

Computer Vision

Object detection

Tae-Hyun Oh (오태현)

전자전기공학과

POSTECH

Slide by Dongmin Choi (최동민)

TAs: { Dongmin Choi , Jongha Kim, Juyong Lee, Sungbin Kim } (in alphabetic order)

1. Object detection

- 1.1 What is object detection?
- 1.2 What are the applications of object detection?

2. Two-stage detector (R-CNN family)

- 2.1 R-CNN
- 2.2 Fast R-CNN
- 2.3 Faster R-CNN

3. Single-stage detector

- 3.1 YOLO
- 3.2 SSD

4. Single-stage detector vs. two-stage detector

- 4.1 Focal loss
- 4.2 RetinaNet

5. Detection with Transformer

1.

Object detection

1.1 What is object detection?

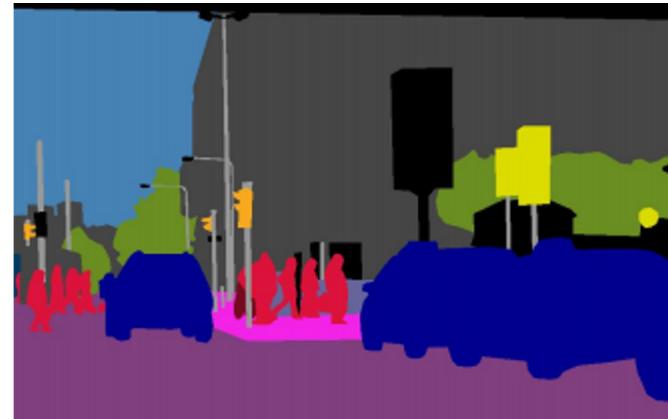
Object detection

Fundamental image recognition tasks

[Kirillov et al., CVPR 2019]



Image



Semantic segmentation



Instance segmentation

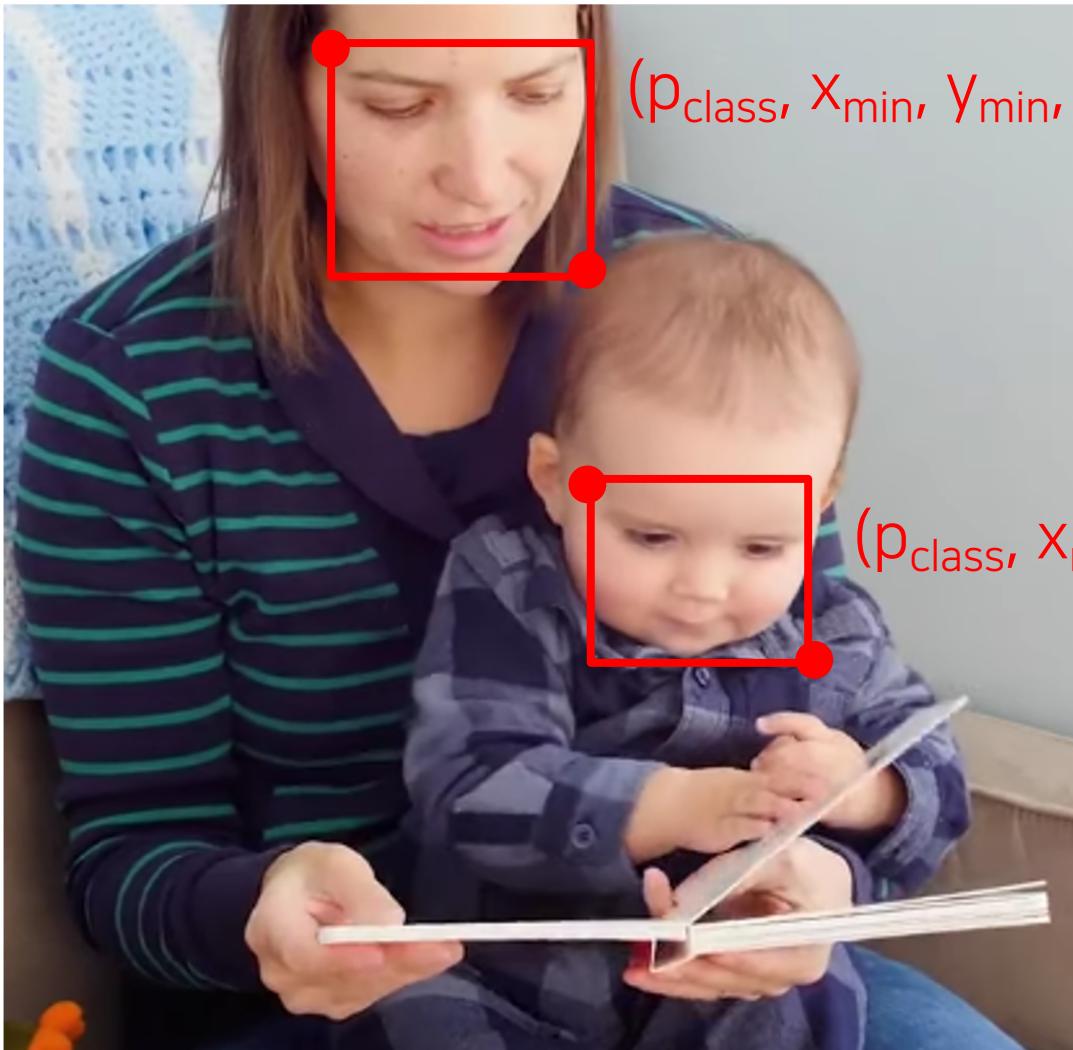


Panoptic segmentation

X

1.1 What is object detection?

Object detection

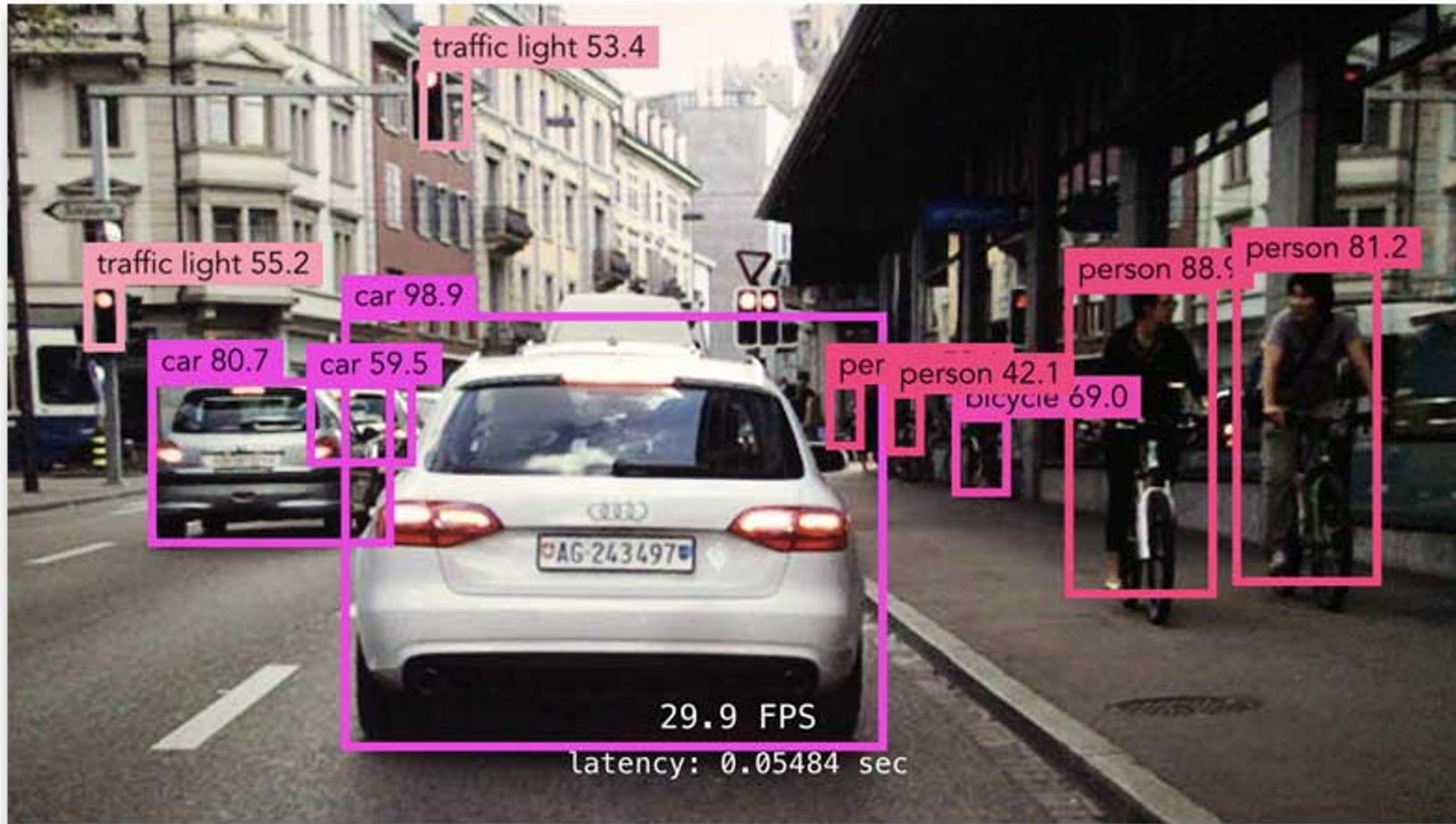


Classification + Box localization

1.2 What are the applications of object detection?

Object detection

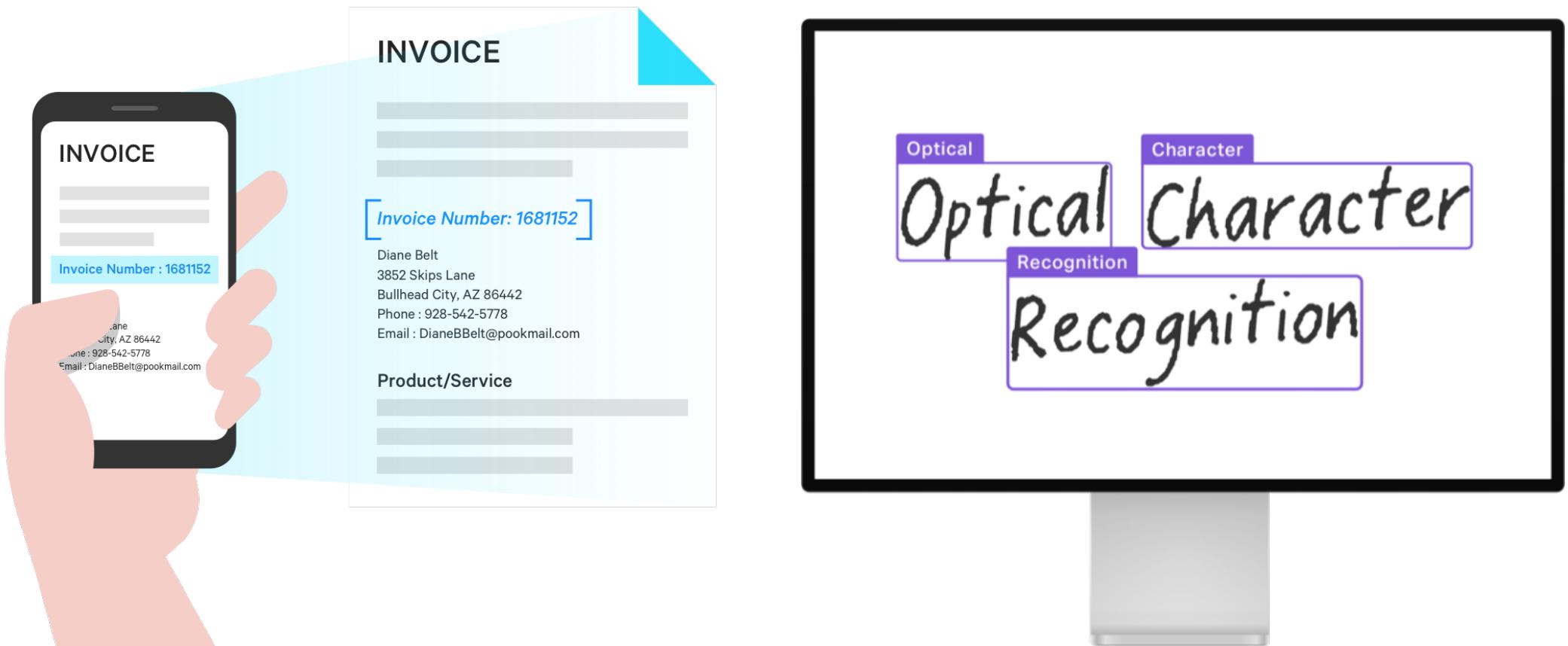
Autonomous driving



1.2 What are the applications of object detection?

Object detection

Optical Character Recognition (OCR)



2.

Two-stage detector

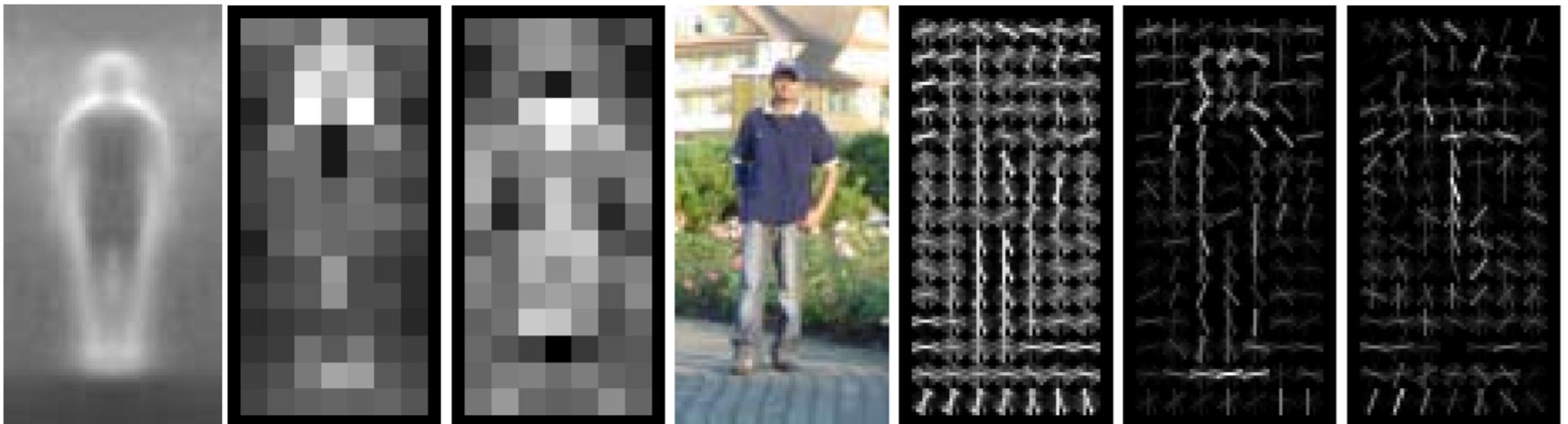
2.0 Traditional methods – hand-crafted techniques

Two-stage detector

Gradient-based detector (e.g., HOG)

[Dalal et al., CVPR 2005]

* HOG = Histogram of Oriented Gradients, SVM = Support Vector Machine



(a)

(b)

(c)

(d)

(e)

(f)

(g)

Average
Gradient

max (+)
SVM weight

max (-)
SVM weight

Image

R-HOG
descriptor

R-HOG
w/ (+) SVM

R-HOG
w/ (-) SVM

2.0 Traditional methods – hand-crafted techniques

Two-stage detector

Selective search

[Uijlings et al., IJCV 2013]

1. Over-segmentation
2. Iteratively merging similar regions
3. Extracting candidate boxes from all remaining segmentations



2.1 R-CNN

Two-stage detector

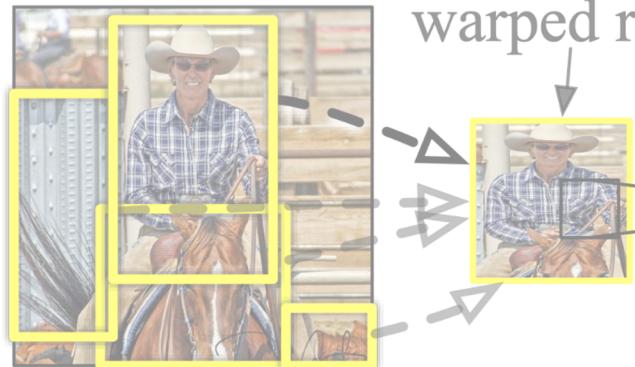
[Girshick et al., CVPR 2014]

Directly leverage image classification networks for object detection

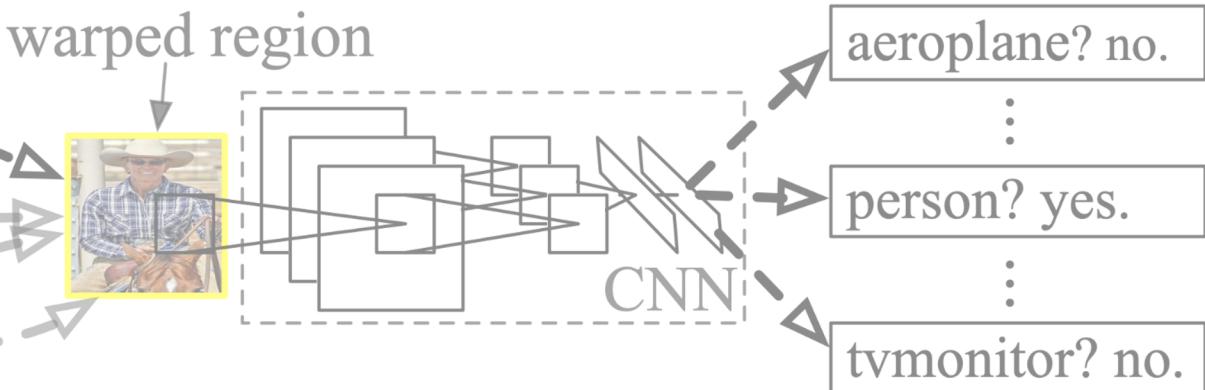
R-CNN: *Regions with CNN features*



1. Input
image



2. Extract region
proposals (~2k)



3. Compute
CNN features

4. Classify
regions

2.1 R-CNN

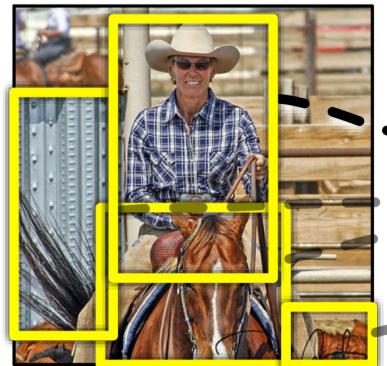
Two-stage detector

[Girshick et al., CVPR 2014]

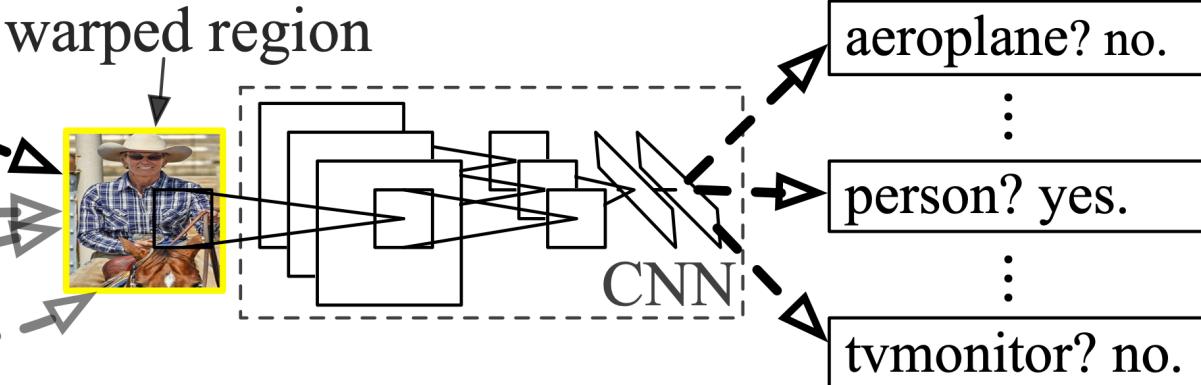
R-CNN: *Regions with CNN features*



1. Input
image



2. Extract region
proposals (~2k)



3. Compute
CNN features

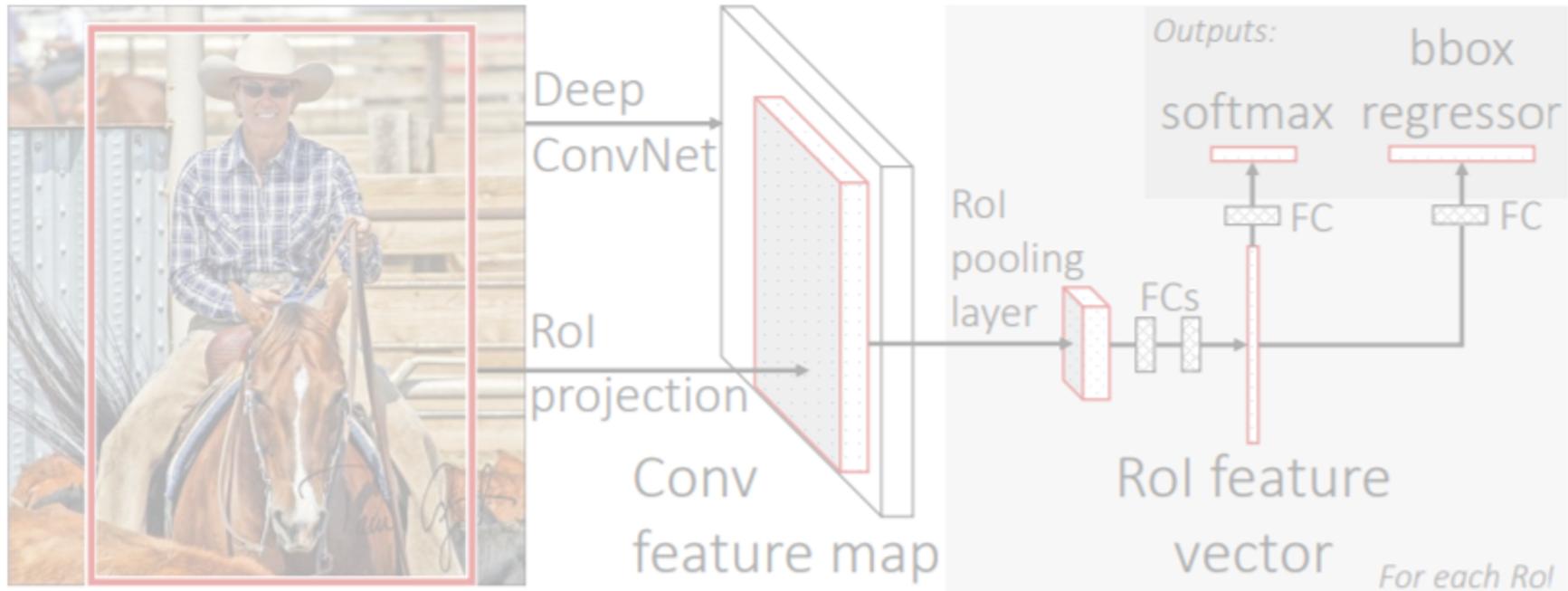
4. Classify
regions

2.2 Fast R-CNN

Two-stage detector

[Girshick et al., ICCV 2015]

Recycle a pre-computed feature for multiple object detection

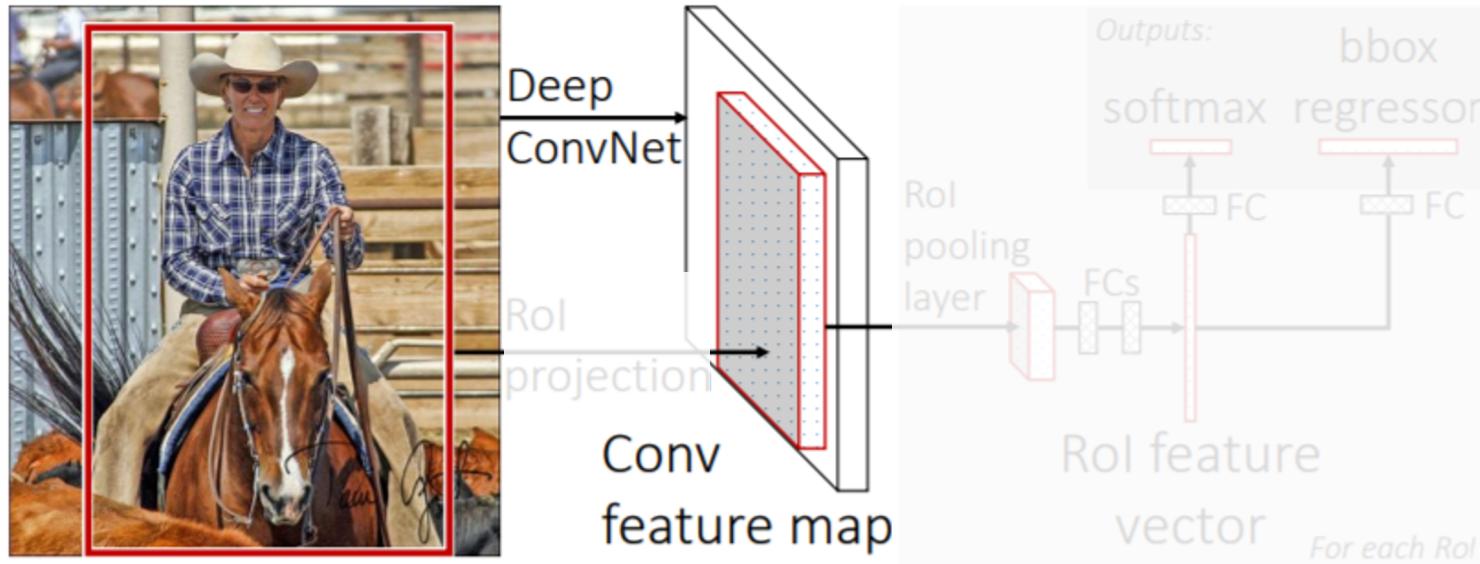


2.2 Fast R-CNN

Two-stage detector

Fast R-CNN

[Girshick et al., ICCV 2015]



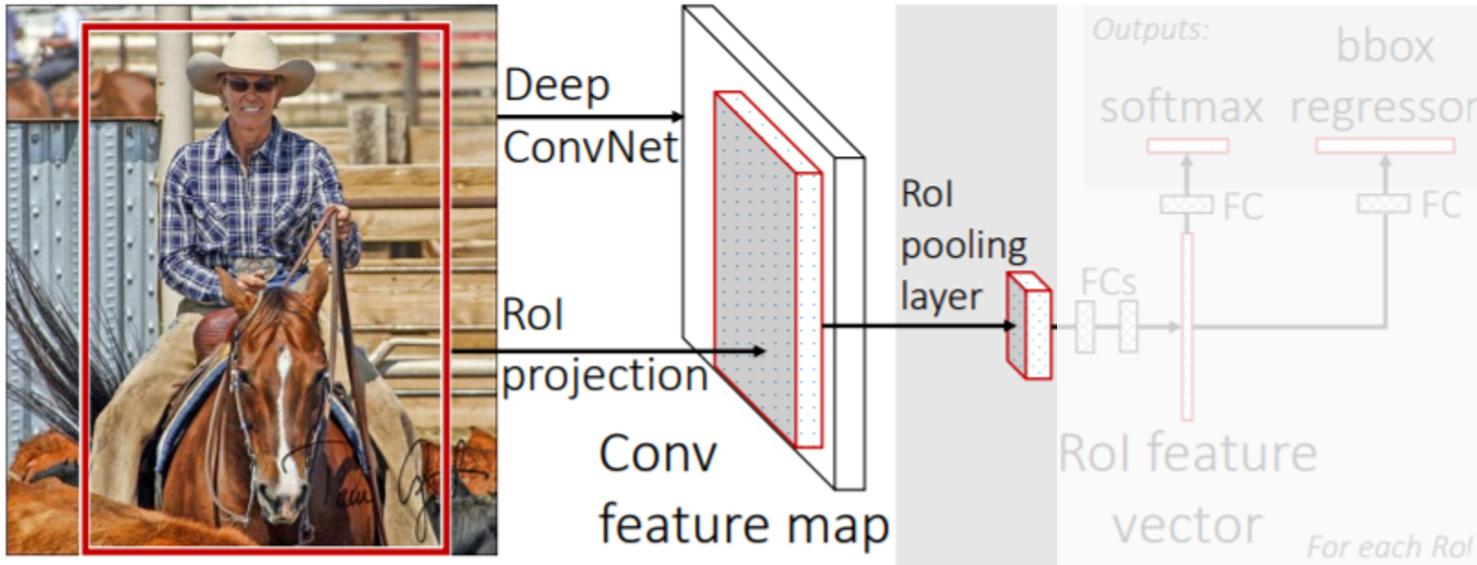
1. Conv. feature map from the original image

2.2 Fast R-CNN

Two-stage detector

[Girshick et al., ICCV 2015]

* ROI = Region of Interest



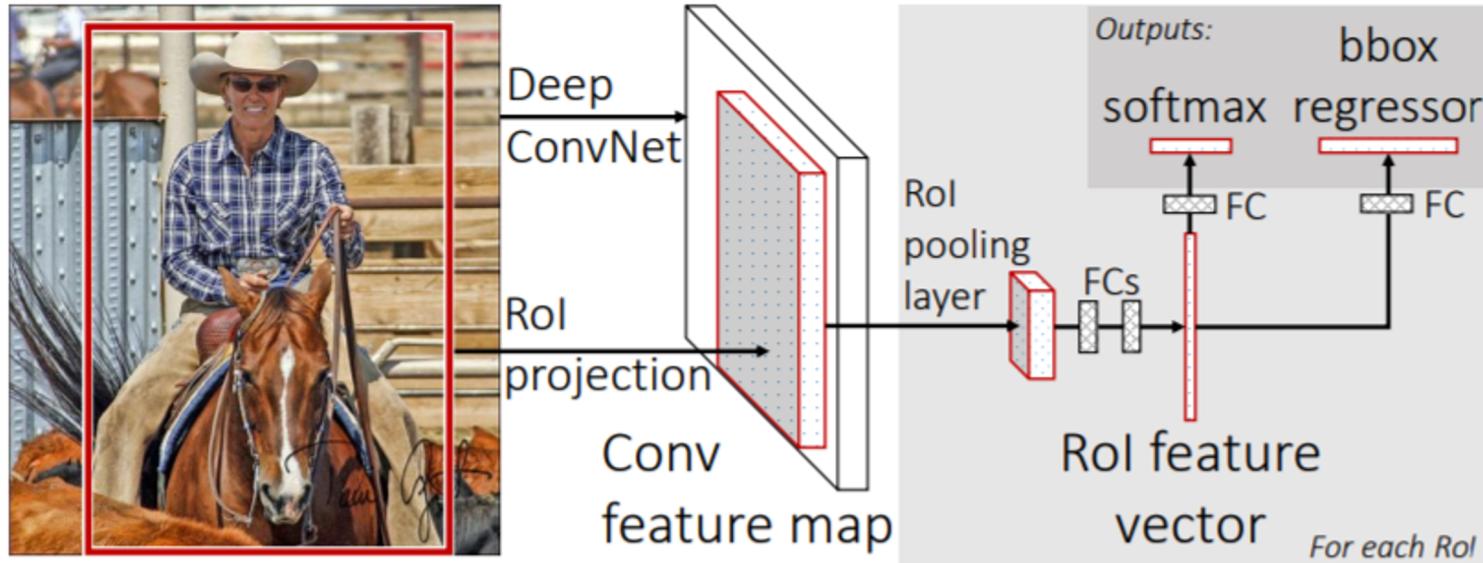
1. Conv. feature map from the original image
2. ROI feature extraction from the feature map through ROI pooling

2.2 Fast R-CNN

Two-stage detector

[Girshick et al., ICCV 2015]

* ROI = Region of Interest



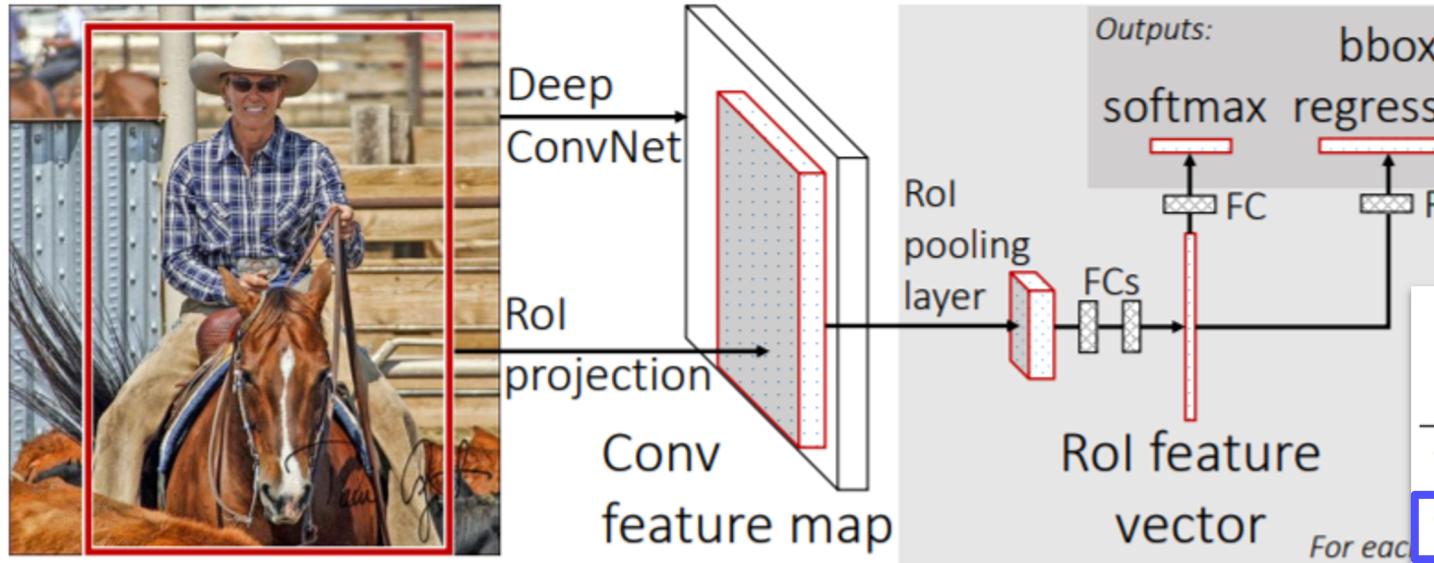
1. Conv. feature map from the original image
2. ROI feature extraction from the feature map through ROI pooling
3. Class and box prediction for each ROI

2.2 Fast R-CNN

Two-stage detector

[Girshick et al., ICCV 2015]

* ROI = Region of Interest



1. Conv. feature map from the original image
2. ROI feature extraction from the feature map through ROI pooling
3. Class and box prediction for each ROI

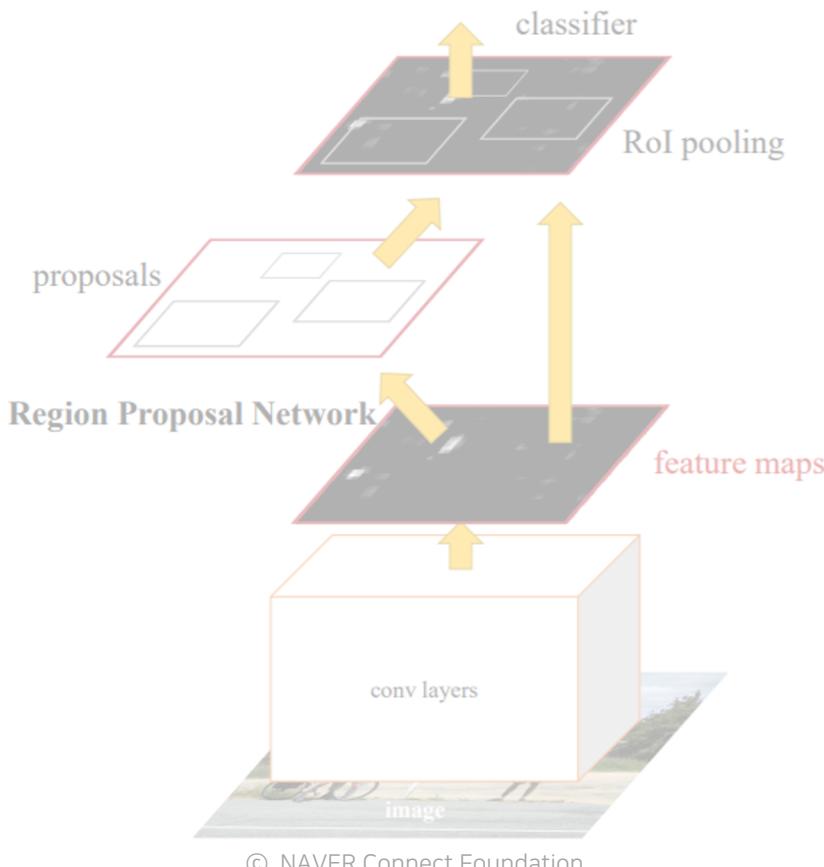
	Fast R-CNN			R-CNN		
	S	M	L	S	M	L
train time (h)	1.2	2.0	9.5	22	28	84
train speedup	18.3×	14.0×	8.8×	1×	1×	1×
test rate (s/im)	0.10	0.15	0.32	9.8	12.1	47.0
▷ with SVD	0.06	0.08	0.22	-	-	-
test speedup	98×	80×	146×	1×	1×	1×
▷ with SVD	169×	150×	213×	-	-	-
VOC07 mAP	57.1	59.2	66.9	58.5	60.2	66.0
▷ with SVD	56.5	58.7	66.6	-	-	-

2.3 Faster R-CNN

Two-stage detector

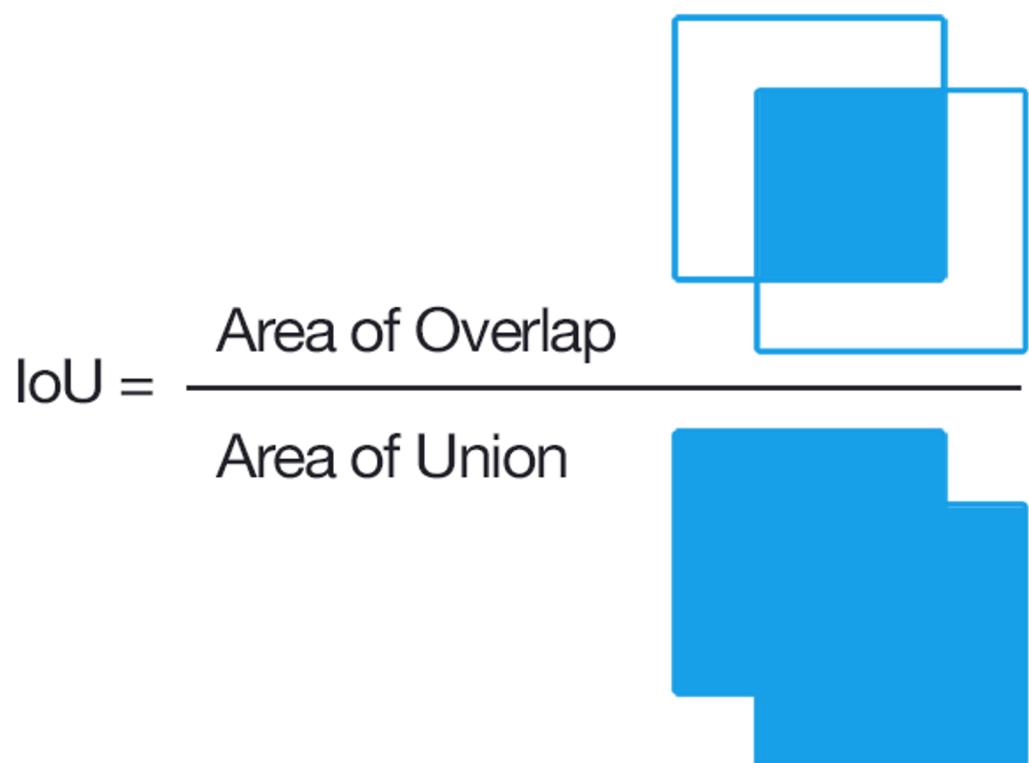
[Ren et al., NeurIPS, 2015]

End-to-end object detection by neural region proposal

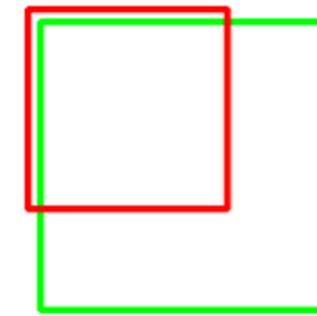


2.3 Faster R-CNN

Two-stage detector

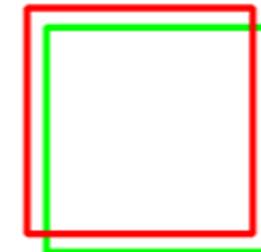


IoU: 0.4034



Poor

IoU: 0.7330



Good

IoU: 0.9264



Excellent

IoU (Intersection over Union)

: A metric commonly used in object detection

2.3 Faster R-CNN

Two-stage detector

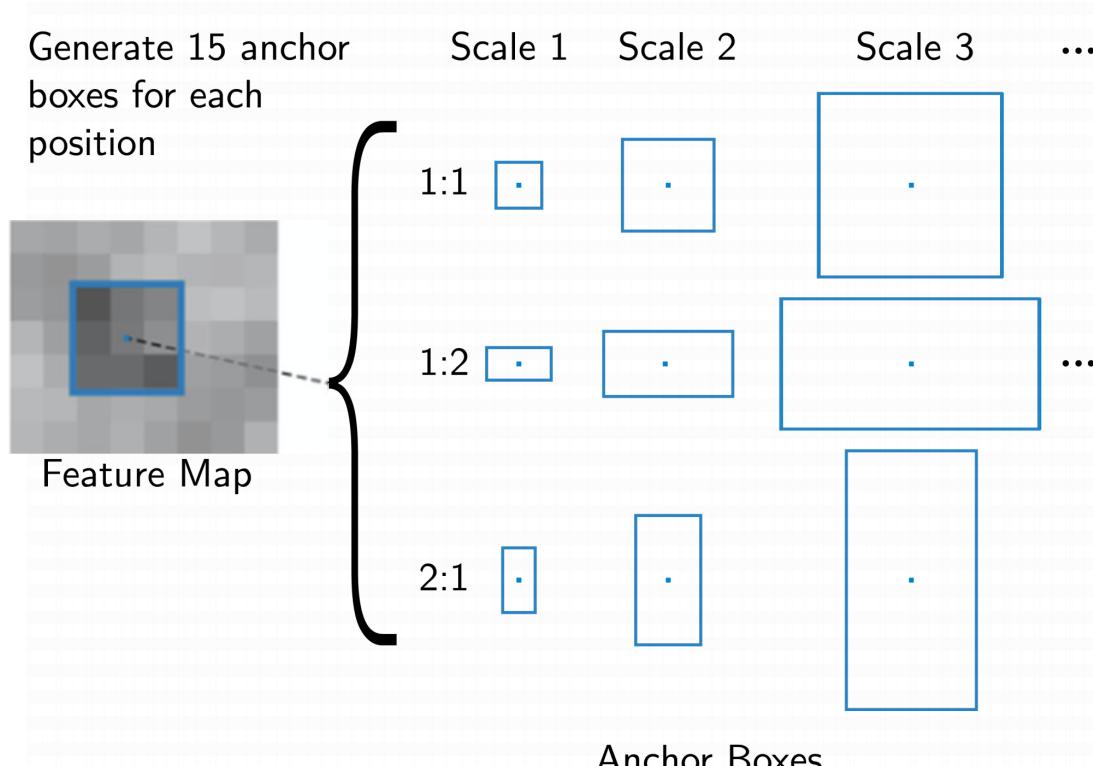


Illustration from [Ferguson et al. 2018]

Anchor boxes

- A set of pre-defined bounding boxes
- IoU with GT $> 0.7 \Rightarrow$ positive sample
- IoU with GT $< 0.3 \Rightarrow$ negative sample

2.3 Faster R-CNN

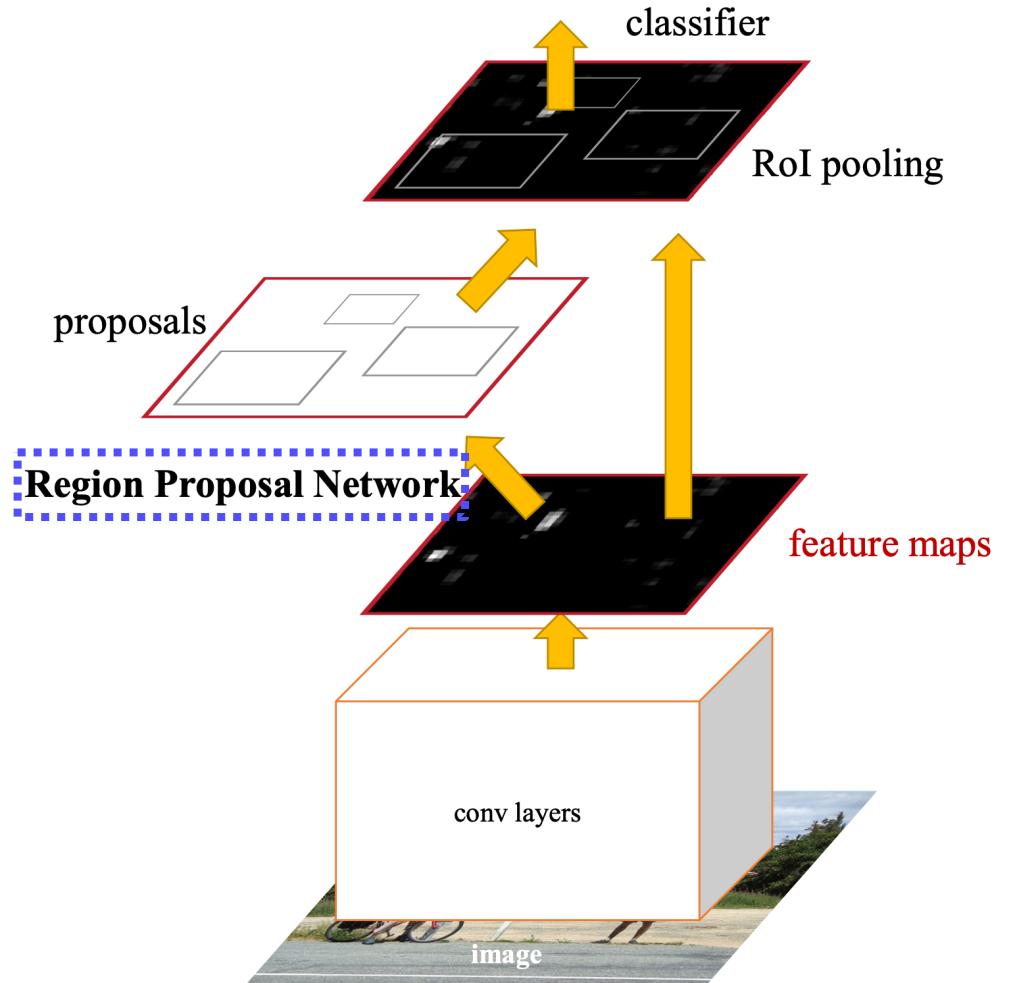
Two-stage detector

[Ren et al., NeurIPS, 2015]

Time-consuming selective search



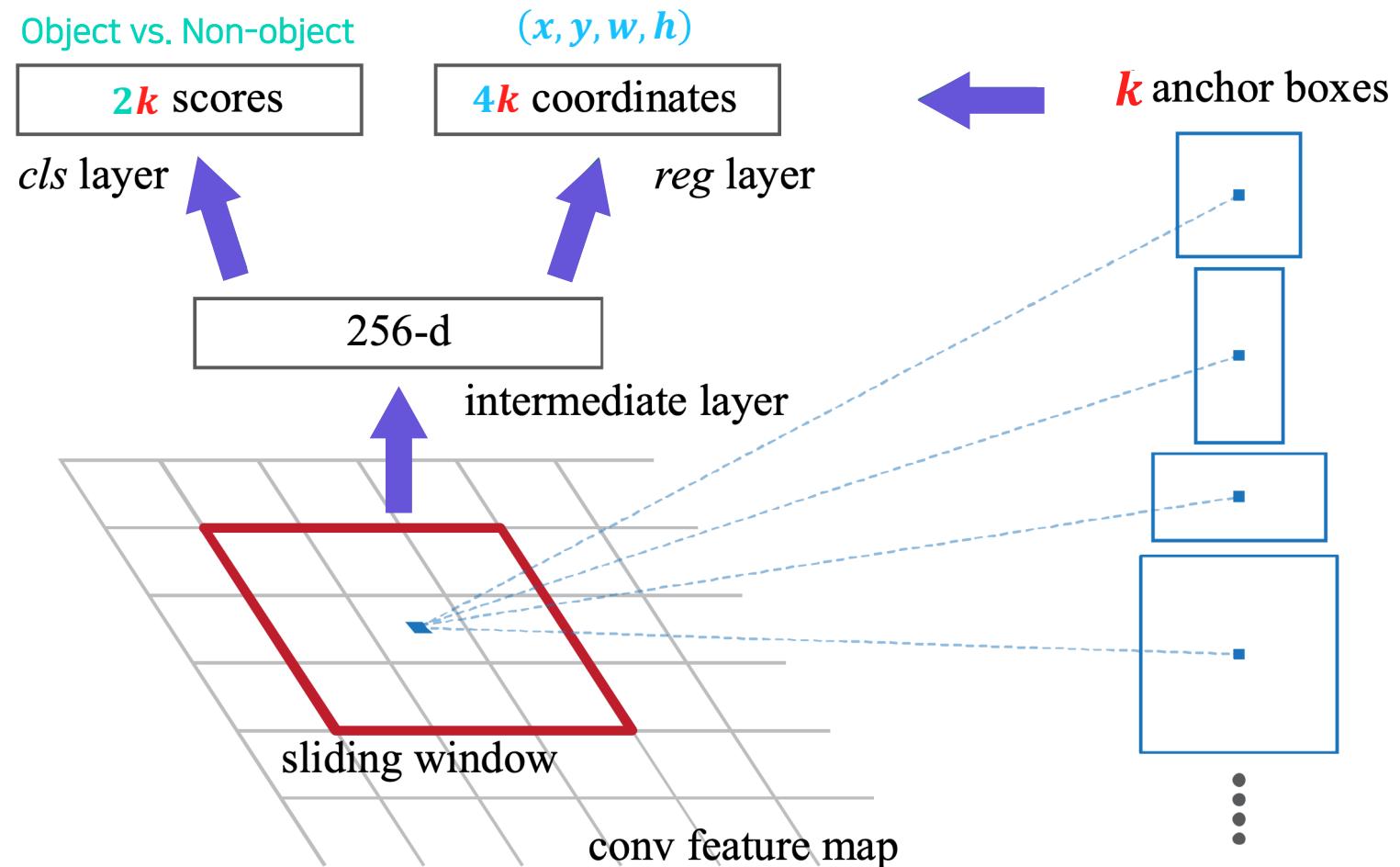
Region Proposal Network (RPN)



2.3 Faster R-CNN

Two-stage detector

[Ren et al., NeurIPS, 2015]



2.3 Faster R-CNN

Two-stage detector

Non-Maximum Suppression (NMS)

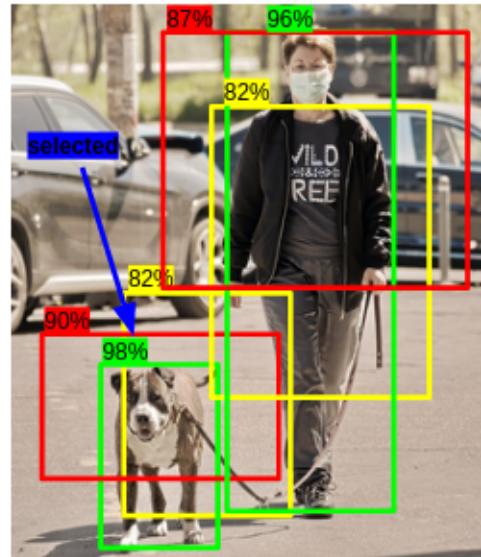
Step 1: Select the box with the highest objectiveness score

Step 2: Compare IoU of this box with other boxes

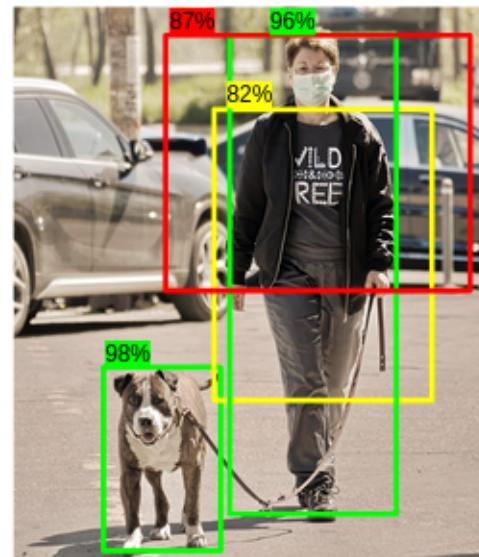
Step 3: Remove the bounding boxes with $\text{IoU} \geq 50\%$

Step 4: Move to the next highest objectiveness score

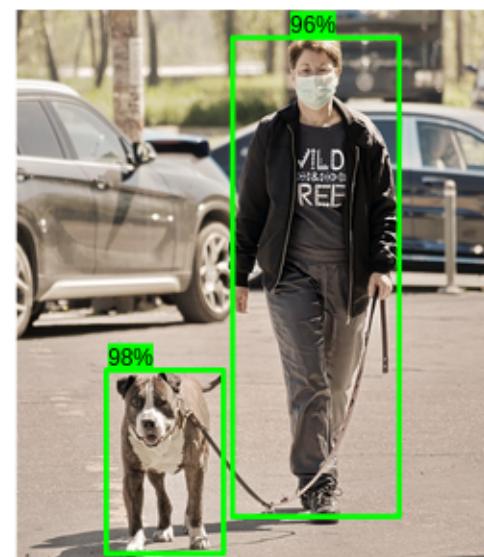
Step 5: Repeat steps 2-4



Step 1: Selecting Bounding box with highest score



Step 3: Delete Bounding box with high overlap

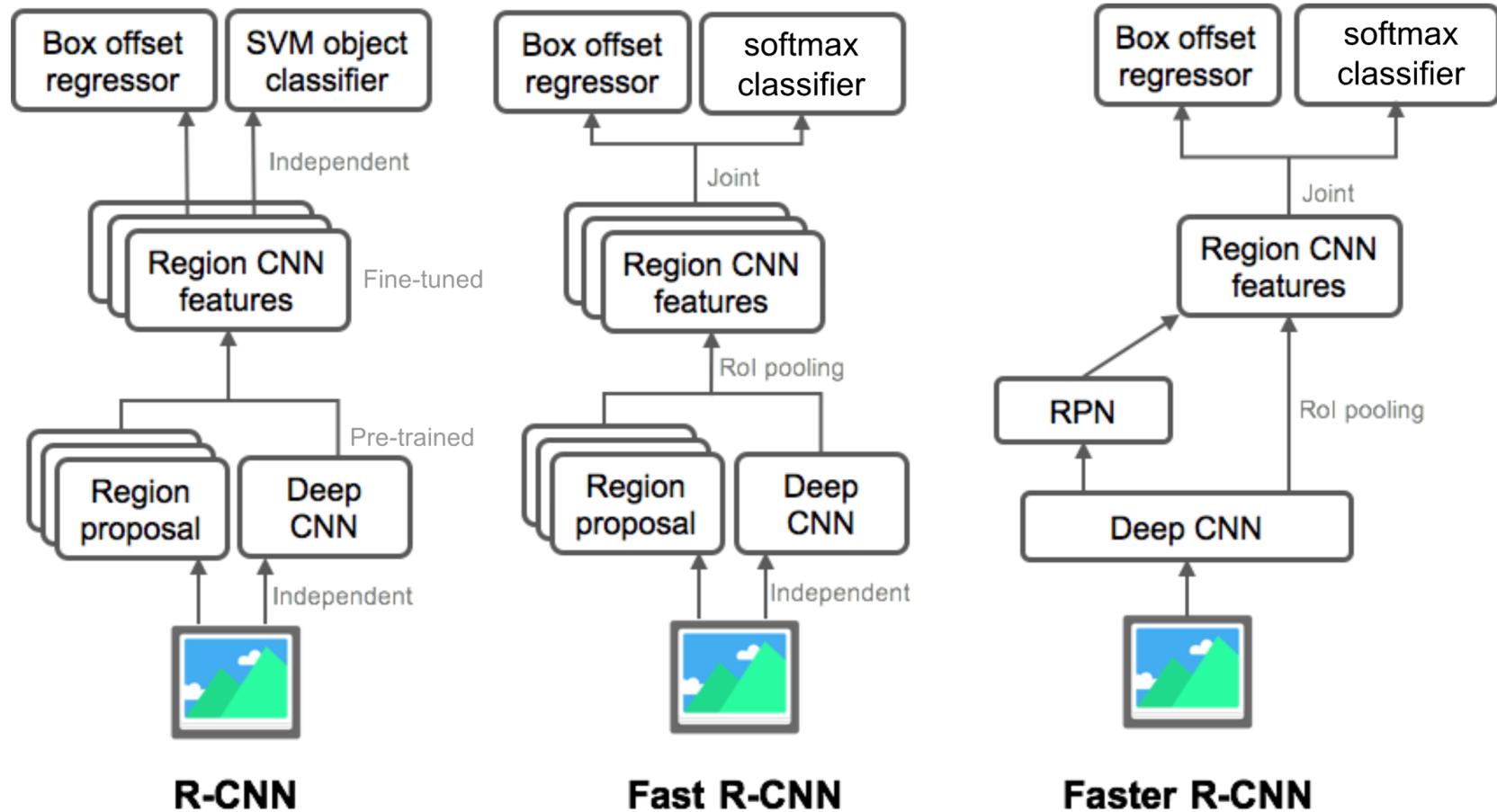


Step 5: Final Output

2.3 Faster R-CNN

Two-stage detector

Summary of the R-CNN family



3.

Single-stage detector

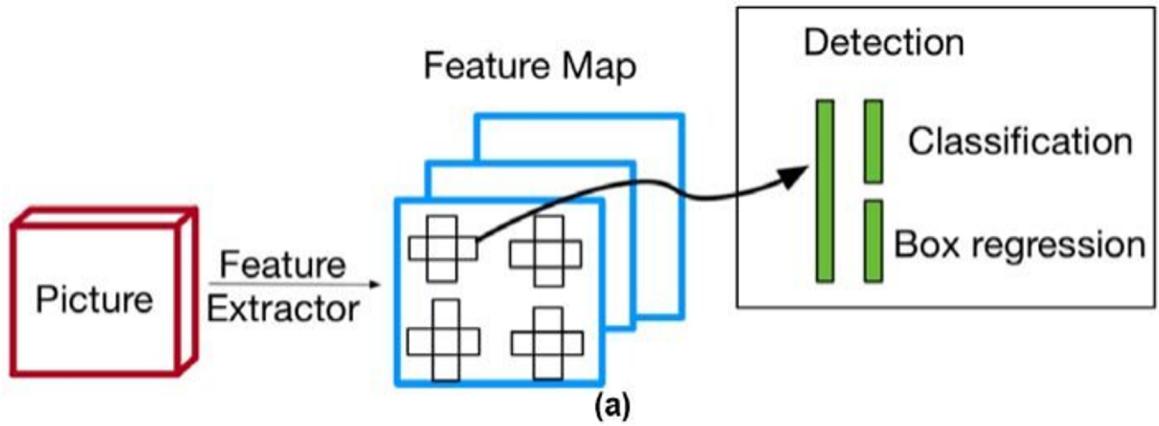
3.0 Comparison with two-stage detectors

One-stage detector

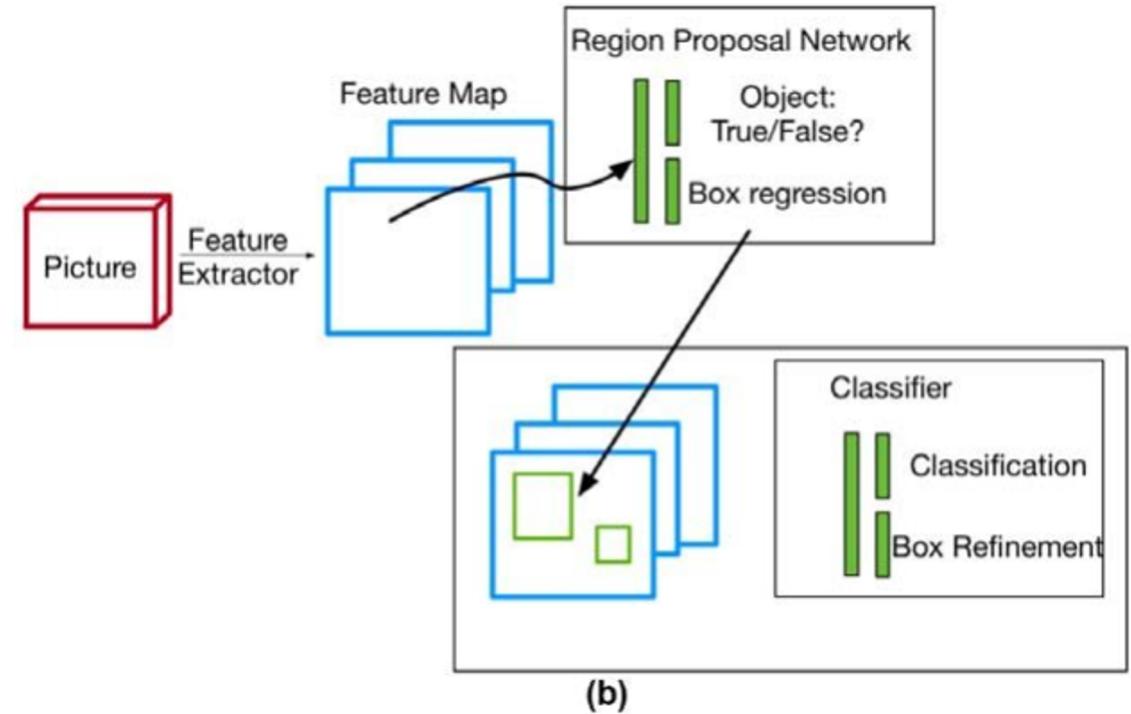
One-stage vs. two-stage

Illustrations from [Ndonhong et al., Offshore Technology Conference 2019]

- No explicit RoI pooling



One-stage detector

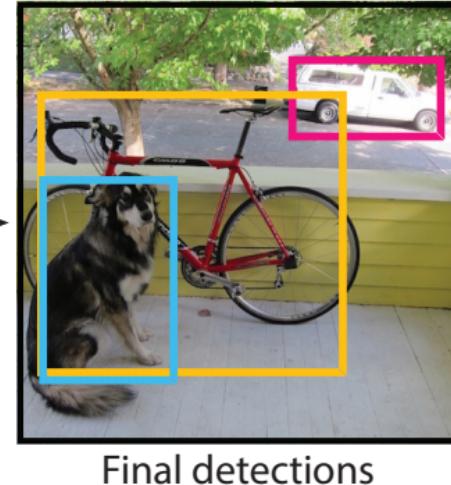
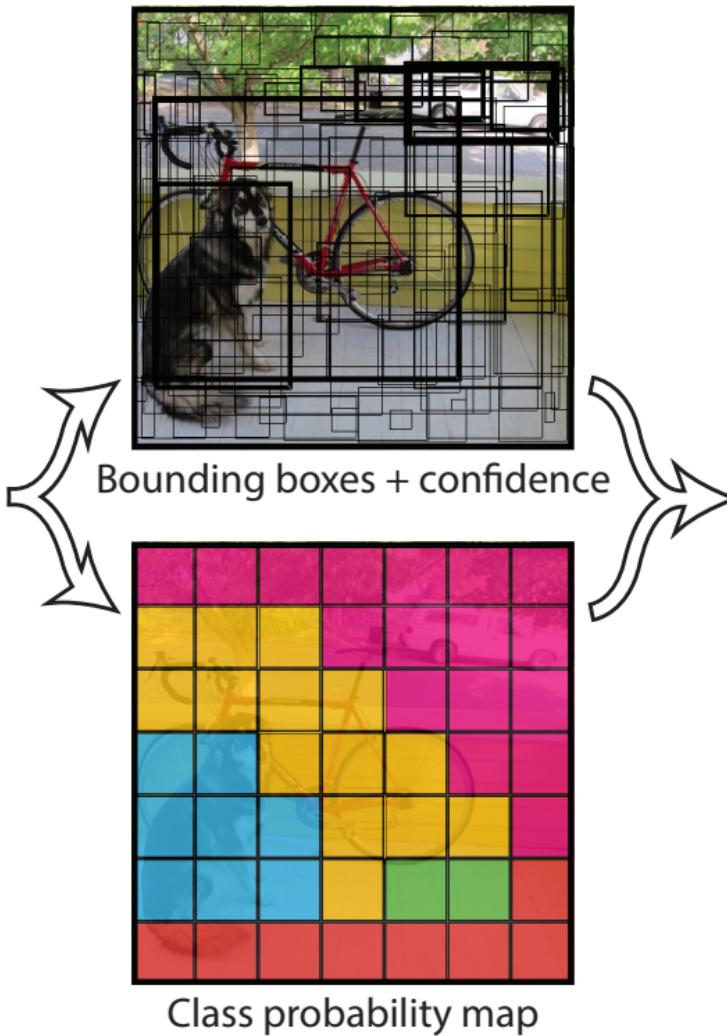
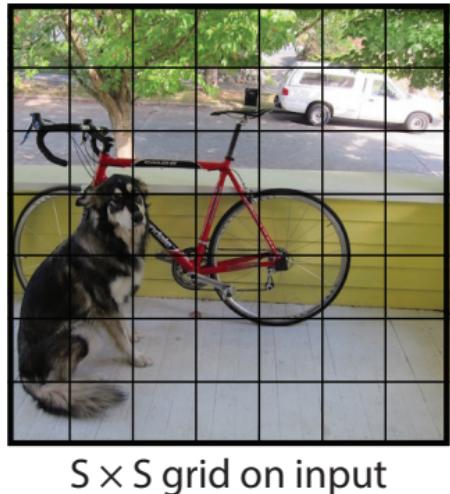


Two-stage detector

3.1 You only look once (YOLO)

One-stage detector

[Redmon et al., CVPR 2016]

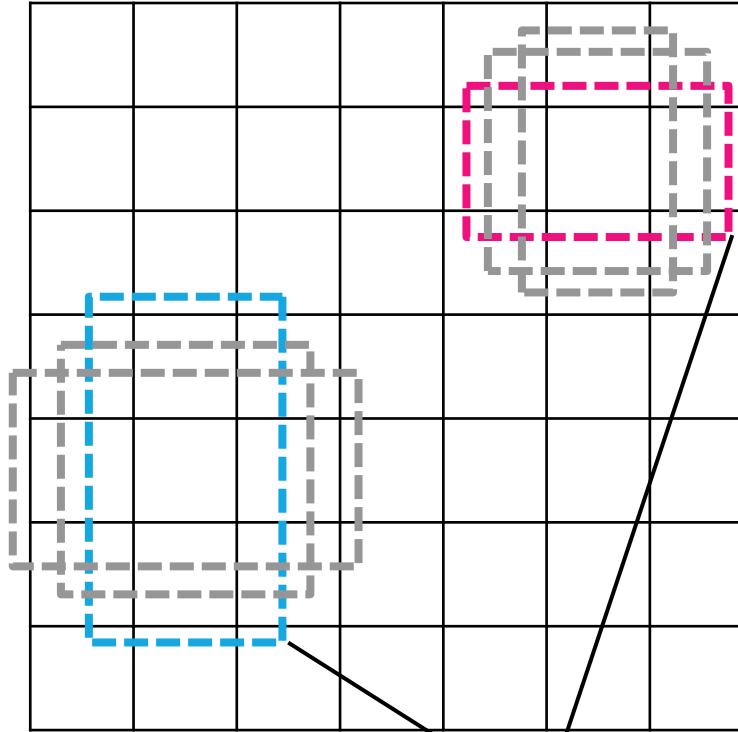
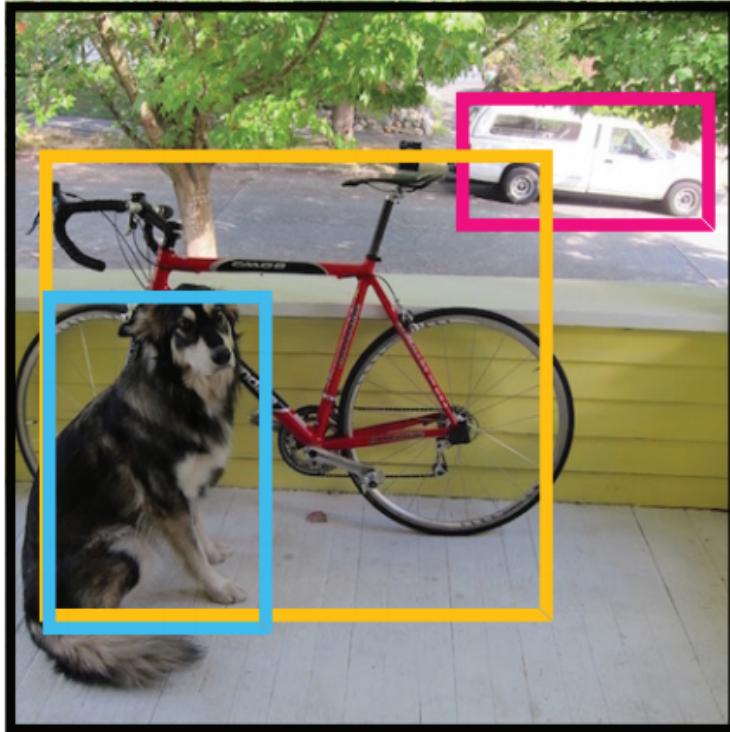


3.1 YOLO

One-stage detector

[Redmon et al., CVPR 2016]

* Anchor boxes for yellow bicycle are omitted

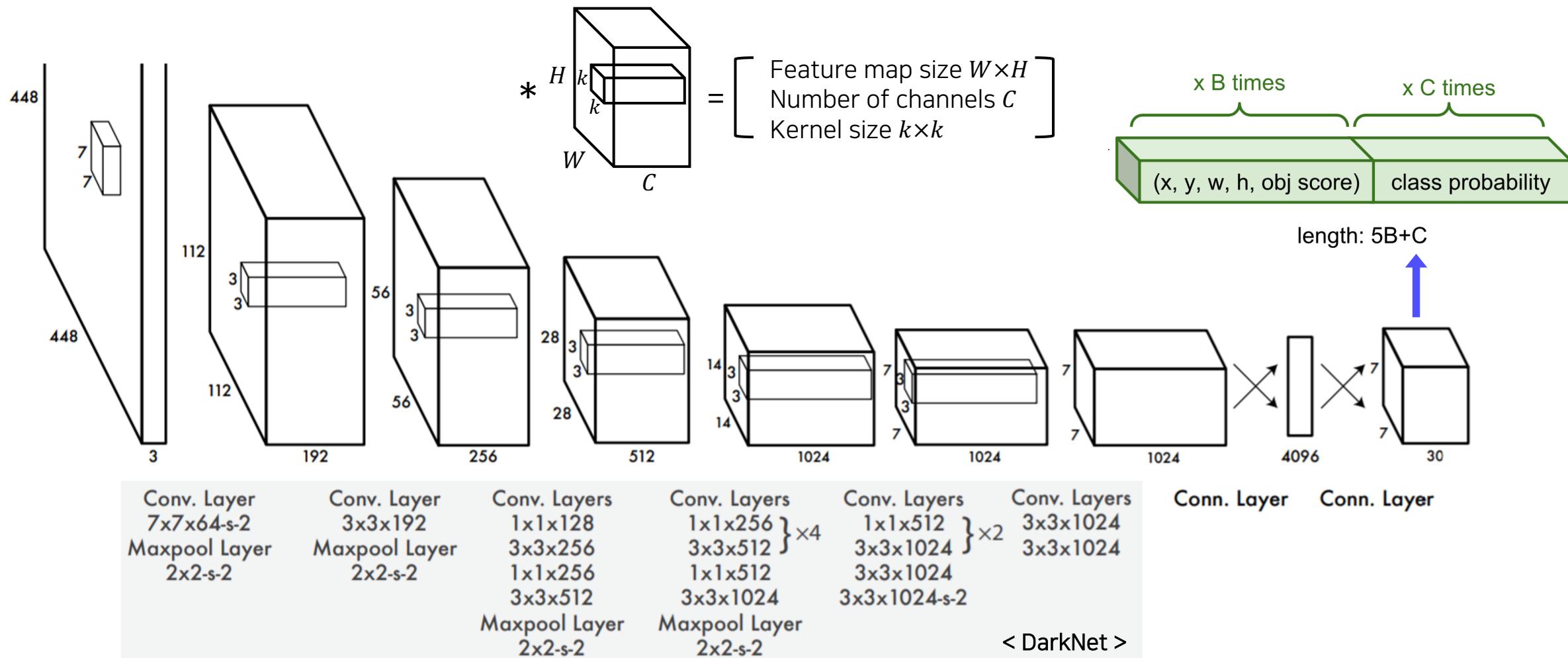


Colored boxes are
matched positive samples

3.1 YOLO

One-stage detector

[Redmon et al., CVPR 2016]



3.1 YOLO

One-stage detector

[Redmon et al., CVPR 2016]

* mAP : mean Average Precision (metrics)

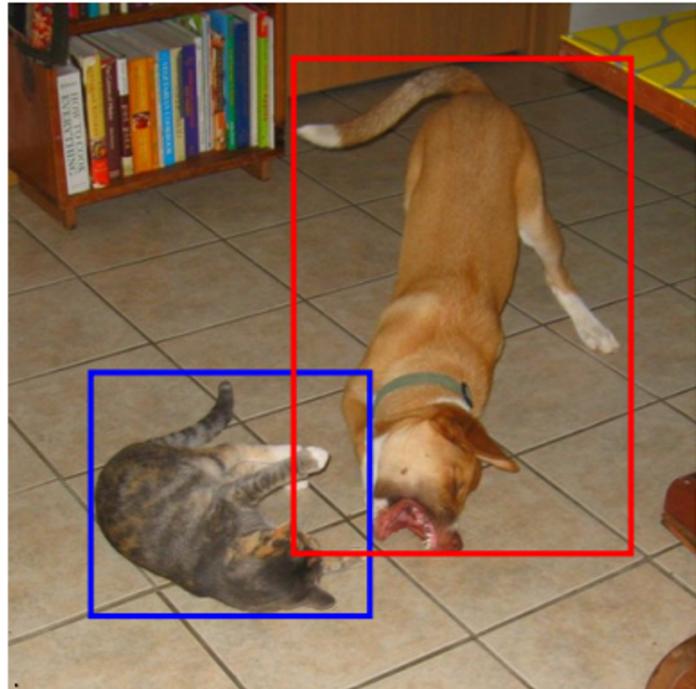
* FPS : Frame Per Second (speed)

Real-Time Detectors	Train	mAP	FPS
100Hz DPM [31]	2007	16.0	100
30Hz DPM [31]	2007	26.1	30
Fast YOLO	2007+2012	52.7	155
YOLO	2007+2012	63.4	45
Less Than Real-Time			
Fastest DPM [38]	2007	30.4	15
R-CNN Minus R [20]	2007	53.5	6
Fast R-CNN [14]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[28]	2007+2012	73.2	7
Faster R-CNN ZF [28]	2007+2012	62.1	18
YOLO VGG-16	2007+2012	66.4	21

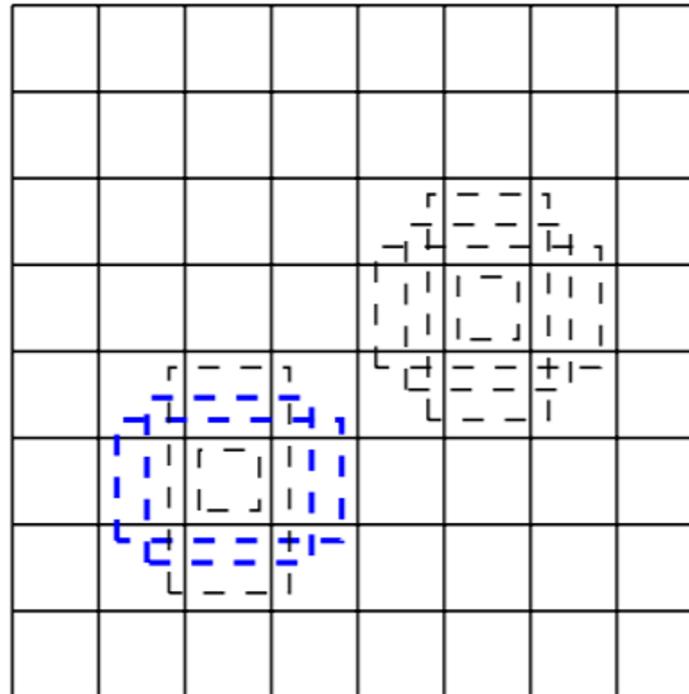
3.2 Single Shot MultiBox Detector (SSD)

One-stage detector

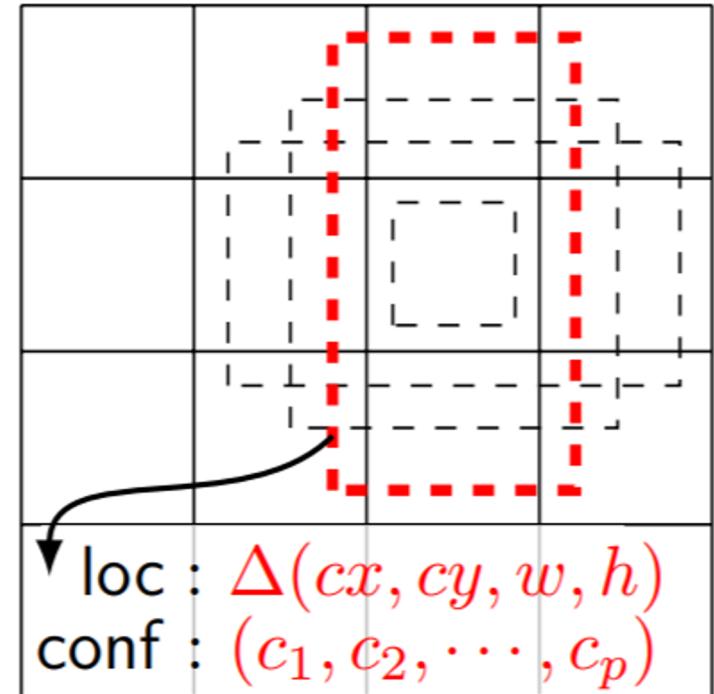
[Liu et al., ECCV 2016]



(a) Image with GT boxes



(b) 8×8 feature map



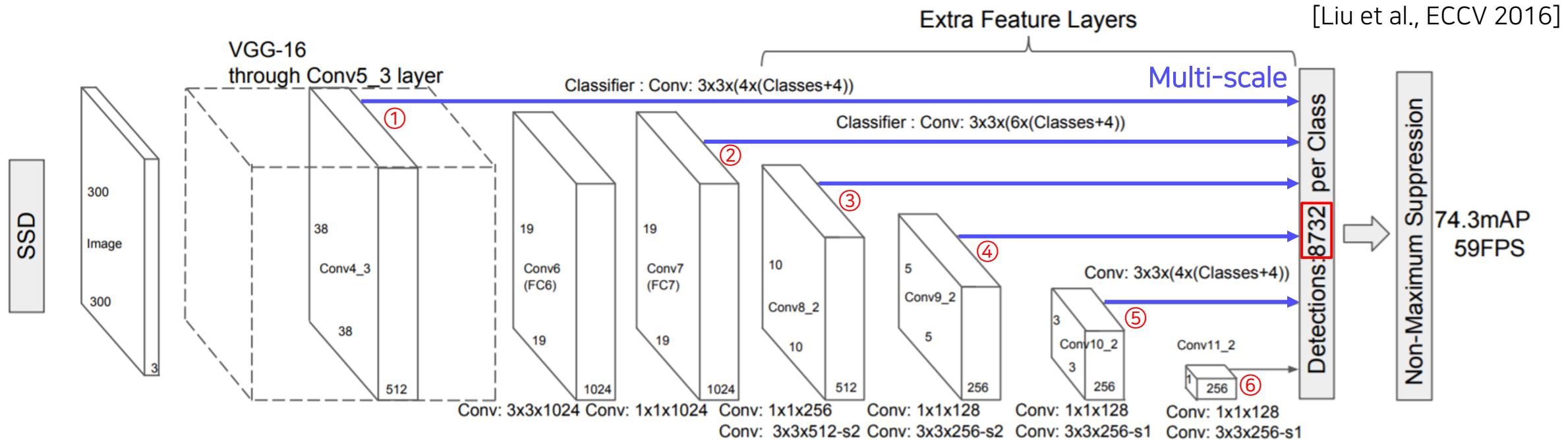
loc : $\Delta(cx, cy, w, h)$
conf : (c_1, c_2, \dots, c_p)

(c) 4×4 feature map

The use of **multi-scale** outputs attached to **multiple feature maps** enable effectively modeling a diverse space of possible box shapes

3.2 SSD

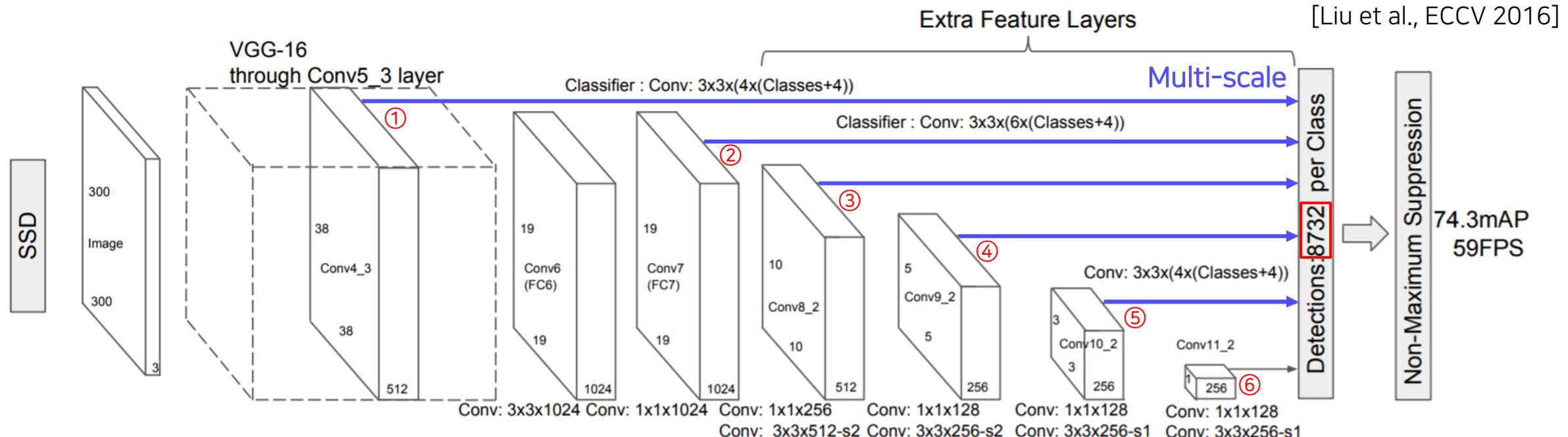
One-stage detector



$$\begin{matrix} * \\ \text{Feature map size } W \times H \\ \text{Number of channels } C \end{matrix} = \left[\begin{matrix} W \\ H \\ C \end{matrix} \right]$$

3.2 SSD

One-stage detector



$$* \begin{matrix} H \\ W \\ C \end{matrix} = \left[\begin{matrix} \text{Feature map size } W \times H \\ \text{Number of channels } C \end{matrix} \right]$$

Number of anchor boxes
for each pixel in feature map

Feature map size
 $* \frac{W \times H}{\uparrow} \times \frac{N}{\uparrow}$

- ① $38 \times 38 \times 4 = 5776$
- ② $19 \times 19 \times 6 = 2166$
- ③ $10 \times 10 \times 6 = 600$
- ④ $5 \times 5 \times 6 = 150$
- ⑤ $3 \times 3 \times 4 = 36$
- ⑥ $1 \times 1 \times 4 = 4$

Total number of anchor boxes
 $\rightarrow 5776 + 2166 + 600 + 150 + 36 + 4 = 8732$

3.2 SSD

One-stage detector

[Liu et al., ECCV 2016]

* SSD300 = SSD for 300×300 input

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 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.

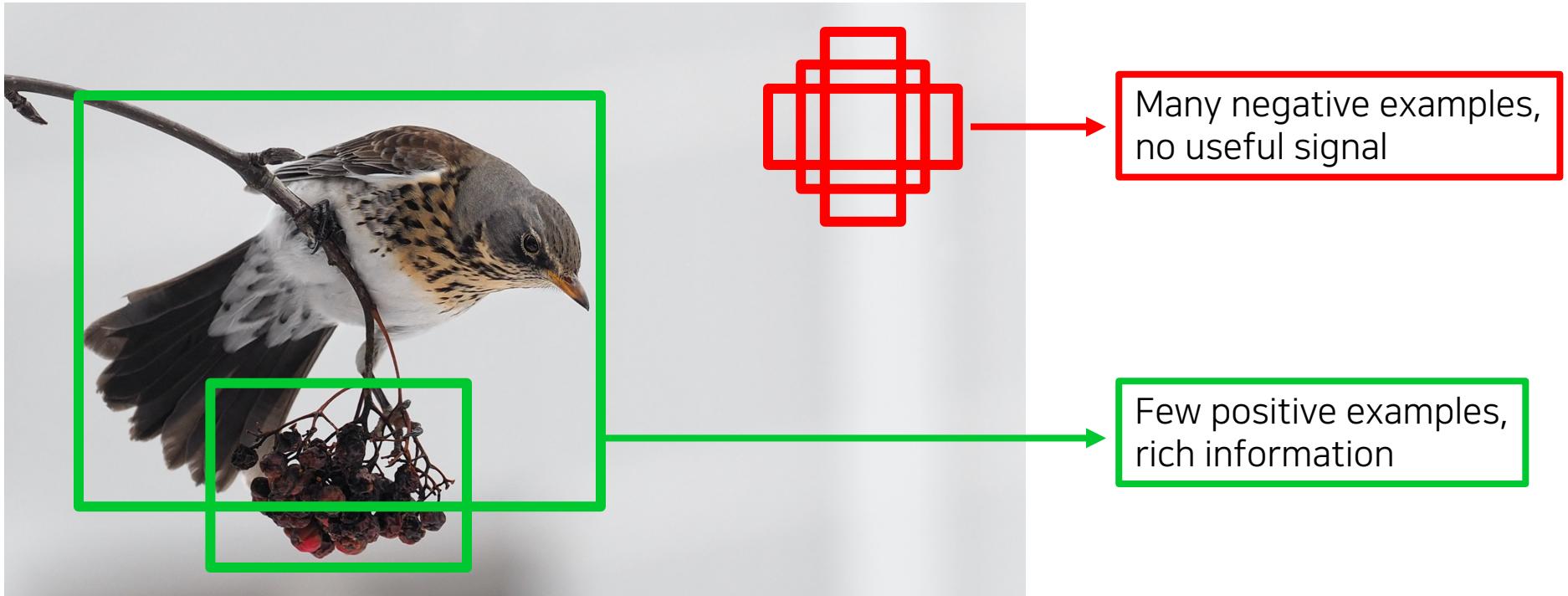
4.

Two-stage detector vs. one-stage detector

4.1 Focal loss

Two-stage detector vs. one-stage detector

Class imbalance problem



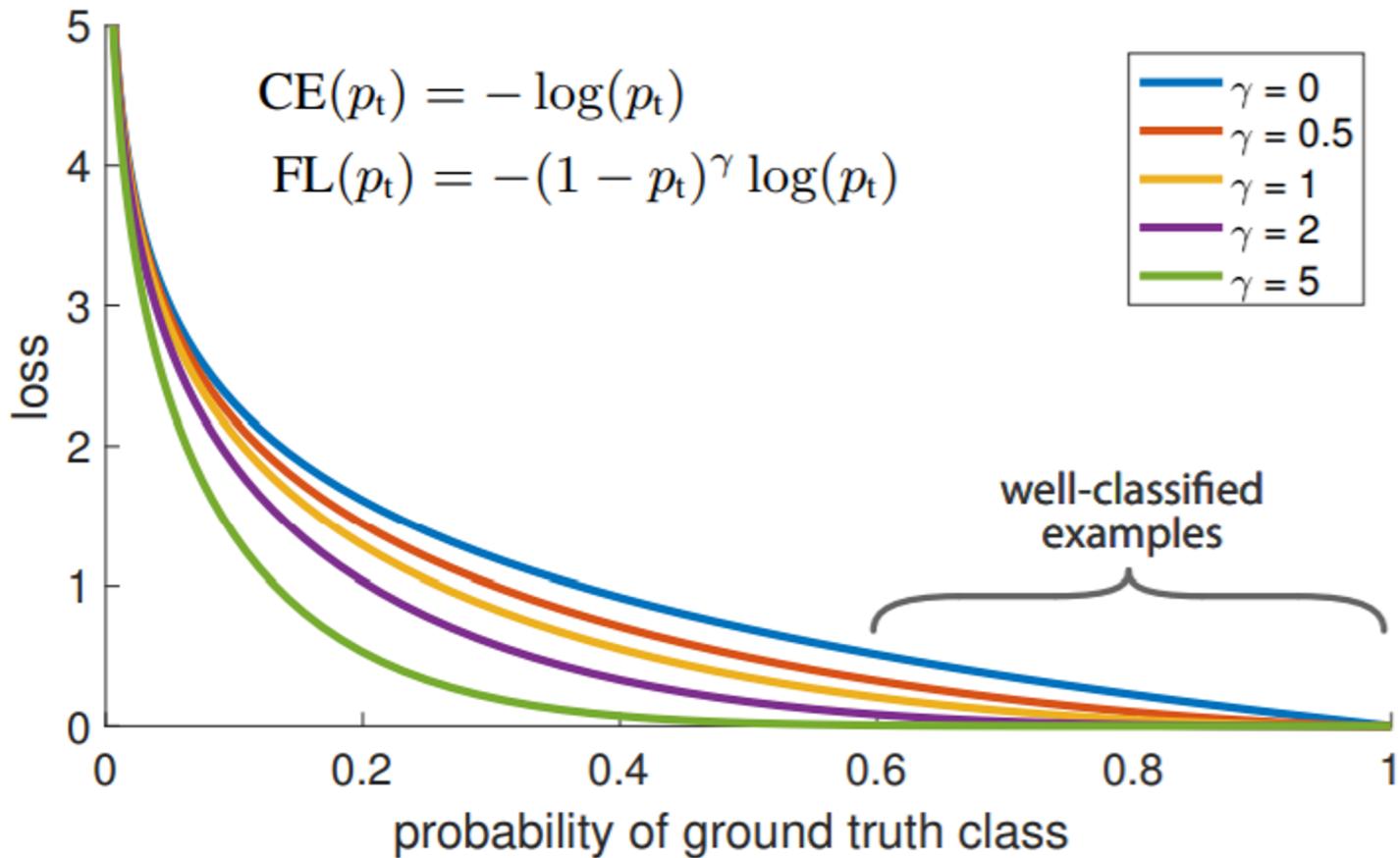
Class imbalance (#neg. anchor boxes \gg #pos. anchor boxes)

Single-stage detectors (e.g., YOLO, SSD) are prone to this problem...

4.1 Focal loss

Two-stage detector vs. one-stage detector

[Lin et al., ICCV 2017]



Focal loss

- Improved cross entropy loss
- Deal with class imbalance
 - Over-weights hard or misclassified examples
 - Down-weights easy examples

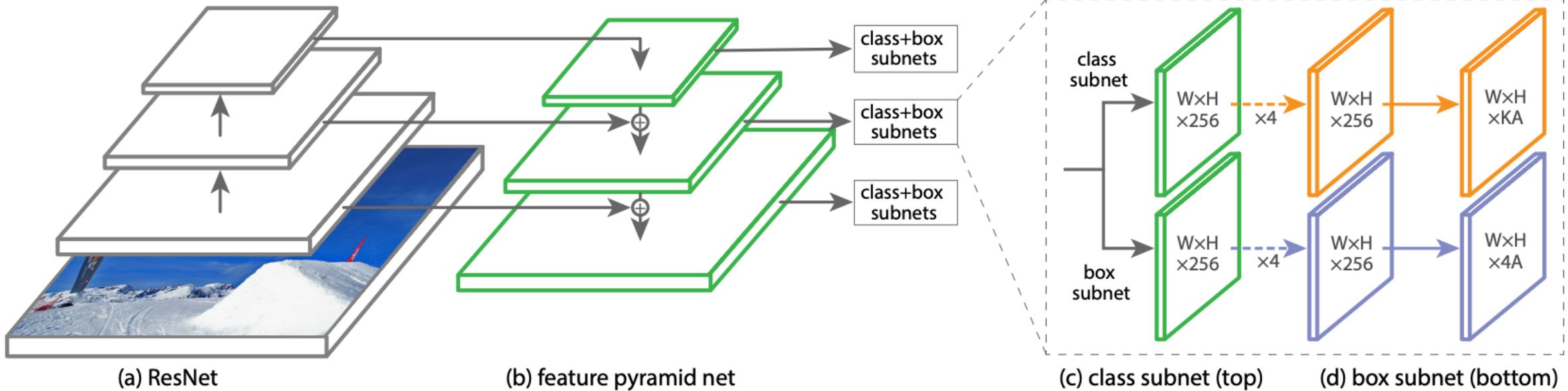
4.2 RetinaNet

Two-stage detector vs. one-stage detector

RetinaNet is a one-stage network

[Lin et al., ICCV 2017]

Feature Pyramid Networks (FPN) + class/box prediction branches

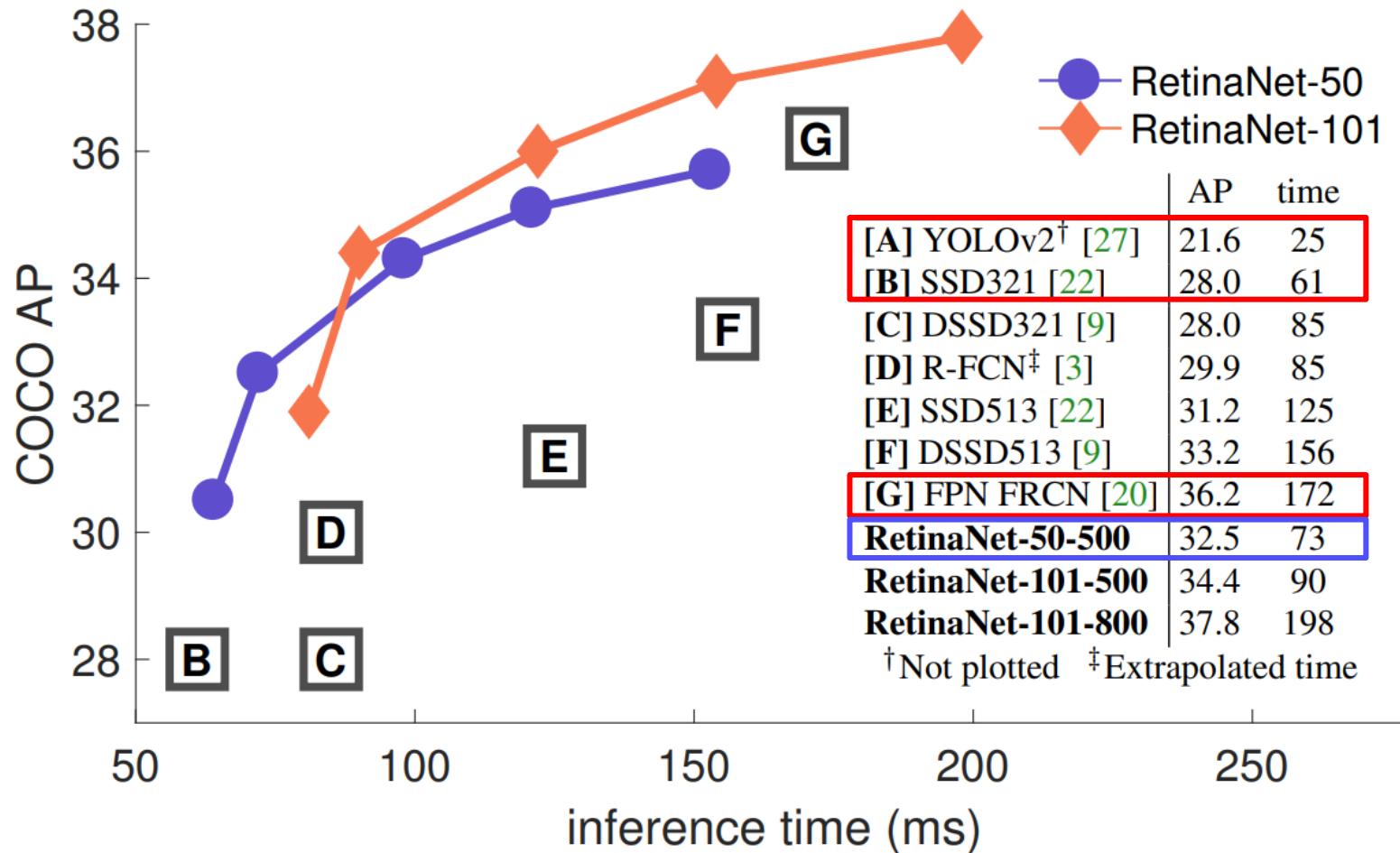


4.2 RetinaNet

Two-stage detector vs. one-stage detector

[Lin et al., ICCV 2017]

* RetinaNet-50-500 = ResNet-50 backbone RetinaNet for 500×500 input



5.

Detection with Transformer

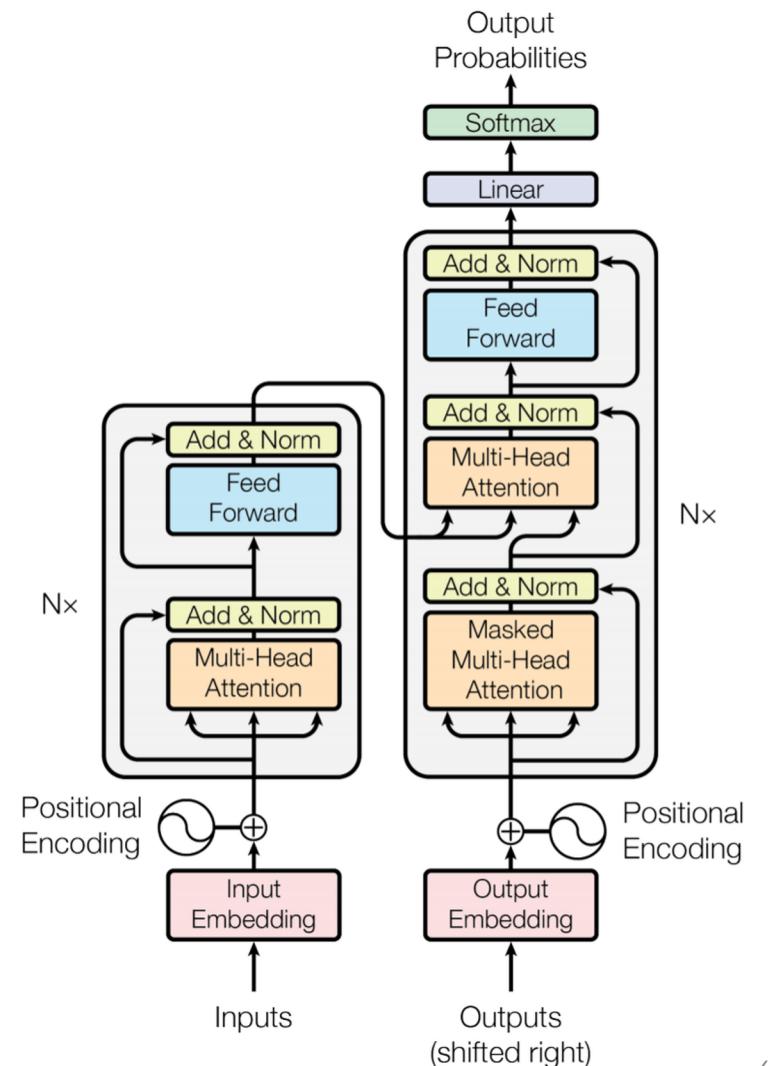
Transformer

Transformer has shown a great success in NLP

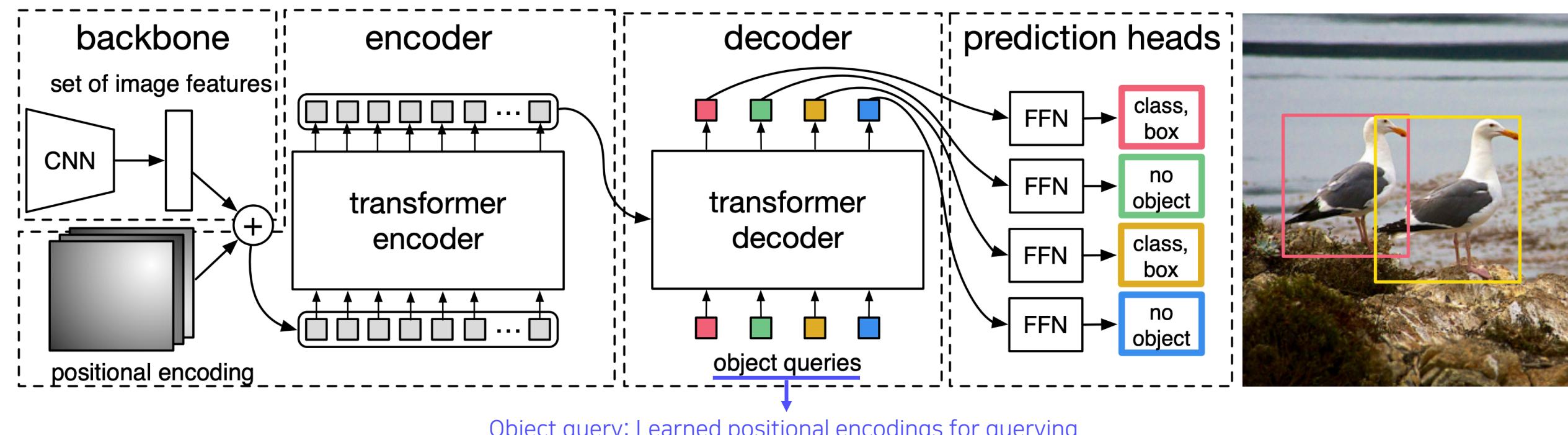
Why not extending Transformer to computer vision tasks!

- ViT (Vision Transformer) by Google
- DeiT (Data-efficient image Transformer) by Facebook
- DETR (DEtection TRansformer) by Facebook

[Vaswani et al., NeurIPS 2017]



[Carion et al., ECCV 2020]



Hyper-parameter:
(N object queries = max. N objects exist in a single image)

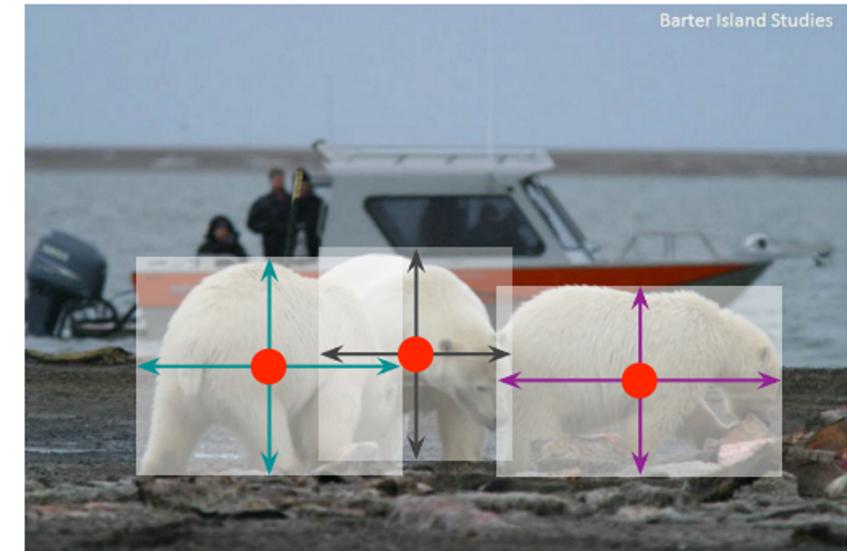
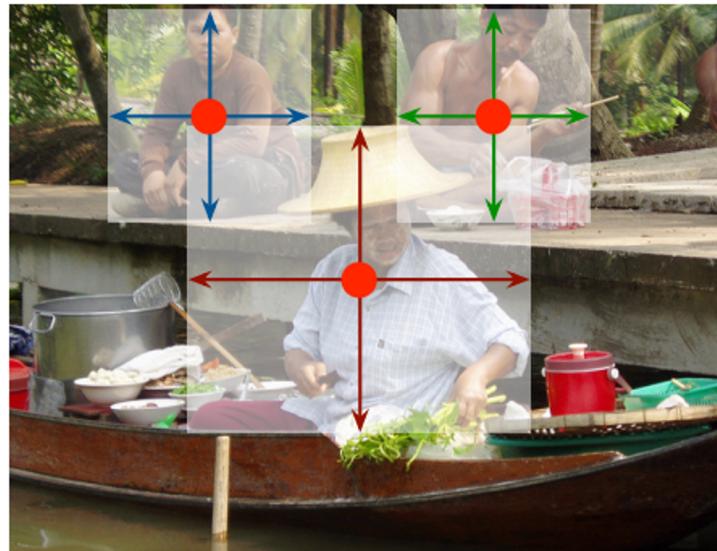
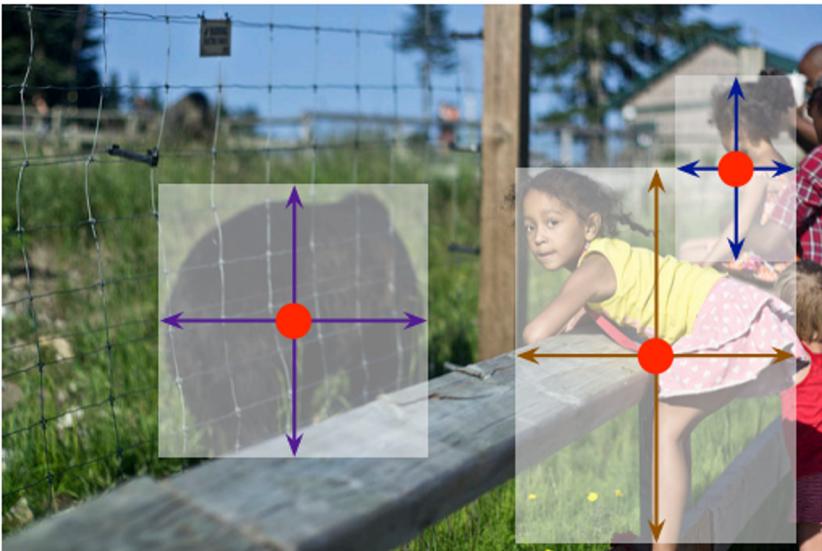
Further reading

Detection with Transformer

Detecting objects as points

[Zhou et al., arXiv 2019]

- Bounding box can be represented by other ways (left-top, right-bottom, centroid & size)
- Idea: Let's detect objects using corresponding points!
- CornerNet/CenterNet will be covered in Lecture 7



Reference

1. Object detection

- Kirillov et al., Panoptic Segmentation, CVPR 2019

2. Two-stage detector (R-CNN Family)

- Dalal et al., Histograms of Oriented Gradients for Human Detection, CVPR 2005
- Uijlings et al., Selective Search for Object Recognition, IJCV 2013
- Girshick et al., Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, CVPR 2014
- Girshick et al., Fast R-CNN, ICCV 2015
- Ferguson et al., Detection and Segmentation of Manufacturing Defects with Convolutional Neural Networks and Transfer Learning, Smart and Sustainable Manufacturing Systems 2018
- Ren et al., Faster R-CNN: Towards Real-Time Object detection with Region Proposal Networks, NeurIPS 2015

3. Single-stage detector

- Ndonhong et al., Wellbore Schematics to Structured Data Using Artificial Intelligence Tools, Offshore Technology Conference 2019
- Redmon et al., You Only Look Once: Unified, Real-Time Object detection, CVPR 2016
- Liu et al., SSD: Single Shot MultiBox Detector, ECCV 2016

×

Reference

4. Single-stage detector vs. two-stage detector

- Lin et al., Focal loss for Dense Object detection, ICCV 2017

5. Detection with Transformer

- Vaswani et al., Attention Is All You Need, NeurIPS 2017
- Carion et al., End-to-end Object Detection with Transformers, ECCV 2020
- Zhou et al., Objects as Points, arXiv 2019

End of Document

Thank You.