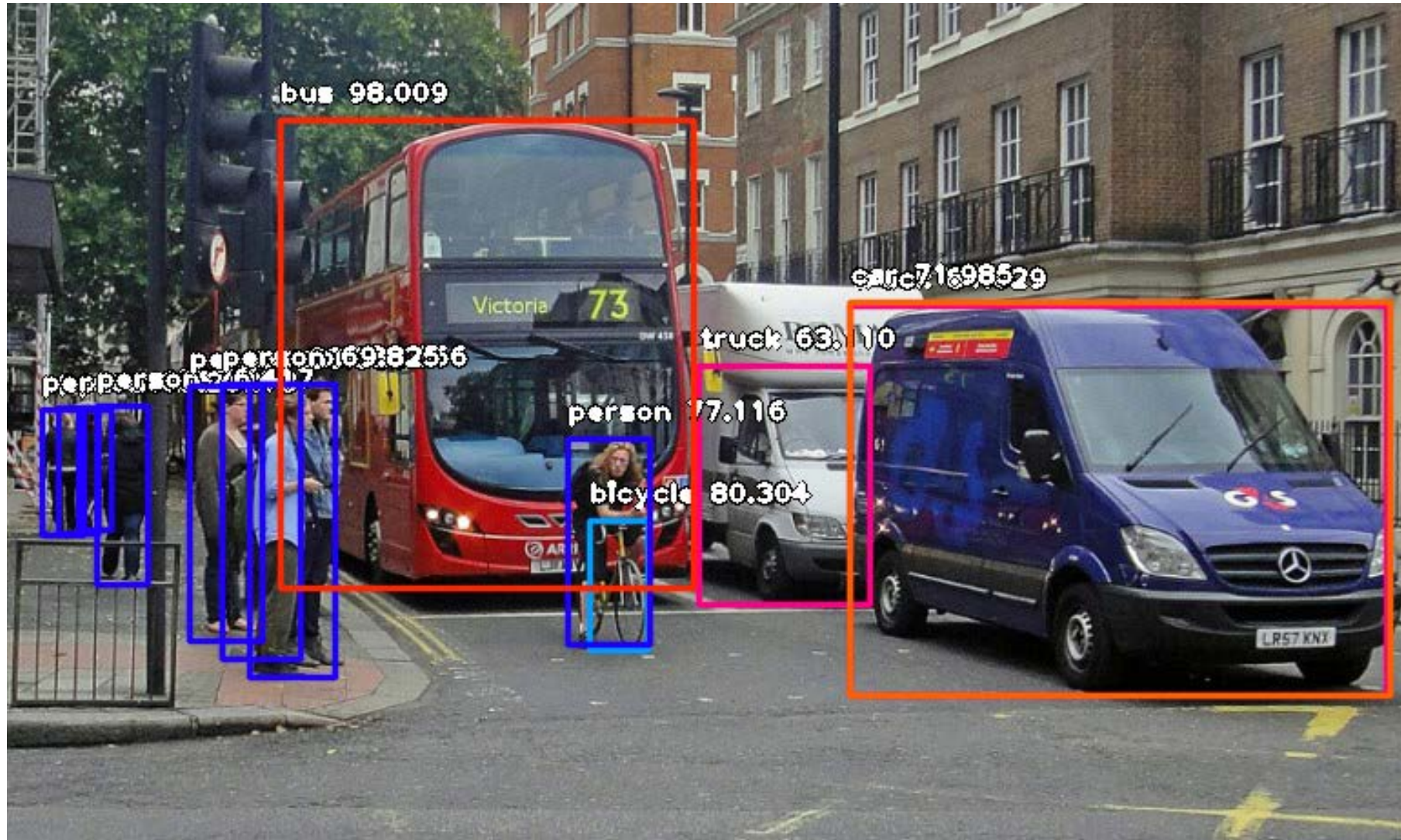# SSD : Single Shot MultiBox Detector
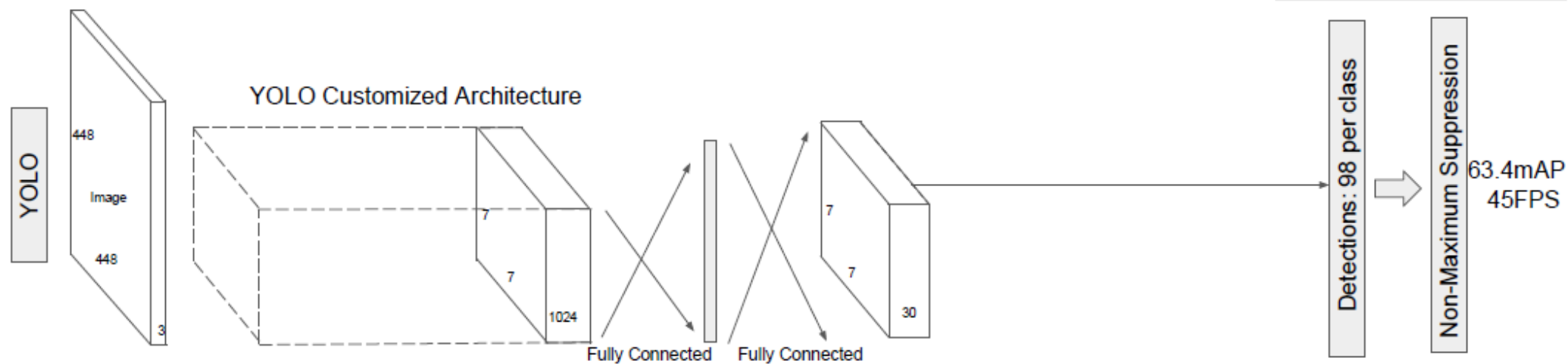
## ECCV 2016

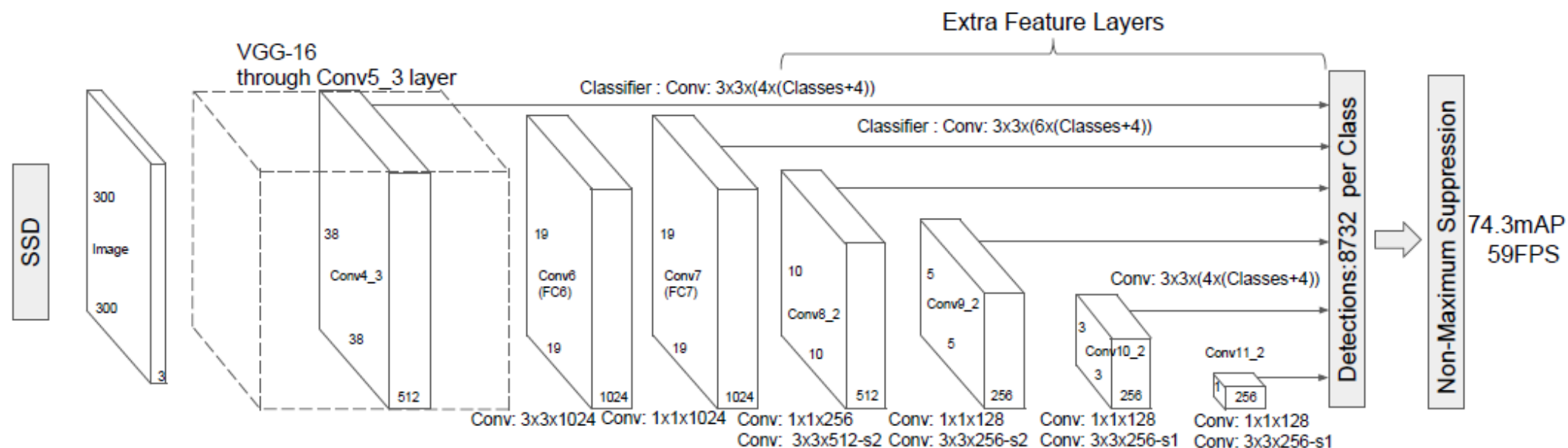# Main Task : Object Detection



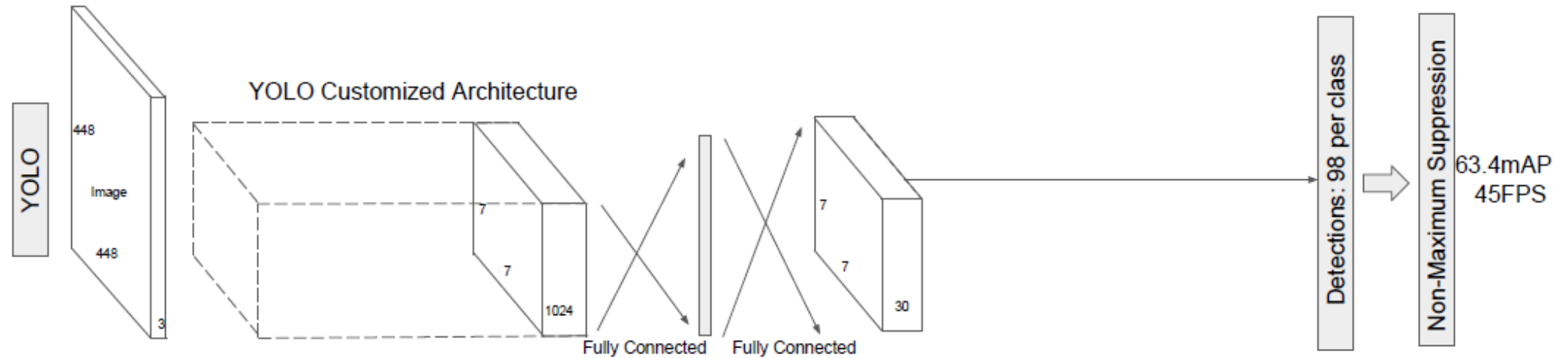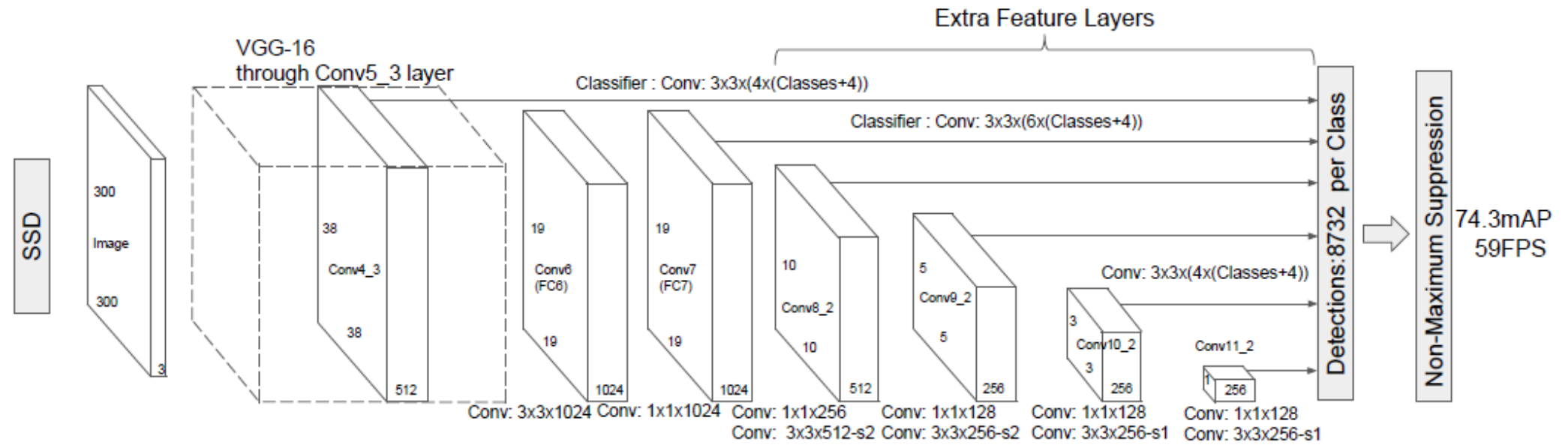Localization + Classification

# Main Problem

- But, conventional object detection network has several problems

- R-CNN series
  - Good performance
  - Too slow

- Yolo
  - So fast
  - Low performance

# Related researches



SSD — VGG-16 through Conv5_3 layer, Extra Feature Layers. Classifier: Conv: 3x3x(4x(Classes+4)), Conv: 3x3x(6x(Classes+4)), Conv: 3x3x(4x(Classes+4)). Conv4_3 38x38x512, Conv6 (FC6) 19x19x1024, Conv7 (FC7) 19x19x1024, Conv8_2 10x10x512, Conv9_2 5x5x256, Conv10_2 3x3x256, Conv11_2 256. Detections: 8732 per Class → Non-Maximum Suppression → 74.3mAP 59FPS. Conv: 3x3x1024, Conv: 1x1x1024, Conv: 1x1x256, Conv: 3x3x512-s2, Conv: 1x1x128, Conv: 3x3x256-s2, Conv: 1x1x128, Conv: 3x3x256-s1, Conv: 1x1x128, Conv: 3x3x256-s1.

YOLO — YOLO Customized Architecture. Image 448x448x3, 7x7x1024, Fully Connected, Fully Connected, 7x7x30. Detections: 98 per class → Non-Maximum Suppression → 63.4mAP 45FPS.
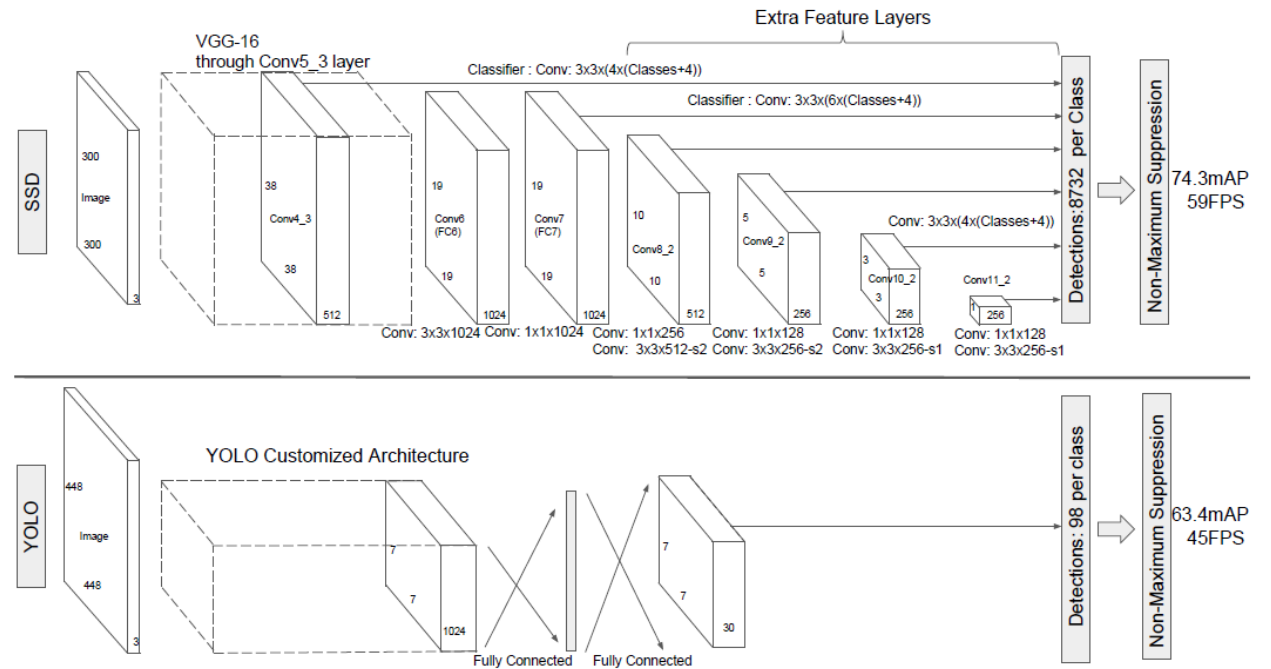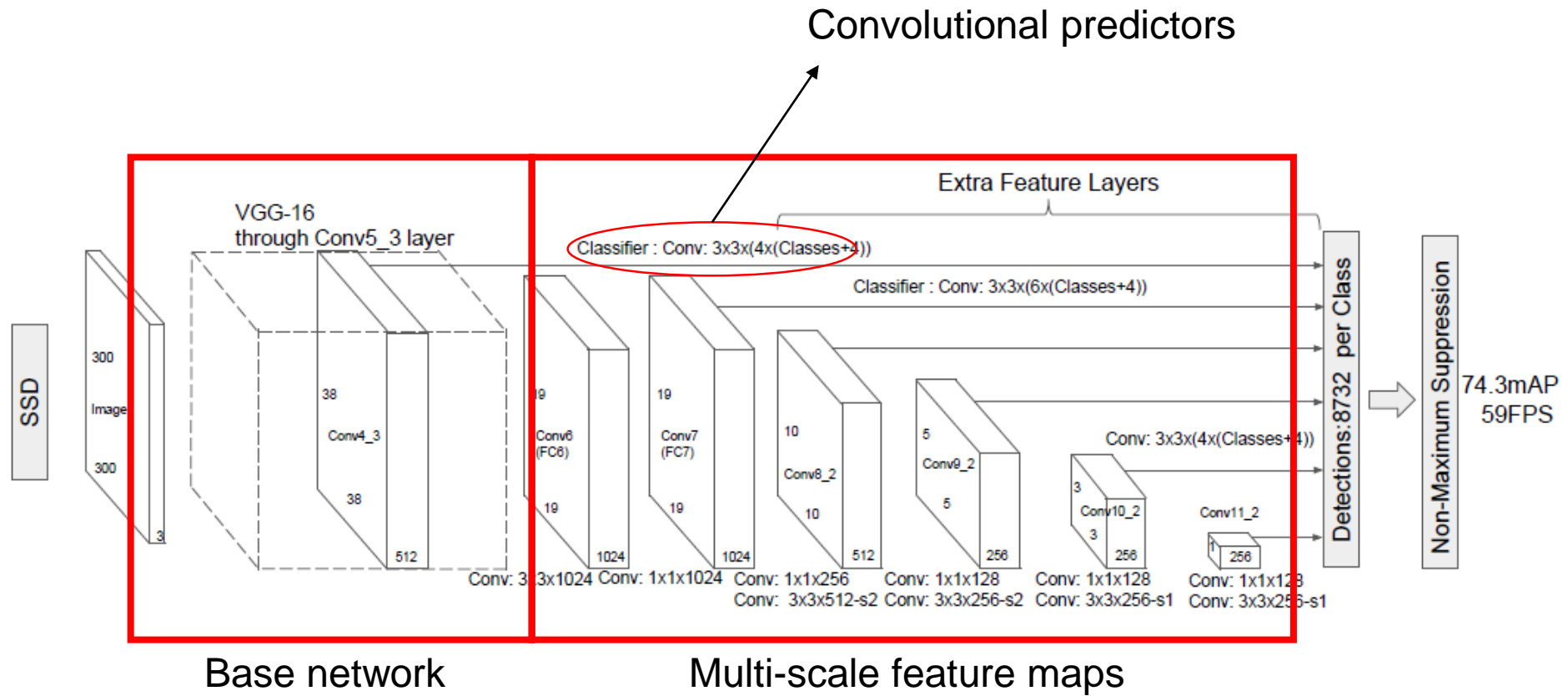
# Solution : SSD architecture
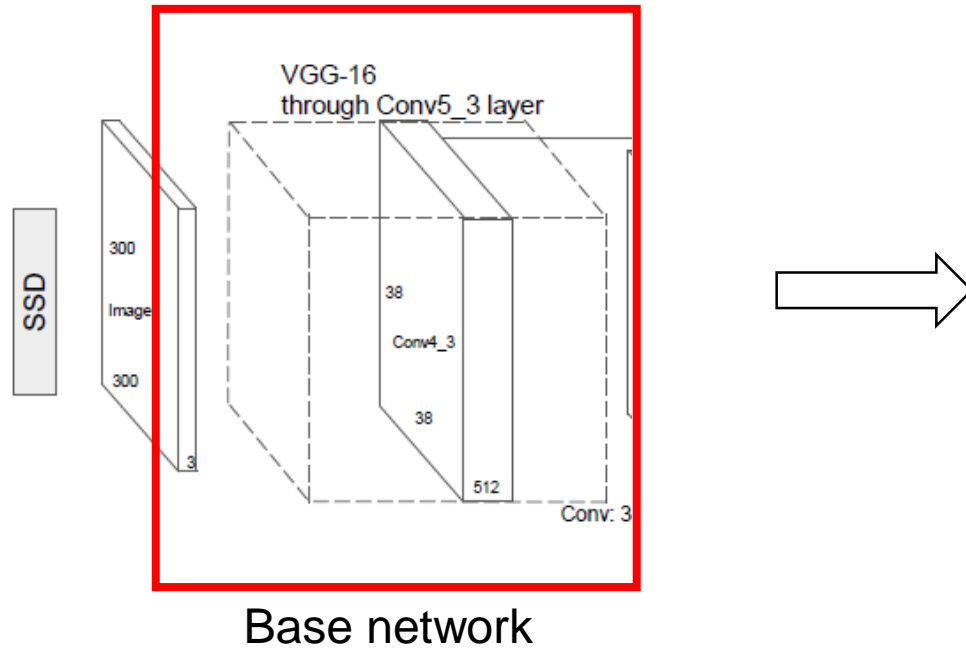
# SSD : Main Contributions

- Propose faster than YOLO, and more accurate, comparable with R-CNN series

- Use box offsets for more faster&accurate prediction

- Use multi-scale prediction scheme

# SSD : Architecture

# SSD : Architecture



VGG-16
through Conv5_3 layer

SSD

300

Image

300

3

38

Conv4_3

38

512

Conv: 3

Base network

- Backbone : VGG16

- Used for high-quality image classification
  - Learned common features of images

- Truncated at Conv5_3 layer
  - 300x300 → 38x38

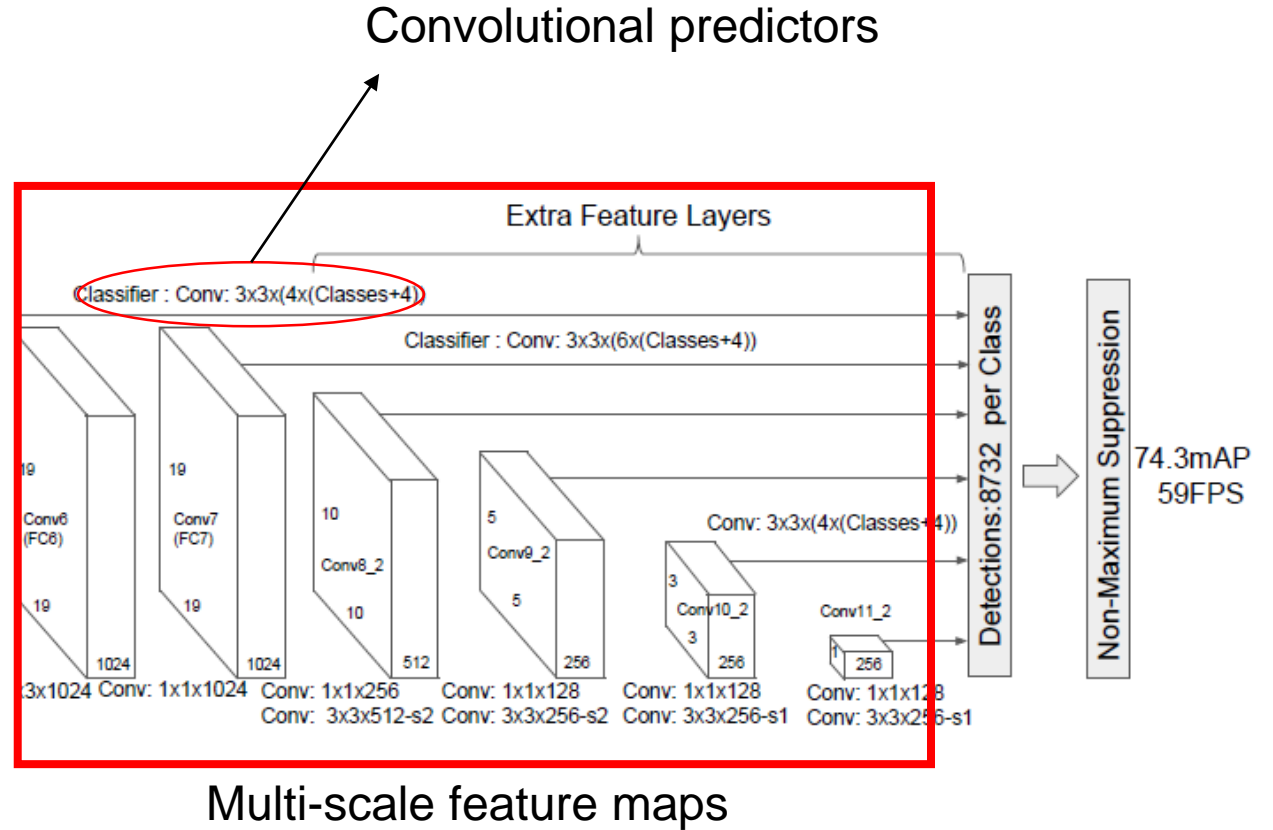# SSD : Architecture

<Multi-scale feature maps>

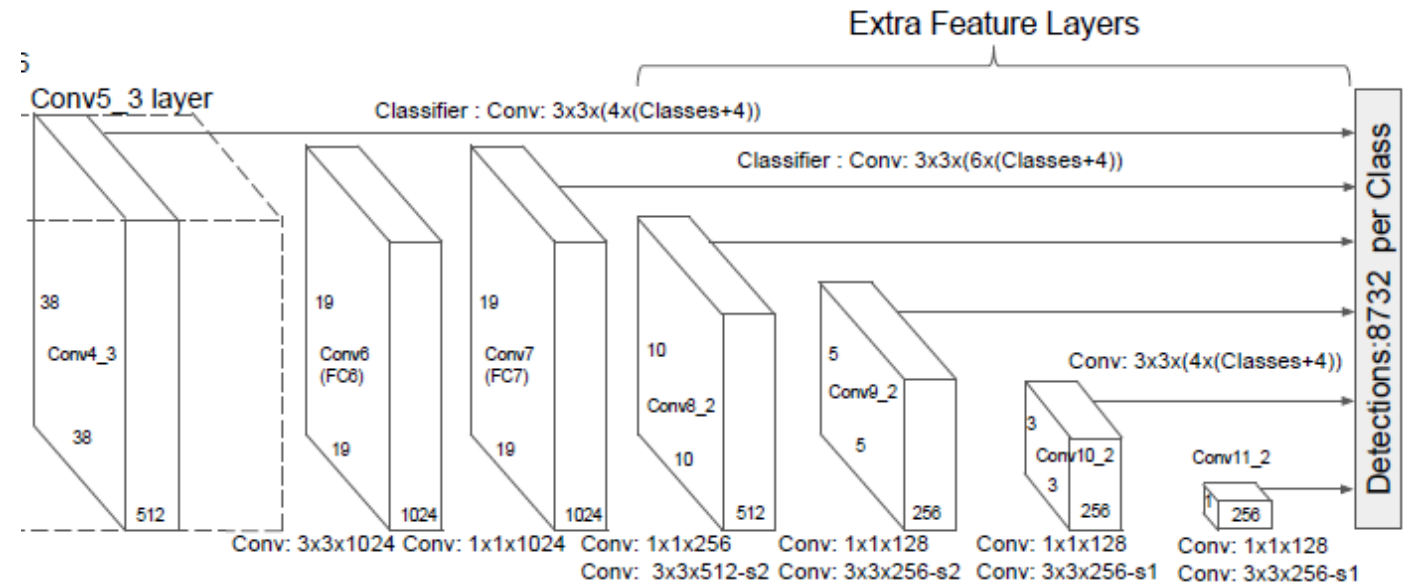- Added auxiliary convolutional layers for Multi-scale detections

<Convolutional predictors>

- Produce a fixed set of detection predictions by using CNN
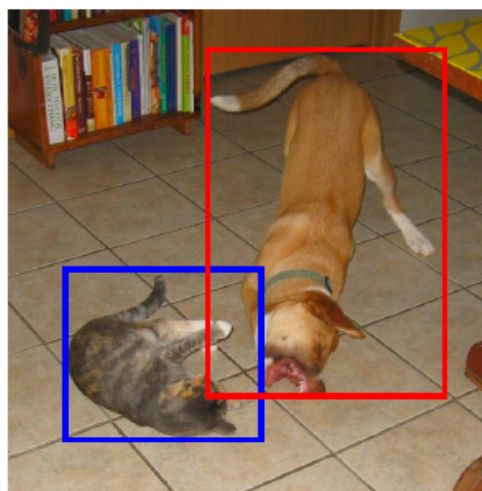
- Applied to feature maps for each scales

Convolutional predictors



Multi-scale feature maps

# SSD : Architecture

- Produce fixed number of predictions for each feature map

- If we predict k bounding box region, we should use (C+4)k channel CNN kernel
    - For mxn feature map, produce (C+4)kmn predictions

- Predict the offsets relative to the default box shapes in the cell



Extra Feature Layers

Conv5_3 layer

Classifier : Conv: 3x3x(4x(Classes+4))

Classifier : Conv: 3x3x(6x(Classes+4))

Conv: 3x3x(4x(Classes+4))

Detections:8732 per Class

38  Conv4_3  38  512

19  Conv6 (FC6)  19  1024

19  Conv7 (FC7)  19  1024

10  Conv8_2  10  512

5  Conv9_2  5  256

3  Conv10_2  3  256

Conv11_2  256

Conv: 3x3x1024  Conv: 1x1x1024  Conv: 1x1x256  Conv: 1x1x128  Conv: 1x1x128  Conv: 1x1x128
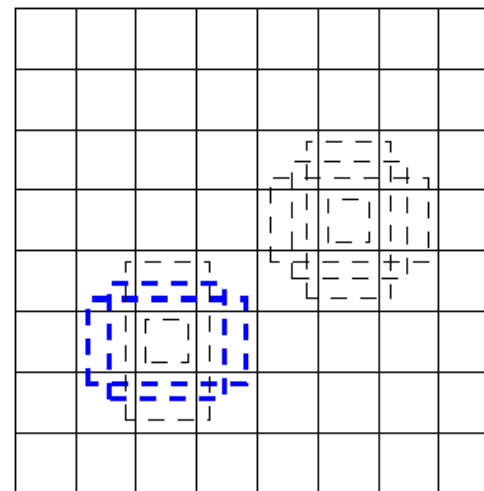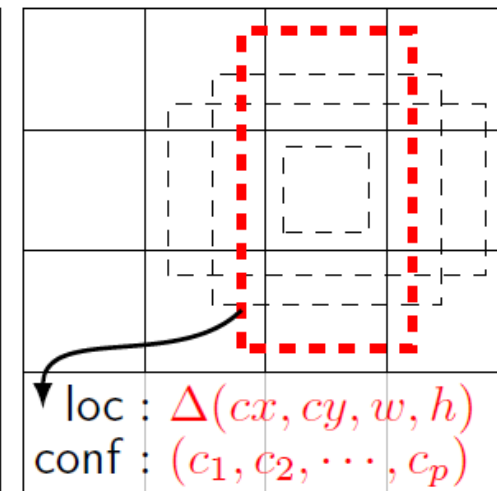Conv: 3x3x512-s2  Conv: 3x3x256-s2  Conv: 3x3x256-s1  Conv: 3x3x256-s1

# SSD : Details

- Discretize space of the possible output box shapes

- Allocated default box set to each scales

- Useful to predict bounding box



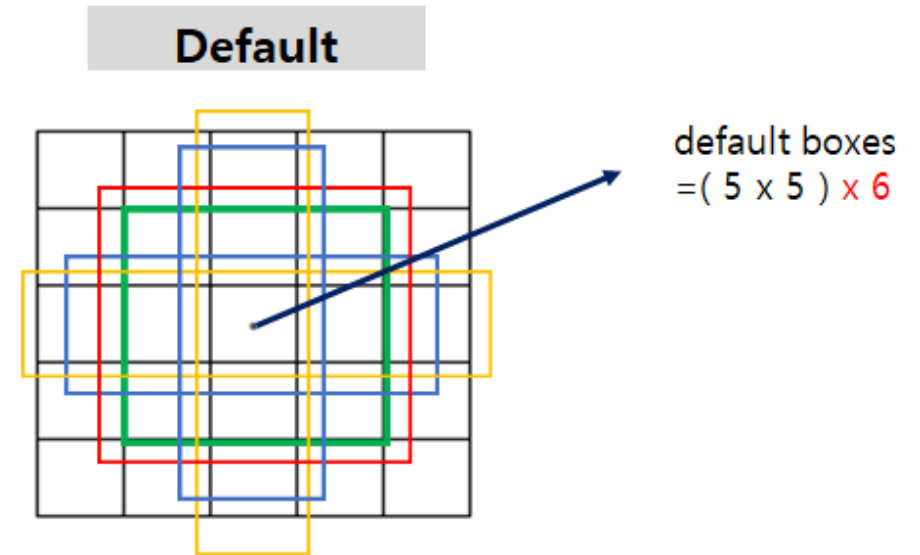(a) Image with GT boxes  (b) $8 \times 8$ feature map  (c) $4 \times 4$ feature map

loc : $\Delta(cx, cy, w, h)$
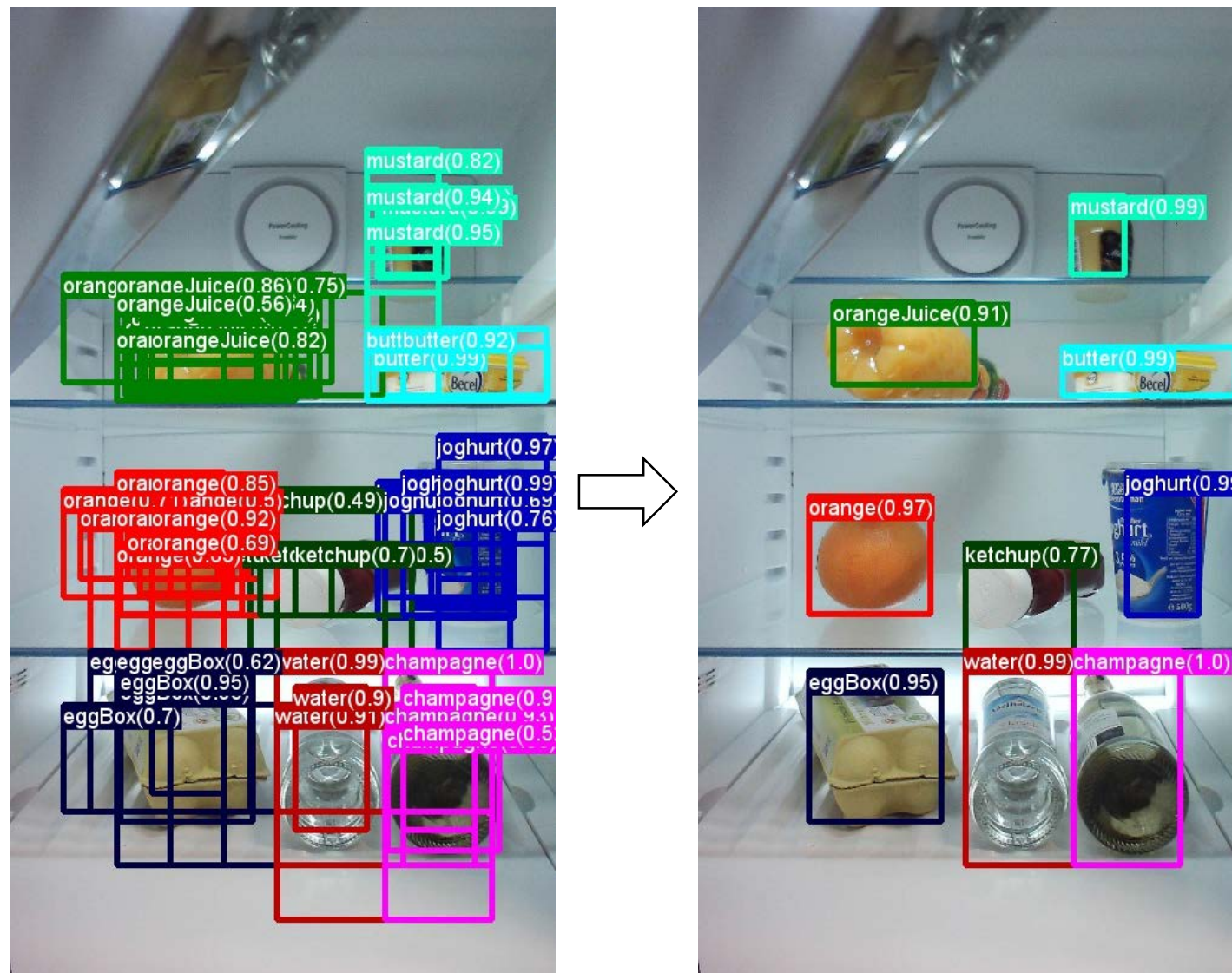conf : $(c_1, c_2, \cdots, c_p)$

# SSD : Details

- Match default boxes to ground truth and select boxes with threshold (IOU 0.5)

- Use multiple boxes, calculate loss function

- Useful to predict bounding box

**Default**

default boxes
=( 5 x 5 ) x 6

# SSD : Non-maximum suppression

- Remove overlapped bounding boxes

- Sorting list of bounding box by confidence

- Eliminate overlapped bounding box, determine by IOU threshold

# SSD : Training

- Summarized form

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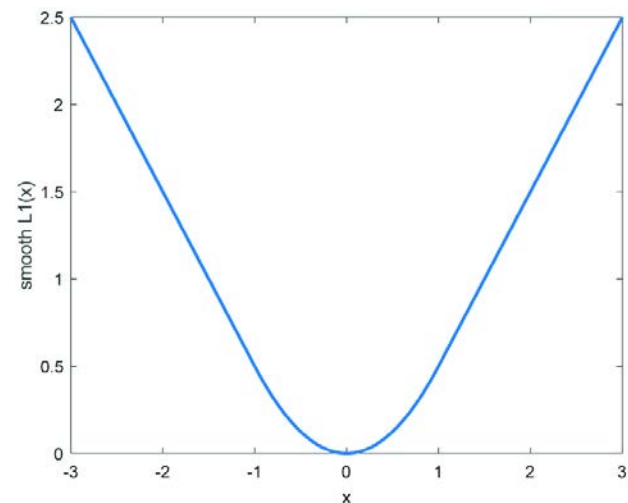
Confidence loss          Localization loss

Number of matched default boxes

# SSD : Training

- Localization loss

$$L_{loc}(x, l, g) = \sum_{i \in Pos}^{N} \sum_{m \in \{cx, cy, w, h\}} x_{ij}^k \text{smooth}_{L1}(l_i^m - \hat{g}_j^m)$$

$$\hat{g}_j^{cx} = (g_j^{cx} - d_i^{cx})/d_i^w \qquad \hat{g}_j^{cy} = (g_j^{cy} - d_i^{cy})/d_i^h$$

$$\hat{g}_j^w = \log\left(\frac{g_j^w}{d_i^w}\right) \qquad \hat{g}_j^h = \log\left(\frac{g_j^h}{d_i^h}\right)$$



- Confidence loss

Classification probability

$$L_{conf}(x, c) = -\sum_{i \in Pos}^{N} x_{ij}^p log(\hat{c}_i^p) - \sum_{i \in Neg} log(\hat{c}_i^0) \quad \text{where} \quad \hat{c}_i^p = \frac{\exp(c_i^p)}{\sum_p \exp(c_i^p)}$$
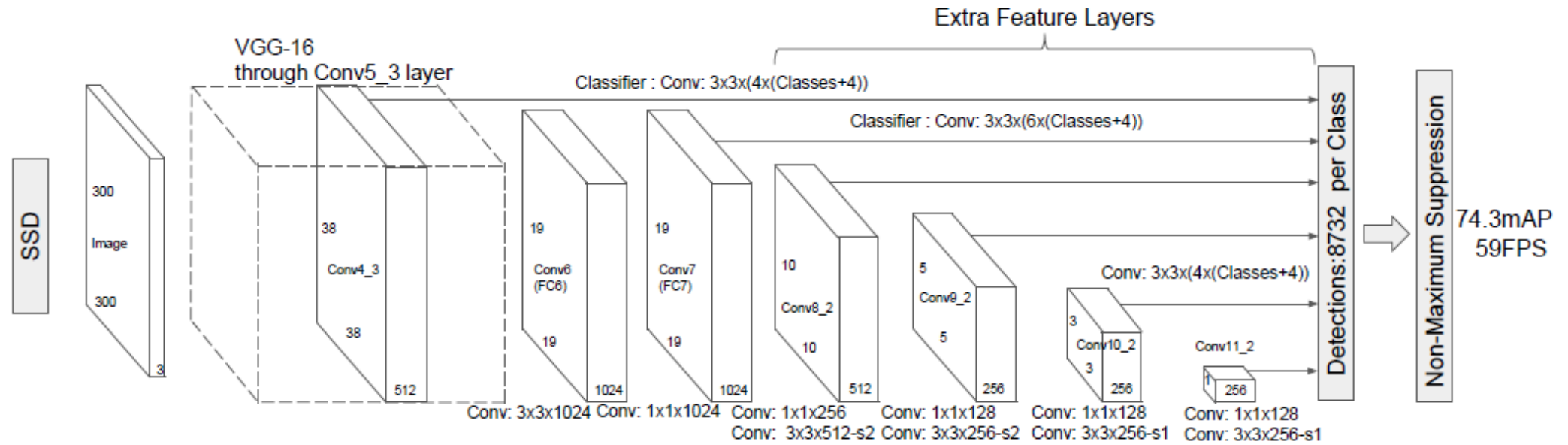
Matching coefficient

# SSD : Training

$$s_k = s_{\min} + \frac{s_{\max} - s_{\min}}{m - 1}(k - 1), \quad k \in [1, m]$$

Number of feature map

- $S \in [0.2, 0.9]$

- $a_r \in \{1, 2, 3, \frac{1}{2}, \frac{1}{3}\} \rightarrow w_k^a = s_k\sqrt{a_r} , h_k^a = s_k/\sqrt{a_r}$

- Determine scale of default boxes

# SSD : Conclusion



- Proposed new brilliant architecture for object detection

- Faster, but more accurate