

Henderson and Davis.
Shape recognition using hierarchical
Constraint Analysis. 1979

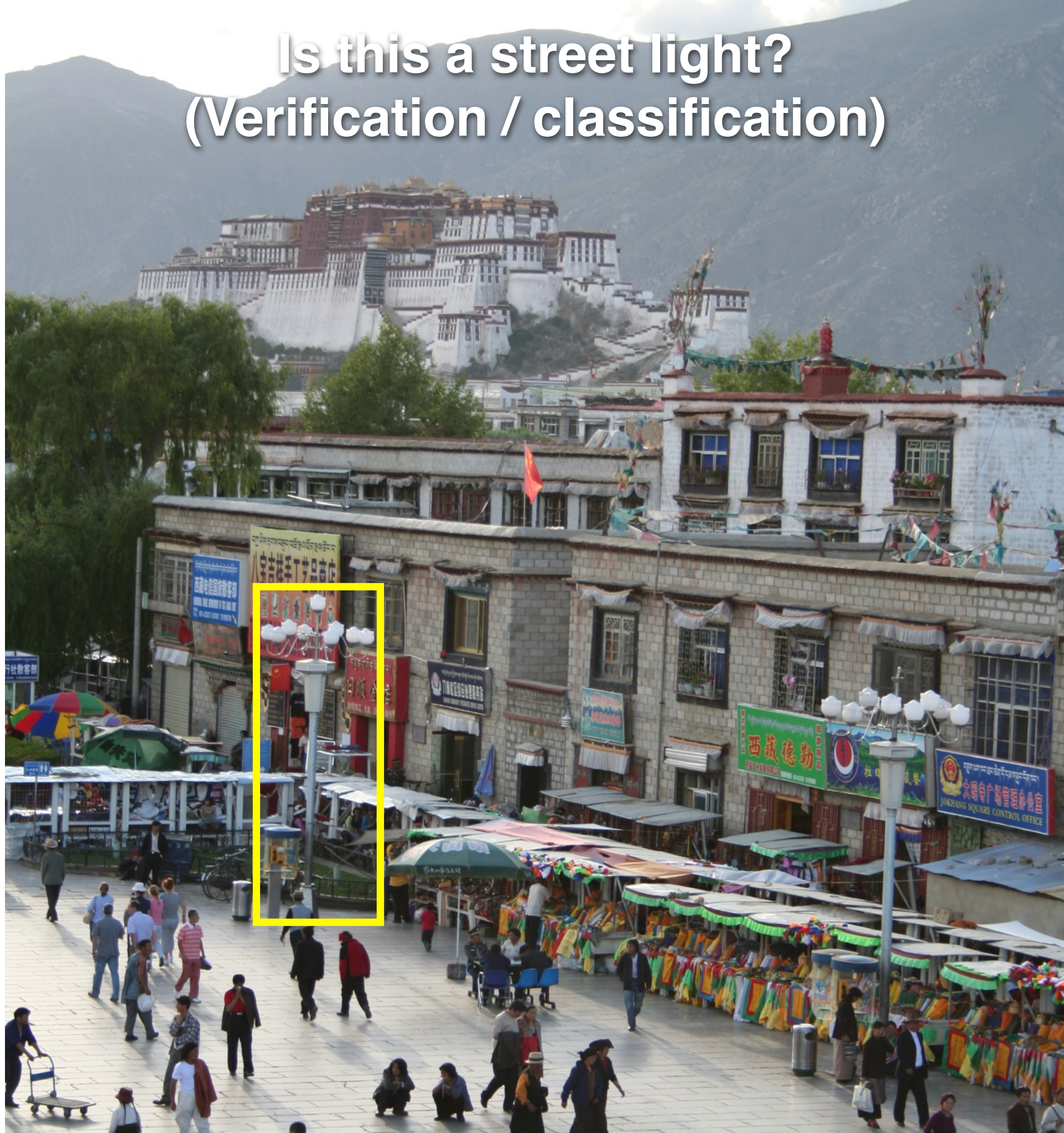
Object Recognition

16-385 Computer Vision (Kris Kitani)

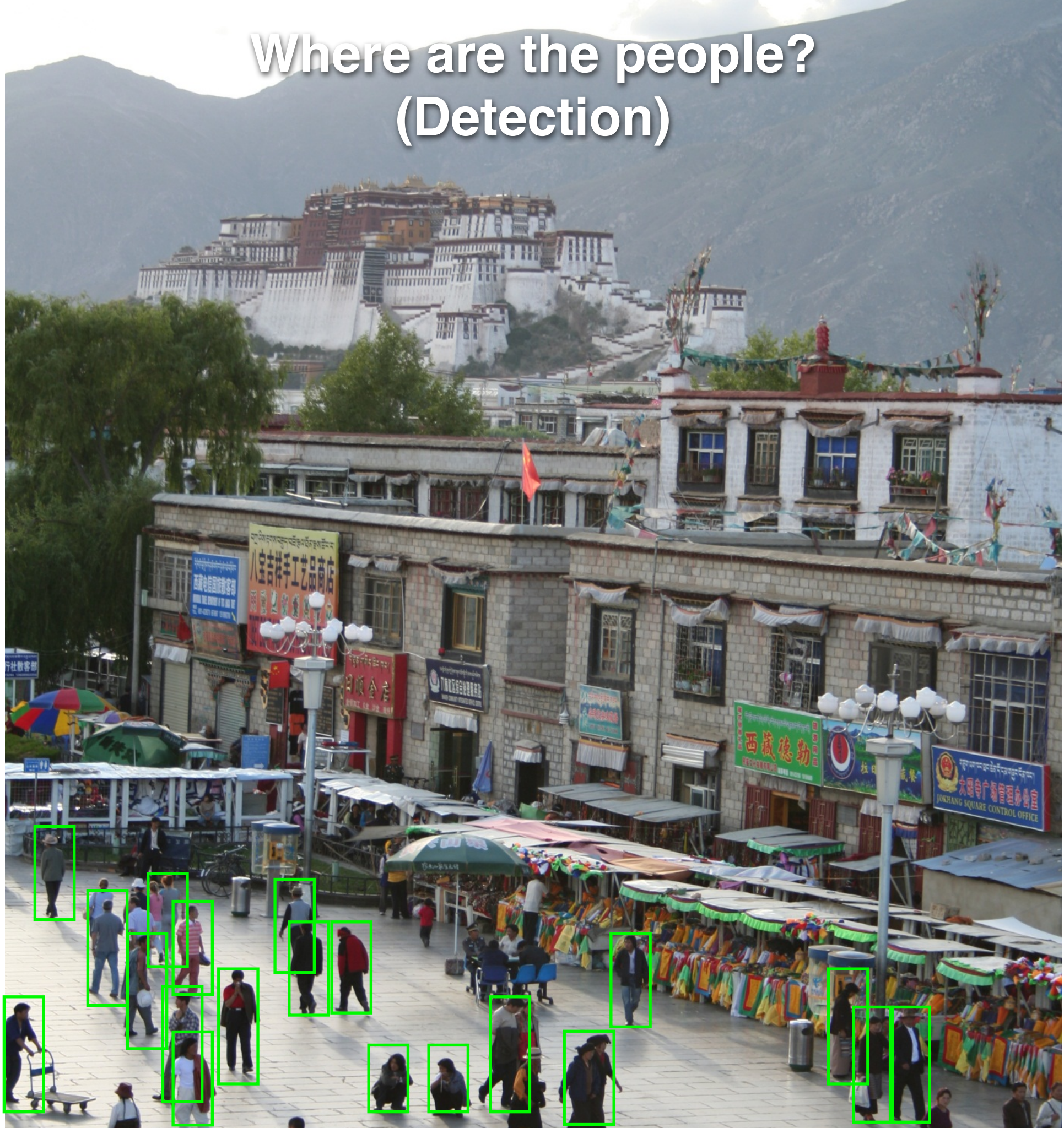
Carnegie Mellon University

What do we mean by
'object recognition'?

Is this a street light?
(Verification / classification)



Where are the people? (Detection)



Is that Potala palace? (Identification)



Sky

What's in the scene? (semantic segmentation)

Mountain

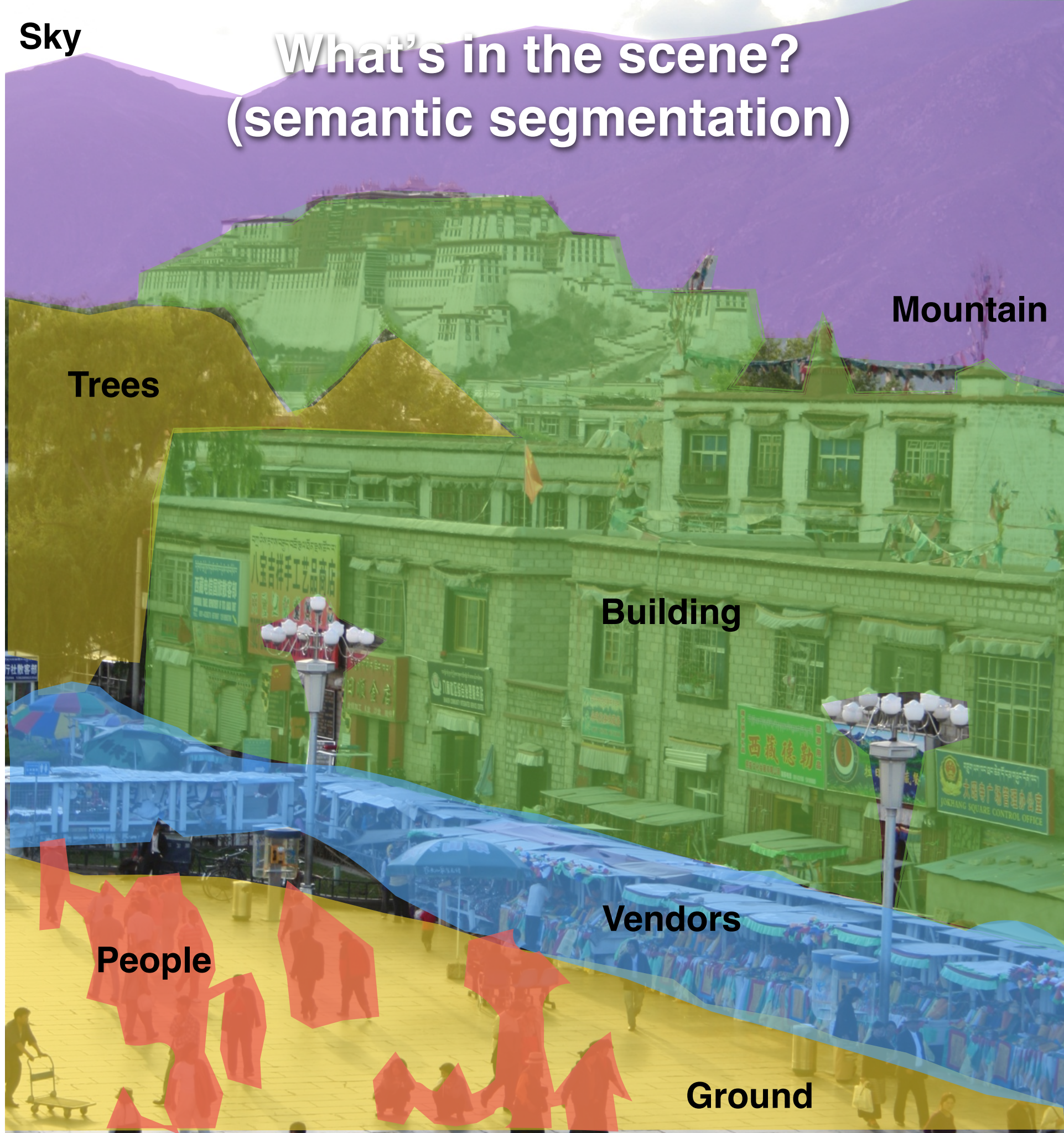
Trees

Building

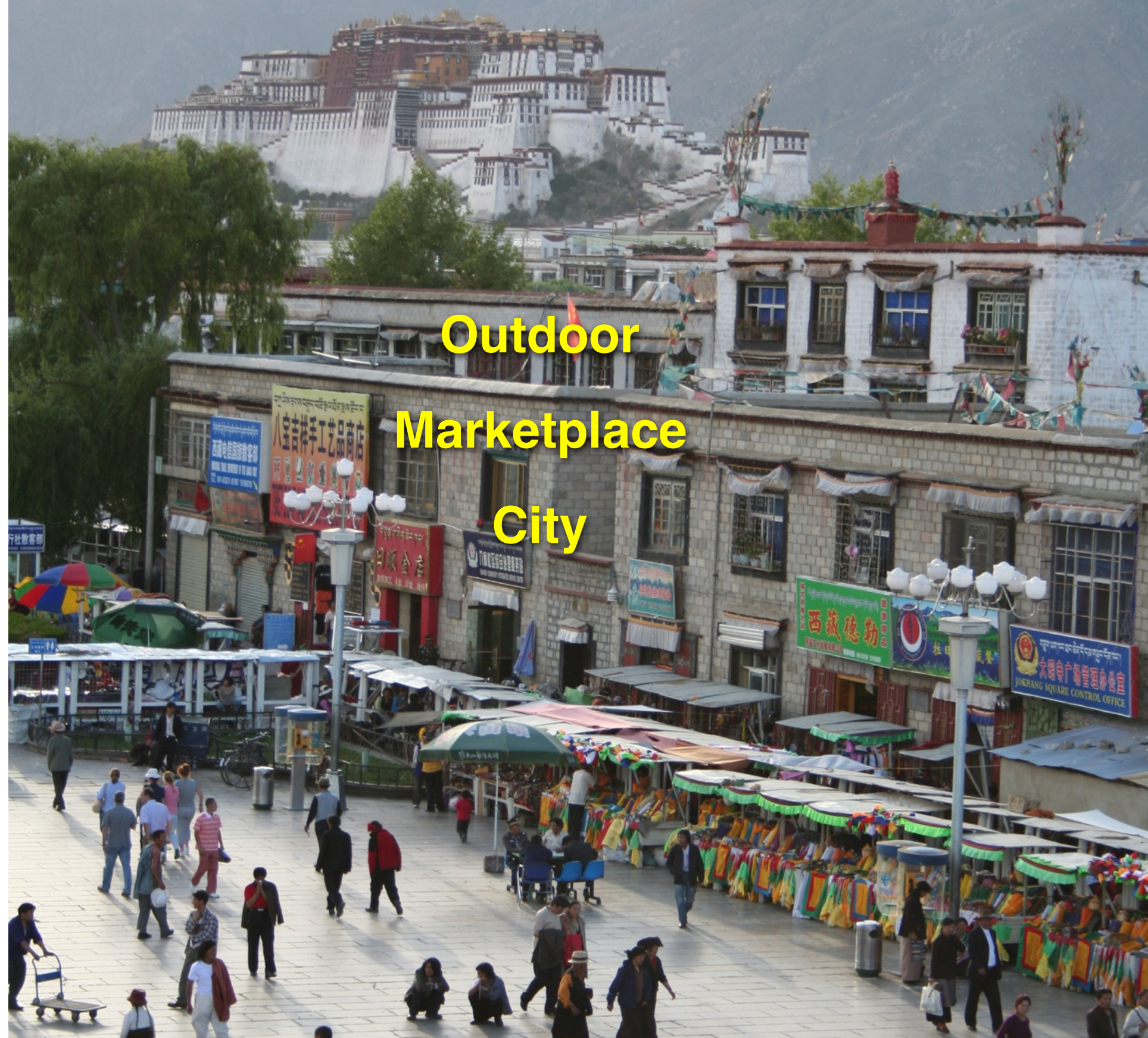
Vendors

People

Ground



What type of scene is it?
(Scene categorization)



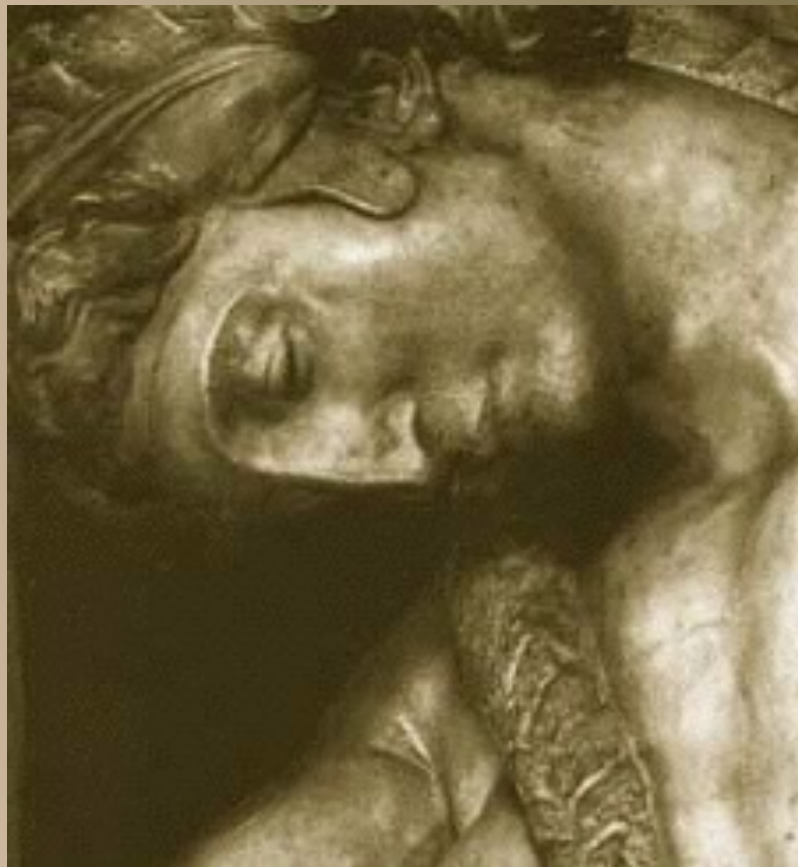
Outdoor

Marketplace

City

Challenges

(Object Recognition)



Viewpoint variation



Illumination variation





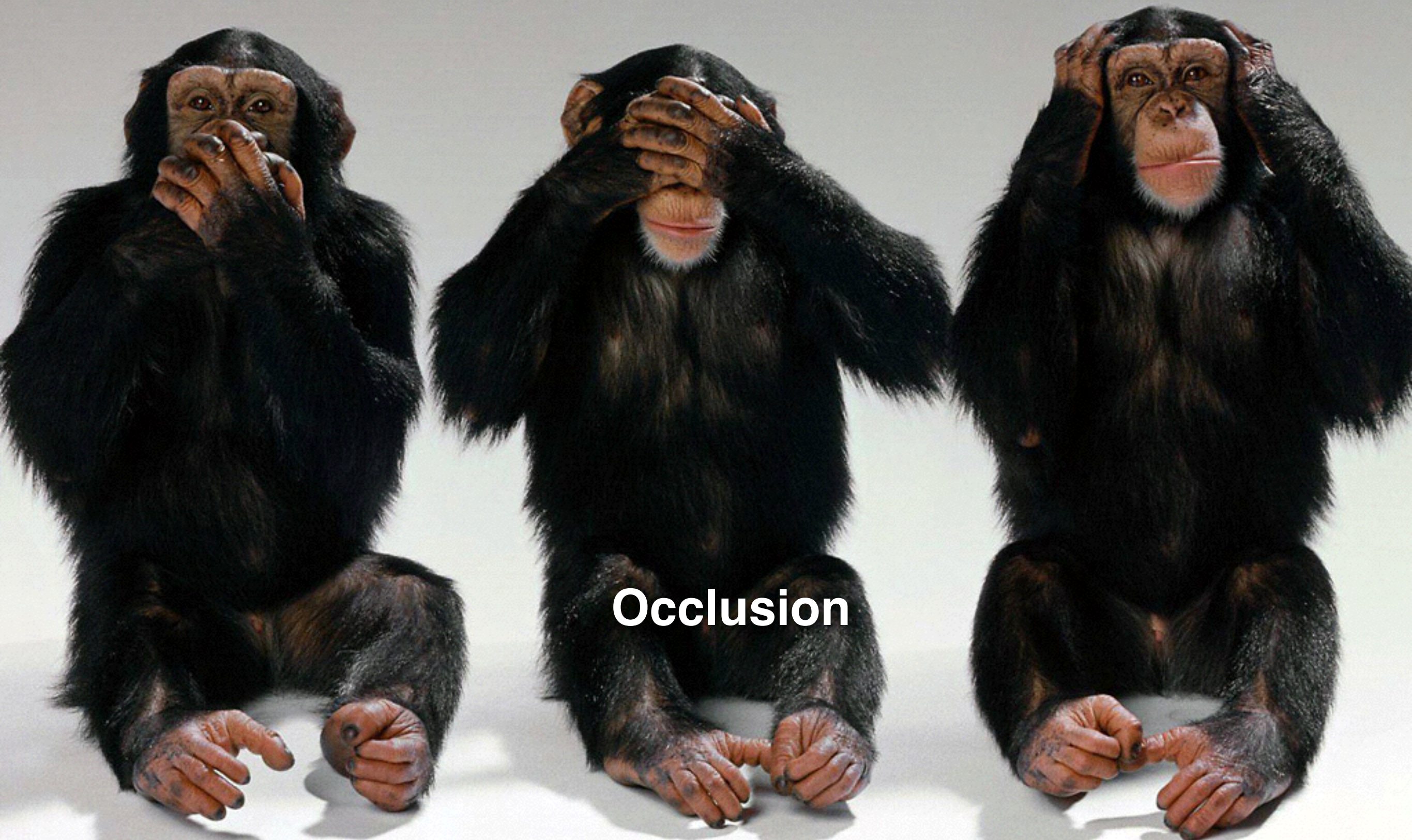
Scale variation



Background clutter



Deformation



Occlusion



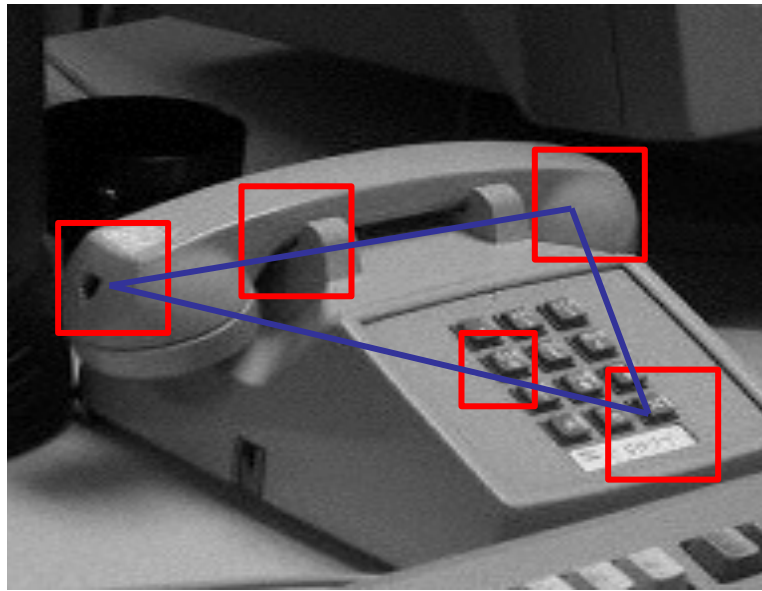
Intra-class variation

Common approaches

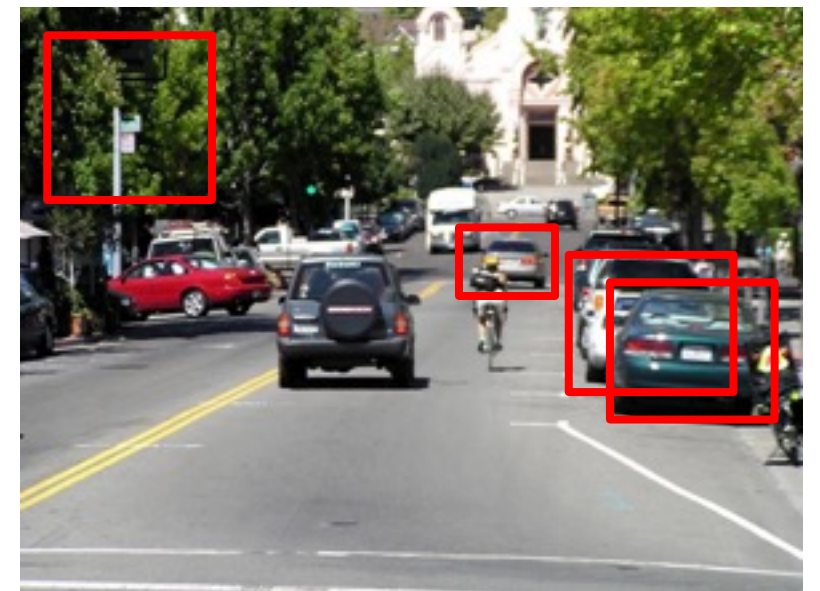
Common approaches: object recognition



Feature
Matching



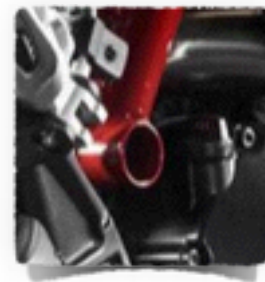
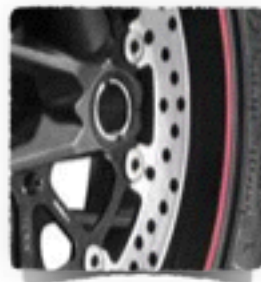
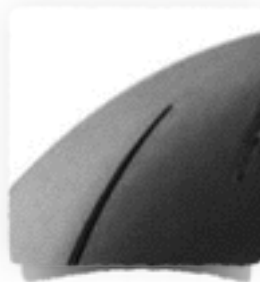
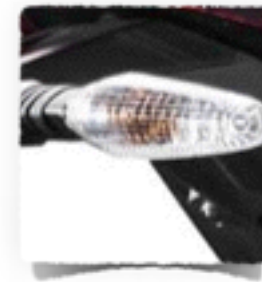
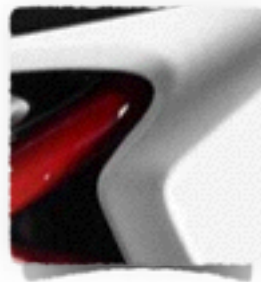
Spatial
reasoning



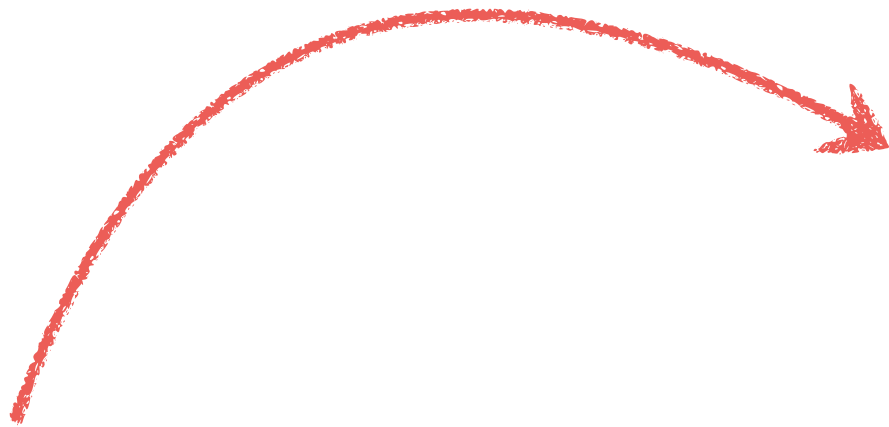
Window
classification

Feature matching

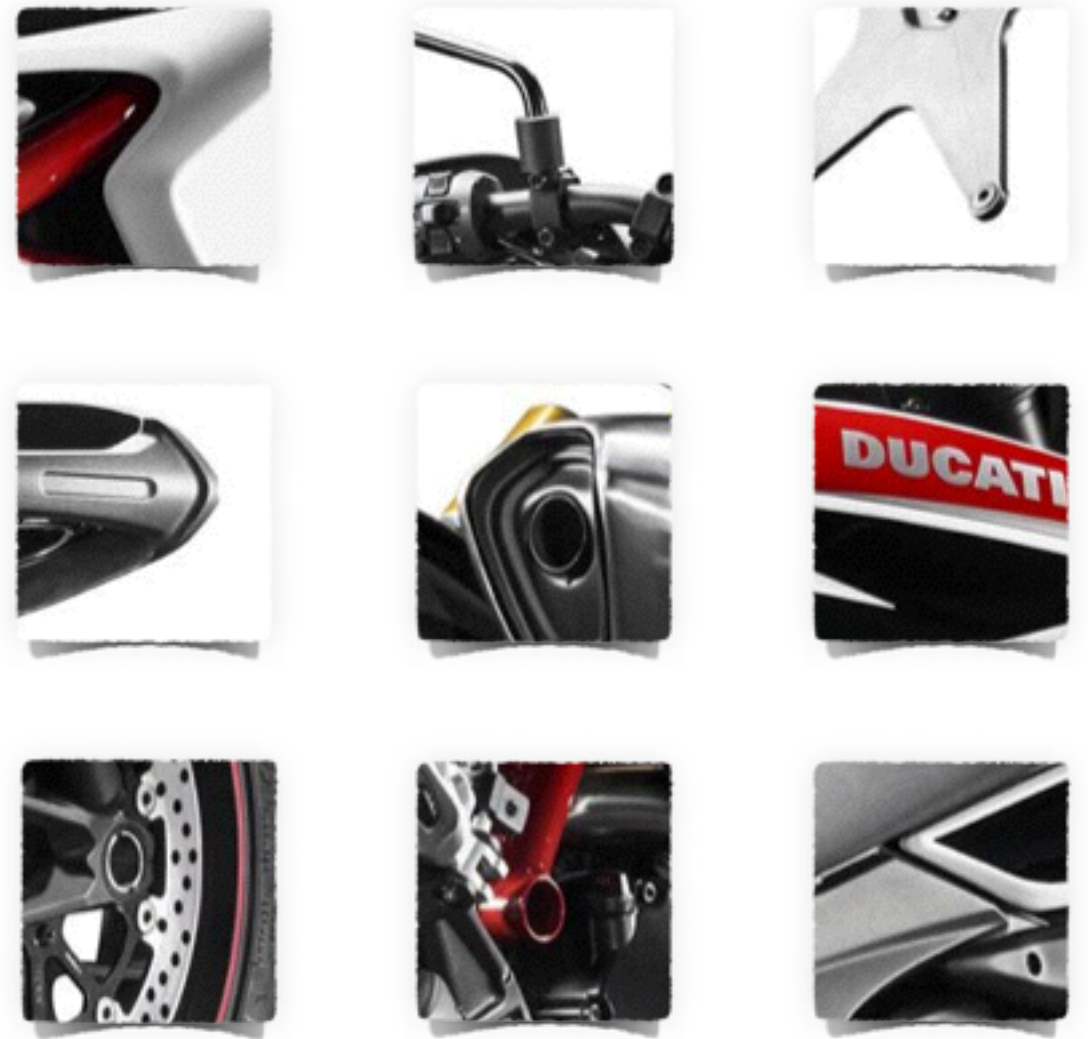
What object do these parts belong to?



Some local feature are
very informative



An object as



a collection of local features
(bag-of-features)

- deals well with occlusion
- scale invariant
- rotation invariant

Are the positions of the parts important?

Pros

- Simple
- Efficient algorithms
- Robust to deformations

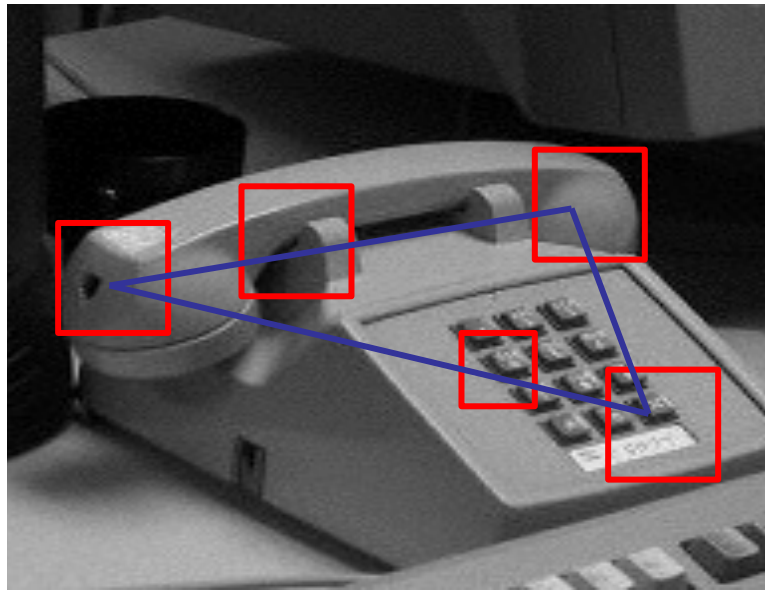
Cons

- No spatial reasoning

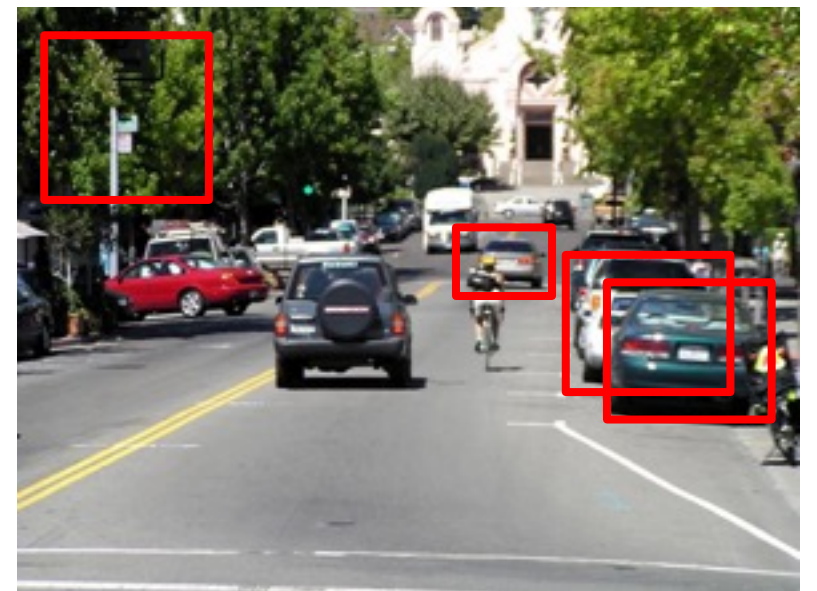
Common approaches: object recognition



Feature
Matching



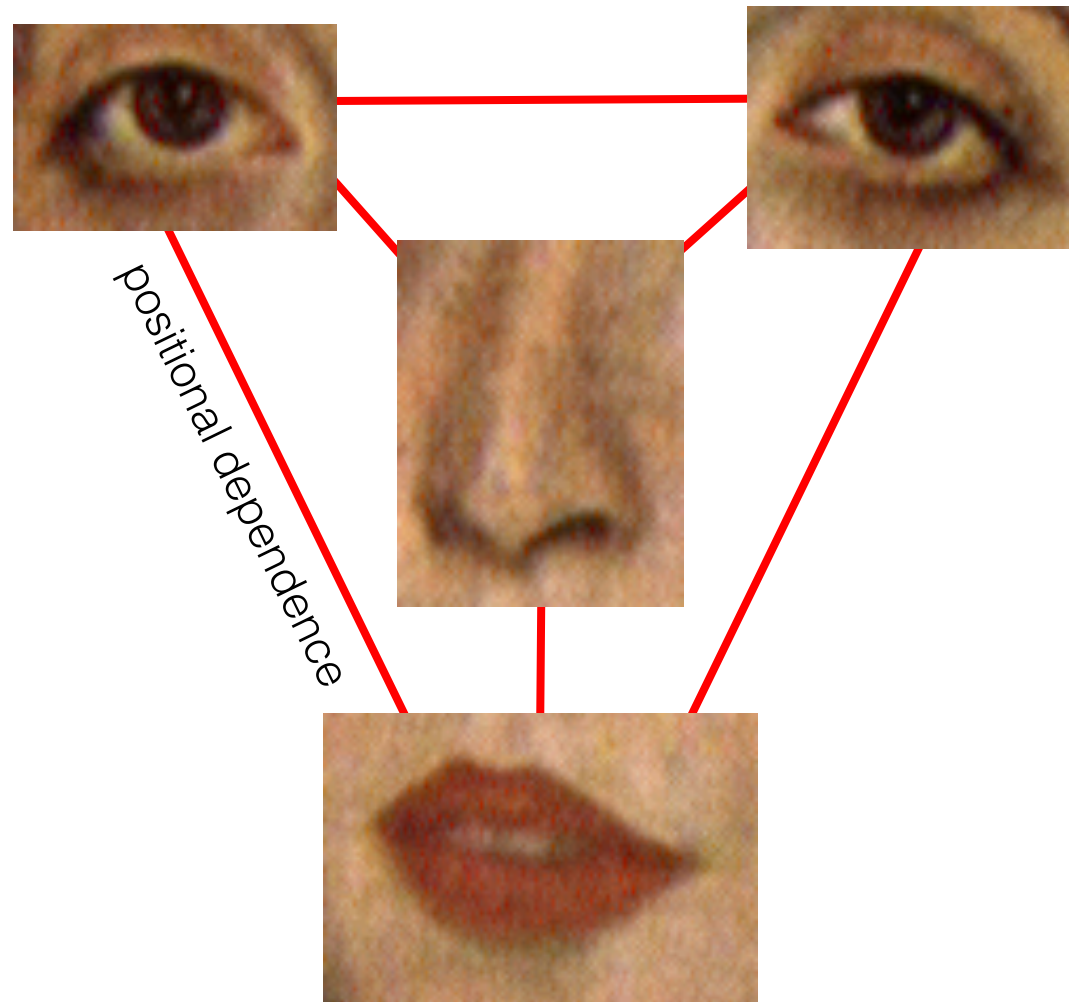
Spatial
reasoning



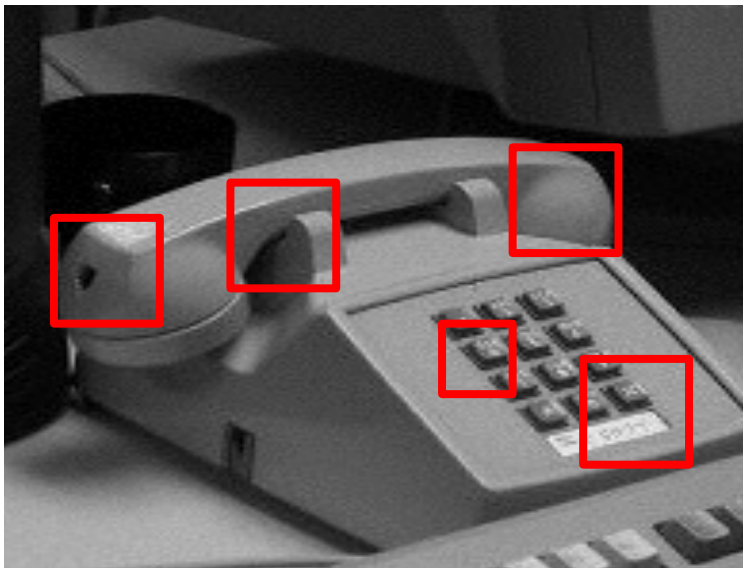
Window
classification

Spatial reasoning

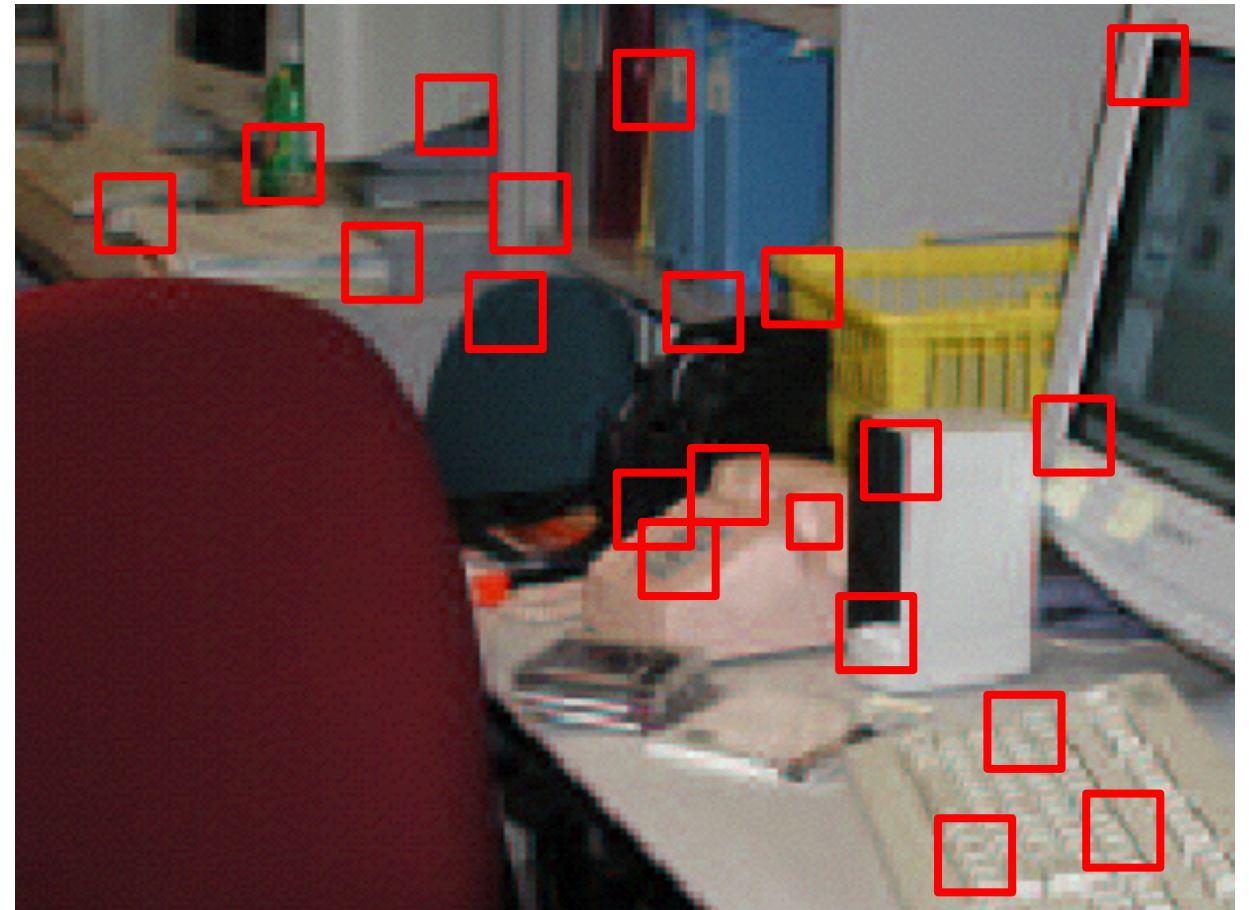
The position of every part depends on the positions of all the other parts



Many parts, many dependencies!

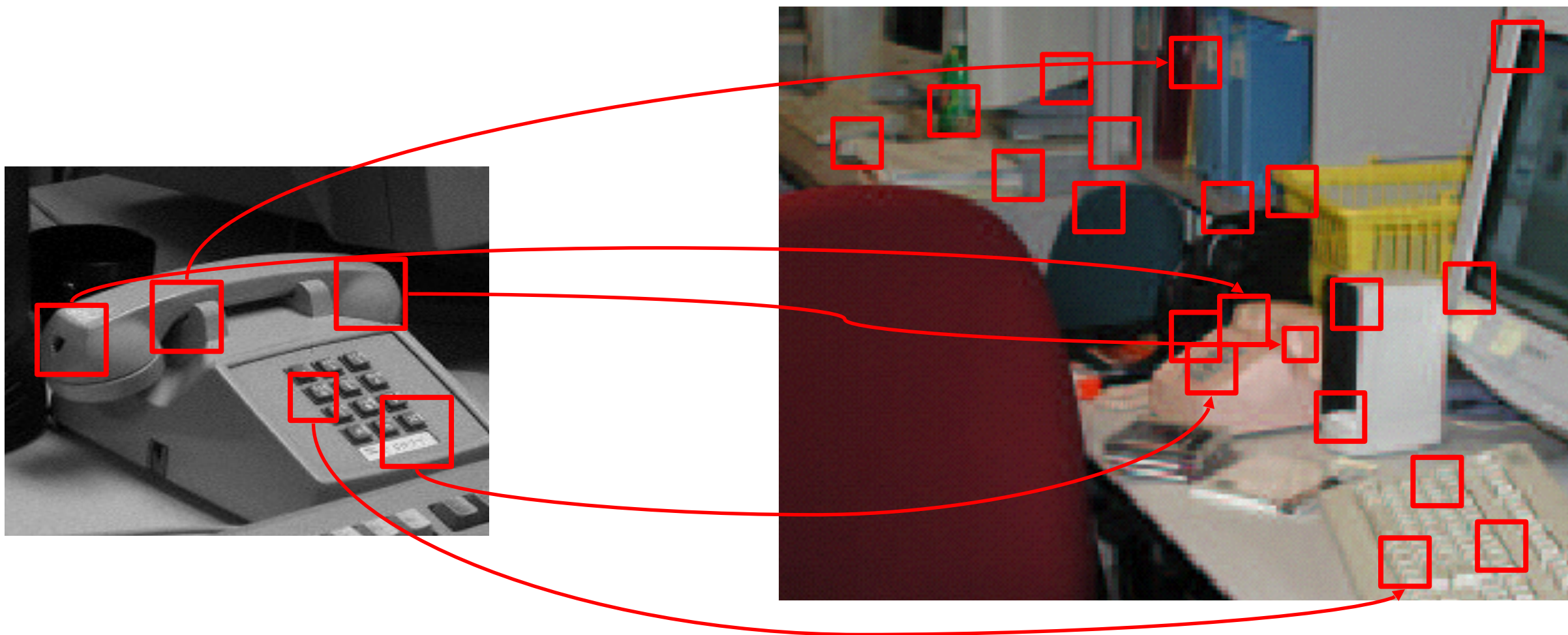


1. Extract features



2. Match features

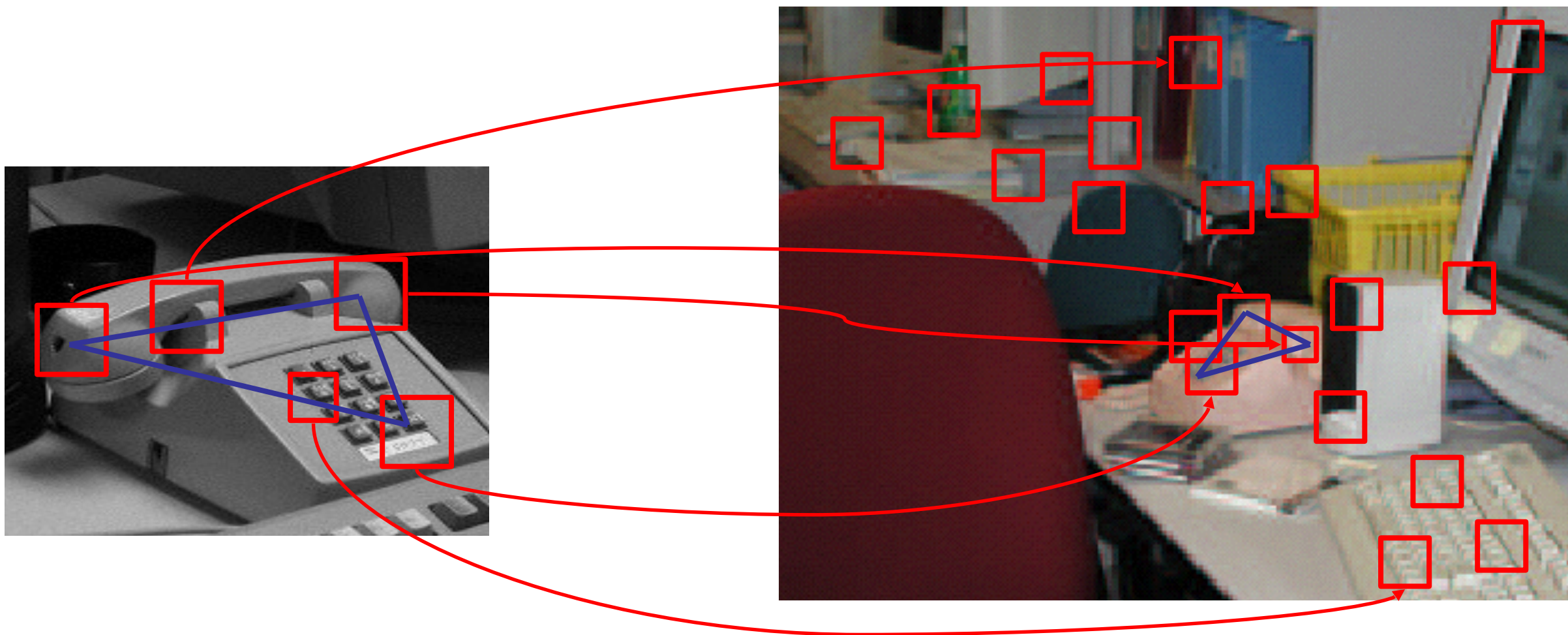
3. Spatial verification



1. Extract features

2. Match features

3. Spatial verification



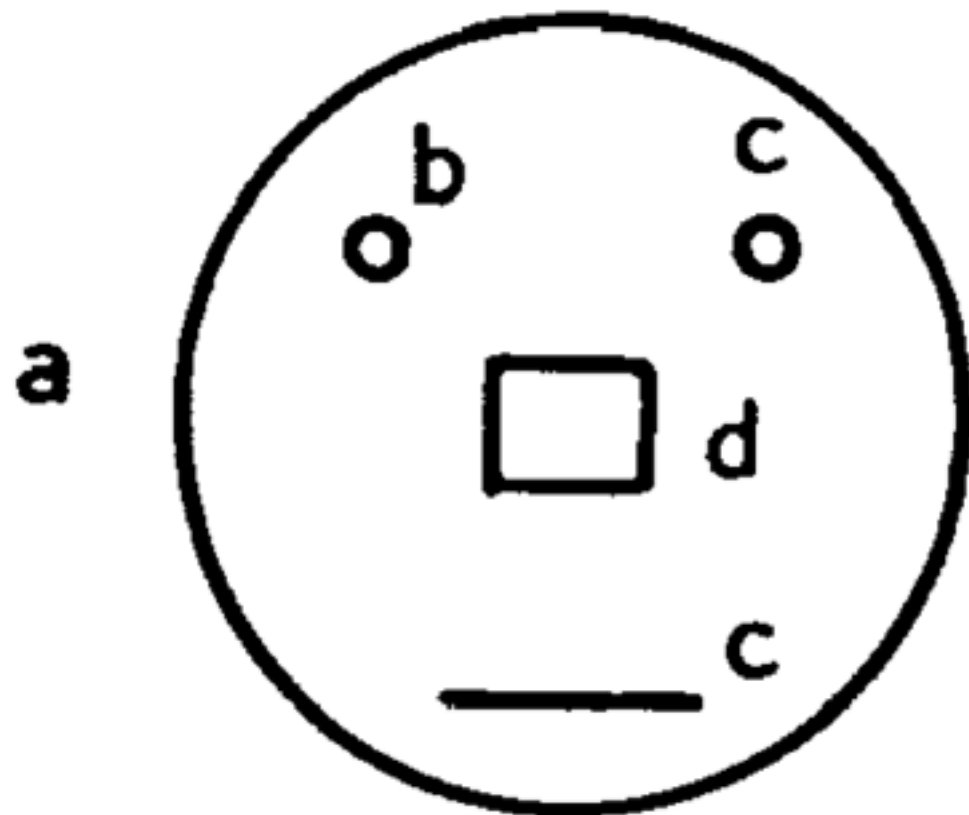
1. Extract features

2. Match features

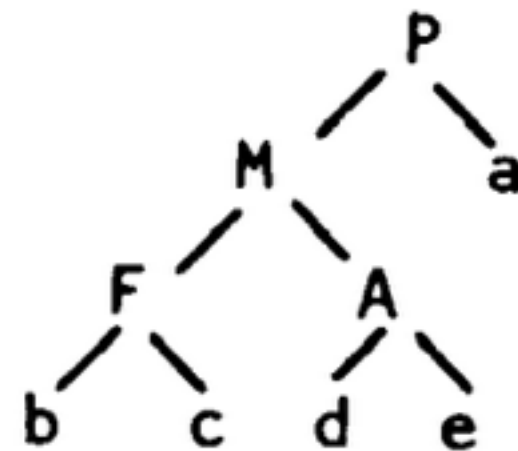
3. Spatial verification

an old idea...

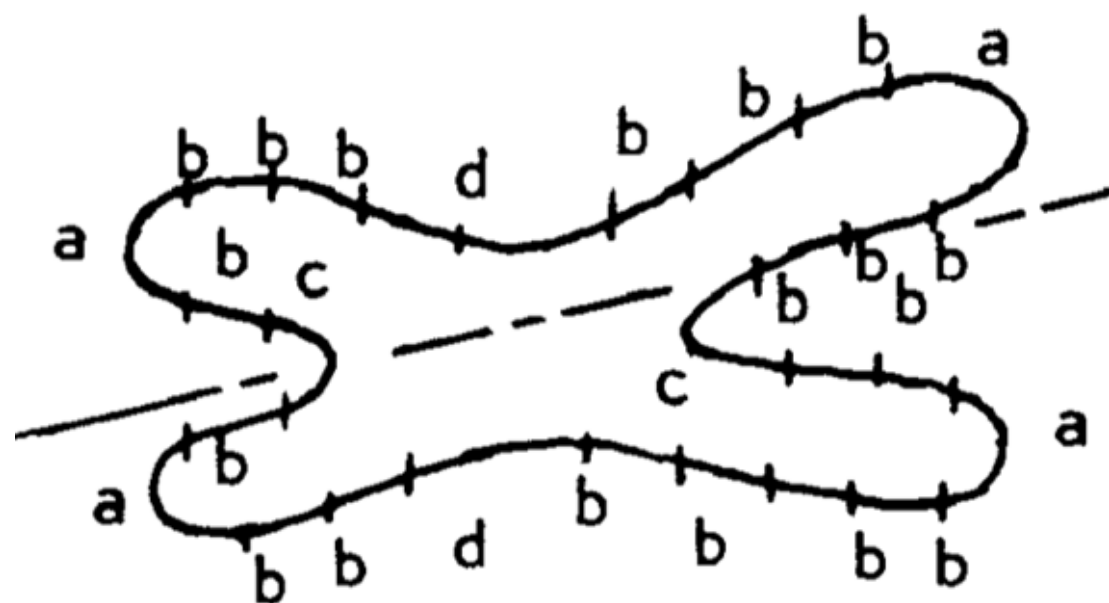
Fu and Booth. Grammatical Inference. 1975



Scene



Structural (grammatical) description

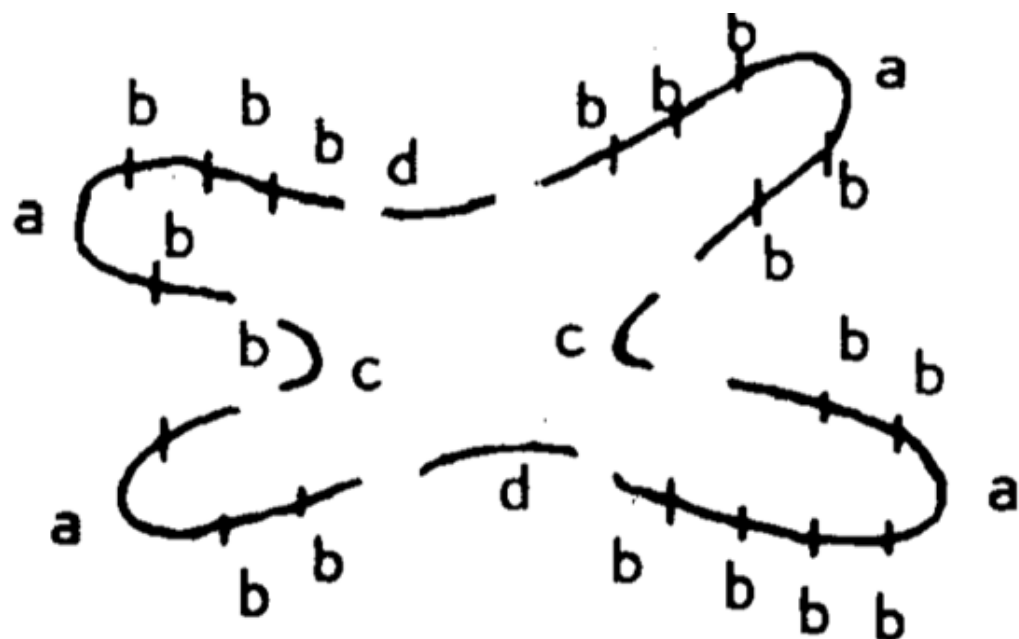


Coded Chromosome

$$V_T = \left\{ \begin{array}{c} \curvearrowright_a, \quad \nearrow_b, \quad \curvearrowright_c, \quad \curvearrowleft_d \end{array} \right\}$$

$$x = cdabbbdbbbabbcbabbabbdbbbabb$$

Substructures of Coded Chromosome



$$S_1 = \{ [b[[[a]b]b]b]; [b[b[b[a]]b]b]; \\ [b[b[[[a]b]b]b]b]; [b[b[a]]b] \}$$

The Representation and Matching of Pictorial Structures

MARTIN A. FISCHLER AND ROBERT A. ELSCHLAGER

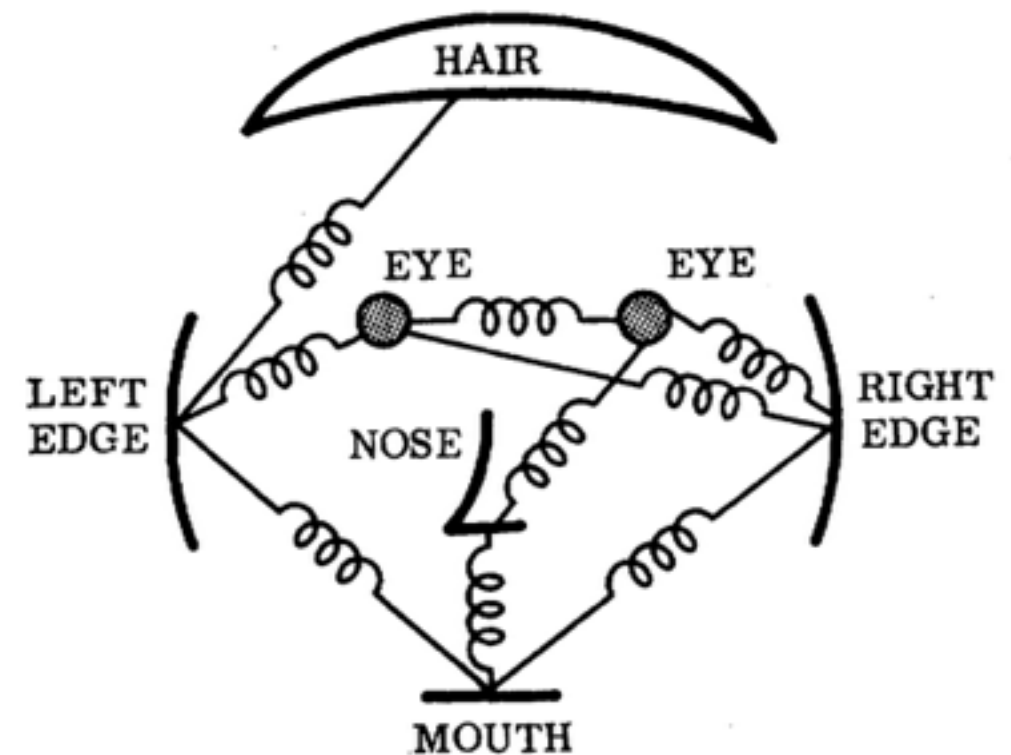
Abstract—The primary problem dealt with in this paper is the following. Given some description of a visual object, find that object in an actual photograph. Part of the solution to this problem is the specification of a descriptive scheme, and a metric on which to base the decision of “goodness” of matching or detection.

We offer a combined descriptive scheme and decision metric which is general, intuitively satisfying, and which has led to promising experimental results. We also present an algorithm which takes the above descriptions, together with a matrix representing the intensities of the actual photograph, and then finds the described object in the matrix. The algorithm uses a procedure similar to dynamic programming in order to cut down on the vast amount of computation otherwise necessary.

One desirable feature of the approach is its generality. A new programming system does not need to be written for every new description; instead, one just specifies descriptions in terms of a certain set of primitives and parameters.



1972



Description for left edge of face

A		E
B		F
C	X	G
D		H

$$\text{VALUE}(X) = (E + F + G + H) - (A + B + C + D)$$

Note: $\text{VALUE}(X)$ is the value assigned to the $L(EV)A$ corresponding to the location X as a function of the intensities of locations A through H in the sensed scene.

A more modern probabilistic approach...

think of locations as random variables (RV)

vector of RVs:
set of part locations

$$\mathbf{L} = \{ \overset{\text{RV}}{L_1}, \overset{\text{RV}}{L_2}, \dots, \overset{\text{RV}}{L_M} \}$$

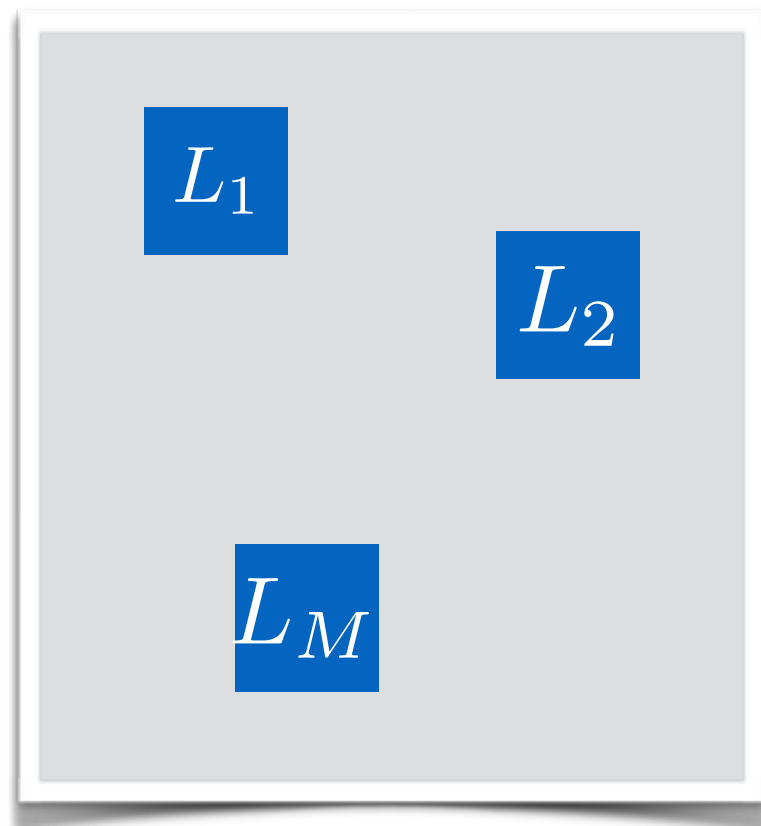
A more modern probabilistic approach...

think of locations as random variables (RV)

vector of RVs:
set of part locations

$$\mathbf{L} = \{ \overset{\text{RV}}{L_1}, \overset{\text{RV}}{L_2}, \dots, \overset{\text{RV}}{L_M} \}$$

image (N pixels)



What are the dimensions of R.V. L ?

How many possible combinations of part locations?

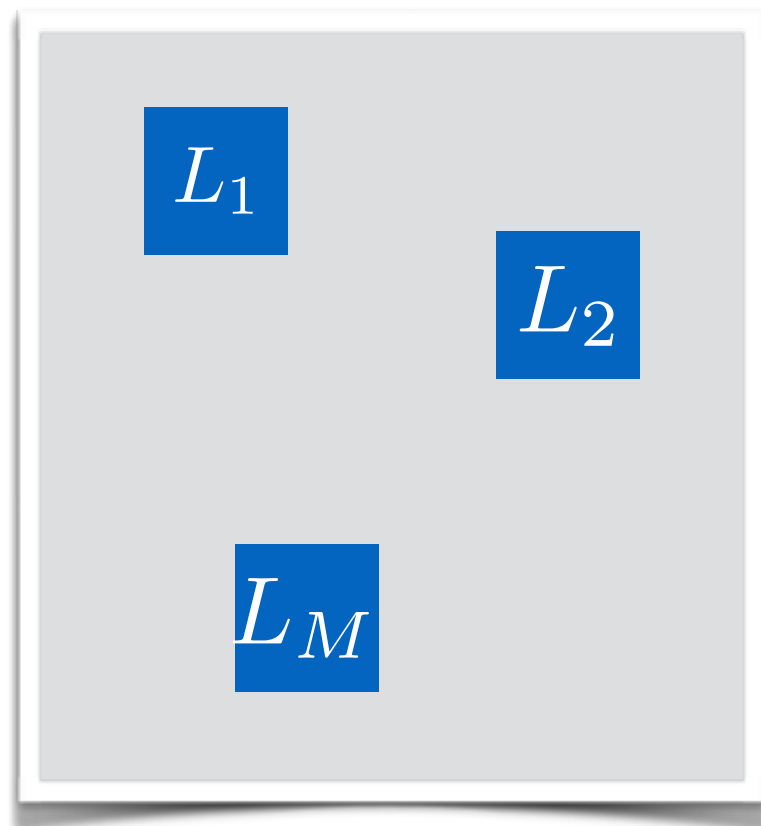
A more modern probabilistic approach...

think of locations as random variables (RV)

vector of RVs:
set of part locations

$$\mathbf{L} = \{ \overset{\text{RV}}{L_1}, \overset{\text{RV}}{L_2}, \dots, \overset{\text{RV}}{L_M} \}$$

image



What are the dimensions of R.V. L ?

$$L_m = [x \ y]$$

How many possible combinations of part locations?

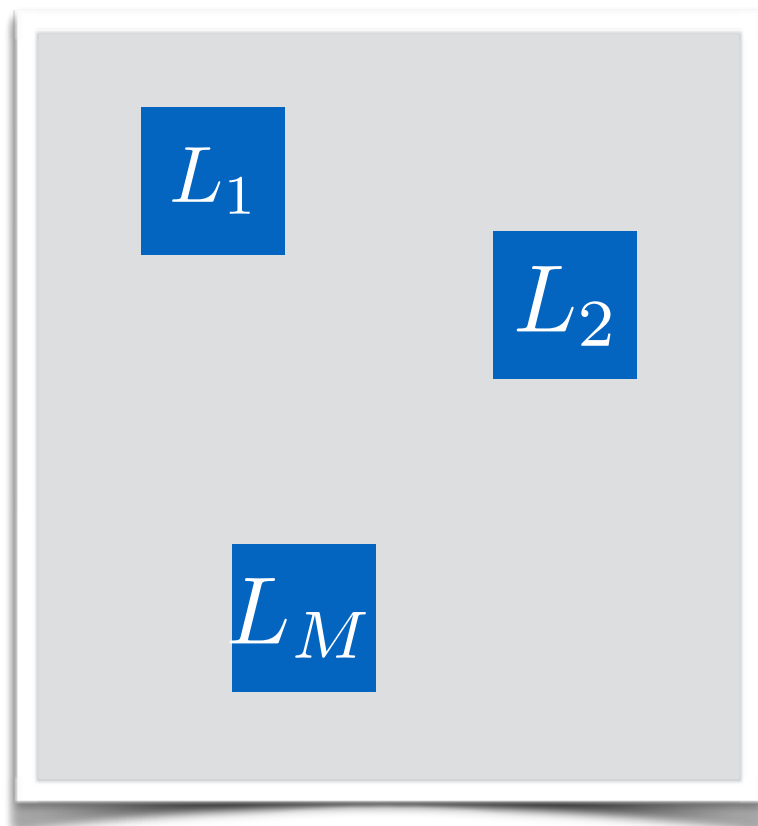
A more modern probabilistic approach...

think of locations as random variables (RV)

vector of RVs:
set of part locations

$$\mathbf{L} = \{ \overset{\text{RV}}{L_1}, \overset{\text{RV}}{L_2}, \dots, \overset{\text{RV}}{L_M} \}$$

image



What are the dimensions of R.V. L ?

$$L_m = [x \ y]$$

How many possible combinations of part locations?

$$N^M$$

Most likely set of locations \mathbf{L} is found by maximizing:

$$p(\overset{\text{part}}{\underset{\text{locations}}{\mathbf{L}}} | \overset{\text{image}}{\mathbf{I}}) \propto p(\mathbf{I} | \mathbf{L}) p(\mathbf{L})$$

Posterior

Likelihood:
How likely it is to observe
image \mathbf{I} given that the M parts
are at locations \mathbf{L}
(scaled output of a classifier)

Prior:
spatial prior controls the
geometric configuration of the
parts

What kind of prior can we formulate?

Given any collection of selfie images,
where would you expect the nose to be?



*What would be an appropriate **prior**?*

$$P(L_{\text{nose}}) = ?$$

A simple factorized model

$$p(\mathbf{L}) = \prod_m p(L_m)$$

Break up the joint probability into
smaller (independent) terms

Independent locations

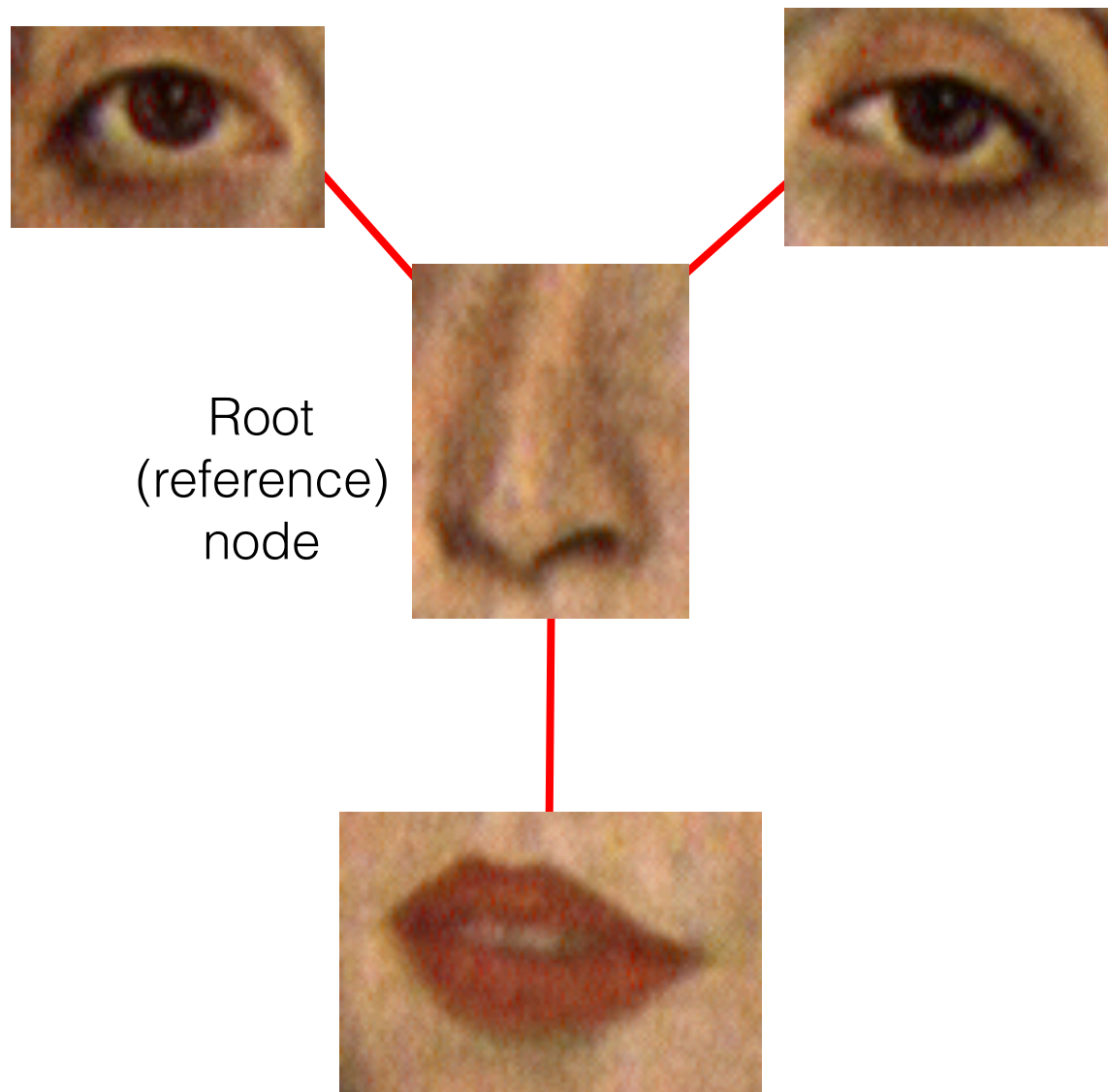


$$p(\mathbf{L}) = \prod_m p(L_m)$$

Each feature is allowed to move independently

Does not model the **relative** location of parts at all

Tree structure (star model)

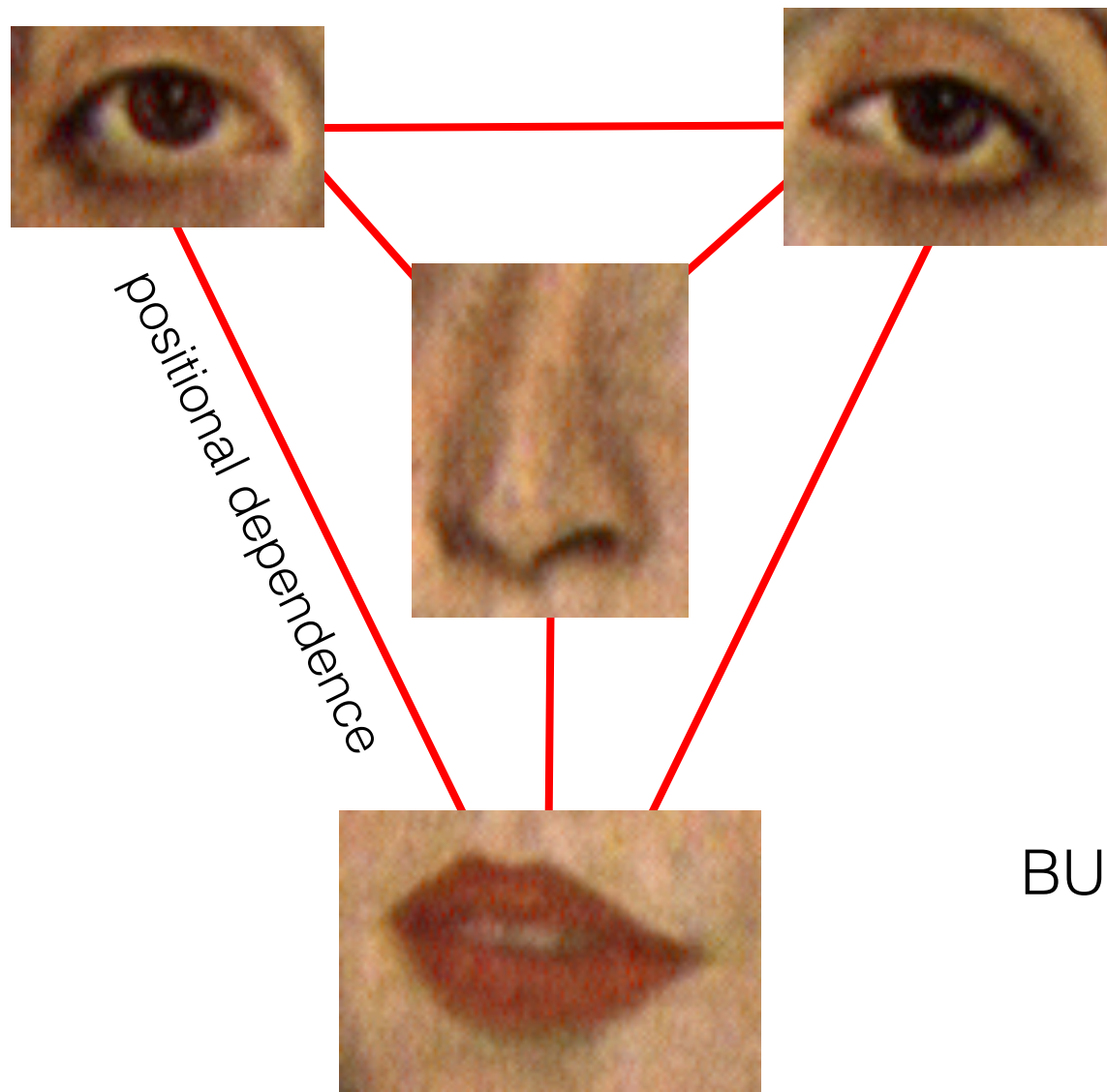


$$p(\mathbf{L}) = p(L_{\text{root}}) \prod_{m=1}^{M-1} p(L_m | L_{\text{root}})$$

Represent the location of
all the parts relative to a single
reference part

Assumes that one
reference part is defined
(who will decide this?)

Fully connected (constellation model)



$$p(L) = p(l_1, \dots, l_N)$$

Explicitly represents the joint distribution of locations

Good model:

Models relative location of parts
BUT Intractable for moderate number of parts

Pros

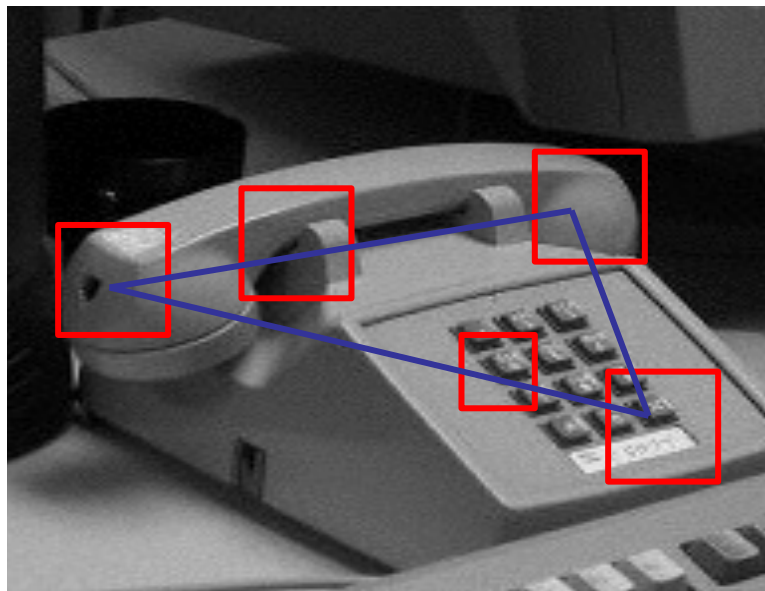
- Retains spatial constraints
- Robust to deformations

Cons

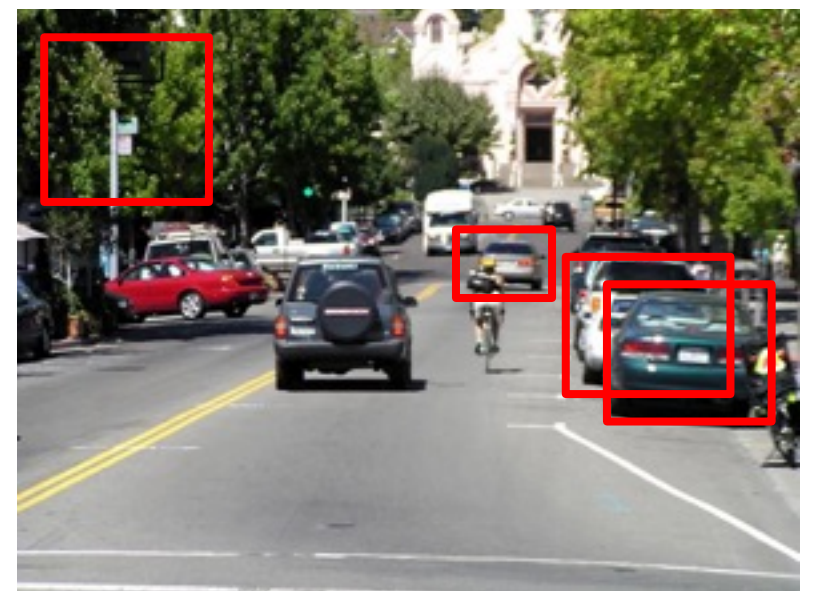
- Computationally expensive
- Generalization to **large** inter-class variation (e.g., modeling chairs)



Feature
Matching



Spatial
reasoning



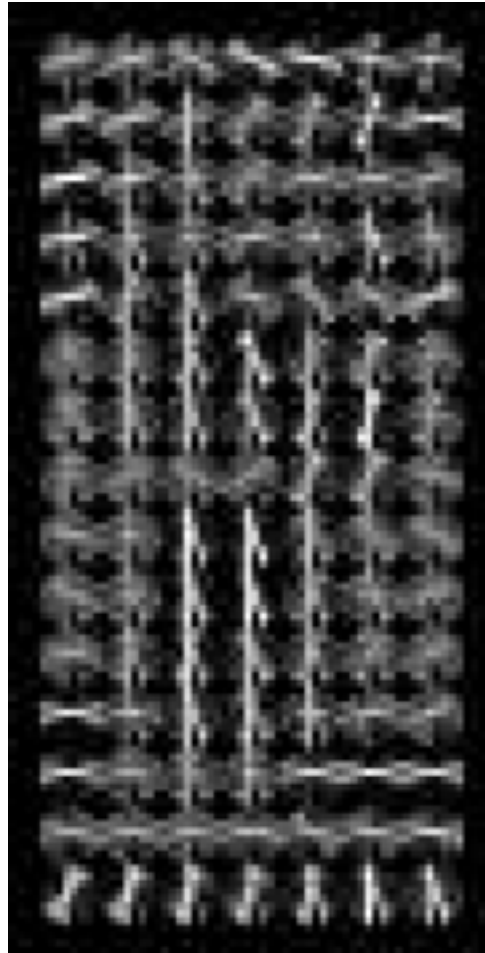
Window
classification

Window-based

Template Matching



1. get image window



2. extract features



3. classify

When does this work and when does it fail?

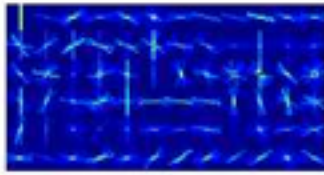
How many templates do you need?

Per-exemplar

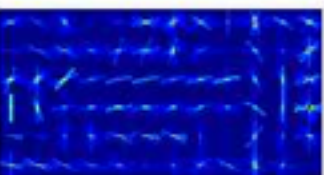
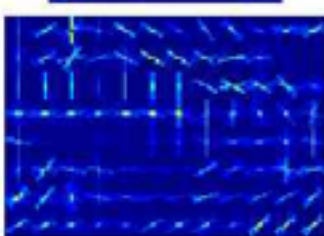
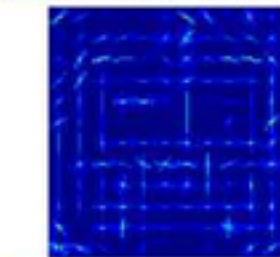
exemplar



template



top hits from test data

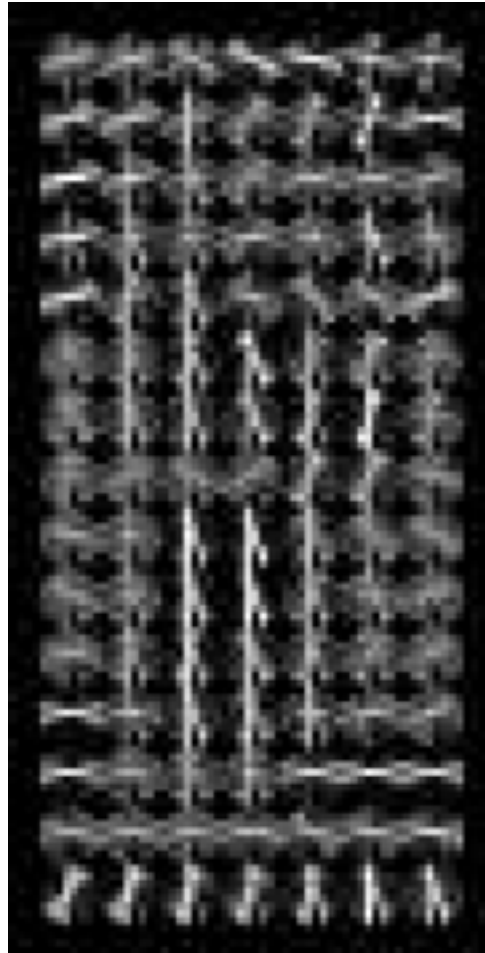


find the 'nearest' exemplar, inherit its label

Template Matching



1. get image window
(or region proposals)



2. extract features

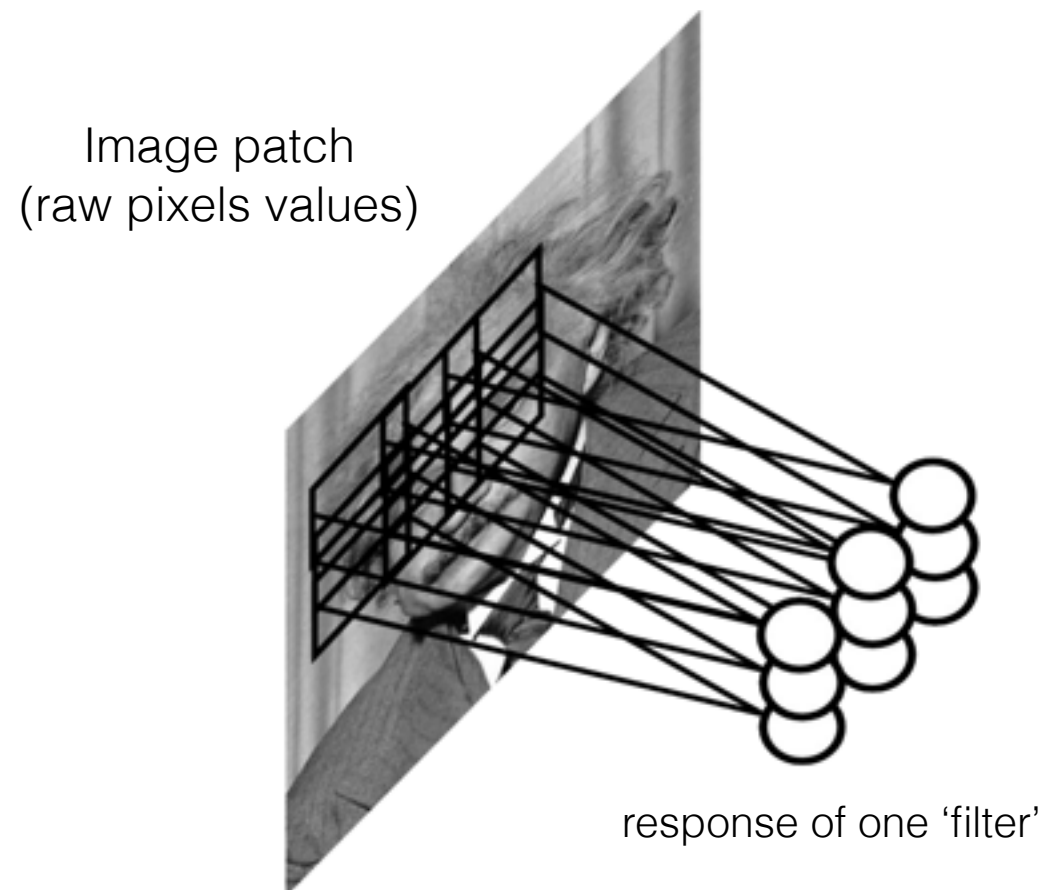


3. compare to template

Do this part with one big classifier
'end to end learning'

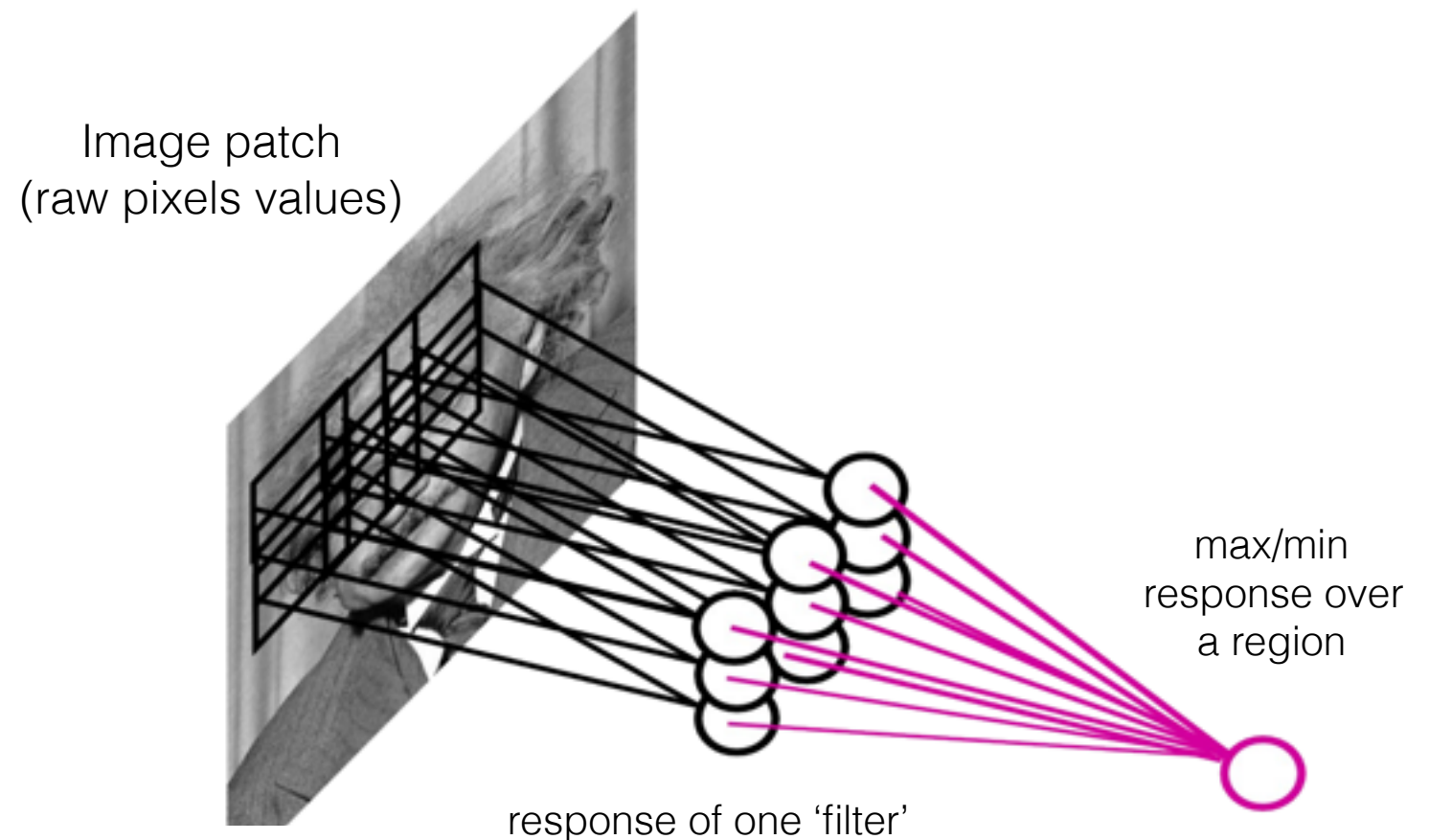
Convolutional Neural Networks

Convolution

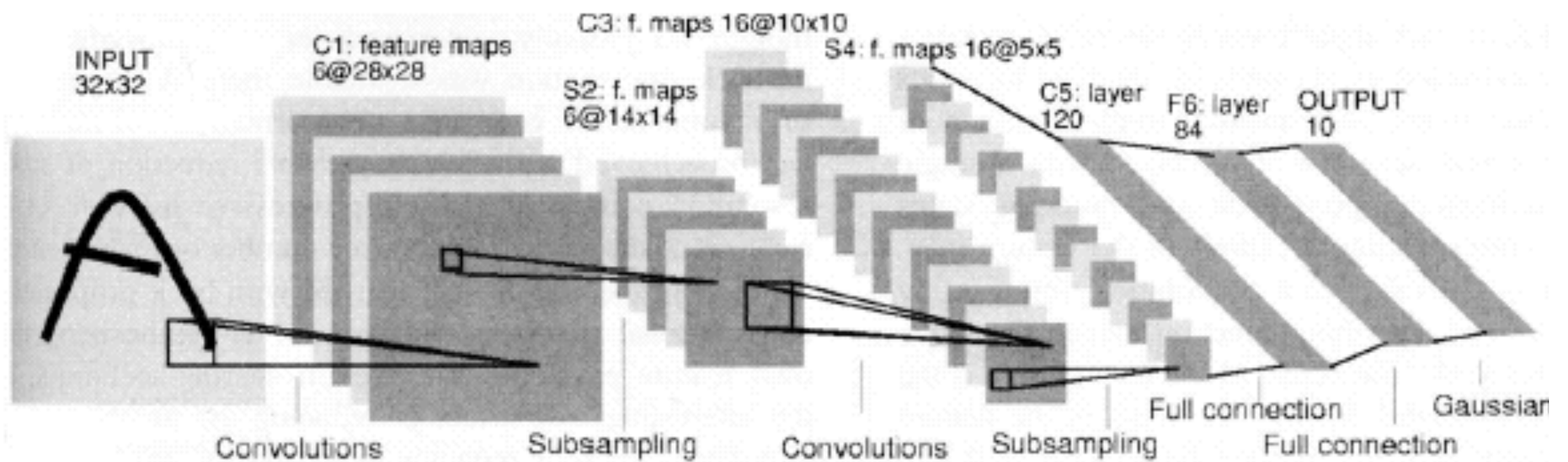
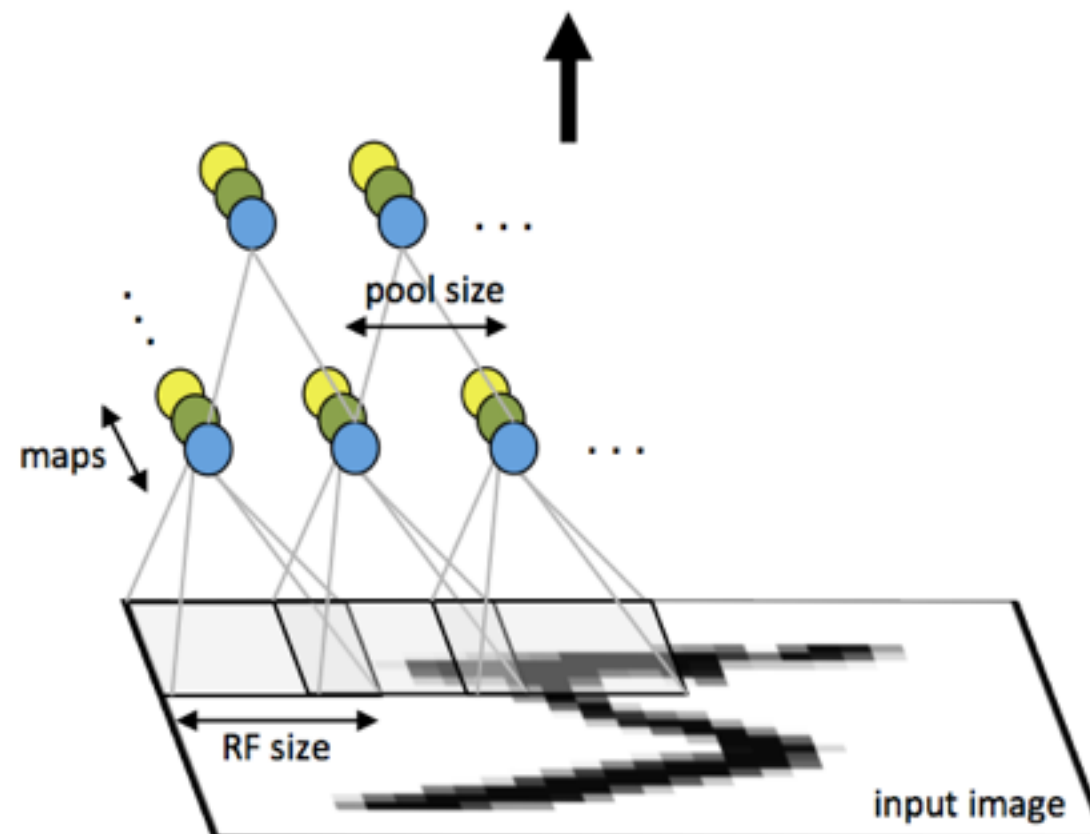


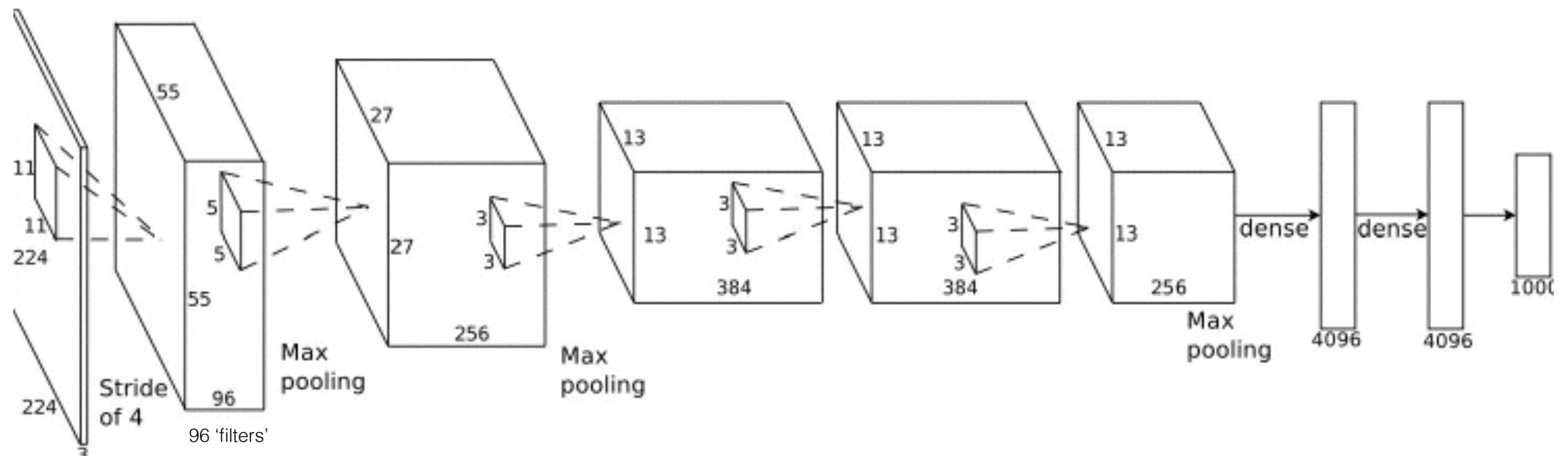
A 96 x 96 image convolved with 400 filters (features) of size 8 x 8 generates about 3 million values ($89^2 \times 400$)

Pooling



Pooling aggregates statistics and lowers the dimension of convolution





630 million connections
60 millions parameters to learn

Krizhevsky, A., Sutskever, I. and Hinton, G. E.
ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012.

Pros

- Retains spatial constraints
- Efficient test time performance

Cons

- Many many possible windows to evaluate
- Requires large amounts of data
- Sometimes (very) slow to train

How to write an effective CV resume

1. **System and Method for Deep Learning.** Deep Learning, Deep Learning , Deep Learning , Deep Learning