

SSD: Single Shot MultiBox Detector

Wei Liu(1), **Dragomir Anguelov(2)**, Dumitru Erhan(3), Christian Szegedy(3),
Scott Reed(4), Cheng-Yang Fu(1), Alexander C. Berg(1)

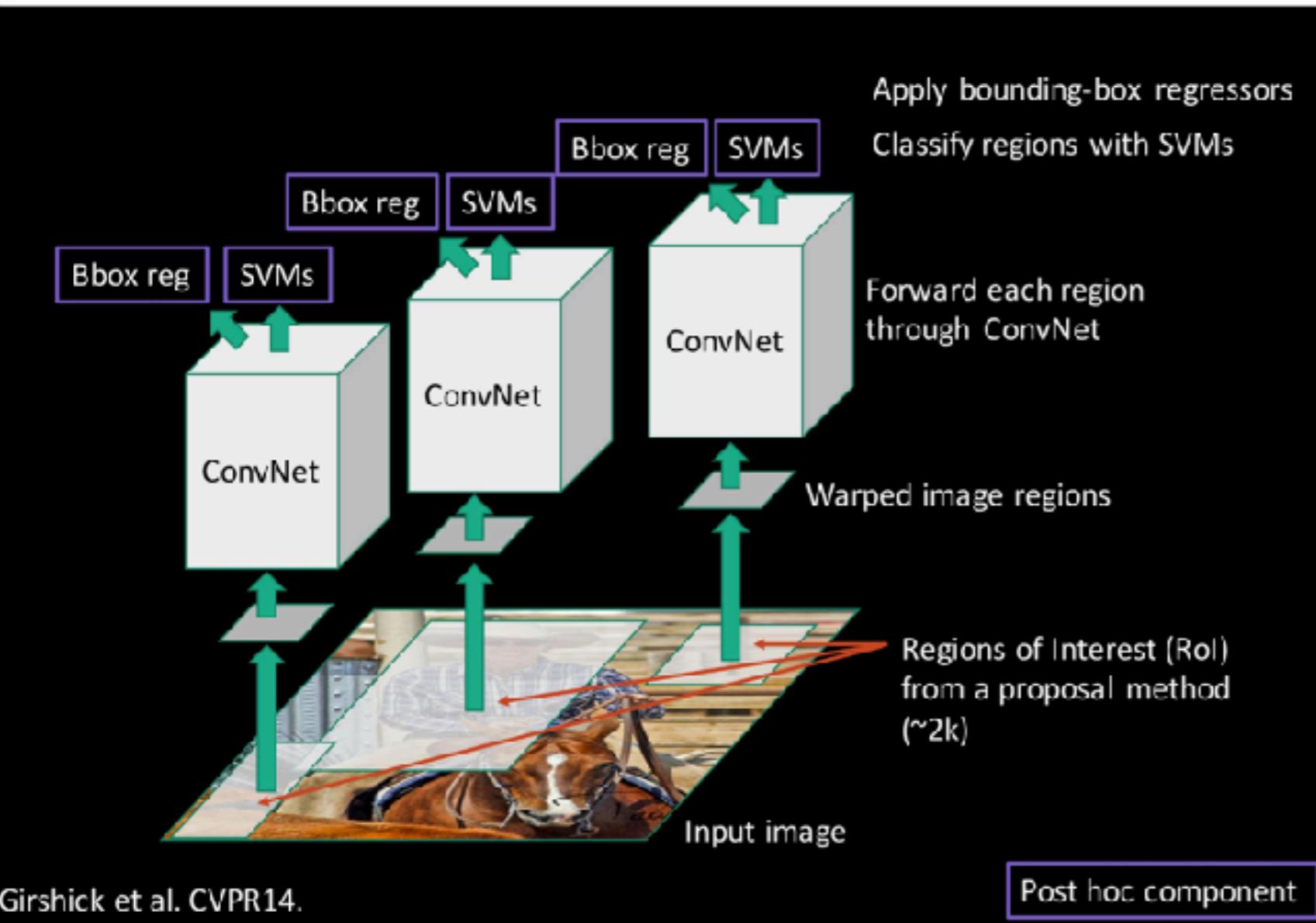
UNC Chapel Hill(1), **Zoox Inc.(2)**, Google Inc.(3),
University of Michigan(4)

Presented by Hongyan Wang and Nathan Watts

Classical Object Detection

- Region selection: Sliding Window
- Feature extraction: SIFT, HOG
- Classification: SVM, Adaboost

Putting it together: R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014

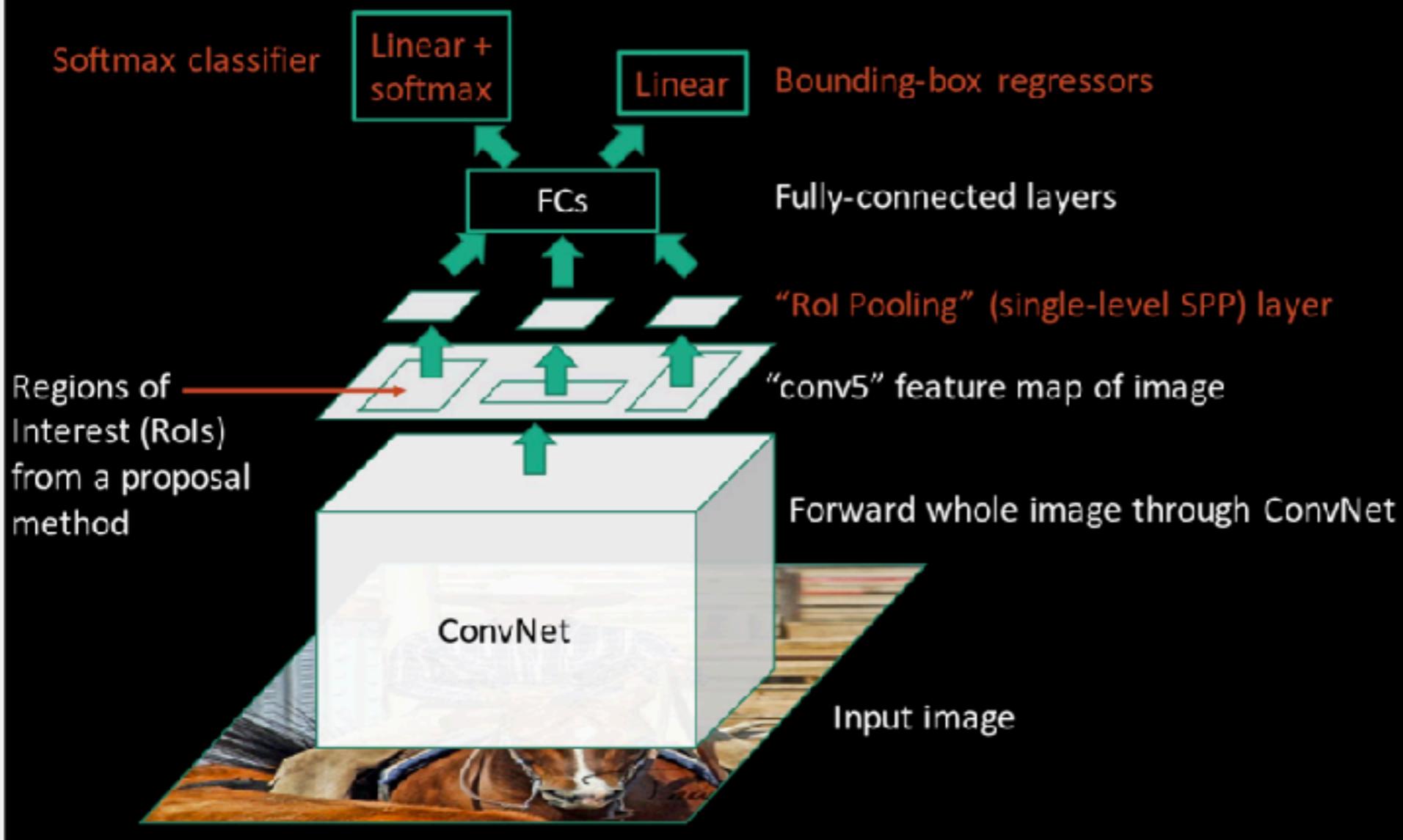
Slide credit: Ross Girshick

R-CNN problems

- Slow at test time: need to run full forward pass of CNN for each region proposal.
- Complex training pipeline

Original slides are from http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

Fast R-CNN (test time)



Girschick, "Fast R-CNN", ICCV 2015

Slide credit: Ross Girschick

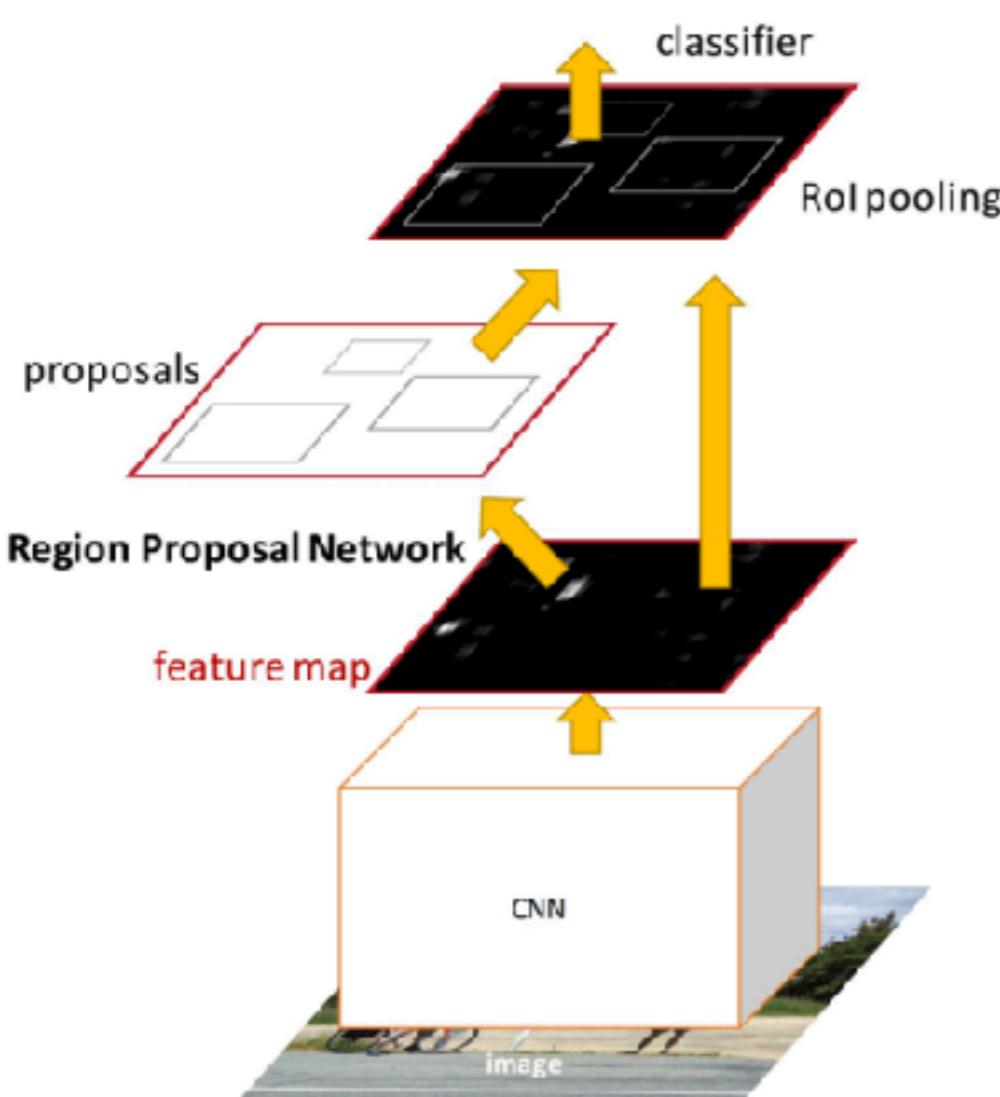
Original slides are from http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

Fast R-CNN problem

- Region proposal is the bottleneck.

Original slides are from http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

Faster R-CNN:



Insert a **Region Proposal Network (RPN)** after the last convolutional layer

RPN trained to produce region proposals directly; no need for external region proposals!

After RPN, use RoI Pooling and an upstream classifier and bbox regressor just like Fast R-CNN

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

Slide credit: Ross Girshick

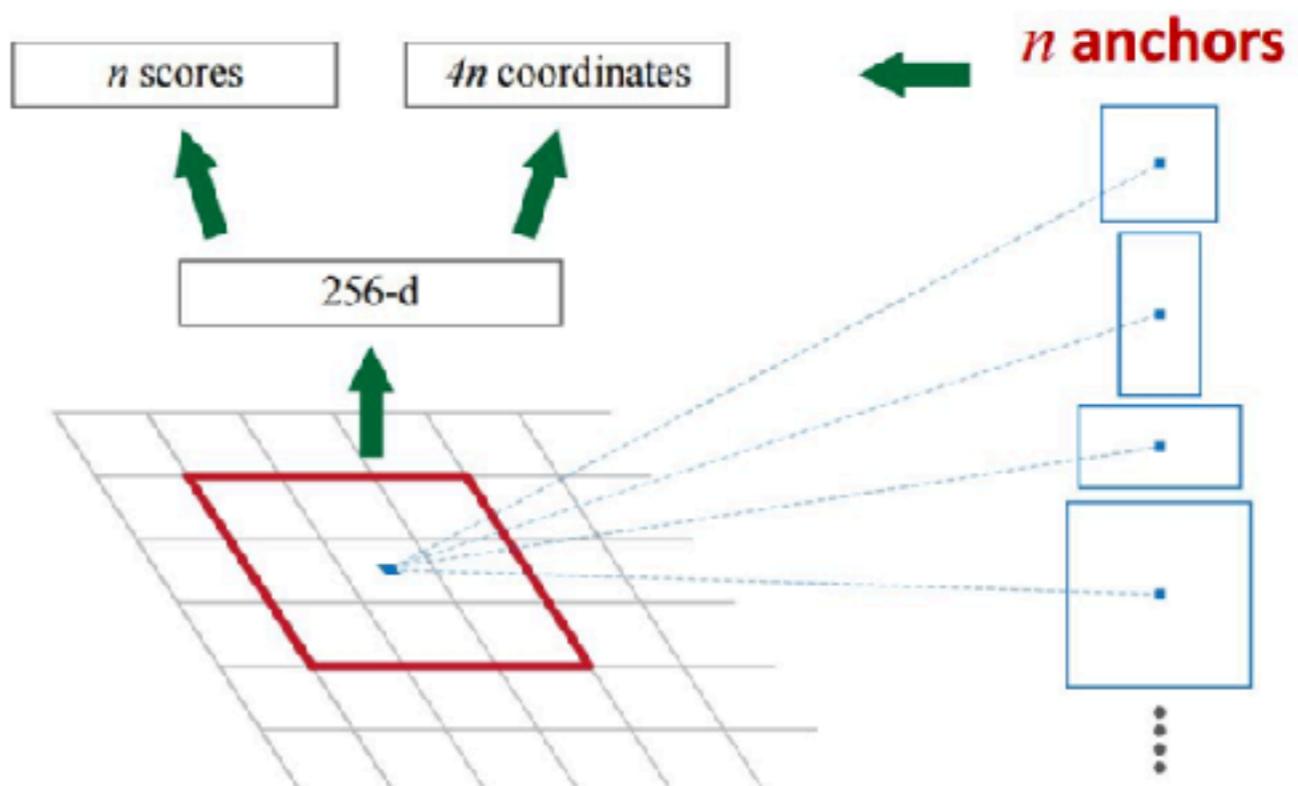
Faster R-CNN: Region Proposal Network

Use **N anchor boxes** at each location

Anchors are **translation invariant**: use the same ones at every location

Regression gives offsets from anchor boxes

Classification gives the probability that each (regressed) anchor shows an object



Original slides are from http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

Faster R-CNN problem

- Still slow for real-time detection.

Original slides are from http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

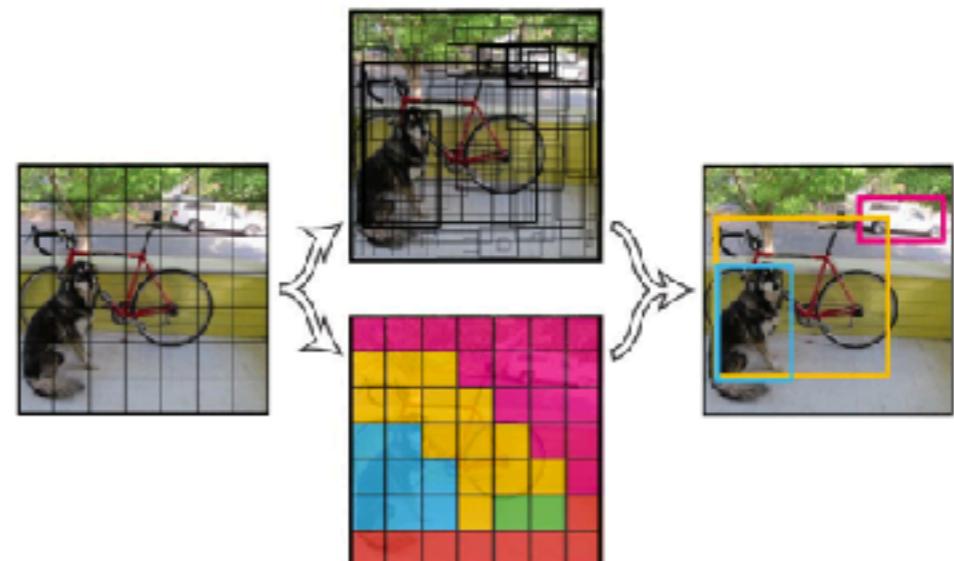
YOLO: You Only Look Once Detection as Regression

Divide image into $S \times S$ grid

Within each grid cell predict:

B Boxes: 4 coordinates + confidence

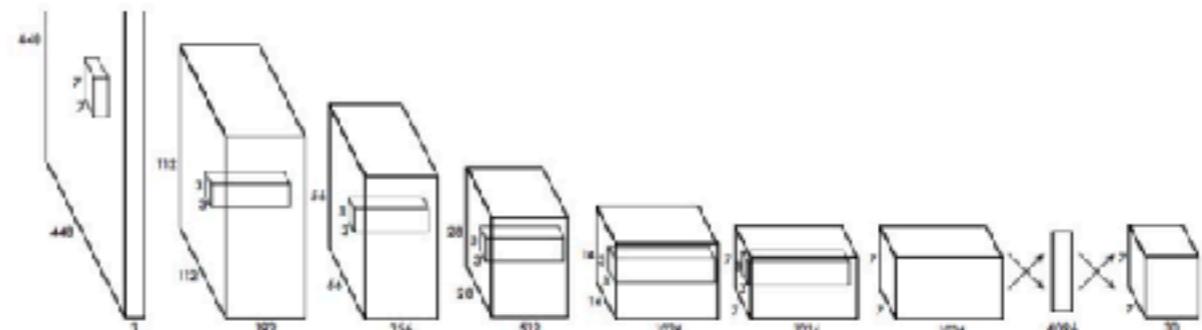
Class scores: C numbers



Regression from image to
 $7 \times 7 \times (5 * B + C)$ tensor

Direct prediction using a CNN

Redmon et al, "You Only Look Once:
Unified, Real-Time Object Detection", arXiv 2015



Original slides are from http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

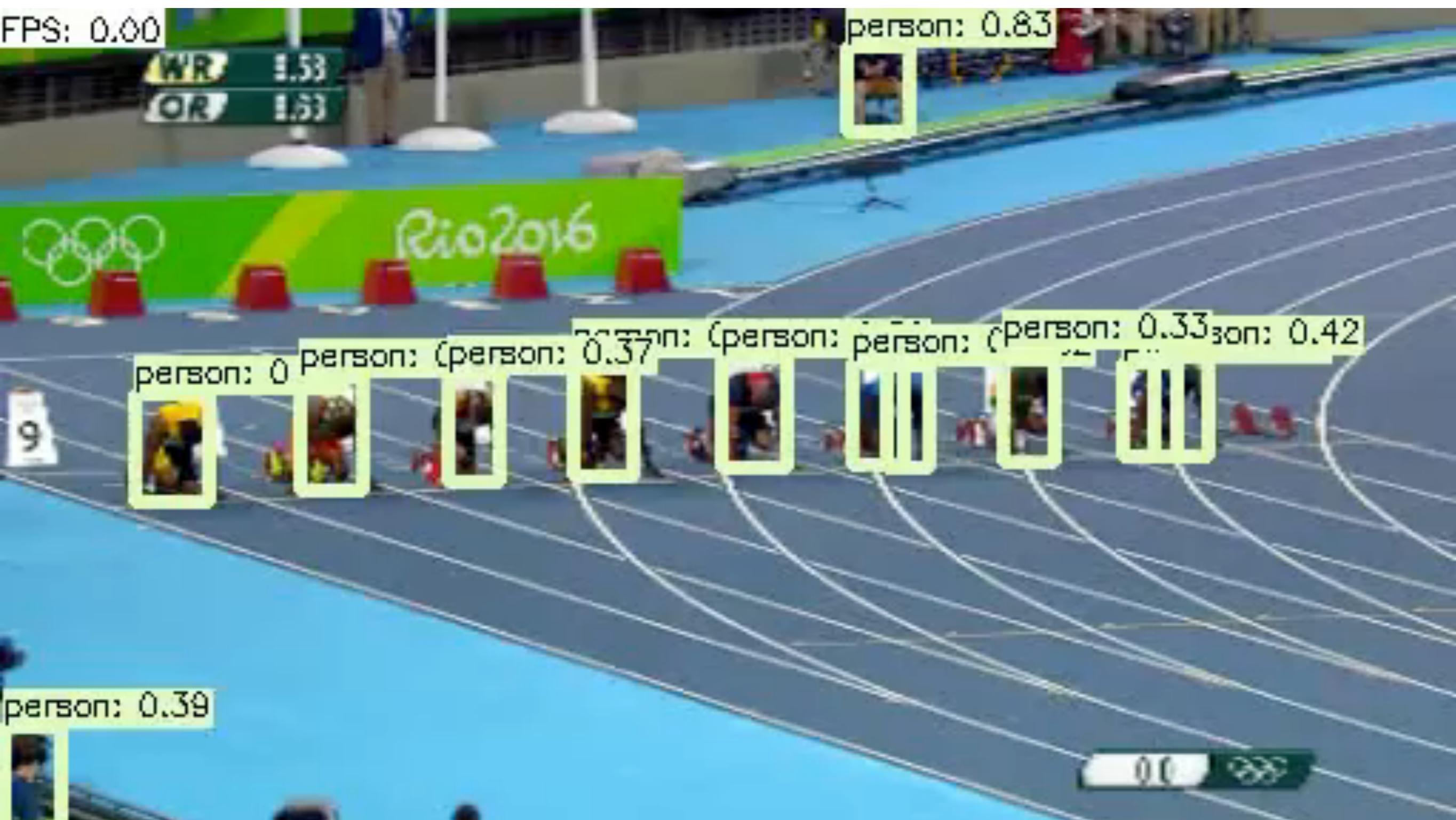
YOLO problems

- Accuracy is much worse than faster R-CNN.
- Not good at small objects.

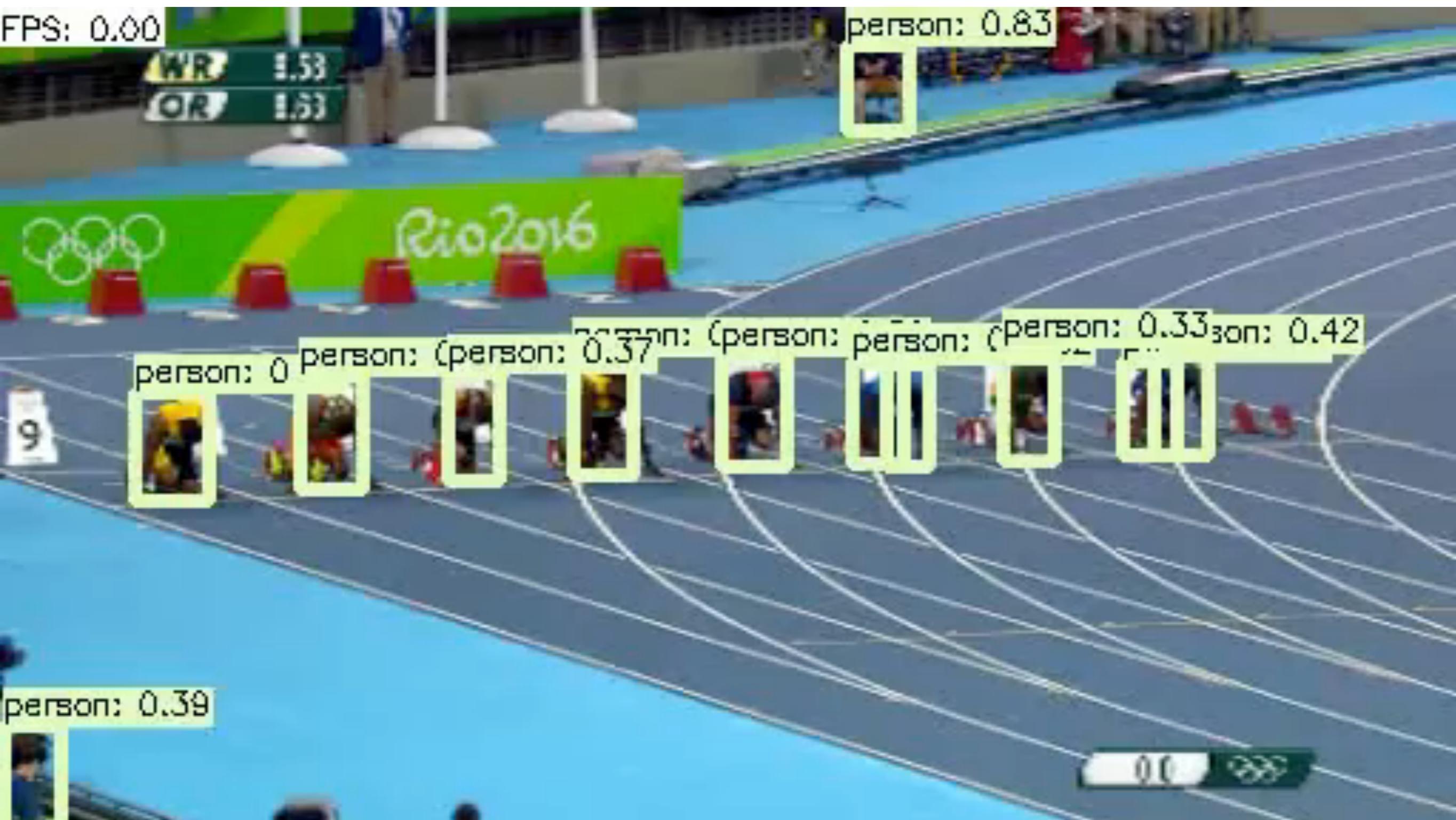
Original slides are from http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

Next ? SSD !

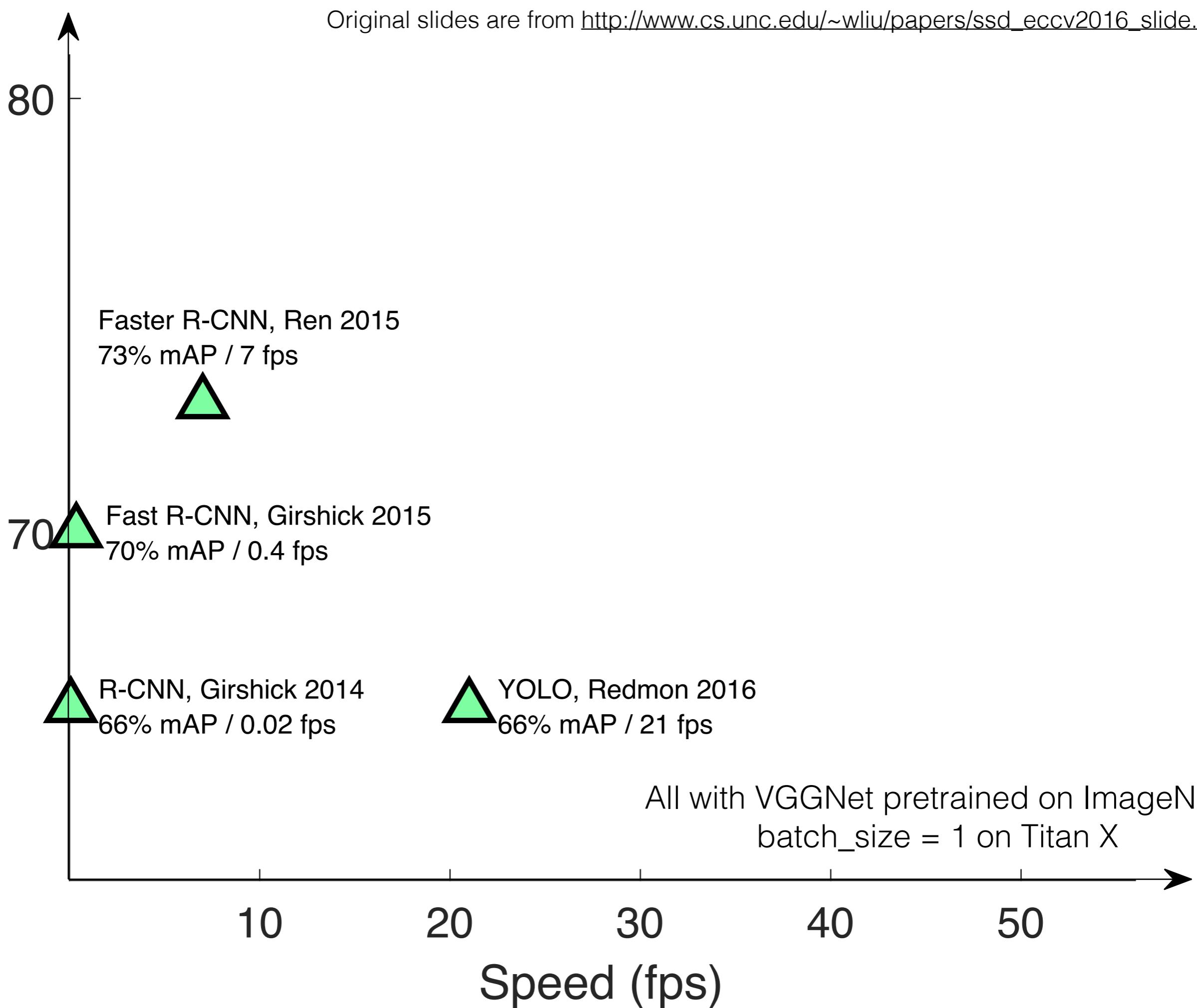
VGGNet
Titan X Pascal



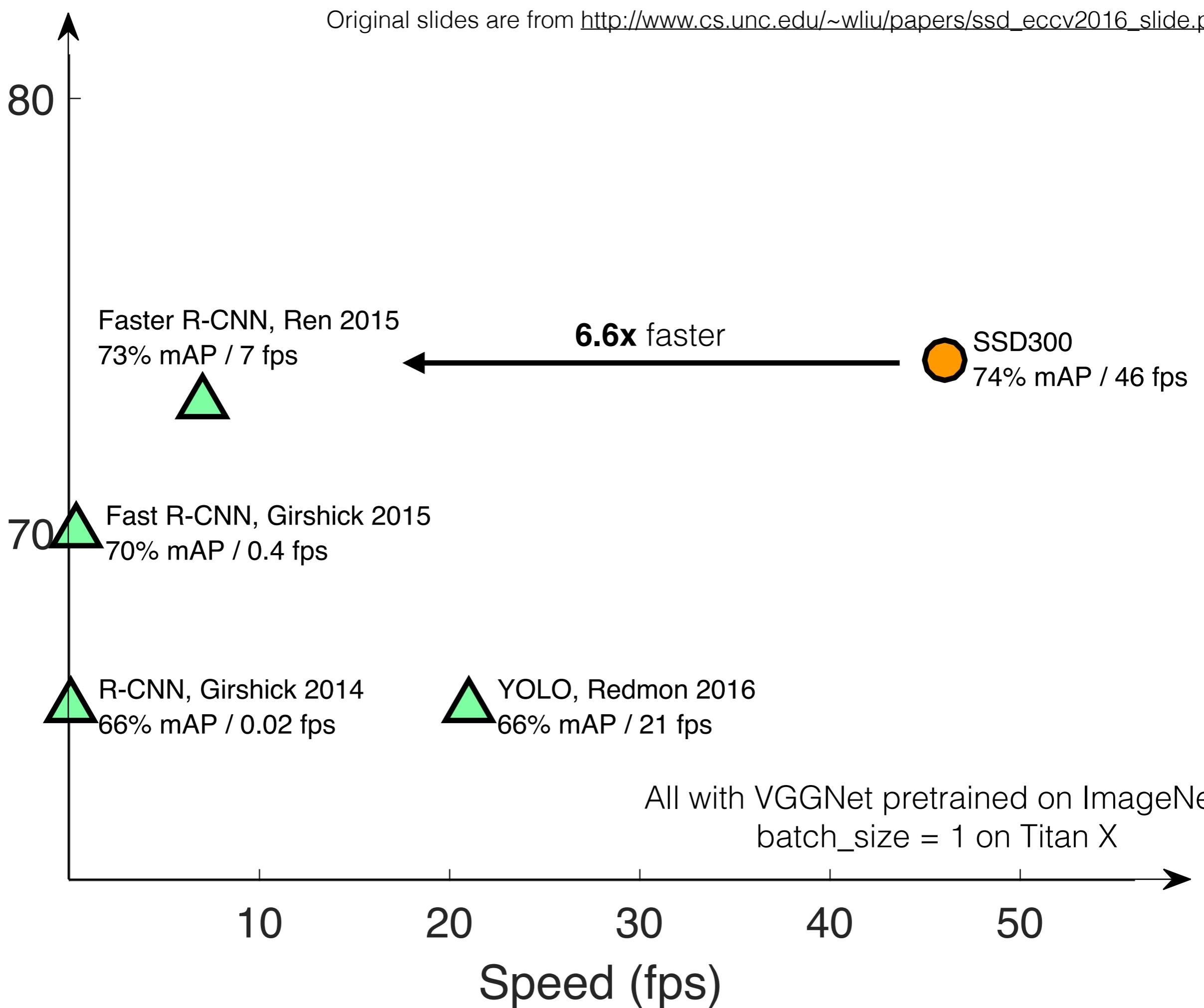
VGGNet
Titan X Pascal

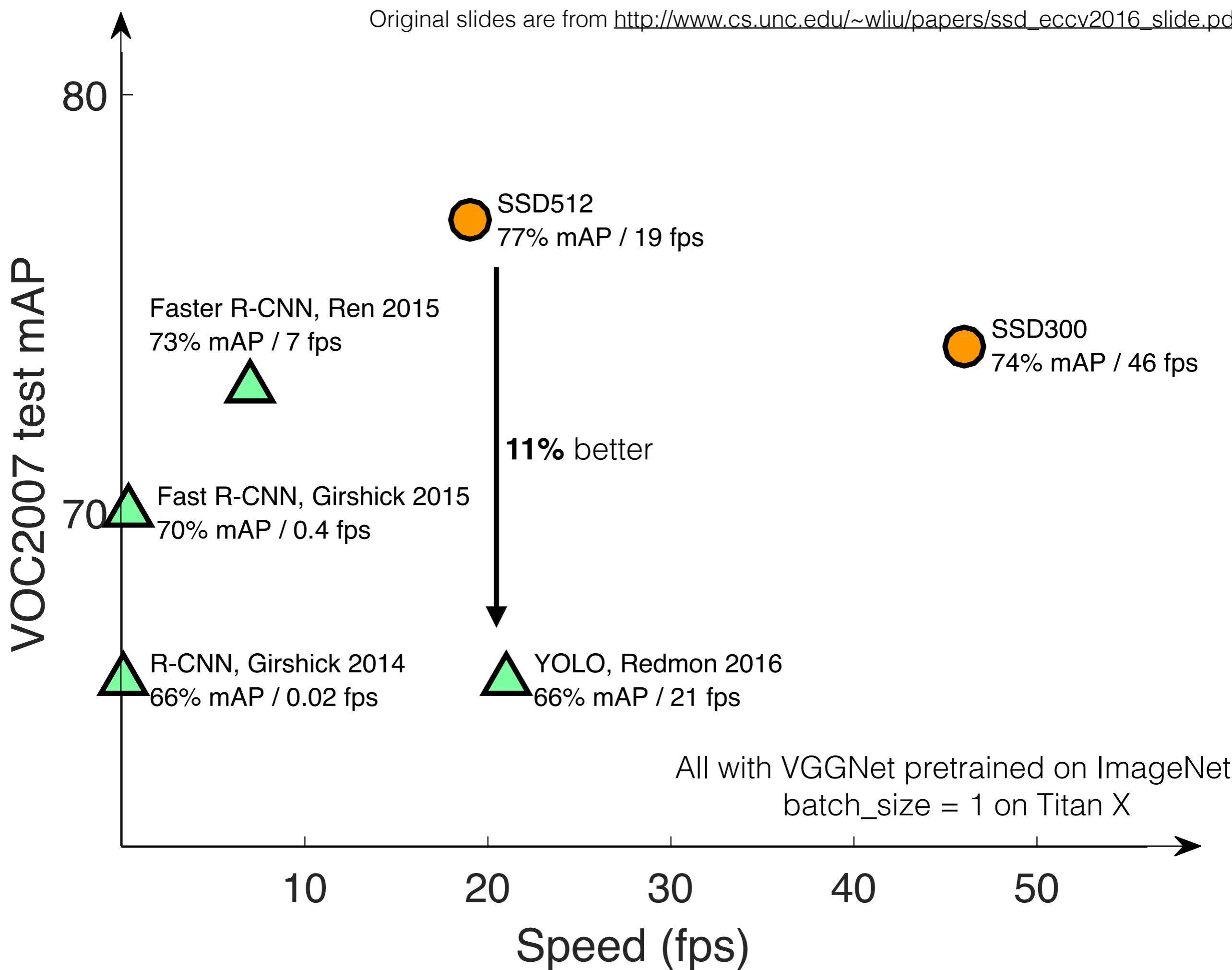


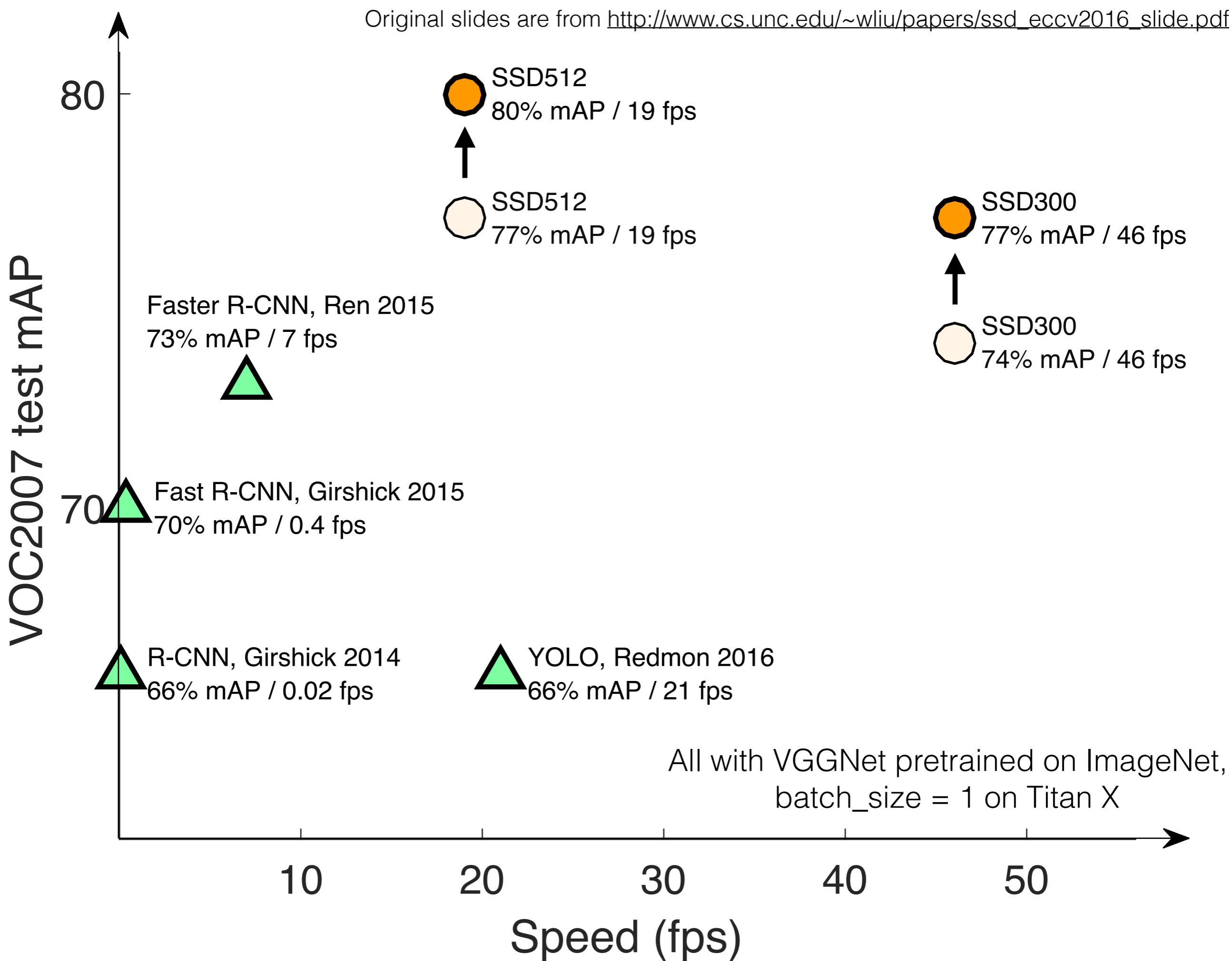
VOC2007 test mAP

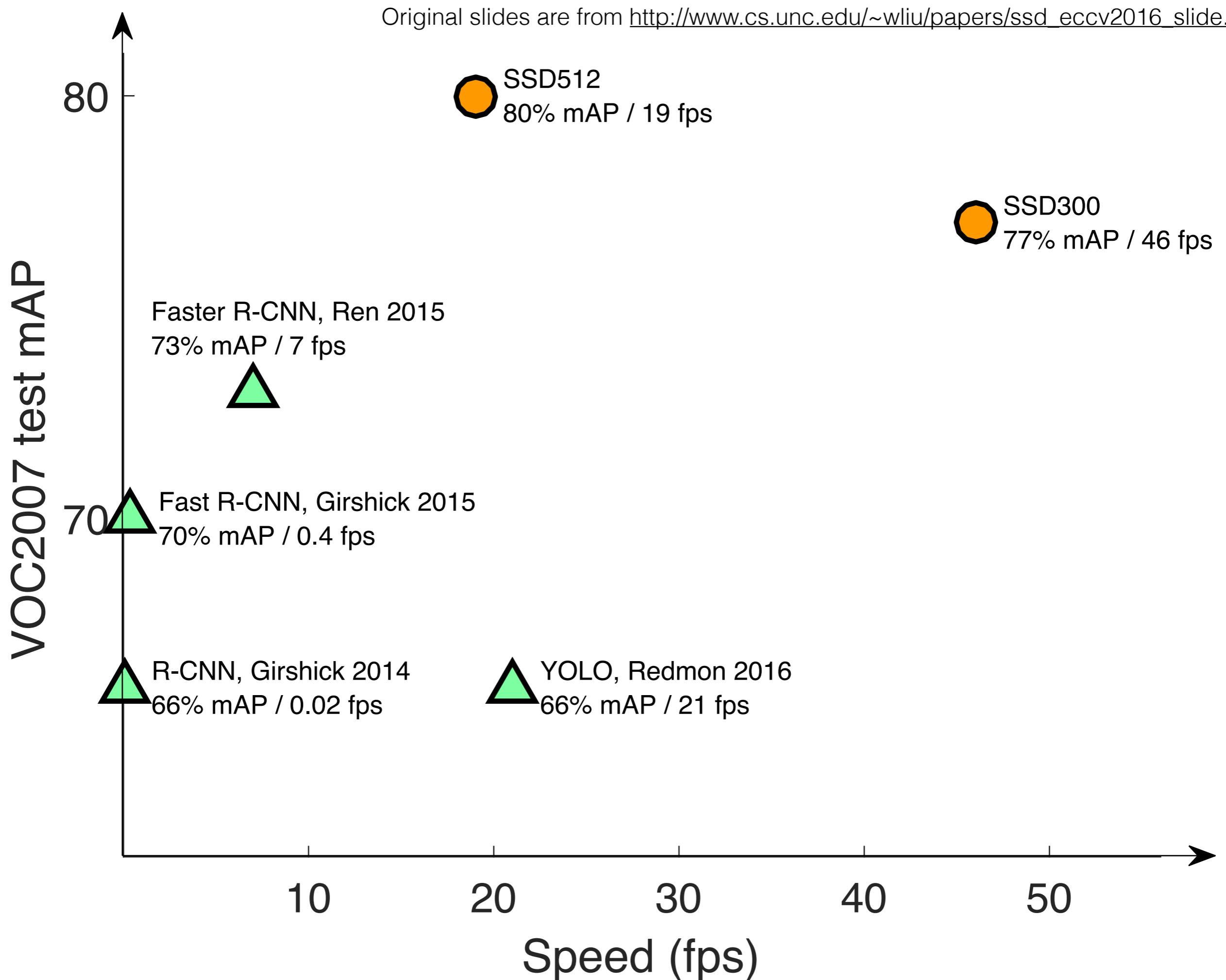


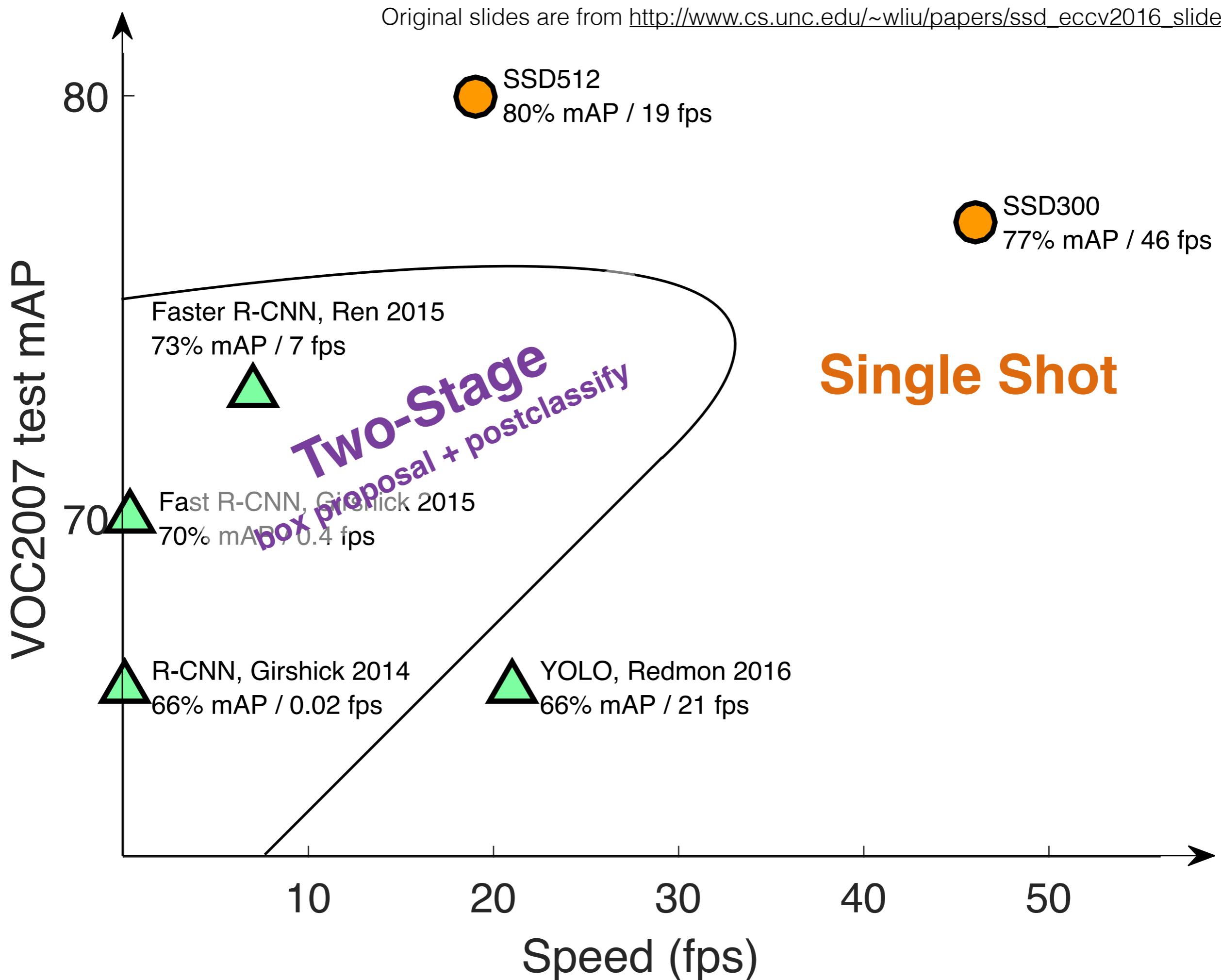
VOC2007 test mAP











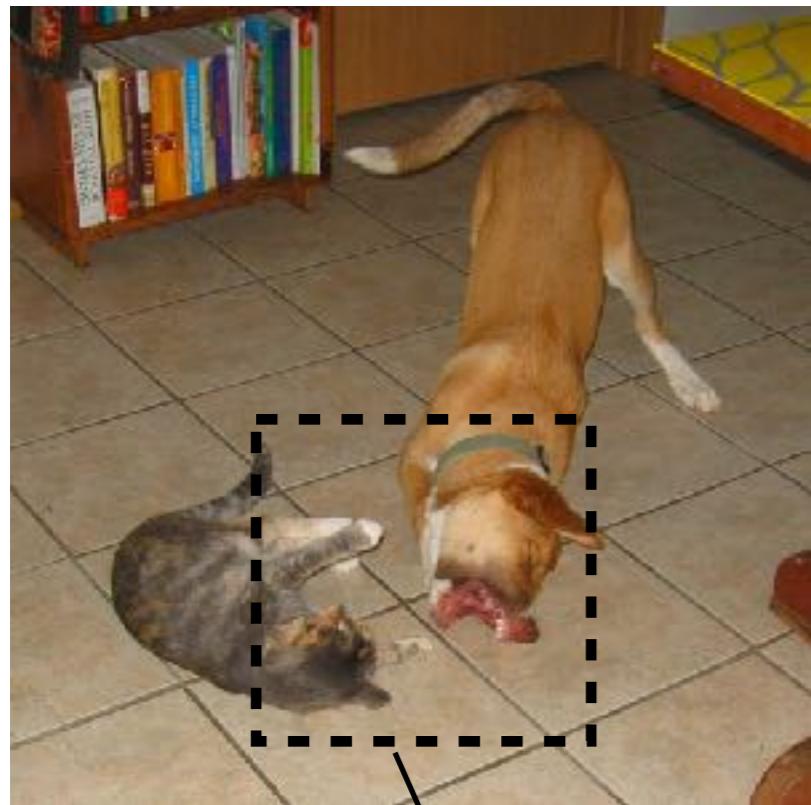
Bounding Box Prediction

Classical sliding
windows



Bounding Box Prediction

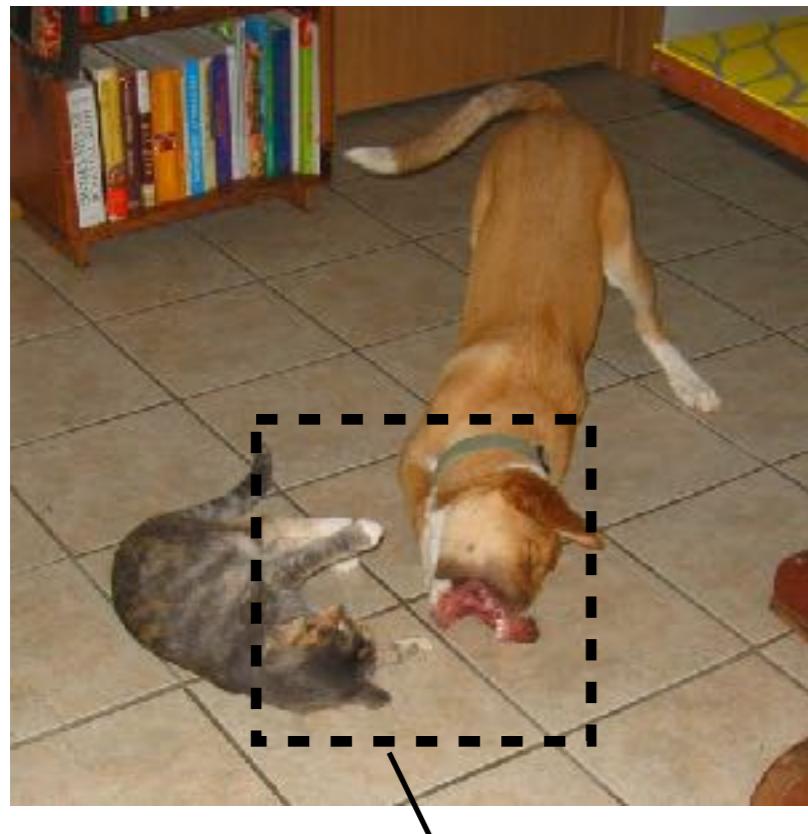
Classical sliding
windows



Is it a cat? **No**

Bounding Box Prediction

Classical sliding
windows

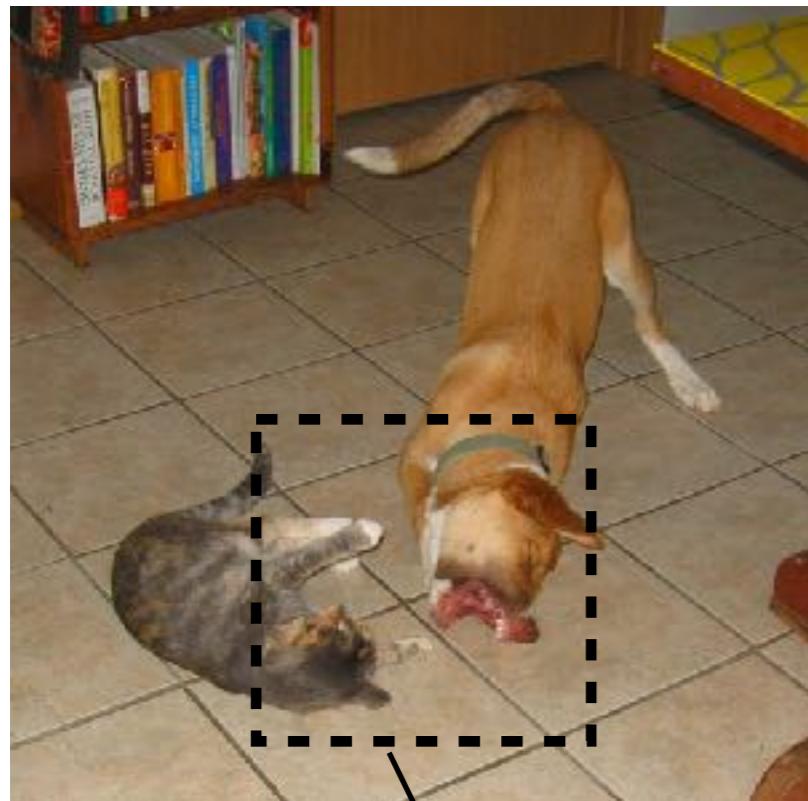


Is it a cat? **No**

Discretize the box space **densely**

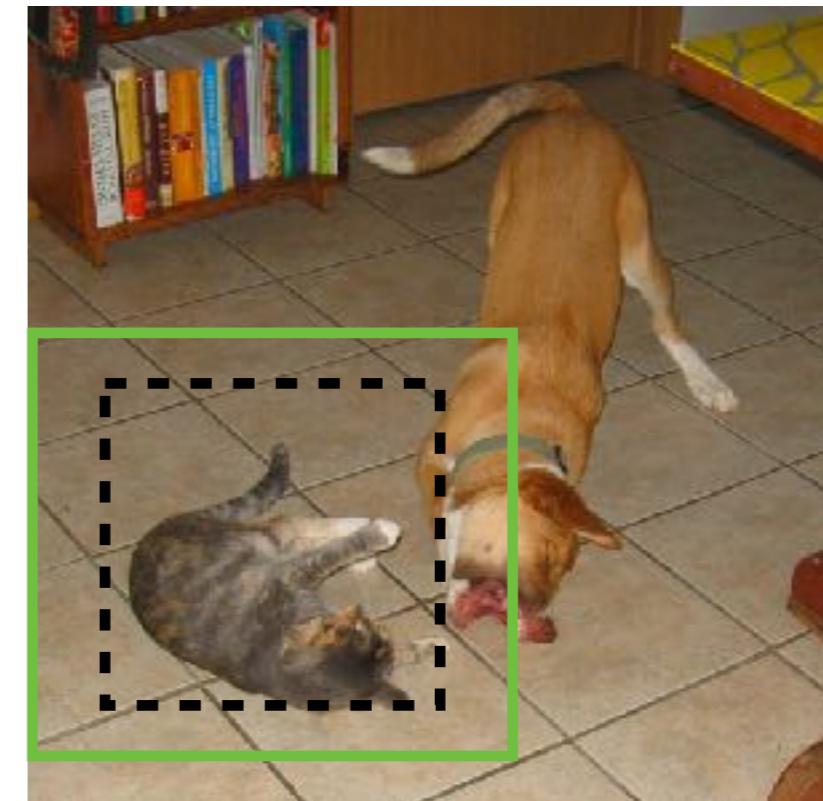
Bounding Box Prediction

Classical sliding windows



Is it a cat? **No**

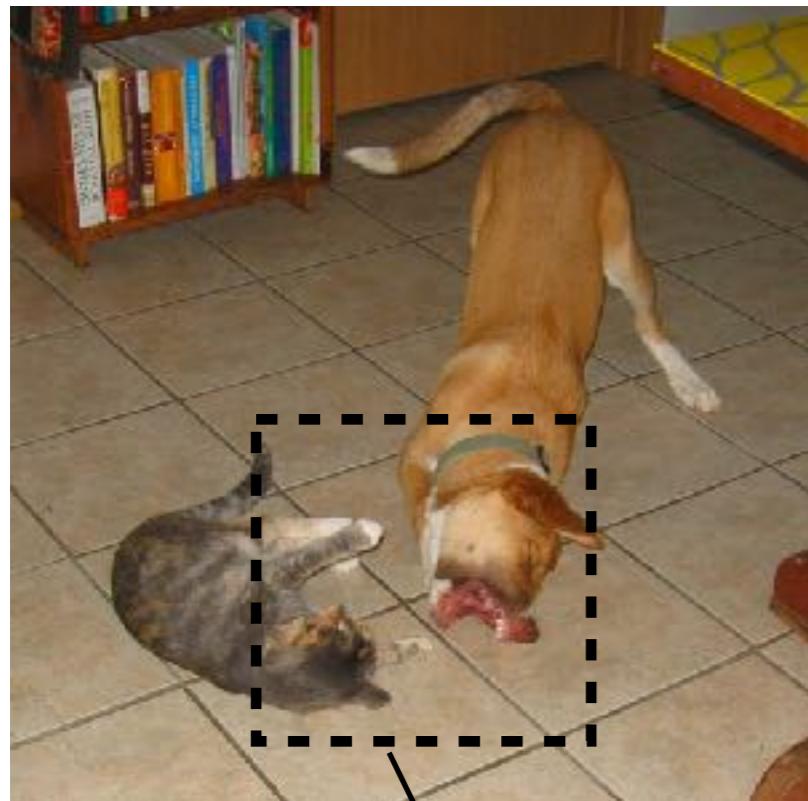
SSD and other deep approaches



Discretize the box space **densely**

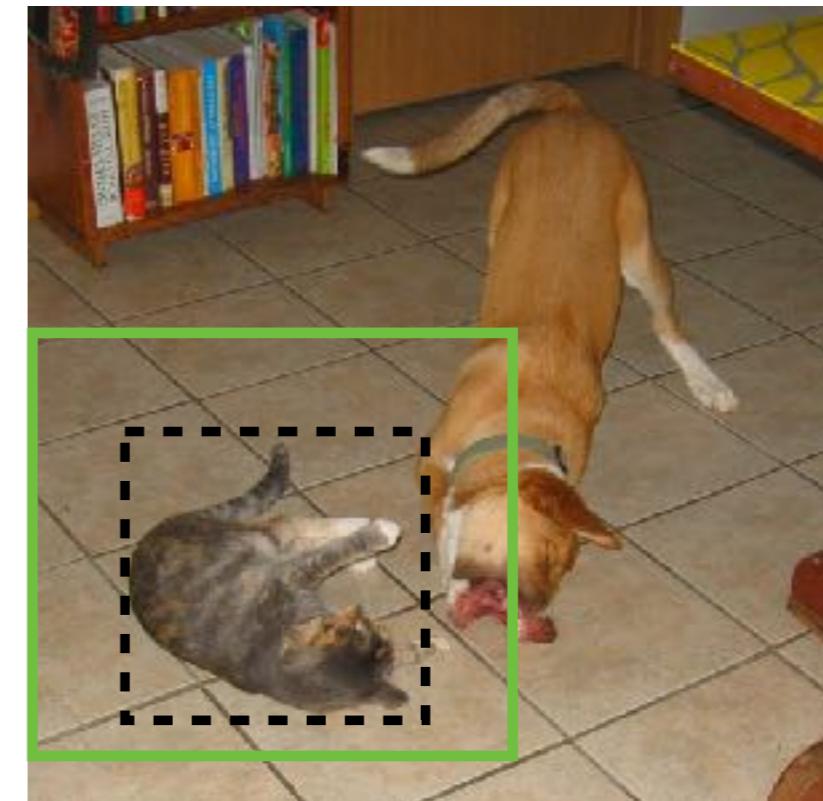
Bounding Box Prediction

Classical sliding windows



Is it a cat? **No**

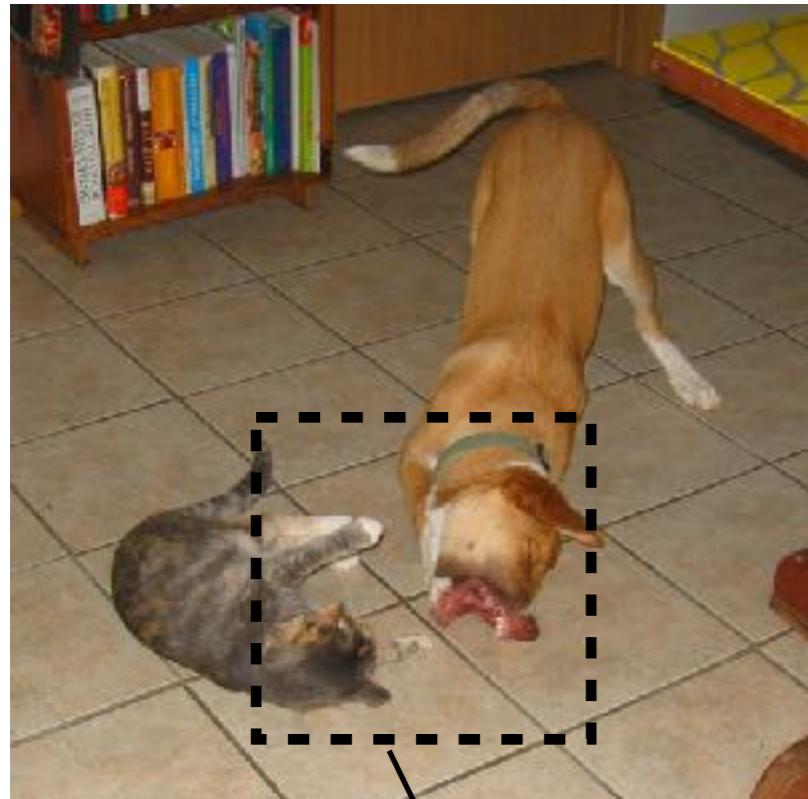
SSD and other deep approaches



Discretize the box space **densely**

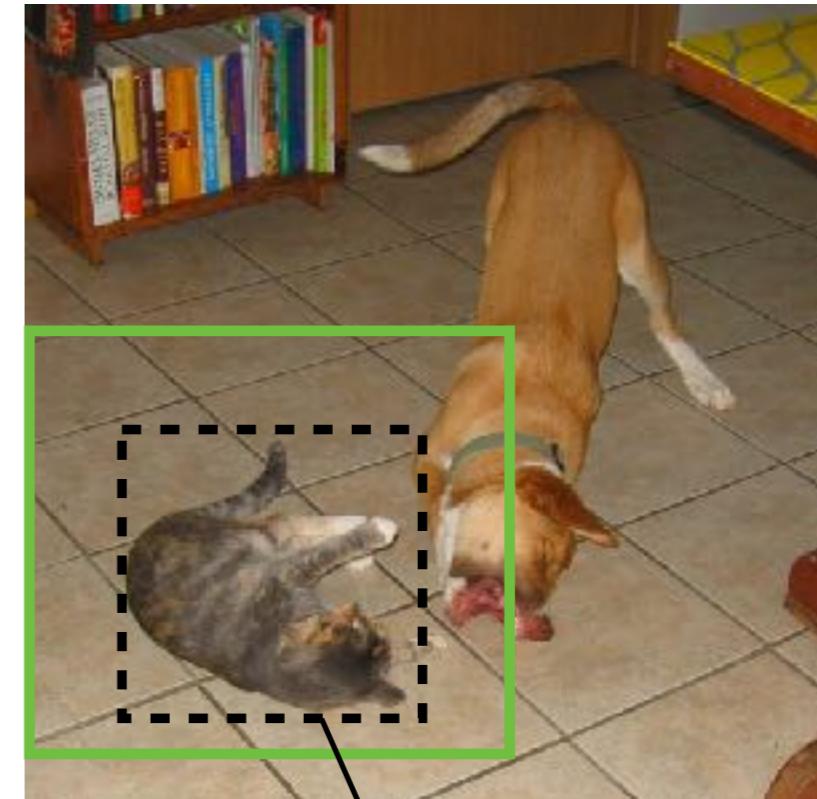
Bounding Box Prediction

Classical sliding windows



Is it a cat? **No**

SSD and other deep approaches

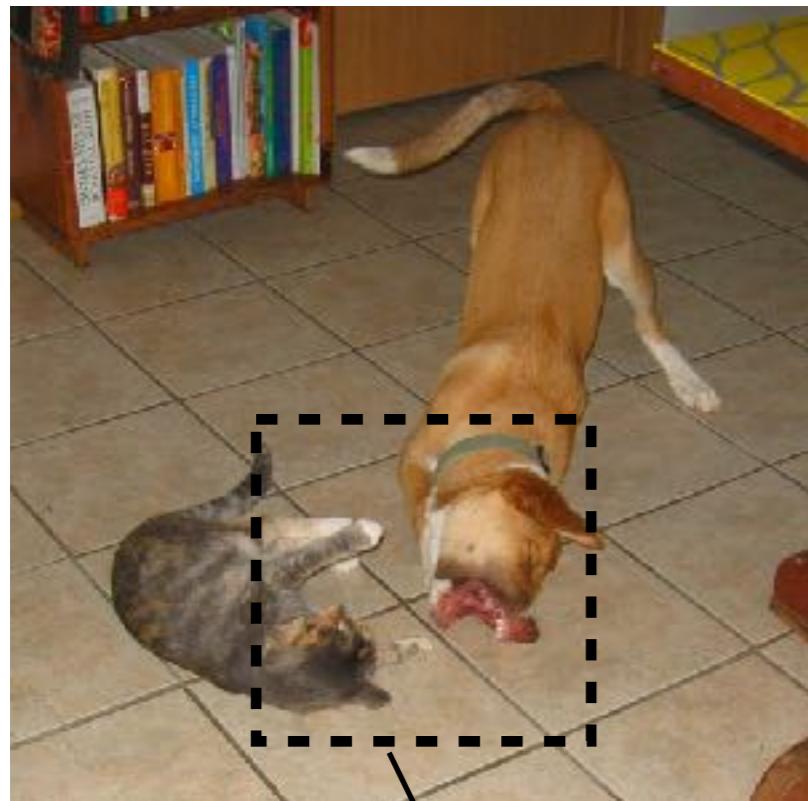


cat: 0.8 dog: 0.1

Discretize the box space **densely**

Bounding Box Prediction

Classical sliding windows



Is it a cat? **No**

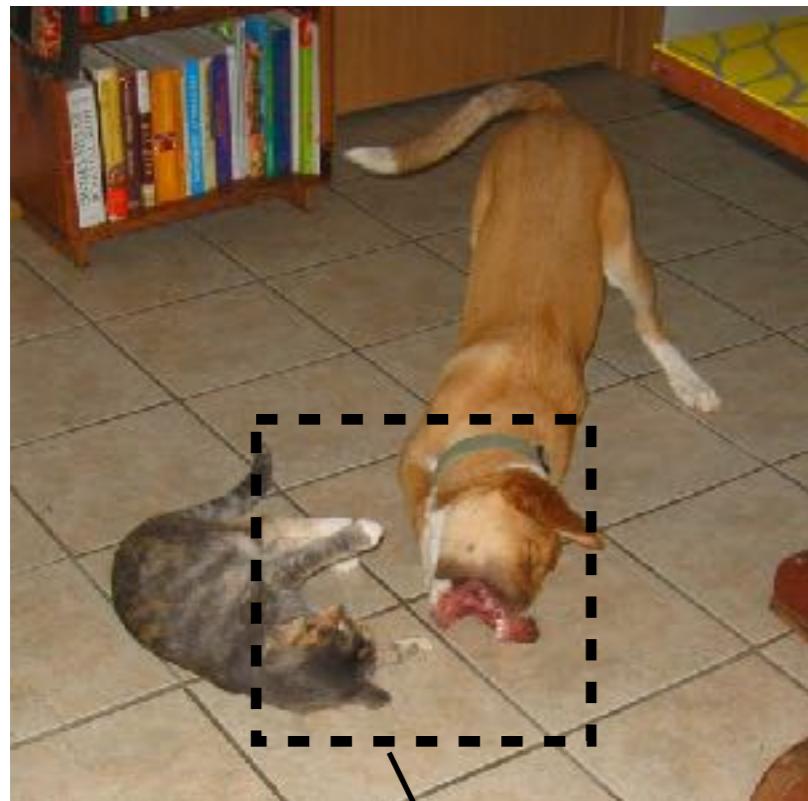
SSD and other deep approaches



Discretize the box space **densely**

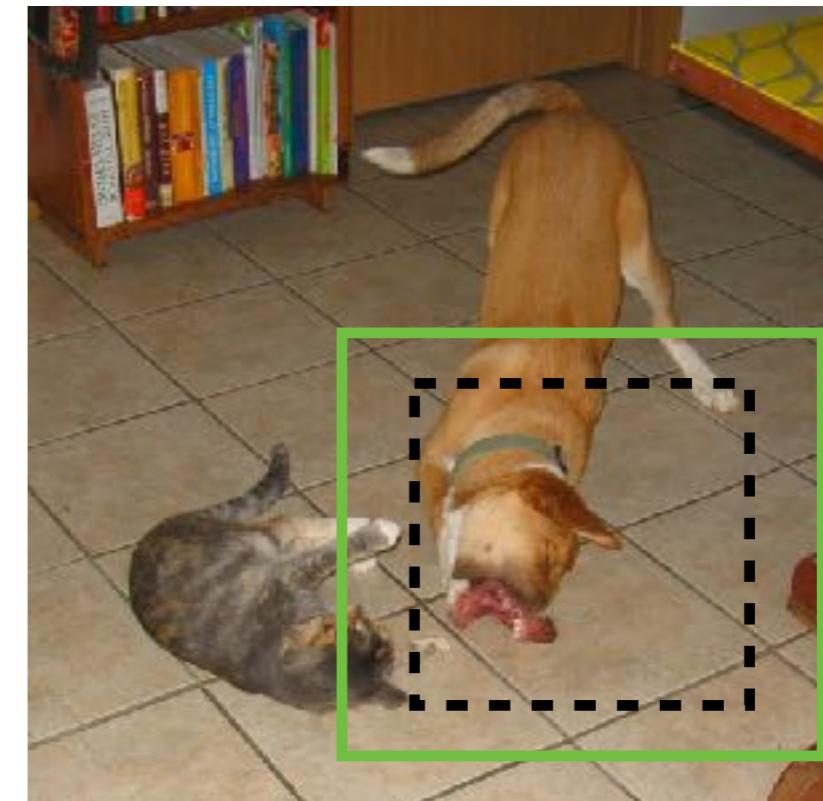
Bounding Box Prediction

Classical sliding windows



Is it a cat? **No**

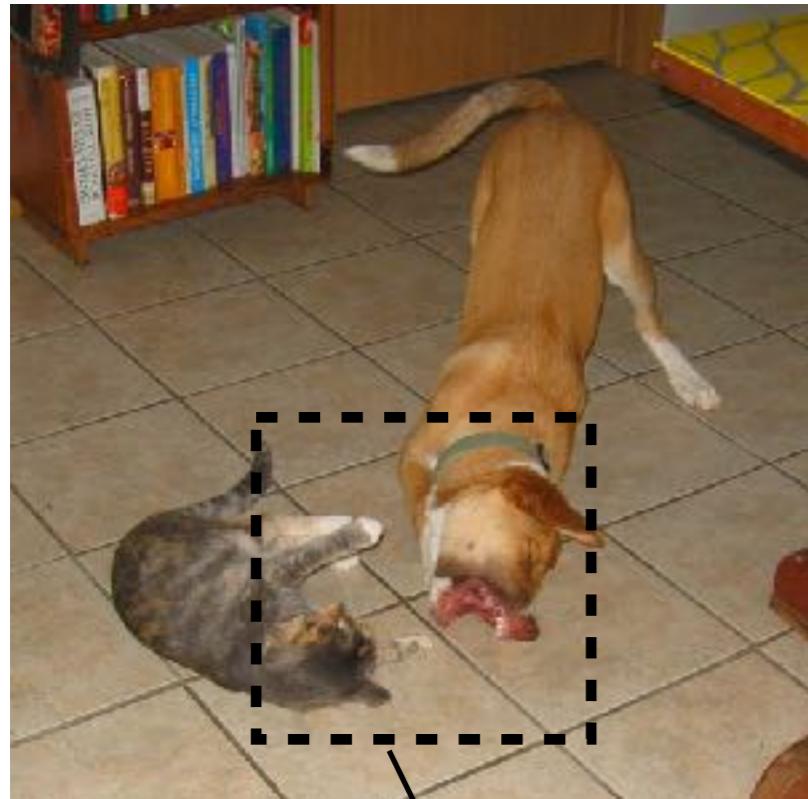
SSD and other deep approaches



Discretize the box space **densely**

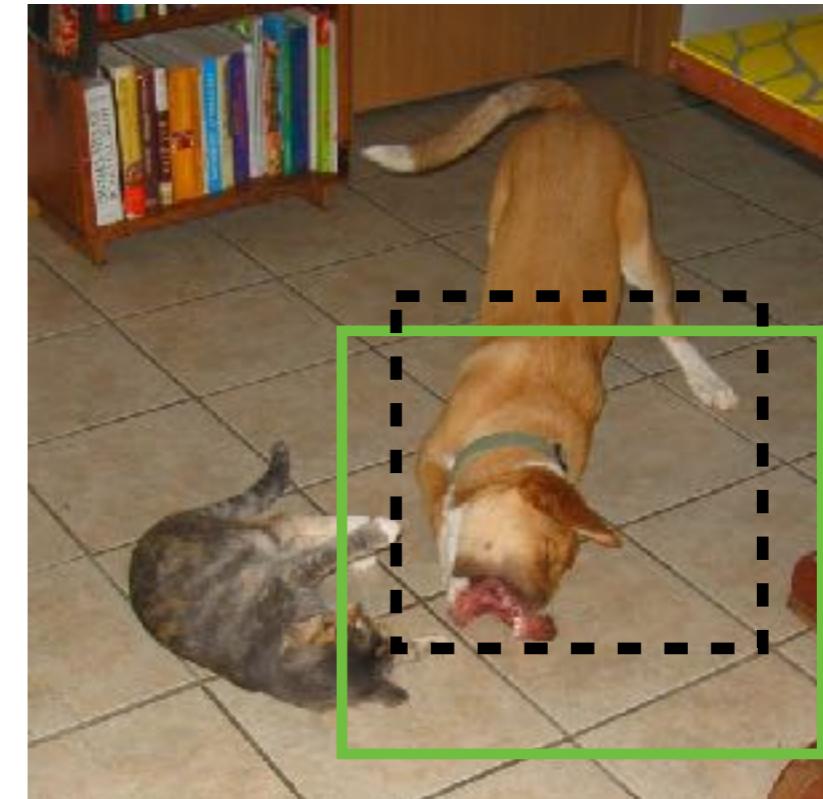
Bounding Box Prediction

Classical sliding windows



Is it a cat? **No**

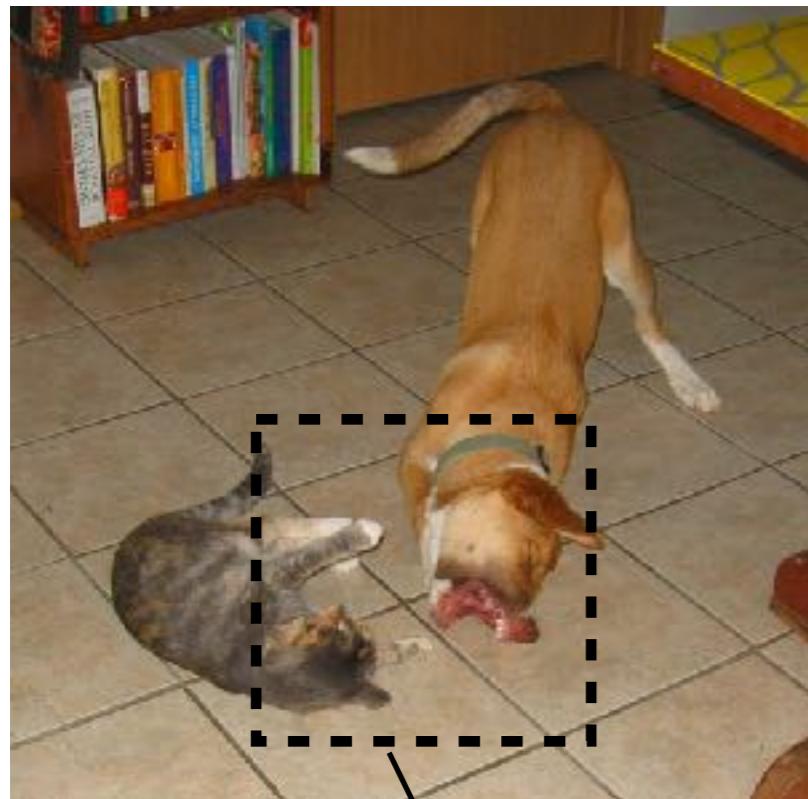
SSD and other deep approaches



Discretize the box space **densely**

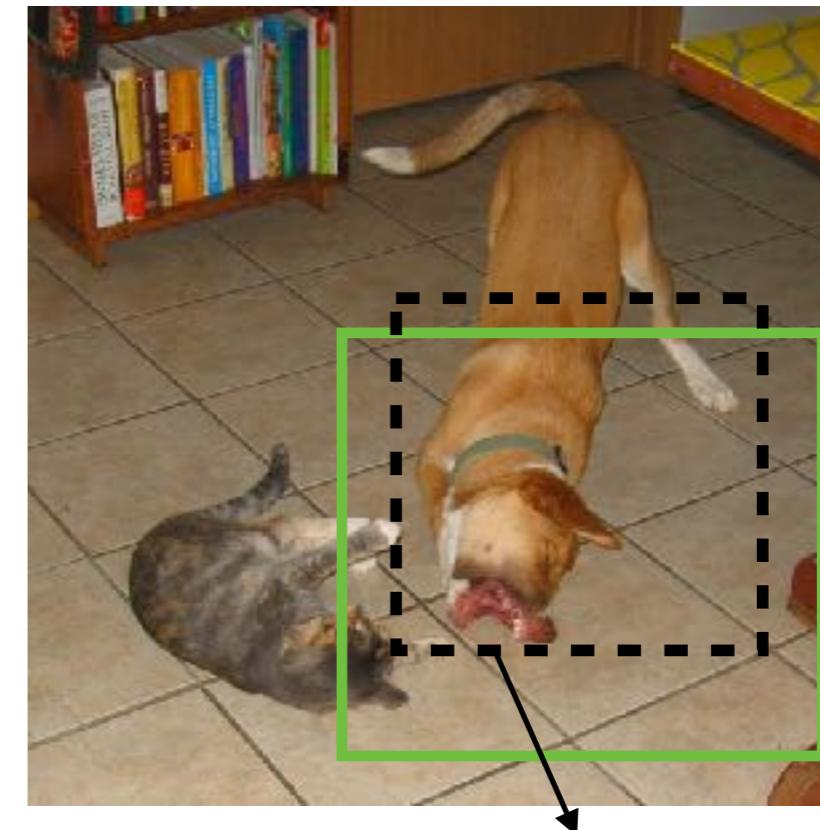
Bounding Box Prediction

Classical sliding windows



Is it a cat? **No**

SSD and other deep approaches

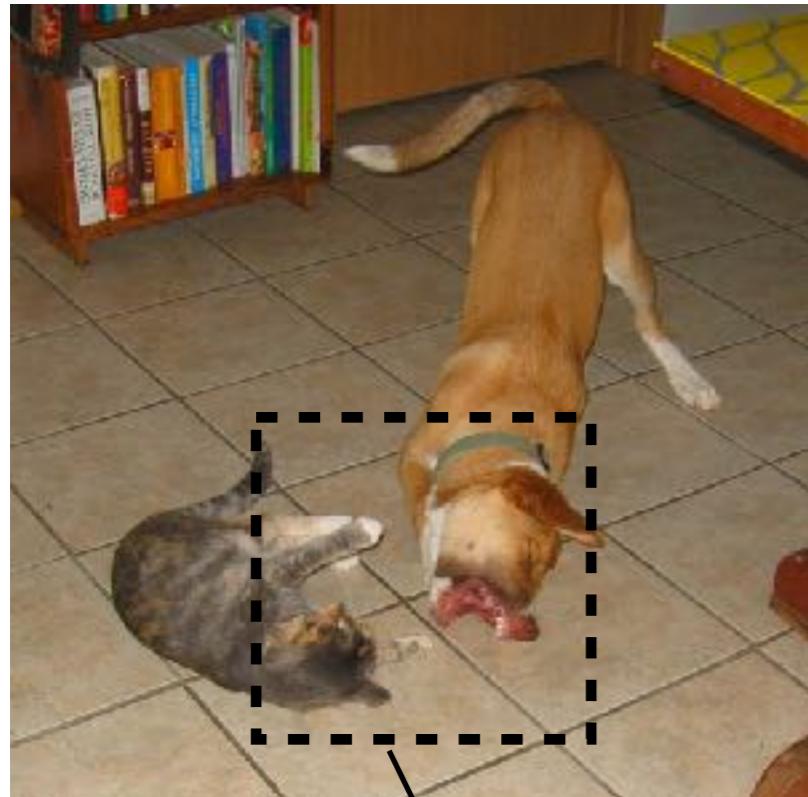


dog: 0.4 cat: 0.2

Discretize the box space **densely**

Bounding Box Prediction

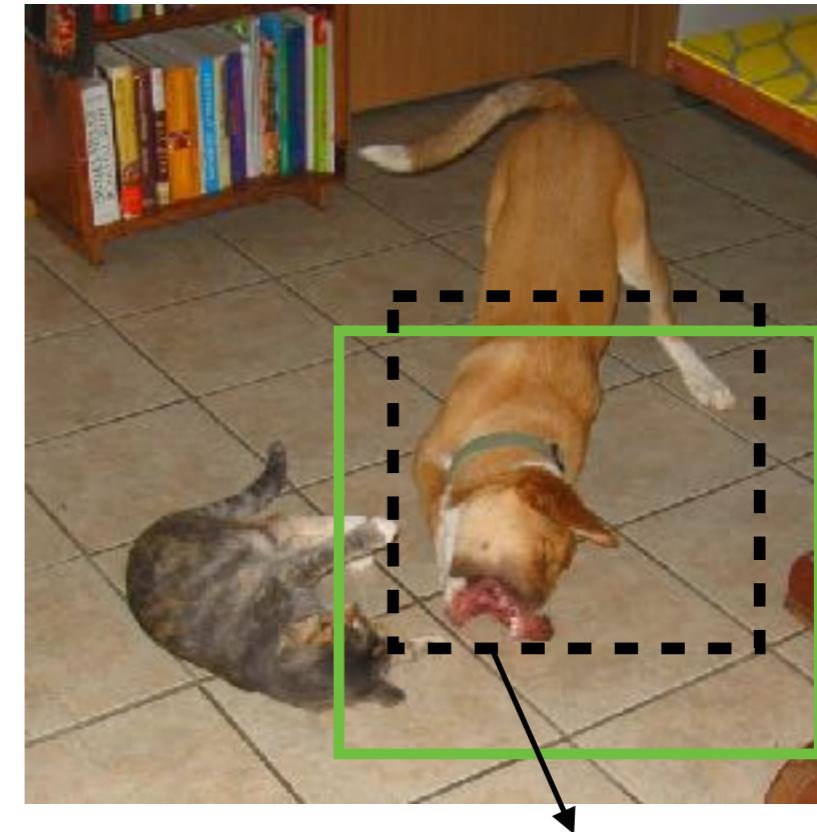
Classical sliding windows



Is it a cat? **No**

Discretize the box space **densely**

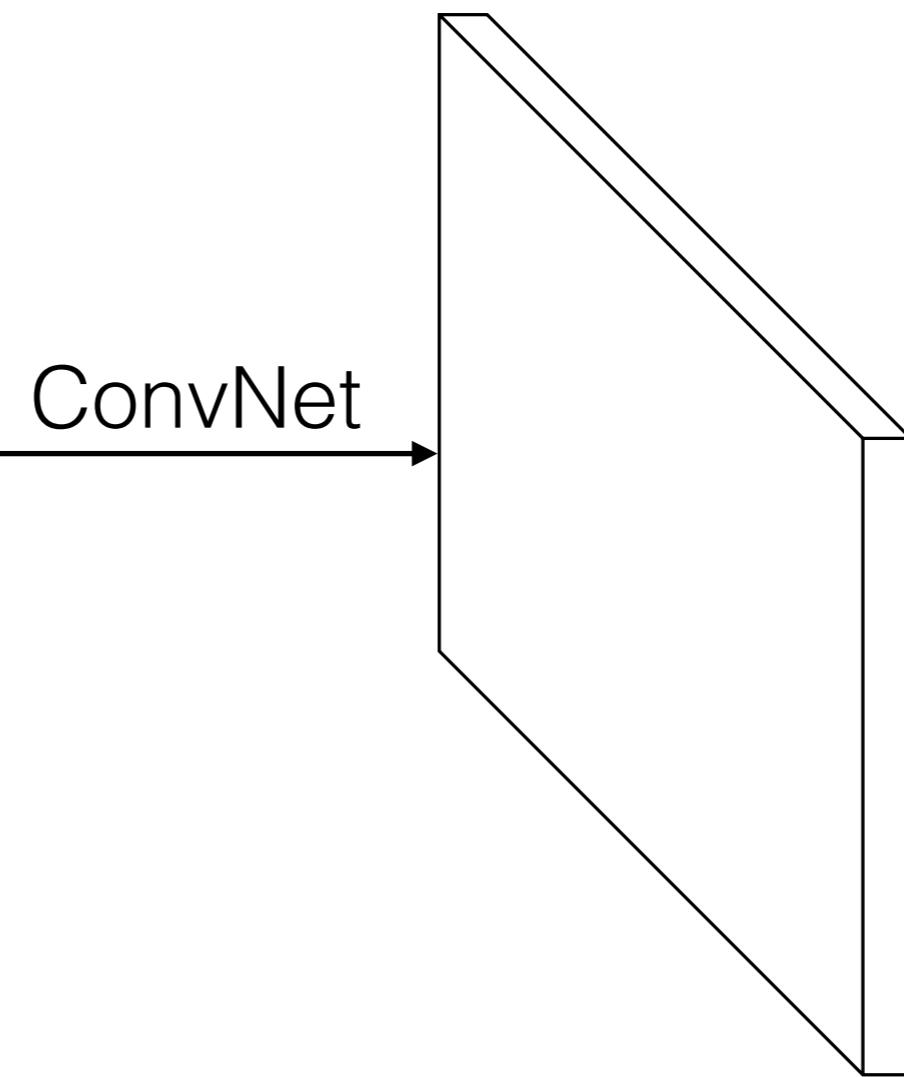
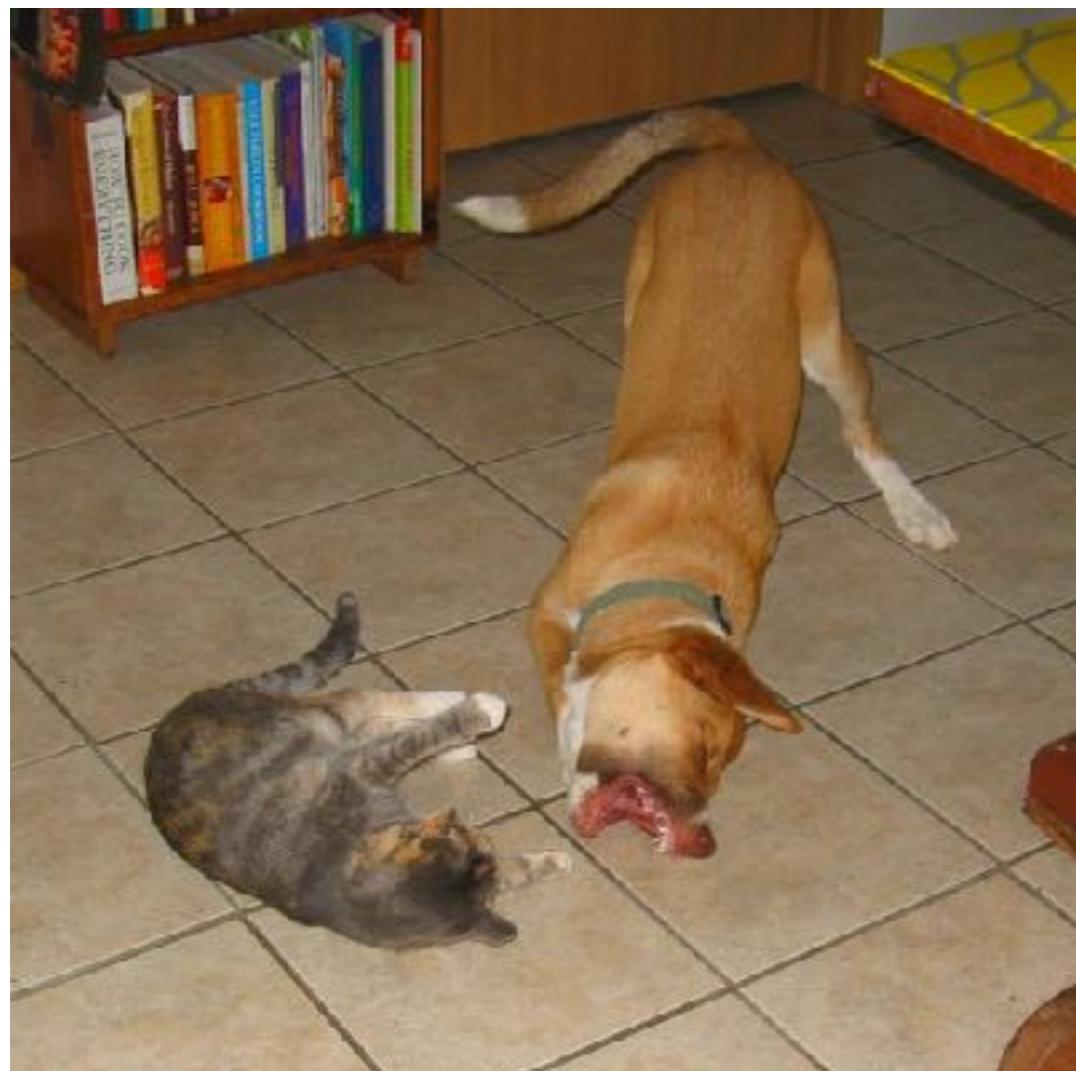
SSD and other deep approaches



dog: 0.4 cat: 0.2

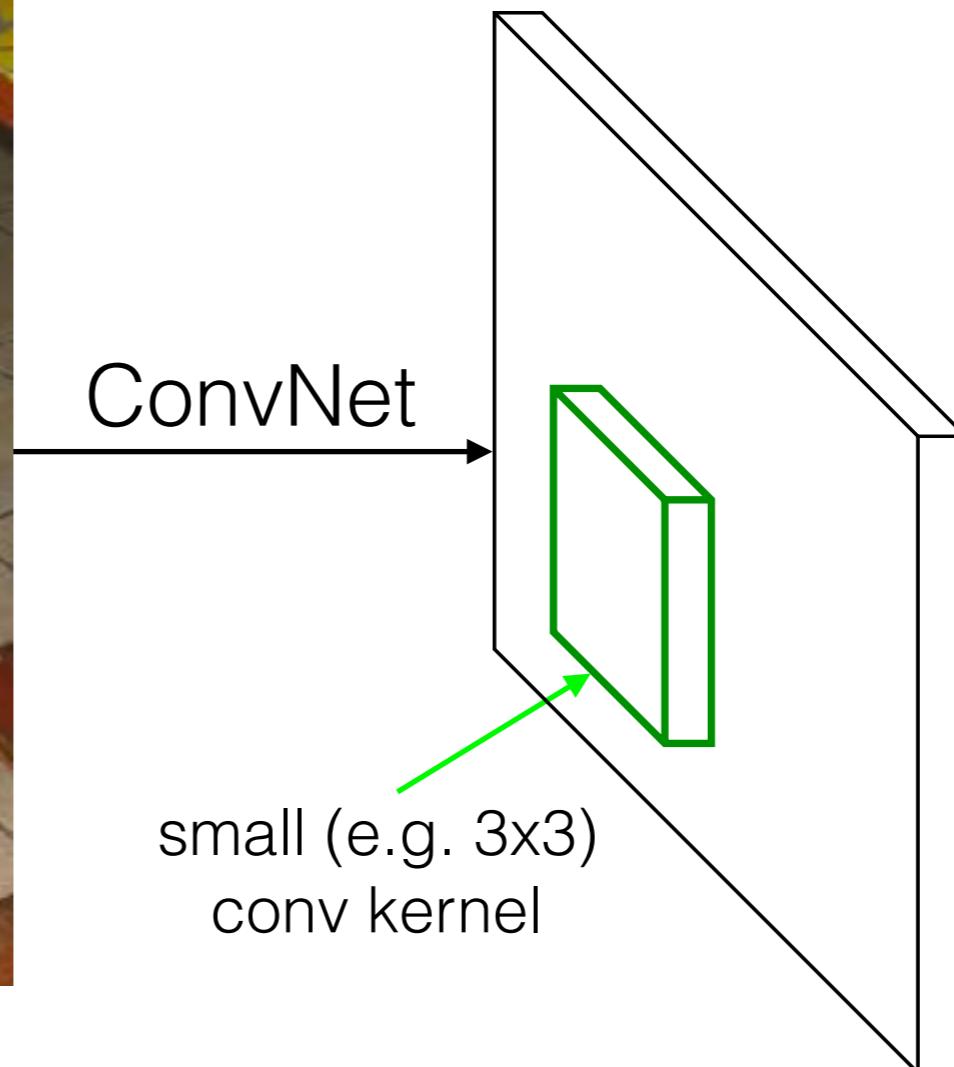
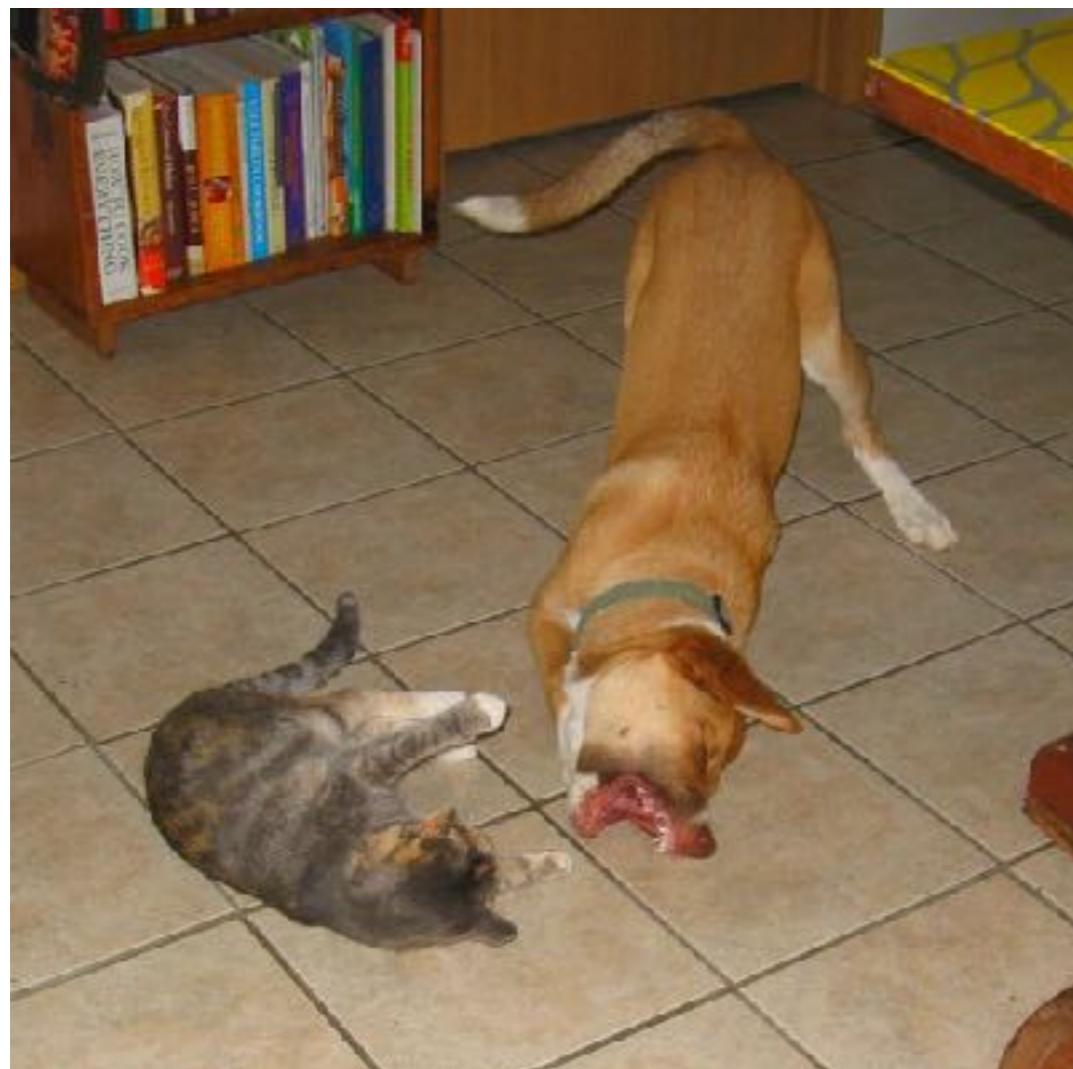
Discretize the box space more **coarsely**
Refine the coordinates of each box

SSD Output Layer



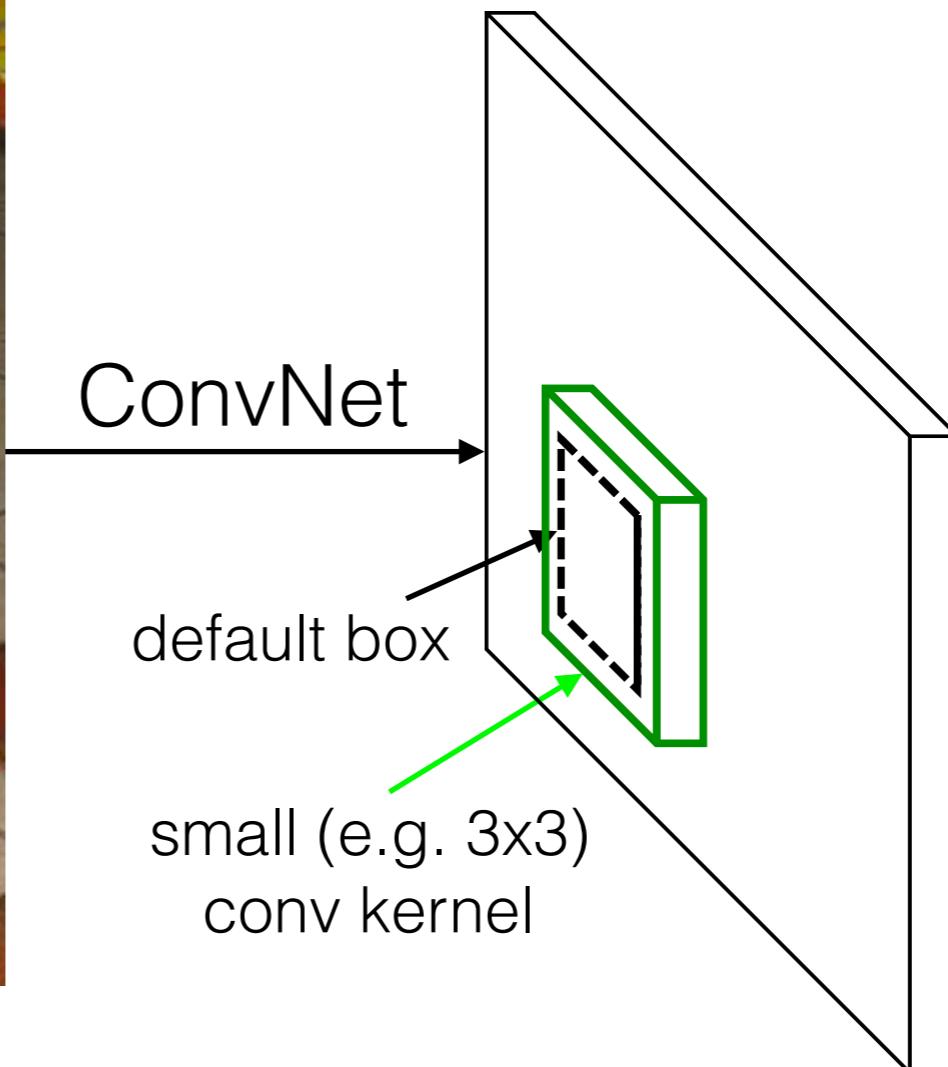
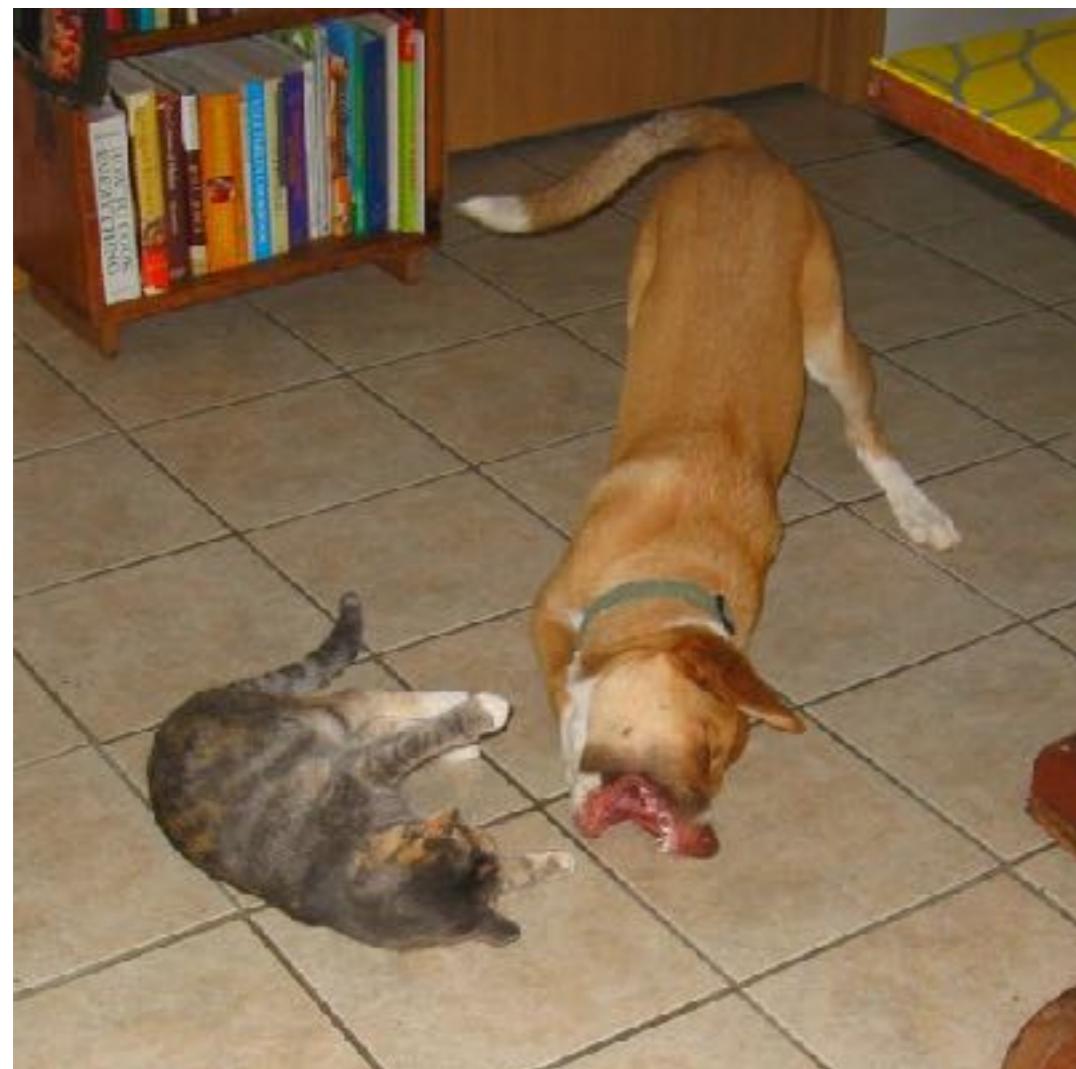
feature map

SSD Output Layer



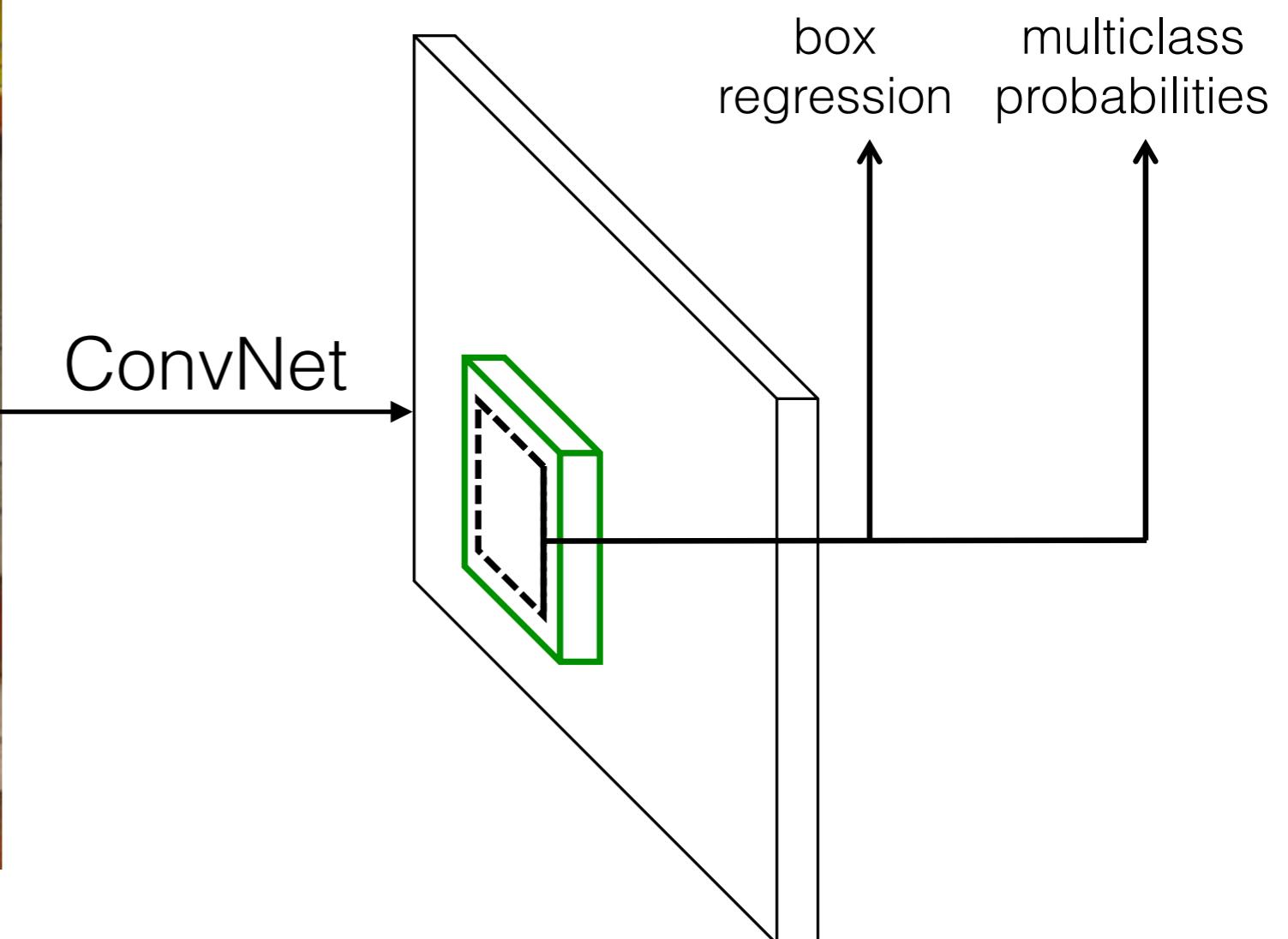
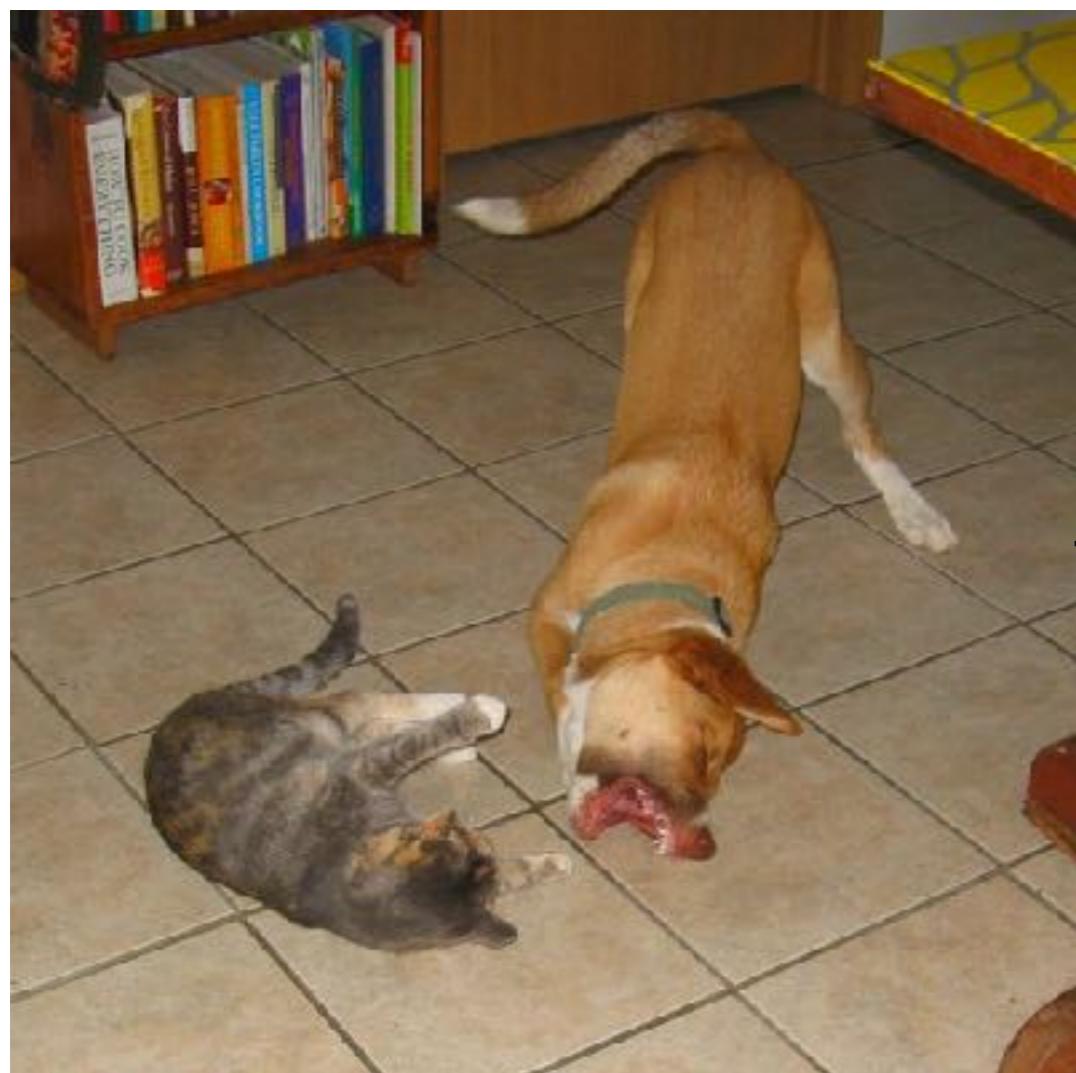
feature map

SSD Output Layer



feature map

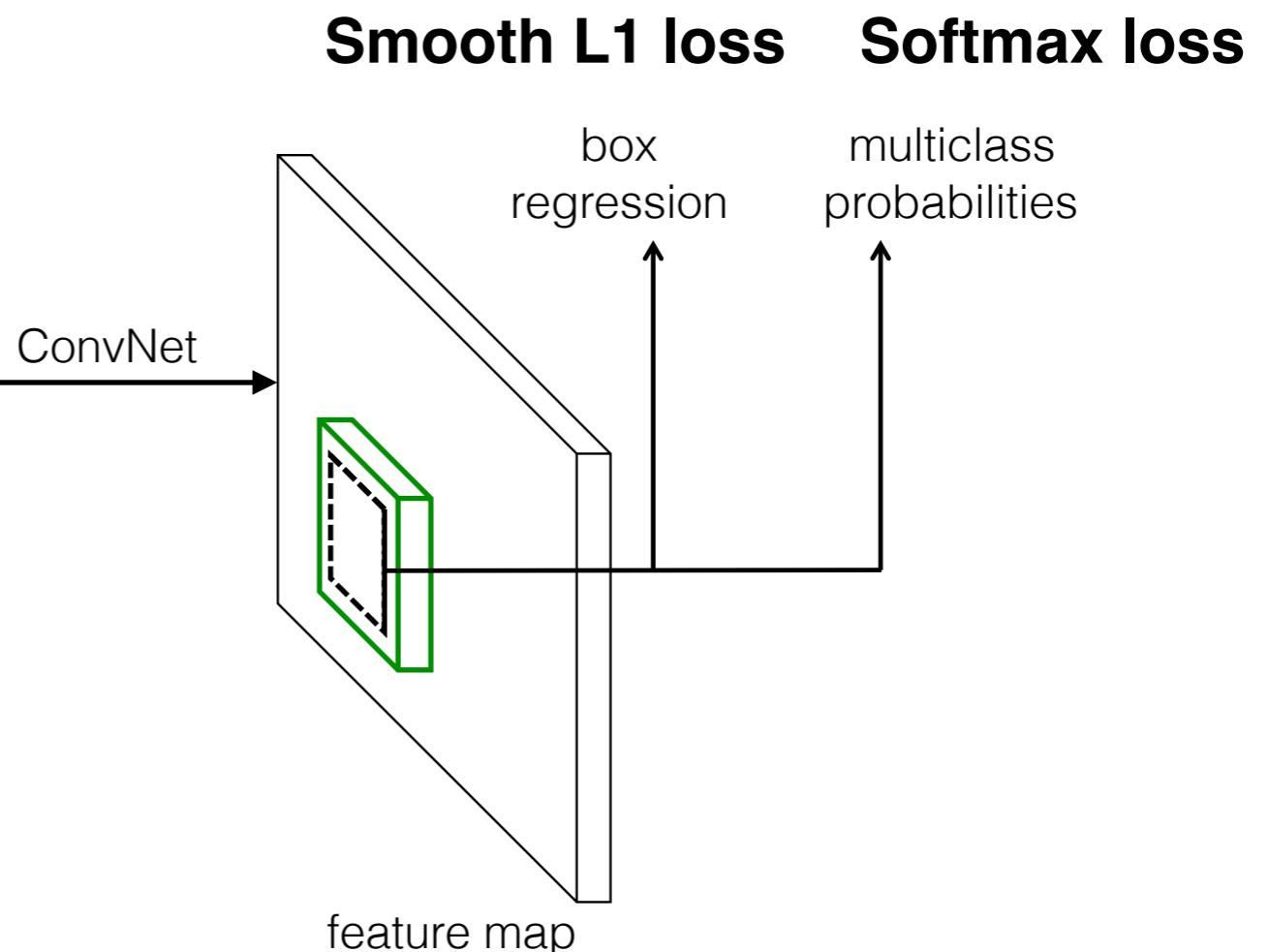
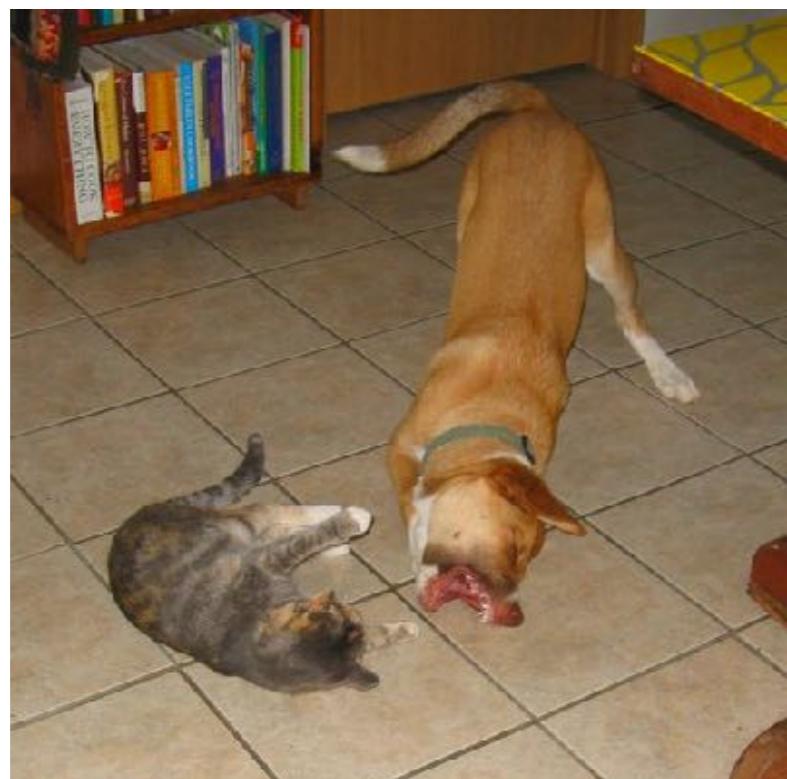
SSD Output Layer



feature map

SSD Training

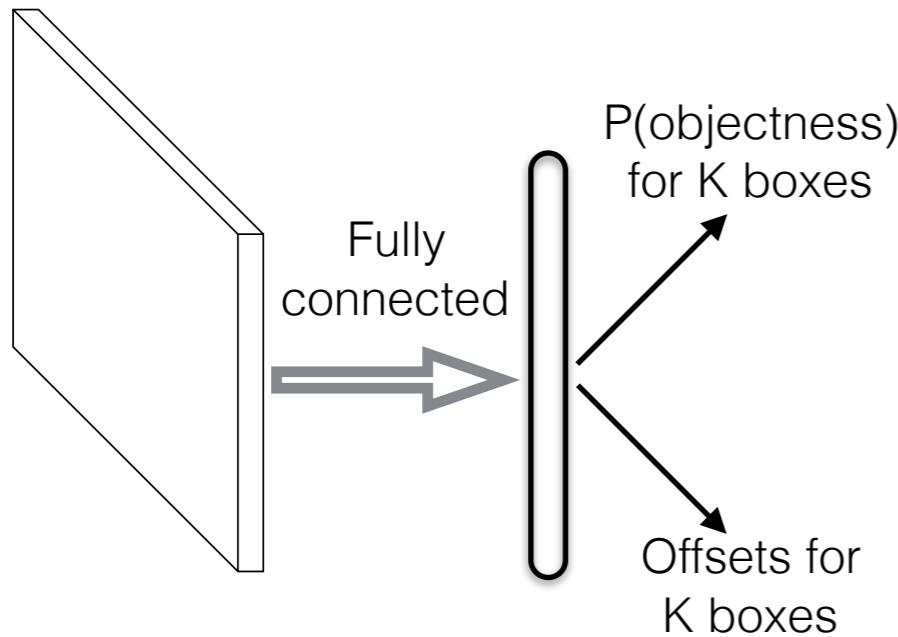
- Match default boxes to ground truth boxes to determine true/false positives.
- Loss = **SmoothL1**(box param) + **Softmax**(class prob)



Related Work

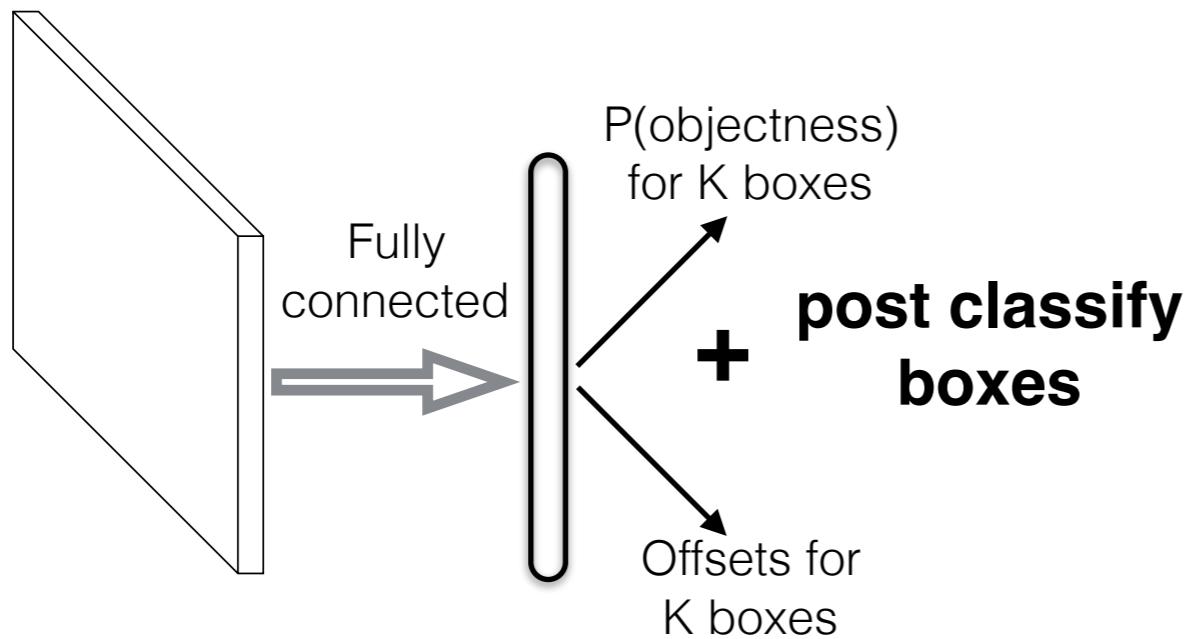
Related Work

MultiBox [Erhan et al. CVPR14]



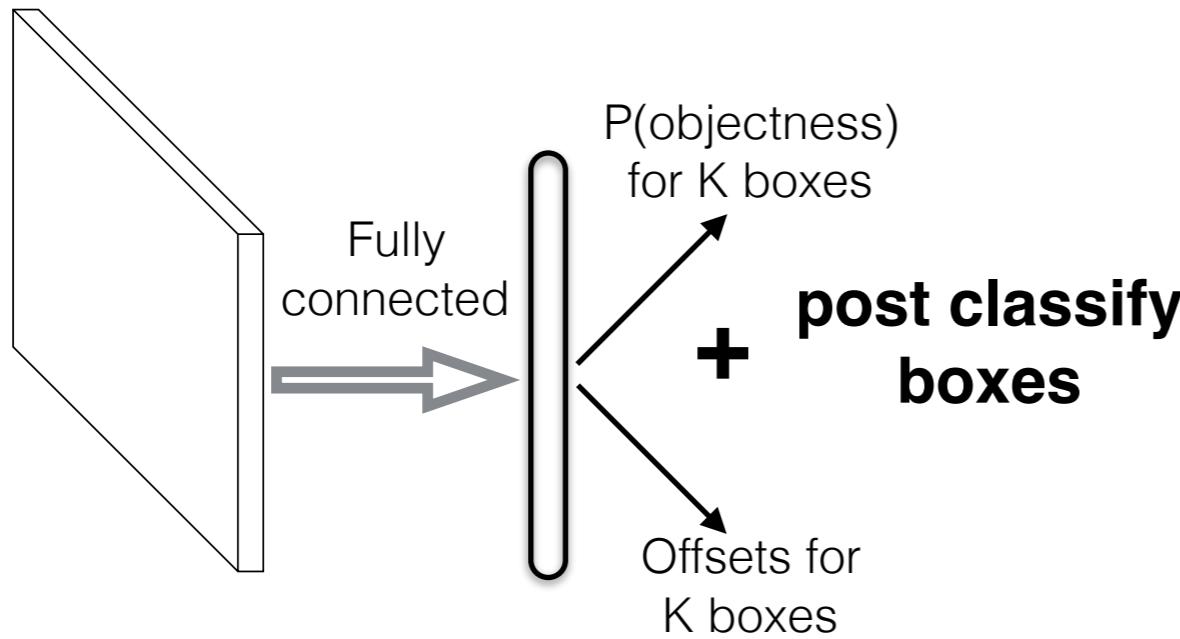
Related Work

MultiBox [Erhan et al. CVPR14]

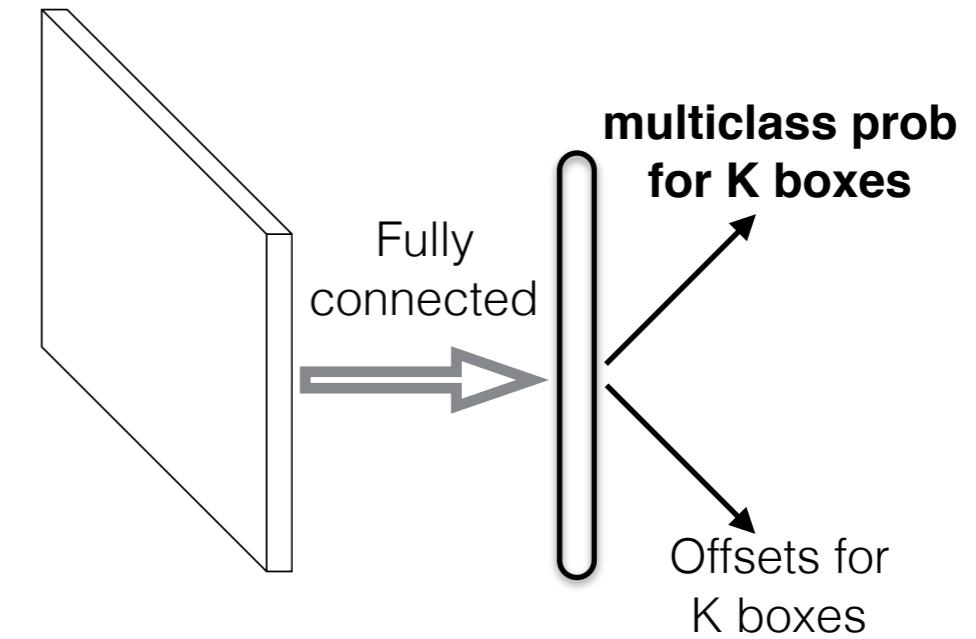


Related Work

MultiBox [Erhan et al. CVPR14]

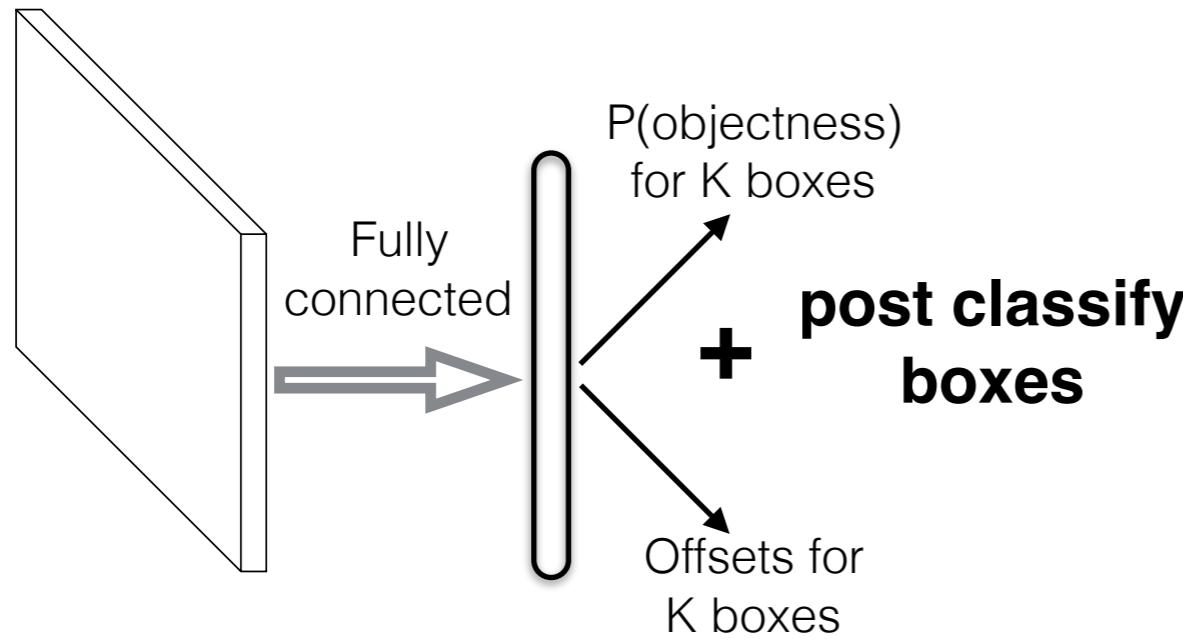


YOLO [Redmon et al. CVPR16]

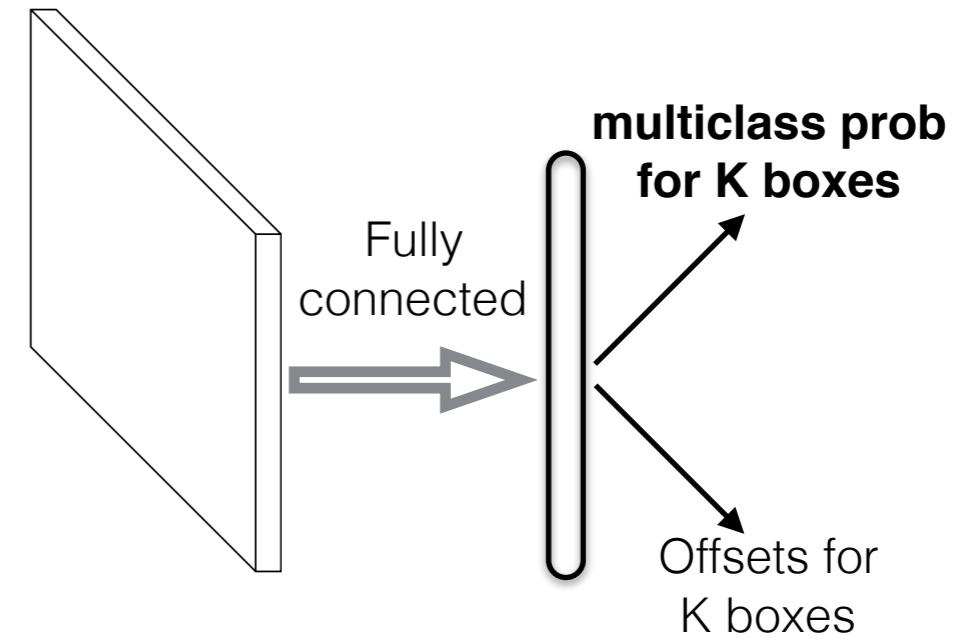


Related Work

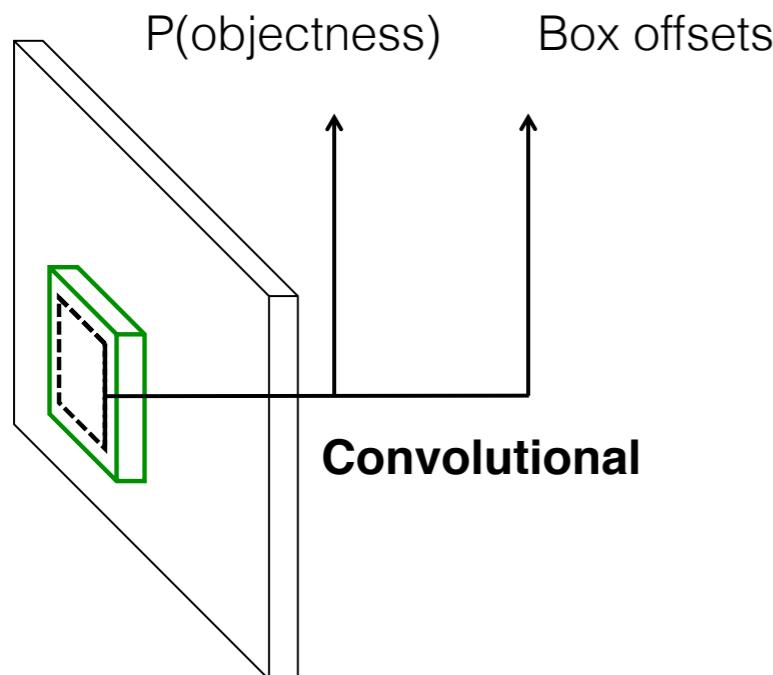
MultiBox [Erhan et al. CVPR14]



YOLO [Redmon et al. CVPR16]

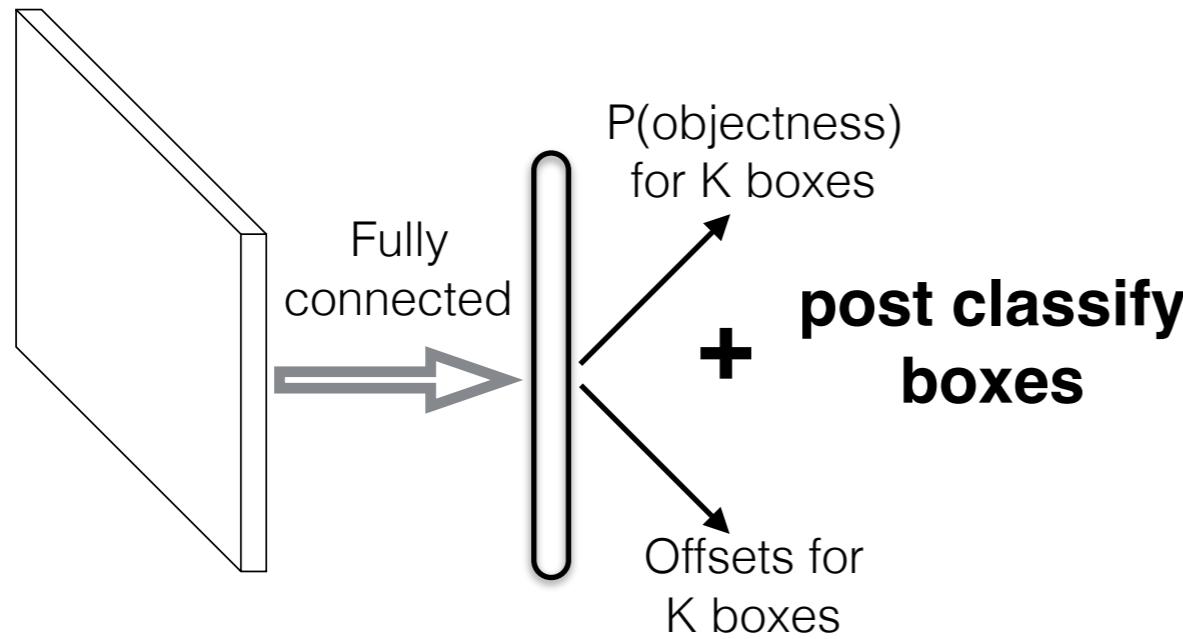


Faster R-CNN [Ren et al. NIPS15]

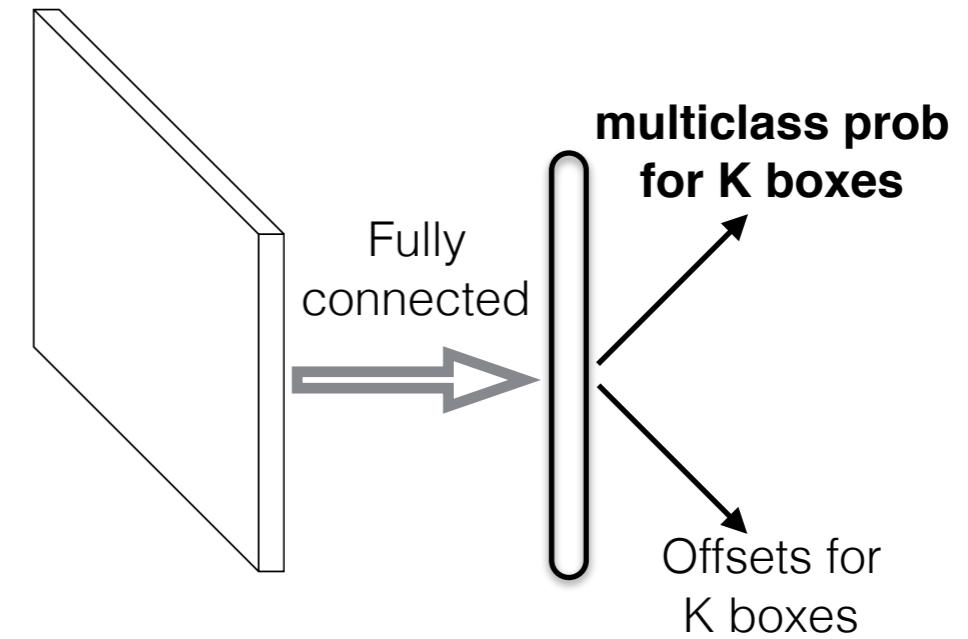


Related Work

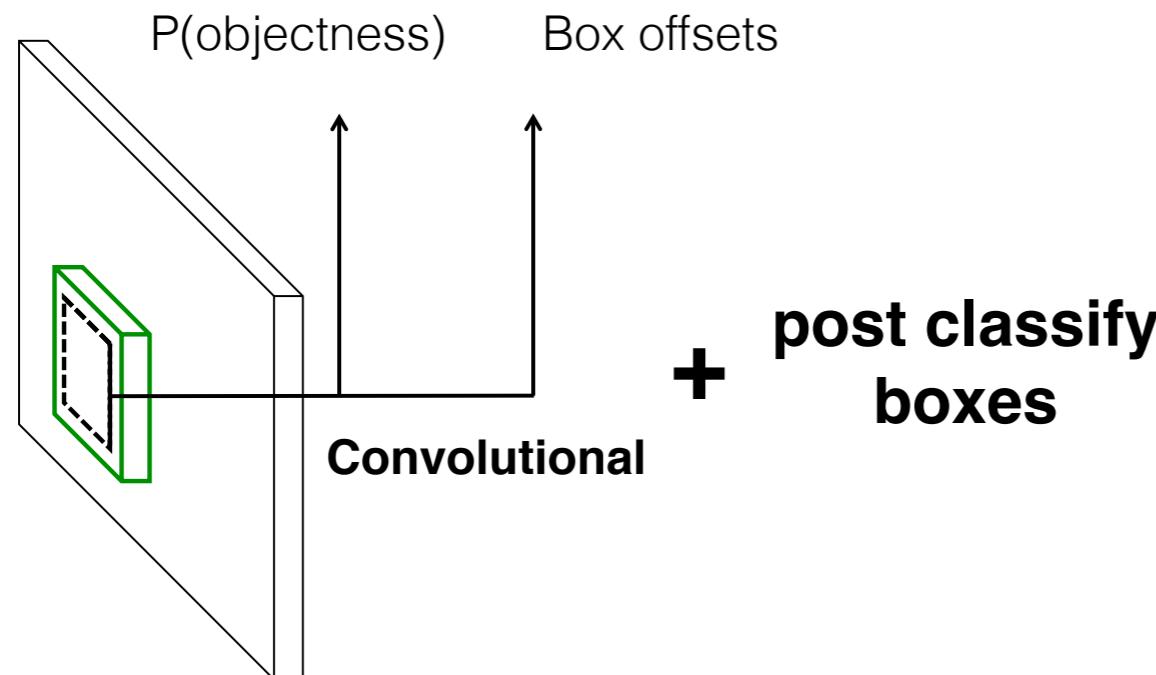
MultiBox [Erhan et al. CVPR14]



YOLO [Redmon et al. CVPR16]

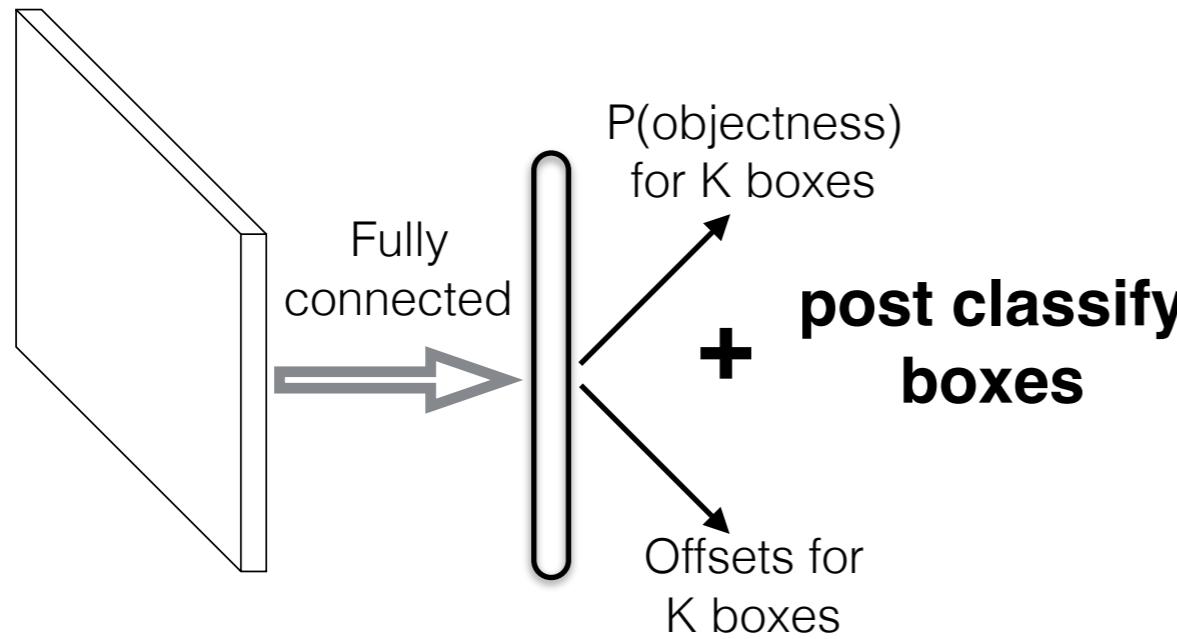


Faster R-CNN [Ren et al. NIPS15]

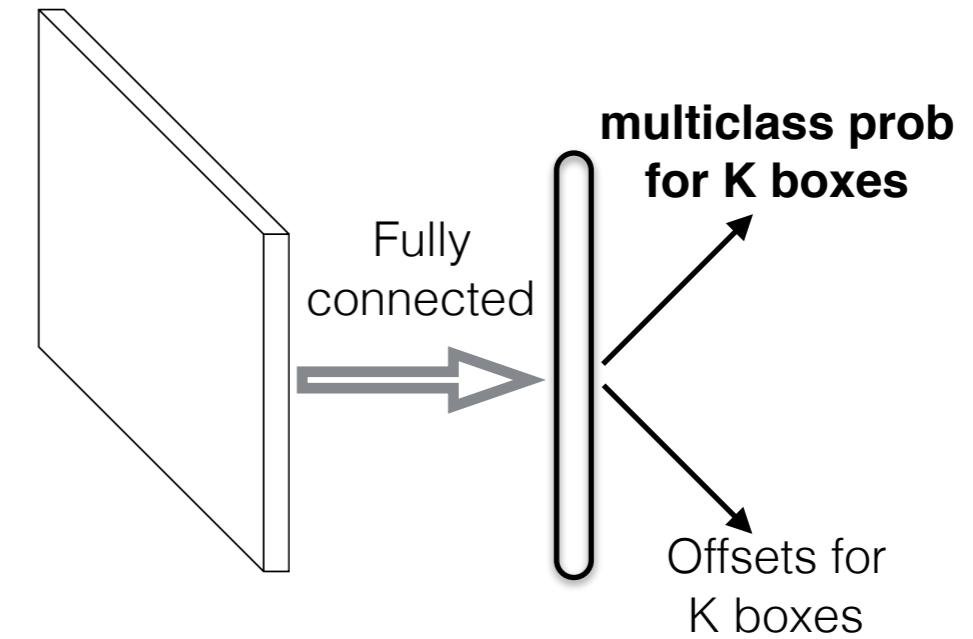


Related Work

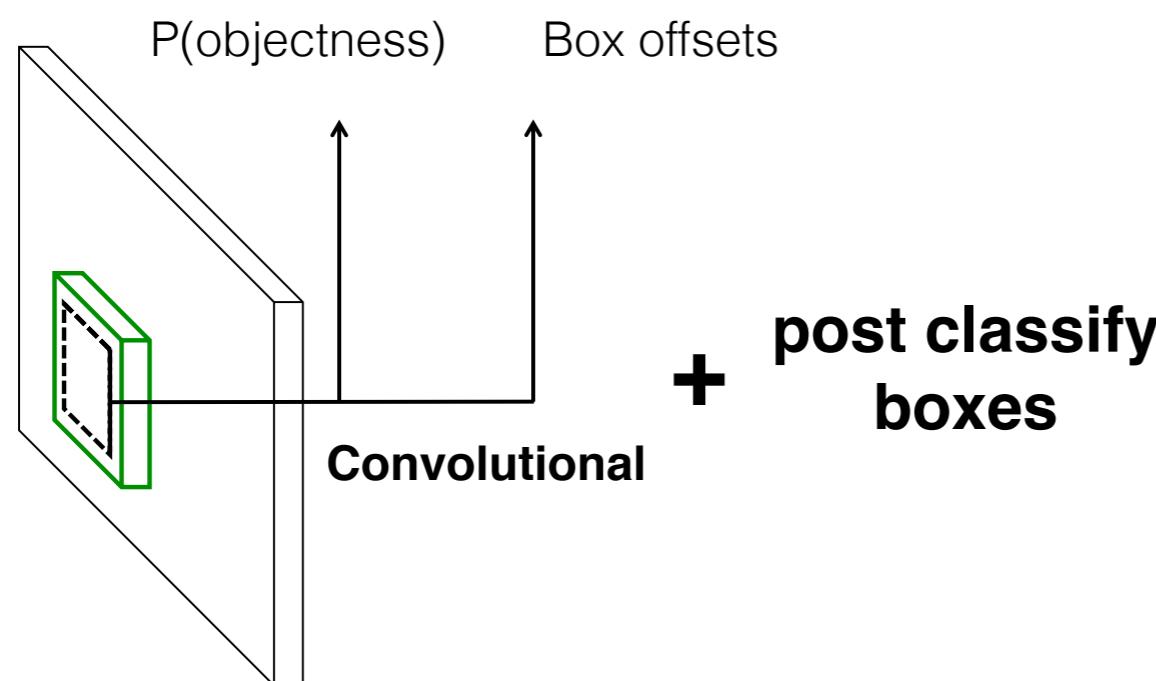
MultiBox [Erhan et al. CVPR14]



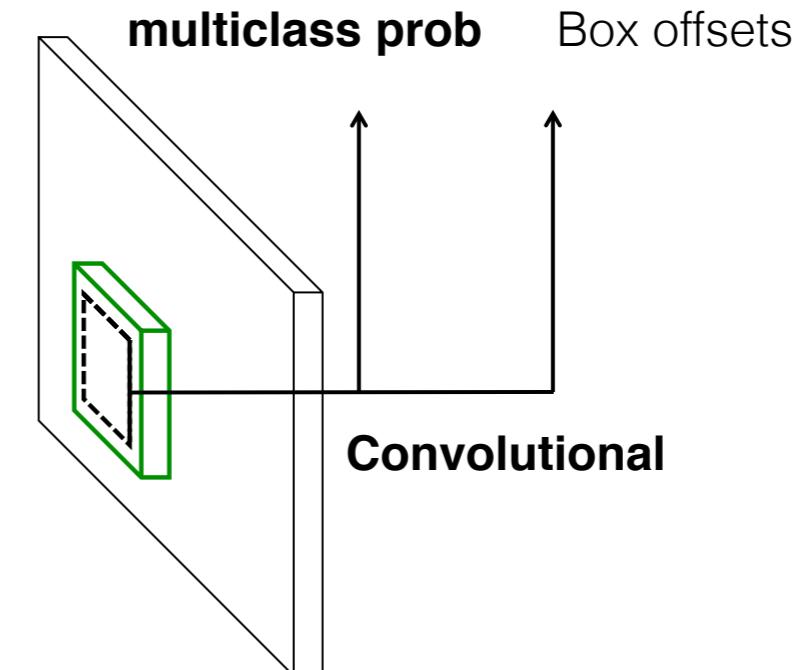
YOLO [Redmon et al. CVPR16]



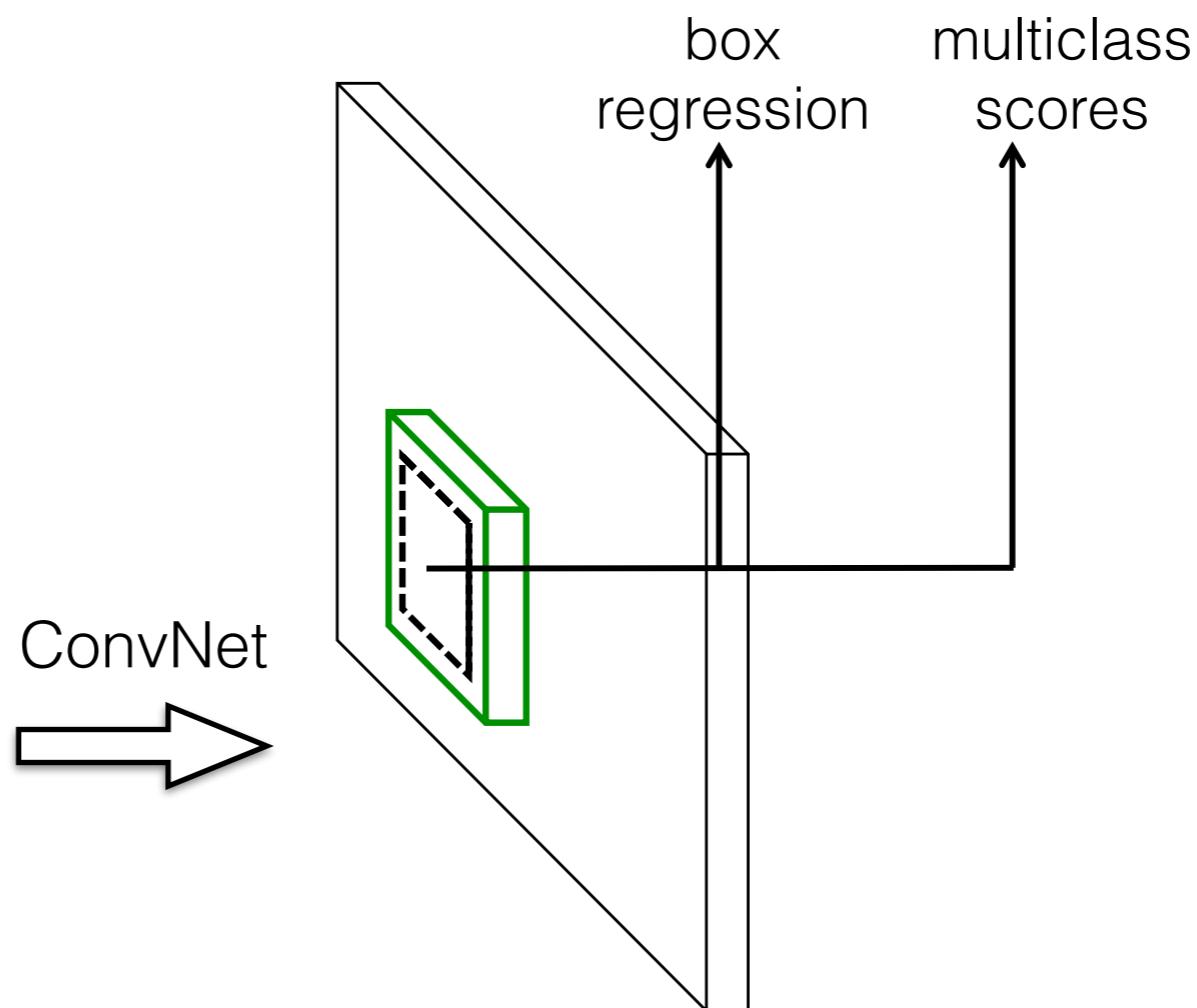
Faster R-CNN [Ren et al. NIPS15]



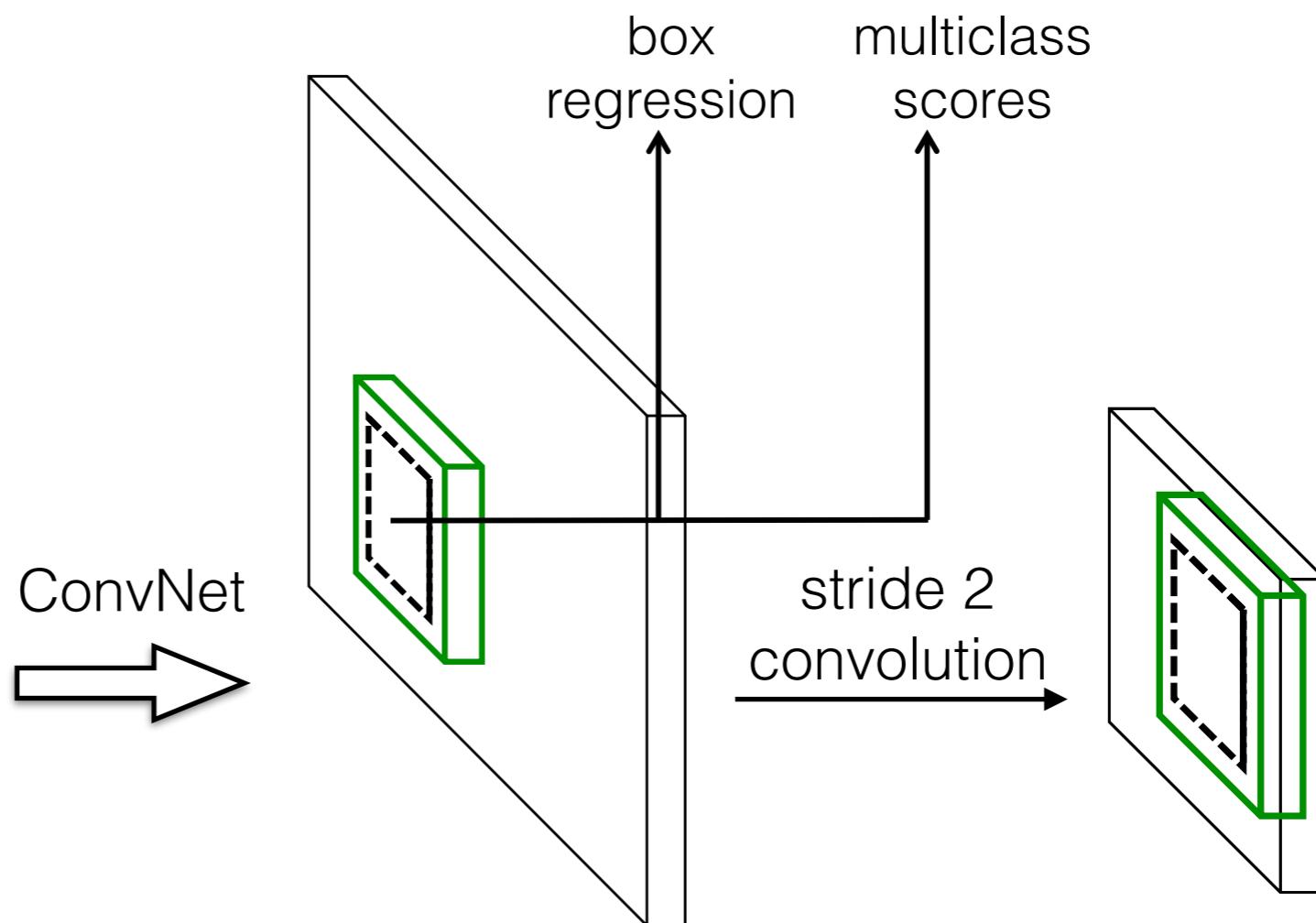
SSD



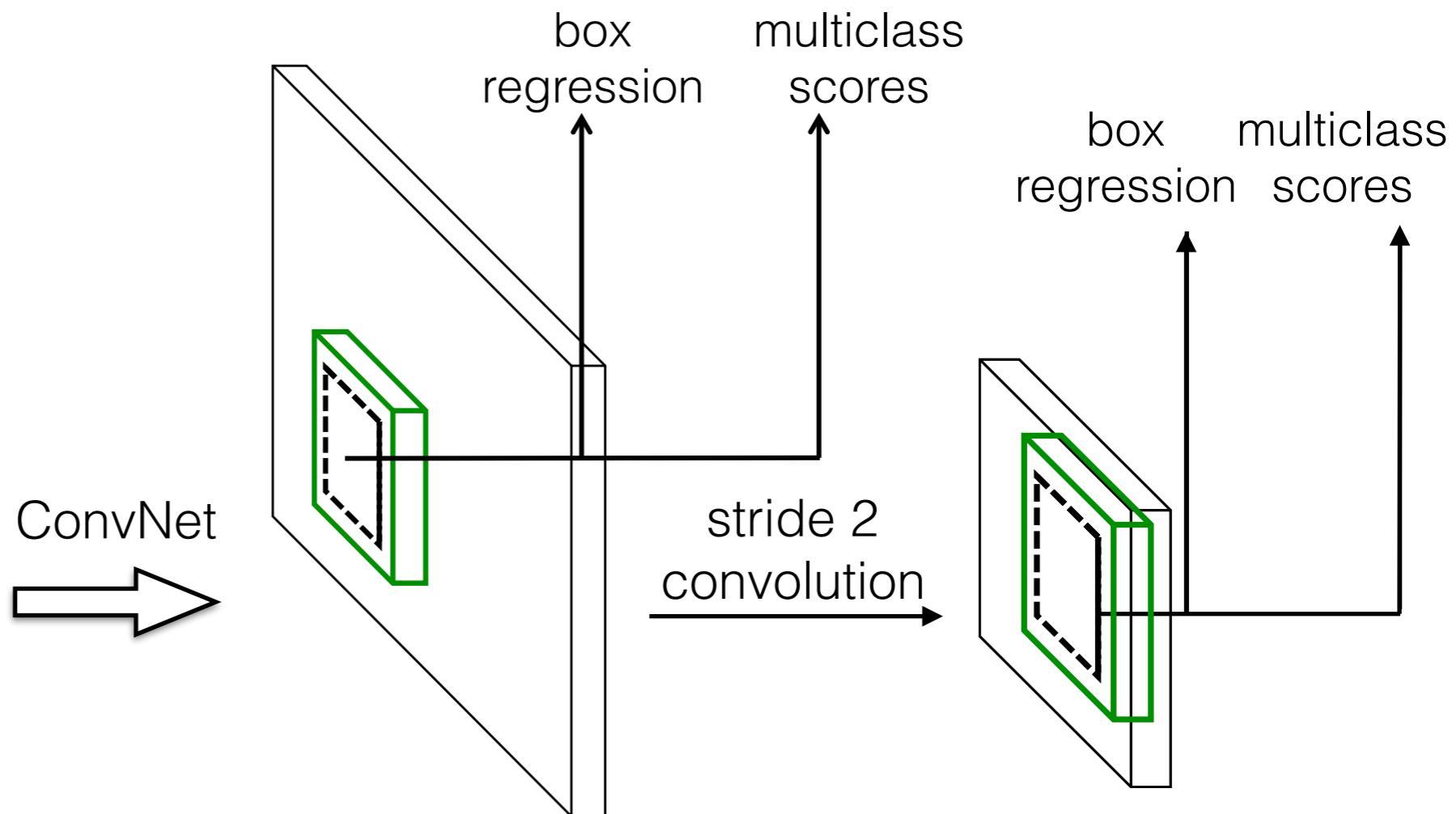
Contribution #1: Multi-Scale Feature Maps



Contribution #1: Multi-Scale Feature Maps

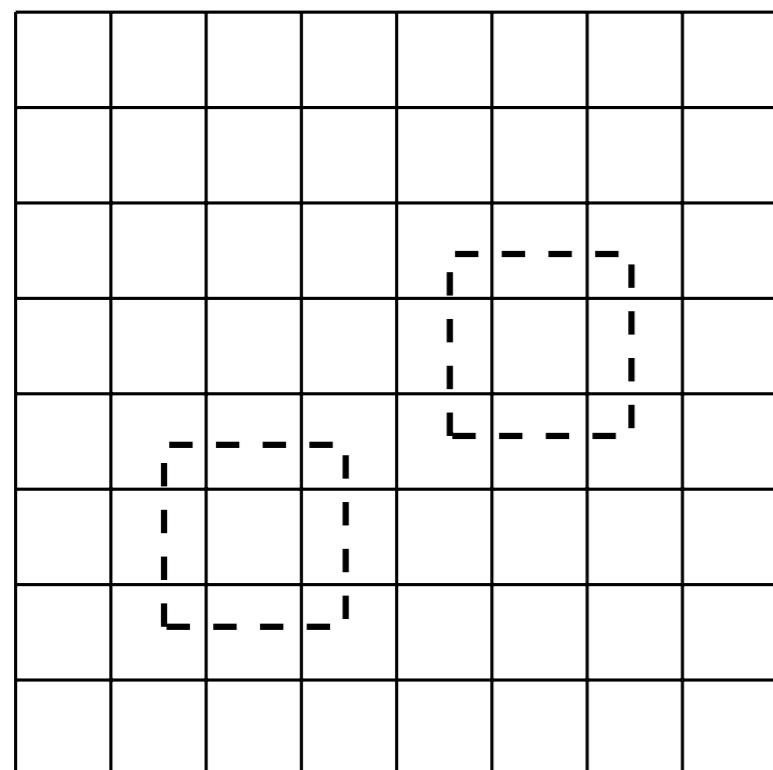


Contribution #1: Multi-Scale Feature Maps

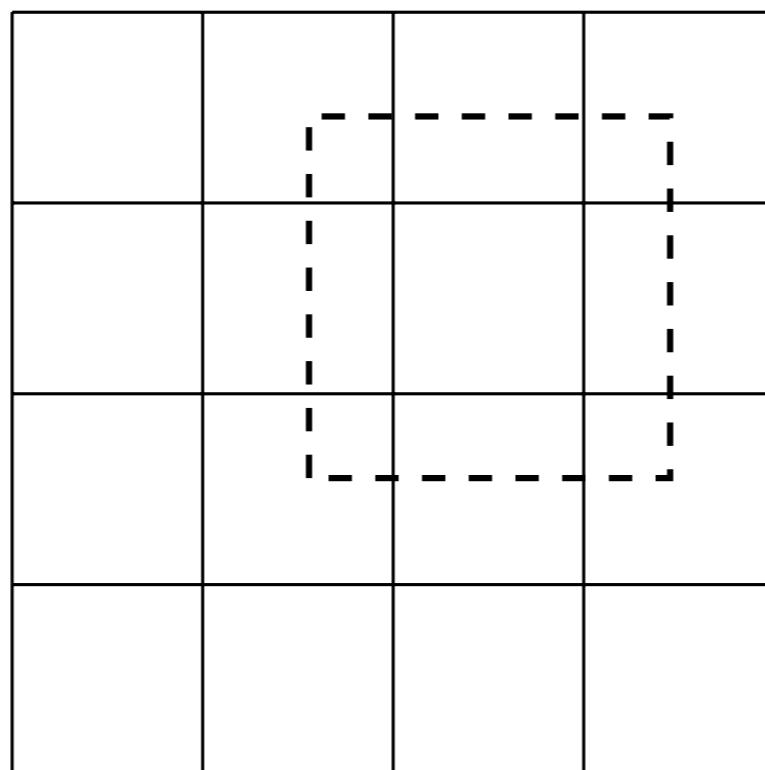


Multi-Scale Feature Maps

SSD



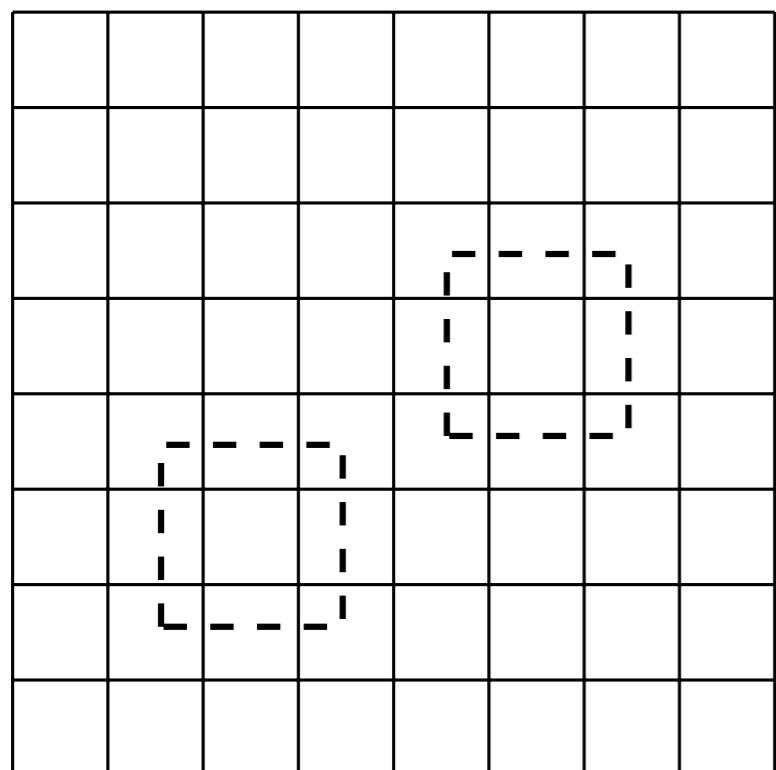
8×8 feature map



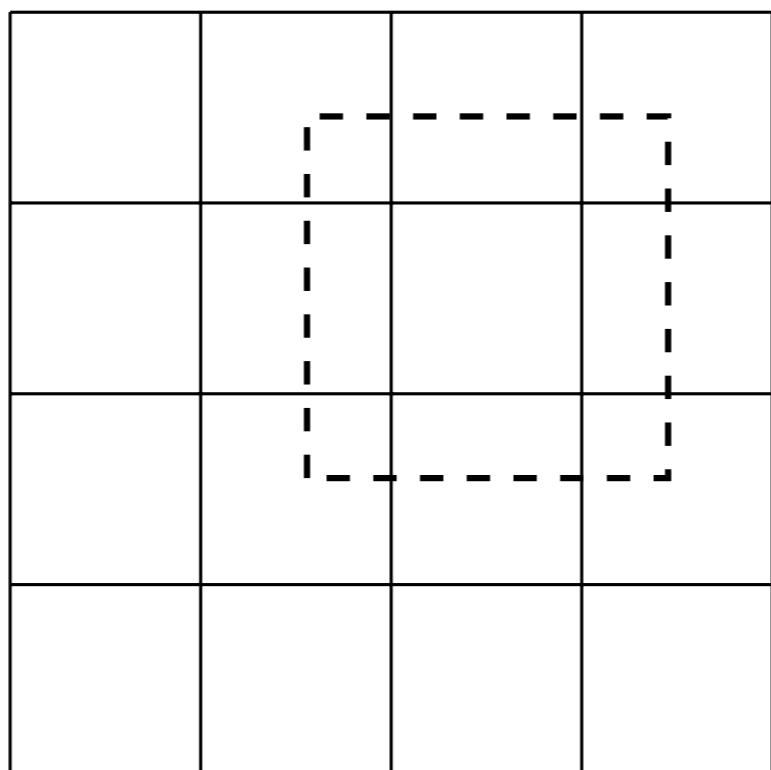
4×4 feature map

Multi-Scale Feature Maps

SSD



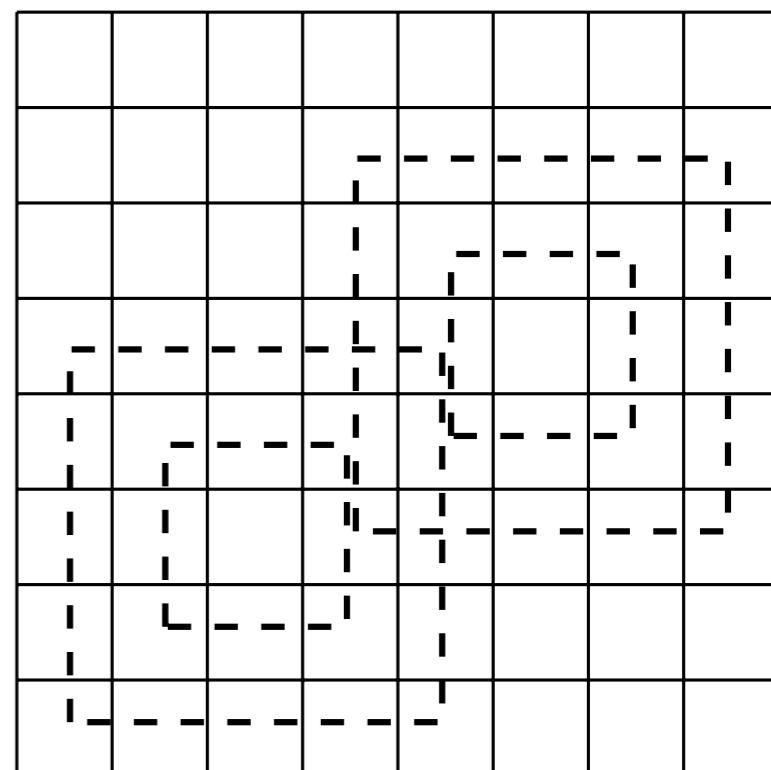
8×8 feature map



4×4 feature map

Faster R-CNN Objectness
Proposal, Ren 2015

vs.



8×8 feature map

Multi-Scale Feature Maps Experiment

Prediction source layers from:						mAP		# Boxes
38 × 38	19 × 19	10 × 10	5 × 5	3 × 3	1 × 1	use boundary boxes?		
✓	✓	✓	✓	✓	✓	Yes	No	
						74.3	63.4	8732
						70.7	69.2	9864
						62.4	64.0	8664

Multi-Scale Feature Maps Experiment

Prediction source layers from:						use boundary boxes?	mAP	# Boxes
38 × 38	19 × 19	10 × 10	5 × 5	3 × 3	1 × 1			
✓	✓	✓	✓	✓	✓	Yes	74.3	8732
✓	✓	✓				No	63.4	
✓						70.7	69.2	9864
						62.4	64.0	8664

Multi-Scale Feature Maps Experiment

Prediction source layers from:						mAP		# Boxes
38 × 38	19 × 19	10 × 10	5 × 5	3 × 3	1 × 1	use boundary boxes?		
✓	✓	✓	✓	✓	✓	Yes	No	
						74.3	63.4	8732
						70.7	69.2	9864
						62.4	64.0	8664

Multi-Scale Feature Maps Experiment

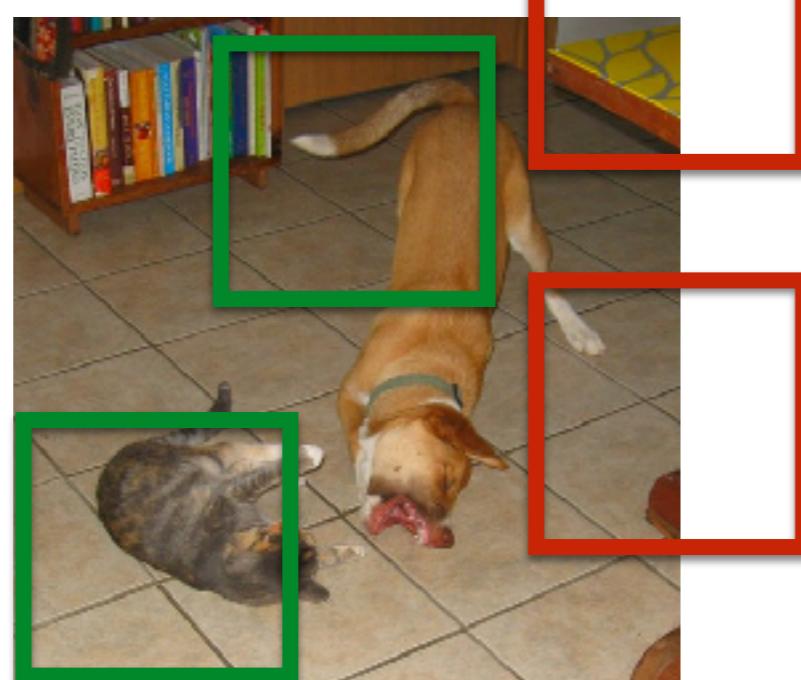
Prediction source layers from:						mAP		# Boxes
38 × 38	19 × 19	10 × 10	5 × 5	3 × 3	1 × 1	use boundary boxes?		
✓	✓	✓	✓	✓	✓	Yes	No	
						74.3	63.4	8732
						70.7	69.2	9864
						62.4	64.0	8664

Multi-Scale Feature Maps Experiment

Prediction source layers from:						mAP		# Boxes
38 × 38	19 × 19	10 × 10	5 × 5	3 × 3	1 × 1	use boundary boxes?		
✓	✓	✓	✓	✓	✓	Yes	No	
						74.3	63.4	8732
						70.7	69.2	9864
						62.4	64.0	8664

Multi-Scale Feature Maps Experiment

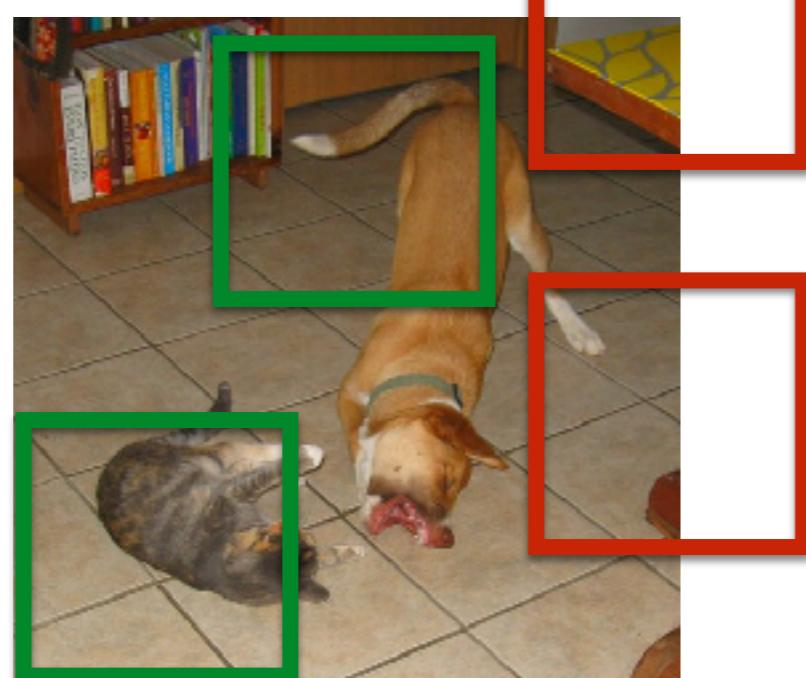
Prediction source layers from:						mAP		# Boxes
38×38	19×19	10×10	5×5	3×3	1×1	use boundary boxes?		
✓	✓	✓	✓	✓	✓	Yes	No	
✓	✓	✓	✓	✓	✓	74.3	63.4	8732
✓	✓	✓				70.7	69.2	9864
✓						62.4	64.0	8664



Original slides are from http://www.cs.unc.edu/~wliu/papers/ssd_eccv2016_slide.pdf

Multi-Scale Feature Maps Experiment

Prediction source layers from:						mAP		# Boxes
38×38	19×19	10×10	5×5	3×3	1×1	use boundary boxes?		
✓	✓	✓	✓	✓	✓	Yes	No	
						74.3	63.4	8732
						70.7	69.2	9864
						62.4	64.0	8664



boundary boxes

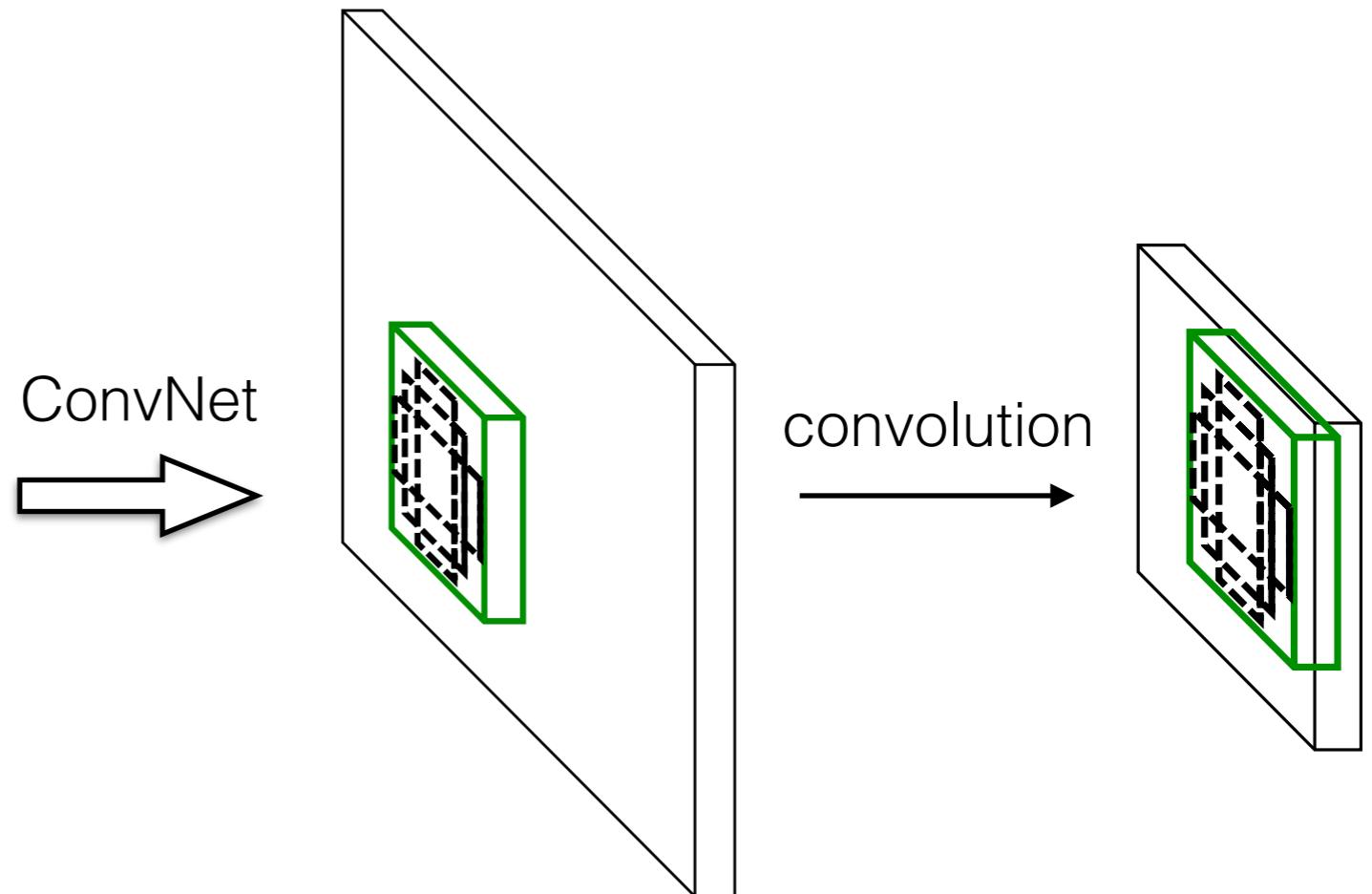
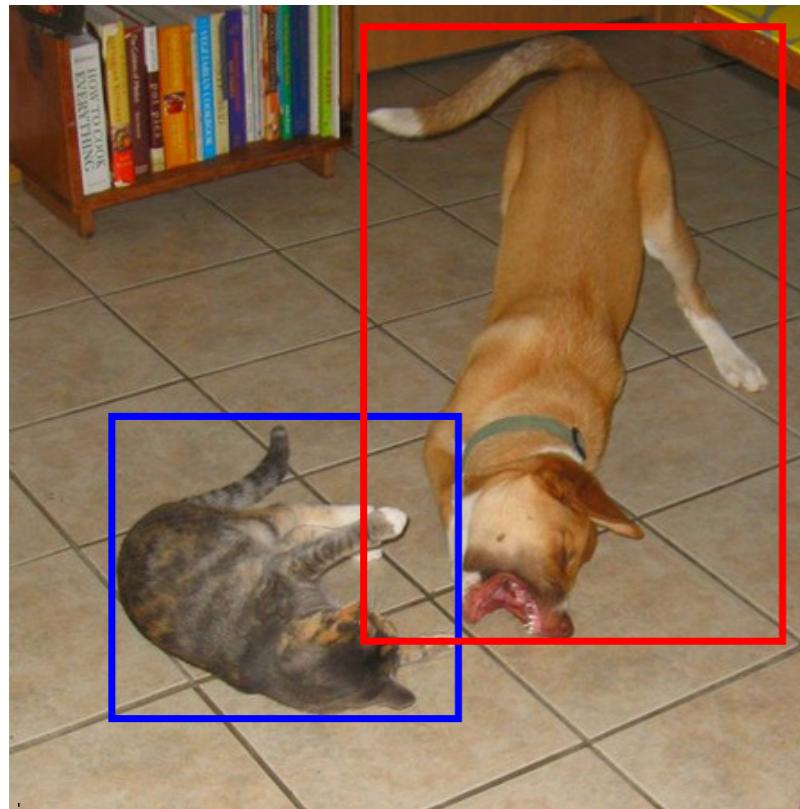
Multi-Scale Feature Maps Experiment

Prediction source layers from:						use boundary boxes?	mAP	# Boxes
38 × 38	19 × 19	10 × 10	5 × 5	3 × 3	1 × 1			
✓	✓	✓	✓	✓	✓	Yes	74.3	63.4
✓	✓	✓				No	70.7	69.2
	✓						62.4	64.0
								8732
								9864
								8664

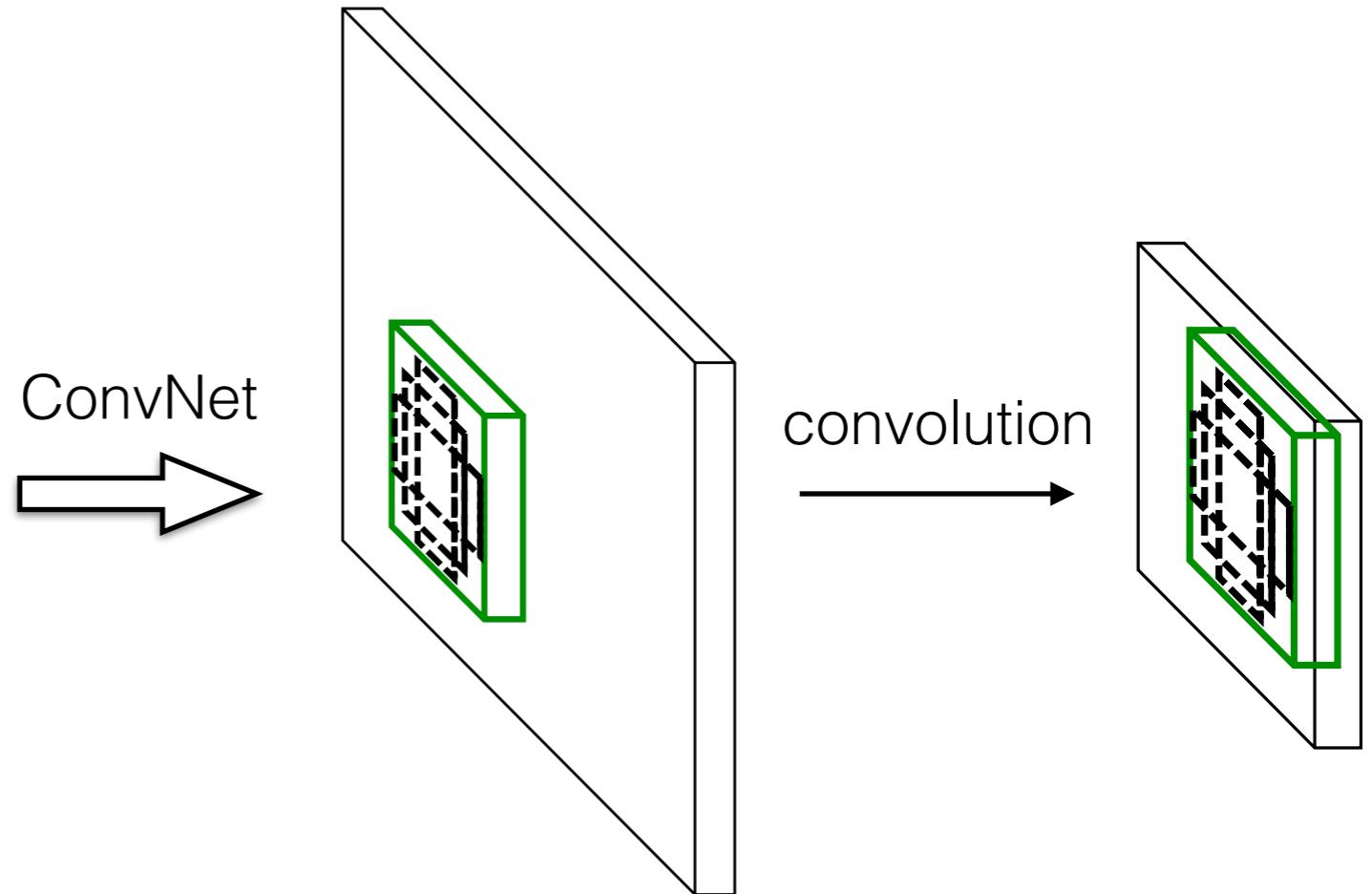
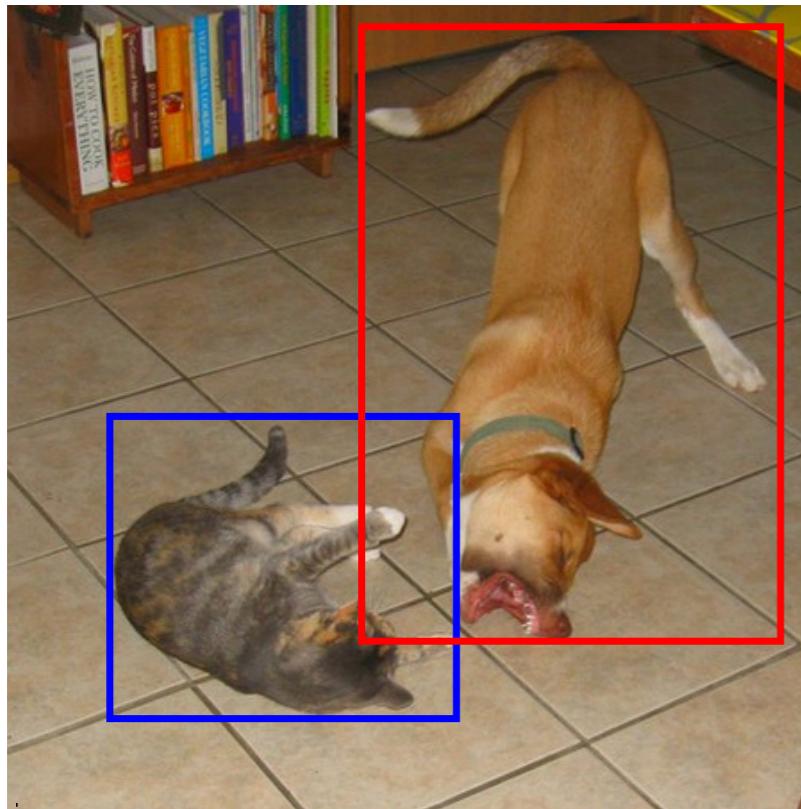
Multi-Scale Feature Maps Experiment

Prediction source layers from:						use boundary boxes?	mAP	# Boxes
38 × 38	19 × 19	10 × 10	5 × 5	3 × 3	1 × 1			
✓	✓	✓	✓	✓	✓	Yes	74.3	8732
✓	✓	✓				No	63.4	
						70.7	69.2	9864
						62.4	64.0	8664

Contribution #2: Splitting the Region Space

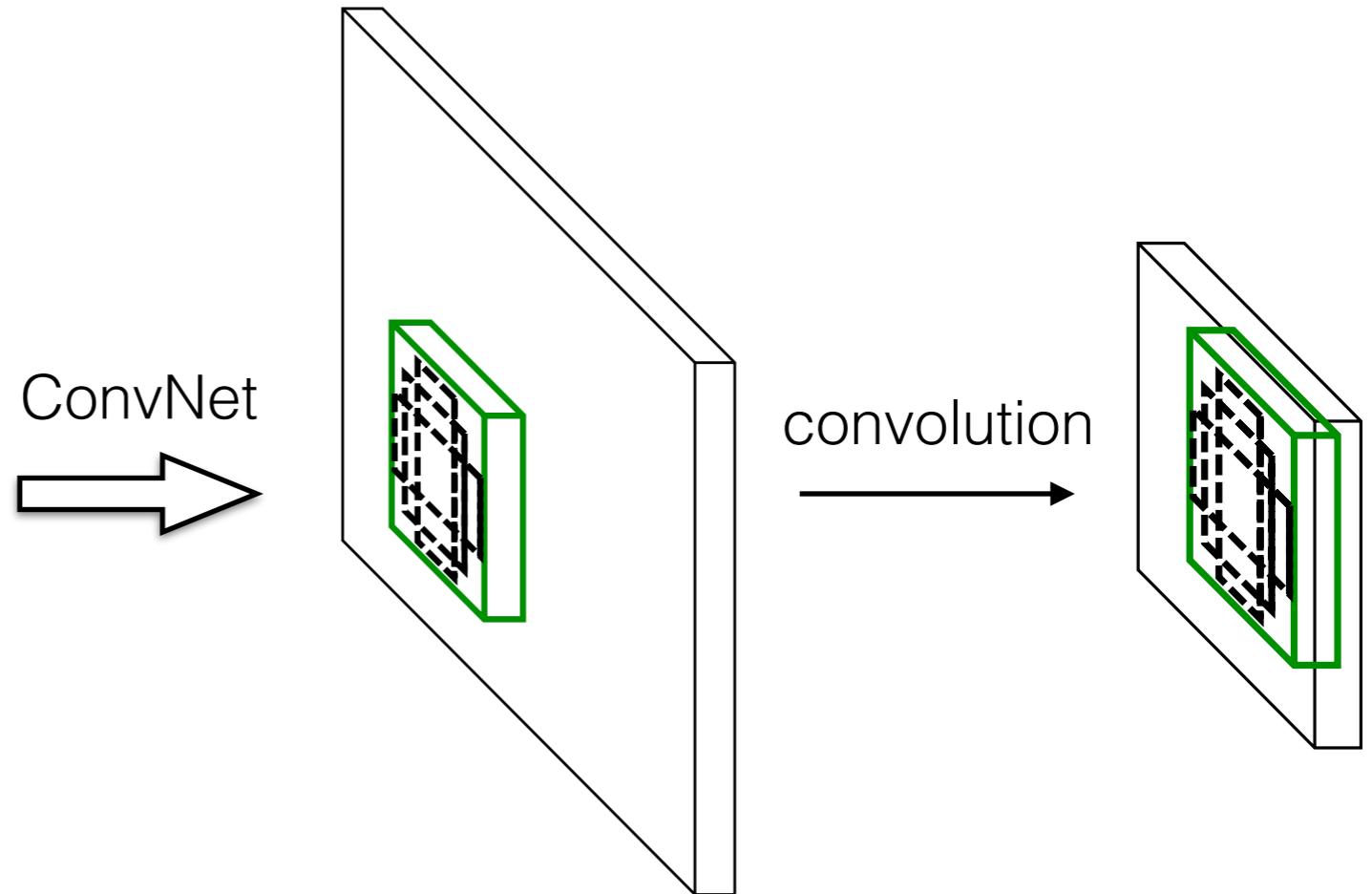
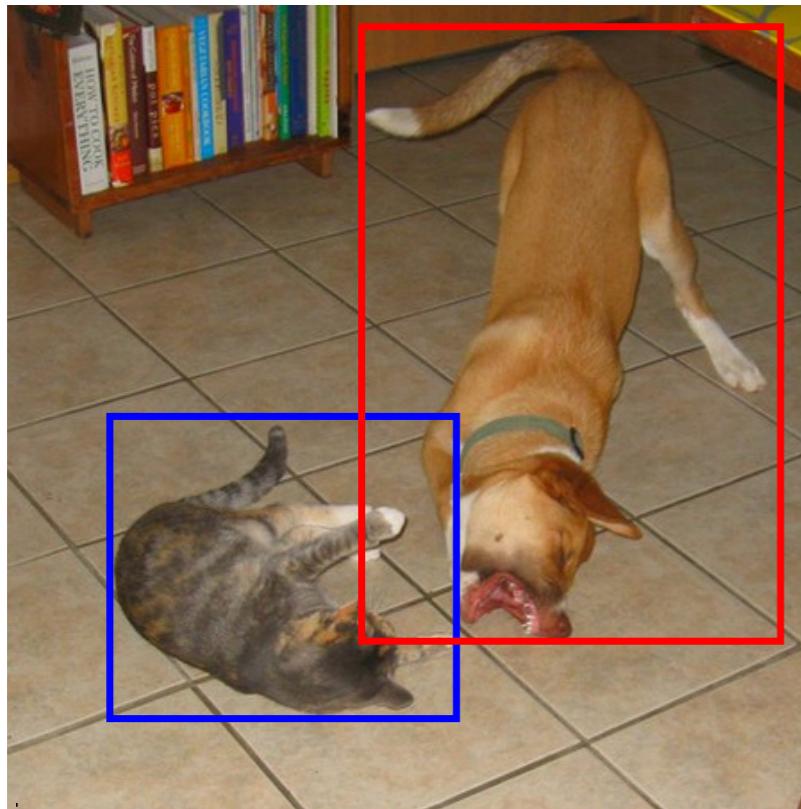


Contribution #2: Splitting the Region Space



	SSD300		
include $\{\frac{1}{2}, 2\}$ box?	✓	✓	
include $\{\frac{1}{3}, 3\}$ box?		✓	
number of Boxes	3880	7760	8732
VOC2007 test mAP	71.6	73.7	74.3

Contribution #2: Splitting the Region Space



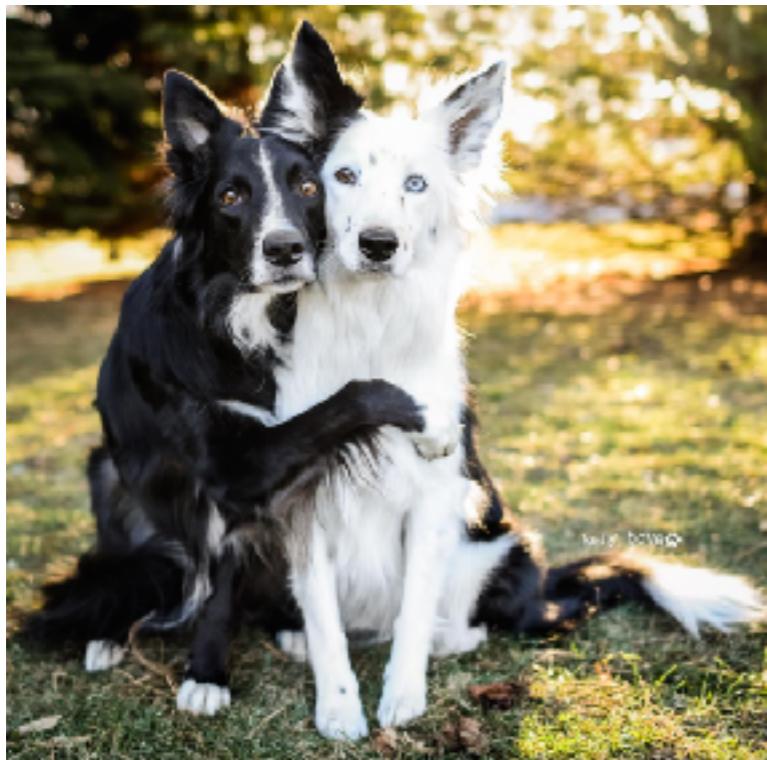
Use 38x38 feature map : **+2.5 mAP**
(conv4_3)

Why So Many Default Boxes?

	Faster R-CNN	YOLO	SSD300	SSD512
# Default Boxes	6000	98	8732	24564
Resolution	1000x600	448x448	300x300	512x512

Why So Many Default Boxes?

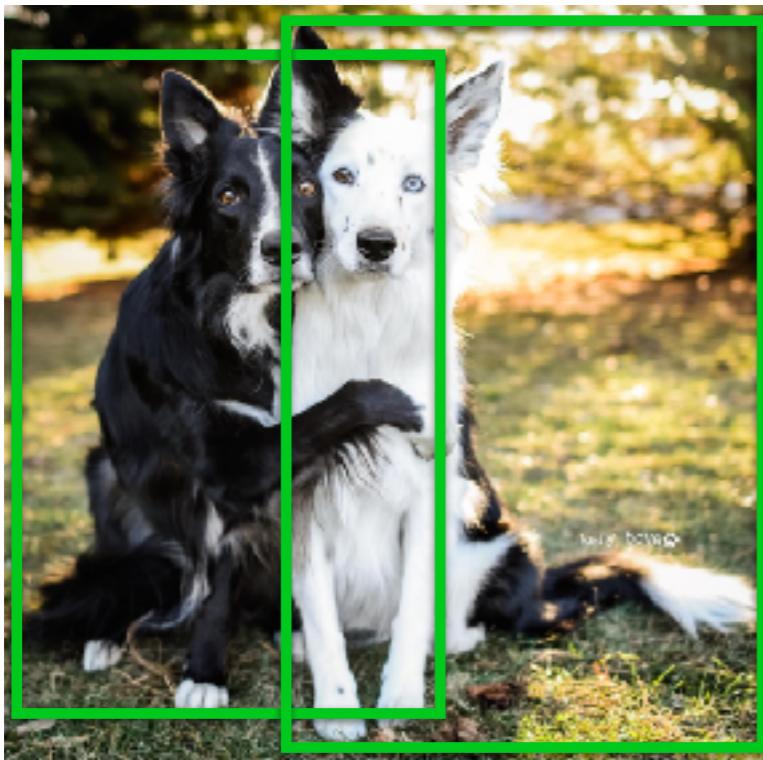
	Faster R-CNN	YOLO	SSD300	SSD512
# Default Boxes	6000	98	8732	24564
Resolution	1000x600	448x448	300x300	512x512



Why So Many Default Boxes?

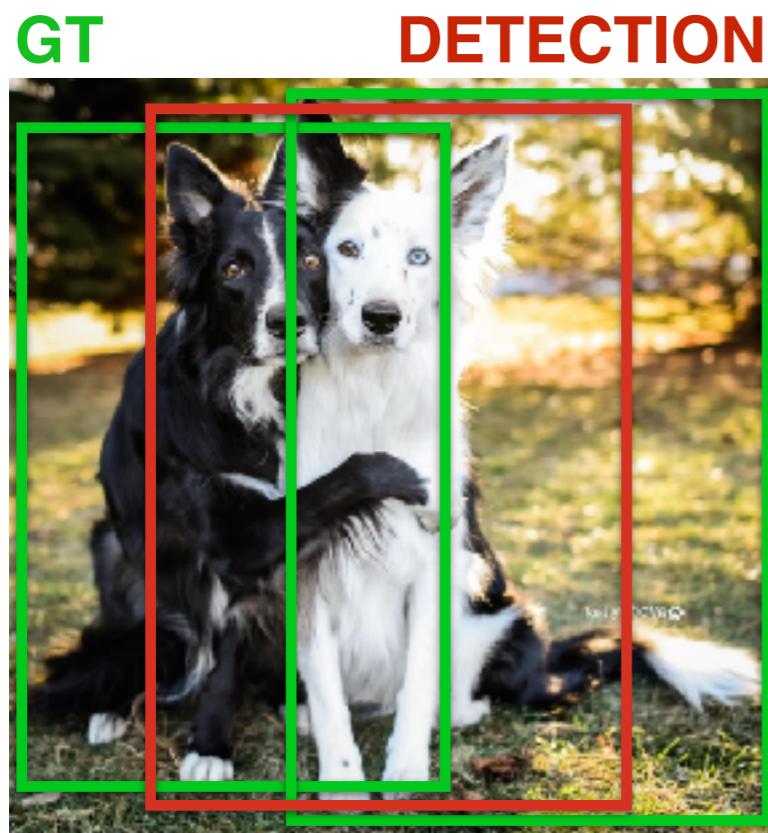
	Faster R-CNN	YOLO	SSD300	SSD512
# Default Boxes	6000	98	8732	24564
Resolution	1000x600	448x448	300x300	512x512

GT



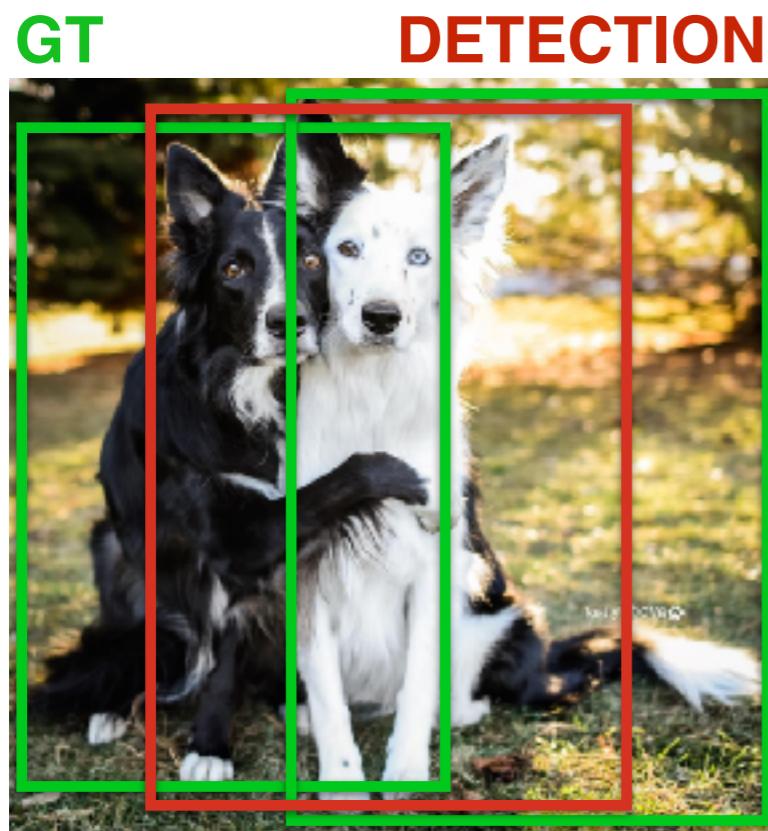
Why So Many Default Boxes?

	Faster R-CNN	YOLO	SSD300	SSD512
# Default Boxes	6000	98	8732	24564
Resolution	1000x600	448x448	300x300	512x512



Why So Many Default Boxes?

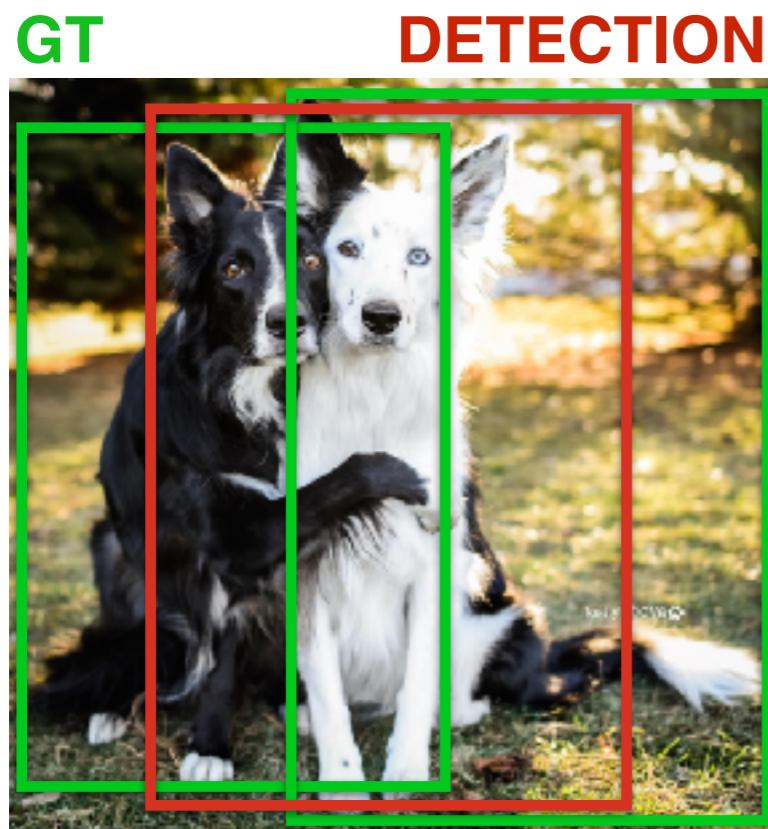
	Faster R-CNN	YOLO	SSD300	SSD512
# Default Boxes	6000	98	8732	24564
Resolution	1000x600	448x448	300x300	512x512



- SmoothL1 or L2 loss for box shape averages among likely hypotheses

Why So Many Default Boxes?

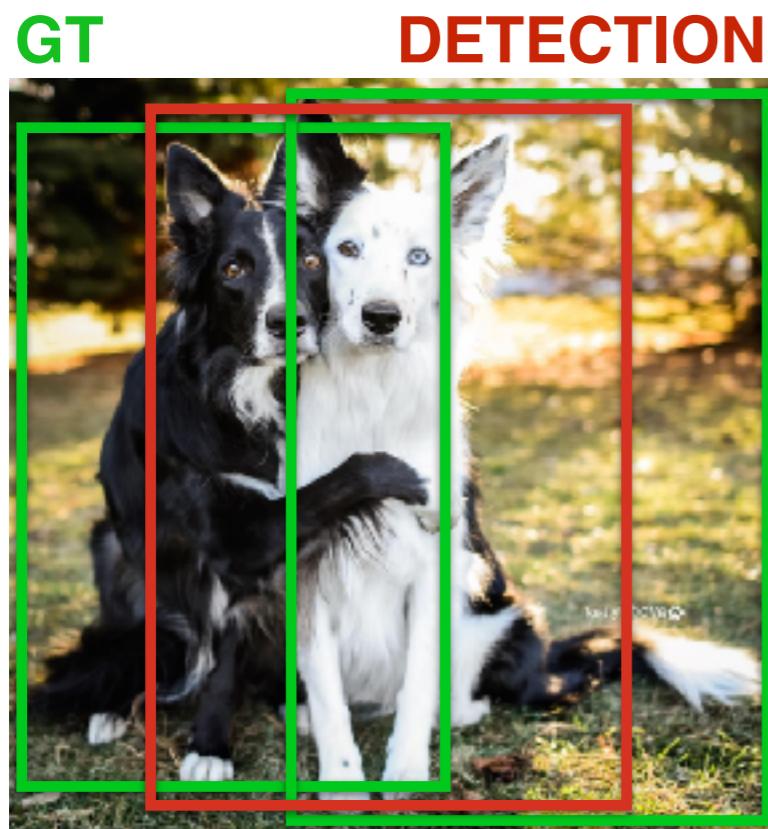
	Faster R-CNN	YOLO	SSD300	SSD512
# Default Boxes	6000	98	8732	24564
Resolution	1000x600	448x448	300x300	512x512



- SmoothL1 or L2 loss for box shape averages among likely hypotheses
- Need to have enough default boxes (discrete bins) to do accurate regression in each

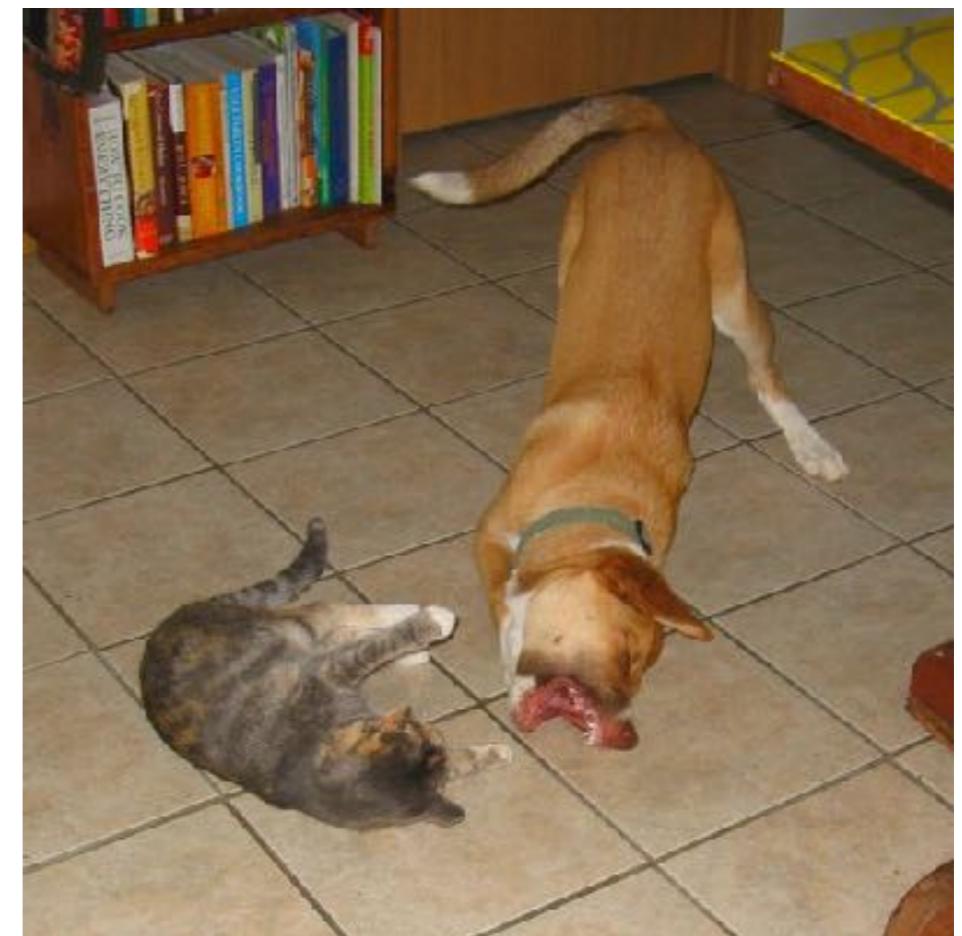
Why So Many Default Boxes?

	Faster R-CNN	YOLO	SSD300	SSD512
# Default Boxes	6000	98	8732	24564
Resolution	1000x600	448x448	300x300	512x512



- SmoothL1 or L2 loss for box shape averages among likely hypotheses
- Need to have enough default boxes (discrete bins) to do accurate regression in each
- General principle for regressing complex continuous outputs with deep nets

Handling Many Default Boxes



Handling Many Default Boxes

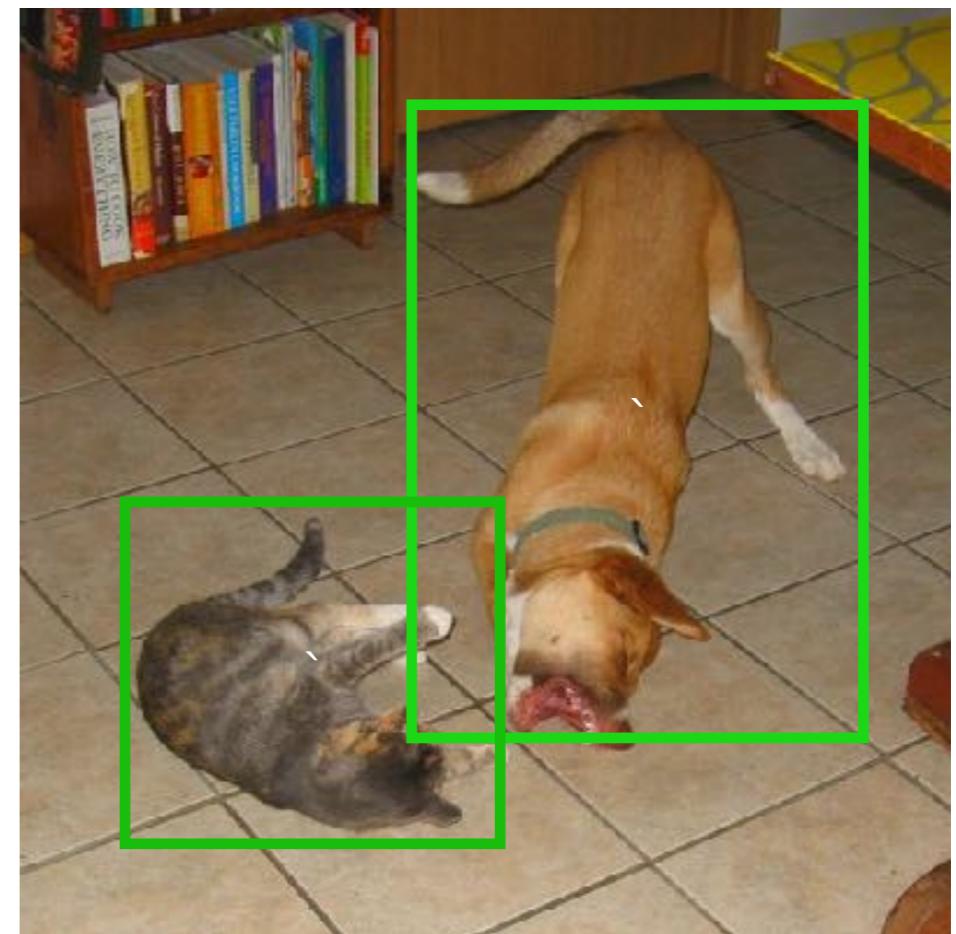
- Matching ground truth and default boxes



Handling Many Default Boxes

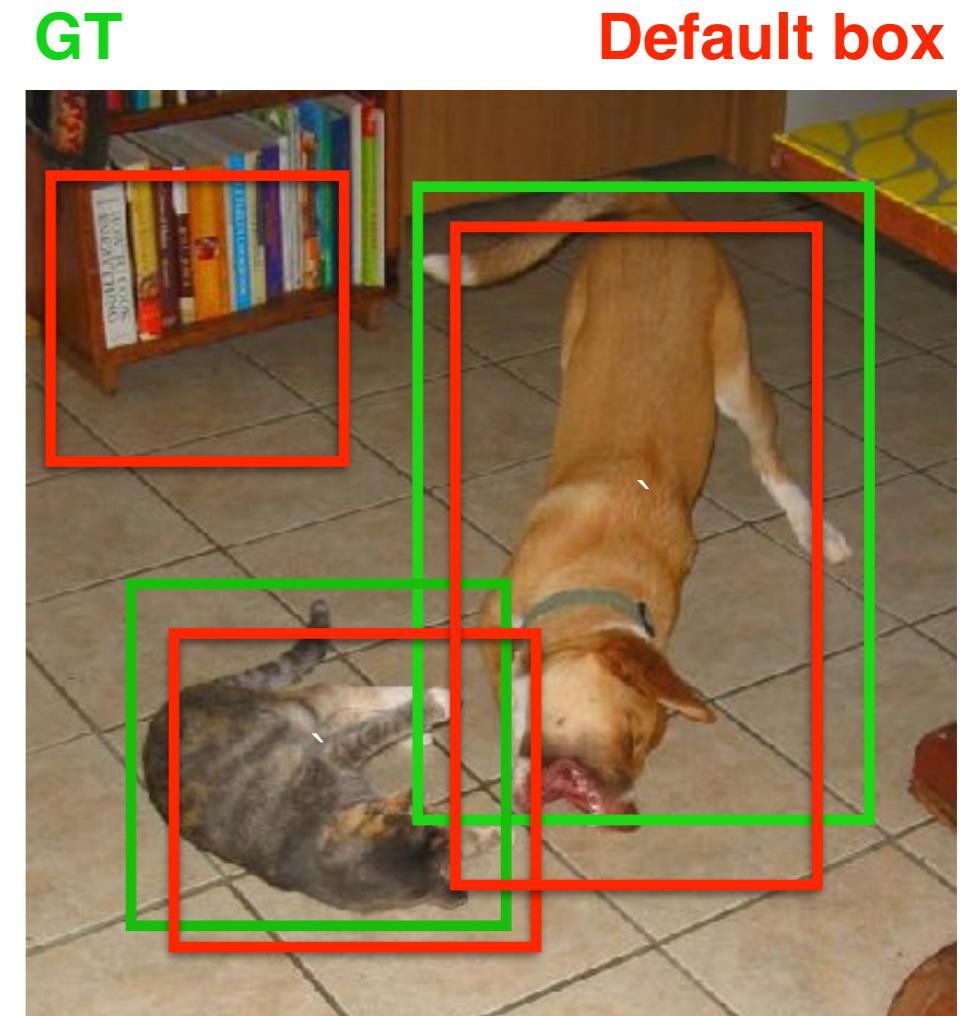
- Matching ground truth and default boxes

GT



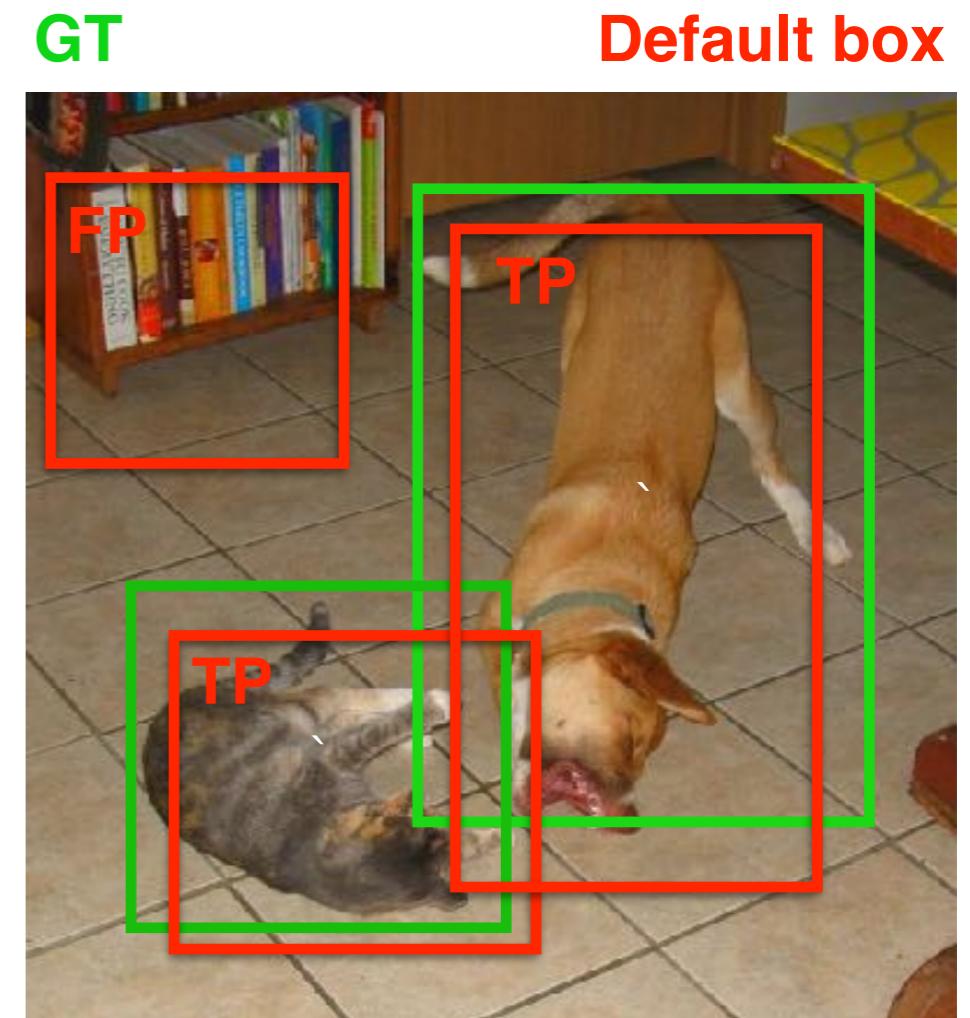
Handling Many Default Boxes

- Matching ground truth and default boxes



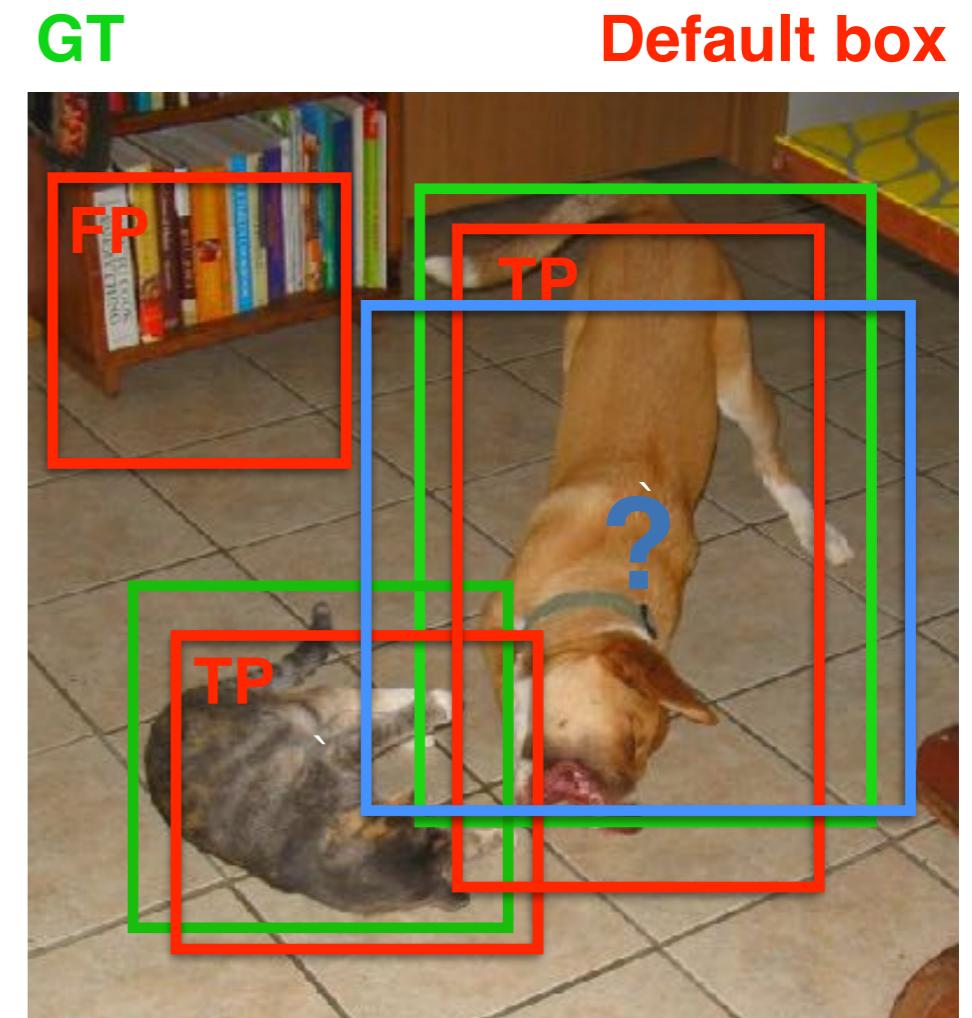
Handling Many Default Boxes

- Matching ground truth and default boxes



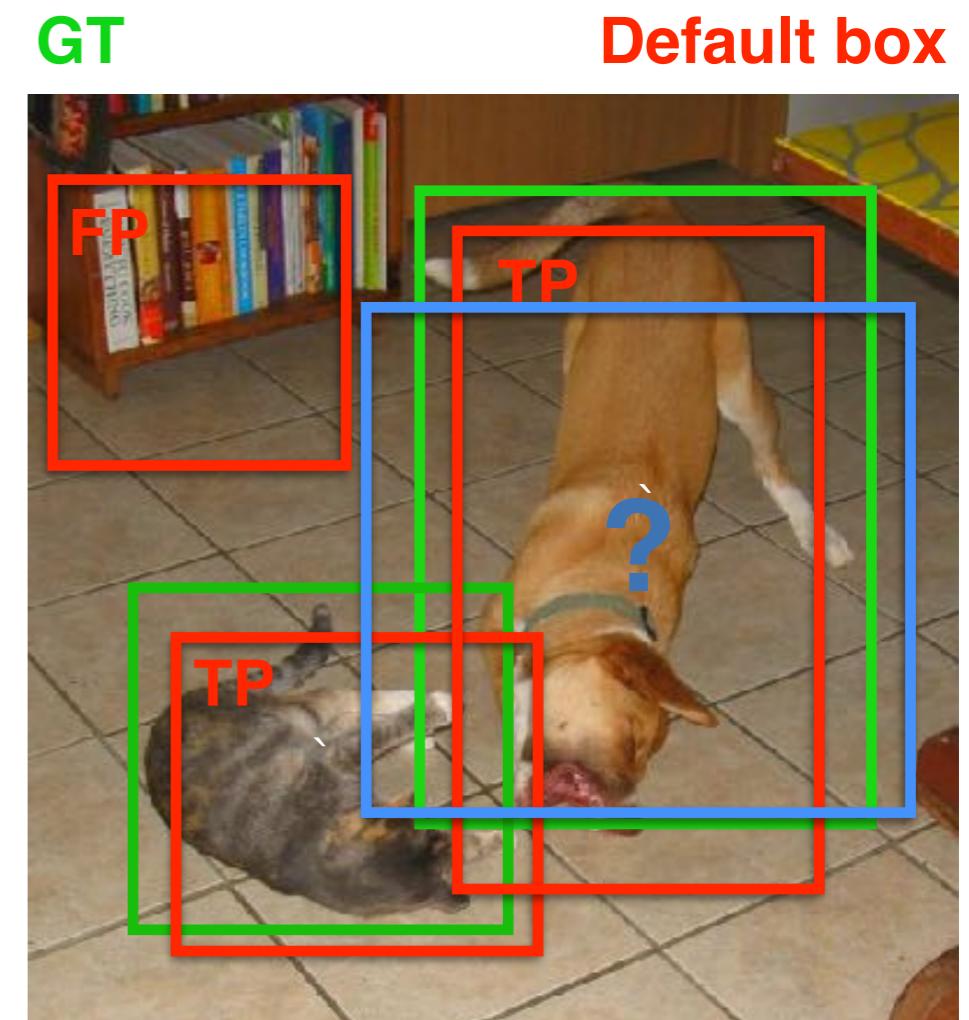
Handling Many Default Boxes

- Matching ground truth and default boxes



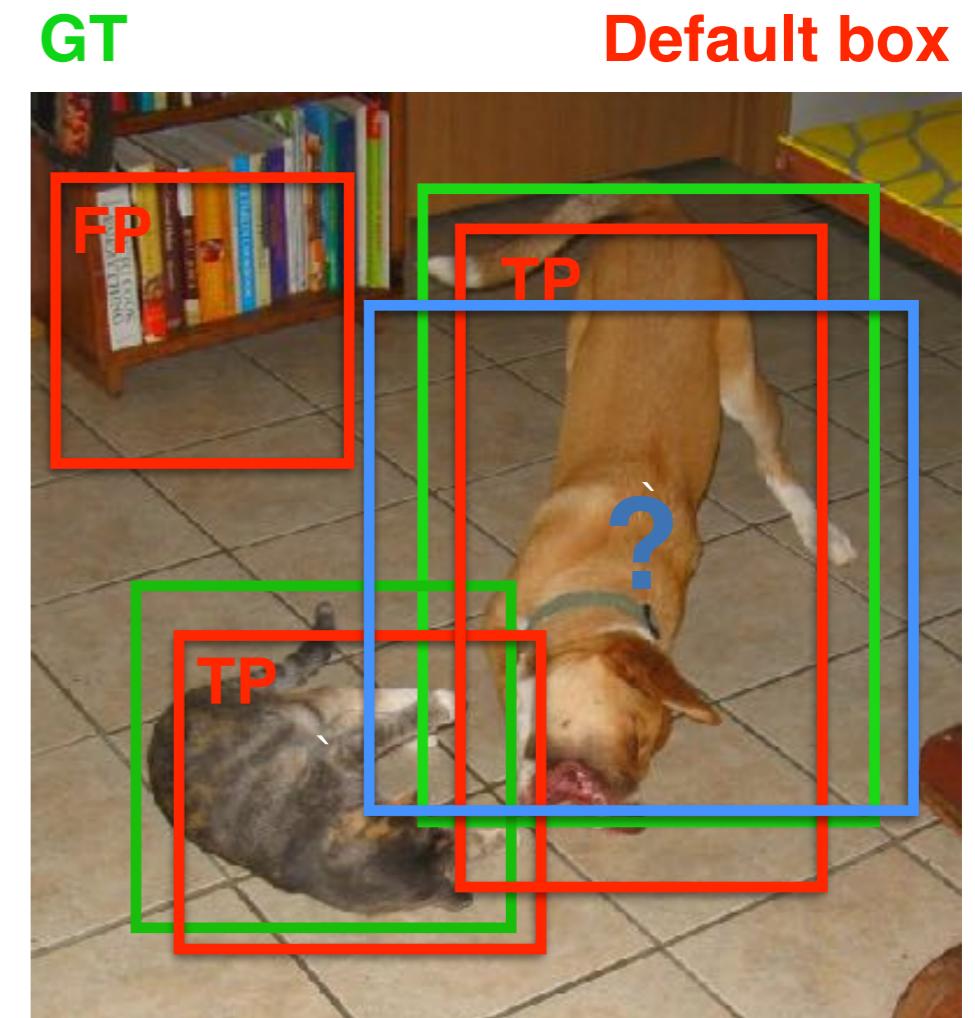
Handling Many Default Boxes

- Matching ground truth and default boxes
 - Match each GT box to closest default box



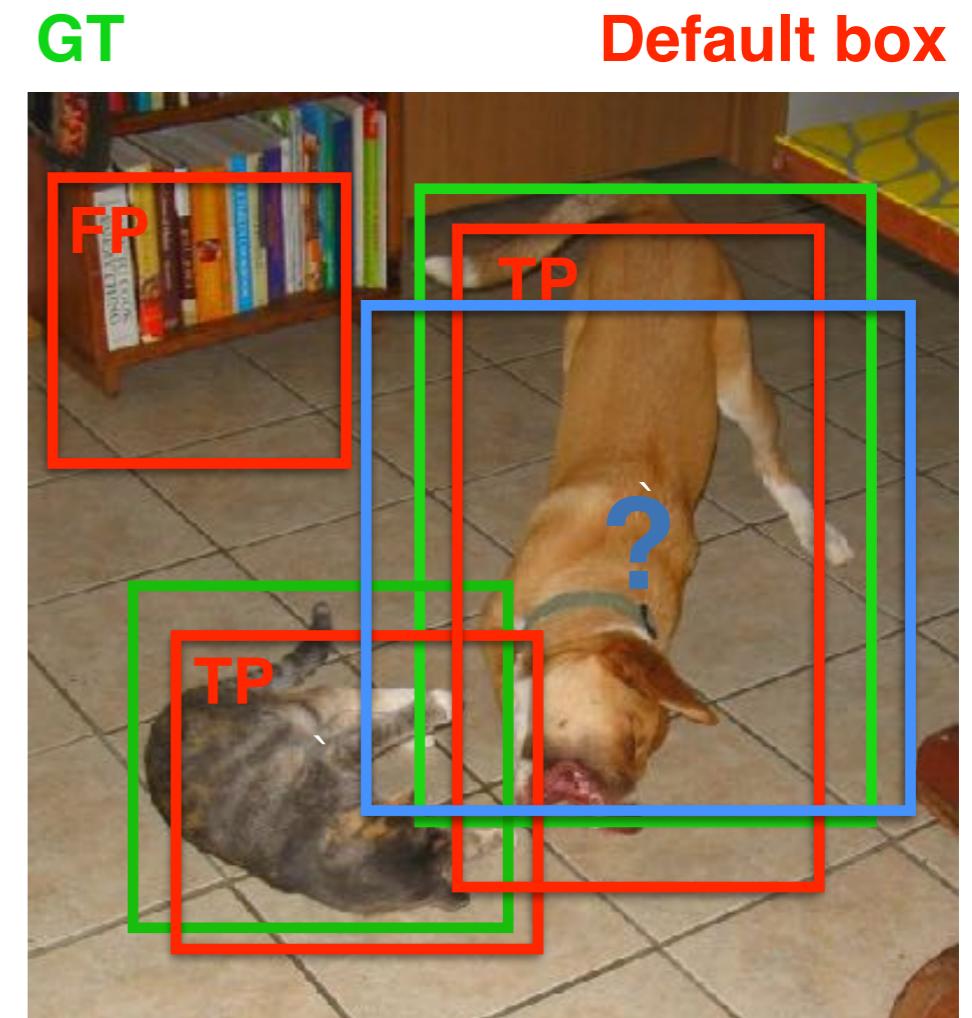
Handling Many Default Boxes

- Matching ground truth and default boxes
 - Match each GT box to closest default box
 - Also match each GT box to all unassigned default boxes with $\text{IoU} > 0.5$



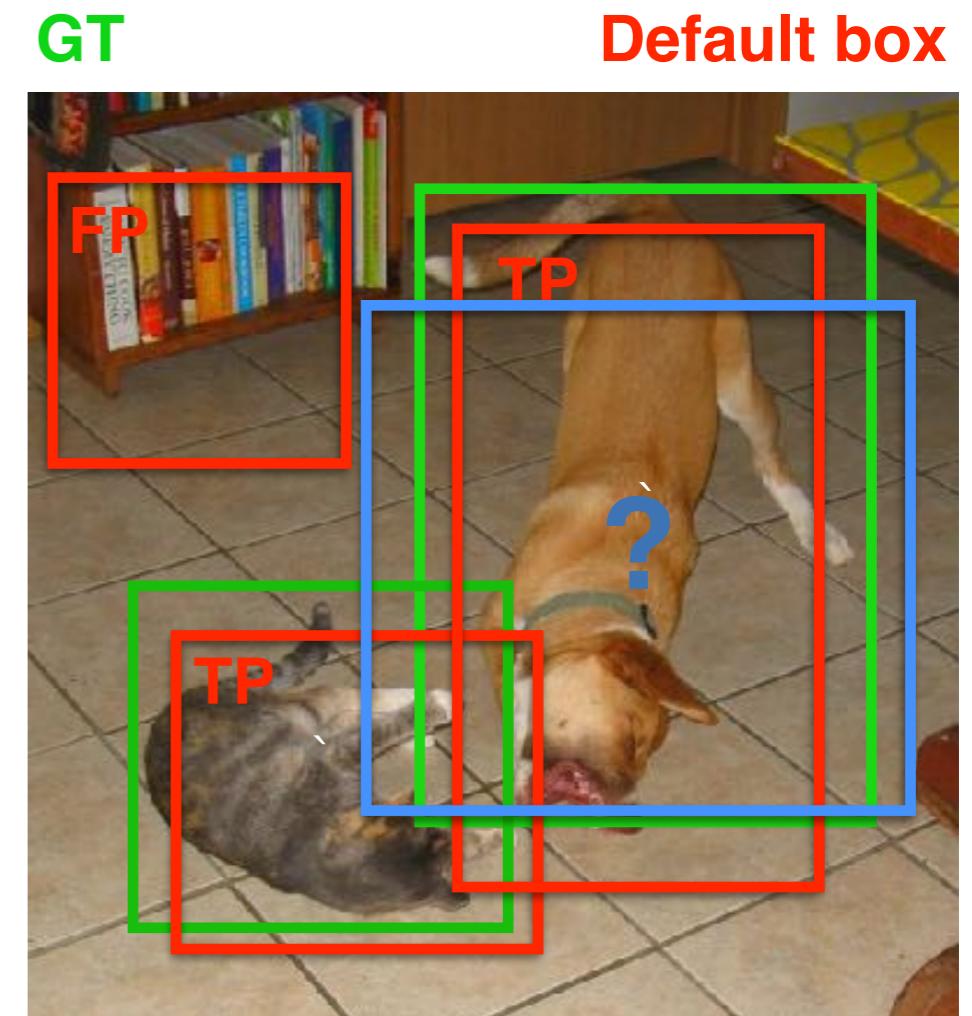
Handling Many Default Boxes

- Matching ground truth and default boxes
 - Match each GT box to closest default box
 - Also match each GT box to all unassigned default boxes with $\text{IoU} > 0.5$
- Hard negative mining



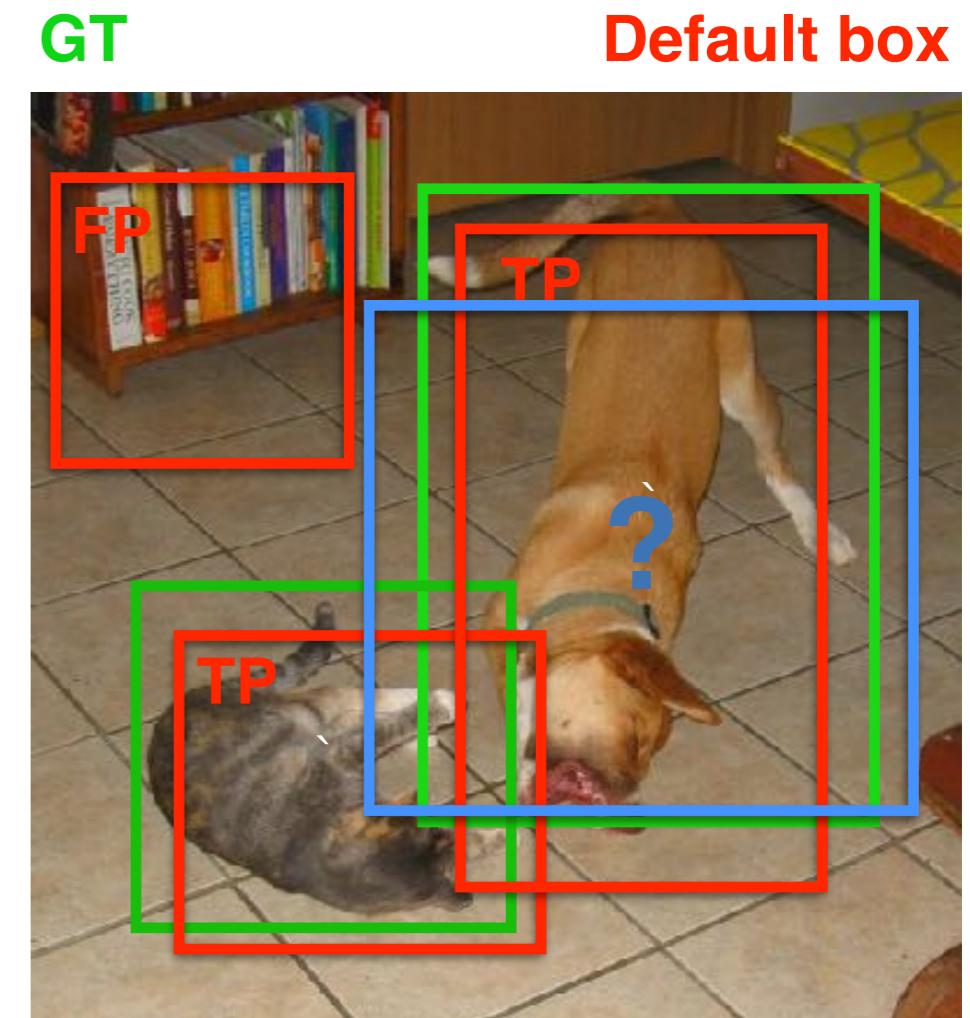
Handling Many Default Boxes

- Matching ground truth and default boxes
 - Match each GT box to closest default box
 - Also match each GT box to all unassigned default boxes with $\text{IoU} > 0.5$
- Hard negative mining
 - Unbalanced training: 1-30 TP, 8k-25k FP

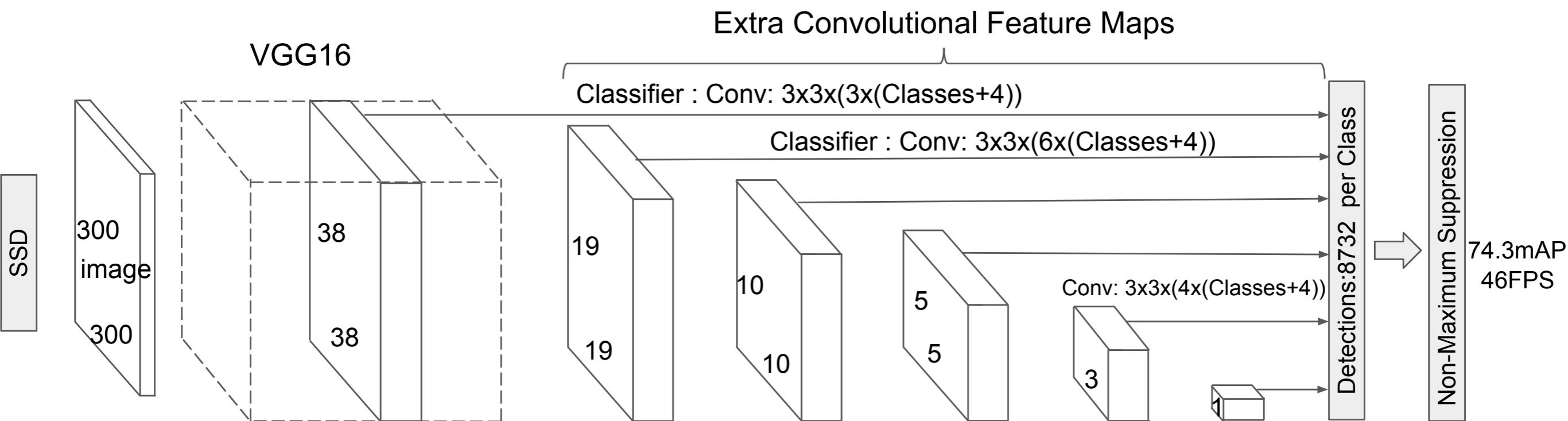


Handling Many Default Boxes

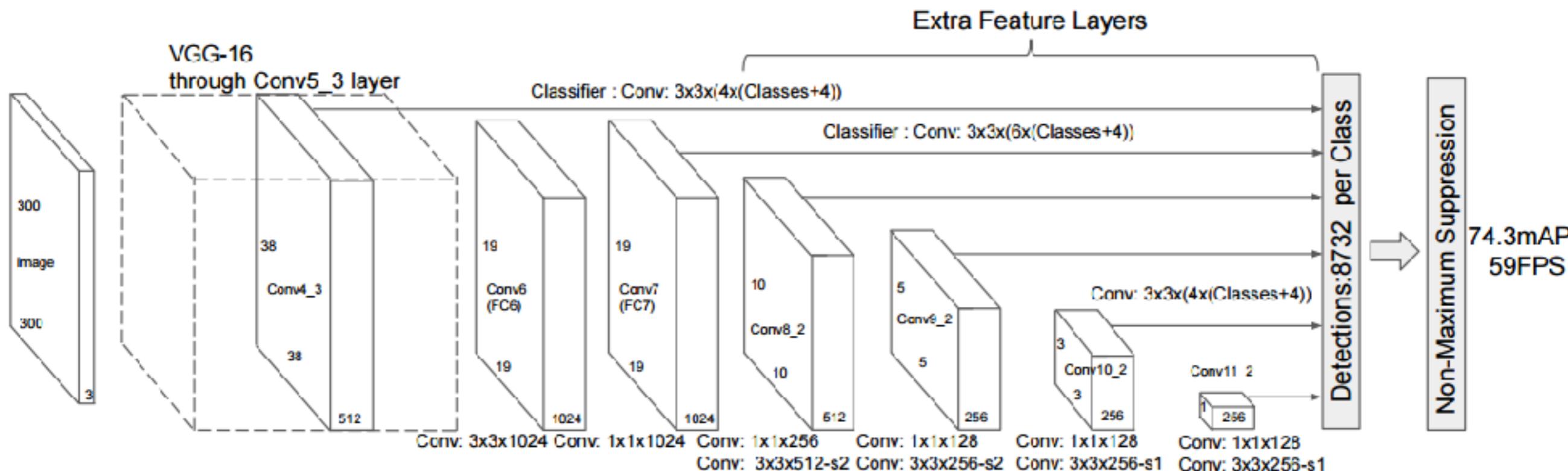
- Matching ground truth and default boxes
 - Match each GT box to closest default box
 - Also match each GT box to all unassigned default boxes with $\text{IoU} > 0.5$
- Hard negative mining
 - Unbalanced training: 1-30 TP, 8k-25k FP
 - Keep TP:FP ratio fixed (1:3), use worst-misclassified FPs.



SSD Architecture



SSD Architecture

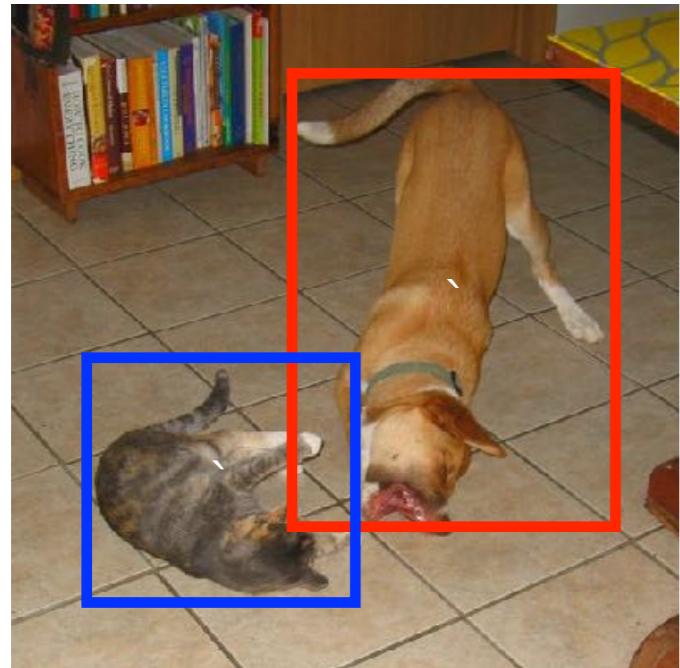


Contribution #3: The Devil is in the Details

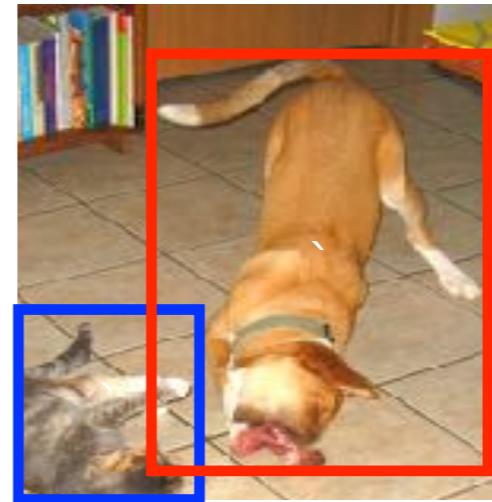
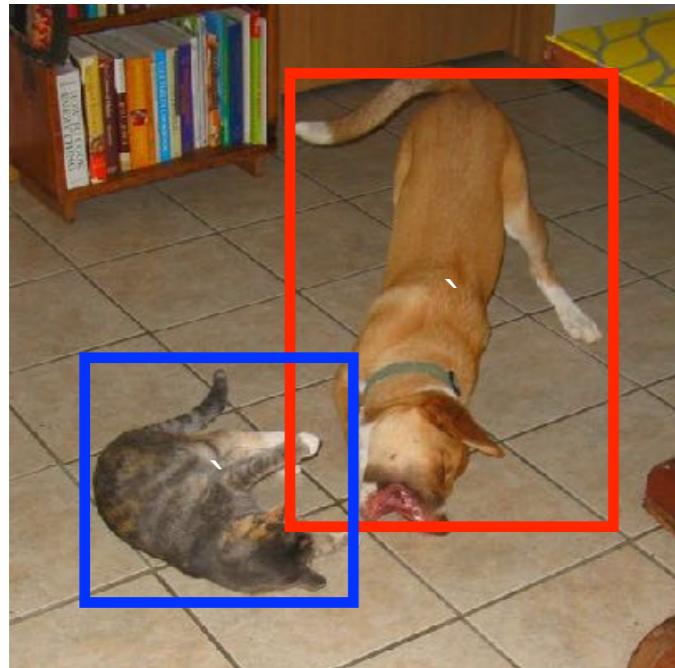


Data Augmentation

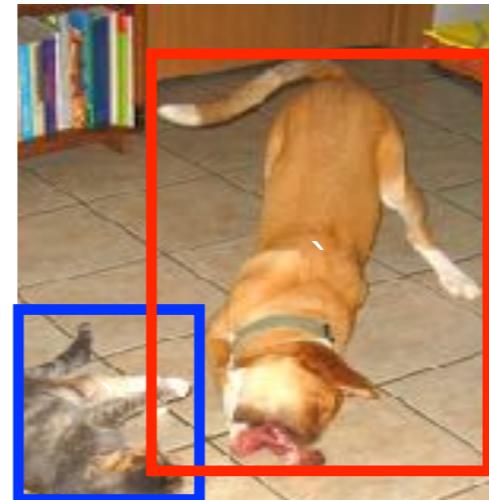
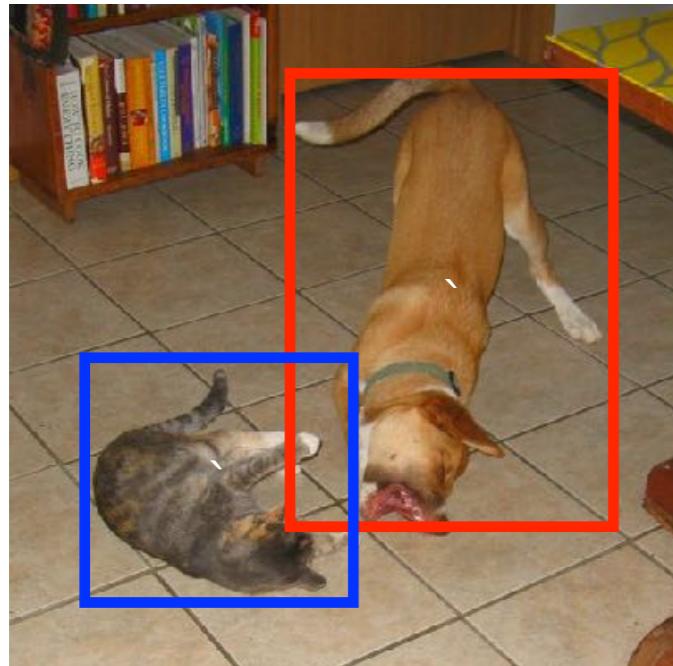
Data Augmentation



Data Augmentation



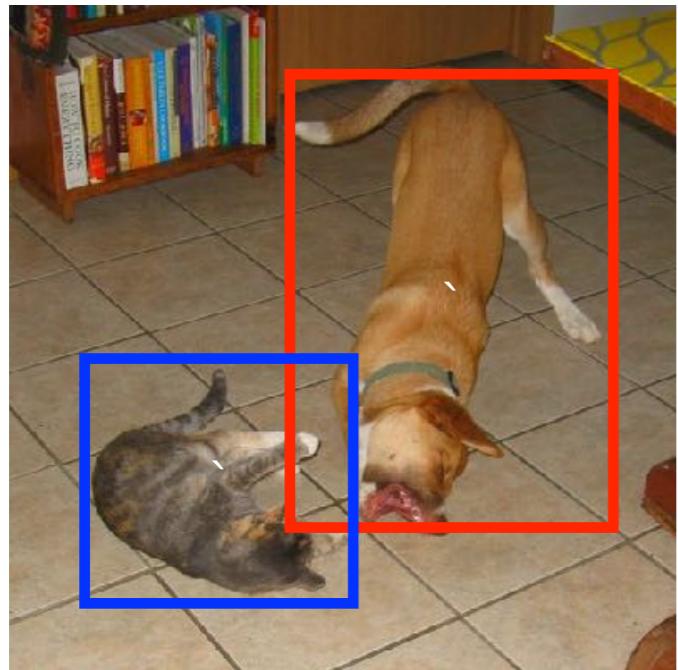
Data Augmentation



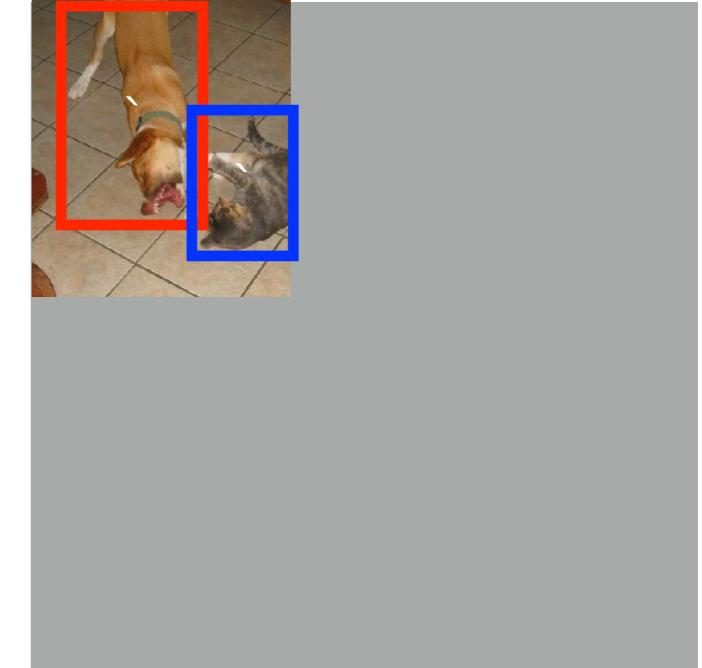
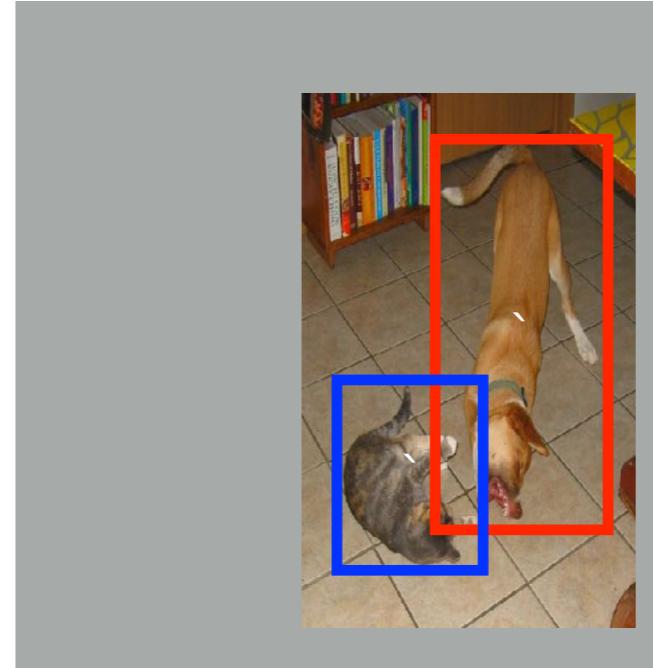
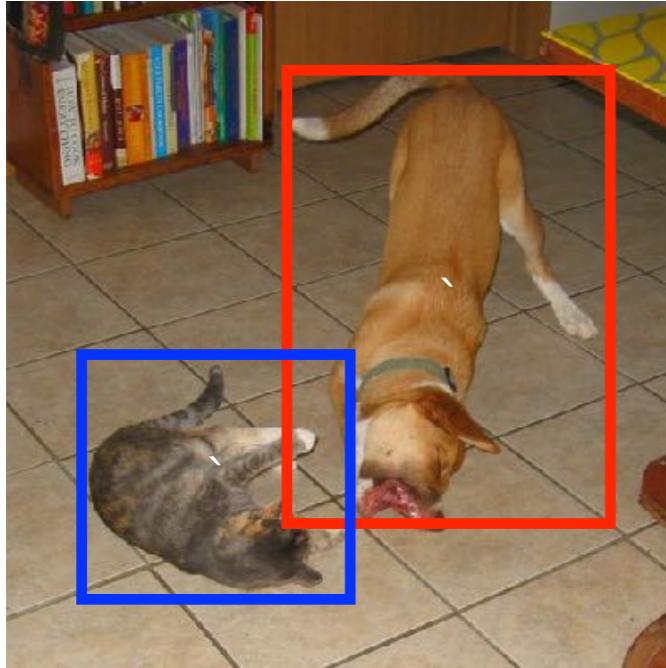
data augmentation	SSD300	
horizontal flip	✓	✓
random crop & color distortion		✓
VOC2007 test mAP	65.5	74.3

Data Augmentation

Data Augmentation

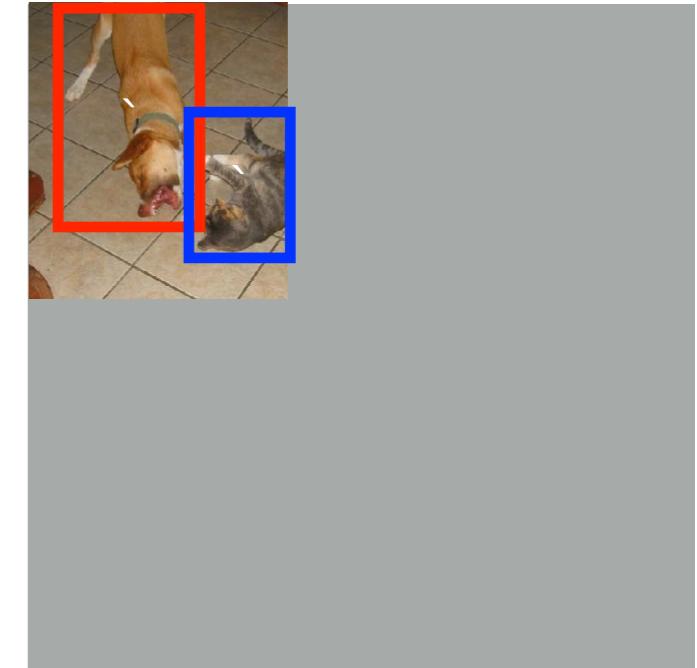
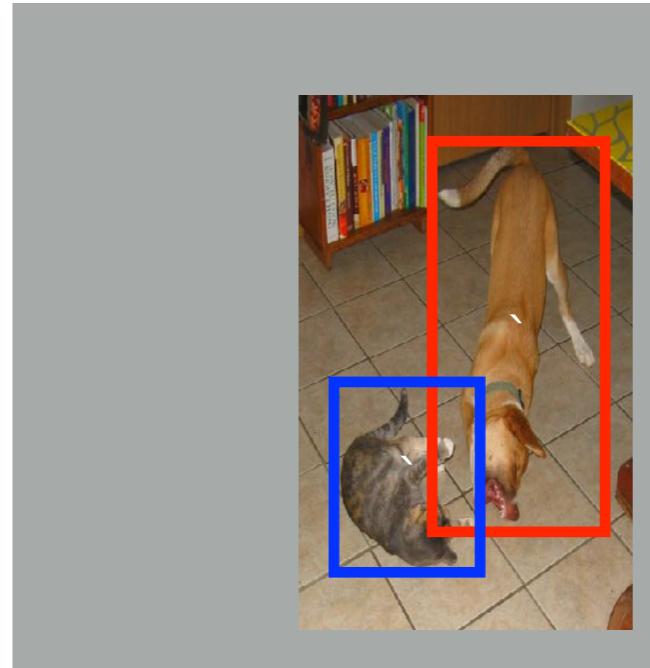
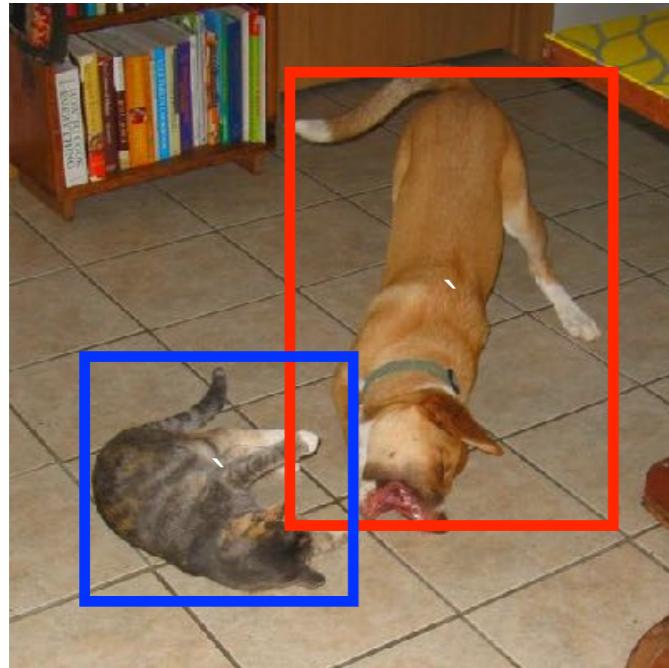


Data Augmentation



Random expansion creates more
small training examples

Data Augmentation



Random expansion creates more
small training examples

data augmentation	SSD300		
horizontal flip	✓	✓	✓
random crop & color distortion		✓	✓
random expansion			✓
VOC2007 test mAP	65.5	74.3	77.2

Results on VOC2007 test

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	74.3	46	1	8732	300 × 300
SSD512	76.8	19	1	24564	512 × 512
SSD300	74.3	59	8	8732	300 × 300
SSD512	76.8	22	8	24564	512 × 512

Results on VOC2007 test

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	6.6x ↑	74.3	46	1	8732
SSD512	10% ↑	76.8	19	1	24564
SSD300	74.3	59	8	8732	300 × 300
SSD512	76.8	22	8	24564	512 × 512

Results on VOC2007 test

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	6.6x ↑	74.3	46	1	8732
SSD512	10% ↑	76.8	19	1	24564
SSD300	74.3	59	8	8732	300 × 300
SSD512	76.8	22	8	24564	512 × 512

Results on VOC2007 test

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	6.6x ↑	74.3	46	1	8732
SSD512	10% ↑	76.8	19	1	24564
SSD300	74.3	59	8	8732	300 × 300
SSD512	76.8	22	8	24564	512 × 512

Results on VOC2007 test

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000×600
Fast YOLO	52.7	155	1	98	448×448
YOLO (VGG16)	66.4	21	1	98	448×448
SSD300	74.3	46	1	8732	300×300
SSD512	76.8	19	1	24564	512×512
SSD300	74.3	59	8	8732	300×300
SSD512	76.8	22	8	24564	512×512

Results on VOC2007 test

Method	mAP	FPS	batch size	# Boxes	Input resolution
Faster R-CNN (VGG16)	73.2	7	1	~ 6000	~ 1000 × 600
Fast YOLO	52.7	155	1	98	448 × 448
YOLO (VGG16)	66.4	21	1	98	448 × 448
SSD300	77.2	74.3	1	8732	300 × 300
SSD512	79.8	76.8	1	24564	512 × 512
SSD300	77.2	74.3	8	8732	300 × 300
SSD512	79.8	76.8	8	24564	512 × 512

Results on More Datasets

Results on More Datasets

Method	VOC2007 test	VOC2012 test	MS COCO test-dev	ILSVRC2014 val2
Fast R-CNN	70.0	68.4	19.7	N/A
Faster R-CNN	73.2	70.4	21.9	N/A
YOLO	63.4	57.9	N/A	N/A

Results on More Datasets

Method	VOC2007 test	VOC2012 test	MS COCO test-dev	ILSVRC2014 val2
Fast R-CNN	70.0	68.4	19.7	N/A
Faster R-CNN	73.2	70.4	21.9	N/A
YOLO	63.4	57.9	N/A	N/A
SSD300	74.3	72.4	23.2	43.4

Results on More Datasets

Method	VOC2007 test	VOC2012 test	MS COCO test-dev	ILSVRC2014 val2
Fast R-CNN	70.0	68.4	19.7	N/A
Faster R-CNN	73.2	70.4	21.9	N/A
YOLO	63.4	57.9	N/A	N/A
SSD300	74.3	72.4	23.2	43.4
SSD512	76.8	74.9	26.8	46.4

Results on More Datasets

Method	VOC2007 test	VOC2012 test	MS COCO test-dev	ILSVRC2014 val2
Fast R-CNN	70.0	68.4	19.7	N/A
Faster R-CNN	73.2	70.4	21.9	N/A
YOLO	63.4	57.9	N/A	N/A
SSD300*	77.2	75.8	25.1	N/A
SSD512*	79.8	78.5	28.8	N/A

COCO Bounding Box precision

COCO Bounding Box precision

mAP @ IoU	0.5	0.75	0.5:0.95
Faster R-CNN	45.3	23.5	24.2
SSD512*	48.5	30.3	28.8
gain	+3.2	+6.8	+4.6

Future Work

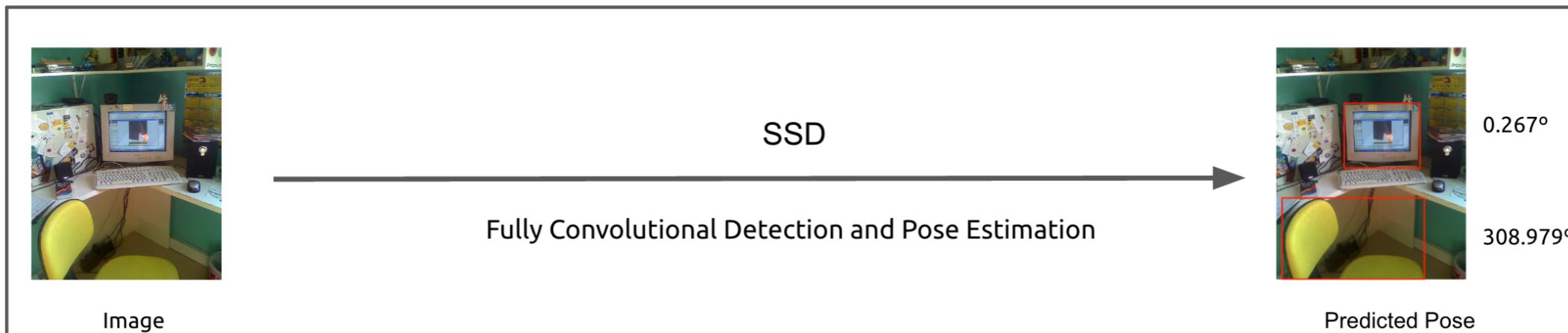
Future Work

- Object detection + pose estimation

Future Work

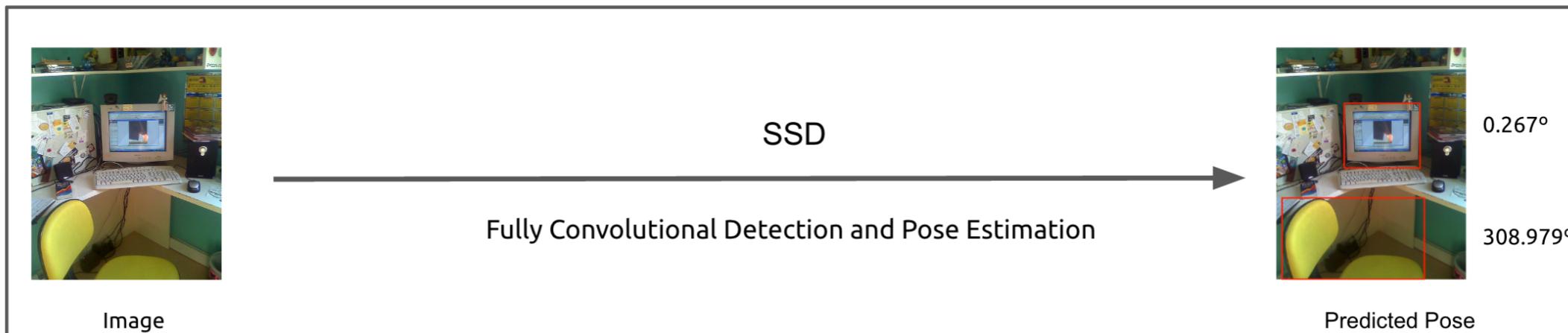
- Object detection + pose estimation

[Poirson et al, coming out at 3DV, 2016]



Future Work

- Object detection + pose estimation
[Poirson et al, coming out at 3DV, 2016]

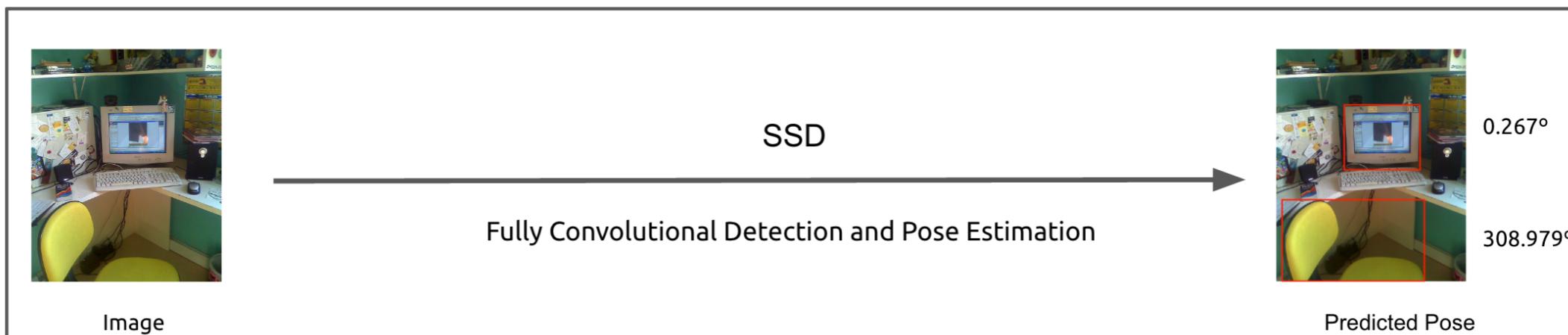


- Single shot 3D bounding box detection

Future Work

- Object detection + pose estimation

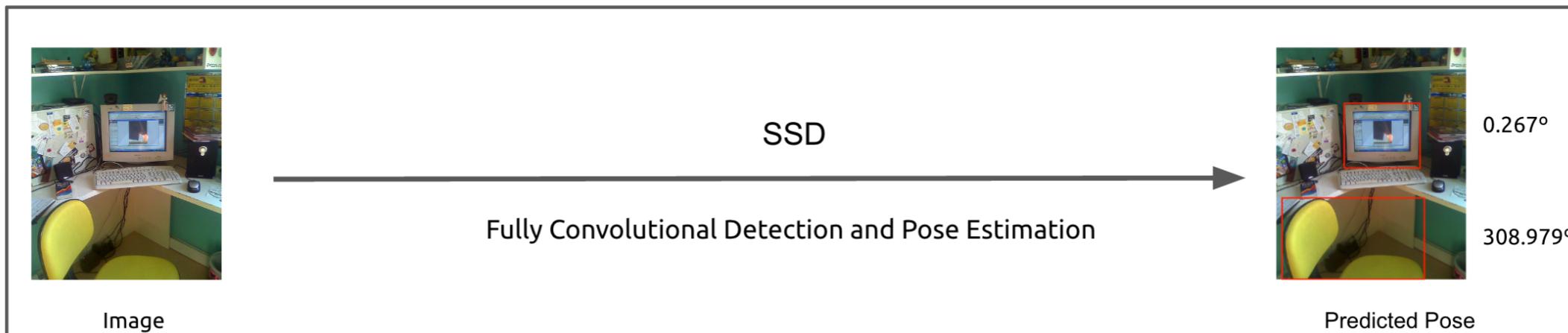
[Poirson et al, coming out at 3DV, 2016]



- Single shot 3D bounding box detection
- Joint object detection + tracking model

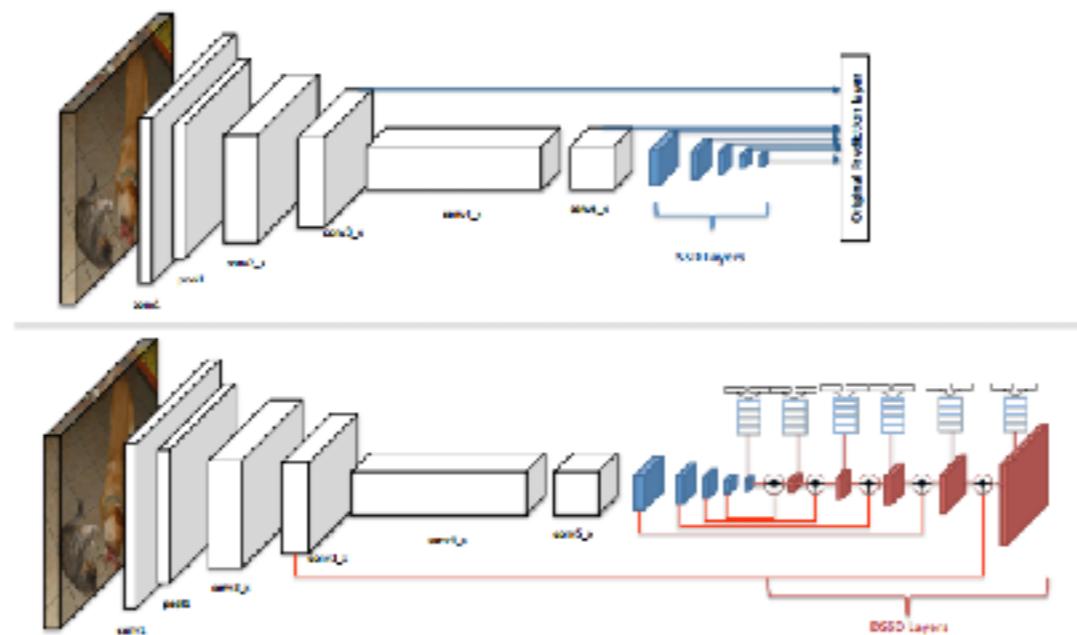
Future Work

- Object detection + pose estimation
[Poirson et al, coming out at 3DV, 2016]



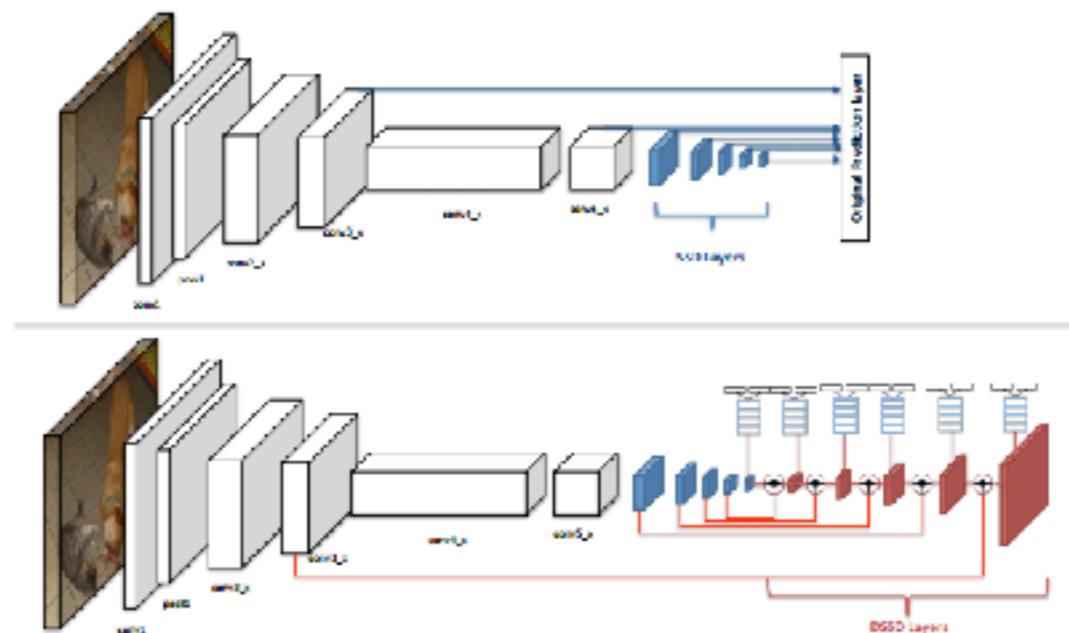
- Single shot 3D bounding box detection
- Joint object detection + tracking model
- <https://arxiv.org/abs/1609.05590>

Future Work



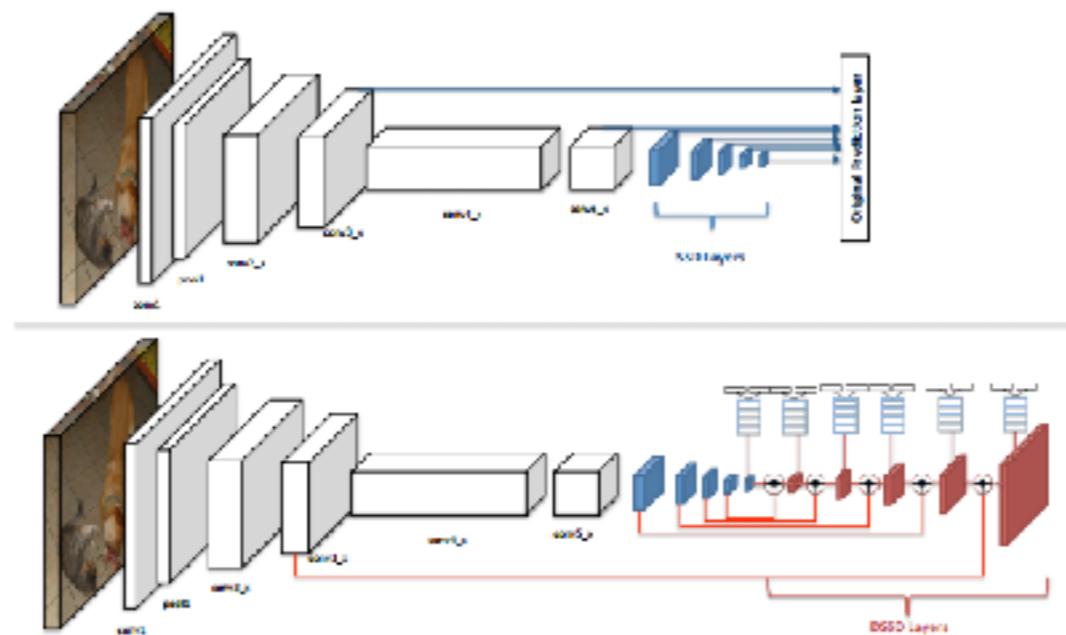
Future Work

- DSSD: Deconvolutional Single-Shot Detector



Future Work

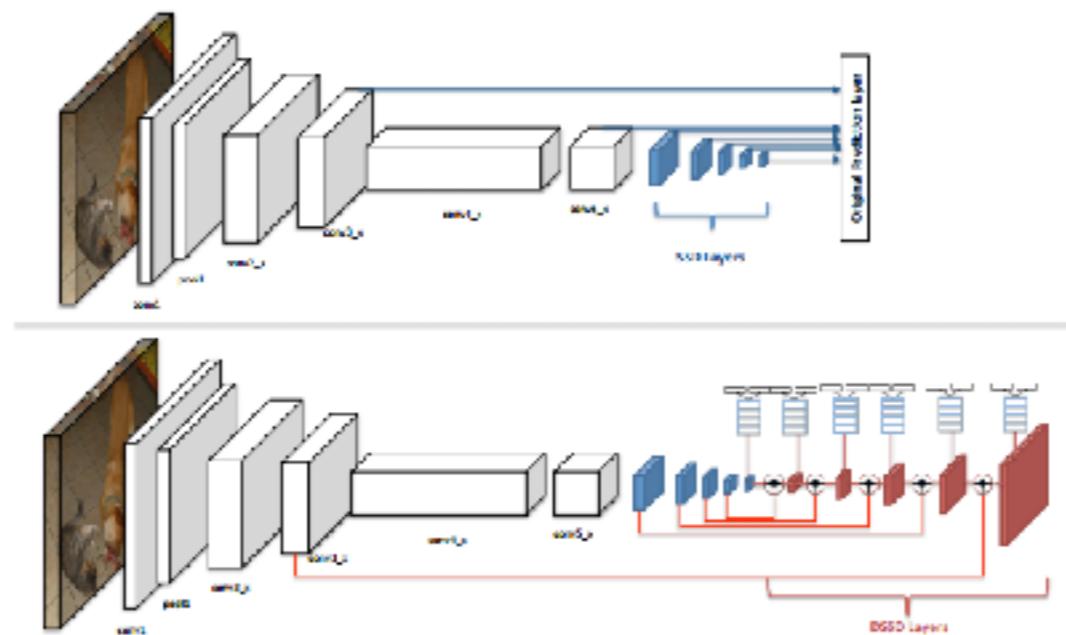
- DSSD: Deconvolutional Single-Shot Detector



- Deconvolution layers added to the end

Future Work

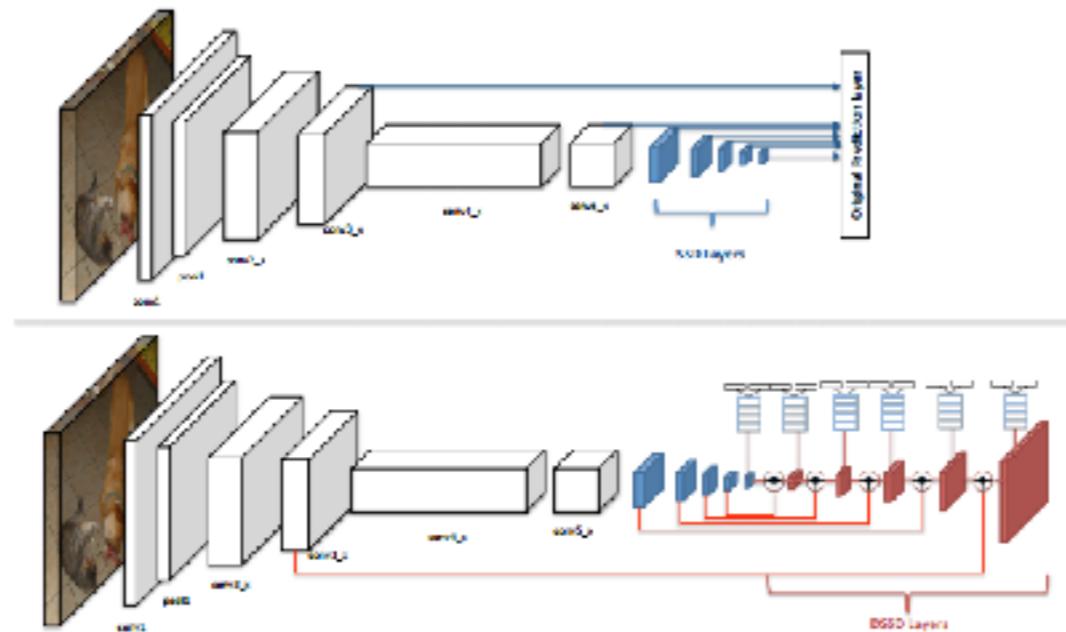
- DSSD: Deconvolutional Single-Shot Detector



- Deconvolution layers added to the end
- Improves performance even further: 81.5% mAP

Future Work

- DSSD: Deconvolutional Single-Shot Detector



- Deconvolution layers added to the end
- Improves performance even further: 81.5% mAP
- <https://arxiv.org/abs/1701.06659>

Check out the code/models



<https://github.com/weiliu89/caffe/tree/ssd>