

Rich feature hierarchies for accurate object detection and semantic segmentation

Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik
UC Berkeley

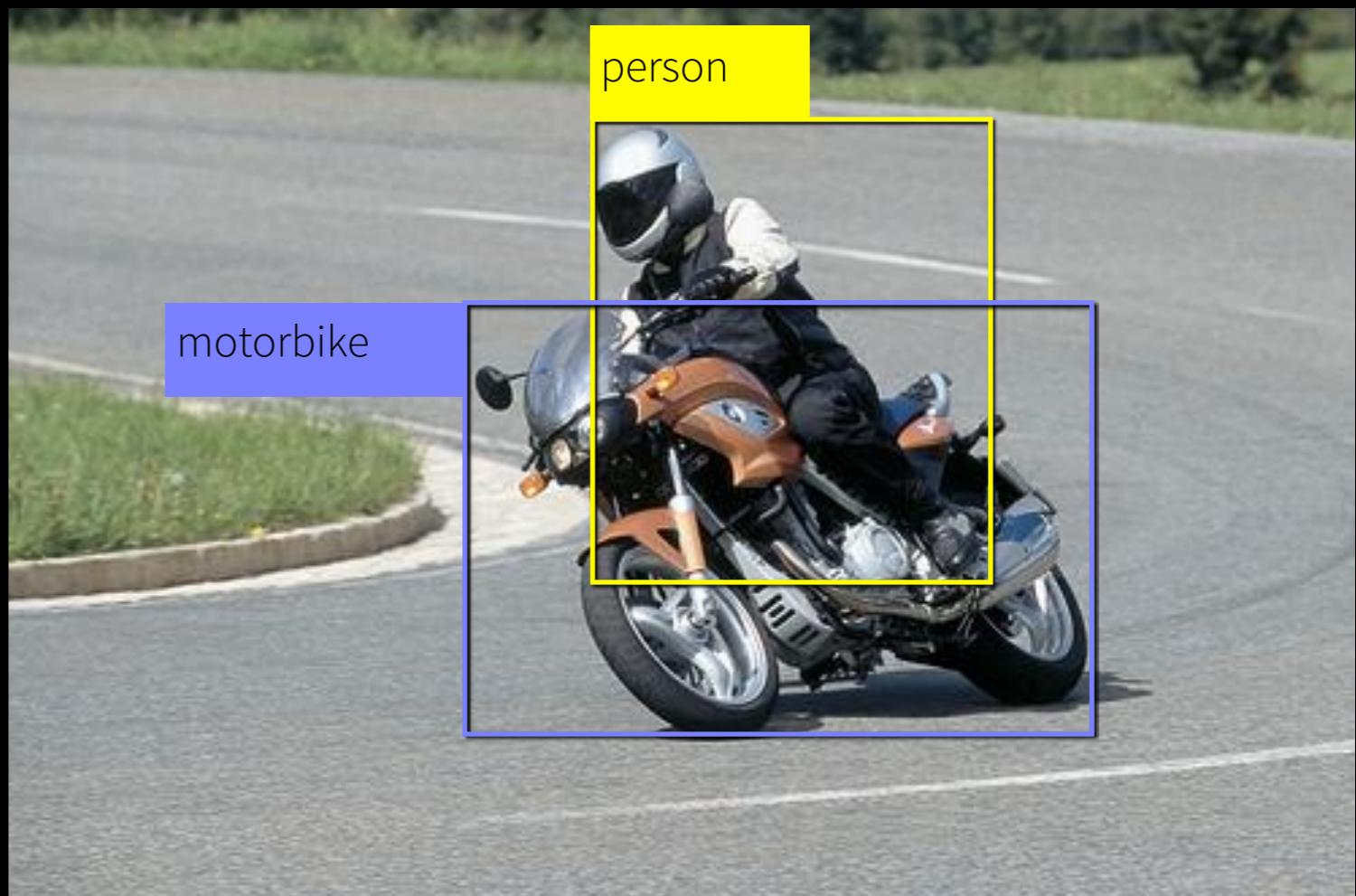
Tech Report @ <http://arxiv.org/abs/1311.2524>

Detection & Segmentation

input



background



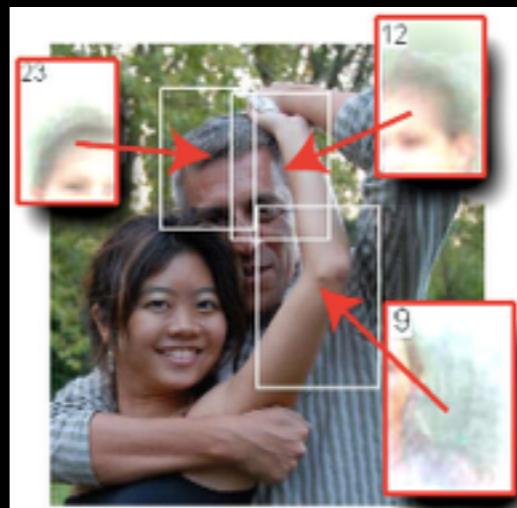
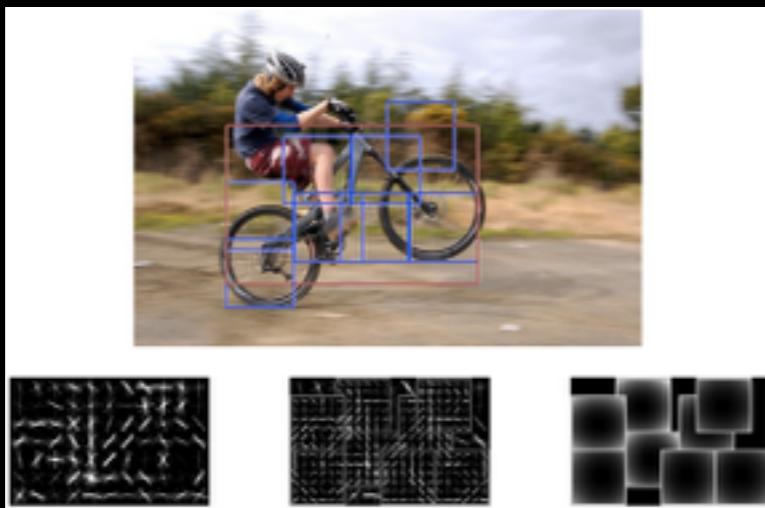
PASCAL VOC



Example PASCAL VOC images

Dominant detection methods

1. Part-based sliding window methods (HOG)



2. Region-proposal classifiers (SIFT++ BoW)



Russell et al. 2006

Gu et al. 2009

van de Sande et al. 2011

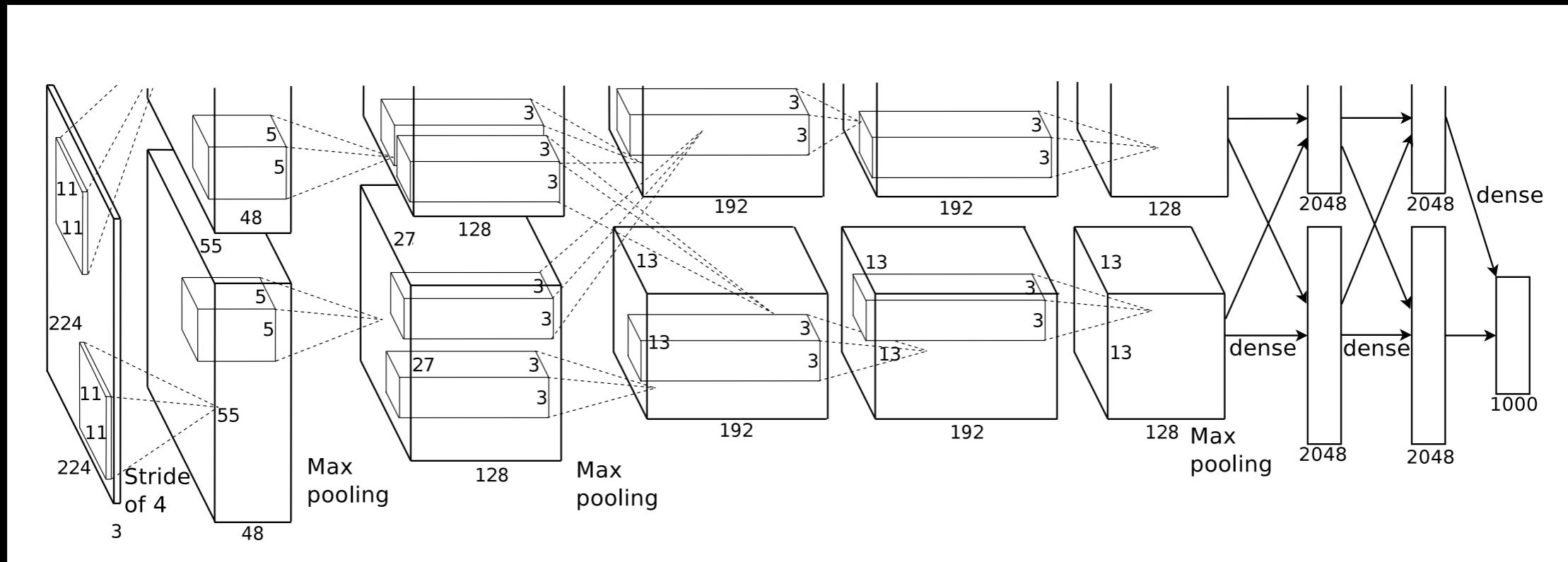
> “selective search”

PASCAL VOC epochs (detection)

- 2007-2010 The Moore's law years
- 2010-2011 The year of kitchen sinks (or the all-too-soon end of Moore's law)
- 2011-2012 Stagnation (no new features left, juice all squeezed from context)
- 2013- *Learning rich features?*

ImageNet LSVRC'12 winner

UToronto “SuperVision” CNN



Krizhevsky, Sutskever, and Hinton.

ImageNet Classification with Deep Convolutional Neural Networks.
NIPS 2012.

cf. LeCun et al. Neural Comp. '89 & Proc. of the IEEE '98

Impressive ImageNet results!

Task: 1000-way whole-image classification

Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
<i>SIFT + FVs [7]</i>	—	—	26.2%
1 CNN	40.7%	18.2%	—
5 CNNs	38.1%	16.4%	16.4%
1 CNN*	39.0%	16.6%	—
7 CNNs*	36.7%	15.4%	15.3%

metric: classification error rate (lower is better)

But... does it generalize to other datasets and tasks?

See: Donahue, Jia, et al. DeCAF Tech Report.

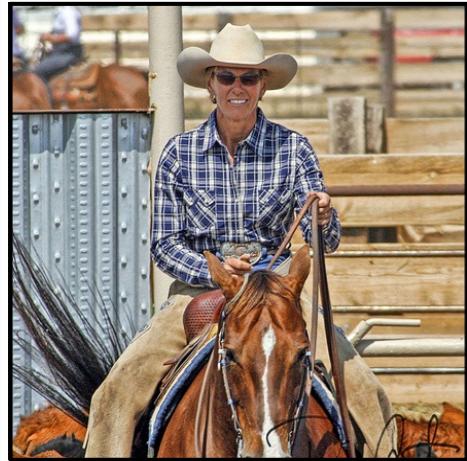
Much debate at ECCV'12

Objective

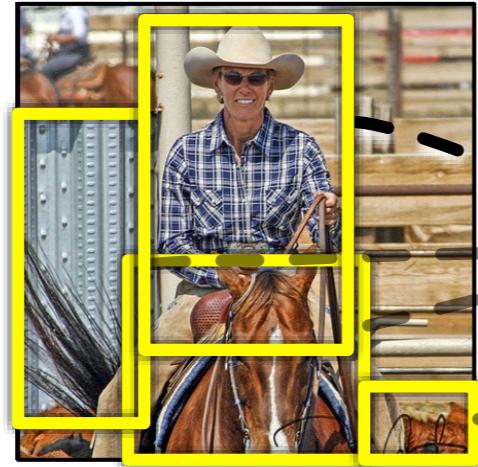
Understand if the SuperVision CNN can be made to work as an object detector.

Object detection system

R-CNN: “Regions with CNN features”

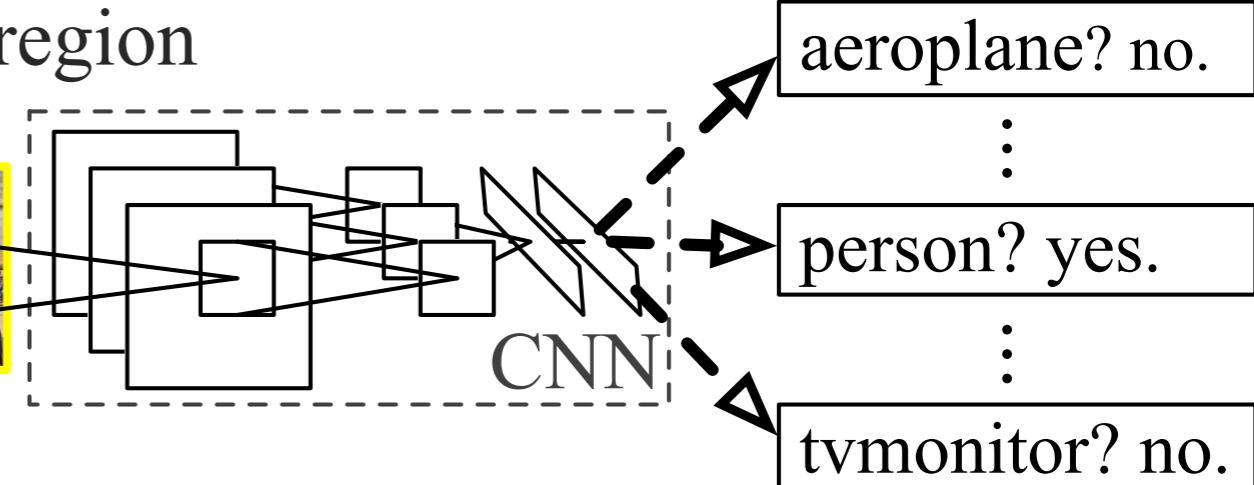


1. Input image



2. Extract region proposals (~2k)

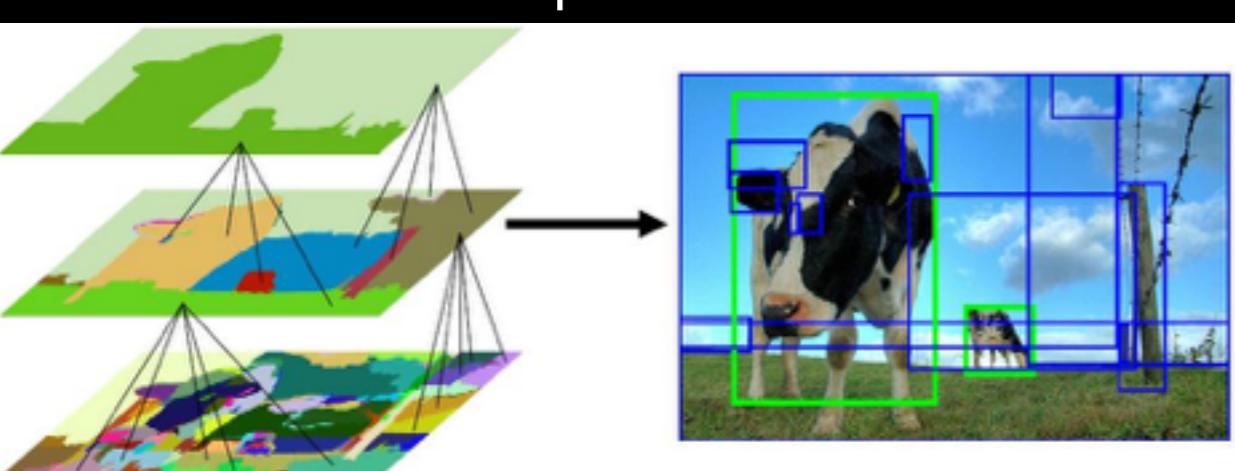
warped region



3. Compute CNN features

4. Classify regions

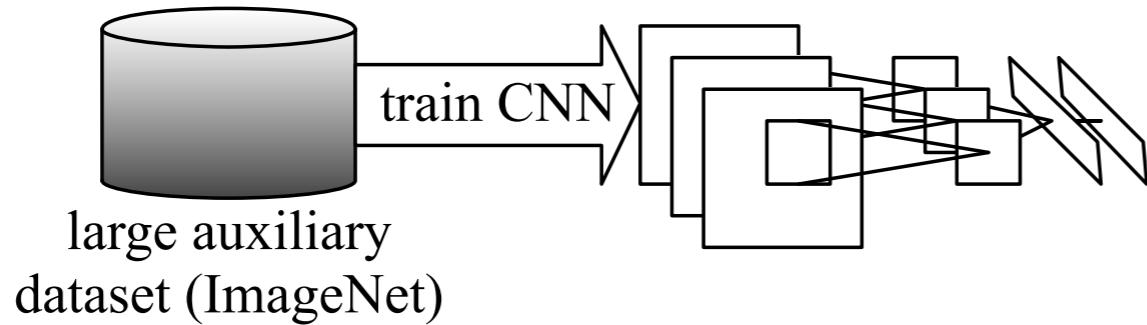
(With a few minor tweaks:
semantic segmentation)



(e.g. selective search)

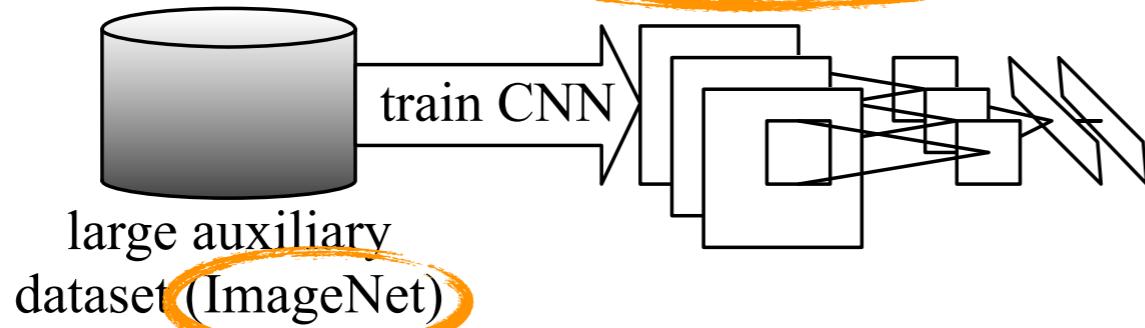
Training

1. Pre-train CNN for **image classification**



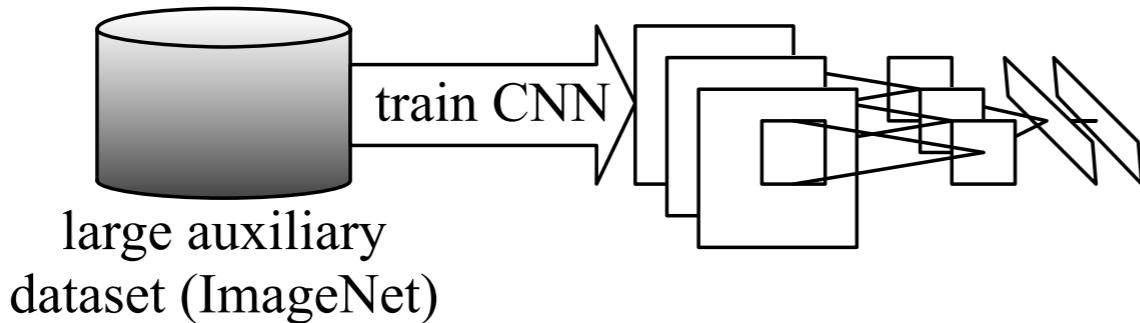
Training

1. Pre-train CNN for **image classification**

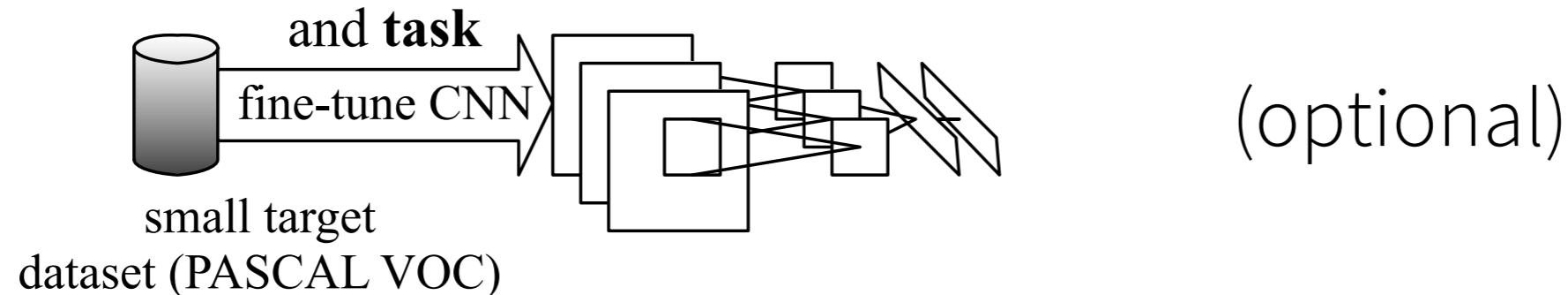


Training

1. Pre-train CNN for **image classification**

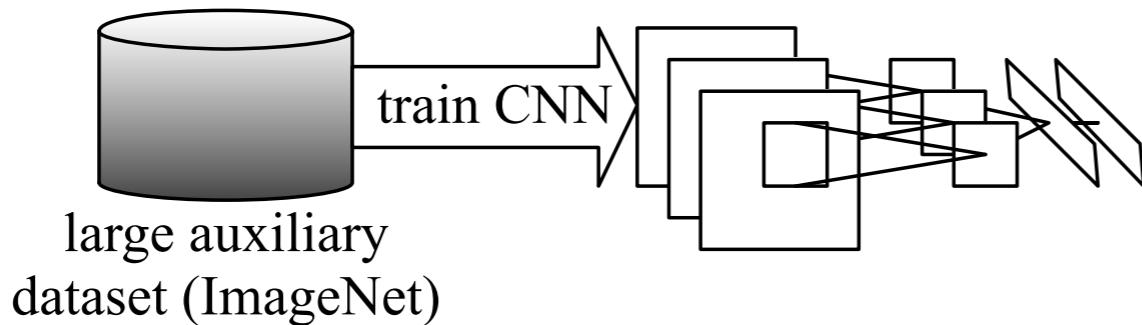


2. Fine-tune CNN on **target dataset** and **task**

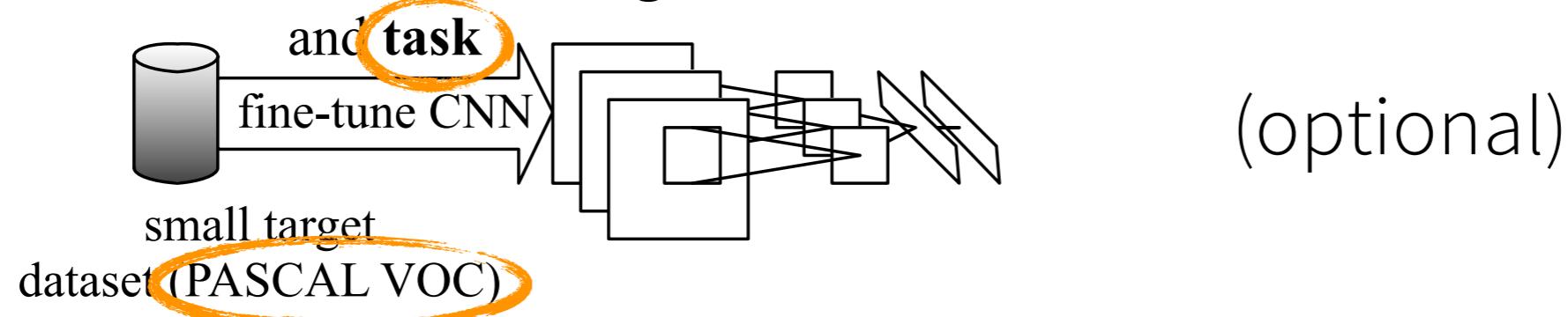


Training

1. Pre-train CNN for **image classification**

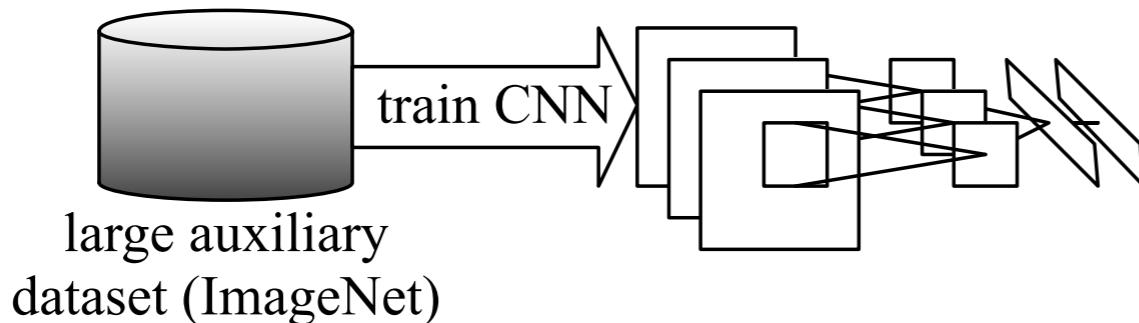


2. Fine-tune CNN on **target dataset**

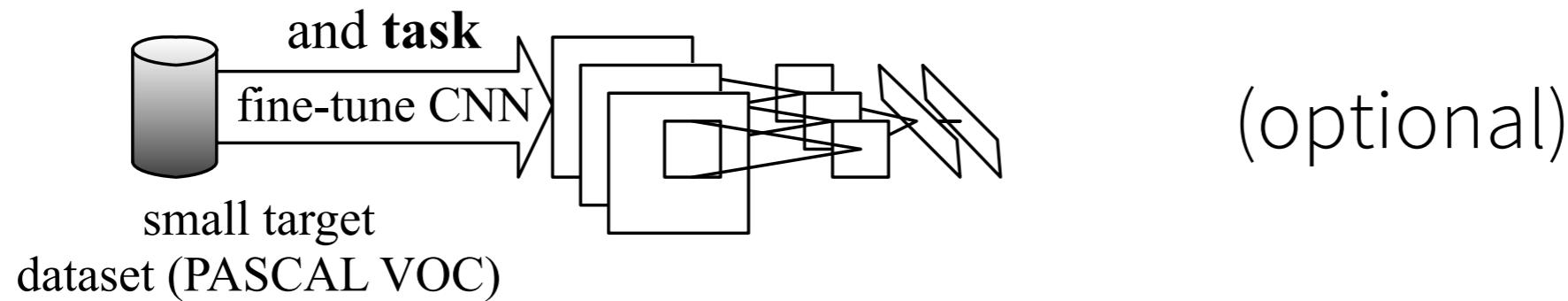


Training

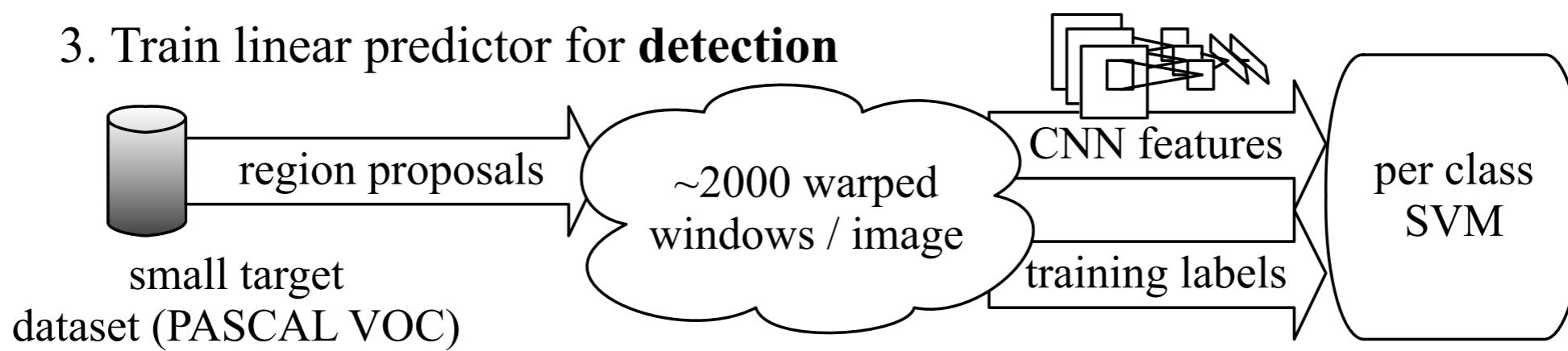
1. Pre-train CNN for **image classification**



2. Fine-tune CNN on **target dataset** and **task**

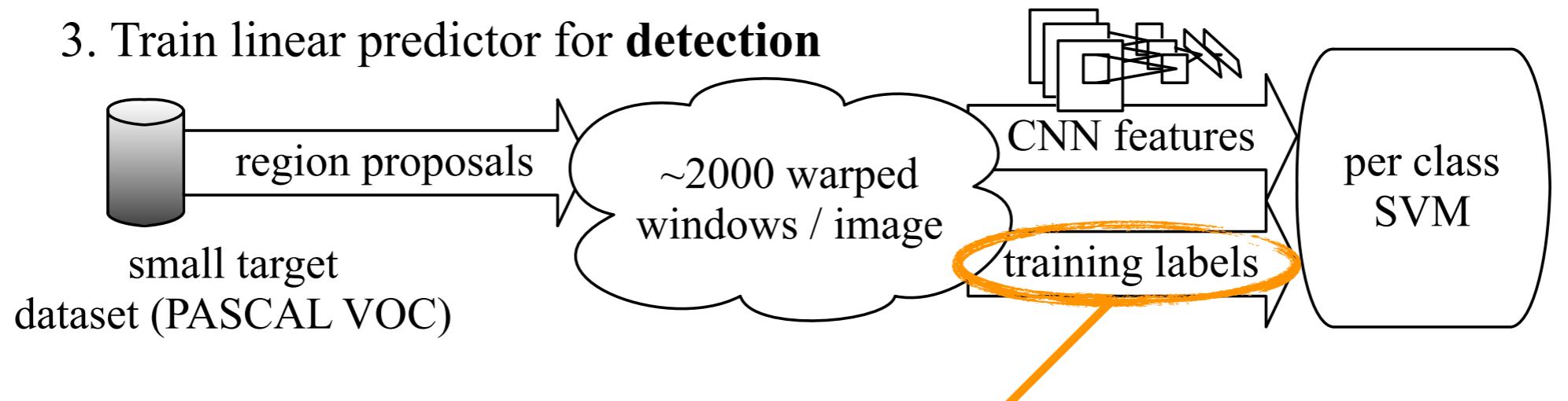


3. Train linear predictor for **detection**



Training labels

3. Train linear predictor for **detection**

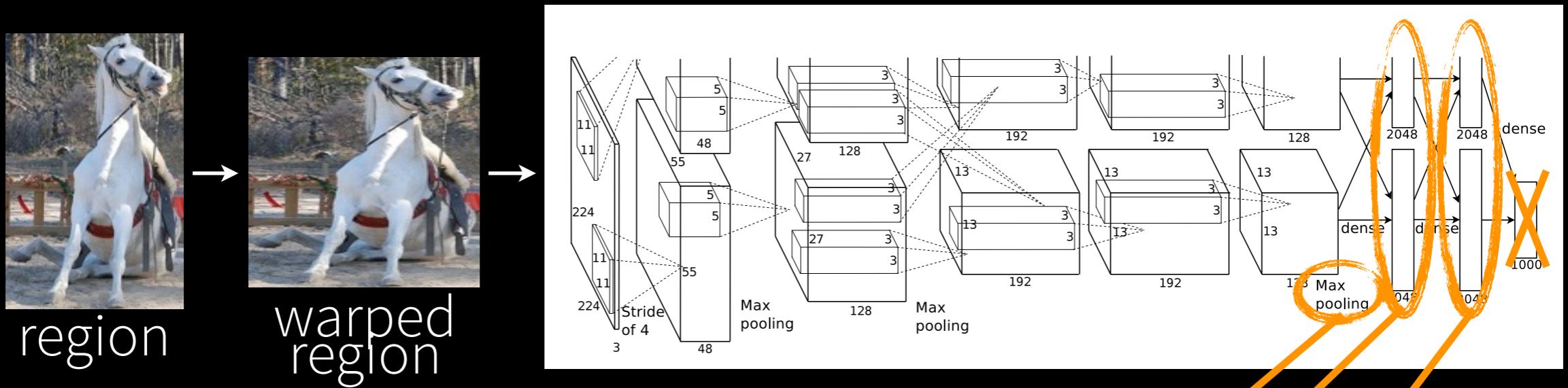


labeling protocol

positives = ground truth

negatives = max IoU < 0.3

CNN features for detection



pool₅: $6 \times 6 \times 256 = 9216$ -dim

6.4% / 15% non-zero

fc₆: 4096-dimensional

71.2% / 20% nz

fc₇: 4096-dimensional

100% / 20% nz

Results

reference	VOC 2007	VOC 2010
DPM v5 (Girshick et al. 2011)	33.7%	29.6%
UVA sel. search (Uijlings et al. 2012)		35.1%
Regionlets (Wang et al. 2013)	41.7%	39.7%

metric: mean average precision (higher is better)

Results

pre-trained
only

	VOC 2007	VOC 2010
DPM v5 (Girshick et al. 2011)	33.7%	29.6%
UVA sel. search (Uijlings et al. 2012)		35.1%
Regionlets (Wang et al. 2013)	41.7%	39.7%
R-CNN pool ₅	40.1%	
R-CNN fc ₆		43.4%
R-CNN fc ₇	42.6%	

metric: mean average precision (higher is better)

Results

	VOC 2007	VOC 2010
DPM v5 (Girshick et al. 2011)	33.7%	29.6%
UVA sel. search (Uijlings et al. 2012)		35.1%
Regionlets (Wang et al. 2013)	41.7%	39.7%
R-CNN pool ₅	40.1%	
R-CNN fc ₆	43.4%	
R-CNN fc ₇	42.6%	
fine-tuned	R-CNN FT pool ₅	42.1%
	R-CNN FT fc ₆	47.2%
	R-CNN FT fc ₇	48%
		43.5%

metric: mean average precision (higher is better)

Results – update

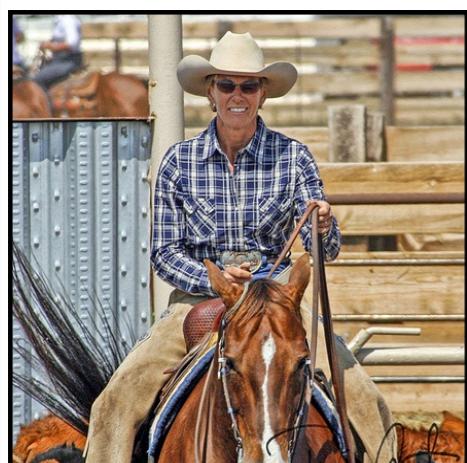
pre-trained
Only

	VOC 2007	VOC 2010
DPM v5 (Girshick et al. 2011)	33.7%	29.6%
UVA sel. search (Uijlings et al. 2012)		35.1%
Regionlets (Wang et al. 2013)	41.7%	39.7%
R-CNN pool ₅	40.1% 44.0%	
R-CNN fc ₆	43.4% 46.2%	
R-CNN fc ₇	42.6% 43.5%	

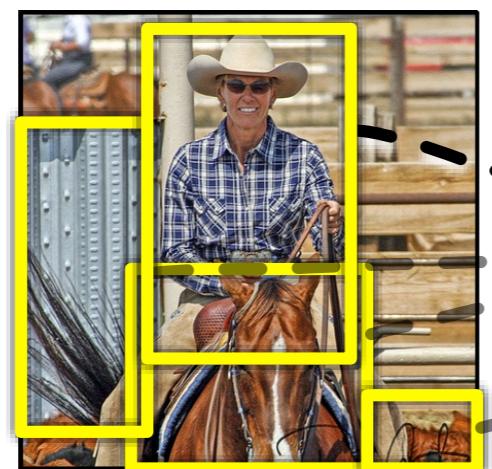
metric: mean average precision (higher is better)

CV and DL together

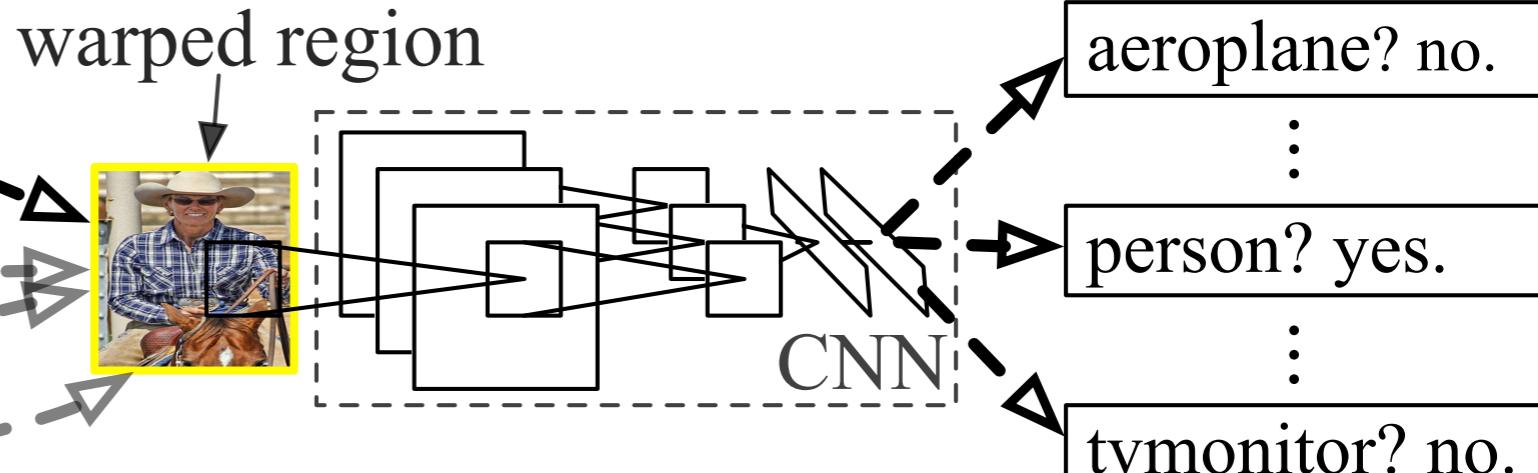
Good features are not enough!



1. Input
image



2. Extract region
proposals (~2k)



3. Compute
CNN features

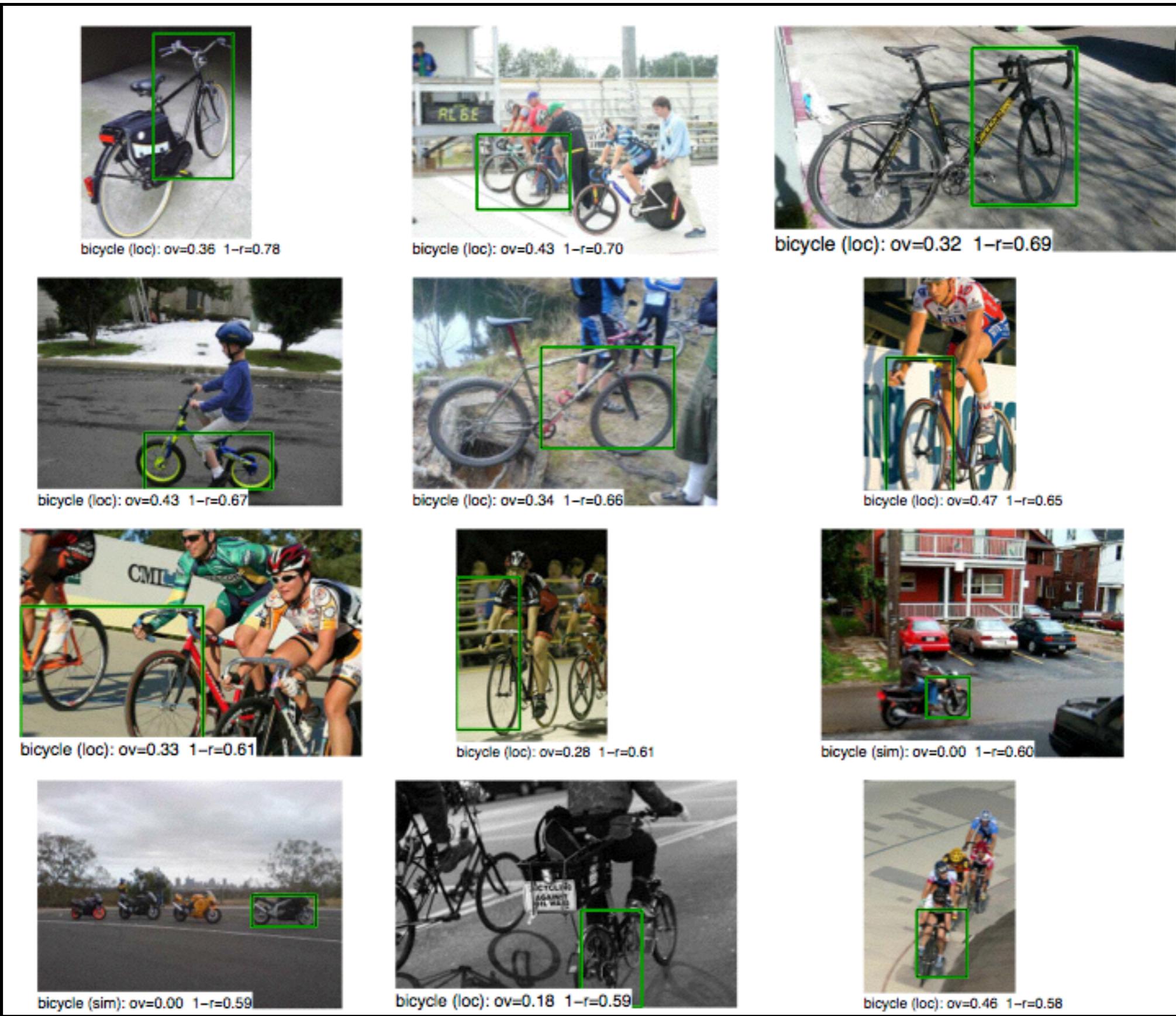
4. Classify
regions

Computer
Vision

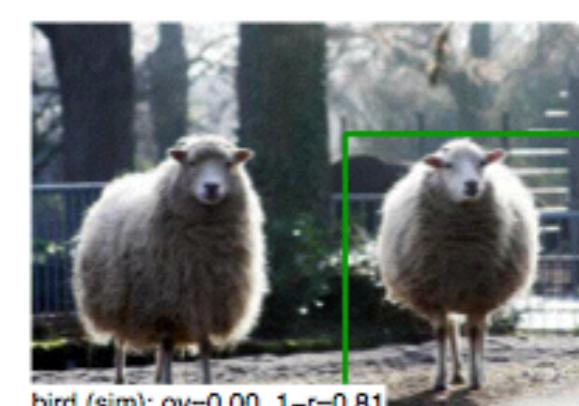
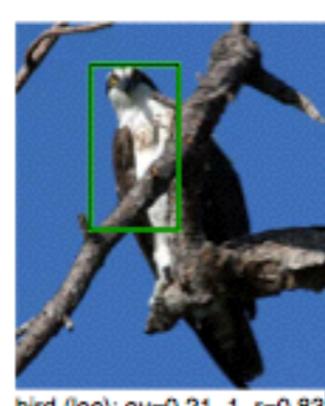
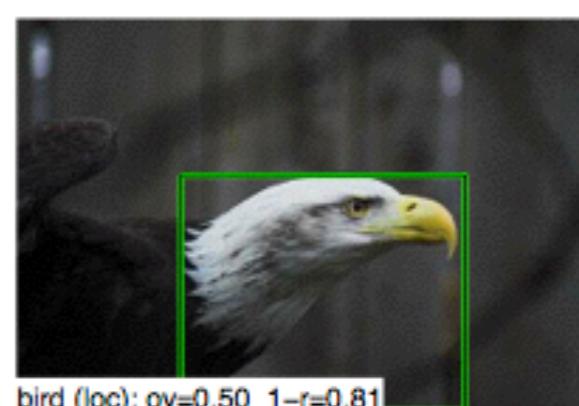
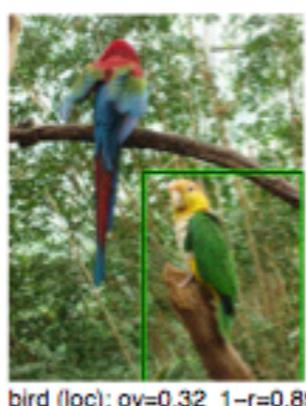
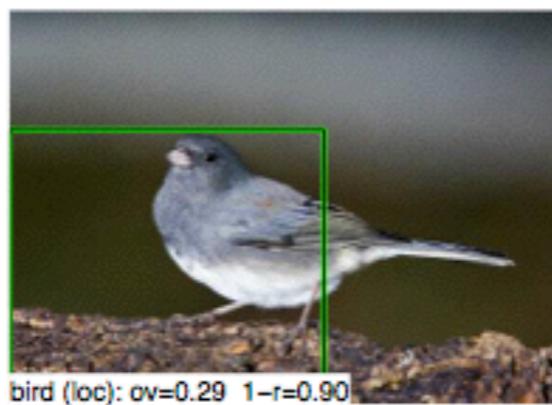
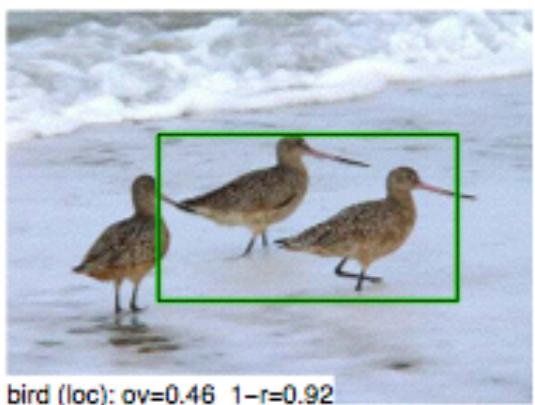
Deep
Learning

Computer
Vision

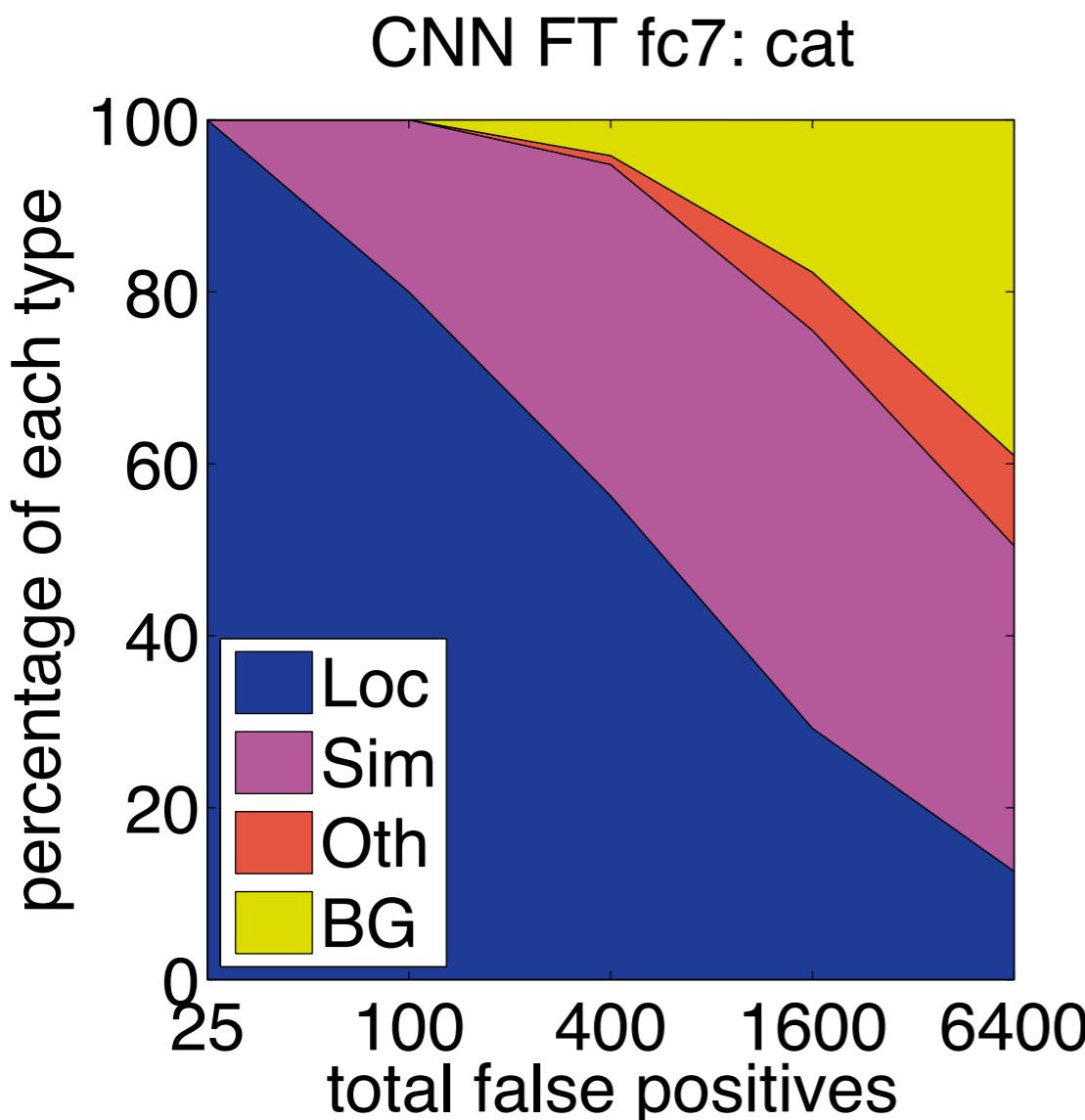
Top bicycle FPs (AP 62.5%)



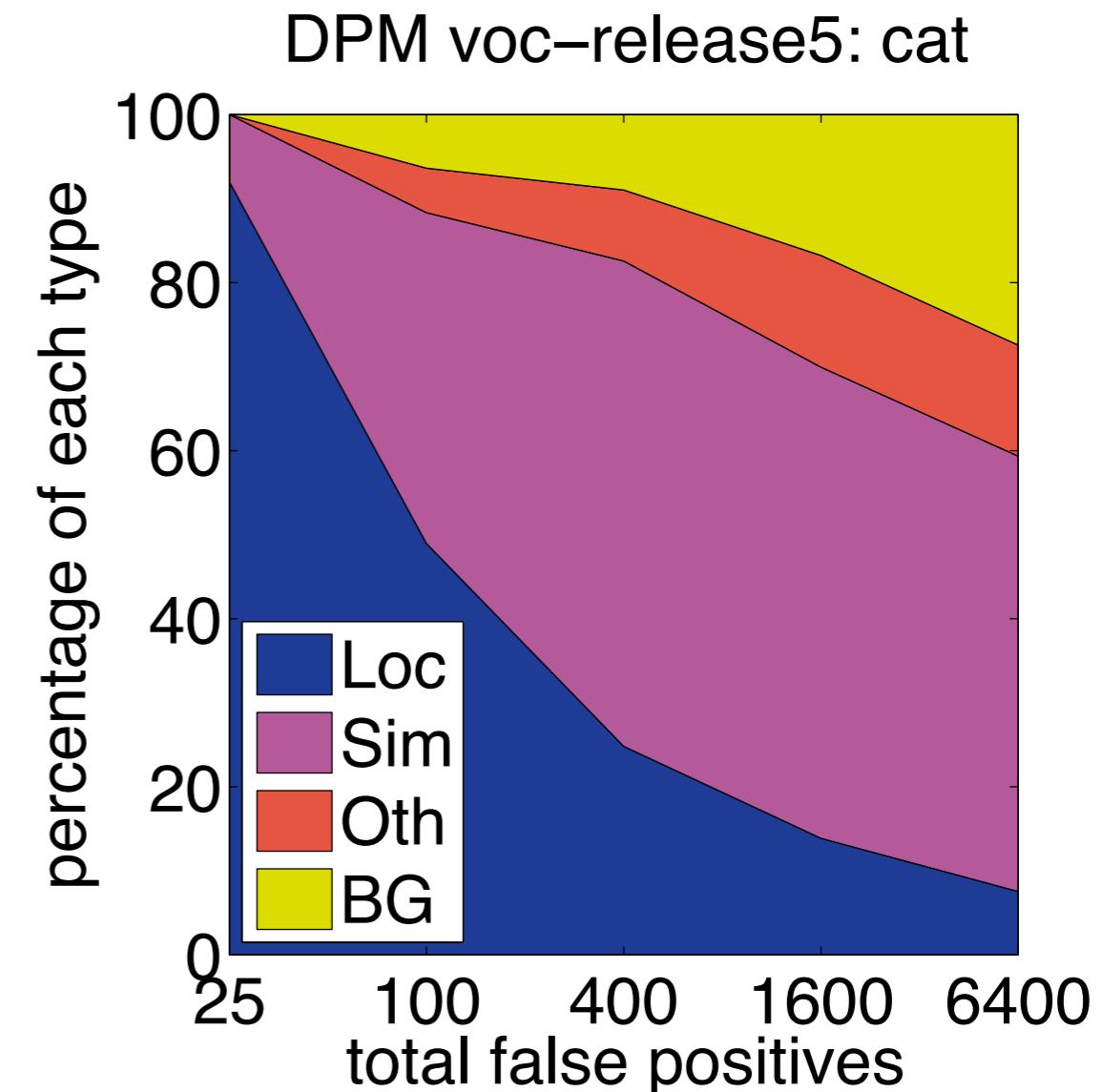
Top bird FPS (AP 41.4%)



False positive types: cat



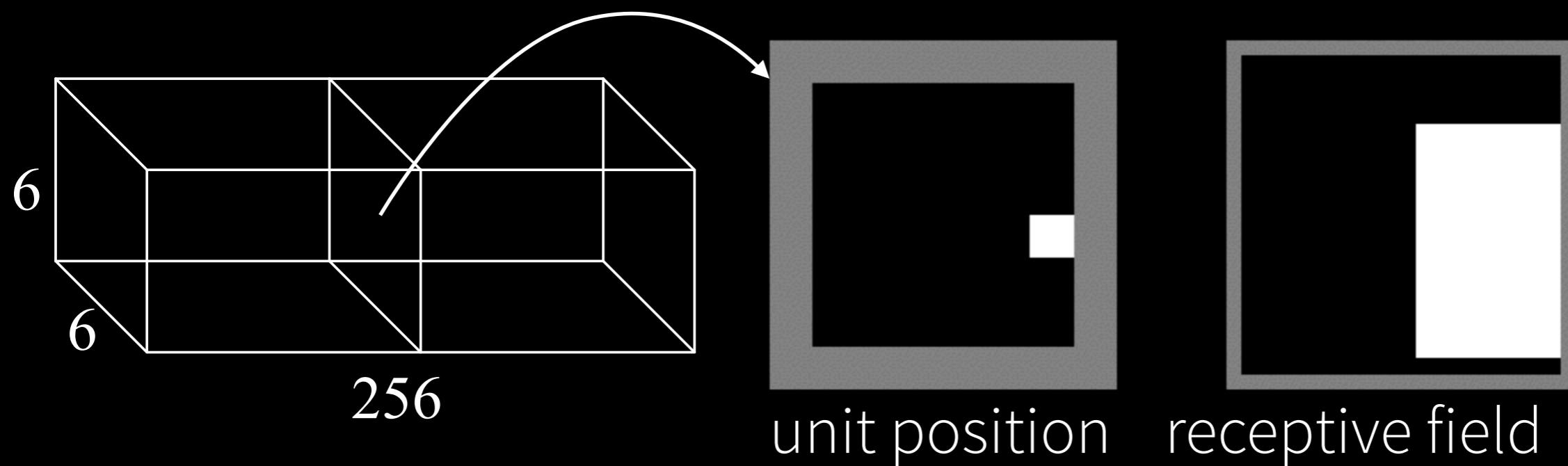
AP 56.3%



AP 23.0%

Visualizing features

- > What does pool₅ learn?
- > Recap:
 - > pool₅: max-pooled output of last conv. layer
 - > 6 × 6 spatial structure (with 256 channels)
 - > receptive field size 163 × 163 (of 224 × 224)

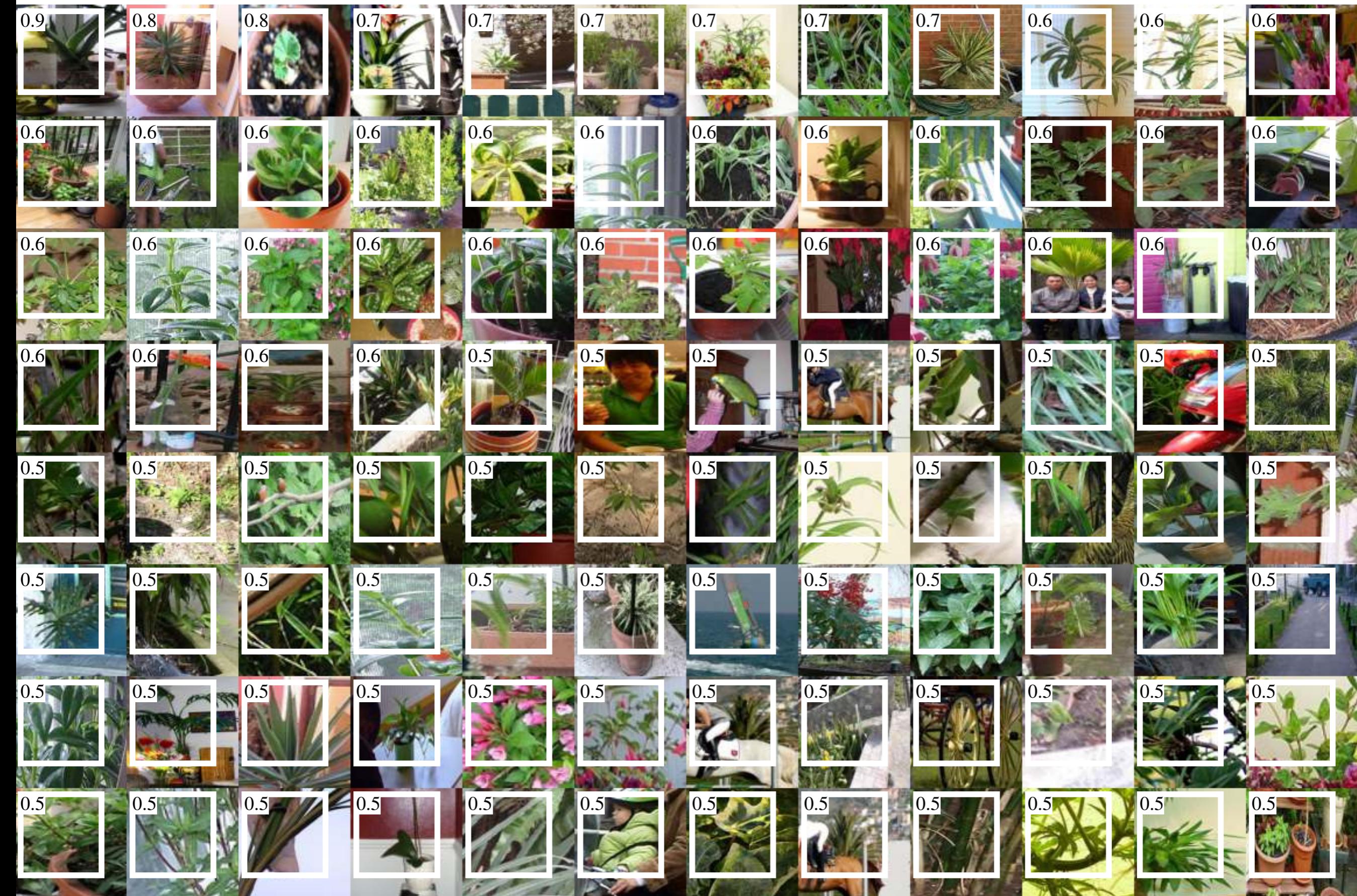


Visualization method

- > Select a unit in pool₅
- > Run it as a detector
- > Show top-scoring regions
- > Non-parametric, lets unit “speak for itself”

(Used ~10 million held-out regions.)

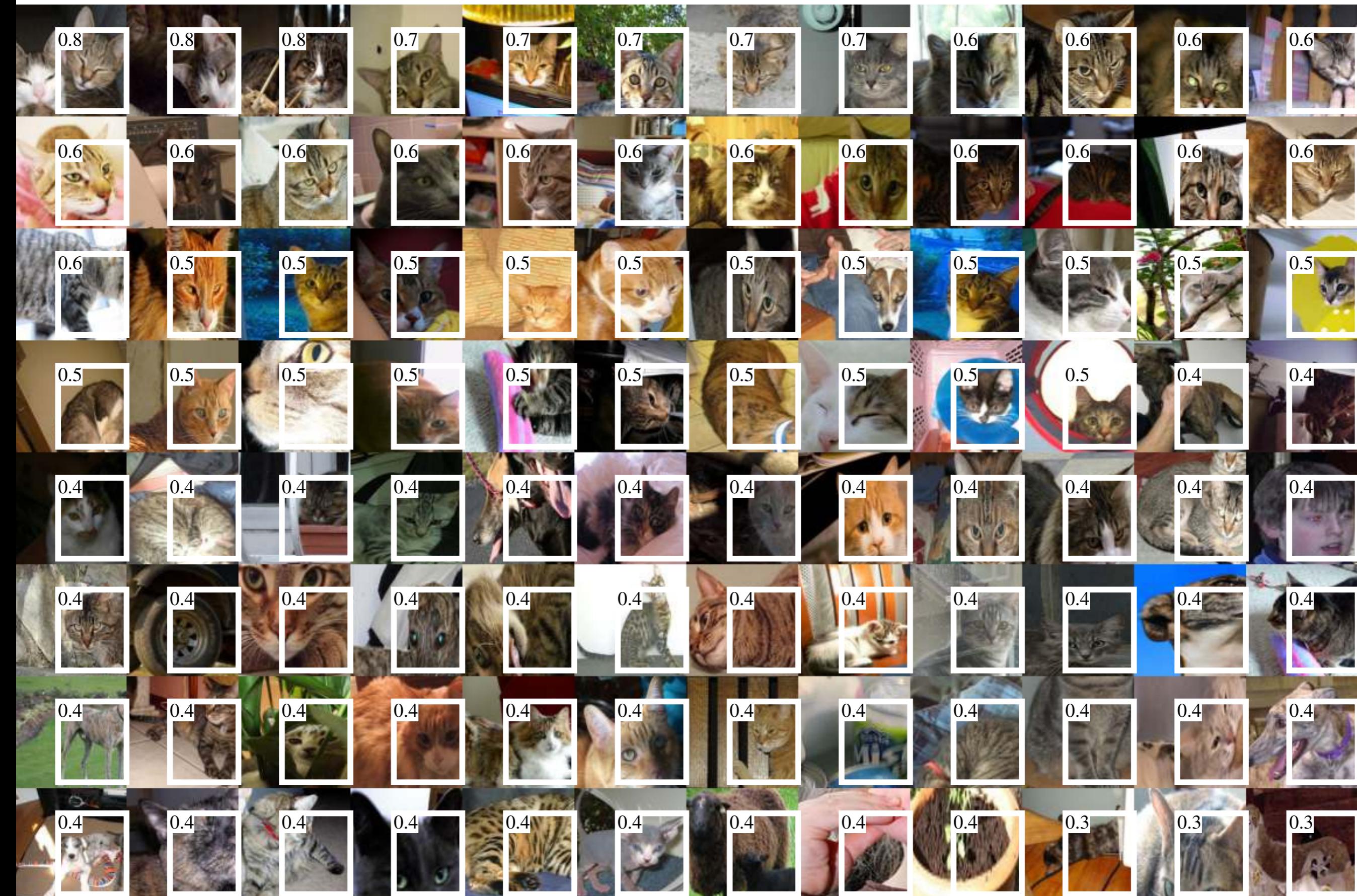
pool5 feature: (3,3,42) (top 1 – 96)



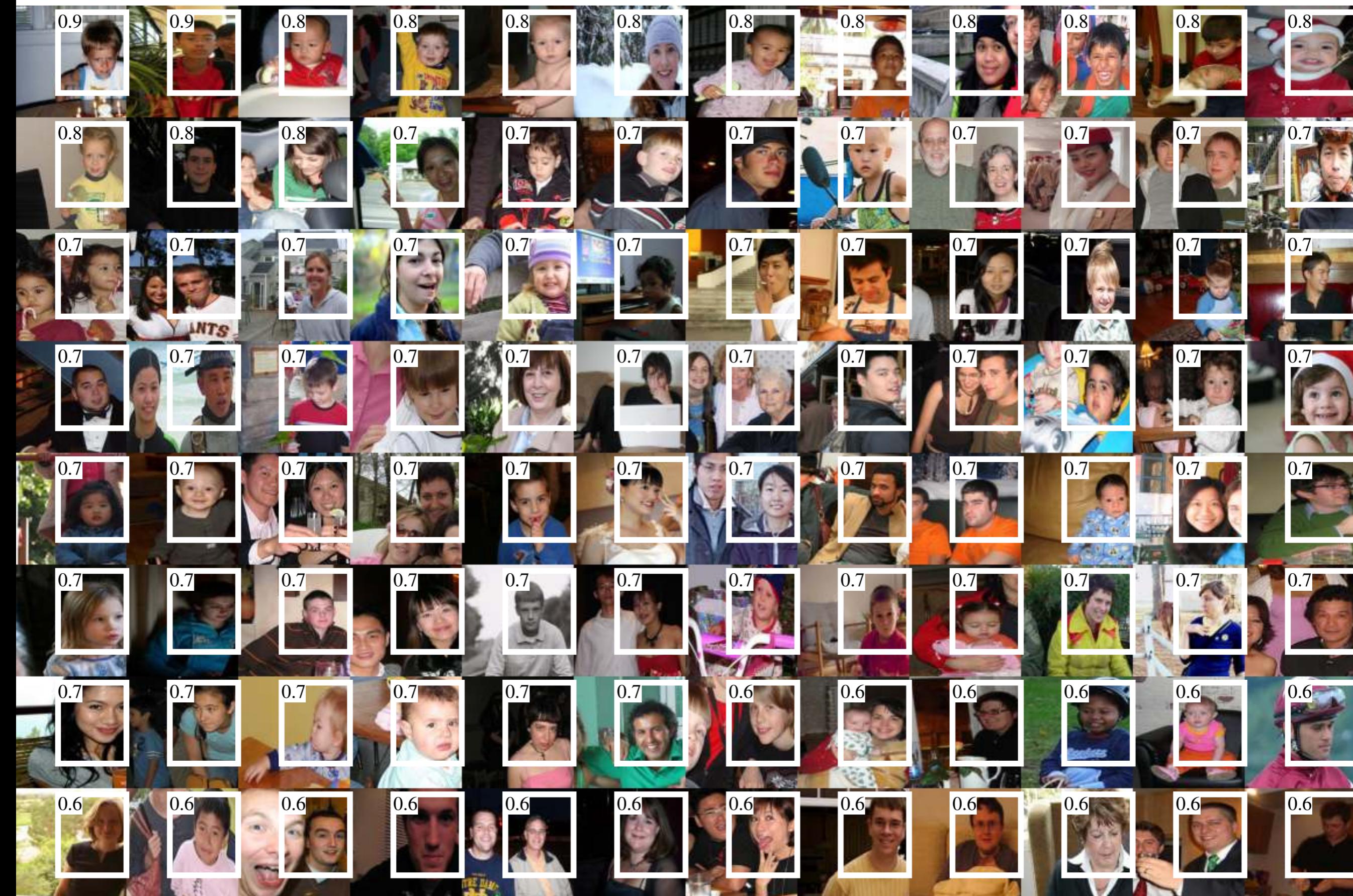
pool5 feature: (3,4,80) (top 1 – 96)



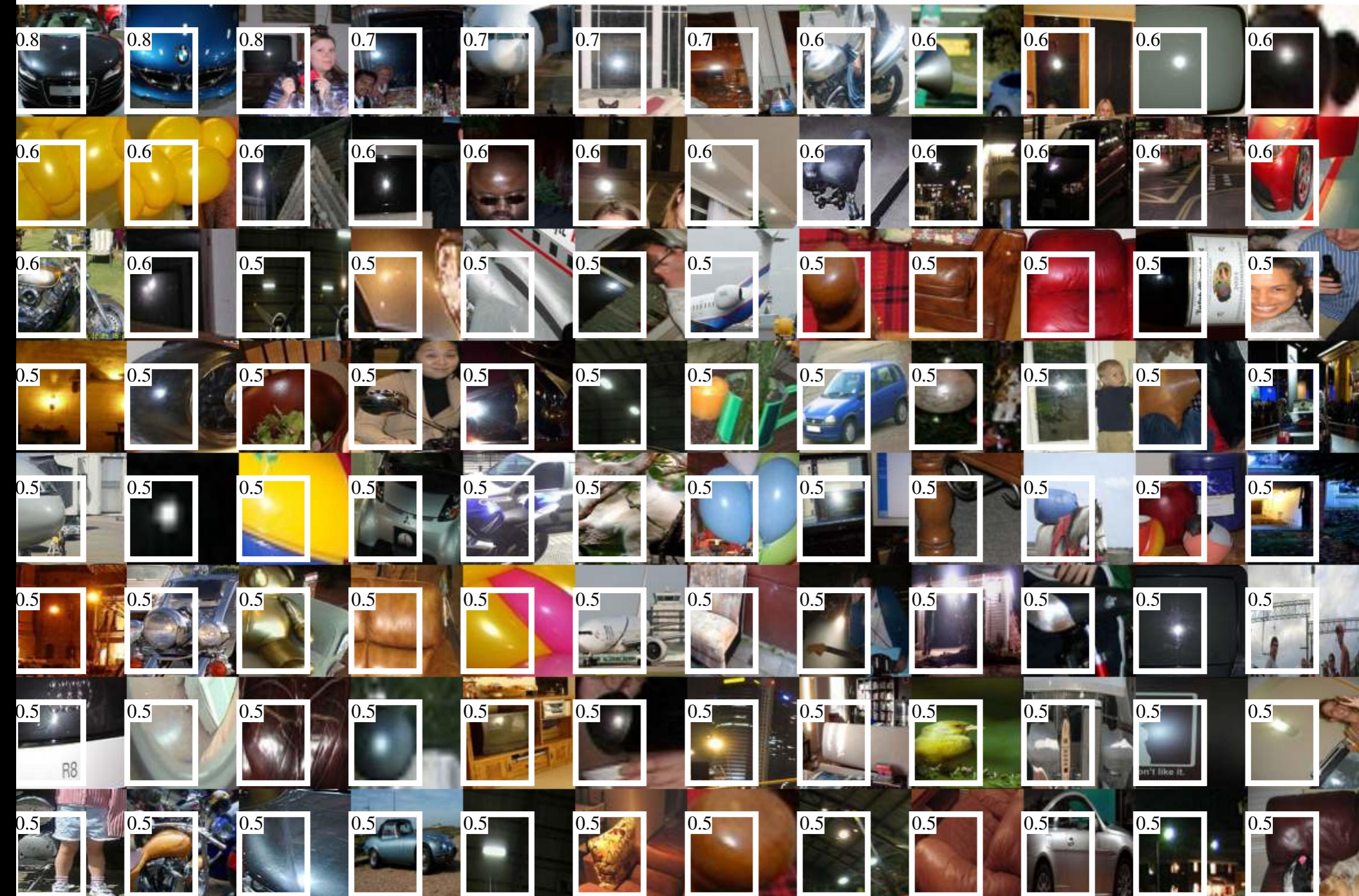
pool5 feature: (4,5,110) (top 1 – 96)



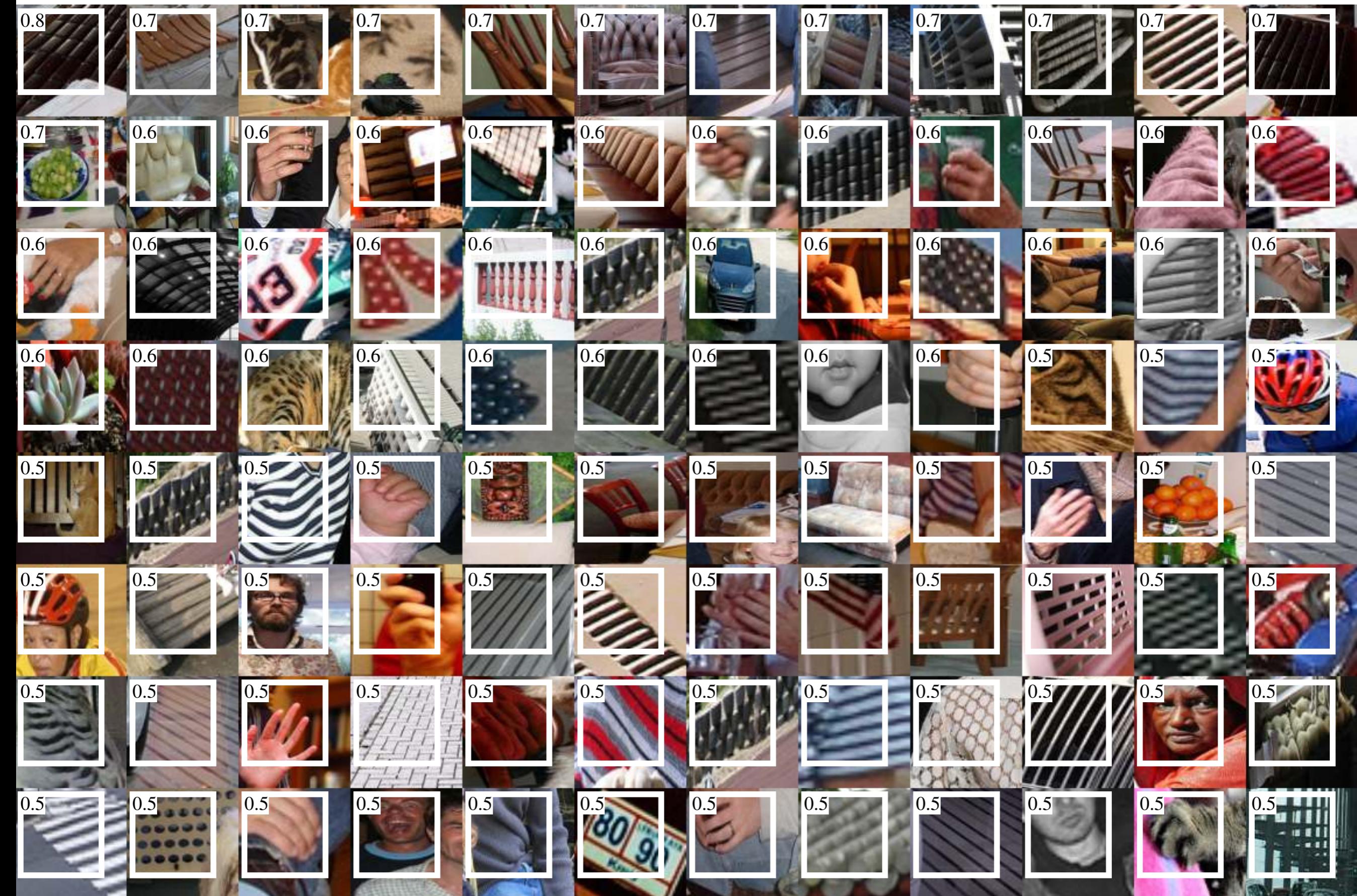
pool5 feature: (3,5,129) (top 1 - 96)



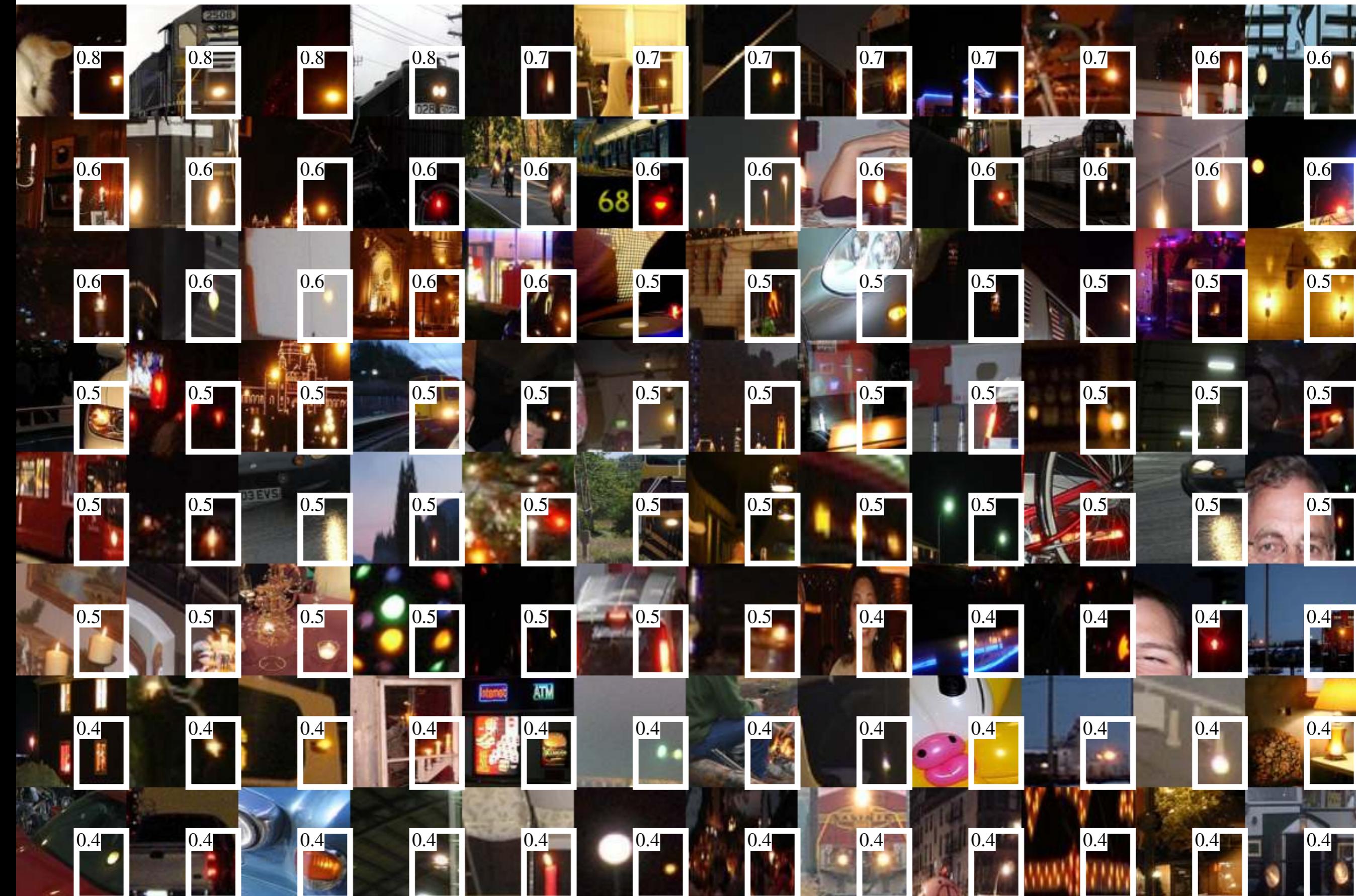
pool5 feature: (4,2,26) (top 1 – 96)



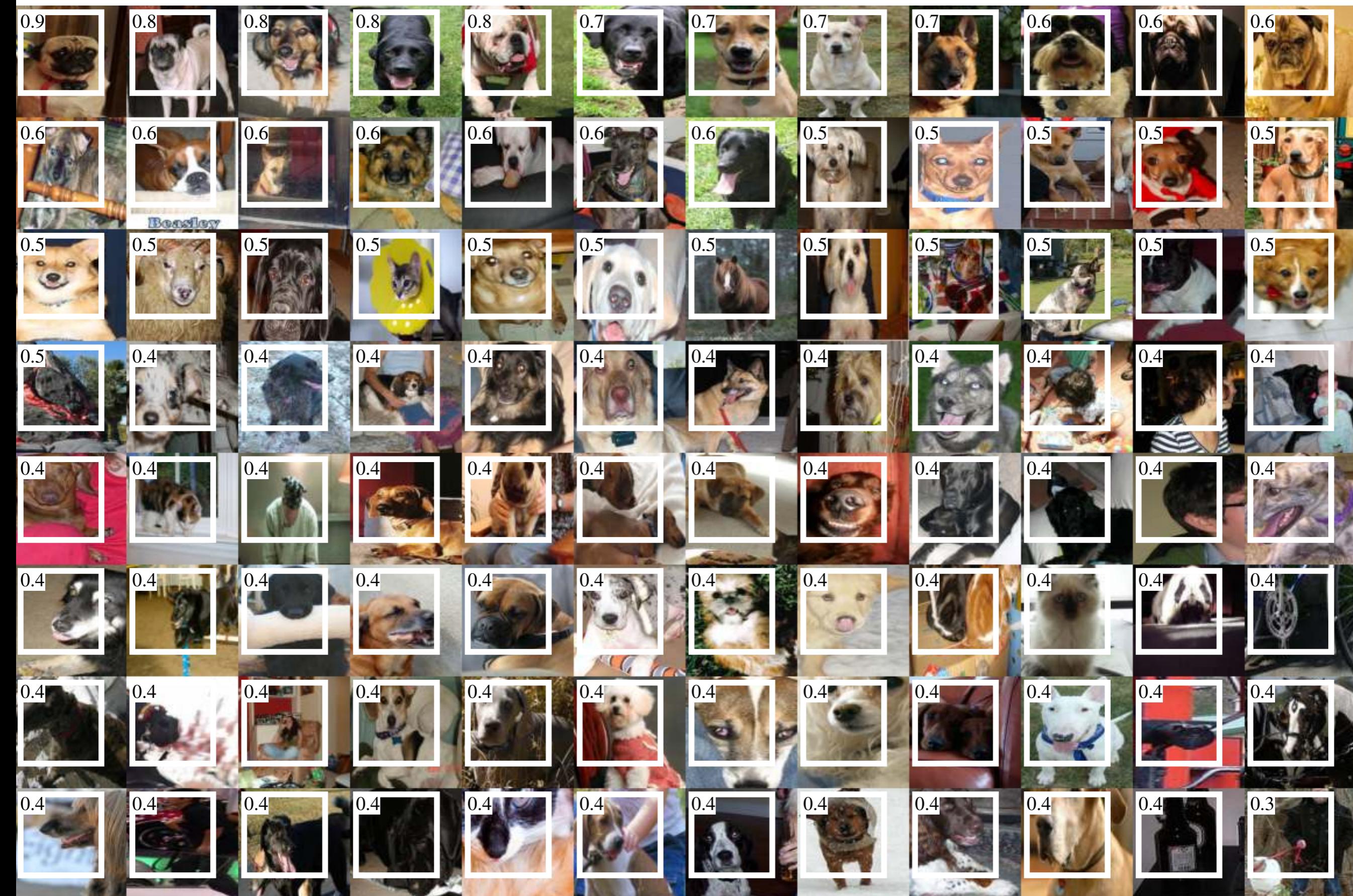
pool5 feature: (3,3,39) (top 1 – 96)



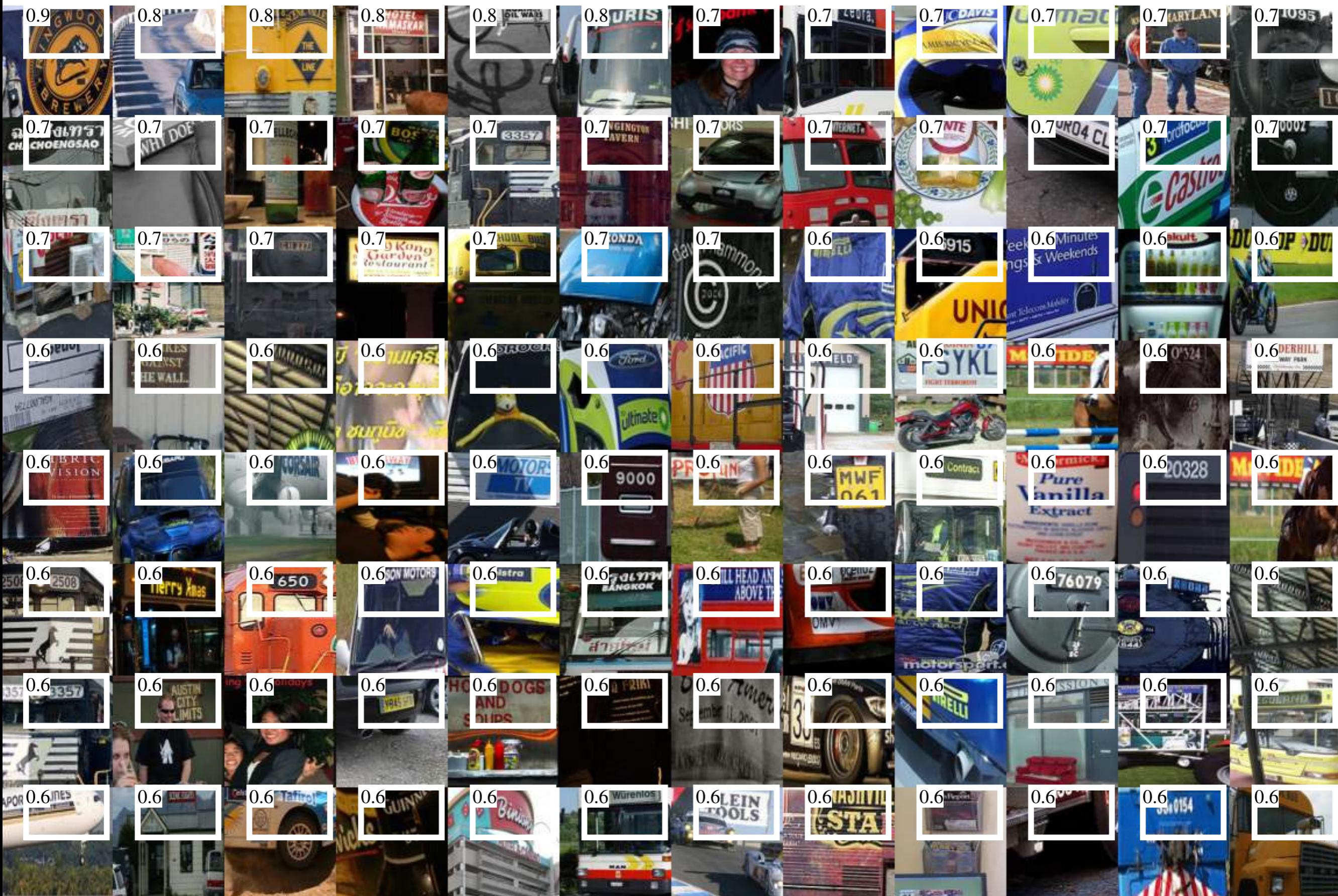
pool5 feature: (5,6,53) (top 1 – 96)



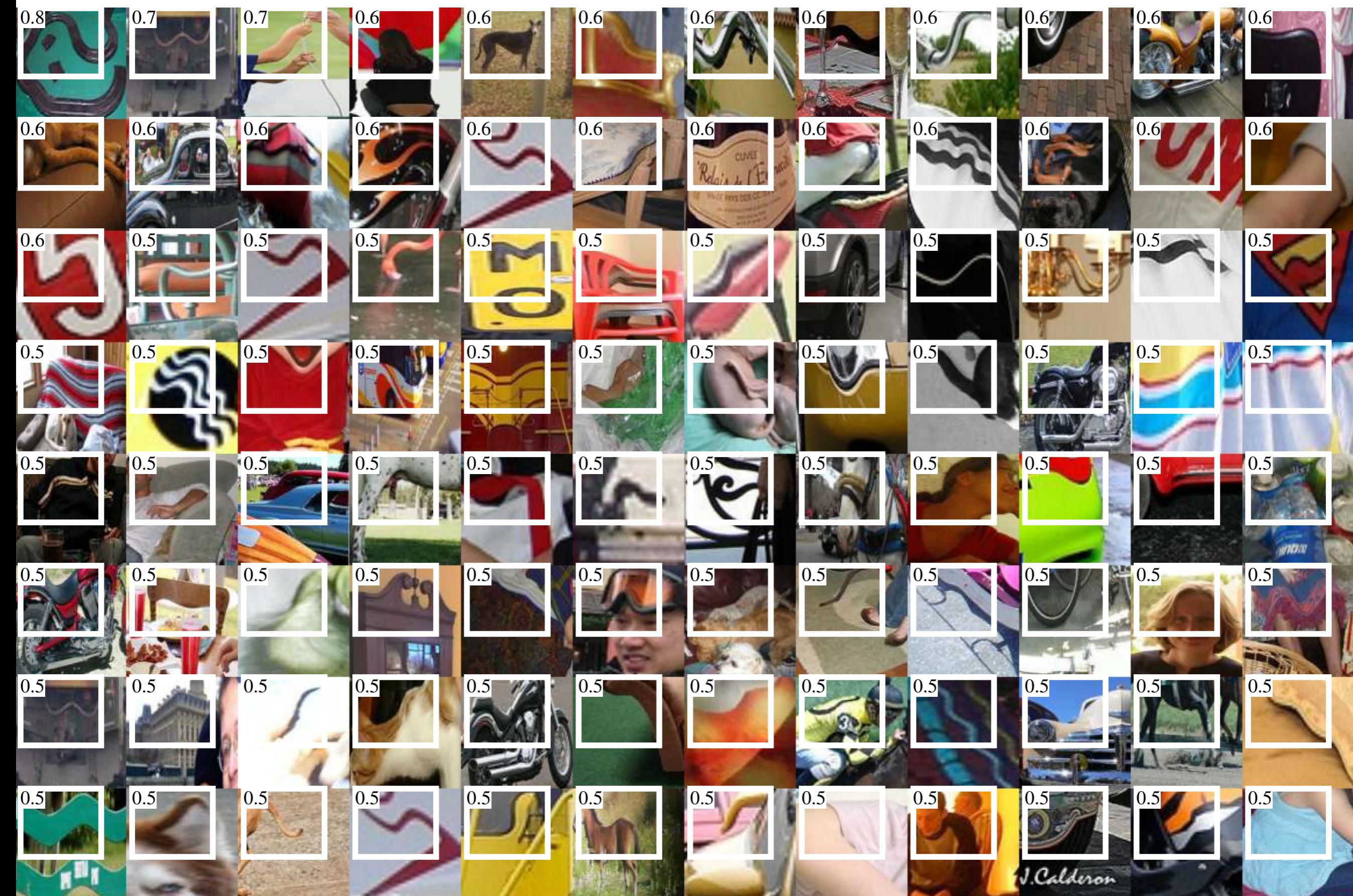
pool5 feature: (3,3,139) (top 1 – 96)



pool5 feature: (1,4,138) (top 1 - 96)



pool5 feature: (2,3,210) (top 1 – 96)



Semantic segmentation

	VOC 2011 test
1. UCB Regions and Parts	40.8%
2. Bonn O ₂ P	47.6%
3. R-CNN full+fg fc ₆	47.9%

metric: mean segmentation accuracy (higher is better)



(Carreira & Sminchisescu)

