

# Detection and Segmentation

## Semantic Segmentation (semantic segmentation)

- input an image and then output a decision of a category for every pixel in the image
- don't differentiate instances, only care about pixels
- IDEA: Sliding Window — window DCT classification



**Fully Convolutional**

— (계산 비용 too expensive) Don't do it  
classification decision per pixel  
put 'cross entropy loss' on every pixel

ground truth category label fit loss 계산  
**downsampling & upsampling**

inside the network  $\Rightarrow$  깊지만 효율적 계산

Pooling, Strided Convolution...

"Unpooling"

1. "Nearest Neighbor"

1	2		
3	4		
		1	1
		3	3

→

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

2. "Bed of Nails"

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

3. Max Unpooling

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

→

5	6
7	8

→ ... →

1	2
3	4

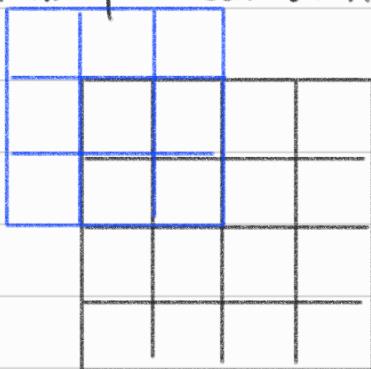
→

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

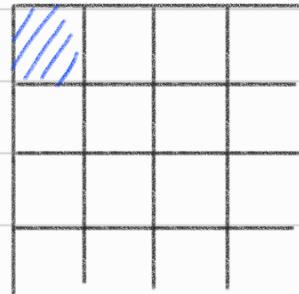
<MAX POOLING>

\* Corresponding pairs of downsampling and upsampling layers

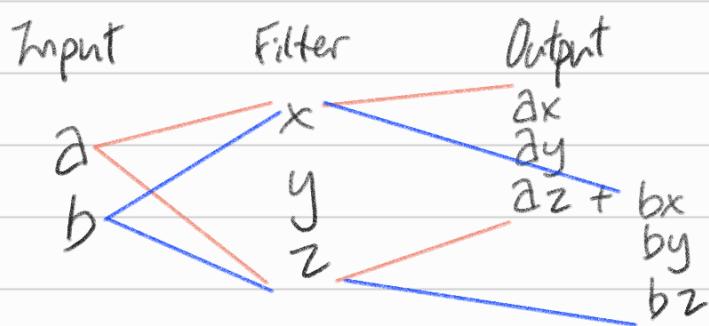
#### 4. Transpose Convolution



Dot product btw. filter and input



Ex.



3, output은 가로로 3개의 결과를 갖고 있고 3개의 행을 갖고 있다

Convolution = Matrix multiplication

Classification + Localization  $\Rightarrow$  Image in the category

have two loss function ① Softmax for ground truth category label

② L2 loss for bounding box

어려움

Human Pose Estimation 14m parts

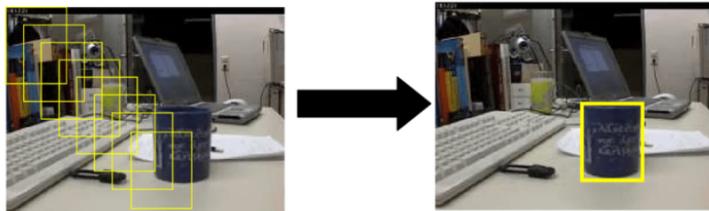


left foot L2 loss  
right foot L2 loss  $\Sigma \Rightarrow$  Loss  
⋮ ⋮

## \* Object Detection (物体検出)

image obj object 개수도 종류도 다르니까 challenging

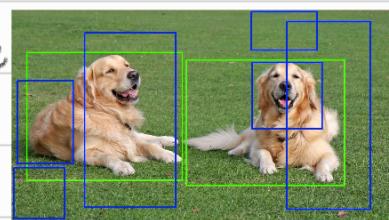
### Sliding Window



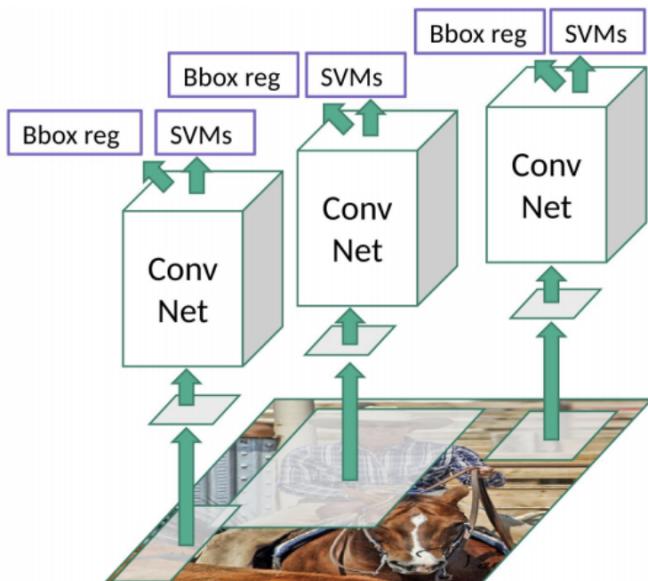
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

problem: need to apply CNN to huge # of locations and scales,  
Very computationally EXPENSIVE

Region Proposal - find blobby image regions that are likely to contain objects (not DL)



### R-CNN (Region Proposal + CNN 결합)



4. Classify regions with SVMs

3. Forward each region through convnet

2. CNN을 같은 크기 input > feature map warp

1.  
} ROI (Regions of Interest) from a proposed method

Problem: 여전히 비효율적인 연산 (too expensive), memory 부족, too slow

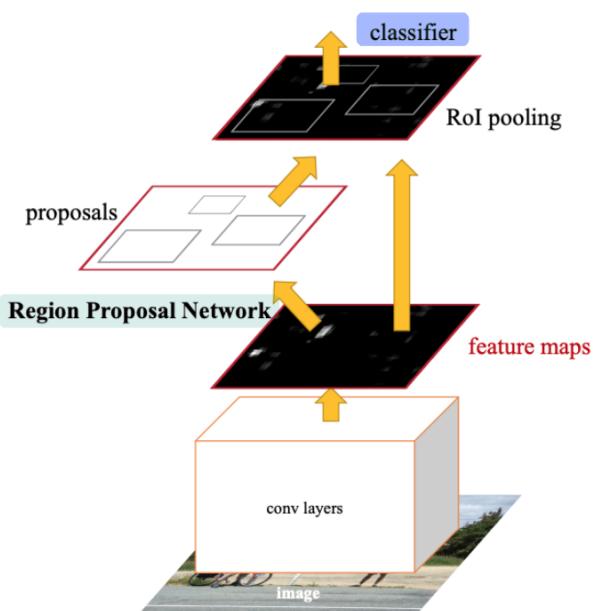
### Fast R-CNN

1. CNN layers로 feature map으로 변환
- + 2. ROI from proposal method

### 3. ROI Pooling layers

4. Softmax classifier , Bounding-box regressors

## Faster R-CNN



Jointly train 4 losses:

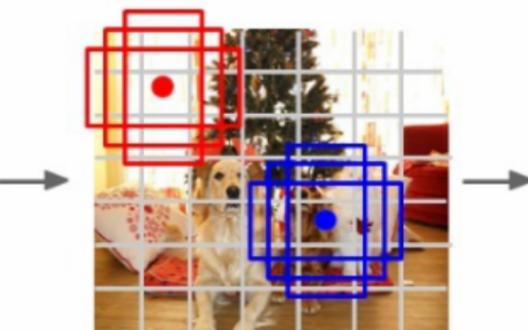
1. RPN classify object / not object (binary)
2. RPN regress box coordinates
3. Final classification score
4. Final box coordinates

→ represents entire high resolution img.

## YOLO (you only look once)



Input image  
 $3 \times H \times W$



Divide image into grid  
 $7 \times 7$

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)

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

Redmon et al, "You Only Look Once:  
Unified, Real-Time Object Detection", CVPR 2016  
Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016

ex.

Divide input image into coarse grid  $(7 \times 7)$

Imagine some set of bounding boxes (tall, square, wide)

## Dense Captioning

## Instance Segmentation

predict the locations and identities of objects in that image,  
predict a whole segmentation mask for each of objects and  
predict which pixels in the input image corresponds to  
each object instance (hybrid of semantic segmentation &  
object detection)

### \* Mask R-CNN

