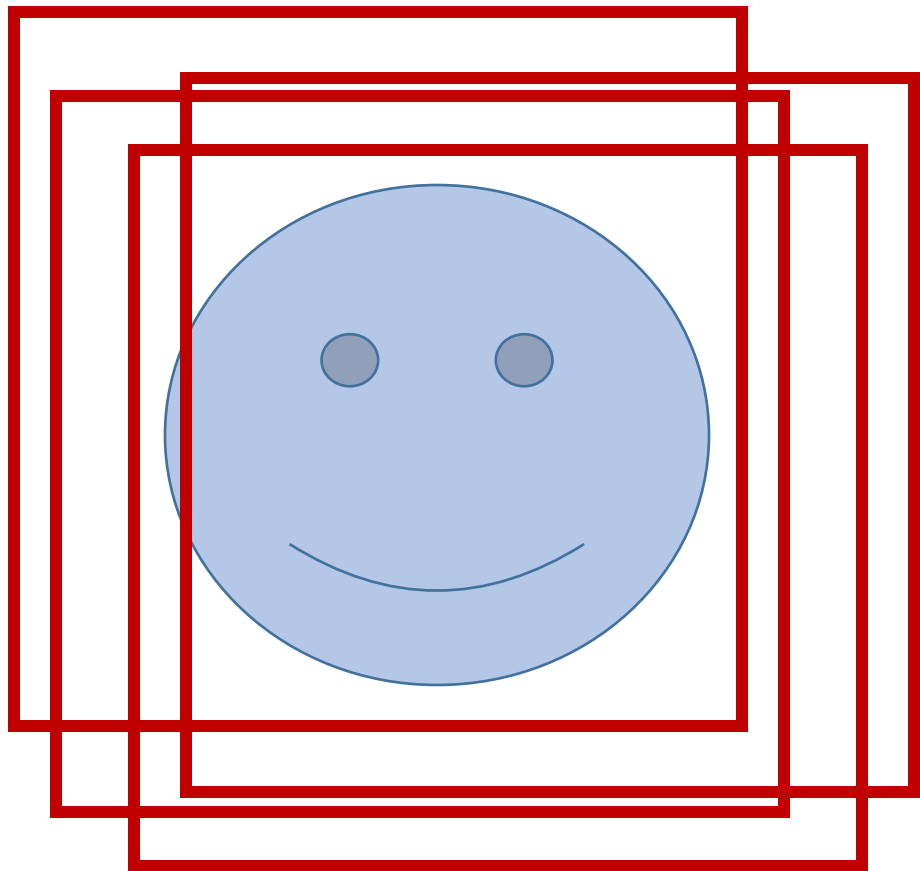


Non-max Suppression

How select only one box?



Non-max Suppression



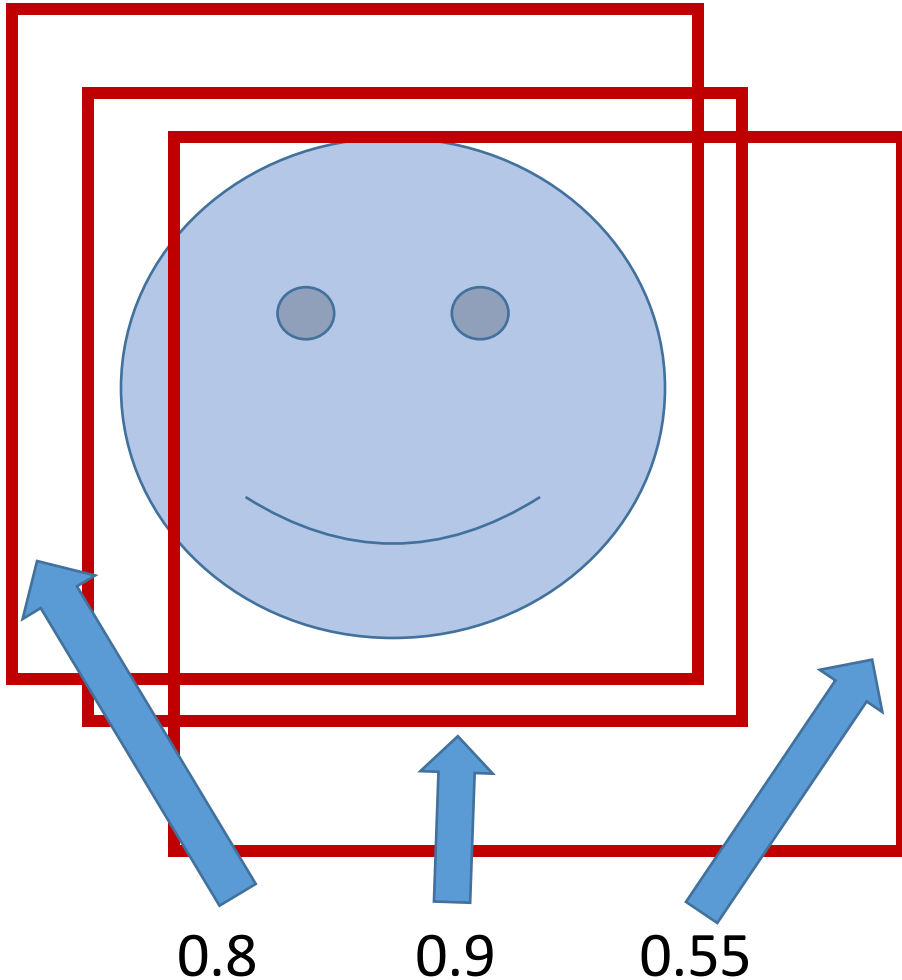
How select only one box?

1. Discard all boxes with confidence smaller or equal to 0.6
2. Select the box with largest confidence
3. Discard all remaining box with IoU greater or equal to 0.5

Non-max Suppression

How select only one box?

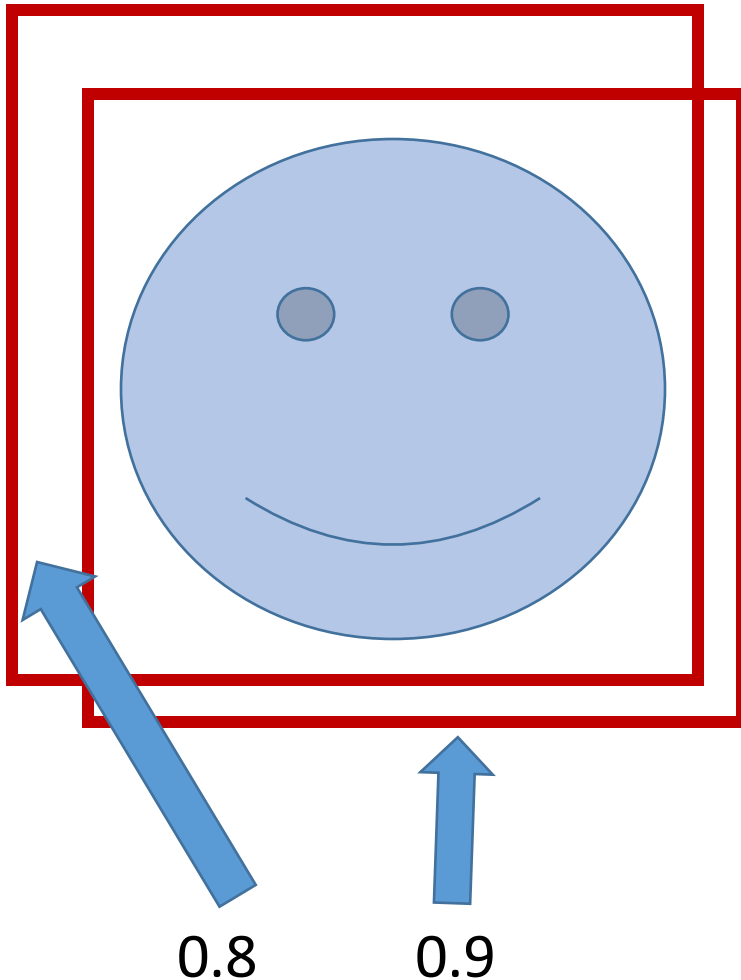
1. Discard all boxes with confidence smaller or equal to 0.6



Non-max Suppression

How select only one box?

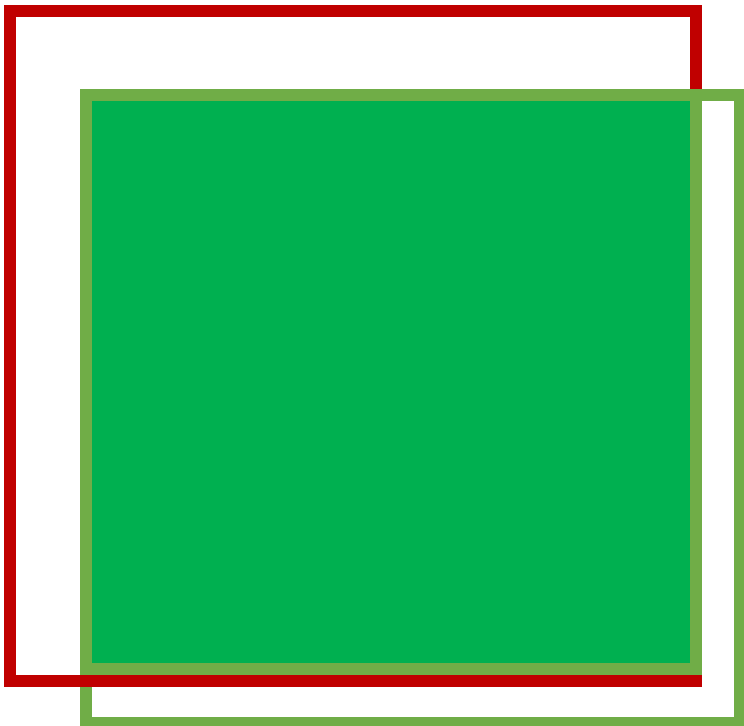
1. Discard all boxes with confidence smaller or equal to 0.6
2. Select the box with largest confidence



Non-max Suppression

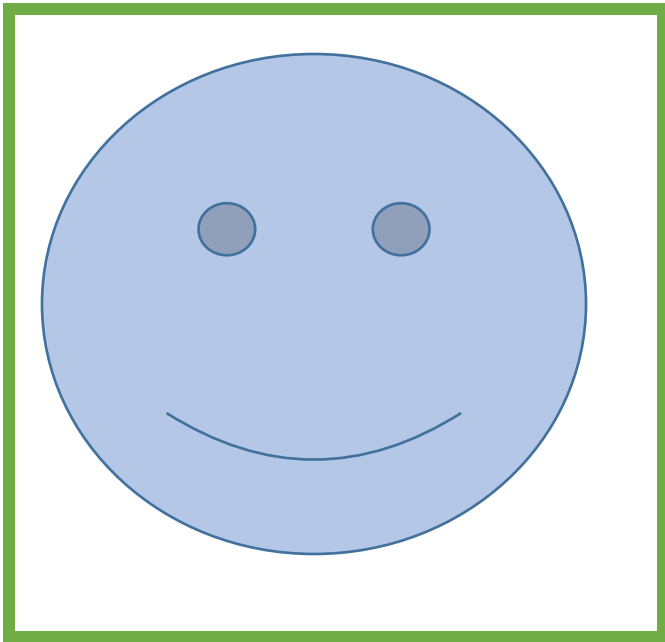
How select only one box?

1. Discard all boxes with confidence smaller or equal to 0.6
2. Select the box with largest confidence
3. Discard all remaining box with IoU greater or equal to 0.5



Non-max Suppression

How select only one box?



1. Discard all boxes with confidence smaller or equal to 0.6
2. Select the box with largest confidence
3. Discard all remaining box with IoU greater or equal to 0.5



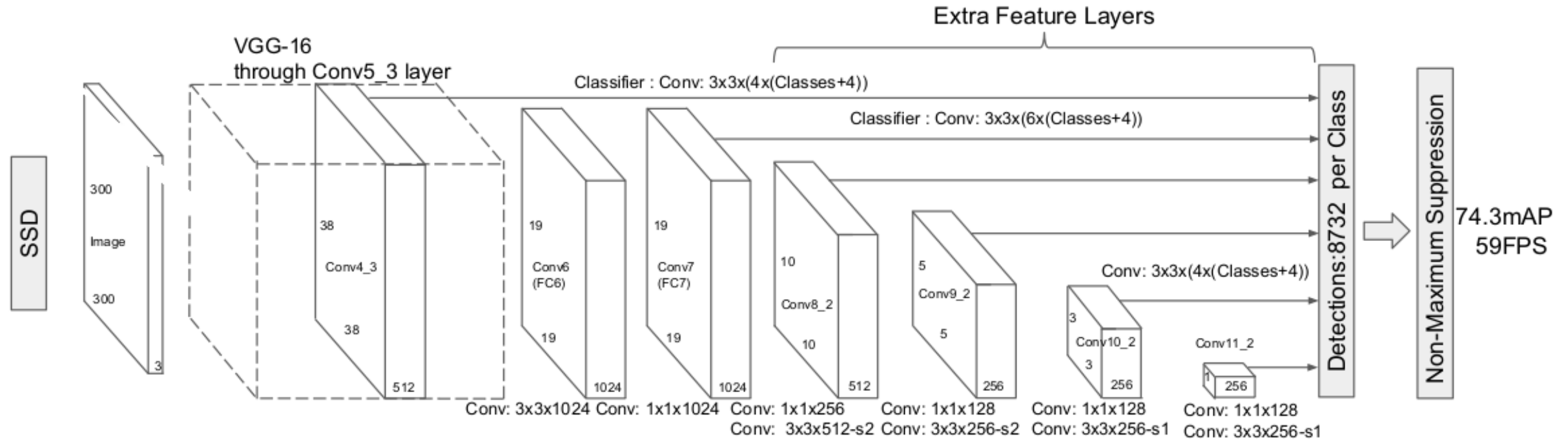
SSD: Single Shot MultiBox Detector

Wei Liu¹, Dragomir Anguelov², Dumitru Erhan³, Christian Szegedy³,
Scott Reed⁴, Cheng-Yang Fu¹, Alexander C. Berg¹

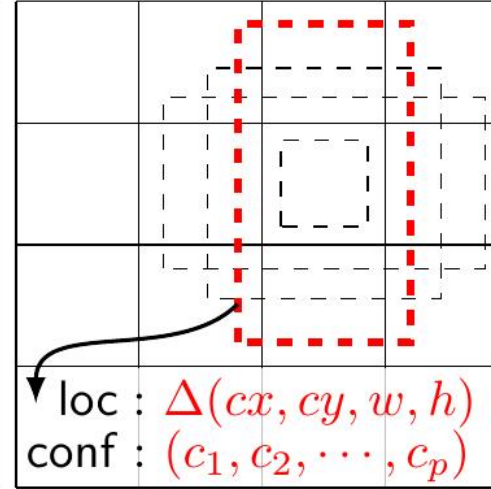
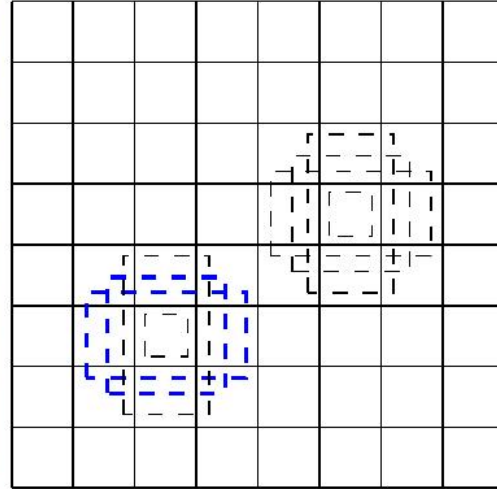
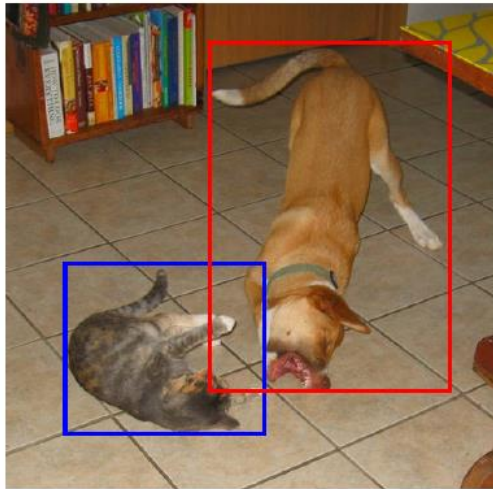
¹UNC Chapel Hill ²Zoox Inc. ³Google Inc. ⁴University of Michigan, Ann-Arbor
¹wliu@cs.unc.edu, ²drago@zoox.com, ³{dumitru,szegedy}@google.com,
⁴reedscot@umich.edu, ¹{cyfu,aberg}@cs.unc.edu

- We introduce SSD, a single-shot detector for multiple categories that is faster than the previous state-of-the-art for single shot detectors (YOLO), and significantly more accurate, in fact as accurate as slower techniques that perform explicit region proposals and pooling (including Faster R-CNN).
- The core of SSD is predicting category scores and box offsets for a fixed set of default bounding boxes using small convolutional filters applied to feature maps.
- To achieve high detection accuracy we produce predictions of different scales from feature maps of different scales, and explicitly separate predictions by aspect ratio.
- These design features lead to simple end-to-end training and high accuracy, even on low resolution input images, further improving the speed vs accuracy trade-off.
- Experiments include timing and accuracy analysis on models with varying input size evaluated on PASCAL VOC, COCO, and ILSVRC and are compared to a range of recent state-of-the-art approaches.

SSD



Anchor Boxes



(a) Image with GT boxes (b) 8×8 feature map (c) 4×4 feature map

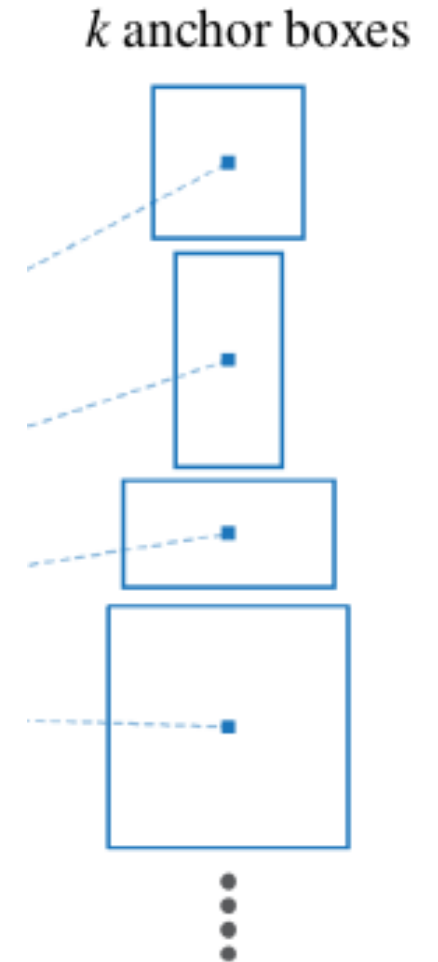
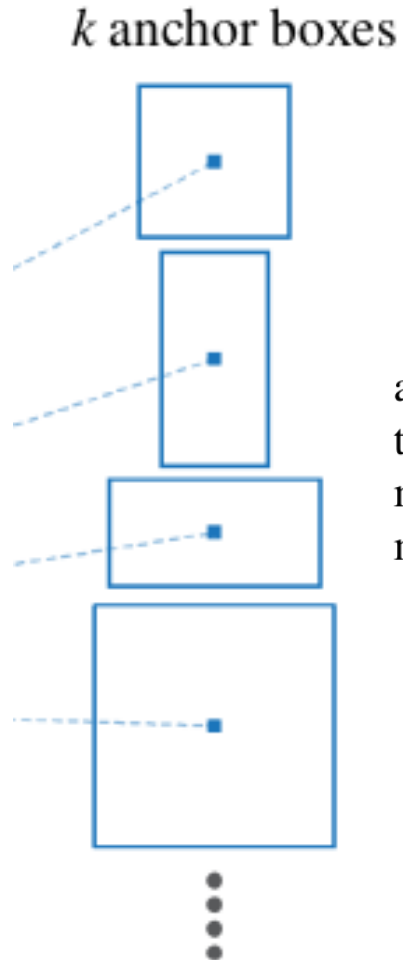


Fig. 1: SSD framework. (a) SSD only needs an input image and ground truth boxes for each object during training. In a convolutional fashion, we evaluate a small set (e.g. 4) of default boxes of different aspect ratios at each location in several feature maps with different scales (e.g. 8×8 and 4×4 in (b) and (c)). For each default box, we predict both the shape offsets and the confidences for all object categories $((c_1, c_2, \dots, c_p))$. At training time, we first match these default boxes to the ground truth boxes. For example, we have matched two default boxes with the cat and one with the dog, which are treated as positives and the rest as negatives. The model loss is a weighted sum between localization loss (e.g. Smooth L1 [6]) and confidence loss (e.g. Softmax).

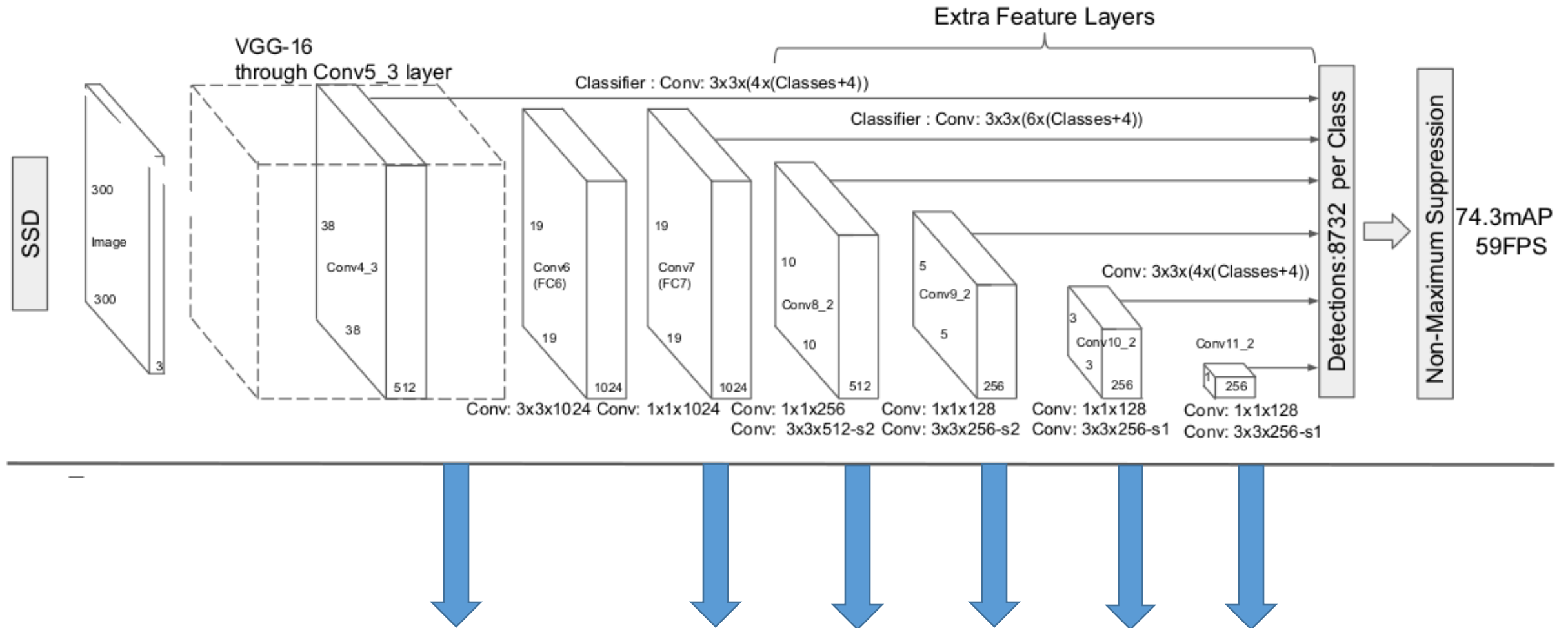
Anchor Boxes

Face
Pedestrian
Car
...
..
..



an illustration of default boxes, please refer to Fig. 1. Our default boxes are similar to the *anchor boxes* used in Faster R-CNN [2], however we apply them to several feature maps of different resolutions. Allowing different default box shapes in several feature maps let us efficiently discretize the space of possible output box shapes.

Sizes can be obtained from dataset



cnn layers make predictions, each with different anchor boxes

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	$\sim 1000 \times 600$
Fast YOLO	52.7	155	1	98	448×448
YOLO (VGG16)	66.4	21	1	98	448×448
SSD300	74.3	46	1	8732	300×300
SSD512	76.8	19	1	24564	512×512
SSD300	74.3	59	8	8732	300×300
SSD512	76.8	22	8	24564	512×512

Table 7: Results on Pascal VOC2007 test. SSD300 is the only real-time detection method that can achieve above 70% mAP. By using a larger input image, SSD512 outperforms all methods on accuracy while maintaining a close to real-time speed.



SSD

