



Recognition with Bag-of-Words

(Borrowing heavily from Tutorial Slides by Li Fei-fei)

Recognition

- So far, we've worked on recognizing edges
- Now, we'll work on recognizing objects
- We will use a bag-of-words approach

Object



Bag of 'words'



Analogy to documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based on the messages that our eyes.

For a long time, the visual image was thought of as a movie. The image is discovered by the eye, the nerve, image Hubel, Wiesel.

Hubel and Wiesel demonstrate that the *message about the image falling on the retina undergoes a step-wise analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.*

**sensory, brain,
visual, perception,
retinal, cerebral cortex,
eye, cell, optical
nerve, image
Hubel, Wiesel**

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports to \$750bn, compared with \$575bn in 2004.

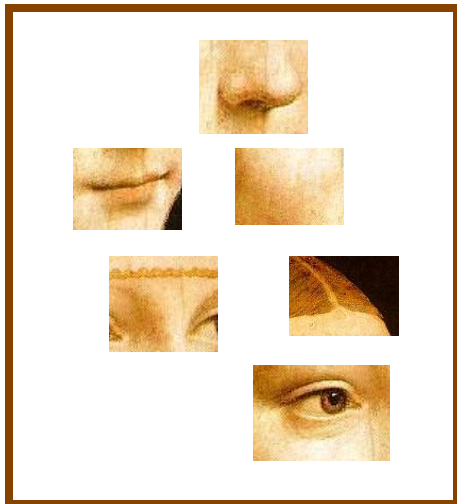
The US government is annoyed by China's deliberate policy of keeping the yuan undervalued. The US government also needs to increase demand for its own country. China's trade surplus is a problem for the US.

China's trade surplus is a problem for the US. The US government also needs to increase demand for its own country. China's trade surplus is a problem for the US.

**China, trade,
surplus, commerce,
exports, imports, US,
yuan, bank, domestic,
foreign, increase,
trade, value**

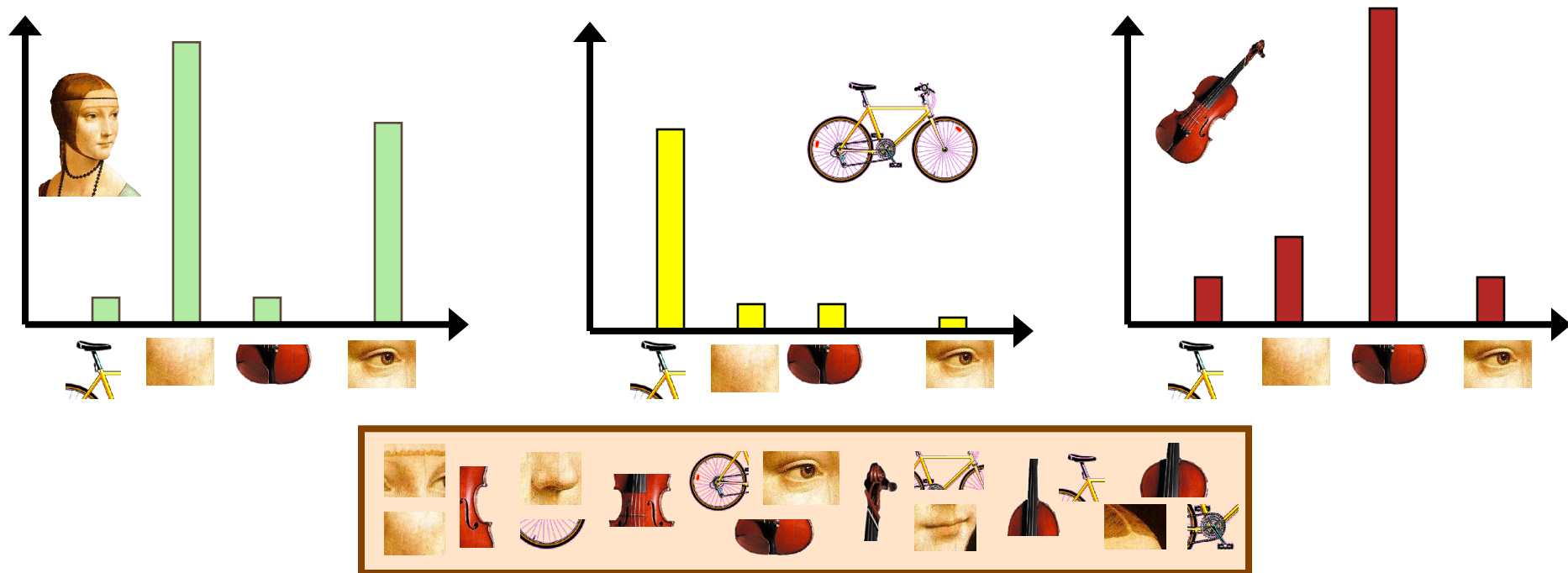
A clarification: definition of “BoW”

- Looser definition
 - Independent features

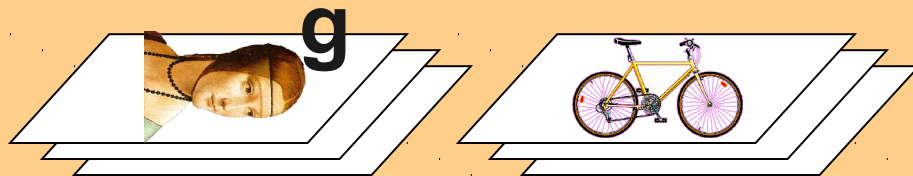


A clarification: definition of “BoW”

- Looser definition
 - Independent features
- Stricter definition
 - Independent features
 - histogram representation



learning



feature detection
& representation

image representation

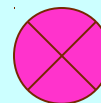
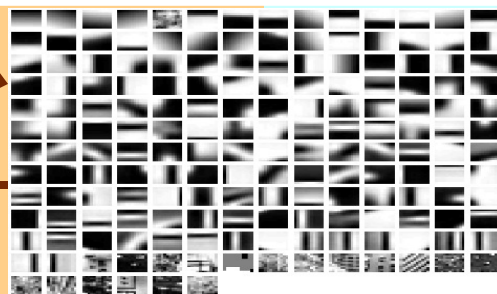


**category models
(and/or) classifiers**

recognition

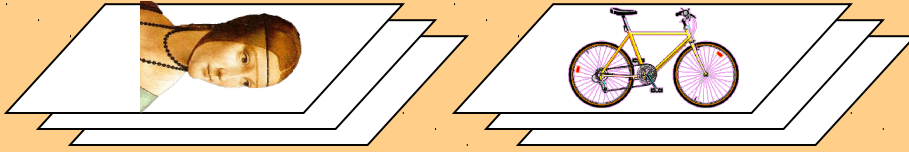


codewords dictionary



**category
decision**

Representation



2.

codewords dictionary

1. feature detection
& representation

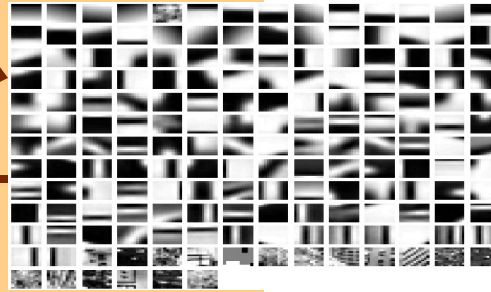


image representation

3.



1.Feature detection and representation

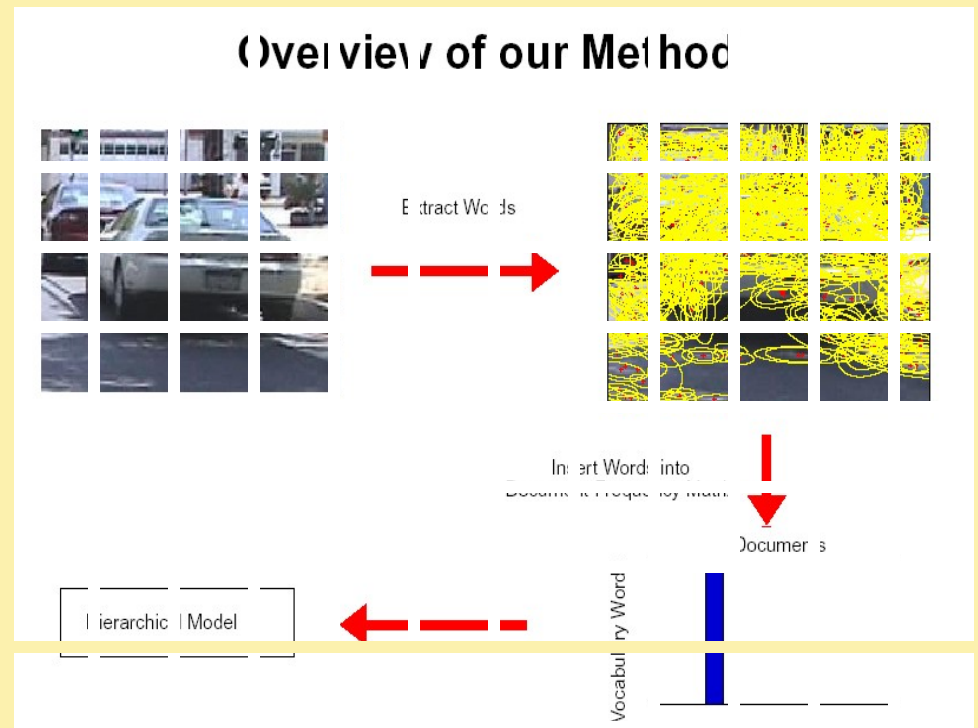


1. Feature detection and representation

Regular grid

Vogel & Schiele, 2003

Fei-Fei & Perona, 2005



1.Feature detection and representation

Regular grid

Vogel & Schiele, 2003

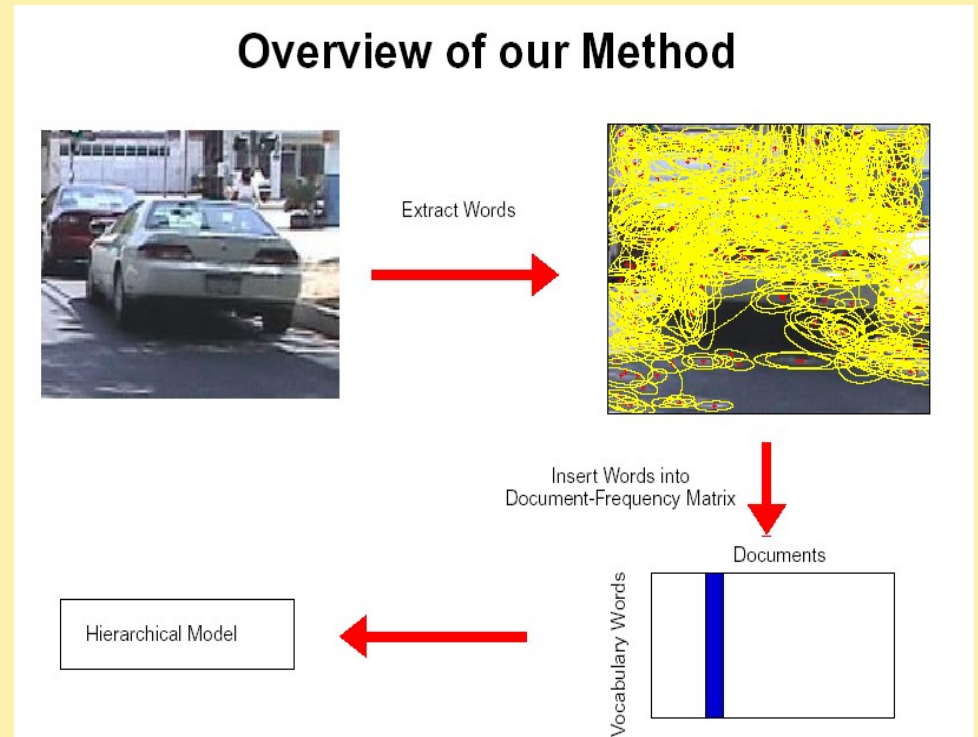
Fei-Fei & Perona, 2005

Interest point detector

Csurka, et al. 2004

Fei-Fei & Perona, 2005

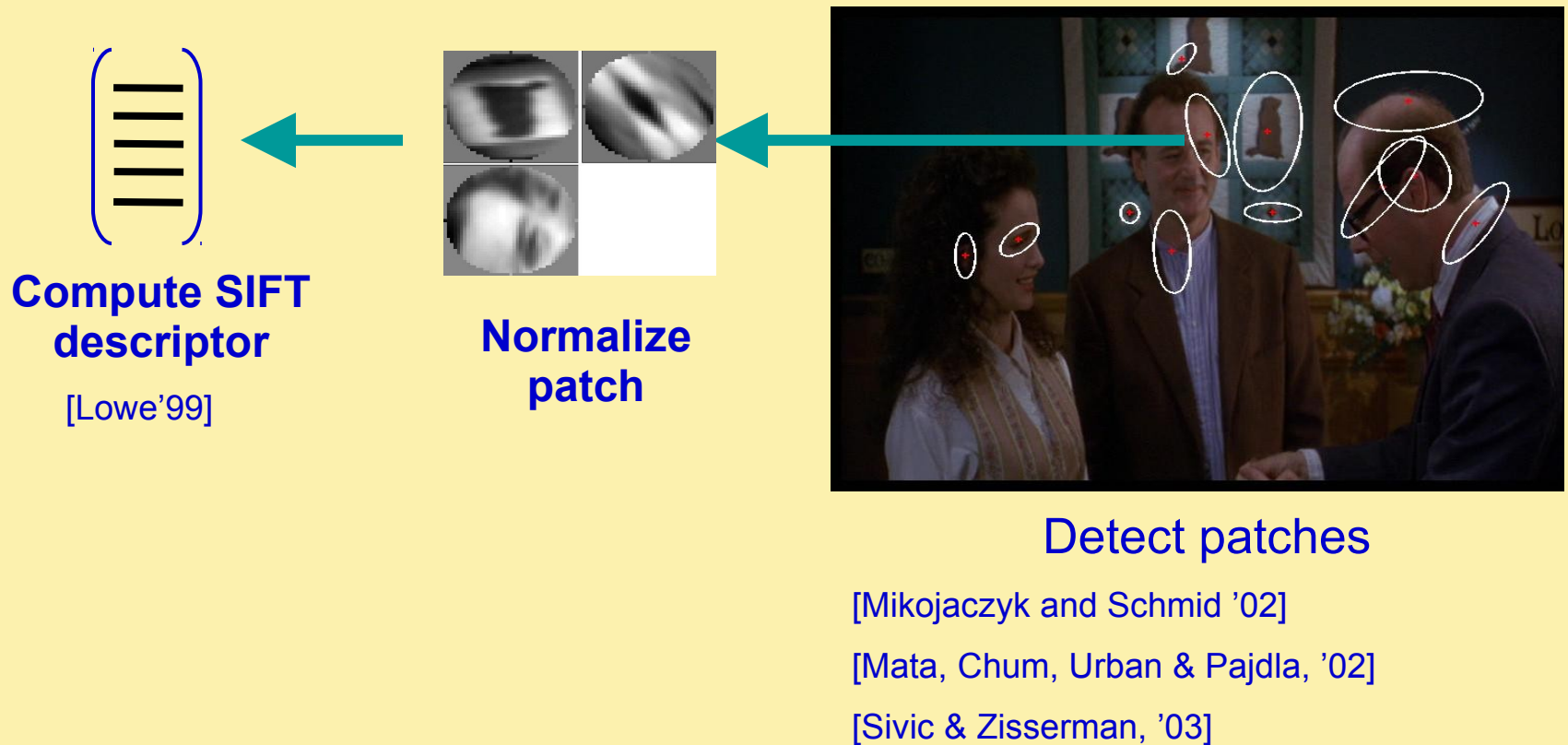
Sivic, et al. 2005



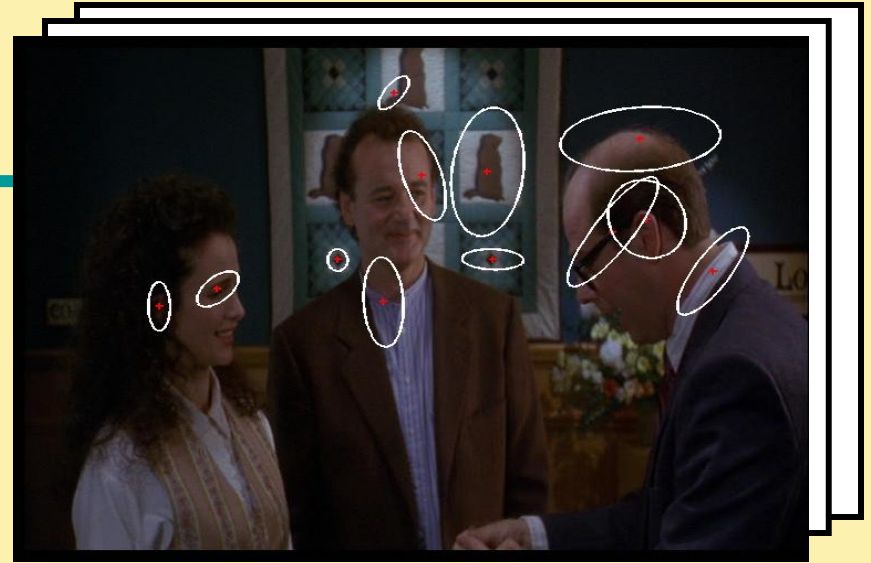
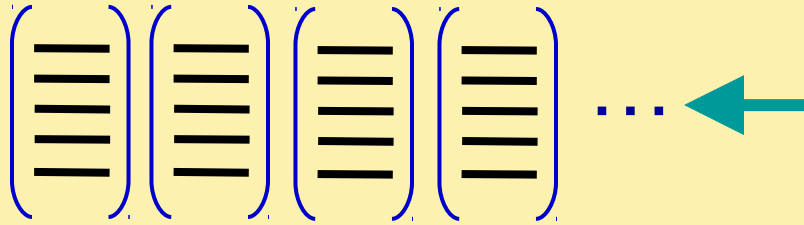
1.Feature detection and representation

- Regular grid
 - Vogel & Schiele, 2003
 - Fei-Fei & Perona, 2005
- Interest point detector
 - Csurka, Bray, Dance & Fan, 2004
 - Fei-Fei & Perona, 2005
 - Sivic, Russell, Efros, Freeman & Zisserman, 2005
- Other methods
 - Random sampling (Vidal-Naquet & Ullman, 2002)
 - Segmentation based patches (Barnard, Duygulu, Forsyth, de Freitas, Blei, Jordan, 2003)

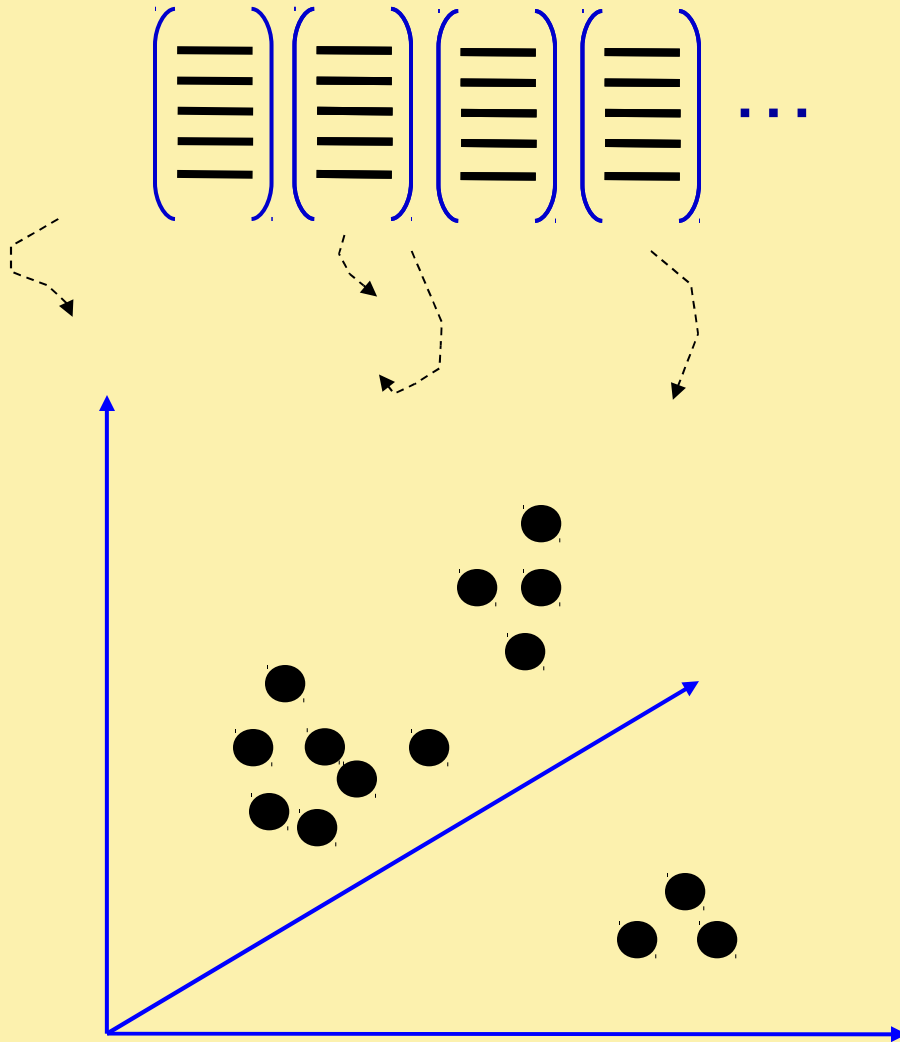
1.Feature detection and representation



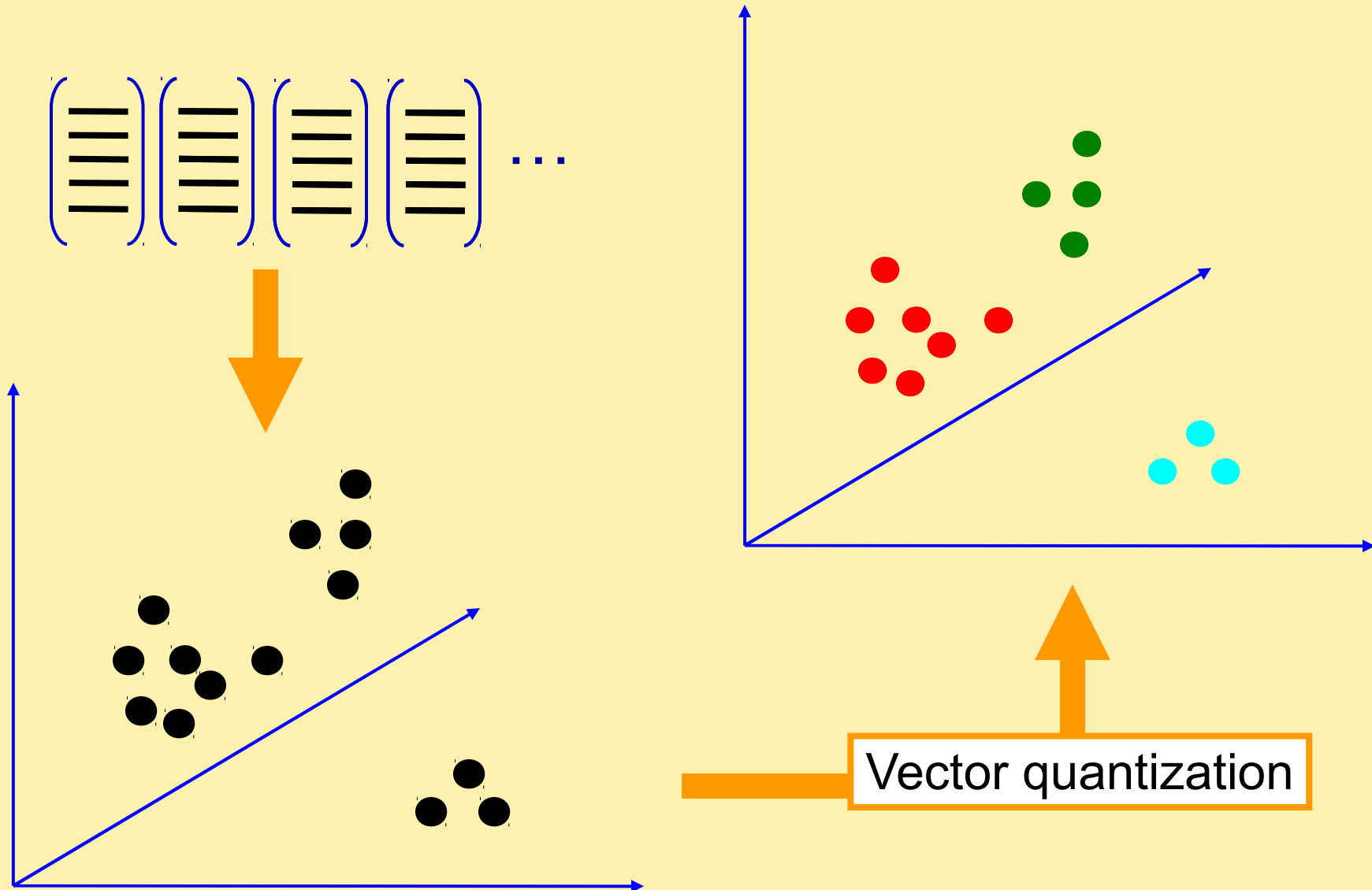
1. Feature detection and representation



2. Codewords dictionary formation



2. Codewords dictionary formation



2. Codewords dictionary formation

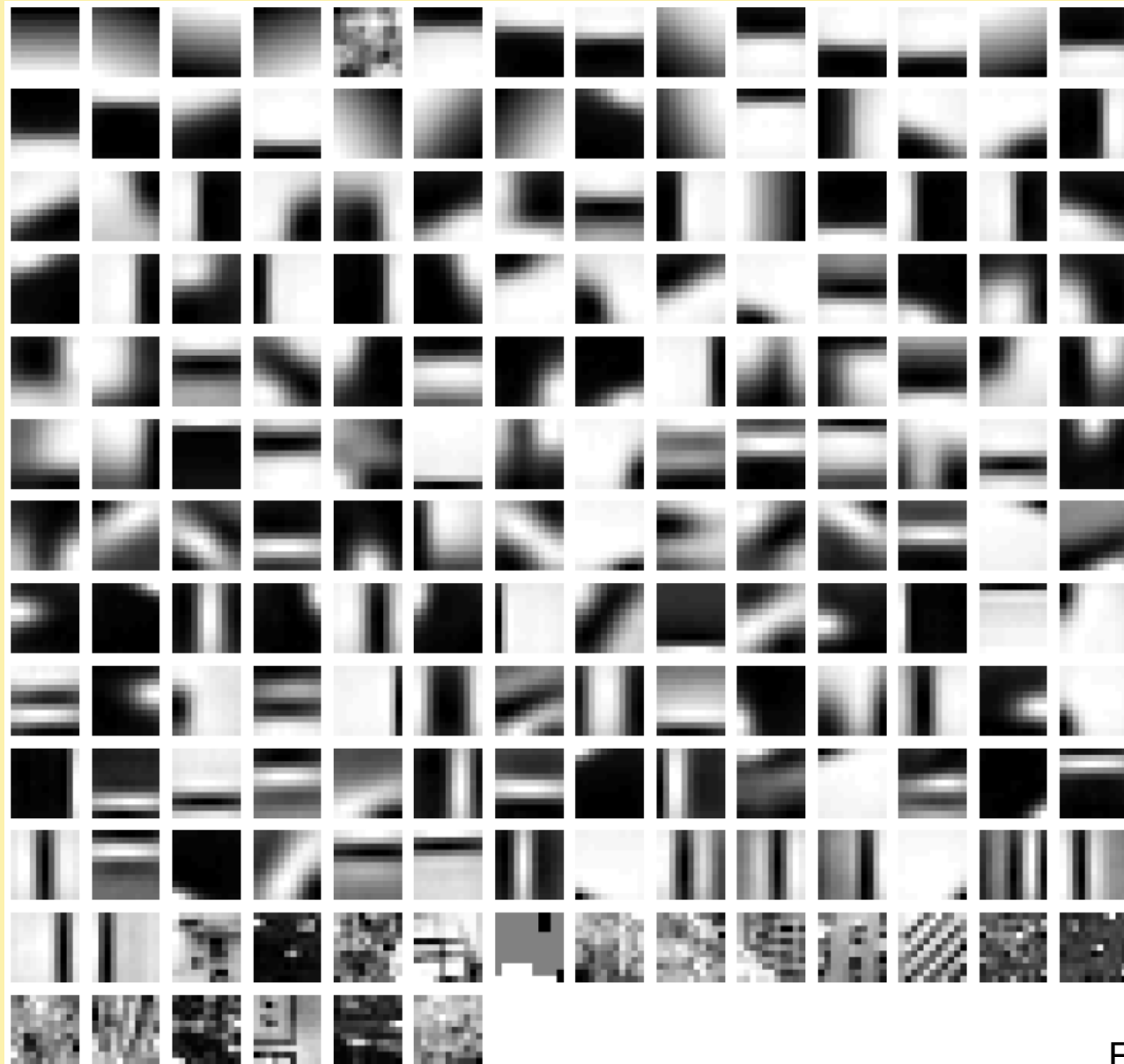
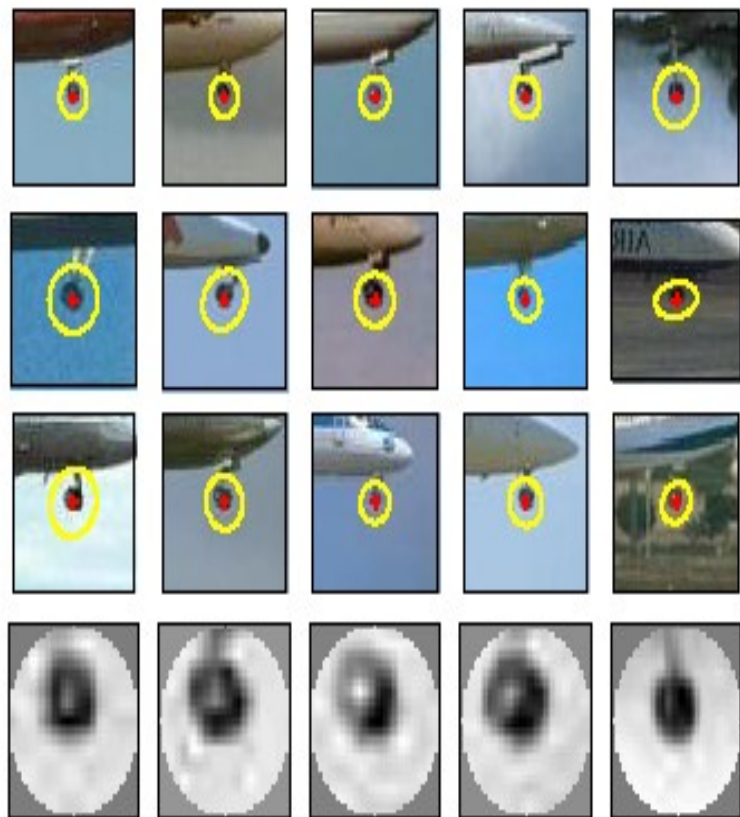


Image patch examples of codewords

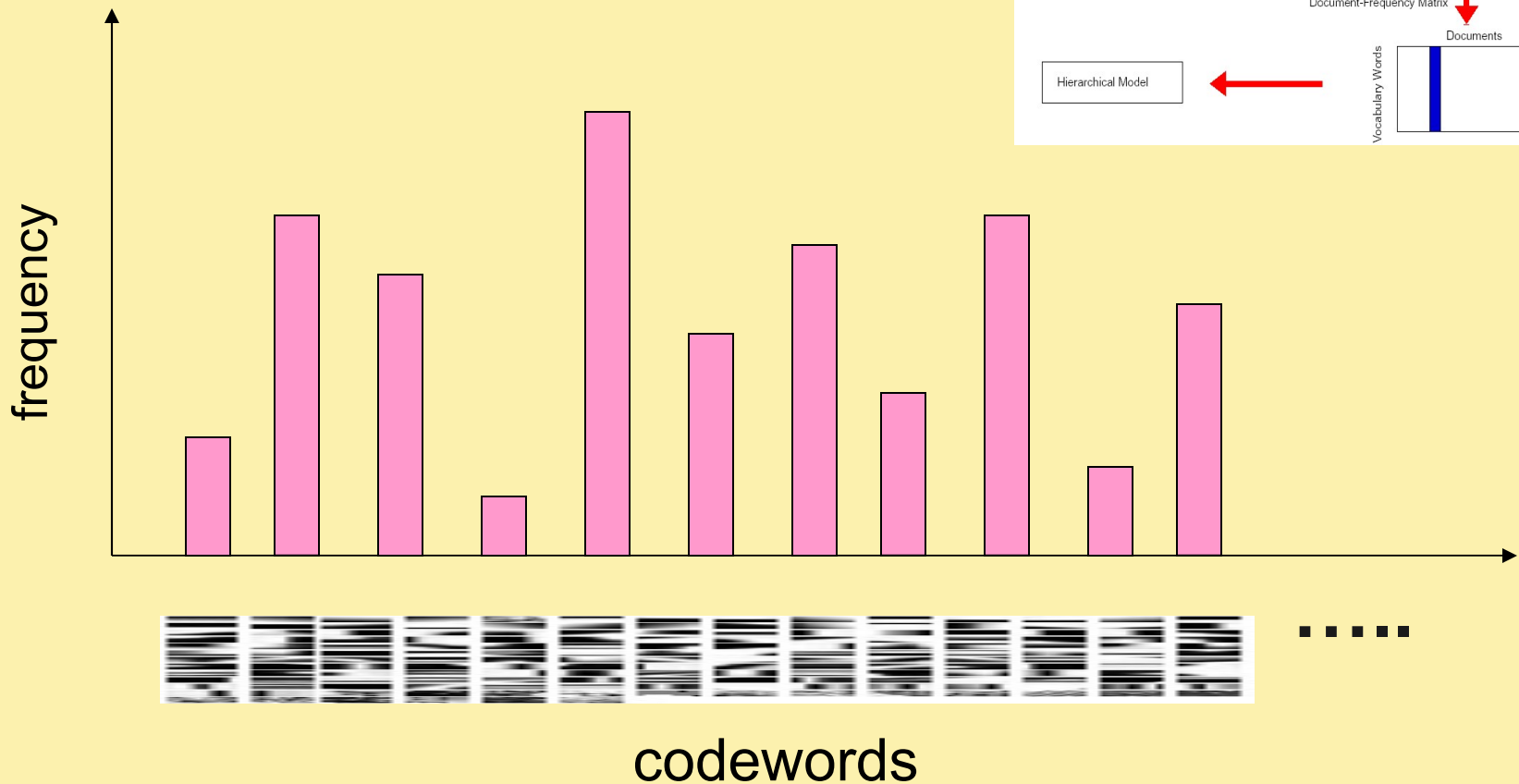
(a)



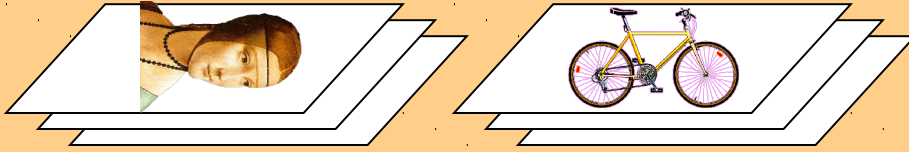
(b)



3. Image representation



Representation



2.

codewords dictionary

1. feature detection
& representation

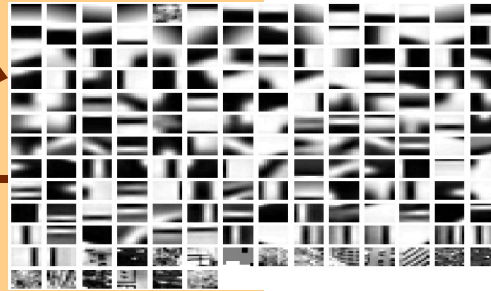


image representation

3.

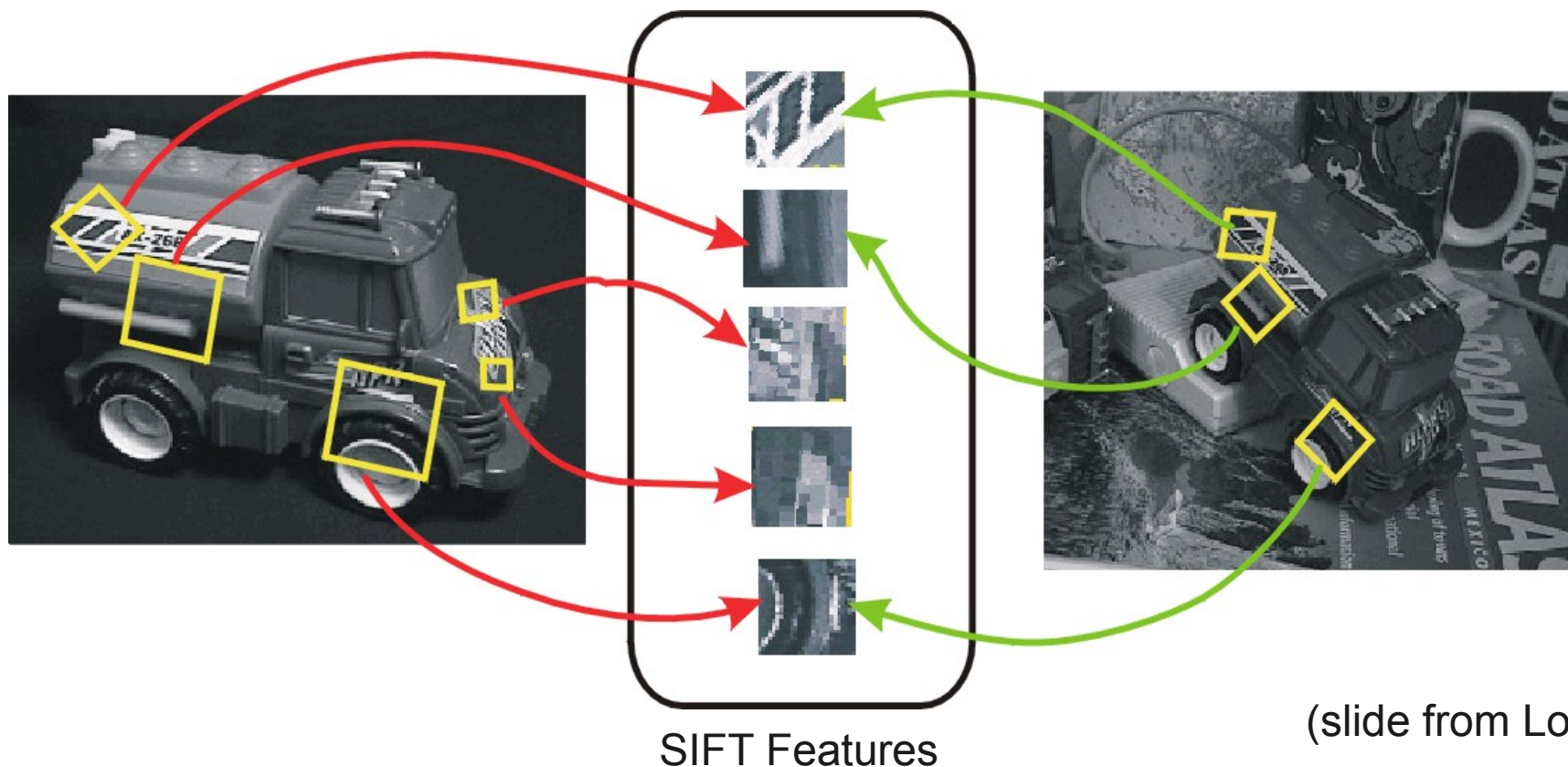


One of the keys to success is a good representation of features

- Just pixels is a bad representation
- Pixel intensities are affected by a lot of different things
 - Rotation, scaling, perspective
 - Illumination changes
 - Reordering of scenes
- We want a good way of characterizing image patches that is somewhat robust to these different effects

Scale-Invariant Local Features

- Image content is transformed into local feature coordinates that are invariant to translation, rotation, scale, and other imaging parameters



Advantages of invariant local features

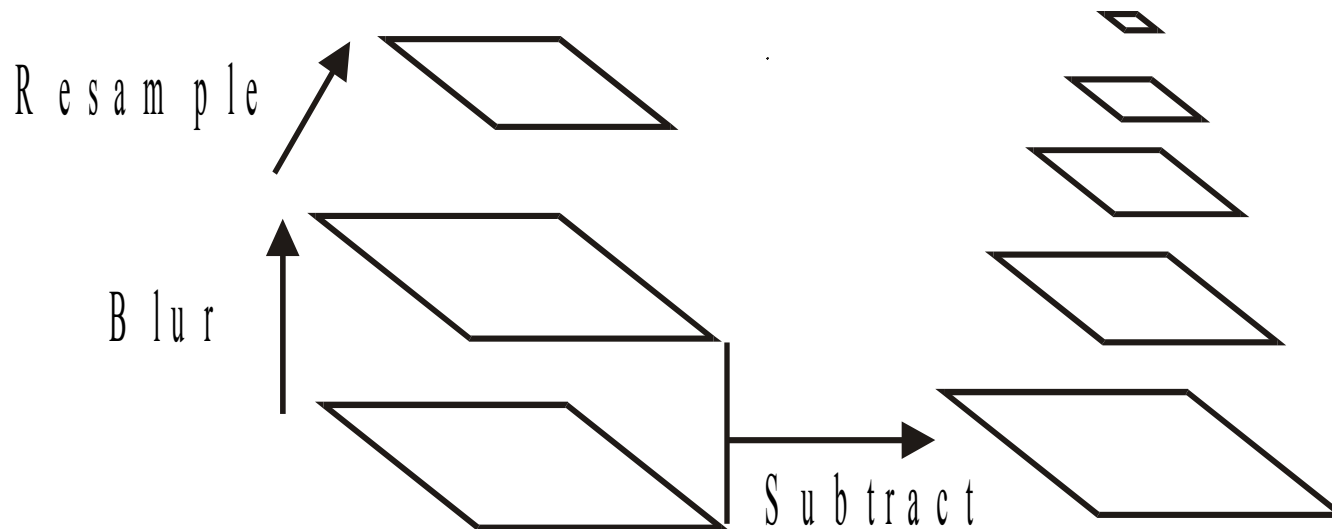
- **Locality:** features are local, so robust to occlusion and clutter (no prior segmentation)
- **Distinctiveness:** individual features can be matched to a large database of objects
- **Quantity:** many features can be generated for even small objects
- **Efficiency:** close to real-time performance
- **Extensibility:** can easily be extended to wide range of differing feature types, with each adding robustness

Think Back to Bag of Words - Two Key Problems

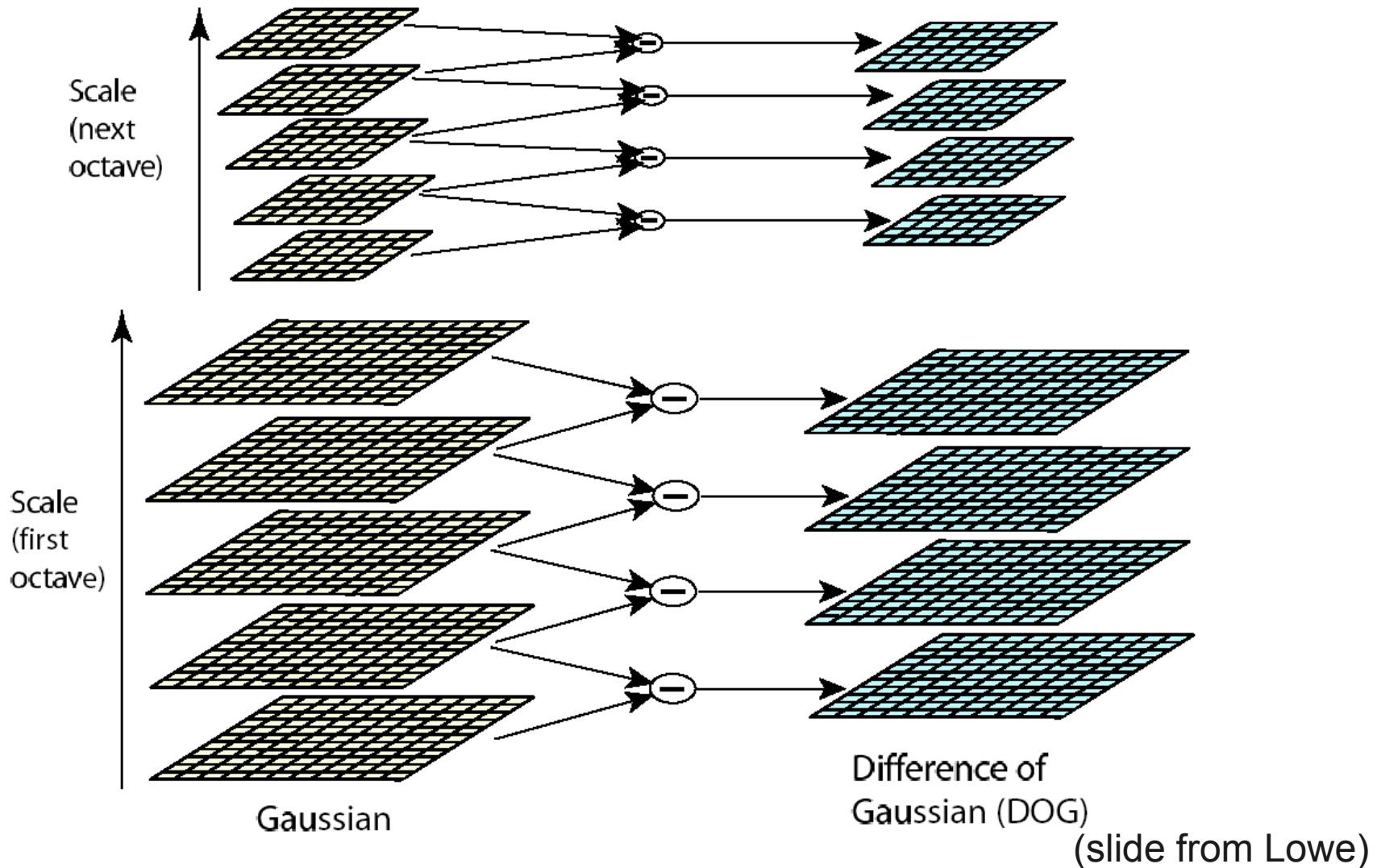
- Problem 1: What parts of the image do I look at?
- Problem 2: How do I represent the patches of pixels

Build Scale-Space Pyramid

- All scales must be examined to identify scale-invariant features
- An efficient function is to compute the Difference of Gaussian (DOG) pyramid (Burt & Adelson, 1983)

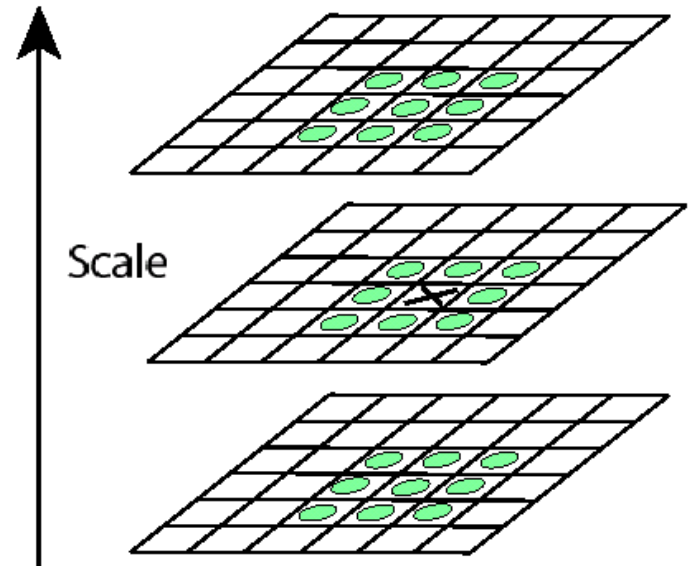


Scale space processed one octave at a time



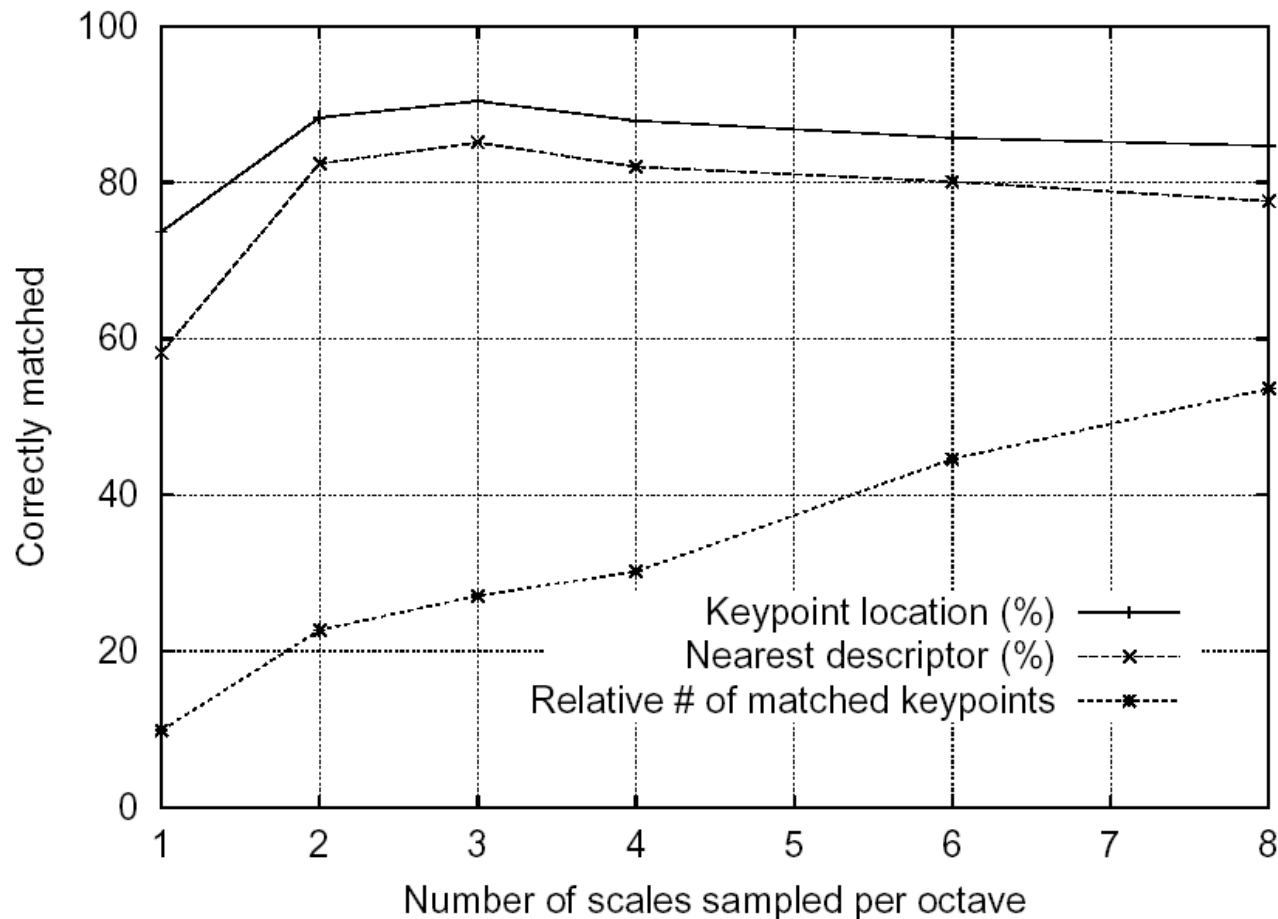
Key point localization

- Detect maxima and minima of difference-of-Gaussian in scale space



Sampling frequency for scale

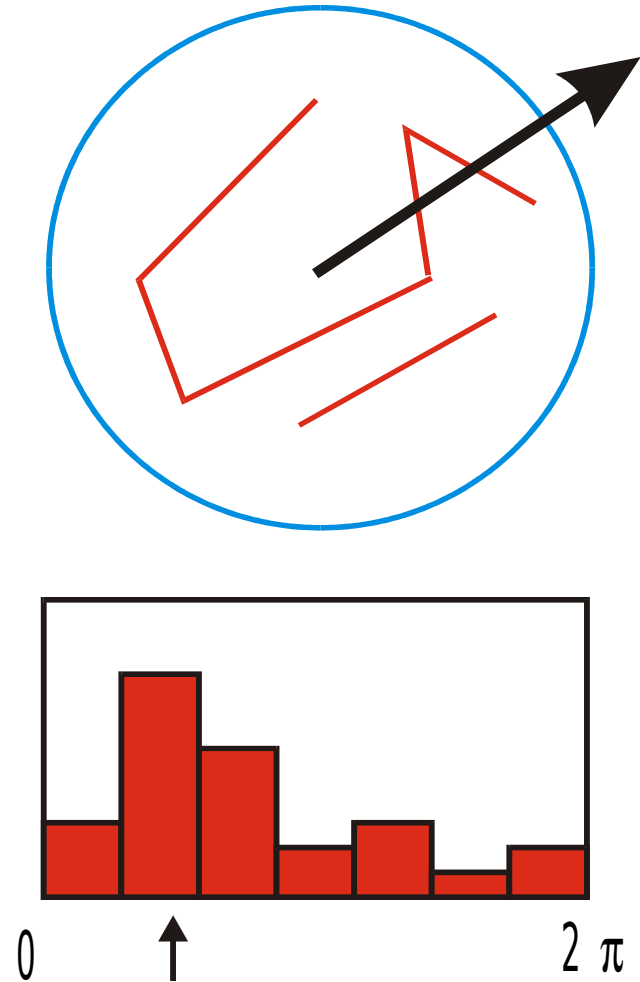
More points are found as sampling frequency increases, but accuracy of matching decreases after 3 scales/octave



(slide from Lowe)

Select canonical orientation

- Create histogram of local gradient directions computed at selected scale
- Assign canonical orientation at peak of smoothed histogram
- Each key specifies stable 2D coordinates (x, y, scale, orientation)



(slide from Lowe)

Example of keypoint detection

Threshold on value at DOG peak and on ratio of principle curvatures (Harris approach)



- (a) 233x189 image
- (b) 832 DOG extrema
- (c) 729 left after peak value threshold
- (d) 536 left after testing ratio of principle curvatures

(slide from Lowe)

Detecting Keypoints is not always better

Descriptor	Grid	Random	Saliency [4]	DoG [7]
11 × 11 Pixel	64.0%	47.5%	45.5%	N/A
128-dim Sift	65.2%	60.7%	53.1%	52.5%

(From L. Fei-Fei and Perona)

SIFT vector formation

- Thresholded image gradients are sampled over 16x16 array of locations in scale space
- Create array of orientation histograms
- 8 orientations x 4x4 histogram array = 128 dimensions

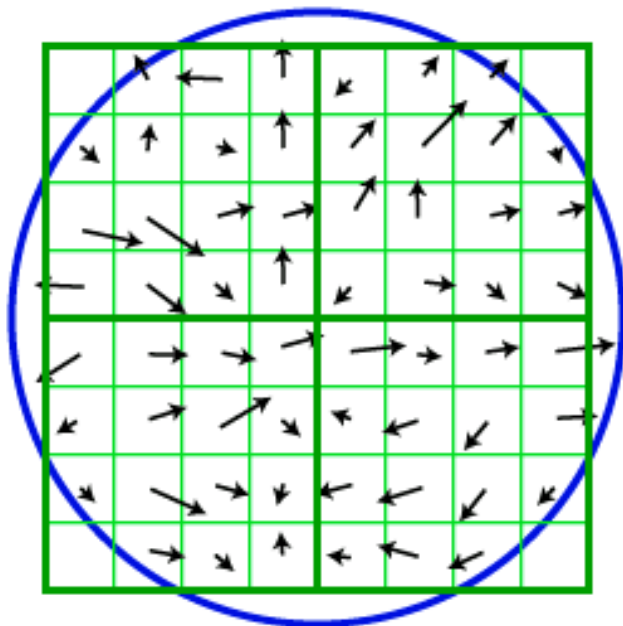
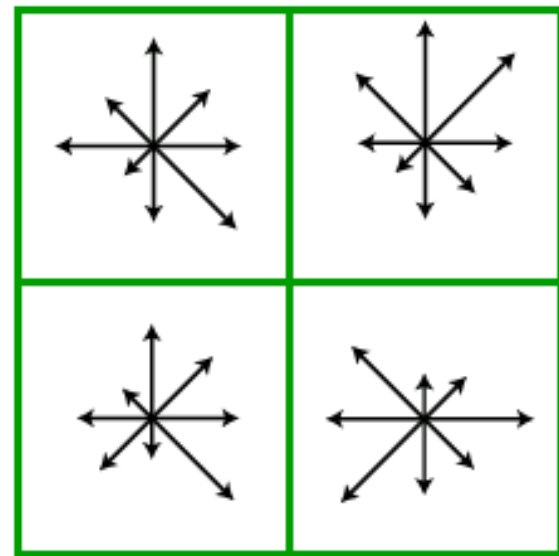


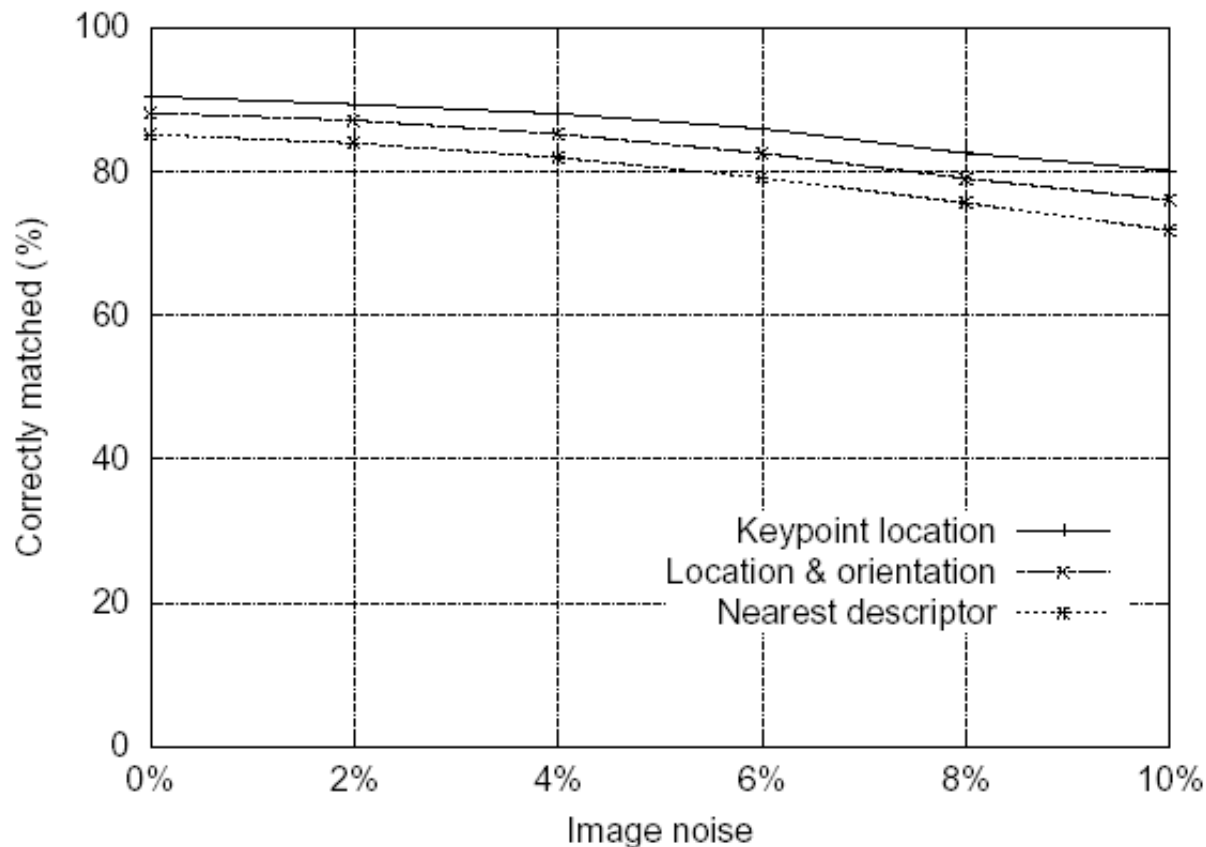
Image gradients



Keypoint descriptor
(slide from Lowe)

Feature stability to noise

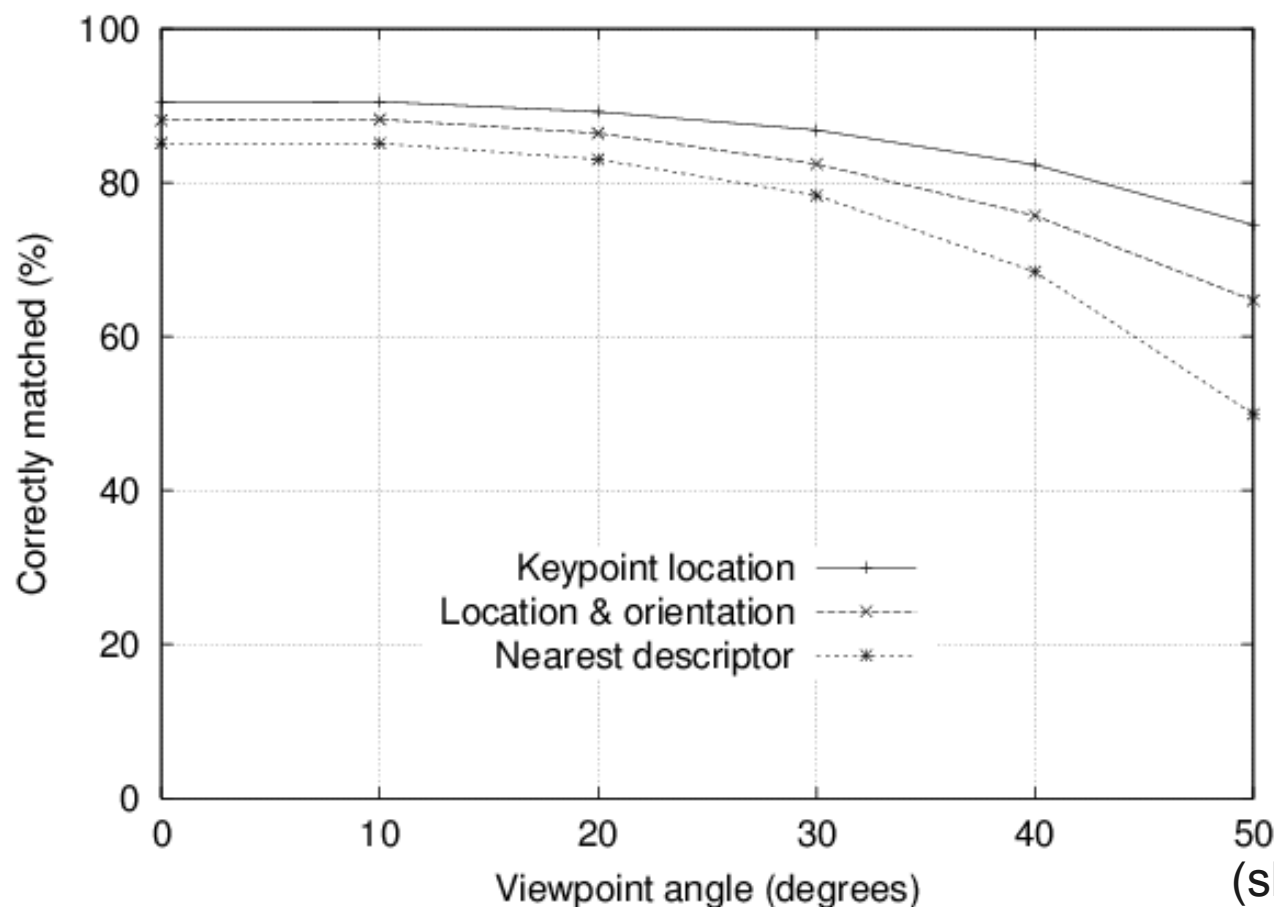
- Match features after random change in image scale & orientation, with differing levels of image noise
- Find nearest neighbor in database of 30,000 features



(slide from Lowe)

Feature stability to affine change

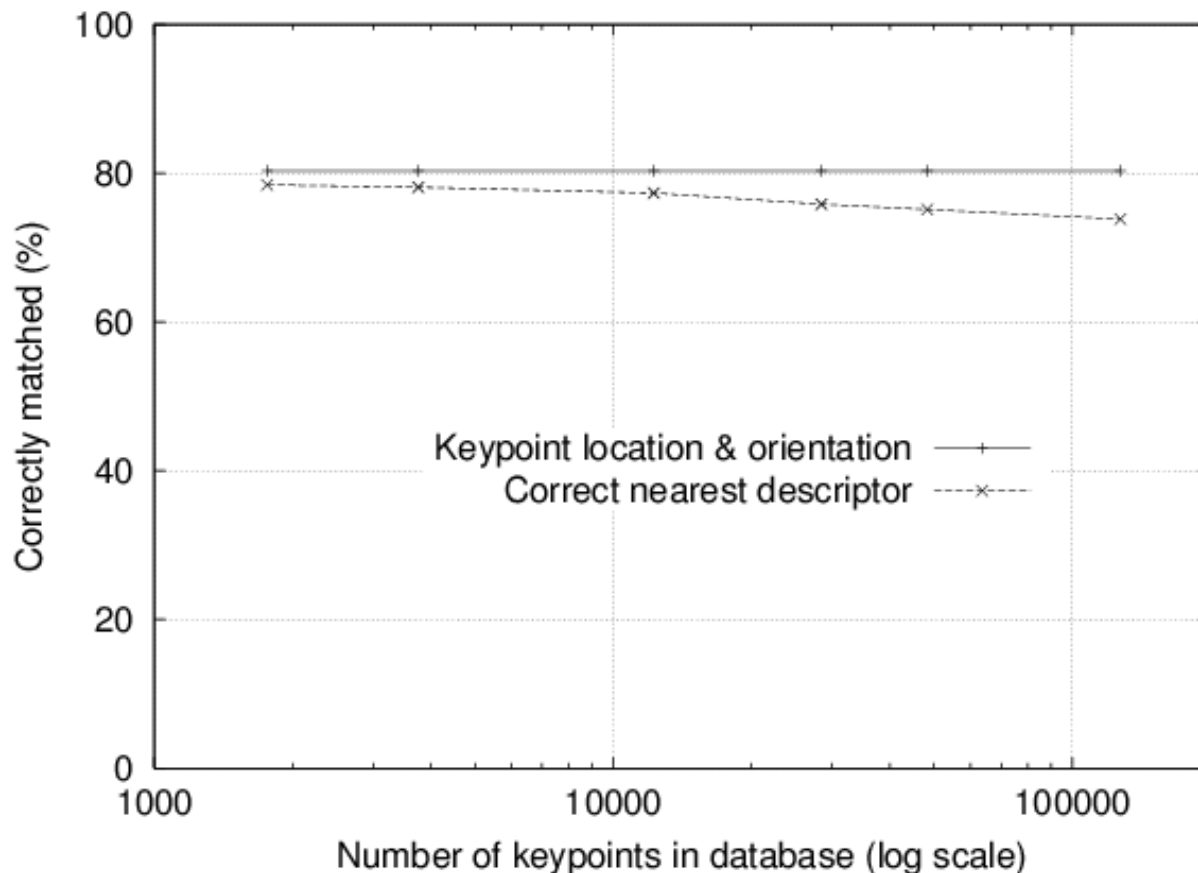
- Match features after random change in image scale & orientation, with 2% image noise, and affine distortion
- Find nearest neighbor in database of 30,000 features



(slide from Lowe)

Distinctiveness of features

- Vary size of database of features, with 30 degree affine change, 2% image noise
- Measure % correct for single nearest neighbor match



(slide from Lowe)

Sony Aibo (Evolution Robotics)

SIFT usage:

- Recognize charging station
- Communicate with visual cards

AIBO® Entertainment Robot

Official U.S. Resources and Online Destinations



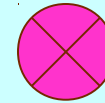
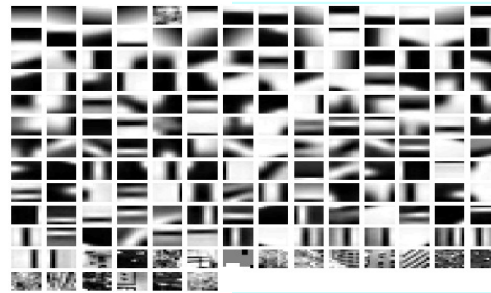
SIFT is just the beginning

- Authors have proposed more feature point detectors
 - Harris-Laplace,....
- Authors have proposed other feature descriptors
 - ColorSIFT
 - SURF
- The Koen executable implements many of this

Learning and Recognition



codewords dictionary

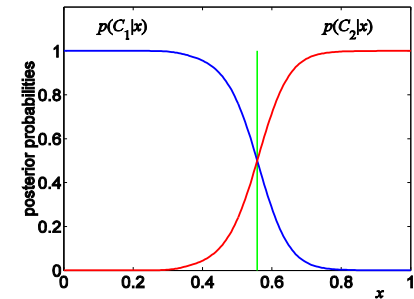
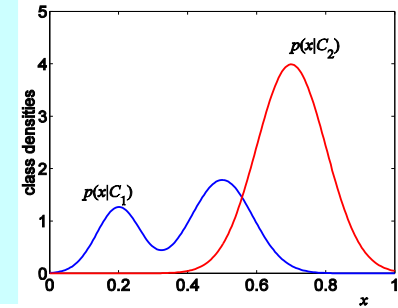


**category models
(and/or) classifiers**

**category
decision**

Learning and Recognition

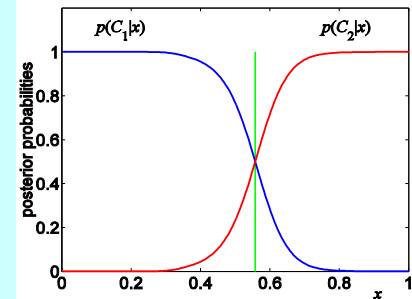
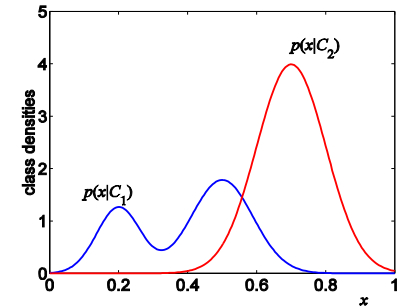
1. Generative method:
 - graphical models
2. Discriminative method:
 - SVM



**category models
(and/or) classifiers**

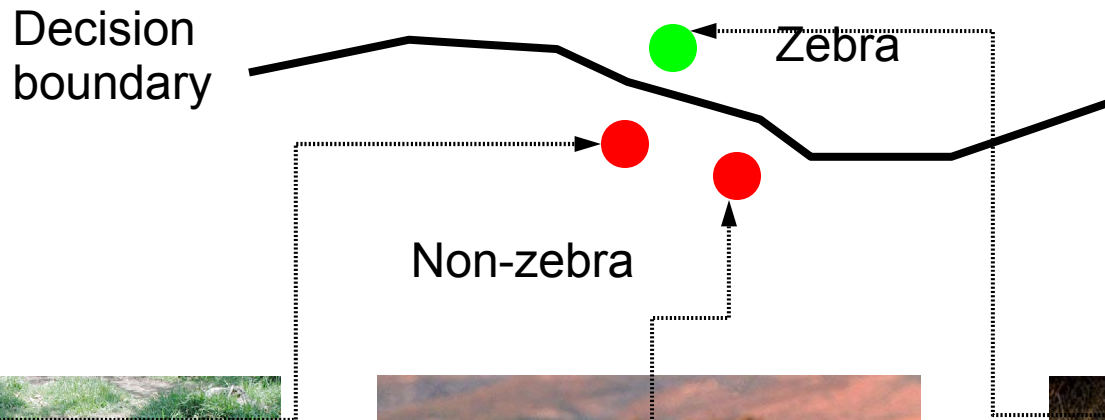
Learning and Recognition

1. Generative method:
 - graphical models
2. Discriminative method:
 - SVM



**category models
(and/or) classifiers**

Discriminative methods based on 'bag of words' representation



Discriminative methods based on 'bag of words' representation

Grauman & Darrell, 2005, 2006:

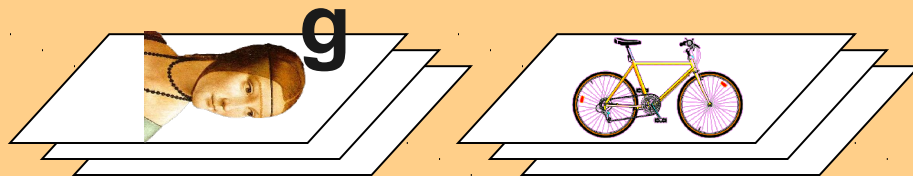
SVM w/ Pyramid Match kernels

Others

Csurka, Bray, Dance & Fan, 2004

Serre & Poggio, 2005

learning



feature detection
& representation

image representation

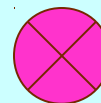
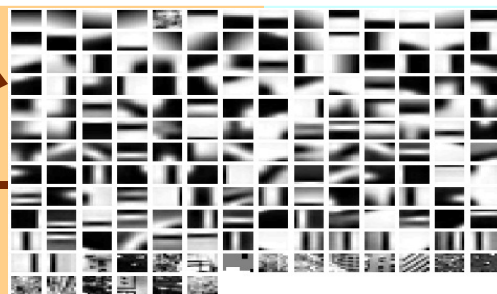


**category models
(and/or) classifiers**

recognition

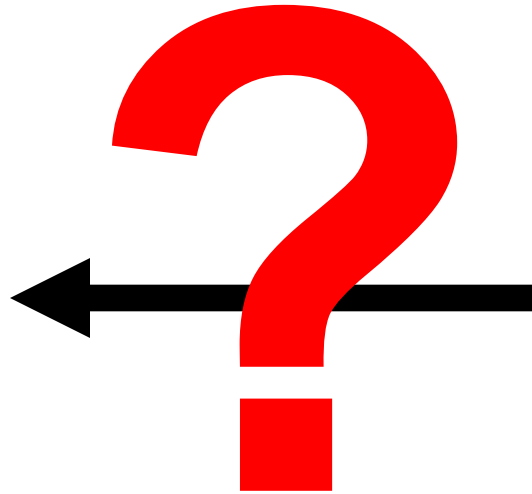


codewords dictionary



**category
decision**

What about spatial info?



What about spatial info?

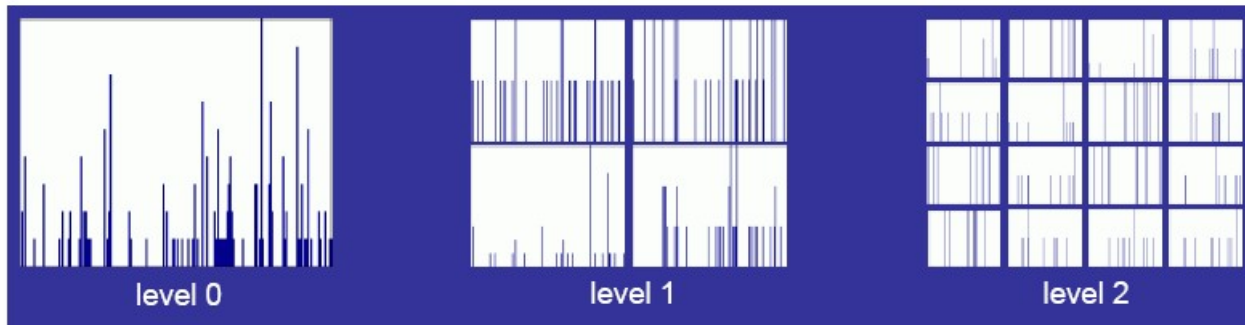
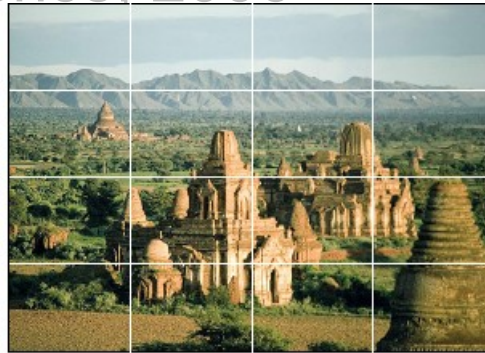


Feature level

Generative models

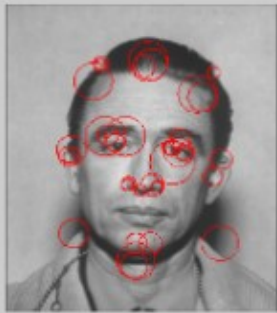
Discriminative methods

Lazebnik, Schmid & Ponce, 2006



Invariance issues

- Scale and rotation
 - Implicit
 - Detectors and descriptors



Invariance issues

- Scale and rotation
- Occlusion
 - Implicit in the models
 - Codeword distribution: small variations
 - (In theory) Theme (z) distribution: different occlusion patterns



Invariance issues

- Scale and rotation
- Occlusion
- Translation
 - Encode (relative) location information
 - Sudderth, Torralba, Freeman & Willsky, 2005, 2006
 - Niebles & Fei-Fei, 2007



Invariance issues

- Scale and rotation
- Occlusion
- Translation
- View point (in theory)
 - Codewords: detector and descriptor
 - Theme distributions: different view points



Model properties



- Intuitive
 - Analogy to documents

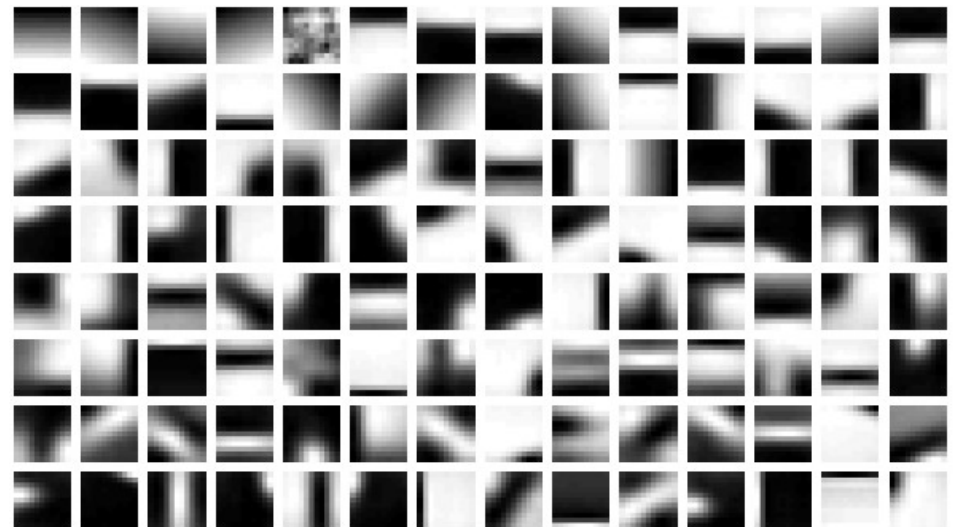
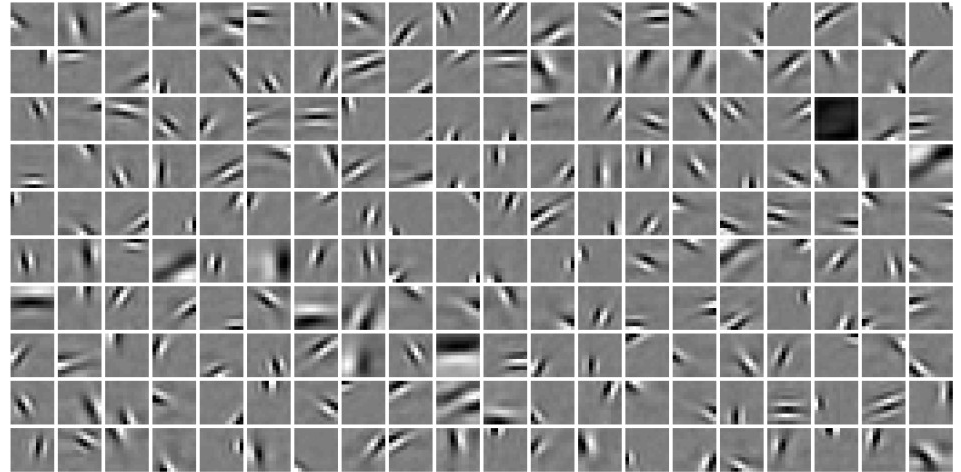
Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes. For a long time, the visual image was considered as a movie image. It was discovered that the visual image is not a movie image but a more complex one. Following the discovery of the path to the various centers of the cortex, Hubel and Wiesel have demonstrated that the message about the image falling on the retina undergoes a point-by-point analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.

**sensory, brain,
visual, perception,
retinal, cerebral cortex,
eye, cell, optical
nerve, image
Hubel, Wiesel**

Model properties

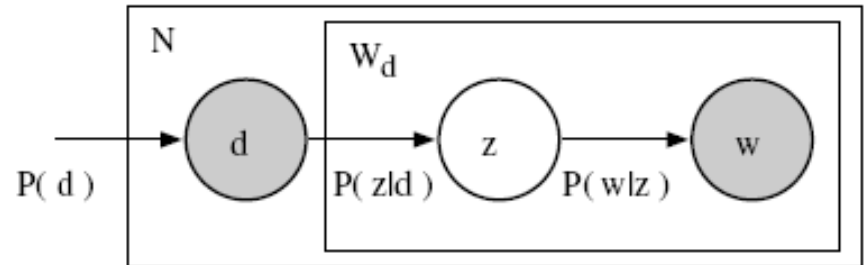


- Intuitive
 - Analogy to documents
 - Analogy to human vision



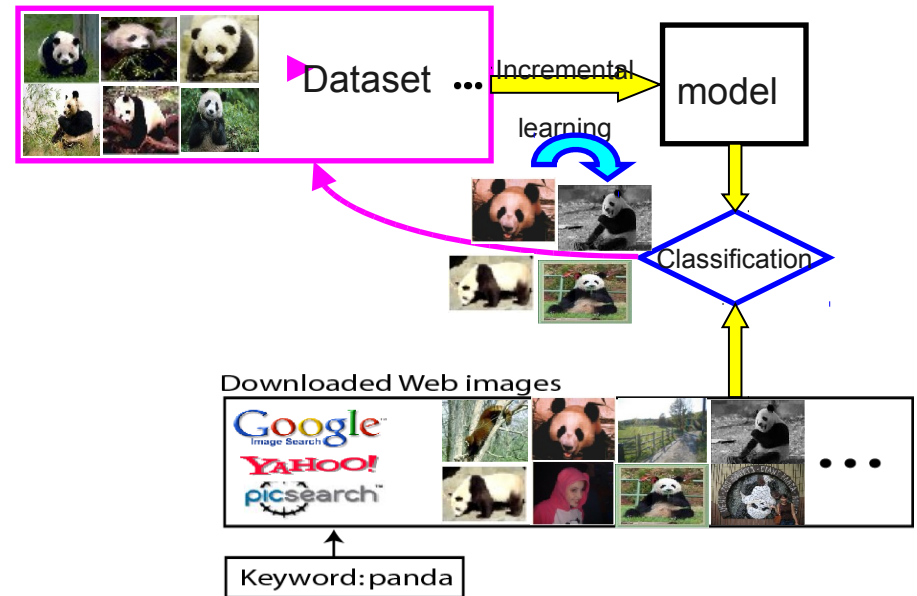


Model properties



Sivic, Russell, Efros, Freeman, Zisserman, 2005

- Intuitive
- generative models
 - Convenient for weakly- or unsupervised, incremental training
 - Prior information
 - Flexibility (e.g. HDP)

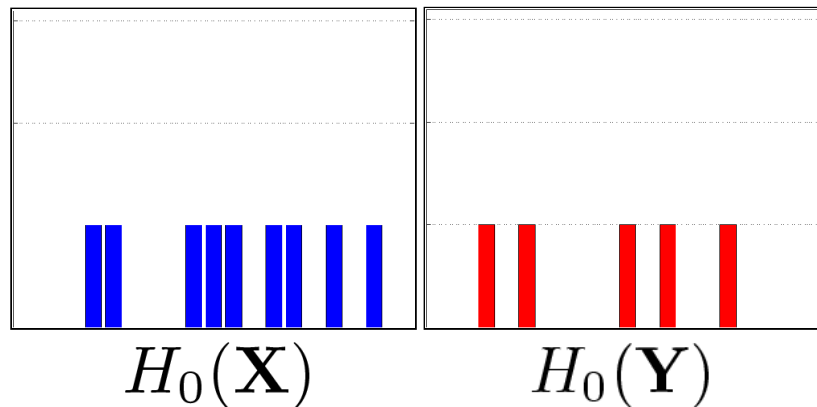
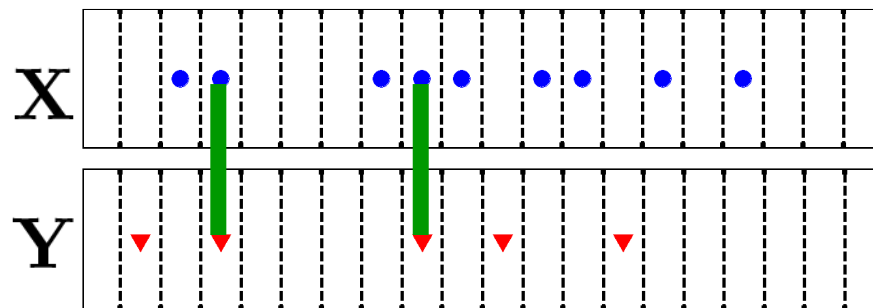


Li, Wang & Fei-Fei, CVPR 2007



Model properties

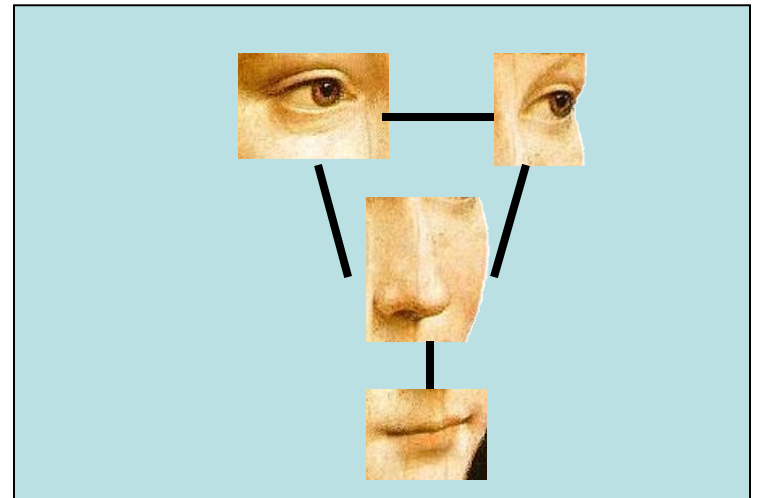
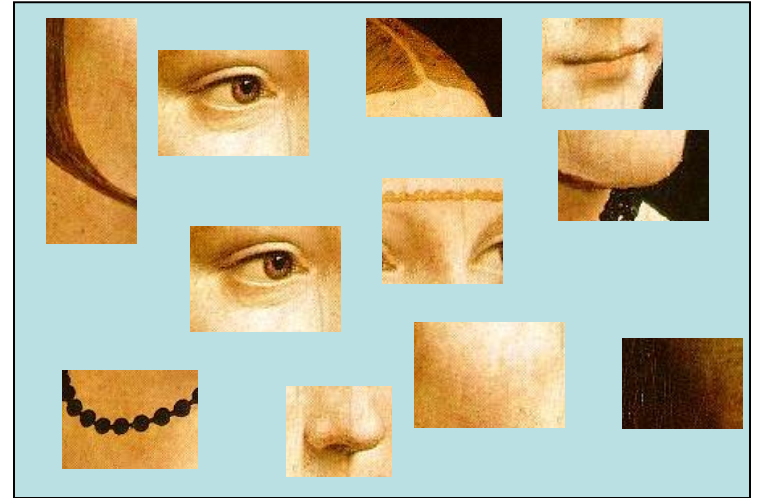
- Intuitive
- generative models
- Discriminative method
 - Computationally efficient



Model properties



- Intuitive
- generative models
- Discriminative method
- Learning and recognition relatively fast
 - Compare to other methods





Weakness of the model

- No rigorous geometric information of the object components
- It's intuitive to most of us that objects are made of parts – no such information
- Not extensively tested yet for
 - View point invariance
 - Scale invariance
- Segmentation and localization unclear