



BIG DATA & AI ANALYTICS
EXPERT COMPANY

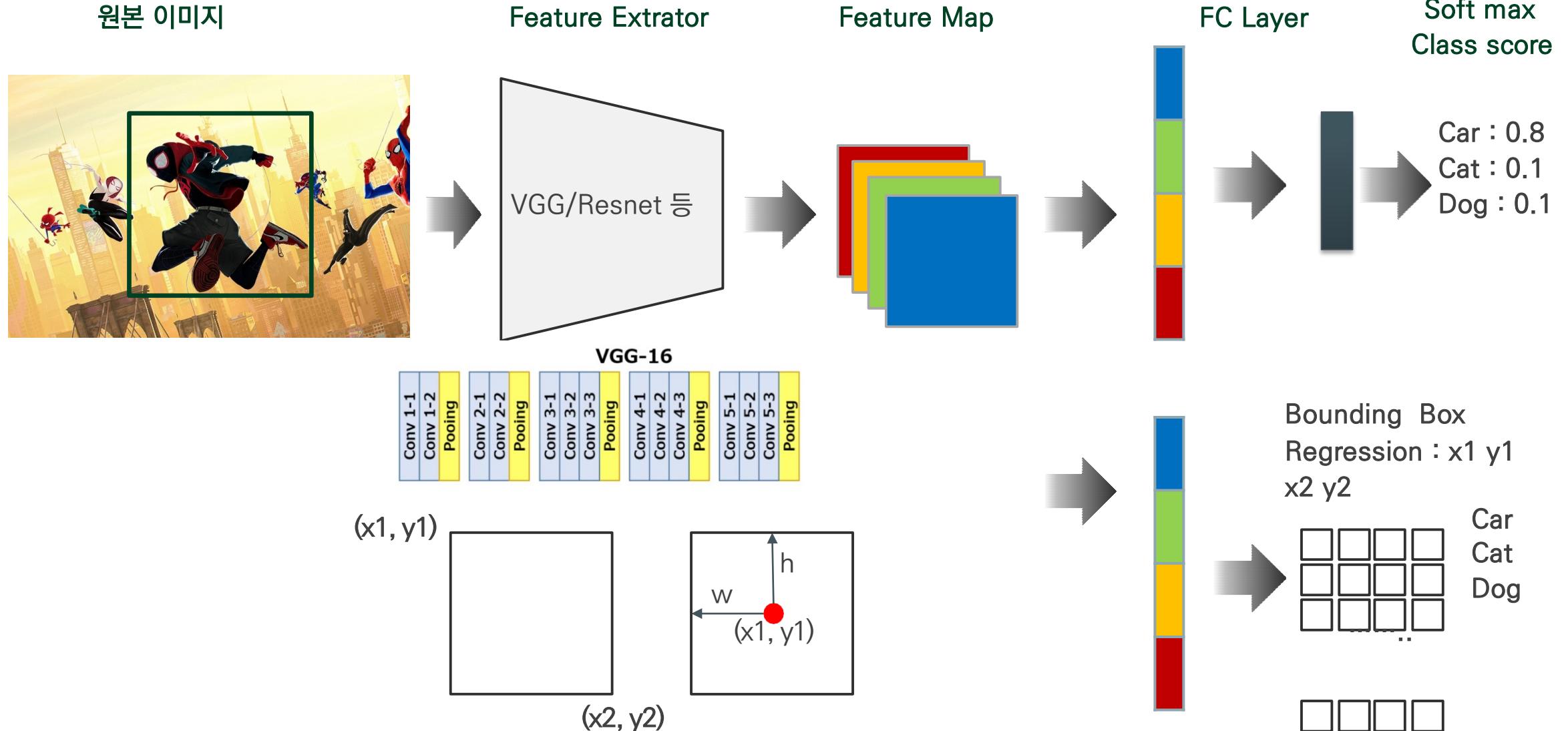
Two-Stage Detection
Faster RCNN



BIG DATA & AI ANALYTICS
EXPERT COMPANY

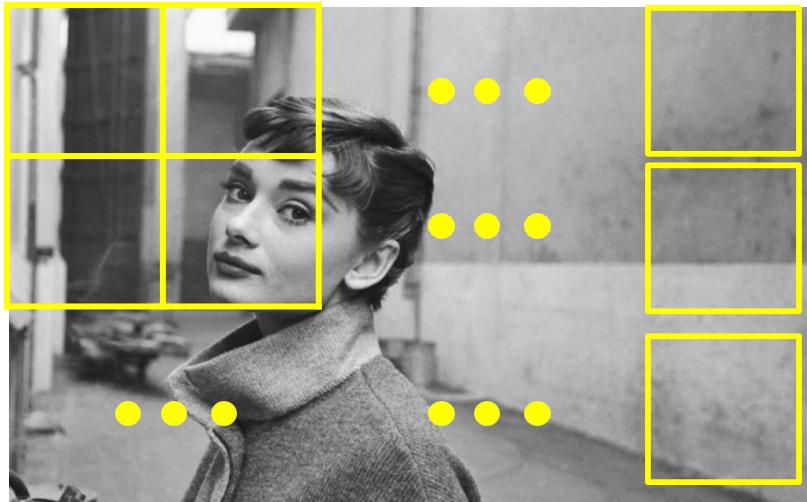
RCNN 개요

Object Detection – bbox Regression



| Sliding Window, Region Proposal 방식

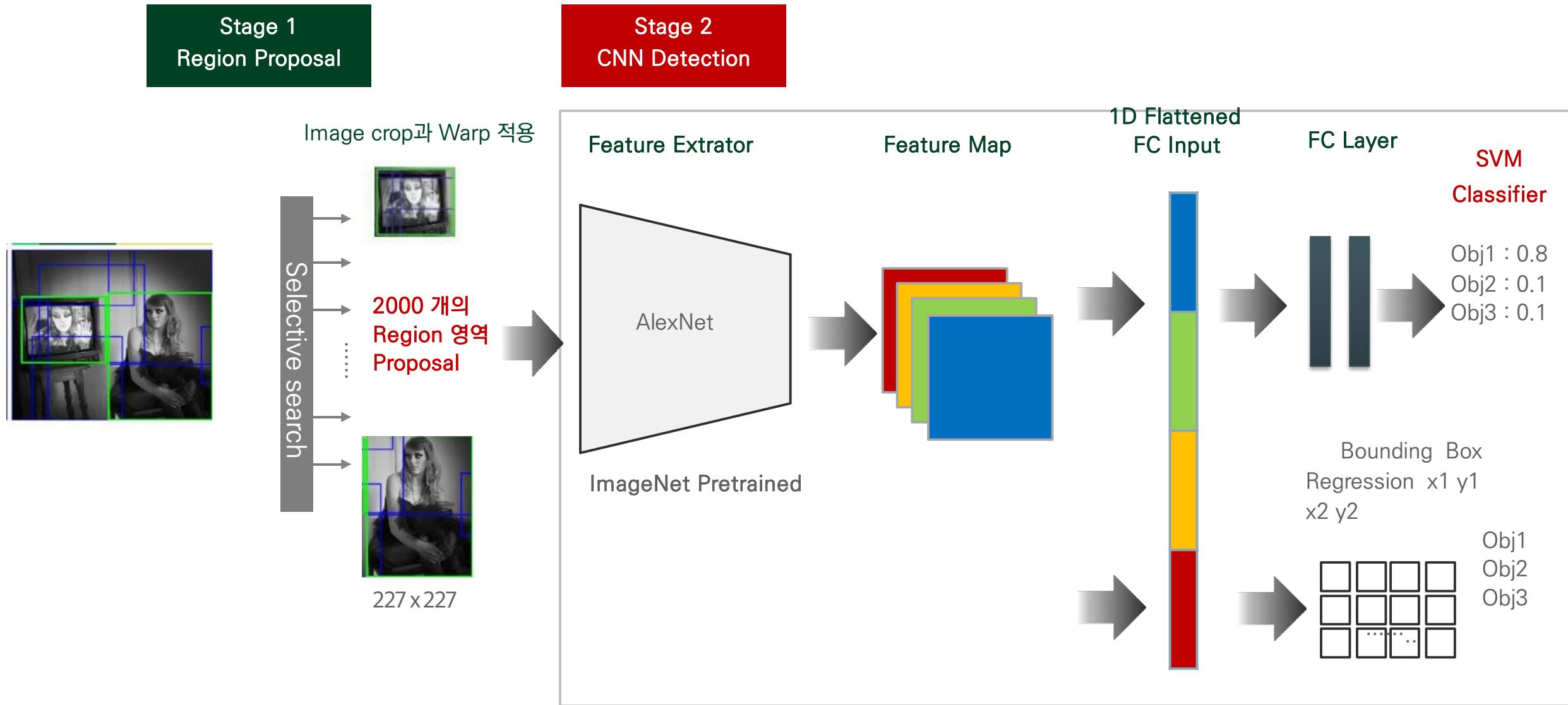
Sliding Window 방식



Region Proposal 방식

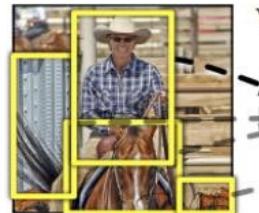


RCNN – Region Proposal 기반

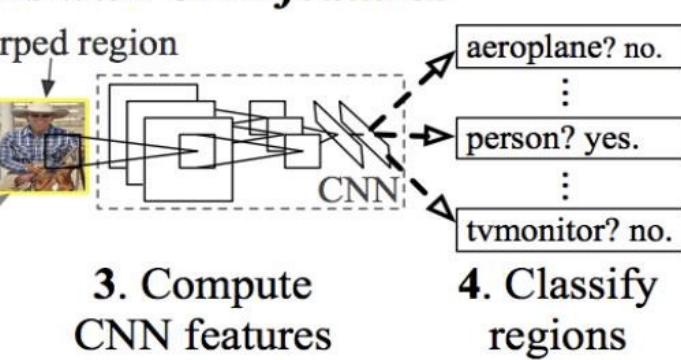




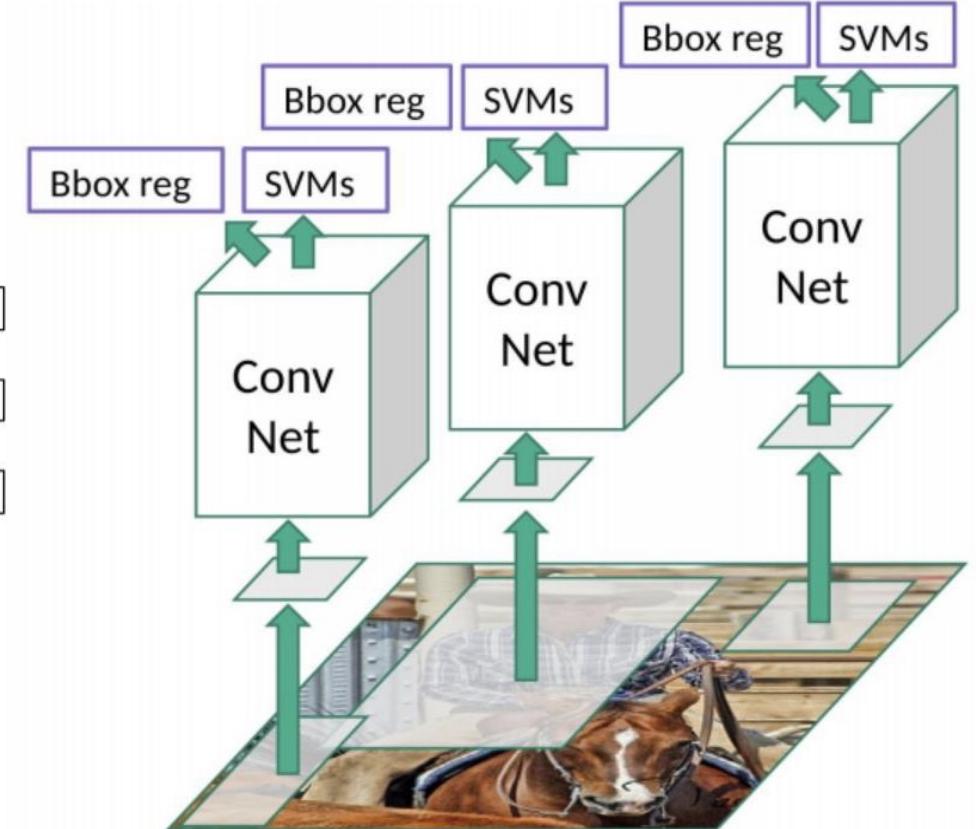
1. Input image



2. Extract region proposals (~2k)



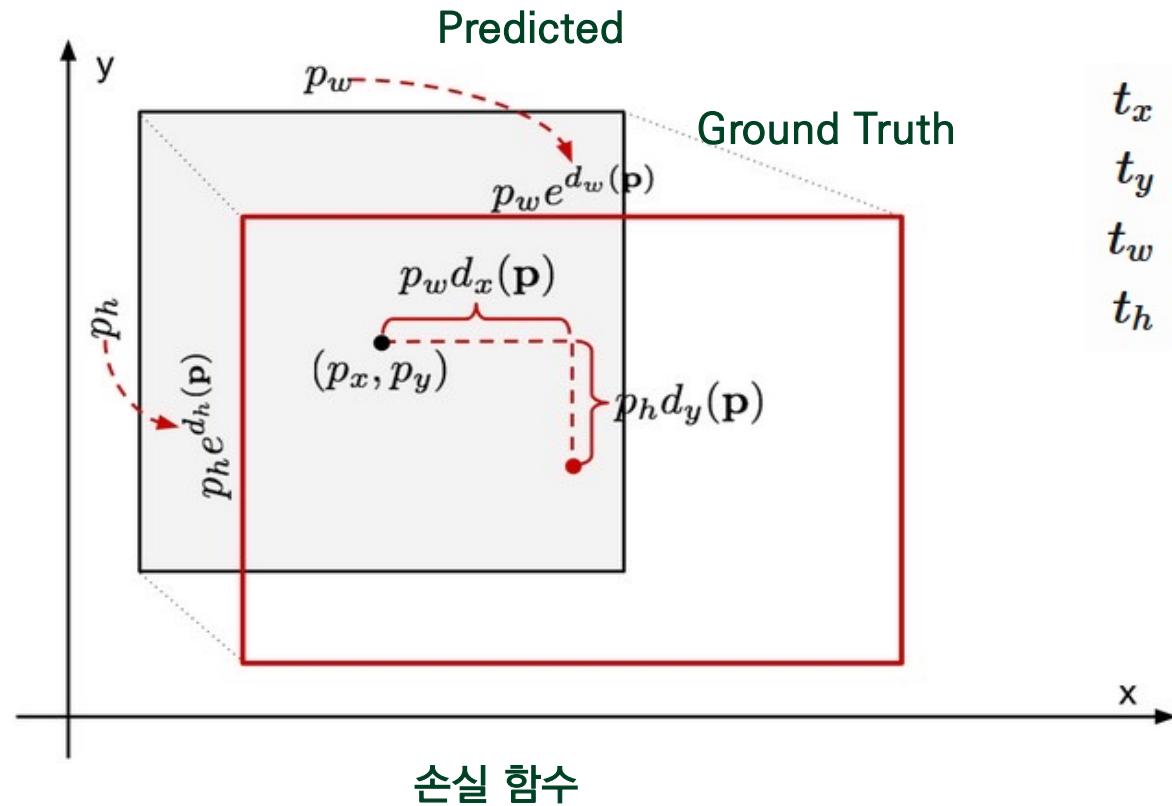
R-CNN



| Bounding Box Regression

수정 예측 값

$$\begin{aligned}\hat{g}_x &= p_w d_x(\mathbf{p}) + p_x \\ \hat{g}_y &= p_h d_y(\mathbf{p}) + p_y \\ \hat{g}_w &= p_w \exp(d_w(\mathbf{p})) \\ \hat{g}_h &= p_h \exp(d_h(\mathbf{p}))\end{aligned}$$



Target

$$\begin{aligned}t_x &= (g_x - p_x)/p_w \\ t_y &= (g_y - p_y)/p_h \\ t_w &= \log(g_w/p_w) \\ t_h &= \log(g_h/p_h)\end{aligned}$$

$$\mathcal{L}_{\text{reg}} = \sum_{i \in \{x, y, w, h\}} (t_i - d_i(\mathbf{p}))^2 + \lambda \|\mathbf{w}\|^2$$

PASCAL VOC 2010 기준

| VOC 2010 test | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv | mAP |
|--------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| DPM v5 [20] [†] | 49.2 | 53.8 | 13.1 | 15.3 | 35.5 | 53.4 | 49.7 | 27.0 | 17.2 | 28.8 | 14.7 | 17.8 | 46.4 | 51.2 | 47.7 | 10.8 | 34.2 | 20.7 | 43.8 | 38.3 | 33.4 |
| UVA [39] | 56.2 | 42.4 | 15.3 | 12.6 | 21.8 | 49.3 | 36.8 | 46.1 | 12.9 | 32.1 | 30.0 | 36.5 | 43.5 | 52.9 | 32.9 | 15.3 | 41.1 | 31.8 | 47.0 | 44.8 | 35.1 |
| Regionlets [41] | 65.0 | 48.9 | 25.9 | 24.6 | 24.5 | 56.1 | 54.5 | 51.2 | 17.0 | 28.9 | 30.2 | 35.8 | 40.2 | 55.7 | 43.5 | 14.3 | 43.9 | 32.6 | 54.0 | 45.9 | 39.7 |
| SegDPM [18] [†] | 61.4 | 53.4 | 25.6 | 25.2 | 35.5 | 51.7 | 50.6 | 50.8 | 19.3 | 33.8 | 26.8 | 40.4 | 48.3 | 54.4 | 47.1 | 14.8 | 38.7 | 35.0 | 52.8 | 43.1 | 40.4 |
| R-CNN | 67.1 | 64.1 | 46.7 | 32.0 | 30.5 | 56.4 | 57.2 | 65.9 | 27.0 | 47.3 | 40.9 | 66.6 | 57.8 | 65.9 | 53.6 | 26.7 | 56.5 | 38.1 | 52.8 | 50.2 | 50.2 |
| R-CNN BB | 71.8 | 65.8 | 53.0 | 36.8 | 35.9 | 59.7 | 60.0 | 69.9 | 27.9 | 50.6 | 41.4 | 70.0 | 62.0 | 69.0 | 58.1 | 29.5 | 59.4 | 39.3 | 61.2 | 52.4 | 53.7 |

장점

- 당시에 높은 Detection 정확도
- 동시대의 다른 알고리즘 대비 매우 높은 Detection 정확도

단점

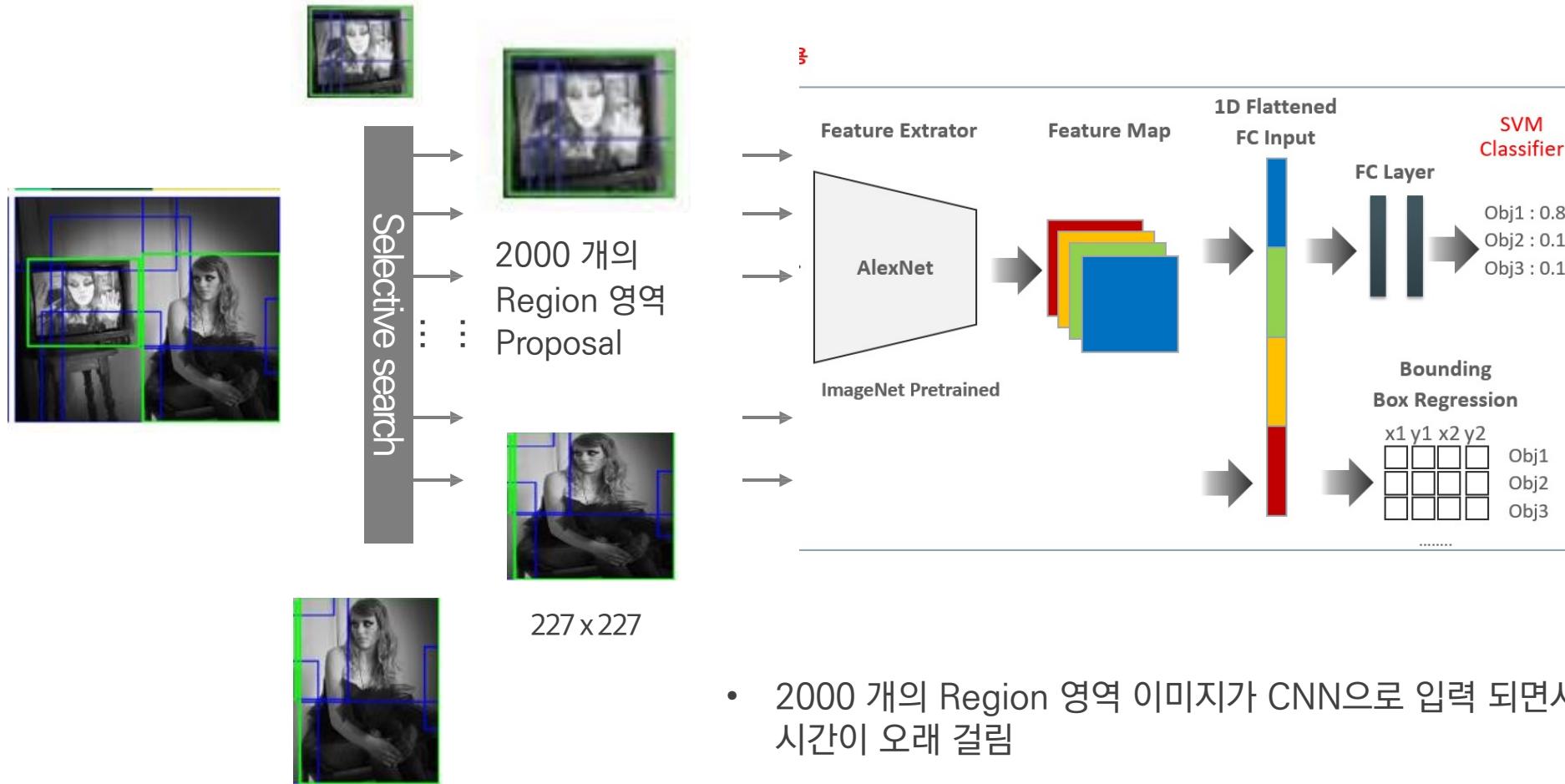
- 1장의 이미지를 Object Detection 하는데 약 50초 소요
- 너무 느린 Detection 시간과 복잡한 아키텍처 및 학습 프로세스
- 하나의 이미지에서 selective search를 수행하여 2000개의 region 영역을 도출
- 2000개씩 생성된 region 이미지를 CNN Feature map 생성
- 복잡한 구성 요소들: Selective search, CNN Feature Extractor, SVM과 Bounding box regressor로 구성



BIG DATA & AI ANALYTICS
EXPERT COMPANY

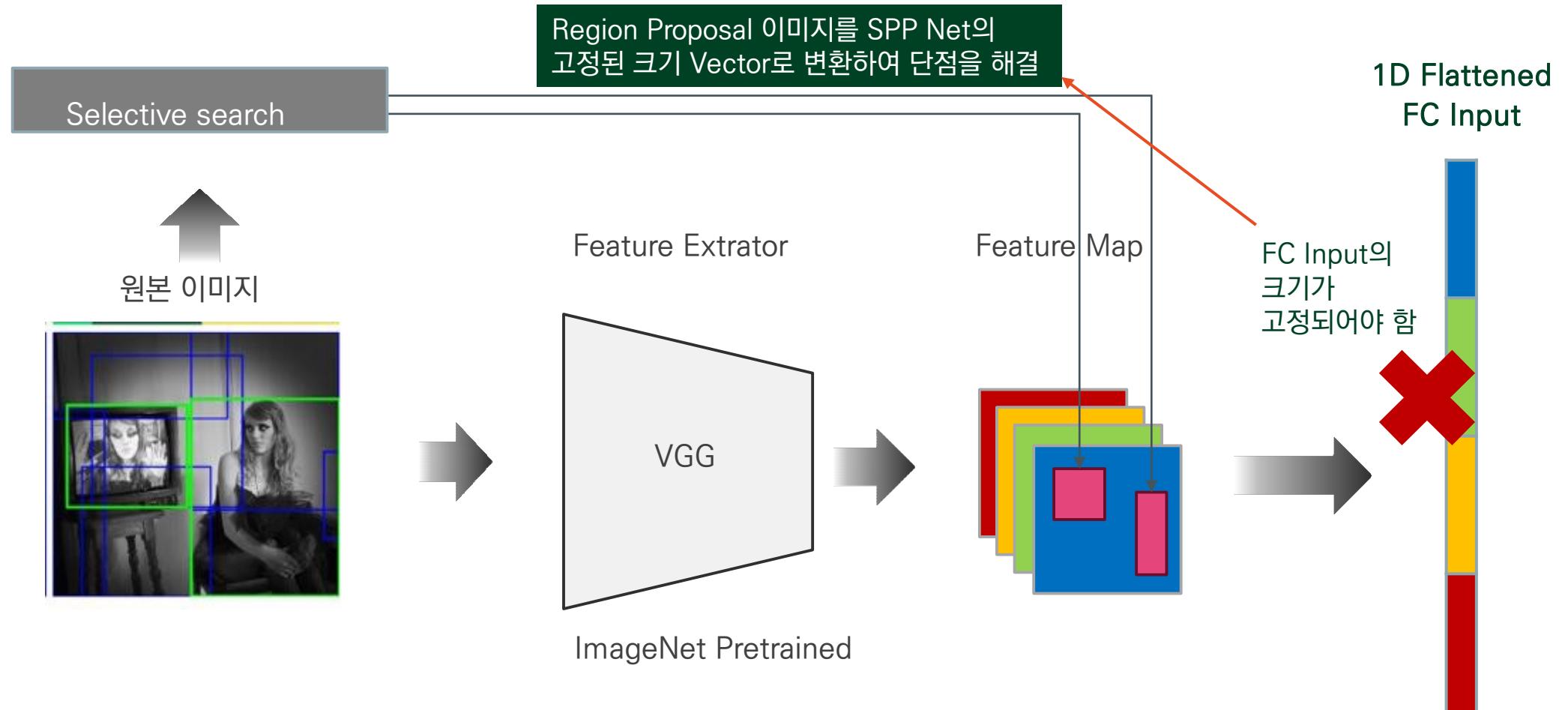
SPP Net
(Spatial Pyramid Pooling)

Image crop과 Warp 적용



- 2000 개의 Region 영역 이미지가 CNN으로 입력 되면서 Object Detection 시간이 오래 걸림
- Region 영역 이미지가 Crop/Warp 되어야 함.

- 2000개의 Region Proposal 이미지를 CNN으로 Feature Extraction 하지 않고 원본 이미지만 CNN으로 Feature Map 생성
- 원본 이미지의 Selective search로 추천된 영역의 이미지만 Feature Map으로 매핑하여 별도 추출



| SPP(Spatial pyramid Pooling)

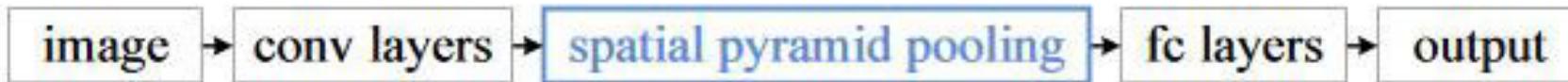
- SPP는 CNN상에서 Image classification에서 서로 다른 이미지의 크기를 고정된 크기로 변환하는 매칭 기법



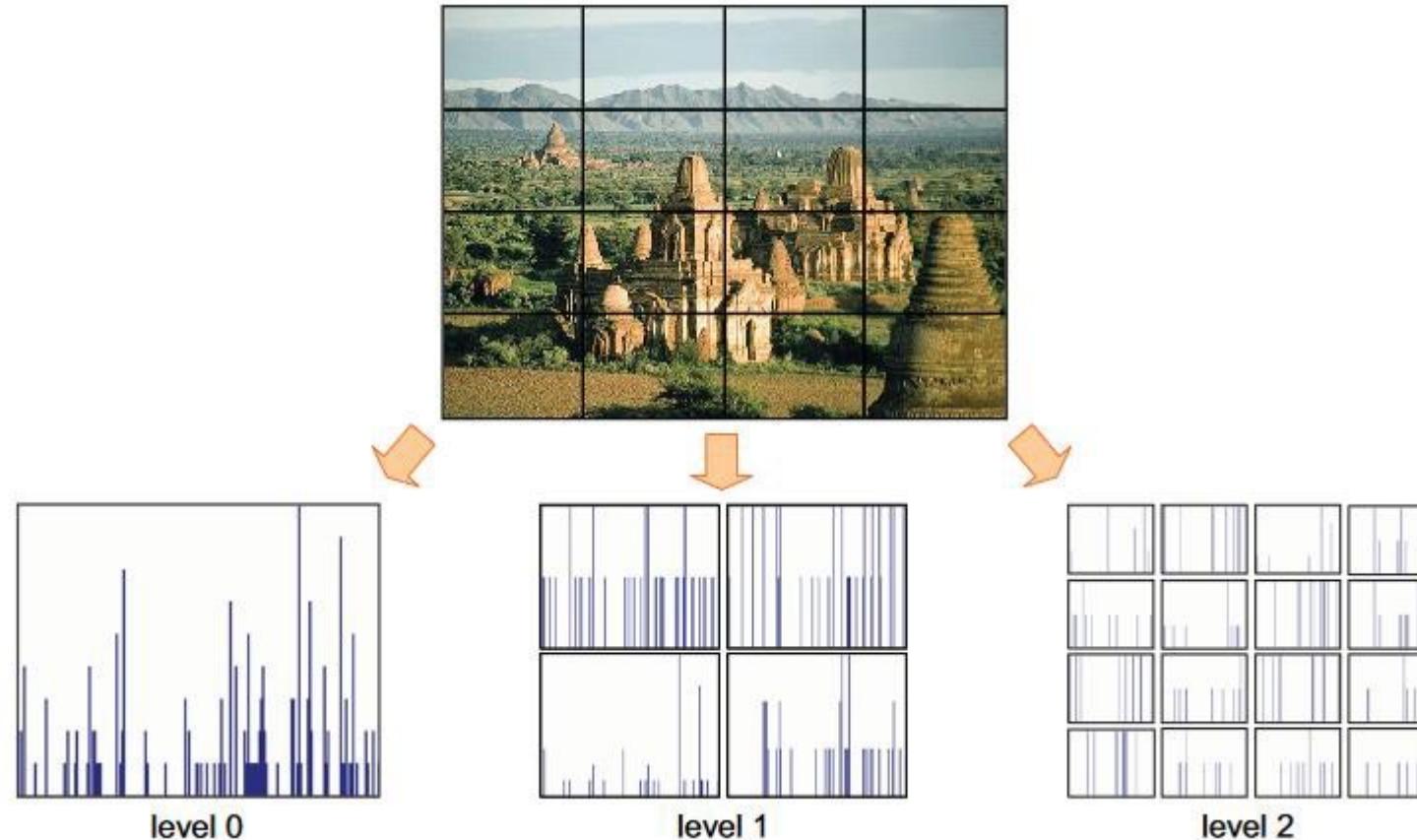
crop



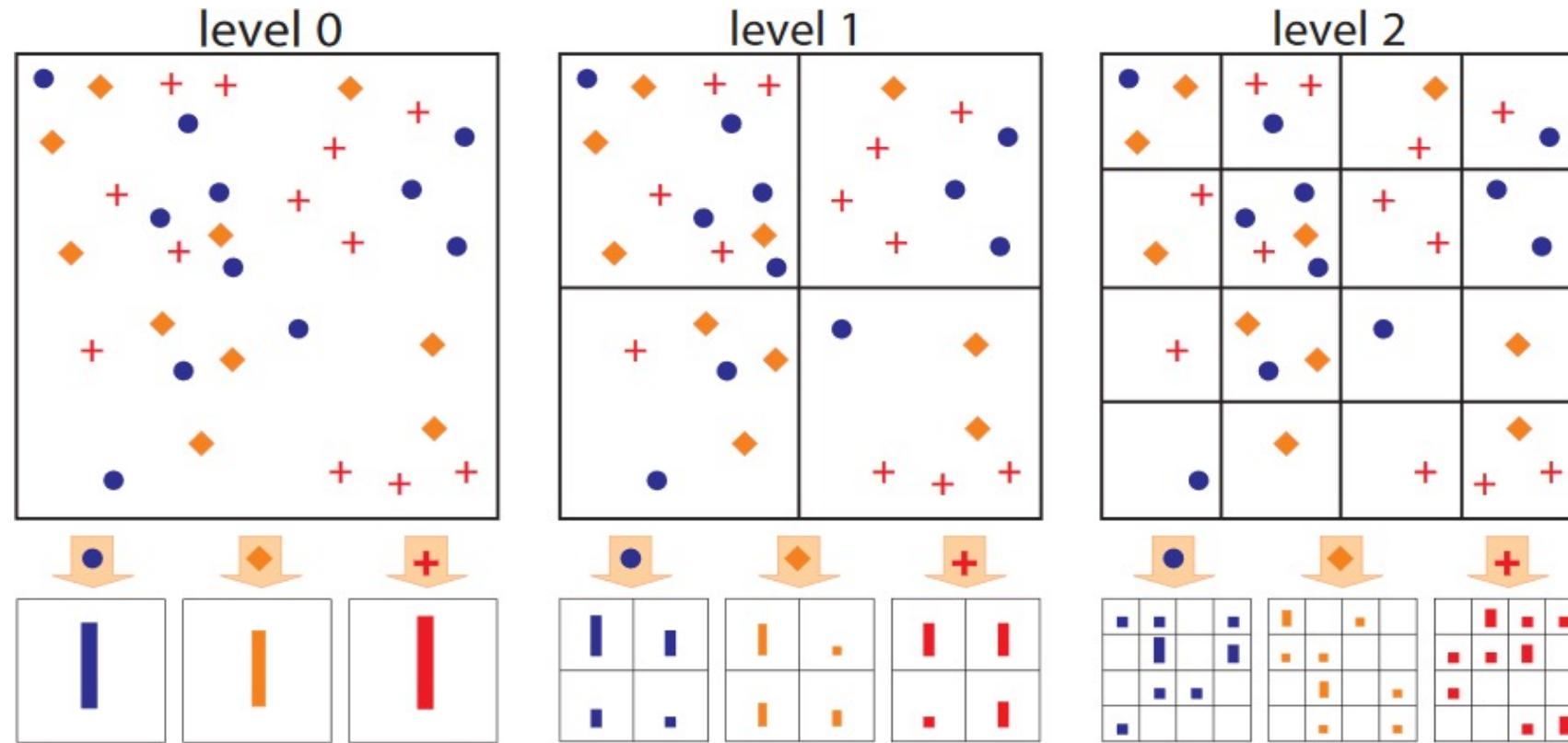
warp



| SPM(Spatial Pyramid Matching) 개요

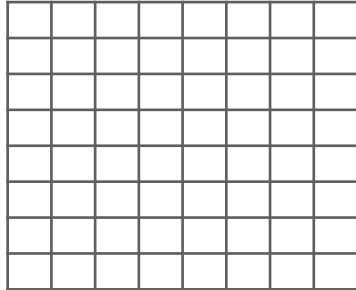


| SPM(Spatial Pyramid Matching) 개요

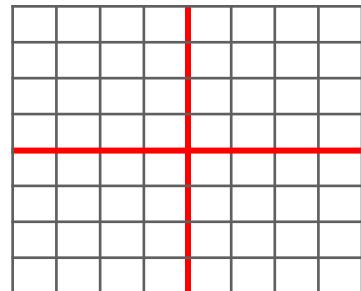


| SPM을 Feature Map을 균일한 Vector 크기로 표현

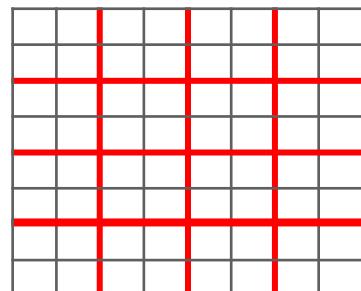
8 x 8 Feature Map



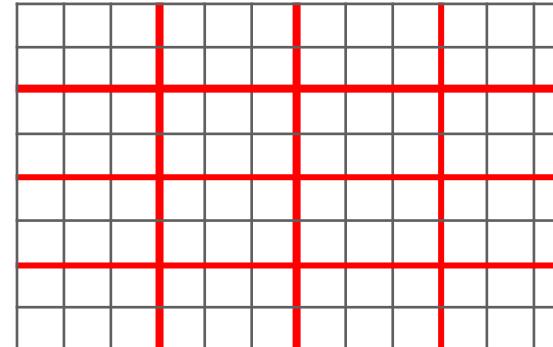
$$3 \times 1 = 3$$



$$3 \times 4 = 12$$



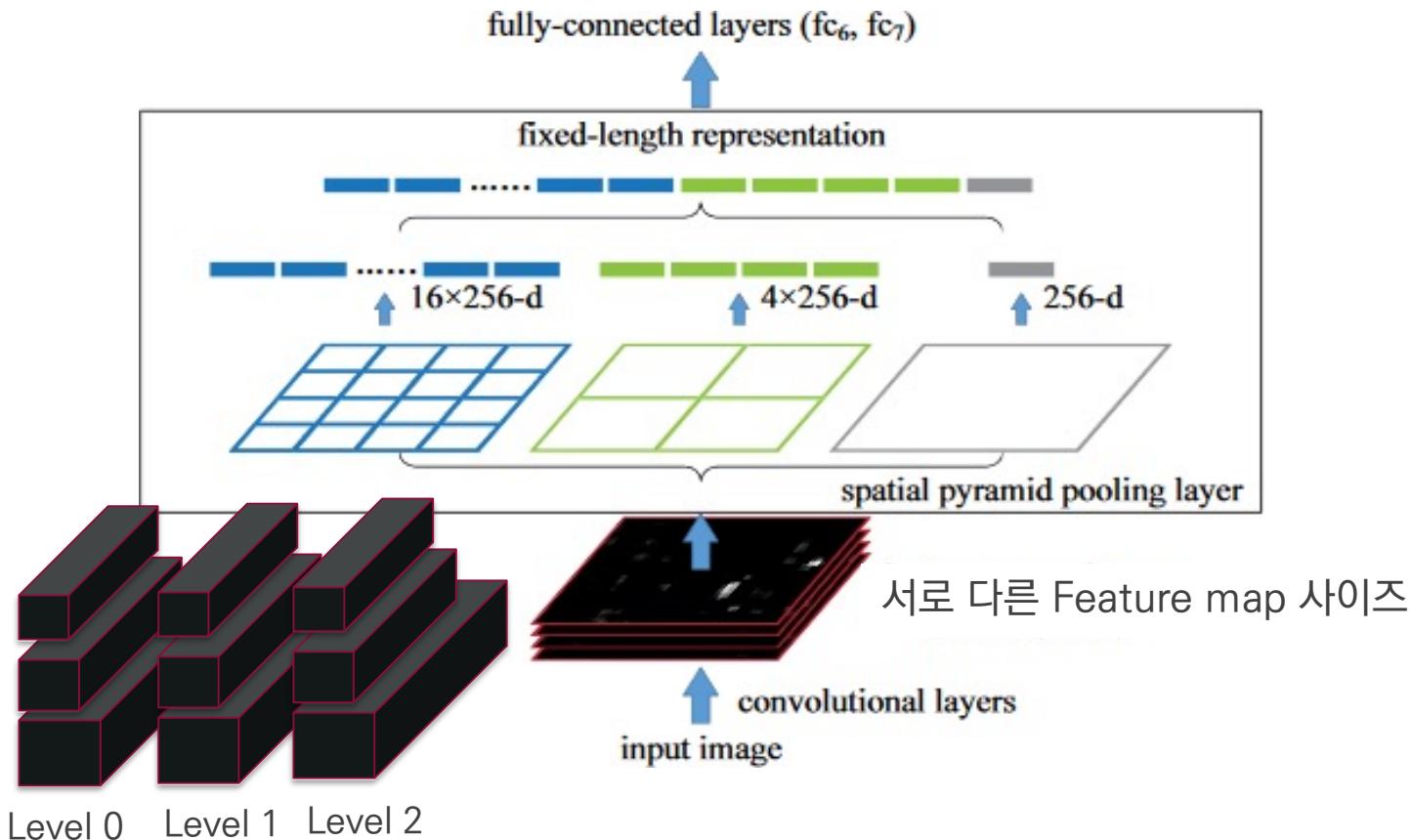
8 x 12 Feature Map



$$3 \times 16 = 48$$

$3 + 12 + 48 = 63$ 개 원소의 vector 값으로 표현 가능

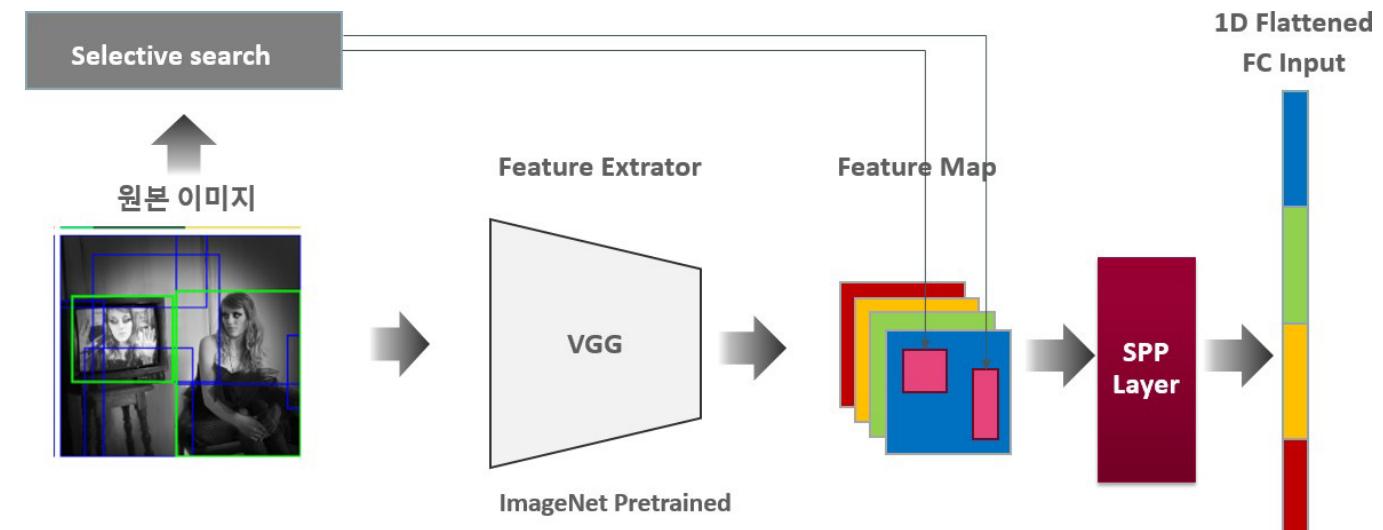
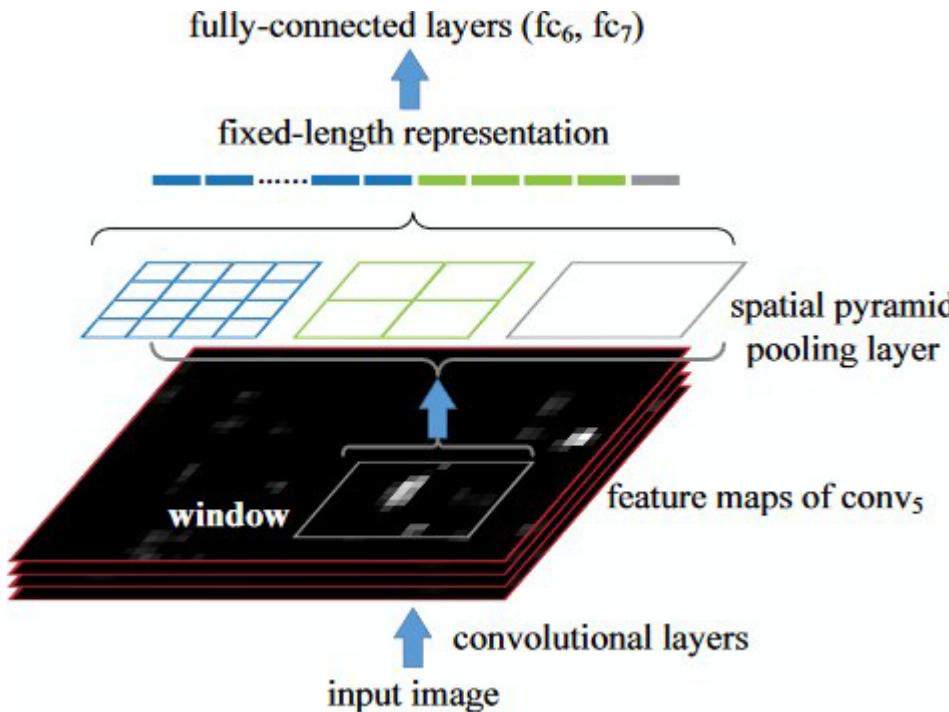
이미지 Classification에서의 SPP-Net 구조



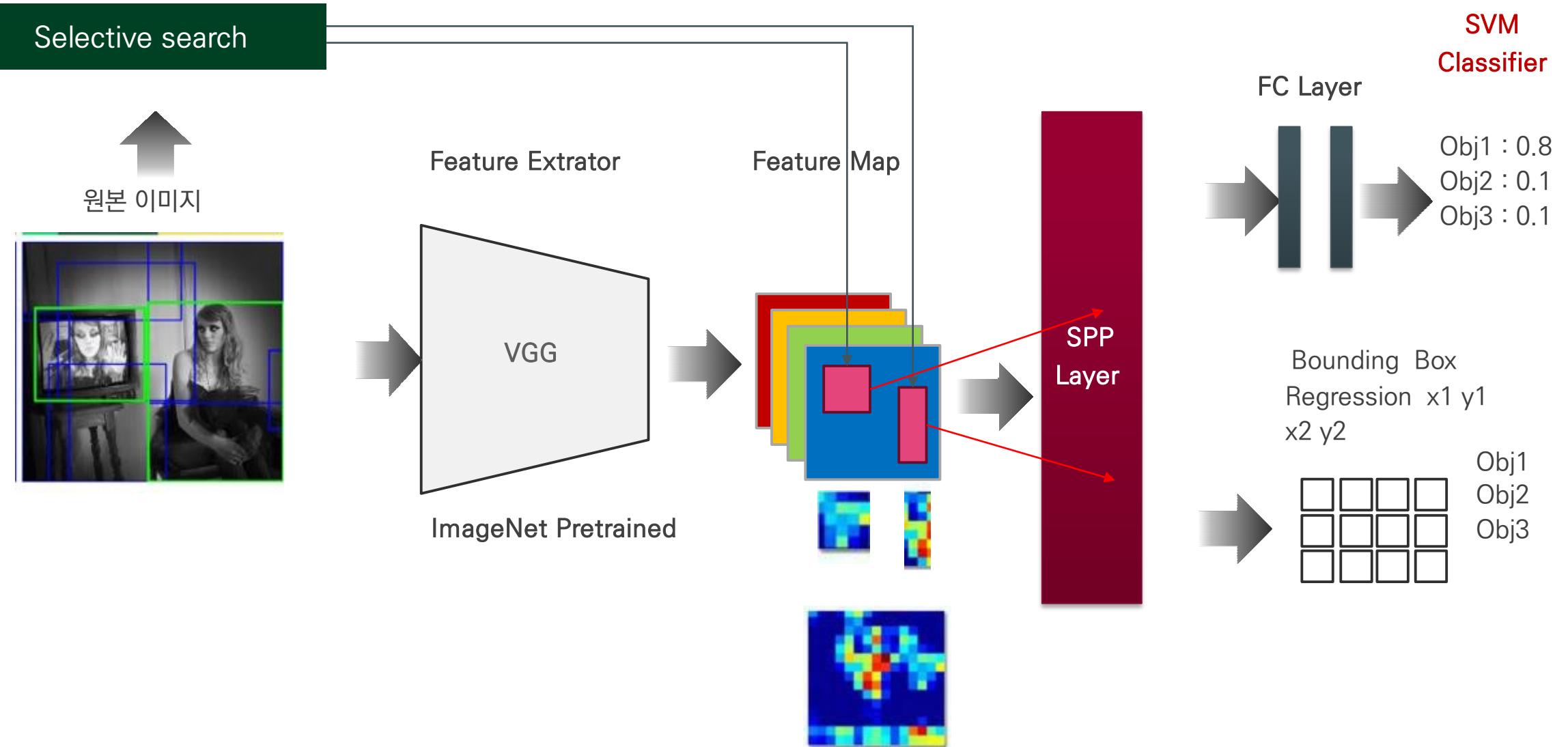
2014 ILSVRC Classification 결과

| rank | team | top-5 test |
|------|----------------|-------------|
| 1 | GoogLeNet [32] | 6.66 |
| 2 | VGG [33] | 7.32 |
| 3 | ours | 8.06 |
| 4 | Howard | 8.11 |
| 5 | DeeperVision | 9.50 |
| 6 | NUS-BST | 9.79 |
| 7 | TTIC_ECP | 10.22 |

| SPP-Net Object Detection

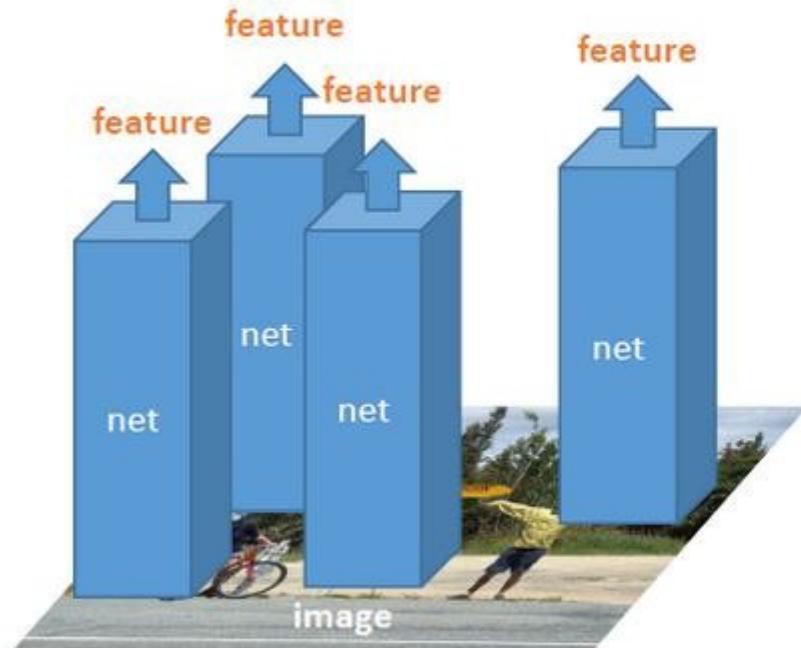


SPP-Net 구조



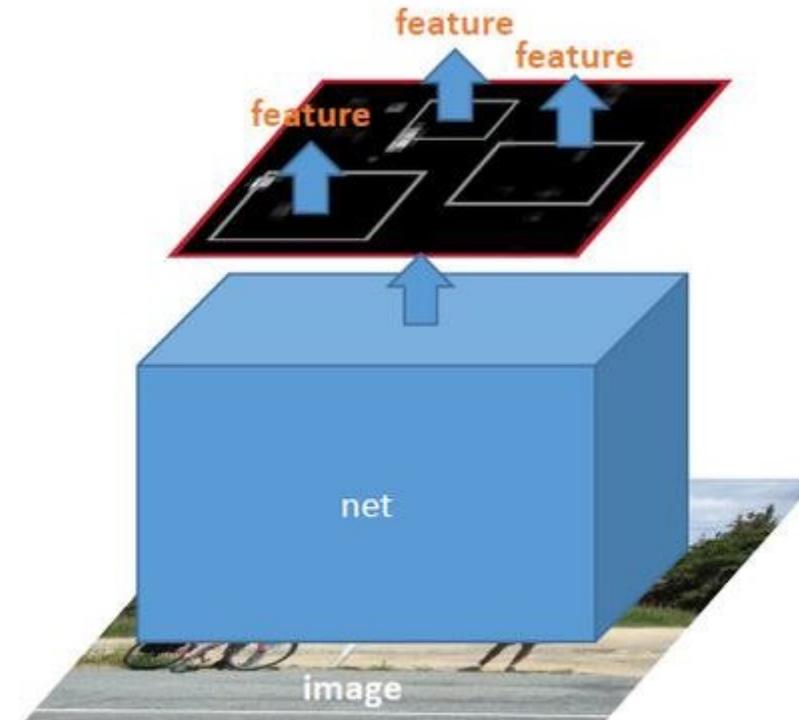
| R-CNN과 SPP-Net 비교

이미지 한 개에 2000번 CNN을 통과



R-CNN

이미지 한 개는 한번만 CNN 통과



SPP-Net

PASCAL VOC 2007 데이터 세트 적용 결과

| method | mAP | areo | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv |
|------------------------------|------|------|------|------|------|--------|------|------|------|-------|------|-------|------|-------|-------|--------|-------|-------|------|-------|------|
| DPM [23] | 33.7 | 33.2 | 60.3 | 10.2 | 16.1 | 27.3 | 54.3 | 58.2 | 23.0 | 20.0 | 24.1 | 26.7 | 12.7 | 58.1 | 48.2 | 43.2 | 12.0 | 21.1 | 36.1 | 46.0 | 43.5 |
| SS [20] | 33.8 | 43.5 | 46.5 | 10.4 | 12.0 | 9.3 | 49.4 | 53.7 | 39.4 | 12.5 | 36.9 | 42.2 | 26.4 | 47.0 | 52.4 | 23.5 | 12.1 | 29.9 | 36.3 | 42.2 | 48.8 |
| Regionlet [39] | 41.7 | 54.2 | 52.0 | 20.3 | 24.0 | 20.1 | 55.5 | 68.7 | 42.6 | 19.2 | 44.2 | 49.1 | 26.6 | 57.0 | 54.5 | 43.4 | 16.4 | 36.6 | 37.7 | 59.4 | 52.3 |
| DetNet [40] | 30.5 | 29.2 | 35.2 | 19.4 | 16.7 | 3.7 | 53.2 | 50.2 | 27.2 | 10.2 | 34.8 | 30.2 | 28.2 | 46.6 | 41.7 | 26.2 | 10.3 | 32.8 | 26.8 | 39.8 | 47.0 |
| RCNN ftfc ₇ (A5) | 54.2 | 64.2 | 69.7 | 50.0 | 41.9 | 32.0 | 62.6 | 71.0 | 60.7 | 32.7 | 58.5 | 46.5 | 56.1 | 60.6 | 66.8 | 54.2 | 31.5 | 52.8 | 48.9 | 57.9 | 64.7 |
| RCNN ftfc ₇ (ZF5) | 55.1 | 64.8 | 68.4 | 47.0 | 39.5 | 30.9 | 59.8 | 70.5 | 65.3 | 33.5 | 62.5 | 50.3 | 59.5 | 61.6 | 67.9 | 54.1 | 33.4 | 57.3 | 52.9 | 60.2 | 62.9 |
| SPP ftfc ₇ (ZF5) | 55.2 | 65.5 | 65.9 | 51.7 | 38.4 | 32.7 | 62.6 | 68.6 | 69.7 | 33.1 | 66.6 | 53.1 | 58.2 | 63.6 | 68.8 | 50.4 | 27.4 | 53.7 | 48.2 | 61.7 | 64.7 |
| RCNN bb (A5) | 58.5 | 68.1 | 72.8 | 56.8 | 43.0 | 36.8 | 66.3 | 74.2 | 67.6 | 34.4 | 63.5 | 54.5 | 61.2 | 69.1 | 68.6 | 58.7 | 33.4 | 62.9 | 51.1 | 62.5 | 64.8 |
| RCNN bb (ZF5) | 59.2 | 68.4 | 74.0 | 54.0 | 40.9 | 35.2 | 64.1 | 74.4 | 69.8 | 35.5 | 66.9 | 53.8 | 64.2 | 69.9 | 69.6 | 58.9 | 36.8 | 63.4 | 56.0 | 62.8 | 64.9 |
| SPP bb (ZF5) | 59.2 | 68.6 | 69.7 | 57.1 | 41.2 | 40.5 | 66.3 | 71.3 | 72.5 | 34.4 | 67.3 | 61.7 | 63.1 | 71.0 | 69.8 | 57.6 | 29.7 | 59.0 | 50.2 | 65.2 | 68.0 |

ILSVRC 2014 Object Detection 결과

| rank | team | mAP |
|------|--------------------|---------|
| 1 | NUS | 37.21 |
| 2 | <u>ours</u> | 35.11 |
| 3 | UvA | 32.02 |
| - | (our single-model) | (31.84) |
| 4 | Southeast-CASIA | 30.47 |
| 5 | 1-HKUST | 28.86 |
| 6 | CASIA_CRPAC_2 | 28.61 |

PASCAL VOC 2007 데이터 세트 기반 수행 시간

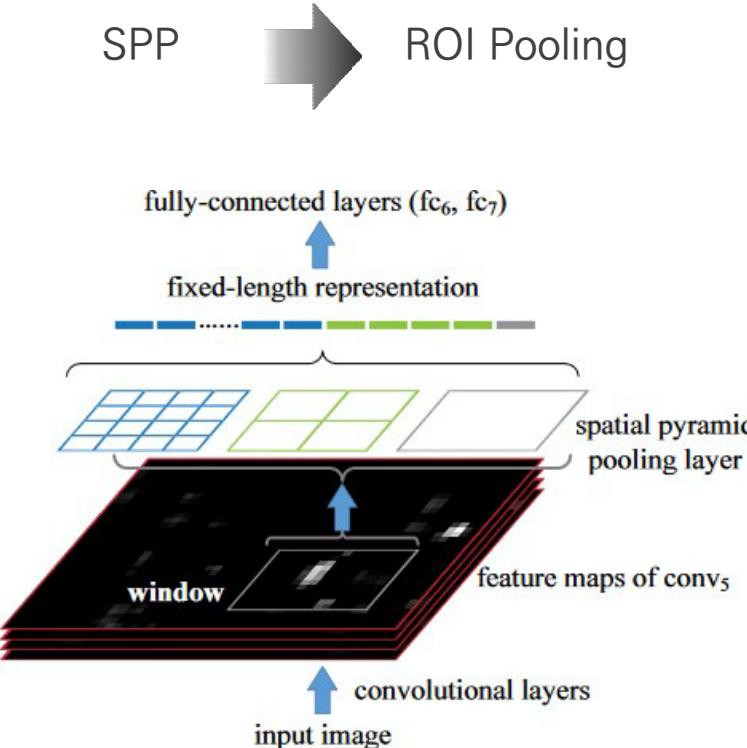
| | SPP-net 1-scale | SPP-net 5-scale | RCNN |
|----------------|--------------------|--------------------|------|
| mAP | 58.0 | 59.2 | 58.5 |
| GPU time / img | 0.14s | 0.38s | 9s |
| speed-up | 64x | 24x | - |



BIG DATA & AI ANALYTICS
EXPERT COMPANY

Fast RCNN

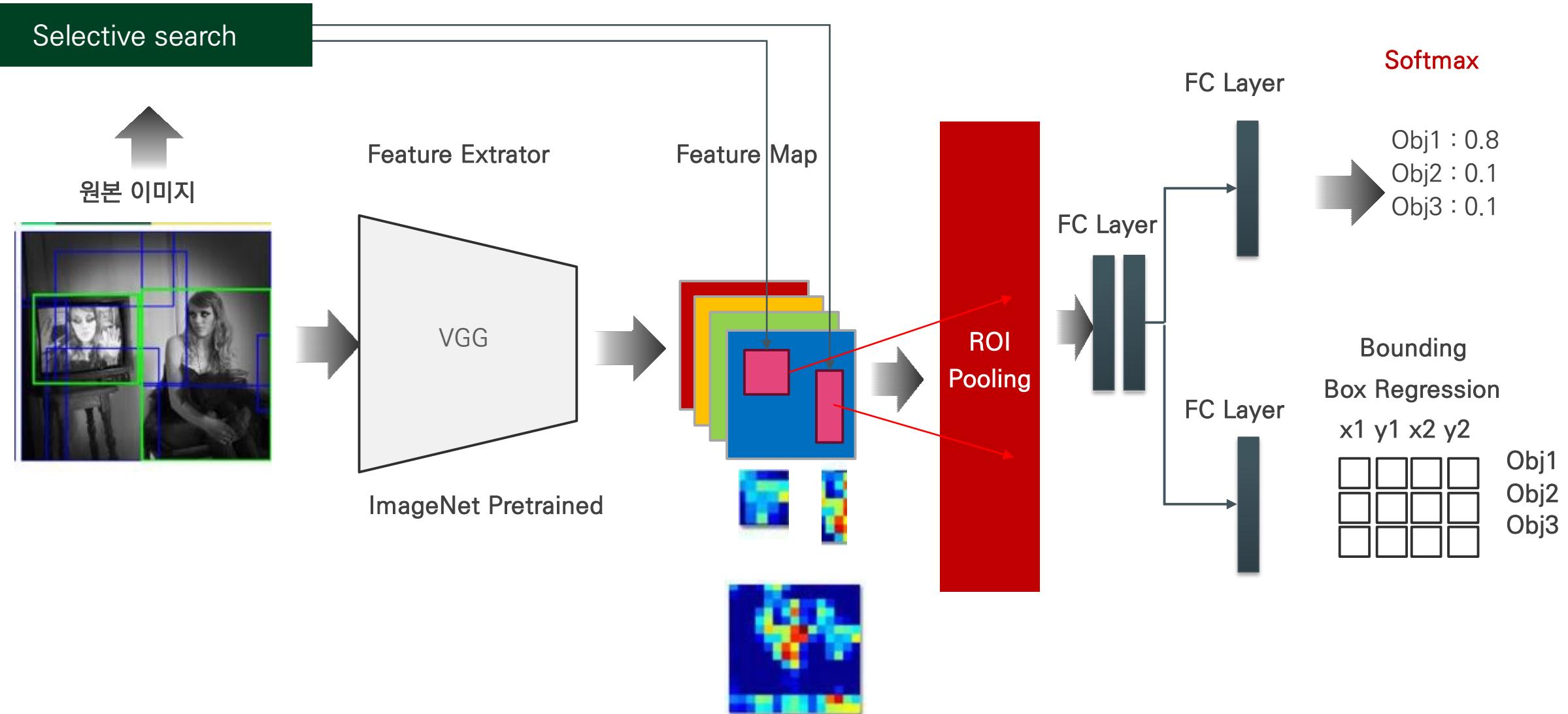
SPP Layer를 ROI Pooling Layer로



End-to-End Network Learning (ROI Proposal은 제외)

- SVM을 Softmax로 변환
- Multi-task loss 함수로 Classification과 Regression을 함께 최적화

| Fast RCNN 구조



Classification과 Regression Loss를 함께 반영한 Loss 함수

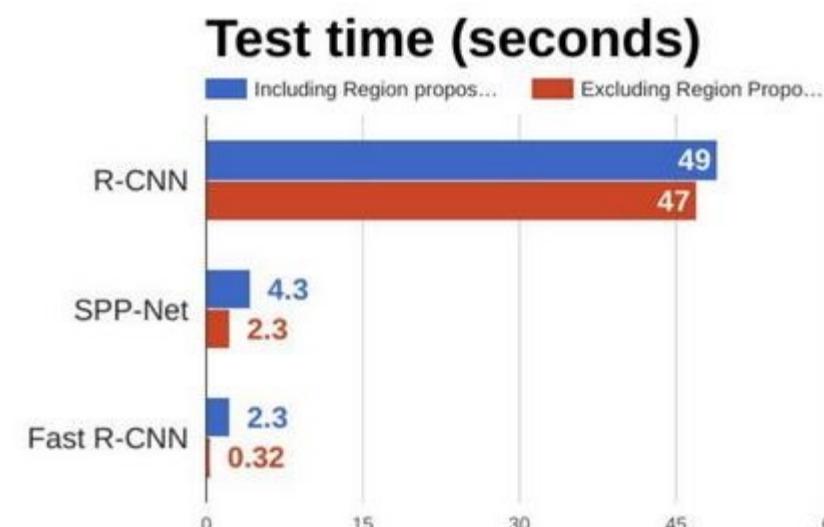
Multi-task Loss 함수:
$$L(p, u, t^u, v) = \underbrace{L_{\text{cls}}(p, u)}_{\text{Classification Loss}} + \lambda[u \geq 1] \underbrace{L_{\text{loc}}(t^u, v)}_{\text{Regression Loss}},$$

Regression Loss 함수:
$$L_{\text{loc}}(t^u, v) = \sum_{i \in \{x, y, w, h\}} \text{smooth}_{L_1}(t_i^u - v_i), \quad \text{smooth}_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise,} \end{cases}$$

PASCAL VOC 2012 데이터 세트 적용 결과

| method | train set | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | persn | plant | sheep | sofa | train | tv | mAP |
|---------------|-----------|------|------|------|------|--------|------|------|------|-------|------|-------|------|-------|-------|-------|-------|-------|------|-------|------|------|
| BabyLearning | Prop. | 78.0 | 74.2 | 61.3 | 45.7 | 42.7 | 68.2 | 66.8 | 80.2 | 40.6 | 70.0 | 49.8 | 79.0 | 74.5 | 77.9 | 64.0 | 35.3 | 67.9 | 55.7 | 68.7 | 62.6 | 63.2 |
| NUS_NIN_c2000 | Unk. | 80.2 | 73.8 | 61.9 | 43.7 | 43.0 | 70.3 | 67.6 | 80.7 | 41.9 | 69.7 | 51.7 | 78.2 | 75.2 | 76.9 | 65.1 | 38.6 | 68.3 | 58.0 | 68.7 | 63.3 | 63.8 |
| R-CNN BB [10] | 12 | 79.6 | 72.7 | 61.9 | 41.2 | 41.9 | 65.9 | 66.4 | 84.6 | 38.5 | 67.2 | 46.7 | 82.0 | 74.8 | 76.0 | 65.2 | 35.6 | 65.4 | 54.2 | 67.4 | 60.3 | 62.4 |
| FRCN [ours] | 12 | 80.3 | 74.7 | 66.9 | 46.9 | 37.7 | 73.9 | 68.6 | 87.7 | 41.7 | 71.1 | 51.1 | 86.0 | 77.8 | 79.8 | 69.8 | 32.1 | 65.5 | 63.8 | 76.4 | 61.7 | 65.7 |
| FRCN [ours] | 07++12 | 82.3 | 78.4 | 70.8 | 52.3 | 38.7 | 77.8 | 71.6 | 89.3 | 44.2 | 73.0 | 55.0 | 87.5 | 80.5 | 80.8 | 72.0 | 35.1 | 68.3 | 65.7 | 80.4 | 64.2 | 68.4 |

수행 시간 비교

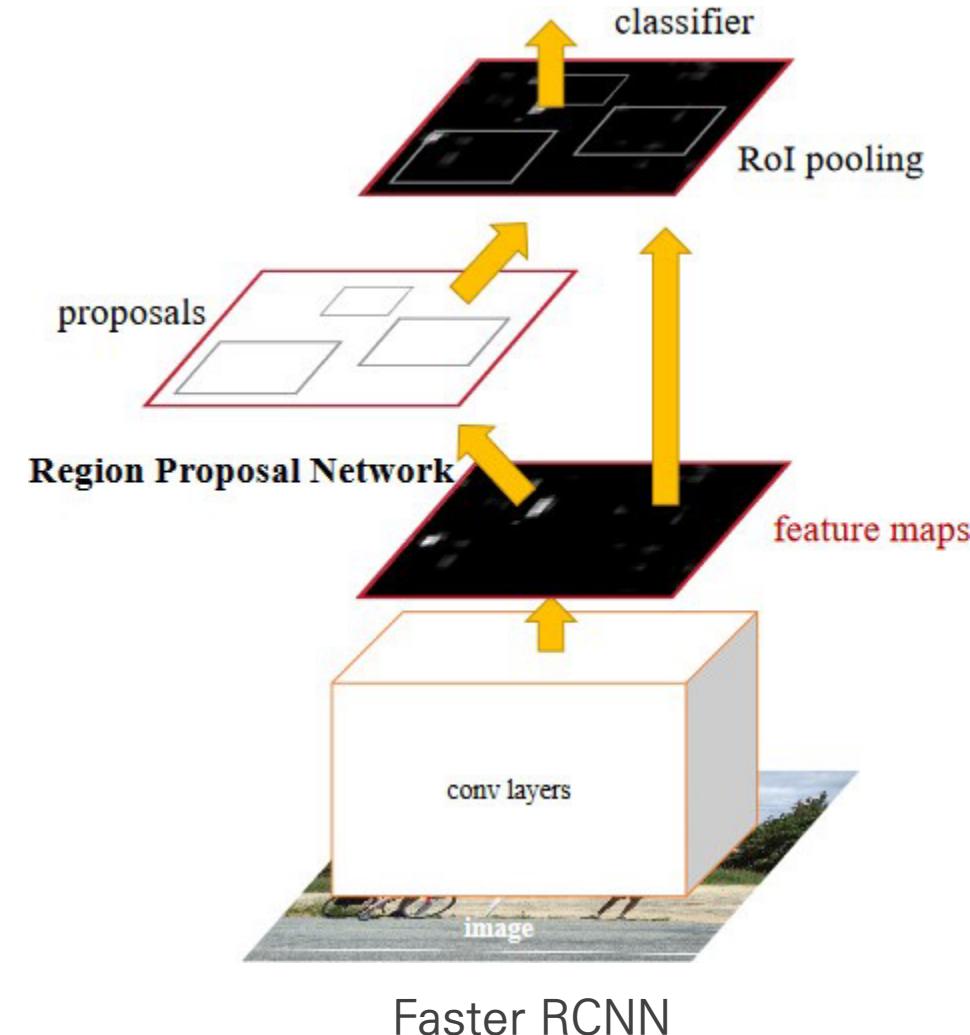


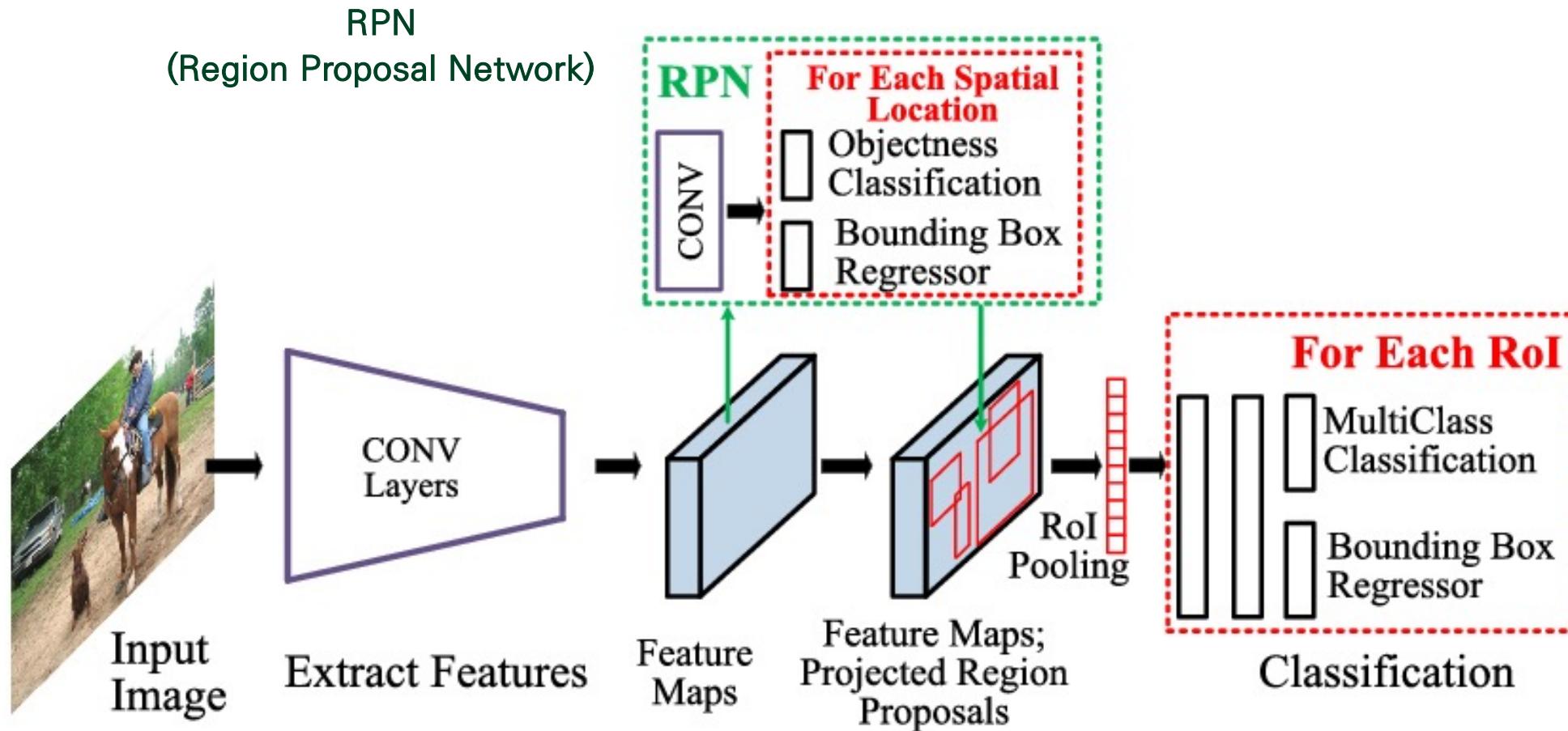


BIG DATA & AI ANALYTICS
EXPERT COMPANY

Faster RCNN

Faster RCNN = RPN + Fast RCNN



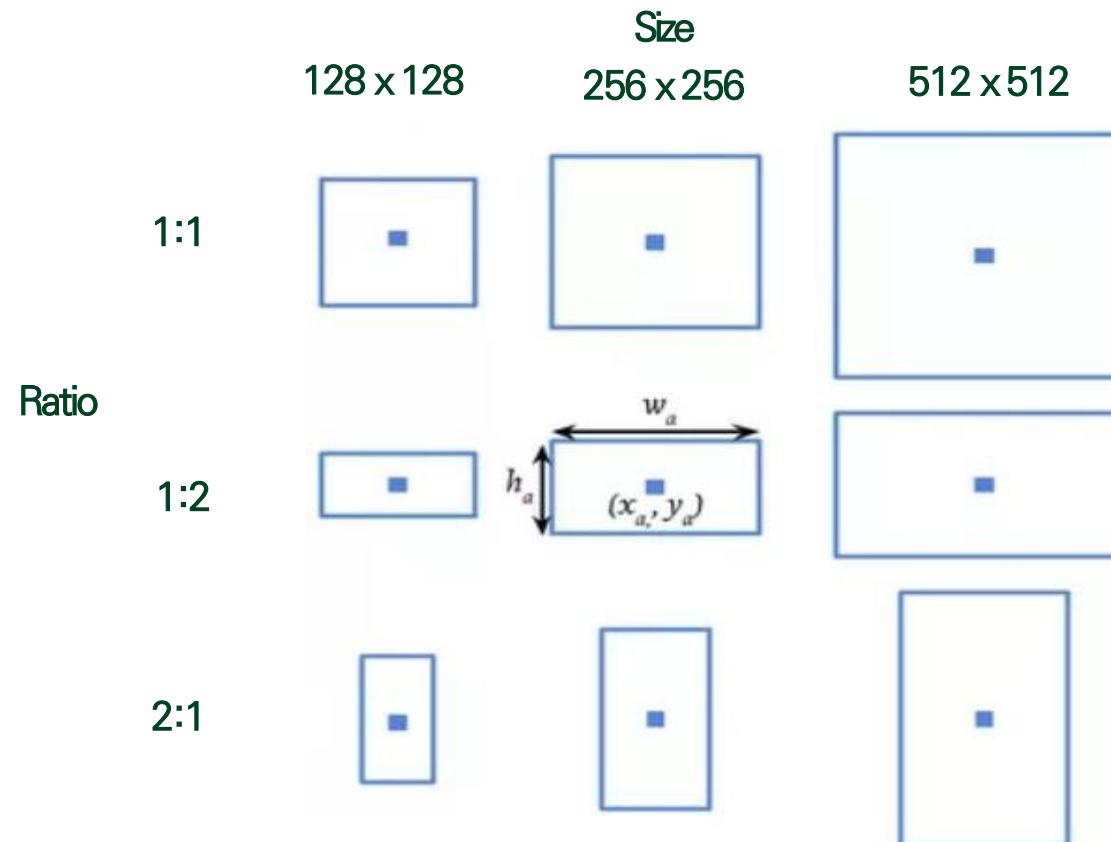


Selective Search를 대체하기 위한 Region Proposal Network 구현 이슈

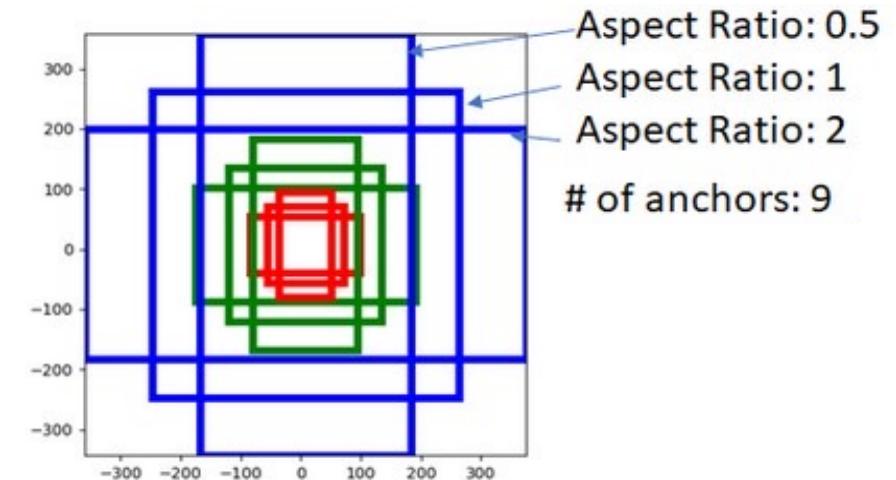
- 데이터로 주어질 피처는 pixel 값, Target은 Ground Truth Bounding Box인데 이를 이용해 어떻게 Selective Search 수준의 Region Proposal을 할 수 있을 것인가?

(Reference) Anchor Box

Object 가 있는지 없는지의 후보 Box



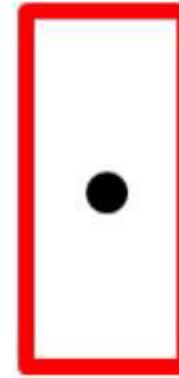
총 9 개의 Anchor box, 3개의 서로 다른 크기,
3개의 서로 다른 ratio로 구성



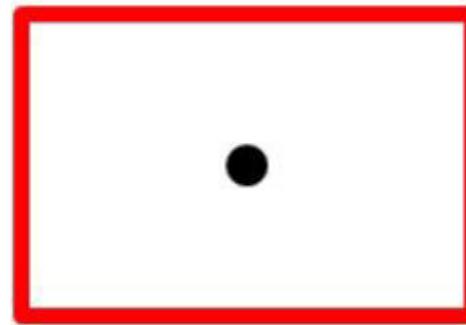
| Anchor Box



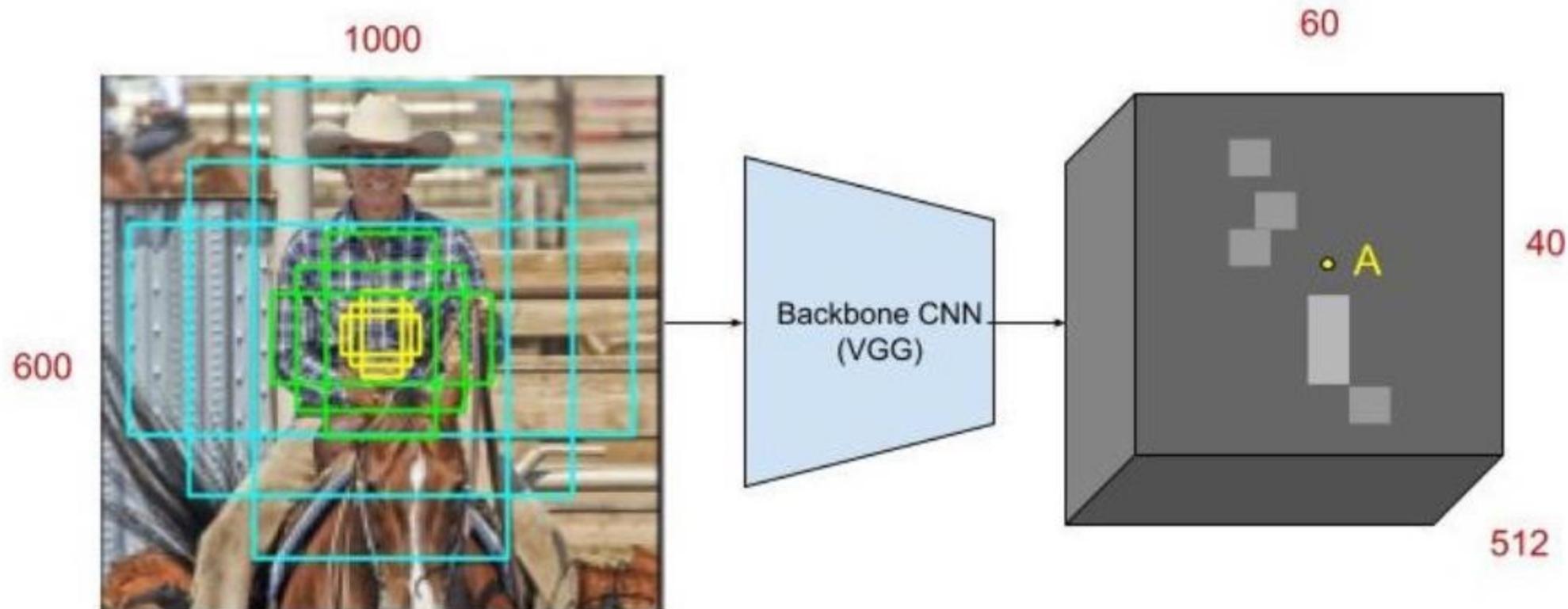
Anchor box 1:



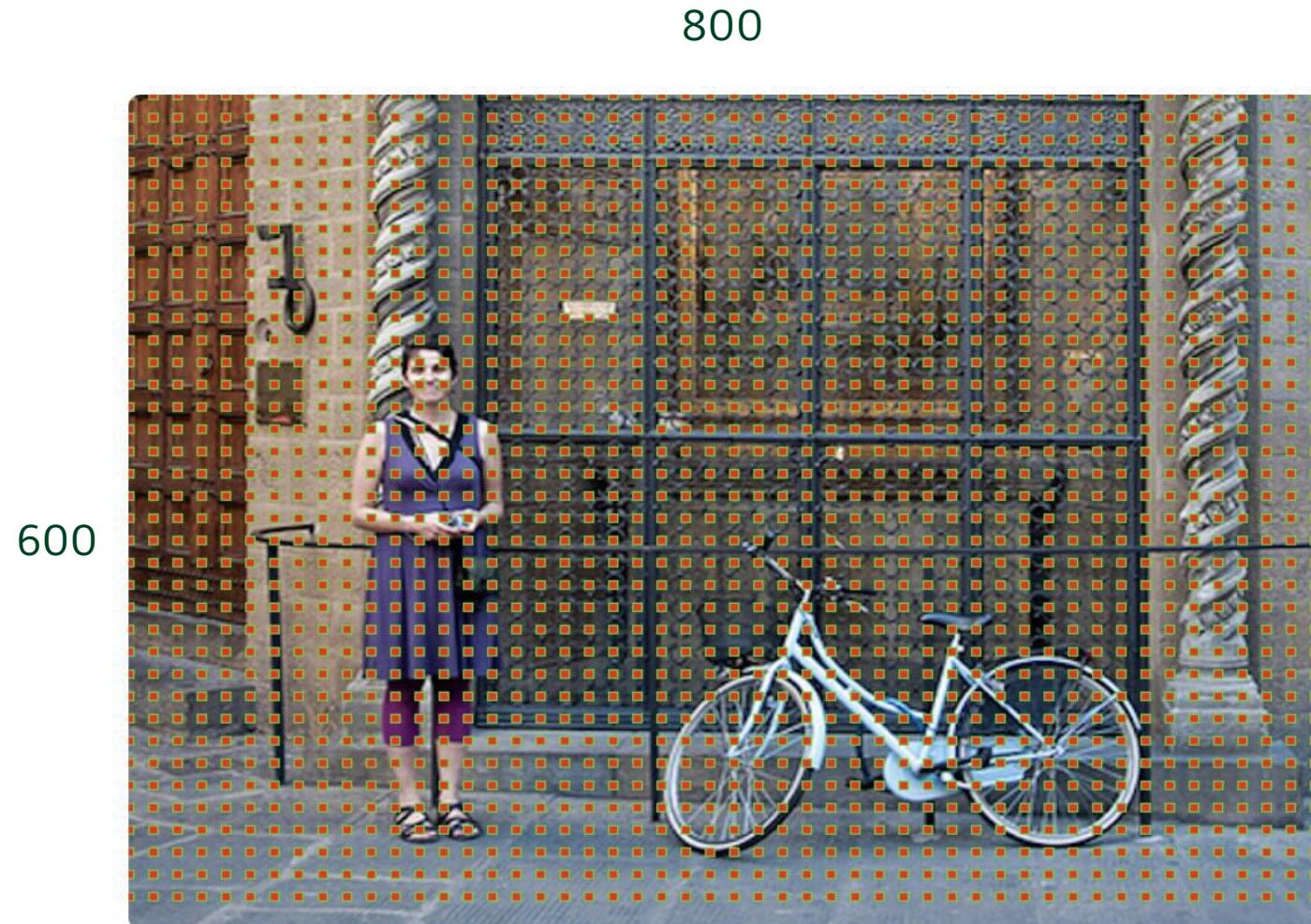
Anchor box 2:



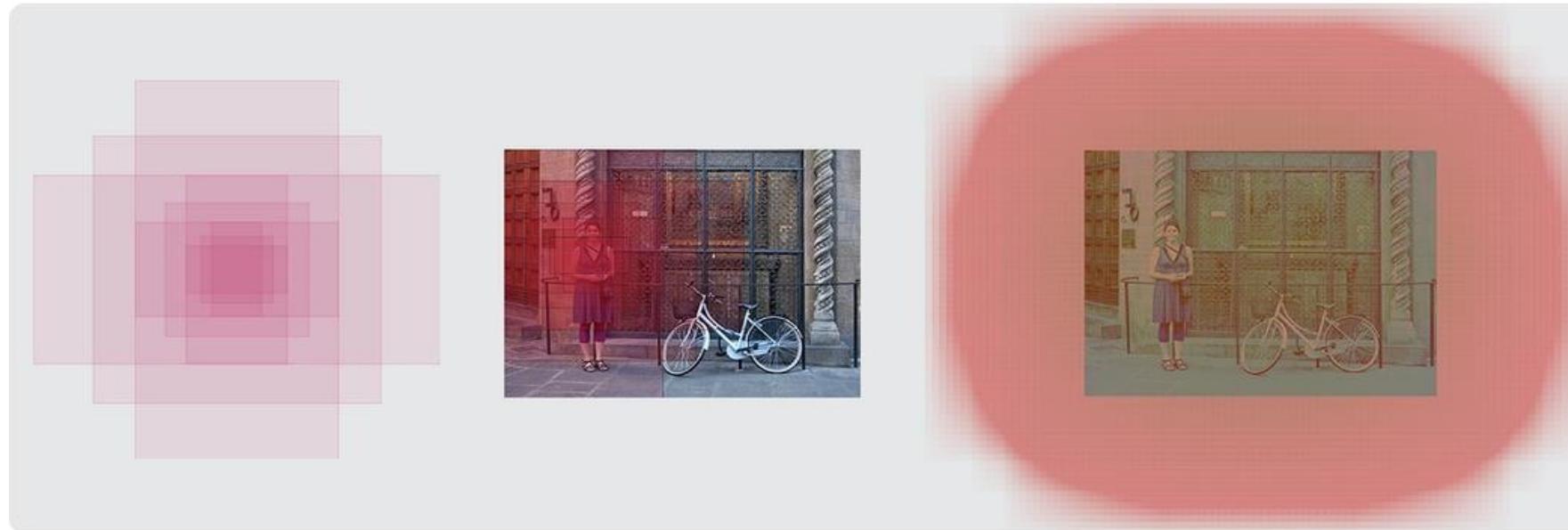
- 이미지와 Feature Map에서 Anchor Box 매핑



- 이미지와 Feature Map에서 Anchor Box 매핑



- 이미지와 Feature Map에서 Anchor Box 매핑



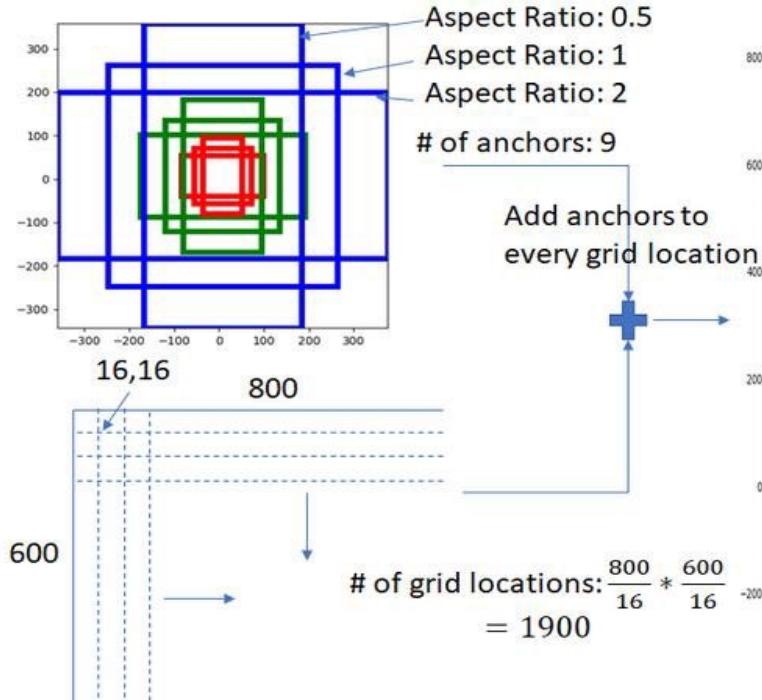
Left: Anchors, Center: Anchor for a single point, Right: All anchors

- 이미지와 Feature Map에서 Anchor Box 매핑

Generate Anchors

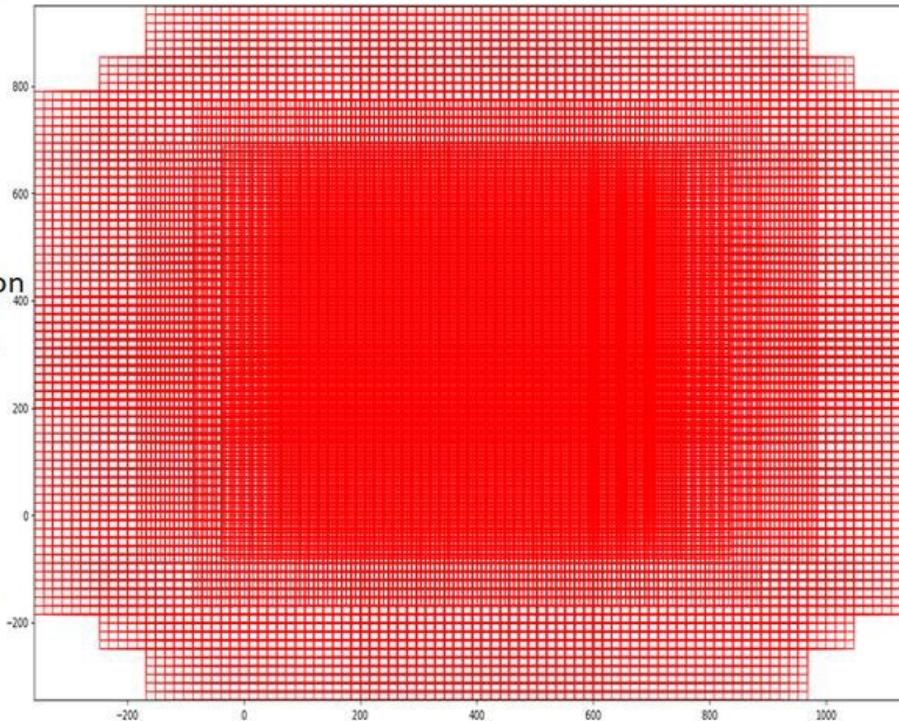
Given:

- Set of aspect ratios (0.5, 1, 2)
- Stride length (downscaling performed by resnet head: 16)
- Anchor Scales (8, 16, 32)

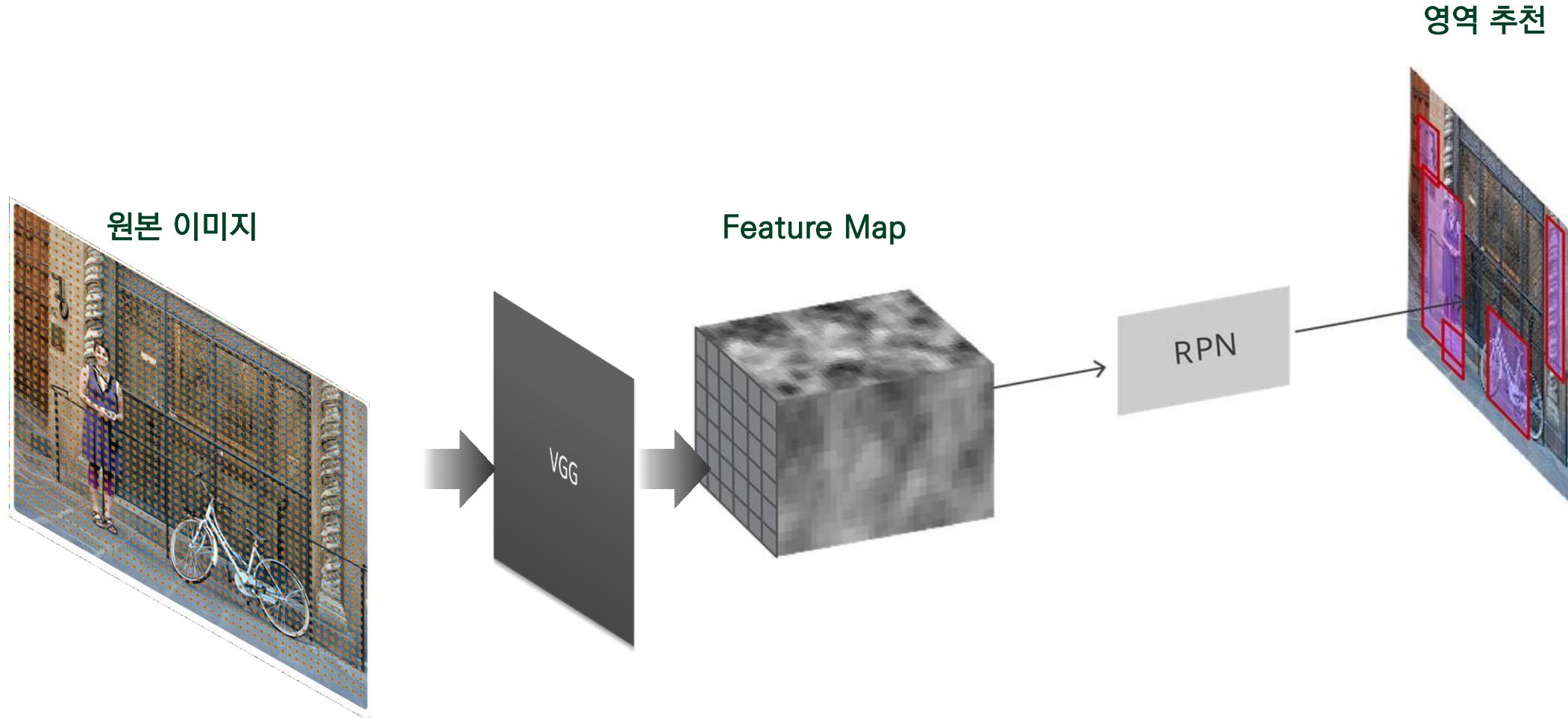


Create uniformly spaced grid with
spacing = stride length

Total number of anchors: $1900 * 9 = 17100$
Some boxes lie outside the image boundary

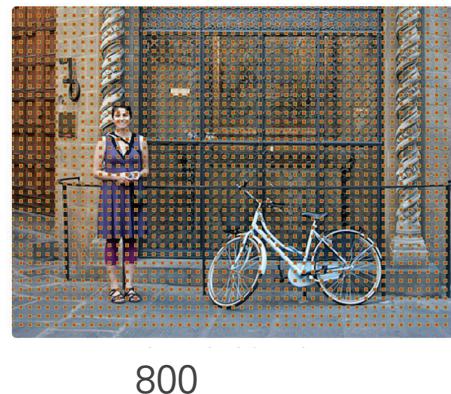


| Region Proposal Network(RPN)

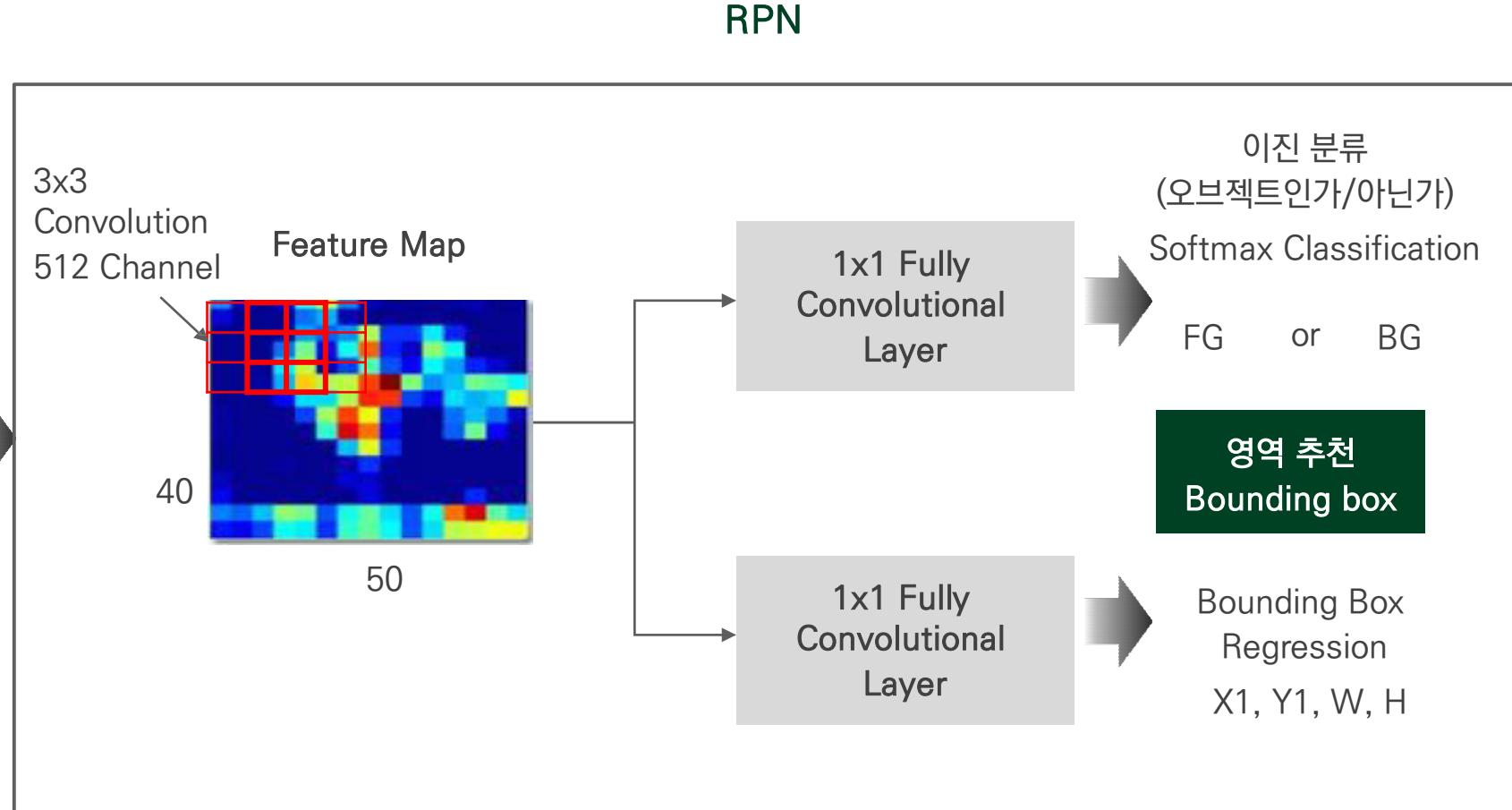


| RPN 구조

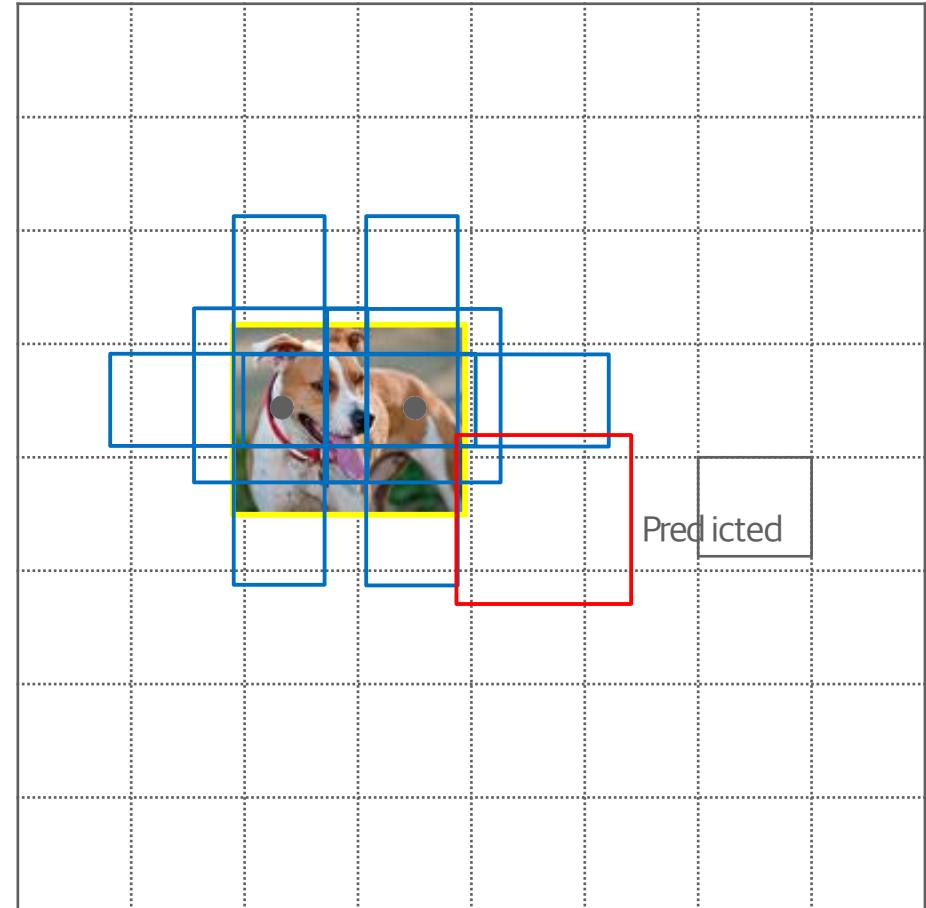
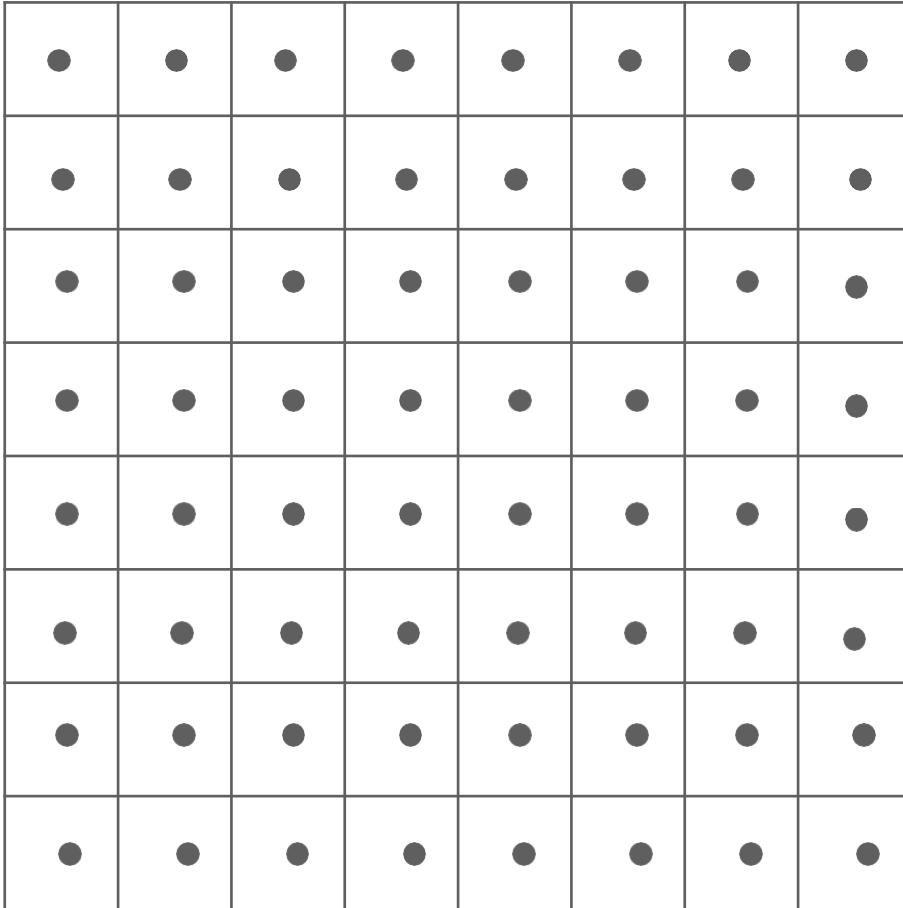
원본 이미지

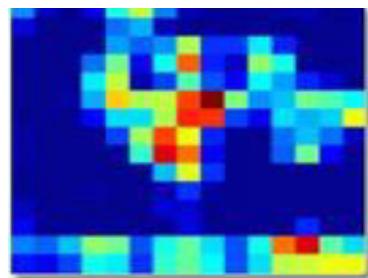


VGG

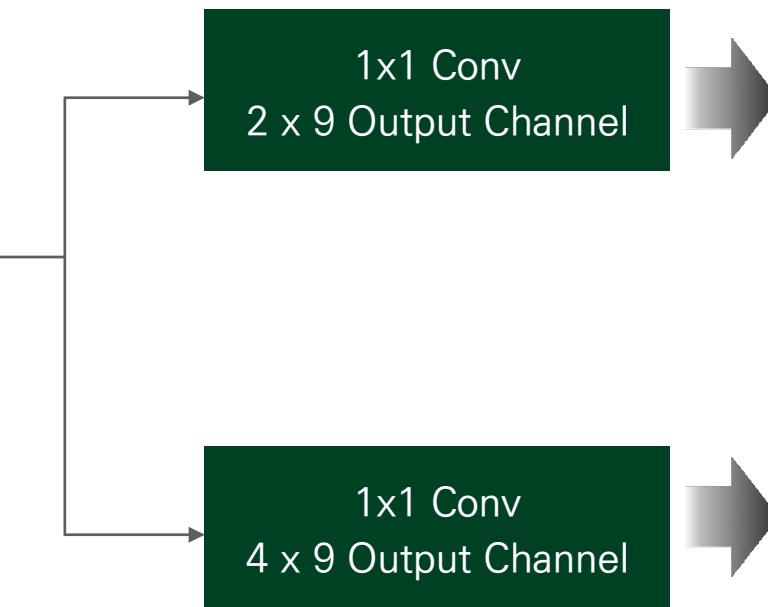


Anchor box 와 Predicted Anchor box





3 x 3 Conv
512 Channel

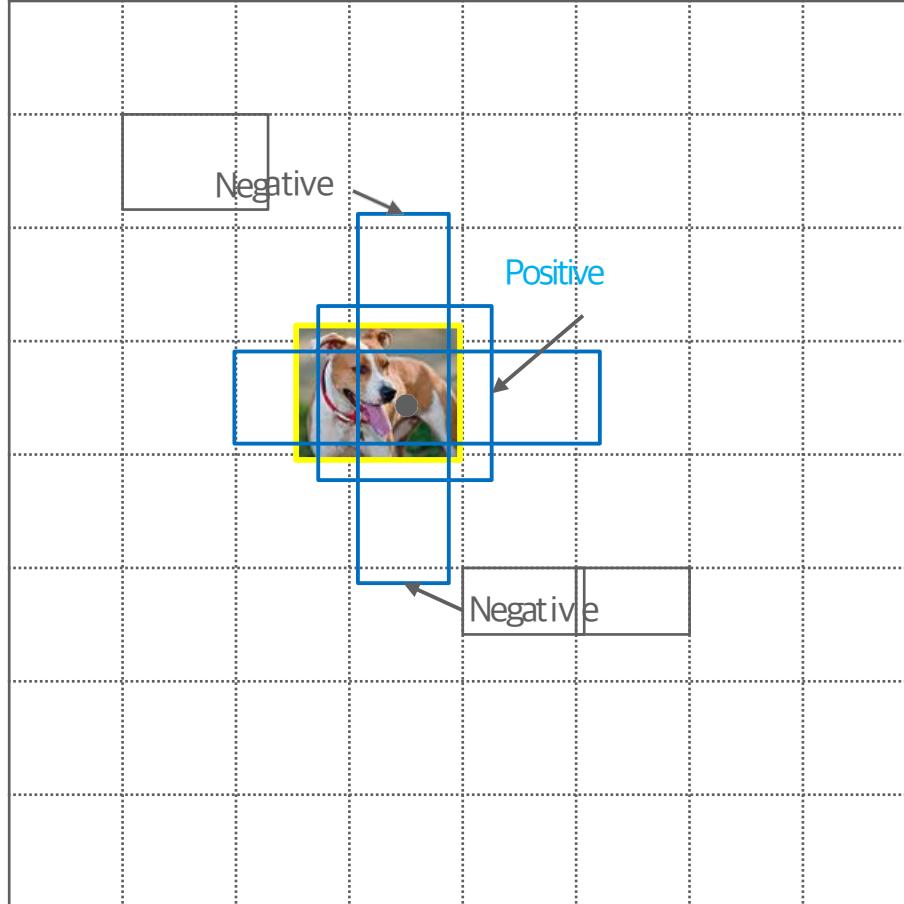


이진 분류 (오브젝트인가/아
닌가) Softmax
Classification
FG or BG

영역 추천
Bounding box

Bounding Box
Regression
X1, Y1, W, H

Positive/Negative Anchor Box

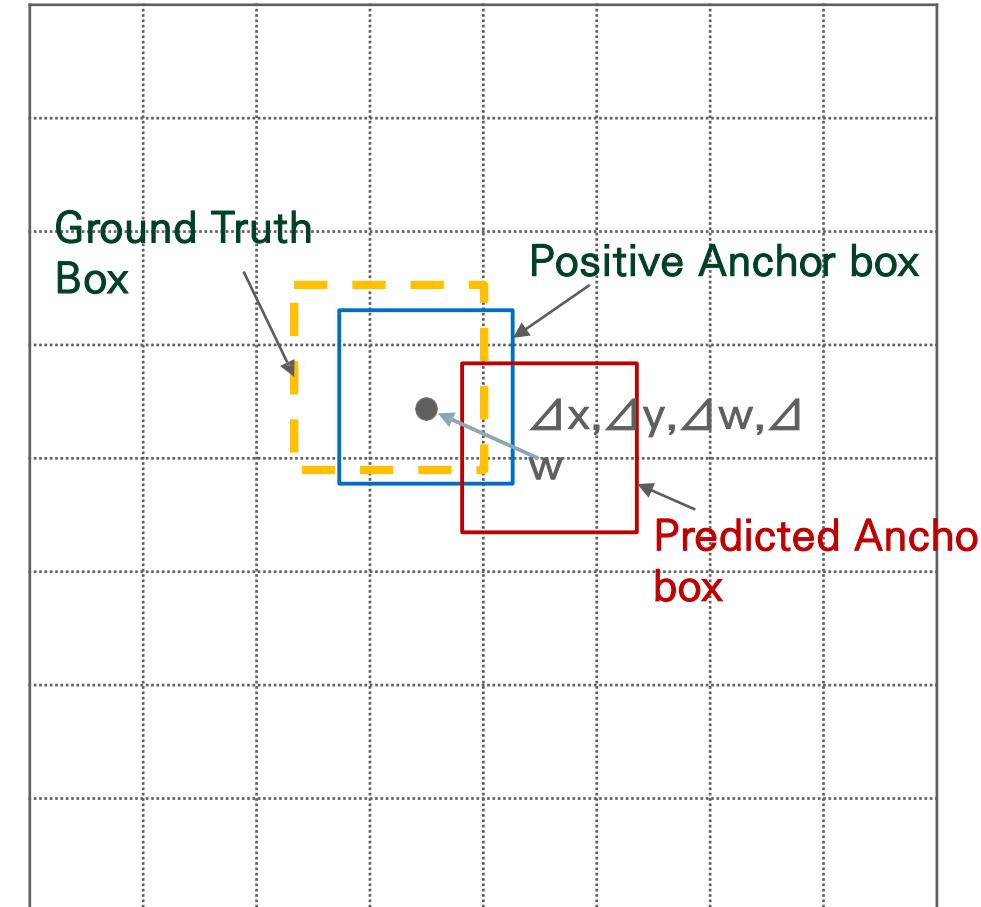


Ground Truth BB 겹치는 IOU 값에 따라 Anchor Box를 Positive Anchor Box, Negative Anchor box로 분류

- IOU가 0.7 이상이면 Positive
- IOU가 가장 높은 Anochor는 Positive
- IOU가 0.3보다 낮으면 Negative

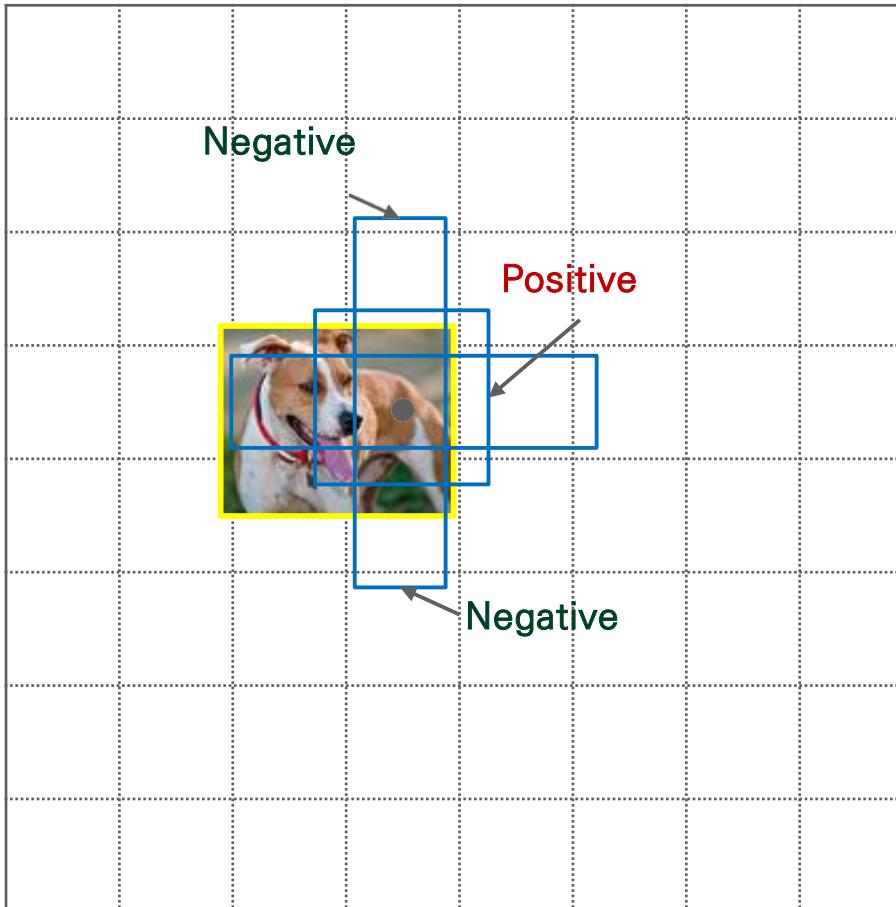
| Bounding Box Regression with Anchor Box

예측 Anchor box는 Positive Anchor box와의 좌표 값 차이를 최소화 할 수 있는
Bounding Box Regression 수행

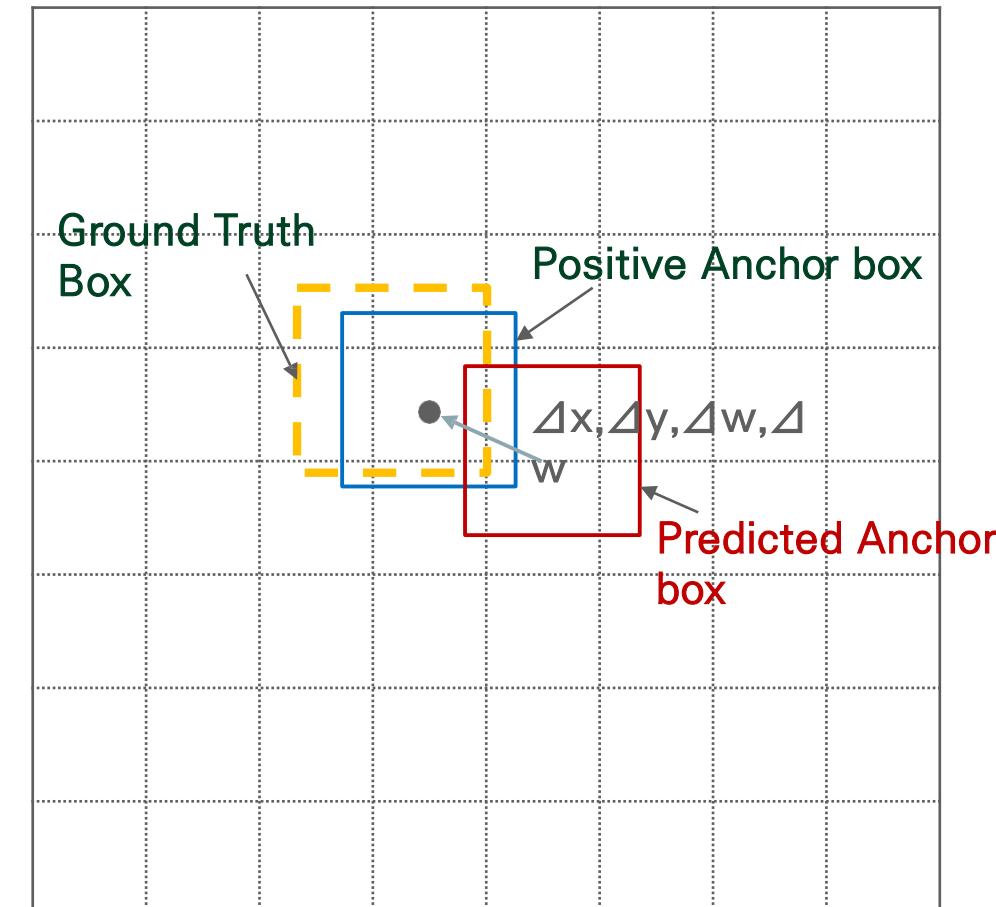


RPN Classification, Bounding Box Regression

Classification

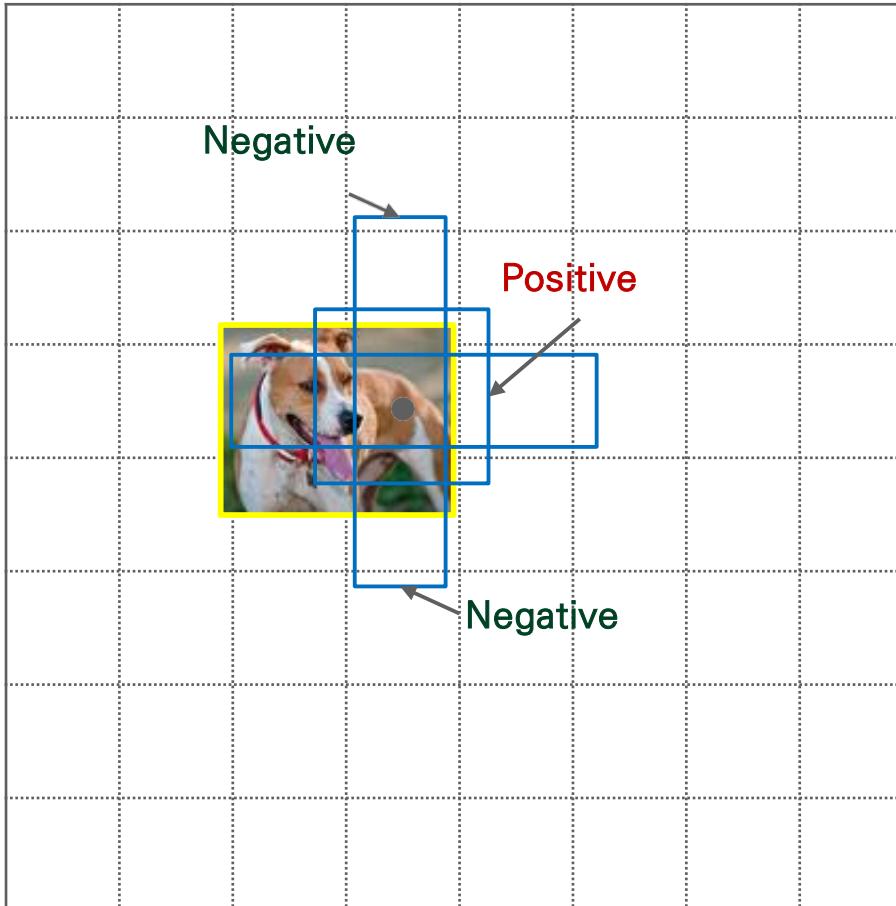


Bounding Box Regression

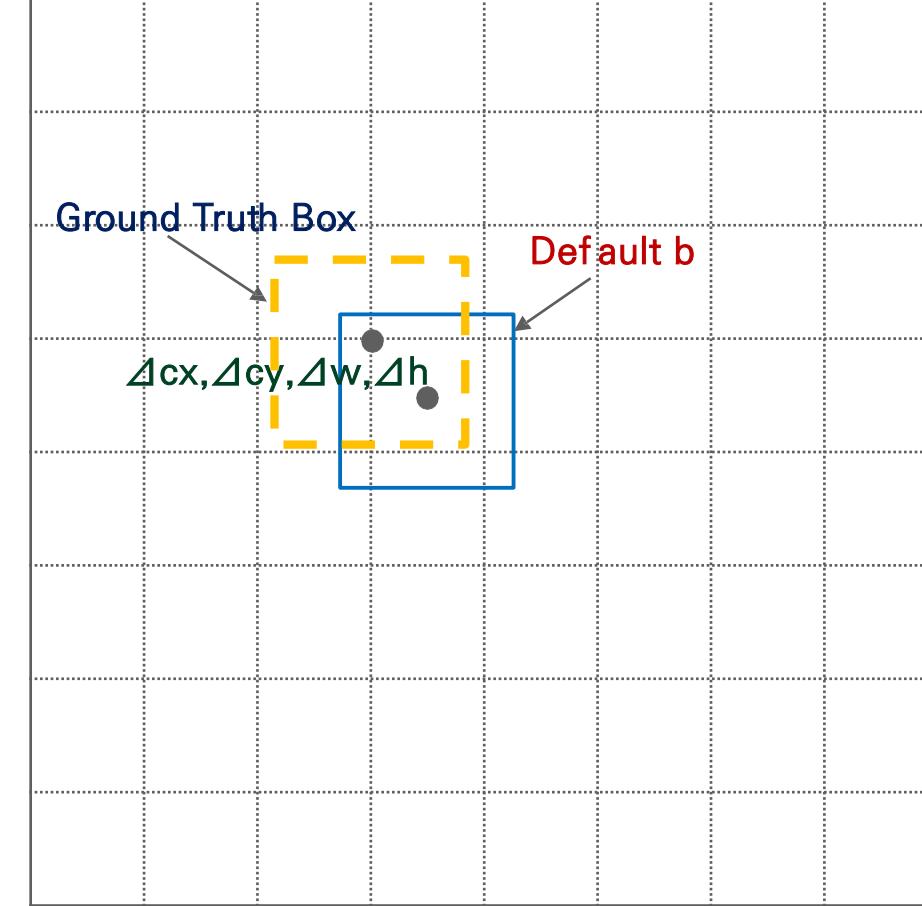


RPN Classification, Bounding Box Regression

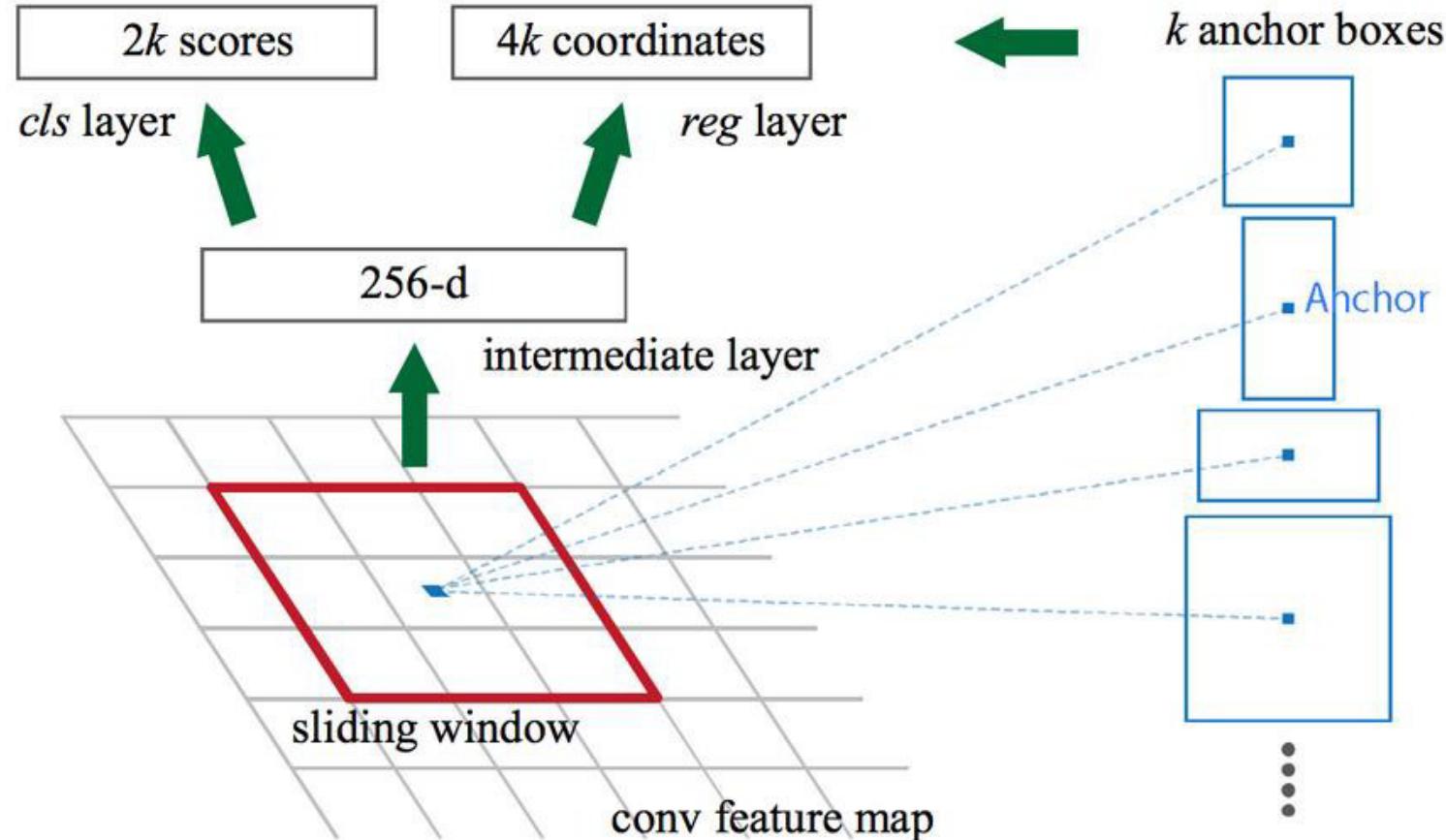
Classification



Bounding Box Regression



| Anchor box에 따른 RPN Output



$$\mathcal{L}(\{p_i\}, \{t_i\}) = \frac{1}{N_{\text{cls}}} \sum_i \mathcal{L}_{\text{cls}}(p_i, p_i^*) + \frac{\lambda}{N_{\text{box}}} \sum_i p_i^* \cdot L_1^{\text{smooth}}(t_i - t_i^*)$$

Symbol Explanation

p_i Anchor i가 오브젝트일 예측 확률

p_i^* Anchor i의 Ground truth Object 여부(Positive 1, Negative 0)

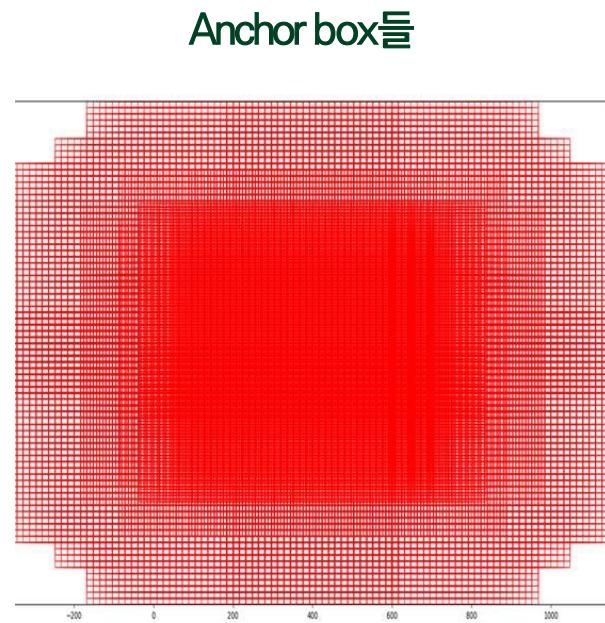
t_i Anchor i에 대한 예측 좌표 4개(x, y, w, h)

t_i^* Anchor i에 실제 좌표 4개(x, y, w, h)

N_{cls} 미니 배치에 따른 정규화 값(256)

N_{box} 박스 개수 정규화 값 (최대 2400)

λ 밸런싱 값: 10



Mini Batch

256 anchors
(128 positive, 128 negative)

.....

.....

256 anchors
(128 positive, 128 negative)

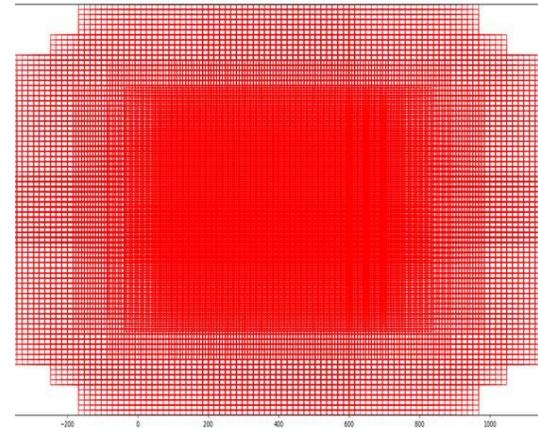
Training



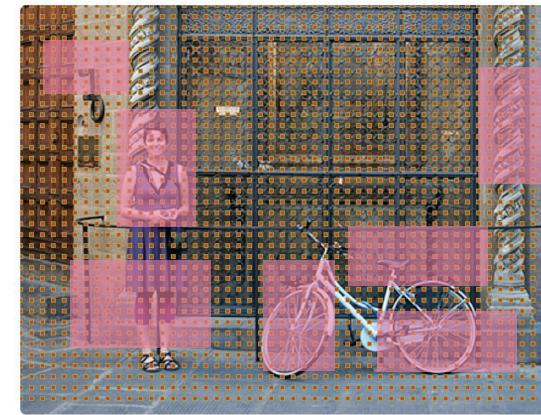
이진 분류 (Object/not Object) Softmax Classification

| 예측 Region Proposal box의 Objectness Score

- Objectness Score: 예측 Box가 Object 일 확률(Softmax 값) * Ground Truth bounding box와의 IOU값
- Objectness Score가 높은 순으로 Region Proposal Box를 추출

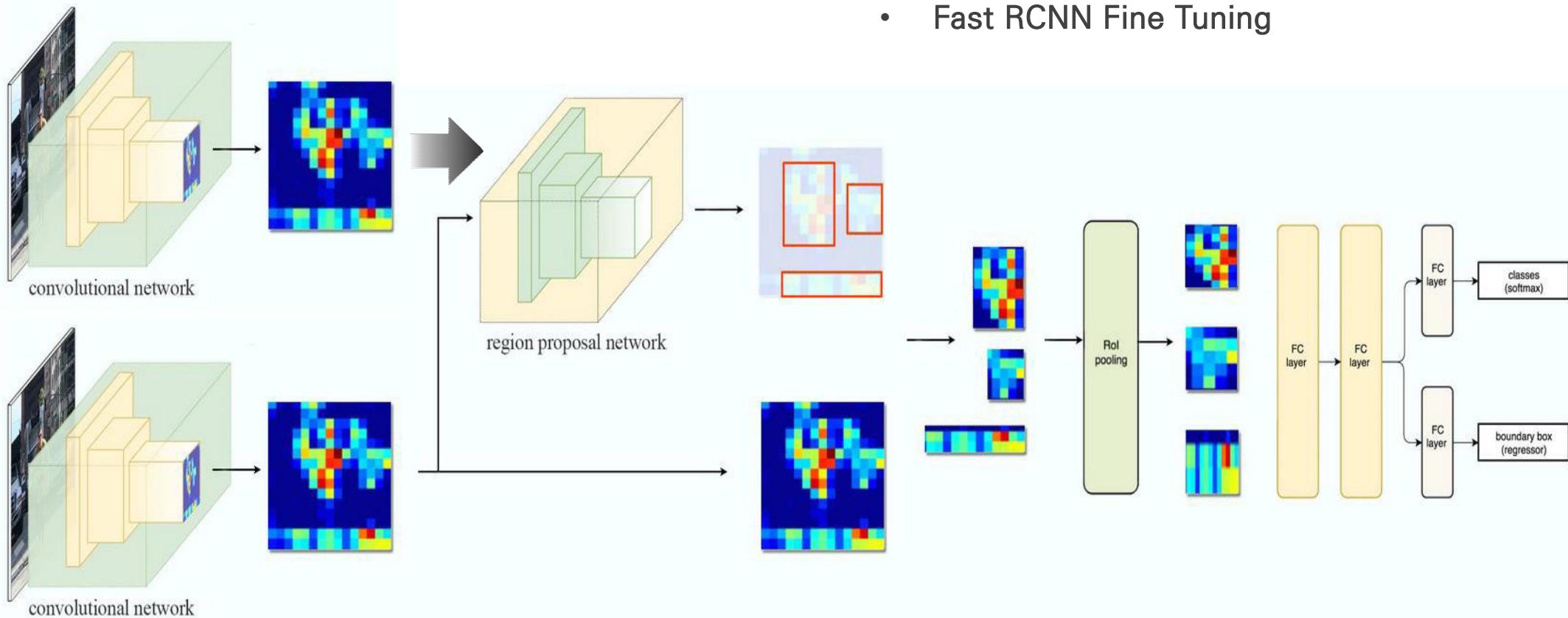


$N < 2000$
.....
 $N < 300$
.....
 $N < 50$



| Faster RCNN 학습 순서

- RPN을 먼저 학습
- Fast RCNN Classification/Regression 학습
- RPN Fine Tuning
- Fast RCNN Fine Tuning



PASCAL VOC 데이터 세트에서 Detection 성능

| method | # box | data | mAP | areo | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv |
|--------|-------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| SS | 2000 | 12 | 65.7 | 80.3 | 74.7 | 66.9 | 46.9 | 37.7 | 73.9 | 68.6 | 87.7 | 41.7 | 71.1 | 51.1 | 86.0 | 77.8 | 79.8 | 69.8 | 32.1 | 65.5 | 63.8 | 76.4 | 61.7 |
| SS | 2000 | 07++12 | 68.4 | 82.3 | 78.4 | 70.8 | 52.3 | 38.7 | 77.8 | 71.6 | 89.3 | 44.2 | 73.0 | 55.0 | 87.5 | 80.5 | 80.8 | 72.0 | 35.1 | 68.3 | 65.7 | 80.4 | 64.2 |
| RPN | 300 | 12 | 67.0 | 82.3 | 76.4 | 71.0 | 48.4 | 45.2 | 72.1 | 72.3 | 87.3 | 42.2 | 73.7 | 50.0 | 86.8 | 78.7 | 78.4 | 77.4 | 34.5 | 70.1 | 57.1 | 77.1 | 58.9 |
| RPN | 300 | 07++12 | 70.4 | 84.9 | 79.8 | 74.3 | 53.9 | 49.8 | 77.5 | 75.9 | 88.5 | 45.6 | 77.1 | 55.3 | 86.9 | 81.7 | 80.9 | 79.6 | 40.1 | 72.6 | 60.9 | 81.2 | 61.5 |
| RPN | 300 | COCO+07++12 | <u>75.9</u> | <u>87.4</u> | <u>83.6</u> | <u>76.8</u> | <u>62.9</u> | <u>59.6</u> | <u>81.9</u> | <u>82.0</u> | <u>91.3</u> | <u>54.9</u> | <u>82.6</u> | <u>59.0</u> | <u>89.0</u> | <u>85.5</u> | <u>84.7</u> | <u>84.1</u> | <u>52.2</u> | <u>78.9</u> | <u>65.5</u> | <u>85.4</u> | <u>70.2</u> |

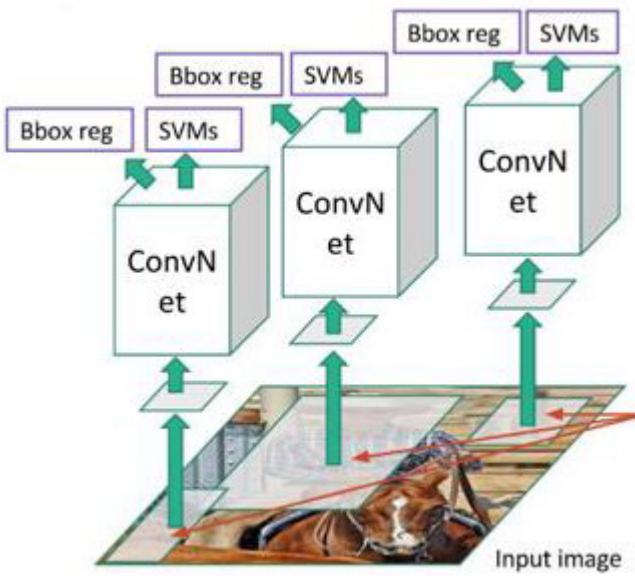
Coco 데이터 세트에서 Detection 성능

| method | proposals | training data | COCO val | | COCO test-dev | |
|----------------------------------|-----------|---------------|----------|---------------|---------------|---------------|
| | | | mAP@.5 | mAP@[.5, .95] | mAP@.5 | mAP@[.5, .95] |
| Fast R-CNN [2] | SS, 2000 | COCO train | - | - | 35.9 | 19.7 |
| Fast R-CNN [impl. in this paper] | SS, 2000 | COCO train | 38.6 | 18.9 | 39.3 | 19.3 |
| Faster R-CNN | RPN, 300 | COCO train | 41.5 | 21.2 | 42.1 | 21.5 |
| Faster R-CNN | RPN, 300 | COCO trainval | - | - | 42.7 | 21.9 |

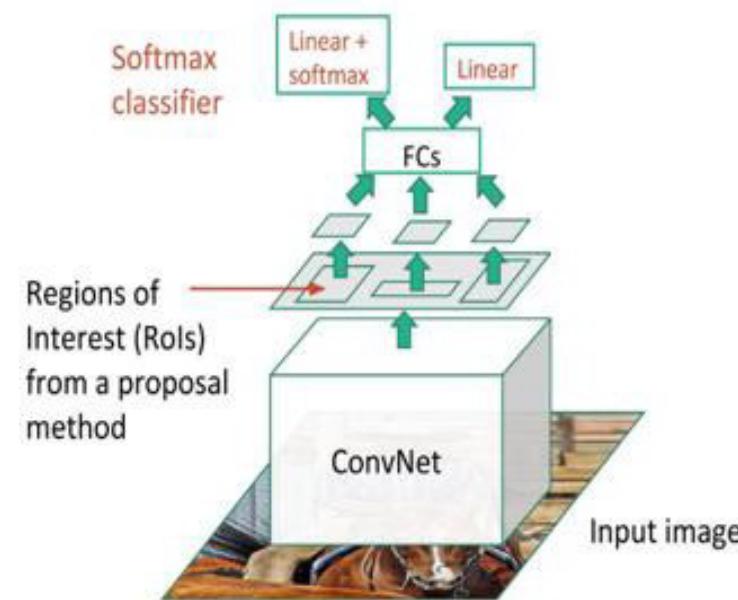
수행 시간 성능

| model | system | conv | proposal | region-wise | total | rate |
|-------|------------------|------|----------|-------------|-------|---------|
| VGG | SS + Fast R-CNN | 146 | 1510 | 174 | 1830 | 0.5 fps |
| VGG | RPN + Fast R-CNN | 141 | 10 | 47 | 198 | 5 fps |
| ZF | RPN + Fast R-CNN | 31 | 3 | 25 | 59 | 17 fps |

RCNN



Fast RCNN



Faster RCNN

