

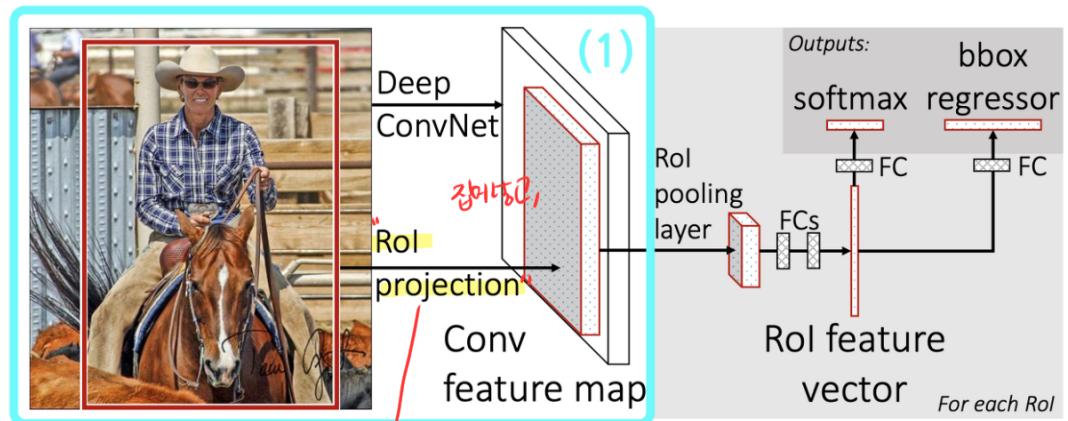
R-CNN : CNN 사용 초기의 Object detection 방식

- * 단점
 1. AlexNet 사용 224×224 깊이 warping \Rightarrow 시간↑
 2. Selective Search 통해 2000개 Image proposal all input to CNN \Rightarrow time↑
 3. " " + SVM이 GPU 사용 짜증 주고 X
 4. 뒤 부족 수반 Computation share X (No back propagation)

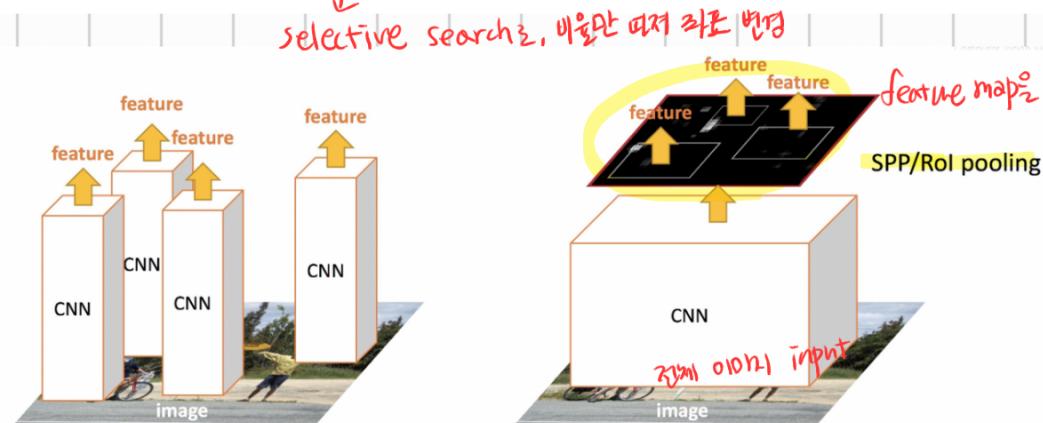


Fast R-CNN

1. CNN model 바탕



- RoI: region of interest



R-CNN

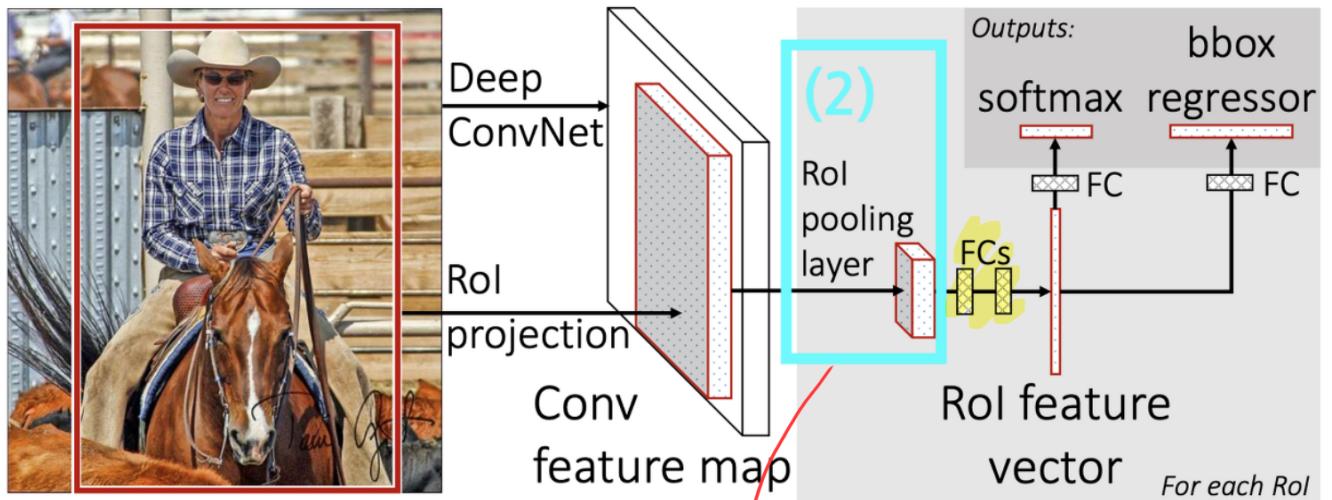
- Extract image regions
- 1 CNN per region (2000 CNNs)
- Classify region-based features
- 높아짐 대상은 crop

SPP-net & Fast R-CNN (the same forward pipeline)

- 1 CNN on the entire image
- Extract features from feature map regions
- Classify region-based features
- Crop \times 256x101x121 \Rightarrow CNN model이 짜증을 주는 feature map의 RoI projection

(주제)

2. RoI Pooling



• Resolution을 고려한 맵의 위치 위상 RoI pooling 흐름

↳ 각각의 feature map의 위치에 Max pooling

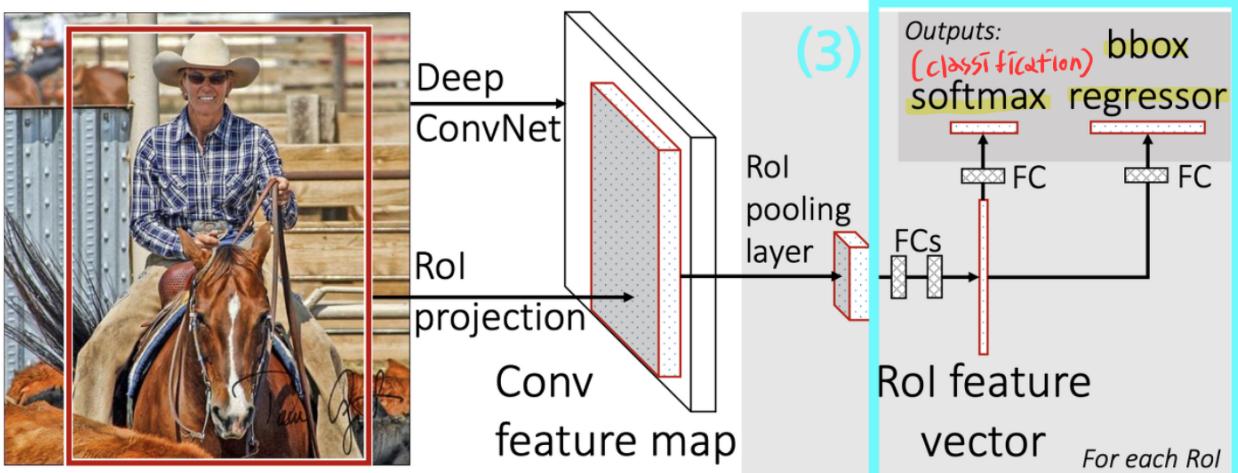
Ex) 7×7 output feature map = fixed length Feature Vector $1 \times n$

$$21 \times 1 \rightarrow 21/3 = 7, 14/2 = 7 \quad \text{stride}(3,2)$$

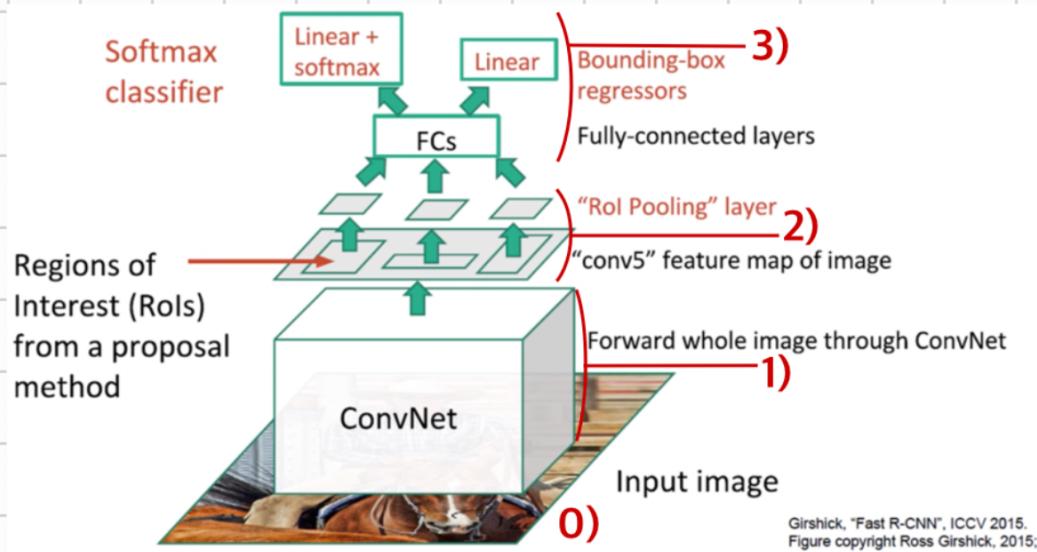
$$42 \times 35 \rightarrow 42/6 = 7, 35/5 = 7 \quad " (6,5)$$

\Rightarrow 각각의 feature map의 위치에 Max pooling

3. Classification & Bounding Box Regression



* Fast R-CNN 2014



- 0) Input image on Selective Search 진행, ROI 영역 미리 뽑고
- 1) CNN 모델에 Image形式 Input, Feature map Output
- 2) 뽑아낸 ROI 영역을 Feature Map에 적용 후, ROI pooling하여 fixed feature vector output
- 3) FC를 거쳐 두 가지 Layer³ class prediction, Bounding Box Regression 진행

- Loss function

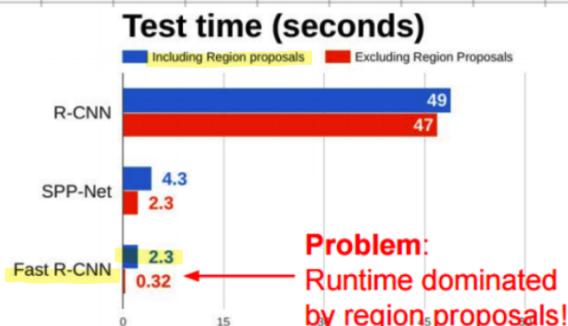
$$L(p, u, t^u, v) = L_{\text{cls}}(p, u) + \lambda[u \geq 1]L_{\text{loc}}(t^u, v)$$

Classification Bounding box Regression

$p = (p_0, \dots, p_K)$: 예측된 Class Score
in $K + 1$ categories (0은 background)
 u = 실제 Class Score
 $t^u = (t_x^u, t_y^u, t_w^u, t_h^u)$: 예측된 tuple
 (v_x, v_y, v_w, v_h) : 실제 Bounding box 좌표 값
 λ 가 0(background)일때 0, 나머지일때 1

Smooth L1 Loss

<문1>



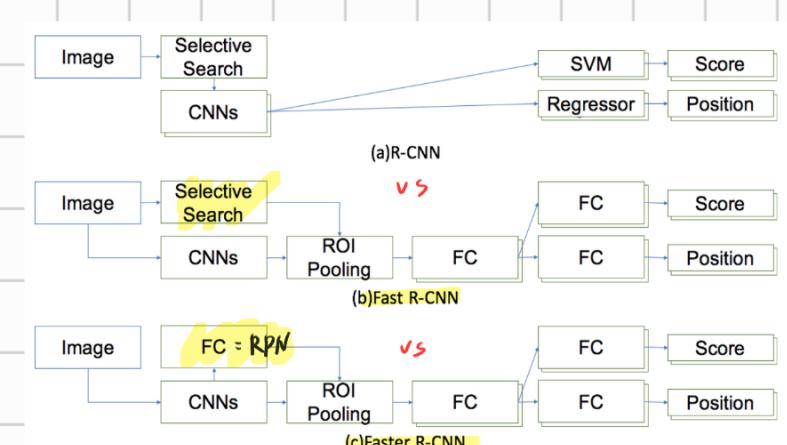
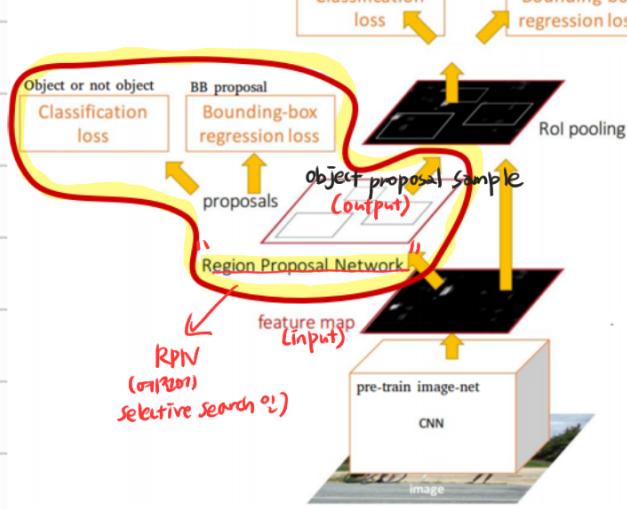
- Region Proposal에 사용되는 selective search가 GPU와 CPU 사용 일고려 때문
→ Faster R-CNN으로 바뀜

Not real-time
object detector

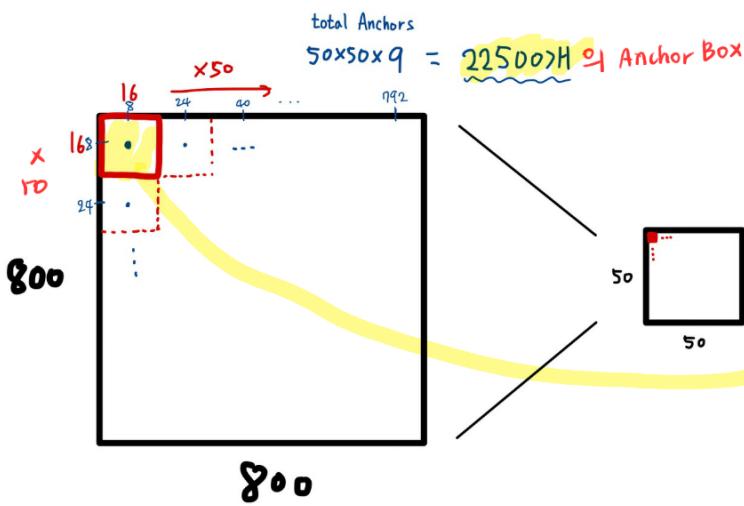
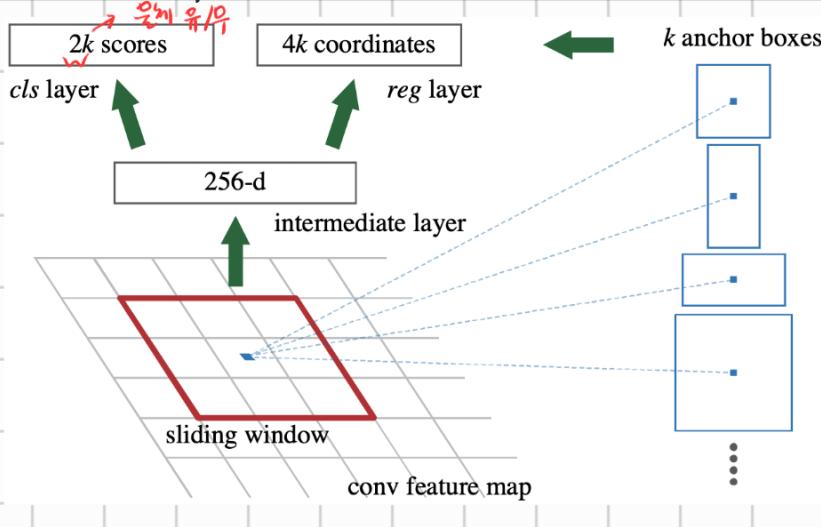
*Faster R-CNN

- Deep Network w/o Region Proposal \Rightarrow RPN (Region Proposal Networks)

Faster R-CNN



1. RPN (Region proposal Network)

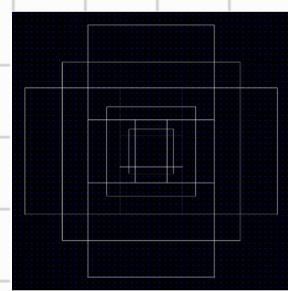
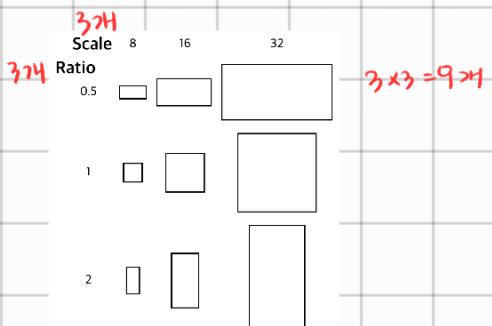


VGG-16의 외각(backbone) 3개

1) Anchor targeting

$800 \times 800 \times 3 \rightarrow \text{CNN} (\text{VGG-16}) \rightarrow 50 \times 50 \times 512$ feature map (Input)

\rightarrow Sliding window $\Rightarrow K=9$ Anchor Box $\times 16^2$



9개 Anchor Box

- 22500개의 Anchor Box 간의 detection score! 위해 Labeling 필요
→ Ground Truth
- (IT Label) 22,500개의 Anchor & Ground Truth Box의 IoU를 모두 계산

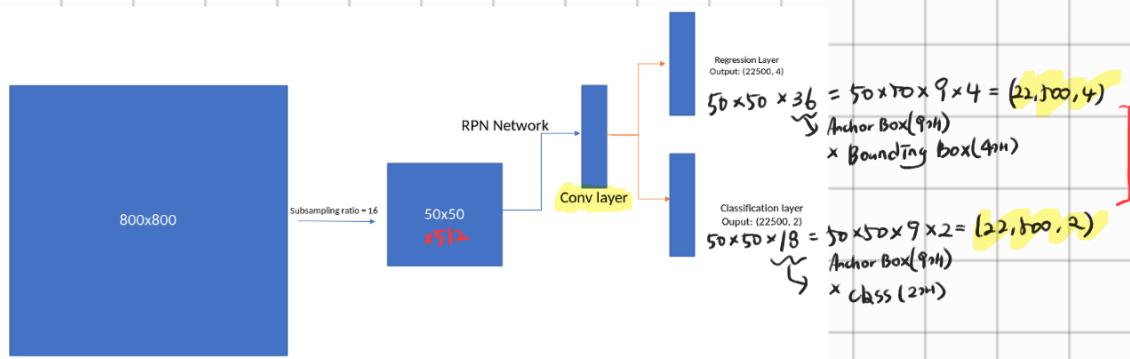
$$\text{IoU} = \frac{\text{Area of Overlap (교집합)}}{\text{Area of Union (합집합)}}$$

0.7 ↑ = 1 = positive

other = -1 true X

0.3 ↓ = 0 = negative

2. prediction



1) prediction에 GT Label 포함 Loss function을 통해 RPN 학습

&
→ P:N = 1:3의 Sampling

2) prediction 같은 NMS 및 ROI Sampling 후 Fast-R-CNN을 쓰임 like selective search

Non-maximum Suppression
: Detection only 2nd best Bounding Box만

3. Loss Function for PRN

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*)$$

Log Loss GT Label (물체가 있으면 1, 없으면 0)
Predicted Probability of Anchor i Smooth L1 Loss
GT box

Smooth L1 Loss: 물체가 없으면 bbox는 고려하지 않음.

4. 성능

	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image (with proposals)	50 seconds	2 seconds	0.2 seconds
(Speedup)	1x	25x	250x
mAP (VOC 2007)	66.0	66.9	69.9