

# CNN for CV

AI for CV Group  
2019



Week 10-11.  
Detection

# Contents:

## I. Two Stage Detection

- A. RCNN: NMS Series
- B. Fast RCNN: ROI Series
- C. Faster RCNN: RPN + Anchor

## II. One Stage Detection

- D. Yolo V1: 1<sup>st</sup> Trial
- E. Yolo V2: Anchor
- F. Yolo V3: FPN
- G. RetinaNet: Focal Loss

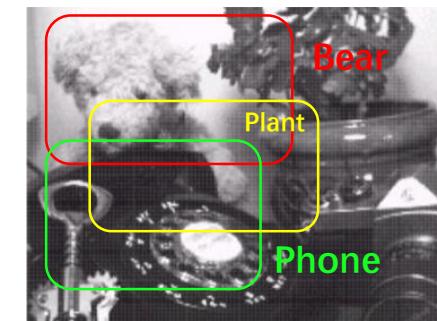
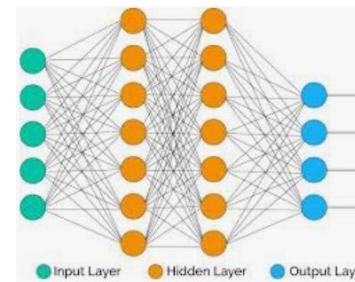
## III. Others

- H. Other methods

# I. Two Stage Detection



# Start from:



Input



Get Manually

Feature



Choose Carefully

Classic ML Tools



Use Hopelessly

Target

**End up with:**



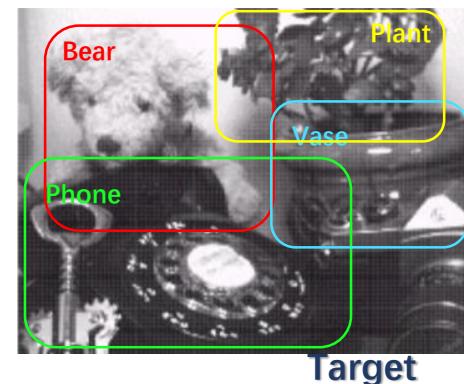
Input



Put In

CNN

Only Focus  
On  
This Part

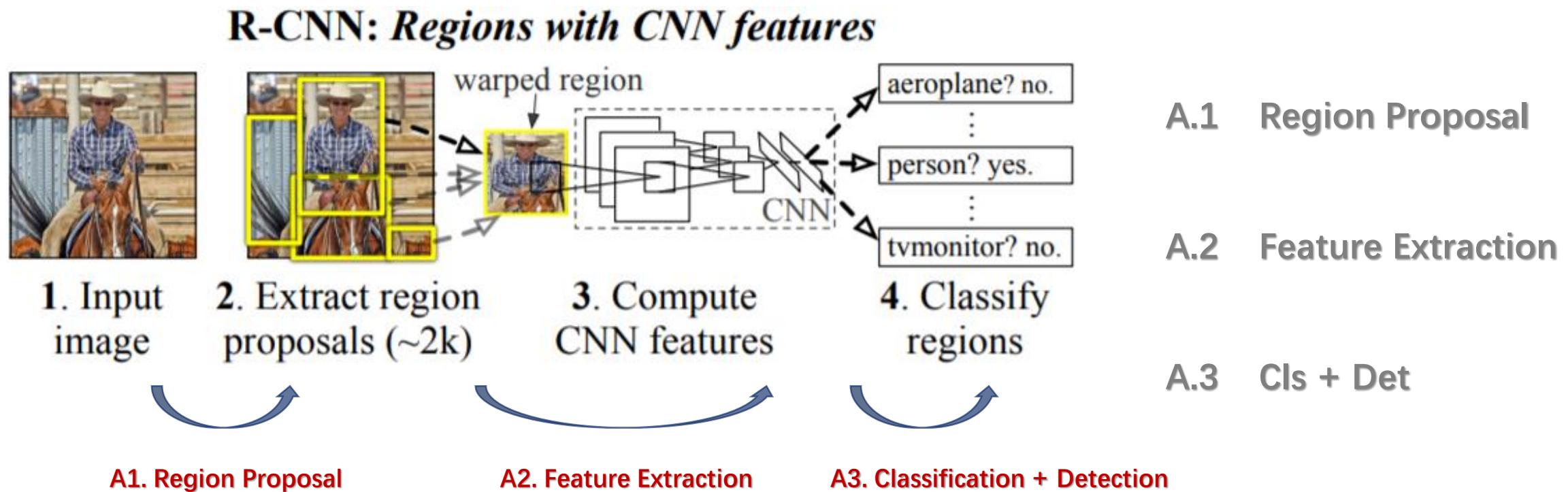


Target

Use It

# I. Two Stage Detection

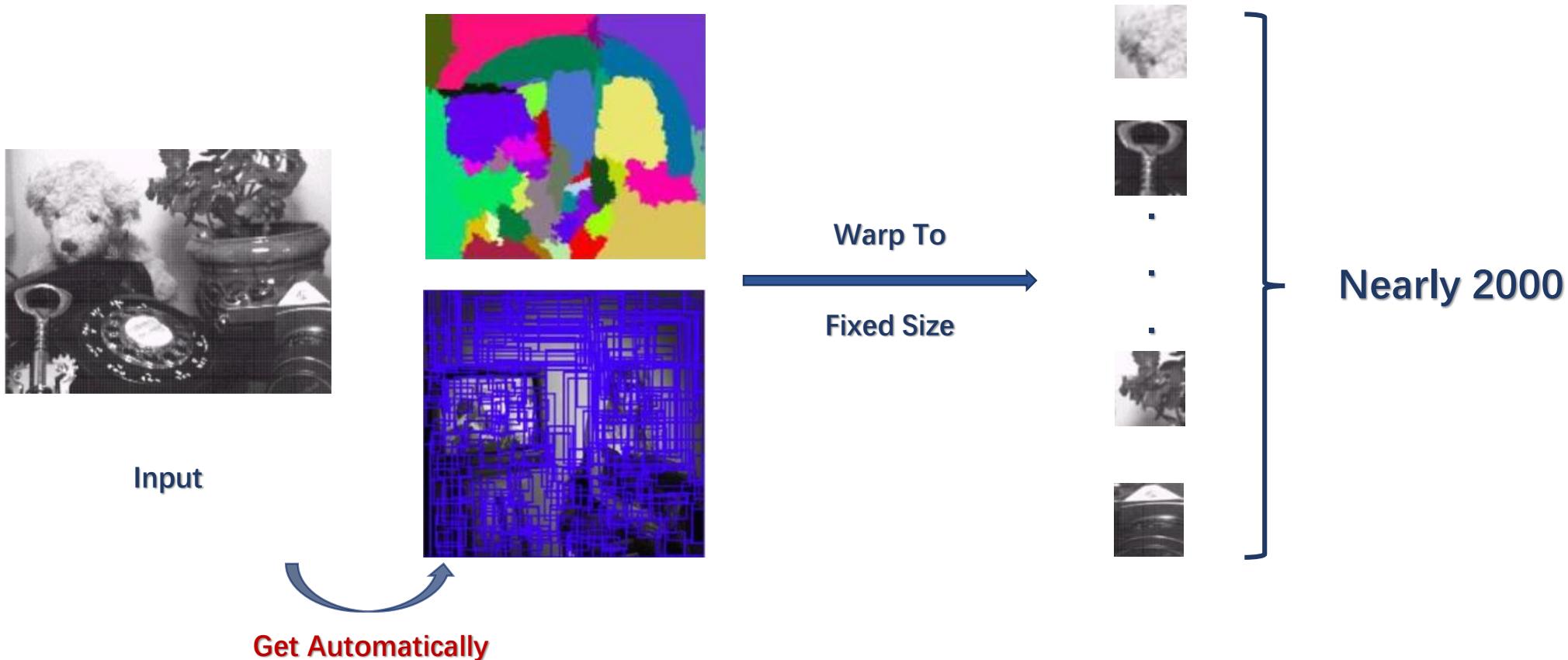
## A. First Trial - RCNN [2014 Ross]:



# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

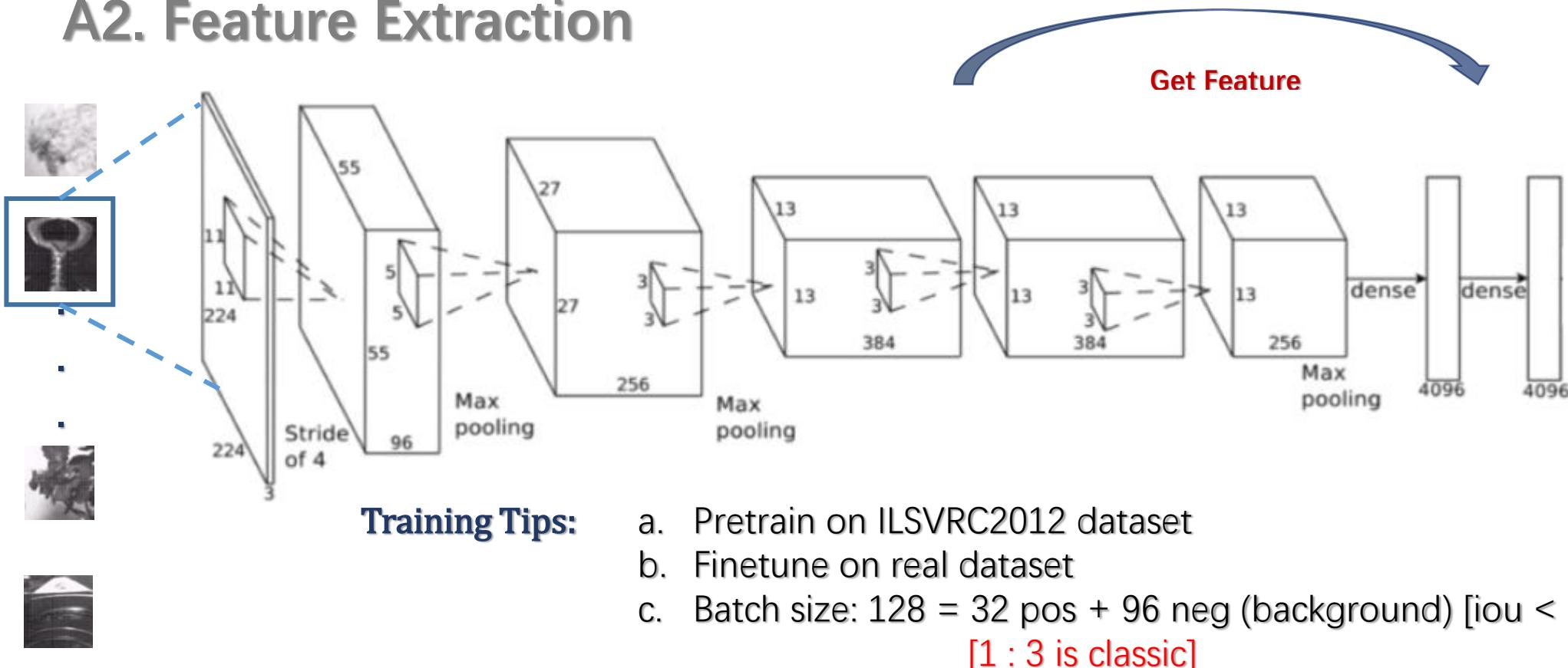
### A1. Region Proposal: Selective Search



# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

### A2. Feature Extraction



# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

### A3. SVM Classification + NMS + BBox Regression

#### a. SVM Classification

- Do classification for each class
  - Pos: Ground truth  
Neg:  $iou < 0.3$  [others ignored]
- [Different from CNN part]

Q: Why you threshold is different here comparing to training AlexNet?

A: We need more data when training CNN. But here more data may be harmful.

Q: Why use SVM? AlexNet has softmax at last, why not use directly?

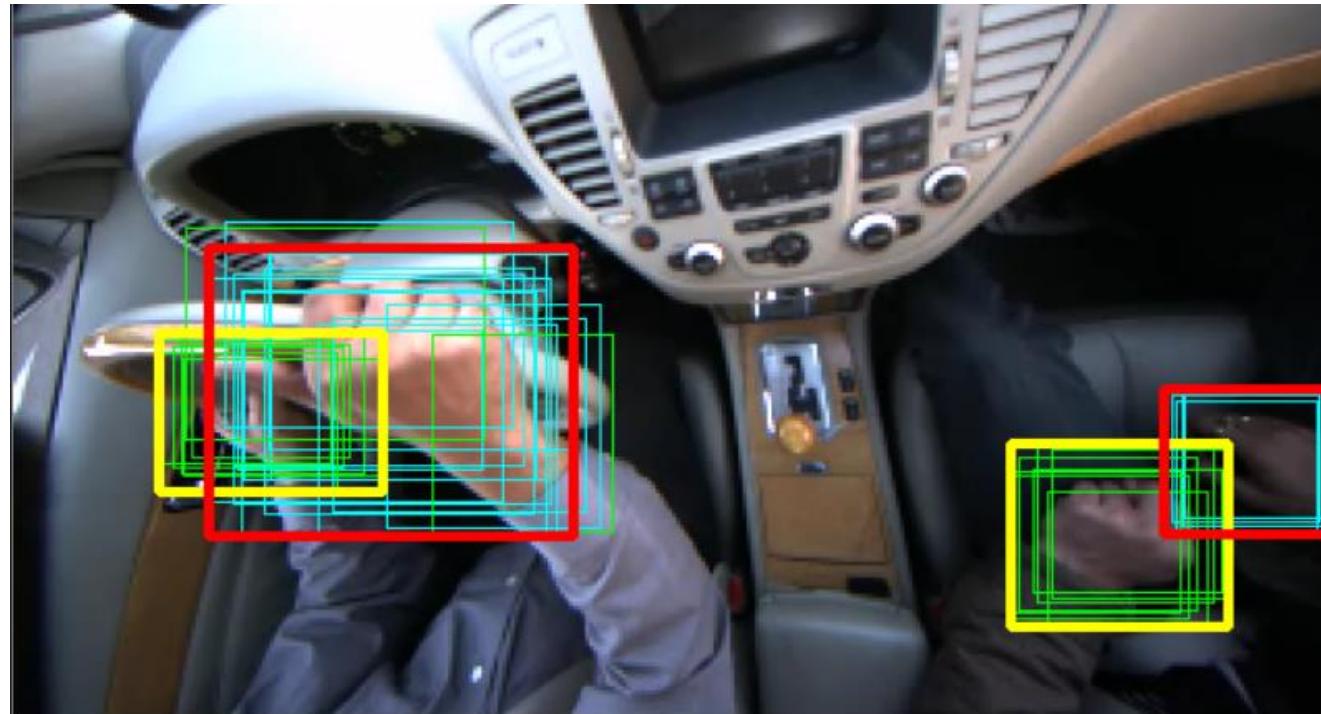
A: AlexNet focuses on classification. Not accurate on localization  
Use hard negative mining when using SVM

# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

**A3. SVM Classification + NMS + BBox Regression**

**b. NMS: Non-Maximum Suppression**

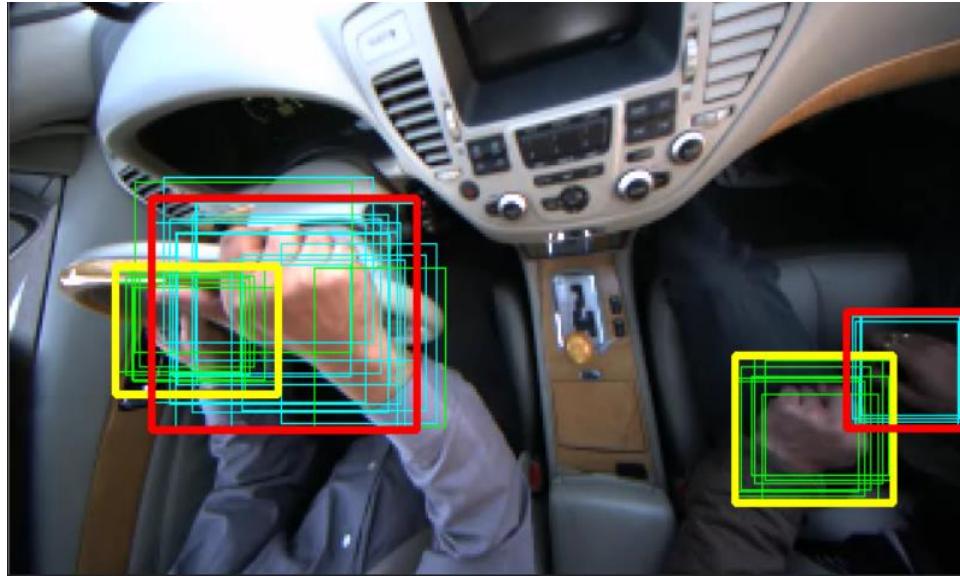


# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

### A3. SVM Classification + NMS + BBox Regression

#### b. NMS: Non-Maximum Suppression



- Sort bbox according to score [here class score]
- Get the 1<sup>st</sup> one. Remove those which has higher iou with it
- Then repeat the last step until there's no bbox

# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

### **A3. SVM Classification + NMS + BBox Regression**

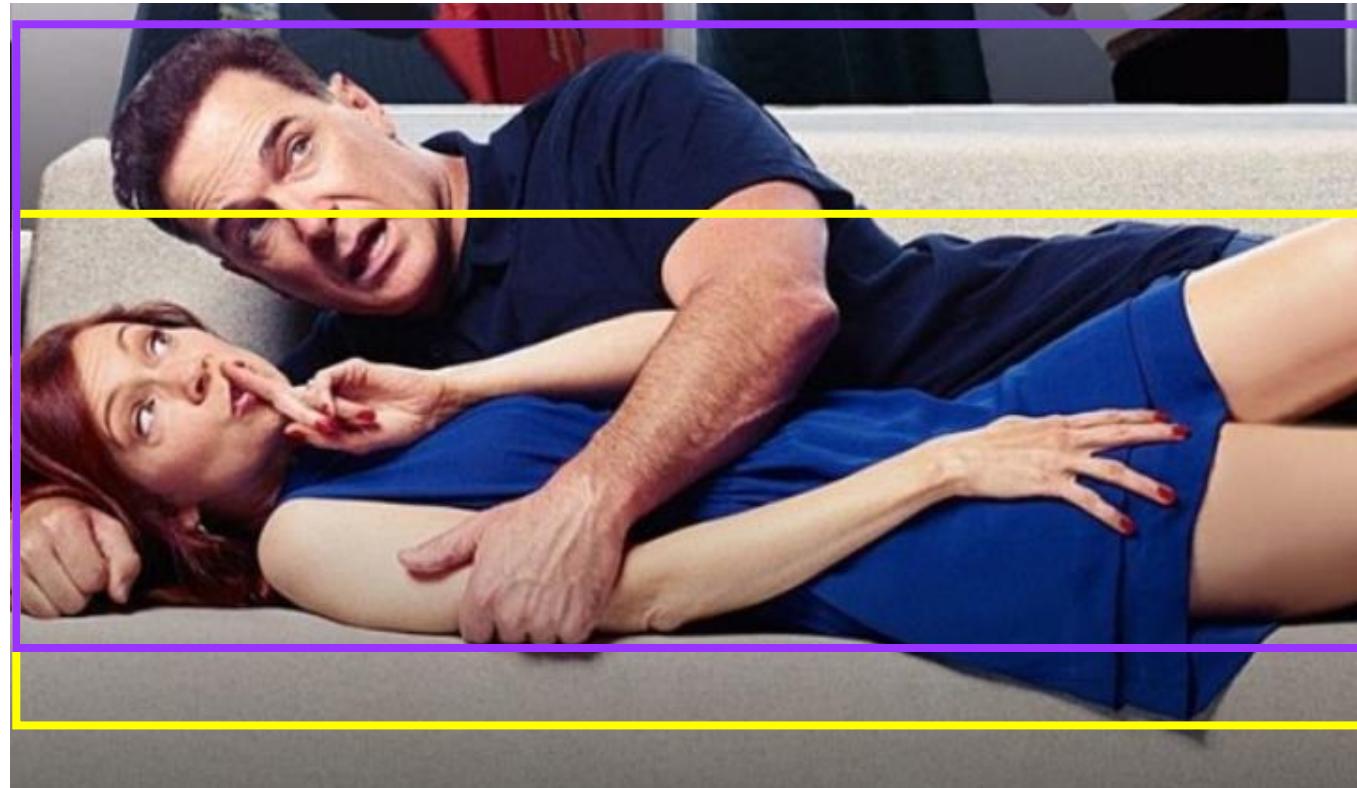
#### **c. Bbox Regression**

- Regression targets of  $(dx, dy, dw, dh)$
- Coordinates need to be normalized
- We'll discuss in detail later

# I. Two Stage Detection

A. First Trial - RCNN [2014 Ross]:

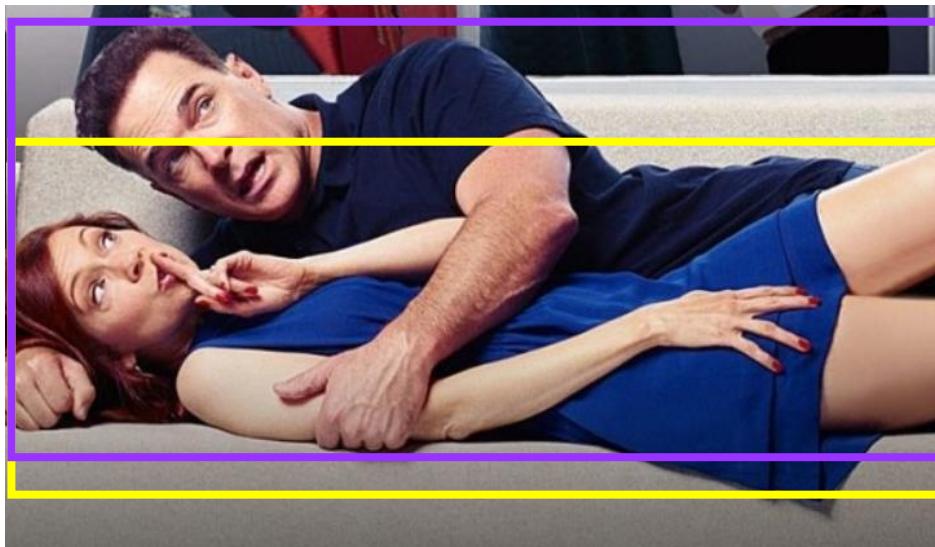
**A\*. Rethinking about NMS**



# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

### A\*. Rethinking about NMS



- If □ has higher score, we'll lose □
- If □ has higher score, we'll lose □
- That's awkward.

**Let's Correct It!**

# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

### A\*. Rethinking about NMS

- Brief history about NMS
  - Traditional NMS [Here we use. Need to know how to code in C++/python]
  - Soft-NMS [2017, better to know]
  - IoU-Net [2018, Localization Confidence + PrROI Pooling, better to know]
  - Softer-NMS [2019, KL-Loss + Softer-NMS, better to know]

Important / Will cover / Remind

# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

### A\*. Rethinking about NMS

- Soft-NMS

**Input :**  $\mathcal{B} = \{b_1, \dots, b_N\}$ ,  $\mathcal{S} = \{s_1, \dots, s_N\}$ ,  $N_t$   
 $\mathcal{B}$  is the list of initial detection boxes  
 $\mathcal{S}$  contains corresponding detection scores  
 $N_t$  is the NMS threshold

```
begin
     $\mathcal{D} \leftarrow \{\}$ 
    while  $\mathcal{B} \neq \text{empty}$  do
         $m \leftarrow \text{argmax } \mathcal{S}$ 
         $\mathcal{M} \leftarrow b_m$ 
         $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{M}; \mathcal{B} \leftarrow \mathcal{B} - \mathcal{M}$ 
        for  $b_i$  in  $\mathcal{B}$  do
            if  $iou(\mathcal{M}, b_i) \geq N_t$  then
                |  $\mathcal{B} \leftarrow \mathcal{B} - b_i; \mathcal{S} \leftarrow \mathcal{S} - s_i$ 
            end
             $s_i \leftarrow s_i f(iou(\mathcal{M}, b_i))$ 
        end
    end
    return  $\mathcal{D}, \mathcal{S}$ 
end
```

NMS

Soft-NMS

# I. Two Stage Detection

## A. First Trial - RCNN [2014 Ross]:

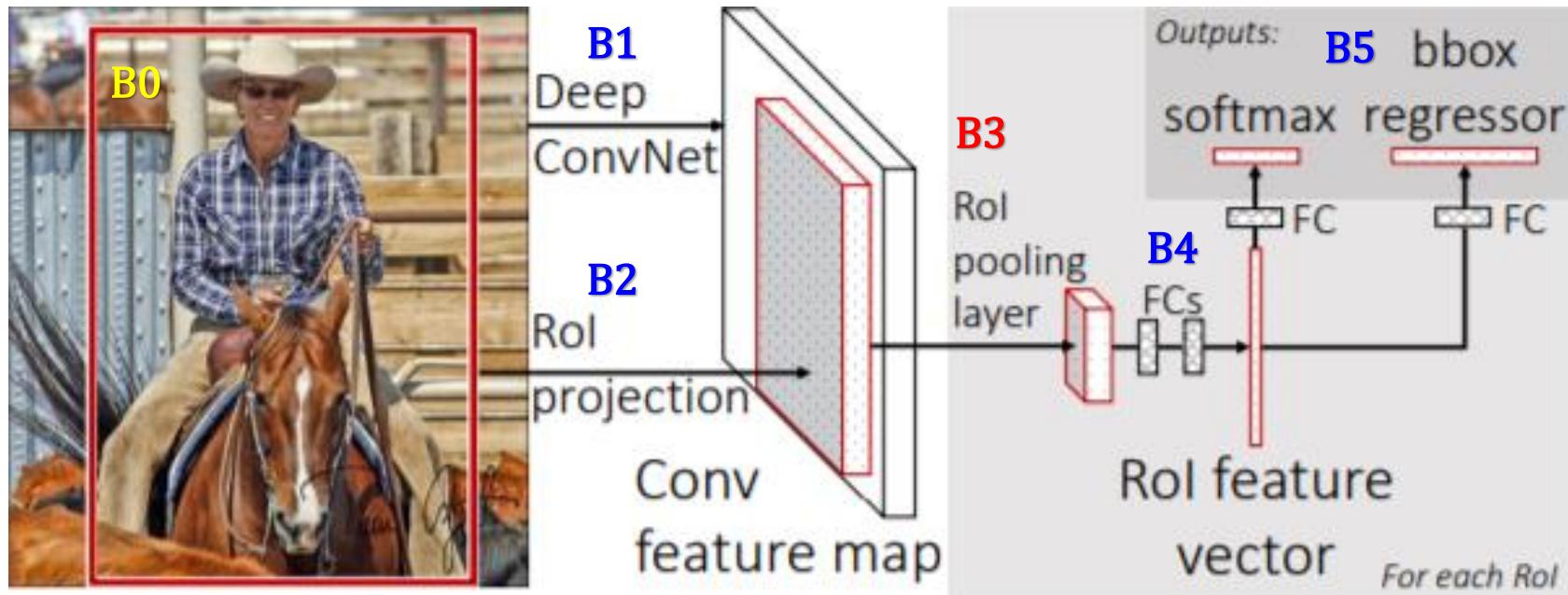
### A4. RCNN Problems

- Slow: need to run full forward pass of CNN for each region proposal
- SVMs and regressors are post-hoc: CNN features not updated in response to SVMs and regressors
- Complex multistage training pipeline

**Let's Fix It, Step by Step**

# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

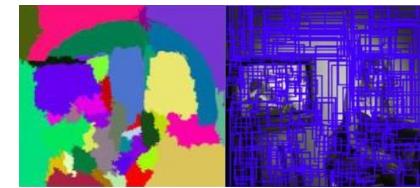


# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

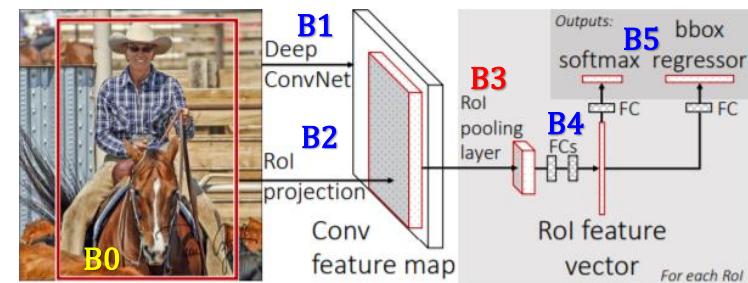
### B0. Region Proposal: Same as RCNN

- 2k per image. Record location for each ROI



### B1 & 2. Convolution & Projection

- Do Conv per ROI. Project location for each ROI
- 3 basic structures provided, use Vgg16 as an E.g.
- 4 max pooling /16



### B3. ROI Pooling

- Grid each ROI in feature map to fixed size, and do max pooling within each grid
- So different size of feature maps can transfer into feature maps with same size.
- $bs = 2; 64$  ROIs / image -> 128 proposals;
- $\frac{1}{4}$  pos &  $\frac{3}{4}$  neg
- IoU > 0.5 pos, [0.1, 0.5) neg, [0.0, 0.1) hard neg
- Train on pos & neg. Test use neg. Retrain bad case

# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

### B4. FC Layers

- FC layer cost a lot. Can use SVD to accelerate

### B5. Multi-task Loss

- Classification:  $n + 1$ (background), softmax + cross entropy,  $L_{cls}$
- Localization: Offset, SmoothL1,  $L_{loc}$
- Total:  $L(p, u, t^u, t^v) = L_{cls}(p, u) + \lambda[u \geq 1] L_{loc}(t^u, v)$

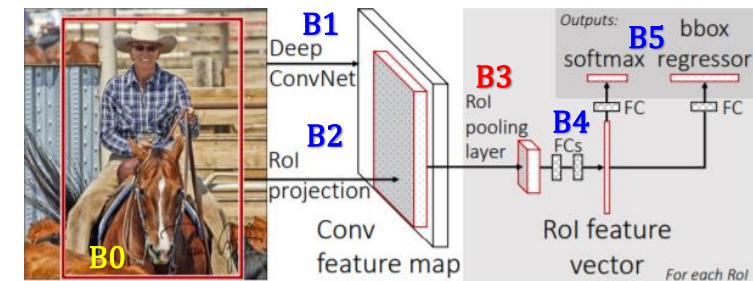
$p$ : predicted class score

$u$ : real class

$t^u$ : generated bbox [offset]

$t^v$ : real bbox

$$\lambda = \begin{cases} 1, & u \geq 1 \\ 0, & u = 0 \end{cases}$$

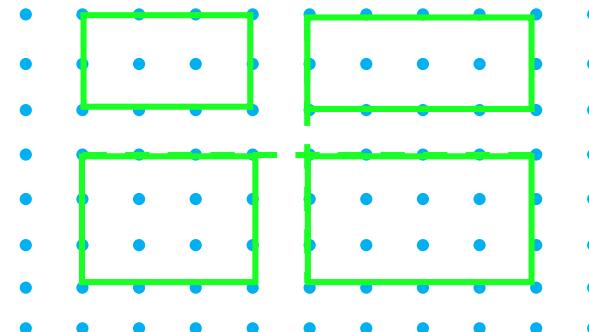
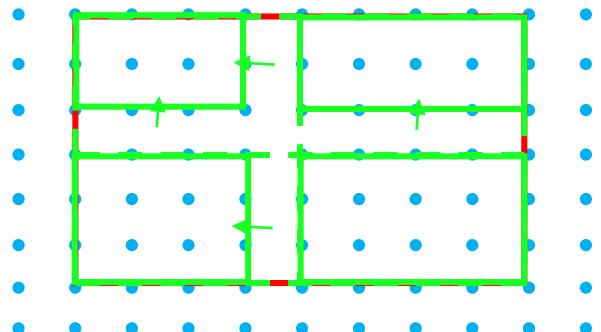
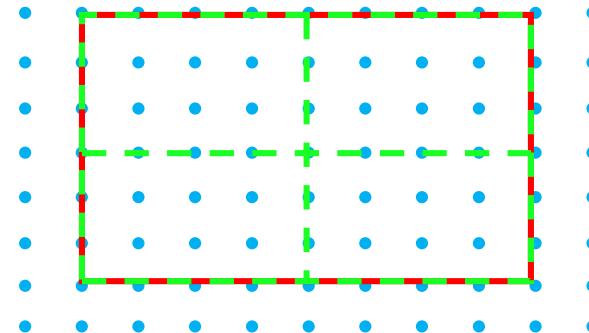
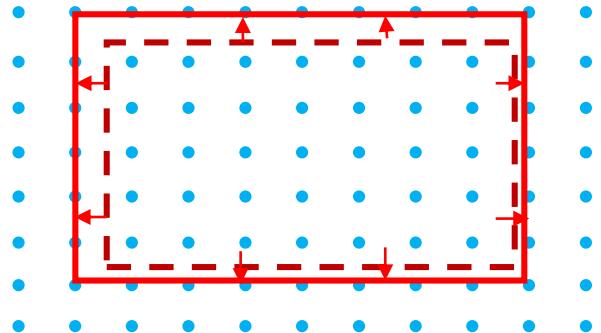


SmoothL1 & Offset regression  
will be discussed in next

# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

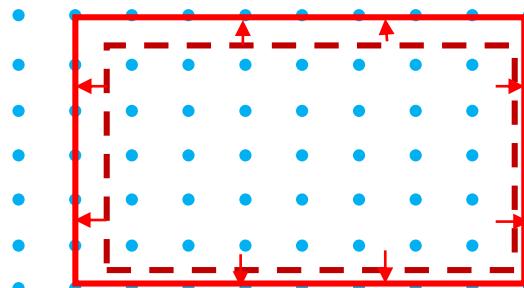
### B\*. Rethinking ROI Pooling



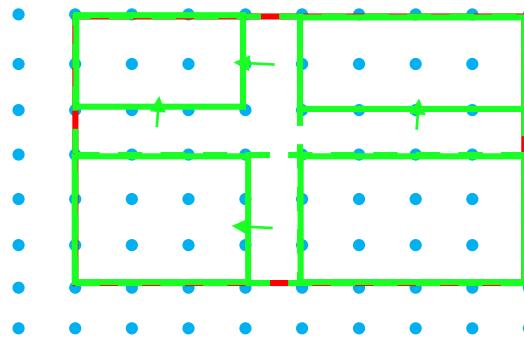
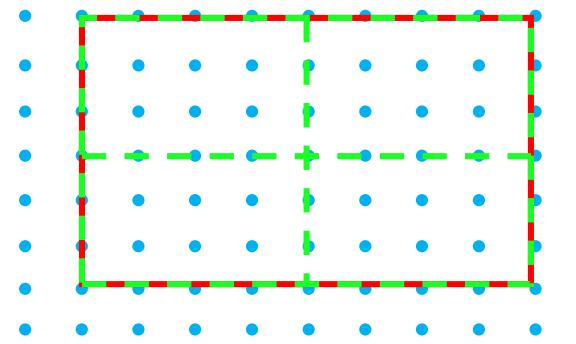
# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

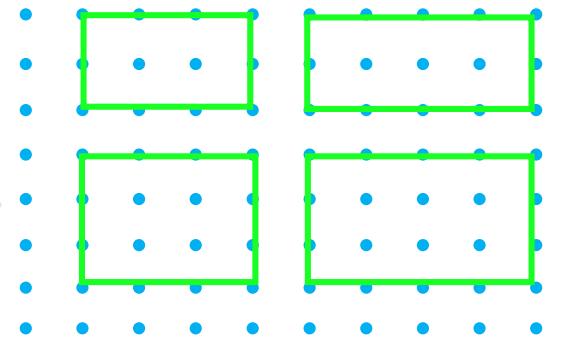
### B\*. Rethinking ROI Pooling



Quantized  
2 Times!



0.8 pixels in feature  
> 10 pixels in image  
Bad for small objects



# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

### **B\*. Rethinking ROI Pooling**

- BP for ROI Pooling

$$\frac{\partial L}{\partial x_i} = \sum_r \sum_j [i = i^*(r, j)] \frac{\partial L}{\partial y_{r,j}}$$

$y_{r,j}$ : rth region, jth feature point

$i^*(r, j)$ : the position where  $y_{r,j}$  comes from

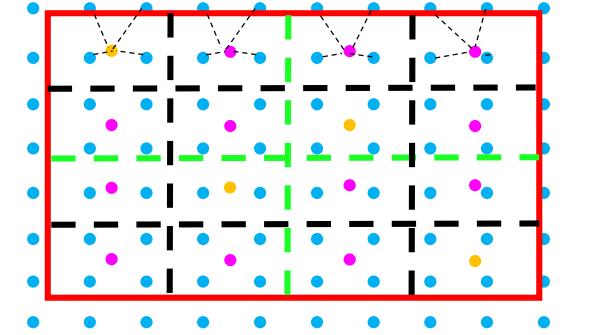
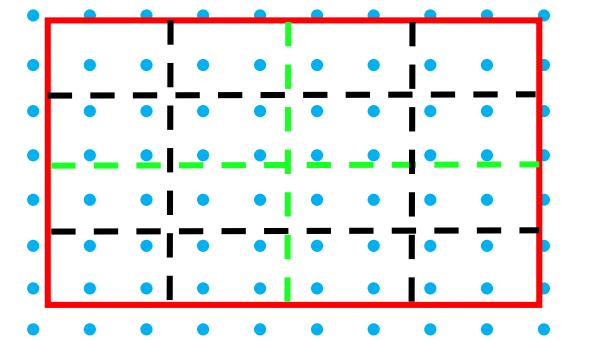
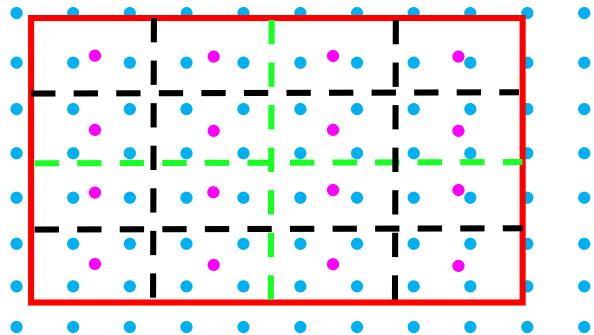
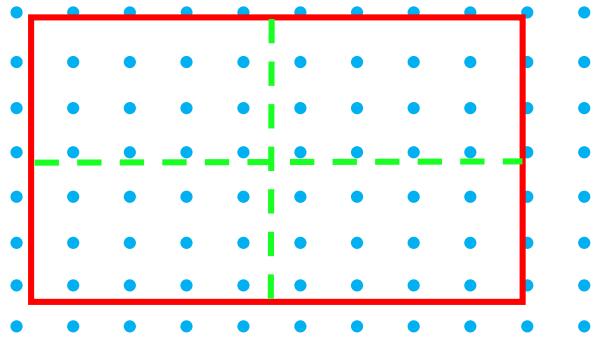
Just one point has loss passed by in each bin.

# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

### **B\*. Rethinking ROI Pooling**

- ROI Align [2017, Kaiming]

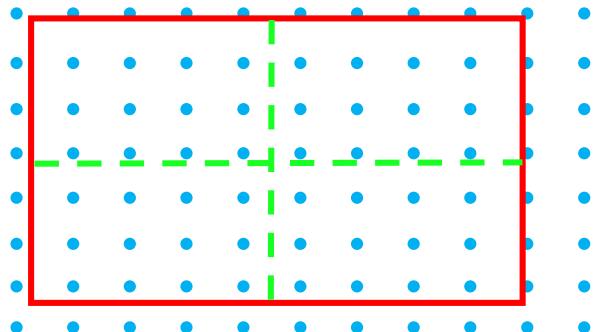


# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

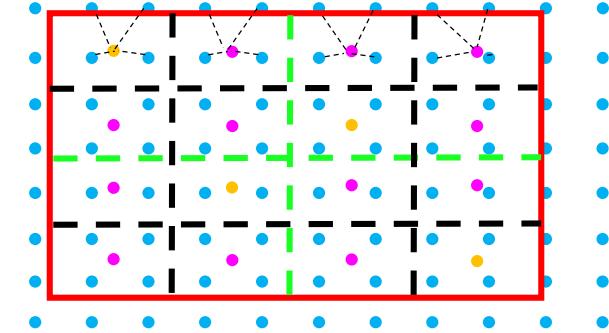
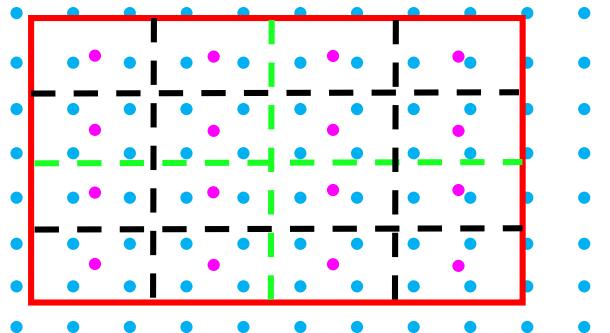
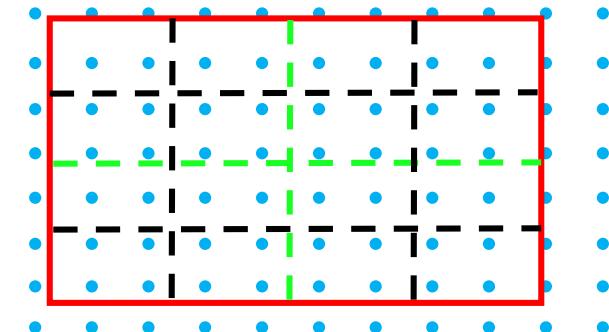
### B\*. Rethinking ROI Pooling

- ROI Align [2017, Kaiming]



N is a param

Not all feature points have contributions



# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

### **B\*. Rethinking ROI Pooling**

- ROI Align [2017, Kaiming] – BP

$$\frac{\partial L}{\partial x_i} = \sum_r \sum_j [d(i, i^*(r, j))] (1 - \Delta h) (1 - \Delta w) \frac{\partial L}{\partial y_{r,j}}$$

$d(x, y)$ : distance between  $x$  &  $y$

$\Delta w, \Delta h$ : width & height difference of  $x_i$  &  $x_{i^*(r,j)}$

$x_{i^*(r,j)}$ : a float value interpolated by the forward pass

$x_i$ : feature point before pooling

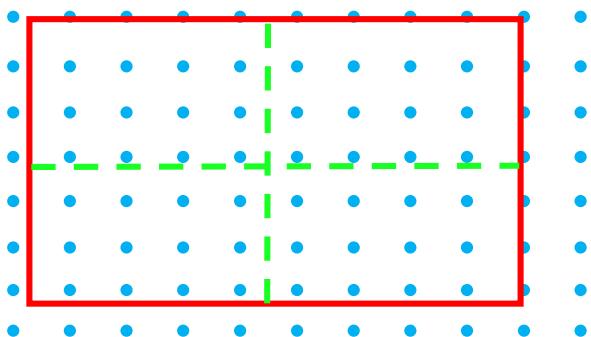
Only several points have losses passed by in each bin.

# I. Two Stage Detection

## B. Fast R-CNN [2015 Ross]:

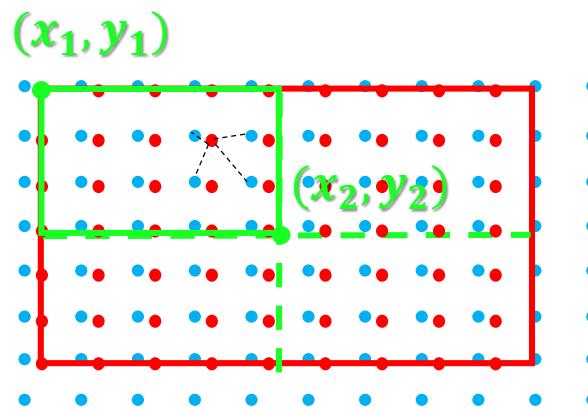
### B\*. Rethinking ROI Pooling

- Precise ROI Pooling [2018, IoU-Net]



$$f(x, y) = \sum_{i,j} IC(x, y, i, j) * w_{i,j}$$

$$IC(x, y, i, j) = \max(0, 1 - |x - i|) + \max(0, 1 - |y - j|)$$



Do Pr ROI Pool for  $\square$ :

$$PrPool(bin, F) = \frac{\int_{y_1}^{y_2} \int_{x_1}^{x_2} f(x, y) dx dy}{(x_2 - x_1) * (y_2 - y_1)}$$

BP:

$$\frac{\partial PrPool(bin, F)}{\partial x_1} = \frac{PrPool(bin, F)}{x_2 - x_1} - \frac{\int_{y_1}^{y_2} f(x, y) dy}{(x_2 - x_1) * (y_2 - y_1)}$$

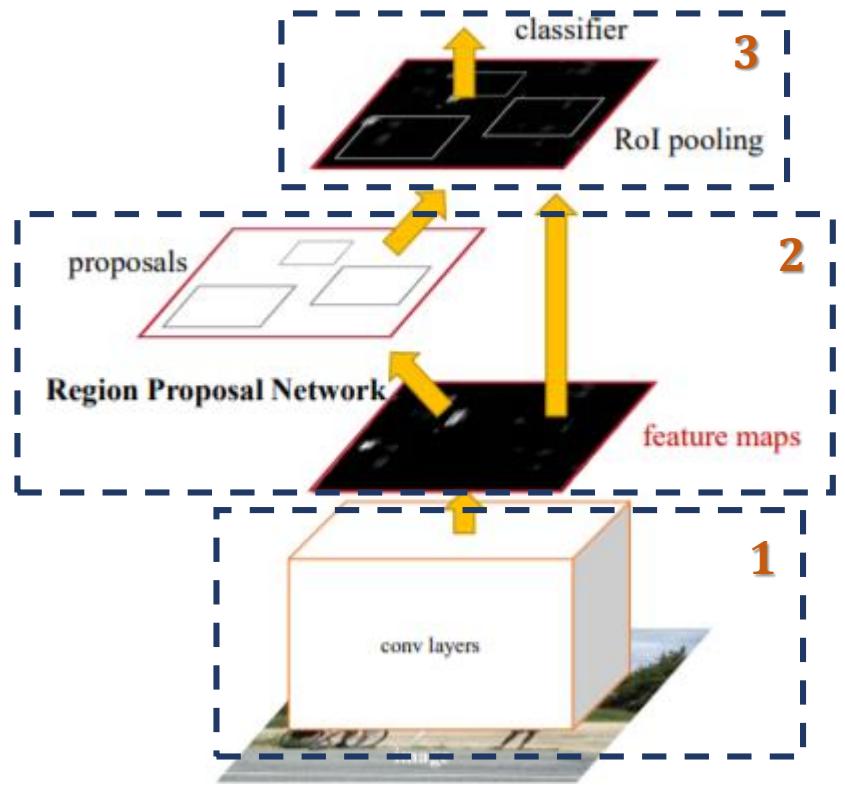
No Quantization

No Parameter

Pass losses for all

# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:



### C0. Structure

3. Fast RCNN: ROI + { Classification  
Regression }

2. RPN

1. Backbone

Input

# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

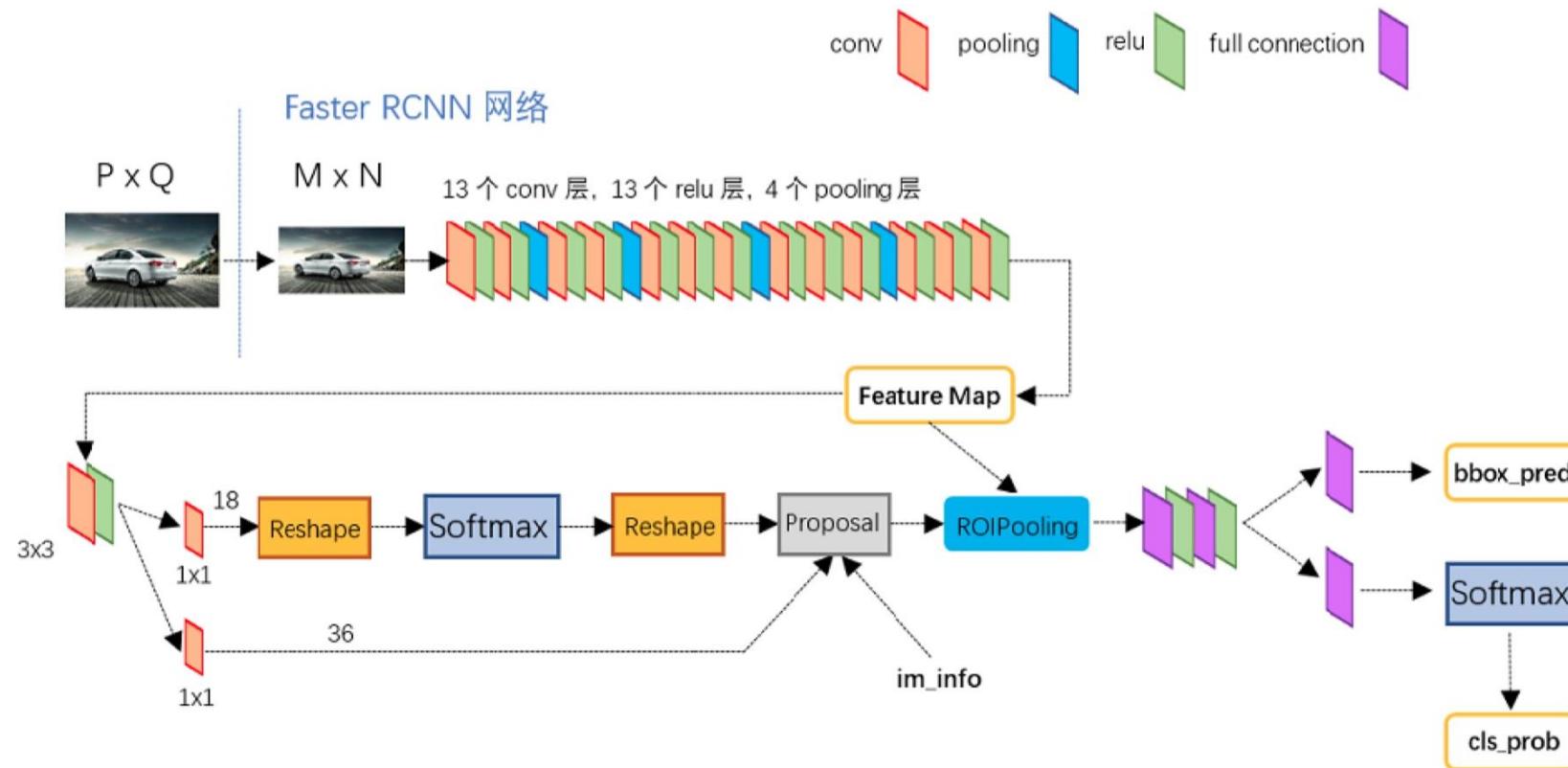
### C1. Backbone

- **Aim:** Extract feature integratedly
  - No need to extract feature per ROI
  - Use generated feature to generate proposals
- **Structure:** ZF / Resnet / VGG16 (13 conv + 13 ReLU + 4 Pooling)  
[Do have the sense of splicing network from now on]
- **Output:**  $B \times C \times H/16 \times W/16$   
 $1 \times 256 \times 38 \times 50$   
[Feature Map]

# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

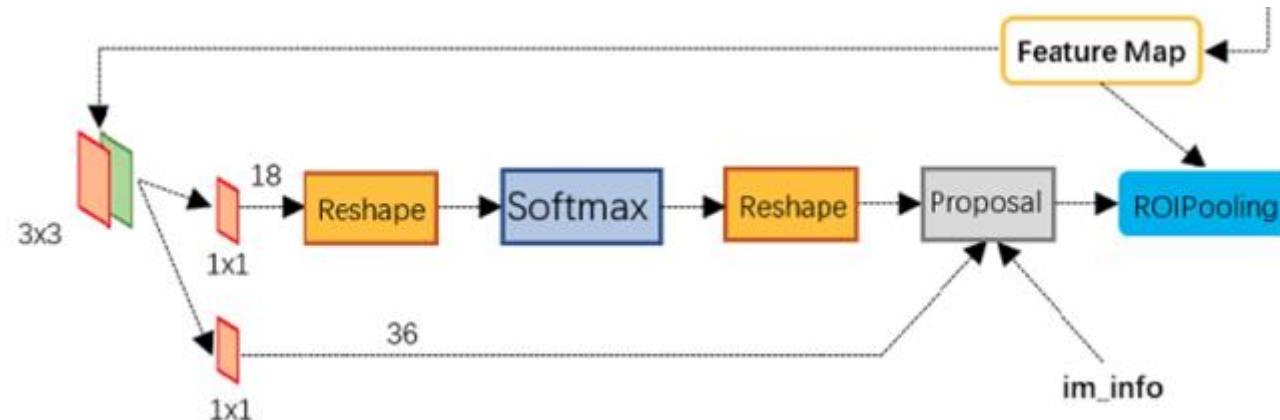
### C1. Backbone



# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN



# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Aim: Generate region proposal  
[Real detected objects are coming from those RPs]
  - a. It's the reason where the name "Two-Stage" coming from:  
RPN + Bbox Regression
  - b. It's the reason why some argue two-stage detection has better results than one-stage methods.

# I. Two Stage Detection

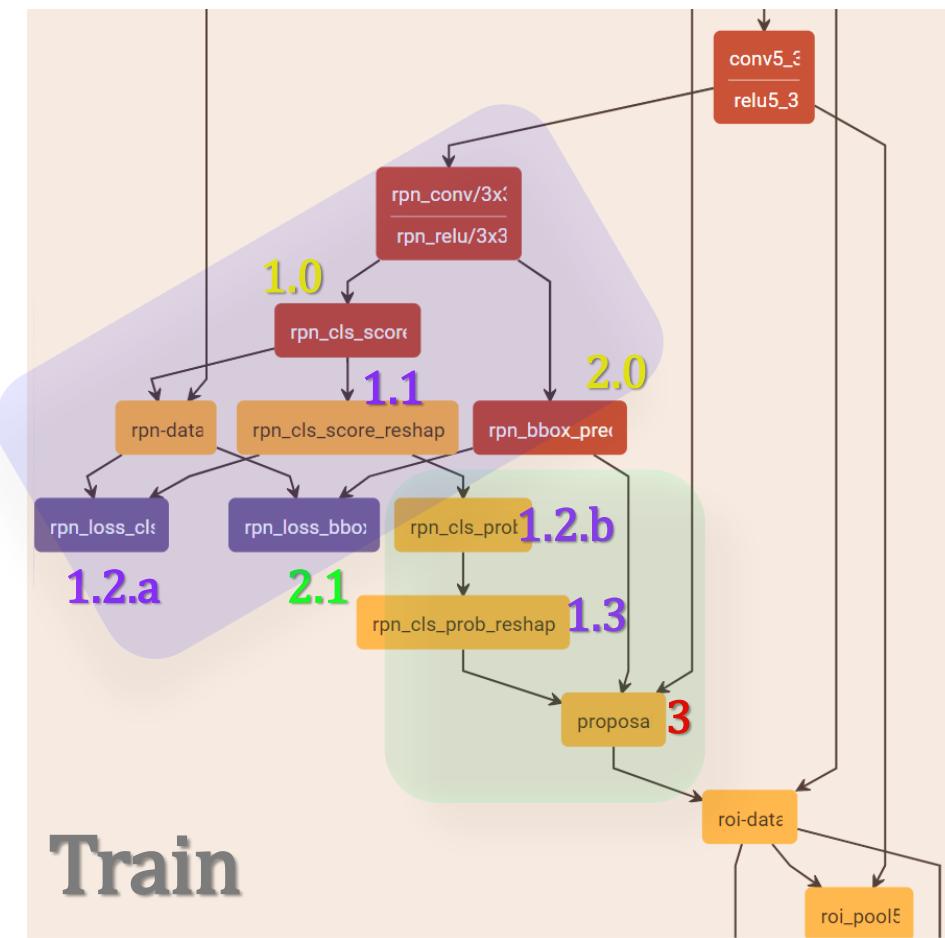
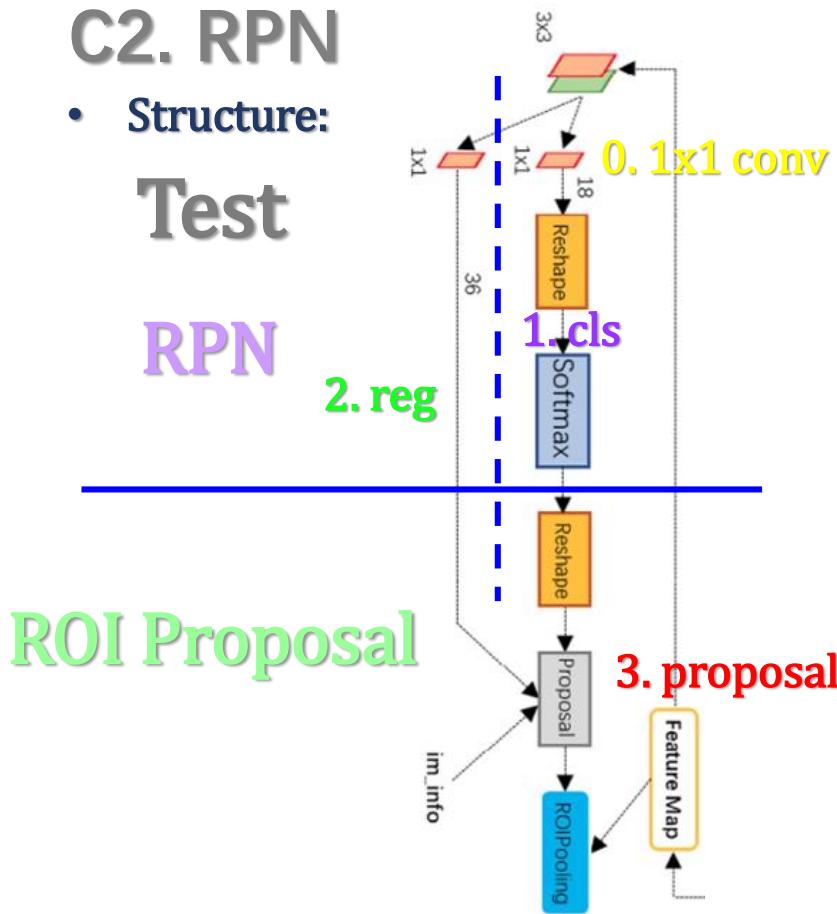
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Output of RPN:
  - a. ROIs:  $128 \times 5$   
[ $0, x_1, y_1, x_2, y_2$ ] -> physical region proposal
  - b. Label:  $128$   
[ $0 \sim 20$ ] -> ROIs' classifications
  - c. bbx\_target:  $128 \times 84$   
[( $20 + 1$ )  $\times 4$ ] -> targets for bounding box regression
  - d. bbx\_weight:  $128 \times 84$   
[0 or 1] -> weights for bbx\_target when box regressing

# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:



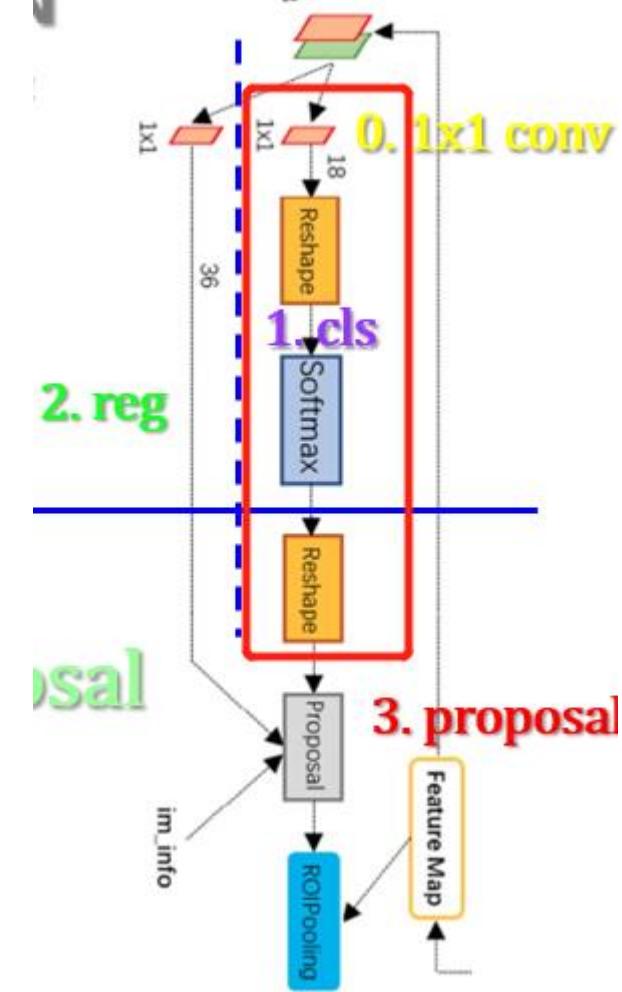
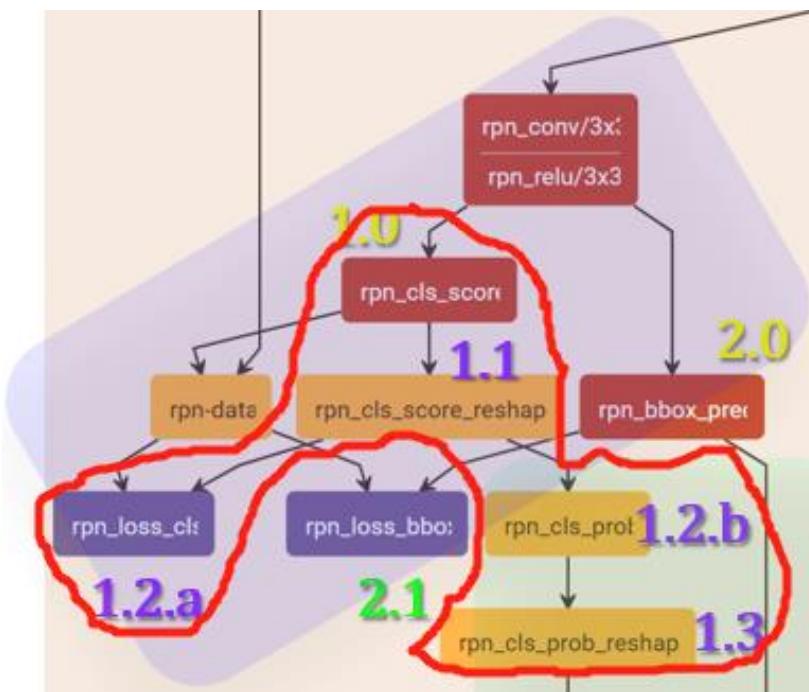
# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

#### C2.1: classification branch



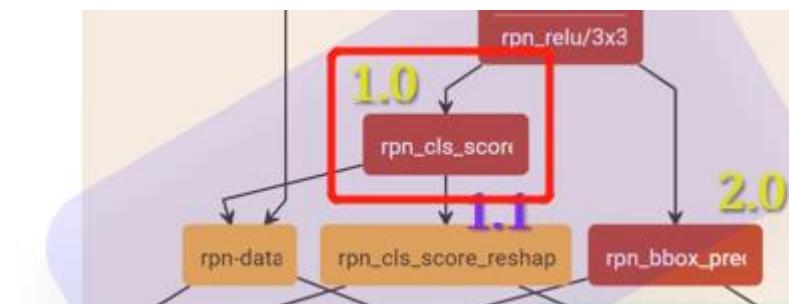
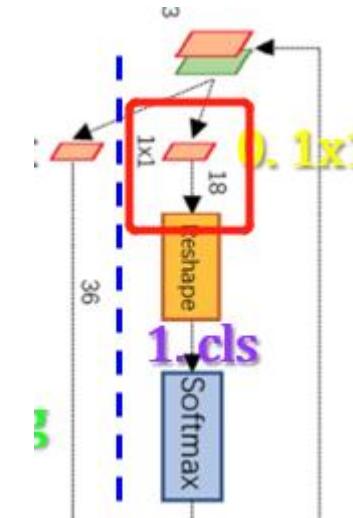
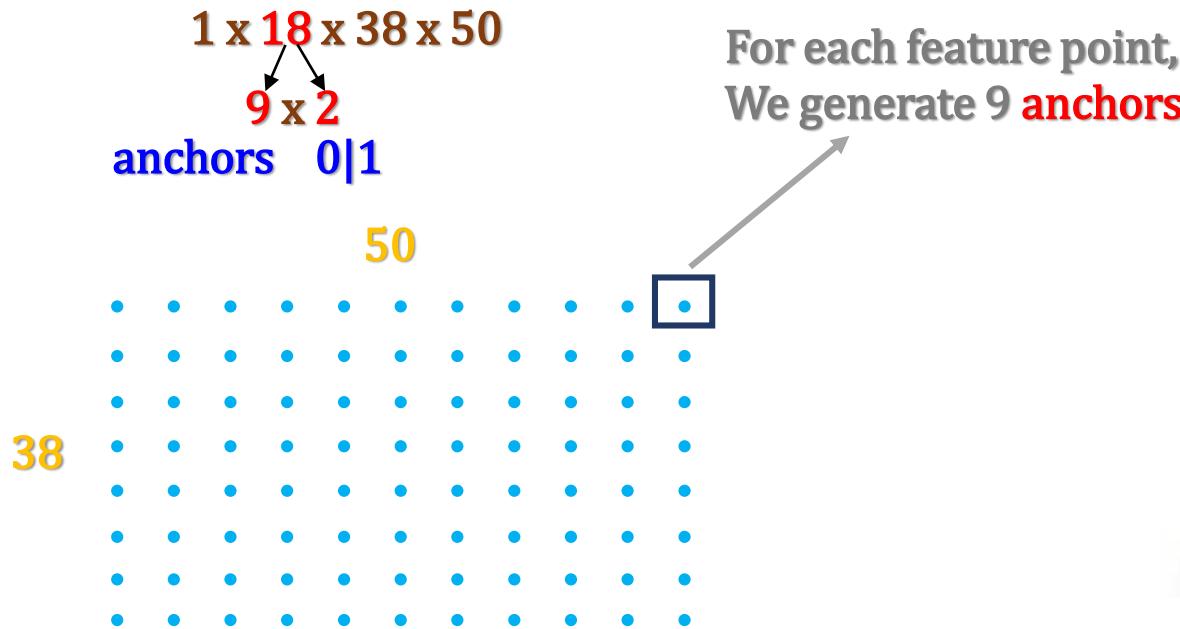
# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

#### C2.1.0: 1x1 conv



# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure - **Anchor:**

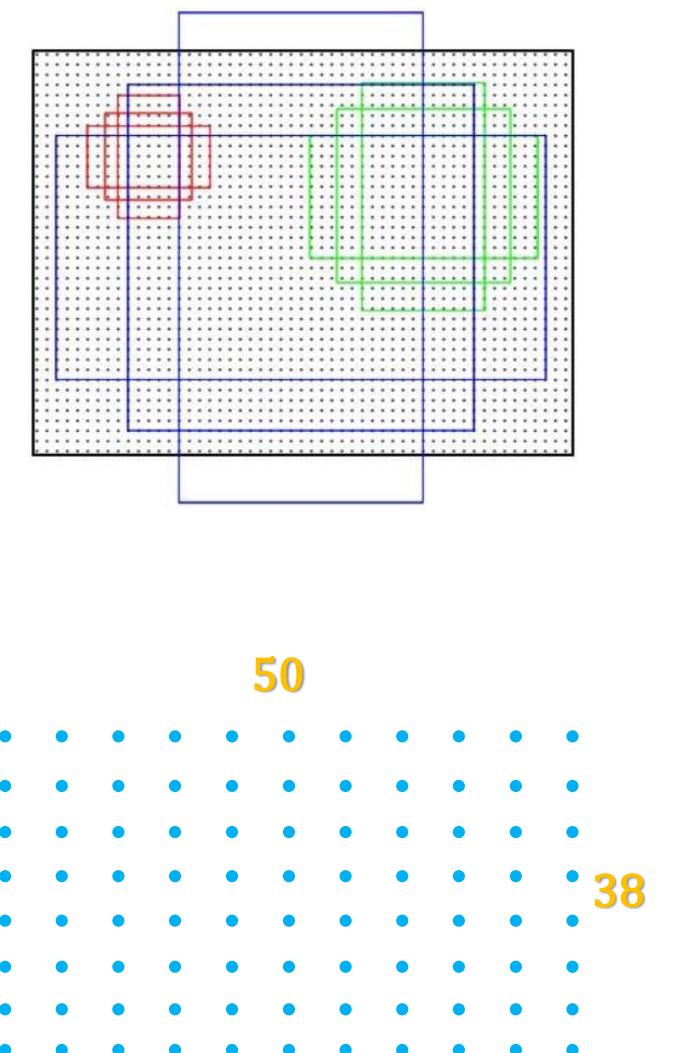
- ❑ Represent an area in original image  
(a.k.a. where objects are)

- ❑ Have different scale & ratio  
to cover all kinds of objects

- ❑ 9 in total: 3 scales x 3 ratios

- ❑  $1 \times 9 \times 38 \times 50 = 17100$  anchors

- ❑ Coords of anchors are of original imgs



# I. Two Stage Detection

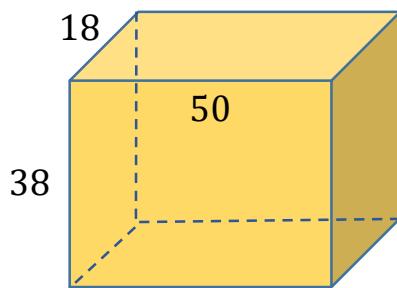
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

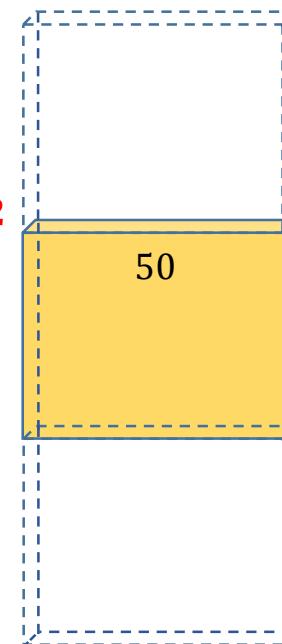
- Structure:

C2.1.1: **cls reshape**

**Input:**

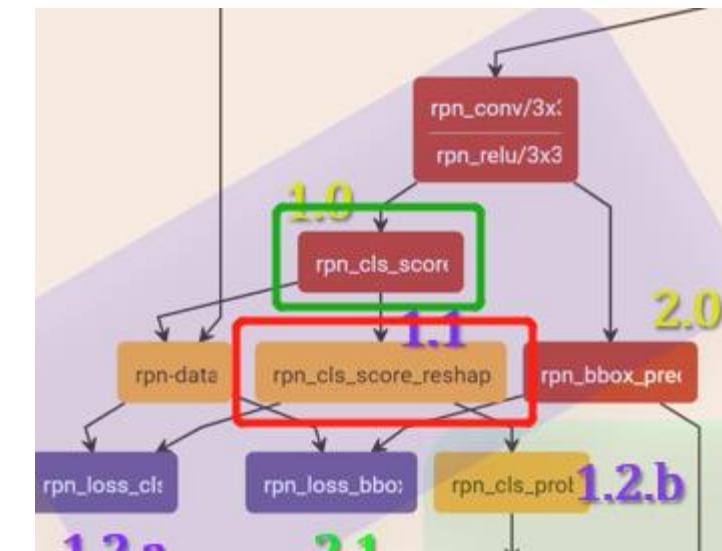
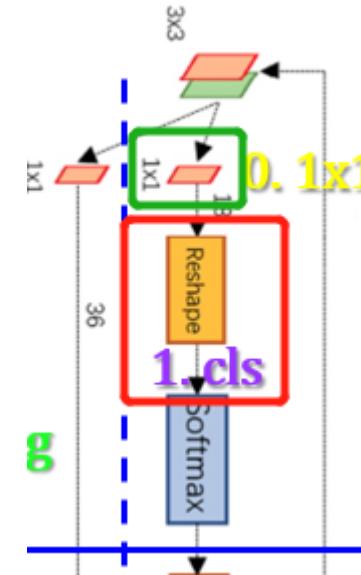


→ reshape



**Output:**

This is for doing **Softmax** to get **fore/background**



# I. Two Stage Detection

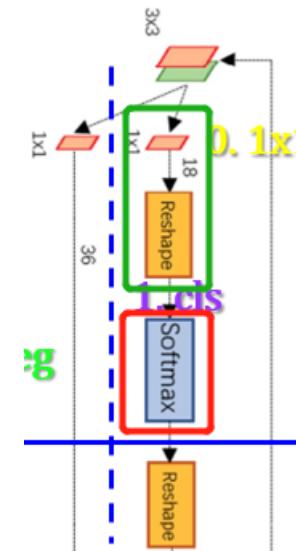
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

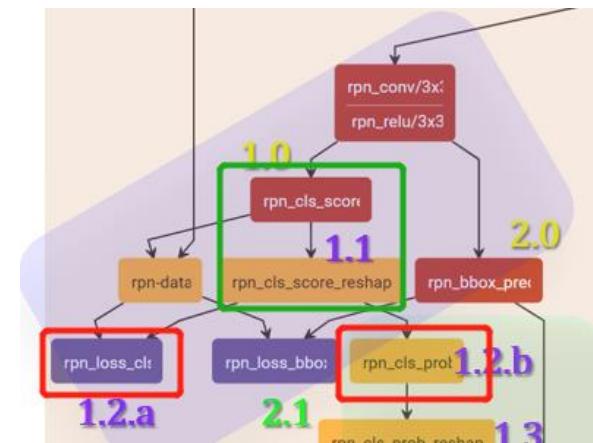
#### C2.1.2: Softmax

This is for doing Softmax  
to get fore/background



2 branches: a. bp loss & b. no loss / get score

- a. 1 foreground & 0 background.  
Get loss with anchor and pass back.  
Use IoU
- b. Just do classification and get the score,  
Use the score to select proposal



# I. Two Stage Detection

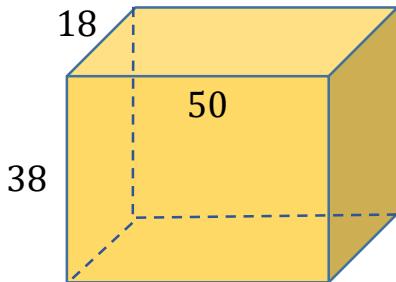
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

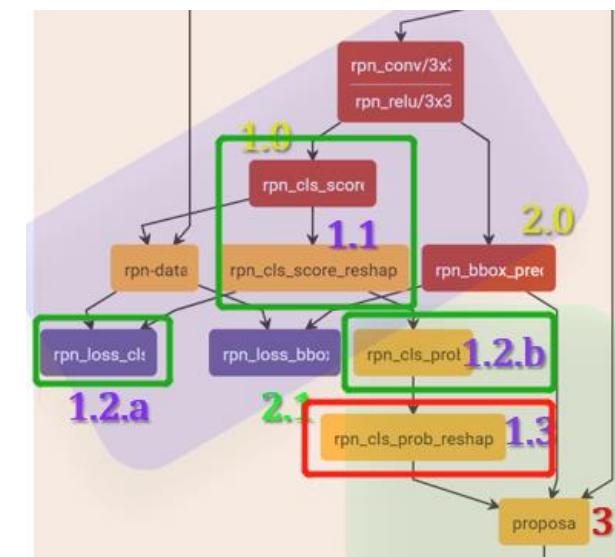
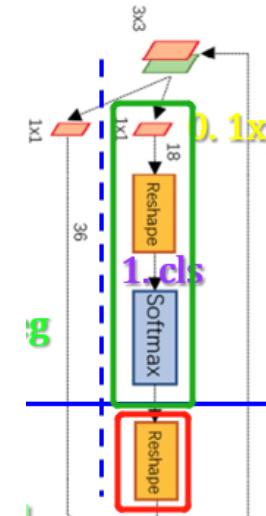
- Structure:

C2.1.3: **cls reshape 2**

**Back to original shape**



Each voxel represents the possibility  
of being a foreground or a background  
of an anchor



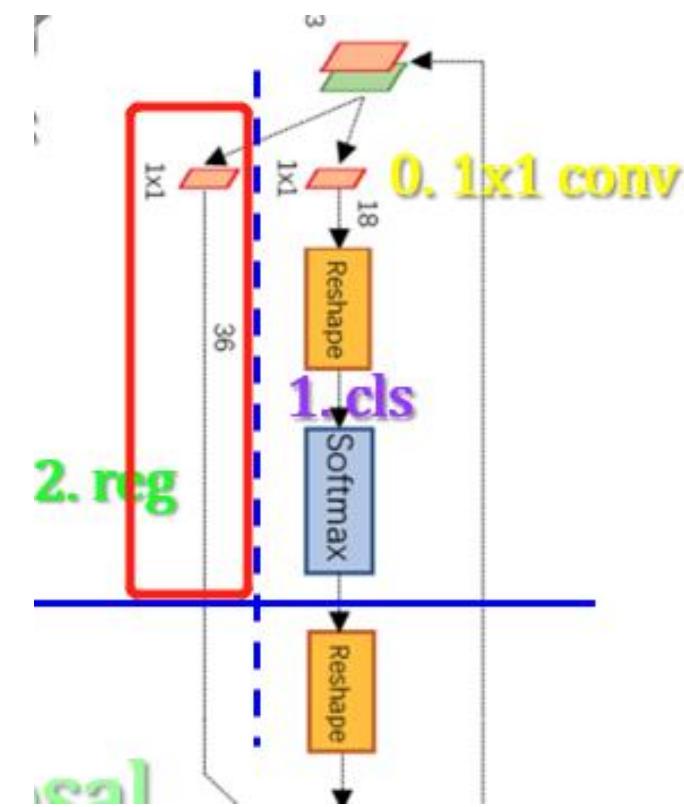
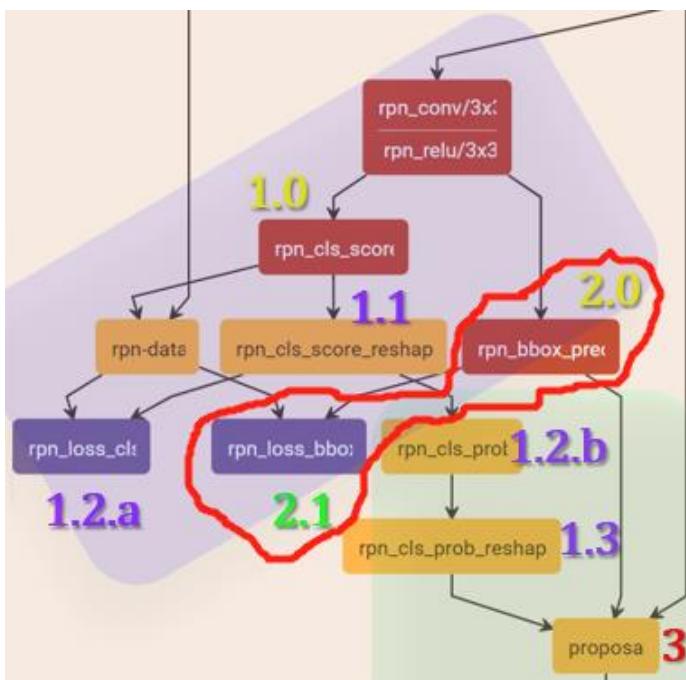
# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

#### C2.2: regression



# I. Two Stage Detection

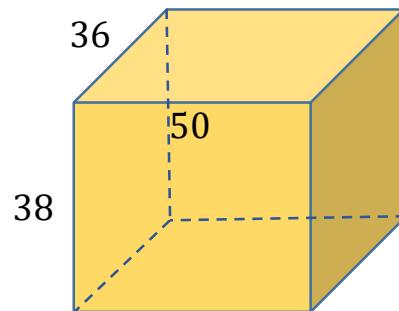
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

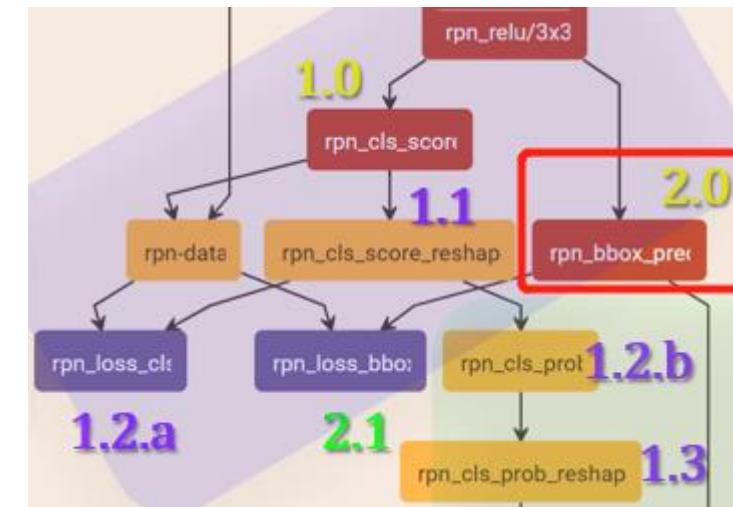
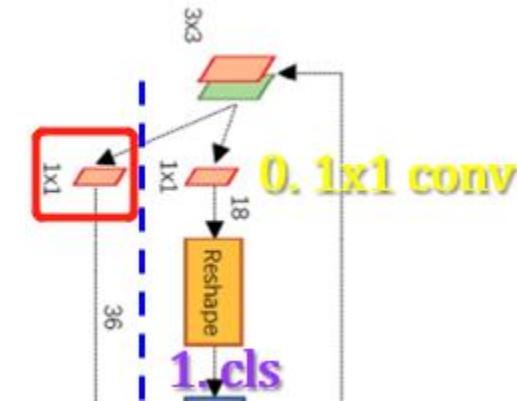
C2.2.0: 1x1 conv

$1 \times 36 \times 38 \times 50$   
9 x 4  
anchors 4 targets



Rgress anchor to gt\_bbox

Targets: proposals



# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

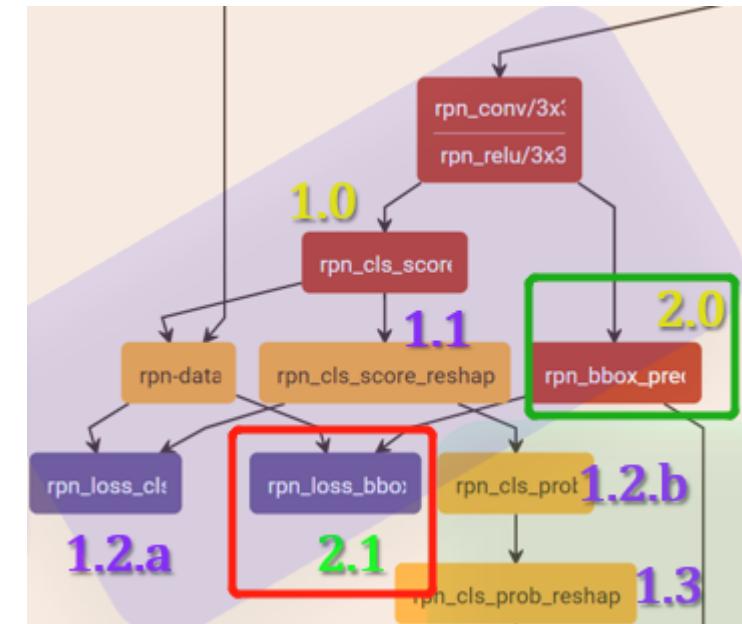
#### C2.2.1: Smooth L1 Loss

- Need to multiply the mask

$$\begin{cases} \text{anchor} = 1, \text{mask} = 1 \\ \text{anchor} = 0, -1, \text{mask} = 0 \end{cases}$$

- Smooth L1 Loss

$$f(x) = \begin{cases} \frac{(\sigma x)^2}{2}, & \text{if } x < 1/\sigma^2 \\ |x| - \frac{0.5}{\sigma^2}, & \text{otherwise} \end{cases}$$



# I. Two Stage Detection

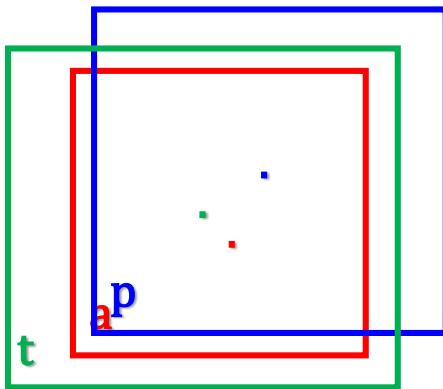
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

C2.2.1: Smooth L1 Loss

➤ Regress (center) offset, NOT COORDINATES



We hope  $a \rightarrow t$ ,  
but actually  $a$  generates  $p$ .

So as long as  $p \rightarrow t$ ,  
We get a good result.

So as long as the offset of  
 $p - a \rightarrow t - a$   
We get a good result

So here we hope  $t \rightarrow t^*$

predicated bbox:  $x, y, w, h$

anchor bbox:  $x_a, y_a, w_a, h_a$

true bbox:  $x^*, y^*, w^*, h^*$

$$f(x) = \begin{cases} t_x = \frac{x - x_a}{w_a}, t_y = \frac{y - y_a}{h_a} \\ t_w = \log\left(\frac{w}{w_a}\right), t_h = \log\left(\frac{h}{h_a}\right) \\ t_x^* = \frac{x^* - x_a}{w_a}, t_y^* = \frac{y^* - y_a}{h_a} \\ t_w^* = \log\left(\frac{w^*}{w_a}\right), t_h^* = \log\left(\frac{h^*}{h_a}\right) \end{cases}$$

# I. Two Stage Detection

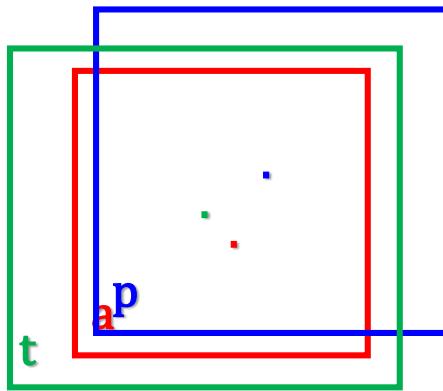
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

C2.2.1: Smooth L1 Loss

➤ Regress (center) offset, NOT COORDINATES



So how to transfer from  
one box to another?

predicated bbox:  $x, y, w, h$

anchor bbox:  $x_a, y_a, w_a, h_a$

true bbox:  $x^*, y^*, w^*, h^*$

$$f(x) = \begin{cases} t_x = \frac{x - x_a}{w_a}, t_y = \frac{y - y_a}{h_a} \\ t_w = \log\left(\frac{w}{w_a}\right), t_h = \log\left(\frac{h}{h_a}\right) \\ t_x^* = \frac{x^* - x_a}{w_a}, t_y^* = \frac{y^* - y_a}{h_a} \\ t_w^* = \log\left(\frac{w^*}{w_a}\right), t_h^* = \log\left(\frac{h^*}{h_a}\right) \end{cases}$$

# I. Two Stage Detection

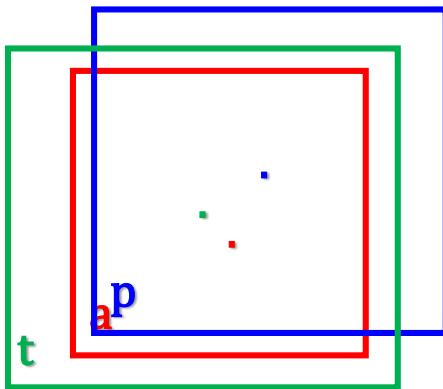
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

C2.2.1: Smooth L1 Loss

➤ Regress (center) offset, NOT COORDINATES



So how to transfer from  
one box to another?

Shift + Scale

predicated bbox:  $x, y, w, h$

anchor bbox:  $x_a, y_a, w_a, h_a$

true bbox:  $x^*, y^*, w^*, h^*$

$$f(x) = \begin{cases} t_x = \frac{x - x_a}{w_a}, t_y = \frac{y - y_a}{h_a} \\ t_w = \log\left(\frac{w}{w_a}\right), t_h = \log\left(\frac{h}{h_a}\right) \\ t_x^* = \frac{x^* - x_a}{w_a}, t_y^* = \frac{y^* - y_a}{h_a} \\ t_w^* = \log\left(\frac{w^*}{w_a}\right), t_h^* = \log\left(\frac{h^*}{h_a}\right) \end{cases}$$

# I. Two Stage Detection

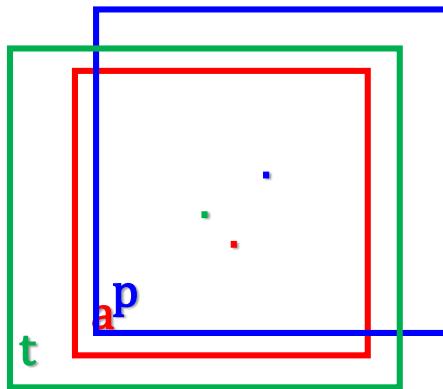
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

C2.2.1: Smooth L1 Loss

➤ Regress (center) offset, NOT COORDINATES



$$x = x_a + \Delta x$$

$$y = y_a + \Delta y$$

$$w = w_a \cdot \Delta w$$

$$h = h_a \cdot \Delta h$$

Shift

Scale

$$f(x) = \begin{cases} t_x = \frac{x - x_a}{w_a}, t_y = \frac{y - y_a}{h_a} \\ t_w = \log\left(\frac{w}{w_a}\right), t_h = \log\left(\frac{h}{h_a}\right) \\ t_x^* = \frac{x^* - x_a}{w_a}, t_y^* = \frac{y^* - y_a}{h_a} \\ t_w^* = \log\left(\frac{w^*}{w_a}\right), t_h^* = \log\left(\frac{h^*}{h_a}\right) \end{cases}$$

predicated bbox:  $x, y, w, h$

anchor bbox:  $x_a, y_a, w_a, h_a$

true bbox:  $x^*, y^*, w^*, h^*$

# I. Two Stage Detection

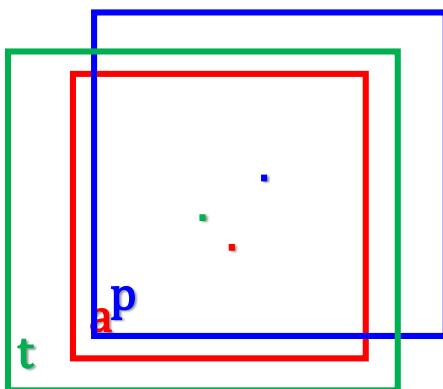
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

C2.2.1: Smooth L1 Loss

➤ Regress (center) offset, NOT COORDINATES



$$\begin{aligned}x &= x_a + \Delta x \\y &= y_a + \Delta y \\w &= w_a \cdot \Delta w \\h &= h_a \cdot \Delta h \\ \Delta x &= w_a \cdot t_x \\ \Delta y &= h_a \cdot t_y \\ \Delta w &= w_a \cdot e^{t_w} \\ \Delta h &= h_a \cdot e^{t_h}\end{aligned}$$

Shift

Scale



$$f(x) = \begin{cases} t_x = \frac{x - x_a}{w_a}, t_y = \frac{y - y_a}{h_a} \\ t_w = \log\left(\frac{w}{w_a}\right), t_h = \log\left(\frac{h}{h_a}\right) \\ t_x^* = \frac{x^* - x_a}{w_a}, t_y^* = \frac{y^* - y_a}{h_a} \\ t_w^* = \log\left(\frac{w^*}{w_a}\right), t_h^* = \log\left(\frac{h^*}{h_a}\right) \end{cases}$$

predicated bbox:  $x, y, w, h$

anchor bbox:  $x_a, y_a, w_a, h_a$

true bbox:  $x^*, y^*, w^*, h^*$

# I. Two Stage Detection

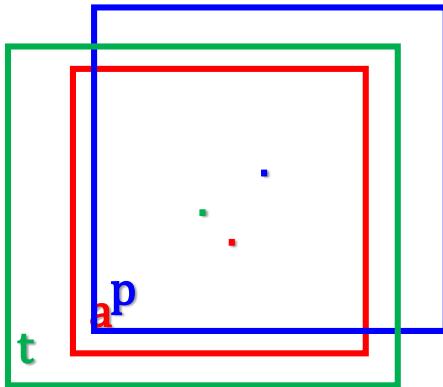
## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

C2.2.1: Smooth L1 Loss

➤ Regress (center) offset, NOT COORDINATES



Q: Why use exp

A: tx & ty is used for scaling  
we need to ensure scale > 0

Q: Why use log

A: 1. reversing procedure of exp  
2. Suppress greater bbox

predicated bbox:  $x, y, w, h$

anchor bbox:  $x_a, y_a, w_a, h_a$

true bbox:  $x^*, y^*, w^*, h^*$

$$f(x) = \begin{cases} t_x = \frac{x - x_a}{w_a}, t_y = \frac{y - y_a}{h_a} \\ t_w = \log\left(\frac{w}{w_a}\right), t_h = \log\left(\frac{h}{h_a}\right) \\ t_x^* = \frac{x^* - x_a}{w_a}, t_y^* = \frac{y^* - y_a}{h_a} \\ t_w^* = \log\left(\frac{w^*}{w_a}\right), t_h^* = \log\left(\frac{h^*}{h_a}\right) \end{cases}$$

# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

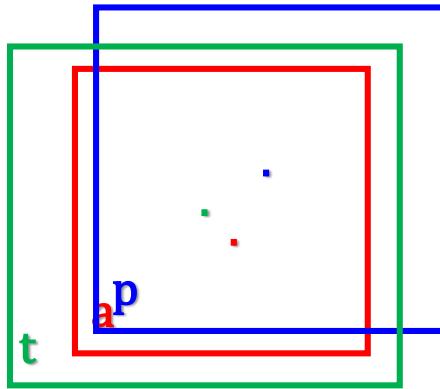
- Structure:

C2.2.1: Smooth L1 Loss

➤ Regress (center) offset, NOT COORDINATES

Q: Why x/y uses linear transferring  
but w/h uses non-linear transferring?  
Is that reasonable?

A: When bbox1->bbox2, then we can.



predicated bbox:  $x, y, w, h$

anchor bbox:  $x_a, y_a, w_a, h_a$

true bbox:  $x^*, y^*, w^*, h^*$

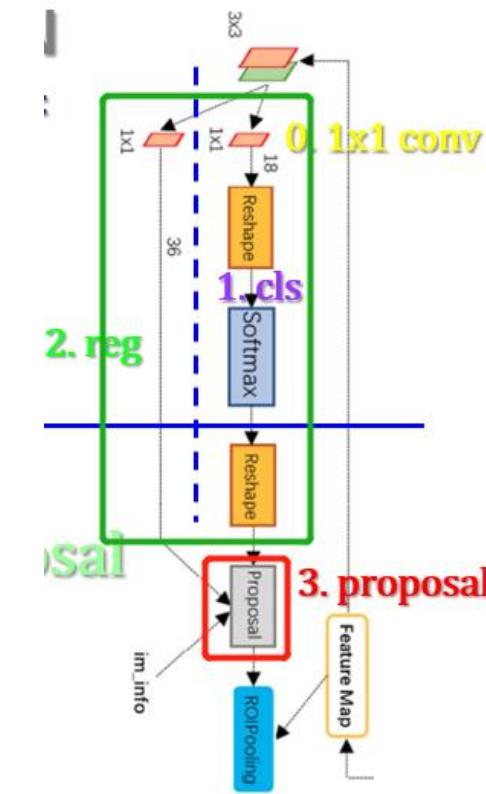
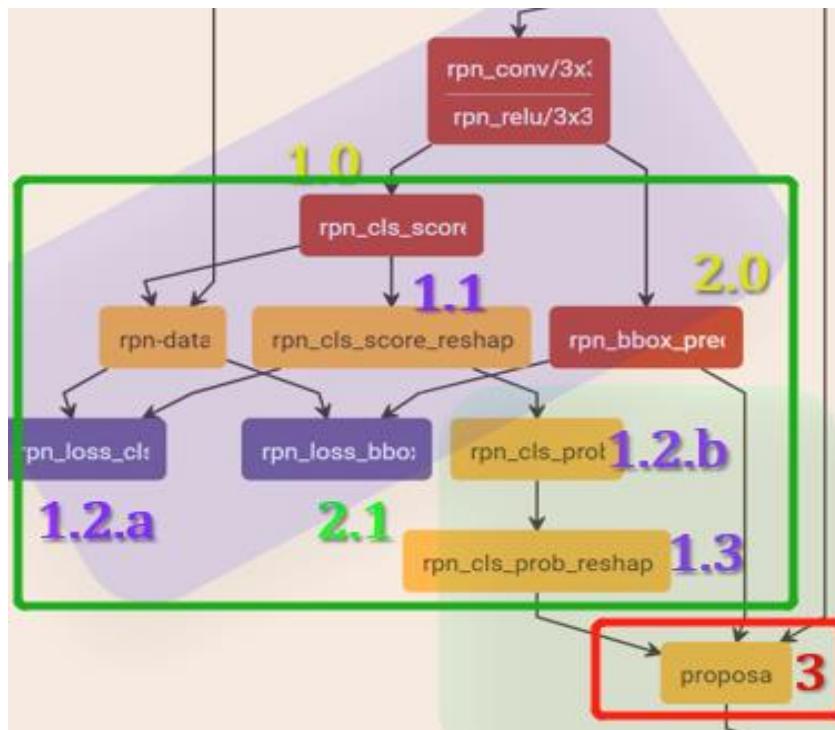
$$f(x) = \begin{cases} t_x = \frac{x - x_a}{w_a}, t_y = \frac{y - y_a}{h_a} \\ t_w = \log\left(\frac{w}{w_a}\right), t_h = \log\left(\frac{h}{h_a}\right) \\ t_x^* = \frac{x^* - x_a}{w_a}, t_y^* = \frac{y^* - y_a}{h_a} \\ t_w^* = \log\left(\frac{w^*}{w_a}\right), t_h^* = \log\left(\frac{h^*}{h_a}\right) \end{cases}$$

# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure: C2.3: Proposal: Get proposal



# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C2. RPN

- Structure:

#### C2.3: Proposal: Get proposal

- Regard anchor as FG when  $\text{IoU} > 0.7$
- Regard anchor as FG when  $\text{IoU} < 0.3$
- Regardless other anchors
- Then just keep 128: 0.25 fg + 0.75 bg anchors

# I. Two Stage Detection

## C. Faster R-CNN [2015 Ren]:

### C3. After RPN

- Same as Fast RCNN

### C4. How to train

- Classic: alternatively 4 steps

- Use ImageNet finetune our RPN
- Use trained proposal from step1 to train an Fast RCNN
- Use detected results to initialize RPN training where we freeze the backbone layers but to train pure RPN-related layers
- Finetune Fast RCNN-related layers only.

- Other method:

- We can also train Faster RCNN as a whole in just one step

## II. One Stage Detection



## II. One Stage Detection

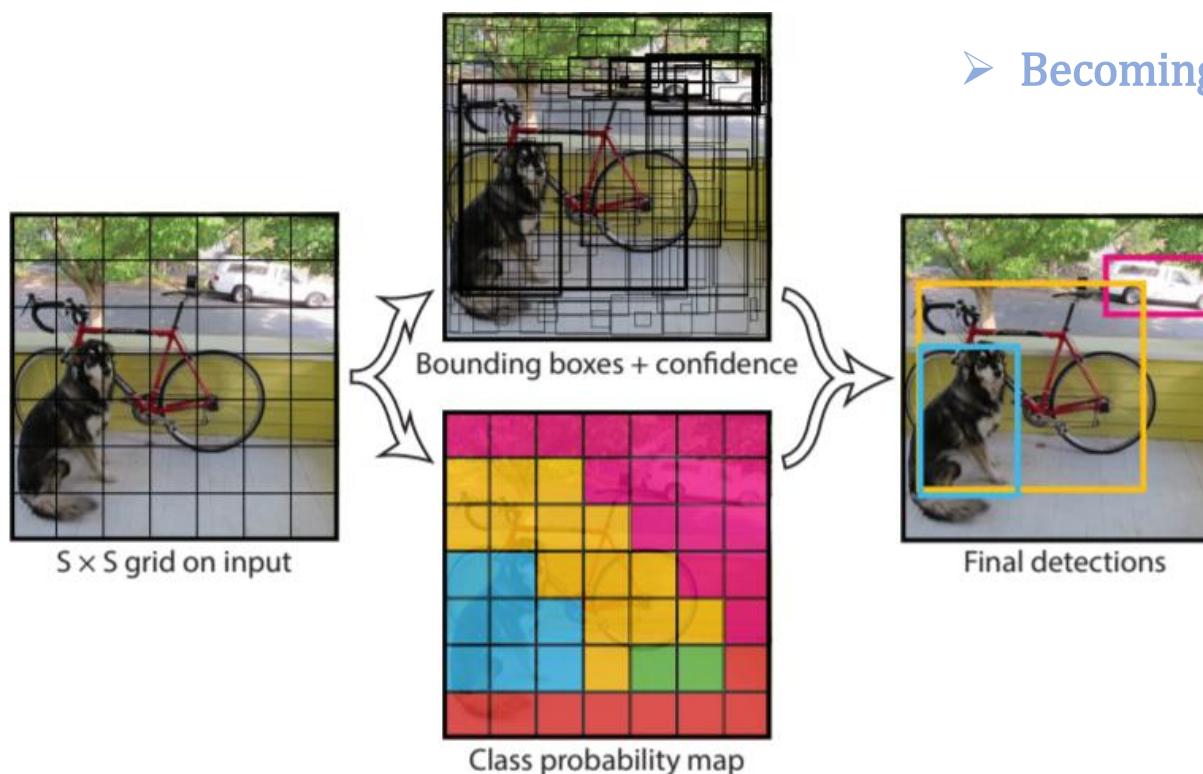
**Why we have to use  
anchors to refine our  
bounding boxes?**

# II. One Stage Detection

## D. Yolo V1 [2015, Joseph]

### D0. You Only Look Once!

- Really fast (18 faster rcnn [ZF] vs 45 yolo )
- Becoming prevalent (62.1 mAP vs 63.4)



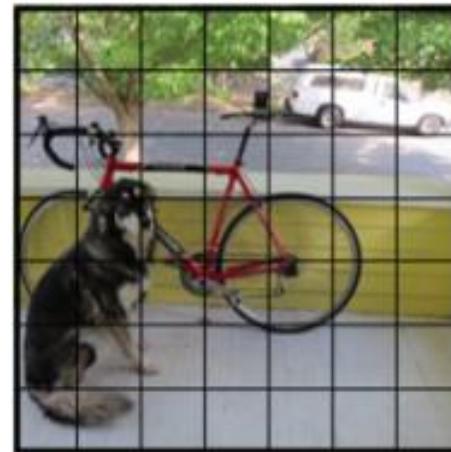
# II. One Stage Detection

## D. Yolo V1 [2015, Joseph]

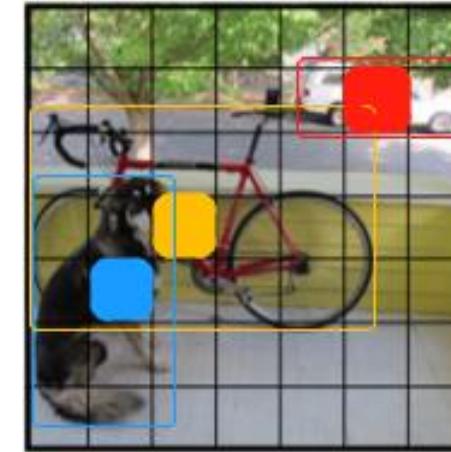
### D1. Procedure

- Grid an image into  $S \times S$  cells [ $448 \times 448 \rightarrow 7 \times 7$ ]

One cell will be responsible for predicting an object as long as an object's center locating in that cell.



$S \times S$  grid on input



$S \times S$  grid on input

# II. One Stage Detection

## D. Yolo V1 [2015, Joseph]

### D1. Procedure

- Each cell predicts **B** bounding box **with a confidence**

**Bounding Box:**  $x, y, w, h$  (center)

**Confidence:**  $P_r(\text{object}) \cdot IoU_{pred}^{truth}$

**Final output tensor:**  $S \times S \times (5 * B + C)$

$[7 \times 7 \times (5 * 2 + 20)] \rightarrow v1$

# II. One Stage Detection

D. Yolo V1 [2015, Joseph]

## D2. Loss Function

$$2.1 \quad \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \quad (1)$$

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \quad (2)$$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \quad (3)$$

$$2.2 \quad + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \quad (4)$$

$$2.3 \quad + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad (5)$$

$i: 0 \sim (S^2 - 1)$  [iterate each grid cell (0~48)]

$j: 0 \sim (B - 1)$  [iterate each bbox (0~2)]

$\mathbb{1}_{ij}^{\text{obj}}$  &  $\mathbb{1}_{ij}^{\text{noobj}}$ :

0	0	0	0	0	0	0	0
0	0	0	0	0	0	<b>1</b>	0
0	0	0	0	0	0	0	0
0	0	<b>1</b>	0	0	0	0	0
0	<b>1</b>	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0

1	1	1	1	1	1	1	1
1	1	1	1	1	1	<b>0</b>	1
1	1	1	1	1	1	1	1
1	1	<b>0</b>	1	1	1	1	1
1	<b>0</b>	1	1	1	1	1	0
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1

# III. Other Methods



# Summary