



Computer Vision

Ch (part 3): Deep Learning for CV

Nguyễn Thị Oanh
oanhnt@soict.hust.edu.vn

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

Classification



CAT

Semantic Segmentation



GRASS, CAT,
TREE, SKY

Object Detection



DOG, DOG, CAT

Instance Segmentation



Multiple Object



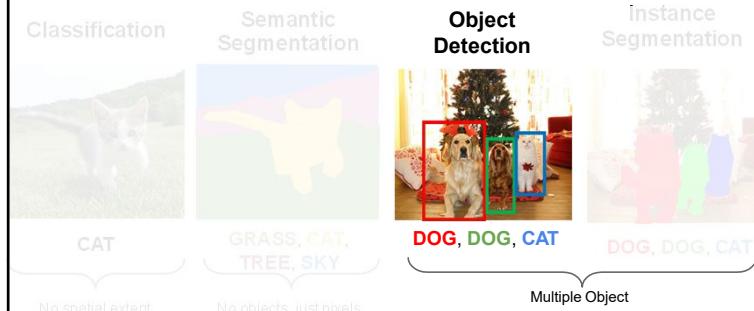
Object detection: sliding windows



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

5

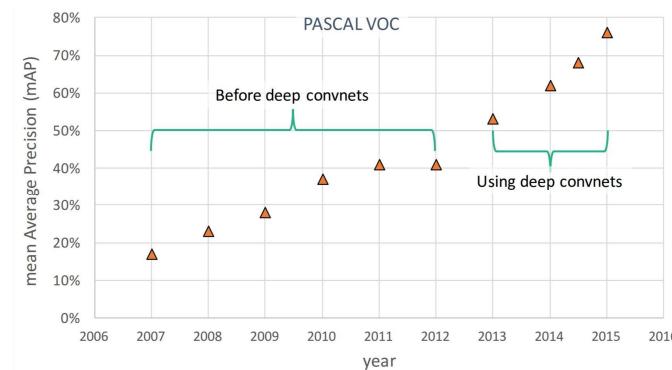
Object Detection



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

6

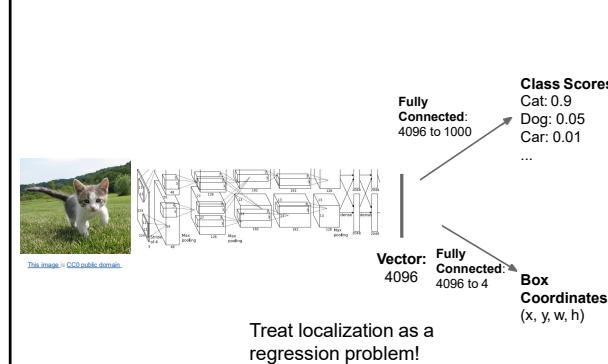
Object Detection: Impact of Deep Learning



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

7

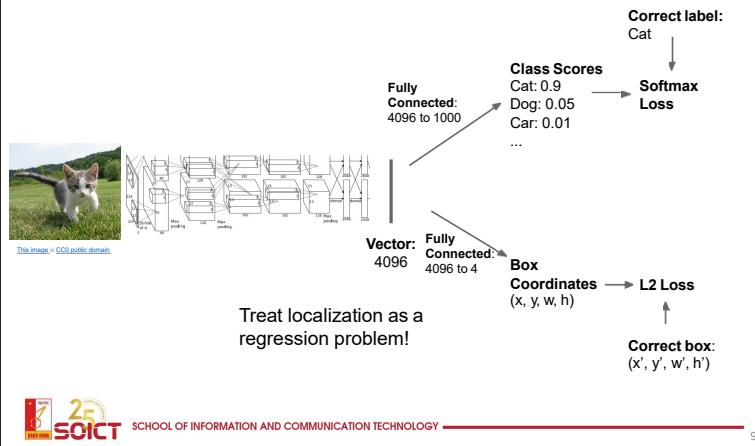
Object Detection: Single Object (Classification + Localization)



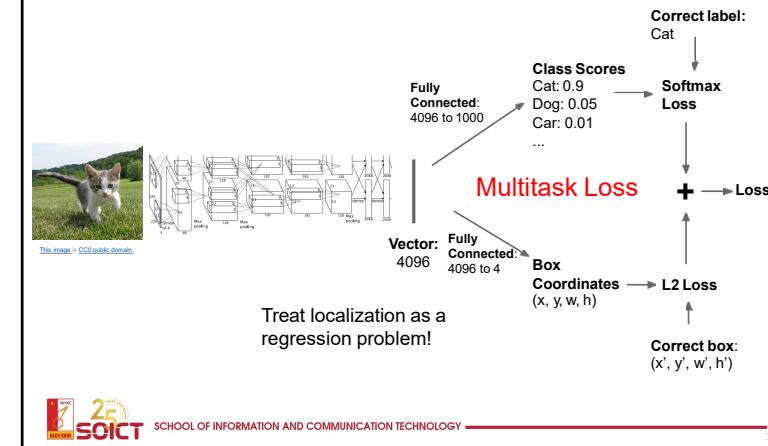
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

8

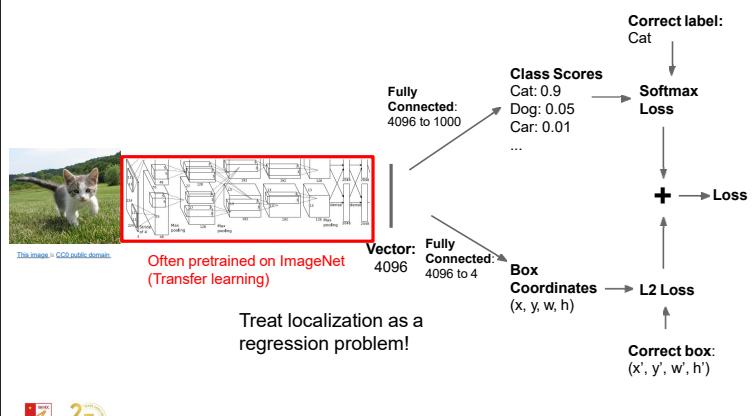
Object Detection: Single Object (Classification + Localization)



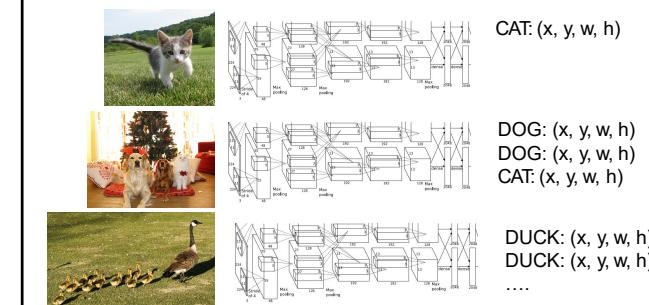
Object Detection: Single Object (Classification + Localization)



Object Detection: Single Object (Classification + Localization)

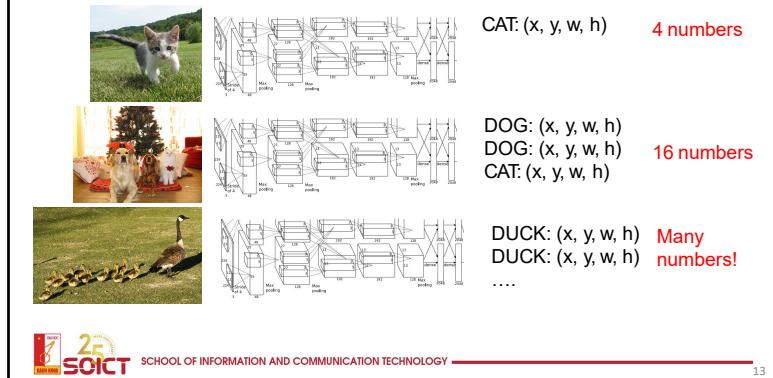


Object Detection: Multiple Objects



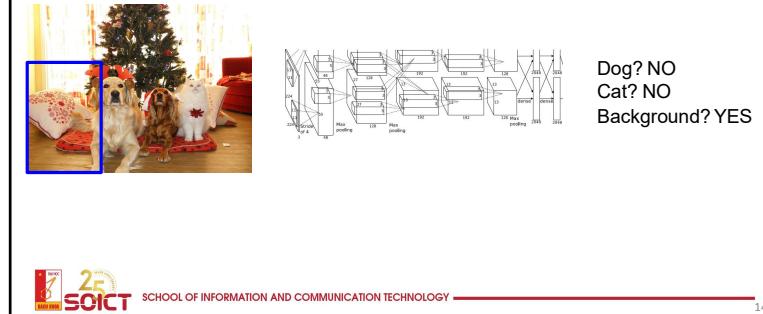
Object Detection: Multiple Objects

Each image needs a different number of outputs!



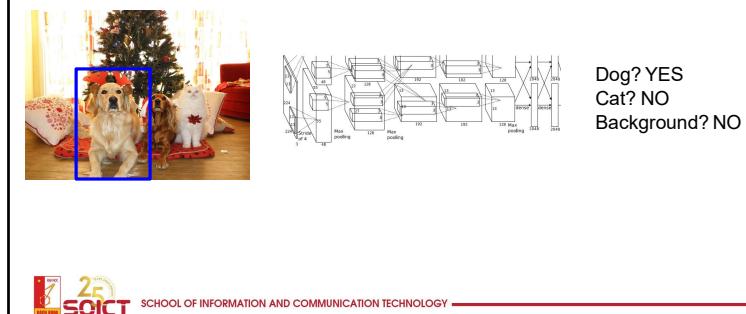
Object Detection: Multiple Objects

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



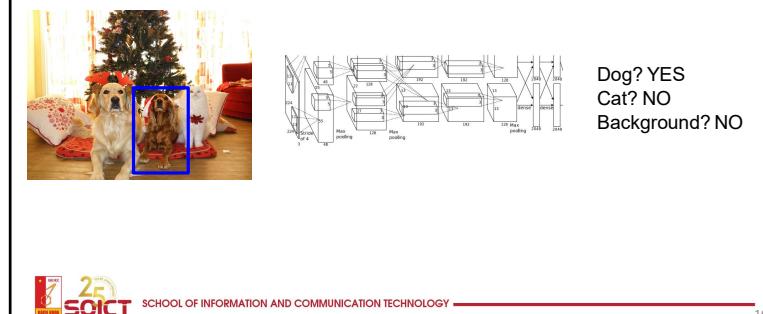
Object Detection: Multiple Objects

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



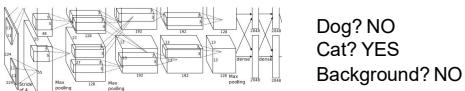
Object Detection: Multiple Objects

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

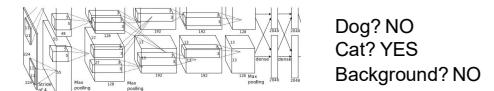
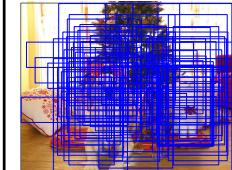


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

17

Object Detection: Multiple Objects

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



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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

18

Two-stage Object detection



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

19

Object Detection

Two Stages

- Propose “objects”
- Classify each candidate

One-Stage

- Sliding window to classify all candidates



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

20

Object Detection

Two Stages

- Propose “objects”
- Classify each candidate

One-Stage

- Sliding window to classify all candidates

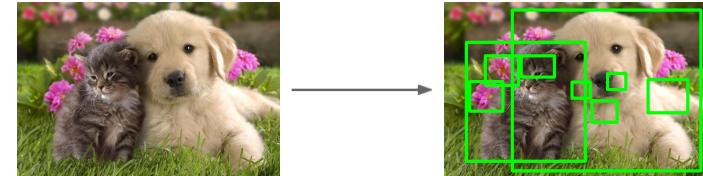


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

21

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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

22

R-CNN



Input image

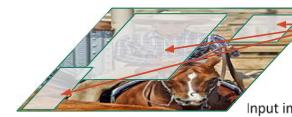
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

23

R-CNN



Regions of Interest (RoI) from a proposal method (~2k)

Input image

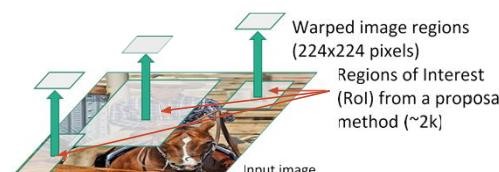
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

24

R-CNN



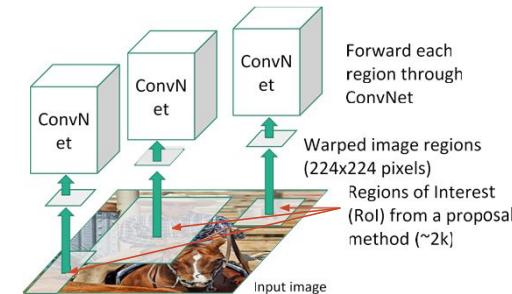
Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

25



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

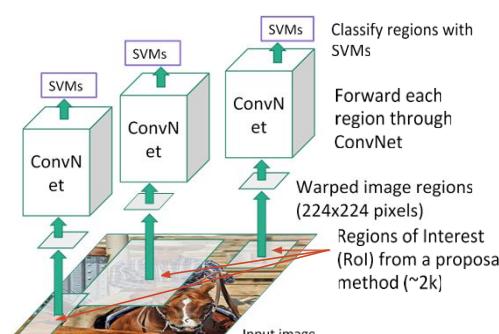
R-CNN



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

26

R-CNN



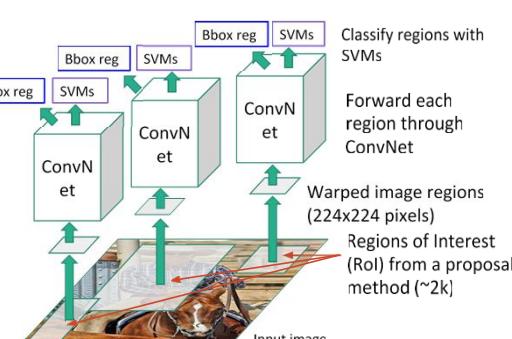
Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

27



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

R-CNN



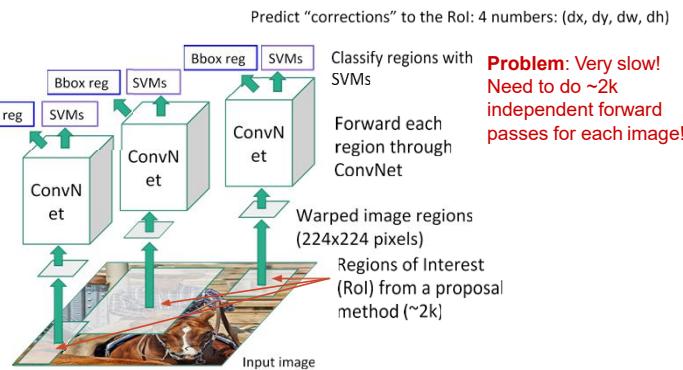
Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.

28



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

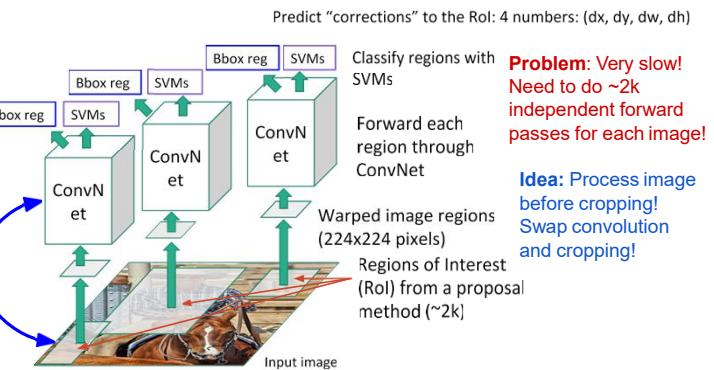
R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

29

"Slow" R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

30

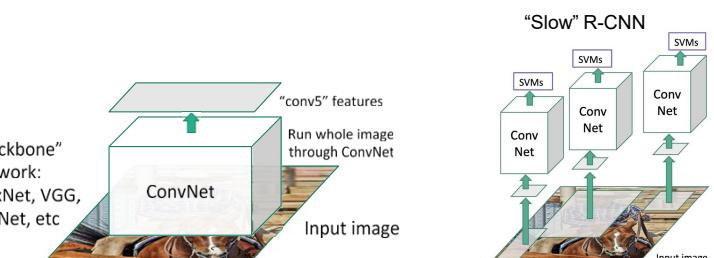
Fast R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

31

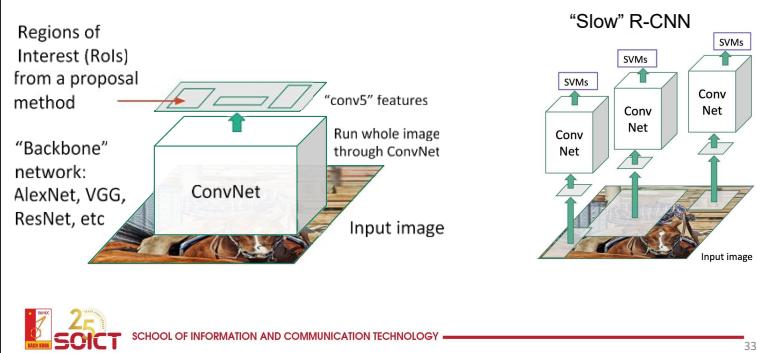
Fast R-CNN



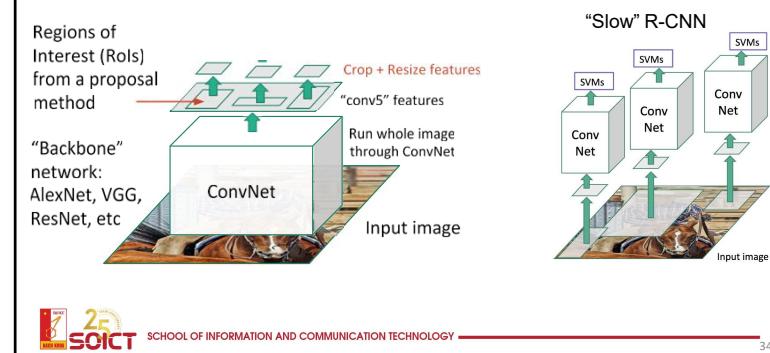
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

32

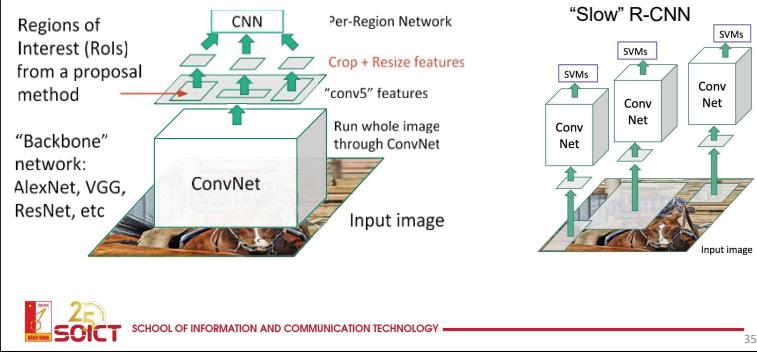
Fast R-CNN



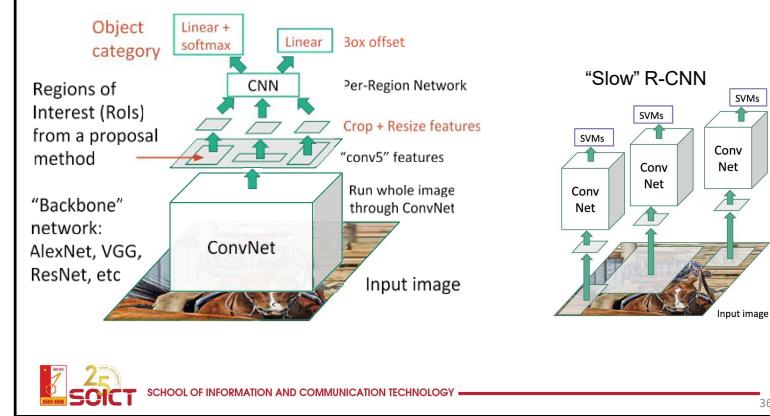
Fast R-CNN



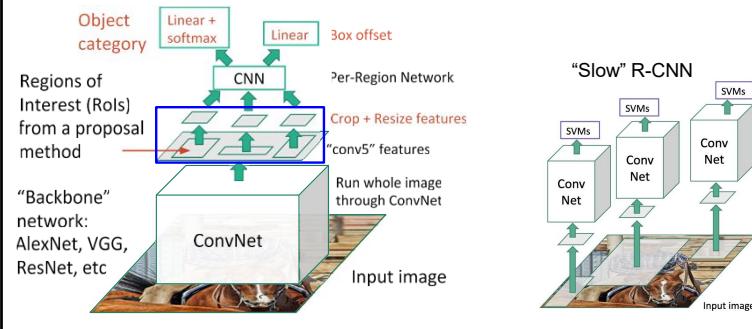
Fast R-CNN



Fast R-CNN



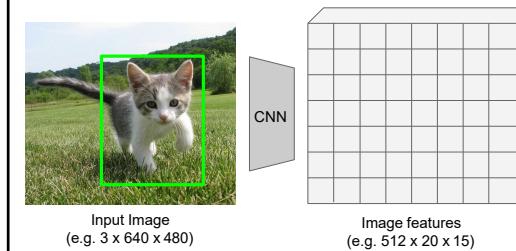
Fast R-CNN



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

37

Cropping Features: RoI Pool



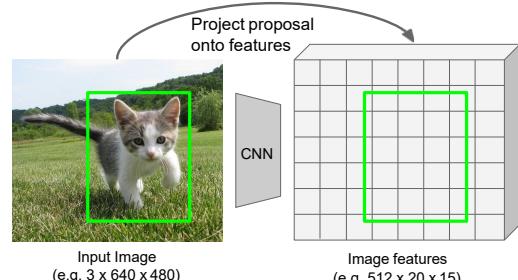
Girshick, “Fast R-CNN”, ICCV 2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

38

Cropping Features: RoI Pool



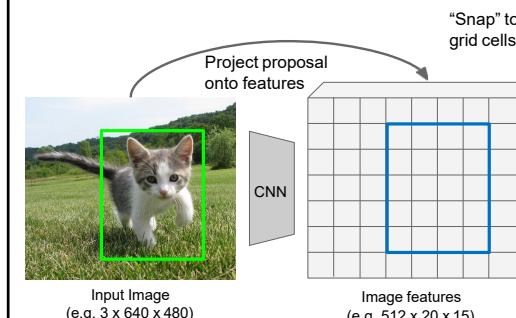
Girshick, “Fast R-CNN”, ICCV 2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

39

Cropping Features: RoI Pool



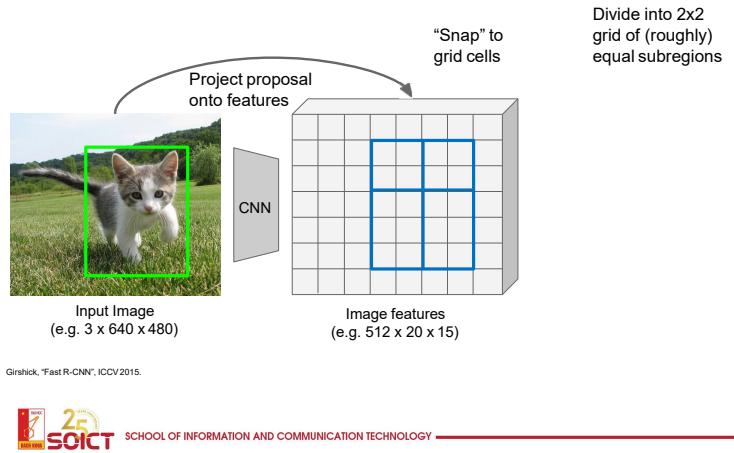
Girshick, “Fast R-CNN”, ICCV 2015.



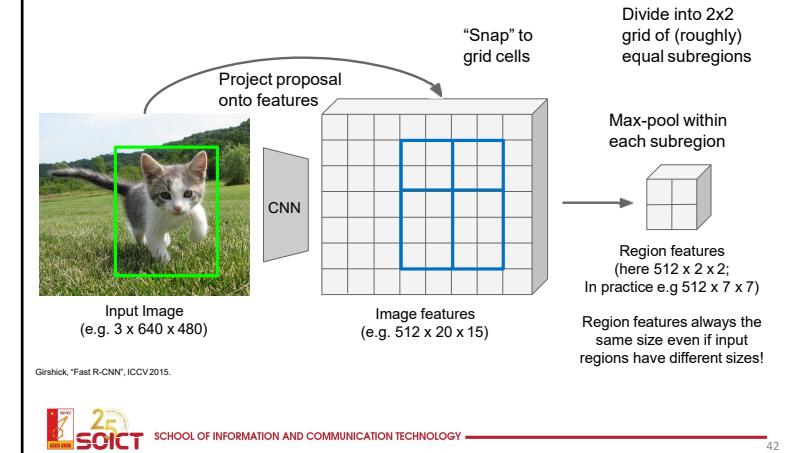
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

40

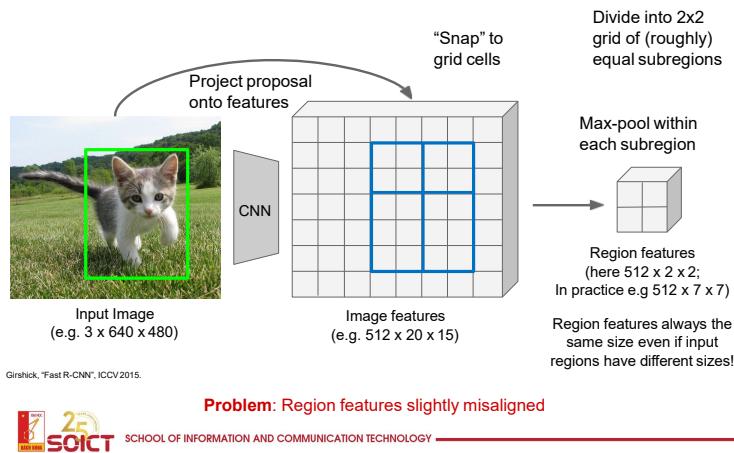
Cropping Features: RoI Pool



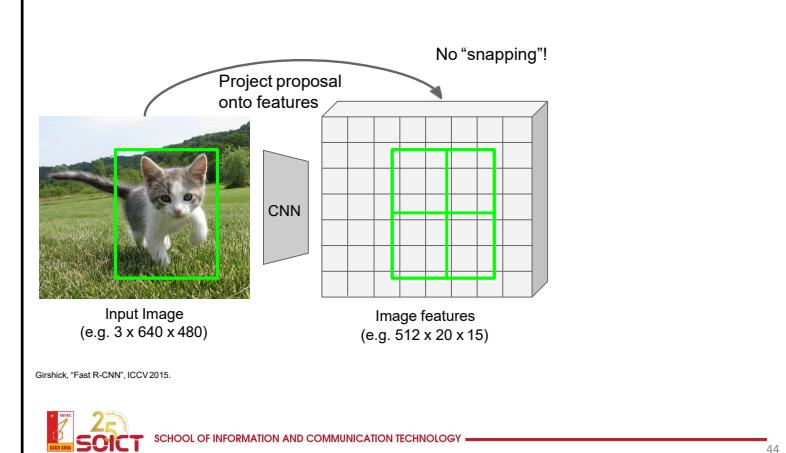
Cropping Features: RoI Pool



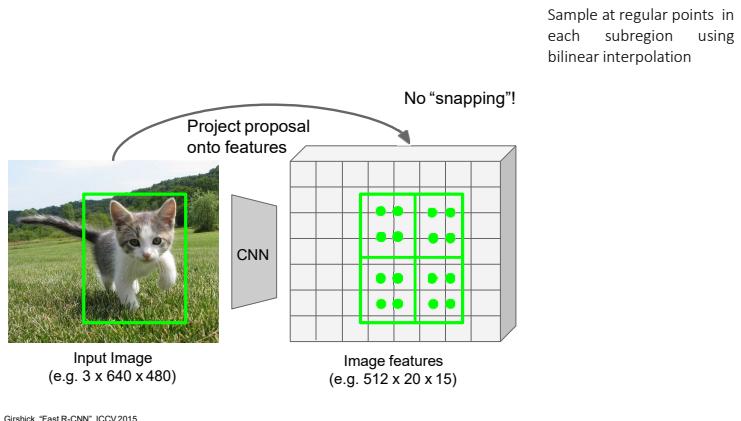
Cropping Features: RoI Pool



Cropping Features: RoI Align

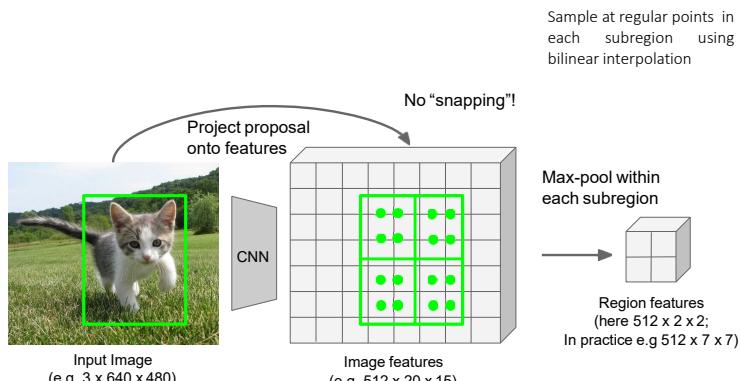


Cropping Features: RoI Align



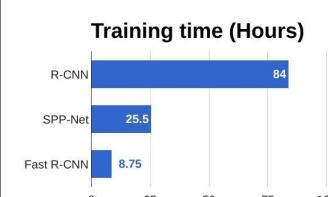
Girshick, "Fast R-CNN", ICCV 2015.

Cropping Features: RoI Align

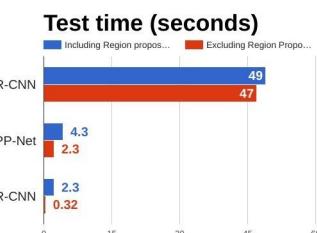


Girshick, "Fast R-CNN", ICCV 2015.

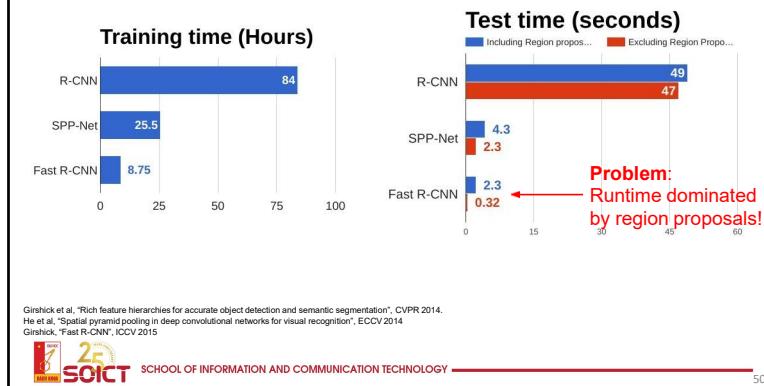
R-CNN vs Fast R-CNN



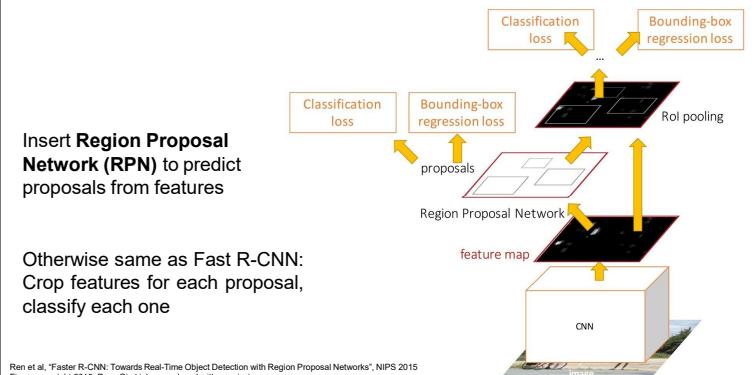
Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
 He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
 Girshick, "Fast R-CNN", ICCV 2015



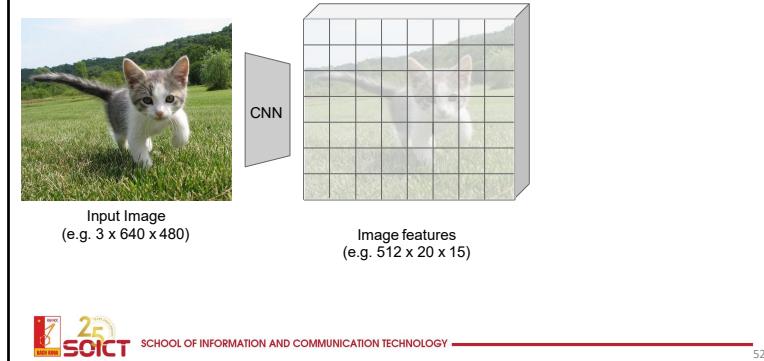
R-CNN vs Fast R-CNN



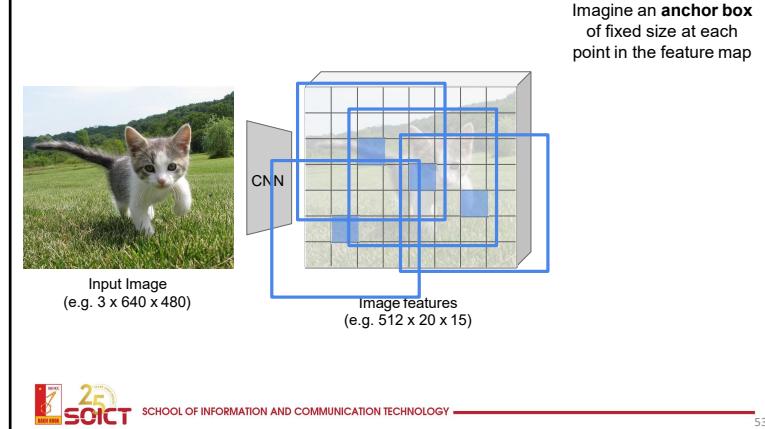
Faster R-CNN: Make CNN do proposals!



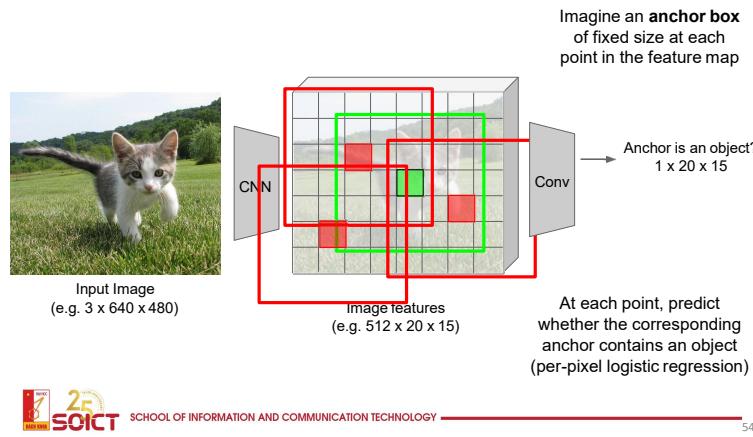
Region Proposal Network



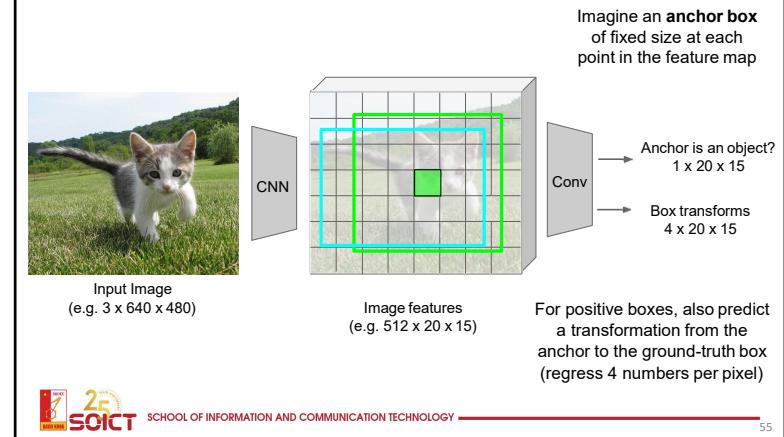
Region Proposal Network



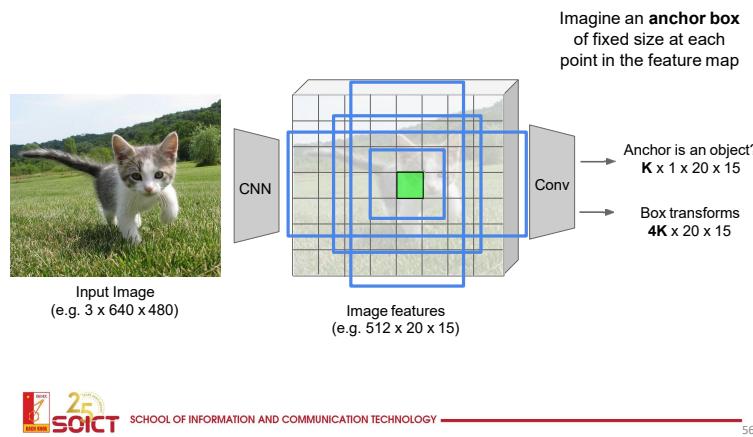
Region Proposal Network



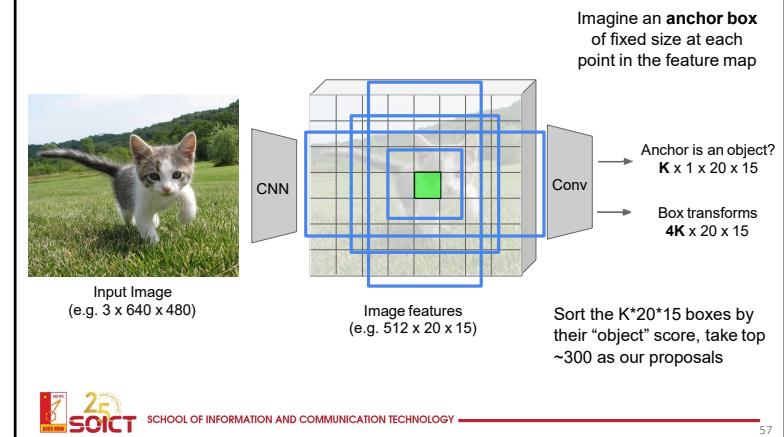
Region Proposal Network



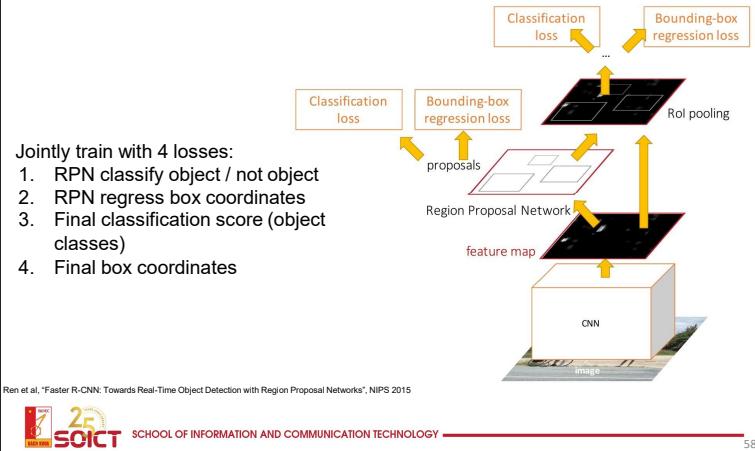
Region Proposal Network



Region Proposal Network



Faster R-CNN: Make CNN do proposals!



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

58

Faster R-CNN: Make CNN do proposals!

R-CNN Test-Time Speed



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

59

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

61

One-stage Object detection Anchor-based



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

63

Object Detection

Two Stages

- Propose “objects”
- Classify each candidate

One-Stage

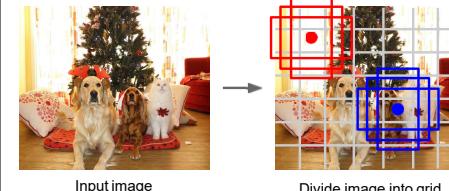
- Sliding window to classify all candidates



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

64

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

- Within each grid cell:
- Regress from each of the B base boxes to a final box with 5 numbers: (dx, dy, dh, dw, confidence)
 - Predict scores for each of C classes (including background as a class)
 - Looks a lot like RPN, but category-specific!

Output:
 $7 \times 7 \times (5 * B + C)$



65

Imbalance

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

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

What happens to CE loss in this case?

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

(Derivative of CE loss : <https://deeplearning�io/softmax-crossentropy#derivative-of-cross-entropy-loss-with-softmax>)

Loss and gradient are dominated by correctly classified negative examples

Training outcome \rightarrow constant “negative” prediction.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

66

Imbalance – Focal Loss

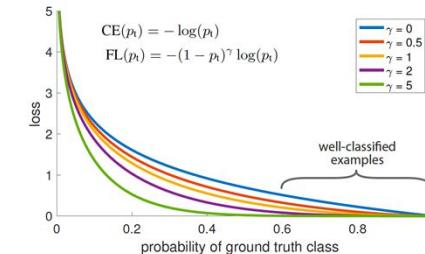


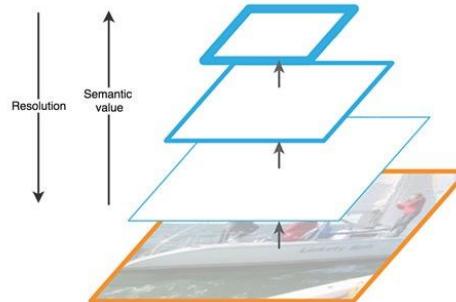
Figure 1: We propose a novel loss we term the *Focal Loss* that adds a factor $(1 - p_i)^\gamma$ to the standard cross entropy criterion. Setting $\gamma > 0$ reduces the relative loss for well-classified examples ($p_i > .5$), putting more focus on hard, misclassified examples. As our experiments will 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)
 SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

67

Feature Pyramid Network (FPN)

- How to handle multiscale predictions?

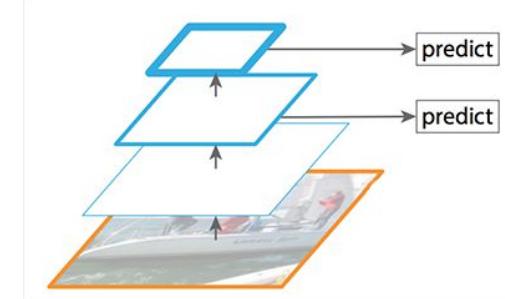


Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. Feature Pyramid Networks for Object Detection (CVPR 2017)

68

Feature Pyramid Network (FPN)

- How to handle multiscale predictions?

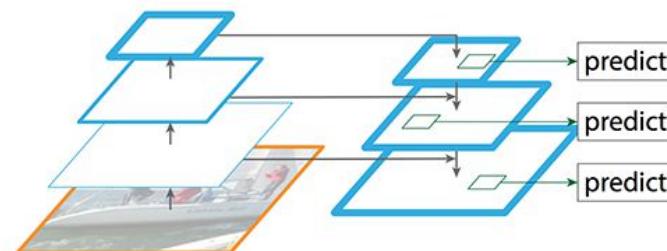


Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. Feature Pyramid Networks for Object Detection (CVPR 2017)

69

Feature Pyramid Network (FPN)

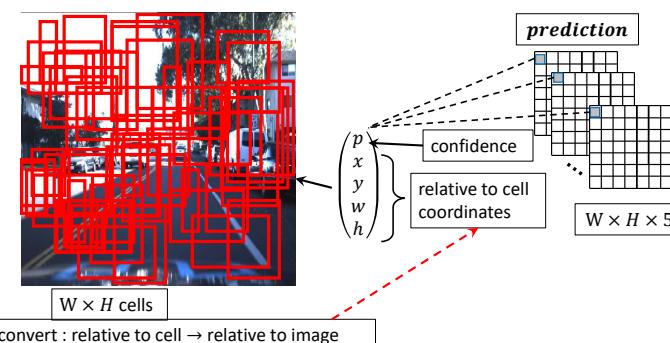
- How to handle multiscale predictions?



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

70

Postprocessing: NMS

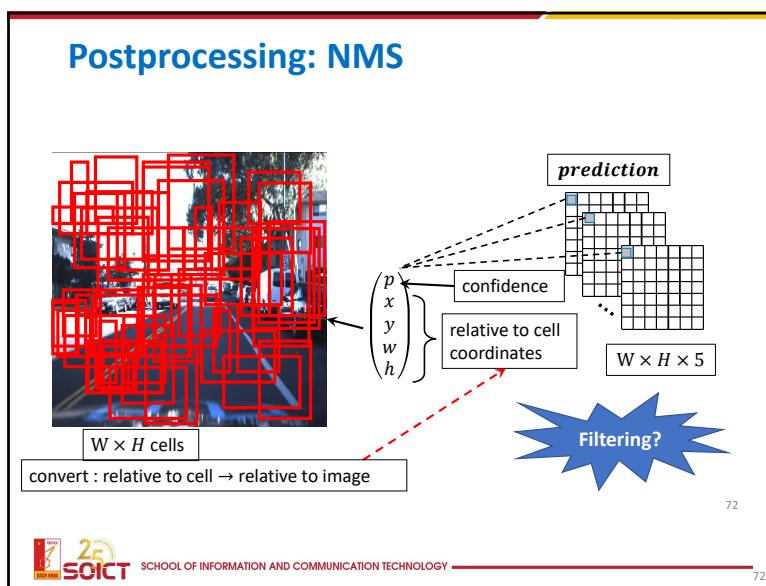


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

71

71

Postprocessing: NMS



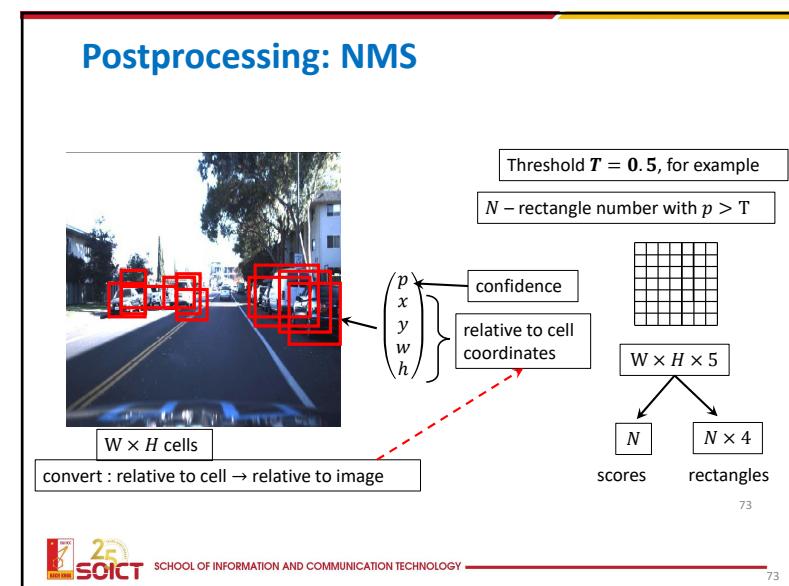
72



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

72

Postprocessing: NMS



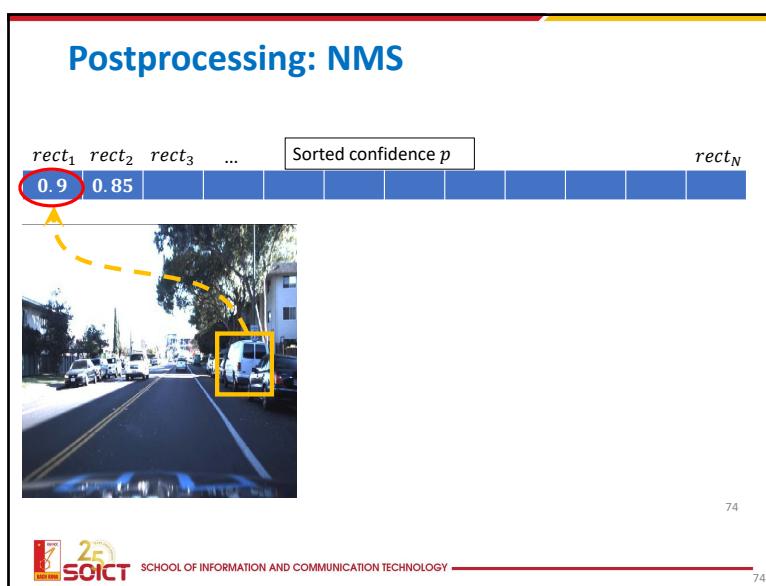
73



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

73

Postprocessing: NMS



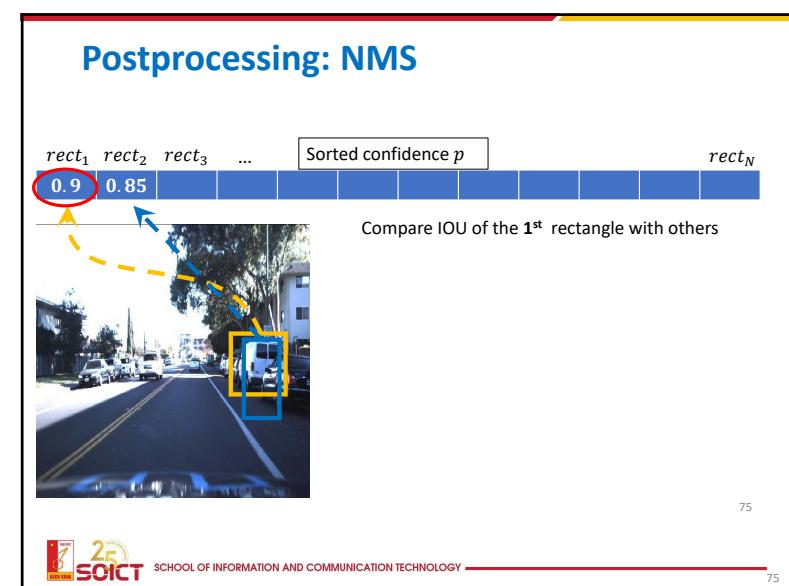
74



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

74

Postprocessing: NMS



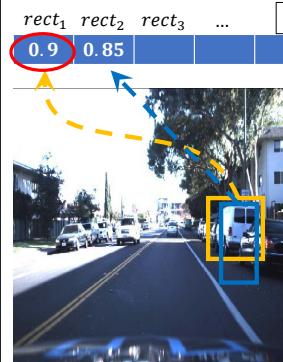
75



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

75

Postprocessing: NMS

Sorted confidence p

$rect_1 \ rect_2 \ rect_3 \ ... \ rect_N$

Compare IOU of the 1st rectangle with others

$$IOU = \frac{\text{Intersection Area}}{\text{Union Area}} > 0.5$$

76



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

76

Postprocessing: NMS

$rect_1 \ rect_2 \ rect_3 \ ... \ rect_N$

Sorted confidence p

~~0.9~~ ~~0.85~~ ~~0.82~~ ... $rect_N$

Compare IOU of the 1st rectangle with others

$$IOU = \frac{\text{Intersection Area}}{\text{Union Area}} > 0.5$$

77



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

77

Postprocessing: NMS

Sorted confidence p

$rect_1 \ rect_2 \ rect_3 \ ... \ rect_N$

Compare IOU of the 1st rectangle with others

$$IOU = \frac{\text{Intersection Area}}{\text{Union Area}} = 0$$

do nothing

78



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

78

Postprocessing: NMS

$rect_1 \ rect_2 \ rect_3 \ ... \ rect_N$

Sorted confidence p

~~0.9~~ ~~0.85~~ ~~0.82~~ ~~0.77~~ ... $rect_N$

Compare IOU of the 1st rectangle with others

$$IOU = \frac{\text{Intersection Area}}{\text{Union Area}} = 0$$

IOU = 0 with the chosen rectangle \rightarrow *do nothing!*

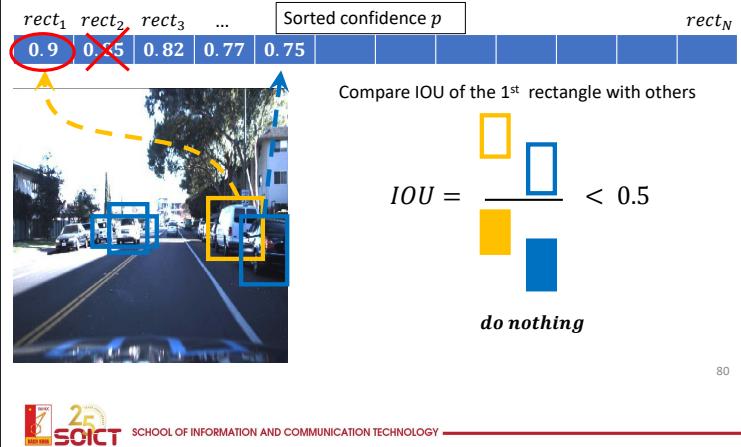
79



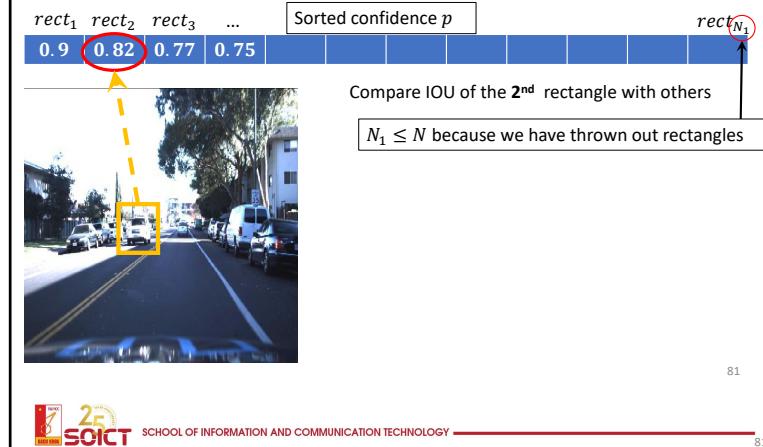
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

79

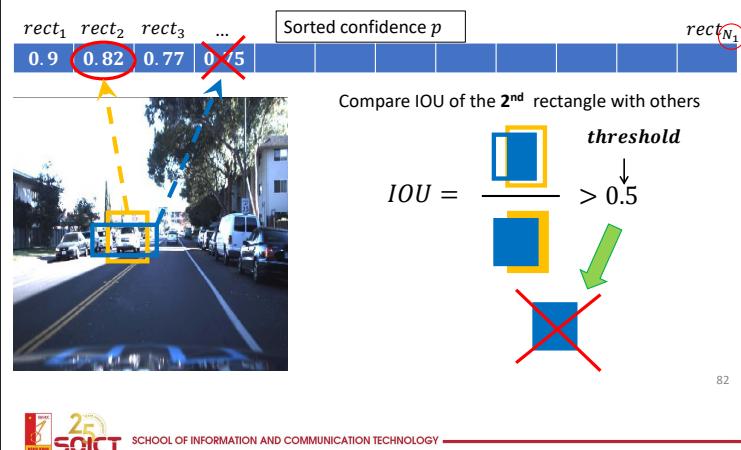
Postprocessing: NMS



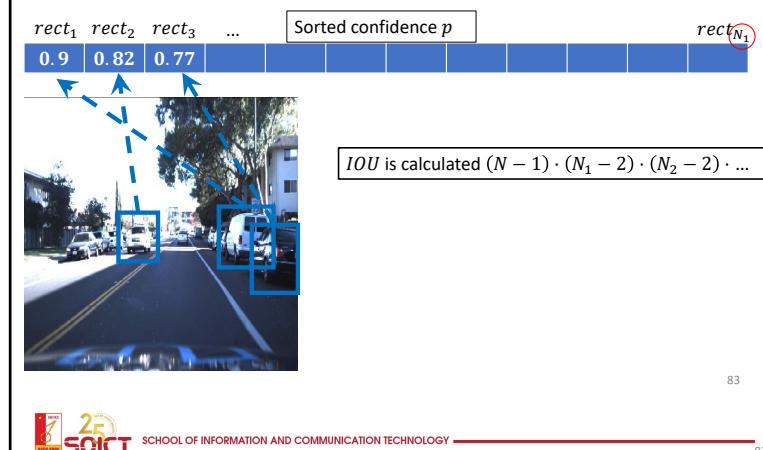
Postprocessing: NMS



Postprocessing: NMS



Postprocessing: NMS



One-stage Object detection

Anchor-free



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

84

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

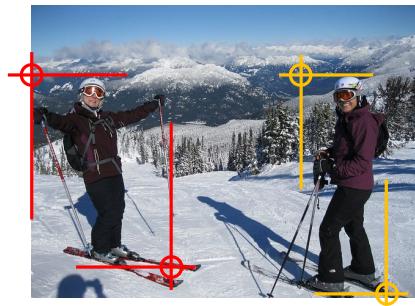
Source: <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

85

CornerNet: Detecting Objects as Paired Keypoints



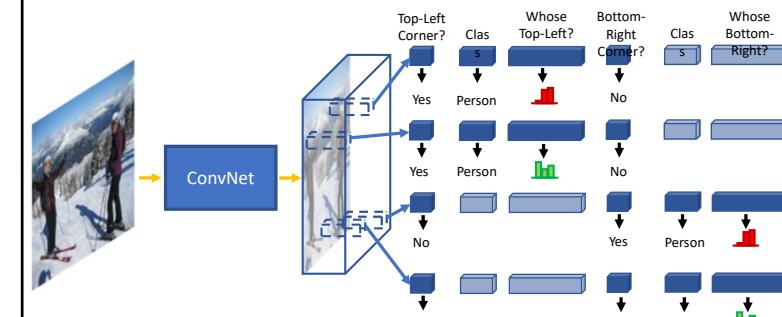
Hei Law, Jia Deng. CornerNet: Detecting Objects as Paired Keypoints – ECCV2018
Source: <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

86

CornerNet: Detecting Objects as Paired Keypoints



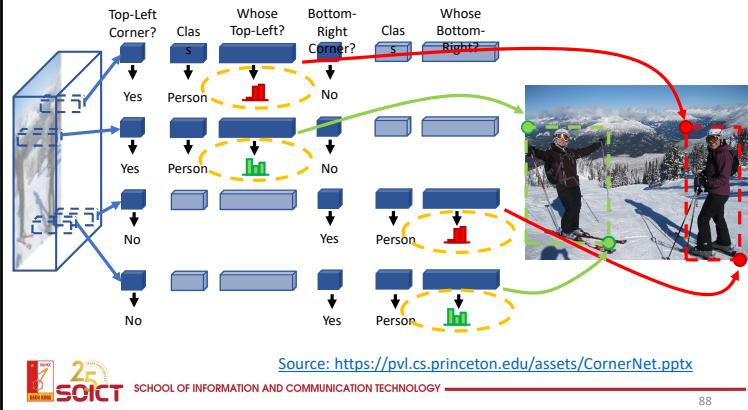
Source: <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>



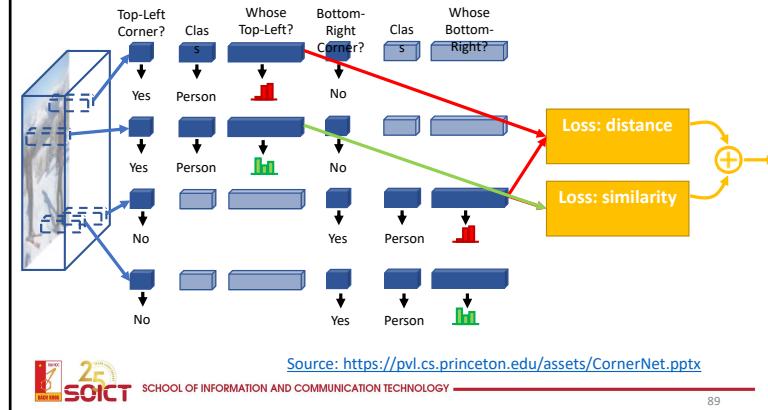
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

87

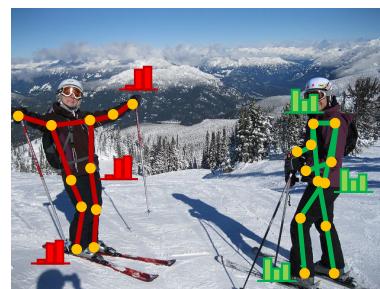
CornerNet: Detecting Objects as Paired Keypoints



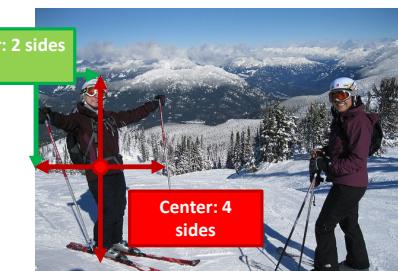
CornerNet: Detecting Objects as Paired Keypoints



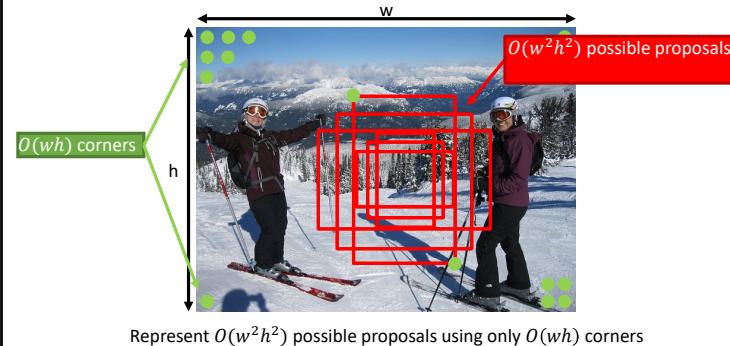
Associative Embedding [Newell et al. NIPS'17]



Advantages of Detecting Corners



Advantages of Detecting Corner

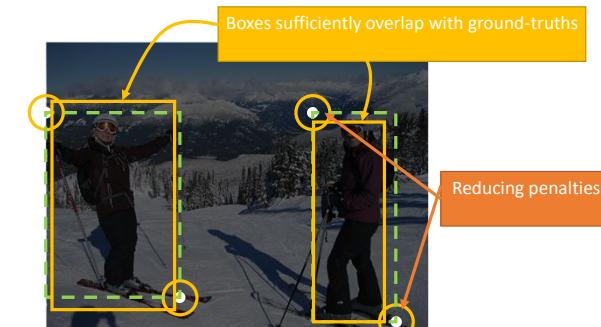


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Source: <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>

92

Supervising Corner Detection



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Source: <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>

93

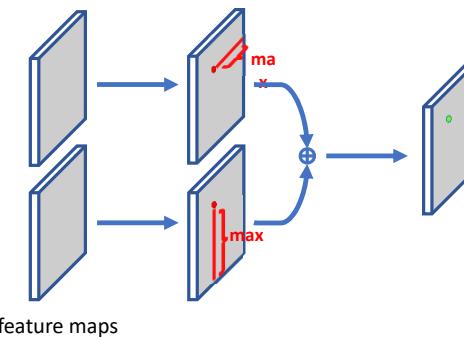
Corner Pooling

Source: <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

94

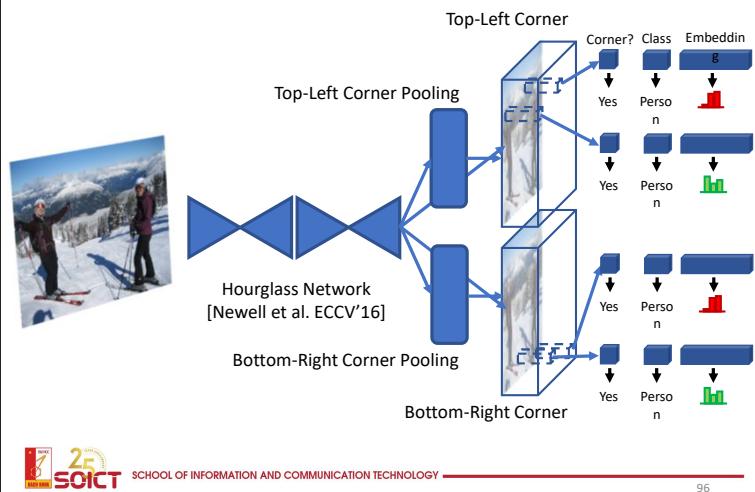
Top-Left Corner Pooling

Source: <https://pvl.cs.princeton.edu/assets/CornerNet.pptx>

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

95

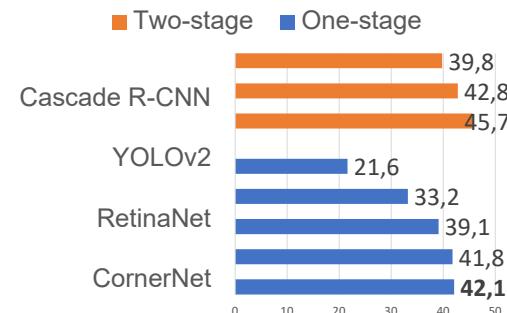
CornerNet



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

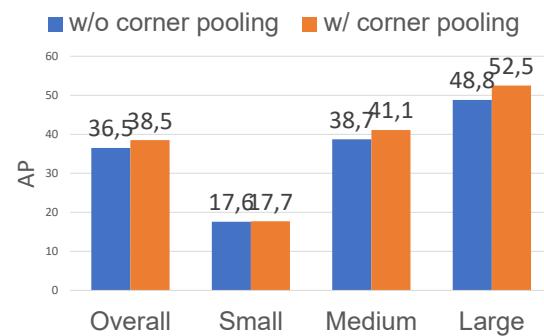
96

Experiment: CornerNet versus Others



97

Experiment: Corner Pooling



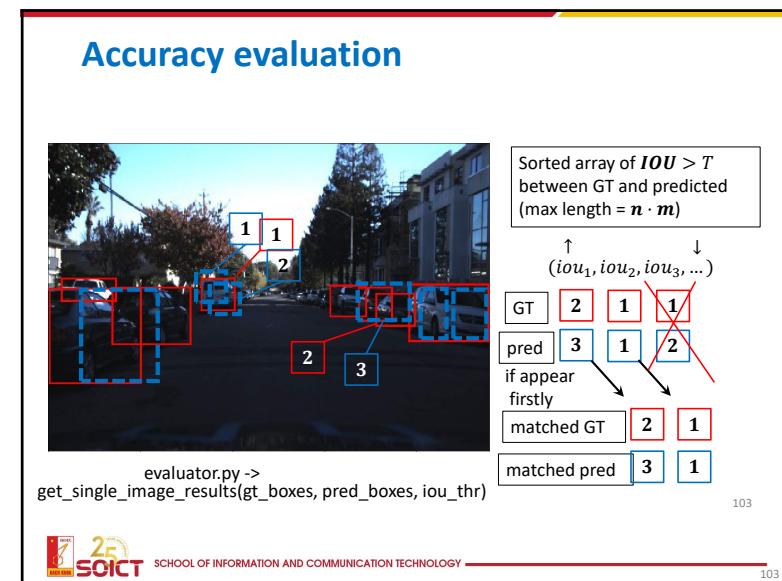
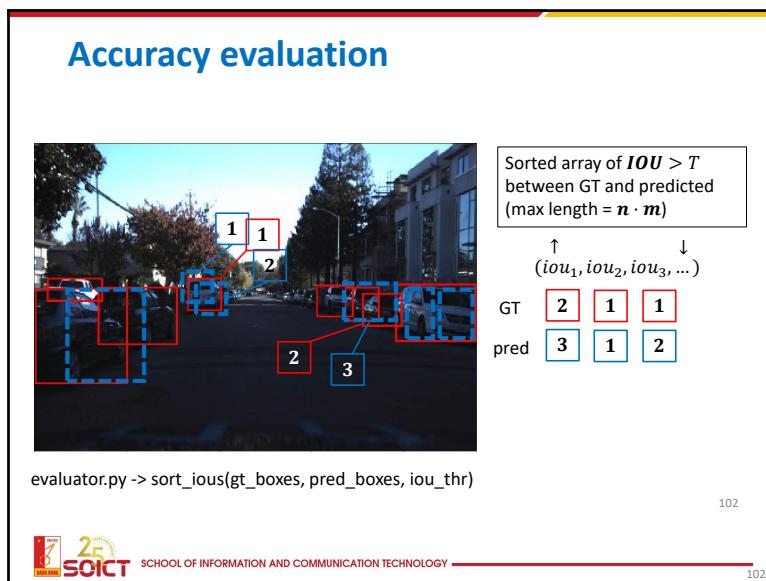
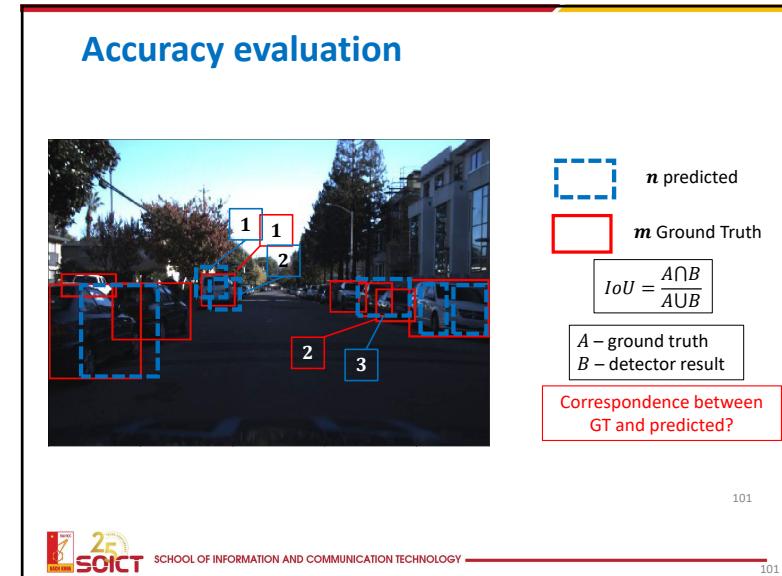
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

98

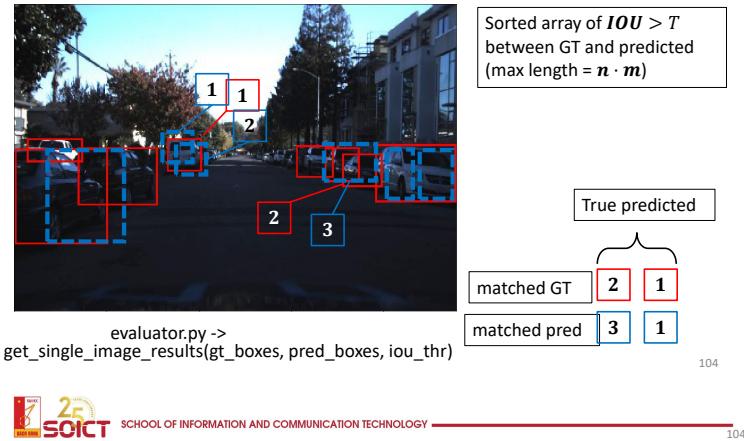


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

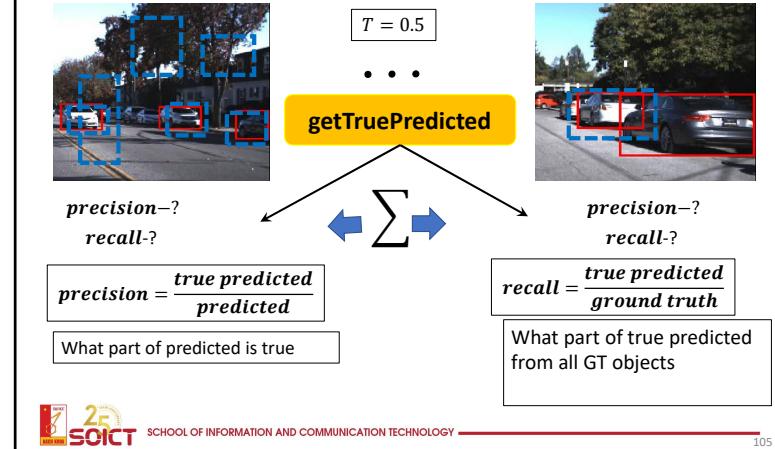
99



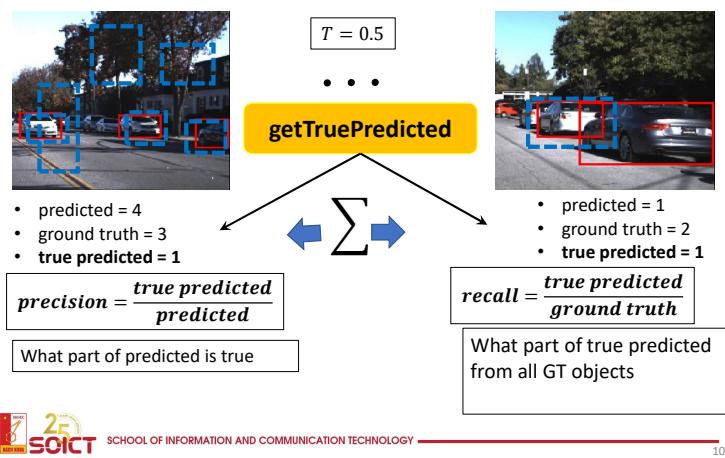
Accuracy evaluation



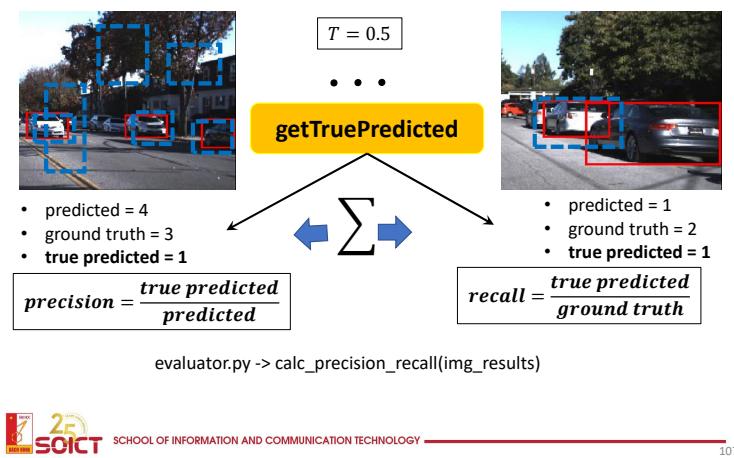
Accuracy evaluation: precision, recall



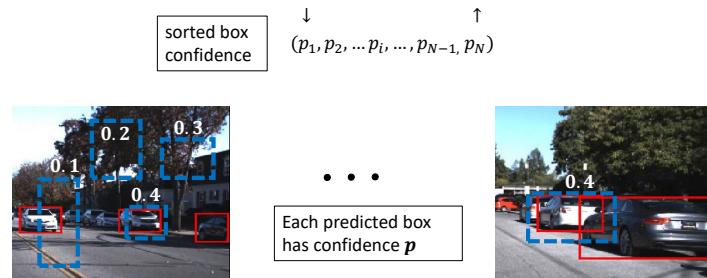
Accuracy evaluation: precision, recall



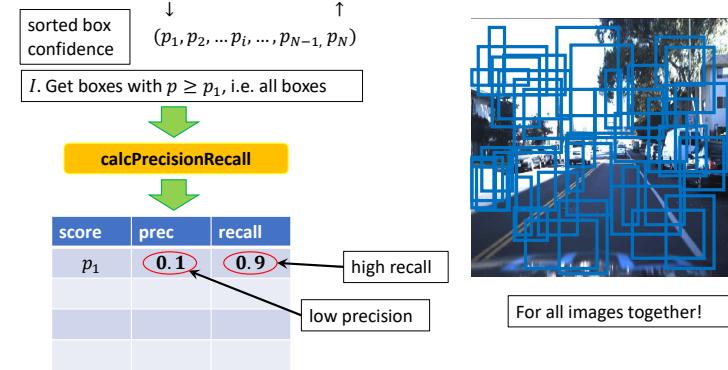
Accuracy evaluation: precision, recall



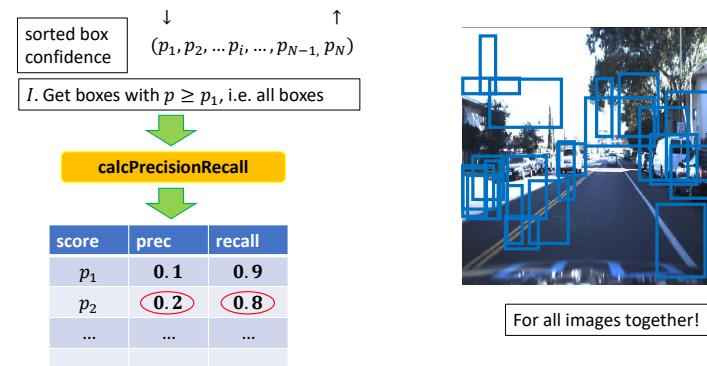
Accuracy evaluation: precision, recall



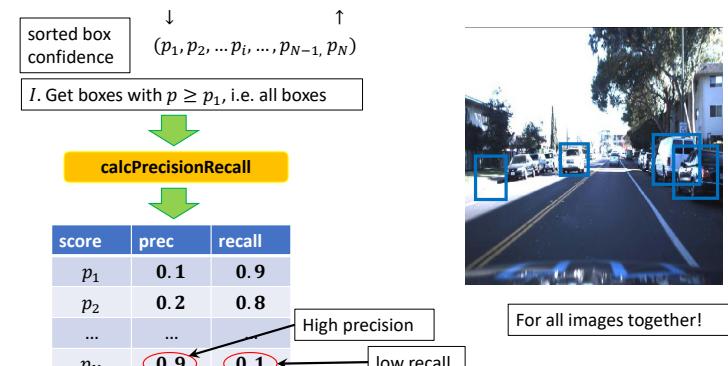
Accuracy evaluation: precision, recall



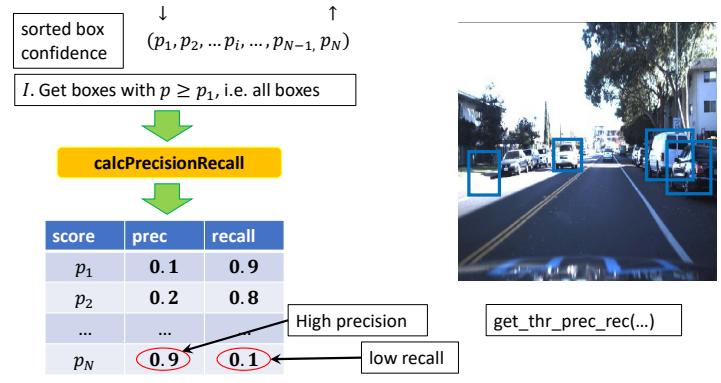
Accuracy evaluation: precision, recall



Accuracy evaluation: precision, recall



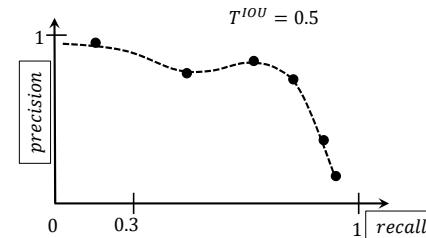
Accuracy evaluation: precision, recall



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

112

Accuracy evaluation: mAP



hundreds of values for real data sets

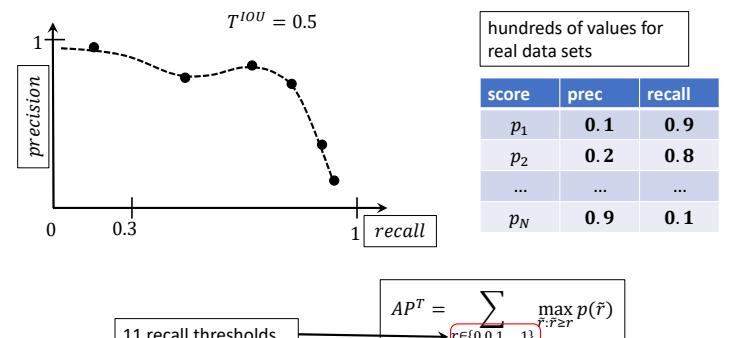
score	prec	recall
p_1	0.1	0.9
p_2	0.2	0.8
...
p_N	0.9	0.1



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

113

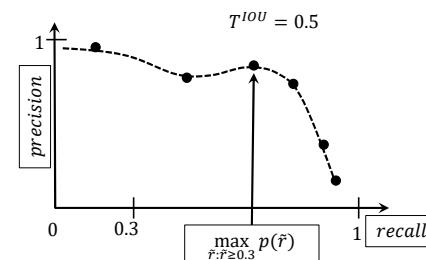
Accuracy evaluation: mAP



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

114

Accuracy evaluation: mAP



hundreds of values for real data sets

score	prec	recall
p_1	0.1	0.9
p_2	0.2	0.8
...
p_N	0.9	0.1

A lot of mAP tutorials has misunderstanding: "area under curve" (AUC) is not the same!

11 recall thresholds

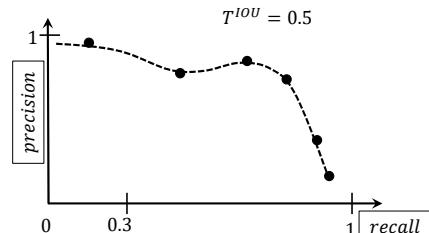
$$AP^T = \sum_{r \in \{0, 0.1, \dots, 1\}} \max_{\tilde{r}: \tilde{r} \geq r} p(\tilde{r})$$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

115

Accuracy evaluation: mAP



hundreds of values for real data sets

score	prec	recall
p_1	0.1	0.9
p_2	0.2	0.8
...
p_N	0.9	0.1

$$AP^T = \sum_{r \in \{0, 0.1, \dots, 1\}} \max_{\tilde{r}: \tilde{r} \geq r} p(\tilde{r})$$

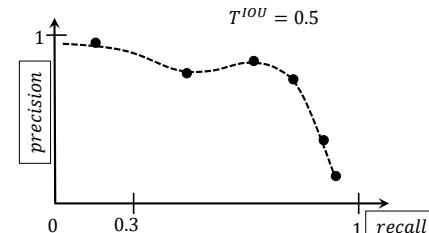
Pascal VOC metric is $AP^{0.5}$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

116

Accuracy evaluation: mAP



hundreds of values for real data sets

score	prec	recall
p_1	0.1	0.9
p_2	0.2	0.8
...
p_N	0.9	0.1

$$AP = \frac{1}{10} \cdot \sum_{T \in \{0.5, 0.55, \dots, 0.95\}} AP^T$$

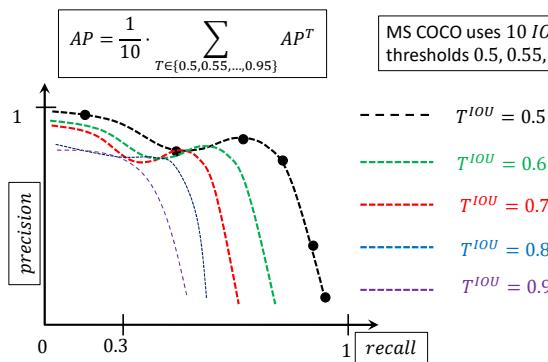
MS COCO has metrics $AP^{0.5}, AP^{0.75}, AP$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

117

Accuracy evaluation: mAP



MS COCO uses 10 IOU thresholds 0.5, 0.55, ... 0.95



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

118

Computer Vision Tasks

Classification



CAT

Semantic Segmentation



GRASS, CAT, TREE, SKY

No spatial extent

Object Detection



DOG, DOG, CAT

No objects, just pixels

Instance Segmentation



Multiple Object



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

119

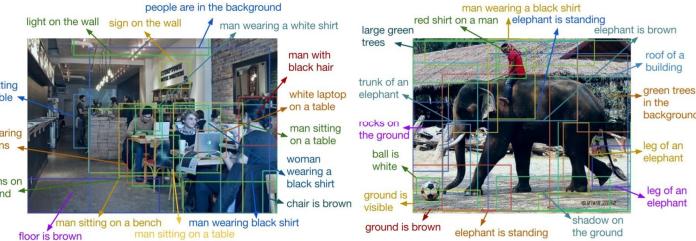
Beyond 2D Object Detection...



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

120

Object Detection + Captioning = Dense Captioning



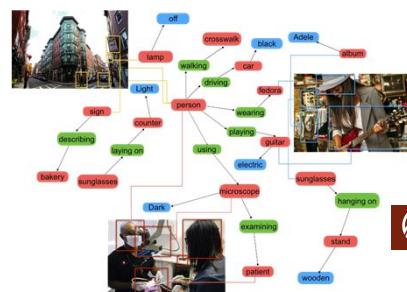
Johnson, Karpathy, and Fei-Fei, "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", CVPR 2016
Figure copyright IEEE, 2016. Reproduced for educational purposes.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

121

Objects + Relationships = Scene Graphs



108,077 Images
5.4 Million Region Descriptions
1.7 Million Visual Question Answers
3.8 Million Object Instances
2.8 Million Attributes
2.3 Million Relationships
Everything Mapped to Wordnet Synsets

VISUALGENOME

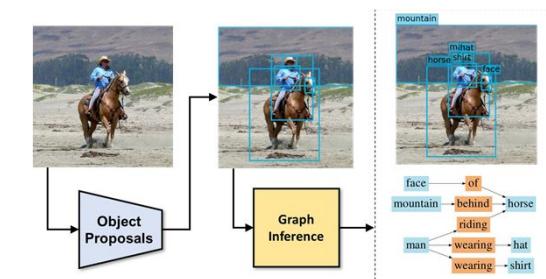
Krishna, Ranjay, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen et al. "Visual genome: Connecting language and vision using crowdsourced dense image annotations." International Journal of Computer Vision 123, no. 1 (2017): 32-73.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

123

Scene Graph Prediction



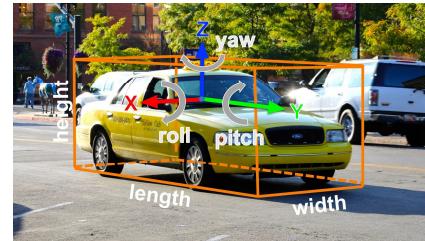
Xu, Zhu, Choy, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", CVPR 2017
Figure copyright IEEE, 2018. Reproduced for educational purposes.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

124

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!

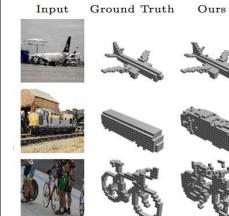


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

125

3D Shape Prediction

Voxel:
D x D x D binary



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

Pointcloud:
V x 3 float



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

Mesh:
V x 3 float, F x 3 int



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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

128

Semantic segmentation

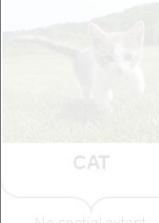


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

129

Semantic Segmentation

Classification



CAT

No spatial extent

**Semantic
Segmentation**



GRASS, CAT, TREE, SKY

No objects, just pixels

Object
Detection



DOG, DOG, CAT

Instance
Segmentation



DOG, DOG, CAT

Multiple Object



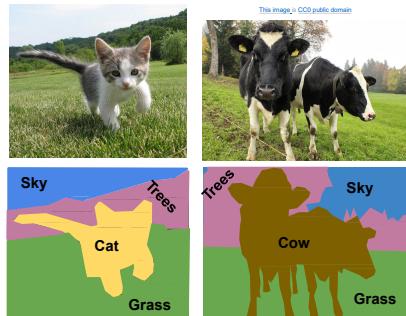
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

130

Semantic Segmentation

Label each pixel in the image with a category label

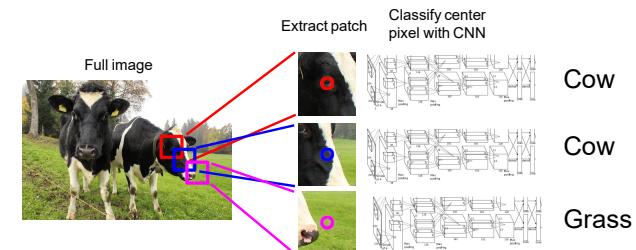
Don't differentiate instances, only care about pixels



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

131

Semantic Segmentation Idea: Sliding Window



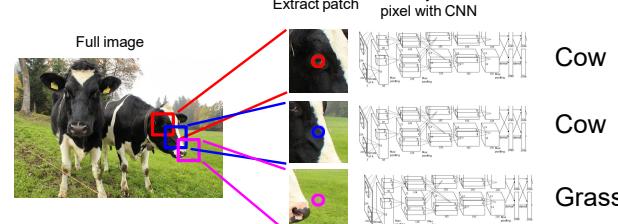
Farabet et al., "Learning Hierarchical Features for Scene Labeling", TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

132

Semantic Segmentation Idea: Sliding Window



Problem: Very inefficient! Not reusing shared features between overlapping patches

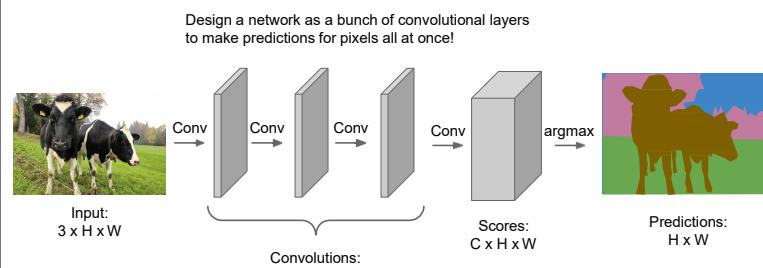
Farabet et al., "Learning Hierarchical Features for Scene Labeling", TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

133

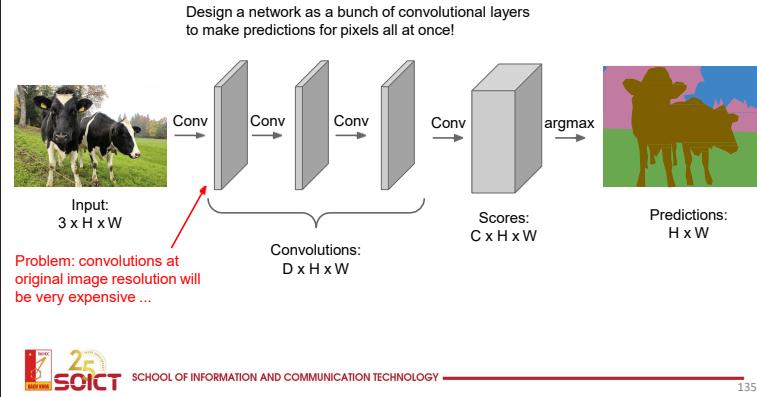
Semantic Segmentation Idea: Fully Convolutional



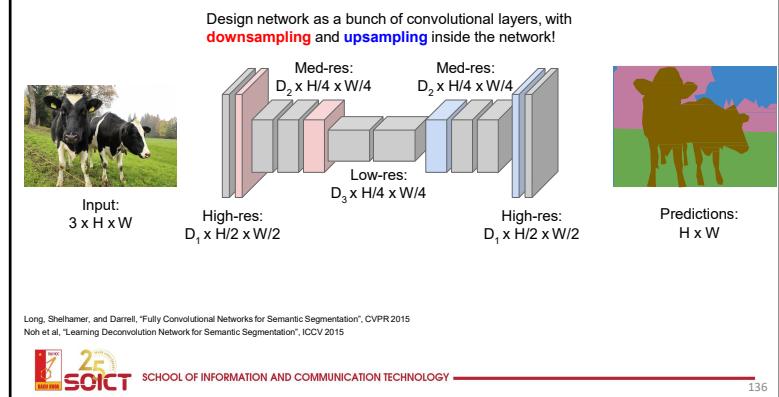
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

134

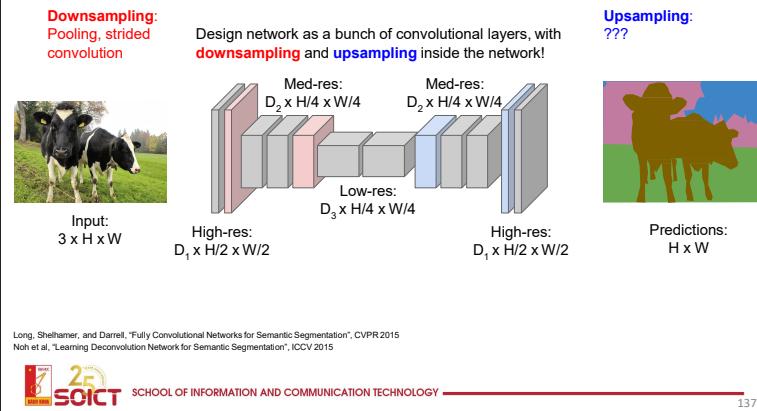
Semantic Segmentation Idea: Fully Convolutional



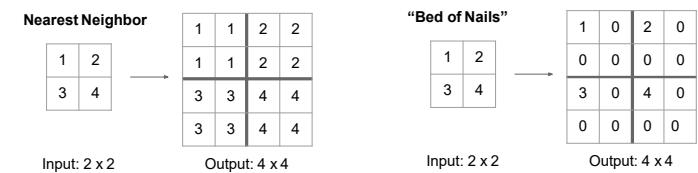
Semantic Segmentation Idea: Fully Convolutional



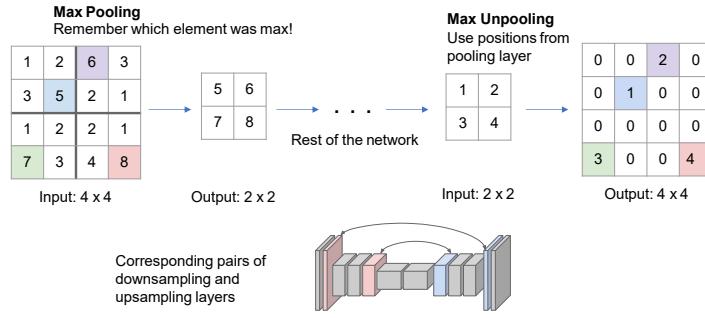
Semantic Segmentation Idea: Fully Convolutional



In-Network upsampling: "Unpooling"



In-Network upsampling: “Max Unpooling”

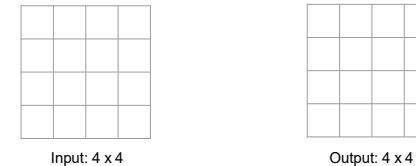


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

139

Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1

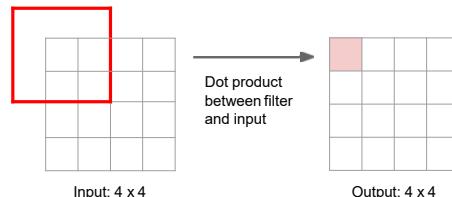


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

140

Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1

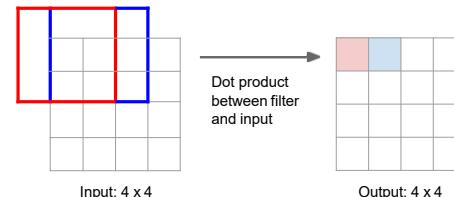


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

141

Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1

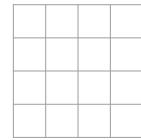


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

142

Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1

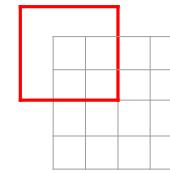


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

143

Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Dot product
between filter
and input

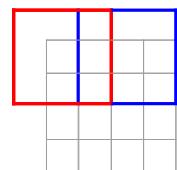


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

144

Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Dot product
between filter
and input



Filter moves 2 pixels in
the input for every one
pixel in the output
Stride gives ratio between
movement in input and
output

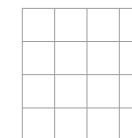


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

145

Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

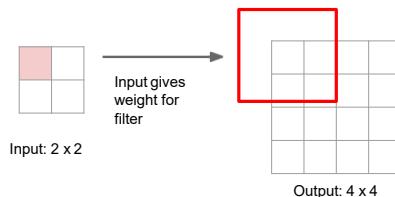


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

146

Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

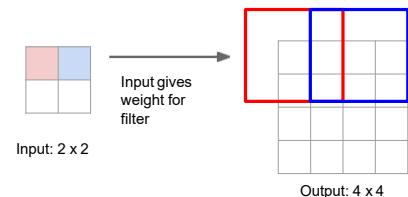


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

147

Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

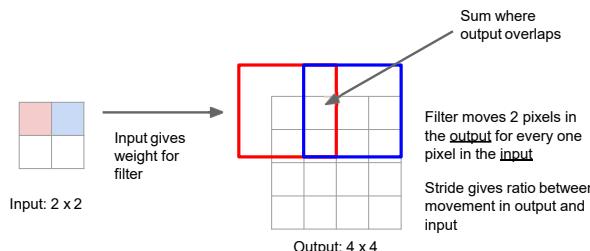


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

148

Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1



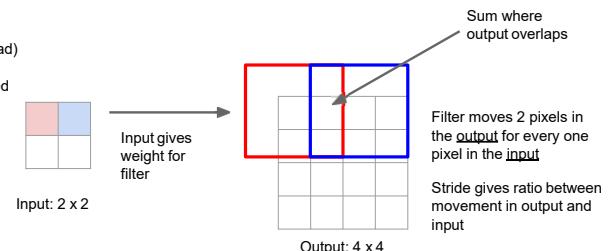
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

149

Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

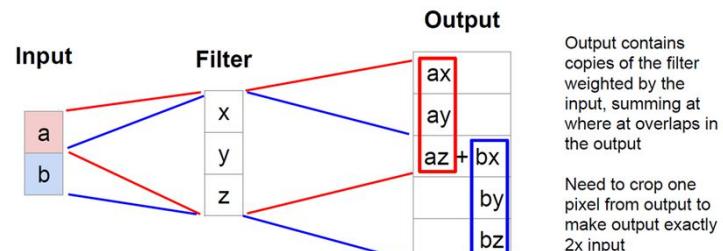
Other names:
 -Deconvolution (bad)
 -Upconvolution
 -Fractionally strided convolution
 -Backward strided convolution



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

150

Learnable Upsampling: Transpose Convolution



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

151

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ ax+by+cz \\ bx+cy+dz \\ cx+dy \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

152

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ ax+by+cz \\ bx+cy+dz \\ cx+dy \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay+bx \\ az+by+cx \\ bz+cy+dx \\ cz+dy \\ dz \end{bmatrix}$$

When stride=1, convolution transpose is just a regular convolution (with different padding rules)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

153

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ ay+bz \\ bx+cy+dz \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

154

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

When stride>1, convolution transpose is no longer a normal convolution!



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

155

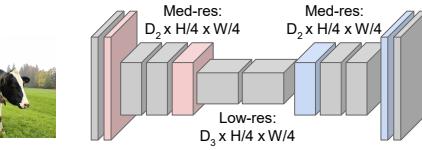
Semantic Segmentation Idea: Fully Convolutional

Downsampling:
Pooling, strided convolution



Input:
3 x H x W

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Upsampling:
Unpooling or strided transpose convolution



Predictions:
H x W



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

156

Metrics for Segmentation Models

• Pixel Accuracy (PA):

$$PA = \frac{\sum_{i=0}^K p_{ii}}{\sum_{i=0}^K \sum_{j=0}^K p_{ij}}$$

— The ratio of pixels properly classified, divided by the total number of pixels.

• Mean Pixel Accuracy (MPA):

$$MPA = \frac{1}{K+1} \sum_{i=0}^K \frac{p_{ii}}{\sum_{j=0}^K p_{ij}}$$

— The ratio of correct pixels is computed in a per-class manner and then averaged over the total number of classes.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

155

Metrics for Segmentation Models

• Intersection over Union (IoU):

$$IoU = J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$IoU = \frac{\text{Area of overlap}}{\text{Area of union}}$$

— Also called **Jaccard Index**

— The most commonly used metrics in semantic segmentation. (mean-IoU/mIoU)

— A denotes the ground truth and B denotes the predicted segmentation maps.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

156

Metrics for Segmentation Models

- Precision / Recall / F1 score:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{F1-score} = \frac{2 \cdot \text{Prec} \cdot \text{Rec}}{\text{Prec} + \text{Rec}}$$

F1: harmonic mean of precision and recall

		Actual Value (as confirmed by experiment)	
		positives	negatives
Predicted Value (predicted by the test)	positives	TP True Positive	FP False Positive
	negatives	FN False Negative	TN True Negative



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Metrics for Segmentation Models

- Dice coefficient:

$$\text{Dice} = \frac{2|A \cap B|}{|A| + |B|}$$

$$\text{Dice} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}} = \text{F1}$$

- Essentially identical to the F1 score.
- The Dice coefficient and IoU are positively correlated.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Instance segmentation



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

161

Instance Segmentation

Classification



CAT

Semantic Segmentation



GRASS, CAT, TREE, SKY

Object Detection



DOG, DOG, CAT

Instance Segmentation



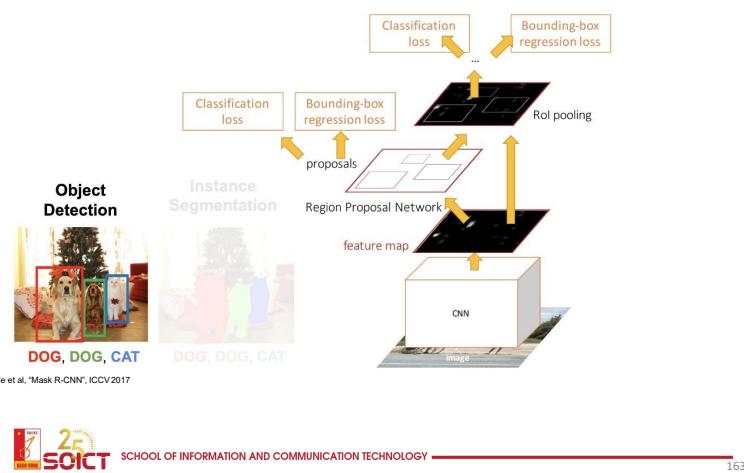
DOG, DOG, CAT



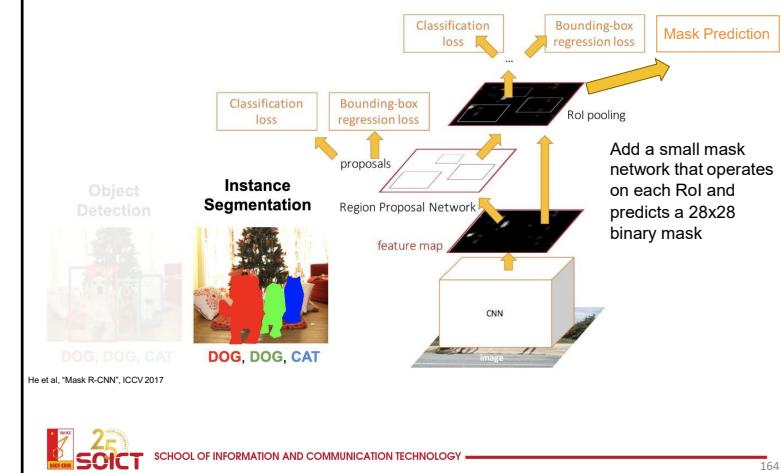
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

162

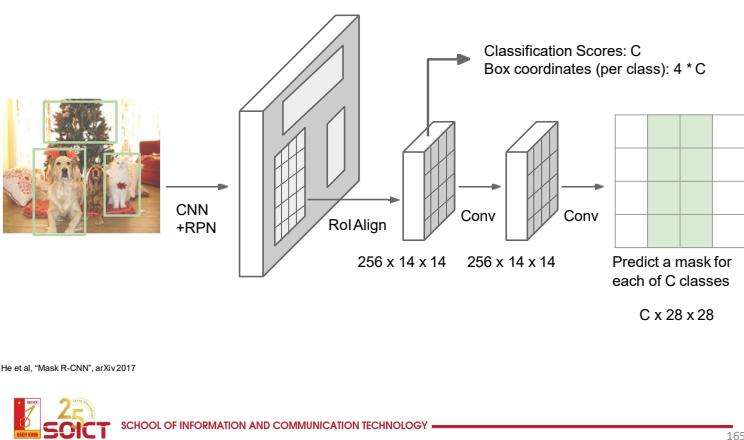
Object Detection: Faster R-CNN



Instance Segmentation: Mask R-CNN



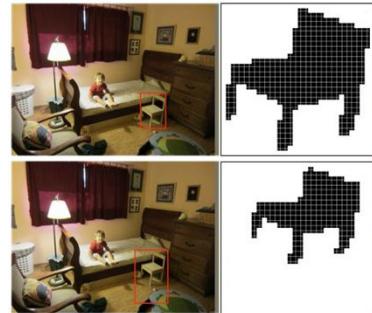
Mask R-CNN



Mask R-CNN: Example Mask Training Targets



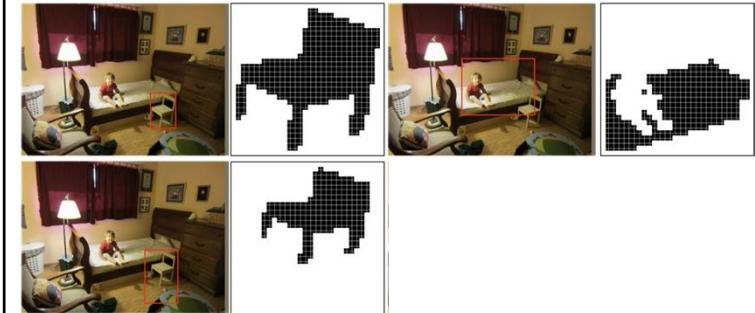
Mask R-CNN: Example Mask Training Targets



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

167

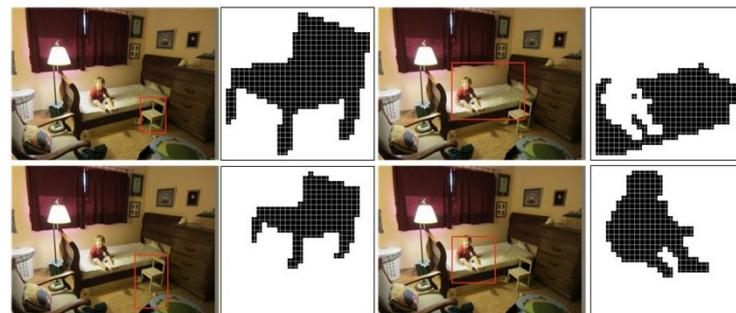
Mask R-CNN: Example Mask Training Targets



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

168

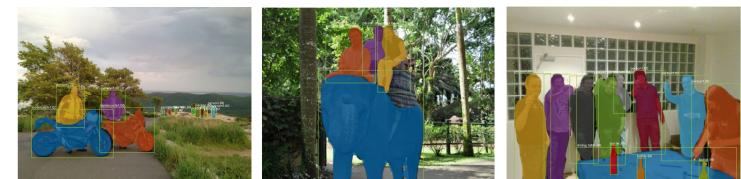
Mask R-CNN: Example Mask Training Targets



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

169

Mask R-CNN: Very Good Results!



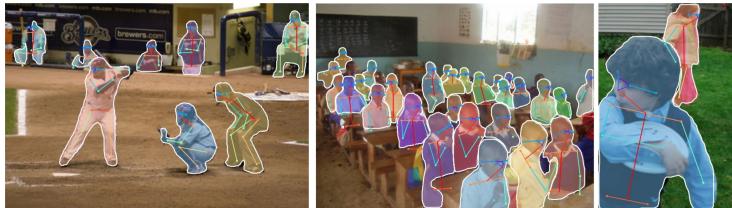
He et al., "Mask R-CNN", ICCV 2017



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

170

Mask R-CNN Also does pose



He et al., "Mask R-CNN". ICCV 2017



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

171

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

172

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

173

Thank
you!



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY



soict.hust.edu.vn/ fb.com/groups/soict