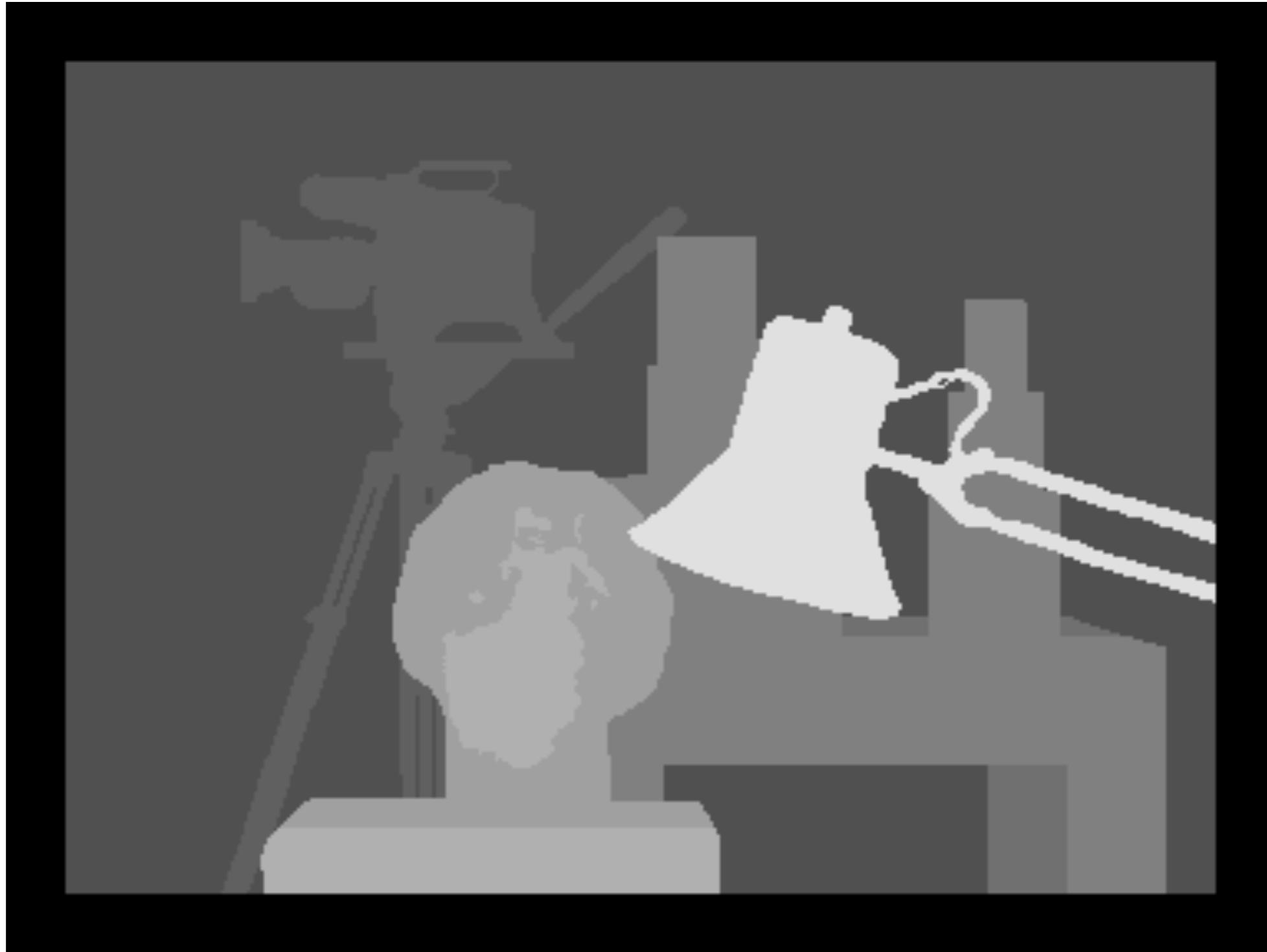


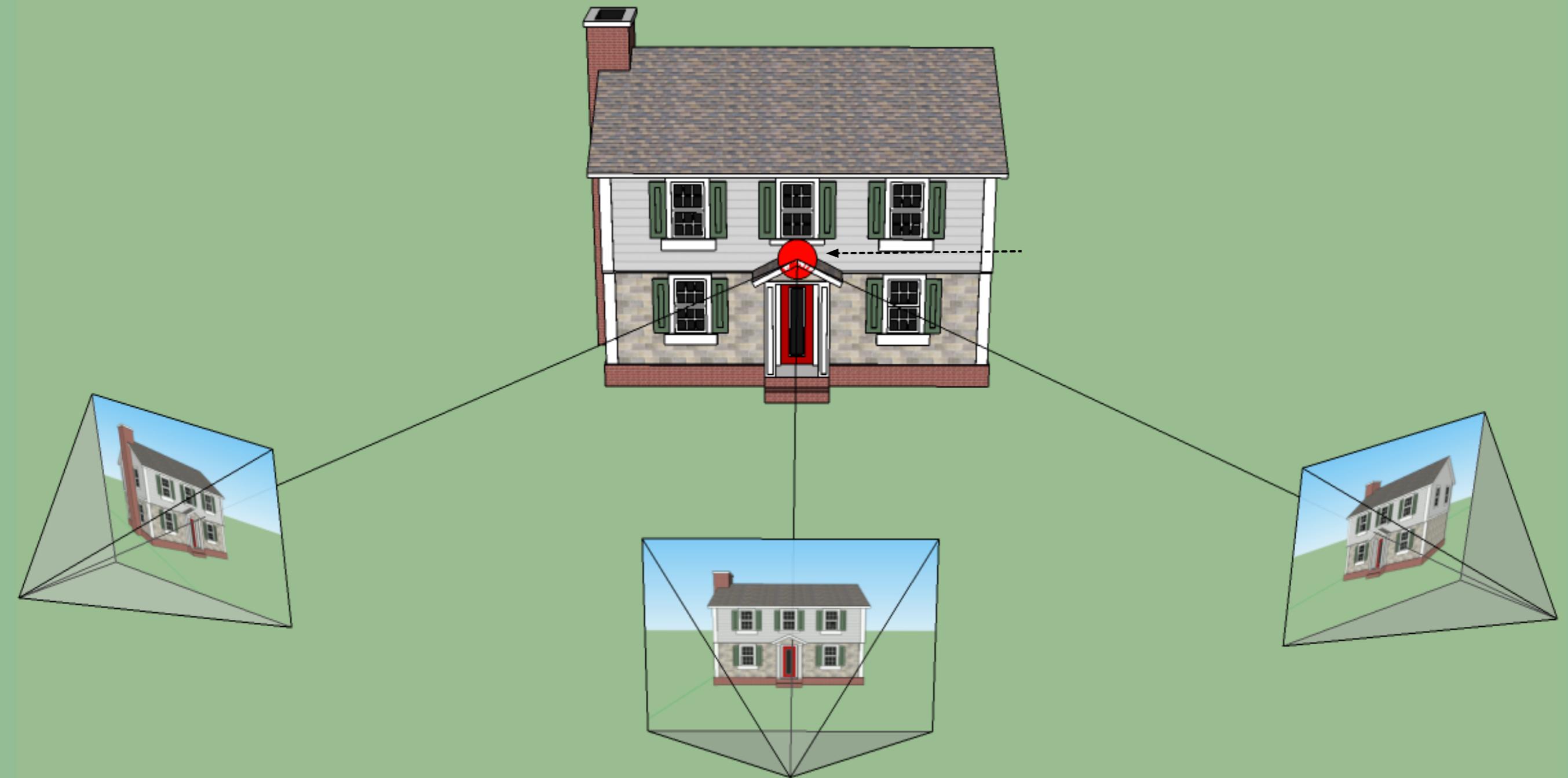
3D Reconstruction with Computer Vision

Meeting 9: Stereo Correspondence



Slides by Kristen Grauman, Richard Szeliski and others
CS 378 Fall 2014, UT Austin, Bryan Klingner, 25 September

Idea: 3D from multiple images of a scene

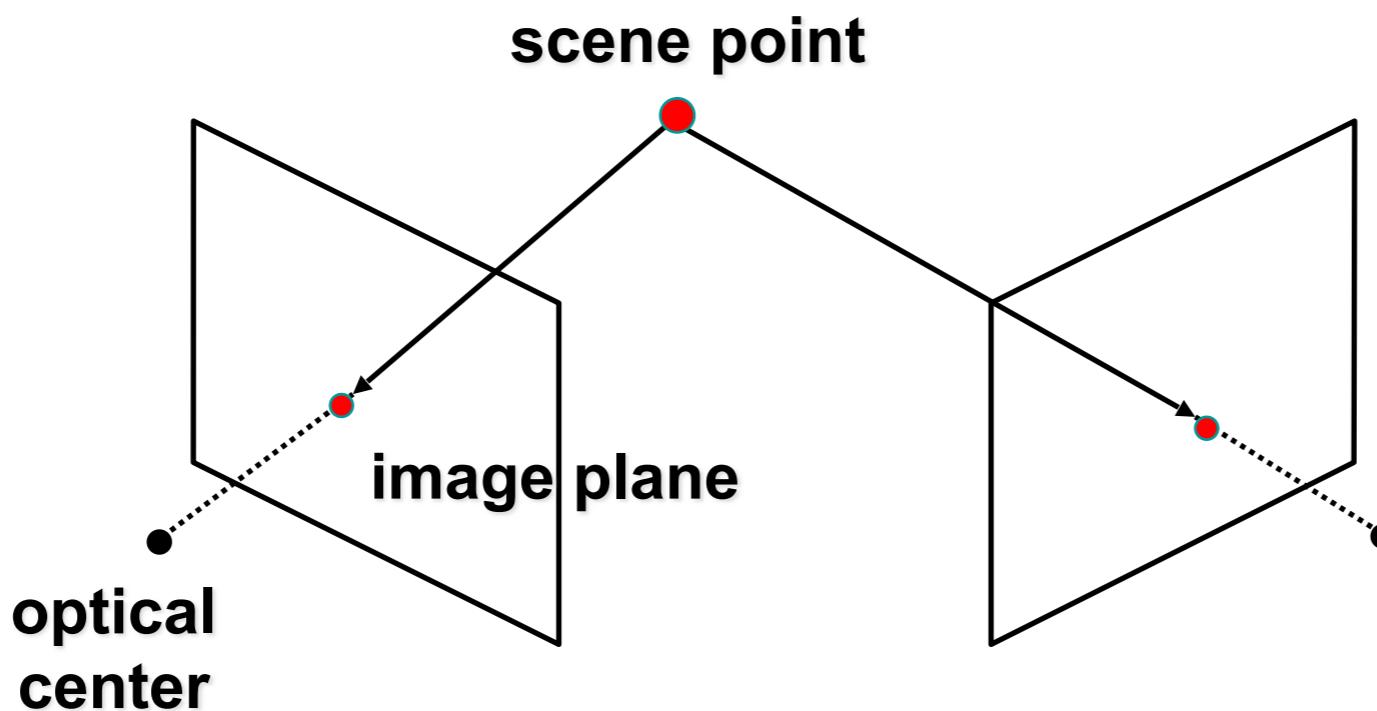


Today

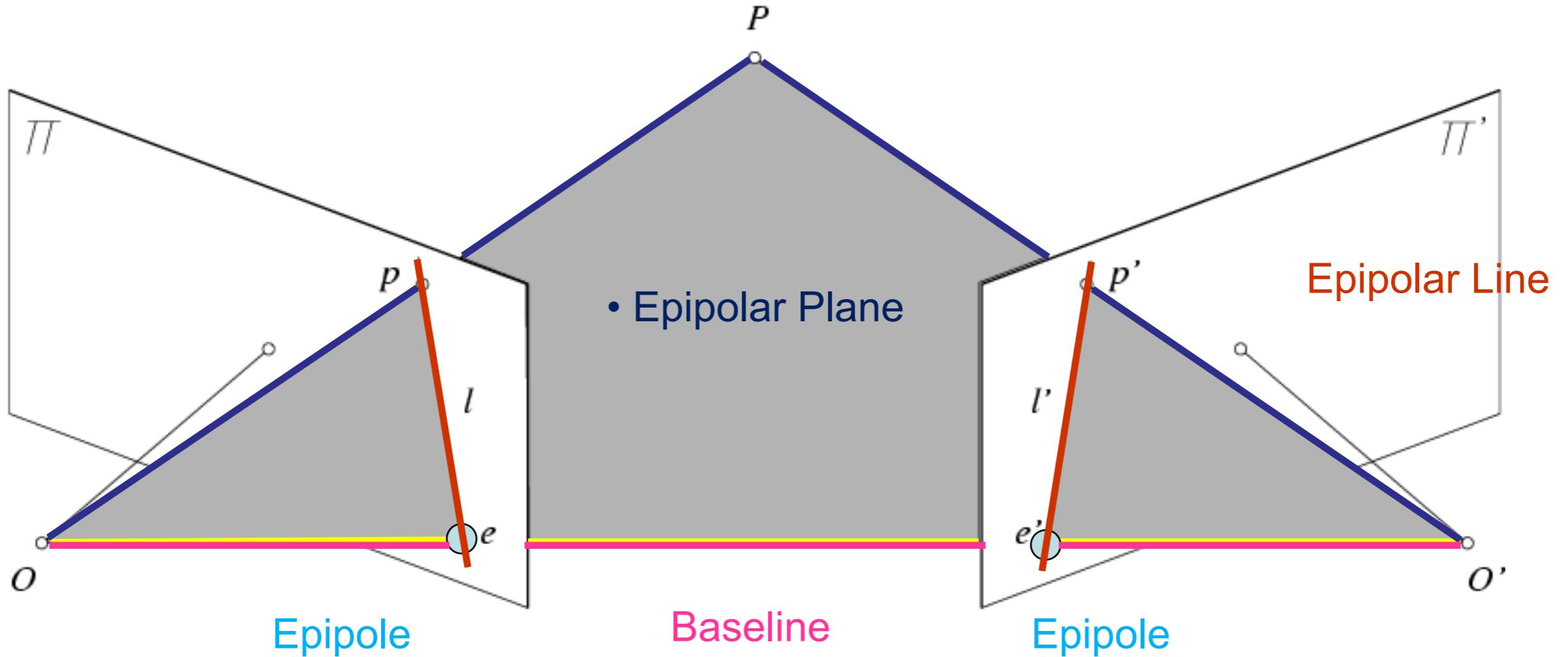
- Recap: epipolar constraint
- Stereo image rectification
- Stereo solutions
 - Computing correspondences
 - Non-geometric stereo constraints
- Calibration
- Example stereo applications

Last time: Estimating depth with stereo

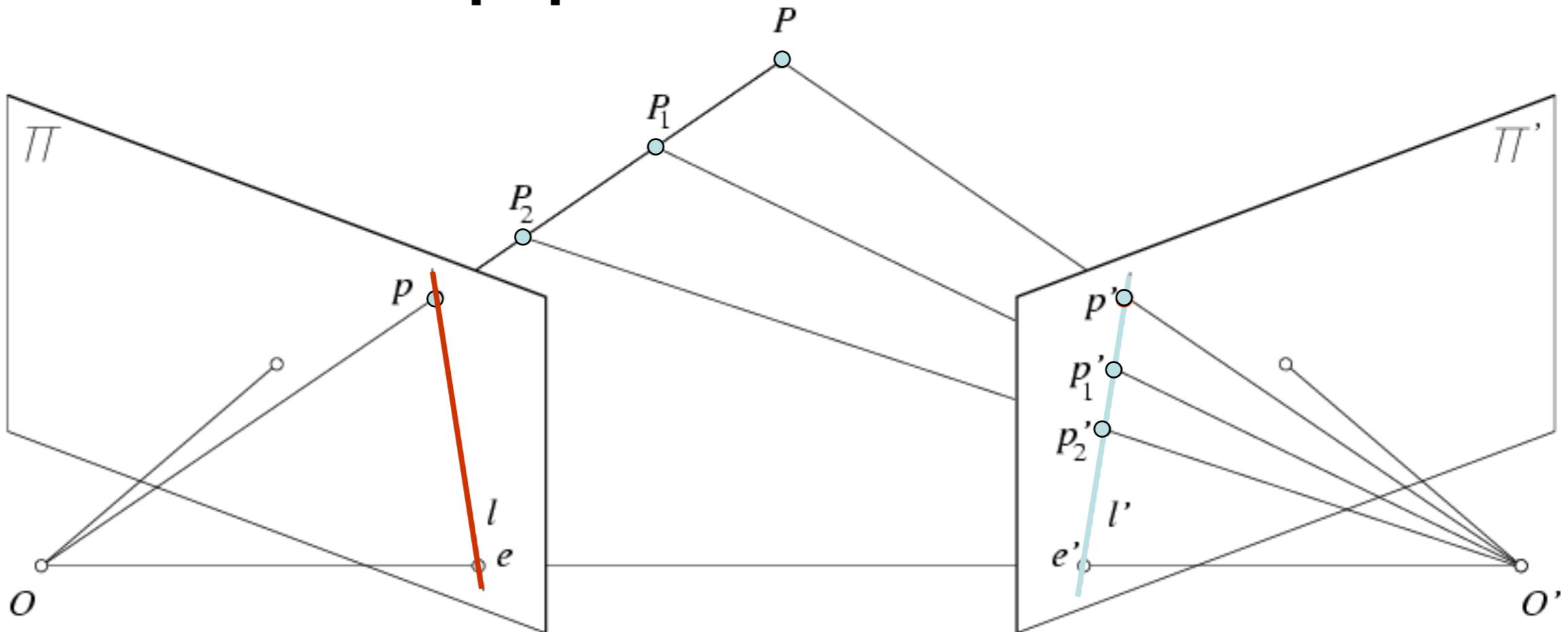
- **Stereo:** shape from “motion” between two views
- We need to consider:
- Info on camera pose (“calibration”)
- Image point correspondences



Last time: Epipolar geometry

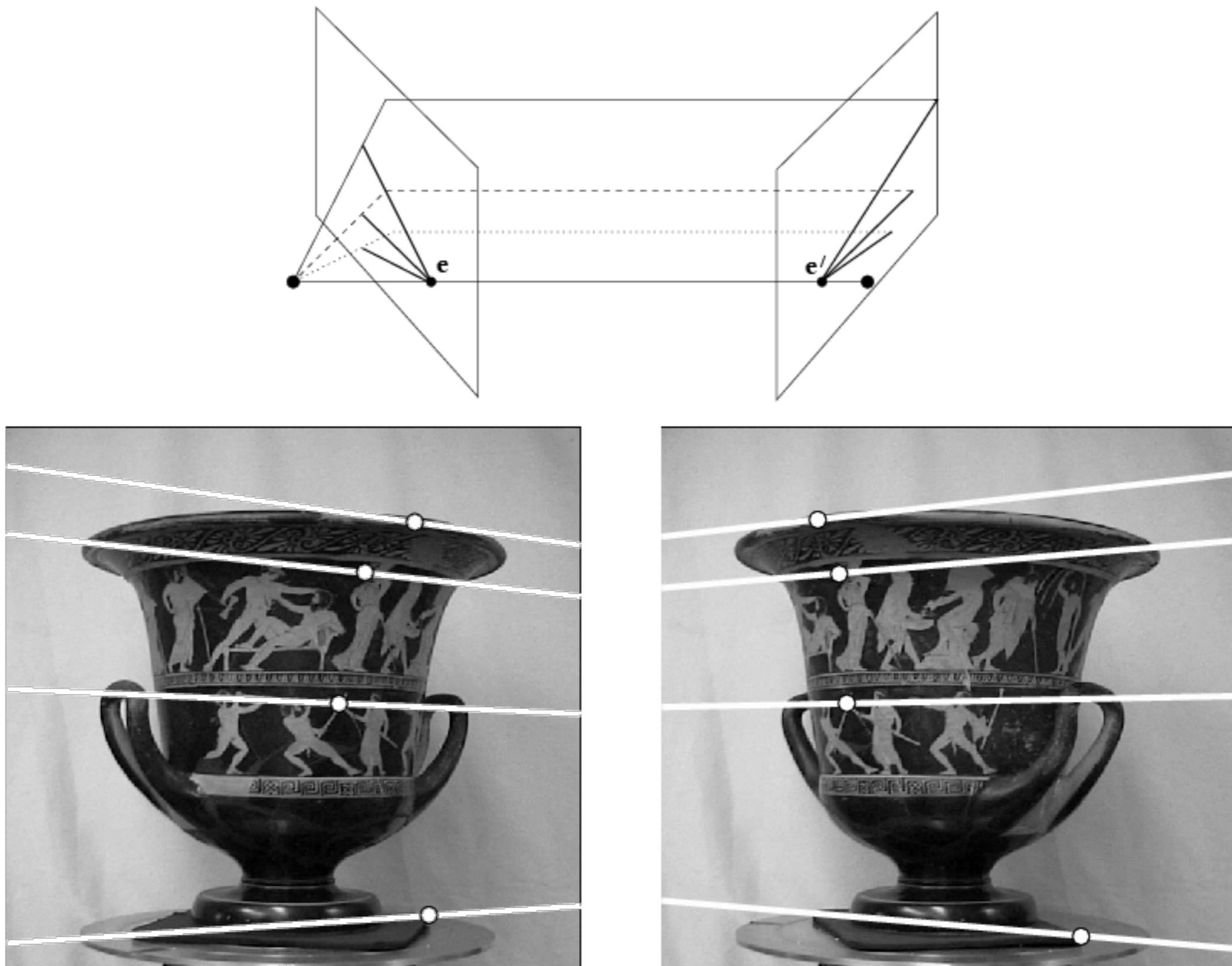


Last time: Epipolar constraint

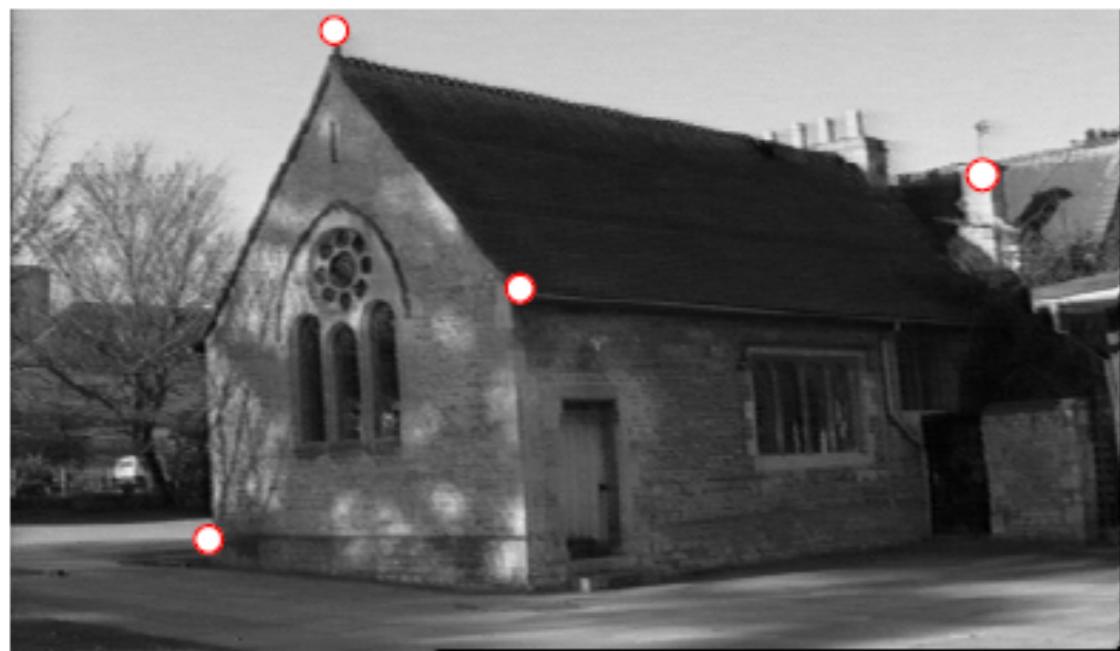
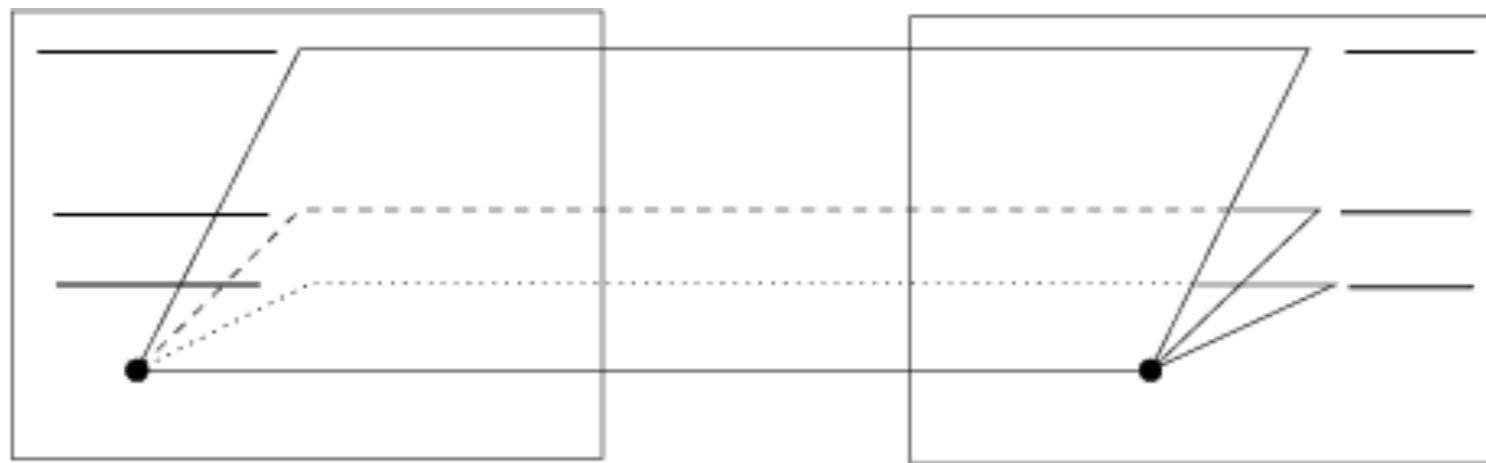


- Potential matches for p have to lie on the corresponding epipolar line l' .
- Potential matches for p' have to lie on the corresponding epipolar line l .

Example: converging cameras



Example: parallel cameras



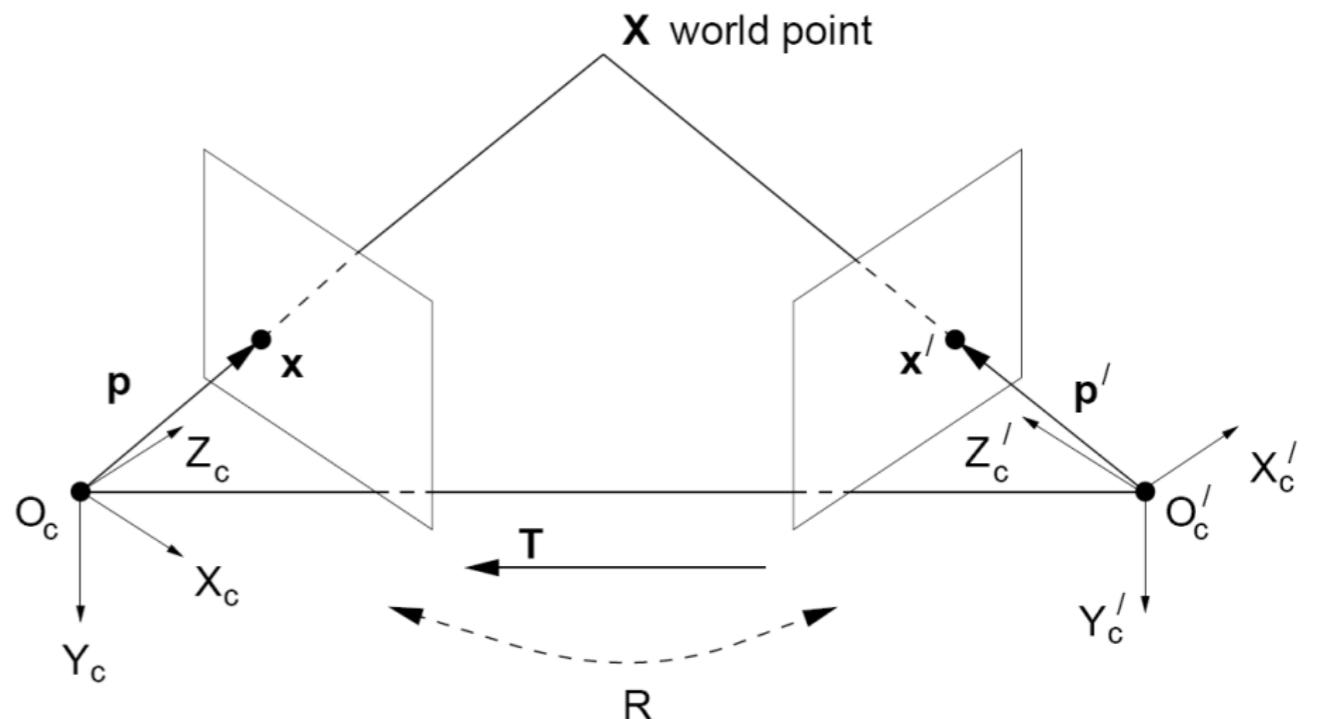
Last time: Essential matrix

$$\mathbf{X}' \cdot (\mathbf{T} \times \mathbf{R}\mathbf{X}) = 0$$

$$\mathbf{X}' \cdot ([\mathbf{T}_x] \mathbf{R}\mathbf{X}) = 0$$

Let $\mathbf{E} = [\mathbf{T}_x] \mathbf{R}$

$$\mathbf{X}'^T \mathbf{E} \mathbf{X} = 0$$



\mathbf{E} is called the **essential matrix**, and it relates corresponding points between both cameras, given the rotation and translation.

If we observe a point in one image, its position in other image is constrained to lie on line defined by above.

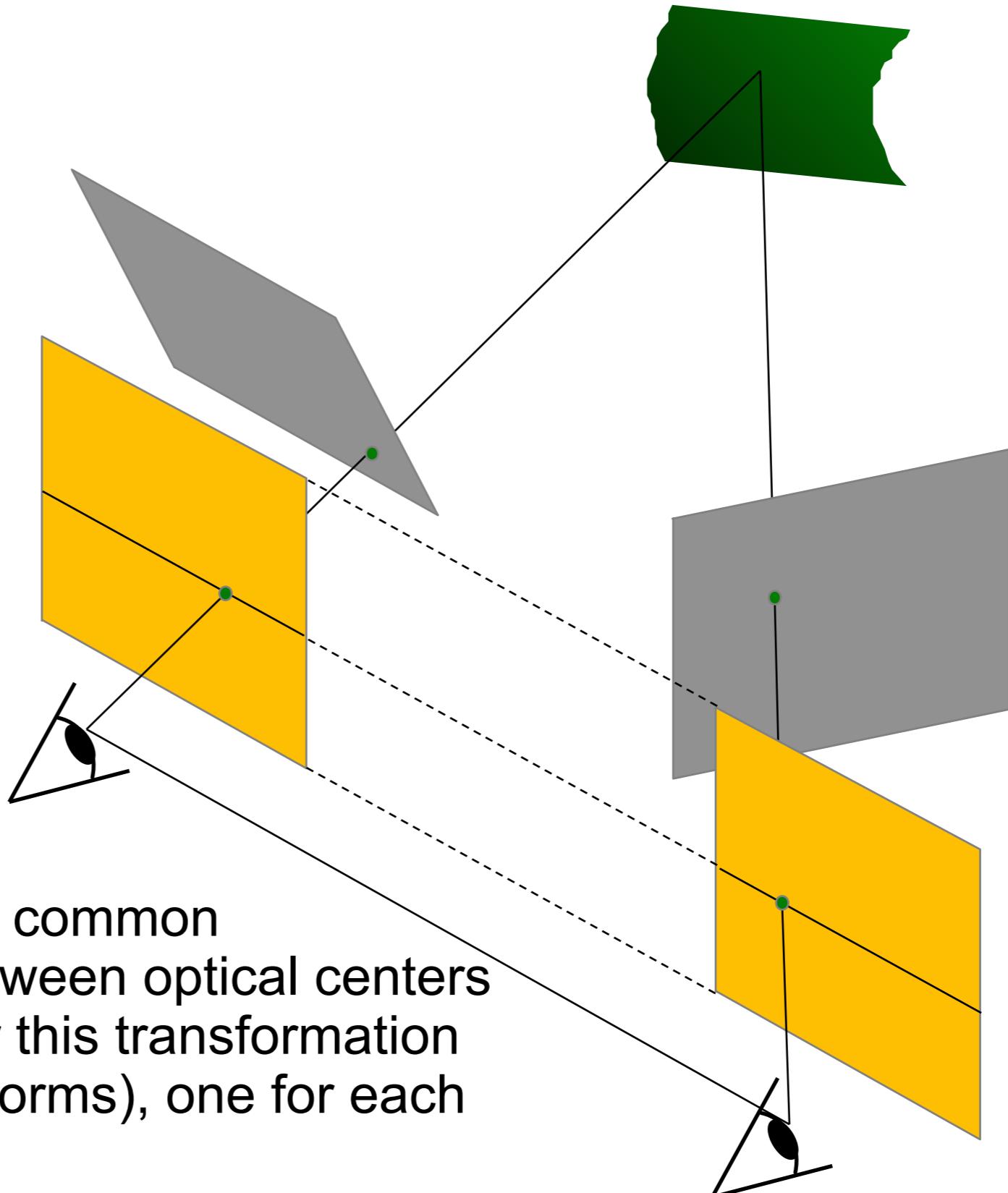
Note: these points are in **camera coordinate systems**.

Today

- Recap: epipolar constraint
- Stereo image rectification
- Stereo solutions
 - Computing correspondences
 - Non-geometric stereo constraints
- Calibration
- Example stereo applications

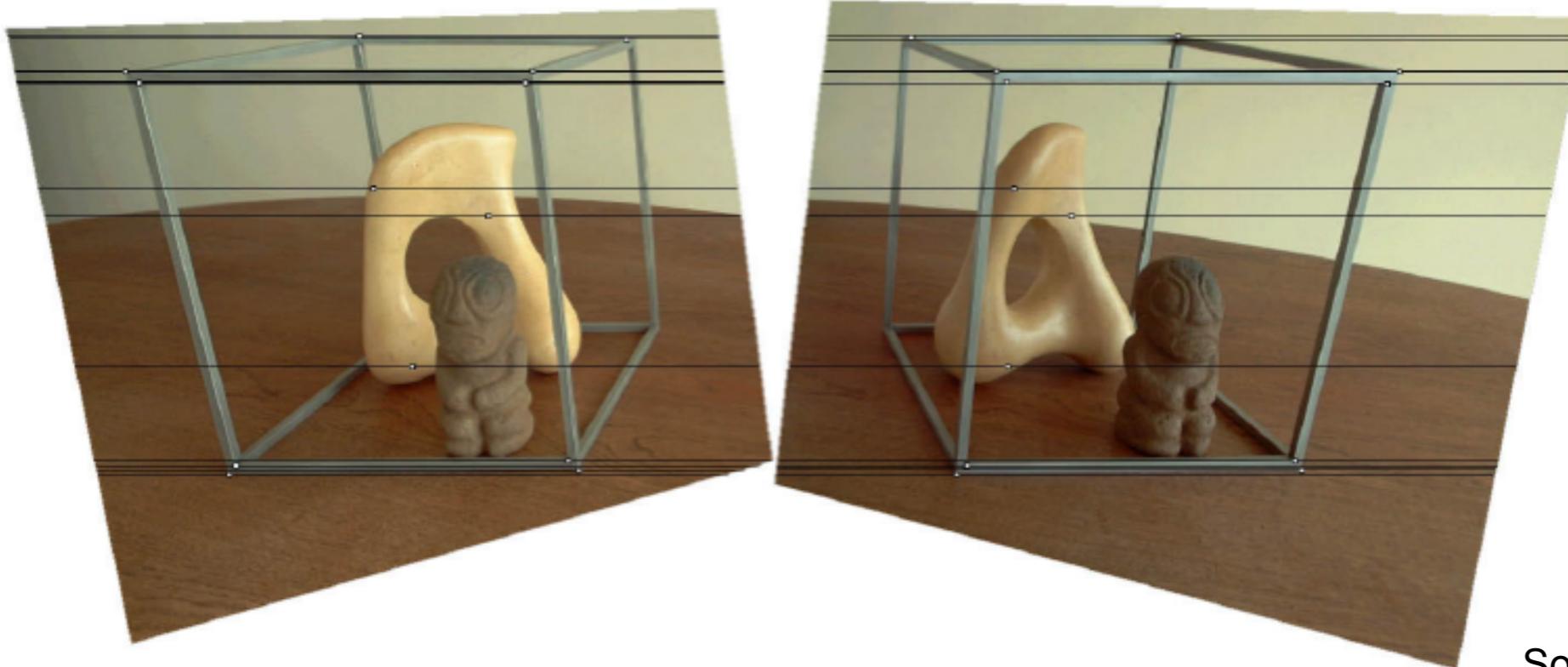
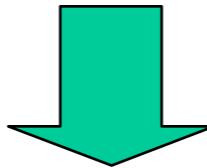
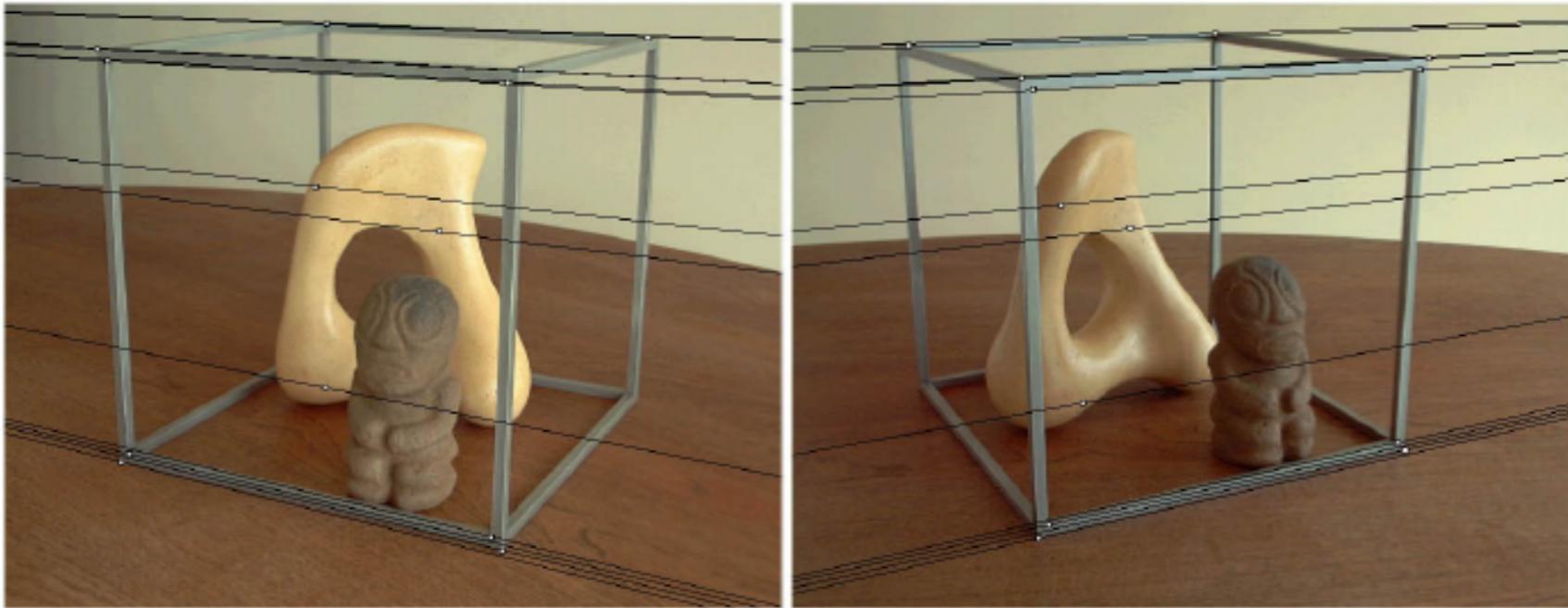
Stereo image rectification

In practice, it is convenient if image scanlines (rows) are the epipolar lines.



reproject image planes onto a common plane parallel to the line between optical centers
pixel motion is horizontal after this transformation
two homographies (3x3 transforms), one for each input image reprojeciton

Stereo image rectification: example

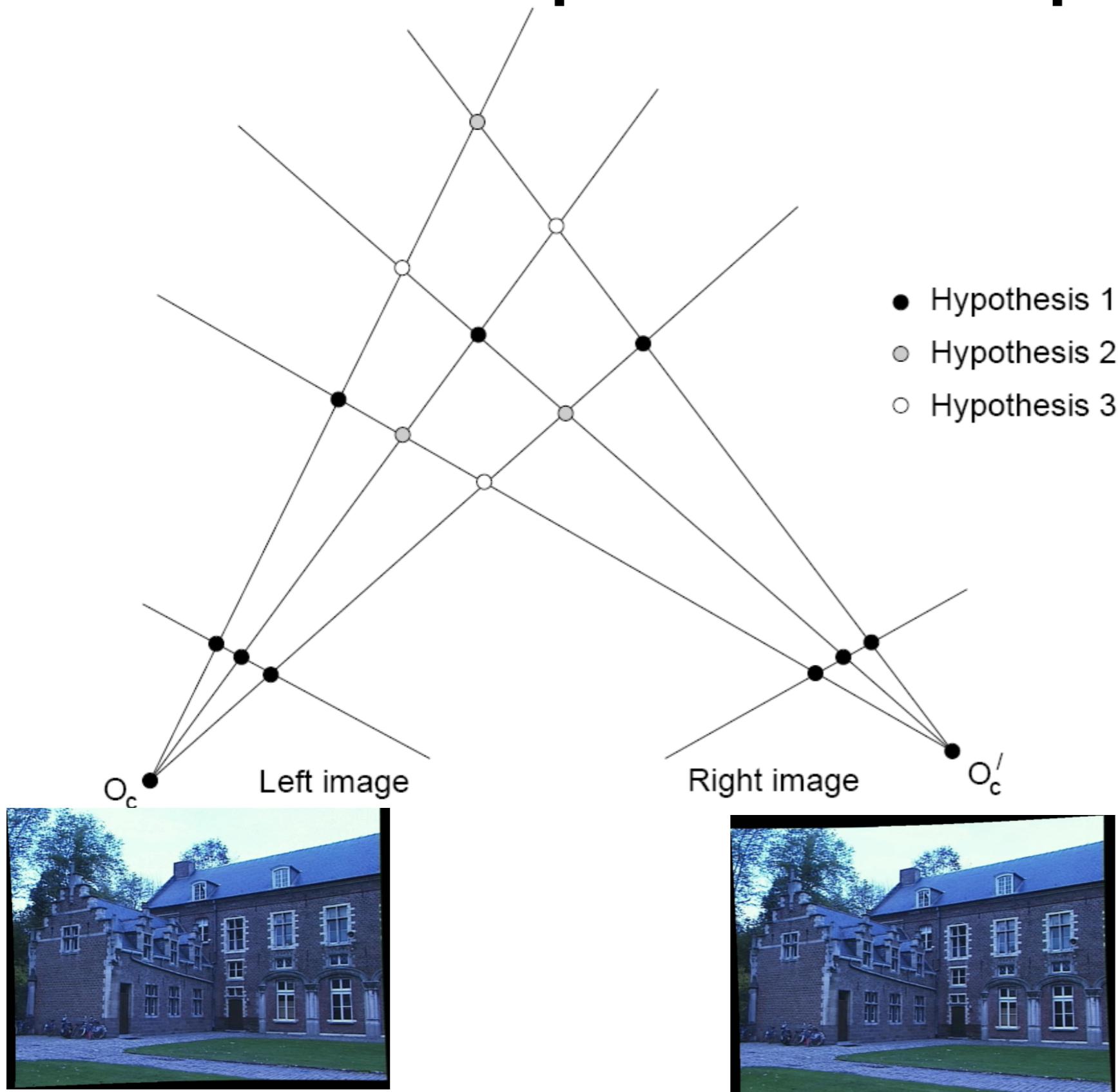


Source: Alyosha Efros

Today

- Recap: epipolar constraint
- Stereo image rectification
- Stereo solutions
 - Computing correspondences
 - Non-geometric stereo constraints
- Calibration
- Example stereo applications

Correspondence problem

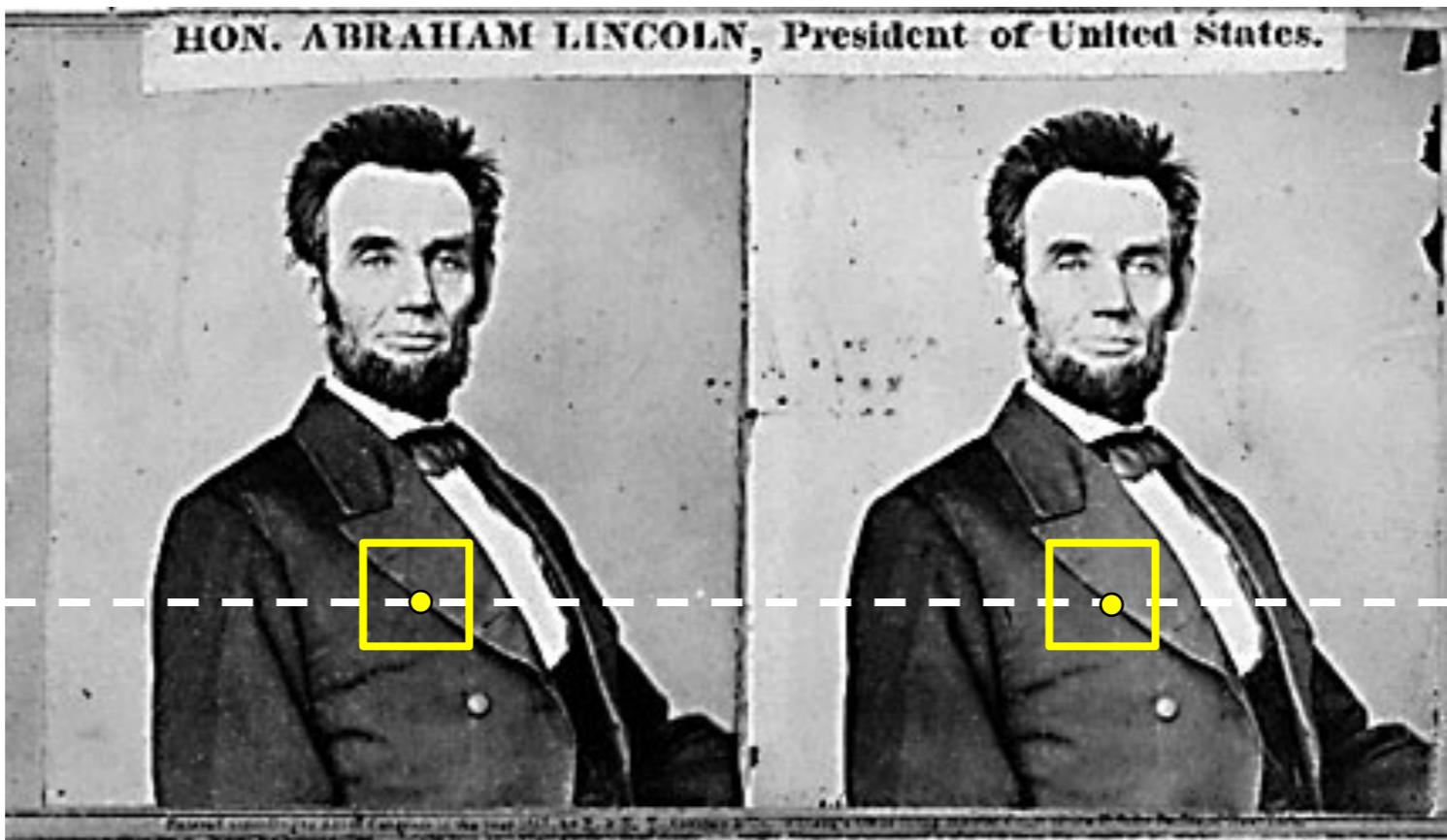


Multiple match hypotheses satisfy epipolar constraint, but which is correct?

Correspondence problem

- Beyond the hard constraint of epipolar geometry, there are “soft” constraints to help identify corresponding points
 - Similarity
 - Uniqueness
 - Ordering
 - Disparity gradient
- To find matches in the image pair, we will assume
 - Most scene points visible from both views
 - Image regions for the matches are similar in appearance

Dense correspondence search

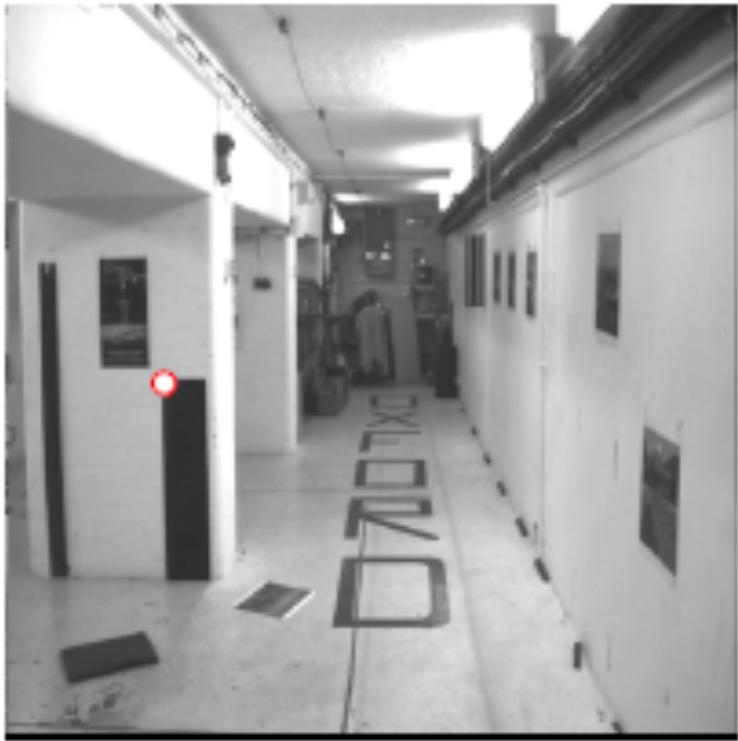


For each epipolar line

For each pixel / window in the left image

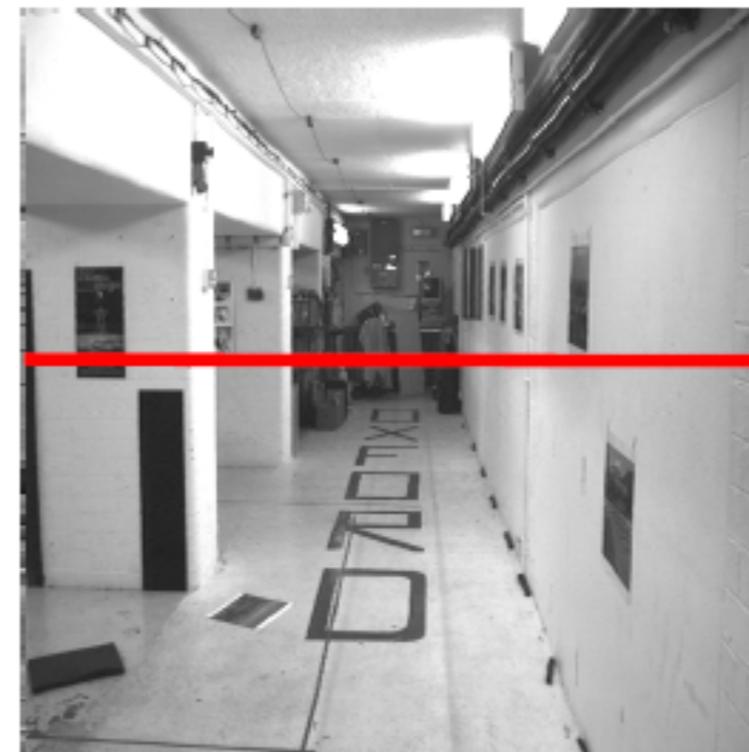
- compare with every pixel / window on same epipolar line in right image
- pick position with minimum match cost (e.g., SSD, correlation)

Correspondence problem

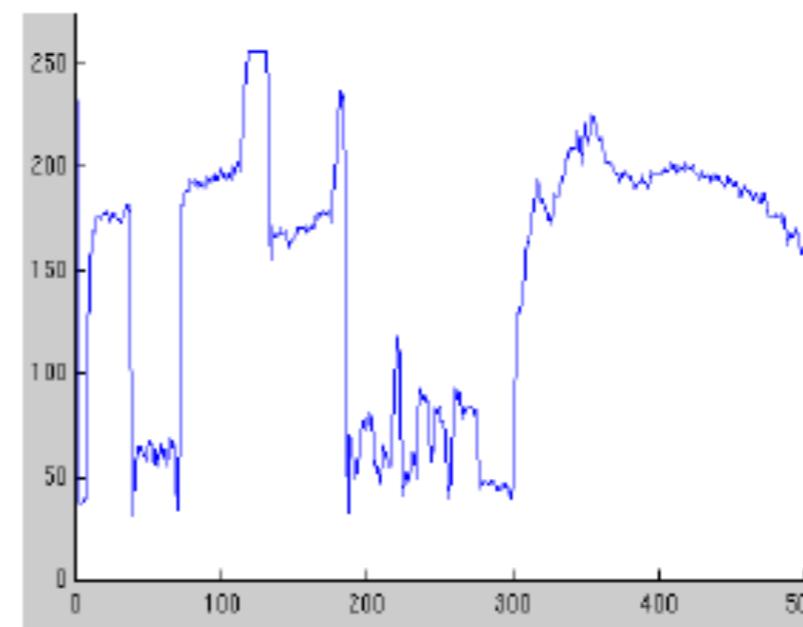
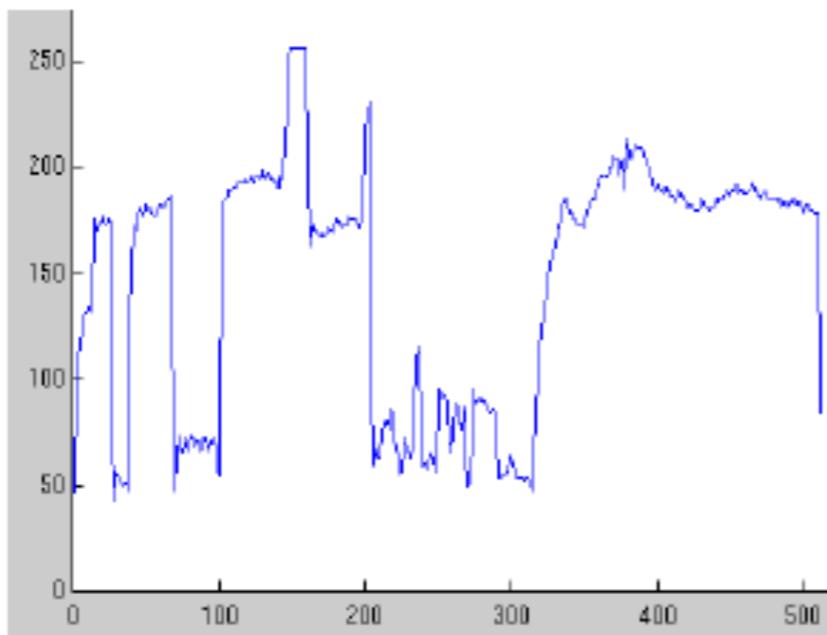


Parallel camera example: epipolar lines are corresponding image scanlines

Correspondence problem

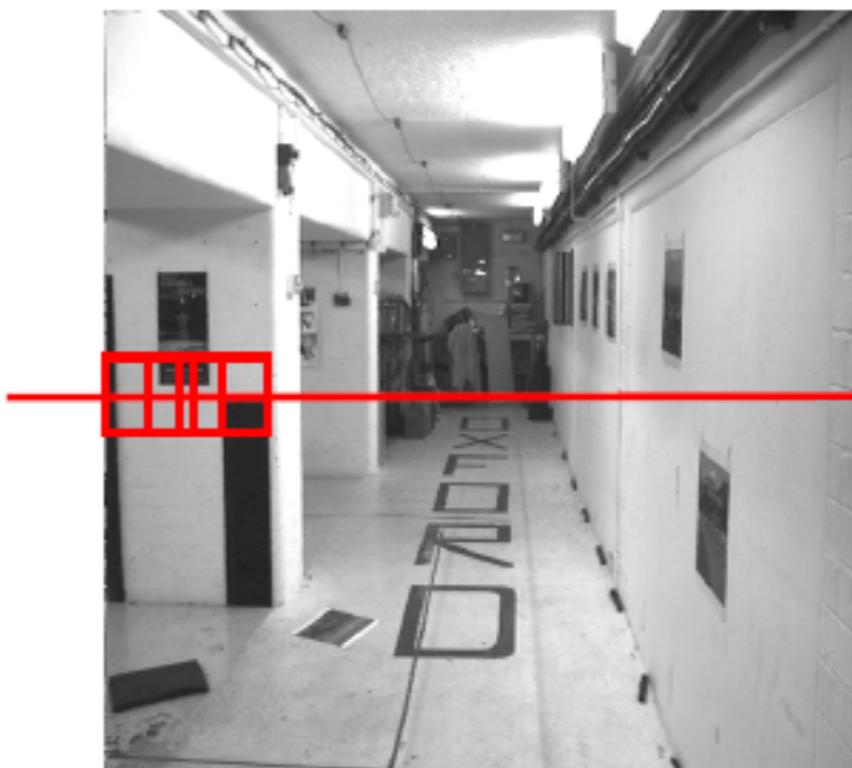


Intensity profiles



- Clear correspondence between intensities, but also noise and ambiguity

Correspondence problem



epipolar
line

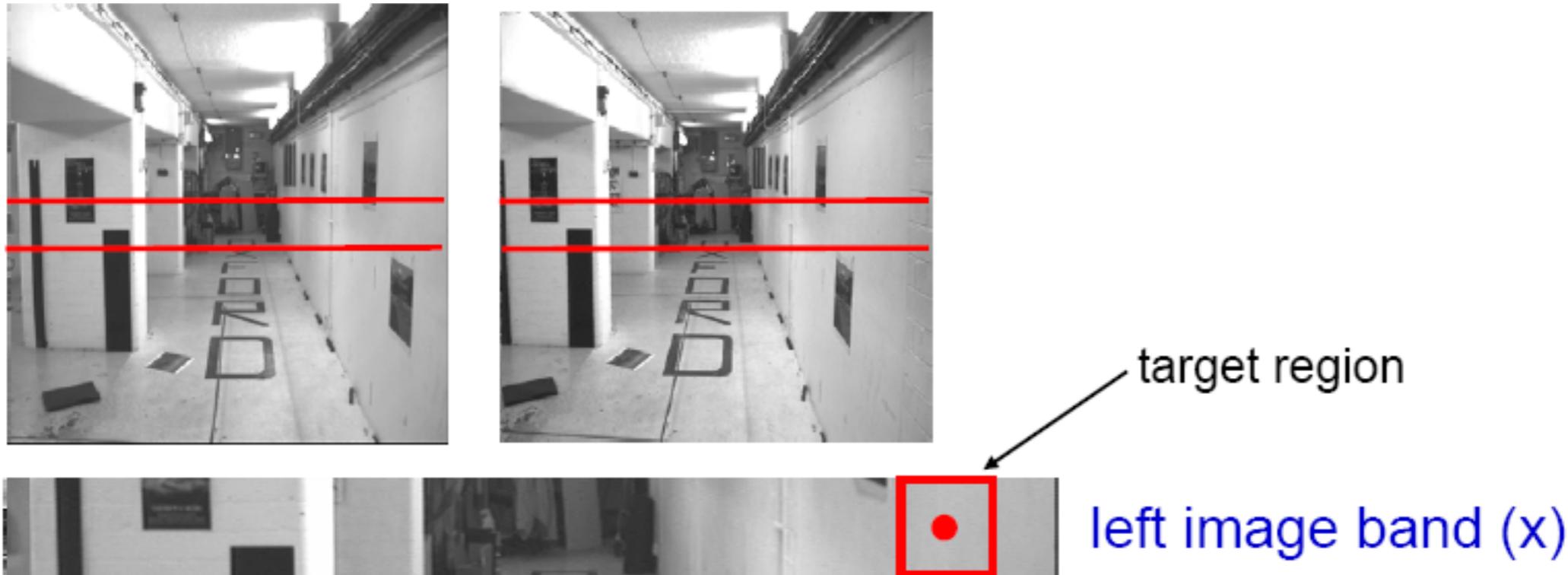
Neighborhoods of corresponding points are similar in intensity patterns.

Correlation-based window matching



left image band (x)

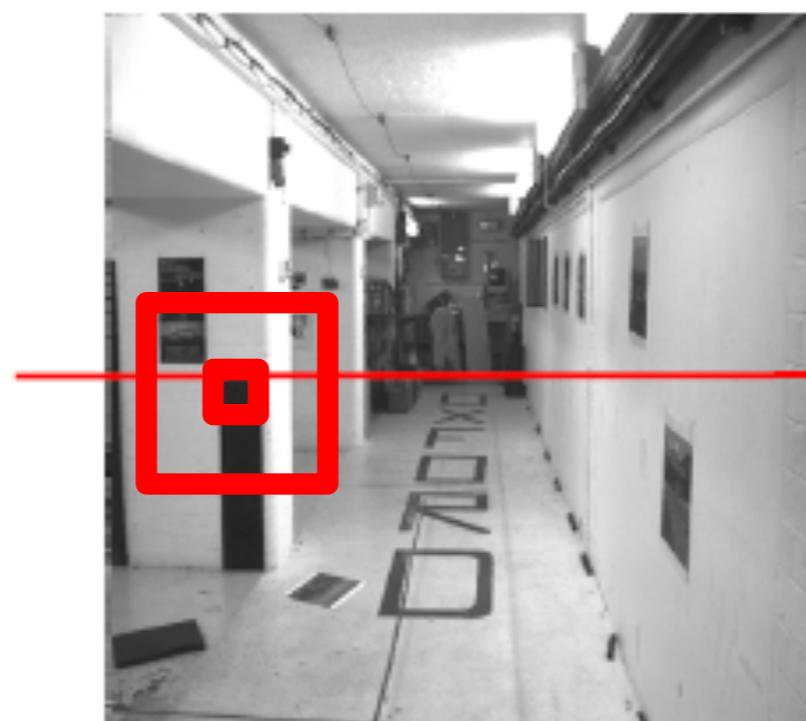
Textureless regions



e

;

Effect of window size

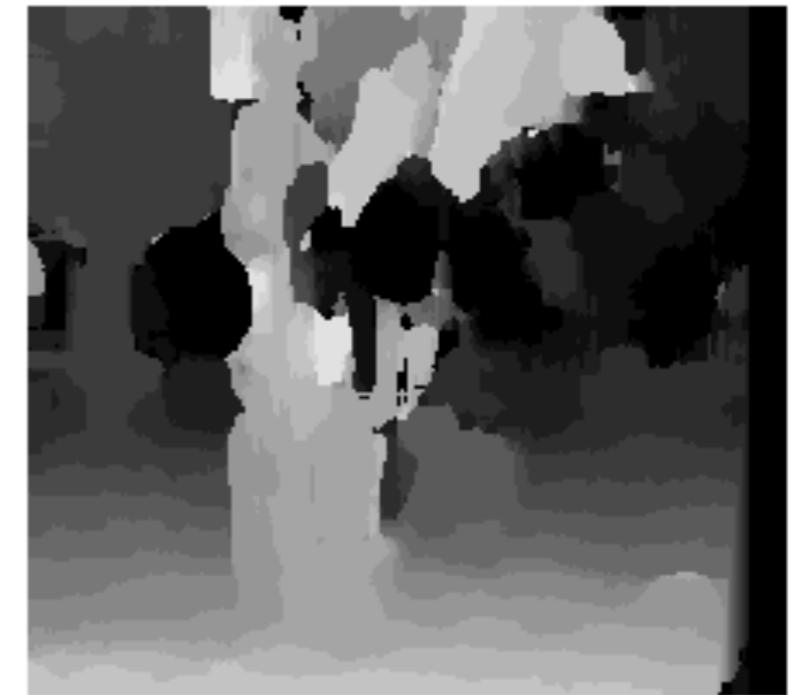


epipolar
line

Effect of window size



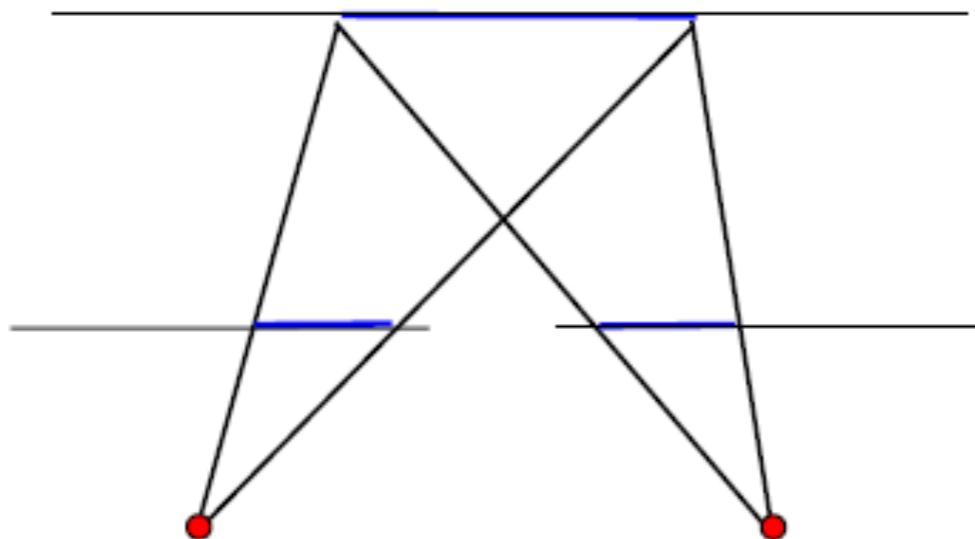
$W = 3$



$W = 20$

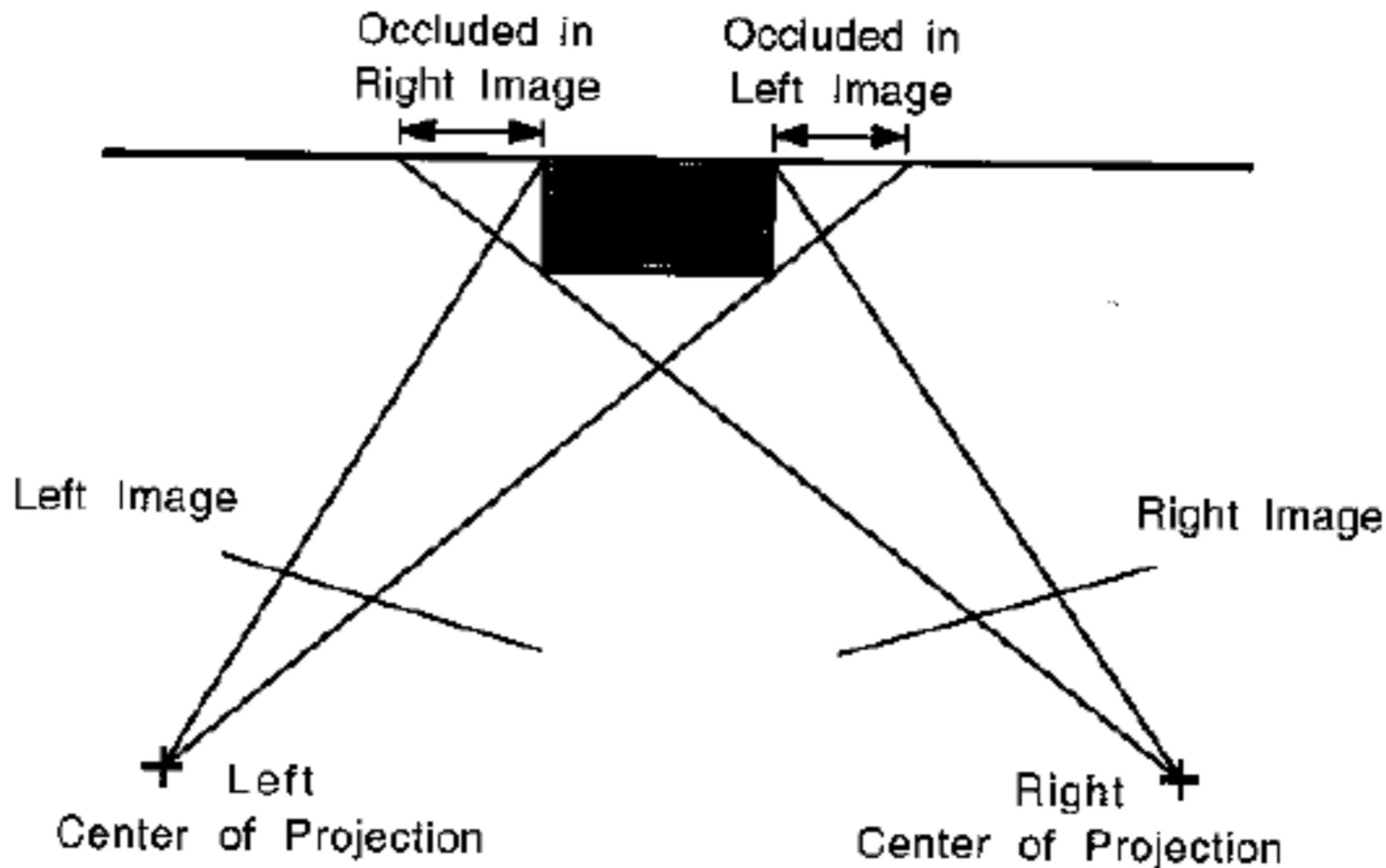
Want window large enough to have sufficient intensity variation, yet small enough to contain only pixels with about the same disparity.

Foreshortening effects

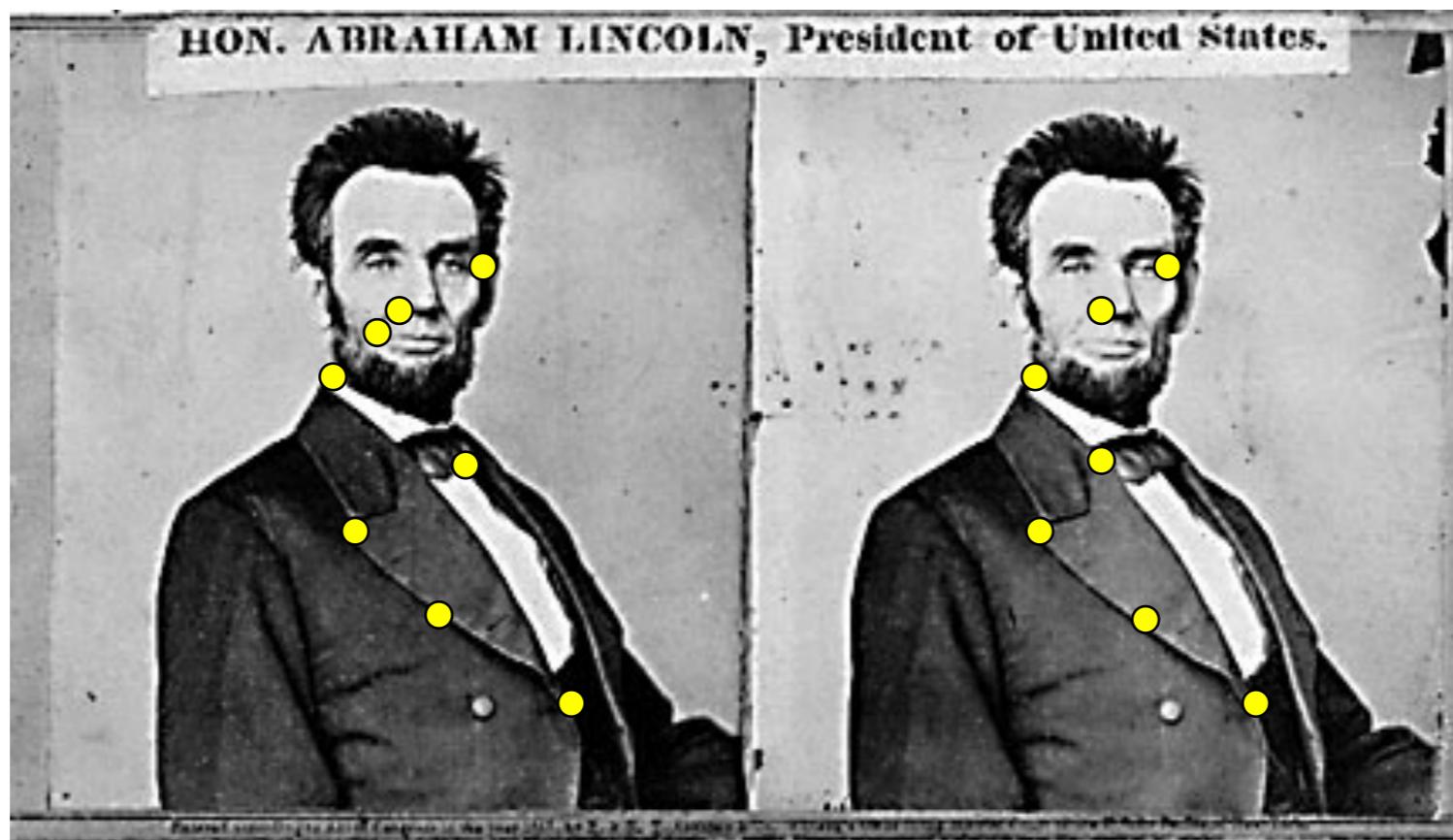


fronto-parallel surface
imaged length the same

Occlusion



Sparse correspondence search



- Restrict search to sparse set of **detected features** (e.g., corners)
- Rather than pixel values (or lists of pixel values) use *feature descriptor* and an associated *feature distance*
- Still narrow search further by epipolar geometry

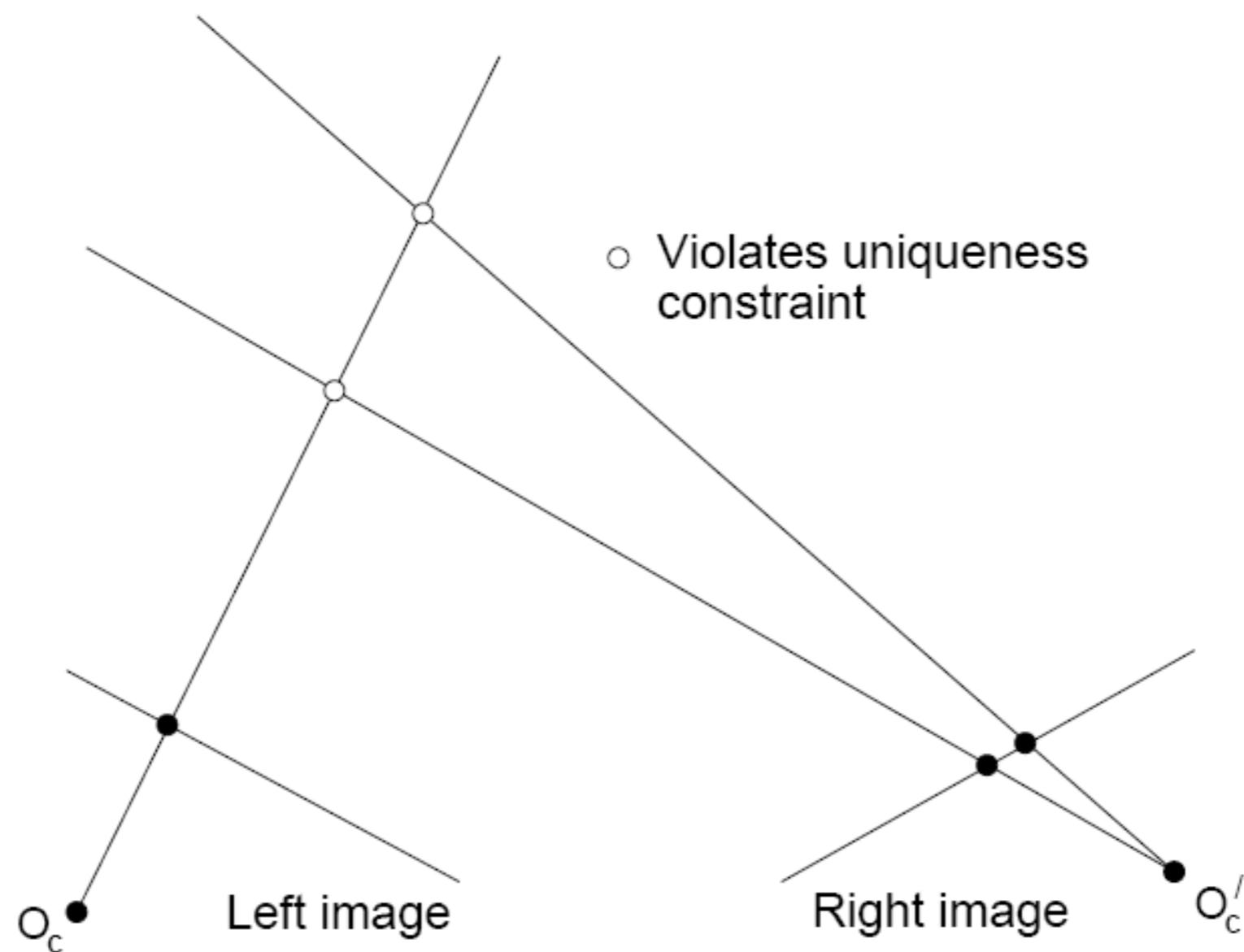
Tradeoffs between dense and sparse search?

Correspondence problem

- Beyond the hard constraint of epipolar geometry, there are “soft” constraints to help identify corresponding points
 - Similarity
 - Uniqueness
 - Disparity gradient
 - Ordering

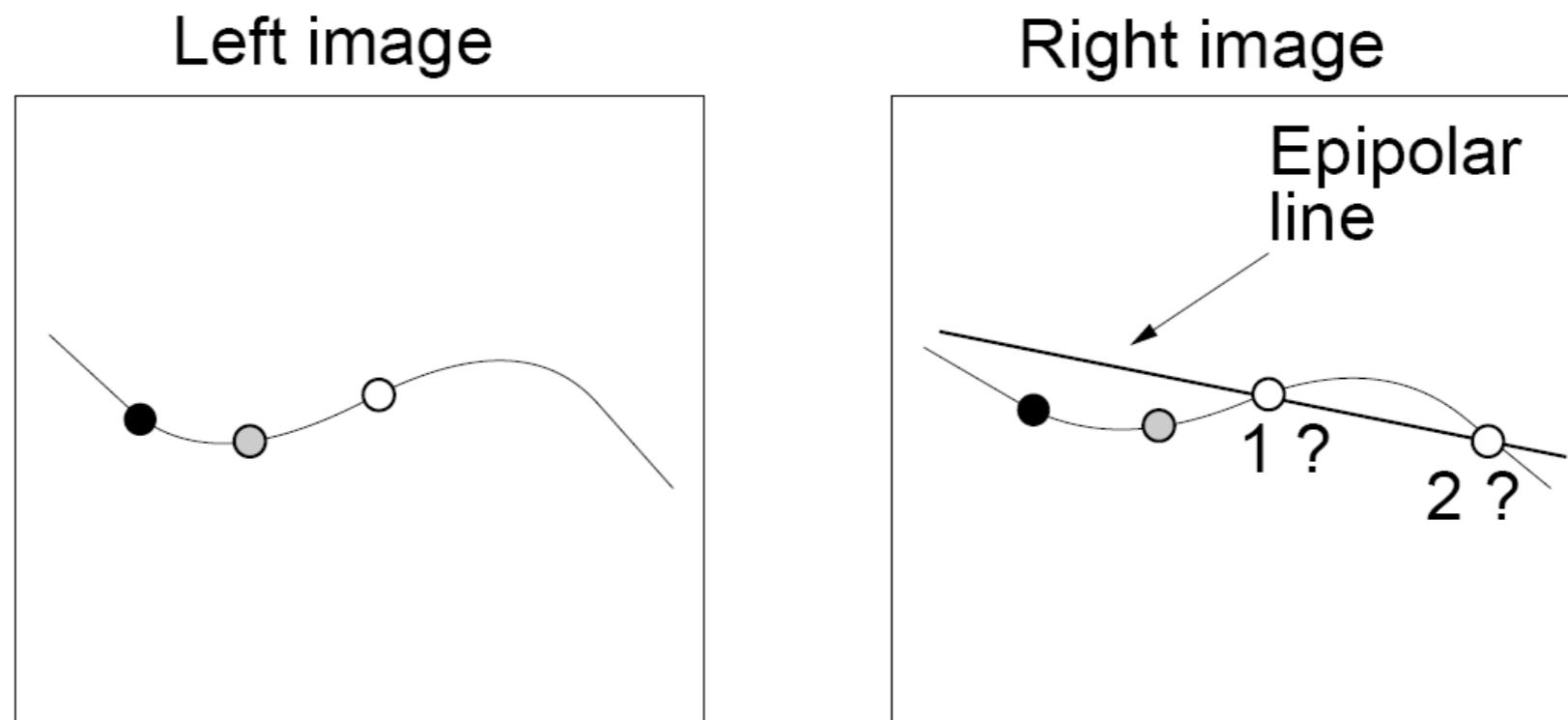
Uniqueness constraint

- Up to one match in right image for every point in left image



Disparity gradient constraint

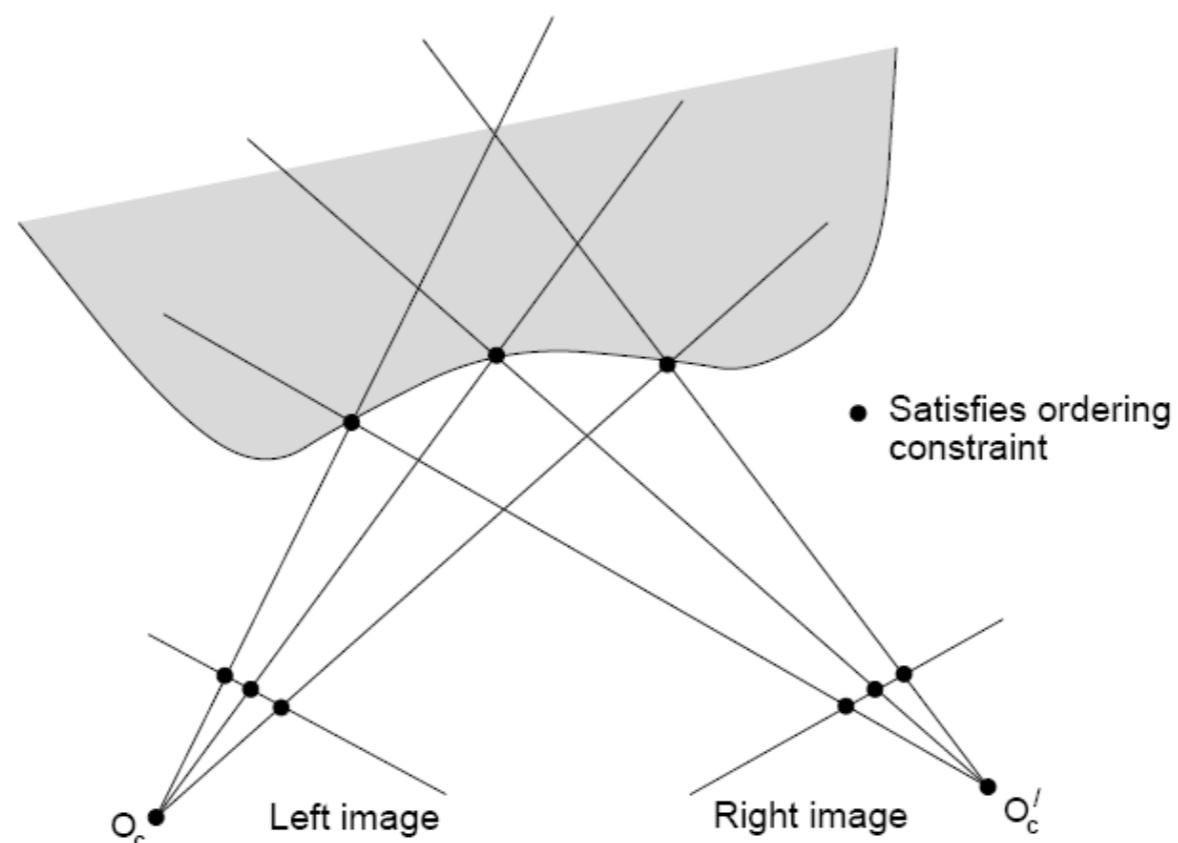
- Assume piecewise continuous surface, so want disparity estimates to be locally smooth



Given matches ● and ○, point ○ in the left image must match point 1 in the right image. Point 2 would exceed the disparity gradient limit.

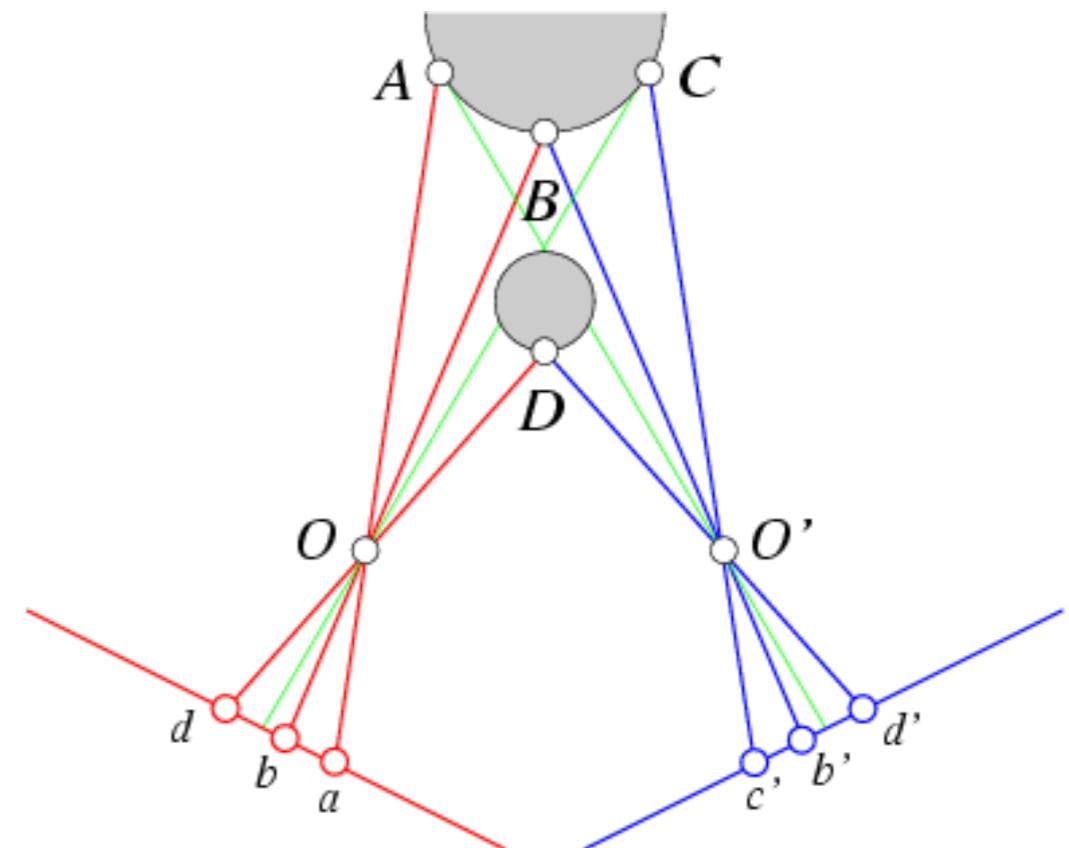
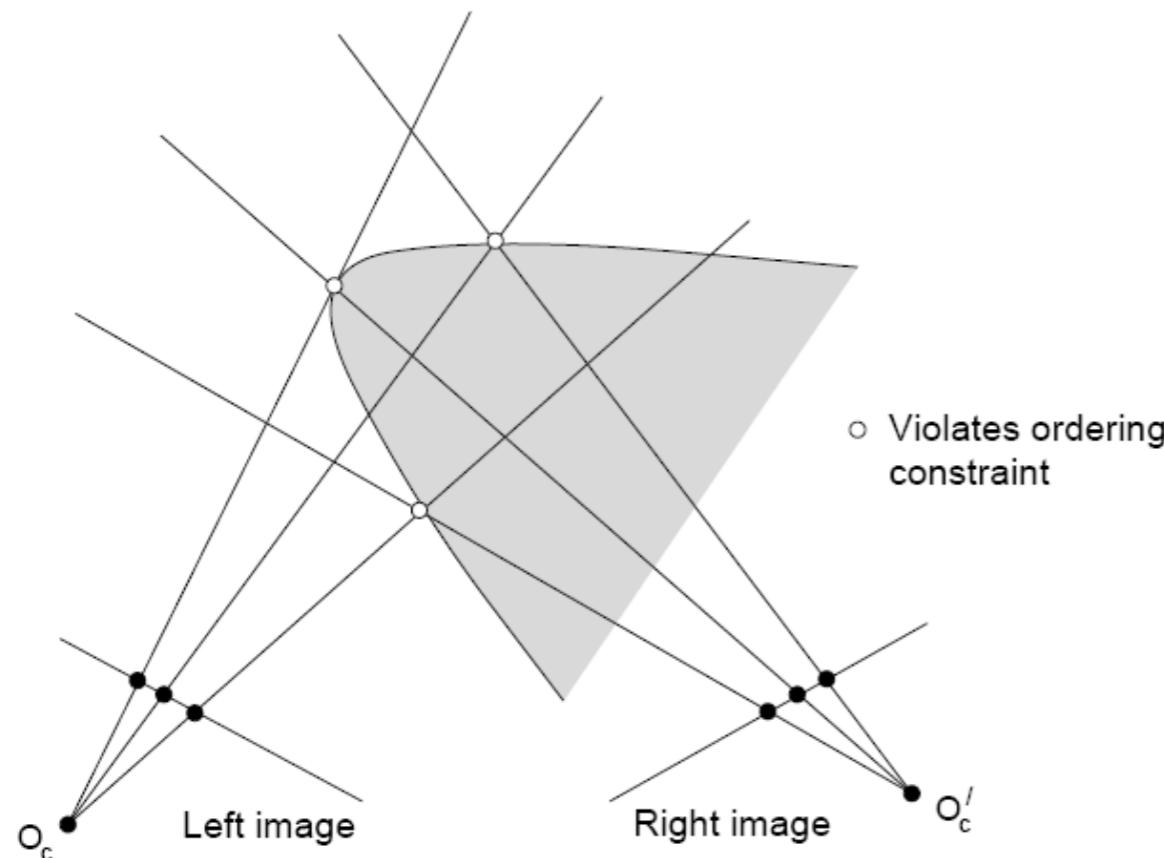
Ordering constraint

- Points on **same surface** (opaque object) will be in same order in both views



Ordering constraint

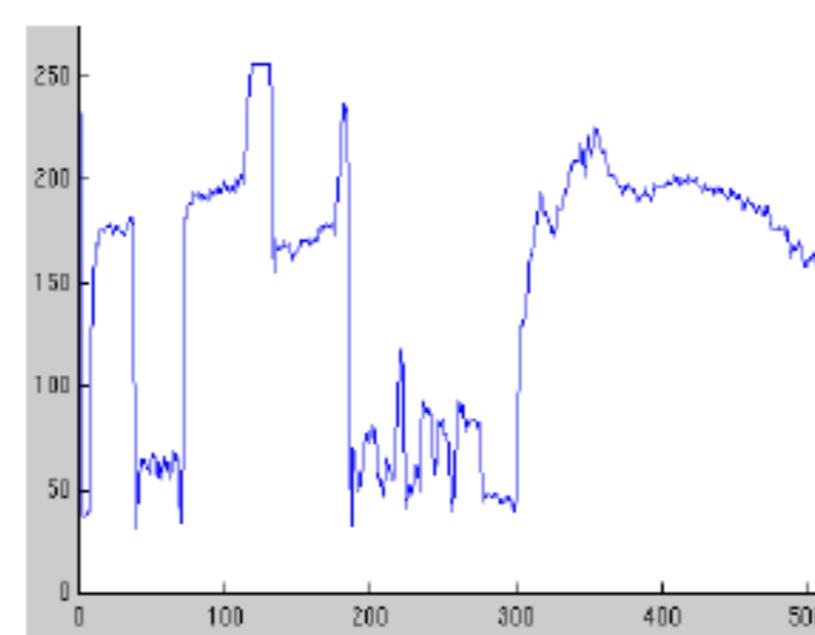
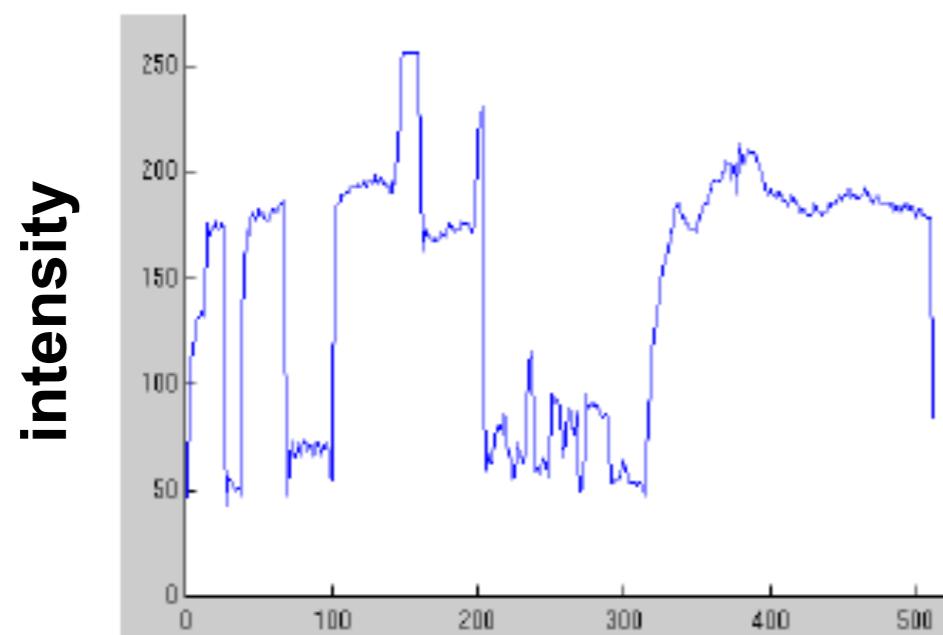
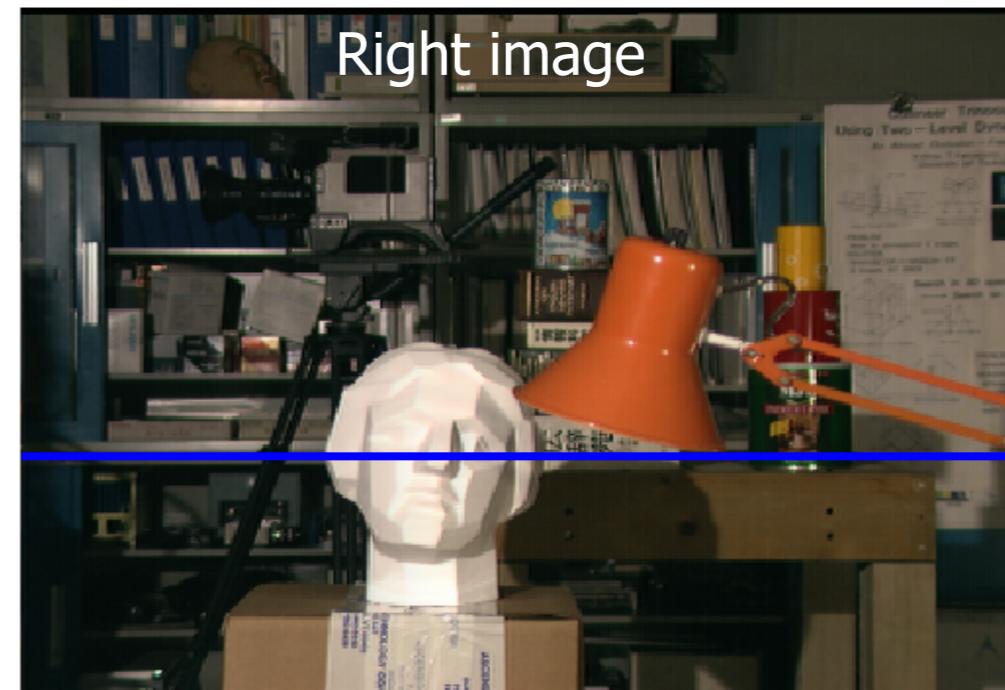
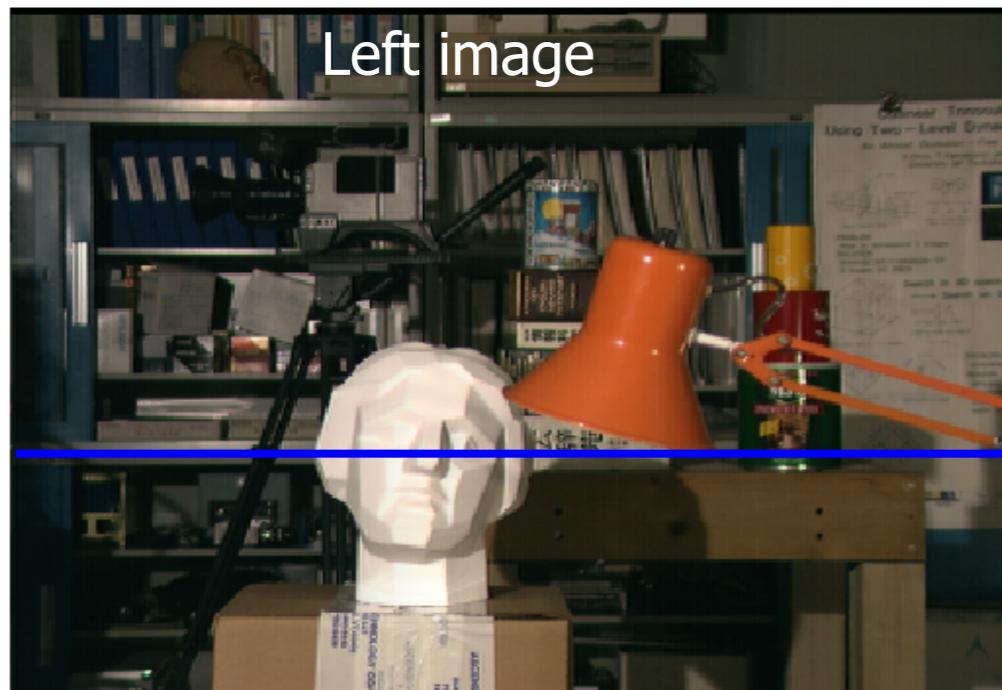
- Won't always hold, e.g. consider transparent object, or an occluding surface



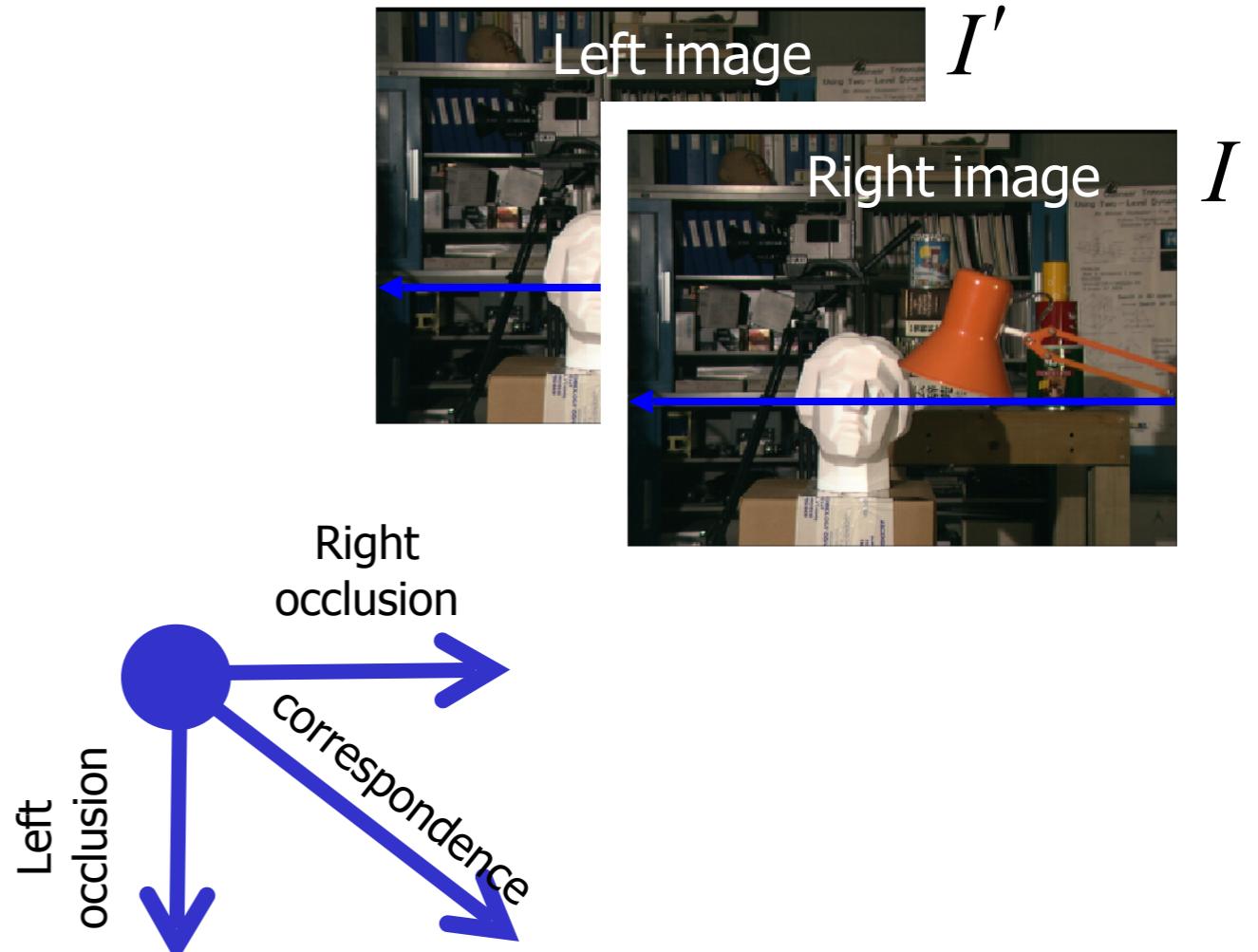
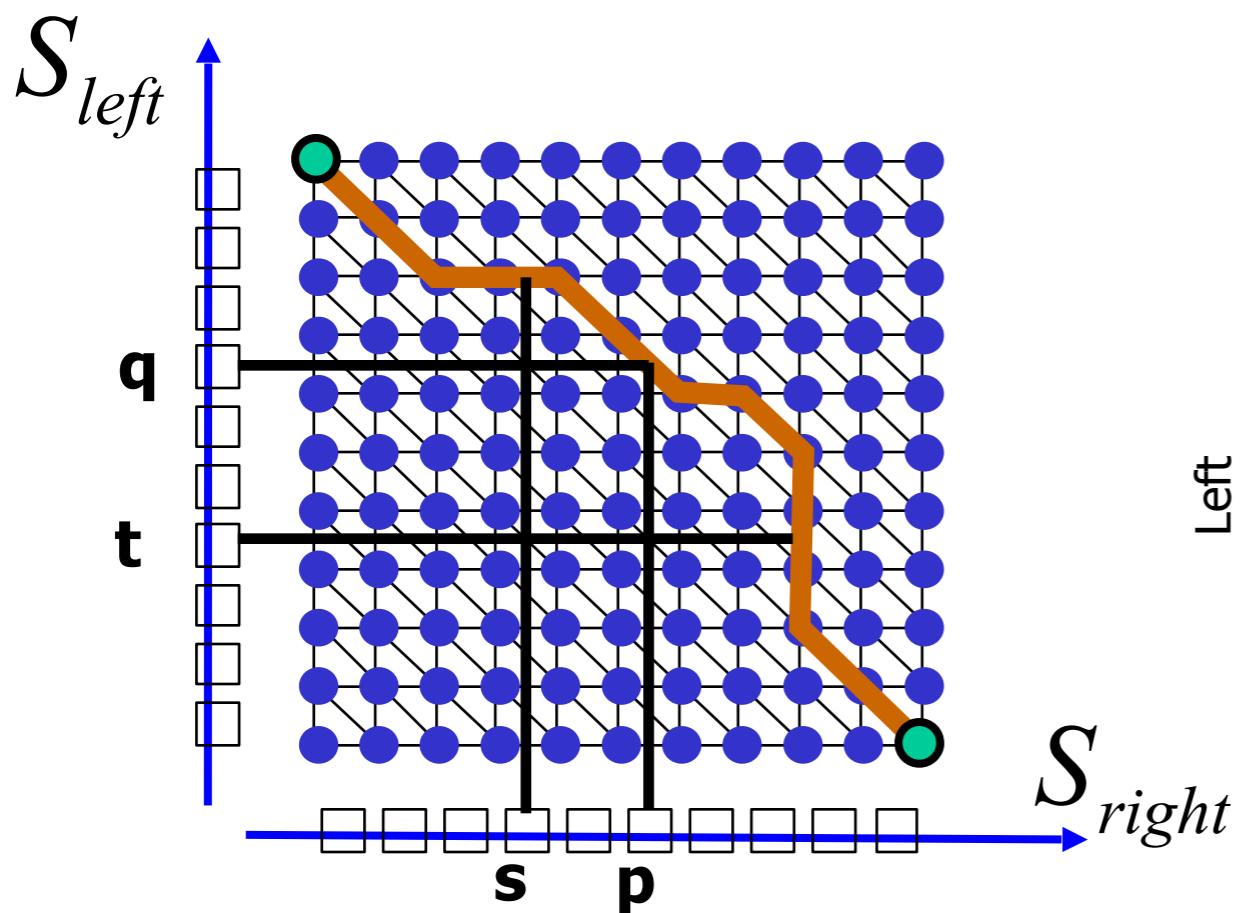
- Beyond individual correspondences to estimate disparities:
- Optimize correspondence assignments jointly
 - Scanline at a time (DP)
 - Full 2D grid (graph cuts)

Scanline stereo

- Try to coherently match pixels on the entire scanline
- Different scanlines are still optimized independently



“Shortest paths” for scan-line stereo

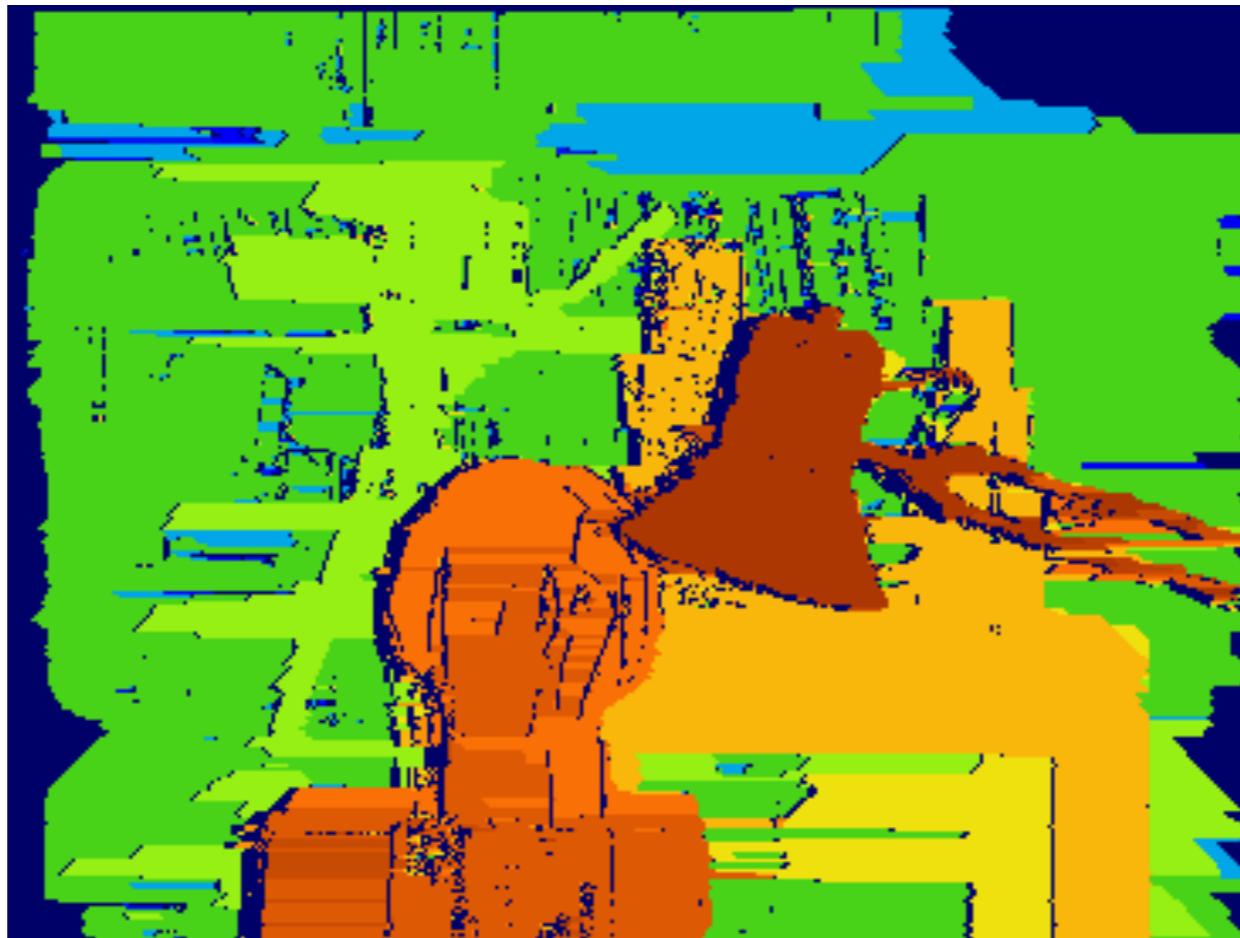


Can be implemented with dynamic programming
Ohta & Kanade '85, Cox et al. '96

Slide credit: Y. Boykov

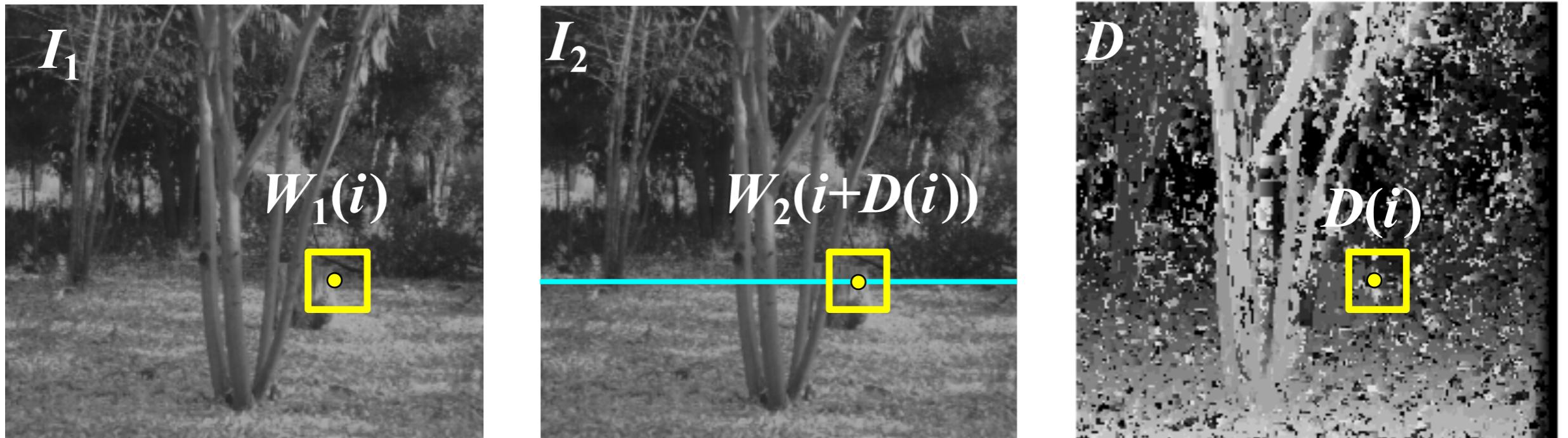
Coherent stereo on 2D grid

- Scanline stereo generates streaking artifacts



- Can't use dynamic programming to find spatially coherent disparities/ correspondences on a 2D grid

Stereo matching as energy minimization



$$E = \alpha E_{\text{data}}(I_1, I_2, D) + \beta E_{\text{smooth}}(D)$$

$$E_{\text{data}} = \sum_i (W_1(i) - W_2(i + D(i)))^2$$

$$E_{\text{smooth}} = \sum_{\text{neighbors } i, j} \rho(D(i) - D(j))$$

- Energy functions of this form can be minimized using *graph cuts*

Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001

Source: Steve Seitz

Recap: stereo with calibrated cameras

- Given image pair, \mathbf{R} , \mathbf{T}
- Detect some features
- Compute essential matrix \mathbf{E}
- Match features using the epipolar and other constraints
- Triangulate for 3d structure



Left



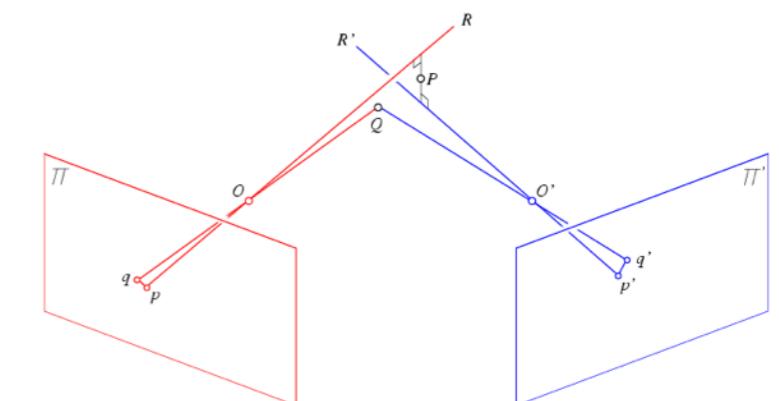
Right



Left



Right



Error sources

- Low-contrast ; textureless image regions
- Occlusions
- Camera calibration errors
- Violations of *brightness constancy* (e.g., specular reflections)
- Large motions

Today

- Recap: epipolar constraint
- Stereo image rectification
- Stereo solutions
 - Computing correspondences
 - Non-geometric stereo constraints
- Calibration
- Example stereo applications

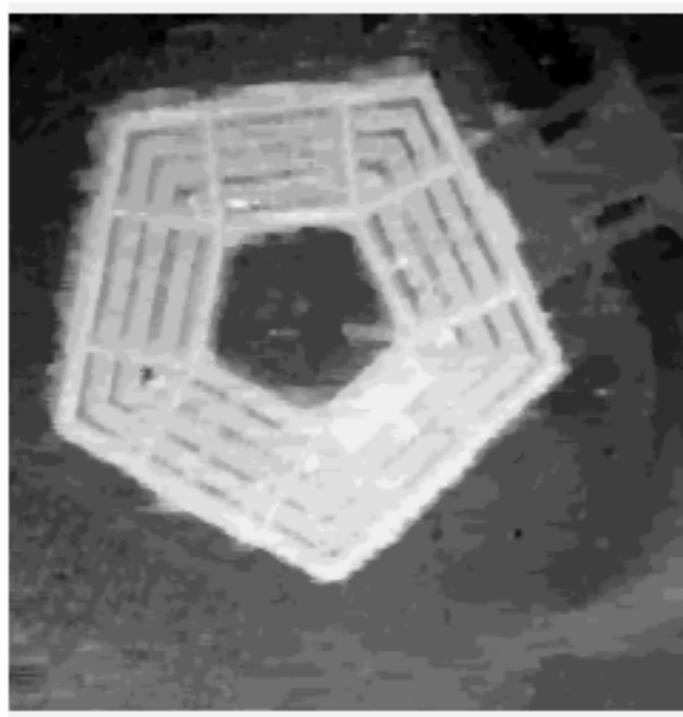
left image



right image



range map



value (4) for the constraint described earlier. The figure below shows the results obtained by both methods. It can be seen that the two methods show almost identical results. The boundaries of buildings in the ground part, the two ellipsoids provide approximately the same results. The most

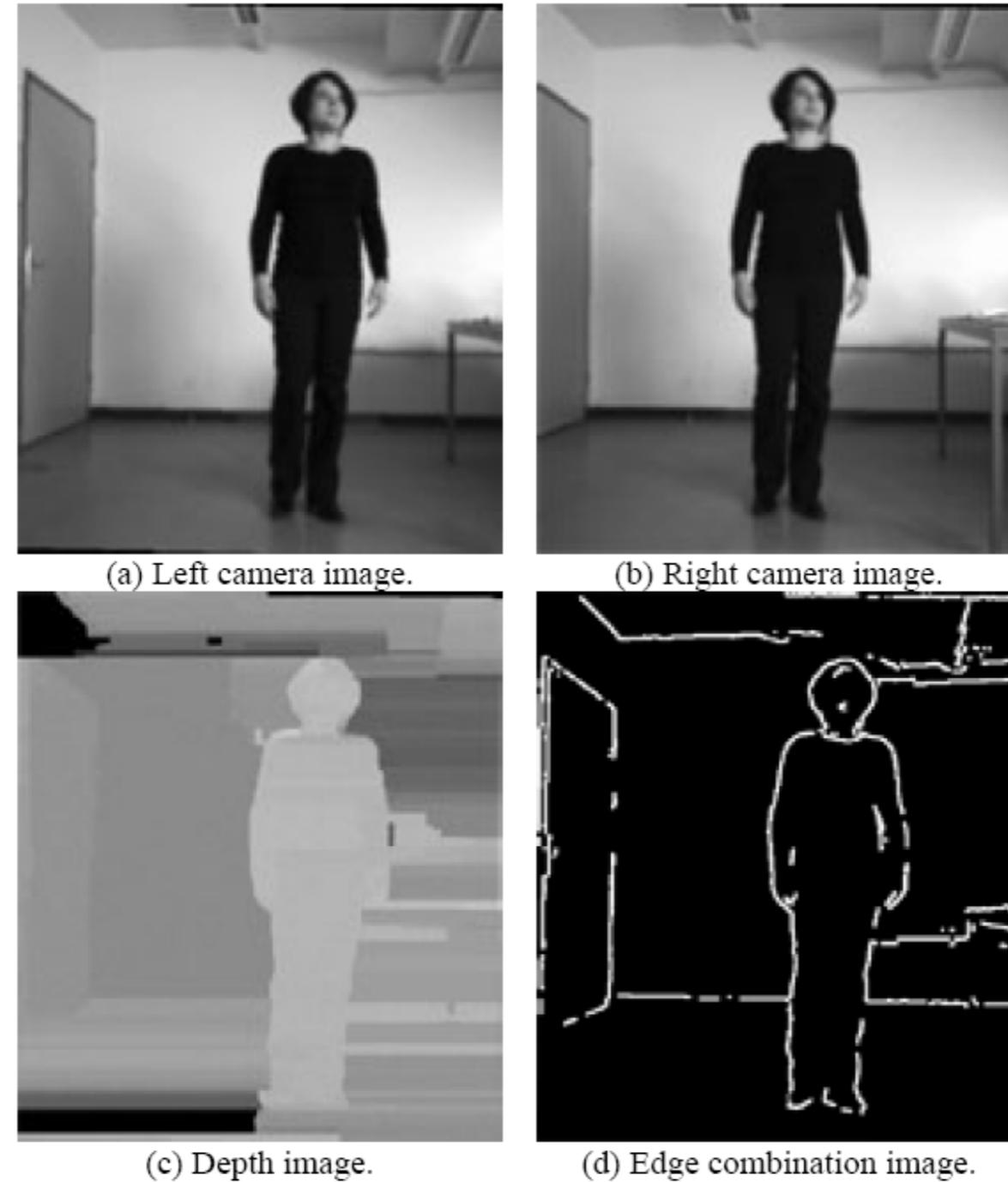
Stereo in machine vision systems



Left : The Stanford cart sports a single camera moving in discrete increments along a straight line and providing multiple snapshots of outdoor scenes

Right : The INRIA mobile robot uses three cameras to map its environment

Depth for segmentation



Edges in disparity in conjunction with image edges enhances contours found

Figure 3 Stereo video frames with computed depth map and edge combination result.

Depth for segmentation



(a) Original image with snake initialization.



(b) Final snake on original image.



(c) Final snake on depth image.



(d) Original image with snake from (c) overlaid.



(e) Final snake on edge combination image.



(f) Original image with snake from (e) overlaid.

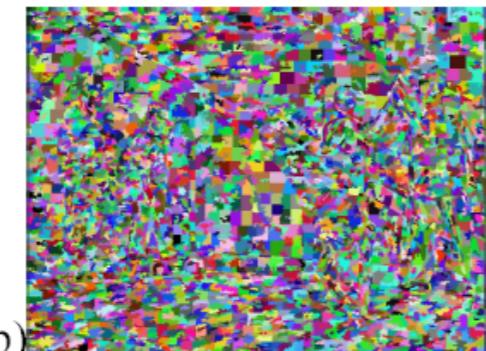
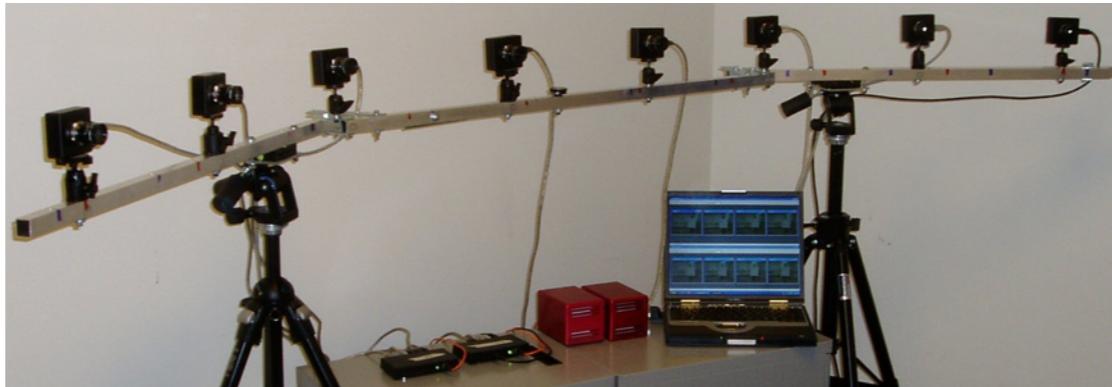
Model-based body tracking, stereo input



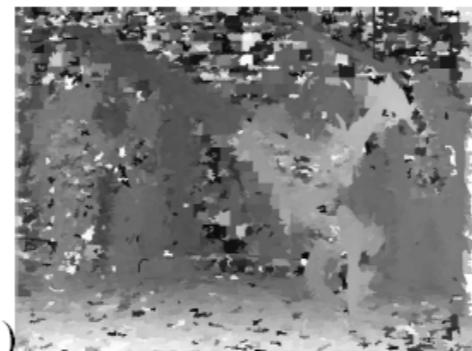
David Demirdjian, MIT Vision Interface Group

<http://people.csail.mit.edu/demirdji/movie/artic-tracker/turn-around.m1v>

Virtual viewpoint video



(a)



(c)



(d)



(e)

Figure 6: Sample results from stereo reconstruction stage: (a) input color image; (b) color-based segmentation; (c) initial disparity estimates \hat{d}_{ij} ; (d) refined disparity estimates; (e) smoothed disparity estimates $d_i(x)$.
(d) A depth-matted object from earlier in the sequence is inserted into the video.

Virtual viewpoint video



Massive Arabesque

Uncalibrated case

- What if we don't know the camera parameters?

Two possibilities:

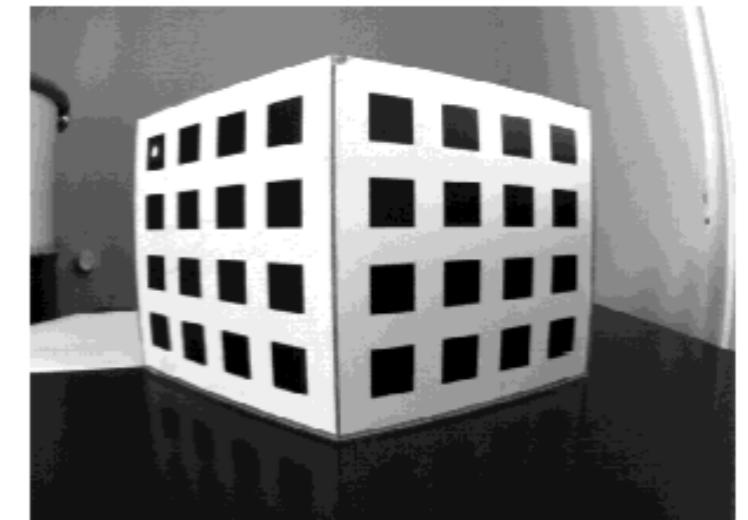
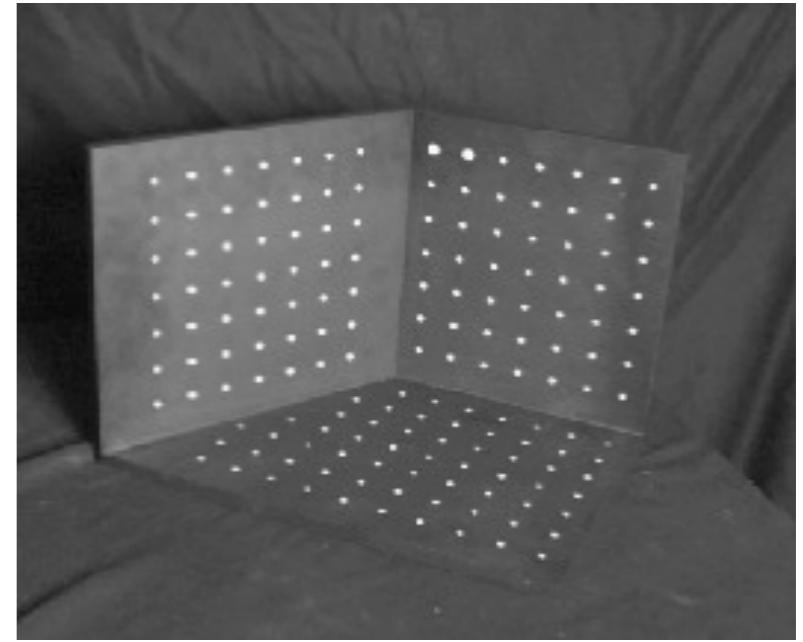
1. Calibrate with a calibration object
2. Weak calibration

Calibrating a camera

- Compute intrinsic and extrinsic parameters using observed camera data

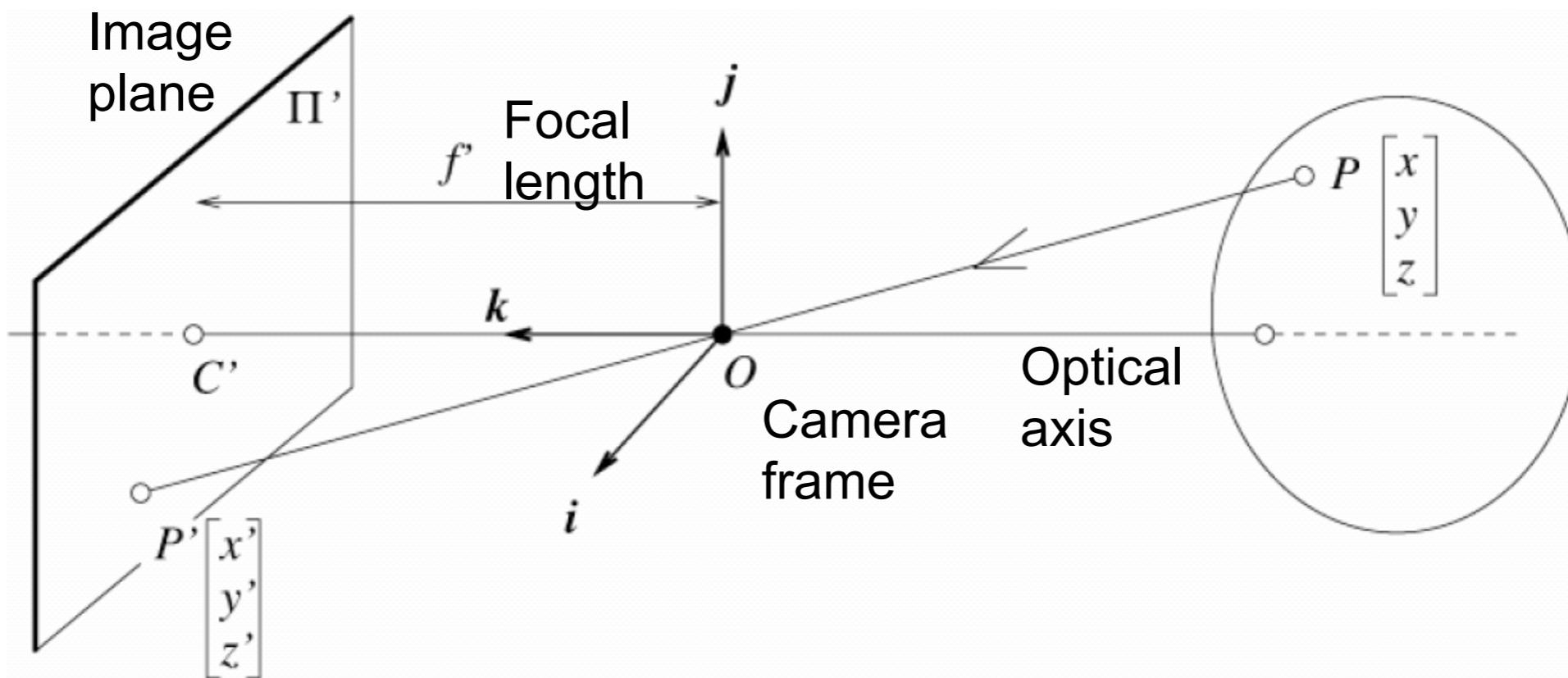
Main idea

- Place “calibration object” with known geometry in the scene
- Get correspondences
- Solve for mapping from scene to image



The Opti-CAL Calibration Target Image

Perspective projection



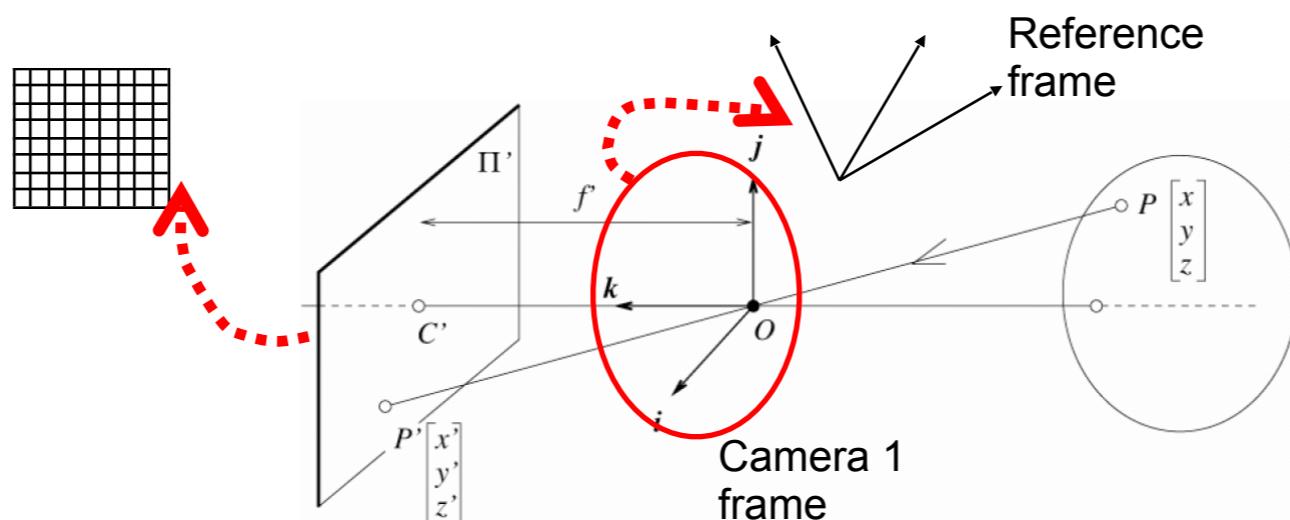
$$(x, y, z) \rightarrow \left(f \frac{x}{z}, f \frac{y}{z} \right)$$

Scene point \rightarrow Image coordinates

Thus far, in **camera's** reference frame only.

Camera parameters

- **Extrinsic:** location and orientation of camera frame with respect to reference frame
- Intrinsic: how to map pixel coordinates to image plane coordinates



Extrinsic camera parameters

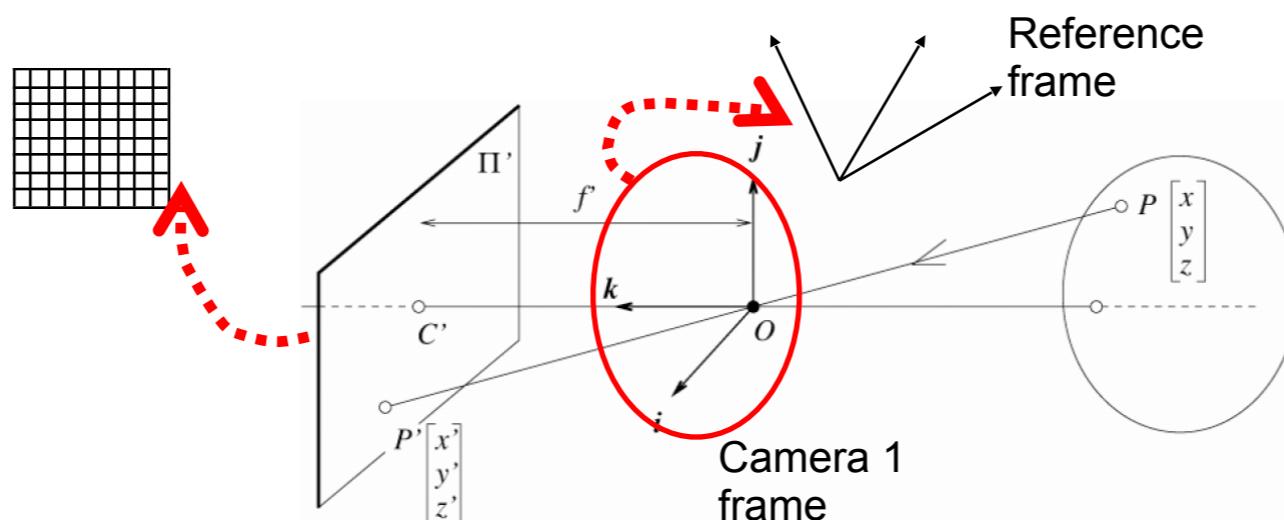
$$\mathbf{P}_c = \mathbf{R}(\mathbf{P}_w - \mathbf{T})$$

↑ ↑
Camera reference World reference
frame frame

$$\mathbf{P}_c = (X, Y, Z)^T$$

Camera parameters

- Extrinsic: location and orientation of camera frame with respect to reference frame
- **Intrinsic: how to map pixel coordinates to image plane coordinates**



Intrinsic camera parameters

- Ignoring any geometric distortions from optics, we can describe them by:

$$x = -(x_{im} - o_x) s_x$$

$$y = -(y_{im} - o_y) s_y$$

Coordinates of projected point in camera reference frame

Coordinates of image point in pixel units

Coordinates of image center in pixel units

Effective size of a pixel (mm)



Camera parameters

- We know that in terms of camera reference frame:

$$x = f \frac{X}{Z} \quad y = f \frac{Y}{Z} \quad \text{and} \quad \mathbf{P}_c = \mathbf{R}(\mathbf{P}_w - \mathbf{T})$$
$$\mathbf{P}_c = (X, Y, Z)^T$$

- Substituting previous eqns describing intrinsic and extrinsic parameters, can relate *pixels coordinates* to *world points*:

$$-(x_{im} - o_x)s_x = f \frac{\mathbf{R}_1 \cdot (\mathbf{P}_w - \mathbf{T})}{\mathbf{R}_3 \cdot (\mathbf{P}_w - \mathbf{T})}$$

\mathbf{R}_i = Row i of rotation matrix

$$-(y_{im} - o_y)s_y = f \frac{\mathbf{R}_2 \cdot (\mathbf{P}_w - \mathbf{T})}{\mathbf{R}_3 \cdot (\mathbf{P}_w - \mathbf{T})}$$

Projection matrix

- This can be rewritten as a matrix product using homogeneous coordinates:

where:

$$\begin{bmatrix} wx_{im} \\ wy_{im} \\ w \end{bmatrix} = \mathbf{M}P_w$$

$$\mathbf{M}_{int} = \begin{bmatrix} -f/s_x & 0 & o_x \\ 0 & -f/s_y & o_y \\ 0 & 0 & 1 \end{bmatrix}$$

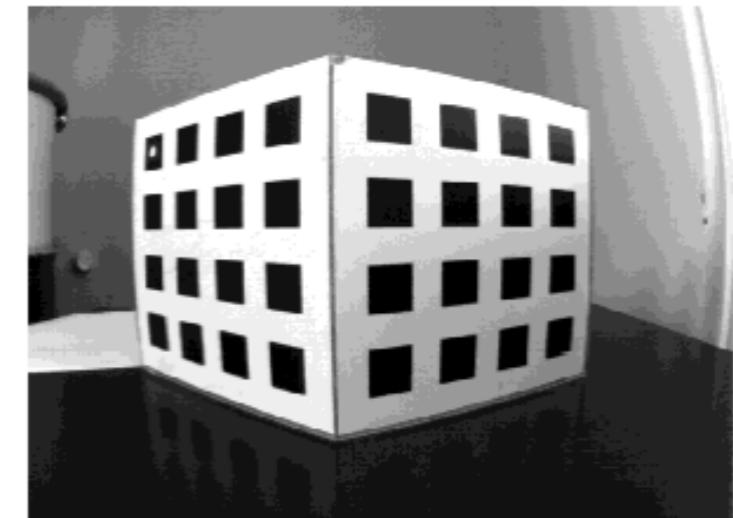
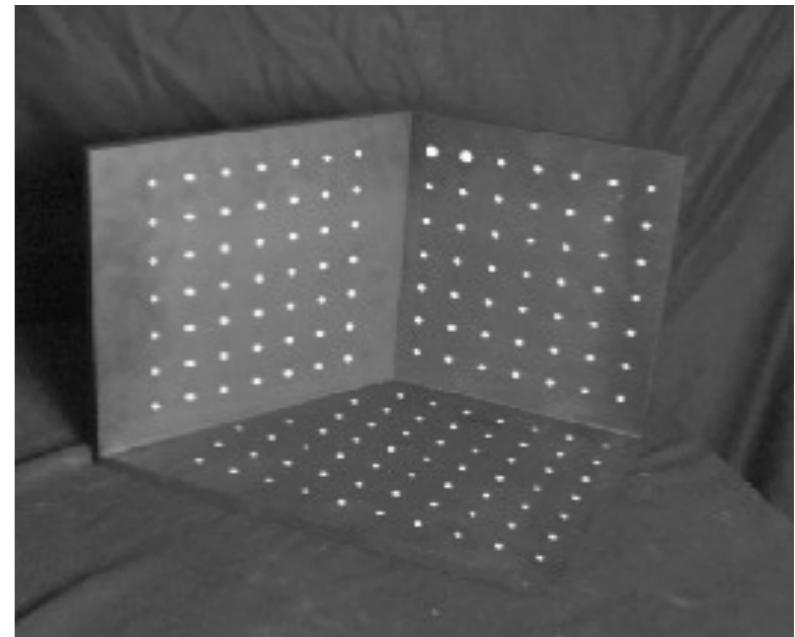
$$\mathbf{M}_{ext} = \begin{bmatrix} r_{11} & r_{12} & r_{13} & -\mathbf{R}_1^T \mathbf{T} \\ r_{21} & r_{22} & r_{23} & -\mathbf{R}_2^T \mathbf{T} \\ r_{31} & r_{32} & r_{33} & -\mathbf{R}_3^T \mathbf{T} \end{bmatrix}$$

Calibrating a camera

- Compute intrinsic and extrinsic parameters using observed camera data

Main idea

- Place “calibration object” with known geometry in the scene
- Get correspondences
- Solve for mapping from scene to image: estimate $\mathbf{M} = \mathbf{M}_{\text{int}} \mathbf{M}_{\text{ext}}$



The Opti-CAL Calibration Target Image

When would we calibrate this way?

- Makes sense when geometry of system is not going to change over time

...when would it change?

Weak calibration

- Want to estimate world geometry without requiring calibrated cameras
 - Archival videos
 - Photos from multiple unrelated users
 - Dynamic camera system
- **Main idea:**
 - Estimate epipolar geometry from a (redundant) set of point correspondences between two uncalibrated cameras

From before: Projection matrix

- This can be rewritten as a matrix product using homogeneous coordinates:

where:

$$\mathbf{M}_{int} = \begin{bmatrix} -f/s_x & 0 & o_x \\ 0 & -f/s_y & o_y \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} wx_{im} \\ wy_{im} \\ w \end{bmatrix} = \mathbf{M}_{int} \mathbf{M}_{ext} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$\mathbf{M}_{ext} = \begin{bmatrix} r_{11} & r_{12} & r_{13} & -\mathbf{R}_1^T \mathbf{T} \\ r_{21} & r_{22} & r_{23} & -\mathbf{R}_2^T \mathbf{T} \\ r_{31} & r_{32} & r_{33} & -\mathbf{R}_3^T \mathbf{T} \end{bmatrix}$$

From before: Projection matrix

- This can be rewritten as a matrix product using homogeneous coordinates:

$$\begin{bmatrix} wx_{im} \\ wy_{im} \\ w \end{bmatrix} = \mathbf{M}_{int} \mathbf{M}_{ext} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$\mathbf{p}_{im} = \mathbf{M}_{int} \mathbf{M}_{ext} \underbrace{\mathbf{P}_w}_{\mathbf{p}_c}$$

$$\mathbf{p}_{im} = \mathbf{M}_{int} \mathbf{p}_c$$

Uncalibrated case

For a given camera:

$$\mathbf{p}_{im} = \mathbf{M}_{int} \mathbf{p}_c$$

So, for two cameras (left and right):

$$\mathbf{p}_{c, left} = \mathbf{M}_{int, left}^{-1} \mathbf{p}_{im, left}$$

$$\mathbf{p}_{c, right} = \mathbf{M}_{int, right}^{-1} \mathbf{p}_{im, right}$$



Internal calibration
matrices, one per camera

$$\begin{aligned}\mathbf{p}_{c, \text{left}} &= \mathbf{M}_{int, \text{left}}^{-1} \mathbf{p}_{im, \text{left}} \\ \mathbf{p}_{c, \text{right}} &= \mathbf{M}_{int, \text{right}}^{-1} \mathbf{p}_{im, \text{right}}\end{aligned}$$

Uncalibrated case

$$\mathbf{p}_{c, \text{right}}^T \mathbf{E} \mathbf{p}_{c, \text{left}} = 0$$

From before, the **essential matrix** \mathbf{E} .

$$(\mathbf{M}_{int, \text{right}}^{-1} \mathbf{p}_{im, \text{right}})^T \mathbf{E} (\mathbf{M}_{int, \text{left}}^{-1} \mathbf{p}_{im, \text{left}}) = 0$$

$$\mathbf{p}_{im, \text{right}}^T \underbrace{(\mathbf{M}_{int, \text{right}}^{-T} \mathbf{E} \mathbf{M}_{int, \text{left}}^{-1})}_{\mathbf{F}} \mathbf{p}_{im, \text{left}} = 0$$

F “Fundamental matrix”

$$\mathbf{p}_{im, \text{right}}^T \mathbf{F} \mathbf{p}_{im, \text{left}} = 0$$

Computing F from correspondences

Each point correspondence generates one constraint on F

$$\mathbf{p}_{im,right}^T \mathbf{F} \mathbf{p}_{im,left} = 0$$

$$\begin{bmatrix} u' & v' & 1 \end{bmatrix} \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = 0$$

Collect n of these constraints

$$\begin{bmatrix} u'_1 u_1 & u'_1 v_1 & u'_1 & v'_1 u_1 & v'_1 v_1 & v'_1 & u_1 & v_1 & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = 0$$

Solve for f , vector of parameters.

Fundamental matrix

- Relates **pixel coordinates** in the two views
- More general form than essential matrix: we remove need to know intrinsic parameters
- If we estimate fundamental matrix from correspondences in *pixel coordinates*, can reconstruct epipolar geometry without intrinsic or extrinsic parameters.

Stereo pipeline with weak calibration

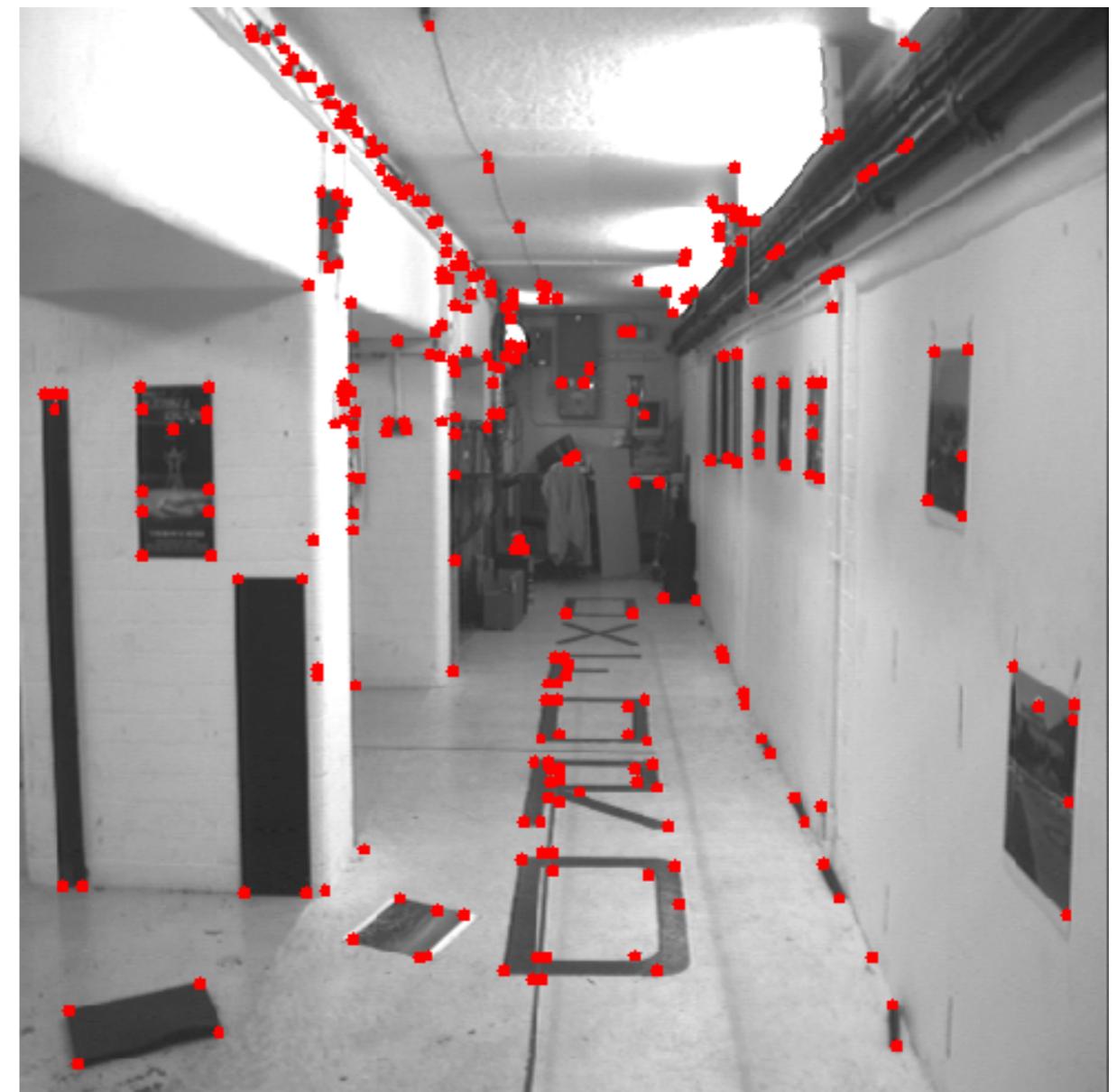
- So, where to start with uncalibrated cameras?
 - Need to find fundamental matrix F **and** the correspondences (pairs of points $(u',v') \leftrightarrow (u,v)$).



- 1) Find interest points in image
- 2) Compute correspondences
- 3) Compute epipolar geometry
- 4) Refine

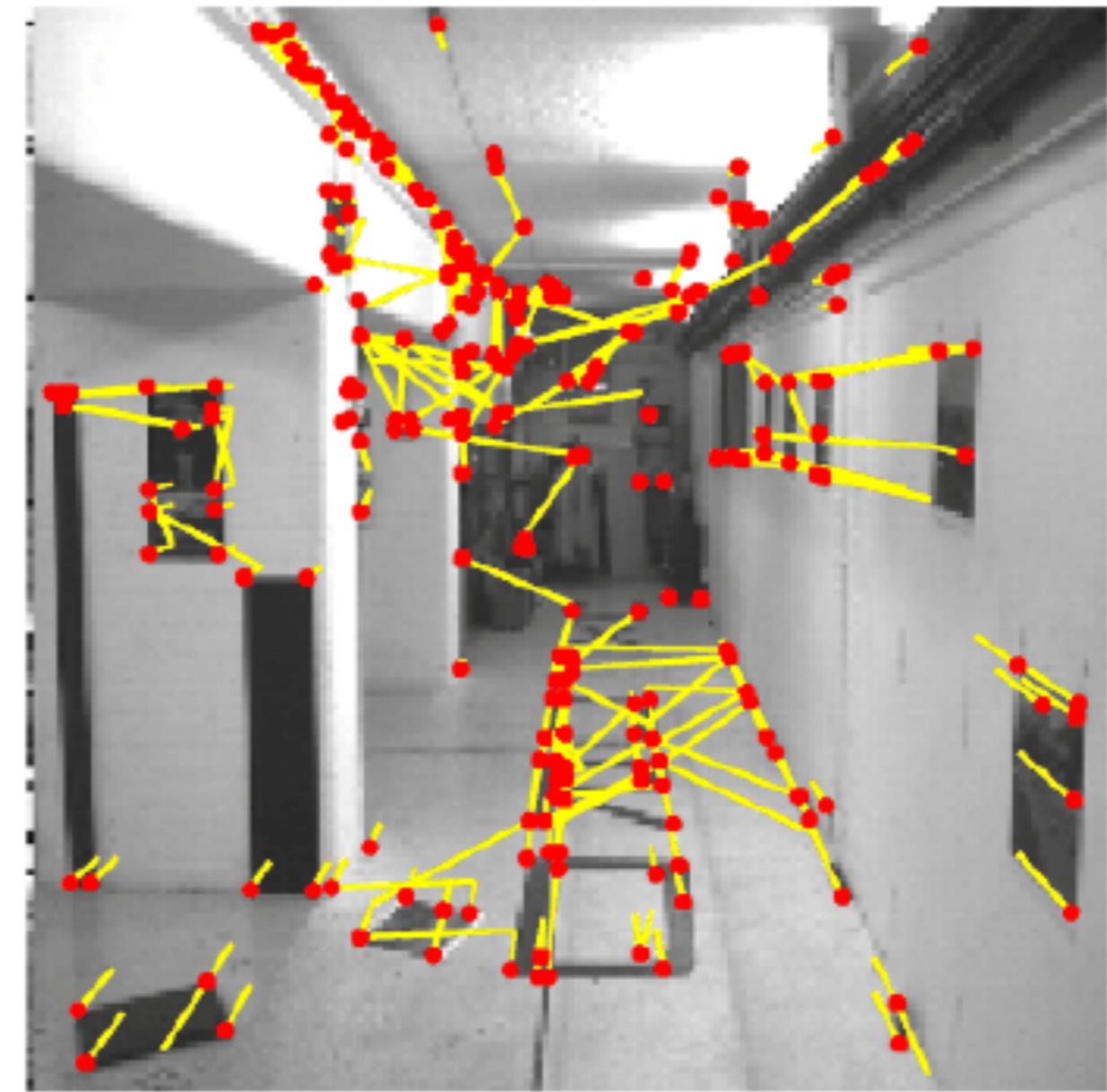
Stereo pipeline with weak calibration

1) Find interest points



Stereo pipeline with weak calibration

2) Match points within proximity to get putative matches



Stereo pipeline with weak calibration

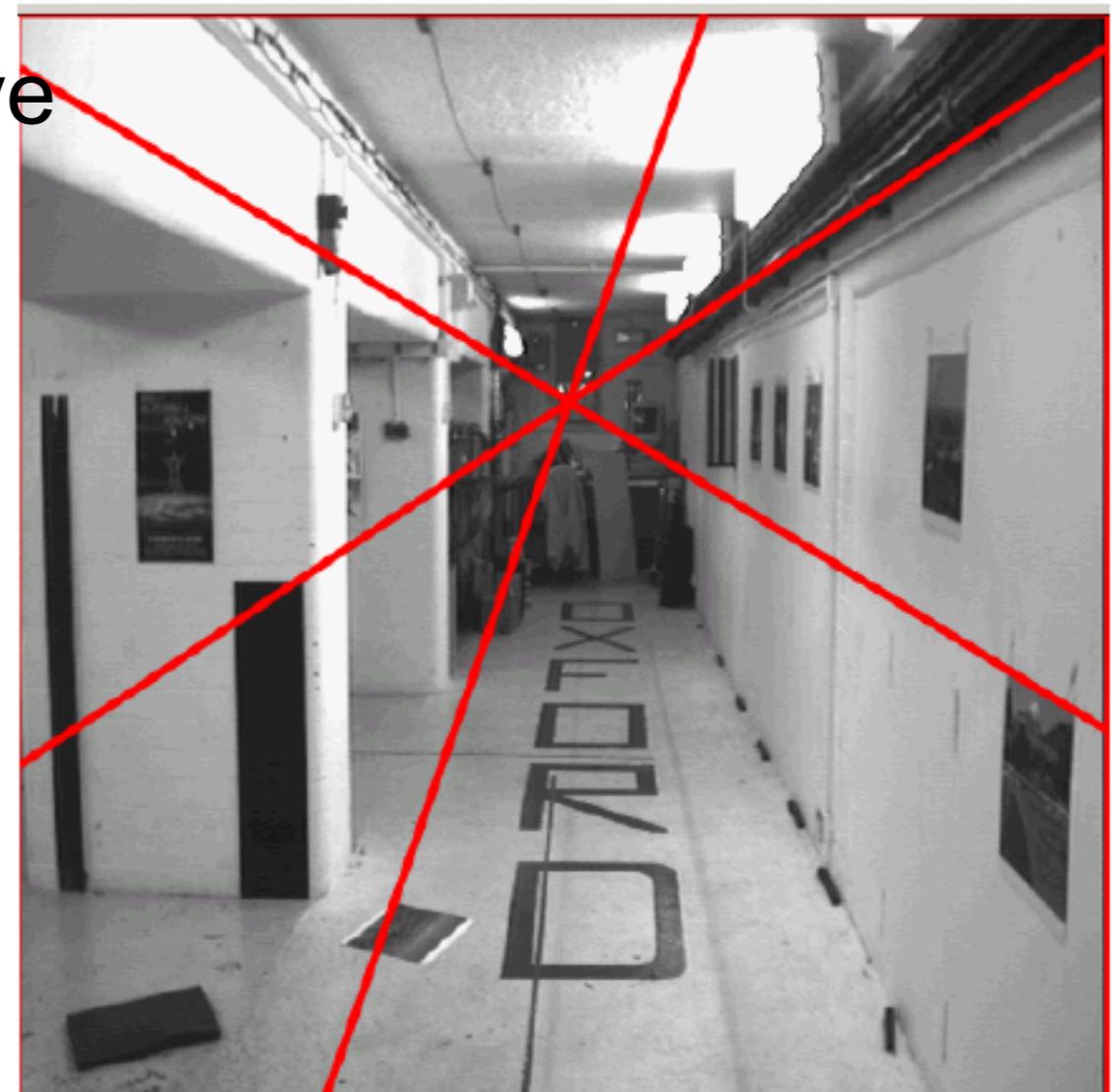
3) Compute epipolar geometry -- robustly with RANSAC

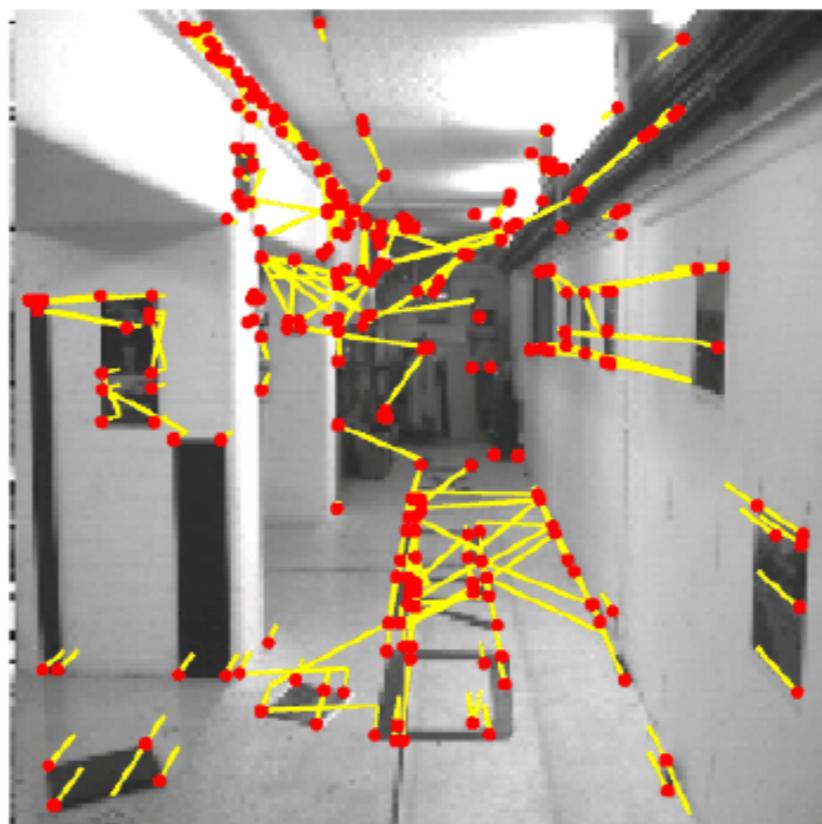
Select random sample of putative correspondences

Compute F using them
- determines epipolar constraint

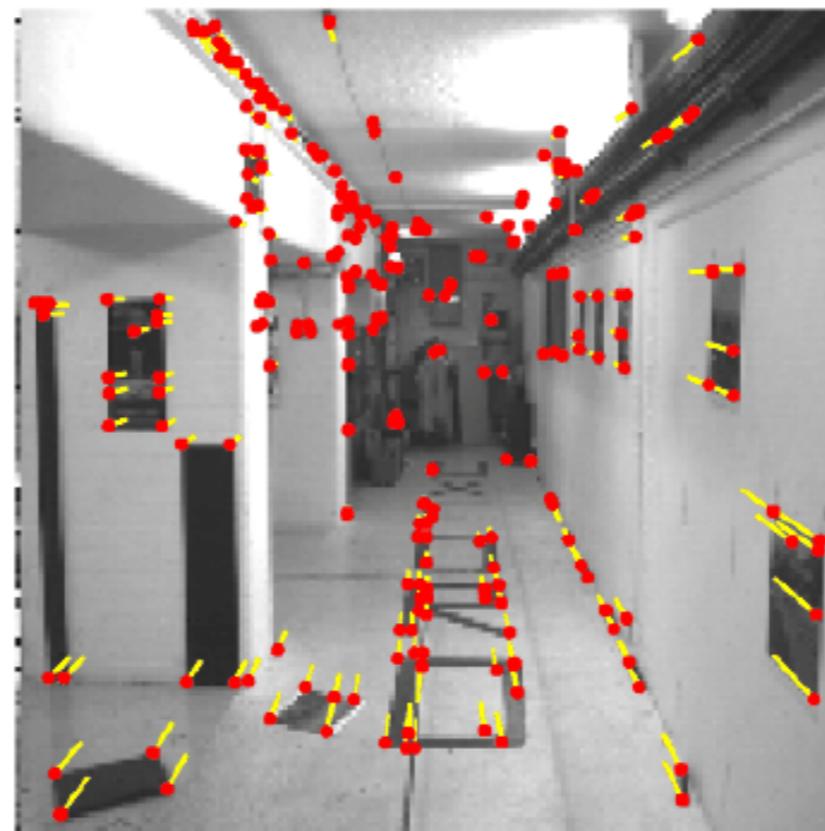
Evaluate amount of support
- inliers within threshold distance of epipolar line

Choose F with most support
(inliers)





Using correlation search to get putative matches: noisy, but enough to compute F using RANSAC



Pruned matches: those consistent with epipolar geometry

Summary

- **Rectification:** make epipolar lines align with scanlines
- Stereo solutions:
 - **Correspondence:** dense, or at interest points
 - **Non-geometric stereo constraints** (e.g., similarity, order, smoothness)
- Calibration
 - **With calibration object** in scene: relate world coordinates to image coordinates
 - **Weak calibration:** solve for fundamental matrix, relate image coordinates to image coordinates

Demo: Code Review