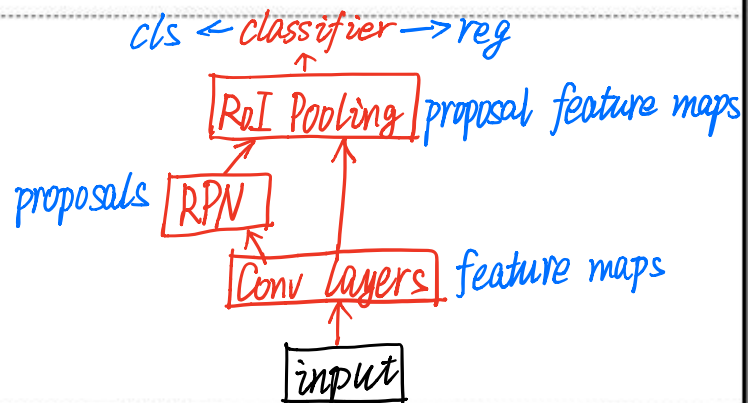# Faster R-CNN (NIPS 2015)

Faster R-CNN: Towards Real-Time Object Detection with Region

Ren et al.

① Problem Description:

Faster R-CNN: region proposal computation.

SPPnet (slide window),

R-CNN (selective search).

② Problem Solution:

RPN (Region Proposal Network):

it shares full-image convolutional features with the detection network.

"attention" mechanisms:
RPN component tells the unified network where to look.

③ Conceptual Understanding

cls ← classifier → reg

RoI Pooling — proposal feature maps

proposals — RPN

Conv layers — feature maps

input

① Conv layers: conv + relu + pooling.  (size=3, pad=1, stride=1)  (2,1,1) (2,0,1)  M×N $\xrightarrow{VGG16}$ $\frac{M}{16} \times \frac{N}{16}$.

② RPN: feature maps → 1×1 (anchors) → 1×1 $\frac{18}{2k}$ → Reshape → softmax → Reshape → proposal     36 (4k) (regression)  (classification)

③ RoI pooling: (M×N) proposals → (w×h) → proposal feature maps   feature maps

④ Classifition: proposal feature maps → softmax → cls    (regression) → bbox_pred   (classification)

# Details of implementation

## RPN:

### ① anchors:

output: $(2k+4k) \times 256$-d.

| 2k scores | 4k coordinates |

cls ↗     ↖ reg
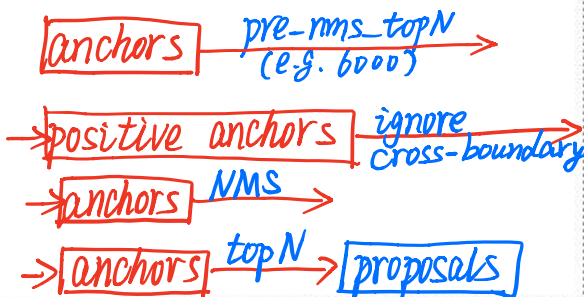
256-d

↑ ZF model

feature map  W×H

k (3×3) anchor boxes

128
256
512

2:1
1:1
1:2

### ② classification + regression:

[ ] → 1×1 $\xrightarrow{\frac{18}{2k}}$ softmax → positive anchors

→ 1×1 $\xrightarrow{\frac{36}{4k}}$ bbox regression

### ③ proposal layers:

anchors $\xrightarrow{\text{pre-nms\_topN (e.g. 6000)}}$

→ positive anchors $\xrightarrow{\substack{\text{ignore} \\ \text{cross-boundary}}}$

→ anchors $\xrightarrow{\text{NMS}}$

→ anchors $\xrightarrow{\text{topN}}$ proposals
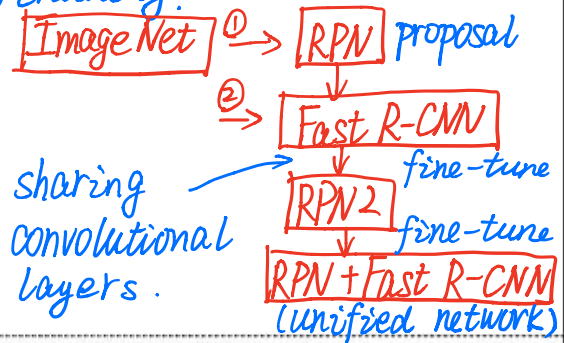
## Experiments:

### ① loss function:

$$L(\{P_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(P_i, P_i^*)$$

$$+ \lambda \frac{1}{N_{reg}} \sum_i P_i^* L_{reg}(t_i, t_i^*).$$

e.g. $N_{cls} = 256$, $N_{reg} = 2400$, $\lambda = \frac{N_{reg}}{N_{cls}} \approx 10$.

$$L_{reg}(t_i, t_i^*) = \sum_{i \in \{x, y, w, h\}} smooth_{L_1}(t_i - t_i^*).$$

### ② training:

ImageNet $\xrightarrow{①}$ RPN proposal

$\xrightarrow{②}$ Fast R-CNN

↓ fine-tune

RPN2

↓ fine-tune

RPN + Fast R-CNN
(unified network)

sharing convolutional layers.

### ③ results:

VOC

COCO

COCO + VOC

## Code: (PyTorch Impletation)

① modified files to suit own machine.

② change custom datasets to train.