

80 Million Tiny Images

*IPAM Workshop on Numerical Tools and Fast Algorithms for
Massive Data Mining, Search Engines and Applications*

October 23rd 2007



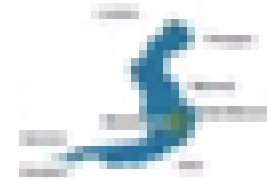
Massachusetts
Institute of
Technology

Antonio Torralba
Rob Fergus
William T. Freeman



Overview

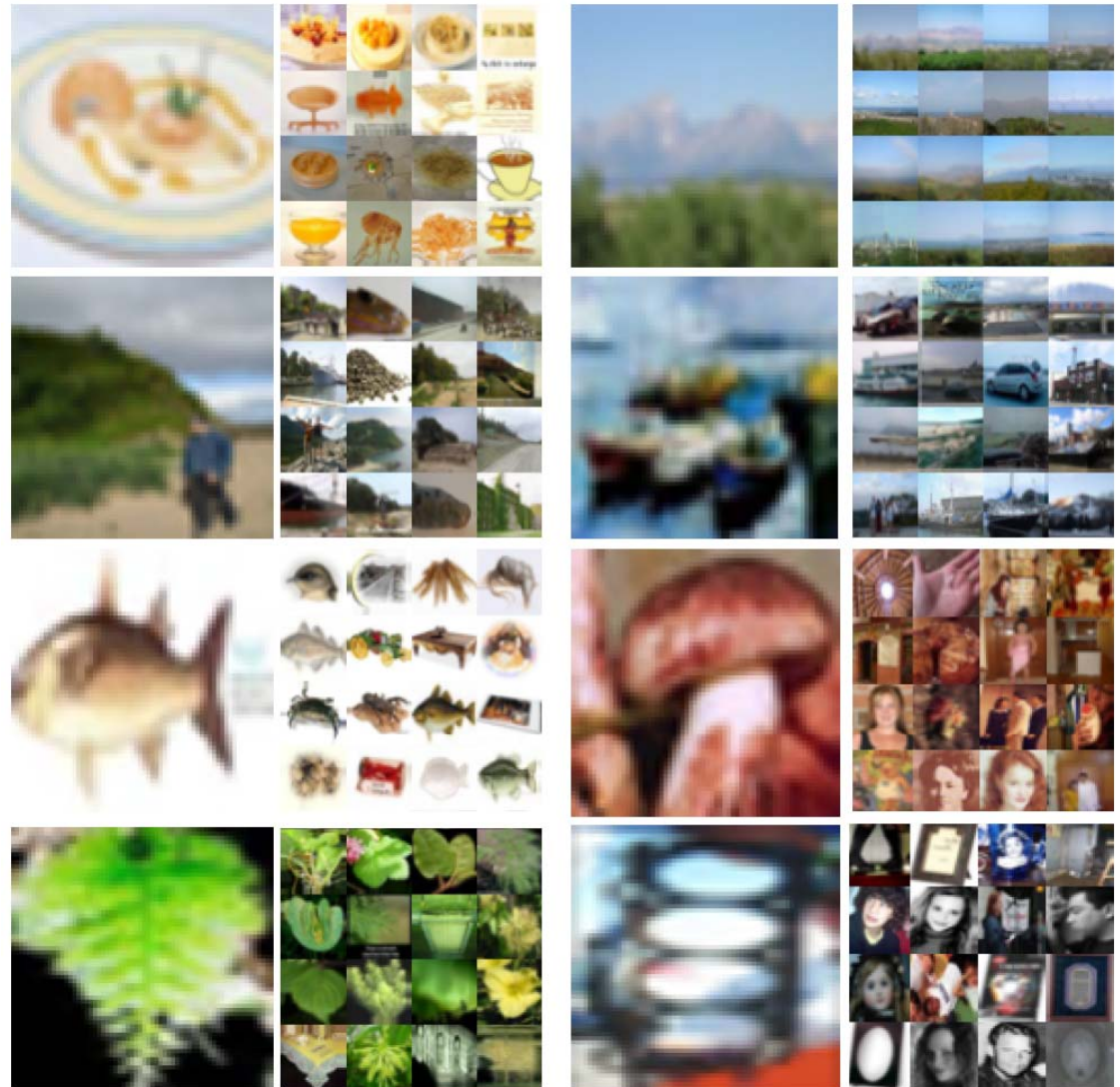
- Non-parametric approach to category-level recognition
- Dataset of 80 million images from Internet



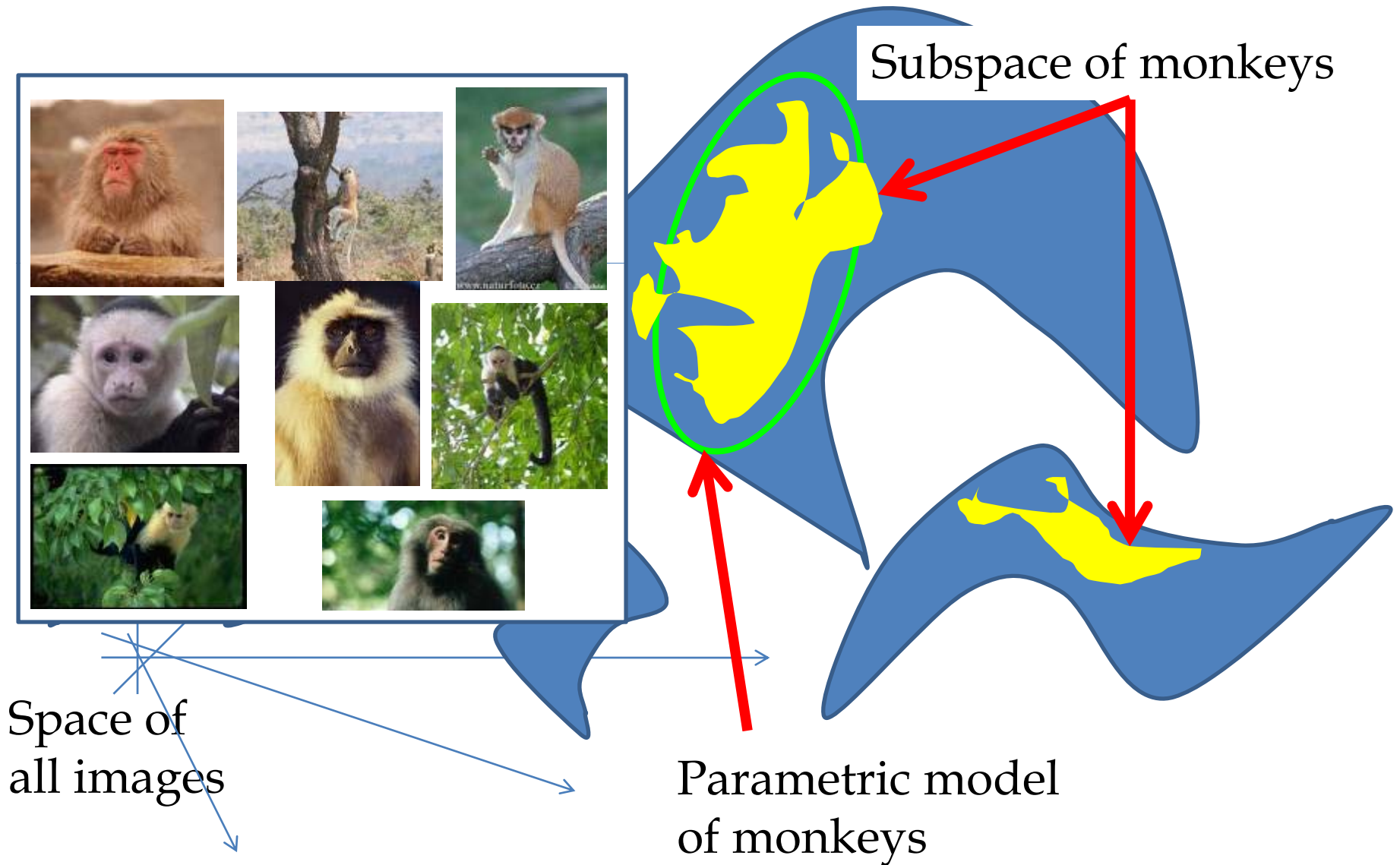
- Use very low resolution images (32x32 color)

Overview

- Use simple algorithms: nearest neighbors



Motivation

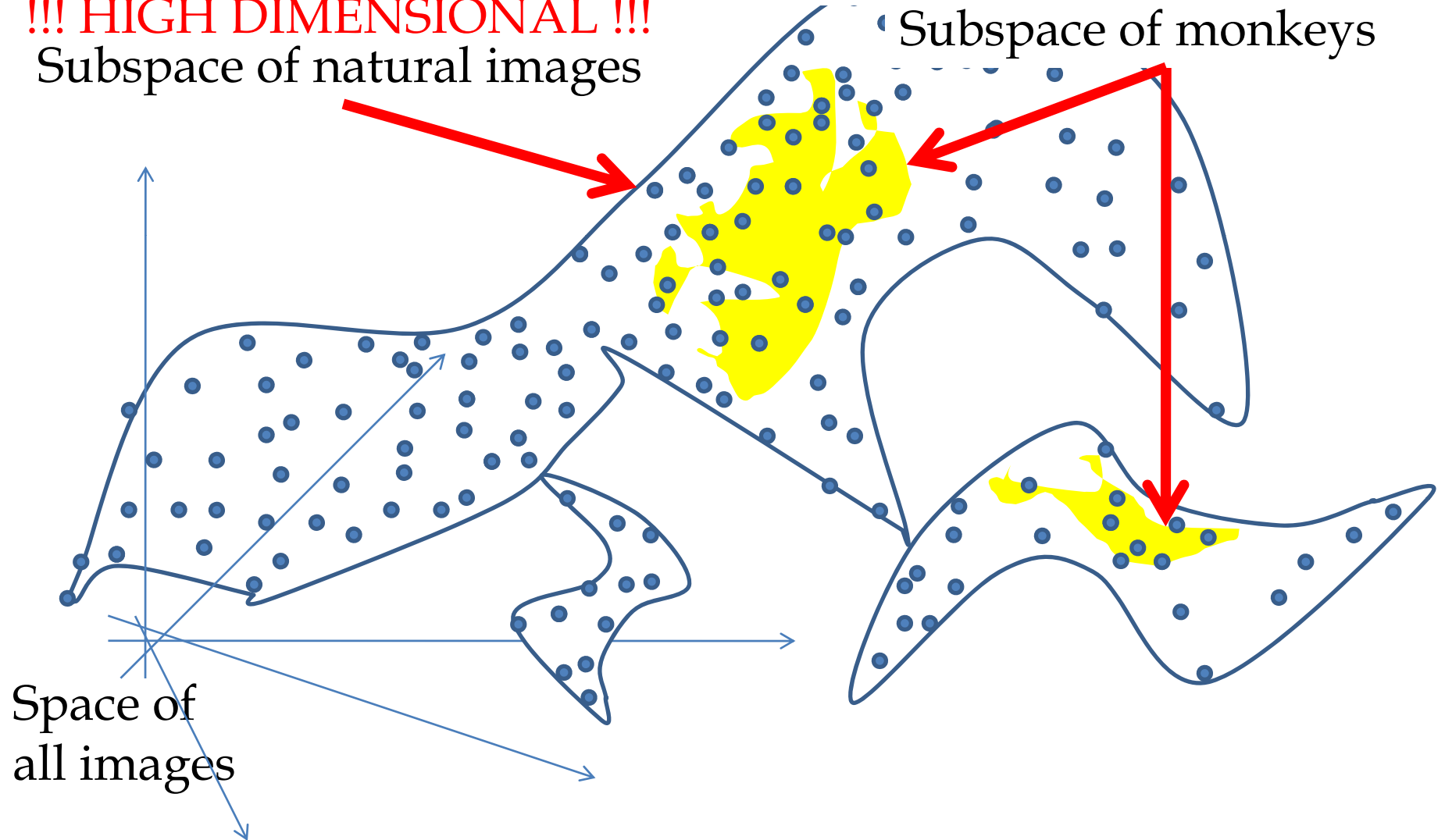


Non-parametric Approach

!!! HIGH DIMENSIONAL !!!
Subspace of natural images

!!! HIGH DIMENSIONAL !!!

Subspace of monkeys

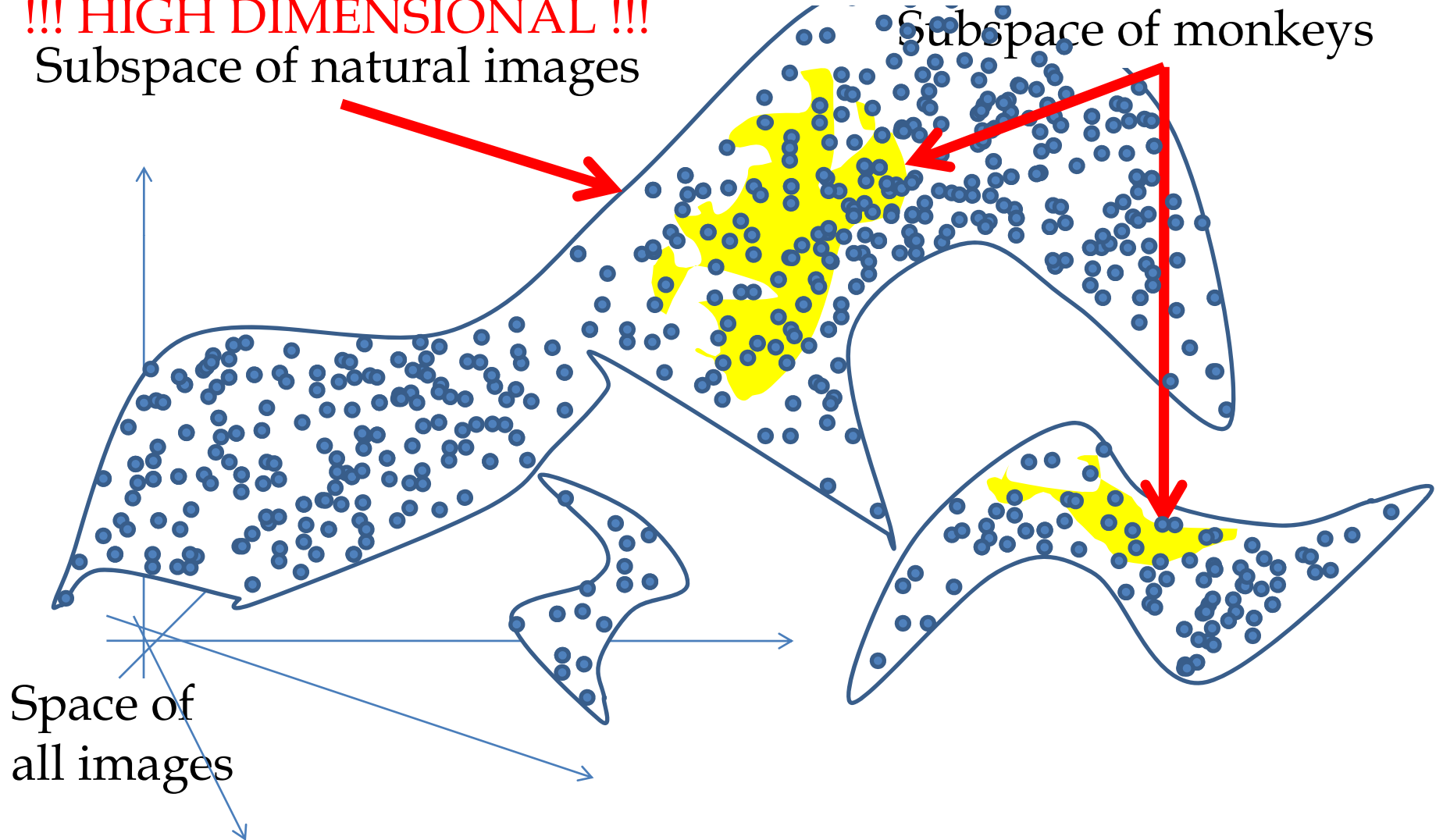


Non-parametric Approach

!!! HIGH DIMENSIONAL !!!
Subspace of natural images

!!! HIGH DIMENSIONAL !!!

Subspace of monkeys



Space of
all images

The Data

Thumbnail Collection Project

- Collect images for ALL objects
 - List obtained from WordNet
 - 75,378 non-abstract nouns in English
- Example first 20:

a-bomb
a-horizon
a._conan_doyle
a._e._burnside
a._e._housman
a._e._kennelly
a.e.
a_battery
a_cappella_singing
a_horizon

a_kempis
aalborg
aalii
aalost
aalto
aar
aardvark
aardwolf
aare
aare_river

Thumbnail Collection

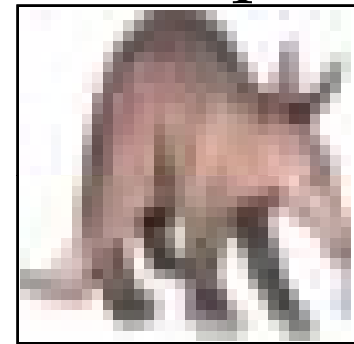
- 7 different search engines



Dataset Statistics

- Overall stats
 - 79,302,017 images
 - 75,062 different words
- Details
 - Two formats: square & rectangular
 - Gathered at 4.5 images/second
 - Downloaded 97,245,098 images
 - 18% duplicate rate
 - Disk usage: ~ 700Gb
 - Collection time: ~ 9 months

32x32 square



32xN rectangular

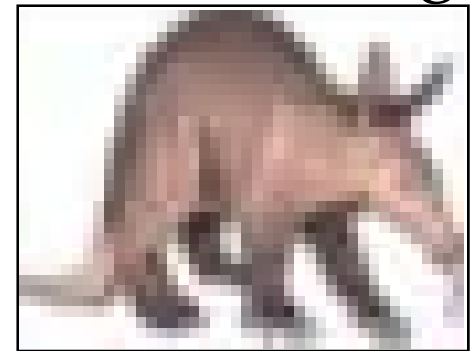
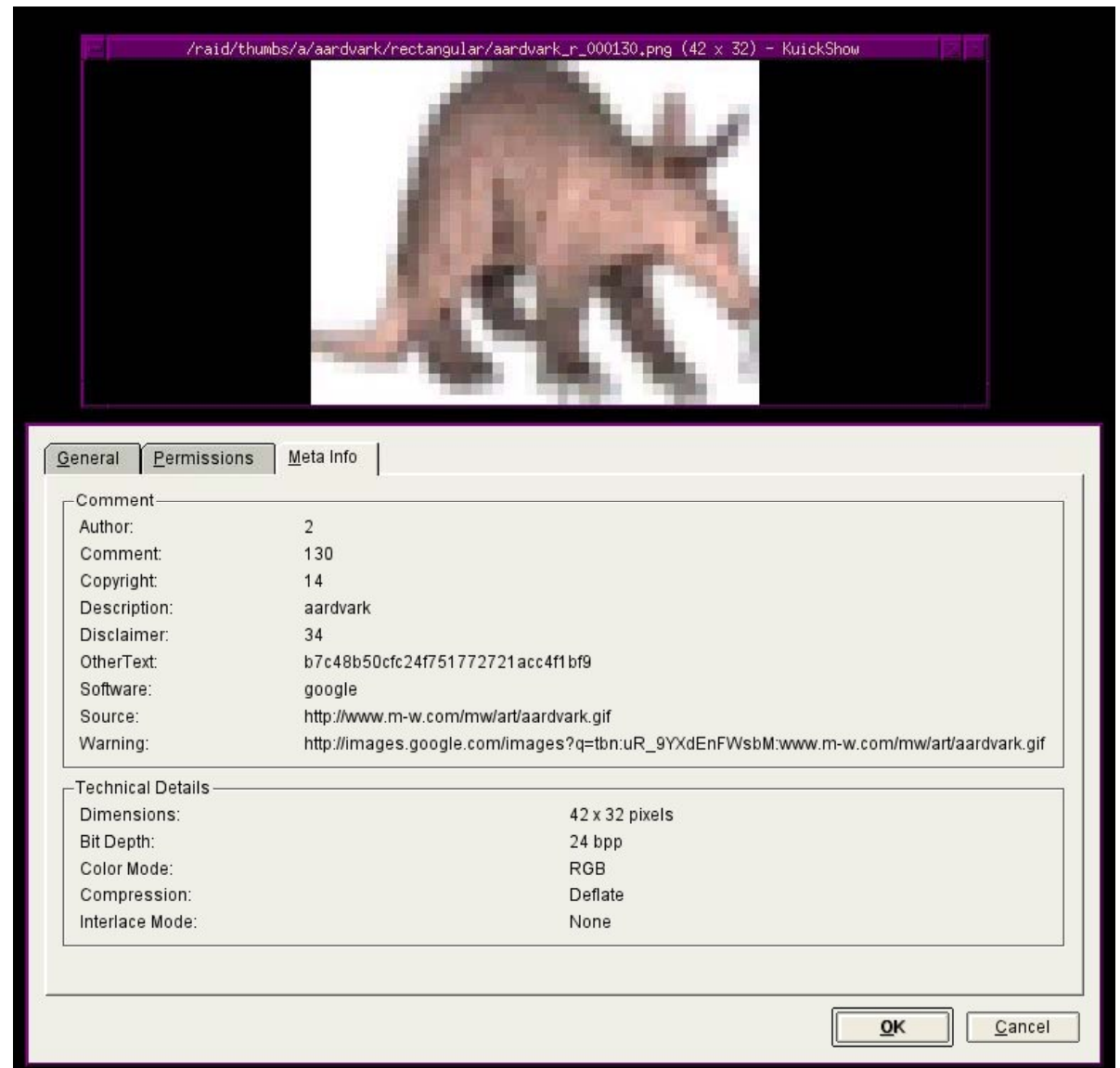
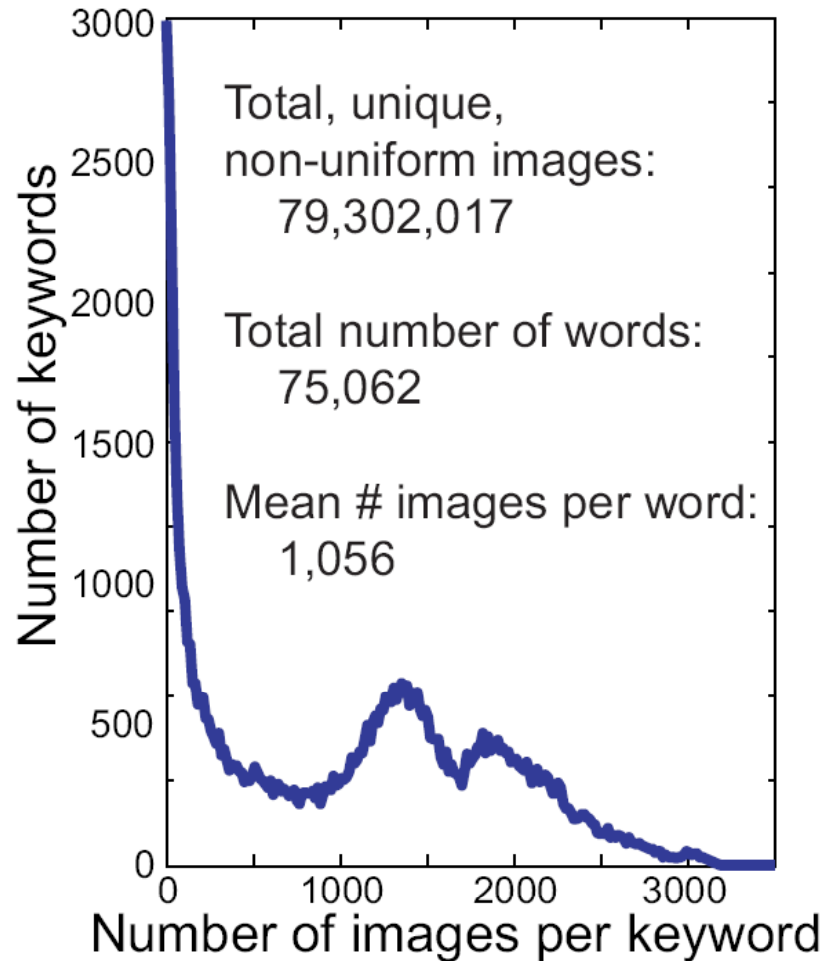


Image Metadata

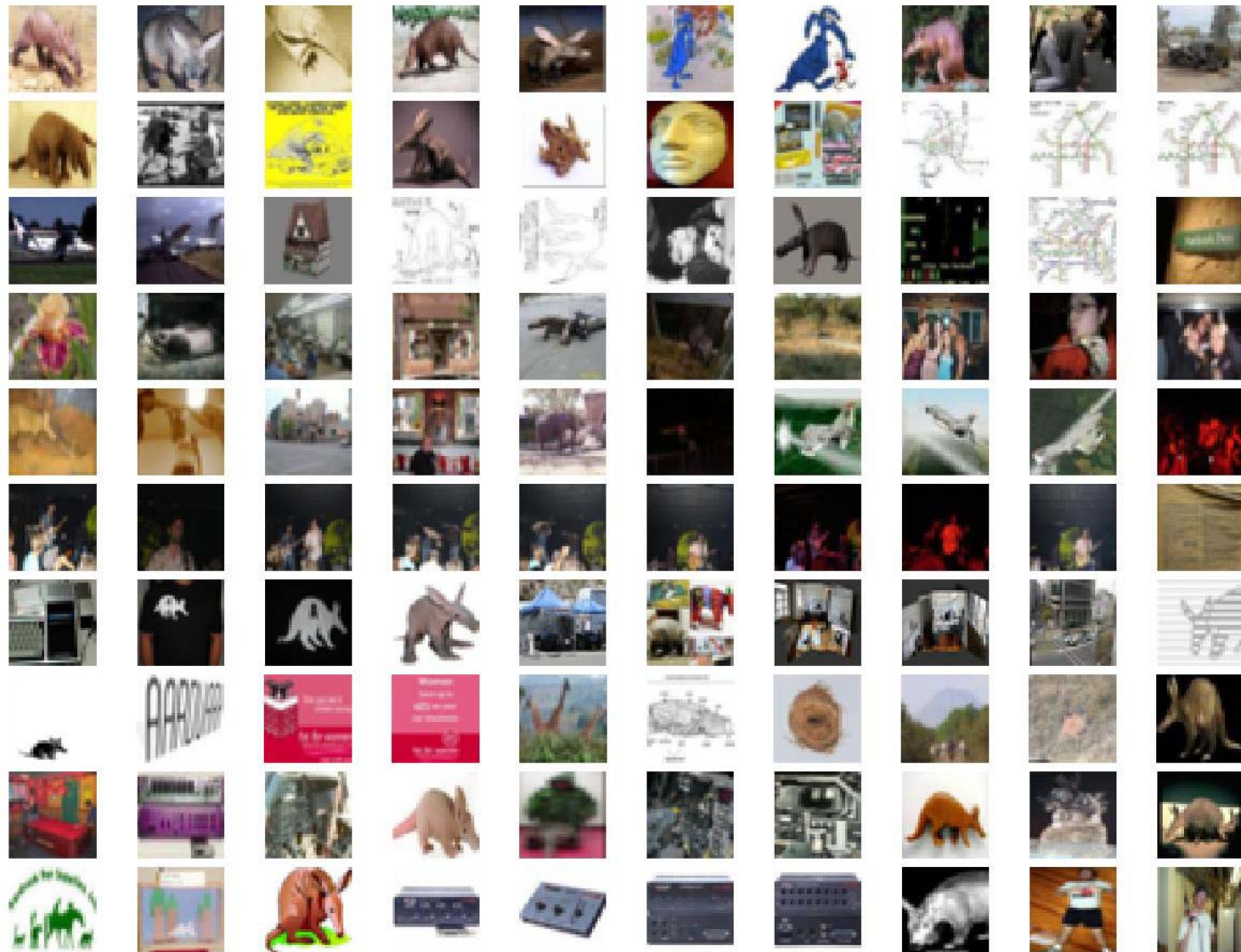
- URL to high-res
- URL of thumbnail
- Engine & Rank



Histogram Images/Word

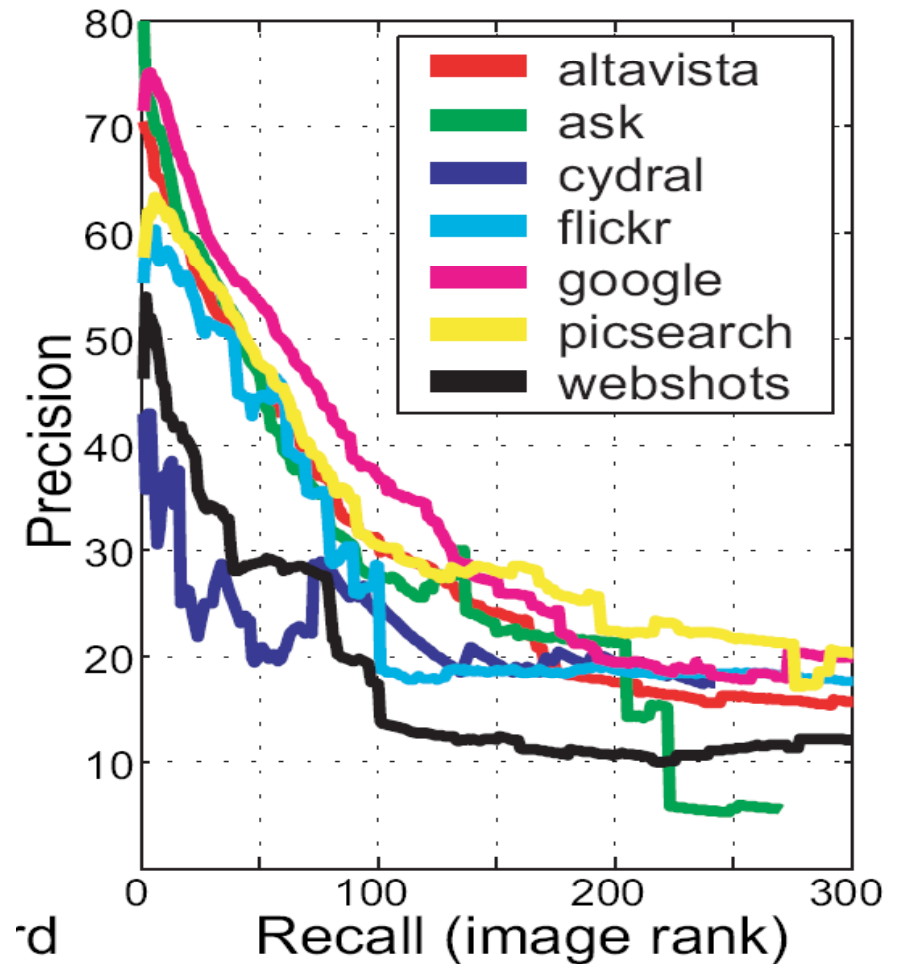


Aardvark Images



Labeling Noise

- Manual labeling of 78 classes
- Best: Google & Altavista
- Worst: Cydral & Webshots

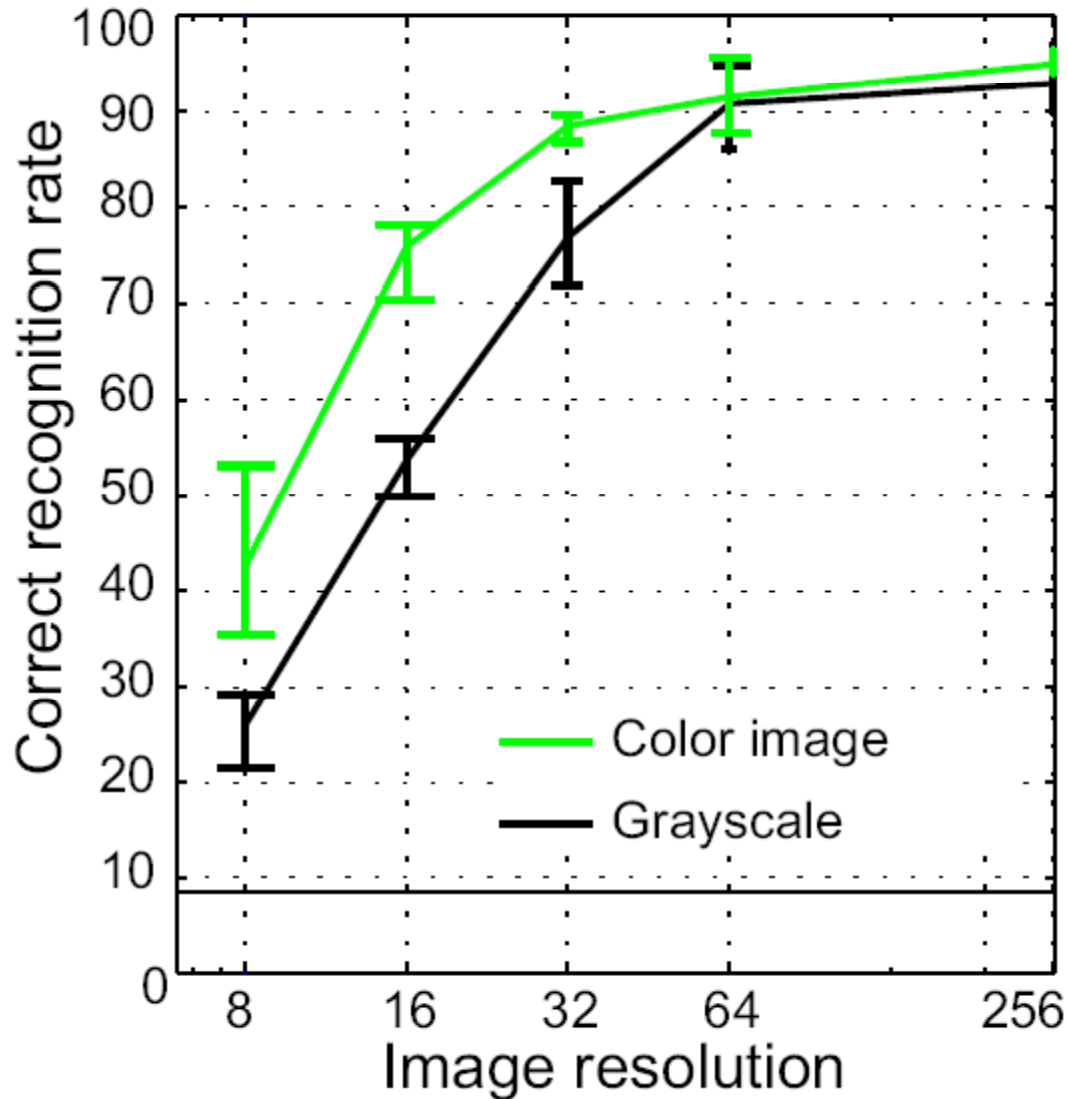


Representation

Suitable Image Representation

- Want minimal representation for task:
 - Classifying scene and dominant objects
- Compact representation has low storage requirements
- We blur & subsample to give low-res image (32x32 color)

Human Performance at Scene Recognition



The role of context in object recognition
A. Oliva, A. Torralba
Trends in Cognitive Sciences, in press.
December 2007.

Human Labeling of Tiny Scenes



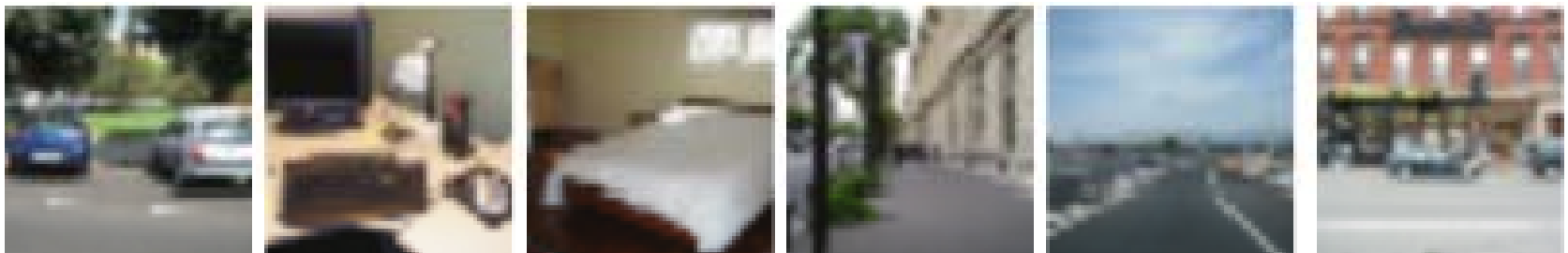
Image Patches vs Tiny Images



a) Patches



b) Object chips



c) Scenes, small thumbnails

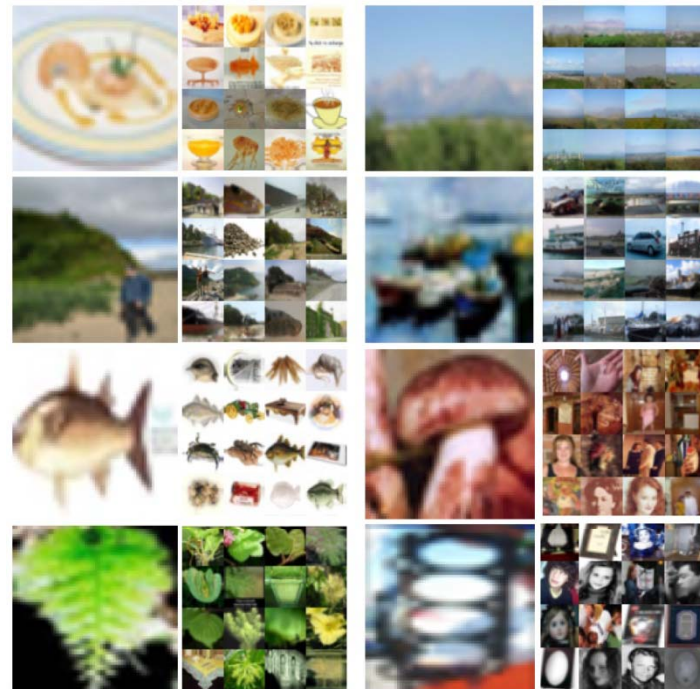
Approach

Non-parametric Classifier

- Nearest-neighbors
- For each query, obtain **sibling set** (neighbors)

- 3 different types of distance metric

- Hand-designed, use whole image



Metric 1 - D_{ssd}

- Sum of squared differences (SSD)

$$D_{ssd}^2 = \sum_{x,y,c} \left[\text{Image 1} - \text{Image 2} \right]^2$$

To give invariance to illumination:
Each image normalized to
be zero mean, unit variance



Target

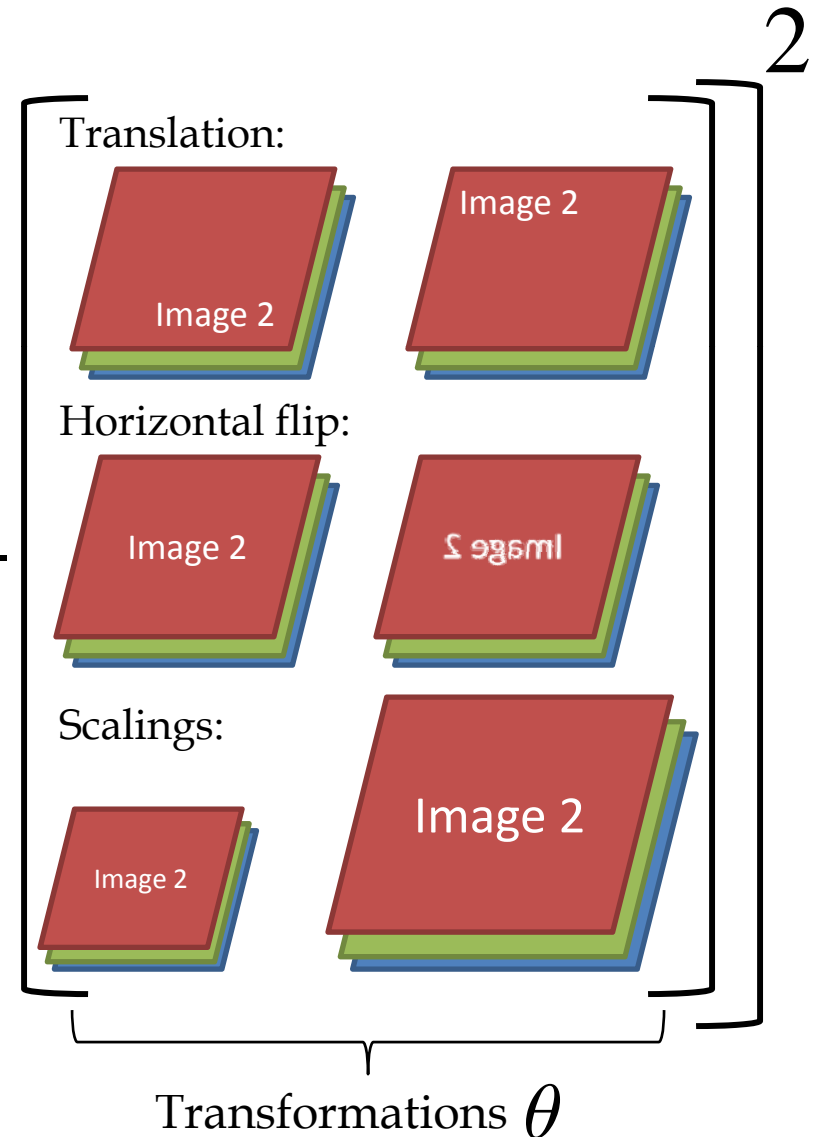
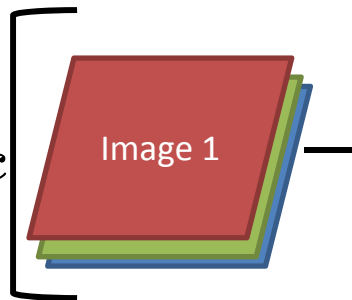


Neighbor

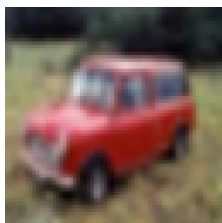
Metric 2 - D_{warp}

- SSD but allow small transformations

$$D_{warp}^2 = \min_{\theta} \sum_{x,y,c}$$



Find min using gradient descent



Target



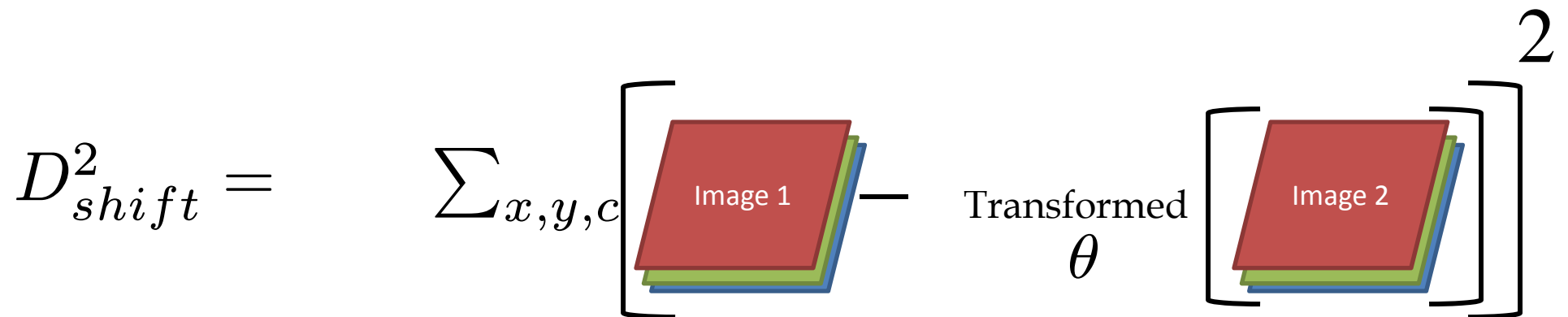
SSD



Warping

Metric 3 - D_{shift}

- As per Warping but also allow sub-window shifts

$$D_{shift}^2 = \sum_{x,y,c} \left[\text{Image 1} - \text{Transformed Image 2} \right]^2$$


Start with warped version of image 2, as per D_{warp}

Metric 3 - D_{shift}

- As per Warping but also allow sub-window shifts

$$D_{shift}^2 = \sum_{x,y,c} \left[\begin{array}{c} \text{Image 1} \\ \text{Image 2} \end{array} - \text{Transformed } \theta \left[\begin{array}{c} \text{Image 2} \\ \text{Image 2} \end{array} \right] \right]^2$$

Start with warped version of image 2, as per D_{warp}

Metric 3 - D_{shift}

- As per Warping but also allow sub-window shifts

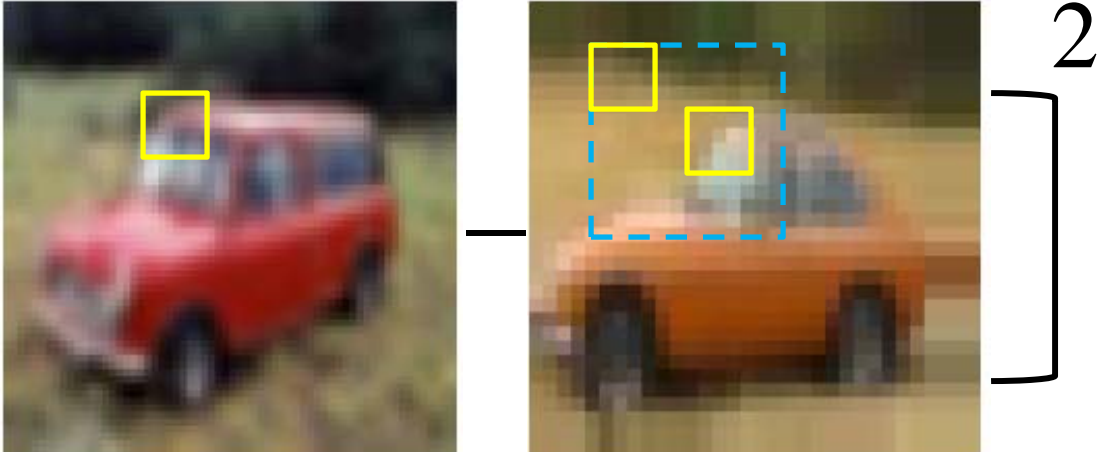
$$D_{shift}^2 = \sum_{x,y,c} \left[\begin{array}{c} \text{Image 1} \\ \text{Image 2} \end{array} \right]^2$$

The equation shows the squared L2 distance between two images, summed over all pixels (x, y, c) . The first image is a red car, and the second image is an orange car. The images are shown in a vertical stack within large square brackets, with a minus sign between them, and a superscript 2 to the right of the closing bracket.

Start with warped version of image 2, as per D_{warp}

Metric 3 - D_{shift}

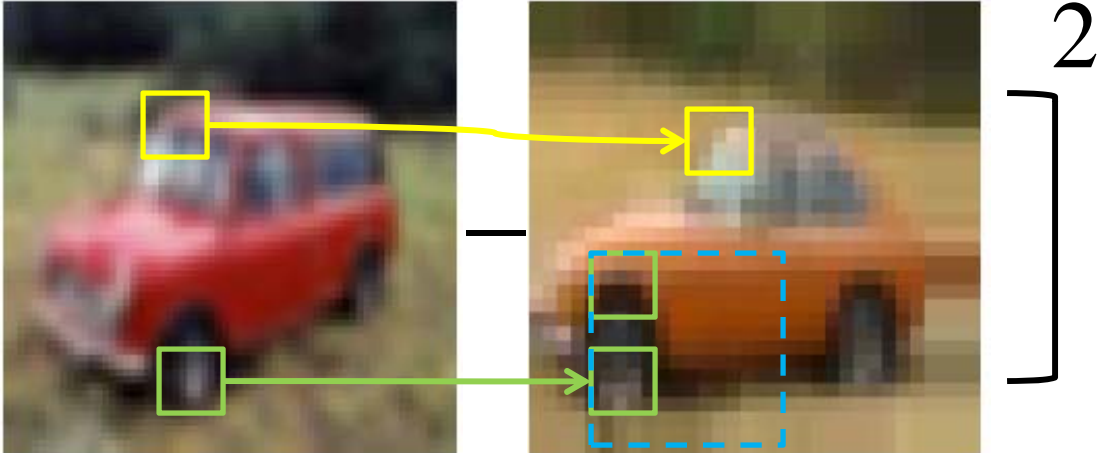
- As per Warping but also allow sub-window shifts

$$D_{shift}^2 = \min_{\text{Local sub-window}} \sum_{x,y,c}$$


The diagram illustrates the D_{shift} metric. It shows two images of a car. The left image is a red car with a yellow bounding box. The right image is an orange car with a yellow bounding box and a blue dashed bounding box. A large bracket on the right side of the images is labeled with the number 2.

Metric 3 - D_{shift}

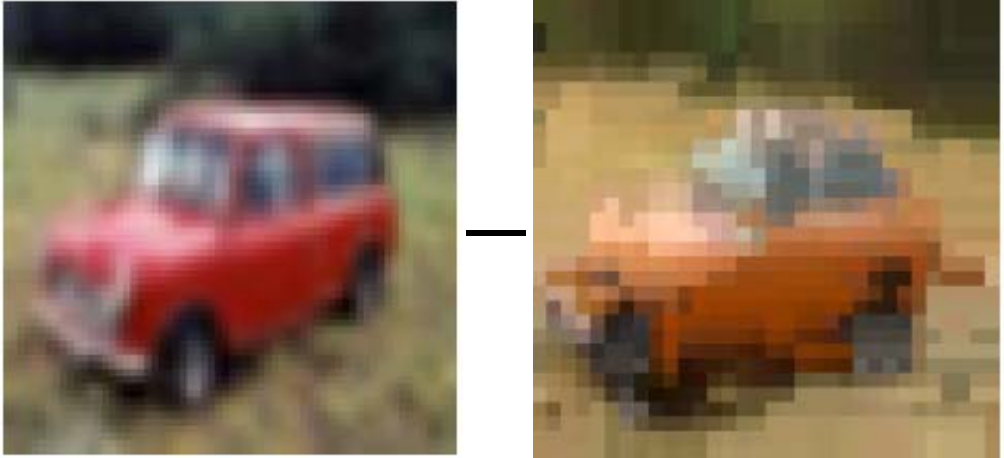
- As per Warping but also allow sub-window shifts

$$D_{shift}^2 = \min_{\text{Local sub-window}} \sum_{x,y,c}$$


- Quick since images are so small

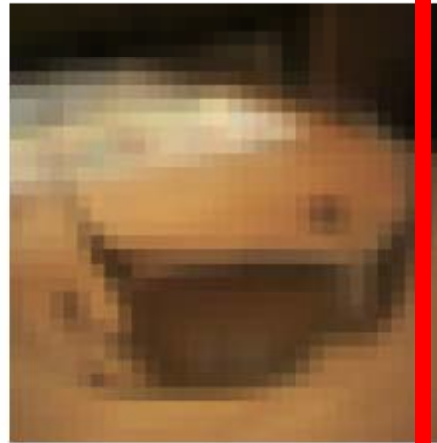
Metric 3 - D_{shift}

- As per Warping but also allow sub-window shifts

$$D_{shift}^2 = \min_{\text{Local sub-window}} \sum_{x,y,c} \left[\begin{array}{c} \text{Image 1} \\ \text{Image 2} \end{array} \right]^2$$


Tried various sizes of sub-window
→ 1x1 (i.e. single pixel) worked best

Comparison of metrics



Target

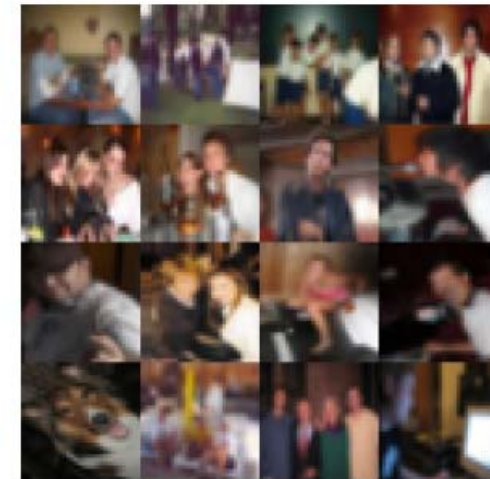
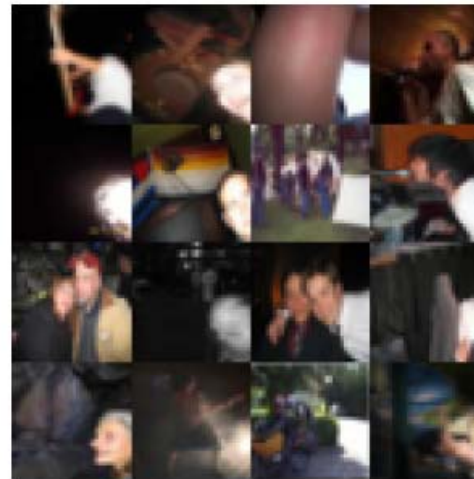
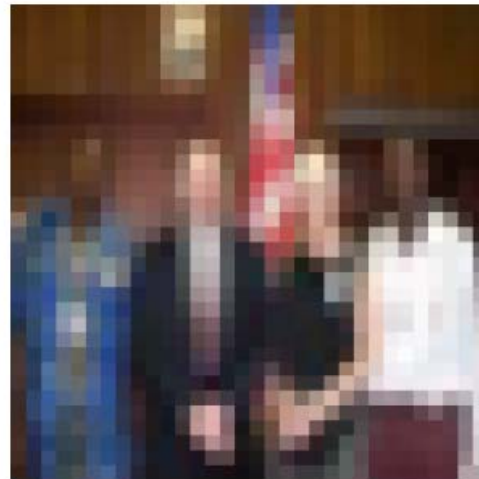
SSD

Warping

Pixel shifting

Sibling Sets with Different Metrics

- Sibling set is 50 images

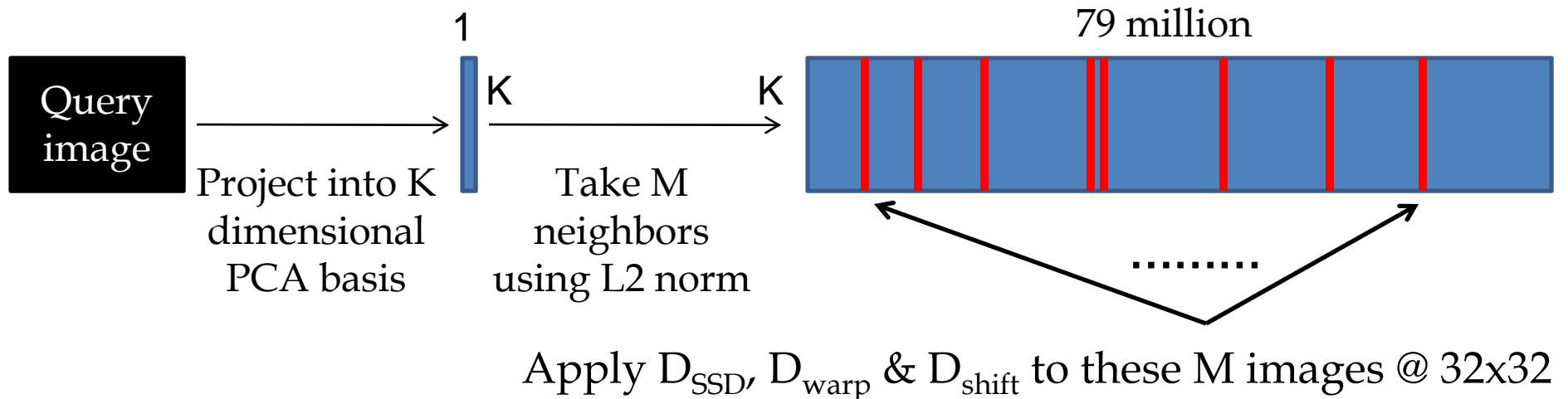


D_{ssd}

D_{shift}

Approximate D_{SSD}

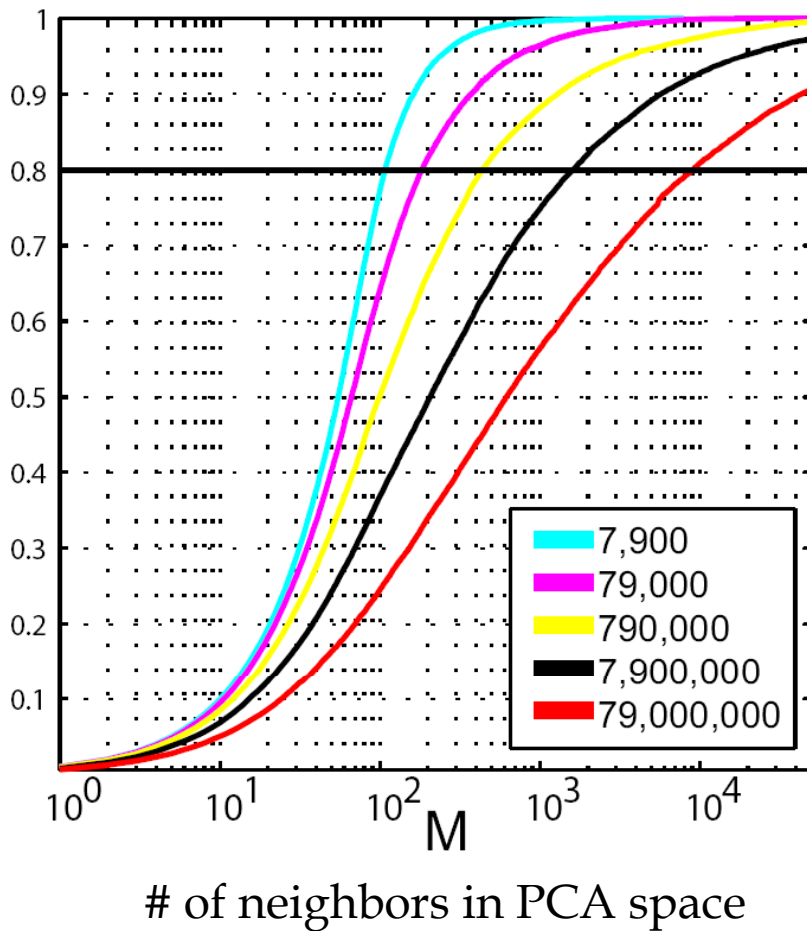
- Exact distance metrics are too expensive to apply to all 79 million images
- Use approximate scheme based on taking first $K=19$ principal components



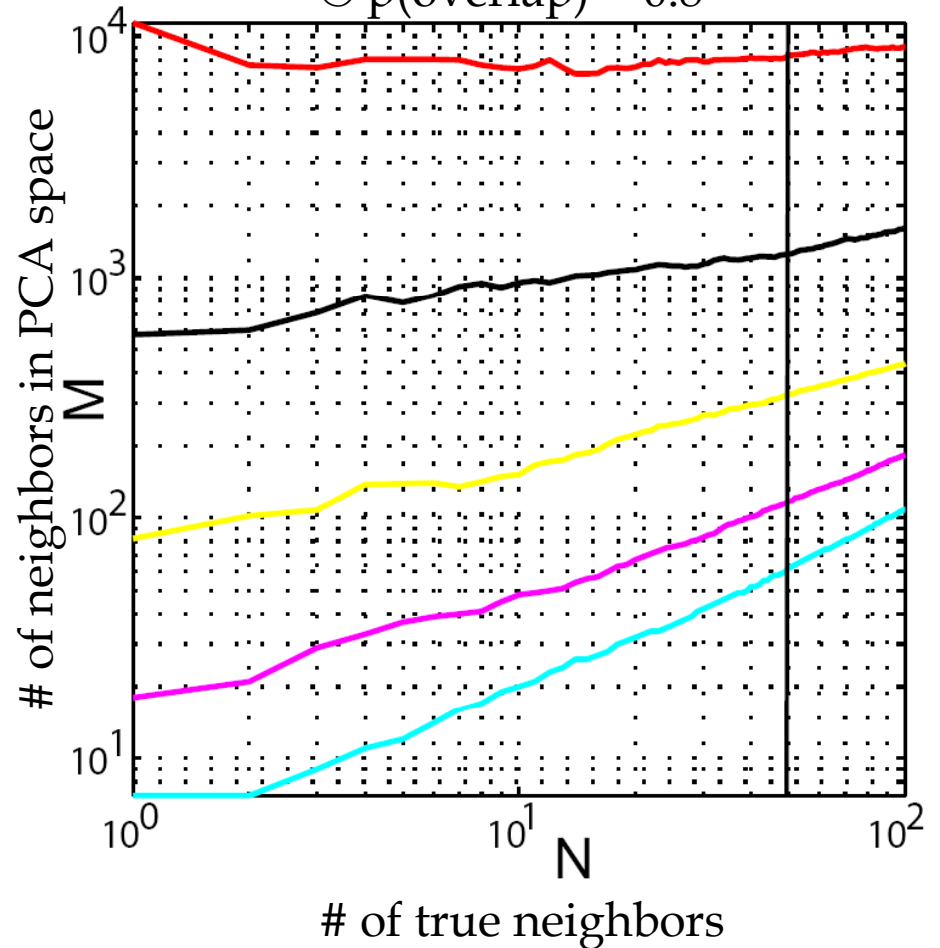
Exact SSD vs Approximate SSD

Overlap between approximate set & 50 true neighbors

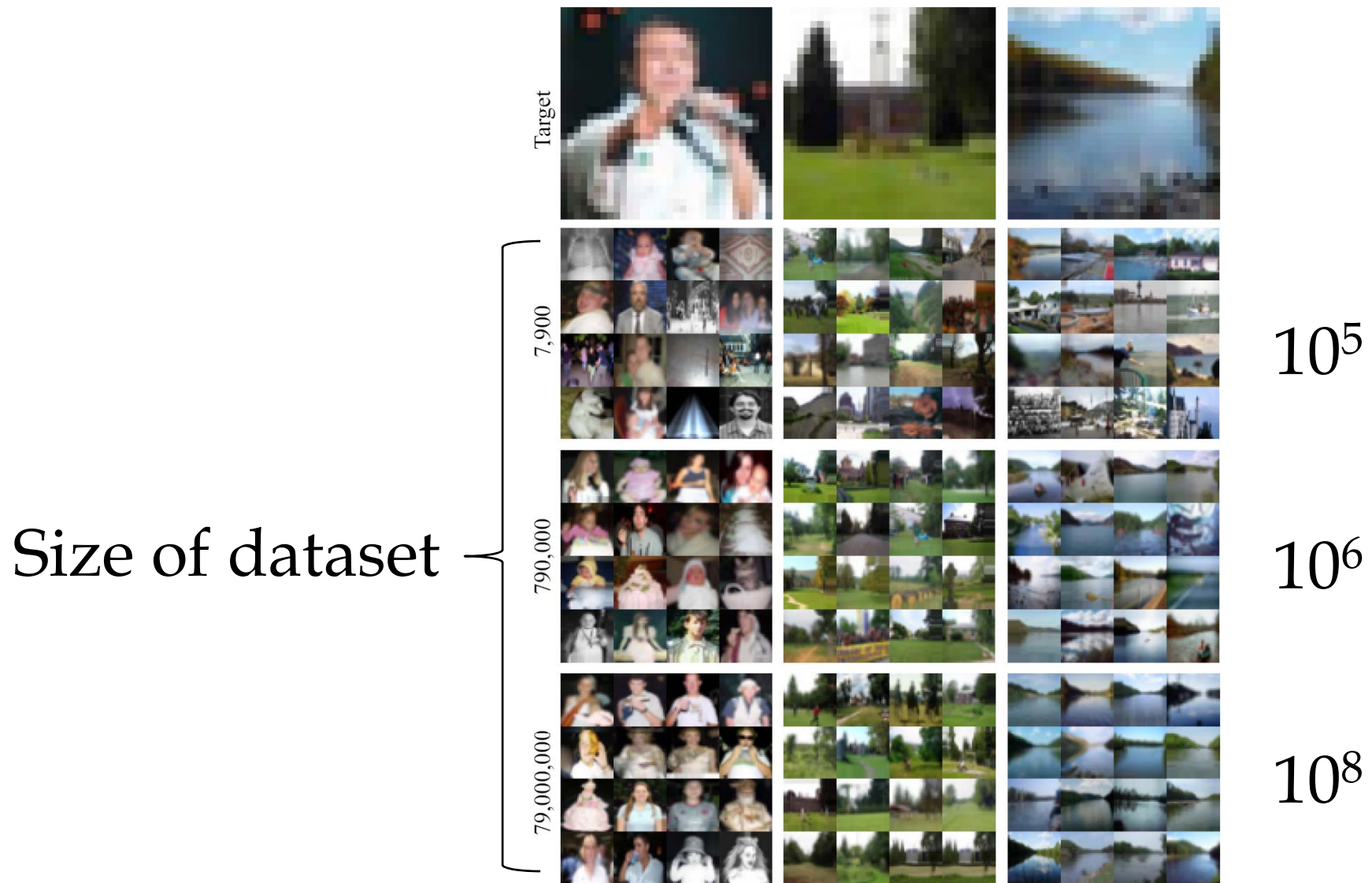
Using N=50 neighbors



@ p(overlap) = 0.8

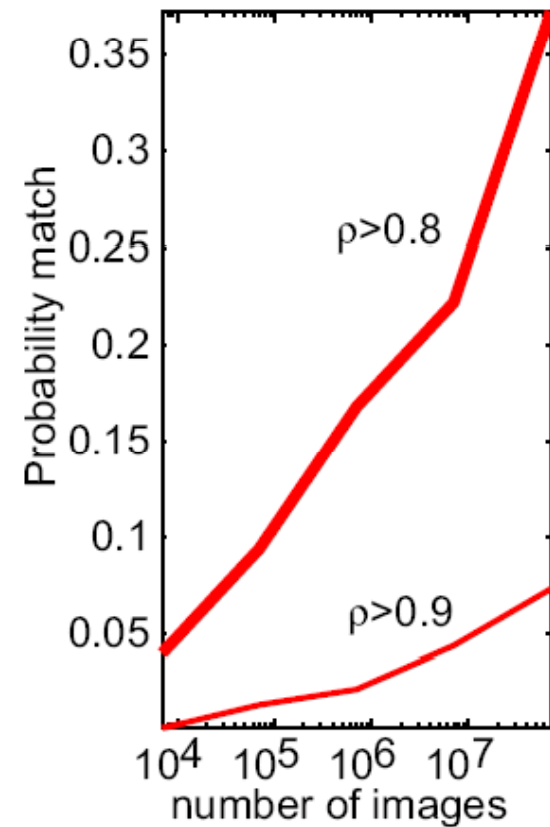
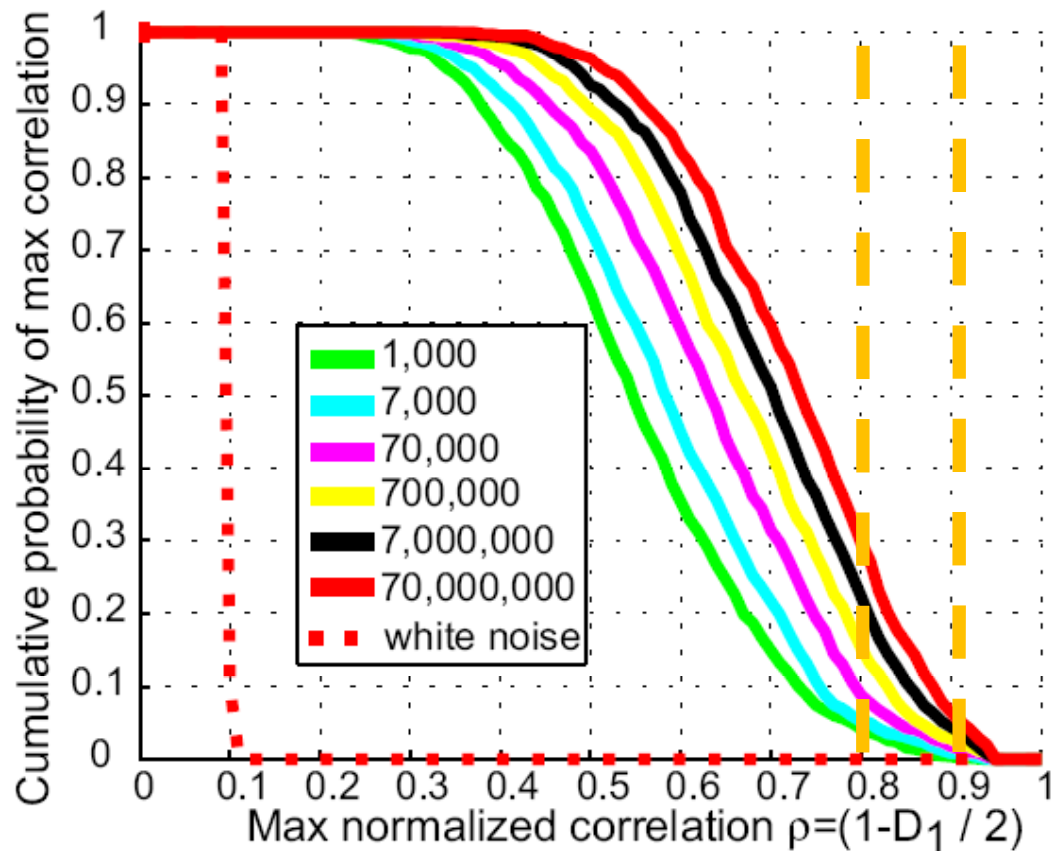


Quality of Sibling Set using D_{shift}



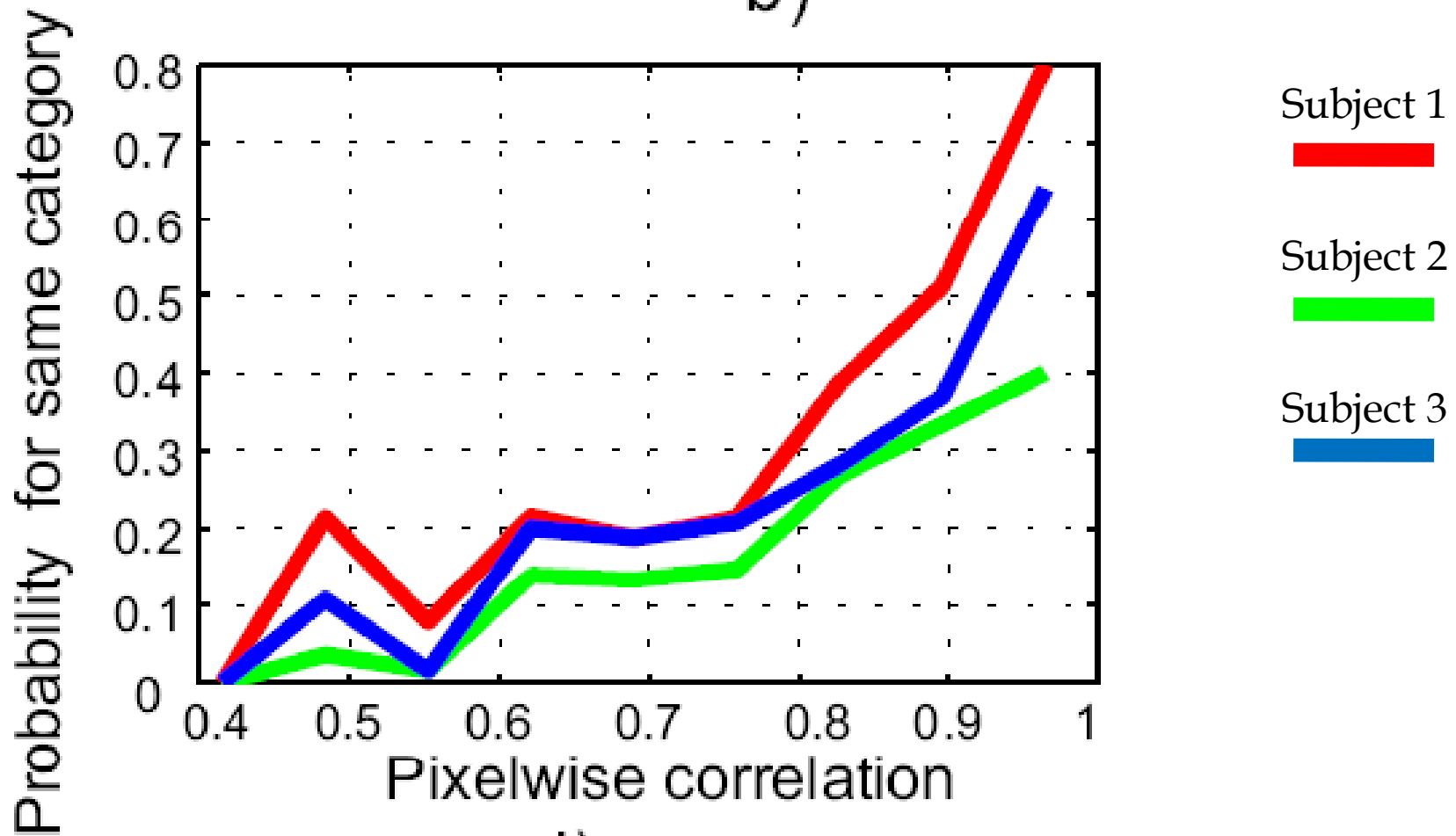
Exploring the Manifold of Images

How Many Images Are There?



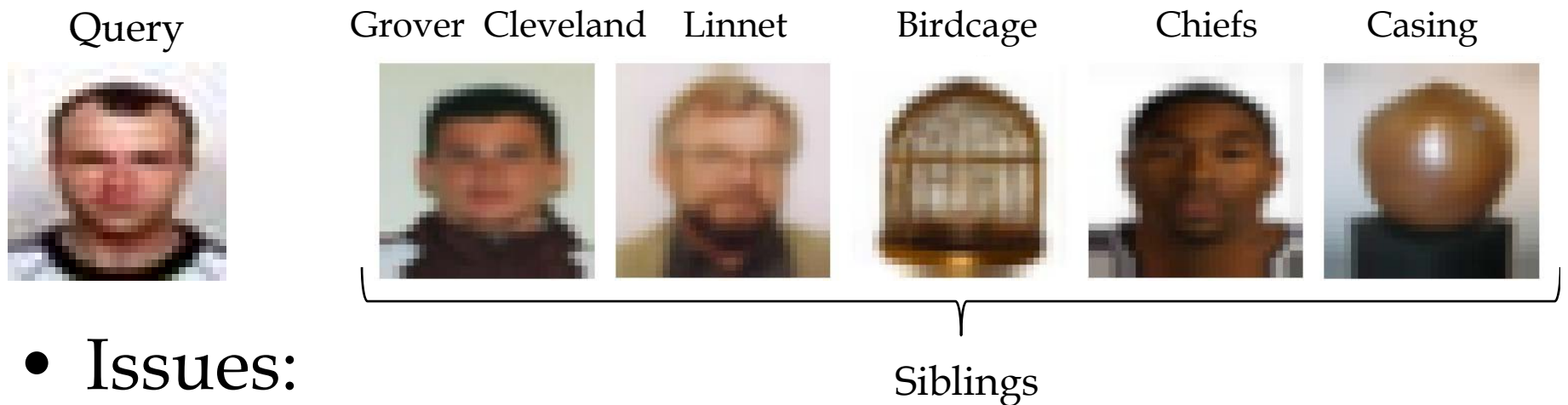
Note: $D_1 = D_{SSD}$

How Does D_{ssd} Relate to Semantic Distance?



Label Assignment

- Distance metrics give set of nearby images
- How to compute label?



- Issues:
 - Labeling noise
 - Keywords can be very specific
 - e.g. yellowfin tuna

Wordnet – a Lexical Dictionary

<http://wordnet.princeton.edu/>

Synonyms/Hypernyms (Ordered by Estimated Frequency) of noun **aardvark**

Sense 1

aardvark, ant bear, anteater, *Orycteropus afer*

=> placental, placental mammal, eutherian, eutherian mammal

=> mammal

=> vertebrate, craniate

=> chordate

=> animal, animate being, beast, brute, creature

=> organism, being

=> living thing, animate thing

=> object, physical object

=> entity

Wordnet Hierarchy

Synonyms/Hypernyms (Ordered by Estimated Frequency) of noun **aardvark**

Sense 1

aardvark, ant bear, anteater, Orycteropus afer

=> **placental**, placental mammal, eutherian, eutherian mammal

=> **mammal**

=> **vertebrate**, craniate

=> **chordate**

=> **animal**, animate being, beast, brute, creature

=> **organism**, being

=> **living thing**, animate thing

=> **object**, physical object

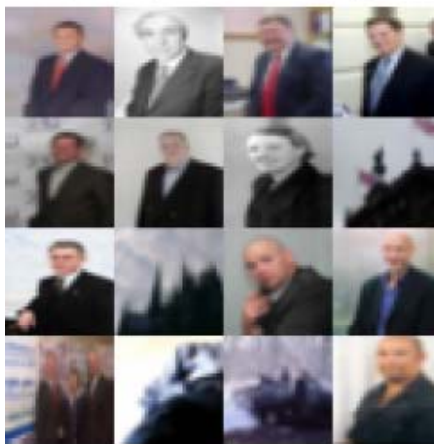
=> **entity**

- Convert graph structure into tree by taking most common meaning

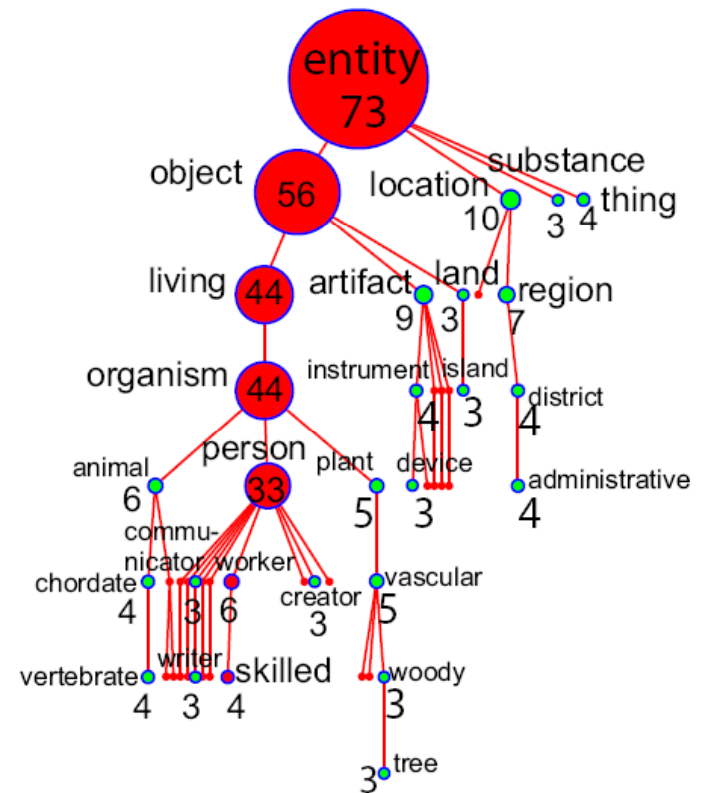
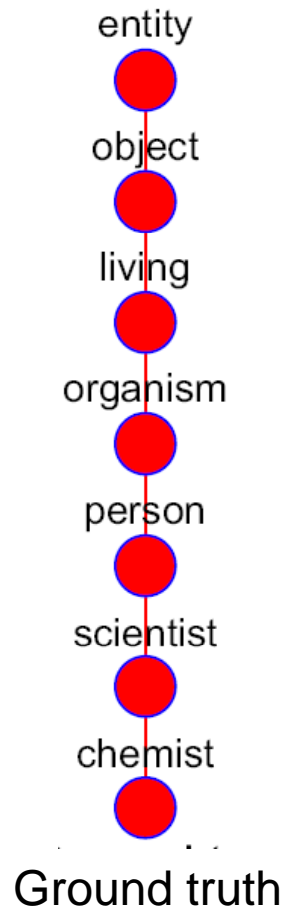
Wordnet Voting Scheme



a) Input image



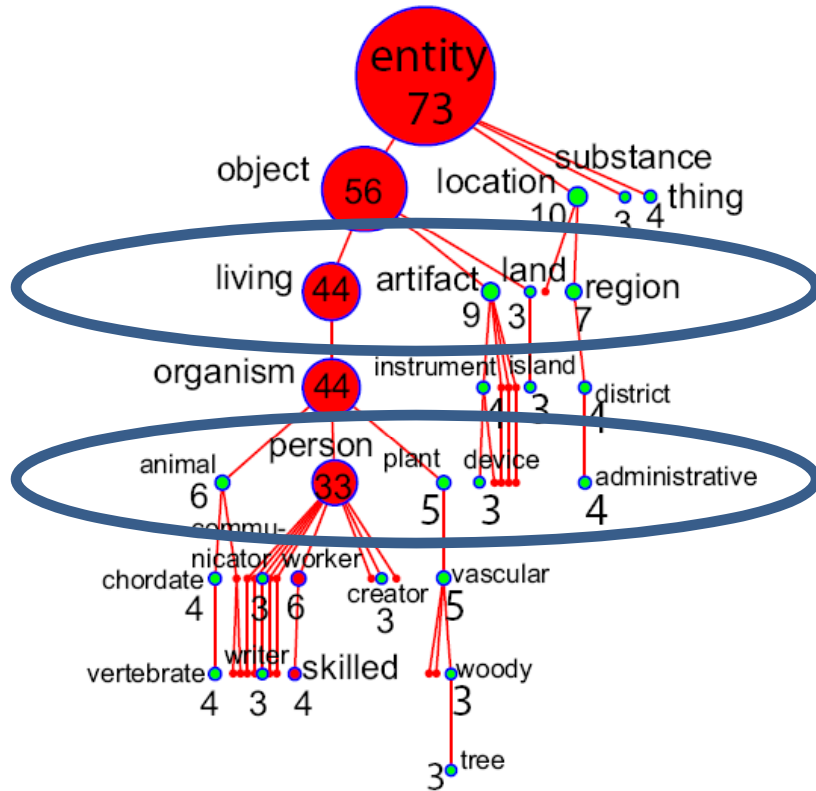
b) Neighbors



d) Wordnet voted branches

One image - one vote

Classification at Multiple Semantic Levels



Votes:

Living	64
Artifact	93
Land	5
Region	3
Administrative	40
Others	22

1) d) Wordnet voted branches

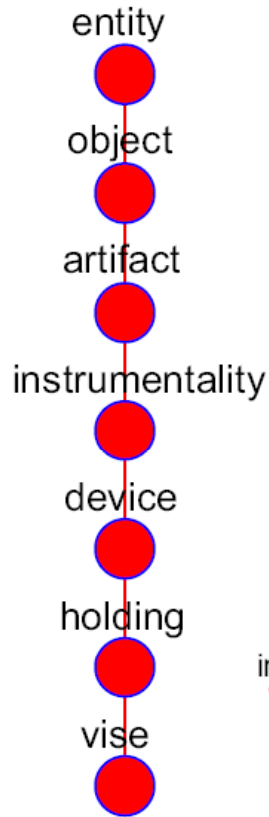
Wordnet Voting Scheme



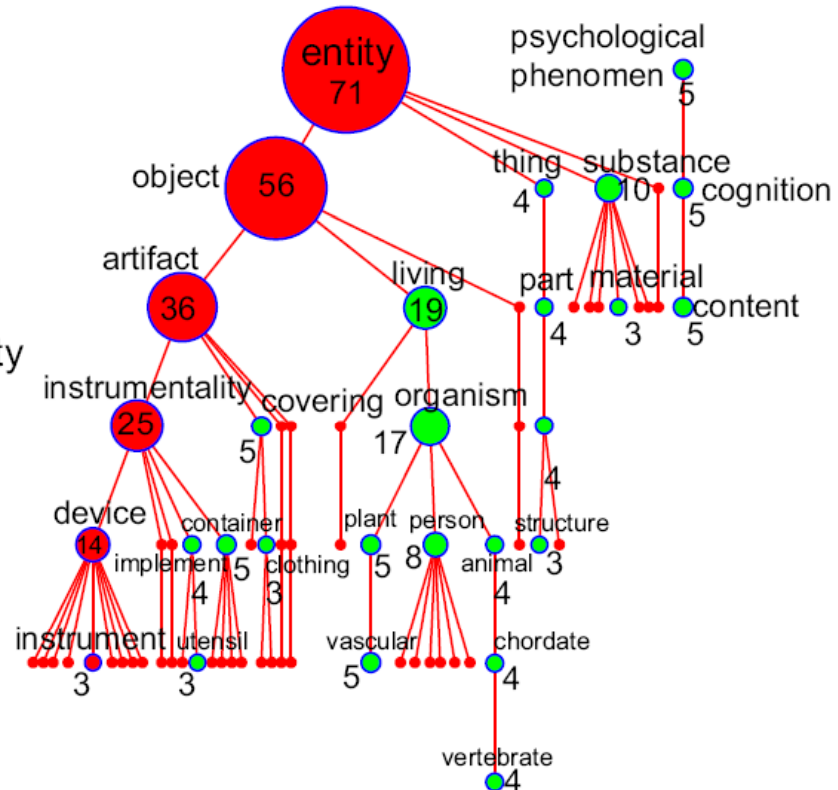
a) Input image



b) Neighbors



c) Ground truth



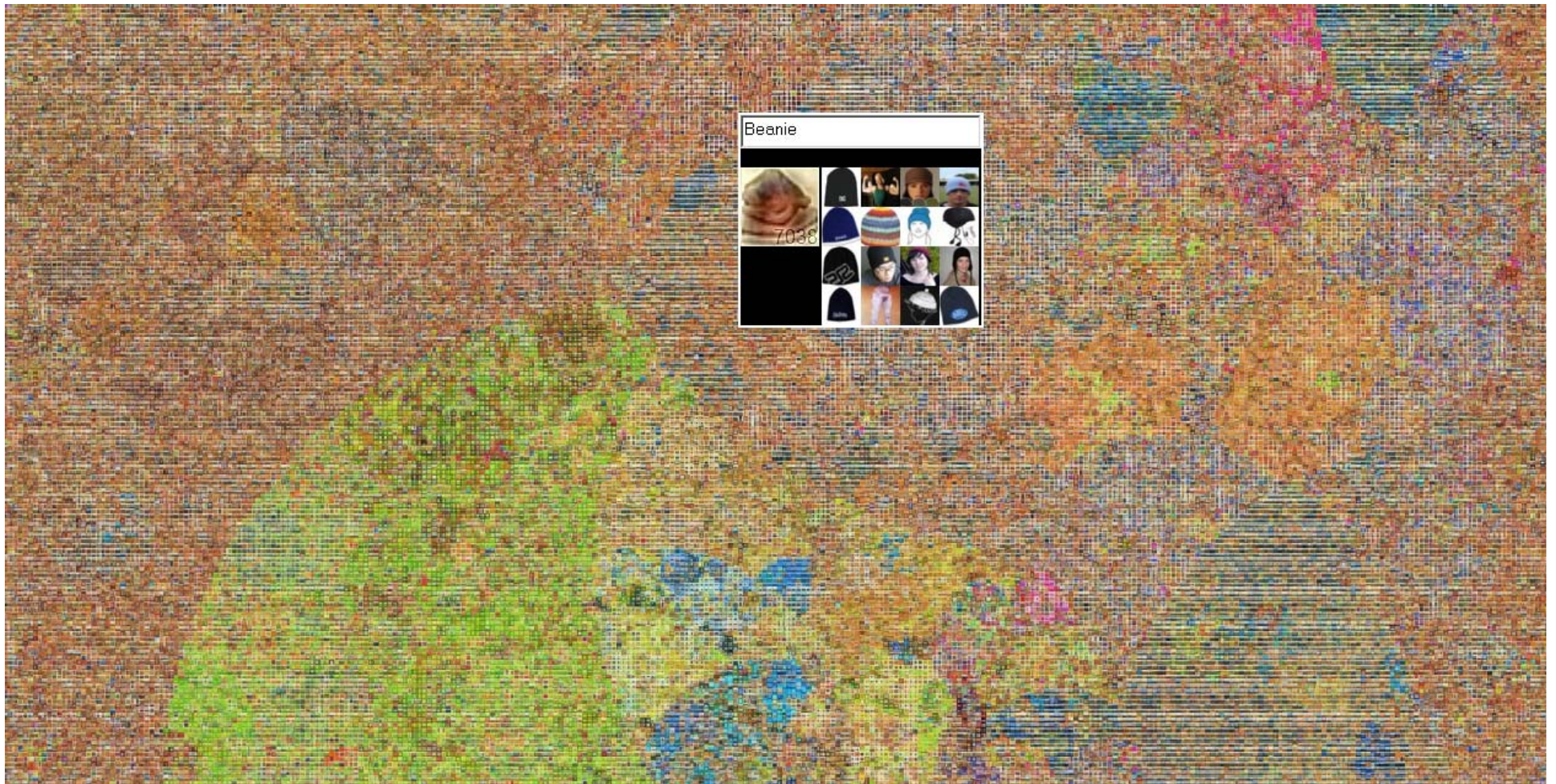
d) Wordnet voted branches

Wordnet Voting

- Overcomes differences in level of semantic labeling:
 - e.g. “person” & “sir arthur conan doyle”
- Totally incorrect labels form hopefully uniform background noise
- Assumes semantic and visual consistency are closely related

Semantic vs Visual Hierarchy

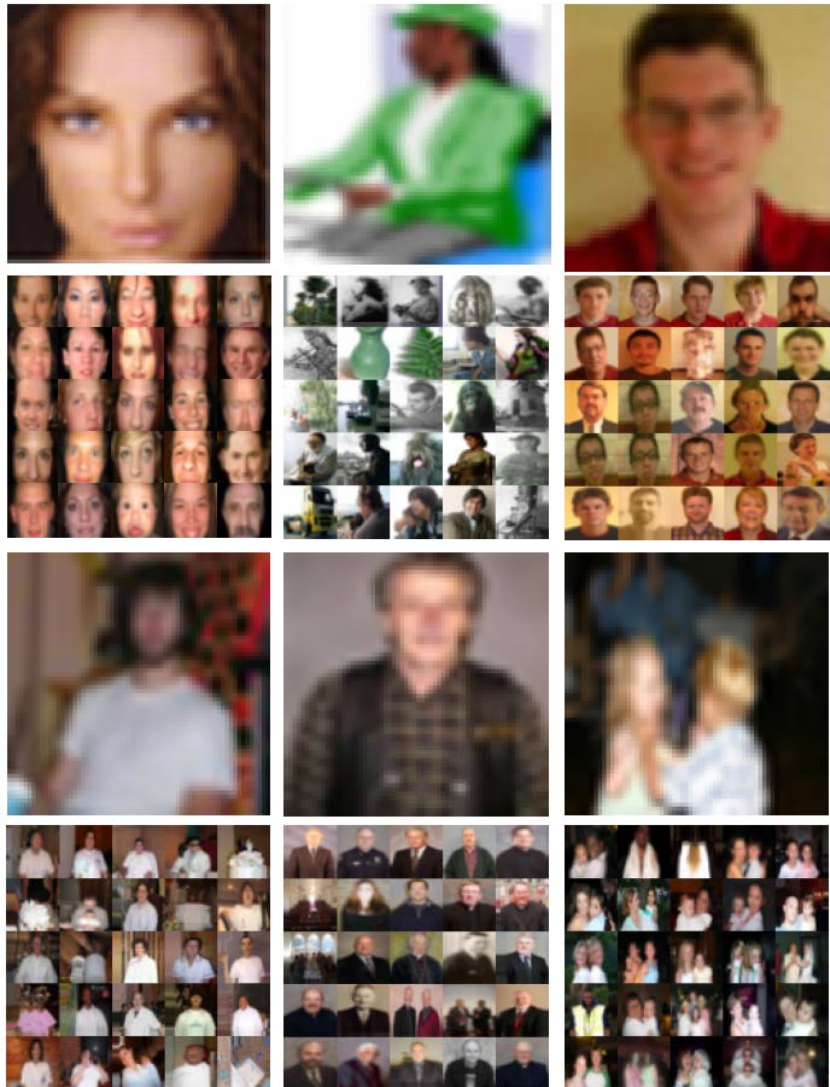
Interactive version at <http://people.csail.mit.edu/torralba/tinyimages>



Recognition Experiments

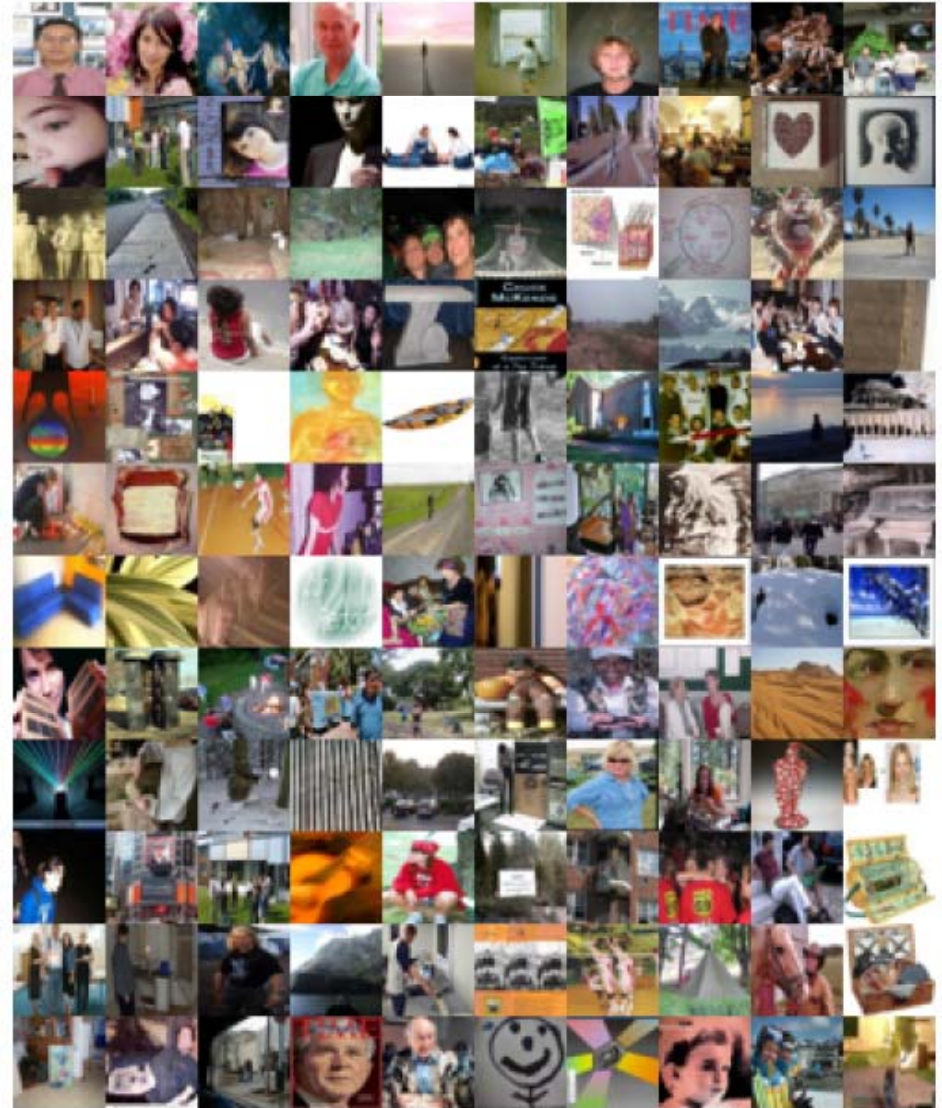
Person Recognition

- 23% of all images in dataset contain people
- Wide range of poses: not just frontal faces



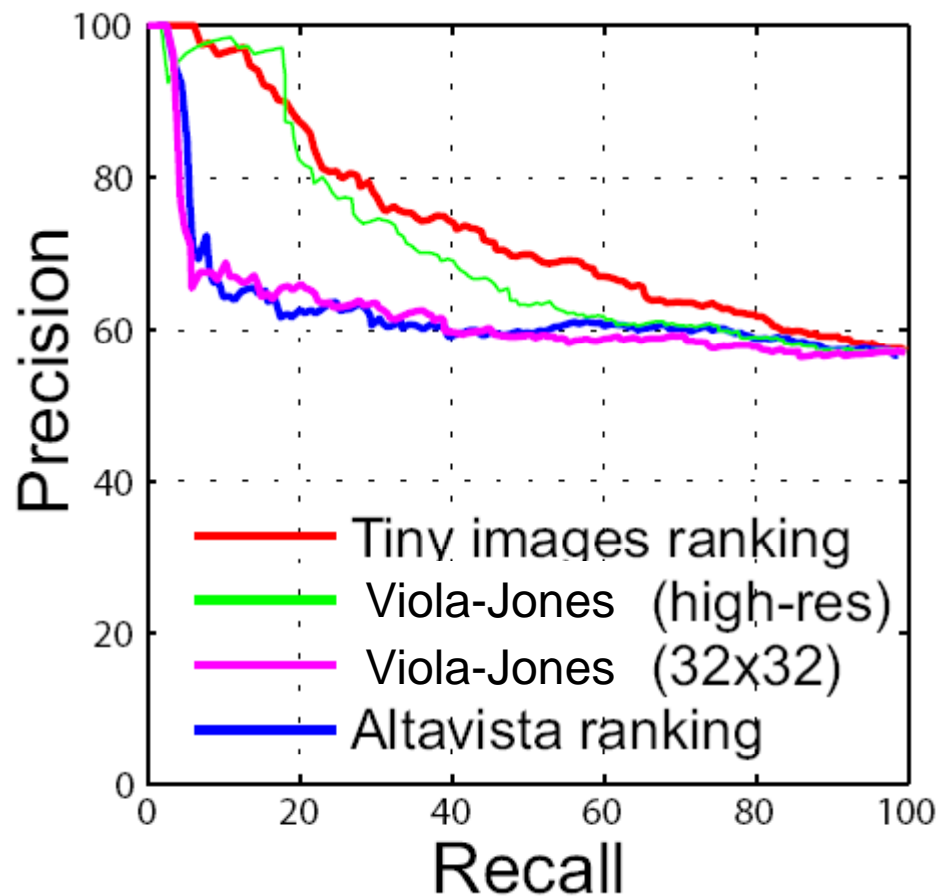
Person Recognition – Test Set

- 1016 images from Altavista using “person” query
- High res and 32x32 available
- Disjoint from 79 million tiny images



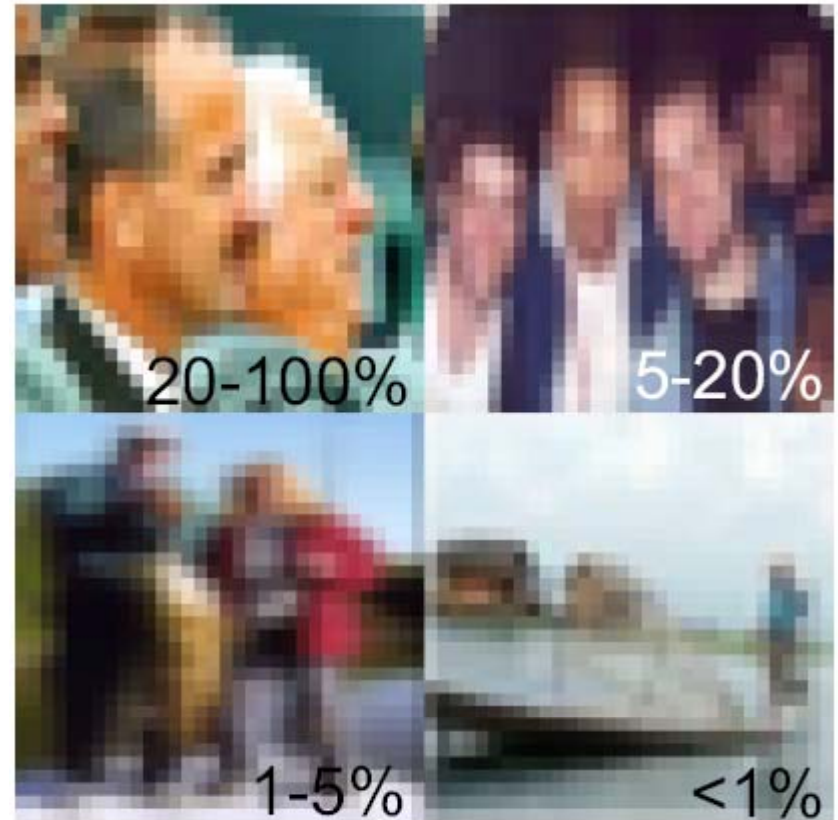
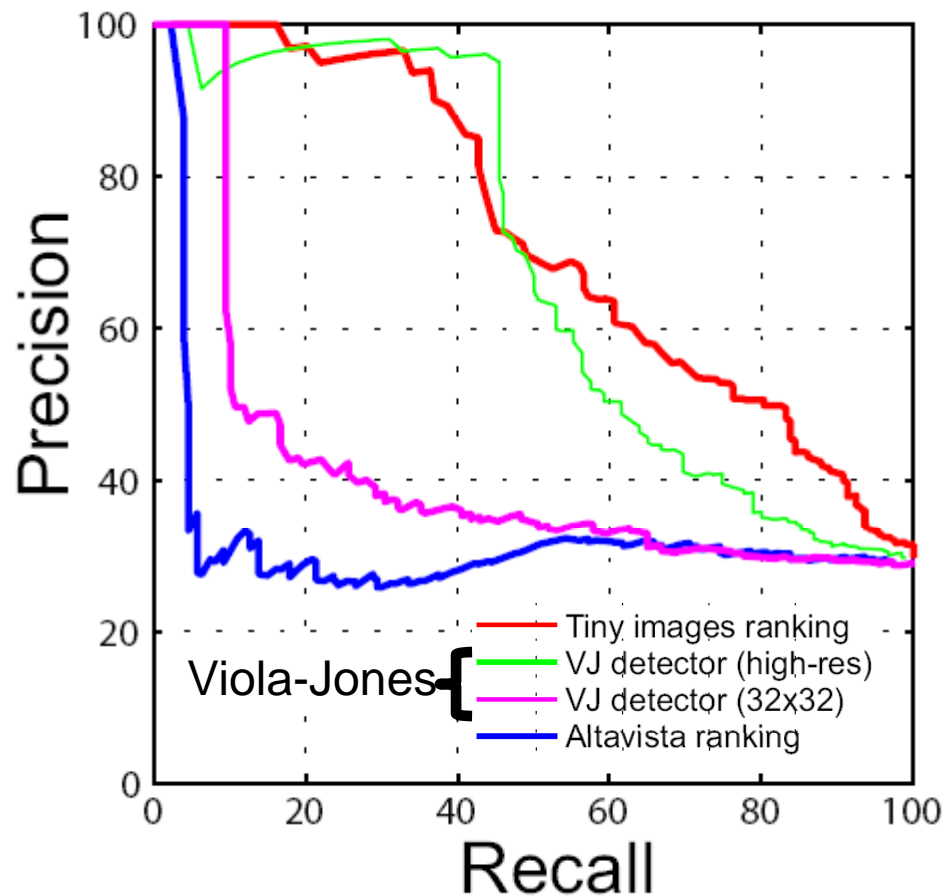
Person Recognition

- Task: person in image or not?



Person Recognition

- Subset where face $>20\%$ of image

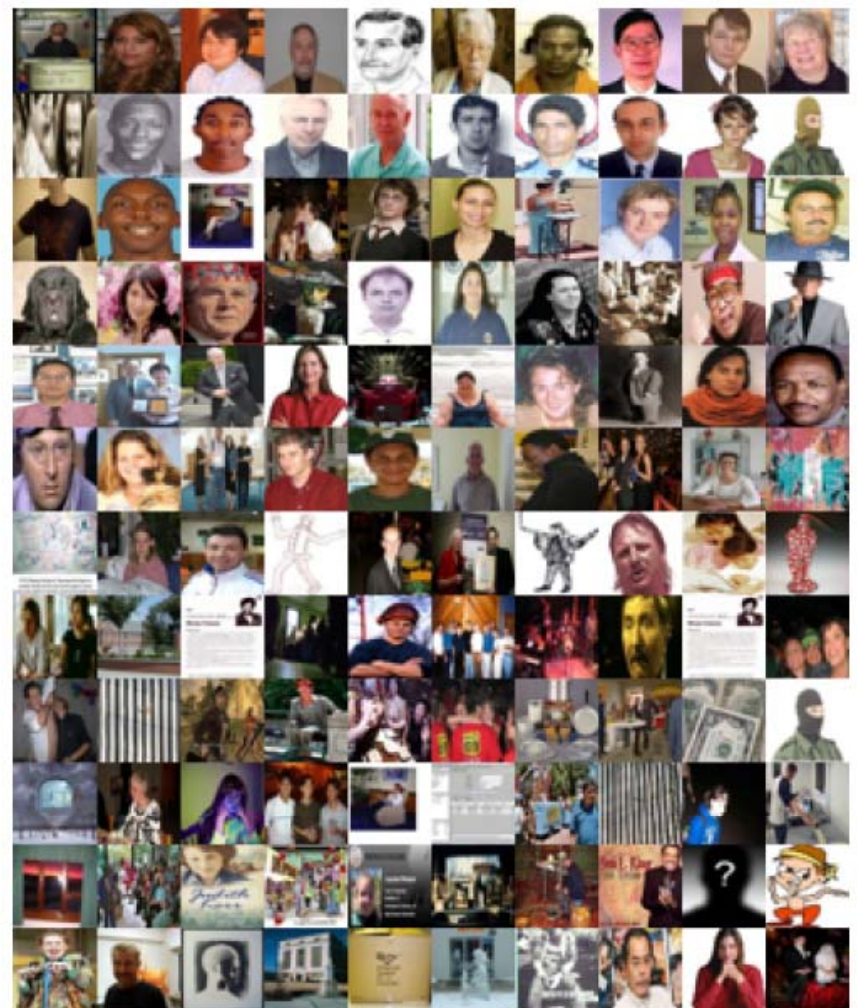


Re-ranked Altavista Images

Original

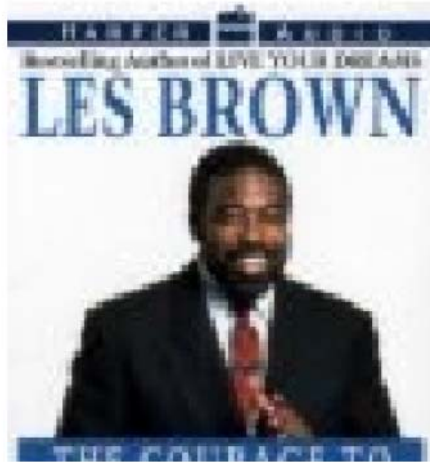


Re-ranked



Person Localization

High-res
Query



a)

Ncuts
Segmentation

Gives putative
crops



b)

Scene Classification

- Test set: 1125 images randomly drawn from 79 million.
- Task: {scene} vs all other classes

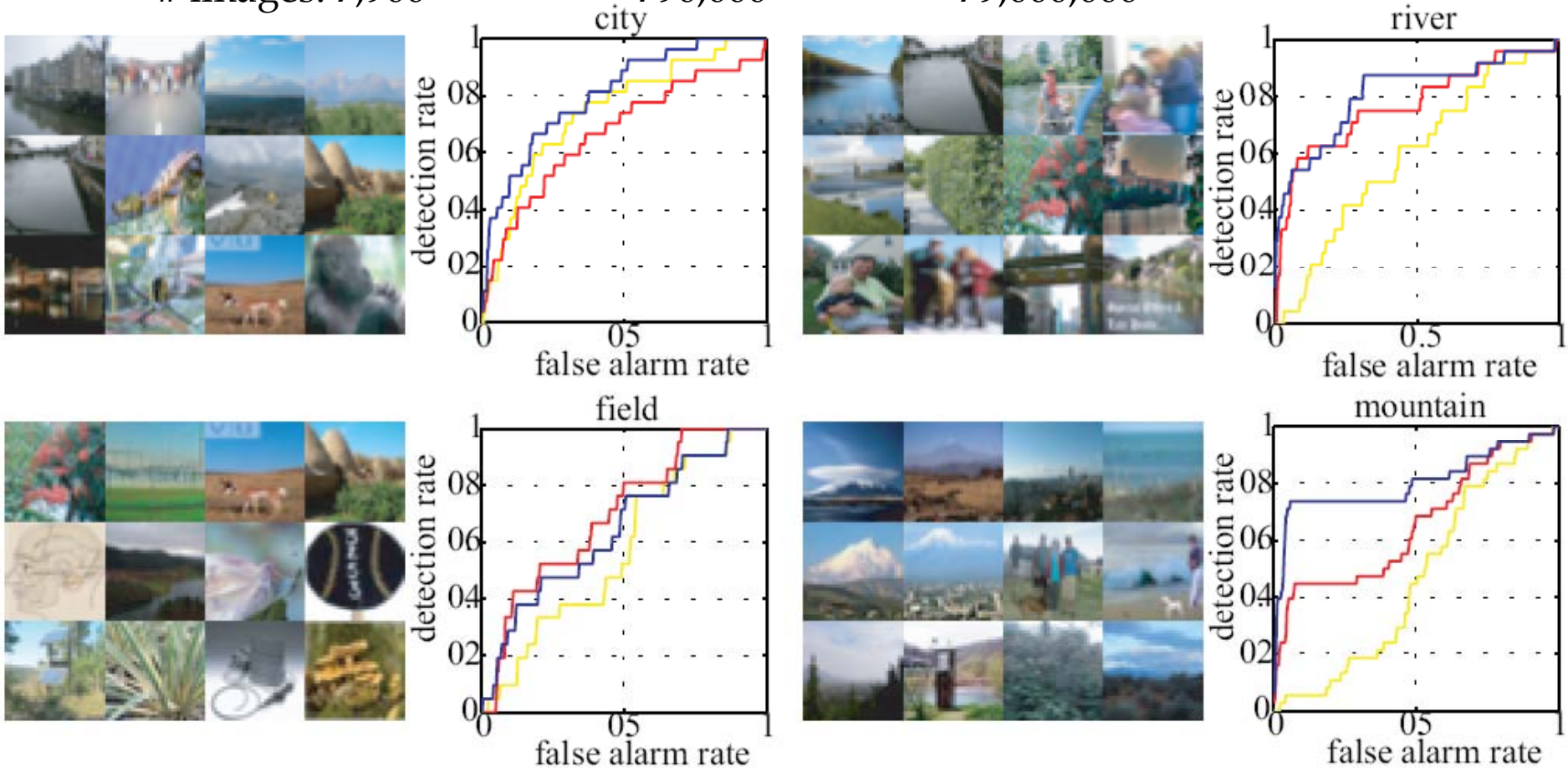
images: 7,900



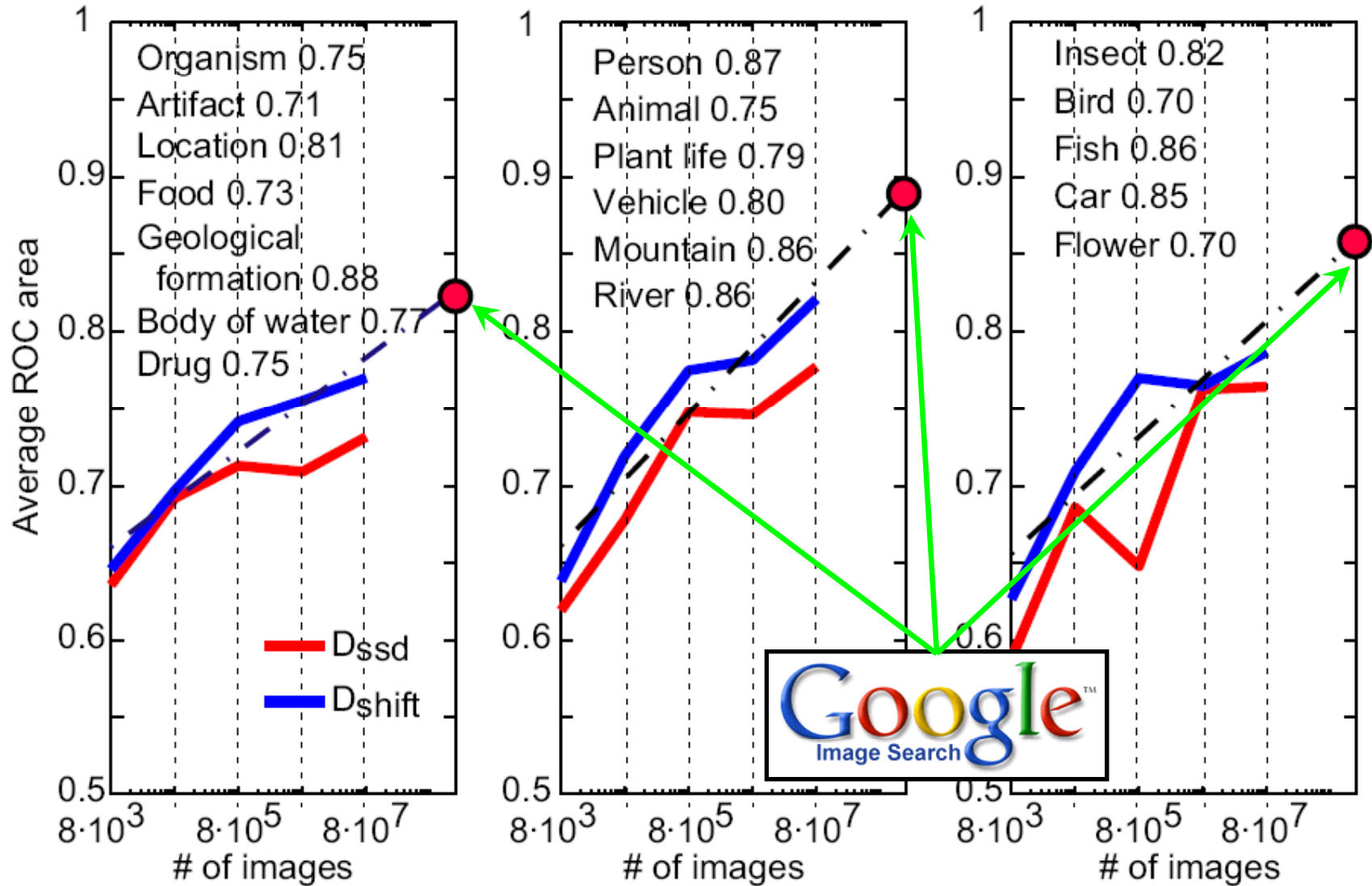
790,000



79,000,000



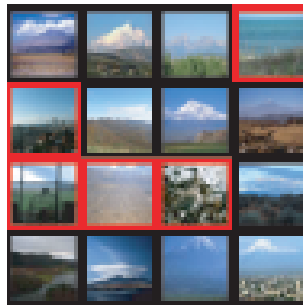
Object Classification



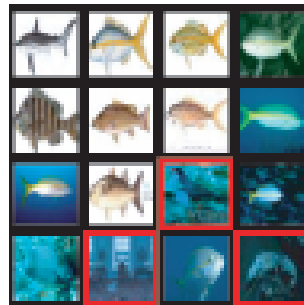
Object Classification

images: 7,900 ■ 790,000 ■ 79,000,000 ■

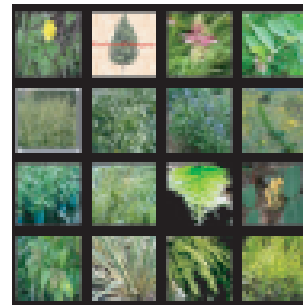
Geological
formation (32)



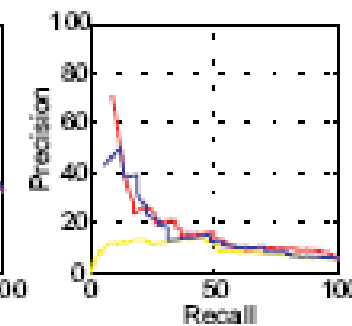
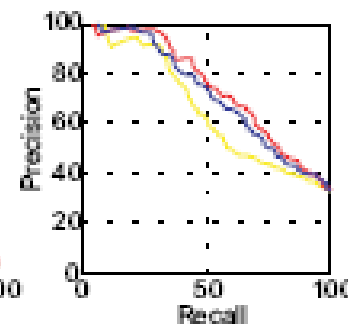
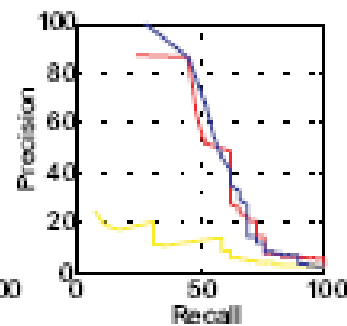
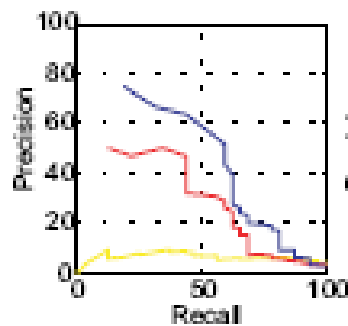
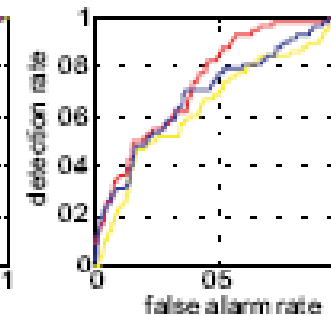
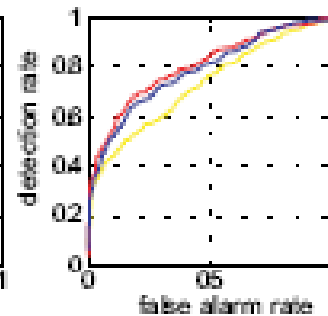
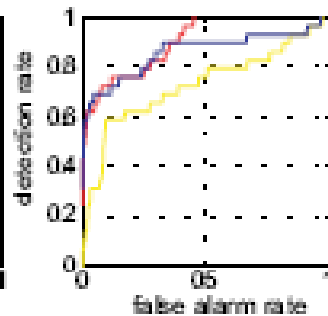
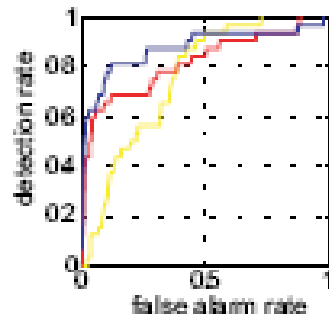
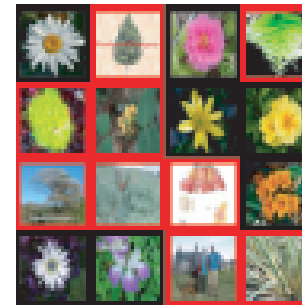
Fish
(29)



Plant life
(335)



Flower
(58)



Other Applications

Automatic Colorization

Grayscale input
High resolution

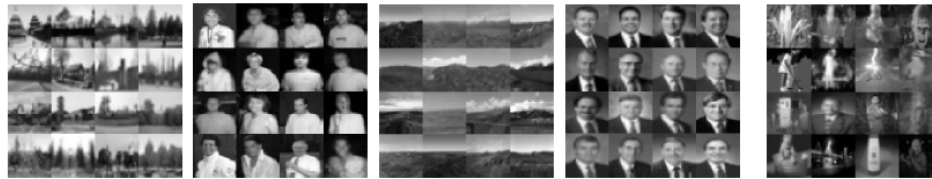


Automatic Colorization

Grayscale input
High resolution



Grayscale
32x32 siblings

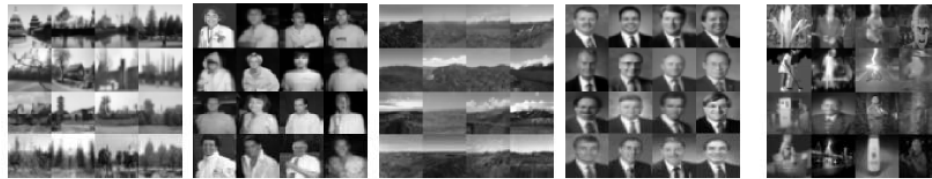


Automatic Colorization

Grayscale input
High resolution



Grayscale
32x32 siblings



Color siblings
high resolution

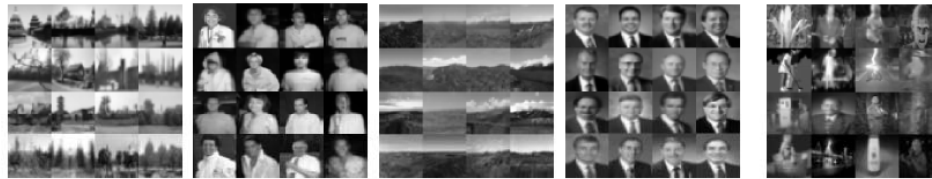


Automatic Colorization

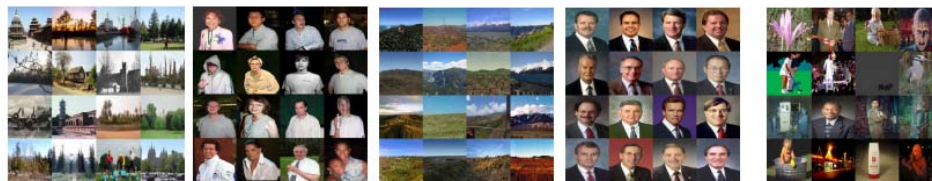
Grayscale input
High resolution



Grayscale
32x32 siblings



Color siblings
high resolution



Average of
color siblings

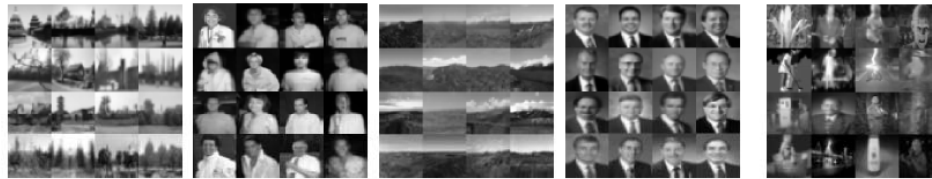


Automatic Colorization

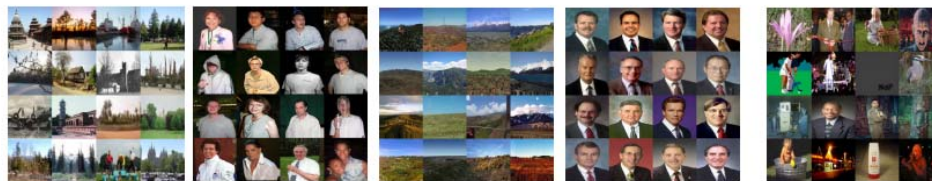
Grayscale input
High resolution



Grayscale
32x32 siblings



Color siblings
high resolution



Average of
color siblings



Colorization of input
using average

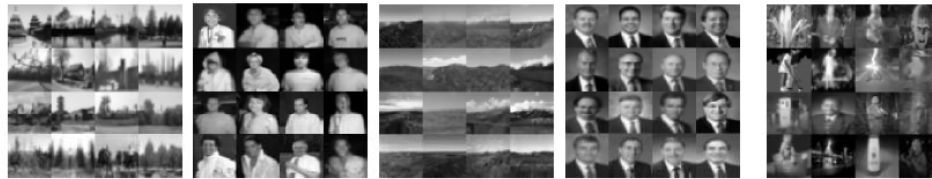


Automatic Colorization

Grayscale input
High resolution



Grayscale
32x32 siblings



Color siblings
high resolution



Average of
color siblings



Colorization of input
using average



Colorization of input
using specific siblings



Automatic Colorization Result

Grayscale input High resolution

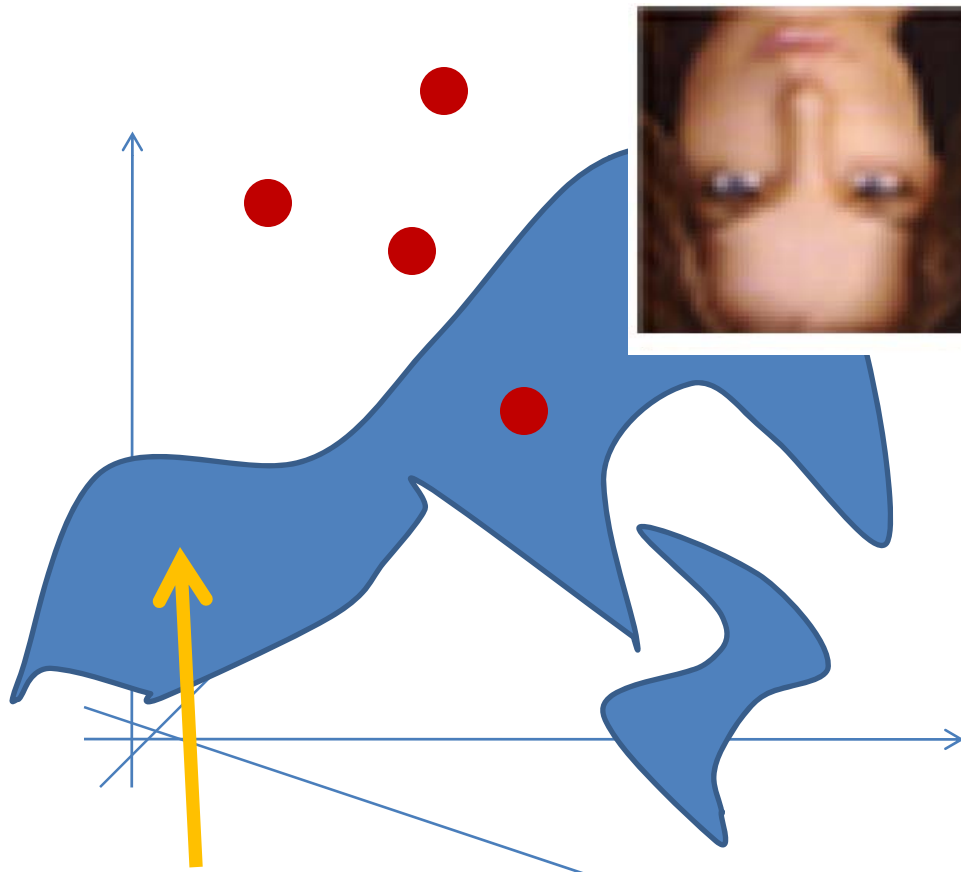


Colorization of input using average

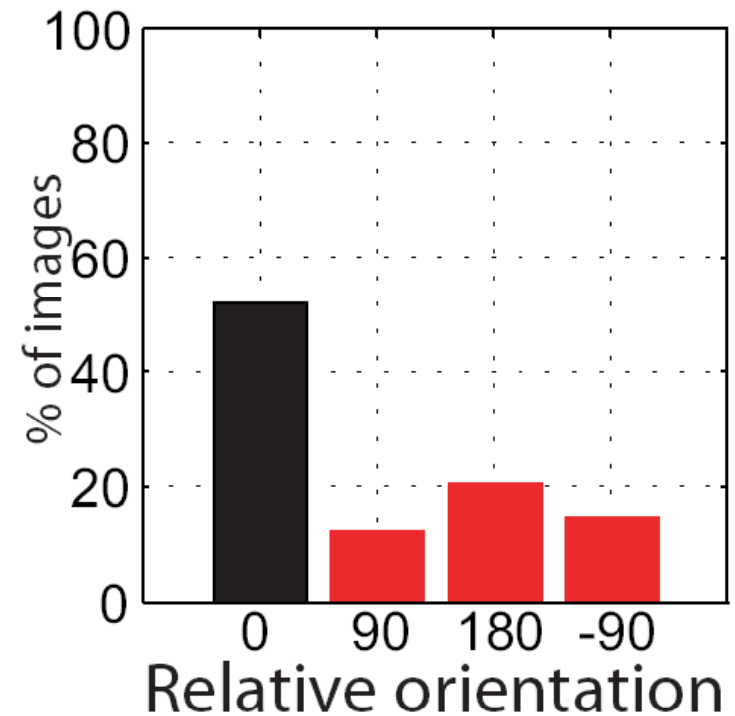


Automatic Orientation

- Look at mean distance to neighbors



Subspace of natural images



Automatic Orientation Examples

0.70



0.64



0.66



0.64



0.86



0.76



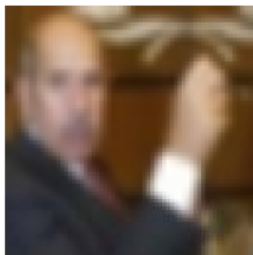
0.79



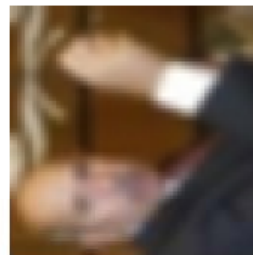
0.77



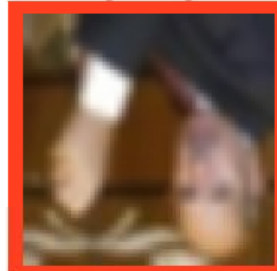
0.66



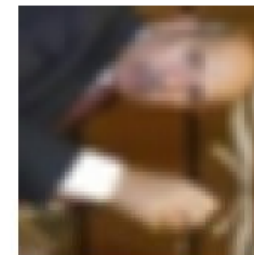
0.62



0.70



0.63

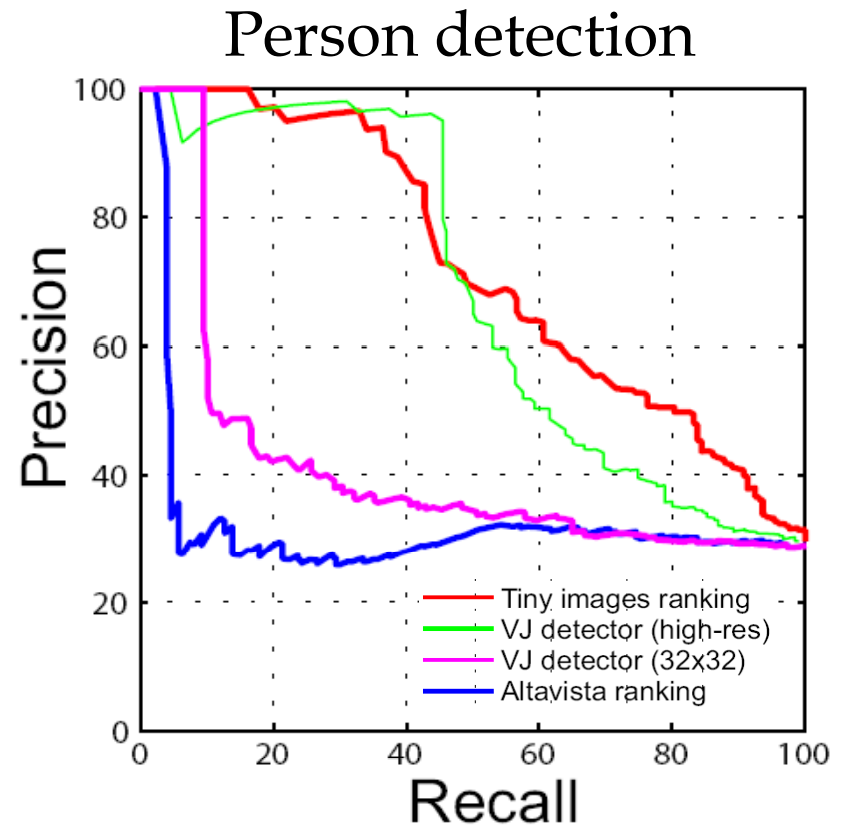


Related Work

- Hayes & Efros, Scene Completion using Millions of photographs, SIGGRAPH 2007.
- Nister & Stewenius. Scalable recognition with a vocabulary tree, CVPR 2006.
- Hoogs & Collins. Object boundary detection in images using a semantic ontology. In *AAAI, 2006*.
- Barnard et al., Matching words and pictures. *JMLR*, 2003.
- Shakhnarovich et al. Fast pose estimation with parameter sensitive hashing, *ICCV 2003*

Conclusions

- Can get good results simple algorithms & lots of data
- Issues with Internet images: labeling noise & image biases.
- Bring in learning:
Distance metrics, text & images



Webpage: <http://people.csail.mit.edu/torr/alba/tinyimages>