

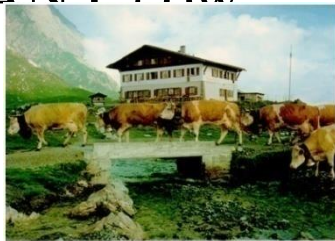
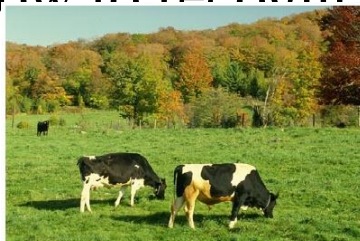
Object Detection and Recognition

Object Categorization

- How to recognize ANY car



- How to recognize ANY cow



Challenges: robustness



Illumination



Object pose



Clutter



Occlusions

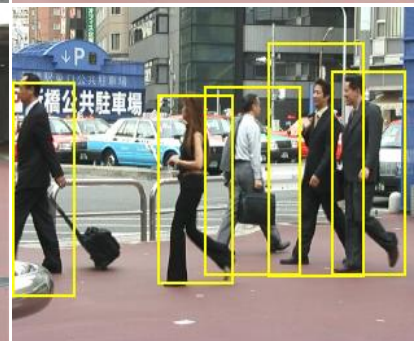
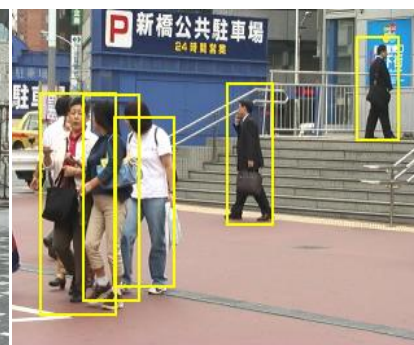
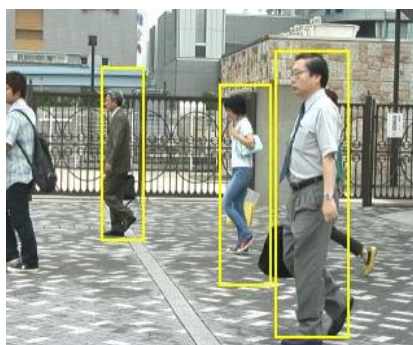


**Intra-class
appearance**



Viewpoint

Challenges: robustness



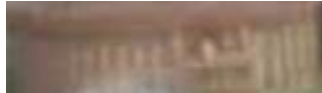
- Detection in Crowded Scenes

- Learn object variability

- ◆ Changes in appearance, scale, and articulation

- Compensate for clutter, overlap, and occlusion

Challenges: context and human experience

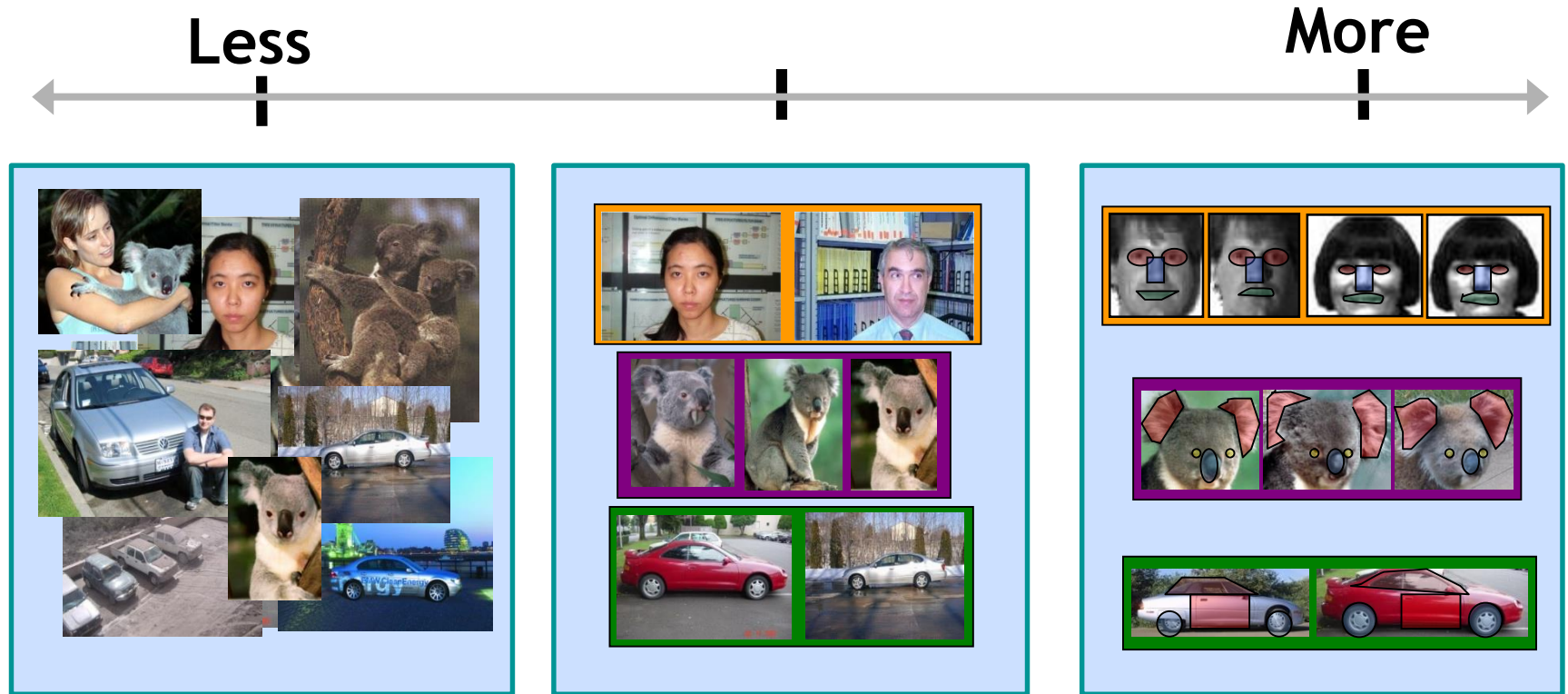


Challenges: context and human experience



Context cues

Challenges: learning with minimal supervision



Unlabeled,
multiple
objects

Classes labeled,
some clutter

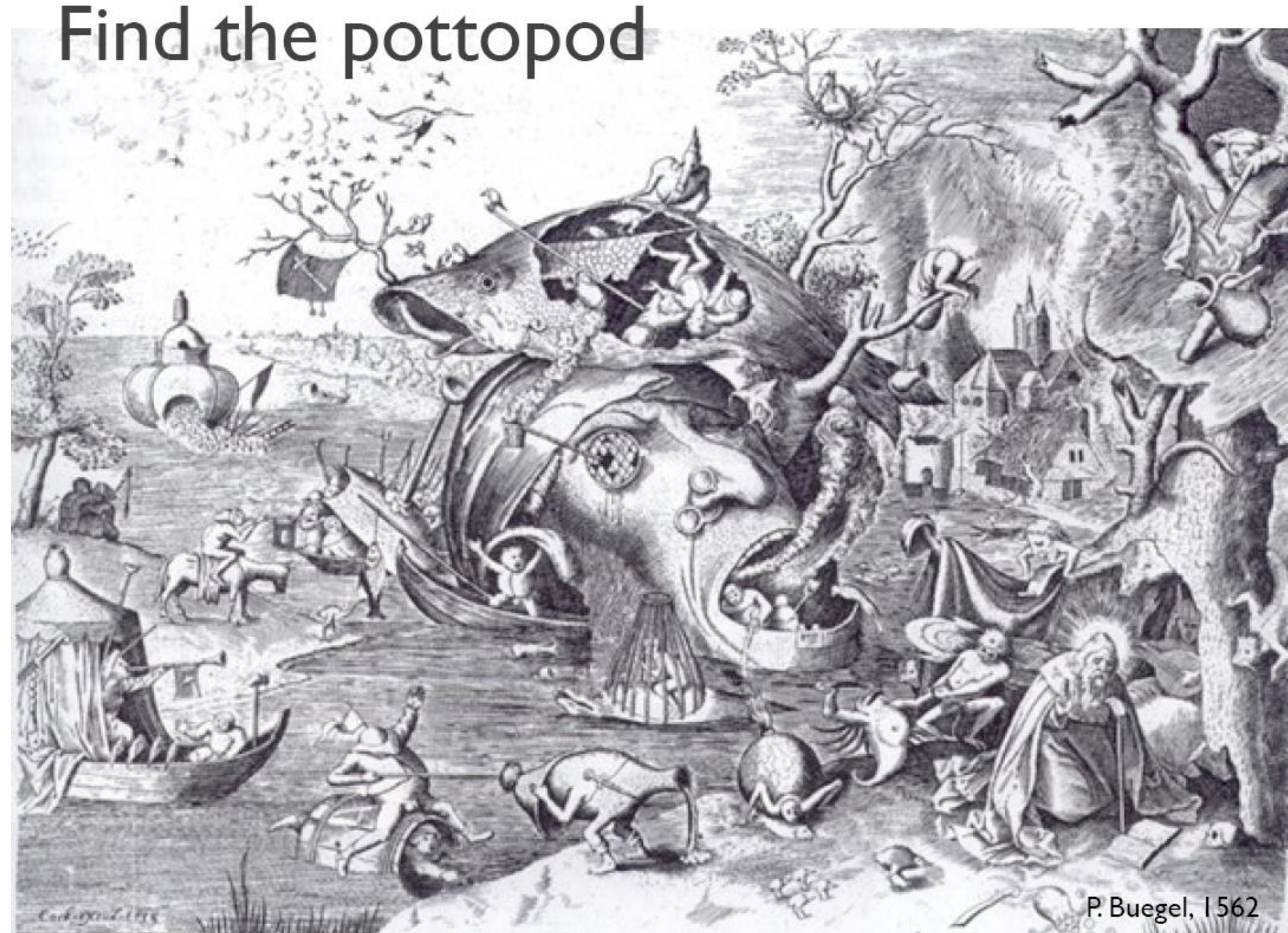
Cropped to
object, parts and
classes labeled



S. Savarese, 2003

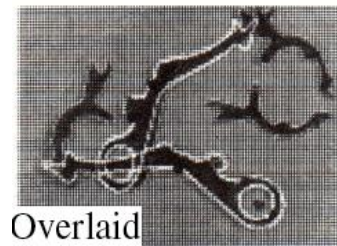
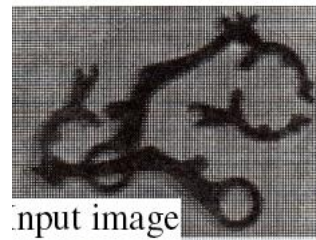
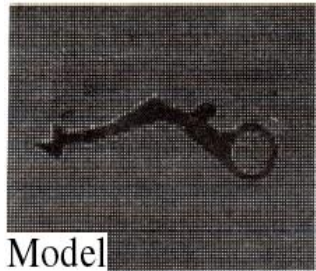
This is a
pottopod

Find the pottopod



P. Buegel, 1562

Rough evolution of focus in recognition research



1980s



7 5 9 2 6 5
2 2 2 2 2 3
0 2 3 8 0 7



1990s to early 2000s

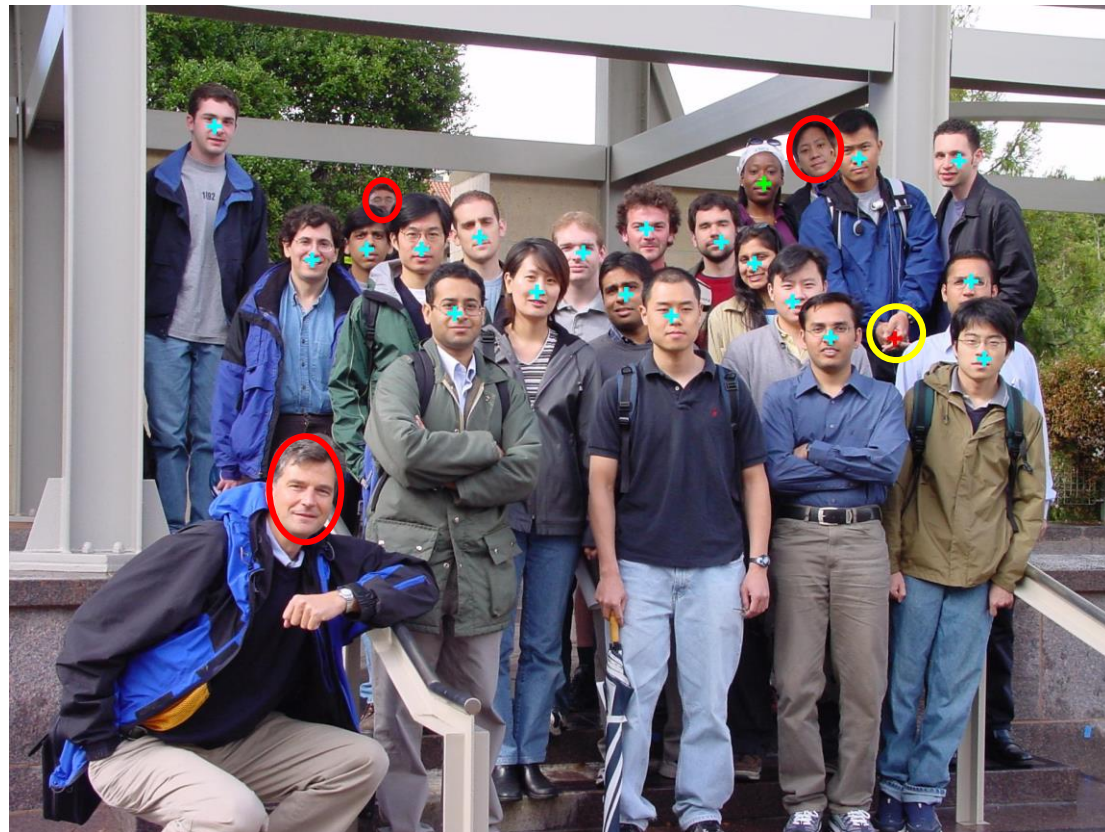
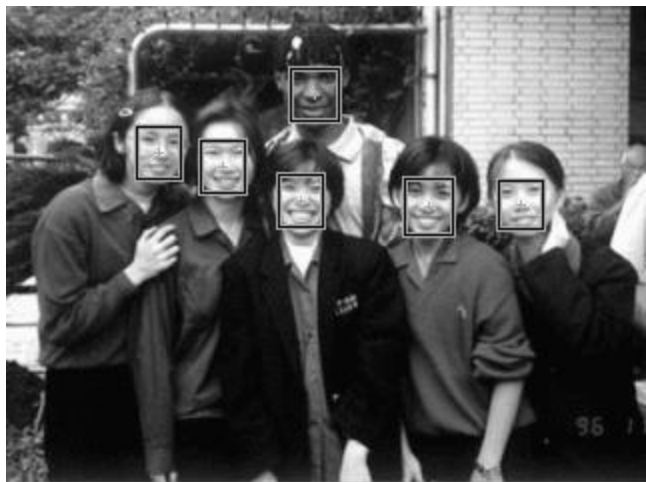


2000-2010...

Detection, recognition, and classification

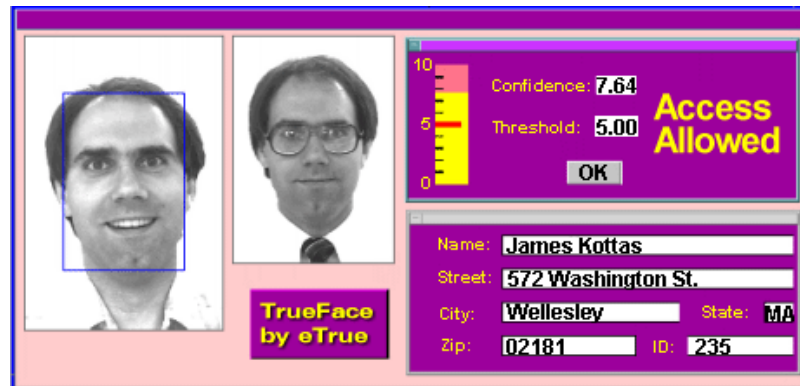
- Detection = 2-class classification problem
 - Object/class or not object/class
 - E.g., detect all the faces in this image
- Recognition of identity = within-class classification problem
 - Within a given class of objects (e.g., faces, logos), identify the object as one particular member of the class (e.g., Joe's face, Nike logo)
- Recognition of class = among-class classification
 - Which class of things is this: sky, cloud, forest, face, ...

Example: Face detection



Found Face at [x=108, y=80]
Found Face at [x=76, y=73]
Found Face at [x=257, y=99]
Found Face at [x=154, y=44]
Found Face at [x=211, y=100]
Found Face at [x=147, y=97]

Example: Face recognition



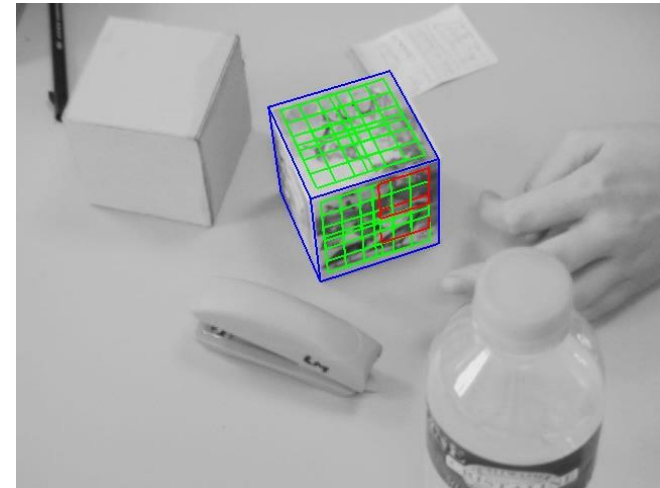
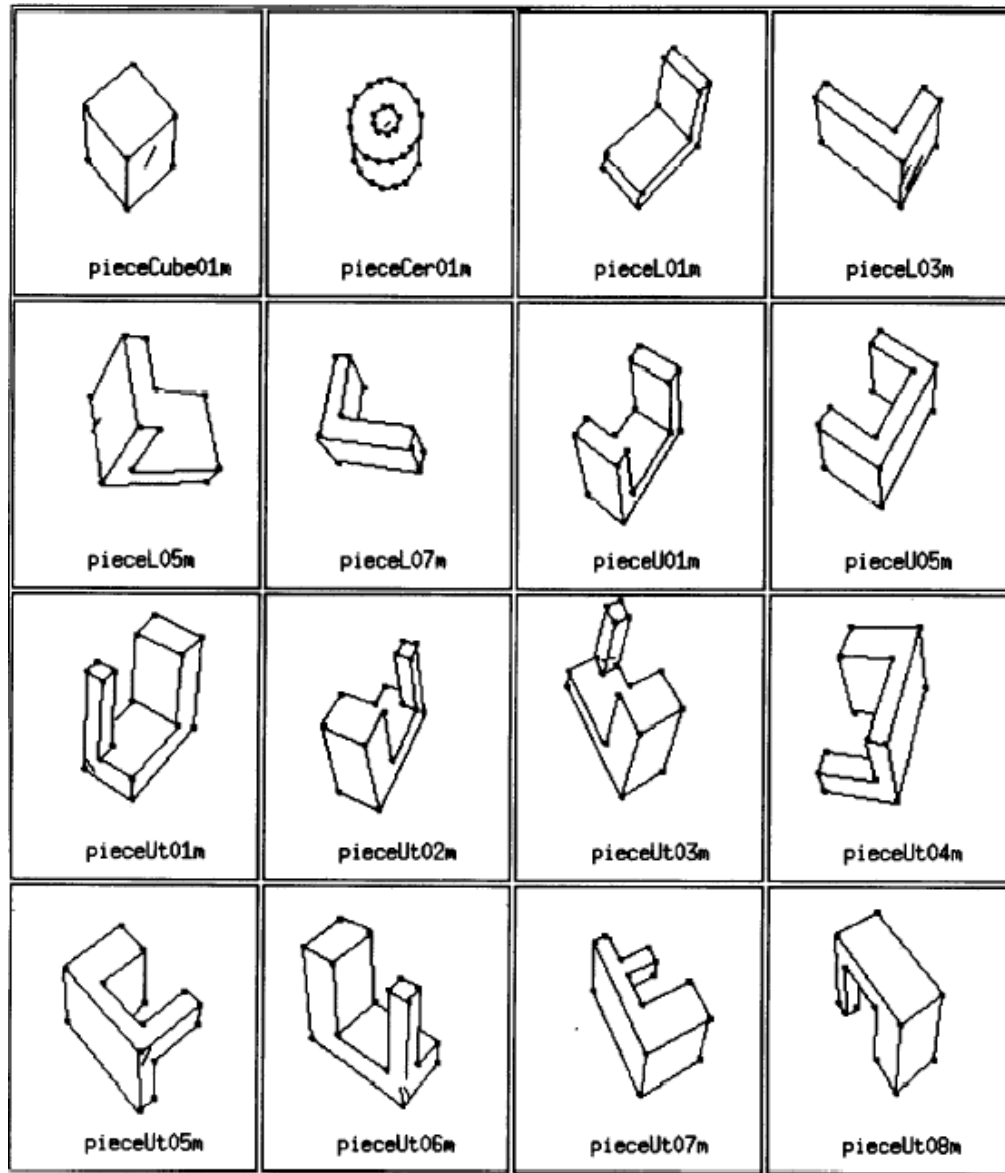
TrueFace by eTrue

Confidence: 7.64
Threshold: 5.00
Access Allowed
OK

Name: James Kottas
Street: 572 Washington St.
City: Wellesley State: MA
Zip: 02181 ID: 235



Example: Polyhedral object recognition



Approaches to detection and classification

- There are many approaches to object detection and recognition, depending on how the object is modeled
 - Template-based: Match an image template (or a family of image templates) of the desired object
 - Feature-based: Derive image features and then match with feature model of object
 - ♦ Colors, texture, edges, corners, ...
 - Shape (2D or 3D): Describe (parameterize) the object contour or full shape, and look for that shape in the image
 - And many more...
- In some sense, all these can be viewed as three steps:
 - Modeling the object(s) (“training”)
 - Preprocessing the image (computing features, shape, ...)
 - Classify based on a comparison or match between model and image data

Template matching and classification

- If we want to detect and recognize (classify) objects in images, one simple technique is to use normalized correlation
 - Provides a measure of how well the correlation template matches the image region
 - I.e., “template matching”
- But in general there is not just one template to match
 - E.g., in face recognition – possibly many example templates (different people, expressions, lighting, rotation, scale, ...)
 - A *classifier* takes an input feature set and produces an output class label

Classifiers

- A classifier assigns a label to any new example
 - E.g., the object name
 - Classes: {Joe, Bob, Mary, Fred, Lisa, unknown}
- A two-class classifier is a detector
 - Classes: {face, no face}
- The classifier is trained from a *training set*
 - Training set: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), \dots$
 - x_i – measurements (image, features, histogram, ...)
 - y_i – labels
 - This is typically framed as a **learning** problem
- Outcome: $(i \rightarrow j)$ means outcome i is labeled as j
 - (*Matthew* \rightarrow *Ralph*) – error
 - (*Matthew* \rightarrow *Matthew*) – correct

Training set

- Training set examples:
 - (image₁, “Joe”), (image₂, “Fred”), (image₃, “Sue”), ...
 - (color₁, “Face”), (color₂, “Hair”), (color₃, “Lips”), ...
 - (template₁, “eye”), (template₂, “eye”), (template₃, “eye”), ...
 - Perhaps *negative* examples also: (image_i, “Not a face”)...
- We want a rule (function) that does
 - $F(\text{new measurement}) = \text{label}$
...with a low error rate
- Errors
 - False positives: **Yes** when the true answer is **No**
 - False negatives: **No** when the true answer is **Yes**
 - Misclassifications: **A** when the true answer is **B**

Classification errors

- For detection (two-class)

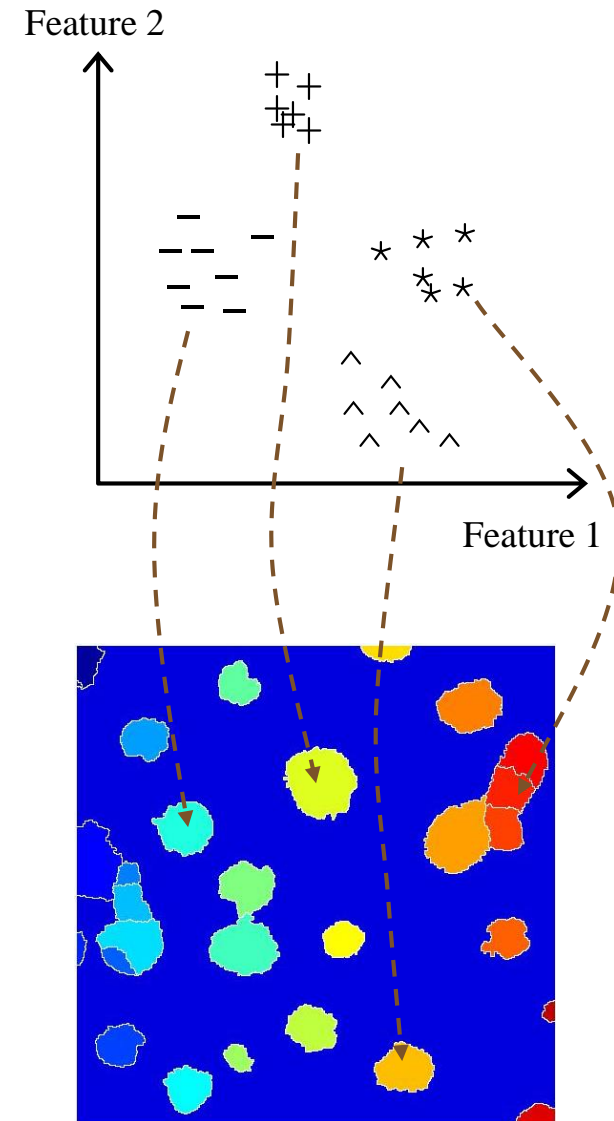
	Absent (0)	Present (1)
Not detected (0)	True negative (0 \rightarrow 0)	False negative (1 \rightarrow 0)
Detected (1)	False positive (0 \rightarrow 1)	True positive (1 \rightarrow 1)

Misclassifications: $(i \rightarrow j)$, where $i \neq j$

Misclassifications = False negatives + False positives

Segmentation and clustering

- Segmentation is about labeling similar pixels as belonging to the same group or segment
 - Pixels that *belong together* = pixels that *cluster*
- Clustering can be done along many dimensions (intensity, color, depth, motion, texture, ...)
 - Individually or combined
- There are some basic clustering methods that do well in certain cases
 - E.g., “k-means clustering”



Clustering/segmenting by k -means

- The “ k -means” algorithm is a fast, simple way to cluster N -dimensional data
 - Given a bunch of data points, group them into k different clusters
 - Each data point is typically a feature vector
 - ♦ But could even be RGB values
- We would like to minimize the *objective function*

$$\Theta(\mathbf{c}_i, \mathbf{x}_j) = \sum_{i \in \text{clusters}} \sum_{j \in \text{cluster}(i)} (\mathbf{x}_j - \mathbf{c}_i)^T (\mathbf{x}_j - \mathbf{c}_i)$$

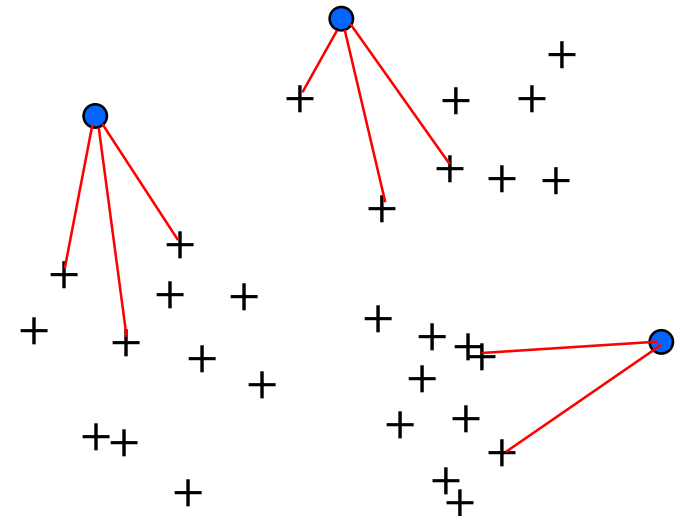
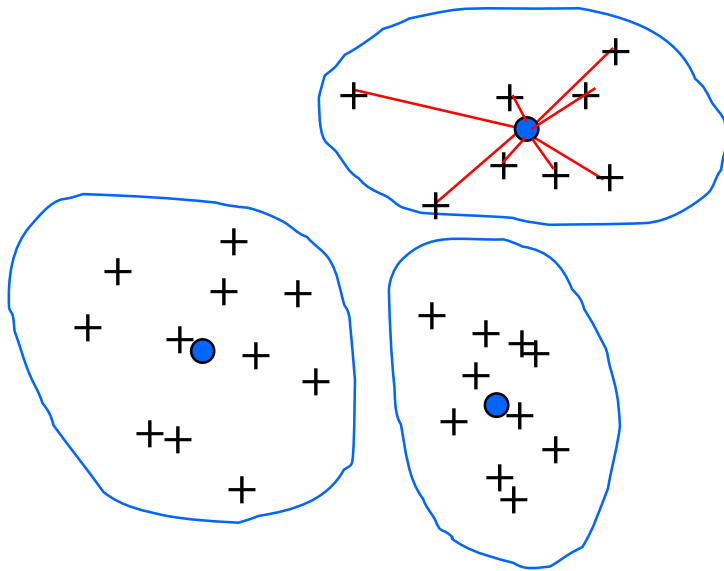
...but this is too expensive to do for lots of data points!

k -means clustering

$$\Theta(\text{clusters}, \text{data}) = \sum_{i \in \text{clusters}} \sum_{j \in \text{cluster}(i)} (x_j - c_i)^T (x_j - c_i)$$

Data points
Cluster centers

= the sum of the squares of the distances to cluster centers (means)



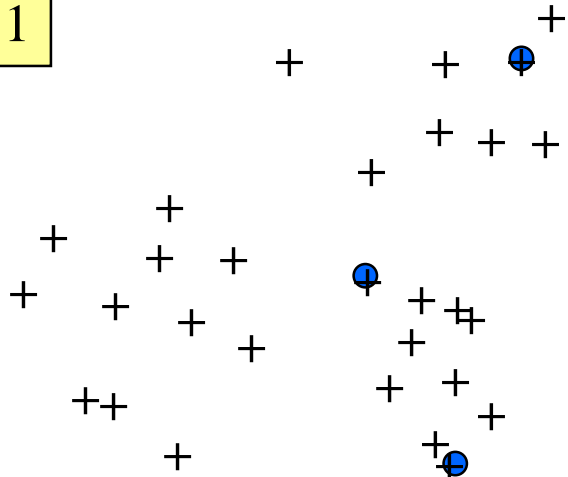
The objective function would be larger in this case

k -means clustering

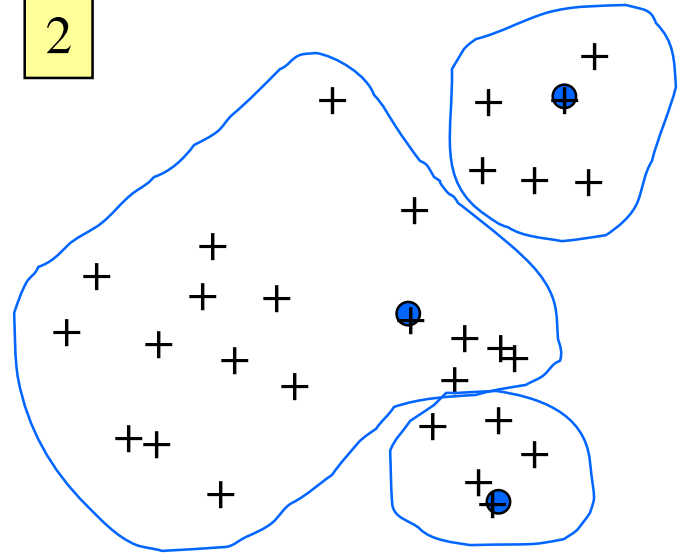
Algorithm

- Randomly choose k data points to be the initial cluster centers
- Iterate until centers are stable:
 - Assign each point to the nearest cluster
 - Recalculate the cluster center (mean)

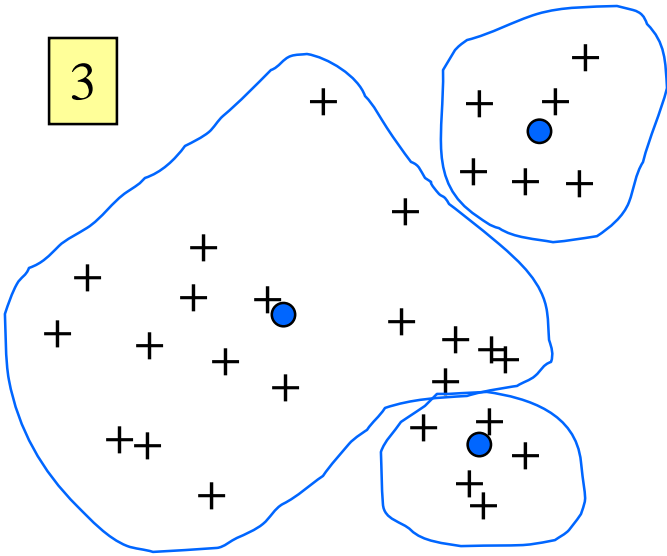
1



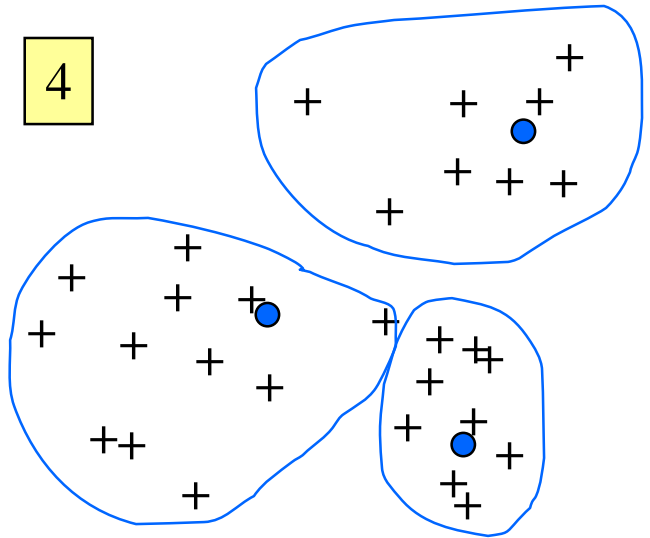
2



3



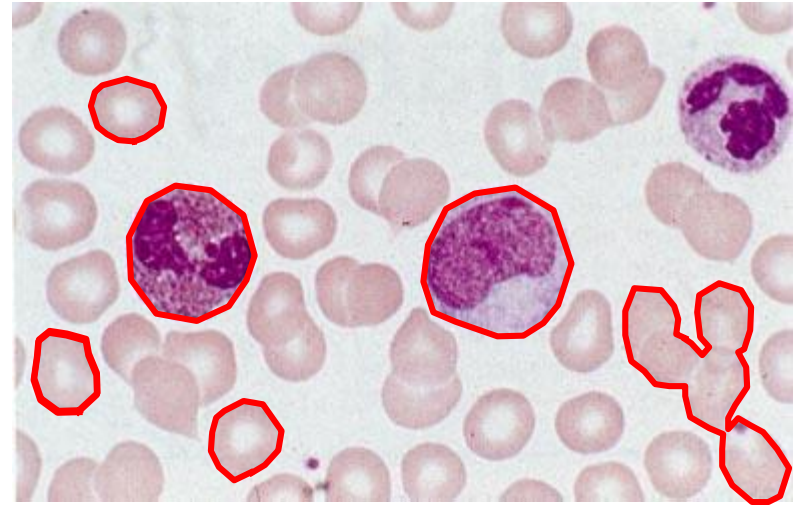
4



Classification/detection example

Task: Automatically detect
abnormal white blood cells

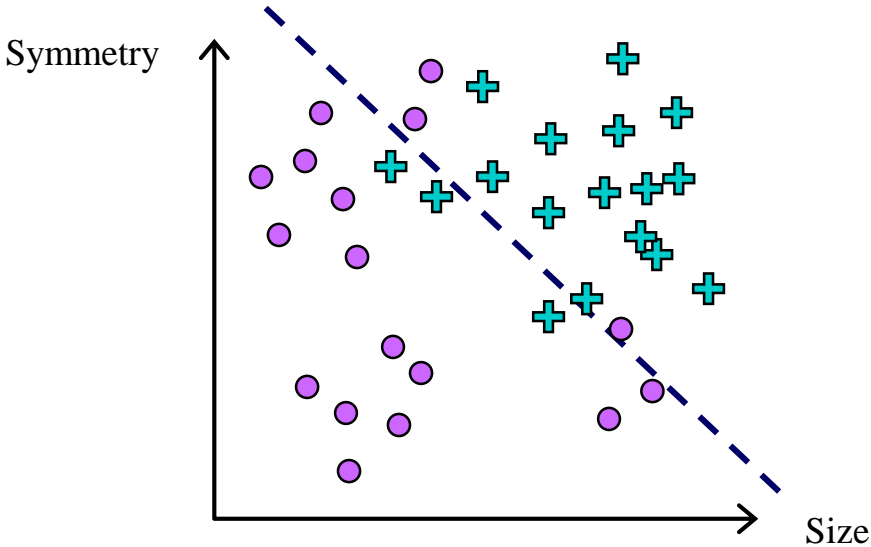
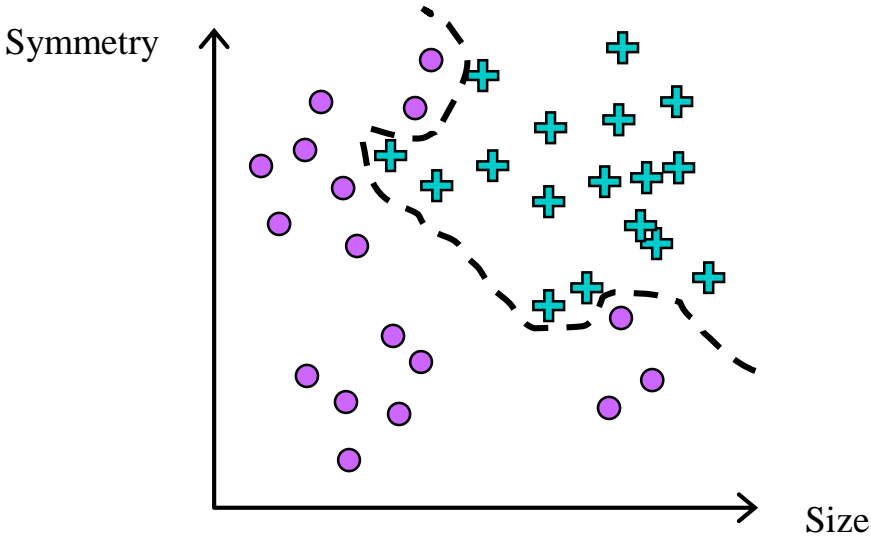
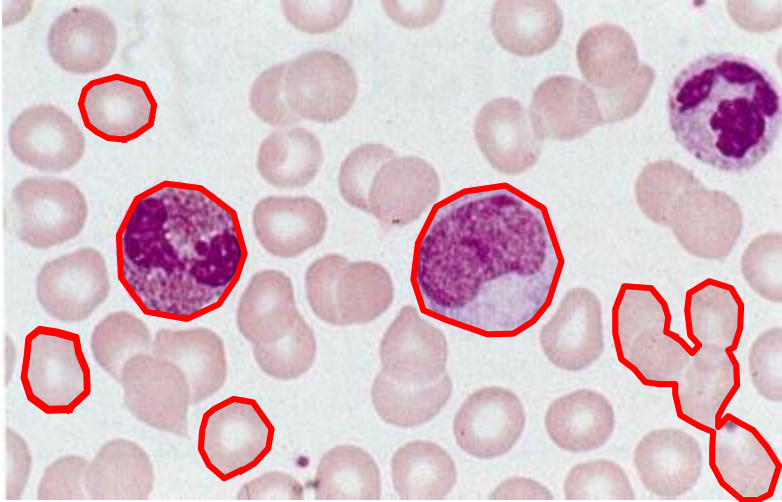
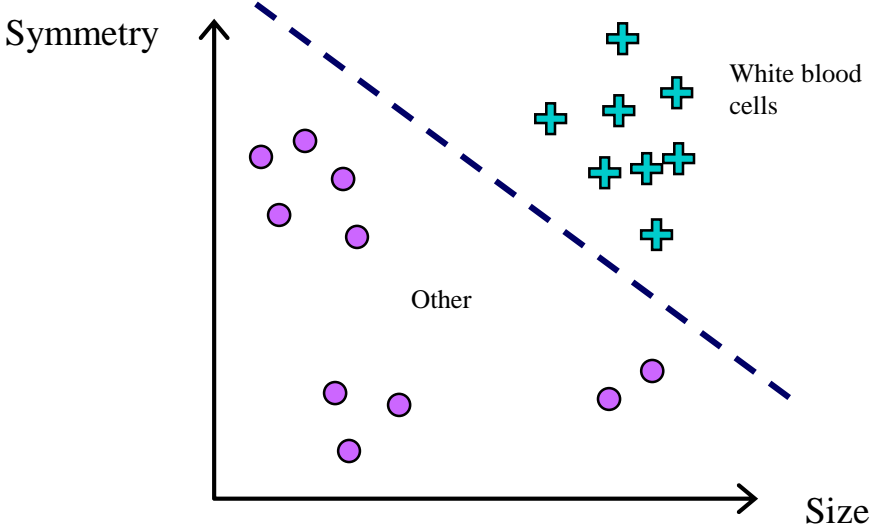
1. Process images to find outlines
2. Count white blood cells
3. Classify abnormal white blood cells



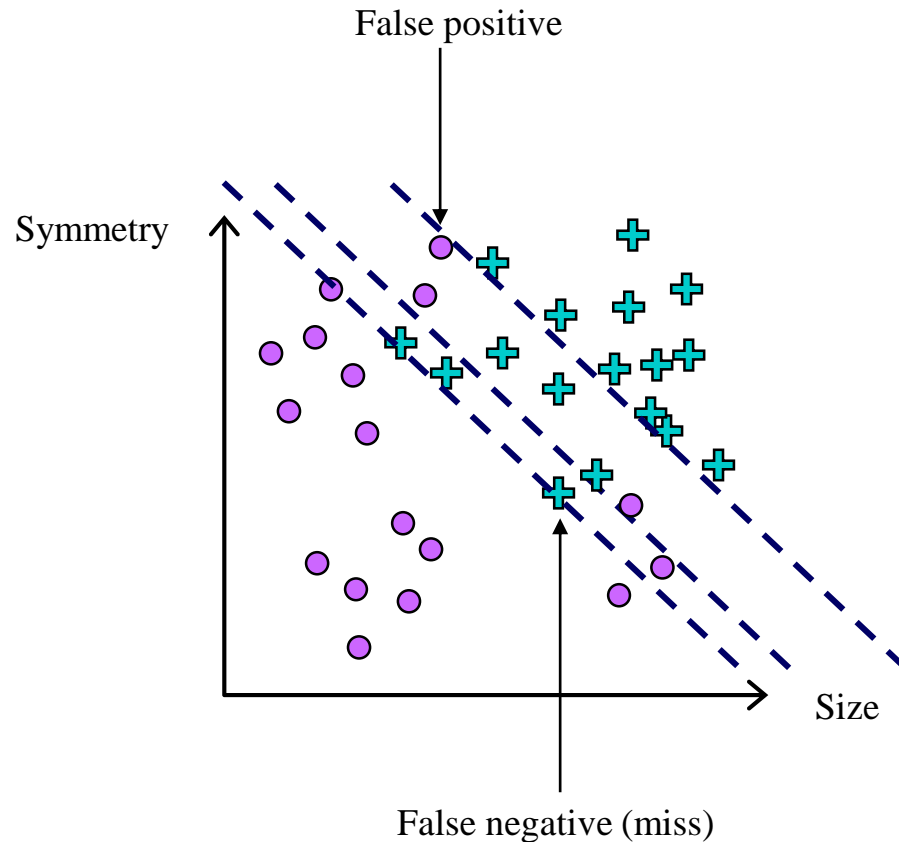
Steps 2 and 3 require *training* – teaching the system how to distinguish between white blood cells and others, and between normal and abnormal white blood cells

It's very important to choose good, discriminating features

Classification/detection example



Where to place the boundary?



We wish to minimize false positives and false negatives

Training vs. Testing

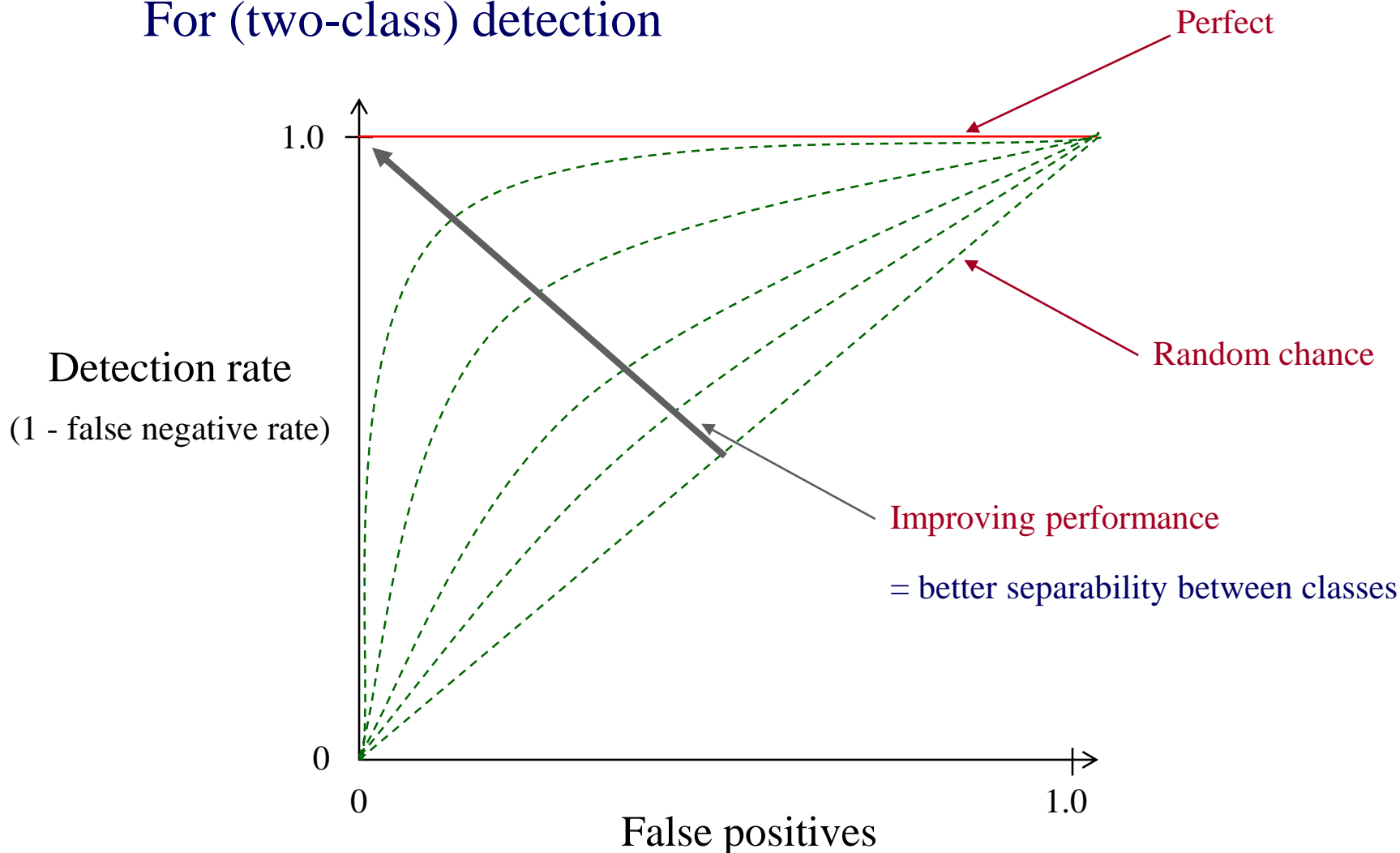
The *training set* (known examples) should be representative of the *testing set* (real data).

Good performance on your training set alone is meaningless...!

In every experiment, keep the *training* and *testing* data separate.

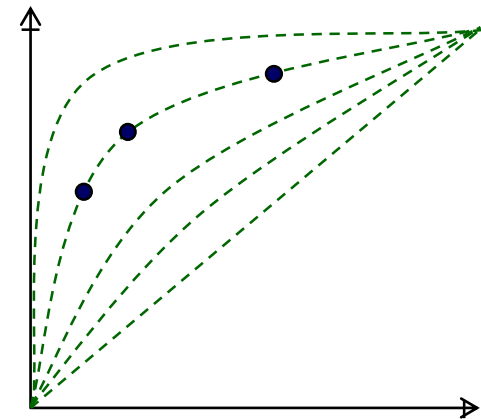
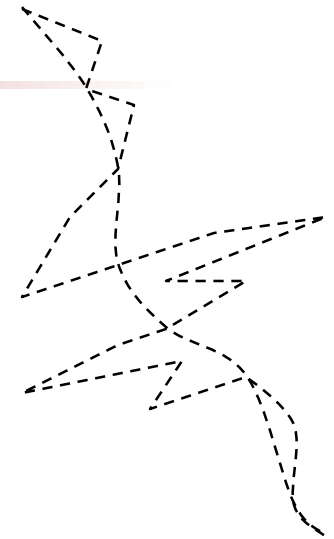
Evaluating performance – the ROC curve

For (two-class) detection

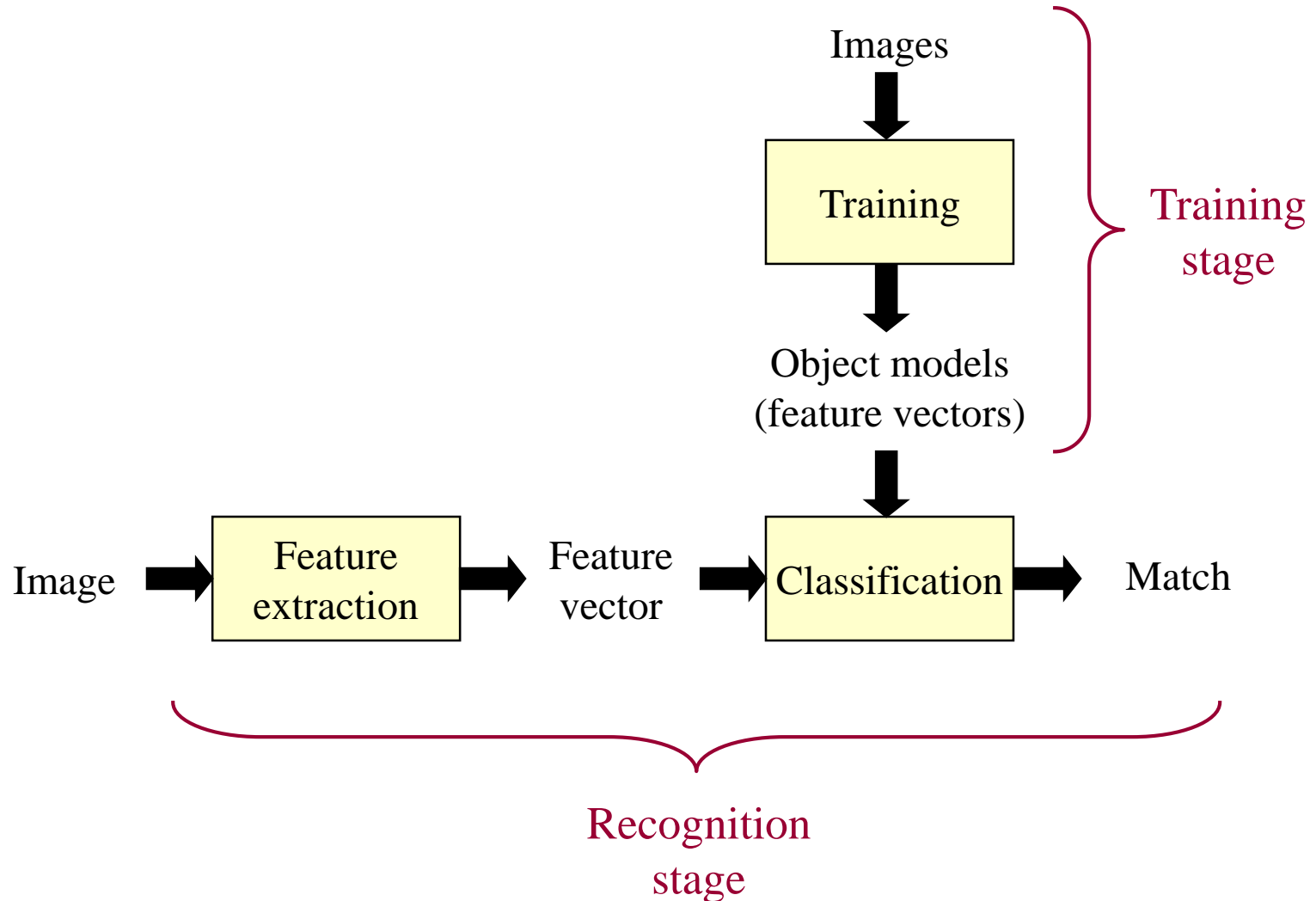


Improving classification/detection

- More training samples
 - So classification strategy is more general
- Don't "overtrain"
 - Don't want to "learn the noise" – keep it simple
- Use better features
 - Good features lead to good class separation
- Don't confuse movement *along* the ROC curve with *improving* the ROC curve



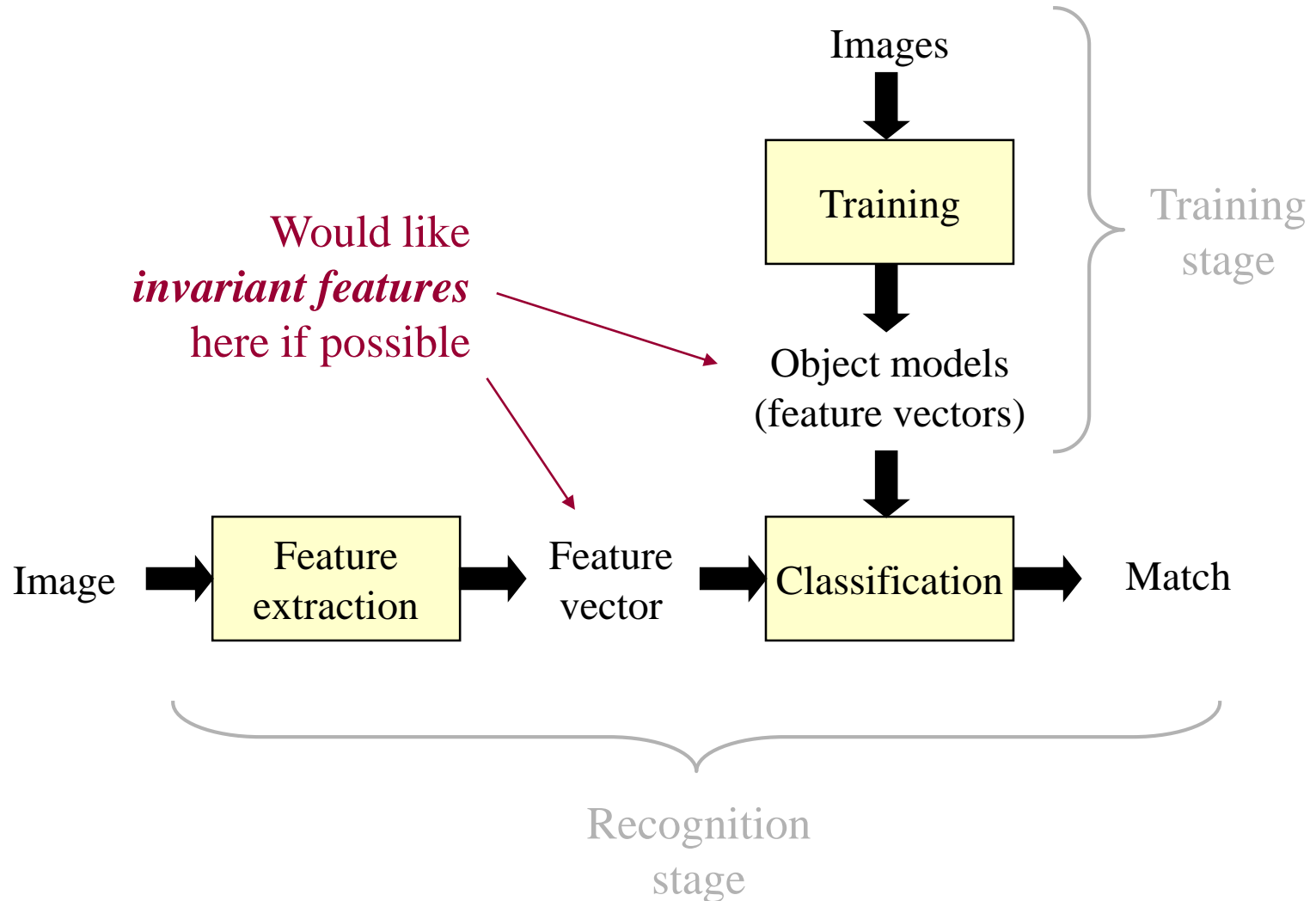
General approach to recognition



Invariant features for recognition

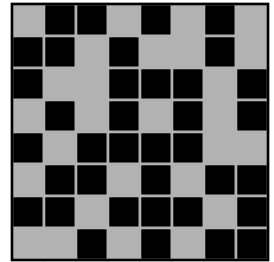
- An *invariant* feature is one that does not change under a certain class of transformations
 - Lengths and angles are invariant under rigid motion
 - Normalized correlation is invariant under scaling of image intensities
 - Brightness/color of a Lambertian surface is invariant under rotation
 - Length and angles *in the image* are **not** invariant under out-of-plane rotation and translation
 - Etc....
- Invariants can be geometric (location, shape) or radiometric (image values)
 - Geometric invariants tend to be much more common and useful

General approach to recognition

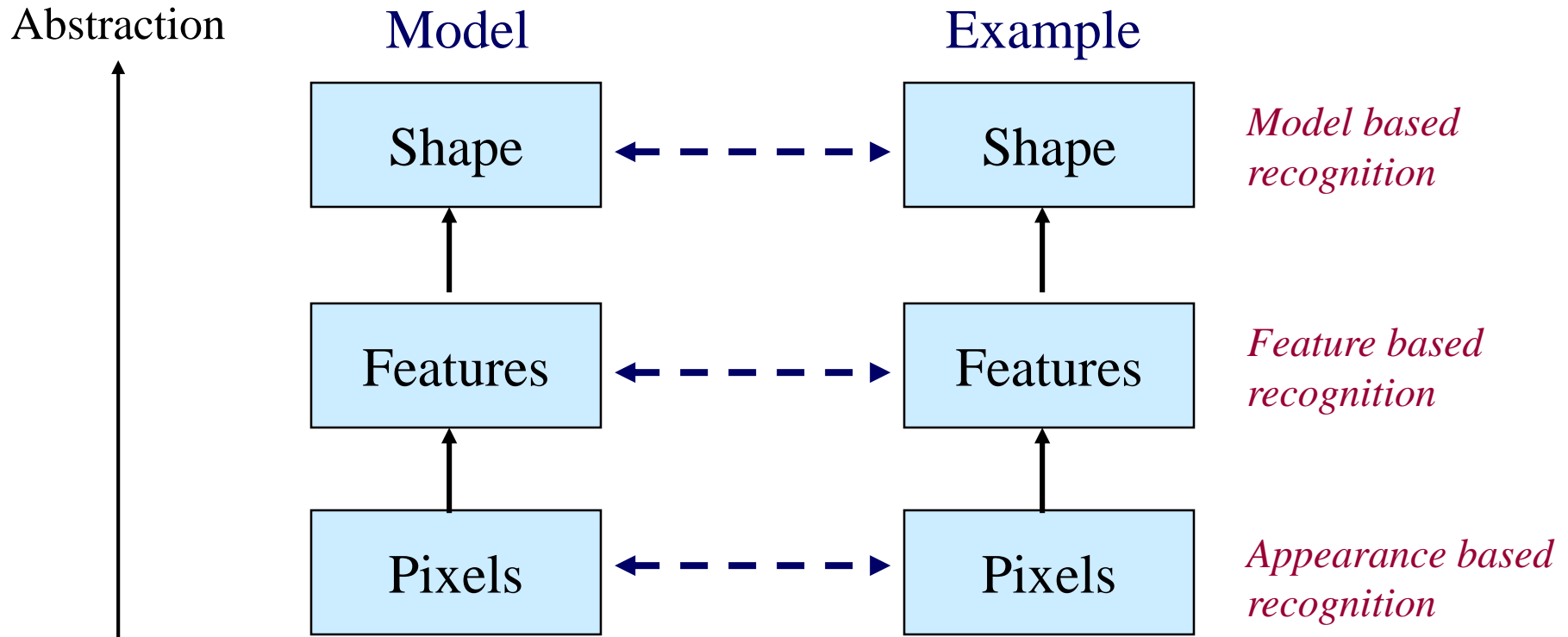


Why not this for face recognition/detection?

1. Ahead of time, search over all possible images to see which ones look like my face, and save these
 2. During recognition, see if the input image is one of these
- Image space is vastly large
 - 8x8 binary image $\rightarrow 2^{64}$ image points (distinct images)
 - 1 billion images per second \rightarrow **600 years**
 - Step #1 would never finish!
 - Not to mention, we'd have to do this for every possible view of my face
 - Range of facial expressions, lighting conditions, poses, etc.



Levels of Recognition/Matching



Reminder: Why computer vision is so hard!

These are all images
of Simon's face!



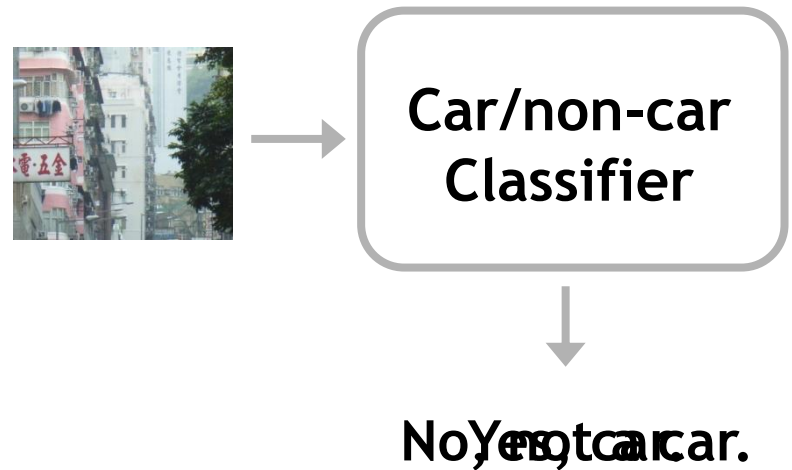
In general, object recognition is difficult because of the immense variability of object appearance. With faces, this is even worse!

How can we reliably detect/recognize Simon???

Scanning windows...

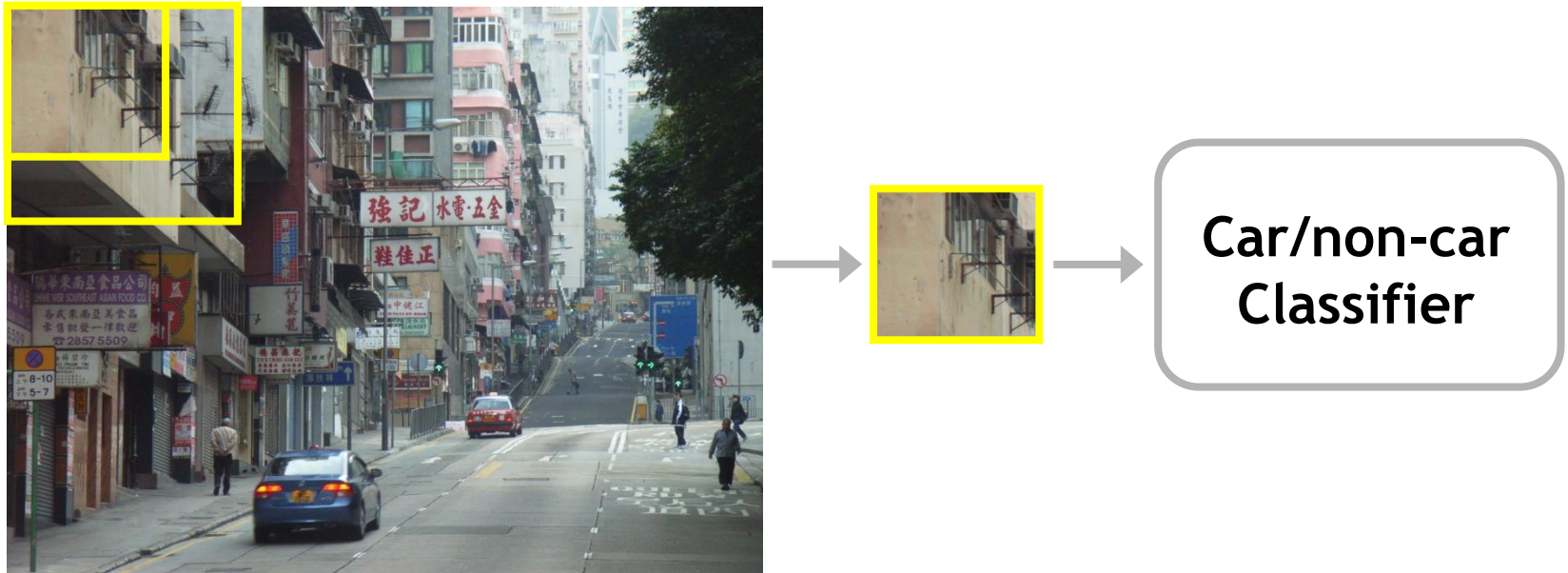
Detection via classification: Main idea

Basic component: a binary classifier



Detection via classification: Main idea

If object may be in a cluttered scene, slide a window around looking for it.

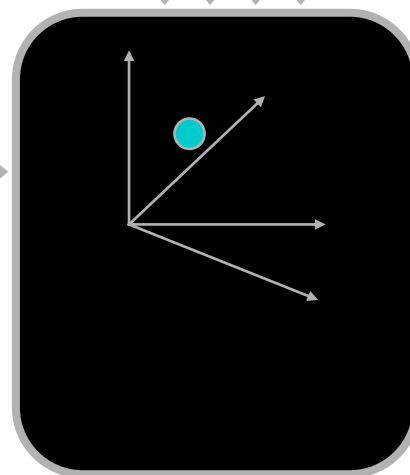
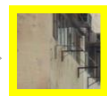


(Essentially, our skin detector was doing this, with a window that was one pixel big.)

Detection via classification: Main idea

Fleshing out this pipeline a bit more, we need to:

1. Obtain training data
2. Define features
3. Define classifier

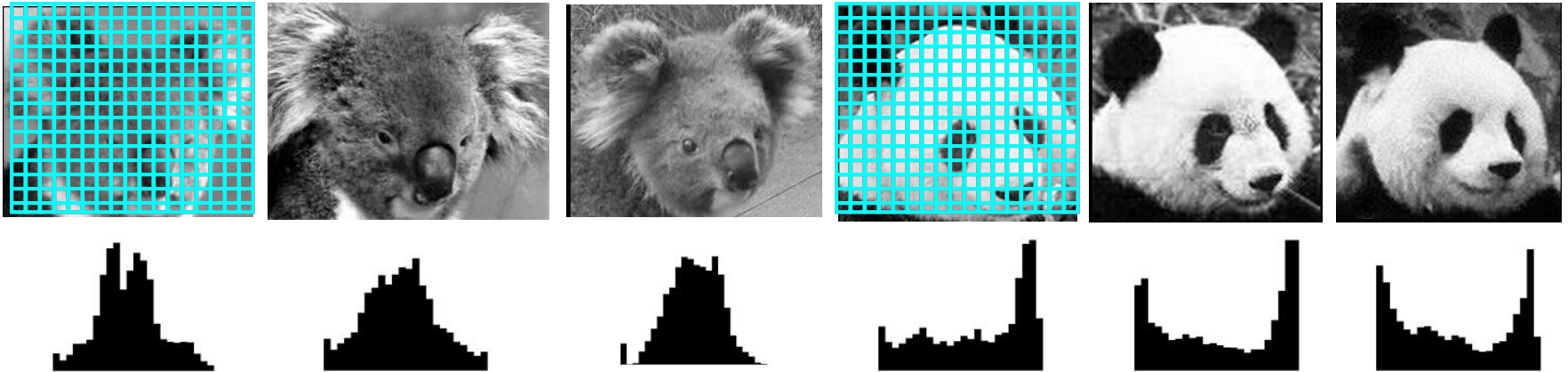
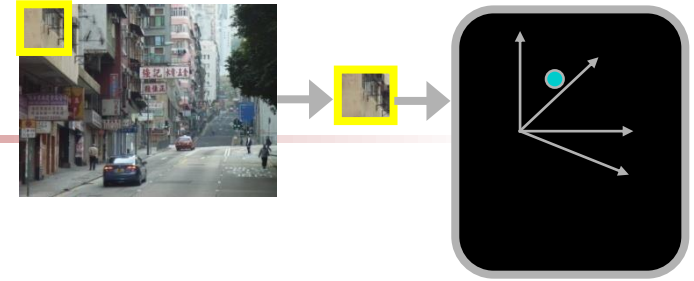


Car/non-car
Classifier

Detection via classification: Main idea

- Consider all subwindows in an image
 - Sample at multiple scales and positions (and orientations)
- Make a decision per window:
 - “Does this contain object category X or not?”

Feature extraction: global appearance



Simple holistic descriptions of image content

- grayscale / color histogram
- vector of pixel intensities

Eigenfaces: global appearance description

An early appearance-based approach to face recognition



Training images



Mean

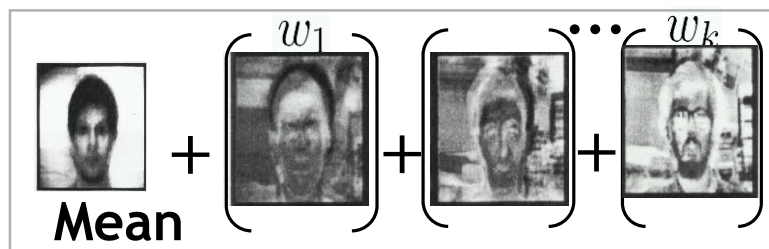


Eigenvectors computed from covariance matrix

Generate low-dimensional representation of appearance with a linear subspace.



\approx

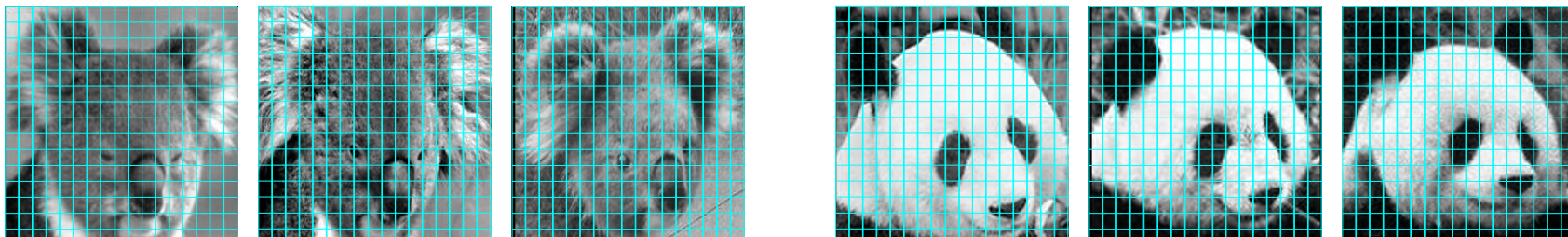


Project new images to “face space”.

Recognition via nearest neighbors in face space

Feature extraction: global appearance

- Pixel-based representations sensitive to small shifts



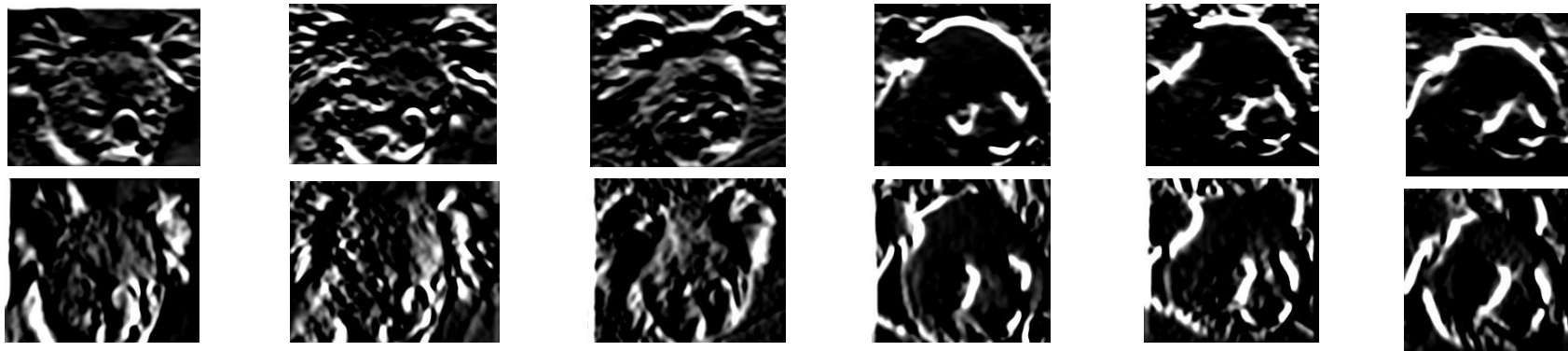
- Color or grayscale-based appearance description can be sensitive to illumination and intra-class appearance variation



**Cartoon example:
an albino koala**

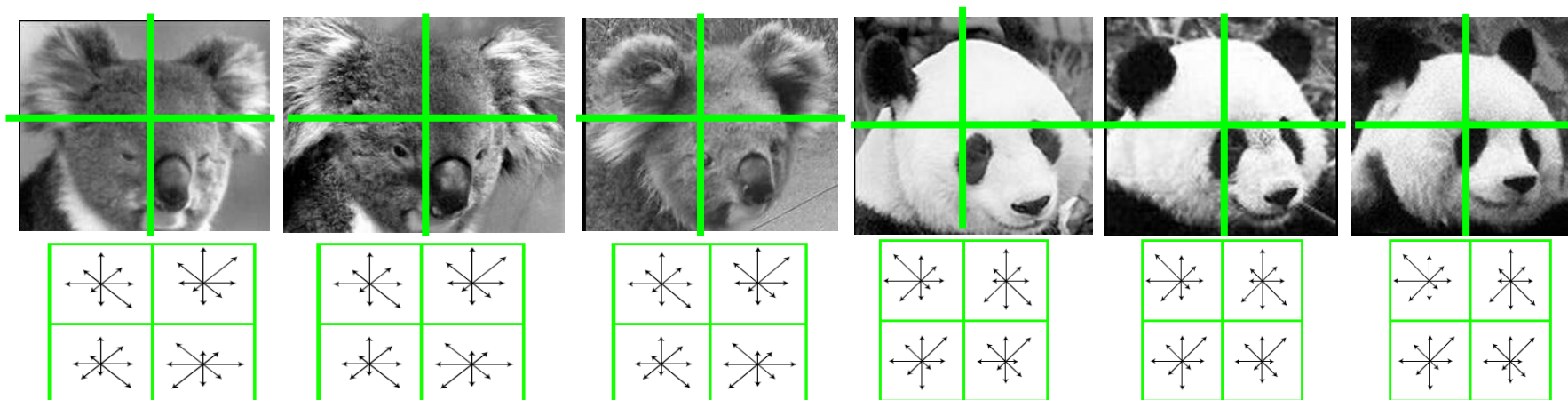
Gradient-based representations

- Consider edges, contours, and (oriented) intensity gradients



Gradient-based representations

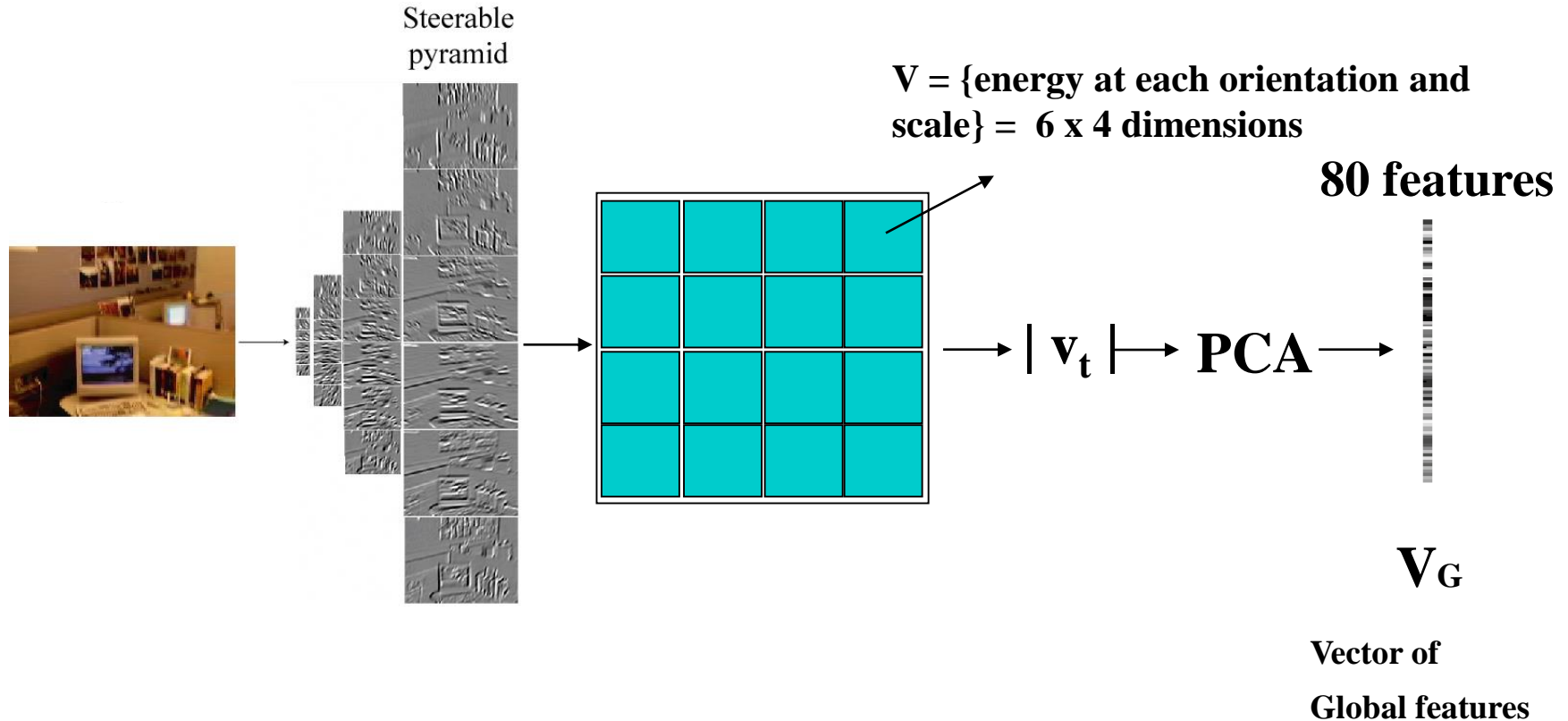
- Consider edges, contours, and (oriented) intensity gradients



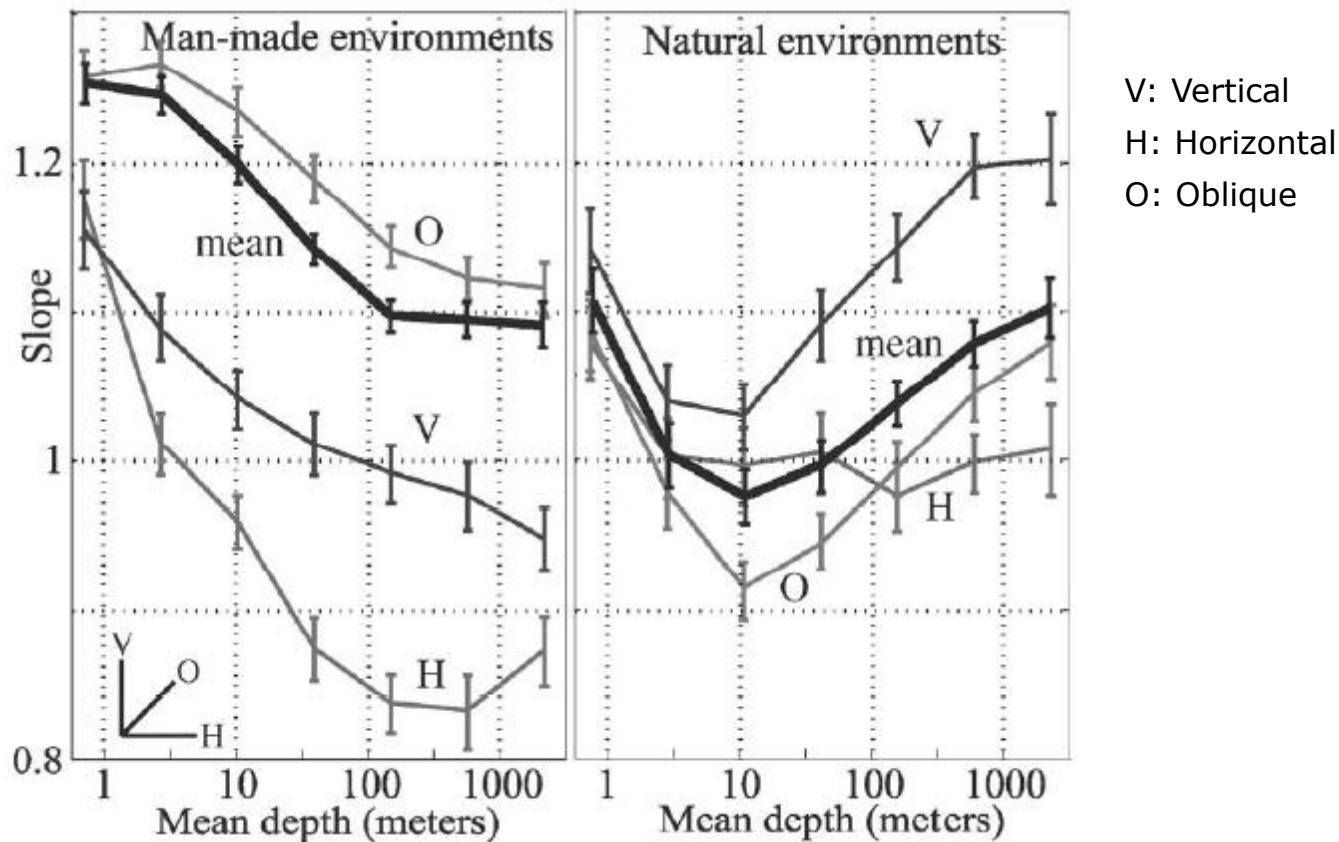
- Summarize local distribution of gradients with histogram
 - Locally orderless: offers invariance to small shifts and rotations
 - Contrast-normalization: try to correct for variable illumination

GIST

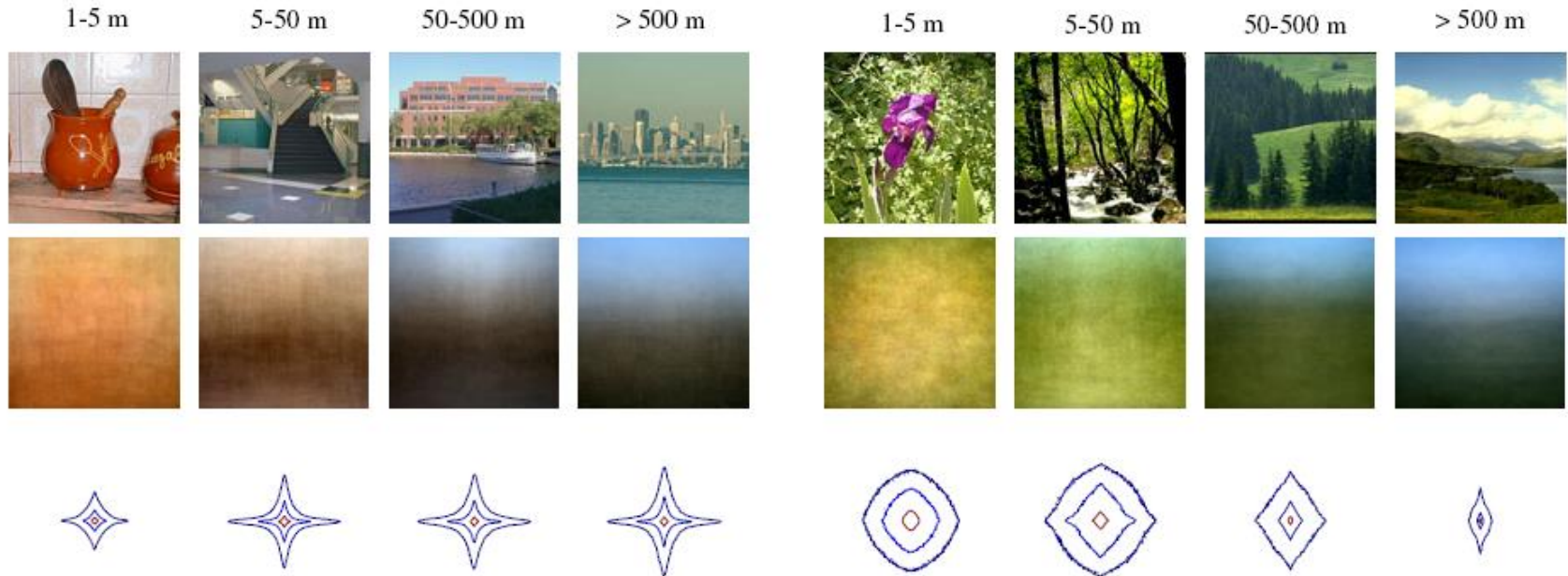
Representing Image Structure with “GIST”



What do Images Statistics say about Depth?



Scene Scale



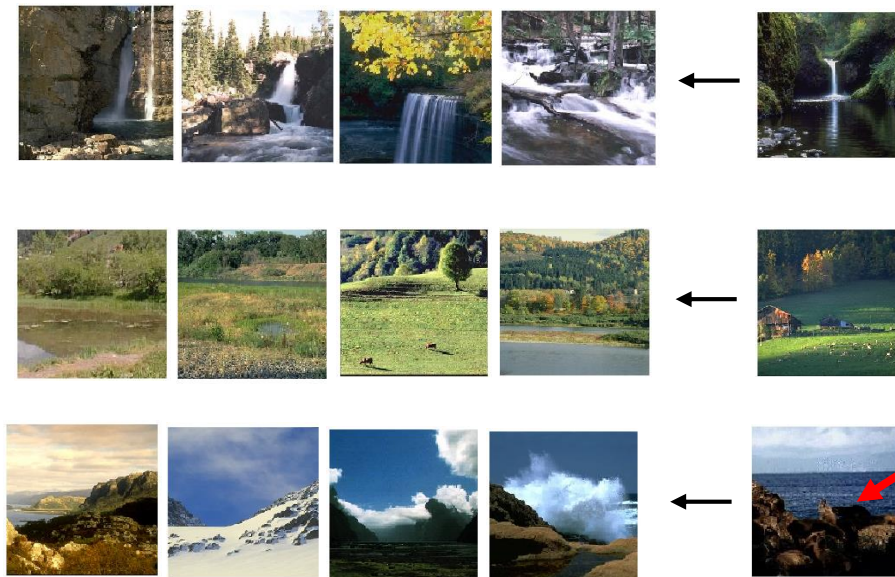
- “The point of view that any given observer adopts on a specific scene is constrained by the volume of the scene.”
- How does the amount of clutter vary against scene scale in man-made environments? In natural environments?

Categorization of Natural Scenes

Confusion Matrix (in % using Layout template) :

Classification of prototypical scenes (400 / category)

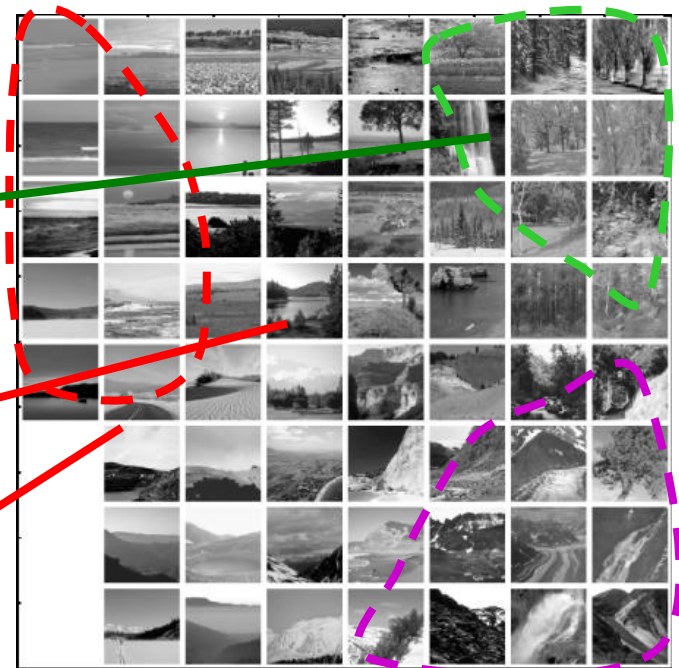
	Coast	Countryside	Forest	Mountain
Coast	88.6	8.9	1.2	1.3
Countryside	9.8	85.2	3.7	1.3
Forest	0.4	3.6	91.5	4.5
Mountain	0.4	4.6	3.8	91.2



Local organization:

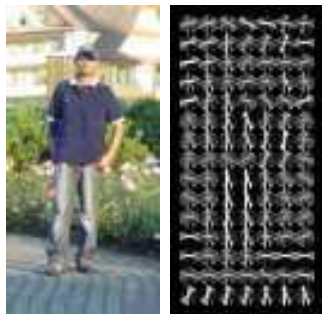
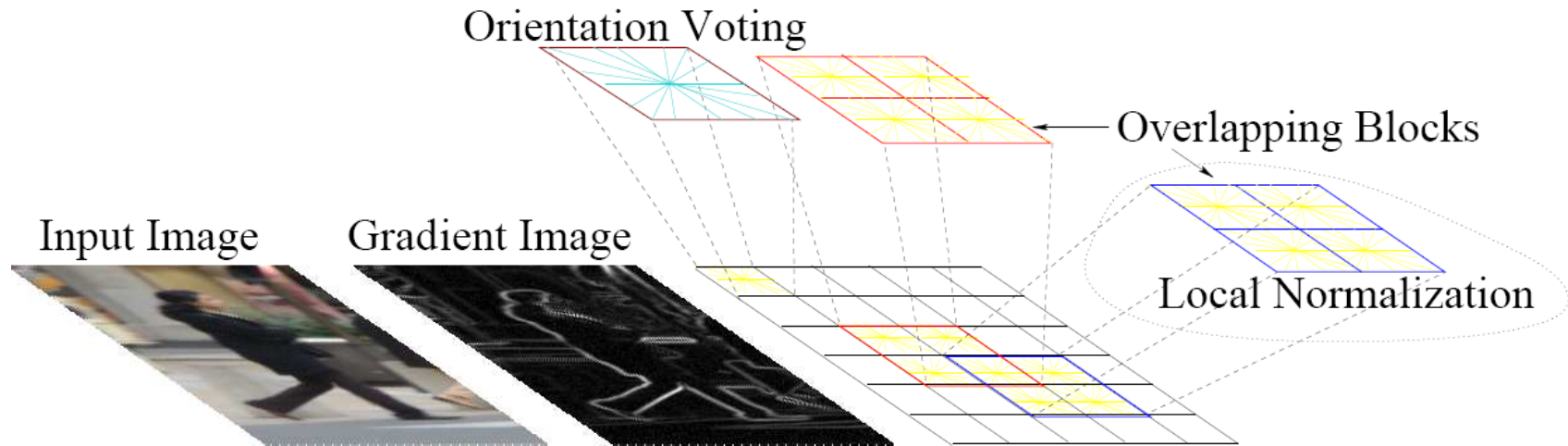
correct for 92 % images

(4 similar images on 7 K-NN)



HOG

Gradient-based representations: Histograms of oriented gradients (HoG)



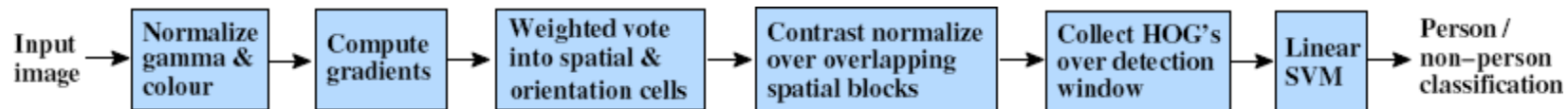
Map each grid cell in the input window to a histogram counting the gradients per orientation.

Code available:

<http://pascal.inrialpes.fr/soft/olt/>



Slide credit: Dalal, Triggs, P. Barnum



Slide credit: Dalal, Triggs, P. Barnum



- Tested with
 - RGB
 - LAB
 - Grayscale
- Gamma Normalization and Compression
 - Square root
 - Log



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

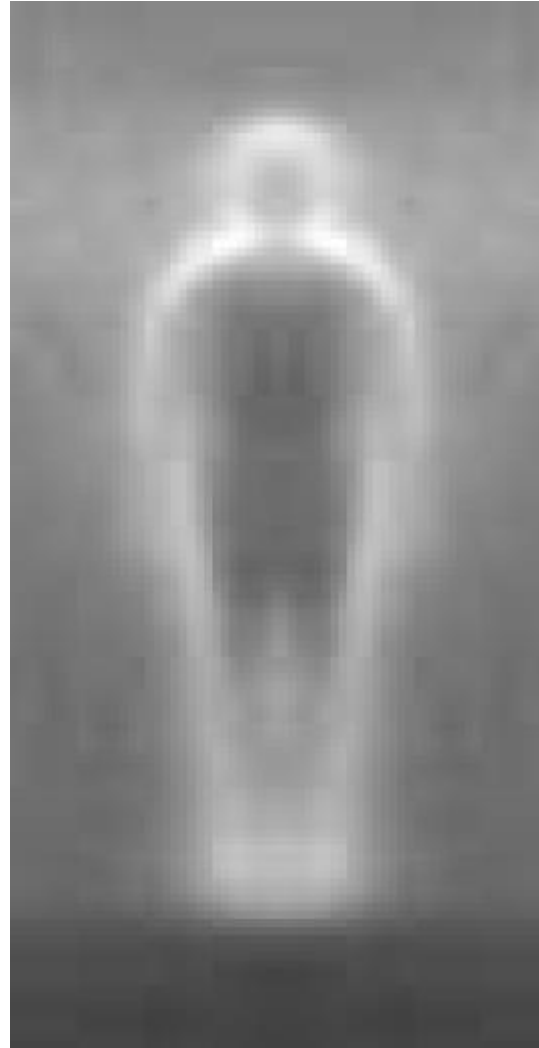
centered

-1	1
----	---

uncentered

1	-8	0	8	-1
---	----	---	---	----

cubic-corrected



0	1
-1	0

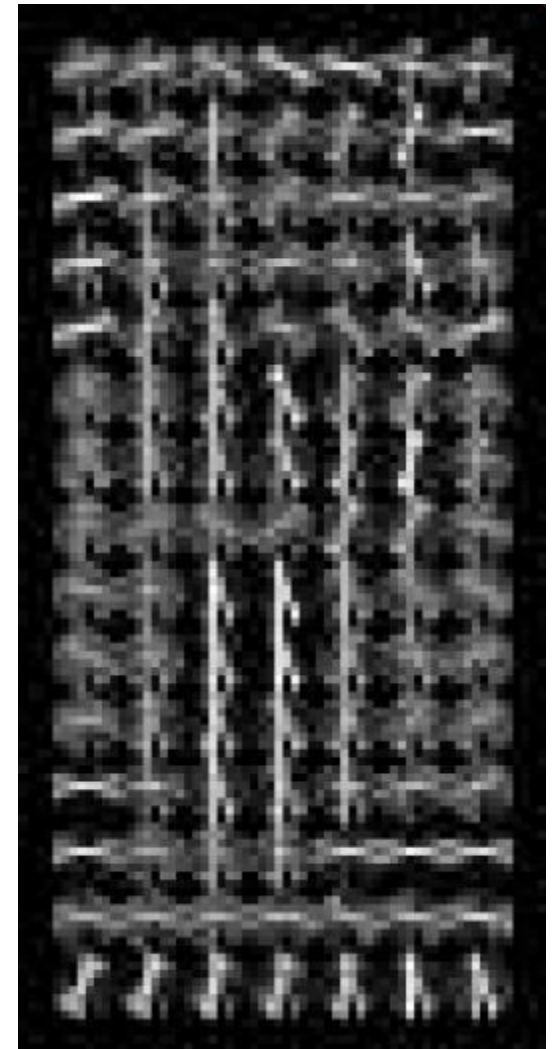
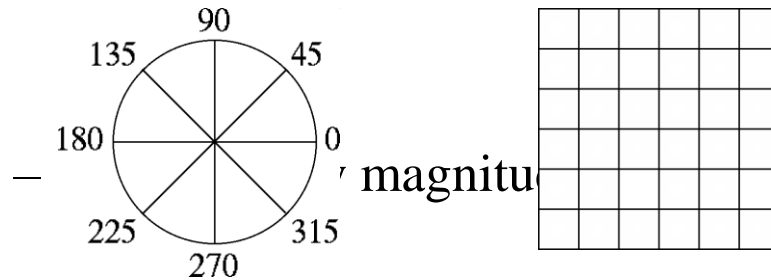
diagonal

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

Sobel

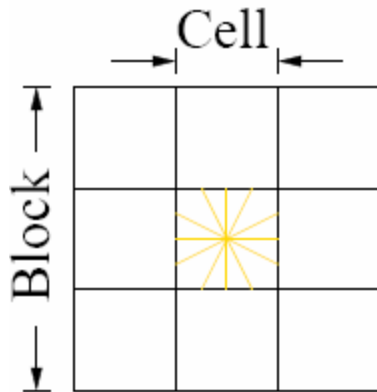


- Histogram of gradient orientations
 - Orientation
 - Position

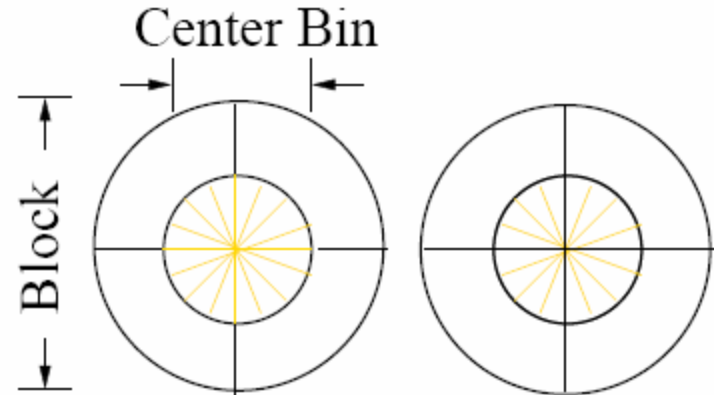




R-HOG



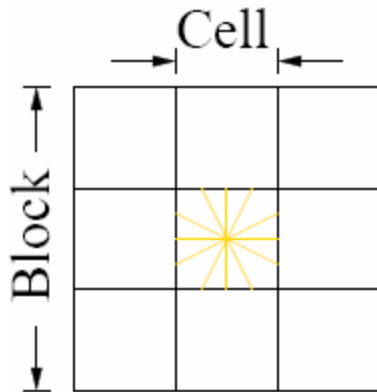
C-HOG



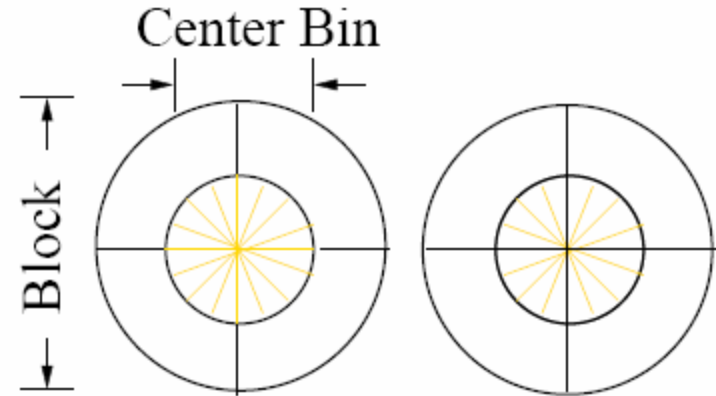
Radial Bins, Angular Bins



R-HOG



C-HOG



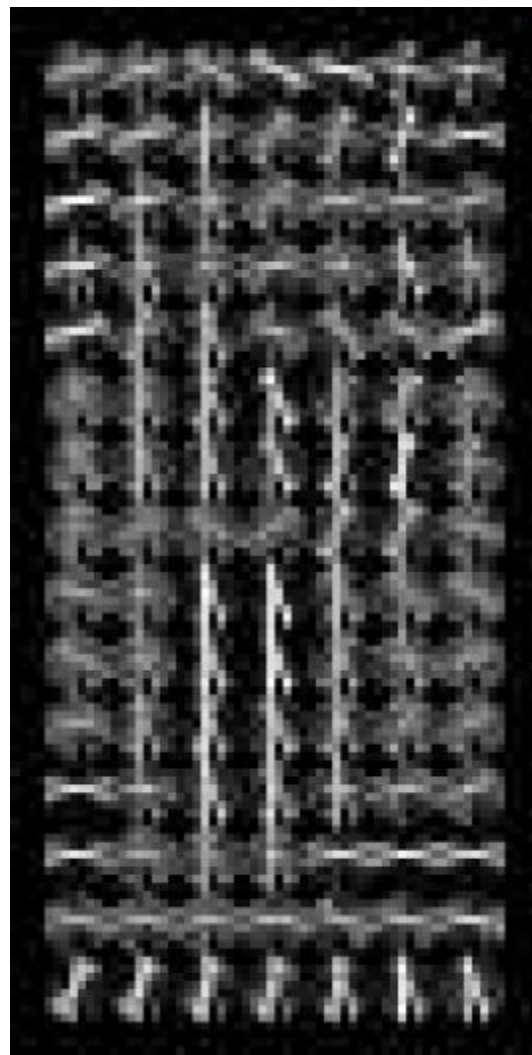
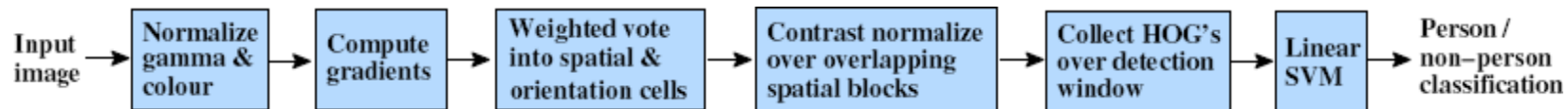
Radial Bins, Angular Bins

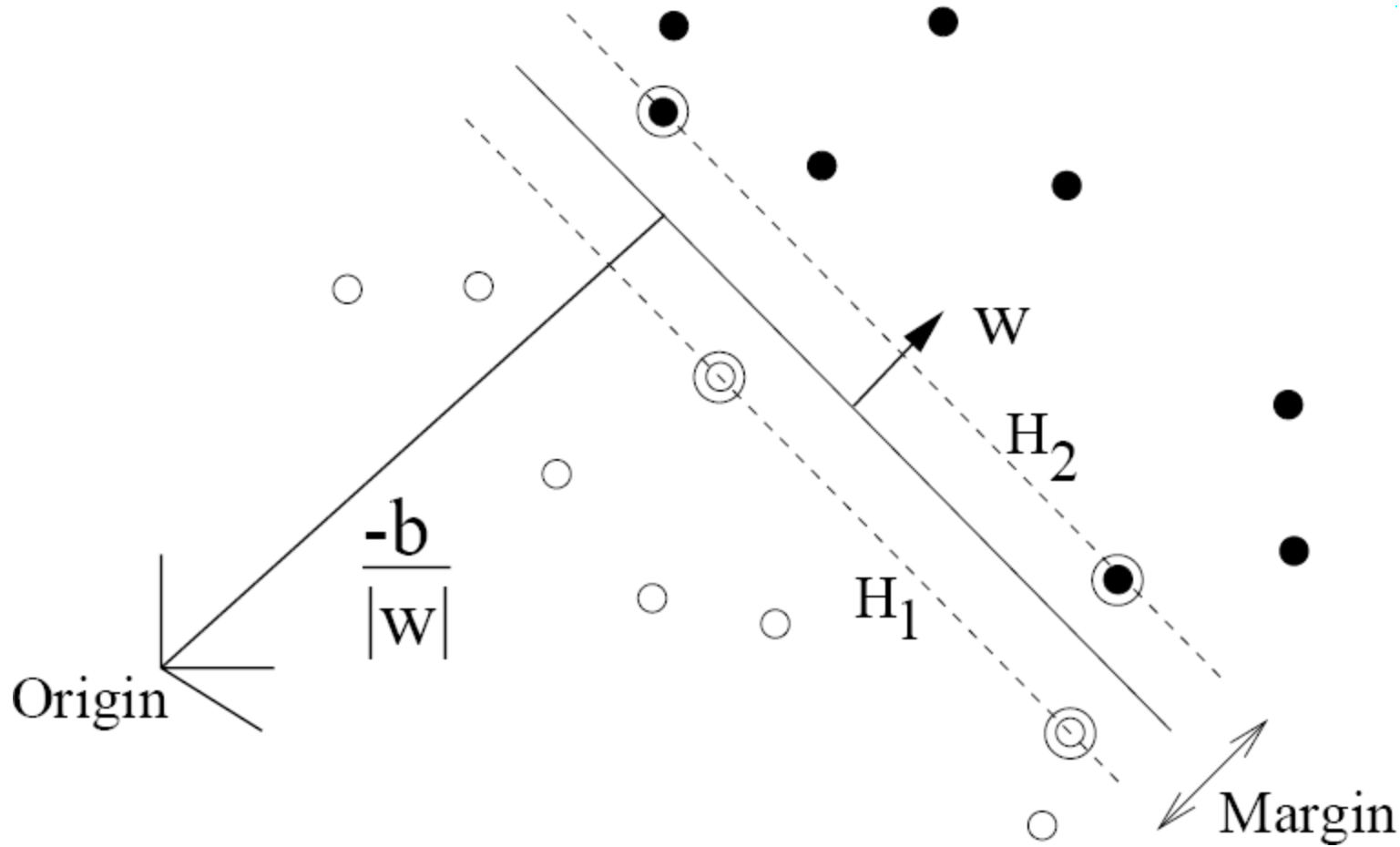
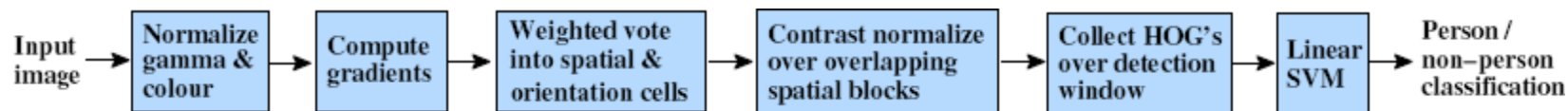
$$L1 - norm : v \longrightarrow v / (\|v\|_1 + \epsilon)$$

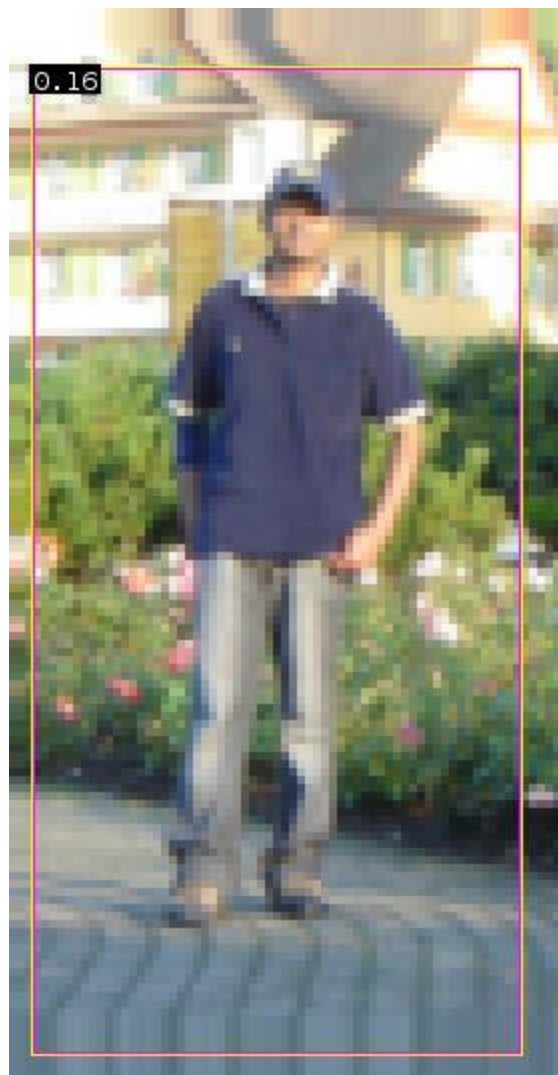
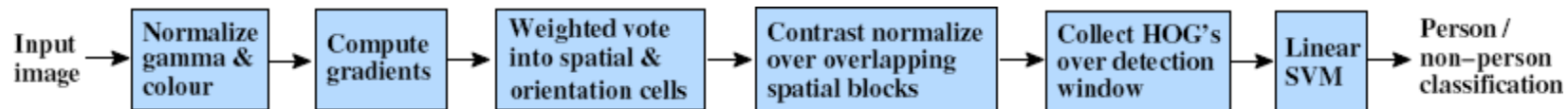
$$L1 - sqrt : v \longrightarrow \sqrt{v / (\|v\|_1 + \epsilon)}$$

$$L2 - norm : v \longrightarrow v / \sqrt{\|v\|_2^2 + \epsilon^2}$$

$L2 - hys$: L2-norm, plus clipping at .2 and renormalizing







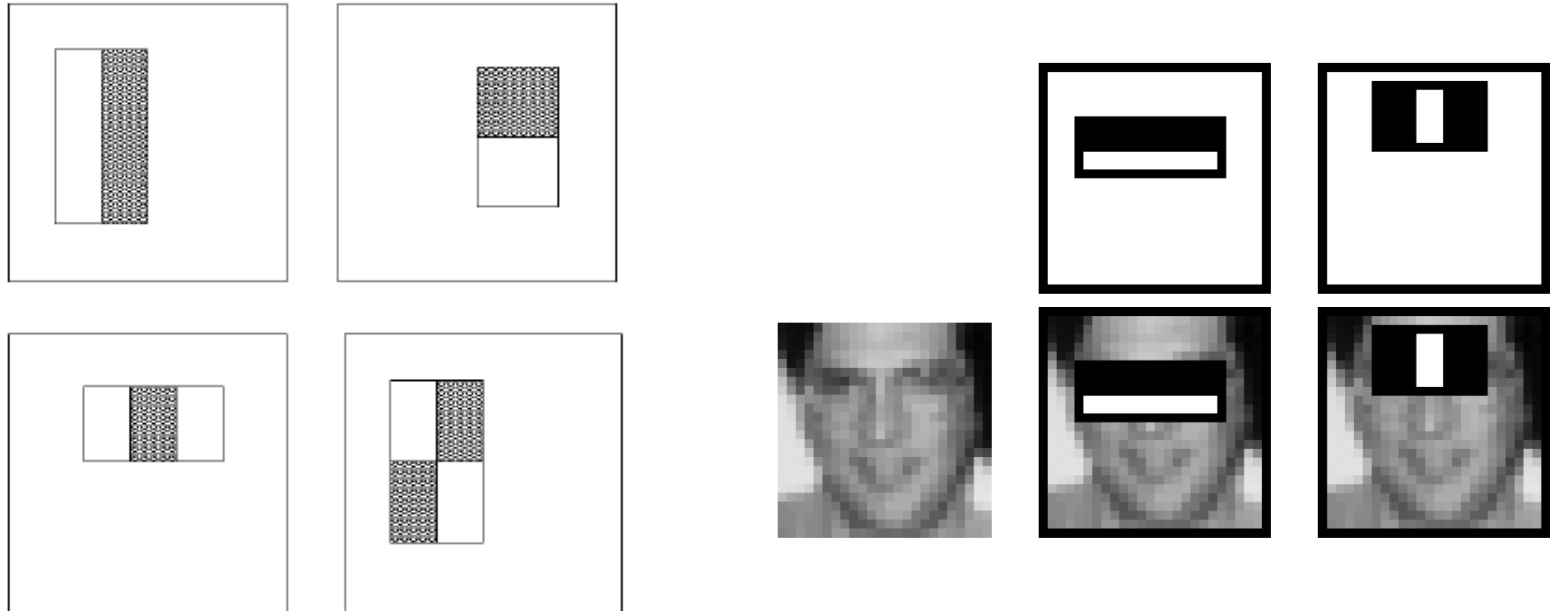
Slide credit: Dalal, Triggs, P. Barnum



Slide credit: Dalal, Triggs, P. Barnum

Boosted Face Detection with Gradient Features

Gradient-based representations: Rectangular features



Compute differences between sums of pixels in rectangles

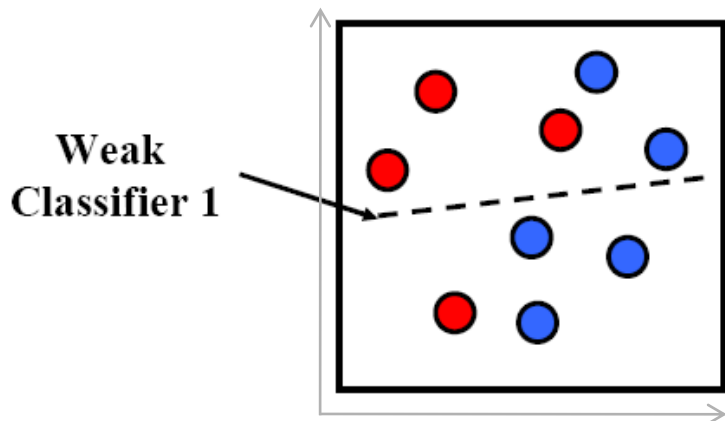
Captures contrast in adjacent spatial regions, efficient to compute

Each feature parameterized by scale, position, type.

Boosting

- Build a strong classifier by combining number of “weak classifiers”, which need only be better than chance
- Sequential learning process: at each iteration, add a weak classifier
- Flexible to choice of weak learner
 - including fast simple classifiers that alone may be inaccurate
- We’ll look at Freund & Schapire’s AdaBoost algorithm
 - Easy to implement
 - Base learning algorithm for Viola-Jones face detector

AdaBoost: Intuition

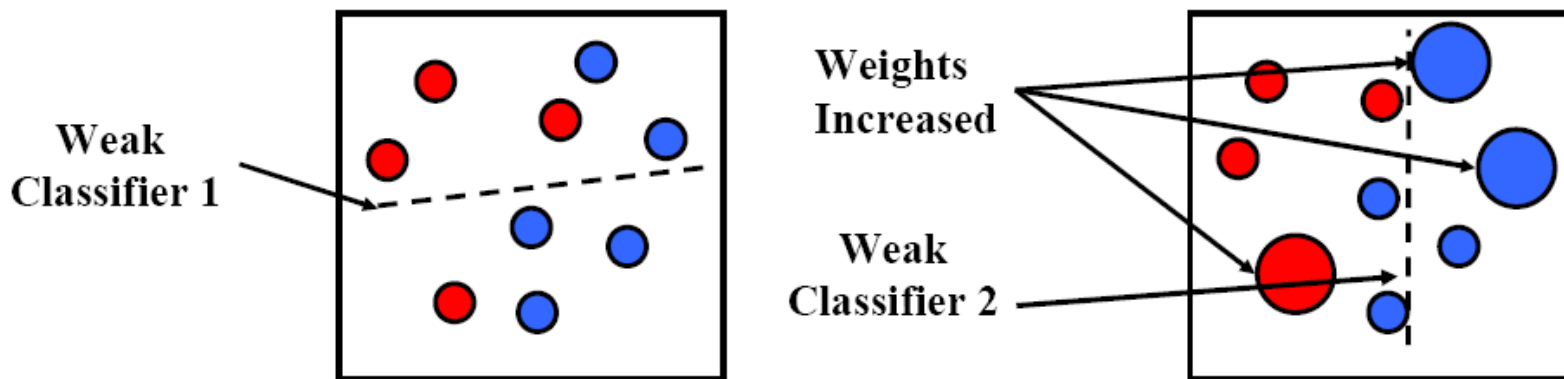


Consider a 2-d feature space with **positive** and **negative** examples.

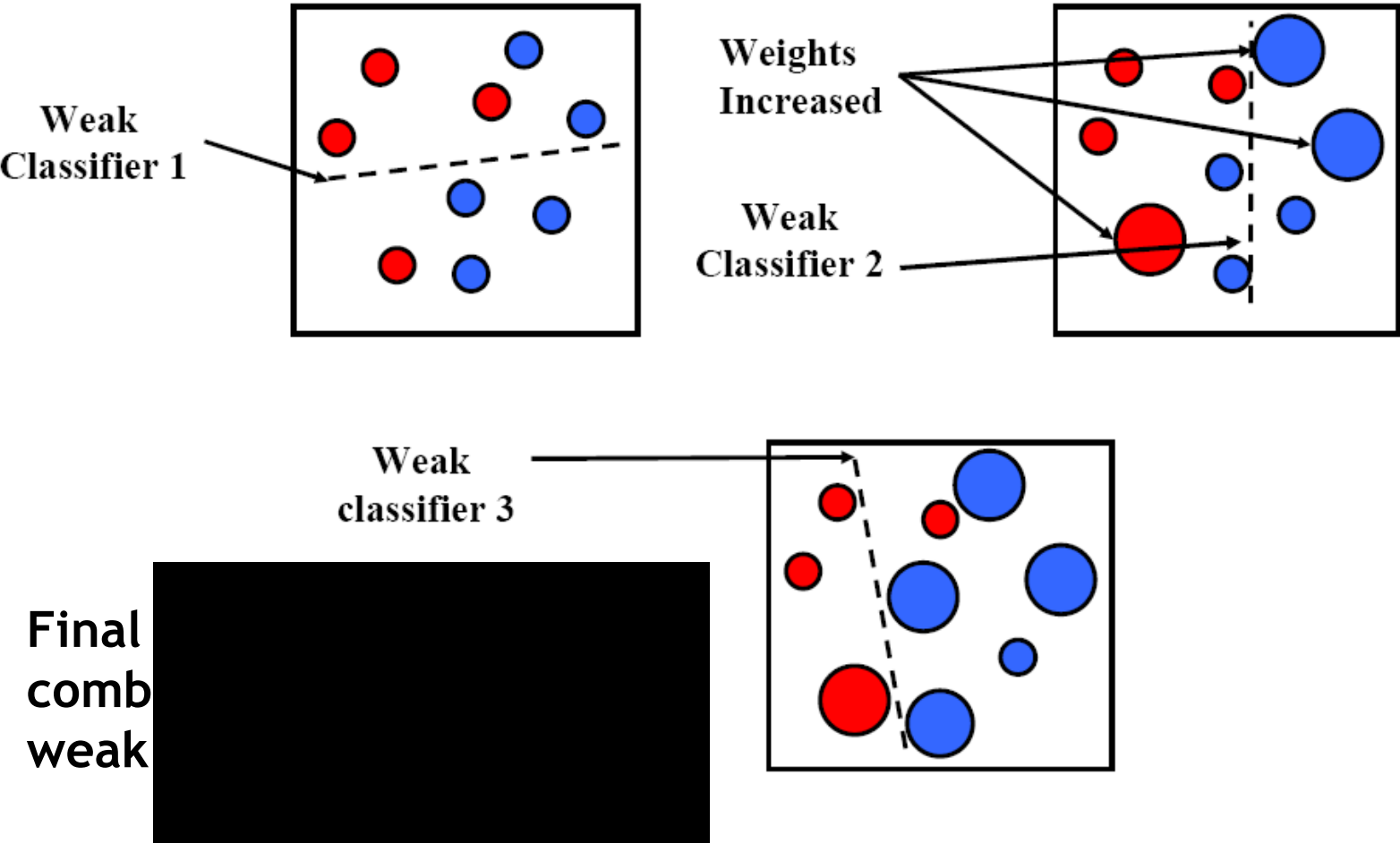
Each weak classifier splits the training examples with at least 50% accuracy.

Examples misclassified by a previous weak learner are given more emphasis at future rounds.

AdaBoost: Intuition



AdaBoost: Intuition



AdaBoost Algorithm

- Given example images $(x_1, y_1), \dots, (x_n, y_n)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{1,i} = \frac{1}{2m}, \frac{1}{2l}$ for $y_i = 0, 1$ respectively, where m and l are the number of negatives and positives respectively.
- For $t = 1, \dots, T$:

1. Normalize the weights,

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$$

so that w_t is a probability distribution.

2. For each feature, j , train a classifier h_j which is restricted to using a single feature. The error is evaluated with respect to w_t , $\epsilon_j = \sum_i w_i |h_j(x_i) - y_i|$.
3. Choose the classifier, h_t , with the lowest error ϵ_t .
4. Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_t^{1-e_i}$$

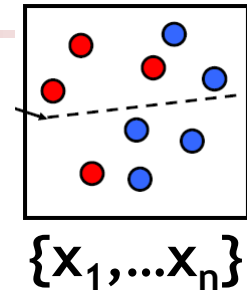
where $e_i = 0$ if example x_i is classified correctly, $e_i = 1$ otherwise, and $\beta_t = \frac{\epsilon_t}{1-\epsilon_t}$.

- The final strong classifier is:

$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

where $\alpha_t = \log \frac{1}{\beta_t}$

Start with
uniform
weights on
training
examples
For T rounds



Evaluate
weighted error
for each
feature, pick
best.

Re-weight the examples:

Incorrectly classified -> more weight

Correctly classified -> less weight

Final classifier is combination of
the weak ones, weighted
according to error they had.

Freund & Schapire 1999

Example: Face detection

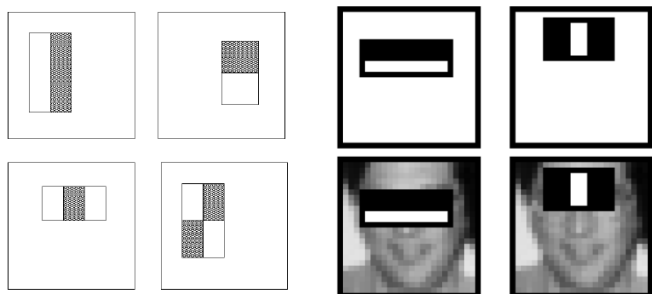
- Frontal faces are a good example of a class where global appearance models + a sliding window detection approach fit well:
 - Regular 2D structure
 - Center of face almost shaped like a “patch”/window



- Now we'll take AdaBoost and see how the Viola-Jones face detector works

Feature extraction

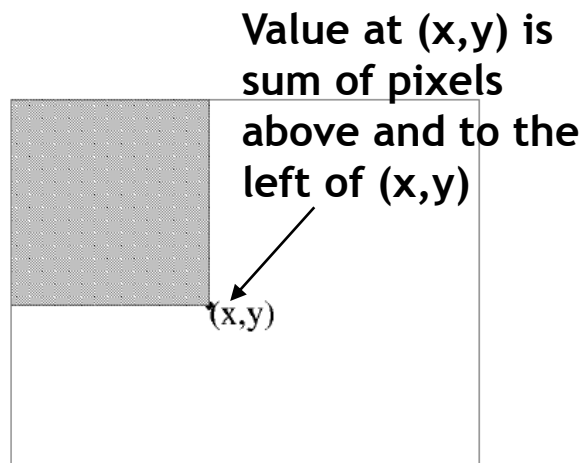
“Rectangular” filters



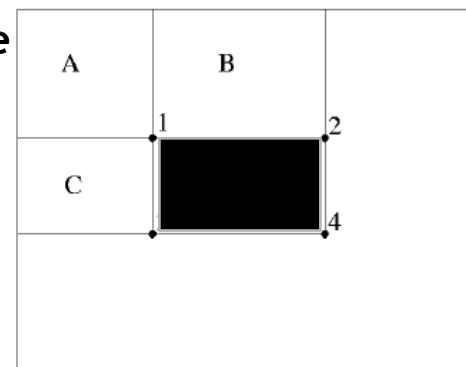
Feature output is difference between adjacent regions

Efficiently computable with integral image: any sum can be computed in constant time

Avoid scaling images → scale features directly for same cost

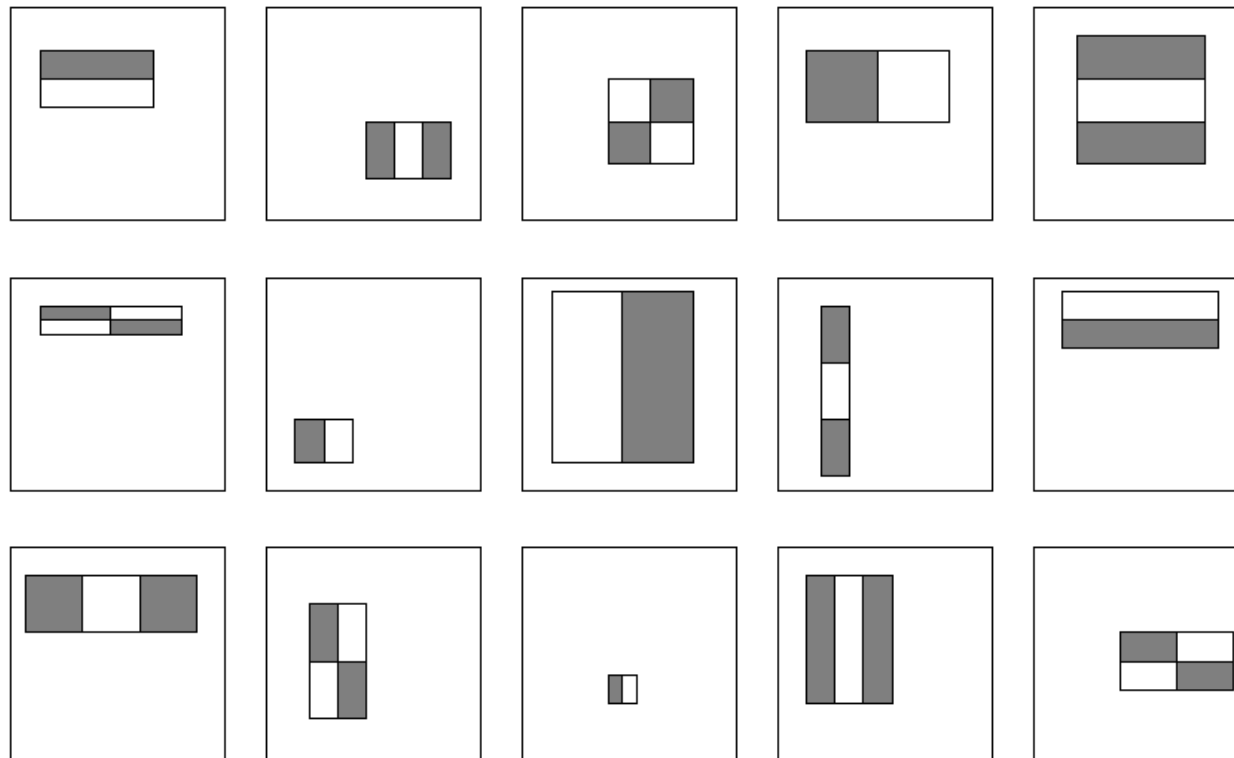


Integral image



$$\begin{aligned} D &= 1 + 4 - (2 + 3) \\ &= A + (A + B + C + D) - (A + C + A + B) \\ &= D \end{aligned}$$

Large library of filters

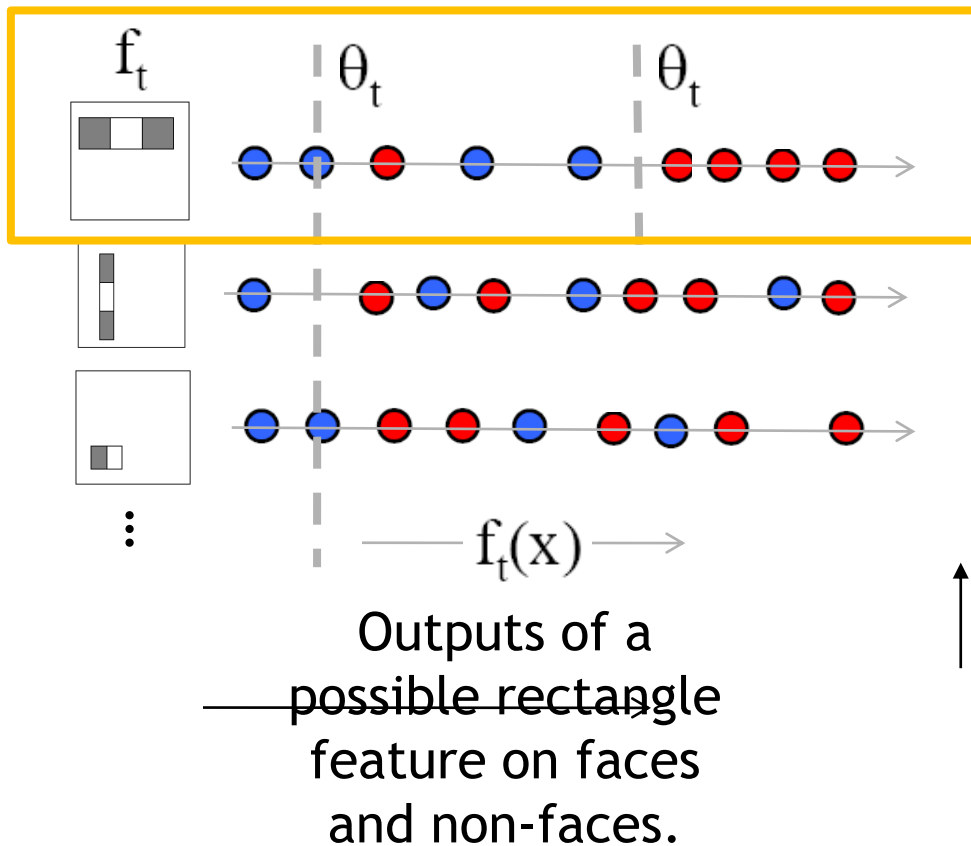


Considering all possible filter parameters:
position, scale, and type:
180,000+ possible features associated with each 24 x 24 window

Use AdaBoost both to select the informative features and to form the classifier

AdaBoost for feature+classifier selection

- Want to select the single rectangle feature and threshold that best separates **positive** (faces) and (non-faces) training examples, in terms of *weighted* error.



Resulting weak classifier:

$$h_t(x) = \begin{cases} +1 & \text{if } f_t(x) > \theta_t \\ -1 & \text{otherwise} \end{cases}$$

For next round, reweight the examples according to errors, choose another filter/threshold combo.

AdaBoost Algorithm

- Given example images $(x_1, y_1), \dots, (x_n, y_n)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{1,i} = \frac{1}{2m}, \frac{1}{2l}$ for $y_i = 0, 1$ respectively, where m and l are the number of negatives and positives respectively.
- For $t = 1, \dots, T$:

1. Normalize the weights,

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$$

so that w_t is a probability distribution.

- For each feature, j , train a classifier h_j which is restricted to using a single feature. The error is evaluated with respect to w_t , $\epsilon_j = \sum_i w_i |h_j(x_i) - y_i|$.
- Choose the classifier, h_t , with the lowest error ϵ_t .
- Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_t^{1-e_i}$$

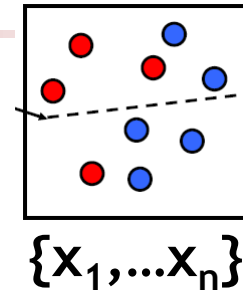
where $e_i = 0$ if example x_i is classified correctly. $e_i = 1$ otherwise. and $\beta_t = \frac{1-\epsilon_t}{1+\epsilon_t}$.

- The final strong classifier is:

$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

where $\alpha_t = \log \frac{1}{\beta_t}$

Start with
uniform
weights on
training
examples
For T rounds



Evaluate
weighted error
for each
feature, pick
best.

Re-weight the examples:

Incorrectly classified -> more weight

Correctly classified -> less weight

Final classifier is combination of
the weak ones, weighted
according to error they had.

Freund & Schapire 1999

AdaBoost for Efficient Feature Selection

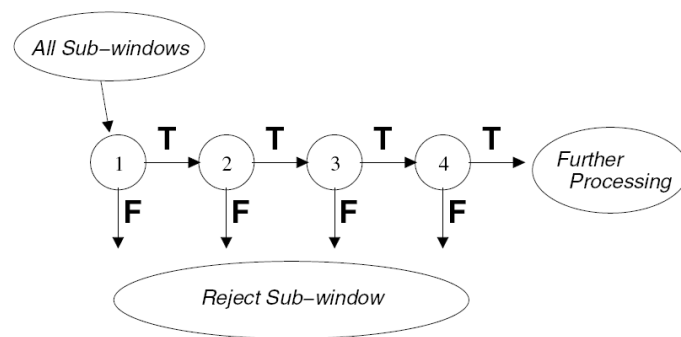
- Image Features = Weak Classifiers
- For each round of boosting:
 - Evaluate each rectangle filter on each example
 - Sort examples by filter values
 - Select best threshold for each filter (min error)
 - Sorted list can be quickly scanned for the optimal threshold
 - Select best filter/threshold combination
 - Weight on this feature is a simple function of error rate
 - Reweight examples

-
- Even if the filters are fast to compute, each new image has a lot of possible windows to search.
 - How to make the detection more efficient?

Cascading classifiers for detection

For efficiency, apply less accurate but faster classifiers first to immediately discard windows that clearly appear to be negative; e.g.,

- Filter for promising regions with an initial inexpensive classifier
- Build a chain of classifiers, choosing cheap ones with low false negative rates early in the chain

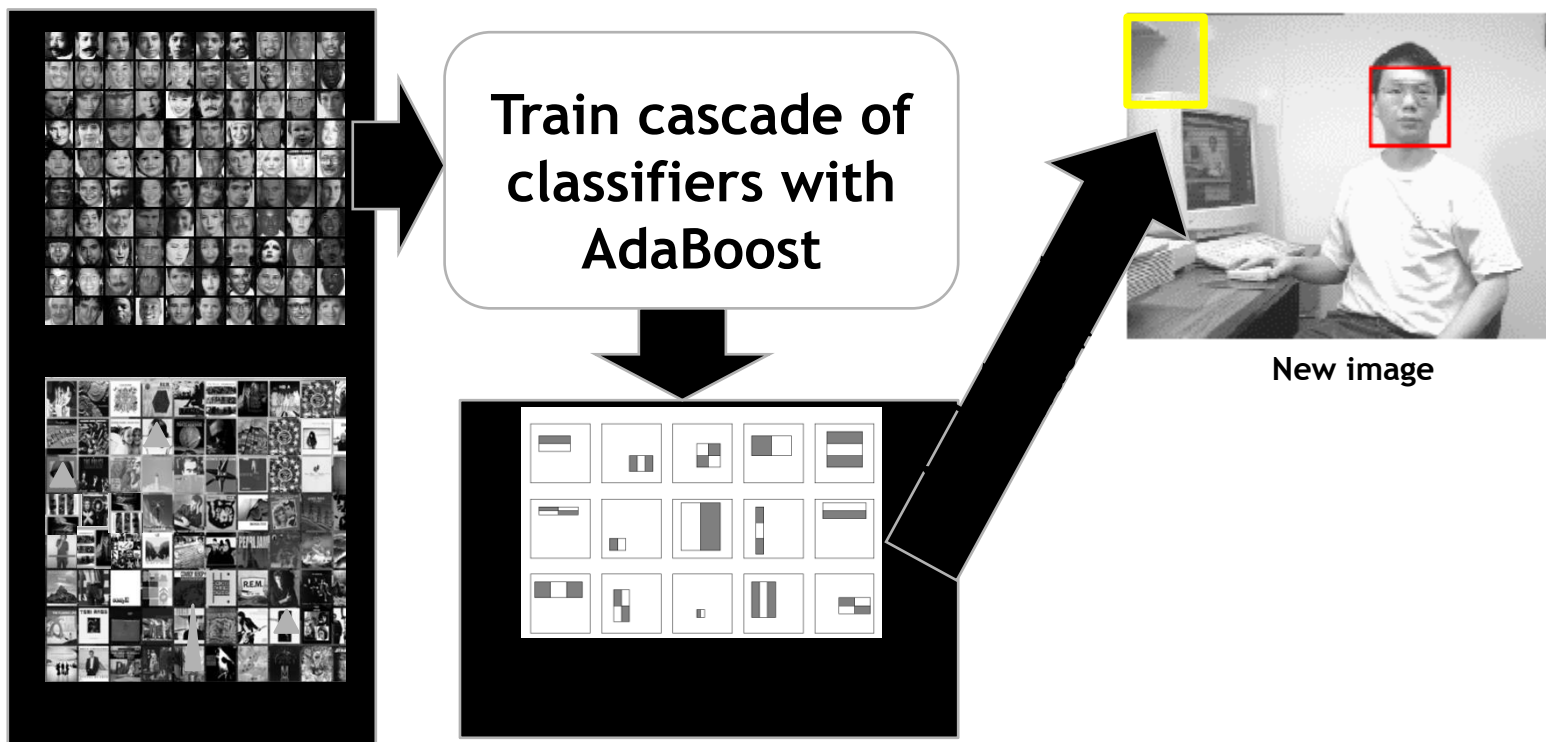


Fleuret & Geman, IJCV 2001

Rowley et al., PAMI 1998

Viola & Jones, CVPR 2001
W. Gorman, B. Leibe

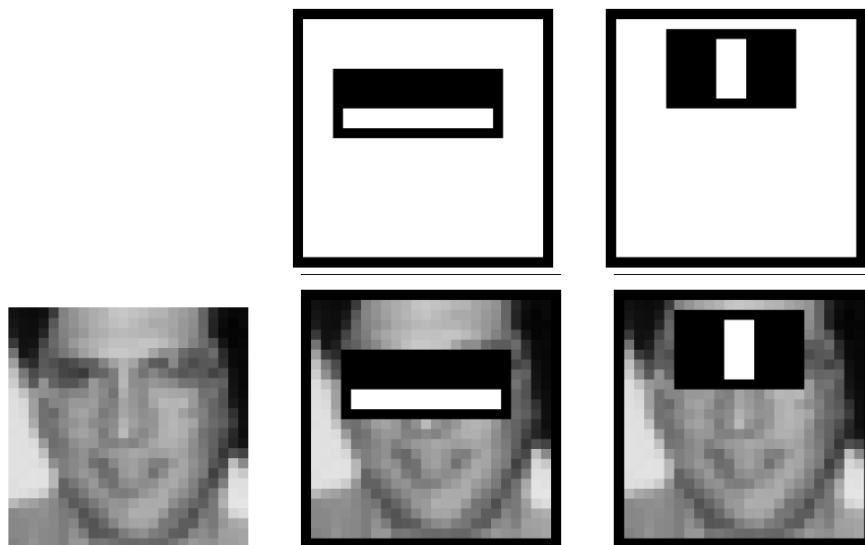
Viola-Jones Face Detector: Summary



- Train with 5K positives, 350M negatives
- Real-time detector using 38 layer cascade
- 6061 features in final layer
- [Implementation available in OpenCV:
<http://www.intel.com/technology/computing/opencv/>]

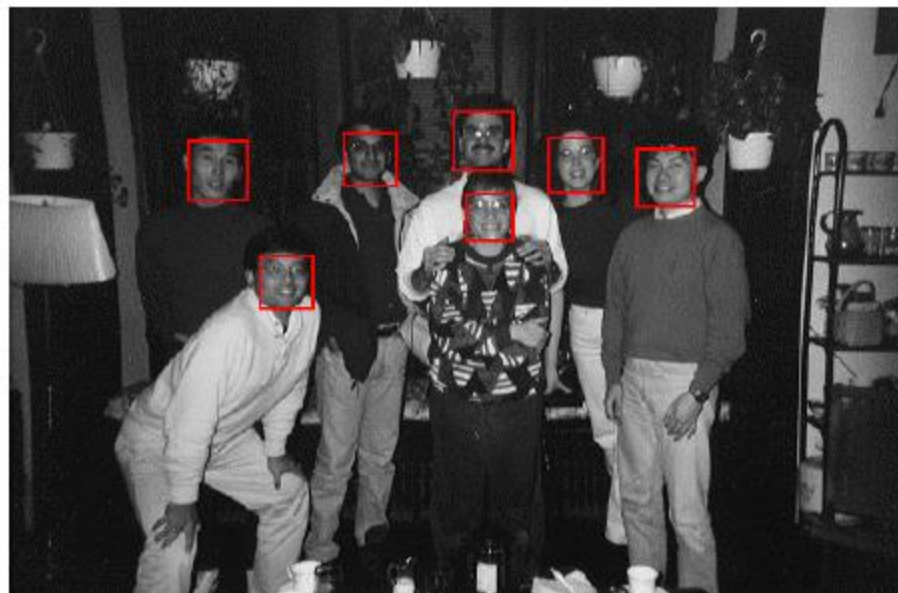
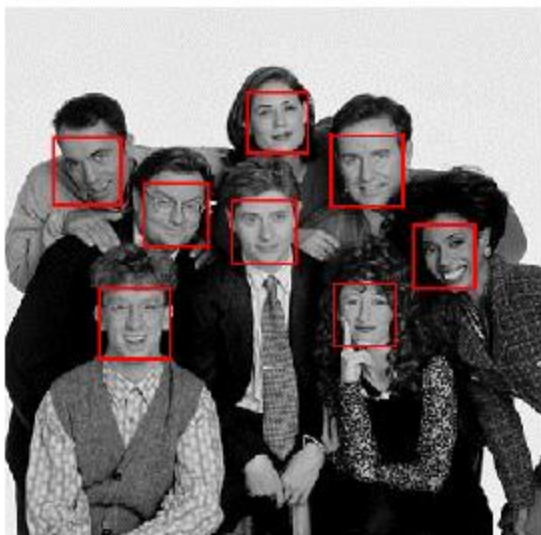
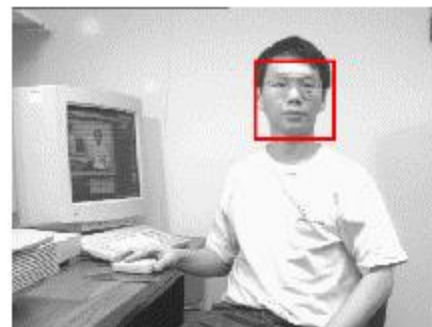
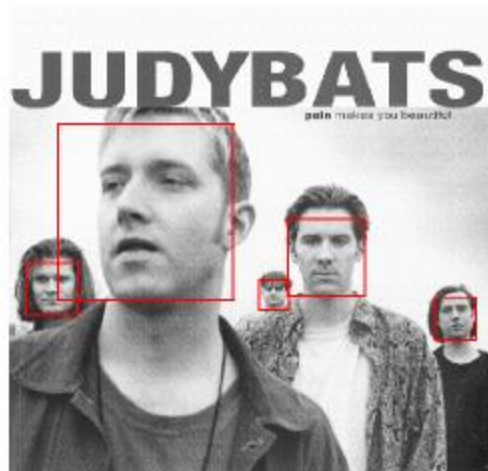
K. Grauman, B. Leibe

Viola-Jones Face Detector: Results

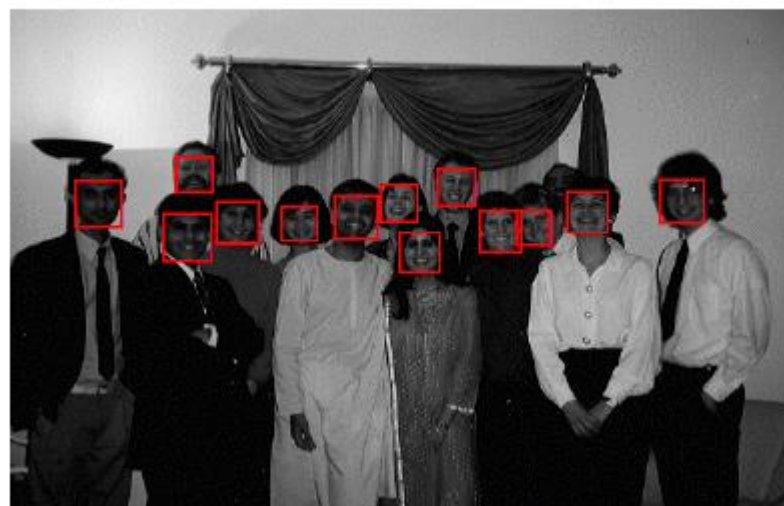
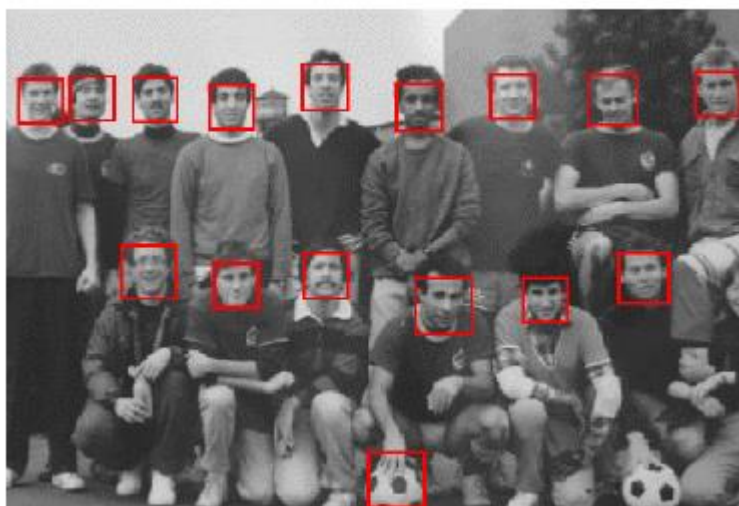
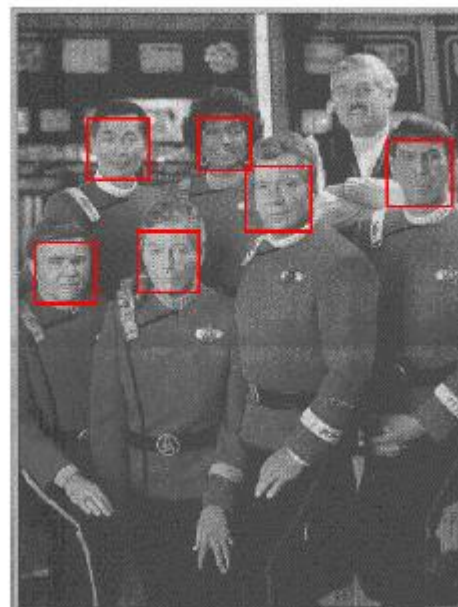
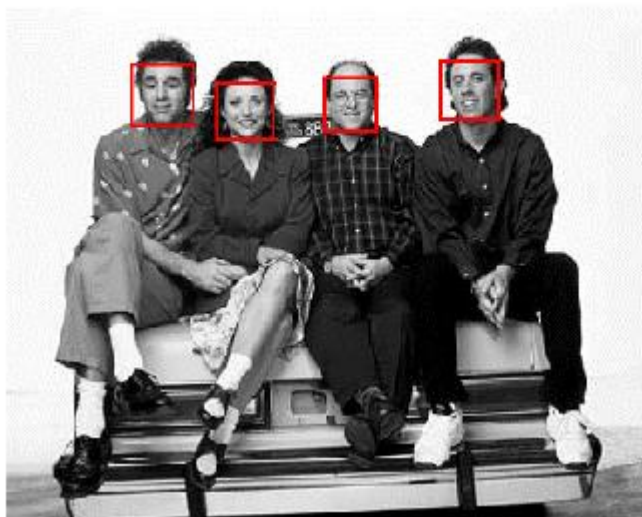


First two features
selected

Viola-Jones Face Detector: Results



Viola-Jones Face Detector: Results

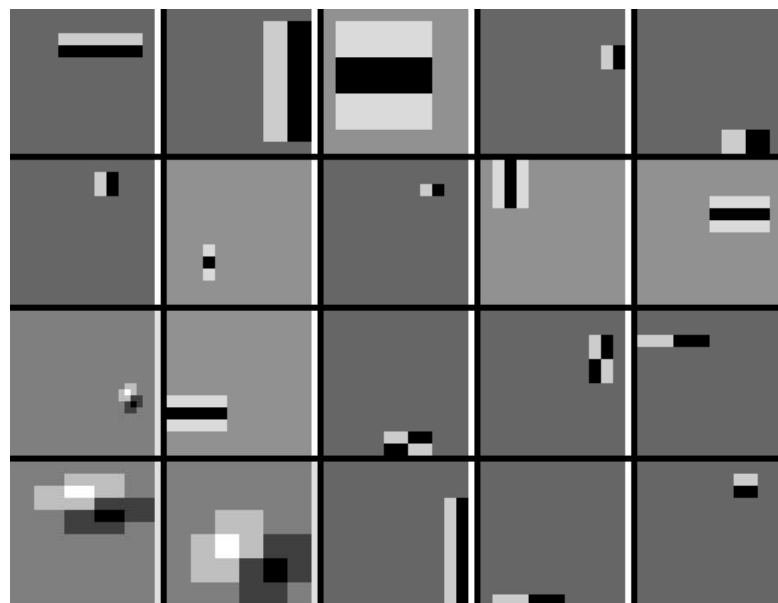


Viola-Jones Face Detector: Results

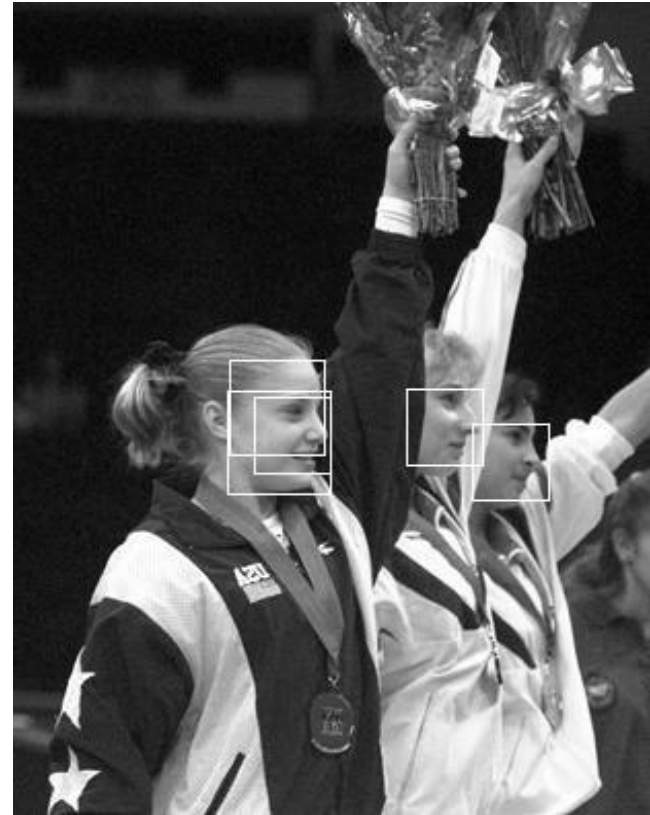


Detecting profile faces?

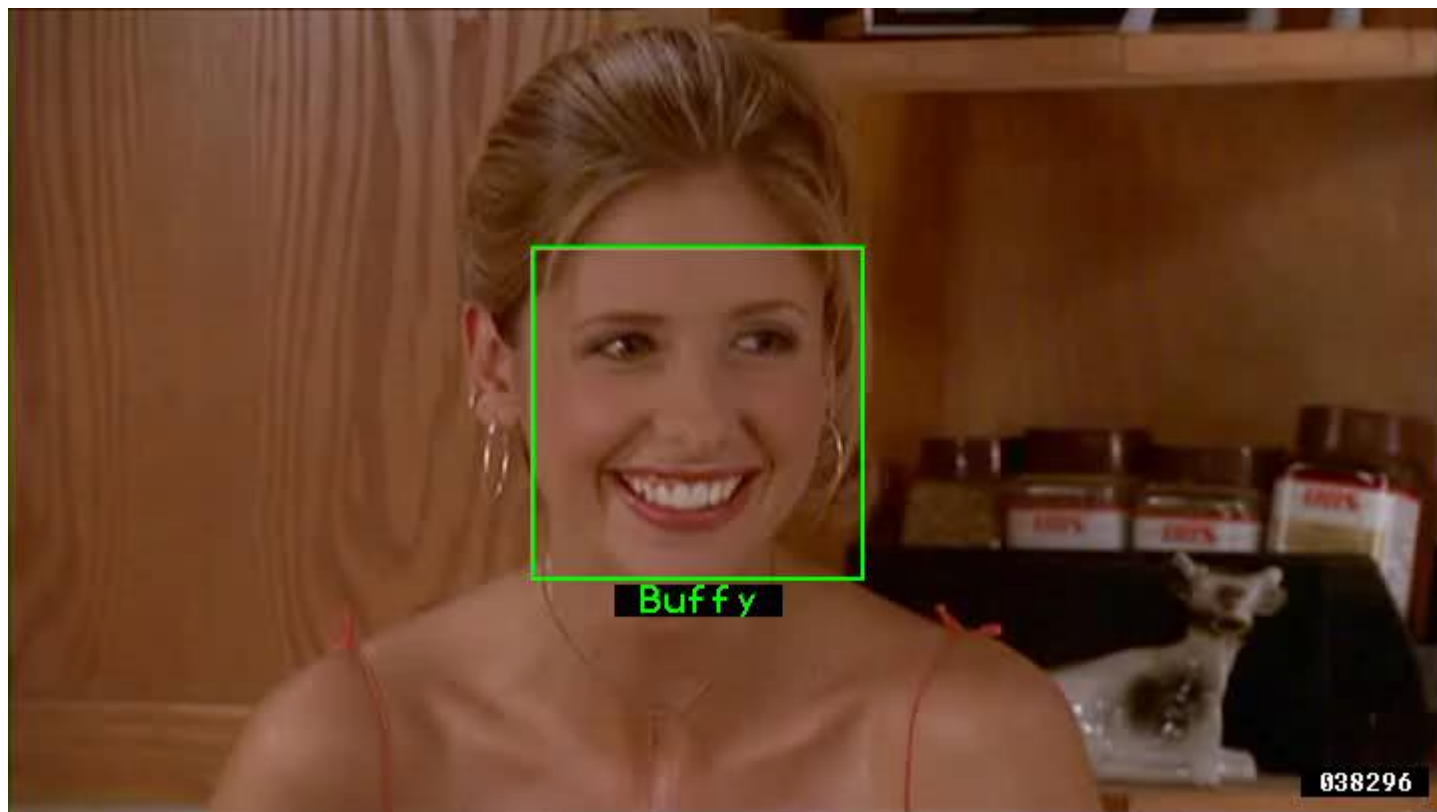
Detecting profile faces requires training separate detector with profile examples.



Viola-Jones Face Detector: Results



Example application



Frontal faces detected and then tracked, character names inferred with alignment of script and subtitles.

Everingham, M., Sivic, J. and Zisserman, A.

"Hello! My name is... Buffy" - Automatic naming of characters in TV video
BMVC 2006.

<http://www.robots.ox.ac.uk/~vgg/research/nface/index.html>

Example application: faces in photos



[All Web](#) [People](#) [Objects](#) [Tags](#) **My Photos**

SEARCH

[Advanced](#)

Riya Personal Search

Use our face recognition and text recognition, to search your personal photos

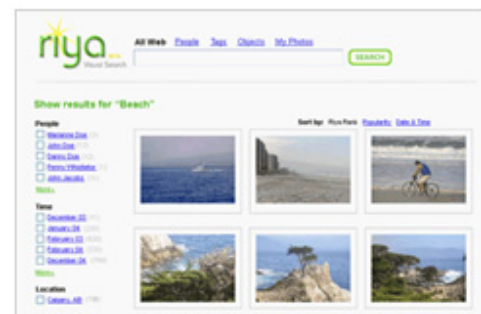
Upload your personal photos
(public or privately)



Use our face and text recognition
to auto tag your photos



Search & share photos
with your friends



Riya's Personal Search lets you upload and search your own photos by name. You can keep them private or make them public and share with all Riya searchers. We allow you to use face and text recognition to search your own photos.

Highlights

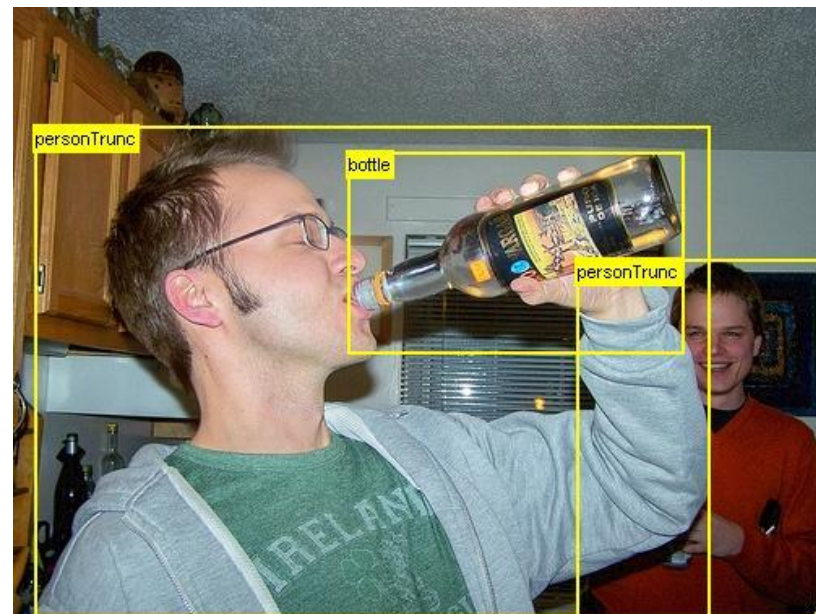
- Sliding window detection and global appearance descriptors:
 - Simple detection protocol to implement
 - Good feature choices critical
 - Past successes for certain classes

Limitations

- High computational complexity
 - For example: 250,000 locations x 30 orientations x 4 scales = 30,000,000 evaluations!
 - If training binary detectors independently, means cost increases linearly with number of classes
- With so many windows, false positive rate better be low

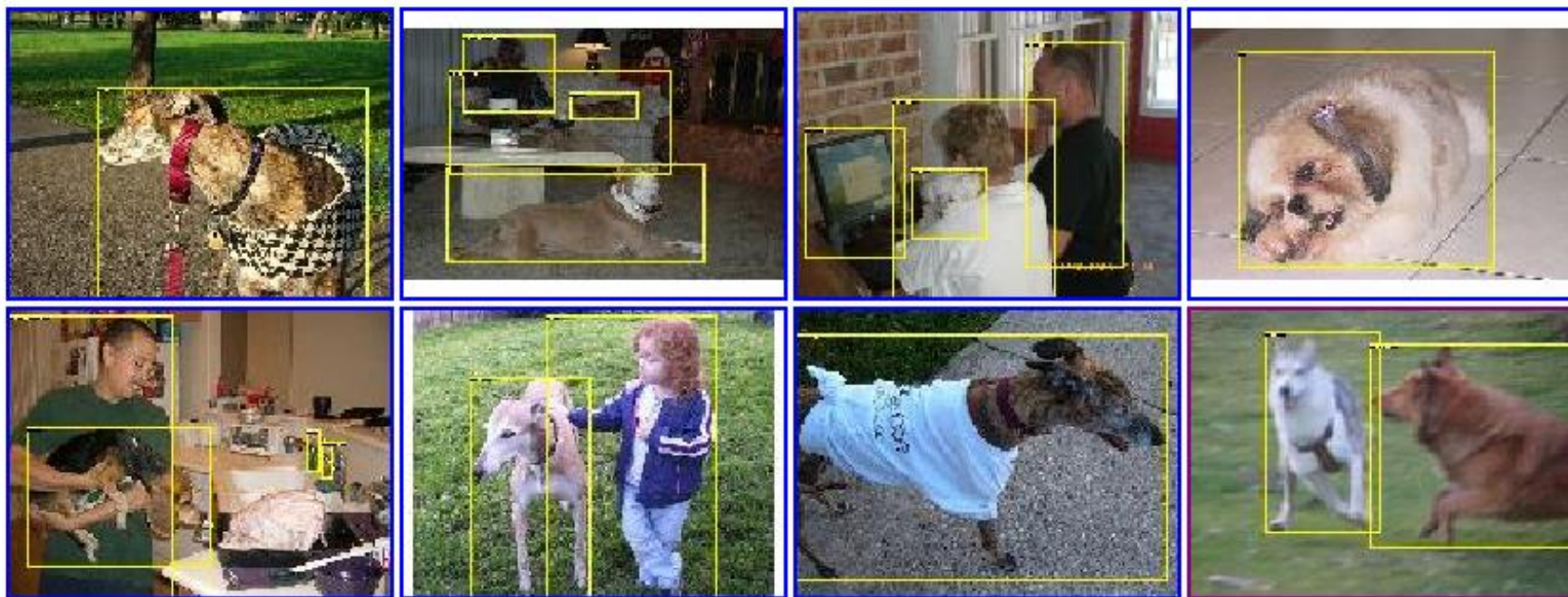
Limitations (continued)

- Not all objects are “box” shaped



Limitations (continued)

- Non-rigid, deformable objects not captured well with representations assuming a fixed 2d structure; or must assume fixed viewpoint
- Objects with less-regular textures not captured well with holistic appearance-based descriptions



Limitations (continued)

- If considering windows in isolation, context is lost

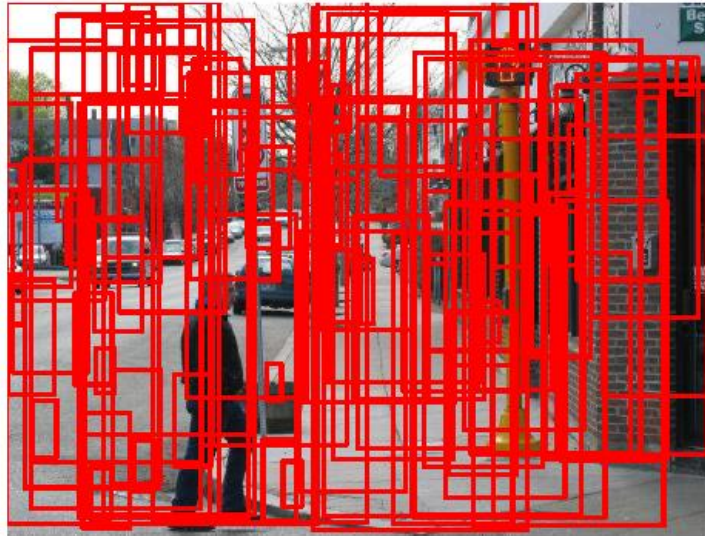


Sliding window

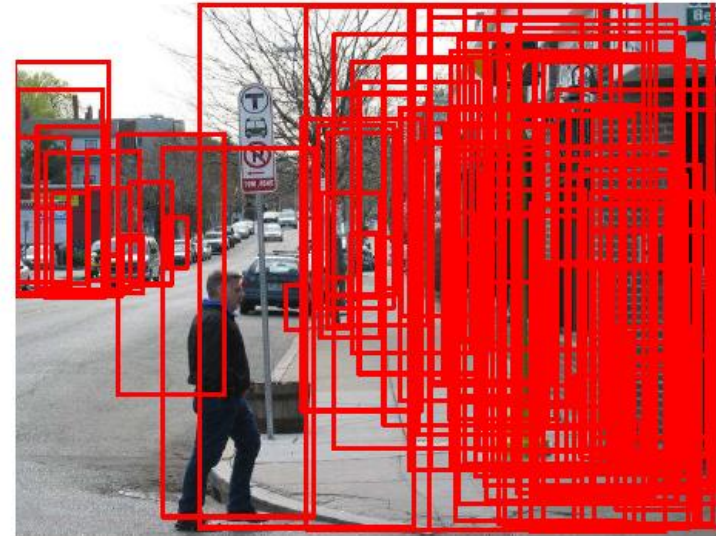


Detector's view

Context can constrain a sliding window search



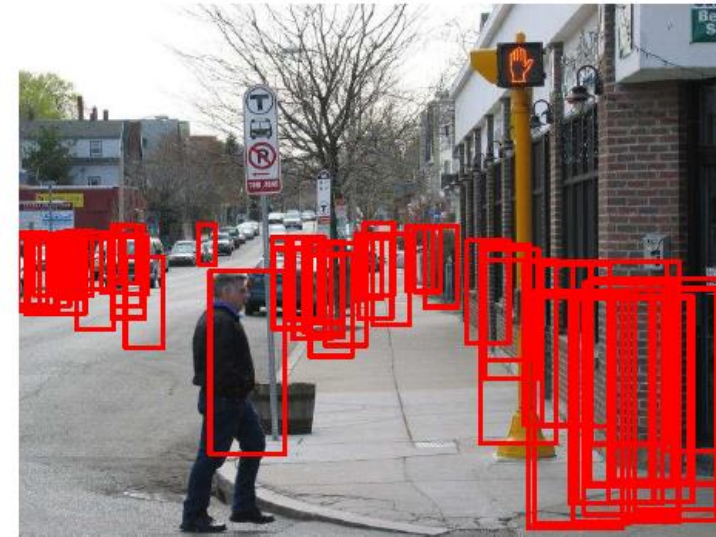
(b) $P(\text{person}) = \text{uniform}$



(d) $P(\text{person} \mid \text{geometry})$



(f) $P(\text{person} \mid \text{viewpoint})$



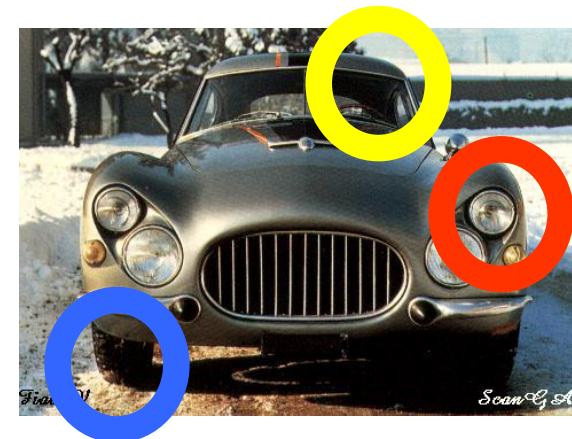
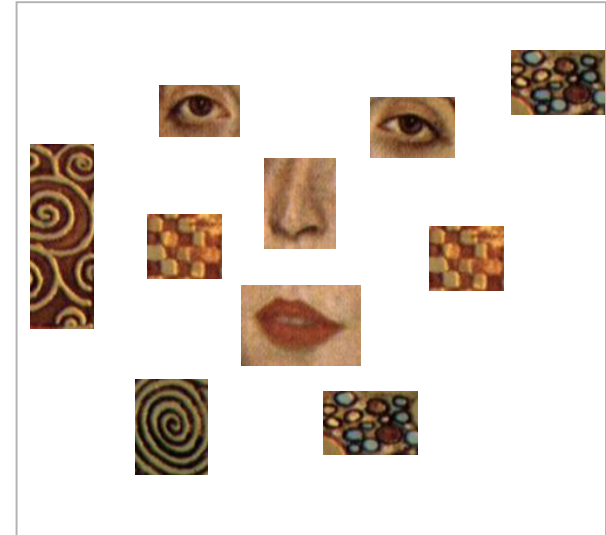
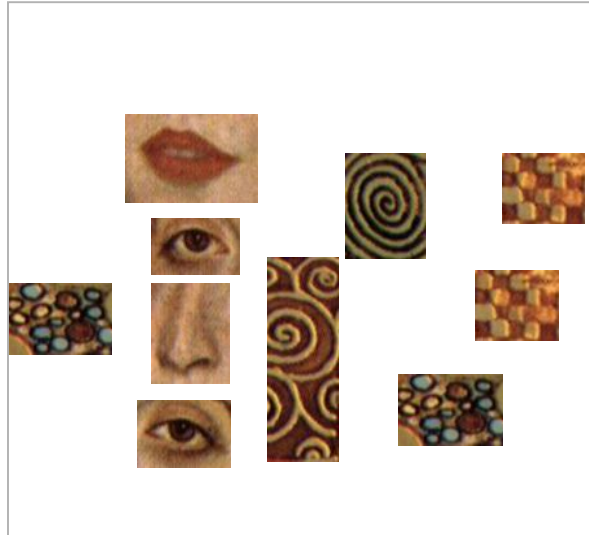
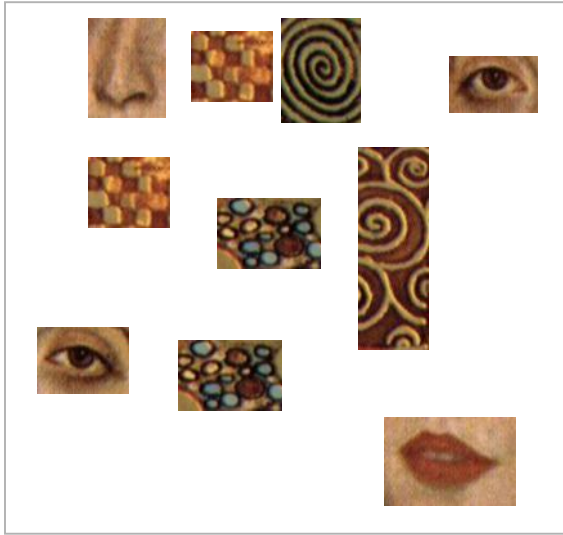
(g) $P(\text{person} \mid \text{viewpoint, geometry})$

Limitations (continued)

- In practice, often entails large, cropped training set (expensive)
- Requiring good match to a global appearance description can lead to sensitivity to partial occlusions



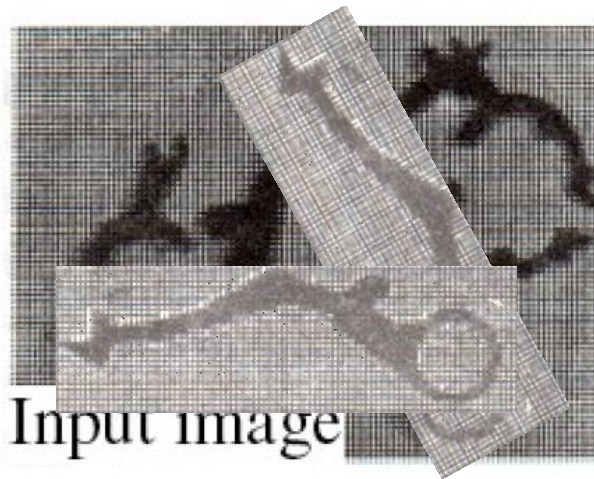
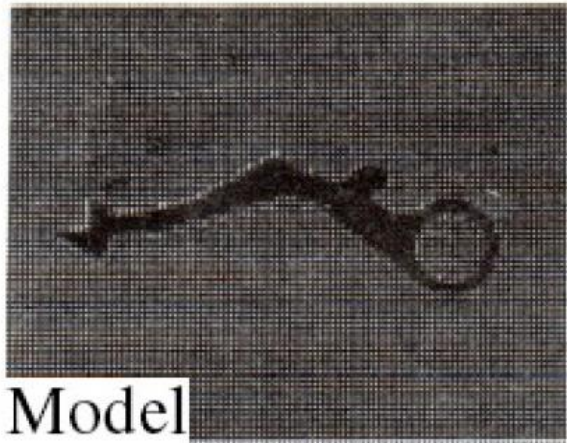
Models based on local features will alleviate some of these limitations...



Local-feature Alignment

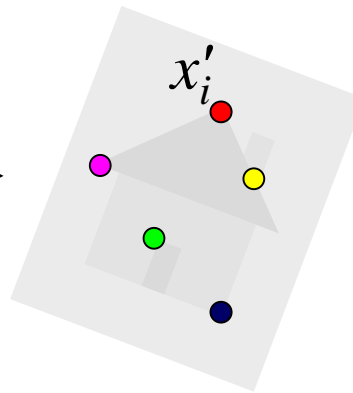
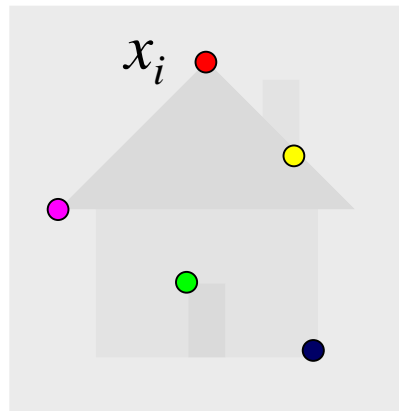
Hypothesize and test: main idea

- Given model of object
- New image: hypothesize object identity and pose
- Render object in camera
- Compare rendering to actual image: if close, good hypothesis.



Recall: Alignment

- Alignment: fitting a model to a transformation between pairs of features (*matches*) in two images

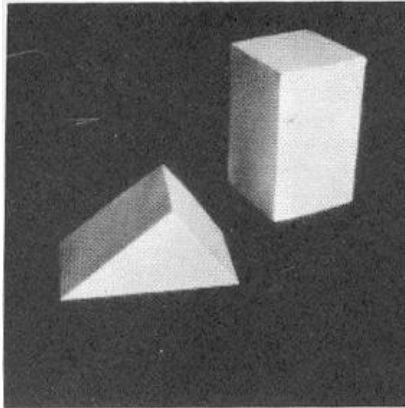


Find transformation T
that minimizes

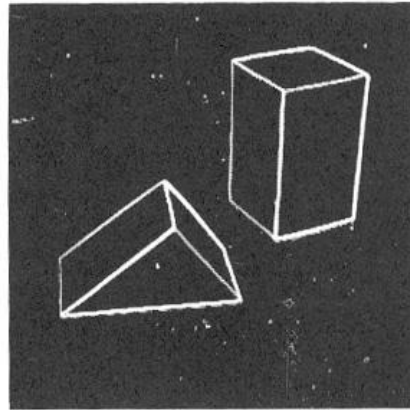
$$\sum_i \text{residual}(T(x_i), x'_i)$$

Alignment-based

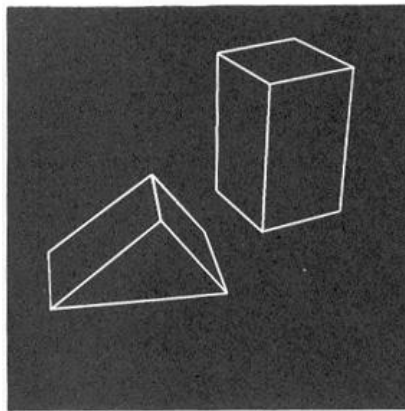
-23-4445(a-d)



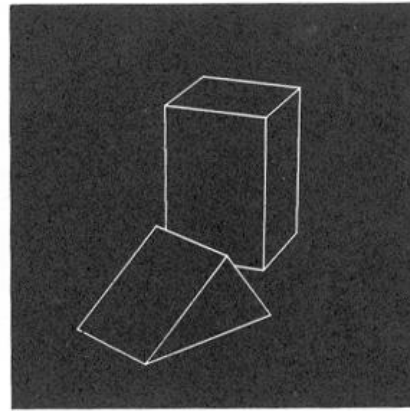
(a) Original picture.



(b) Differentiated picture.



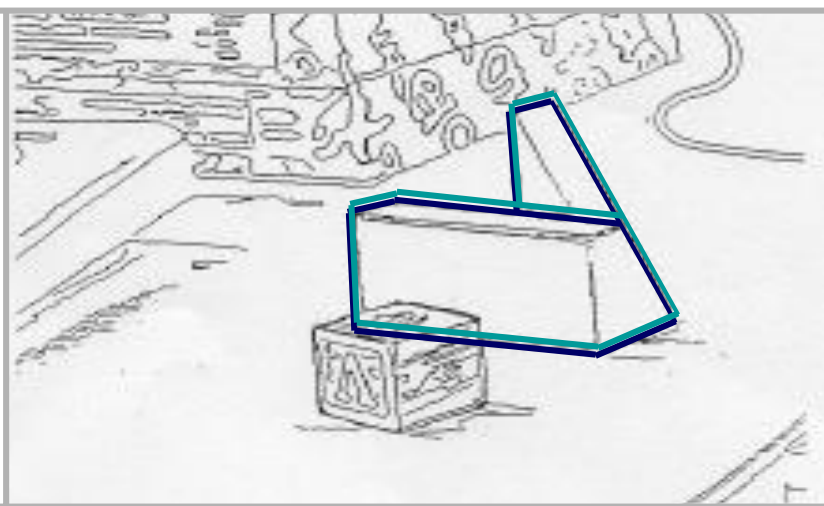
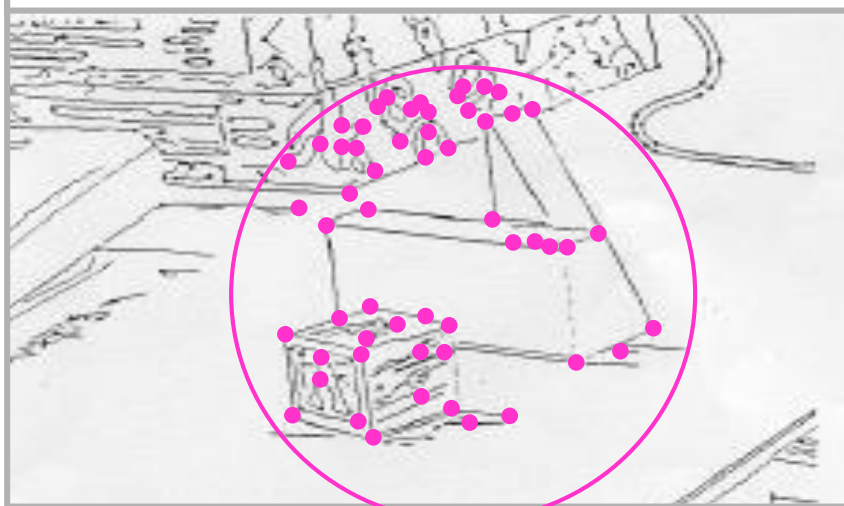
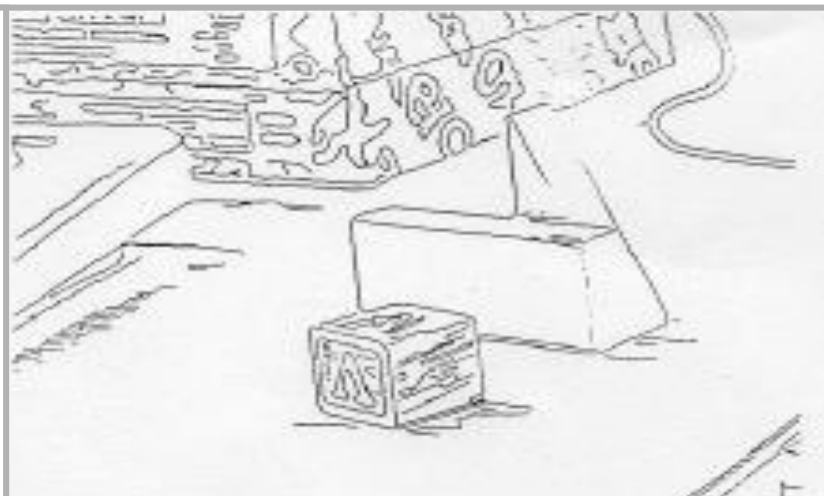
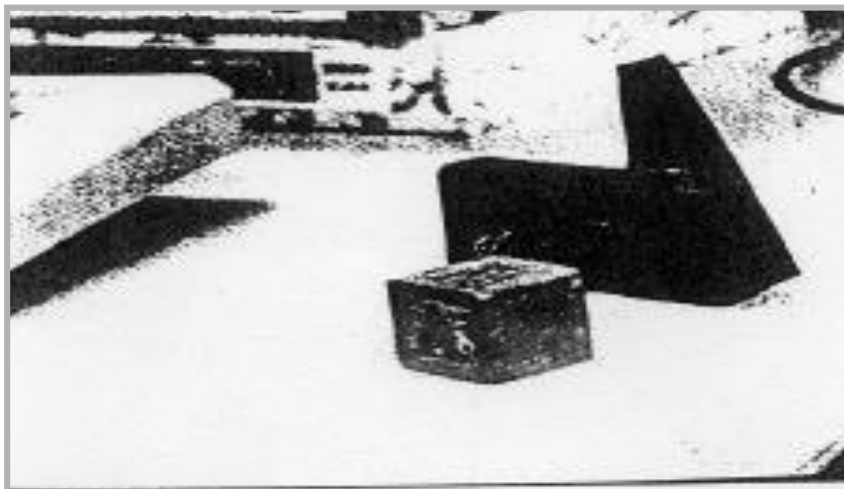
(c) Line drawing.



(d) Rotated view.

L. G. Roberts, *Machine Perception of Three Dimensional Solids*, Ph.D. thesis, MIT Department of Electrical Engineering, 1963.

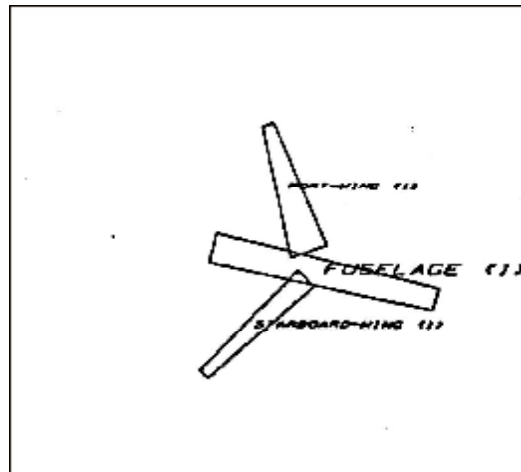
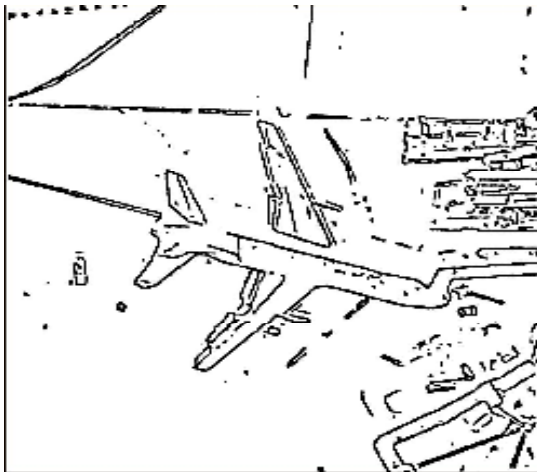
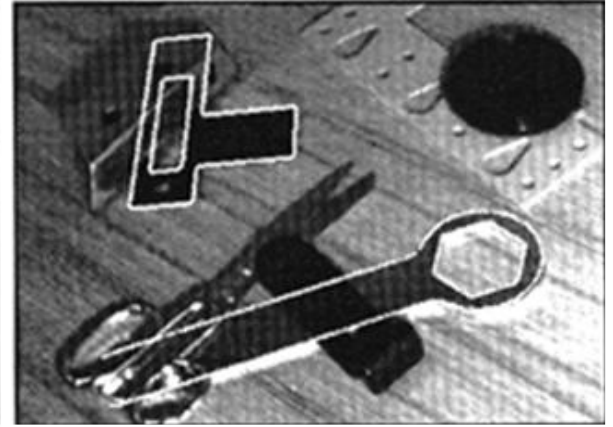
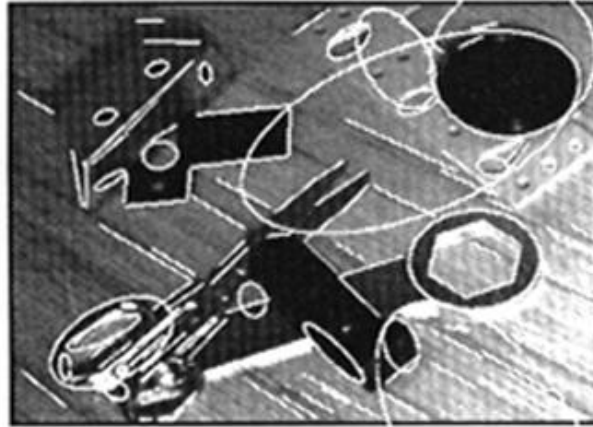
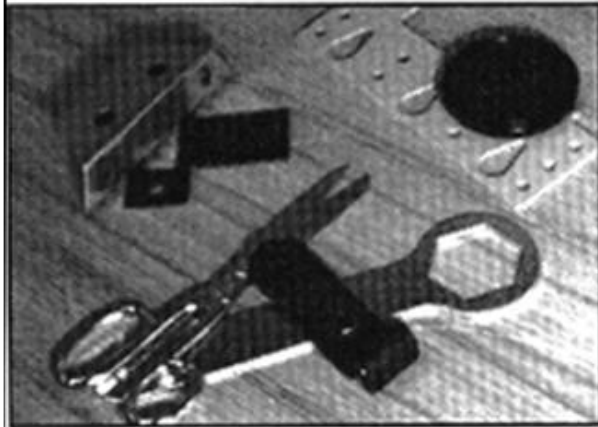
Alignment-based



Huttenlocher & Ullman (1987)

Source: Lana Lazebnik

Alignment-based



ACRONYM (Brooks and Binford, 1981)

How to form a hypothesis?

Given a particular model object, we can estimate the *correspondences* between image and model features

Use correspondence to estimate model pose relative to object coordinate frame

Generating hypotheses

We want a good correspondence between model features and image features.

- Brute force?

Generating hypotheses

We want a good correspondence between model features and image features.

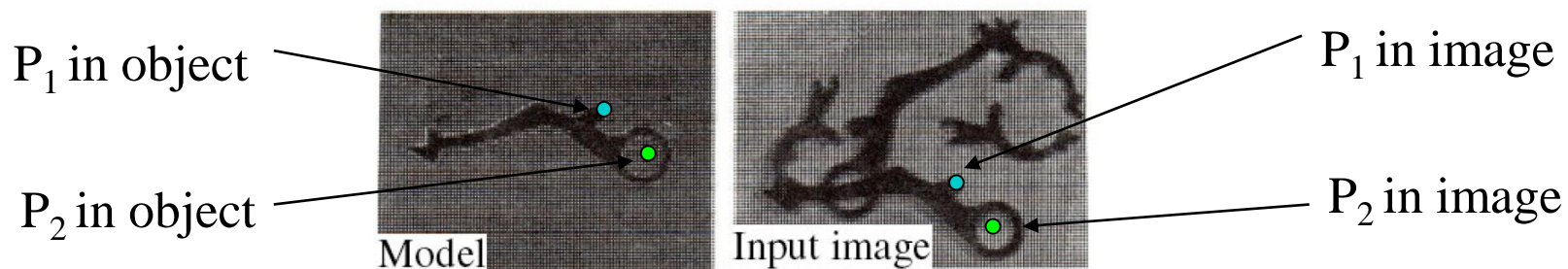
- Brute force?
- **Pose consistency**, alignment: use subsets of features to estimate larger correspondence
- **Voting**, pose clustering

Pose consistency / alignment

- Key idea:
 - If we find good correspondences for a small set of features, it is easy to obtain correspondences for a much larger set.
- Strategy:
 - Generate hypotheses using small numbers of correspondences
 - Backproject: transform *all* model features to image features
 - Verify

Example: 2d affine mappings

- Say camera is looking down perpendicularly on planar surface



- We have two coordinate systems (object and image), and they are related by some affine mapping (rotation, scale, translation, shear).

Alignment: verification

- Given the back-projected model in the image:
 - Check if image edges coincide with predicted model edges
 - May be more robust if also require edges to have the same orientation
 - Consider texture in corresponding regions
- Possible issues?

Alignment: verification

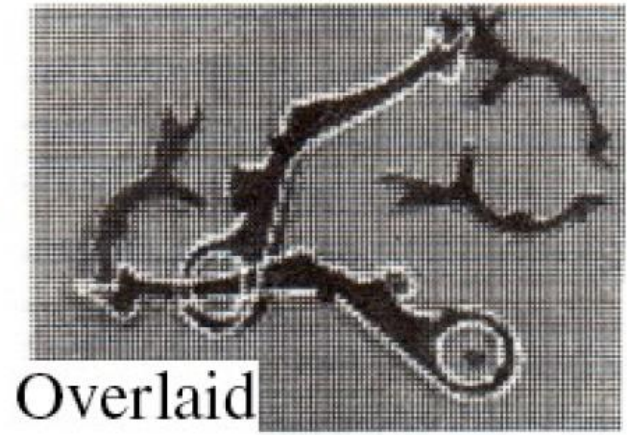
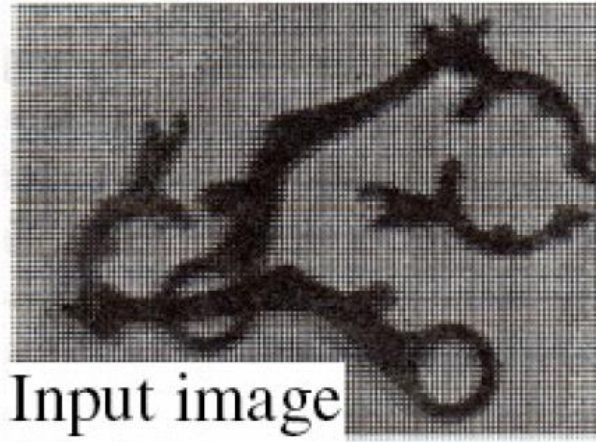
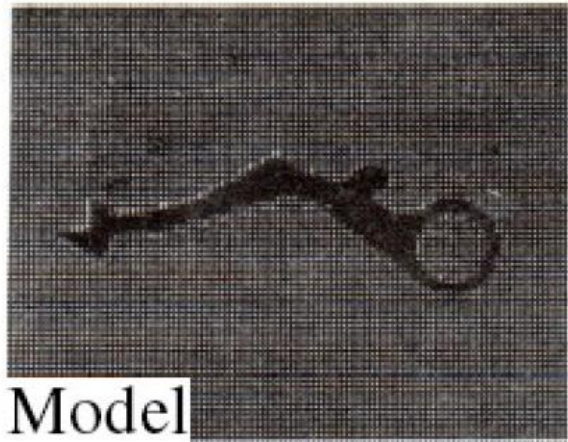
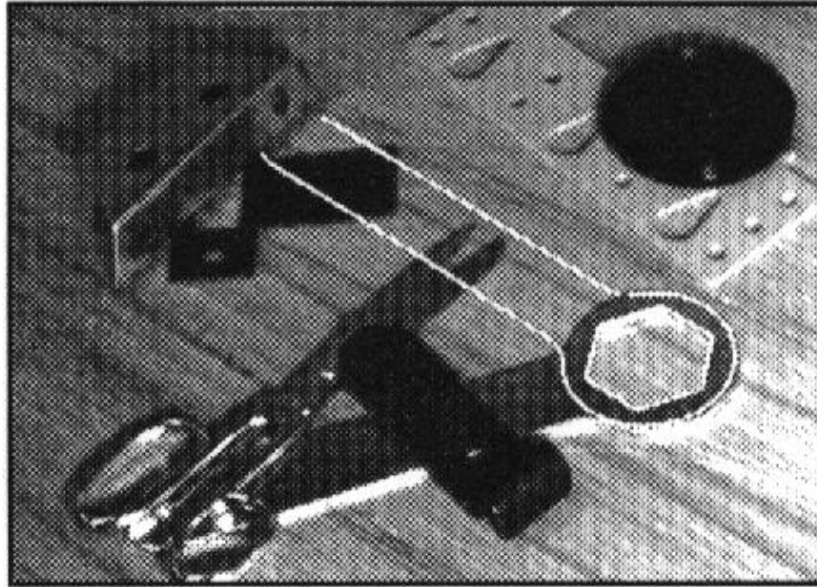


Figure from “Object recognition using alignment,” D.P. Huttenlocher and S. Ullman, Proc. Int. Conf. Computer Vision, 1986, copyright IEEE, 1986

Alignment: verification



Issue with hypothesis & test approach

- May have false matches
 - We want *reliable* features to form the matches
 - ◆ **Local invariant features** useful to find matches, and to verify hypothesis
(*SIFT, etc.*)
- May be too many hypotheses to consider
 - We want to look at the *most likely* hypotheses first
 - ◆ **Pose clustering (voting)**: Narrow down number of hypotheses to verify by letting features *vote* on model parameters.

Pose clustering (voting)

- Narrow down the number of hypotheses to verify: identify those model poses that a lot of features agree on.
 - Use each group's correspondence to estimate pose
 - Vote for that object pose in accumulator array (one array per object if we have multiple models)
- Local invariant features can give more reliable matches and means of verification

Pose clustering and verification with SIFT [Lowe]

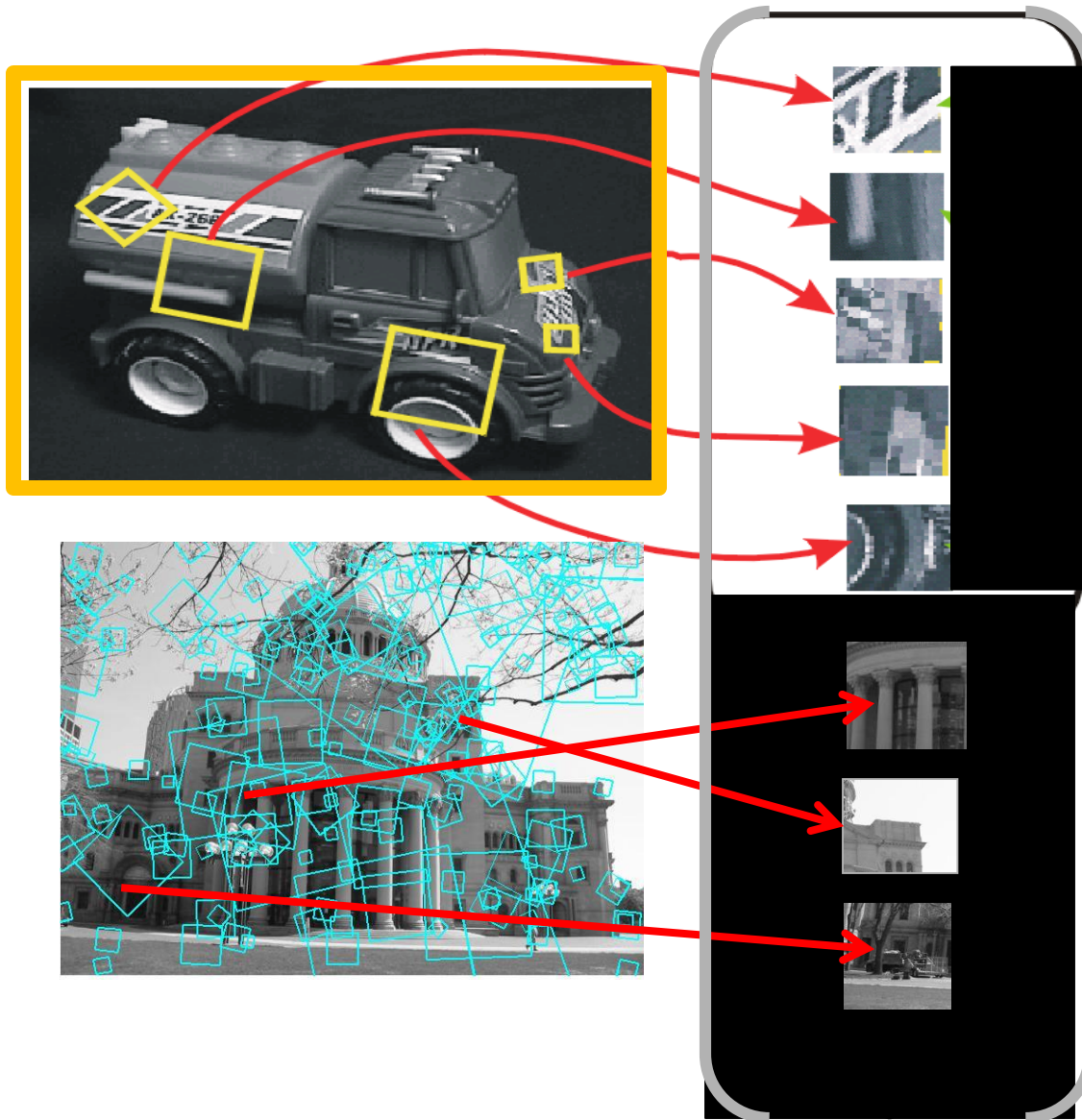
To detect **instances** of objects from a model base:



- 1) Index descriptors (distinctive features narrow possible matches)



Indexing local features



Pose clustering and verification with SIFT [Lowe]

To detect **instances** of objects from a model base:



- 1) Index descriptors (distinctive features narrow possible matches)
- 2) Generalized Hough transform to vote for poses (keypoints have record of parameters relative to model coordinate system)
- 3) Affine fit to check for agreement between model and image features (approximates perspective projection for planar objects)

Planar objects

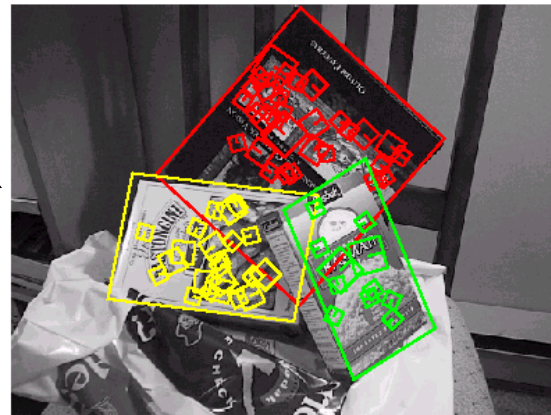


Model images and their SIFT keypoints

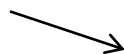


Input image

Model keypoints that were used to recognize, get least squares solution.



Recognition result



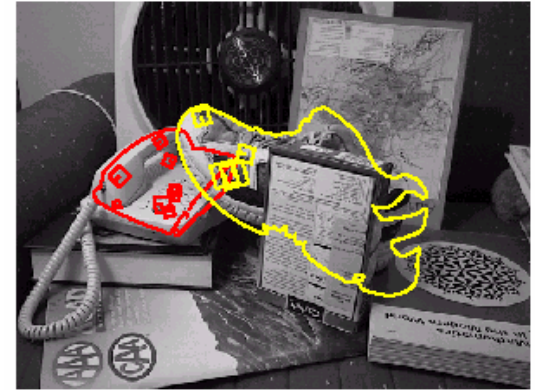
3d objects



Background subtract for
model boundaries



Objects recognized,
though affine model not
as accurate.



Recognition in spite
of occlusion

Recall: difficulties of voting

- Noise/clutter can lead to as many votes as true target
- Bin size for the accumulator array must be chosen carefully
- (Recall Hough Transform)
- In practice, good idea to make broad bins and spread votes to nearby bins, since verification stage can prune bad vote peaks.

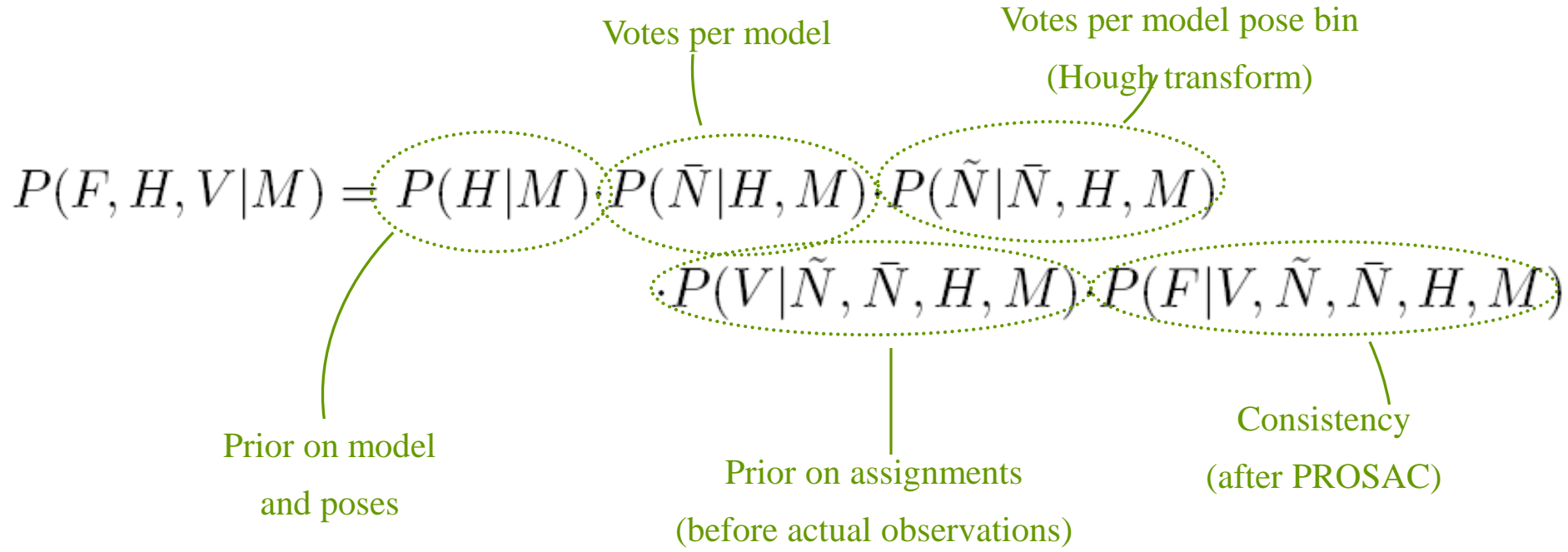
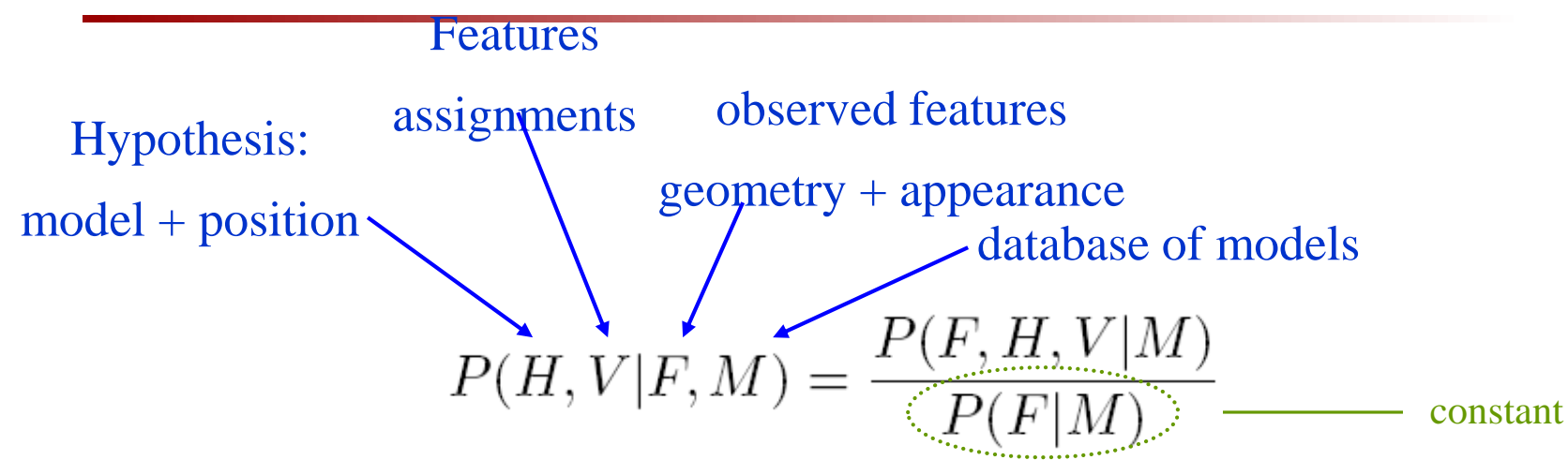
A probabilistic interpretation (and re-tuning) of Lowe's system:

P. Moreels and P. Perona, "A probabilistic cascade of detectors for individual object recognition," European Conference on Computer Vision, 2008.

Coarse-to-Fine detection

- Progressively narrow down focus on correct region of hypothesis space
- Reject with little computation cost irrelevant regions of search space
- Use first information that is easy to obtain
- Simple building blocks organized in a cascade
- Probabilistic interpretation of each step

Score of an extended hypothesis

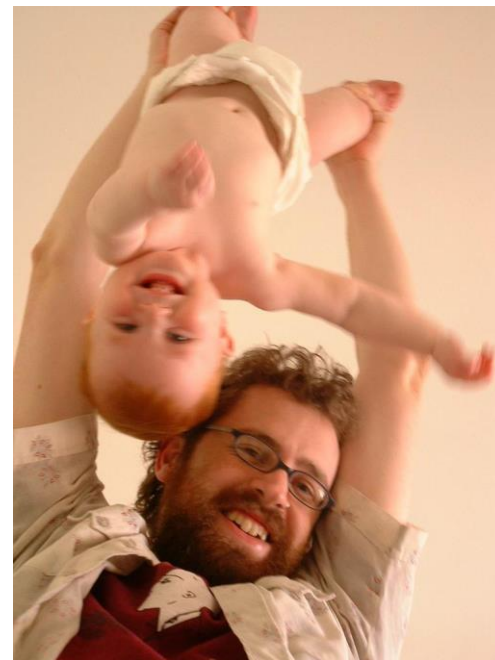


Coarse data : prior knowledge

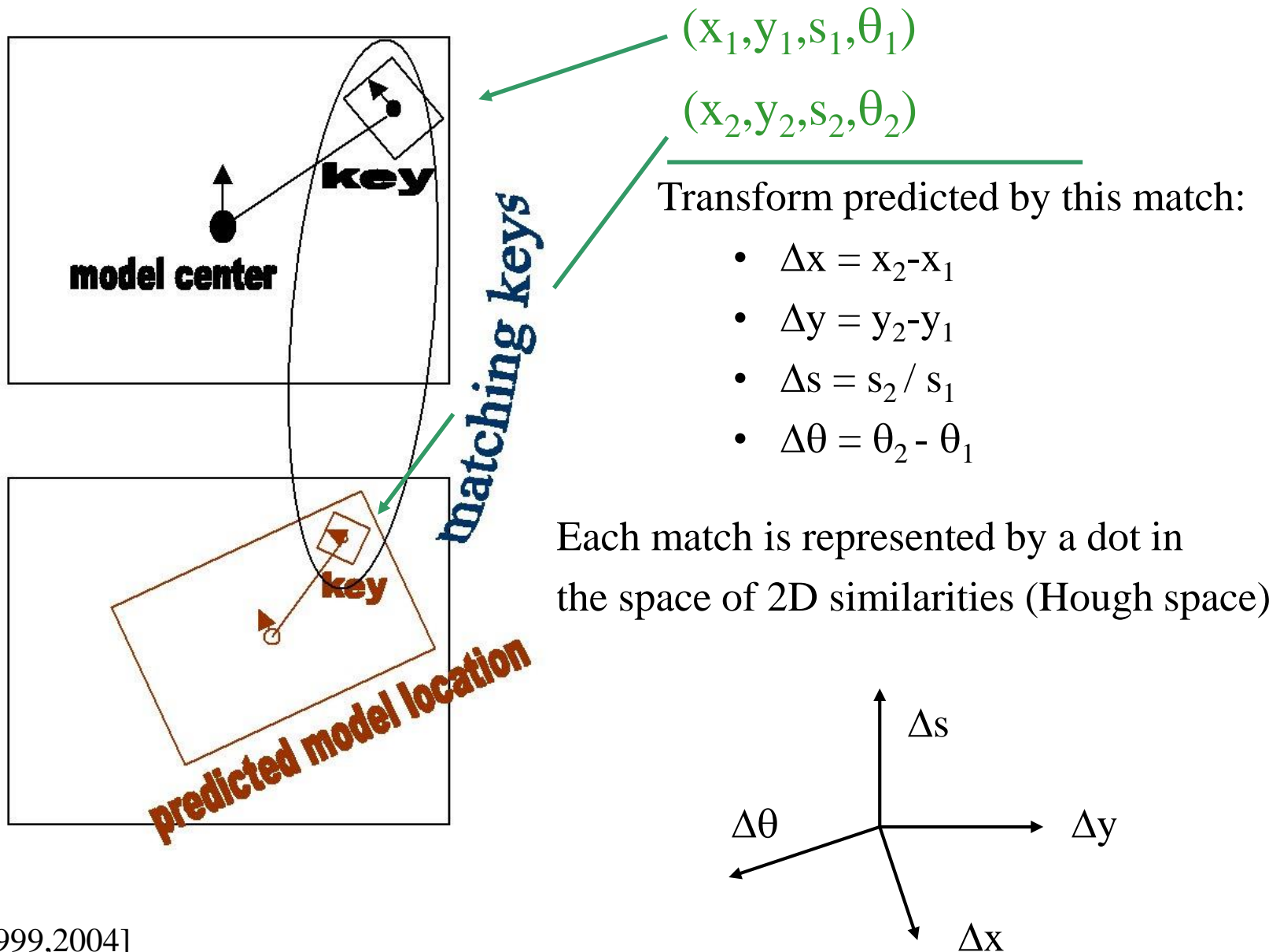
- Which objects are likely to be there, which pose are they likely to have ?



unlikely
situations

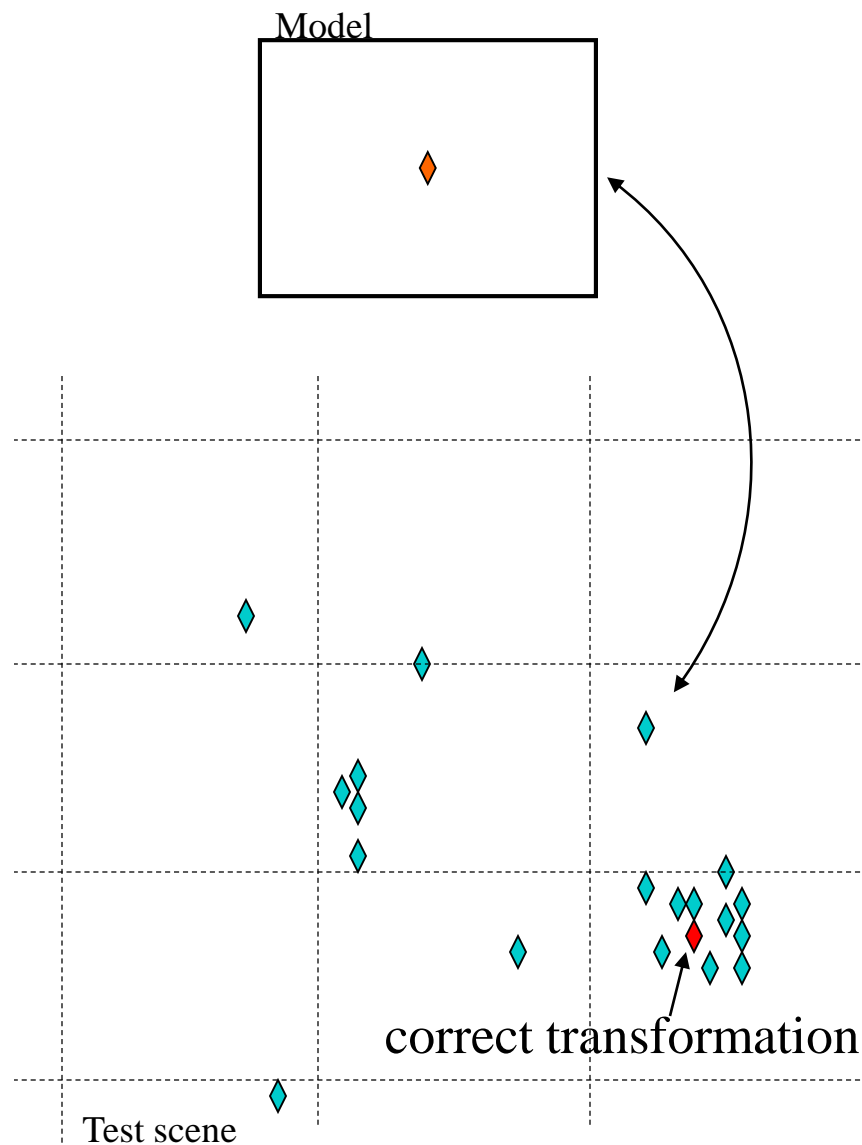


Coarse Hough transform

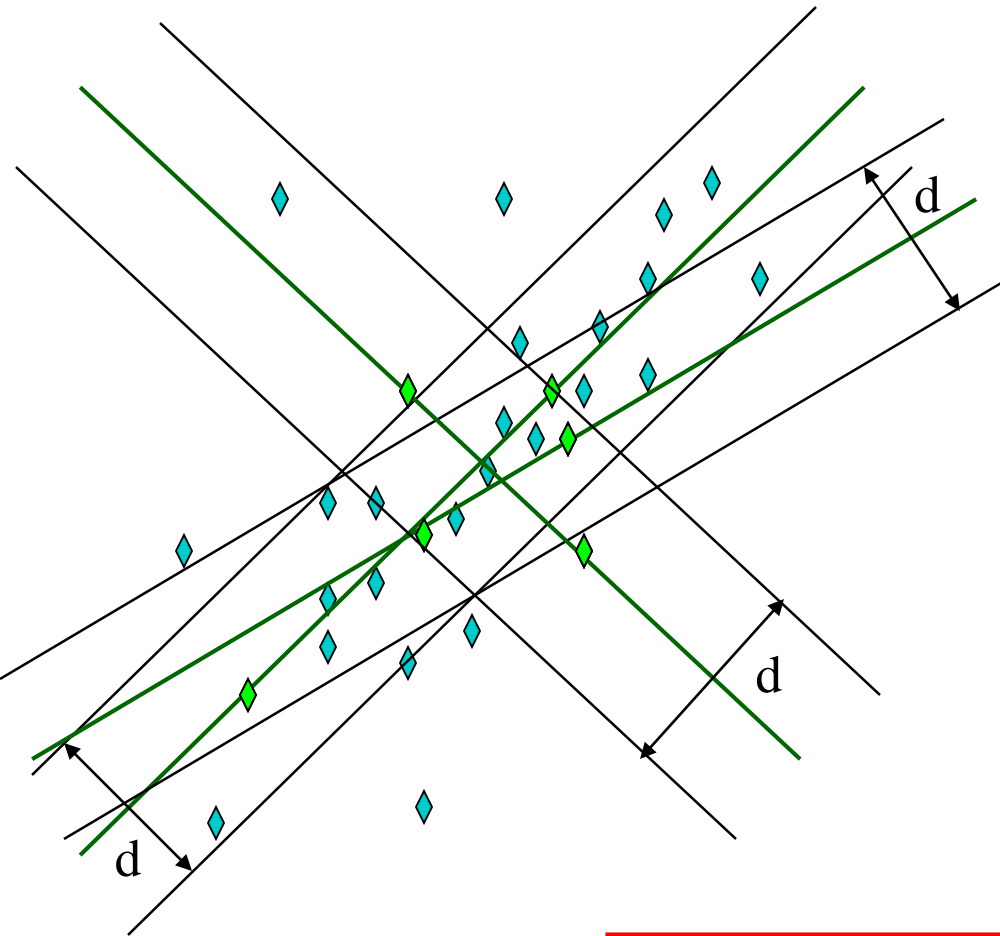


Coarse Hough transform

- Prediction of position of model center after transform
- The space of transform parameters is discretized into ‘bins’
- Coarse bins to limit boundary issues and have a low false-alarm rate for this stage
- We count the number \tilde{N} of votes collected by each bin.



Correspondence or clutter ? PROSAC



- Similar to RANSAC – robust statistic for parameter estimation
- Priority to candidates with good **quality** of appearance match
- 2D affine transform : 6 parameters
⇒ each sample contains 3 candidate correspondences.

[Fischler 1973]

[Chum&Matas 2005]

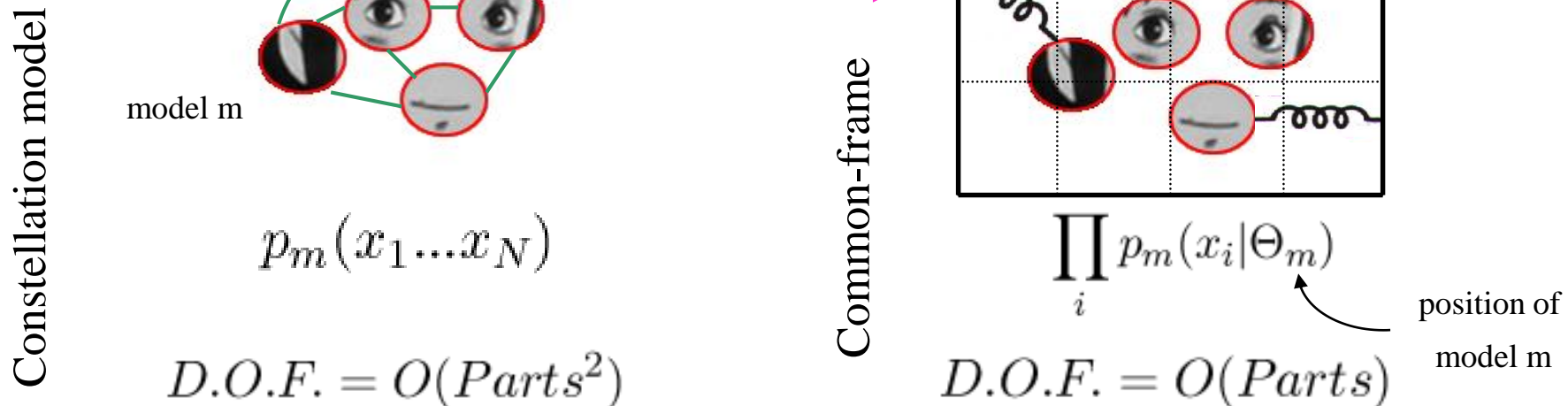
Output of PROSAC : pose transformation
+ set of features correspondences

Consistency

Consistency between observations and predictions from hypothesis

$$P(F|V, \tilde{N}, \bar{N}, H, M) = \prod_{V(i) \neq 0} p_{fg}(f_i|H, f_{V(i)}) \cdot \prod_{V(i)=0} p_{bg}(f_i)$$

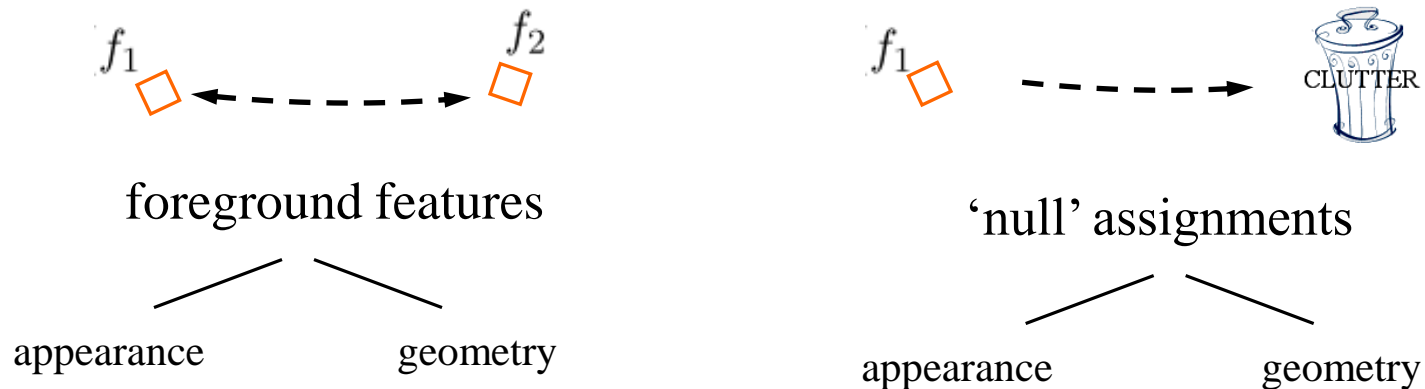
Common-frame approximation : parts are conditionally independent once reference position of the object is fixed. [Lowe1999, Huttenlocher90, Moreels04]



Consistency

Consistency between observations and predictions from hypothesis

$$P(F|V, \tilde{N}, \bar{N}, H, M) = \prod_{V(i) \neq 0} p_{fg}(f_i|H, f_{V(i)}) \cdot \prod_{V(i)=0} p_{bg}(f_i)$$

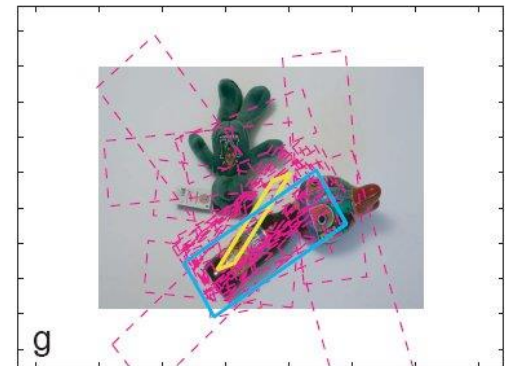
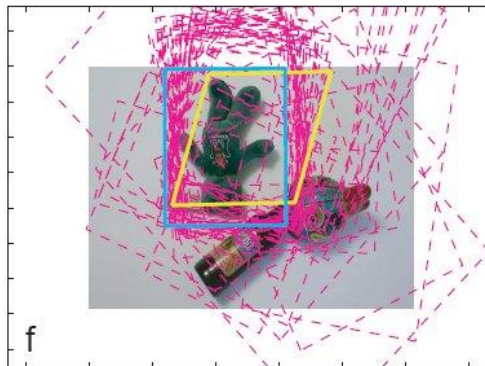
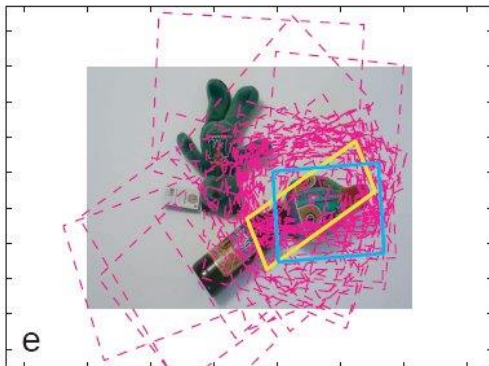
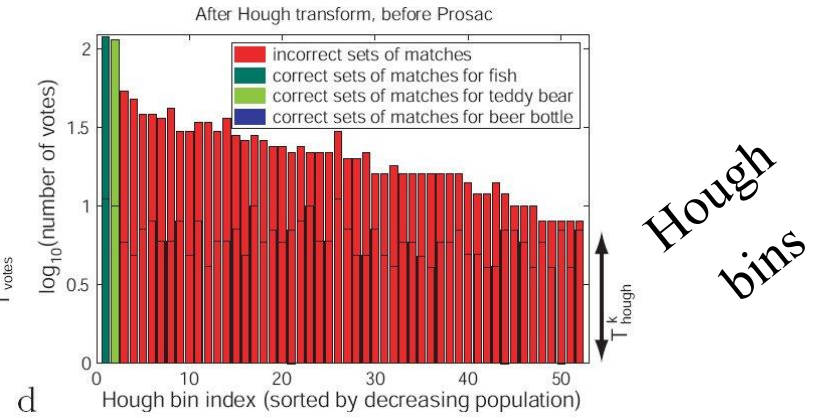
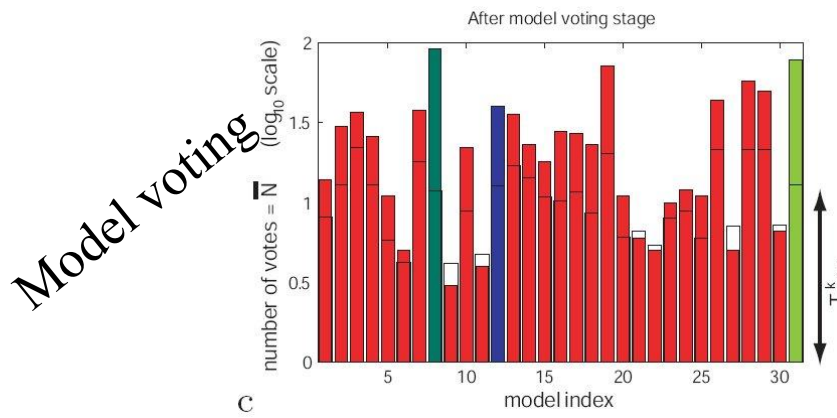
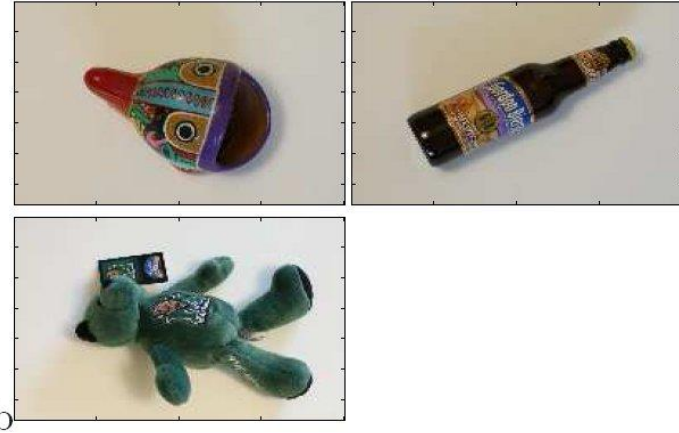


$$p_{fg}(f_i|H, f_{V(i)}) = p_{fg, \mathcal{A}}(\mathcal{A}|H, \mathcal{A}_{V(i)}) \cdot p_{fg, \mathcal{X}}(\mathcal{X}|H, \mathcal{X}_{V(i)})$$

Consistency - appearance

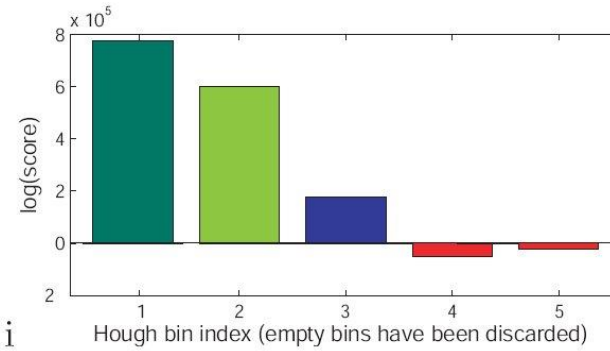
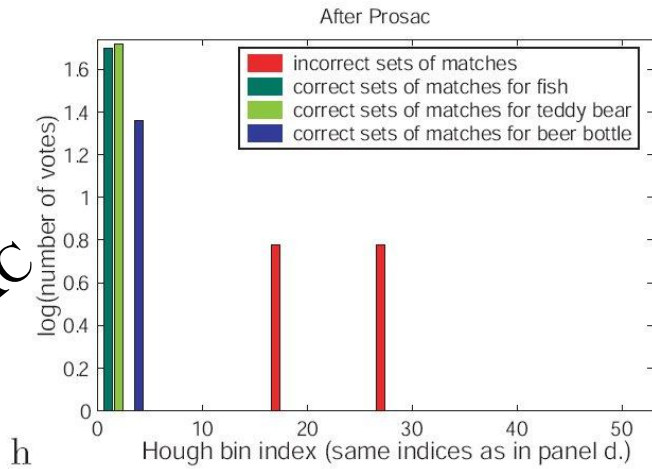
Consistency - geometry

An example

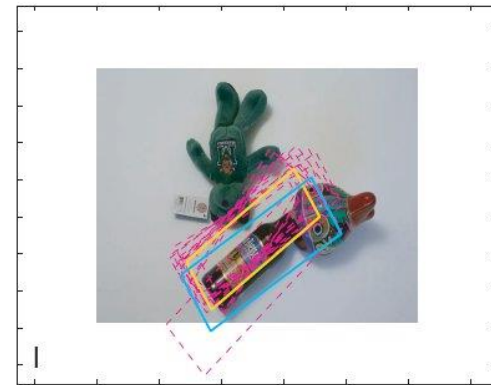
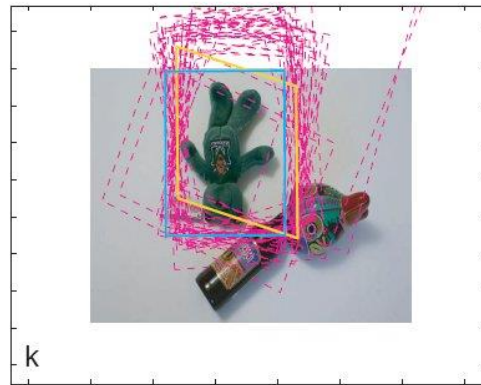
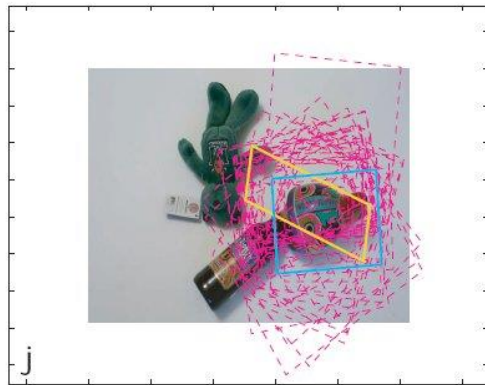


An example

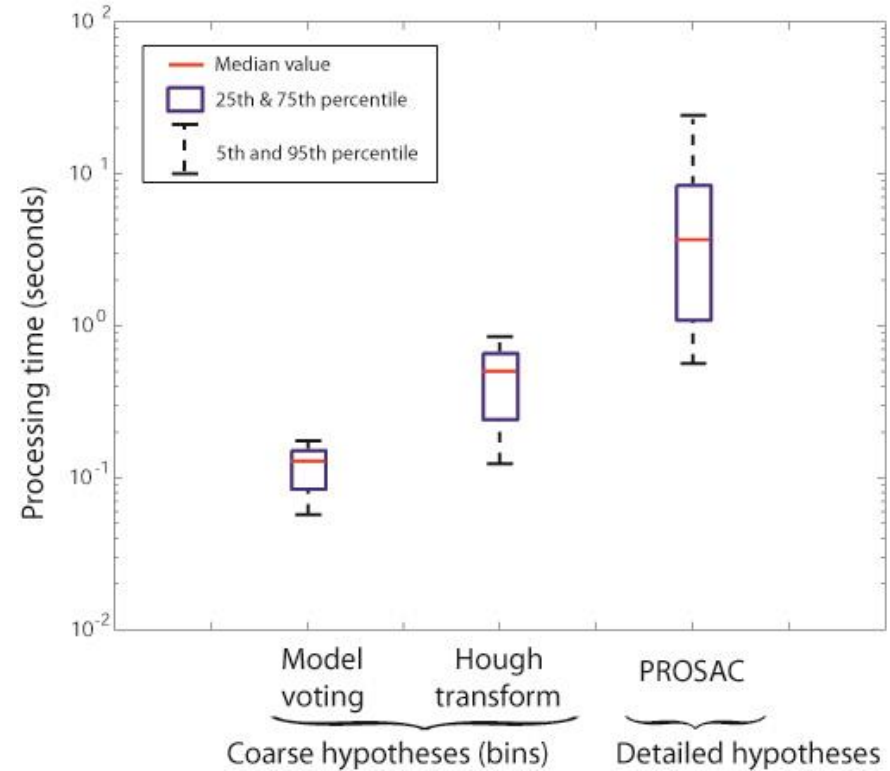
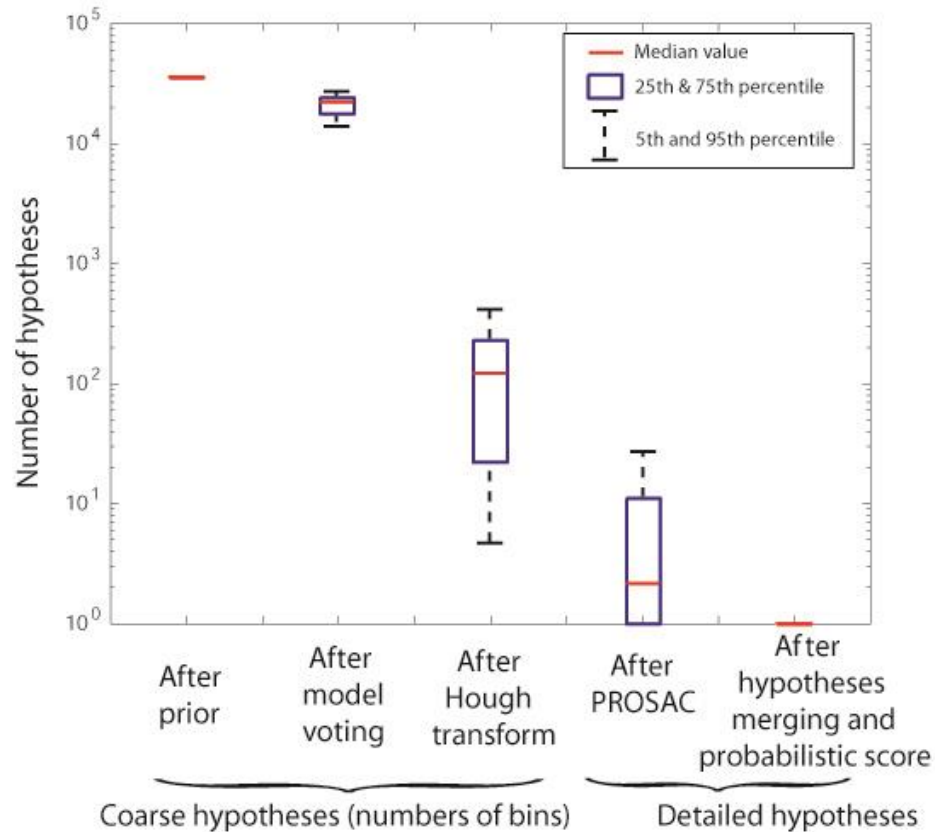
After
PROSAC



Probabilistic
scores



Efficiency of coarse-to-fine processing



Giuseppe Toys database – Models



61 objects, 1-2 views/object

Giuseppe Toys database – Test scenes



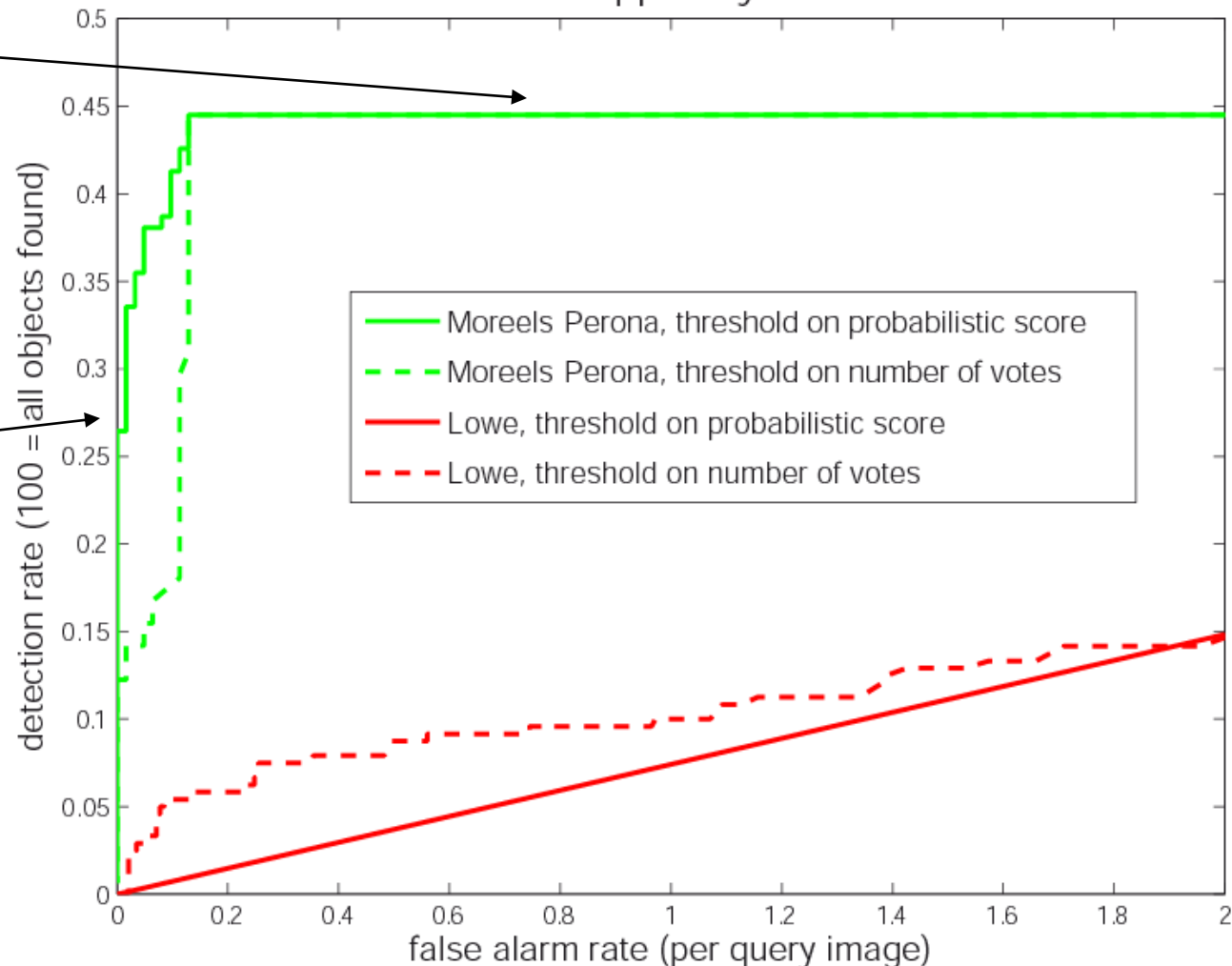
Results – Giuseppe Toys database

–
undetected objects:
features with poor
appearance distinctiveness
index to incorrect models

+
Lower false alarm
rate

- more systematic
verification of
geometry consistency
- more consistent
verification of
geometric consistency

ROC curve Giuseppe Toys database

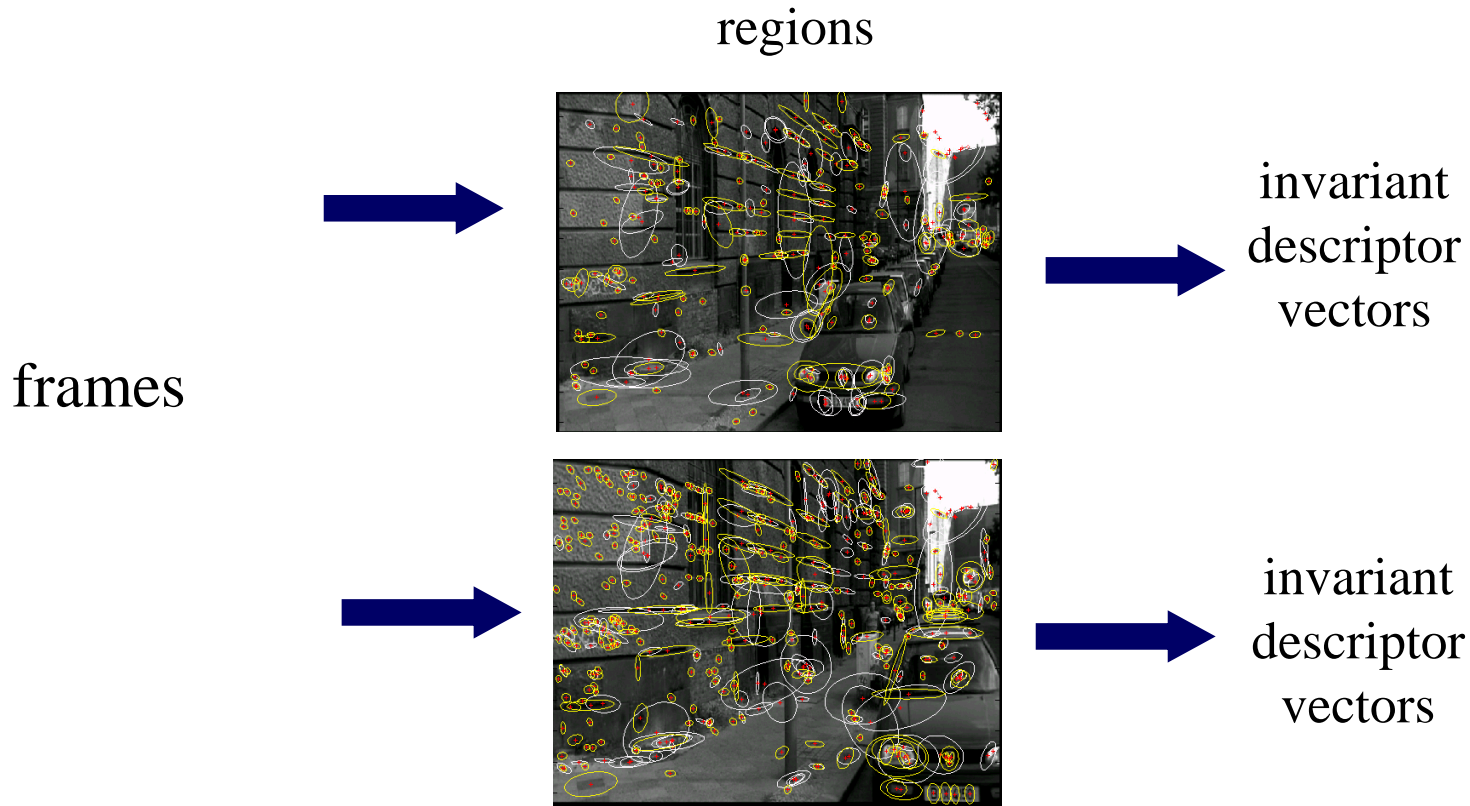


Conclusions – Moreels and Perona

- Coarse-to-fine strategy prunes irrelevant search branches at early stages.
- Probabilistic interpretation of each step.
- Higher performance than Lowe, especially in cluttered environment.
- Front end (features) needs more work for smooth or shiny surfaces.

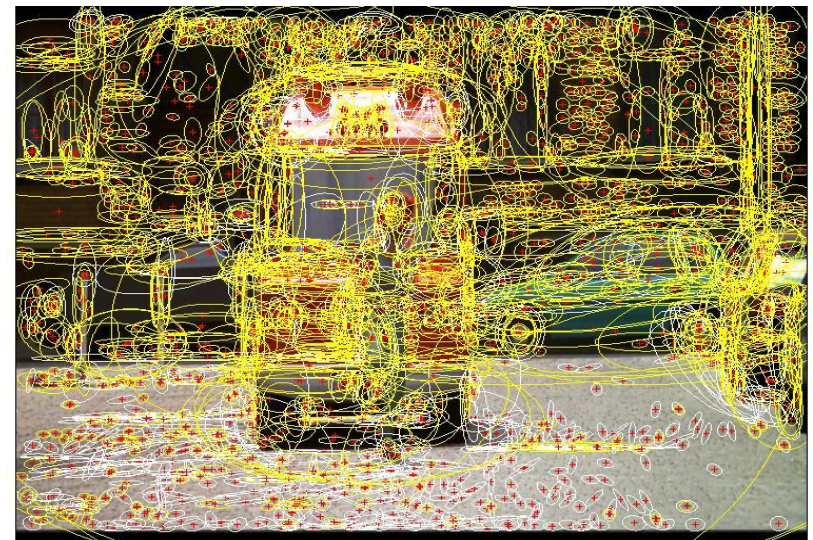
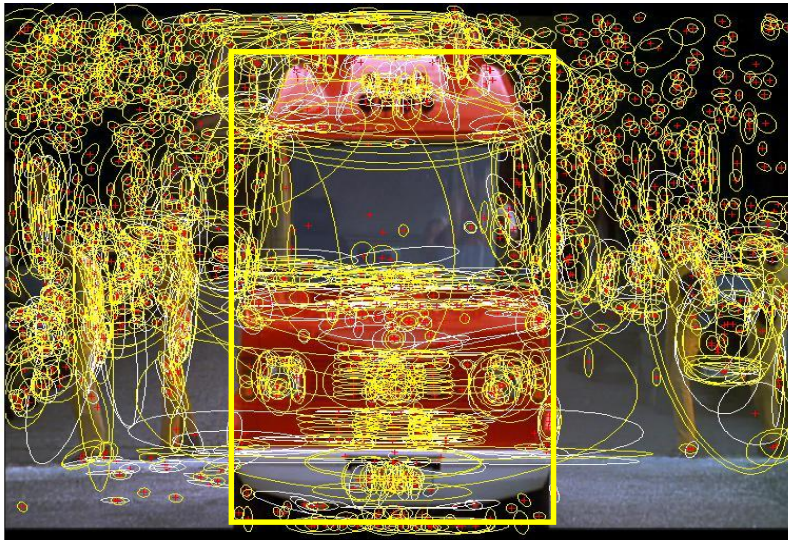
Scaling up: BOW Indexing

Outline of a large-scale retrieval strategy



1. Compute affine covariant regions in each frame independently
2. “Label” each region by a vector of descriptors based on its intensity
3. Finding corresponding regions is transformed to **finding nearest neighbour vectors**
4. Rank retrieved frames by number of corresponding regions
5. Verify retrieved frame based on **spatial consistency**

Example of object recognition



1000+ descriptors per frame

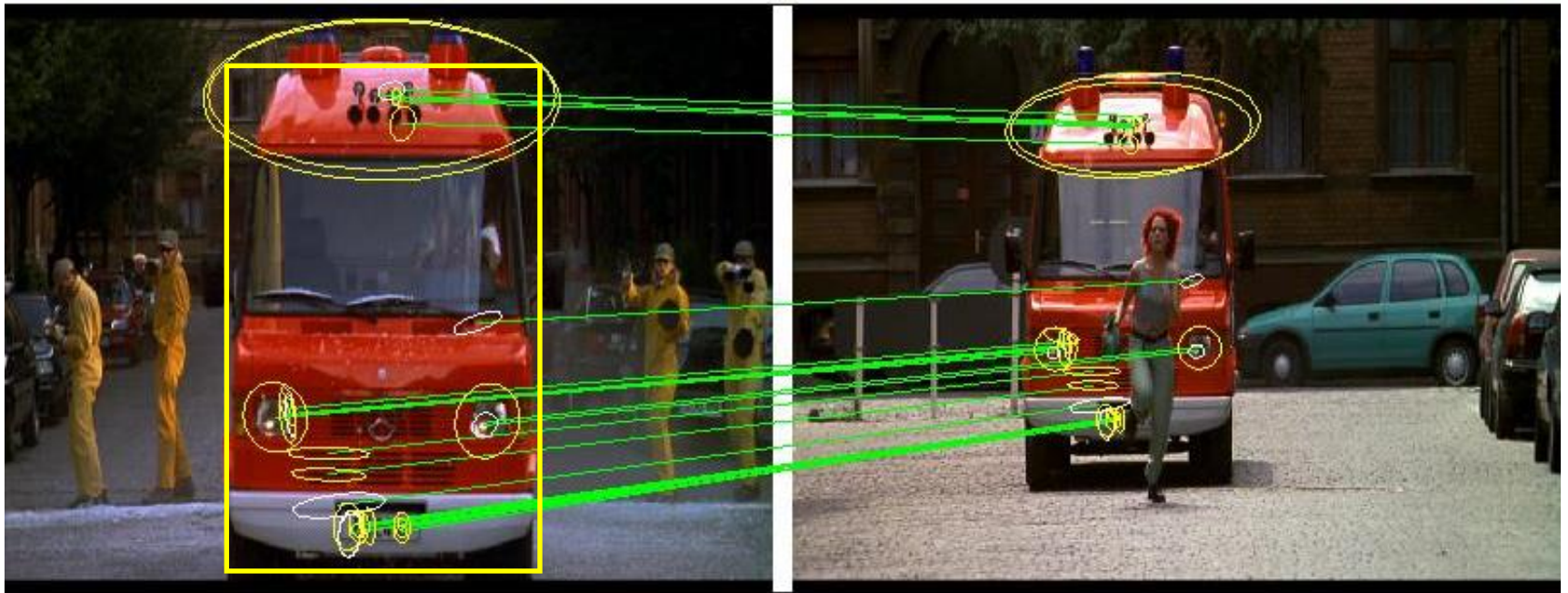


Shape adapted regions





Maximally stable regions

Match regions between frames using SIFT descriptors and spatial consistency



Multiple regions overcome problem of partial occlusion

-  Shape adapted regions
-  Maximally stable regions

Visual search using local regions

Schmid and Mohr '97	– 1k images
Sivic and Zisserman'03	– 5k images
Nister and Stewenius'06 (1M)	– 50k images
Philbin et al.'07	– 100k images
Chum et al.'07 + Jegou and Schmid'07	– 1M images
Chum et al.'08	– 5M images

Index 1 billion (10^9) images

– 200 servers each indexing 5M images?



Beyond Nearest Neighbors...

Indexing local features using inverted file index

Index

"Along I-75," From Detroit to Florida; *inside back cover*
"Drive I-95," From Boston to Florida; *inside back cover*
1929 Spanish Trail Roadway; 101-102,104
511 Traffic Information; 83
A1A (Barrier Isl) - I-95 Access; 86
AAA (and CAA); 83
AAA National Office; 88
Abbreviations,
 Colored 25 mile Maps; cover
 Exit Services; 196
 Travelogue; 85
Africa; 177
Agricultural Inspection Stns; 126
Ah-Tah-Thi-Ki Museum; 160
Air Conditioning, First; 112
Alabama; 124
Alachua; 132
 County; 131
Alafia River; 143
Alapaha, Name; 126
Alfred B Maclay Gardens; 106
Alligator Alley; 154-155
Alligator Farm, St Augustine; 169
Alligator Hole (definition); 157
Alligator, Buddy; 155
Alligators; 100,135,138,147,156
Anastasia Island; 170
Anhaica; 108-109,146
Apalachicola River; 112
Appleton Mus of Art; 136
Aquifer; 102
Arabian Nights; 94
Art Museum, Ringling; 147
Aruba Beach Cafe; 183
Aucilla River Project; 106
Babcock-Web WMA; 151
Bahia Mar Marina; 184
Baker County; 99
Barefoot Mailmen; 182
Barge Canal; 137
Bee Line Expy; 80
Belz Outlet Mall; 89
Bernard Castro; 136
Big "I"; 165
Big Cypress; 155,158
Big Foot Monster; 105
Billie Swamp Safari; 160
Blackwater River SP; 117
Blue Angels
 Butterfly Center, McGuire; 134
 CAA (see AAA)
 CCC, The; 111,113,115,135,142
 Ca d'Zan; 147
 Caloosahatchee River; 152
 Name; 150
 Canaveral Natnl Seashore; 173
 Cannon Creek Airpark; 130
 Canopy Road; 106,169
 Cape Canaveral; 174
 Castillo San Marcos; 169
 Cave Diving; 131
 Cayo Costa, Name; 150
 Celebration; 93
 Charlotte County; 149
 Charlotte Harbor; 150
 Chautauqua; 116
 ChIPLEY; 114
 Name; 115
 Choctawatchee, Name; 115
 Circus Museum, Ringling; 147
 Citrus; 88,97,130,136,140,180
 CityPlace, W Palm Beach; 180
 City Maps,
 Fl Lauderdale Expwys; 194-195
 Jacksonville; 163
 Kissimmee Expwys; 192-193
 Miami Expressways; 194-195
 Orlando Expressways; 192-193
 Pensacola; 26
 Tallahassee; 191
 Tampa-St. Petersburg; 63
 St. Augustine; 191
 Civil War; 100,108,127,138,141
 Clearwater Marine Aquarium; 187
 Collier County; 154
 Collier, Barron; 152
 Colonial Spanish Quarters; 168
 Columbia County; 101,128
 Coquina Building Material; 165
 Corkscrew Swamp, Name; 154
 Cowboys; 85
 Crab Trap II; 144
 Cracker, Florida; 88,95,132
 Crosstown Expy; 11,35,98,143
 Cuban Bread; 184
 Dade Battlefield; 140
 Dade, Maj. Francis; 139-140,161
 Dania Beach Hurricane; 184
 Daniel Boone, Florida Walk; 117
 Daytona Beach; 172-173
 De Land; 87
 Driving Lanes; 85
 Duval County; 163
 Eau Gallie; 175
 Edison, Thomas; 152
 Eglin AFB; 116-118
 Eight Reale; 176
 Ellenton; 144-145
 Emanuel Point Wreck; 120
 Emergency Callboxes; 83
 Epiphytes; 142,148,157,159
 Escambia Bay; 119
 Bridge (I-10); 119
 County; 120
 Estero; 153
 Everglade,90,95,139-140,154-160
 Draining of; 156,181
 Wildlife MA; 160
 Wonder Gardens; 154
 Falling Waters SP; 115
 Fantasy of Flight; 95
 Fayer Dykes SP; 171
 Fires, Forest; 168
 Fires, Prescribed ; 148
 Fisherman's Village; 151
 Flagler County; 171
 Flagler, Henry; 97,165,167,171
 Florida Aquarium; 186
 Florida,
 12,000 years ago; 187
 Cavern SP; 114
 Map of all Expressways; 2-3
 Mus of Natural History; 134
 National Cemetery ; 141
 Part of Africa; 177
 Platform; 187
 Sheriff's Boys Camp; 126
 Sports Hall of Fame; 130
 Sun 'n Fun Museum; 97
 Supreme Court; 107
 Florida's Turnpike (FTP), 178,189
 25 mile Strip Maps; 66
 Administration; 189
 Coin System; 190
 Exit Services; 189
 HEFT; 76,161,190
 History; 189
 Names; 189
 Service Plazas; 190
 Spur SR91; 76
 Ticket System; 190
 Toll Plazas; 190
 Ford, Henry; 152

For text documents, an efficient way to find all *pages* on which a *word* occurs is to use an index...

We want to find all *images* in which a *feature* occurs.

To use this idea, we'll need to map our features to "visual words".

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 visual impressions which reach the brain from the eye.

thought that the point by which the cerebral cortex is connected upon which the visual image is formed.

Through the eye, cell, optical nerve, image

more complex. Hubel, Wiesel
the visual input. Hubel and Wiesel
cell layers of the optic tectum. Hubel and Wiesel
have been able to demonstrate that the visual image
*about the image falling on the retina undergoes a
step-wise analysis in a system of nerve cells organized
in columns. In this system each cell has its own
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 \$100bn, a 30% jump in exports to \$100bn, a 30% rise in

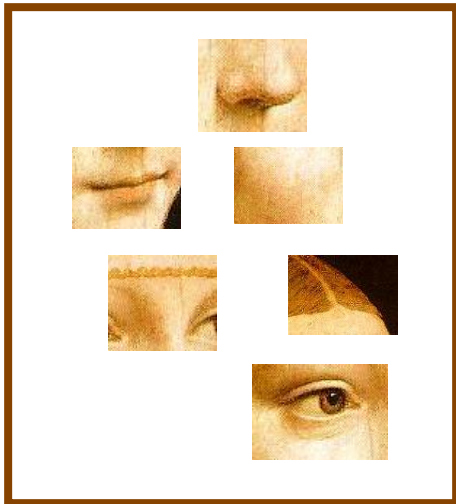
imports to \$100bn, a 30% rise in further a 30% rise in at China's 30% rise in at deliber 30% rise in surplus 30% rise in factor. He 30% rise in said the 30% rise in domestic 30% rise in the country. China 30% rise in against the dollar 30% rise in it to trade within a narrow 30% rise in wants the yuan to be allowed to trade freely. 30% rise in er, Beijing has made it clear that it will take 30% rise in and tread carefully before allowing the yuan 30% rise in 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



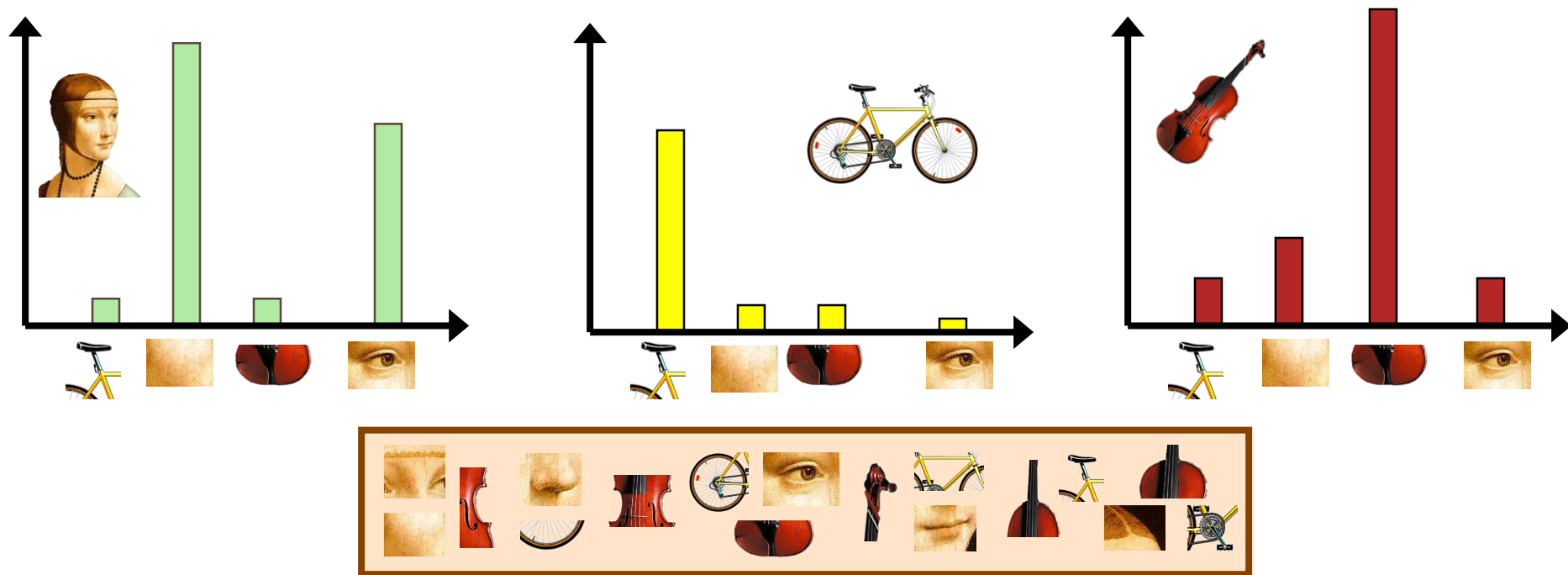
A clarification: definition of “BoW”

Looser definition

- Independent features

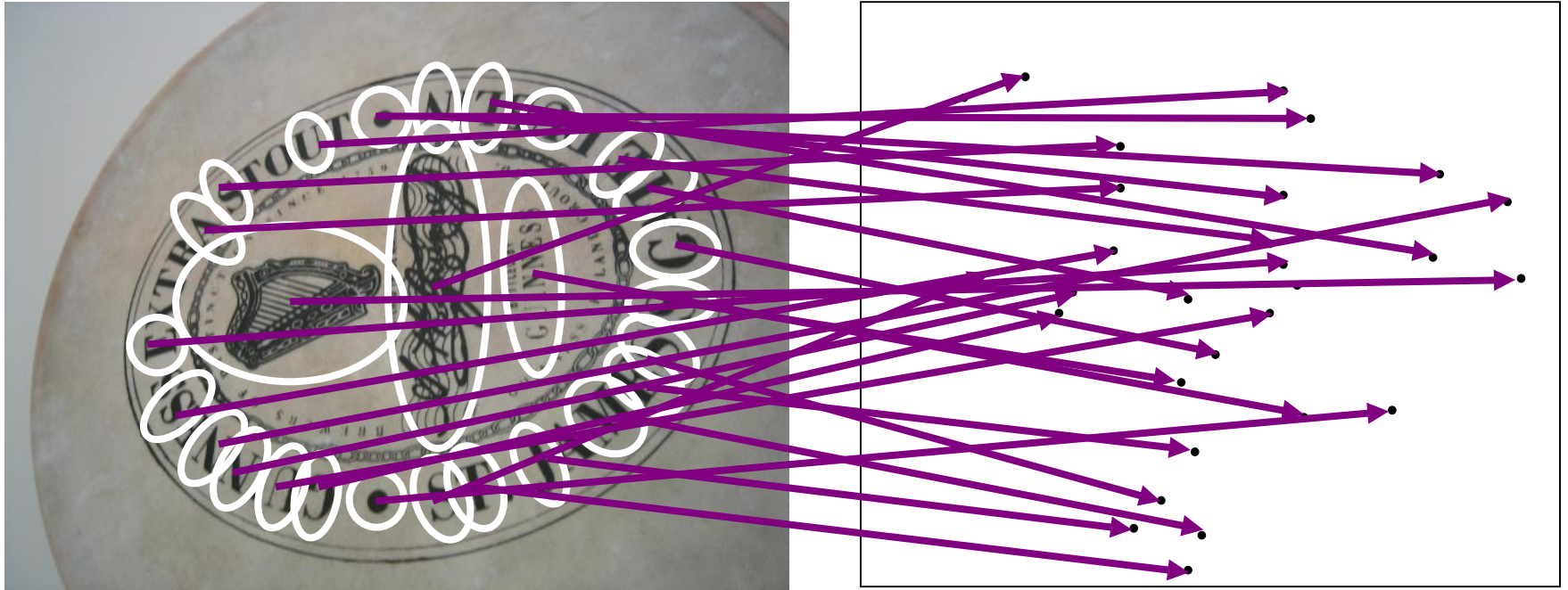
Stricter definition

- Independent features
- histogram representation



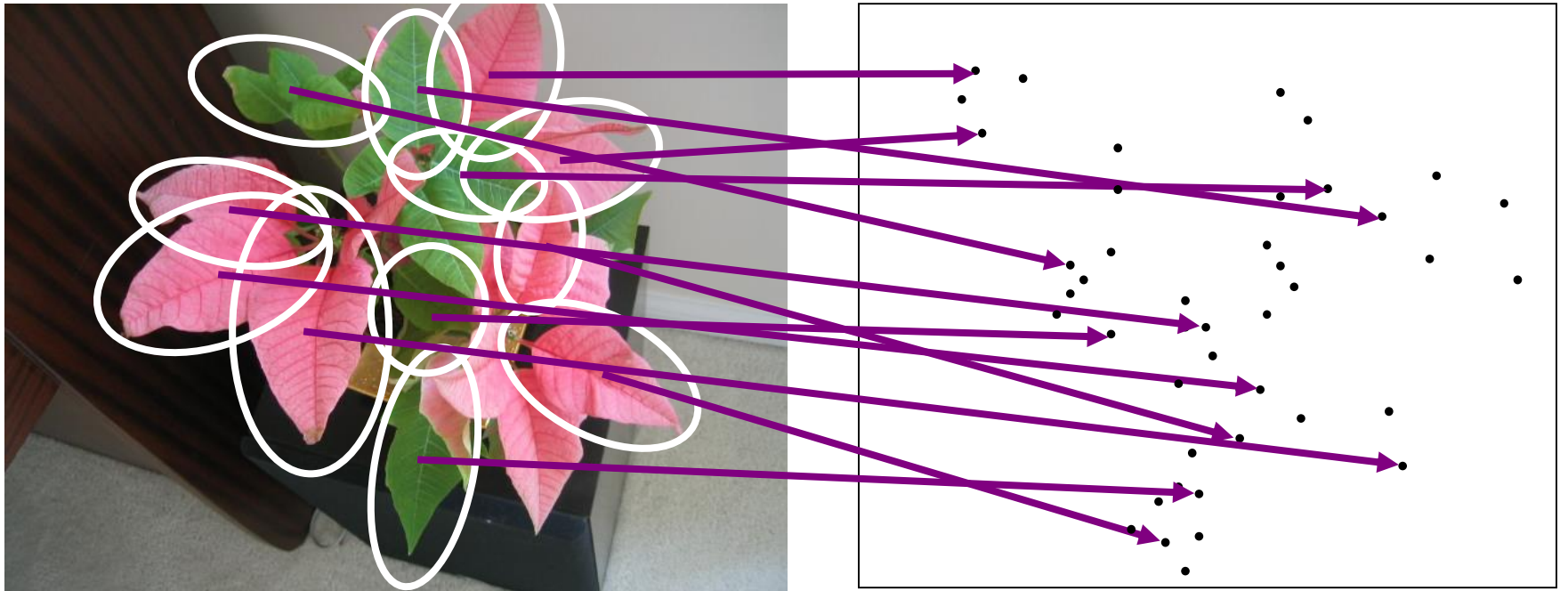
Visual words: main idea

Extract some local features from a number of images ...

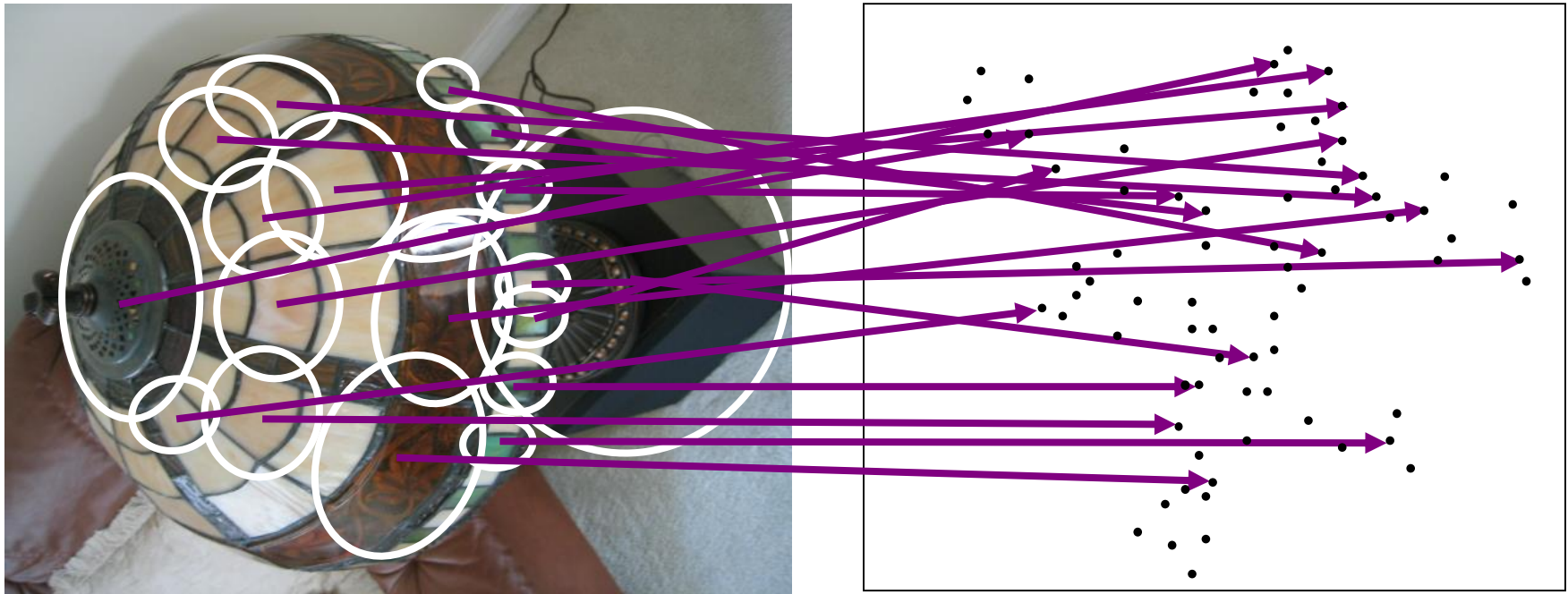


e.g., SIFT descriptor space:
each point is 128-dimensional

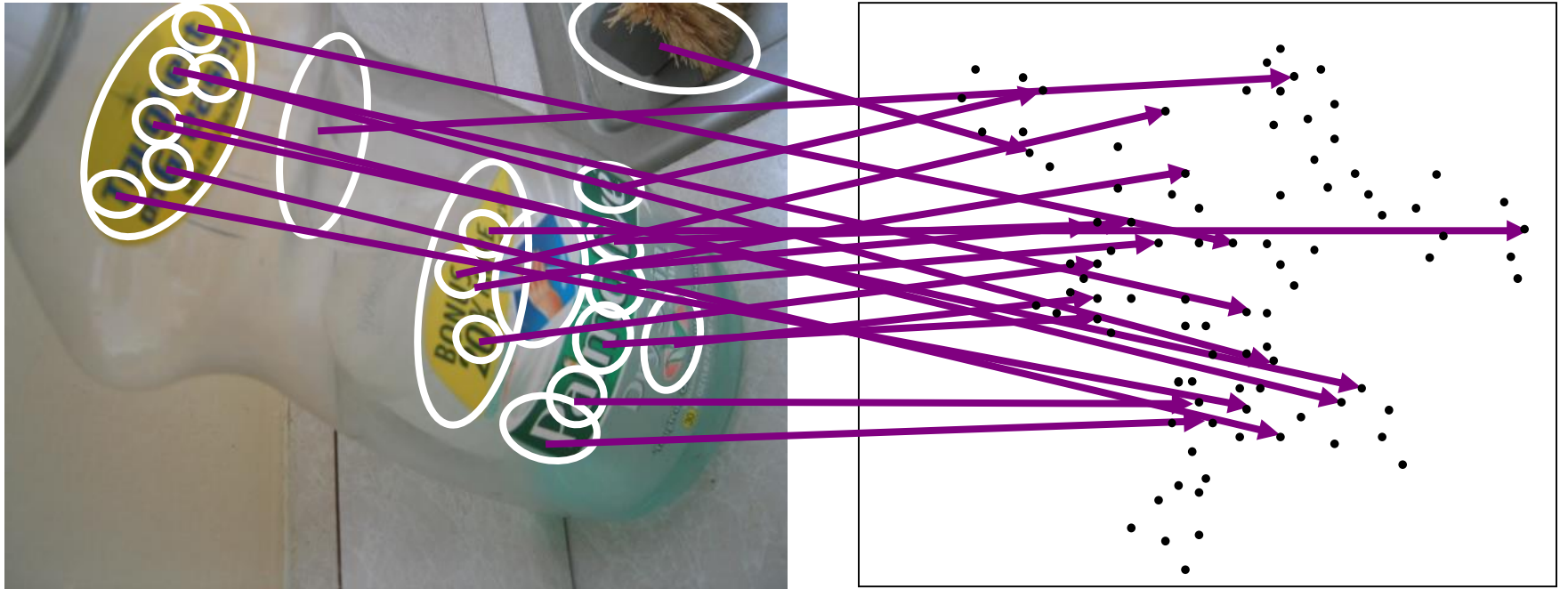
Visual words: main idea

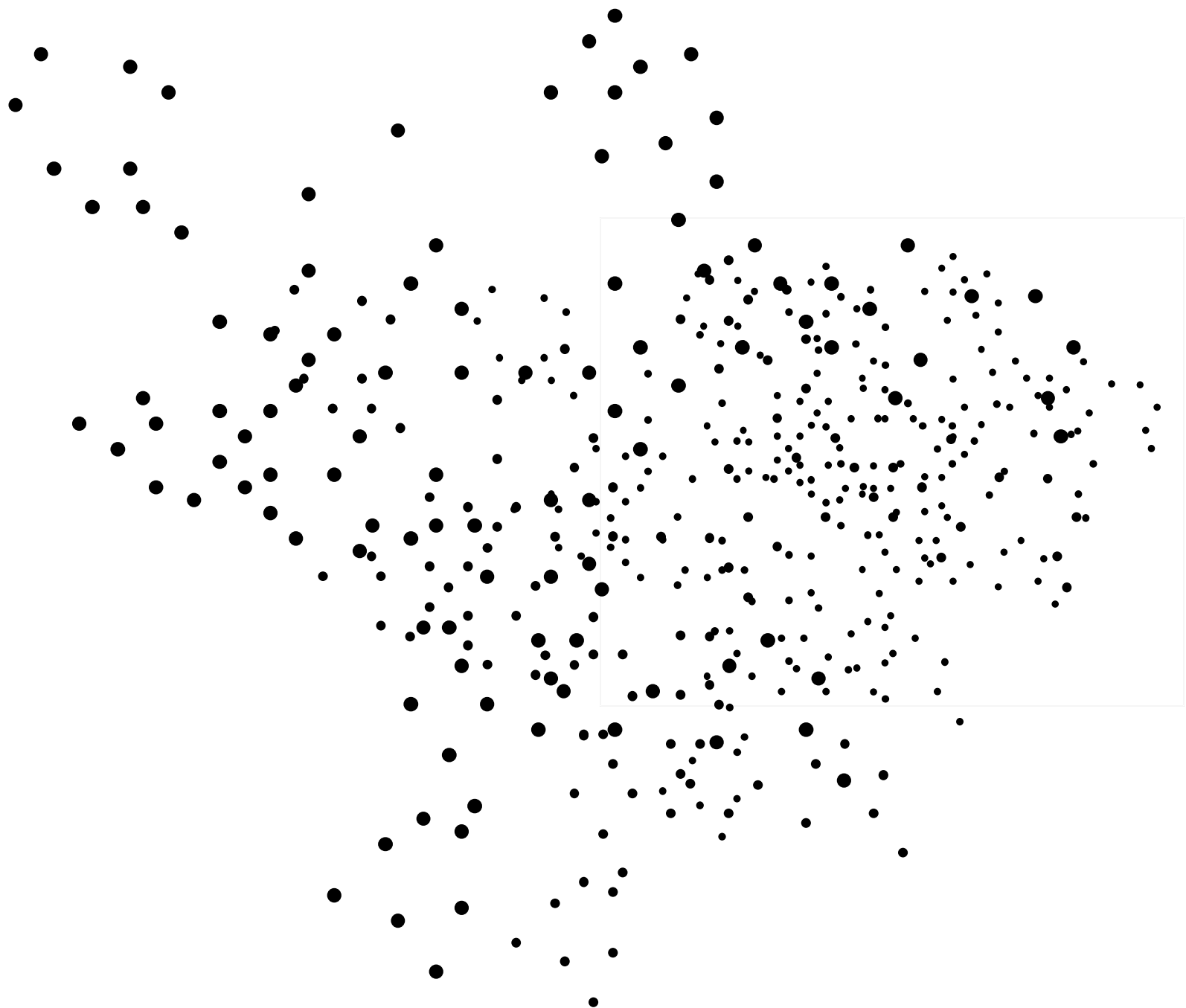


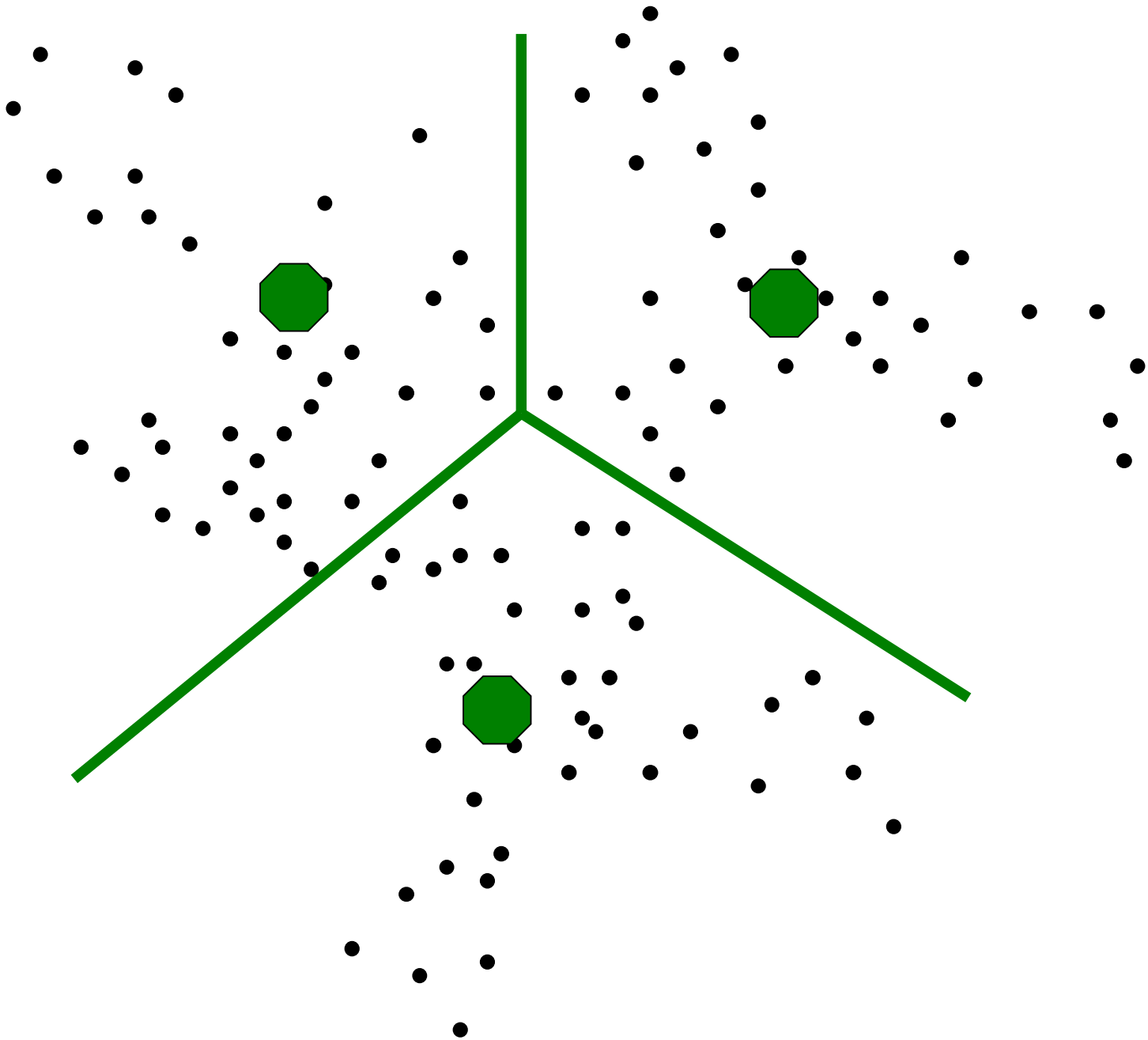
Visual words: main idea



Visual words: main idea

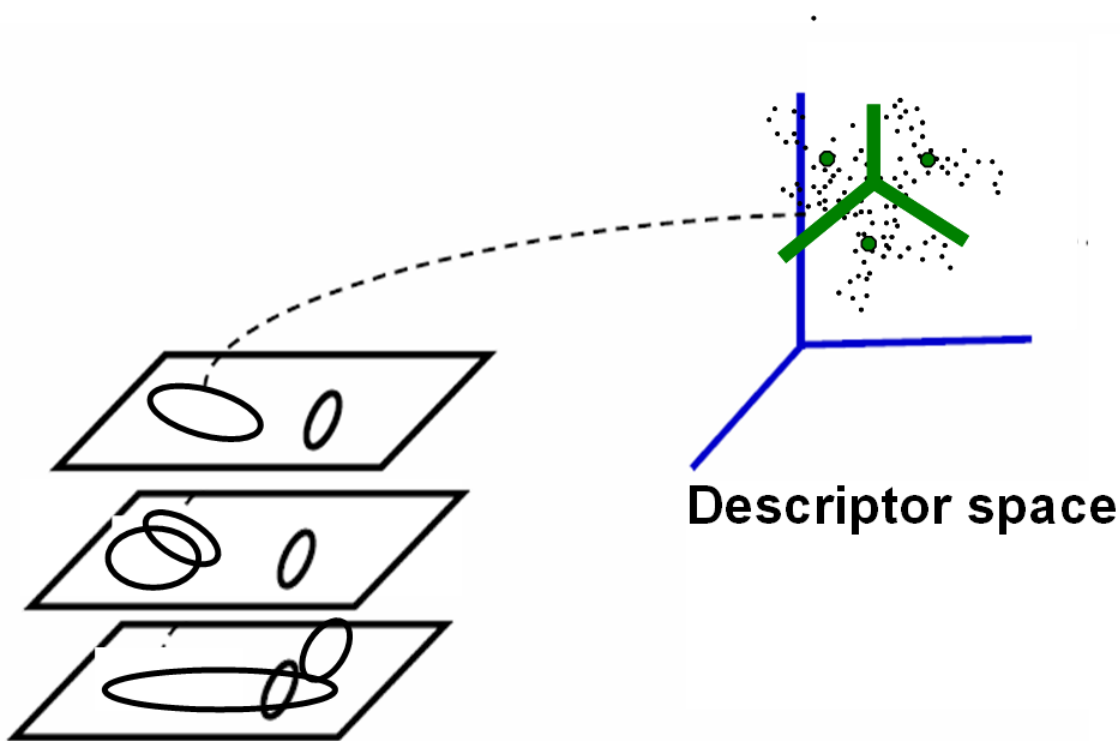






Visual words: main idea

Map high-dimensional descriptors to tokens/words by quantizing the feature space

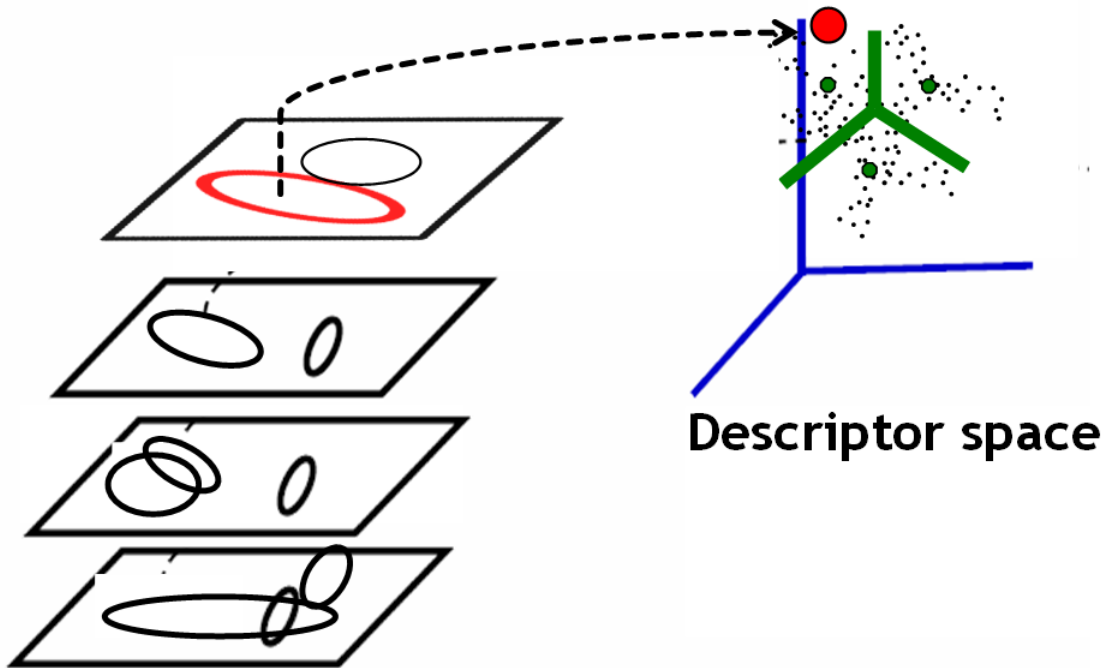


- Quantize via clustering, let cluster centers be the prototype “words”



Visual words: main idea

Map high-dimensional descriptors to tokens/words by quantizing the feature space

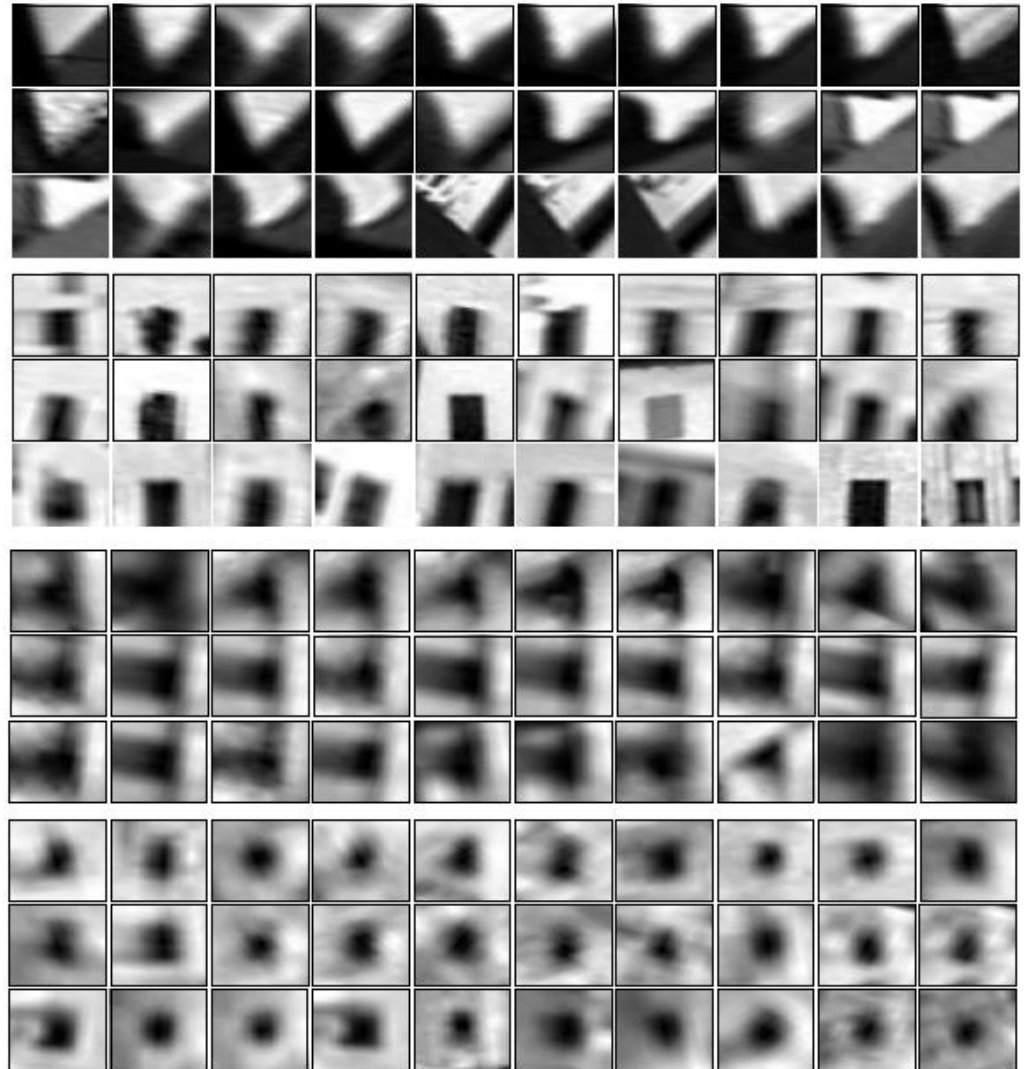


- Determine which word to assign to each new image region by finding the closest cluster center.



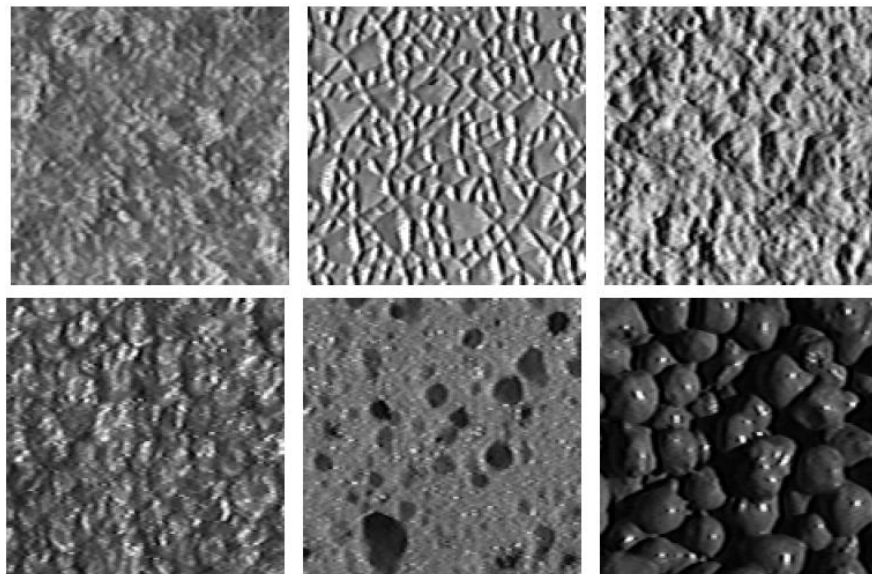
Visual words

Example: each group of patches belongs to the same visual word



Visual words

- First explored for texture and material representations
- *Texon* = cluster center of filter responses over collection of images
- Describe textures and materials based on distribution of prototypical texture elements.



Leung & Malik 1999; Varma & Zisserman, 2002; Lazebnik, Schmid & Ponce, 2003;

Inverted file index for images comprised of visual words



frame #5



frame #10

Word number	List of image numbers
-------------	-----------------------

1	→ 5, 10, ...
2	→ 10, ...
...	...

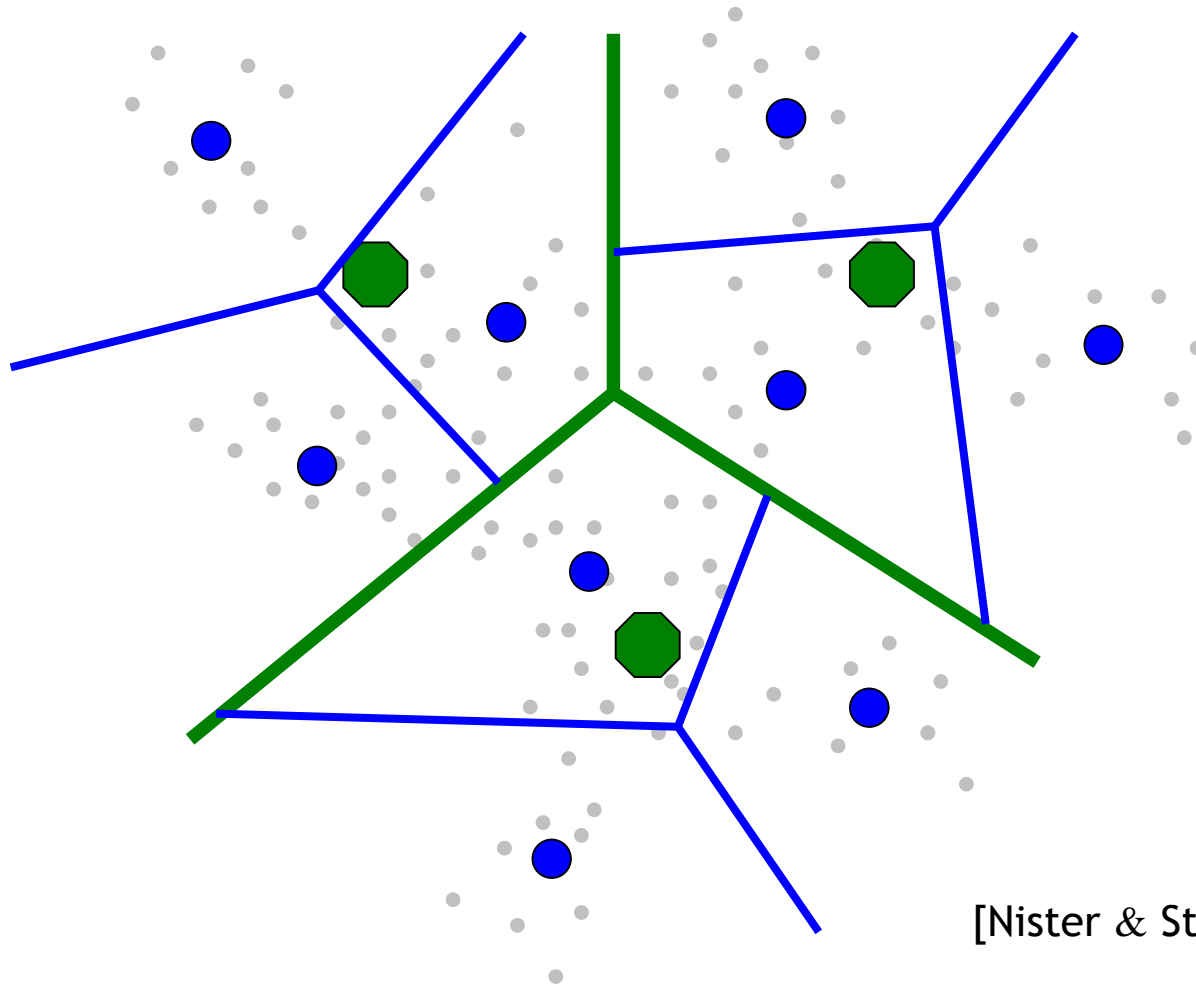
- Score each image by the number of common visual words (tentative correspondences)
- But: does not take into account spatial layout of regions

Clustering / quantization methods

- k-means (typical choice), agglomerative clustering, mean-shift,...
- **Hierarchical clustering: allows faster insertion / word assignment while still allowing large vocabularies**
 - **Vocabulary tree [Nister & Stewenius, CVPR 2006]**

Example: Recognition with Vocabulary Tree

Tree construction:

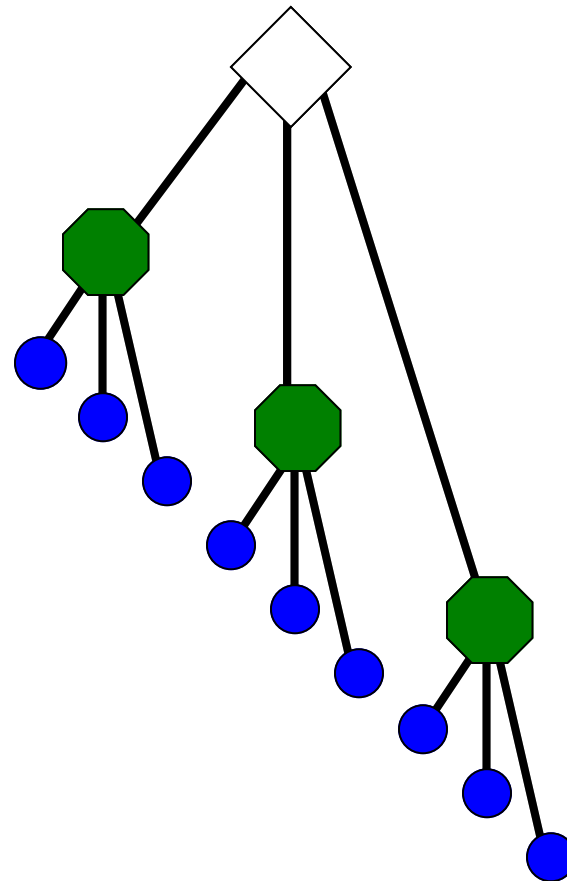


[Nister & Stewenius, CVPR'06]

Slide credit: David Nister

Vocabulary Tree

Training: Filling the tree

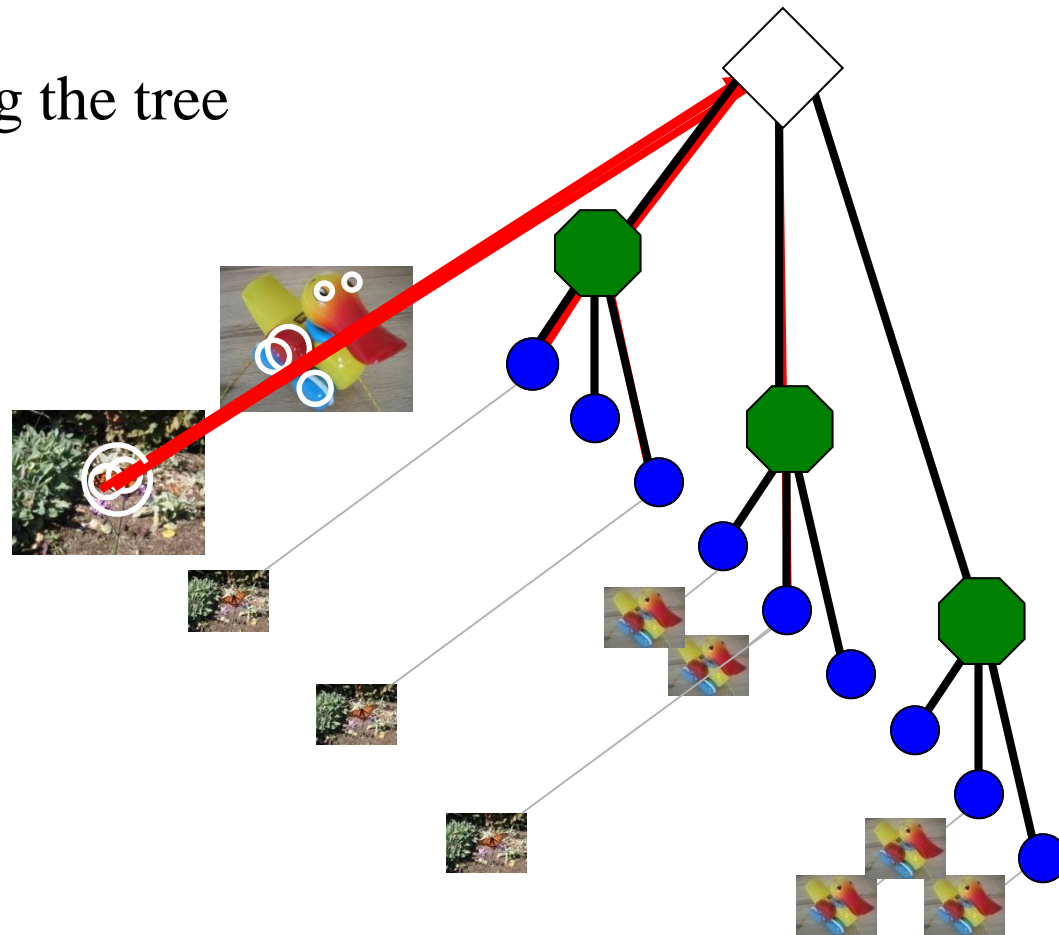


[Nister & Stewenius, CVPR'06]

Slide credit: David Nister

Vocabulary Tree

Training: Filling the tree

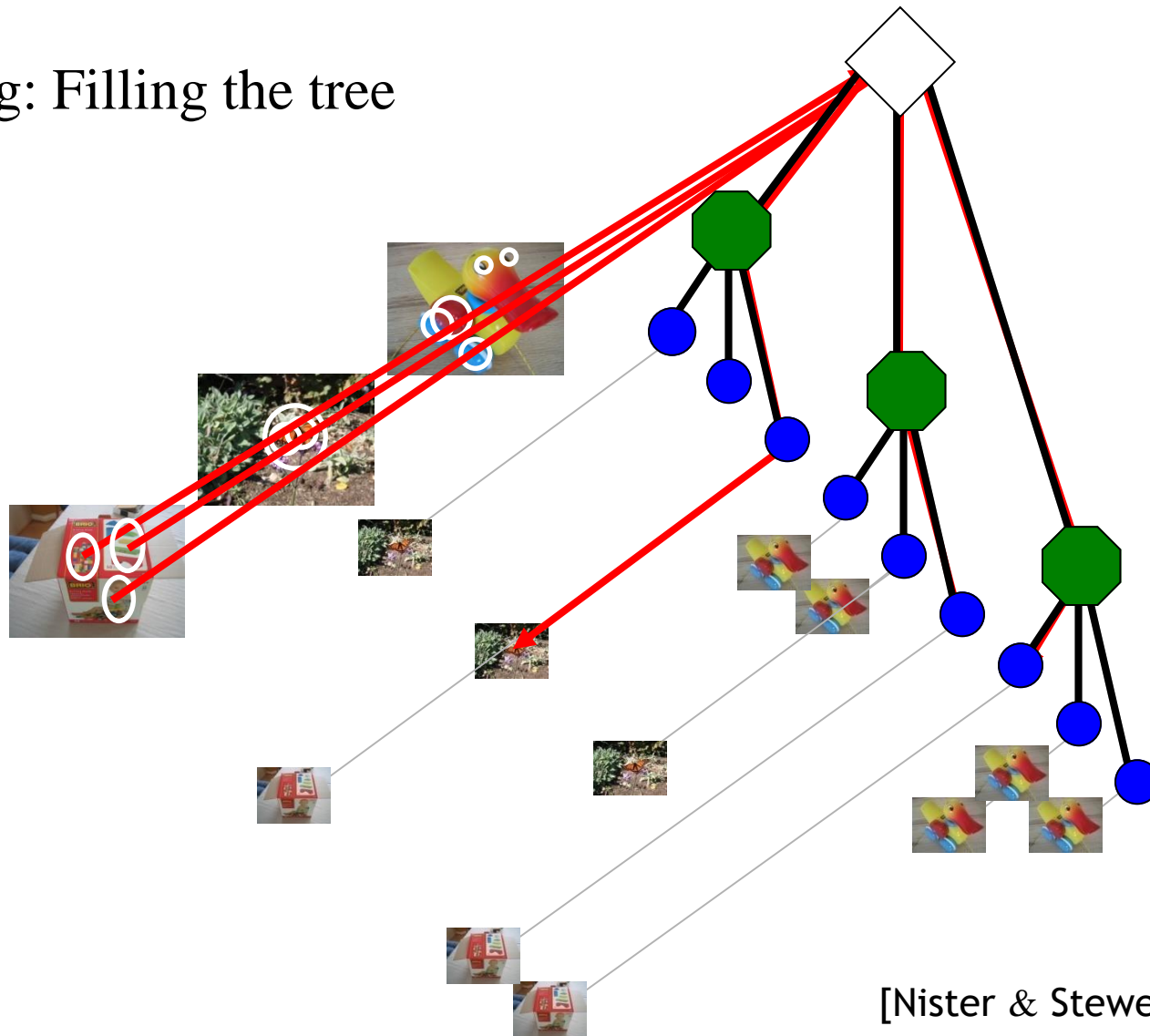


[Nister & Stewenius, CVPR'06]

Slide credit: David Nister

Vocabulary Tree

Training: Filling the tree

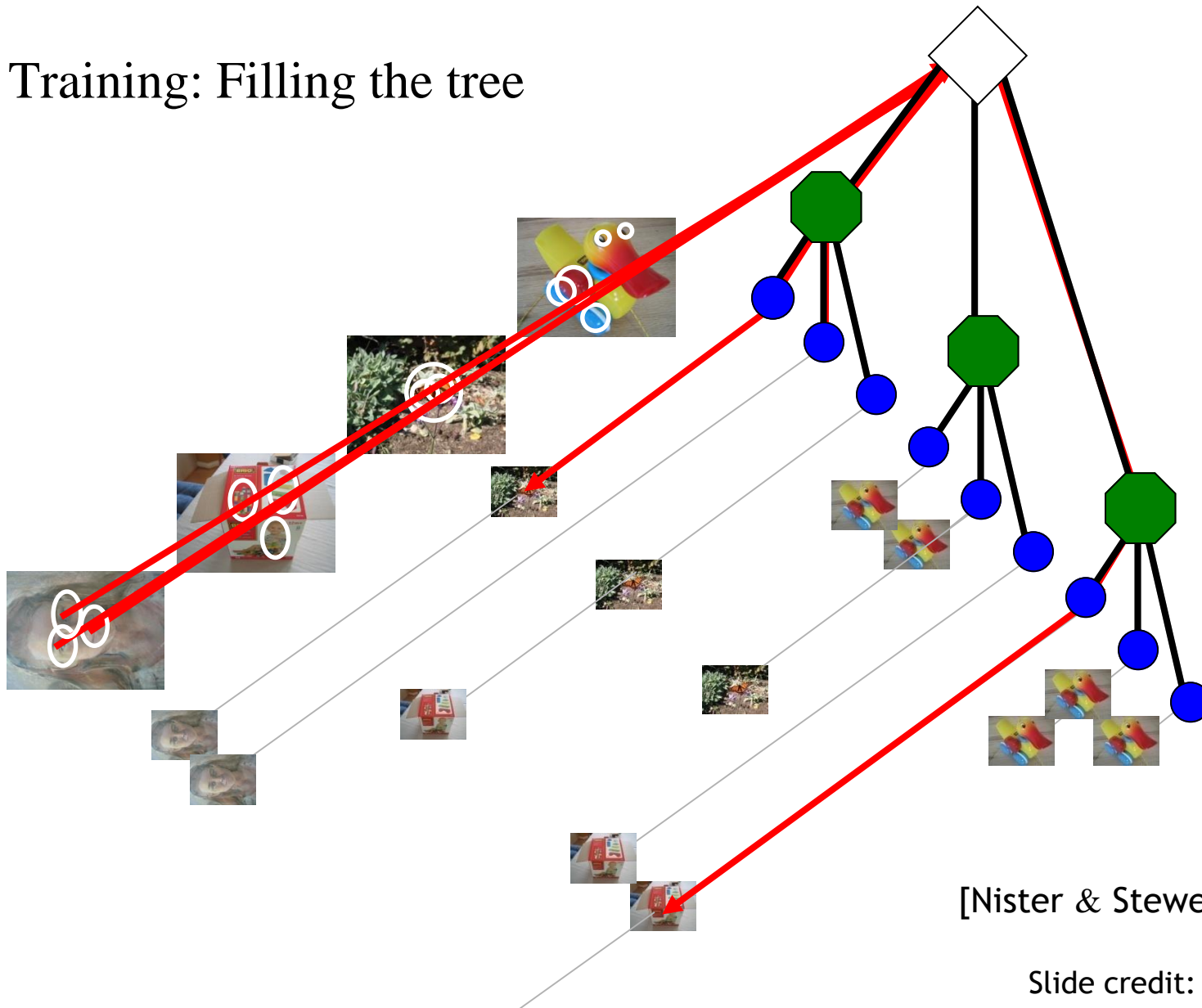


[Nister & Stewenius, CVPR'06]

Slide credit: David Nister

Vocabulary Tree

Training: Filling the tree



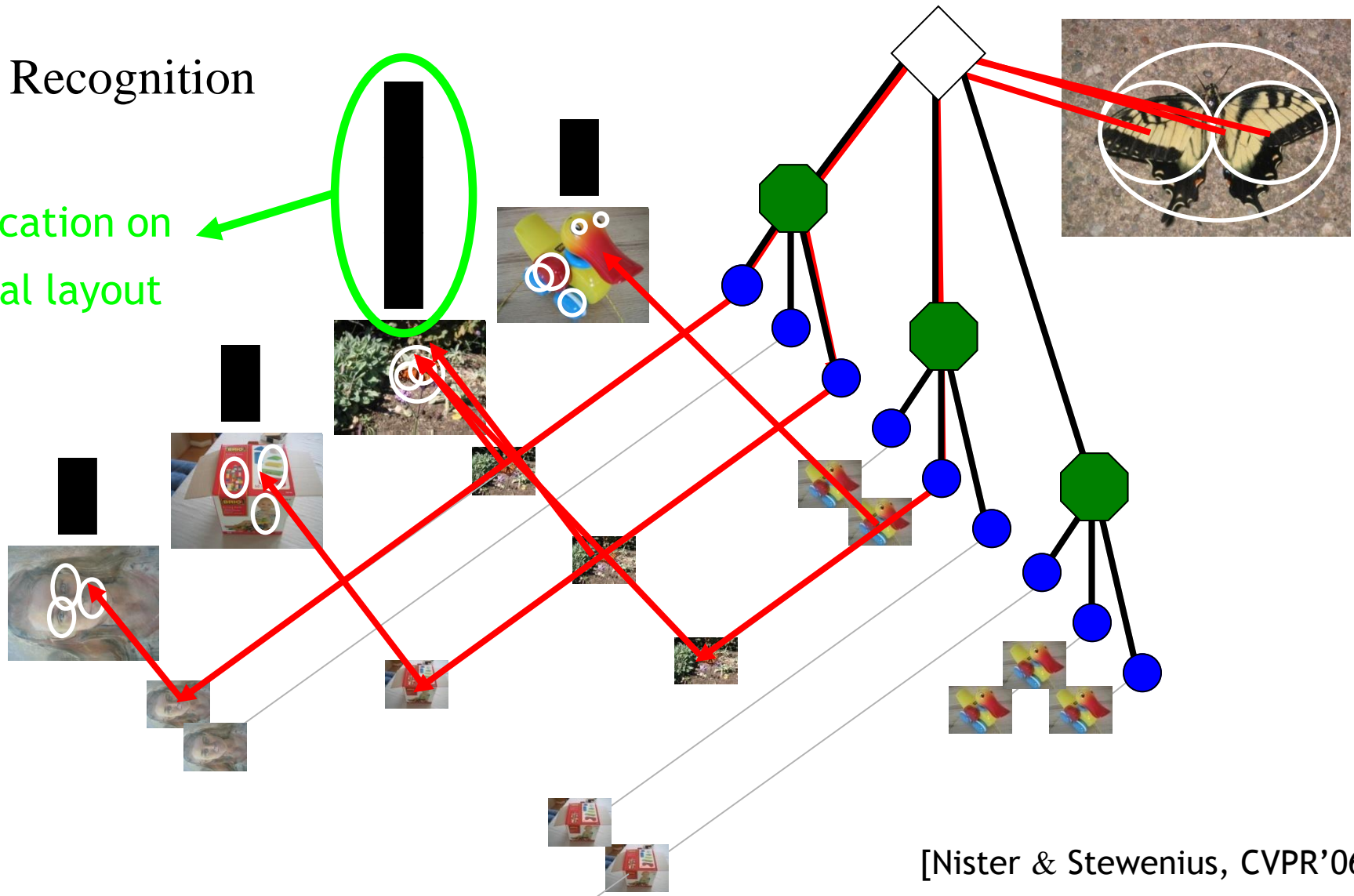
[Nister & Stewenius, CVPR'06]

Slide credit: David Nister

Vocabulary Tree

Recognition

Verification on spatial layout



[Nister & Stewenius, CVPR'06]

Slide credit: David Nister

Vocabulary Tree: Performance

Evaluated on large databases

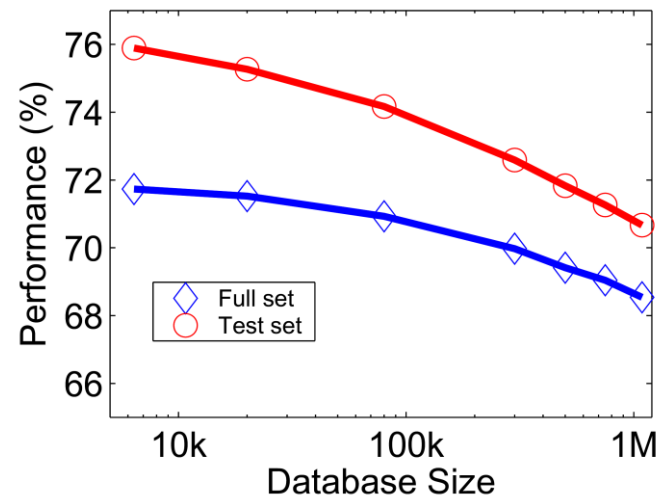
- Indexing with up to 1M images

Online recognition for database of 50,000 CD covers

- Retrieval in ~1s

Find experimentally that large vocabularies can be beneficial for recognition

[Nister & Stewenius, CVPR'06]



“Bag of visual words”

