# Object Categorization

# Bag-of-words models

**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 essentially on the messages that reach the brain from our eyes. For a long time it was thought that the retinal image was transmitted point by point to visual centers in the brain; the cerebral cortex was a movie screen, so to speak, upon which the image in the eye was projected. Through the discoveries of Hubel and Wiesel we now know that behind the origin of the visual perception in the brain there is a considerably more complicated course of events. By following the visual impulses along their path to the various cell layers of the optical cortex, Hubel and Wiesel have been able to 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**
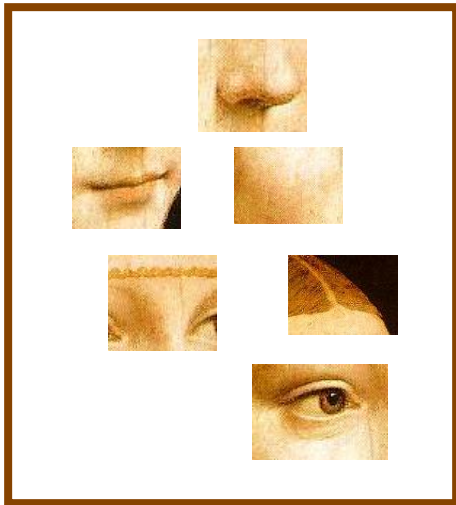
China is forecasting a trade surplus of $90bn (£51bn) to $100bn this year, a threefold increase on 2004's $32bn. The Commerce Ministry said the surplus would be created by a predicted 30% jump in exports to $750bn, compared with a 18% rise in imports to $660bn. The figures are likely to further annoy the US, which has long argued that China's exports are unfairly helped by a deliberately undervalued yuan. Beijing agrees the surplus is too high, but says the yuan is only one factor. Bank of China governor Zhou Xiaochuan said the country also needed to do more to boost domestic demand so more goods stayed within the country. China increased the value of the yuan against the dollar by 2.1% in July and permitted it to trade within a narrow band, but the US wants the yuan to be allowed to trade freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

**China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value**
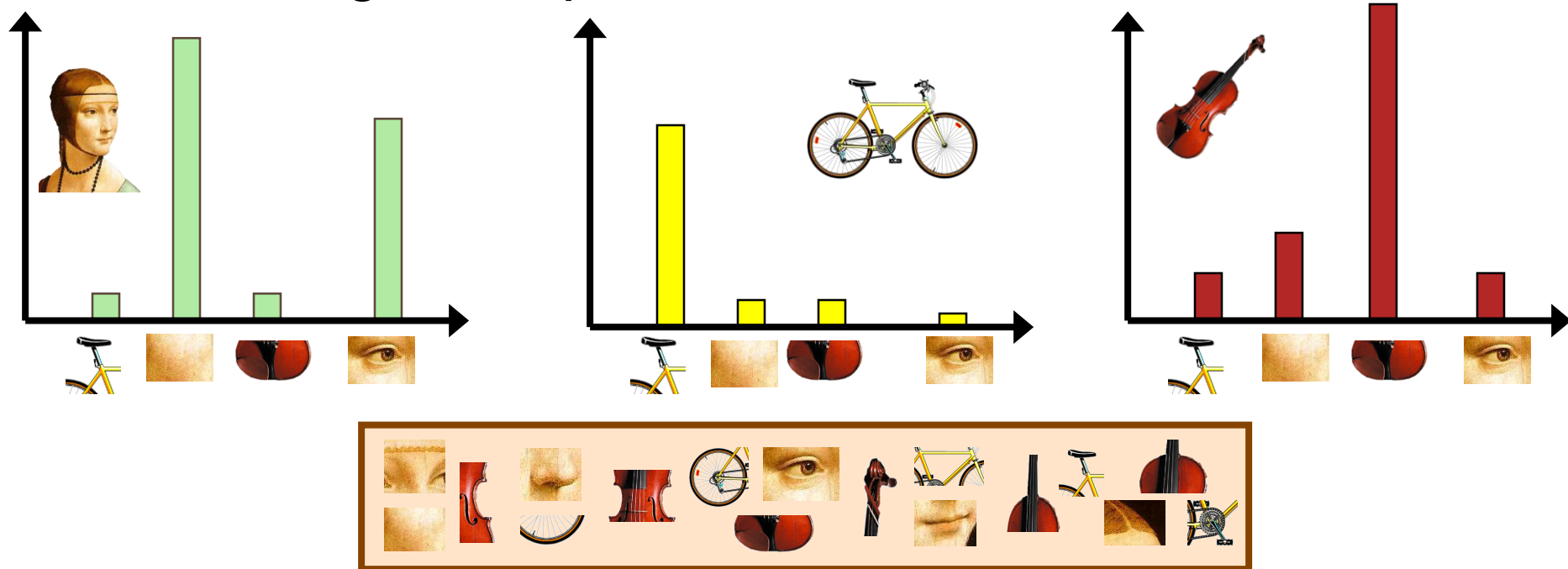
# A clarification: definition of "BoW"

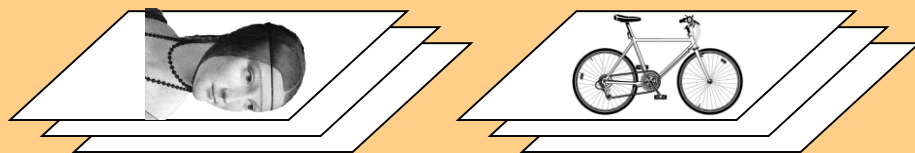- Looser definition
  - Independent features

# A clarification: definition of "BoW"

- Looser definition
  - Independent features

- Stricter definition
  - Independent features
  - histogram representation

**learning**

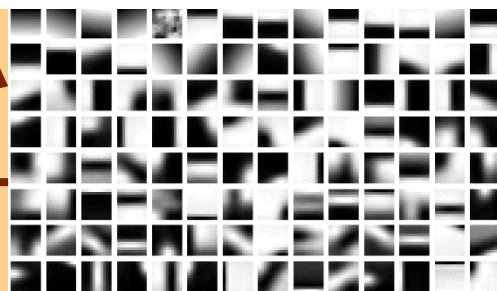**recognition**

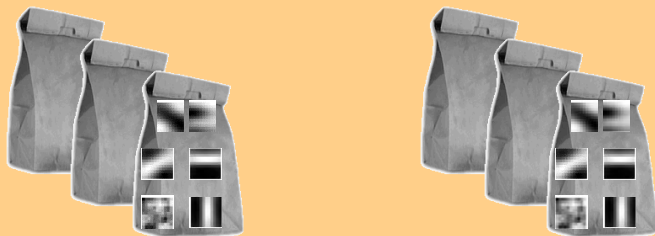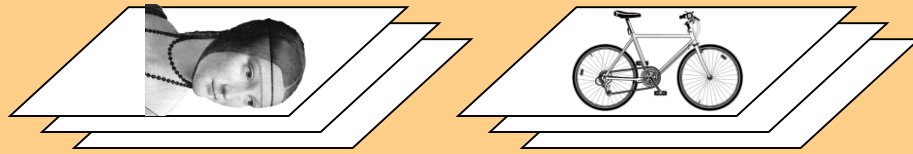feature detection & representation

**codewords dictionary**

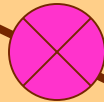image representation

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

**category decision**

# Representation



**1.** feature detection & representation

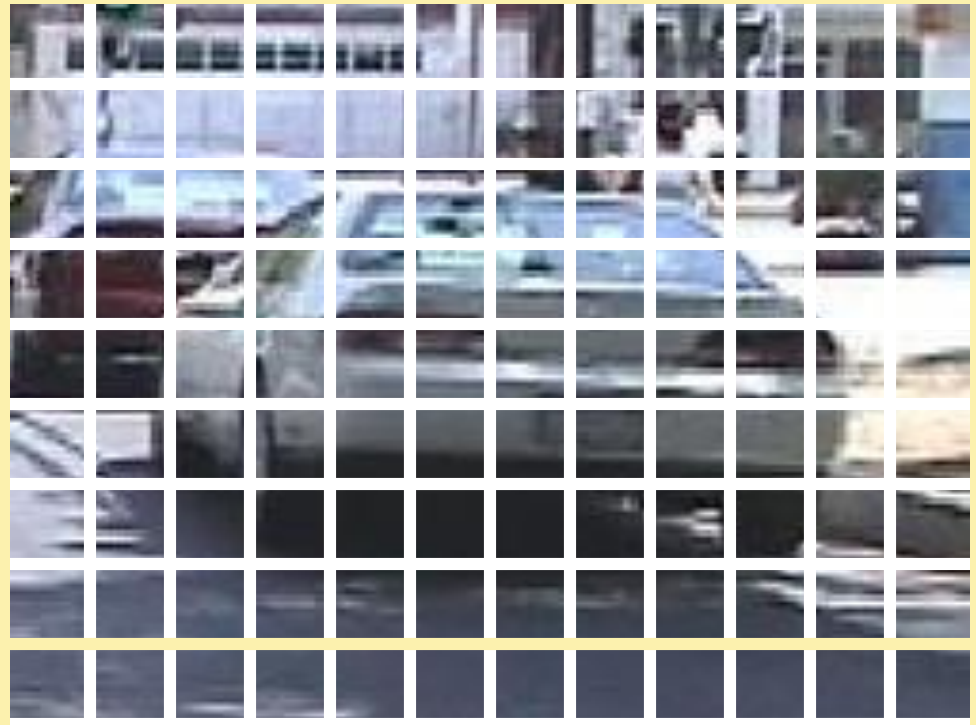**2.** codewords dictionary

image representation

**3.**

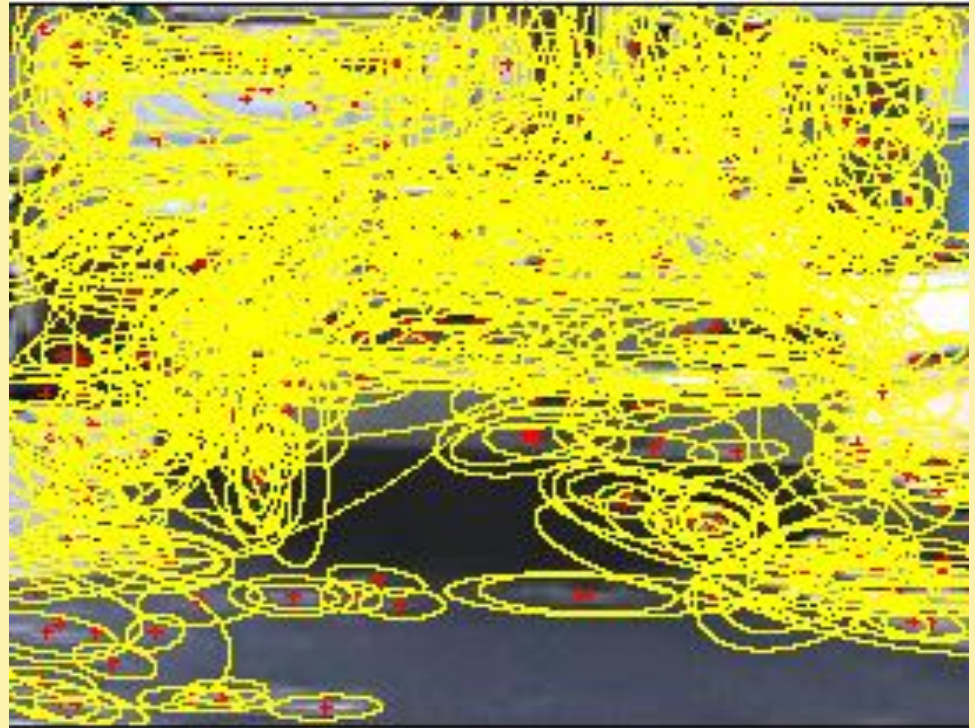# 1.Feature detection and representation

# 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
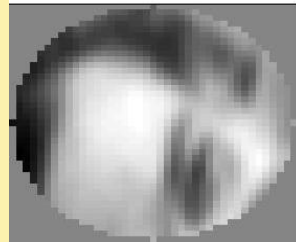  - Fei-Fei & Perona, 2005
  - Sivic, et al. 2005

# 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)
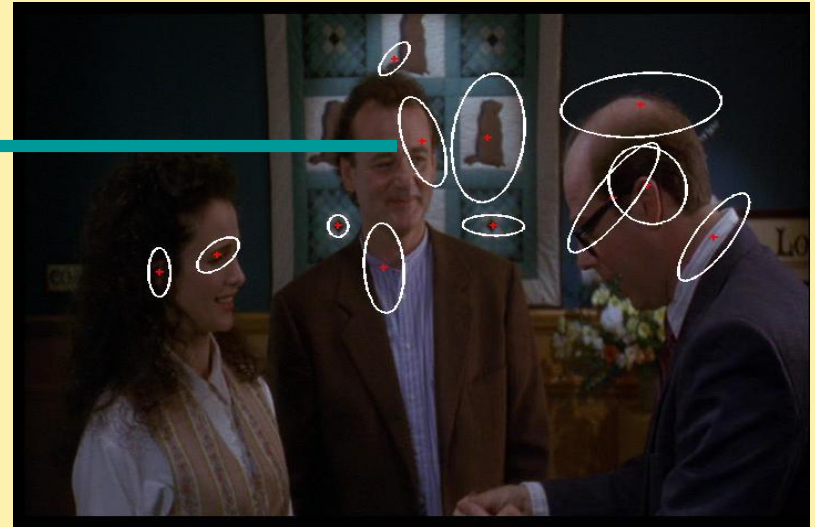
# 1.Feature detection and representation



**Compute
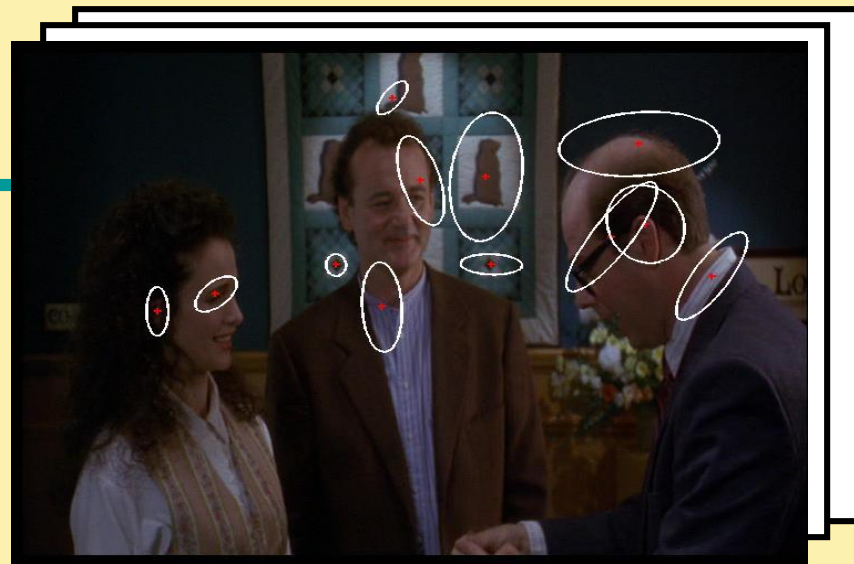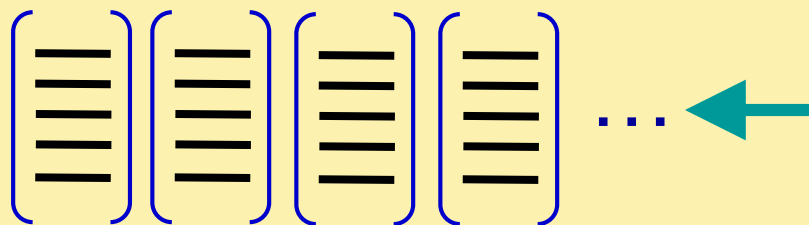SIFT
descriptor**

[Lowe'99]

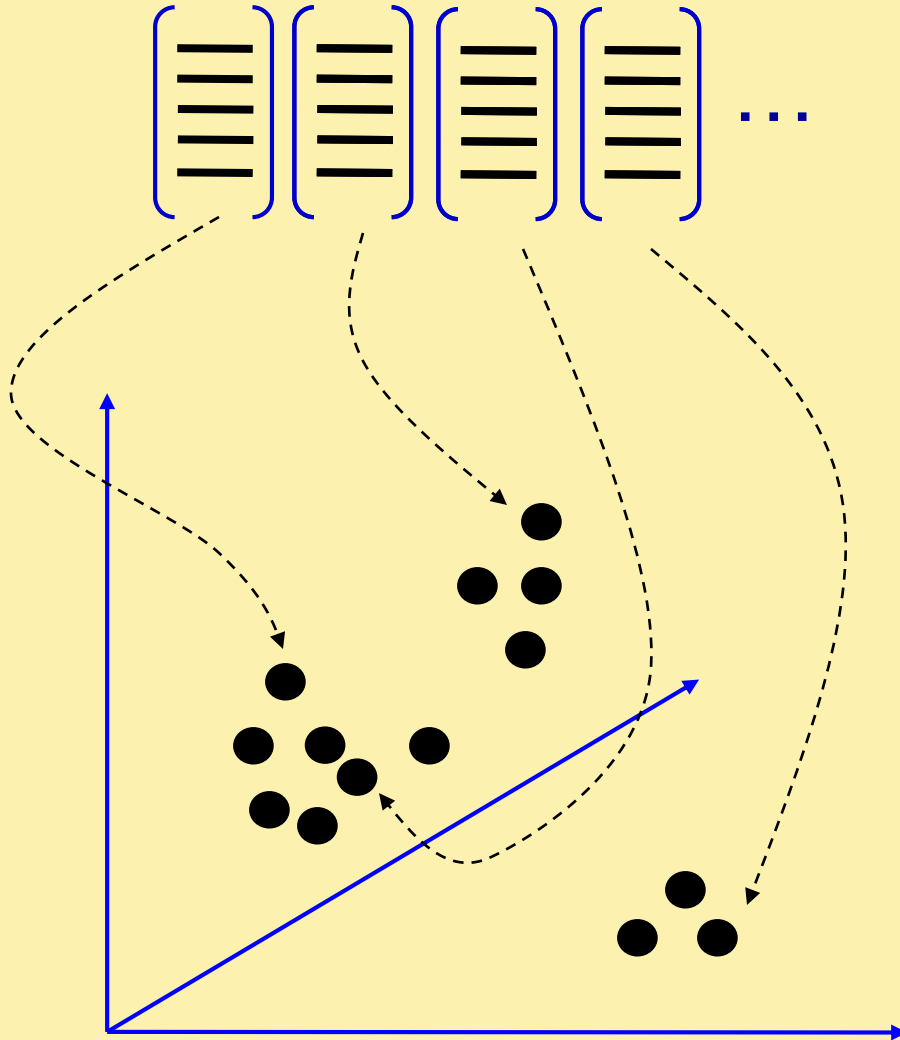**Normalize
patch**

Detect patches

[Mikojaczyk and Schmid '02]

[Mata, Chum, Urban & Pajdla, '02]

[Sivic & Zisserman, '03]

Slide credit: Josef Sivic

# 1.Feature detection and representation

# 2. Codewords dictionary formation

# 2. Codewords dictionary formation



**Clustering**

Vector quantization

# 2. Codewords dictionary formation



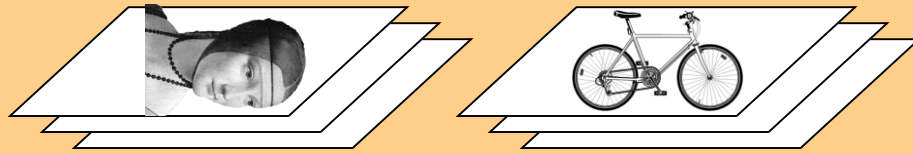Regular grid

Fei-Fei et al. 2005
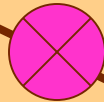
# Image patch examples of codewords

# 3. Image representation

# Representation



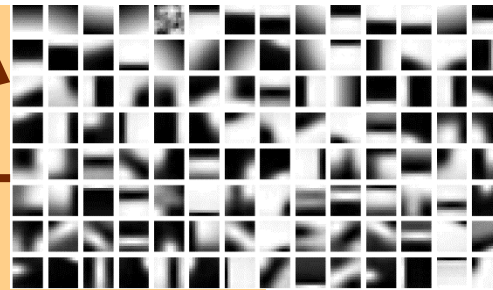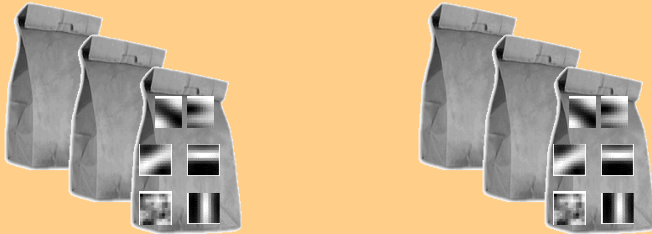**1.** feature detection & representation

**2.** codewords dictionary

image representation

**3.**

# Learning and Recognition

**codewords dictionary**

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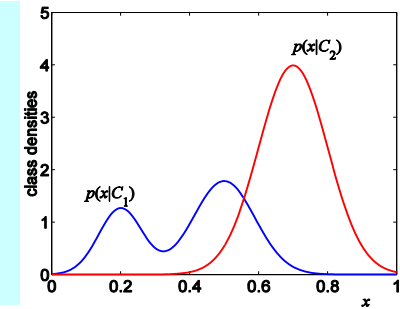
**category decision**

# BoW-based Object Categorization
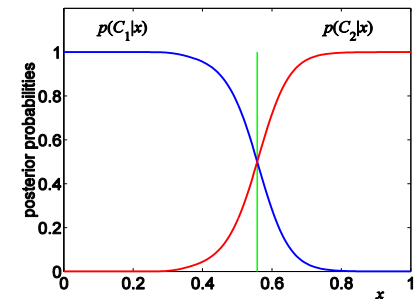
- Basic steps
  - Feature extraction and representation
  - Building codebook (codewords dictionary) from training samples with clustering
  - Represent an image with histogram of codebook (i.e. Bag-of-words of an image)
  - Classify an unknown image with its BoW.

1. Generative method:
   - graphical models

2. Discriminative method:
   - SVM

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

# 2 generative models

1. Naïve Bayes classifier
   – Csurka Bray, Dance & Fan, 2004

2. Hierarchical Bayesian text models (pLSA and LDA)
   – Background: Hoffman 2001, Blei, Ng & Jordan, 2004
   – Object categorization: Sivic et al. 2005, Sudderth et al. 2005
   – Natural scene categorization: Fei-Fei et al. 2005

# Demo

- Course website



**A demonstration of bag-of-words classifiers - Microsoft Internet Explorer provided by Insight Broadban**

File    Edit    View    Favorites    Tools    Help

Back    Search    Favorites

Address: http://people.csail.mit.edu/fergus/iccv2005/bagwords.html

Google    Search    100 blocked    Check    AutoLink    AutoF

## Two bag-of-words classifiers

ICCV 2005 short courses on
**Recognizing and Learning Object Categories**

A simple approach to classifying images is to treat them as a collection of regions, describing only their appearance and ignoring their s... have been successfully used in the text community for analyzing documents and are known as "bag-of-words" models, since each docu... distribution over fixed vocabulary(s). Using such a representation, methods such as probabalistic latent semantic analysis (pLSA) [1] (LDA) [2] are able to extract coherent topics within document collections in an unsupervised manner.

Recently, Fei-Fei et al. [3] and Sivic et al. [4] have applied such methods to the visual domain. The demo code implements pLSA, incl... For comparison, a Naive Bayes classifier is also provided which requires labelled training data, unlike pLSA.

The code consists of Matlab scripts (which should run under both Windows and Linux) and a couple of 32-bit Linux binaries for doing... representation. Hence the whole system will need to be run on Linux. The code is for teaching/research purposes only. If you find a b... where csail point mit point edu.

## Download

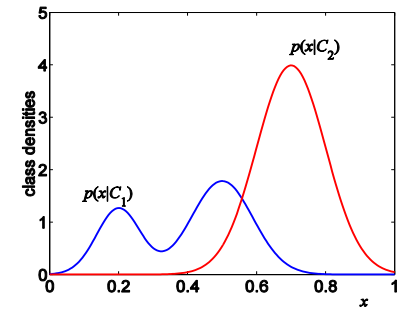Download the code and datasets (32 Mbytes)

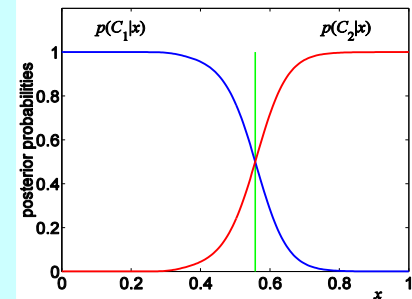## Operation of code

To run the demos:

start    Microsoft Outlook We...    未名空间(mitbbs.co...    A demonstration of b...    ICCV200

# Learning and Recognition
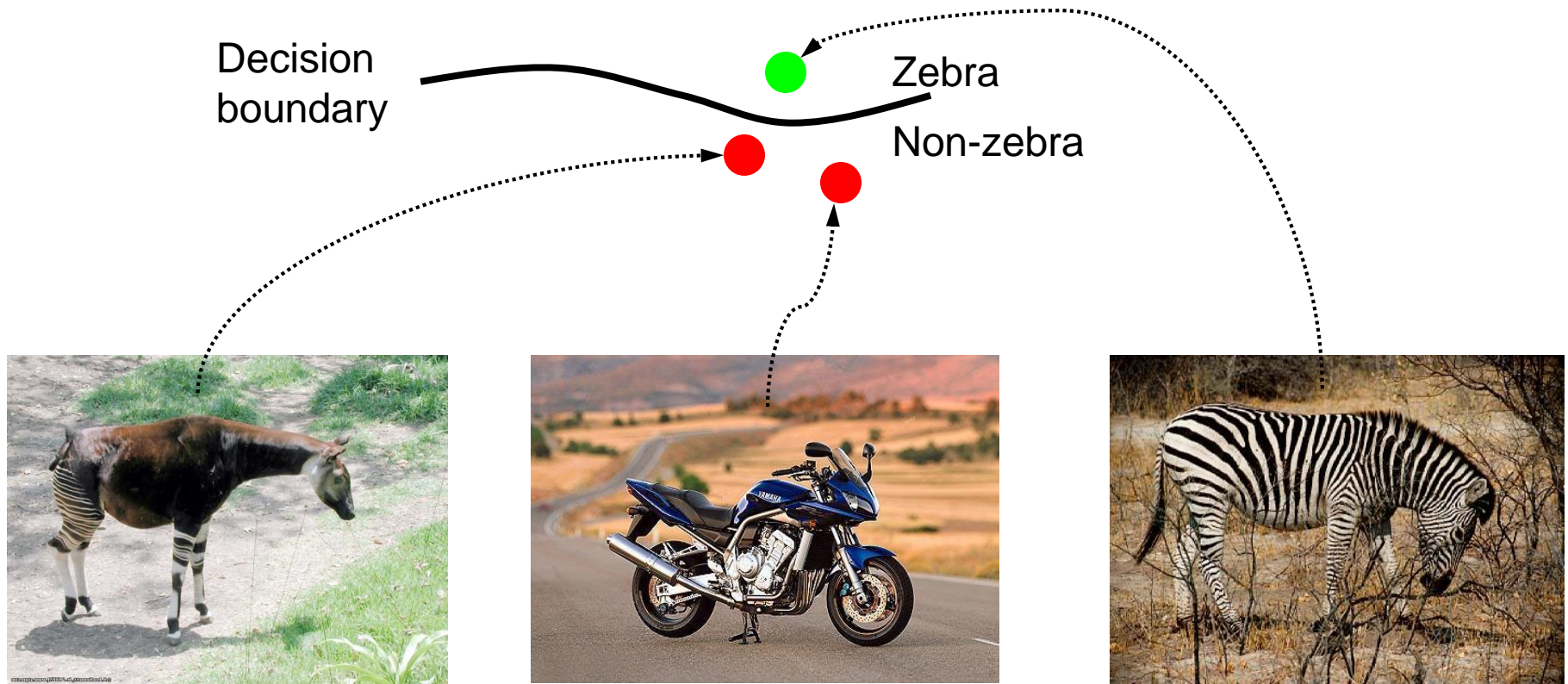
1. Generative method:
   - graphical models

2. Discriminative method:
   - SVM

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

# Discriminative methods based on 'bag of words' representation



Decision boundary

Zebra

Non-zebra

# 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

# Object recognition results

- ETH-80 database
  8 object classes
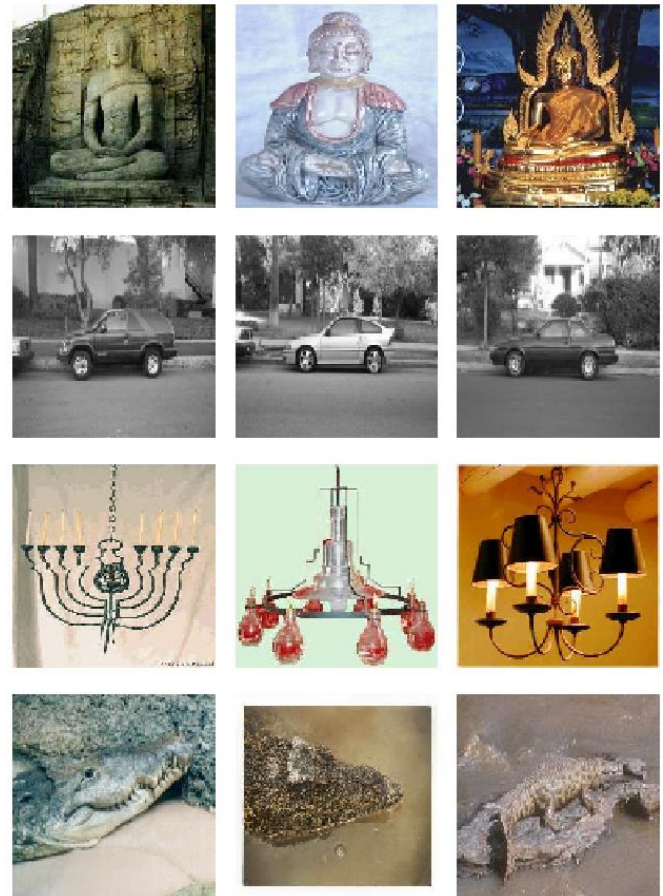  (*Eichhorn and Chapelle 2004*)

- Features:
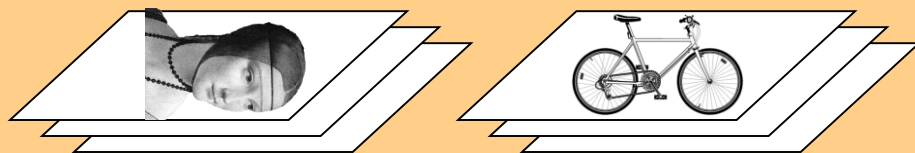  - Harris detector
  - PCA-SIFT descriptor, *d*=10



| Kernel | Complexity | Recognition rate |
|---|---|---|
| Match *[Wallraven et al.]* | $O(dm^2)$ | 84% |
| Bhattacharyya affinity *[Kondor & Jebara]* | $O(dm^3)$ | 85% |
| Pyramid match | $O(dmL)$ | 84% |

# Object recognition results

- Caltech objects database 101 object classes
- Features:
  - SIFT detector
  - PCA-SIFT descriptor, $d$=10
- 30 training images / class
- 43% recognition rate (1% chance performance)
- 0.002 seconds per match

# learning

# recognition
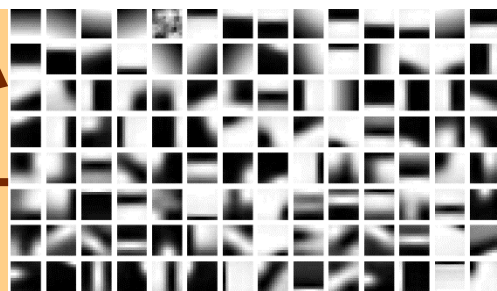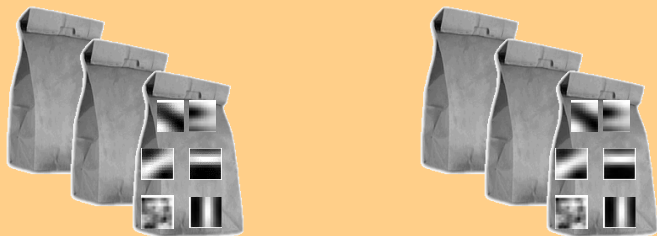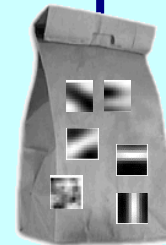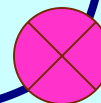
feature detection & representation

**codewords dictionary**

image representation

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

**category decision**
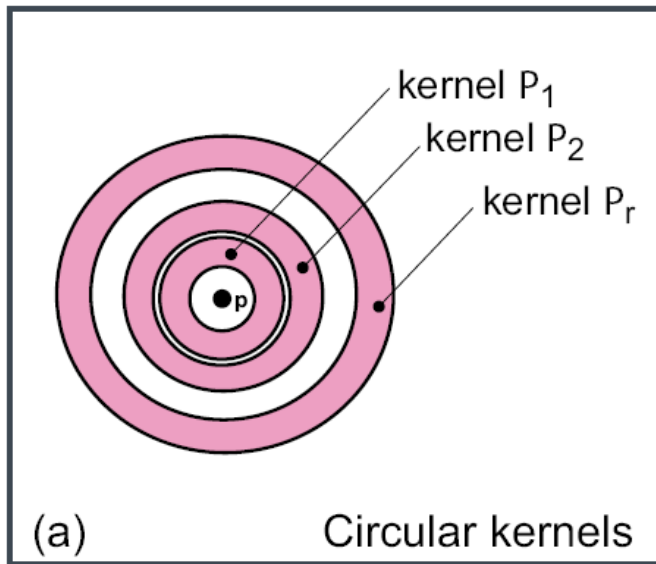
# What about spatial info?

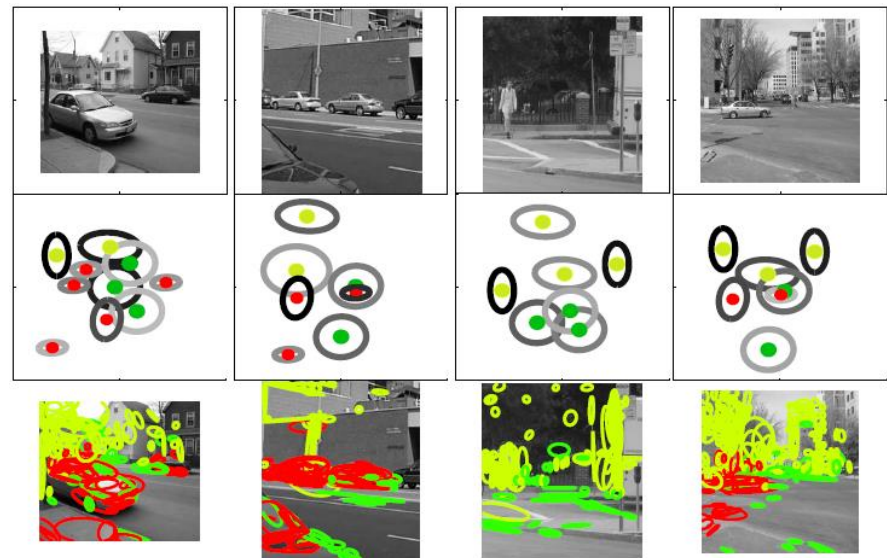- ## Feature level
  - Spatial influence through correlogram features: Savarese, Winn and Criminisi, CVPR 2006



kernel $P_1$
kernel $P_2$
kernel $P_r$

(a)   Circular kernels



frequency

(a)       10       20       30       40   correlaton label

# What about spatial info?

- Feature level

- Generative models
  - Sudderth, Torralba, Freeman & Willsky, 2005, 2006
  - Niebles & Fei-Fei, CVPR 2007

# What about spatial info?

- Feature level

- Generative models
  - Sudderth, Torralba, Freeman & Willsky, 2005, 2006
  - Niebles & Fei-Fei, CVPR 2007

# What about spatial info?
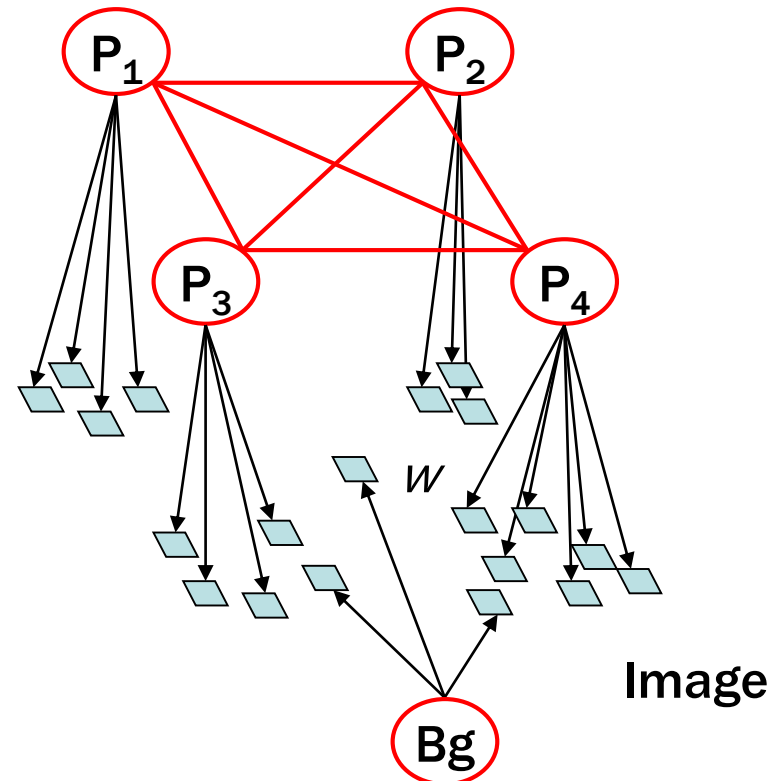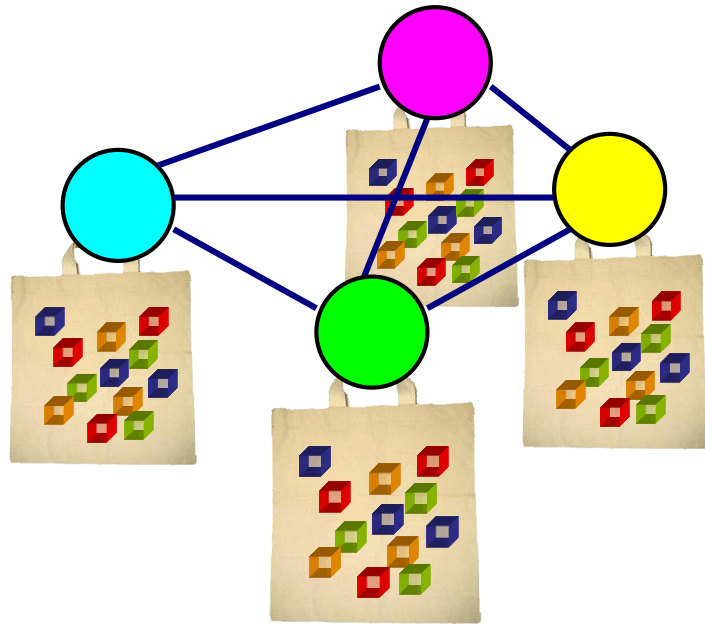
- Feature level
- Generative models
- Discriminative methods
  - Lazebnik, Schmid & Ponce, 2006



level 0                level 1                level 2

# Invariance issues

- Scale and rotation
  - Implicit
  - Detectors and descriptors



Kadir and Brady. 2003

# **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
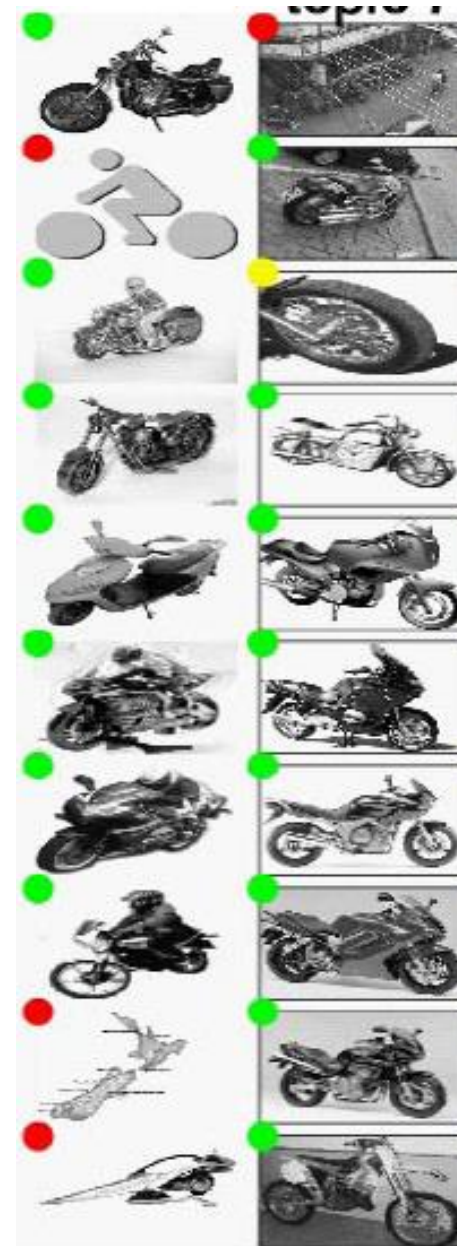
- Translation

- View point (in theory)

  – Codewords: detector and descriptor

  – Theme distributions: different view points

Fergus, Fei-Fei, Perona & Zisserman, 2005

# Model properties

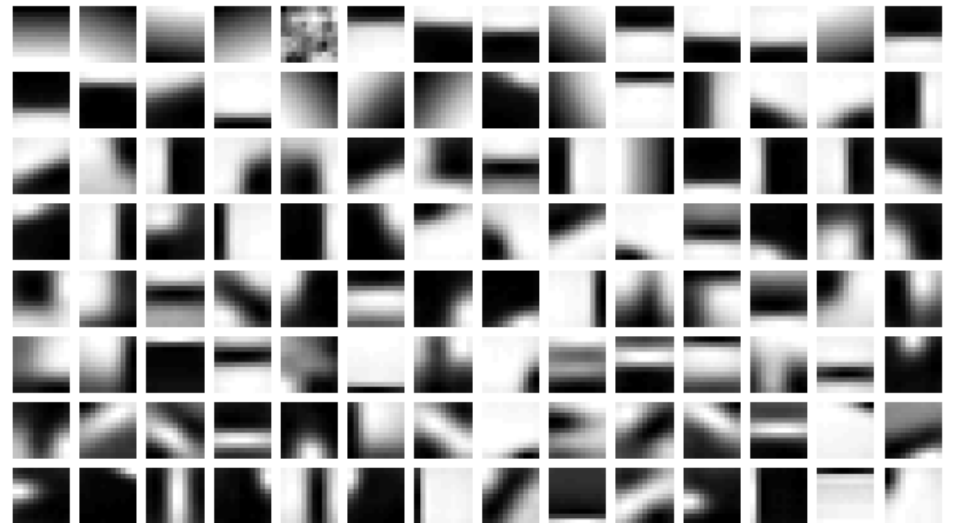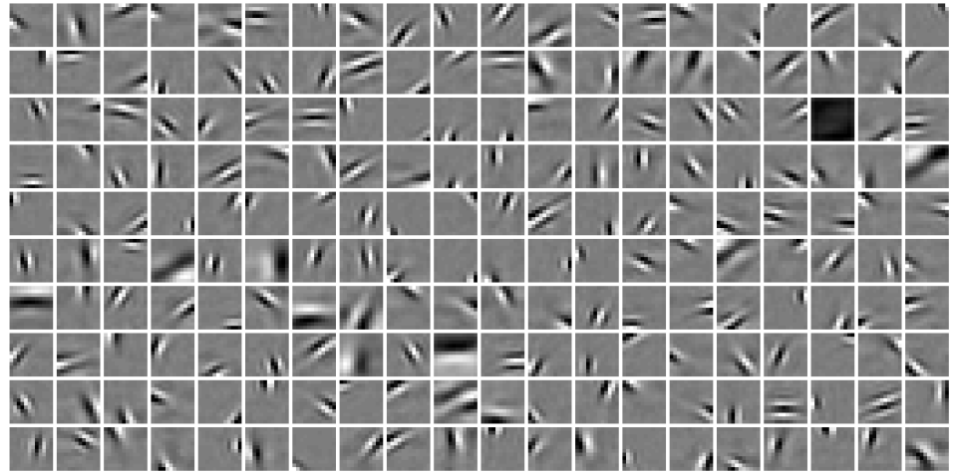- Intuitive
  - Analogy to documents

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 the brain from our eyes. For a long time it was thought that the retinal image was transmitted point by point to visual centers in the brain; the cerebral cortex was a movie screen, so to speak, upon which the image in the eye was projected. Through the discoveries of Hubel and Wiesel we now know that behind the origin of the visual perception in the brain there is a considerably more complicated course of events. By following the visual impulses along their path to the various cell layers of the optical cortex, Hubel and Wiesel have been able to 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**
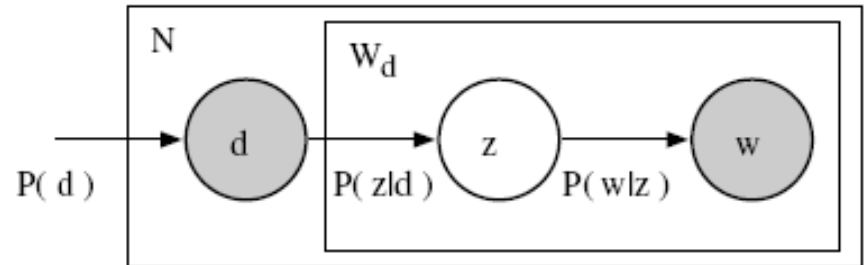
# **Model properties**



- Intuitive
  - Analogy to documents
  - Analogy to human vision

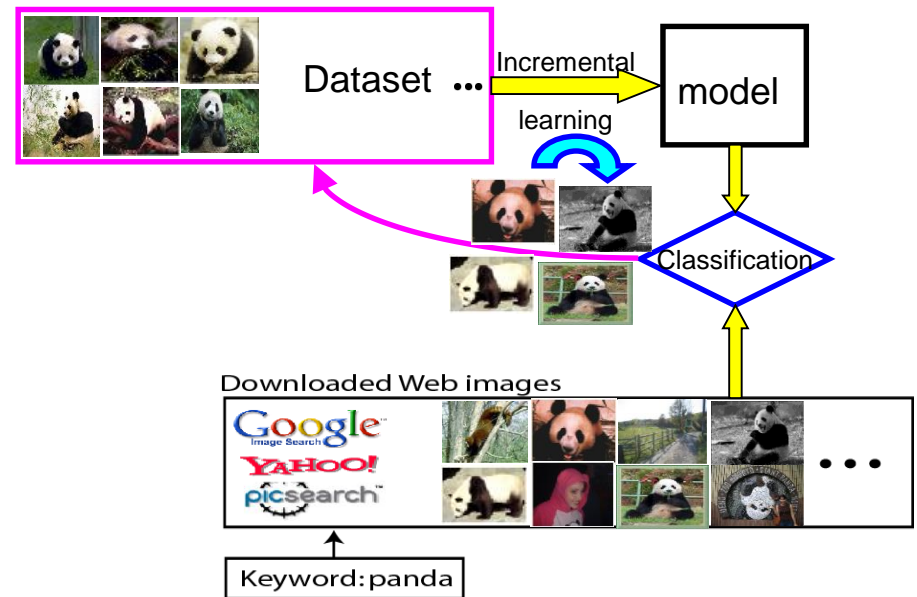Olshausen and Field, 2004, Fei-Fei and Perona, 2005

# **Model properties**



Sivic, Russell, Efros, Freeman, Zisserman, 2005

- Intuitive

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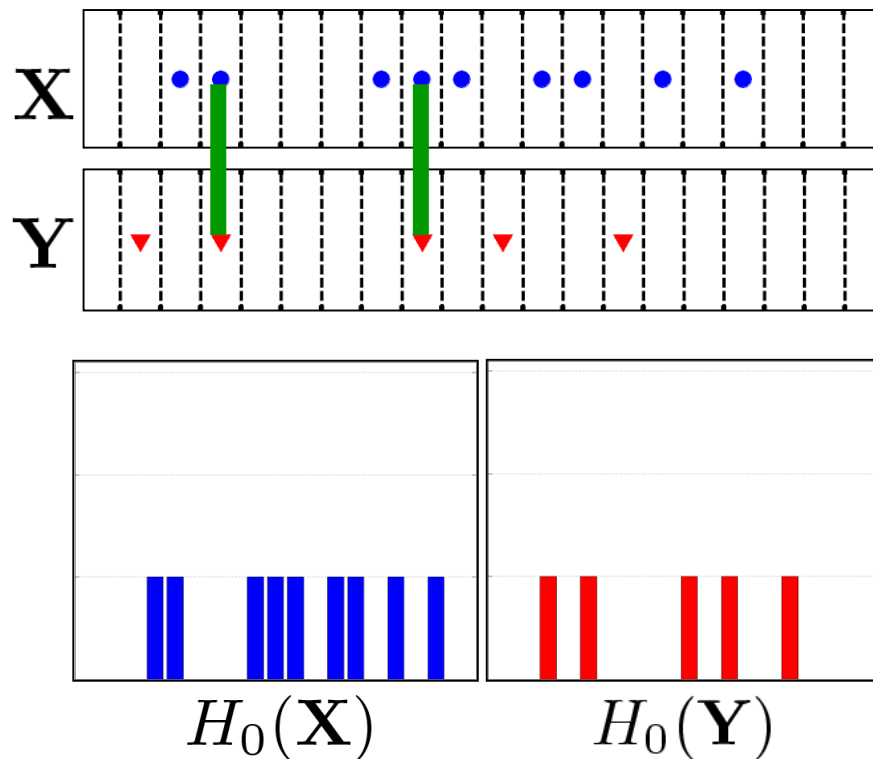  - Prior information
  - Flexibility (e.g. HDP)
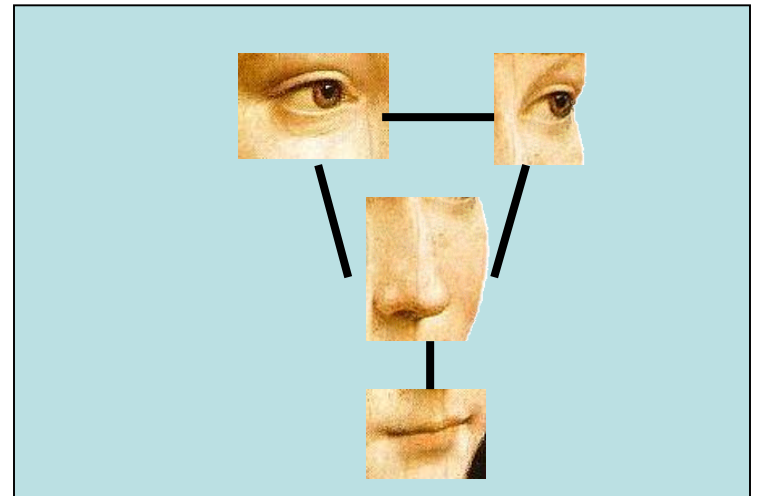

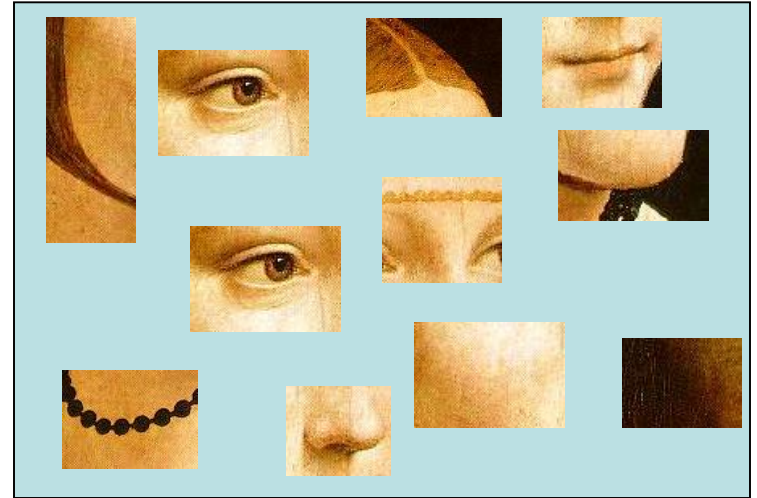
Li, Wang & Fei-Fei, CVPR 2007

# **Model properties**



- Intuitive

- generative models

- Discriminative method
  - Computationally efficient



$$H_0(\mathbf{X}) \qquad H_0(\mathbf{Y})$$

Grauman et al. CVPR 2005

# 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