



S A I R

Spatial AI & Robotics Lab

CSE 473/573-A

L21: DETECTION

Chen Wang

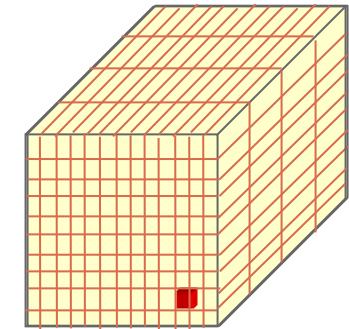
Spatial AI & Robotics Lab

Department of Computer Science and Engineering



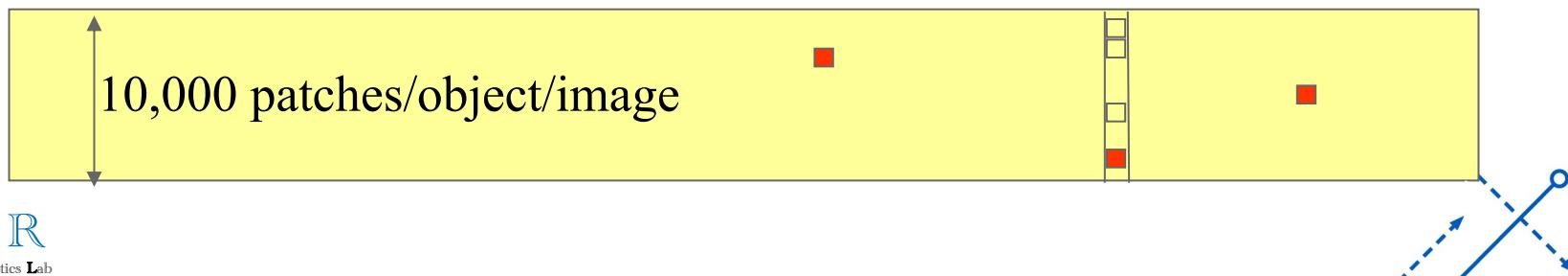
University at Buffalo The State University of New York

Why is detection hard?



We want to do this
for ~ 1000 objects

1,000,000 images/day



Why is detection hard?

If we have 1000 categories (detectors), and each detector produces 1 false every 10 images, we will have 100 false alarms per image... pretty much garbage...



Local information is helpful



Slide credit: A. Torralba

Is local information enough?



Information

Local features

Contextual features

Distance

Slide credit: A. Torralba

Is local information even enough?

We know there is a keyboard present in this scene even if we cannot see it clearly.



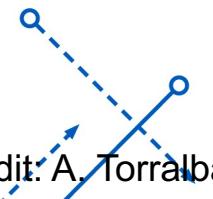
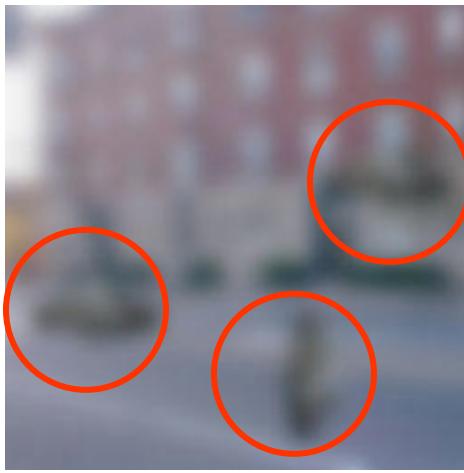
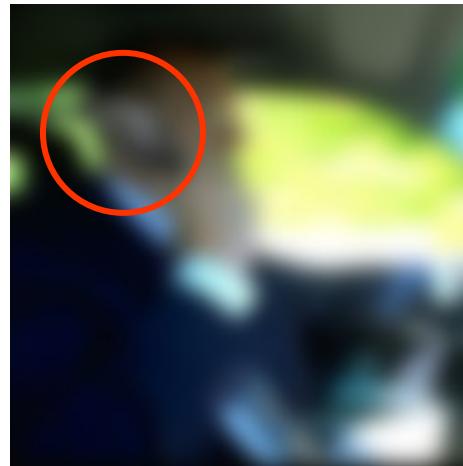
We know there is no keyboard present in this scene



... even if there is one indeed.

Slide credit: A. Torralba

The multiple personalities of a blob



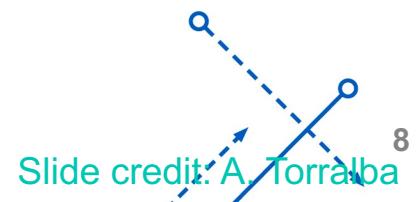
Slide credit: A. Torralba

The multiple personalities of a blob

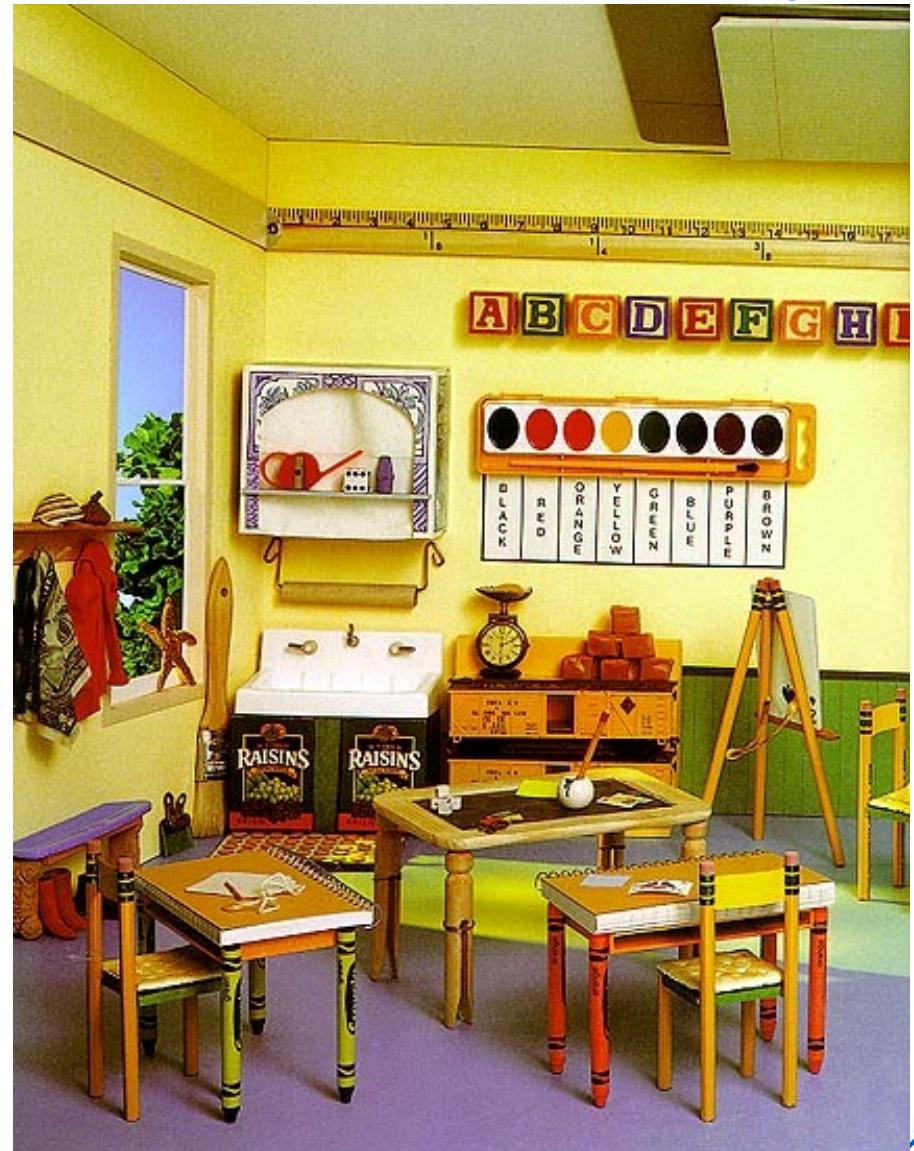
A B C

12
B
14

12
A B C
14



The multiple personalities of a blob



The context challenge

- What are the hidden objects?



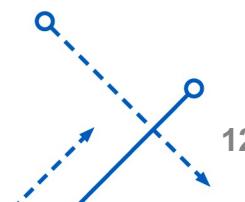
The context challenge

- What are the hidden objects?



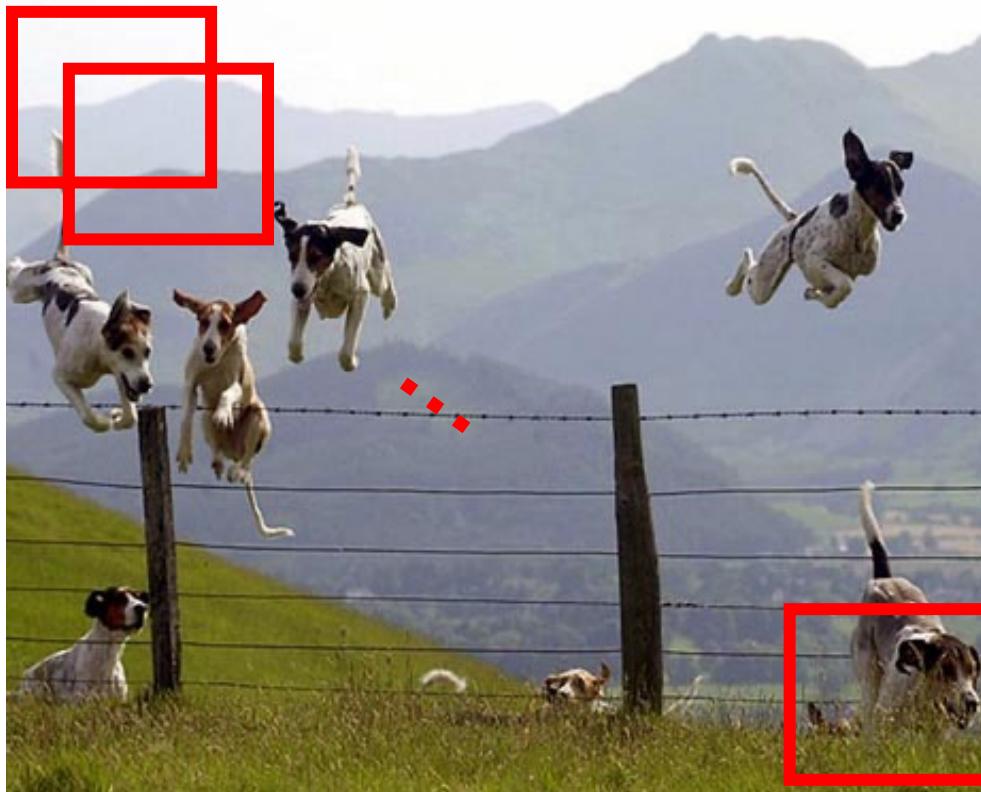
Object detection vs Scene Recognition

- Objects (even if deformable and articulated) probably have more **consistent shapes** than scenes.
- Scenes can be defined by distribution of “stuff” – materials and surfaces with **arbitrary shape**.
- Objects are “**things**” that own their **boundaries**
- Bag of words (BoW) were less popular for object detection because they **throw away shape** info.

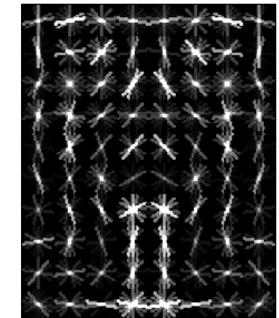


Object Category Detection

- Focus on object search: “Where is it?”
- Build templates that quickly differentiate object patch from background patch



Dog Model



Object or
Non-Object?

Challenges in modeling the object class



Illumination



Object pose



Clutter



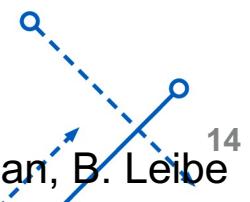
Occlusions



Intra-class
appearance



Viewpoint



Challenges in modeling the non-object class



True
Detections



Bad
Localization



Confused with
Similar Object



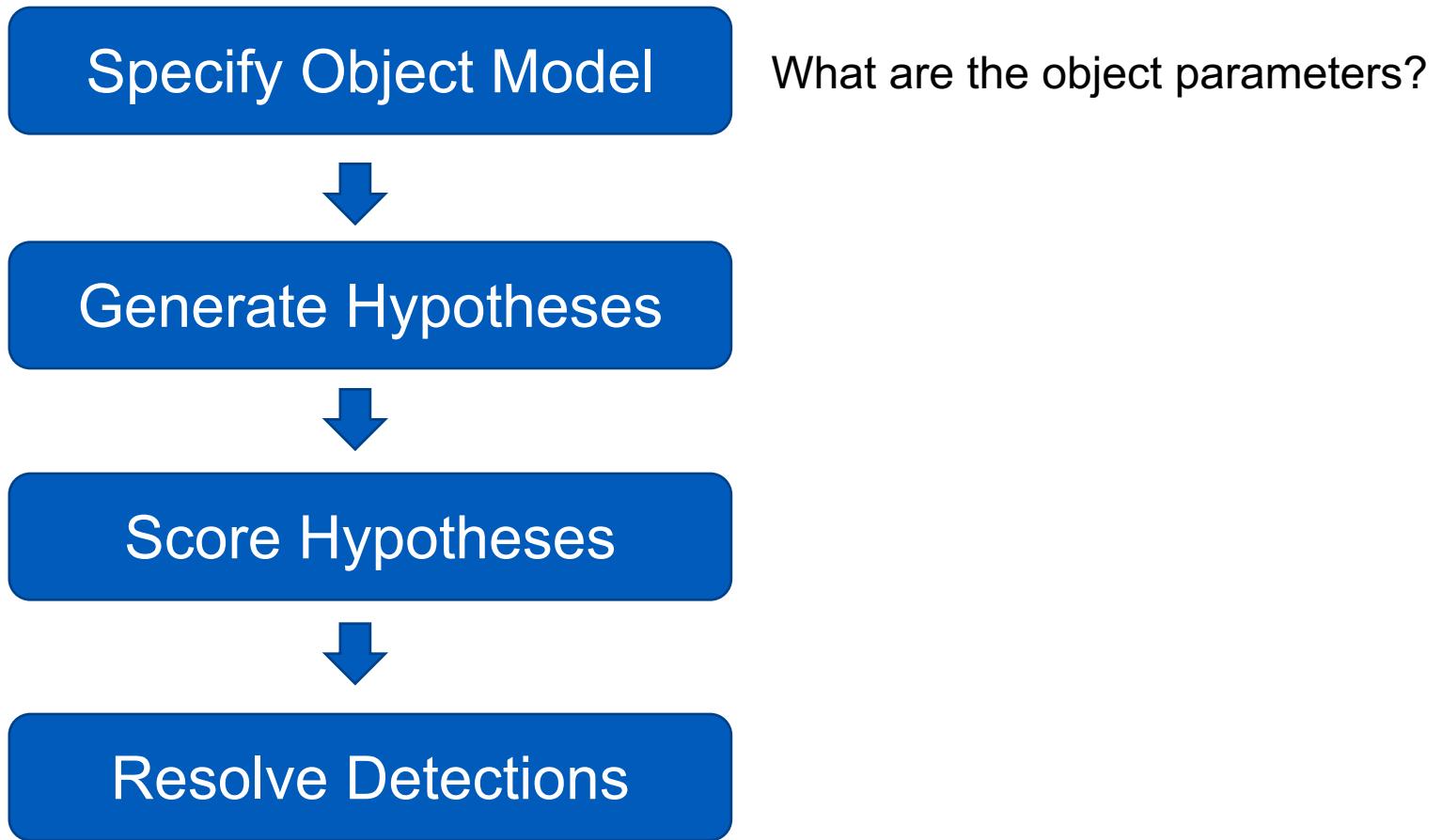
Misc. Background



Confused with
Dissimilar Objects



General Process of Object Recognition



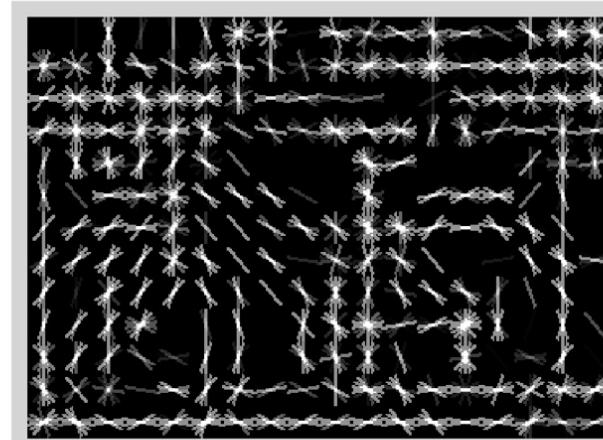
Specifying an object model

1. Statistical Template in Bounding Box

- Object is somewhere (x, y, w, h) in image
- Features defined wrt bounding box coordinates



Image

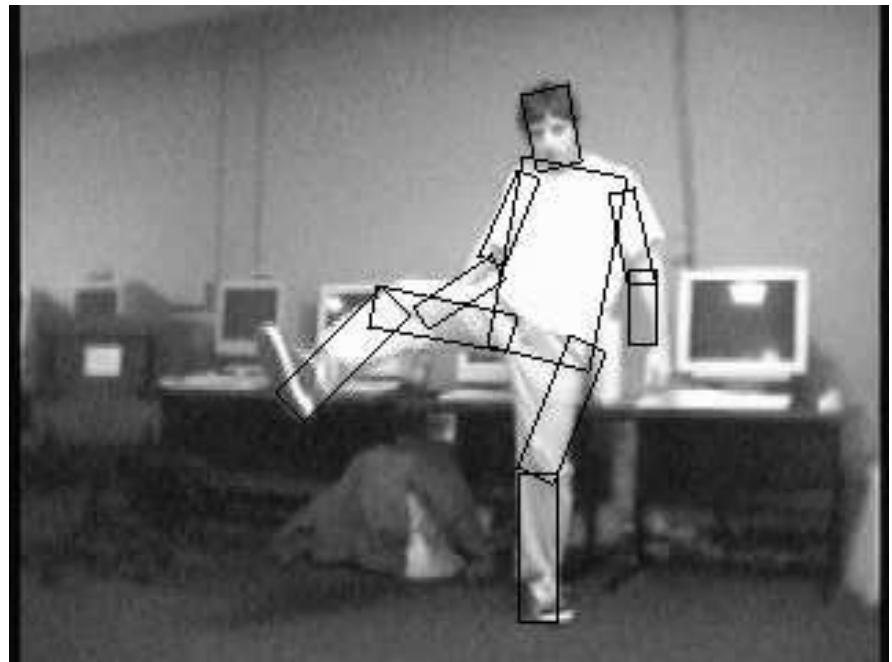
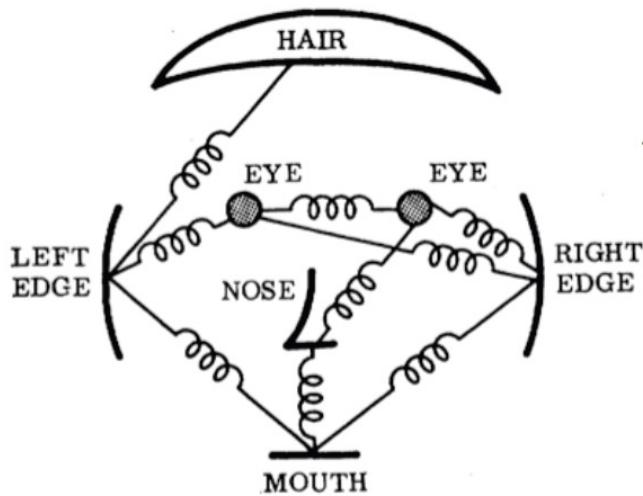


Template Visualization
(HoG features)

Specifying an object model

2. Articulated parts model

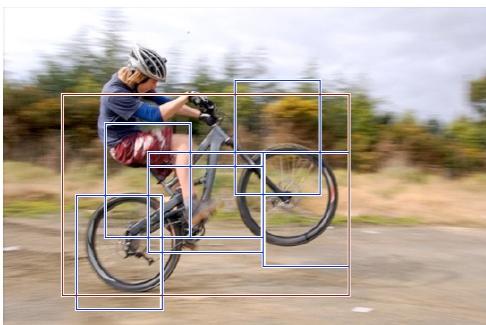
- Object is configuration of parts
- Each part is detectable



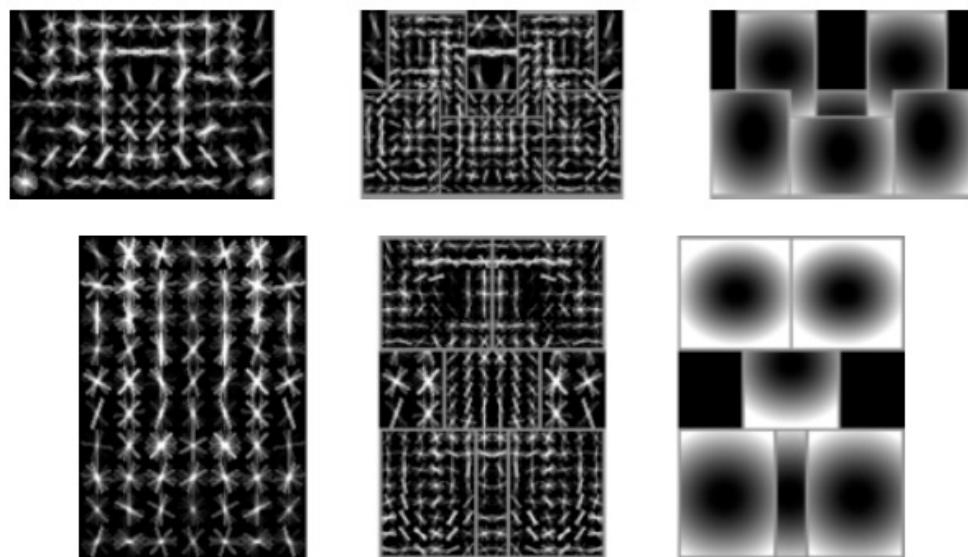
Specifying an object model

3. Hybrid template/parts model

Detections



Template Visualization
HoG features



root filters
coarse resolution

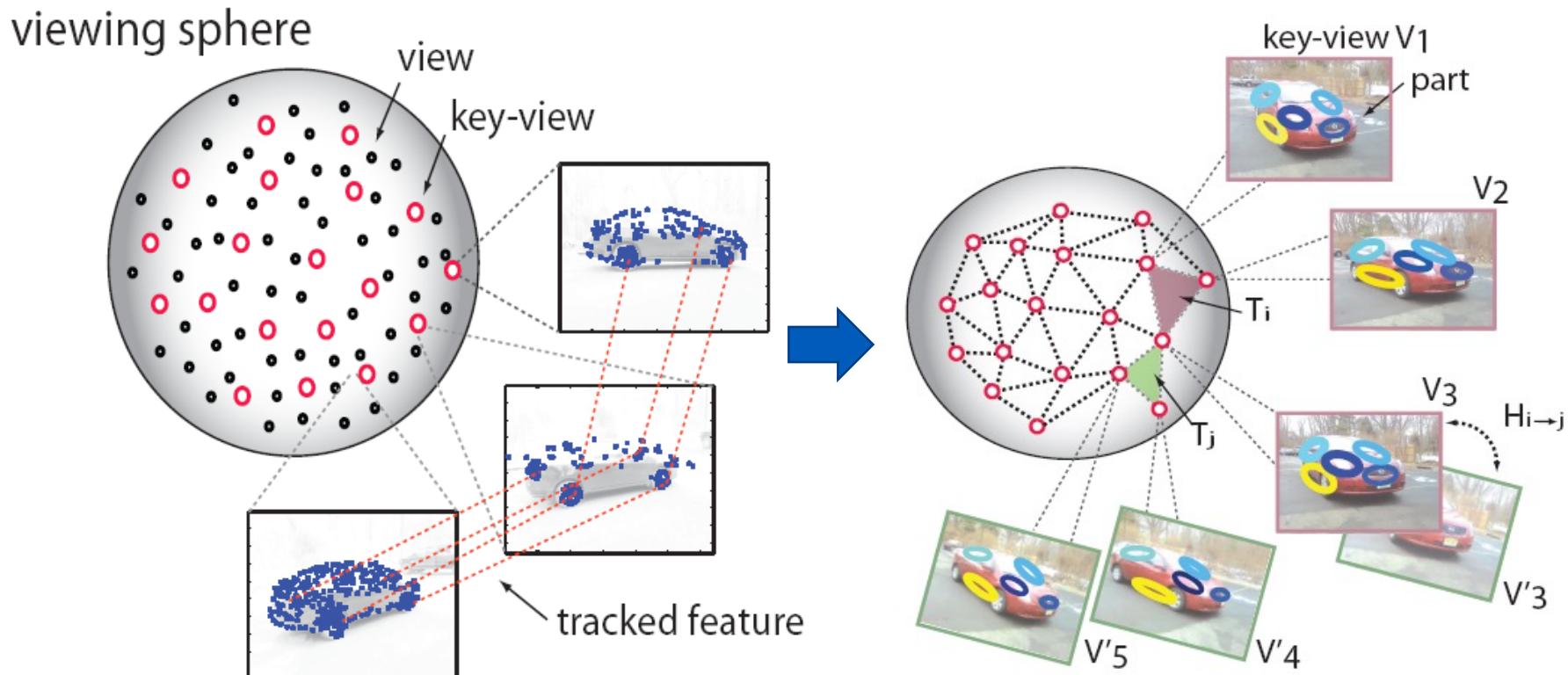
part filters
finer resolution

deformation
models

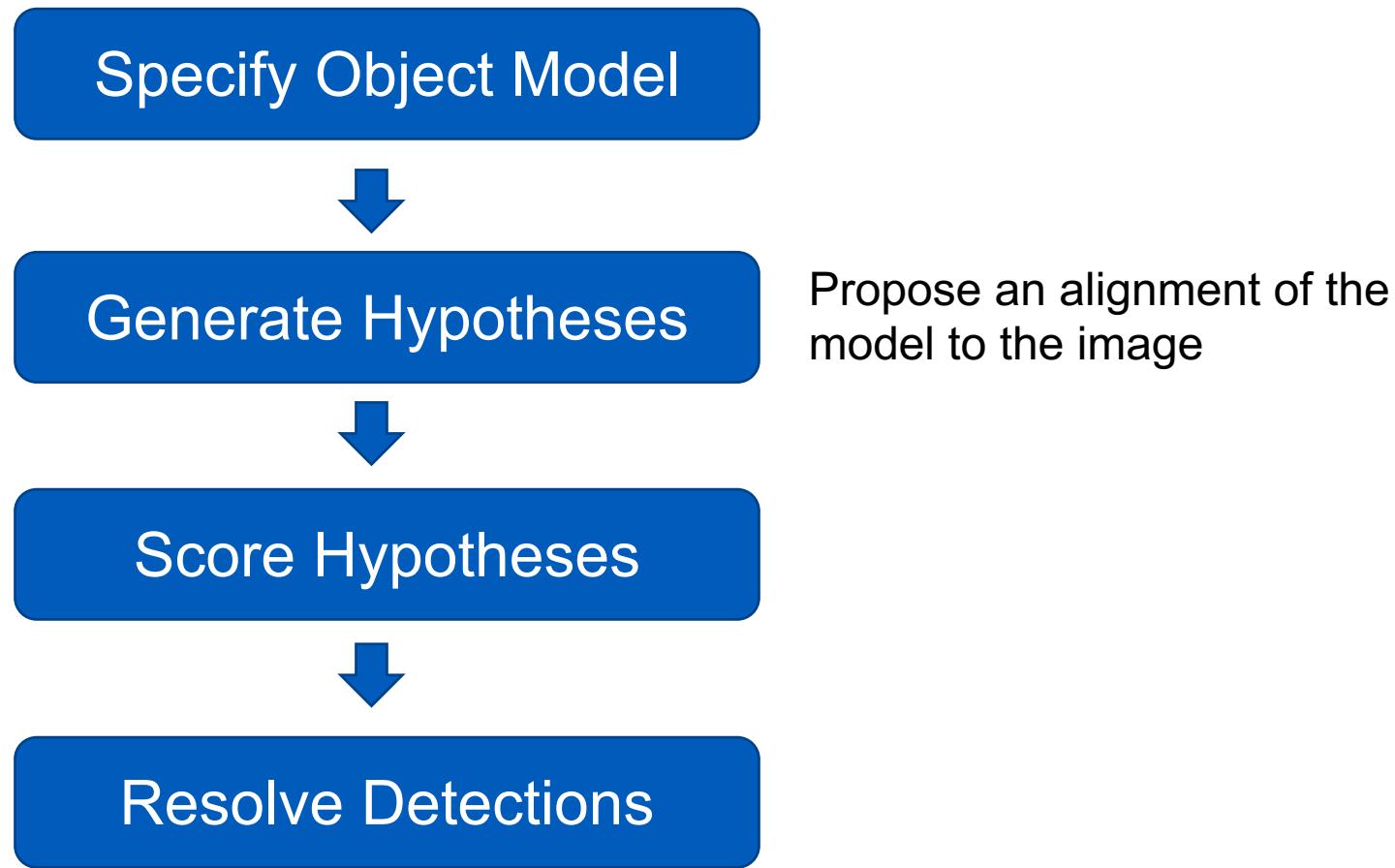
Specifying an object model

4. 3D-ish model

- Object is collection of 3D planar patches under affine transformation



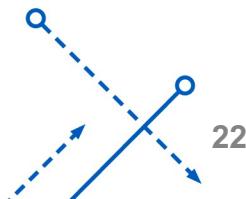
General Process of Object Recognition



Generating hypotheses

1. Sliding window

- Test patch at each **location** and **scale**



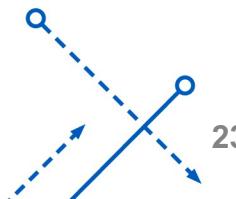
Generating hypotheses

1. Sliding window

- Test patch at each **location** and **scale** (**pyramids**)



Note – Template did not change size



Each window is separately classified



Generating hypotheses

2. Voting from patches/keypoints

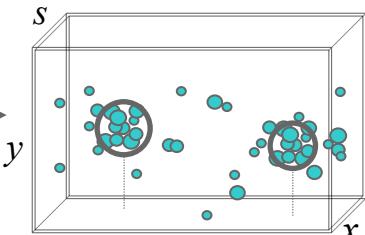
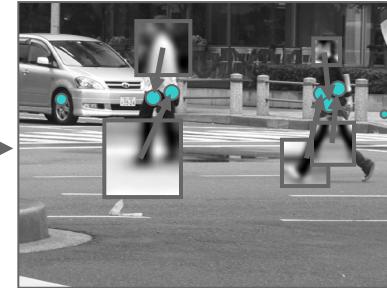
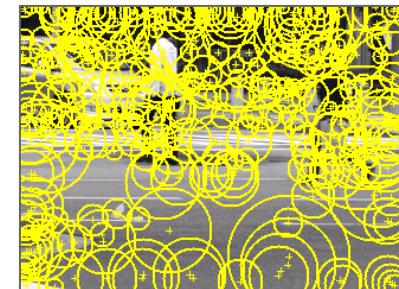


Interest Points

Matched Codebook
Entries

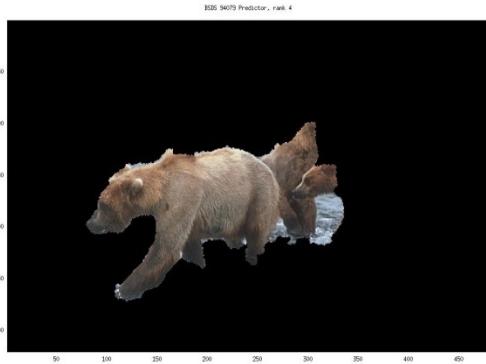
Probabilistic
Voting

3D Voting Space
(continuous)



Generating hypotheses

3. Region-based proposal



General Process of Object Recognition

Specify Object Model



Generate Hypotheses



Score Hypotheses



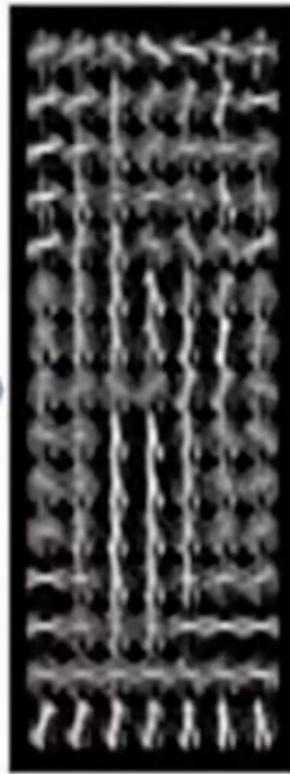
Resolve Detections

Mainly-gradient based features,
usually based on summary
representation, many classifiers.

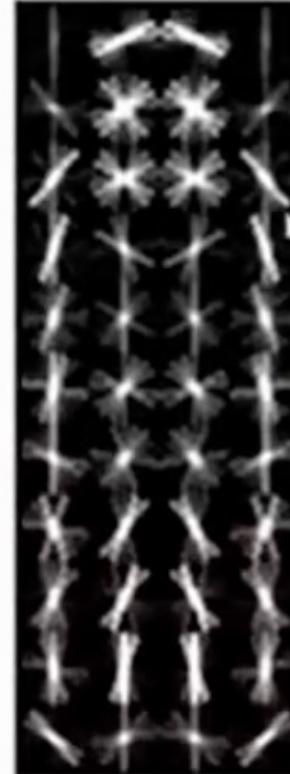
HOG + SVM for Object Detection



Image



HOG Features



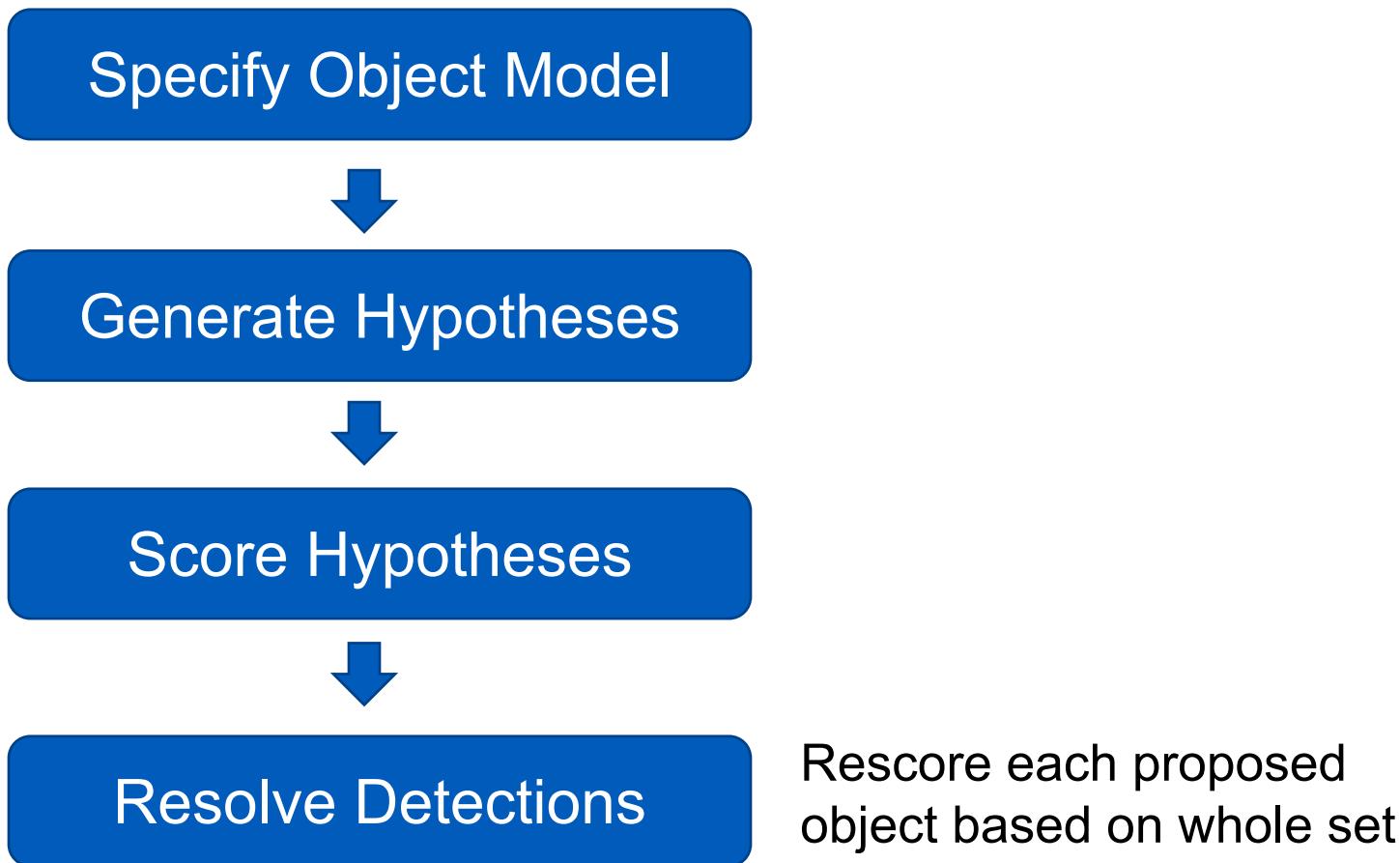
SVM Coefs
Human Model



Human or Not

<https://youtu.be/sDByl84n5mY>

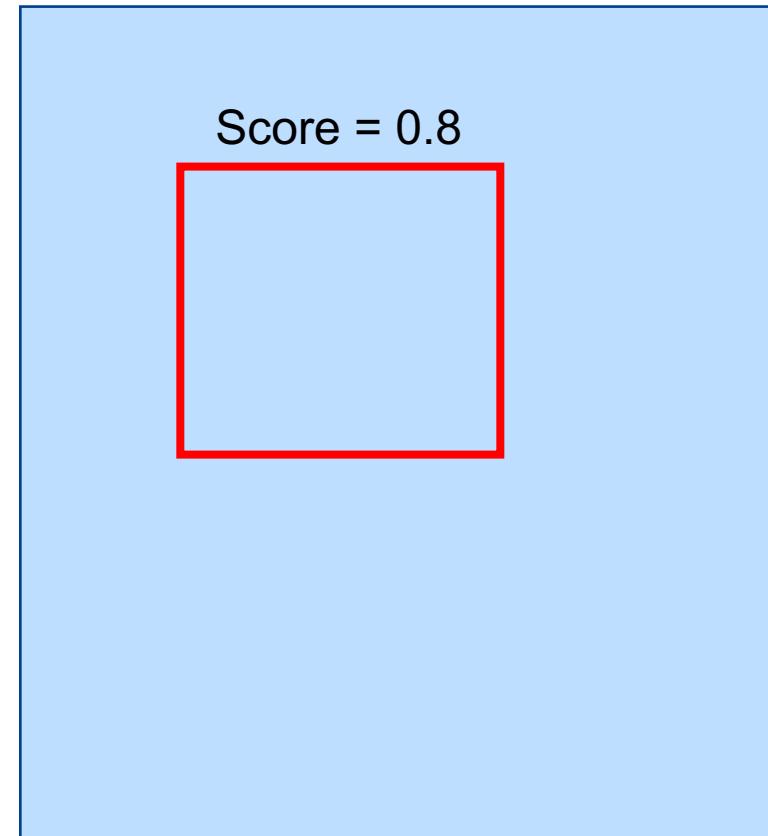
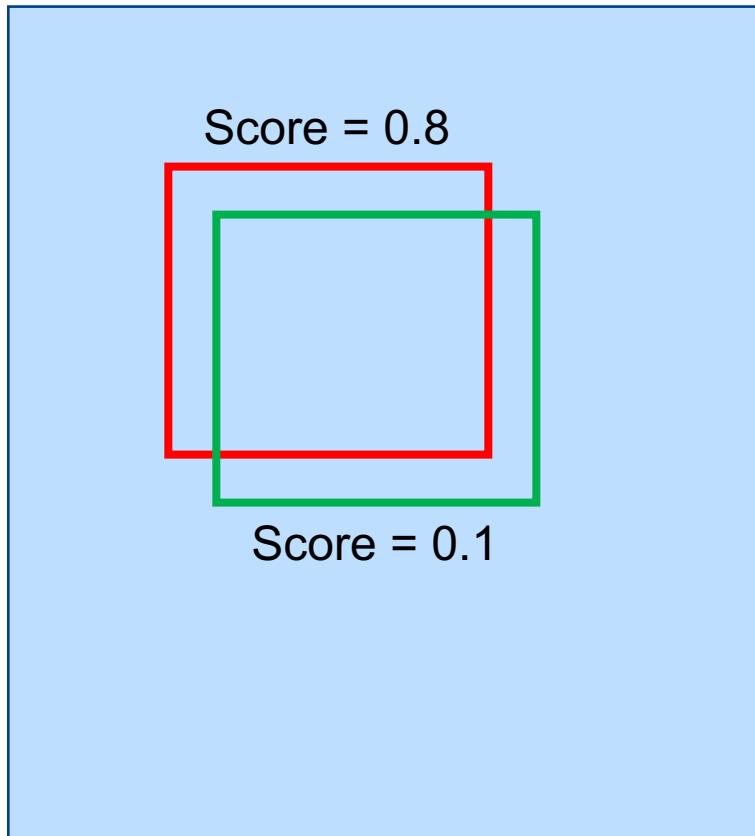
General Process of Object Recognition



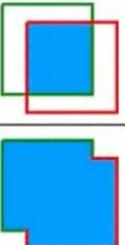
Resolving detection scores

$$IOU = \frac{\text{area of overlap}}{\text{area of union}} = \frac{\text{blue} \cap \text{red}}{\text{blue} \cup \text{red}}$$

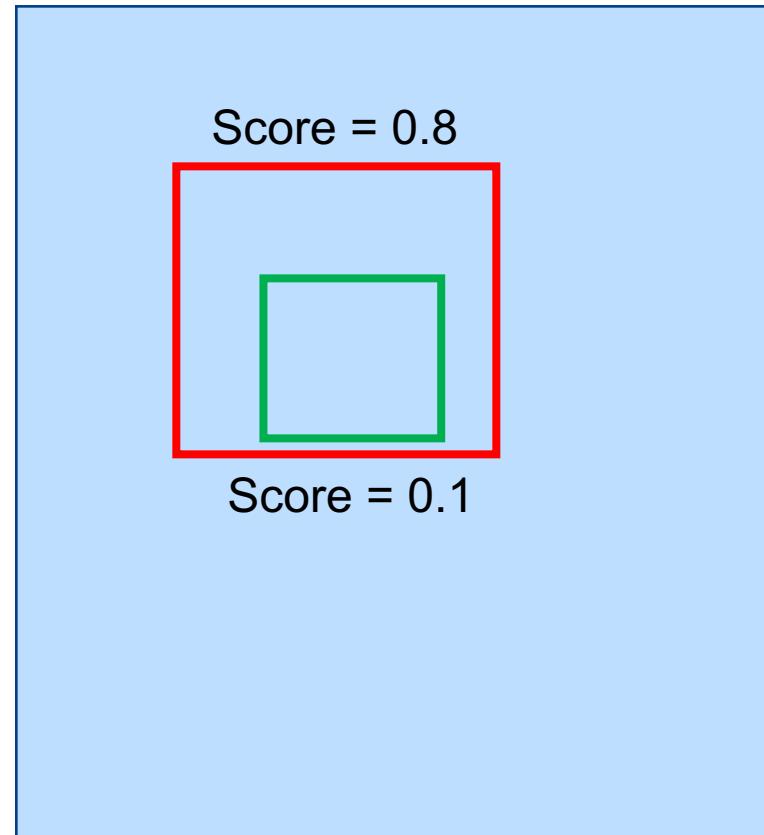
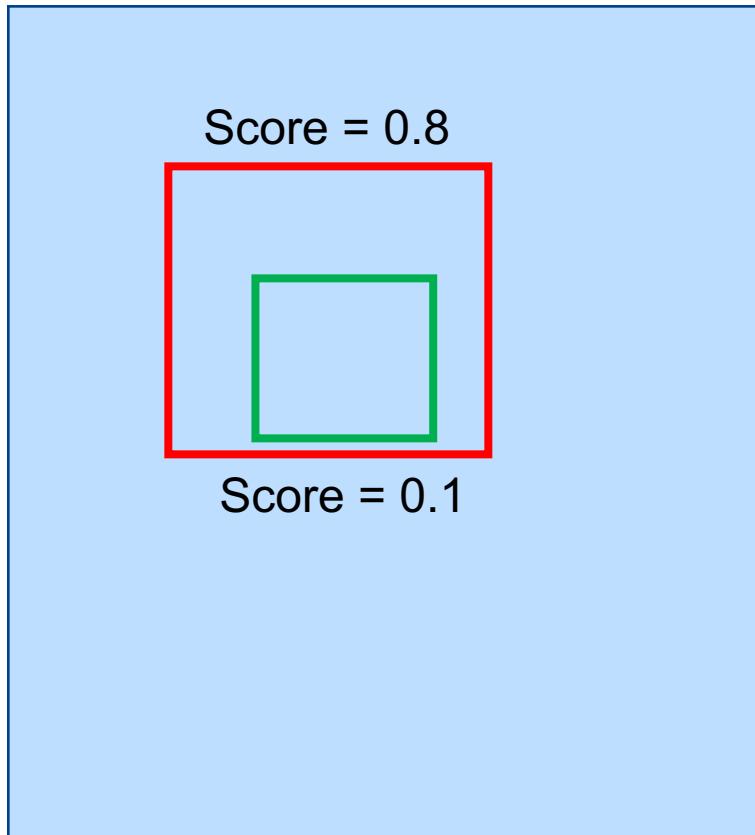
1. Non-max suppression



Resolving detection scores

$$IOU = \frac{\text{area of overlap}}{\text{area of union}} = \frac{\text{red} \cap \text{blue}}{\text{red} \cup \text{blue}}$$


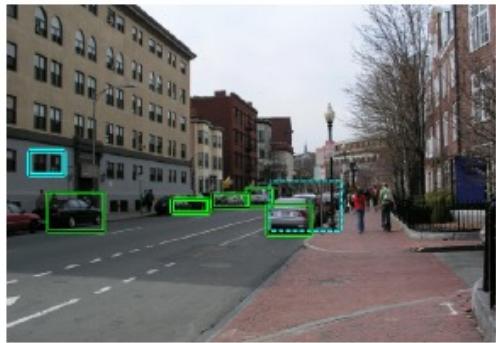
1. Non-max suppression



“Overlap” score is below some threshold

Resolving Detection Scores

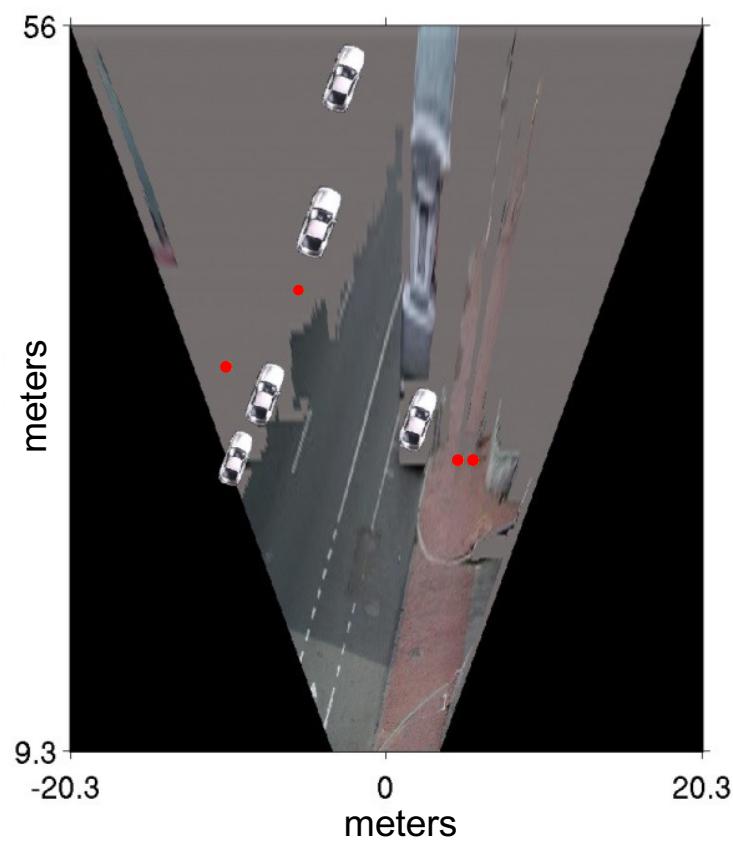
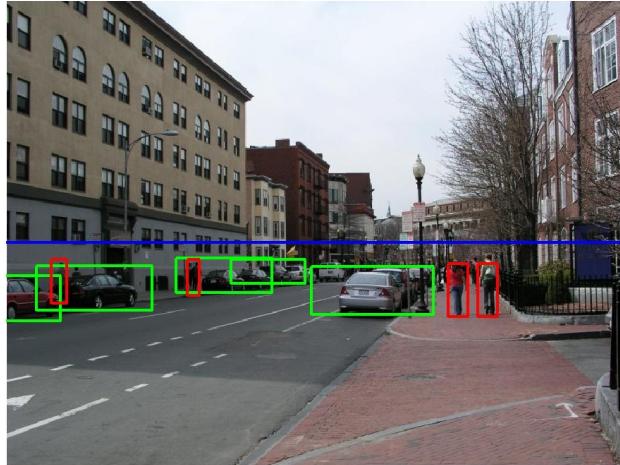
2. Context/reasoning



(g) Car Detections: Local



(h) Ped Detections: Local



Hoiem et al., 2006

Non-Max Suppression

H N



Confidence Scores

0 - 100

5 15 20 20 15 10 5 1 2 3 1

70	75	70	65
30	85	80	
90	95	90	
80	85	80	
75	75	70	

Set a high confidence score threshold of 95%



How to measure performance?

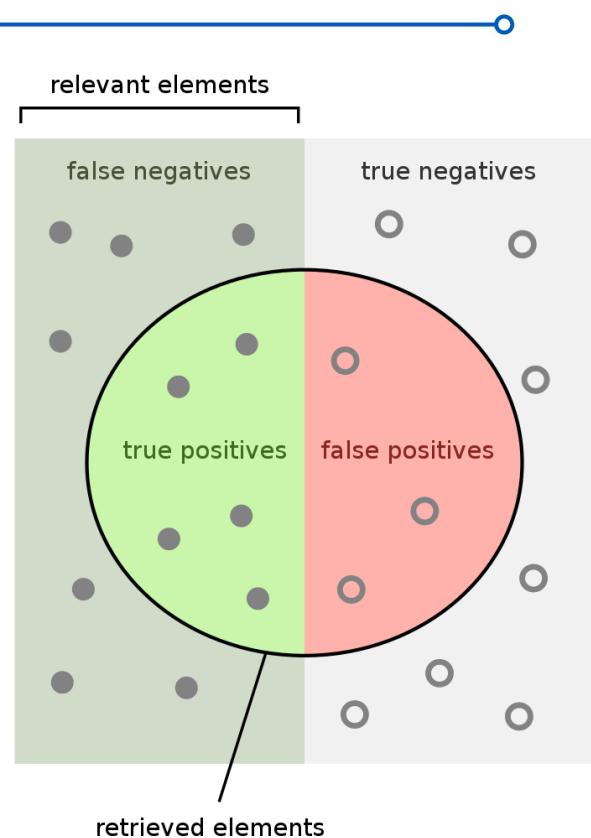
- Precision – How precise are the answers you gave?
 - $(\# \text{ Relevant}) / (\# \text{ Total Returned})$
- Recall – How many did you find of the ones that could be found?
 - $(\# \text{ Relevant}) / (\# \text{ Total Relevant})$
- F Measure – Harmonic Mean of Precision and Recall
 - $(2 * \text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$

Precision, Recall, F1 (Recap)

$$\text{Precision} = \frac{\text{True Positive}}{\text{Predicted Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{Actual Positive}}$$

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$



How many retrieved items are relevant?

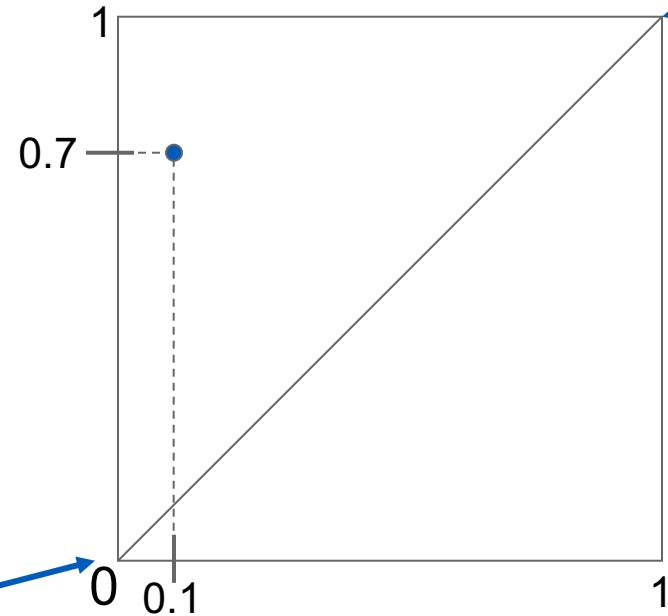
$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are retrieved?

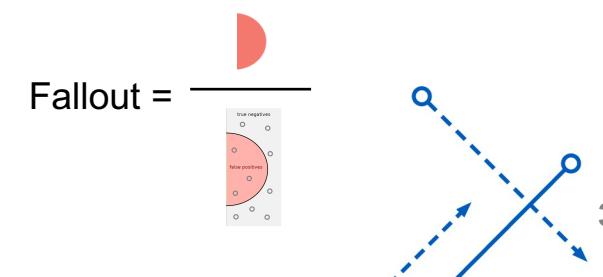
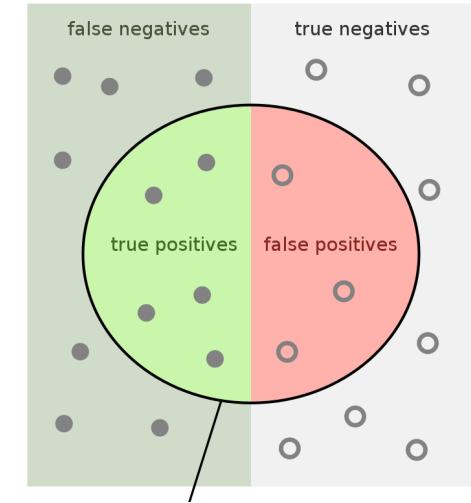
$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

ROC curve (Recap)

- How can we measure the performance of a feature matcher?



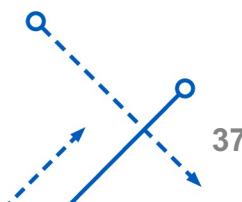
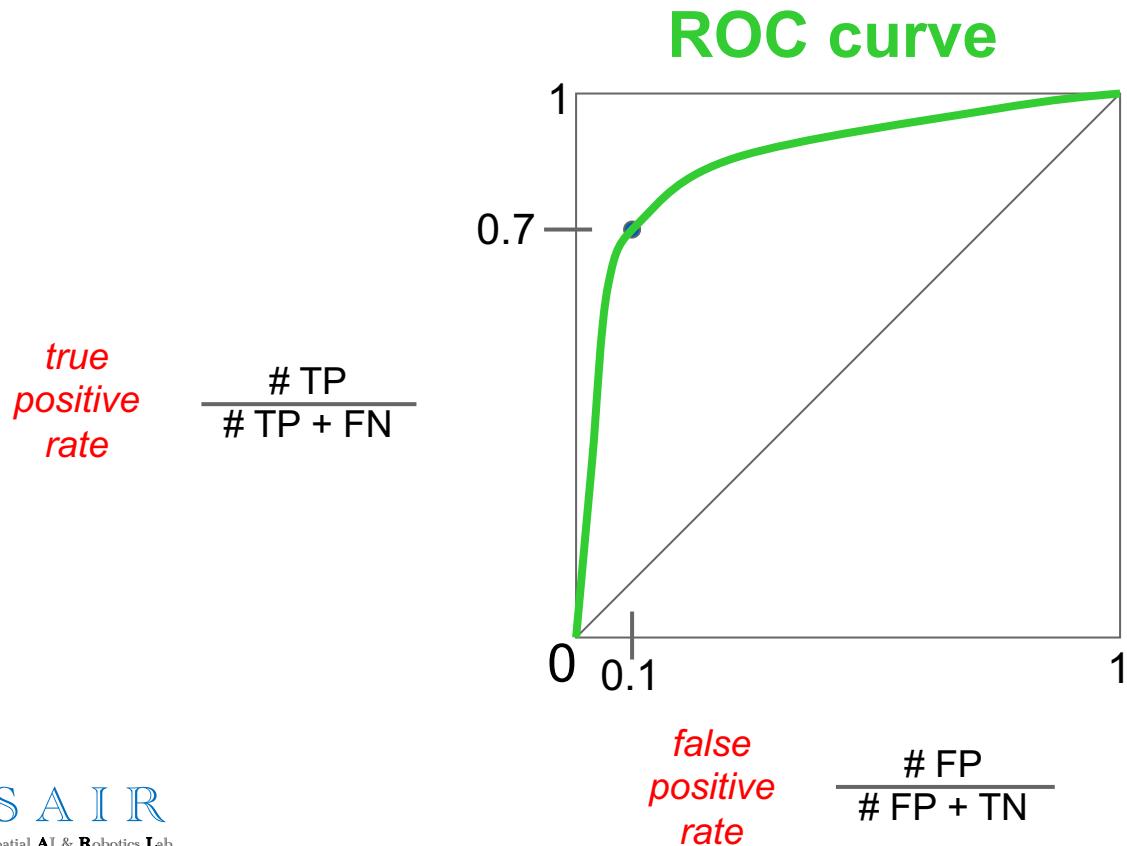
Lowest threshold
all pairs as match



ROC curve (Recap)

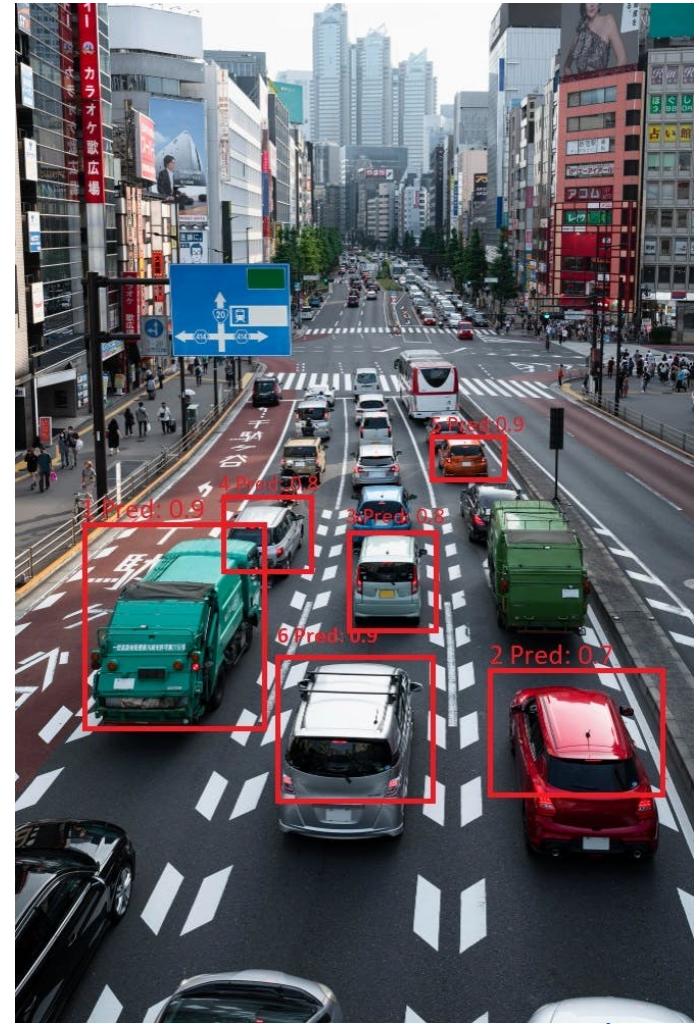
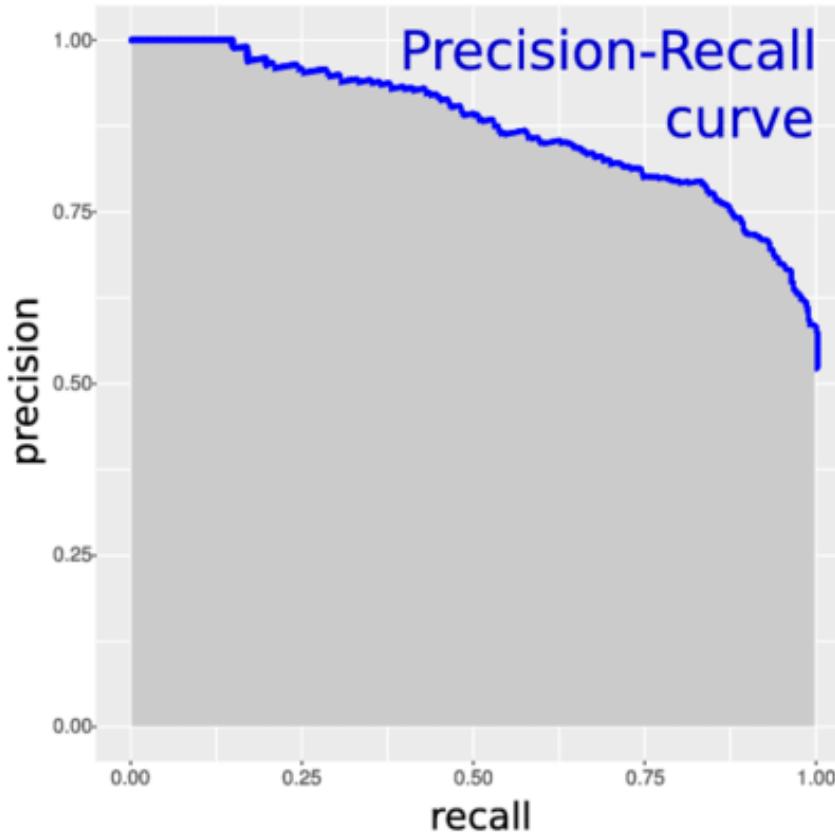
ROC Curves (Receiver Operator Characteristic)

- Generated by counting # correct/incorrect matches, for different thresholds.
- Want to maximize area under the curve (AUC)
- Useful for comparing different feature matching methods.



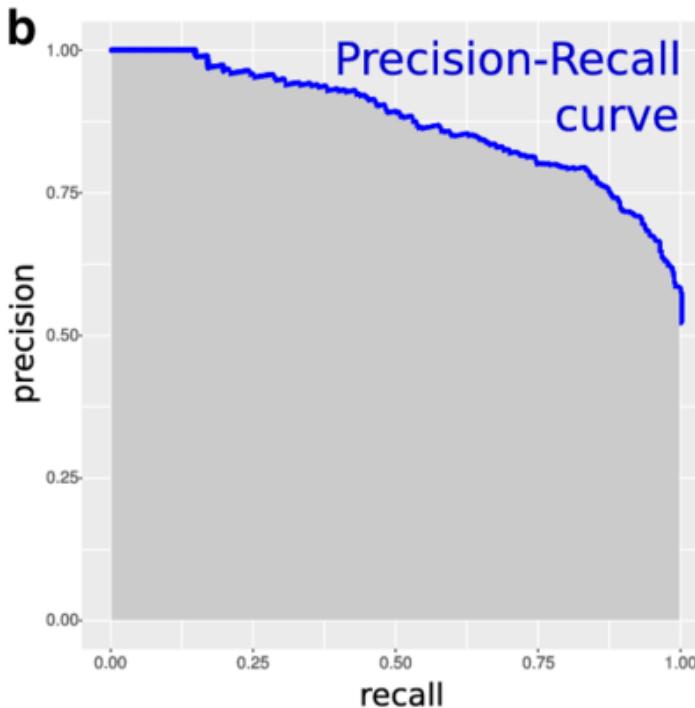
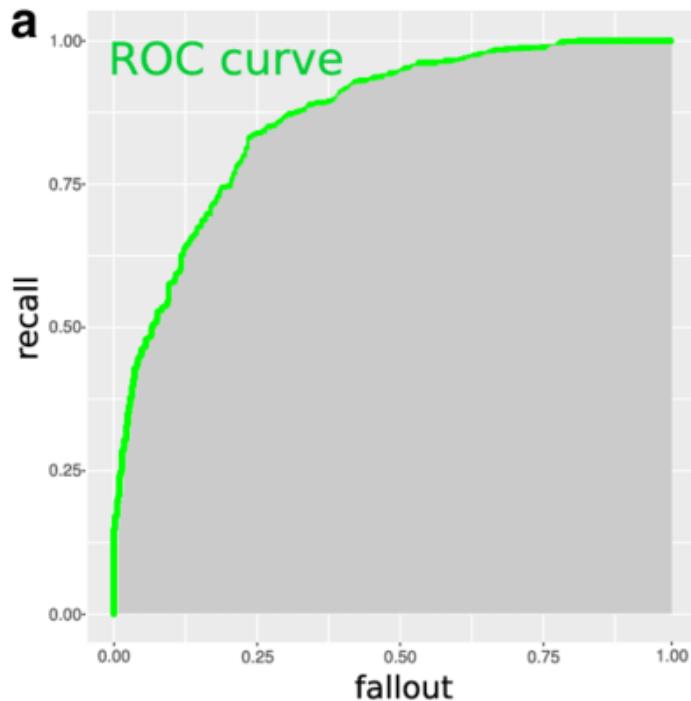
Precision and Recall Curve (PR curve)

precision = #relevant / #returned
recall = #relevant / #total relevant



Slide credit: Ondrej Chum 38

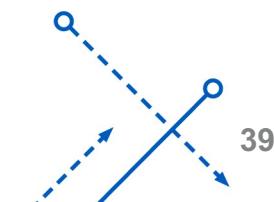
ROC curve VS. PR curve



Fallout =

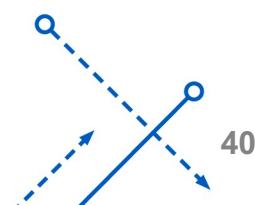
A 2x2 confusion matrix diagram. The top-left cell is labeled "True negative". The top-right cell is labeled "False positive". The bottom-left cell is labeled "True positive". The bottom-right cell is labeled "False negative". A red arrow points from the "False positive" cell to the "Fallout" equation above it.

Recall =



ROC curve VS. PR curve

- ROC Curves (Lecture 4)
 - summarize the trade-off between the **true positive rate** and **false positive rate** for a predictive model using different probability thresholds.
- PR curves
 - summarize the trade-off between the **true positive rate** and the **positive predictive value** for a predictive model using different probability thresholds.



ROC curve VS. PR curve

- ROC: good for balanced datasets/classes.
- PR curves: appropriate for imbalanced datasets.
- If the proportion of positive to negative instances changes in a test set, the ROC curves will not change, so do not depend on class distributions.
- ROC “weights” equally true positives and true negatives (not expected sometimes, e.g., place recognition).

