# MIL-UT at ILSVRC2015

## Masataka Yamaguchi, Qishen Ha, Katsunori Ohnishi, Masatoshi Hidaka, Yusuke Mukuta, Tatsuya Harada
### The University of Tokyo

MIL: Machine Intelligence Laboratory
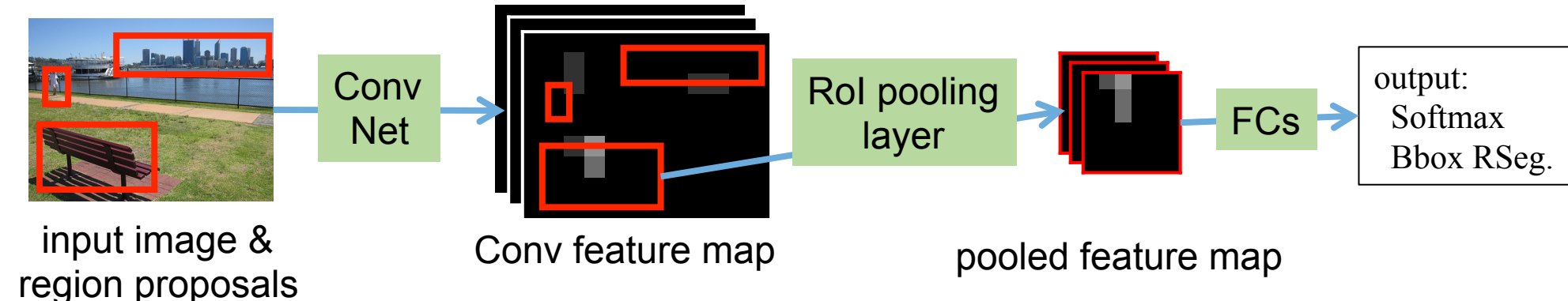
東京大学 THE UNIVERSITY OF TOKYO

## Overview

We use **Fast-RCNN** [Girshick, 2015] as the base detection system and **VGG-16** [Simonyan and Zisserman, 2014] as the base model. We improve the detection accuracy by **concatenating the whole image features** with the fc7-layer output and using it as the input to the inner product layer before Softmax.

In addition, we demonstrate that **replacing pool4 layer** rather than pool5 layer **with RoI pooling layer** improves mAP.
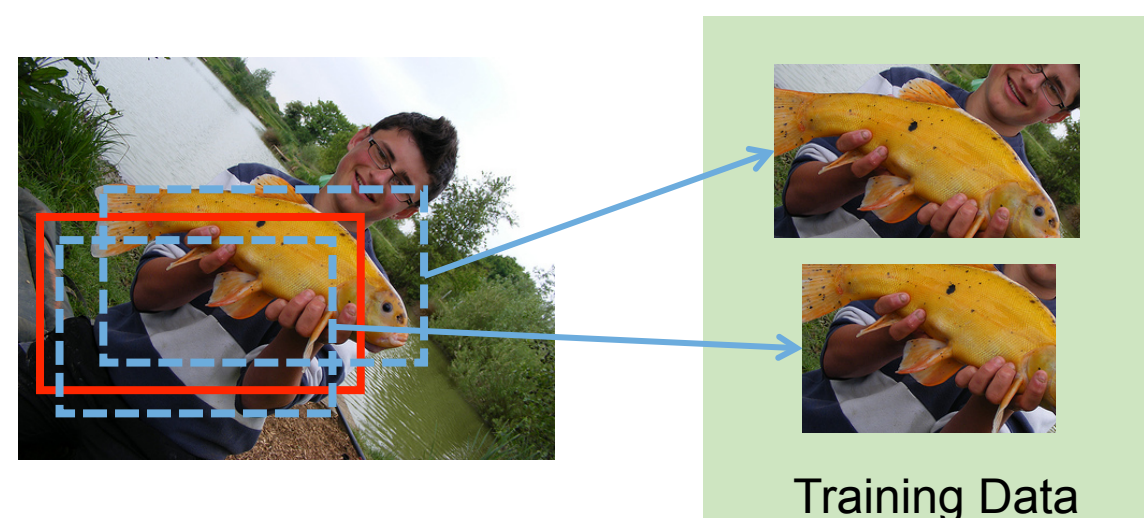
We submitted two results. One is obtained by model fusion using the same weights for all models and the other is obtained by **model fusion using weights learned by Bayesian optimization** on the val2 dataset.

## Fast R-CNN



input image & region proposals — Conv Net — Conv feature map — RoI pooling layer — pooled feature map — FCs — output: Softmax Bbox RSeg.
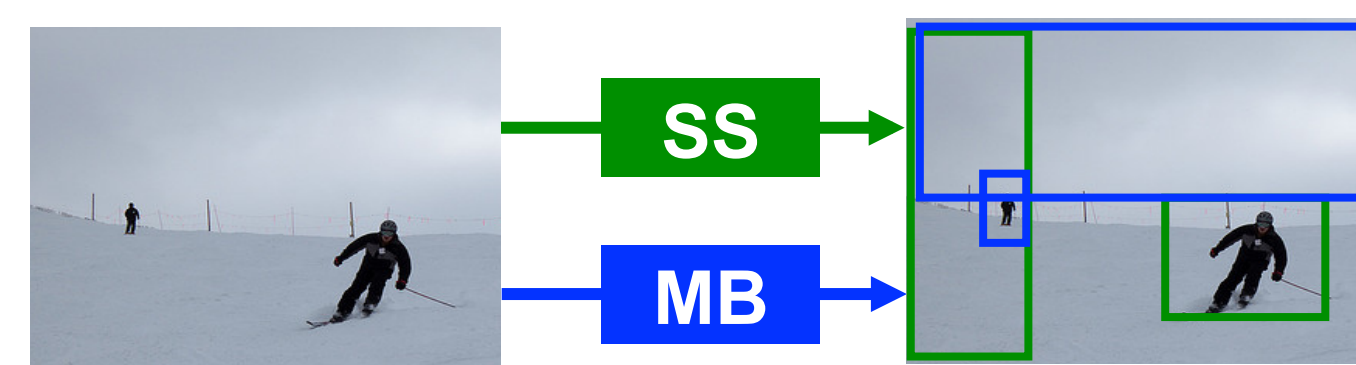
- Ross Girshick, "Fast R-CNN", ICCV 2015
- CNN model
  - VGG-16 (Simonyan and Zisserman, "Very deep convolutional networks for large-scale image recognition." arXiv 2014)
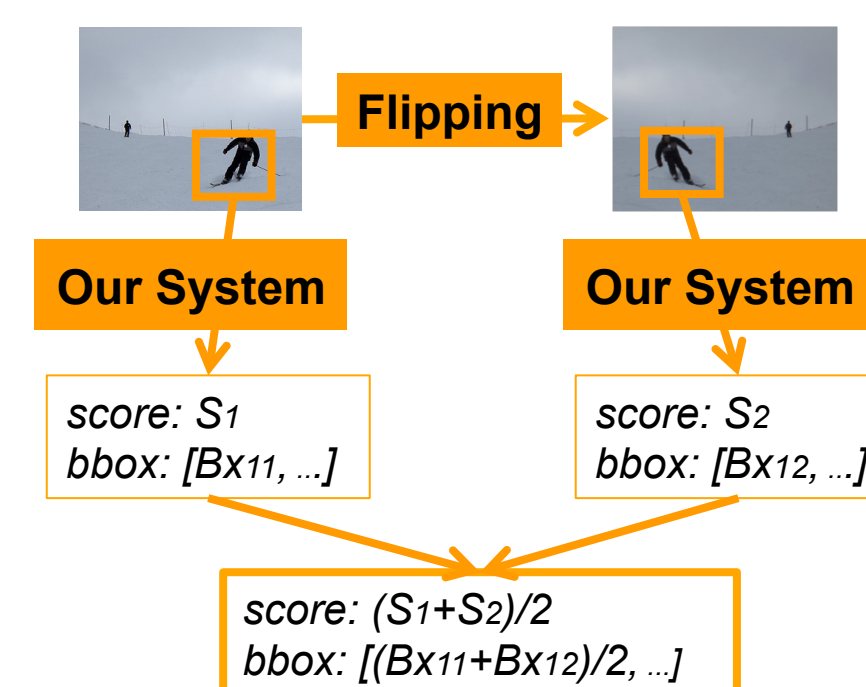- Train and test on the single-scale mode

## Techniques to Improve mAP



Training Data

- Retrain VGG-16 on annotated bounding boxes in CLS-LOC and DET dataset
  - DeepID-net [OuyangTODO]
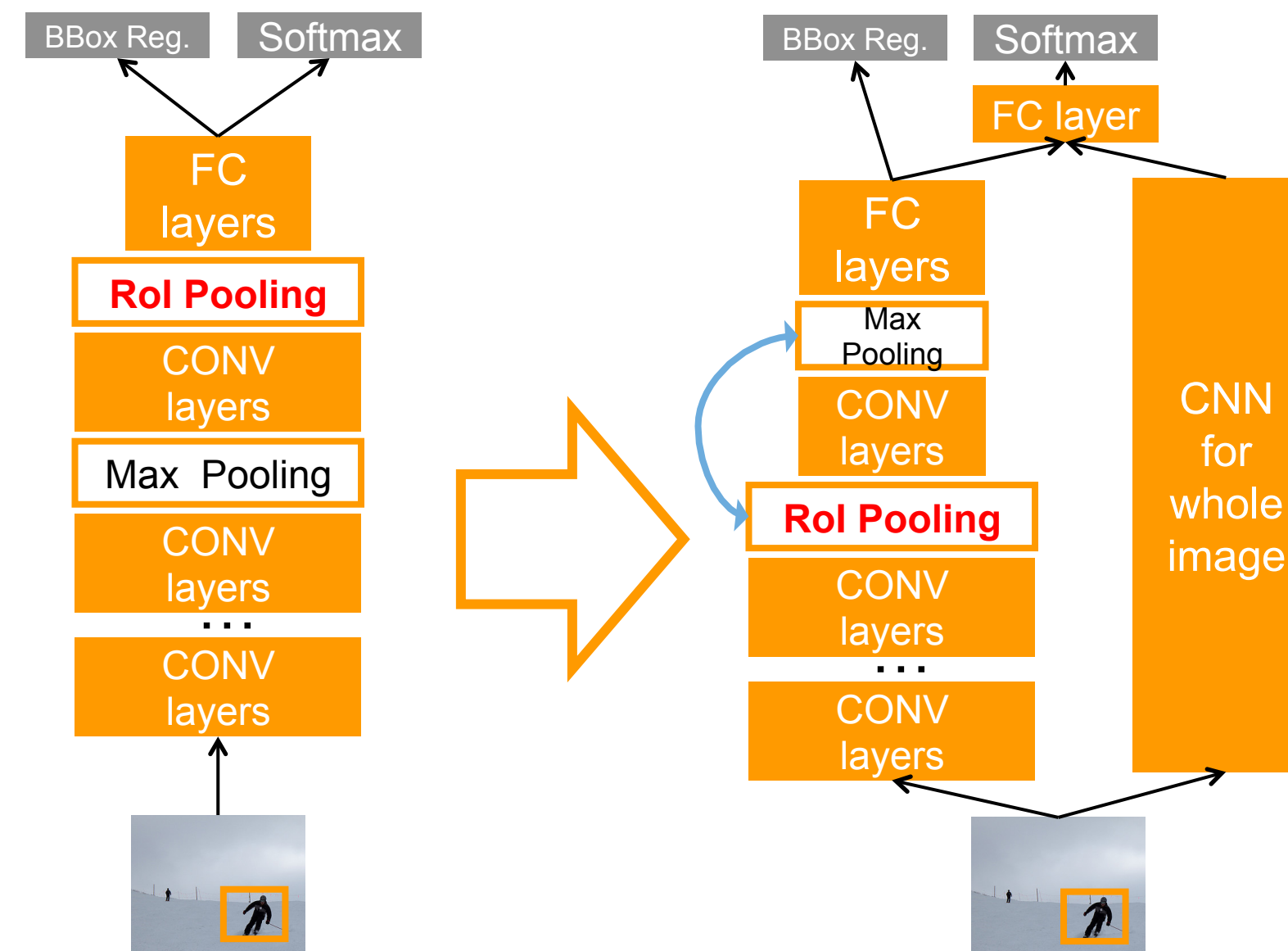  - We follow the RCNN framework[??]
- Use multiple region proposal methods when testing
  - Selective Search [TODO refer]
    - when training and testing
  - Multibox [TODO refer]
    - when testing
- Test not only the original images but also horizontally-flipped ones and combine them

score: $S_1$ / bbox: [$Bx_{11}$, ...]
score: $S_2$ / bbox: [$Bx_{12}$, ...]
score: ($S_1$+$S_2$)/2 / bbox: [($Bx_{11}$+$Bx_{12}$)/2, ...]

## Improvement on the Network Architecture



- Context Modeling
  - Concatenate the whole image features with the fc7-layer output and use it as the input to the inner product layer before Softmax
  - Fix the weights of CNN for whole image except its last FC layer for simplicity
- Replacing "pool4" layer with RoI pooling layer
  - Computational time per iteration gets 1.5 times slower

| VOC2007 mAP | Pool5→RoI Pool | Pool4→RoI Pool |
|---|---|---|
| w/o whole image feature | 66.7 | 68.9 |
| with whole image feature | 67.8 | 70.1 |

## Model Ensemble

- 1. Ensemble using same weights for all models
- 2. Ensemble using weights learned by Bayesian Optimization
  - At step t+1, choose weights $w_{t+1}$ as follows:

$$w_{t+1} = \underset{w \in D}{\operatorname{argmax}} \int \max\left(0, y - \max_{i=1,\cdots,t} AP(w_i)\right) P(y \mid w, U_t)\, dy$$

where
$$D = \left\{ w \mid \sum_i^N w_i = 1, 0 \le w_i \le 1 \right\}, U_t = \{(w_i, AP(w_i)) \mid i = 1, 2, \cdots, t\}$$
$N$ : the number of models

  - Learn weights separately for each class on Val2 (see R-CNN paper)

## Results of DET Task



Example Results

- mAP Using Single Model (on Val2)

| | Pool5→RoI | Pool4→RoI |
|---|---|---|
| original VGG-16 | 42.9 | 44.2 |
| pre-trained on annotated boxes | 45.6 | 45.6 |

- mAP Using Single Model (on Val2 and Test)

| | mAP on Val2 | mAP on Test |
|---|---|---|
| by averaging | 48.1 | todo |
| using weights learned by BO | 50.6 | todo |

## Acknowledgement

IBM    NVIDIA