# Computer Vision by Learning

Cees Snoek, UvA

Arnold W.M. Smeulders, UvA

Efstratios Gavves, UvA

Laurens van der Maaten Facebook

# Tomorrow

Invited tutorial by **Laurens van der Maaten**

- Understanding and Improving Convolutional Networks
- From Visual Recognition to Visual Reasoning

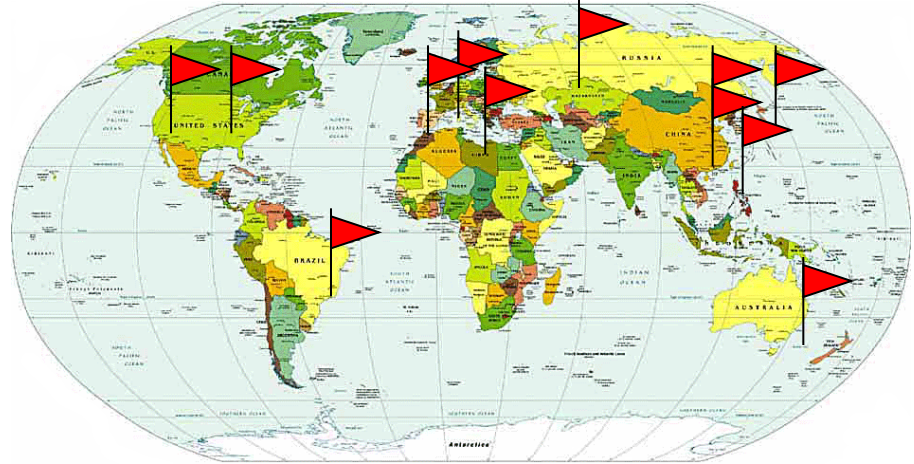Note change of location

- **CWI, Z009 Eulerzaal**

# CWI, Z009 Eulerzaal

# Overview

1. Image benchmarks, PASCAL, ImageNet, MSCOCO
2. Video benchmarks, TRECVID, ActivityNet
3. Labels from humans, experts, volunteers, crowdsourcing
4. Labels from similarity, nearest neighbor, simple features
5. Weakly-supervised computer vision

6. Event recognition by learning

# Evaluation of computer vision



Situation in 2000

- – Various video concept definitions
- – *Specific* and *small* data sets
- – Hard to compare methodologies

= Researchers

For object tracking still the case in 2013

# 1. Image benchmarks

The PASCAL Visual Object Classes (VOC) challenge is a benchmark in visual object category recognition and detection, which provides challenging images and high quality annotation, together with a standard evaluation methodology. Measured the state-of-the-art on a yearly basis from 2005 to 2012. It has been succeeded by the ImageNet challenge which evaluates algorithms for object detection and image classification at large scale.

# Pascal Dataset Collection

500K Images downloaded from **flickr** and random subset selected for annotation
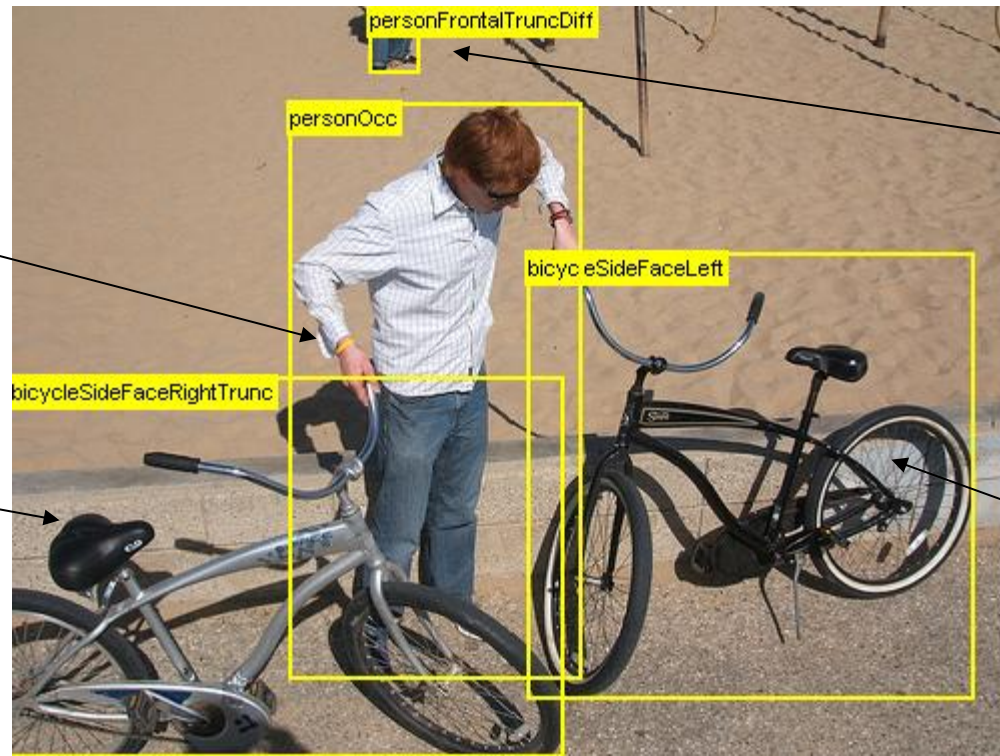
Complete annotation of all objects from 20 categories



**Occluded**
Object is significantly occluded within BB

**Truncated**
Object extends beyond BB

**Difficult**
Not scored in evaluation

**Pose**
Facing left

# Examples

**Dining Table**



**Dog**



**Horse**



**Motorbike**



**Person**



**Potted Plant**



**Sheep**



**Sofa**



**Train**



**TV/Monitor**

# 2010 Dataset Statistics

| | Training | | Testing | |
|---|---|---|---|---|
| **Images** | 10,103 | (7,054) | 9,637 | (6,650) |
| **Objects** | 23,374 | (17,218) | 22,992 | (16,829) |

VOC2009 counts shown in brackets

Minimum ~500 training objects per category

   ~1700 cars, 1500 dogs, 7000 people

~Equal distribution across training and test sets

# PASCAL VOC Challenges

## Object classification
– Does the image contain an airplane?



## Object deteciton
– Where is the airplane, (if any)?



## Object segmentation
– Which pixels are part of an airplane, (if any)?

# ImageNet Challenge

Yearly competition

Automatically label 1.4M images with 1K objects

Measure top-5 classification error



**Output**
Scale
T-shirt
Steel drum ✔
Drumstick
Mud turtle

**Output**
Scale
T-shirt
Giant panda ✘
Drumstick
Mud turtle

# Some highlights

## Year 2010

NEC-UIUC



Dense grid descriptor: HOG, LBP

↓

Coding: local coordinate, super-vector

↓

Pooling, SPM

↓

Linear SVM

Lin *et al.* CVPR11

## Year 2012

SuperVision



Krizhevsky *et al.* NIPS12

## Year 2014

GoogLeNet



Convolution
Pooling
Softmax
Other

Szegedy *et al.* CVPR15

VGG



image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096
FC-1000
softmax

Simonyan *et al.* ICLR15

# Progress in ImageNet



*Machine makes less mistakes than human*

# Progress: Classification & Detection

# ImageNet object detection

**Modeled after PASCAL VOC**

Algorithm outputs a list of bounding box detections with confidences

A detection is considered correct if intersection over union (IoU) overlap with ground truth > threshold (0.5)

Evaluated by average precision per object class

Winner is the team that wins the most object categories

# ImageNet detection challenge

| Statistics | | PASCAL VOC 2012 | | ILSVRC 2013 |
|---|---|---|---|---|
| Object classes | | 20 | **10x** | **200** |
| Training | Images | 5.7K | | **395K** |
| | Objects | 13.6K | **25x** | **345K** |
| Validation | Images | 5.8K | | **20.1K** |
| | Objects | 13.8K | **4x** | **55.5K** |
| Testing | Images | 11.0K | | **40.1K** |
| | Objects | --- | | --- |



Person
Car
Motorcycle
Helmet

# Progress

**Mean average precision on test set**



| Model | Value |
|---|---|
| CUHK ('16) | 66.3 |
| ResNet ('15) | 62.1 |
| GoogLeNet ('14) | 44 |
| UvA/Euvision ('13) | 23 |

# MSCOCO

80 object categories

200k images

1.2M instances (350k people)

106k people with keypoints


Dataset examples

# Instance segmentations



Every instance segmented in MSCOCO

# Challenges in 2016

# Segmentation winner

Fully convolutional end-to-end for instance segmentation
Based on ResNet-101



fg

conv

bg

translation-aware fg/bg score maps

position-sensitive
RoI pooling

position-sensitive
RoI pooling

fg likelihood

bg likelihood

pixel-wise softmax

pixel-wise
max

"person"
aggregate
& vote

# Some results

# Some more

# 2. Video benchmarks

Crucial drivers for progress in large-scale computer vision are international search engine benchmarks. The National Institute of Standards and Technology's TRECVID (TREC Video Retrieval) benchmark has played a significant role. The main goal of TRECVID is to promote progress in content-based analysis of and retrieval from digital video via open, metrics-based evaluation. TRECVID is a laboratory-style evaluation that attempts to model real world situations or significant component tasks involved in such situations.

# International competition

NIST TRECVID Benchmark

Promote progress in video retrieval research

Open data, tasks, evaluation *and* innovation

http://trecvid.nist.gov/

# Video data sets

US TV news (`03/`04)



International TV news (`05/`06)



Dutch TV infotainment (`07/`08/`09)



Web video (since 2010)

# NIST TRECVID evolution

# Task: concept detection

Goal

– Build benchmark collection for visual concept detection methods

Secondary goals

– encourage generic (scalable) methods for detector development
– semantic annotation is important for search/browsing

**Aircraft**

**Beach**

Note the variety in visual appearance

**Mountain**

# De facto evaluation standard

# Annotation efforts

# Measuring performance

Results

1.
2.
3.
4.
5.

Set of relevant items

Set of retrieved items

Set of relevant retrieved items

Precision

Recall

inverse relationship

# Evaluation measure

Average Precision

- – Combines precision and recall
- – Averages precision after relevant shot
- – Top of ranked list most important

$$AP = \frac{\sum_{r=1}^{N}(P(r) \times \mathrm{rel}(r))}{\text{number of relevant documents}}$$

$$AP = \frac{1/1 + 2/3 + 3/4 + \dots}{\text{number of relevant documents}}$$

Results

1. 
2. 
3. 
4. 
5. 

# Progress in video concept search



• = 1000+ others

* = UvA / Euvision / Qualcomm

# 2010: Bag-of-words

Color SIFT, soft assignment and  kernel approximations.



**Software available for download at http://colordescriptors.com**

# Benchmarking is compute intensive

Distributed ASCI super computer: *priceless*

# Performance doubled in 3 years

• 36 concept detectors



Mean Average Precision

2006    2009

# 2013: AlexNet-variant



224

224

11×11      5×5      3×3      3×3      3×3

Layer 6      Dropout      Layer 7      Dropout      Loss

4,096      4,096

Convolution      Non-linearity      Pooling

# Latest jump due to deep learning

# MediaMill video search engine

CrossBrowser combines query results and time



time

ranked results

**2010 Version**

# Other challenge: Instance search

Given a single query example, including a segmentation mask, find similar occurrences of the named instance in a collection of video.

instance "Eiffel tower"

instance "a circular 'no smoking' logo"

instance "Stephen Colbert"

instance "an Audi logo"

instance "this man"

# Other challenge: event recognition

Given 100, 10 or 0 training example videos, recognize and recount videos in a huge test collection containing the event of interest.

Working on a metal project



Cleaning an appliance

# Goal

Recognize all activities in daily life

# ActivityNet

# Challenges

- ## Task I: Untrimmed Video Classification

input: long untrimmed video

output

activity presence (binary)

- ## Task II: Activity Detection

input: long untrimmed video

output

activity temporal location

http://activity-net.org/challenges/2016/

# 3. Labels from humans

The most precious resource in computer vision by learning is data.

The most traditional source for obtaining labeled examples is to rely on human experts. The Internet has launched the trend to let volunteers label visual content, either for fun, for winning a game or for a small compensation. ImageNet is a labeled image database organized according to the WordNet hierarchy in which each node of the hierarchy is depicted by hundreds of images.

# Labeling by library experts

LSCOM (Large Scale Concept Ontology for Multimedia)

Provides manual annotations for 449 concepts

– In international broadcast TV news

Connection to Cyc ontology

http://www.lscom.org/

# Labeling by volunteers

# Polygon quality

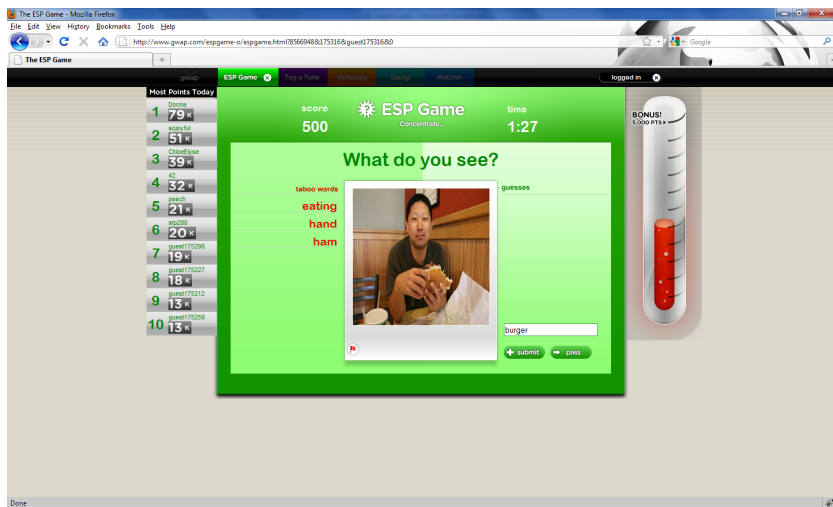# Online hooligans

# Testing



**Most common labels:**

test

adksdsa

woiieiie

# Quiz: downside of volunteers?

Lack of incentive

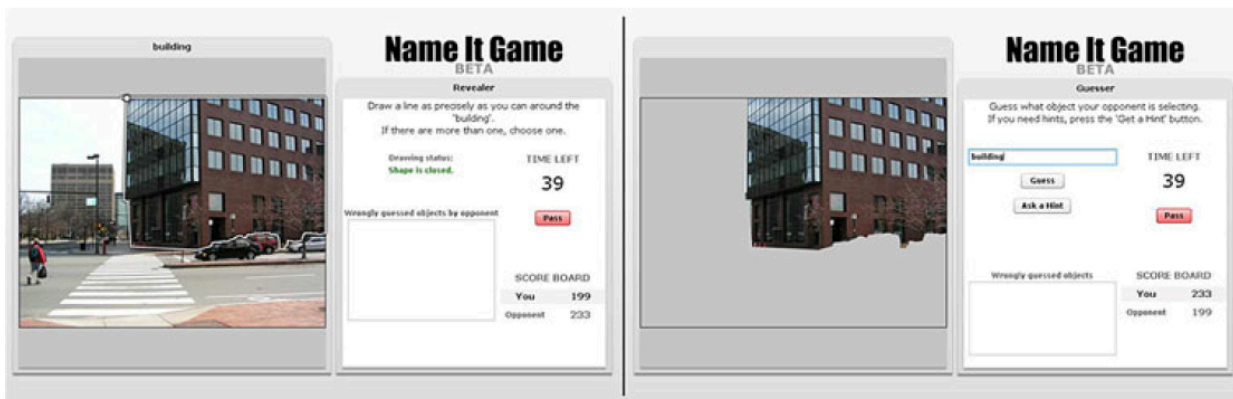Limited quality control

Limited number of labels

# Labels from games



von Ahn, ESP Game



Bubble sizes as proportions of image

Deng CVPR 2013



Steggink, MM Sys, 2011

# Labels from games

Games are a fun way to motivate volunteers
- Words are often too abstract
- Requires some sort of label validation

More descriptive labels by
- Adding semantic structure
- Linking labels to regions

Any game suffers from lack of popularity

# Labels from micro-payments

ImageNet (11M images)

- 4000 categories

- > 100 examples

SUN (130K images)

- 397 scene categories

- > 100 examples



Deng et al, CVPR 2009



Xiao et al, CVPR 2010

# IM GENET demo

# Constructing ImageNet

**Step 1:**
Collect candidate images via the Internet

→

**Step 2:**
Clean up the candidate Images by humans

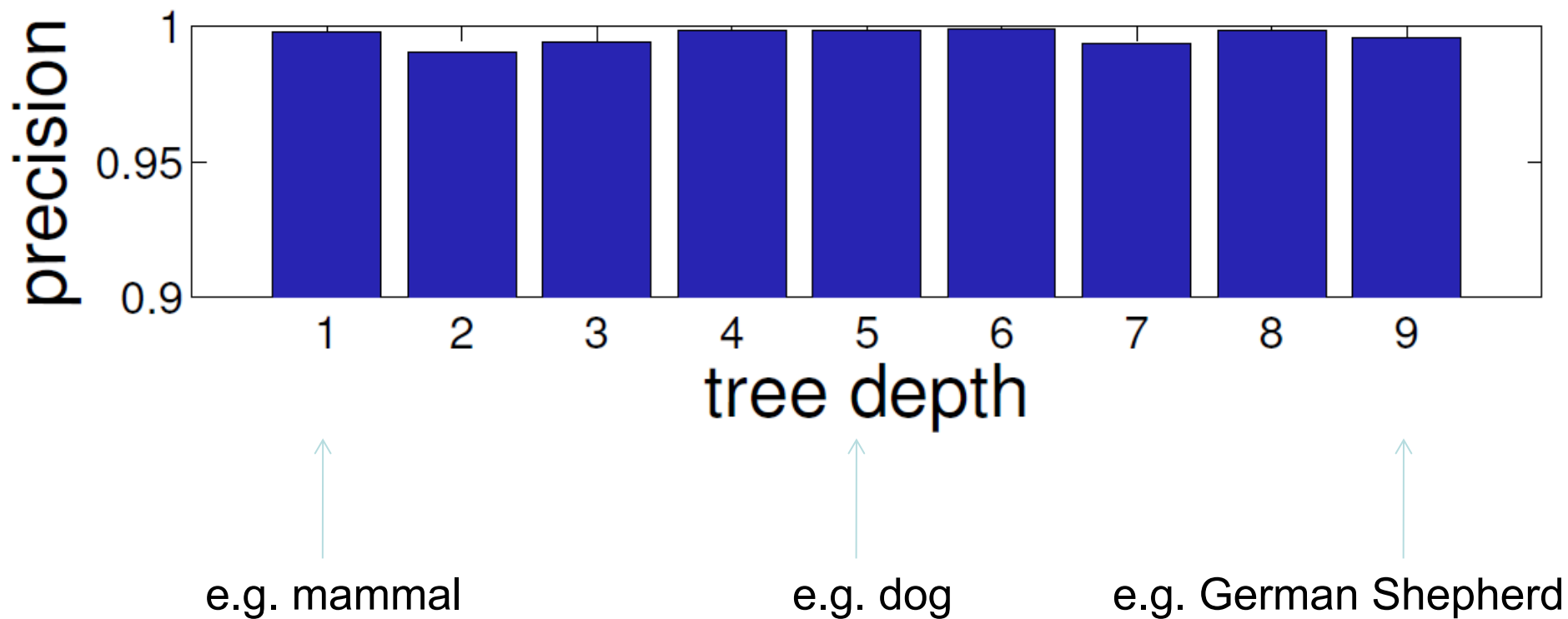# IMAGENET  is  built by crowdsourcing
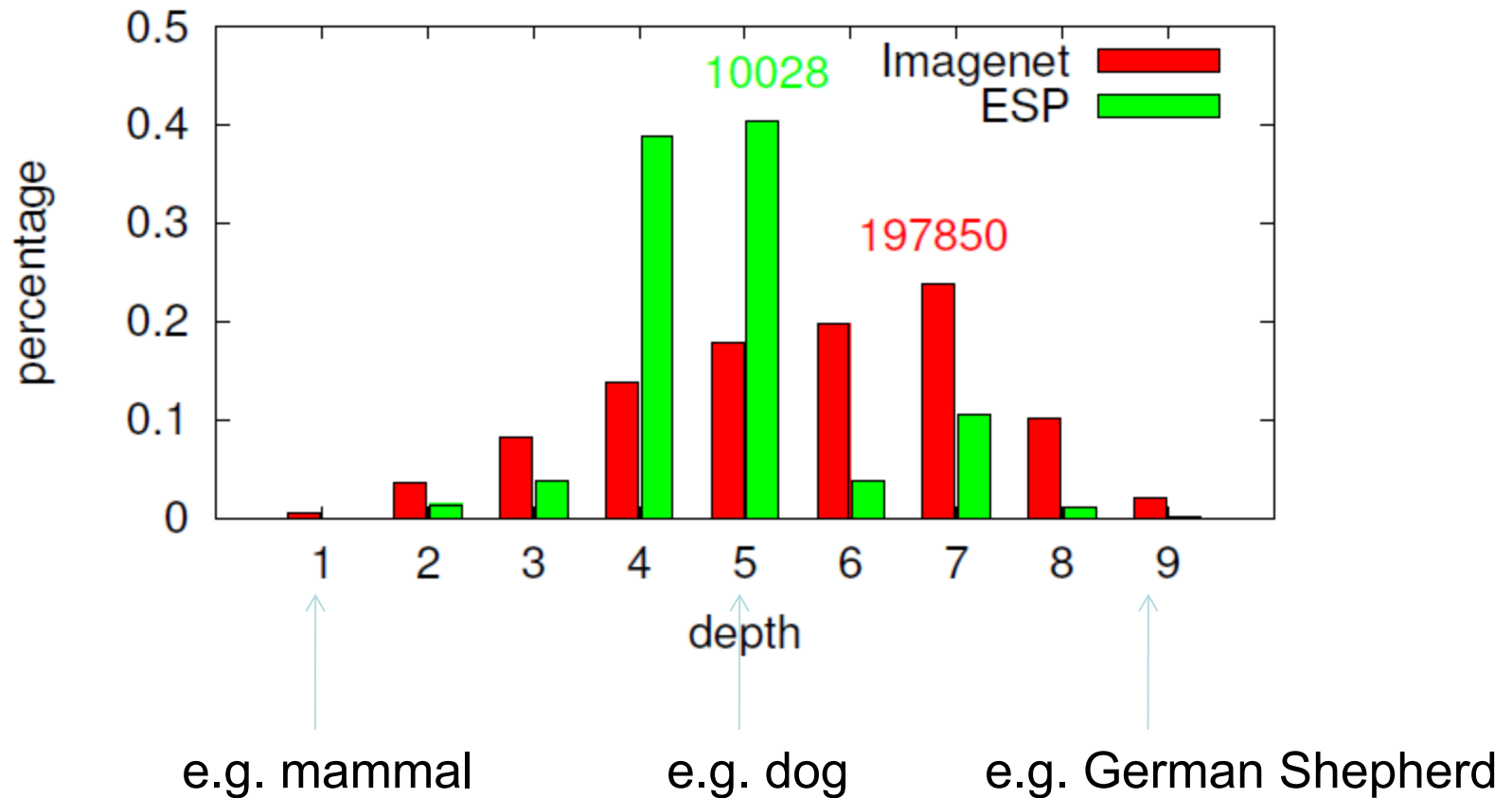
July 2008: 0 images

Dec 2008: 3 million images, 6K+ synsets

April 2010: 11 million images, 15K+ synsets

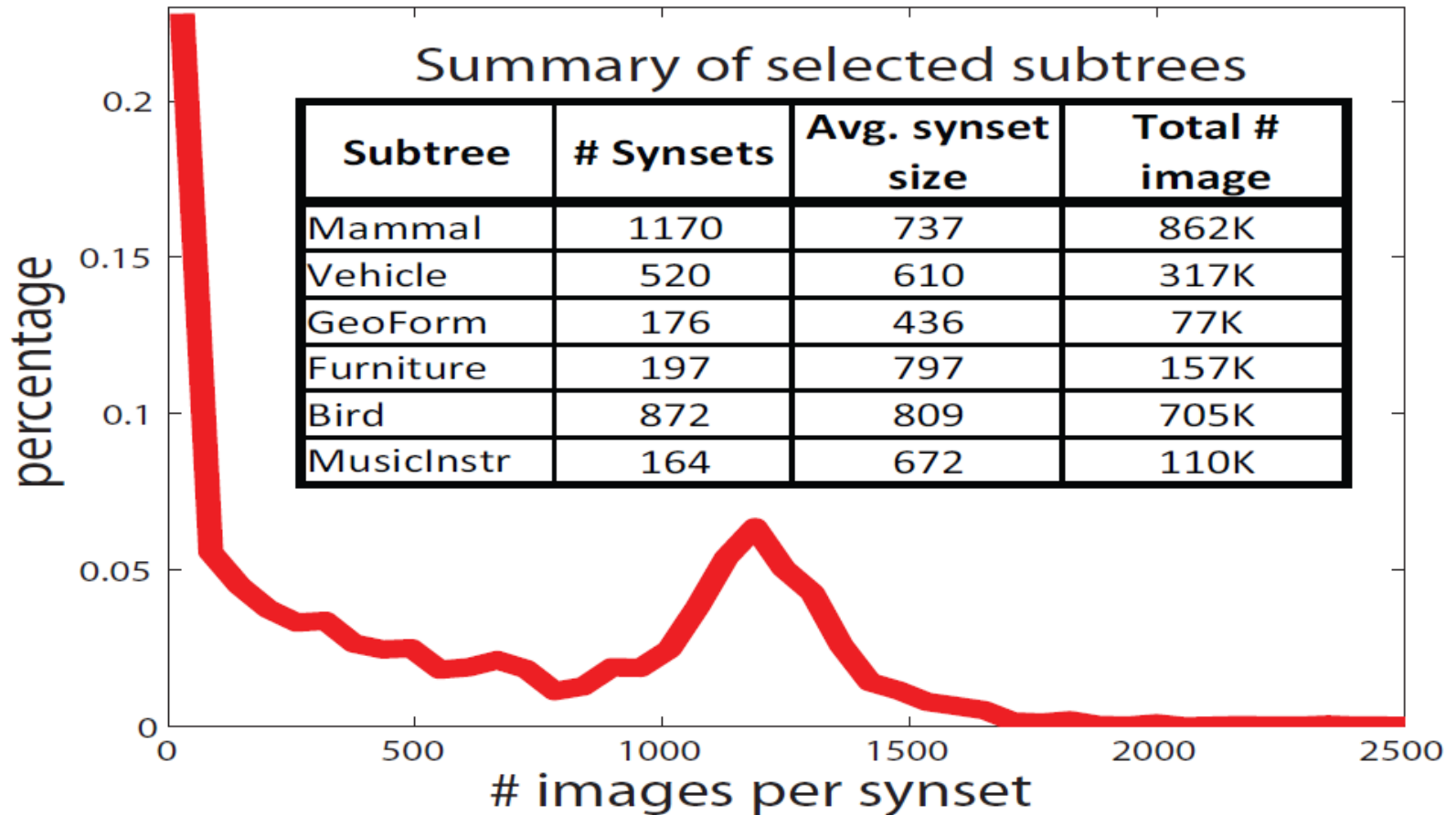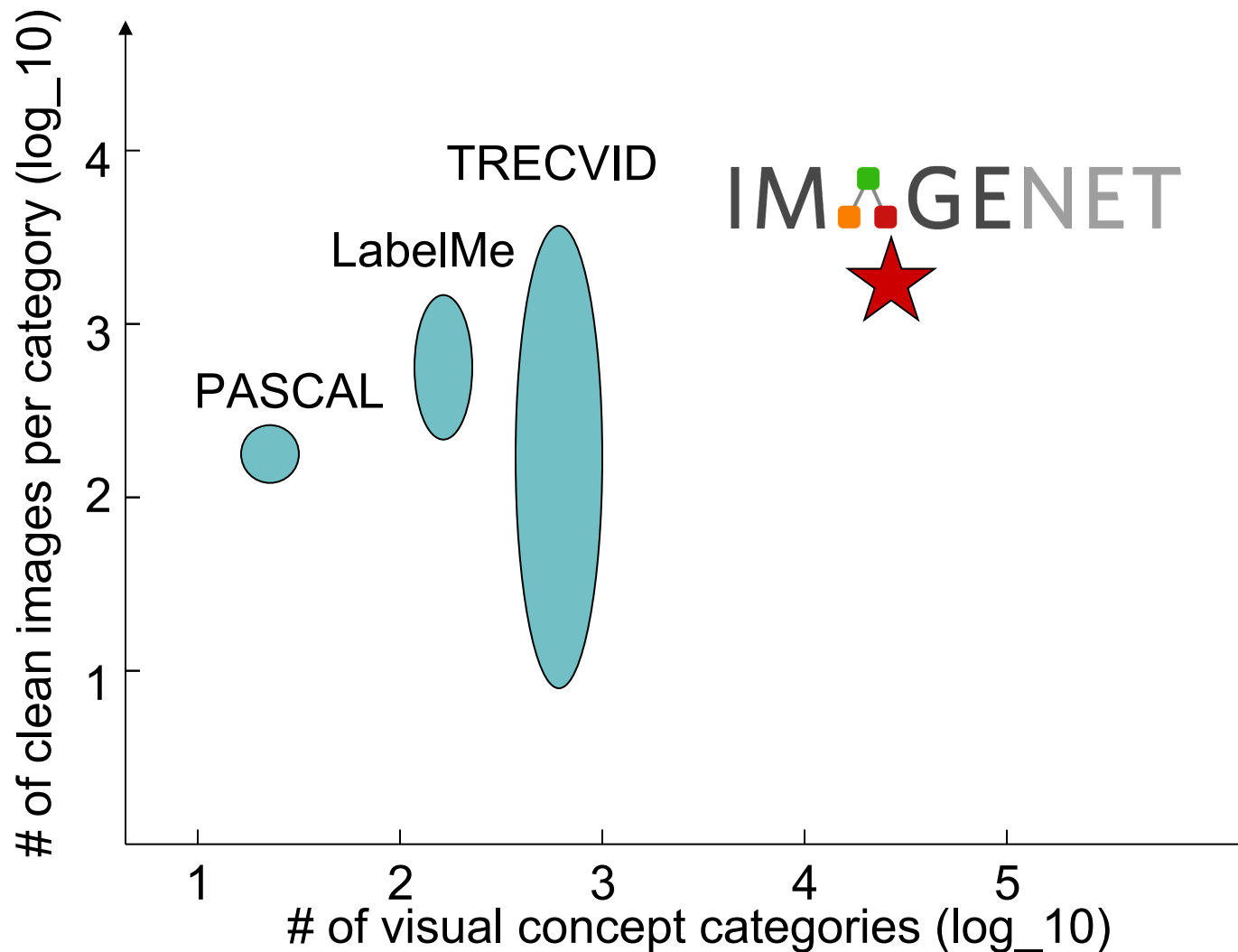Yesterday: 14 million images, 21K synsets indexed

# Accuracy

# Diversity

# Scale



| Subtree | # Synsets | Avg. synset size | Total # image |
|---|---|---|---|
| Mammal | 1170 | 737 | 862K |
| Vehicle | 520 | 610 | 317K |
| GeoForm | 176 | 436 | 77K |
| Furniture | 197 | 797 | 157K |
| Bird | 872 | 809 | 705K |
| MusicInstr | 164 | 672 | 110K |

Summary of selected subtrees

# Datasets comparison

# Constructing ImageNet

**Step 1:**
Collect candidate images via the Internet

→

**Step 2:**
Clean up the candidate Images by humans

# Constructing ImageNet

**Free**

Step 2:
Clean up the candidate
Images by humans

YAHOO!

picsearch™

flickr™

Live Search

Google™
Image Search
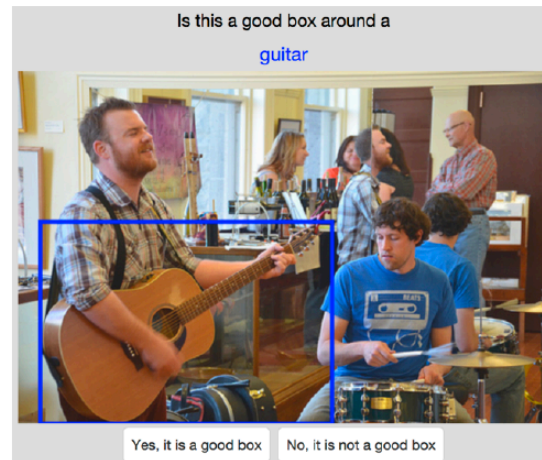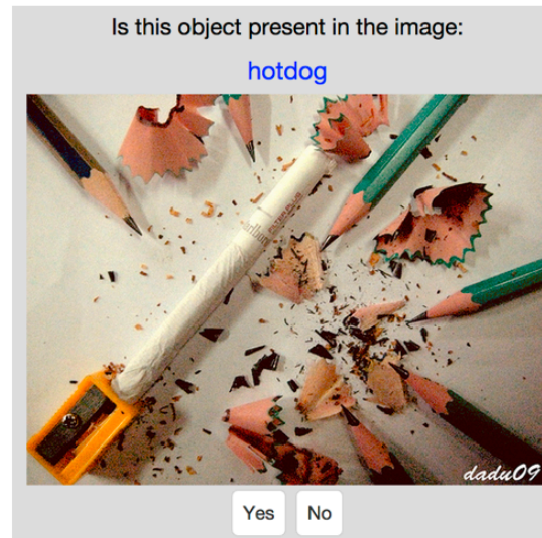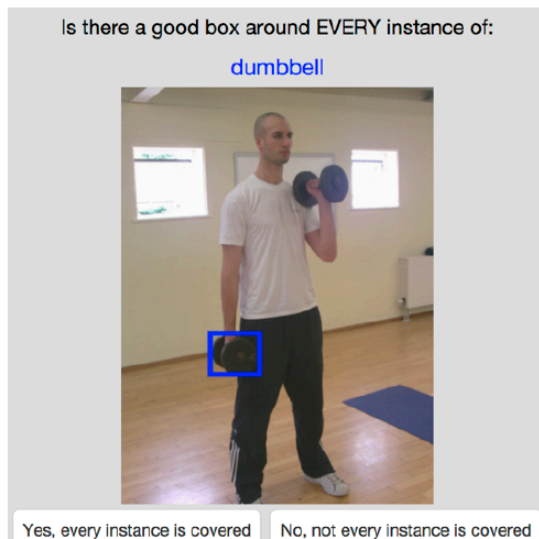
amazon.com

amazonmechanical turk
Artificial Artificial Intelligence
beta

# Constructing ImageNet

Free → $$$

# User interfaces

For image labeling

# 4. Labels from similarities

The most precious resource in computer vision by learning is data.

Huge amounts of weakly labeled images and videos are available online. How reliable are these tags? Can we use them for learning classifiers, segment images, or localize distinctive parts? It turns out that 'good old' nearest neighbor with simple visual features provides a free, scalable and effective means to collect valuable data.

**YouTube** — **72** HOURS of Video uploaded

**Google** — **2** MILLION Searches

**Spotify** — **14** New Songs Added

DOMAINS — **70** New Registered

**facebook** — **41** THOUSAND posts every second

**1.8** MILLION likes

**350GB** of data

**skype** — **1.4** MILLION Minutes Connecting with Each other

**15** THOUSAND Tracks downloaded from iTunes

WORDPRESS — **347** New Blog Posts

**flickr** — **20** MILLION Photo Views

WEBSITES — **571** New Created

**twitter** — **278** THOUSAND Tweets

**17** THOUSAND Transactions — Walmart

**104** THOUSAND Photos Shared

snapchat

New Tumblr Photos — **20** THOUSAND — tumblr.

Amazon Sales **$83,000**

amazon.com

Professional searches — **11** THOUSAND

Linked in

Active Users **11** THOUSAND

Pinterest

photos every second **3,600**

Instagram

emails sent **204** MILLION — EMAIL

# Fundamental problem

Social tags for image and video were never meant to meet professional standards, consequently they are

- – subjective

- – ambiguous,

- – overly personalized, and

- – limited.

Tagged images are notoriously difficult to find.

# Searching for 'tiger'

# Searching for 'classroom'

kindergarten classroom layout of ... ✗

stone age stone age ... ✗

365days me of me ... ✗

view details

cfc minnesotawoods minnesota forest ... ✗

12 favorite classroom incubator ✓

tour tampere church upload ... ✓

view details

# Quiz

*What image tags in this example are suited as training label?*



bicycle
perfect
bridge
MyWinners

# Computer vision is essential

## Free text



## User tags

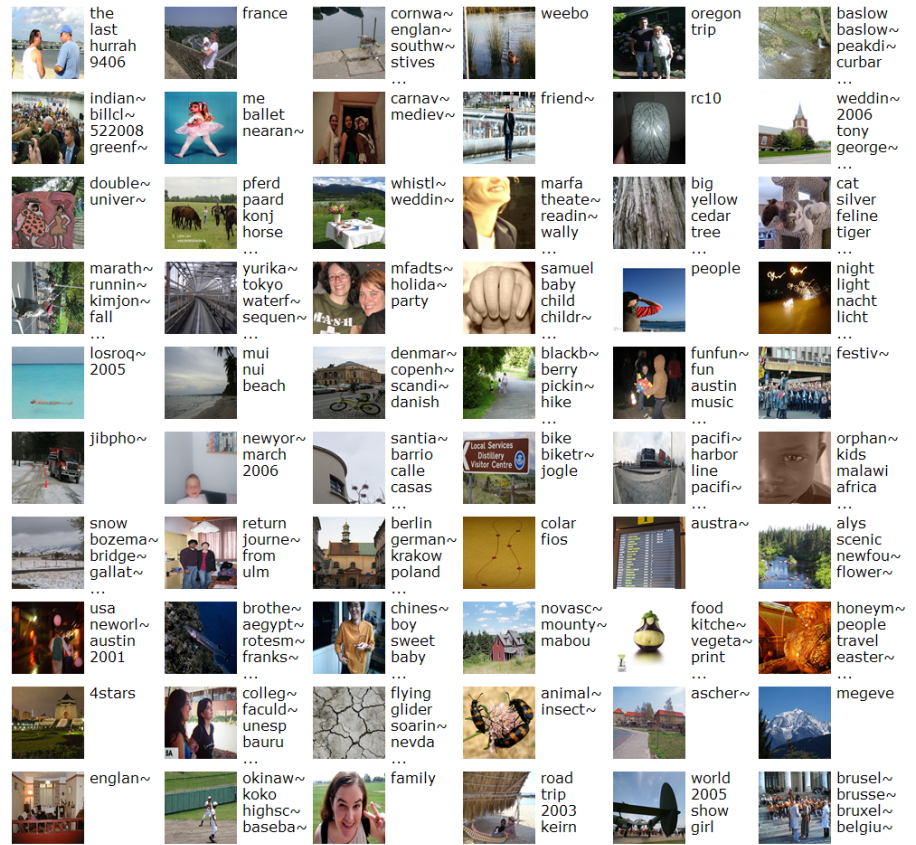

bridge
bicycle
perfect
MyWinners

?

bridge
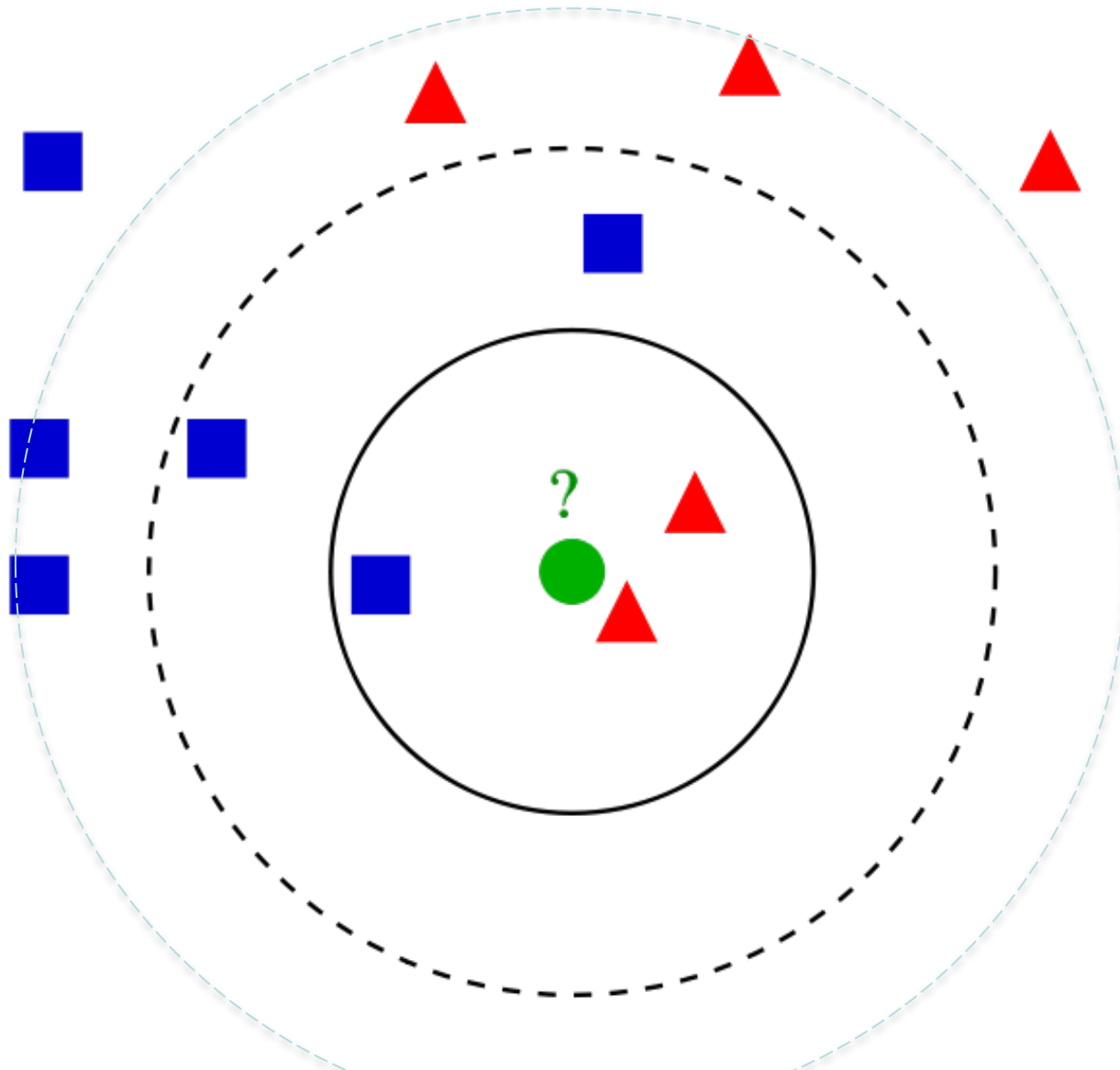bicycle
perfect
MyWinners

# Challenges

Many tags & many images

A prospective algorithm
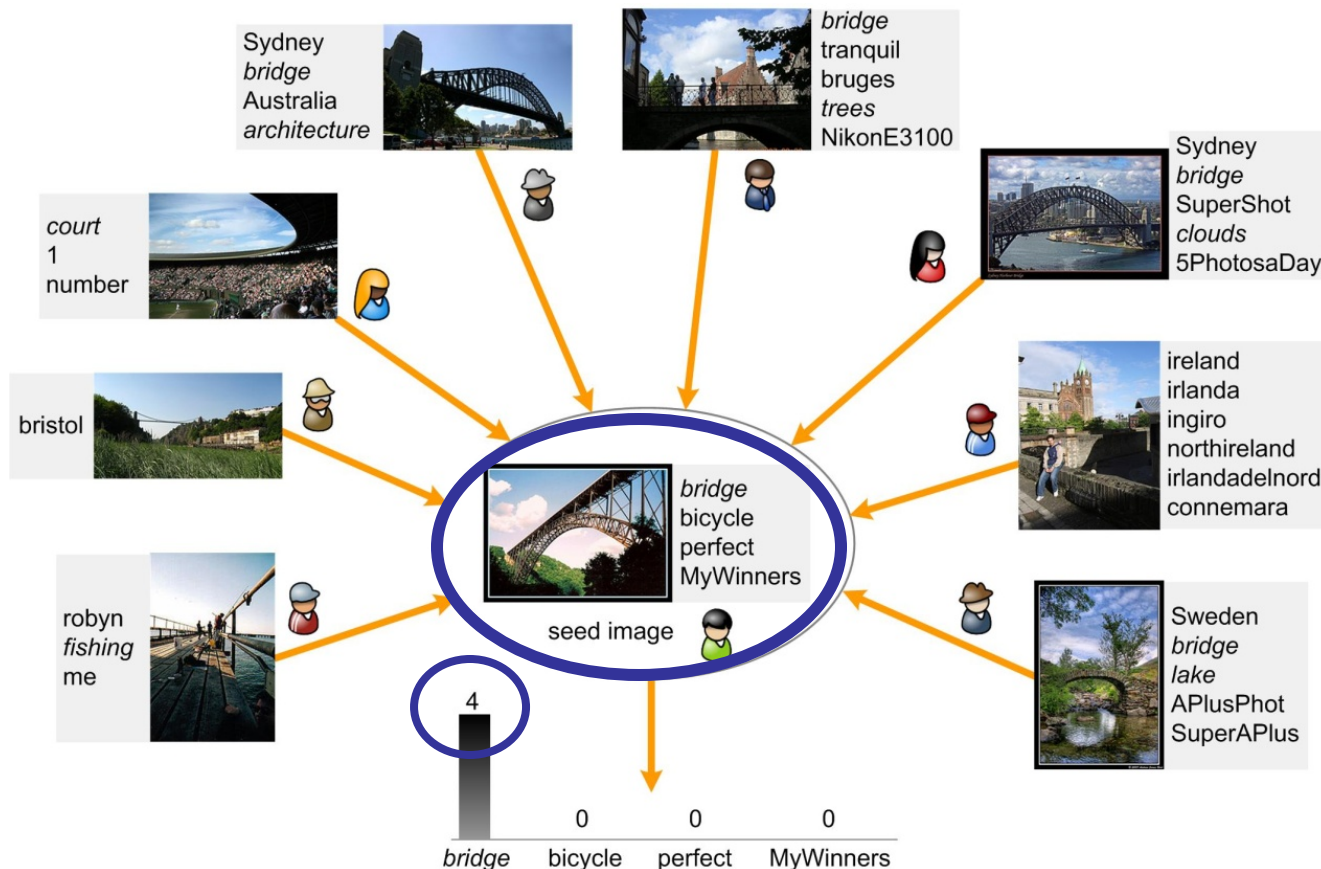
    scalable

      unsupervised

# Nearest neighbor

# Intuition for tagged images

Similar images with similar tags are reliable
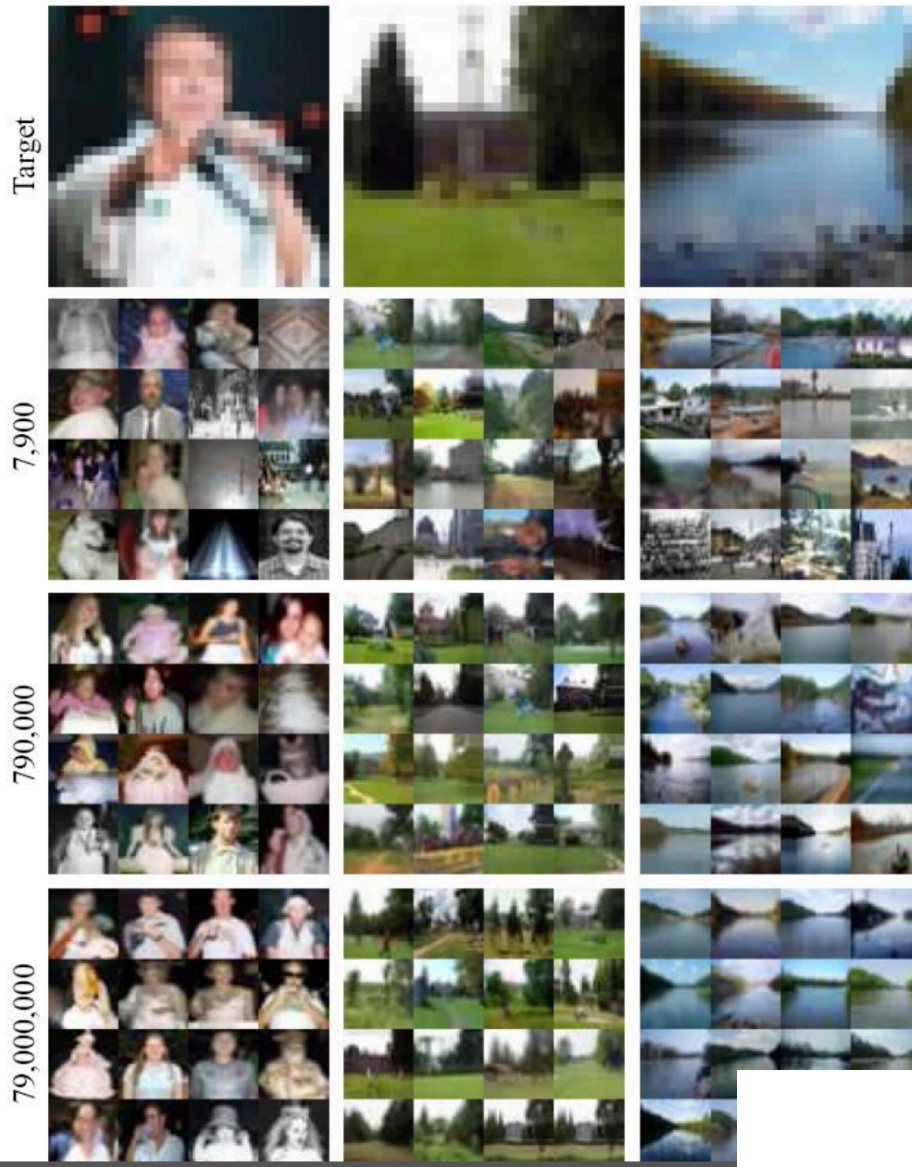
# Nearest neighbor for tag relevance
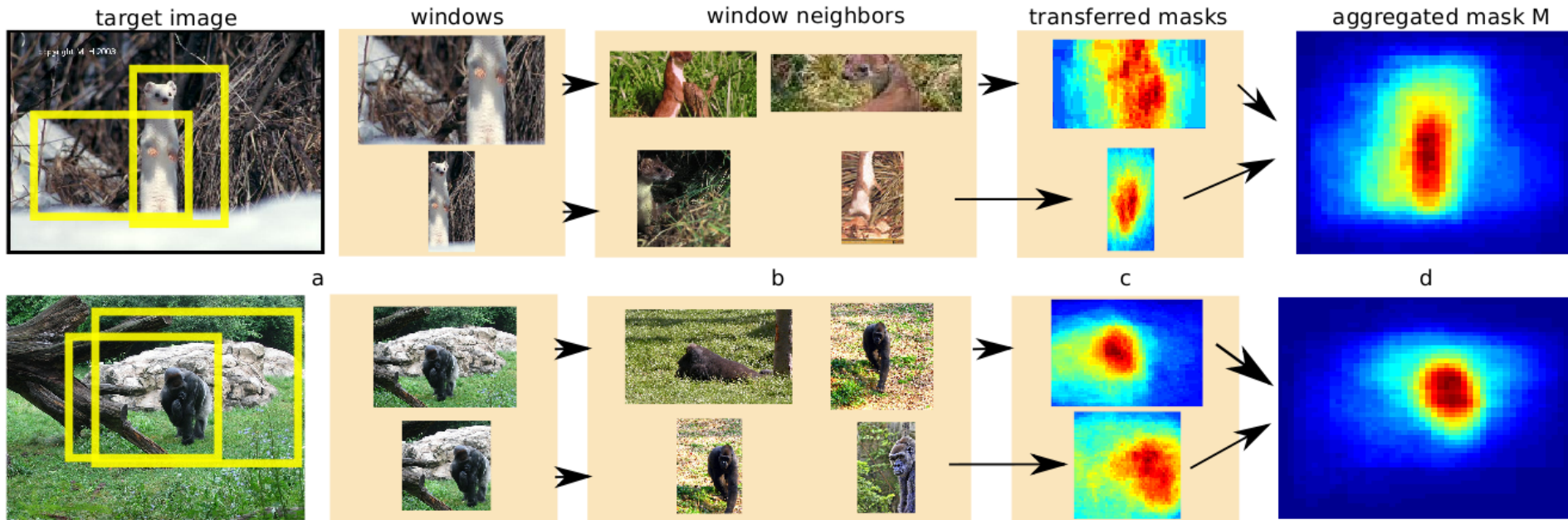
Objective tags are identified



Based on 3.5 Million images downloaded from Flickr

# Even more efficient with tiny images
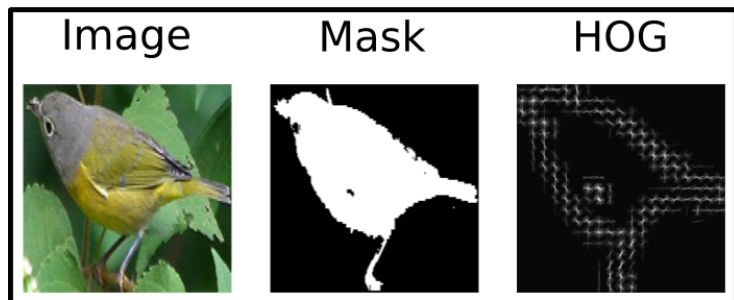


32x32 resolution

80M images

Nearest neighbor

Torralba, PAMI 2008

# Nearest neighbor for segments



target image     windows     window neighbors     transferred masks     aggregated mask M

Annotates many classes with accurate segmentations

Scales efficiently

Segmentations available

Kuettel, ECCV 2012, best paper

# Nearest neighbor for parts



Gavves, ICCV 2013

# Nearest neighbor localized actions?

Write paper.

# Take home message

Nearest neighbor with simple visual features provides a free, scalable and effective means to collect valuable data for many computer vision by learning problems.

# 5. Weakly-supervised vision

In this Chapter we consider computer vision by weakly-supervised learning. In such scenarios some limited supervision is available at train time, typically an object or action class label. The goal is then to enrich this label, for example by predicting bounding boxes, segments or spatio-temporal tubes.

# Weakly-supervised object detection

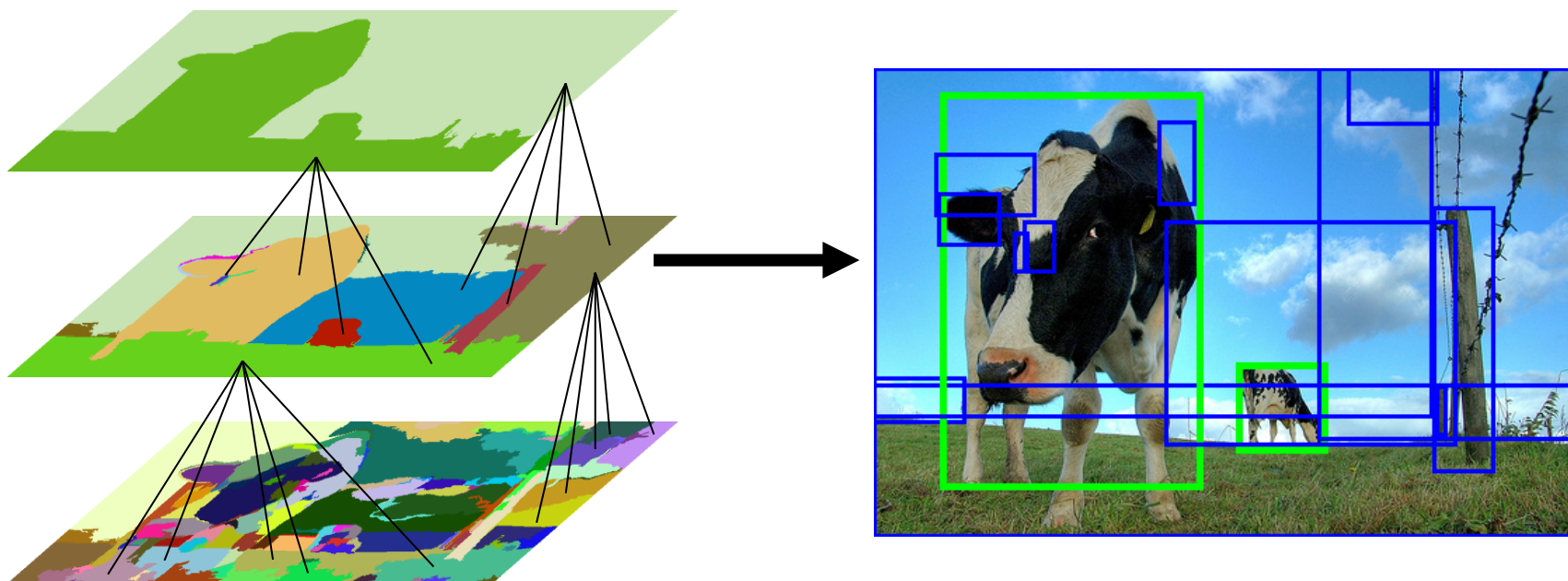Typically casted as Multiple Instance Learning problem

Each image is considered as a "bag" of examples given by object proposals.

Positive images are assumed to contain at least one positive object proposal

The object detector is obtained by alternating detector training, and using the detector to select the single most likely object instance in each positive image.

# Object proposals

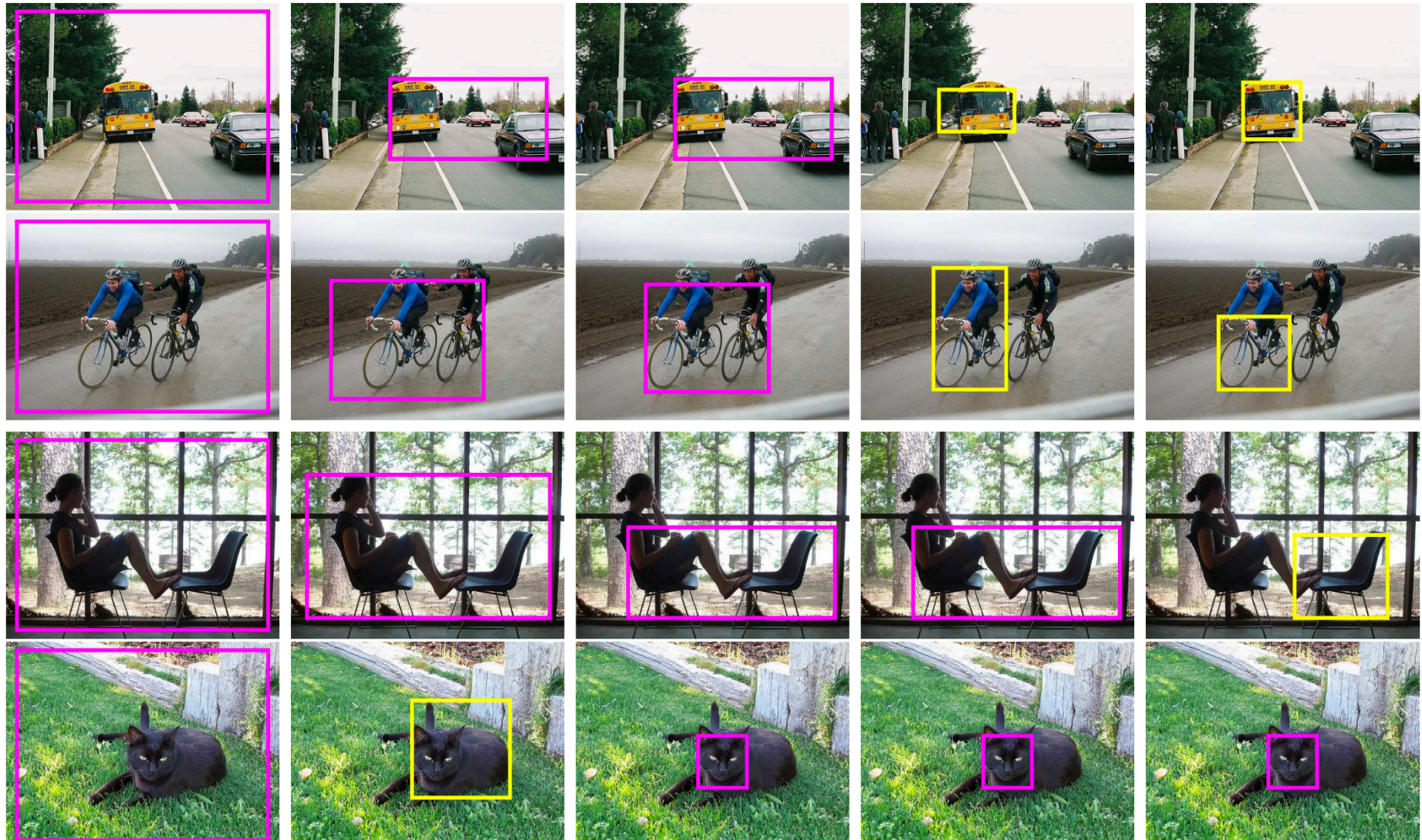Hypotheses from hierarchical grouping of super-pixels
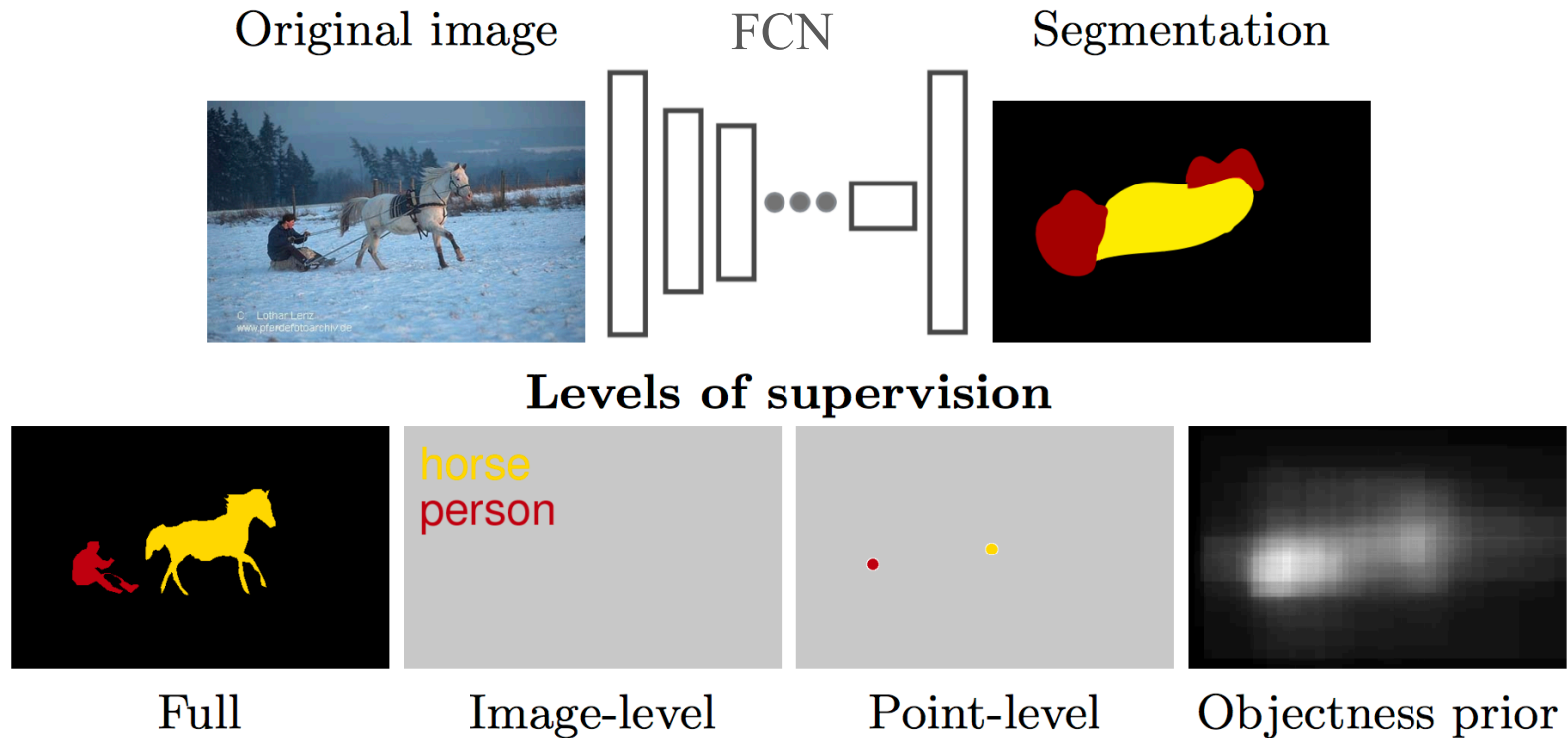
# Multi-fold multiple instance learning

**Algorithm 1** — Multi-fold weakly supervised training

1. Initialization: positive and negative windows are set to entire images up to a 4% border.

2. For iteration $t = 1$ to $T$

   (a) Divide positive images randomly into $K$ folds.

   (b) For $k = 1$ to $K$

      i. Train using positives in all folds but $k$.

      ii. Re-localize positives in fold $k$ using this detector.

   (c) Train detector using positive windows from all folds.

   (d) Perform hard-negative mining using this detector.

3. Return final detector and object windows in train data.

Cinbis *et al.* CVPR 2014 / PAMI 2017

# Re-localization process

# Weakly-supervised segmentation



Original image     FCN     Segmentation

**Levels of supervision**

horse
person

Full     Image-level     Point-level     Objectness prior

Adapt loss function depending on supervision scheme
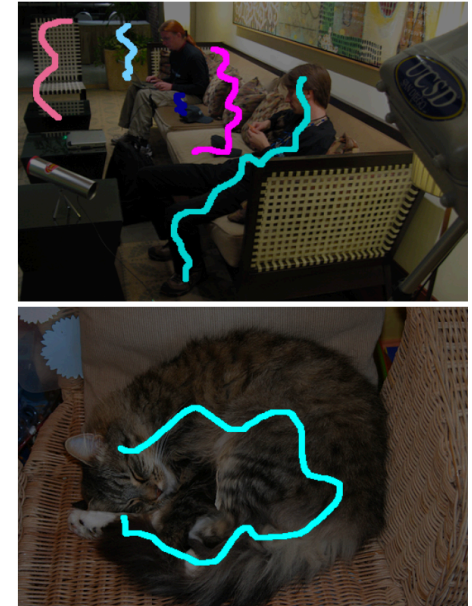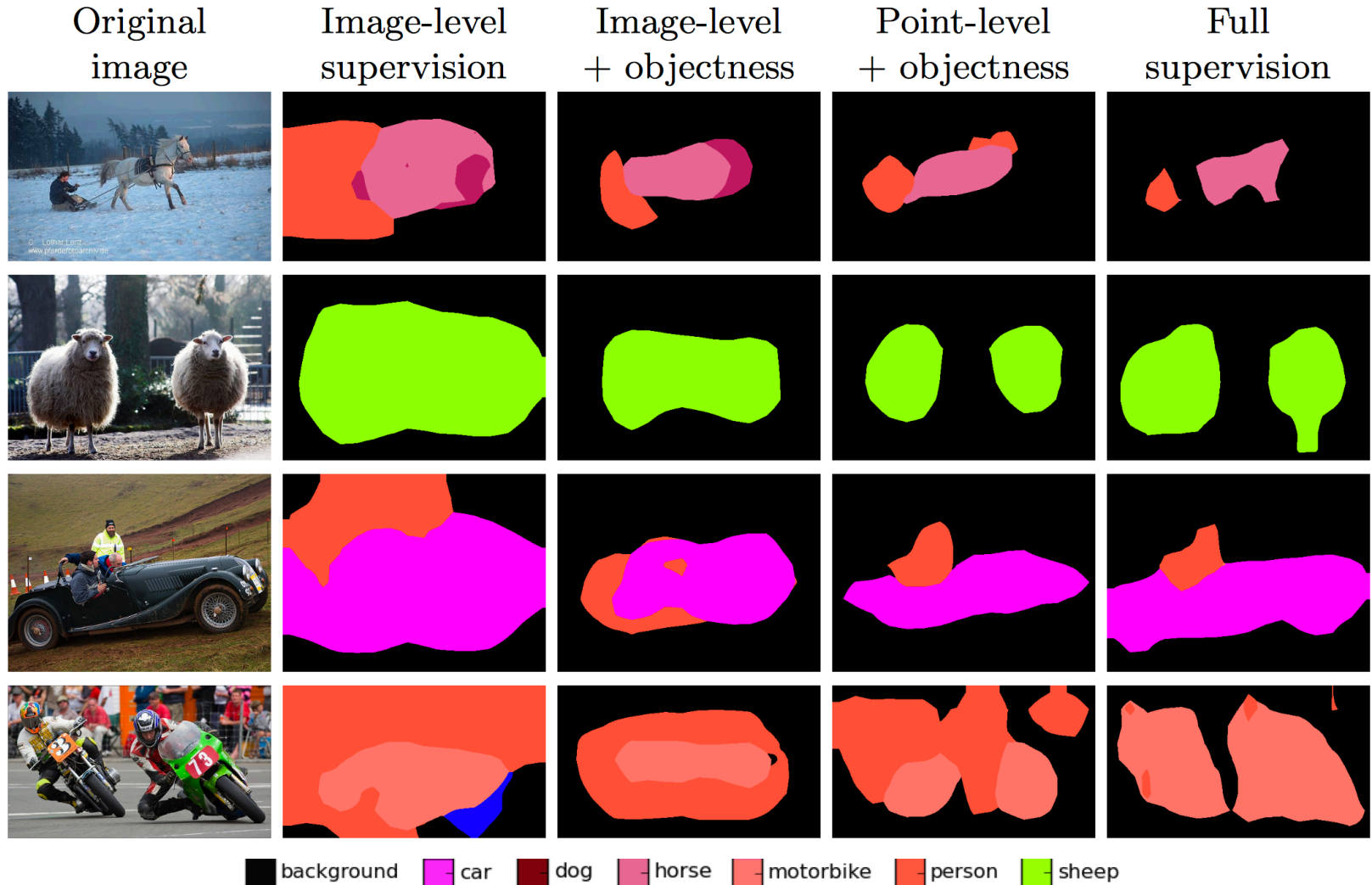
# Crowdsourcing point annotations



Image-level labels: 20.0 sec/image

Points: **22.1** sec/image

Squiggles: 34.9 sec/image
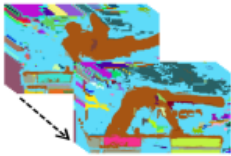
Full supervision: **239.7** sec/image

# Some results



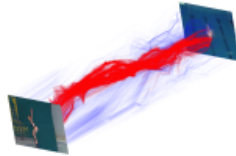| | Original image | Image-level supervision | Image-level + objectness | Point-level + objectness | Full supervision |
|---|---|---|---|---|---|

Legend: background, car, dog, horse, motorbike, person, sheep

Bearman *et al.* ECCV 2016

# Recap: Action proposals

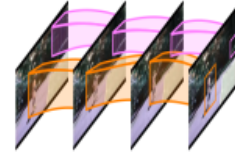| **Supervoxels** | **Trajectories** | **Detect & Track** |
|---|---|---|
| Jain et al. *CVPR'14*<br>Oneata et al. *ECCV'14* | van Gemert et al. *BMVC'15*<br>Puskas et al. *ICCV'15* | Yu et al. *CVPR'15*<br>Weinzaepfel et al. *ICCV'15* |

**Action proposals**

. . . .          . . . .

# Action localization with proposals

**At train time**

Annotate spatiotemporal tubes with class labels

Extract video representation from tubes

Train favorite classifier

**At test time**

Extract action proposals

Extract video representation from each proposal

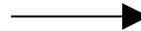Classify all proposals, select proposal with maximum response

# Hypothesis

Training on bounding boxes not required.
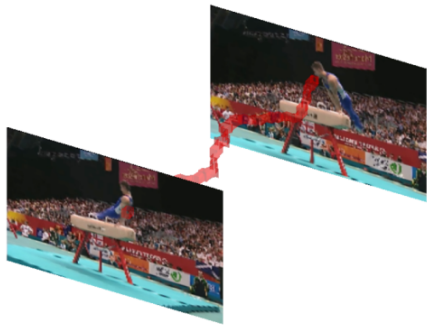Training on proposals with fast point annotations is as effective.
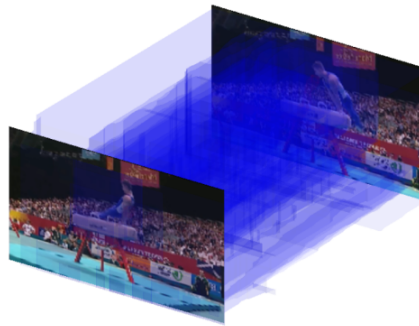


*Annotation time for video:*
5 min. 11 sec.

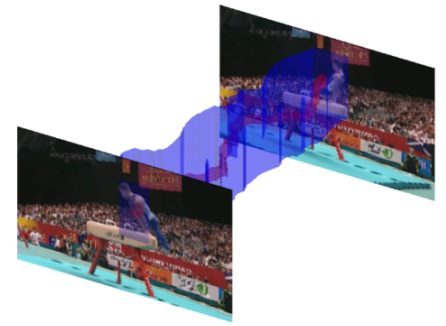*Annotation time for video:*
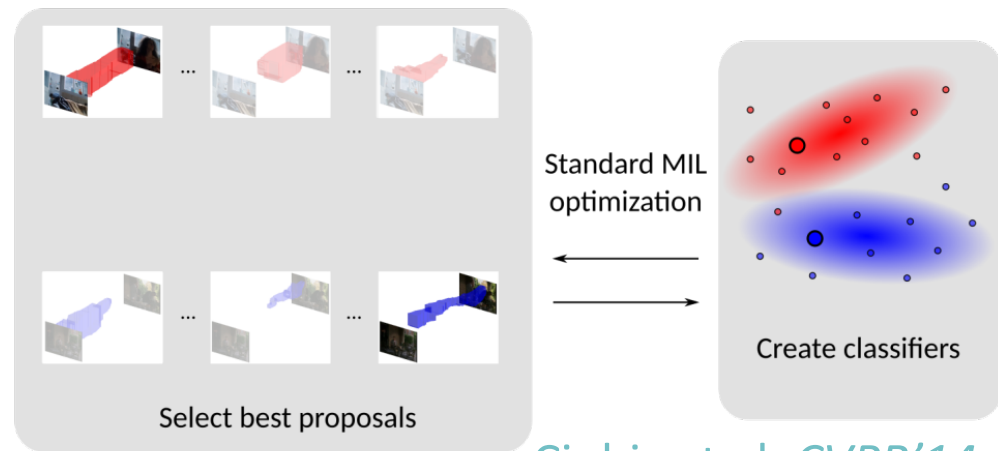25 sec.

# Idea



Human point supervision      Compute proposal affinity      Mine best proposal

Mettes *et al.* ECCV 2016
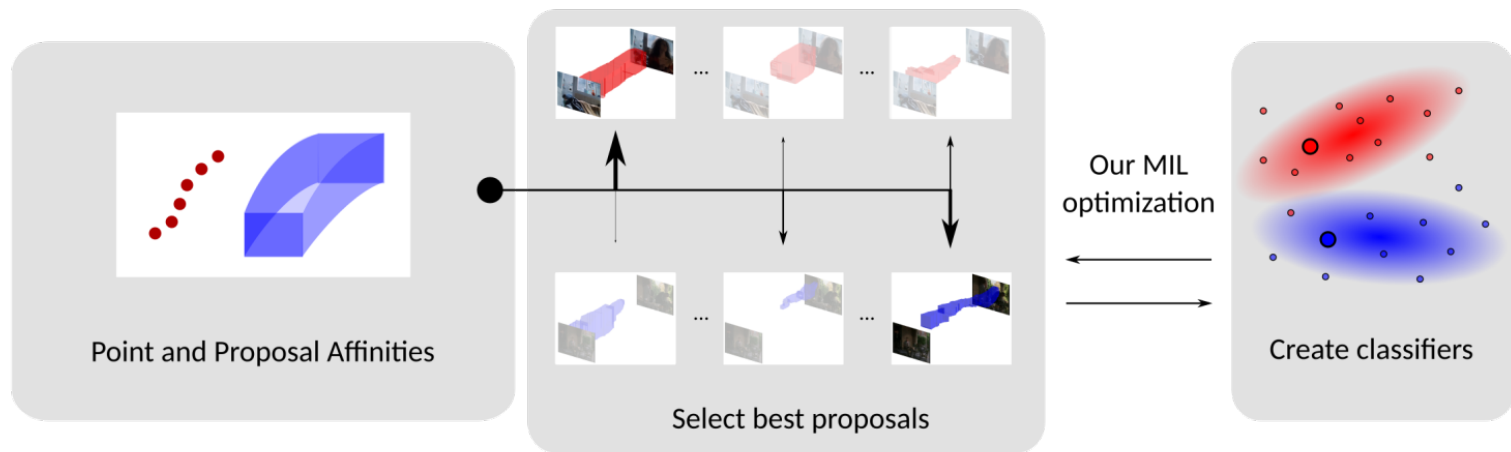
# Mining the best proposal

Train action classifiers using best proposals only.
Casted as a Multiple Instance Learning problem.



Cinbis et al. *CVPR'14*

# Mining the best proposal

Train action classifiers using best proposals only.
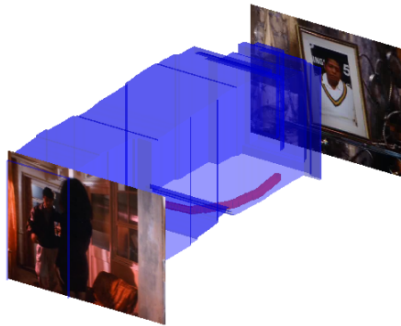Casted as a Multiple Instance Learning problem.



Use affinity with point annotations to guide the mining.
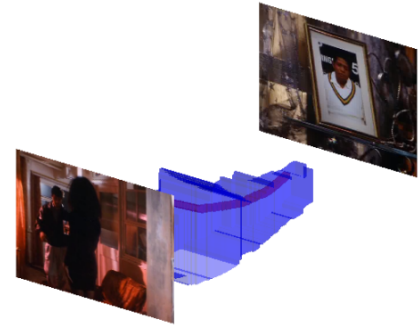
# Proposal affinity

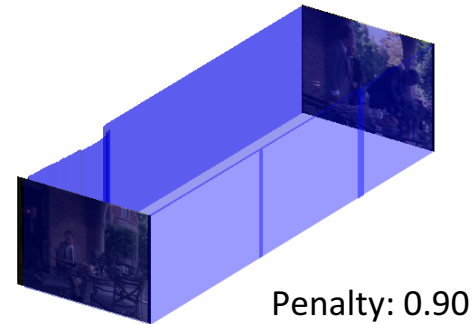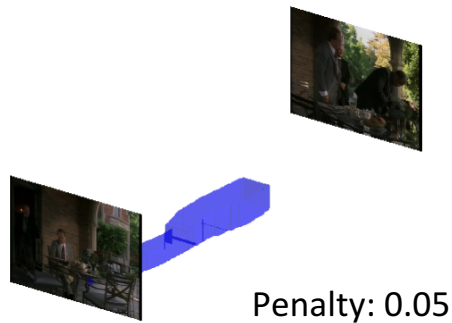Novel overlap measure between point annotations and proposals.



No overlap          Small overlap          High overlap
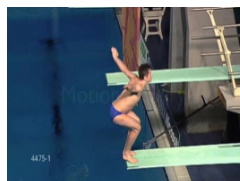
# Mind the center bias

Subtract the size of the proposal from the match.
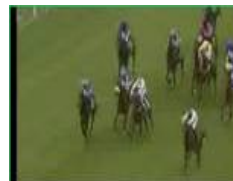To alleviate center bias of large proposals.



Penalty: 0.05
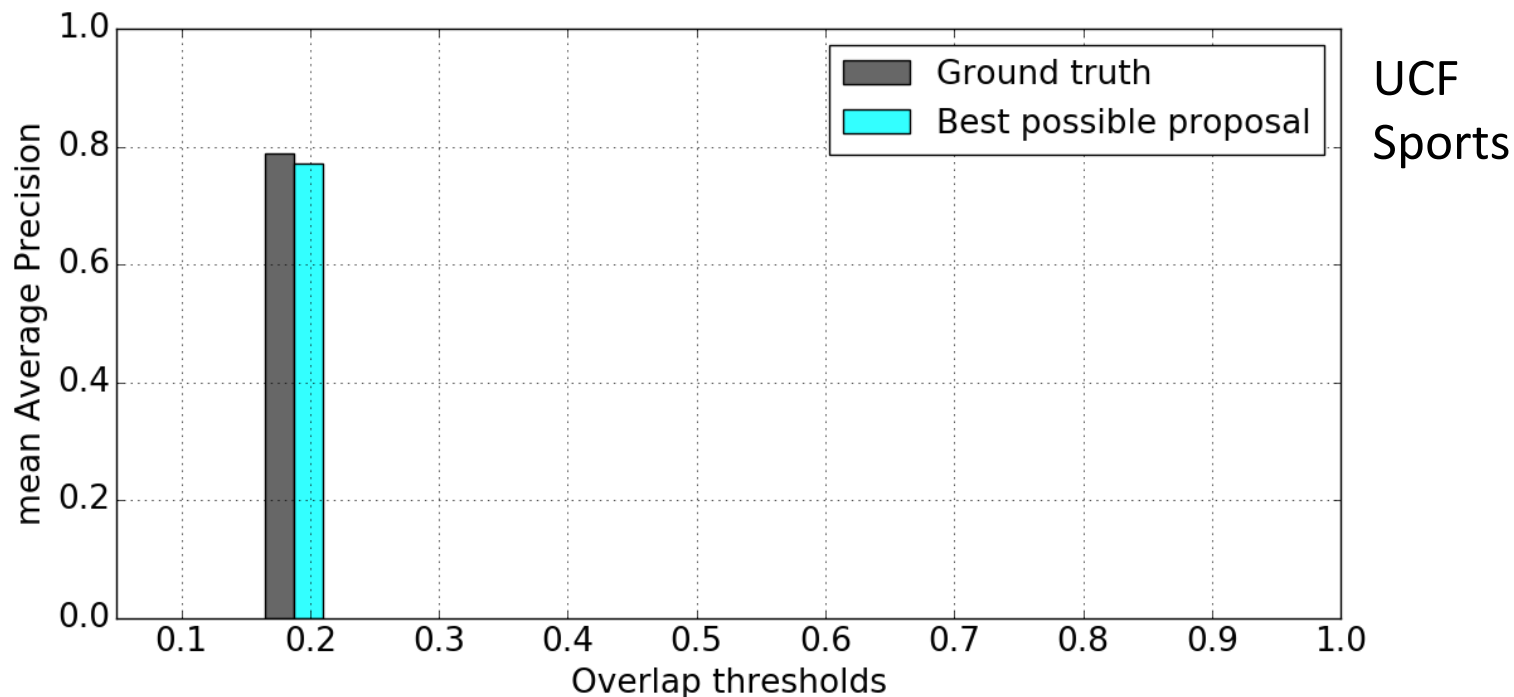
Penalty: 0.90

# Experiments

**UCF Sports**



**UCF 101** (in paper)



Unsupervised proposals from clustered trajectory features.
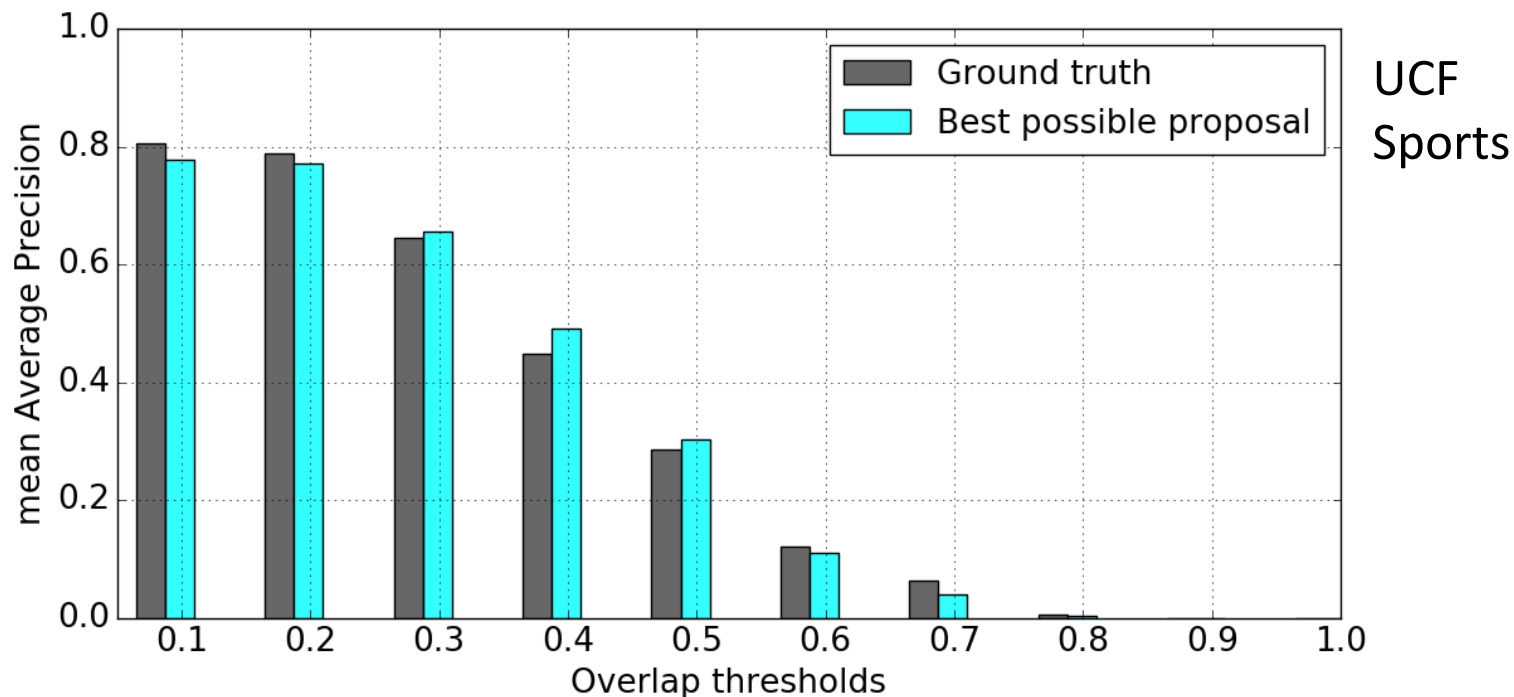Evaluated with Fisher Vectors and SVMs.

van Gemert *et al.* BMVC'15
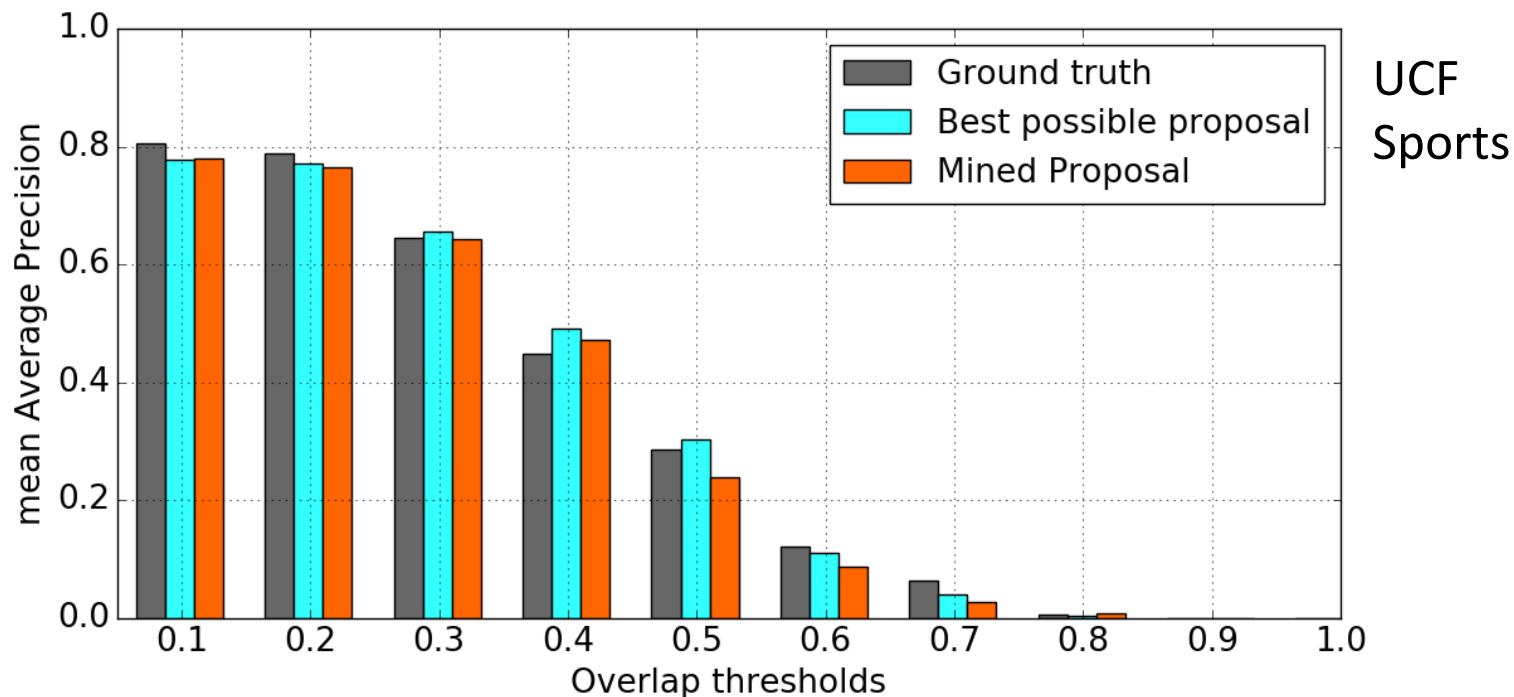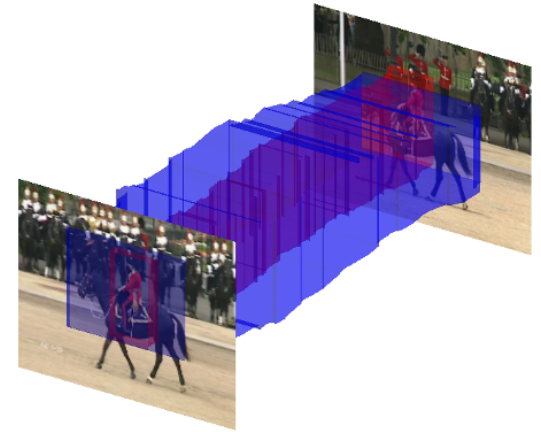
# Training without ground truth boxes



UCF Sports

Best possible proposal performs as well as ground truth boxes.

# Training without ground truth boxes
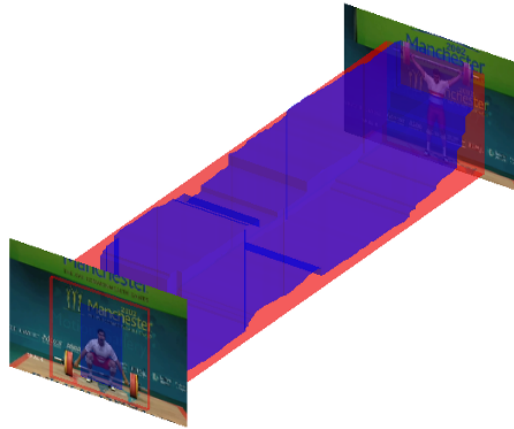


UCF
Sports

Best possible proposal performs as well as ground truth boxes.

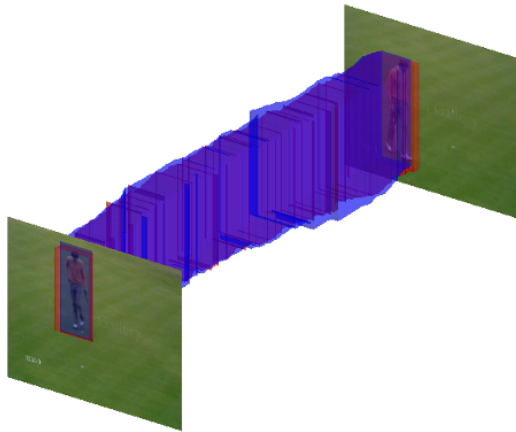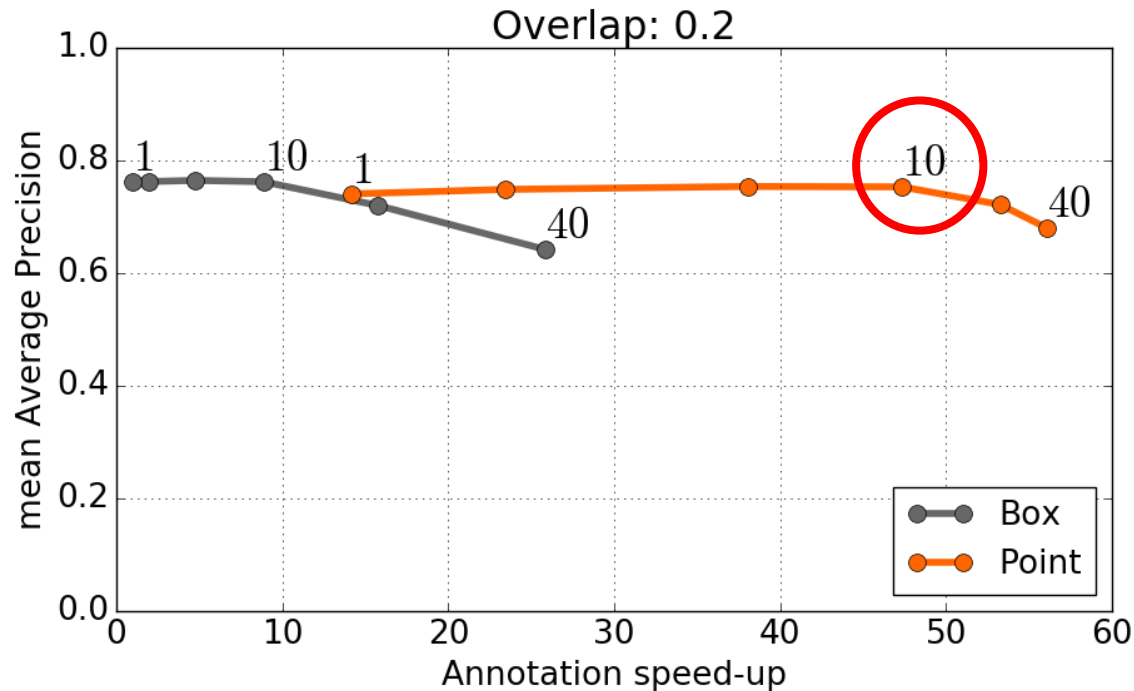# Training without ground truth boxes



UCF Sports

Mean AP maintained using our mined proposals.

# Qualitative results



Ground truth boxes
Mined proposal

# Lowering annotation frame-rate



Up to 50 times speed-up at similar performance.

# Hollywood2Tubes

Dataset to demonstrate how easy action annotation becomes. Contains actions and instances new to action localization.



a. Multi-label videos.

b. Contextual actions.

c. Group interactions.

Download:

**tinyurl.com/hollywood2tubes**

# Take home message

Weakly-supervised computer vision is aided by reasonable proposals for objects, segments and/or actions.

Proposals are further refined with point annotations. Especially useful for precise annotations, like segments and actions.

Facilitates dataset construction and/or enrichment.

# Overview

1. Image benchmarks, PASCAL, ImageNet, MSCOCO
2. Video benchmarks, TRECVID, ActivityNet
3. Labels from humans, experts, volunteers, crowdsourcing
4. Labels from similarity, nearest neighbor, simple features
5. Weakly-supervised computer vision

6. Event recognition by learning