# Single Shot Detector (SSD300)

Bahdah Shin

Matt Witman

Dan Moore

# What is Object Detection?

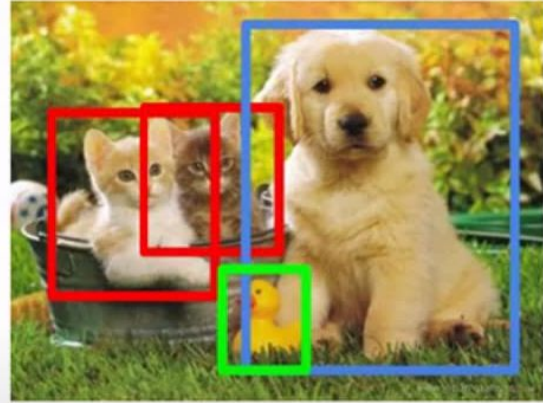# What is SSD?

# Base Convolutions: part 1 (vgg-16 architecture)

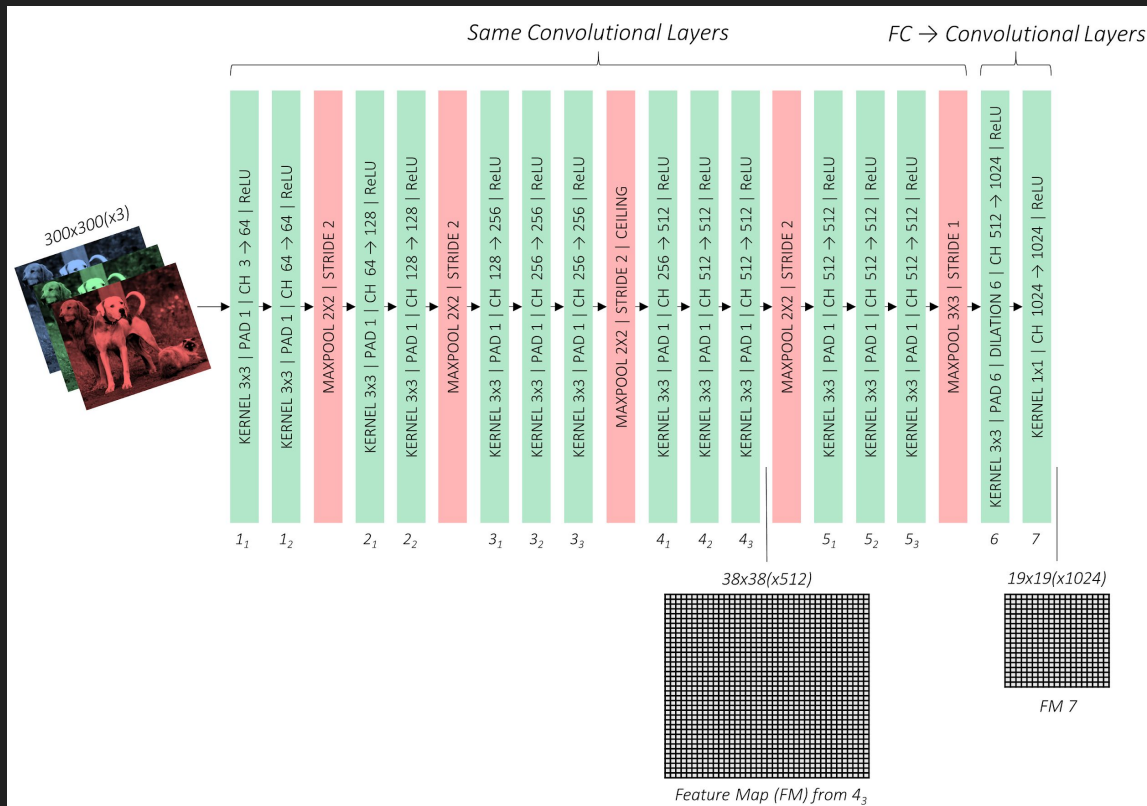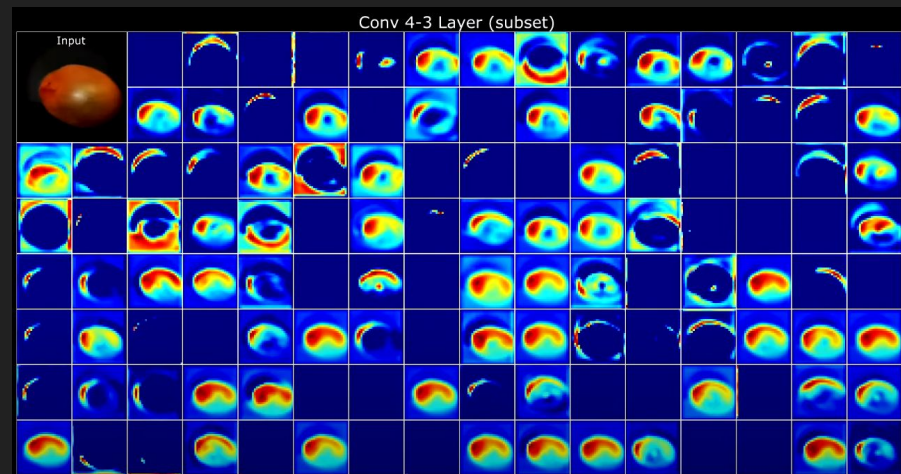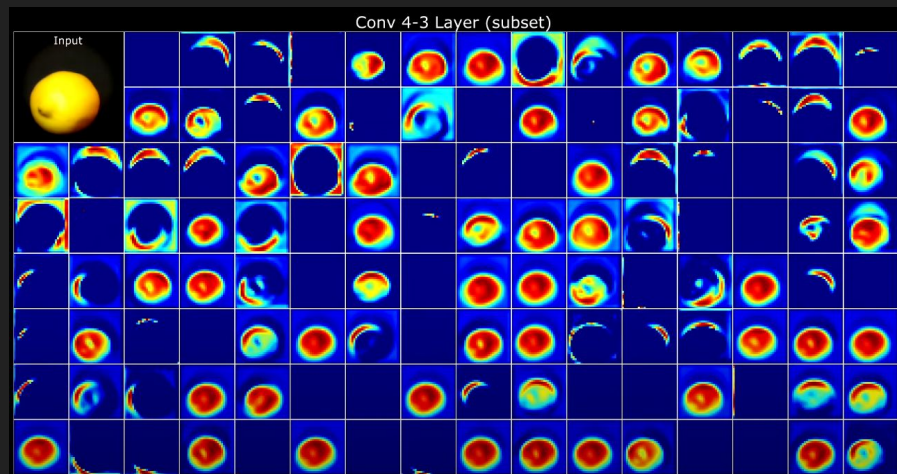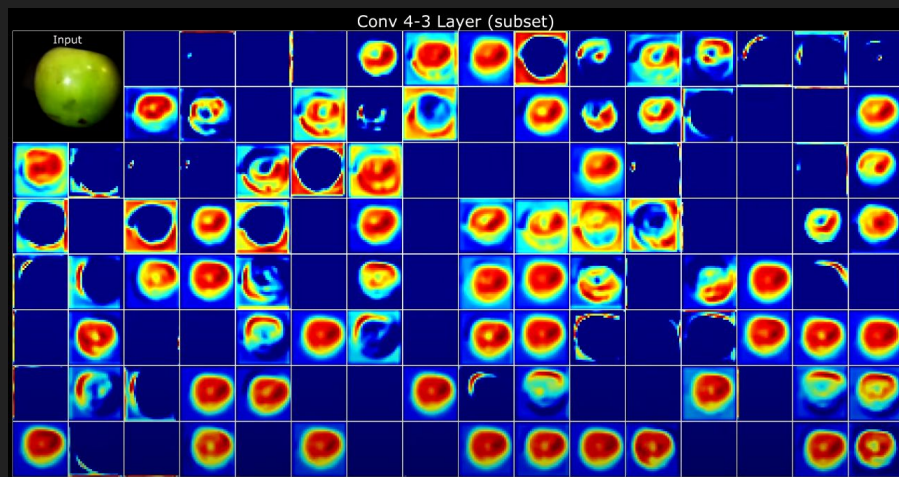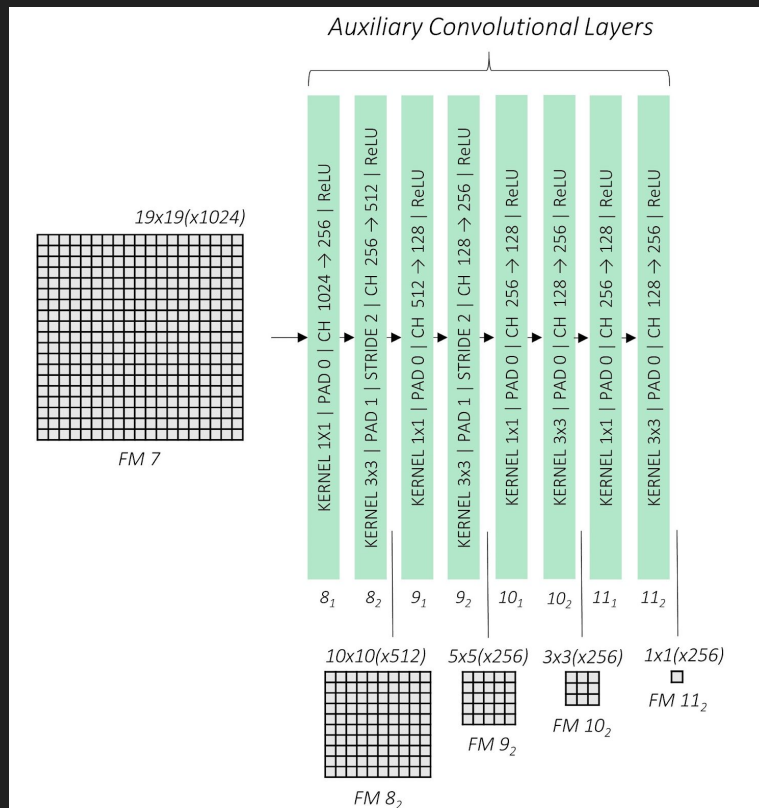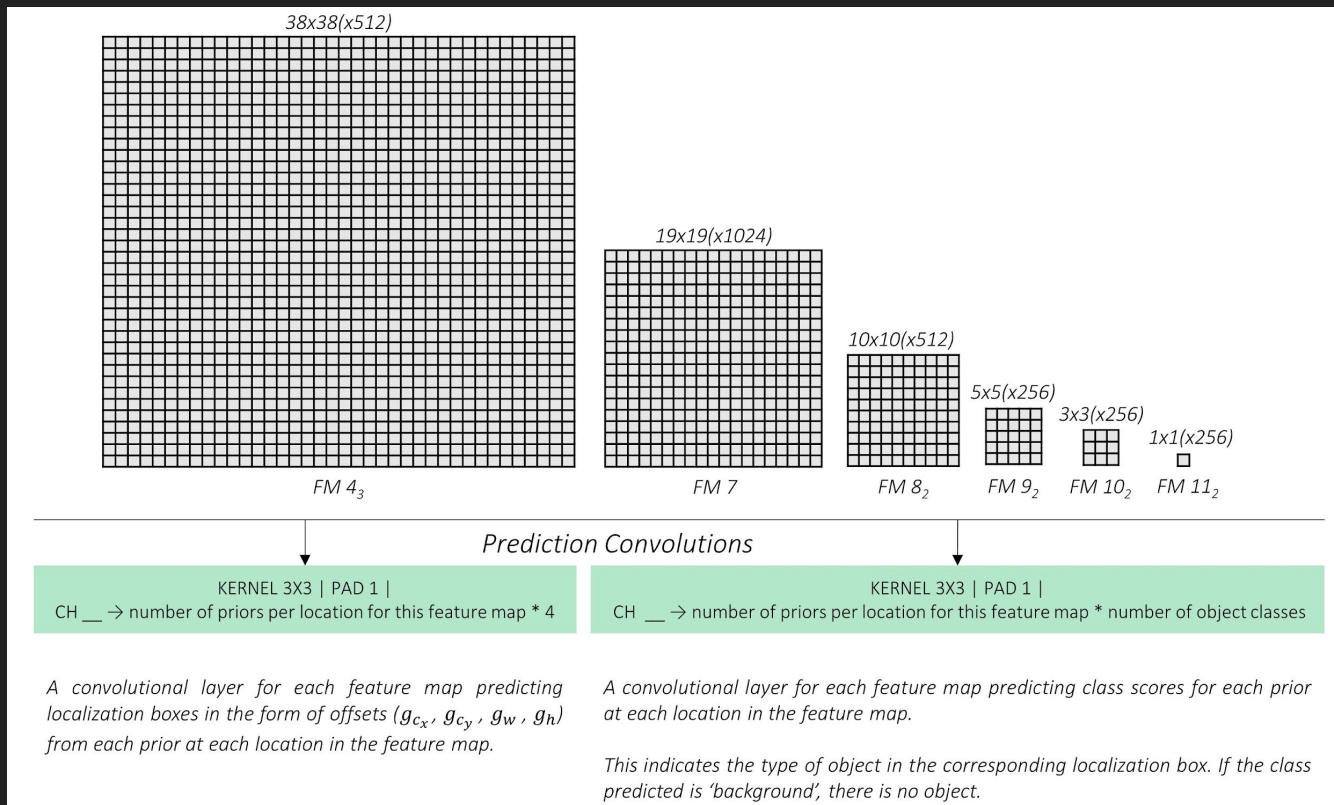# Base Convolutions: Part 2 (modified VGG-16)



Same Convolutional Layers

FC → Convolutional Layers

300x300(x3)

KERNEL 3x3 | PAD 1 | CH 3 → 64 | ReLU — $1_1$
KERNEL 3x3 | PAD 1 | CH 64 → 64 | ReLU — $1_2$
MAXPOOL 2X2 | STRIDE 2
KERNEL 3x3 | PAD 1 | CH 64 → 128 | ReLU — $2_1$
KERNEL 3x3 | PAD 1 | CH 128 → 128 | ReLU — $2_2$
MAXPOOL 2X2 | STRIDE 2
KERNEL 3x3 | PAD 1 | CH 128 → 256 | ReLU — $3_1$
KERNEL 3x3 | PAD 1 | CH 256 → 256 | ReLU — $3_2$
KERNEL 3x3 | PAD 1 | CH 256 → 256 | ReLU — $3_3$
MAXPOOL 2X2 | STRIDE 2 | CEILING
KERNEL 3x3 | PAD 1 | CH 256 → 512 | ReLU — $4_1$
KERNEL 3x3 | PAD 1 | CH 512 → 512 | ReLU — $4_2$
KERNEL 3x3 | PAD 1 | CH 512 → 512 | ReLU — $4_3$
MAXPOOL 2X2 | STRIDE 2
KERNEL 3x3 | PAD 1 | CH 512 → 512 | ReLU — $5_1$
KERNEL 3x3 | PAD 1 | CH 512 → 512 | ReLU — $5_2$
KERNEL 3x3 | PAD 1 | CH 512 → 512 | ReLU — $5_3$
MAXPOOL 3X3 | STRIDE 1
KERNEL 3x3 | PAD 6 | DILATION 6 | CH 512 → 1024 | ReLU — 6
KERNEL 1x1 | CH 1024 → 1024 | ReLU — 7

38x38(x512)

19x19(x1024)

FM 7

Feature Map (FM) from $4_3$

Conv 4-3 Layer (subset)

# Auxiliary Convolutions

# Prediction Convolutions



38x38(x512)

19x19(x1024)

10x10(x512)

5x5(x256)

3x3(x256)

1x1(x256)

FM $4_3$    FM 7    FM $8_2$    FM $9_2$    FM $10_2$    FM $11_2$

*Prediction Convolutions*

KERNEL 3X3 | PAD 1 |
CH ___ → number of priors per location for this feature map * 4

KERNEL 3X3 | PAD 1 |
CH ___ → number of priors per location for this feature map * number of object classes

*A convolutional layer for each feature map predicting localization boxes in the form of offsets ($g_{c_x}$, $g_{c_y}$, $g_w$, $g_h$) from each prior at each location in the feature map.*

*A convolutional layer for each feature map predicting class scores for each prior at each location in the feature map.*

*This indicates the type of object in the corresponding localization box. If the class predicted is 'background', there is no object.*
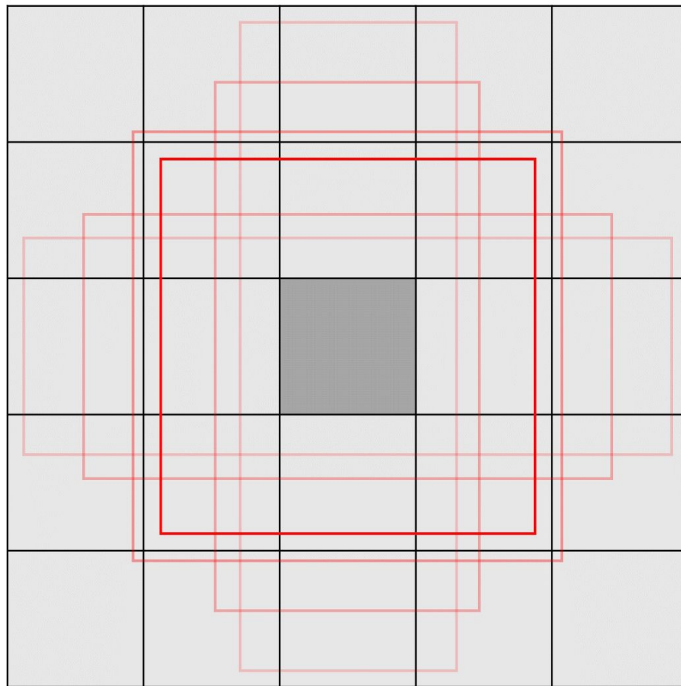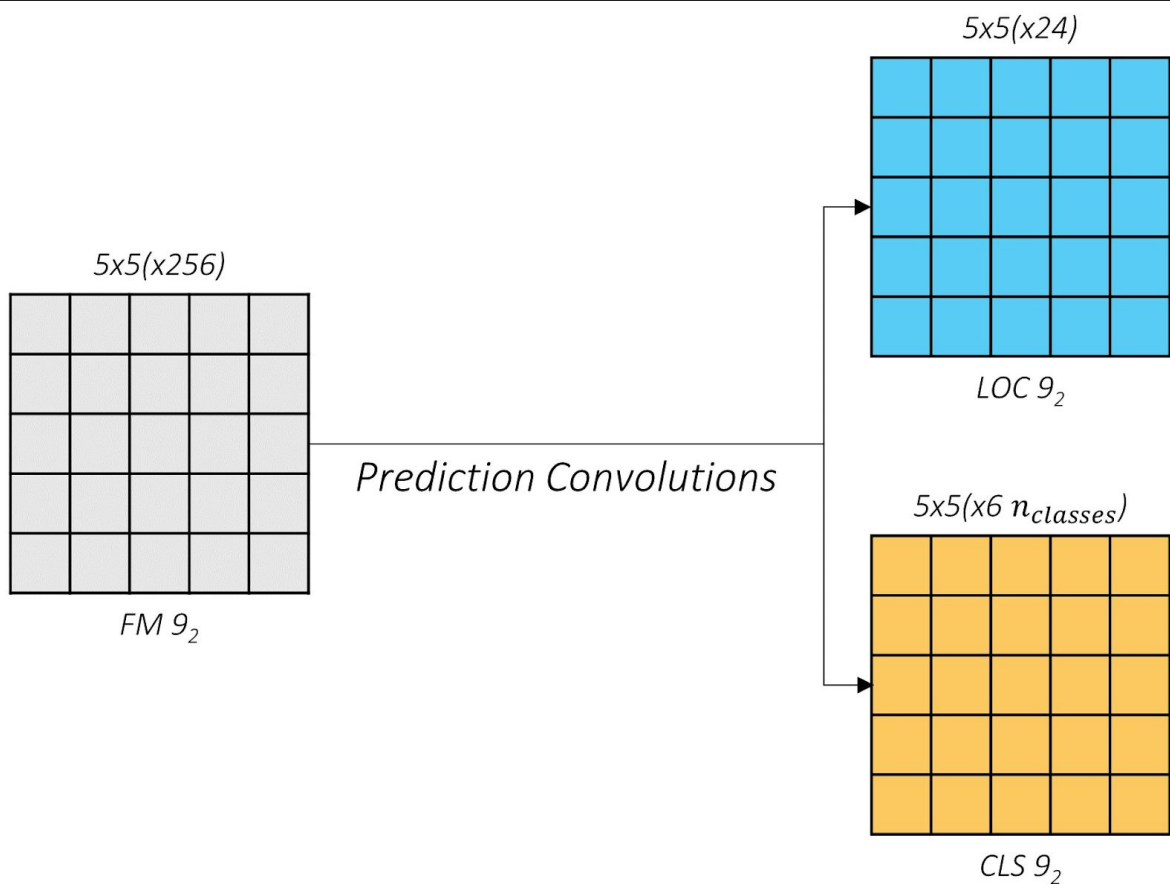
# Model: Anchor boxes (default box) (prior box)

FM 9₂
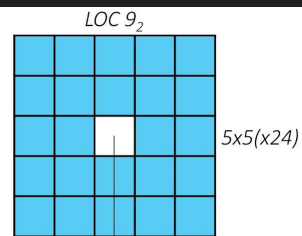
At each location, there are 5 priors with aspect ratios 1, 2, 3, ½, ⅓ and areas equal to that of a square of side 0.55

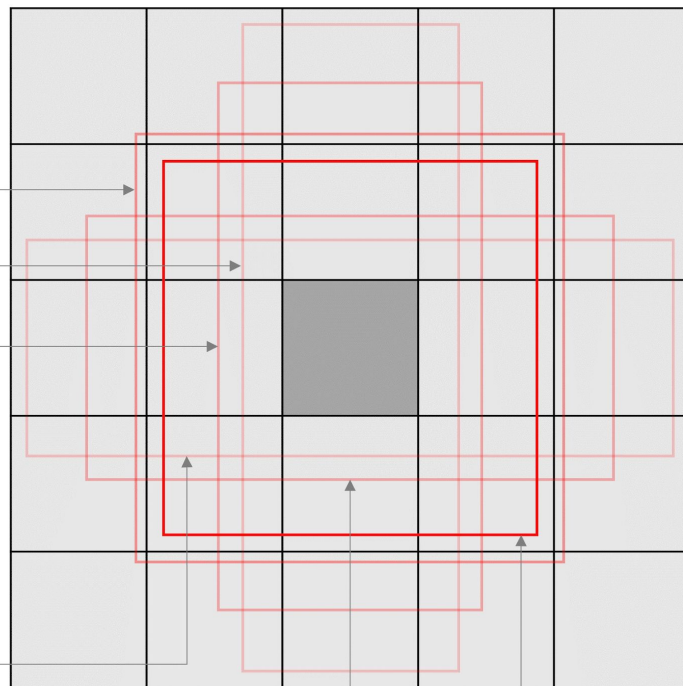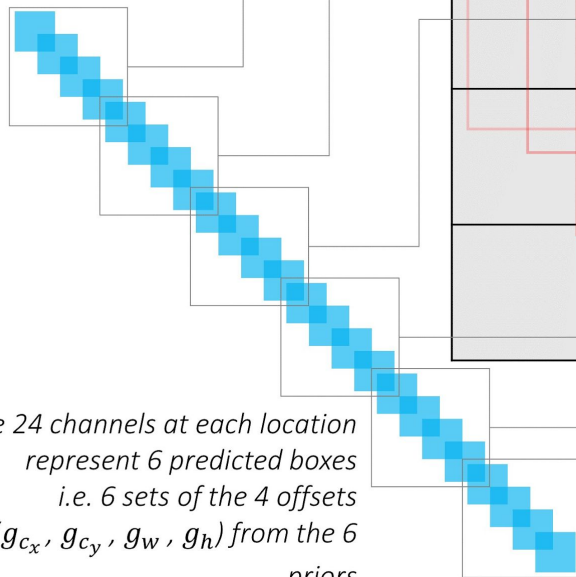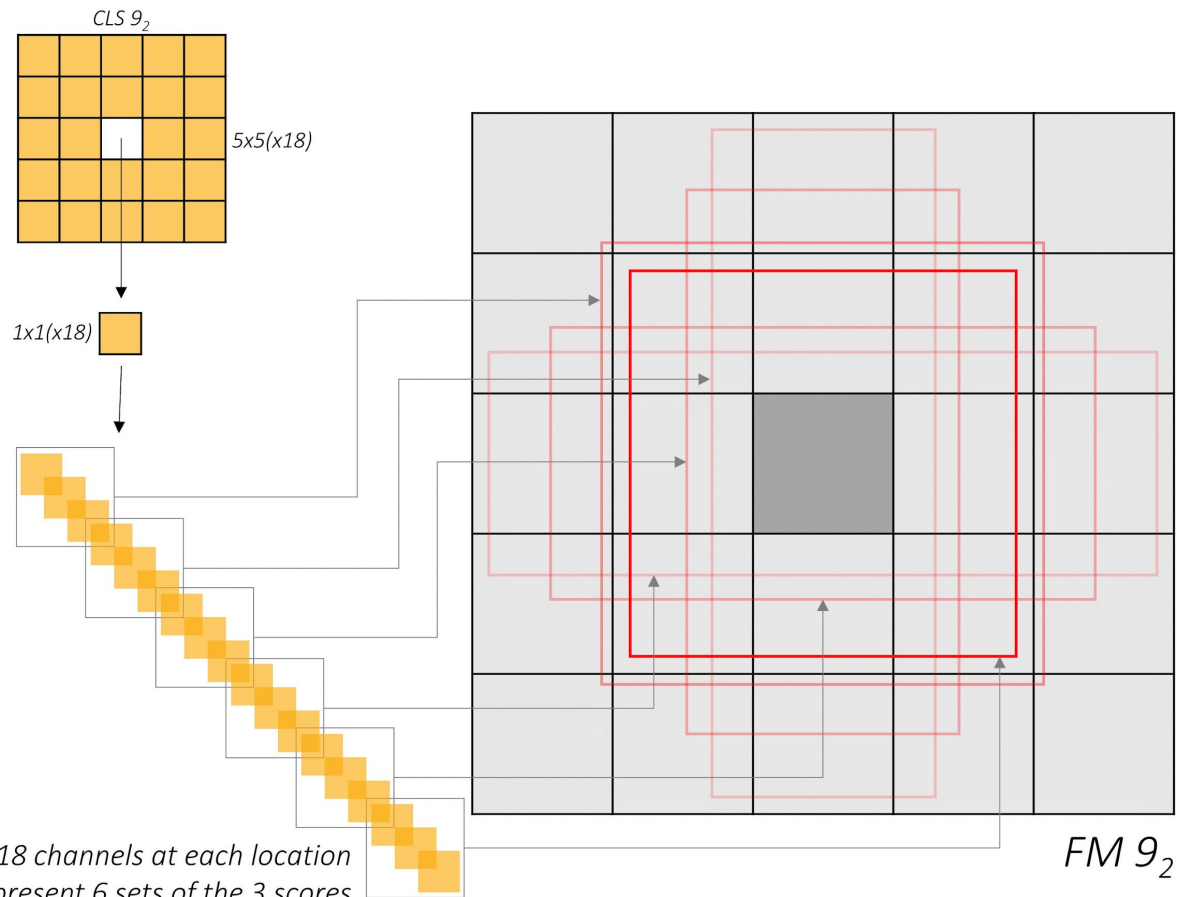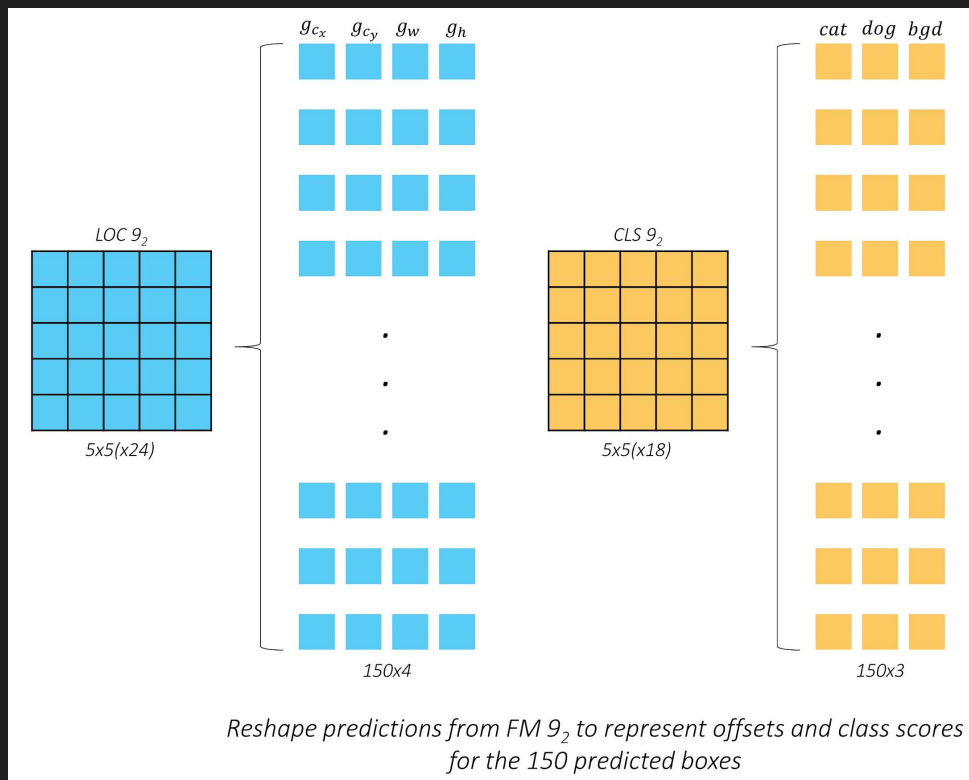Also, a 6ᵗʰ prior with aspect ratio 1 and of side 0.63

LOC $9_2$

5x5(x24)

1x1(x24)

*The 24 channels at each location represent 6 predicted boxes i.e. 6 sets of the 4 offsets $(g_{c_x}, g_{c_y}, g_w, g_h)$ from the 6 priors*

FM $9_2$

Assume $n_{classes} = 3$ $(cat, dog, background)$

CLS $9_2$

5x5(x18)

1x1(x18)

The 18 channels at each location
represent 6 sets of the 3 scores
$(cat, dog, bgd)$ for the 6 priors

FM $9_2$

$g_{c_x}$ $g_{c_y}$ $g_w$ $g_h$  cat dog bgd

LOC 9₂

5x5(x24)

CLS 9₂

5x5(x18)

150x4

150x3

Reshape predictions from FM 9₂ to represent offsets and class scores
for the 150 predicted boxes

$g_{c_x}$ $g_{c_y}$ $g_w$ $g_h$  cat dog bgd

Reshaped predictions
from all feature maps
stacked together

8732x4

8732x3

A total of 8732 predicted boxes

Conv4_3: 38 * 38 * 4 = 5776

Conv7: 19 * 19 * 6 = 2166

Conv8_2: 10 * 10 * 6 = 600
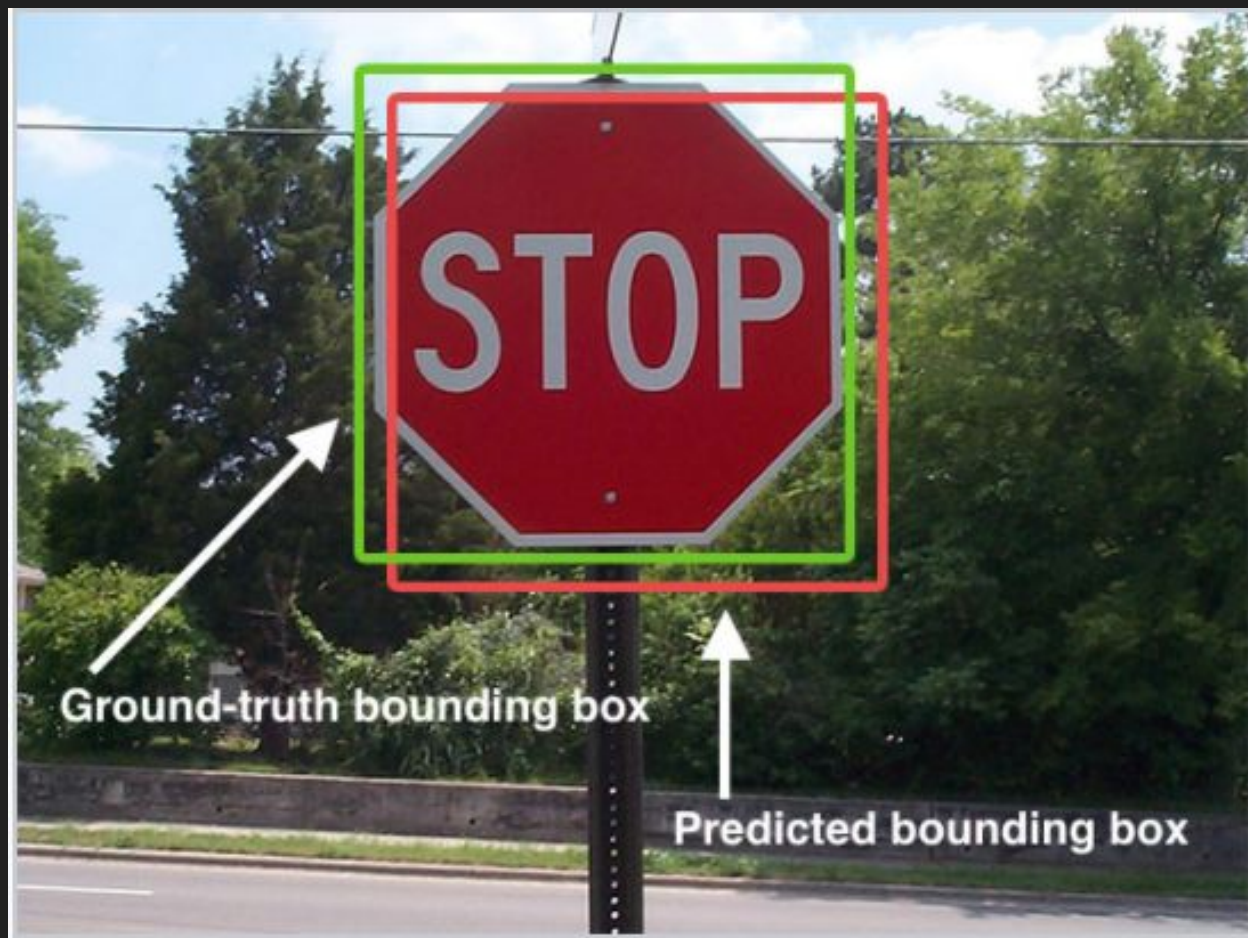
Conv9_2: 5 * 5 * 6 = 150

Conv10_2: 3 * 3 * 4 = 36

Conv11_2: 4

Total: 5776+ 2166 + 600 + 150 + 36 + 4 = 8732 priors

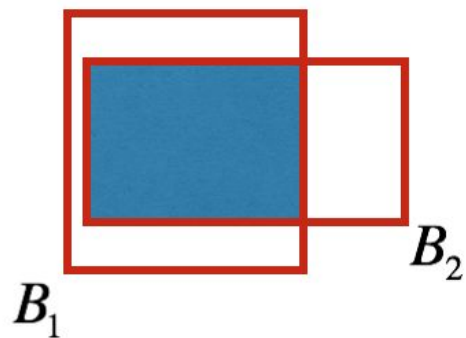# Non-Maximum Suppression

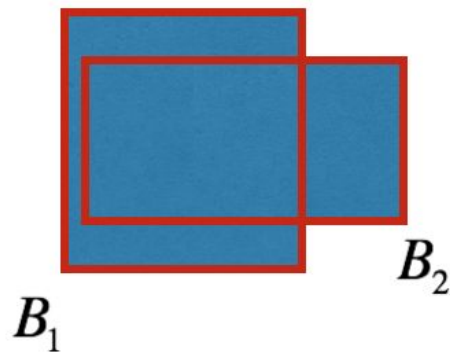# Training

**Intersection**

**Union**

**Intersection over Union**

$$IoU = \frac{B_1 \cap B_2}{B_1 \cup B_2} =$$

$B_1$    $B_2$

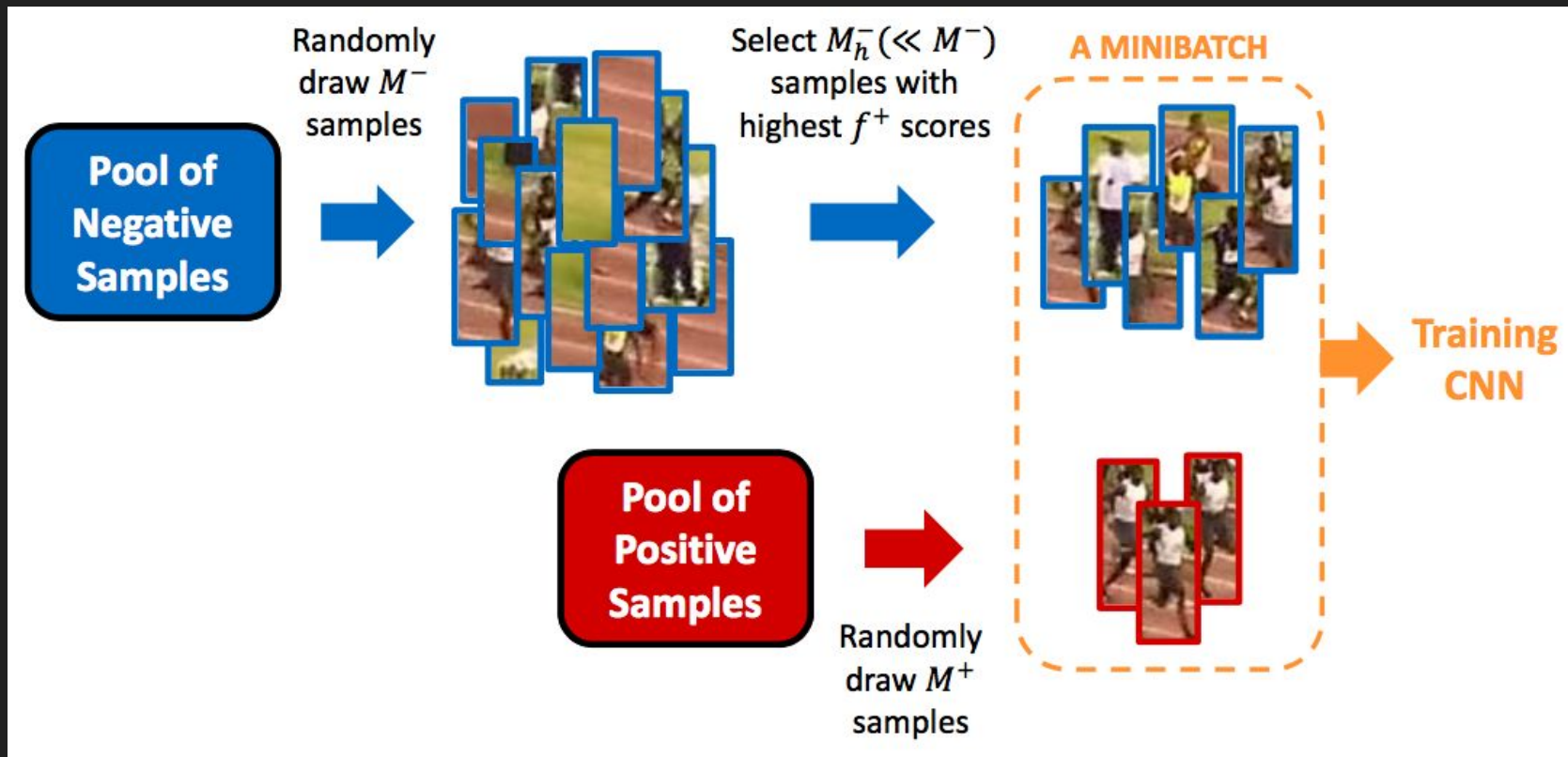IoU: 0.4034     IoU: 0.7330     IoU: 0.9264

**Poor**      **Good**      **Excellent**

# Hard Negative Mining

# Not always accurate

# Sources

1. https://towardsdatascience.com/understanding-ssd-multibox-real-time-object-detection-in-deep-learning-495ef744fab

2. https://d2l.ai/chapter_computer-vision/ssd.html

3. https://towardsdatascience.com/object-detection-with-neural-networks-a4e2c46b4491

4. https://www.youtube.com/watch?v=RNnKtNrsrmg

5. https://arxiv.org/pdf/1512.02325.pdf

6. https://arxiv.org/pdf/1409.1556.pdf

7. https://arxiv.org/pdf/1409.1556.pdf