# Project Admin

Someone from every team must come and see me!

If you chose:
- Depth Prediction
- Neural Style Transfer
- Image Captioning
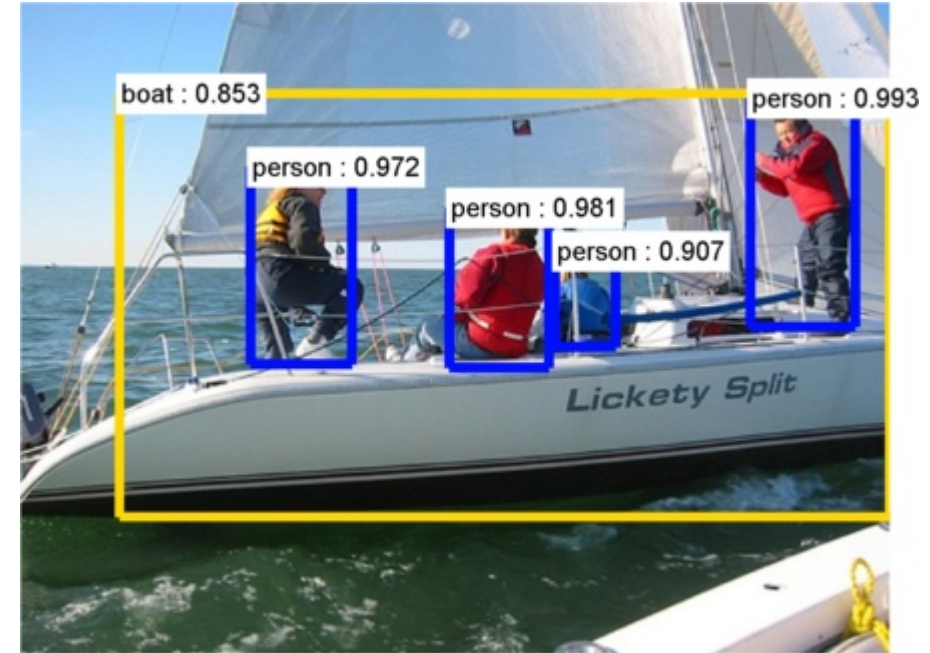
I will outline what project involves IMMEADIATELY after this class

For other projects, come and see me during office hours.

# Object Detection



Image Classification
(what?)

Object Detection
(what + where?)

# Detection with ConvNets

- So far, all about classification

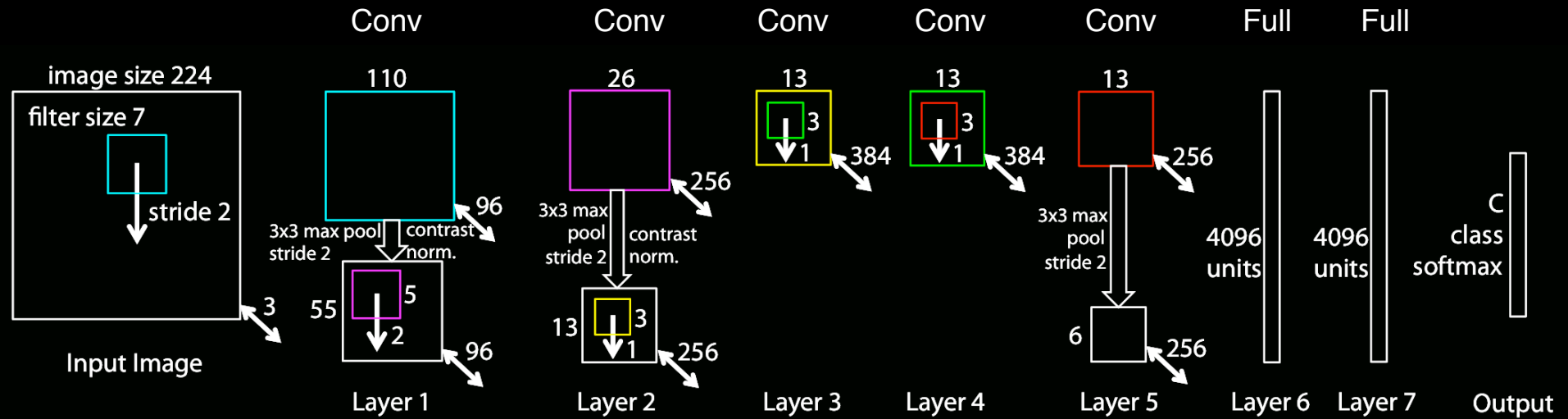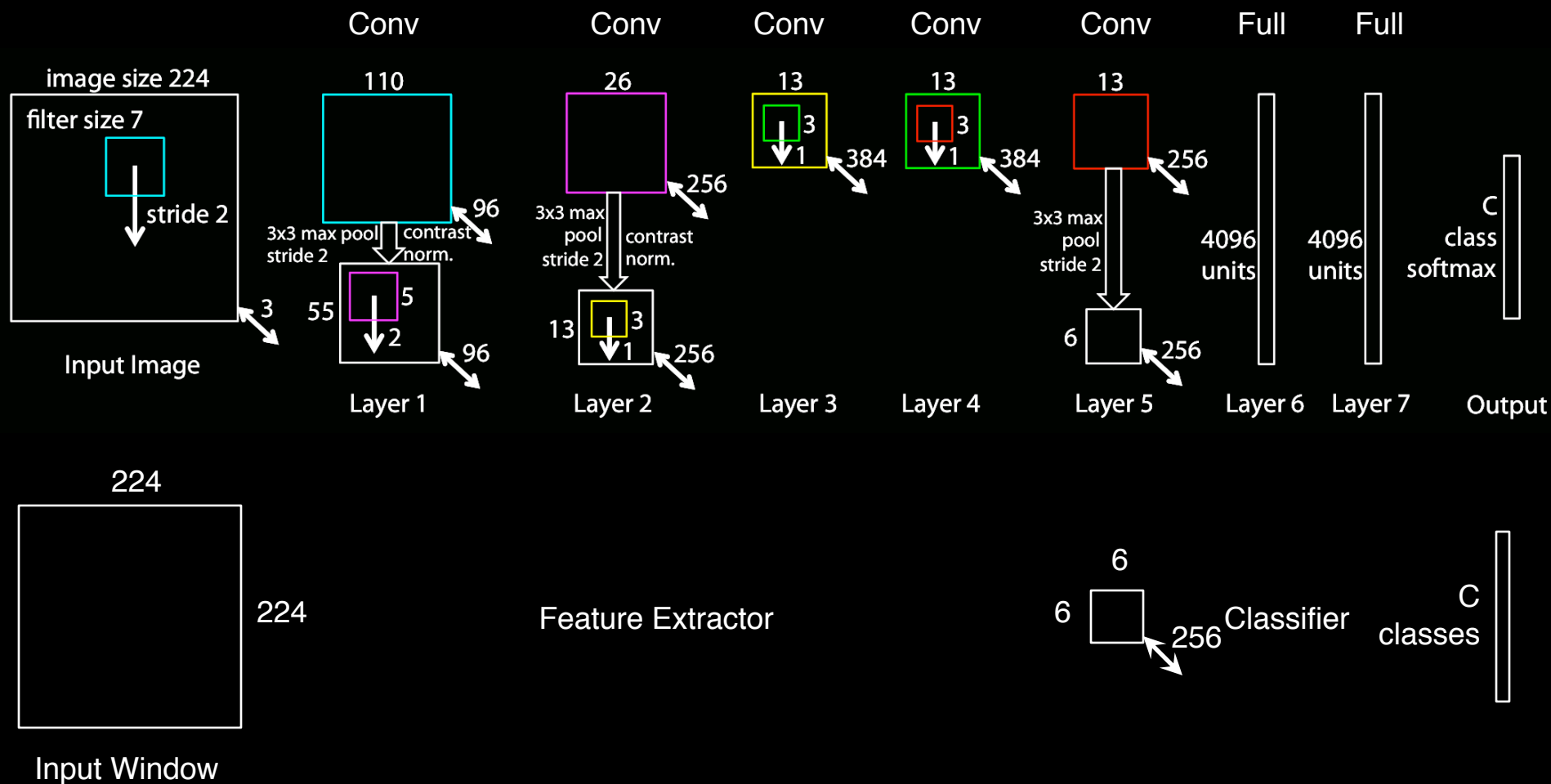- What about localizing objects within the scene?

# Two General Approaches

1. Examine very position / scale
   - E.g. Overfeat: Integrated recognition, localization and detection using convolutional networks, Sermanet et al., ICLR 2014

2. Use some kind of proposal mechanism to attend to a set of possible regions
   - E.g. Region-CNN [Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al., CVPR 2014]
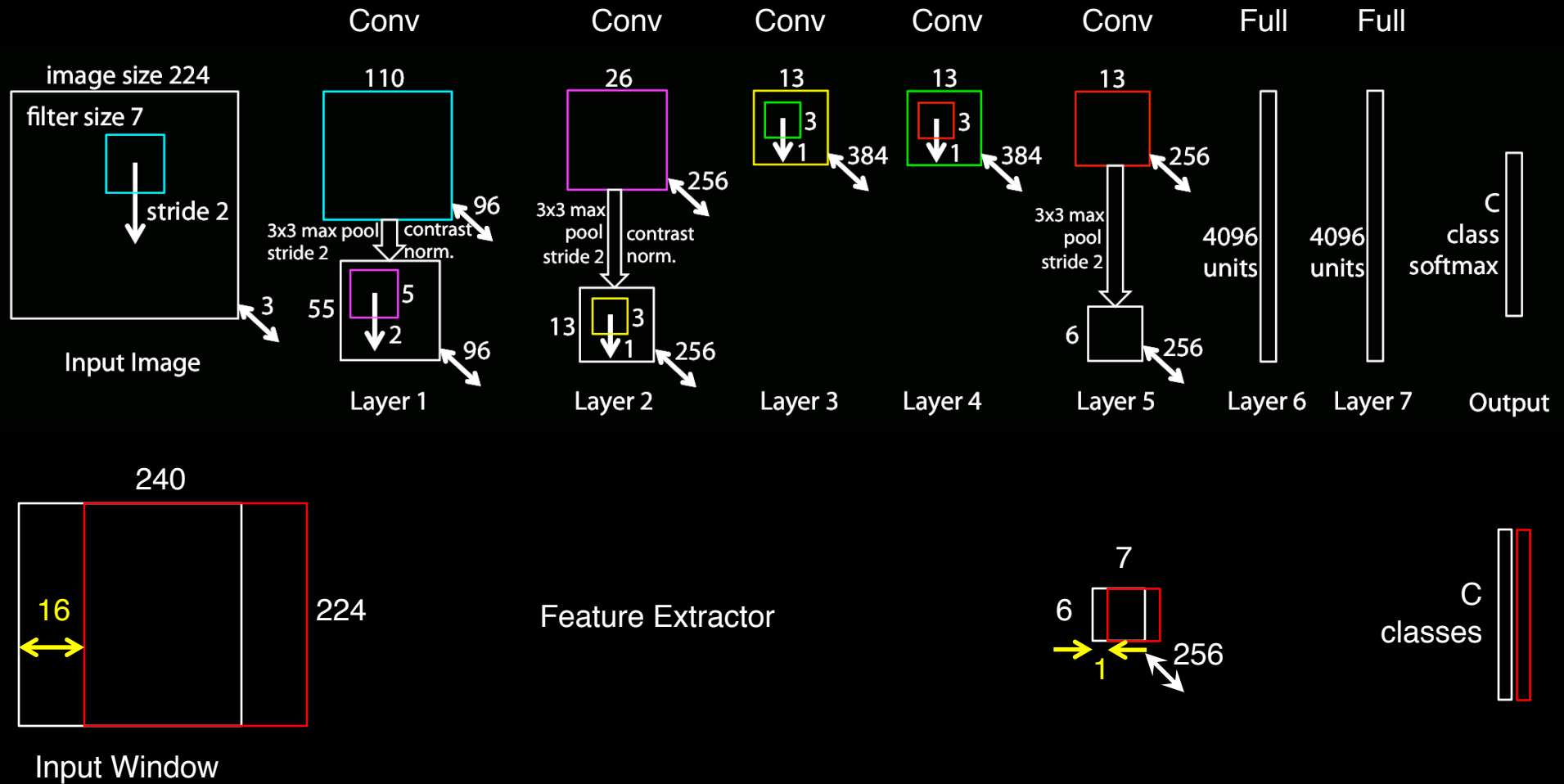
# Sliding Window with ConvNet

# Sliding Window with ConvNet
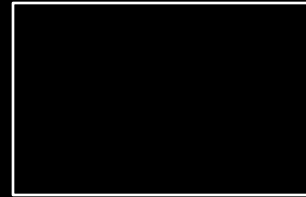
# Sliding Window with ConvNet



No need to compute two separate windows --- Just one big input window
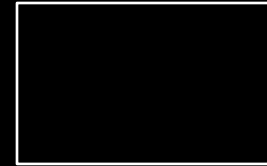
# Multi-Scale Sliding Window ConvNet

Feature
Maps

Class
Maps

256

C=1000

Feature
Extractor

256

Classifier

C=1000

256

C=1000

256
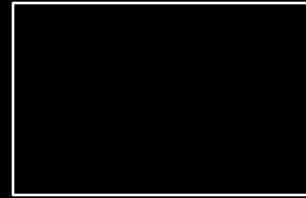
C=1000
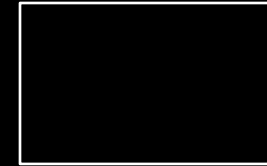
# Multi-Scale Sliding Window ConvNet



Feature
Maps

Bounding Box
Maps

256

4

Feature
Extractor

256
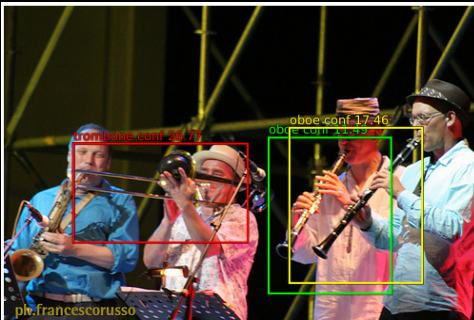
Regression
Network

4

256

4

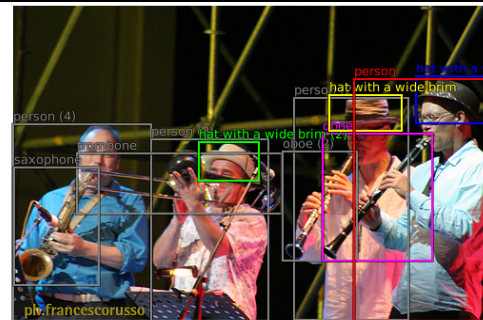256

4

# OverFeat – Output before NMS

# Overfeat Detection Results
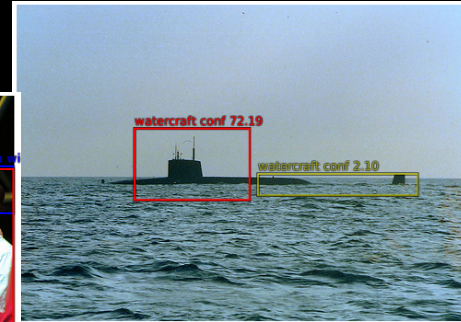
[Sermanet et al. ICLR 2014]

**Top predictions:**
trombone (confidence 26.8)
oboe (confidence 17.5)
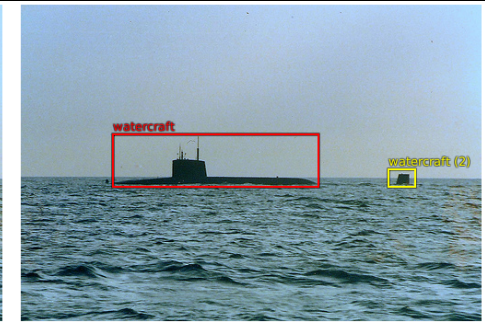oboe (confidence 11.5)

ILSVRC2012_val_00000614.JPEG

**Groundtruth:**
person
hat with a wide brim
hat with a wide brim (2)
hat with a wide brim (3)
oboe
oboe (2)
saxophone
trombone
person (2)
person (3)
person (4)

**Top predictions:**
watercraft (confidence 72.2)
watercraft (confidence 2.1)

ILSVRC2012_val_00000623.JPEG

**Groundtruth:**
watercraft
watercraft (2)

**Top predictions:**
tennis ball (confidence 3.5)
banana (confidence 2.4)
banana (confidence 2.1)
hotdog (confidence 2.0)
banana (confidence 1.9)

ILSVRC2012_val_00000320.JPEG

**Groundtruth:**
strawberry
strawberry (2)
strawberry (3)
strawberry (4)
strawberry (5)
strawberry (6)
strawberry (7)
strawberry (8)
strawberry (9)
strawberry (10)
apple
apple (2)
apple (3)

**Top predictions:**
microwave (confidence 5.6)
refrigerator (confidence 2.5)

ILSVRC2012_val_00000519.JPEG

**Groundtruth:**
bowl
microwave

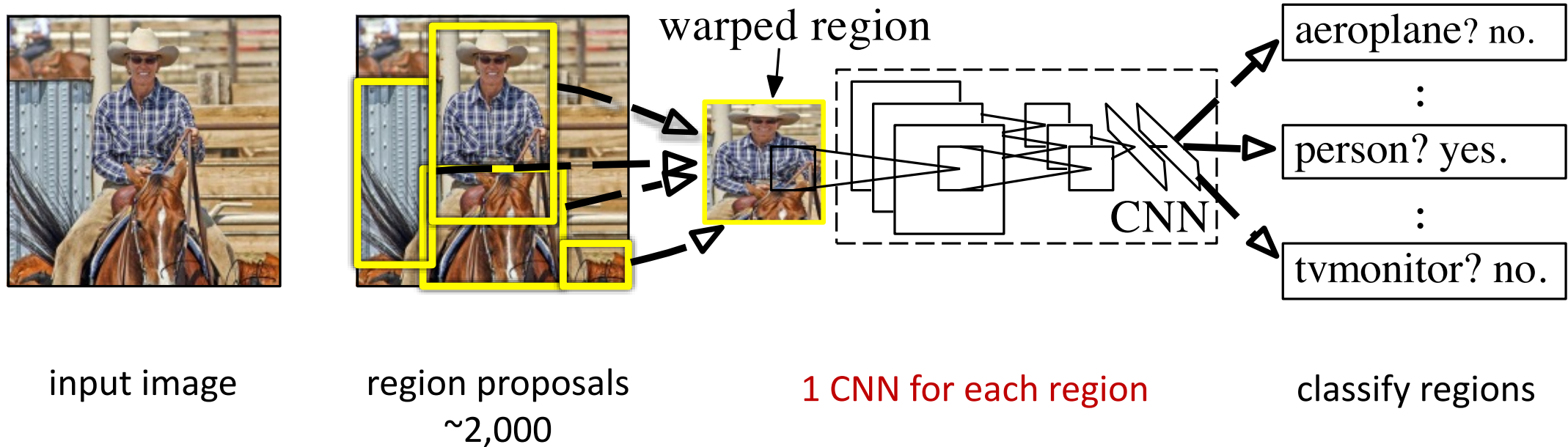# Two General Approaches

1. Examine very position / scale
   – E.g. Overfeat: Integrated recognition, localization and detection using convolutional networks, Sermanet et al., ICLR 2014

2. Use some kind of proposal mechanism to attend to a set of possible regions
   – E.g. Region-CNN [Rich feature hierarchies for accurate object detection and semantic segmentation, Girshick et al., CVPR 2014]

# Object Detection: R-CNN

warped region

CNN

aeroplane? no.

:

person? yes.

:

tvmonitor? no.

input image

region proposals
~2,000

1 CNN for each region

classify regions

**R**egion-based **CNN** pipeline

Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014

# Object Detection: R-CNN

- R-CNN



feature

feature

feature

feature

CNN

CNN

CNN

CNN

pre-computed
Regions-of-Interest
(RoIs)

image

End-to-End
training

Girshick, Donahue, Darrell, Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. CVPR 2014
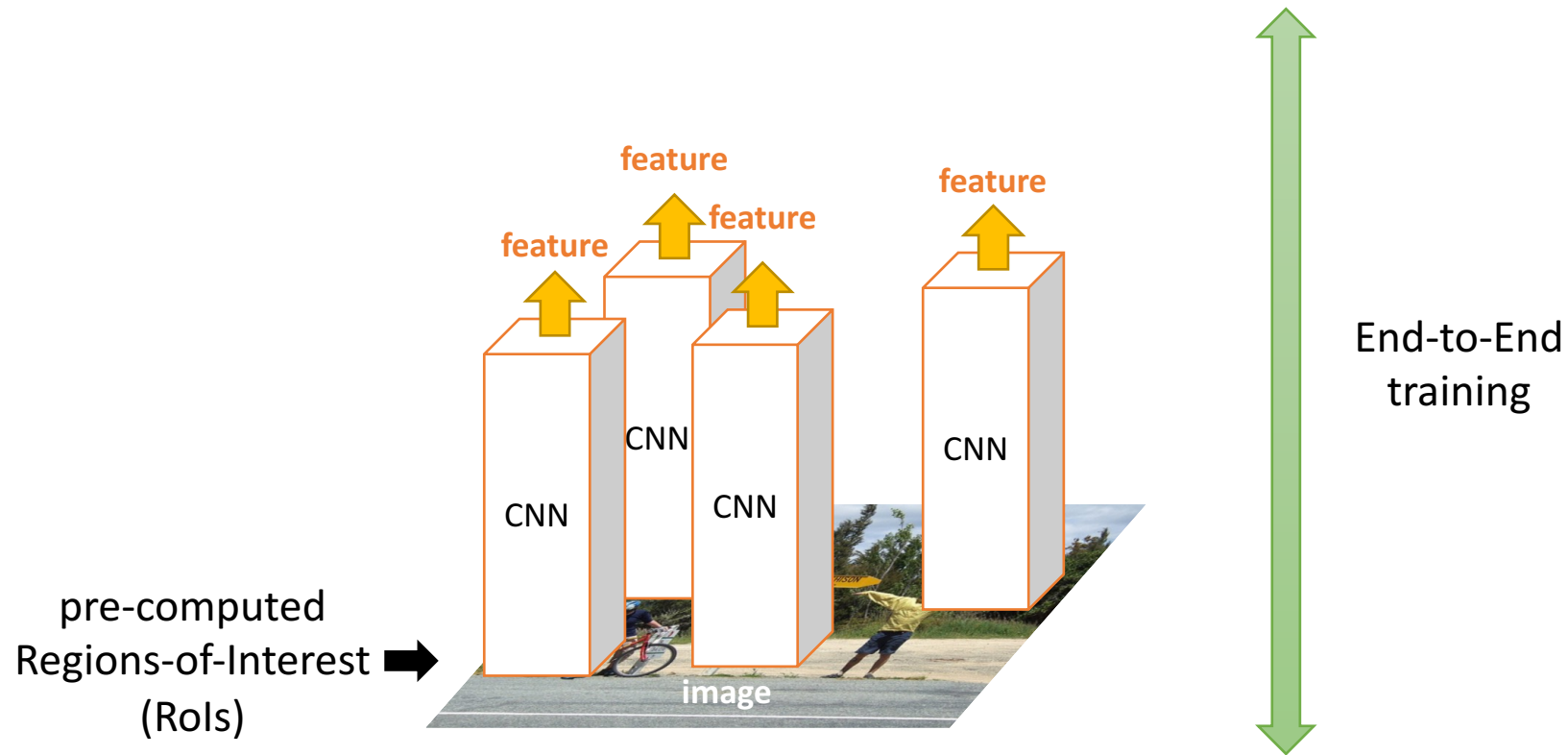
# Object Detection: Fast R-CNN

- Fast R-CNN



pre-computed
Regions-of-Interest
(RoIs) ➡

**feature**   **feature**   **feature**

RoI pooling

shared conv
layers

CNN

image
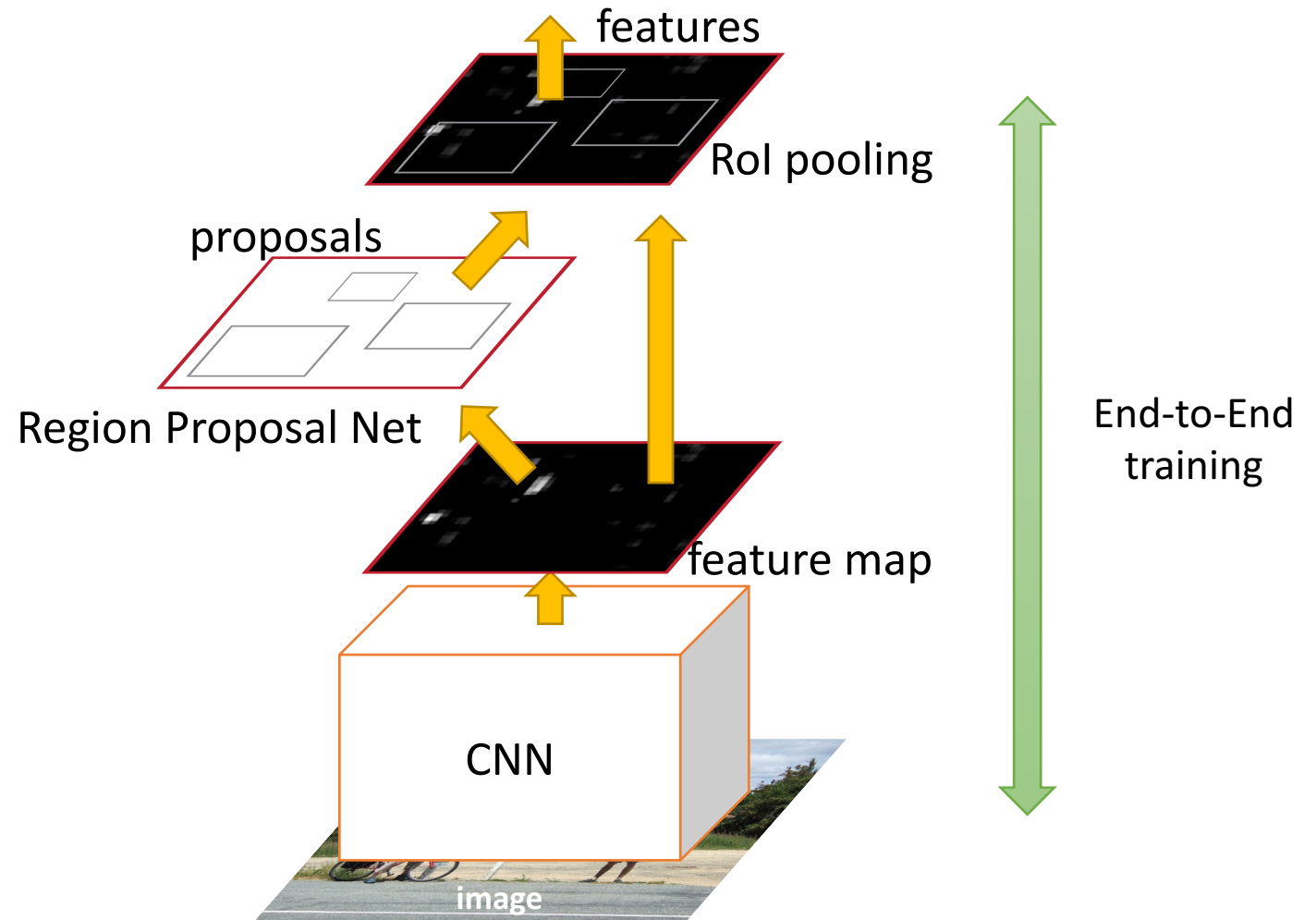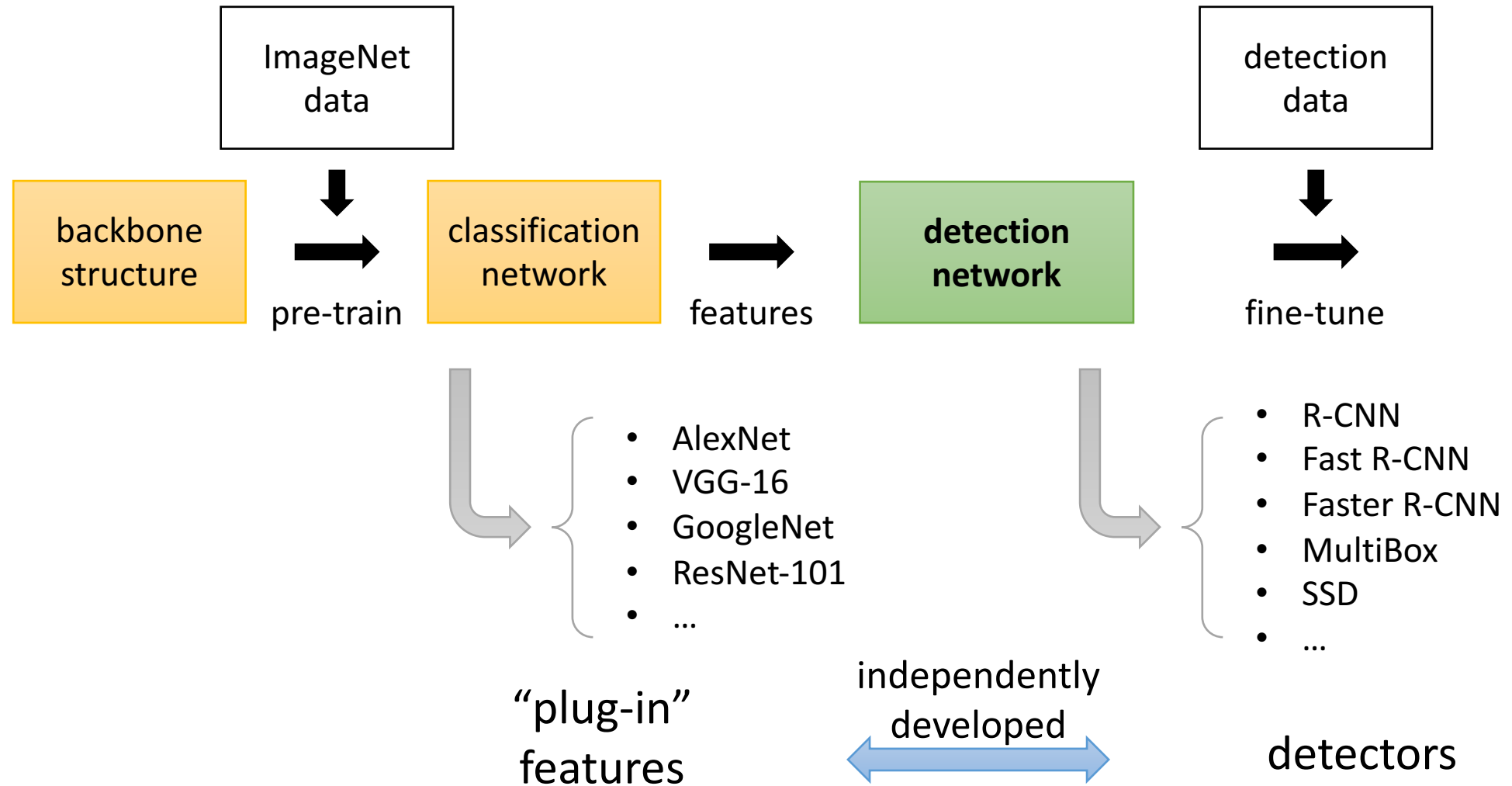
End-to-End
training

Girshick. Fast R-CNN. ICCV 2015

# Object Detection: Faster R-CNN

- Faster R-CNN
  - Solely based on CNN
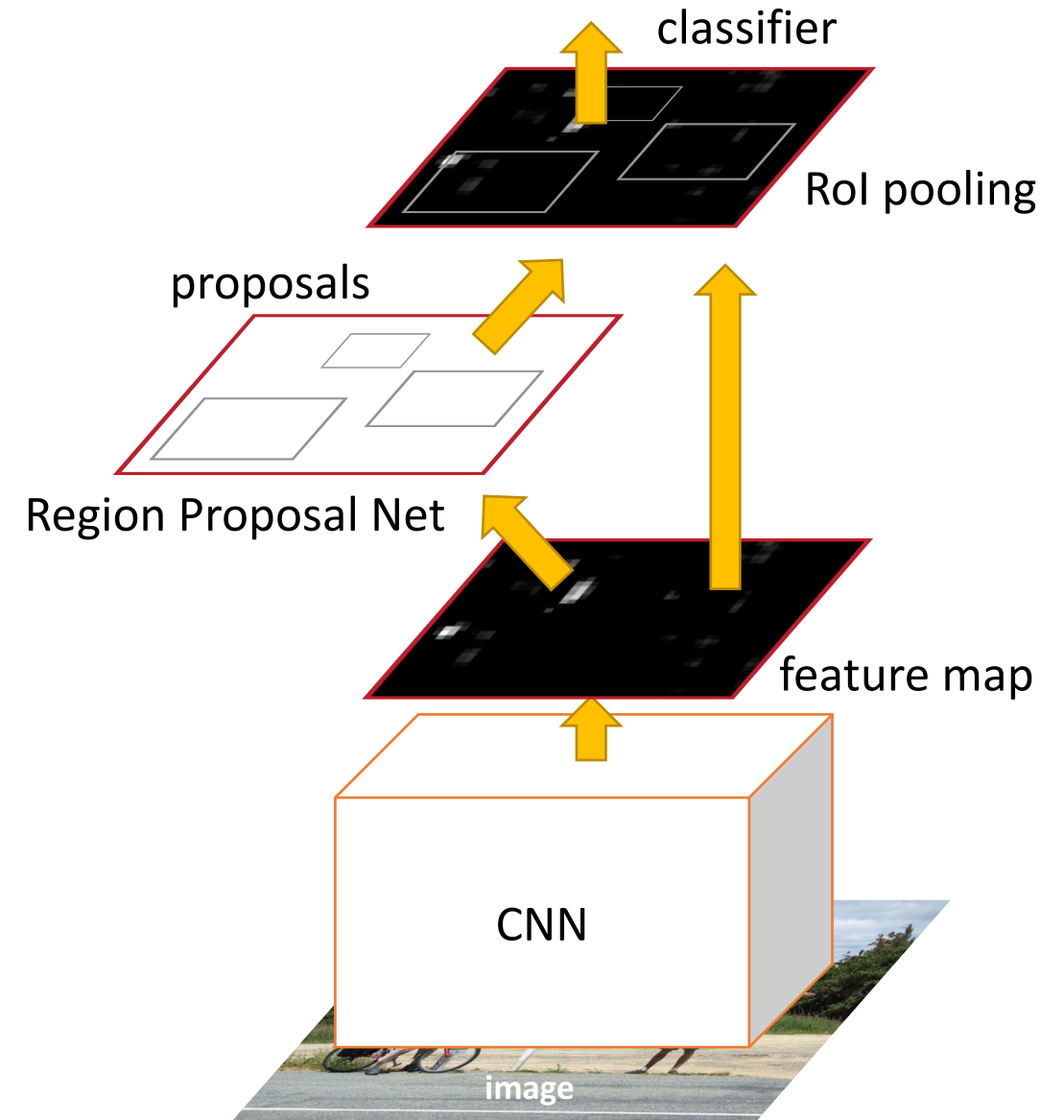  - No external modules
  - Each step is end-to-end

features

RoI pooling

proposals

Region Proposal Net

feature map

CNN

image

End-to-End training

Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

# Object Detection

# Object Detection

- Simply "Faster R-CNN + ResNet"

| Faster R-CNN baseline | mAP@.5 | mAP@.5:.95 |
|---|---|---|
| VGG-16 | 41.5 | 21.5 |
| ResNet-101 | **48.4** | **27.2** |

COCO detection results

**ResNet-101 has 28% relative gain vs VGG-16**



classifier

RoI pooling

proposals

Region Proposal Net

feature map

CNN

image

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

# Object Detection

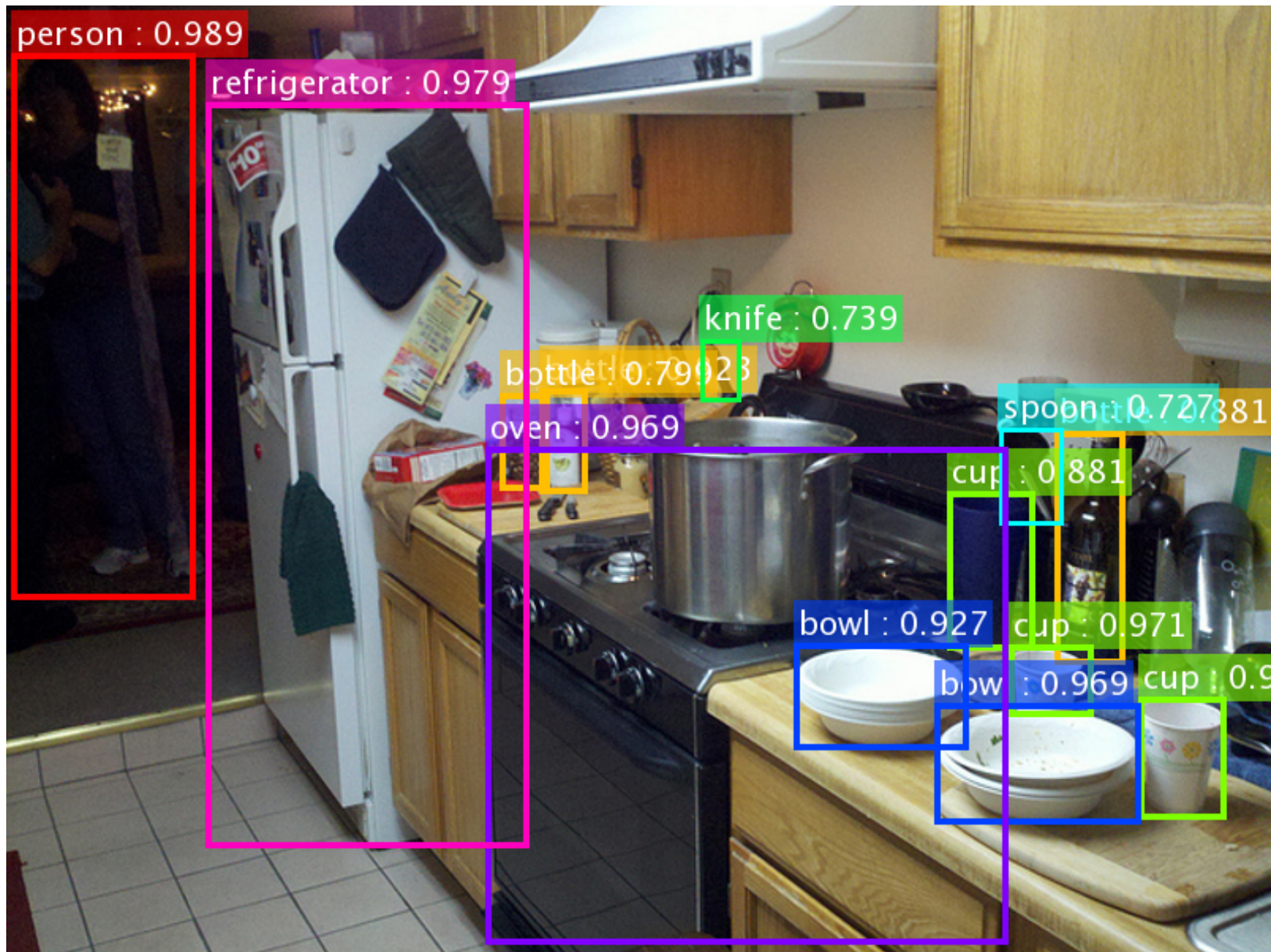- RPN <span style="color:red">learns</span> proposals by extremely deep nets
  - We use <span style="color:red">only 300 proposals</span> (no hand-designed proposals)

- Add components:
  - Iterative localization
  - Context modeling
  - Multi-scale testing

- All components are based on CNN features; all steps are end-to-end

- All benefit <span style="color:red">more</span> from <span style="color:red">deeper</span> features – cumulative gains!

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.
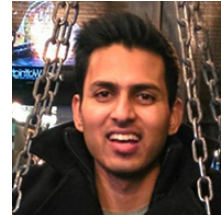
ResNet's object detection result on COCO

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.
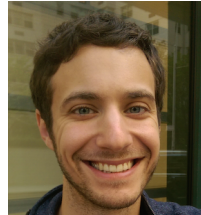
Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

Results on real video. Models trained on MS COCO (80 categories).
(frame-by-frame; no temporal processing)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

# FAIR COCO Object Detection

Sergey Zagoruyko*, Tsung-Yi Lin*, Pedro Pinheiro*, Adam Lerer, Sam Gross, Soumith Chintala, Piotr Dollár

(*equal contribution)

# Results

| | AP bbox | AP small | AP medium | AP large | AR max=100 | AP segm |
|---|---|---|---|---|---|---|
| MSRA | 0.373 | 0.183 | 0.419 | 0.524 | 0.491 | 0.282 |
| *FAIRCNN* | 0.335 | 0.139 | 0.378 | 0.477 | 0.485 | 0.251 |
| ION | 0.310 | 0.123 | 0.332 | 0.447 | 0.457 | |
| FastRCNN | 0.197 | 0.035 | 0.188 | 0.346 | 0.298 | |

66% improvement over FastRCNN baseline

# Overview

I.  DeepMask segmentation proposals [Pinheiro NIPS 15]

    + iterative localization

    + top-down refinement

II. Fast R-CNN object detector [Girshick ICCV 15]

    + foveal context regions

    + modified loss function

    + skip connections

    + ensembling

# I. Deep MASK Object Proposals

# DeepMask Framework

**Model:**



x: 3x224x224

VGG

512x14x14

1x1 conv

512x14x14    512x1x1    56x56

$f_{segm}(x)$: 224x224

2x2 pool

512x7x7    512x1x1    1024x1x1

$f_{score}(x)$: 1x1

# DeepMask Framework

**Model:**



x: 3x224x224

VGG

512x14x14

Segmentation Mask

1x1 conv

512x14x14    512x1x1    56x56

$f_{segm}(x)$: 224x224

2x2 pool

512x7x7    512x1x1    1024x1x1

$f_{score}(x)$: 1x1

# DeepMask Framework

**Model:**



x: 3x224x224

VGG

512x14x14

1x1 conv

512x14x14    512x1x1    56x56

f_segm(x): 224x224

2x2 pool

512x7x7    512x1x1    1024x1x1

f_score(x): 1x1

'Objectness' score

# DeepMask Framework

**Model:**

# DeepMask Framework

**Model:**

# DeepMask Framework

**Model:**

# DeepMask Framework

**Model:**

# Single Scale Inference

image



scores



masks

# Single Scale Inference

image



scores



masks

# *New:* Iterative Localization (+1.0 AP)

# *New:* Top-Down Refinement (+0.7 AP)

# Proposal Quality (boxes)

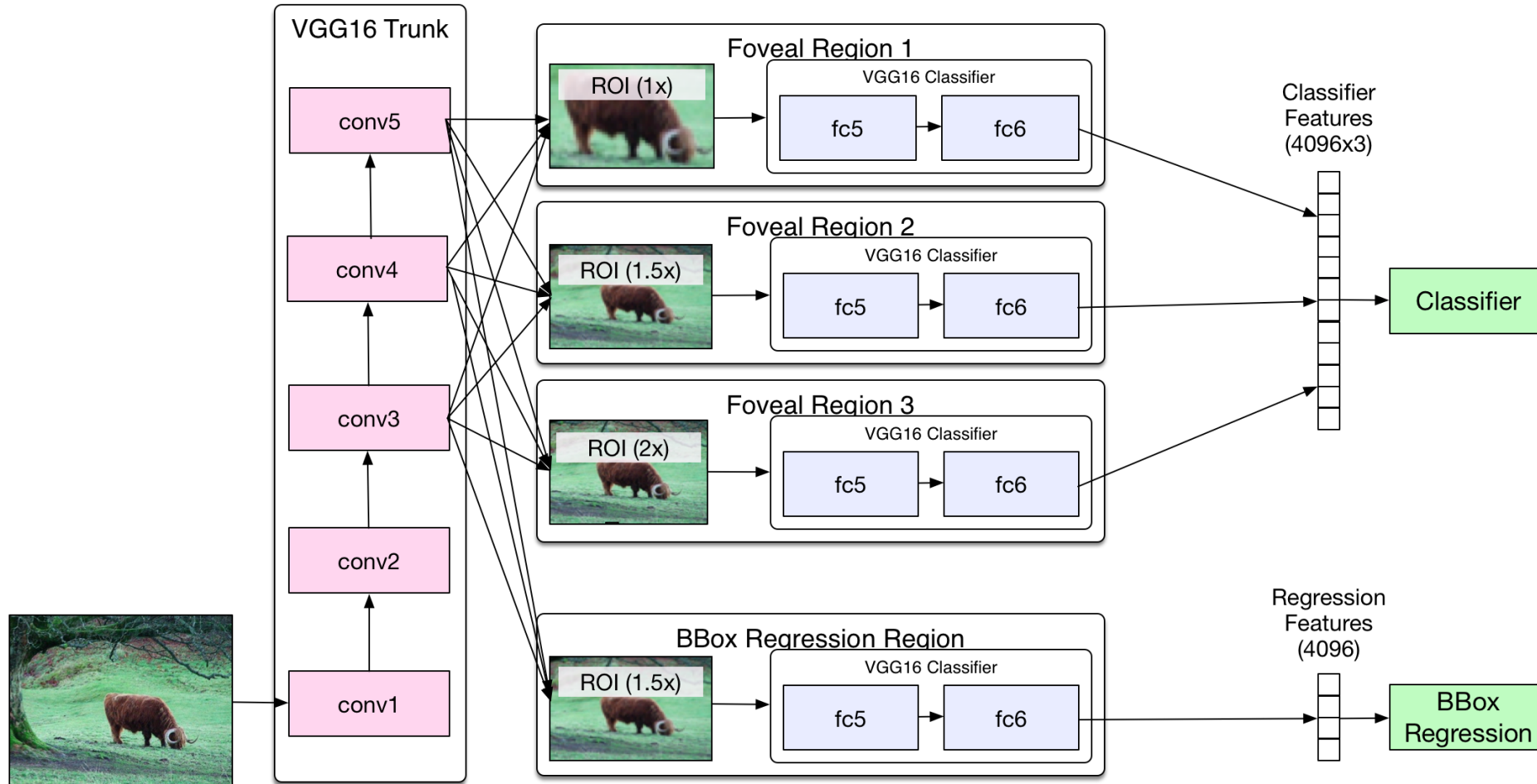# DeepMask Object Proposals

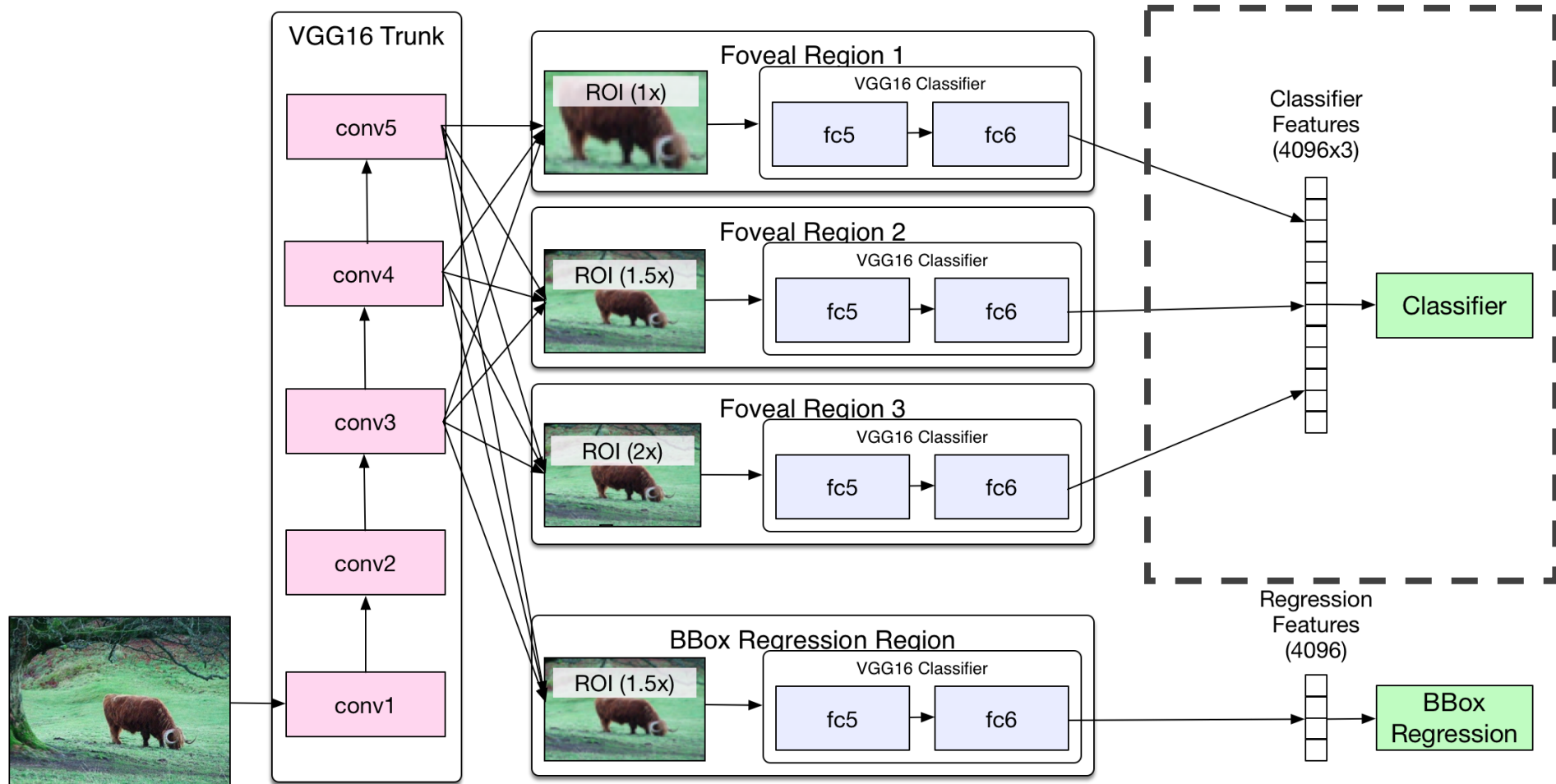# II. Classification framework

- Fast R-CNN setup [Girshick, ICCV15]

- Fast R-CNN setup [Girshick, ICCV15]
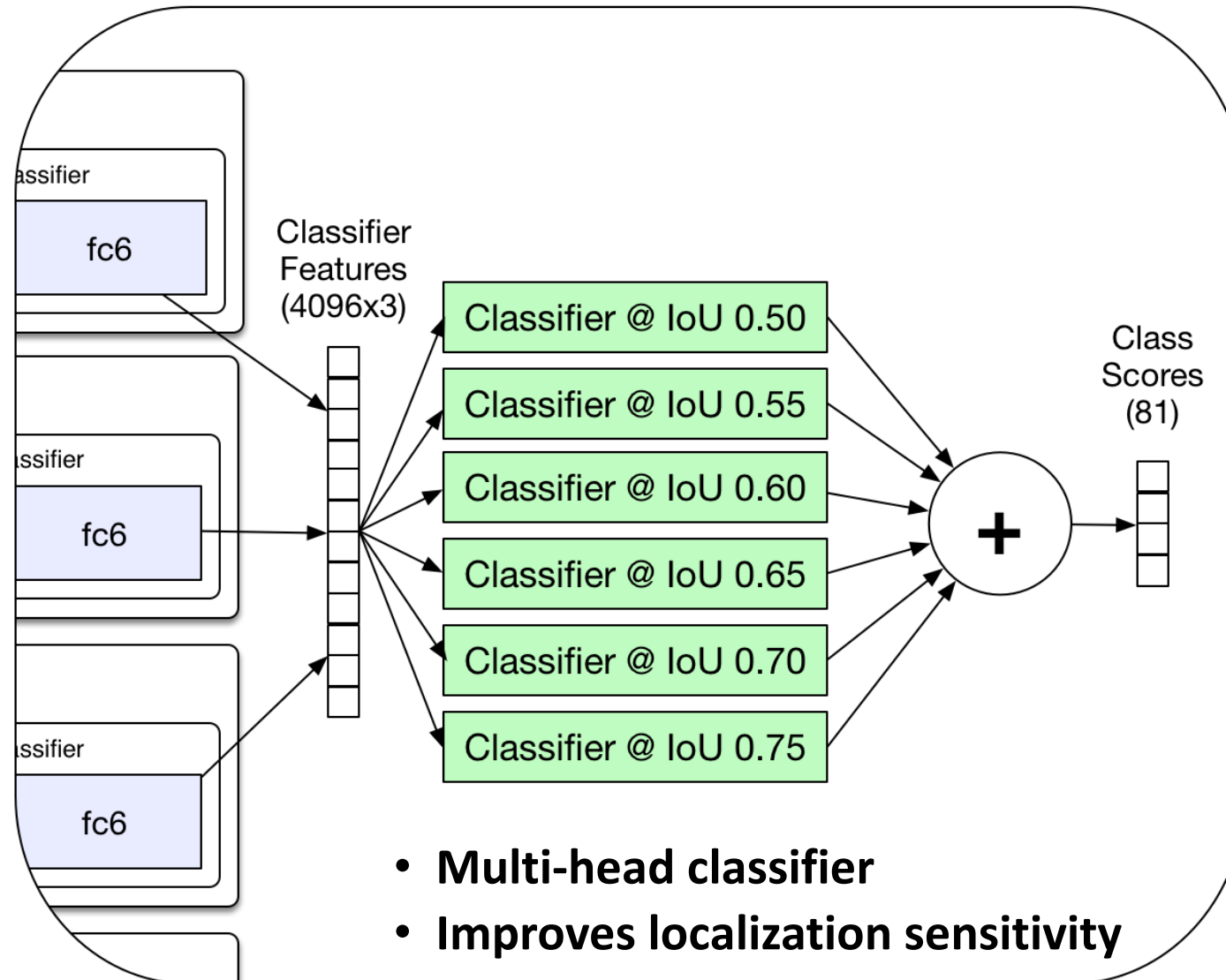- Foveal structure [inspired by Gidaris & Komodakis, ICCV15] (+2 AP)

- Fast R-CNN setup [Girshick, ICCV15]
- Foveal structure [inspired by Gidaris & Komodakis, ICCV15] (+2 AP)
- Skip connections (+1 AP)

- Fast R-CNN setup [Girshick, ICCV15]
- Foveal structure [inspired by Gidaris & Komodakis, ICCV15] (+2 AP)
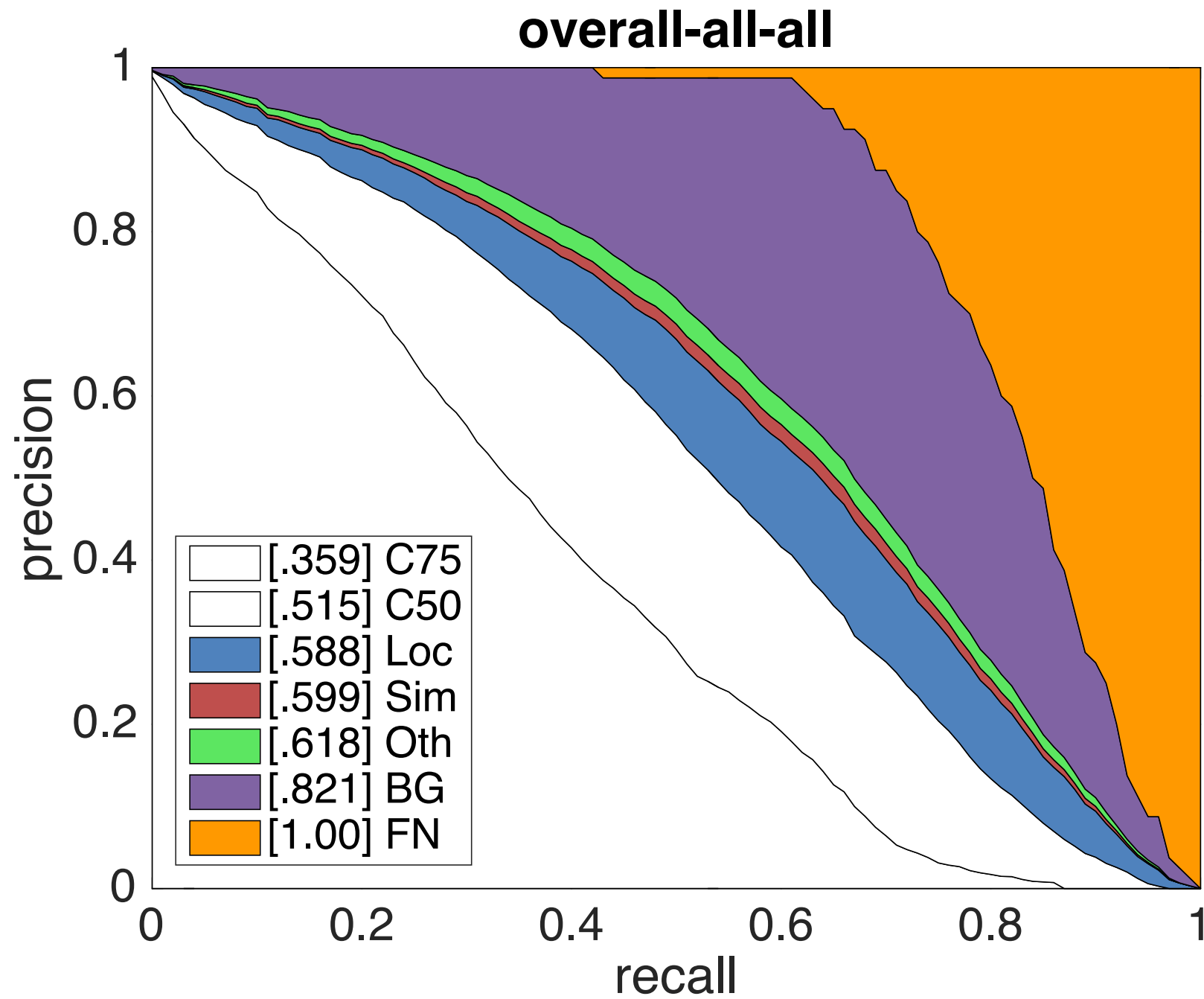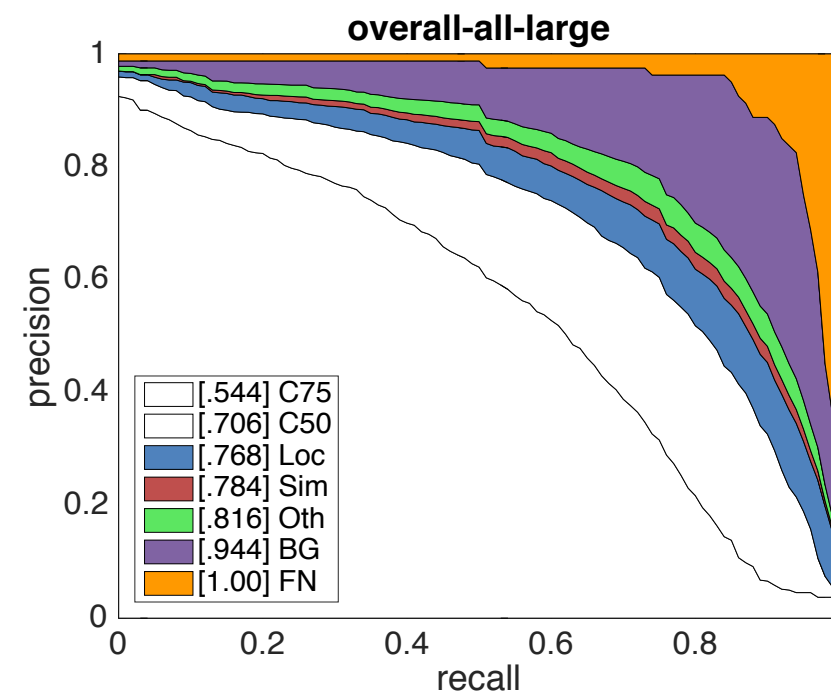- Skip connections (+1 AP)
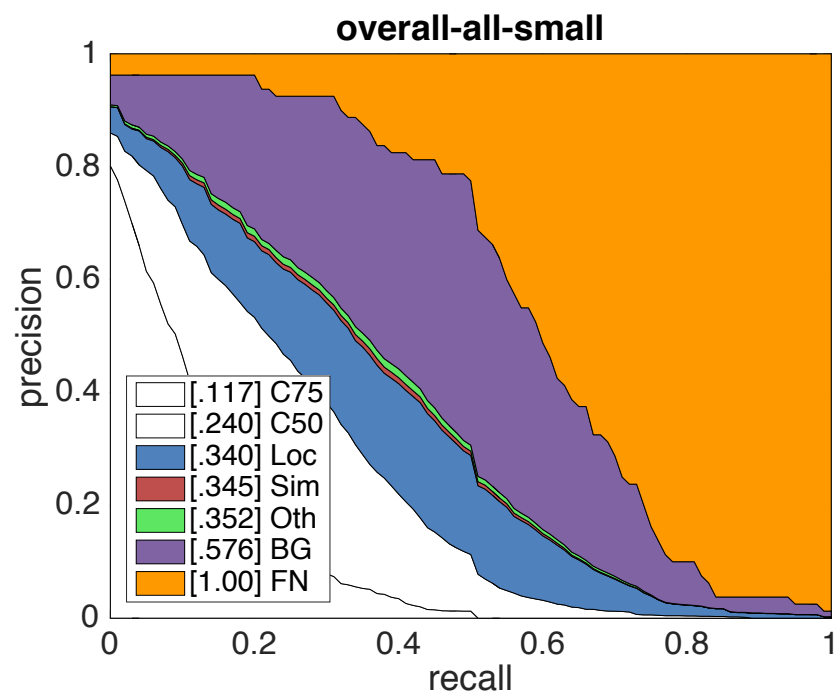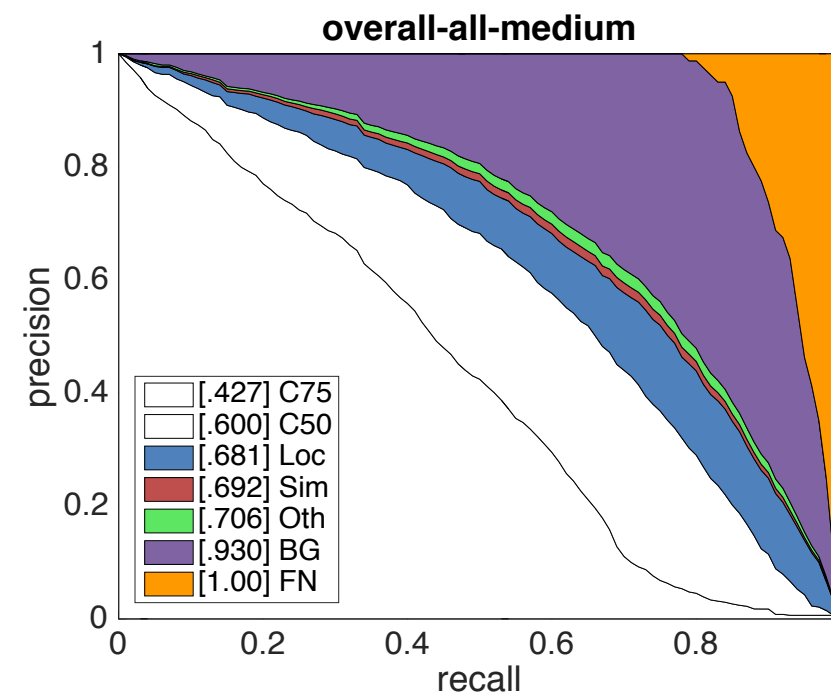
# Multi-threshold Loss (+1.5 AP)



- **Multi-head classifier**
- **Improves localization sensitivity**

# Inference

|  |  |
|---|---|
| Base Model | **30.1 AP** |
| + horizontal flip | **31.1 AP** |
| + ROI Pooling '2 crop' | **32.1 AP** |
| + 7-model Ensemble | **33.5 AP** |

**overall-all-all**

Legend:
- [.359] C75
- [.515] C50
- [.588] Loc
- [.599] Sim
- [.618] Oth
- [.821] BG
- [1.00] FN

x-axis: recall
y-axis: precision

| AP | |
|---|---|
| Small | 0.139 |
| Medium | 0.378 |
| Large | 0.477 |

**overall-all-medium**

[.427] C75
[.600] C50
[.681] Loc
[.692] Sim
[.706] Oth
[.930] BG
[1.00] FN

**overall-all-small**

[.117] C75
[.240] C50
[.340] Loc
[.345] Sim
[.352] Oth
[.576] BG
[1.00] FN

**overall-all-large**

[.544] C75
[.706] C50
[.768] Loc
[.784] Sim
[.816] Oth
[.944] BG
[1.00] FN

# Segmentation Examples



DeepMask

*Proposal BBoxes*

FastRCNN

*Scored BBoxes*

DeepMask

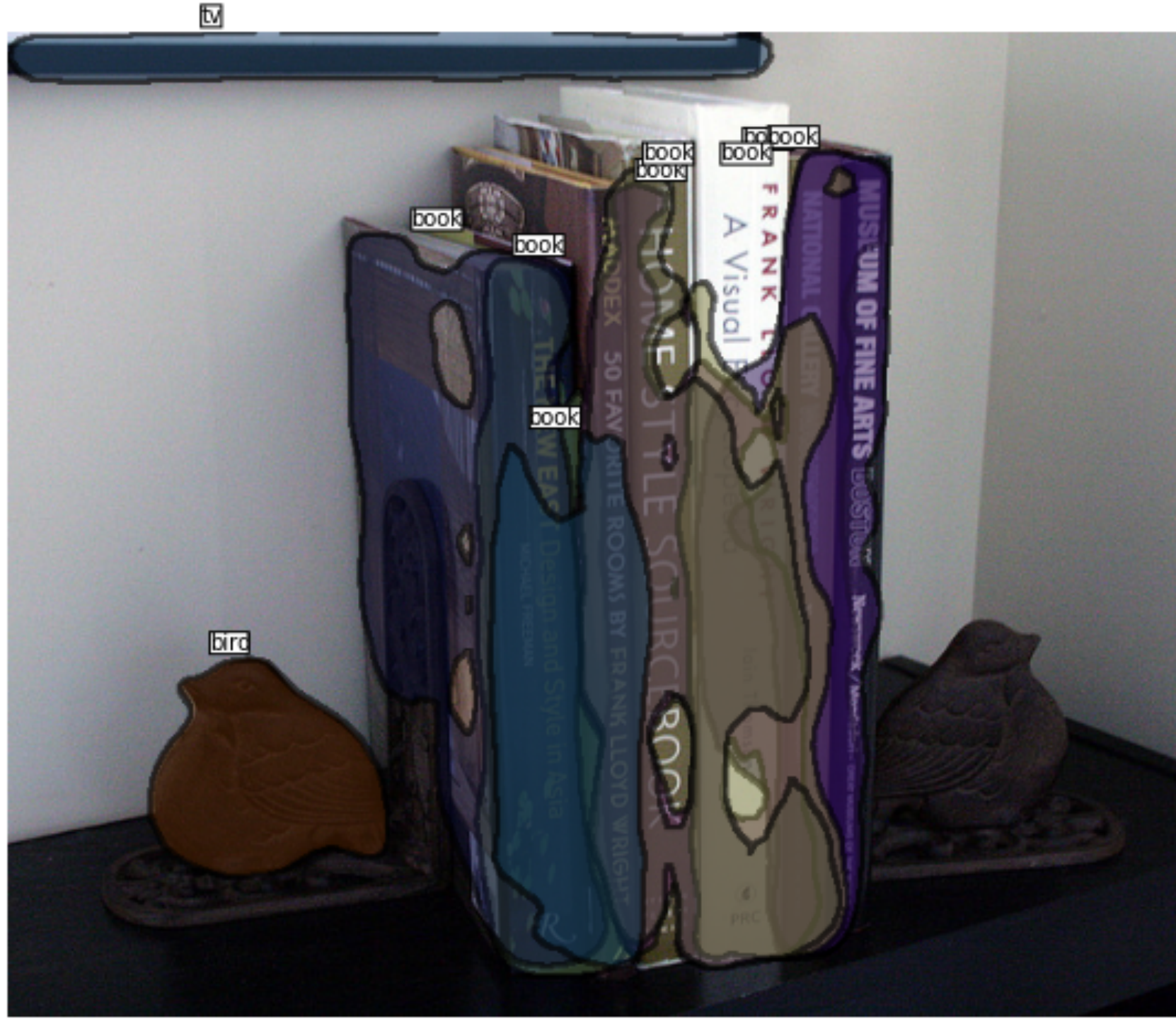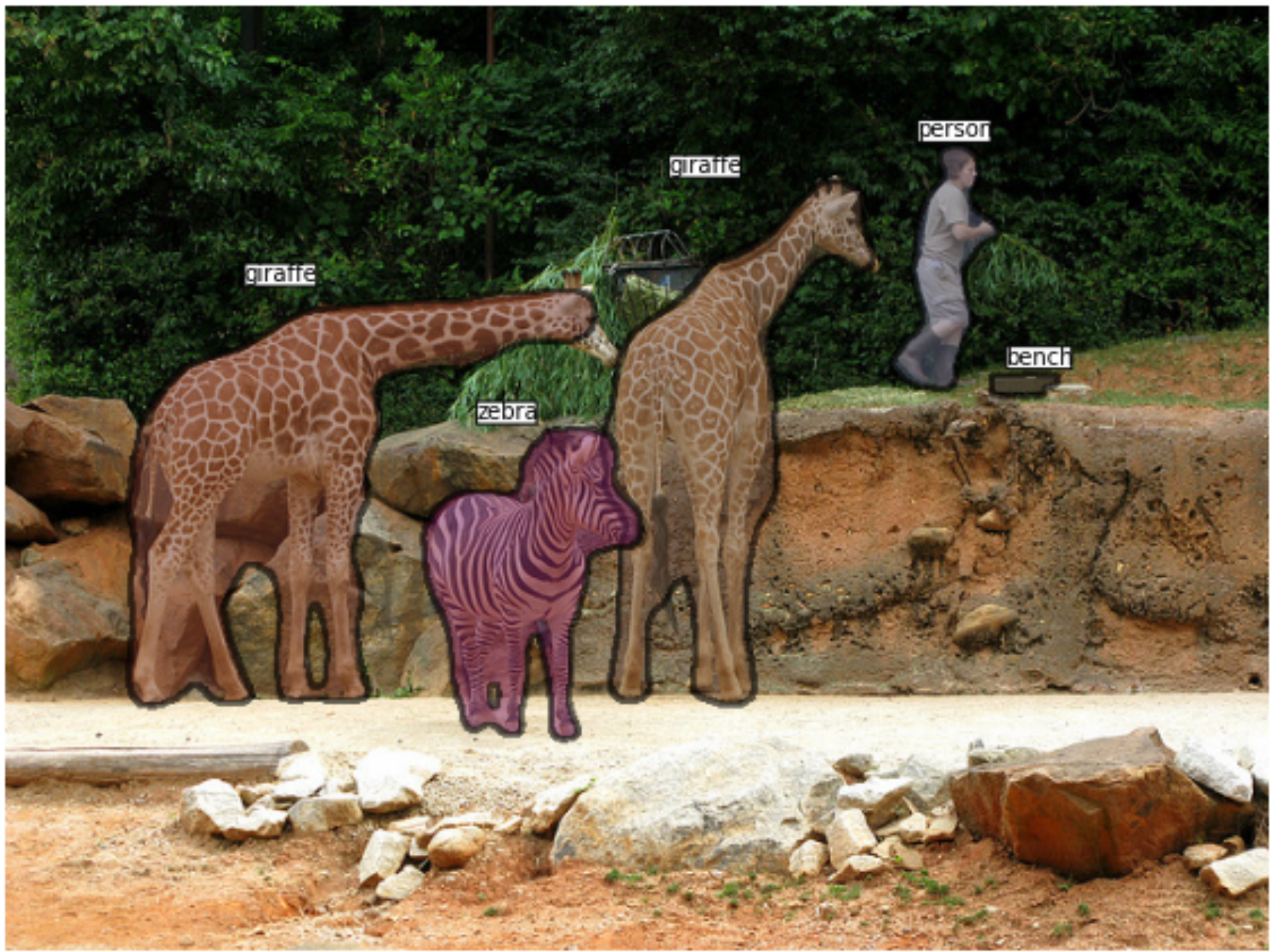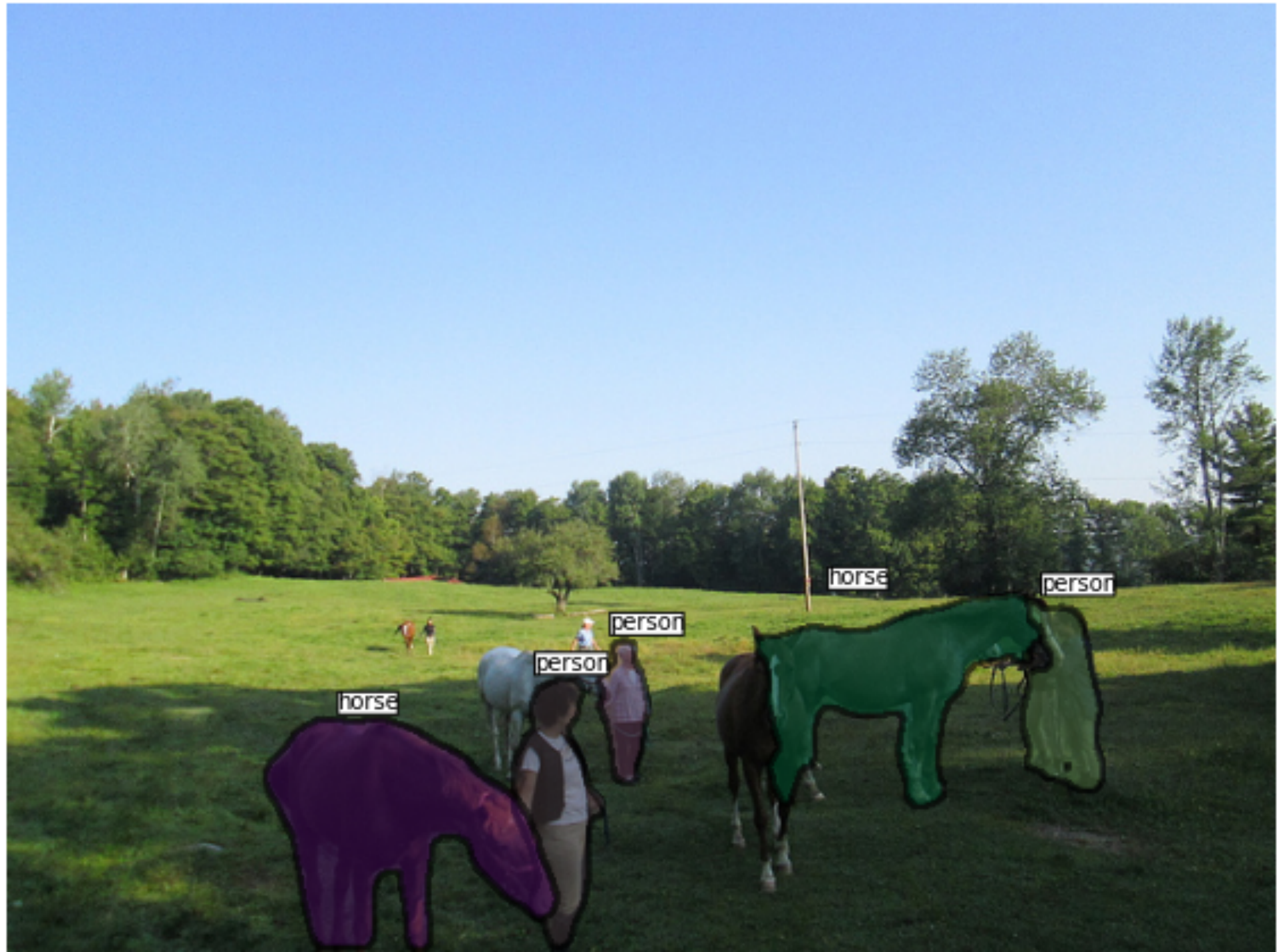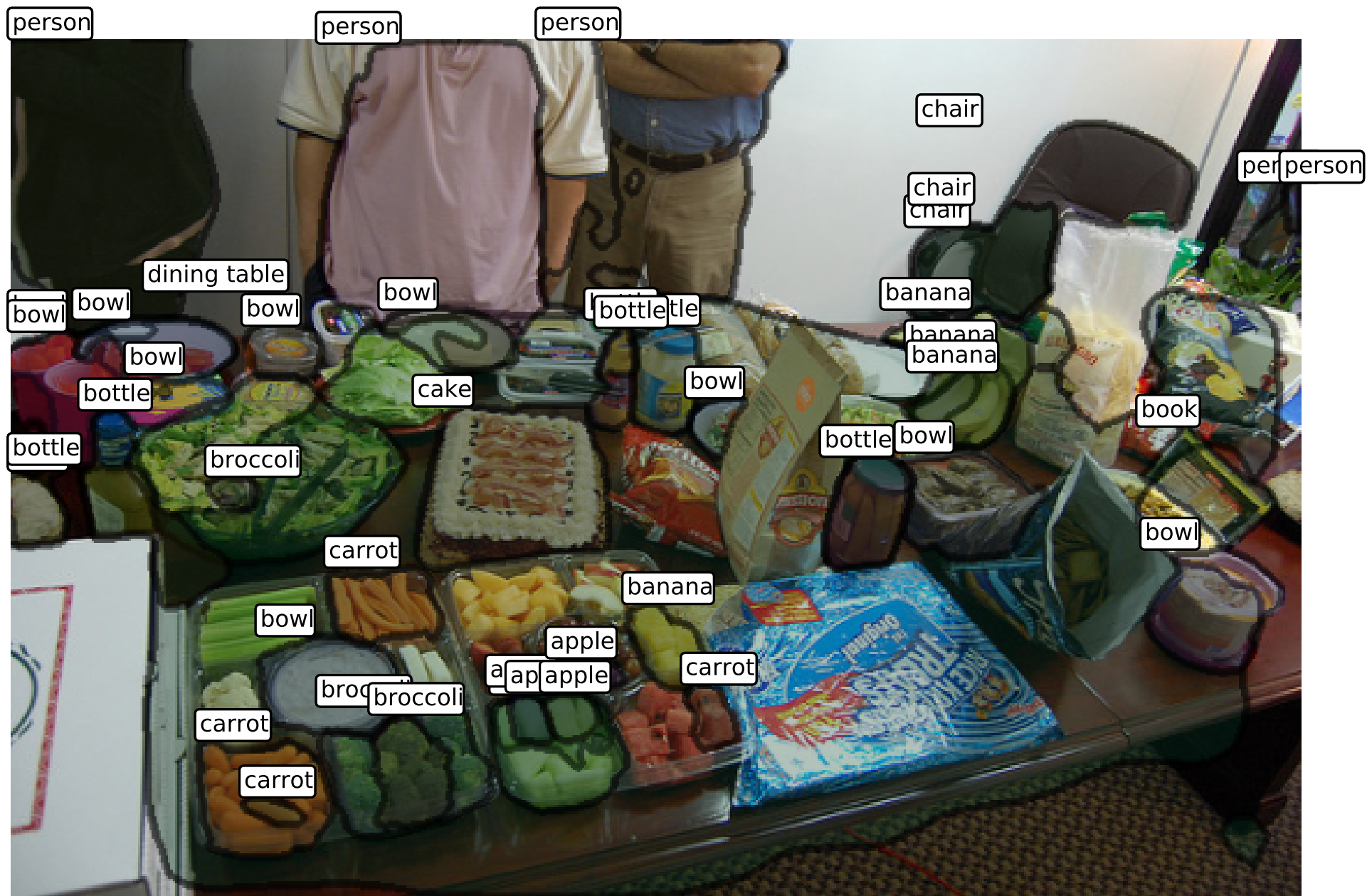*Scored Segments*

557556

sink

sink

# Future Directions

- most room for improvement:
  - background confusion (FP/FN)
  - small objects
- more effective use of context
- fast / proposal-free detection

# Questions?