

Localization-Aware Active Learning for Object Detection

Chieh-Chi Kao

University of California, Santa Barbara

`chiehchi.kao@gmail.com`

Teng-Yok Lee

Mitsubishi Electric Research Laboratories

`tlee@merl.com`

Pradeep Sen

University of California, Santa Barbara

`psen@ece.ucsb.edu`

Ming-Yu Liu

Mitsubishi Electric Research Laboratories

`seanmingyuliu@gmail.com`

Contents

- Introduction
- Background
- Method
- Experiment

Introduce

- Active learning has been shown to be effective at annotating data for image classification, but largely unexplored for object detection.
- We present two metrics for measuring the informativeness of an object hypothesis.
 - Localization Tightness (TL).
 - Localization Stability (LS).

Background

- Active learning.

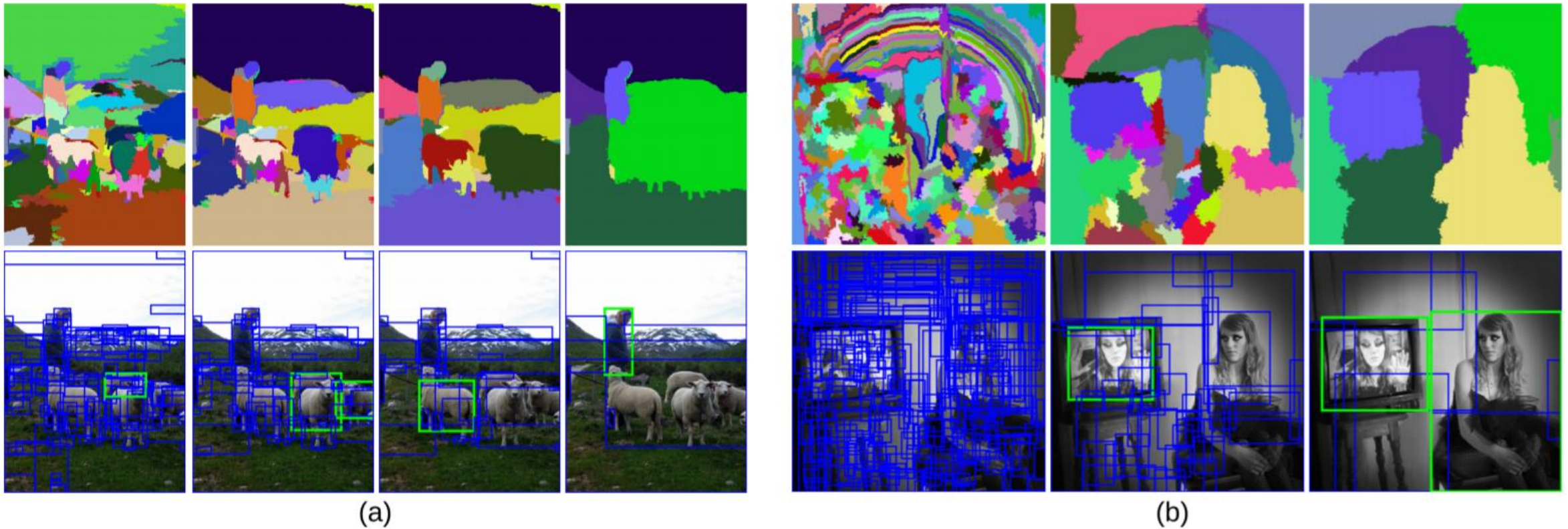
region proposal



- Object Detection.

- Two stage: $img \rightarrow Feature\ extraction \rightarrow RP \rightarrow classification/ Locating\ regression$
 - **R-CNN**、**Fast R-CNN**、**Faster R-CNN**
 - SPP-NET
- One stage: $img \rightarrow Feature\ extraction \rightarrow classification/ Locating\ regression$
 - OverFeat
 - **YOLOv1**、**YOLOv2**、**YOLOv3**
 - SSD

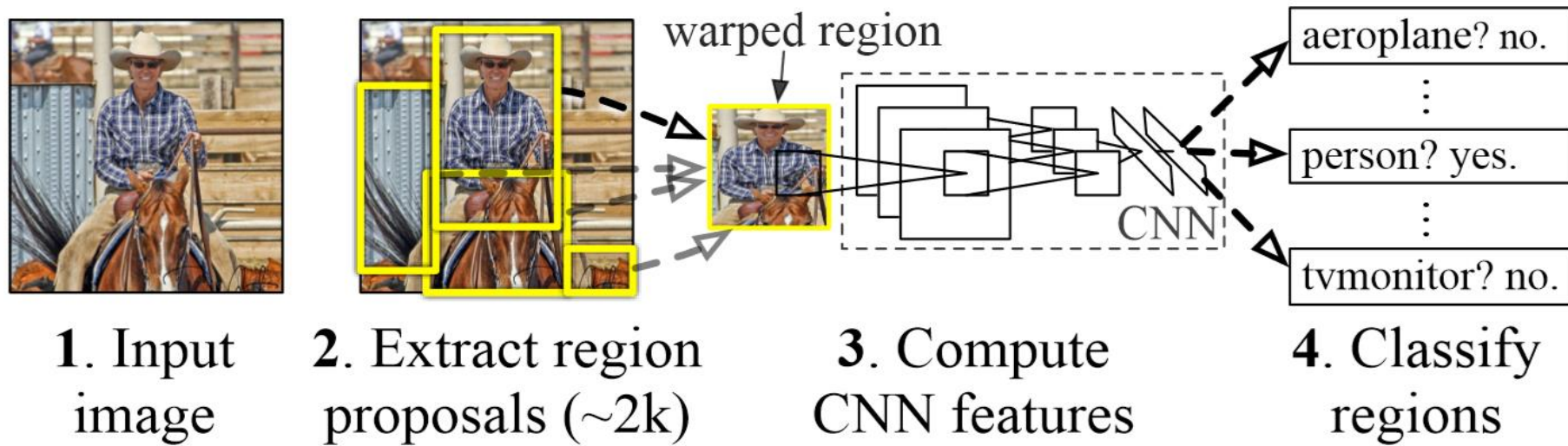
Background – RCNN (regions with cnn)



Selective Search for Object Recognition (IJCV 2013)

Background – RCNN (Regions with CNN)

R-CNN: *Regions with CNN features*

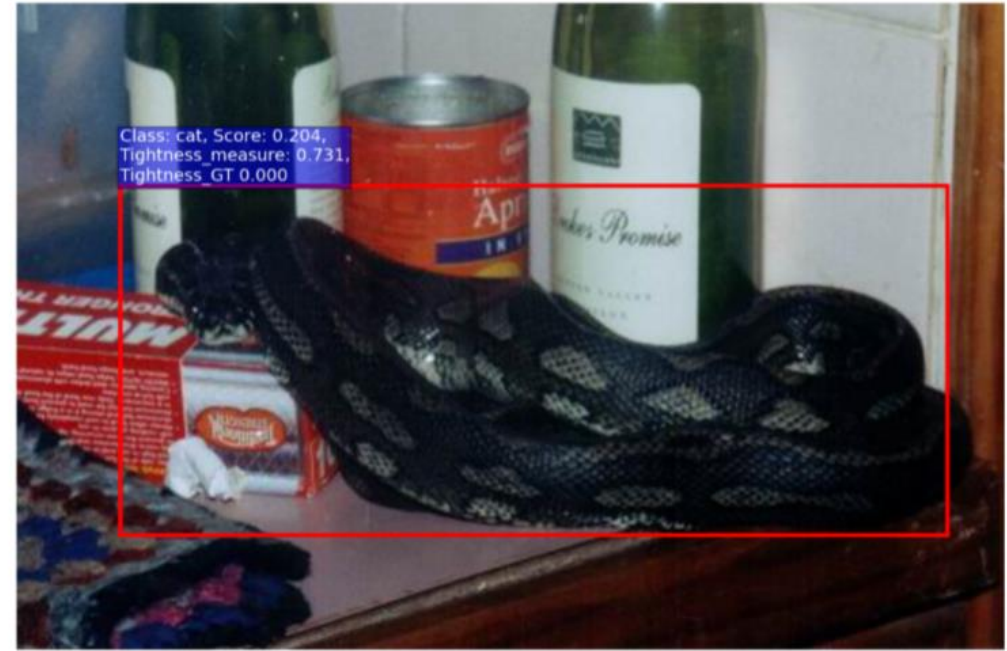


- Process

- Using SS to extract region proposals
- Score each region proposals
- Rank the scores of the candidate boxes, then remove some regions according to some rules. (NMS)
- Make a border regression of the last remaining boxes. (Position、 Size)



(a)



(b)

- Given a predicted box that is absolutely certain about its classification result ($P_{max} = 1$), but it cannot tightly enclose a true object ($T = 0$).
- Reversely, if the predicted box can tightly enclose a true object ($T = 1$), but the classification result is uncertain (*low* P_{max}).

$$J(B_0^j) = |T(B_0^j) + P_{max}(B_0^j) - 1| \quad T_I(I_i) = \min_j J(B_0^j) \quad \downarrow \text{select}$$

Method

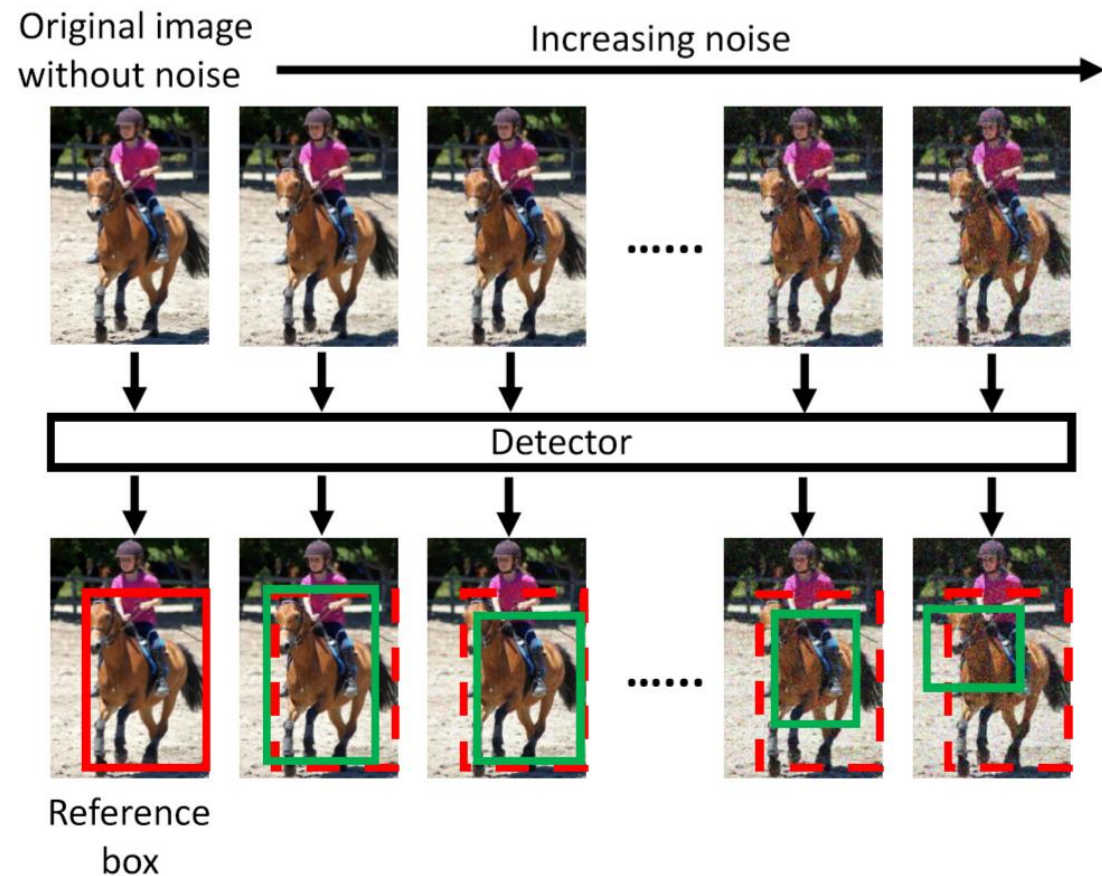
- Localization Stability

If the current model is stable to noise, meaning that the detection result does not dramatically change even if the input unlabeled image is corrupted by noise, the current model already understands this unlabeled image well so there is no need to annotate this unlabeled image.

$$S_B(B_0^j) = \frac{\sum_{n=1}^N IoU(B_0^j, C_n(B_0^j))}{N}$$

select ↓

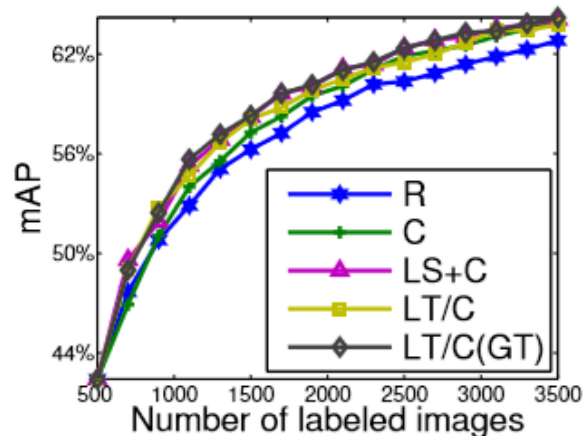
$$S_I(I_i) = \frac{\sum_{j=1}^M P_{max}(B_0^j) S_B(B_0^j)}{\sum_{j=1}^M P_{max}(B_0^j)}$$



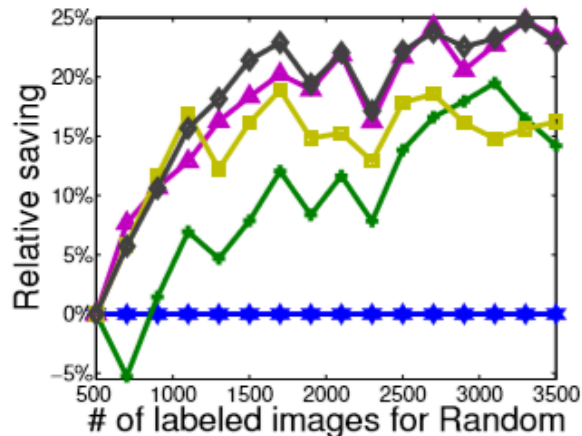
Experiment

- Two baselines:
 - Random (R)
 - Classification Uncertainty (C)
 - Two baselines:
- Ours:
 - Localization stability only (LS)
 - Localization stability + Classification Uncertainty (LS+C)
 - Localization tightness + Classification Uncertainty (LT/C)
 - Localization tightness + Classification Uncertainty + Ground-truth boxes instead of the estimate used in LT/C (LT/C(GT))

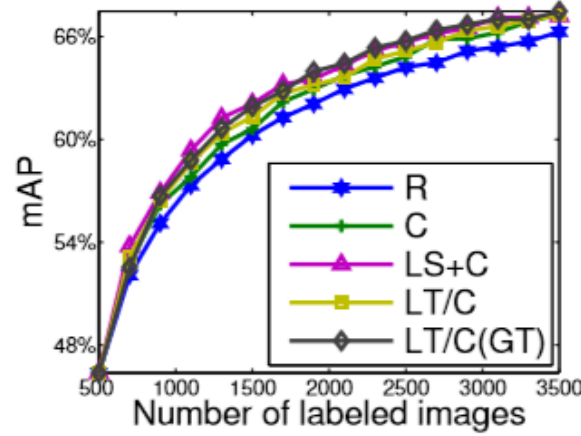
Experiment – FRCNN on PASCAL 2012 and PASCAL 2007



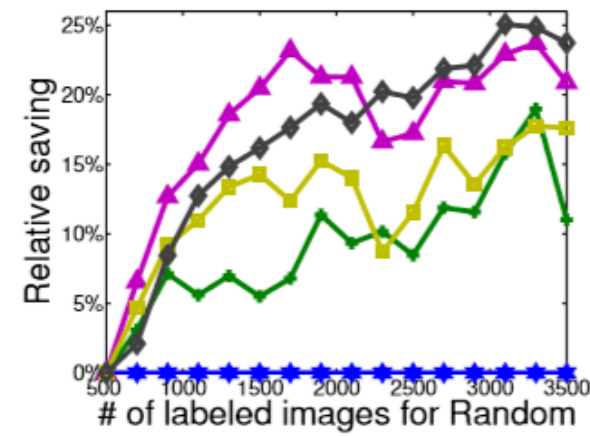
(a) mAP



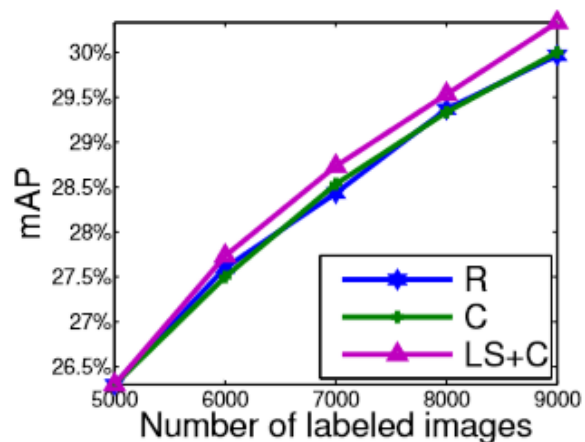
(b) Saving



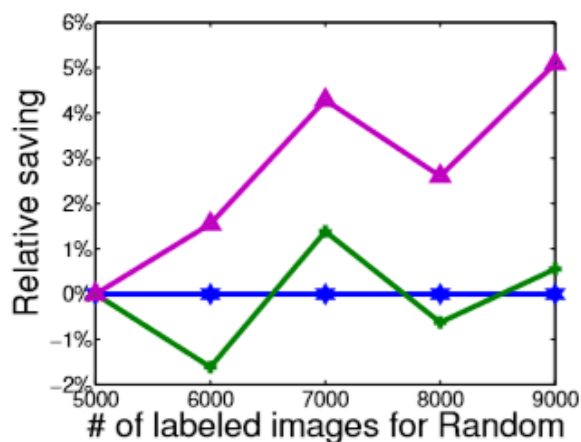
(a) mAP



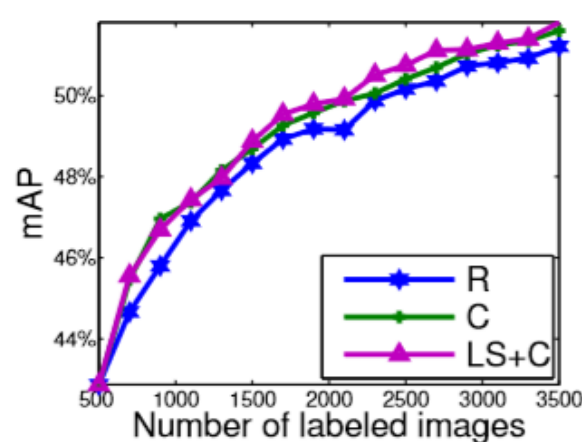
(b) Saving



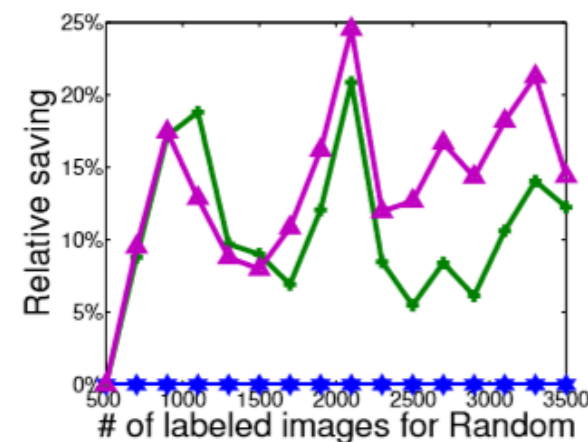
(a) mAP



(b) Saving



(a) mAP



(b) Saving

Experiment – FRCNN on PASCAL 2012 and PASCAL 2007

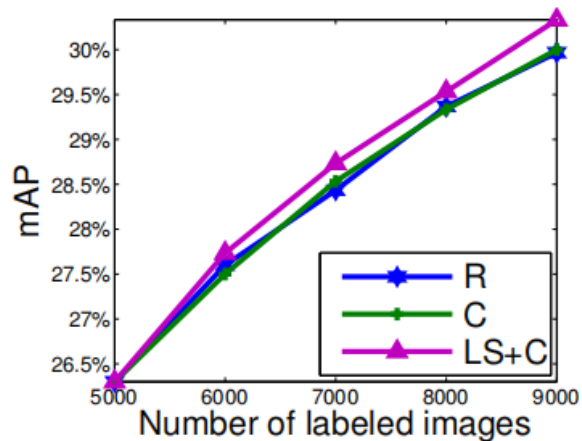
method	aero	bike	bird	boat*	bottle*	bus	car	cat	chair*	cow	table*	dog	horse	mbike	persn	plant*	sheep	sofa	train	tv	mAP
R	<u>71.1</u>	61.5	<u>54.7</u>	28.4	32.0	<u>68.1</u>	57.9	75.4	25.8	44.2	36.4	73.0	61.9	67.3	68.1	21.6	51.9	41.0	65.5	51.7	52.9
C	70.7	62.9	54.7	25.5	30.8	66.1	56.2	78.1	26.4	54.5	36.7	76.9	68.3	<u>67.7</u>	67.4	22.5	<u>57.7</u>	40.8	63.6	52.5	54.0
LS+C	73.9	<u>63.7</u>	56.9	29.6	35.2	66.5	<u>58.5</u>	<u>77.9</u>	31.3	<u>50.8</u>	<u>40.7</u>	<u>73.8</u>	<u>65.4</u>	66.9	<u>68.4</u>	24.8	58.0	44.9	64.2	<u>53.9</u>	55.3
LT/C	69.8	64.6	54.6	<u>29.5</u>	<u>33.8</u>	70.3	59.7	75.5	<u>29.5</u>	46.3	41.8	73.0	62.5	69.0	70.8	<u>23.2</u>	56.5	<u>42.8</u>	<u>64.3</u>	55.9	<u>54.7</u>

PASCAL 2012

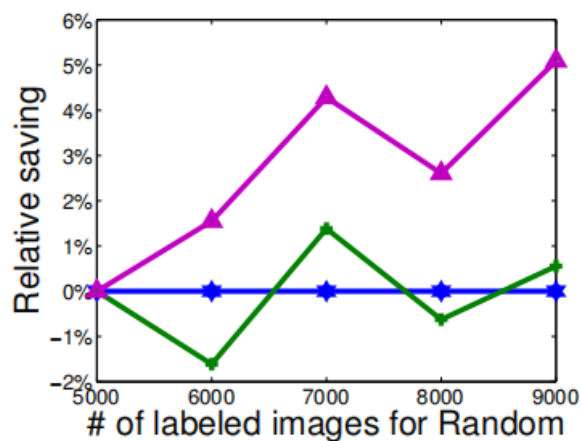
method	aero	bike	bird	boat*	bottle*	bus	car	cat	chair*	cow	table	dog	horse	mbike	persn	plant*	sheep	sofa	train	tv	mAP
R	61.6	67.2	54.1	40.0	33.6	64.5	73.0	73.9	34.5	60.8	52.2	69.3	74.7	<u>66.6</u>	<u>67.1</u>	25.9	52.1	54.2	66.1	54.9	57.3
C	56.9	<u>68.0</u>	<u>54.9</u>	36.8	34.4	<u>68.1</u>	71.7	75.5	34.0	68.6	51.0	<u>71.4</u>	<u>74.7</u>	65.2	65.9	24.9	<u>60.0</u>	53.9	63.0	<u>57.4</u>	57.8
LS+C	<u>61.5</u>	64.4	55.8	<u>40.2</u>	38.7	66.3	<u>73.8</u>	<u>74.7</u>	39.6	<u>68.0</u>	<u>56.3</u>	71.5	73.8	67.2	66.7	<u>27.7</u>	61.3	57.0	<u>65.6</u>	57.4	59.4
LT/C	<u>57.6</u>	69.7	52.9	41.1	<u>38.4</u>	69.7	74.4	71.8	<u>36.4</u>	61.2	58.1	69.5	74.3	66.2	67.8	28.0	55.5	<u>56.3</u>	65.5	58.2	<u>58.6</u>

PASCAL 2007

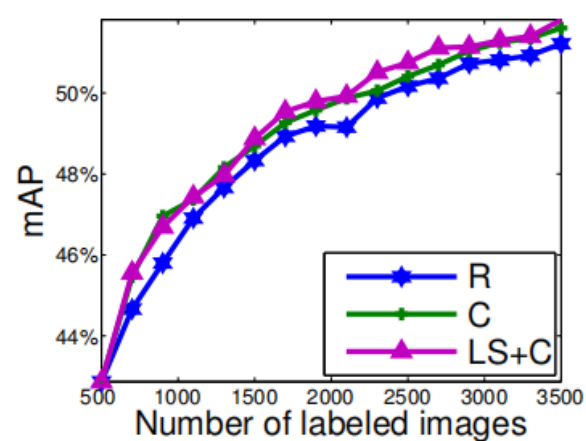
Experiment – FRCNN on MS COCO 2012 and SSD on PASCAL 2007



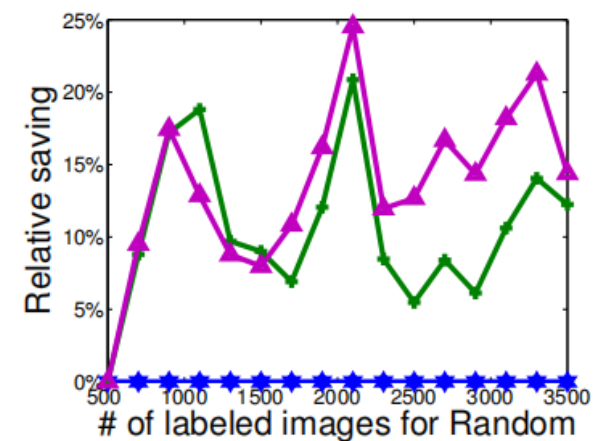
(a) mAP



(b) Saving



(a) mAP



(b) Saving

FRCNN on MS COCO 2012

SSD on PASCAL 2007