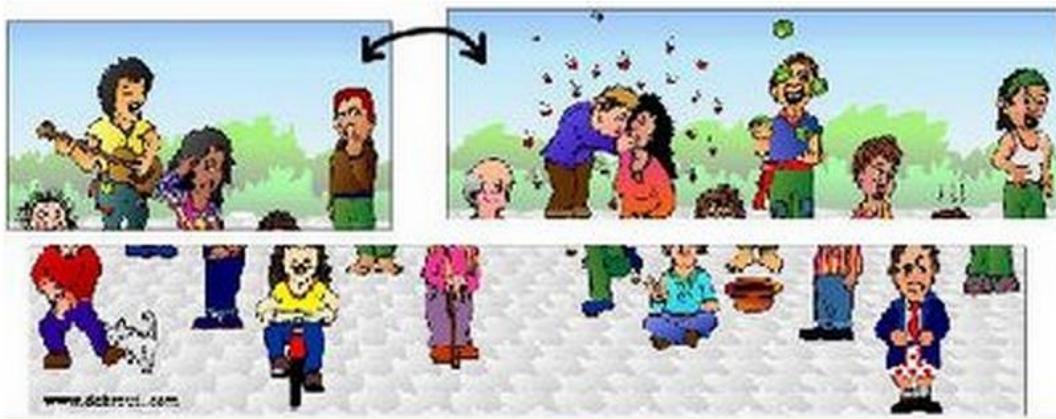


Recognize Objects



Gluing Pixels

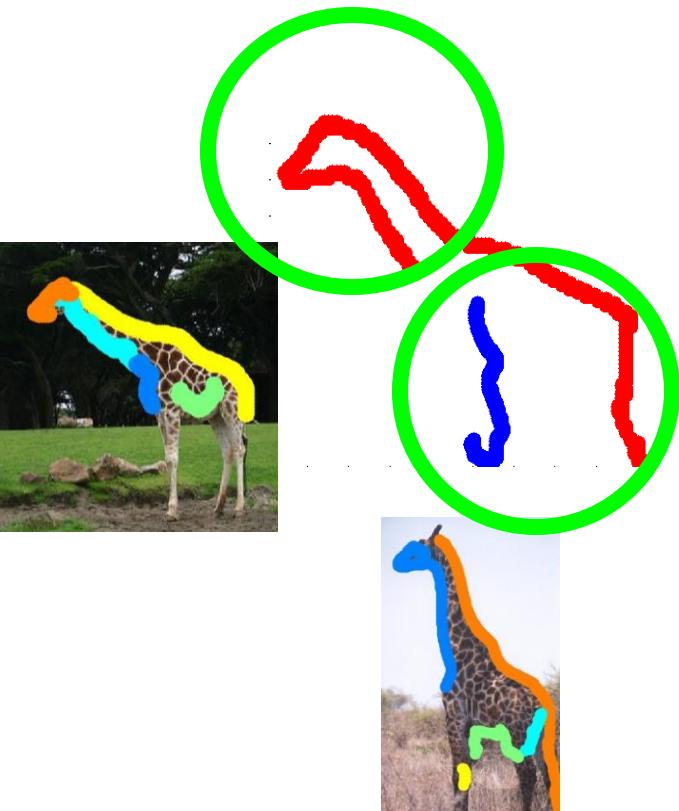


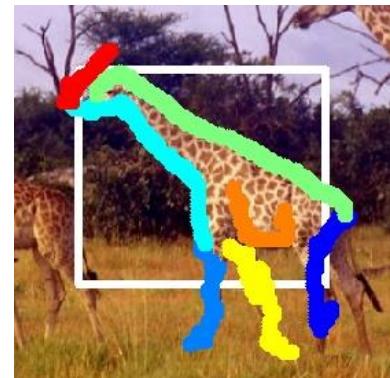
CIS 680, Spring '17



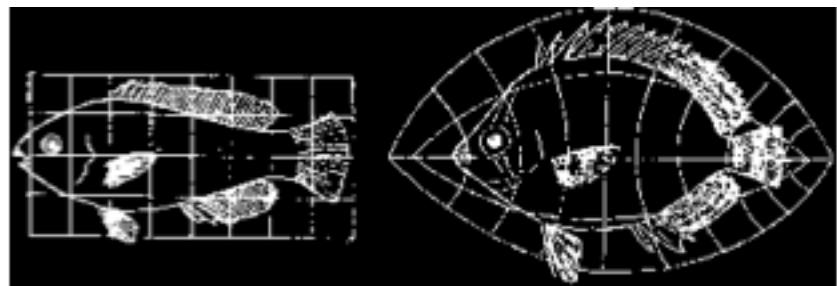
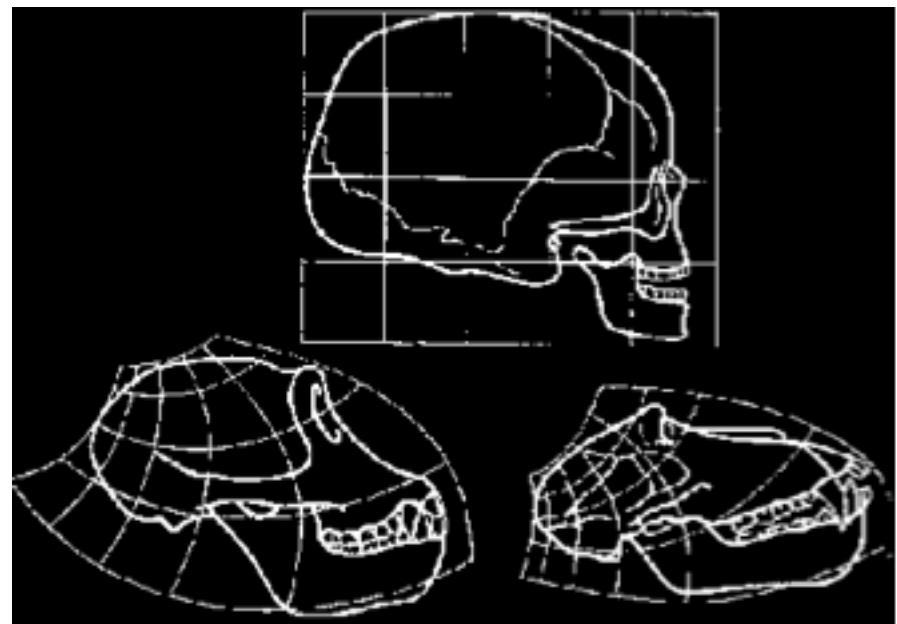
recognize objects: Describe and Discriminate

- **Describe positives** by commonalities of shape
- **Discriminate negatives** by uniqueness of shape

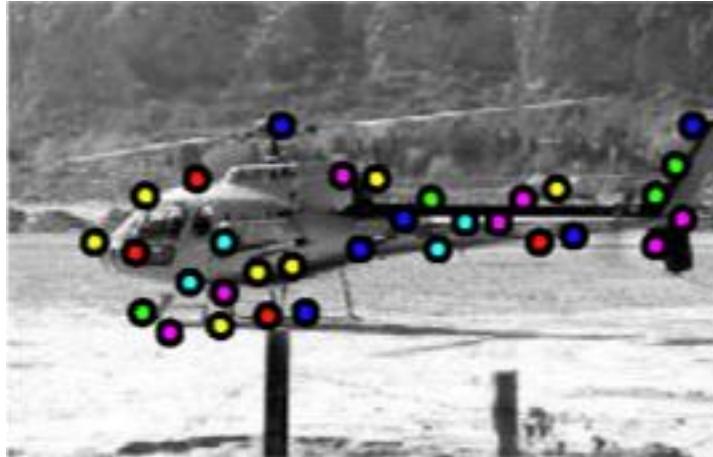




recognize objects: Frame 'it' right

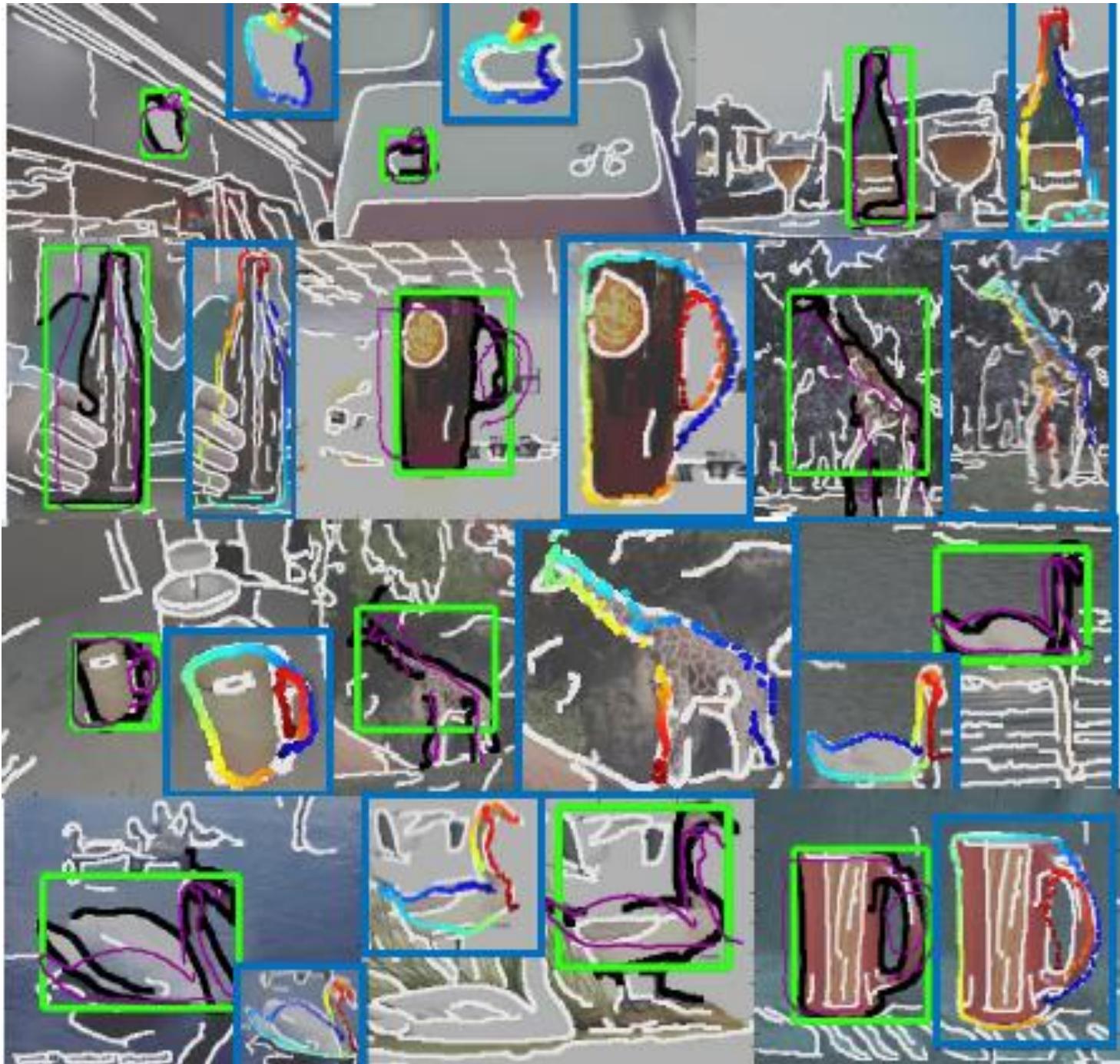


Challenge 1: maintain discriminative under deformation



Similarity of appearance near feature points

Similarity in configuration of the feature points

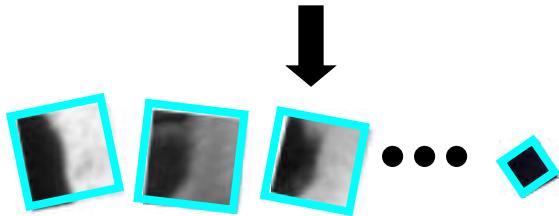


maintain discriminative under deformation

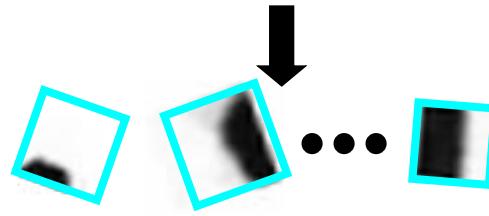
The Pyramid Match Kernel:
Discriminative Classification with
Sets of Image Features

Kristen Grauman
Trevor Darrell

Sets of features

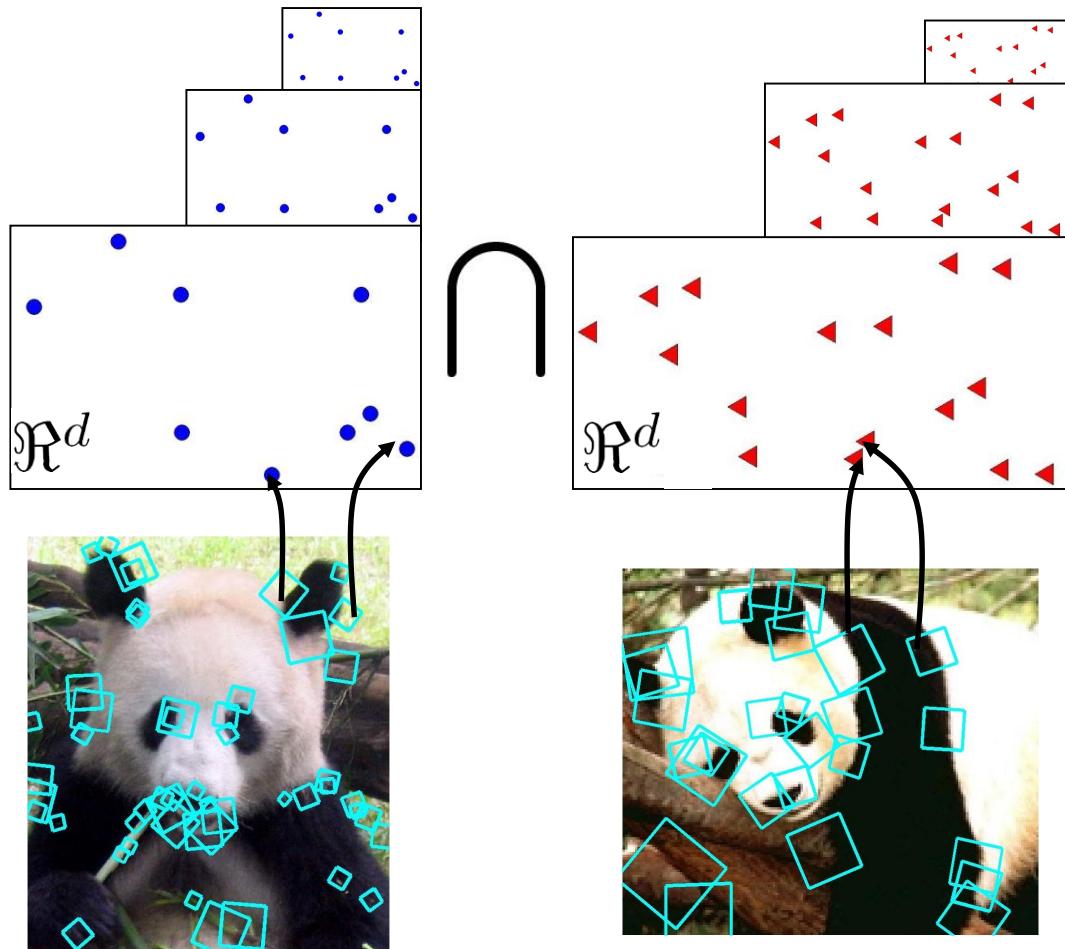


$$X = \{\vec{x}_1, \dots, \vec{x}_m\}$$



$$Y = \{\vec{y}_1, \dots, \vec{y}_n\}$$

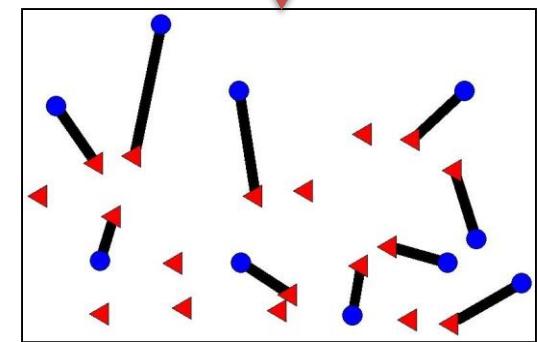
Pyramid match



$$\mathbf{X} = \{\vec{\mathbf{x}}_1, \dots, \vec{\mathbf{x}}_m\} \quad \vec{\mathbf{x}}_i \in \Re^d$$

$$\mathbf{Y} = \{\vec{\mathbf{y}}_1, \dots, \vec{\mathbf{y}}_n\} \quad \vec{\mathbf{y}}_i \in \Re^d$$

Related:
Earth Moving Distance
(EMD)



optimal partial
matching

$$\max_{\pi: \mathbf{X} \rightarrow \mathbf{Y}} \sum_{\mathbf{x}_i \in \mathbf{X}} \mathcal{S}(\mathbf{x}_i, \pi(\mathbf{x}_i))$$

Pyramid match overview

Pyramid match kernel measures similarity of a partial matching between two sets:

- Place multi-dimensional, multi-resolution grid over point sets
- Consider points matched at finest resolution where they fall into same grid cell
- Approximate similarity between matched points with worst case similarity at given level

No explicit search for matches!

Pyramid match kernel

Approximate
partial match
similarity

$$K_{\Delta} = \sum_{i=0}^L w_i N_i$$

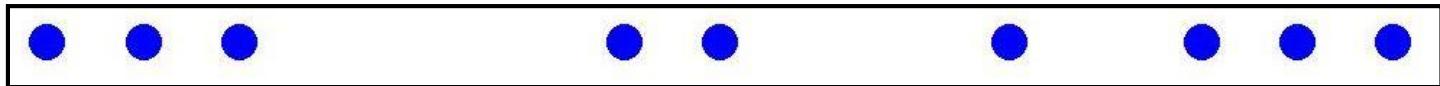
Number of newly matched pairs at level i

Measure of difficulty of a match at level i

```
graph TD; A["Number of newly matched pairs at level i"] --> B["wi Ni"]; C["Measure of difficulty of a match at level i"] --> B;
```

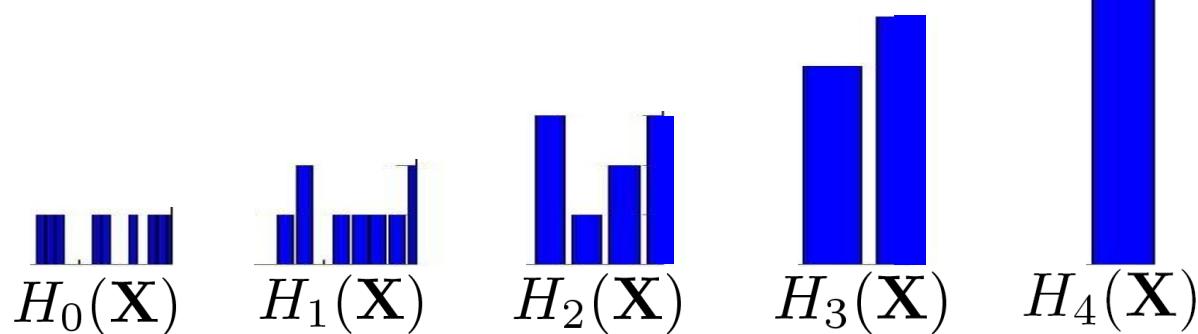
Feature extraction

$$\mathbf{X} = \{\vec{\mathbf{x}}_1, \dots, \vec{\mathbf{x}}_m\}, \quad \vec{\mathbf{x}}_i \in \Re^d$$


$$d = 1$$



Histogram pyramid: level i has
bins of size 2^i

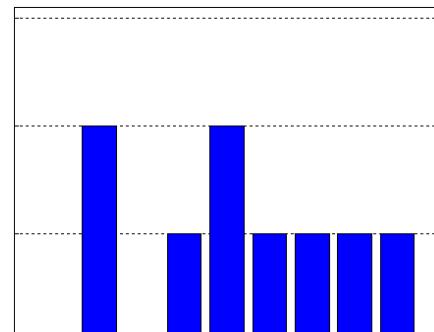
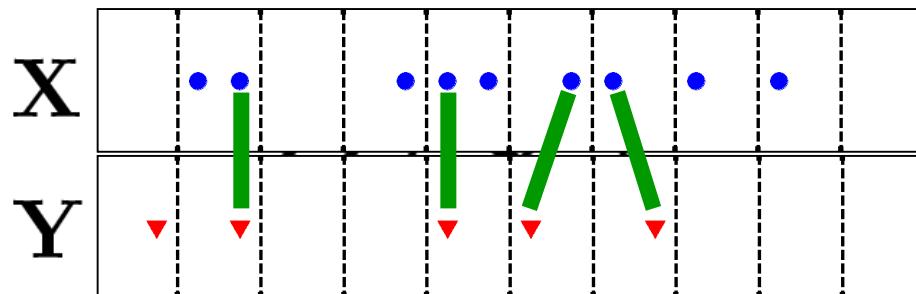


$$\Psi(\mathbf{X}) = [H_0(\mathbf{X}), \dots, H_L(\mathbf{X})]$$

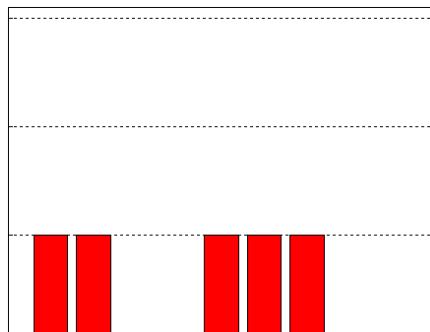
Counting matches

Histogram
intersection

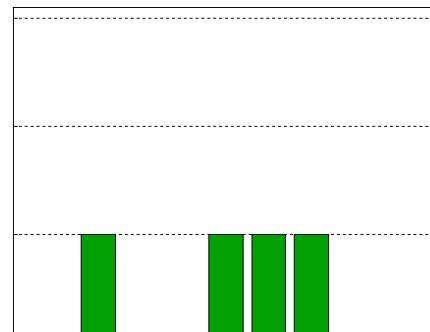
$$\mathcal{I}(H(\mathbf{X}), H(\mathbf{Y})) = \sum_{j=1}^r \min(H(\mathbf{X})_j, H(\mathbf{Y})_j)$$



$$H(\mathbf{X})$$



$$H(\mathbf{Y})$$



$$\mathcal{I}(H(\mathbf{X}), H(\mathbf{Y})) = 4$$

Counting new matches

Histogram
intersection

$$\mathcal{I}(H(\mathbf{X}), H(\mathbf{Y})) = \sum_{j=1}^r \min(H(\mathbf{X})_j, H(\mathbf{Y})_j)$$

$$N_i = \overbrace{\mathcal{I}(H_i(\mathbf{X}), H_i(\mathbf{Y})) - \mathcal{I}(H_{i-1}(\mathbf{X}), H_{i-1}(\mathbf{Y}))}^{\text{matches at this level}} - \overbrace{\mathcal{I}(H_{i-1}(\mathbf{X}), H_{i-1}(\mathbf{Y}))}^{\text{matches at previous level}}$$



Difference in histogram intersections across
levels counts *number of new pairs matched*

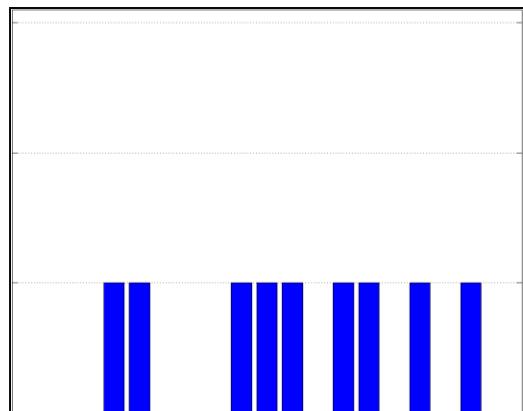
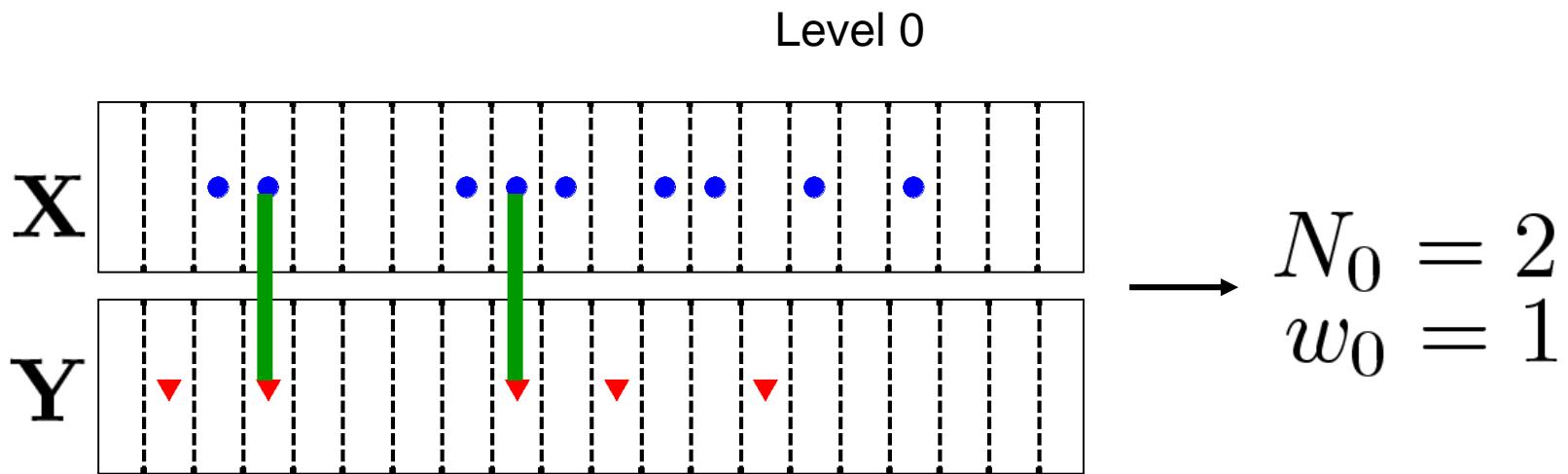
Pyramid match kernel

$$K_{\Delta} (\overbrace{\Psi(\mathbf{X}), \Psi(\mathbf{Y})}^{\text{histogram pyramids}}) = \sum_{i=0}^L \frac{1}{2^i} \left(\underbrace{\mathcal{I}(H_i(\mathbf{X}), H_i(\mathbf{Y})) - \mathcal{I}(H_{i-1}(\mathbf{X}), H_{i-1}(\mathbf{Y}))}_{\text{number of newly matched pairs at level } i} \right)$$

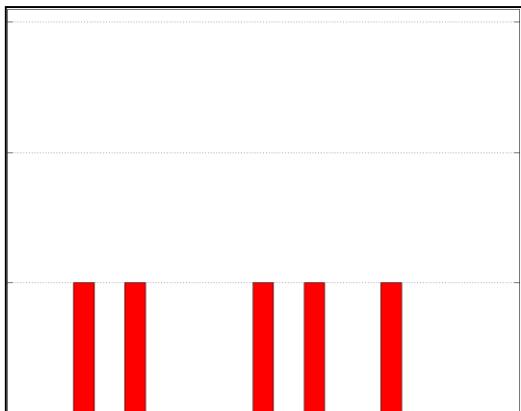
↑
measure of difficulty of
a match at level i

- Weights inversely proportional to bin size
- Normalize kernel values to avoid favoring large sets

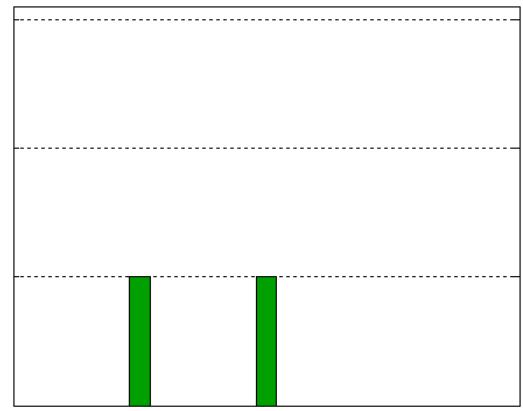
Example pyramid match



$$H_0(\mathbf{X})$$



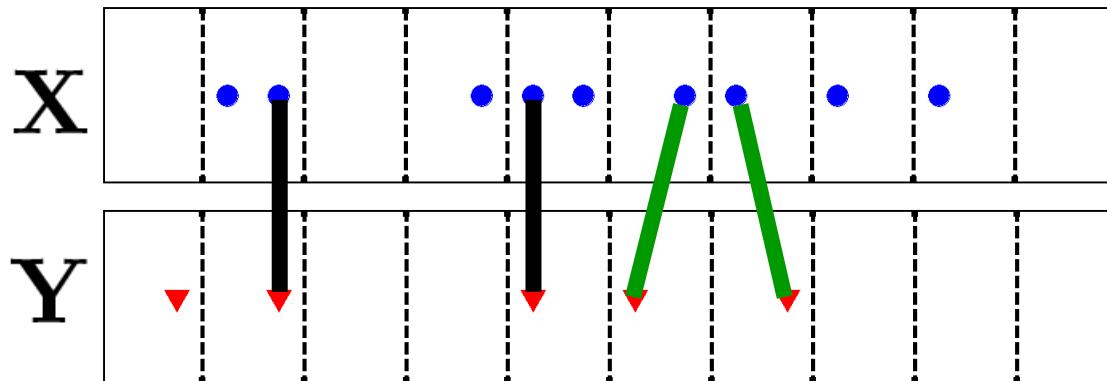
$$H_0(\mathbf{Y})$$



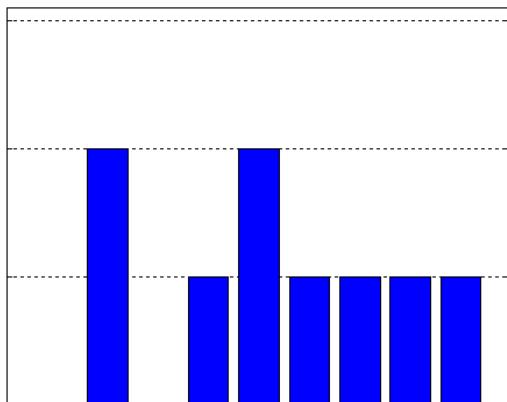
$$\mathcal{I}_0 = 2$$

Example pyramid match

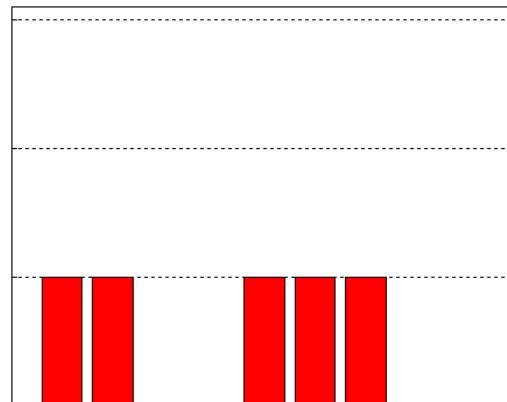
Level 1



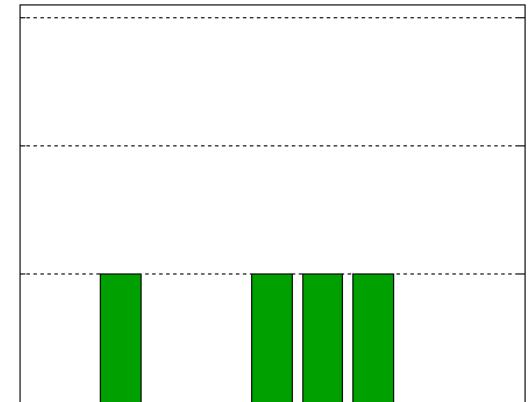
$$\rightarrow N_1 = 4 - 2 = 2$$
$$w_1 = \frac{1}{2}$$



$H_1(\mathbf{X})$



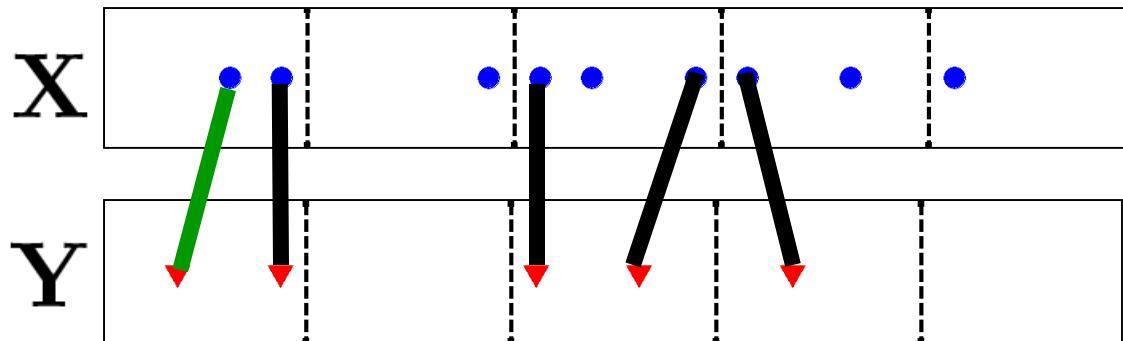
$H_1(\mathbf{Y})$



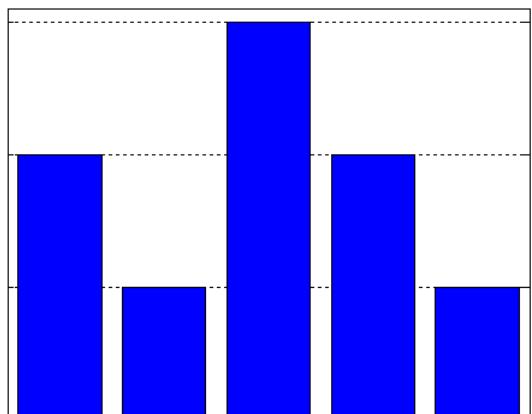
$\mathcal{I}_1 = 4$

Example pyramid match

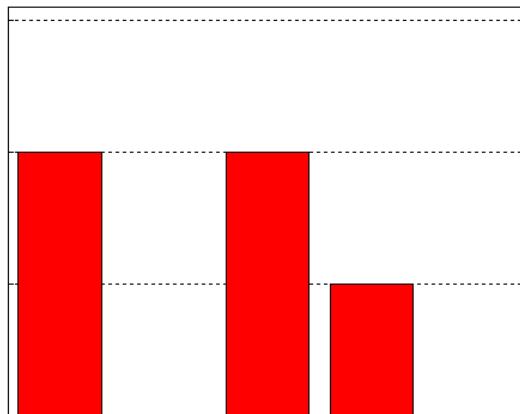
Level 2



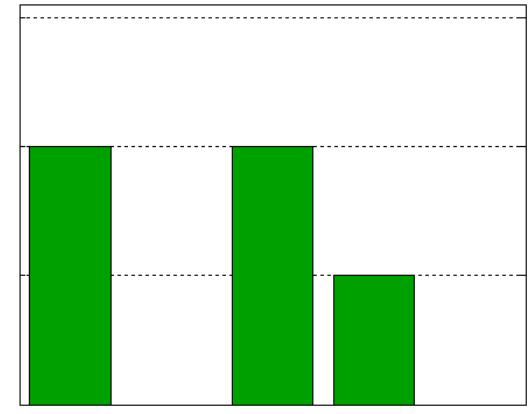
$$\rightarrow N_2 = 5 - 4 = 1$$
$$w_2 = \frac{1}{4}$$



$H_2(\mathbf{X})$



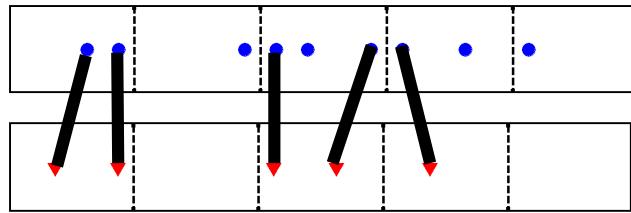
$H_2(\mathbf{Y})$



$\mathcal{I}_2 = 5$

Example pyramid match

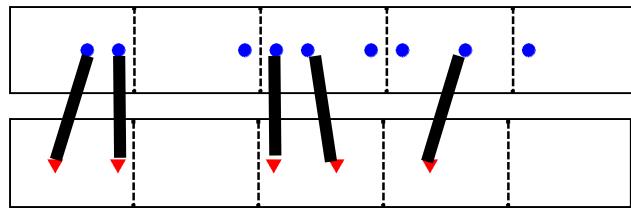
pyramid match



$$K_{\Delta} = \sum_{i=0}^L w_i N_i$$

$$= 1(2) + \frac{1}{2}(2) + \frac{1}{4}(1) = 3.25$$

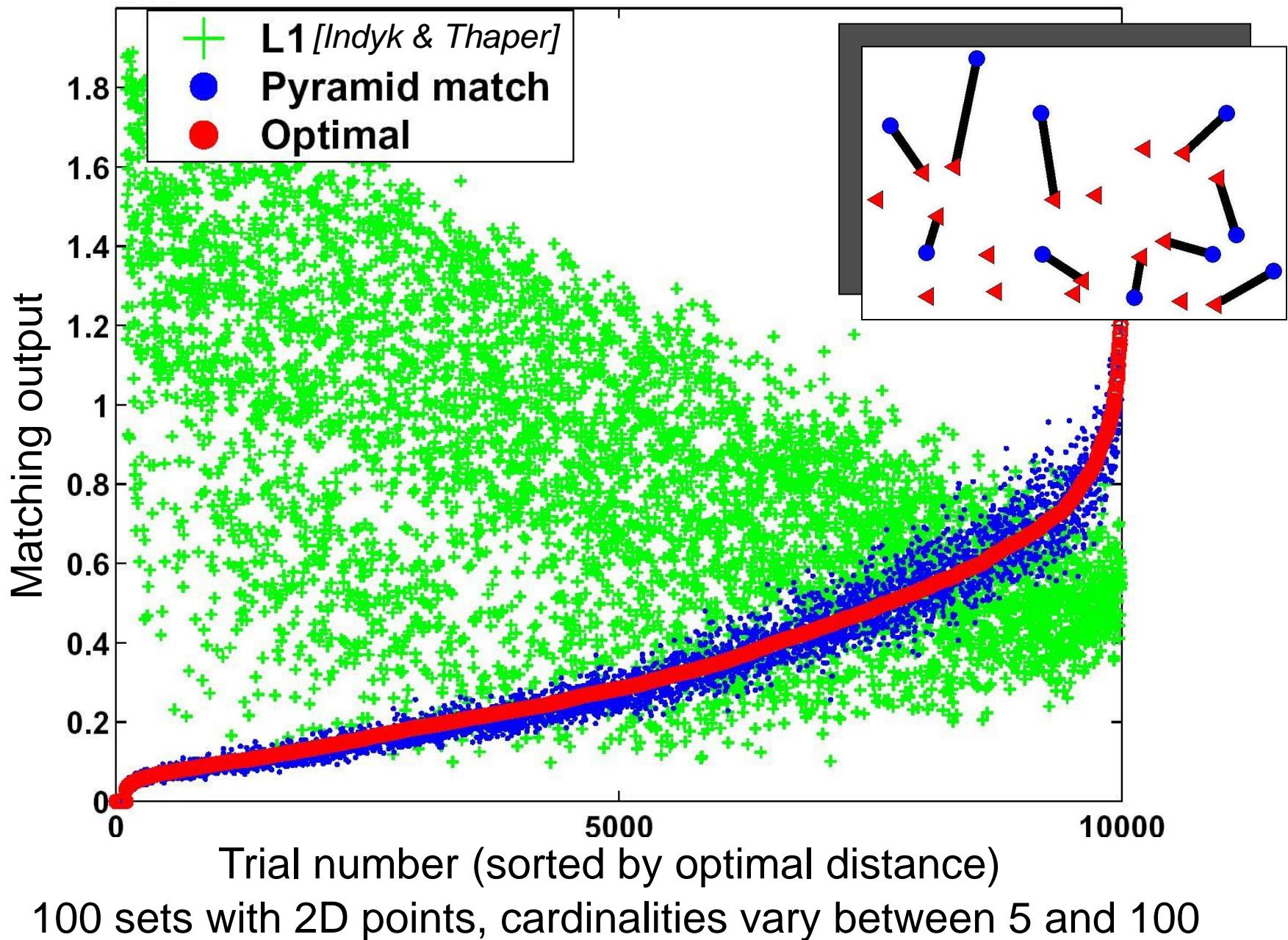
optimal match



$$K = \max_{\pi: \mathbf{X} \rightarrow \mathbf{Y}} \sum_{\mathbf{x}_i \in \mathbf{X}} \mathcal{S}(\mathbf{x}_i, \pi(\mathbf{x}_i))$$

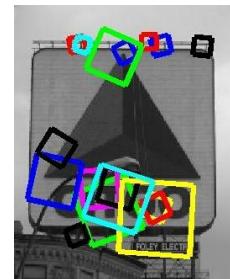
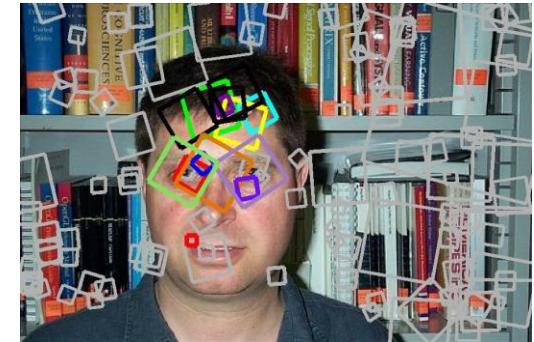
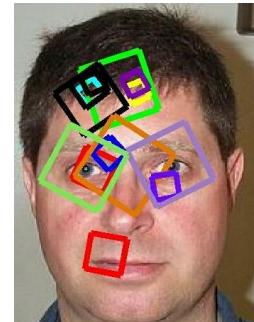
$$= 1(2) + \frac{1}{2}(3) = 3.5$$

Approximation of the optimal partial matching

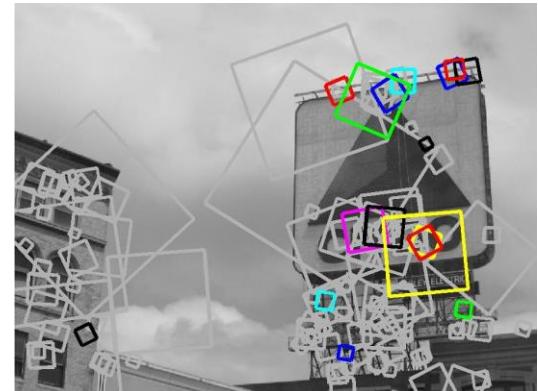


Localization

- Inspect intersections to obtain correspondences between features
- Higher confidence correspondences at finer resolution levels



target



observation

Challenge: maintain discriminative under occlusion



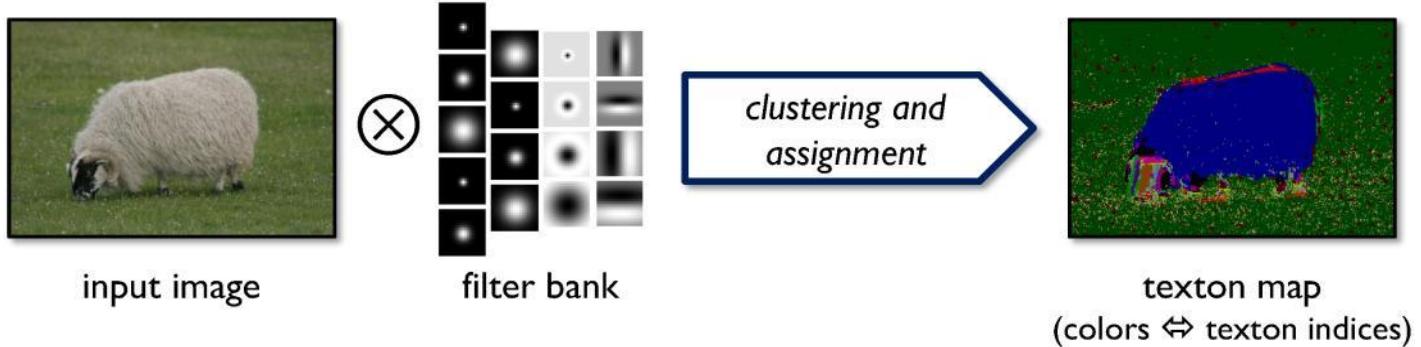
Challenge 2: maintain discriminative under occlusion

Per Pixel Classification

**Texton Boost for
pixel level semantic classification**

Jamie Shotton

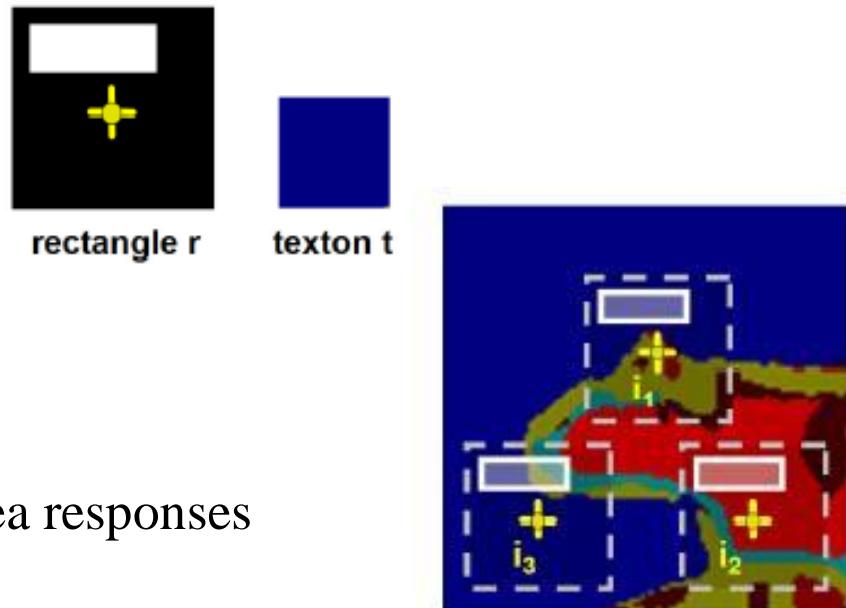
❑ Step 1: Texton Map generation (17 filters, K=400)



❑ Step 2: Shape Filter

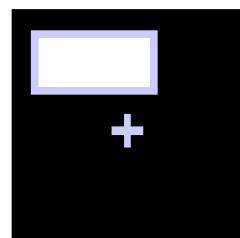
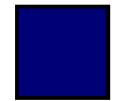
- For each texton t
 - ✓ Inputs
 - Texton Map
 - (Rectangle mask r , texton query t)
 - Pixel location i
 - ✓ Output
 - Area in rectangle mask that match t
- End result is a texton histogram of area responses

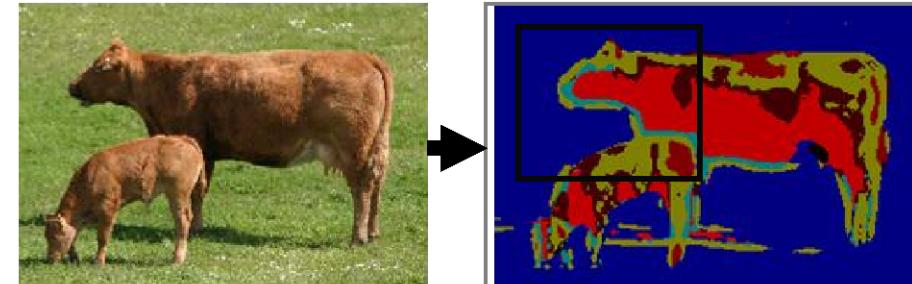
❑ How does this capture shape?



Shape Filters

up to 200 pixels

- ☐ Pair: $($  $,$  $)$
- ☐ rectangle r
- ☐ texton t

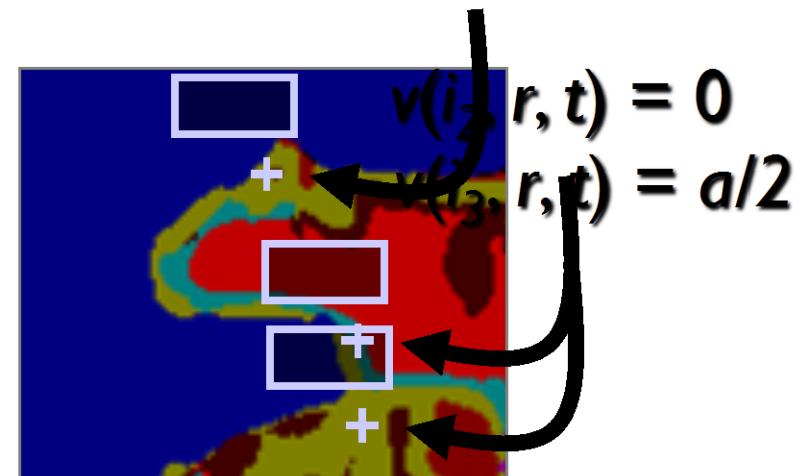


- ☐ Feature responses $v(i, r, t)$

- ☐ Large bounding boxes enable
long range interactions

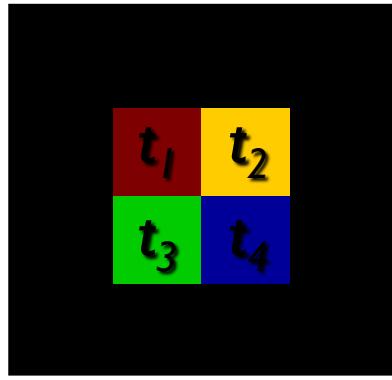
- ☐ Integral images

$$v(i_1, r, t) = a$$

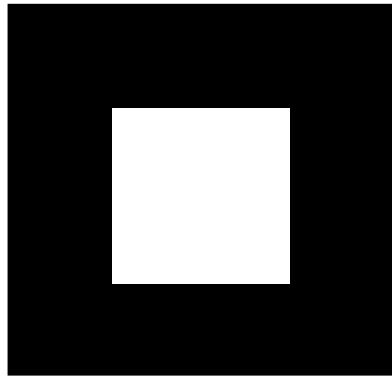


appearance context

Shape as Texton Layout

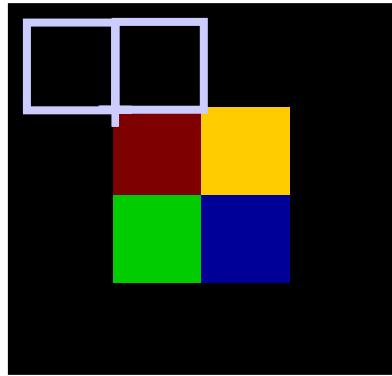


texton map

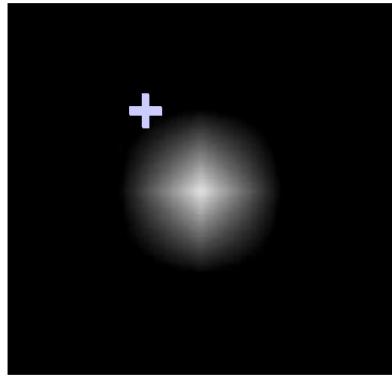


ground truth

$$(r_1, t_1) = \left(\begin{array}{c} \text{white square} \\ \text{black background} \end{array}, \text{red square} \right)$$
$$(r_2, t_2) = \left(\begin{array}{c} \text{white square} \\ \text{black background} \end{array}, \text{yellow square} \right)$$



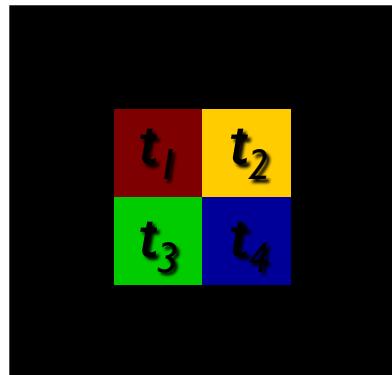
texton map



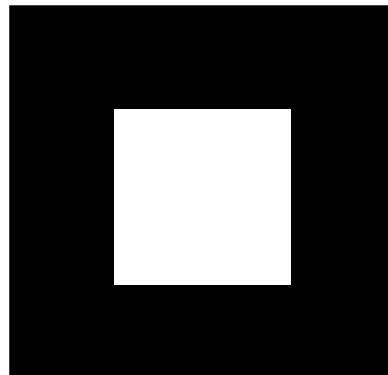
feature response image

$$v(i, r_2, t_2)$$

Shape as Texton Layout



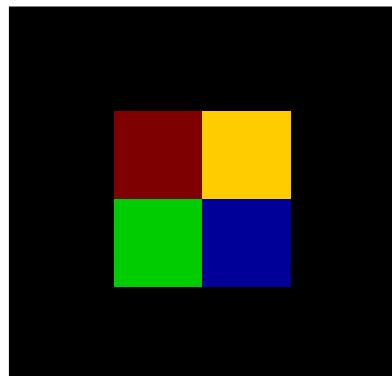
texton map



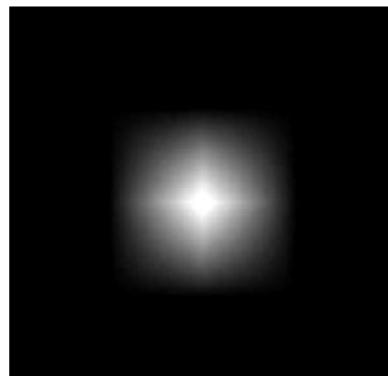
ground truth

$$(r_1, t_1) = \left(\begin{array}{c} \text{white square} \\ \text{black background} \end{array}, \text{red square} \right)$$

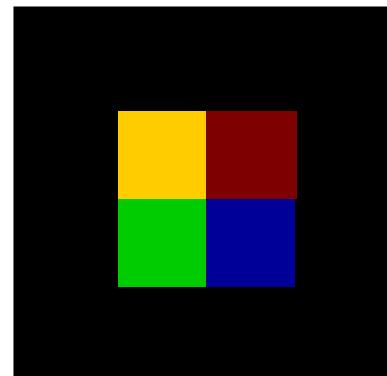
$$(r_2, t_2) = \left(\begin{array}{c} \text{white square} \\ \text{black background} \end{array}, \text{yellow square} \right)$$



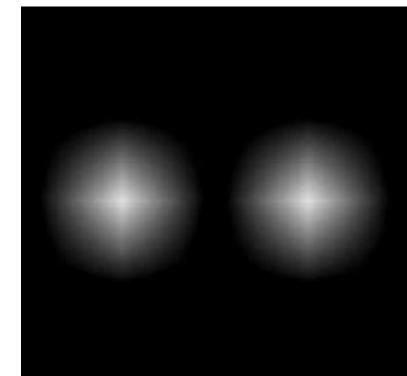
texton map



summed response images



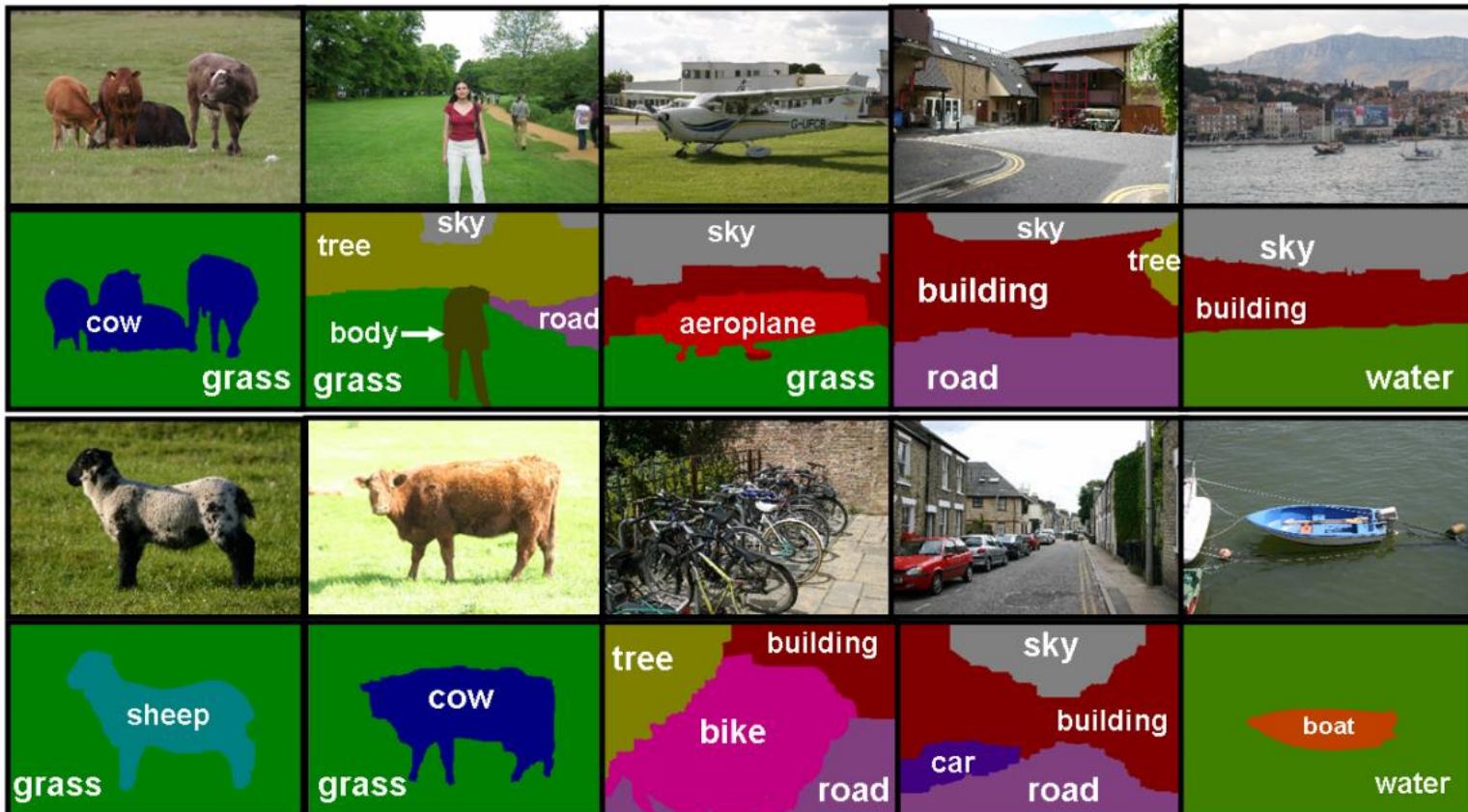
texton map



summed response images

$$v(i, r_1, t_1) + v(i, r_2, t_2)$$

$$v(i, r_1, t_1) + v(i, r_2, t_2)$$



<i>Object classes</i>	Building	Grass	Tree	Cow	Sheep	Sky	Aeroplane	Water	Face	Car
Bike	Flower	Sign	Bird	Book	Chair	Road	Cat	Dog	Body	Boat

❑ Failures



Benchmark Challenges

VOC: 20 classes



COCO: 200 classes

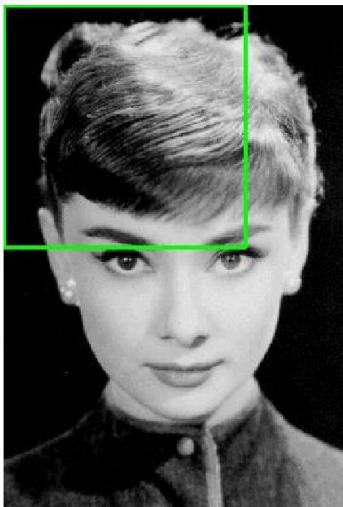


Challenge 3: Dense vs Sparse

Learning Pictorial Structure (DPM)

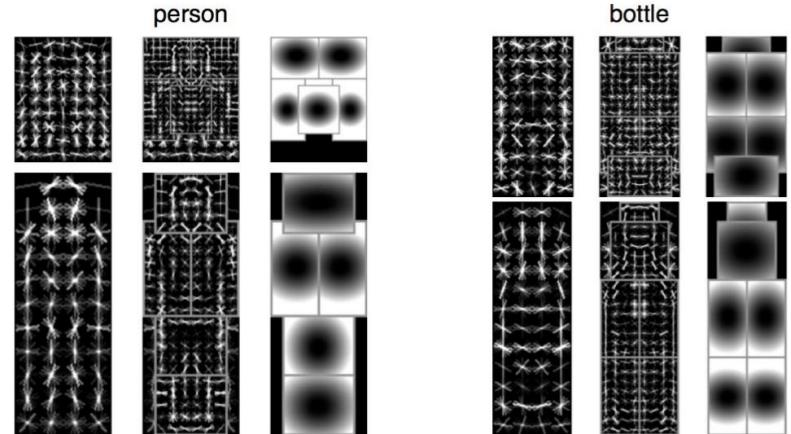
Sliding windows.

- Score every subwindow.

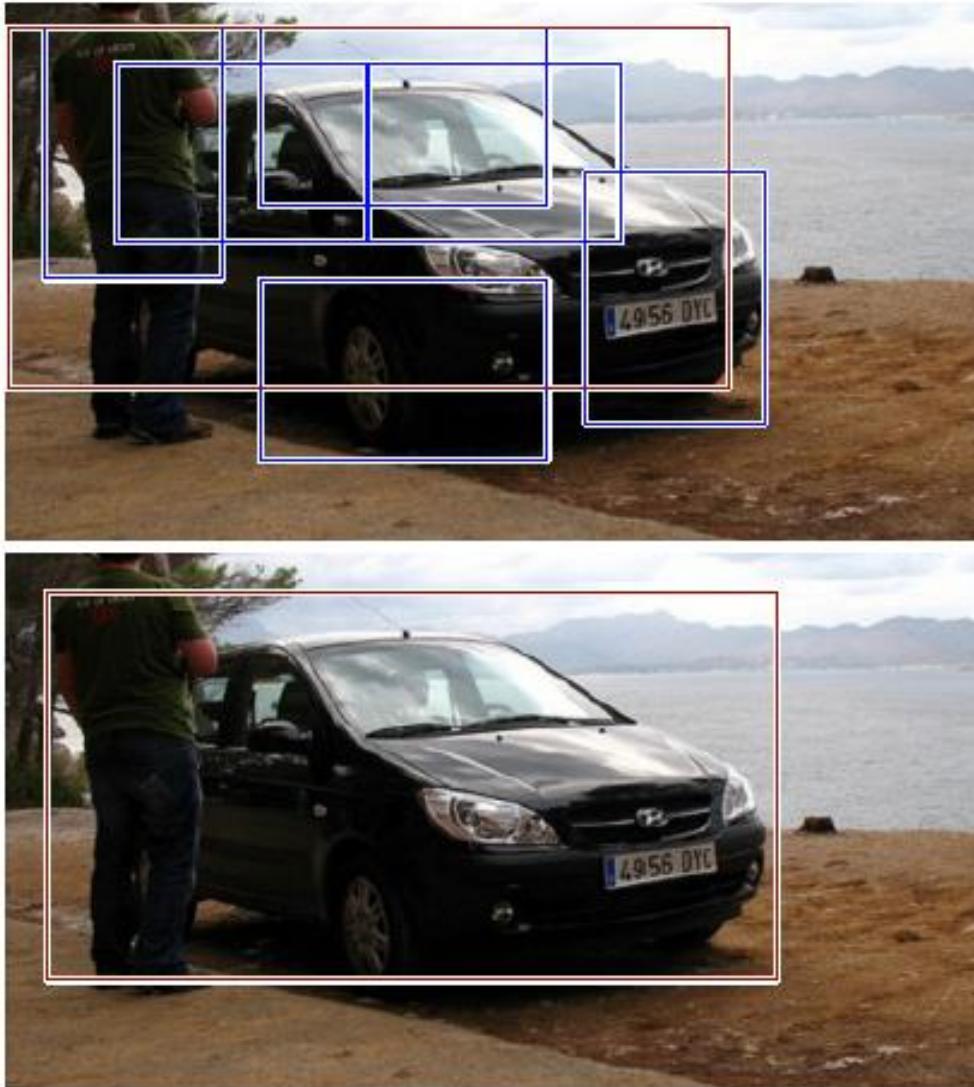


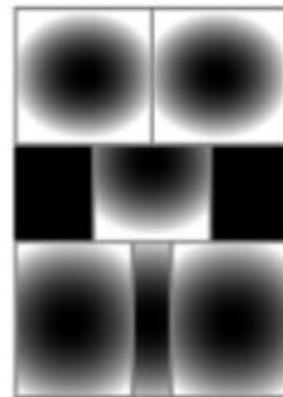
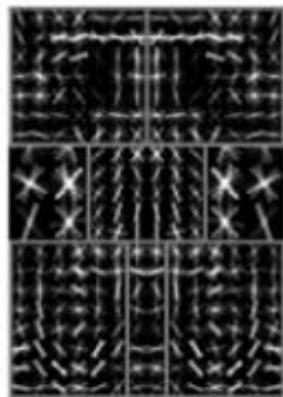
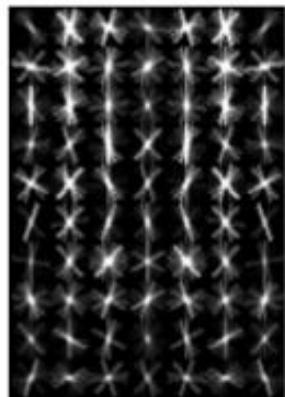
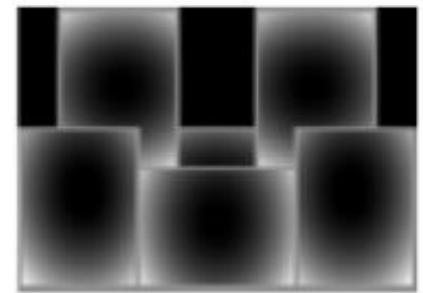
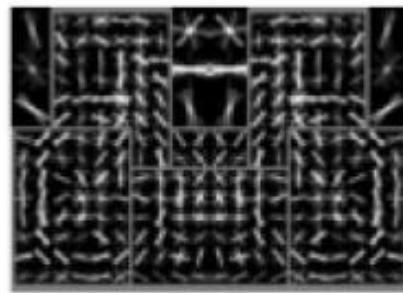
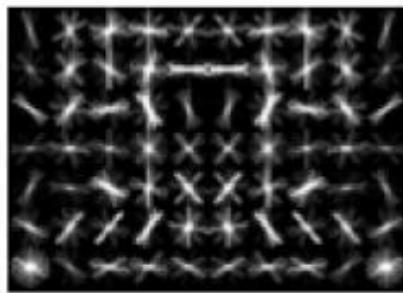
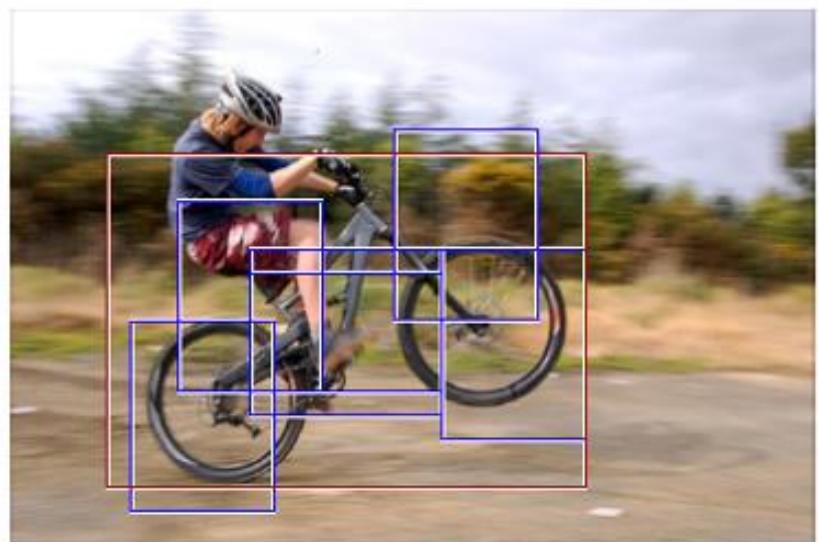
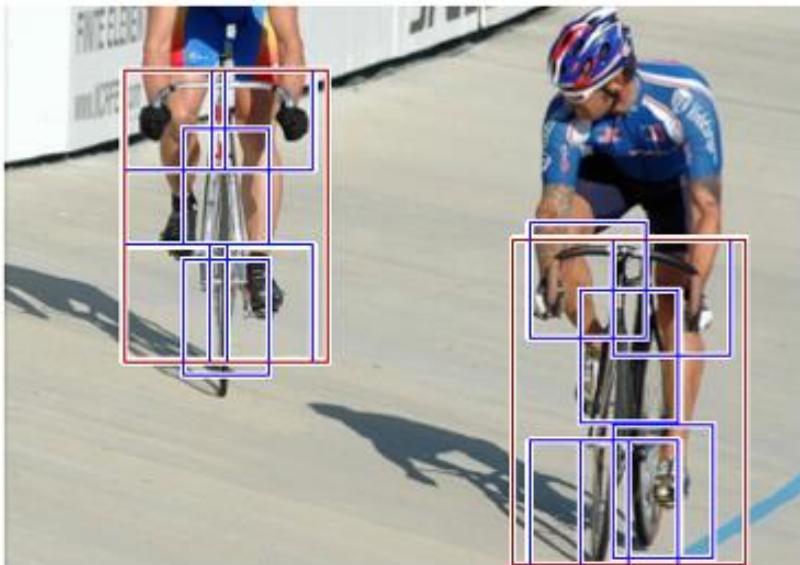
Deformable part models (DPM)

- Uses HOG features
- Very fast

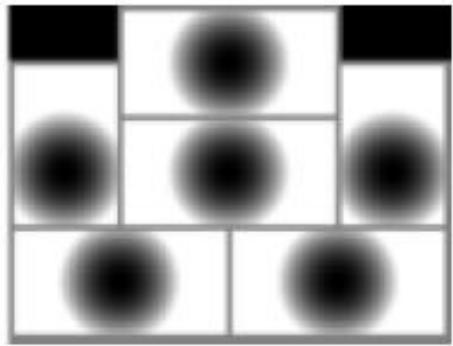
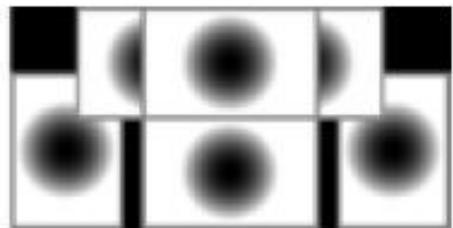
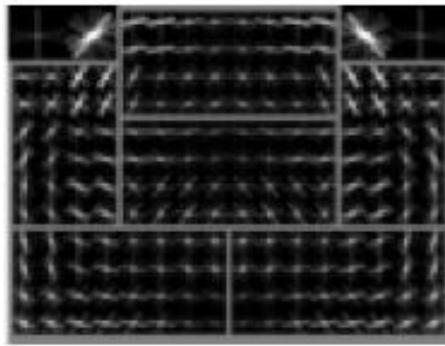
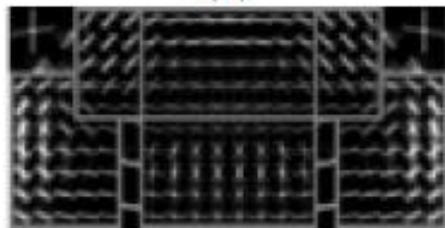
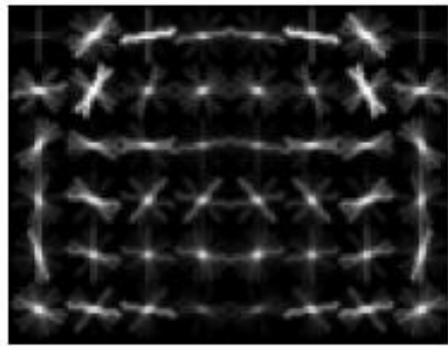
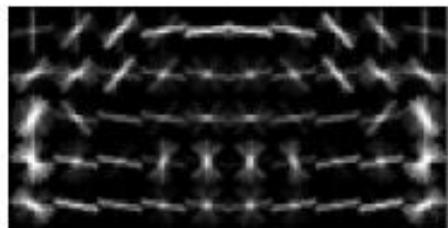


- 1) fine level with deformable parts
- 2) coarse level with a fixed template model

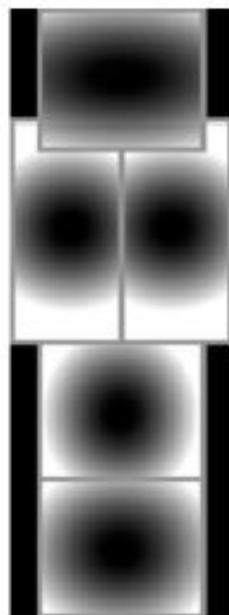
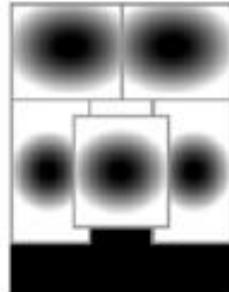
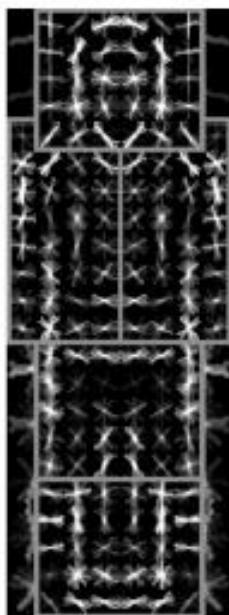
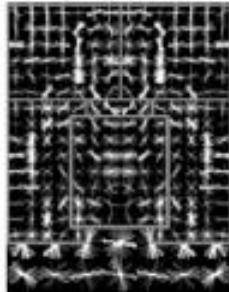
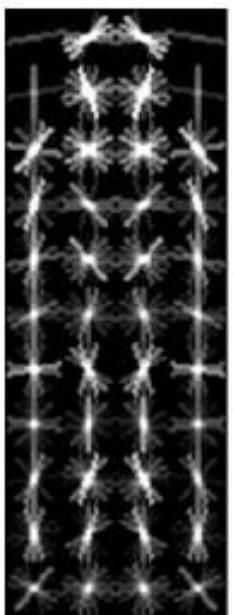
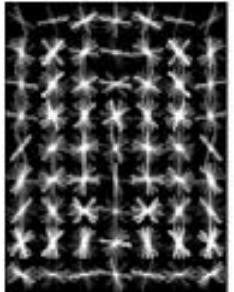




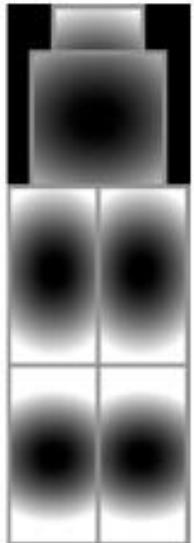
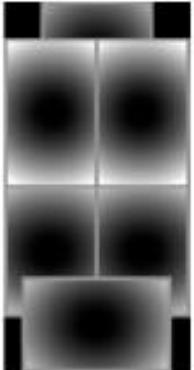
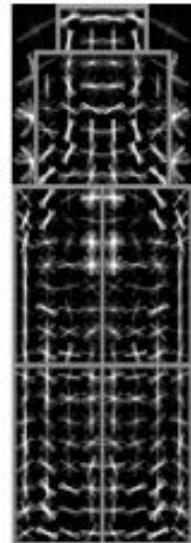
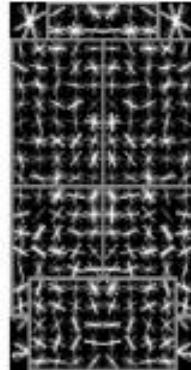
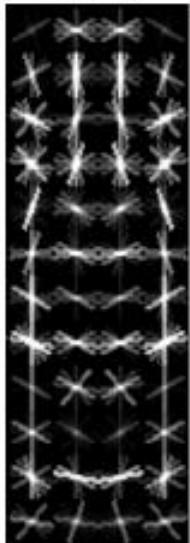
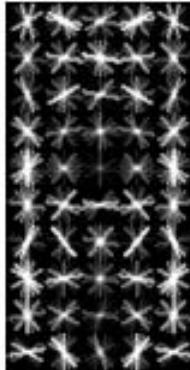
(c)

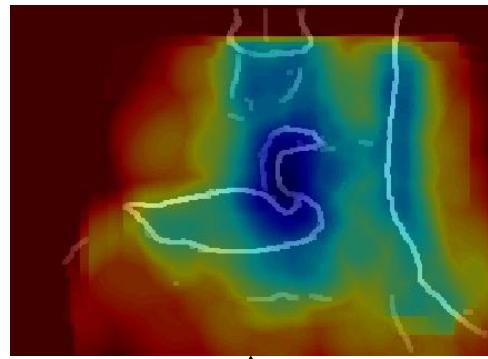
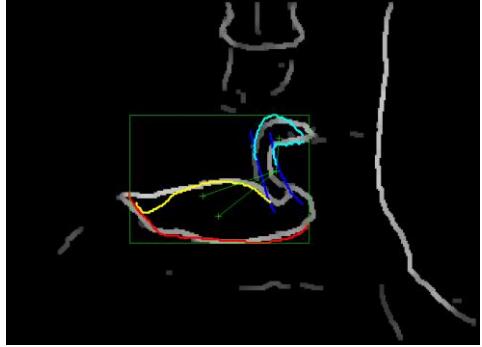


person

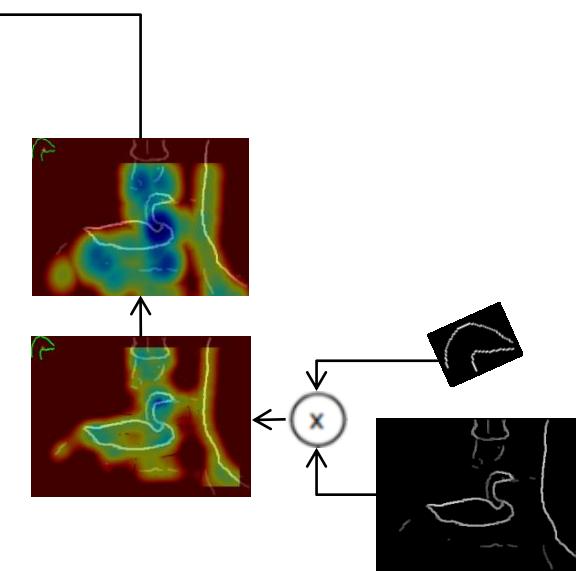
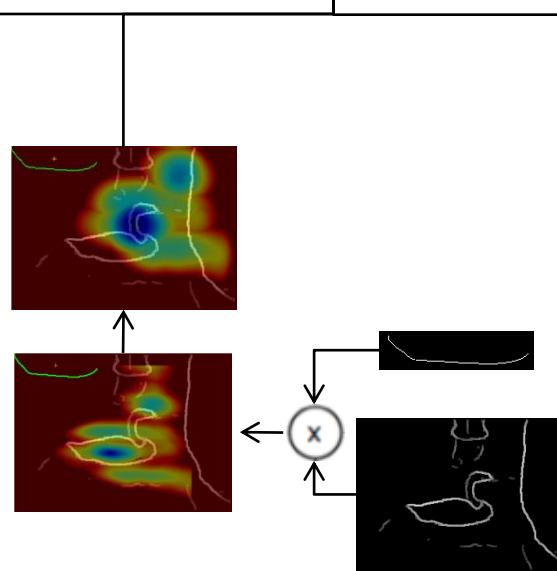
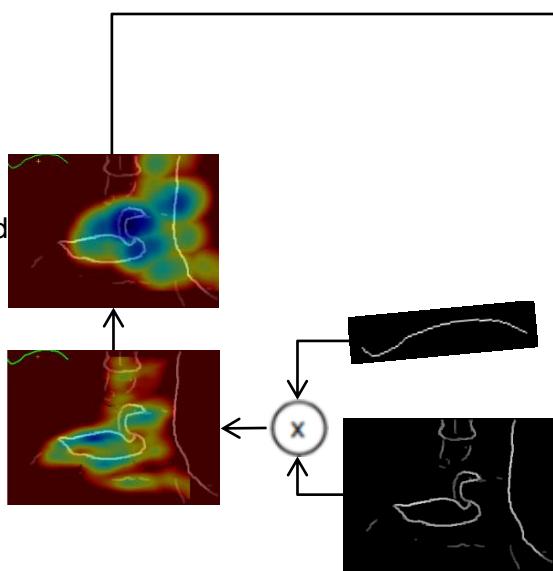
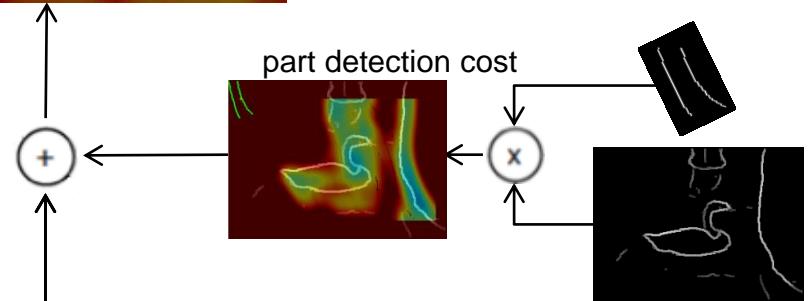
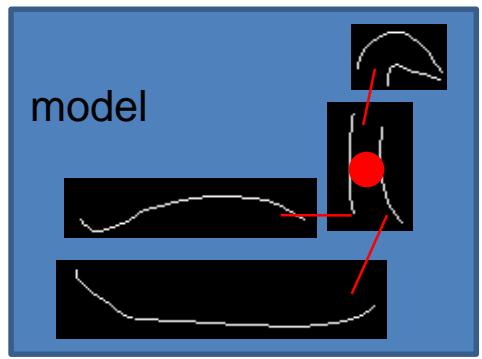


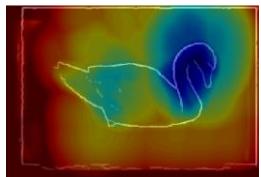
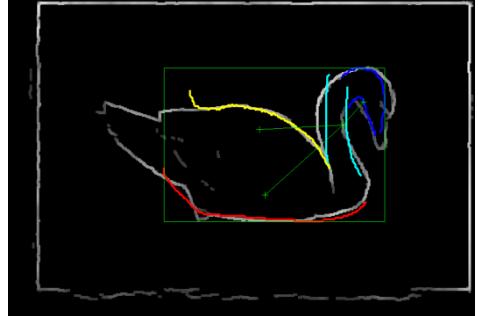
bottle



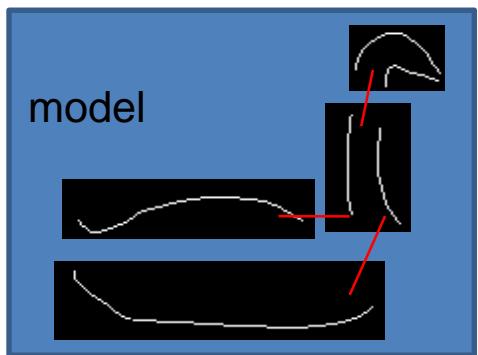
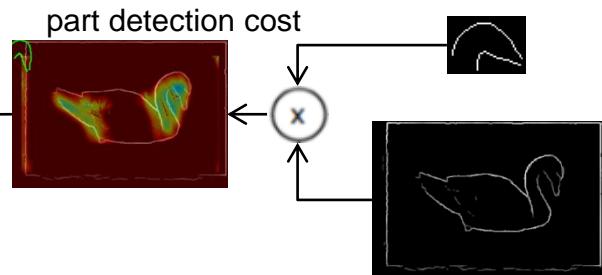


combined cost of root (neck) locations

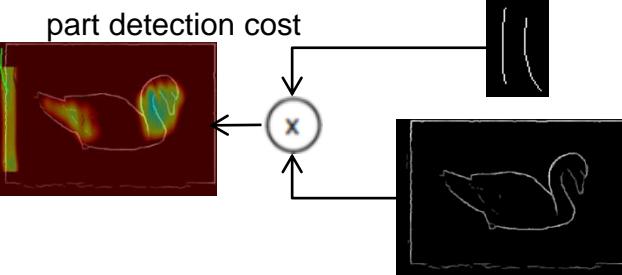
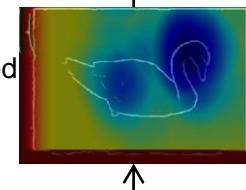




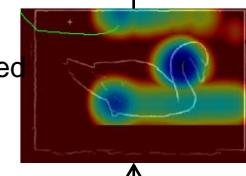
combined cost of root (head) locations



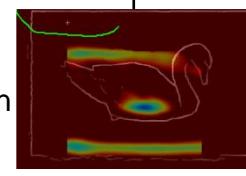
transformed cost

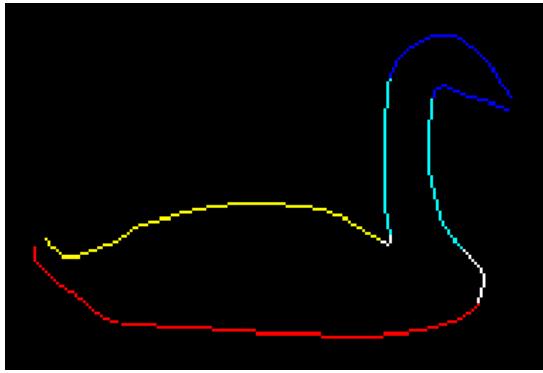


transformed cost

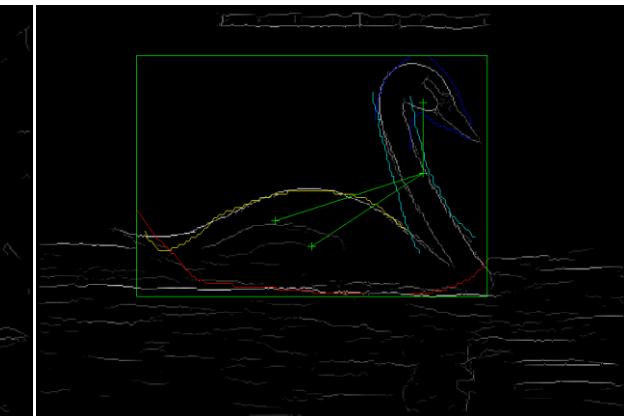
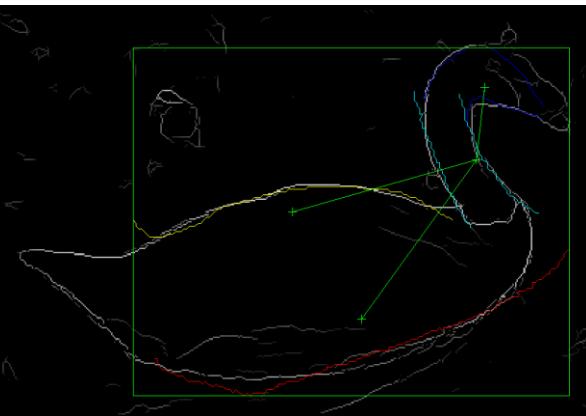
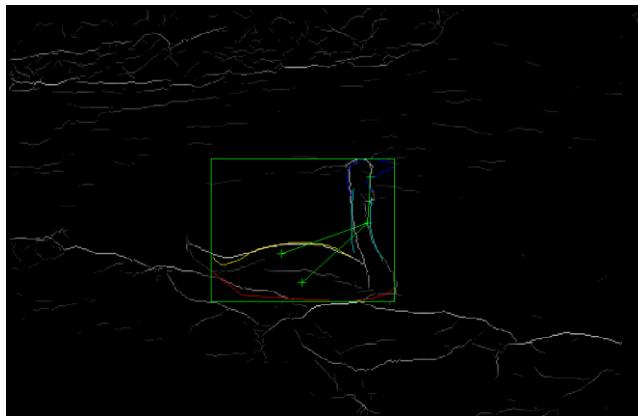


part detection cost

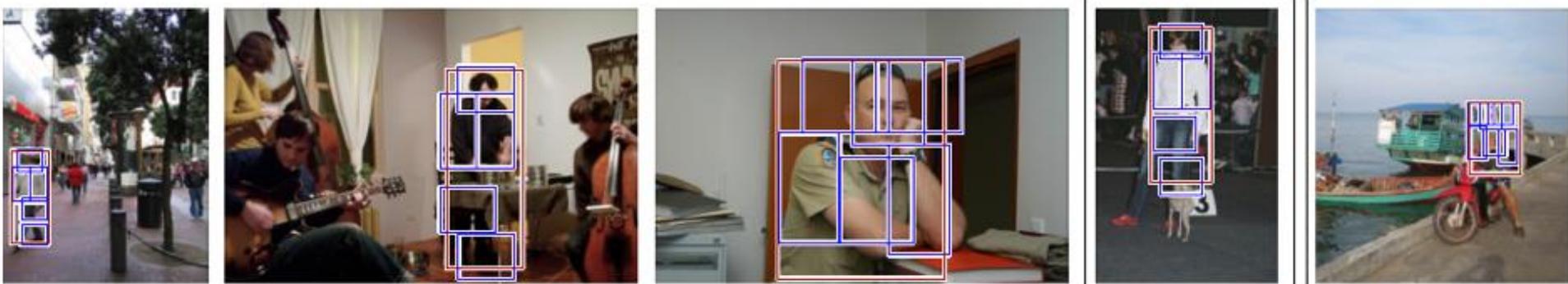




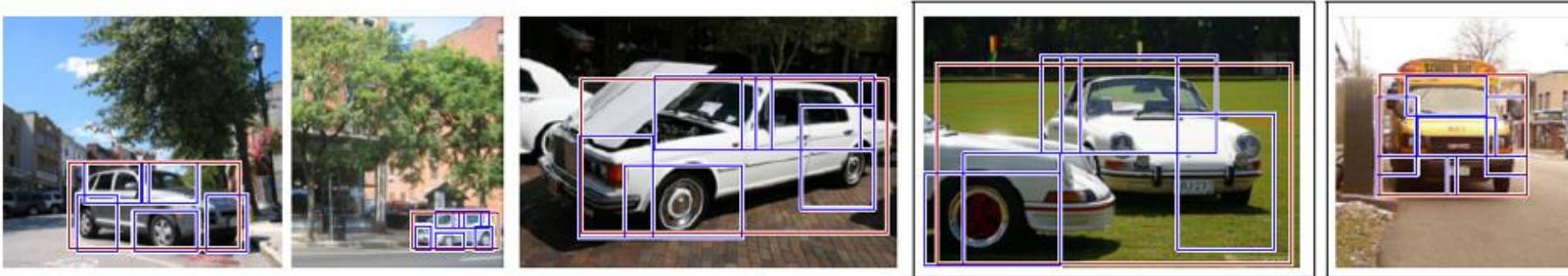
mod



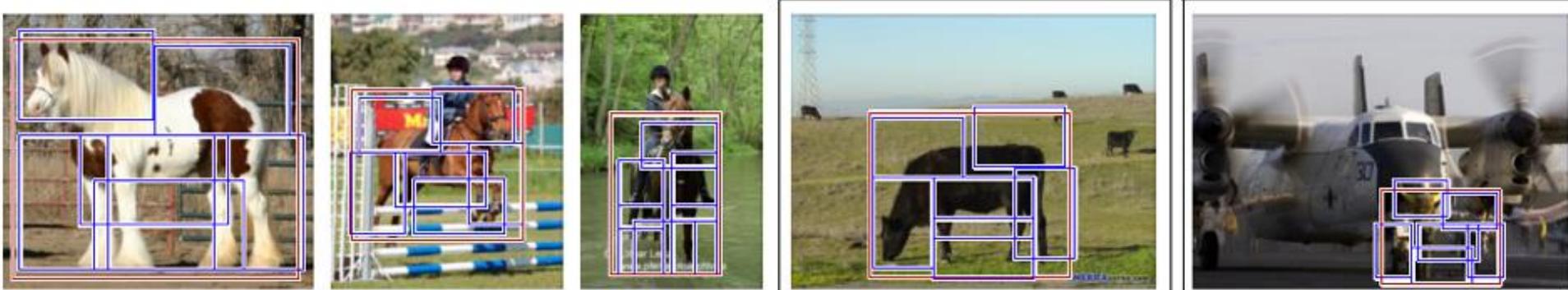
person



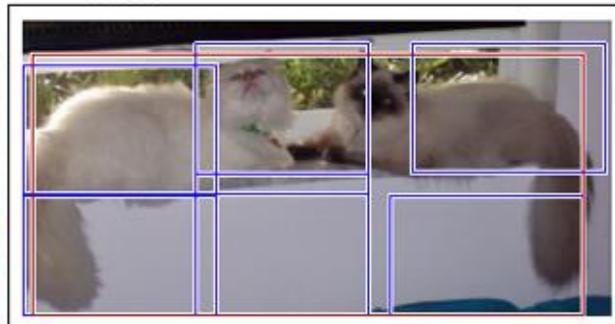
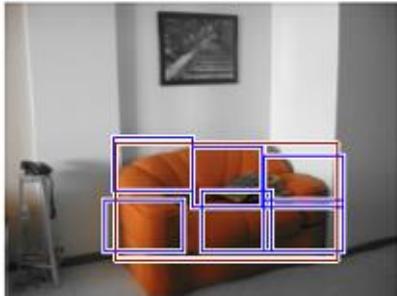
car



horse



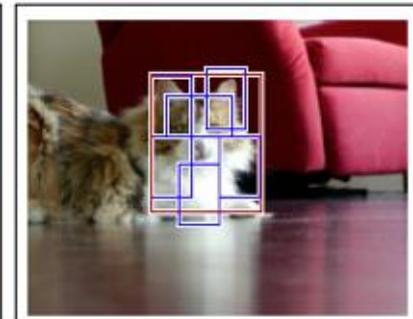
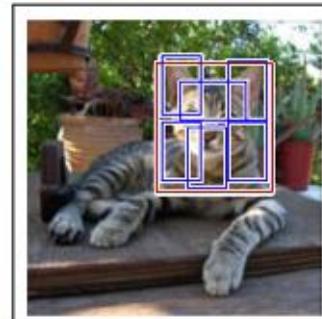
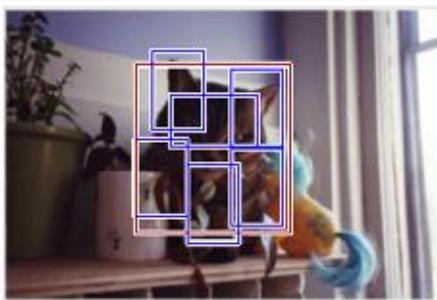
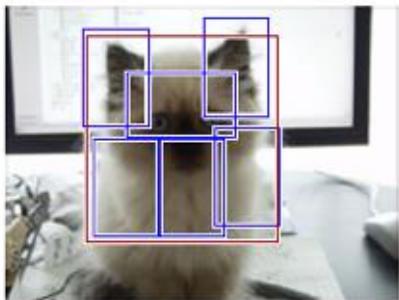
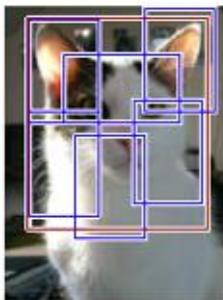
sofa



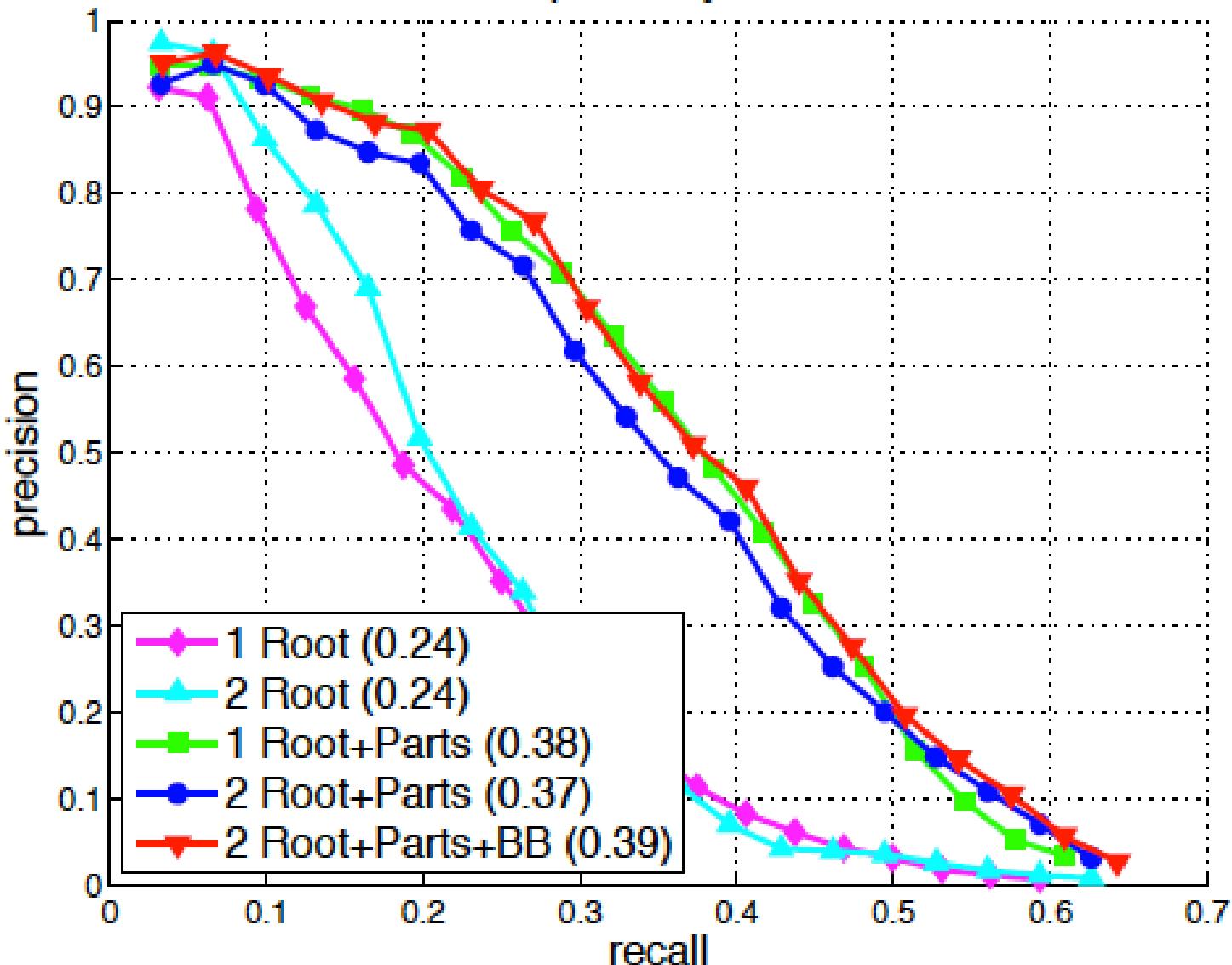
bottle



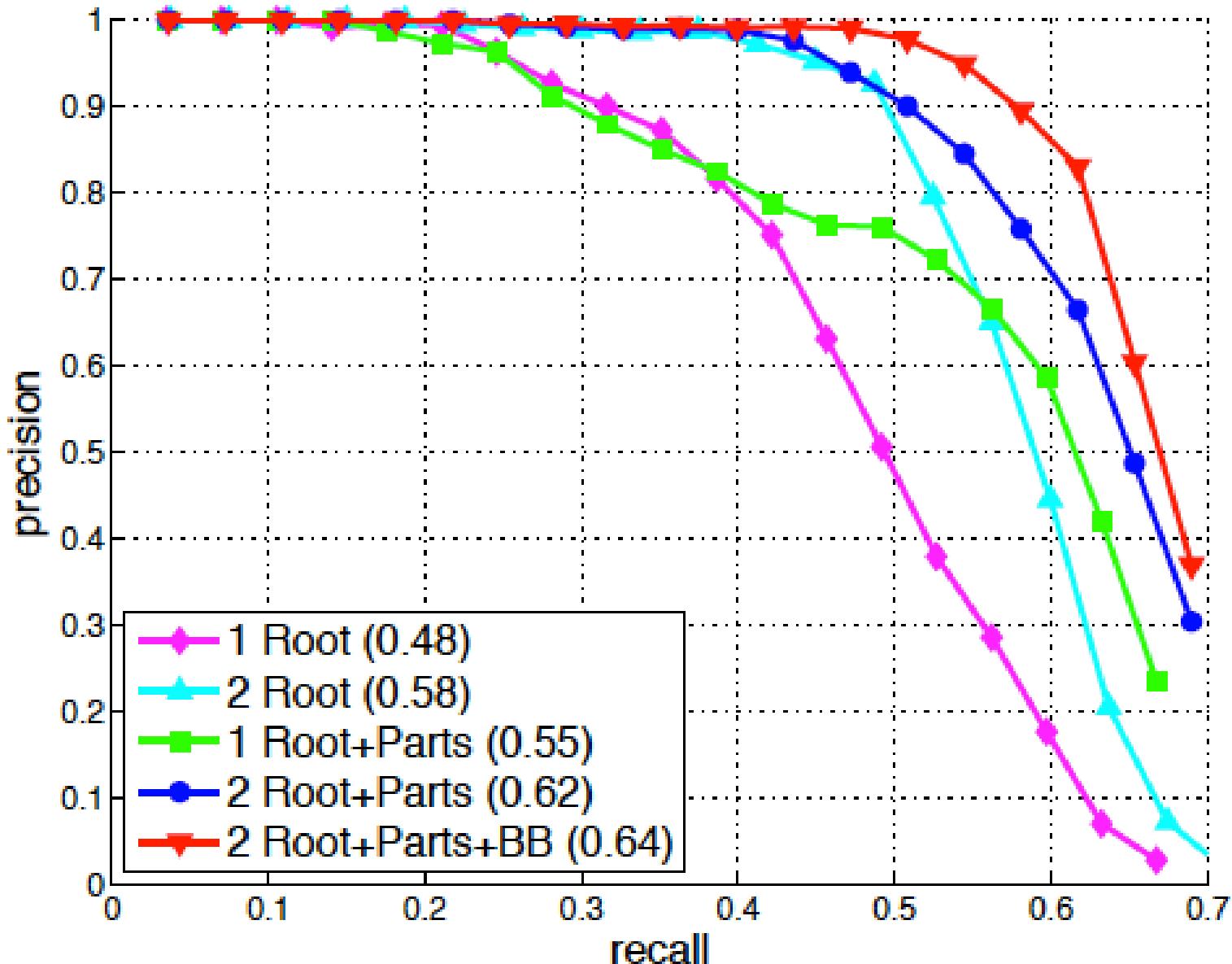
cat



class: person, year 2006



class: car, year 2006



Challenge 3: Dense vs Sparse

Selective Search for Object Recognition

J.R.R. Uijlings^{*1,2}, K.E.A. van de Sande^{†2}, T. Gevers², and A.W.M. Smeulders²

¹University of Trento, Italy

²University of Amsterdam, the Netherlands

Technical Report 2012, submitted to IJCV

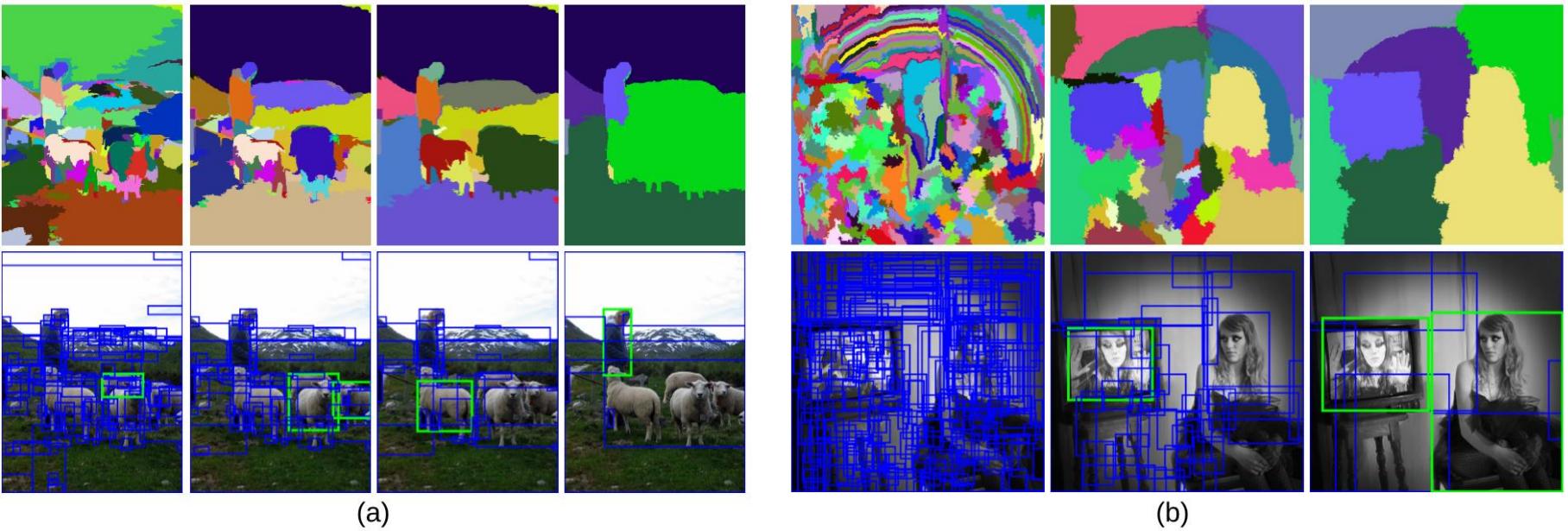


Figure 2: Two examples of our selective search showing the necessity of different scales. On the left we find many objects at different scales. On the right we necessarily find the objects at different scales as the girl is contained by the tv.

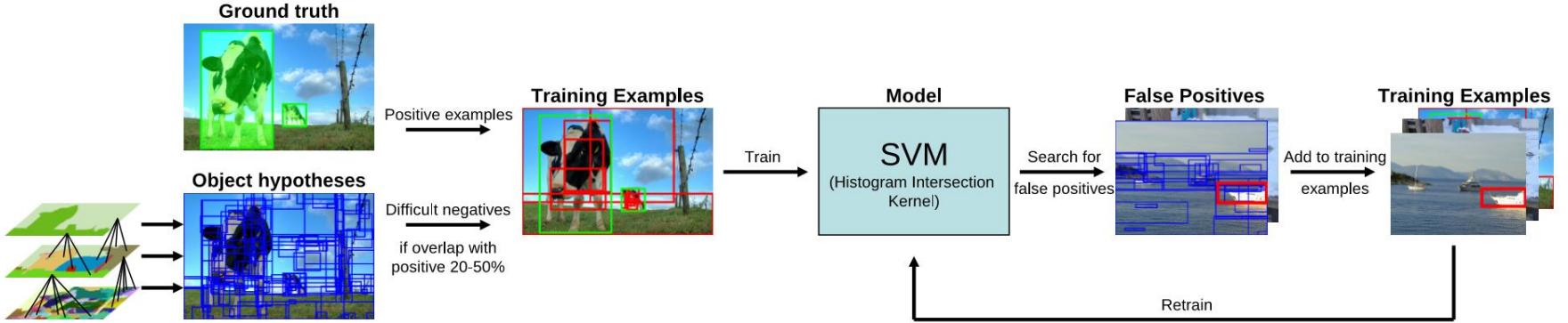


Figure 3: The training procedure of our object recognition pipeline. As positive learning examples we use the ground truth. As negatives we use examples that have a 20-50% overlap with the positive examples. We iteratively add hard negatives using a retraining phase.

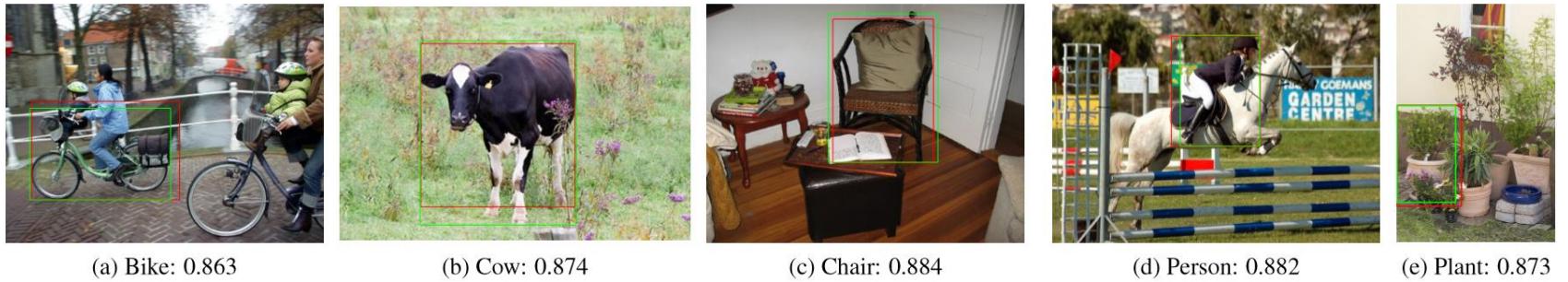


Figure 5: Examples of locations for objects whose Best Overlap score is around our Mean Average Best Overlap of 0.879. The green boxes are the ground truth. The red boxes are created using the “Quality” selective search.

Lessons and Hard Cases:

Hard Negative Mining

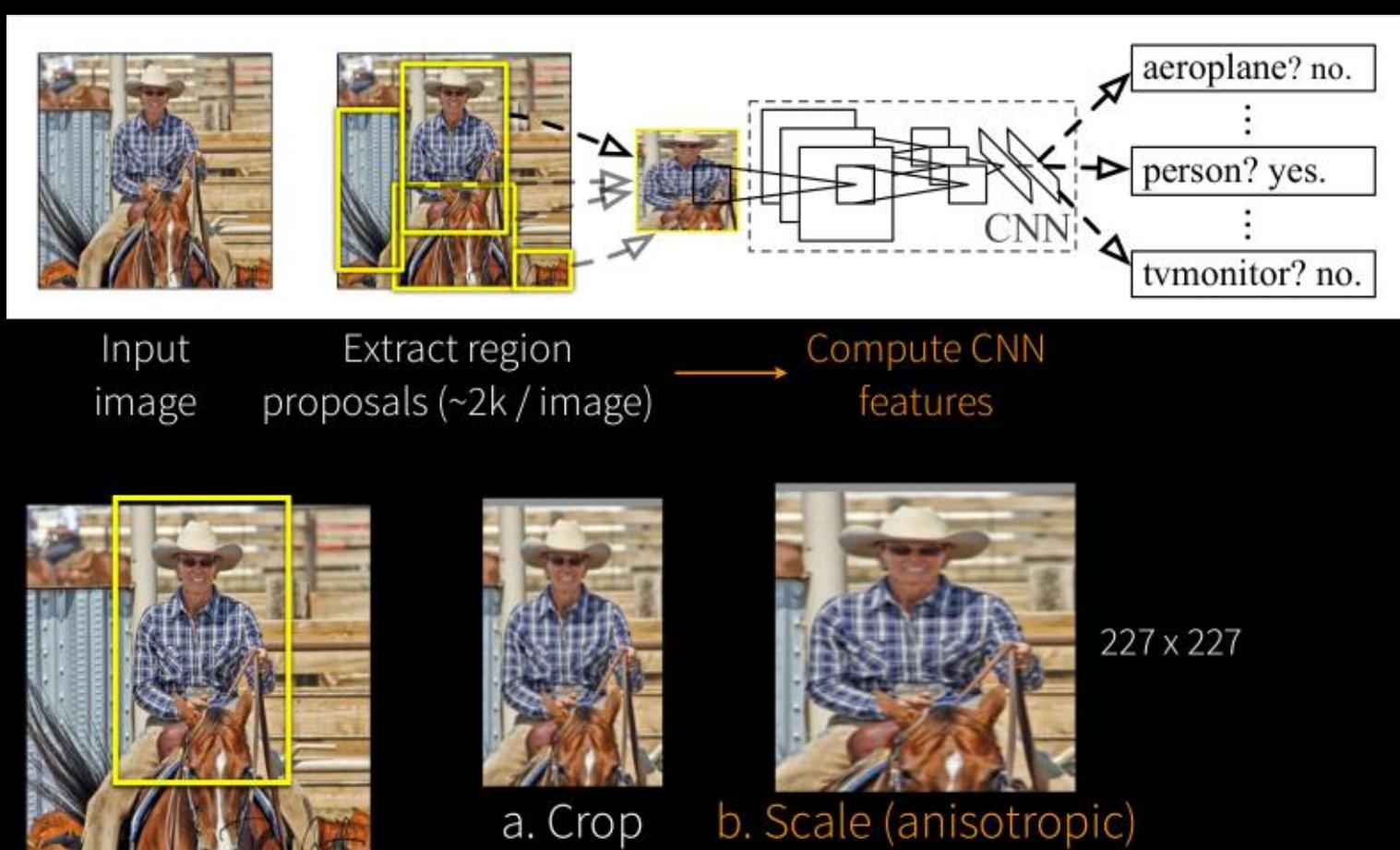
Imbalance between positive and negative examples.

Use negative examples with higher confidence score.

Non Maximum Suppression

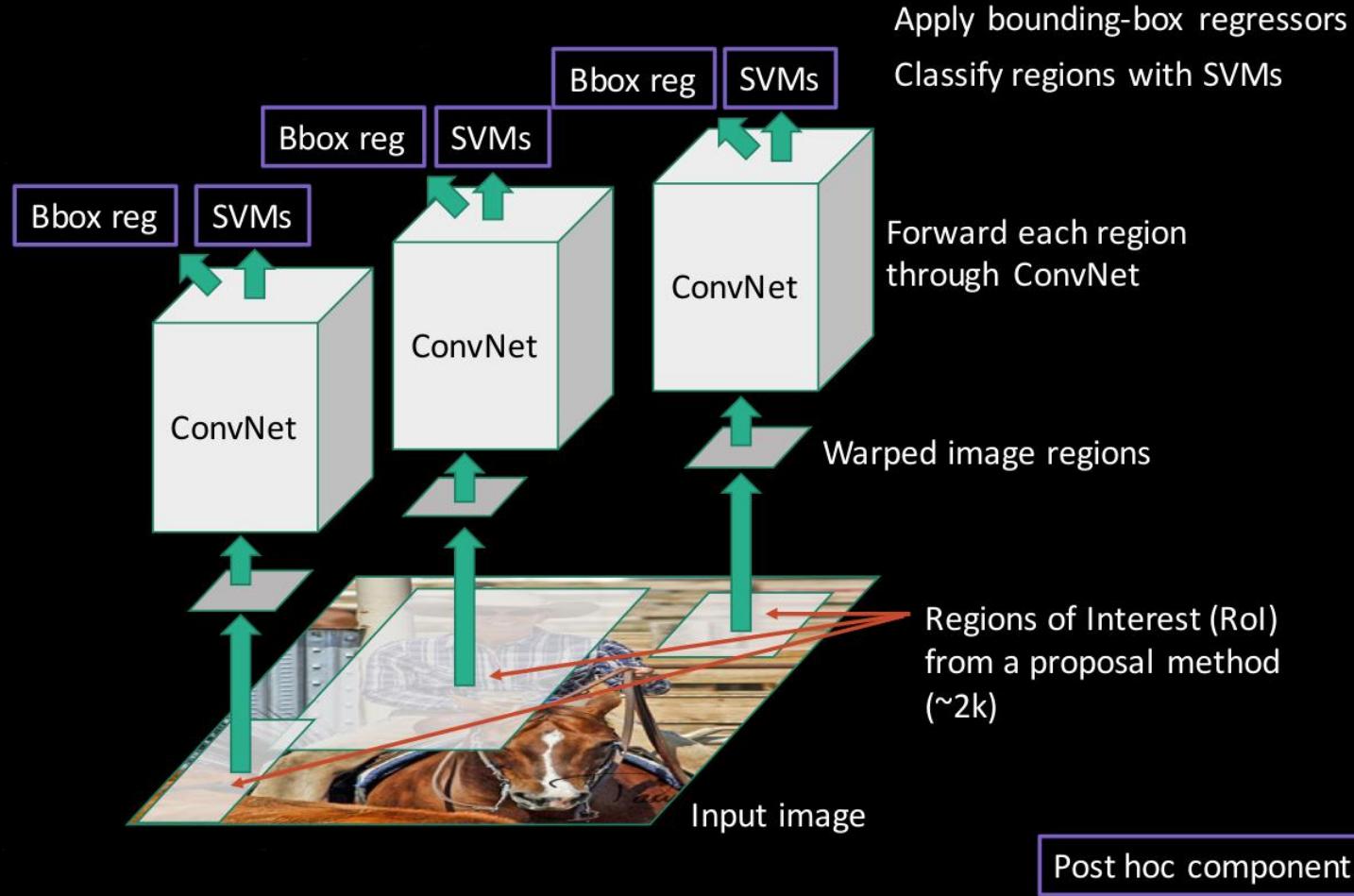
If output boxes overlap, only consider the most confident.

Region based classification: Proposal (where) and Checking (what)



R-CNN: Regions with CNN Features (Girshick et al., 2014)

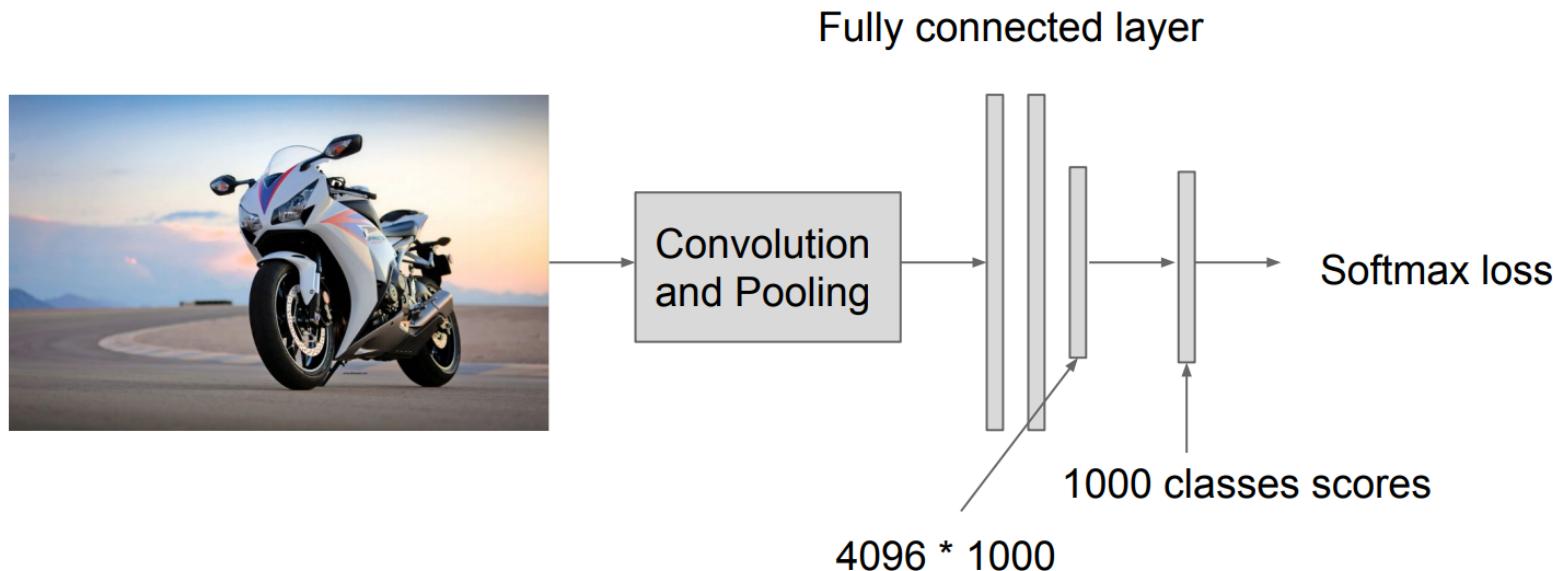
Object Detection



R-CNN: Regions with CNN Features (Girshick et al., 2014)

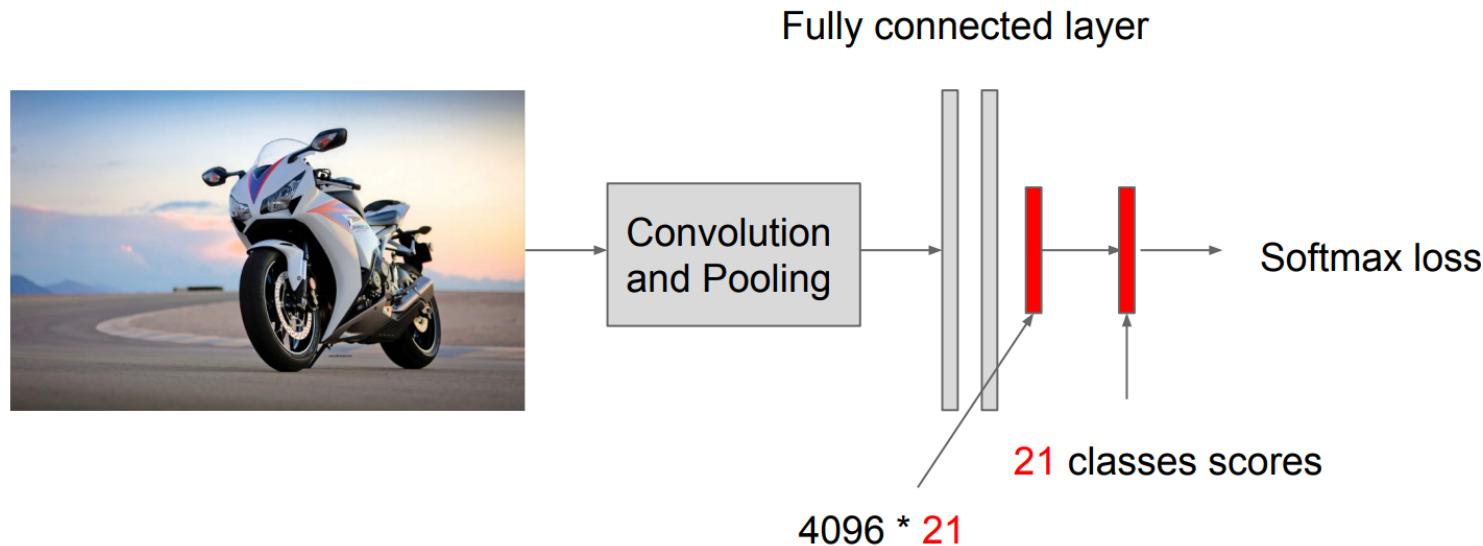
Training RCNN

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

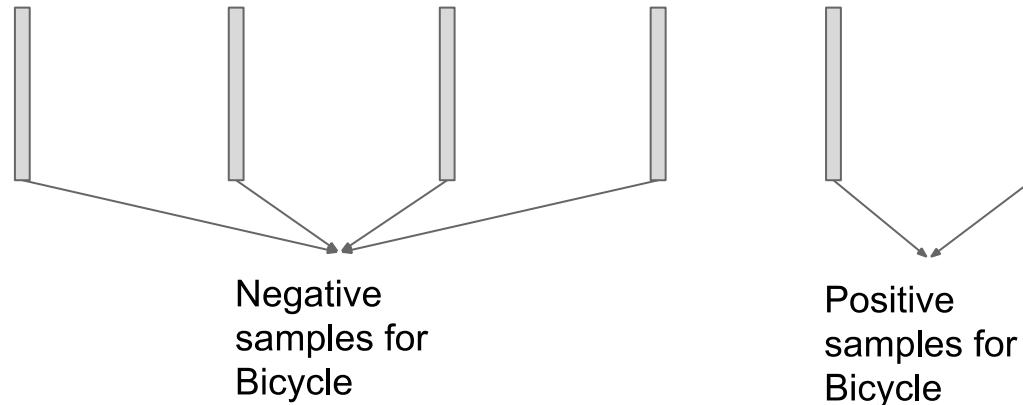
Training RCNN

Step4: train SVM for each class

Crop /
Warp
image



Features
from last
step



Classify regions by SVM

- linear SVM per class
(With the softmax classifier from fine-tuning mAP decreases from 54% to 51%)
- greedy NMS(non-maximum suppression) per class : rejects a region if it has an intersection-overunion (IoU) overlap with a higher scoring selected region larger than a learned threshold

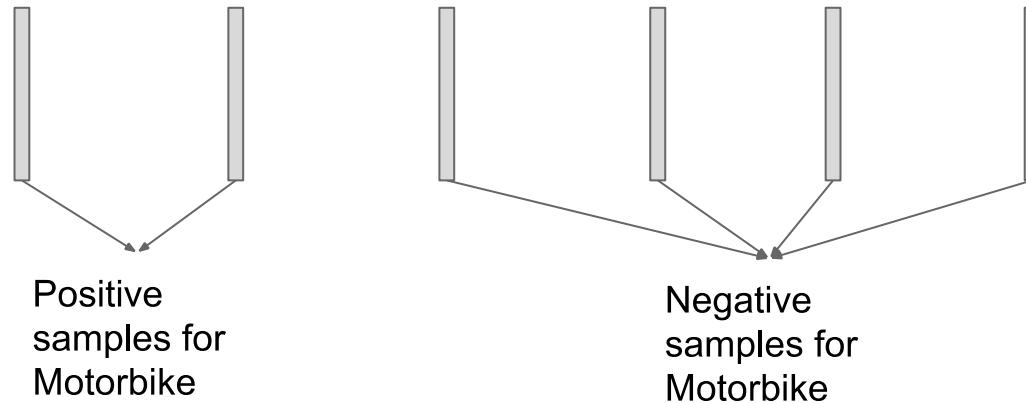
Training RCNN

Step4: train SVM for each class

Crop /
Warp
image



Features
from last
step



Object proposal refinement

- Linear bounding-box regression on CNN features (pool_5 feature: mAP ~4% up)

Training R-CNN

- Bounding-box labeled detection data is scarce
- Use supervised pre-training on a data-rich auxiliary task and transfer to detection

Supervised pre-training

Pre-train CNN on ILSVRC2012(1.2 million 1000-way image classification) using image-level annotations only

Domain-specific fine-tuing

Adapt to new task(detection) and new domain(warped proposal)

- random initialize $(N+1)$ -way classification layer (N classes + background)
- Positives: 0.5 IoU overlap with a ground-truth box. Negative: o.w.
- SGD: learning rate: 0.001 (1/10 of original) mini-batch: 32 pos & 96 neg

Train binary SVM

- IoU overlap threshold: grid search over {0, 0.1, ... 0.5}
 - IoU = 0.5 : mAP ~5% down
 - IoU = 0.0 : mAP ~4% down

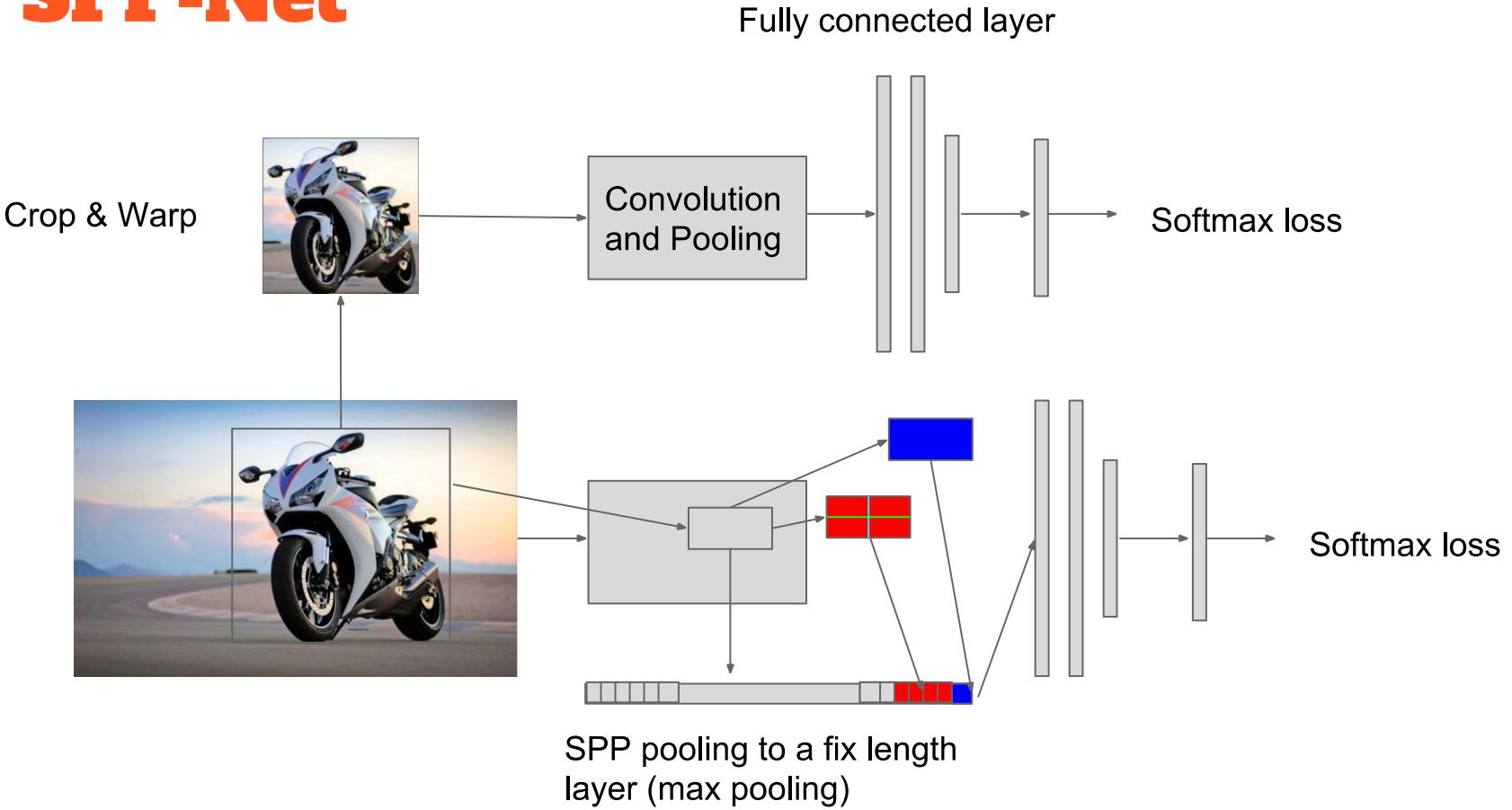
SPP-Net



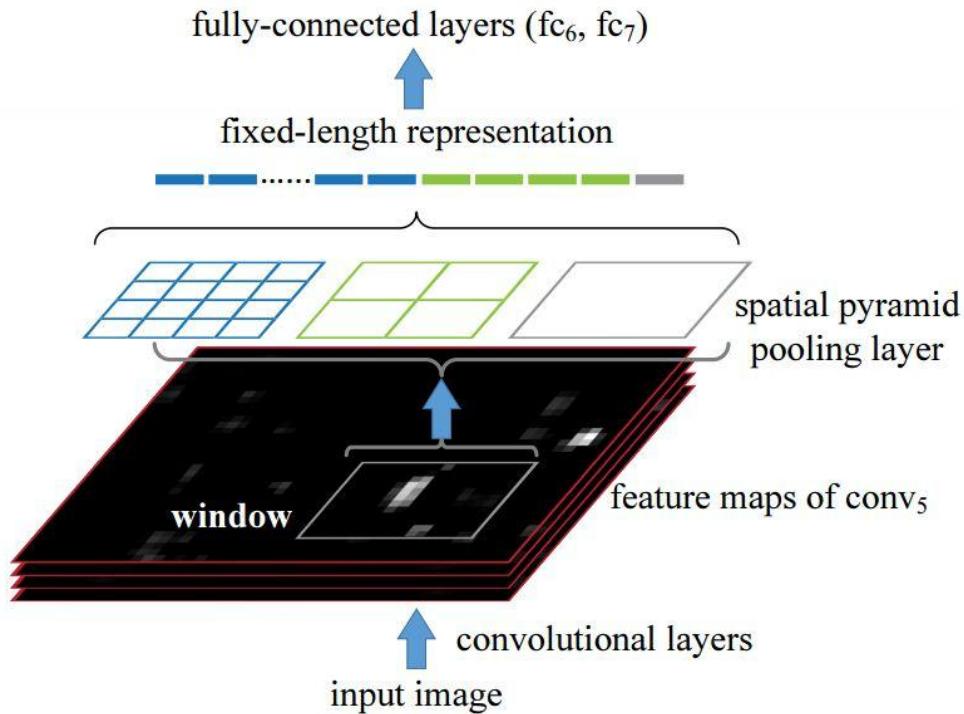
Figure 1: Top: cropping or warping to fit a fixed size. Middle: a conventional CNN. Bottom: our spatial pyramid pooling network structure.

<https://arxiv.org/pdf/1406.4729v4.pdf>

SPP-Net

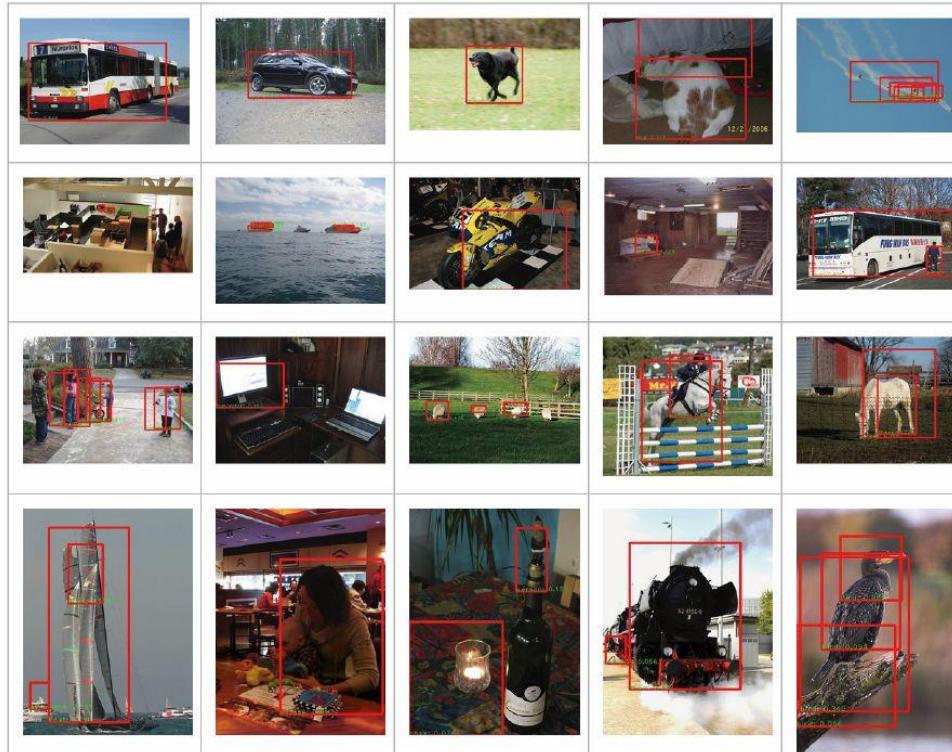


SPP-Net



<https://arxiv.org/pdf/1406.4729v4.pdf>

Multi-Box



Goal: Achieve a class-agnostic scalable object detection by predicting a set of bounding boxes.

Train a neural network to directly predict:

- The upper-left and lower-right coordinates of each bounding box
- The confidence score for the box containing an object

Fast RCNN

Share convolution layers for proposals from the same image

Faster and More accurate than RCNN

ROI Pooling

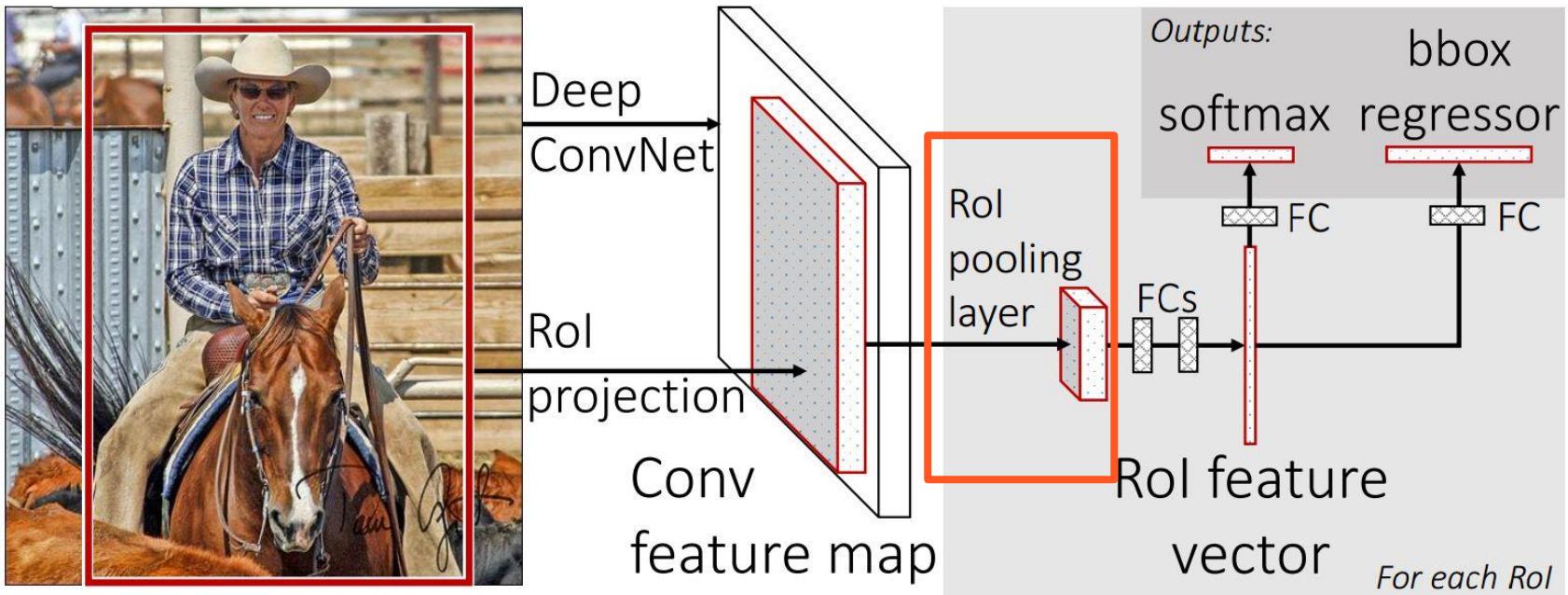
Fast R-CNN

Ross Girshick

Apr 2015

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

Fast RCNN

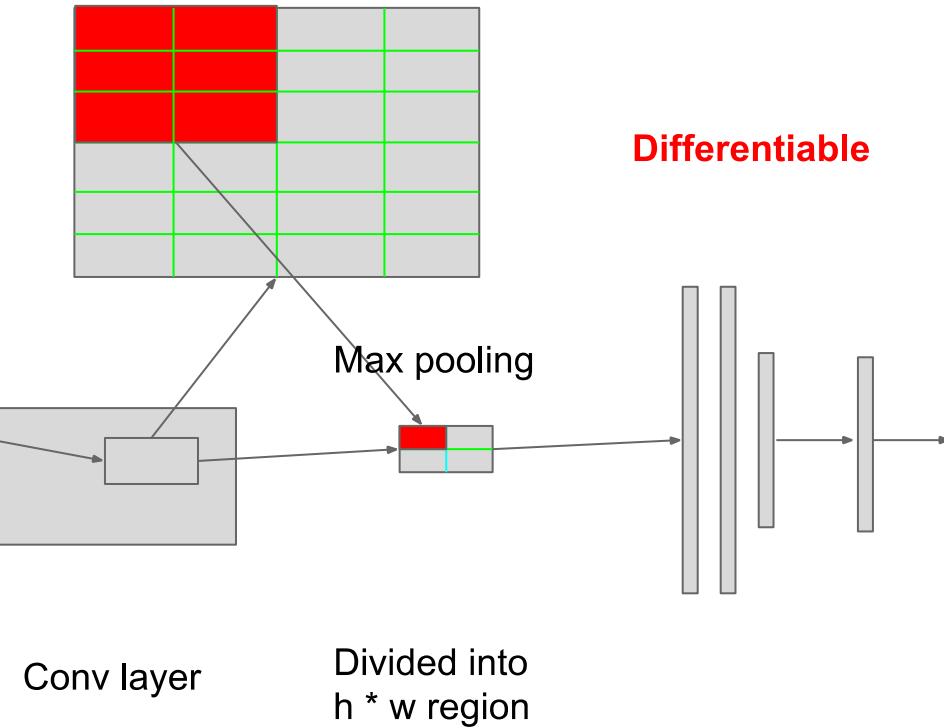


<https://arxiv.org/pdf/1504.08083v2.pdf>

RoI Pooling

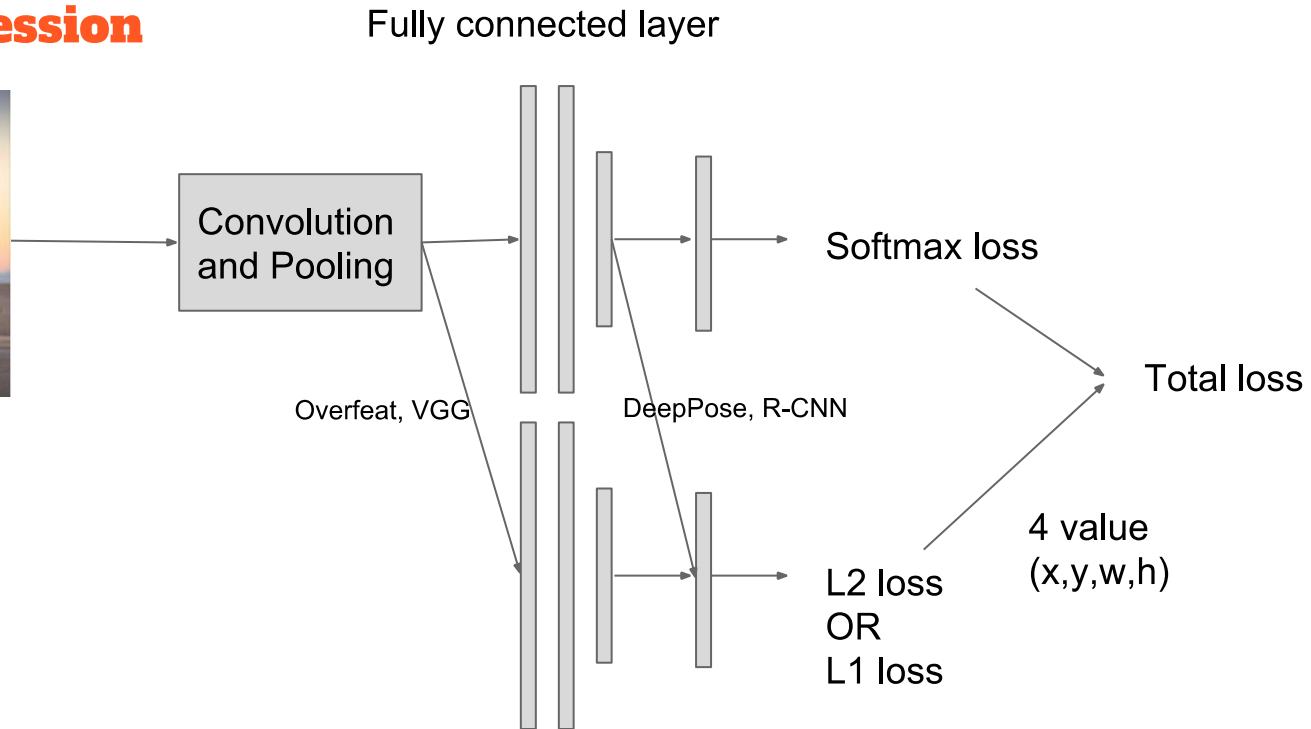


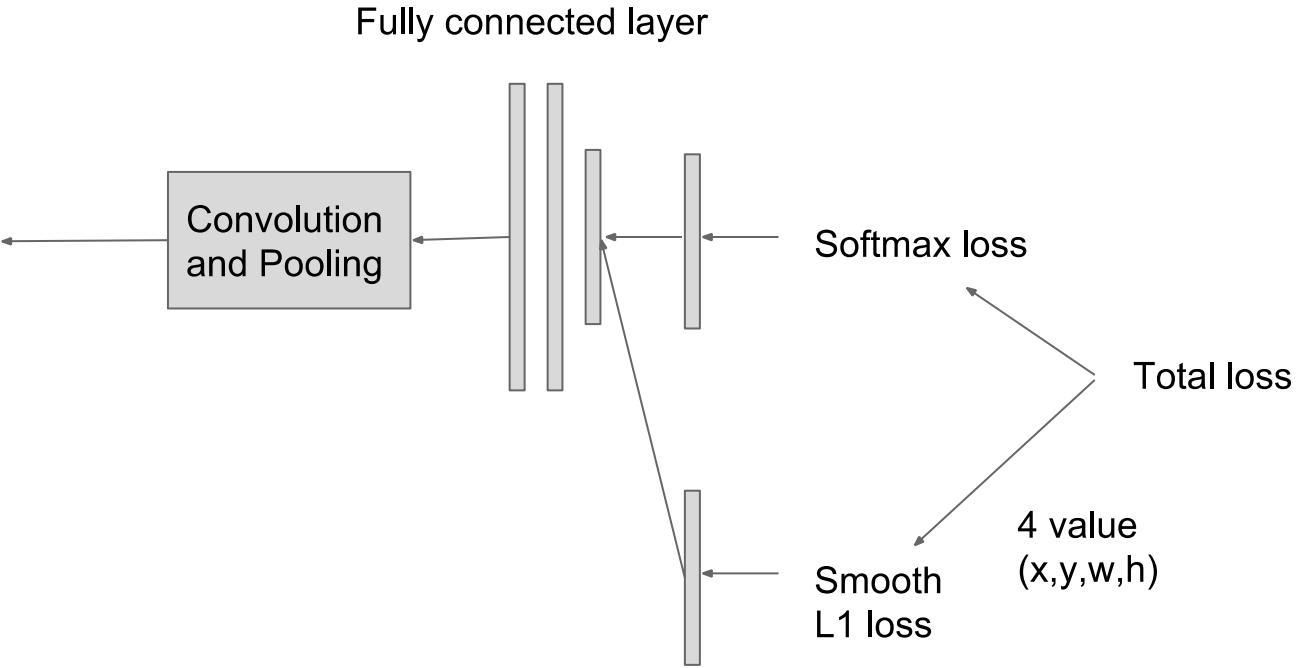
image



What is bbox-regressor?

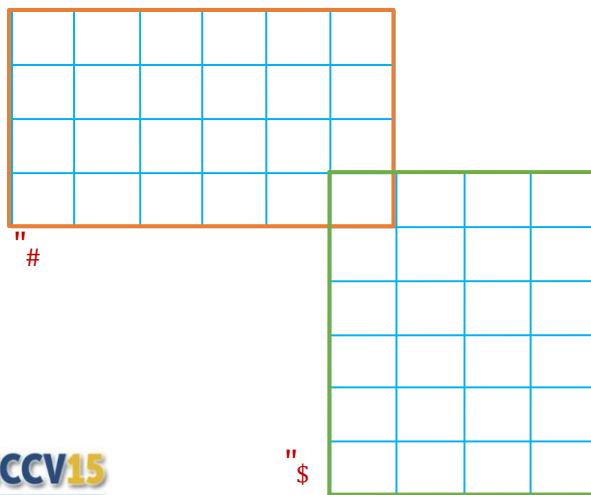
Bounding box regression





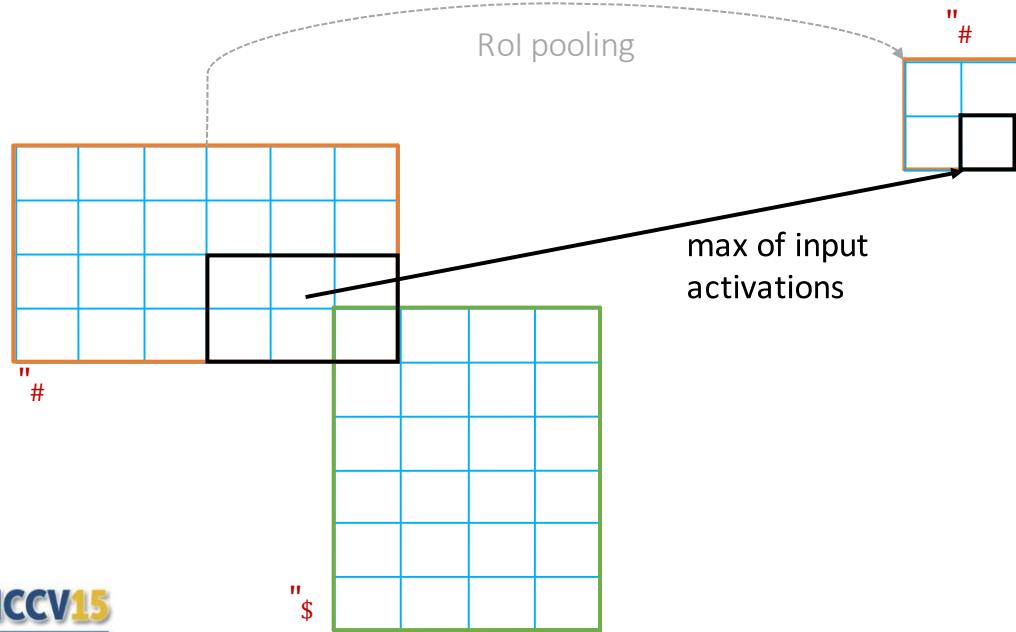
Obstacle '#1:'Differentiable'RoI pooling

RoI pooling'/'SPP'is'just'like'max'pooling,'**except'that'**pooling'regions'
overlap



Obstacle '#1: 'Differentiable' RoI pooling

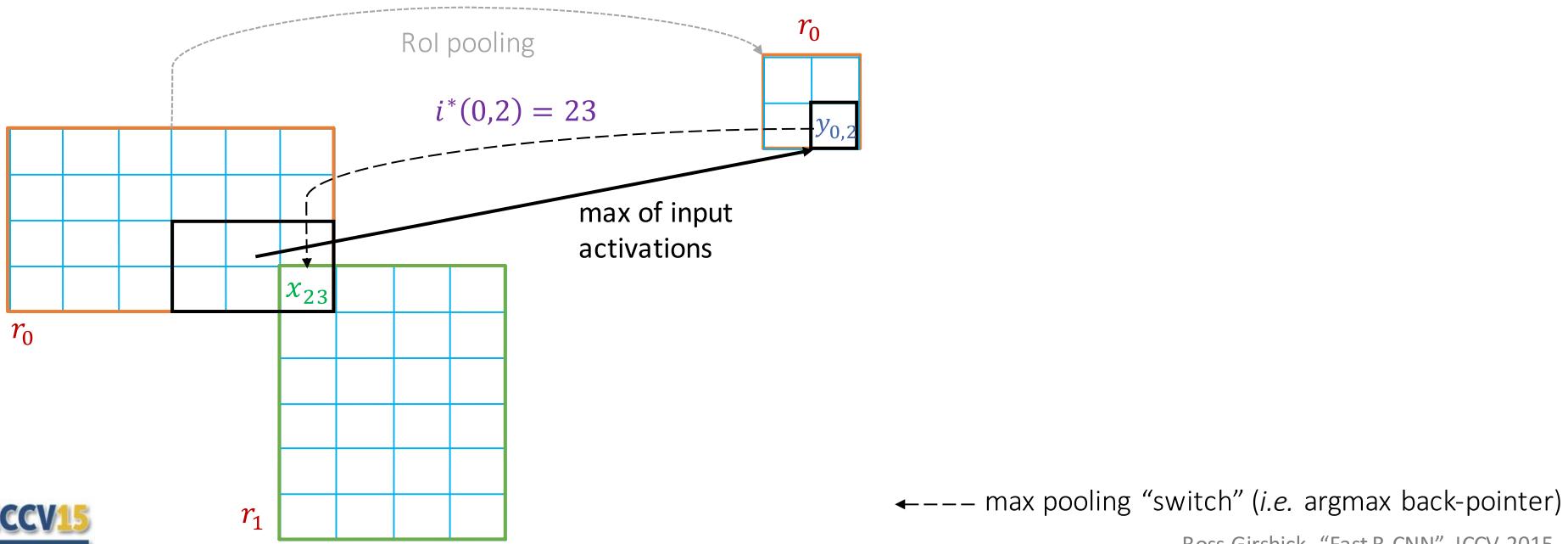
RoI pooling / 'SPP' is just like 'max' pooling, except that 'pooling' regions overlap



Ross Girshick. "Fast R-CNN". ICCV 2015

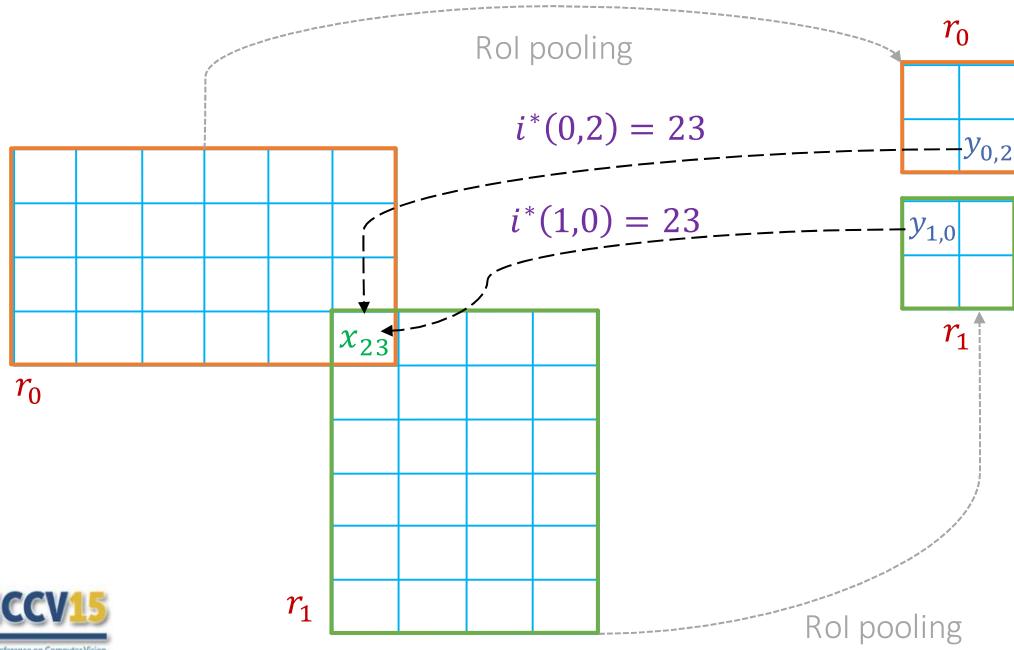
Obstacle #1: Differentiable RoI pooling

RoI pooling / SPP is just like max pooling, except that pooling regions overlap



Obstacle #1: Differentiable RoI pooling

RoI pooling / SPP is just like max pooling, except that pooling regions overlap



$$\frac{\partial L}{\partial x_i} = \sum_r \sum_j [i = i^*(r, j)] \frac{\partial L}{\partial y_{rj}}$$

1 if r, j “pooled”
input i ; 0 o/w

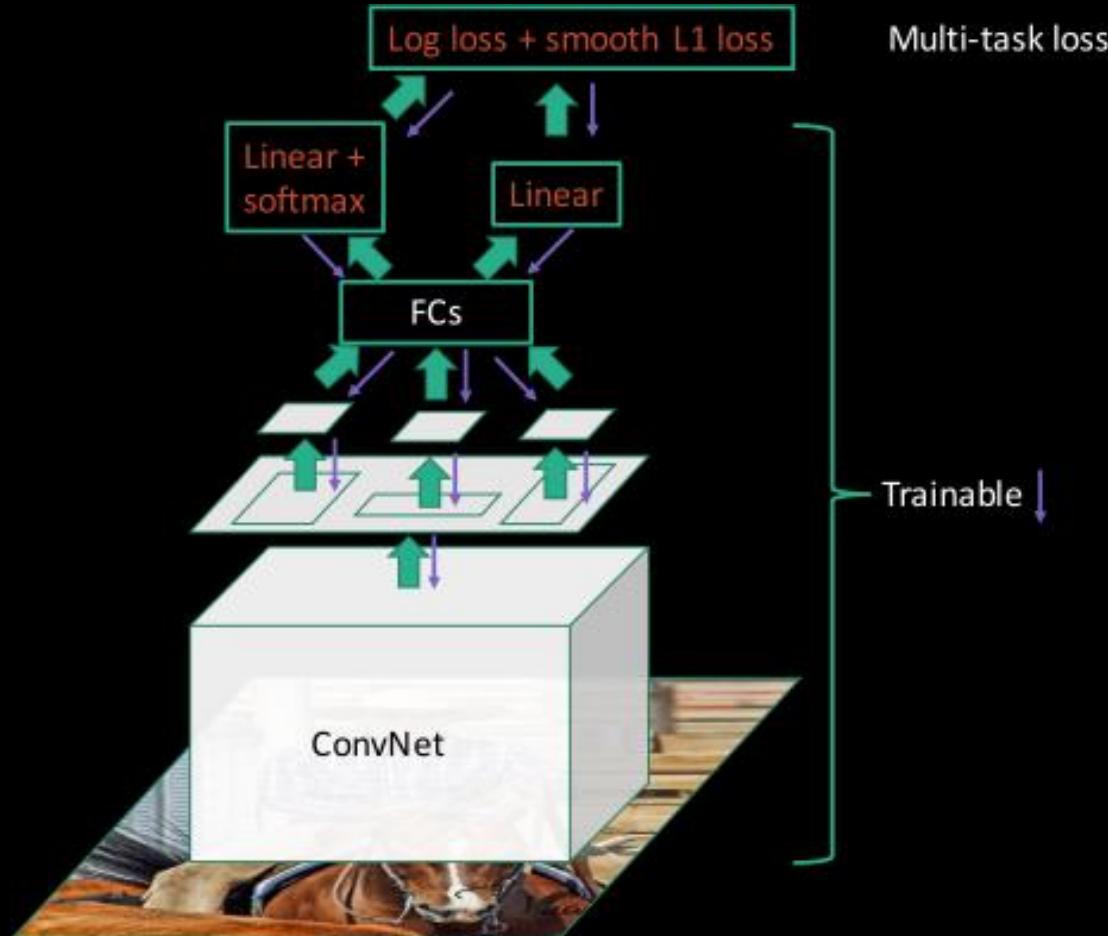
Partial Over regions r ,
for x_i locations j

Partial from next layer

←---- max pooling “switch” (i.e. argmax back-pointer)

Ross Girshick. “Fast R-CNN”. ICCV 2015.

Object Detection

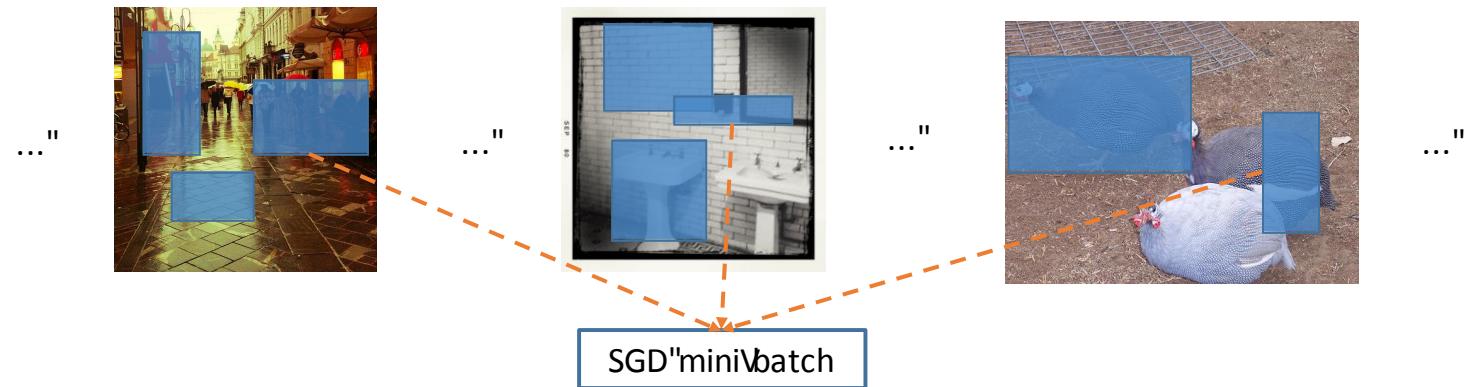


Fast R-CNN (Girshick, 2015)

Obstacle '#2: 'Making' SGD's steps' efficient

Slow'R)CNN'and'SPP)net'use'region)wise'sampling'to'make'mini)batches

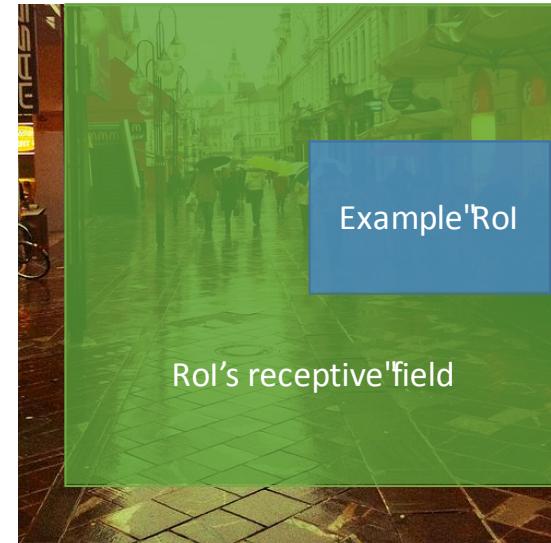
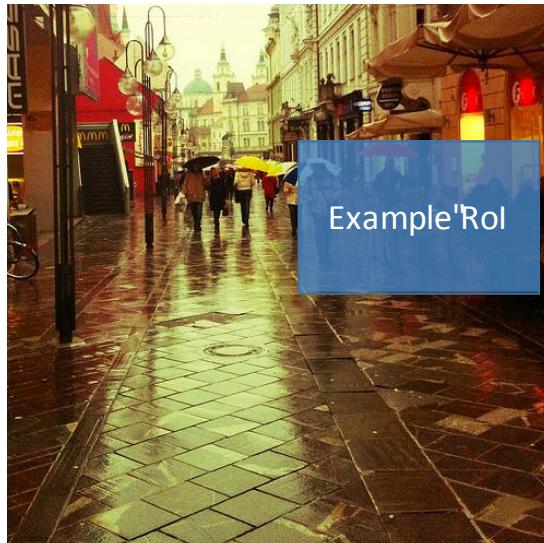
- Sample'128'example'Rolls uniformly'at'random
- Examples'will'come'from'different'images'with'high'probability



Obstacle #2: Making SGD steps efficient

Note the receptive field for one example ROI is often very large

- Worst case: the receptive field is the entire image



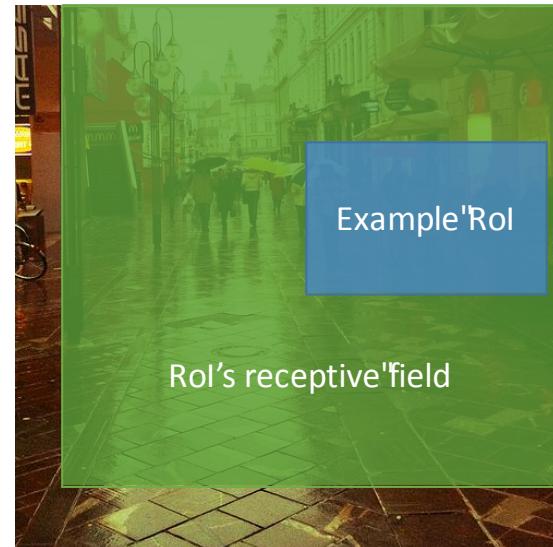
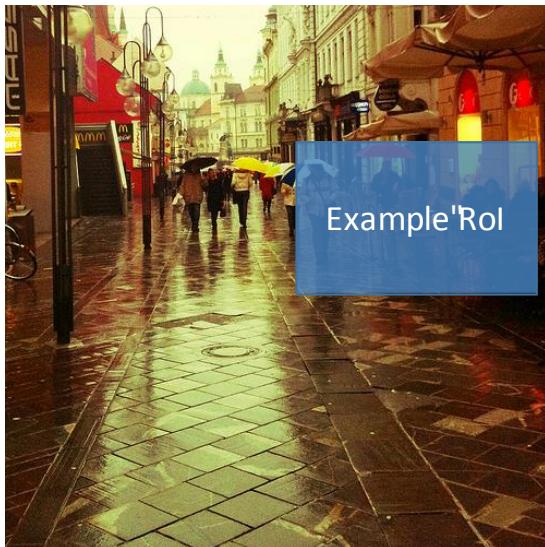
Obstacle #2: Making SGD steps efficient

Worst case cost per mini-batch (crude model of computational complexity)

input size for Fast R-CNN

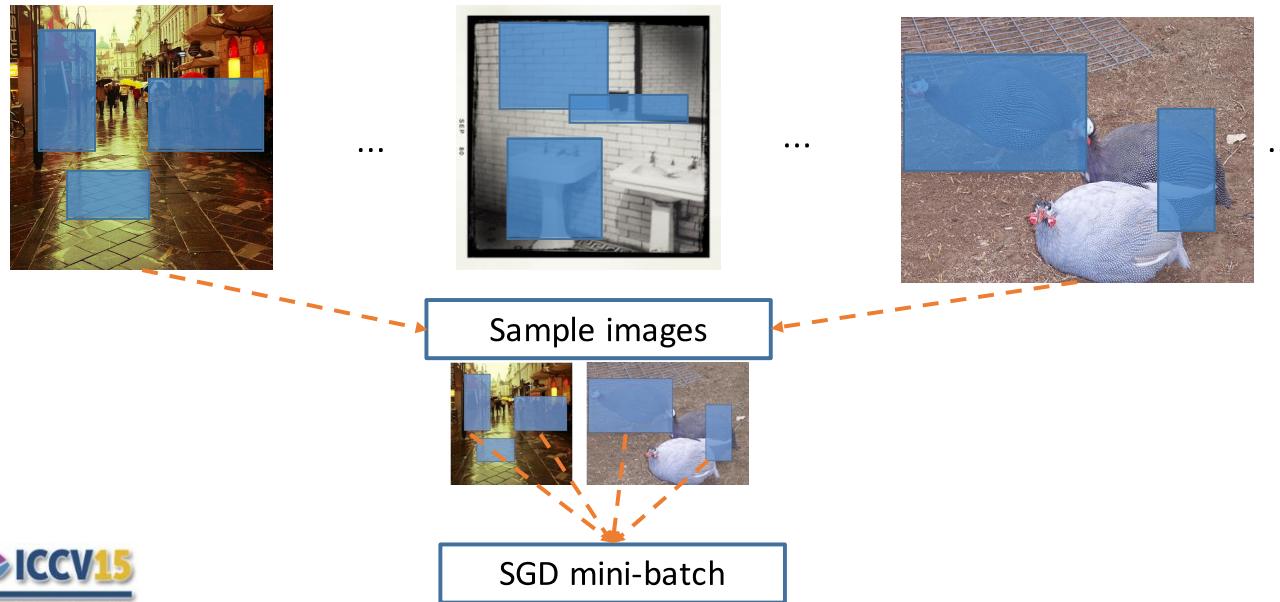
input size for slow R-CNN

- $128*600*1000 / (128*224 *224) = 12x >$ computation than slow R-CNN



Obstacle #2: Making SGD steps efficient

Solution: use hierarchical sampling to build mini-batches

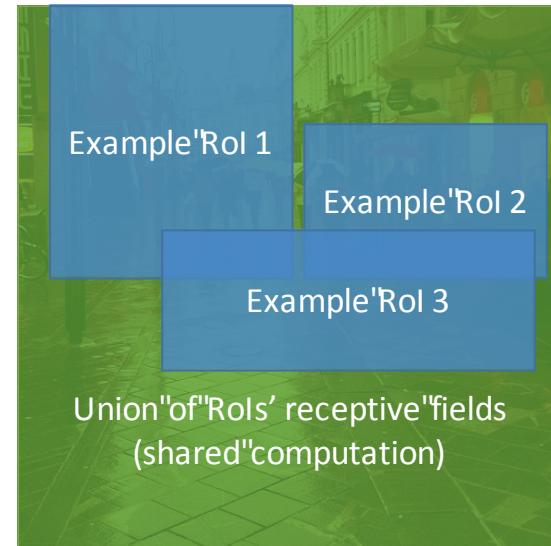


- Sample a **small number of images** (2)
- Sample **many examples from each image** (64)

Obstacle '#2: 'Making' SGD's steps' efficient

Use 'the' test) time 'trick' from 'SPP) net' during 'training

- Share 'computation' between 'overlapping' examples 'from' the 'same' image



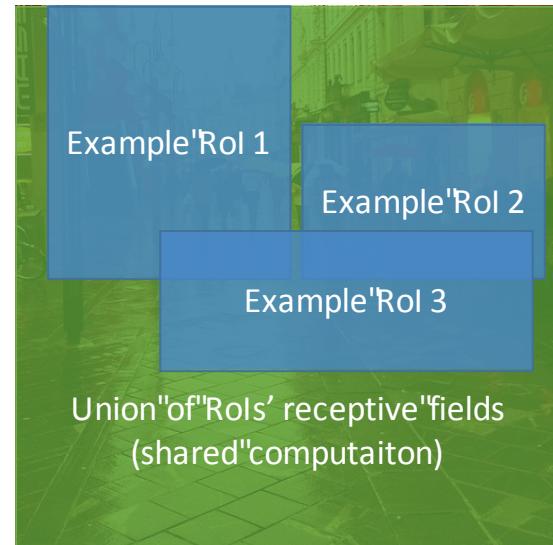
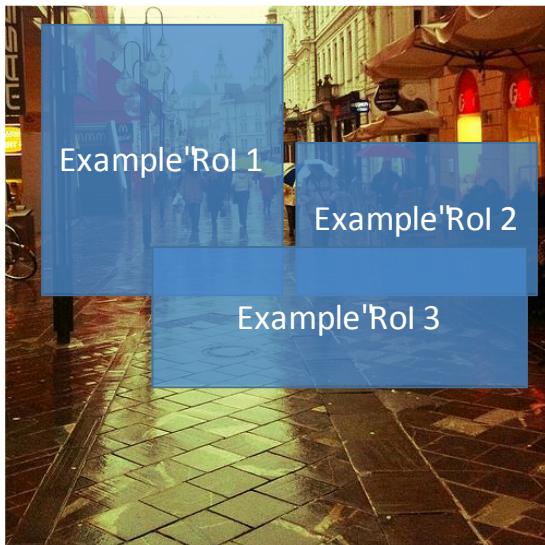
Obstacle #2: Making SGD steps efficient

Cost per mini-batch compared to slow R-CNN (same crude cost model)

input size for Fast R-CNN

input size for slow R-CNN

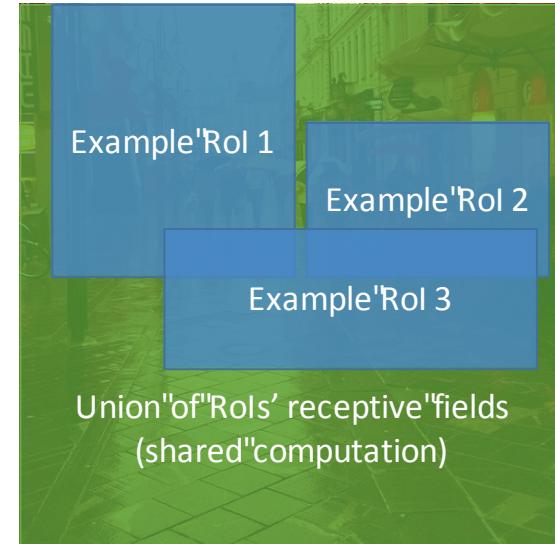
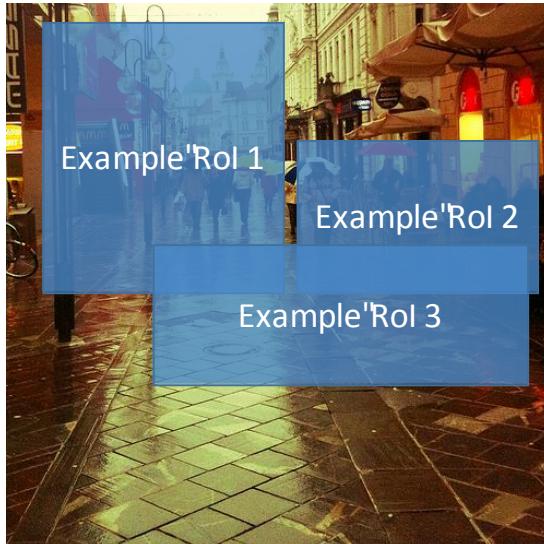
- $2*600*1000 / (128*224*224) = 0.19x < \text{computation than slow R-CNN}$



Obstacle '#2: Making SGD steps efficient

Are the examples from just 2 images diverse enough?

- Concern: examples from the sample image may be too **correlated**



Fast R-CNN outcome

Better training time and testing time with better accuracy than slow R-CNN or SPP-net

- Training time: 84 hours / 25.5 hours / 8.75 hours (Fast R-CNN)
- VOC07 test mAP: 66.0% / 63.1% / 68.1%
- Testing time per image: 47s / 2.3s / 0.32s
 - Plus 0.2 to >2s per image depending on proposal method
 - With selective search: 49s / 4.3s / 2.32s

Updated numbers from the ICCV paper based on implementation improvements

Experimental findings

- End-to-end training is important for very deep networks
- Softmax is a fine replacement for SVMs
- Multi-task training is beneficial
- Single-scale testing is a good tradeoff (noted by Kaiming)
- Fast training and testing enables new experiments
 - Comparing proposals

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

Source: R. Girshick

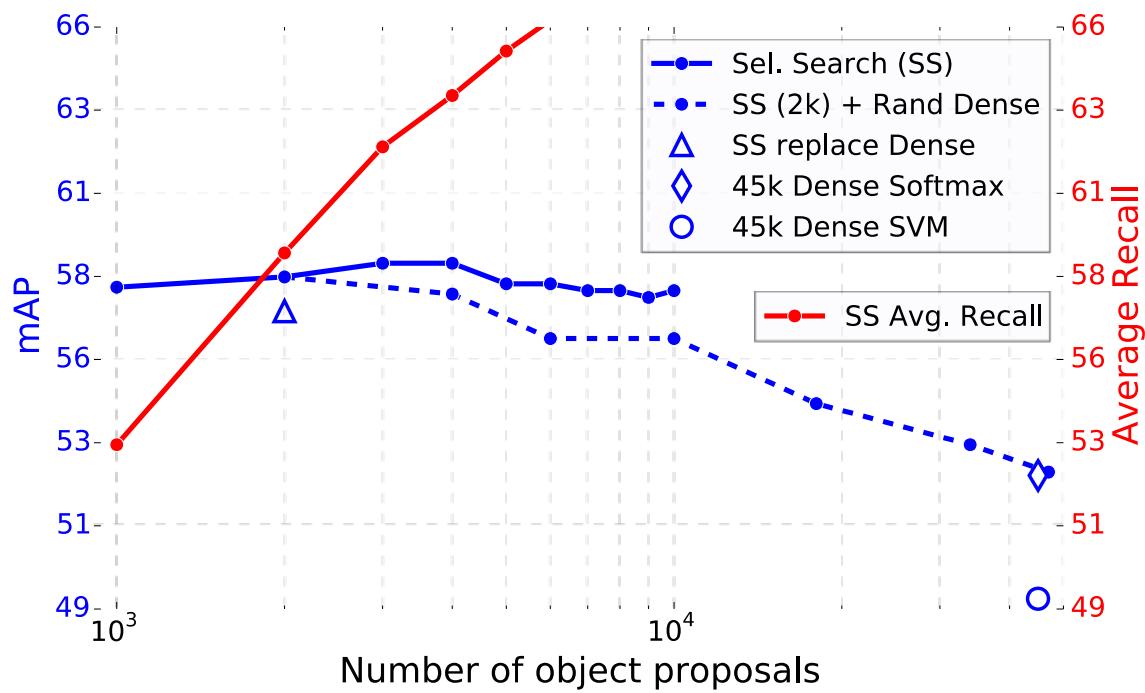
Result compare

	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

Source: cs231 standford

Direct 'region' proposal' evaluation

- VGG_CNN_M_1024
- Training takes < 2 hours
- Fast training makes these experiments possible



Ross Girshick. "Fast R-CNN". ICCV 2015.

Drawback of R-CNN and the modification:

1. Training is a multi-stage pipeline. -> End-to-end joint training.
2. Training is expensive in space and time. -> Convolutional layer sharing. Classification in memory.

For SVM and regressor training, features are extracted from each warped object proposal in each image and written to disk.(VGG16, 5k VOC07 trainval images : 2.5 GPU days). Hundreds of gigabytes of storage.

3. Test-time detection is slow. -> Single scale testing, SVD fc layer.

At test-time, features are extracted from each warped proposal in each img. (VGG16: 47s / image).

Contributions:

1. Higher detection quality (mAP) than R-CNN
2. Training is single-stage, using a multi-task loss
3. All network layers can be updated during training
4. No disk storage is required for feature caching

•Mini-batch Sampling

128: 2 randomly sampled images with 64 Pol sampled from each image

25% positive: IoU > 0.5

75% background:IoU in [0.1, 0.5)

horizontally flipped with prob = 0.5

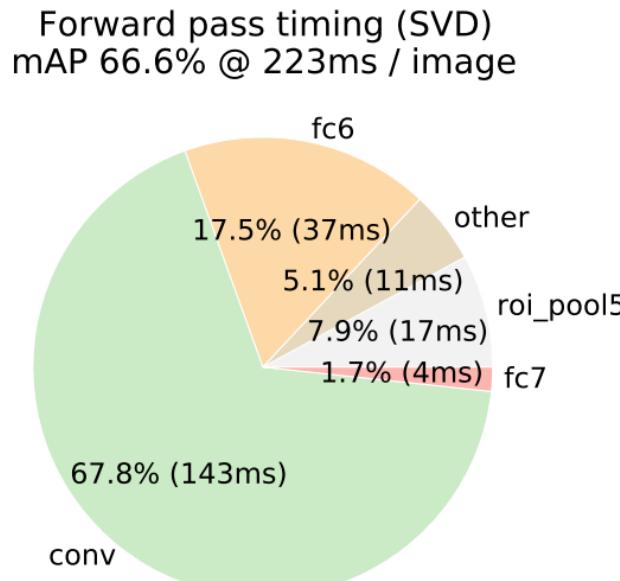
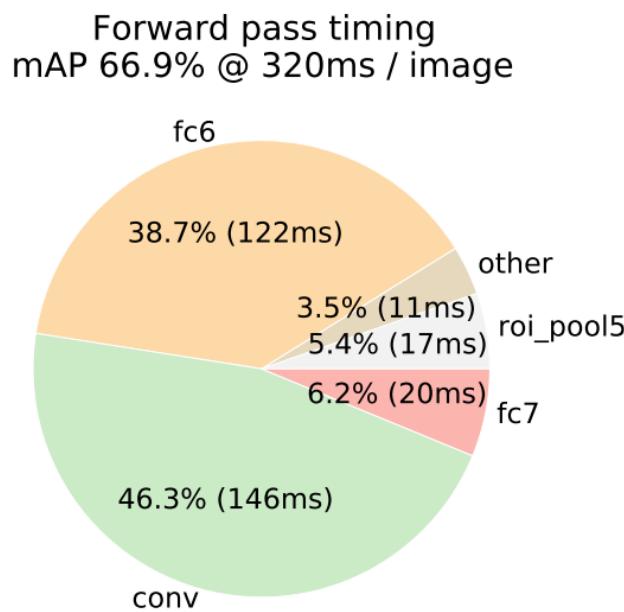
Truncated SVD for faster detection

mAP ~ 0.3% down; speed ~ 30% up

number of RoI for detection is large -> time spent on fc

$$W \sim U \Sigma_t V^T \quad (U : u^*t, \text{Sigma}_t : t^*t, V : v^*t)$$

$$\text{Compression} : (Wx + b) \text{ fc} \rightarrow (\Sigma_t V^T x) \text{ fc} + (Ux + b) \text{ fc}$$



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

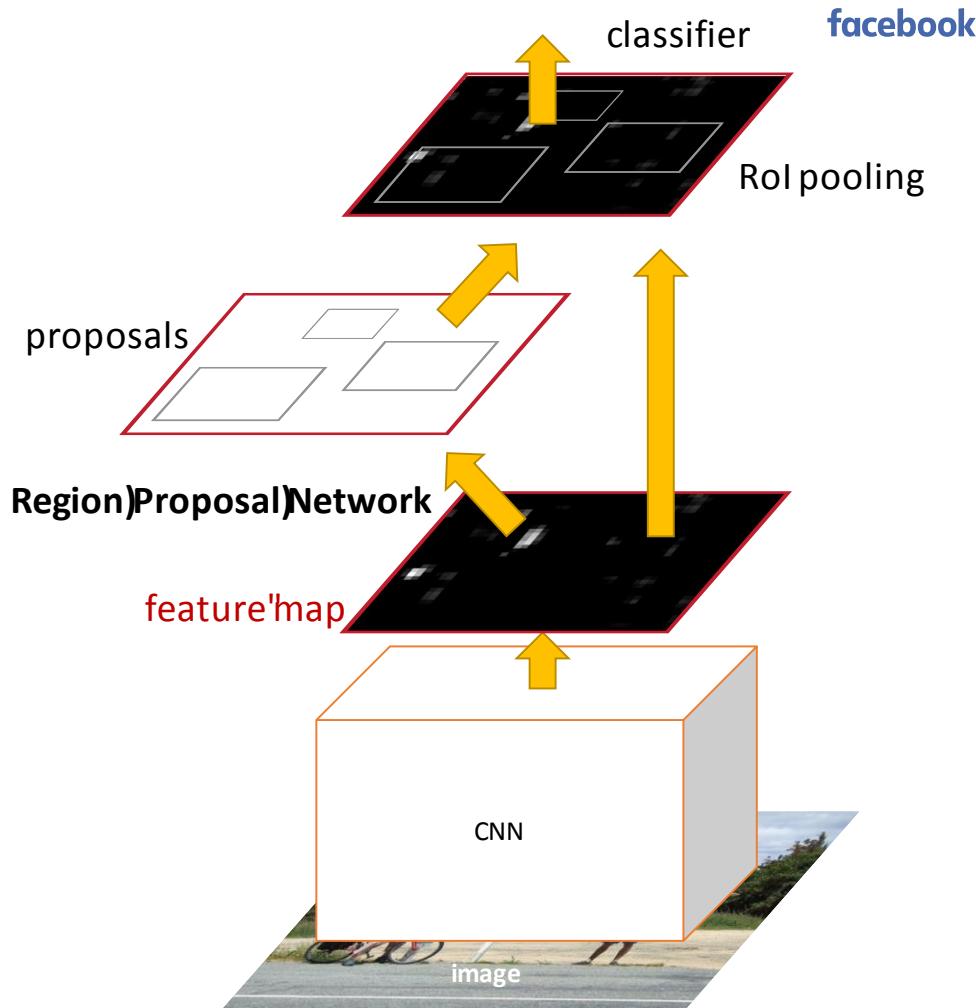
Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun

Jun 2015

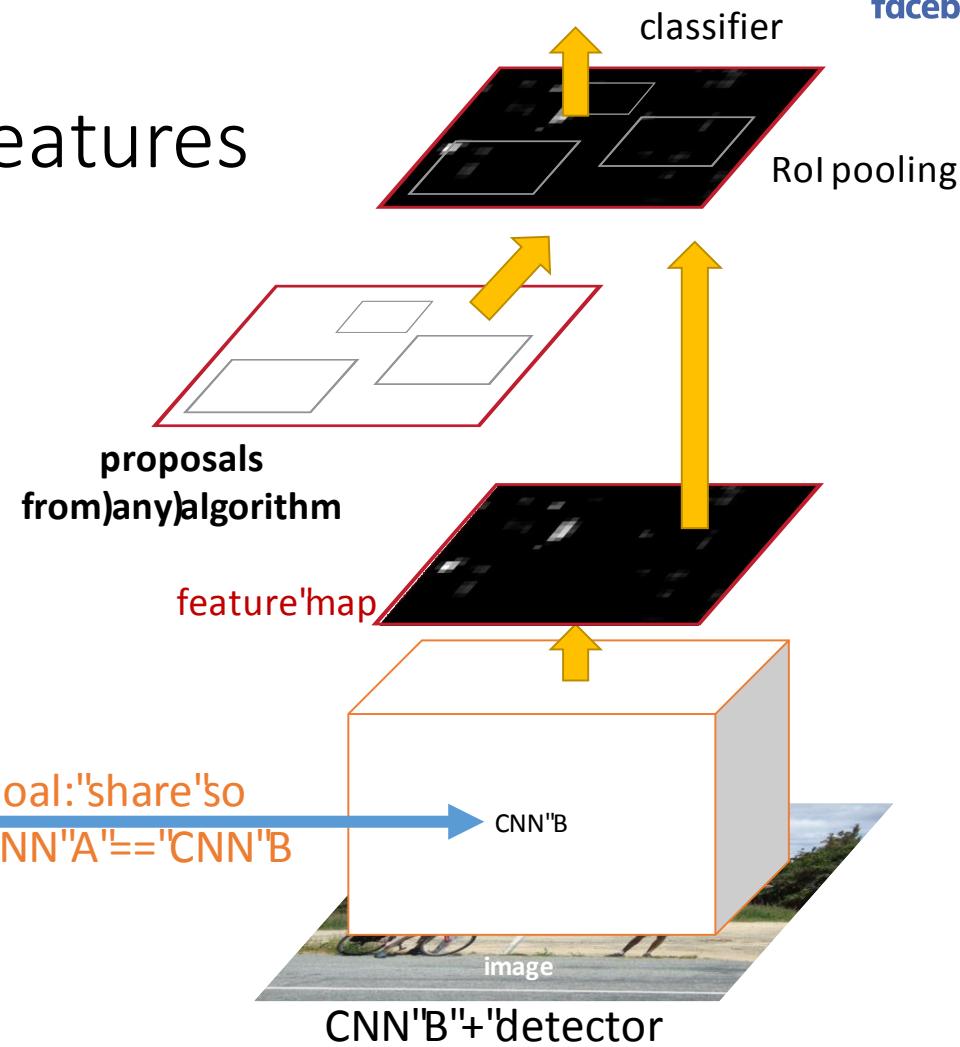
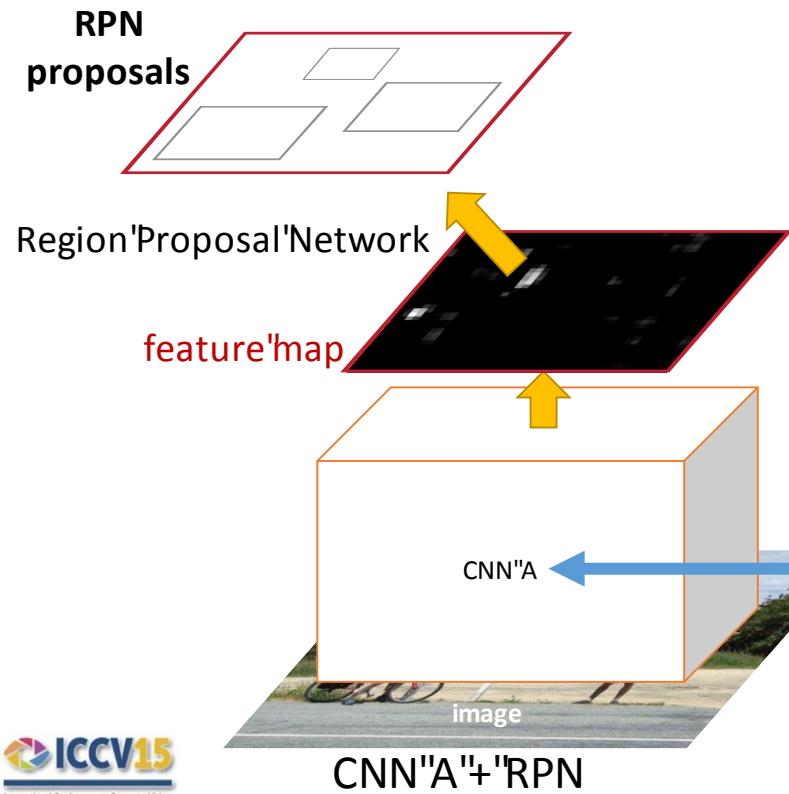
<https://arxiv.org/pdf/1506.01497v3.pdf>

What's Faster R)CNN?

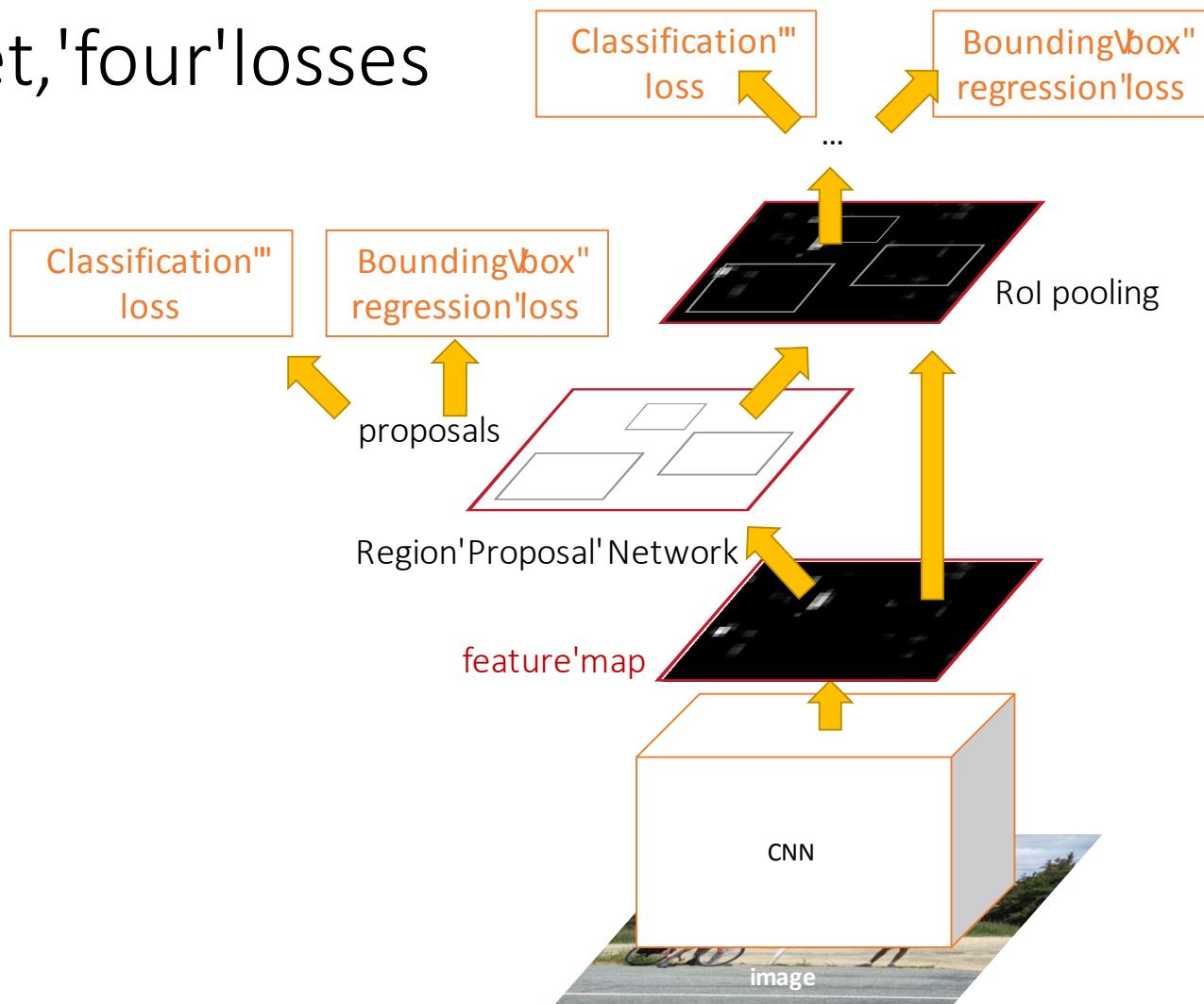
- Presented in Kaiming's section
- Review:
Faster R)CNN = Fast R)CNN + Region Proposal Networks
 - Does not depend on an external region proposal algorithm
 - Does object detection in a single forward pass



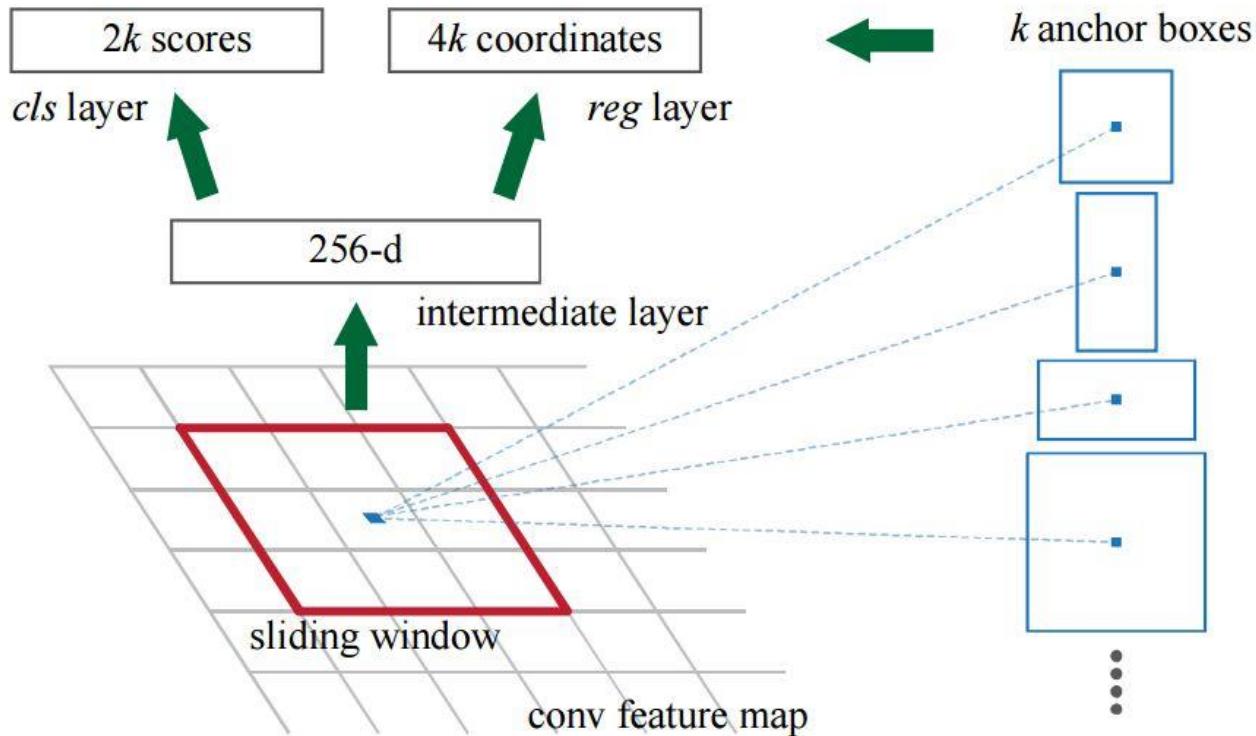
Training goal: Share features



One'net,'four'losses



Faster RCNN



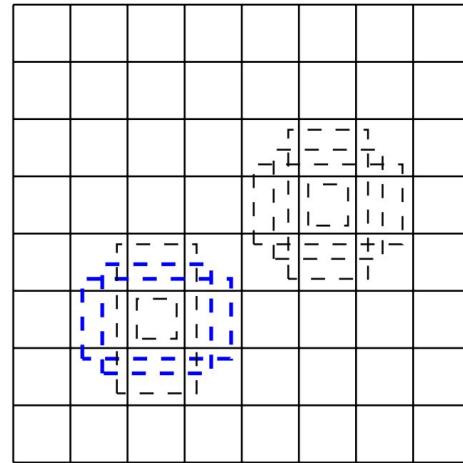
<https://arxiv.org/pdf/1506.01497v3.pdf>

- Translation-Invariant Anchors

At each sliding window loc, predict k proposal: 4k outputs for reg layer, 2k outputs for cls layer (binary softmax).

Anchor: centered at sliding window with scale and aspect ratio: $(128^2, 256^2, 512^2; 1:2, 2:1, 1:1)$

For a conv feature map: $W * H * k$ ($k=9$ anchors) $(2+4)*9$ output layer



- Optimization

fcn trained by end-to-end by bp and sgd

image-centric sampling strategy, sample 256 anchors in an image(Pos:neg = 1:1)

new layer initialization $\sim N(0, 0.01)$

tune ZFnet and conv3_1 and up for VGGnet, lr=0.001 for 60k batches, 0.0001 for 20k on PASCAL

- **Loss function for Learning Region Proposal**

positive label: the anchor has highest IoU with a gt-box or has an IoU>0.7 with any gt-box

negative label: IoU<0.3 for all gt-box

Objective function with multi-task loss: Similar to Fast R-CNN.

$$L(p_i, t_i) = L_{cls(p_i, p_i^*)} + \lambda p_i^* L_{reg}(t_i, t_i^*)$$

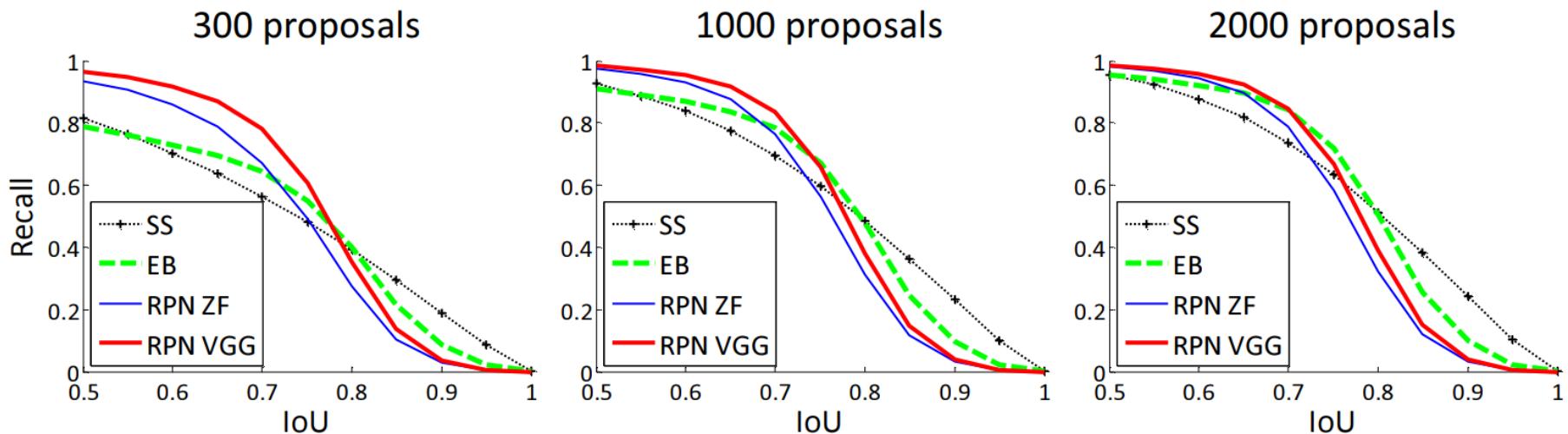
where p_i^* is 1 if the anchor is labeled positive, and is 0 if the anchor is negative.

$\lambda = 10$ bias towards better box location

- Share Convolutional Features for Region Proposal and Objection Detection

Four-step training algorithm:

1. Train RPN, initialized with ImageNet pre-trained model
 2. Train a separate detection network by Fast R-CNN using proposals generated by step-1 RPN, initialized by ImageNet pre-trained model
 3. Fix conv layer, fine-tune unique layers to RPN, initialized by detector network in Step2
 4. Fix conv layer, fine-tune fc-layers of Fast R-CNN



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

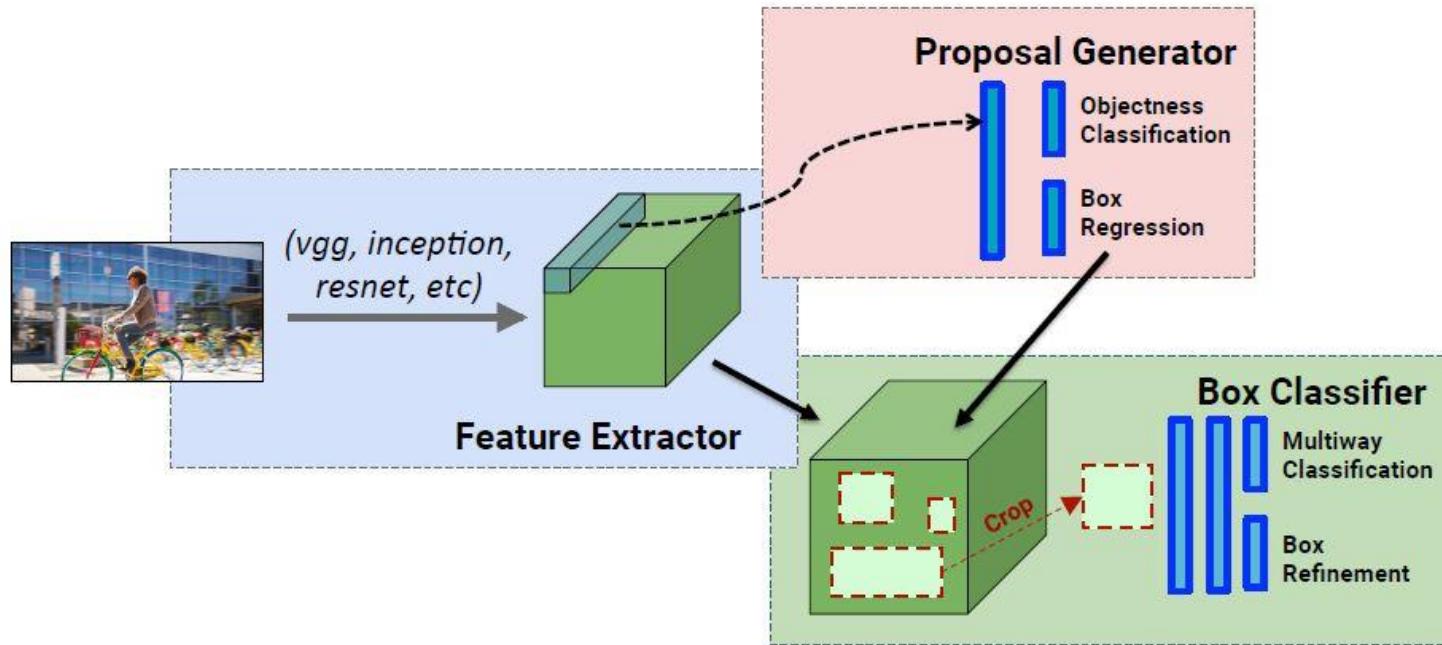
Trained using Pascal VOC 2007 dataset

Source: cs231 standford

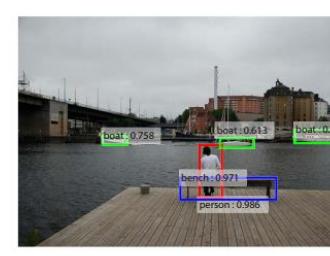
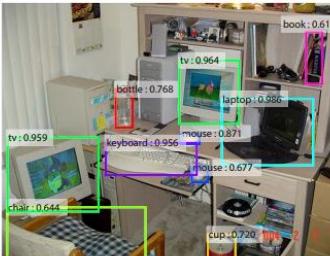
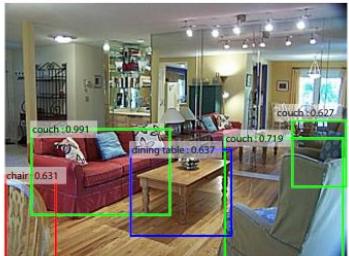
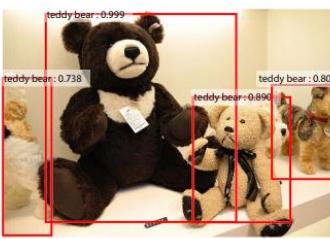
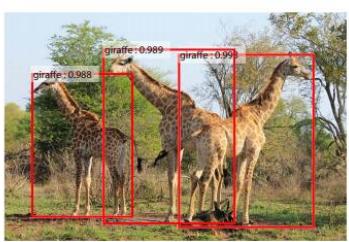
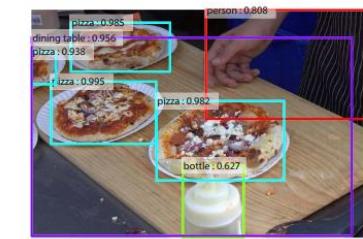
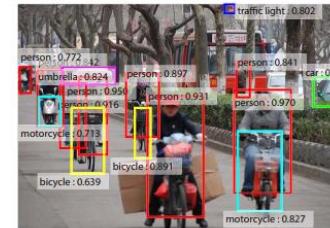
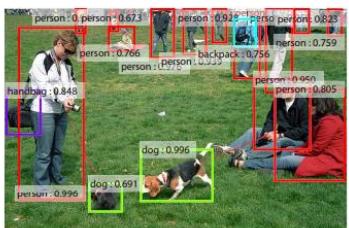
Table 5: **Timing** (ms) on a K40 GPU, except SS proposal is evaluated in a CPU. “Region-wise” includes NMS, pooling, fully-connected, and softmax layers. See our released code for the profiling of running time.

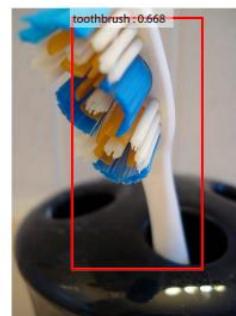
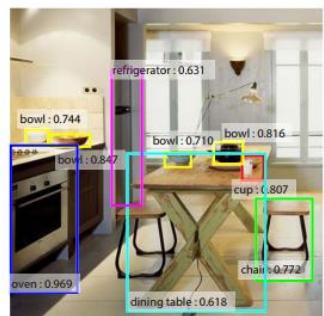
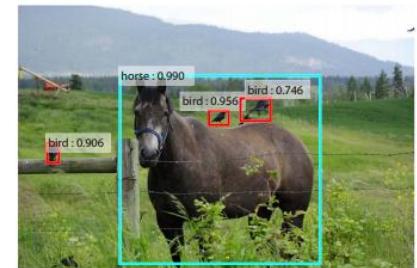
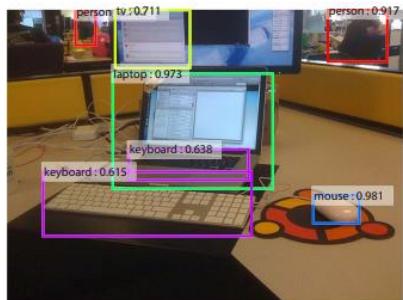
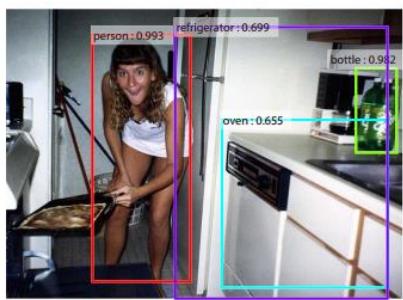
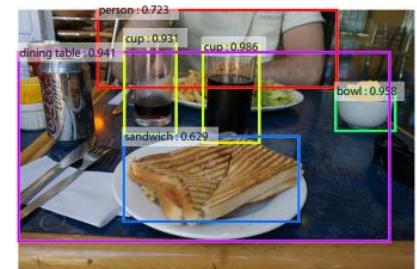
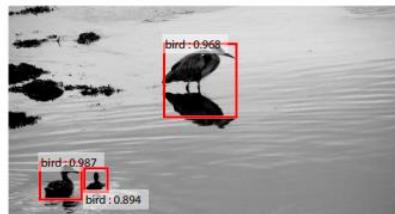
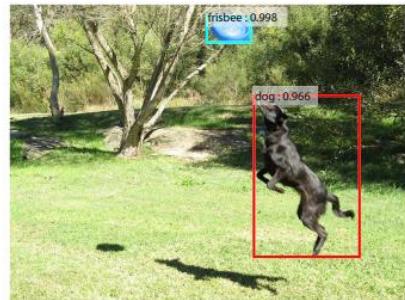
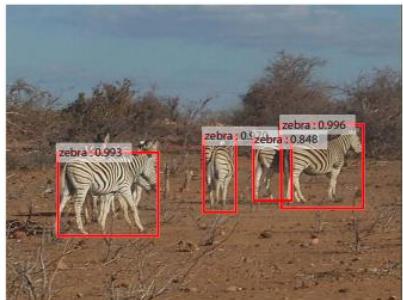
model	system	conv	proposal	region-wise	total	rate
VGG	SS + Fast R-CNN	146	1510	174	1830	0.5 fps
VGG	RPN + Fast R-CNN	141	10	47	198	5 fps
ZF	RPN + Fast R-CNN	31	3	25	59	17 fps

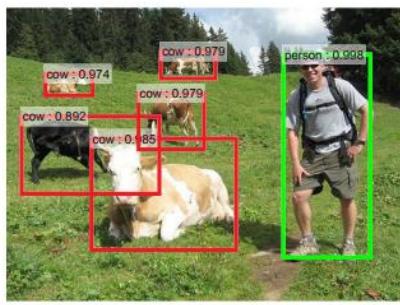
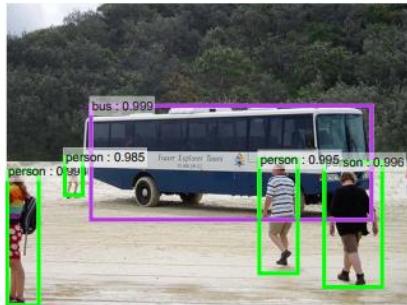
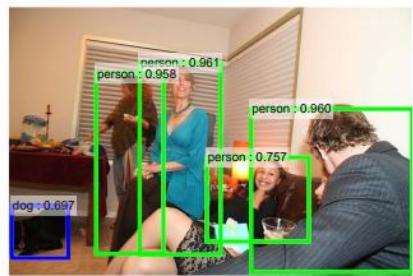
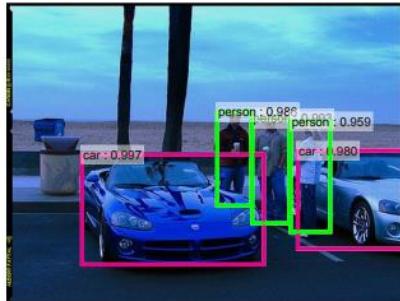
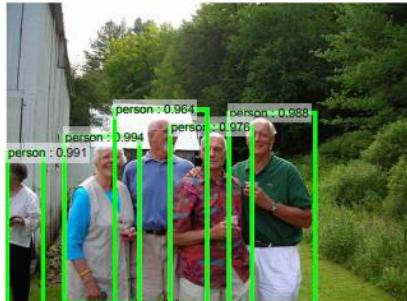
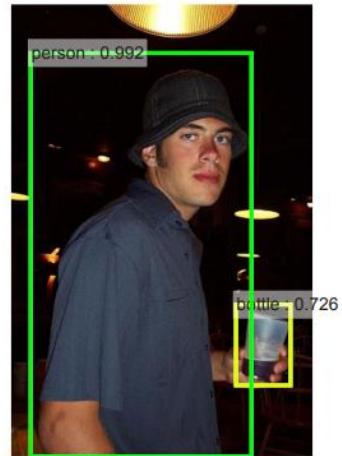
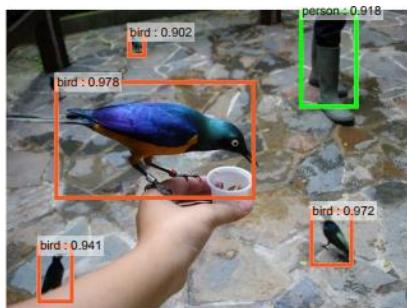
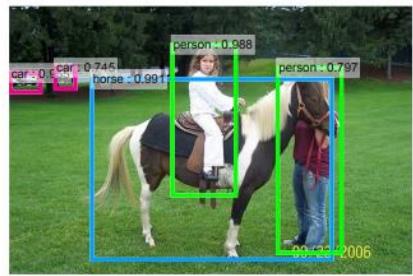
Recap: Faster RCNN

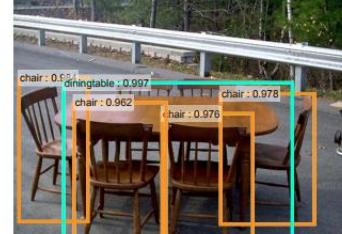
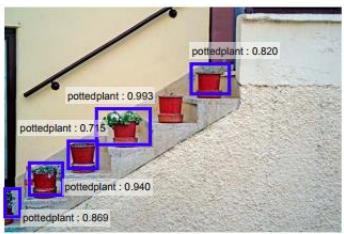
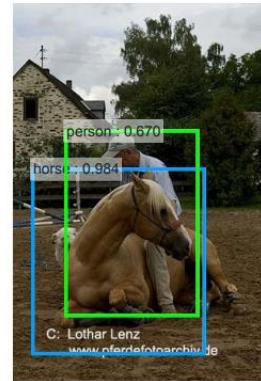
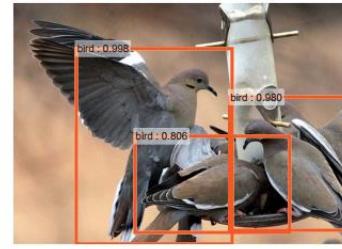
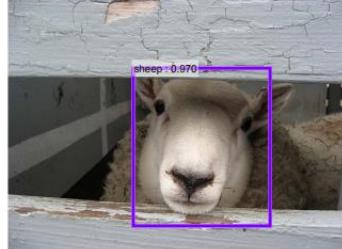
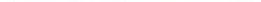
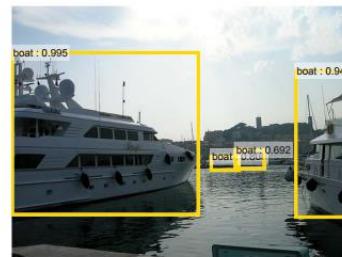
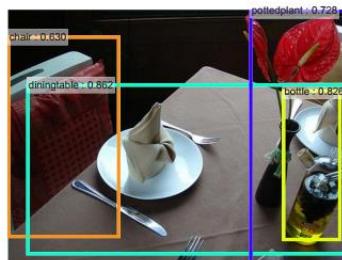
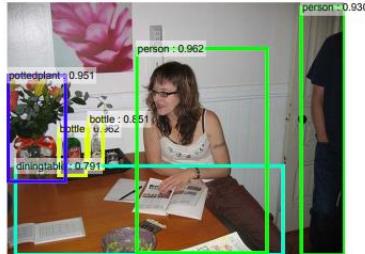
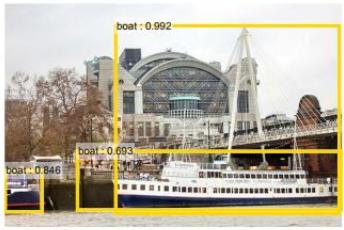
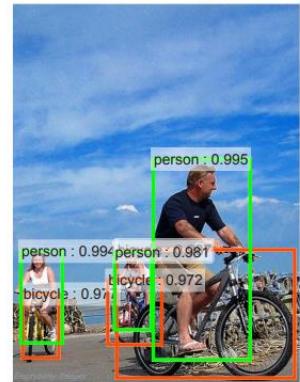
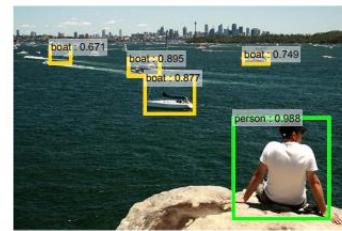


<https://pan.baidu.com/s/1pKIKIIB>









- Recognition = Description + Discrimination. We are doing well on Discrimination
- Challenges of Deformation and Occlusion. This is also known as hard negative mining and non-maximum suppression. The past solutions assume the features are hard coded. CNN solutions extract ‘deformable’ features. Occlusion is dealt with per pixel labels over many object classes.
- R-CNN, Fast R-CNN, Faster R-CNN: sparse sampling of ‘where’ jointly with ‘what’. Using one shared feature extraction front-end, and multiple regression objectives at the FC layers. Enough information from the front-end CNN features allows both detection (of location) and classification (of classes).
- Memory allocation and flow is important for CNN learning and speed up.