

From Visual Subcategory to Webly-supervised Visual Recognition

Moin Nabi

Santosh Divalla (Allen Institute for Artificial Intelligence)

Ali Farhadi (University of Washington)

Massimiliano Pontil (UCL)

Vittorio Murino (Italian Institute of Technology)

Outline

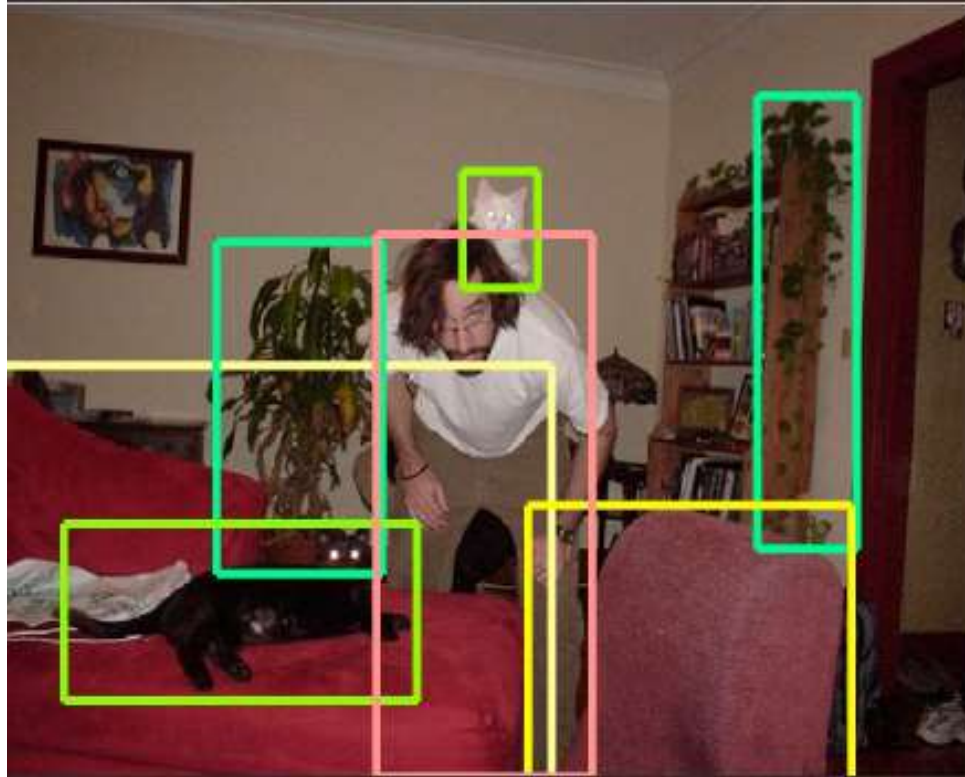
- Introduction on object detection
- Visual subcategory
- Evolution of DPM
- web-supervised visual recognition
- Webly-supervised discriminative patch
- Subcategory and dataset bias

Question



What objects are **where**?

Goal: detecting objects in cluttered images



person, plant, cat, dog chair, sofa, car, bicycle, motorbike, table, plane, ...

Application

- Applications

- Semantic image and video search
- Human-computer interaction (e.g., Kinect)
- Automotive safety
- Camera focus-by-detection
- Surveillance
- Semantic image and video editing
- Assistive technologies
- Medical imaging

Challenges



Variation in illumination



Variation in pose, viewpoint



Variation in appearance



Occlusion and clutter

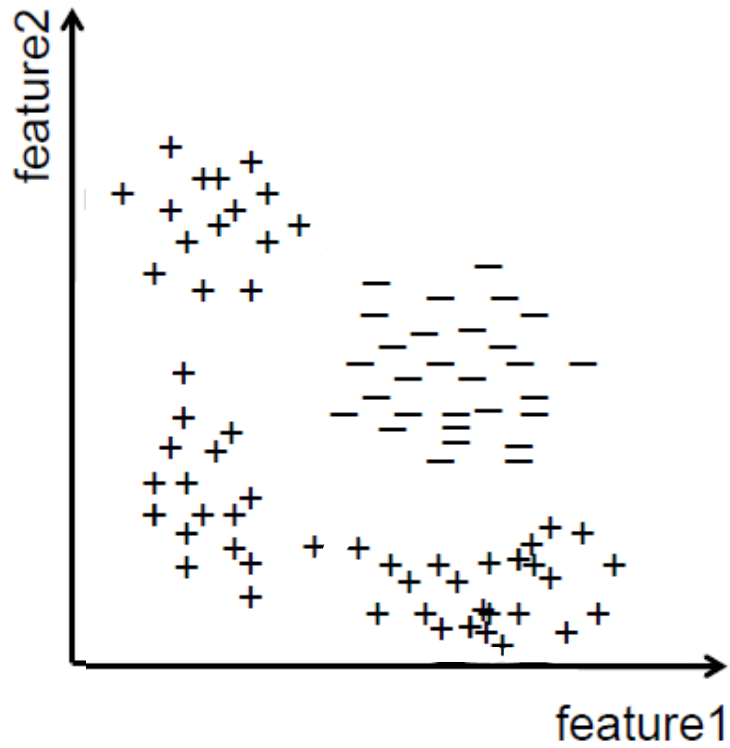
Intra Category Diversity



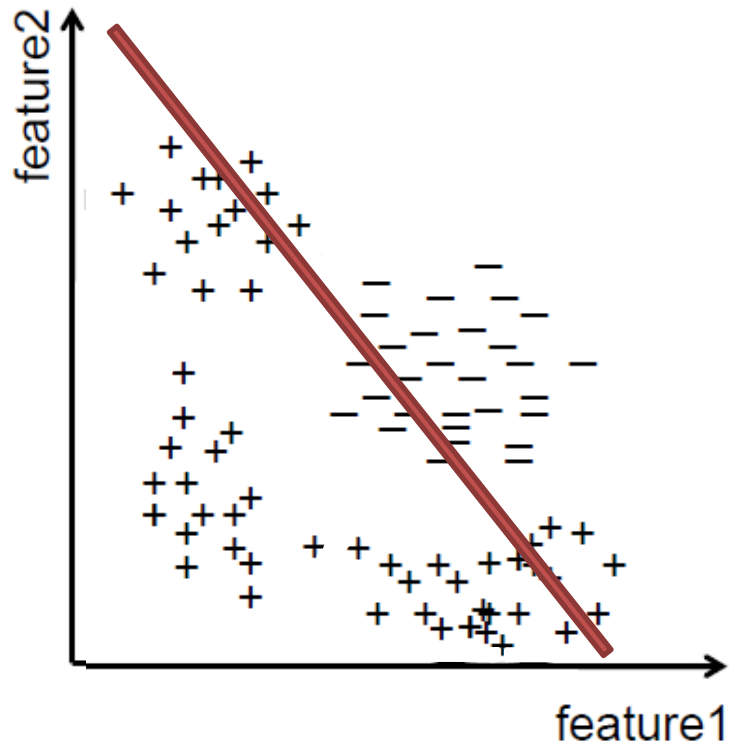
Example images for “Horse” from PASCAL VOC

Variation due to change in camera viewpoint, object pose, and occlusion

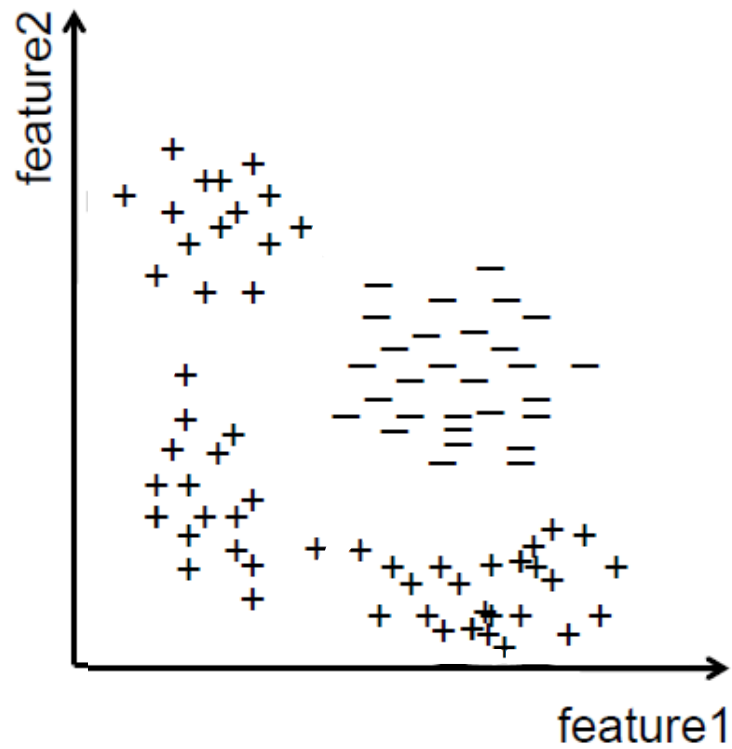
Subcategories



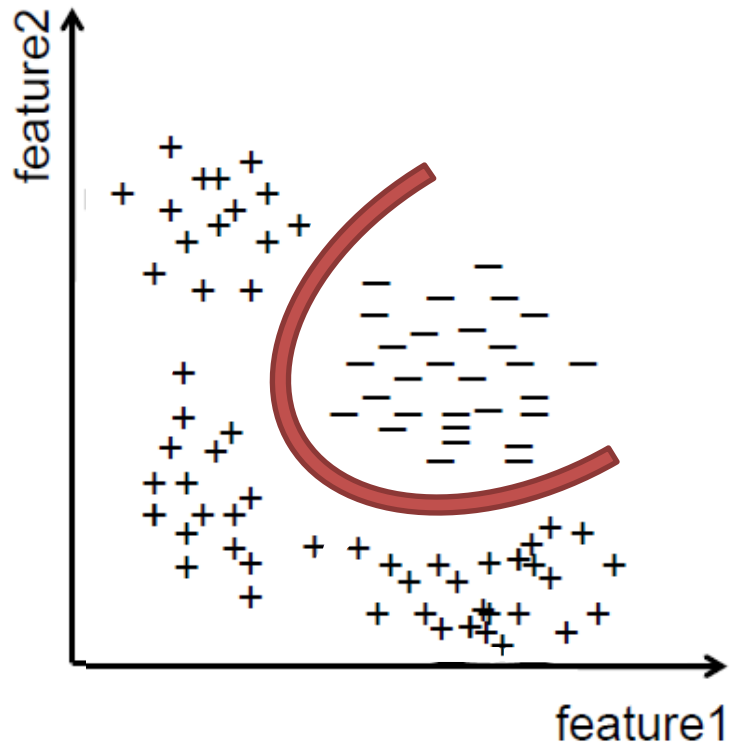
Subcategories



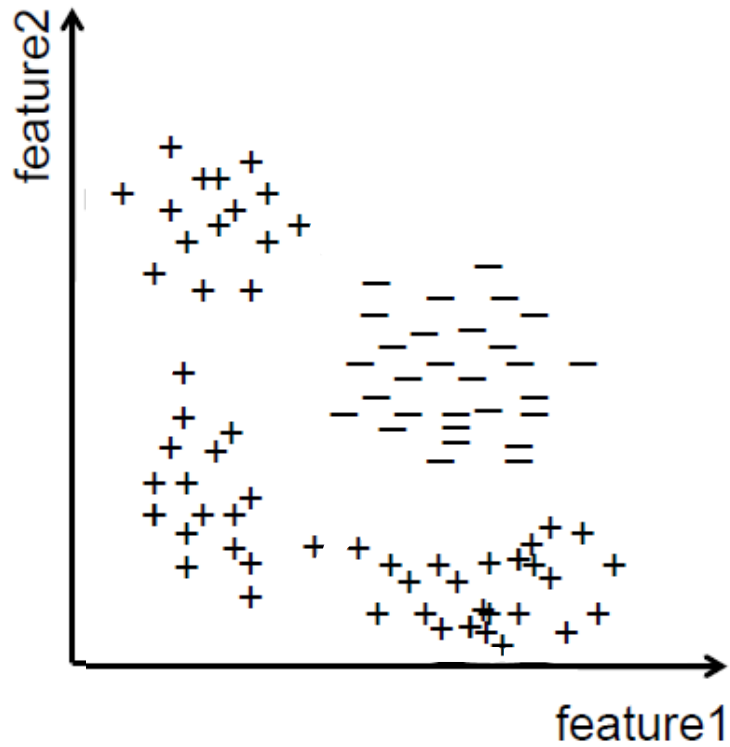
Subcategories



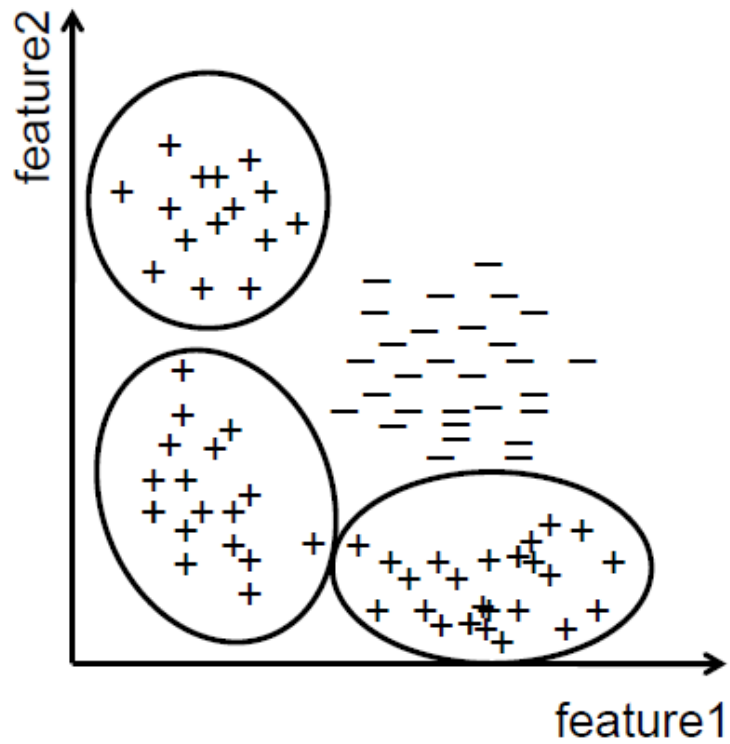
Subcategories



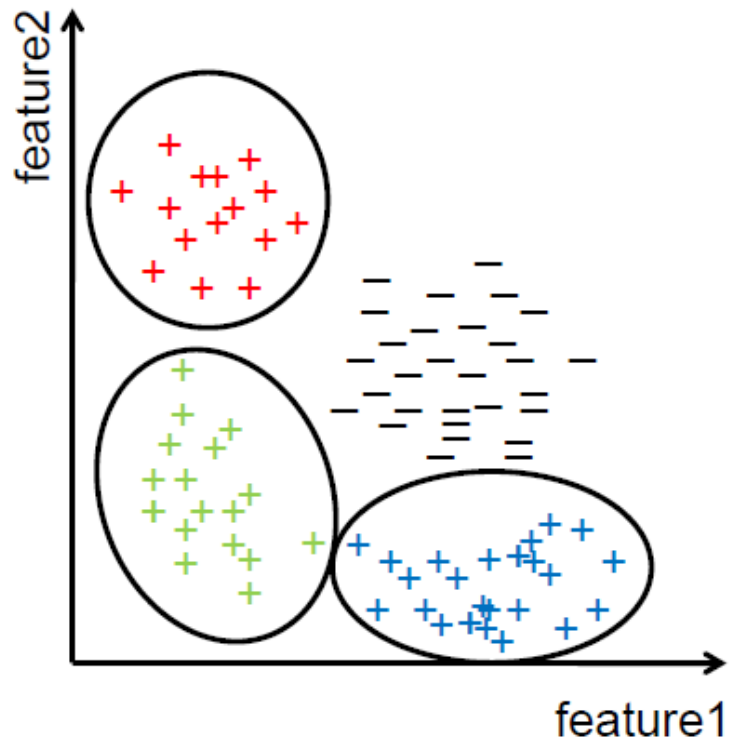
Subcategories



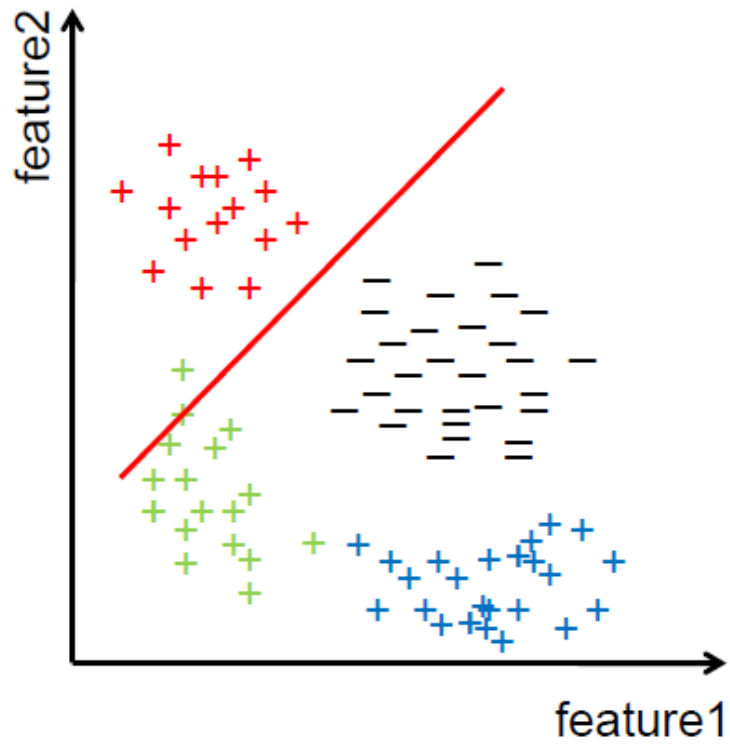
Subcategories



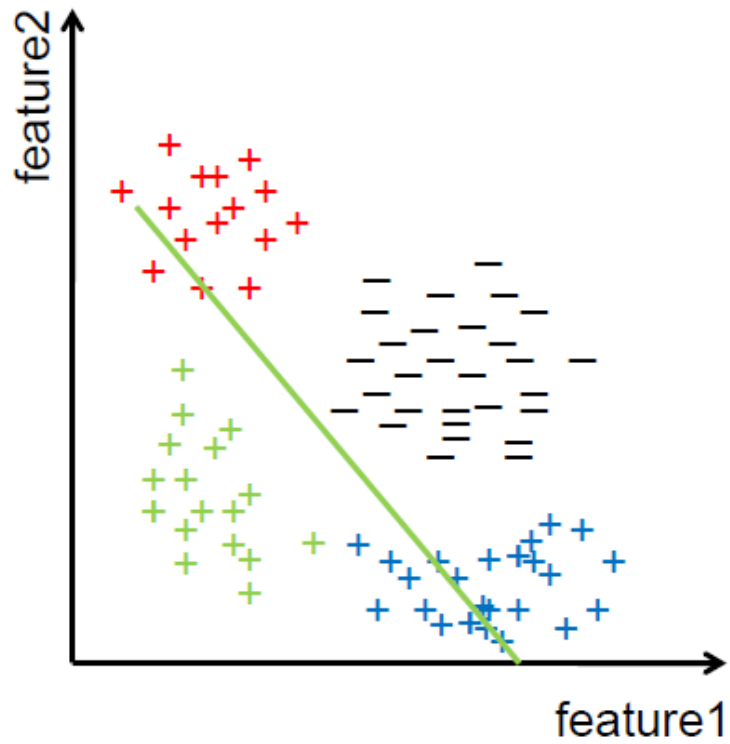
Subcategories



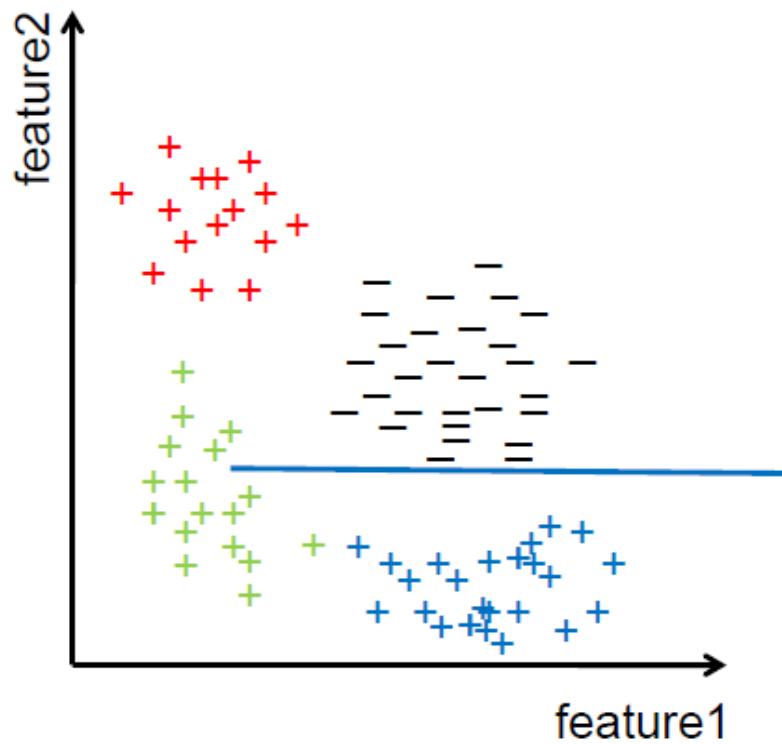
Subcategories



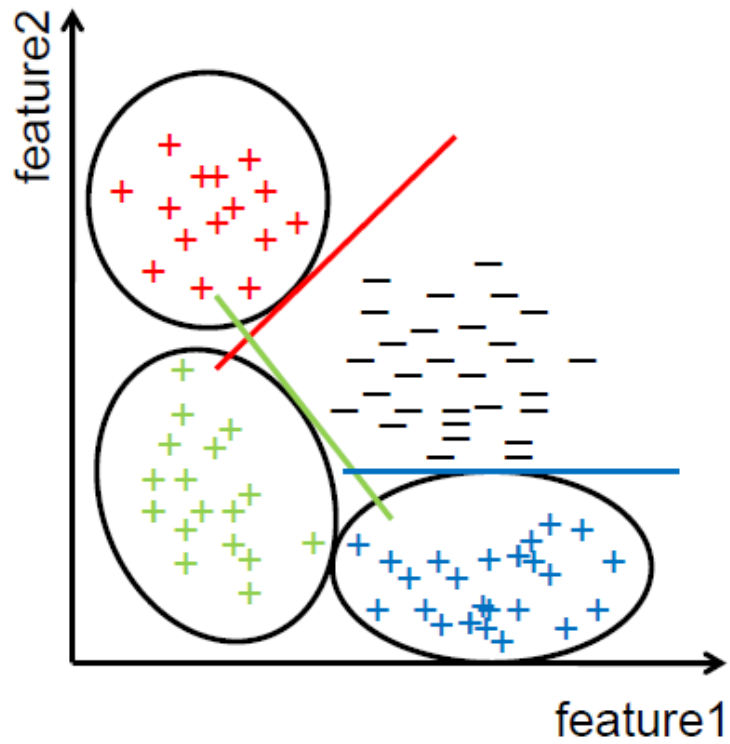
Subcategories



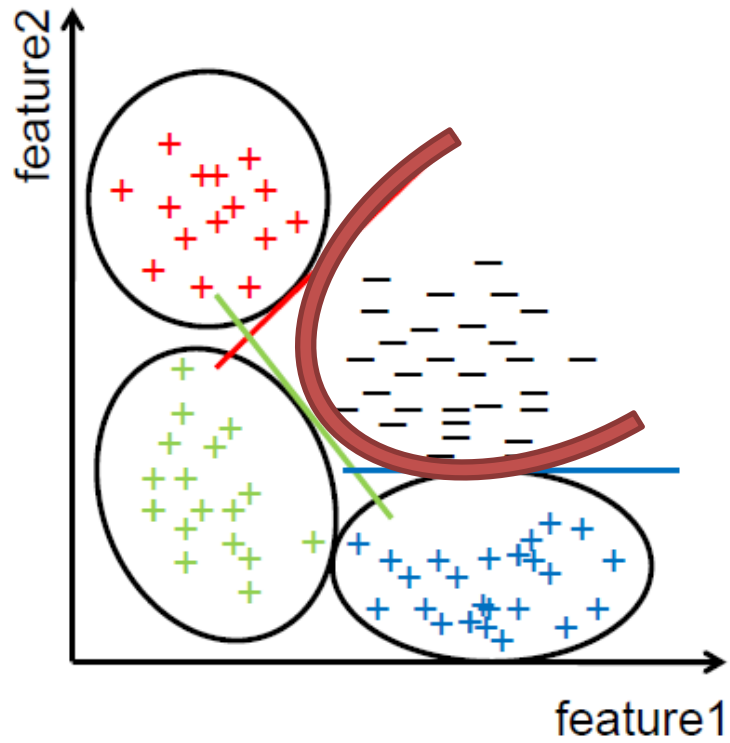
Subcategories



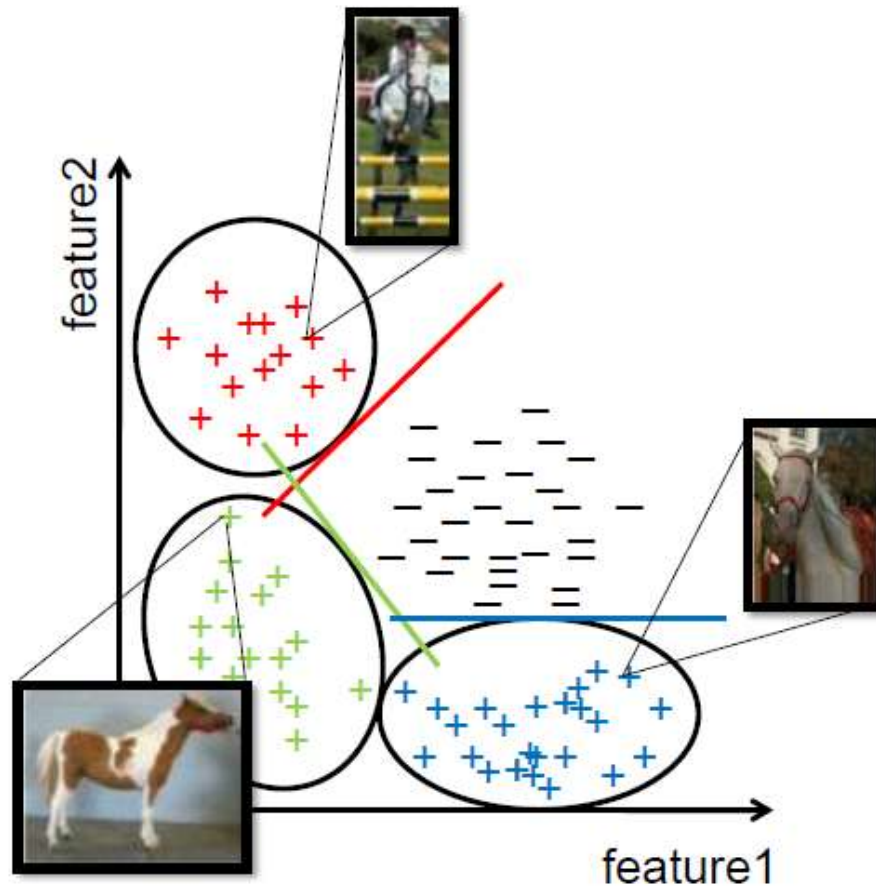
Subcategories



Subcategories



Subcategories



Subcategories



The Evolution of DPM

| |
|----------------------------|
| K=1 no Parts [DT'05] |
| mAP = 0.17 |



HOG (Histogram of Oriented Gradients)

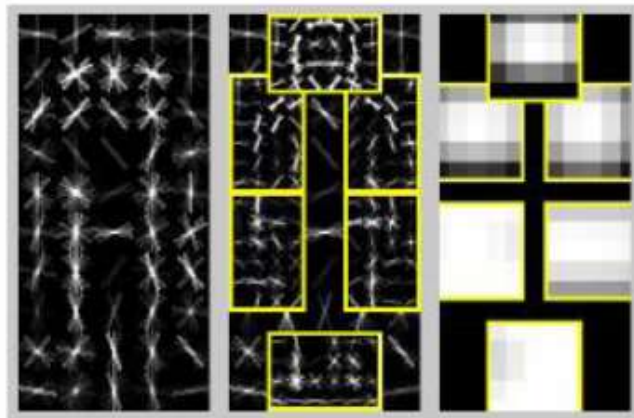
The Evolution of DPM

K=1
6 Parts
[PFF'08]
mAP = 0.21

K=1
no Parts
[DT'05]
mAP = 0.17



Image



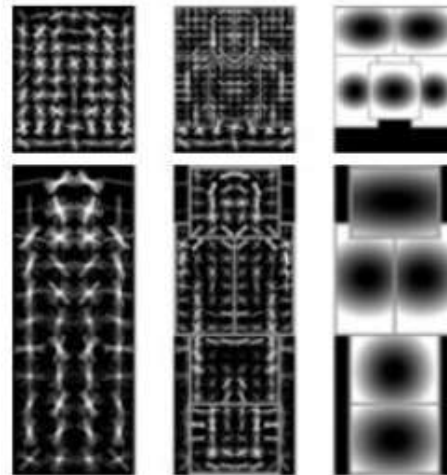
Root filter (Coarse resolution)
Part filters (Fine resolution)
Deformation Models

The Evolution of DPM

| |
|----------------------------|
| K=1 6 Parts [PFF'08] |
| mAP = 0.21 |

| |
|---------------------------------|
| K=2 (ar) 6 Parts [PFF'10] |
| mAP = 0.26 |

| |
|----------------------------|
| K=1 no Parts [DT'05] |
| mAP = 0.17 |



2 DPM (PAMI10)

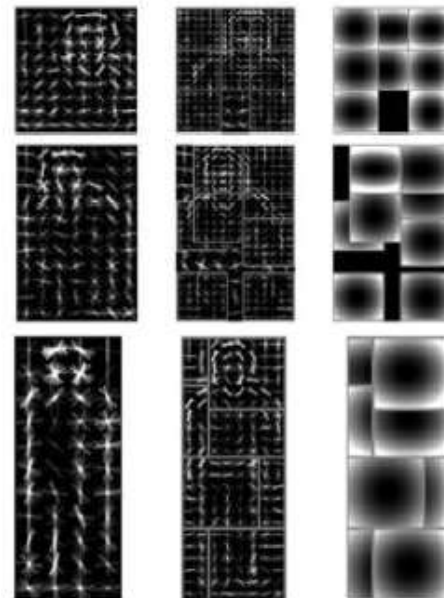
The Evolution of DPM

K=1
6 Parts
[PFF'08]
mAP = 0.21

K=2 (ar)
6 Parts
[PFF'10]
mAP = 0.26

K=6 (ar)
8 Parts
[PFF'11]
mAP = 0.32

K=1
no Parts
[DT'05]
mAP = 0.17



6 DPM (voc-release4)

Santosh's Experiment

K=1
6 Parts
[PFF'08]
mAP = 0.21

K=2 (ar)
6 Parts
[PFF'10]
mAP = 0.26

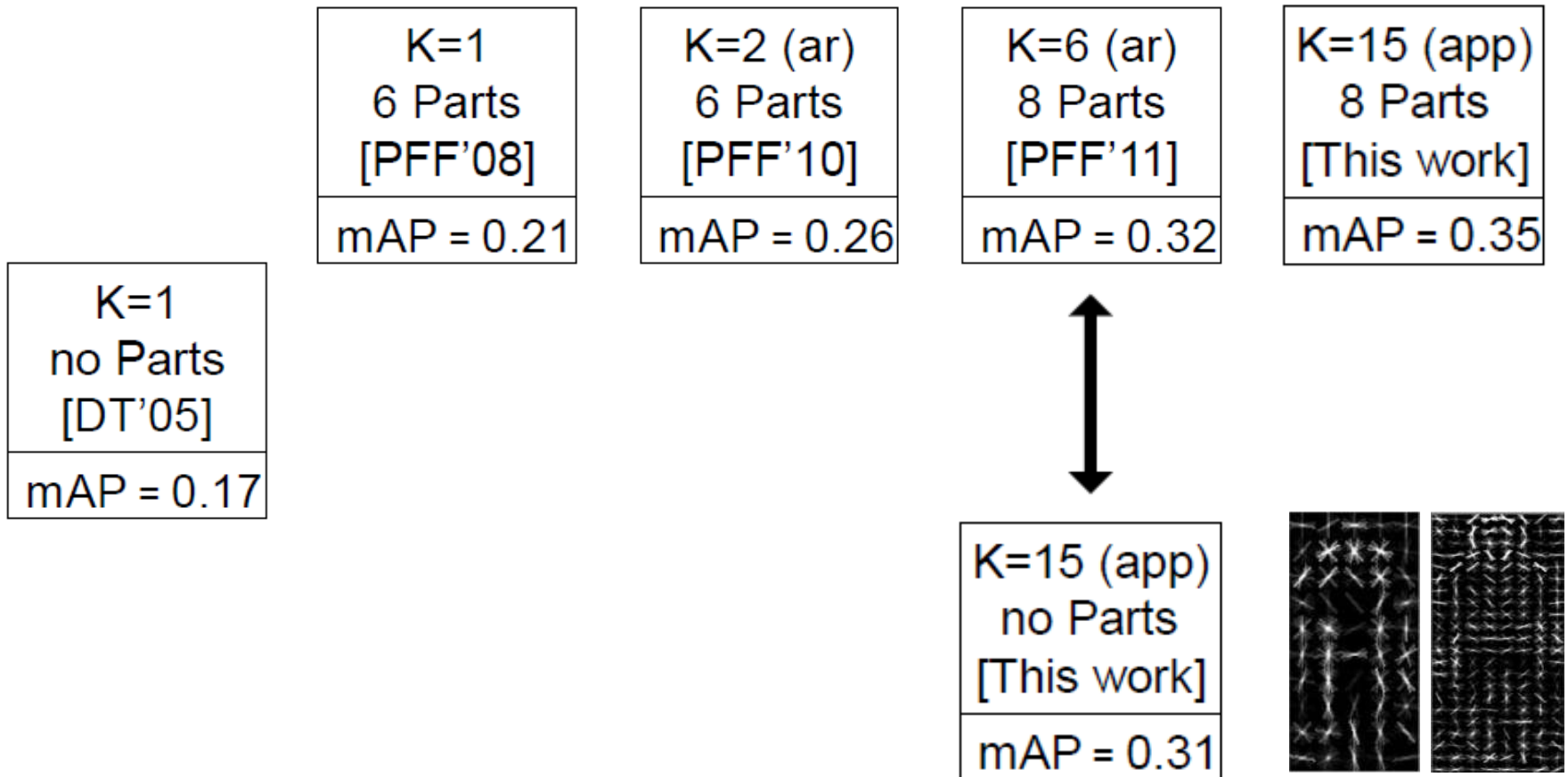
K=6 (ar)
8 Parts
[PFF'11]
mAP = 0.32

K=15 (app)
8 Parts
[This work]
mAP = 0.35

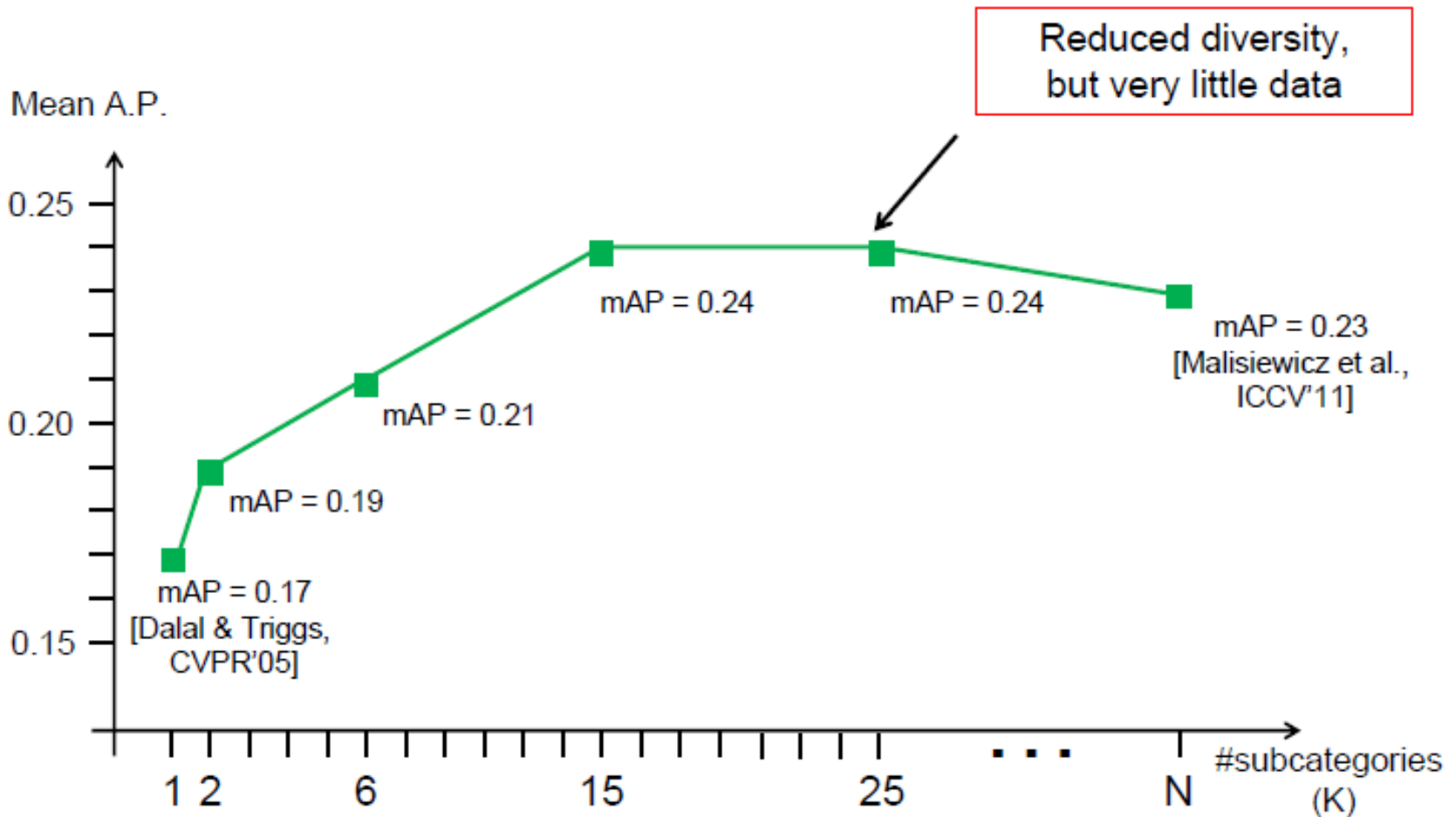
K=1
no Parts
[DT'05]
mAP = 0.17



Santosh's Experiment

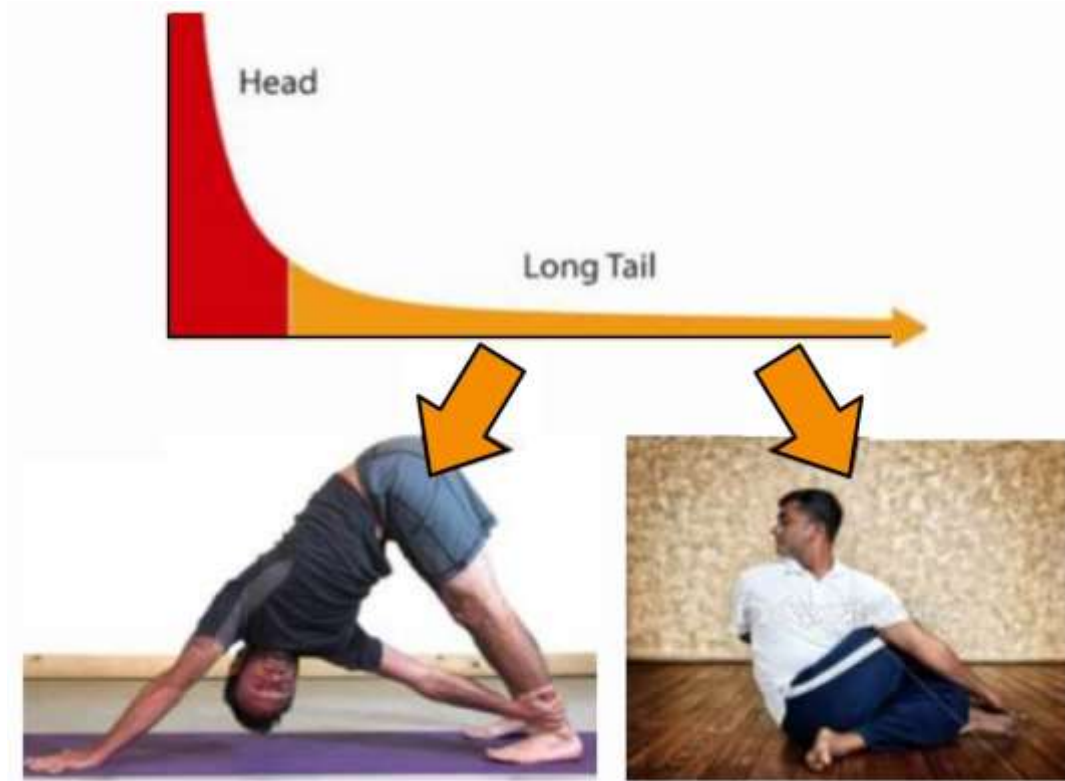


Effect of varying # subcategories



* S. Divvala, A. Effros, M. Hebert "Object Instance Sharing by Enhanced Bounding Box Correspondence", *BMVC'12*.

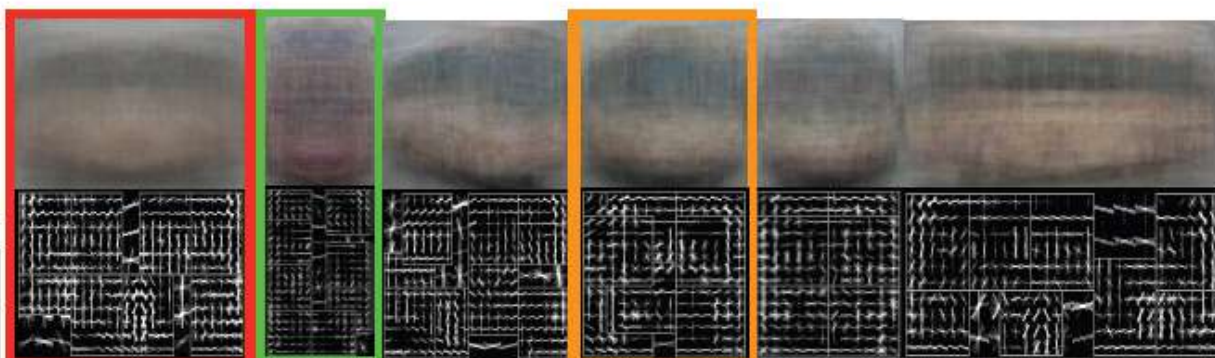
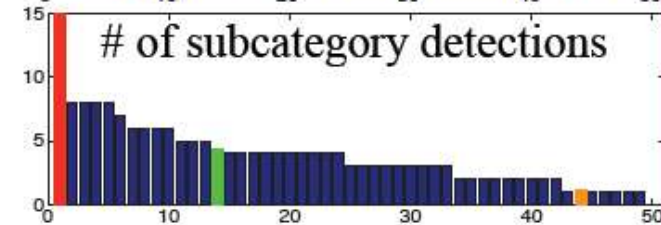
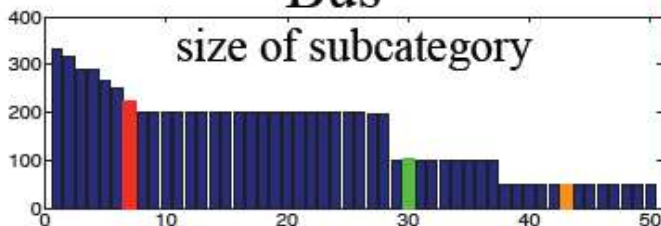
Long-tail distribution



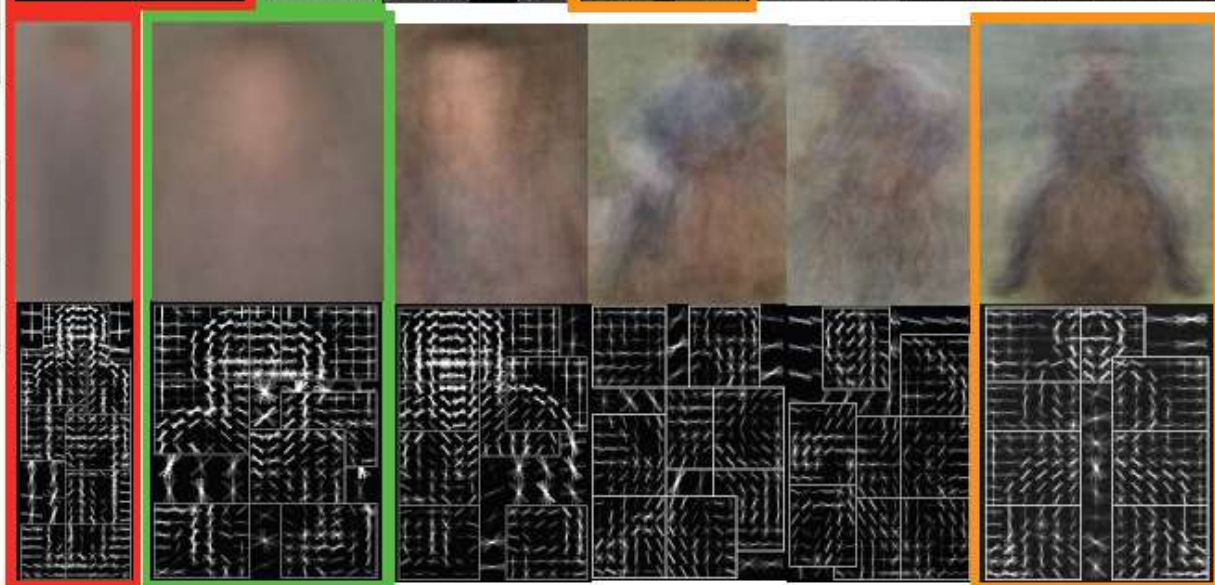
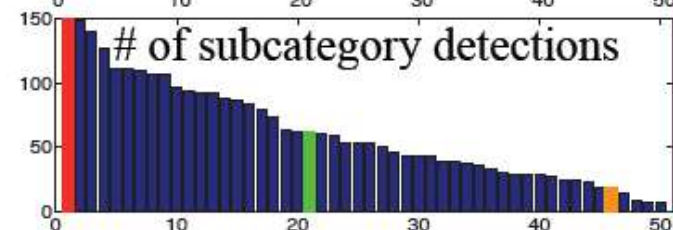
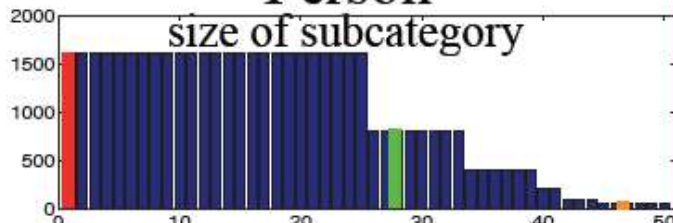
We need lots of templates, have little data of *'yoga twist'* poses

Long-tail distribution

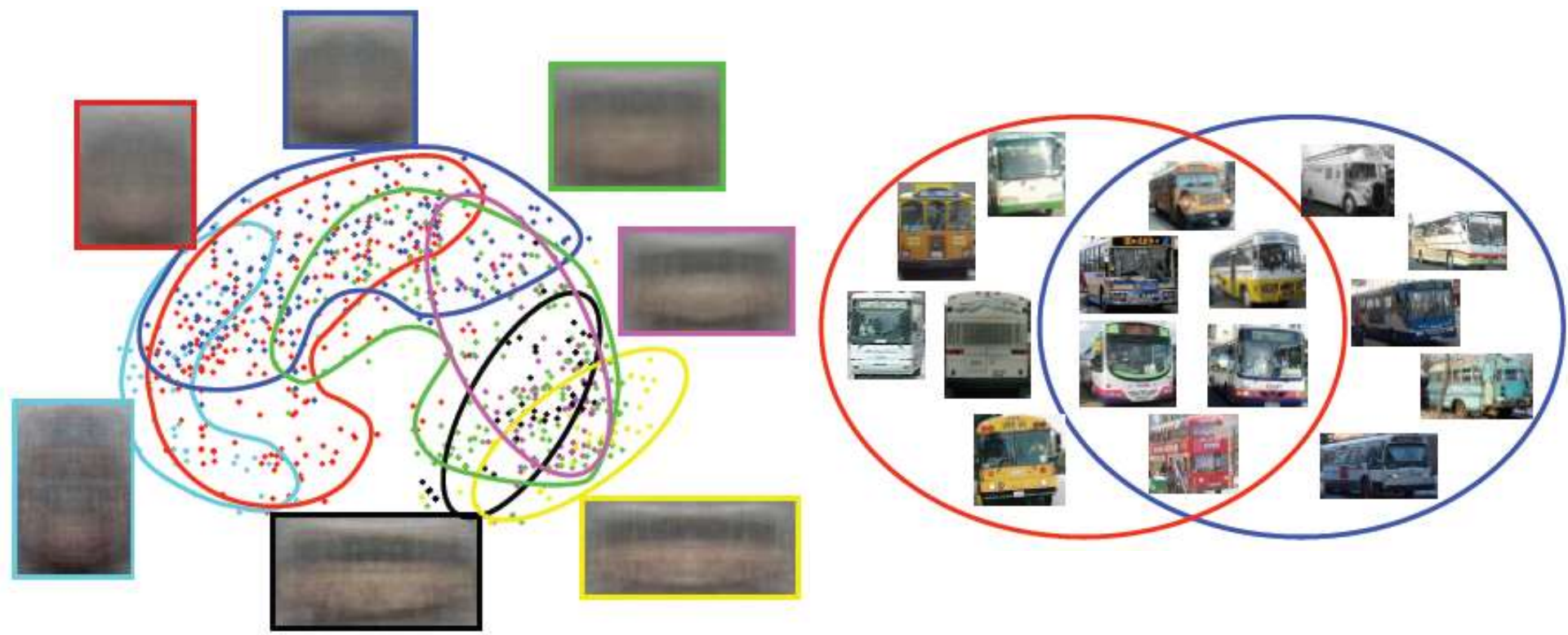
Bus



Person



Instance sharing across subcategories



* X. Zhu, D. Ramanan, "Capturing long-tail distribution of object subcategories", In CVPR 2014.

Web-supervision

Web-supervision

- **GOAL:** Use Internet contents instead of explicit **human supervision**.
- Use internet contents (texts/images) for :
 - *Subcategory discovery:*
 - leveraging vast resources of **online books** to discover the vocabulary of variance.
 - *Enrichment of poor subcategories:*
 - Using gigantic amount of **unlabeled images** on Internet.

The PASCAL Visual Object Classes (VOC) Challenge

Mark Everingham · Luc Van Gool ·
Christopher K. I. Williams · John Winn ·
Andrew Zisserman

Received: 30 July 2008 / Accepted: 16 July 2009 / Published online: 9 September 2009
© Springer Science+Business Media, LLC 2009

Abstract The PASCAL Visual Object Classes (VOC) challenge is a benchmark in visual object category recognition and detection, providing the vision and machine learning communities with a standard dataset of images and annotation, and standard evaluation procedures. Organised annually from 2005 to present, the challenge and its associated dataset has become accepted as *the* benchmark for object detection.

This paper describes the dataset and evaluation procedure. We review the state-of-the-art in evaluated methods for both classification and detection, analyse whether the methods are statistically different, what they are learning from the images (e.g. the object or its context), and what the methods find easy or confuse. The paper concludes with lessons learnt in the three year history of the challenge, and proposes directions for future improvement and extension.

Keywords Database · Benchmark · Object recognition · Object detection

M. Everingham (✉)
University of Leeds, Leeds, UK
e-mail: m.everingham@leeds.ac.uk

L. Van Gool
KU Leuven, Leuven, Belgium

C.K.I. Williams
University of Edinburgh, Edinburgh, UK

J. Winn
Microsoft Research, Cambridge, UK

A. Zisserman
University of Oxford, Oxford, UK

1 Introduction

The PASCAL¹ Visual Object Classes (VOC) Challenge consists of two components: (i) a publicly available *dataset* of images and annotation, together with standardised evaluation software; and (ii) an annual *competition* and workshop. The VOC2007 dataset consists of annotated consumer photographs collected from the flickr² photo-sharing web-site. A new dataset with ground truth annotation has been released each year since 2006. There are two principal challenges: *classification*—"does the image contain any instances of a particular object class?" (where the object classes include cars, people, dogs, etc.), and *detection*—"where are the instances of a particular object class in the image (if any)?" In addition, there are two subsidiary challenges ("tasters") on pixel-level segmentation—assign each pixel a class label, and "person layout"—localise the head, hands and feet of people in the image. The challenges are issued with deadlines each year, and a workshop held to compare and discuss that year's results and methods. The datasets and associated annotation and software are subsequently released and available for use at any time.

The objectives of the VOC challenge are twofold: first to provide challenging images and high quality annotation, together with a standard evaluation methodology—a "plug and play" training and testing harness so that performance of algorithms can be compared (the dataset component); and second to measure the state of the art each year (the competition component).

¹PASCAL stands for pattern analysis, statistical modelling and computational learning. It is an EU Network of Excellence funded under the IST Programme of the European Union.

²<http://www.flickr.com/>

Table 1 Queries used to retrieve images from flickr. Words in bold show the “targeted” class. Note that the query terms are quite general—including the class name, synonyms and scenes or situations where the class is likely to occur

- **horse**, gallop, jump, buck, equine, foal, cavalry, saddle, canter, buggy, mare, neigh, dressage, trial, racehorse, steeplechase, thoroughbred, cart, equestrian, paddock, stable, farrier
- **motorbike**, motorcycle, minibike, moped, dirt, pillion, biker, trials, motorcycling, motorcyclist, engine, motocross, scramble, sidecar, scooter, trail
- **person**, people, family, father, mother, brother, sister, aunt, uncle, grandmother, grandma, grandfather, grandpa, grandson, granddaughter, niece, nephew, cousin
- **sheep**, ram, fold, fleece, shear, baa, bleat, lamb, ewe, wool, flock
- **sofa**, chesterfield, settee, divan, couch, bolster
- **table**, dining, cafe, restaurant, kitchen, banquet, party, meal
- **potted plant**, pot plant, plant, patio, windowsill, window sill, yard, greenhouse, glass house, basket, cutting, pot, cooking, grow
- **train**, express, locomotive, freight, commuter, platform, subway, underground, steam, railway, railroad, rail, tube, underground, track, carriage, coach, metro, sleeper, railcar, buffet, cabin, level crossing
- **tv/monitor**, television, plasma, flatscreen, flat screen, lcd, crt, watching, dvd, desktop, computer, computer monitor, PC, console, game

Query expansion

User
Input
"horse"



Search
Google books
Ngrams



- horse
- grazing horse
- jumping horse
- rolling horse
- reining horse
- sledge horse
- fighting horse
- crazy horse
- horse ears
- last horse
- sleigh horse
- eating horse
- horse head
- ...
- +20K variations

- horse
 - grazing horse
 - jumping horse
 - rolling horse
 - reining horse
 - sledge horse
 - fighting horse
 - crazy horse
 - horse ears
 - last horse
 - sleigh horse
 - eating horse
 - horse head
 - ...
- +20K variations

- horse
 - grazing horse
 - jumping horse
 - rolling horse
 - reining horse
 - sledge horse
 - fighting horse
 - crazy horse
 - horse ears
 - last horse
 - sleigh horse
 - eating horse
 - horse head
 - ...
- +20K variations

Pruning non visual ngrams

- horse
- grazing horse
- jumping horse
- rolling horse
- reining horse
- sledge horse
- fighting horse
- crazy horse
- horse ears
- last horse
- sleigh horse
- eating horse
- horse head
- ...
- +20K variations



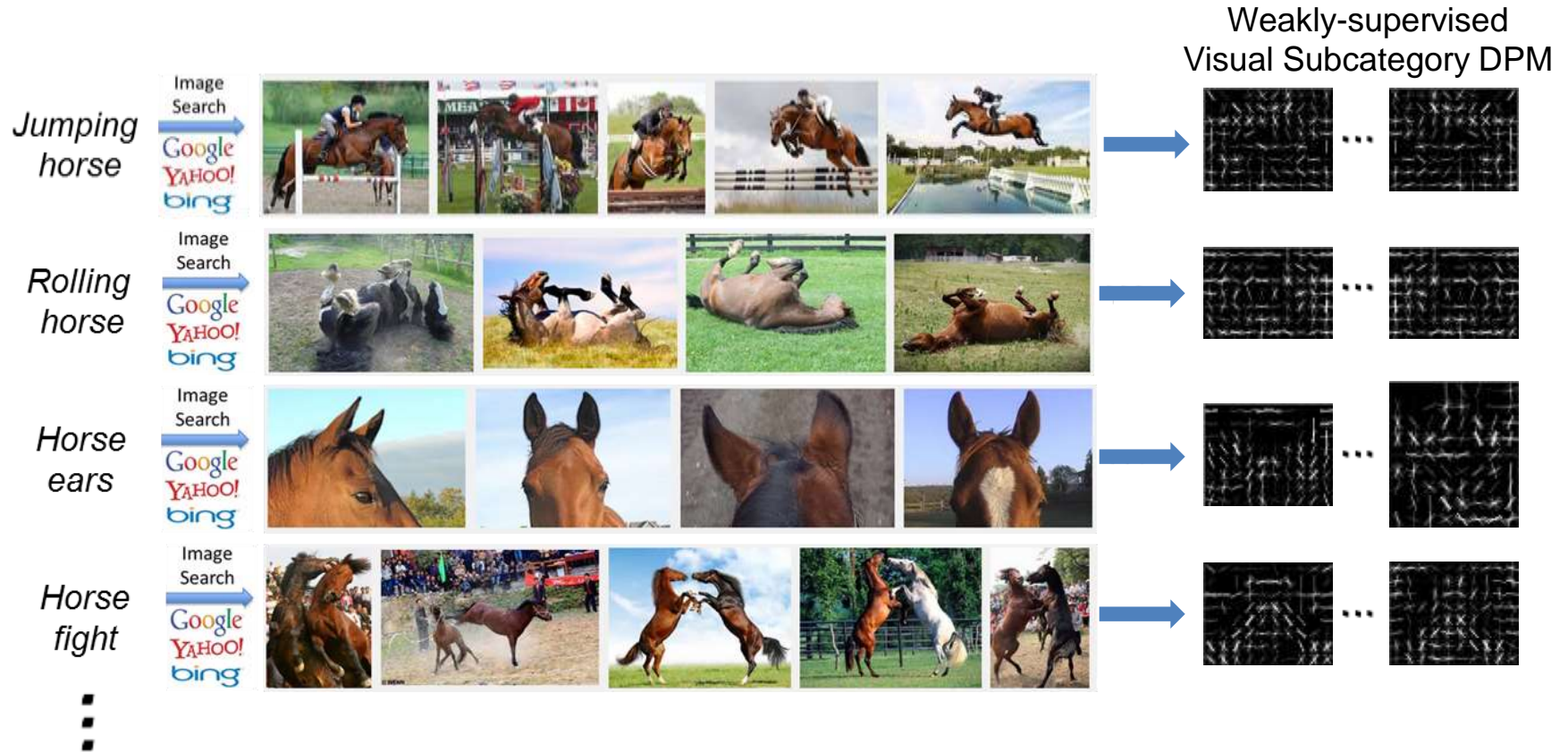
Identifying Visual Synonyms

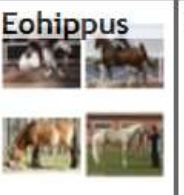
- horse
- grazing horse
- jumping horse
- rolling horse
- reining horse
- sledge horse
- fighting horse
- ~~• crazy horse~~
- horse ears
- ~~• last horse~~
- sleigh horse
- eating horse
- horse head
- ...

+20K variations



Taming Intra-class Variance



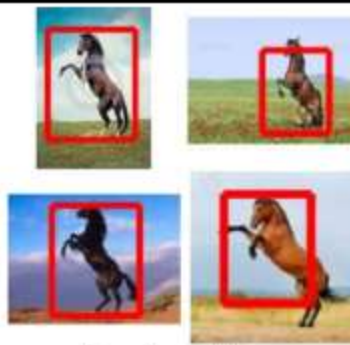




Swimming Horse



Rolling Horse



Rearing Horse



Fighting Horse



Bucking Horse



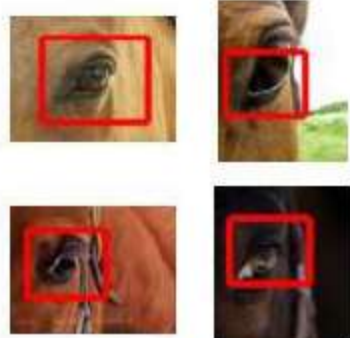
Reining Horse



Barrel Horse



Horse Tram



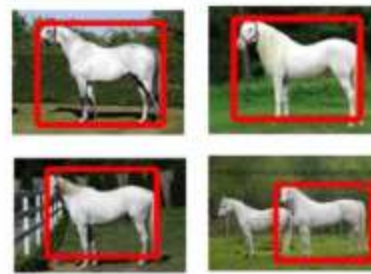
Horse Eye



Front Horse



Bridled Horse

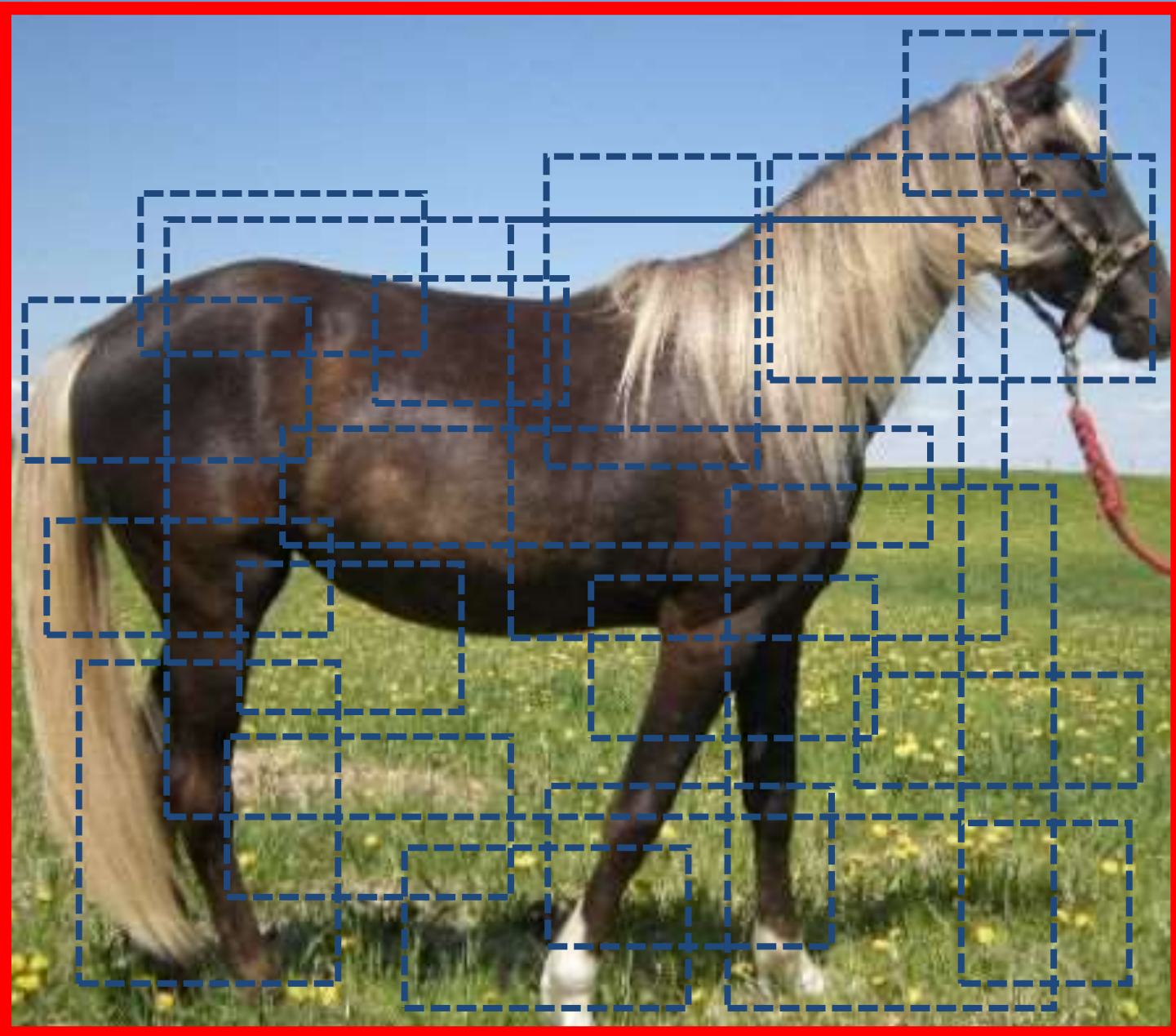


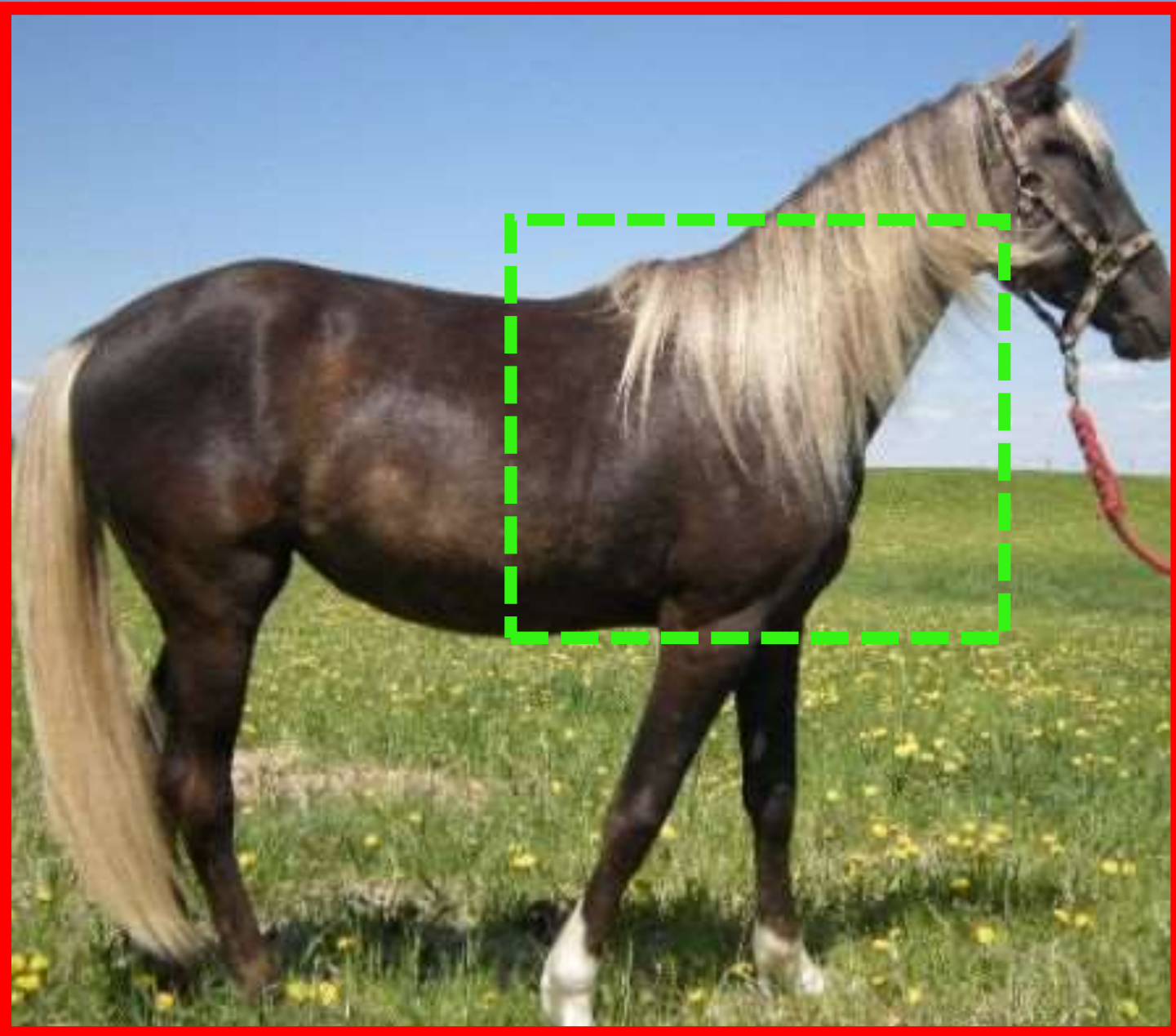
Jennet horse

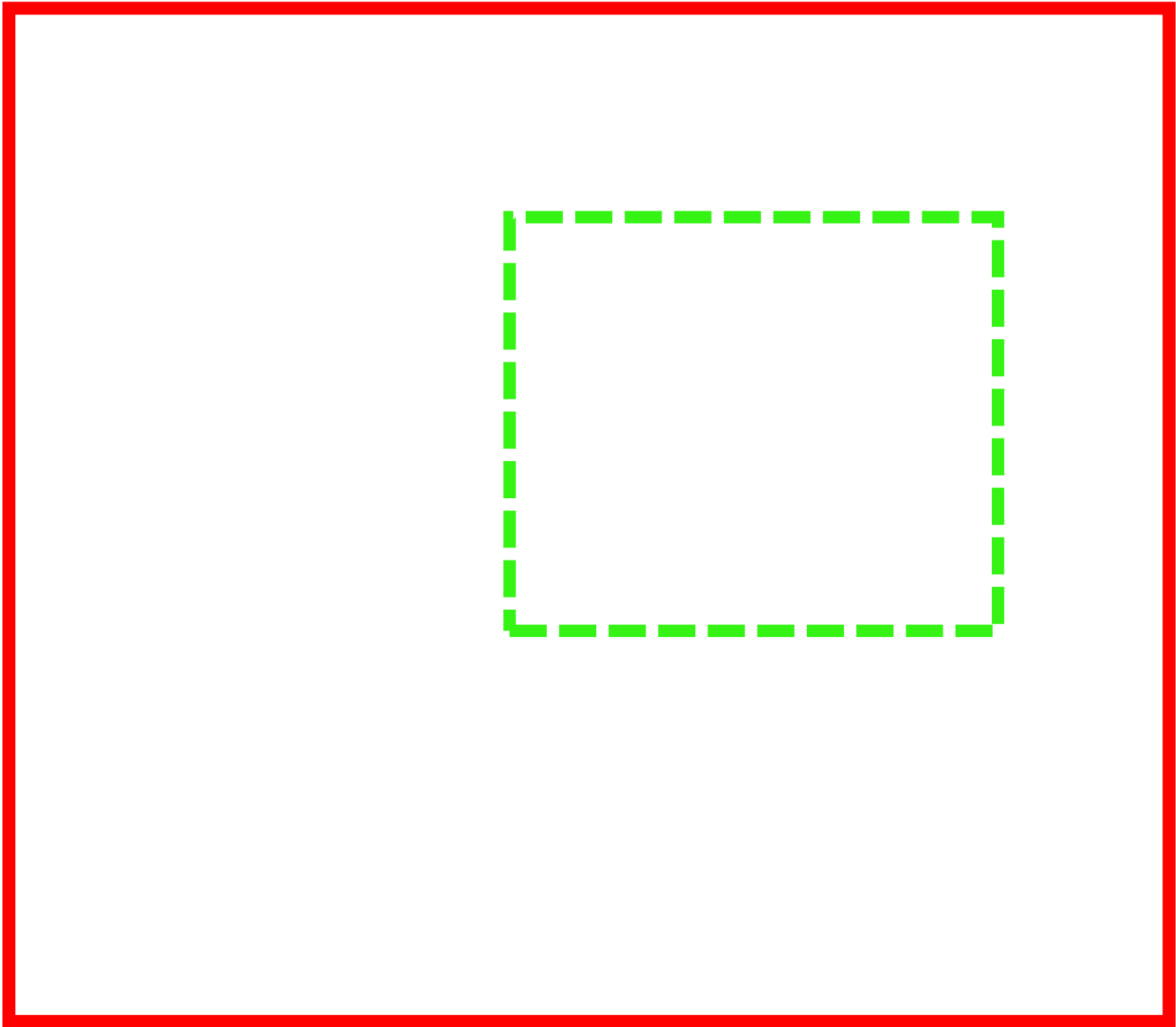
Webly-supervised discriminative patch

- Can we go **inside the box** and find the discriminative patch?
- **Subcategory-aware** discriminative patches.
- **Fixed-position** patches.

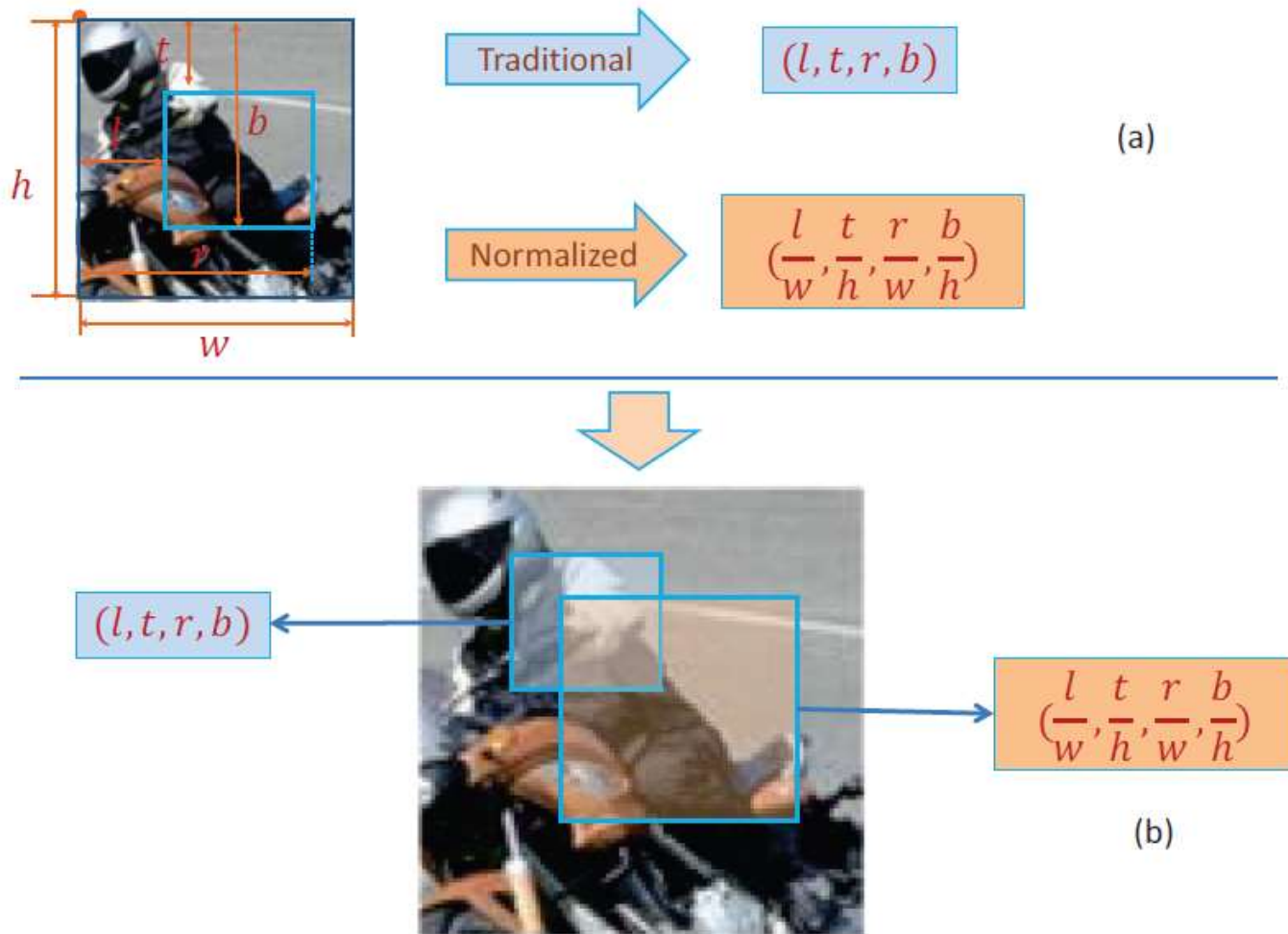


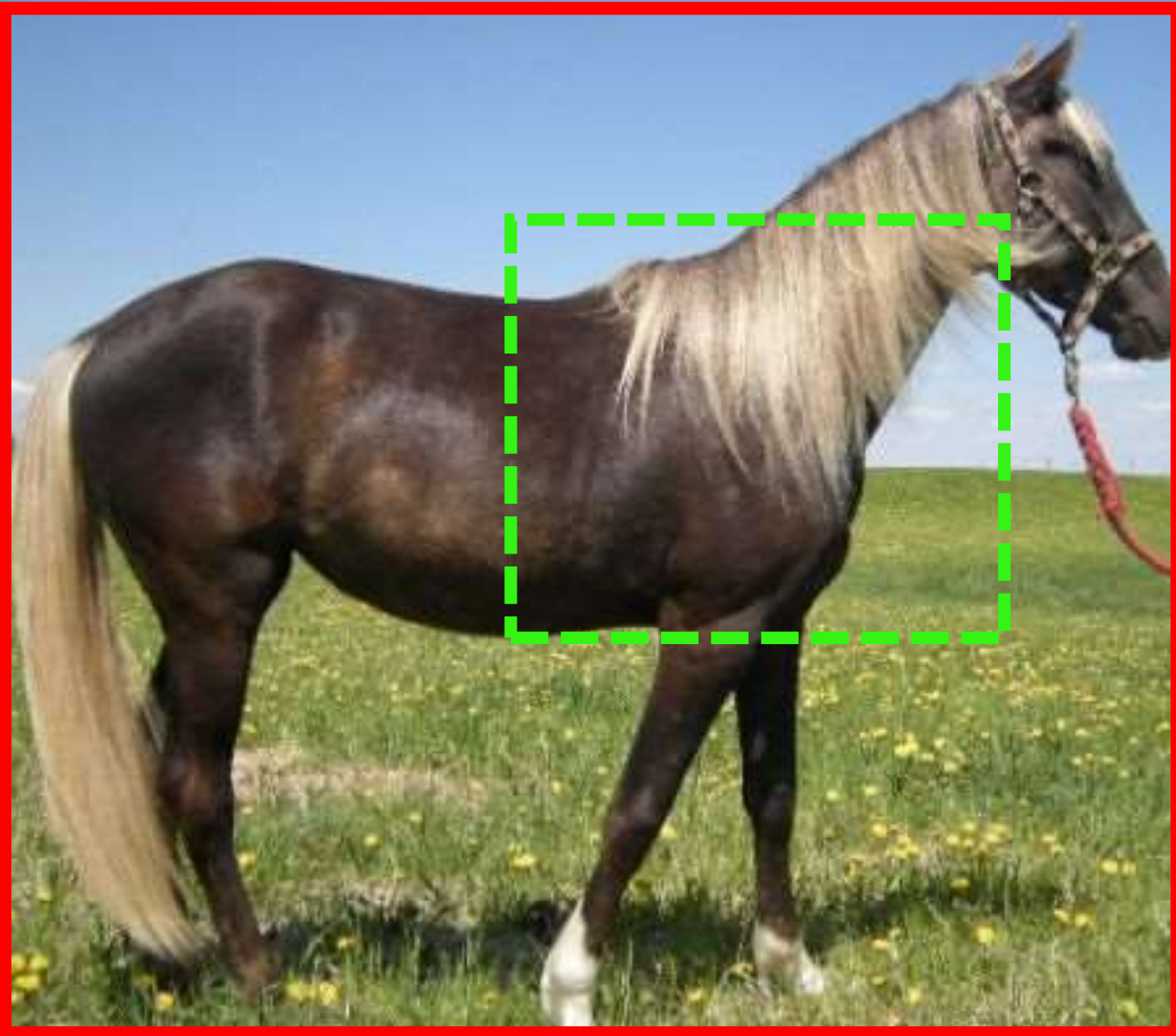






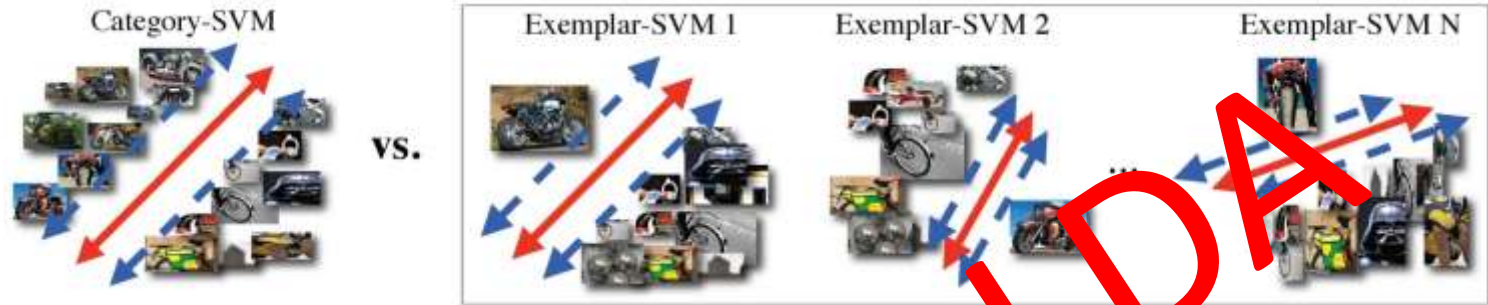
Relative Normalized Position



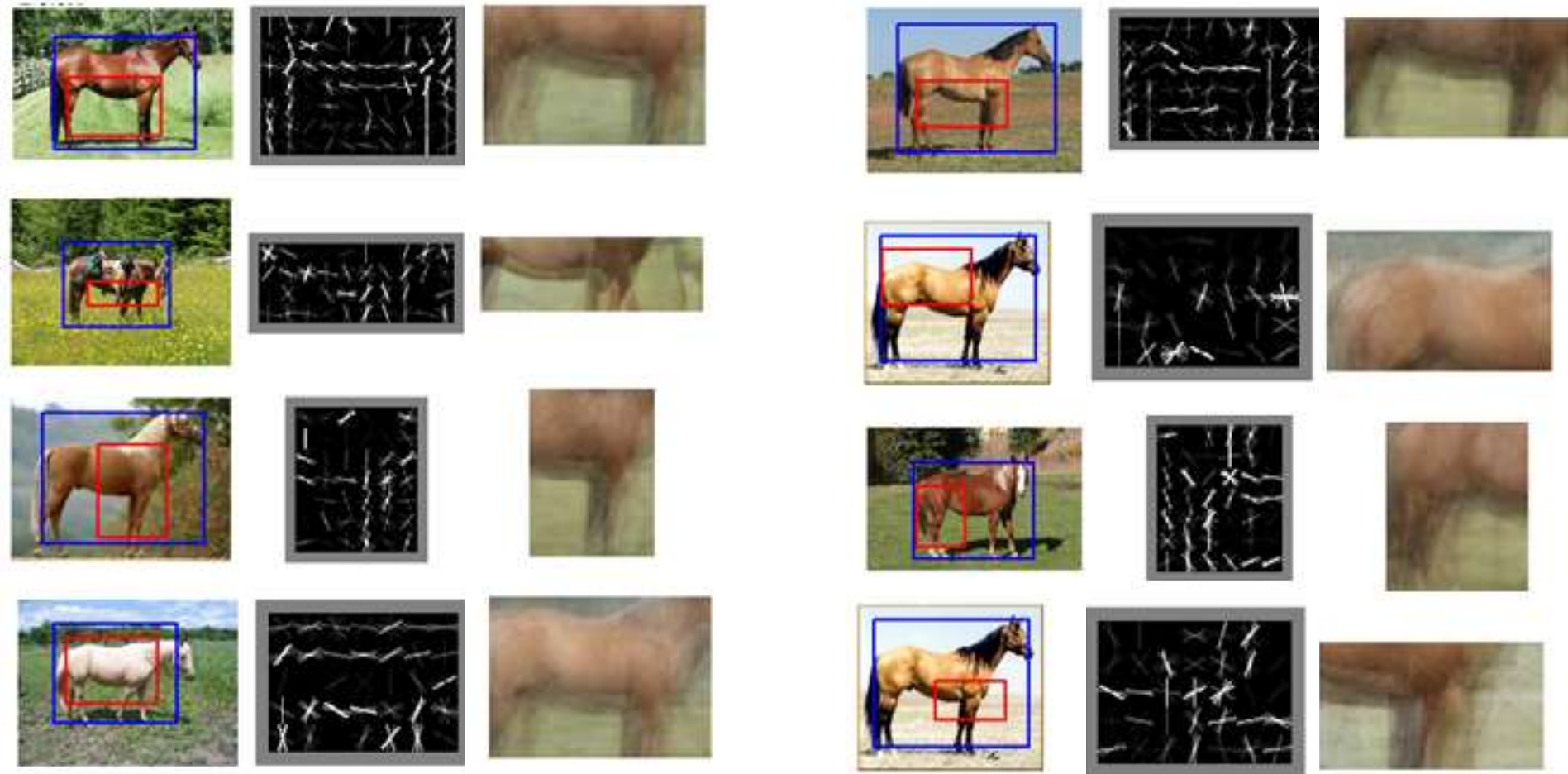




Train Exemplar-SVM

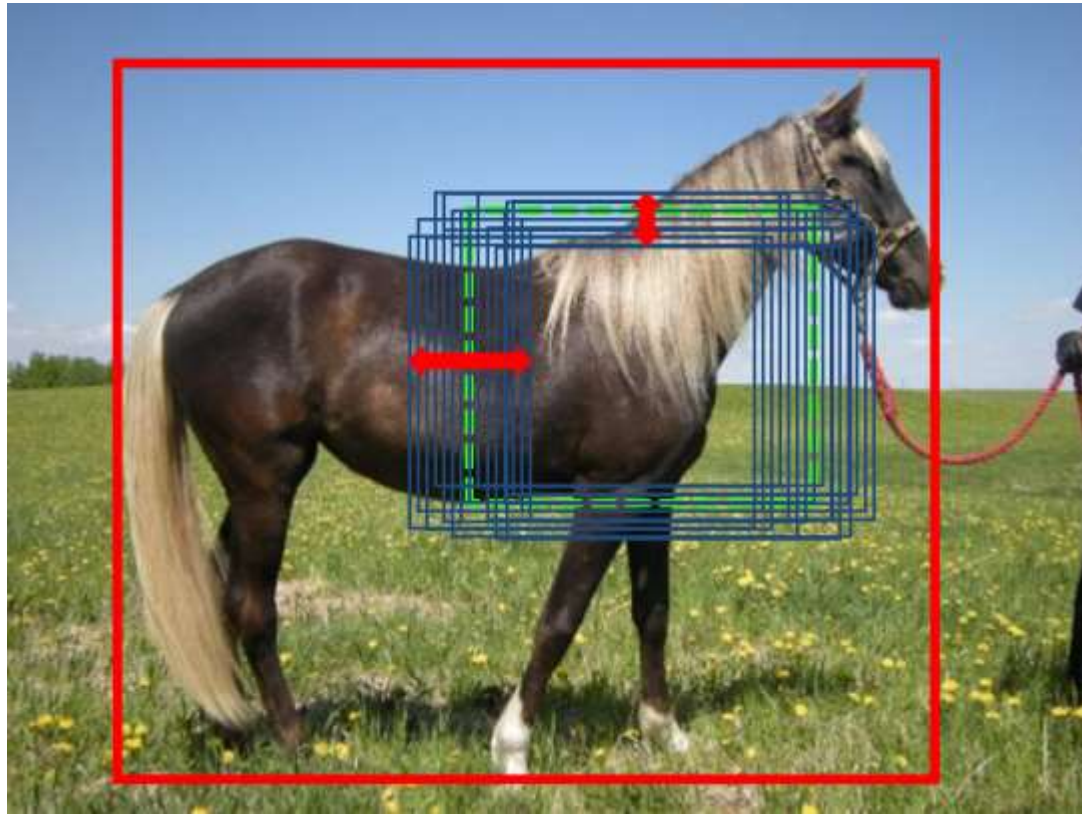


Initial Patch Models



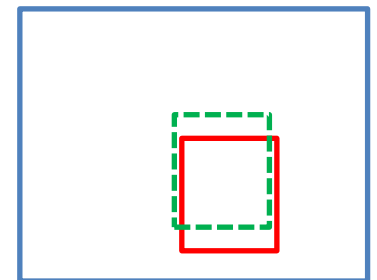
Patch Deformability

- Patch should NOT be fully fixed-position
 - Use NMS to find deformation of the patch



Patch Selection

- What are **good** patch?
 - Appearance consistency score
 - Repetitive visual pattern
 - Confidence score of E-LDA patch model
 - Spatial consistency score
 - Spatially consistent
 - Patch activation



$$Activation(a, b) = \frac{(a \cap b)}{(a \cup b)}$$

Patch Selection

- Representativeness criteria
 - Intra-subcategory consistency

$$rep(p, I_s) = \frac{1}{|I_s|} \sum_{i=1}^{|I_s|} score(x_i, p)$$

- Discrimination criteria
 - Inter-Category discriminativity
 - Normalized median rank on a mixed set of subcategory images and a huge PASCAL negative set.

$$disc(p, I_s, \bar{I}) = \frac{median(rank(p, I_s, \bar{I}))}{|I_s|}$$

$$rank(p, I_s, \bar{I}) : \mathbb{R}^{|I_s \cup \bar{I}|} \mapsto \mathbb{N}^{|I_s|}$$

RP=0.94



Score:1.15



Score:1.68



Score:1.9



Score:1.87



Score:1.73



Score:2.68



Score:1.83



Score:1.56



Score:1.56



Score:1.32



Score:2.37



Score:2.37



Score:1.2



Score:1.23



Score:1.14



Score:1.14



Score:1.82



Score:1.11



Score:1.68



Score:1



RP=1.75



Score:4.14



Score:1.75



Score:1.3



Score:1.89



Score:1.96



Score:1.75



Score:1.7



Score:1.64



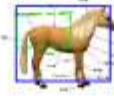
Score:1.82



Score:1.81



Score:1.45



Score:1.33



Score:1.58



Score:1.34



Score:1.28



Score:1.41



RP=0.94



Score:1.46



Score:1.12



Score:1.96



Score:1.87



Score:1.74



Score:2.56



Score:1.55



Score:1.5



Score:1.41



Score:1.4



Score:1.33



Score:1.75



Score:1.35



Score:2.35



Score:1.23



Score:1.22



Score:1.25



Score:1.25



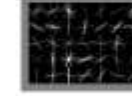
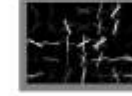
Score:1.11



Score:1.17



RP=1.75



Score:1.33



Score:1.96



Score:1.62



Score:1.47



Score:1.25



Score:1.34



Score:1.12



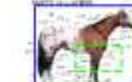
Score:1.06



Score:1.95

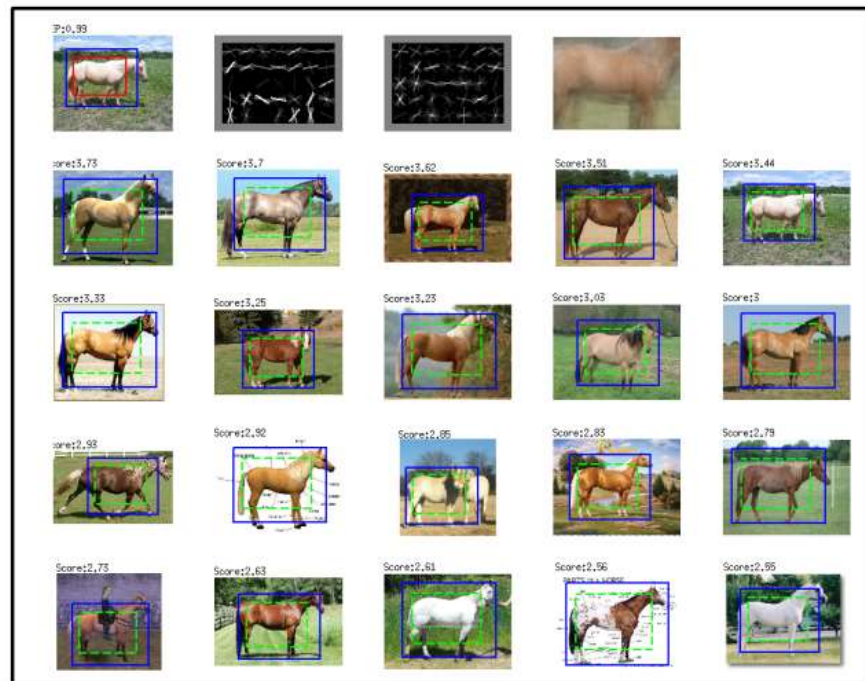
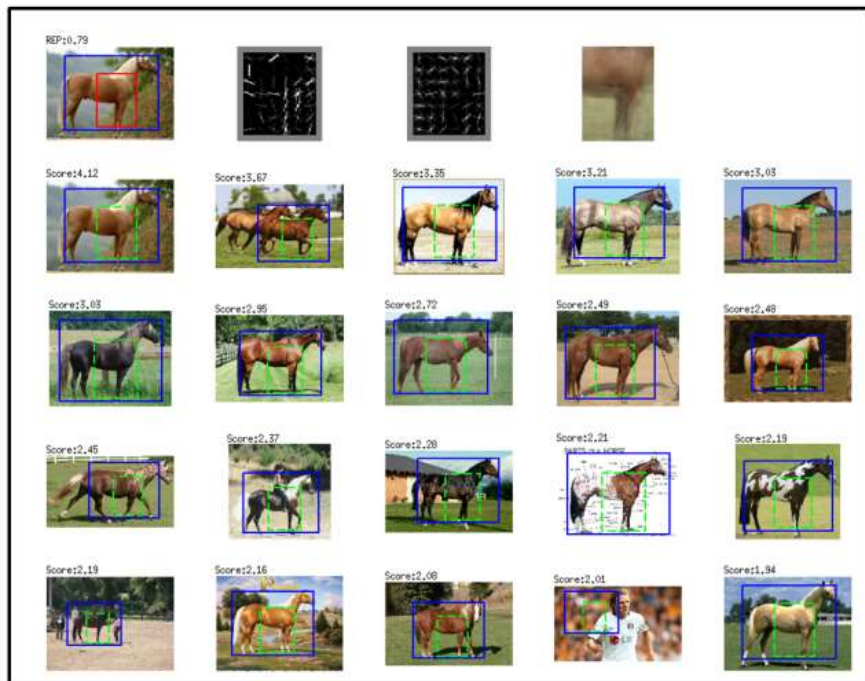


Score:1.88



Score:1.44

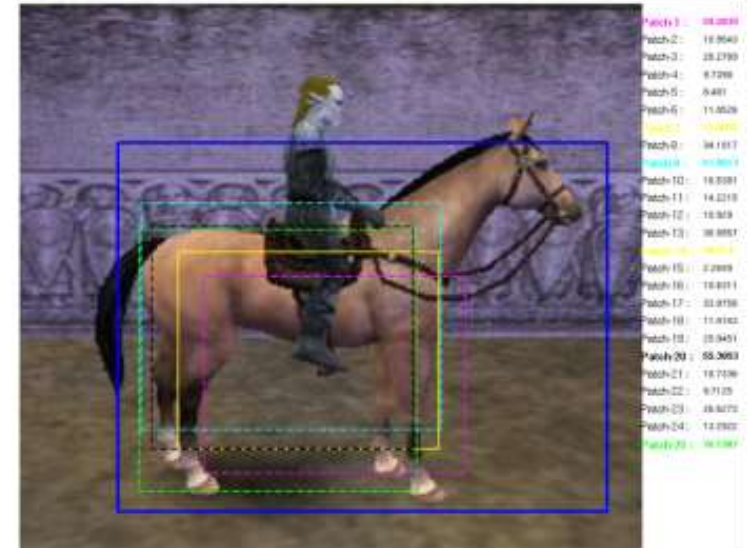
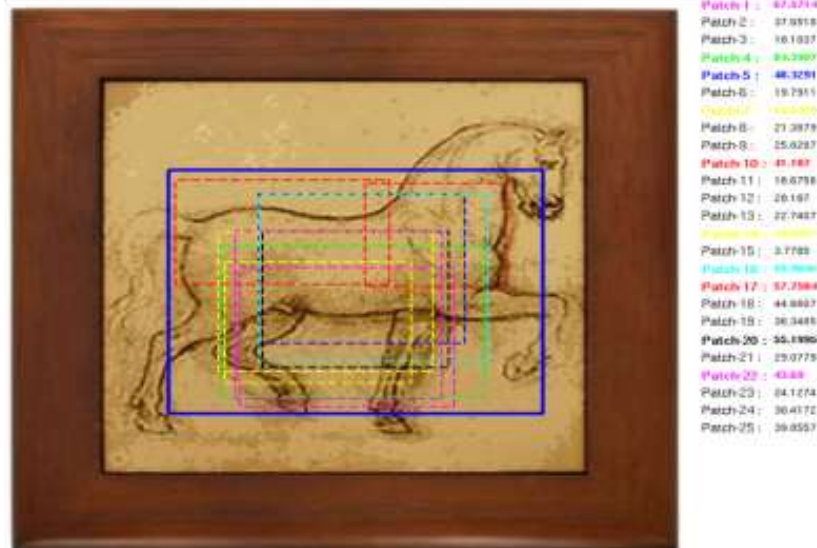
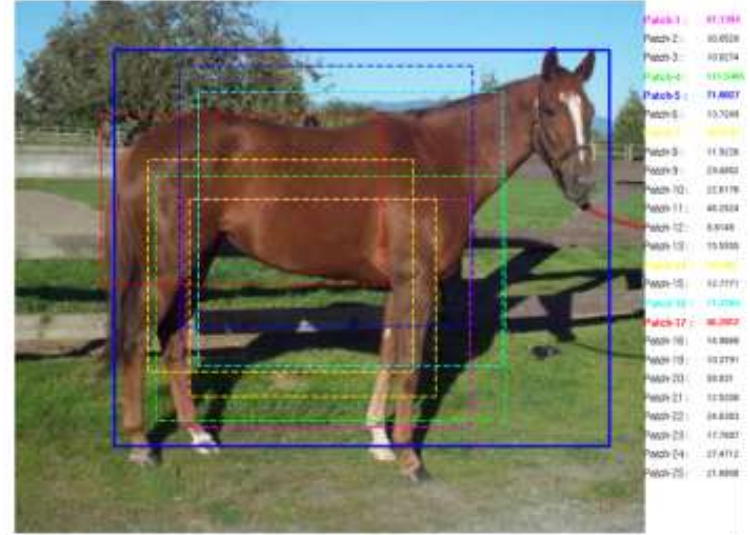
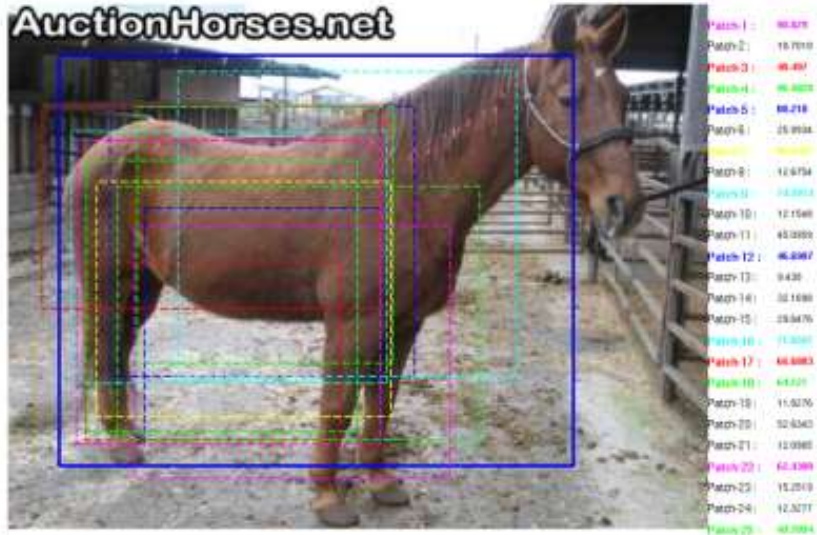




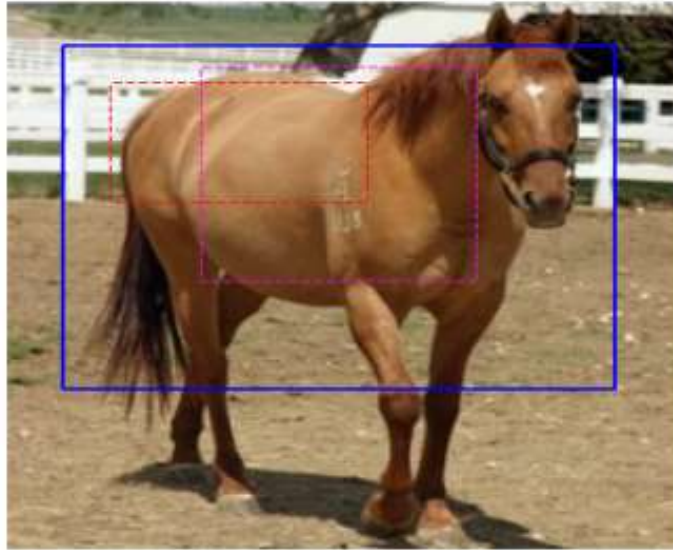
Patch re-training

- **Why?**
 - E-LDA are **shallow models**
 - LDA and SVM models are **highly correlated**
- **Patch retraining:**
 - Train initial patch models with LDA
 - Patch selection
 - Example selection
 - Then re-train the expensive models only for the selected Patches using Latent-LDA

Select good example



Prune noisy images



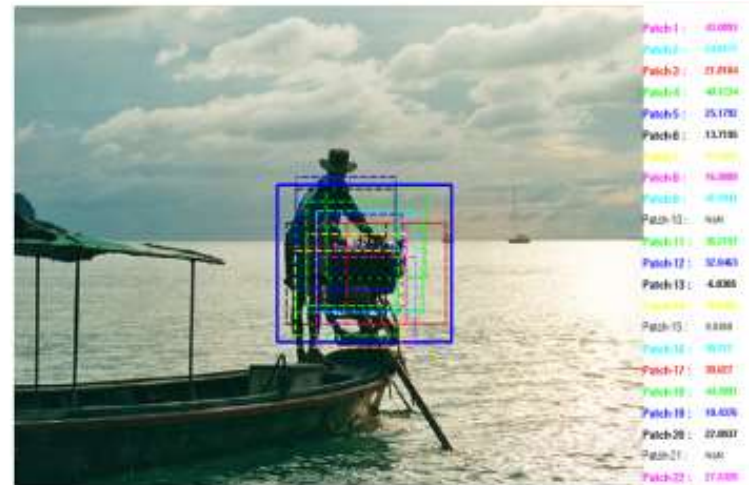
Patch 1: 11.6876
 Patch 2: 29.047
 Patch 3: 25.3779
 Patch 4: 13.2609
 Patch 5: 24.4323
 Patch 6: 11.7513
 Patch 7: 43.7264
 Patch 8: 64.7664
 Patch 9: 7.866
 Patch 10: 26.8171
 Patch 11: 10.3114
 Patch 12: 4.6674
 Patch 13: 18.2499
 Patch 14: 18.6689
 Patch 15: 9.9081
 Patch 16: 9.0933
 Patch 17: 38.7758
 Patch 18: 34.2654
 Patch 19: 27.2745
 Patch 20: 20.2684
 Patch 21: 24.2536
 Patch 22: 35.7097
 Patch 23: 22.2019
 Patch 24: 6.6987
 Patch 25: 6.1479



Patch 1: NaN
 Patch 2: NaN
 Patch 3: NaN
 Patch 4: NaN
 Patch 5: NaN
 Patch 6: NaN
 Patch 7: NaN
 Patch 8: NaN
 Patch 9: NaN
 Patch 10: NaN
 Patch 11: NaN
 Patch 12: NaN
 Patch 13: NaN
 Patch 14: NaN
 Patch 15: 17.476
 Patch 16: NaN
 Patch 17: NaN
 Patch 18: NaN
 Patch 19: NaN
 Patch 20: NaN
 Patch 21: NaN
 Patch 22: NaN
 Patch 23: NaN
 Patch 24: NaN
 Patch 25: 1.4735

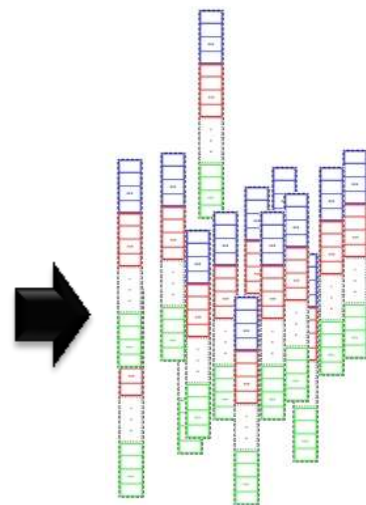


Patch 1: 64.7129
 Patch 2: 25.8256
 Patch 3: 25.2322
 Patch 4: 23.2803
 Patch 5: 34.0489
 Patch 6: 24.0245
 Patch 7: 24.0245
 Patch 8: 23.087
 Patch 9: 25.1173
 Patch 10: 77.4896
 Patch 11: 64.6266
 Patch 12: 36.8424
 Patch 13: 46.626
 Patch 14: 24.0245
 Patch 15: 4.8102
 Patch 16: 32.969
 Patch 17: 23.5829
 Patch 18: 32.969
 Patch 19: 39.0607
 Patch 20: 43.4472
 Patch 21: 24.0245
 Patch 22: 13.8416
 Patch 23: 14.6434
 Patch 24: 22.2712
 Patch 25: 28.4341



Patch 1: 40.0893
 Patch 2: 44.0155
 Patch 3: 27.8144
 Patch 4: 40.1294
 Patch 5: 25.1782
 Patch 6: 13.1746
 Patch 7: 24.0245
 Patch 8: 54.2689
 Patch 9: 47.0191
 Patch 10: NaN
 Patch 11: 30.2757
 Patch 12: 32.0463
 Patch 13: 4.8365
 Patch 14: 24.0245
 Patch 15: 4.8365
 Patch 16: 30.2757
 Patch 17: 38.827
 Patch 18: 44.0155
 Patch 19: 18.4026
 Patch 20: 27.8144
 Patch 21: NaN
 Patch 22: 27.8144
 Patch 23: NaN
 Patch 24: 4.8365
 Patch 25: 40.1294

Patch Calibration



Context

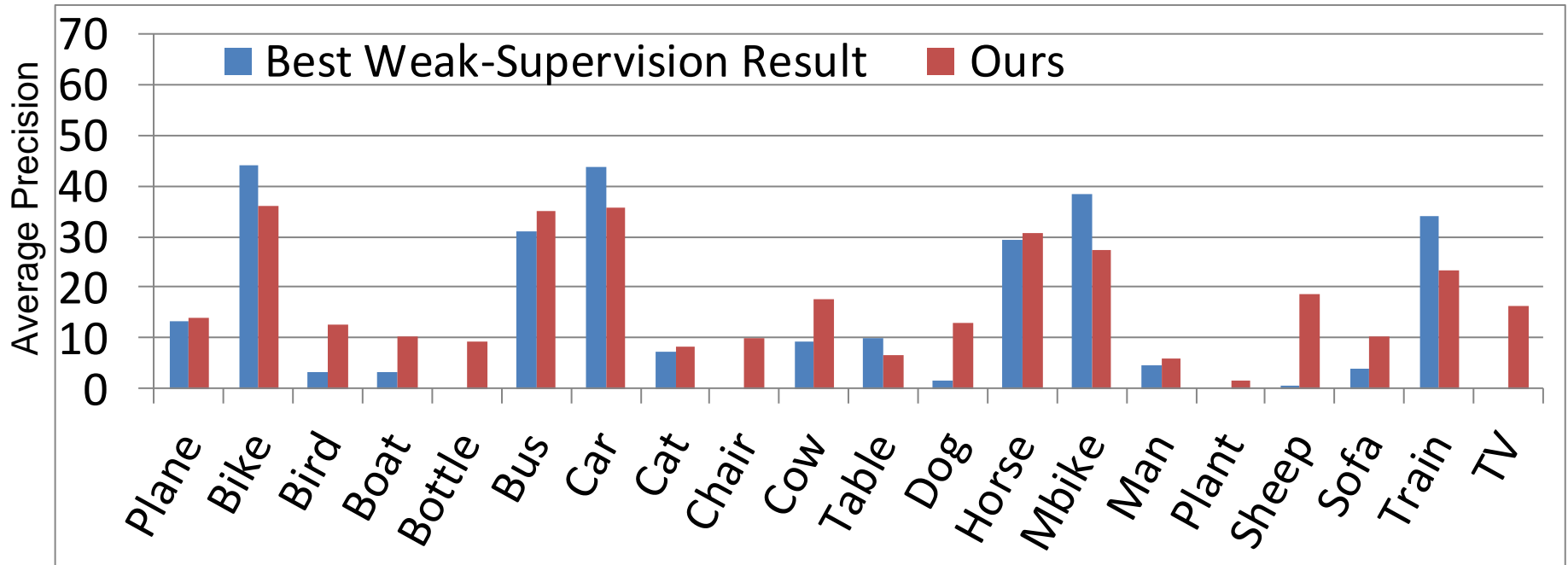


Patch-based detection



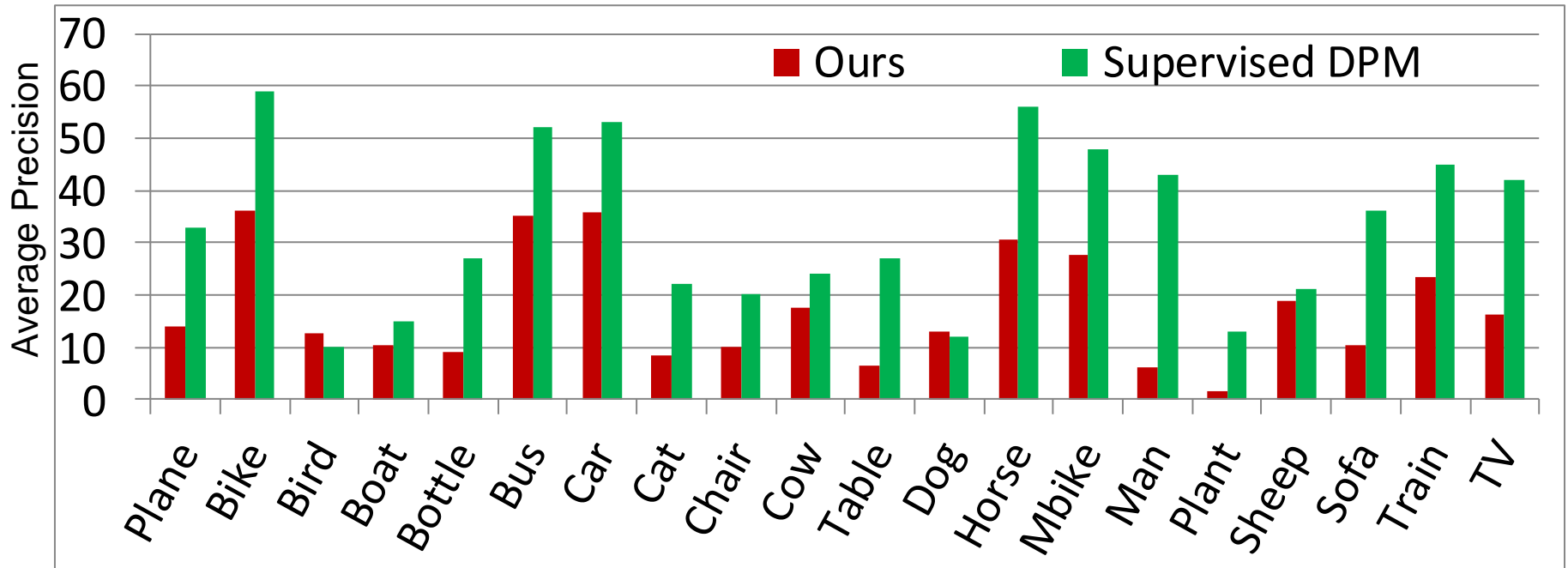
Results

PASCAL VOC 2007 Object Detection



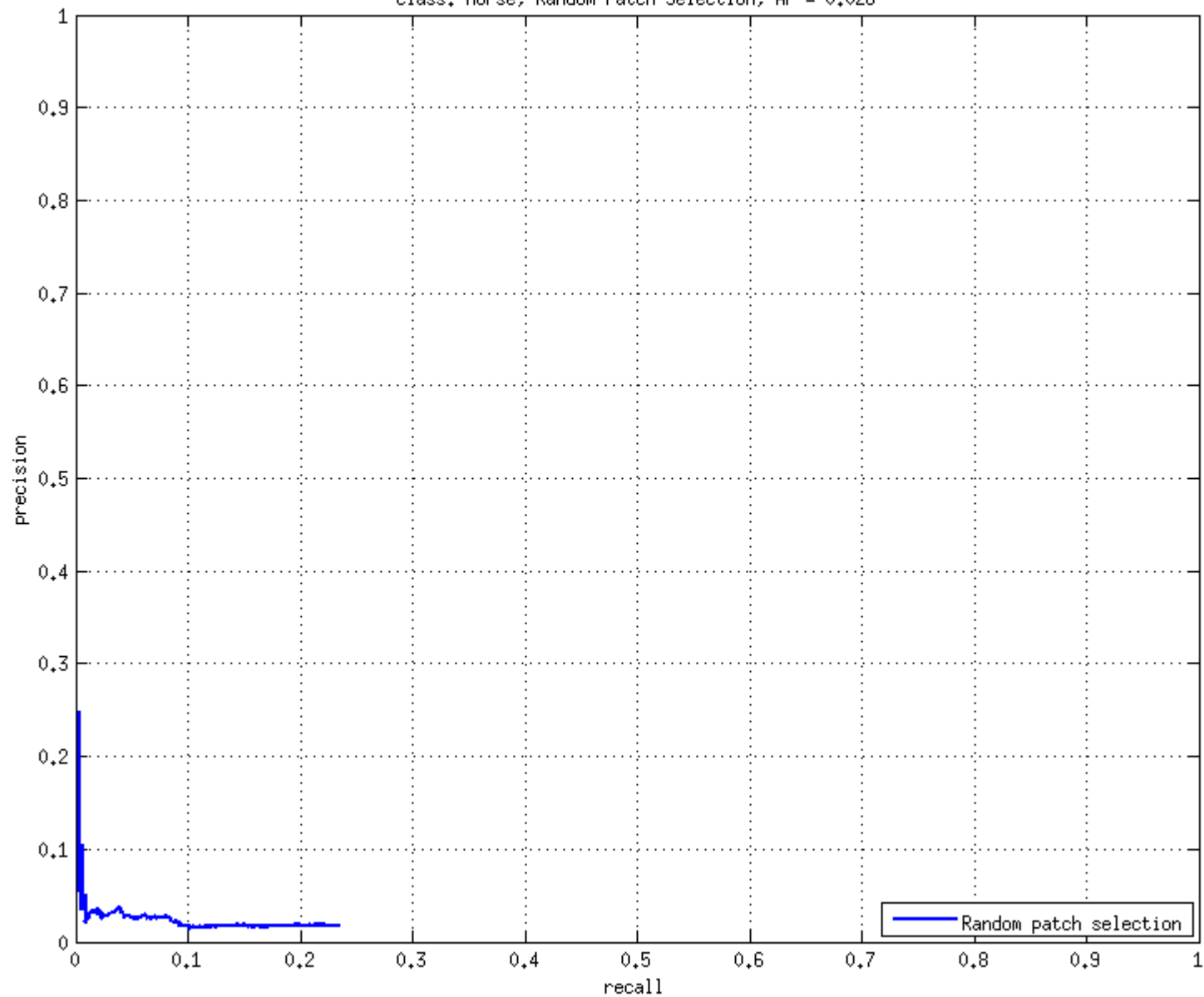
- Holistic model beats the previous best on 13 classes
- Previous best [Siva_ICCV11, Prest_CVPR12] uses weak human-supervision i.e., image/video labels, and Objectness

PASCAL VOC 2007 Object Detection

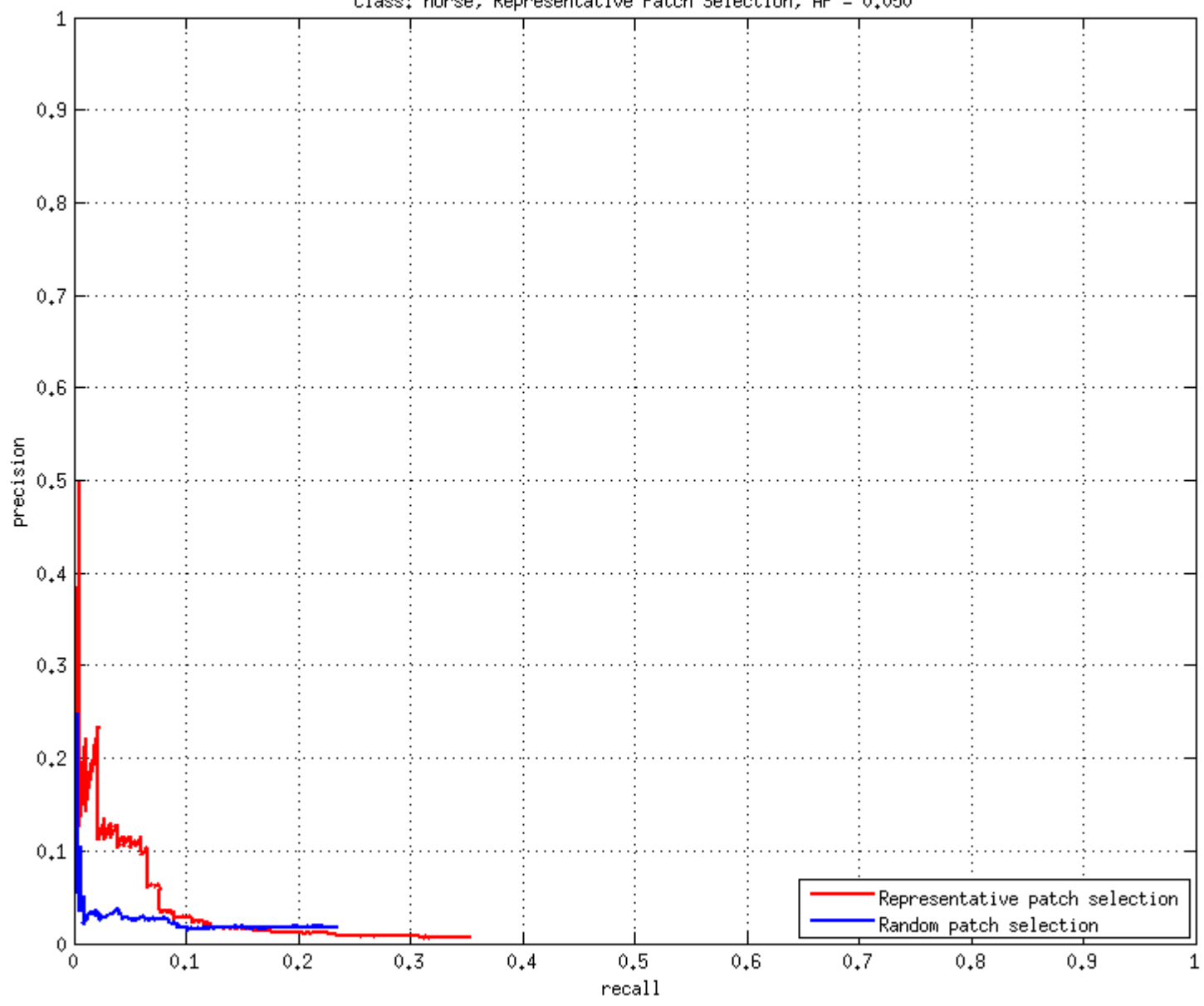


Holistic webly-supervised model is almost on par with supervised DPM on 4 classes

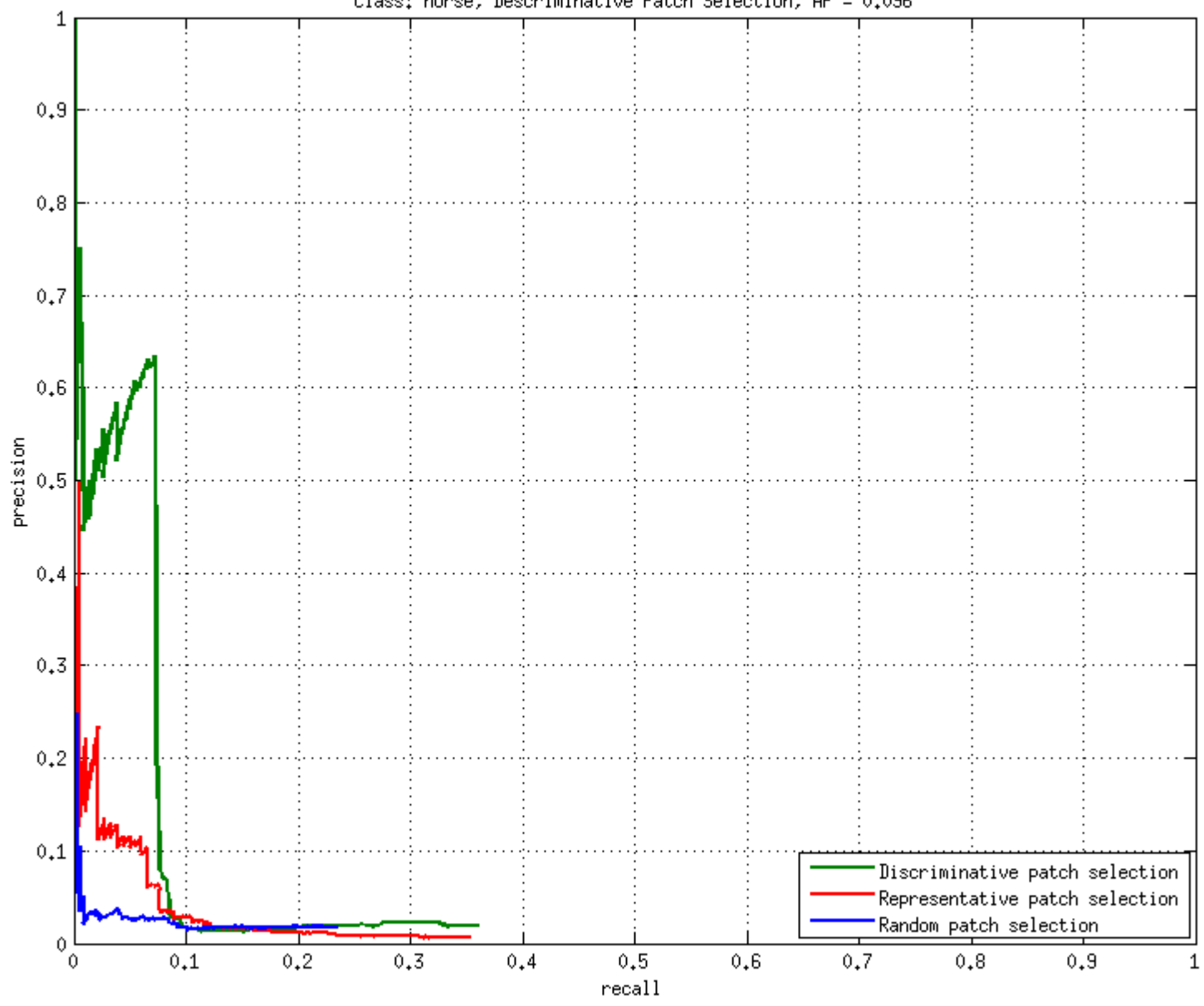
Class: horse, Random Patch Selection, AP = 0.026



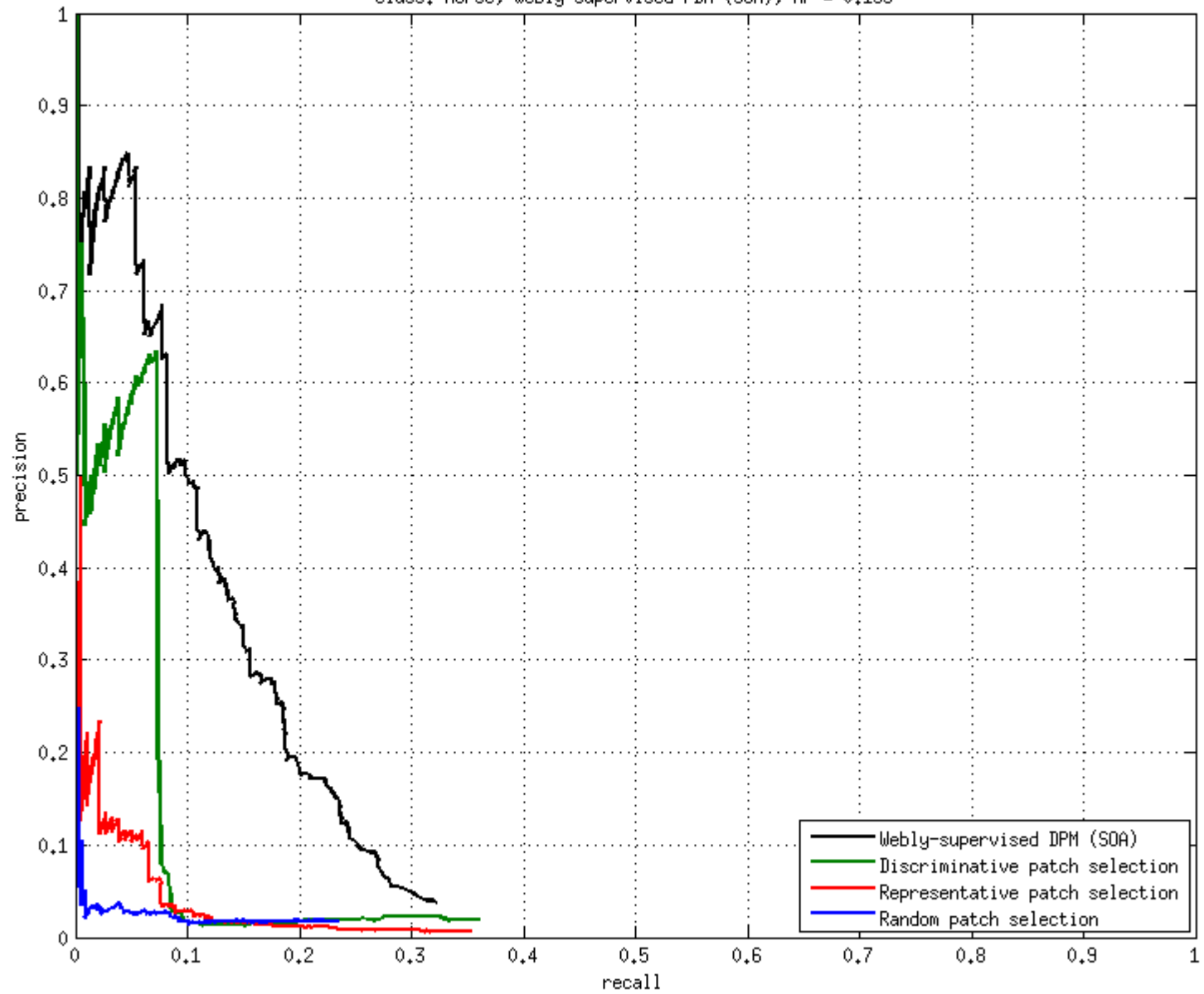
Class: horse, Representative Patch Selection, AP = 0.050



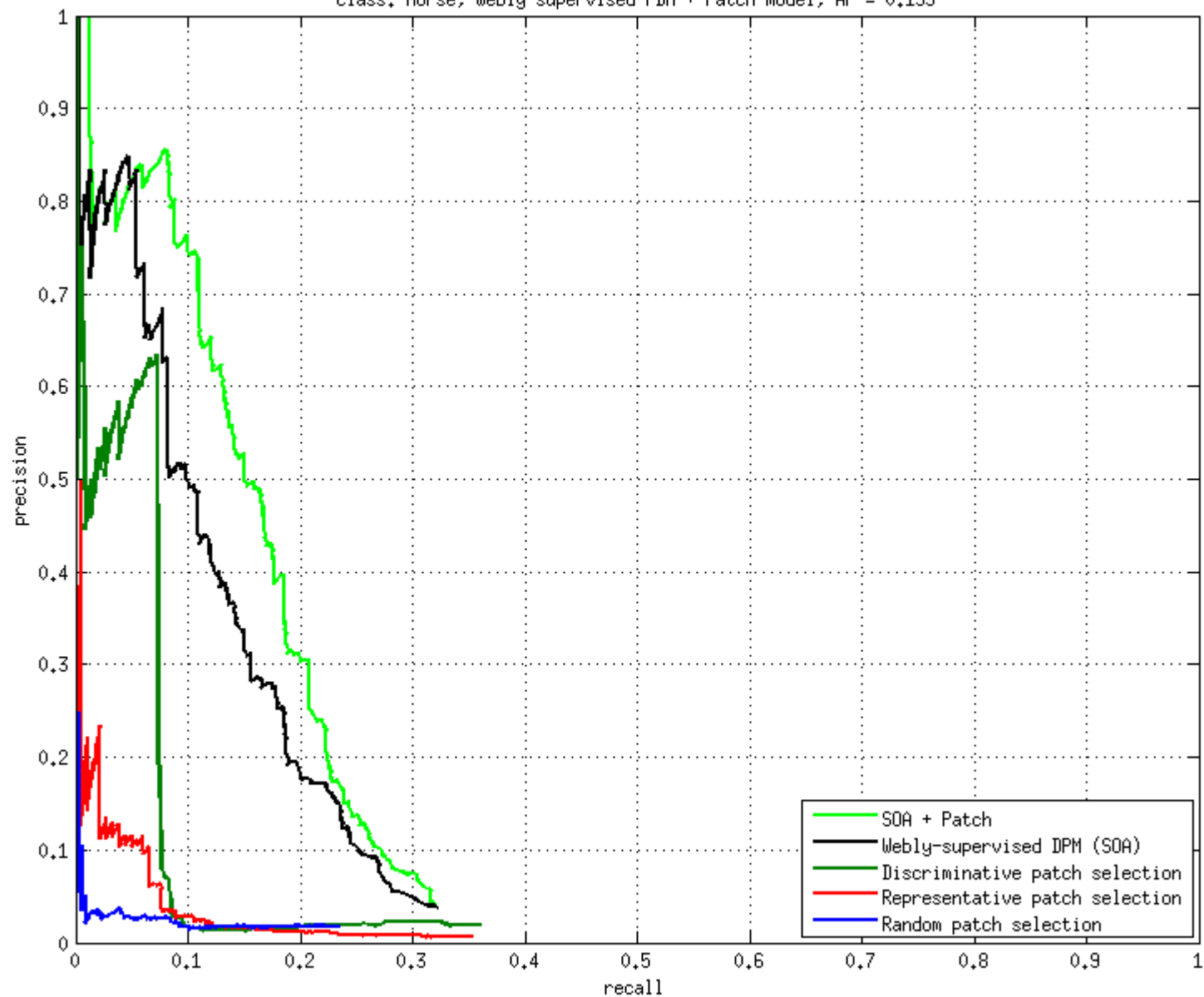
Class: horse, Discriminative Patch Selection, AP = 0,096



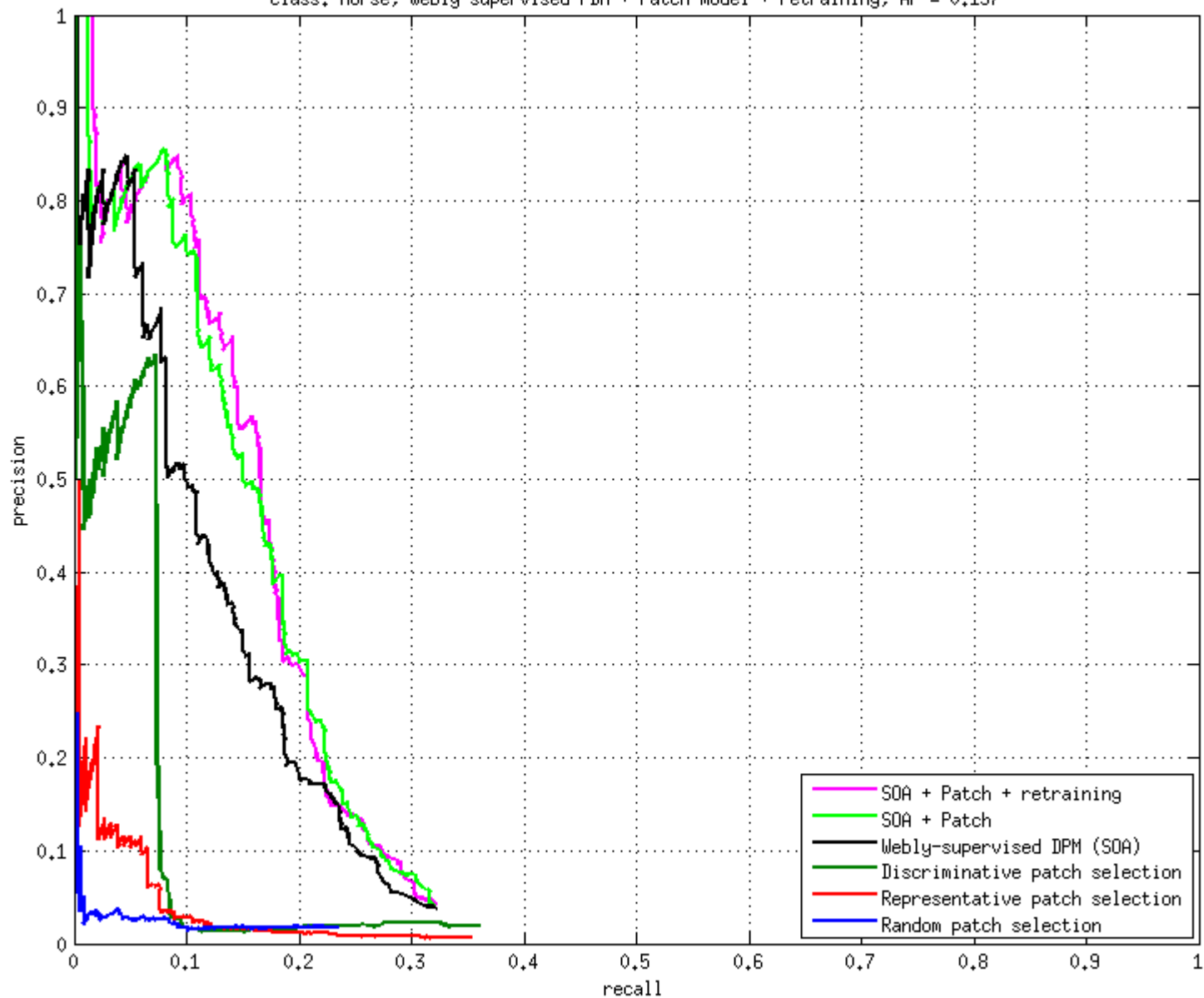
Class: horse, Webly-supervised PDM (SOA), AP = 0.159



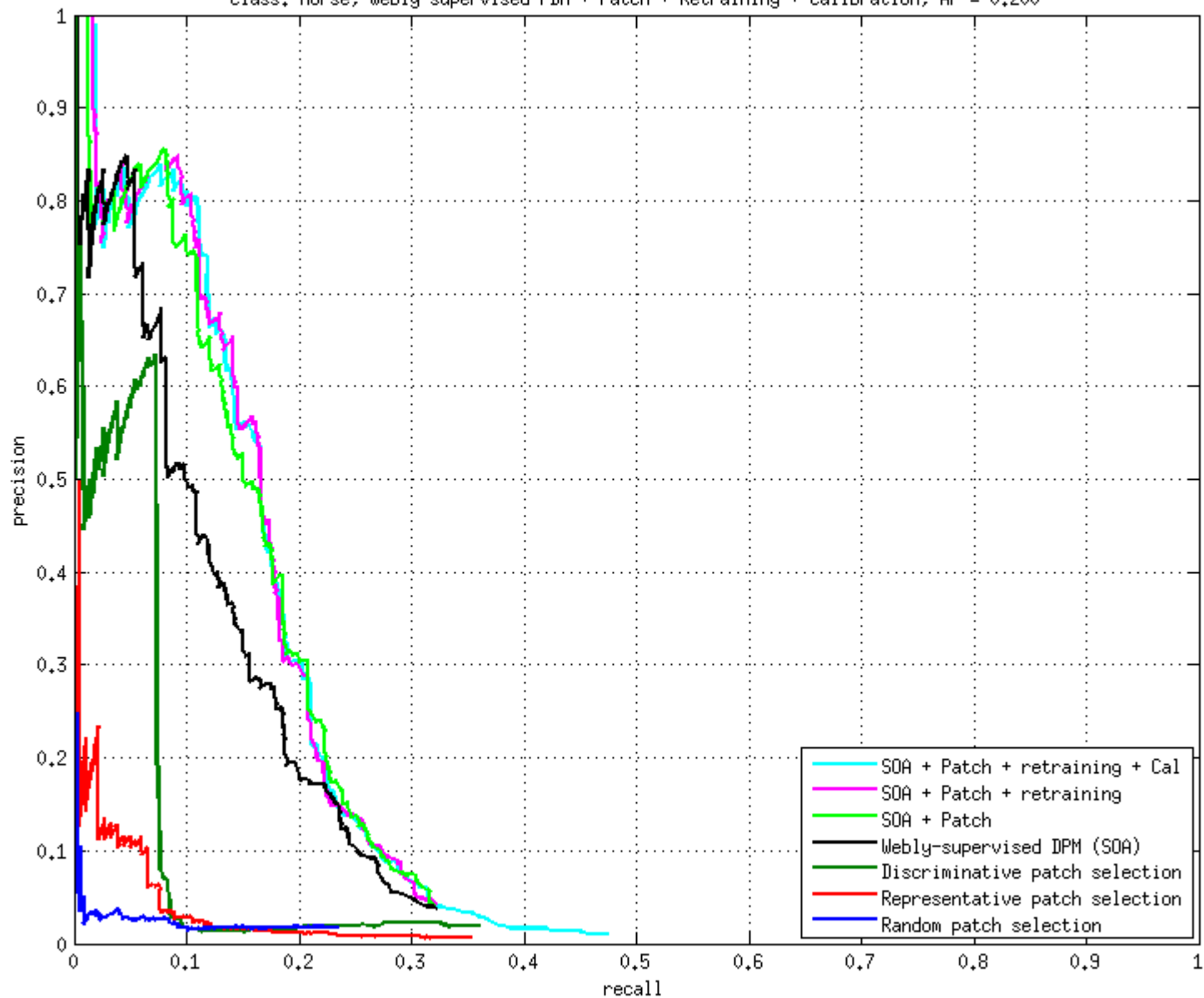
Class: horse, Webly-supervised PDM + Patch model, AP = 0,195



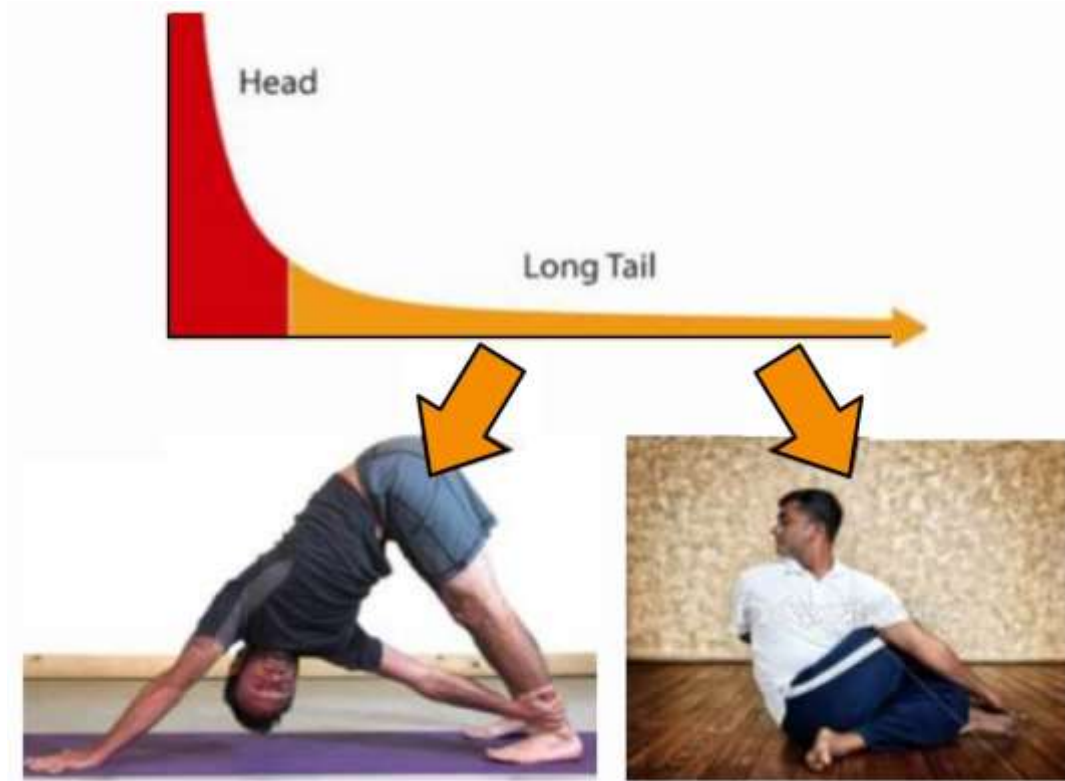
Class: horse, Webly-supervised PDM + Patch model + retraining, AP = 0,197



Class: horse, Webly-supervised PDM + Patch + Retraining + Calibration, AP = 0,200



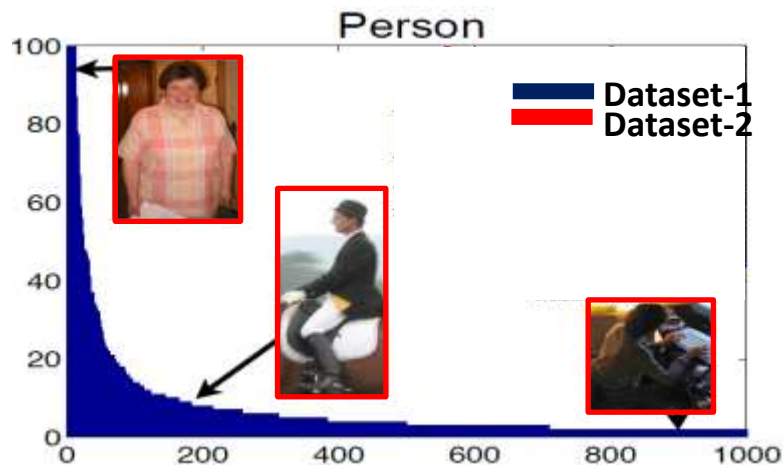
Long-tail distribution



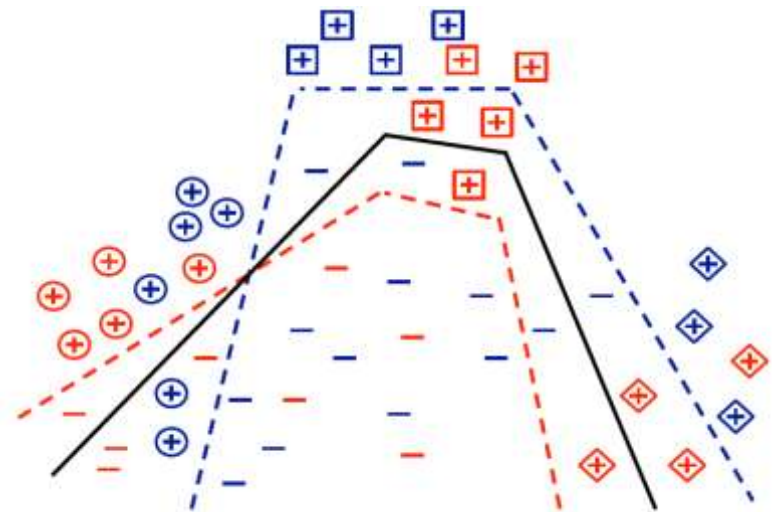
We need lots of templates, have little data of *'yoga twist'* poses

Sharing across datasets

- Solution: Enrich poor subcategory models with *statistical strength* borrowed from other datasets



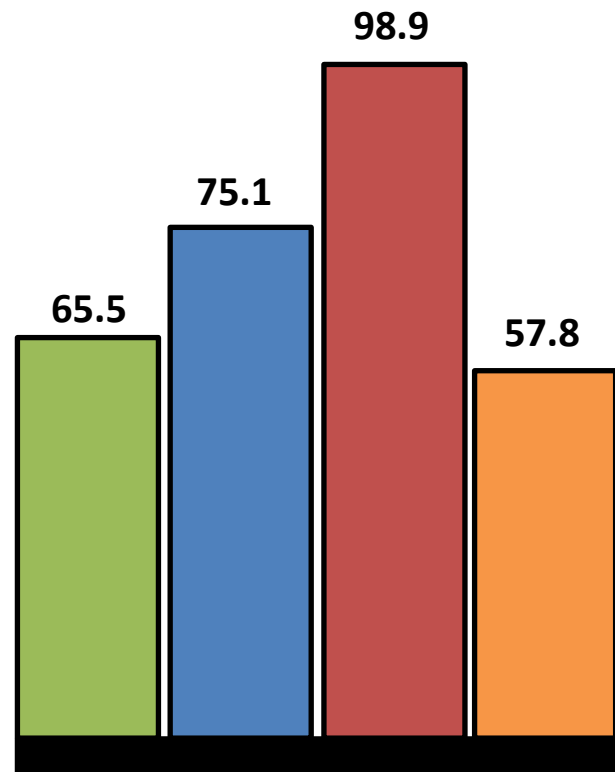
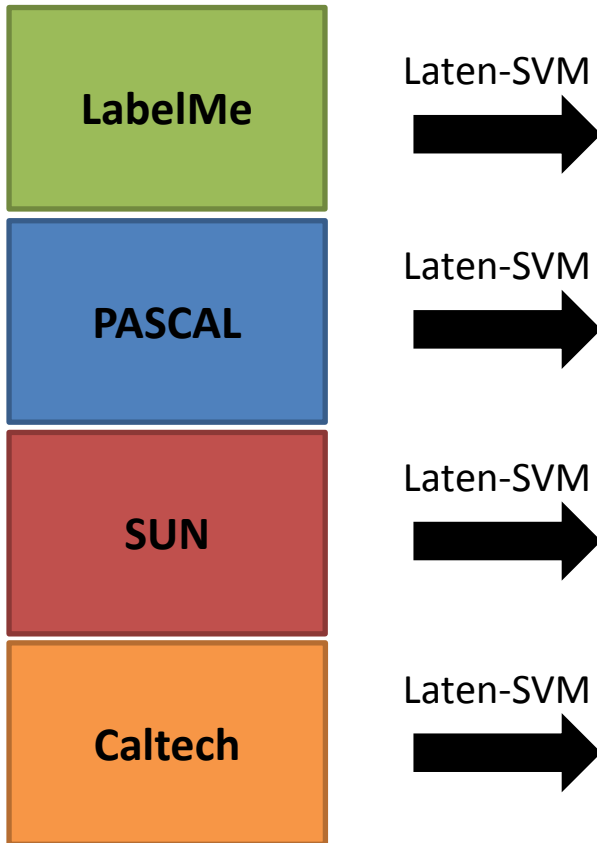
○ Sample sharing



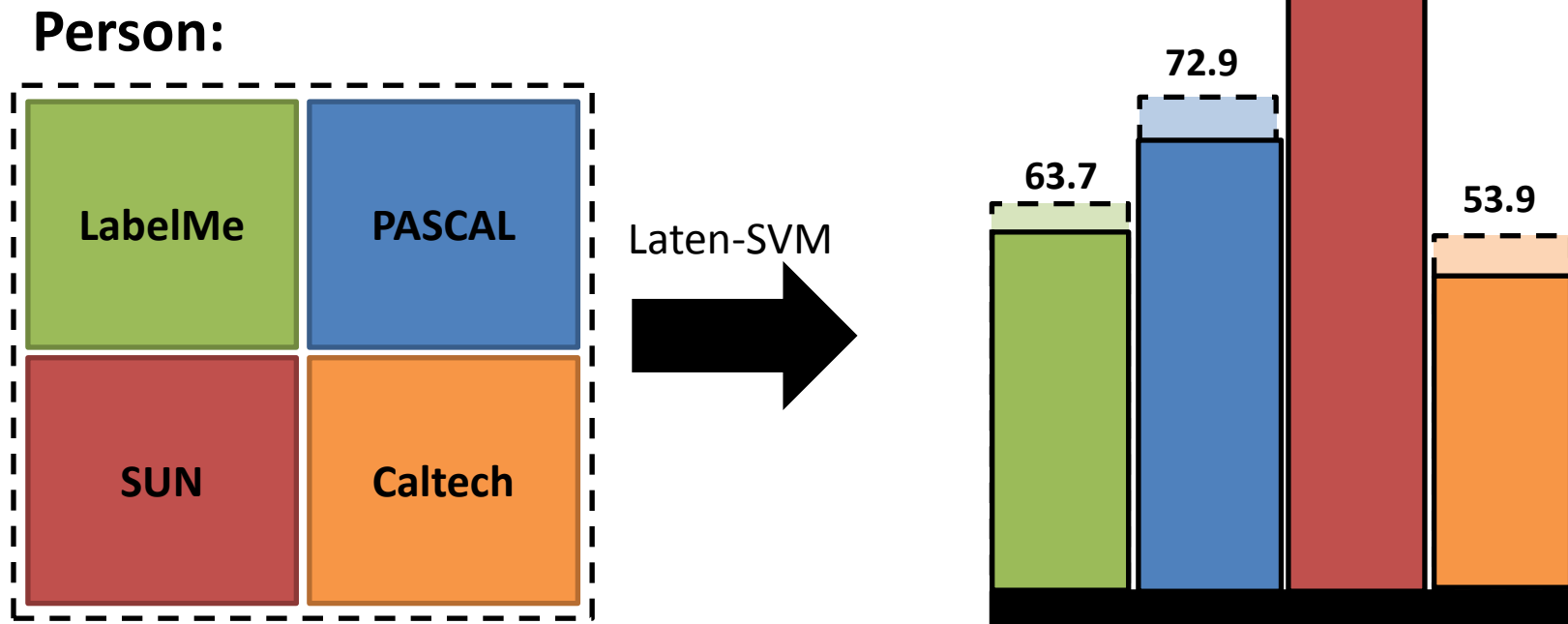
● Parameter sharing

Training on dataset individually

Person:

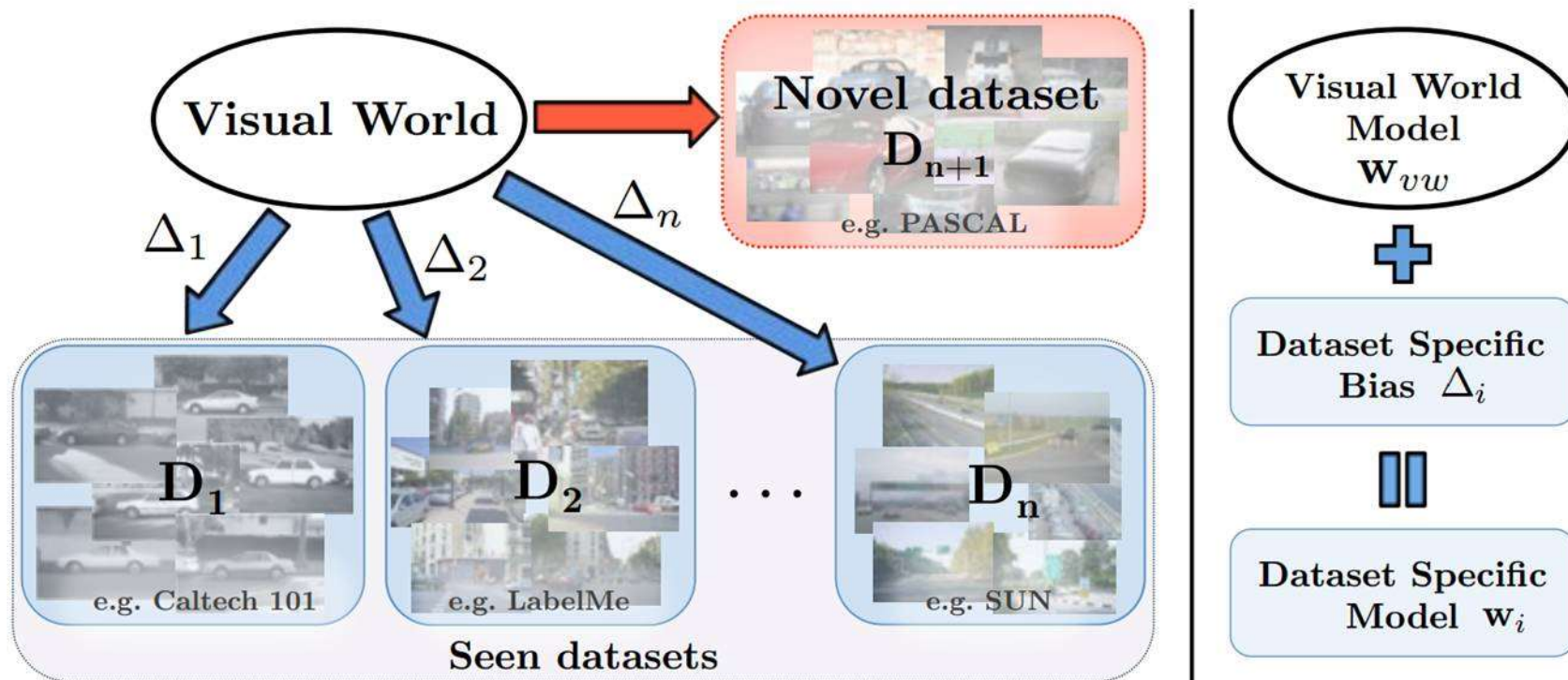


Concatenation of Datasets



Subcategory-based undoing bias

Extend the regularized multi-task learning framework* to the Latent Subcategory setting

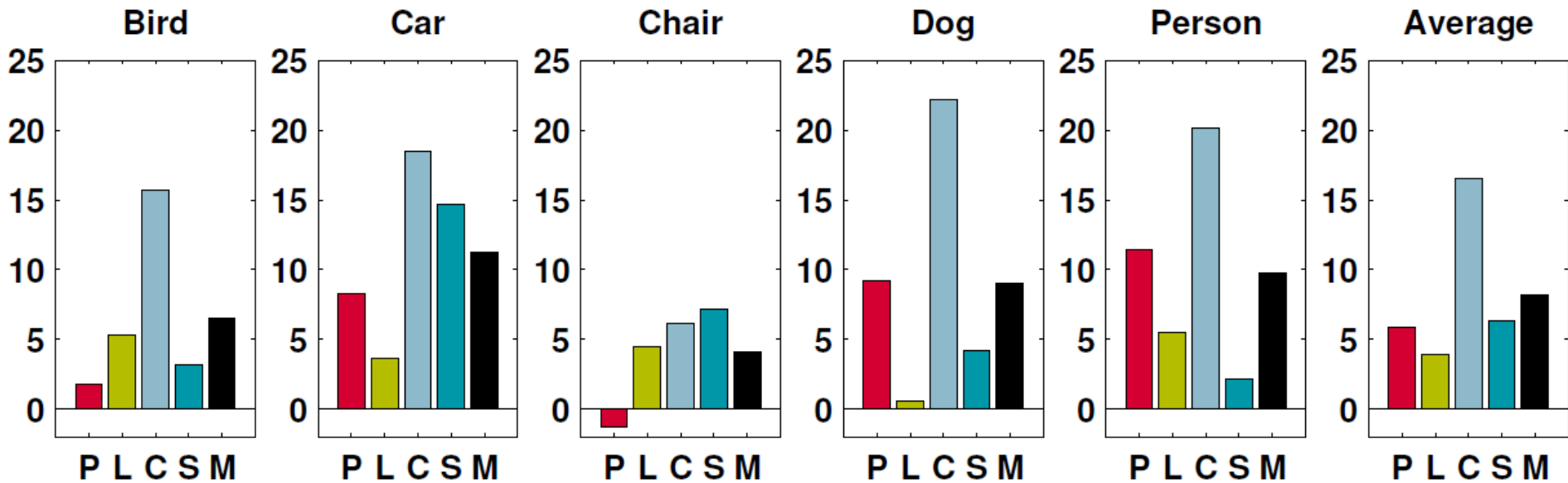


* A. Khosla, "Undoing the damage of dataset bias", In *ECCV* 2012.

Experiments

- **Leave-one-dataset-out**

- Compare to LSVM on concatenated dataset

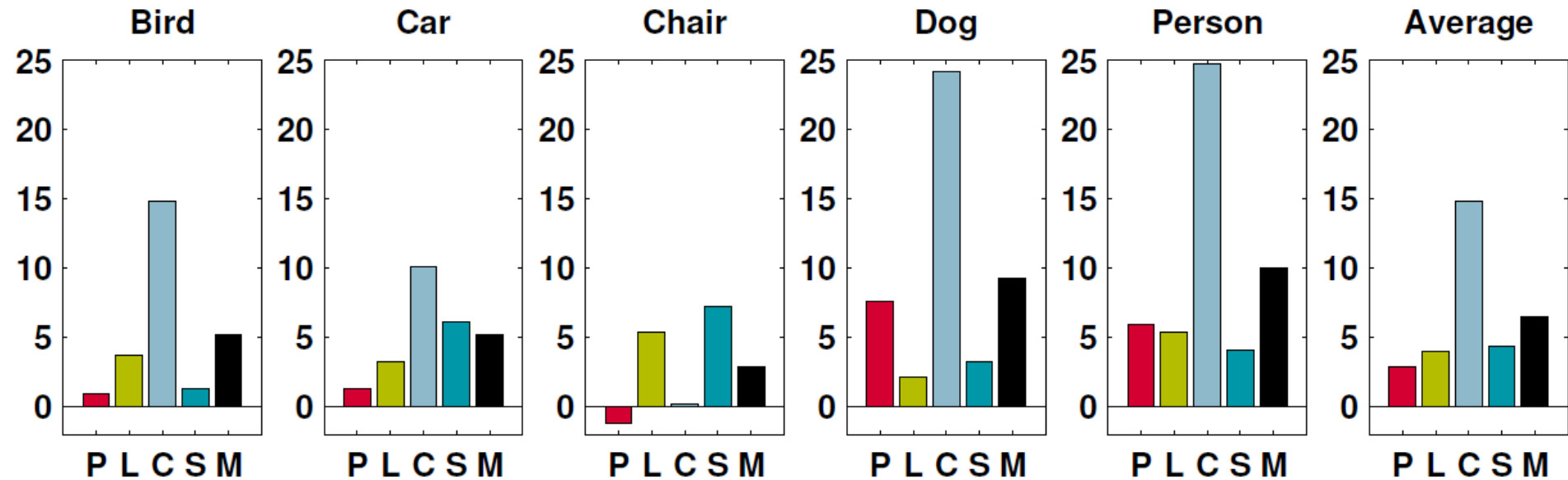


Average improvement: 8.5%

Experiments

- **Leave-one-dataset-out**

- Compare to SVM-based undo bias (SOA)



Average improvement: 6.5%

LabelMe



PASCAL



SUN



Caltech-101

