

What is object detection?



이름 : 김지원

전공 : 항공우주공학 / 컴퓨터공학

연구분야 : **Object Detection**

연구실 : [지능기전공학부 RCV연구실](#)

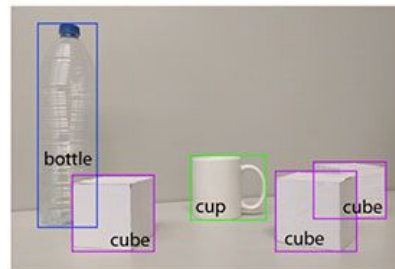
메일 : jwkim@rcv.sejong.ac.kr

Github : <https://github.com/socome>

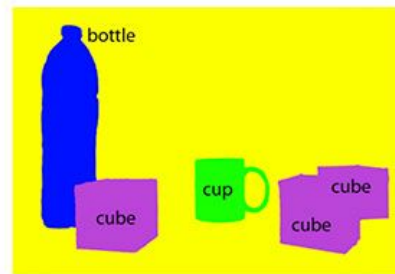
Blog : <https://socome.github.io/about/>



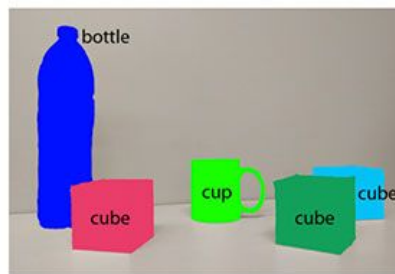
(a) Image classification



(b) Object localization



(c) Semantic segmentation



(d) Instance segmentation

R-CNN → **OverFeat** → MultiBox → SPP-Net → MR-CNN → DeepBox → AttentionNet →
2013.11 ICLR' 14 CVPR' 14 ECCV' 14 ICCV' 15 ICCV' 15 ICCV' 15

Fast R-CNN → DeepProposal → **Faster R-CNN** → **OHEM** → **YOLO v1** → G-CNN → AZNet →
ICCV' 15 ICCV' 15 NIPS' 15 CVPR' 16 CVPR' 16 CVPR' 16 CVPR' 16

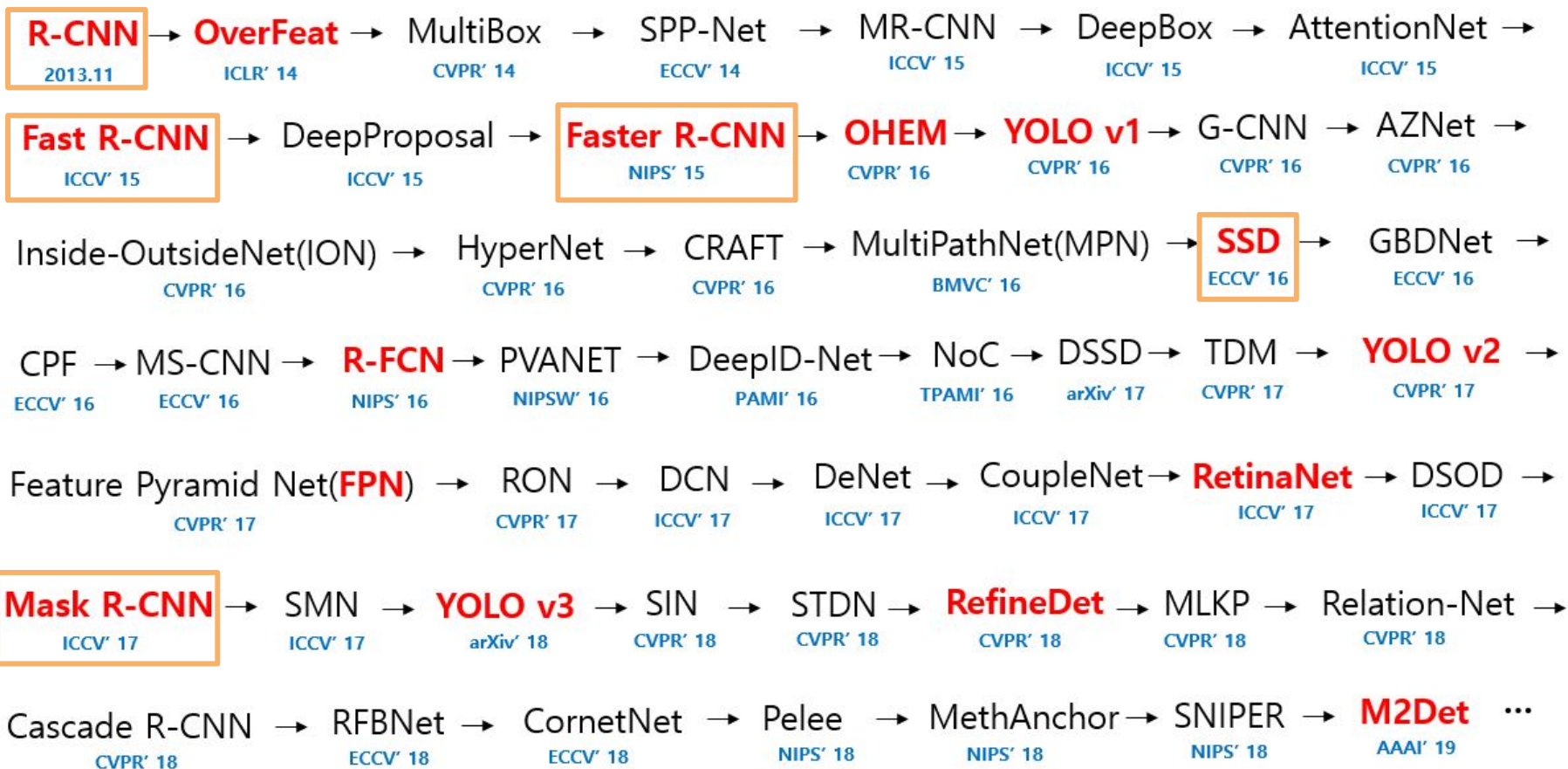
Inside-OutsideNet(ION) → HyperNet → CRAFT → MultiPathNet(MPN) → **SSD** → GBDNet →
CVPR' 16 CVPR' 16 CVPR' 16 BMVC' 16 ECCV' 16 ECCV' 16

CPF → MS-CNN → **R-FCN** → PVANET → DeepID-Net → NoC → DSSD → TDM → **YOLO v2** →
ECCV' 16 ECCV' 16 NIPS' 16 NIPS' 16 PAMI' 16 TPAMI' 16 arXiv' 17 CVPR' 17 CVPR' 17

Feature Pyramid Net(**FPN**) → RON → DCN → DeNet → CoupleNet → **RetinaNet** → DSOD →
CVPR' 17 CVPR' 17 ICCV' 17 ICCV' 17 ICCV' 17 ICCV' 17 ICCV' 17

Mask R-CNN → SMN → **YOLO v3** → SIN → STDN → **RefineDet** → MLKP → Relation-Net →
ICCV' 17 ICCV' 17 arXiv' 18 CVPR' 18 CVPR' 18 CVPR' 18 CVPR' 18 CVPR' 18

Cascade R-CNN → RFBNet → CornetNet → Pelee → MethAnchor → SNIPER → **M2Det** ...
CVPR' 18 ECCV' 18 ECCV' 18 NIPS' 18 NIPS' 18 NIPS' 18 NIPS' 18 AAAI' 19



R-CNN → **OverFeat** → MultiBox → SPP-Net → MR-CNN → DeepBox → AttentionNet →
2013.11 ICLR' 14 CVPR' 14 ECCV' 14 ICCV' 15 ICCV' 15 ICCV' 15

Fast R-CNN → DeepProposal → **Faster R-CNN** → **OHEM** → **YOLO v1** → G-CNN → AZNet →
ICCV' 15 ICCV' 15 NIPS' 15 CVPR' 16 CVPR' 16 CVPR' 16 CVPR' 16

Inside-OutsideNet(ION) → HyperNet → CRAFT → MultiPathNet(MPN) → **SSD** → GBDNet →
CVPR' 16 CVPR' 16 CVPR' 16 BMVC' 16 ECCV' 16 ECCV' 16

CPF → MS-CNN → **R-FCN** → PVANET → DeepID-Net → NoC → DSSD → TDM → **YOLO v2** →
ECCV' 16 ECCV' 16 NIPS' 16 NIPSW' 16 PAMI' 16 TPAMI' 16 arXiv' 17 CVPR' 17 CVPR' 17

Feature Pyramid Net(**FPN**) → RON → DCN → DeNet → CoupleNet → **RetinaNet** → DSOD →
CVPR' 17 CVPR' 17 ICCV' 17 ICCV' 17 ICCV' 17 ICCV' 17 ICCV' 17

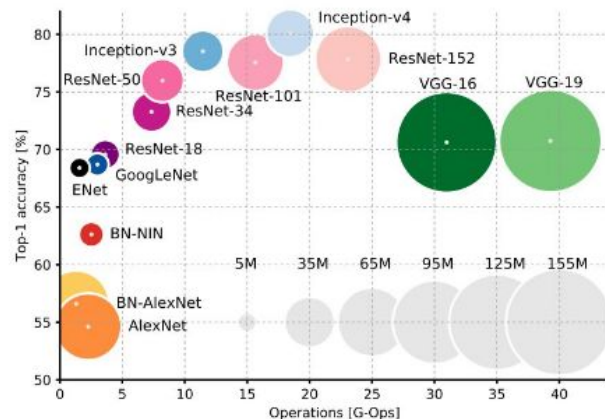
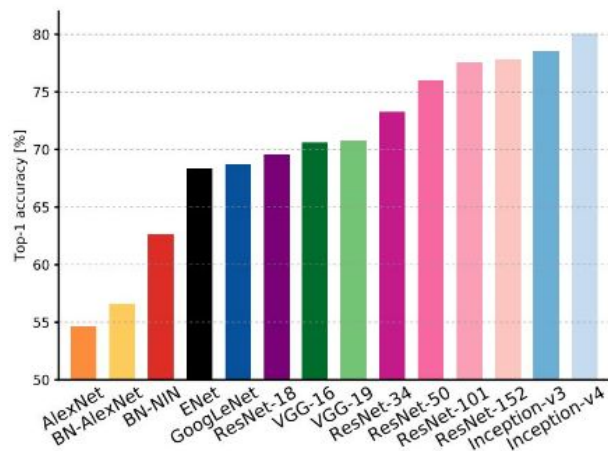
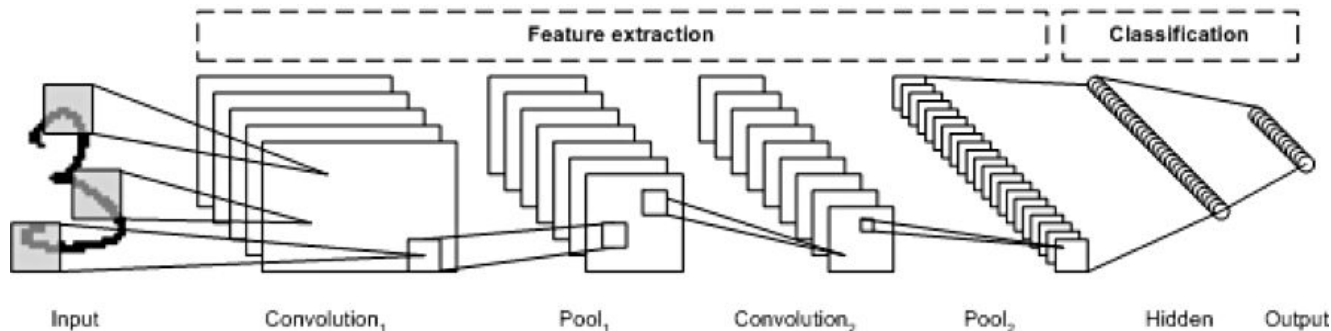
Mask R-CNN → SMN → **YOLO v3** → SIN → STDN → **RefineDet** → MLKP → Relation-Net →
ICCV' 17 ICCV' 17 arXiv' 18 CVPR' 18 CVPR' 18 CVPR' 18 CVPR' 18 CVPR' 18

Cascade R-CNN → RFBNet → CornetNet → Pelee → MethAnchor → SNIPER → **M2Det** ...
CVPR' 18 ECCV' 18 ECCV' 18 NIPS' 18 NIPS' 18 NIPS' 18 NIPS' 18 AAAI' 19

About CNN



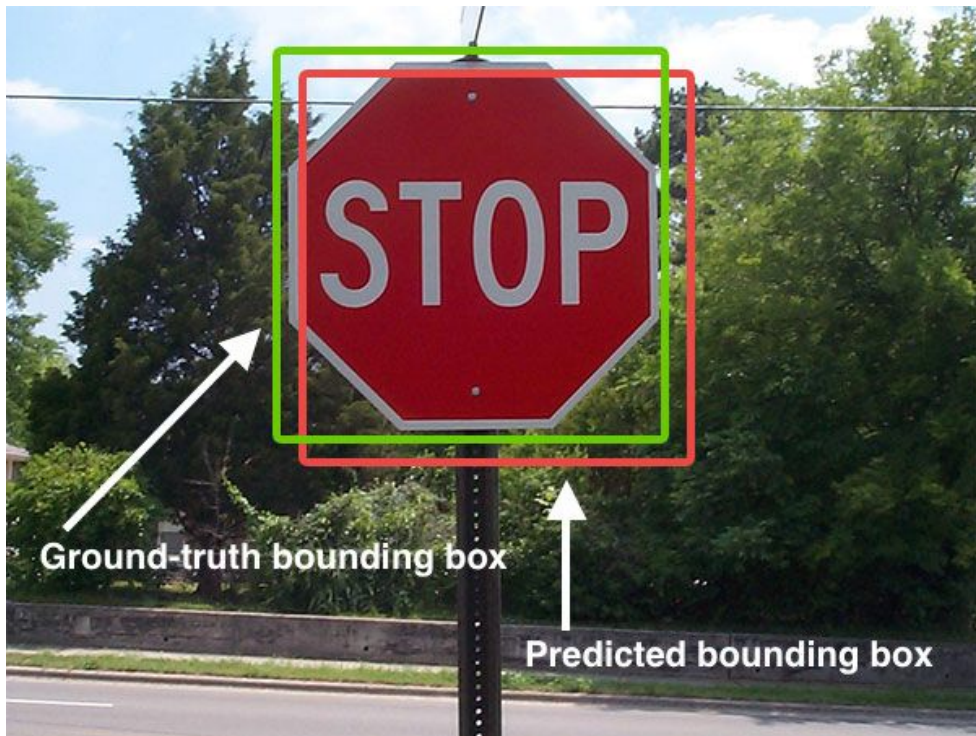
(a) Image classification



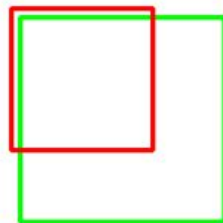
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Object Detection

IOU

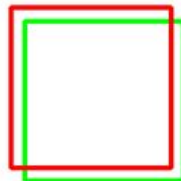


IoU: 0.4034



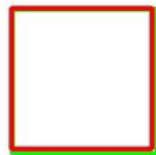
Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

Box Regression

R-CNN: Pipeline

R-CNN in detail

- Bounding box regression*

- Proposal $P = (P_x, P_y, P_w, P_h)$ / Ground-truth $G = (G_x, G_y, G_w, G_h)$

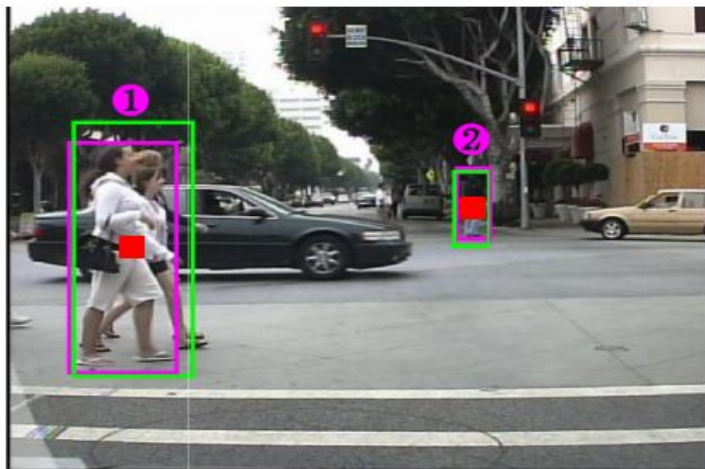
- **For box-scale invariance**, regression target t for the training pair (P, G)

$$t_x = (G_x - P_x) / P_w$$

$$t_y = (G_y - P_y) / P_y$$

$$t_w = \log(G_w / P_w)$$

$$t_h = \log(G_h / P_h)$$



① $G = (57.4, 141.0, 97.2, 237.0)$
 $P = (61.6, 120.0, 107.3, 261.7)$

$$t_x = (57.4 - 61.6) / 107.3 = -0.039$$

$$t_y = (141.0 - 120.0) / 261.7 = 0.080$$

$$t_w = \log(97.2 / 107.3) = -0.043$$

$$t_h = \log(237.0 / 261.7) = -0.043$$

② $G = (408.0, 167.0, 29.9, 73.0)$
 $P = (405.8, 170.8, 31.5, 76.9)$

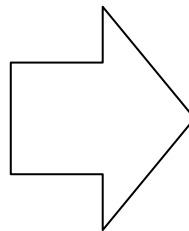
$$t_x = (408.0 - 405.8) / 31.5 = 0.070$$

$$t_y = (167.0 - 170.8) / 76.9 = -0.049$$

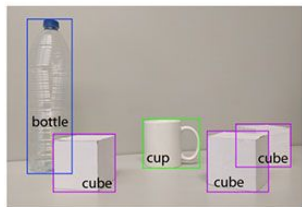
$$t_w = \log(29.9 / 31.5) = -0.023$$

$$t_h = \log(73.0 / 76.9) = -0.023$$

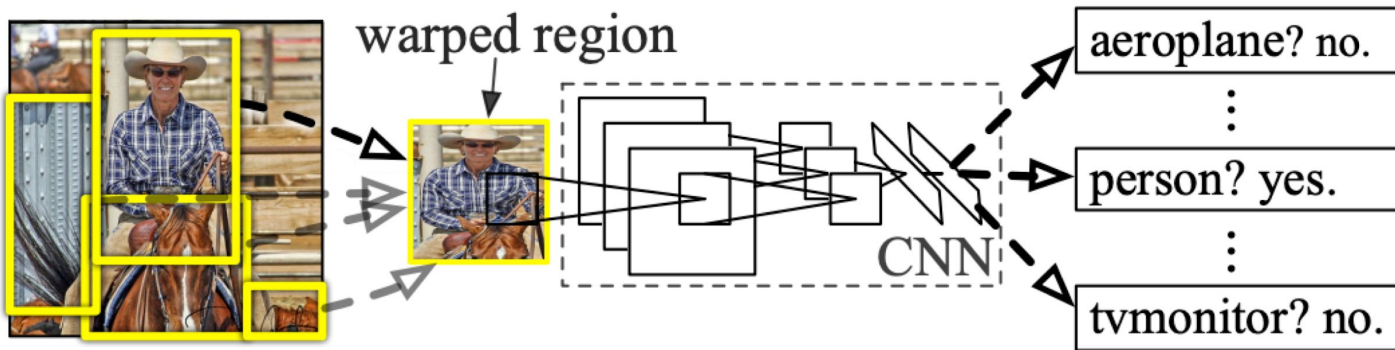
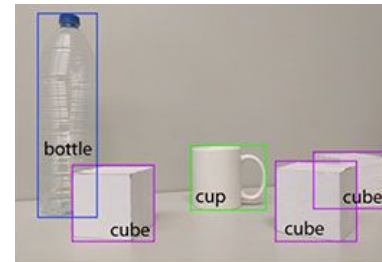
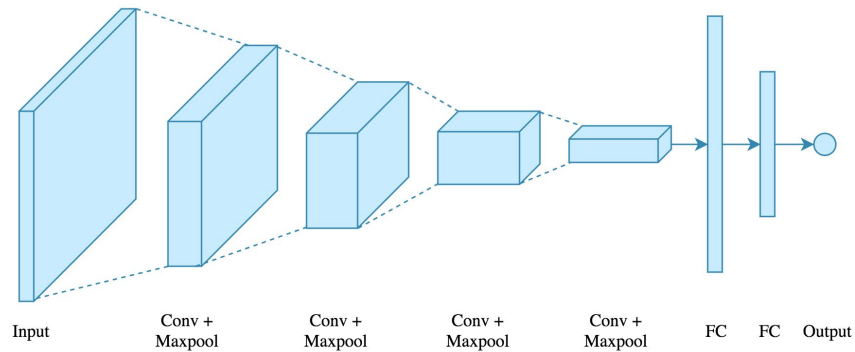
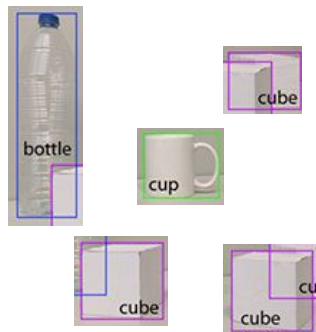
NMS



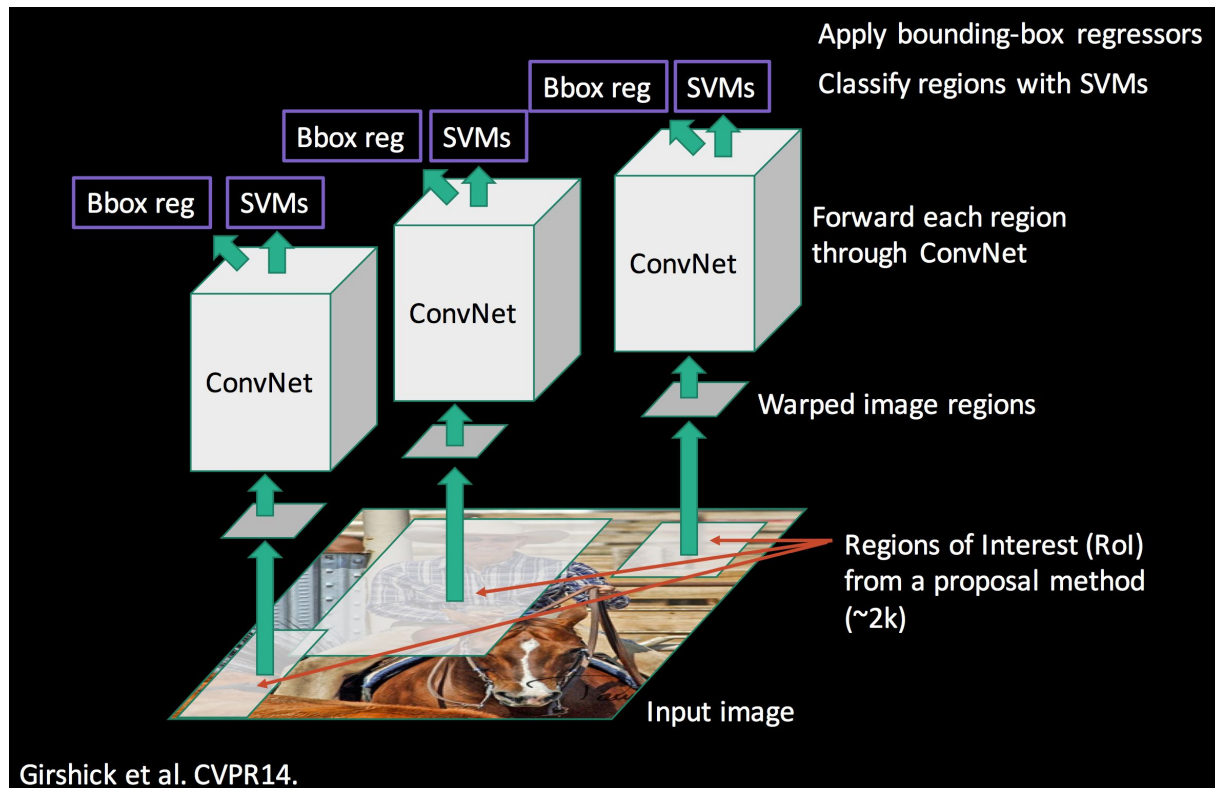
RCNN , Fast RCNN, Faster RCNN, Mask RCNN



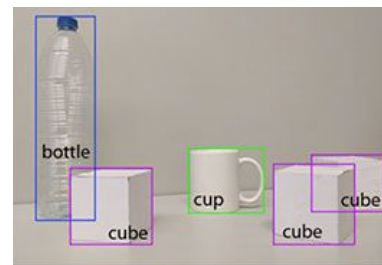
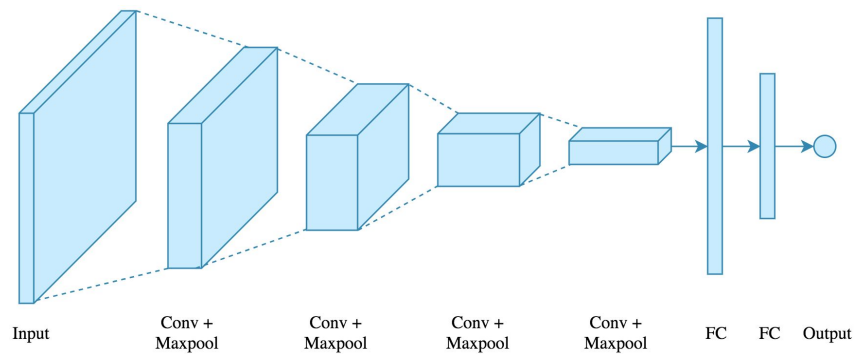
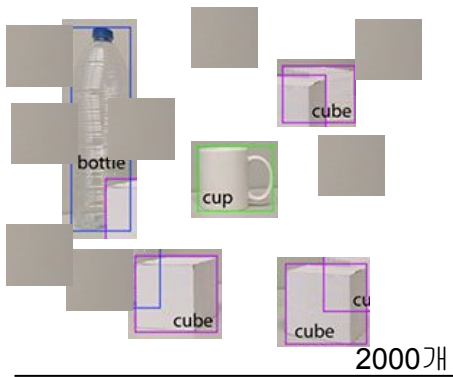
(b) Object localization



RCNN



RCNN , Fast RCNN, Faster RCNN, Mask RCNN



2000개의 모든 **proposal**을 처리하기에는 **너무 많은 시간 소모**

Fast R-CNN (training)

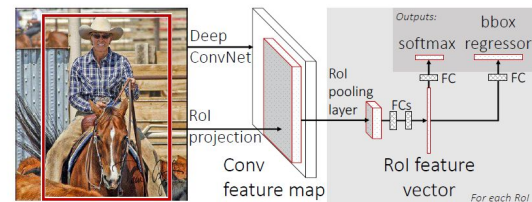
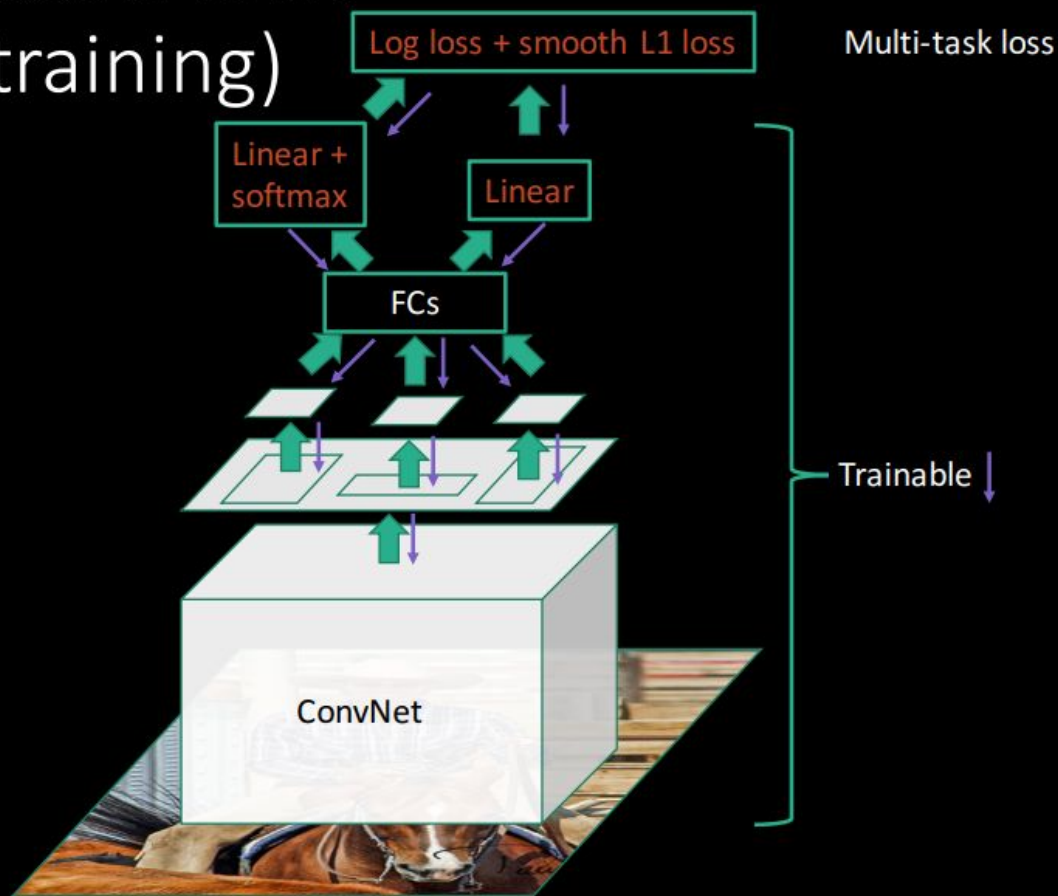
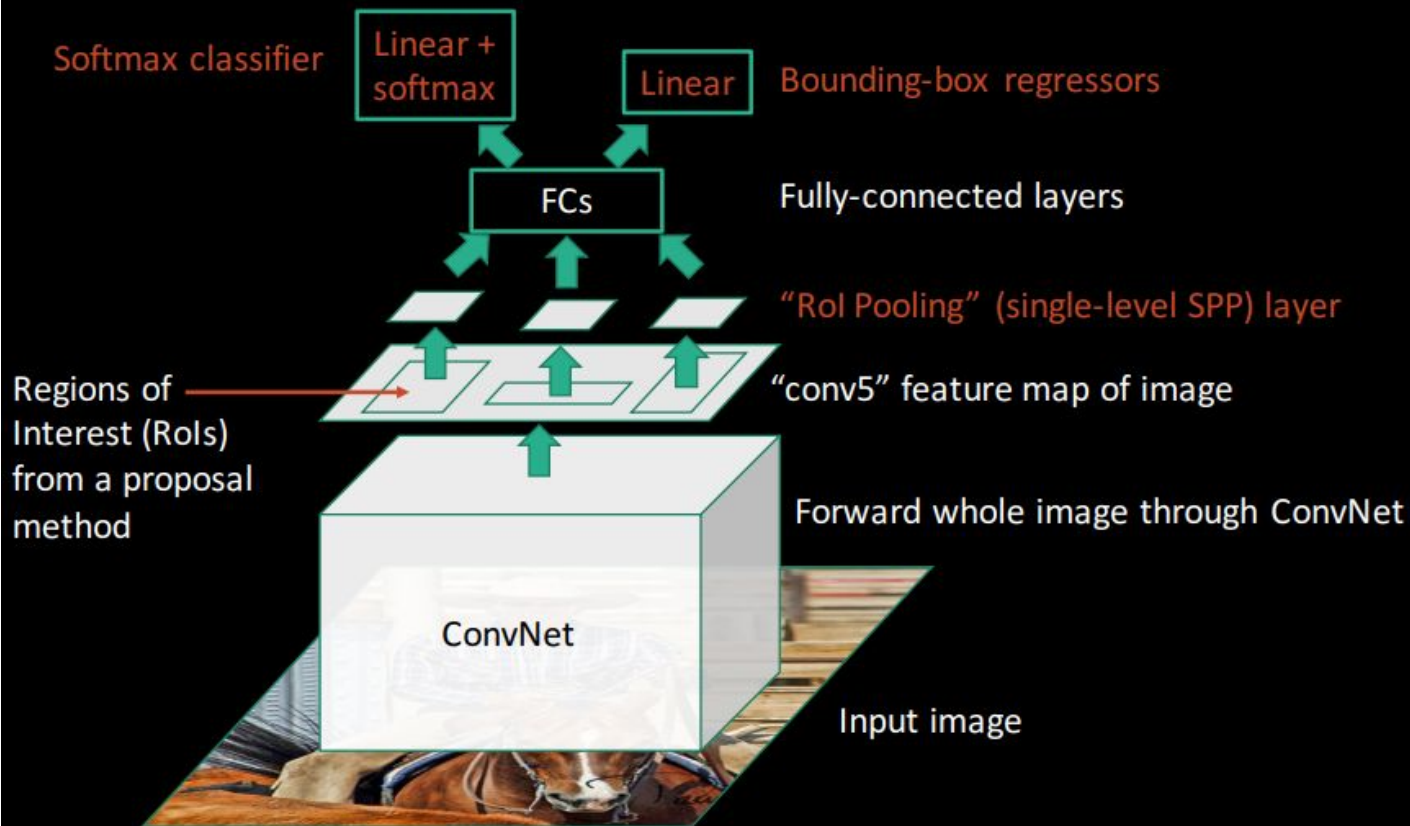
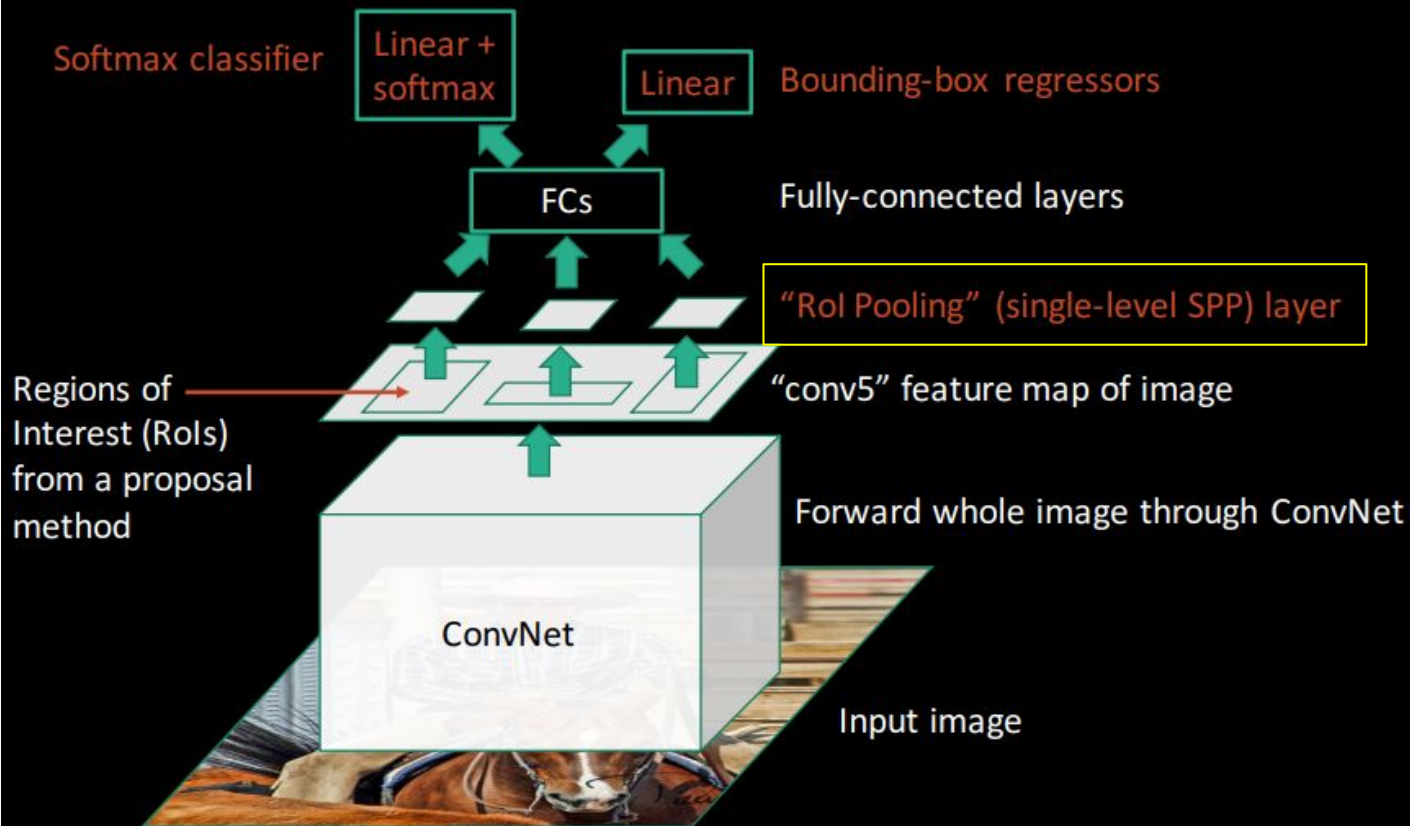


Figure 1. Fast R-CNN architecture. An input image and multiple regions of interest (RoIs) are input into a fully convolutional network. Each RoI is pooled into a fixed-size feature map and then mapped to a feature vector by fully connected layers (FCs). The network has two output vectors per RoI: softmax probabilities and per-class bounding-box regression offsets. The architecture is trained end-to-end with a multi-task loss.

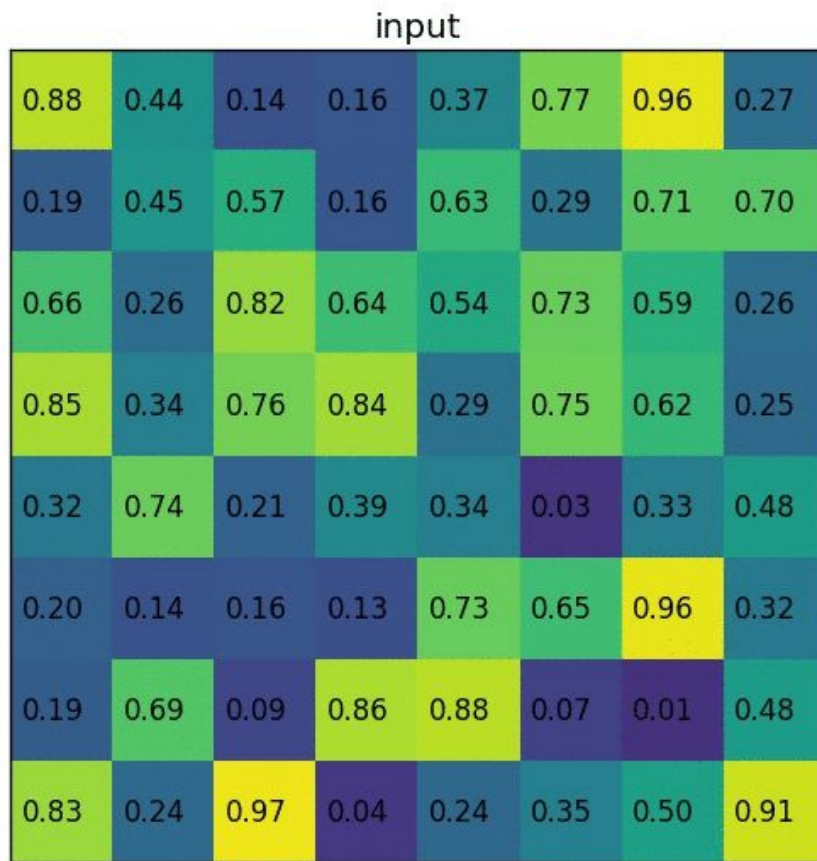
Fast R-CNN (test time)



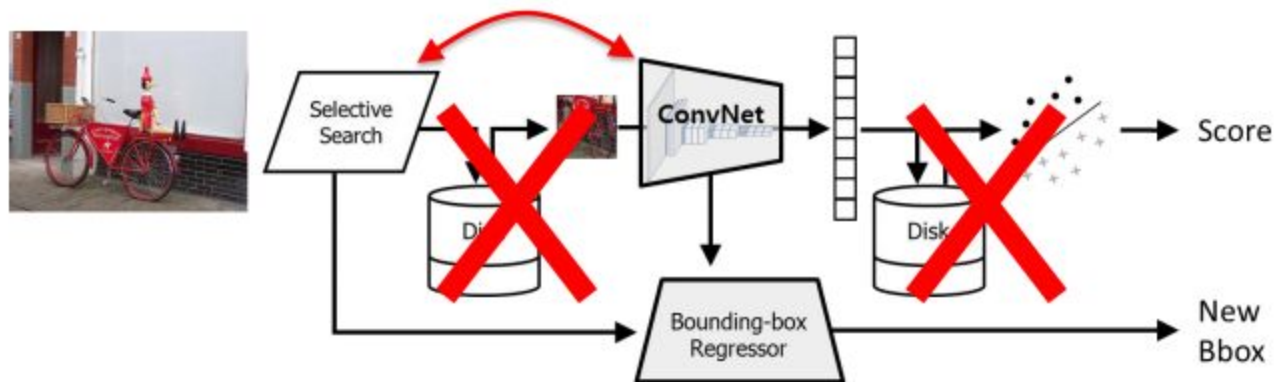
Fast R-CNN (test time)



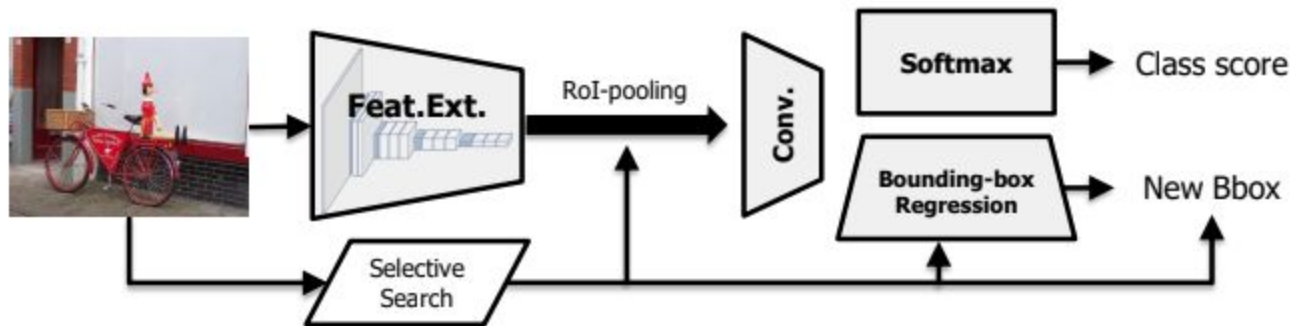
ROI Pooling



R-CNN

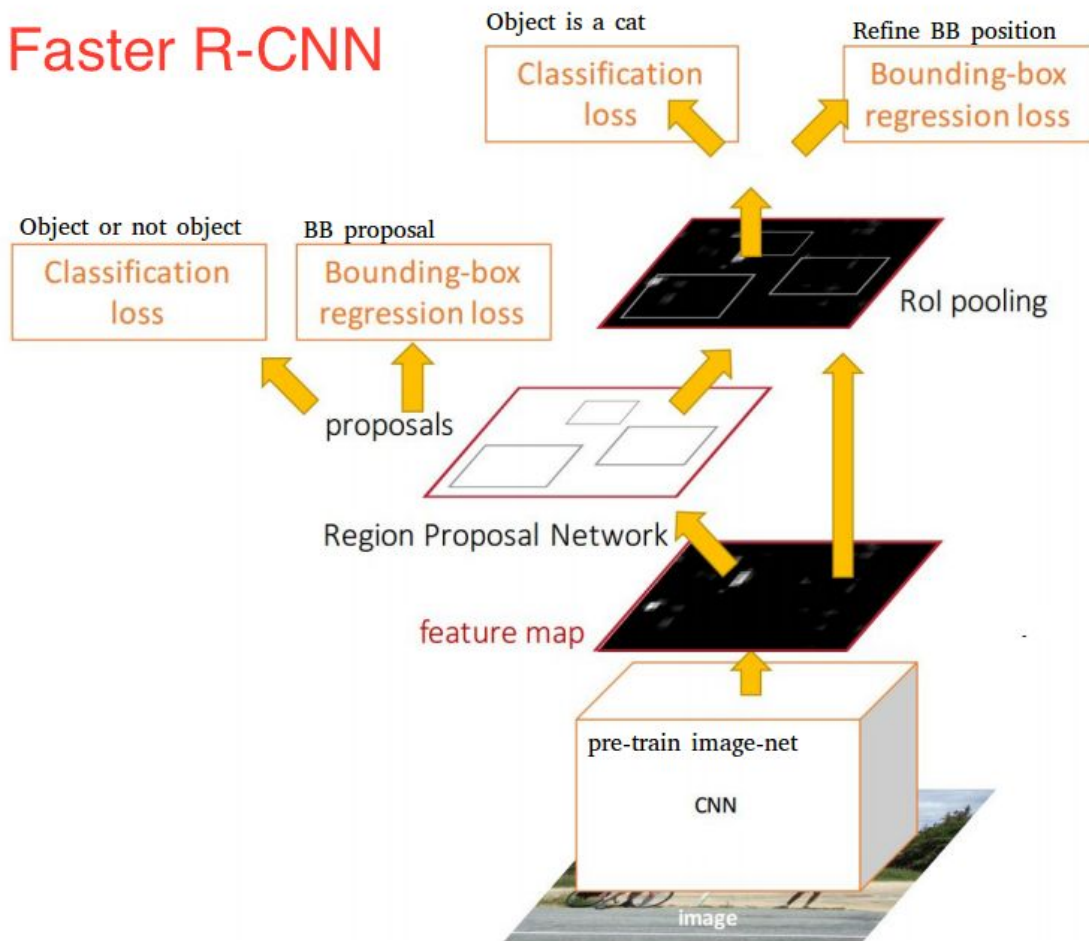


Fast R-CNN

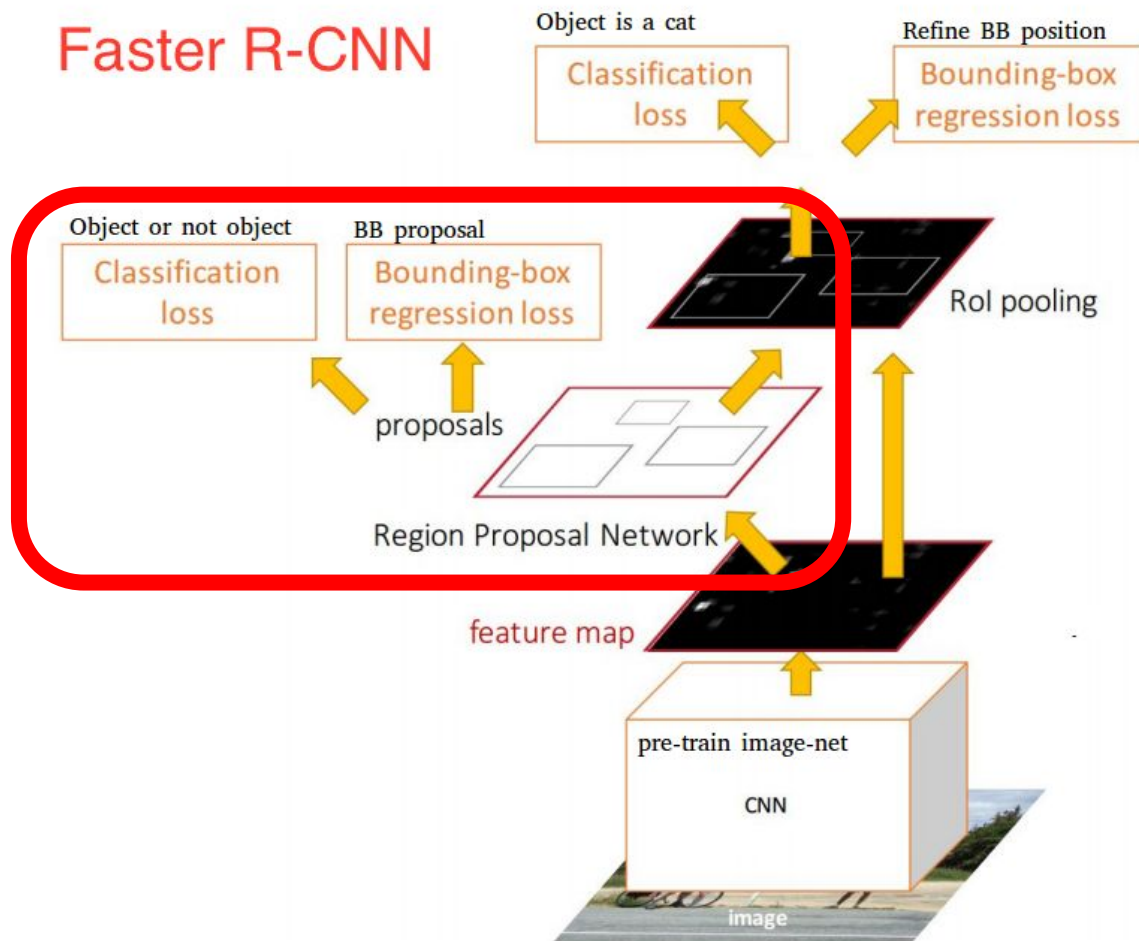


RCNN , Fast RCNN, Faster RCNN, Mask RCNN

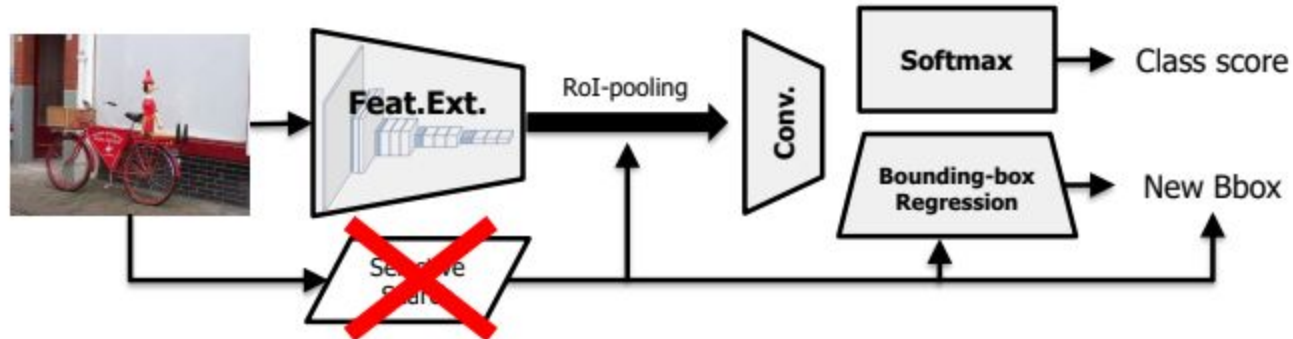
Faster R-CNN



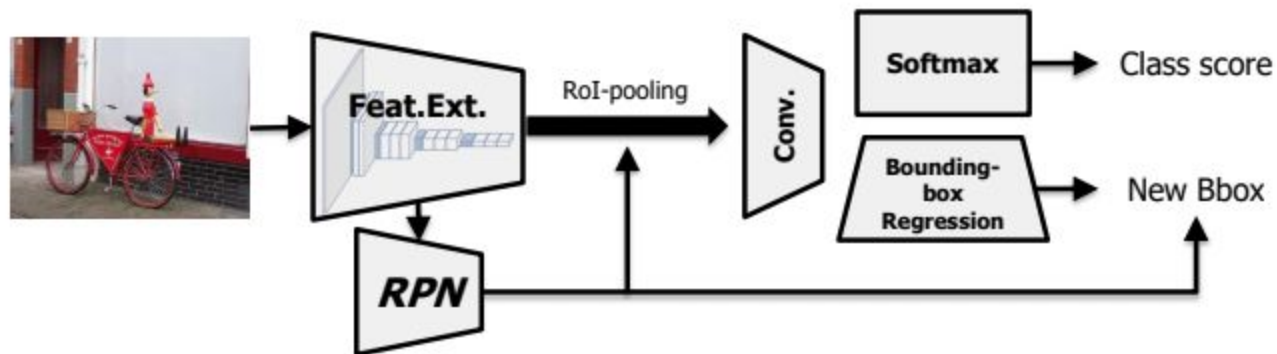
Faster R-CNN



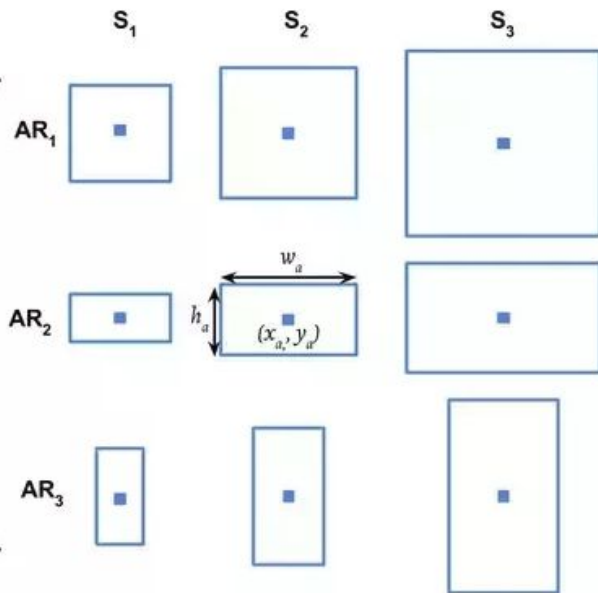
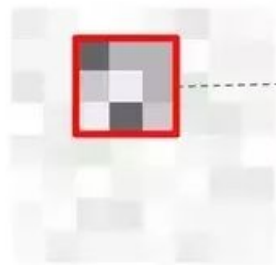
Fast R-CNN



Faster R-CNN

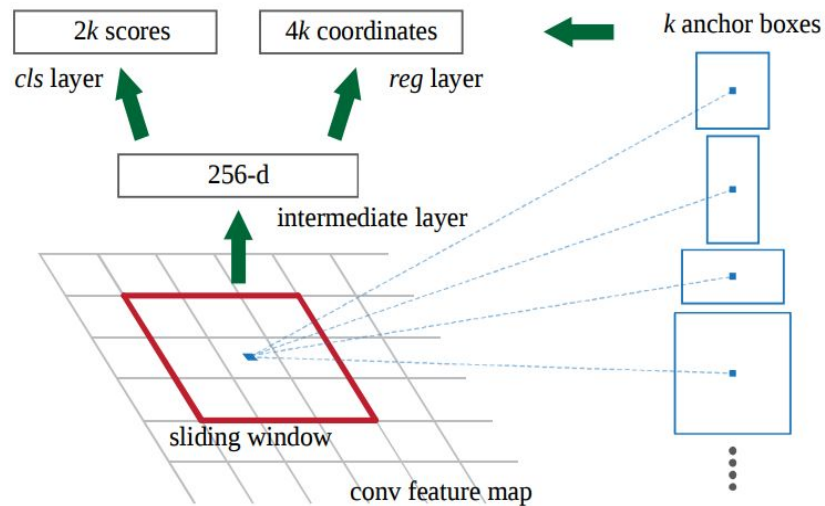


Generate 9 anchors for each **sliding window** on conv. feature map



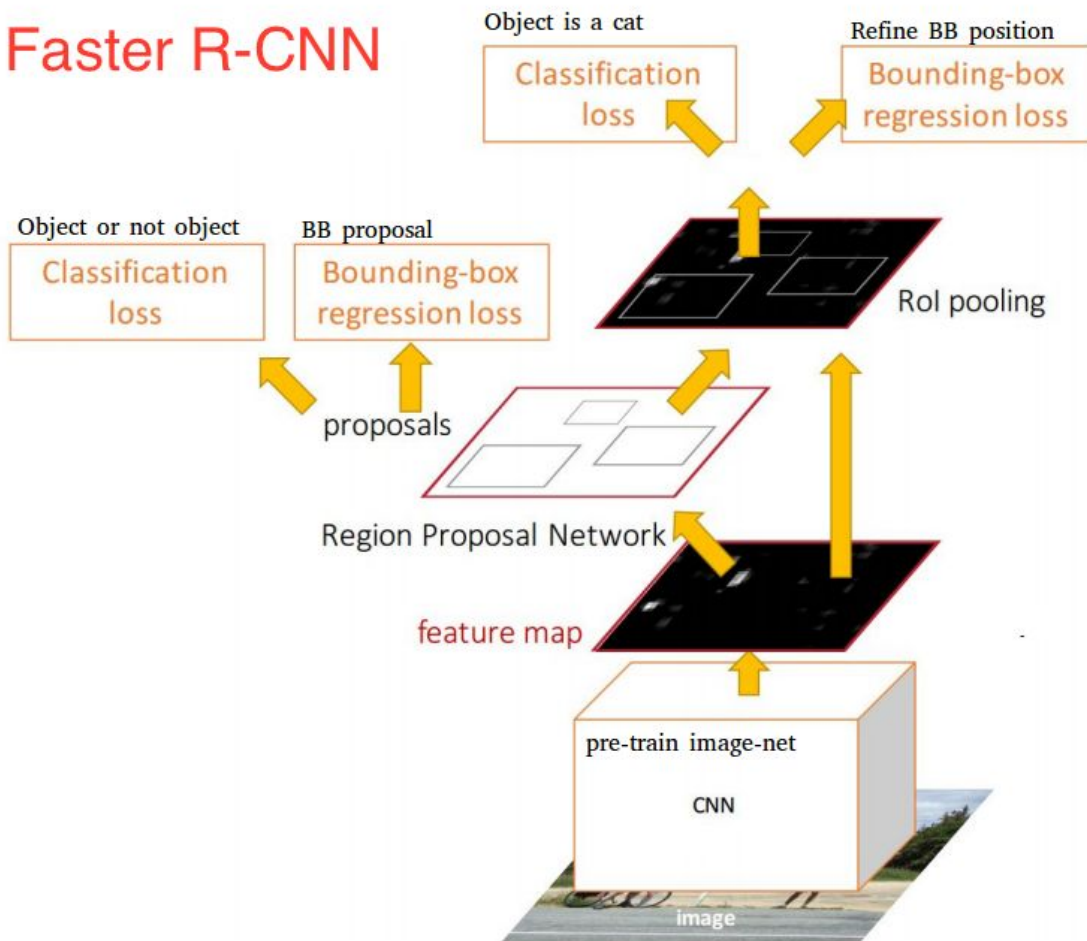
w_a : anchor's width
 h_a : anchor's height
 x_a, y_a : anchor's center

@vmirly

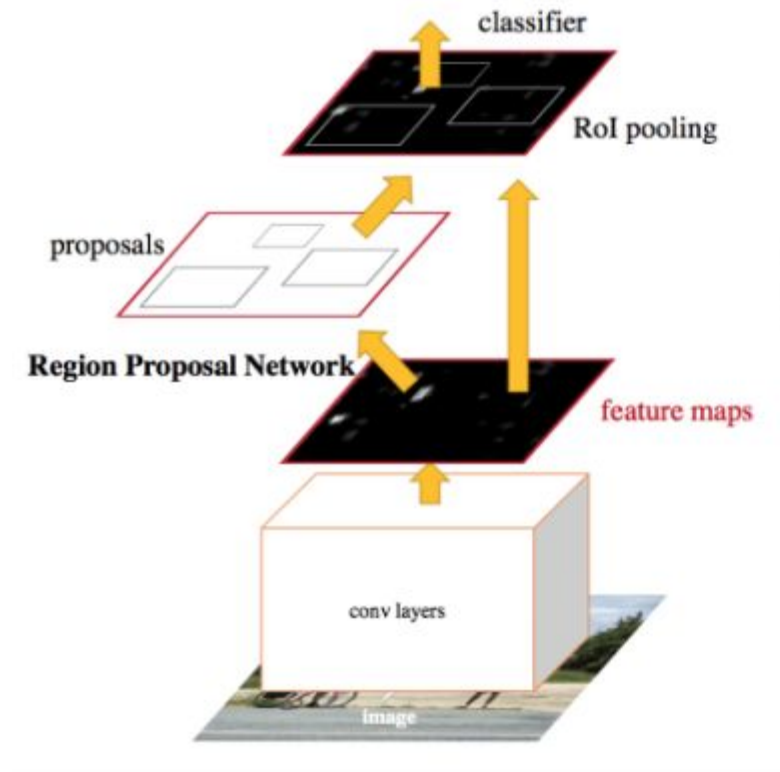


Train

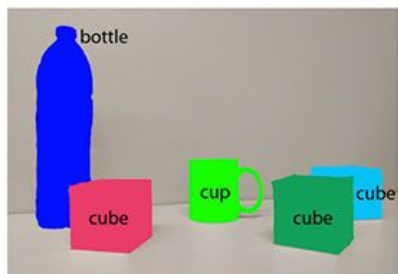
Faster R-CNN



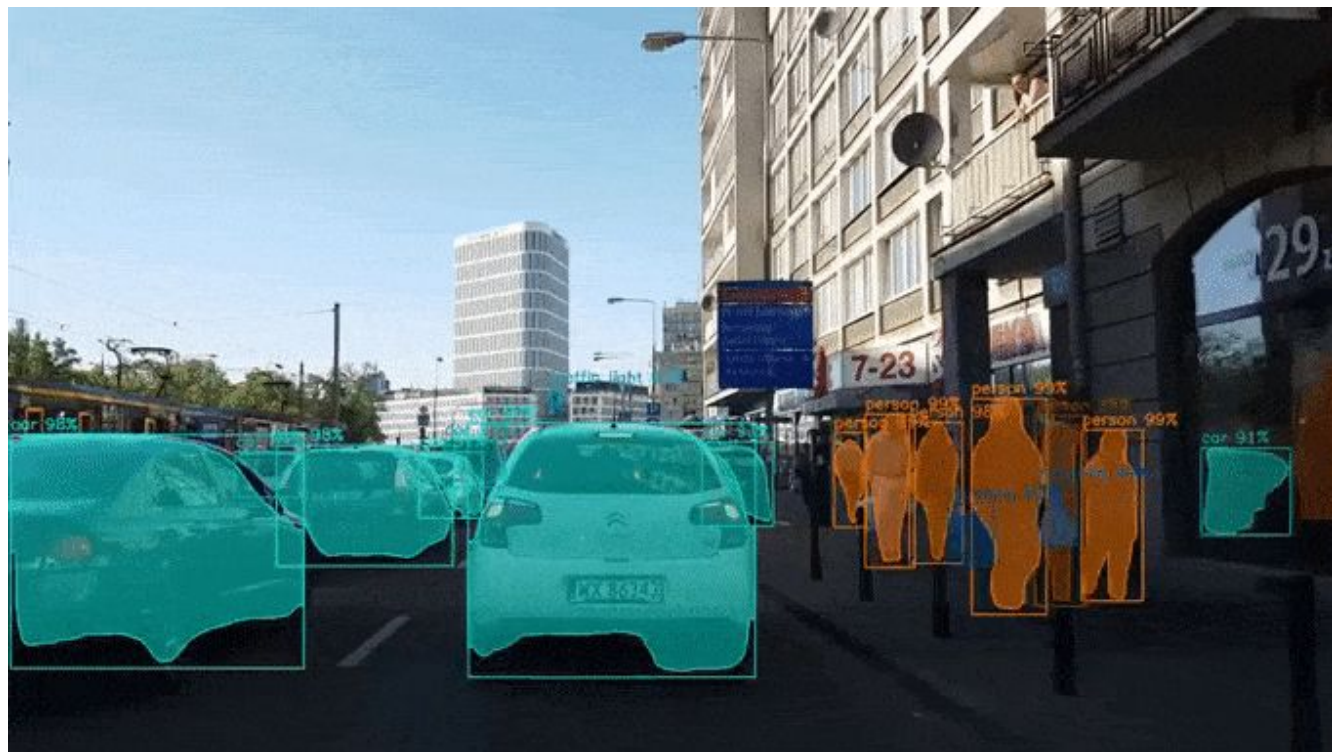
Test



RCNN , Fast RCNN, Faster RCNN, Mask RCNN



(d) Instance segmentation



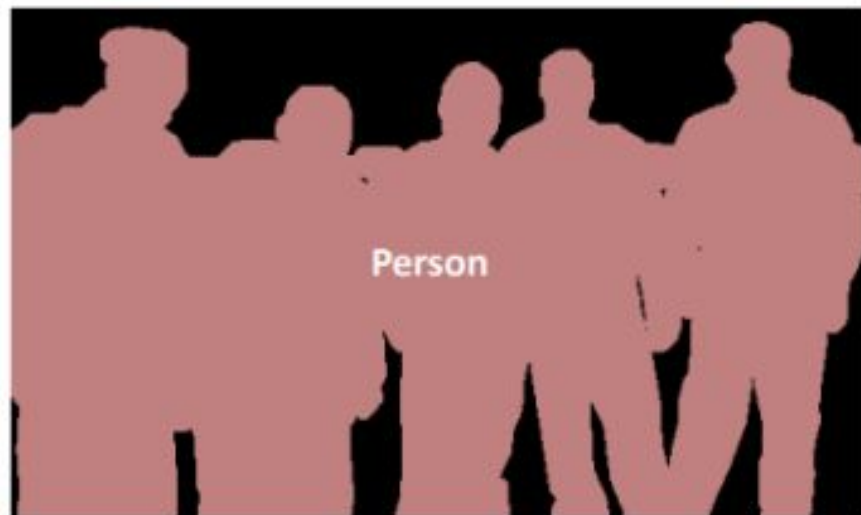


Image 1

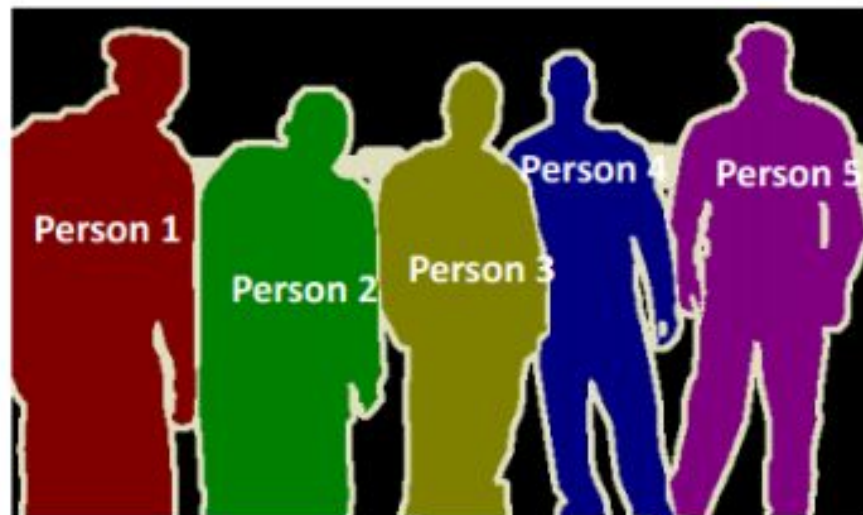
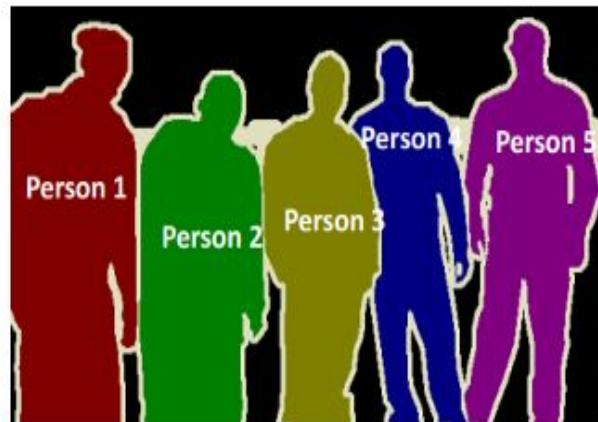
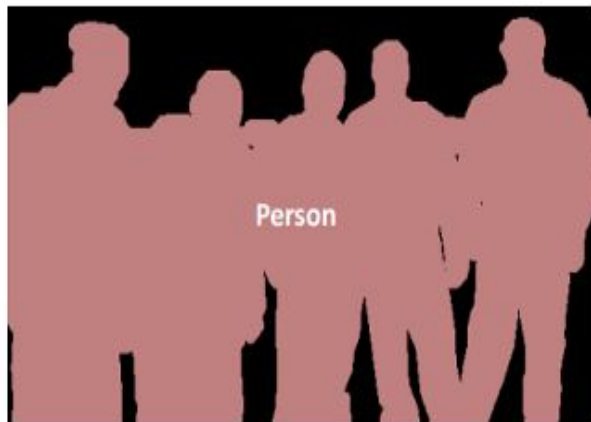
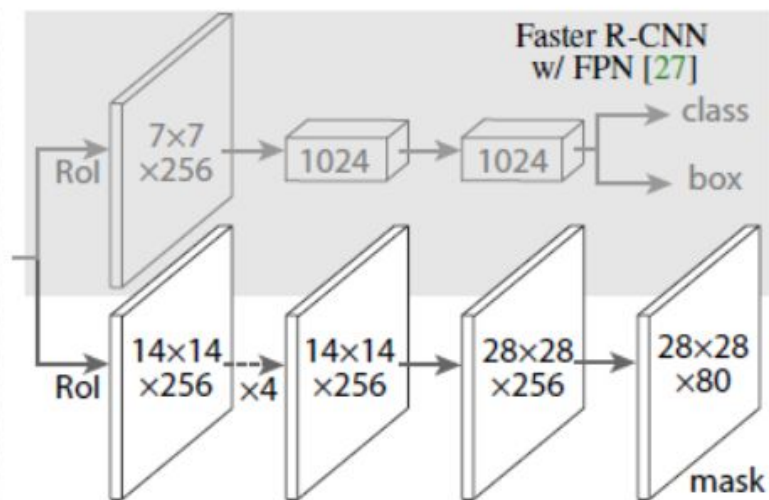
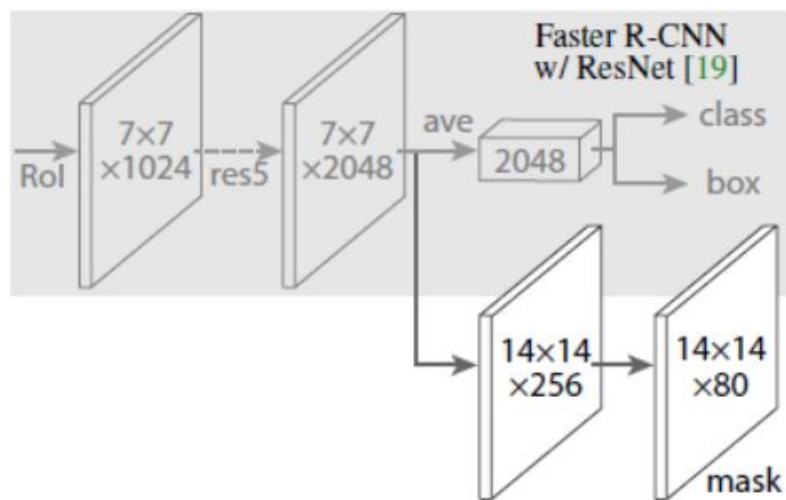
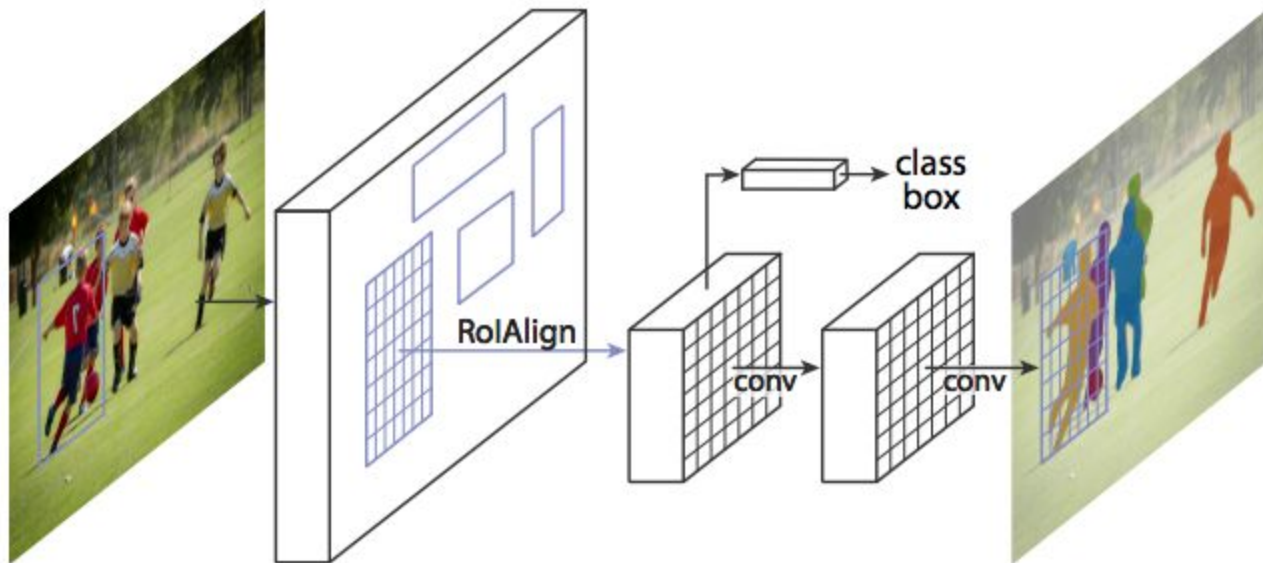


Image 2





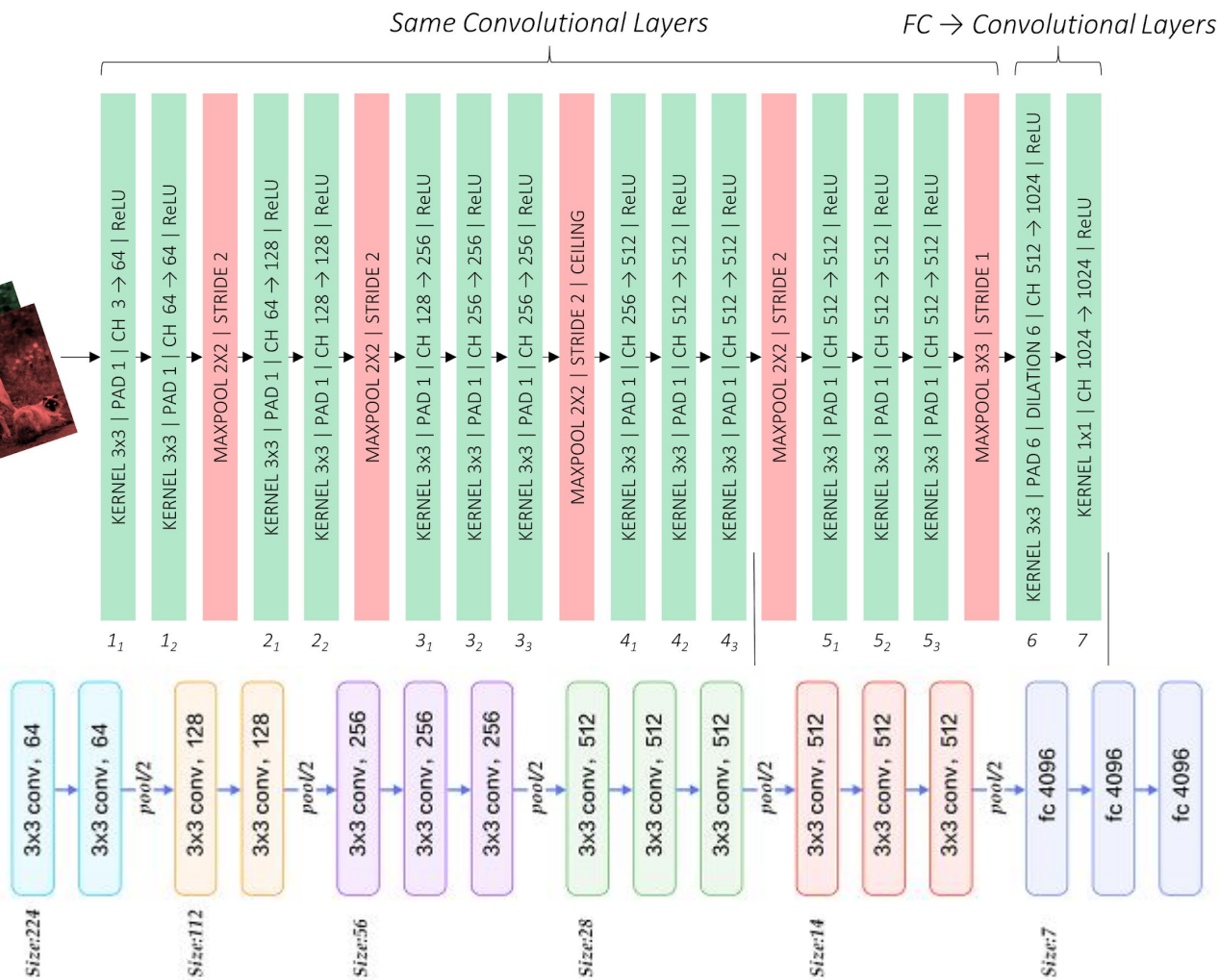
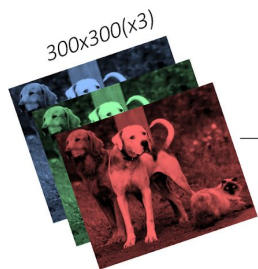


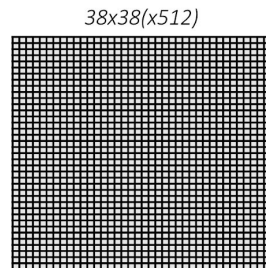
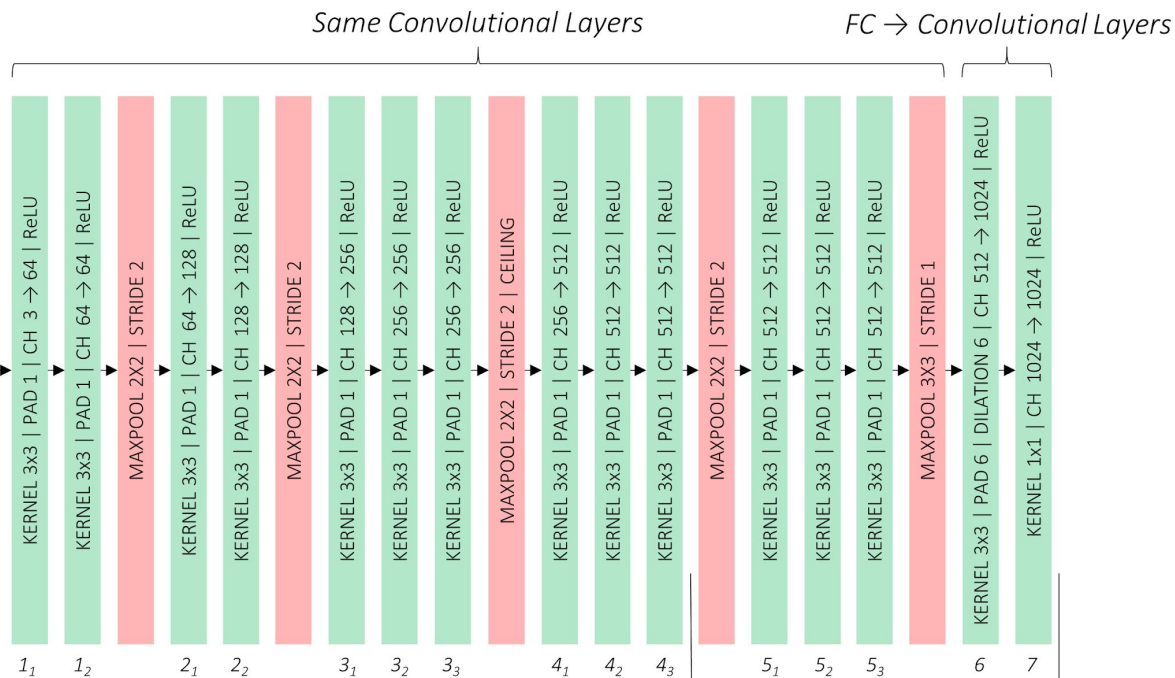
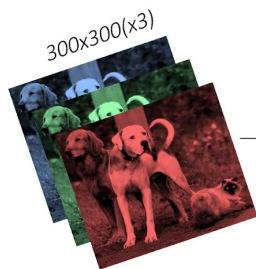
<https://github.com/chenyuntc/simple-faster-rcnn-pytorch>

<https://github.com/facebookresearch/maskrcnn-benchmark>

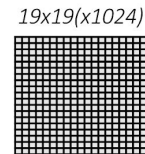
Single ShotmultiBox Detector





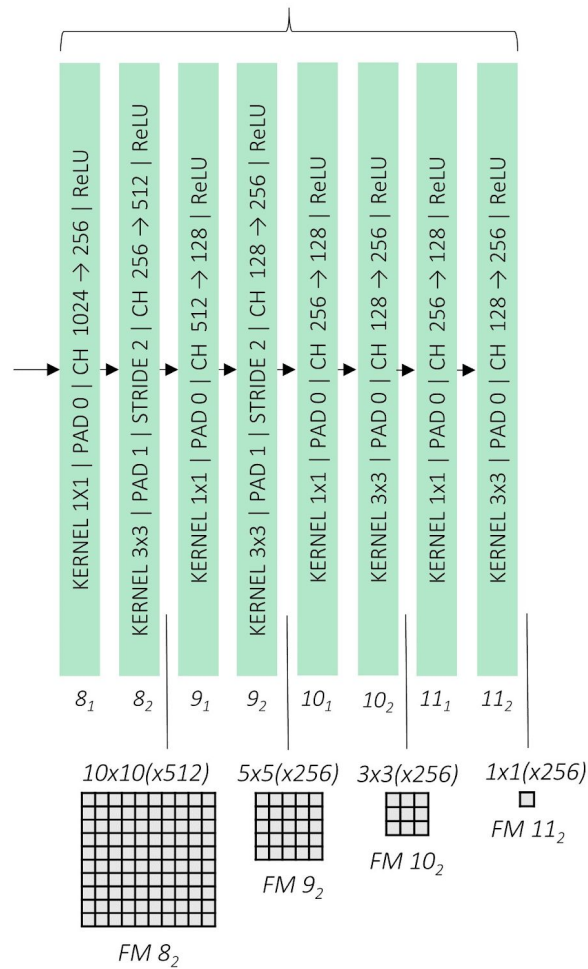
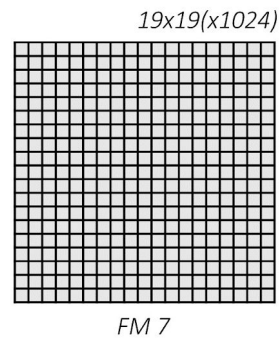


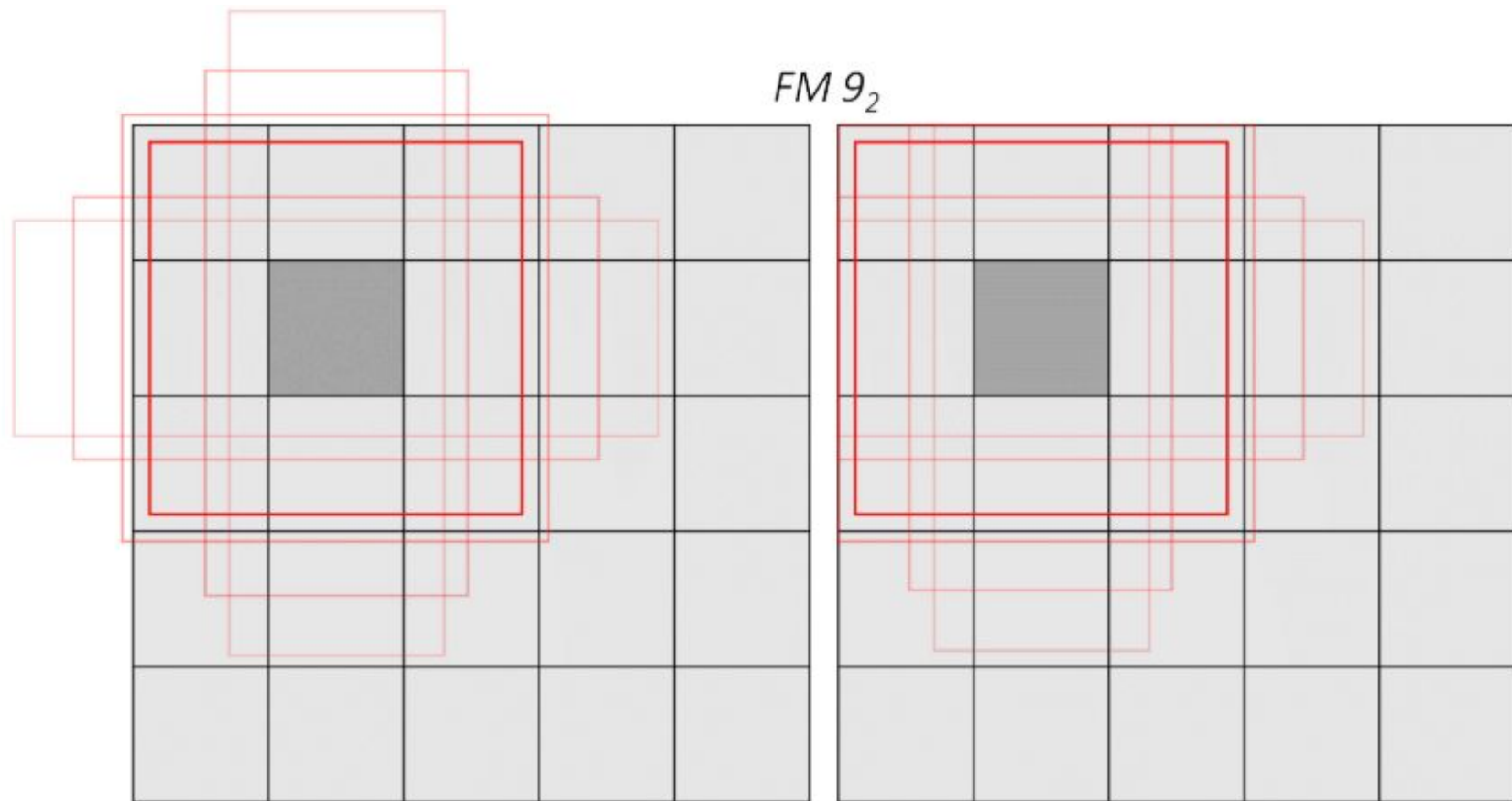
Feature Map (FM) from 4₃



FM 7

Auxiliary Convolutional Layers

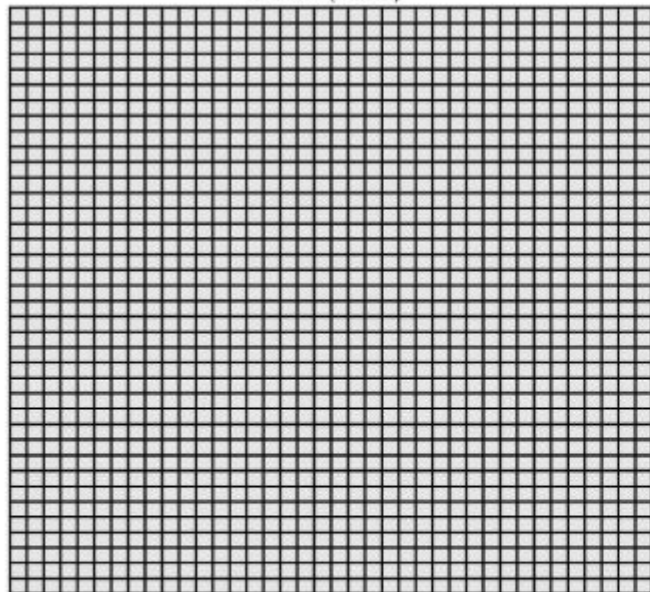




When priors at a location overshoot the edges of the feature map, they are clipped

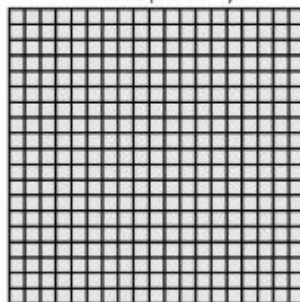
Feature Map From	Feature Map Dimensions	Prior Scale	Aspect Ratios	Number of Priors per Position	Total Number of Priors on this Feature Map
conv4_3	38, 38	0.1	1:1, 2:1, 1:2 + an extra prior	4	5776
conv7	19, 19	0.2	1:1, 2:1, 1:2, 3:1, 1:3 + an extra prior	6	2166
conv8_2	10, 10	0.375	1:1, 2:1, 1:2, 3:1, 1:3 + an extra prior	6	600
conv9_2	5, 5	0.55	1:1, 2:1, 1:2, 3:1, 1:3 + an extra prior	6	150
conv10_2	3, 3	0.725	1:1, 2:1, 1:2 + an extra prior	4	36
conv11_2	1, 1	0.9	1:1, 2:1, 1:2 + an extra prior	4	4
Grand Total	–	–	–	–	8732 priors

$38 \times 38 (x512)$



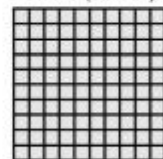
$FM 4_3$

$19 \times 19 (x1024)$



$FM 7$

$10 \times 10 (x512)$



$FM 8_2$

$5 \times 5 (x256)$



$FM 9_2$

$3 \times 3 (x256)$



$FM 10_2$

$1 \times 1 (x256)$



$FM 11_2$

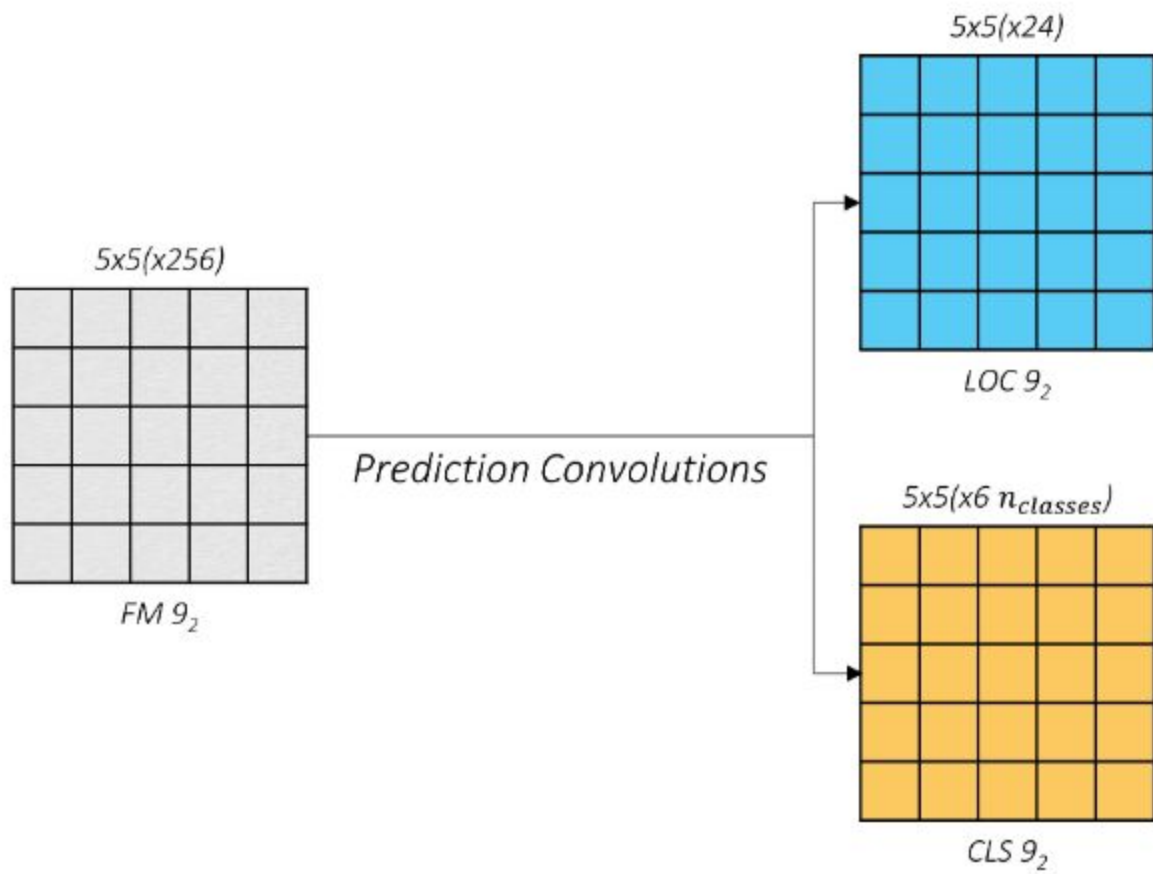
Prediction Convolutions

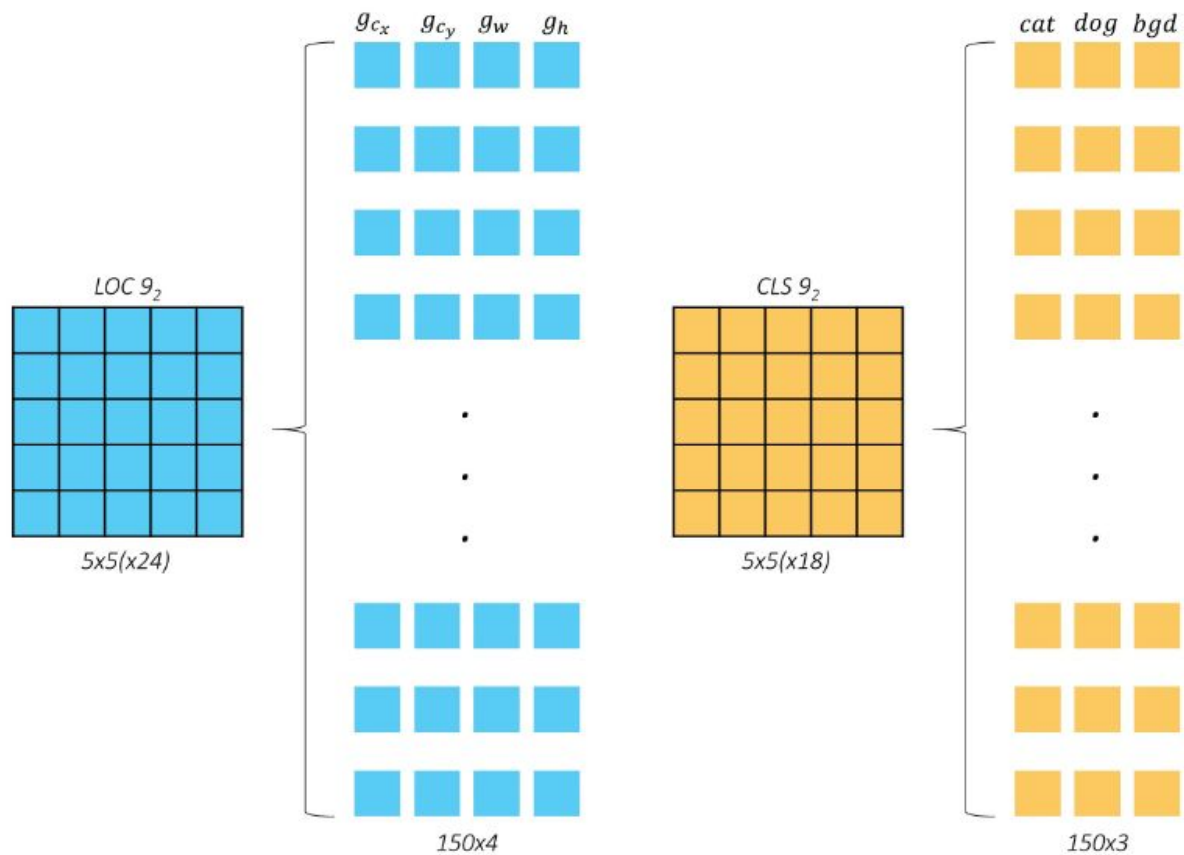
KERNEL 3×3 | PAD 1 |

CH __ \rightarrow number of priors per location for this feature map * 4

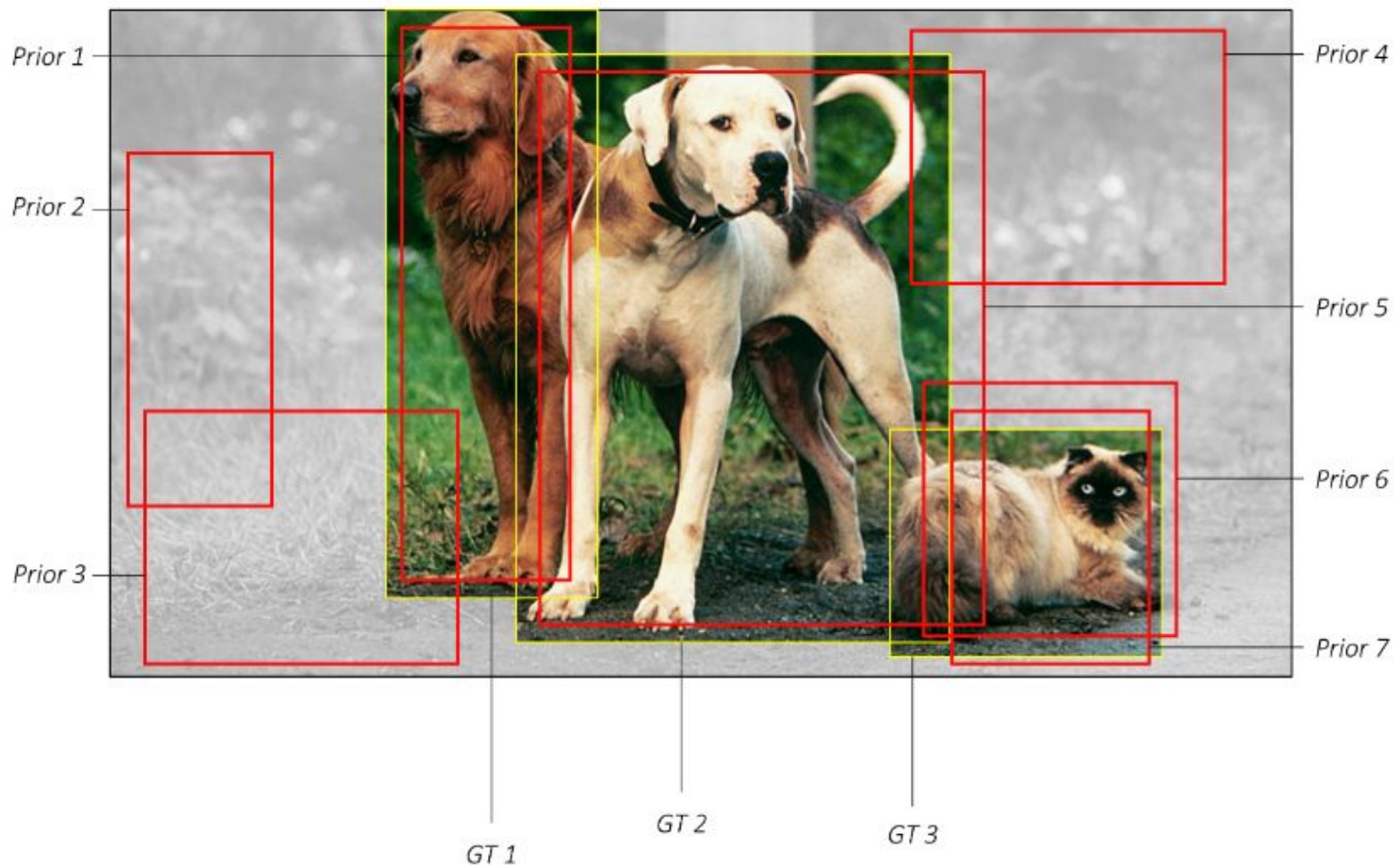
KERNEL 3×3 | PAD 1 |

CH __ \rightarrow number of priors per location for this feature map * number of object classes





Reshape predictions from $FM\ 9_2$ to represent offsets and class scores for the 150 predicted boxes



<i>IoU</i>	<i>GT 1</i>	<i>GT 2</i>	<i>GT 3</i>
<i>Prior 1</i>	0.75	0.07	0
<i>Prior 2</i>	0	0	0
<i>Prior 3</i>	0.05	0	0
<i>Prior 4</i>	0	0.03	0
<i>Prior 5</i>	0.05	0.88	0.06
<i>Prior 6</i>	0	0.03	0.65
<i>Prior 7</i>	0	0	0.68

	<i>Matched GT</i>	<i>Type</i>	<i>Label</i>	<i>Coordinates</i>
<i>Prior 1</i>	GT 1	Positive	<i>dog</i>	of GT 1
<i>Prior 2</i>	GT 1	Negative	<i>bgd</i>	-
<i>Prior 3</i>	GT 1	Negative	<i>bgd</i>	-
<i>Prior 4</i>	GT 2	Negative	<i>bgd</i>	-
<i>Prior 5</i>	GT 2	Positive	<i>dog</i>	of GT 2
<i>Prior 6</i>	GT 3	Positive	<i>cat</i>	of GT 3
<i>Prior 7</i>	GT 3	Positive	<i>cat</i>	of GT 3