

RCNN, Fast RCNN, Faster RCNN



Presented by: Roi Shikler & Gil Elbaz
Advisor: Prof. Michael Lindenbaum

Topics of the lecture:

- Problem statement
- Review of slow R-CNN
- Review of Fast R-CNN
- Review of Faster R-CNN
- Compare with other methods
- Take away

Topics of the lecture:

→ Problem statement

- Review of slow R-CNN
- Review of Fast R-CNN
- Review of Faster R-CNN
- Compare with other methods
- Take away

PASCAL Visual Object Classes



"The main goal of this challenge is to **recognize objects** from a number of visual object classes in realistic scenes"

[from the challenge homepage]

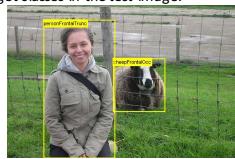


Desired output

PASCAL Visual Object Classes

The challenge face the participants in front of two problems:

- **Classification:** For each of the twenty classes, predicting presence/absence of an example of that class in the test image.
- **Detection (localization):** Predicting the bounding box and label of each object from the twenty target classes in the test image.

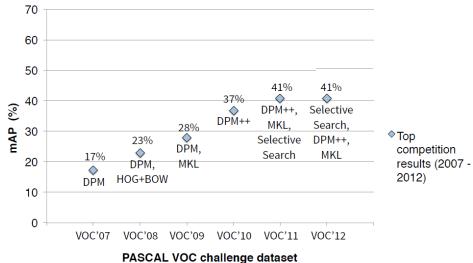


PASCAL Visual Object Classes

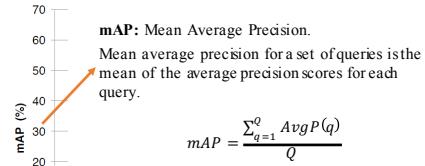
The twenty object classes that have been selected are:

- *Person*: person
- *Animal*: bird, cat, cow, dog, horse, sheep
- *Vehicle*: airplane, bicycle, boat, bus, car, motorbike, train
- *Indoor*: bottle, chair, dining table, potted plant, sofa, TV/monitor

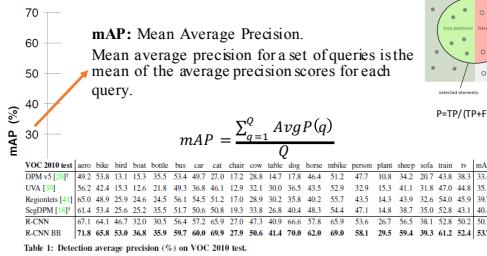
PASCAL VOC detection history



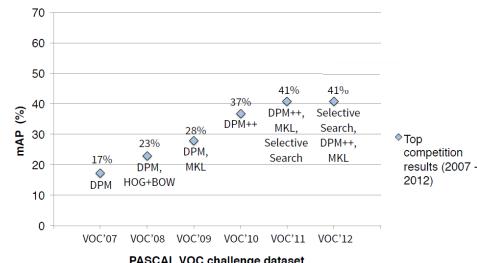
PASCAL VOC detection history



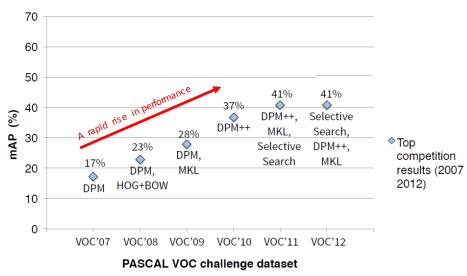
PASCAL VOC detection history



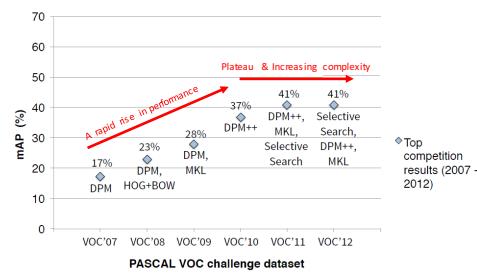
PASCAL VOC detection history



PASCAL VOC detection history



PASCAL VOC detection history



Topics of the lecture:

✓ Problem statement

➡ Review of slow R-CNN

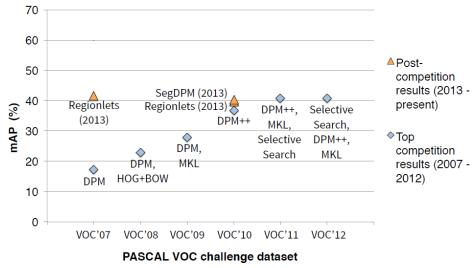
- Review of Fast R-CNN
- Review of Faster R-CNN
- Compare with other methods
- Take away

R-CNN - Article (2013)

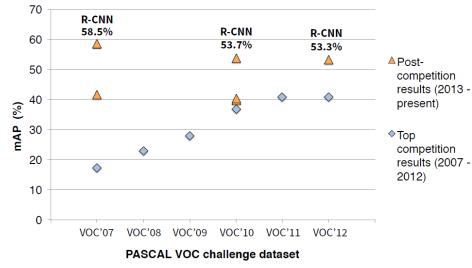
Rich feature hierarchies for accurate object detection and semantic segmentation
Tech report

Ross Girshick¹ Jeff Donahue^{1,2} Trevor Darrell^{1,2} Jitendra Malik¹
¹UC Berkeley and ²ICSI
{rgb, jdonahue, trevor, malik}@eecs.berkeley.edu

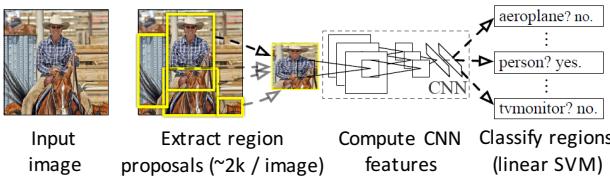
PASCAL VOC detection history



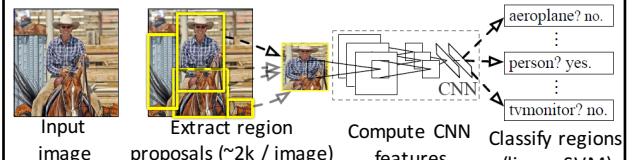
PASCAL VOC detection history



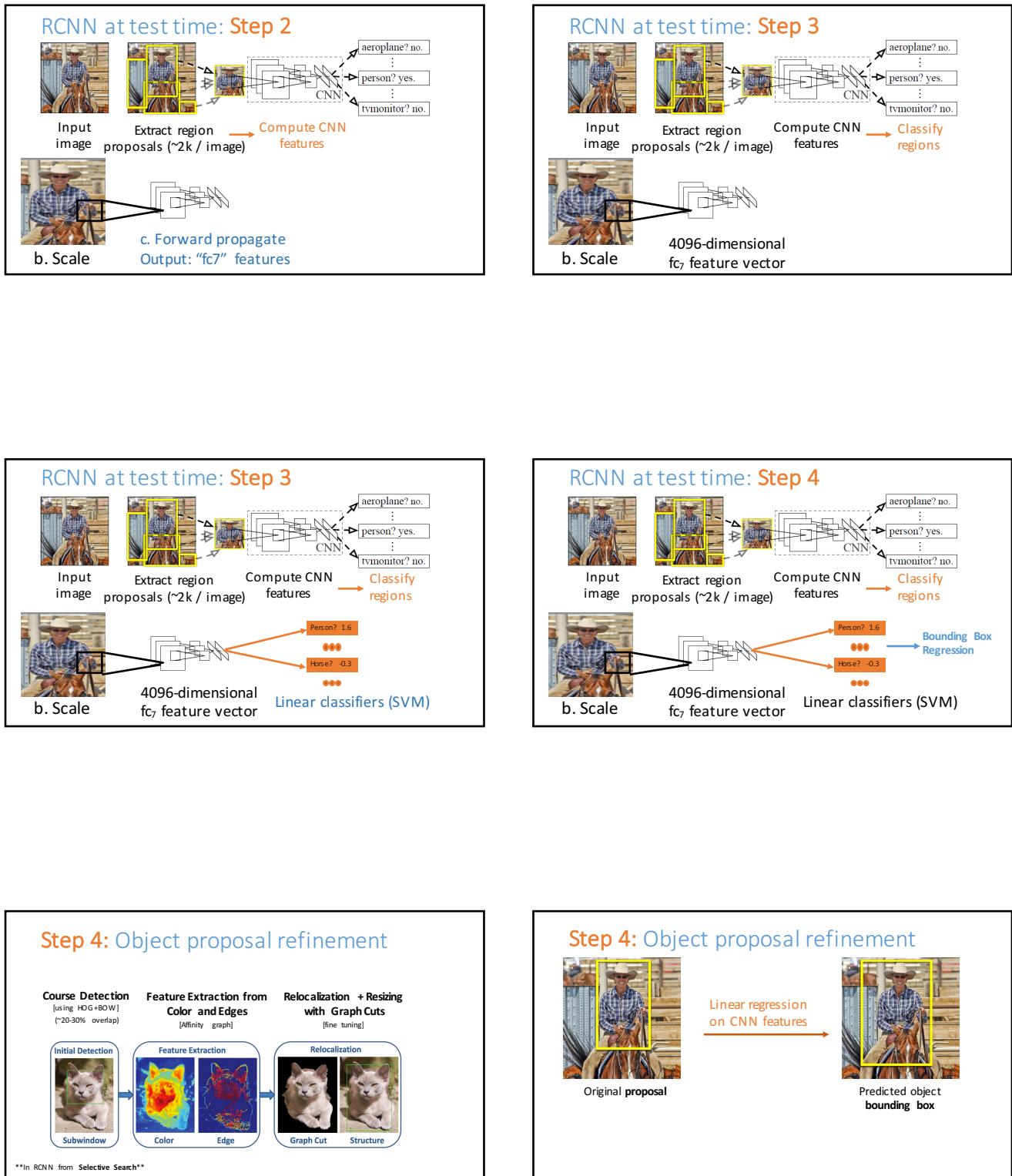
R-CNN: Regions with CNN features



RCNN at test time: Step 1







Review of the slow R-CNN training pipeline

Steps for training a slow R-CNN detector:

1. [offline] $M \leftarrow$ Pre-train a ConvNet for ImageNet classification

Review of the slow R-CNN training pipeline

Steps for training a slow R-CNN detector:

1. [offline] $M \leftarrow$ Pre-train a ConvNet for ImageNet classification
2. $M' \leftarrow$ Fine-tune M for object detection (softmax classifier)

Review of the slow R-CNN training pipeline

Steps for training a slow R-CNN detector:

1. [offline] $M \leftarrow$ Pre-train a ConvNet for ImageNet classification
2. $M' \leftarrow$ Fine-tune M for object detection (softmax classifier)
3. $F \leftarrow$ Cache feature vectors to disc using M'

Review of the slow R-CNN training pipeline

Steps for training a slow R-CNN detector:

1. [offline] $M \leftarrow$ Pre-train a ConvNet for ImageNet classification
2. $M' \leftarrow$ Fine-tune M for object detection (softmax classifier)
3. $F \leftarrow$ Cache feature vectors to disc using M'
4. Train post hoc linear SVMs on F for object classification

* “post hoc” means the parameters are learned after the ConvNet is fixed

Review of the slow R-CNN training pipeline

Steps for training a slow R-CNN detector:

1. [offline] $M \leftarrow$ Pre-train a ConvNet for ImageNet classification
2. $M' \leftarrow$ Fine-tune M for object detection (softmax classifier)
3. $F \leftarrow$ Cache feature vectors to disc using M'
4. Train post hoc linear SVMs on F for object classification
5. Train post hoc linear bounding-box regressors on F

* “post hoc” means the parameters are learned after the ConvNet is fixed

Review slow R-CNN

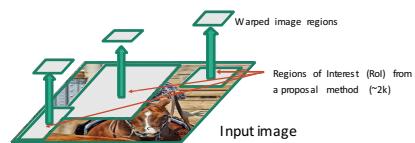


Input image

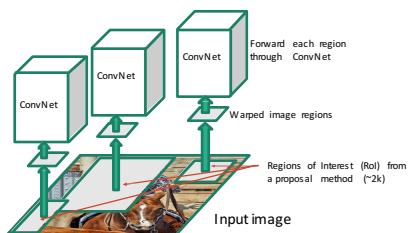
Review slow R-CNN



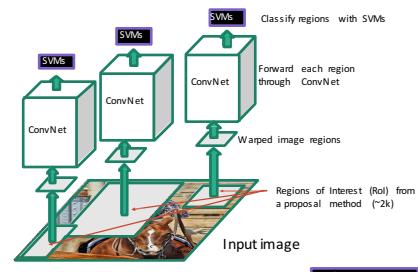
Review slow R-CNN



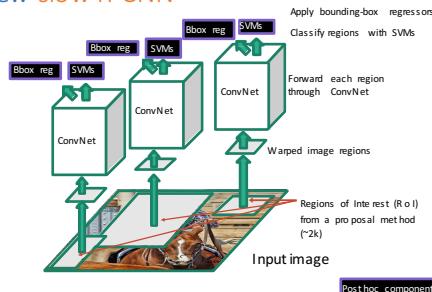
Review slow R-CNN



Review slow R-CNN



Review slow R-CNN



What's wrong with slow R-CNN?

What's wrong with slow R-CNN?

- Ad hoc training objectives

- Fine tune network with softmax classifier (log loss)
- Train post-hoc linear SVMs (hinge loss)
- Train post hoc bounding-box regressors (squared loss)

What's wrong with slow R-CNN?

- Ad hoc training objectives

- Fine tune network with softmax classifier (log loss)
- Train post-hoc linear SVMs (hinge loss)
- Train post hoc bounding-box regressors (squared loss)

- Training is slow, takes a lot of disk space

- The training process takes about 84h

What's wrong with slow R-CNN?

- Ad hoc training objectives

- Fine tune network with softmax classifier (log loss)
- Train post-hoc linear SVMs (hinge loss)
- Train post hoc bounding-box regressors (squared loss)

- Training is slow, takes a lot of disk space

- The training process takes about 84h

- Detection (inference) is slow

- About 47s per image
- ~2k ConvNet forward passes per image

What's wrong with slow R-CNN?

- Ad hoc training objectives

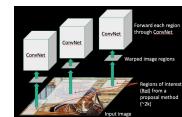
- Fine tune network with softmax classifier (log loss)
- Train post-hoc linear SVMs (hinge loss)
- Train post hoc bounding-box regressors (squared loss)

- Training is slow, takes a lot of disk space

- The training process takes about 84h

- Detection (inference) is slow

- About 47s per image
- ~2k ConvNet forward passes per image



Topics of the lecture:

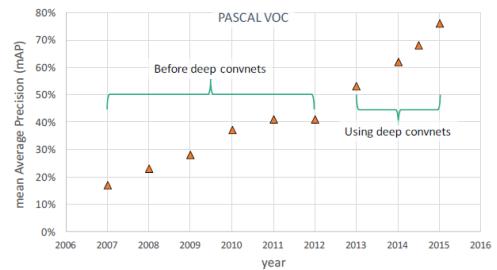
✓ Problem statement

✓ Review of slow R-CNN

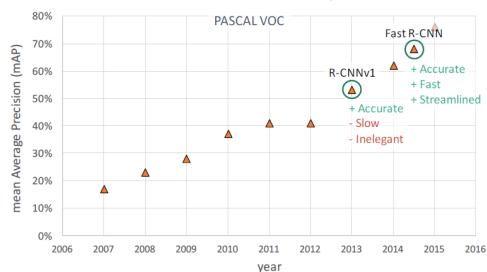
➡ Review of Fast R-CNN

- Review of Faster R-CNN
- Compare with other methods
- Take away

PASCAL VOC detection history



PASCAL VOC detection history



Fast R-CNN objectives

Fix most of what's wrong with slow R-CNN and SPP-net

- Train the detector in a **single stage**, end-to-end
 - No caching features to disk
 - No post hoc training steps
- Train **all layers** of the network
 - [Training the conv layers is important for very deep networks]

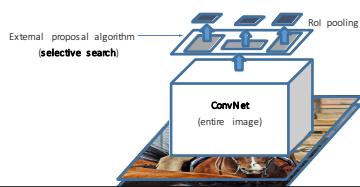
Review of the fast R-CNN training pipeline



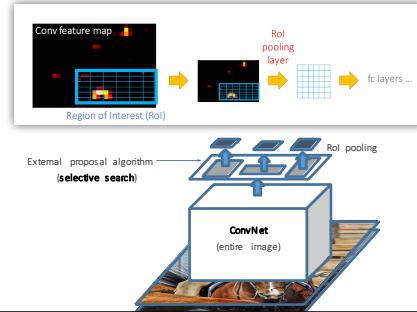
Review of the fast R-CNN training pipeline



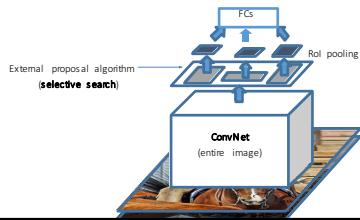
Review of the fast R-CNN training pipeline



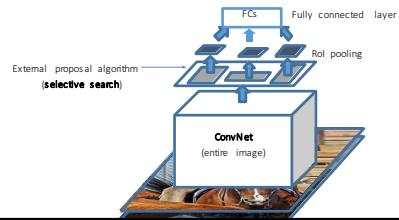
Review of the fast R-CNN training pipeline



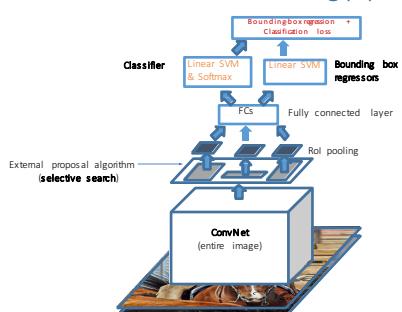
Review of the fast R-CNN training pipeline



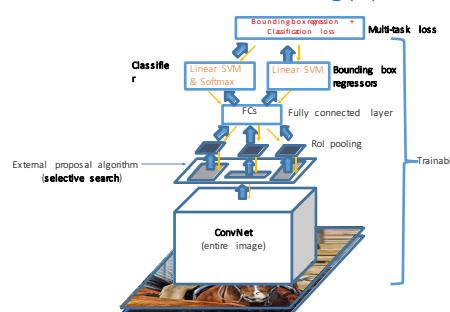
Review of the fast R-CNN training pipeline



Review of the fast R-CNN training pipeline



Review of the fast R-CNN training pipeline



Benefits of end-to-end training

- Faster training
 - No reading/writing features from/to disk
 - No training post hoc SVMs and bounding-box regressors
- Verified empirically: optimizing a single **multi-task objective** is **more accurate** than optimizing objectives independently

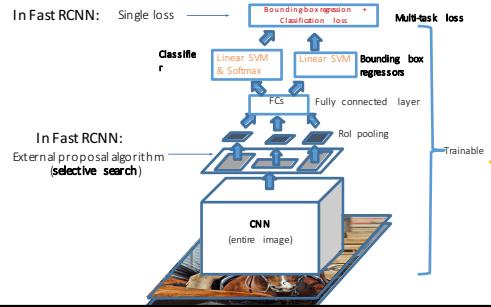
Slow R-CNN vs. Fast R-CNN

- Training time: **84 hours** / **8.75 hours**
- VOC07 test mAP: **66.0%** / **68.1%**
- Testing time per image: **47s** / **0.32s**
 - With selective search: **49s** / **2.32s** (+2s per image)

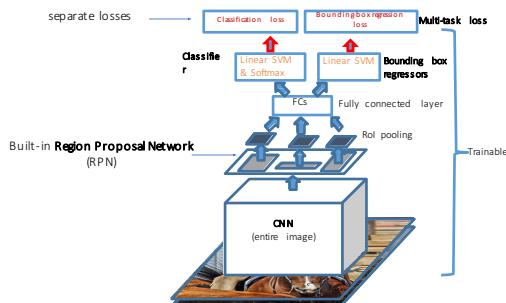
Topics of the lecture:

- ✓ Problem statement
- ✓ Review of slow R-CNN
- ✓ Review of Fast R-CNN
- ➡ **Review of Faster R-CNN**
- Compare with other methods
- Take away

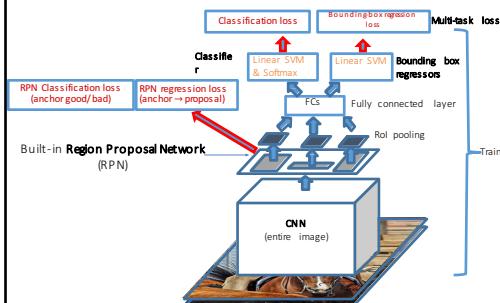
Review of the faster R-CNN



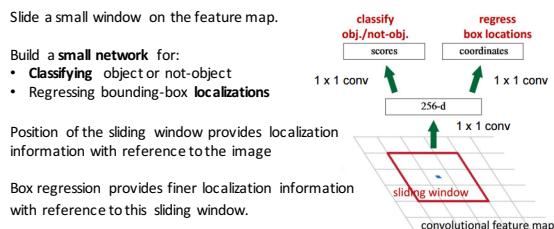
Review of the faster R-CNN



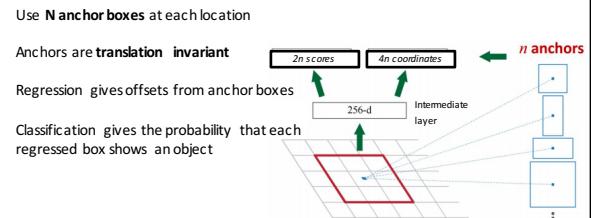
Review of the faster R-CNN



Faster R-CNN: Region Proposal Network



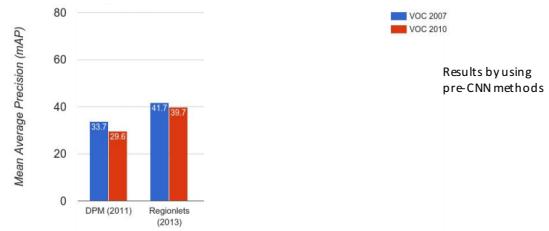
Faster R-CNN: Region Proposal Network



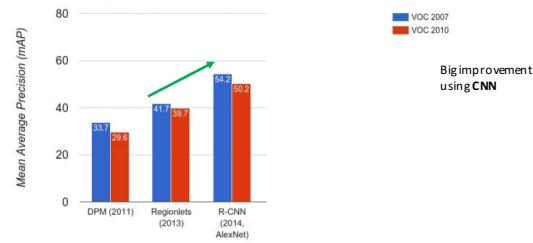
Topics of the lecture:

- ✓ Problem statement
- ✓ Review of slow R-CNN
- ✓ Review of Fast R-CNN
- ✓ Review of Faster R-CNN
- ➡ Compare with other methods
- Take away

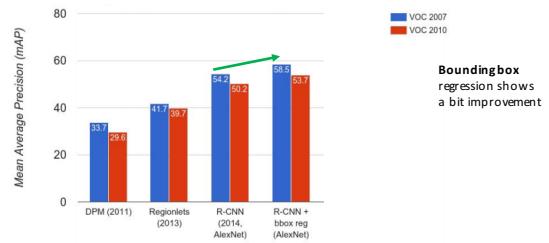
R-CNN Results



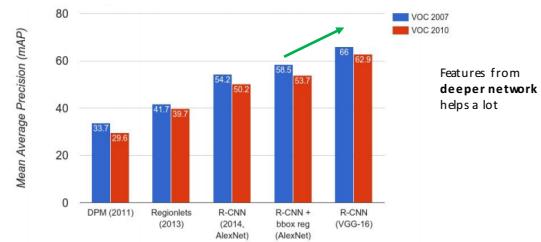
R-CNN Results



R-CNN Results



R-CNN Results



Fast R-CNN Results

	R-CNN	Fast R-CNN
Training Time:	84 hours	9.5 hours
(Speedup)	1x	8.8x
Test time per image	47 seconds	0.32 seconds
(Speedup)	1x	146x
mAP (VOC 2007)	66.0	66.9
Test time per Image with Selective Search	50 seconds	2 seconds
(Speedup)	1x	25x

Faster R-CNN Results

	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image (with proposals)	50 seconds	2 seconds	0.2 seconds
(Speedup)	1x	25x	250x
mAP (VOC 2007)	66.0	66.9	66.9

Faster R-CNN Results

YOLO – You Only Look Once
Detection as Regression

	Real-Time Detectors	Train	mAP	FPS
100Hz DPM [30]	2007	16.0	100	
30Hz DPM [30]	2007	26.1	30	
Fast YOLO	2007+2012	52.7	155	
YOLO	2007+2012	63.4	45	

Faster than Faster R-CNN, but not as good

	Less Than Real-Time	Train	mAP	FPS
Fastest DPM [37]	2007	30.4	15	
R-CNN Minus R [20]	2007	53.5	6	
Fast R-CNN [14]	2007+2012	70.0	0.5	
Faster R-CNN VGG-16[27]	2007+2012	73.2	7	
Faster R-CNN ZF [27]	2007+2012	62.1	18	

Faster R-CNN Results

Selective Search vs. Region Proposal Network

Table 7: Results on PASCAL VOC 2012 test set with Fast R-CNN detectors and VGG-16. For RPN, the train-time proposals for Fast R-CNN are 2000.

method	# box	data	laptop	motorcycle	bicycle	bottle	bus	car	motorcycle	chair	cow	table	dog	horse	motorcycle	person	plant	sheep	sofa	train	tv		
SS	2000	07++12	65.7	80.3	74.7	65.9	44.9	37.7	71.9	68.0	47.7	71.1	51.1	86.0	77.6	69.8	62.1	65.5	63.0	76.4	61.9		
SS	2000		68.4	82.3	79.4	70.8	52.3	38.7	77.8	71.6	69.3	44.2	73.0	55.0	87.5	80.5	80.8	72.0	55.1	68.3	65.7	80.4	64.2
RPN	300	12	67.0	82.3	76.4	71.0	48.4	45.2	72.1	72.3	87.3	42.2	73.7	50.0	86.8	78.7	76.4	77.4	34.5	70.1	57.1	77.1	58.9
RPN	300	07++12	70.4	84.9	79.8	74.3	53.9	49.8	77.5	75.9	88.5	45.6	77.1	55.3	86.9	81.7	80.9	79.6	40.1	72.6	60.9	81.2	61.5
RPN	300	COCO++07++12	75.5	87.4	85.6	82.8	59.6	51.9	81.0	81.3	84.9	52.6	80.0	50.0	85.5	84.7	84.1	53.2	78.9	65.5	84.4	70.2	

Topics of the lecture:

- ✓ Problem statement
- ✓ Review of slow R-CNN
- ✓ Review of Fast R-CNN
- ✓ Review of Faster R-CNN
- ✓ Compare with other methods
- ➡ Take away

Take Away

- Classification and detection of objects in images is approaching **high quality** and **real-time frame rate** results.
- The **major breakthroughs** came from:
 1. Utilizing **CNN** instead of classical computer vision methods
 2. Making **deeper** neural networks with **multiple** objectives
 3. Switching all stages of algorithm with **one unified** network (end-to-end)
- These algorithms are all **available** and **applicable** to real world problems! Once you know how they work you can change and adjust them to your specific research needs.

Object Detection code links:

R-CNN (Caffe+MATLAB):
<https://github.com/rbgirshick/rcnn>
[Probably don't use this; too slow]

Fast R-CNN (Caffe+MATLAB):
https://github.com/rbgirshick/fast_rcnn

Faster R-CNN (Caffe+MATLAB):
https://github.com/ShaoqingRen/faster_rcnn

Faster R-CNN (Caffe+Python):
https://github.com/bgirshick/pyfaster_rcnn

YOLO:
<http://pjreddie.com/darknet/yolo/>

Any Questions?



T·H·A·N·K·
YOU