

Image Features and Categorization



Computer Vision

Jia-Bin Huang, Virginia Tech

Administrative stuffs

- Final project
 - Got your proposals! Thanks!
 - Will reply with feedbacks this week.
- HW 4
 - Due 11:59pm on Wed, November 8
- Happy Halloween!

Review: Interpreting Intensity

- **Light and color**

- What an image records

- **Filtering in spatial domain**

- Filtering = weighted sum of neighboring pixels
 - Smoothing, sharpening, measuring texture

- **Filtering in frequency domain**

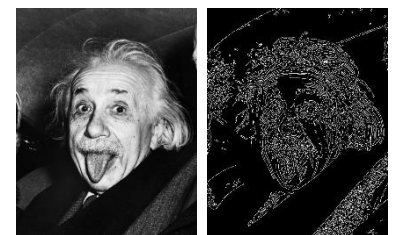
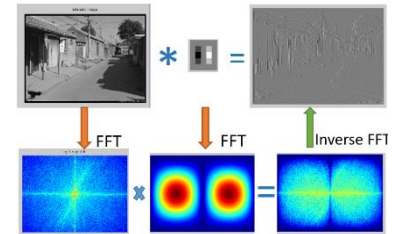
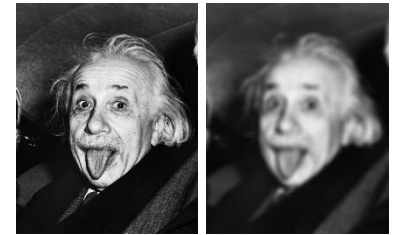
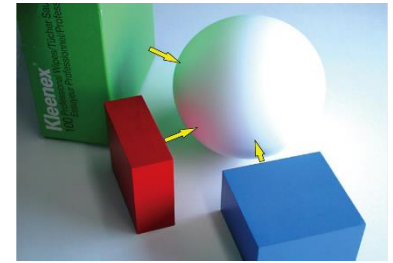
- Filtering = change frequency of the input image
 - Denoising, sampling, image compression

- **Image pyramid and template matching**

- Filtering = a way to find a template
 - Image pyramids for coarse-to-fine search and multi-scale detection

- **Edge detection**

- Canny edge = smooth -> derivative -> thin -> threshold -> link
 - Finding straight lines, binary image analysis



Review: Correspondence and Alignment

• Interest points

- Find *distinct* and *repeatable* points in images
- Harris-> corners, DoG -> blobs
- SIFT -> feature descriptor

• Feature tracking and optical flow

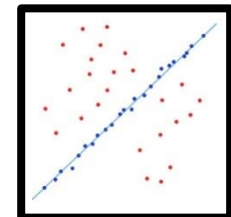
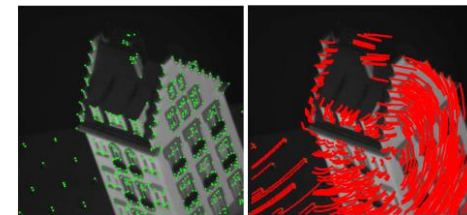
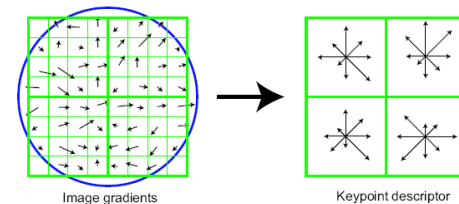
- Find motion of a keypoint/pixel over time
- Lucas-Kanade:
 - brightness consistency, small motion, spatial coherence
- Handle large motion:
 - iterative update + pyramid search

• Fitting and alignment

- find the transformation parameters that best align matched points

• Object instance recognition

- Keypoint-based object instance recognition and search



Review: Perspective and 3D Geometry

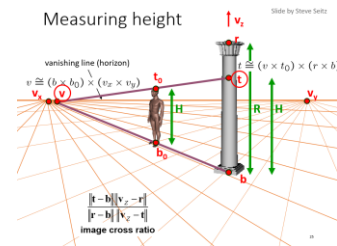
- **Projective geometry and camera models**

- What's the mapping between image and world coordinates?

$$\mathbf{x} = \mathbf{K}[\mathbf{R} \quad \mathbf{t}] \mathbf{X}$$

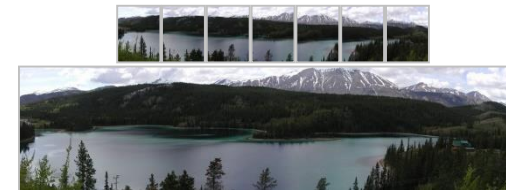
- **Single view metrology and camera calibration**

- How can we measure the size of 3D objects in an image?
- How can we estimate the camera parameters?



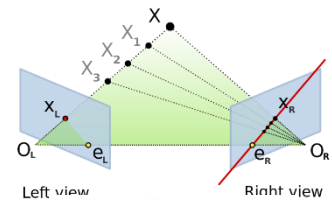
- **Photo stitching**

- What's the mapping from two images taken without camera translation?



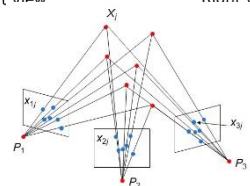
- **Epipolar Geometry and Stereo Vision**

- What's the mapping from two images taken with camera translation?



- **Structure from motion**

- How can we recover 3D points from multiple images?



Review: Perspective and 3D Geometry

• Grouping and Segmentation

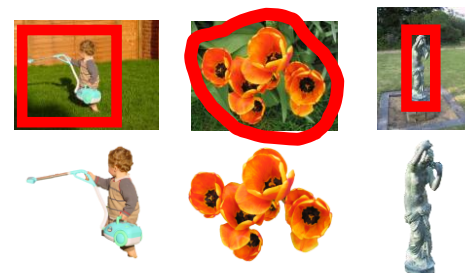
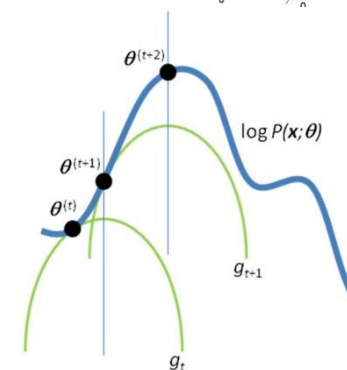
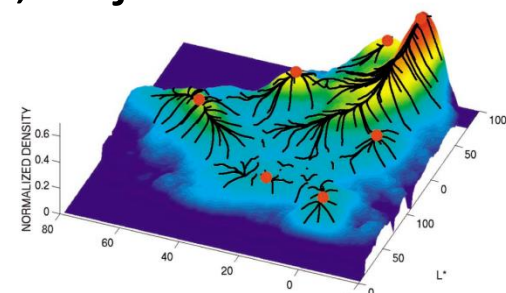
- How do we group pixels into meaningful regions?
- Use of segmentation: efficiency, better features, object region proposal, wanted the segmented object

• EM Algorithm, Mixture of Gaussians

- How do we deal with missing data?
- Maximum likelihood estimation
- Probabilistic inference
- Expectation-Maximization algorithm

• MRFs and Graph Cut

- How do we encode pixel dependencies?
- Markov Random Fields
- Graph Cuts



Recognition and Learning

- Image Features and Categorization
- Convolutional Neural Networks
- Object Detection
- Part and Pixel Labeling
- Action Recognition
- Vision and Language

AIPOLY

Mashable

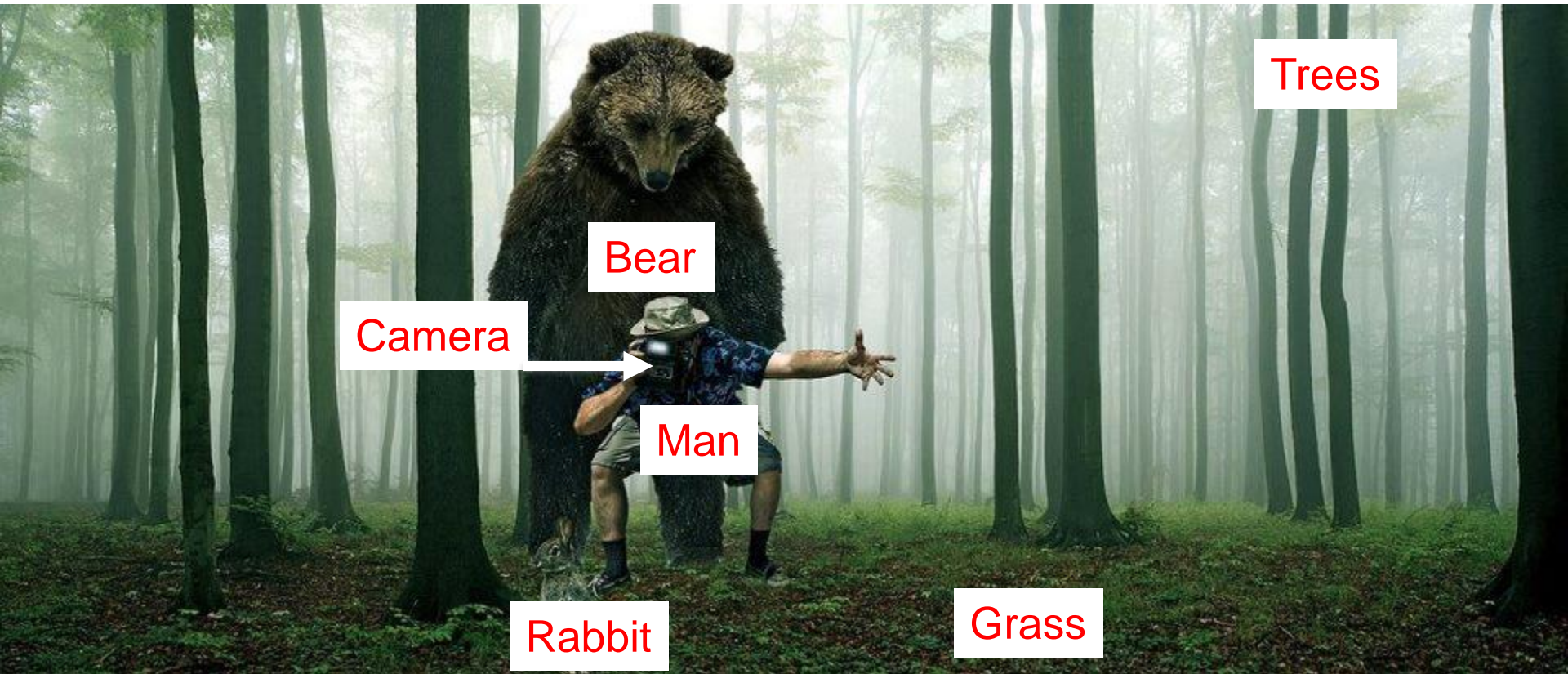
patio
Patio.



Today: Image features and categorization

- General concepts of categorization
 - Why? What? How?
- Image features
 - Color, texture, gradient, shape, interest points
 - Histograms, feature encoding, and pooling
 - CNN as feature
- Image and region categorization

What do you see in this image?



Forest

Describe, predict, or interact with the object based on visual cues



Is it **dangerous**?

How **fast** does it run?

Is it **alive**?

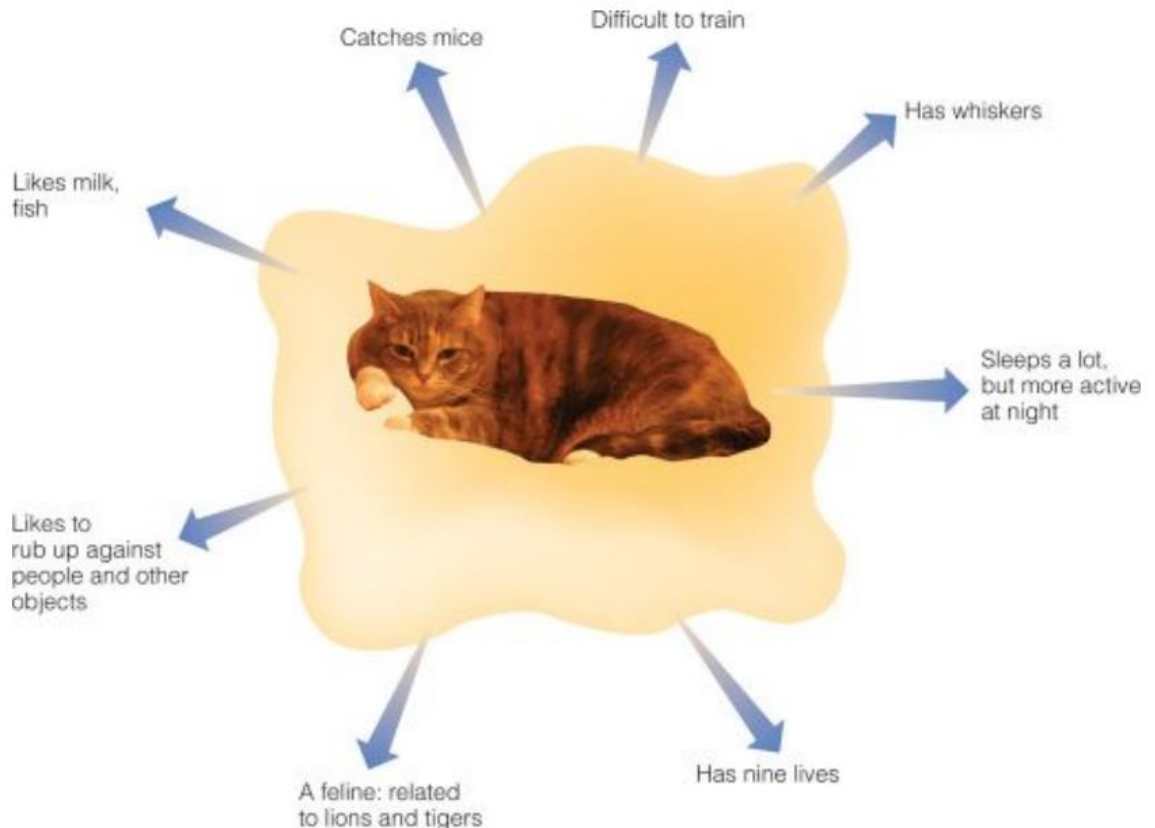
Is it **soft**?

Does it have a **tail**?

Can I **poke with it**?

Why do we care about categories?

- From an object's category, we can make predictions about its behavior in the future, beyond of what is immediately perceived.
- Pointers to knowledge
 - Help to understand individual cases not previously encountered
- Communication



Theory of categorization

How do we determine if something is a member of a particular category?

- Definitional approach
- Prototype approach
- Exemplar approach

Definitional approach: classical view of categories

- Plato & Aristotle

- Categories are defined by a list of properties shared by all elements in a category
- Category membership is binary
- Every member in the category is equal

[The Categories \(Aristotle\)](#)



Aristotle by Francesco Hayez

Prototype or sum of exemplars ?

Prototype Model

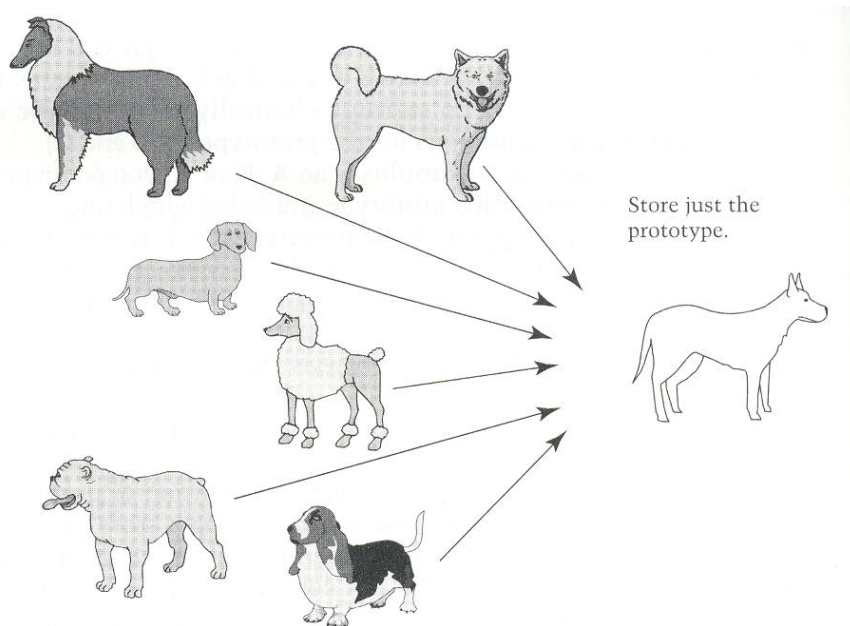


Figure 7.3. Schematic of the prototype model. Although many exemplars are seen, only the prototype is stored. The prototype is updated continually to incorporate more experience with new exemplars.

Category judgments are made by comparing a new exemplar to the prototype.

Exemplars Model

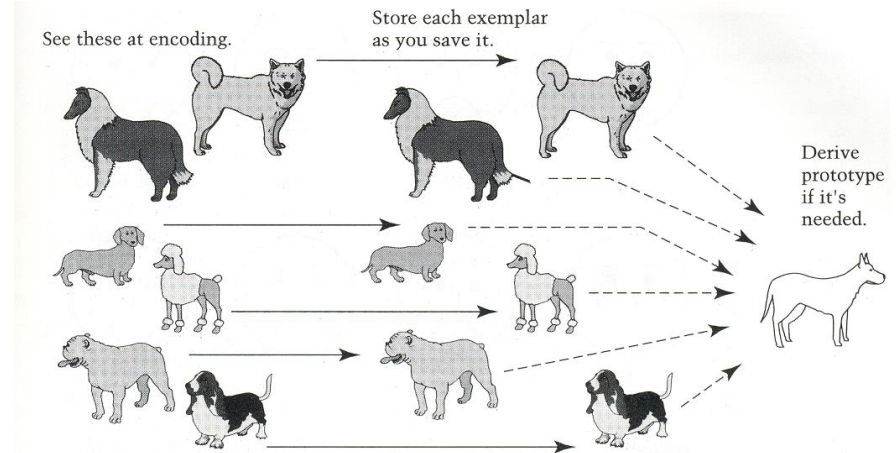


Figure 7.4. Schematic of the exemplar model. As each exemplar is seen, it is encoded into memory. A prototype is abstracted only when it is needed, for example, when a new exemplar must be categorized.

Category judgments are made by comparing a new exemplar to all the old exemplars of a category or to the exemplar that is the most appropriate

Levels of categorization [Rosch 70s]



Definition of Basic Level:

- **Similar shape:** Basic level categories are the highest-level category for which their members have similar shapes.
- **Similar motor interactions:** ... for which people interact with its members using similar motor sequences.
- **Common attributes:** ... there are a significant number of attributes in common between pairs of members.

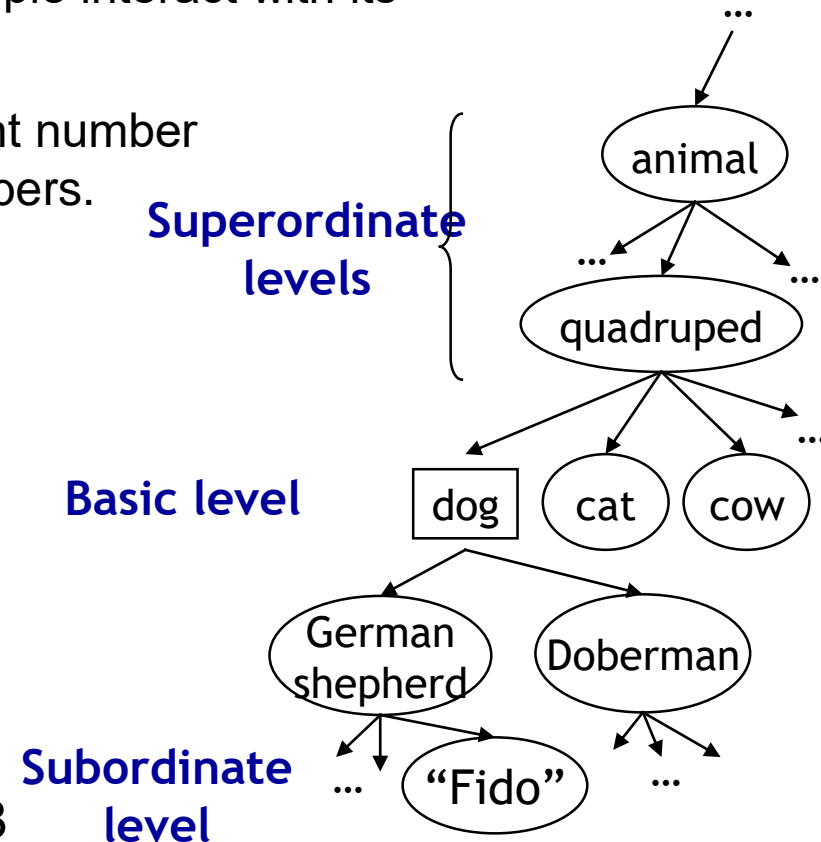
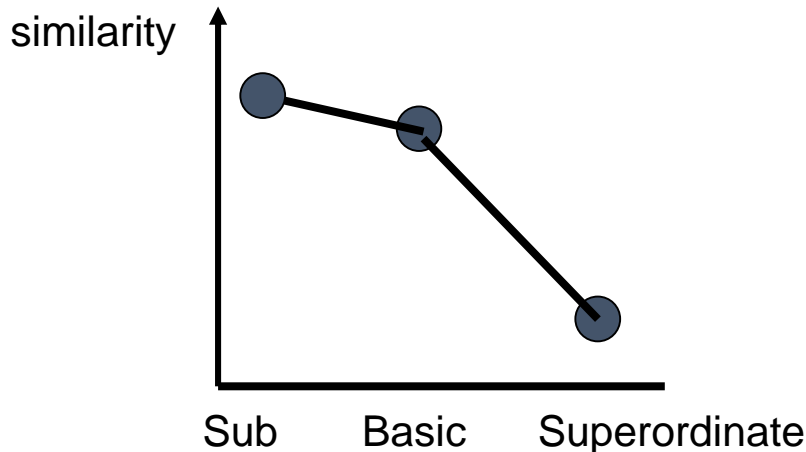


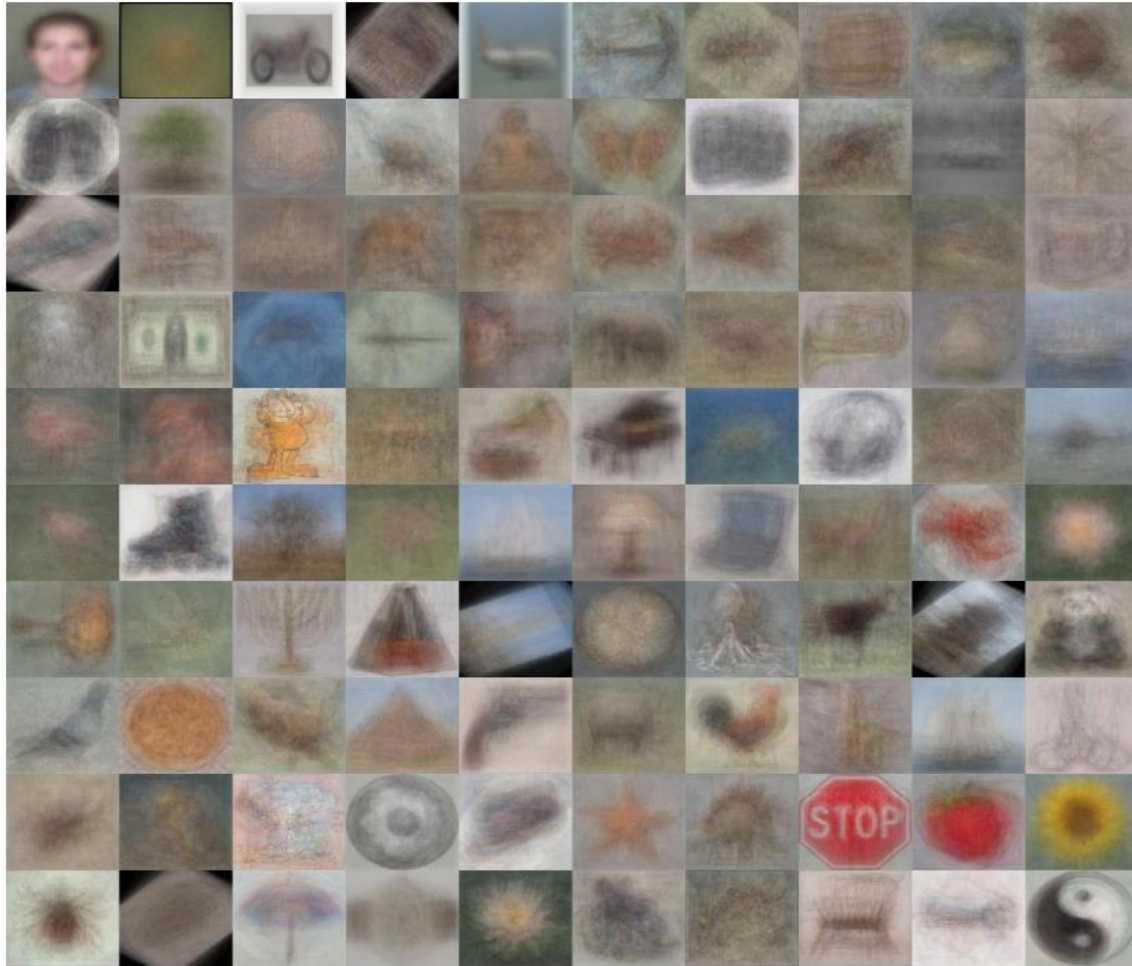
Image categorization

- Cat vs Dog



Image categorization

- Object recognition

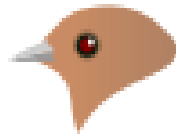


Caltech 101 Average Object Images

Image categorization



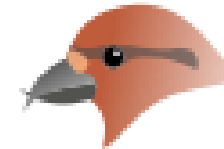
Generalist



Insect catching



Grain eating



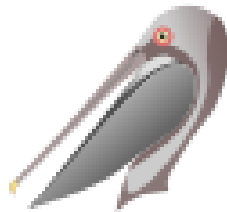
Coniferous-seed eating



Nectar feeding



Chiseling



Dip netting



Surface skimming



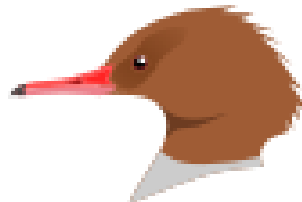
Scything



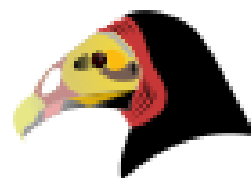
Probing



Aerial fishing



Pursuit fishing



Scavenging



Raptorial



Filter feeding

Image categorization

- Place recognition



Places Database [[Zhou et al. NIPS 2014](#)]

Image categorization

- Visual font recognition

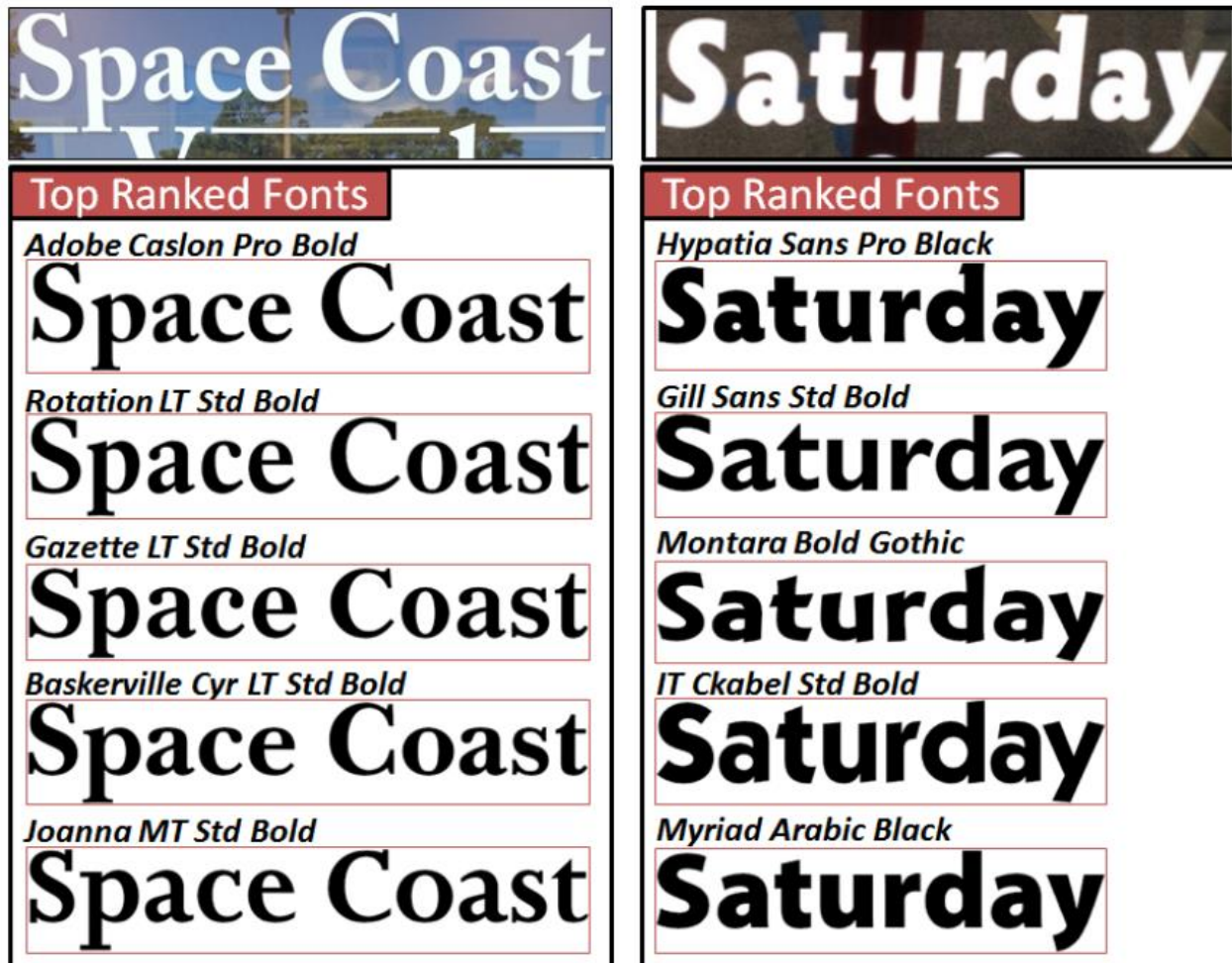


Image categorization

- Dating historical photos



1940



1953



1966



1977

[[Palermo et al. ECCV 2012](#)]

Image categorization

- Image style recognition



HDR



Macro



Baroque



Roccoco



Vintage



Noir



Northern Renaissance



Cubism



Minimal



Hazy



Impressionism



Post-Impressionism



Long Exposure



Romantic



Abs. Expressionism



Color Field Painting

Flickr Style: 80K images covering 20 styles.

Wikipaintings: 85K images for 25 art genres.

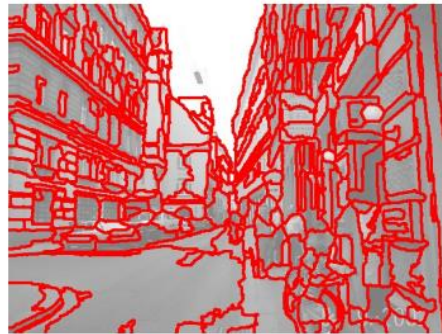
[[Karayev et al. BMVC 2014](#)]

Region categorization

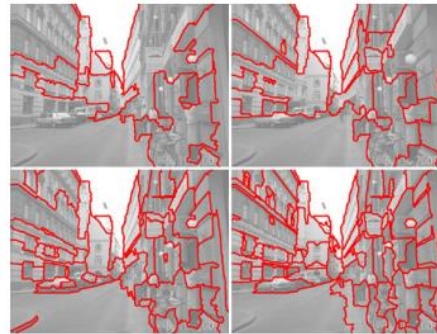
- Layout prediction



Input



Superpixels



Multiple Segmentations



Surface Layout

Assign regions to orientation

Geometric context [[Hoiem et al. IJCV 2007](#)]

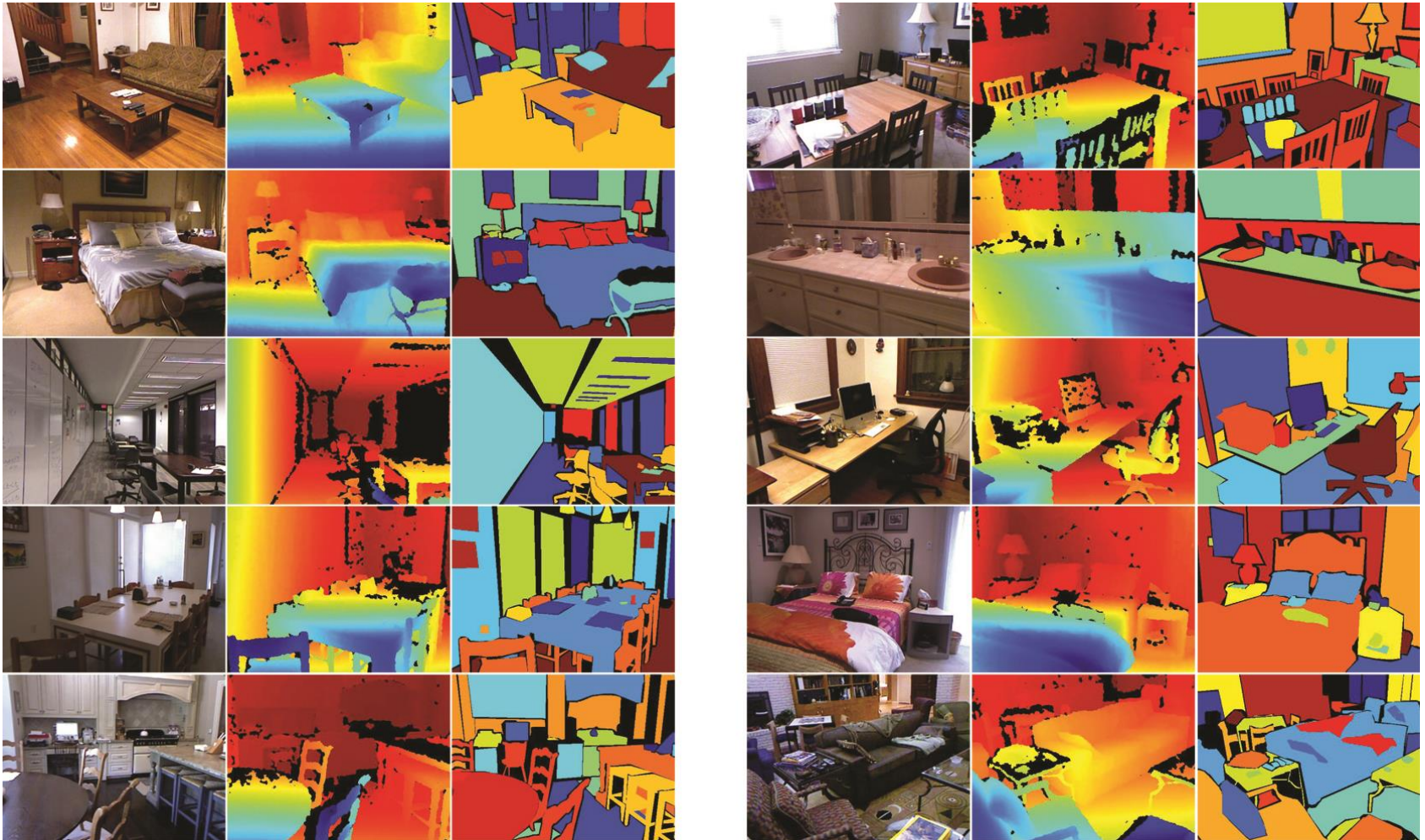


Assign regions to depth

Make3D [[Saxena et al. PAMI 2008](#)]

Region categorization

- Semantic segmentation from RGBD images



[Silberman et al. ECCV 2012]

Region categorization

- Material recognition



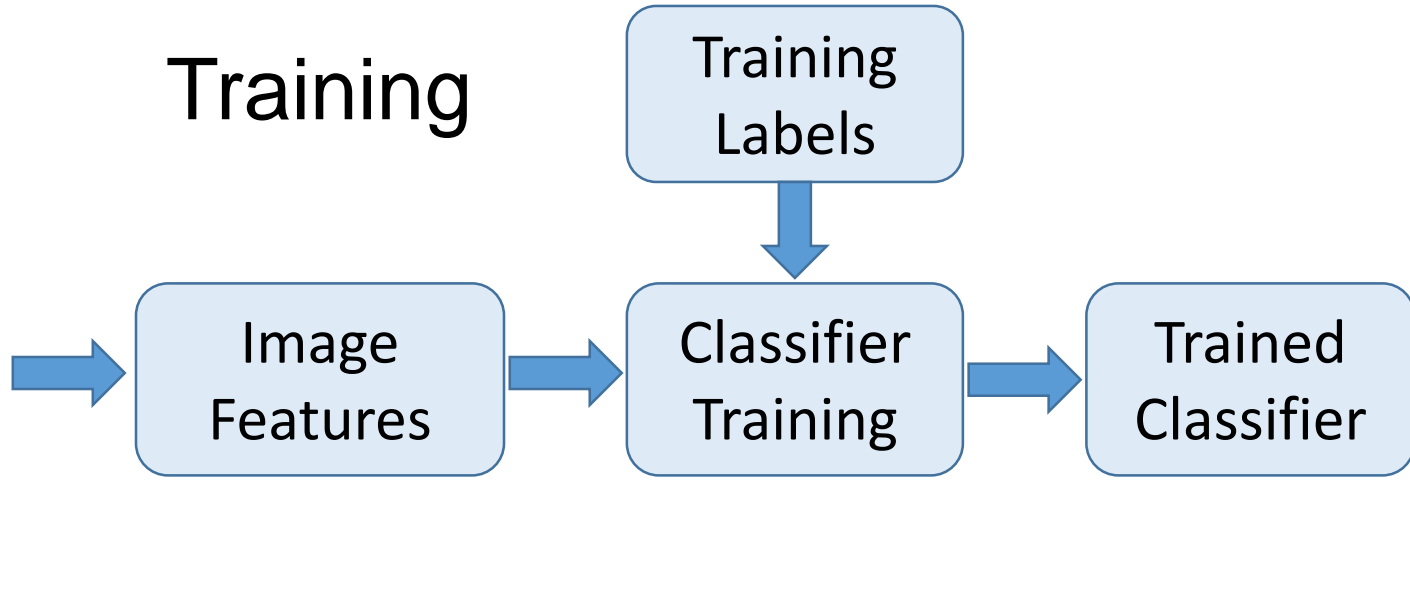
[[Bell et al. CVPR 2015](#)]

Training phase

Training
Images



Training

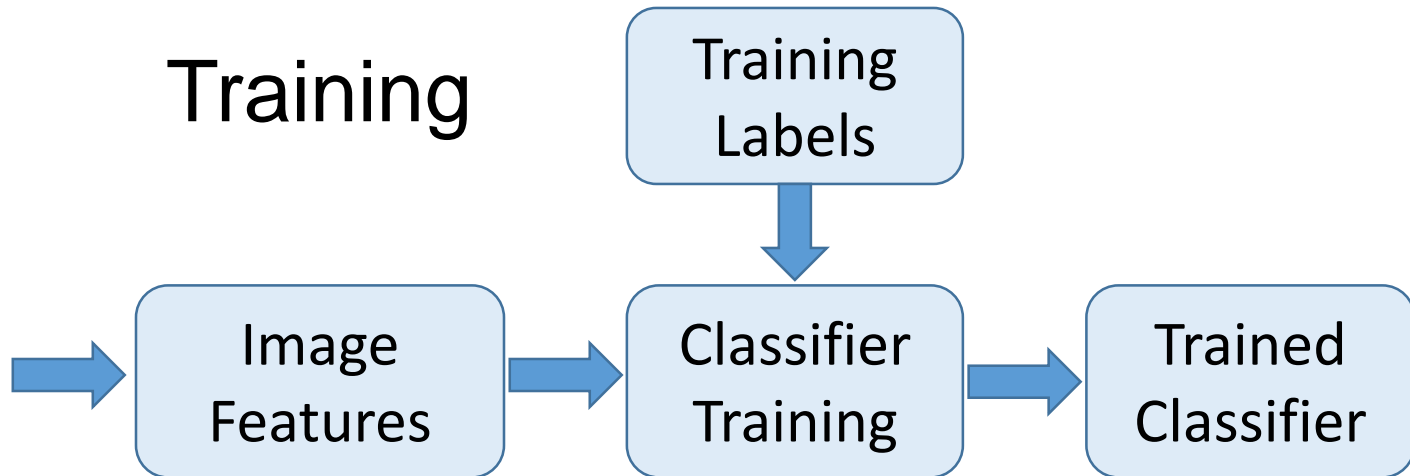


Testing phase

Training Images



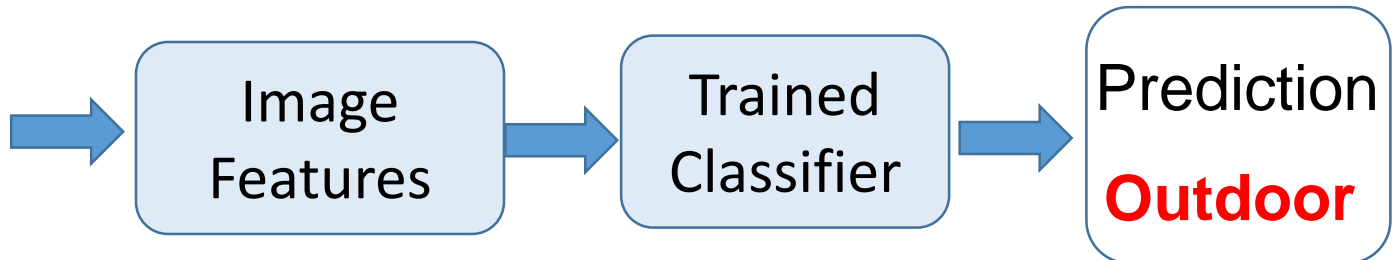
Training



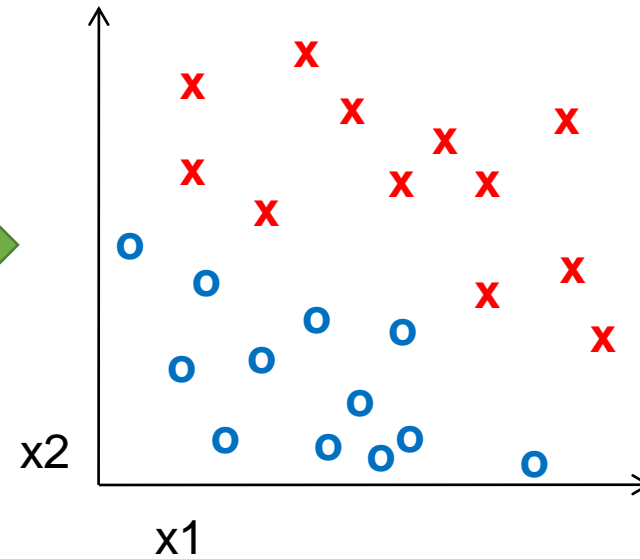
Testing



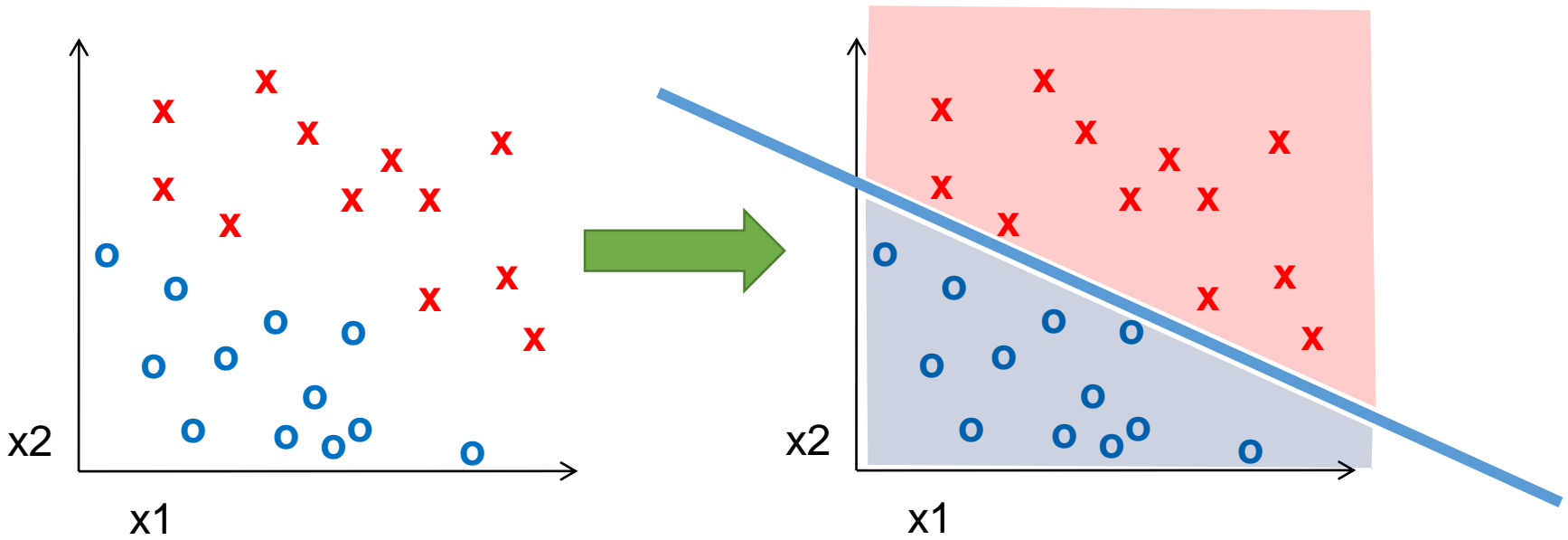
Test Image



- **Image features:** map images to feature space



- **Classifiers:** map feature space to label space



Different types of classification

- **Exemplar-based:** transfer category labels from examples with most similar features
 - What similarity function? What parameters?
- **Linear classifier:** confidence in positive label is a weighted sum of features
 - What are the weights?
- **Non-linear classifier:** predictions based on more complex function of features
 - What form does the classifier take? Parameters?
- **Generative classifier:** assign to the label that best explains the features (makes features most likely)
 - What is the probability function and its parameters?

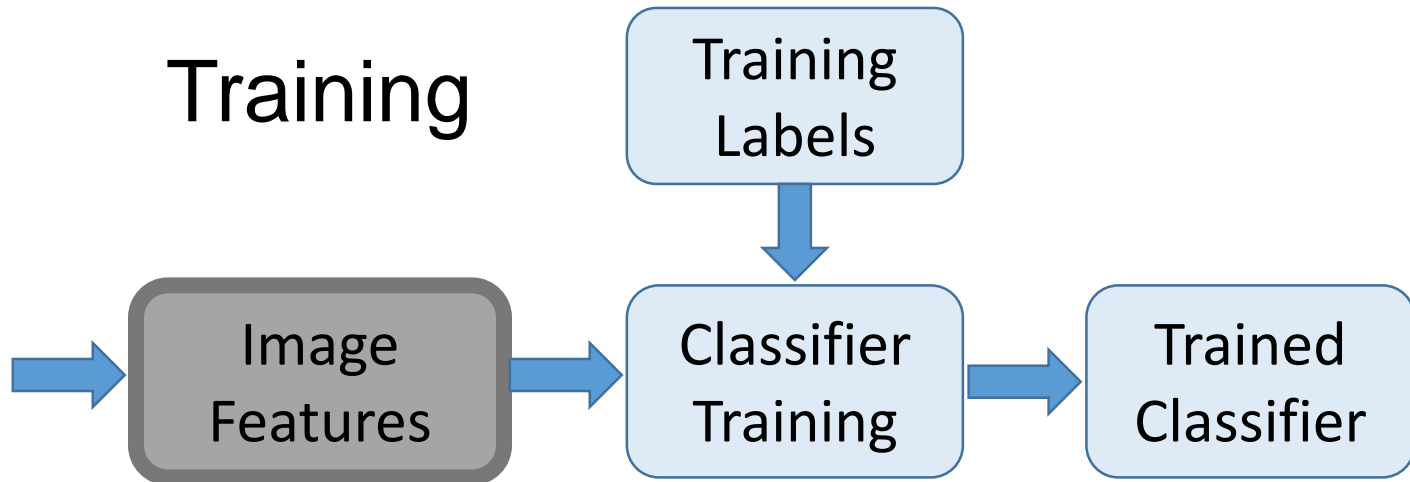
Note: You can always fully design the classifier by hand, but usually this is too difficult. Typical solution: learn from training examples.

Testing phase

Training Images



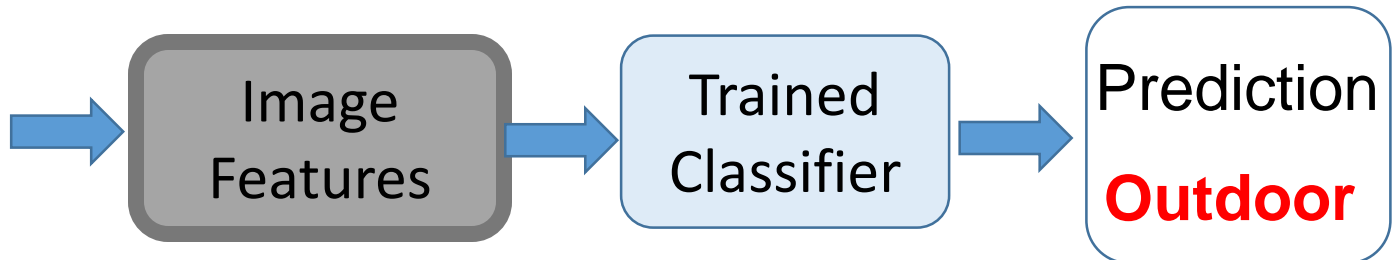
Training



Testing



Test Image



Q: What are good features for...

- recognizing a beach?



Q: What are good features for...

- recognizing cloth fabric?



Q: What are good features for...

- recognizing a mug?



What are the right features?

Depend on what you want to know!

- Object: shape
 - Local shape info, shading, shadows, texture
- Scene : geometric layout
 - linear perspective, gradients, line segments
- Material properties: albedo, feel, hardness
 - Color, texture
- Action: motion
 - Optical flow, tracked points

General principles of representation

- **Coverage**

- Ensure that all relevant info is captured

- **Concision**

- Minimize number of features without sacrificing coverage

- **Directness**

- Ideal features are independently useful for prediction

Image representations

- Templates
 - Intensity, gradients, etc.
- Histograms
 - Color, texture, SIFT descriptors, etc.
- Average of features



Image
Intensity



Gradient
template

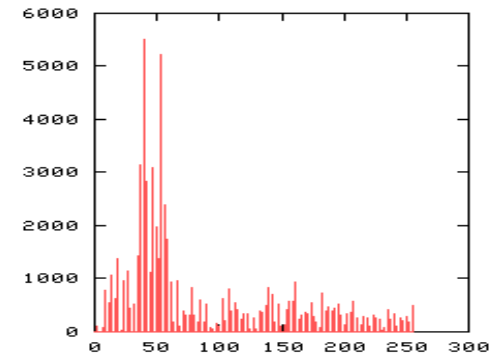
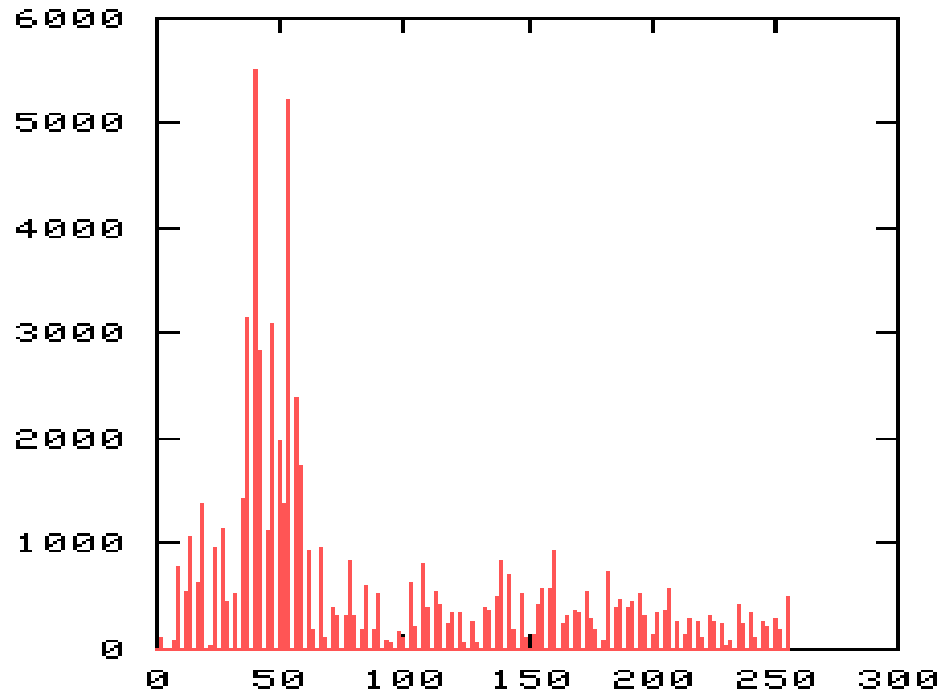


Image representations: histograms



Global histogram

- Represent distribution of features
 - Color, texture, depth, ...

Image representations: histograms

- Data samples in 2D

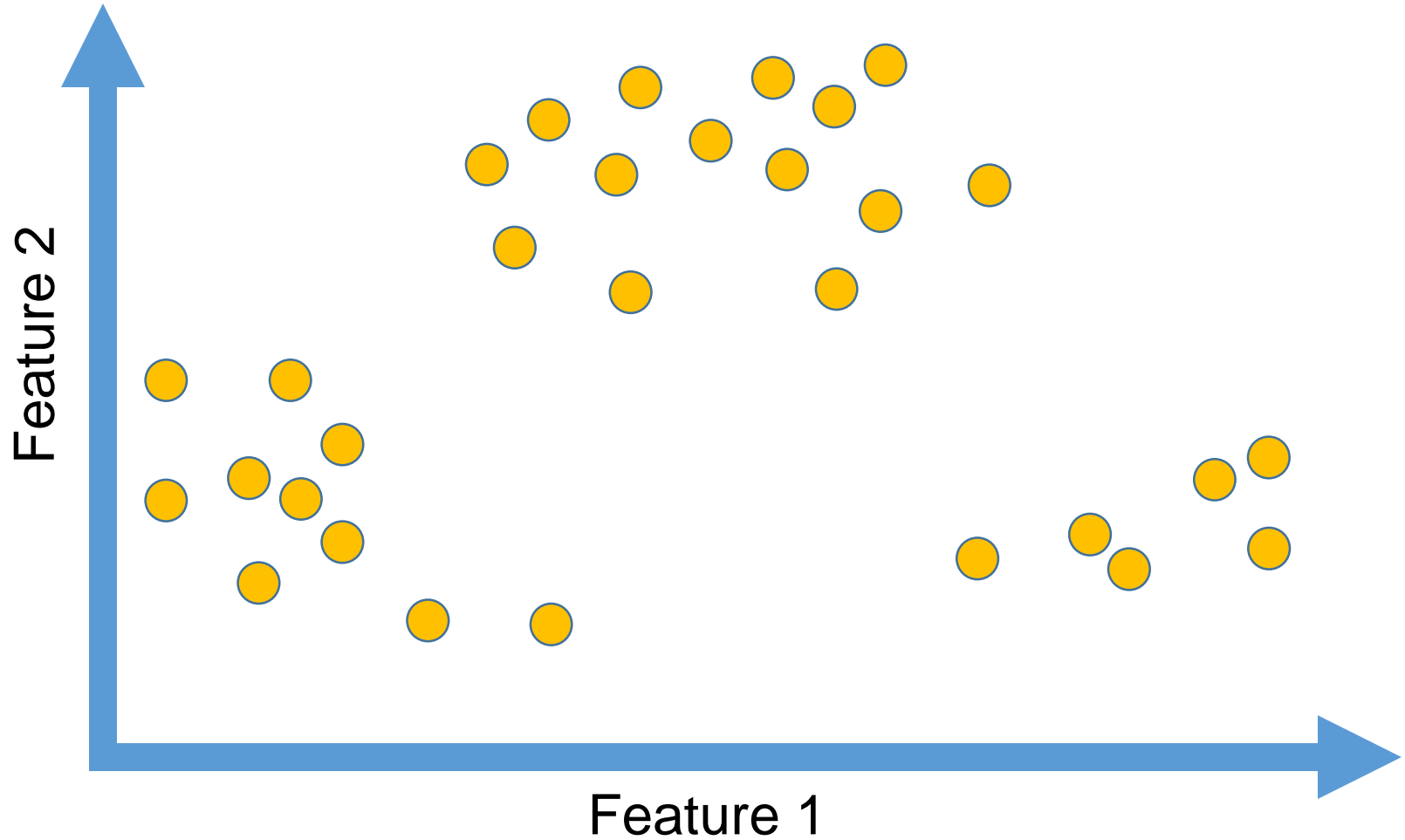


Image representations: histograms

- Probability or count of data in each bin
- Marginal histogram on feature 1

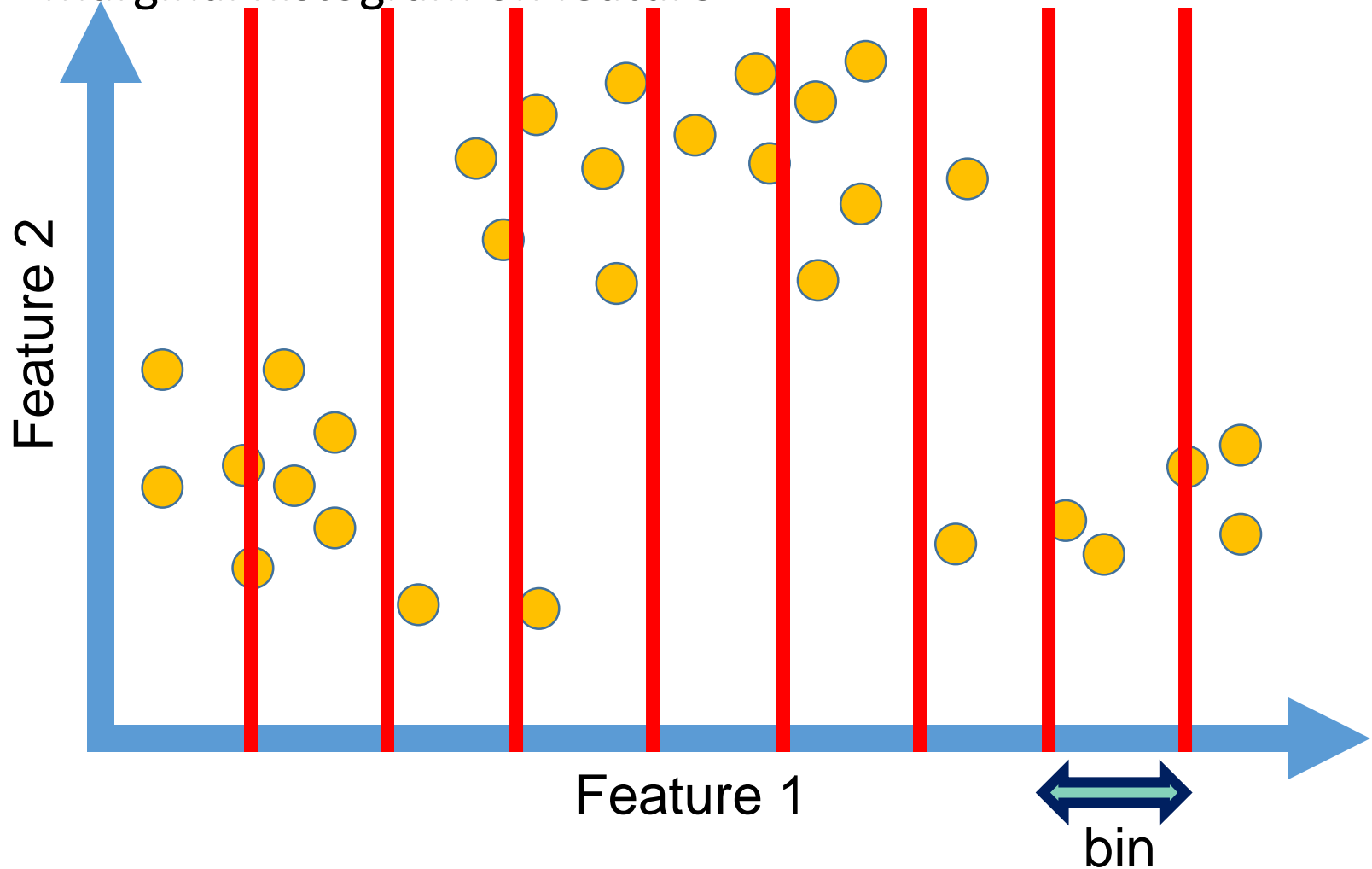


Image representations: histograms

- Marginal histogram on feature 2

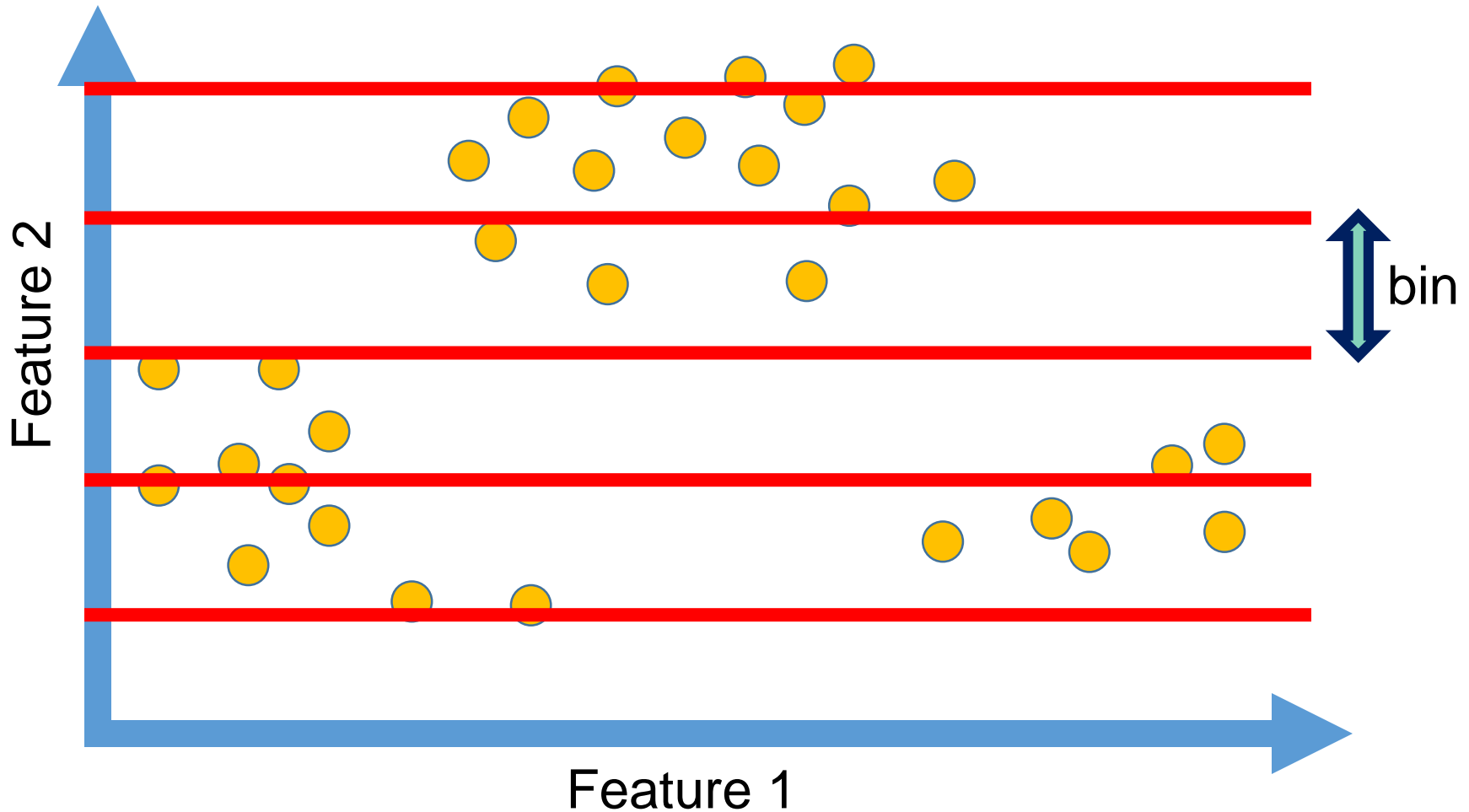
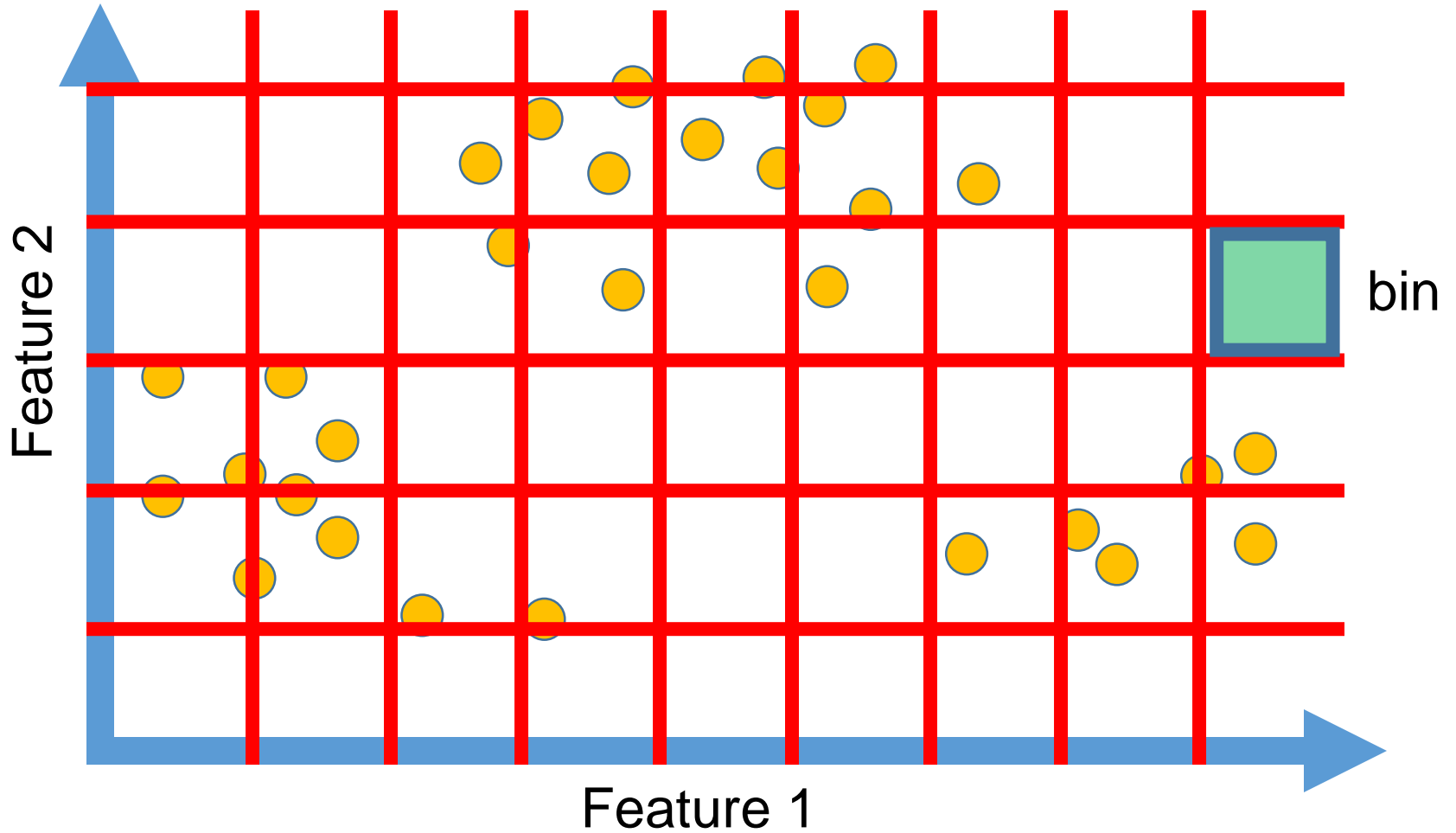
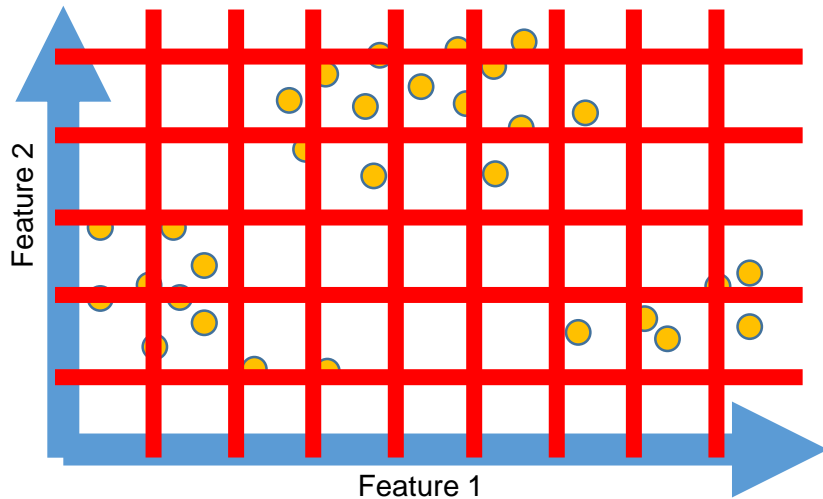


Image representations: histograms

- Joint histogram

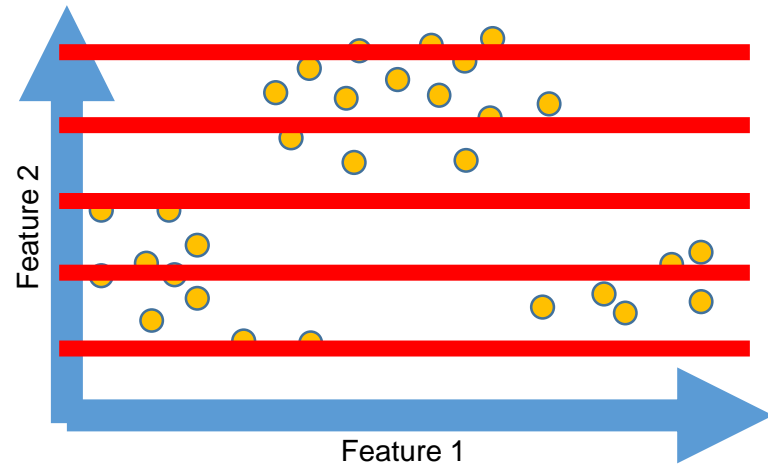
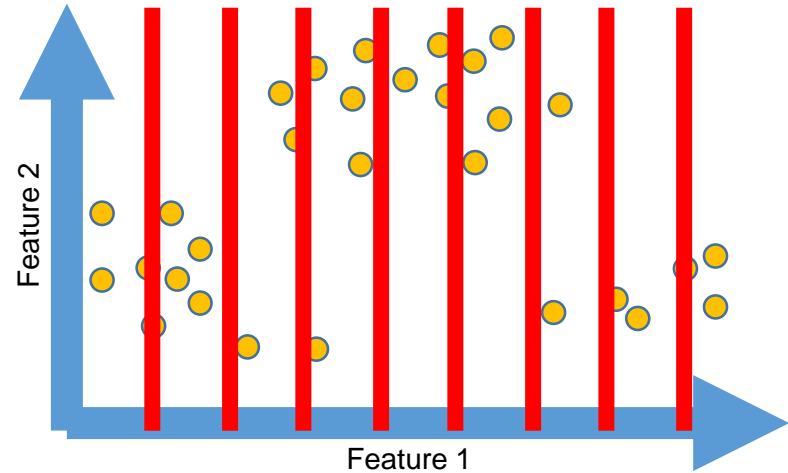


Modeling multi-dimensional data



Joint histogram

- Requires lots of data
- Loss of resolution to avoid empty bins

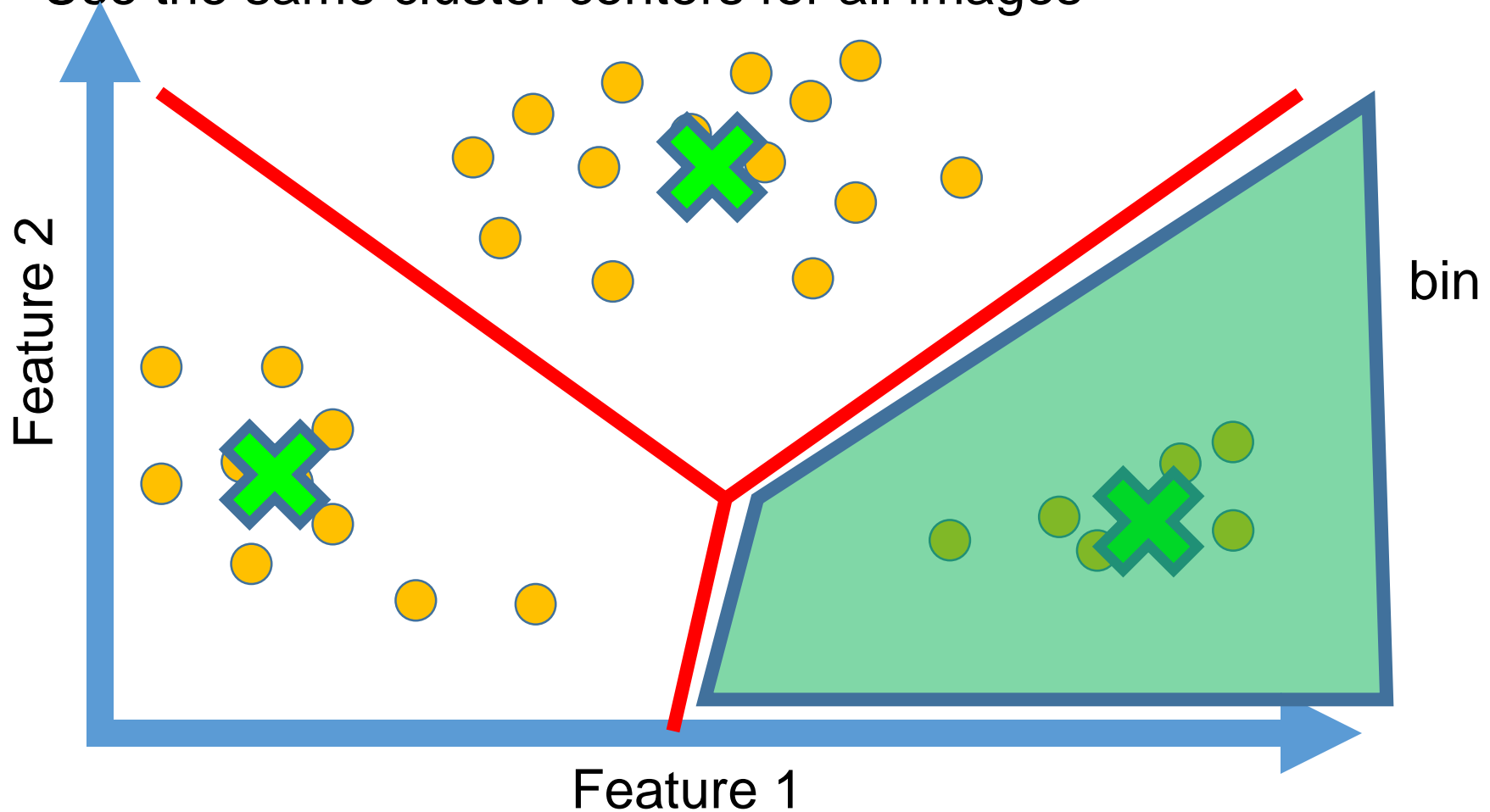


Marginal histogram

- Requires independent features
- More data/bin than joint histogram

Modeling multi-dimensional data

- Clustering
- Use the same cluster centers for all images



Computing histogram distance

- Histogram intersection

$$\text{histint}(h_i, h_j) = 1 - \sum_{m=1}^K \min(h_i(m), h_j(m))$$

- Chi-squared Histogram matching distance

$$\chi^2(h_i, h_j) = \frac{1}{2} \sum_{m=1}^K \frac{[h_i(m) - h_j(m)]^2}{h_i(m) + h_j(m)}$$

- Earth mover's distance
(Cross-bin similarity measure)
 - minimal cost paid to transform one distribution into the other

Histograms: implementation issues

- Quantization

- Grids: fast but applicable only with few dimensions
- Clustering: slower but can quantize data in higher dimensions



Few Bins

Need less data

Coarser representation

Many Bins

Need more data

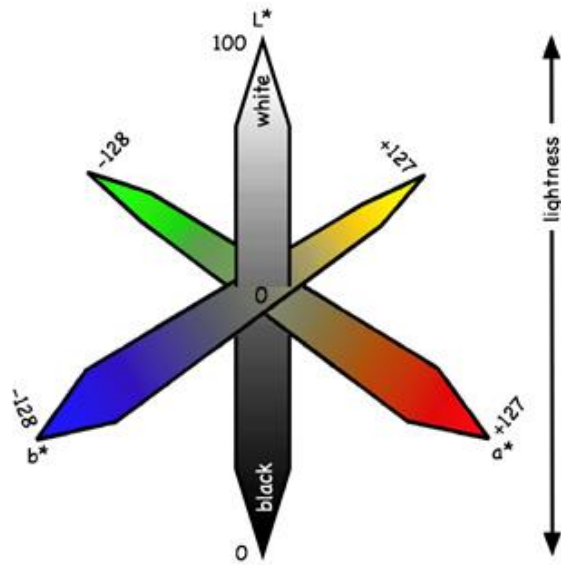
Finer representation

- Matching

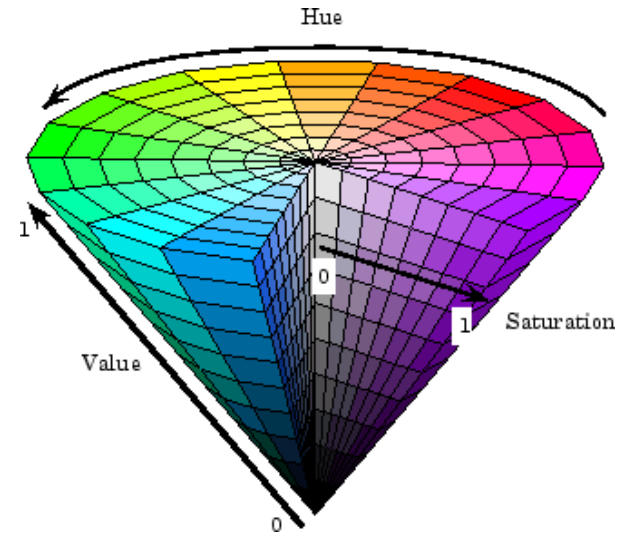
- Histogram intersection or Euclidean may be faster
- Chi-squared often works better
- Earth mover's distance is good for when nearby bins represent similar values

What kind of things do we compute histograms of?

- Color

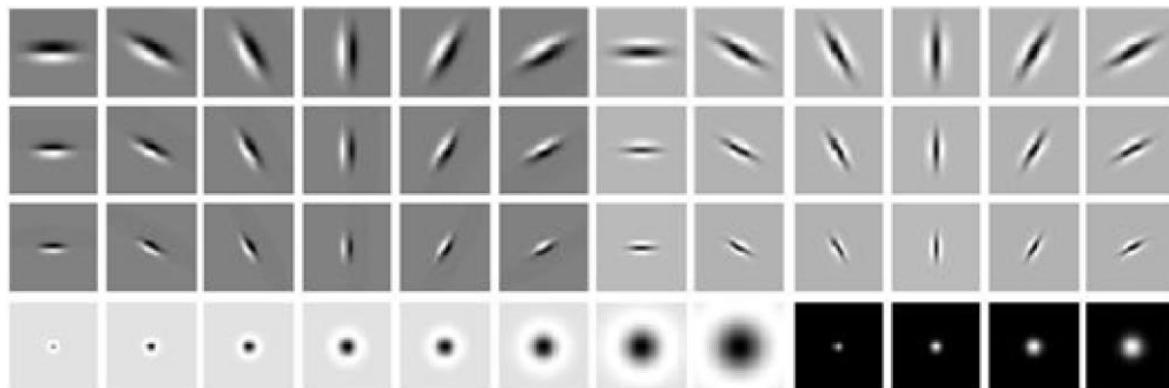


L*a*b* color space



HSV color space

- Texture (filter banks or HOG over regions)



What kind of things do we compute histograms of?

- Histograms of descriptors

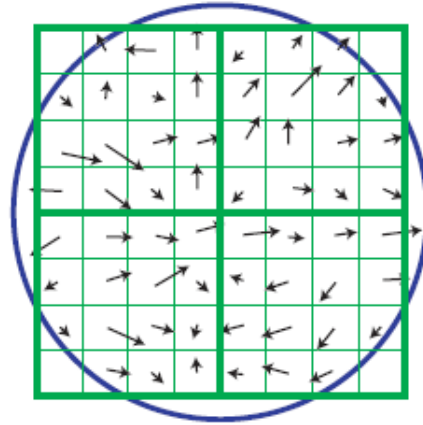
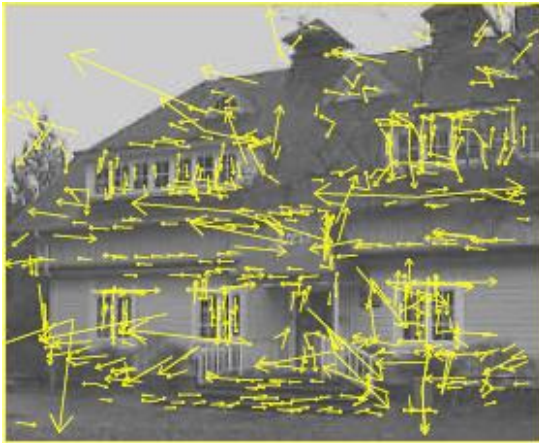
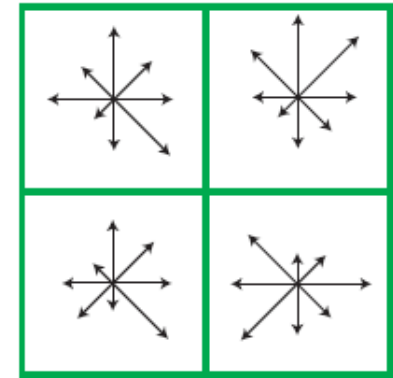


Image gradients



Keypoint descriptor

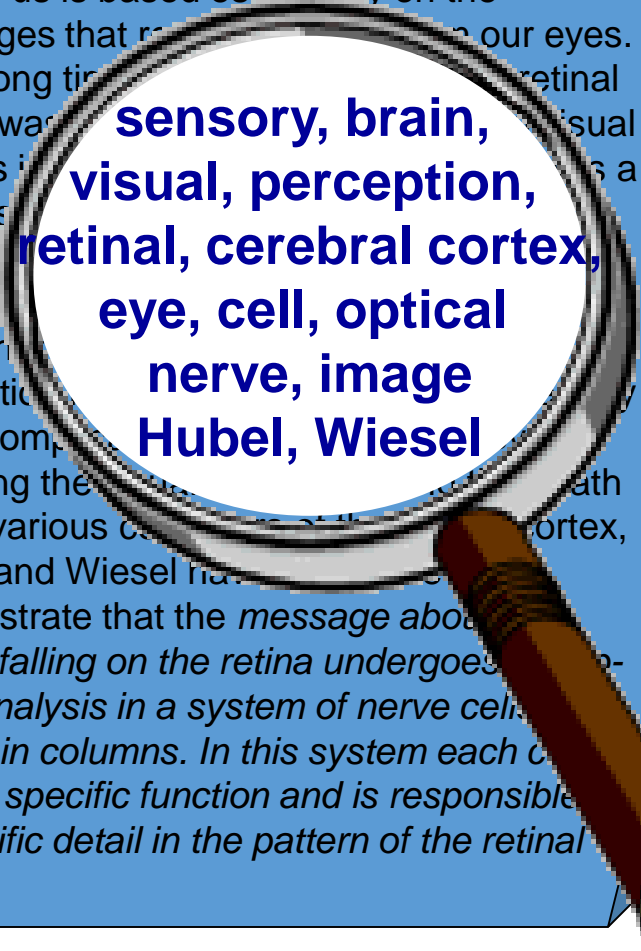
SIFT – [Lowe IJCV 2004]

- “Bag of visual 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 our eyes.

For a long time, the retinal image was considered as a movie screen. It is now known that the image is processed in a more complex way following 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**

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 \$560bn in 2004.

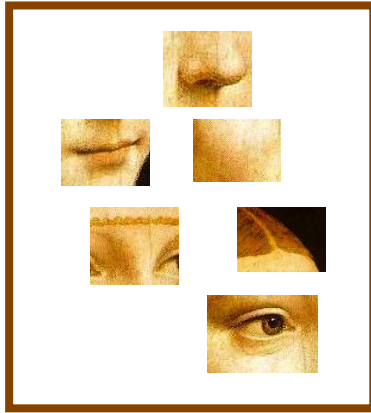
The increase in exports will be partly due to a deliberate policy to encourage exports. China's government has agreed to allow the yuan to rise in value against the dollar. The government also needs to increase the demand for the yuan in the country. China has allowed the yuan against the dollar to rise 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.



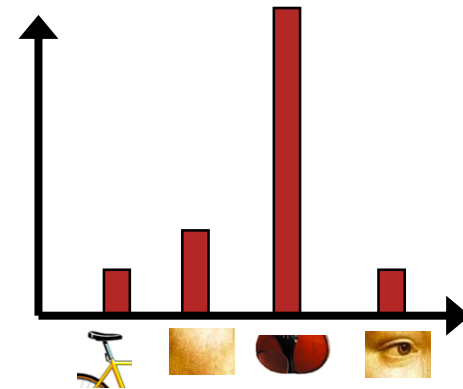
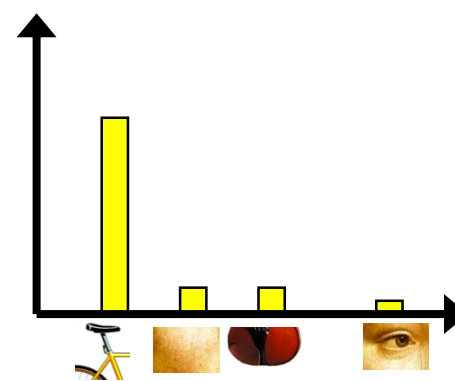
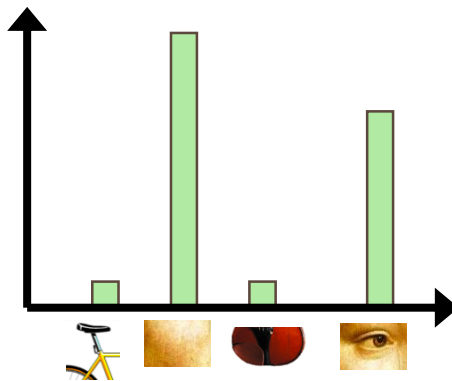
**China, trade,
surplus, commerce,
exports, imports, US,
yuan, bank, domestic,
foreign, increase,
trade, value**

Bag of *visual words*

- Image patches



- BoW histogram



- Codewords

Image categorization with bag of words

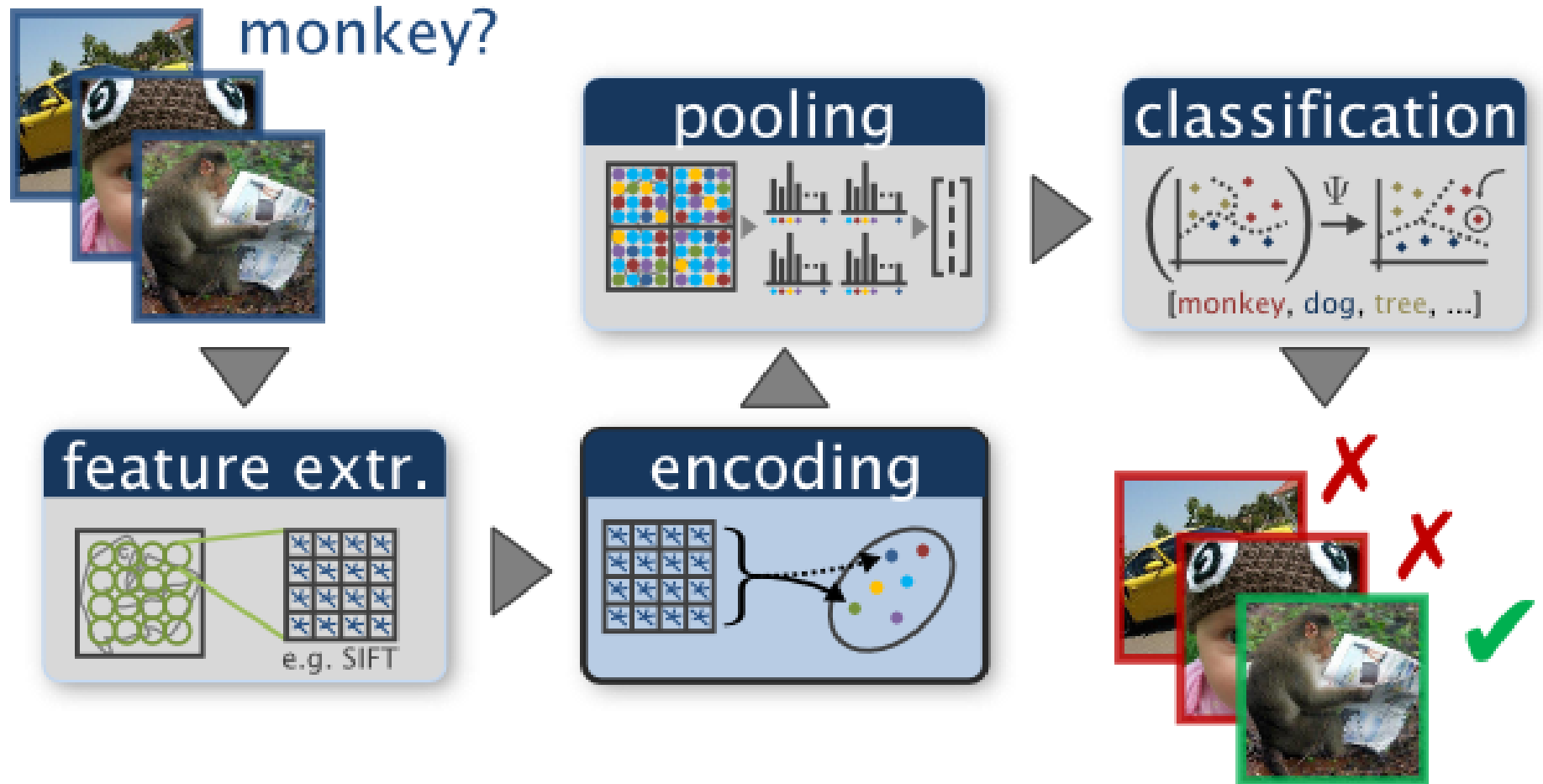
Training

1. Extract keypoints and descriptors for all training images
2. Cluster descriptors
3. Quantize descriptors using cluster centers to get “visual words”
4. Represent each image by normalized counts of “visual words”
5. Train classifier on labeled examples using histogram values as features

Testing

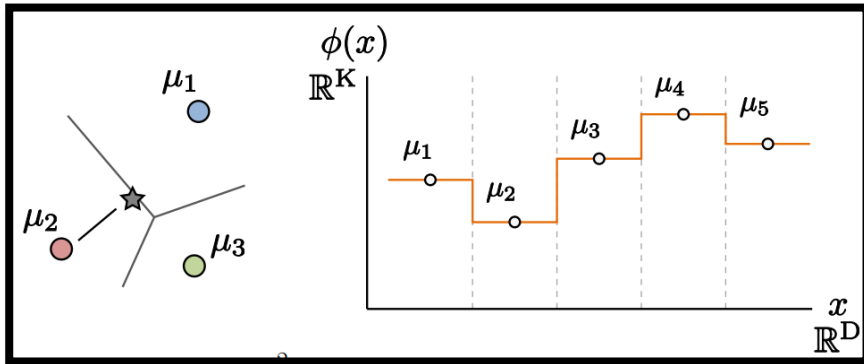
1. Extract keypoints/descriptors and quantize into visual words
2. Compute visual word histogram
3. Compute label or confidence using classifier

Bag of visual words image classification

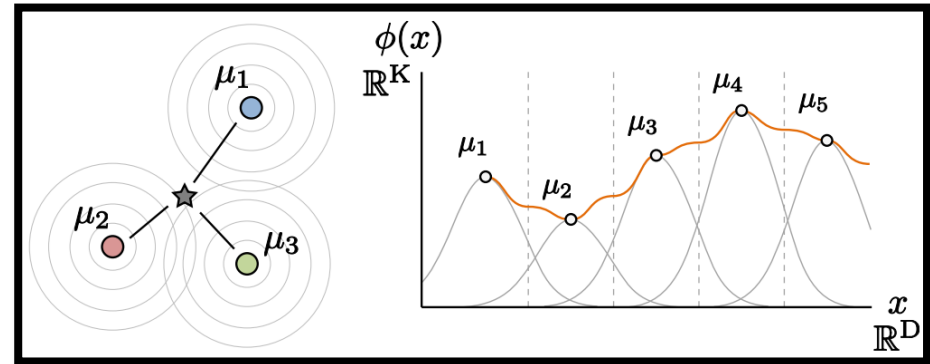


Feature encoding

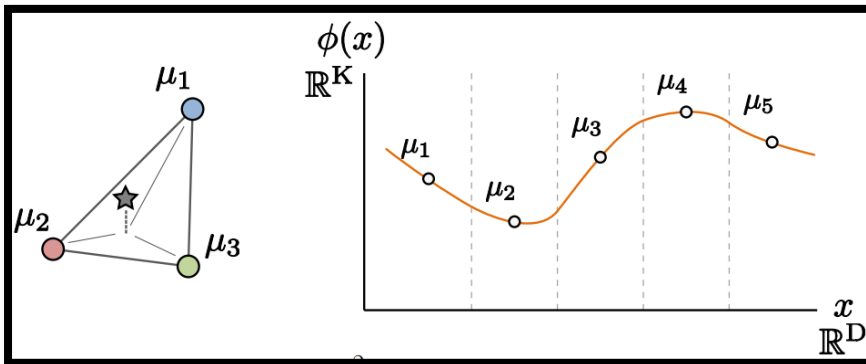
- Hard/soft assignment to clusters



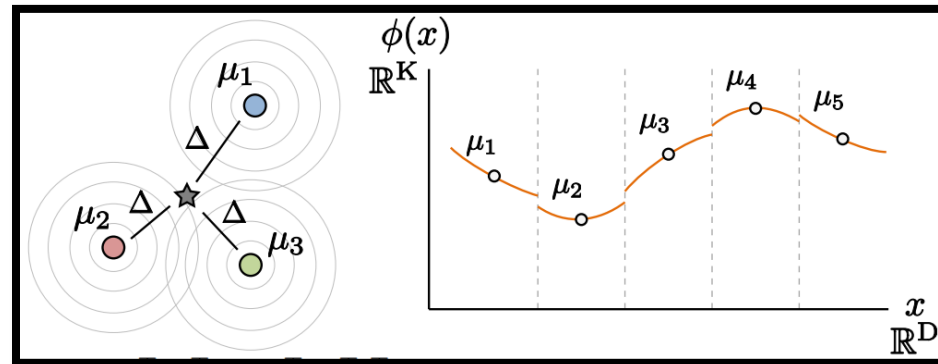
Histogram encoding



Kernel codebook encoding



Locality constrained encoding



Fisher encoding

Fisher vector encoding

- Fit Gaussian Mixture Models

$$\Theta = (\mu_k, \Sigma_k, \pi_k : k = 1, \dots, K)$$

- Posterior probability

$$q_{ik} = \frac{\exp \left[-\frac{1}{2} (\mathbf{x}_i - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_k) \right]}{\sum_{t=1}^K \exp \left[-\frac{1}{2} (\mathbf{x}_i - \mu_t)^T \Sigma_k^{-1} (\mathbf{x}_i - \mu_t) \right]}$$

- First and second order differences to cluster k

$$u_{jk} = \frac{1}{N \sqrt{\pi_k}} \sum_{i=1}^N q_{ik} \frac{x_{ji} - \mu_{jk}}{\sigma_{jk}},$$

$$v_{jk} = \frac{1}{N \sqrt{2\pi_k}} \sum_{i=1}^N q_{ik} \left[\left(\frac{x_{ji} - \mu_{jk}}{\sigma_{jk}} \right)^2 - 1 \right]$$

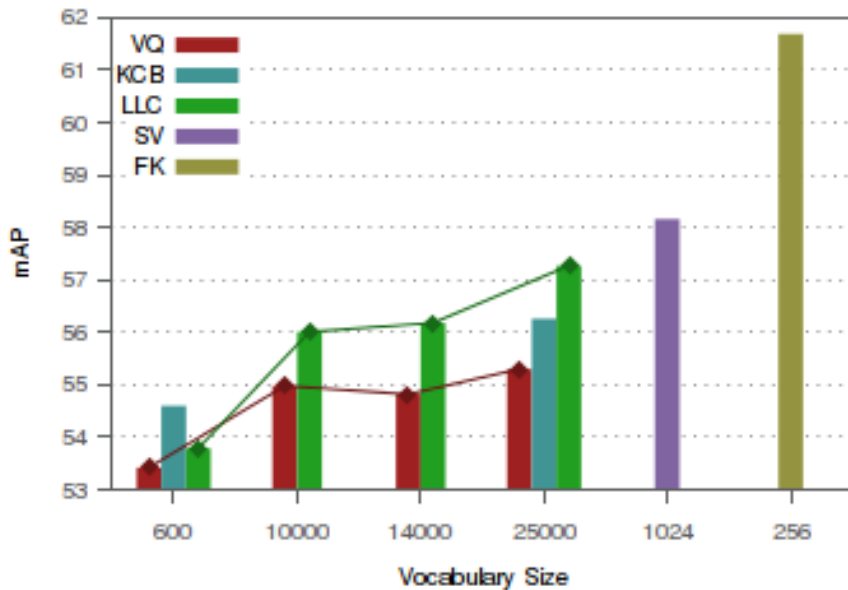
$$\Phi(I) = \begin{bmatrix} \vdots \\ \mathbf{u}_k \\ \vdots \\ \mathbf{v}_k \\ \vdots \end{bmatrix}$$

[\[Perronnin et al. ECCV 2010\]](#)

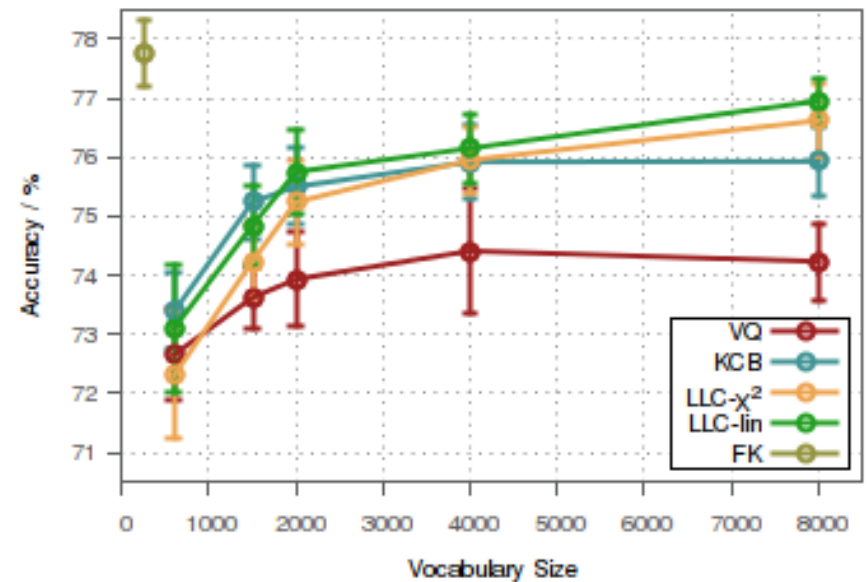
Performance comparisons

- Fisher vector encoding outperforms others
- Higher-order statistics helps

Performance over PASCAL VOC 2007

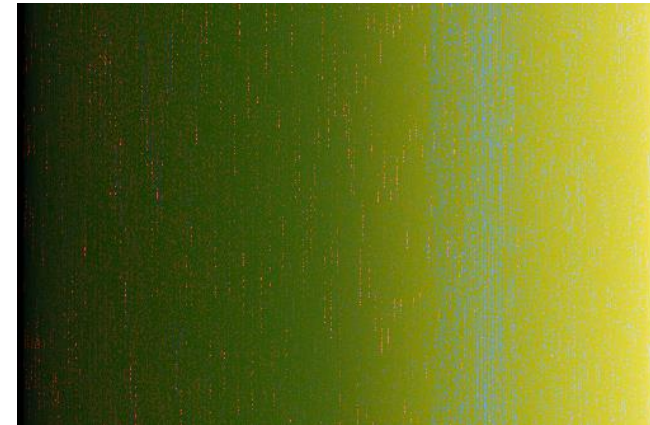
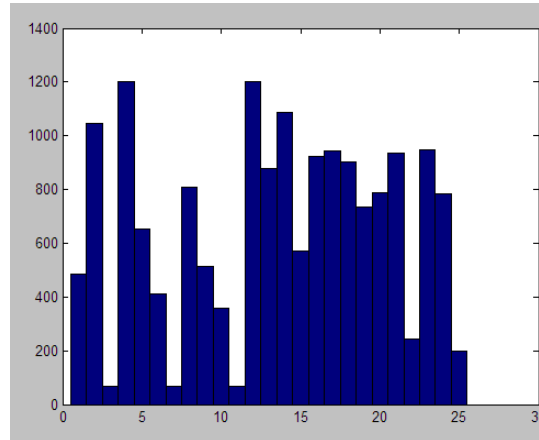


Vocabulary Size vs. Accuracy over Caltech 101



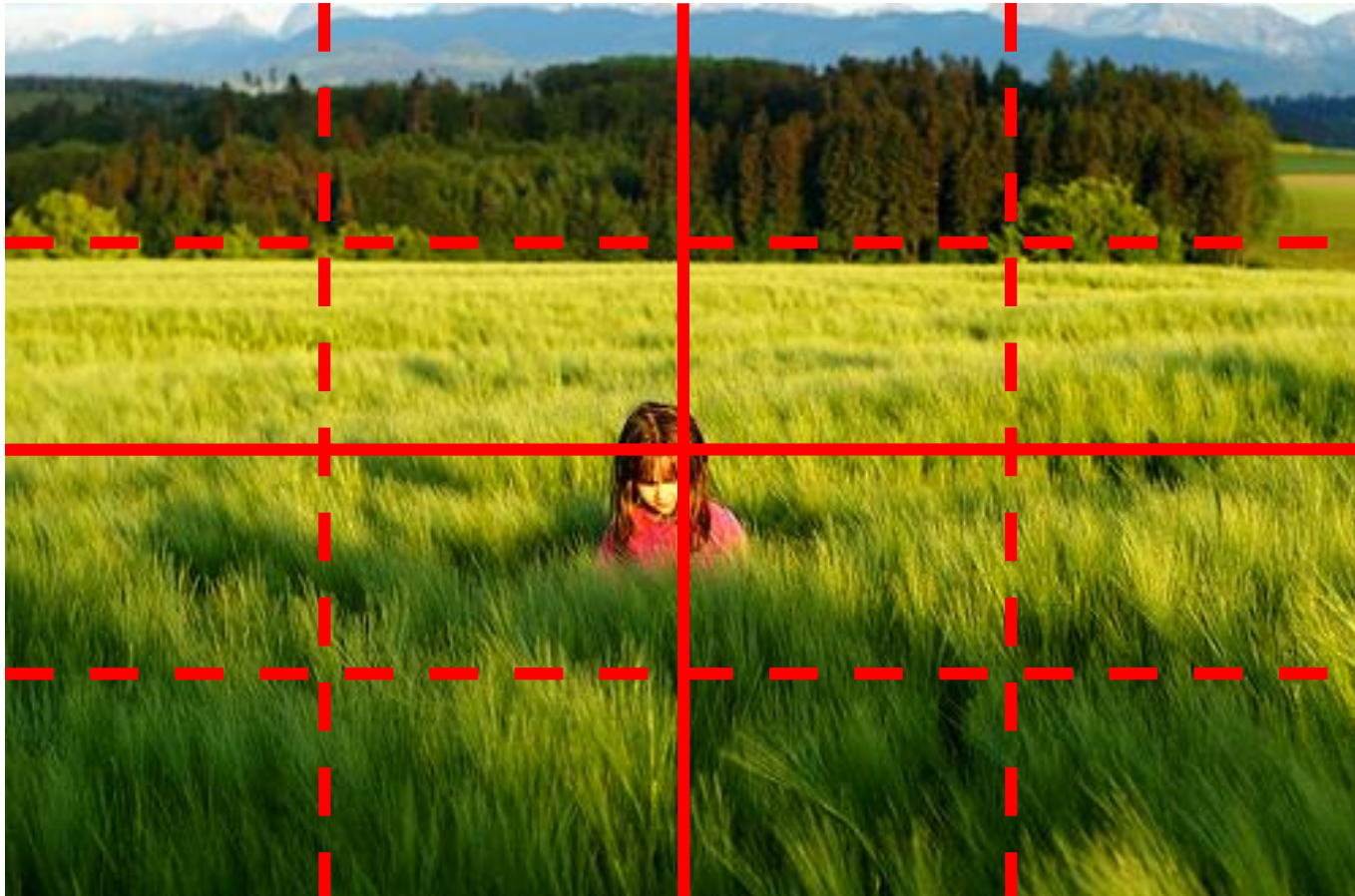
[[Chatfield et al. BMVC 2011](#)]

But what about spatial layout?



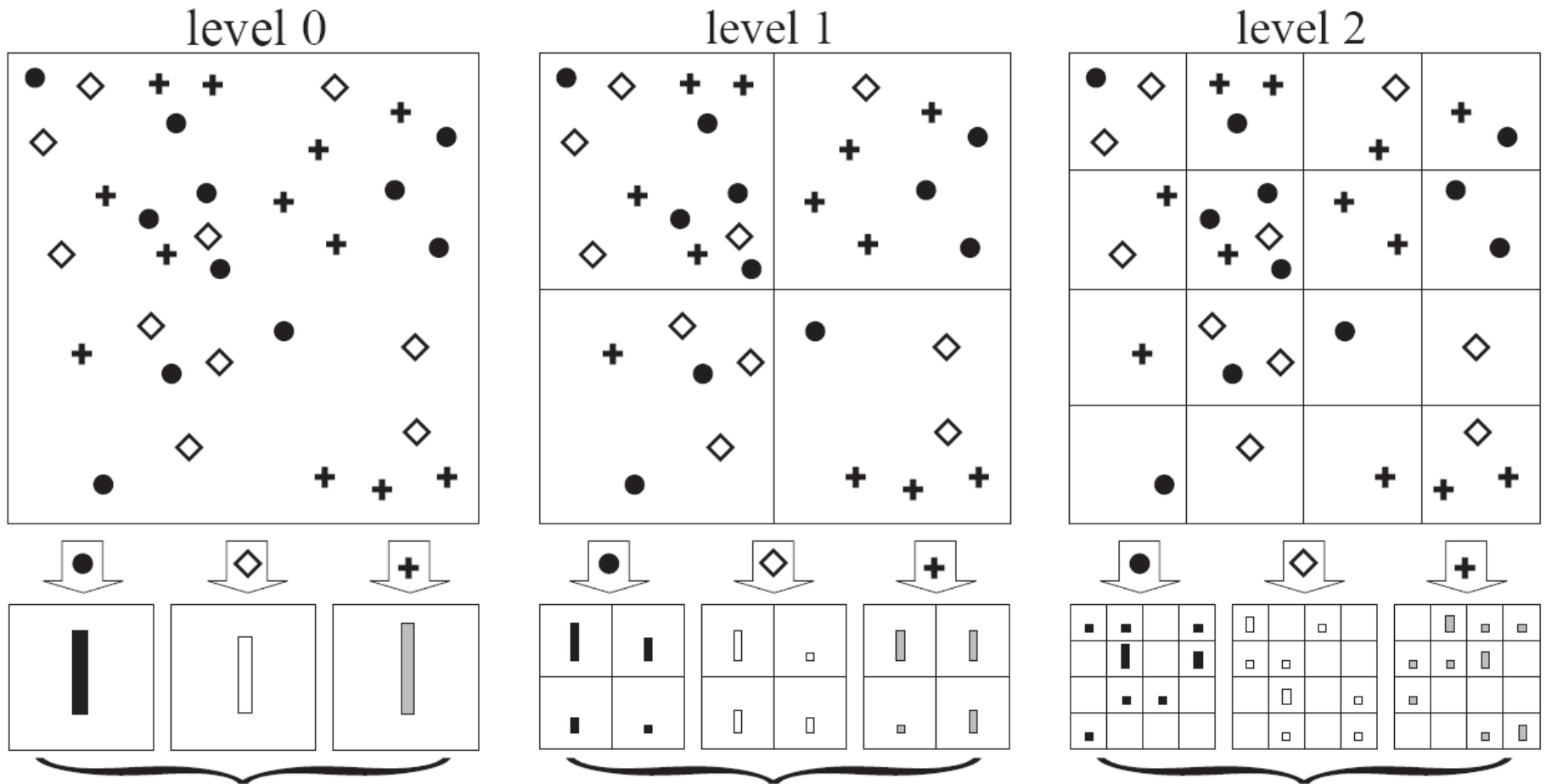
All of these images have the same color histogram

Spatial pyramid



Compute histogram in each spatial bin

Spatial pyramid

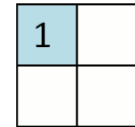
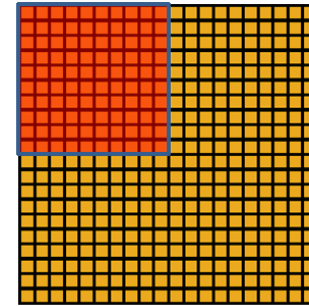
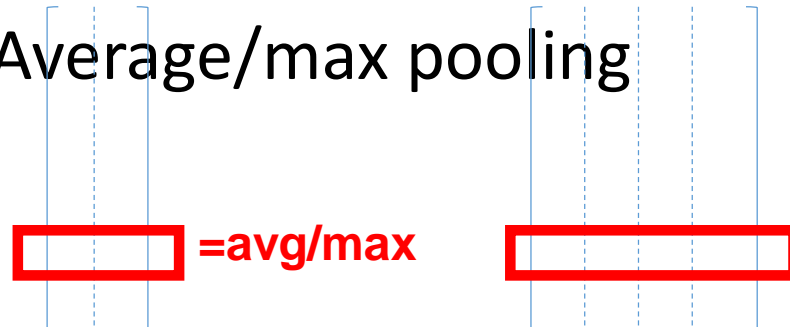


High number of features – PCA to reduce dimensionality

[[Lazebnik et al. CVPR 2006](#)]

Pooling

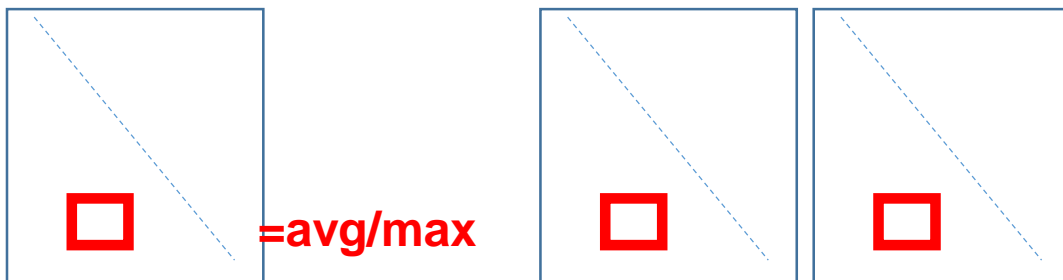
- Average/max pooling



Convolved feature Pooled feature

[Source: Unsupervised Feature Learning and Deep Learning](#)

- Second-order pooling
[[Joao et al. PAMI 2014](#)]

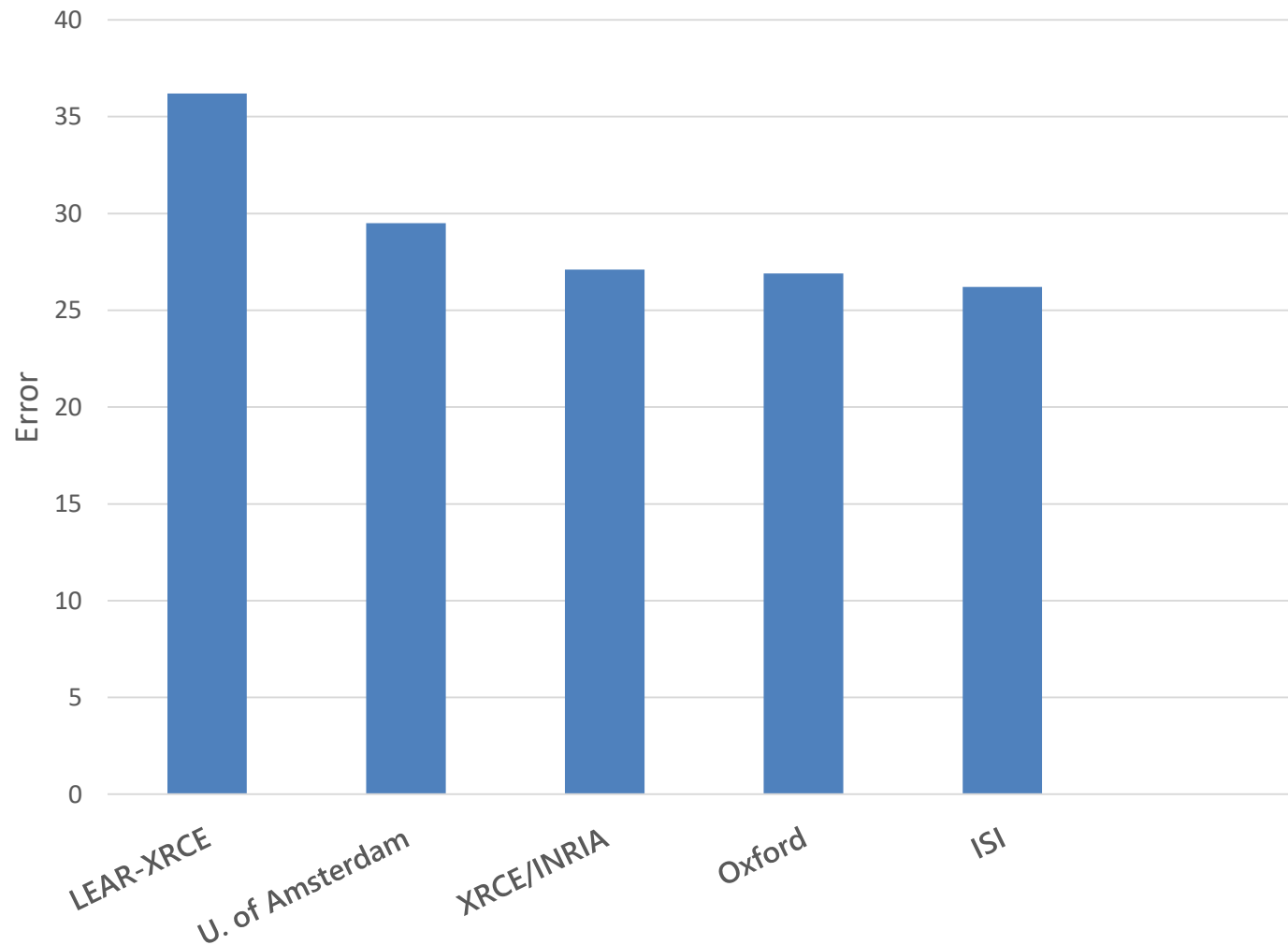


$$\mathbf{G}_{avg}(R_j) = \frac{1}{|F_{R_j}|} \sum_{i:(\mathbf{f}_i \in R_j)} \mathbf{x}_i \cdot \mathbf{x}_i^\top$$

$$\mathbf{G}_{max}(R_j) = \max_{i:(\mathbf{f}_i \in R_j)} \mathbf{x}_i \cdot \mathbf{x}_i^\top$$

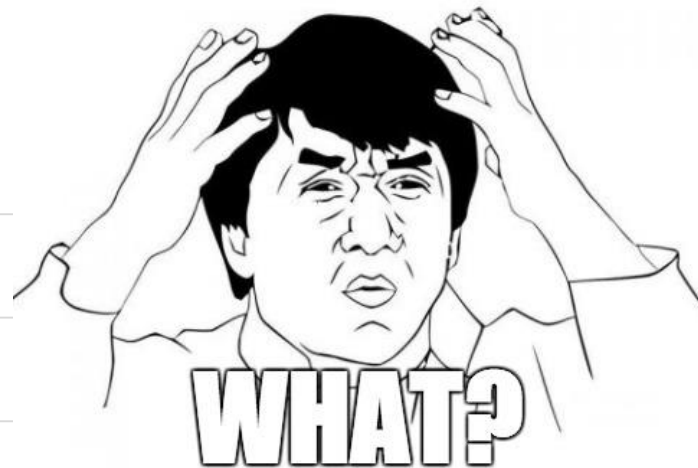
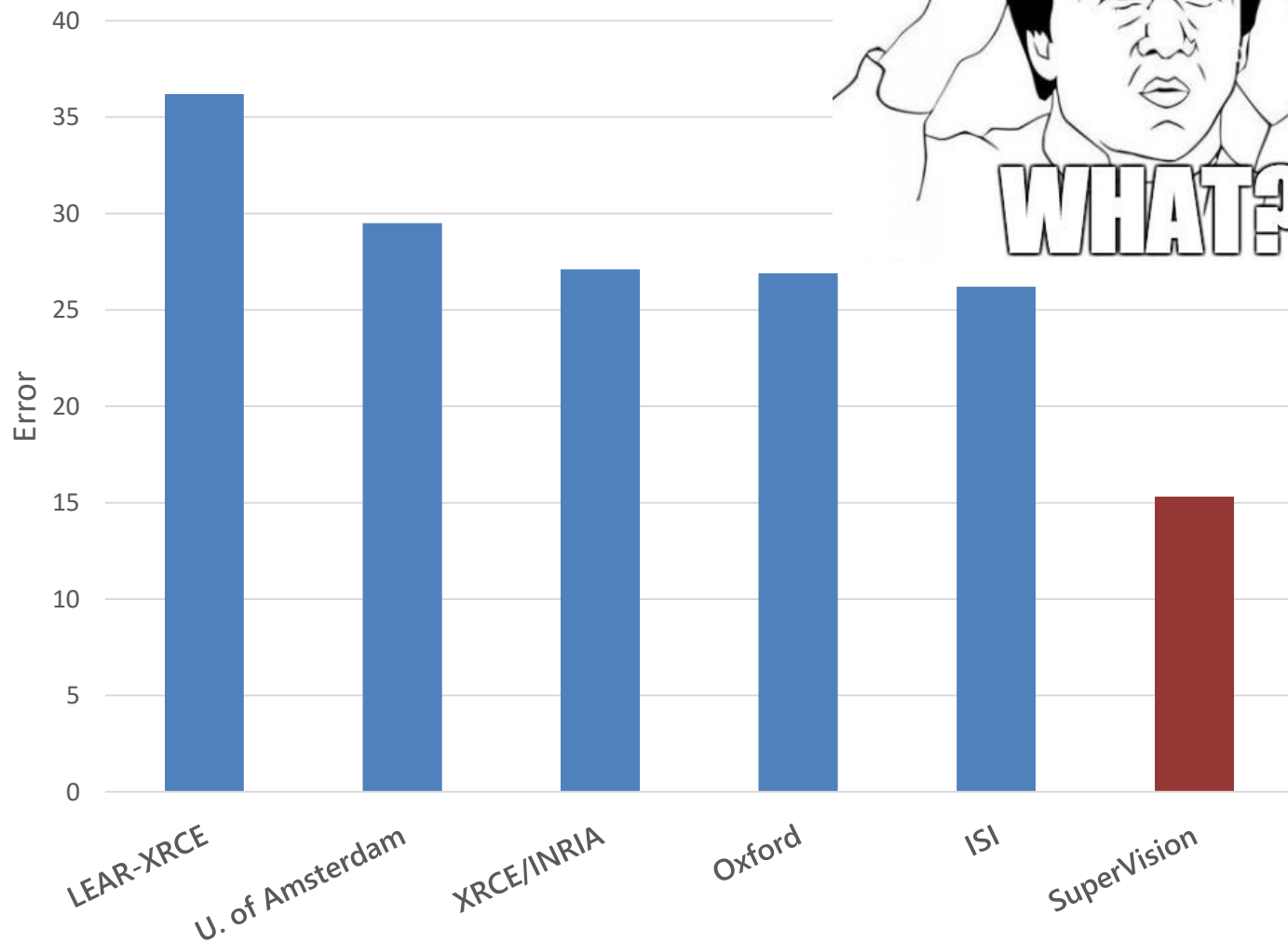
2012 ImageNet 1K

(Fall 2012)



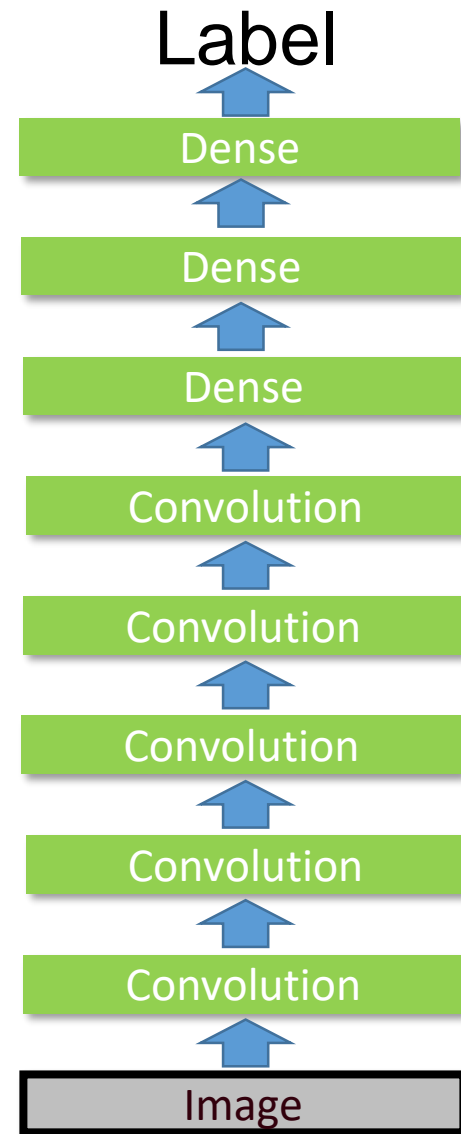
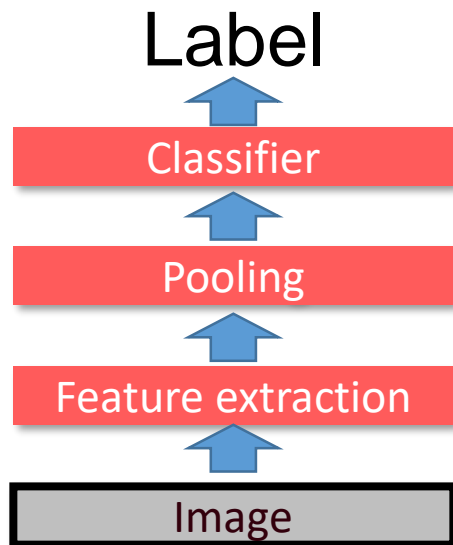
2012 ImageNet 1K

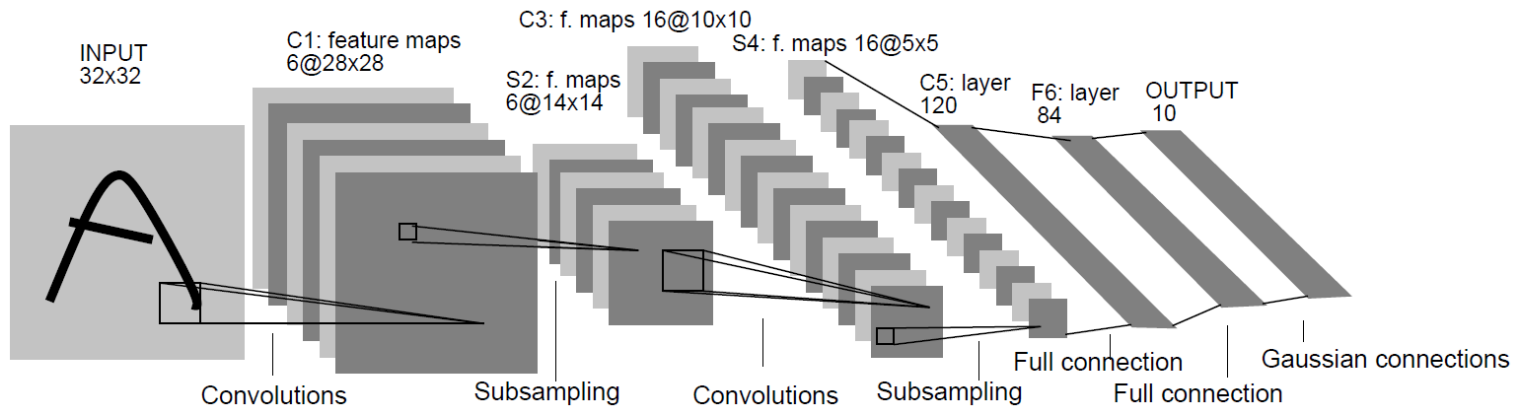
(Fall 2012)



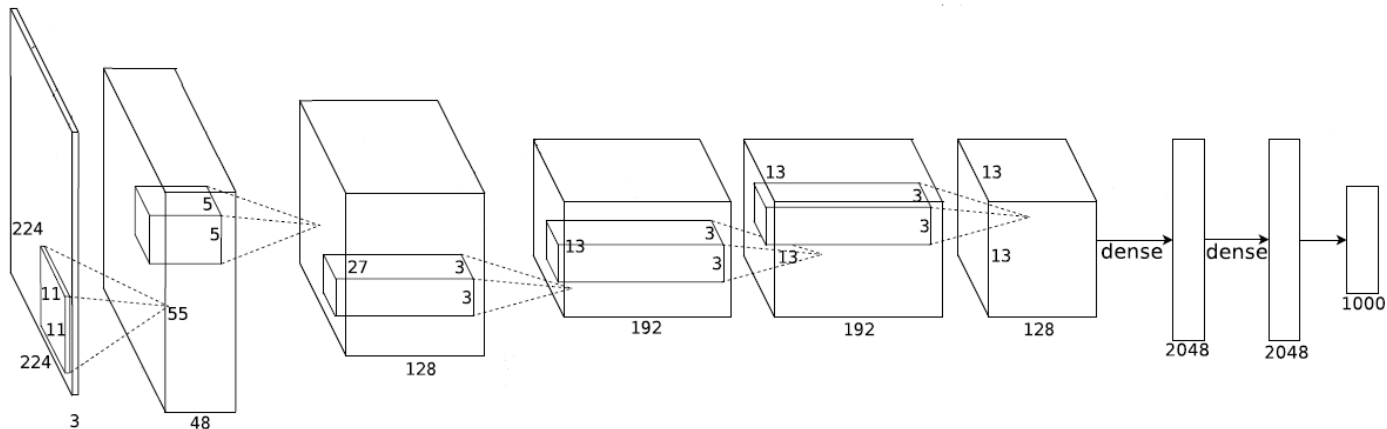
Shallow vs. deep learning

- Engineered vs. learned features

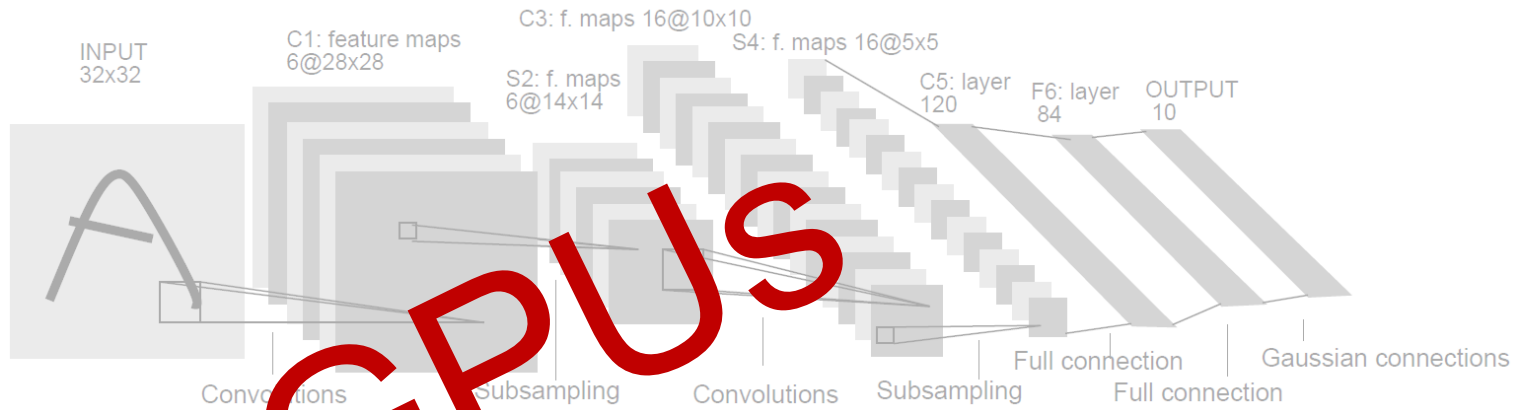




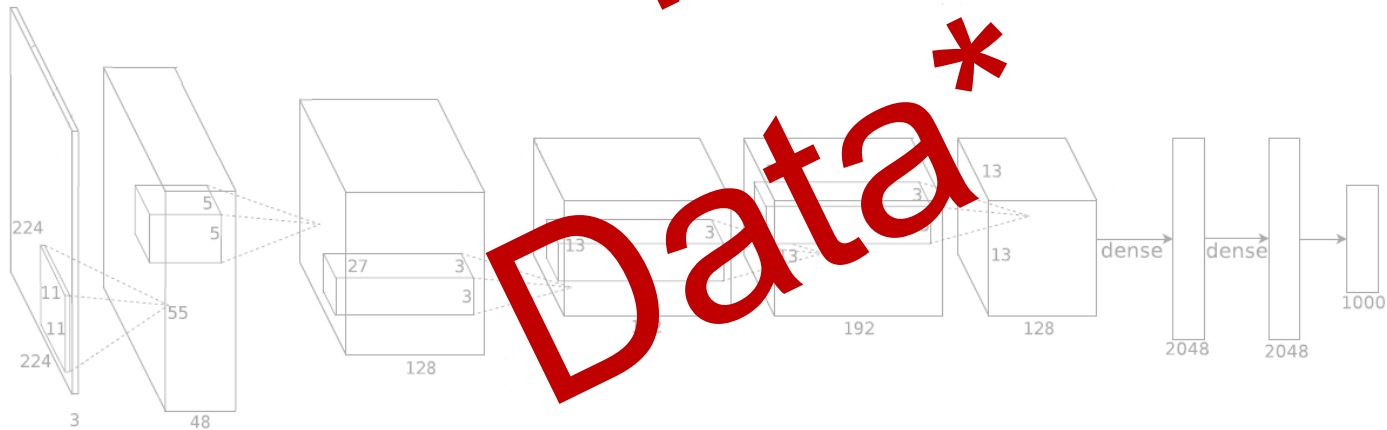
Gradient-Based Learning Applied to Document Recognition, LeCun, Bottou, Bengio and Haffner, Proc. of the IEEE, **1998**



Imagenet Classification with Deep Convolutional Neural Networks, Krizhevsky, Sutskever, and Hinton, NIPS **2012**



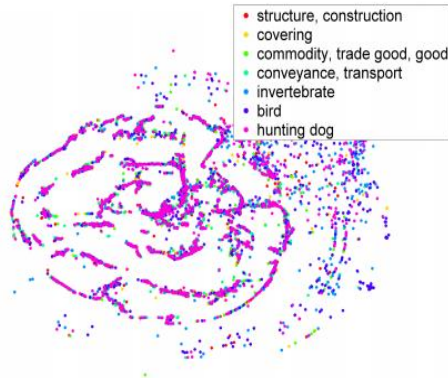
Gradient-Based Learning Applied to Document Recognition, LeCun, Bottou, Bengio and Haffner, Proc. of the IEEE, 1998



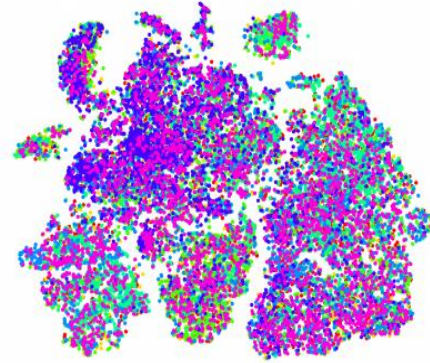
Imagenet Classification
Networks, Krizhevsky et al.

* Rectified activations and dropout

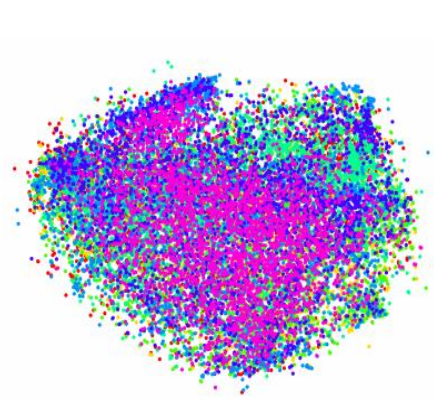
Convolutional activation features



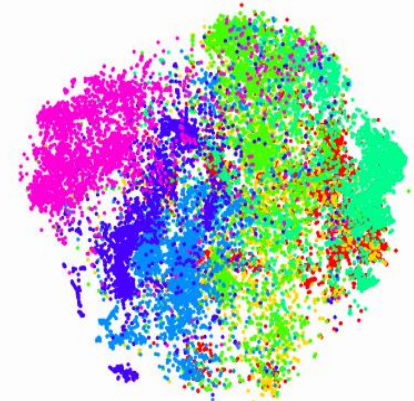
(a) LLC



(b) GIST

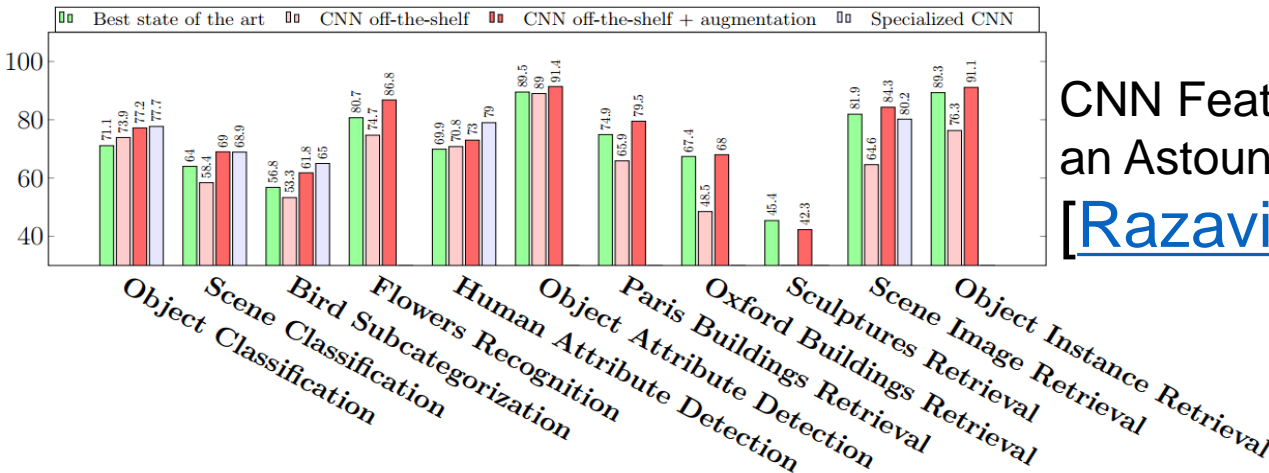
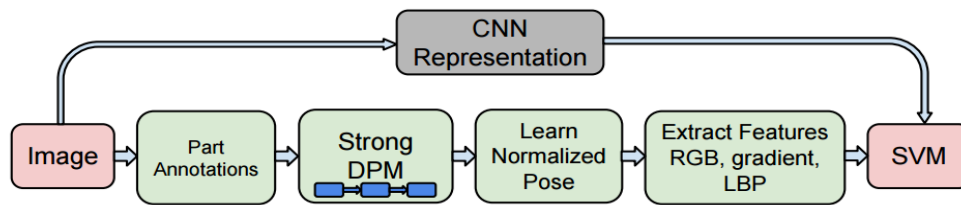


(c) DeCAF₁



(d) DeCAF₆

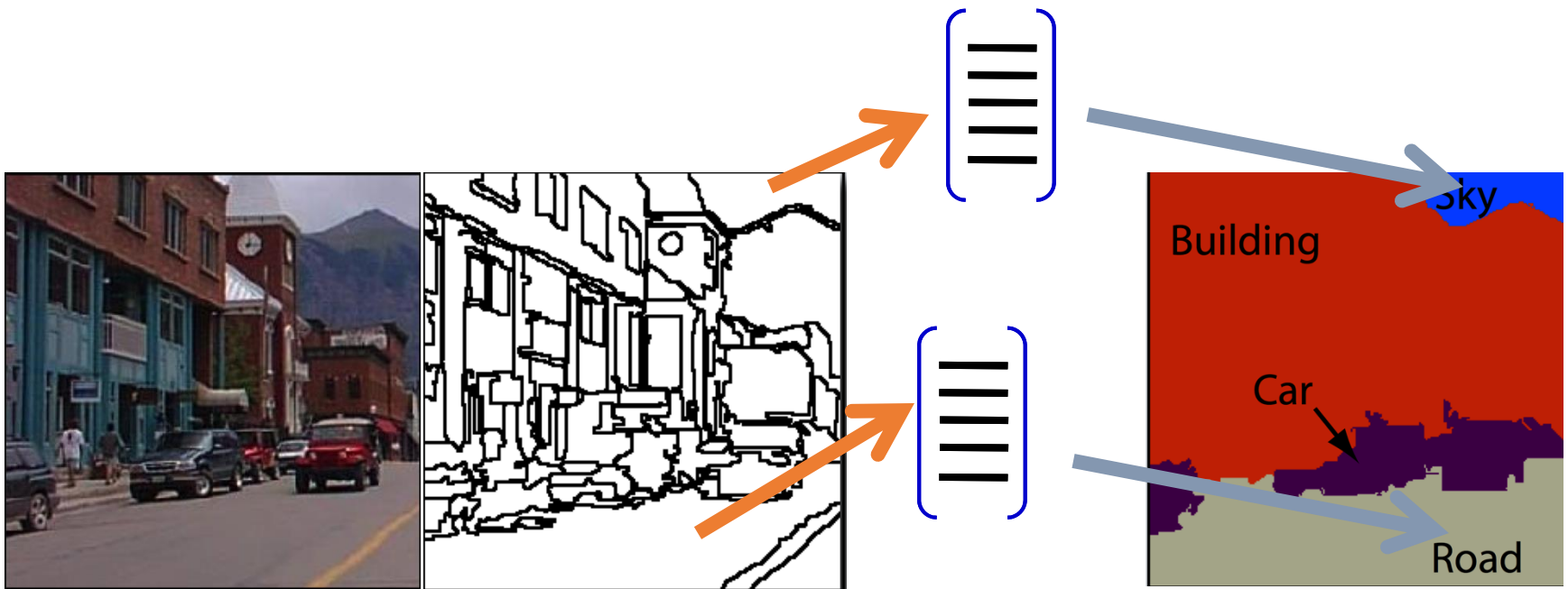
[[Donahue et al. ICML 2013](#)]



CNN Features off-the-shelf:
an Astounding Baseline for Recognition
[[Razavian et al. 2014](#)]

Region representation

- Segment the image into superpixels
- Use features to represent each image segment



Region representation

- Color, texture, BoW
 - Only computed within the local region
- Shape of regions
- Position in the image

Working with regions

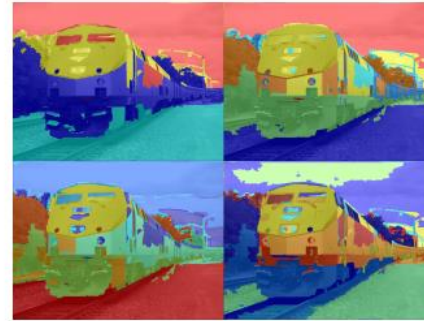
- Spatial support is important – multiple segmentation



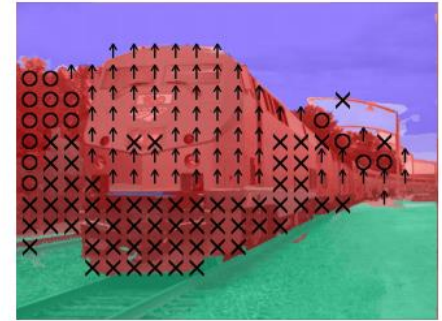
(a) Input



(b) Superpixels



(c) Multiple Hypotheses



(d) Geometric Labels

Geometric context [[Hoiem et al. ICCV 2005](#)]

- Spatial consistency – MRF smoothing

Things to remember

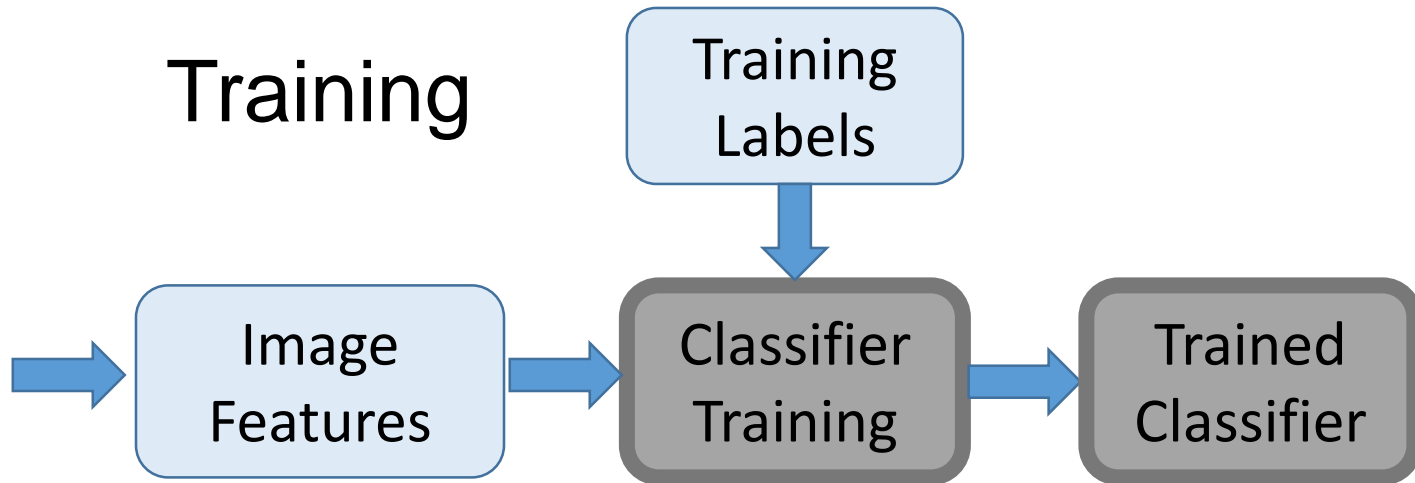
- Visual categorization help transfer knowledge
- Image features
 - Coverage, concision, directness
 - Color, gradients, textures, motion, descriptors
 - Histogram, feature encoding, and pooling
 - CNN as features
- Image/region categorization

Next lecture - Classifiers

Training Images



Training



Testing



Test Image

