

 HA NOI UNIVERSITY OF SCIENCE AND TECHNOLOGY  
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

# Computer Vision

## Chapter 7 : Deep Learning for CV

1

## Chapter 7: Deep Learning for CV (Part 3)

### Content

- Object detection: sliding-windows
- Two-stage object detection
- One-stage object detection: Anchor-based
- One-stage object detection: Anchor-free
- Semantic segmentation
- Instance segmentation

 SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

2

### Computer Vision Tasks

Classification	Semantic Segmentation	Object Detection	Instance Segmentation
			
CAT	GRASS, CAT, TREE, SKY	DOG, DOG, CAT	DOG, DOG, CAT
No spatial extent	No objects, just pixels	Multiple Object	

 SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

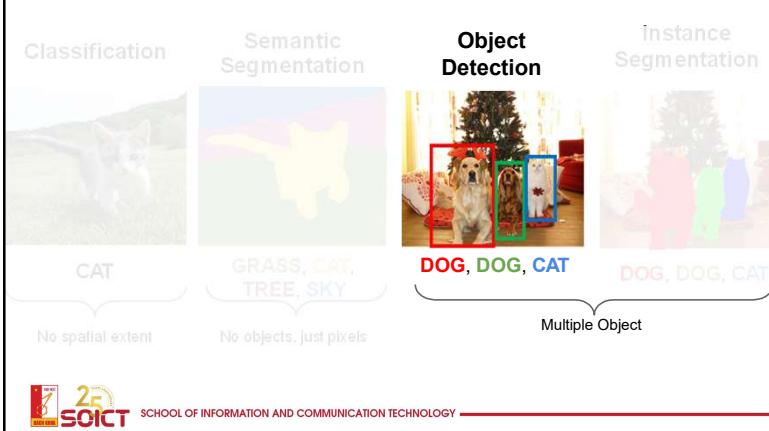
3

### Object detection: Sliding windows

 SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

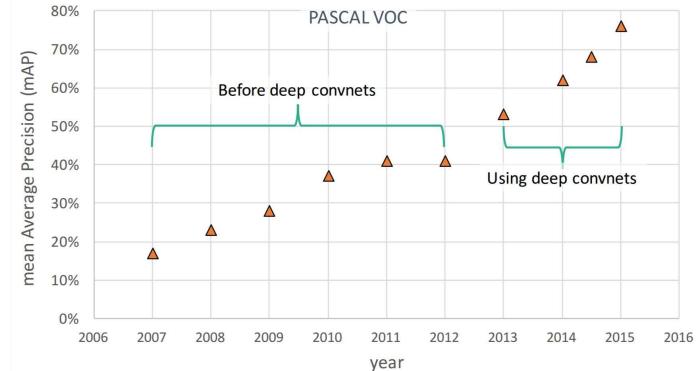
4

## Object Detection



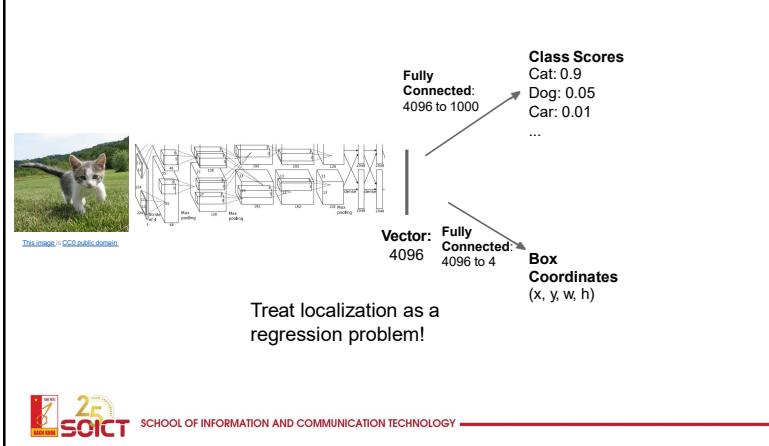
5

## Object Detection: Impact of Deep Learning



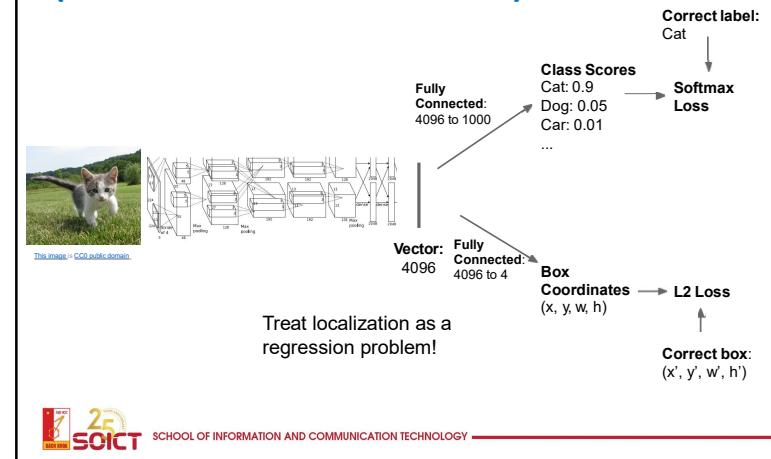
6

## Object Detection: Single Object (Classification + Localization)



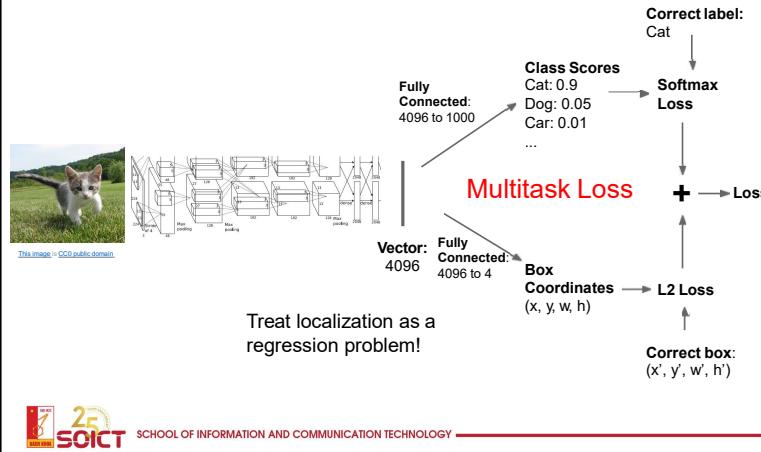
7

## Object Detection: Single Object (Classification + Localization)



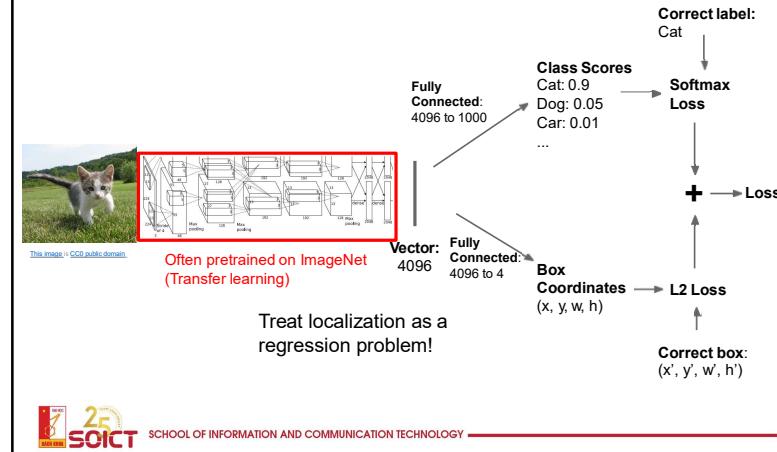
8

## Object Detection: Single Object (Classification + Localization)



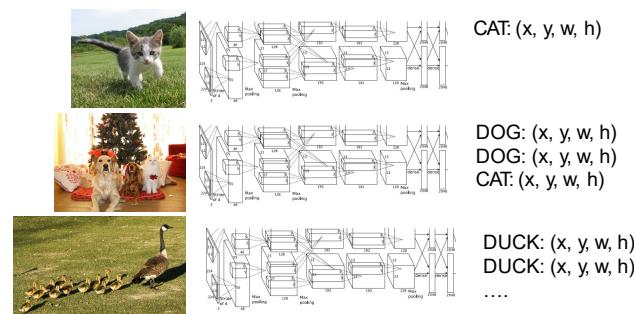
9

## Object Detection: Single Object (Classification + Localization)



10

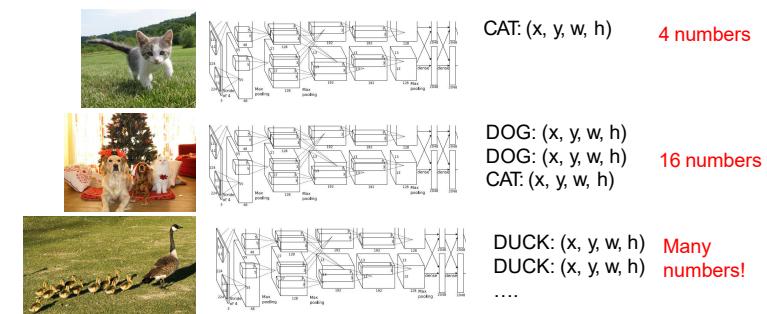
## Object Detection: Multiple Objects



11

## Object Detection: Multiple Objects

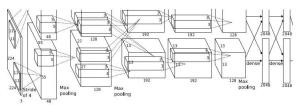
Each image needs a different number of outputs!



12

## 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? NO  
Background? YES

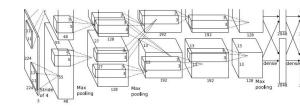


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

13

## Object Detection: Multiple Objects

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



Dog? YES  
Cat? NO  
Background? NO

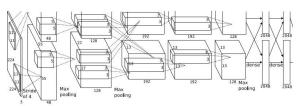


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

14

## Object Detection: Multiple Objects

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



Dog? YES  
Cat? NO  
Background? NO

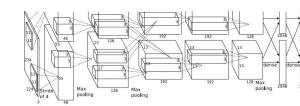


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

15

## 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

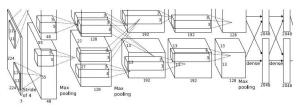
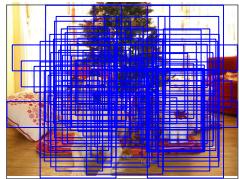


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

16

## 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

## Two-stage Object detection



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

18

## Object Detection

### Two Stages

- Propose “objects”
- Classify each candidate

### One-Stage

- Sliding window to classify all candidates



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

19

## Object Detection

### Two Stages

- Propose “objects”
- Classify each candidate

### One-Stage

- Sliding window to classify all candidates

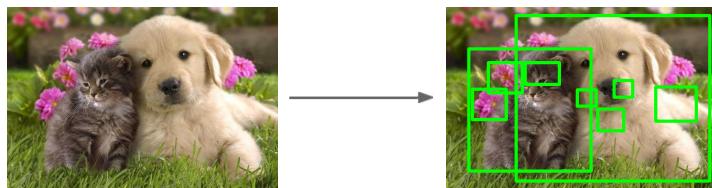


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

20

## 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

21

## R-CNN



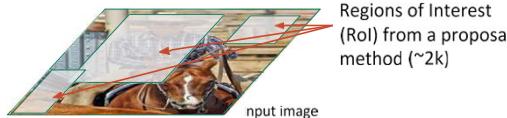
Girshick et al., “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

22

## R-CNN



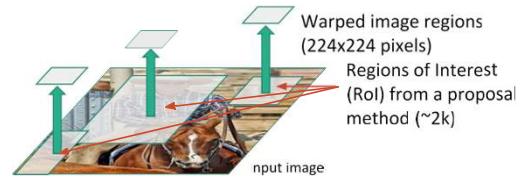
Girshick et al., “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

23

## R-CNN

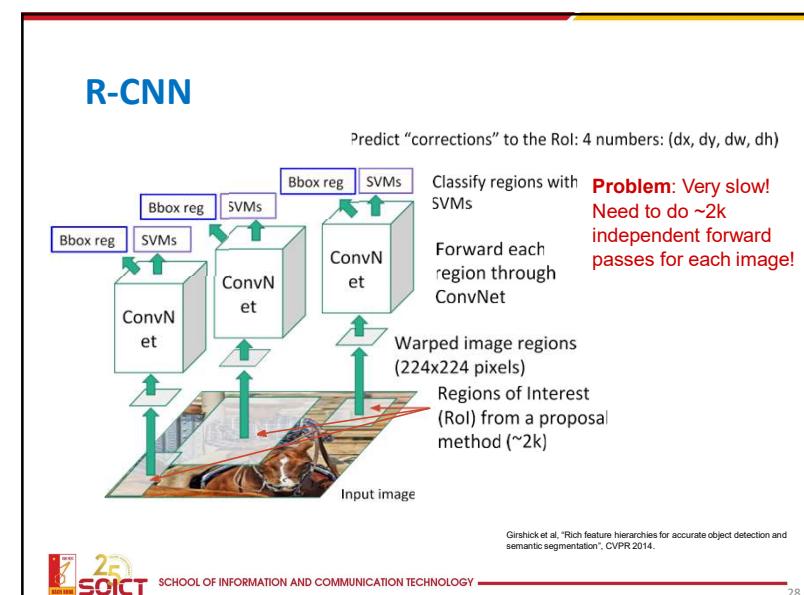
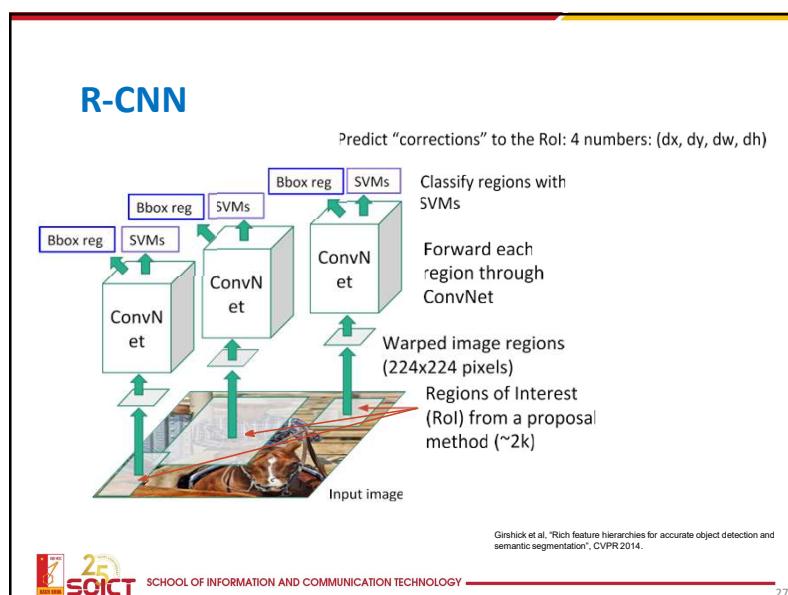
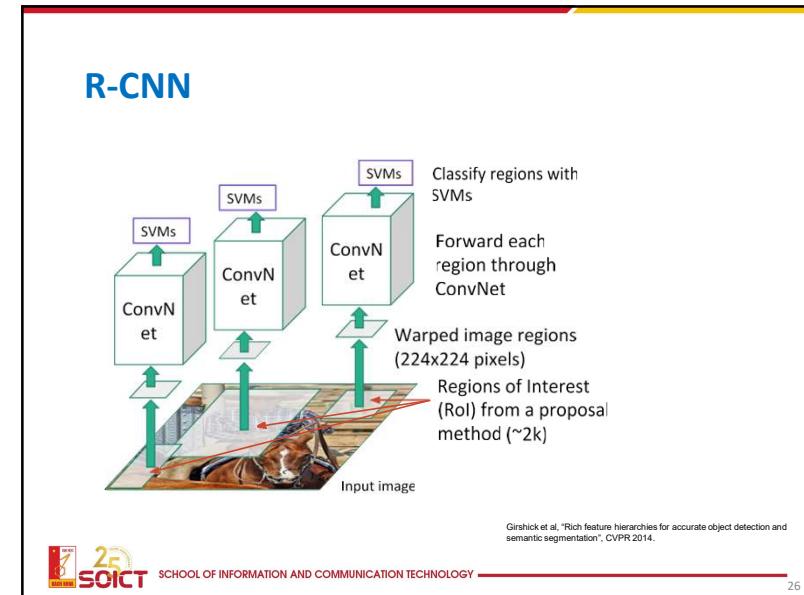
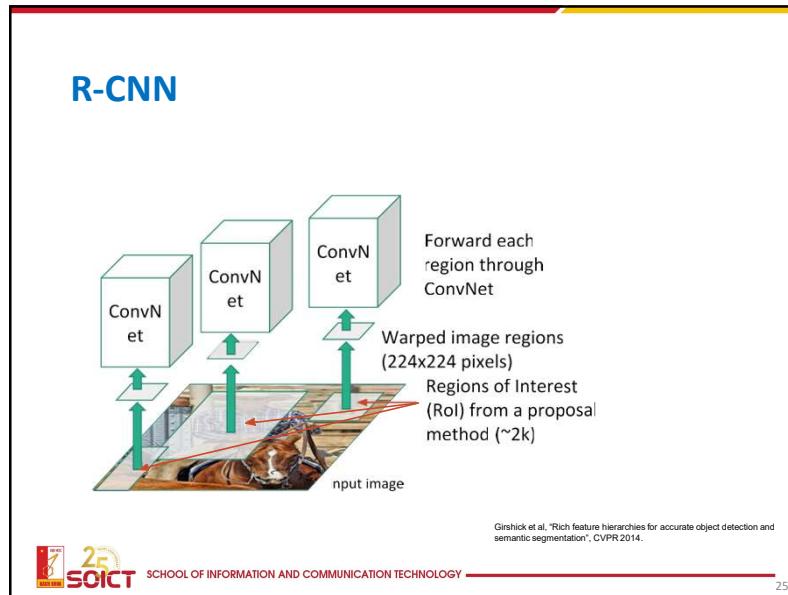


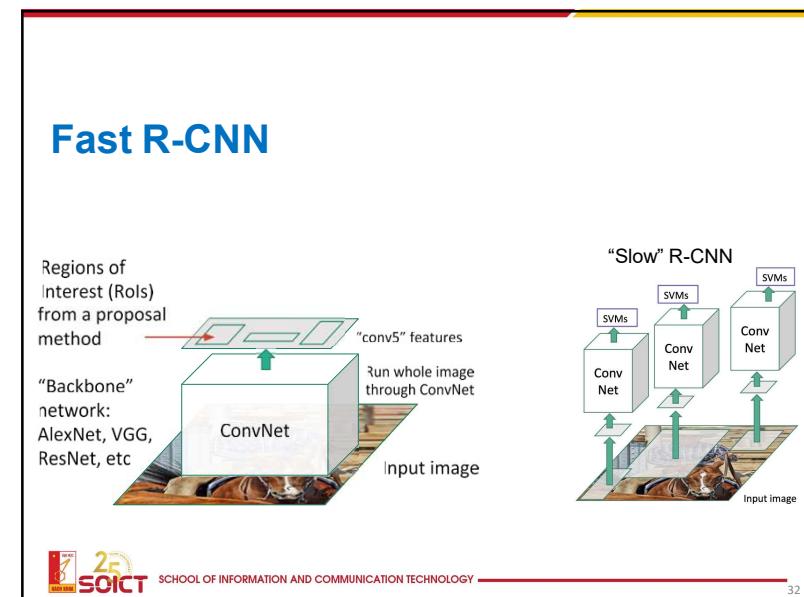
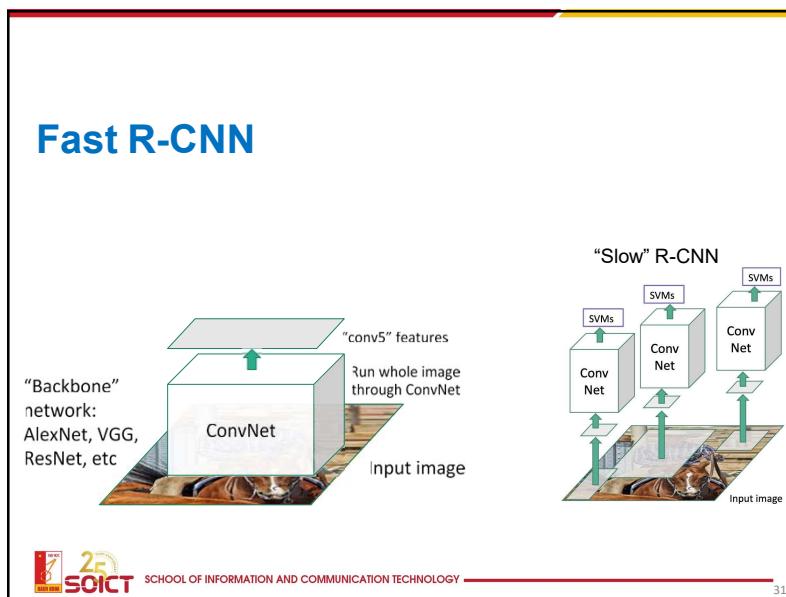
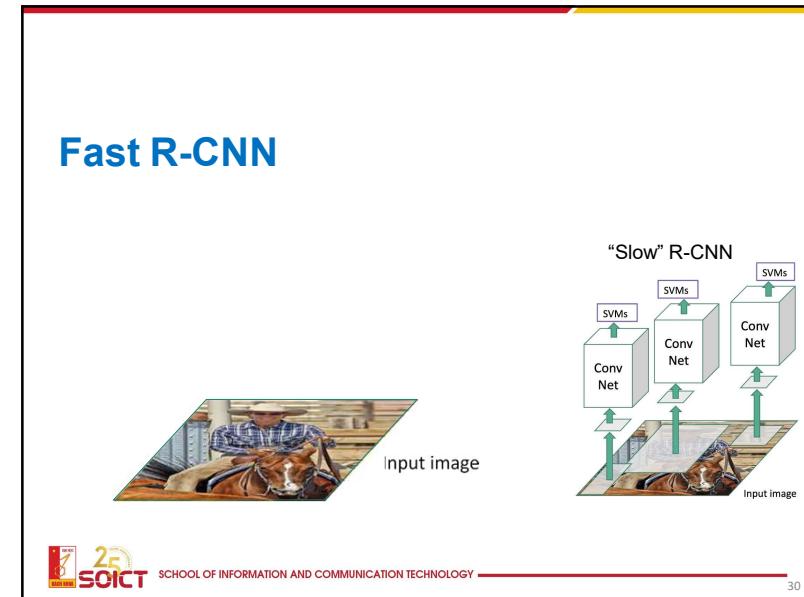
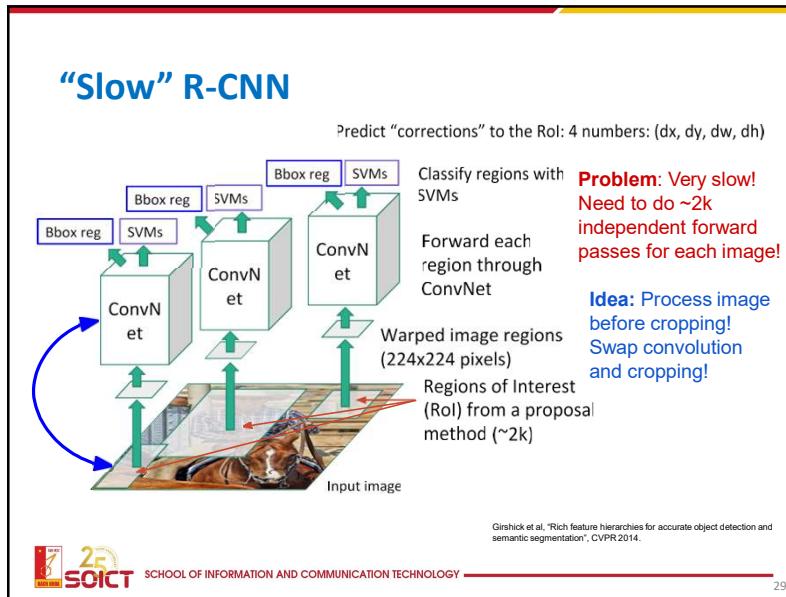
Girshick et al., “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.



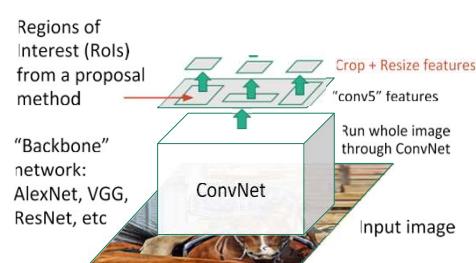
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

24





## Fast R-CNN

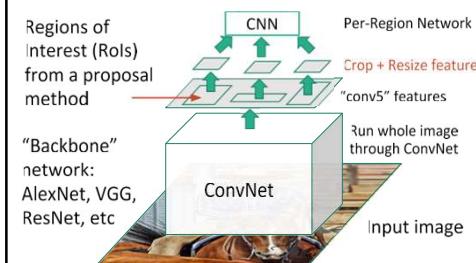


25 SOICT

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

33

## Fast R-CNN

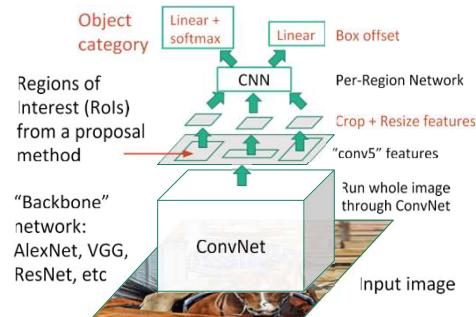


25 SOICT

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

34

## Fast R-CNN

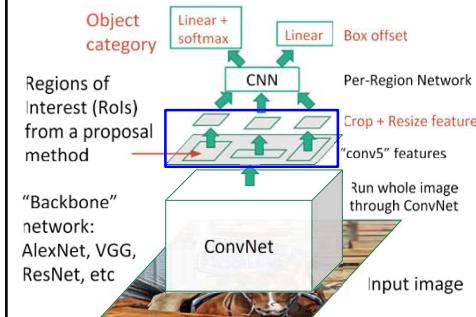


25 SOICT

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

35

## Fast R-CNN



25 SOICT

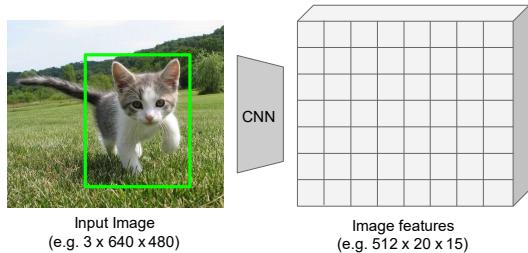
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

36

35

36

## Cropping Features: RoI Pool



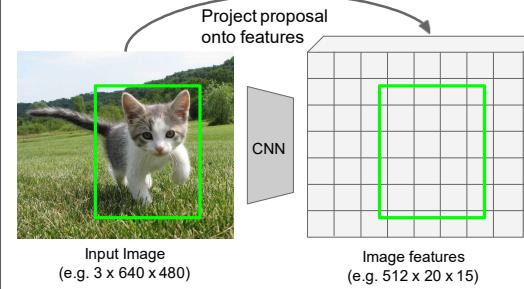
Girshick, "Fast R-CNN", ICCV 2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

37

## Cropping Features: RoI Pool



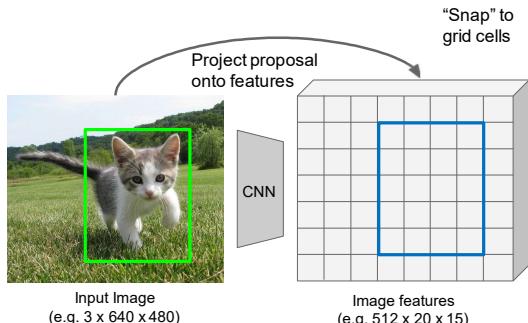
Girshick, "Fast R-CNN", ICCV 2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

38

## Cropping Features: RoI Pool



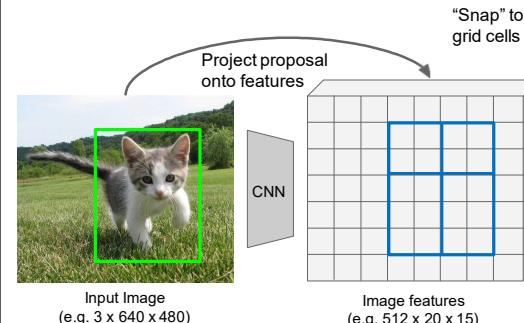
Girshick, "Fast R-CNN", ICCV 2015.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

39

## Cropping Features: RoI Pool

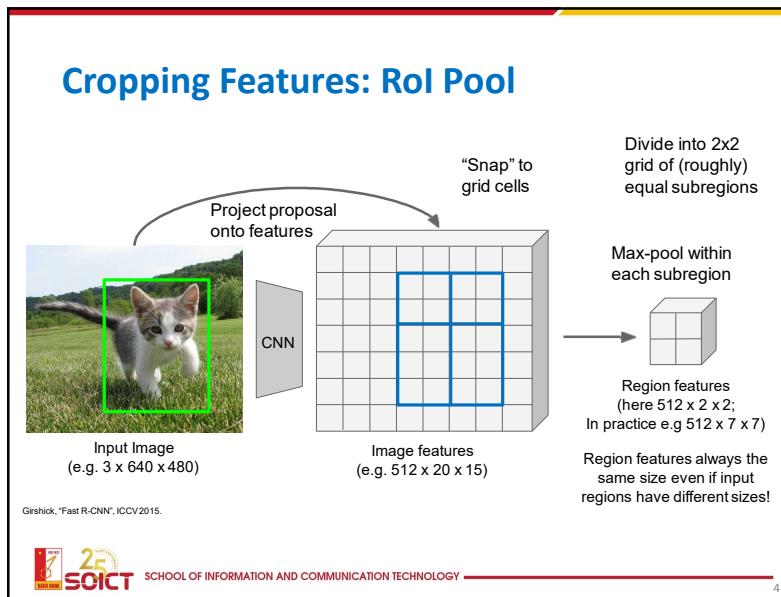


Girshick, "Fast R-CNN", ICCV 2015.

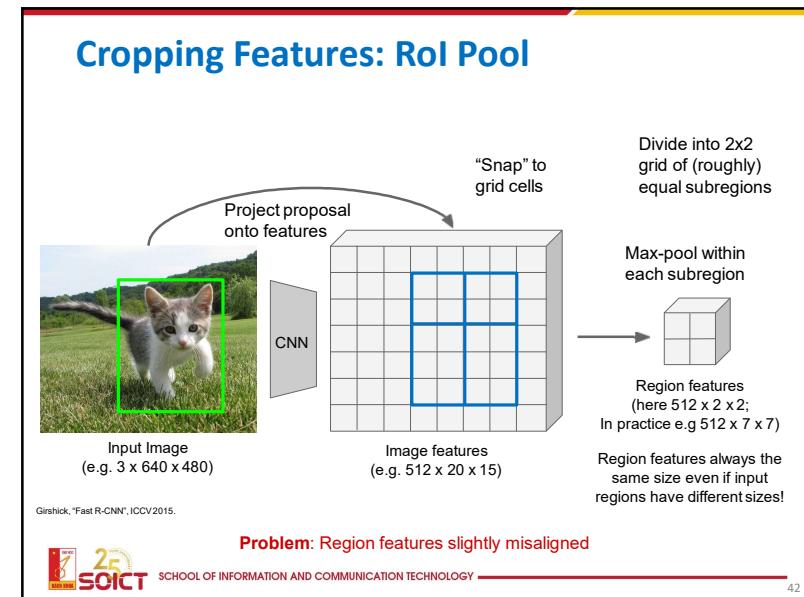


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

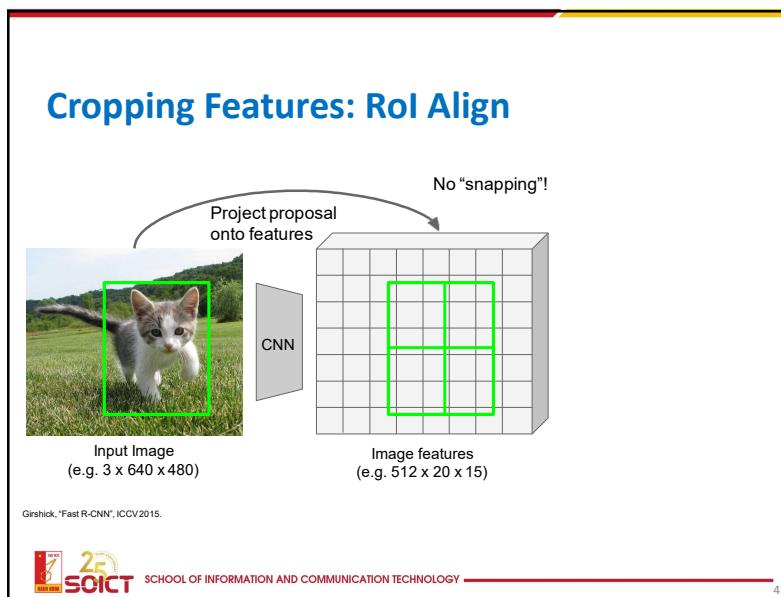
40



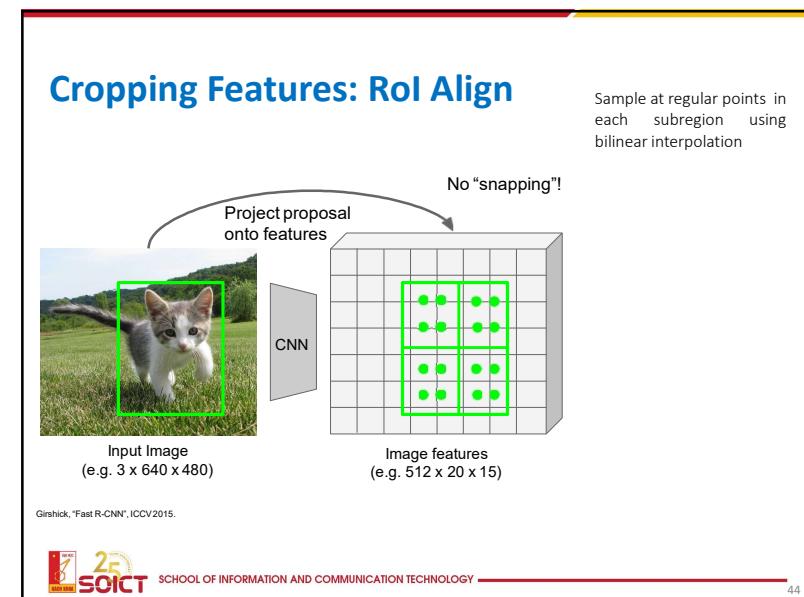
41



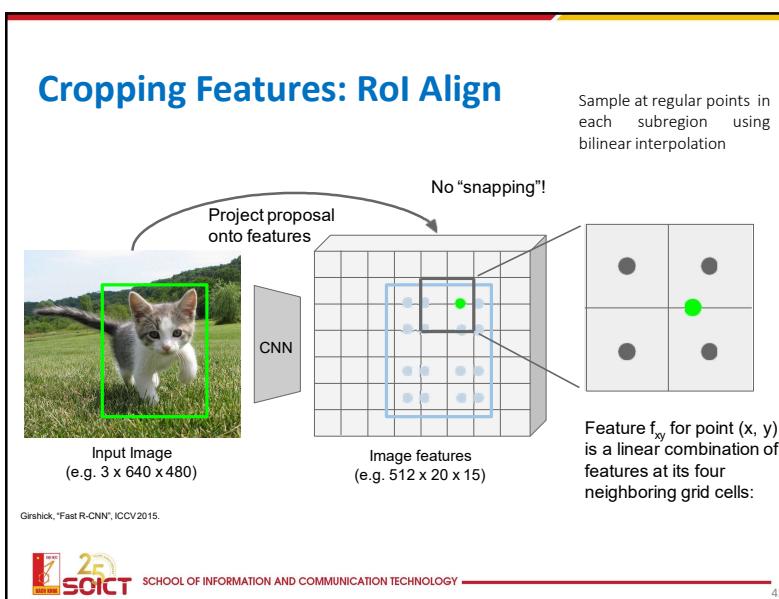
42



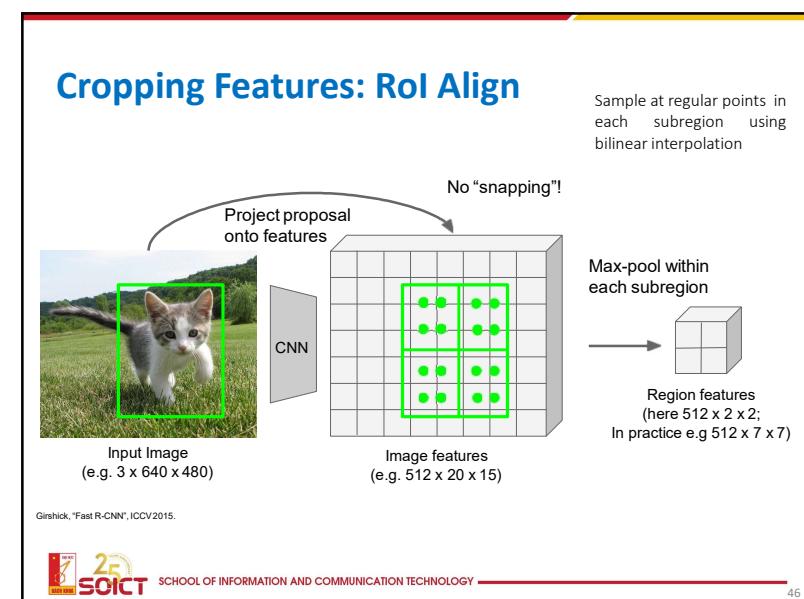
43



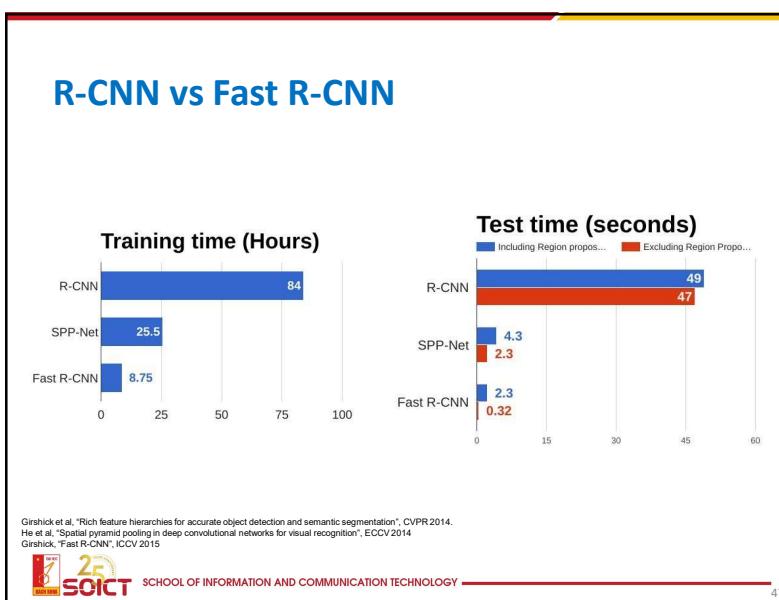
44



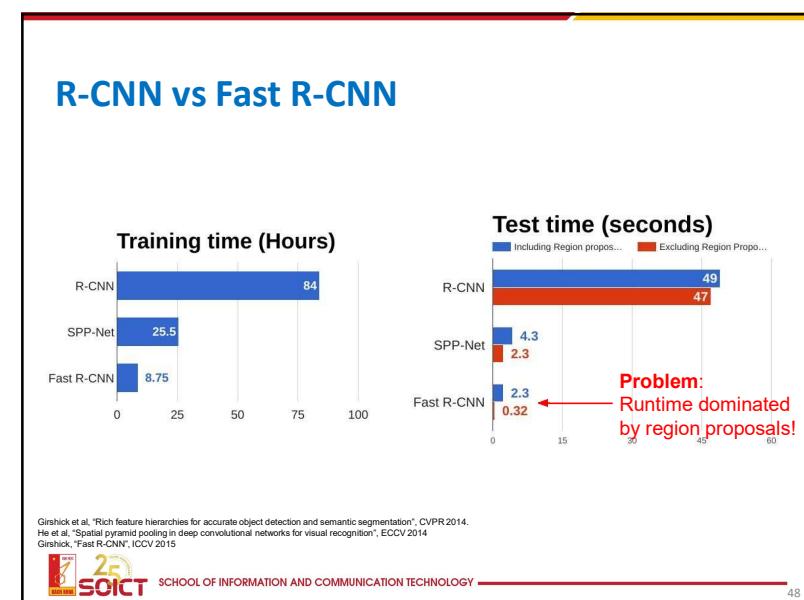
45



46



47

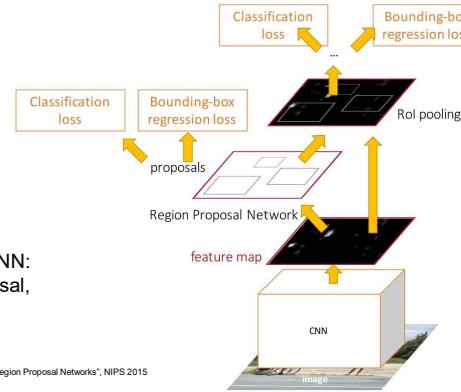


48

## Faster R-CNN: Make CNN do proposals!

Insert Region Proposal Network (RPN) to predict proposals from features

Otherwise same as Fast R-CNN:  
Crop features for each proposal,  
classify each one



Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015  
Figure copyright 2015, Ross Girshick; reproduced with permission



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

49

## Region Proposal Network

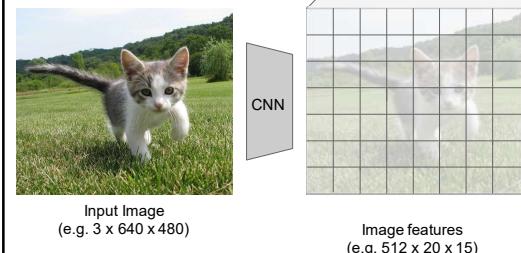


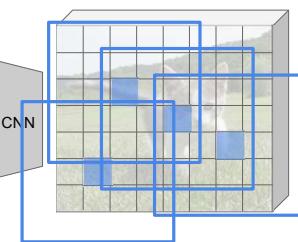
Image features  
(e.g. 512 x 20 x 15)

50

## Region Proposal Network



Input Image  
(e.g. 3 x 640 x 480)



Imagine an anchor box  
of fixed size at each  
point in the feature map



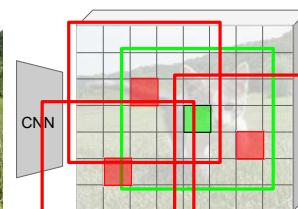
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

51

## Region Proposal Network



Input Image  
(e.g. 3 x 640 x 480)



Imagine an anchor box  
of fixed size at each  
point in the feature map

At each point, predict  
whether the corresponding  
anchor contains an object  
(per-pixel logistic regression)

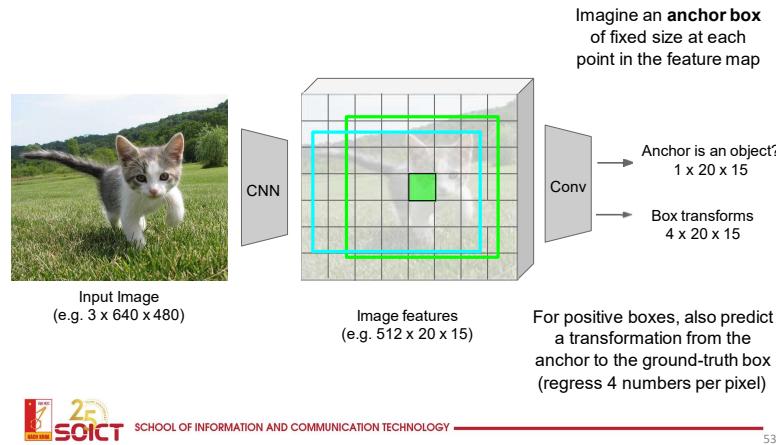
Anchor is an object?  
1 x 20 x 15



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

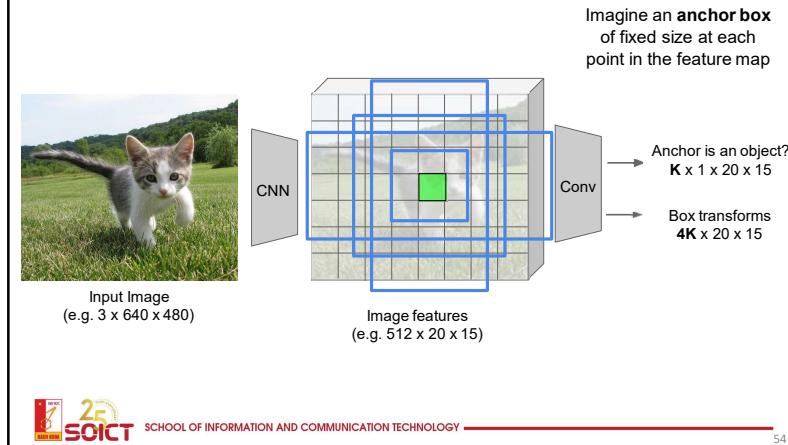
52

## Region Proposal Network



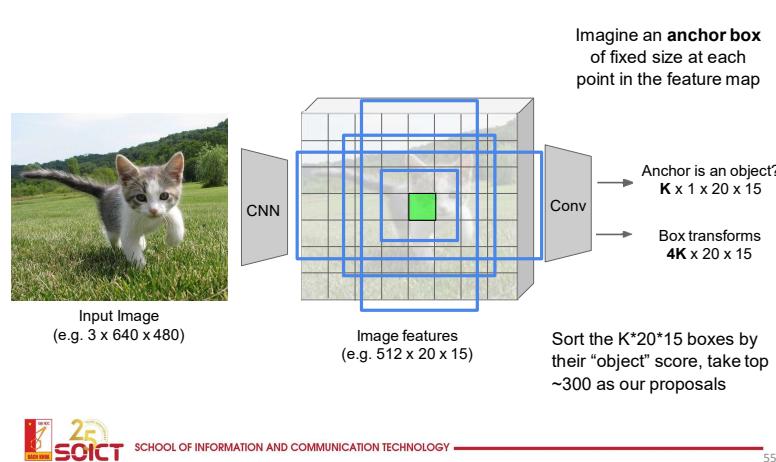
53

## Region Proposal Network



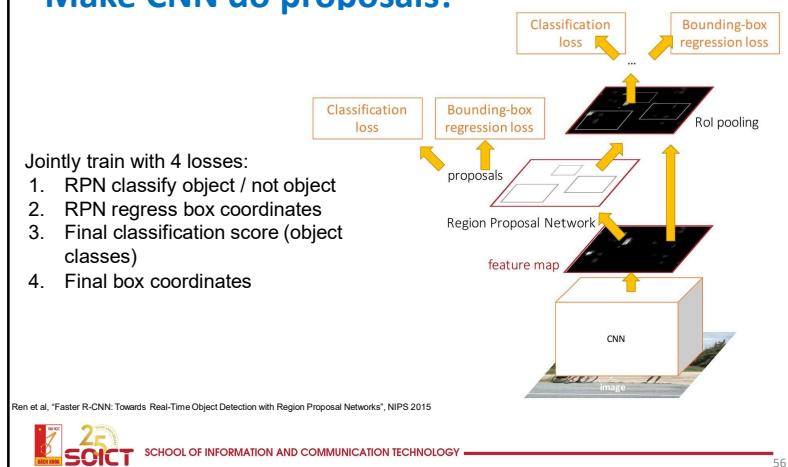
54

## Region Proposal Network



55

## Faster R-CNN: Make CNN do proposals!



56

## Faster R-CNN: Make CNN do proposals!

R-CNN Test-Time Speed



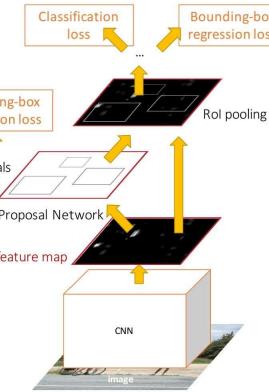
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

57

## Faster R-CNN: Make CNN do proposals!

Glossing over many details:

- Ignore overlapping proposals with **non-max suppression**
- How to determine whether a proposal is positive or negative?
- How many positives / negatives to send to second stage?
- How to parameterize bounding box regression?



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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

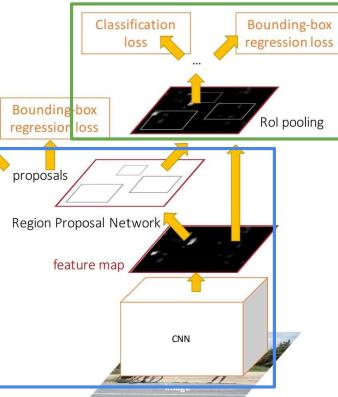
58

## 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

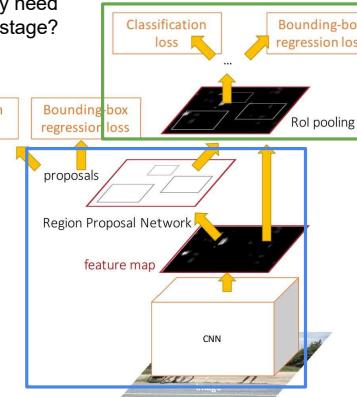
## Faster R-CNN: Make CNN do proposals!

Do we really need  
the second stage?

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

60

## One-stage Object detection

Anchor-based



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

61

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

62

## Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Input image  
3 x H x W



Divide image into grid  
7 x 7  
Image a set of **base boxes**  
centered at each grid cell  
Here B = 3

- 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 x 7 x (5 \* B + C)

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

63

## Imbalance

Number of “negative” anchors ~O(10K)  
Number of “positive” anchors ~O(10)

### What happens to CE loss in this case?

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

On Board

**Loss and gradient are dominated by correctly classified negative examples**

Training outcome → constant “negative” prediction.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

64

64

## Imbalance – Focal Loss

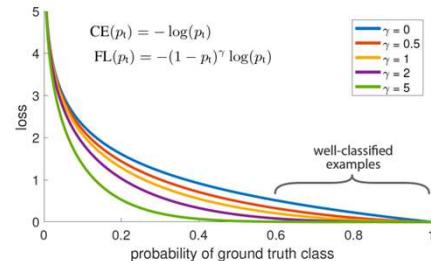


Figure 1: We propose a novel loss we term the *Focal Loss* that adds a factor  $(1 - p_t)^\gamma$  to the standard cross entropy criterion. Setting  $\gamma > 0$  reduces the relative loss for well-classified examples ( $p_t > .5$ ), putting more focus on hard, misclassified examples. As our experiments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

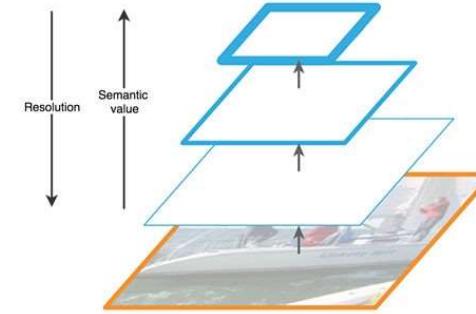


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

65

## Feature Pyramid Network (FPN)

- How to handle multiscale predictions?

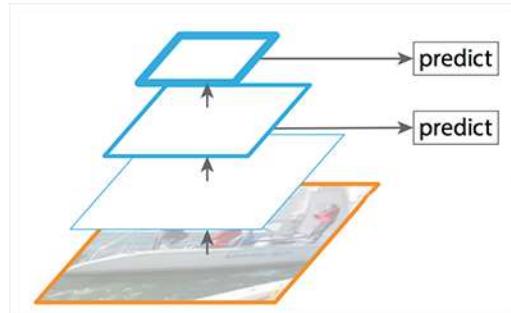


Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. **Feature Pyramid Networks for Object Detection** (CVPR 2017)

66

## Feature Pyramid Network (FPN)

- How to handle multiscale predictions?

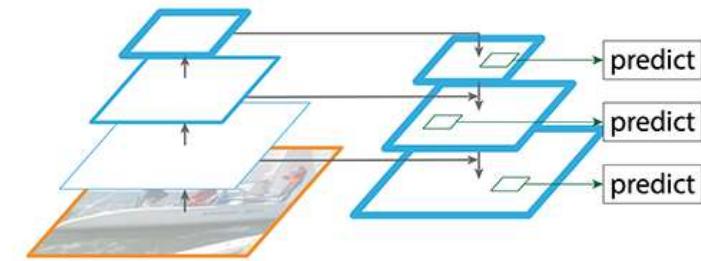


Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. **Feature Pyramid Networks for Object Detection** (CVPR 2017)

67

## Feature Pyramid Network (FPN)

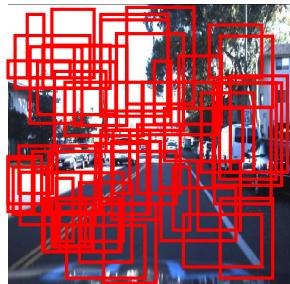
- How to handle multiscale predictions?



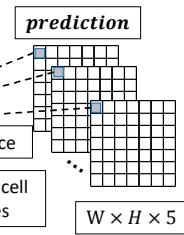
Tsung-Yi, Dollár, Girshick, He, Hariharan and Belongie. **Feature Pyramid Networks for Object Detection** (CVPR 2017)

68

## Postprocessing: NMS



$W \times H$  cells  
convert : relative to cell → relative to image



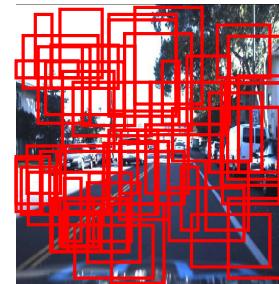
69



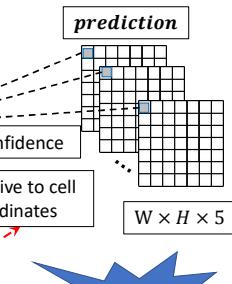
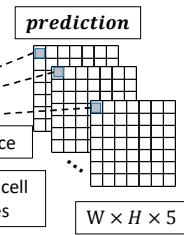
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

69

## Postprocessing: NMS



$W \times H$  cells  
convert : relative to cell → relative to image



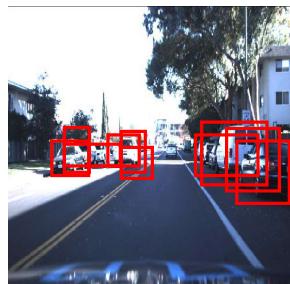
70



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

70

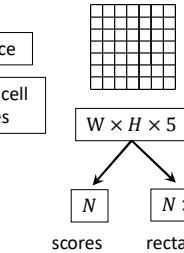
## Postprocessing: NMS



$W \times H$  cells  
convert : relative to cell → relative to image

Threshold  $T = 0.5$ , for example

$N$  – rectangle number with  $p > T$



71



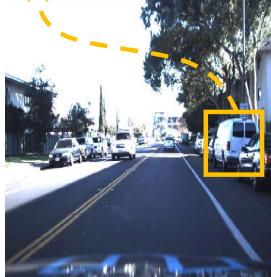
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

71

## Postprocessing: NMS

$rect_1, rect_2, rect_3, \dots, rect_N$   
Sorted confidence  $p$

0.9 0.85 ...



72



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

72

## Postprocessing: NMS

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

Sorted confidence  $p$

0.9    0.85

Compare IOU of the 1<sup>st</sup> rectangle with others



73



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

73

## Postprocessing: NMS

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

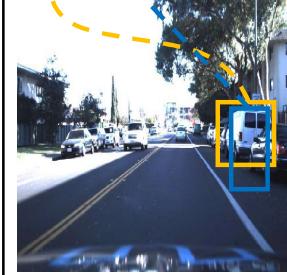
Sorted confidence  $p$

0.9    0.85

Compare IOU of the 1<sup>st</sup> rectangle with others

$threshold$

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



74



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

74

## Postprocessing: NMS

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

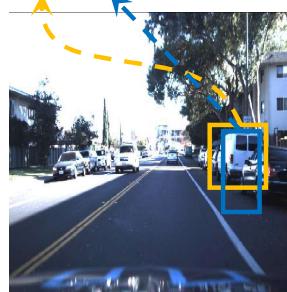
Sorted confidence  $p$

0.9    0.5    0.82

Compare IOU of the 1<sup>st</sup> rectangle with others

$threshold$

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



75



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

75

## Postprocessing: NMS

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

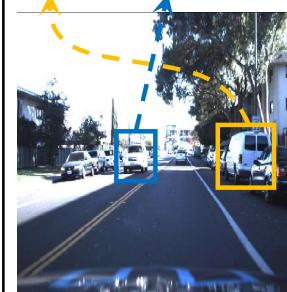
Sorted confidence  $p$

0.9    0.5    0.82

Compare IOU of the 1<sup>st</sup> rectangle with others

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

*do nothing*



76



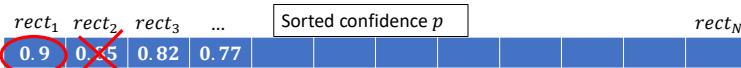
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

76

75

76

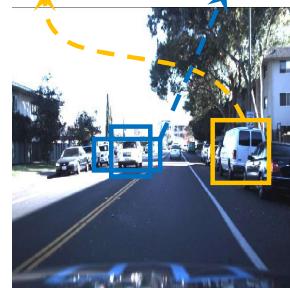
## Postprocessing: NMS



Compare IOU of the 1<sup>st</sup> rectangle with others

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

IOU=0 with the chosen rectangle → **do nothing!**



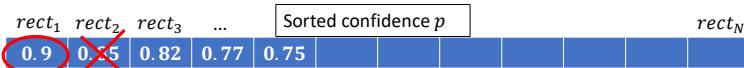
77



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

77

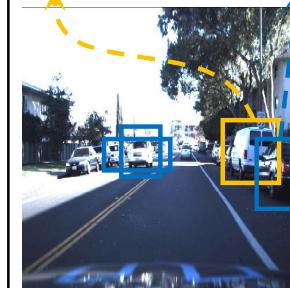
## Postprocessing: NMS



Compare IOU of the 1<sup>st</sup> rectangle with others

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

**do nothing**



78



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

78

## Postprocessing: NMS



Compare IOU of the 2<sup>nd</sup> rectangle with others

$N_1 \leq N$  because we have thrown out rectangles



79



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

79

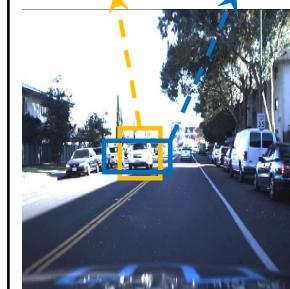
## Postprocessing: NMS



Compare IOU of the 2<sup>nd</sup> rectangle with others

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

threshold



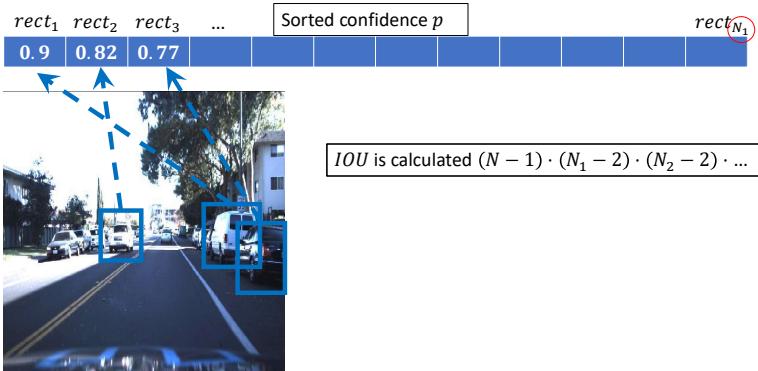
80



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

80

## Postprocessing: NMS



81



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

81

## One-stage Object detection

Anchor-free



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

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]

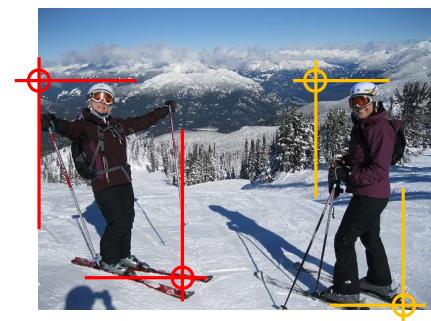
### 2. Extra hyperparameters – sizes and aspect ratios



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

83

## CornerNet: Detecting Objects as Paired Keypoints

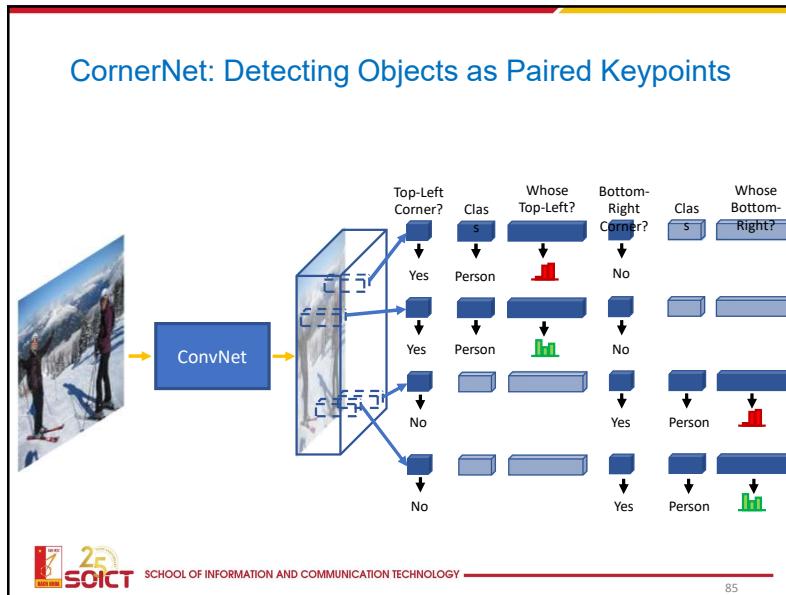


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

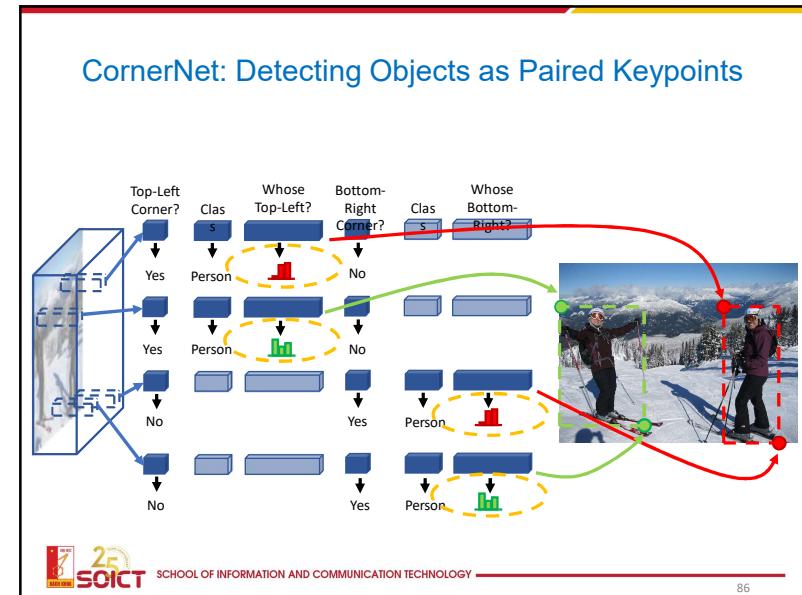
84

83

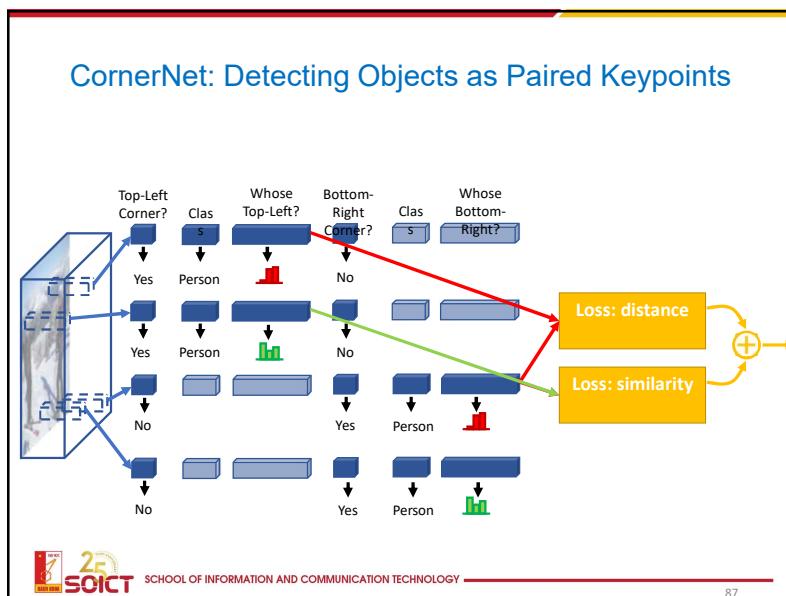
84



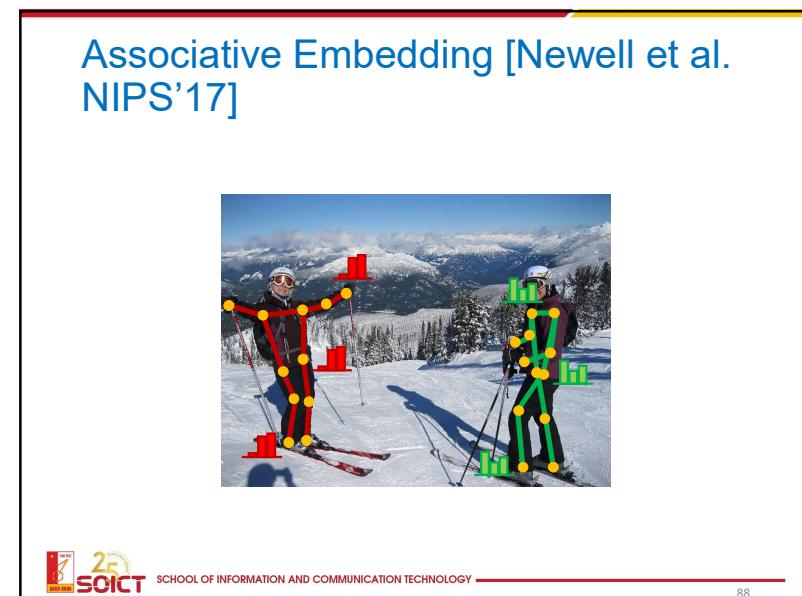
85



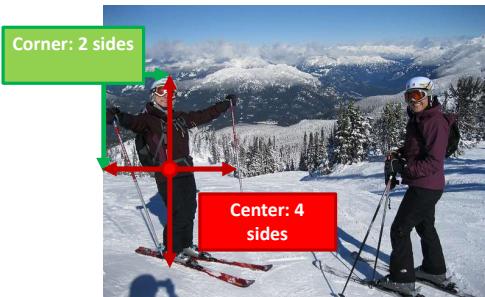
86



87



## Advantages of Detecting Corners



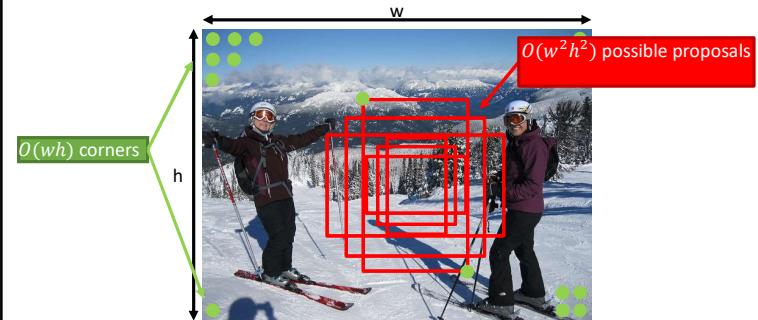
Detecting corner is easier than detecting center



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

89

## Advantages of Detecting Corner



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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

90

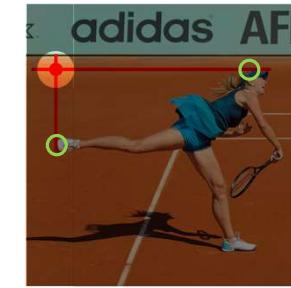
## Supervising Corner Detection



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

91

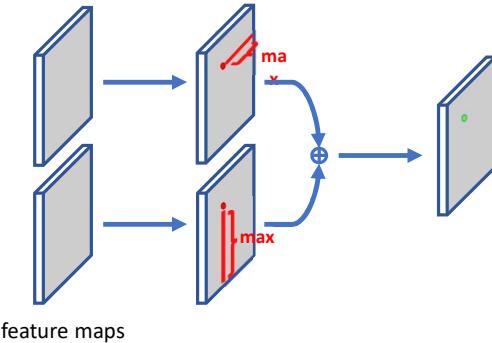
## Corner Pooling



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

92

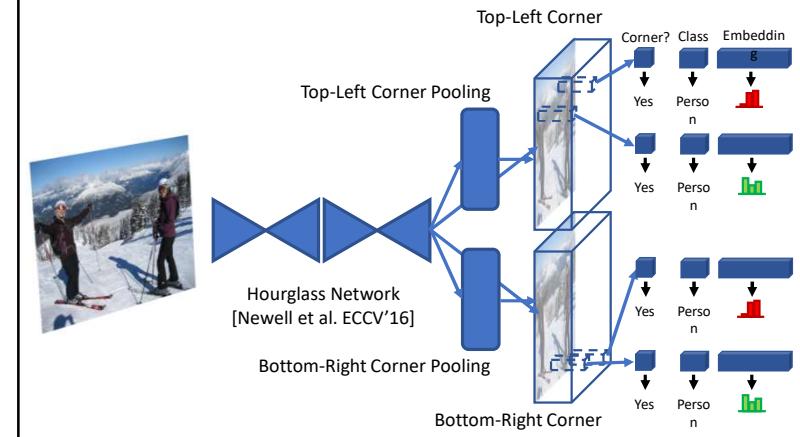
## Top-Left Corner Pooling



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

93

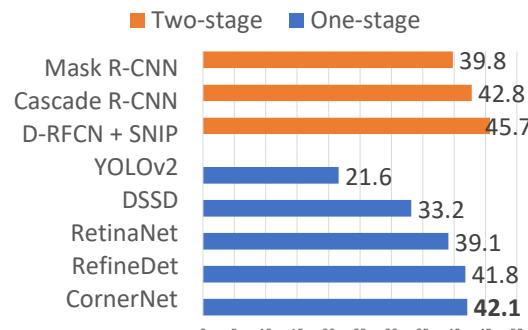
## CornerNet



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

94

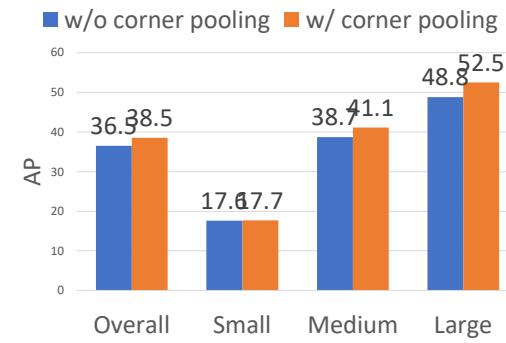
## Experiment: CornerNet versus Others



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

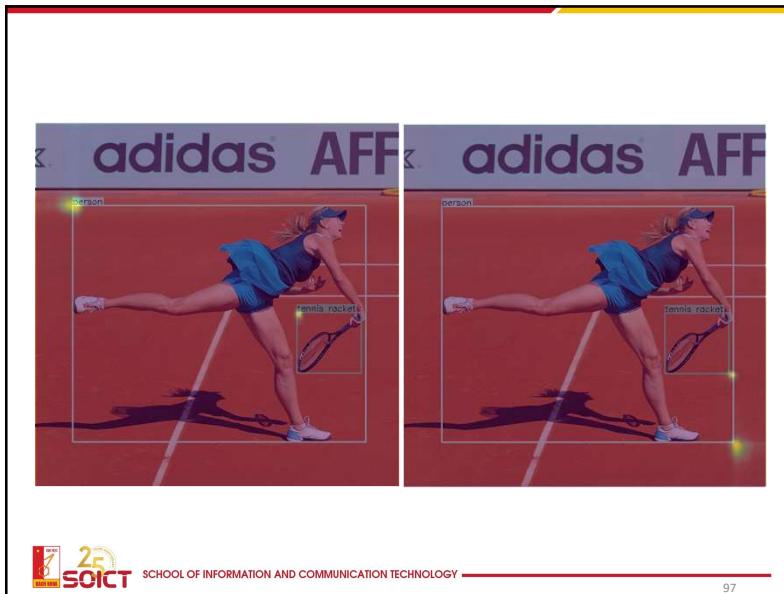
95

## Experiment: Corner Pooling



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

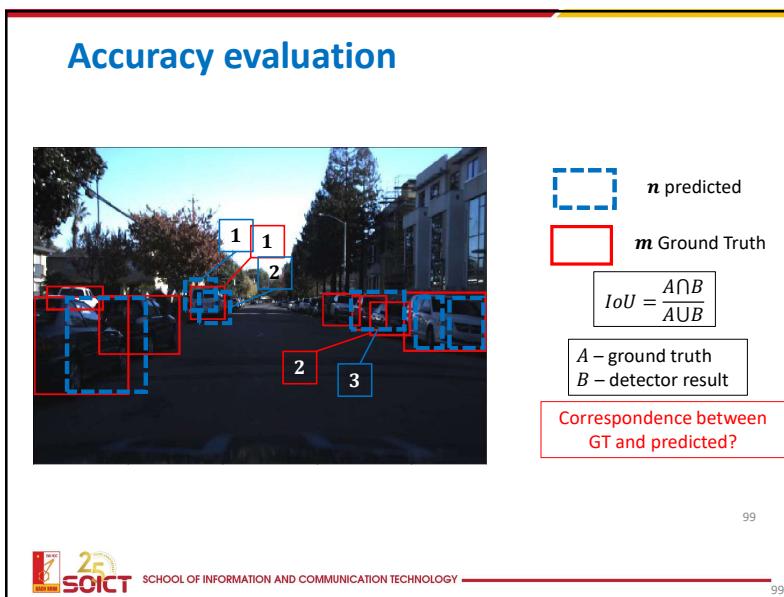
96



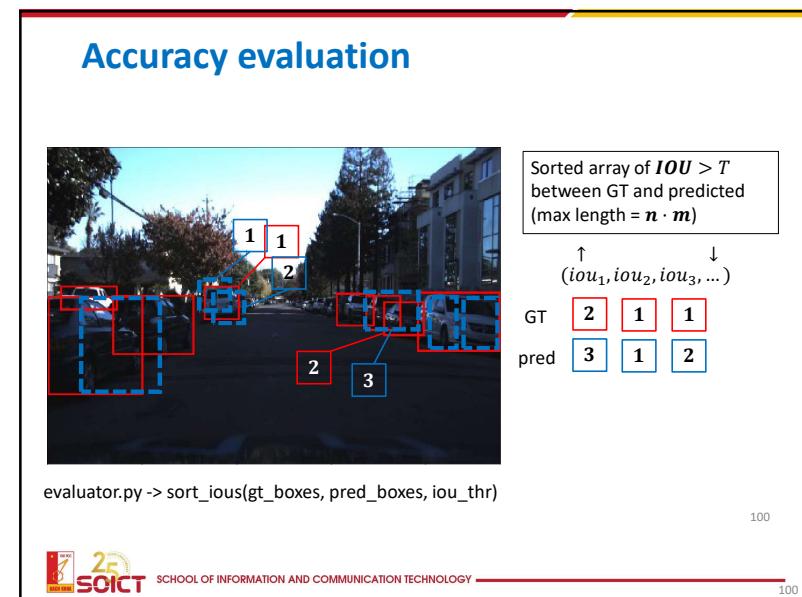
97



98

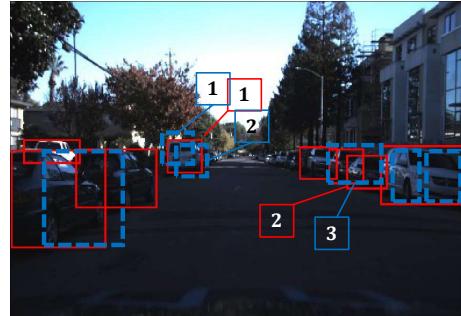


99



100

## Accuracy evaluation



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

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

GT    2    1    1  
pred 3    1    2

if appear  
firstly

matched GT 2    1

matched pred 3    1

101

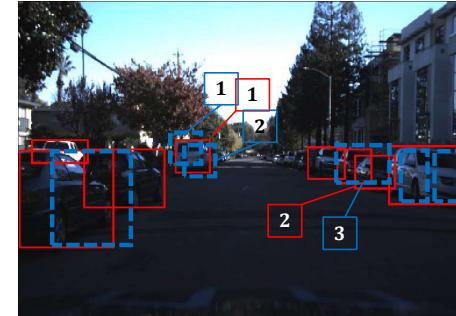
evaluator.py ->  
get\_single\_image\_results(gt\_boxes, pred\_boxes, iou\_thr)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

101

## Accuracy evaluation



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

True predicted

matched GT 2    1

matched pred 3    1

102

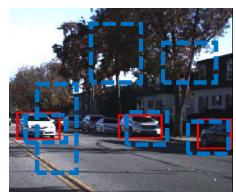
evaluator.py ->  
get\_single\_image\_results(gt\_boxes, pred\_boxes, iou\_thr)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

102

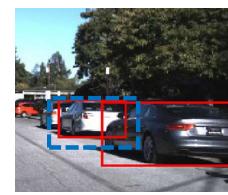
## Accuracy evaluation: precision, recall



$T = 0.5$

...

getTruePredicted



precision -?  
recall -?

$$\text{precision} = \frac{\text{true predicted}}{\text{predicted}}$$

What part of predicted is true

precision -?  
recall -?

$$\text{recall} = \frac{\text{true predicted}}{\text{ground truth}}$$

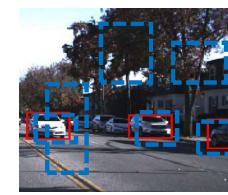
What part of true predicted  
from all GT objects



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

103

## Accuracy evaluation: precision, recall



$T = 0.5$

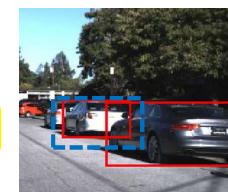
...

getTruePredicted

- predicted = 4
- ground truth = 3
- true predicted = 1

$$\text{precision} = \frac{\text{true predicted}}{\text{predicted}}$$

What part of predicted is true



- predicted = 1
- ground truth = 2
- true predicted = 1

$$\text{recall} = \frac{\text{true predicted}}{\text{ground truth}}$$

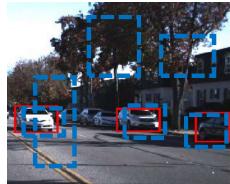
What part of true predicted  
from all GT objects



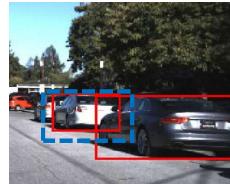
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

104

## Accuracy evaluation: precision, recall

 $T = 0.5$  $\dots$ 

getTruePredicted



- predicted = 4
- ground truth = 3
- true predicted = 1

$$\text{precision} = \frac{\text{true predicted}}{\text{predicted}}$$

- predicted = 1
- ground truth = 2
- true predicted = 1

$$\text{recall} = \frac{\text{true predicted}}{\text{ground truth}}$$

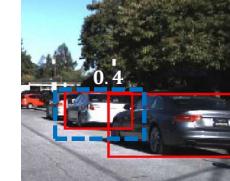
evaluator.py -&gt; calc\_precision\_recall(img\_results)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

105

## Accuracy evaluation: precision, recall

sorted box  
confidence $\downarrow (p_1, p_2, \dots, p_i, \dots, p_{N-1}, p_N) \uparrow$ Each predicted box  
has confidence  $p$ 

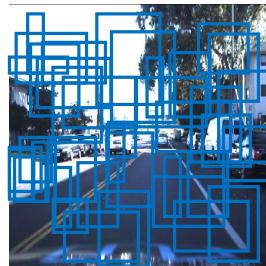
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

106

## Accuracy evaluation: precision, recall

sorted box  
confidence  
 $\downarrow (p_1, p_2, \dots, p_i, \dots, p_{N-1}, p_N) \uparrow$ I. Get boxes with  $p \geq p_1$ , i.e. all boxes

calcPrecisionRecall



score	prec	recall
$p_1$	0.1	0.9
		high recall
		low precision
		For all images together!



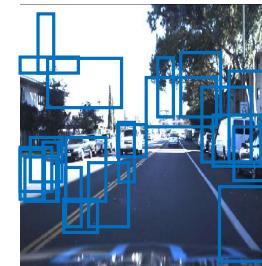
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

107

## Accuracy evaluation: precision, recall

sorted box  
confidence  
 $\downarrow (p_1, p_2, \dots, p_i, \dots, p_{N-1}, p_N) \uparrow$ I. Get boxes with  $p \geq p_1$ , i.e. all boxes

calcPrecisionRecall



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

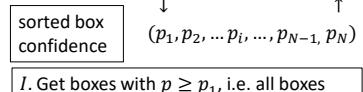
For all images together!



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

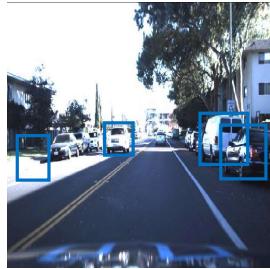
108

## Accuracy evaluation: precision, recall



**calcPrecisionRecall**

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



High precision

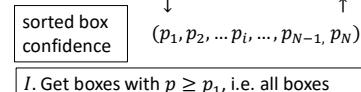
For all images together!



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

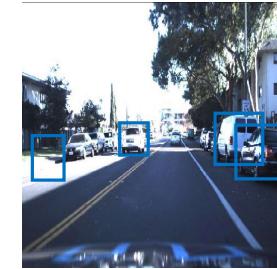
109

## Accuracy evaluation: precision, recall



**calcPrecisionRecall**

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



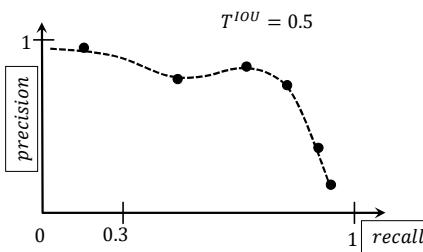
High precision  
get\_thr\_prec\_rec(...)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

110

## Accuracy evaluation: mAP



hundreds of values for real data sets

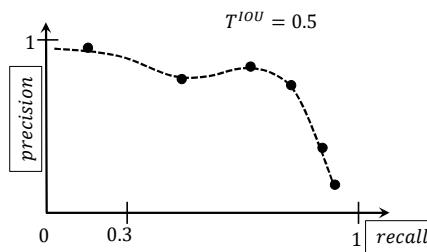
score	prec	recall
$p_1$	0.1	0.9
$p_2$	0.2	0.8
...	...	...
$p_N$	0.9	0.1



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

111

## 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

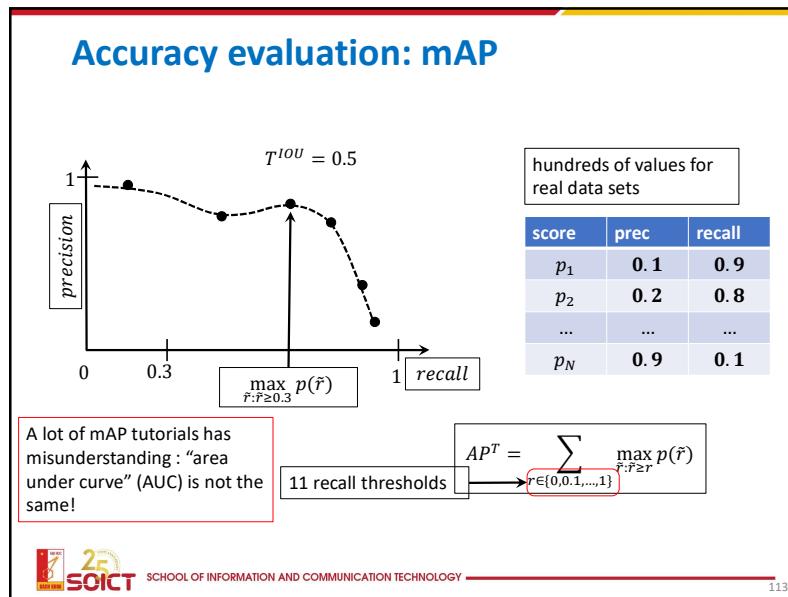
11 recall thresholds

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

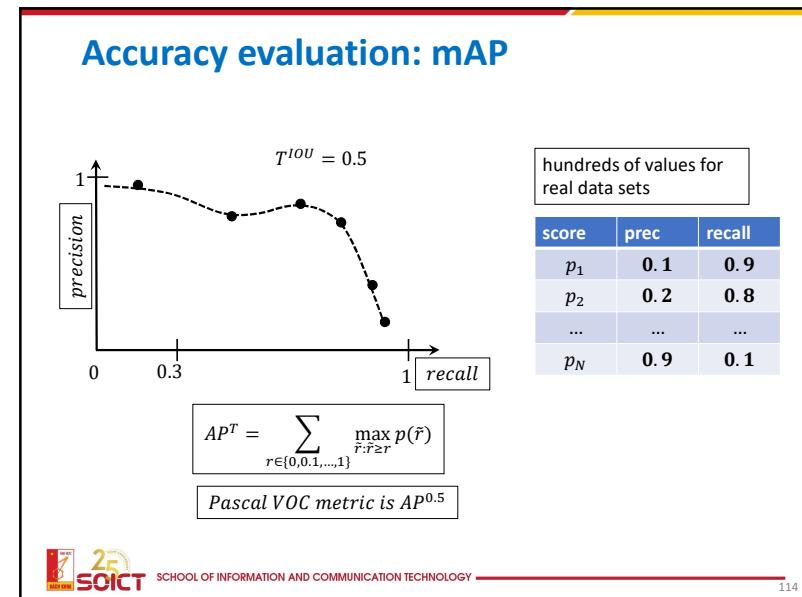


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

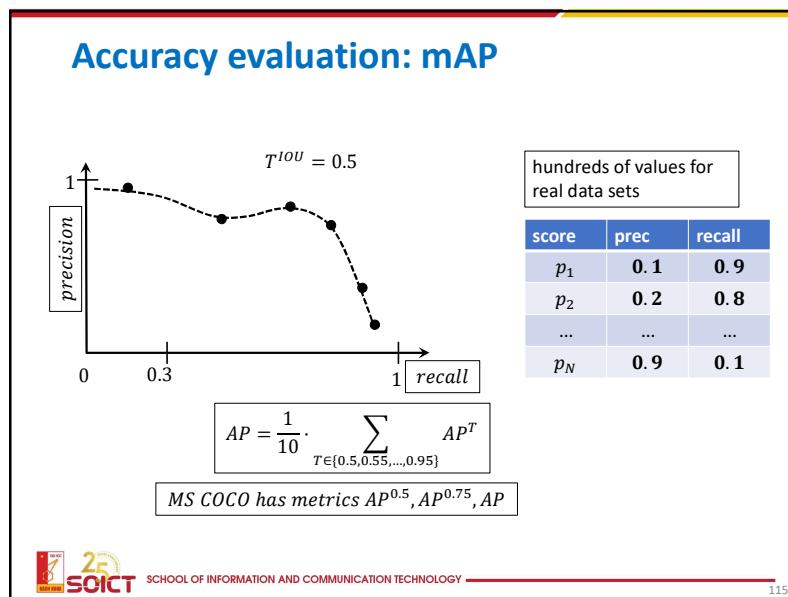
112



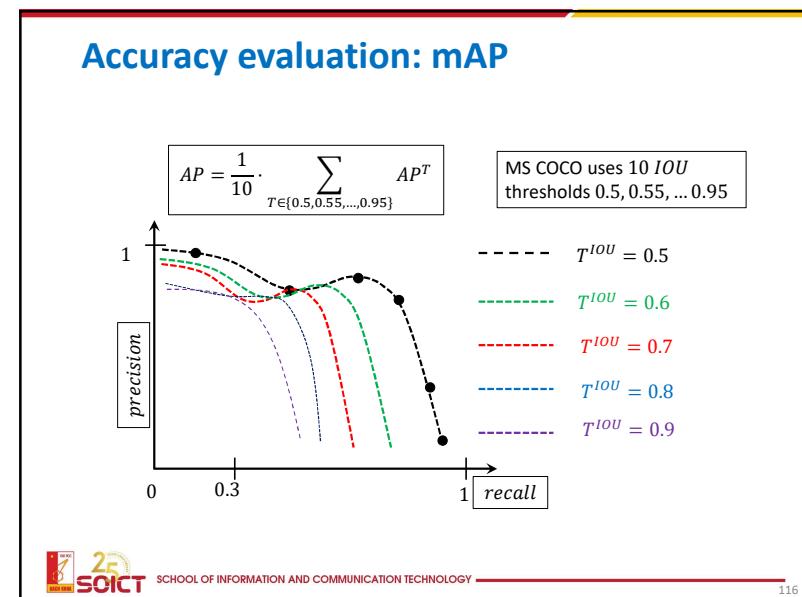
113



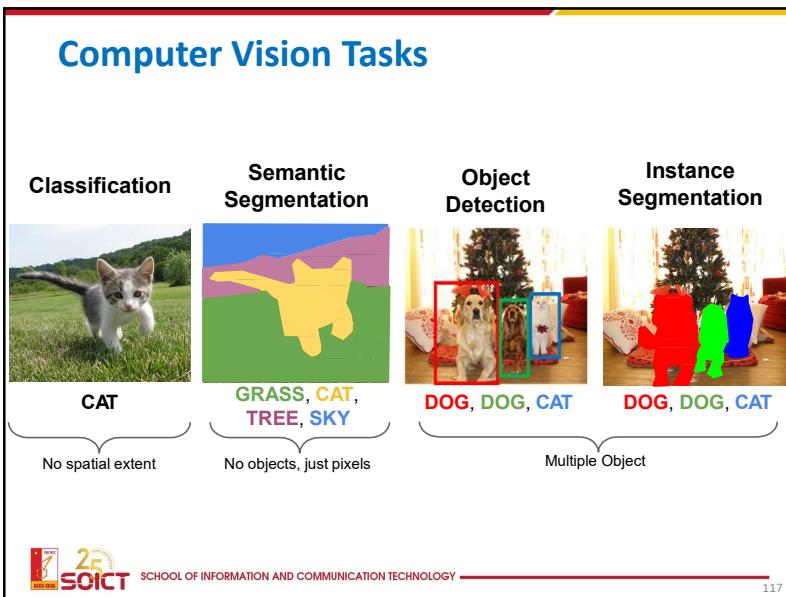
114



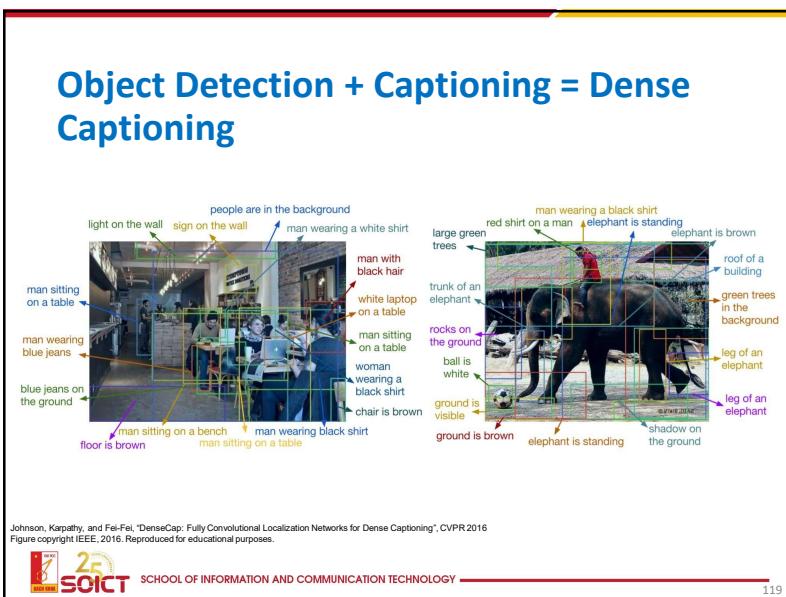
115



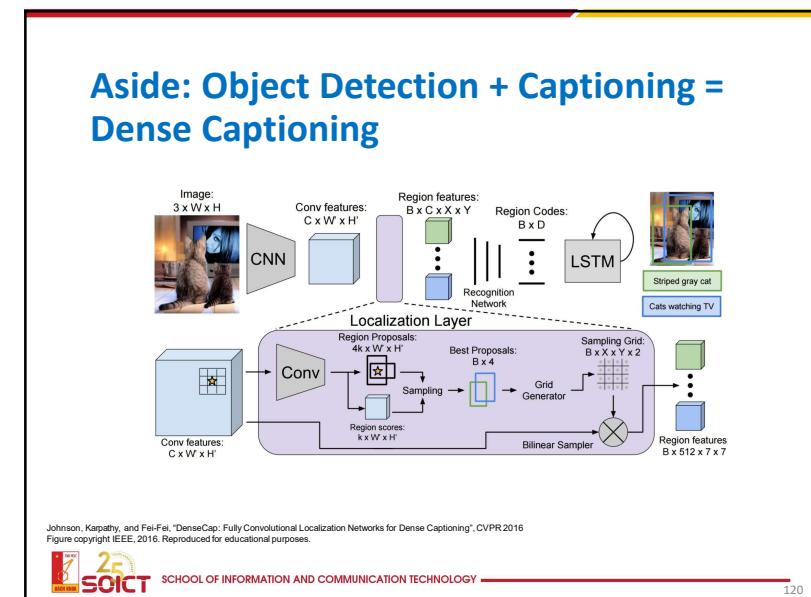
116



117



119



120

## 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.

**SOICT** 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.

**SOICT** 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!

**SOICT** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

123

## 3D Object Detection: Simple Camera Model

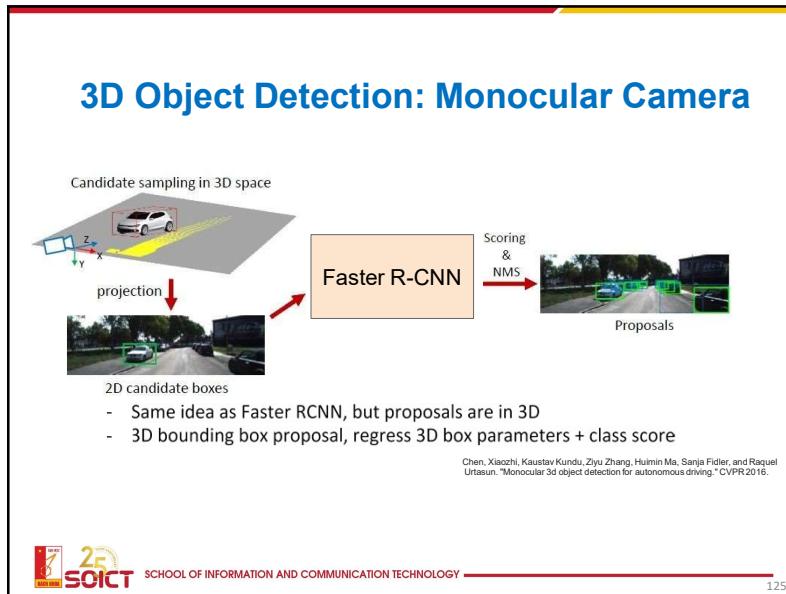
A point on the image plane corresponds to a **ray** in the 3D space

A 2D bounding box on an image is a **frustum** in the 3D space

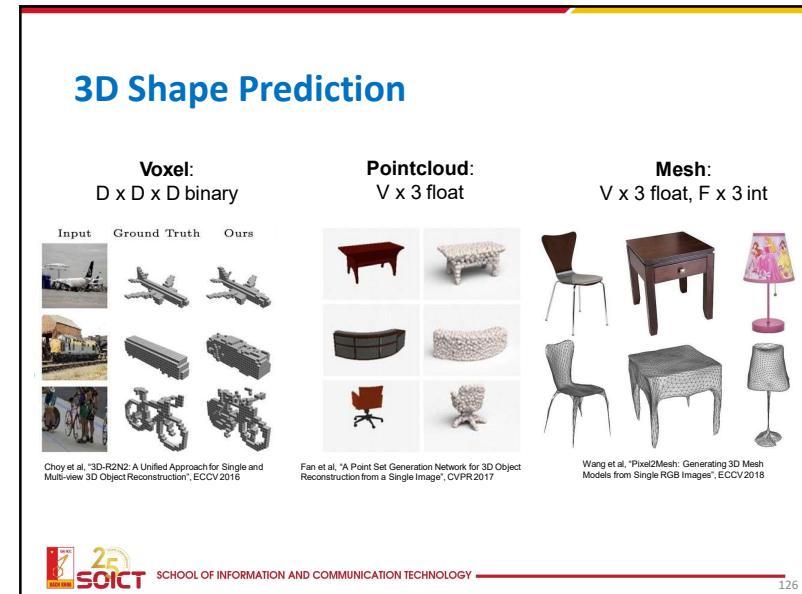
Localize an object in 3D:  
The object can be anywhere in the **camera viewing frustum**!

**SOICT** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

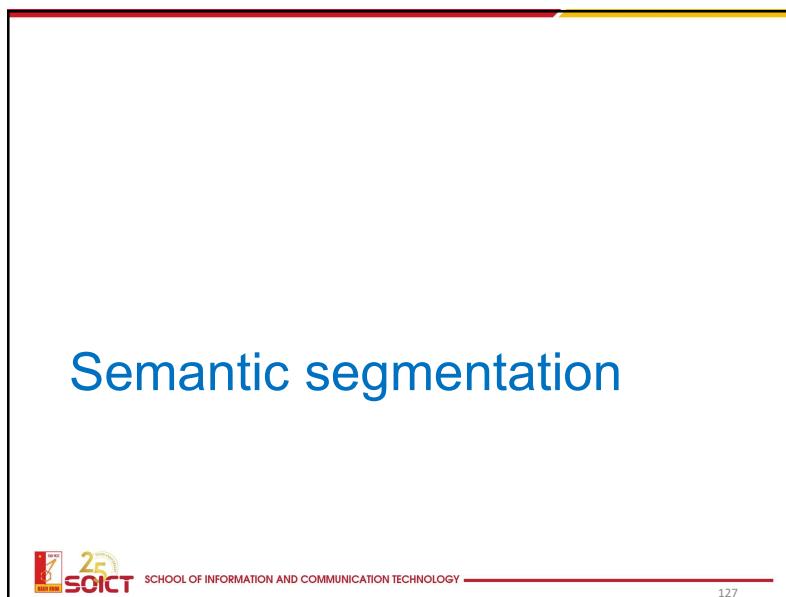
124



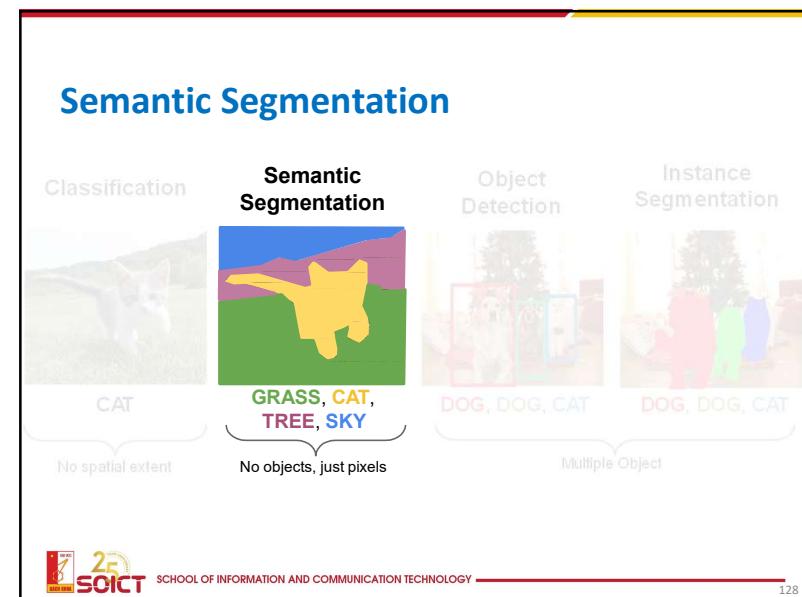
125



126



127



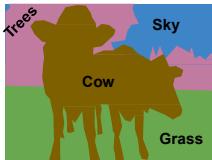
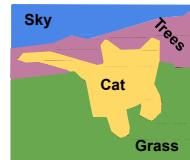
128

## Semantic Segmentation

Label each pixel in the image with a category label



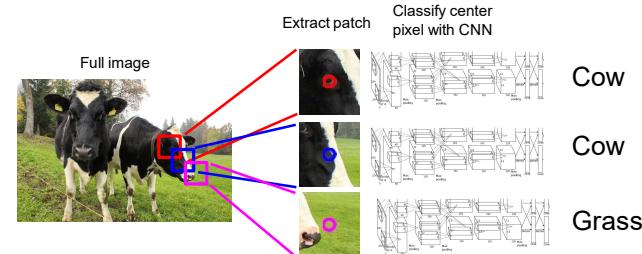
Don't differentiate instances, only care about pixels



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

129

## Semantic Segmentation Idea: Sliding Window



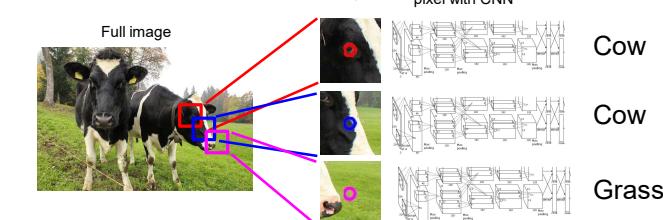
Farabet et al., "Learning Hierarchical Features for Scene Labeling", TPAMI 2013  
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

130

## Semantic Segmentation Idea: Sliding Window



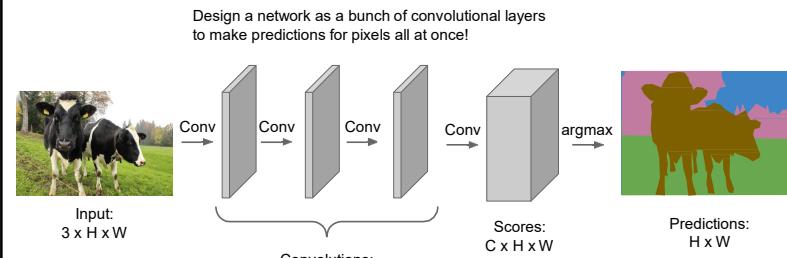
Farabet et al., "Learning Hierarchical Features for Scene Labeling.", TPAMI 2013  
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

131

## Semantic Segmentation Idea: Fully Convolutional

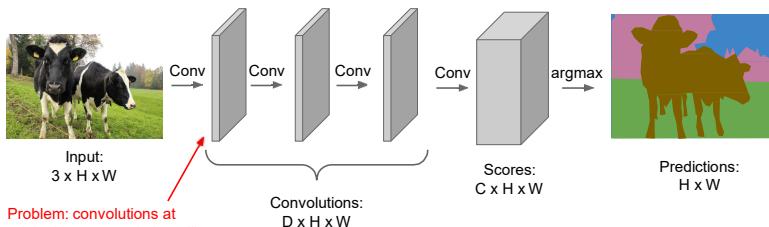


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

132

## Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

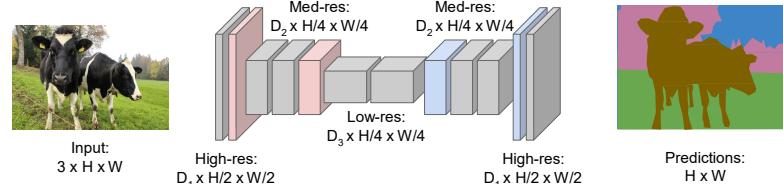


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

133

## Semantic Segmentation Idea: Fully Convolutional

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



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015  
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015



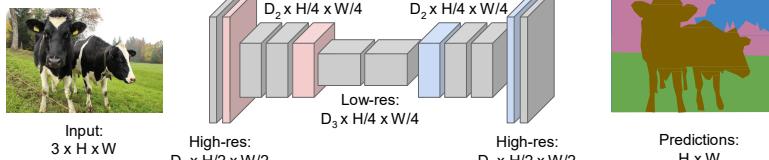
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

134

## Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided convolution

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



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015  
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

135

## In-Network upsampling: "Unpooling"

Nearest Neighbor

1	2
1	1
3	4
3	3

Input:  $2 \times 2$

1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Output:  $4 \times 4$

"Bed of Nails"

1	2
0	0
3	0
0	0

Input:  $2 \times 2$

1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

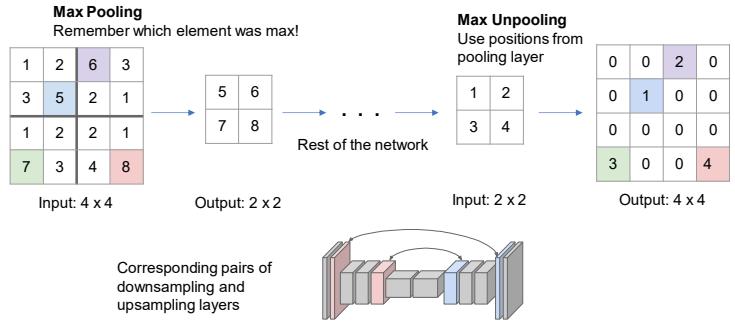
Output:  $4 \times 4$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

136

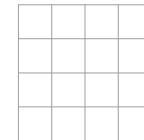
## In-Network upsampling: “Max Unpooling”



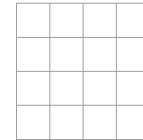
137

## Learnable Upsampling: Transpose Convolution

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



Input: 4 x 4



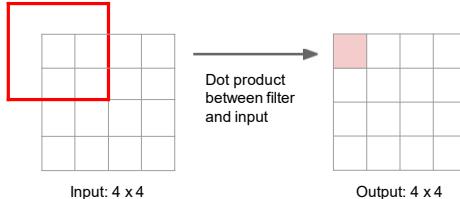
Output: 4 x 4



138

## Learnable Upsampling: Transpose Convolution

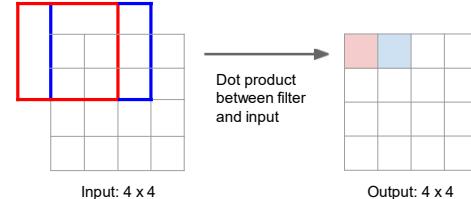
**Recall:** Normal 3 x 3 convolution, stride 1 pad 1



139

## Learnable Upsampling: Transpose Convolution

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

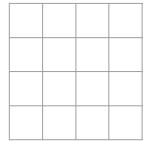


140

140

## Learnable Upsampling: Transpose Convolution

Recall: Normal  $3 \times 3$  convolution, stride 2 pad 1



Input:  $4 \times 4$



Output:  $2 \times 2$

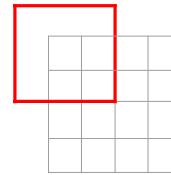


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

141

## Learnable Upsampling: Transpose Convolution

Recall: Normal  $3 \times 3$  convolution, stride 2 pad 1



Input:  $4 \times 4$

Dot product  
between filter  
and input



Output:  $2 \times 2$

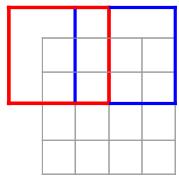


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

142

## Learnable Upsampling: Transpose Convolution

Recall: Normal  $3 \times 3$  convolution, stride 2 pad 1



Input:  $4 \times 4$

Dot product  
between filter  
and input



Output:  $2 \times 2$

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

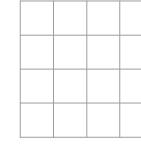
143

## Learnable Upsampling: Transpose Convolution

$3 \times 3$  transpose convolution, stride 2 pad 1



Input:  $2 \times 2$



Output:  $4 \times 4$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

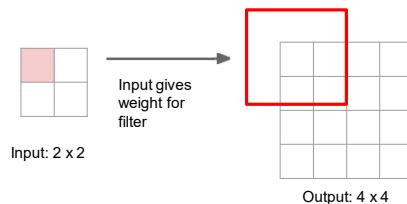
144

143

144

## Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

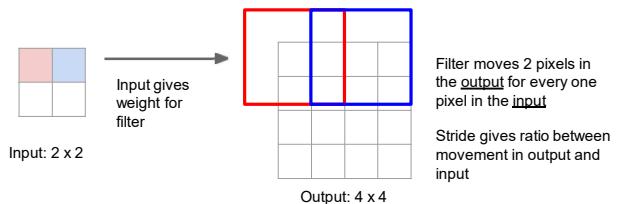


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

145

## Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

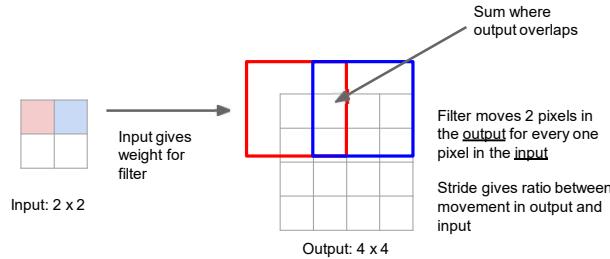


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

146

## Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1



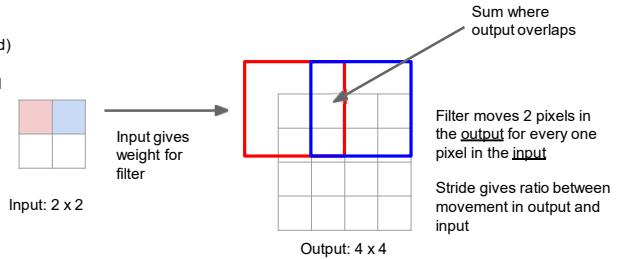
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

147

## Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

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



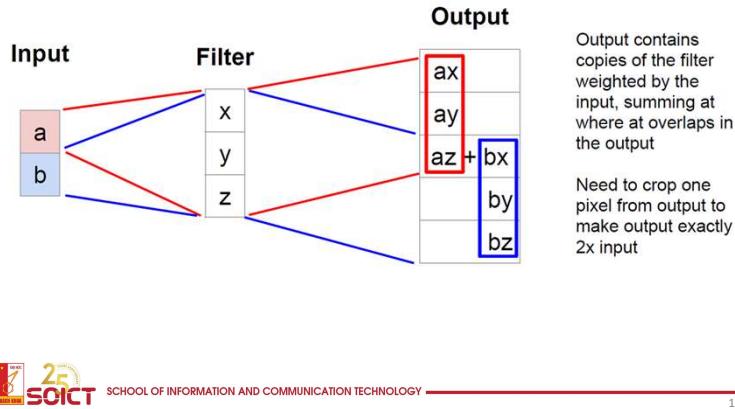
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

148

147

148

## Learnable Upsampling: Transpose Convolution



149

## Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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



150

## Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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

151

## Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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



152

151

152

## Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

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

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

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

Convolution transpose multiplies by the transpose of the same matrix:

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

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

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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

153

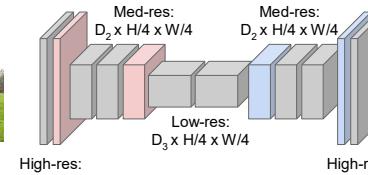
## Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided convolution



Input:  
 $3 \times H \times W$

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



**Upsampling:**  
Unpooling or strided transpose convolution



Predictions:  
 $H \times W$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

154

## 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	<b>TP</b> True Positive	<b>FP</b> False Positive
	negatives	<b>FN</b> False Negative	<b>TN</b> True Negative



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

157

## 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

158

## Instance segmentation



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

159

## Instance Segmentation

Classification



CAT

Semantic Segmentation



GRASS, CAT,  
TREE, SKY

Object Detection



DOG, DOG, CAT

Instance Segmentation



DOG, DOG, CAT

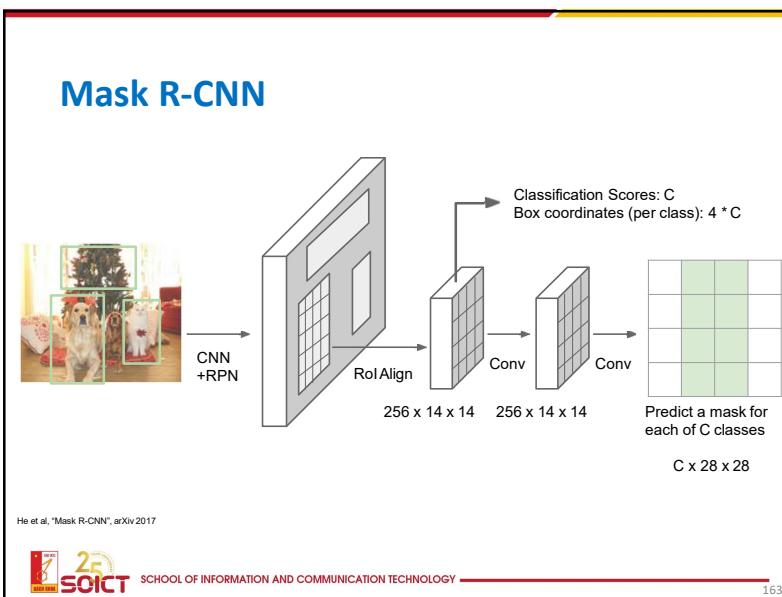
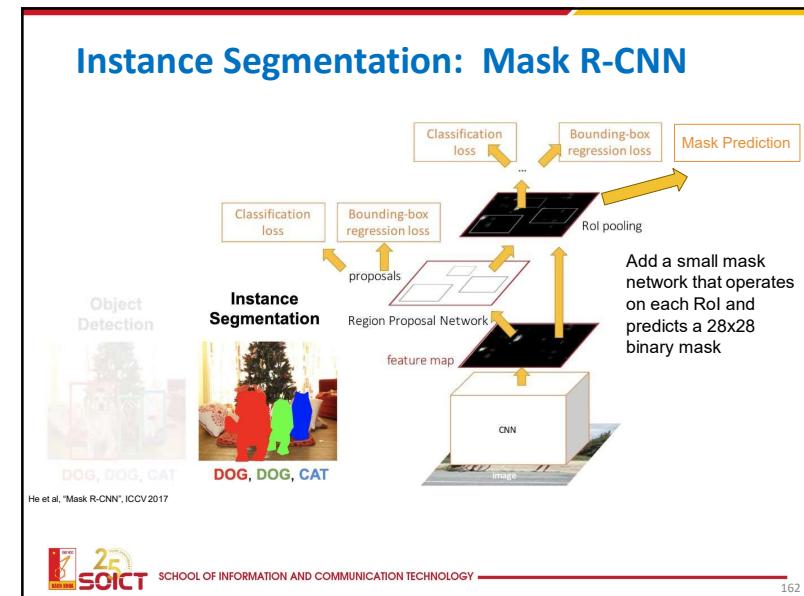
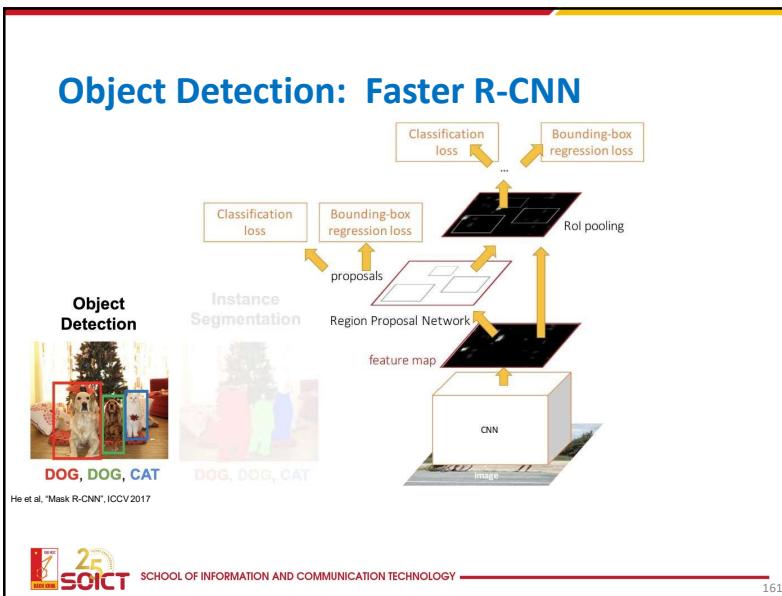
Multiple Object



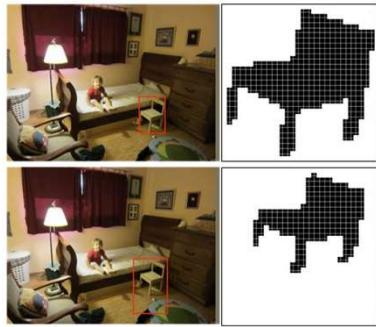
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

160

160



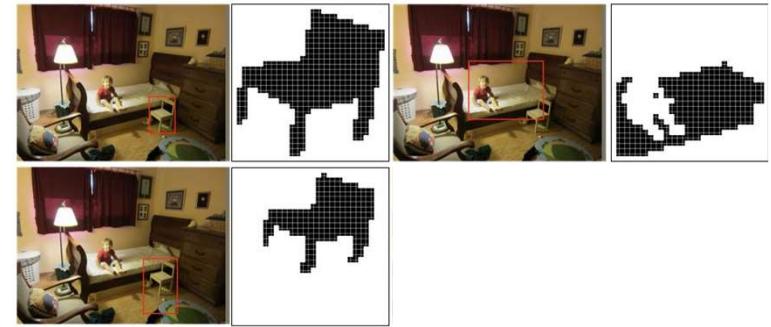
## Mask R-CNN: Example Mask Training Targets



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

165

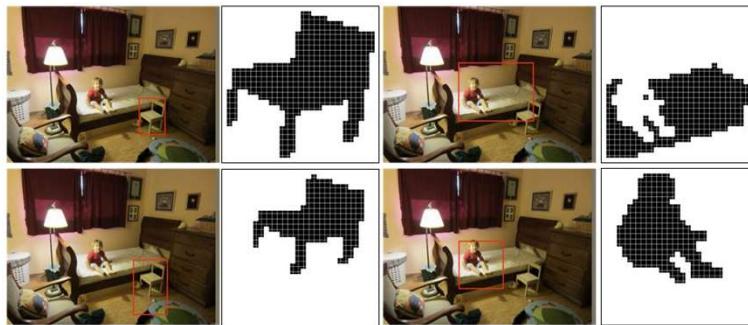
## Mask R-CNN: Example Mask Training Targets



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

166

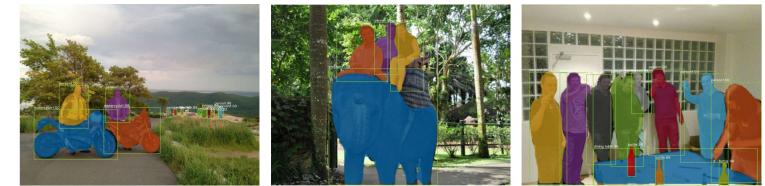
## Mask R-CNN: Example Mask Training Targets



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

167

## Mask R-CNN: Very Good Results!



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



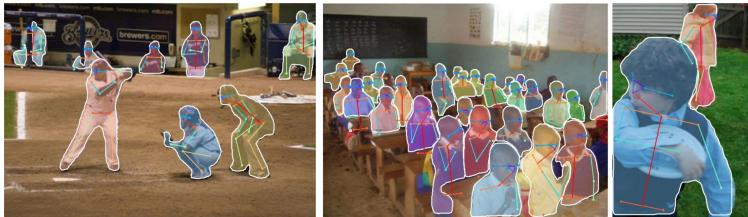
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

168

167

168

## Mask R-CNN Also does pose



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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

169

## Open Source Frameworks

- Lots of good implementations on GitHub!
- TensorFlow Detection API:  
  - [https://github.com/tensorflow/models/tree/master/research/object\\_detection](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

170

## 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

171

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG  
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Thank you!

[soict.hust.edu.vn/](http://soict.hust.edu.vn/) [fb.com/groups/soict](https://fb.com/groups/soict)



172