# Retrieval in multimedia

# Image retrieval systems

- Images can be queried using
  - Metadata (text)
  - User annotations
  - Image features (content)
- Problems
  - Metadata is not complete/informative/available
  - User annotations not supported, unreliable

My pet *Tiger*

# Images and text queries

- Images in web documents
  - Use text around image (URL element name, neighborhood)
  - Same principles as in text retrieval systems
- Example of searching for images with word »Sunset«



Sunset at Rocky Point

Frank Smiles at Sunset

Sunset Beach

# Images and text queries

## Query: »tiger in woods«
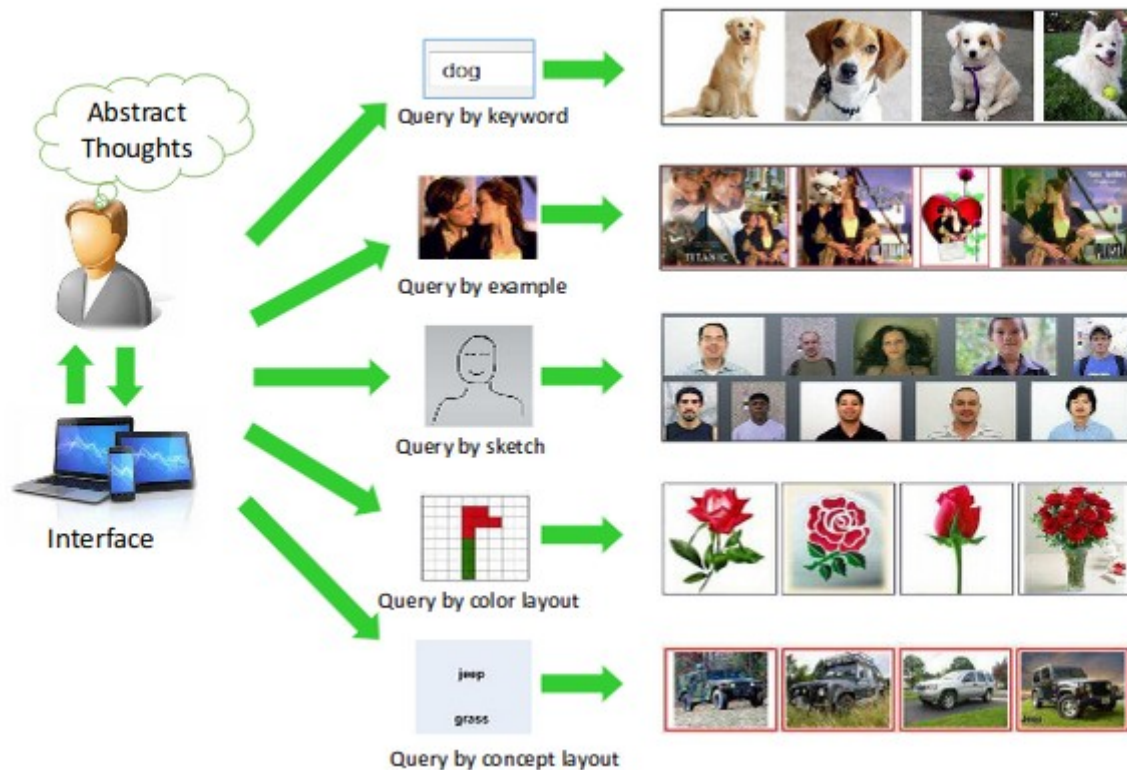
Desired result



Obtained result

# Problems with text queries

- Avoid using image content
  - Annotation bias
  - Metadata ambiguity
- Perceptual relevance
  - Impossible to describe composition
  - Abstract shapes

Development of retrieval systems that

encode image content directly

# Image retrieval systems

# Querying image content

- Extract image content
  - Detecting object and categories
  - Describing relations, actions
  - Ambiguous problem
- Low-level features
  - Color
  - Texture
  - Shape
  - Structural elements

# Image retrieval system

# Querying by color

- Average color – no information about the distribution around average



$$\mu_1 \qquad d = \mu_1 - \mu_2 \qquad \mu_2$$

- Parametric distribution (Gaussian)



$$(\mu_1, \sigma_1) \qquad\qquad (\mu_2, \sigma_2)$$

Bhattacharryya distance: $\quad d = \frac{1}{8}(\mu_1 - \mu_2)^T \boldsymbol{\Sigma}^{-1}(\mu_1 - \mu_2) + \frac{1}{2}\left(\frac{|\boldsymbol{\Sigma}|}{\sqrt{|\boldsymbol{\Sigma_1}||\boldsymbol{\Sigma_2}|}}\right) \qquad \boldsymbol{\Sigma} = \frac{1}{2}(\boldsymbol{\Sigma_1} + \boldsymbol{\Sigma_2})$

# Color histograms

- General non-parametric model
  - Gaussian distribution is single-modal
  - Images are usually multi-modal

# Histogram properties

- Robustness
  - Scale change, rotation
  - Resolution change
  - Partial occlusions
- No spatial information
- Sensitivity to illumination variation
  - Remove the value part

# What is a texture?

- No exact definition

  »Texture is a description of the spatial arrangement of color or intensities in an image or a selected region of an image.«

- Shape and texture

- Level of detail

# Querying using texture

- **Low-level description**
  - Spatial properties
  - Frequency properties
- **Perceptual properties**
  - periodicity, coarseness, dominant orientation, complexity

repeatability

stochasticity

combination

fractals

# Coocurrence matrix

- How many times does pixel of value V1 appear next to pixel of value V2?
  - Displacement vector d=[dy,dx]
  - C(i,j) contains number of times values i an j appear on image in relation d
  - Cooccurence matrix is normalized

d=[dy,dx]=[100,-20]

dx=-20

i: 26

dy=100

j: 50

j=50

i=26

C

v(**C**) = [f1,f2,f3,f4,f5,...]

# Extracting features

Various features can be computed from cooccurence matrix

$$Energy = \sum_{i,j} C_A(i,j)^2$$

$$Entropy = -\sum_{i,j} C_A(i,j) log_2 C(i,j)$$

$$Contrast = \sum_{i,j} C_A(i,j)(i-j)^2$$

$$Homogenity = \sum_{i,j} \frac{C_A(i,j)}{1+|i-j|}$$

$$Correlation = \frac{\sum_{i,j}(i-\mu_i)(j-\mu_j)C_A(i,j)}{\sigma_i \sigma_j}$$

**C**

v(**C**) = [f1,f2,f3,f4,f5,...]

Comparison: Euclidean distance

# Local Binary Pattern

- Describe global texture with local descriptors

- For each pixel p compute 8-bit number

- Texture represented as histogram of these local numbers



T. Ahonen, A. Hadid, M. Pietikinen, "Face Description with Local Binary Patterns: Application to Face Recognition", TPAMI2006

# Auto-correlation

- Normalized scalar product between image and its shifted version

- Shape of response function describes

  - Texture regularity

  - Texture coarseness

$$\rho(x,y) = \frac{\sum_{u,v} I(u,v)I(u+x,v+y)}{\sum_{u,v} I(u,v)^2}$$

# Fourier transform

- Description of image with complex basis functions
  - Energy of spectrum: |F(u,v)|
  - If I is WxH, then F is WxH

$$F(u,v) = \mathcal{F}\{I(x,y)\}(u,v) = \frac{1}{WH} \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I(x,y) e^{-i2\pi(\frac{ux}{W} + \frac{vy}{H})}$$

$I(x,y)$

(0,0)

$F(u,v) = 2dFFT[I(x,y)]$

$\log(abs(F(u,v))+1)$

(W,H)

$F(u,v)$

Low frequencies

(0,0)

$(u_{max}, v_{max})$

# Spectrum features

How much energy is contained in various parts of spectrum



2D FFT

**v** = [f1, f2,f3,f4]

# Query by shape

- Edge detection, threshold
- Vector of features
  - Region moments
  - Freeman differential codes
- Transformation distance
  - Amount of transformation

# Comparing histograms

- Euclidean distance

$$D = \sqrt{\sum (h_1(i) - h_2(i))^2}$$

- Hellinger distance

$$H = \left(\frac{1}{2} \sum_{i=1:N_{bins}} (h_1(i)^{\frac{1}{2}} - h_2(i)^{\frac{1}{2}})^2\right)^{\frac{1}{2}}$$

- Chi-square distance

$$\chi^2 = \frac{1}{2} \sum_{i=1:N_{bins}} \frac{(h_1(i) - h_2(i))^2}{h1(i) + h2(i) + \varepsilon_0}$$

- Histogram intersection

$$I = 1 - \sum_{i=1:N_{bins}} \min(h_1(i), h_2(i))$$

# Including spatial information

- Divide image into subregions
- Stack histograms

# Bag of words

- Inspired by text retrieval systems
- General object categories
  - No clear spatial consistency
  - Objects composed of important parts - words
- Ignoring relationships between parts
  - Dictionary – list of known parts
  - Descriptor – histogram of part occurrences

Object

Bag of words

# Visual words

| | |
|---|---|
| Word | Feature |
| Token | Centroid/Cluster |
| Document | Image/Frame |
| Corpus | Video/Collection |

# Local regions

- Detecting stable regions
  - Robustness
  - Corners, blobs
- Describing neighborhood
  - Invariance (illumination, rotation, scale)



rotate →

scale →

# SIFT features

- Scale invariant feature transform

  - Divide region into 4x4 sub-regions: 16 cells

  - Compute gradients in each sub-region

  - Discretize orientation (8 directions)

  - Compute orientation histogram based on magnitude

  - Stack histograms and normalize: 4x4x8 = 128

# Building a dictionary

- Unsupervised learning
  - Large number of different local descriptors
  - Finite amount of words
  - Clustering



Fei-Fei Li; Perona, P. "A Bayesian Hierarchical Model for Learning Natural Scene Categories".IEEE CVPR 2005

# Example of visual words



Sivic, Josef, and Andrew Zisserman. "Video Google: A text retrieval approach to object matching in videos." IEEE CVPR, 2003

# Hierarchy of parts

- Learn complex shape features
  - Gabor features – edges
  - Cooccurence

- Hierarchical composition

- Histogram of parts



categorical layer $\Omega^{\mathcal{C}}$ car cow swan motorbike bottle giraffe

object layer $\Omega^{\mathcal{O}}$ Layer 6 object layer

Layer 5

Layer 4

Layer 3

$\Omega^2$ Layer 2

fixed layer $\Omega^1$ Layer 1 fixed layer

**Hierarchically organized vocabulary**



$\{\mathcal{V}, \mathcal{E}\}$

**Example of a parse tree at detection**

# Towards high-level categories

- Objects in images
- Scanning image
  - Sliding window
  - Region proposals
- Categorization
  - Features + SVM
  - CNN



sliding window       extract features    classify

region proposals

extract features    classify

# Deep learning

# CNN example – VGG16



$I$   $l_1$   $p_1$   $l_2$   $p_2$   $l_3$   $p_3$   $l_4$   $p_4$   $l_5$   $p_5$   $f_1$ $f_2$ $f_3$   $m_1$

| | convolutional layer | | max-pooling layer | | fully-connected layer | | soft-max layer |
|---|---|---|---|---|---|---|---|

Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. arXiv 2014

# Image retrieval with inverted index

- Multi-object detector (semantic tokens)
- Use Boolean queries to per-process database



A. Popescu, A, Ginsca, H. Le Borgne, "Scale-free content based image retrieval (or nearly so)", ICCV 2017 Workshops

# Efficient retrieval of dense descriptors

- Most descriptors are dense
  - Inverted index not efficient
  - Comparison is slow
- Structure the space
  - Hierarhical clustering
  - Traverse a tree (log n)
- Locality-sensitive hashing
  - Similar descriptors have the same hash value

# Towards image understanding

- Semantic segmentation

- Spatial relationships

- Describing scene



"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

car under elephant

person in cart

person ride dog

person on top of traffic light

cs.stanford.edu/people/karpathy/deepimagesent/

www.di.ens.fr/willow/research/unrel/

# Why decompose images?

- Retrieval with specific queries (e.g. horses)

- Describe entire image
  - Which descriptors belong to object?

- Describe only parts of images
  - How many, what shape?

# Superpixels

- Over-segmentation

- Describe each pixel in CIE Lab and (x,y): $C_i = [l_i \ a_i \ b_i \ x_i \ y_i]^T$

- Manually set number of clusters (superpixes)

- Modified K-means (fast, spatial restrictions)



R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, S. Süsstrunk, SLIC Superpixels Compared to State-of-the-art Superpixel Methods, PAMI2012

# Automatic decomposition examples

# Texton segmentation



- Texton descriptor learning
  - Each pixel described with responses to a bank of filters (e.g. 24 filters)
  - Find textons by clustering responses of filters

- Assign each pixel a texton

- Describe texture around pixel as a histogram of textons

- Segmentation - cluster histograms

filter bank (24)

training images
different light/directions

concatenate responses

response features

clustering into textons

T. Leung and J. Malik. Representing and recognizing the visual appearance of materials using three-dimensional textons. IJCV, 2001

# Segmentation using texture



Original image     k-means (k=5), feature: rgb     k-means (k=5), feature: texton

Multiclass segmentation using textons

# Semantic segmentation

- Segments have semantic meaning
- "Bag-of-textons"
  - Texton features
  - Classifier



- Convolutional neural networks
  - Train network for per-pixel classification
  - Encoding context

# Handling subsampling in CNNs

- Pooling/subsampling
  - Reduce parameter count
  - Increase spatial robustness
- Approaches
  - Interpolation
  - MRF
  - Deconvolution



activations for "cat"

$W \times H \times C$

$I \quad l_1 \quad p_1 \quad l_2 \quad p_2 \quad l_3 \quad p_3 \quad l_4 \quad p_4 \quad l_5 \quad p_5 \quad l_6 \quad l_7 \quad l_8 \ m_1$

convolutional layer — max-pooling layer — soft-max layer

# Avoiding pooling

- Dilated convolution
- MRF

# Encoder-decoder

- Deconvolution produces coarse segments

- Skip connections
  - Information from hi-res features

V. Badrinarayanan, A. Kendall, R. Cipolla, SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation, TPAMI 2017

# Describing video content

- Structure: frame, shot, scene

- Content
  - Dynamics: still, moving objects, camera movement
  - Activity in a frame interval, e.g. jumping, robbery, horse race
  - Categories, e.g. cats, horses, cars
  - Object instances: e.g. Harry Potter, Jack Sparrow, Han Solo

# MPEG-7

- Efficient access and manipulation of multimedia content

- Complementary to MPEG-4

- Standardized text-less object retrieval
  - D – Object descriptors (audio and video)
  - DS – Description schemes
  - DDL – Description definition language (XML)
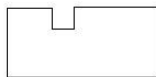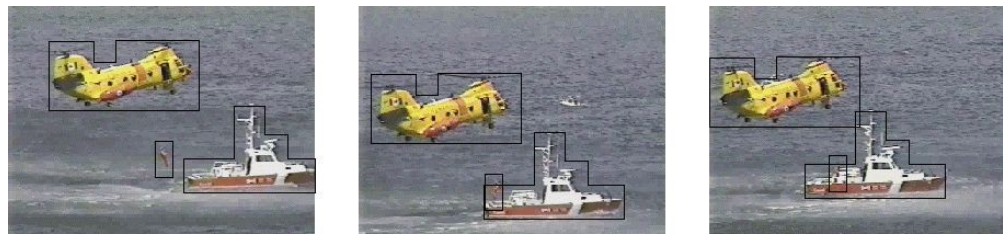
# Examples of descriptors

- Color
  - Color space
  - Color layout
  - Dominant color
  - Color structure
  - GoP color
- Texture
  - Homogenous
  - Non-homogenous

- Shape
  - Shape descriptor
  - Contour
  - 2D-3D shape
- Motion
  - Activity
  - Camera motion
  - Warping parameters
  - Trajectory
  - Parametric motion
- Localization
  - Spatio-temporal
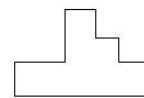  - Region

# Structure description

Describing content at the level of video segment



Moving Region: Helicopter    Moving Region: Person    Moving Region: Boat

Example: three moving objects, describe relations ...

# Applications

- Digital library (Image/video/music catalogue)
- Broadcast media (Radio channel, TV channel)
- Multimedia authoring
- E-business: Searching for products
- Cultural services (art-galleries, museums)
- Educational applications
- Biomedical applications