

Object Detection

Nov. 26-27

What is object detection

Classification



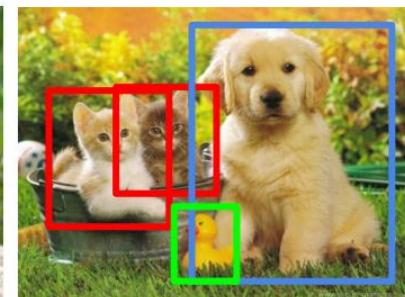
CAT

Classification + Localization



CAT

Object Detection



CAT, DOG, DUCK

Instance Segmentation



CAT, DOG, DUCK

Single object

Multiple objects

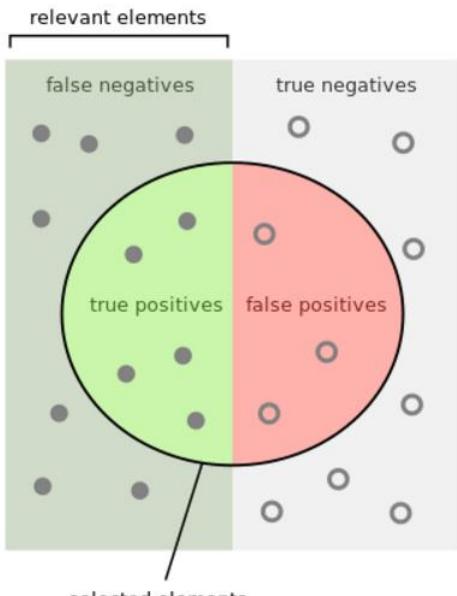
Terms

Recall

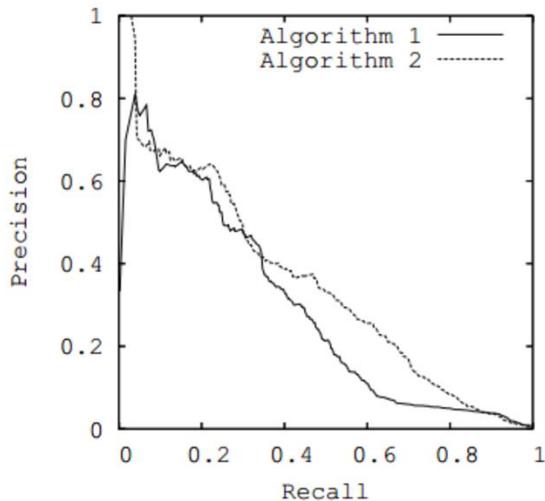
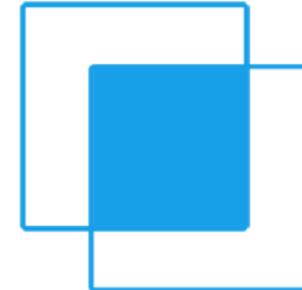
Precision

mAP

IoU



$$\text{Precision} = \frac{\text{How many selected items are relevant?}}{\text{How many selected items are selected?}}$$
$$\text{Recall} = \frac{\text{How many relevant items are selected?}}{\text{How many relevant items are there?}}$$



Detection Competitions

Pascal VOC

COCO

ImageNet ILSVRC

VOC: 20 classes



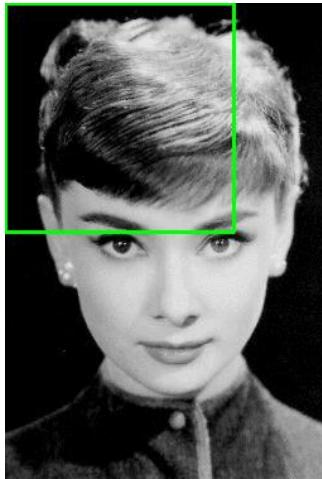
COCO: 200 classes



Before Deep Learning

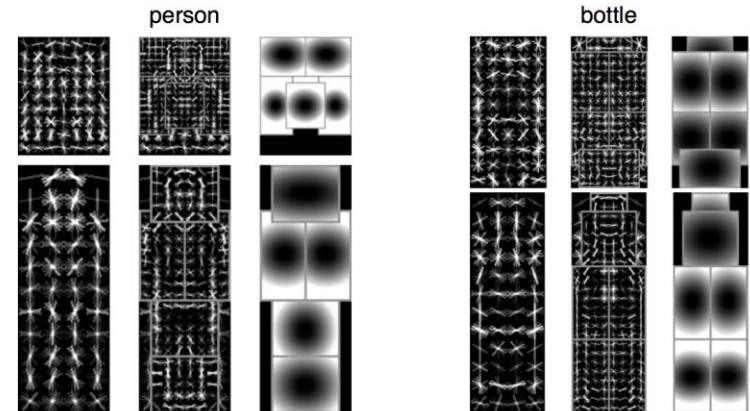
Sliding windows.

- Score every subwindow.



Deformable part models (DPM)

- Uses HOG features
- Very fast



Selective Search



Hard Negative Mining

Imbalance between positive and negative examples.

Use negative examples with higher confidence score.

Non Maximum Suppression

If output boxes overlap, only consider the most confident.

Bounding Box Regression

Regression used to find bounding box parameters

Applied in one of two ways

- Bounding box refinement
- Complete object detection

Example Loss Function

$$\sum_{i \text{ in } I} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 + (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2]$$

I is the set of all matching bounding boxes (Highest IoU with ground truth)

Two Stage Detection

Part I

RCNN : Region Proposal + CNN

Use selective search to come up with regional proposal

First object detection method using CNN

Rich feature hierarchies for accurate object detection and semantic segmentation

Ross Girshick Jeff Donahue Trevor Darrell Jitendra Malik

Nov 2013

<https://people.eecs.berkeley.edu/~rbg/papers/r-cnn-cvpr.pdf>

RCNN :

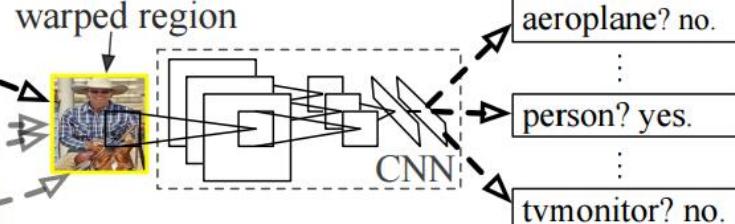
R-CNN: *Regions with CNN features*



1. Input image



2. Extract region proposals (~2k)

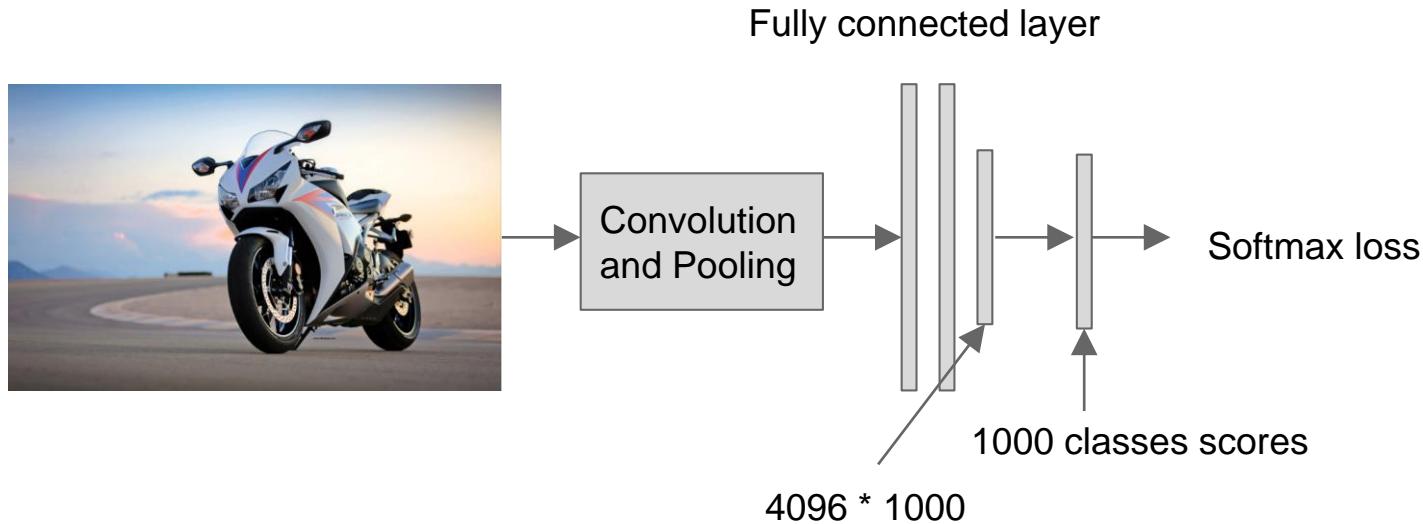


3. Compute CNN features

4. Classify regions

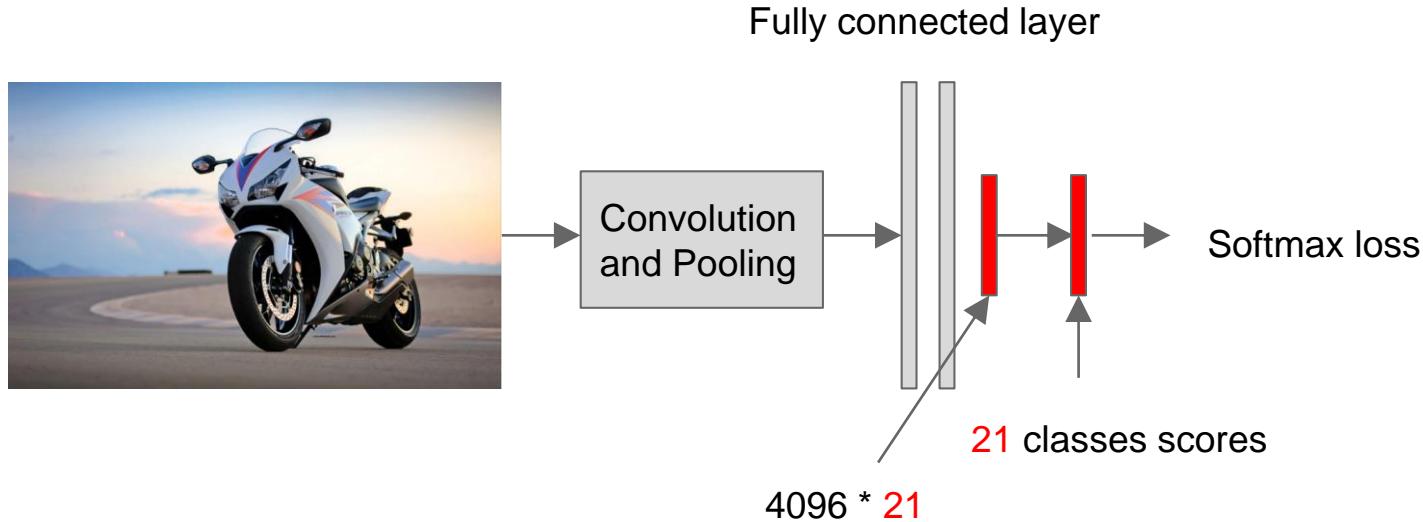
Training RCNN

Step1: train your own CNN model for classification (or use existing model), using ImageNet dataset.



Training RCNN

Step2: focus on 20 classes + 1 background. Remove the last FC layer and replace it with a smaller layer and fine-tune the model using PASCAL VOC dataset



Training RCNN

Step3: extract feature



Crop & Warp



Convolution
and Pooling

Store all the features
after pool 5 layer
and save to disk

It is about ~ 200G
features

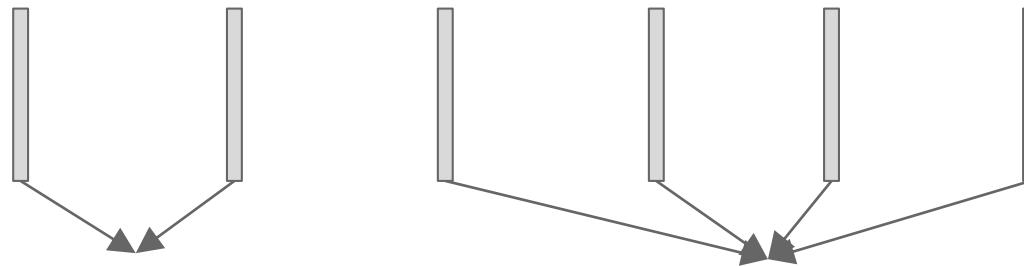
Training RCNN

Step4: train SVM for each class

Crop /
Warp
image



Features
from last
step



Positive
samples for
Motorbike

Negative
samples for
Motorbike

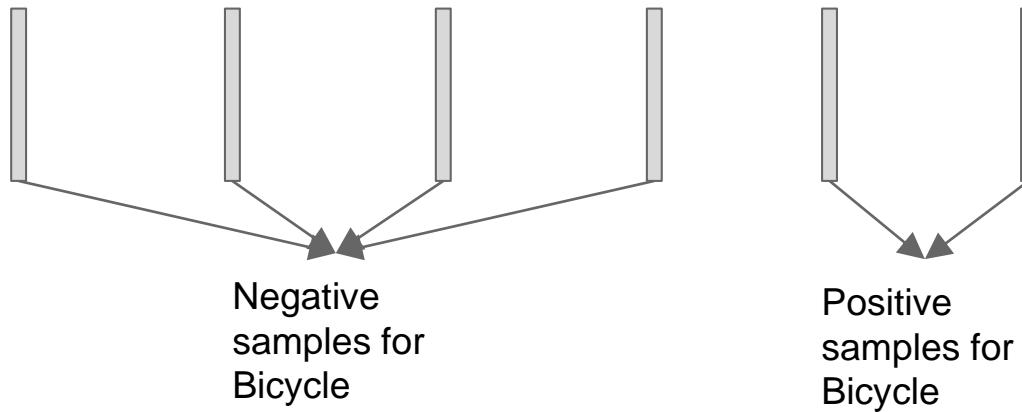
Training RCNN

Step4: train SVM for each class

Crop /
Warp
image



Features
from last
step



Fast RCNN

Share convolution layers for proposals from the same image

Faster and More accurate than RCNN

ROI Pooling

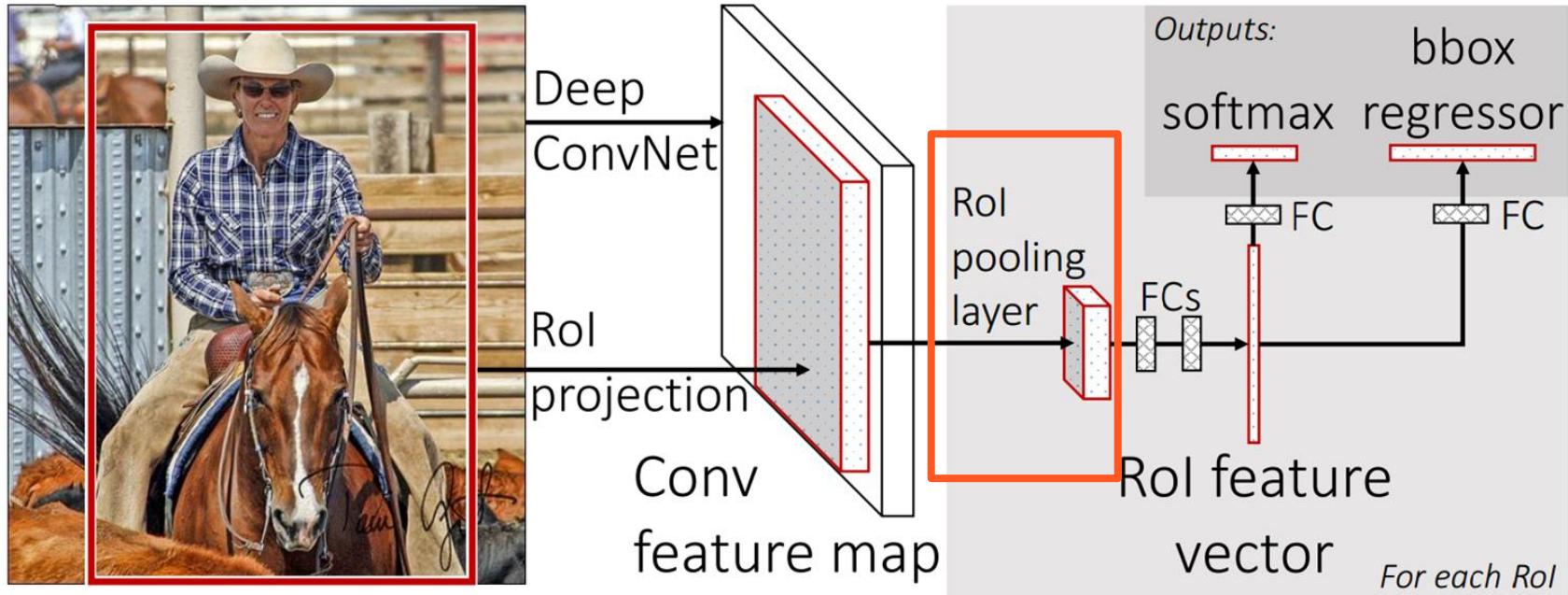
Fast R-CNN

Ross Girshick

Apr 2015

<https://arxiv.org/pdf/1504.08083v2.pdf>

Fast RCNN



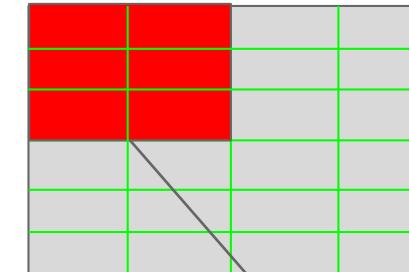
RoI Pooling



image

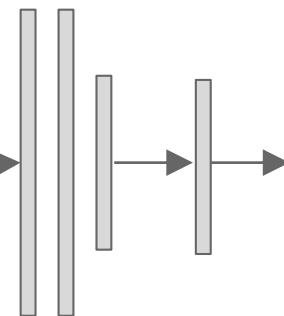
Conv layer

Divided into
 $h * w$ region

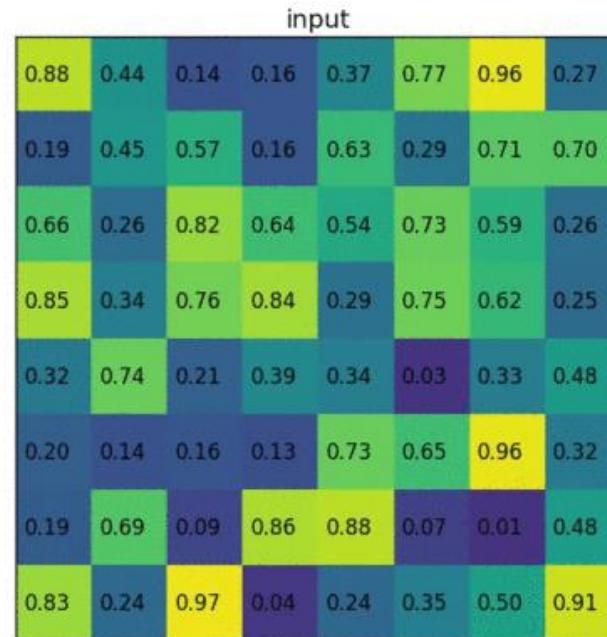


Max pooling

Differentiable



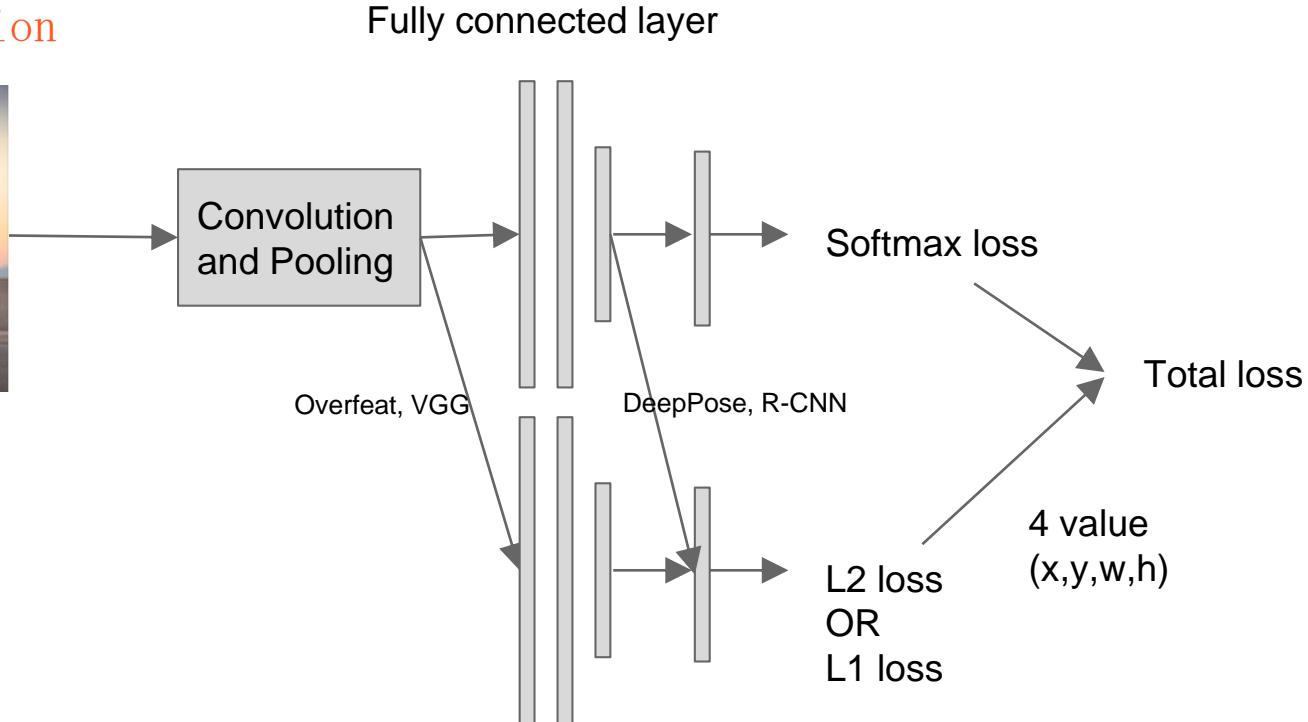
ROI Pooling

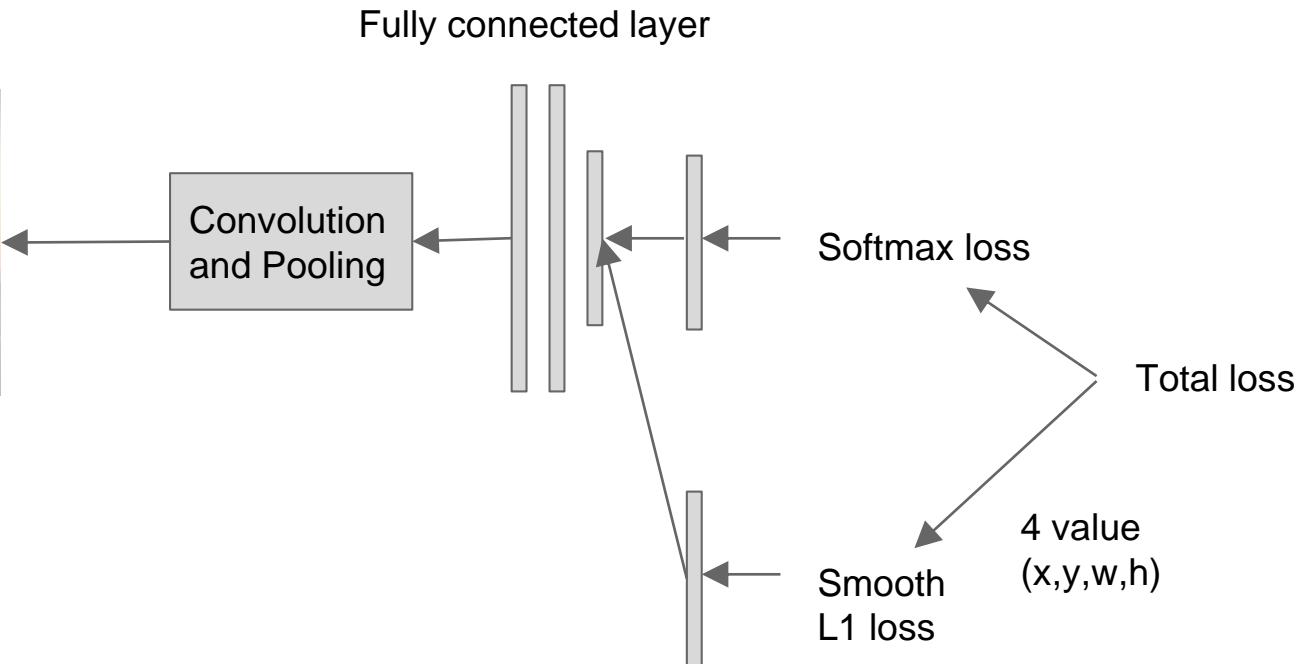


<https://blog.deepsense.ai/region-of-interest-pooling-explained/>

What is bbox-regressor?

Bounding box regression





Result compare

	Fast R-CNN	R-CNN
Train time (h)	9.5	84
-speedup	8.8x	1x
Test time/image	0.32s	47.00 s
-test speedup	146x	1x
mAP	66.9	66

Trained using VGG 16 on Pascal VOC 2007 dataset
Not including proposal time

Result compare

	Fast R-CNN	R-CNN
Test time/image	0.32s	47.00 s
-test speedup	146x	1x
Test time/image with proposal	2s	50 s
-test speedup	25x	1x

Faster RCNN

Don't need to have external regional proposals

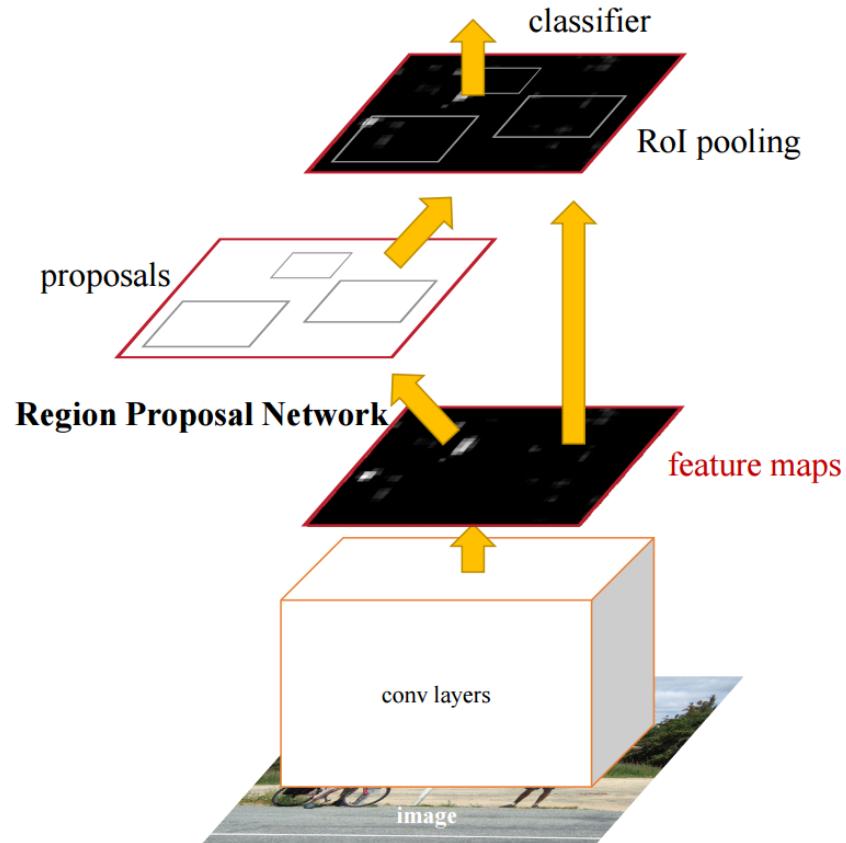
RPN - Regional Proposal Network

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

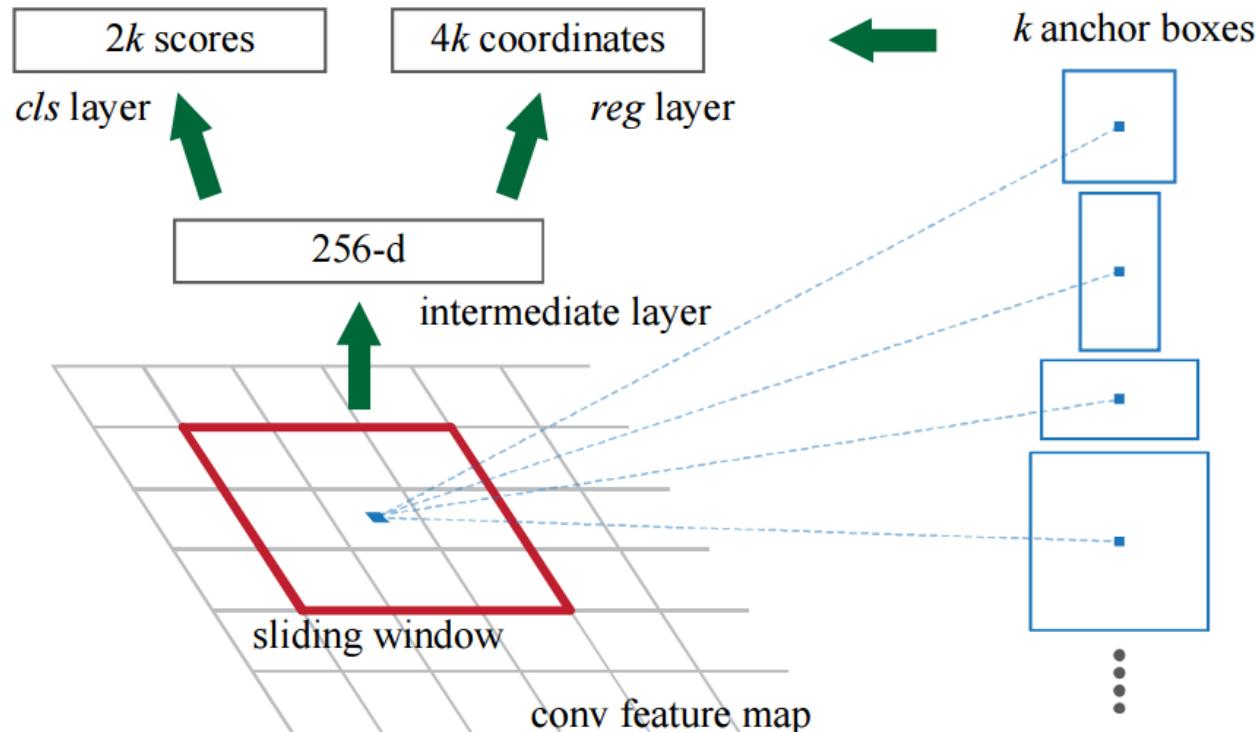
Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun

Jun 2015

Faster RCNN



Faster RCNN



Result compare

	Faster R-CNN	Fast R-CNN	R-CNN
Test time/image With proposal	0.2S	2s	50s
-test speedup	250x	25x	1x
mAP	66.9	66.9	66

Trained using Pascal VOC 2007 dataset

R-FCN :Region-based Fully Convolutional Networks

Use position sensitive score map

Share all conv and fc layers between all proposals for the same image

R-FCN: Object Detection via Region-based Fully Convolutional Networks

Jifeng Dai, Yi Li, Kaiming He, Jian Sun

May 2016

<https://arxiv.org/pdf/1605.06409v2.pdf>

R-FCN

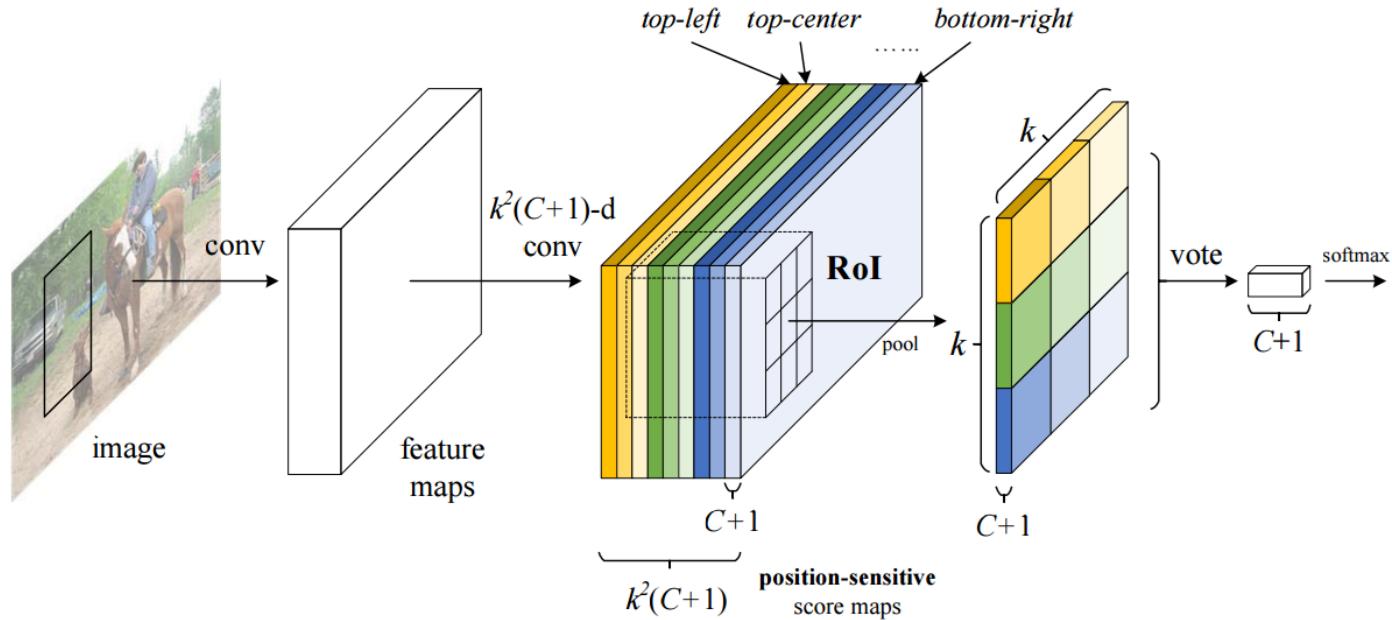


Figure 1: Key idea of **R-FCN** for object detection. In this illustration, there are $k \times k = 3 \times 3$ position-sensitive score maps generated by a fully convolutional network. For each of the $k \times k$ bins in an RoI, pooling is only performed on one of the k^2 maps (marked by different colors).

R-FCN

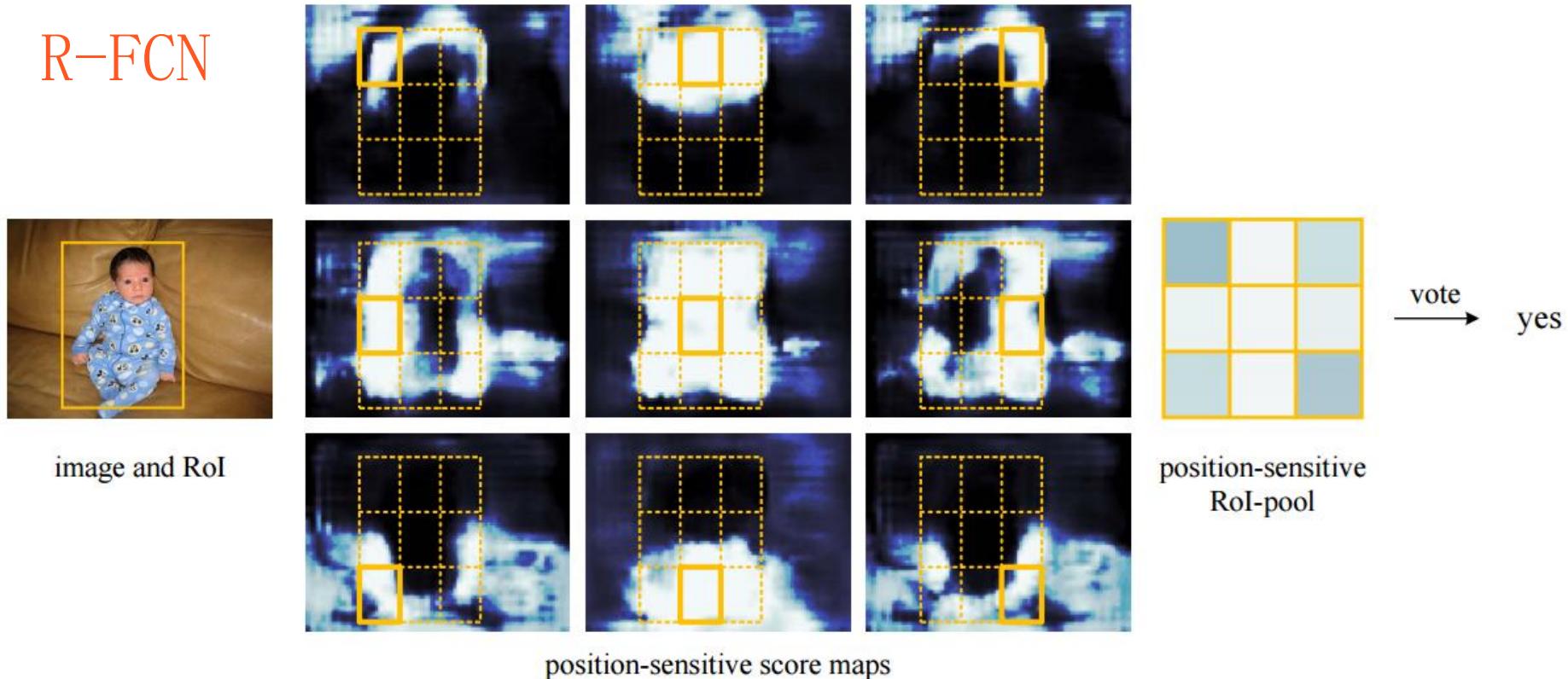


Figure 3: Visualization of R-FCN ($k \times k = 3 \times 3$) for the *person* category.

R-FCN

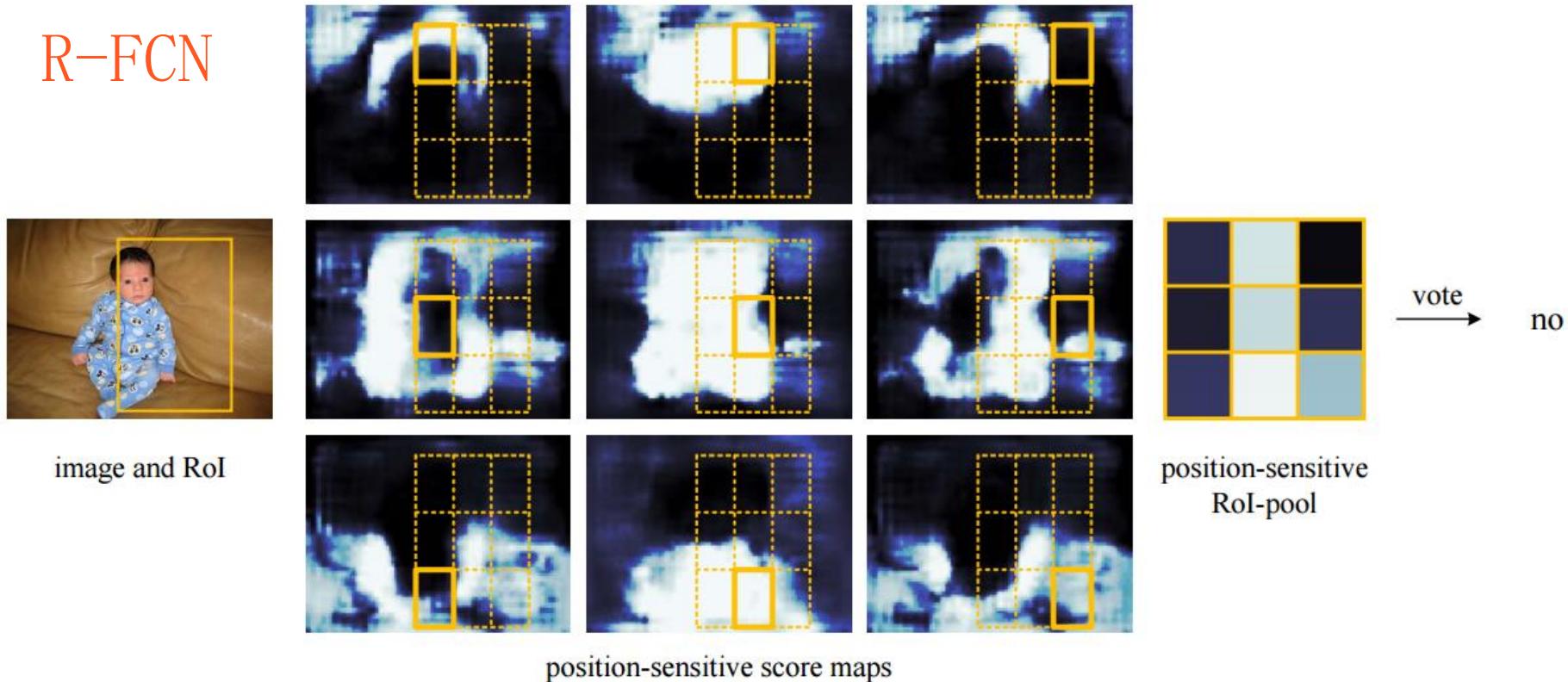


Figure 4: Visualization when an ROI does not correctly overlap the object.

Methodologies Compare

	RCNN	Faster RCNN	RFCN
Depth of shared convolutional subnetwork	0	91	101
Depth of ROI-wise subnetwork	101	10	0

Trained using ResNet 101

Result compare

	Faster R-CNN	R-FCN
Test time/image With proposal	0.42S	0.2s
mAP	76.4	76.6

Trained using ResNet 101 on Pascal VOC 2007 dataset

Unified Detection

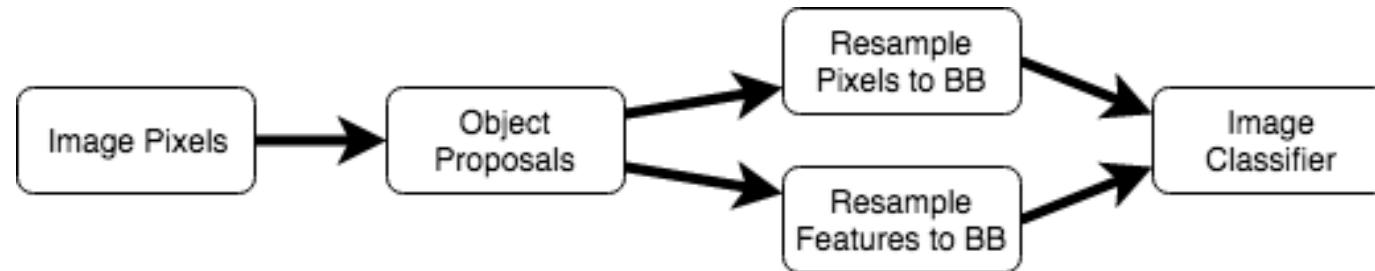
Part II

Problems with 2 step detection.

Complex Pipeline

Slow (Cannot run in real time)

Hard to optimize each component



Yolo: You Only Look Once

Consider detection a regression problem

Use a single ConvNet

Runs once on entire image. Very Fast!

You only look once: Unified, real-time object detection.

Joseph Redmon, Santosh Divvala, Ross Girshick, Ali Farhadi

June 2015

How it works

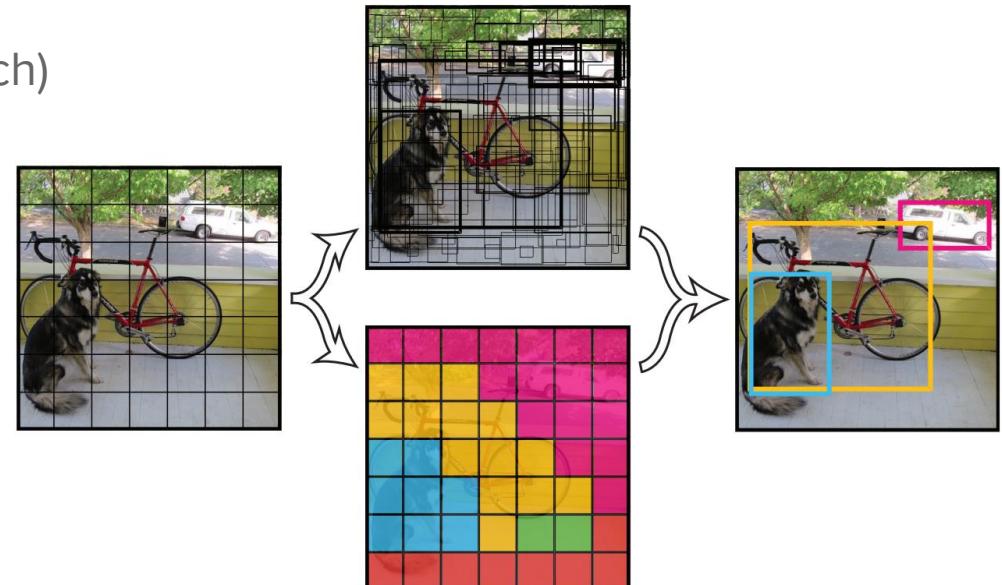
The following predictions are made for each cell in an $S \times S$ grid.

C conditional class probabilities $\text{Pr}(\text{Class}_i \mid \text{Obj})$

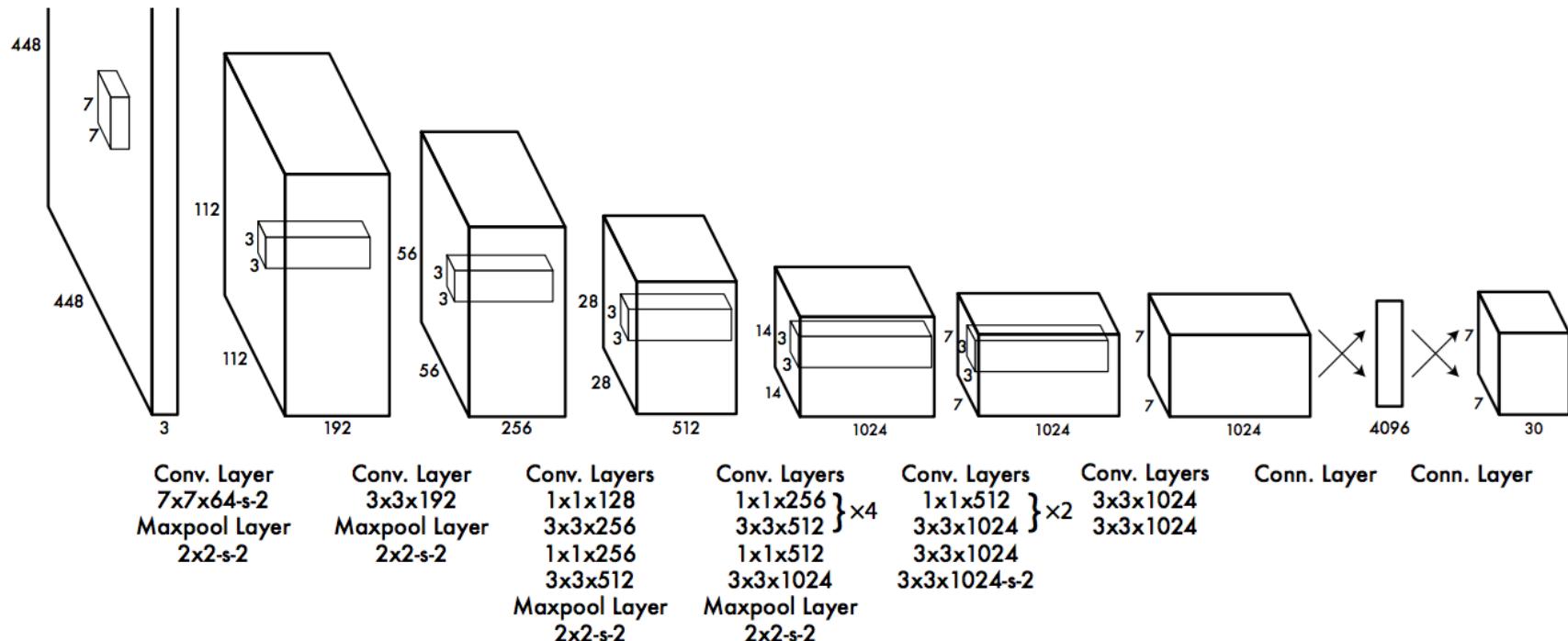
B bounding boxes (4 parameters each)

B confidence scores $\text{Pr}(\text{Obj}) * \text{IoU}$

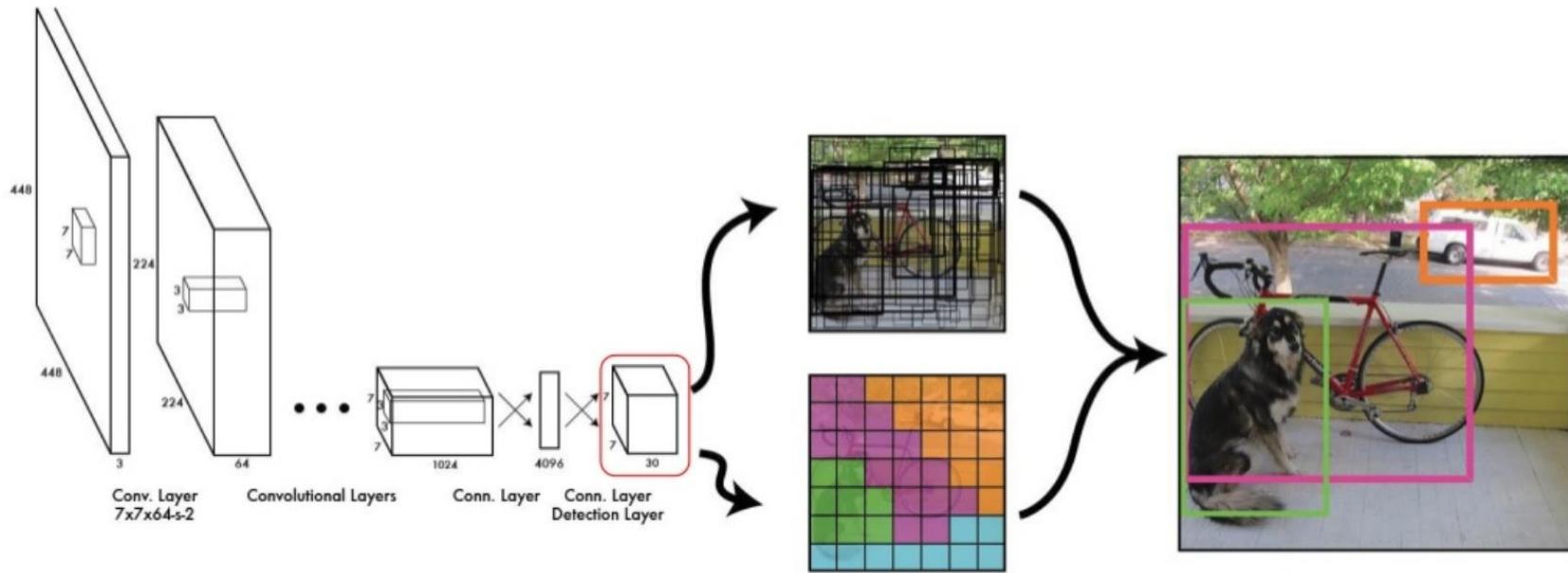
Output is $S \times S \times (5B+C)$ tensor



Architecture

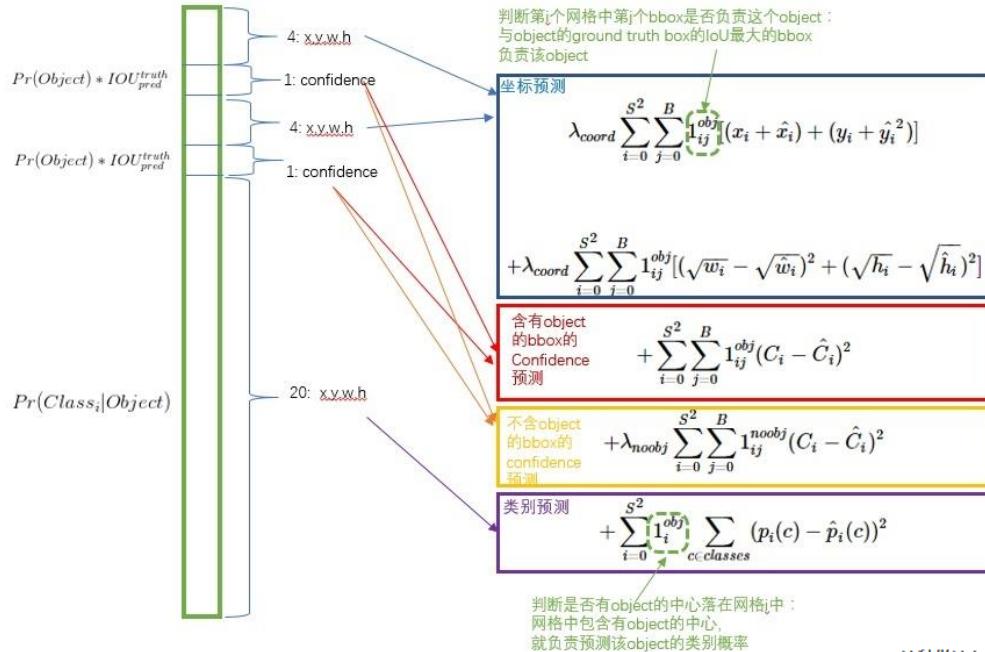


Architecture



$$\Pr(\text{Class}_i \mid \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

Architecture



这种做法存在以下几个问题：

- 第一，8维的localization error和20维的classification error同等重要显然是不合理的；
- 第二，如果一个网格中没有object（一幅图中这种网格很多），那么就会将这些网格中的box的confidence push到0，相比于较少的有object的网格，这种做法是overpowering的，这会导致网络不稳定甚至发散。

解决办法：

- 更重视8维的坐标预测，给这些损失前面赋予更大的loss weight，记为 λ_{coord} 在pascal VOC训练中取5。
- 对没有object的box的confidence loss，赋予小的loss weight，记为 λ_{noobj} 在pascal VOC训练中取0.5。
- 有object的box的confidence loss和类别的loss的loss weight正常取1。

<https://arxiv.org/pdf/1506.02640v5.pdf>

Performance (VOC 2007)

	Yolo	Faster R-CNN (VGG-16)
mAP	63.4	73.2
FPS	45	7

Trained on Pascal VOC 2007 + 2012 dataset

Limitations of Yolo

Struggles with small objects

Struggles with unusual aspect ratios

Poor localization

SSD Single Shot Detector

Faster than Yolo, as accurate as Faster R-CNN

Predicts categories and box offsets

Uses small convolutional filters applied to feature maps

Makes predictions using feature maps of different scales

SSD: Single shot multibox detector

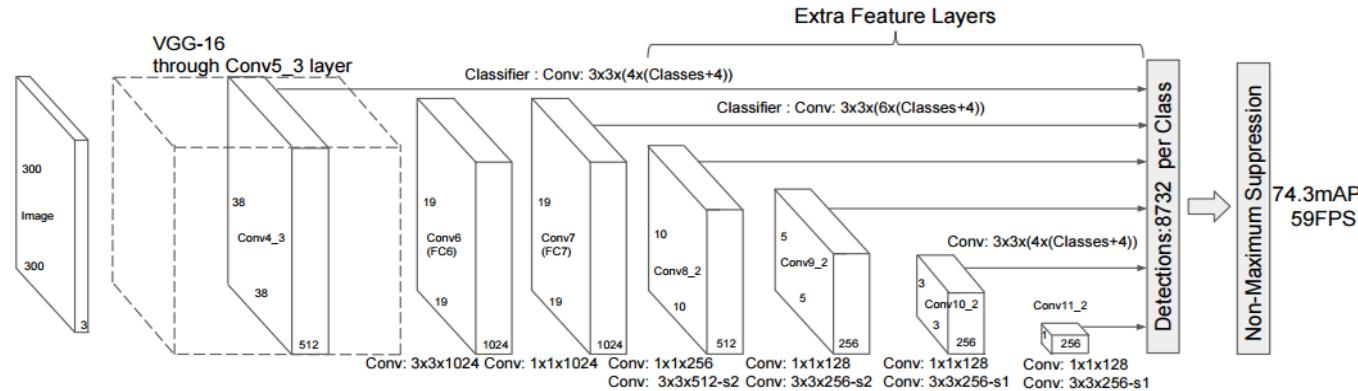
Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, Alexander C. Berg

Dec 2015

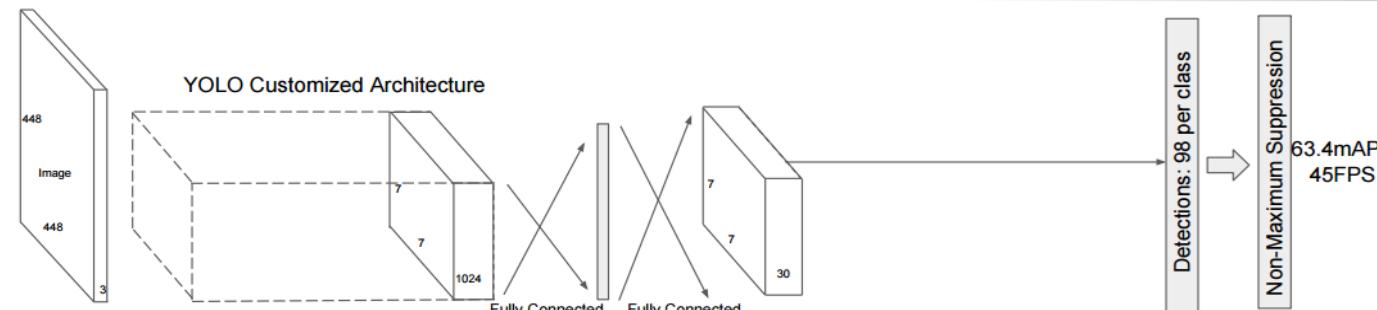
Comparison to Yolo

$$38*38*4+19*19*6+10*10*6+5*5*6+3*3*4+1*1*4=8732 \text{ default boxes}$$

SSD



Yolo

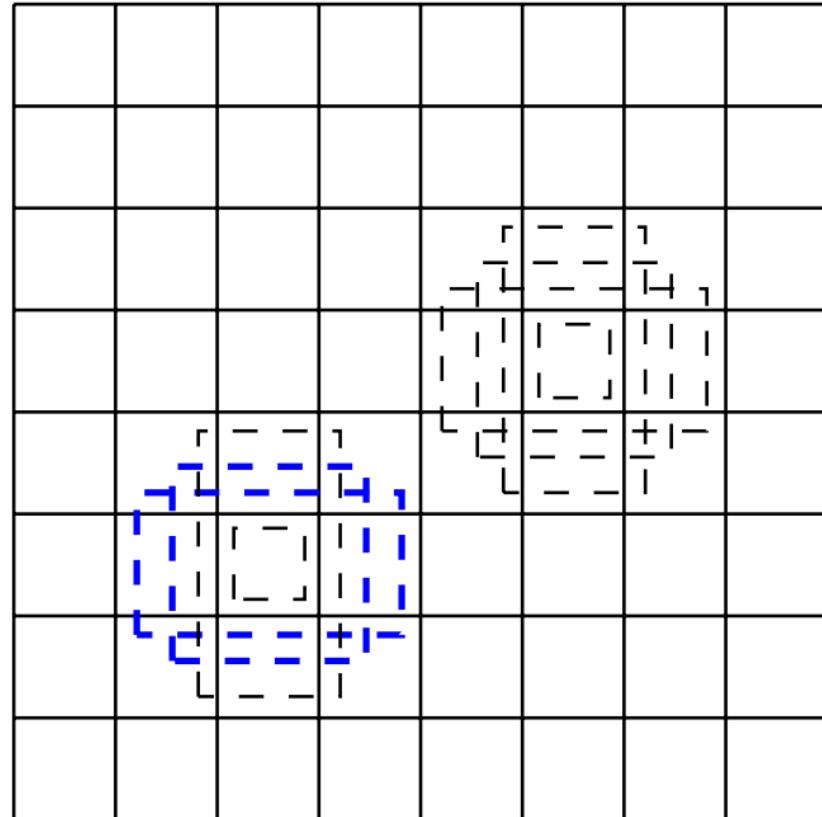


Default Boxes

Multiple aspect ratios per cell.

Similar to Faster R-CNN Anchor Boxes.

- Applied to many feature maps.



SSD Detector

Detectors are convolutional filters.

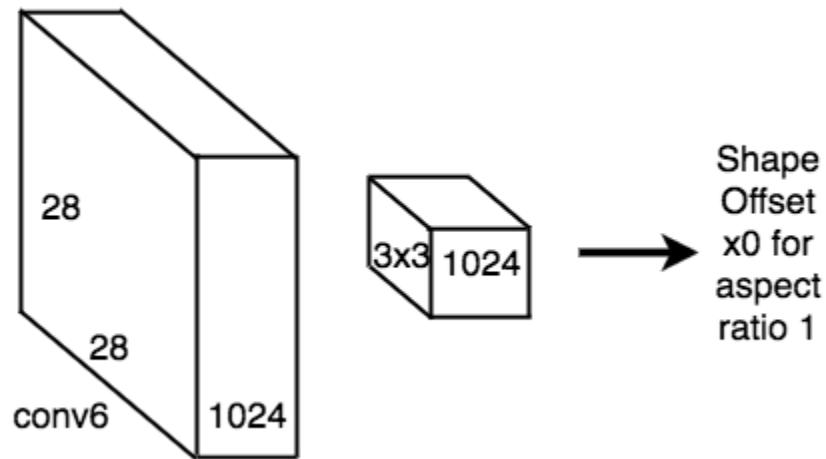
Each detector outputs a single value.

Predict class probabilities.

Predict bounding box offsets.

(classes + 4) detectors are needed for a detection.

(classes + 4) x (#default boxes) x m x n outputs for a mxn feature map.



Results (VOC 2007)

	SSD*	Yolo*	Faster R-CNN*
mAP	74.3	66.4	73.2
FPS	46	21	7

*VGG16

Trained on Pascal VOC 2007 + 2012 dataset

Others

Part III

Feature Pyramid Networks

Uses ConvNet as Feature Pyramid

Includes low level feature maps to detect small objects

Top down pathway provides contextual information

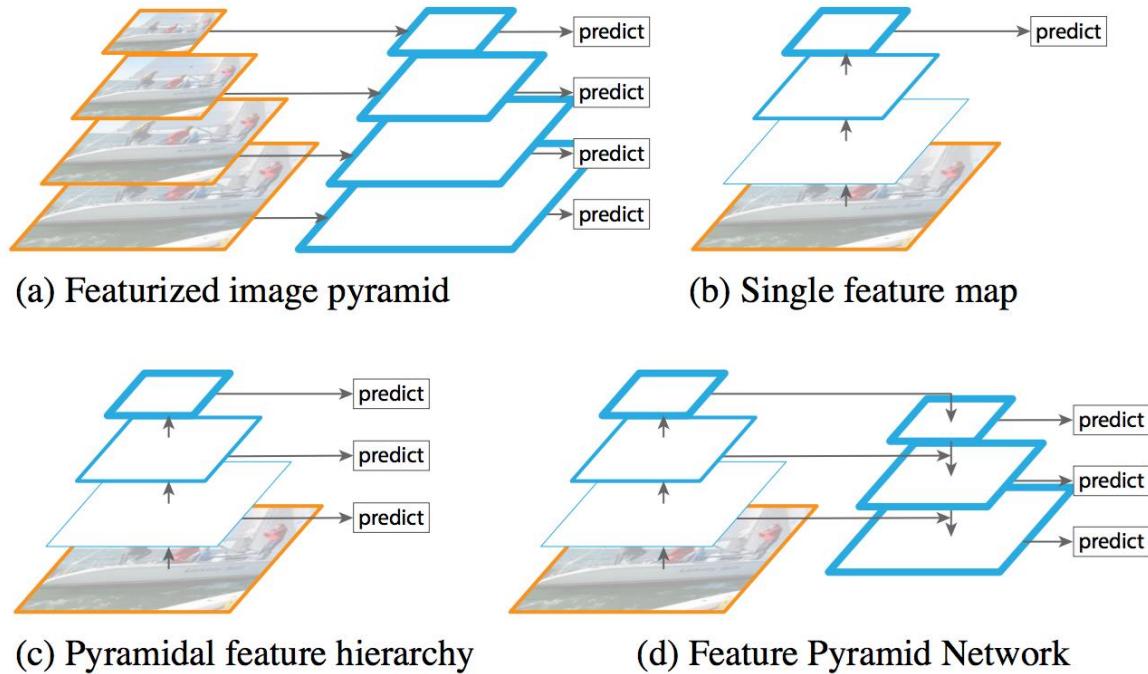
Feature Pyramid Networks for Object Detection.

[Tsung-Yi Lin](#), Piotr Dollár, [Ross Girshick](#), [Kaiming He](#), [Bharath Hariharan](#), Serge Belongie.

Dec 2016

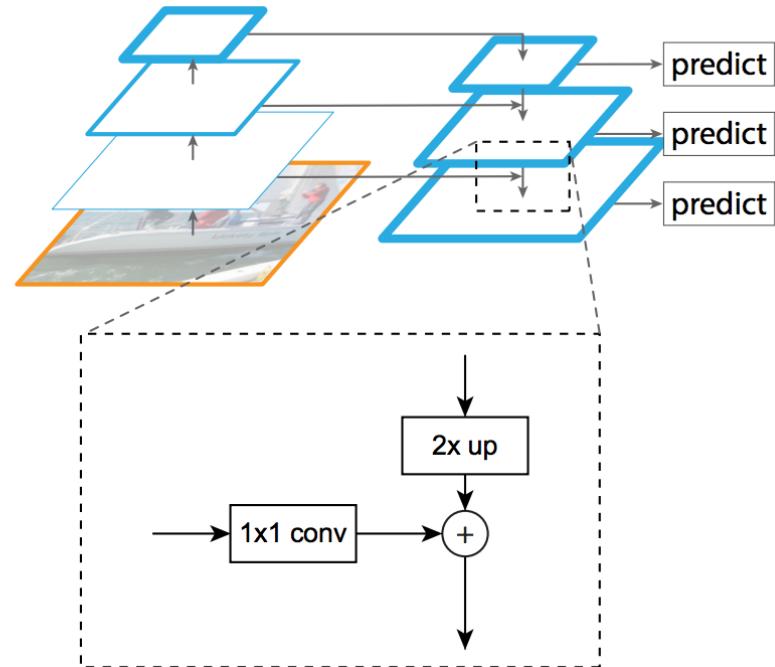
Comparison to prior methods

- a. Accurate but slow.
- b. Misses low level information. (Yolo)
- c. Misses context in low level predictions. (SSD)
- d. Accurate and fast.



Top down pathway

- High level (semantically strong) feature maps are upsampled.
- Lateral connections merge feature maps from bottom up pathway.



Feature Pyramid Networks for RPN

- Replace single scale feature map with FPN.
- Single scale anchors at each level.

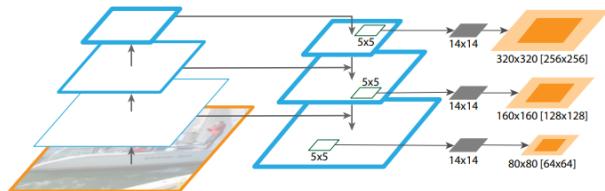
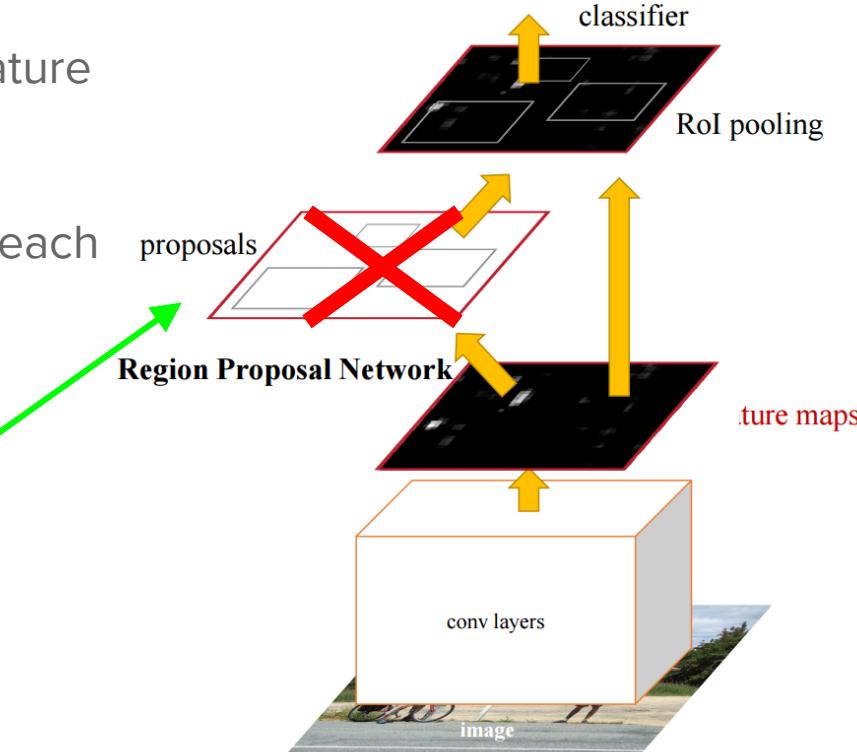


Figure 4. FPN for object segment proposals. The feature pyramid is constructed with identical structure as for object detection. We apply a small MLP on 5×5 windows to generate dense object segments with output dimension of 14×14 . Shown in orange are the size of the image regions the mask corresponds to for each pyramid level (levels P_{3-5} are shown here). Both the corresponding image region size (light orange) and canonical object size (dark orange) are shown. Half octaves are handled by an MLP on 7×7 windows ($7 \approx 5\sqrt{2}$), not shown here. Details are in the appendix.



Results

- COCO
- Previous State of the art single-model

	Faster R-CNN on FPN*	Faster R-CNN +++*	ION**
mAP	36.2	34.9	31.2
mAP (small images)	18.2	15.6	12.8

* ResNet-101

** VGG-16

ION : Inside–Outside Network

Use multi-scale Conv features for inside region

Use four direction RNN features for outside region

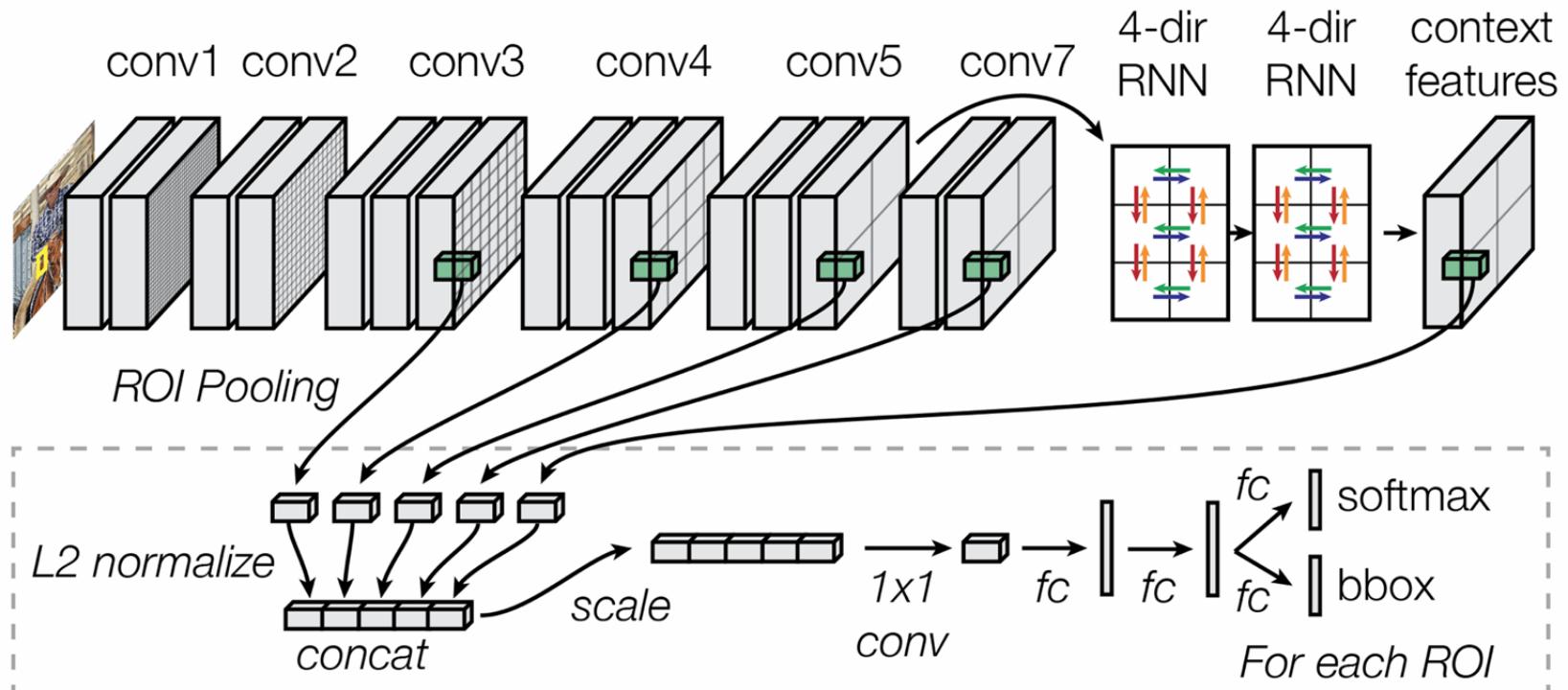
Inside-Outside Net: Detecting Objects in Context with Skip Pooling and Recurrent Neural Networks

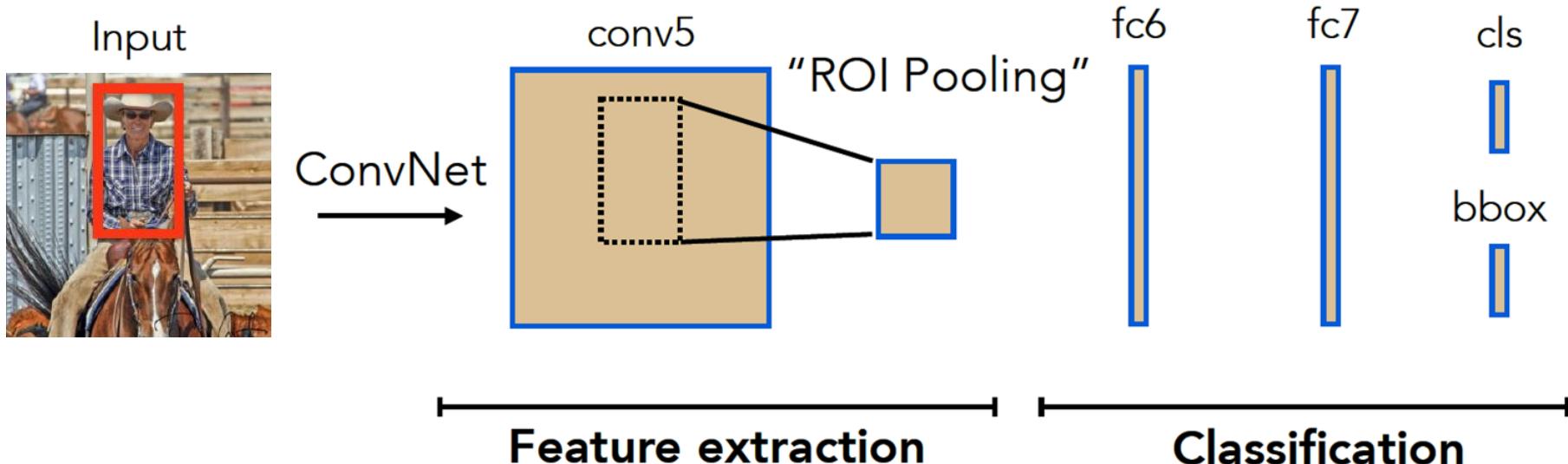
Sean Bell, C.Lawrence Zitnick, Kavita Bala, Ross Girshick

Cornell University, Microsoft Research

CVPR 2016

ION: Inside–Outside Net

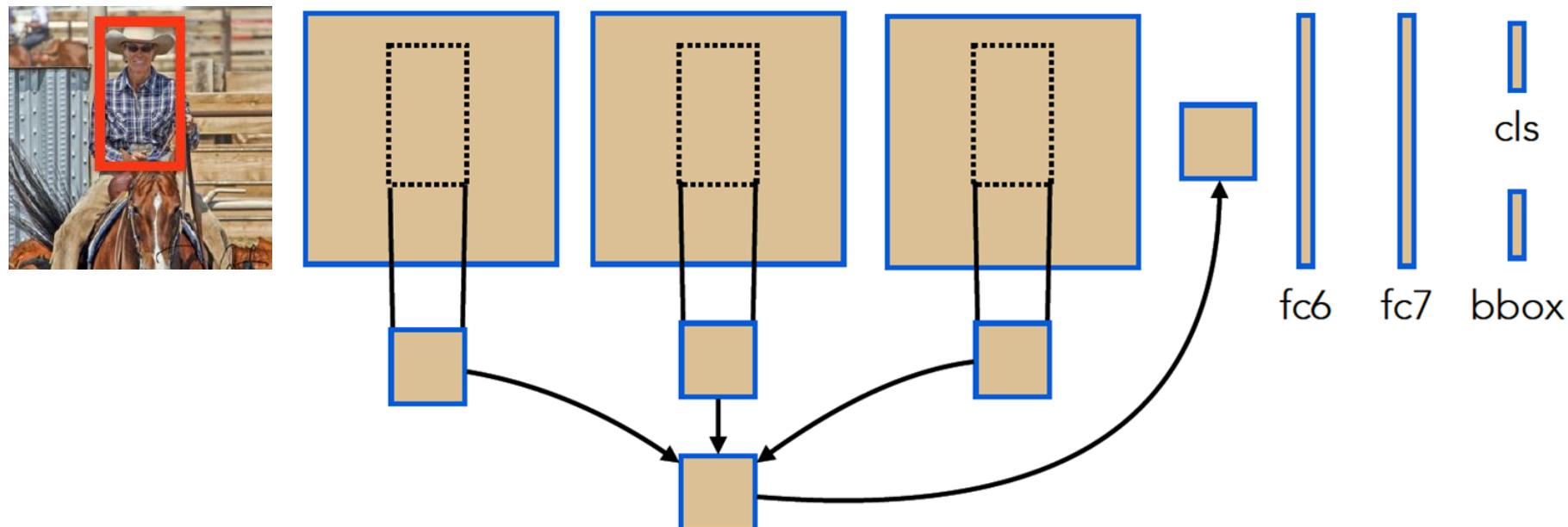




Can we improve on feature extraction?

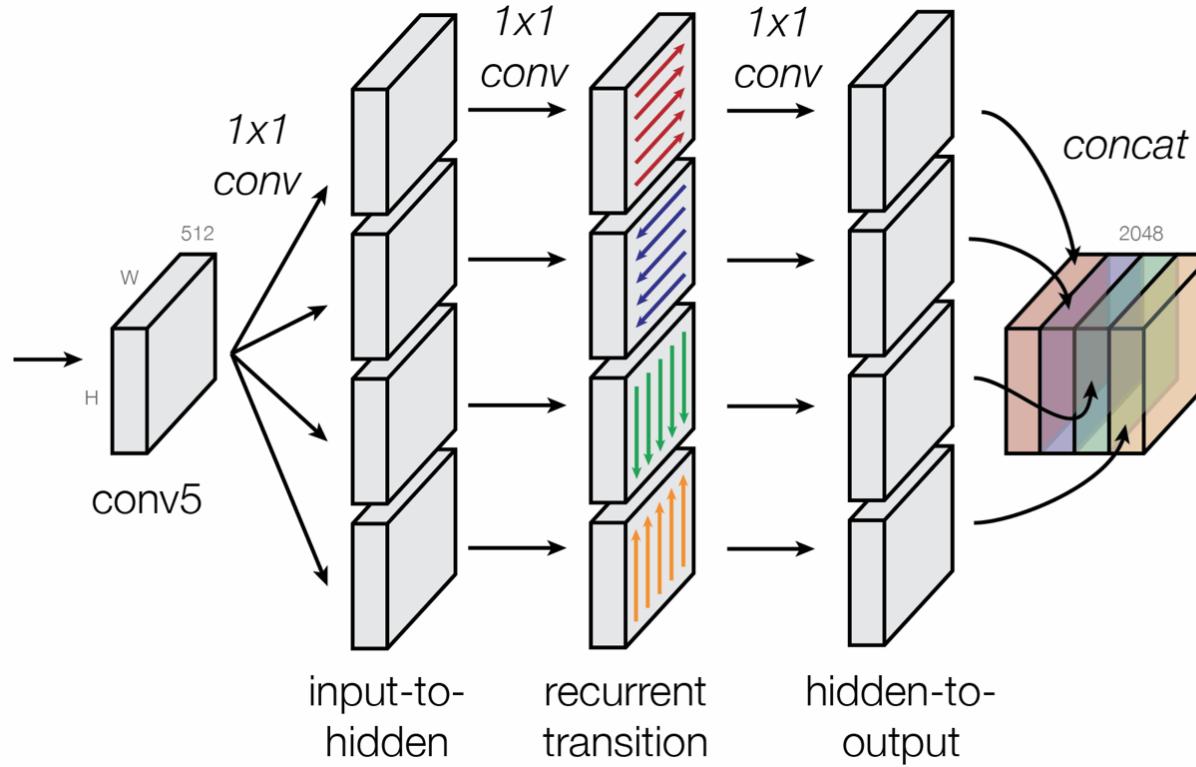
- For small objects, the footprint on conv5 might only cover a 1x1 cell, which gets upsampled to 7x7
- Only local features (inside the ROI) are used for classification

COMBINING ACROSS LAYERS

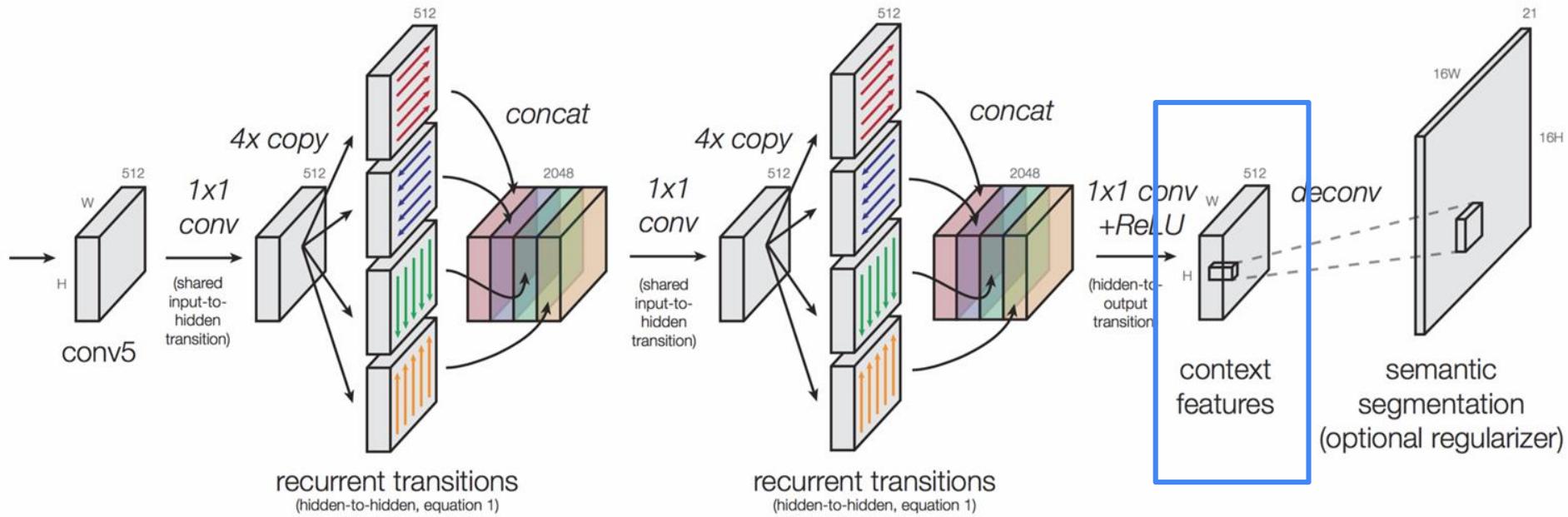


normalize, concatenate, re-scale

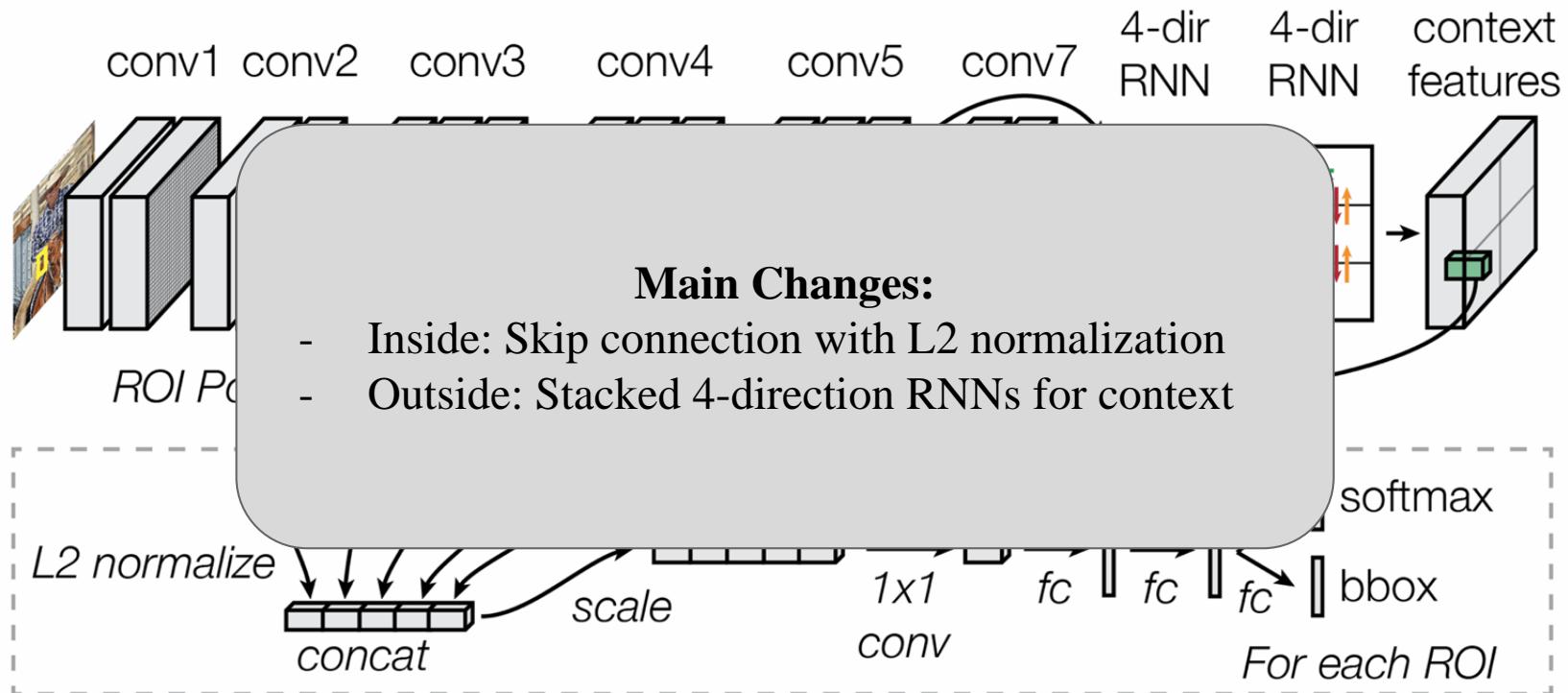
Uses RNNs to capture context info from outside bounding box.



Stack 2 RNNs together



ION: Inside–Outside Net



Results (VOC 2007)

METHOD	MAP
FAST R-CNN [GIRSHICK 2014]	70.0
FASTER R-CNN [GIRSHICK 2015]	73.2
CONV3+CONV4+CONV5	75.6
+ RNN + SEGMENTATION LOSS	76.5
+ SECOND BBOX REGRESSION + WEIGHTED VOTING	78.5
— DROPOUT	79.2

<https://www.robots.ox.ac.uk/~vgg/rg/slides/ion-coco.pdf>

Trained on Pascal VOC 2007 + 2012 dataset

Summary and Comparison

Part IV

Speed / Accuracy Trade-off

Unified tensorflow architecture

Compare speed, accuracy and memory usage

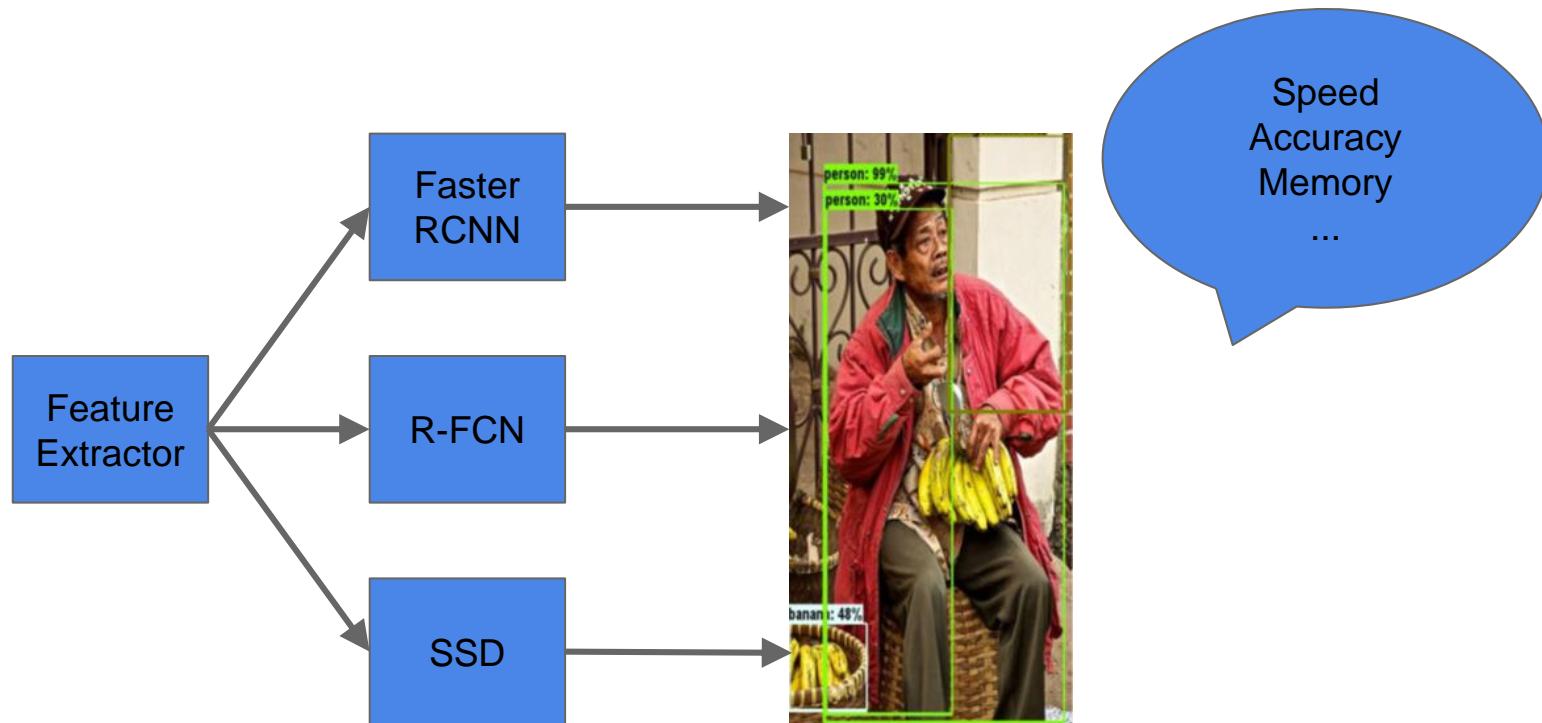
Speed/Accuracy trade-offs for modern convolutional object detectors

Jonathan Huang, Kevin Murphy et al.

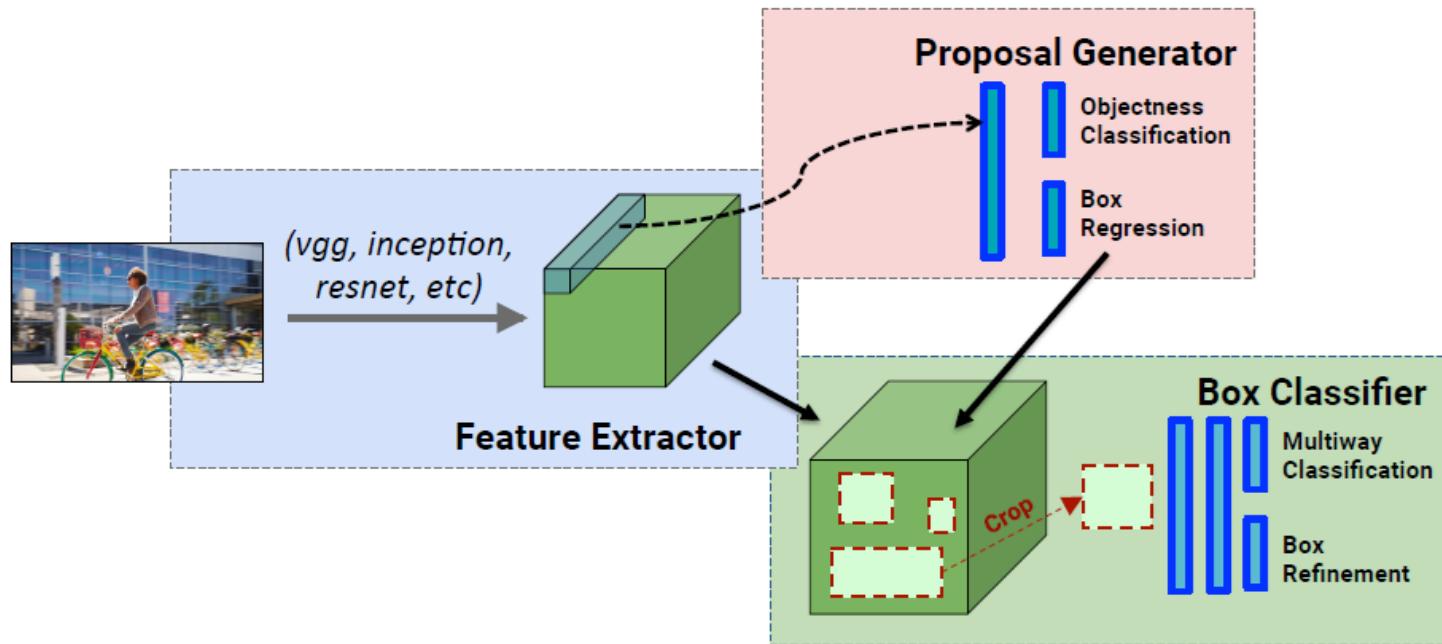
Nov 2016

<https://people.eecs.berkeley.edu/~rbg/papers/r-cnn-cvpr.pdf>

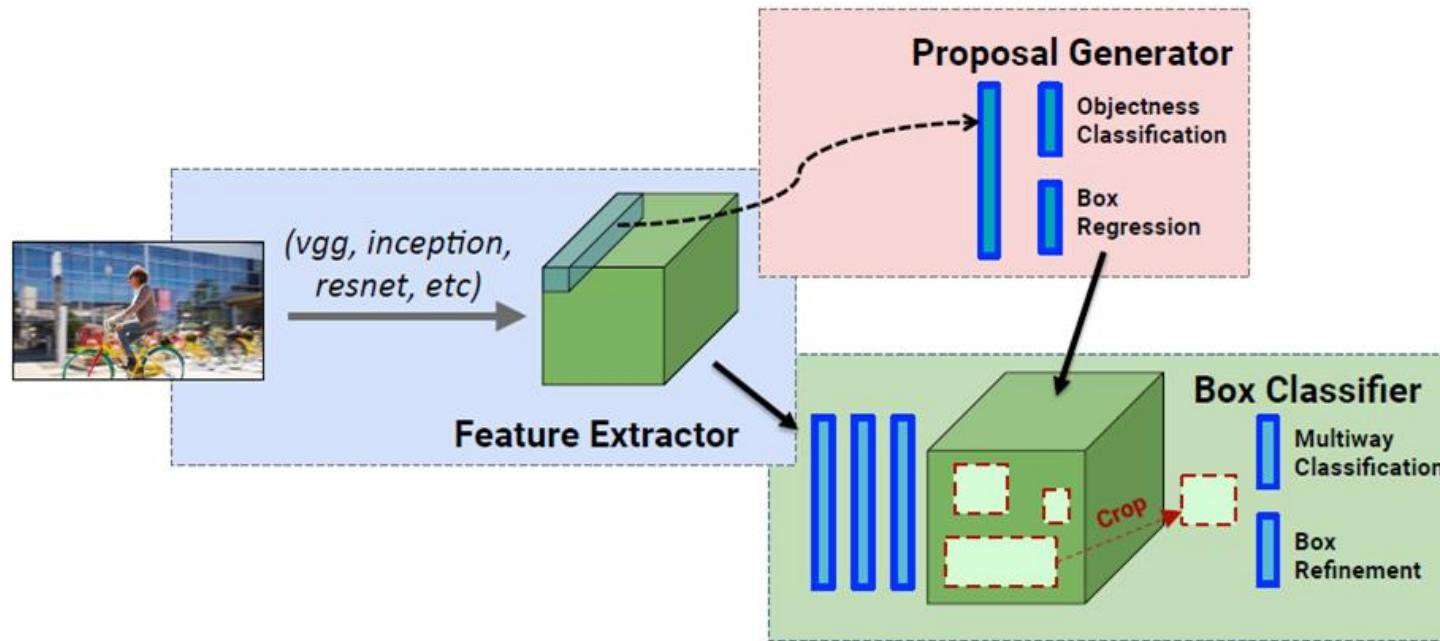
Same Architectures



Recap: Faster RCNN

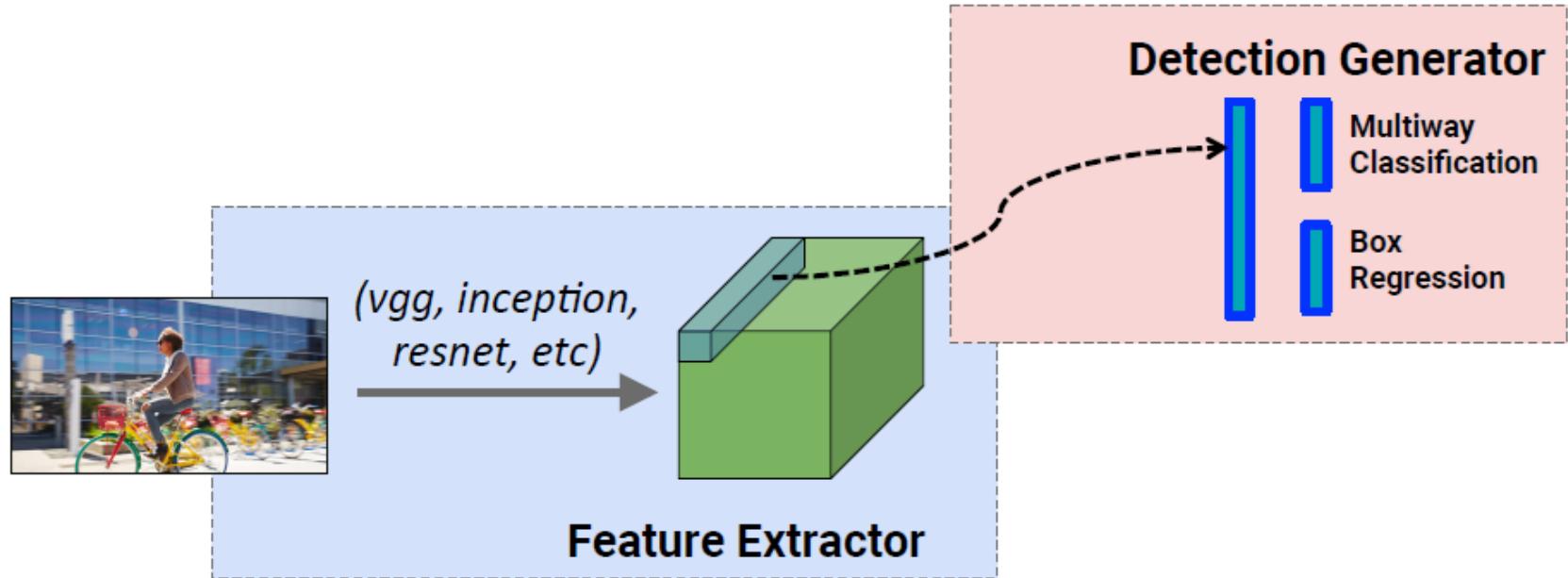


Recap: R-FCN

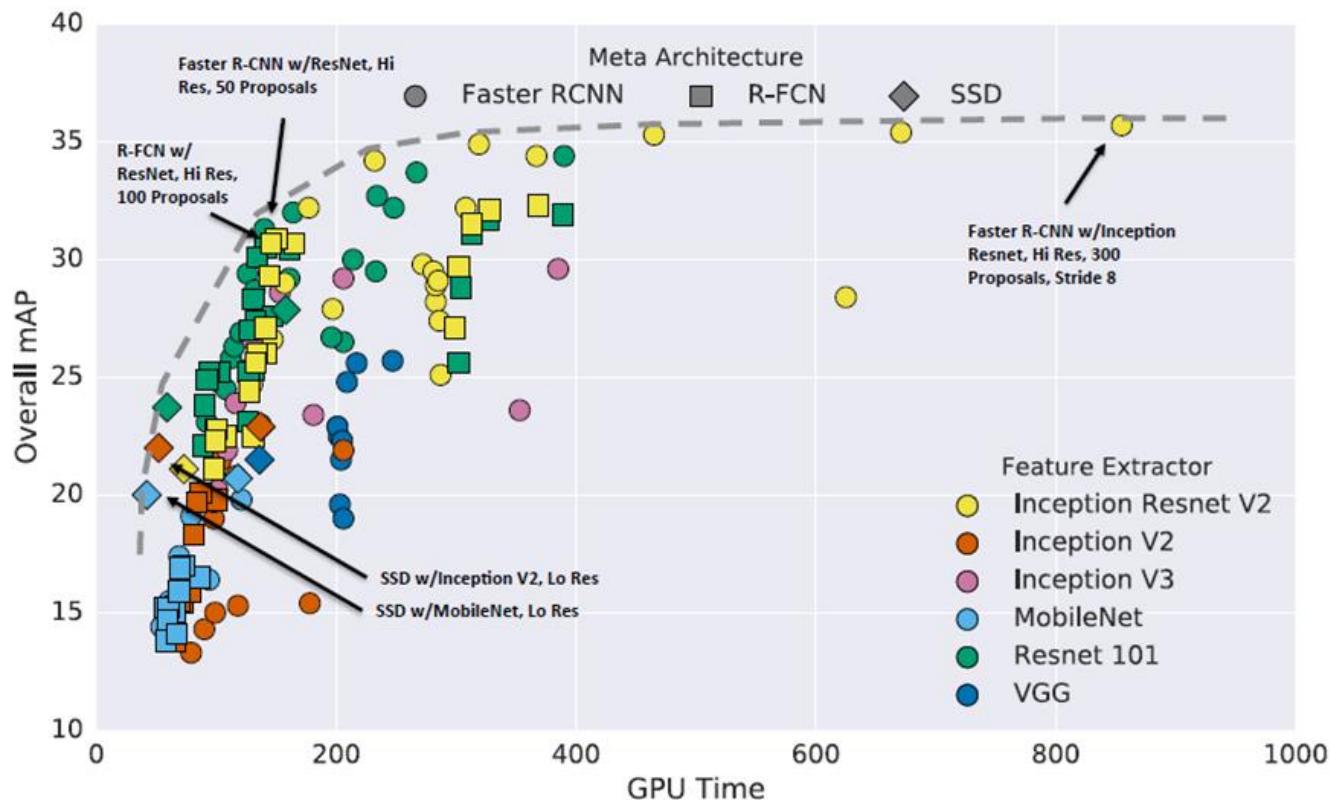


<https://pan.baidu.com/s/1pKIKIIB>

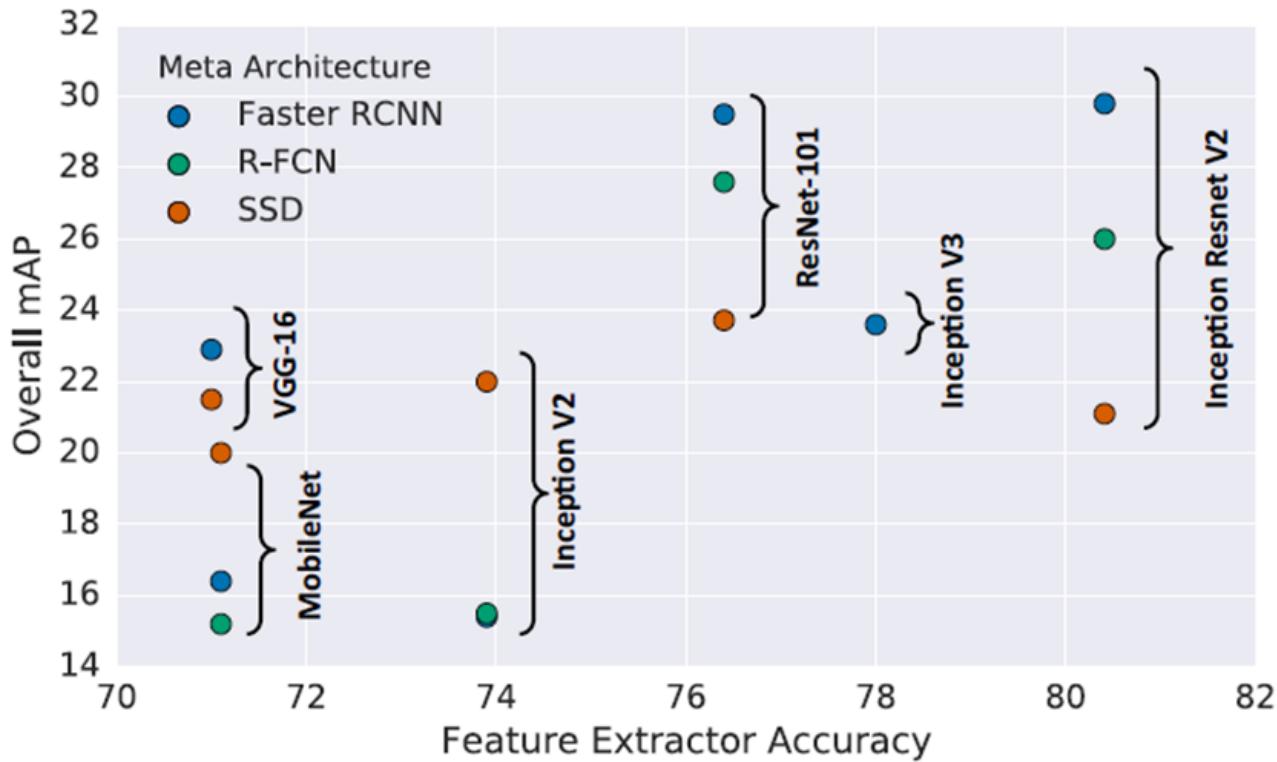
Recap: SSD



Accuracy VS Speed



1. Different Feature Extractor



2. Detect Object Size

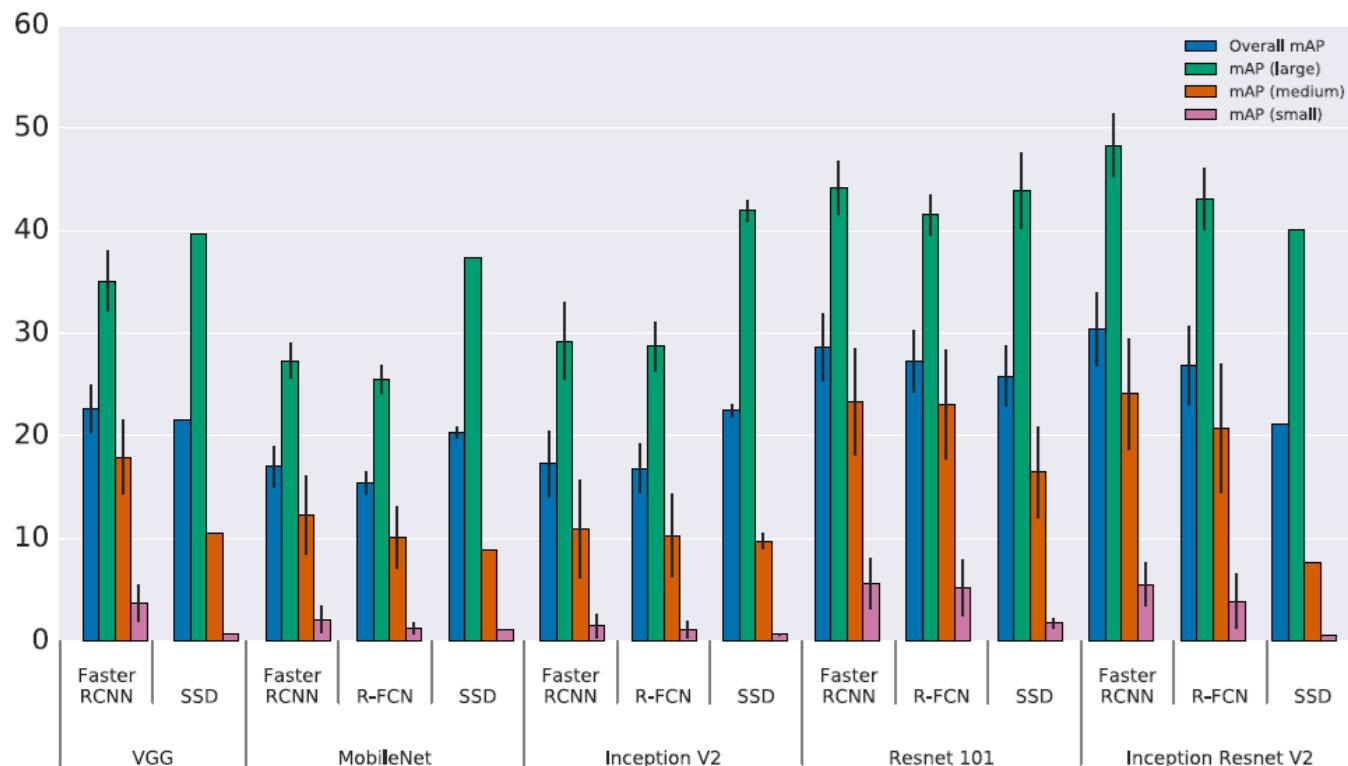
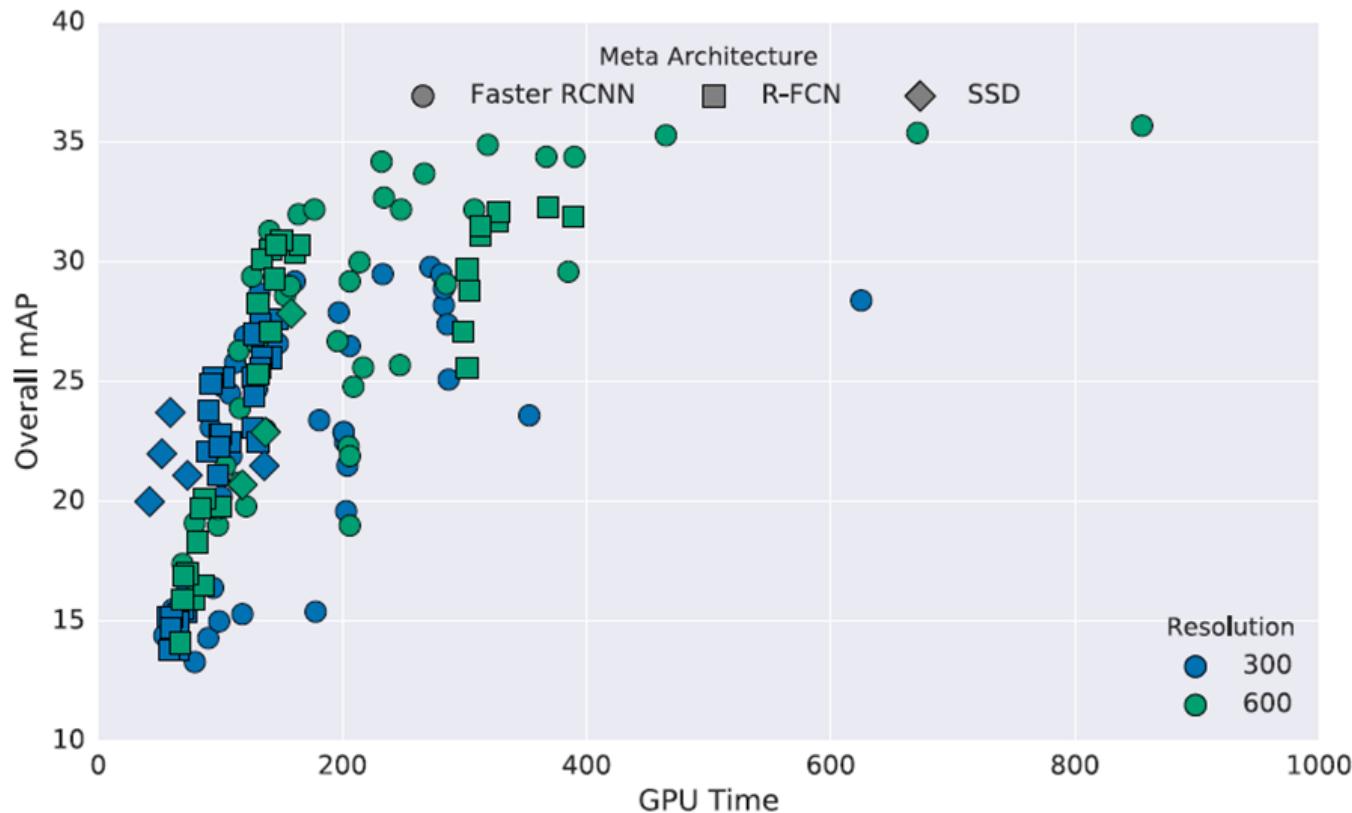
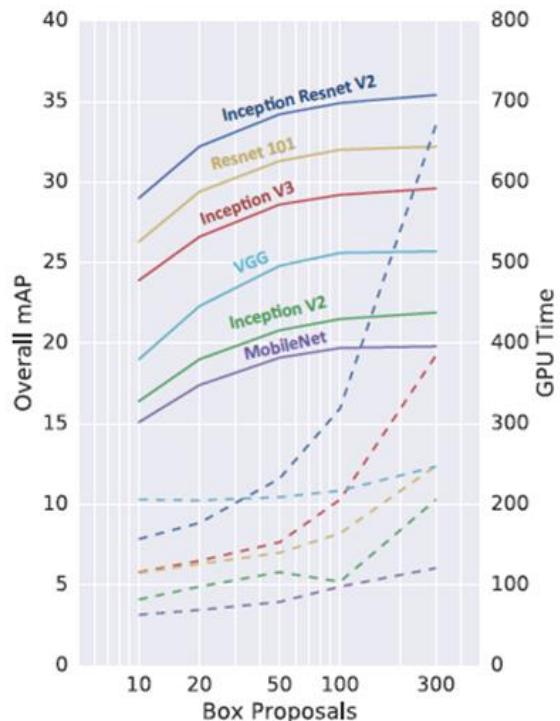


Figure 4: Accuracy stratified by object size, meta-architecture and feature extractor. We fix the image resolution to 300.

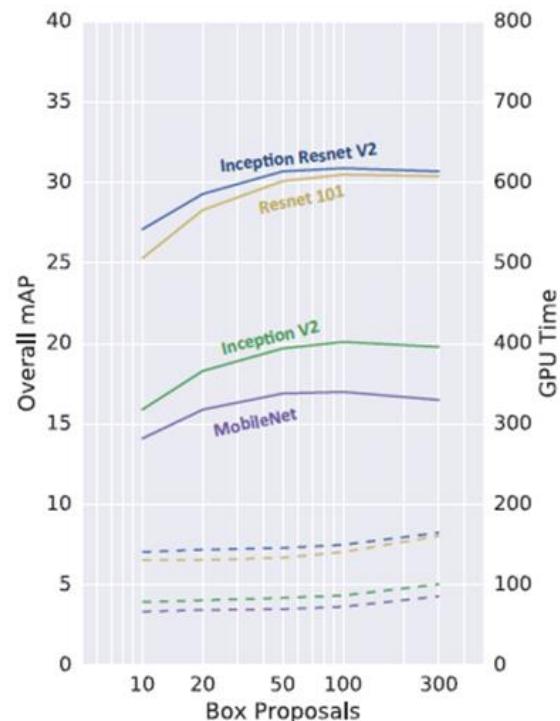
3. Image Resolution



4. Region Proposal Number



(a) FRCNN

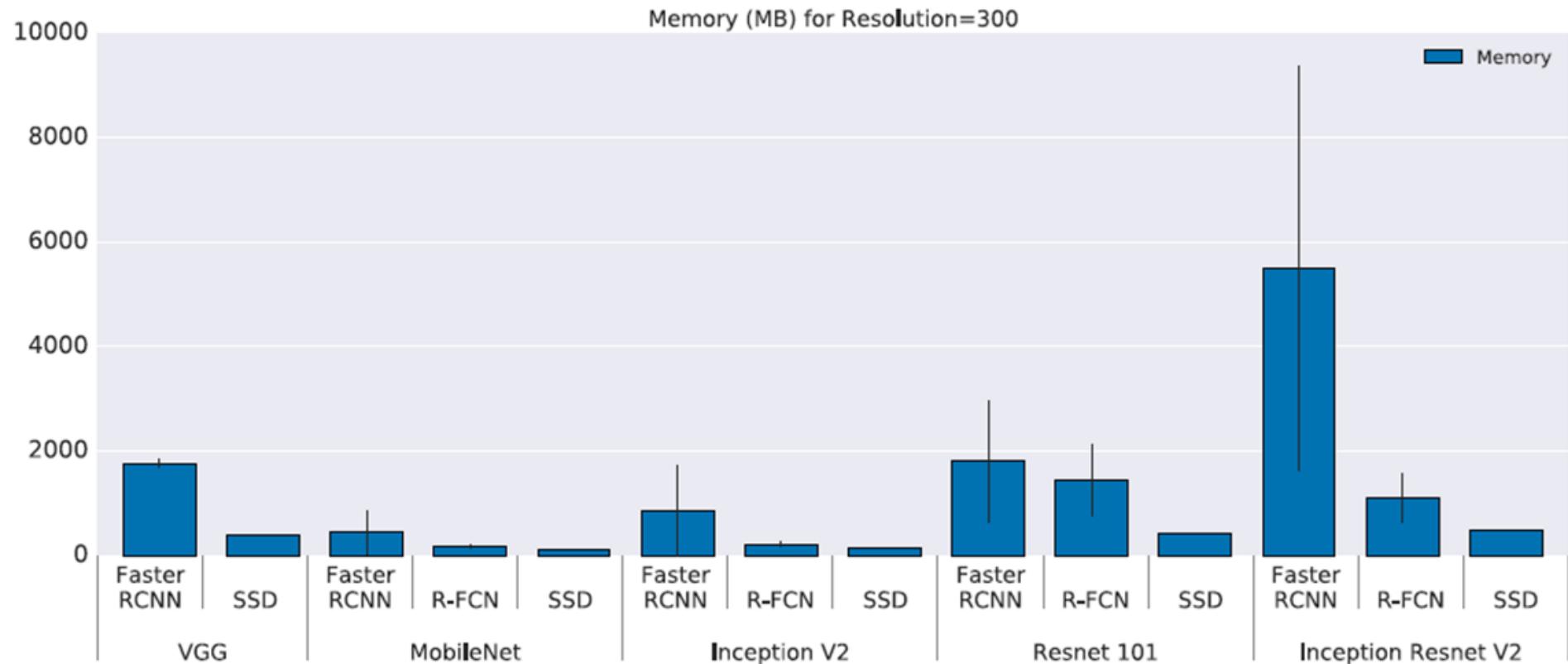


(b) RFCN

5. GPU Time



6. Memory



Summary

Speed First: SSD

Balance Speed and Accuracy: R-FCN

Accuracy First: Faster RCNN (Reduce proposals, it can speed up a lot with some accuracy loss)

Thanks!

Questions?



References

Girshick, Ross and Donahue, Jeff and Darrell, Trevor and Malik, Jitendra, [Rich feature hierarchies for accurate object detection and semantic segmentation](#), CVPR 2014

He, Kaiming and Zhang, Xiangyu and Ren, Shaoqing and Sun, Jian, [Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition](#), ECCV 2014

Girshick, Ross, [Fast R-CNN](#), ICCV 2015

Ren, Shaoqing and He, Kaiming and Girshick, Ross and Sun, Jian, [Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks](#), CVPR 2015

[Jifeng Dai, Yi Li, Kaiming He, Jian Sun](#) R-FCN: [Object Detection via Region-based Fully Convolutional Networks](#), NIPS 2016

Erhan, Dumitru and Szegedy, Christian and Toshev, Alexander and Anguelov, Dragomir, [Scalable Object Detection using Deep Neural Networks](#), CVPR 2014

Bell, Sean and Lawrence Zitnick, C and Bala, Kavita and Girshick, Ross, [Inside-Outside Net: Detecting Objects in Context with Skip Pooling and Recurrent Neural Networks](#), CVPR 2016

Redmon, Joseph and Divvala, Santosh and Girshick, Ross and Farhadi, Ali, [You Only Look Once: Unified, Real-Time Object Detection](#), CVPR 2016

Liu, Wei and Anguelov, Dragomir and Erhan, Dumitru and Szegedy, Christian and Reed, Scott and Fu, Cheng-Yang and Berg, Alexander C, [SSD: Single Shot MultiBox Detector](#), ECCV 2016

Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, Serge Belongie, [Feature Pyramid Networks for Object Detection](#), arXiv 2016

Huang, Jonathan and Rathod, Vivek and Sun, Chen and Zhu, Menglong and Korattikara, Anoop and Fathi, Alireza and Fischer, Ian and Wojna, Zbigniew and Song, Yang and Chen, Qizhe and Yu, Yuxin and Liu, Ming and Guo, Jun and Wang, Zhen and Chen, Yuxin and Yu, Fei-Fei and Yuille, Alan L, [Context Encoding for Semantic Segmentation](#), arXiv 2016