

# Computer vision and Robotics

Roberto Cipolla

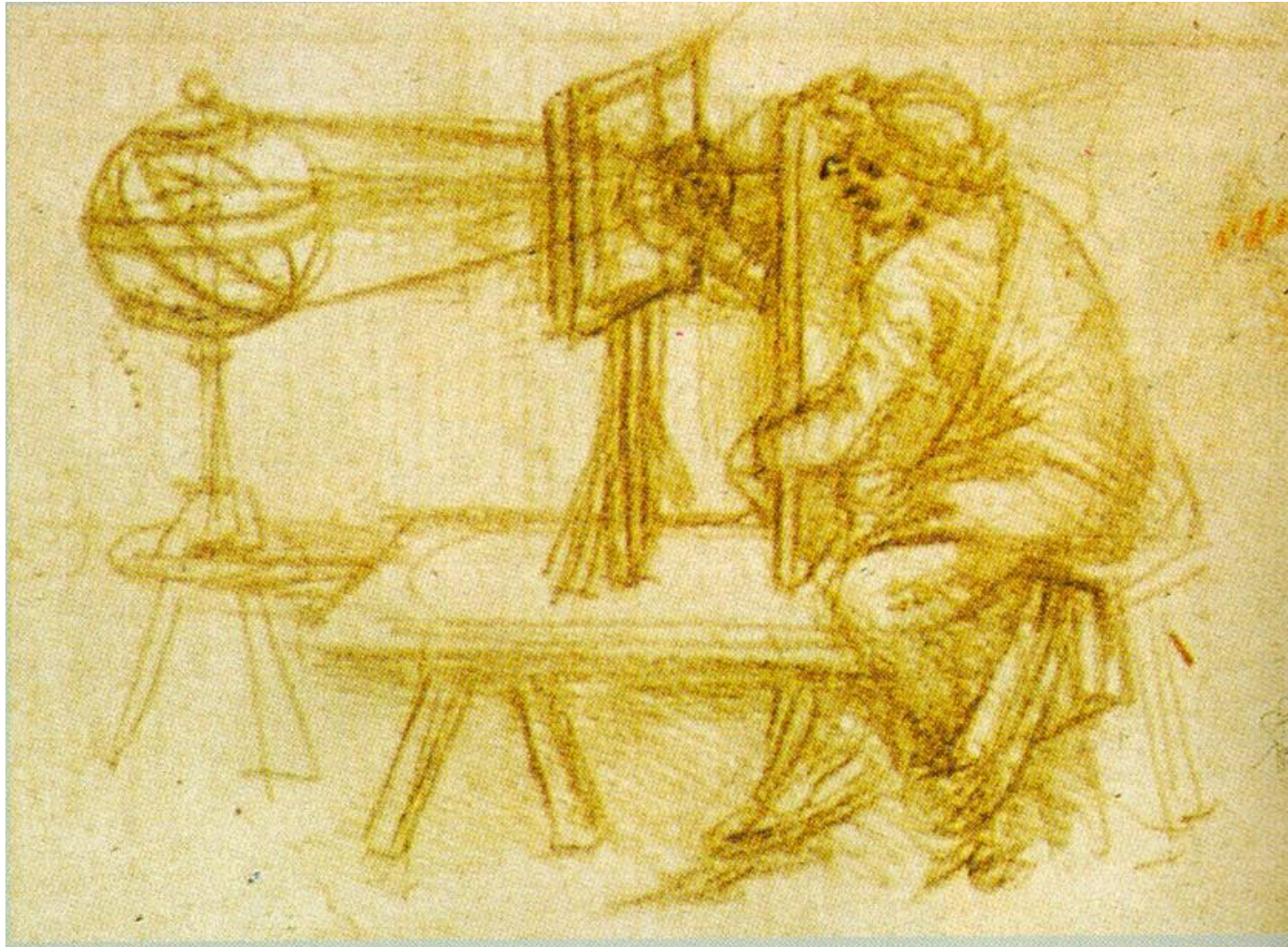
Department of Engineering

Research team

<http://www.eng.cam.ac.uk/~cipolla/people.html>

# Perspective

---



# Making machines see

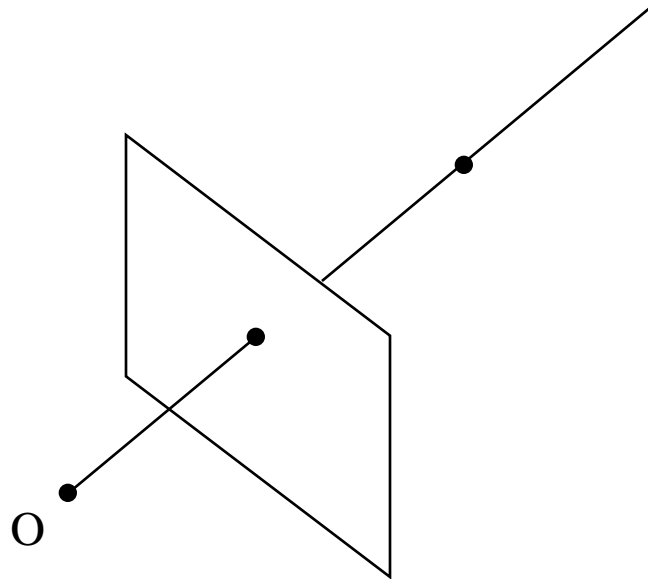
---

- What is vision and how to duplicate it?
- 3D shape: making digital copies of sculpture from photographs from multiple viewpoints
- Image matching and localisation from a single photo using a mobile (camera) phone
- Detection and tracking of objects: hands, faces and people
- Machine learning – object categorization and recognition

# Stereo vision and 3D shape

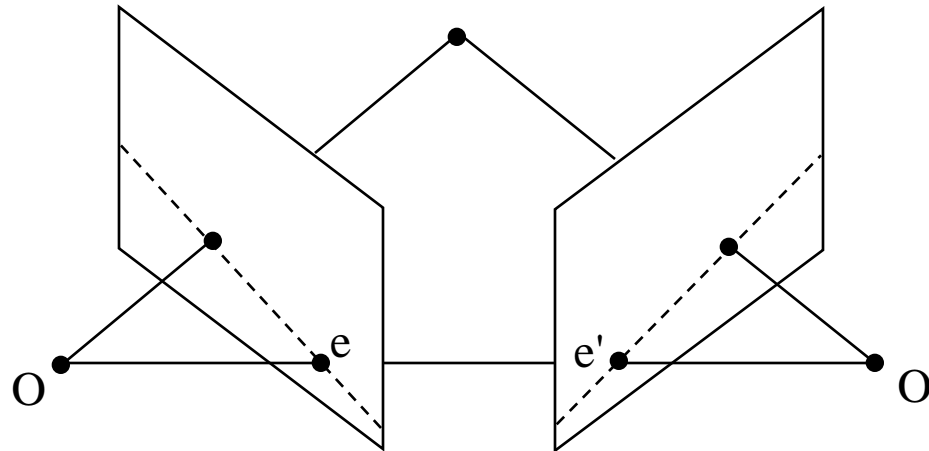
# Ambiguity in a single view

---



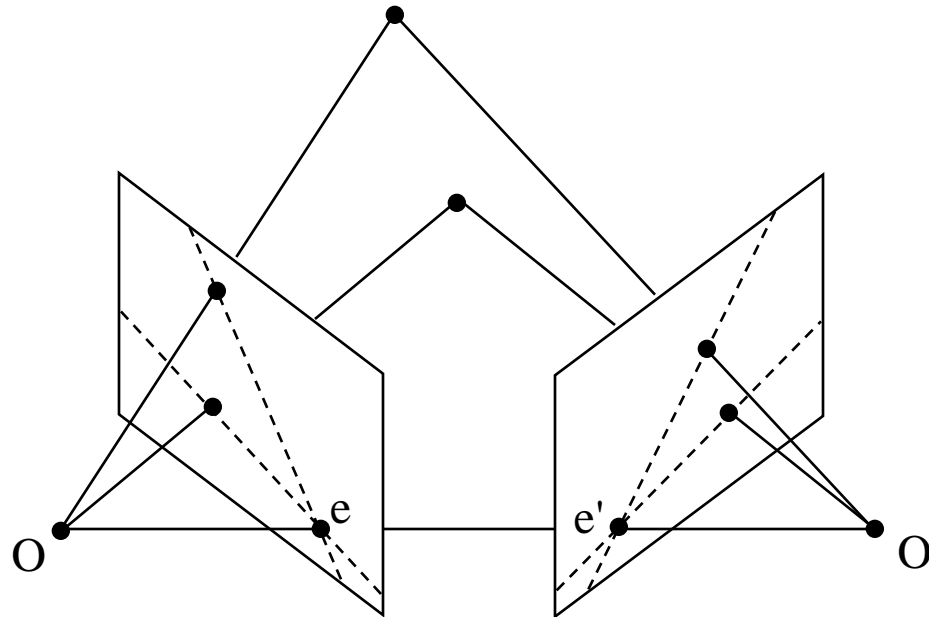
# Stereo vision

---



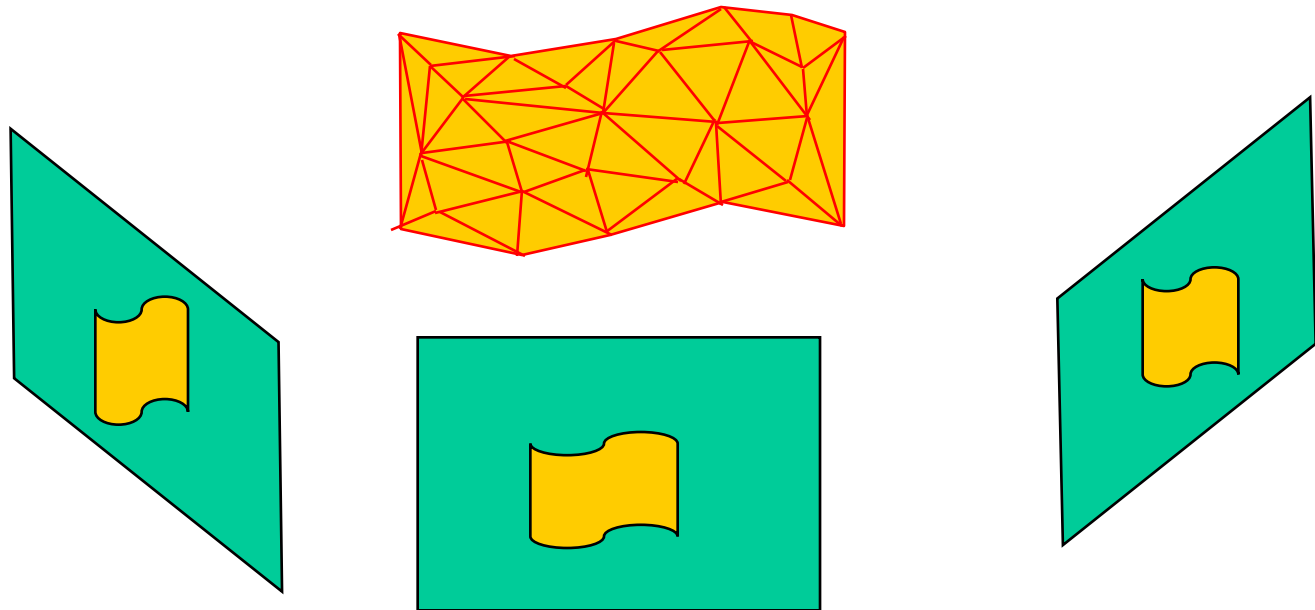
# Stereo vision

---



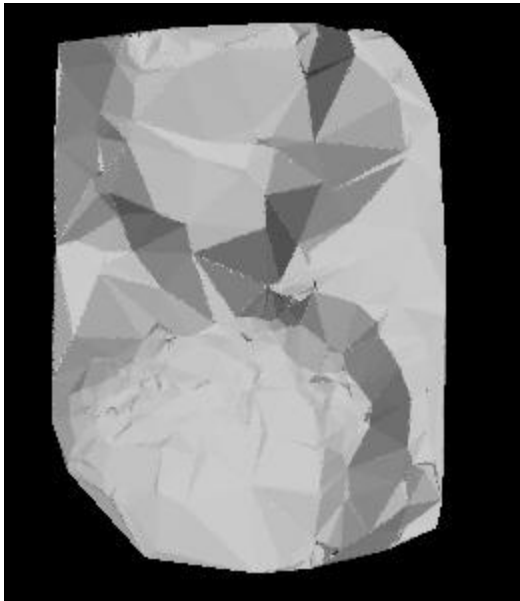
# Shape recover problem:

---



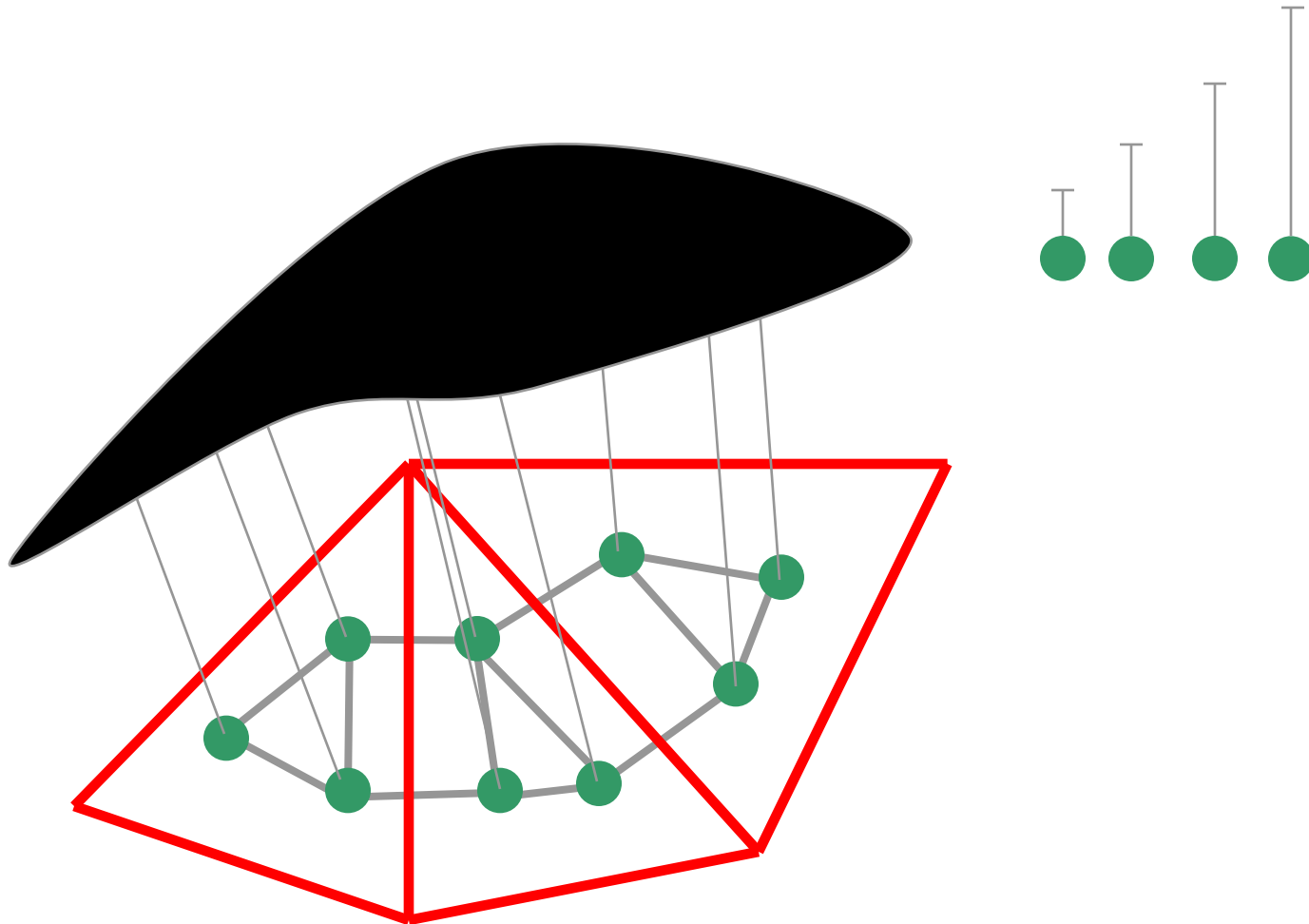


# Dense stereo



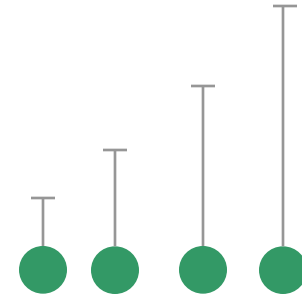
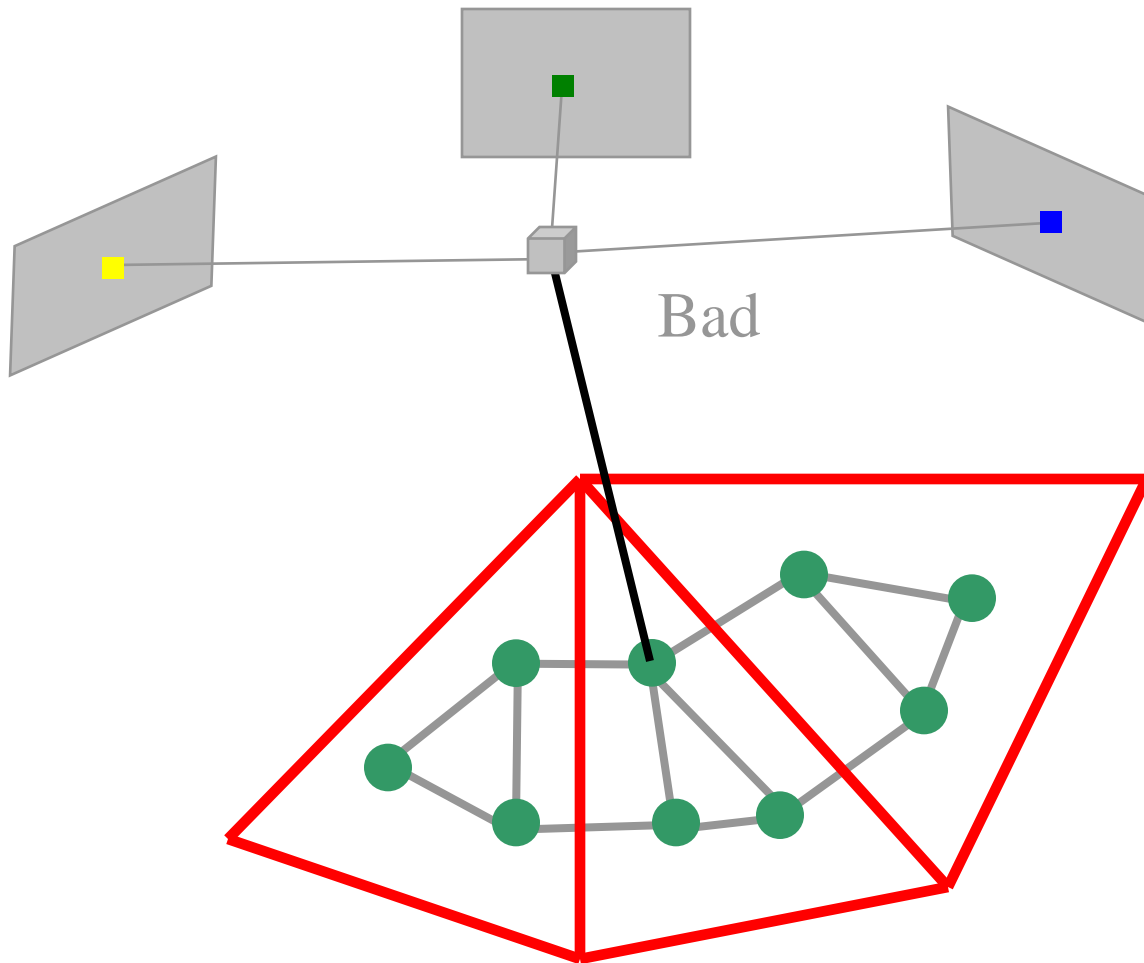
# Surface + height

‘colours’ :



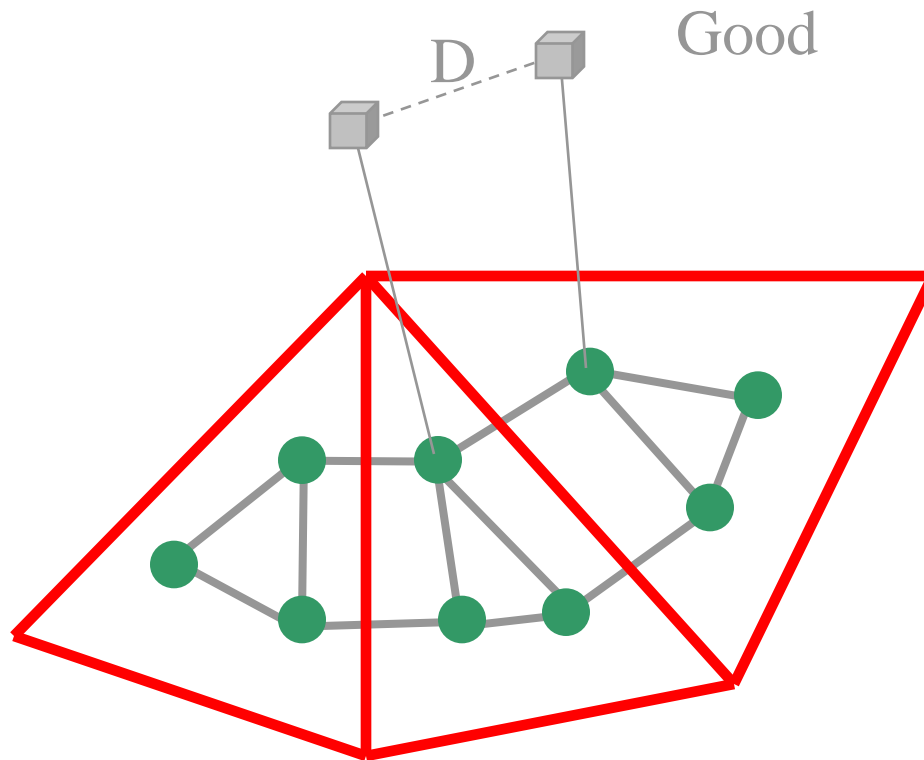
# Surface + height

‘colour’ cost :



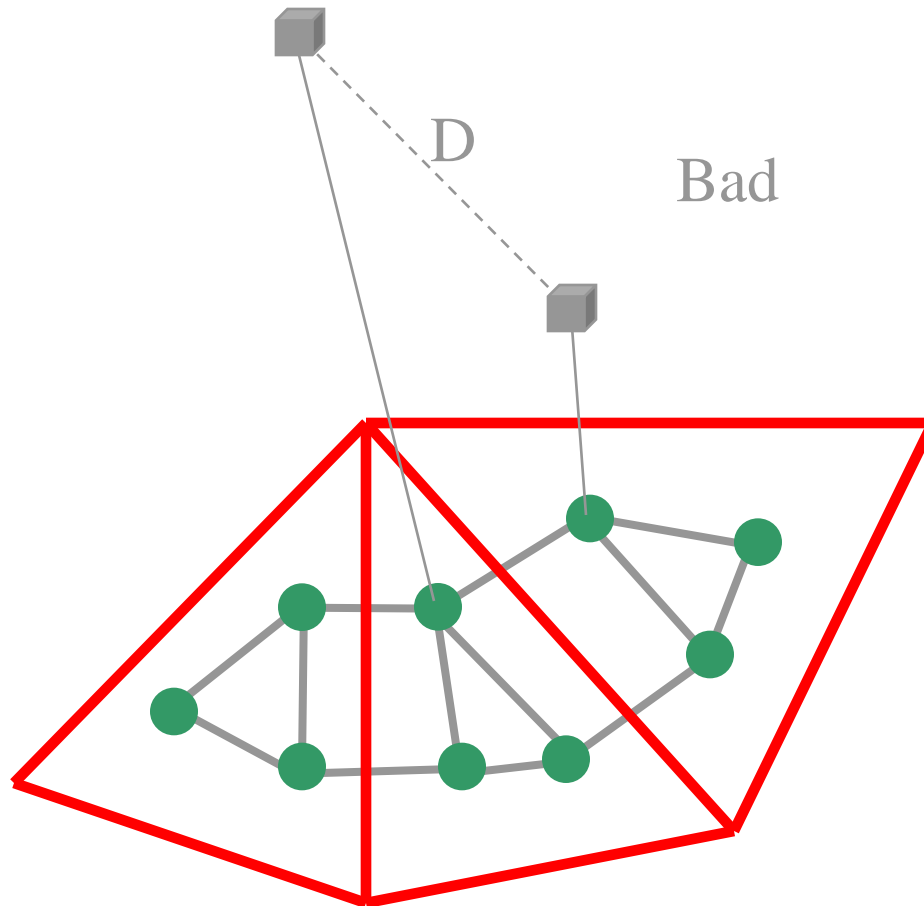
# Surface + height

Neighbour cost :



# Surface + height

Neighbour cost :



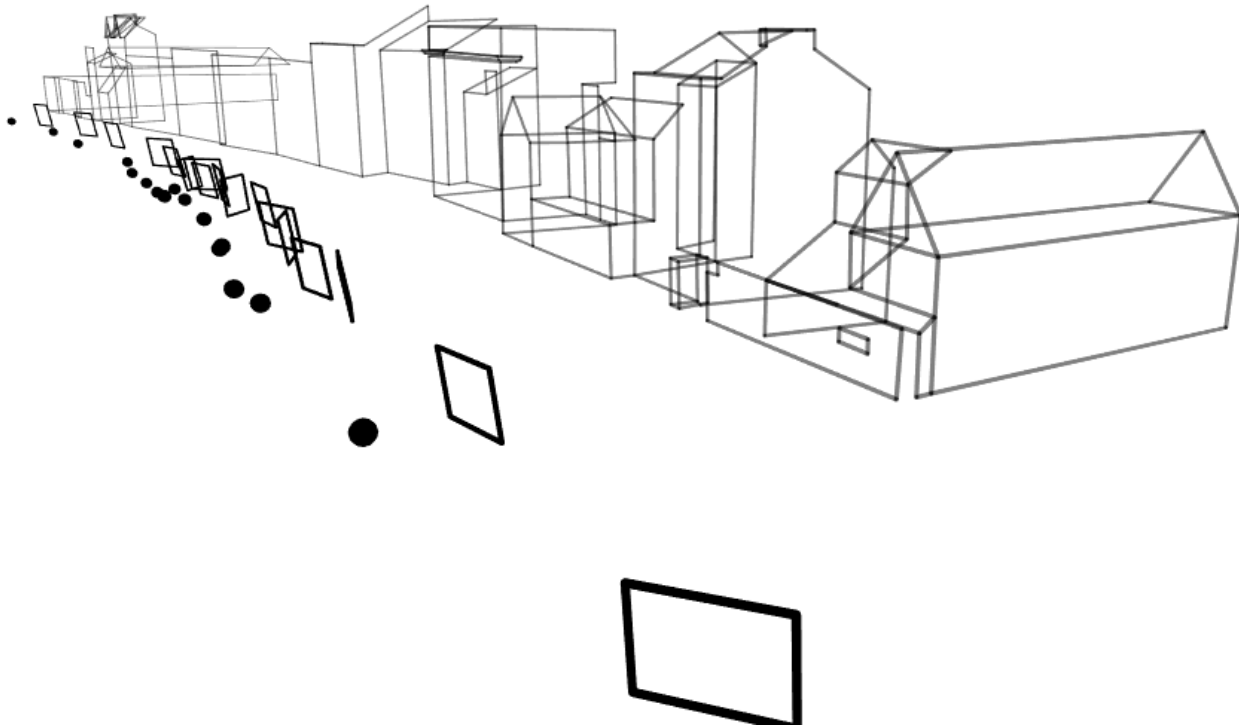
# Digital copies of sculpture – Digital Pygmalion

# Input images

---



# Building 3D models of cities





# Trumpington Street Data

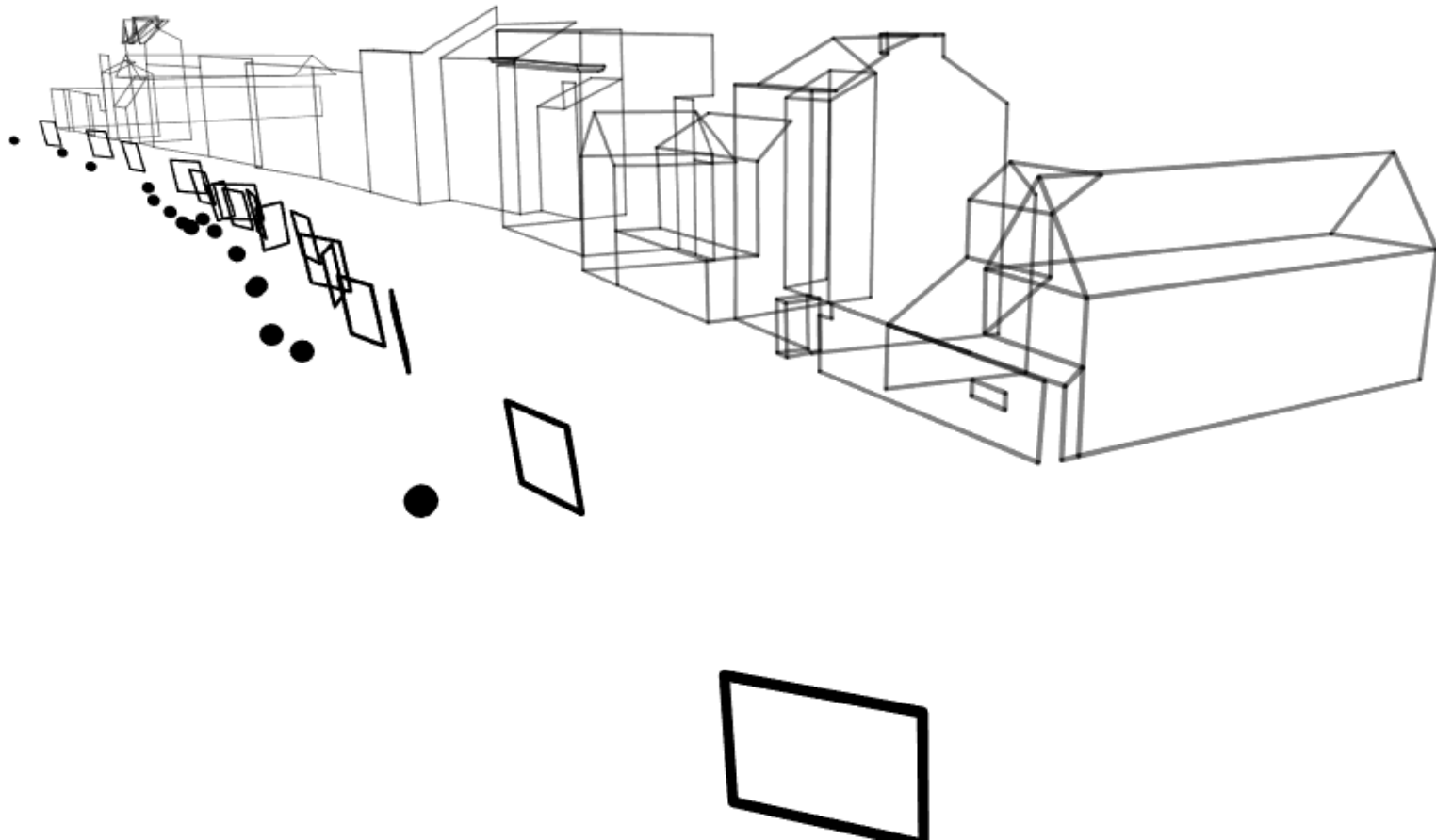


# Camera pose determination



# 3D reconstruction

---



# Reconstruction texture mapped



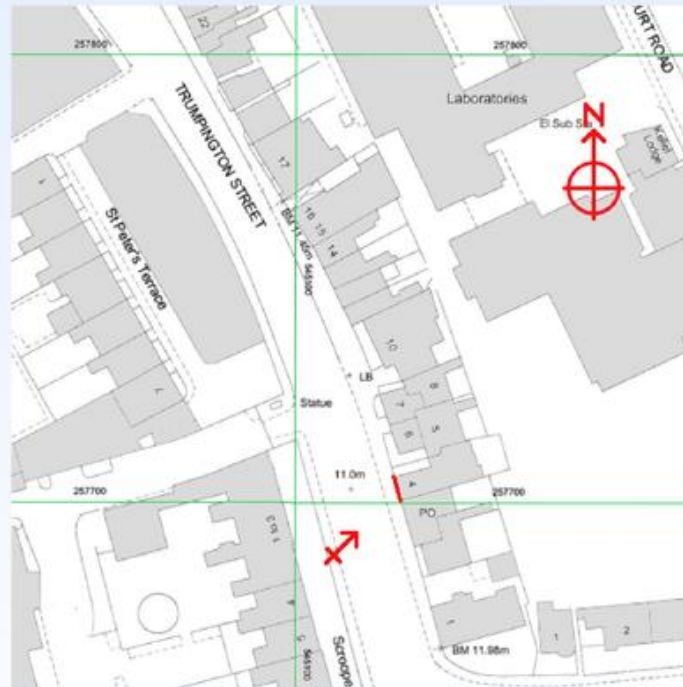
# Image-Based Localisation

## Where am I?

# The goal – where am I?



User takes a picture of a nearby building. System tells you what you are looking at and exactly where you are on a map.

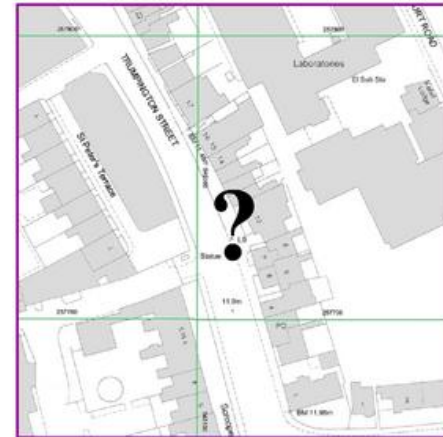


# The problem

---



?  
=



# Why difficult?

---



Extreme perspective distortion



Differences in colour / lighting conditions



Occlusion



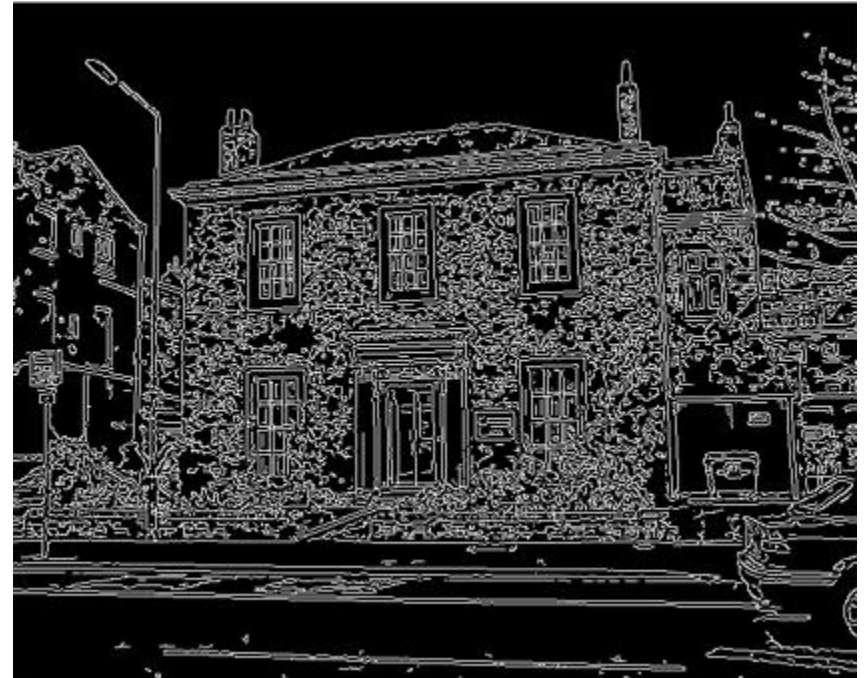
# Constrained matching

---



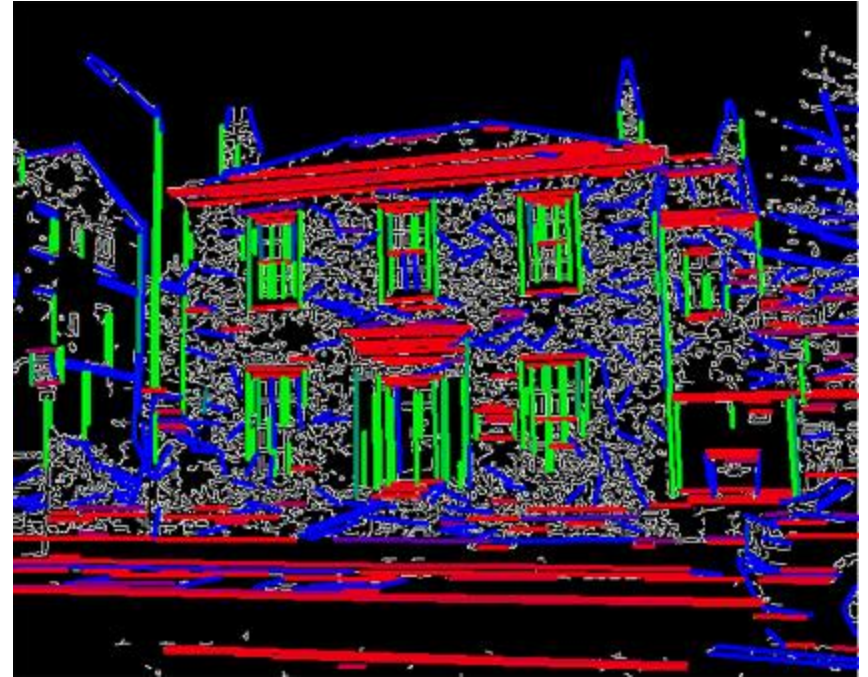
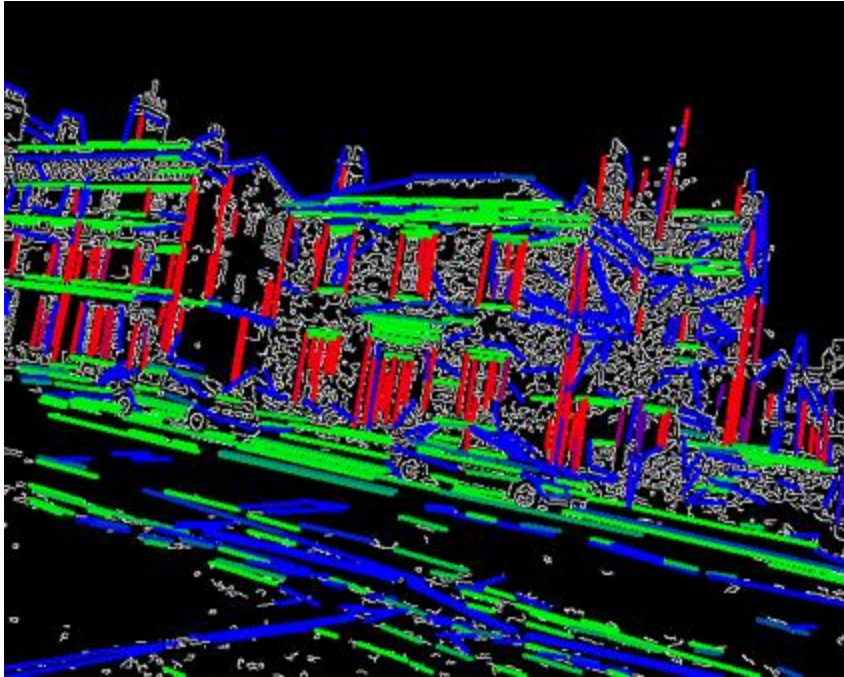
# Constrained matching

---



# Constrained matching

---



# Constrained matching

---

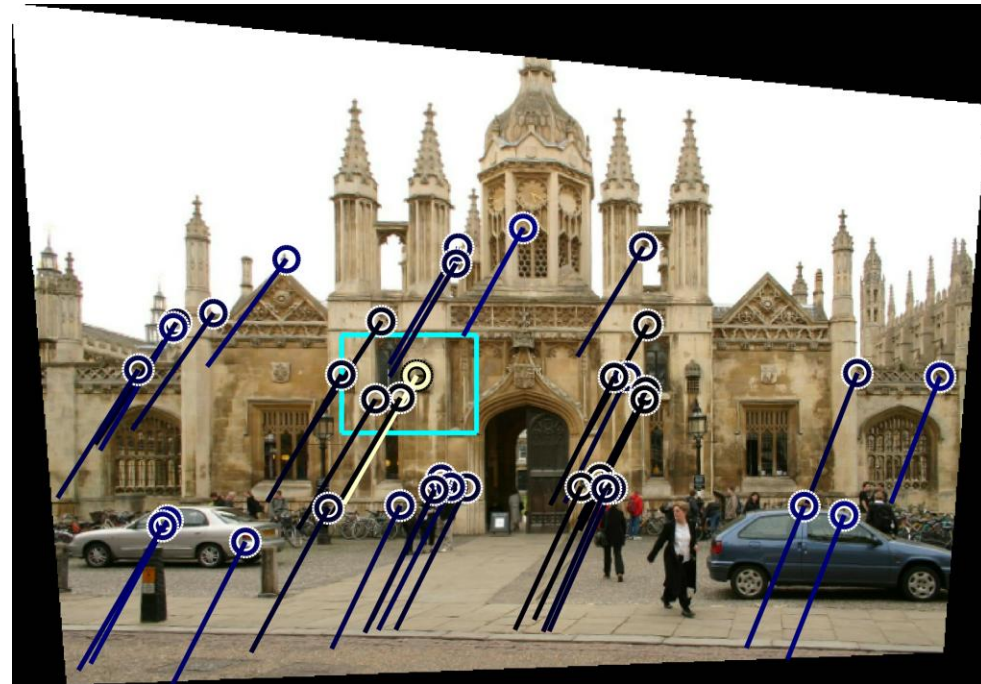


# Constrained matching

---

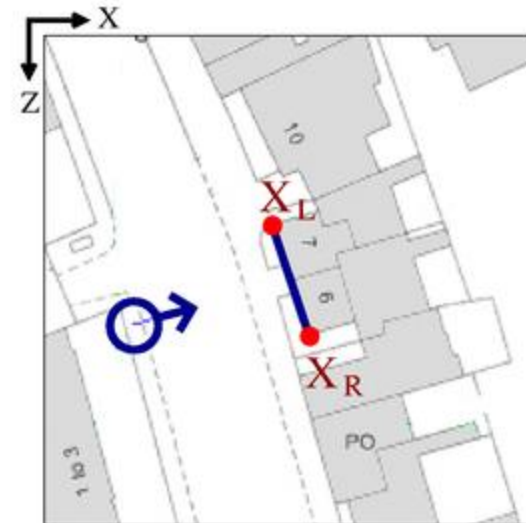
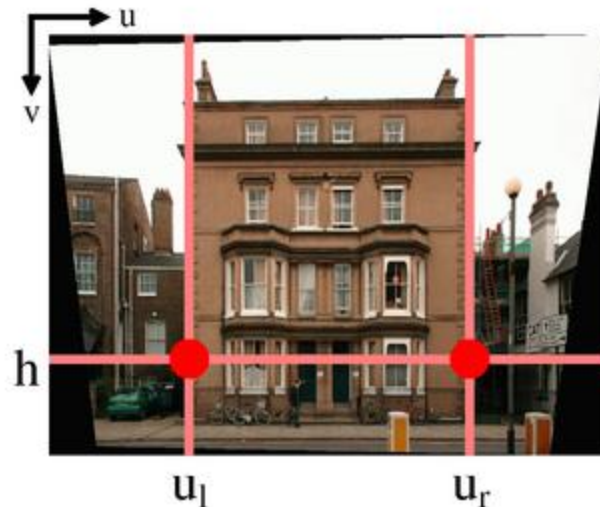


# Matching



# Register database view

First align database view to map



# Localisation

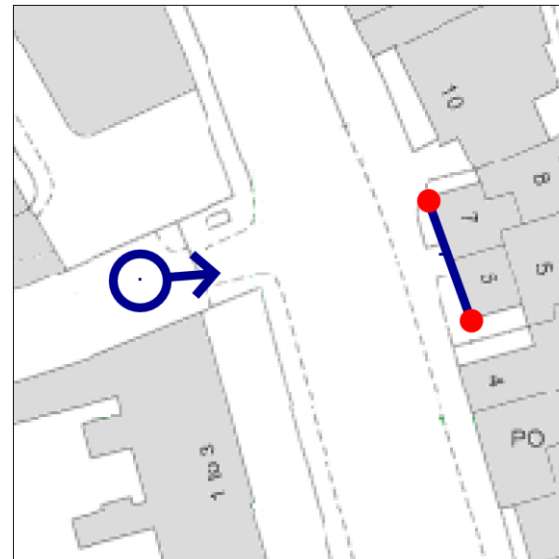
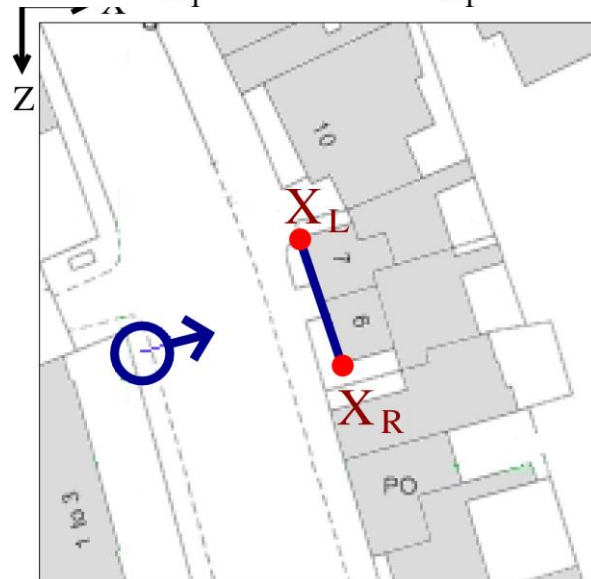
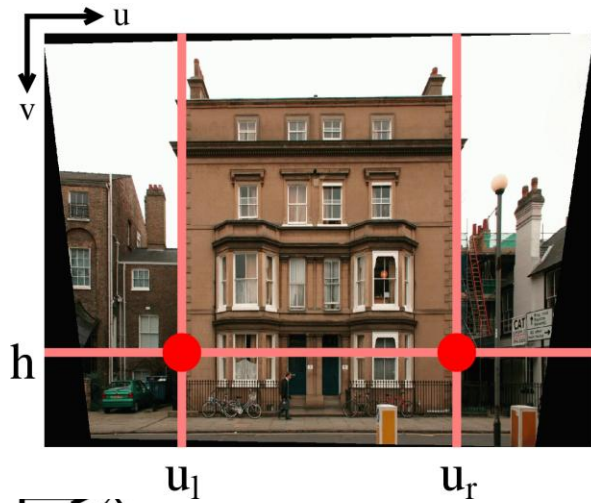
Knowing the rectifying homography ( $H_{\perp}$ ), the alignment ( $H_A$ ), and the database view registration, can work backwards to find user:



Rectifying rotation  $R_{\perp}$  gives the angle from perpendicular and focal length the distance to camera.



# Localisation of query view

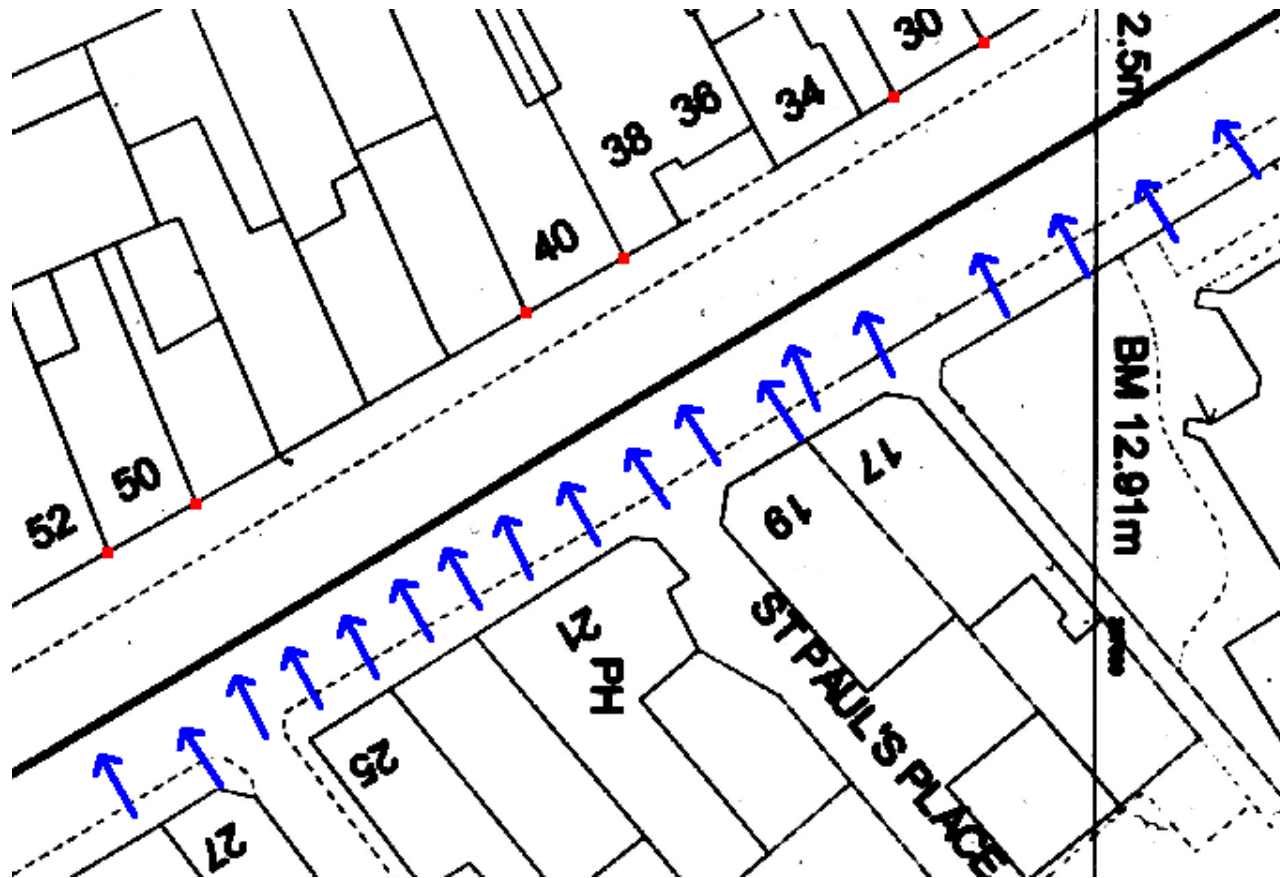


# Image-based localisation

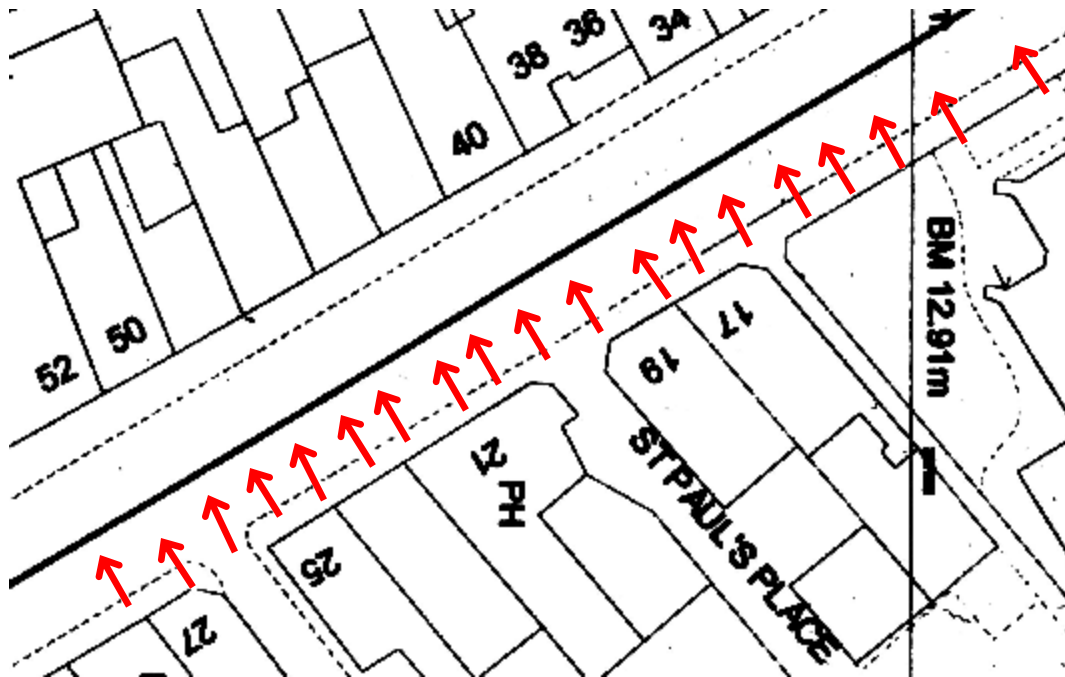
...



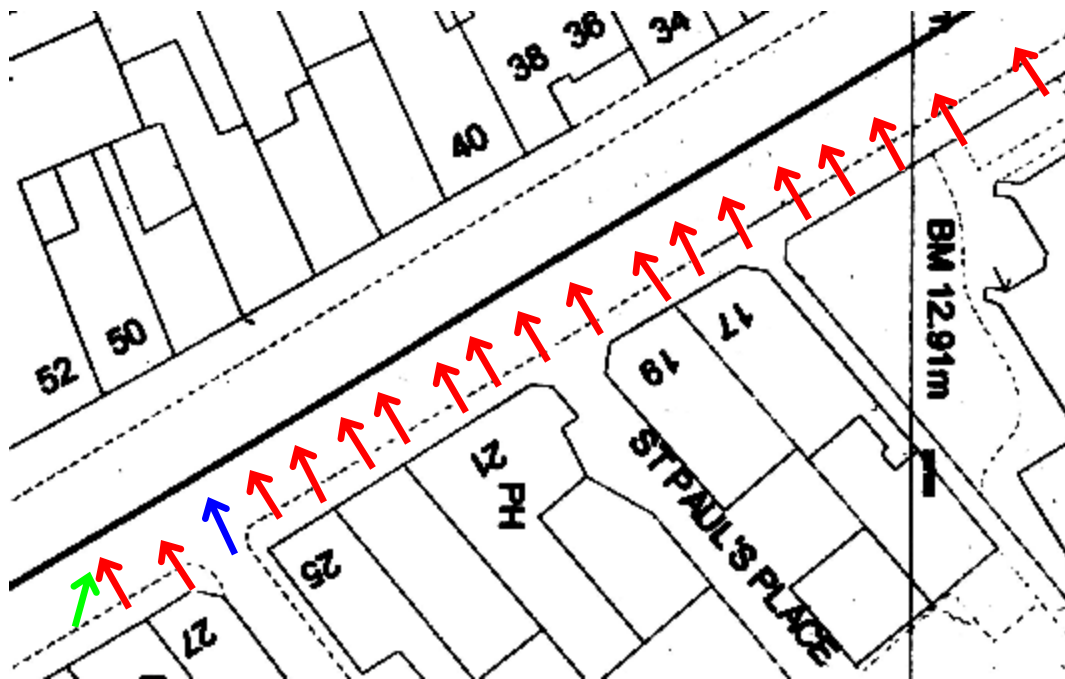
...



# Image-based localisation



# Image-based localisation



# Object detection and tracking

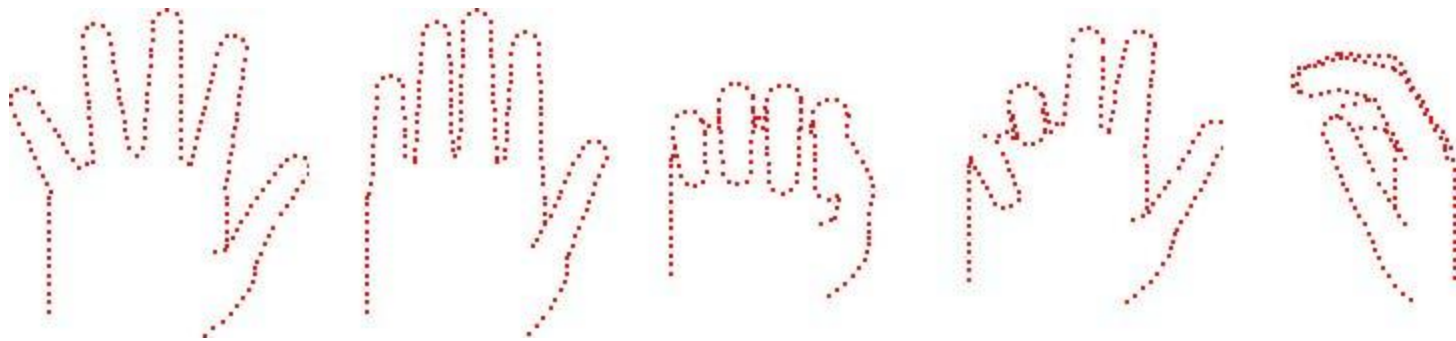
# Hand detection system

---



# Template-based Detection

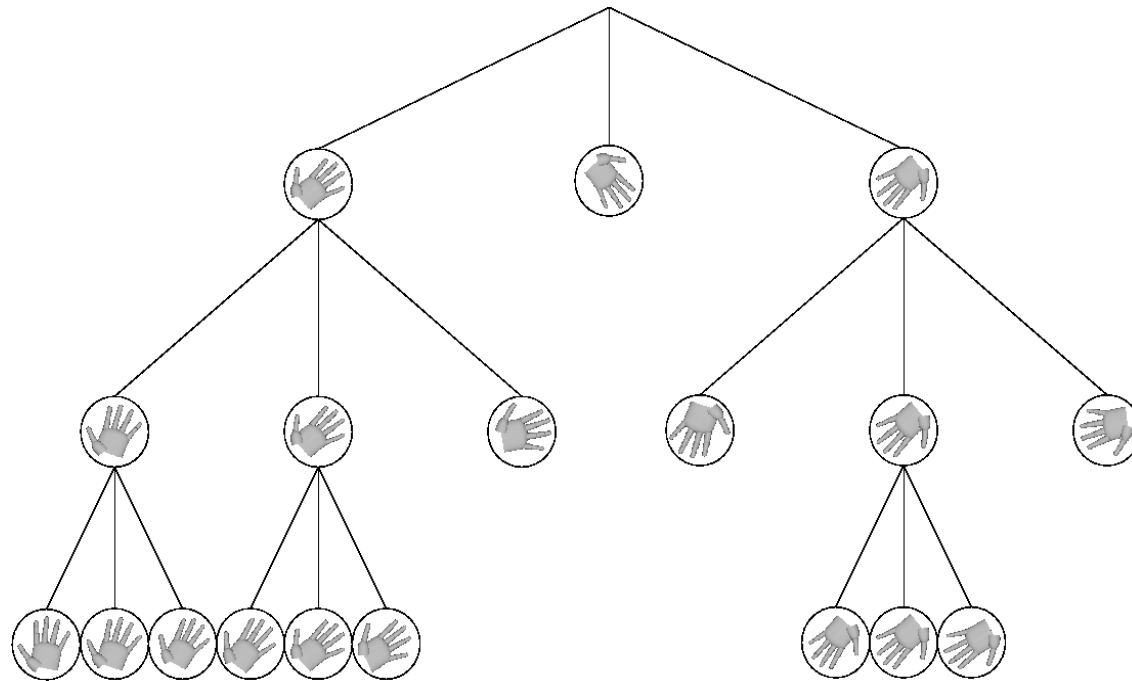
---



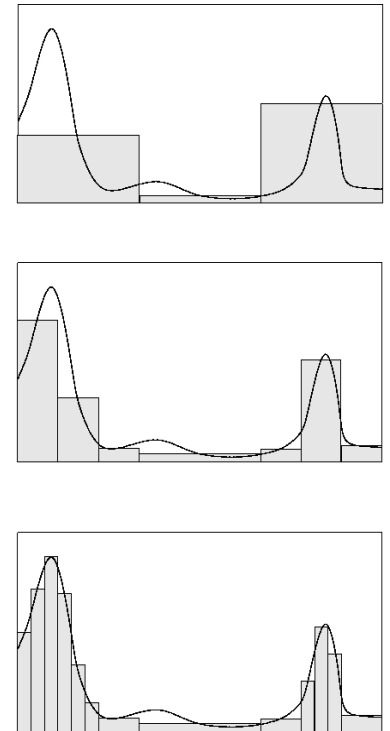
- Large number of templates are generated off-line to handle global motion and finger articulation.
- Need for
  - Inexpensive template-matching function
    - Distance Transform and Chamfer Matching
  - Efficient search structure
    - Bayesian Tree structure

# Bayesian-Tree

State space partitioning



Estimation of posterior pdf



- The search-tree is brought into a Bayesian framework by adding the prior knowledge from previous frame.
- The Bayesian-Tree can be thought as approximating the posterior probability at different resolutions.



# Tracking - 3D mouse

---

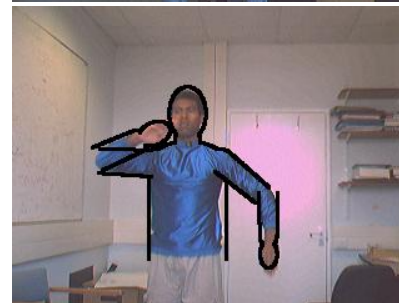
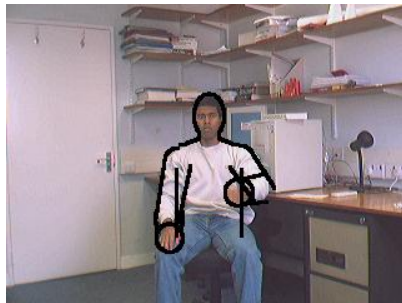


# Opening and closing

---

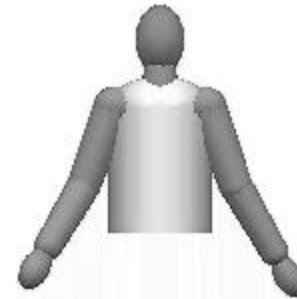


# People and pose detection



# A Tracked Sequence

---



# Detecting and tracking people in crowds

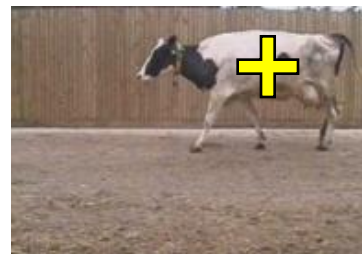
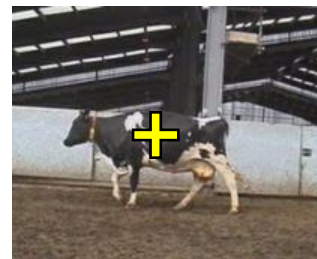
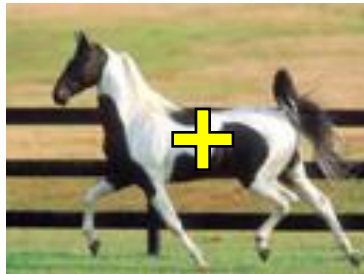
# Learning object categories

# Machine learning

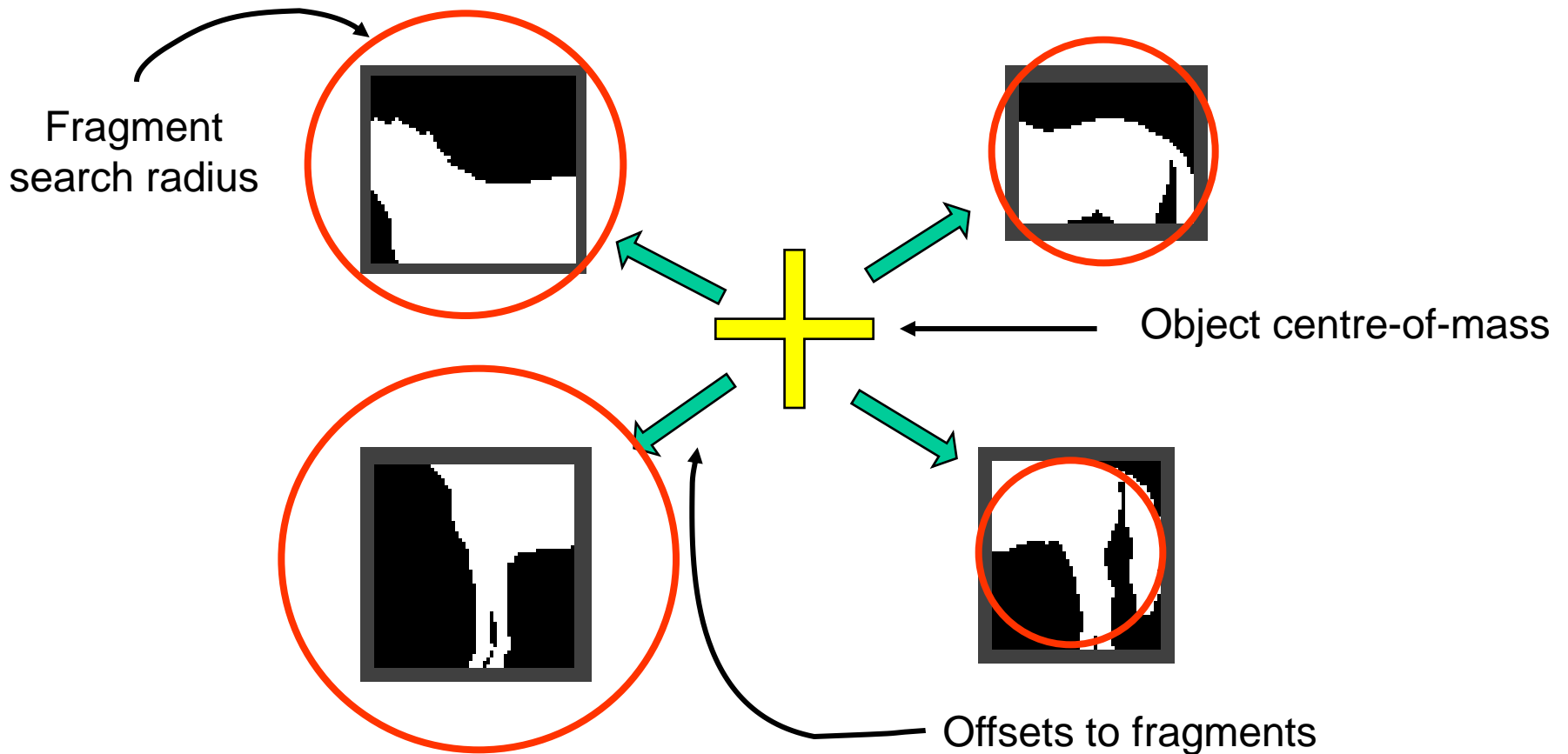
---

- Learn to recognise images of a particular class, localised in space and scale
- i.e. find the horse/cow/car etc!

Desired  
Results



# Learn Object Model





# Learning a Detector

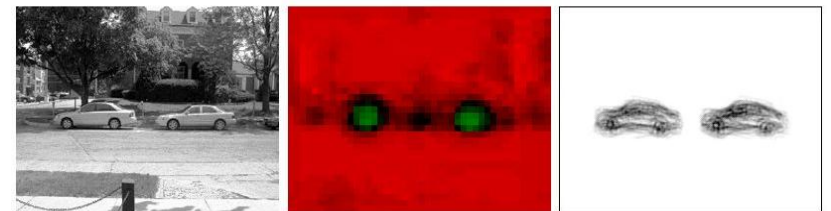
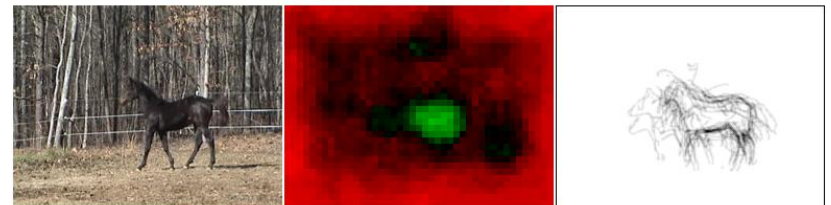
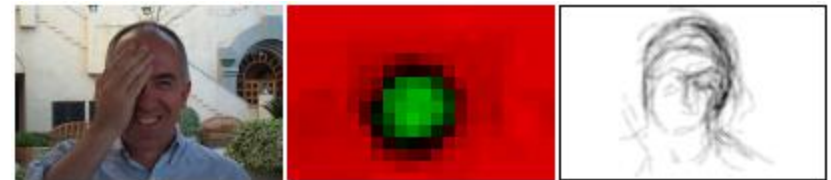
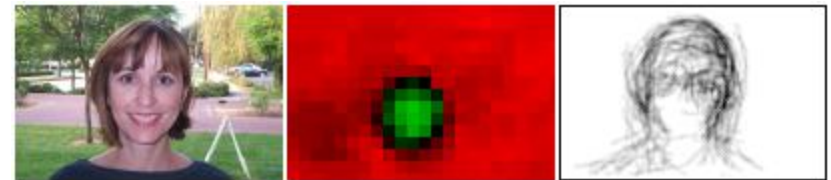
---

- Take large number ( $\sim 1000$ ) of candidate fragments of contour
  - randomly chosen from training set of masks
- Calculate chamfer scores for each fragment
  - over training set of images with known centroids
- Boosting algorithm selects a discriminative subset of fragments ( $\sim 100$ ) and learns their model parameters

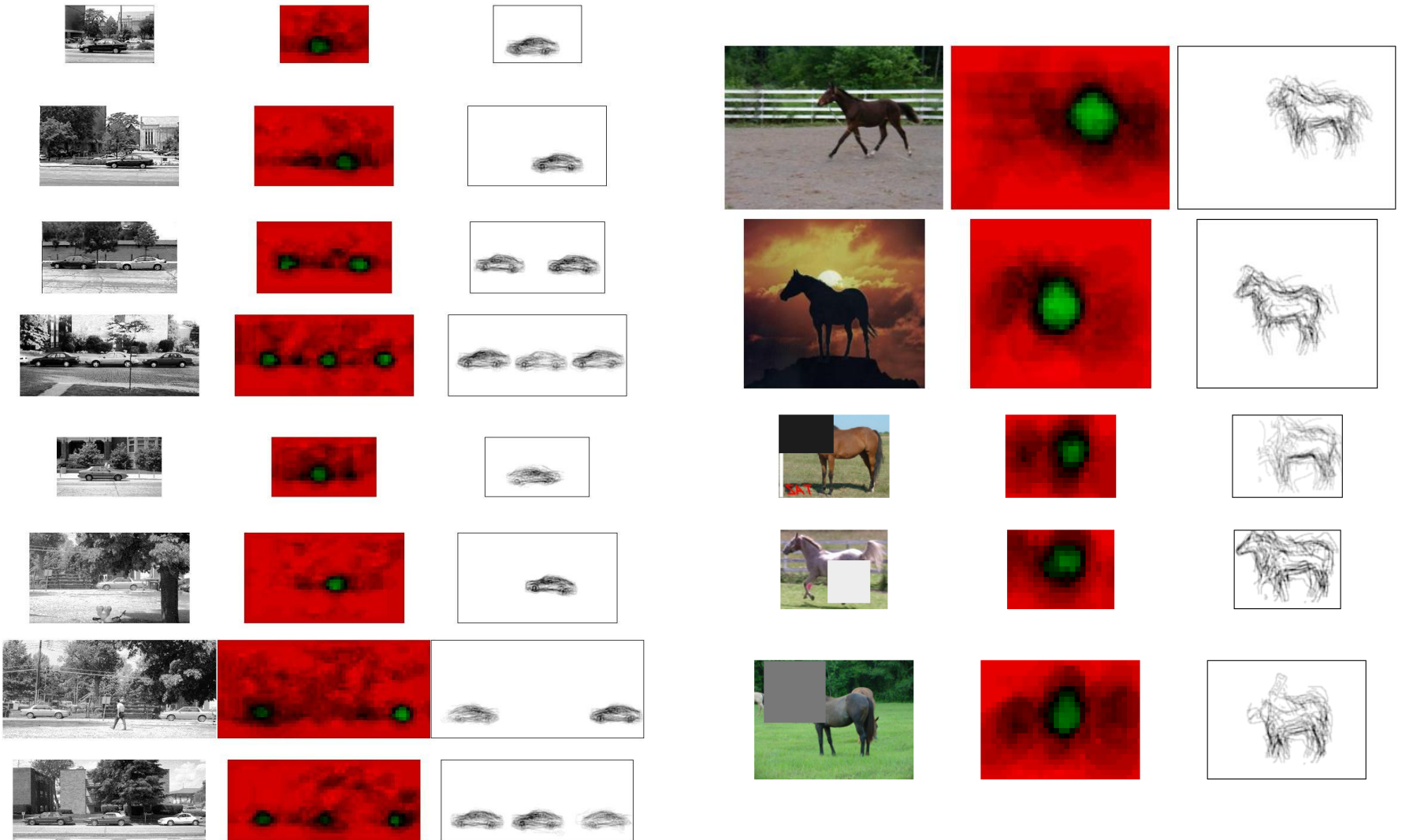
# Object Detection

---

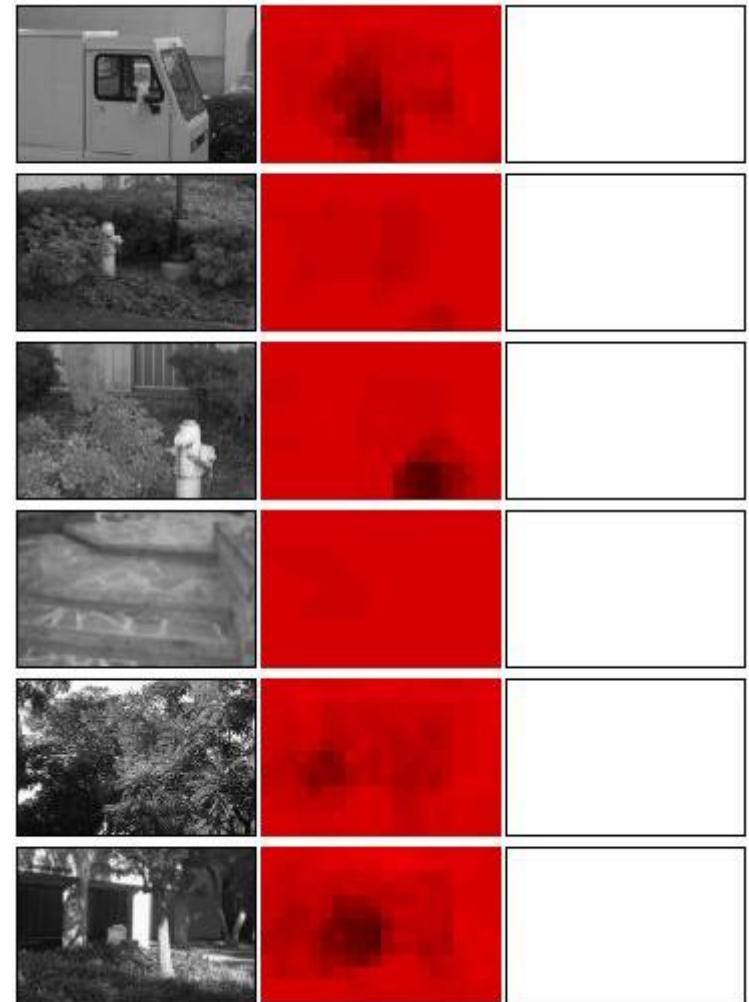
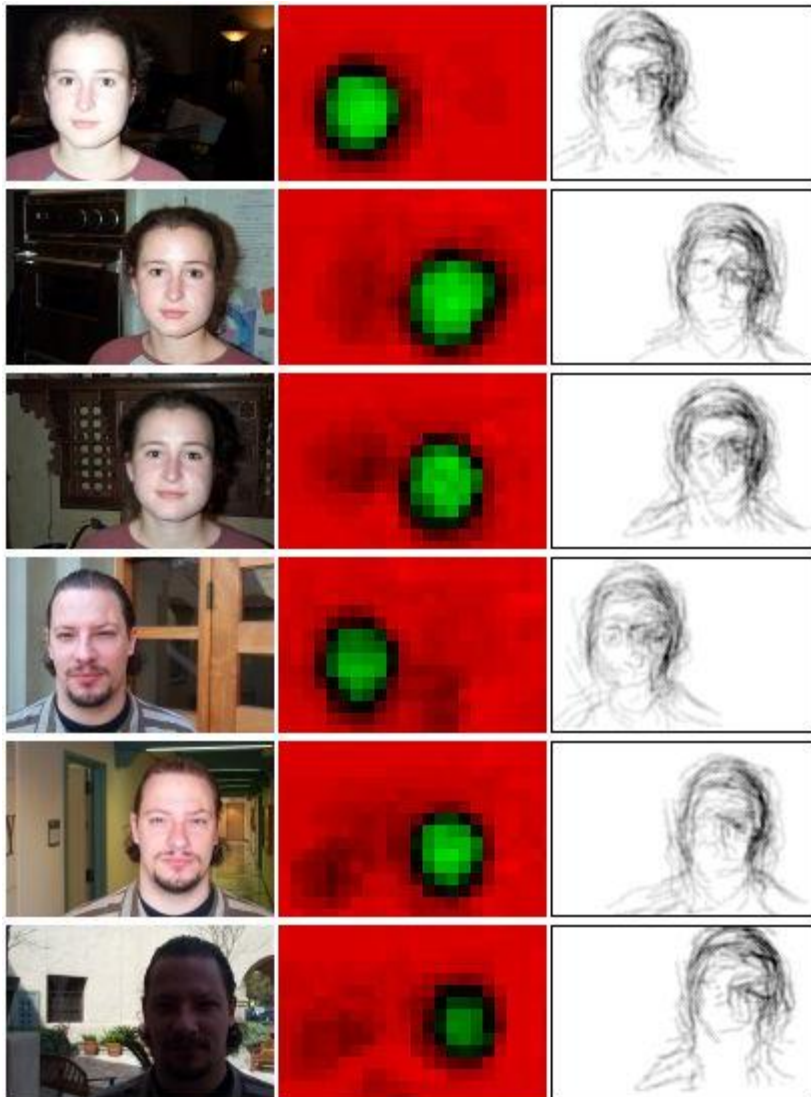
- Given a model, we construct a classification function  $K(\mathbf{c})$ 
  - additive model of feature responses
  - returns confidence value as function of position
    - +ve (green) meaning object present
    - -ve (red) meaning no object
- Evaluate for all centroids in test image gives classification map



# Results

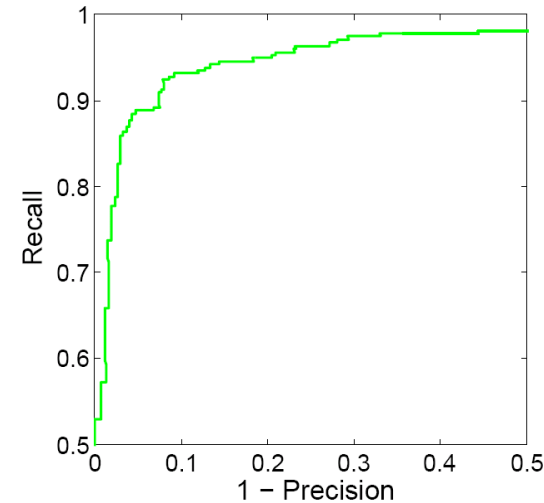


# Results

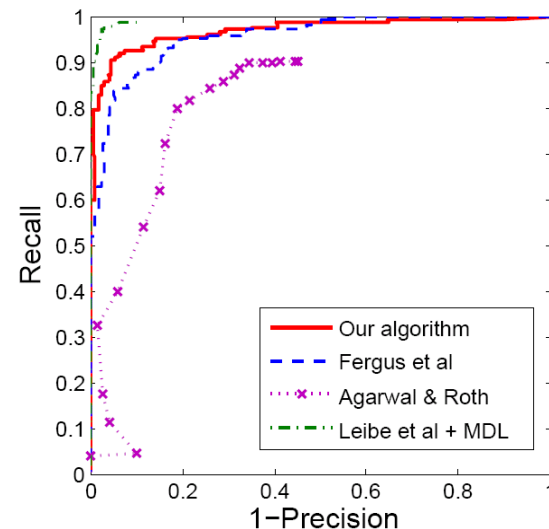


# Results

- Quantification with *recall-precision* curves
  - illustrates trade-off between:
    - correct detection rate
    - proportion of all detections that are correct
  - as a global detection threshold is changed
- A perfect detector would give  $\text{recall}=1$  at  $\text{precision}=1$



Horses



Cars

# Making machines see

---

- 3D shape: making digital copies of sculpture from photographs from multiple viewpoints
- Recognition of a painting/picture from a single photo using a mobile (camera) phone
- Realtime detection of objects: hands, faces and people
- Machine Learning for object categorization and recognition