

Tamara Berg

Representing Images

Announcements

- Alex will release HW this weekend
- Make-up class for last week will be scheduled later this semester

Representing Images



Keep all the pixels!

Pros? Cons?

2nd Try



Compute average pixel

Pros? Cons?

3rd Try



Photo by: [marielito](#)



Represent the image as a spatial grid of average pixel colors

Pros?

Cons?

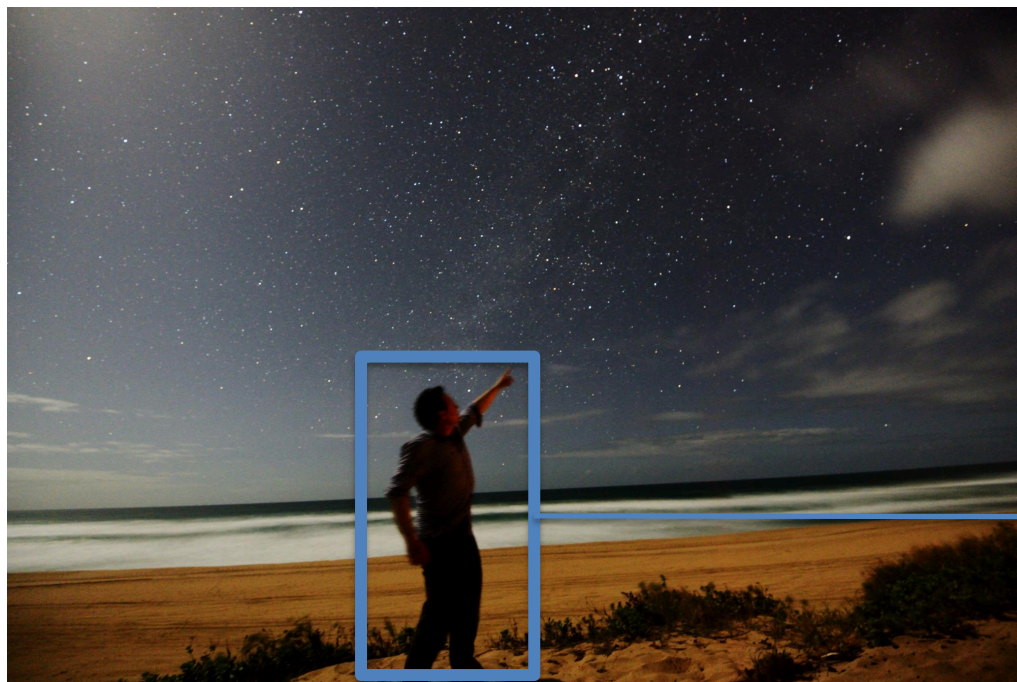
QBIC system



- First content based image retrieval system
 - Query by image content (QBIC)
 - IBM 1995
 - QBIC interprets the virtual canvas as a grid of coloured areas, then matches this grid to other images stored in the database.

Feature Choices

Global vs Local?



Beach

Person

Depends whether you want
whole or local image similarity

Feature types!



The representation of these two umbrella's should be similar....
Under a color based representation they look completely different!

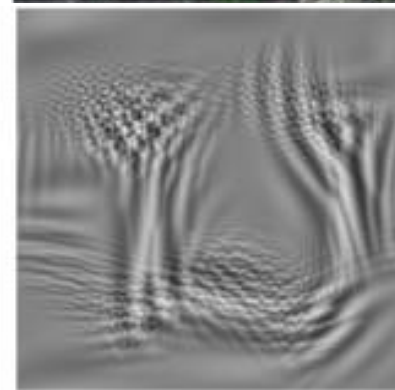
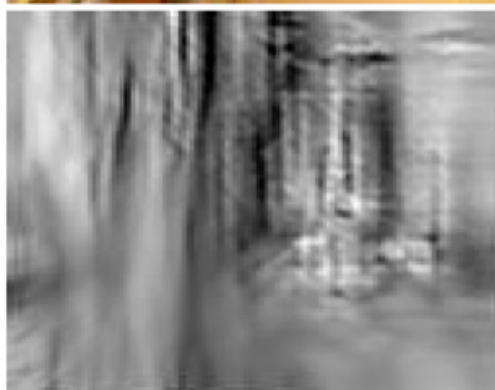
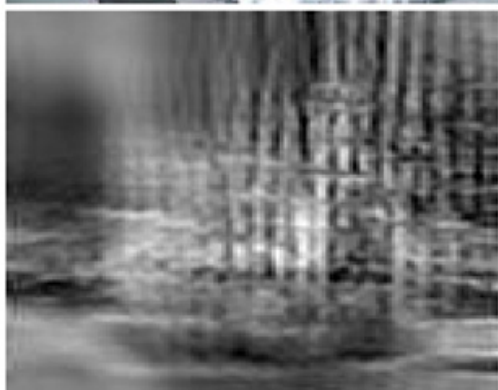
Outline

Image Representations

Global vs Local – depends on the task

Feature Type: color, shape, texture

Global Features



The “gist” of a scene: Oliva & Torralba, 2001

Limitations of global appearance models

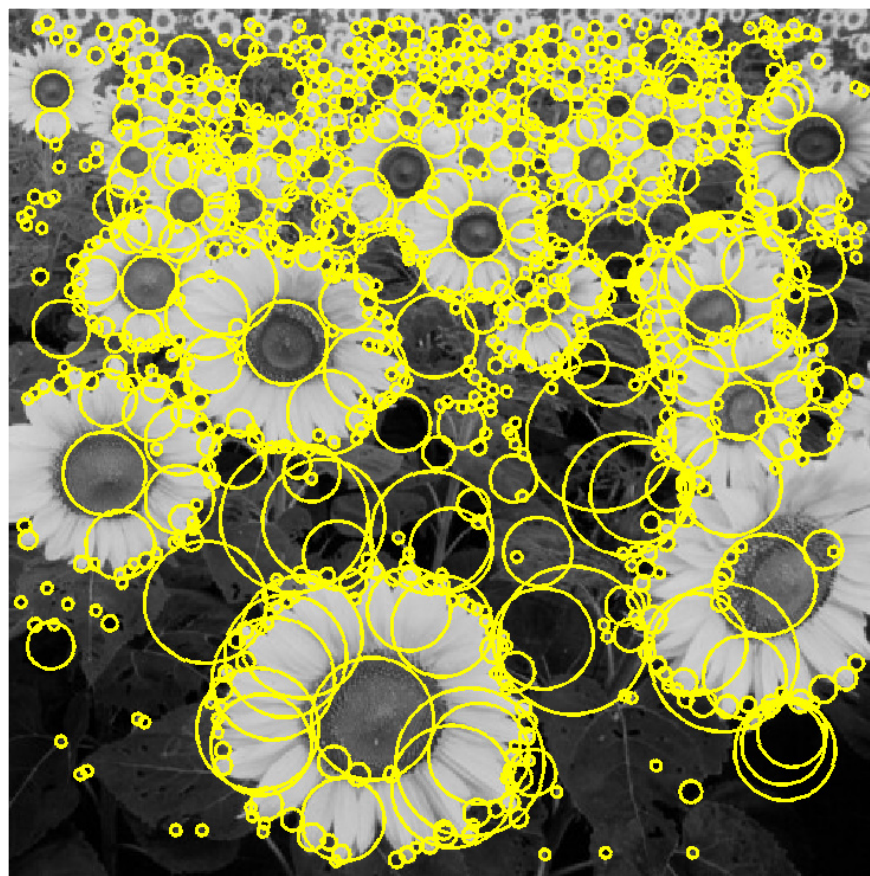
- Can work on relatively simple patterns



- Not robust to clutter, occlusion

Local Features

Feature points (locations) + feature descriptors



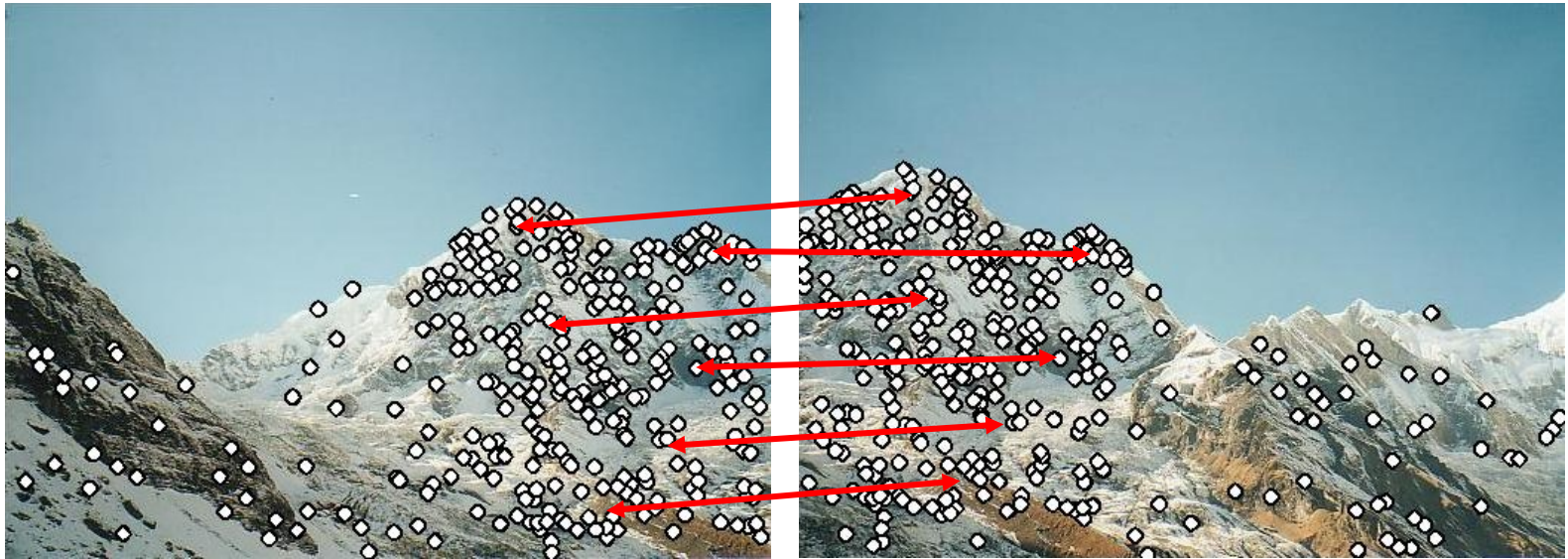
Why extract features?

- Motivation: panorama stitching
 - We have two images – how do we combine them?



Why extract features?

- Motivation: panorama stitching
 - We have two images – how do we combine them?



Step 1: extract features

Step 2: match features

Why extract features?

- Motivation: panorama stitching
 - We have two images – how do we combine them?



Step 1: extract features

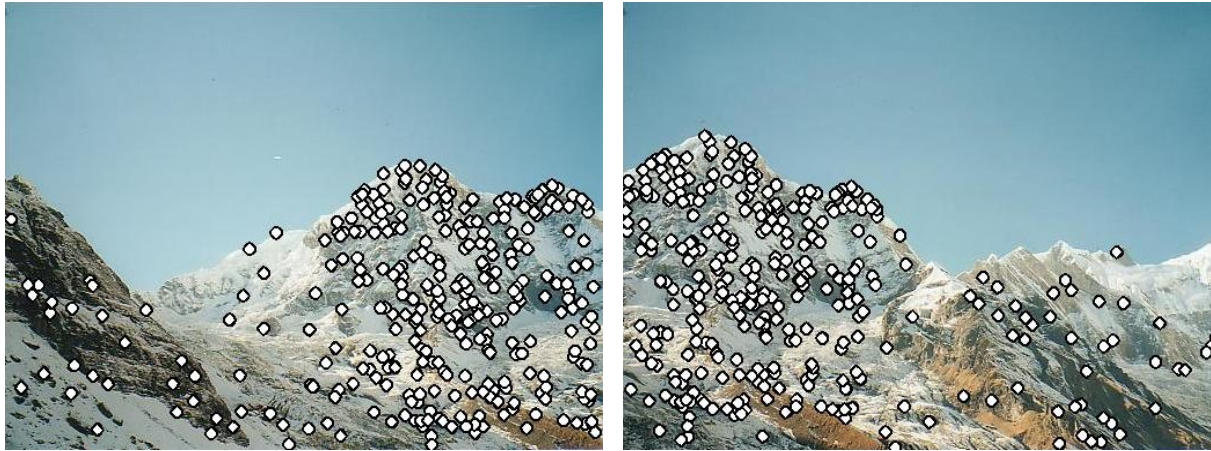
Step 2: match features

Step 3: align images

Local Features for object detection



Characteristics of good features



- **Repeatability**
 - The same feature can be found in several images despite geometric and photometric transformations
- **Saliency**
 - Each feature has a distinctive description
- **Compactness and efficiency**
 - Many fewer features than image pixels
- **Locality**
 - A feature occupies a relatively small area of the image; robust to clutter and occlusion

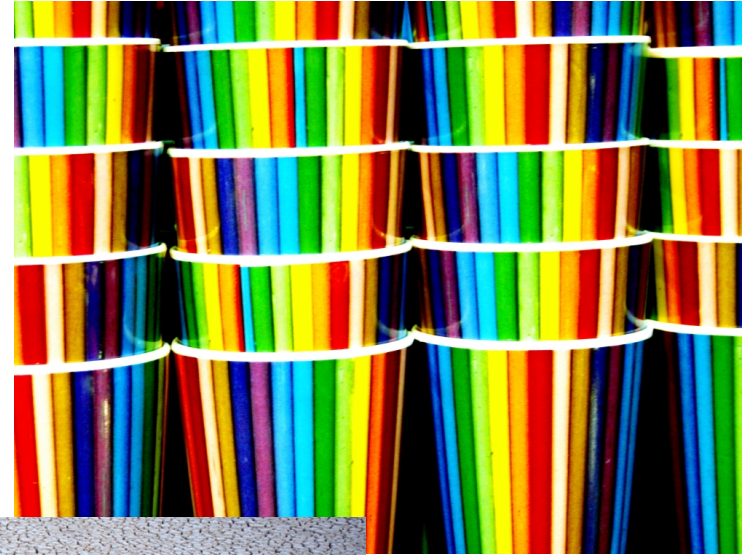
Local features - Applications

- Feature points are used for:
 - Retrieval
 - Tracking
 - Image alignment
 - 3D reconstruction
 - Object recognition
 - Object detection
 - Attribute recognition
 - Indexing and search
 - Robot navigation

Feature Types



Shape!



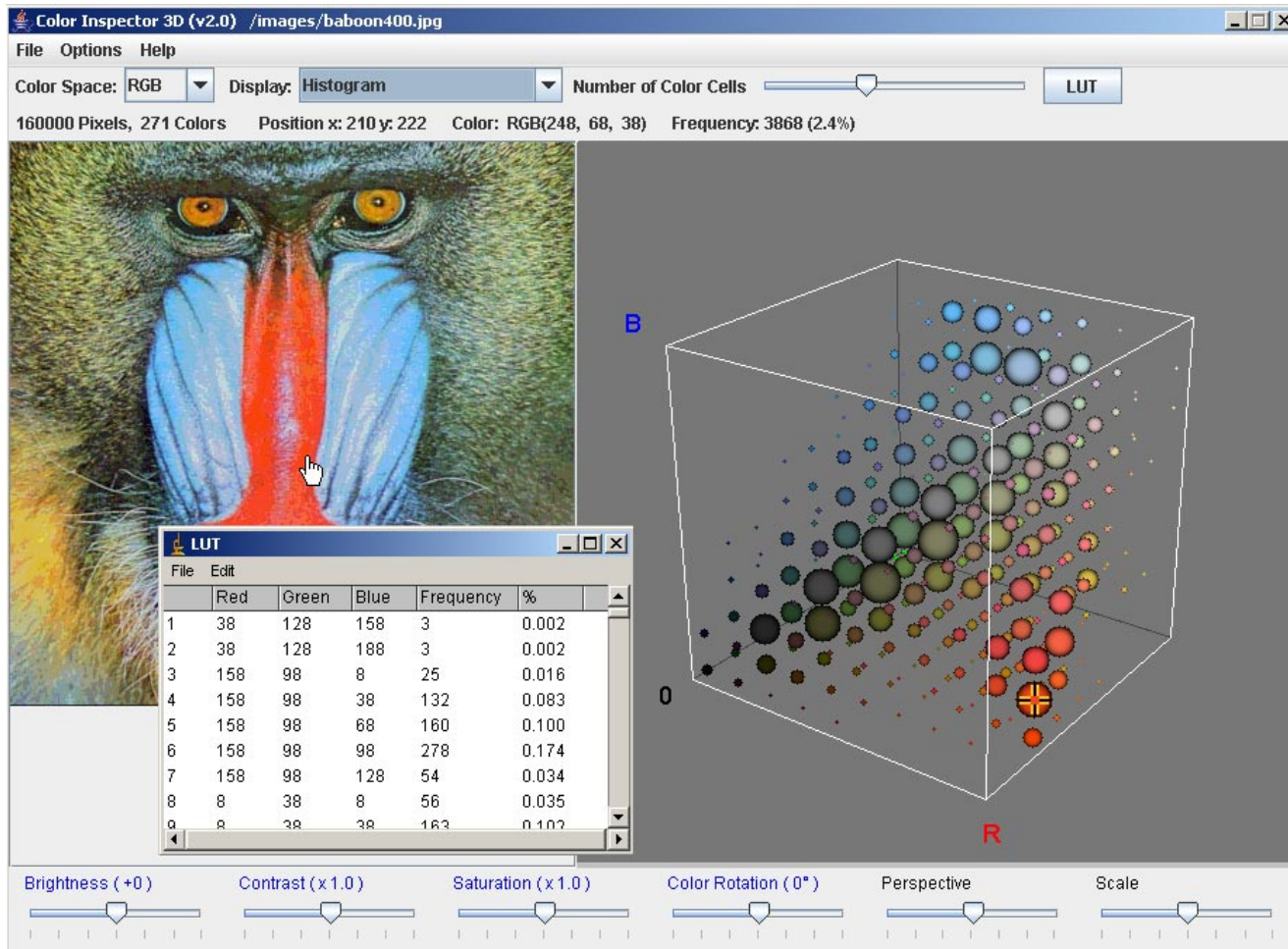
Color!



Texture!

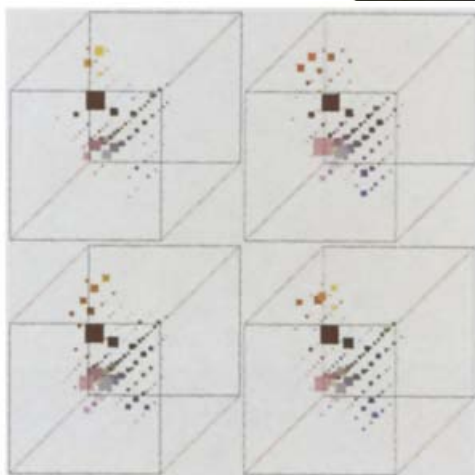
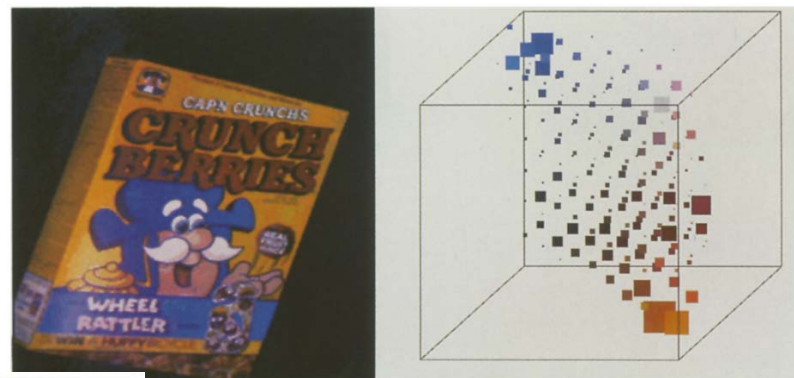
Color Features

Color Histograms



Uses of color in computer vision

Color histograms for indexing and retrieval

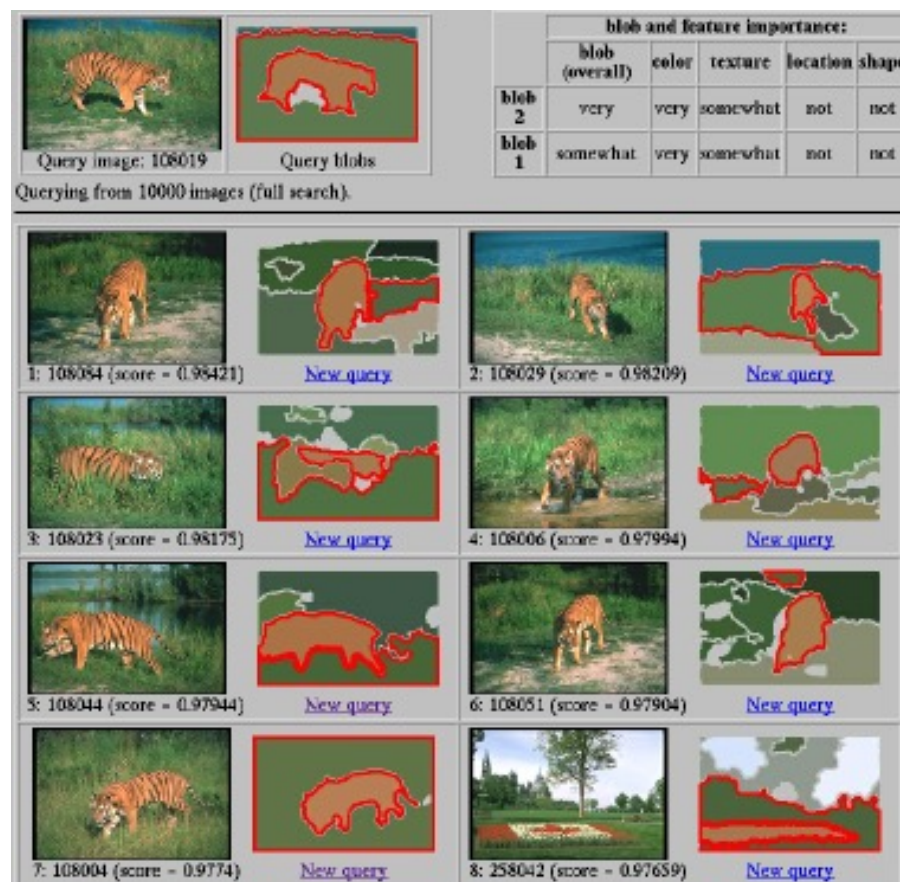


Swain and Ballard, [Color Indexing](#), IJCV 1991.

source: Svetlana Lazebnik

Uses of color in computer vision

Image segmentation and retrieval



C. Carson, S. Belongie, H. Greenspan, and Ji. Malik, Blobworld: Image segmentation using Expectation-Maximization and its application to image querying, ICVIS 1999.

Uses of color in computer vision

Skin detection



M. Jones and J. Rehg, [Statistical Color Models with Application to Skin Detection](#), IJCV 2002.

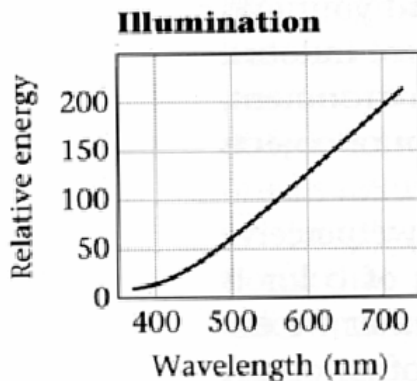
Color features

Pros/Cons?

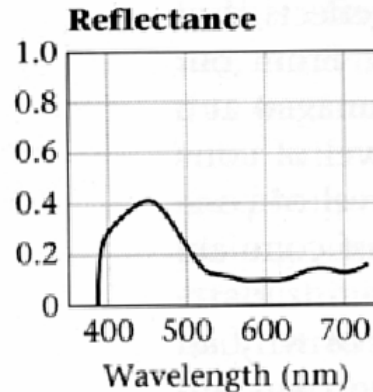
Interaction of light and surfaces



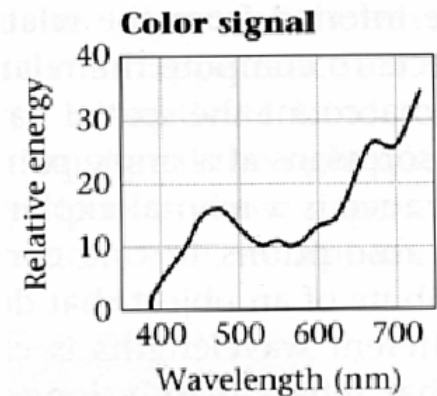
- Reflected color is the result of interaction of light source spectrum with surface reflectance



• *



=



Interaction of light and surfaces



Olafur Eliasson, *Room for one color*, 1997

“In Olafur Eliasson’s *Room for one color* (1997) we actually only see one color. The room is illuminated by yellow monofrequency light. In contrast to ordinary white light containing the full colour spectrum, it consists of a single wavelength at the yellow end of the spectrum. The room appears colorless because the light reduces all the other colours to a scale comprised of yellow. On the other hand, everything stands out with crystal sharpness because there is far less information than usual for the eye to process. After a while, the eye compensates for this deficiency by generating an excess of the missing primaries of red and blue – which together yield violet. So when we move on into the adjacent space we begin by seeing everything as imbued with violet.”

Shape Features

Invariance

- We might want our features to be robust to variations that can be present in images.

Types of invariance

- Illumination



Types of invariance

- Illumination
- Scale



Slide source: Tom Duerig

Types of invariance

- Illumination
- Scale
- Rotation



Slide source: Tom Duerig

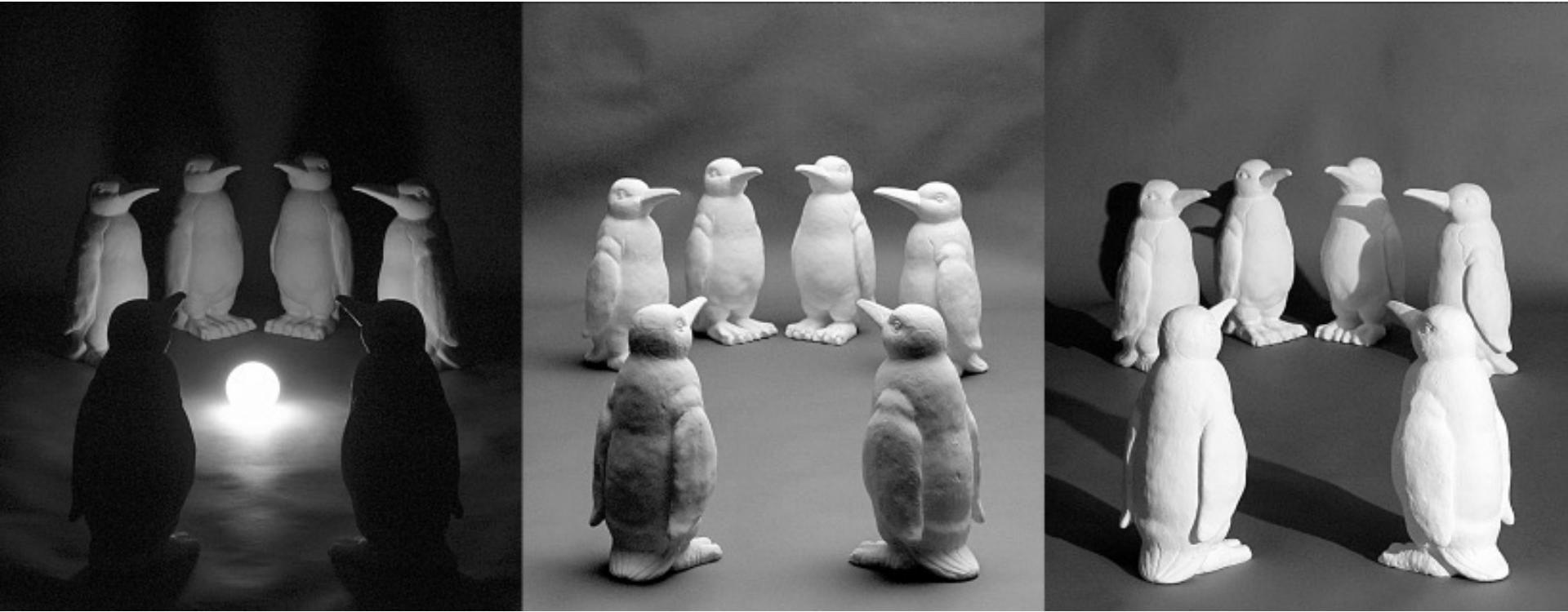
Types of invariance

- Illumination
- Scale
- Rotation
- Affine



Slide source: Tom Duerig

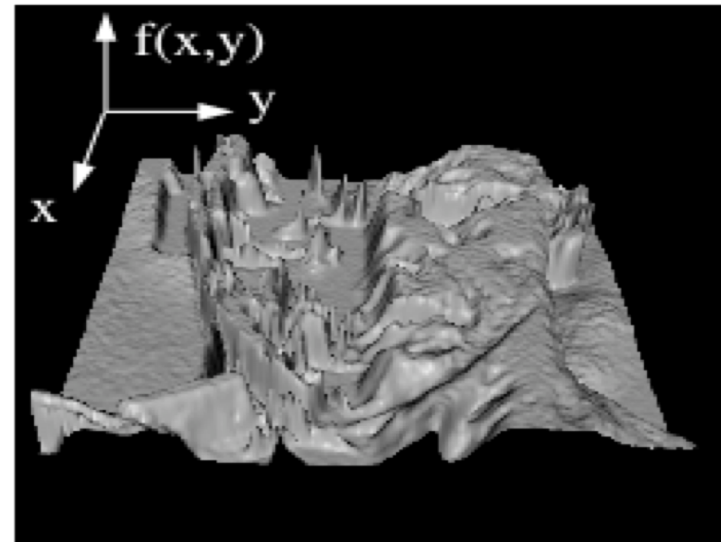
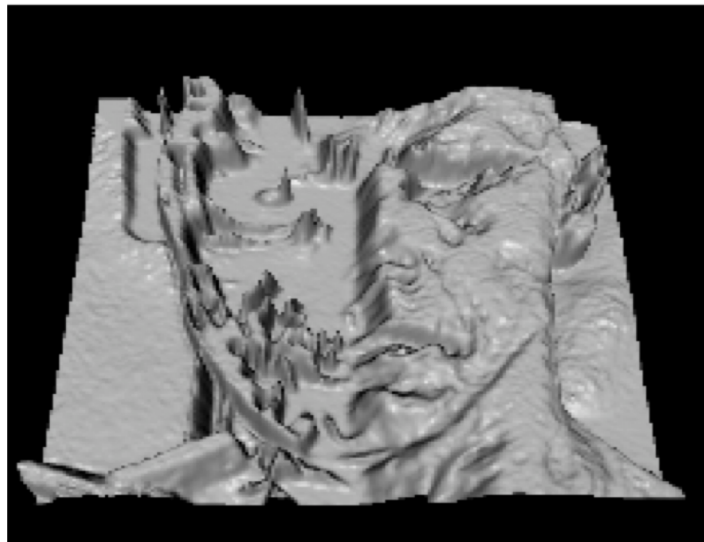
Challenges: illumination



Illumination changes brightness of pixels. What remains unchanged?

Modeling images as functions

Images as functions



Intensity given position

Source: S. Seitz

Images as functions

- We can think of an **image** as a function, f , from \mathbb{R}^2 to \mathbb{R} :
 - $f(x, y)$ gives the **intensity** at position (x, y)
- A color image is just three functions pasted together. We can write this as a “vector-valued” function:

$$f(x, y) = \begin{bmatrix} r(x, y) \\ g(x, y) \\ b(x, y) \end{bmatrix}$$

Motivation 1: Noise reduction

Given a camera and a still scene, how can you reduce noise?



Take lots of images and average them!

What's the next best thing?

First attempt at a solution

- Let's replace each pixel with an average of all the values in its neighborhood

Moving Average In 2D

$F[x, y]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$G[x, y]$

	0								

Moving Average In 2D

$F[x, y]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$G[x, y]$

	0	10							

Moving Average In 2D

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20						

Moving Average In 2D

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20	30					

Moving Average In 2D

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20	30	30				

Moving Average In 2D

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20	30	30	30	20	10	
	0	20	40	60	60	60	40	20	
	0	30	60	90	90	90	60	30	
	0	30	50	80	80	90	60	30	
	0	30	50	80	80	90	60	30	
	0	20	30	50	50	60	40	20	
	10	20	30	30	30	30	20	10	
	10	10	10	0	0	0	0	0	

Generalization of moving average (Convolution)

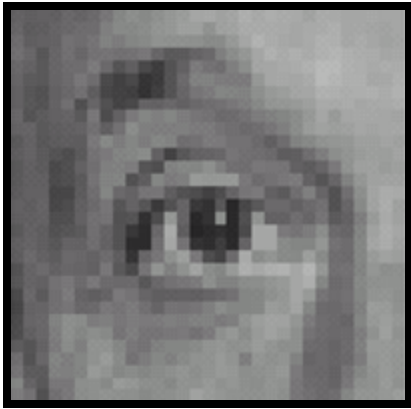
- Let's replace each pixel with a *weighted* average of its neighborhood
- The weights are called the *filter kernel*
- What are the weights for a 3x3 moving average?

Generalization of moving average (Convolution)

- Let's replace each pixel with a *weighted* average of its neighborhood
- The weights are called the *filter kernel*
- What are the weights for a 3x3 moving average?

1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

Practice with linear filters

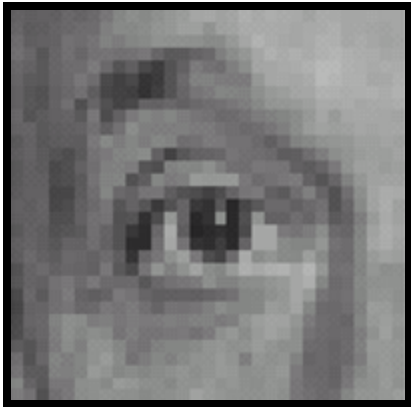


Original

0	0	0
0	1	0
0	0	0

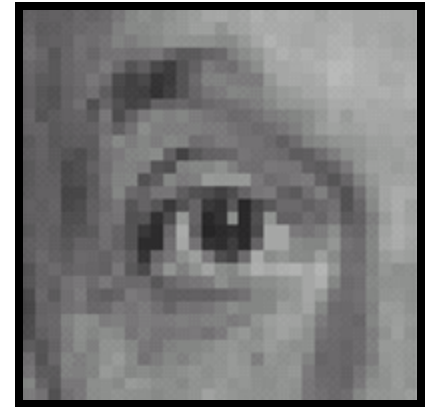
?

Practice with linear filters



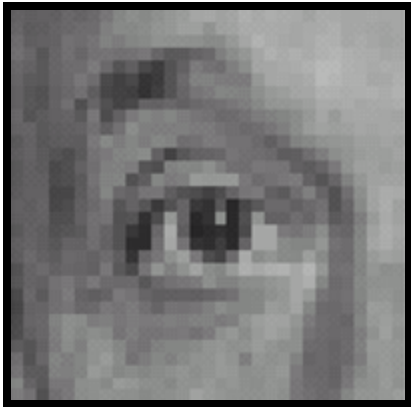
Original

0	0	0
0	1	0
0	0	0



Filtered
(no change)

Practice with linear filters

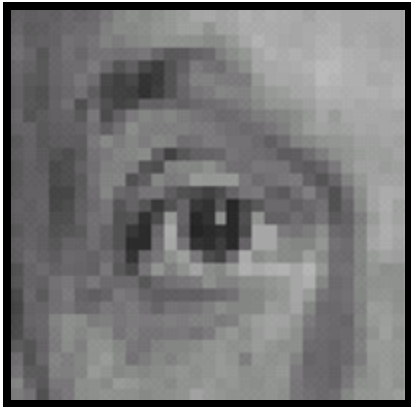


Original

0	0	0
0	0	1
0	0	0

?

Practice with linear filters



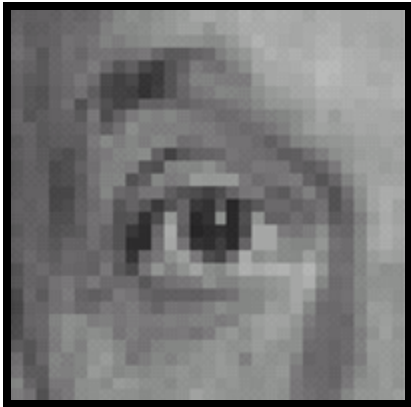
Original

0	0	0
0	0	1
0	0	0



Shifted left
By 1 pixel

Practice with linear filters



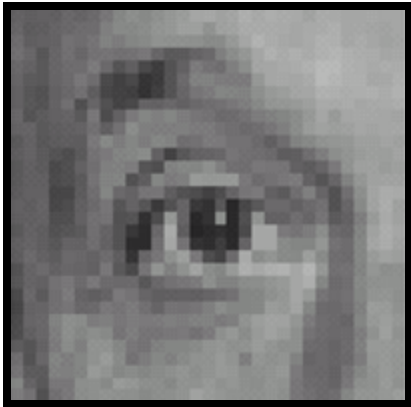
Original

 $\frac{1}{9}$

1	1	1
1	1	1
1	1	1

?

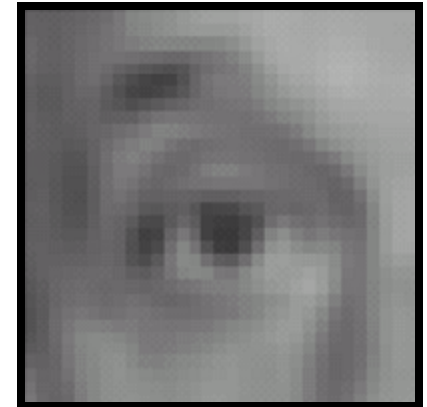
Practice with linear filters



Original

$$\frac{1}{9}$$

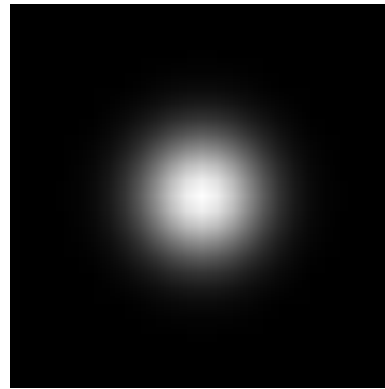
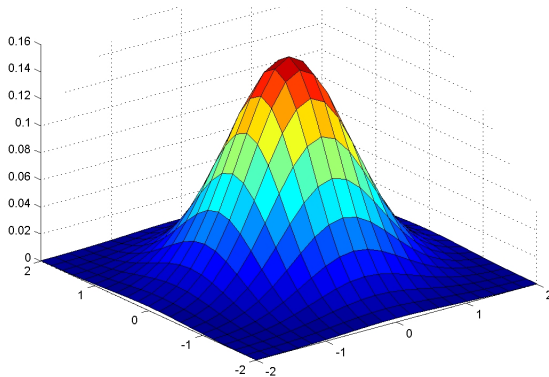
1	1	1
1	1	1
1	1	1



Blur (with a
box filter)

Gaussian Kernel

$$G_{\sigma} = \frac{1}{2\pi\sigma^2} e^{-\frac{(x^2+y^2)}{2\sigma^2}}$$

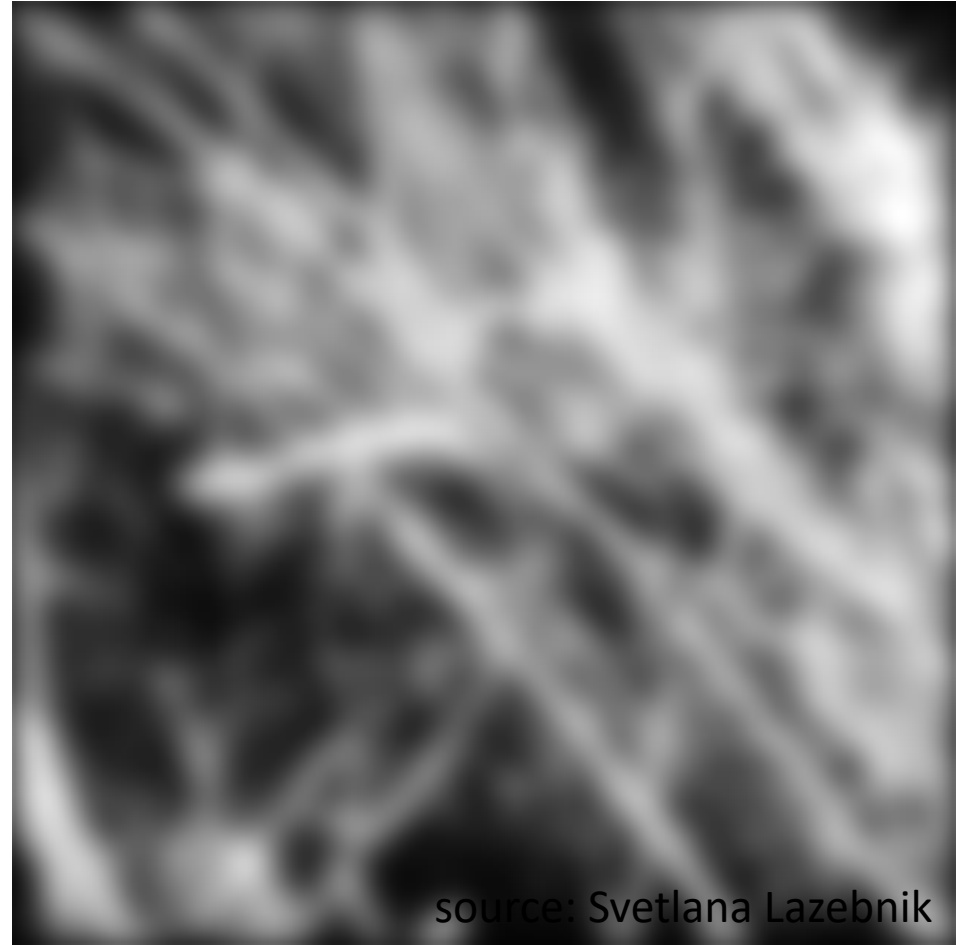
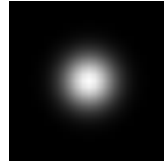


0.003	0.013	0.022	0.013	0.003
0.013	0.059	0.097	0.059	0.013
0.022	0.097	0.159	0.097	0.022
0.013	0.059	0.097	0.059	0.013
0.003	0.013	0.022	0.013	0.003

5 x 5, $\sigma = 1$

- Constant factor at front makes volume sum to 1 (can be ignored, as we should re-normalize weights to sum to 1 in any case)

Example: Smoothing with a Gaussian



source: Svetlana Lazebnik

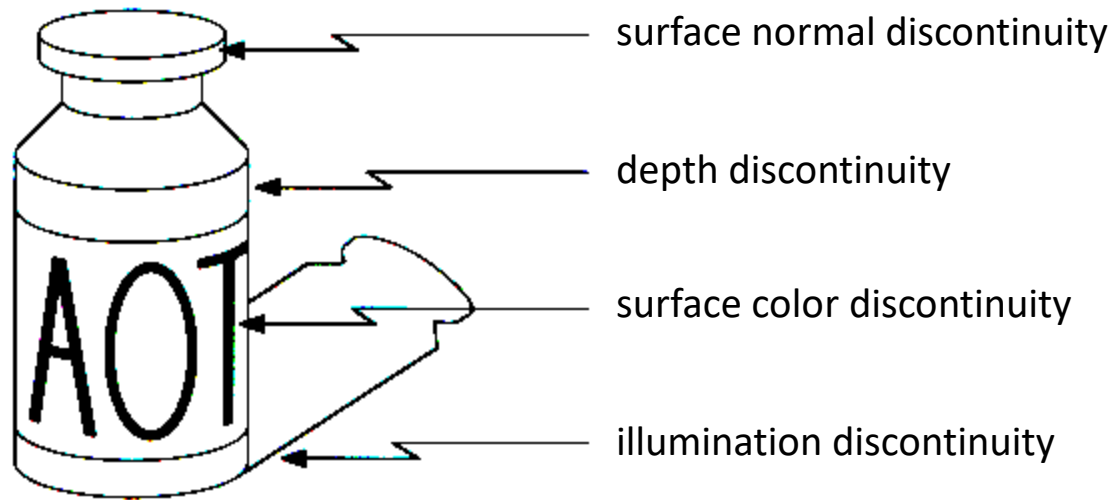
Edges

Edge detection

- **Goal:** Identify sudden changes (discontinuities) in an image
 - Intuitively, most semantic and shape information from the image can be encoded in the edges
 - More compact than pixels
- **Ideal:** artist's line drawing (but artist is also using object-level knowledge)



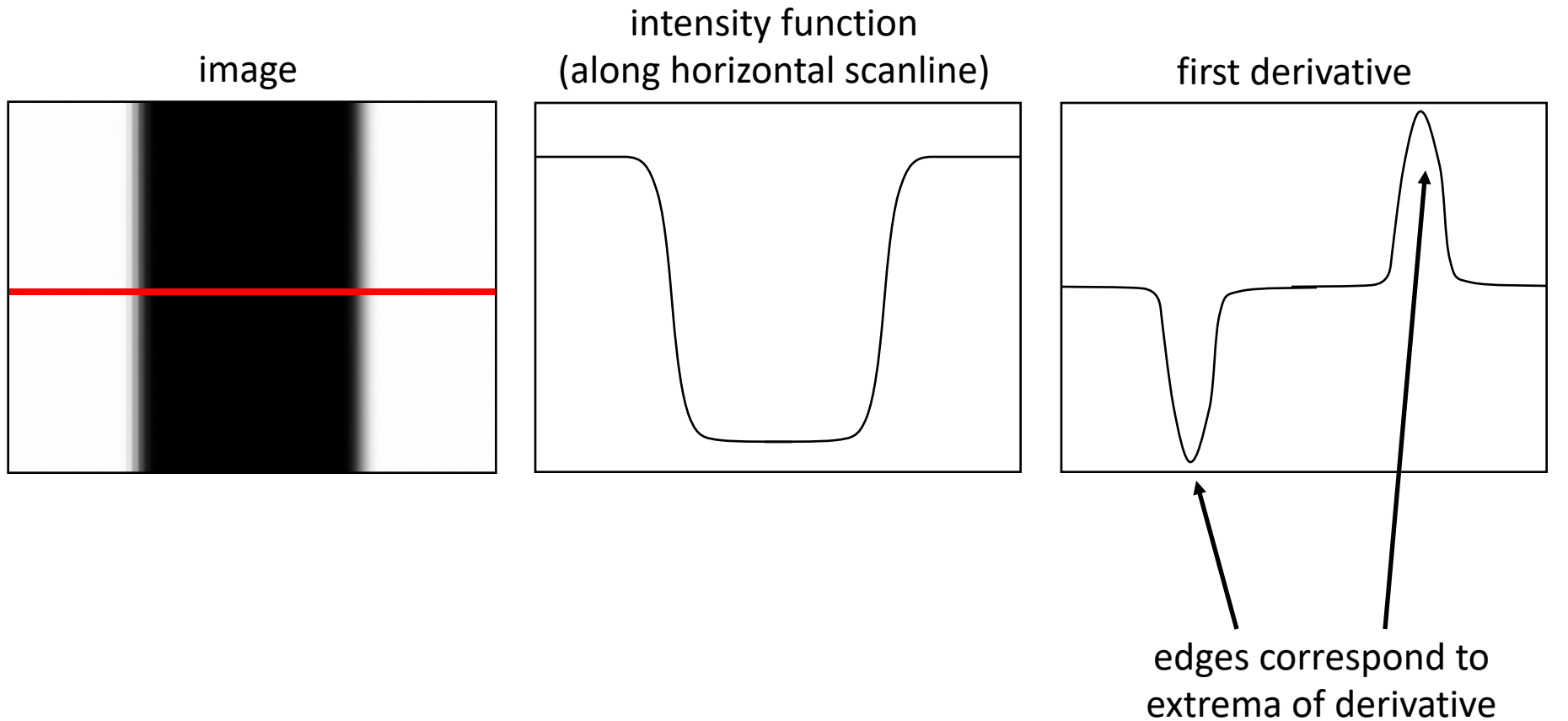
Origin of Edges



Edges are caused by a variety of factors

Characterizing edges

- An edge is a place of rapid change in the image intensity function



Edge filters

Approximations of derivative filters:

Prewitt: $M_x =$

-1	0	1
-1	0	1
-1	0	1

 ; $M_y =$

1	1	1
0	0	0
-1	-1	-1

Sobel: $M_x =$

-1	0	1
-2	0	2
-1	0	1

 ; $M_y =$

1	2	1
0	0	0
-1	-2	-1

Roberts: $M_x =$

0	1
-1	0

 ; $M_y =$

1	0
0	-1

Convolve filter with image to get edge map

Edge filters

Approximations of derivative filters:

Prewitt:

$$M_x = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

;

$$M_y = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

Sobel:

$$M_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

;

$$M_y = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

Roberts:

$$M_x = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$

;

$$M_y = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

Edge filters

Approximations of derivative filters:

Prewitt:

$$M_x = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

;

$$M_y = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

Sobel:

$$M_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

;

$$M_y = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

Roberts:

$$M_x = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$

;

$$M_y = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

Respond highly to vertical edges

Edge filters

Approximations of derivative filters:

Prewitt: $M_x = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$; $M_y = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

Sobel: $M_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$; $M_y = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$

Roberts: $M_x = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$; $M_y = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$

Edge filters

Approximations of derivative filters:

Prewitt: $M_x =$

-1	0	1
-1	0	1
-1	0	1

;

$$M_y =$$

1	1	1
0	0	0
-1	-1	-1

Sobel: $M_x =$

-1	0	1
-2	0	2
-1	0	1

;

$$M_y =$$

1	2	1
0	0	0
-1	-2	-1

Roberts: $M_x =$

0	1
-1	0

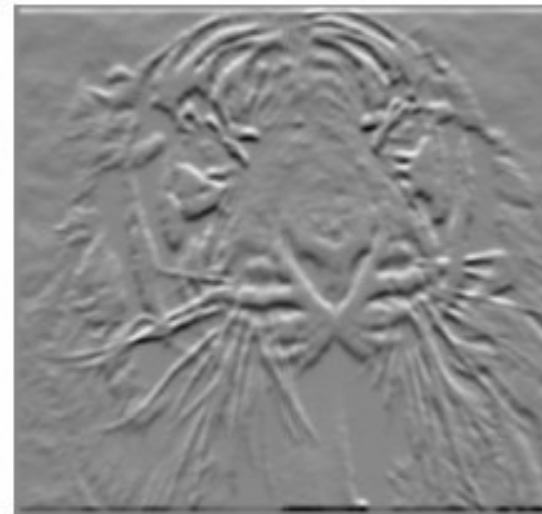
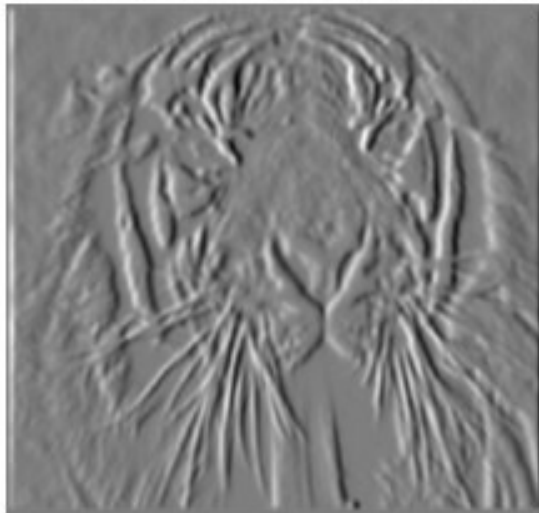
;

$$M_y =$$

1	0
0	-1

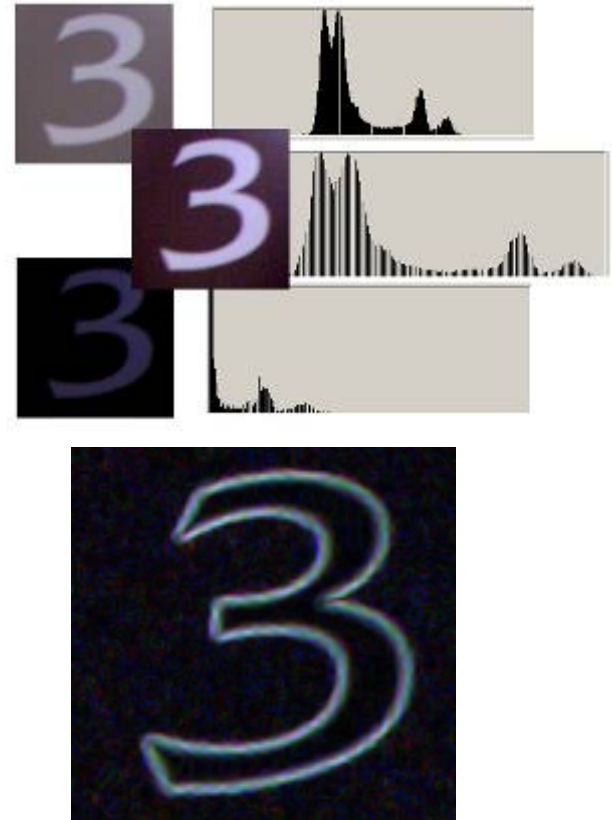
Respond highly to horizontal edges

Edges: example



How to achieve illumination invariance

- Use edges instead of raw values



Feature types!



The representation of these two umbrella's should be similar....
Under a color based representation they look completely different!

Edges



Red umbrella



Gray umbrella

Edges extracted using convolution with Prewitt filter

Edges



Edges overlaid from red and gray umbrellas.

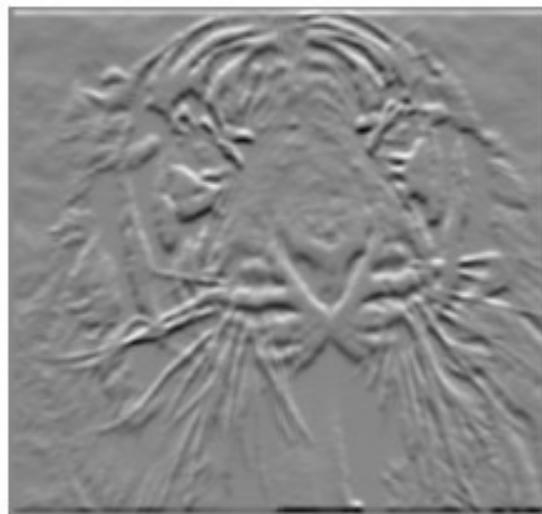
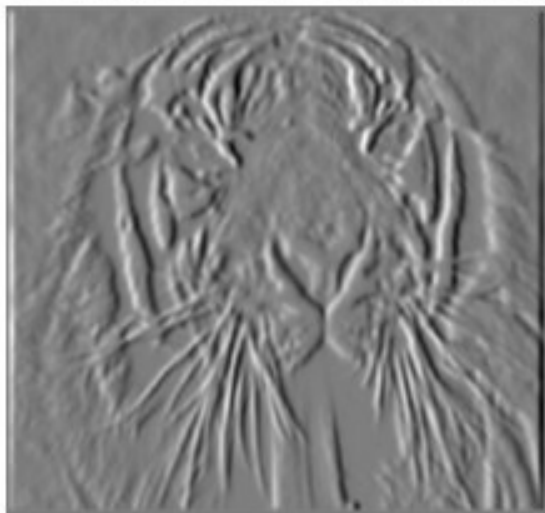
Local Features

Feature points (locations) + feature descriptors:

- 1) Where should we put the features?
- 2) How should we describe the features

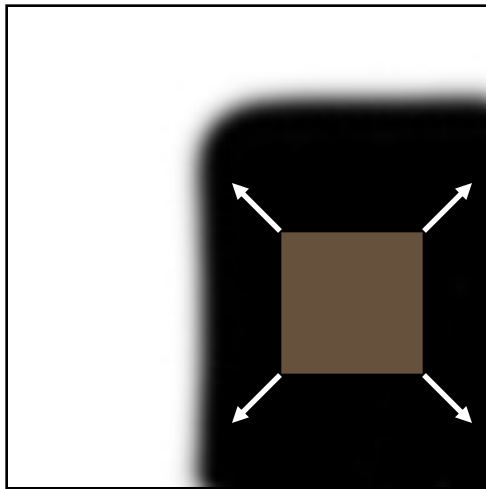
Where should we put features?

Where to put features?

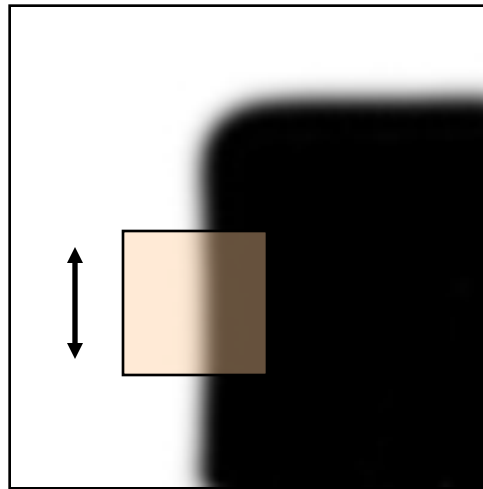


Distinctiveness in x,y

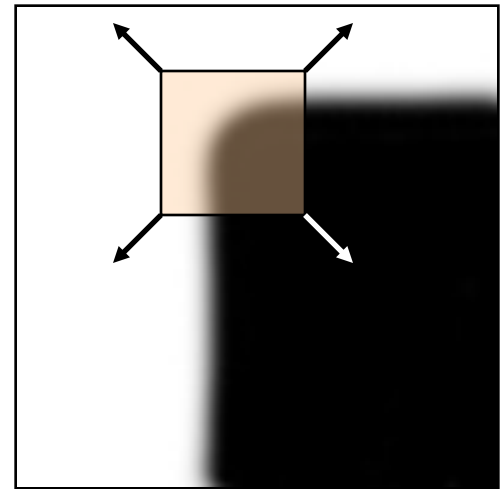
- We should easily recognize the point by looking through a small window
- Shifting a window in *any direction* should give a *large change* in intensity



“flat” region:
no change in
all directions



“edge”:
no change
along the edge
direction



“corner”:
significant
change in all
directions

Corner Detections

C.Harris and M.Stephens. ["A Combined Corner and Edge Detector."](#)
Proceedings of the 4th Alvey Vision Conference: pages 147--151.



source: Svetlana Lazebnik

Toward denser sampling

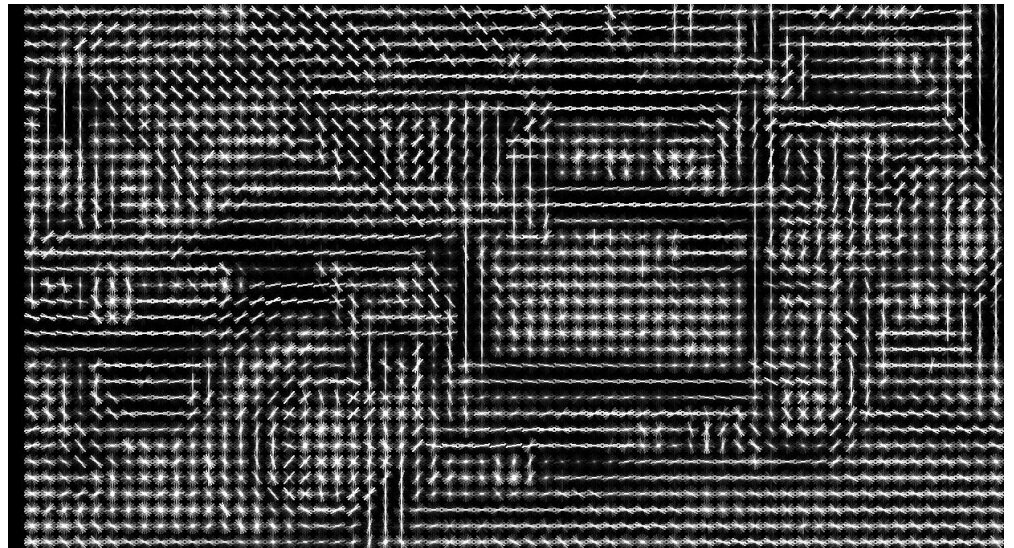


Just sample points randomly on edges

Fully dense sampling



Or sample everywhere!



How should we describe features?

SIFT

- SIFT (Scale Invariant Feature Transform) - stable robust and distinctive local features
- One of the most popular shape (edge) based features

Where to put them:

- Originally sampled points distinctive in both position (x,y) and scale (SI – Scale Invariant)
- Often now used with sampling on a dense grid

Feature descriptor

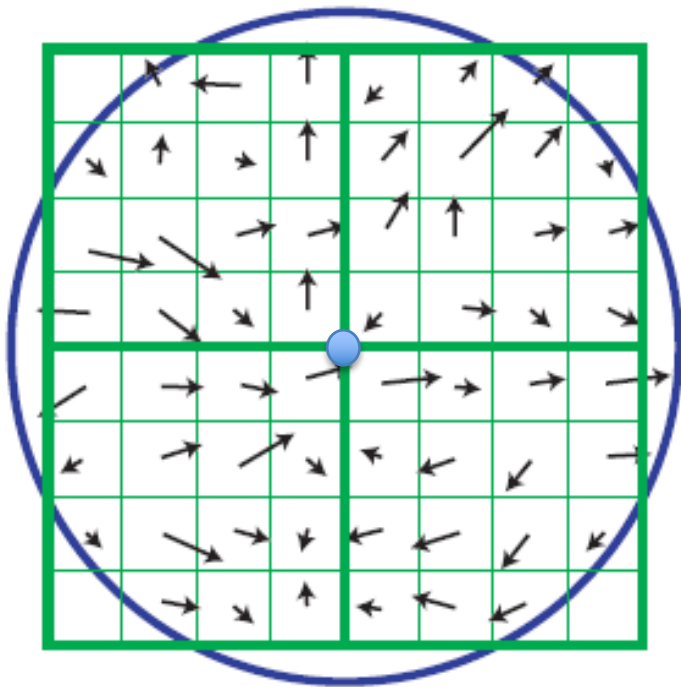
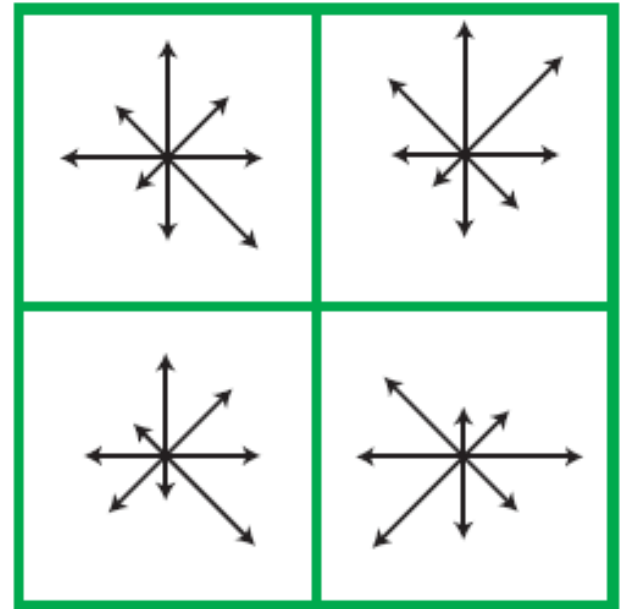


Image gradients



Keypoint descriptor

Feature descriptor

- Based on 16×16 patches
- 4×4 subregions
- 8 orientation bins in each subregion
- $4 \times 4 \times 8 = 128$ dimensions in total

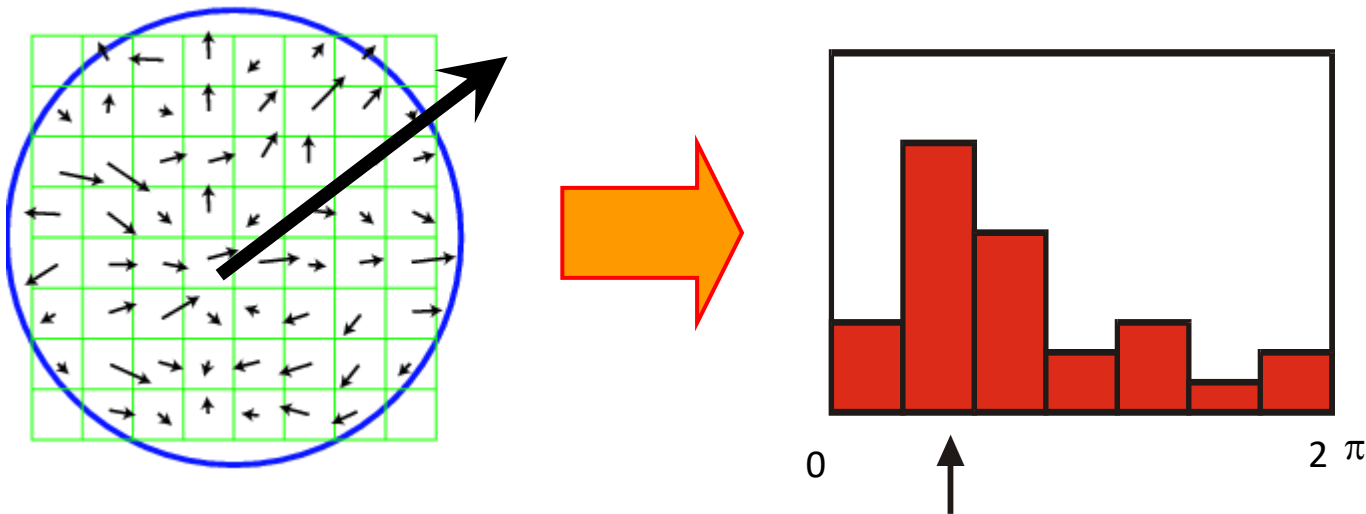
Rotation Invariance

- Rotate all features to go the same way in a determined manner
- A histogram is formed by quantizing the orientations into bins
- Features are rotated according to peak of orientation histogram

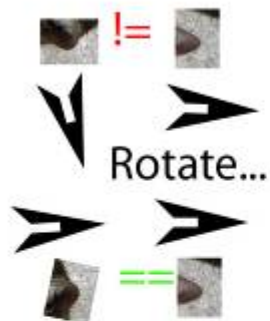
Eliminating rotation ambiguity

To assign a unique orientation:

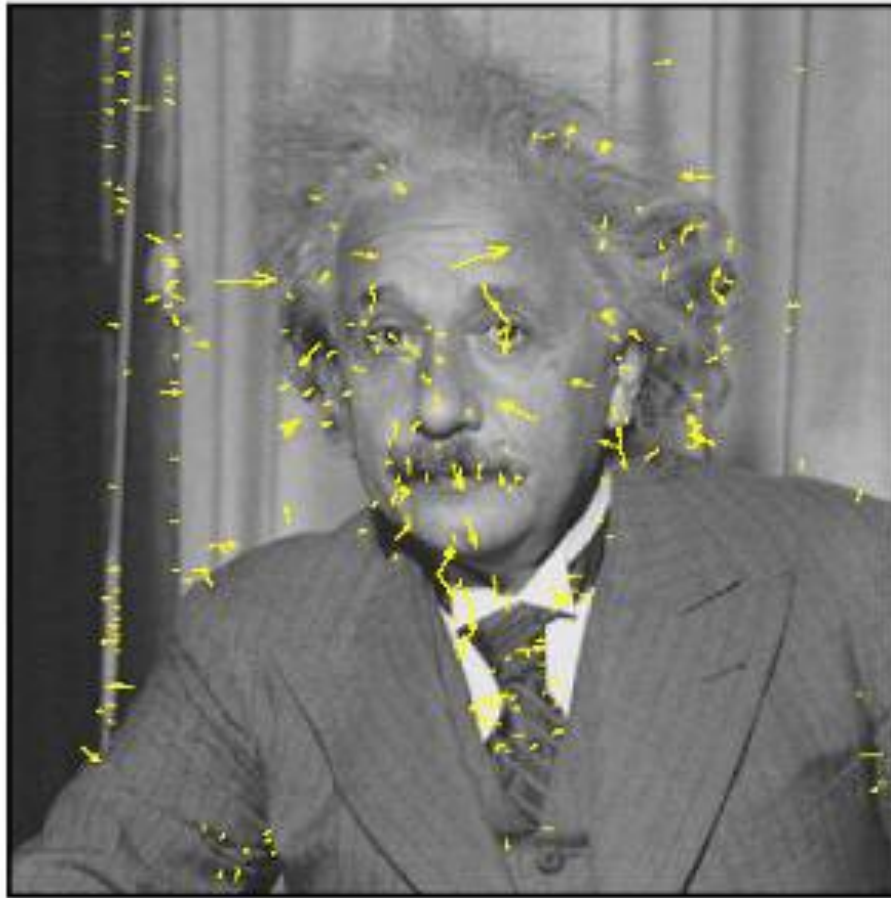
- Create histogram of local gradient directions in the patch
- Assign canonical orientation at peak of smoothed histogram



Rotation Invariance

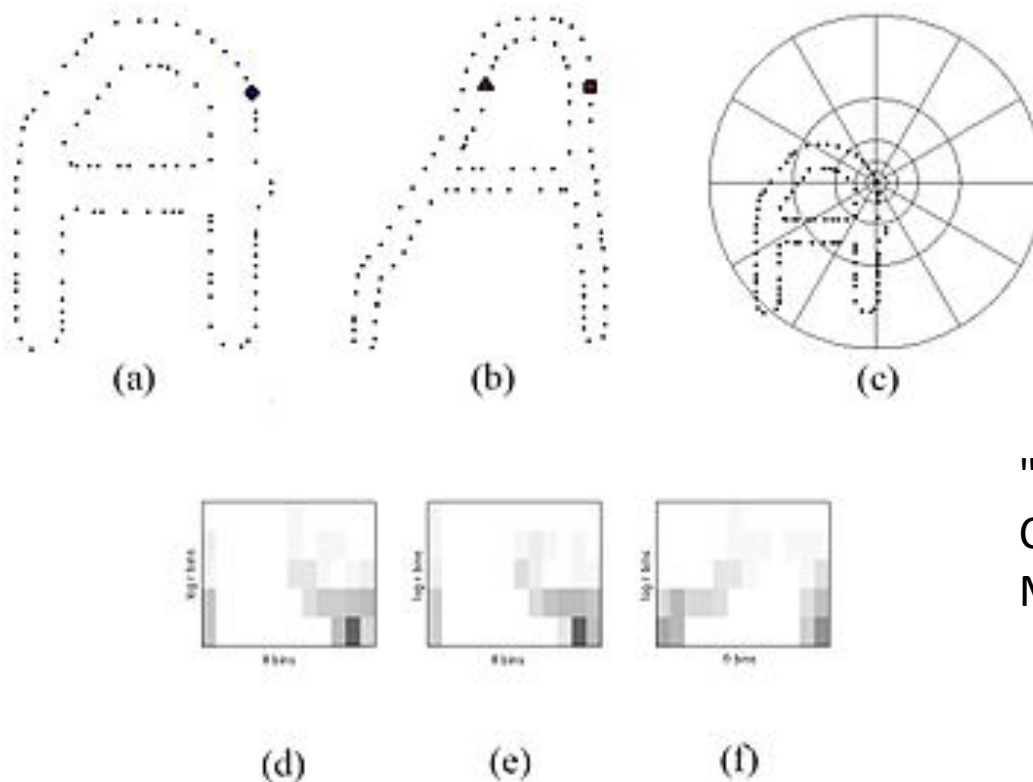


SIFT output



Other shape features

Shape Context

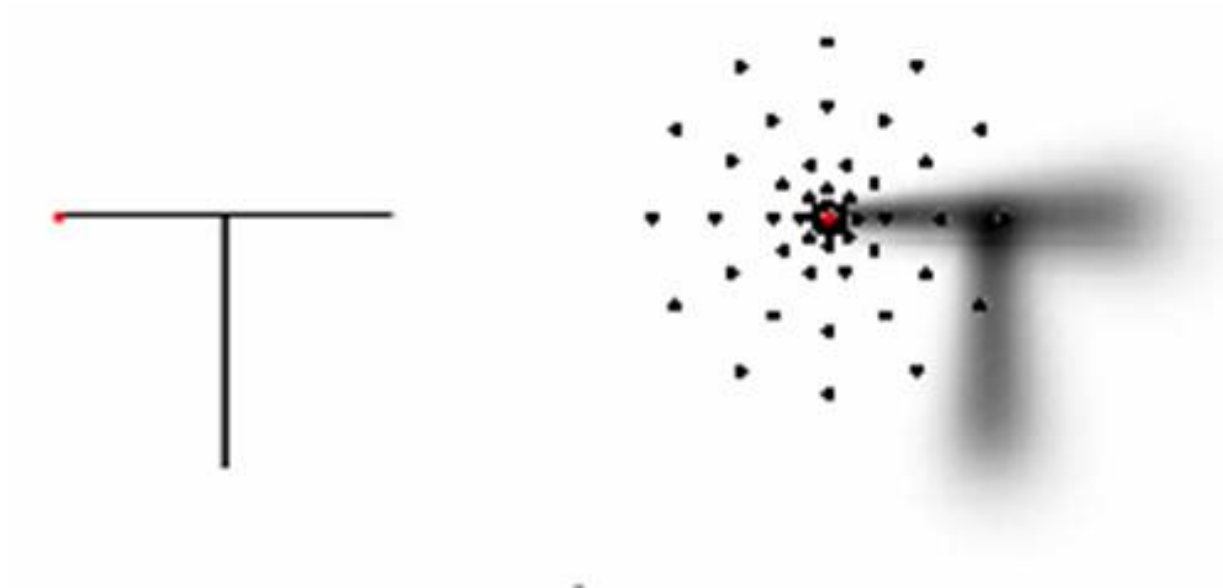


"Matching with Shape Contexts": Belongie, Malik 2000

- (a) and (b) are the sampled edge points of the two shapes. (c) is the diagram of the log-polar bins used to compute the shape context. (d) is the shape context for the circle, (e) is that for the diamond, and (f) is that for the triangle. As can be seen, since (d) and (e) are the shape contexts for two closely related points, they are quite similar, while the shape context in (f) is very different.

Geometric Blur

“Geometric Blur for Template Matching” A. Berg & Malik, 2001



Geometric Blur

“Geometric Blur for Template Matching” A. Berg & Malik, 2001



**Extract
Sparse
Channels**



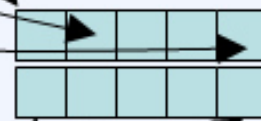
Sparse
Non-Negative
Channels
(eg. oriented edge
energy)

**Apply
Spatially
Varying
Blur**



Geometric Blur
of each
channel

Subsample

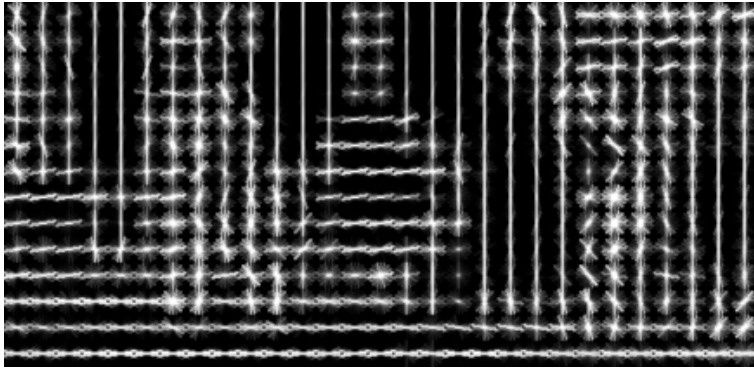


Geometric Blur
descriptor

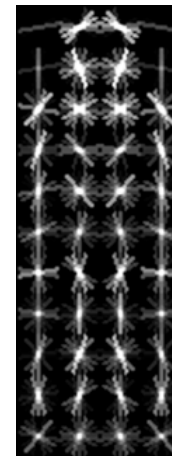
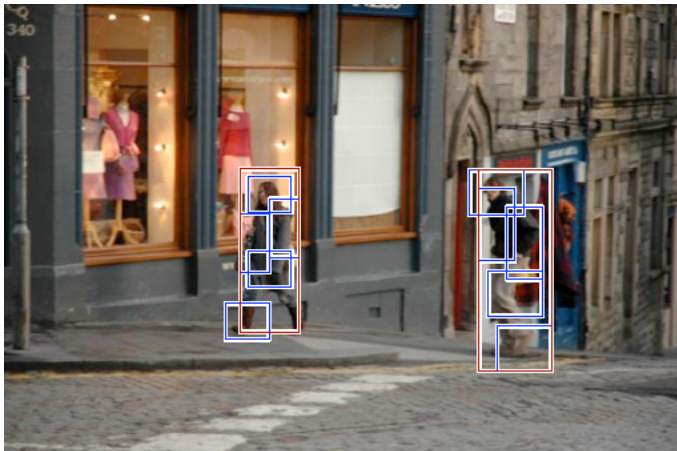
Histograms of Oriented Gradients

Navneet Dalal and Bill Triggs, 2005

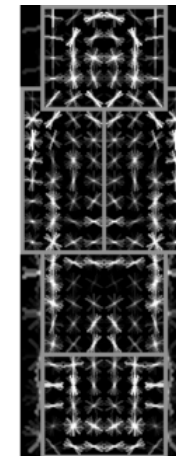
HOG feature map



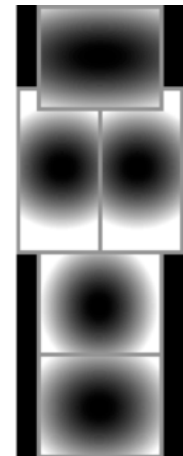
- Similar to SIFT
- Computed on dense grid of uniformly spaced cells
- Good exploration of parameter space (gradient scale, orientation binning, spatial binning, contrast normalization)



Root
filter



Part
filters

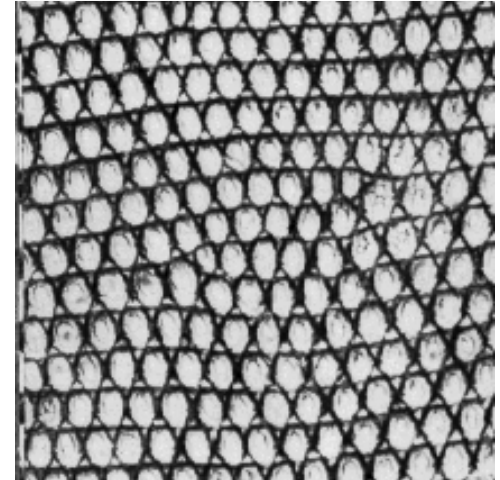
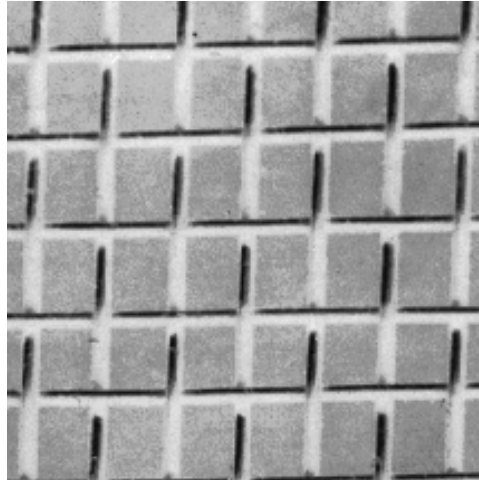
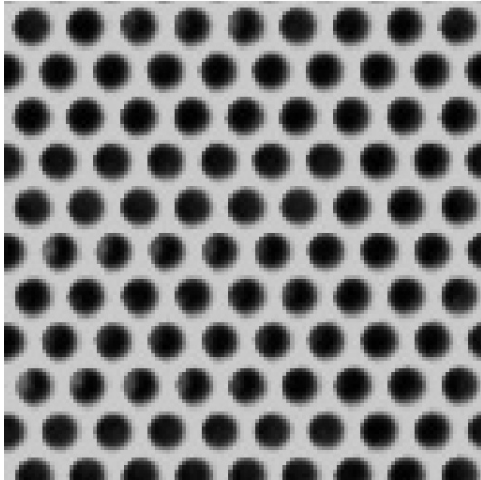


Deformation
weights

Texture Features

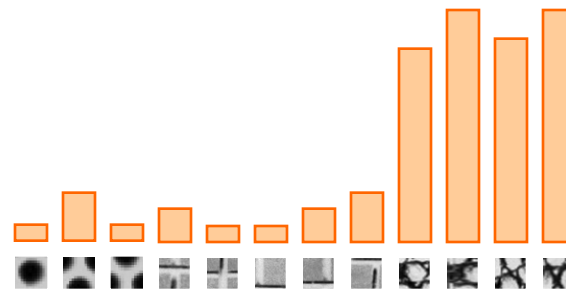
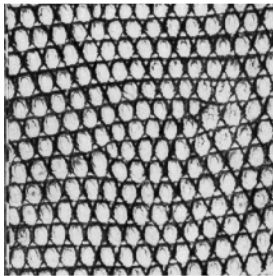
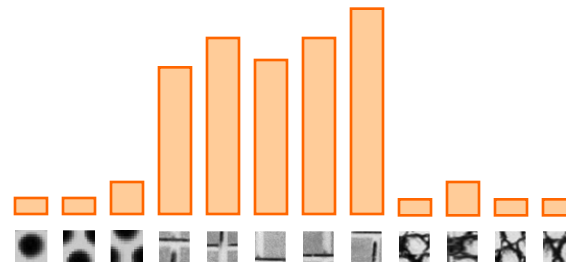
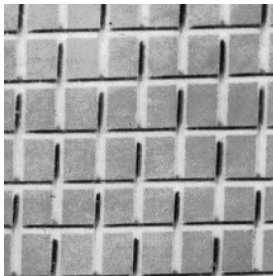
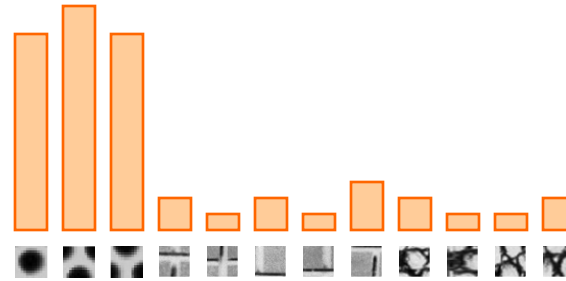
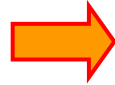
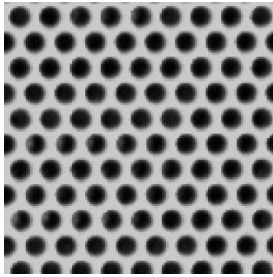
Texture Features

- Texture is characterized by the repetition of basic elements or *textons*
- Generally, it is the identity of the textons, not their spatial arrangement, that matters



Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

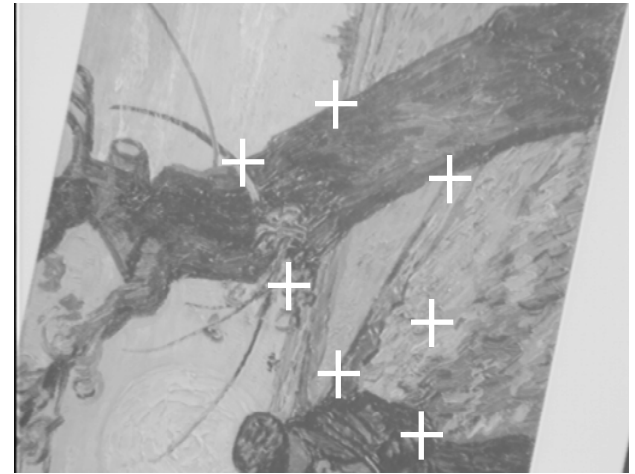
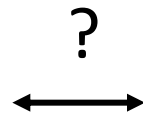
Texture Histograms



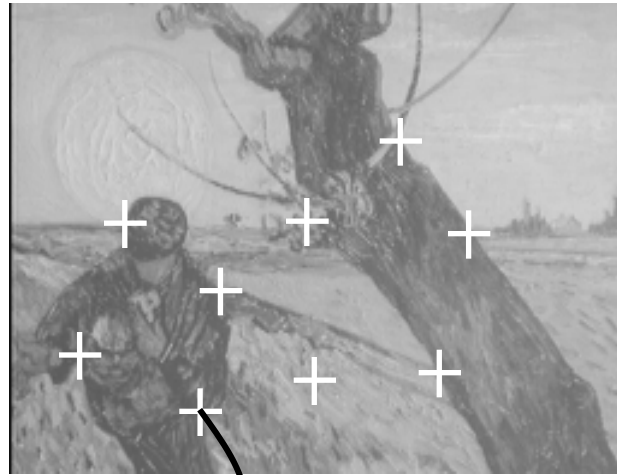
Can be computed locally or globally

What can we do with features?

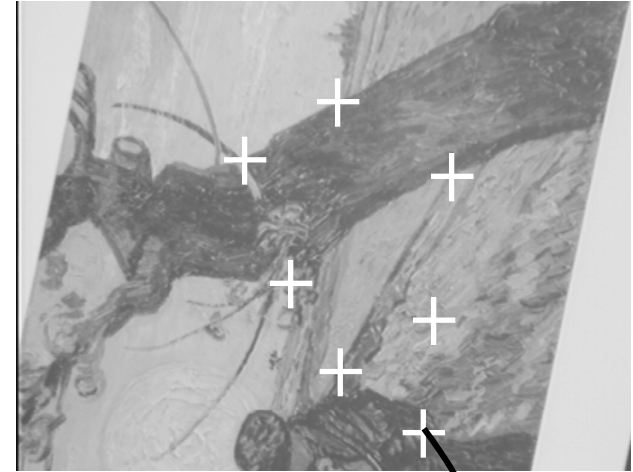
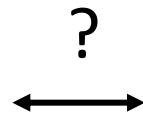
Finding correspondences



Finding correspondences




feature
descriptor



?
=


feature
descriptor

- Need to compare *feature descriptors* of local patches surrounding interest points

Comparing Feature descriptors

How to compare two such vectors?

- Sum of squared differences (SSD)

$$\text{SSD}(u, v) = \sum_i (u_i - v_i)^2$$

- Not invariant to intensity change

- Normalized correlation

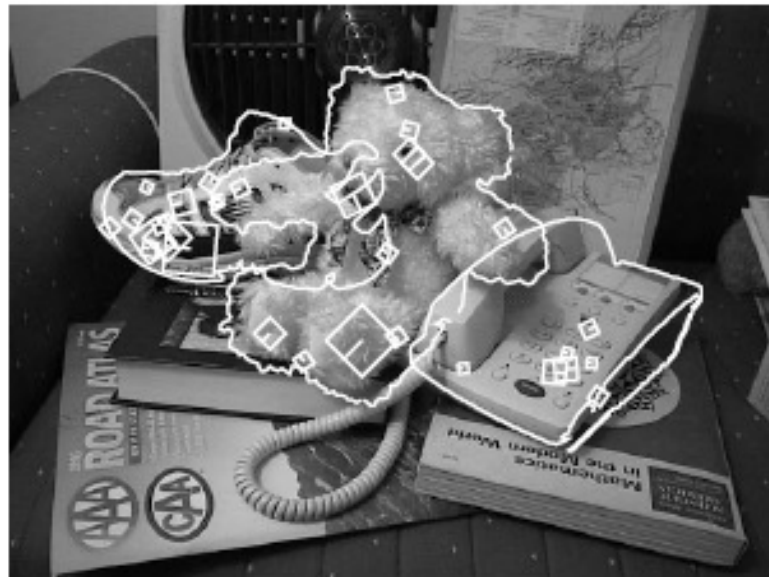
$$\rho(u, v) = \frac{\sum_i (u_i - \bar{u})(v_i - \bar{v})}{\sqrt{\left(\sum_j (u_j - \bar{u})^2 \right) \left(\sum_j (v_j - \bar{v})^2 \right)}}$$

- Invariant to intensity change

Example: object recognition

- The SIFT features of training images are extracted and stored
- For a query image
 1. Extract SIFT feature
 2. Find nearest neighbor match
 3. Given 3 keypoint matches, perform geometric verification





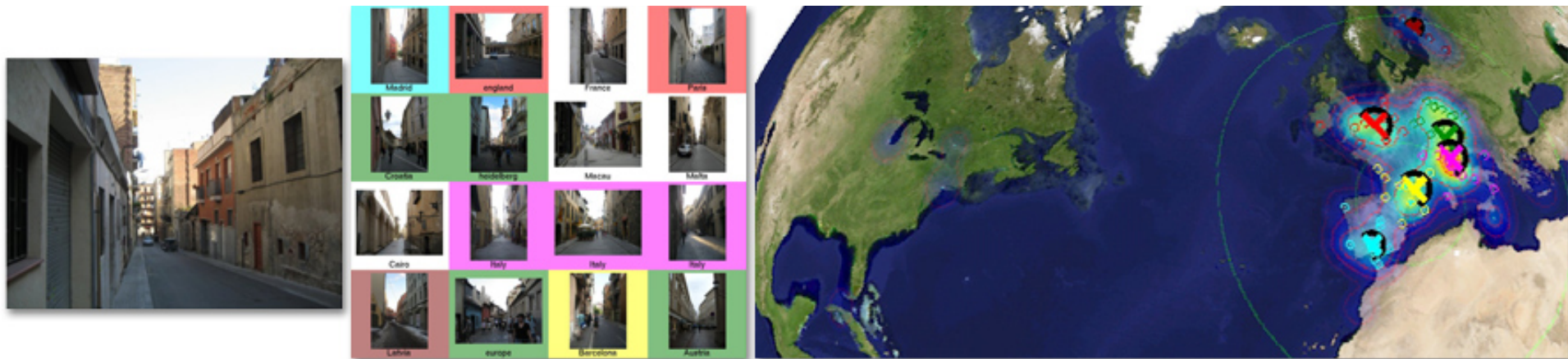
Useful Code

VLFEAT (<http://www.vlfeat.org/index.html>)

open source library implementing various feature descriptors,
clustering, matching algorithms

Some fun applications of
features + matching

NN for estimating location



J. Hays and A. Efros, IM2GPS: estimating geographic information from a single image, CVPR 2008

Where?



What can you say about where these photos were taken?

How?

Collect a large collection of geo-tagged photos

6.5 million images with both GPS coordinates and geographic keywords, removing images with keywords like birthday, concert, abstract, ...

Test set – 400 randomly sampled images from this collection. Manually removed abstract photos and photos with recognizable people – 237 test photos.

Nearest Neighbor Matching

For each input image compute features (color, texture, shape)

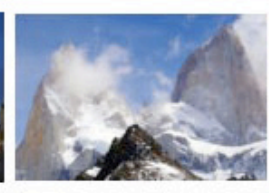
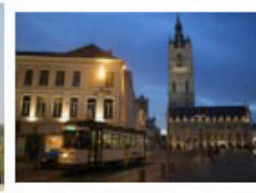
Compute distance in feature space to all 6 million images in the database (each feature contributes equally).

Label the image with GPS coordinates of:

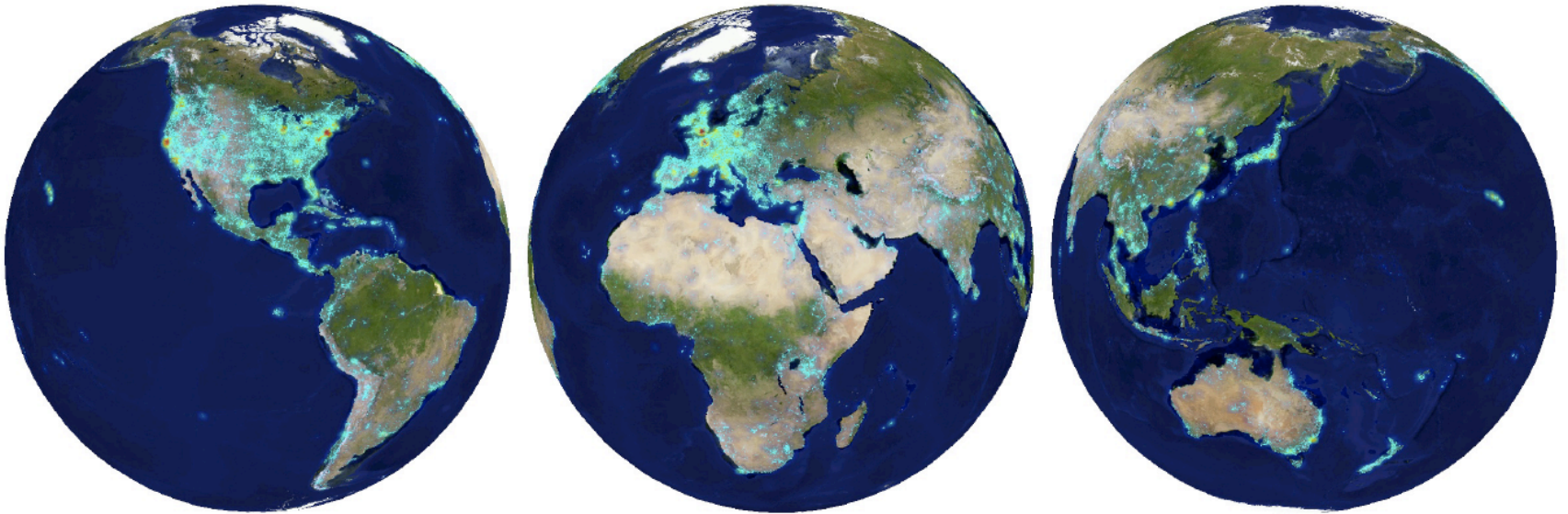
- 1 nearest neighbor

- k=120 nearest neighbors – probability map over entire globe.

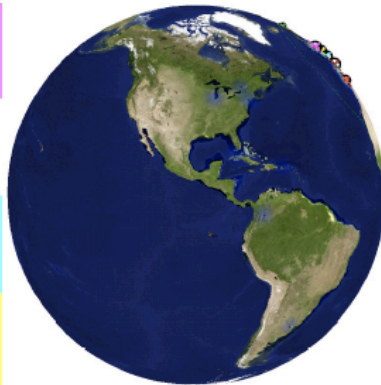
Test Images



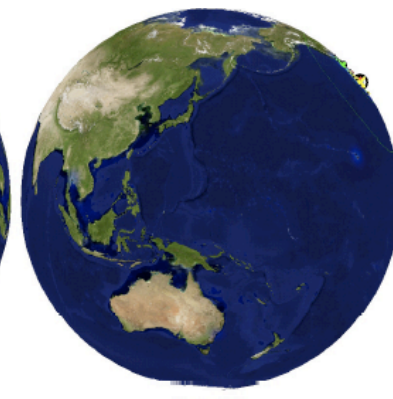
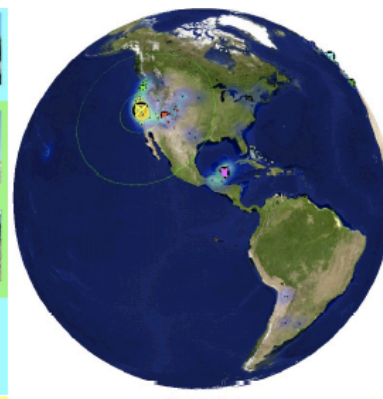
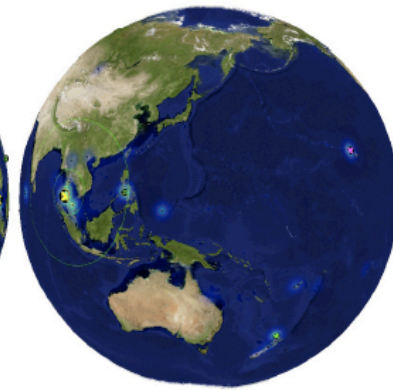
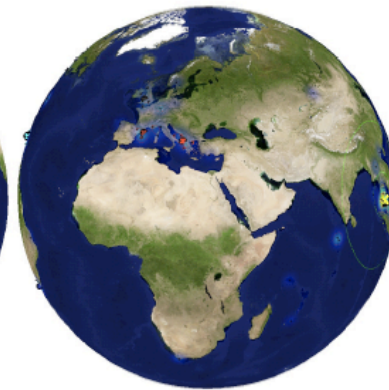
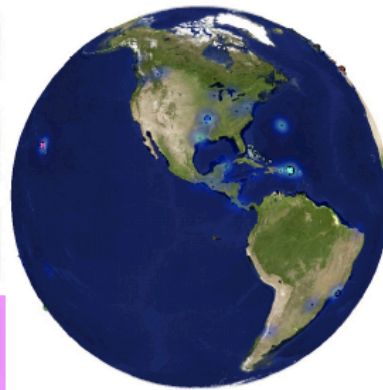
Distribution of photos



Results



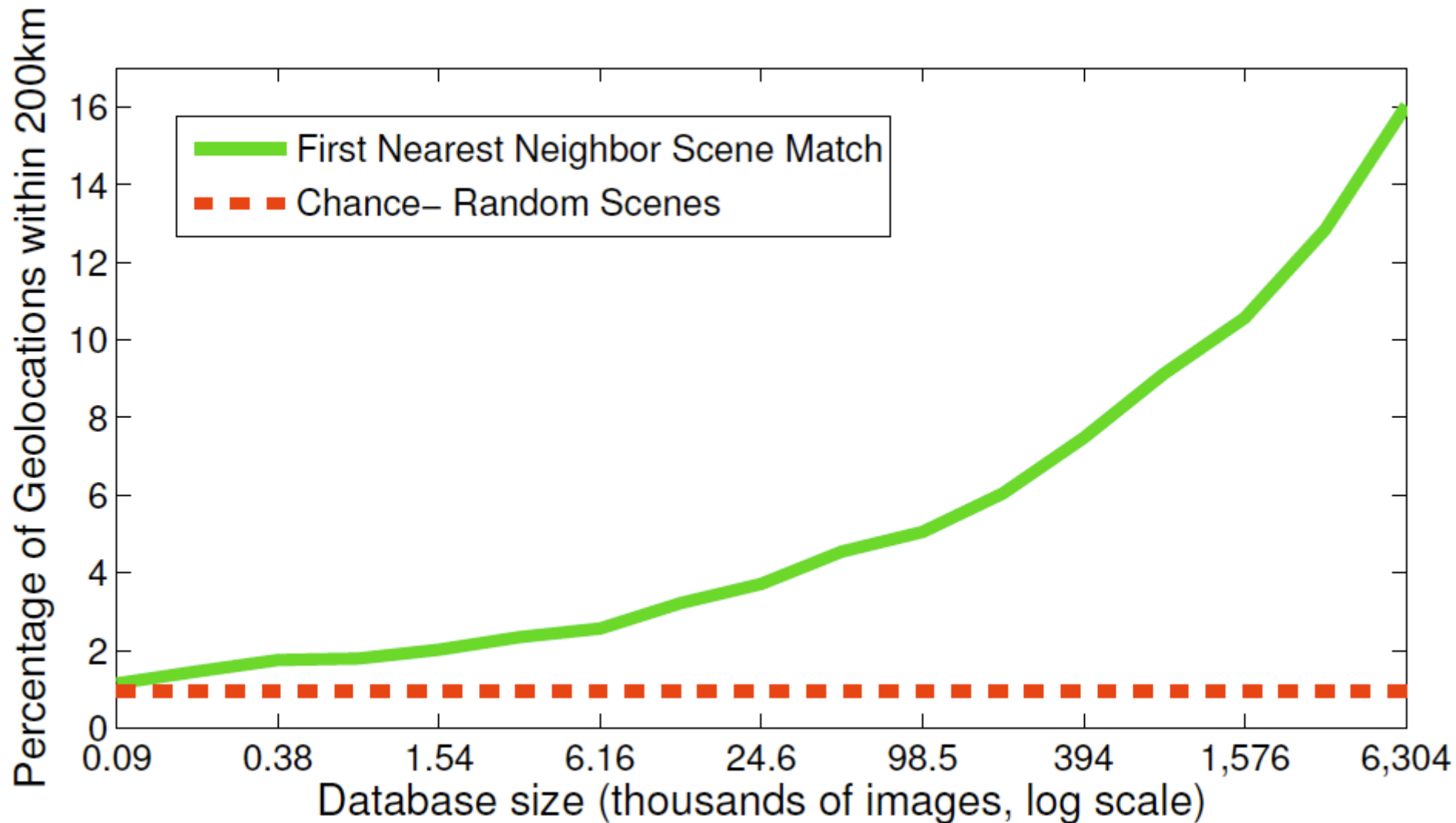
Results



Results



Performance across database size

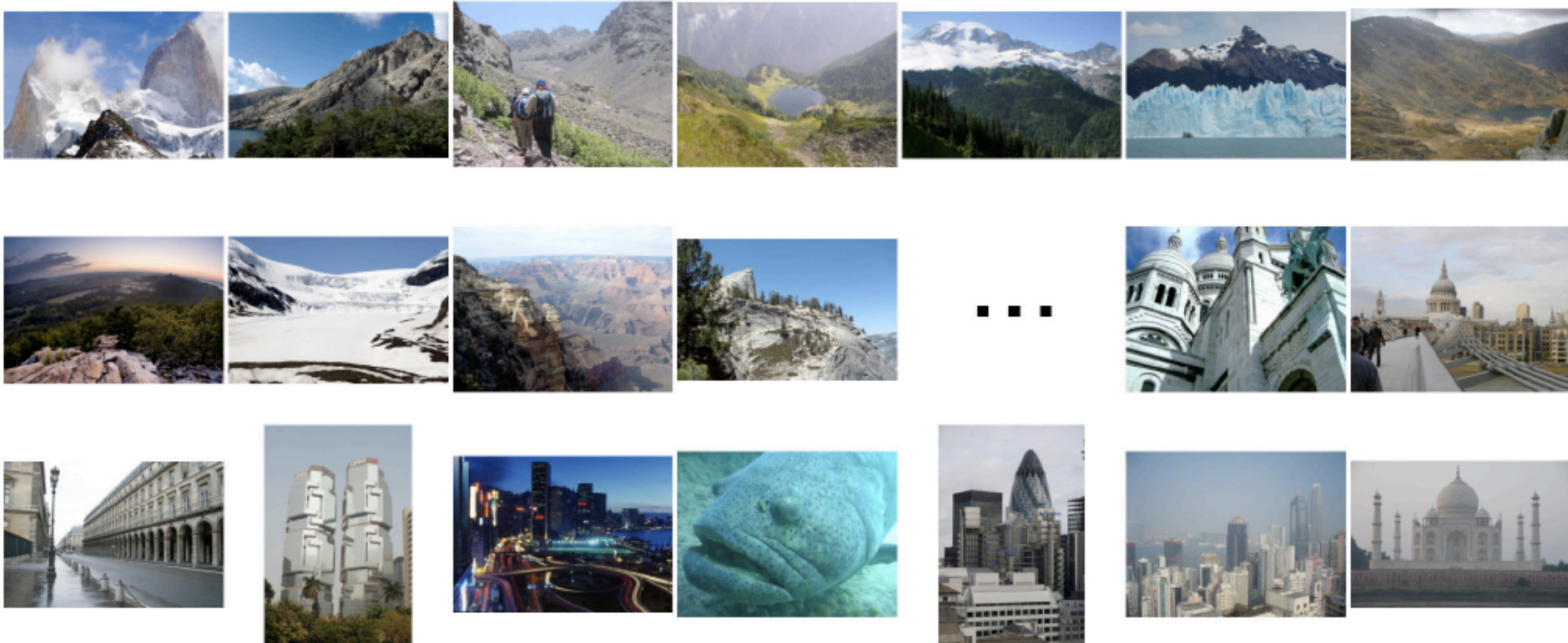


Estimating population density



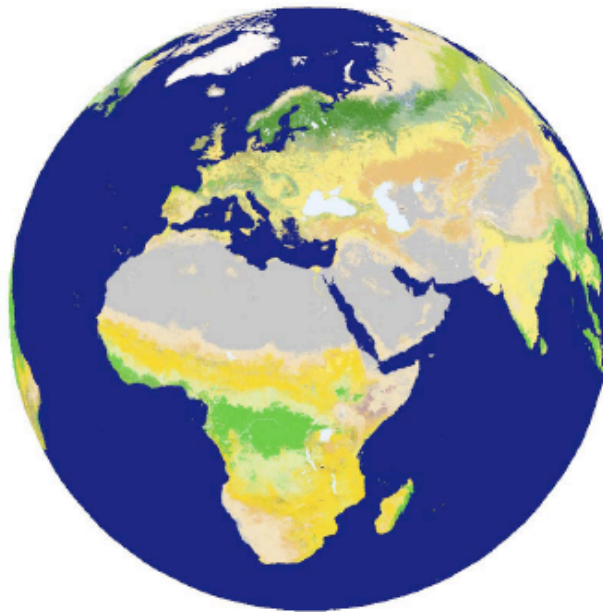
Given a population density map of world, estimate the population of an image by sampling at the estimated location. Images ranked by predicted population density.

Estimating elevation



Given an elevation map of world, they can predict the elevation of an image according to its location. Images ranked by their estimated elevation.

Landcover Classification



Forests



Evergreen Needleleaf Forest



Evergreen Broadleaf Forest



Deciduous Needleleaf Forest



Deciduous Broadleaf Forest



Mixed Forests

Shrublands, Grasslands, and Wetlands



Closed Shrublands



Open Shrublands



Woody Savannas



Savannas



Grasslands



Permanent Wetlands

Agriculture, Urban, and Barren



Croplands



Urban and Built up



Cropland/Natural Vegetation Mosaic



Snow and Ice



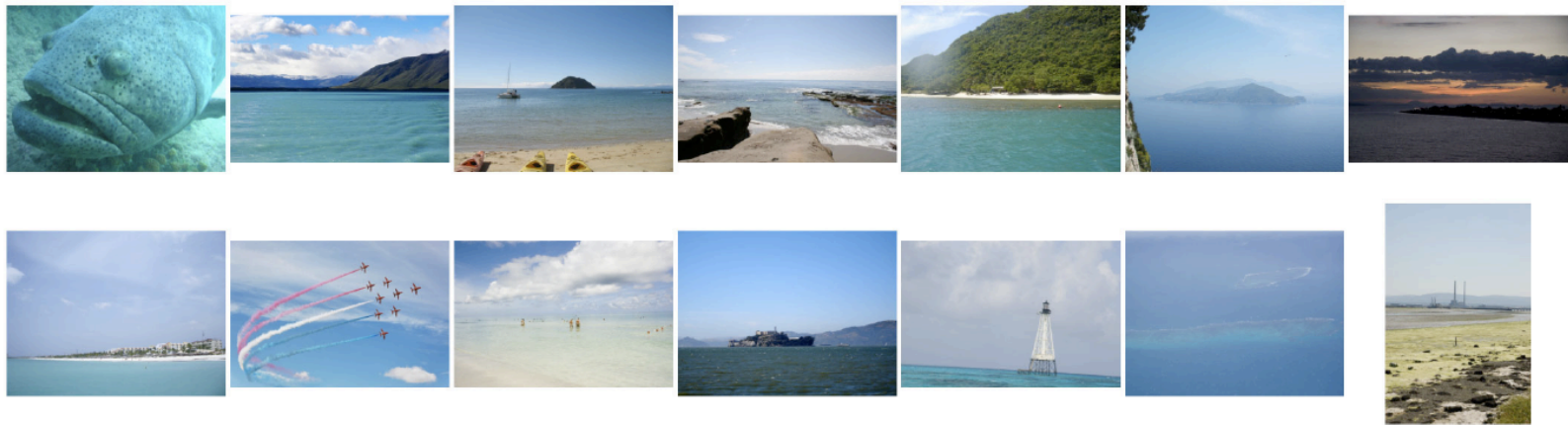
Barren or Sparsely Vegetated

Given landcover map, predict which images are most likely to be examples of each category

Forest



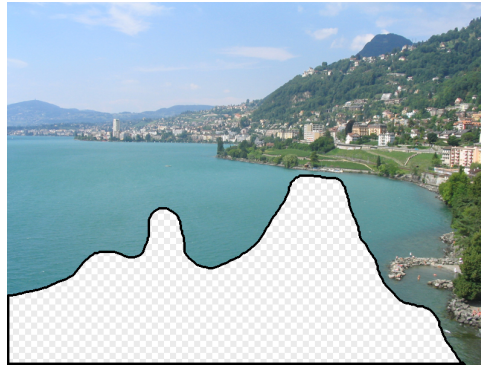
Water



Savanna



NN for Scene Completion Using Millions of Photographs



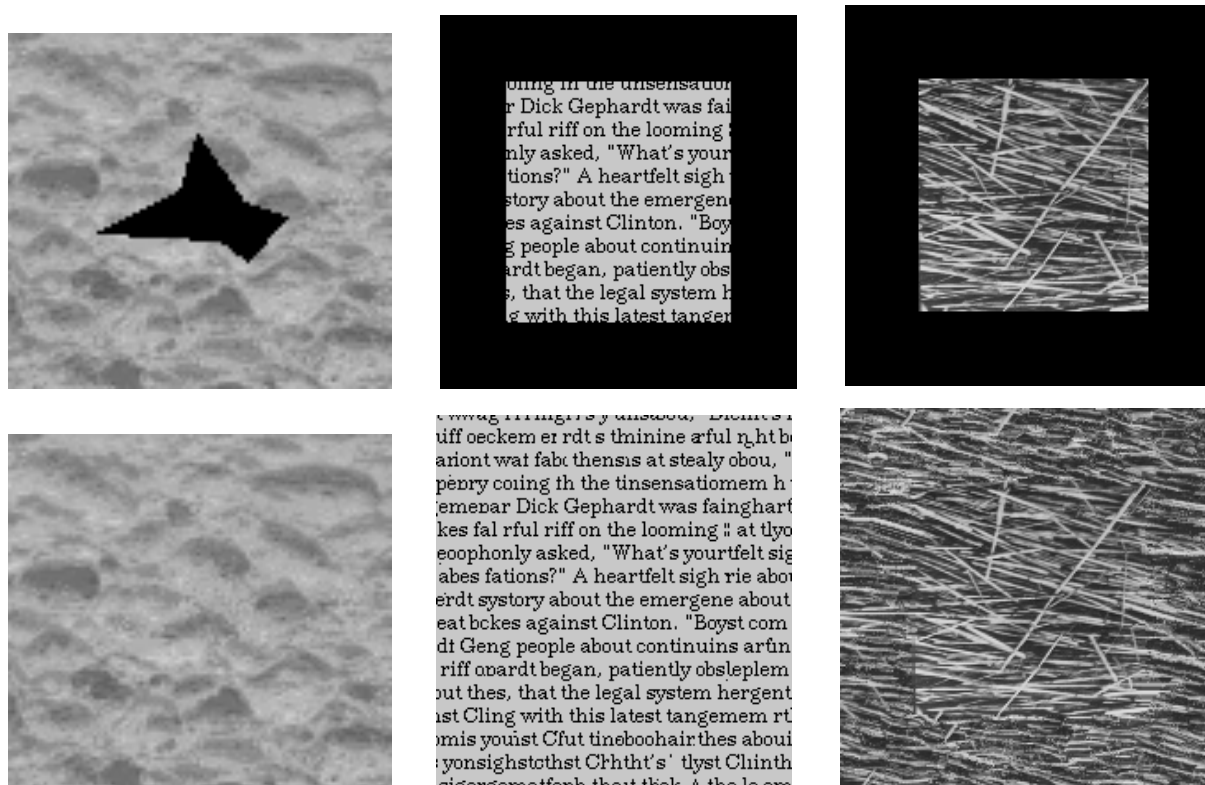
James Hays and Alexei A. Efros
Carnegie Mellon University



Thanks to James & Alyosha for slides!







Efros and Leung. Texture synthesis by non-parametric sampling. ICCV 1999.

Fill in unknown region from source image parts



Efros and Leung result – no notion of semantics, also assumes necessary data is present elsewhere in the image

Scene Matching for Image Completion







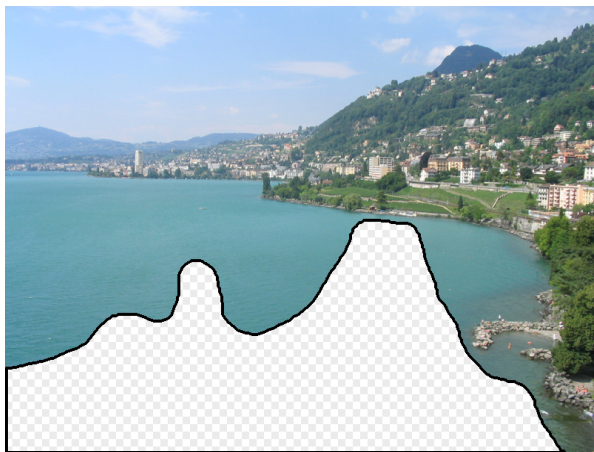
Challenges:

Computational costs of
searching lots of
images

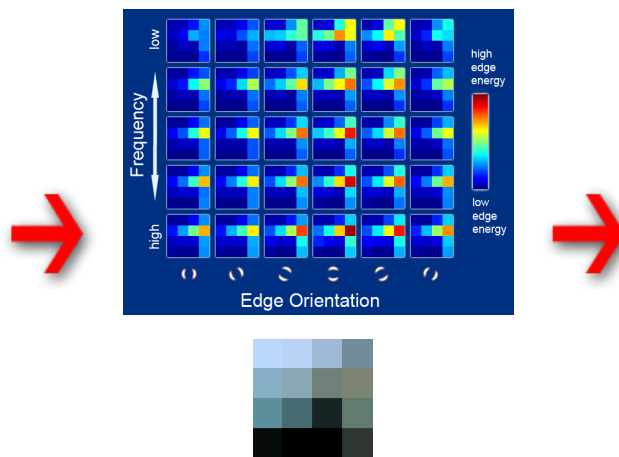
Should fill in missing
regions with
semantically valid
fragments

Scene Completion Result

The Algorithm



Input image



Scene Descriptor



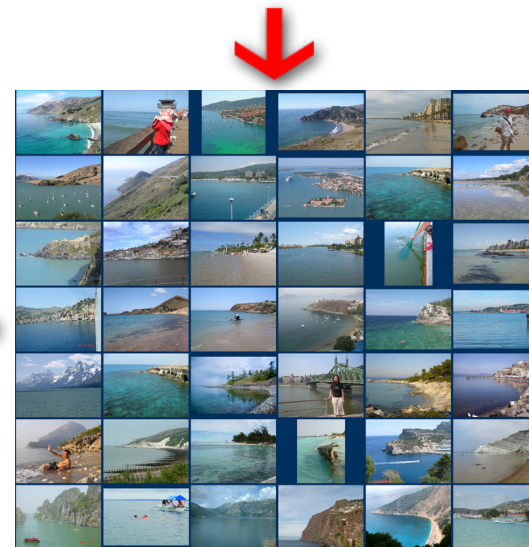
Image Collection



20 completions



Context matching
+ blending



200 matches

Data

They downloaded **2.3 Million** unique images from Flickr groups and keyword searches.

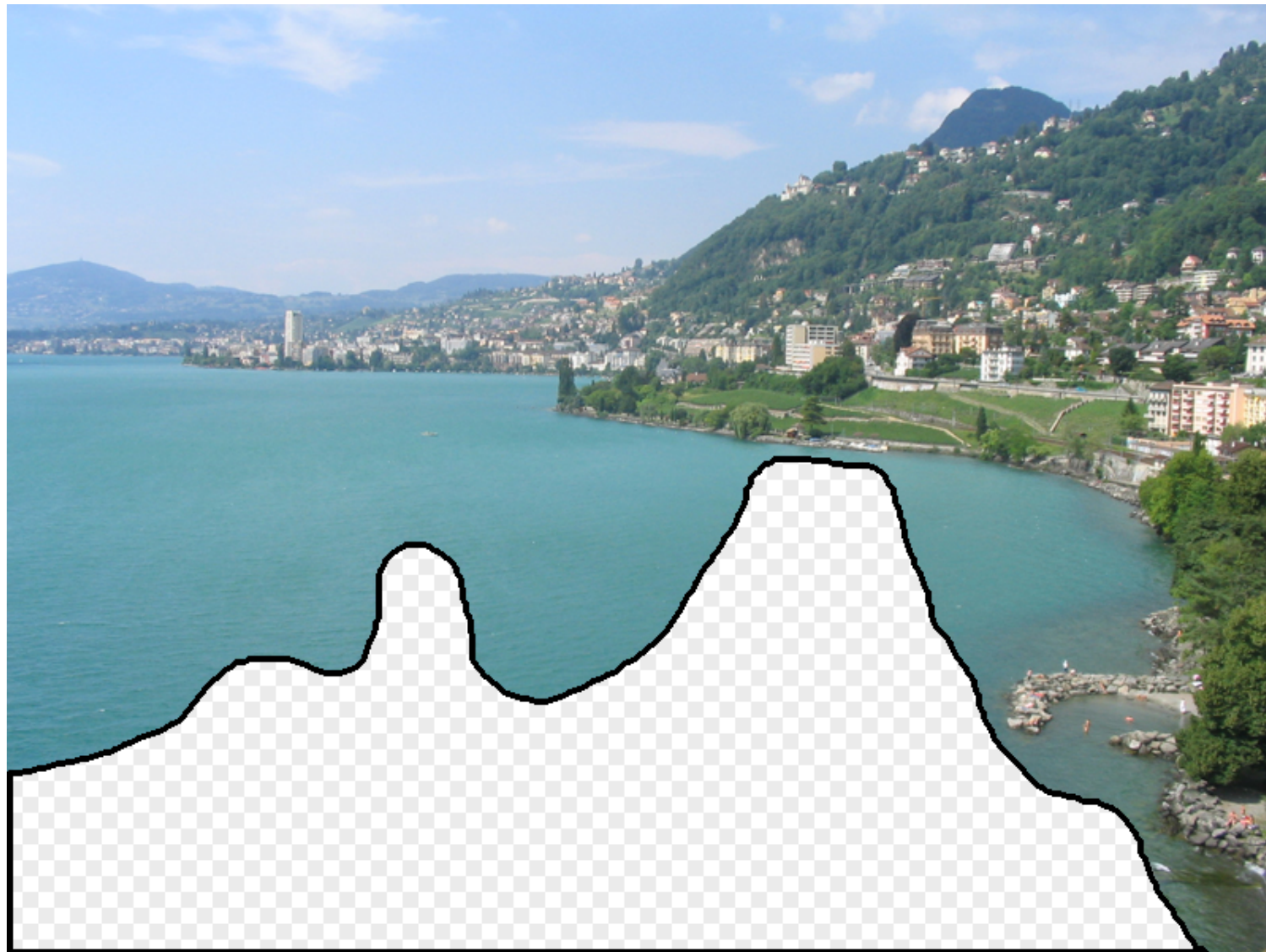


Groups: lonelyplanet, urban-fragments, ruraldecay ...

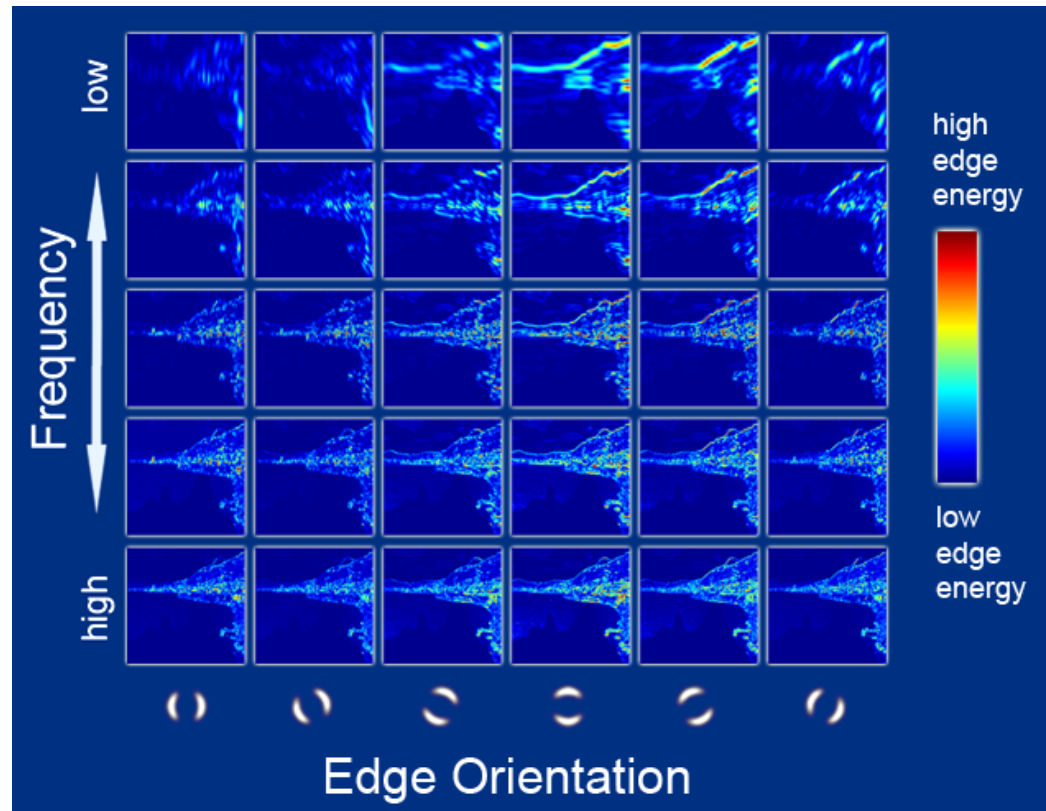
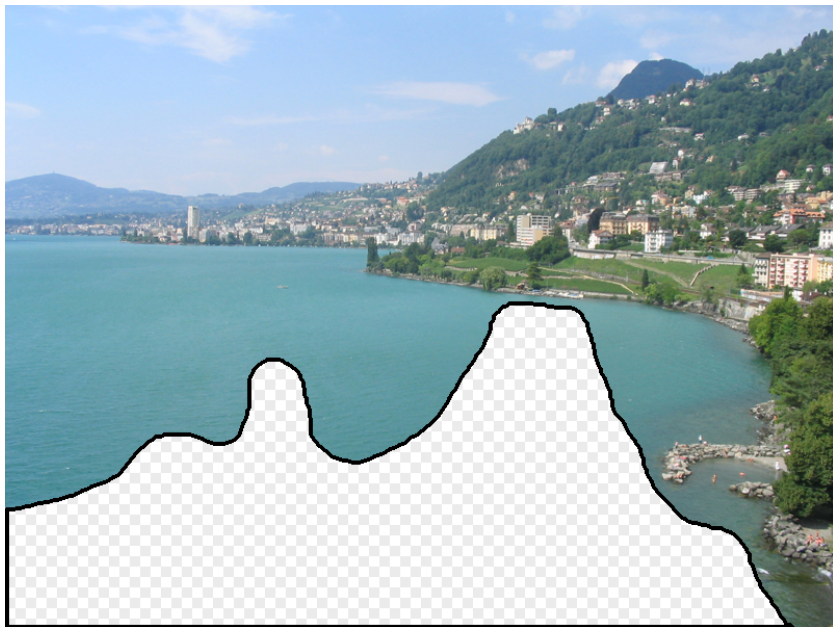
Keywords: outdoors, vacation, river...

Discard duplicates and small images

Scene Matching

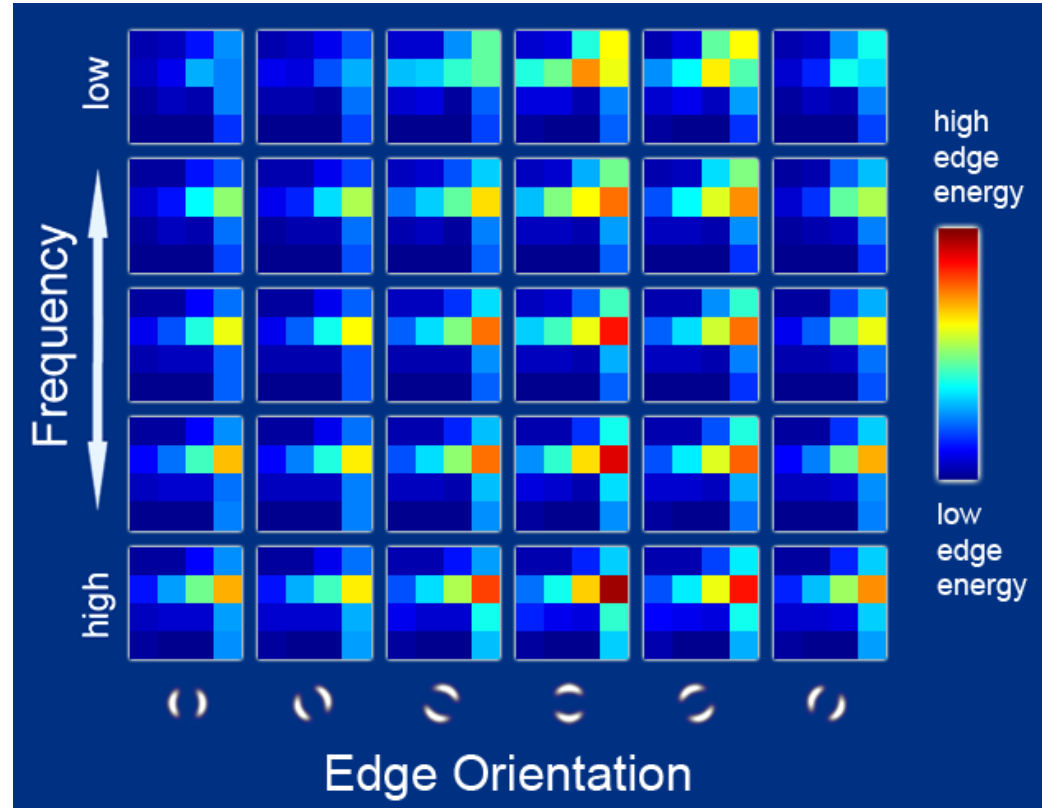
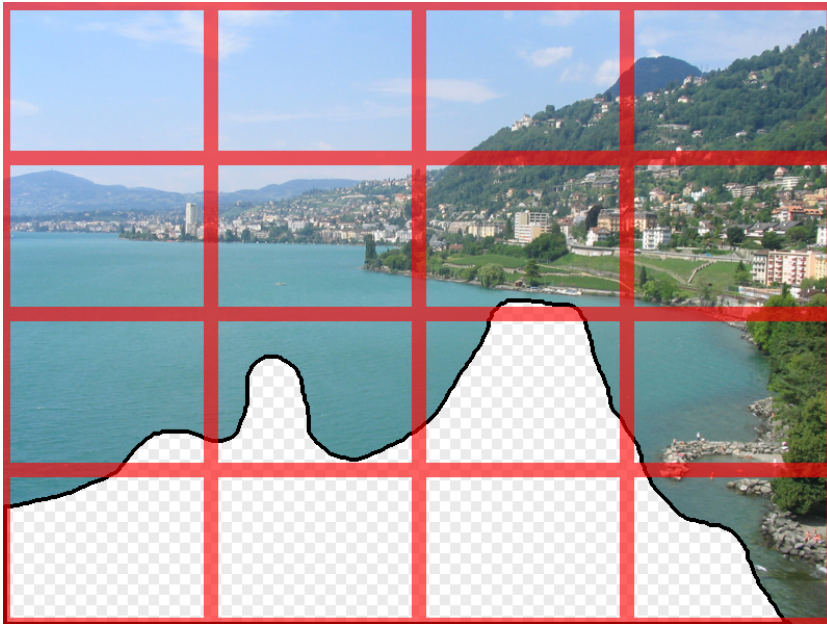


Scene Descriptor



Compute oriented edge response at multiple scales (5 spatial scales, 6 orientations)

Scene Descriptor



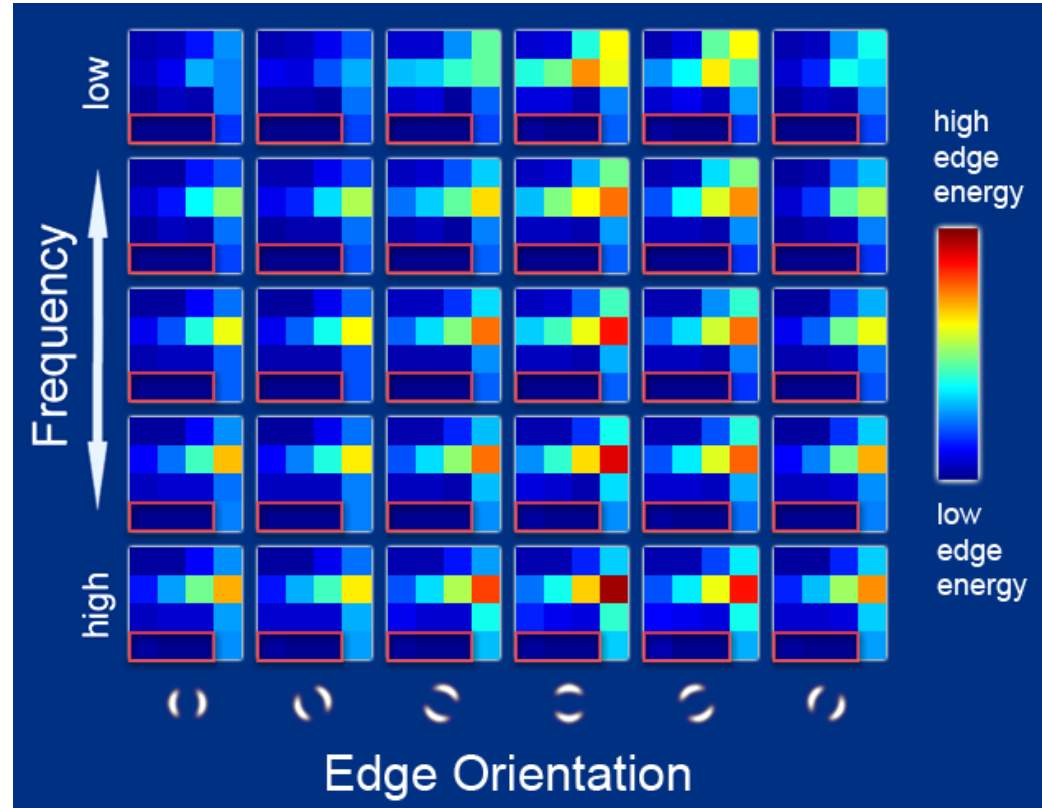
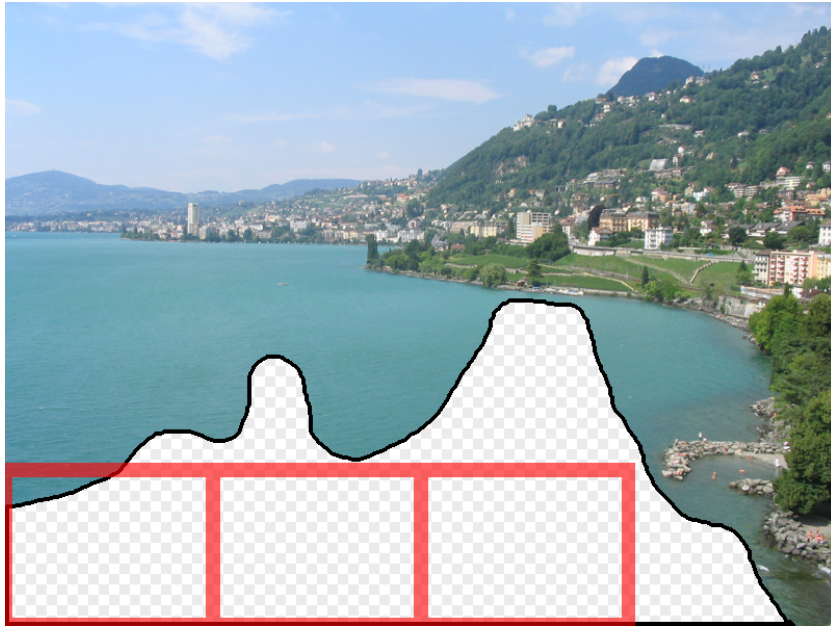
Gist scene descriptor (Oliva and Torralba 2001)

“semantic” descriptor of image composition

aggregated edge responses over 4x4 windows

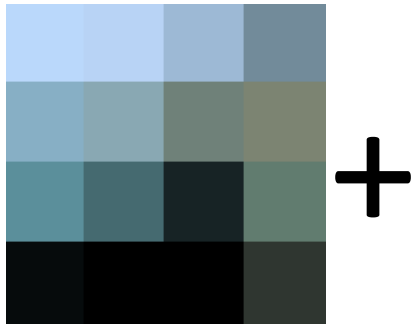
scenes tend to be semantically similar under this descriptor if very close

Scene Descriptor



Gist scene descriptor - with missing regions masked (weighted based on percentage of valid pixels)

Scene Descriptor

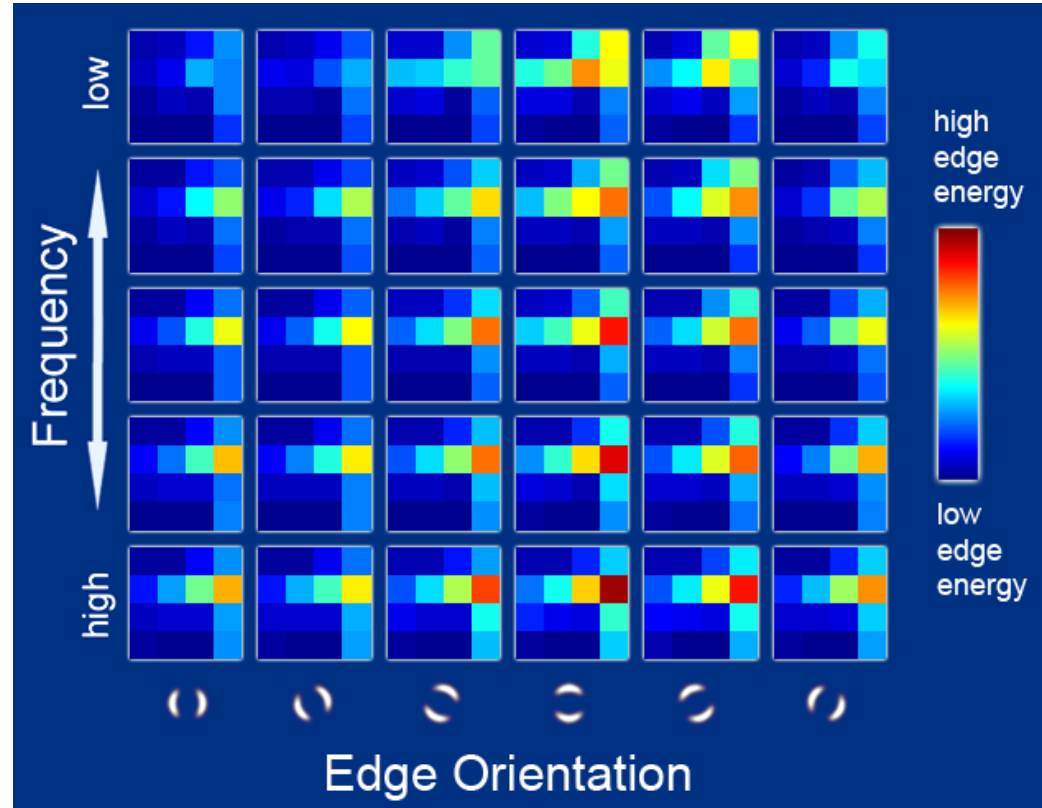


+

Color descriptor – color of the query image downsampled to 4x4

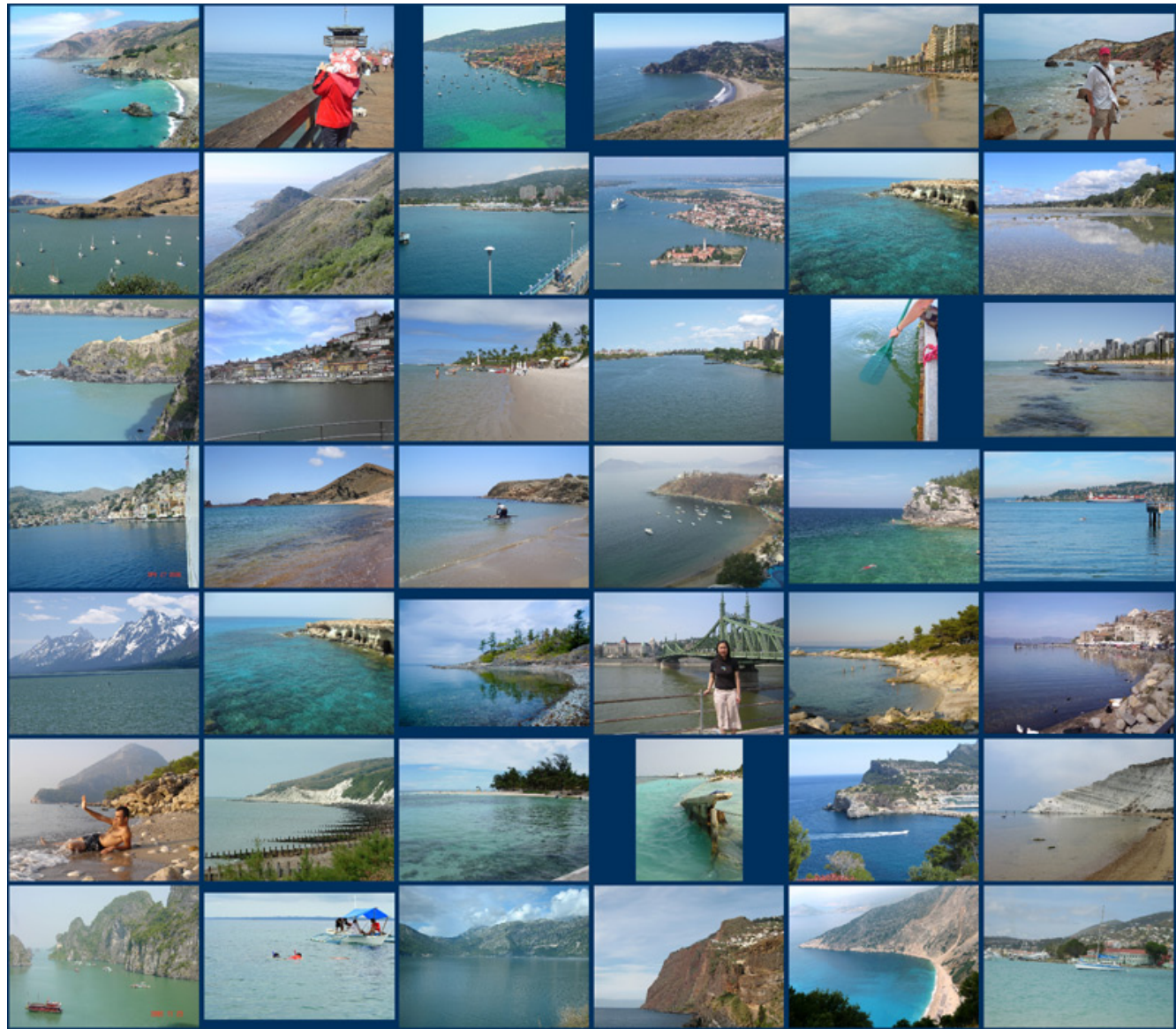
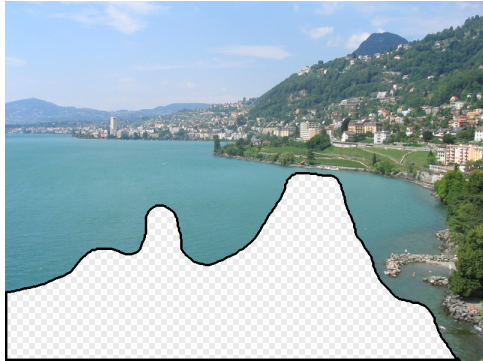
Distances calculated by SSD between query image descriptors & imgs in database

Total Dist = color dist + 2*gist dist



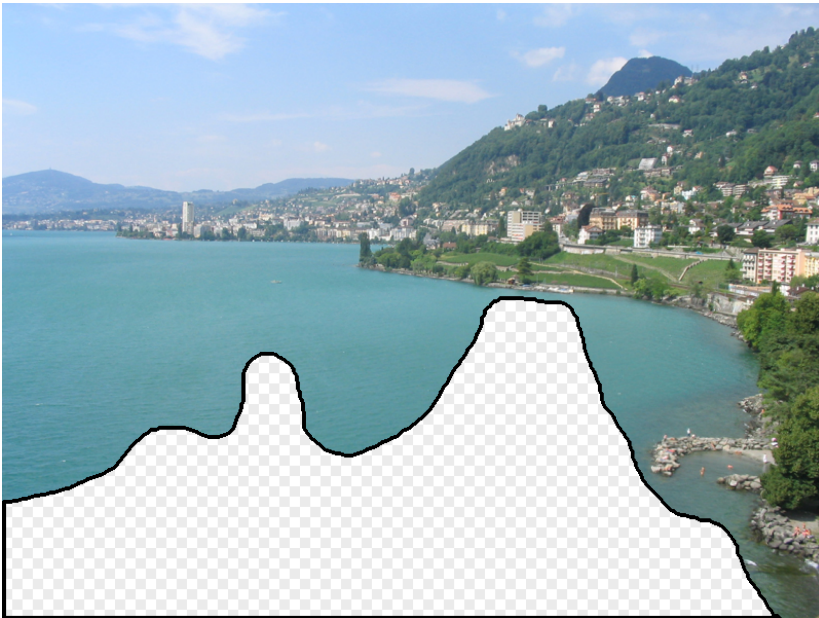
Gist scene descriptor
(Oliva and Torralba 2001)





... 200 total

Context Matching



Need to more precisely align matching scenes to local img context around missing region

local context = all pixels within 80 pixel radius of hole's boundary

Compute pixel-wise error of 200 best scene matches across all valid translations and 3 scales

Compute texture similarity of proposed fill-in to removed region



Graph cut



Final result – blended between the two images along the cut to merge seamlessly







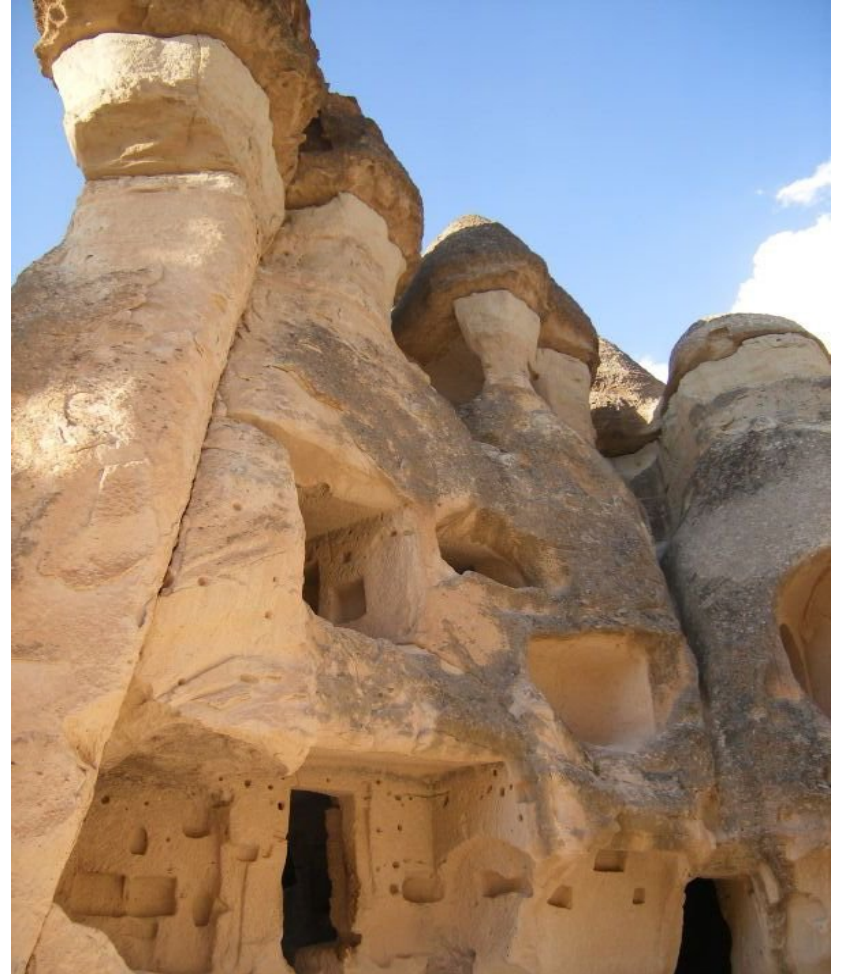
Failures



Failures



Failures



Cause of failure – atypical scene caused lack of good matches

Failures



Related work in Deep Learning



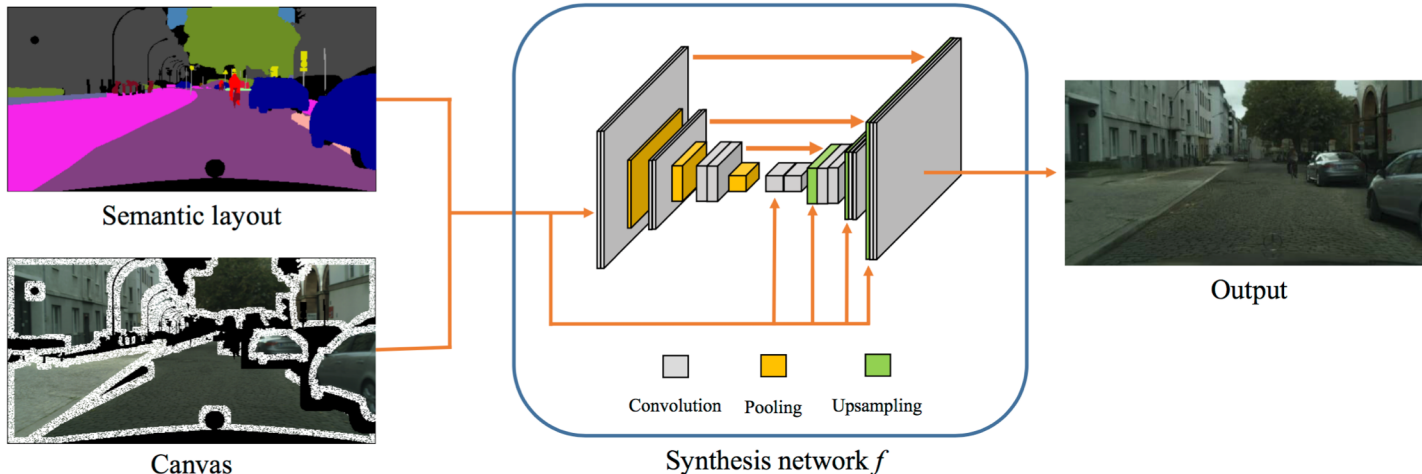
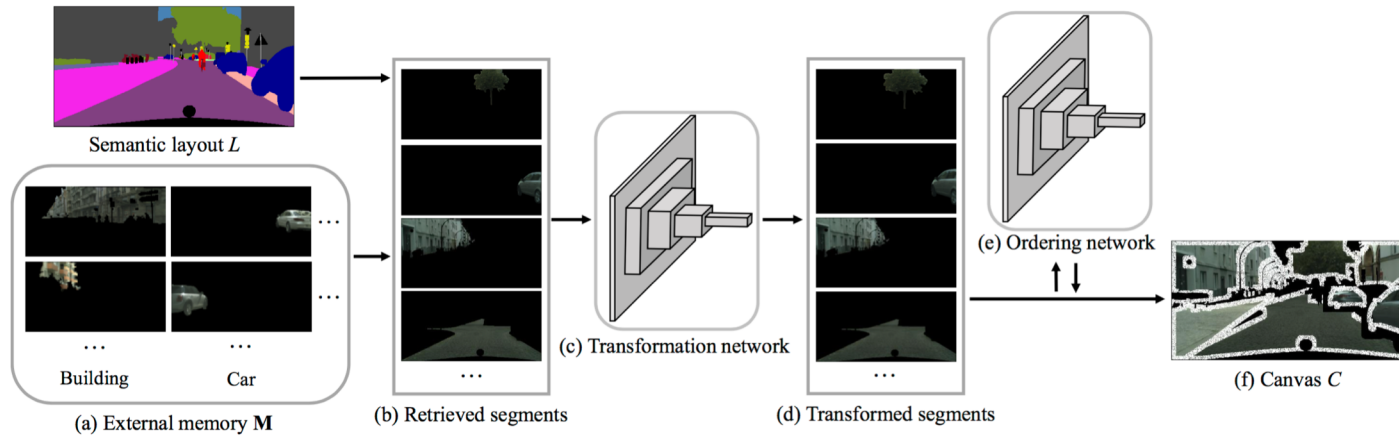
Semantic layout L



Output

Goal: Image synthesis from semantic layout

Related work in Deep Learning



NN for Image Captioning

Traditional Recognition...



→ person



→ shoe



→ car

Human centric recognition outputs



car

Human centric recognition outputs



→ pink car

Human centric recognition outputs



→ car on road

Human centric recognition outputs



Little pink smart car
parked on the side
of a road in a
London shopping
district.

Telling the “*story of an image*”

NN Captioning: Compose descriptions
from recognized content
+ existing descriptions

(Captioning circa 2012-2013)



Through the smoke



Duna Portrait #5



Mirror and gold



the cat lounging in the sink

Data exists, but buried in junk!

Captions in the Wild

<http://tamaraberg.com/sbucaptions>



The Egyptian cat statue by the floor clock and perpetual motion machine in the pantheon



Little girl and her dog in northern Thailand. They both seemed interested in what we were doing



Interior design of modern white and brown living room furniture against white wall with a lamp hanging.



Man sits in a rusted car buried in the sand on Waitarere beach



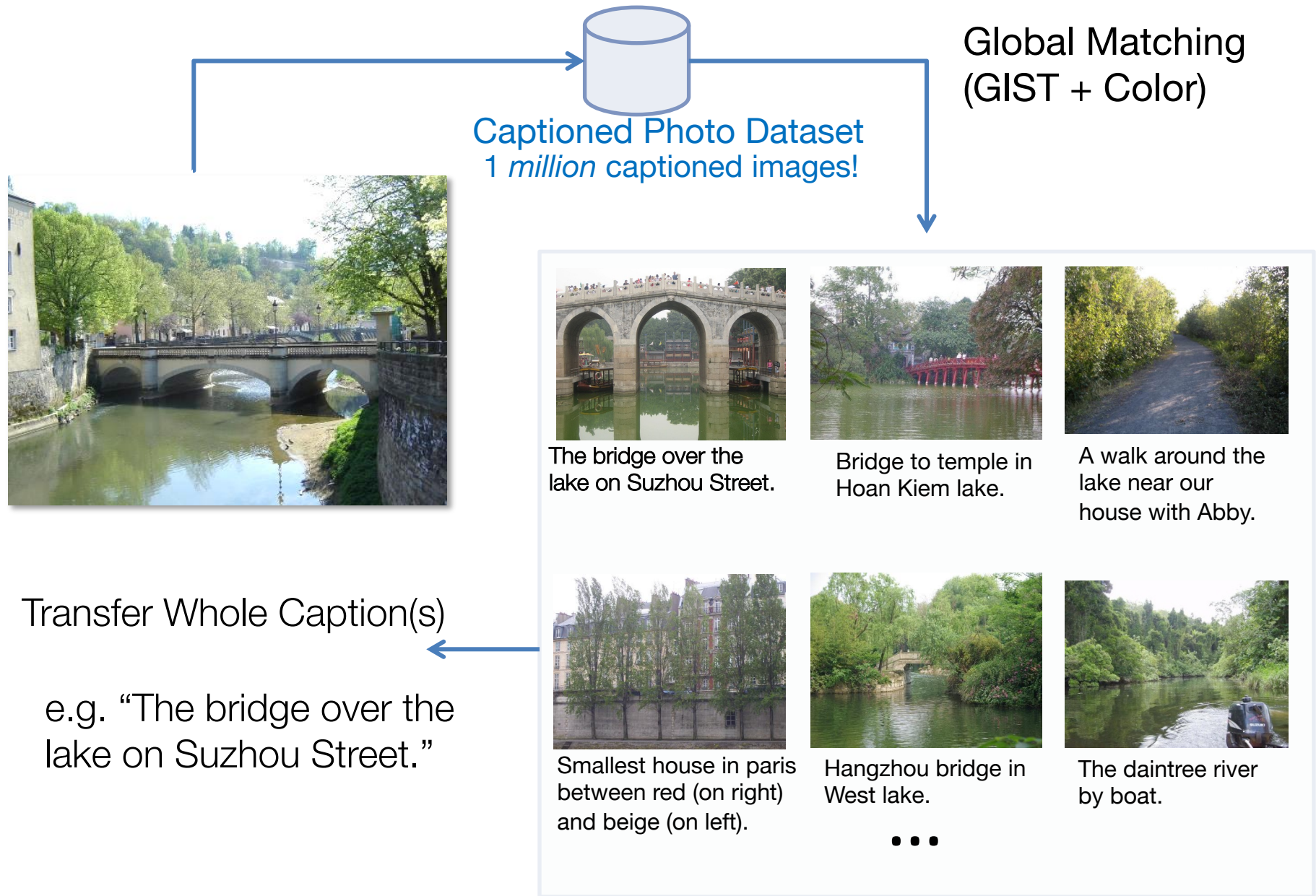
Our dog Zoe in her bed



Emma in her hat looking super cute

Harness the Web

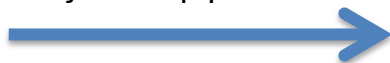
Ordonez et al, NIPS 2011



Transfer pieces of Captions



Object appearance



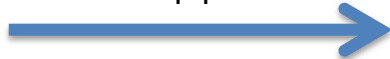
NP: the dirty sheep

Object pose



VP: meandered along a desolate road

Scene appearance



PP: in the highlands of Scotland

Region
appearance &
relationship



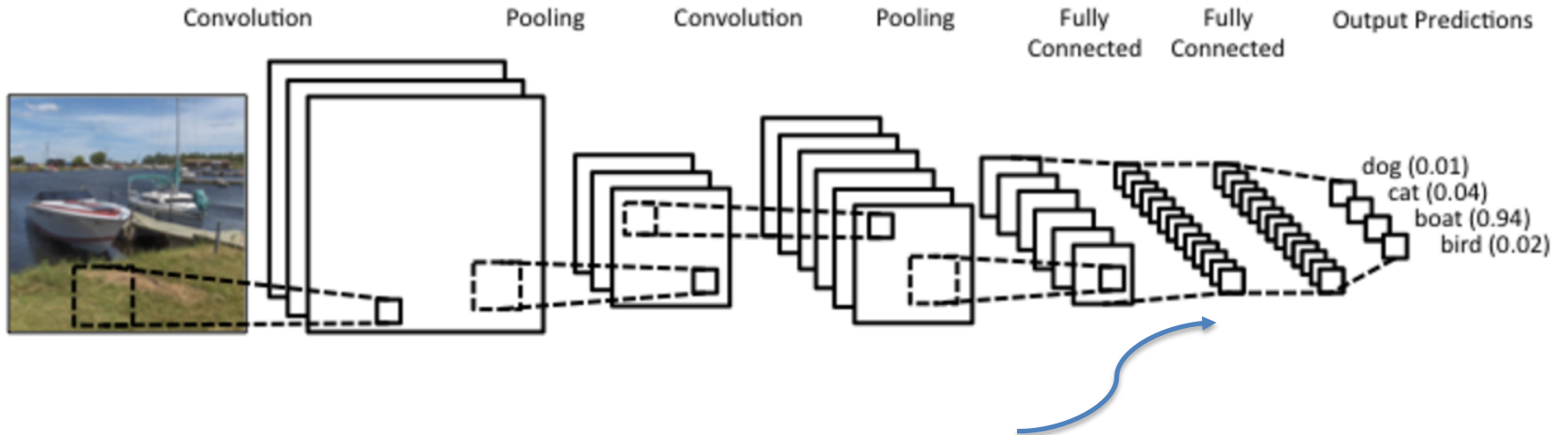
PP: through frozen grass



Example Composed Description:

the dirty sheep meandered
along a desolate road in the
highlands of Scotland
through frozen grass

CNNs (coming soon)



Feature representation extracted from some level of the CNN

With deep learning now features can be learned rather than hand-crafted.

Better features == better captions



A bedroom with a bed and a couch.

A hotel room with two beds and a table.



A train is stopped at a train station.

A red and white train parked in a train station.



A group of people sitting around in a living room.

A group of people sitting on a couch in a living room.



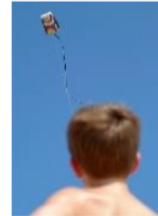
A baseball player holding a bat on a field.

A baseball player holding a bat on a field.



A man riding a wave on a surfboard.

A man riding a wave on a surfboard in the ocean.



A person flying a kite in the sky.

A person flying a kite in the sky.



A cat sitting in a bathroom sink.

A black and white cat sitting in a bathroom sink.



A building with a clock on the top.

A clock tower on the top of a building.

Brainstorming Session

- Form small groups and come up with applications where you could use feature matching for some task