



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

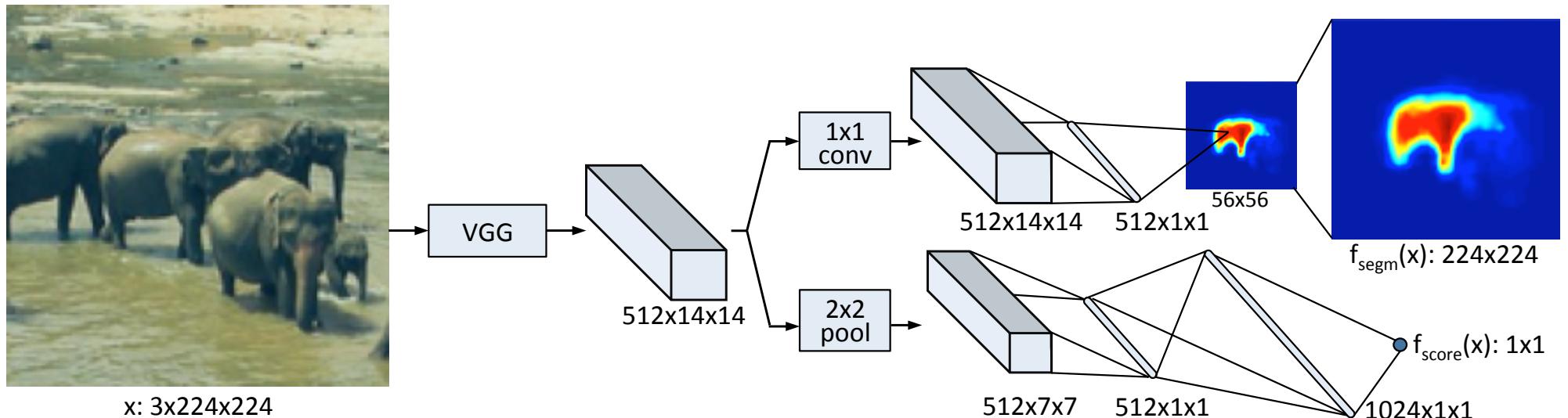
- I. DeepMask segmentation proposals [Pinheiro NIPS 15]
 - + iterative localization
 - + top-down refinement

- II. Fast R-CNN object detector [Girshick ICCV 15]
 - + foveal context regions
 - + modified loss function
 - + skip connections
 - + ensembling

I. DEEP MASK OBJECT PROPOSALS

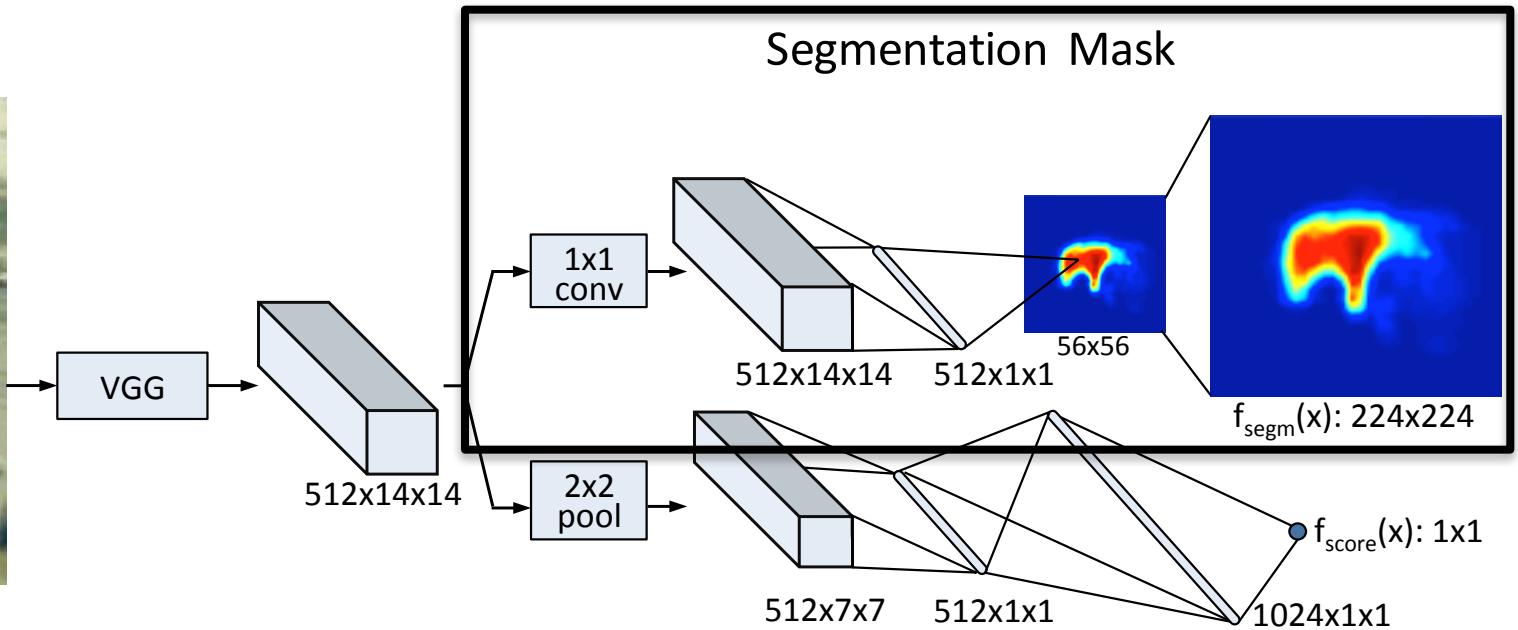
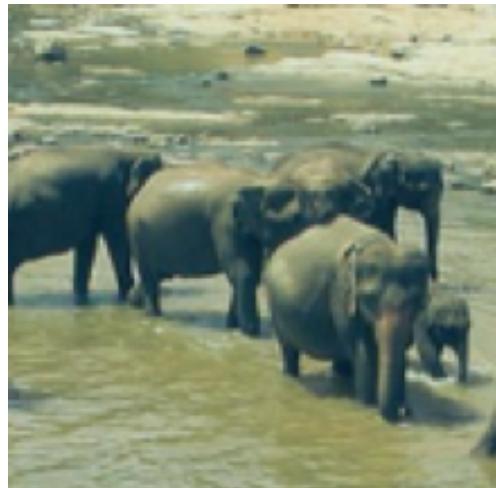
DeepMask Framework

Model:



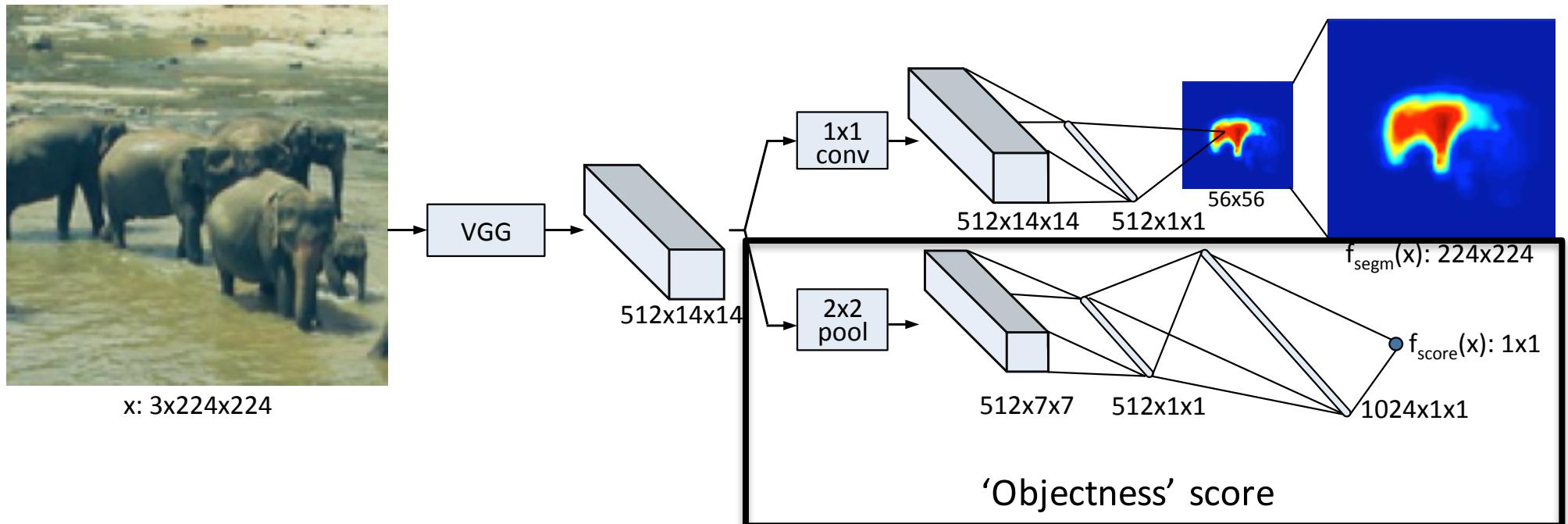
DeepMask Framework

Model:



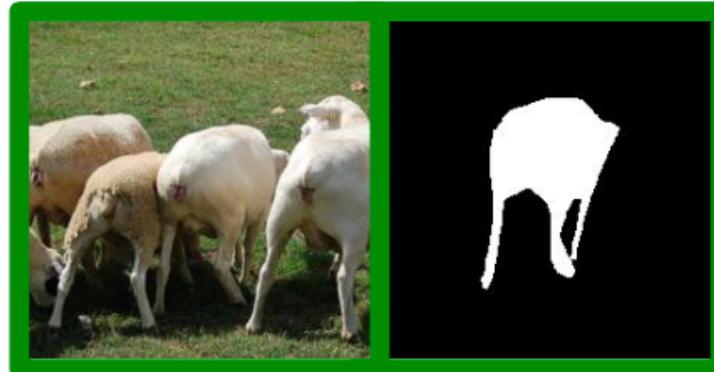
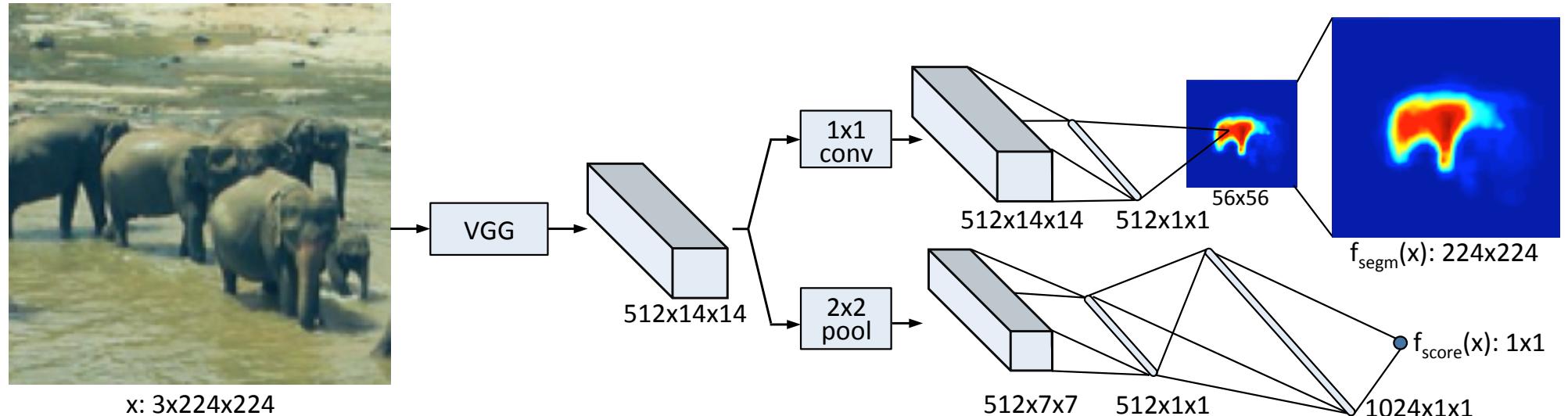
DeepMask Framework

Model:



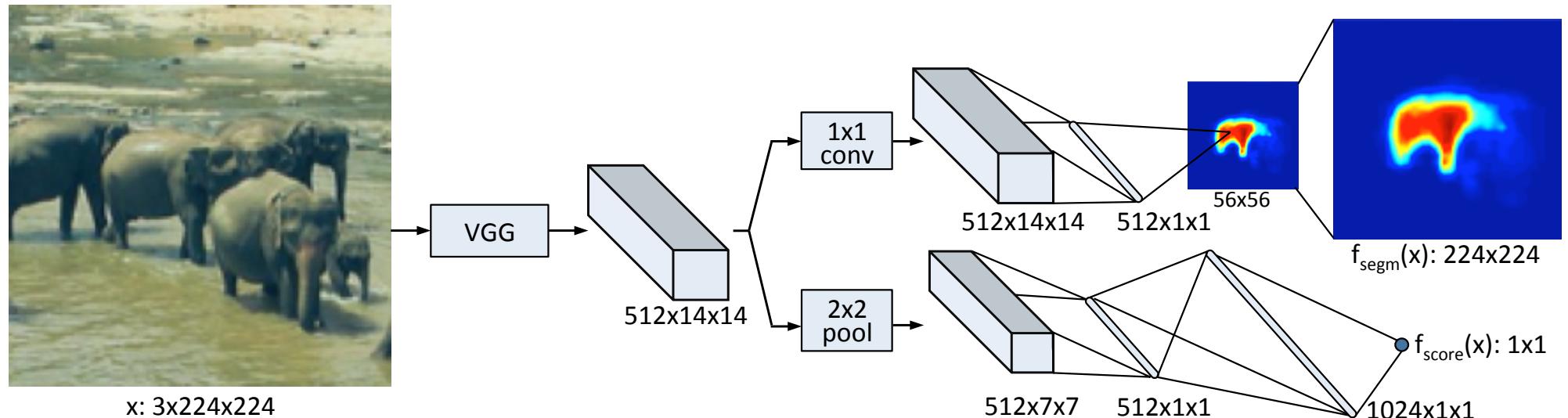
DeepMask Framework

Model:



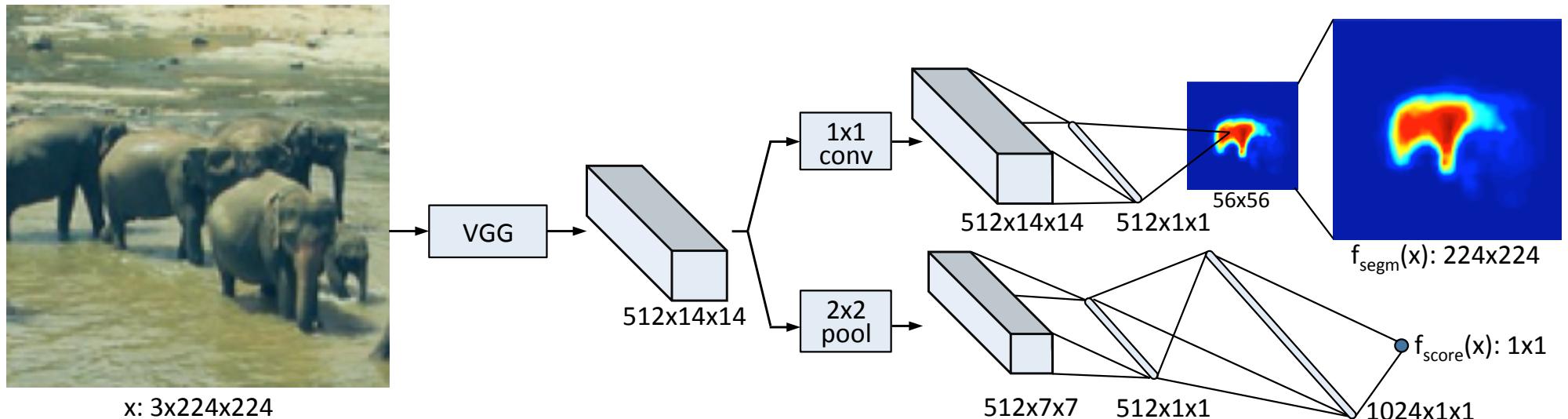
DeepMask Framework

Model:



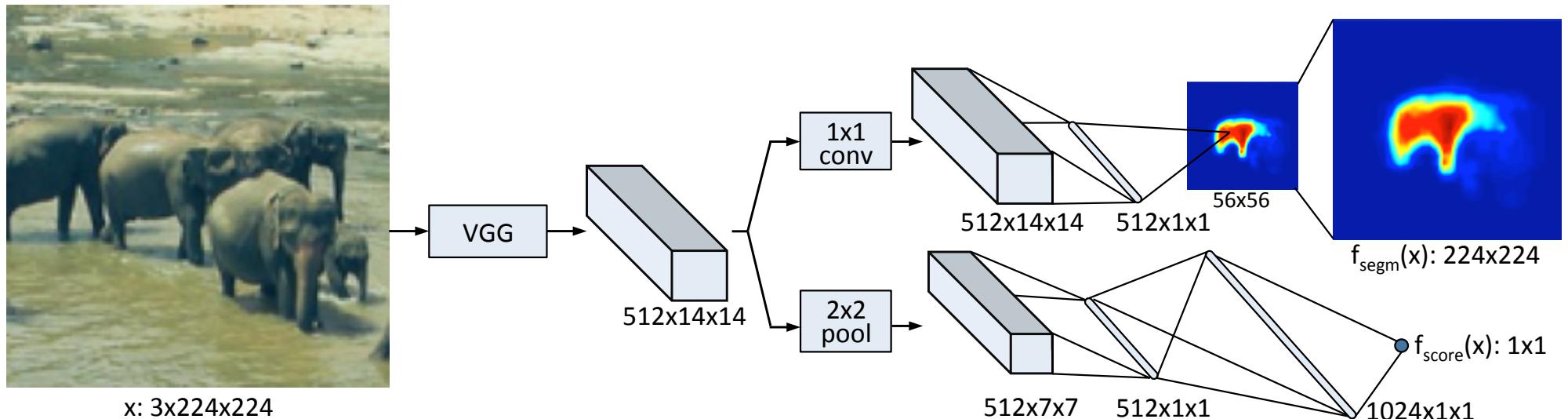
DeepMask Framework

Model:



DeepMask Framework

Model:

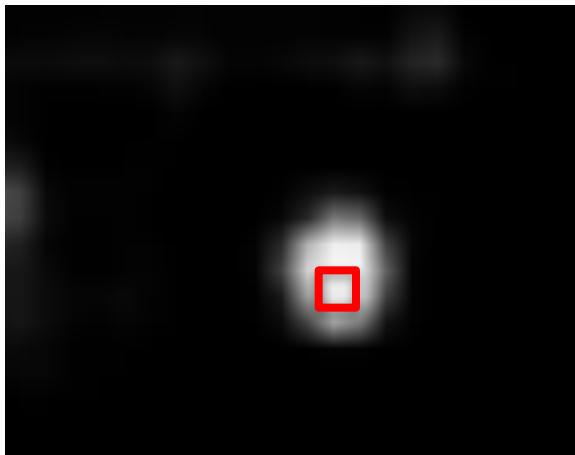


Single Scale Inference

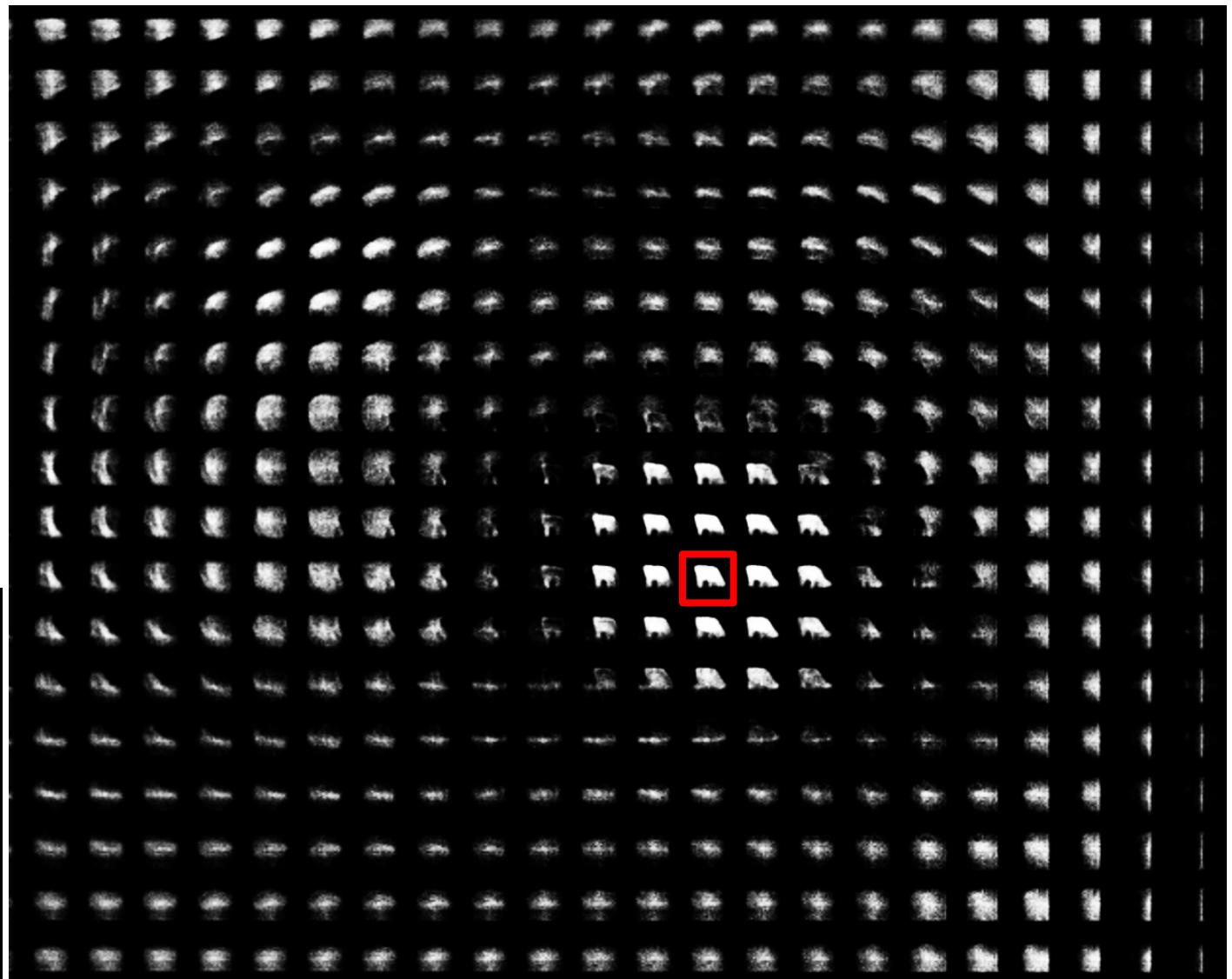
image



scores

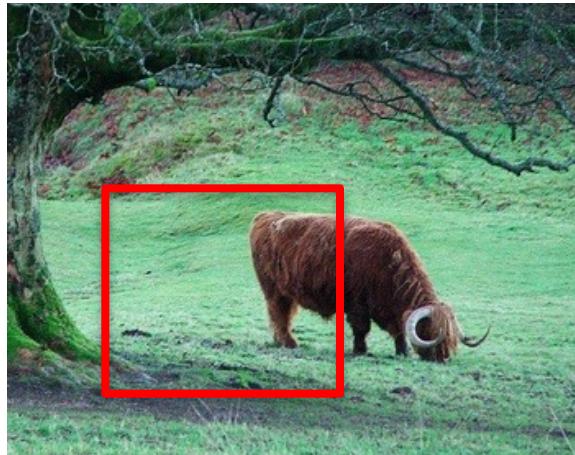


masks

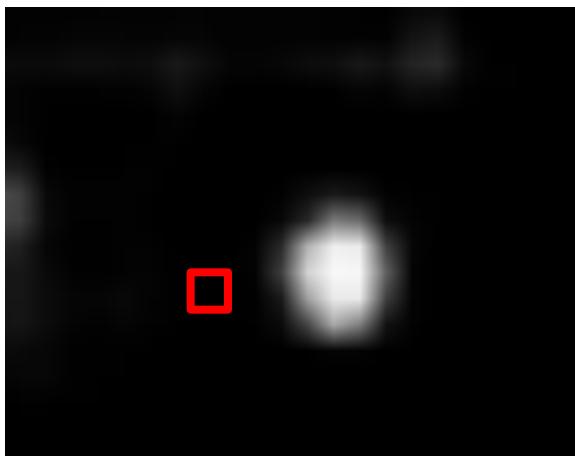


Single Scale Inference

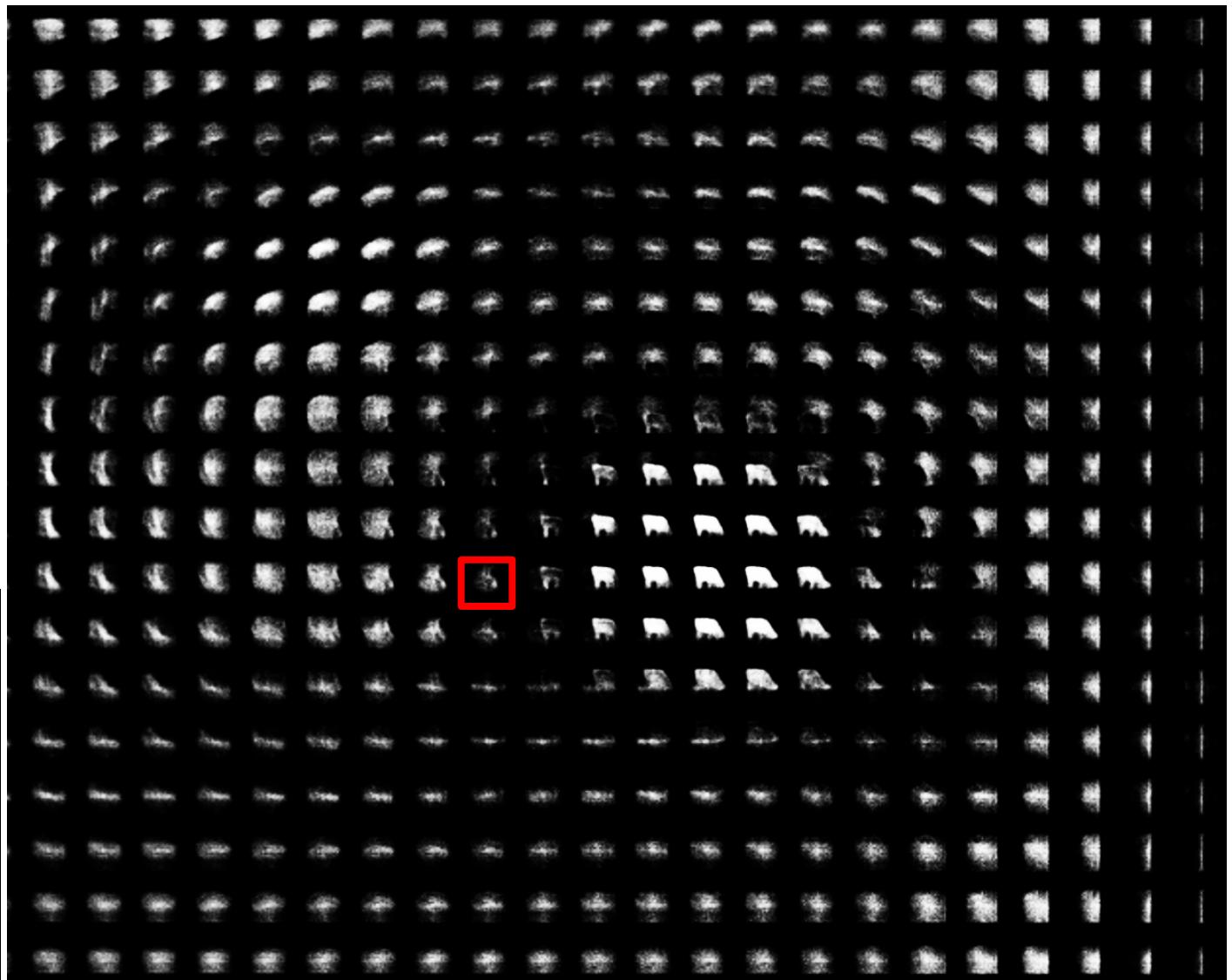
image



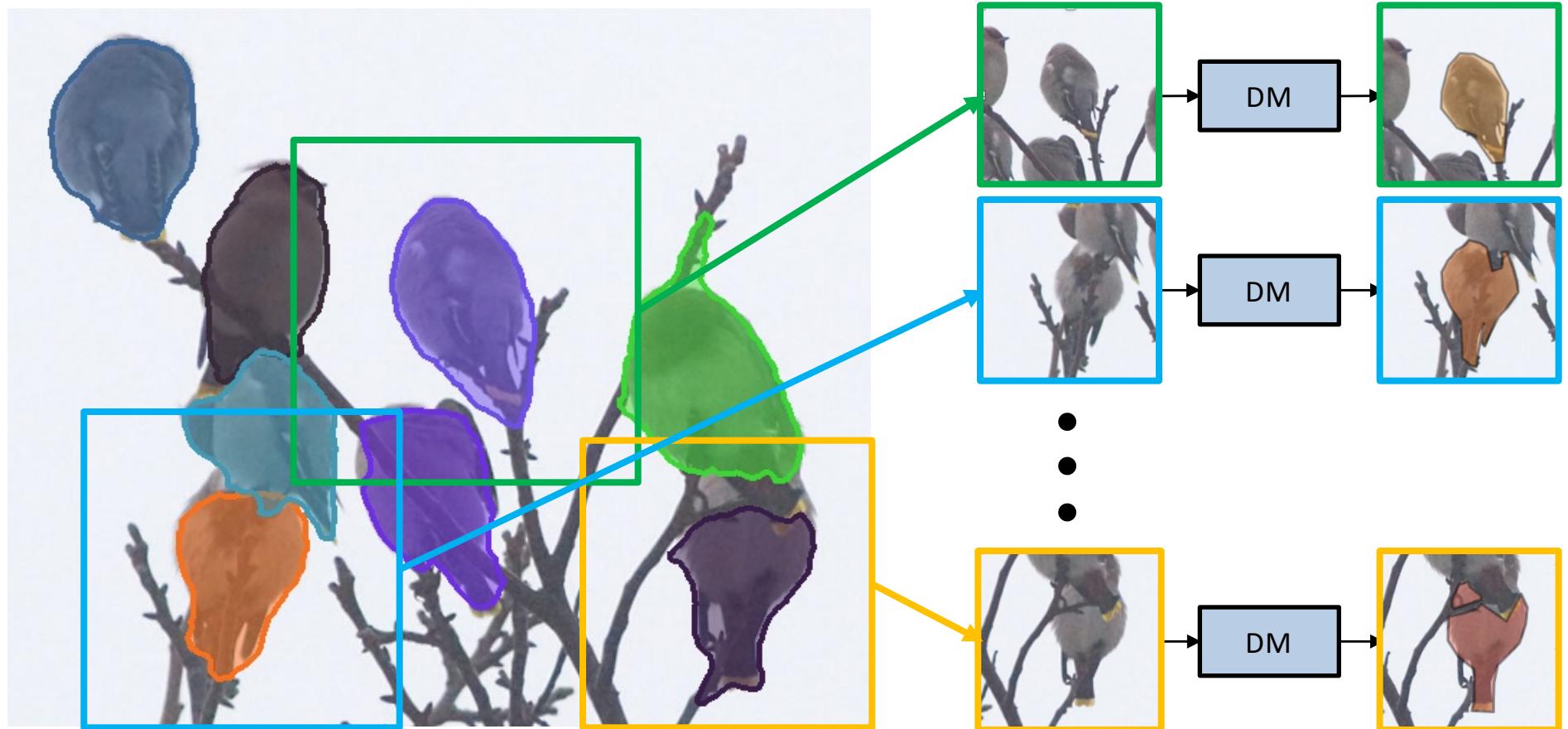
scores



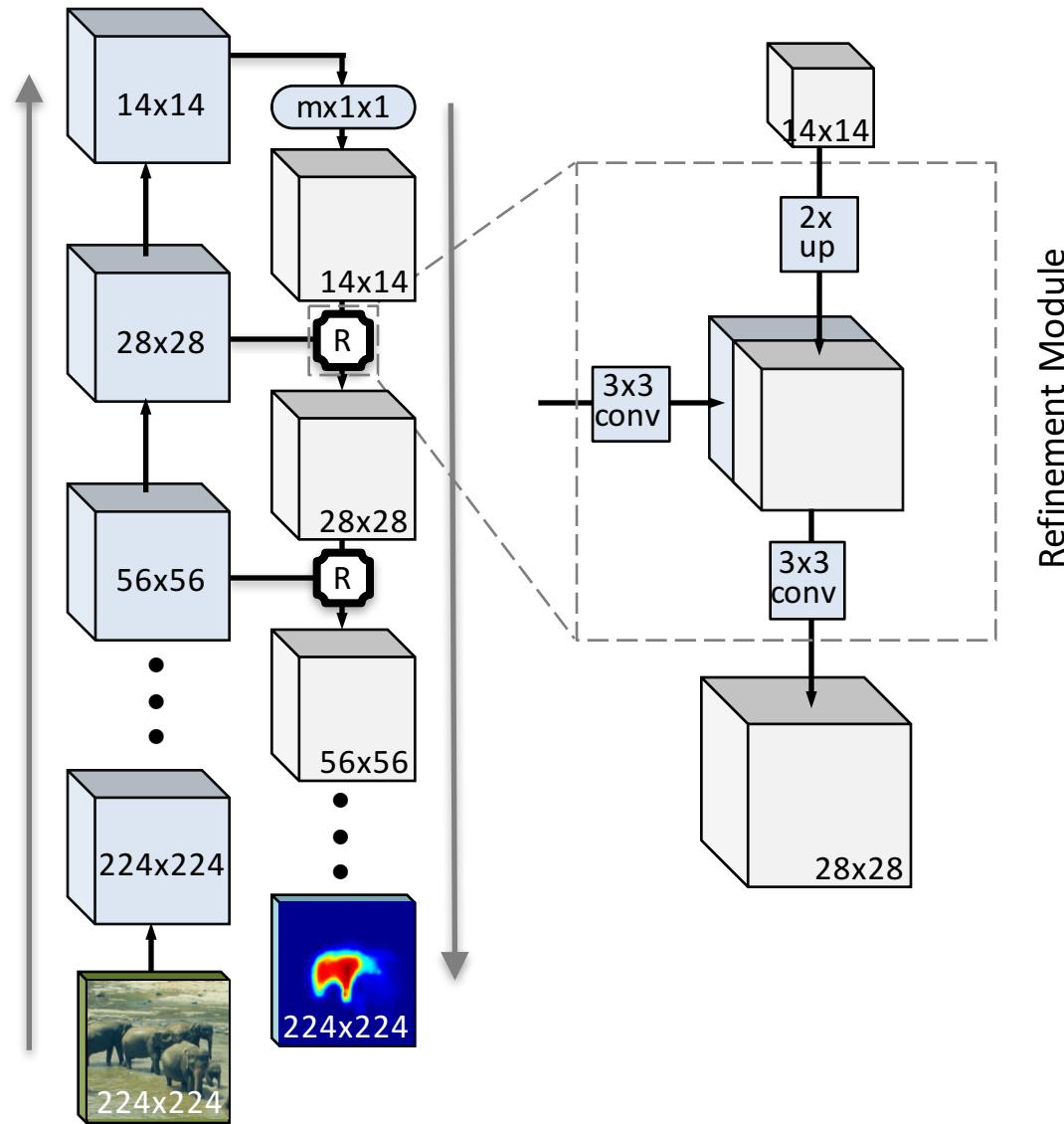
masks



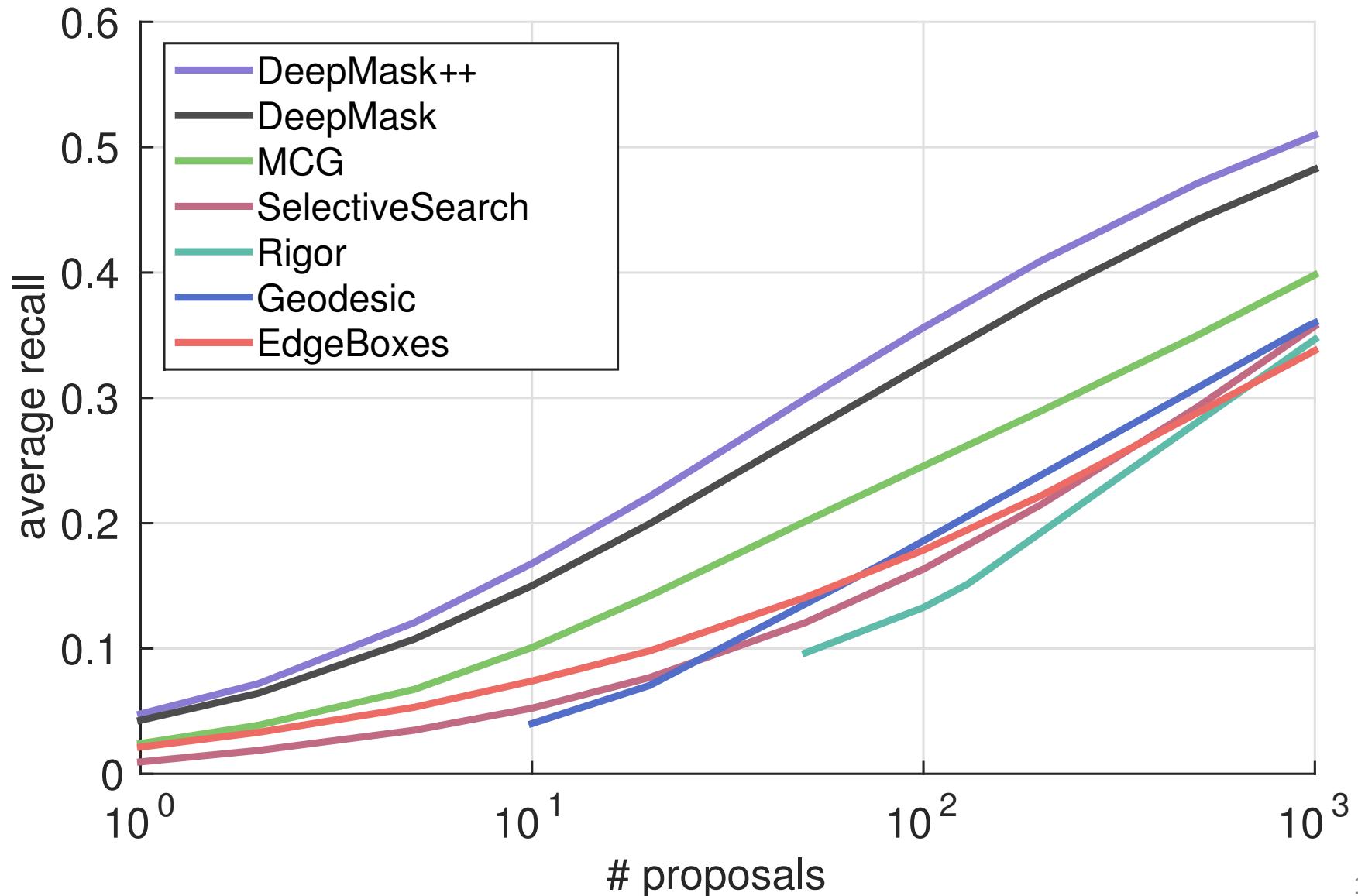
New: Iterative Localization (+1.0 AP)



New: Top-Down Refinement (+0.7 AP)



Proposal Quality (boxes)

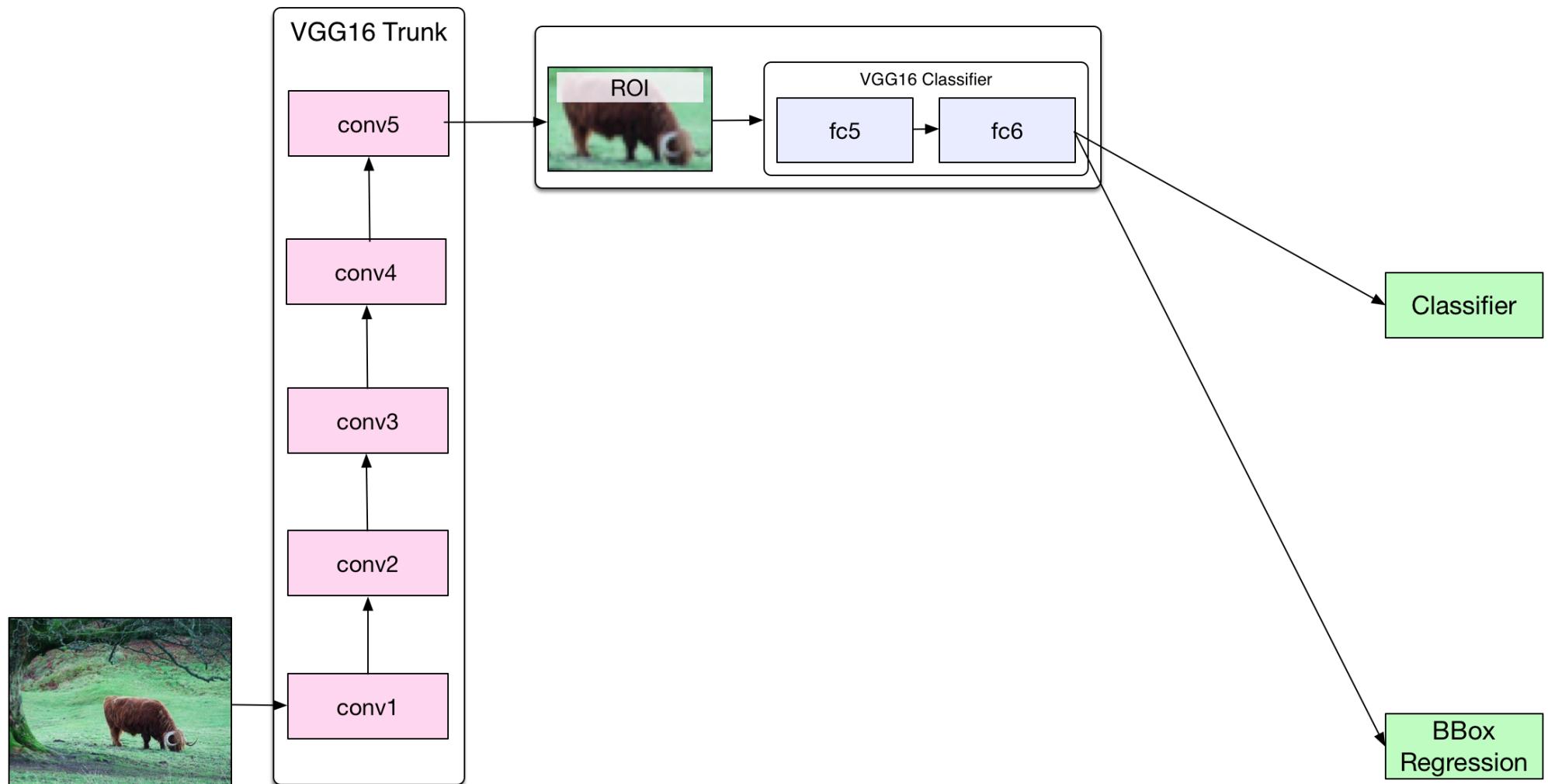


DeepMask Object Proposals

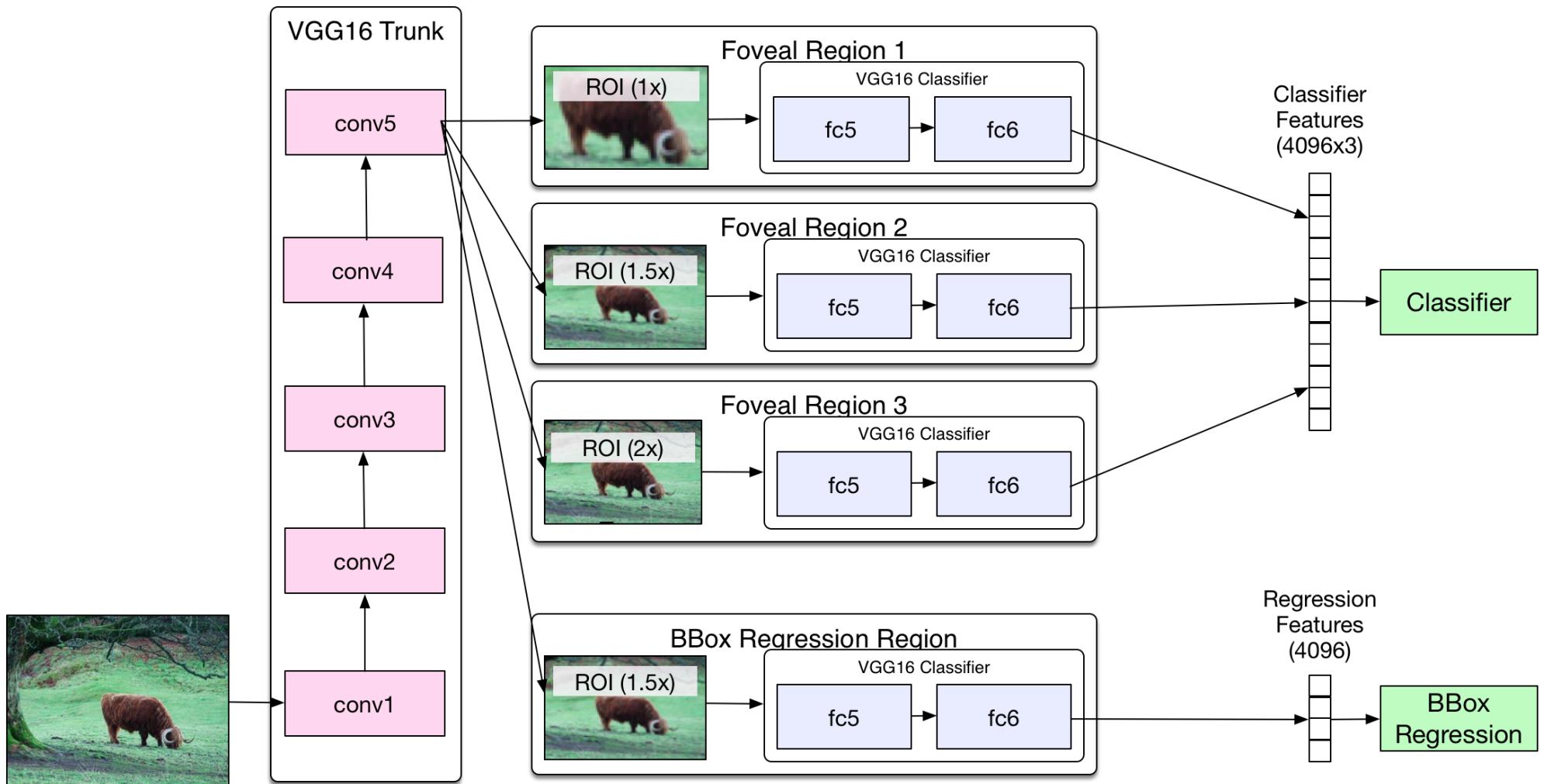


II. CLASSIFICATION FRAMEWORK

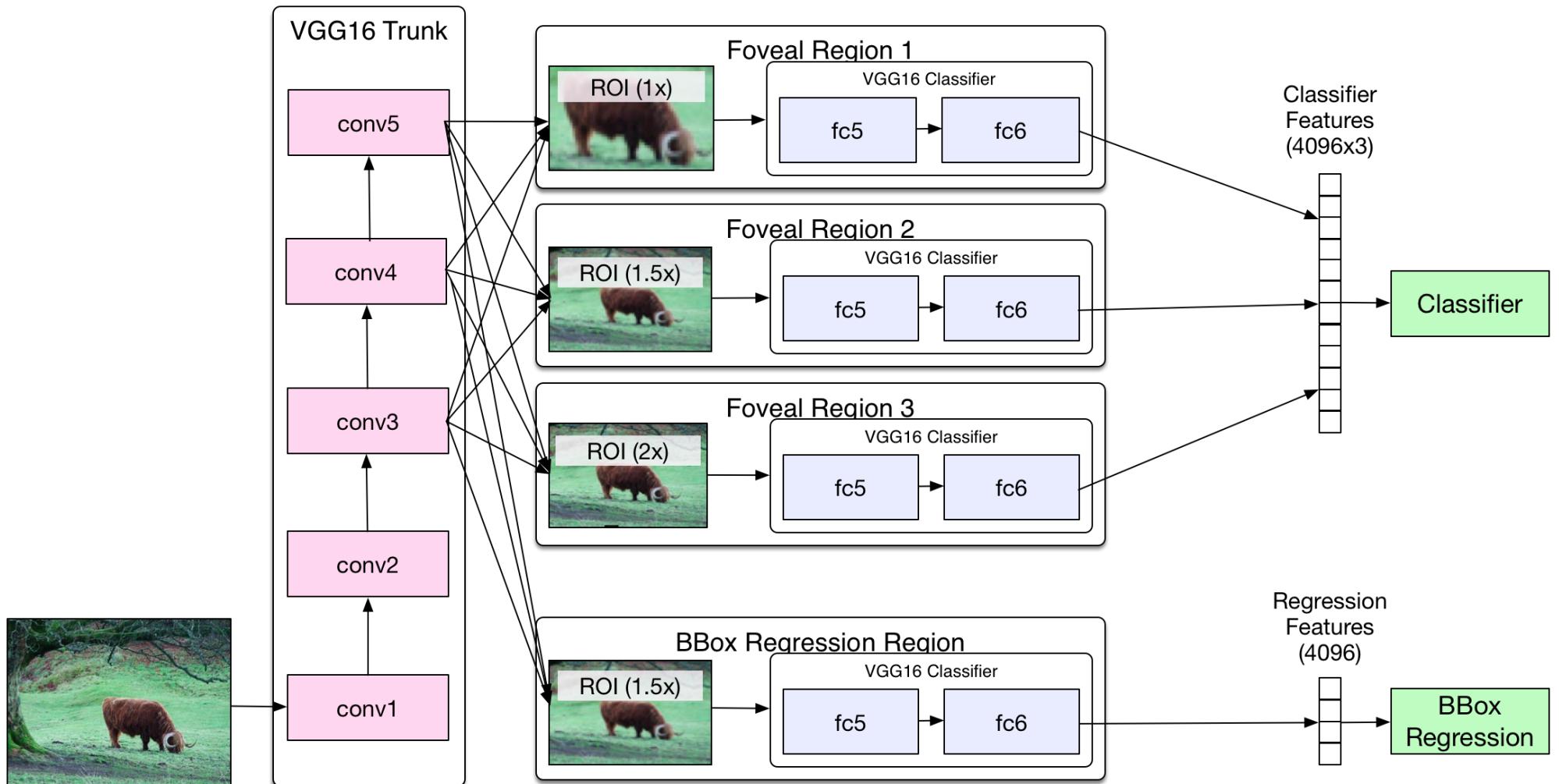
- Fast R-CNN setup [Girshick, ICCV15]



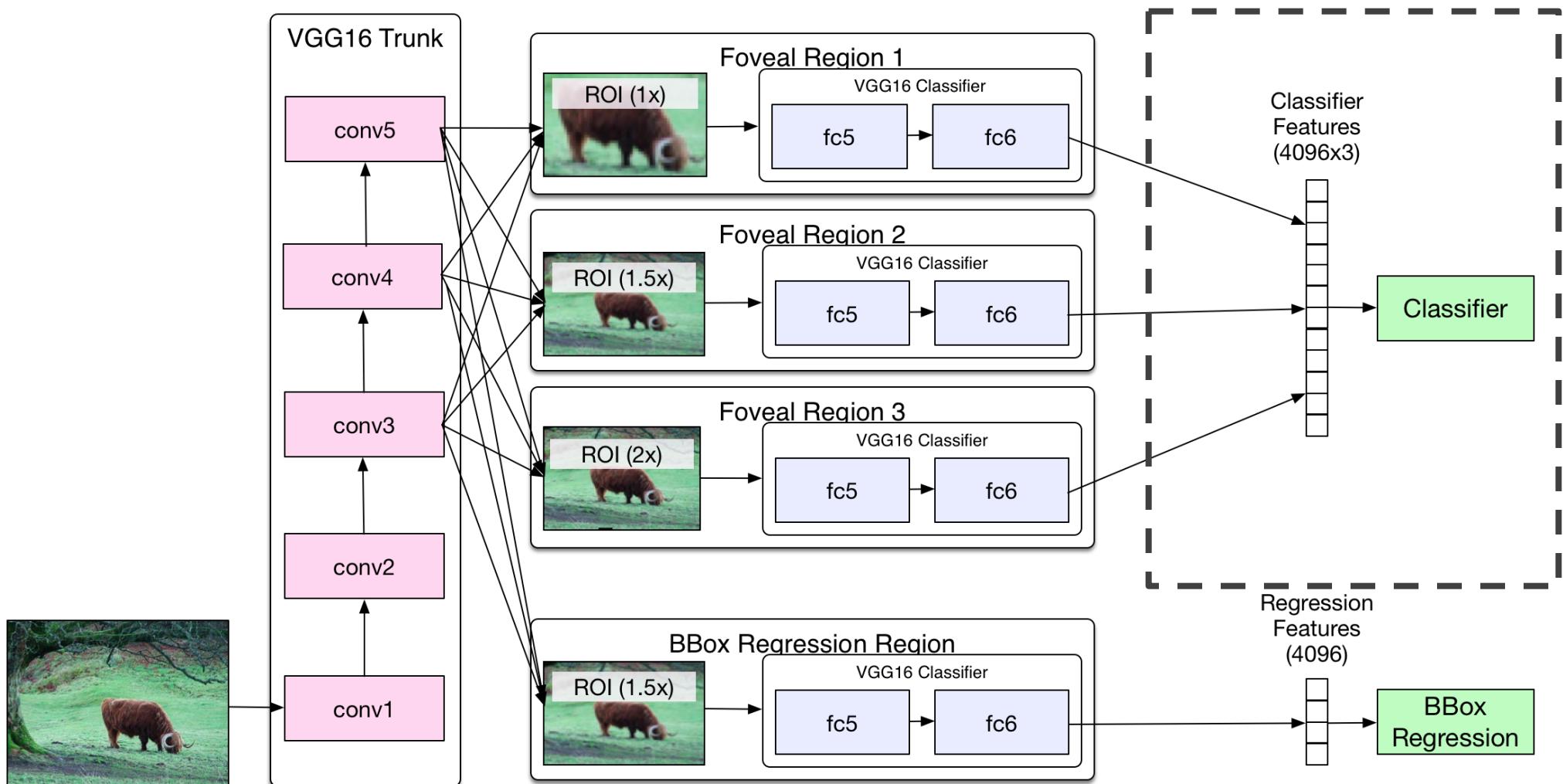
- Fast R-CNN setup [Girshick, ICCV15]
- Foveal structure [inspired by Gidaris & Komodakis, ICCV15] (+2 AP)



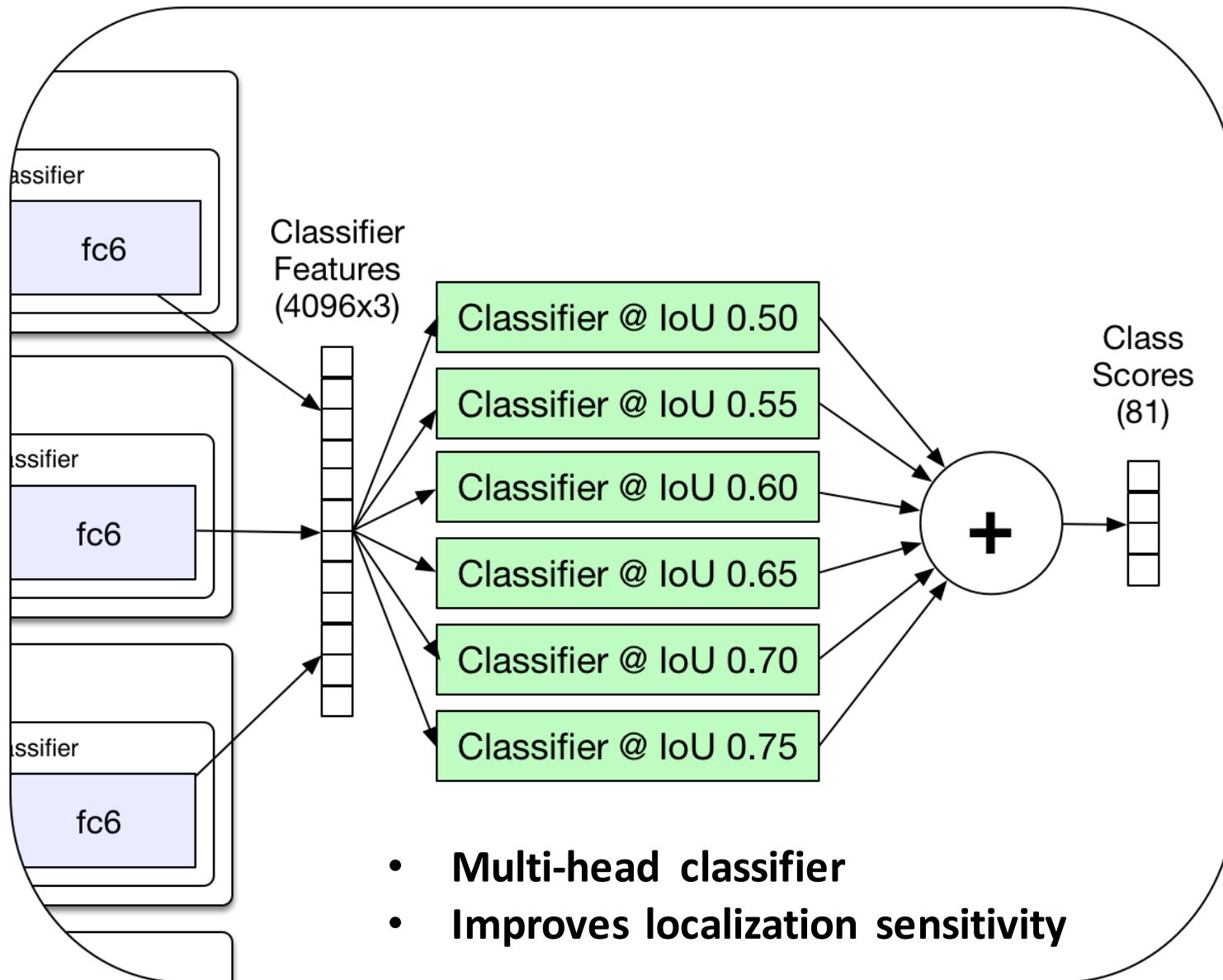
- Fast R-CNN setup [Girshick, ICCV15]
- Foveal structure [inspired by Gidaris & Komodakis, ICCV15] (+2 AP)
- Skip connections (+1 AP)



- 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)



Training

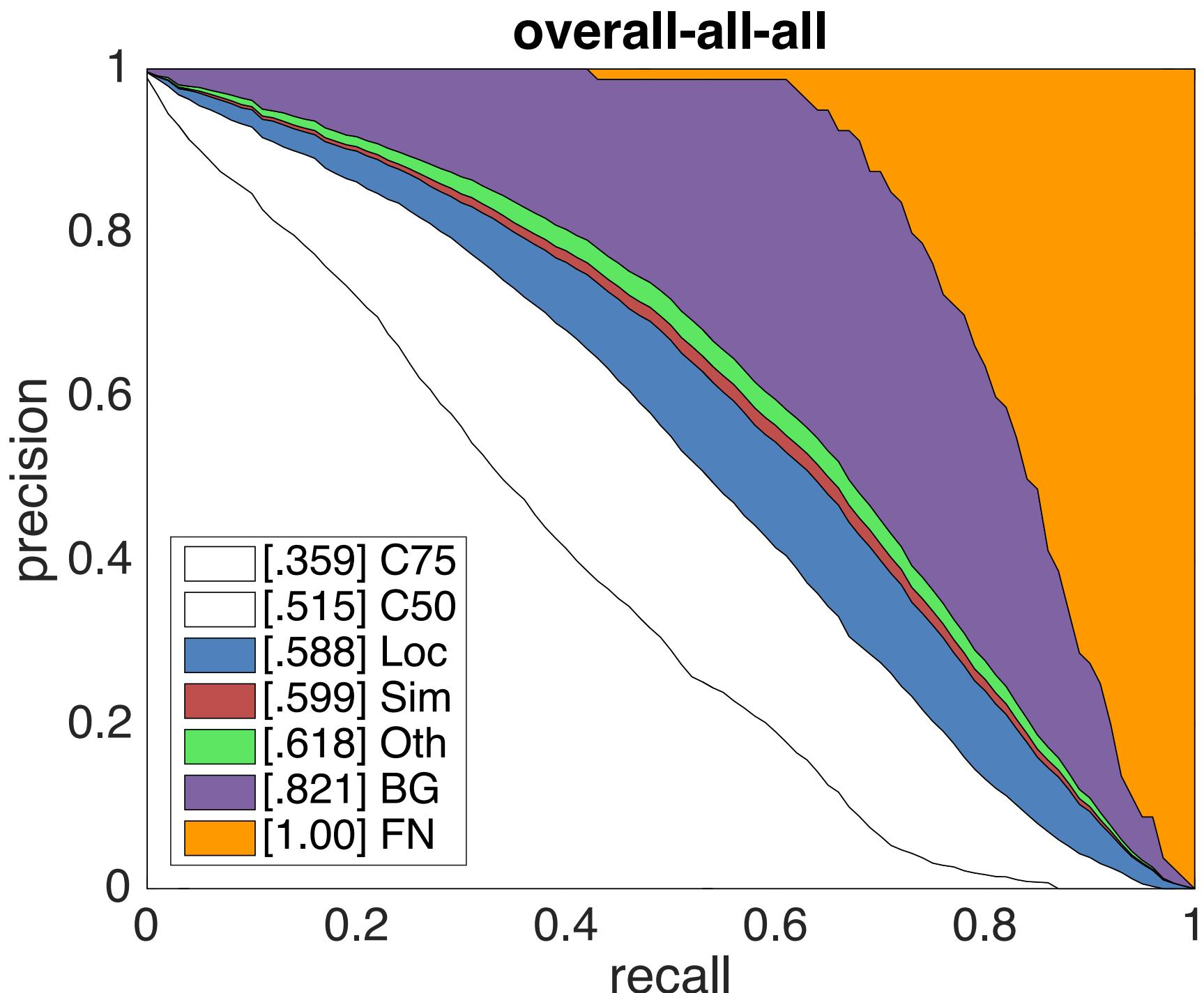
Two strategies:

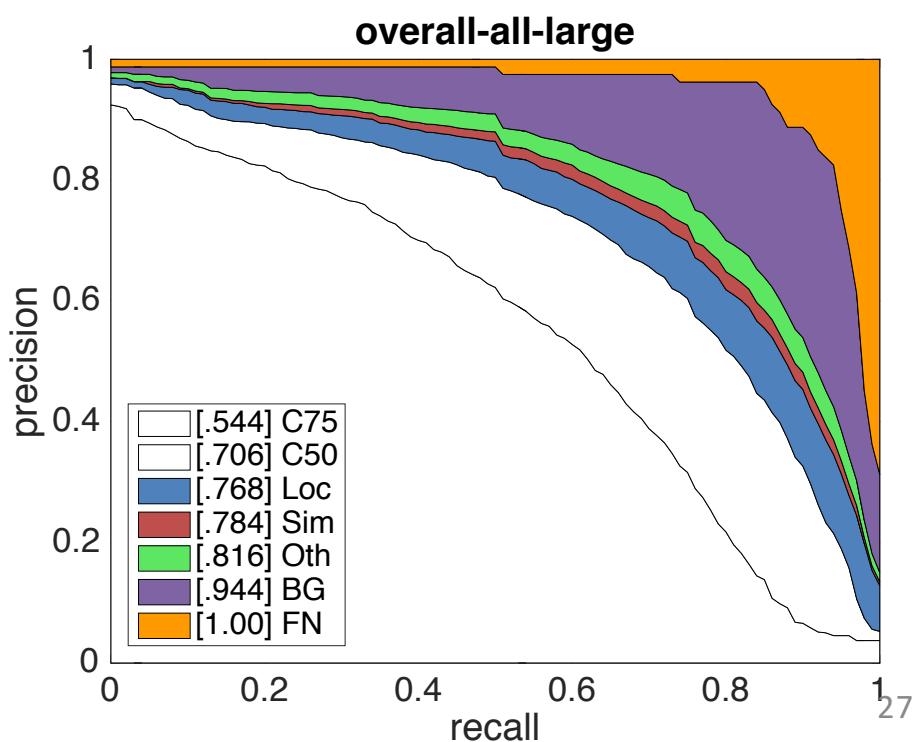
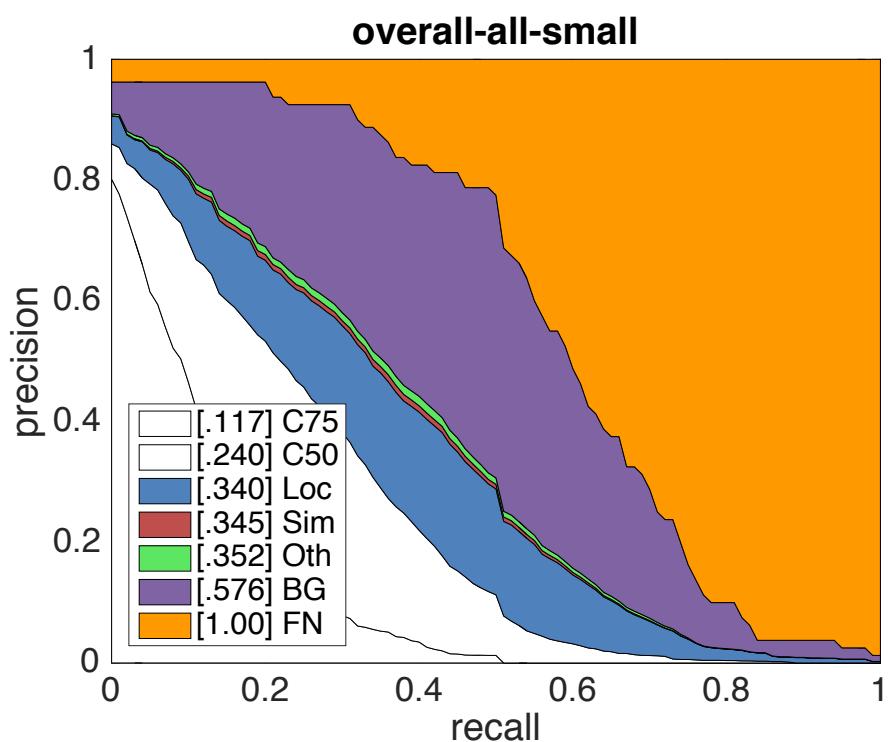
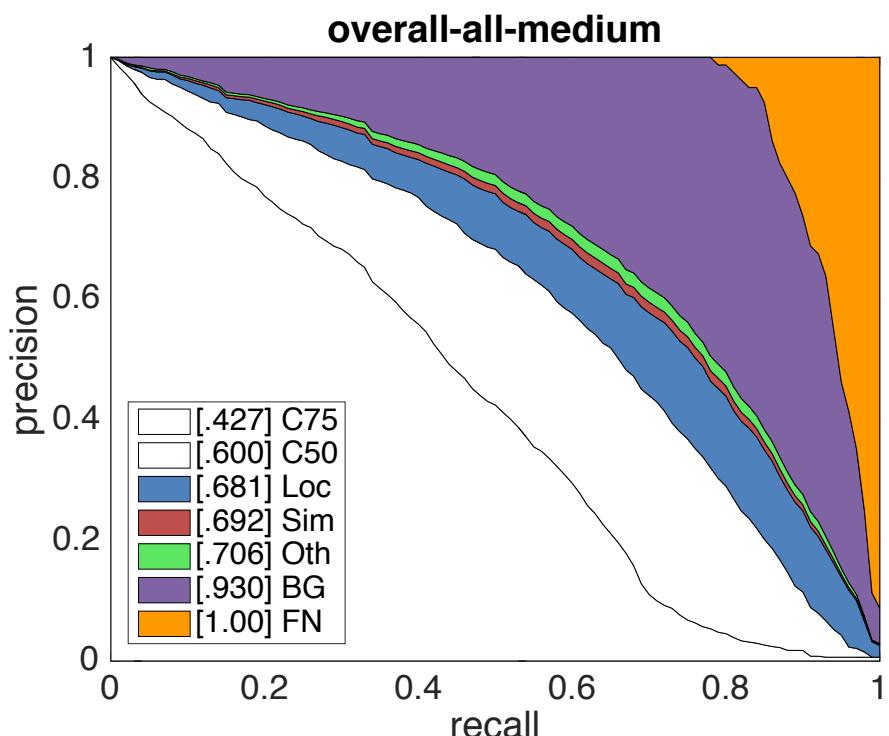
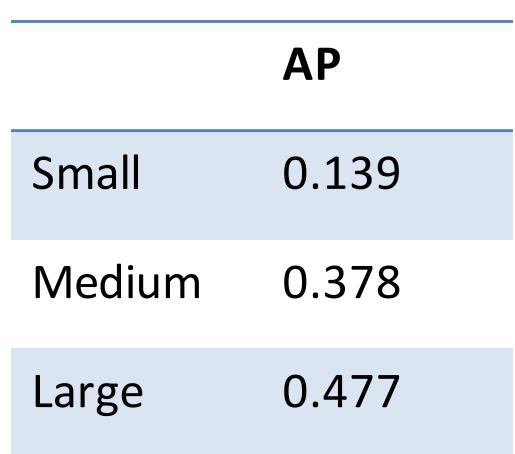
- 4 days on 4 Maxwell GPUs [Big Sur]
- 2.5 days on 8x4 Kepler with Elastic Averaging SGD [Zhang, NIPS 15]



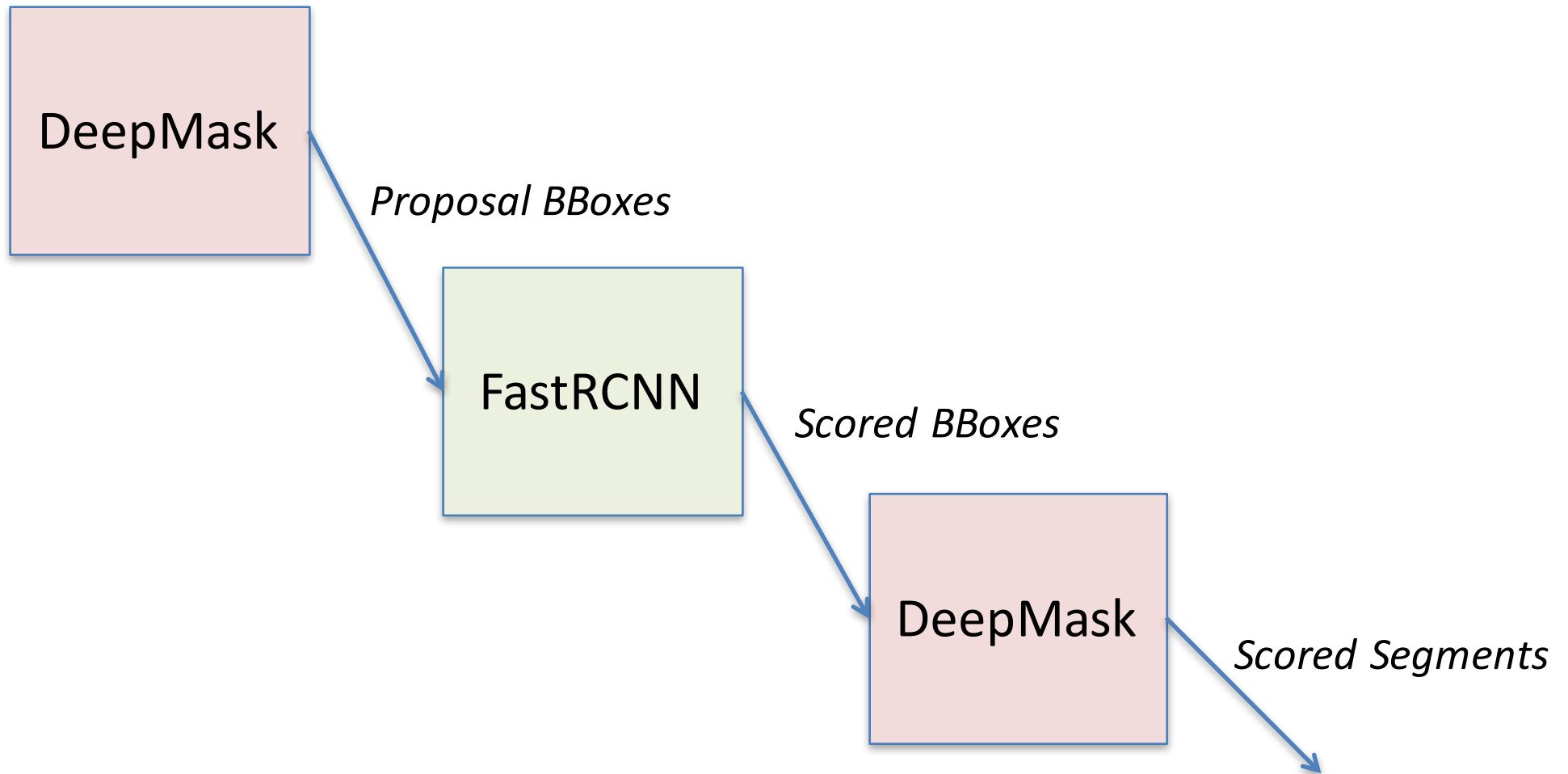
Inference

Base Model	30.1 AP
+ horizontal flip	31.1 AP
+ ROI Pooling ‘2 crop’	32.1 AP
+ 7-model Ensemble	33.5 AP





Segmentation Examples





person

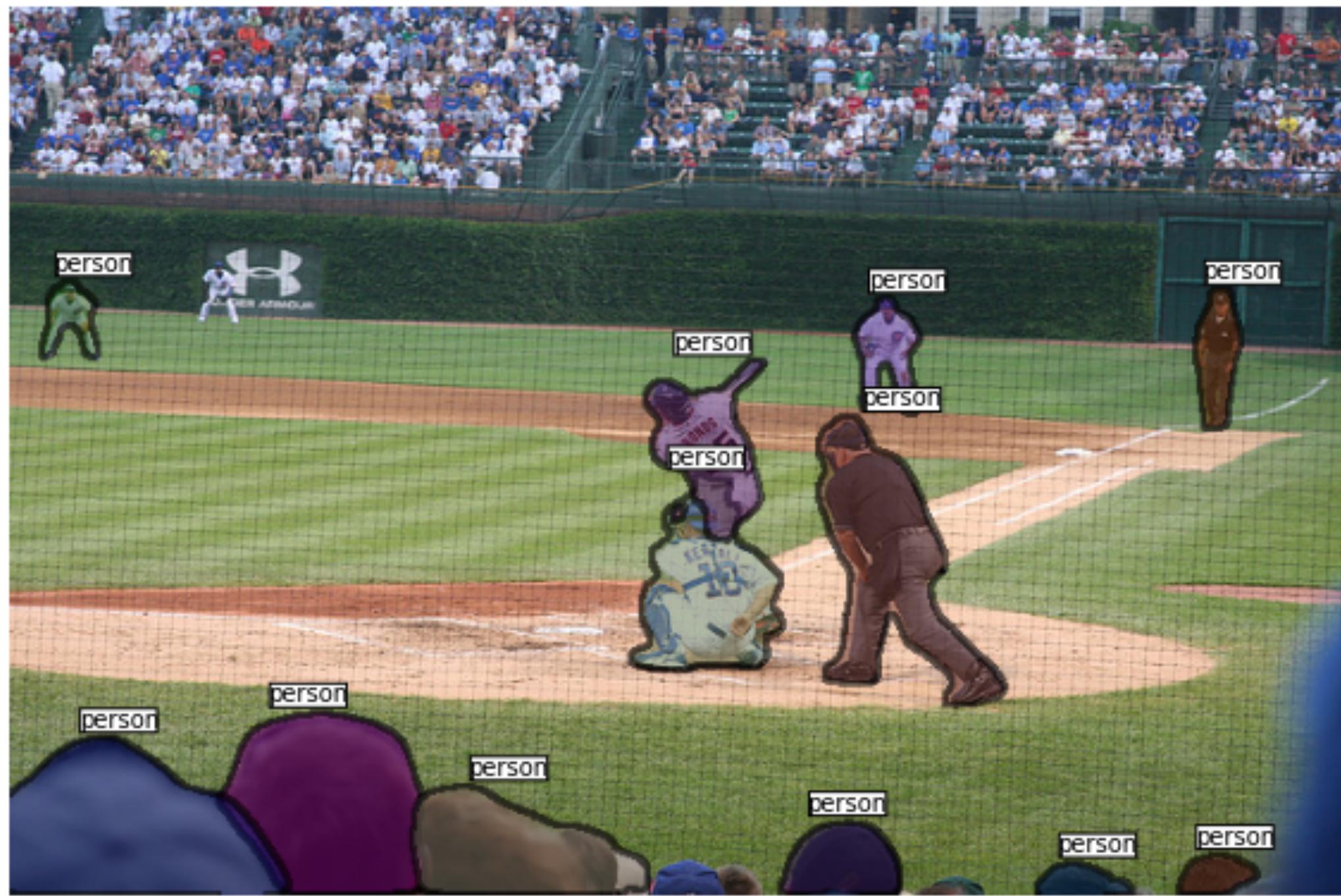
cell phone

cell phone

aptop

aptop

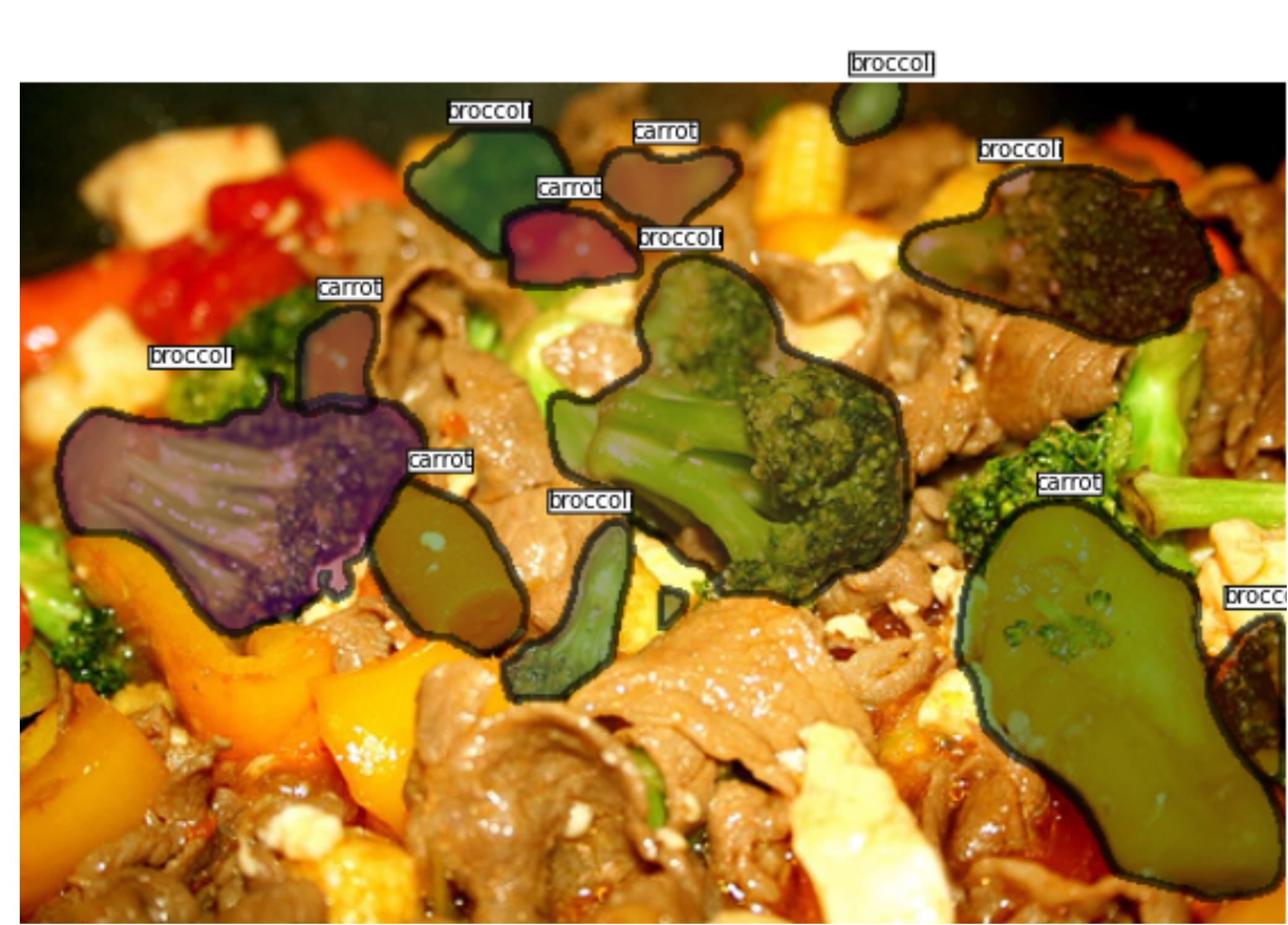




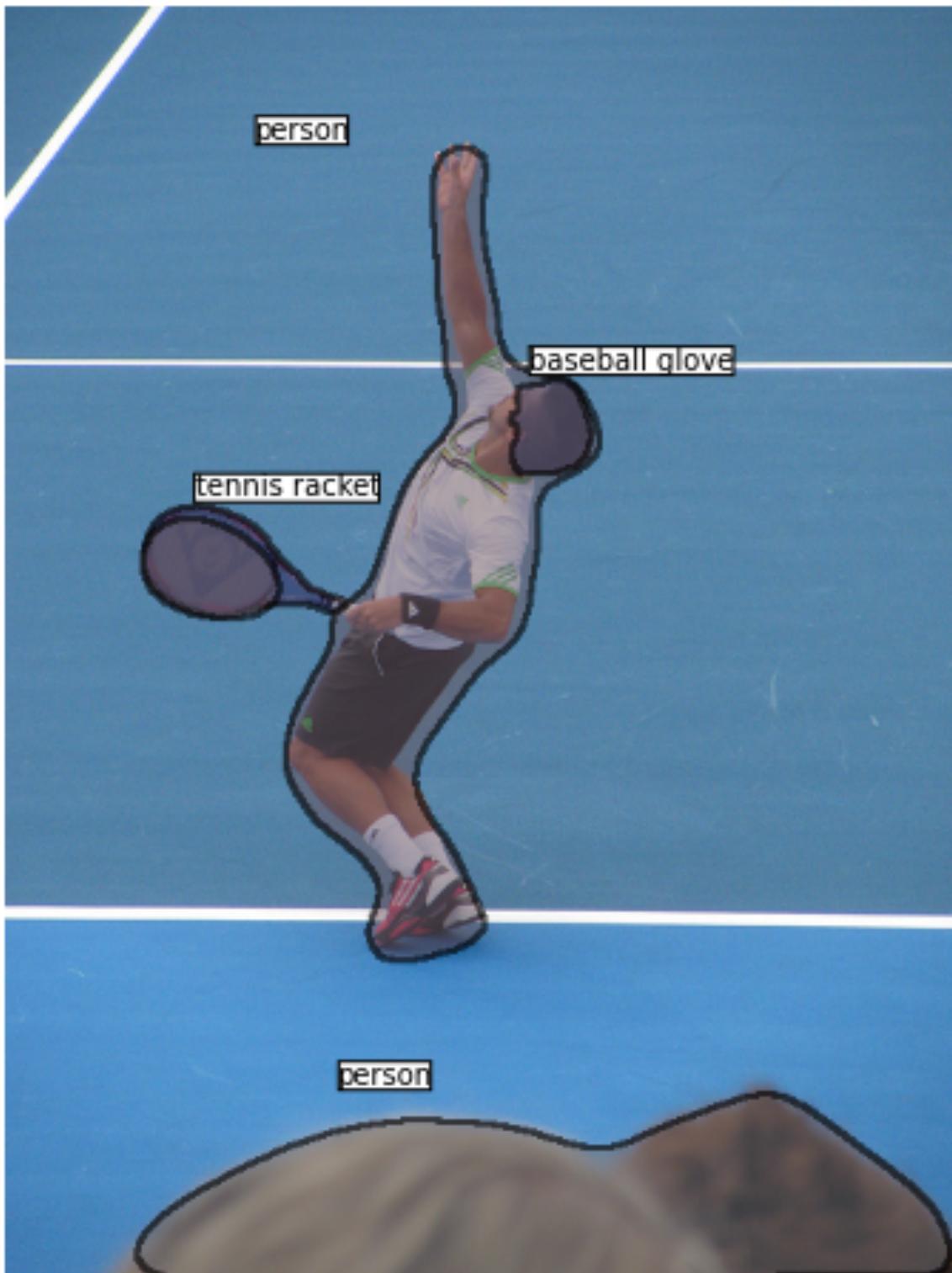


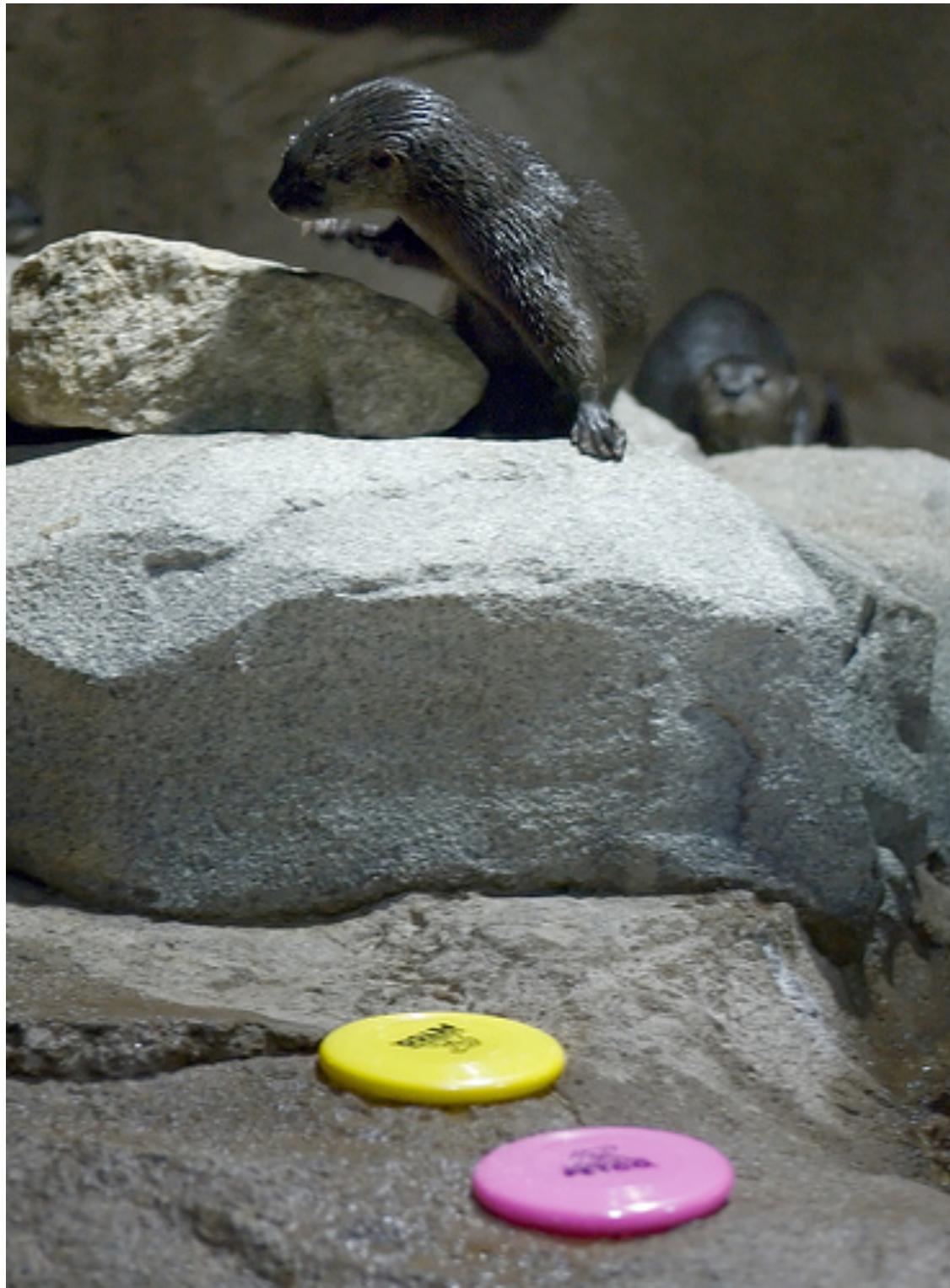


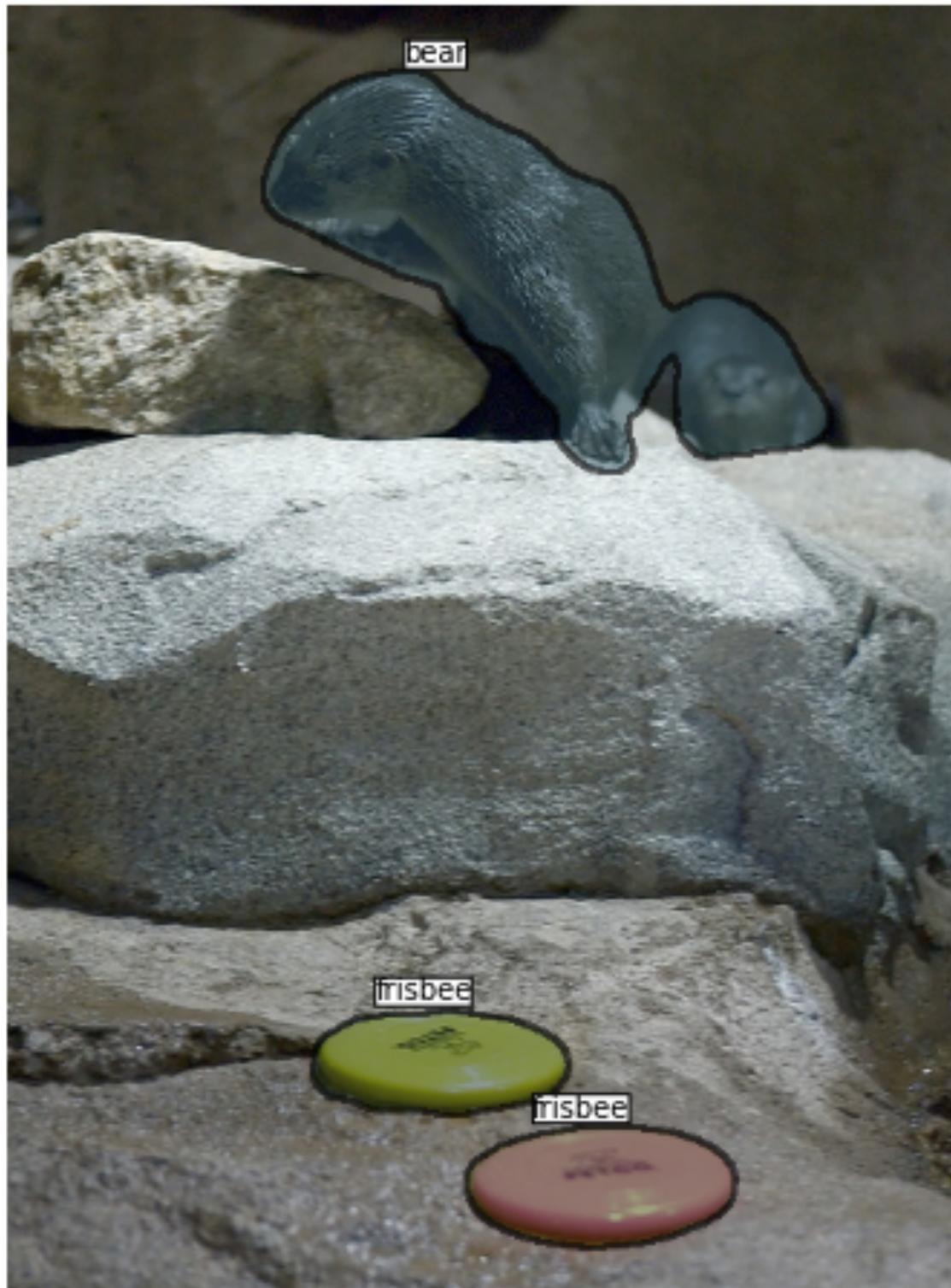


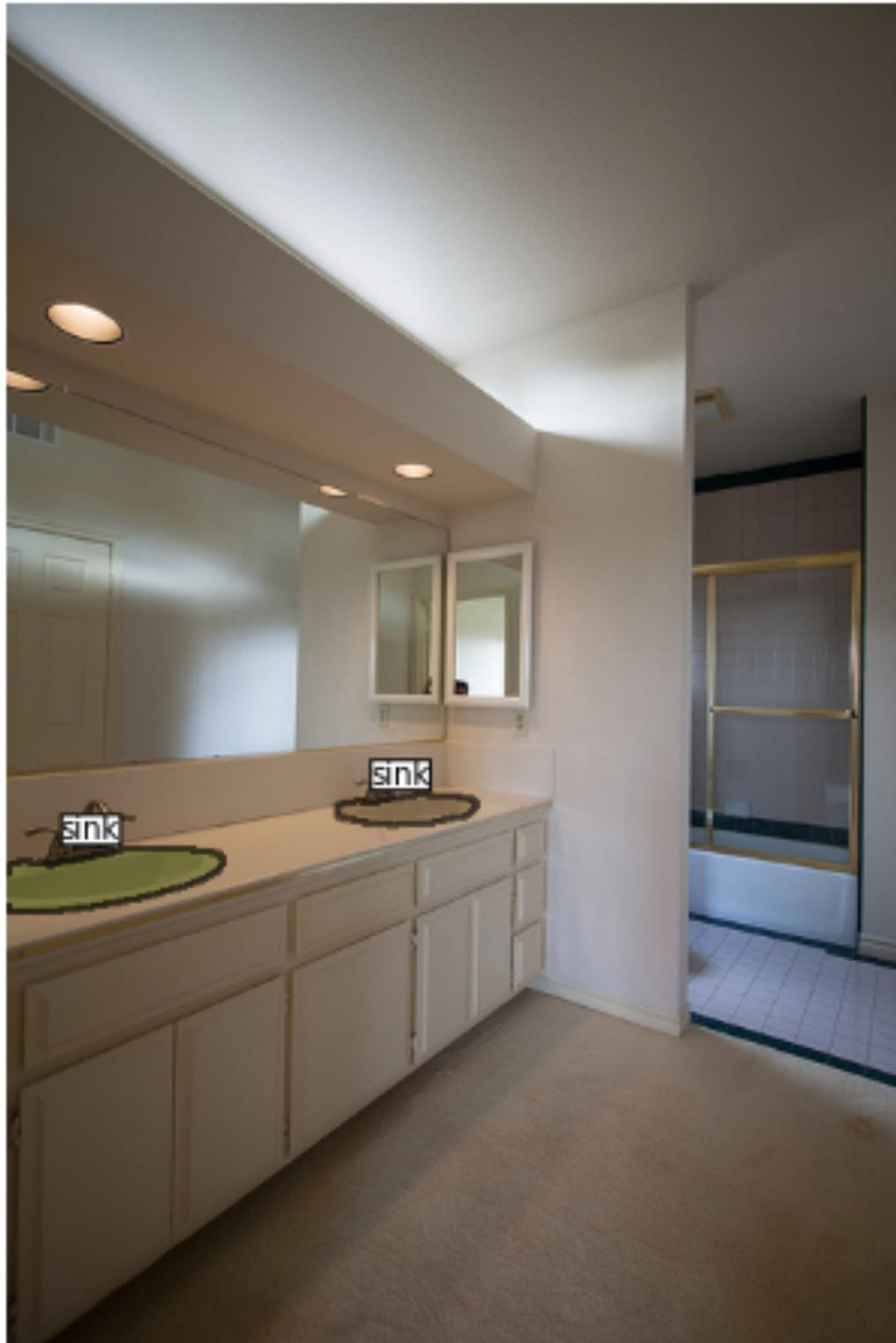
























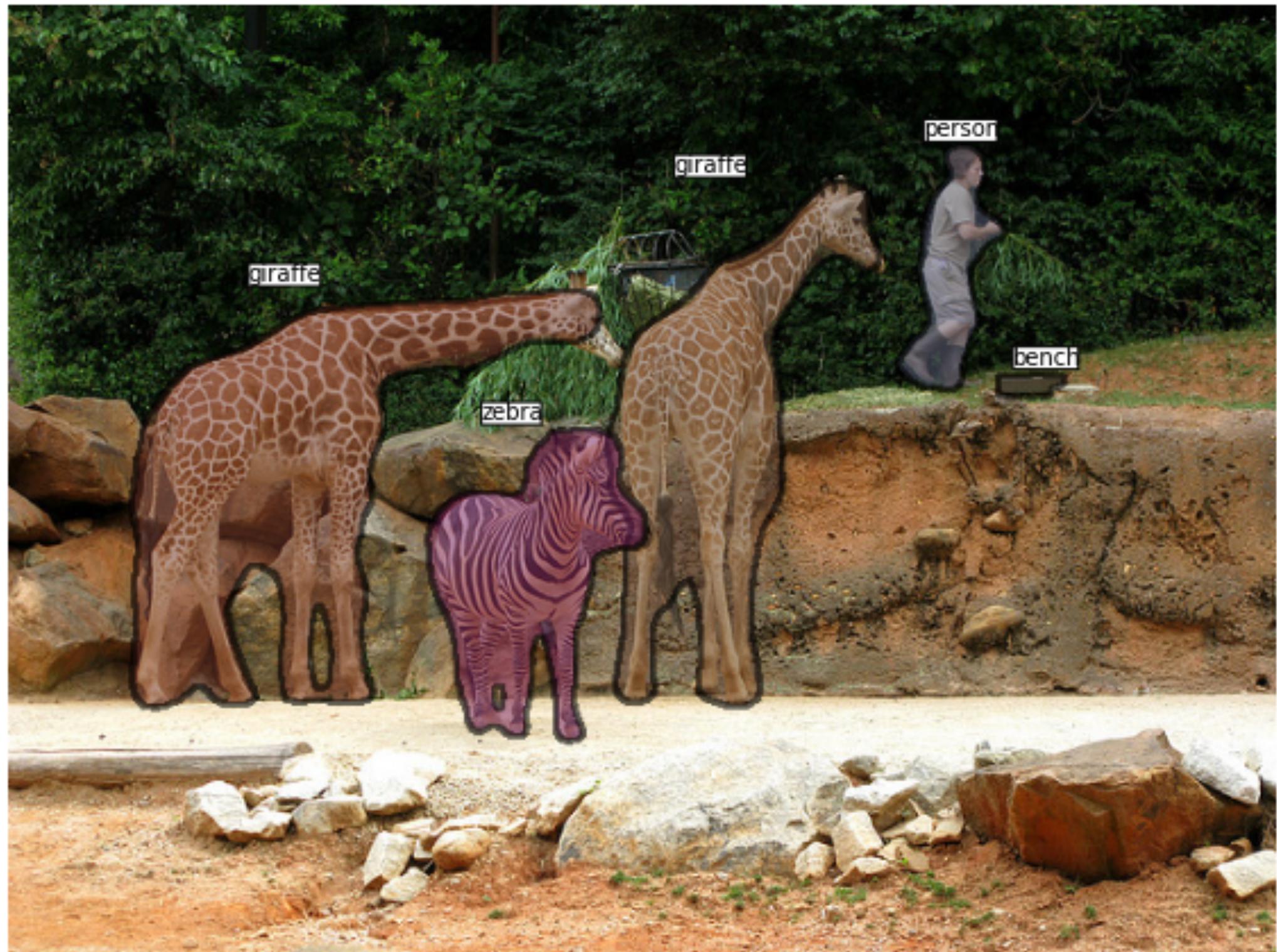


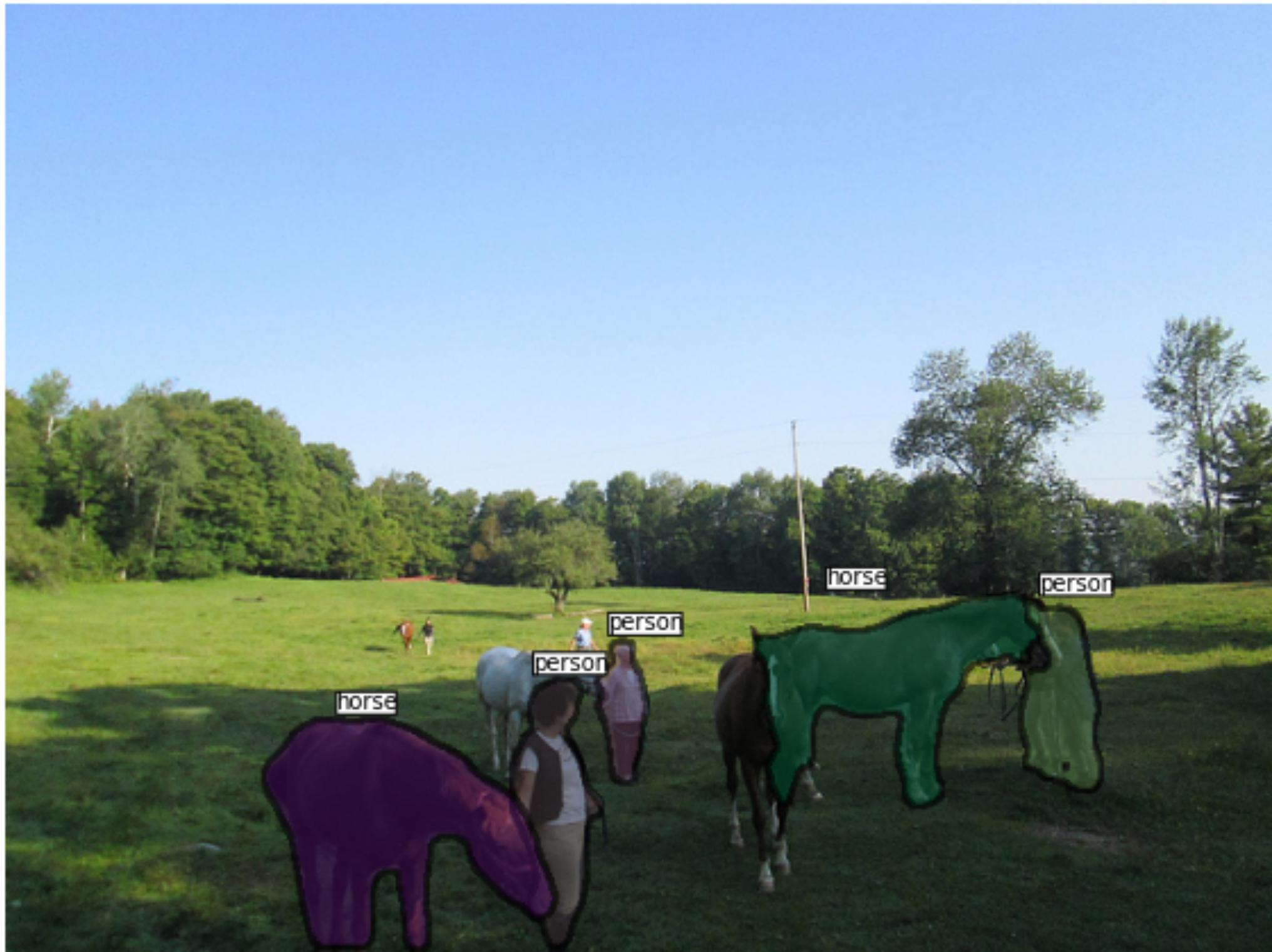
zebra

zebra













person

person

person

chair

chair
chair

per
pers

bowl bowl dining table

bowl

bottle

bowl

bowl

bottle

bottle

broccoli

cake

bowl

banana

banana
banana

book

bowl

carrot

banana

bowl

broccoli
broccoli

apple

a [ap] apple

carrot

carrot

carrot

Future Directions

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

