



Computer Vision

Chapter 7 (part 3):
Deep Learning for CV

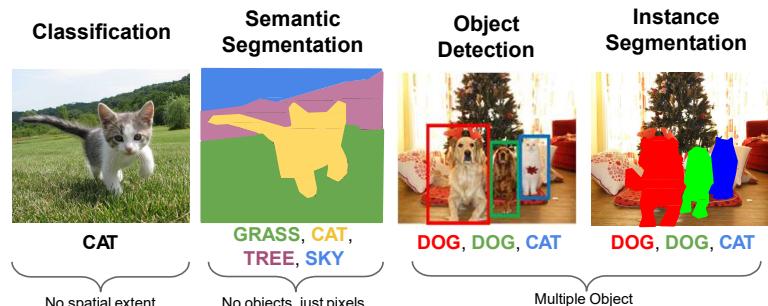
Deep learning for CV

Contents

1. Object detection: sliding-windows
2. Two-stage object detection
3. One-stage object detection: Anchor-based
4. One-stage object detection: Anchor-free
5. Semantic segmentation
6. Instance segmentation



Computer Vision Tasks



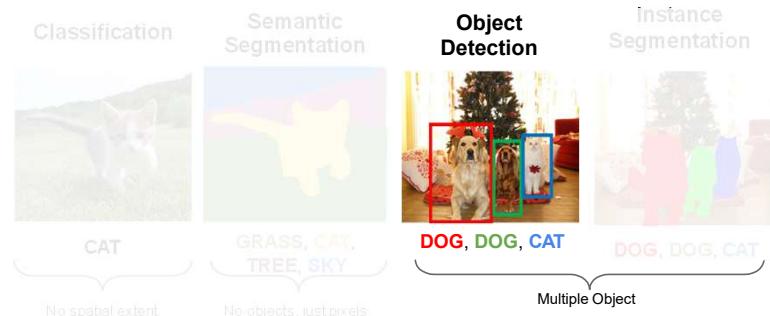
Object detection: sliding windows



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

5

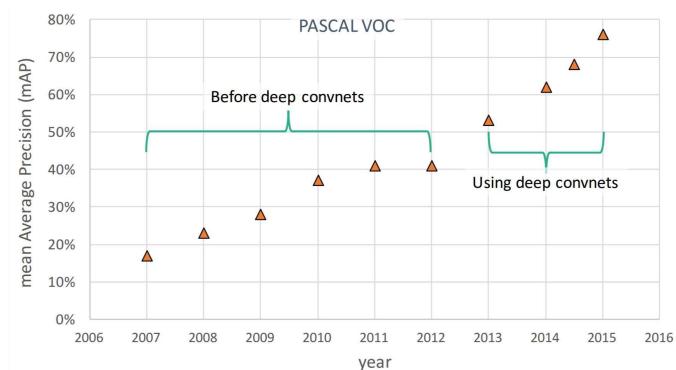
Object Detection



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

6

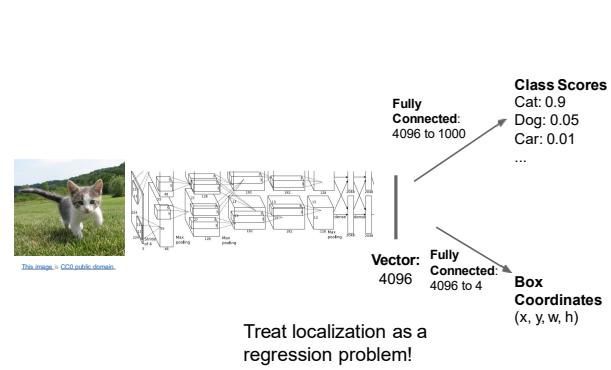
Object Detection: Impact of Deep Learning



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

7

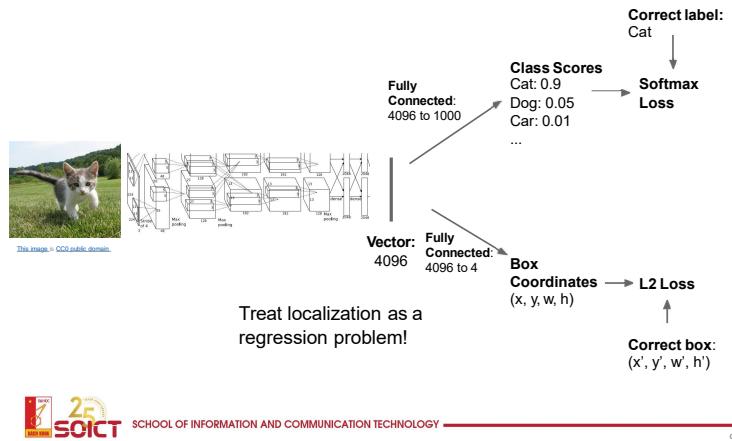
Object Detection: Single Object (Classification + Localization)



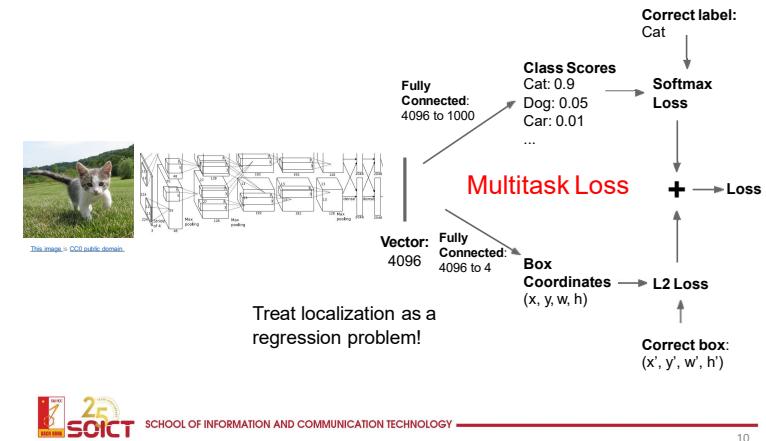
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

8

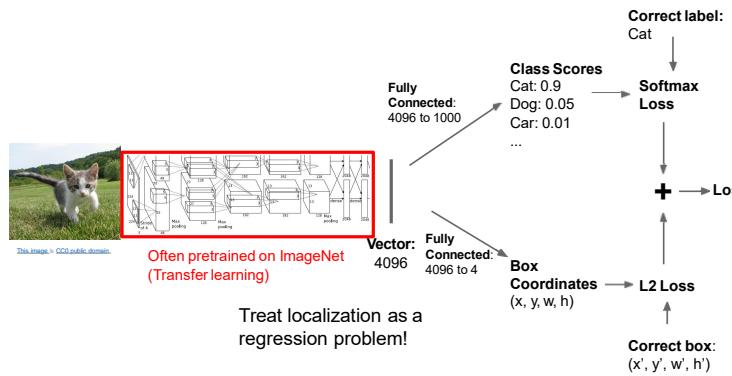
Object Detection: Single Object (Classification + Localization)



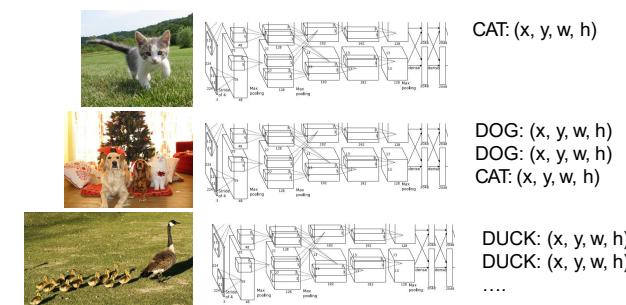
Object Detection: Single Object (Classification + Localization)



Object Detection: Single Object (Classification + Localization)

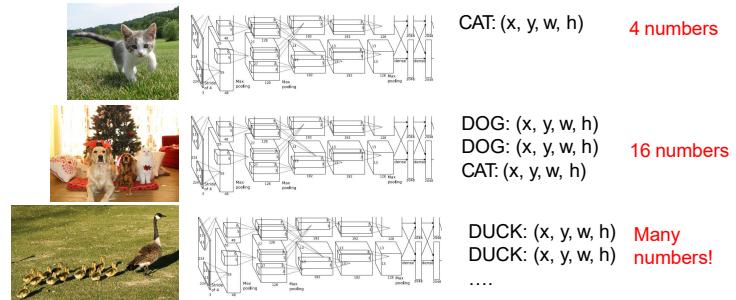


Object Detection: Multiple Objects



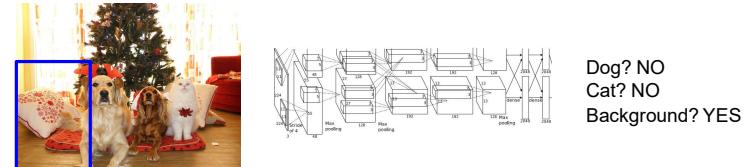
Object Detection: Multiple Objects

Each image needs a different number of outputs!



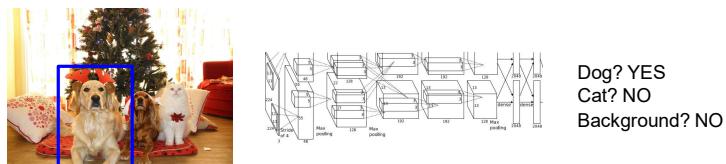
Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



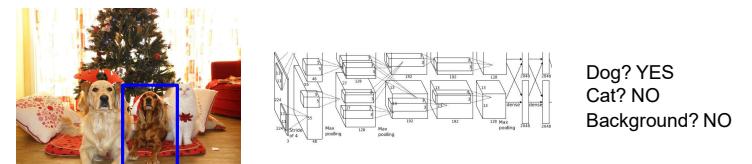
Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



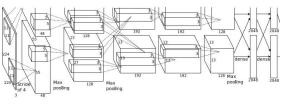
Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Object Detection: Multiple Objects

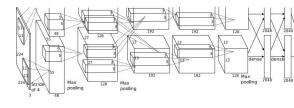
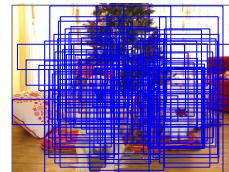
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? YES
Background? NO

Object Detection: Multiple Objects

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? YES
Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

17



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

18

Object Detection

Two Stages

- Propose “objects”
- Classify each candidate

One-Stage

- Sliding window to classify all candidates

Object Detection

Two Stages

- Propose “objects”
- Classify each candidate

One-Stage

- Sliding window to classify all candidates



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

19

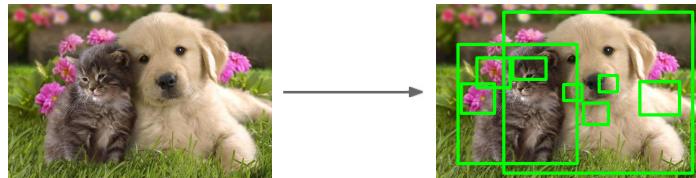


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

20

Region Proposals: Selective Search

- Find “blobby” image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



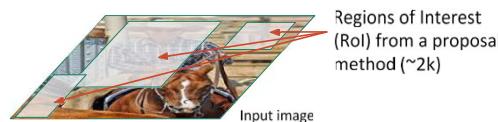
Alexe et al., "Measuring the objectness of image windows", TPAMI 2012 Uijlings et al., "Selective Search for Object Recognition", IJCV 2013
 Cheng et al., "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014 Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

R-CNN



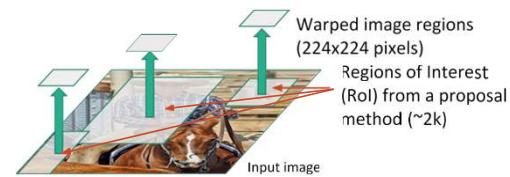
Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

R-CNN



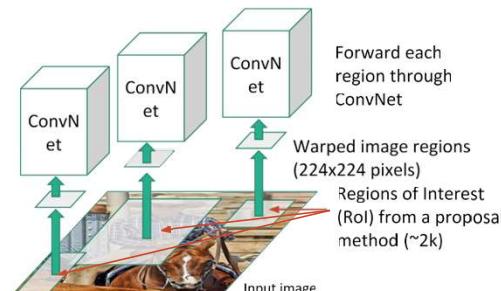
Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

R-CNN



Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

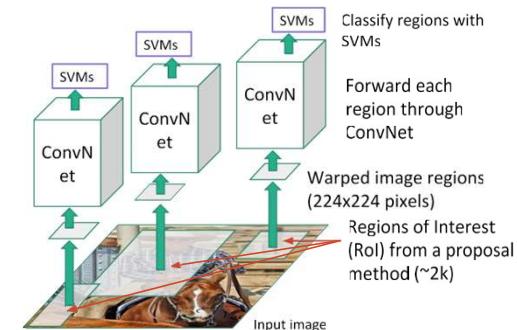
R-CNN



Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

25

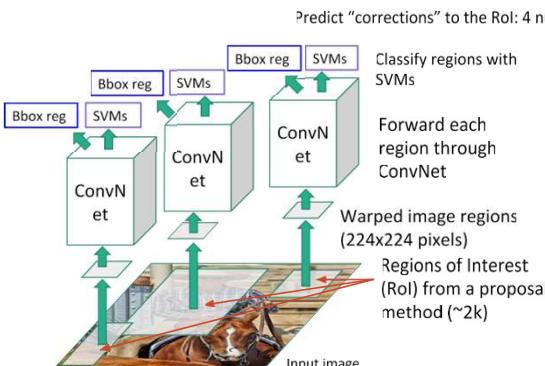
R-CNN



Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

26

R-CNN

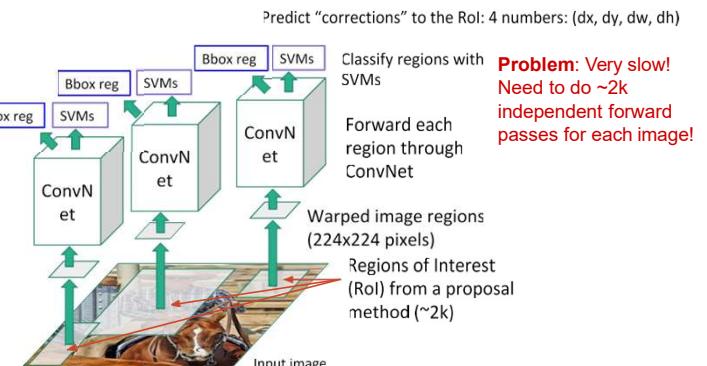


Predict "corrections" to the RoI: 4 numbers: (dx, dy, dw, dh)

Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

27

R-CNN

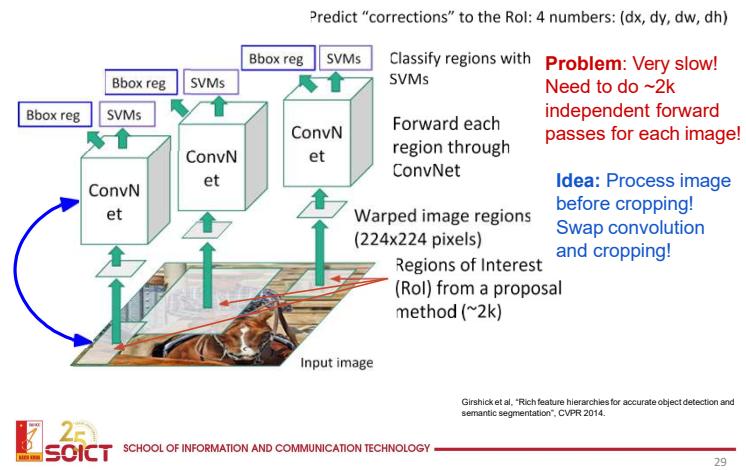


Predict "corrections" to the RoI: 4 numbers: (dx, dy, dw, dh)

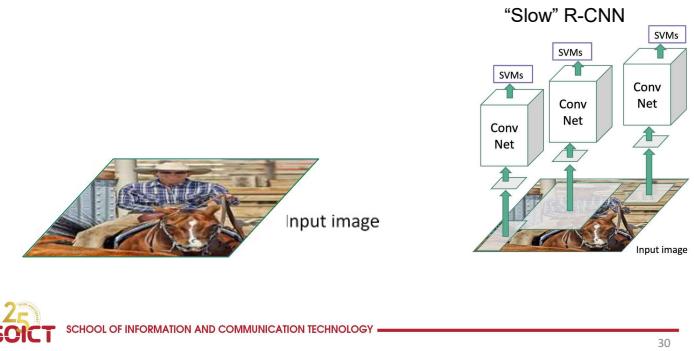
Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

28

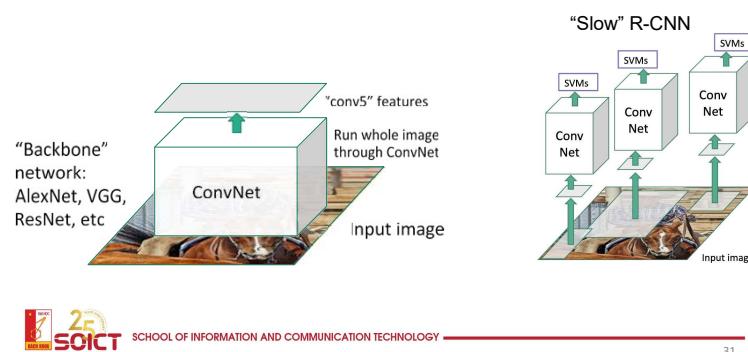
“Slow” R-CNN



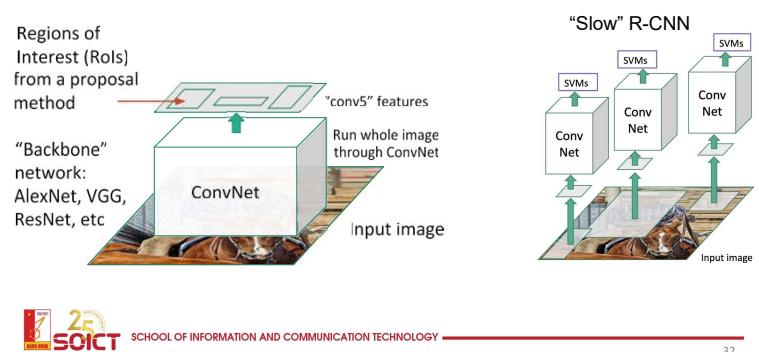
Fast R-CNN



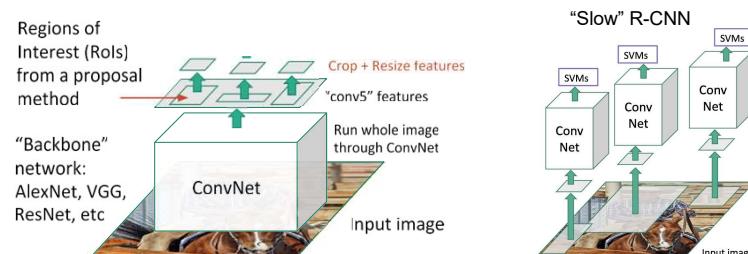
Fast R-CNN



Fast R-CNN



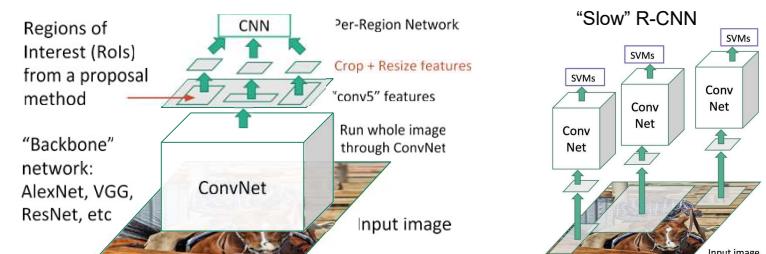
Fast R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

33

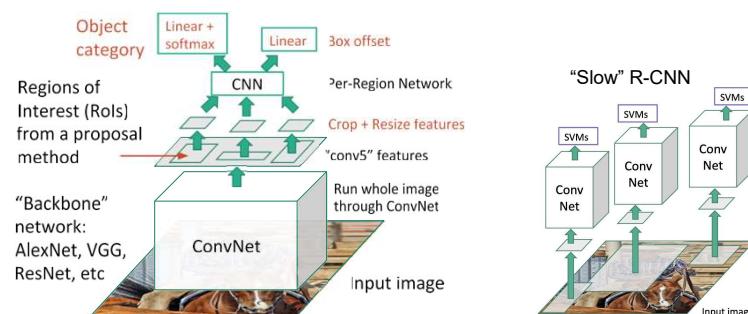
Fast R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

34

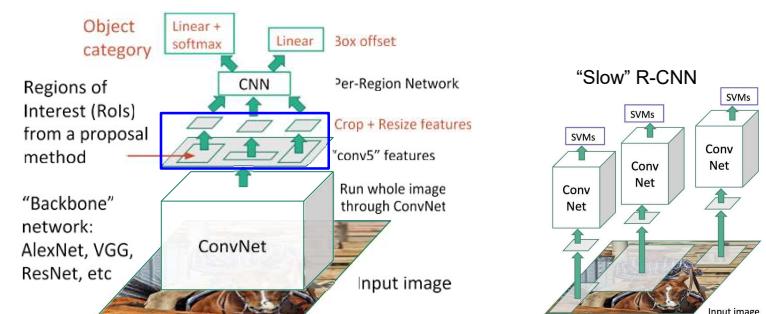
Fast R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

35

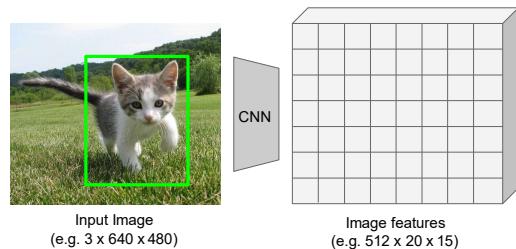
Fast R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

36

Cropping Features: RoI Pool



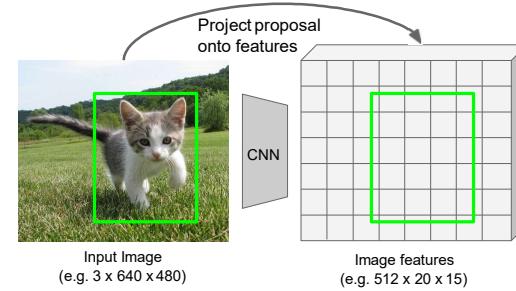
Girshick, "Fast R-CNN", ICCV2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

37

Cropping Features: RoI Pool



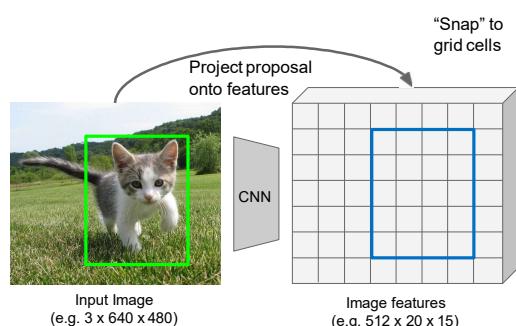
Girshick, "Fast R-CNN", ICCV2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

38

Cropping Features: RoI Pool



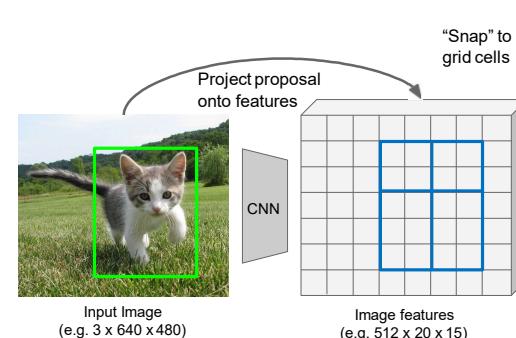
Girshick, "Fast R-CNN", ICCV2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

39

Cropping Features: RoI Pool



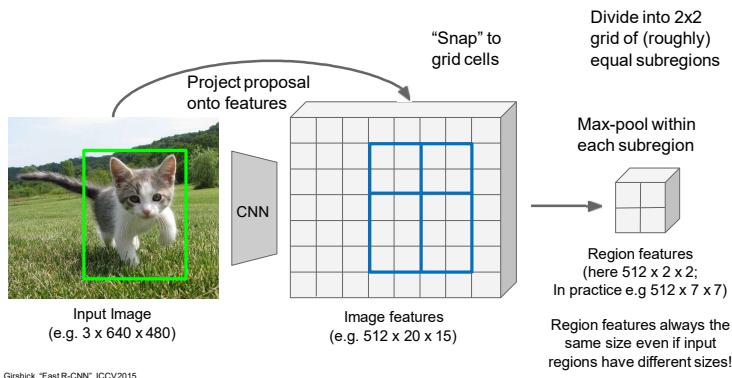
Girshick, "Fast R-CNN", ICCV2015.



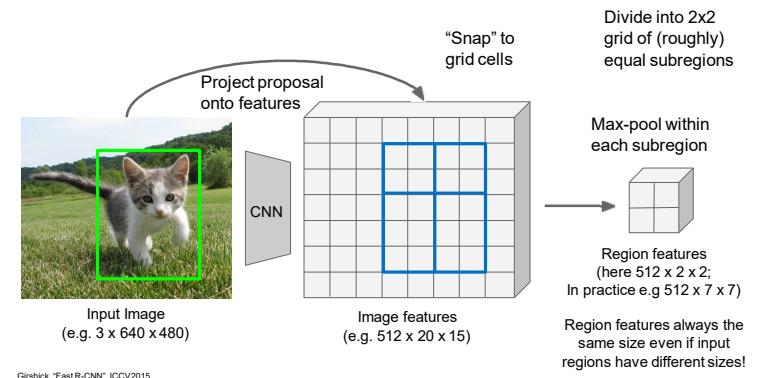
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

40

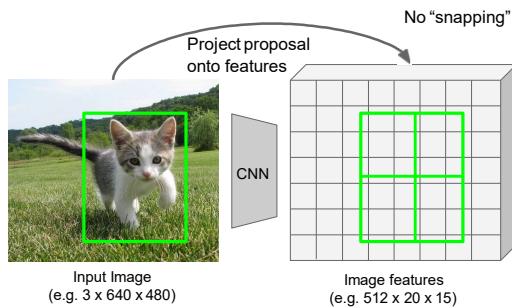
Cropping Features: RoI Pool



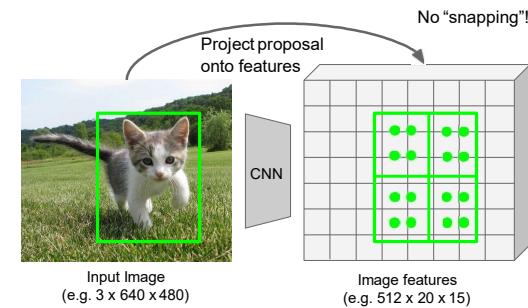
Cropping Features: RoI Pool



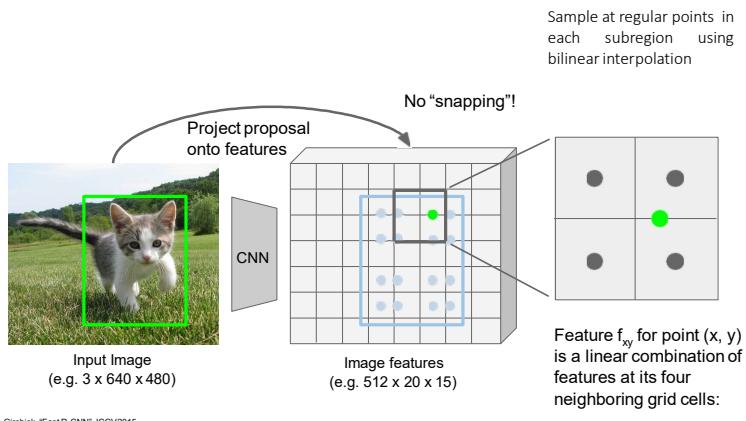
Cropping Features: RoI Align



Cropping Features: RoI Align



Cropping Features: RoI Align



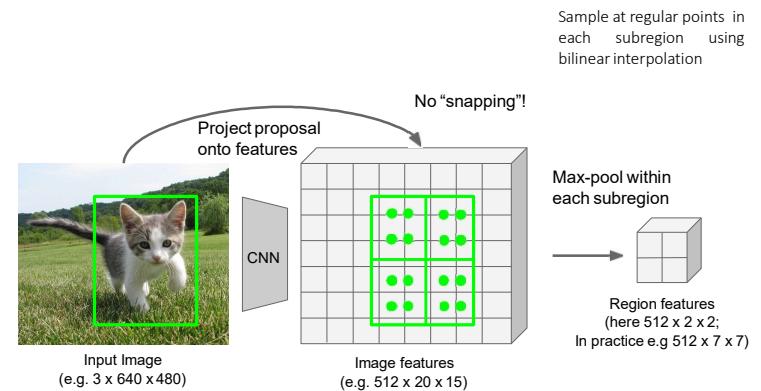
Girshick, "Fast R-CNN", ICCV2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

45

Cropping Features: RoI Align



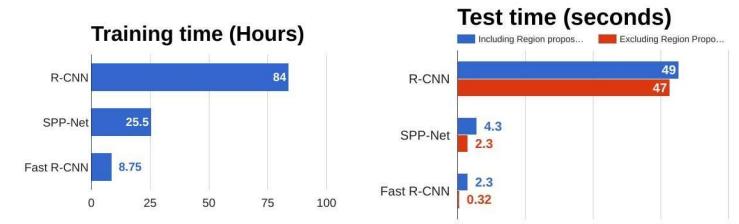
Girshick, "Fast R-CNN", ICCV2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

46

R-CNN vs Fast R-CNN



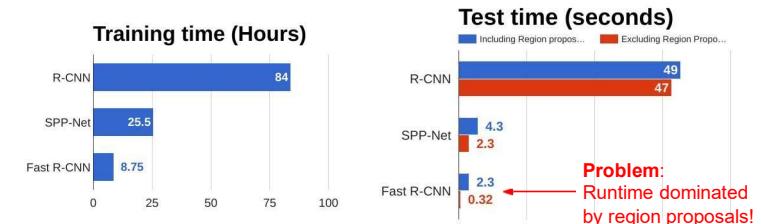
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
 He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
 Girshick, "Fast R-CNN", ICCV 2015



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

47

R-CNN vs Fast R-CNN



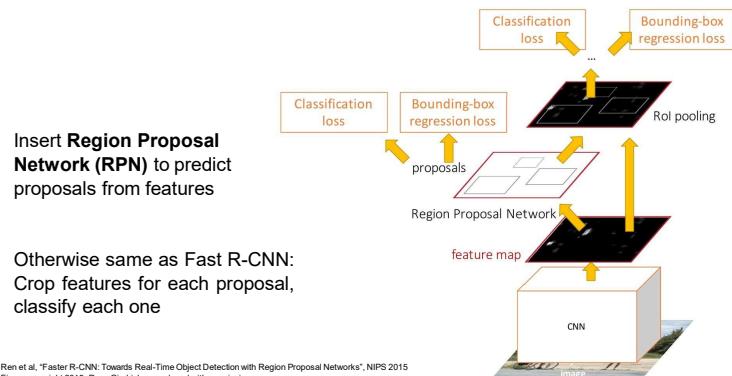
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
 He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
 Girshick, "Fast R-CNN", ICCV 2015



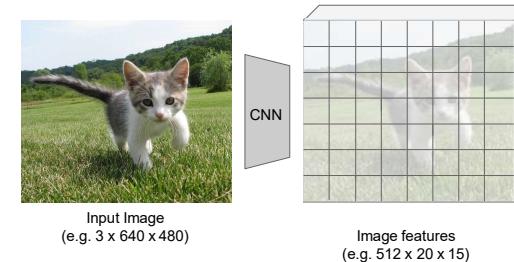
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

48

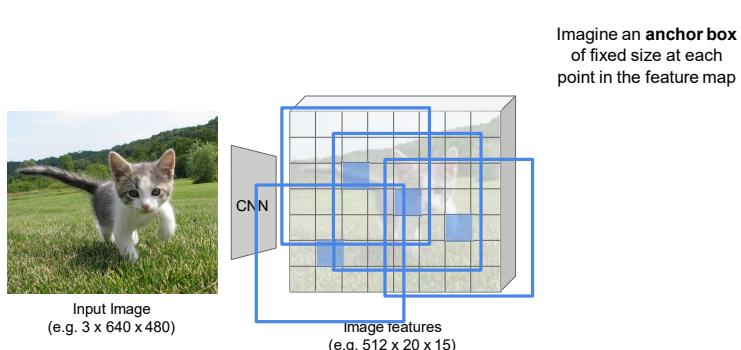
Faster R-CNN: Make CNN do proposals!



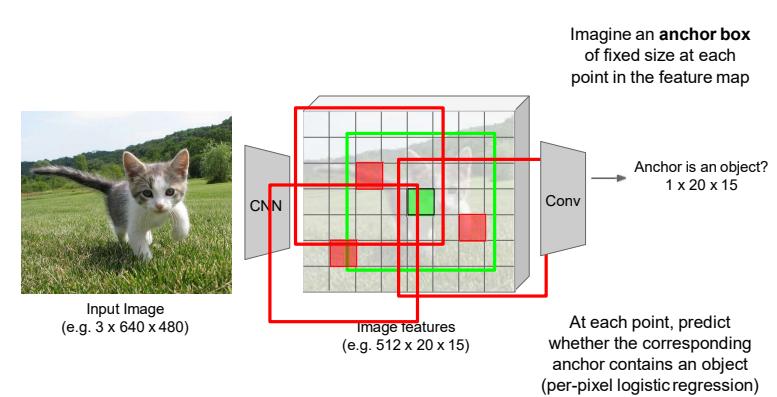
Region Proposal Network



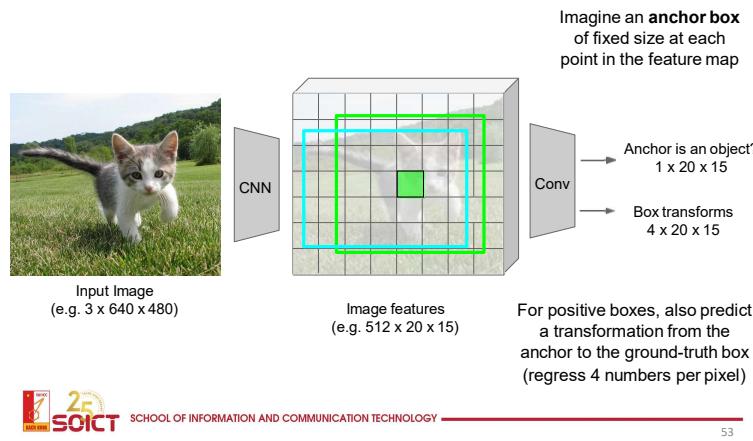
Region Proposal Network



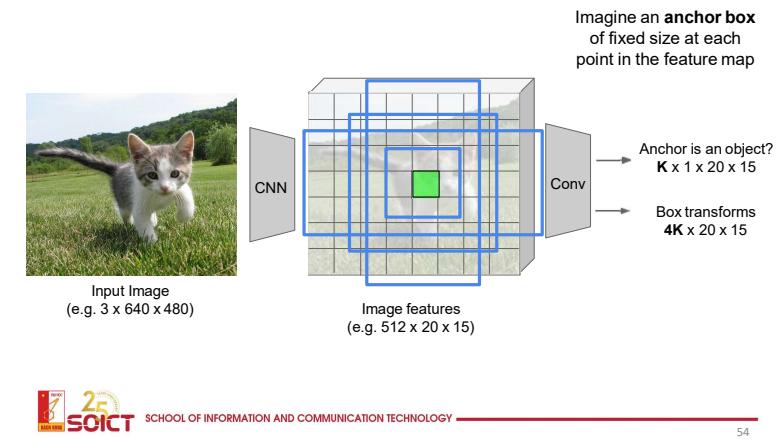
Region Proposal Network



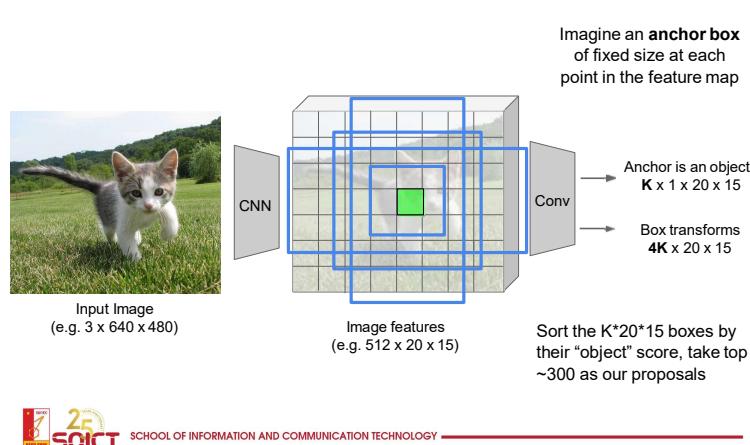
Region Proposal Network



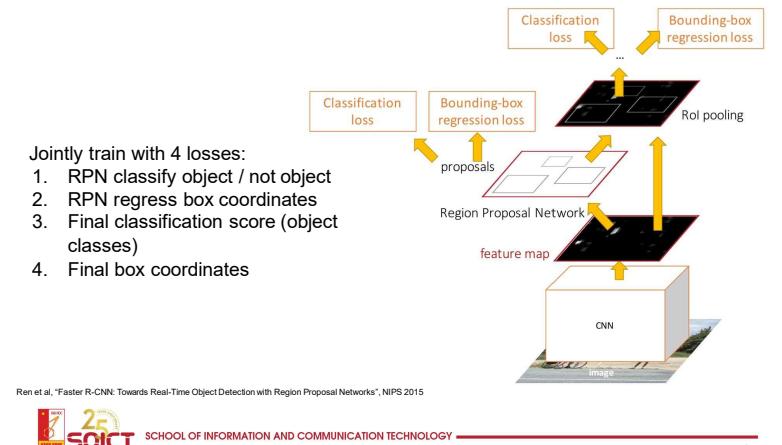
Region Proposal Network



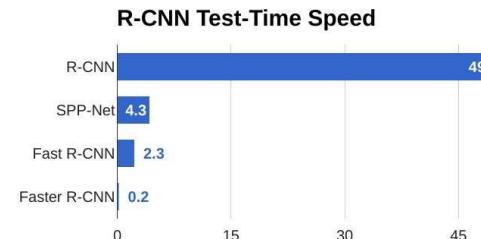
Region Proposal Network



Faster R-CNN: Make CNN do proposals!



Faster R-CNN: Make CNN do proposals!



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

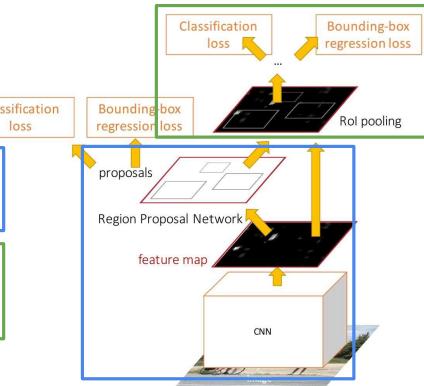
57

Faster R-CNN: Make CNN do proposals!

Faster R-CNN is a
Two-stage object detector

- First stage: Run once per image
 - Backbone network
 - Region proposal network

- Second stage: Run once per region
 - Crop features: RoI pool / align
 - Predict object class
 - Prediction bbox offset



Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

59

One-stage Object detection

Anchor-based



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

61

Object Detection

Two Stages

- Propose “objects”
- Classify each candidate

One-Stage

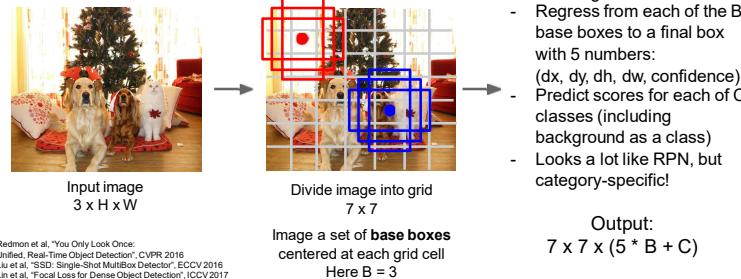
- Sliding window to classify all candidates



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

62

Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Redmon et al., "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016
 Liu et al., "SSD: Single-Shot MultiBox Detector", ECCV 2016
 Lin et al., "Focal Loss for Dense Object Detection", ICCV 2017



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

63

Imbalance – Focal Loss

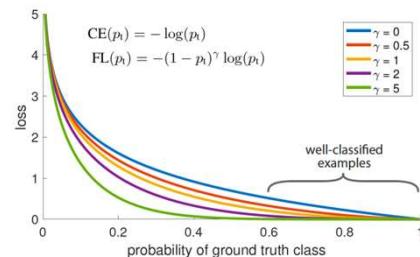


Figure 1: We propose a novel loss we term the *Focal Loss* that adds a factor $(1 - p_t)^\gamma$ to the standard cross entropy criterion. Setting $\gamma > 0$ reduces the relative loss for well-classified examples ($p_t > .5$), putting more focus on hard, misclassified examples. As our experiments will demonstrate, the proposed focal loss enables training highly accurate dense object detectors in the presence of vast numbers of easy background examples.



Lin, Goyal, Girshick, He, and Dollár
Focal loss for dense object detection (PAMI 2018)

65

Imbalance

Number of “negative” anchors $\sim O(10K)$
 Number of “positive” anchors $\sim O(10)$

What happens to CE loss in this case?

$$\mathcal{L} = -\frac{1}{N} \sum_i y_i \log(p_i), \quad \frac{\partial \mathcal{L}}{\partial p_i} = \begin{cases} p_{i,l} & \text{if } y_i \neq l \\ p_{i,l} - 1 & \text{if } y_i = l \end{cases}$$

Loss and gradient are dominated by correctly classified negative examples

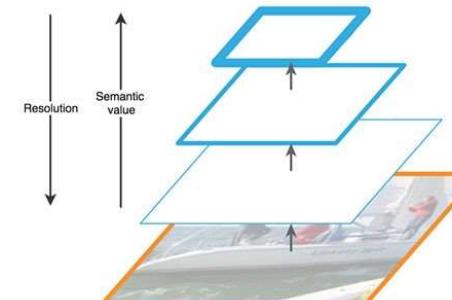
Training outcome \rightarrow constant “negative” prediction.



64

Feature Pyramid Network (FPN)

- How to handle multiscale predictions?

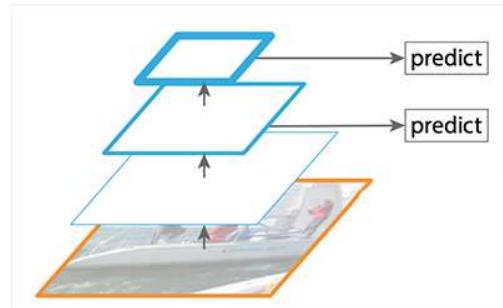


Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. **Feature Pyramid Networks for Object Detection** (CVPR 2017)

66

Feature Pyramid Network (FPN)

- How to handle multiscale predictions?

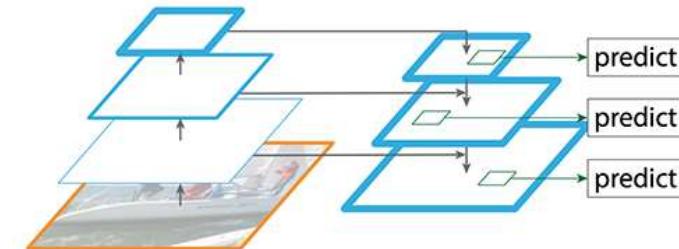


Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. Feature Pyramid Networks for Object Detection (CVPR 2017)

67

Feature Pyramid Network (FPN)

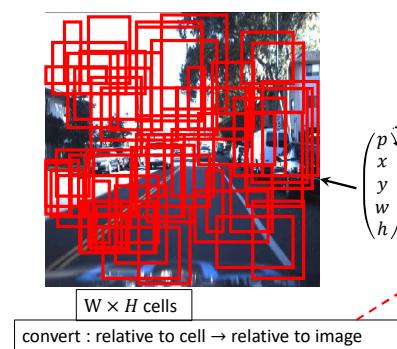
- How to handle multiscale predictions?



Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. Feature Pyramid Networks for Object Detection (CVPR 2017)

68

Postprocessing: NMS



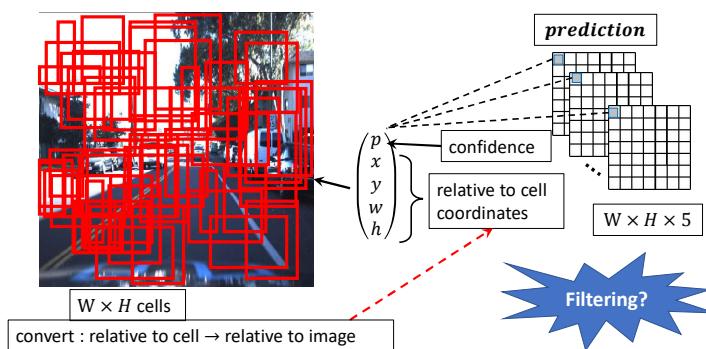
69



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

69

Postprocessing: NMS



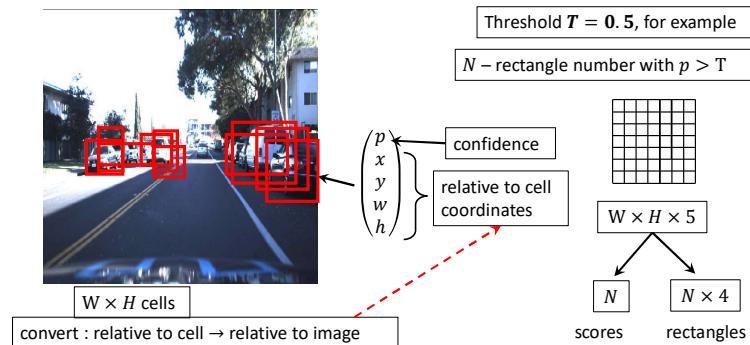
70



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

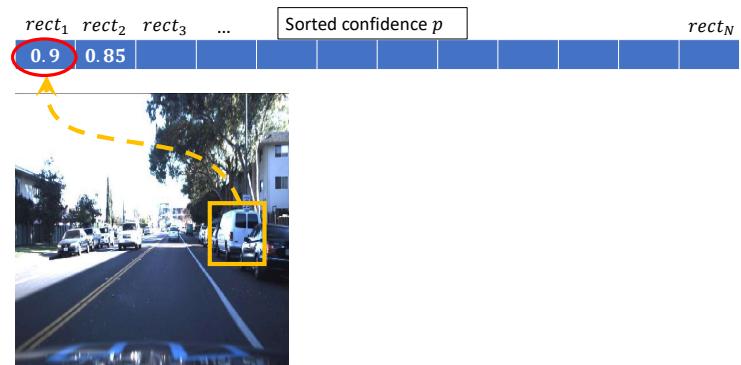
70

Postprocessing: NMS



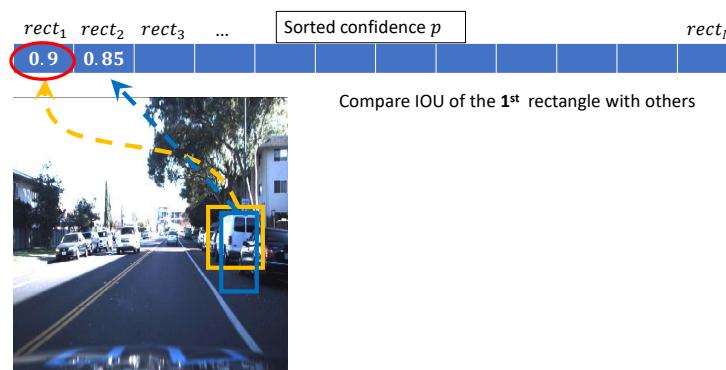
71

Postprocessing: NMS



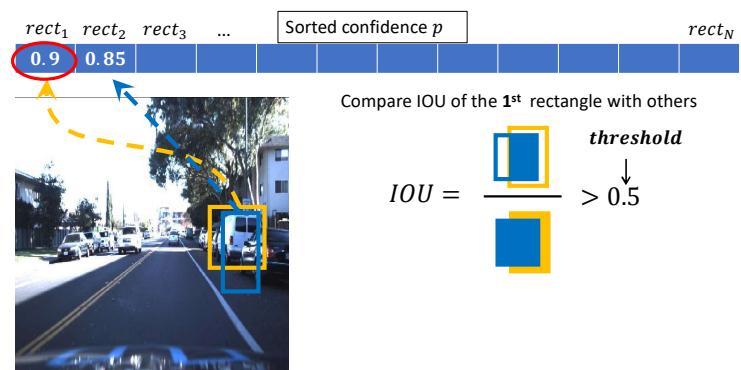
72

Postprocessing: NMS



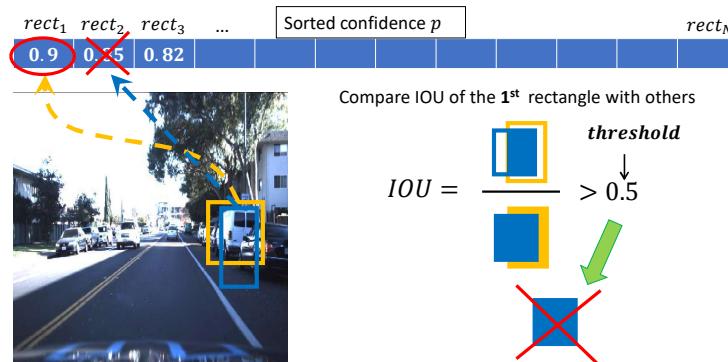
73

Postprocessing: NMS



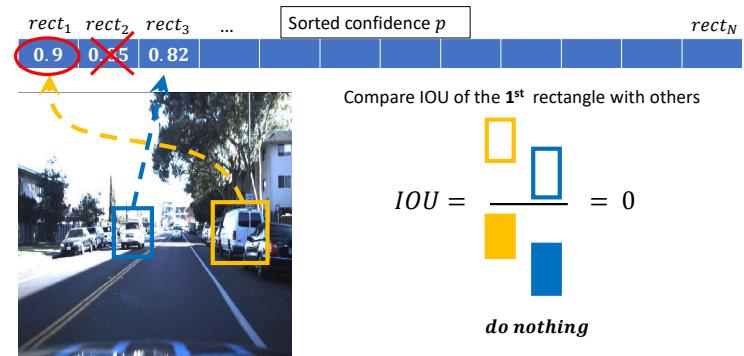
74

Postprocessing: NMS



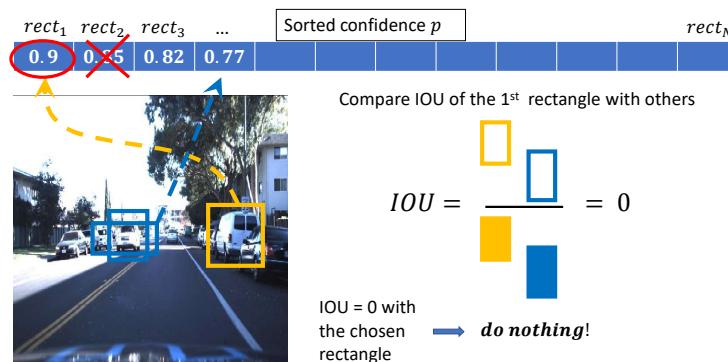
75

Postprocessing: NMS



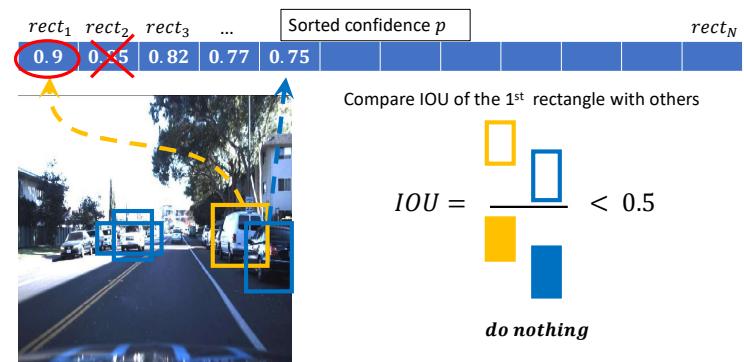
76

Postprocessing: NMS



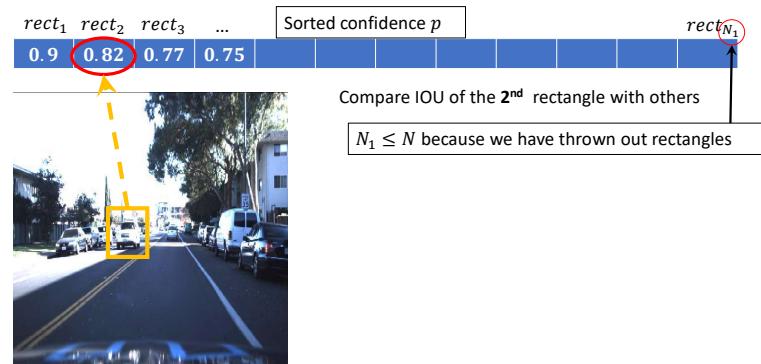
77

Postprocessing: NMS



78

Postprocessing: NMS



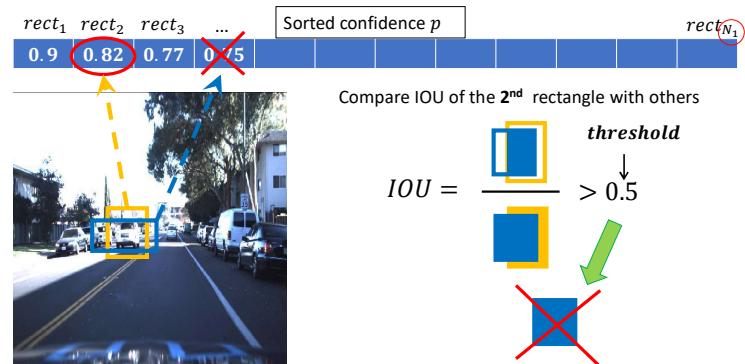
79



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

79

Postprocessing: NMS



I

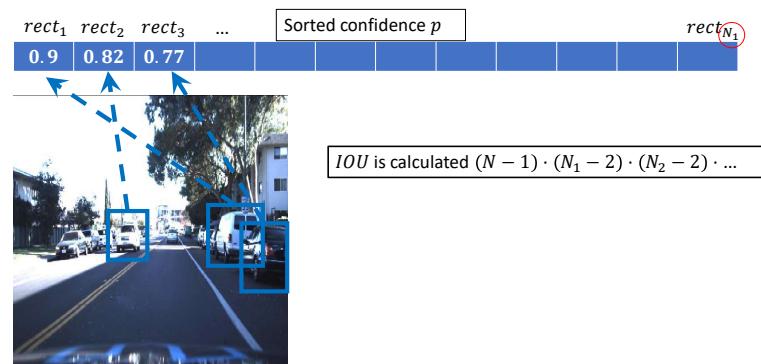
80



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

80

Postprocessing: NMS



81



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

81

One-stage Object detection

Anchor-free



82

Drawbacks of Anchor Boxes

1. Need a large number of anchors



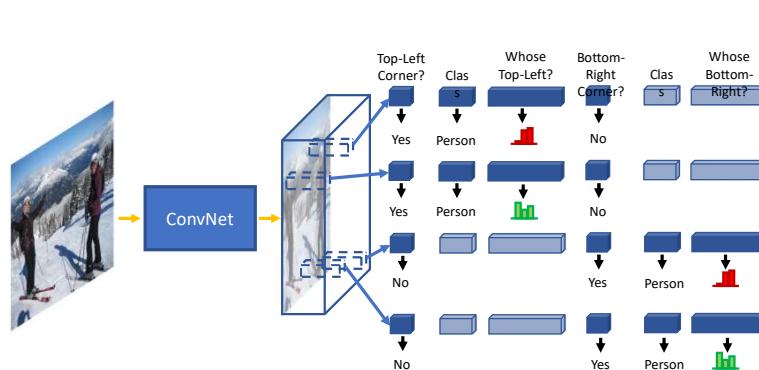
- A tiny fraction of anchors are positive examples
- Slow down training [Lin et al. ICCV'17]

2. Extra hyperparameters – sizes and aspect ratios

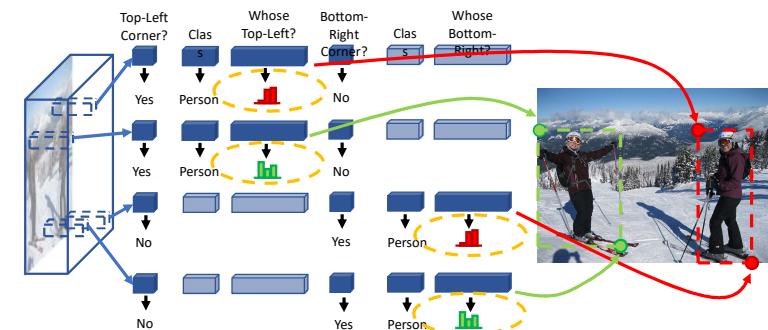
CornerNet: Detecting Objects as Paired Keypoints



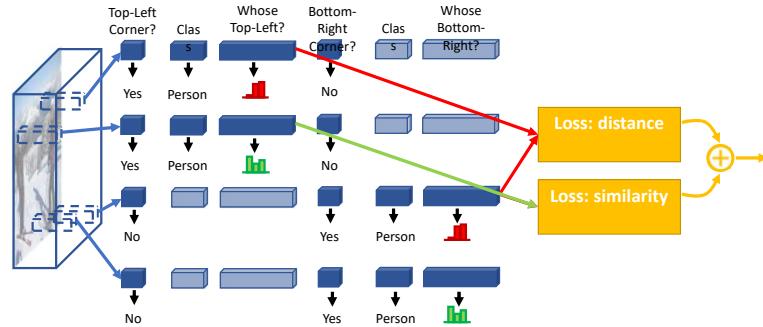
CornerNet: Detecting Objects as Paired Keypoints



CornerNet: Detecting Objects as Paired Keypoints



CornerNet: Detecting Objects as Paired Keypoints



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

87

Associative Embedding [Newell et al. NIPS'17]

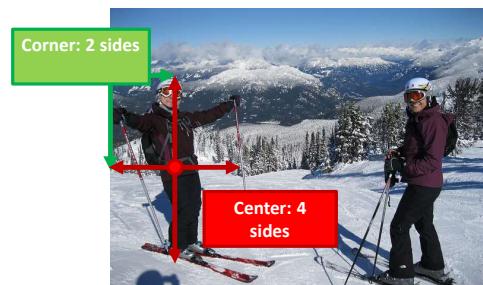
- CornerNet is inspired by the idea of associative embedding from work by Newell et al. [NIPS'17]



<https://proceedings.neurips.cc/paper/2017/file/8edd72158ccd2a879f79cb2538568fdc.pdf>

88

Advantages of Detecting Corners



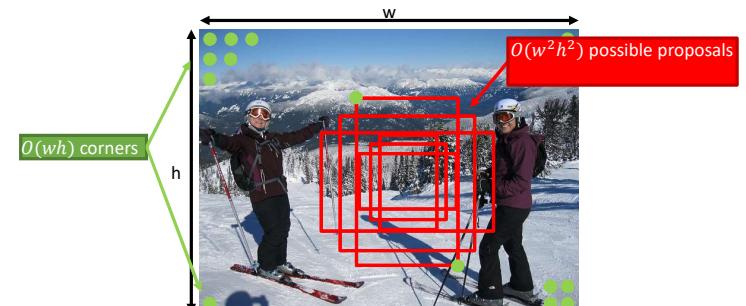
Detecting corner is easier than detecting center



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

89

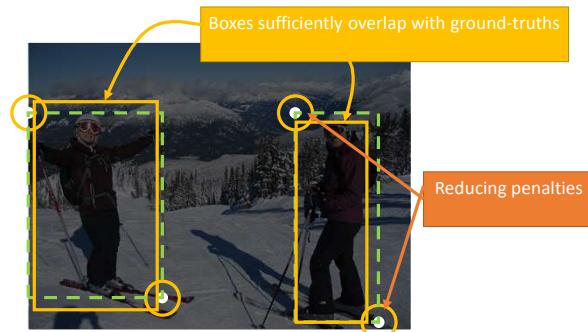
Advantages of Detecting Corner

Represent $O(w^2 h^2)$ possible proposals using only $O(wh)$ corners

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

90

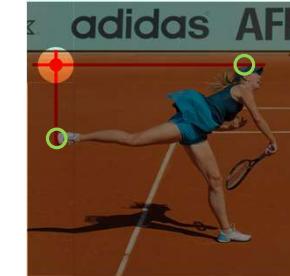
Supervising Corner Detection



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

91

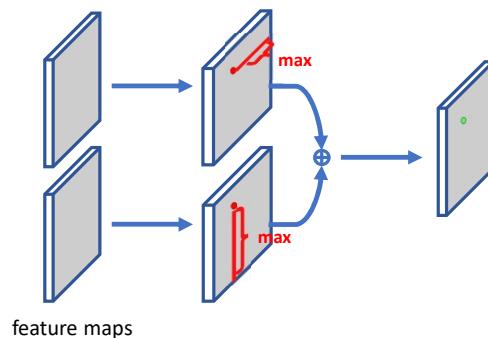
Corner Pooling



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

92

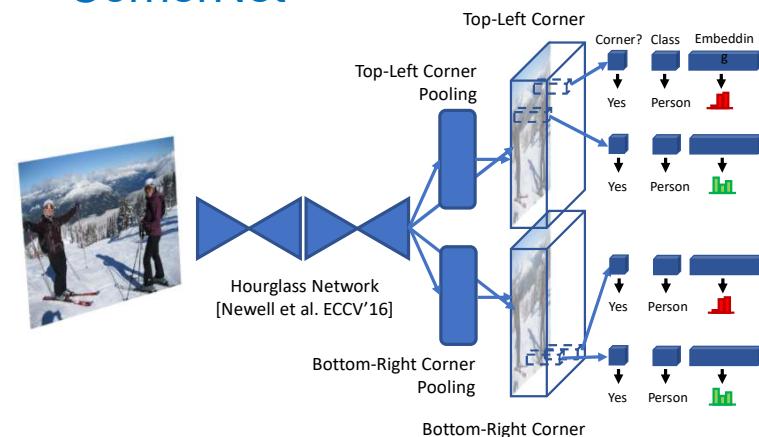
Top-Left Corner Pooling



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

93

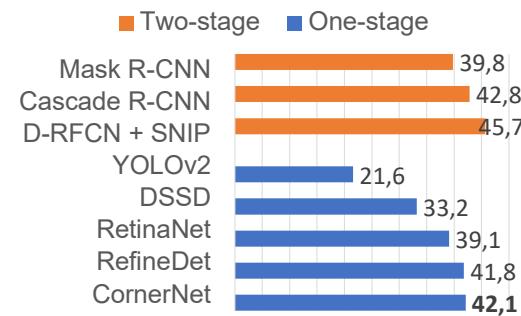
CornerNet



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

94

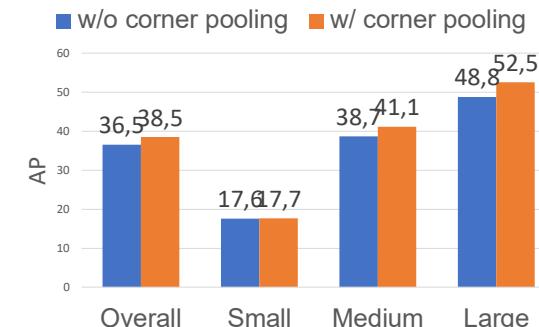
Experiment: CornerNet versus Others



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

95

Experiment: Corner Pooling



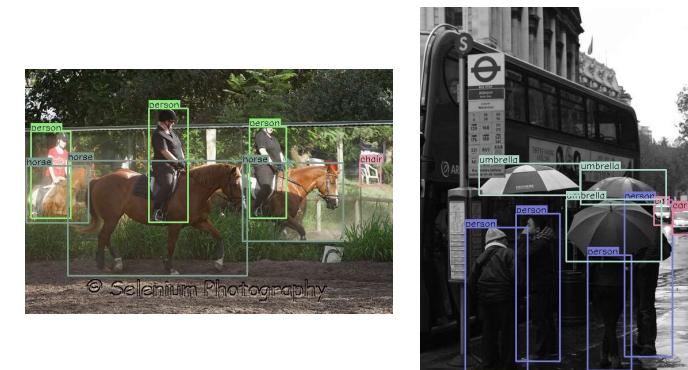
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

96



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

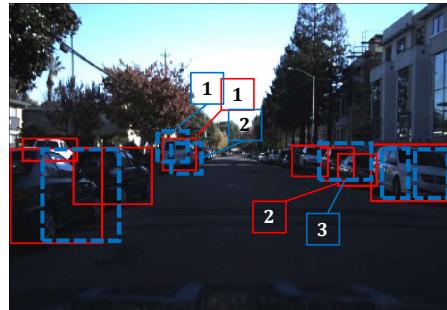
97



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

98

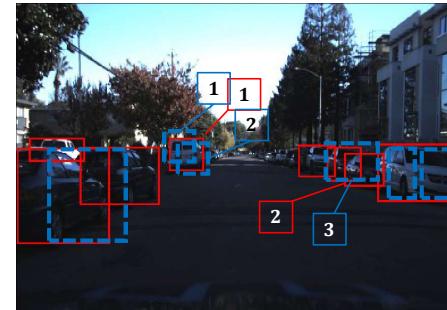
Accuracy evaluation



n predicted
 m Ground Truth
 $IoU = \frac{A \cap B}{A \cup B}$
 A – ground truth
 B – detector result
 Correspondence between GT and predicted?

99

Accuracy evaluation



Sorted array of $IoU > T$ between GT and predicted (max length = $n \cdot m$)

$(iou_1, iou_2, iou_3, \dots)$

GT 2 1 1
pred 3 1 2

if appear firstly

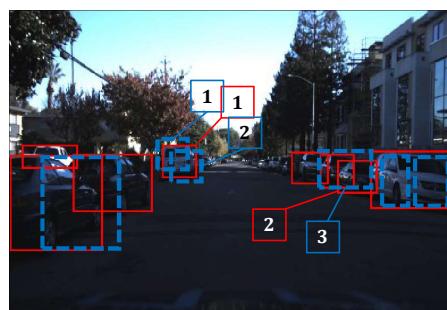
matched GT 2 1

matched pred 3 1

evaluator.py ->
get_single_image_results(gt_boxes, pred_boxes, iou_thr)

101

Accuracy evaluation



Sorted array of $IoU > T$ between GT and predicted (max length = $n \cdot m$)

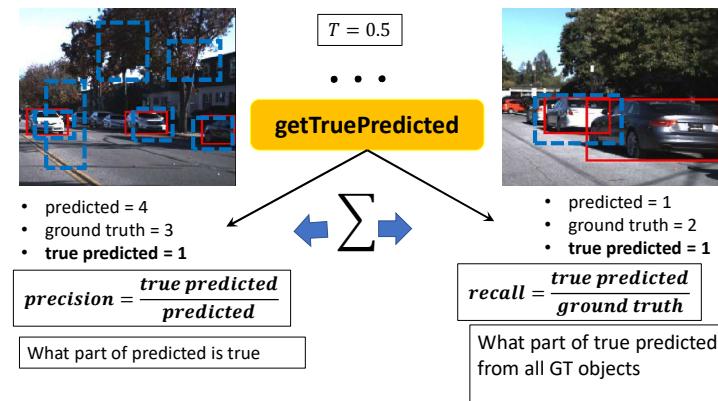
True predicted

$\begin{matrix} \text{matched GT} & 2 & 1 \\ \text{matched pred} & 3 & 1 \end{matrix}$

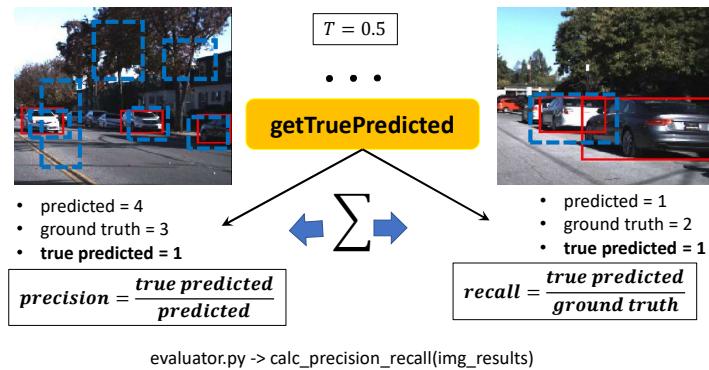
evaluator.py ->
get_single_image_results(gt_boxes, pred_boxes, iou_thr)

102

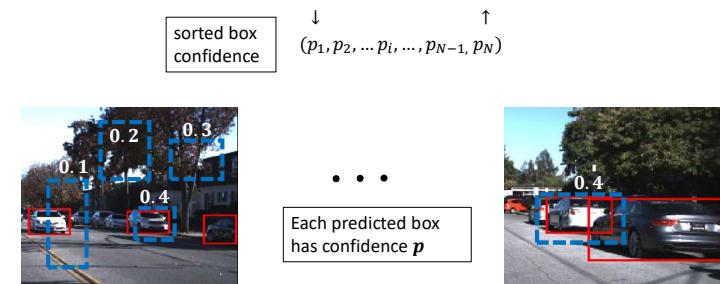
Accuracy evaluation: precision, recall



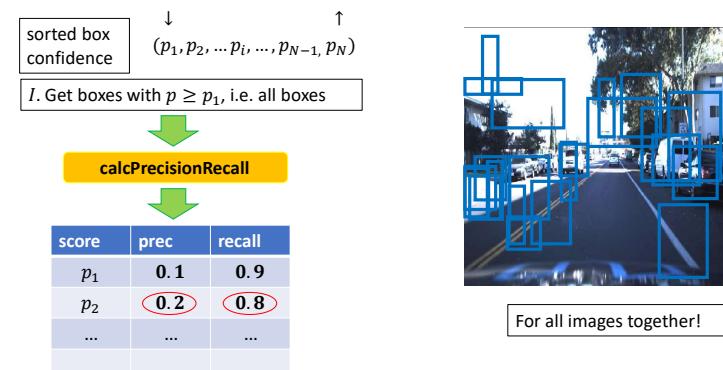
Accuracy evaluation: precision, recall



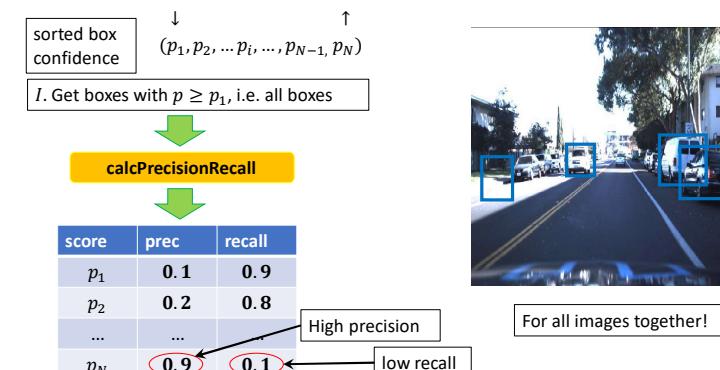
Accuracy evaluation: precision, recall



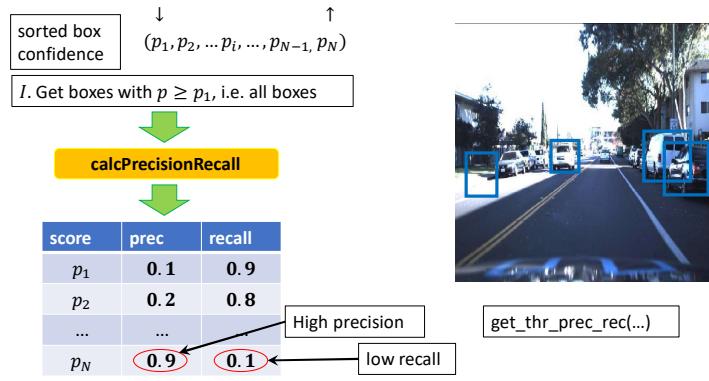
Accuracy evaluation: precision, recall



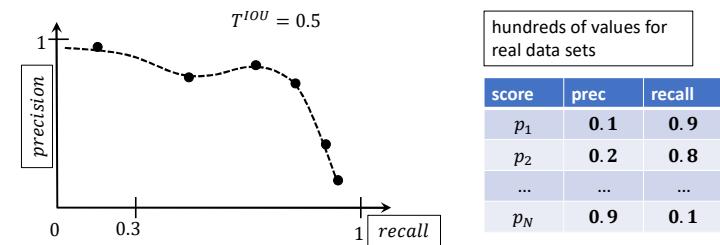
Accuracy evaluation: precision, recall



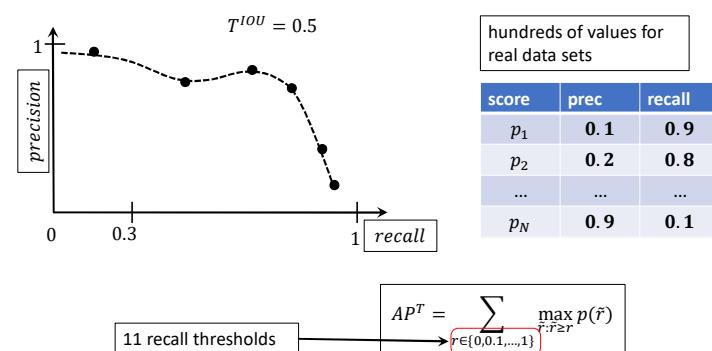
Accuracy evaluation: precision, recall



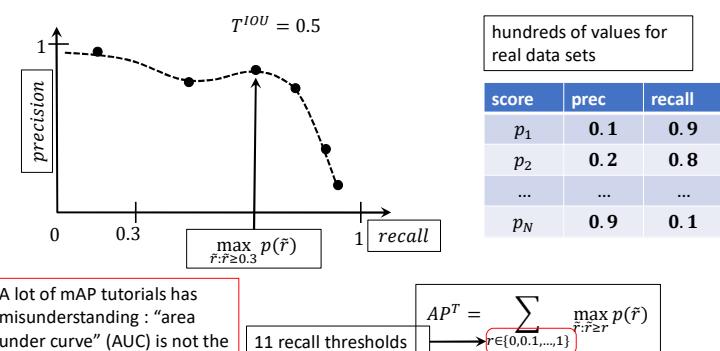
Accuracy evaluation: mAP



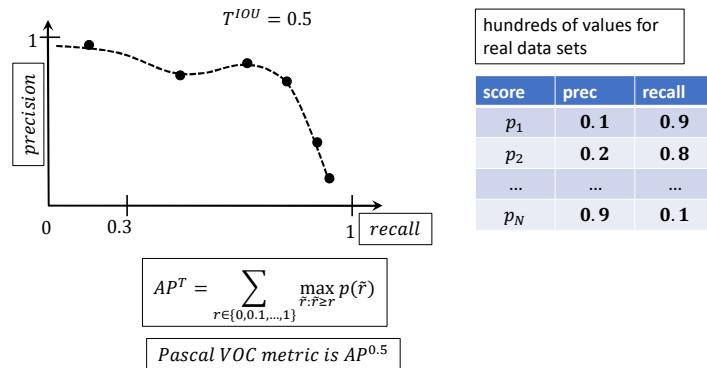
Accuracy evaluation: mAP



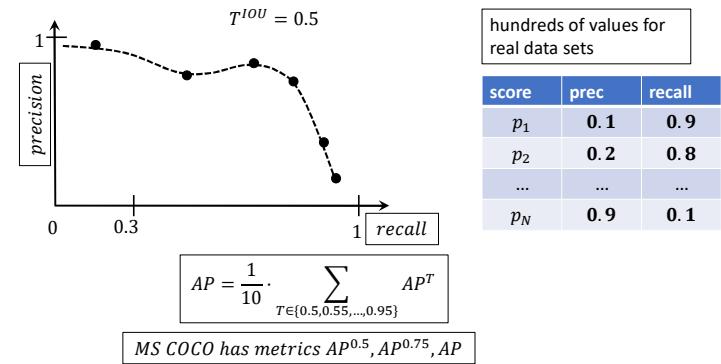
Accuracy evaluation: mAP



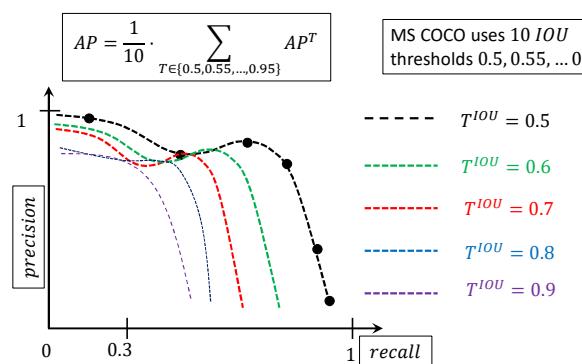
Accuracy evaluation: mAP



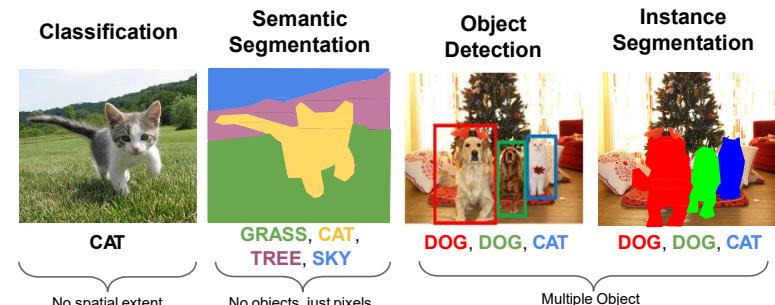
Accuracy evaluation: mAP



Accuracy evaluation: mAP



Computer Vision Tasks



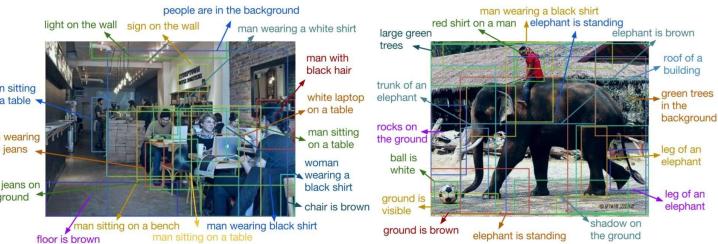
Beyond 2D Object Detection...



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

118

Object Detection + Captioning = Dense Captioning



Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016
Figure copyright IEEE, 2016. Reproduced for educational purposes.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

119

Objects + Relationships = Scene Graphs



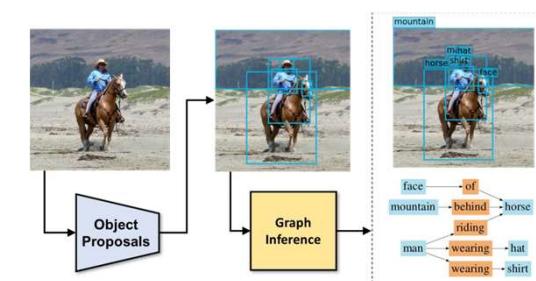
Krishna, Ranjay, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations." International Journal of Computer Vision 123, no. 1 (2017): 32-73.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

121

Scene Graph Prediction



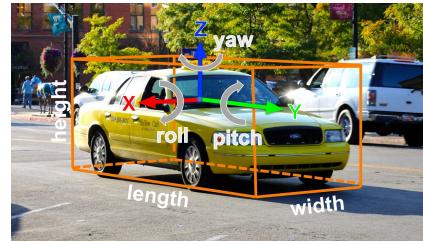
Xu, Zhu, Choy, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", CVPR 2017
Figure copyright IEEE, 2018. Reproduced for educational purposes.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

122

3D Object Detection



2D Object Detection:
2D bounding box
 $[x, y, w, h]$

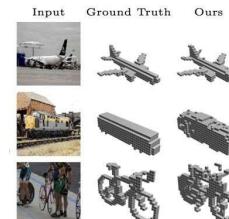
3D Object Detection:
3D oriented bounding box
 $[x, y, z, w, h, l, r, p, y]$

Simplified bbox: no roll & pitch

Much harder problem than 2D object detection!

3D Shape Prediction

Voxel:
 $D \times D \times D$ binary



Choy et al., "3D-R2N2: A Unified Approach for Single and Multi-view 3D Object Reconstruction", ECCV 2016

Pointcloud:
 $V \times 3$ float



Fan et al., "A Point Set Generation Network for 3D Object Reconstruction from a Single Image", CVPR 2017

Mesh:
 $V \times 3$ float, $F \times 3$ int



Wang et al., "Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images", ECCV 2018



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

123



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

126

Open Source Frameworks

- Lots of good implementations on GitHub!
- TensorFlow Detection API:
 - https://github.com/tensorflow/models/tree/master/research/object_detection
 - Faster RCNN, SSD, RFCN, Mask R-CNN
- Caffe2 Detectron:
 - <https://github.com/facebookresearch/Detectron>
 - Mask R-CNN, RetinaNet, Faster R-CNN, RPN, Fast R-CNN, R-FCN
- Finetune on your own dataset with pre-trained models



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

127

References

1. CS231n: Convolutional Neural Networks for Visual Recognition
 - <http://cs231n.stanford.edu/>
2. Object Detection Creation from Scratch. Samsung R&D Institute Ukraine. Vitaliy Bulygin
 <https://aiukraine.com/wp-content/uploads/2018/08/Vitalij-Bulygin-.pptx>
3. CornerNet: Detecting Objects as. Paired Keypoints
 <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

128

