



國立高雄科技大學

National Kaohsiung University of Science and Technology

Faster R-CNN

Towards Real-Time Object Detection with Region Proposal Networks

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun

*IEEE Trans. on Pattern Analysis and Machine
Intelligence (PAMI), 2017*

Speaker: Shih-Shinh Huang

September 24, 2020





Outline

- Introduction
 - About Object Detection
 - Background
 - Architecture Overview
- Network Ingredient
 - Region Proposal Network (RPN)
 - RoI Pooling
 - Detection Network
- Faster R-CNN Training
 - Training Overview
 - RPN Loss
 - Detection Loss

Introduction

- About Object Detection

- Input:

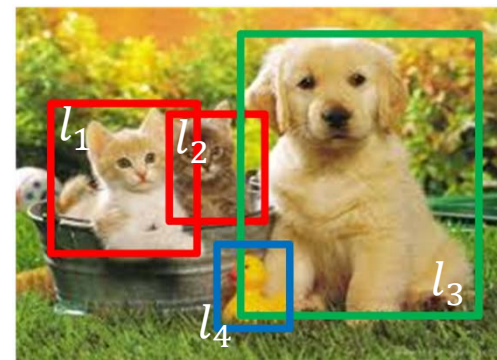
- I : input image
 - $\{c_1, c_2, \dots, c_n\}$: object classes to be detected

- Output:

- $\{r_1, r_2, \dots, r_m\}$: bounding boxes of m detected objects
 - $\{l_1, l_2, \dots, l_m\}$: class labels of all detected objects



$C = \{\text{cat, dog, duck}\}$



$l_1 = \text{cat}$

$l_2 = \text{cat}$

$l_3 = \text{dog}$

$l_4 = \text{duck}$



Introduction

- Background

- **Region-based CNN (R-CNN)** has made remarkable advance in object detection.

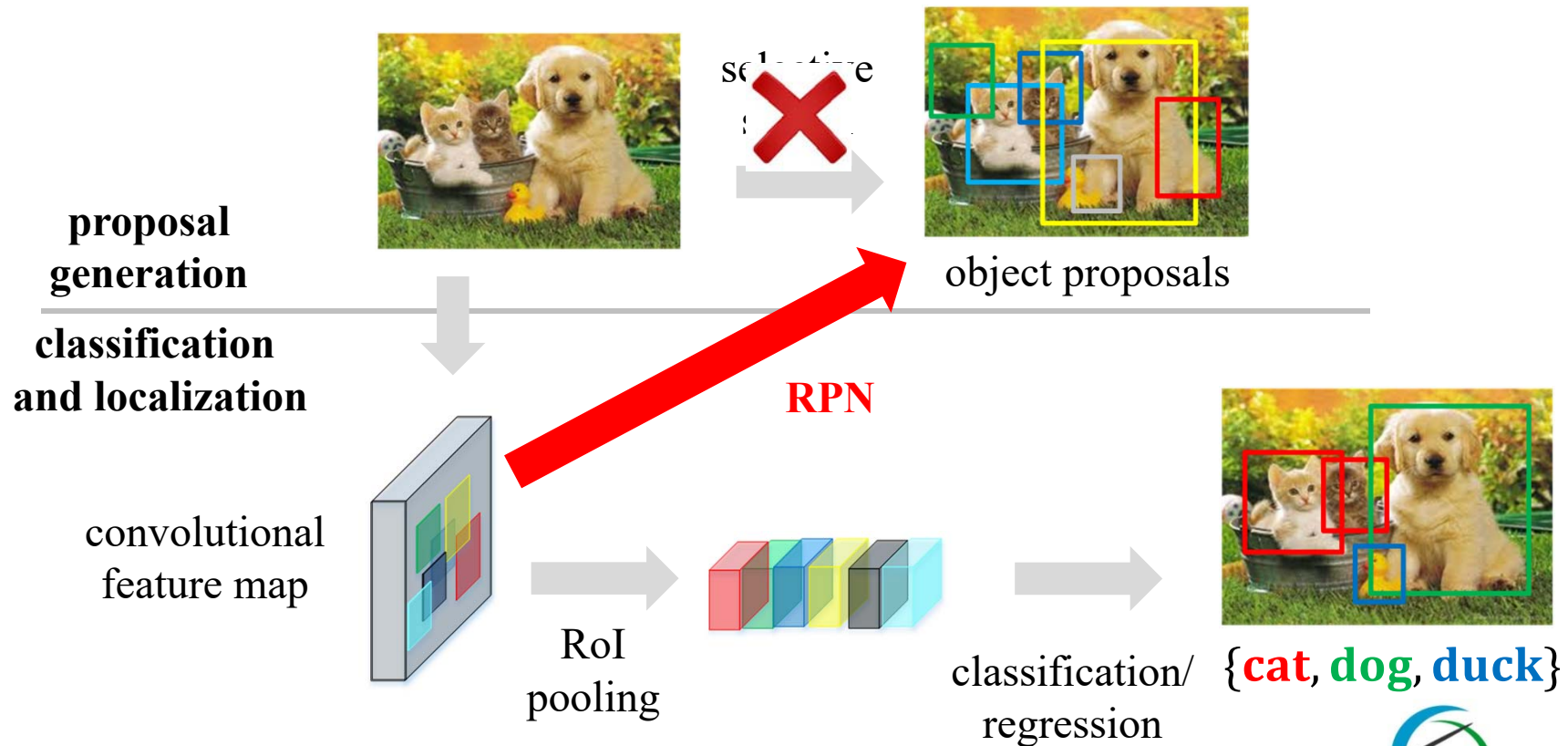
R. Girshick, et. al. “Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation,” *CVPR*, 2014

- **Fast R-CNN** improves detection speed by sharing convolution feature map across proposals.

R. Girshick, “Fast R-CNN,” *CVPR*, 2015

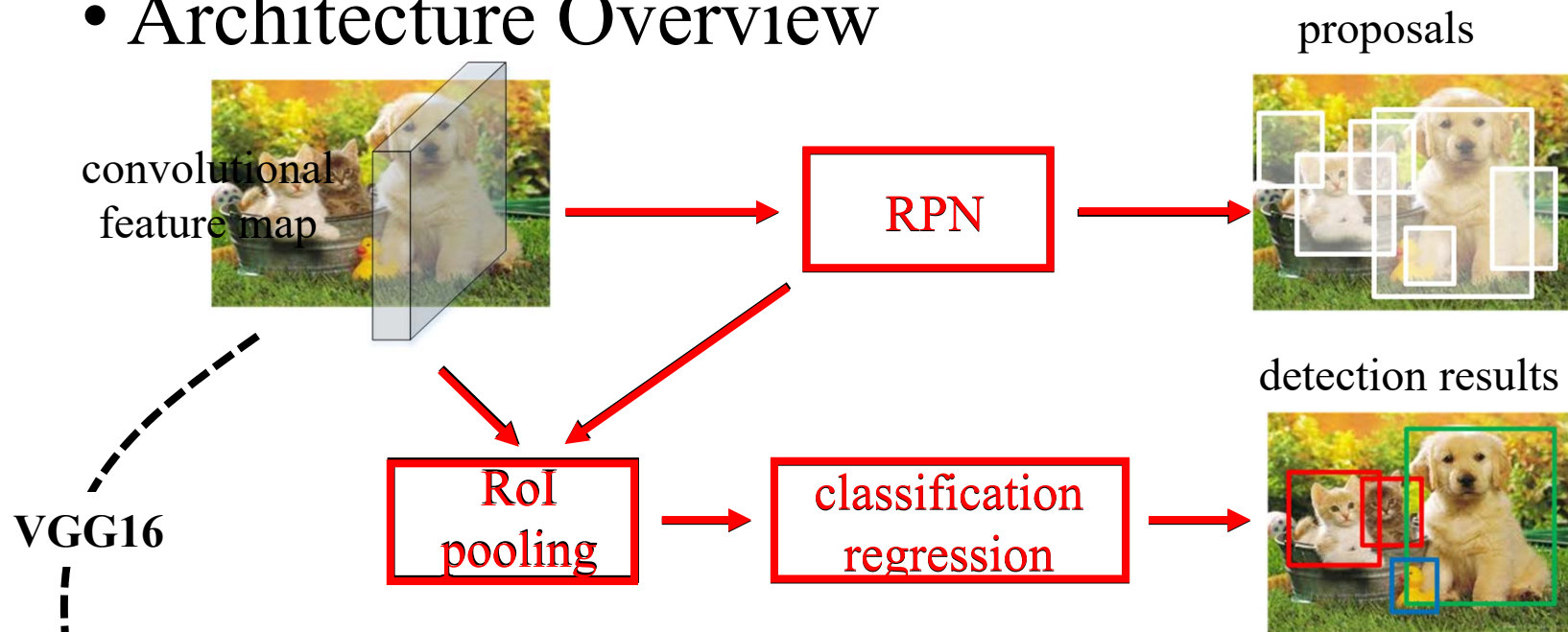
Introduction

- Background: ~~Fast~~R-CNN



Introduction

- Architecture Overview

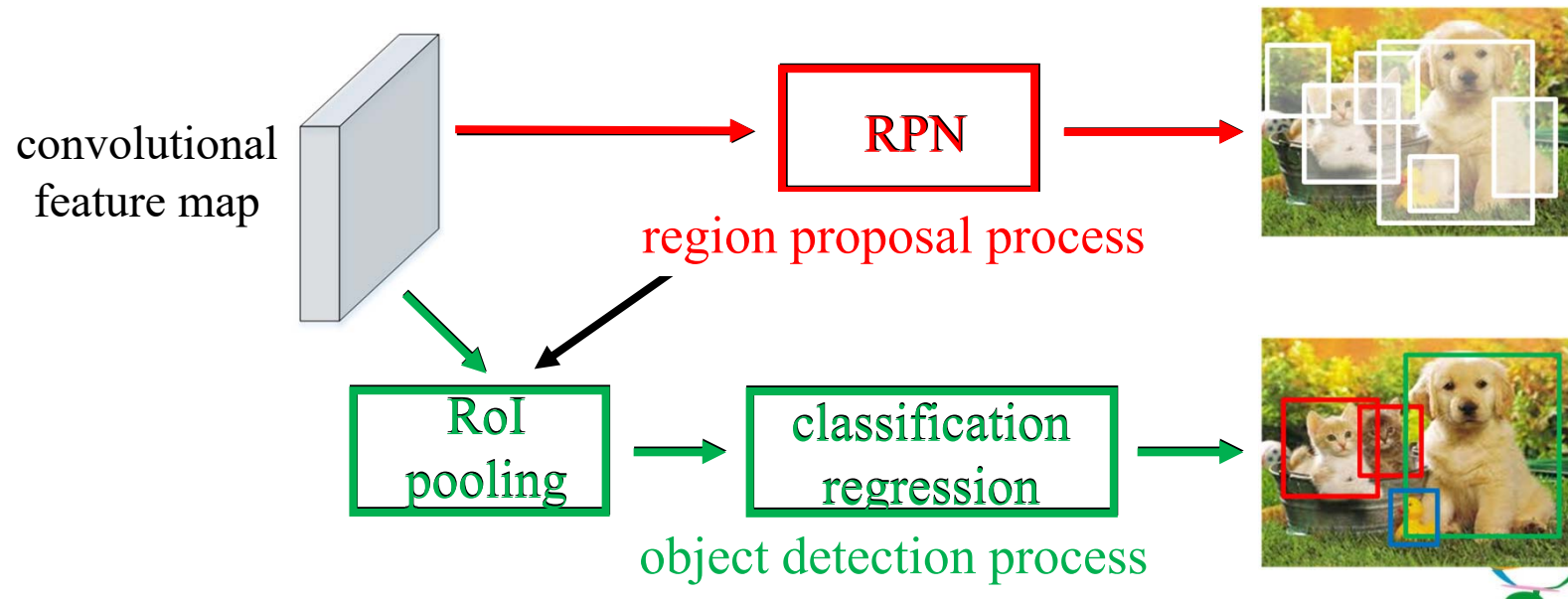


RPN finds proposals by max, pooling. RoI pooling converts proposals into a fixed size. The proposals are then processed by classification and regression.

K. Simonyan et. al. "Very Deep Convolutional Networks for Large-Scale Image Recognition," *ICLR*, 2015.

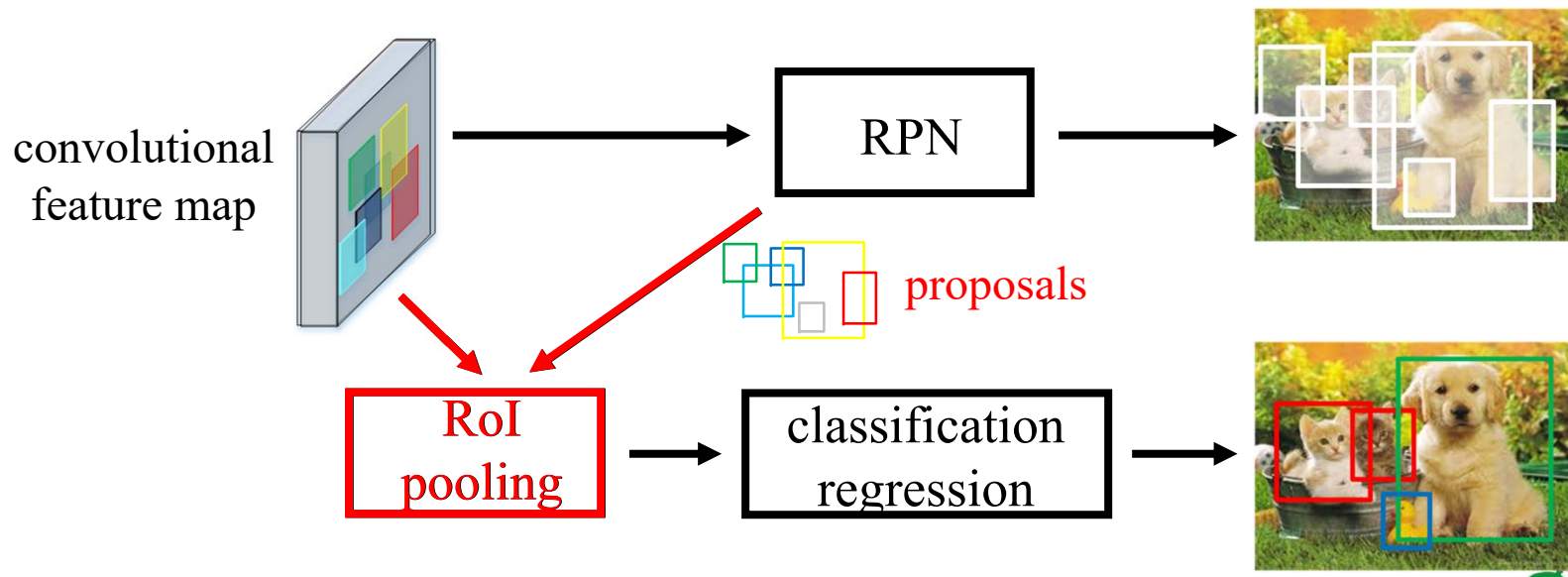
Introduction

- Architecture Overview
 - **Efficiency**: two-level feature sharing
 - **first-level**: region proposal and detection
 - **second-level**: among generated proposals



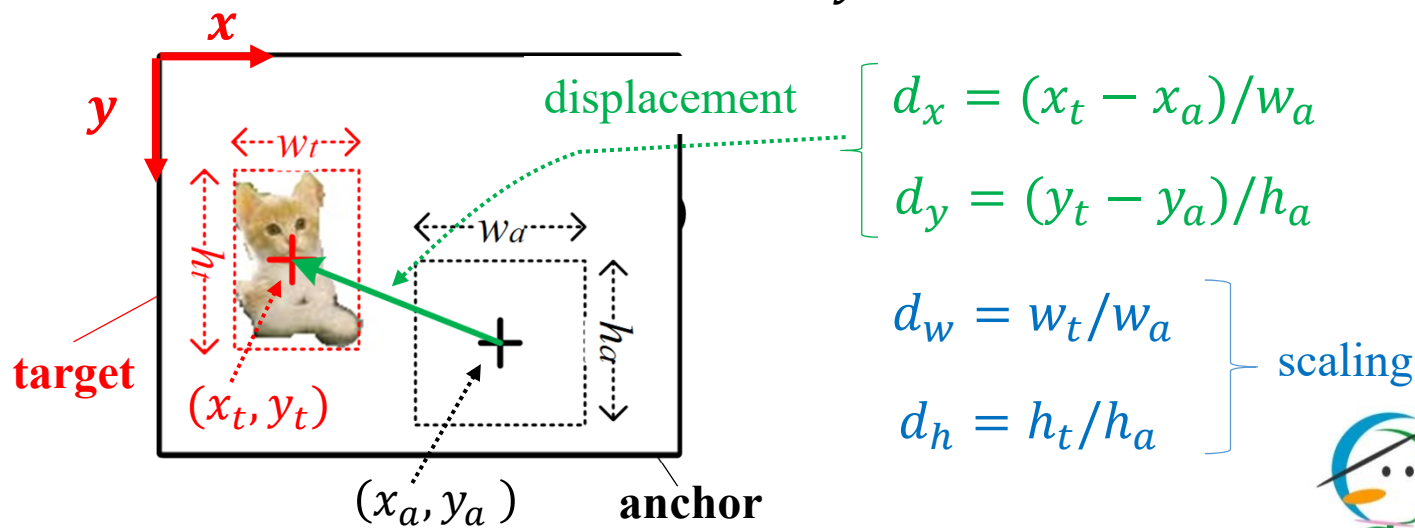
Introduction

- Architecture Overview
 - **Efficiency**: two-level feature sharing
 - **first-level**: region proposal and detection
 - **second-level**: among generated proposals



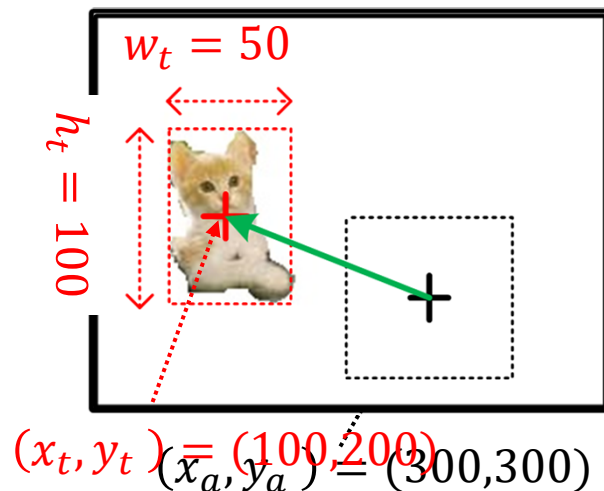
Introduction

- Architecture Overview
 - Effectiveness: two-level anchor mechanism.
 - Anchor is a reference box (x_a, y_a, w_a, h_a)
 - Anchor parameterizes a target box (x_t, y_t, w_t, h_t) by regression offset (d_x, d_y, d_w, d_h)



Introduction

- Architecture Overview
 - Effectiveness: two-level anchor mechanism.
 - Anchor is a reference box (x_a, y_a, w_a, h_a)
 - Anchor parameterizes a target box (x_t, y_t, w_t, h_t) by regression offset (d_x, d_y, d_w, d_h)



$$(x_a, y_a, w_a, h_a) = (100, 200, 50, 100)$$

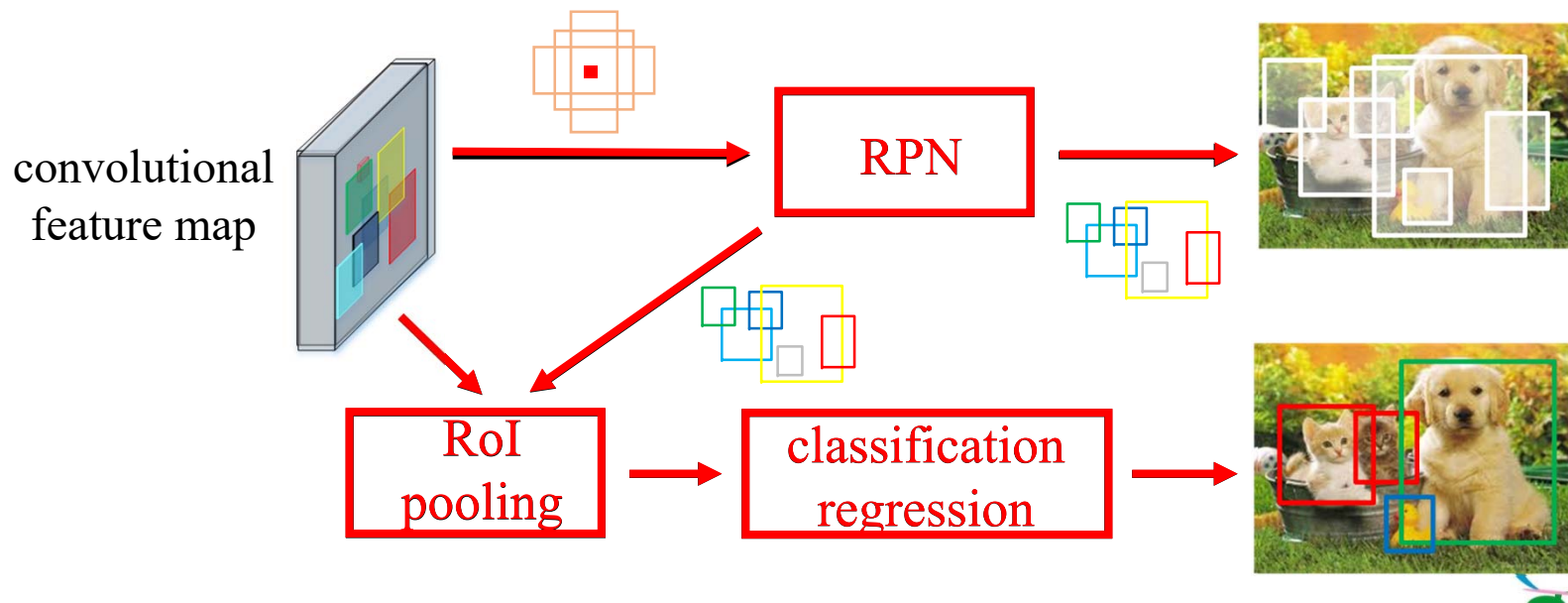
$$(d_x, d_y, d_w, d_h) = (-2.0, -1.0, 0.5, 1.0)$$

$$x_t = 100 = x_a - 3 \cdot w_t = 300 - 3 \cdot 50 \quad d_x = -2.0 = (x_t - x_a) / w_a = (100 - 300) / 50$$

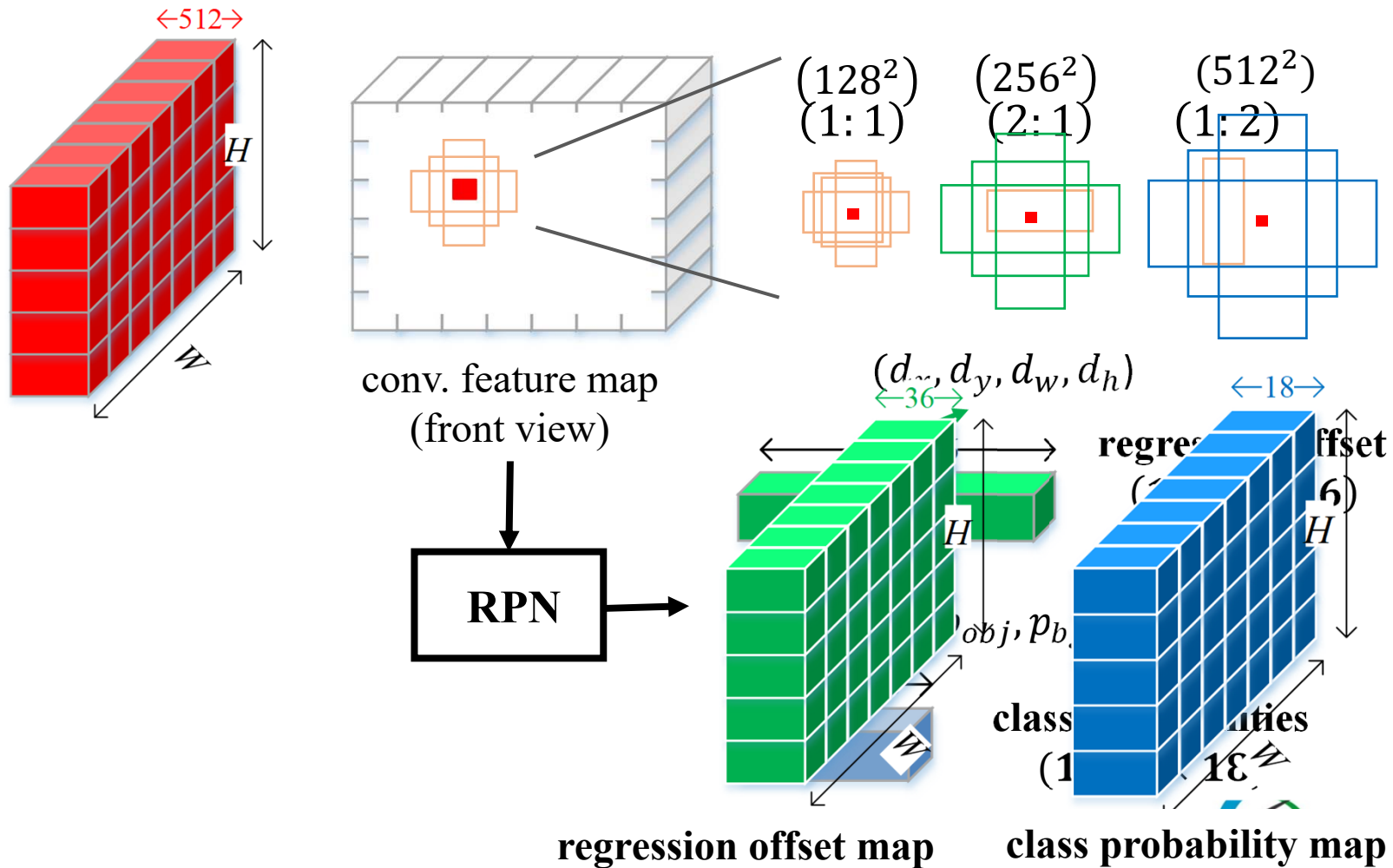
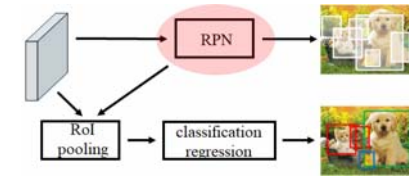
$$y_t = 200 = y_a - 3 \cdot h_t = 300 - 3 \cdot 100 \quad d_y = -1.0 = (y_t - y_a) / h_a = (200 - 300) / 100$$

Introduction

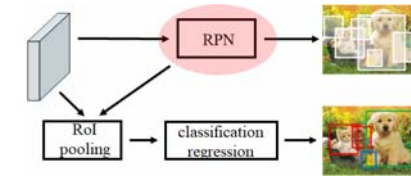
- Architecture Overview
 - Effectiveness: two-level **anchor** mechanism
 - first-level: put several anchors at every point (dense)
 - second-level: take proposals as anchors (sparse)



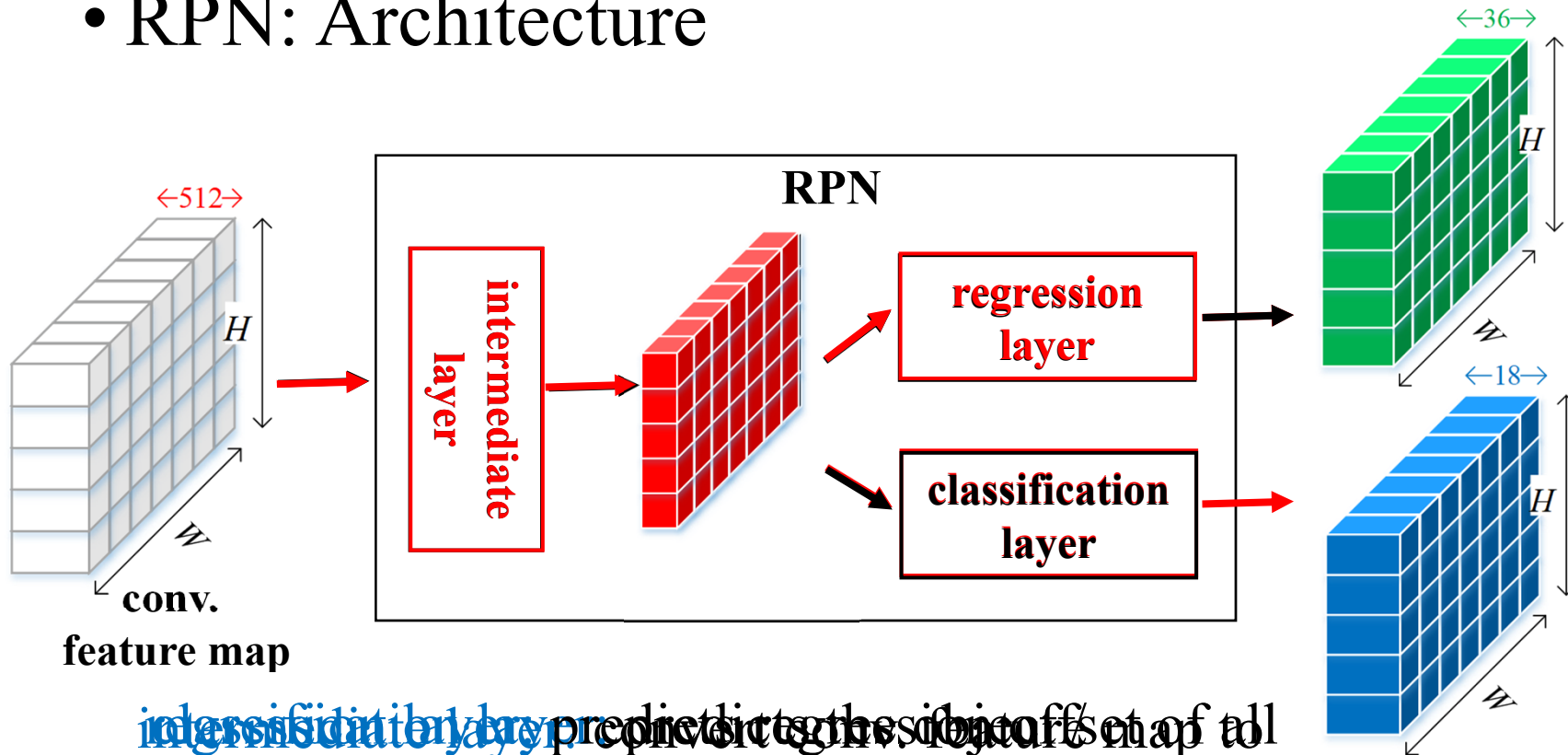
Network Ingredient



Network Ingredient

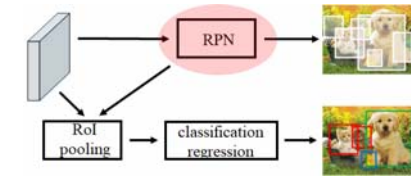


- RPN: Architecture

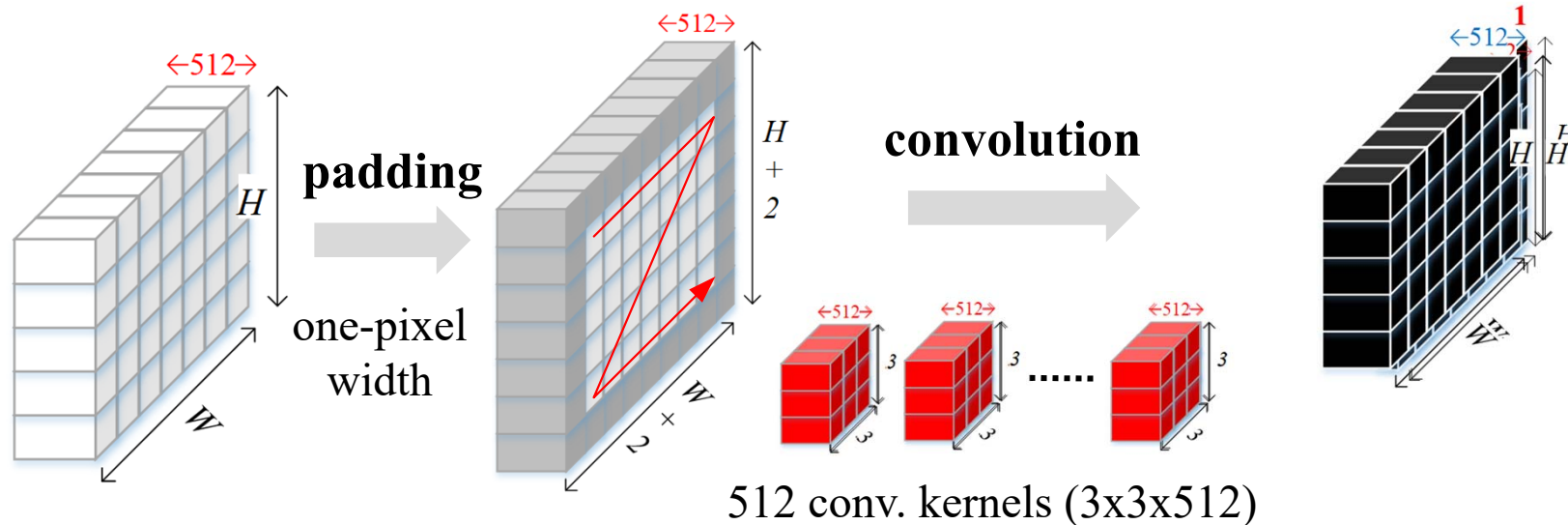


intermediate layer produces the score of all the proposals. regression layer produces the offsets of all the proposals.

Network Ingredient

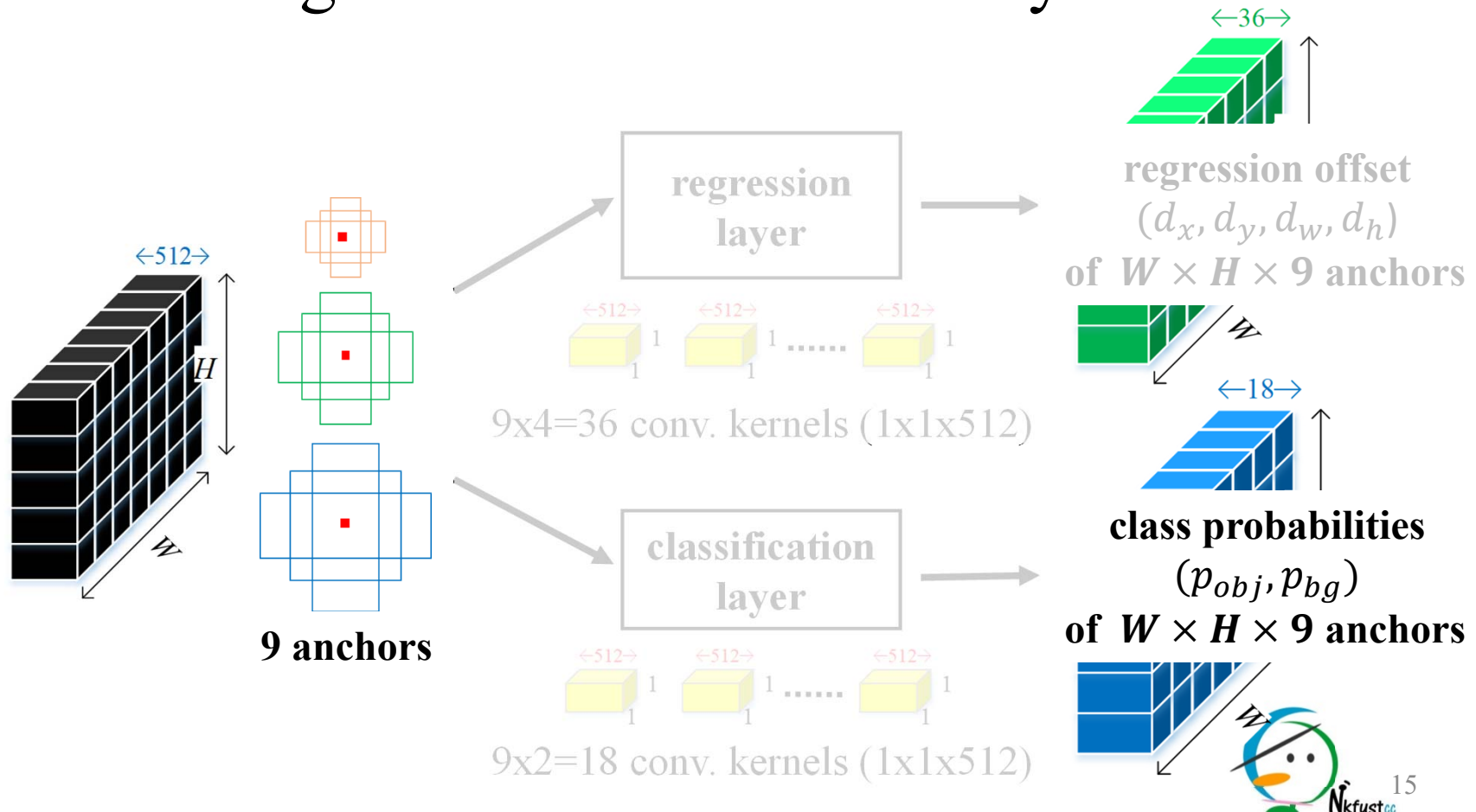


- RPN: intermediate layer
 - extract feature map for proposal generation.



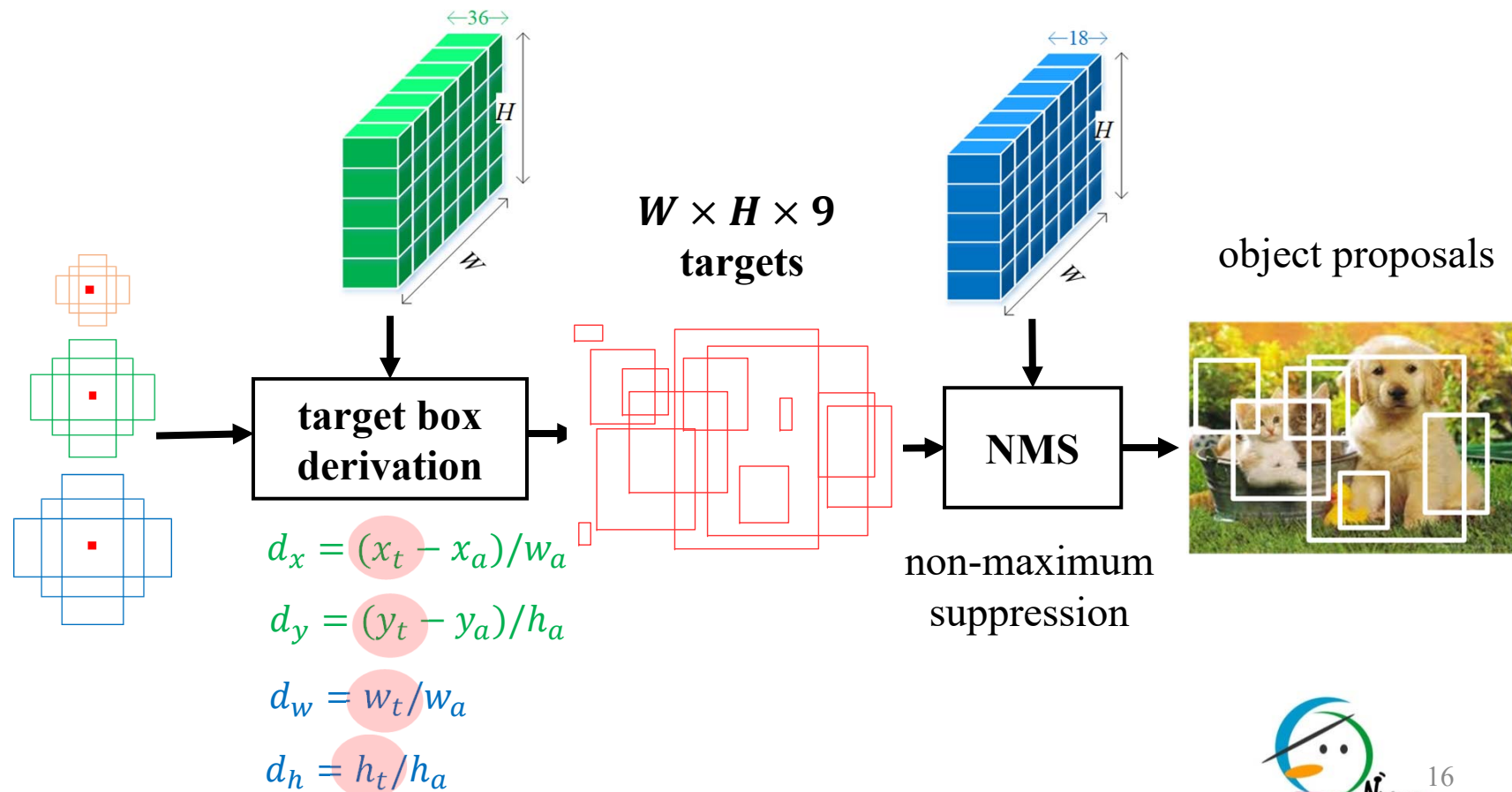
Network Ingredient

- RPN: regression/classification layers

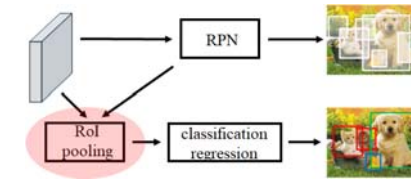


Network Ingredient

- RPN: regression/classification layers

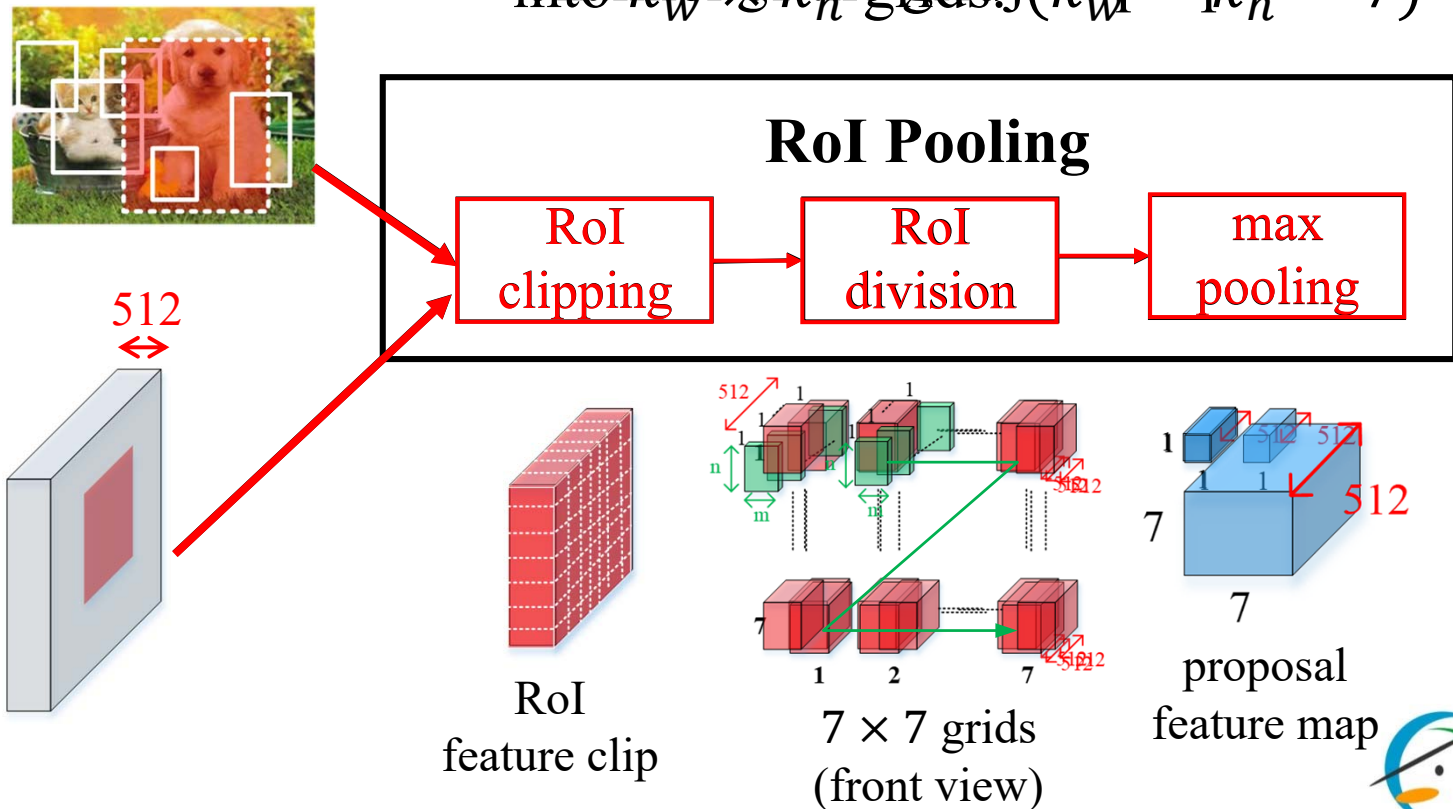


Network Ingredient

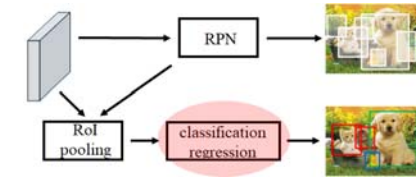


- RoI Pooling

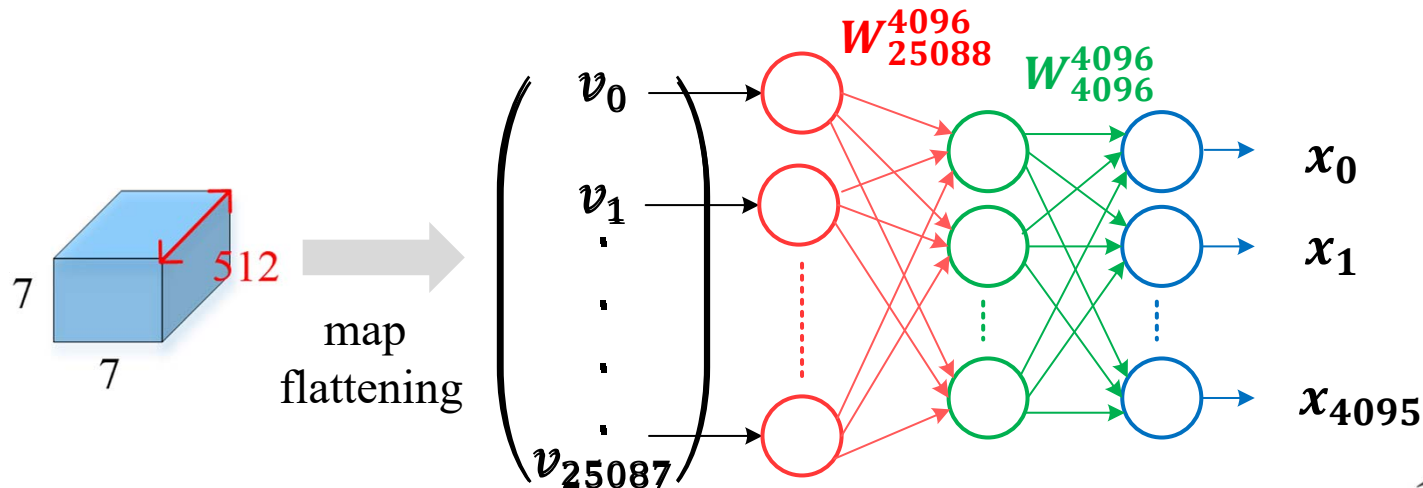
RoI Pooling: apply the max pooling to every element in the grid (e.g. proposal 7×7 grids) ($k_w = k_h = 7$)



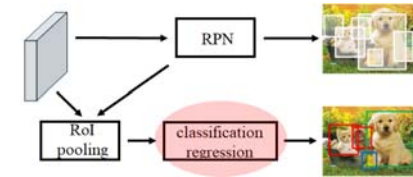
Network Ingredient



- Detection Network
 - flatten proposal feature map to a $7 \times 7 \times 512$ (= 25088) feature vector
 - use two fully connected layers for reducing feature dimension to 4096



Network Ingredient

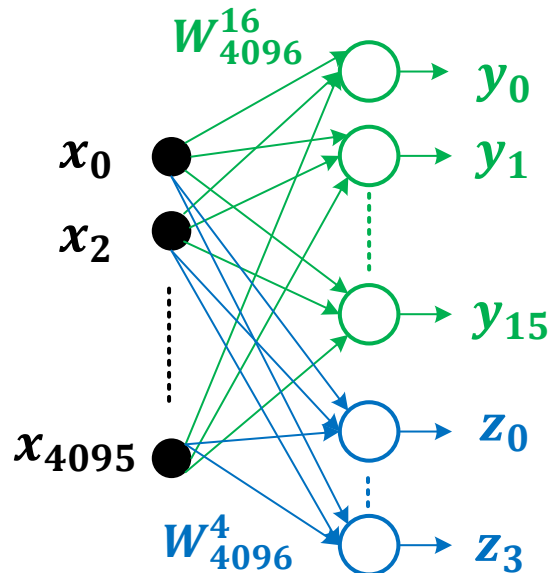


• Detection Network

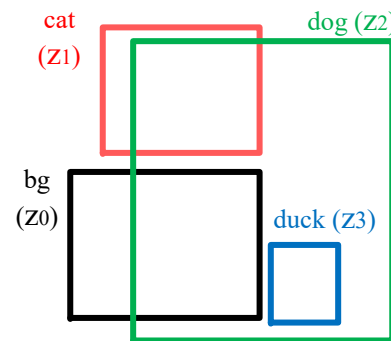
• apply two fully connected branches for

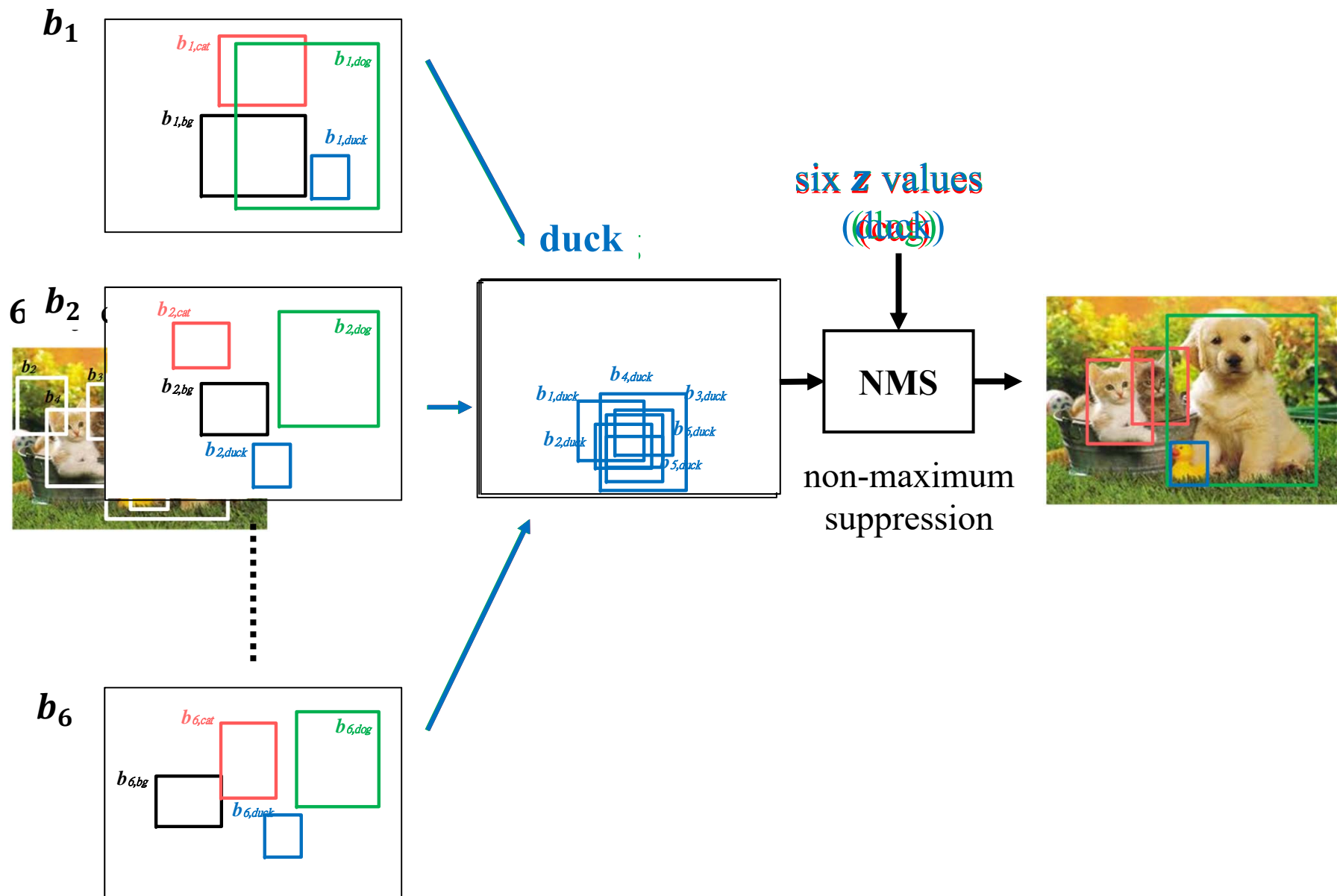
- regression: $4096 \rightarrow (N_{class} + 1) \times 4$ background (bg)

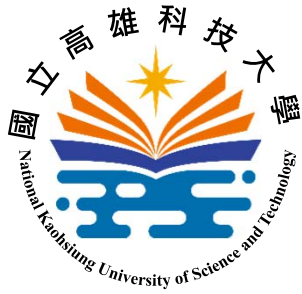
- classification: $4096 \rightarrow (N_{class} + 1)$ offset: (d_x, d_y, d_w, d_h)



$$\hat{C} = \{\mathbf{bg}, \mathbf{cat}, \mathbf{dog}, \mathbf{duck}\} \quad n_{class} = 3$$

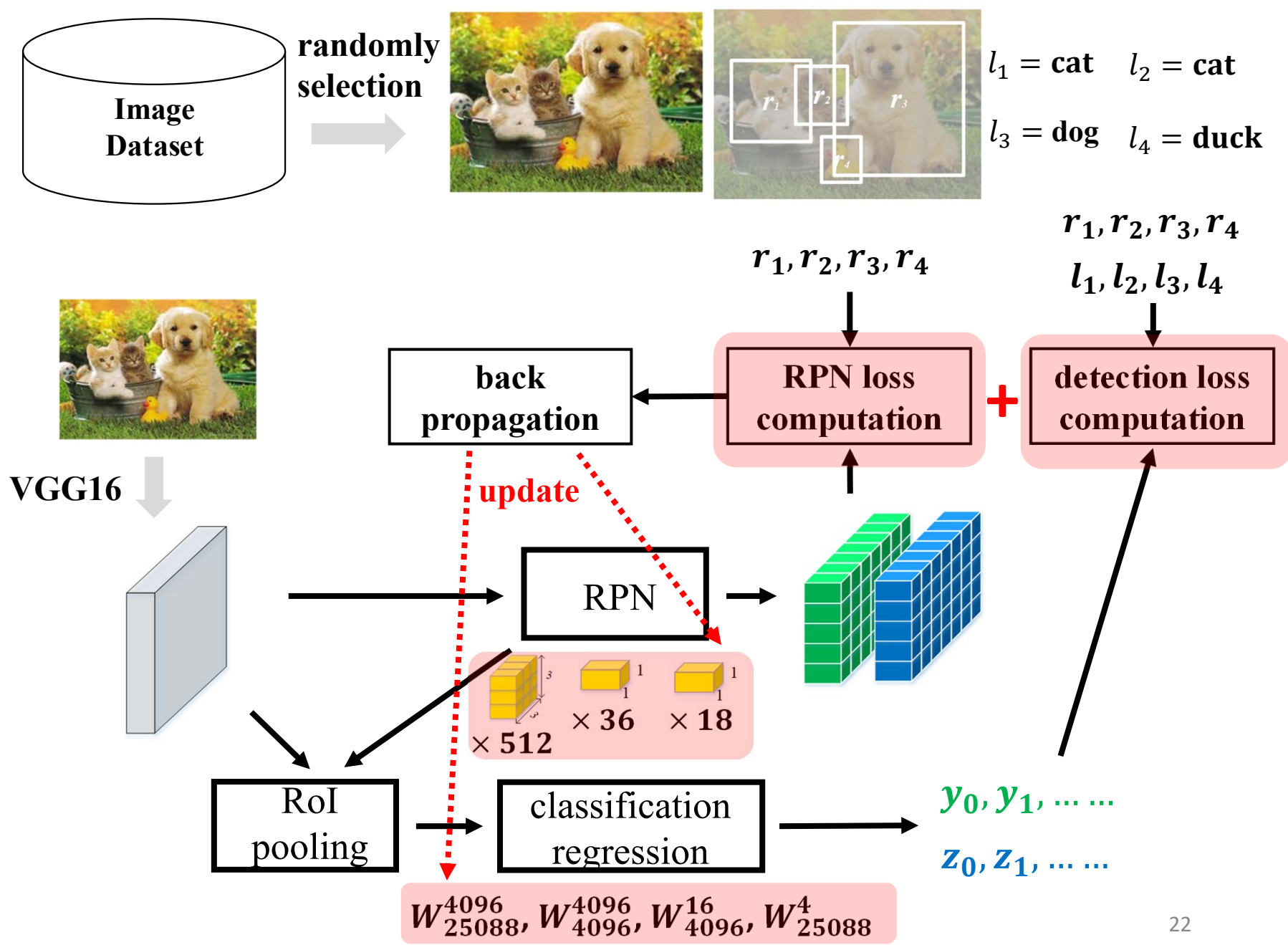






Faster R-CNN Training

- Training Overview
 - select an image with its labels from dataset
 - flow through Faster R-CNN network to obtain
 - **RPN**: regression offset map and classification map
 - **Detection Network**: object regression offset $\{y_0, y_1, \dots\}$ and class probabilities $\{z_0, z_1, \dots\}$
 - compute prediction loss for updating parameters
 - **RPN**: convolutional kernels
 - **Detection Network**: weight matrices





Faster R-CNN Training

- RPN Loss $L^{(r)}(.)$
 - The computation of RPN loss is on a **mini-batch**
 - 128 positive (object) anchors
 - 128 negative (background) anchors

Quarter Unit: Region Proposal Network (PRN)

(00:13:40)

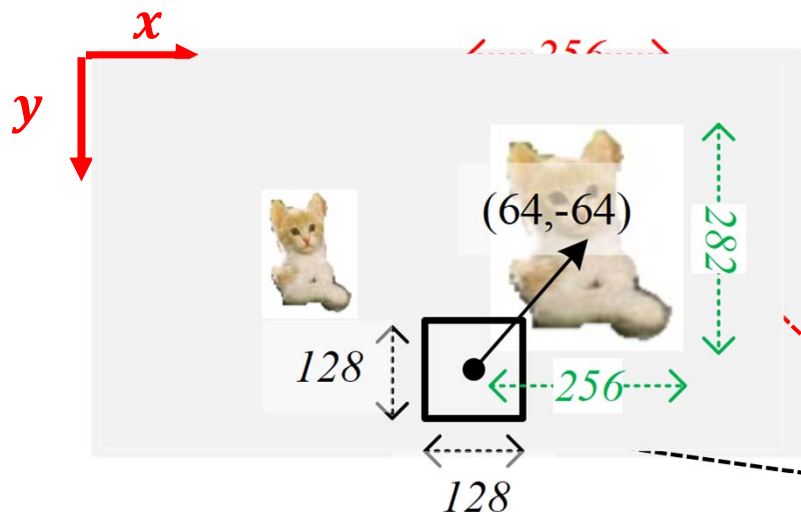
<https://youtu.be/0tBhRfEzUWs>

Faster R-CNN Training

- RPN Loss $L^{(r)}(.)$

- $D = \{\hat{d}^{(i)}, \hat{p}^{(i)}\}_{i=1}^{256}$: selected anchors

- $\hat{p}^{(i)} = (\hat{p}_{obj}^{(i)}, \hat{p}_{bg}^{(i)})$: class probabilities of each anchor-truth box with respect to i th anchor



$$\hat{d}^{(1)} = \begin{pmatrix} 32 & 32 & 256 & 282 \\ 0.125 & 0.125 & 1.0 & 1.1 \\ 256 & 256 & 256 & 256 \end{pmatrix}$$

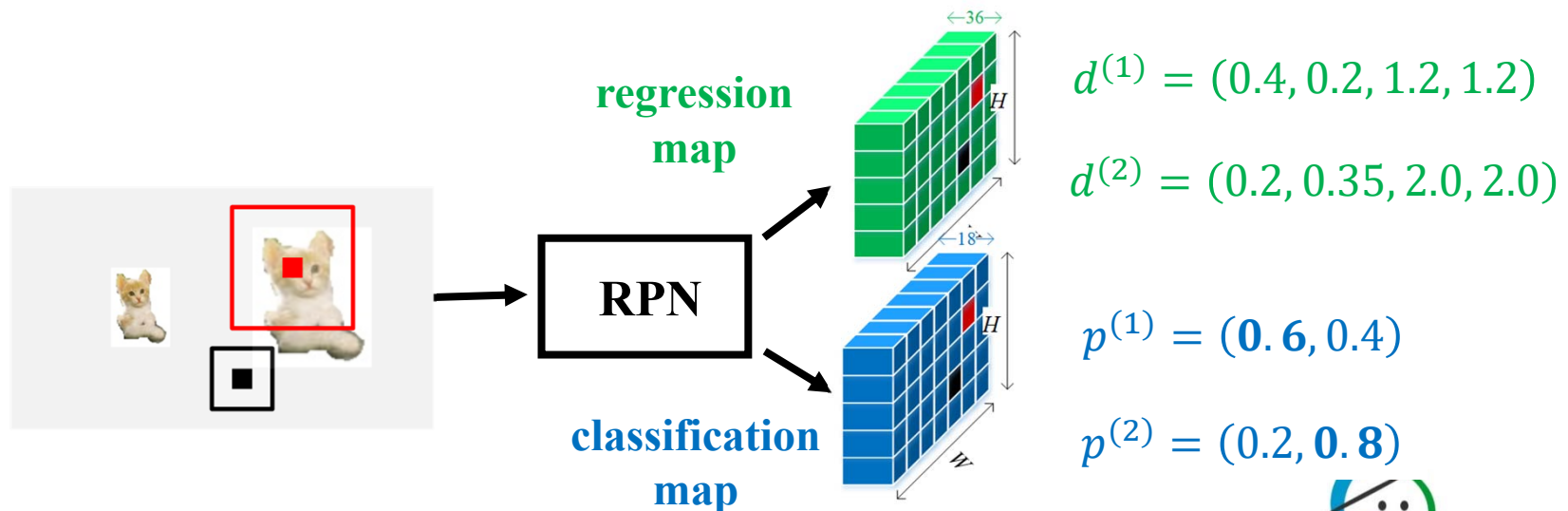
$$\hat{d}^{(2)} = \begin{pmatrix} 64 & -64 & 256 & 282 \\ 0.25 & 0.25 & 2.0 & 2.2 \\ 256 & 256 & 128 & 128 \end{pmatrix}$$

$$\hat{p}^{(1)} = (\hat{p}_{obj}, \hat{p}_{bg}) = (1.0, 0.0)$$

$$\hat{p}^{(2)} = (0.0, 1.0)$$

Faster R-CNN Training

- RPN Loss $L^{(r)}(\cdot)$
 - $\{d^{(i)}, p^{(i)}\}$: predicted parameters from RPN
 - $L^{(r)}(\cdot)$ evaluates fitness of $\{d^{(i)}, p^{(i)}\}$ to $\{\hat{d}^{(i)}, \hat{p}^{(i)}\}$



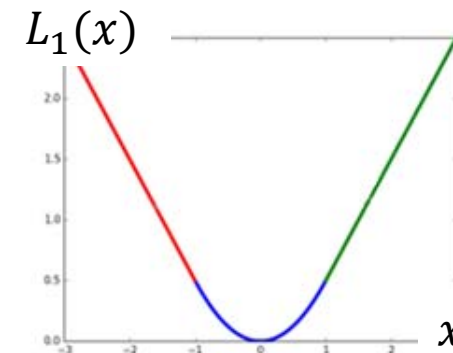


Faster R-CNN Training

- RPN Loss $L^{(r)}(\cdot)$
 - regression term
 - classification term

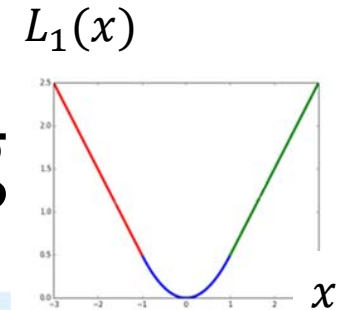
$$L^{(r)}(\{\mathbf{d}^{(i)}, \mathbf{p}^{(i)}\}) = \sum_i \hat{\mathbf{p}}_{obj}^{(i)} \times L_{reg}^{(r)}(\mathbf{d}^{(i)}, \hat{\mathbf{d}}^{(i)}) + \sum_i L_{cls}^{(r)}(\mathbf{p}^{(i)}, \hat{\mathbf{p}}^{(i)})$$

- regression term
 - depend only on **positive** anchors
 - $L_{reg}^{(r)}$: **smooth L_1** function
- classification term
 - $L_{cls}^{(r)}$: **log loss** function



$$L_1(x) = \begin{cases} 0.5x^2, & \text{if } |x| < 1 \\ |x| - 0.5, & \text{otherwise} \end{cases}$$

Faster R-CNN Training



- RPN Loss $L^{(r)}(\cdot)$:
$$\sum_i \hat{p}_{obj}^{(i)} \times L_{reg}^{(r)}(\mathbf{d}^{(i)}, \hat{\mathbf{d}}^{(i)})$$

$$L_{reg}^{(r)}(\mathbf{d}^{(i)}, \hat{\mathbf{d}}^{(i)}) = L_1(\mathbf{d}_x^{(i)} - \hat{\mathbf{d}}_x^{(i)}) + L_1(\mathbf{d}_y^{(i)} - \hat{\mathbf{d}}_y^{(i)}) + L_1(\log(\mathbf{d}_w^{(i)}) - \log(\hat{\mathbf{d}}_w^{(i)})) + L_1(\log(\mathbf{d}_h^{(i)}) - \log(\hat{\mathbf{d}}_h^{(i)}))$$

$$\begin{aligned} \mathbf{d}^{(1)} &= \begin{pmatrix} \mathbf{d}_x & \mathbf{d}_y & \mathbf{d}_w & \mathbf{d}_h \\ 0.4 & 0.2 & 1.2 & 1.2 \end{pmatrix} \\ \hat{\mathbf{d}}^{(1)} &= \begin{pmatrix} 0.125 & 0.125 & 1.0 & 1.1 \end{pmatrix} \\ &= L_1(0.275) + L_1(0.075) + L_1(0.2 - 0.125) \\ &\quad + L_1(\log(1.2) - \log(1.0)) \\ &= (0.5 \times (0.275)^2) + (0.5 \times (0.075)^2) + (0.005) \\ &\quad + (0.5 \times (0.18)^2) + (0.5 \times (0.09)^2) \end{aligned}$$



Faster R-CNN Training

- RPN Loss $L^{(r)}(\cdot)$: $\sum_i L_{cls}^{(r)}(\mathbf{p}^{(i)}, \hat{\mathbf{p}}^{(i)})$

$$L_{cls}^{(r)}(\mathbf{p}^{(i)}, \hat{\mathbf{p}}^{(i)}) = -\left(\hat{\mathbf{p}}_{obj}^{(i)} \times \log(\mathbf{p}_{obj}^{(i)}) + \hat{\mathbf{p}}_{bg}^{(i)} \times \log(\mathbf{p}_{bg}^{(i)})\right)$$

$$\mathbf{p}^{(1)} = \begin{pmatrix} p_{obj} & p_{bg} \\ 0.6 & 0.4 \end{pmatrix}$$

$$\hat{\mathbf{p}}^{(1)} = \begin{pmatrix} 1.00 & 0.0 \end{pmatrix}$$

$$\mathbf{p}^{(2)} = \begin{pmatrix} 0.2 & 0.8 \end{pmatrix}$$

$$\hat{\mathbf{p}}^{(2)} = \begin{pmatrix} 0.00 & 1.0 \end{pmatrix}$$

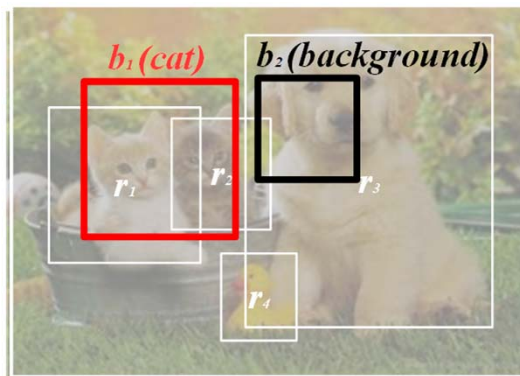
$$\sum_i L_{cls}^{(r)}(\mathbf{p}^{(i)}, \hat{\mathbf{p}}^{(i)})$$

$$= -\left(\hat{\mathbf{p}}_{obj}^{(1)} \times \log(\mathbf{p}_{obj}^{(1)}) + \hat{\mathbf{p}}_{bg}^{(1)} \times \log(\mathbf{p}_{bg}^{(1)})\right)$$

$$- \left(\hat{\mathbf{p}}_{obj}^{(2)} \times \log(\mathbf{p}_{obj}^{(2)}) + \hat{\mathbf{p}}_{bg}^{(2)} \times \log(\mathbf{p}_{bg}^{(2)})\right)$$

Faster R-CNN Training

- Detection Loss $L^{(d)}(.)$
 - A mini-batch for computing detection loss consists of 128 proposals (anchors)
 - 32 object proposals (max. IoU ≥ 0.5)
 - 96 background proposals ($0.1 \leq \text{max. IoU} < 0.5$)

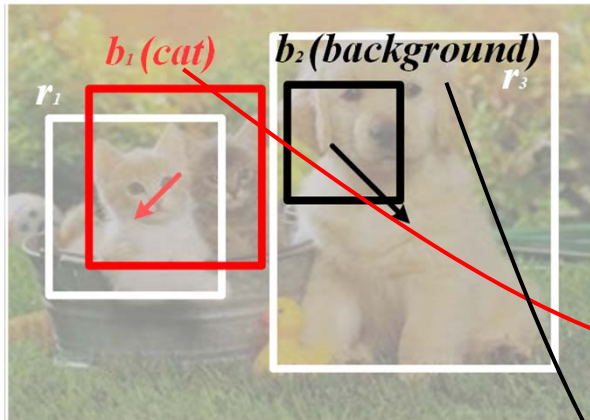


	<i>cat</i>	<i>cat</i>	<i>dog</i>	<i>duck</i>		
	<i>r</i> ₁	<i>r</i> ₂	<i>r</i> ₃	<i>r</i> ₄		
<i>b</i> ₁	0.8	0.6	0.0	0.0	≥ 0.5	<i>cat</i> proposal
<i>b</i> ₂	0.0	0.1	0.2	0.0	< 0.5	<i>bg</i> proposal

IoU: Intersection over Union

Faster R-CNN Training

- Detection Loss $L^{(d)}(.)$
 - $D = \{\hat{d}^{(i)}, \hat{p}^{(i)}\}_{i=1}^{128}$: selected proposals



$$\hat{d}^{(1)} = \begin{pmatrix} d_x & d_y & d_w & d_h \\ -0.2 & 0.2 & 1.0 & 1.0 \end{pmatrix}$$

$$\hat{d}^{(2)} = \begin{pmatrix} 1.0 & 1.0 & 2.0 & 2.0 \end{pmatrix}$$

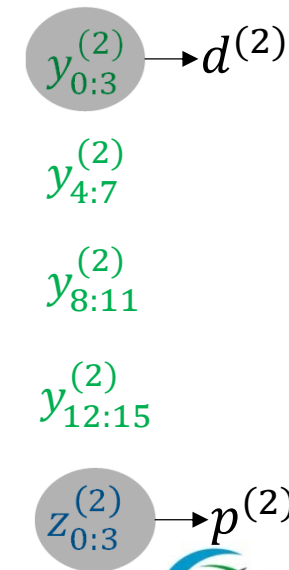
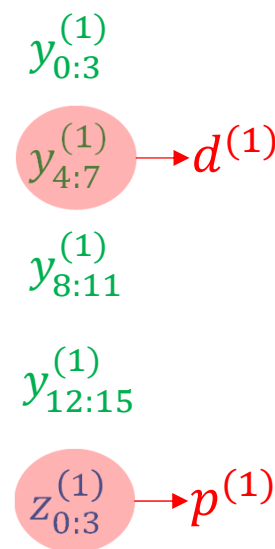
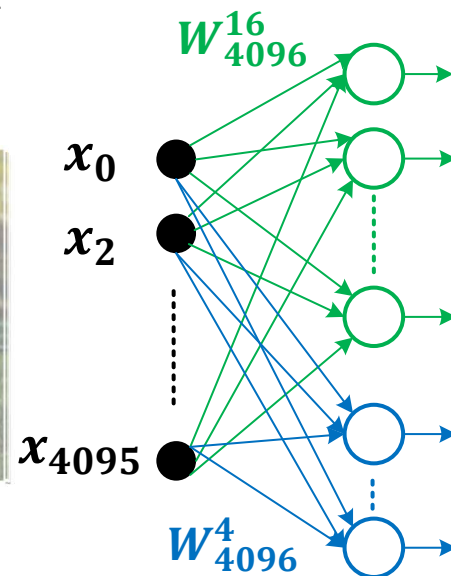
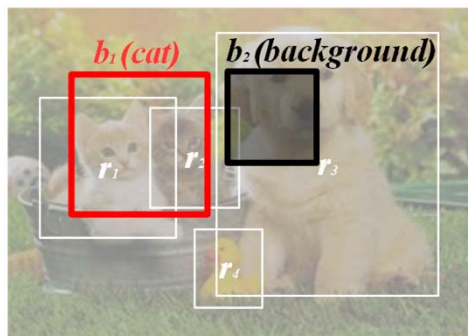
$$\hat{p}^{(1)} = \begin{pmatrix} p_{bg} & p_{cat} & p_{dog} & p_{duck} \\ 0.0 & \underline{1.0} & 0.0 & 0.0 \end{pmatrix}$$

$$\hat{p}^{(2)} = \begin{pmatrix} \underline{1.0} & 0.0 & 0.0 & 0.0 \end{pmatrix}$$

Faster R-CNN Training

- Detection Loss $L^{(d)}(.)$
 - $\{d^{(i)}, p^{(i)}\}$: predicted parameters from detection network.

$\{\text{bg, cat, dog, duck}\}$





Faster R-CNN Training

- Detection Loss $L^{(d)}(\{d^{(i)}, p^{(i)}\})$

$$= \sum_i \underbrace{(1.0 - \hat{p}_{bg}^{(i)}) \times L_{reg}^{(d)}(d^{(i)}, \hat{d}^{(i)})}_{\text{regression term}} + \sum_i \underbrace{L_{cls}^{(d)}(p^{(i)}, \hat{p}^{(i)})}_{\text{classification term}}$$

$$y_{4:7}^{(1)} \hat{p}^{(1)} = d^{(1)} = \begin{pmatrix} p_{bg}^{(1)} & d_x & p_{cat} & d_y & p_{dog} & d_w & p_{duck} & d_h \\ (-0.1 & 0.3 & 0.8 & 1.2) \\ (0.0 & 1.0 & 0.0 & 0.0) \end{pmatrix}$$

$$\hat{p}^{(2)} = \uparrow \underbrace{(L_{reg}^{(d)}(d^{(i)}, \hat{d}^{(i)}))}_{\text{RPN regression}} \times \text{RPN regression}$$

$$\hat{d}^{(1)} = (-0.2 \quad 0.2 \quad 1.0 \quad 1.0)$$



Faster R-CNN Training

- Detection Loss $L^{(d)}(\{d^{(i)}, p^{(i)}\})$

$$= \sum_i (\mathbf{1.0} - \hat{p}_{bg}^{(i)}) \times L_{reg}^{(d)}(d^{(i)}, \hat{d}^{(i)}) + \sum_i L_{cls}^{(d)}(p^{(i)}, \hat{p}^{(i)})$$

classification term

$$L_{cls}^{(d)}(p^{(i)}, \hat{p}^{(i)}) = -\hat{p}_{bg}^{(i)} \times \log(p_{bg}^{(i)}) + \hat{p}_{cat}^{(i)} \times \log(p_{cat}^{(i)})$$

$$z_{0:3}^{(1)} = p^{(1)} = \begin{pmatrix} 0.1 & 0.7 & 0.1 & 0.1 \end{pmatrix}$$

$$z_{0:3}^{(2)} = p^{(2)} = \begin{pmatrix} 0.6 & 0.0 & 0.2 & 0.2 \end{pmatrix}$$

$$\hat{p}^{(1)} = \begin{pmatrix} \hat{p}_{dog}^{(1)} & \times & \log(p_{dog}^{(1)}) & + & \hat{p}_{duck}^{(1)} & \times & \log(p_{duck}^{(1)}) \end{pmatrix}$$

$$\hat{p}^{(2)} = \begin{pmatrix} 1.0 & 0.0 & 0.0 & 0.0 \end{pmatrix}$$

