

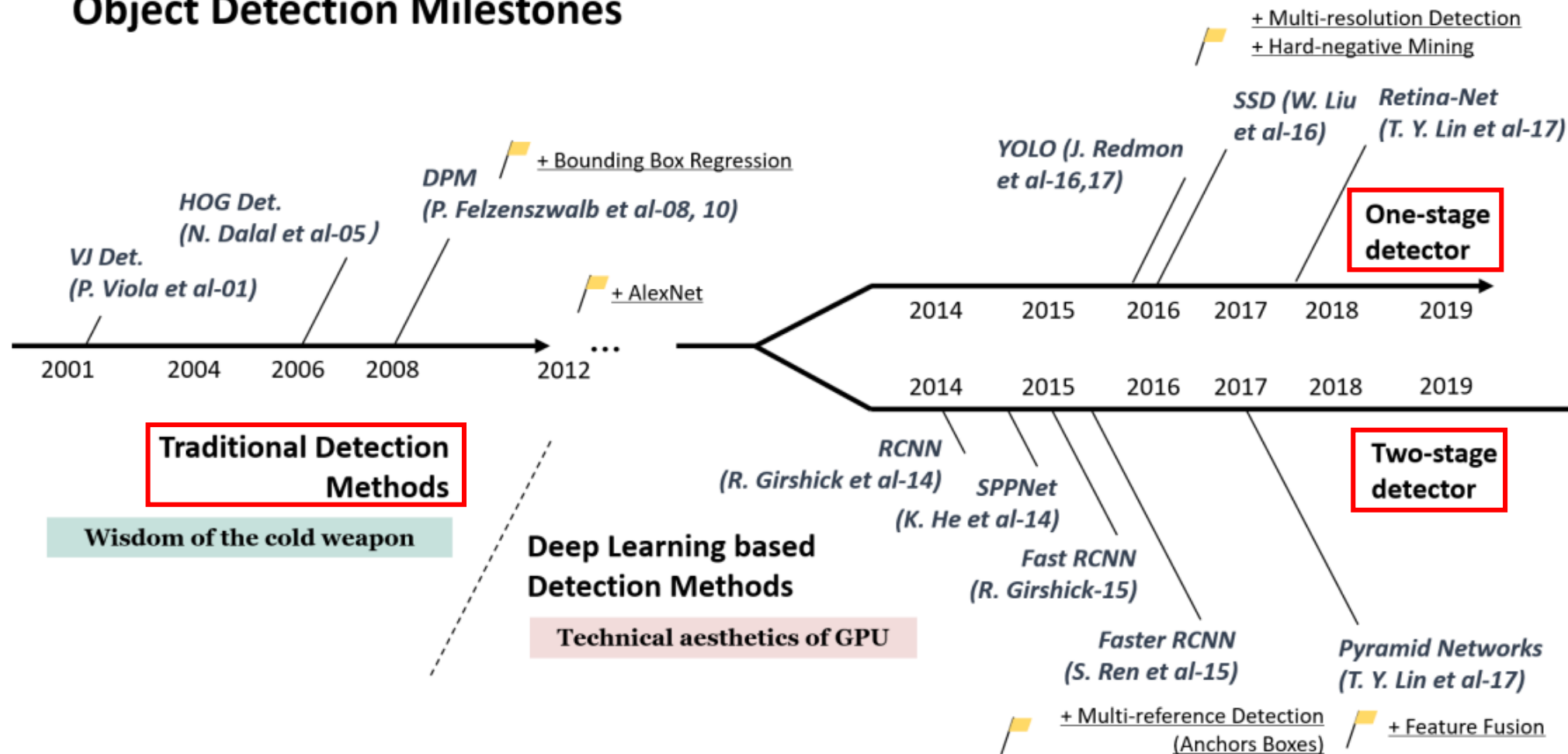
객체 검출 I

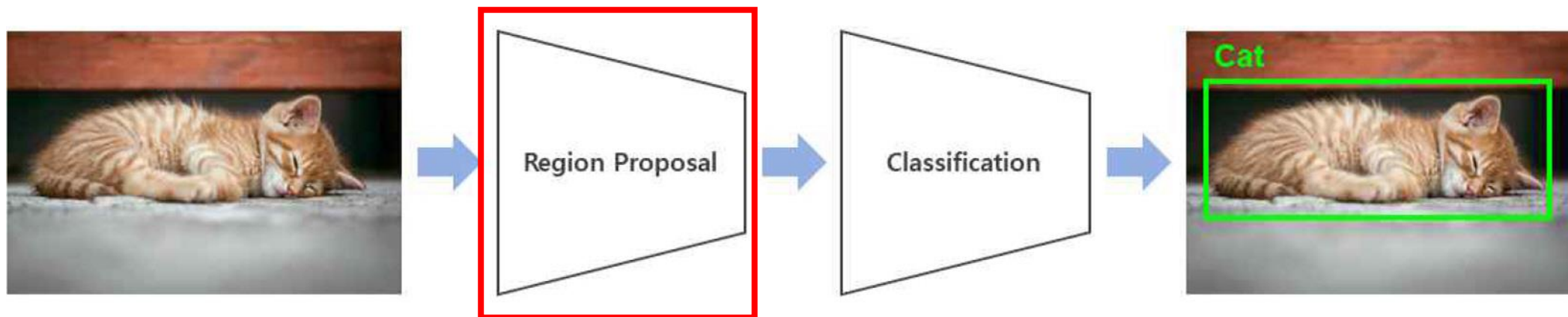
(Two stage detector)

강사: 김 남 범 교수

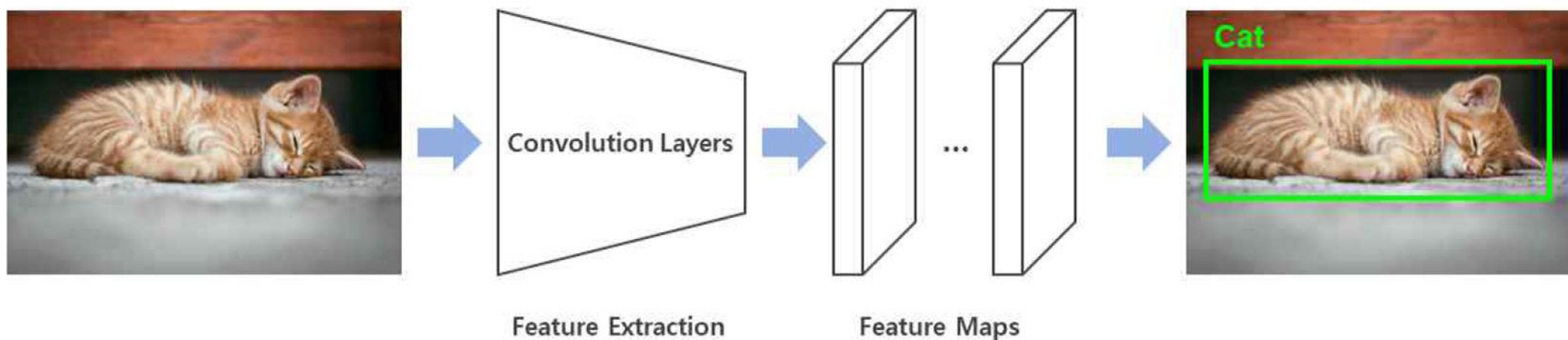
Object detection Milestones

Object Detection Milestones





(a) 2-Stage detector



(b) 1-Stage detector

Stage에 따른 모델분류

- Two stage detector
 - R-CNN (2014)
 - Fast R-CNN
 - Faster R-CNN
 - Mask R-CNN (Instance segmentation)
- One stage detector
 - Yolo series
 - Yolo1 – Yolo11 (Object detection)
 - SSD series
 - SSD
 - RetinaNet

R-CNN: *Regions with CNN features*

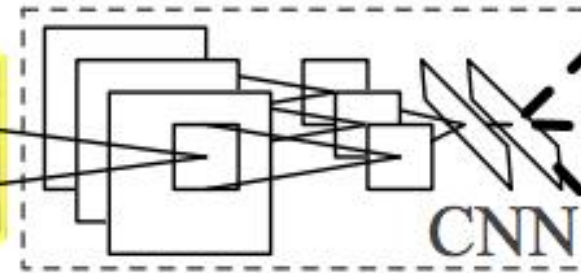


1. Input image

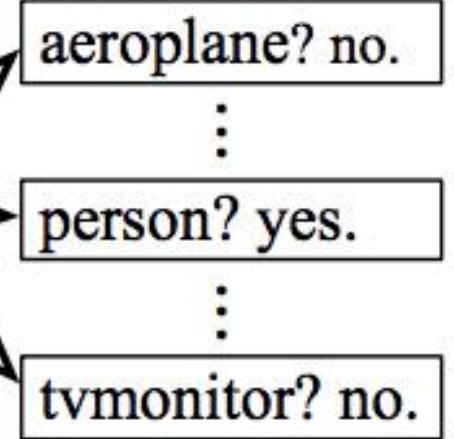


2. Extract region proposals (~2k)

warped region

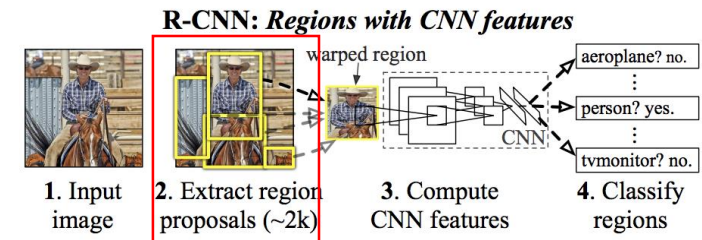


3. Compute CNN features



4. Classify regions

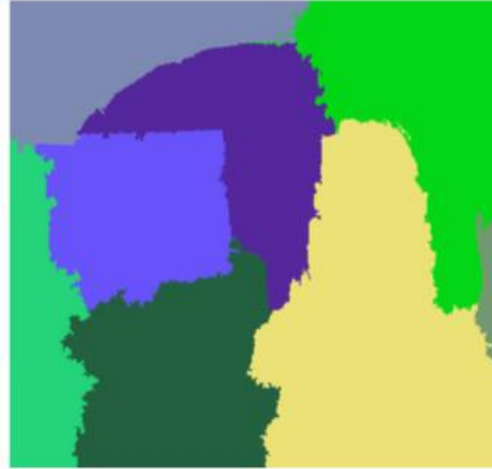
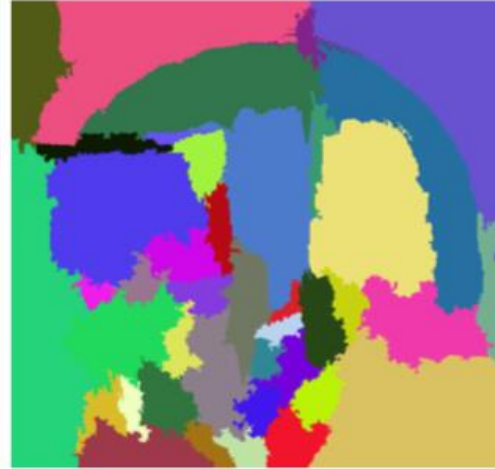
R-CNN: Region proposal



1

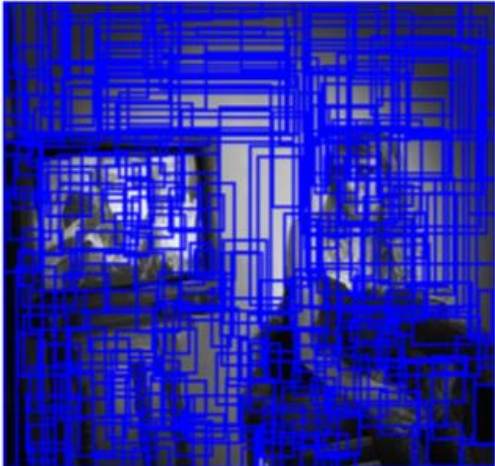


2

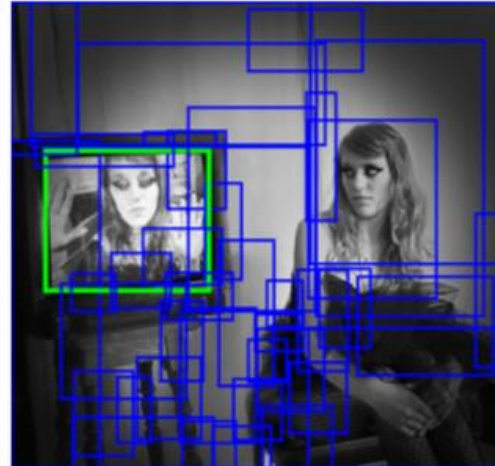


Selective search:

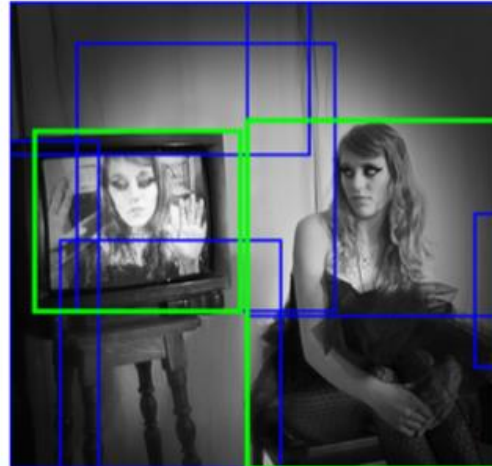
1. Image segmentation
(Graph-based segmentation)
 - felzenszwalb segmentation
2. Merge regions
 - Greedy algorithm
3. Candidate object location (~2)



Initial segmentation

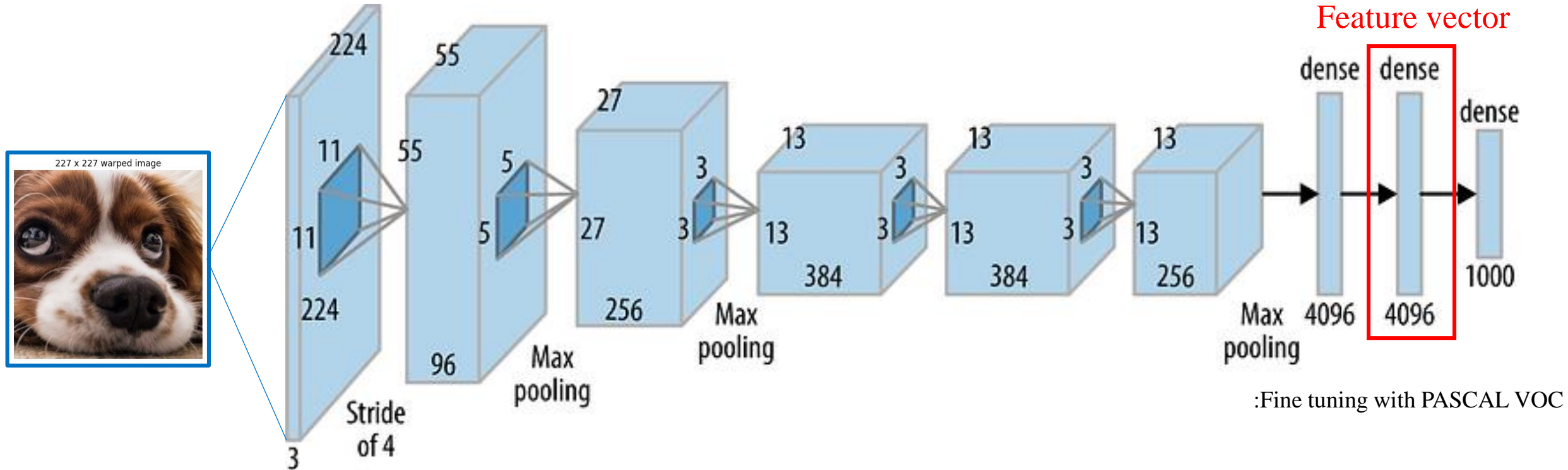
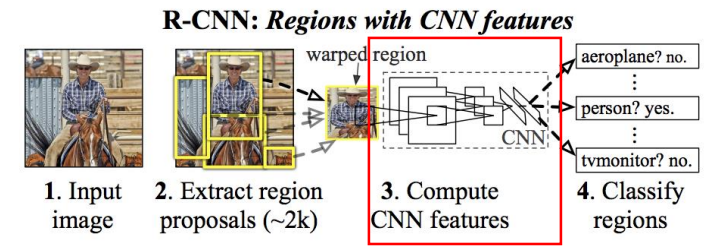


After some iterations



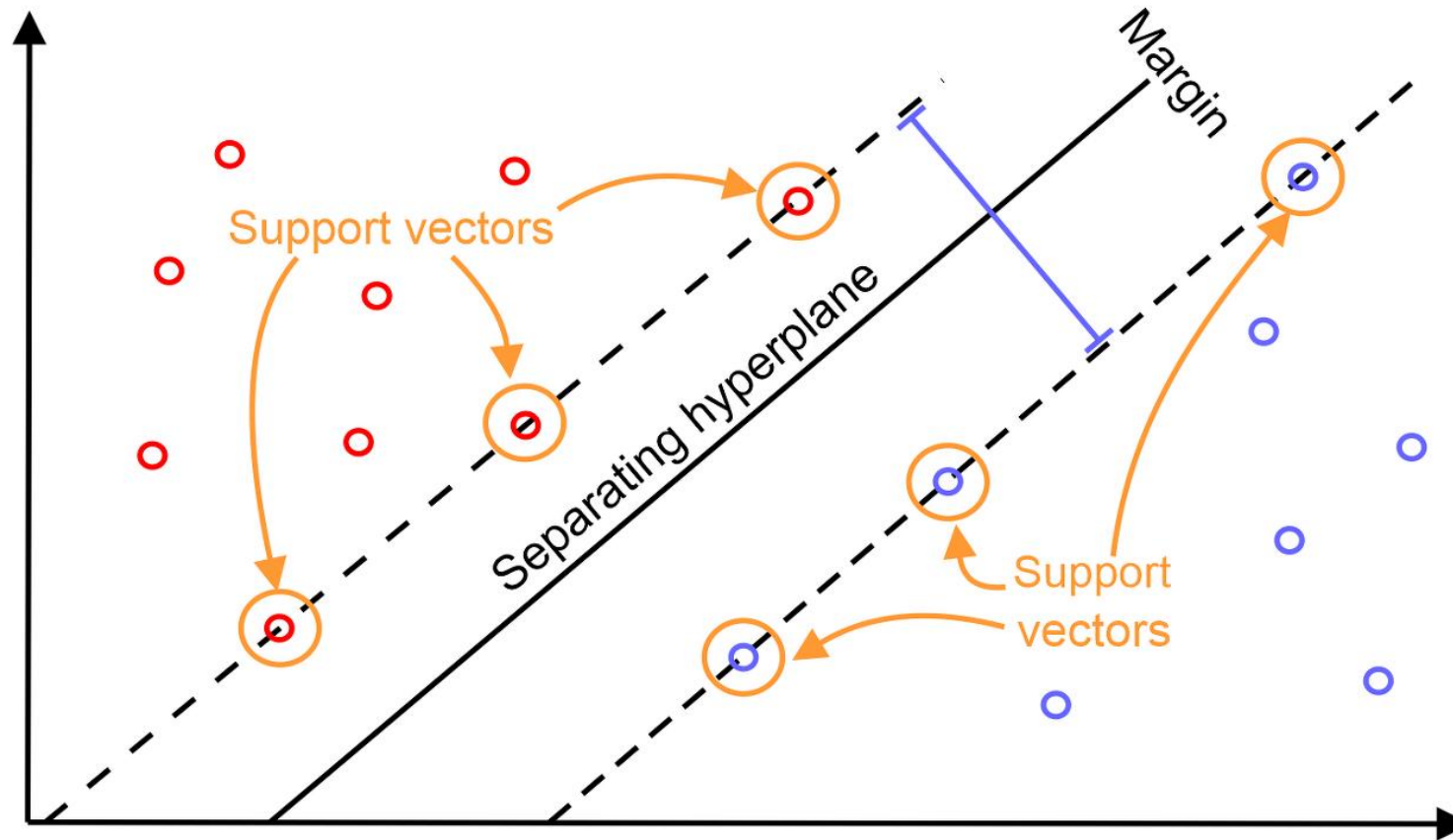
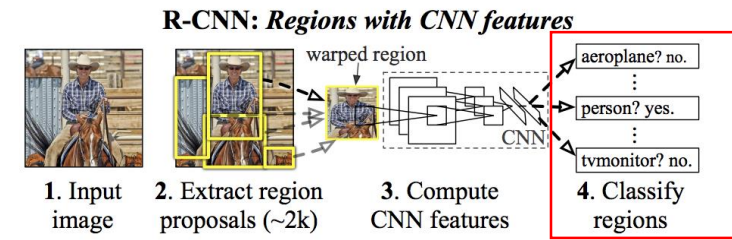
After more iterations

R-CNN: Feature extraction

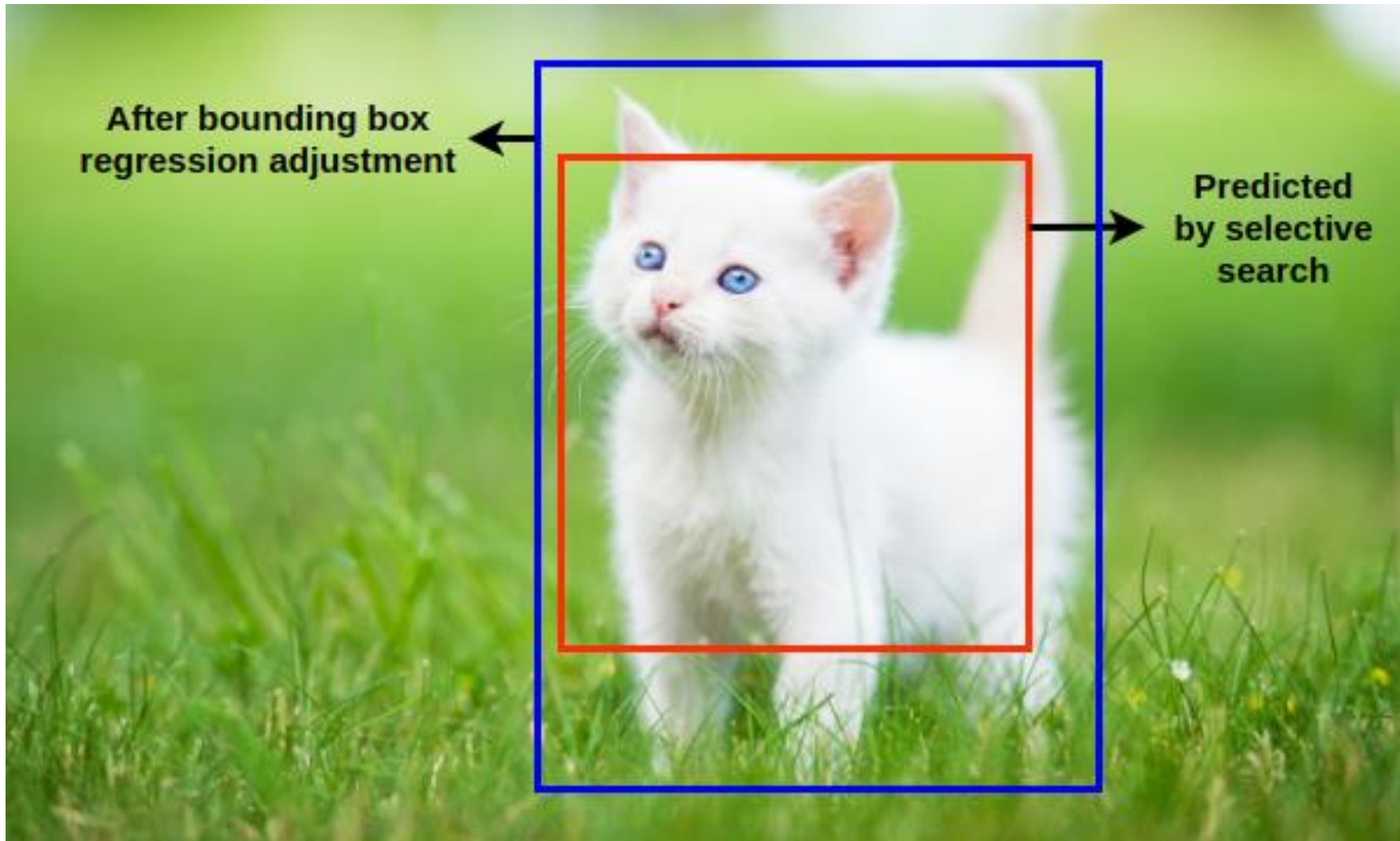
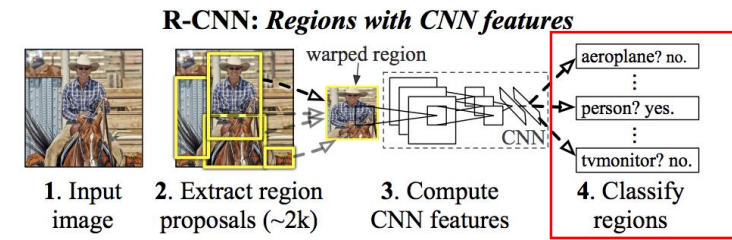


Alexnet (2012)

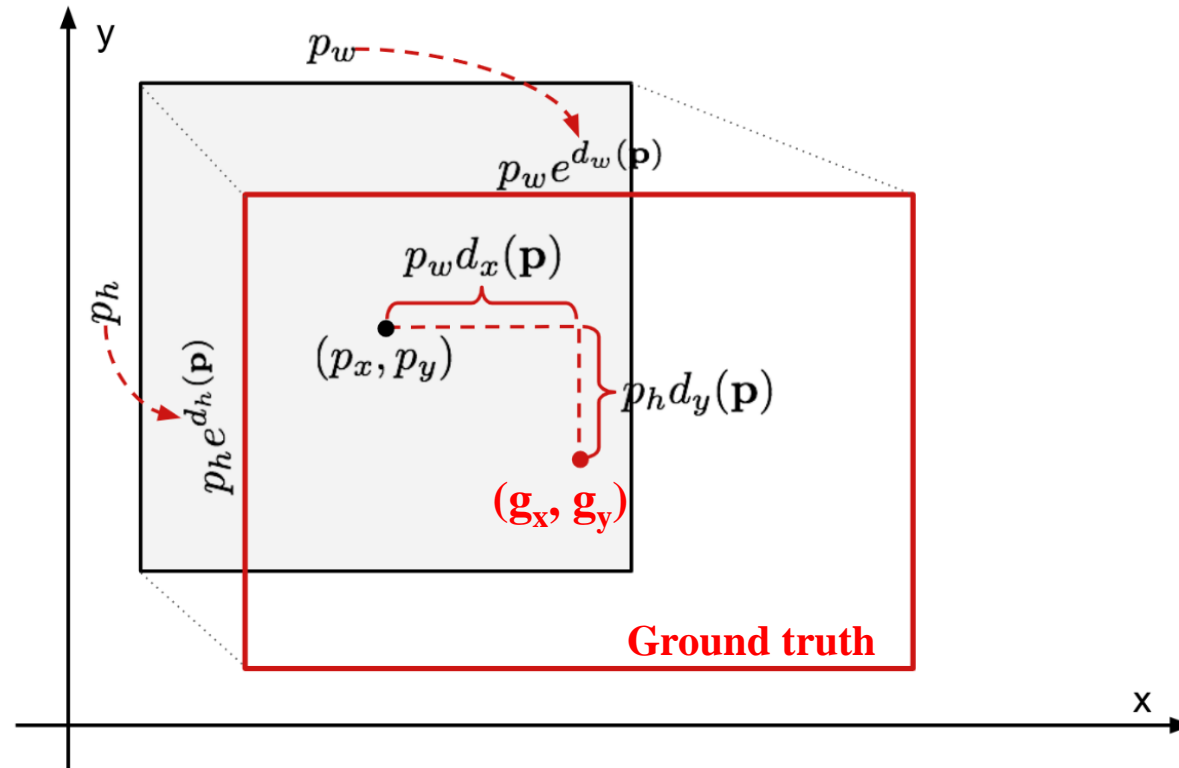
R-CNN: SVM Classification



R-CNN: Bounding box regression



Bounding box regression



$$\hat{g}_x = p_w d_x(\mathbf{p}) + p_x$$

$$\hat{g}_y = p_h d_y(\mathbf{p}) + p_y$$

$$\hat{g}_w = p_w \exp(d_w(\mathbf{p}))$$

$$\hat{g}_h = p_h \exp(d_h(\mathbf{p}))$$

$$t_x = (g_x - p_x) / p_w$$

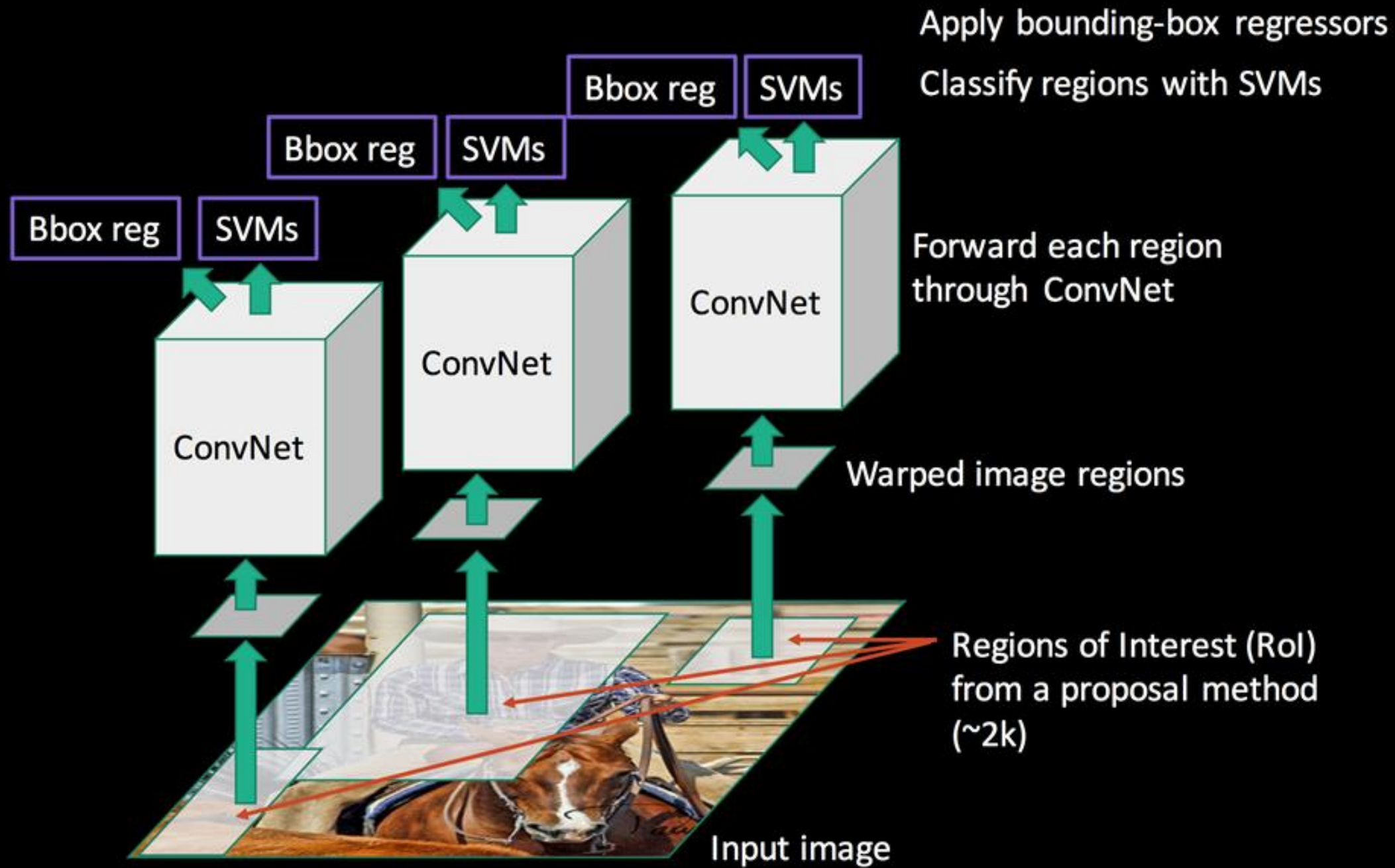
$$t_y = (g_y - p_y) / p_h$$

$$t_w = \log(g_w / p_w)$$

$$t_h = \log(g_h / p_h)$$

Minimizing SSE loss =

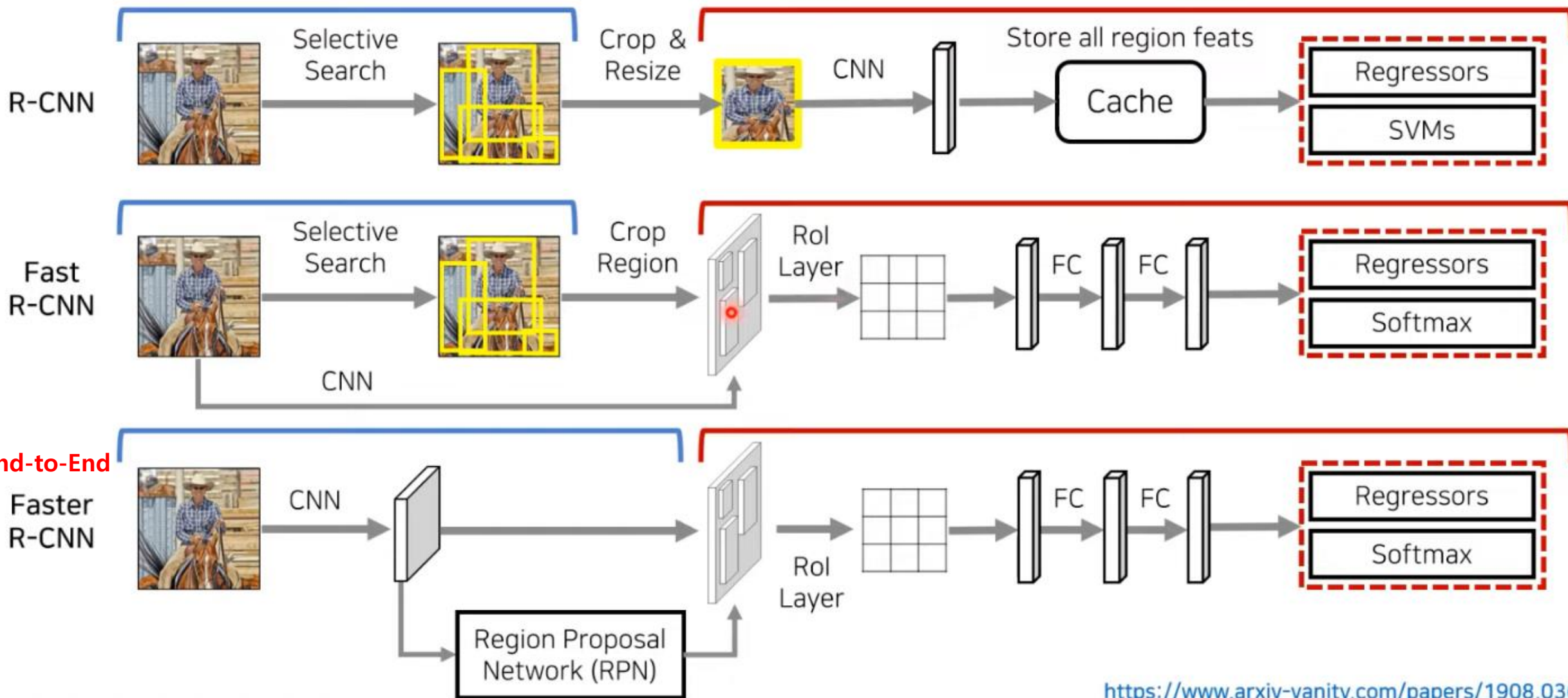
$$\mathcal{L}_{\text{reg}} = \sum_{i \in \{x, y, w, h\}} (t_i - d_i(\mathbf{p}))^2 + \lambda \|\mathbf{w}\|^2$$

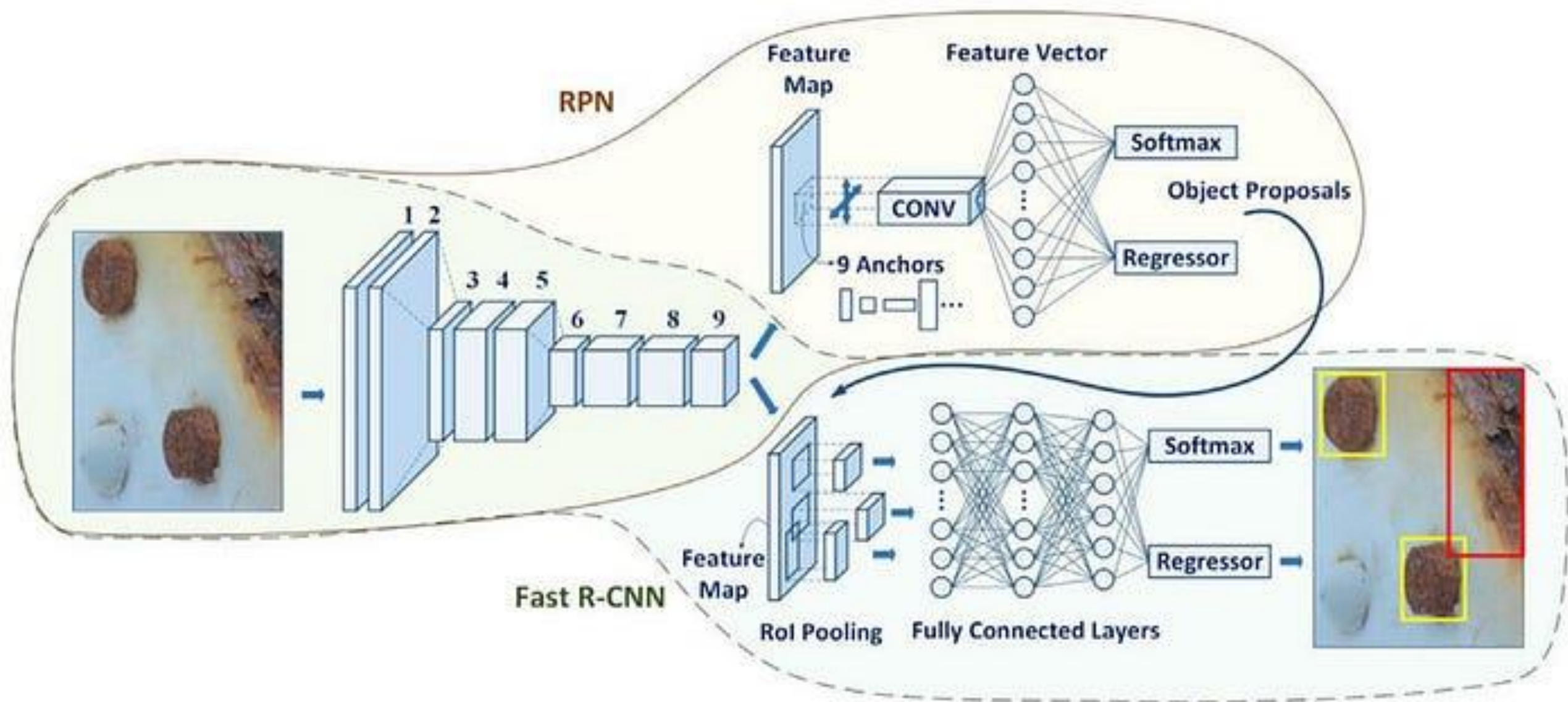


R-CNN의 장단점

- 높은 Detection 정확도
- 2000개씩 생성된 region 이미지를 CNN Feature map 생성
 - 너무 느림
 - 이미지 한 장당 50초 소요
- End-to-End 학습이 안됨
 - CNN Feature Extractor + SVM and Bounding box regressor

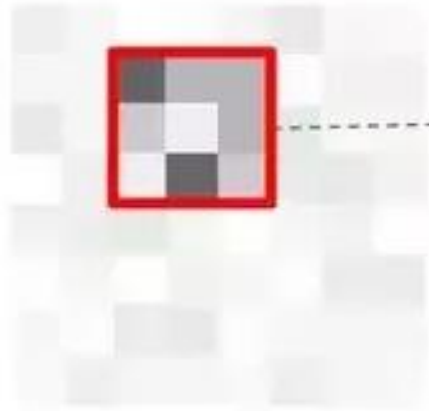
R-CNN, Fast R-CNN, Faster R-CNN



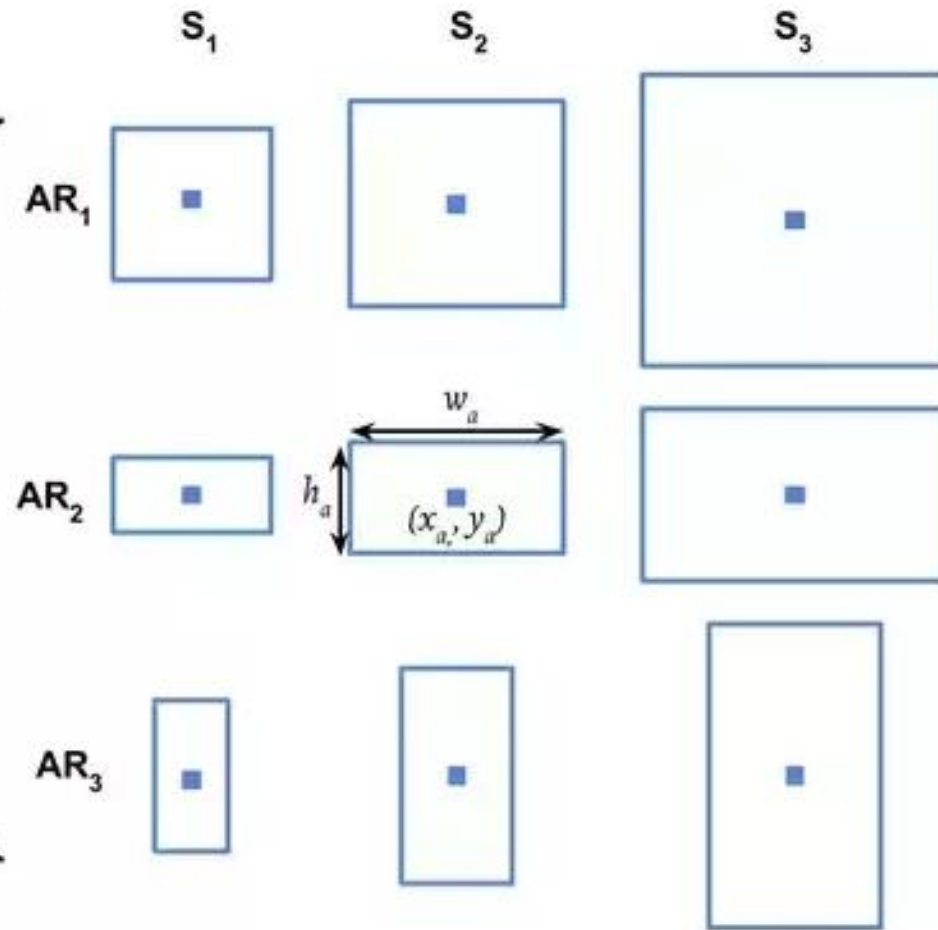


RPN (Region proposal network)

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



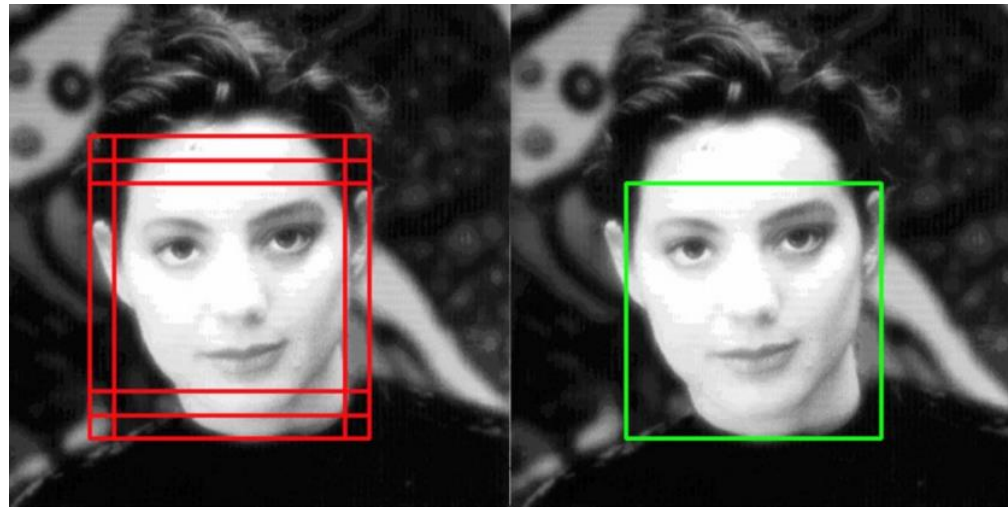
w_a : anchor's width
 h_a : anchor's height
 x_a, y_a : anchor's center



- Anchor target generation
- Calculate the IOU of GT boxes
- Proposal generation

@vmirly

Non Maximum suppression

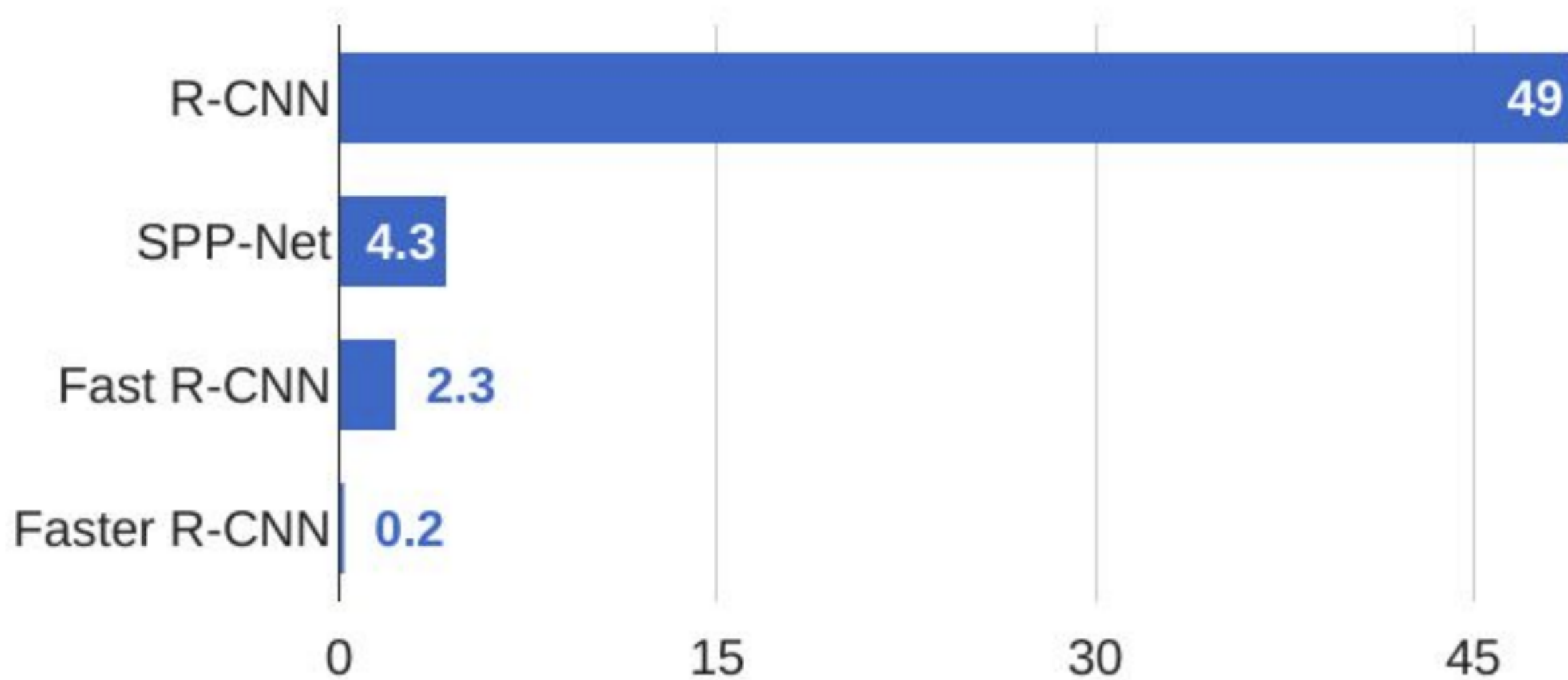


R-CNN 계열 성능비교

System	Time	07 data	07 + 12 data
R-CNN	~ 50s	66.0	-
Fast R-CNN	~ 2s	66.9	70.0
Faster R-CNN	~ 198ms	69.9	73.2

Detection mAP on PASCAL VOC 2007 and 2012, with VGG-16 pre-trained on ImageNet Dataset

R-CNN Test-Time Speed

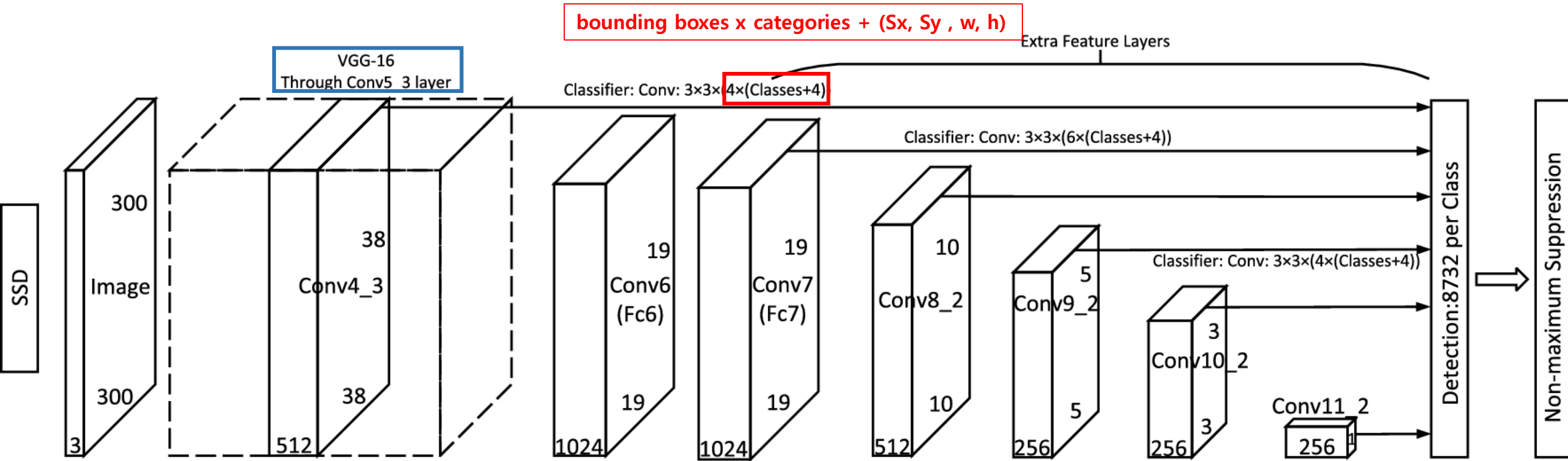


객체 검출 Ⅱ

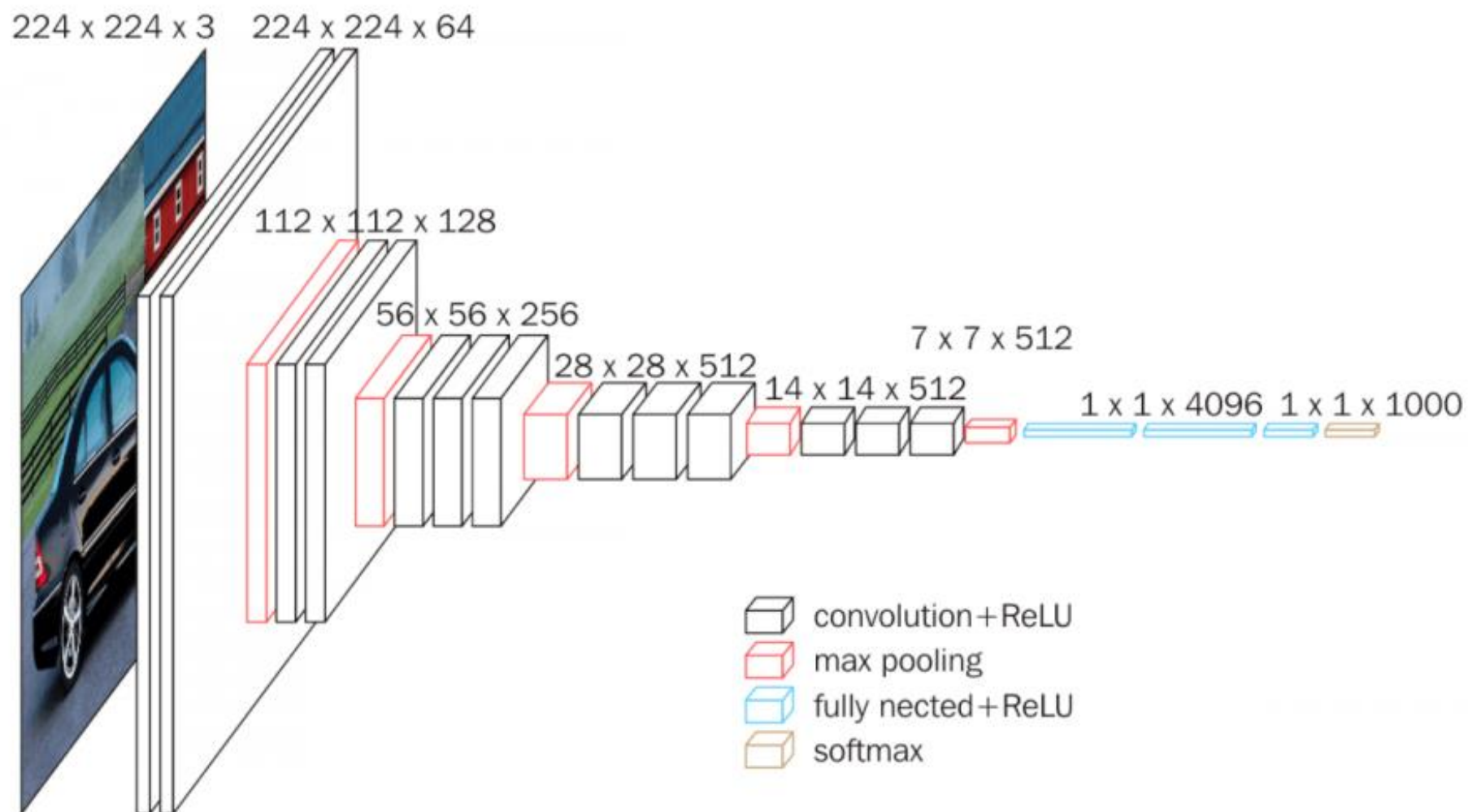
(Real time detector)

강사: 김 남 범 교수

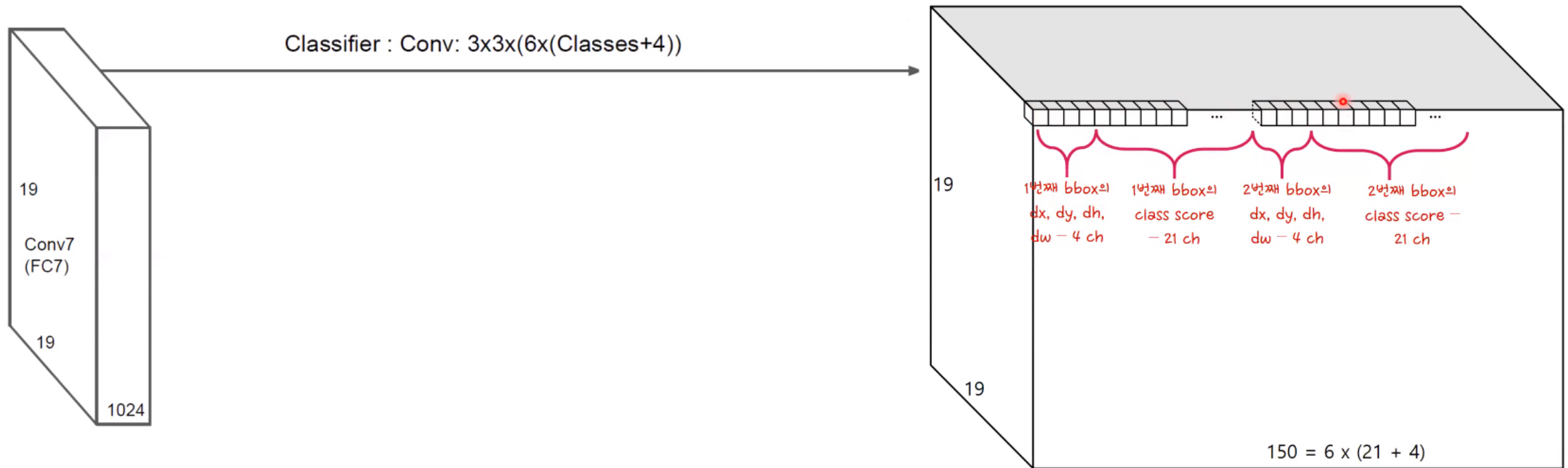
Single shot multibox detector (SSD, 2016)



VGG16 (2014)



Convolutional predictors



Jaccard overlap

True Positive



IoU = 0.922

False Positive



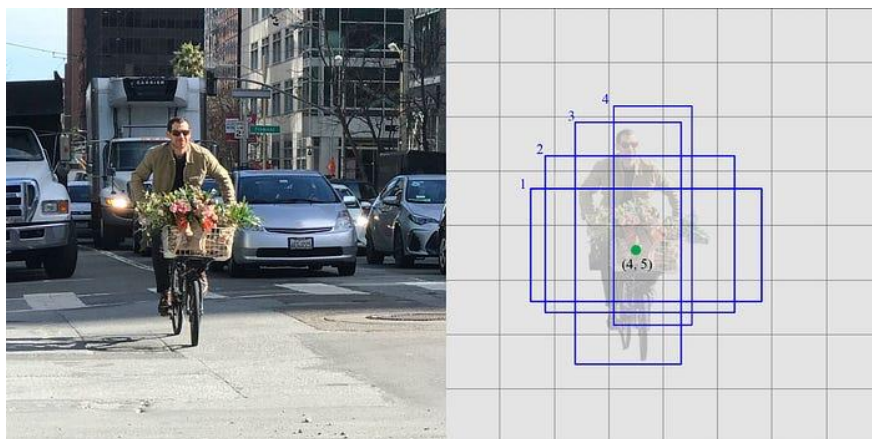
IoU = 0.258

False Negative

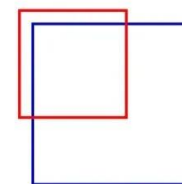
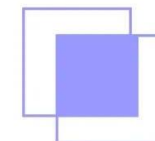


IoU = 0.00

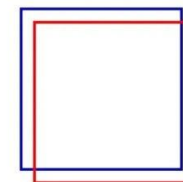
IoU Threshold = 0.5



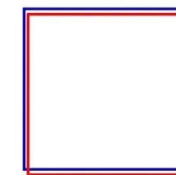
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



Poor

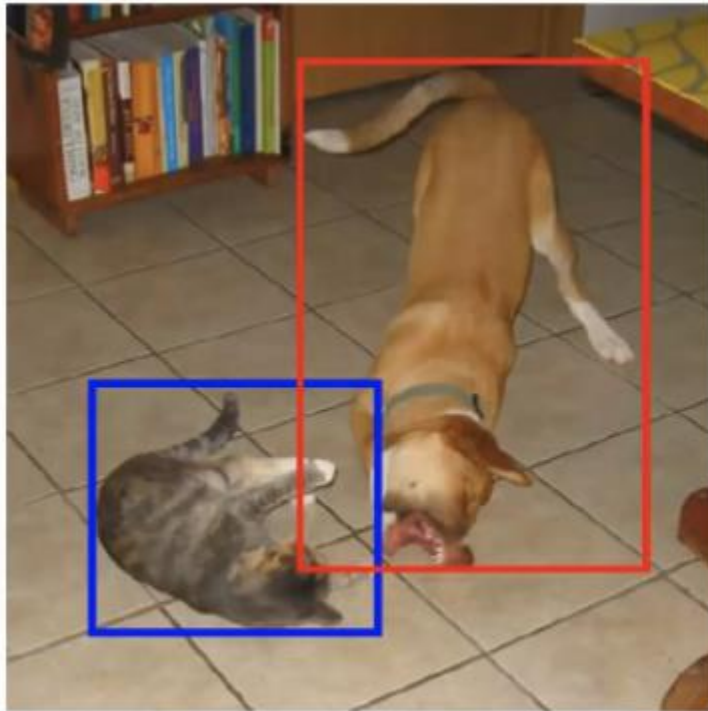


Good

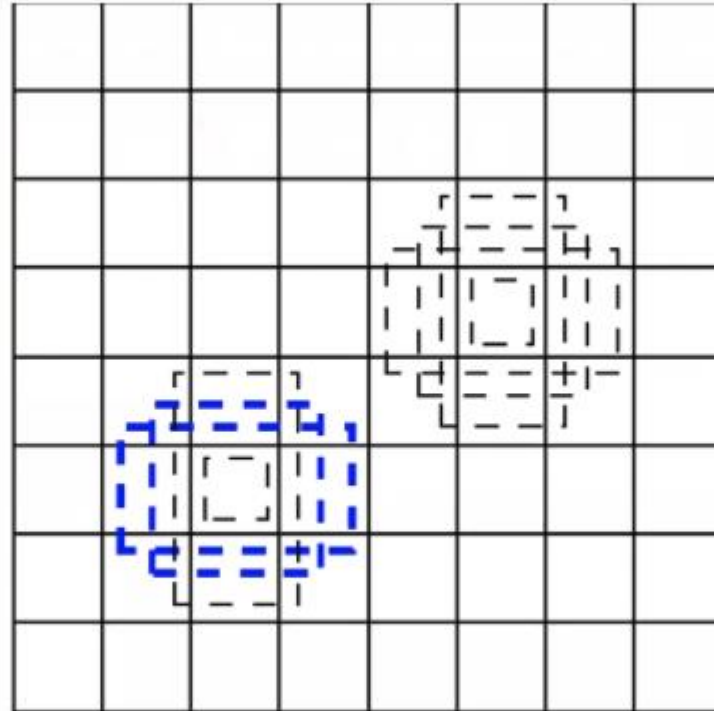


Excellent

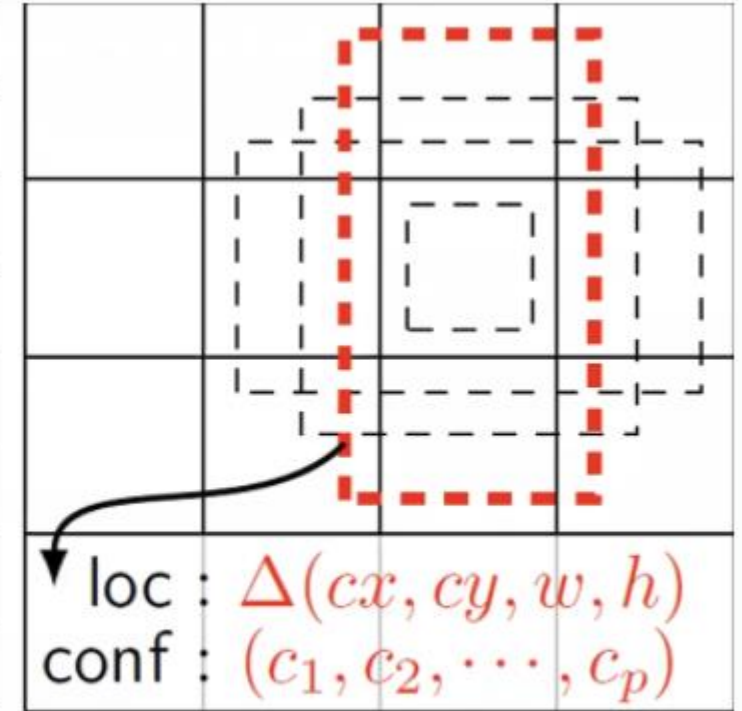
Scales and aspect ratios for default boxes



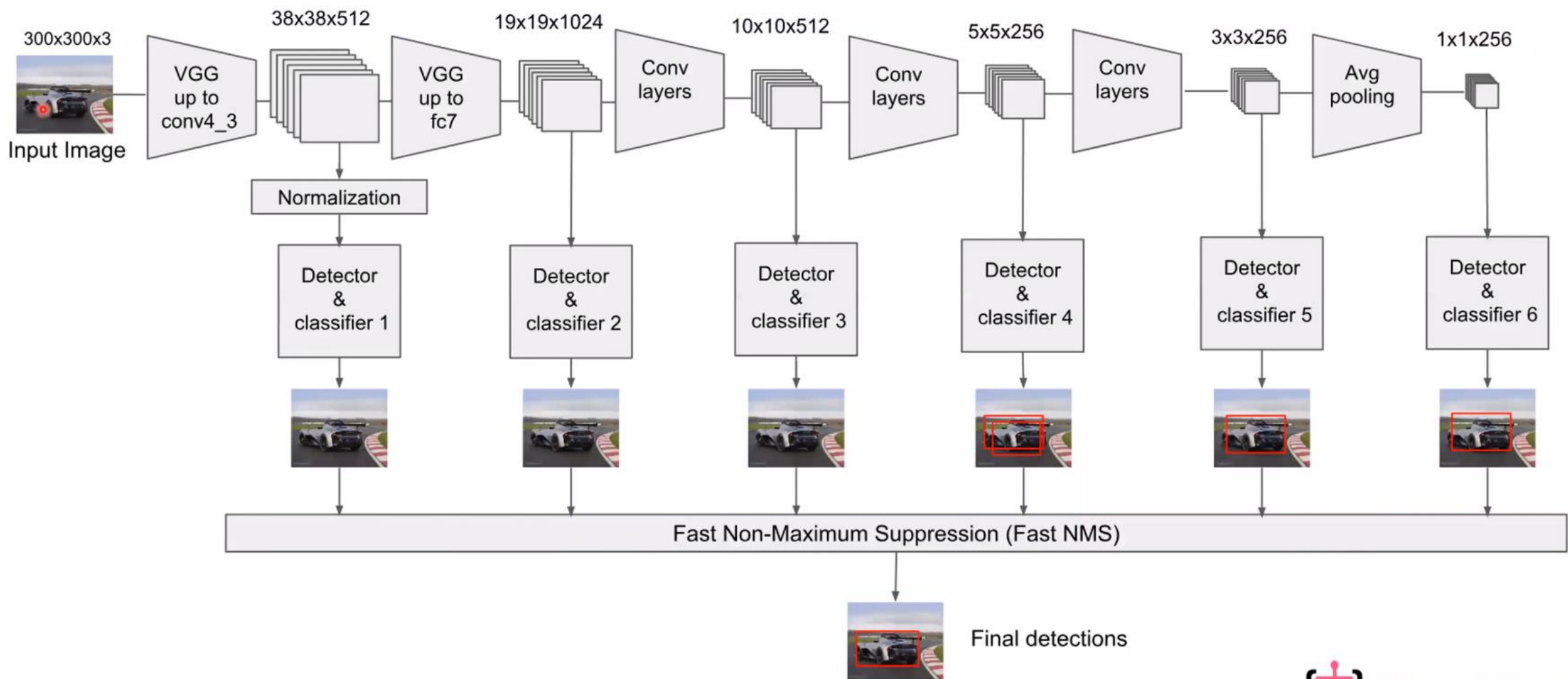
(a) Image with GT boxes

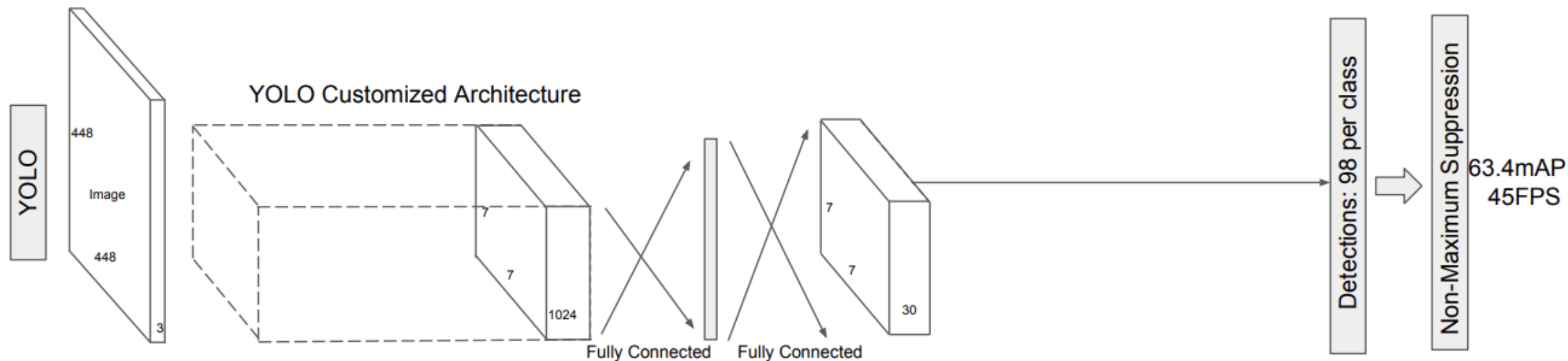
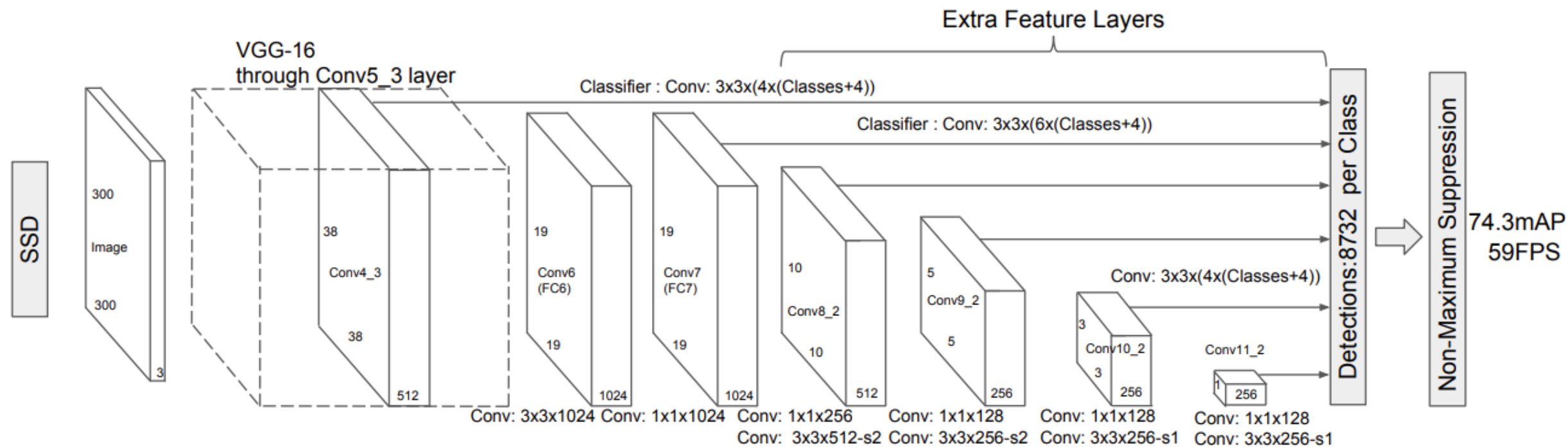


(b) 8×8 feature map



(c) 4×4 feature map





Yolo architecture (You Only Look Once, 2015)

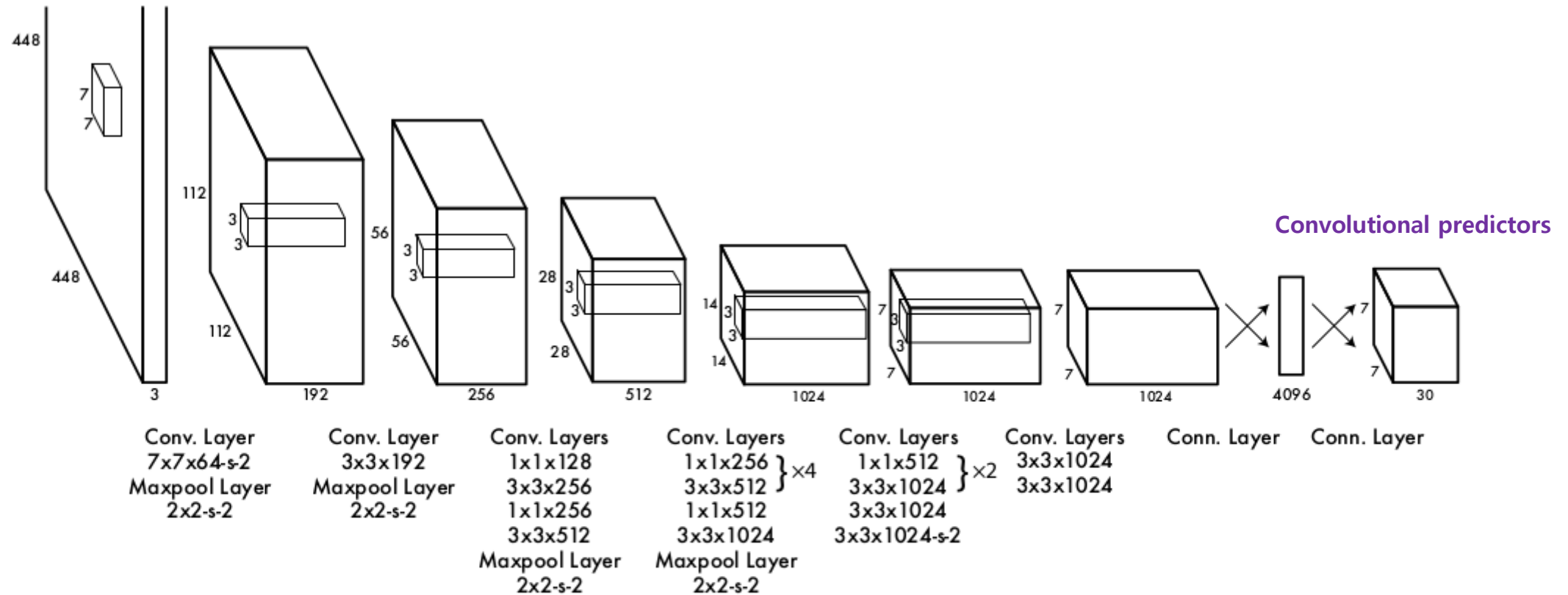


Figure 3: The Architecture. Our detection network has 24 convolutional layers followed by 2 fully connected layers. Alternating 1×1 convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution (224×224 input image) and then double the resolution for detection.

Unified detection

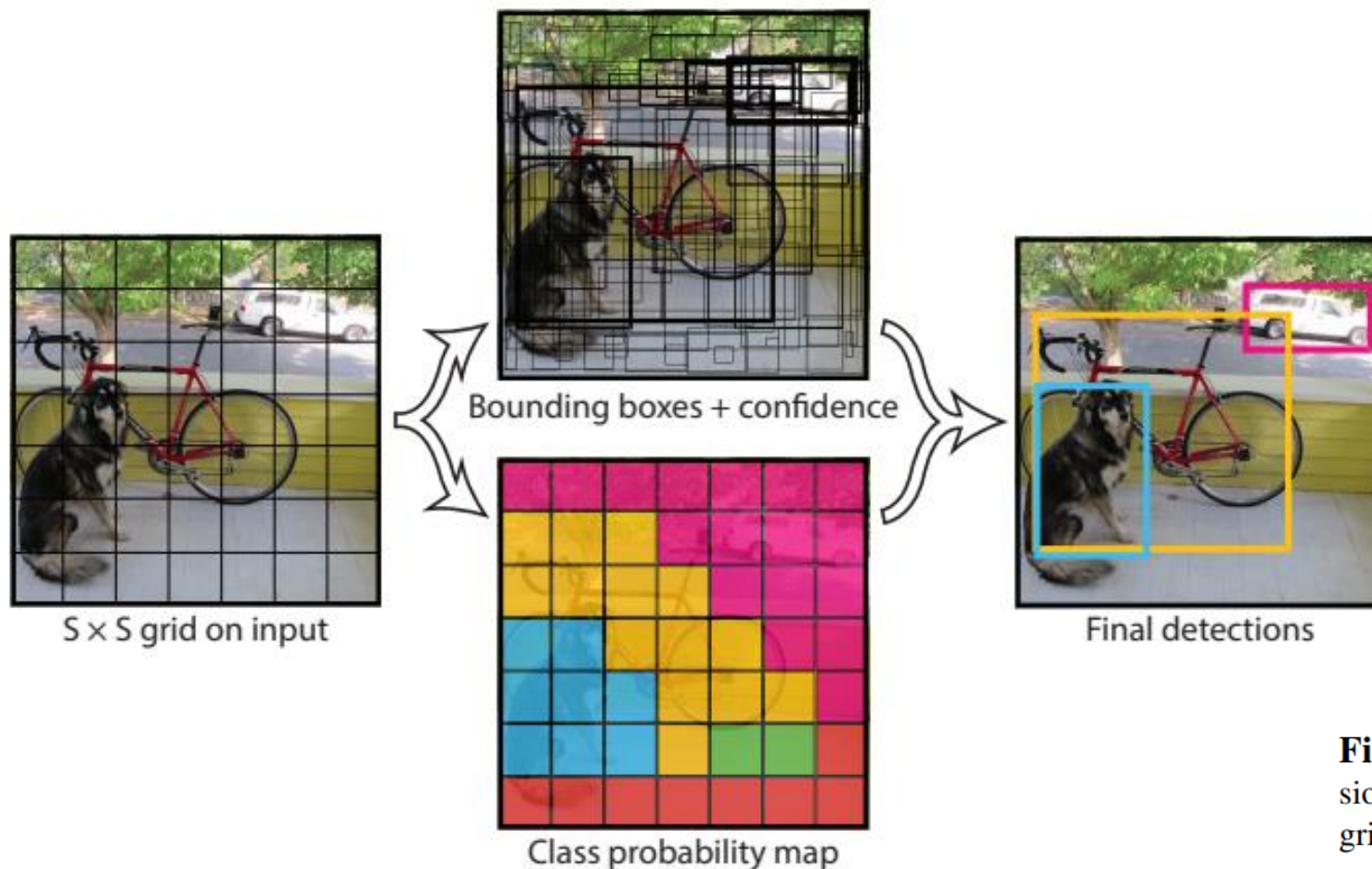
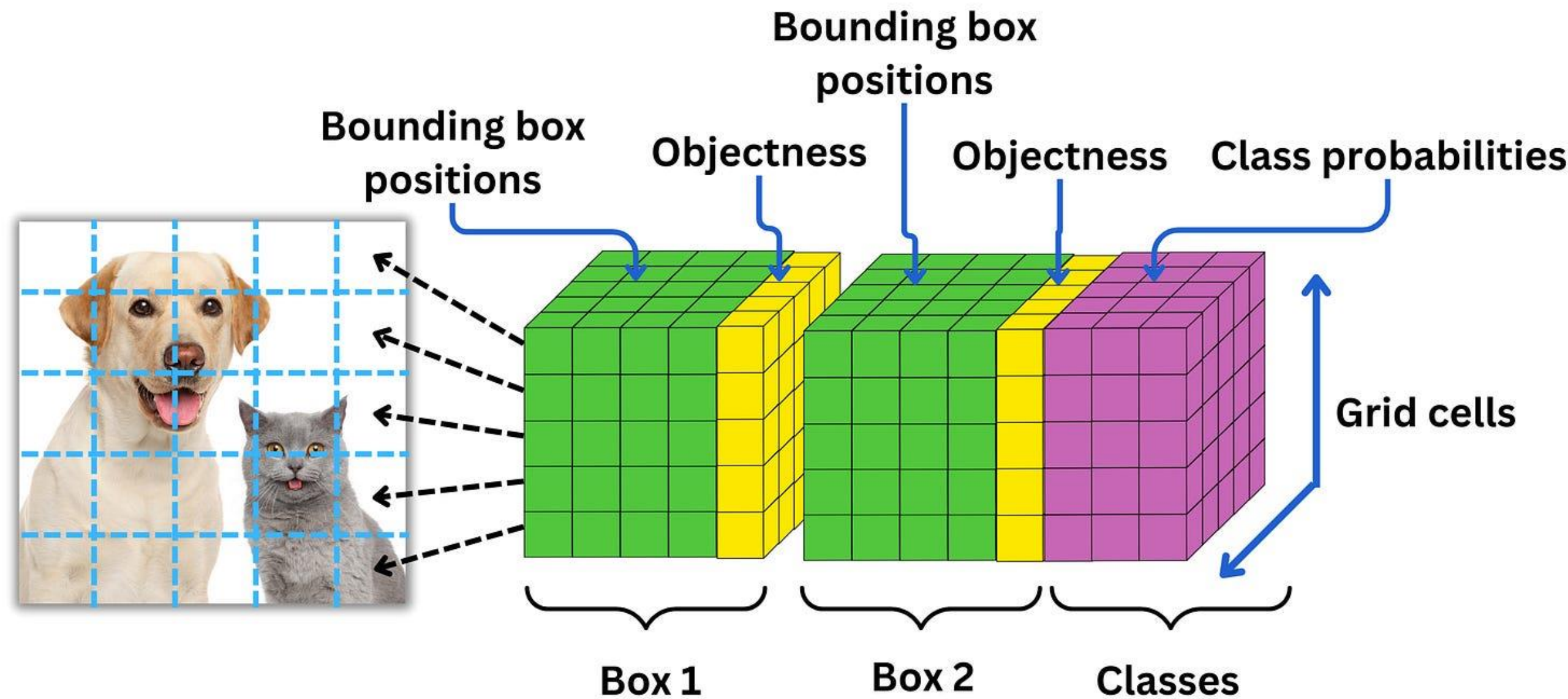
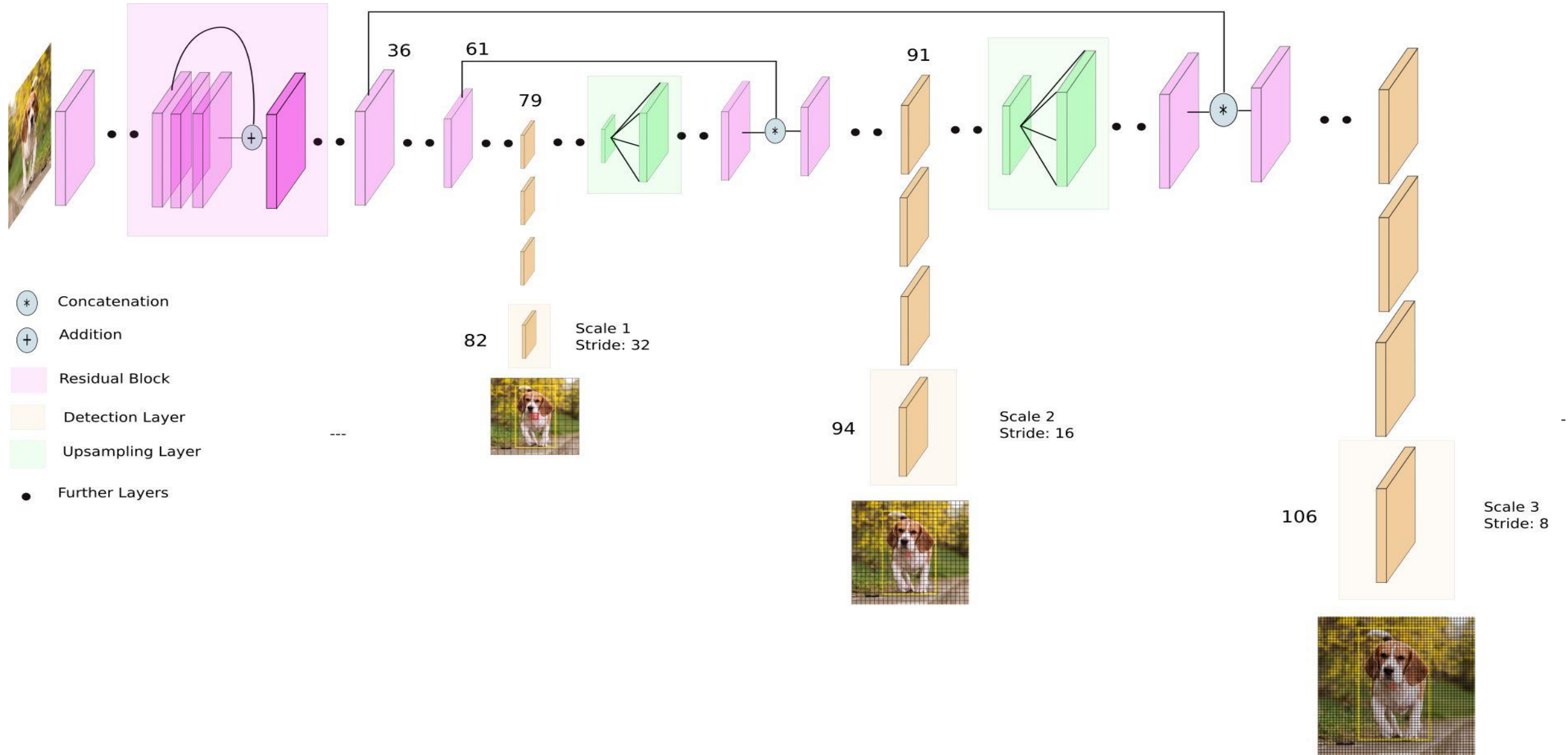


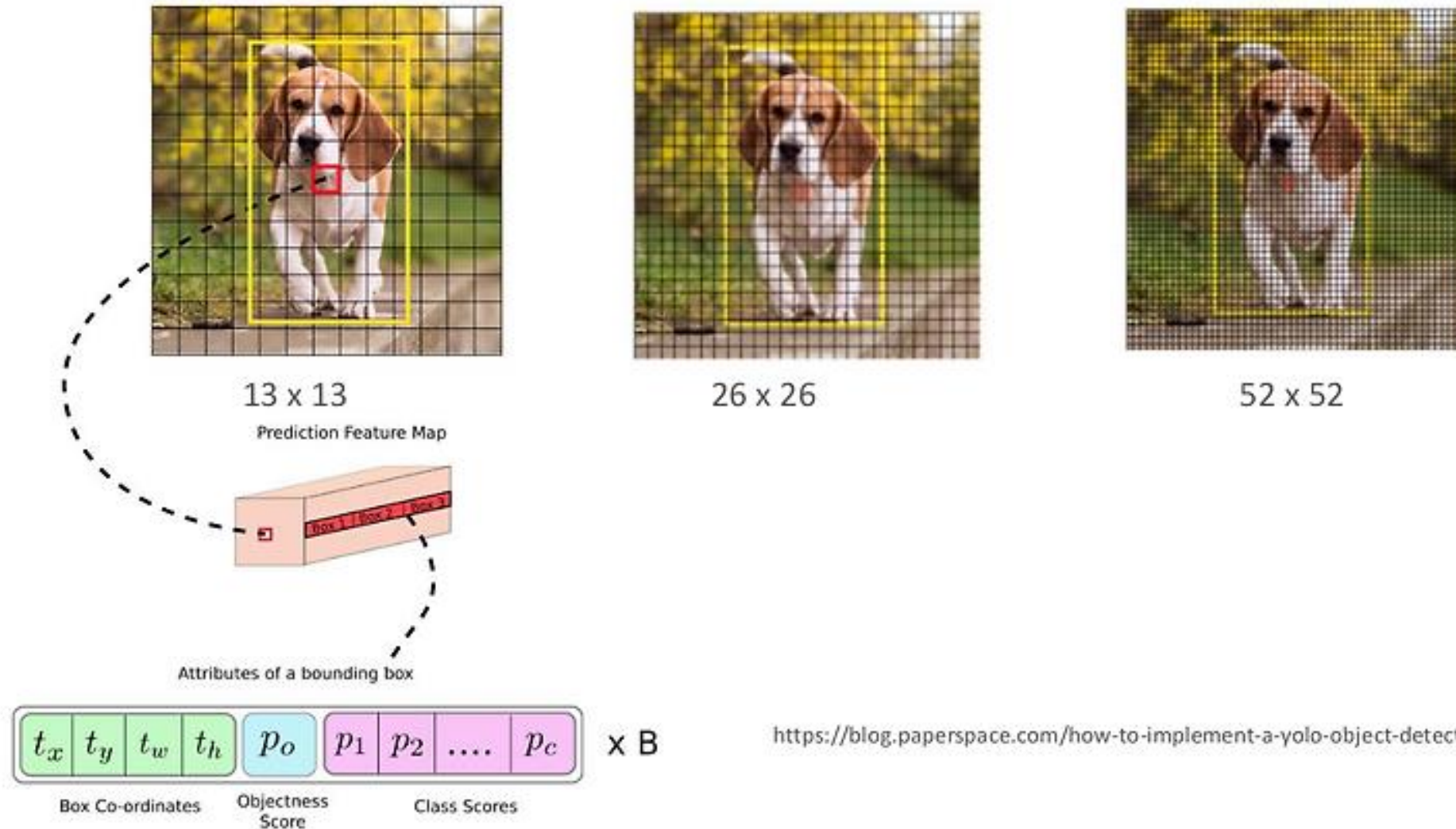
Figure 2: The Model. Our system models detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.



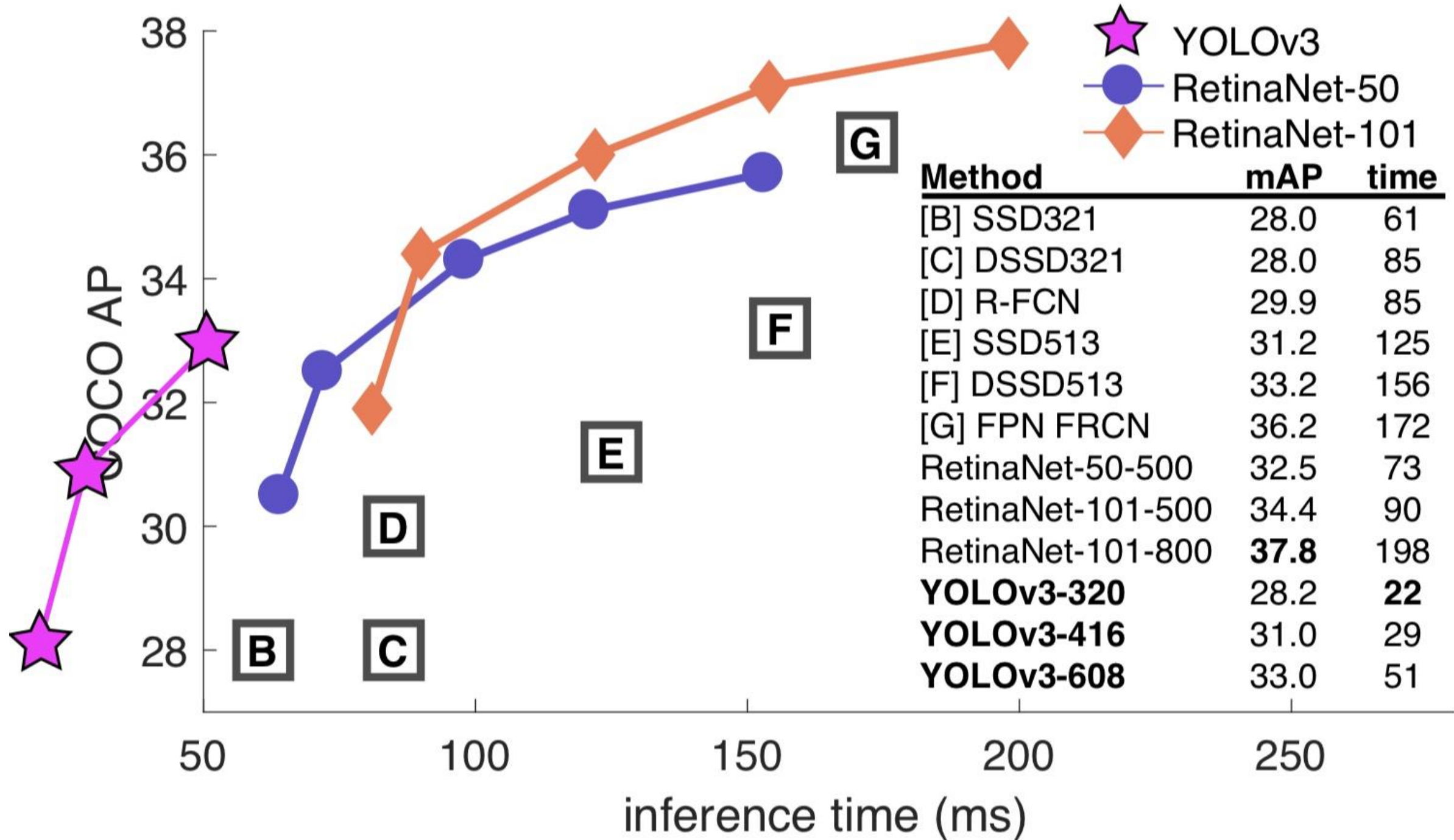
Yolo-3 architecture (2018)



Yolo-3 output layers



<https://blog.paperspace.com/how-to-implement-a-yolo-object-detector-in-pytorch/>



	저자	발표 시점	단체	출간처	특징	정확도 (mAP)	속도 (fps)	논문/github	코드 베이스	라이선스
Yolov1 ~ YOLOv3	Joseph Redmon, Ali Farhadi	2018 년 4 월	University of Washington	arXiv	Feature Pyramid Network(FPN), Multi-Scale Prediction	33.0% (COCO)	35 fps (Titan X)	paper github	Darknet	"YOLO" license
YOLOv4	Alexey Bochkovskiy, et al.	2020 년 4 월	(2,3저자) National Taiwan University	arXiv (YOLOv4)/ CVPR 2021 (Scaled-YOLOv4)	Bag of Freebies (BoF), Bag of Specials (BoS)	43.5% (COCO)	65 fps (V100)	paper paper(scaled) github (darknet) github(pytorch)	Darknet	"YOLO" license
YOLOv5	Glenn Jocher	2020 년 6 월	Ultralytics	-	PyTorch Implementation, CSP backbone	50% (COCO)	93 fps (RTX4090)	github	PyTorch	AGPL 3.0
YOLOv6	Chuyi Li, et al.	2022 년 6 월	Meituan	arXiv	Bi-directional Concatenation(BiC) module, Anchor-aided training(AAT) strategy	52.8% (COCO)	116 fps (T4)	paper(v1.0) paper(v3.0) github	PyTorch	GPL 3.0
YOLOv7	Chien-Yao Wang*, et al. (*v4 논문의 2 저자)	2022 년 7 월	Academia Sinica	CVPR 2023	Extended-ELAN, Model Reparameterization	55.9% (COCO)	56 fps (V100)	paper github	YOLOv5 (PyTorch)	GPL 3.0
YOLOv8	Glenn Jocher, et al.	2023 년 1 월	Ultralytics	-	Anchor-free, Multi-task application	53.9% (COCO)	78 fps (T4)	website	PyTorch	AGPL 3.0
YOLOv9	Chien-Yao Wang, et al.	2024 년 2 월	Academia Sinica	arXiv	Programmable Gradient Information (PGI), Generalized ELAN	55.6% (COCO)	73 fps (T4)	paper github	YOLOv5	(PyTorch)
YOLOv10	Ao Wang, et al.	2024 년 5 월	Tsinghua University	arXiv / NeurIPS 2024 poster	NMS-free training, spatial-channel separate downsampling	54.4% (COCO)	94 fps(T4)	paper github	Ultralytics(YOLOv8/PyTorch)	AGPL 3.0

Yolo version comparison

