

Computer Vision I

From Single to Two-View Geometry -
03.07.2013



TECHNISCHE
UNIVERSITÄT
DARMSTADT





Announcements

- ◆ Class next week
 - ◆ Will be held by Stephan Richter
- ◆ Class in two weeks
 - ◆ is cancelled (hence double lecture today)
- ◆ Q&A Session
 - ◆ Wednesday 24.7.2013, 9:50
 - ◆ Chance to ask detailed questions on the entire material
 - ◆ Participation optional (otherwise office hours, email, etc.)



Outline

- ◆ Part 1: Camera calibration
- ◆ Part 2: Single view geometry
- ◆ Homography
 - ◆ Perspective image of a plane
 - ◆ Views from same camera centre
- ◆ (Homography) estimation
 - ◆ Numerical conditioning of coordinates
 - ◆ Robust fitting with RANSAC
- ◆ Application: Panorama stitching

Projecting a planar object



- ◆ Why is this useful?
 - ◆ undistort perspective images of planar structures
 - ◆ map from one plane to another



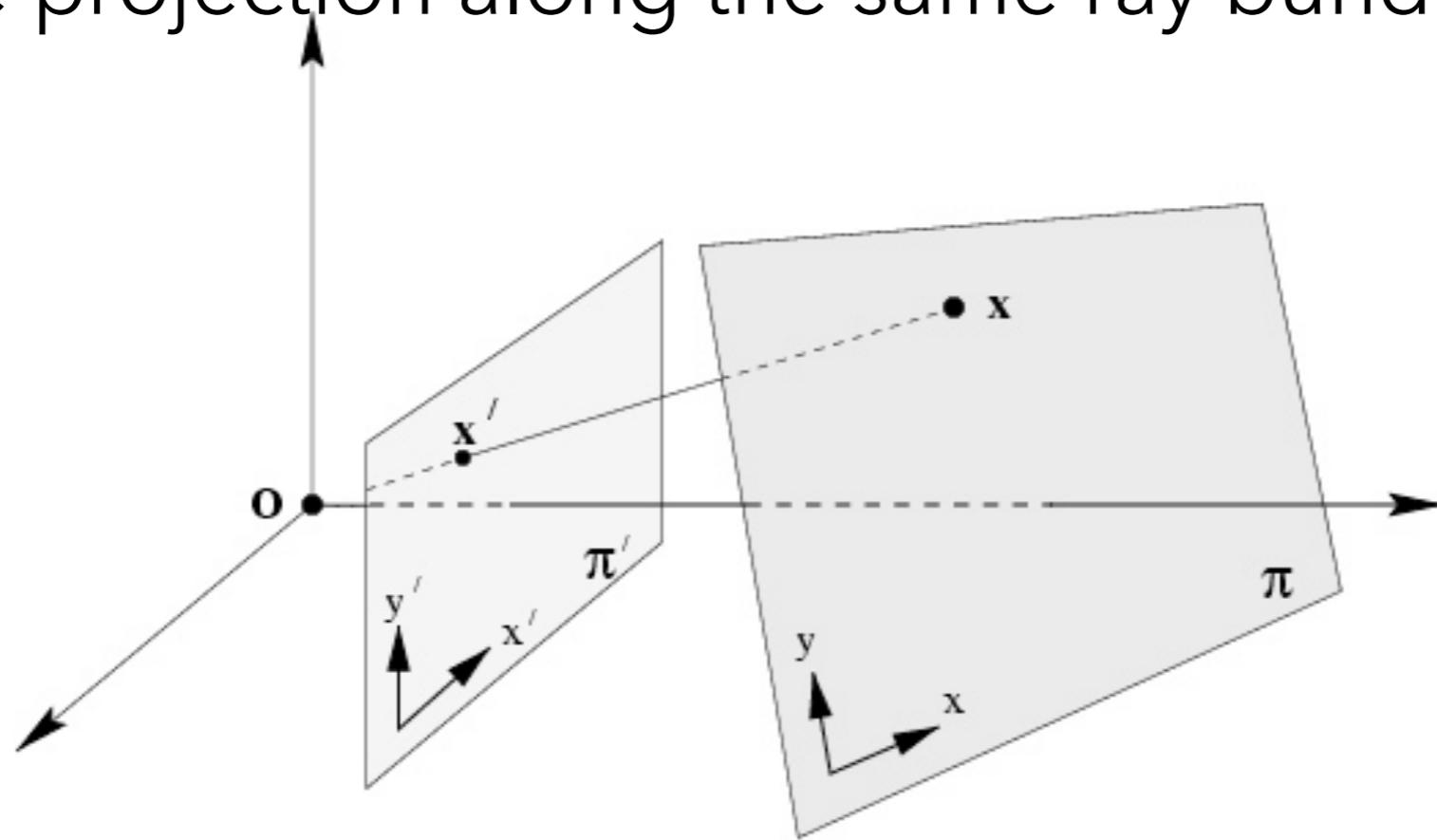
Projecting a planar object

- ◆ Why is this useful?
 - ◆ undistort perspective images of planar structures
 - ◆ map from one plane to another
- ◆ geometrically: How does a pinhole camera map a plane?
- ◆ algebraically: How does a projection matrix act on co-planar points?



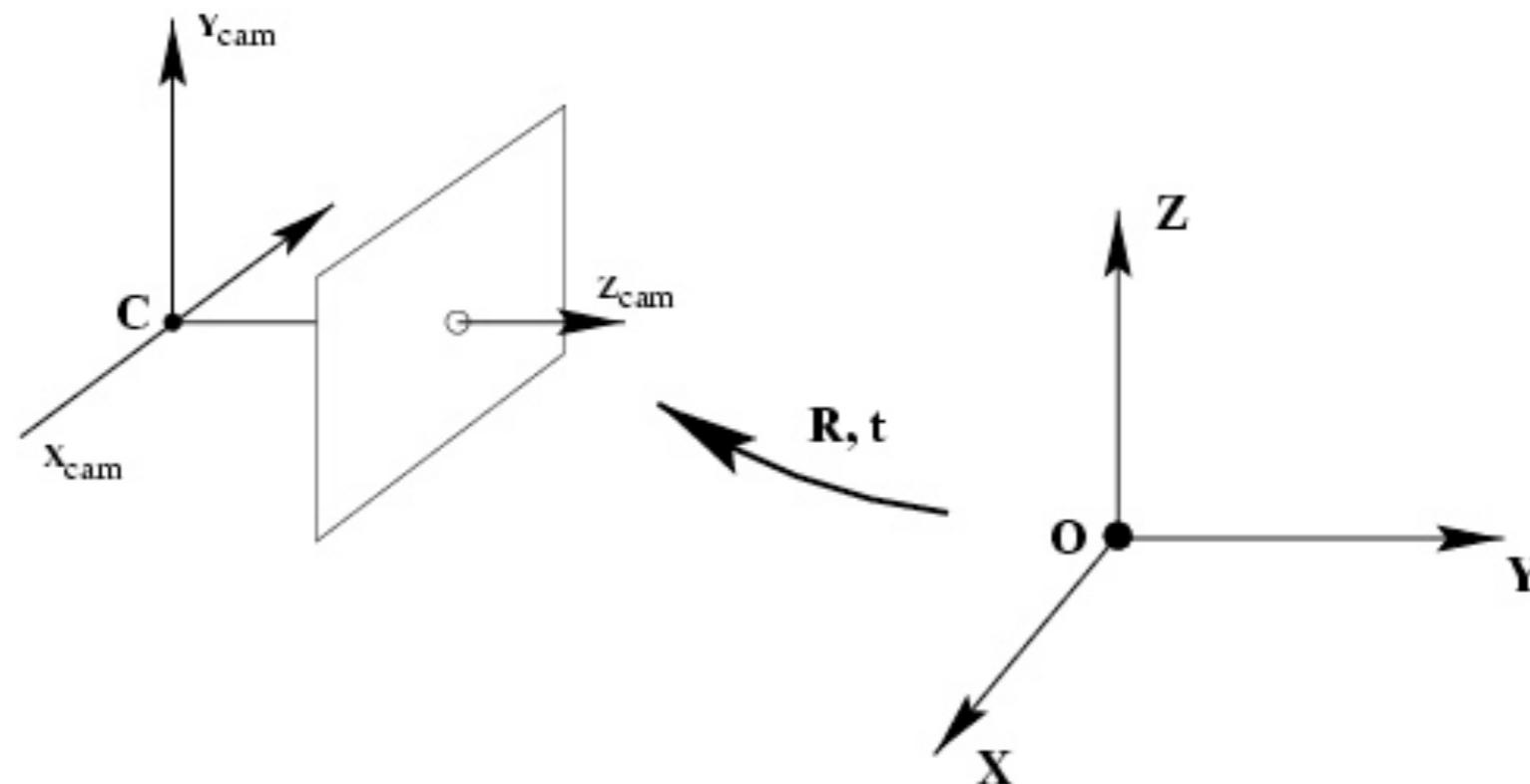
Geometric intuition

- ◆ Projection with a straight ray from a 2D plane to another 2D plane (the retina of the camera)
- ◆ No dimensionality reduction
- ◆ Ray intersects each plane **exactly once**
- ◆ Inverse projection along the same ray bundle



Pinhole camera model

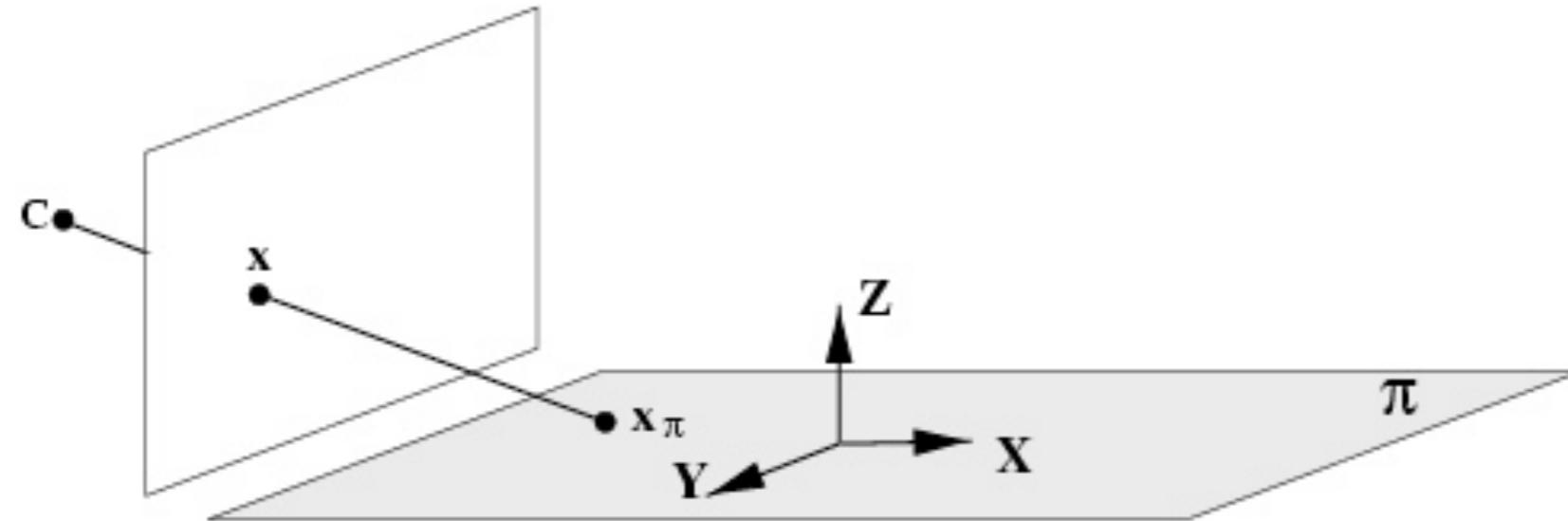
- ◆ Remember: a pinhole camera is represented by a (3x4) homogeneous projection matrix P



$$x = PX = K[R| - R\tilde{C}]X$$

Projecting a plane

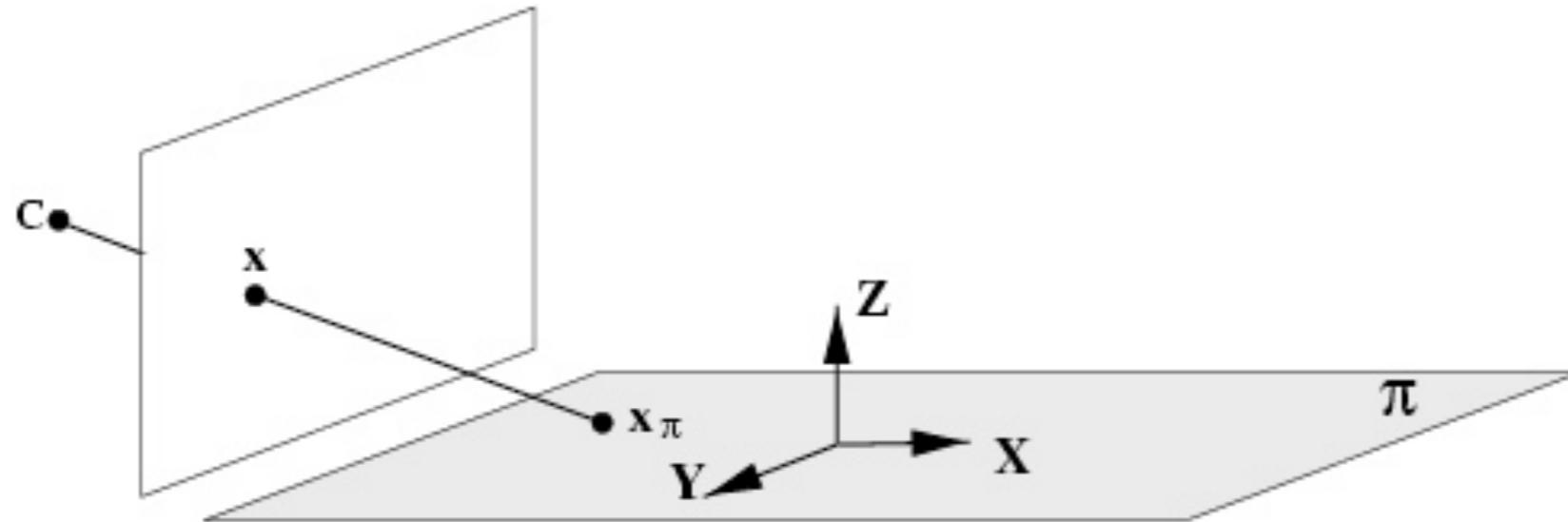
- ◆ Coordinate system in the plane: $Z=0$



$$x = ?x_\pi$$

Projecting a plane

- ◆ Coordinate system in the plane: $Z=0$



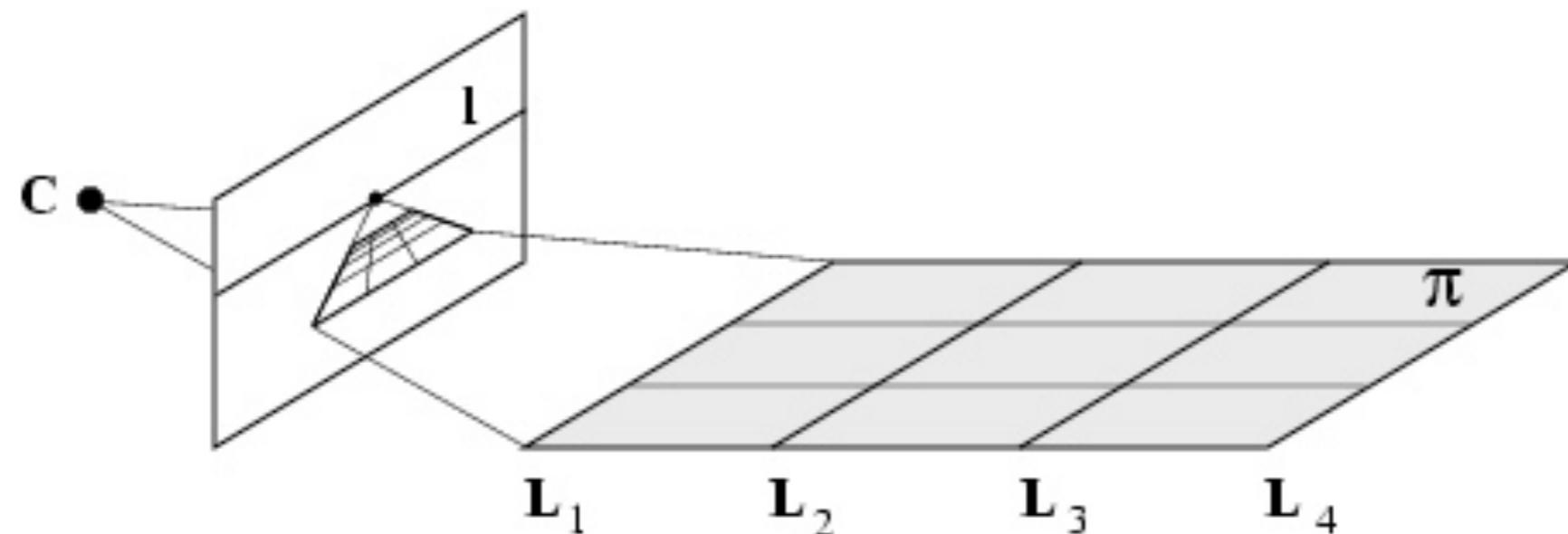
$$x = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ 0 \\ 1 \end{bmatrix} =$$

$$= \begin{bmatrix} p_{11} & p_{12} & p_{14} \\ p_{21} & p_{22} & p_{24} \\ p_{31} & p_{32} & p_{34} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ 1 \end{bmatrix} = Hx_\pi$$

Projecting a plane

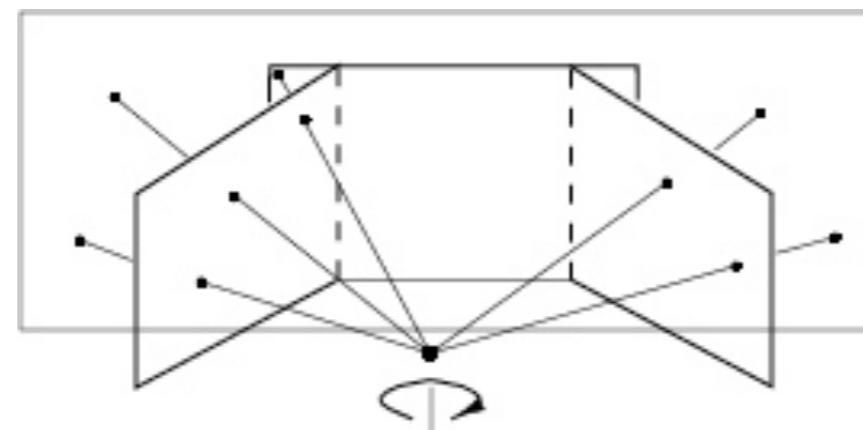
- ◆ The mapping is given by a (3x3) matrix H
 - ◆ H is called a **homography** (or projectivity)
 - ◆ the most general projective 2D transformation
 - ◆ maps a square to a general quadrangle
 - ◆ has 8 degrees of freedom
 - ◆ H is a one-to-one mapping, hence invertible

$$x = Hx_\pi \quad x_\pi = H^{-1}x$$



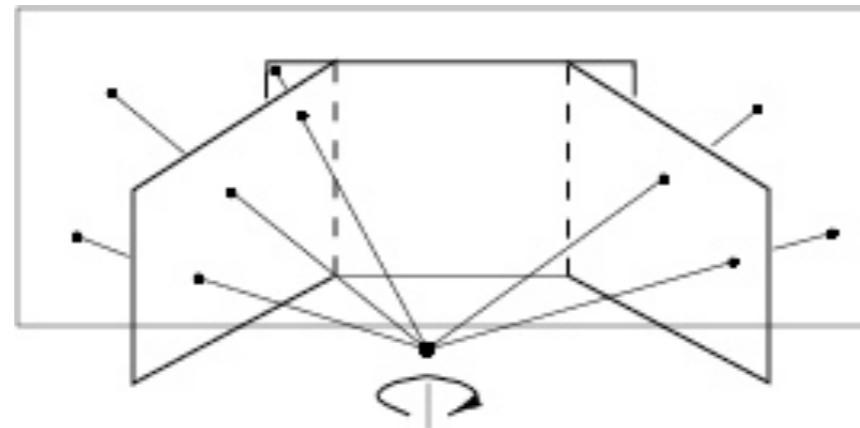
Rotating camera

- ◆ Relation between two images, if the camera only rotates (baseline = 0)



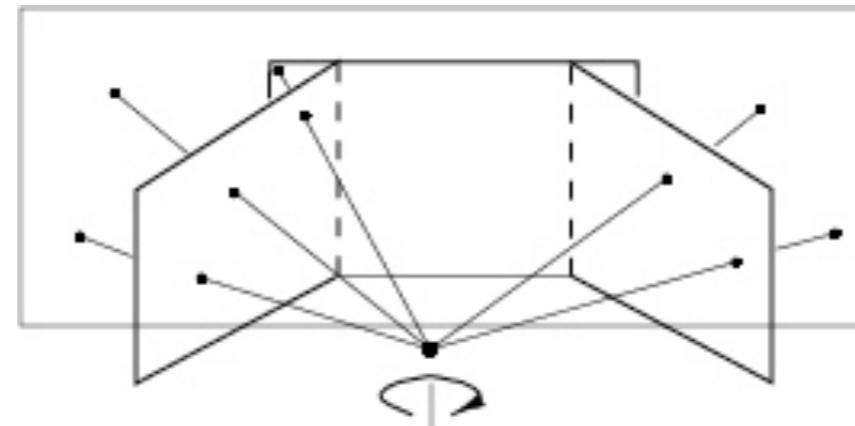
Rotating camera

- ◆ Relation between two images, if the camera only rotates (baseline = 0)
 - ◆ Projection rays are identical
 - ◆ ray maps directly from one image plane to the other
 - ◆ the two images are related by a [homography](#)



Rotating camera

- ◆ Relation between two images, if the camera only rotates (baseline = 0)
 - ◆ Projection rays are identical
 - ◆ ray maps directly from one image plane to the other
 - ◆ the two images are related by a homography



$$x = K[I|0]X = K\tilde{X}$$

$$\tilde{X} = K^{-1}x$$

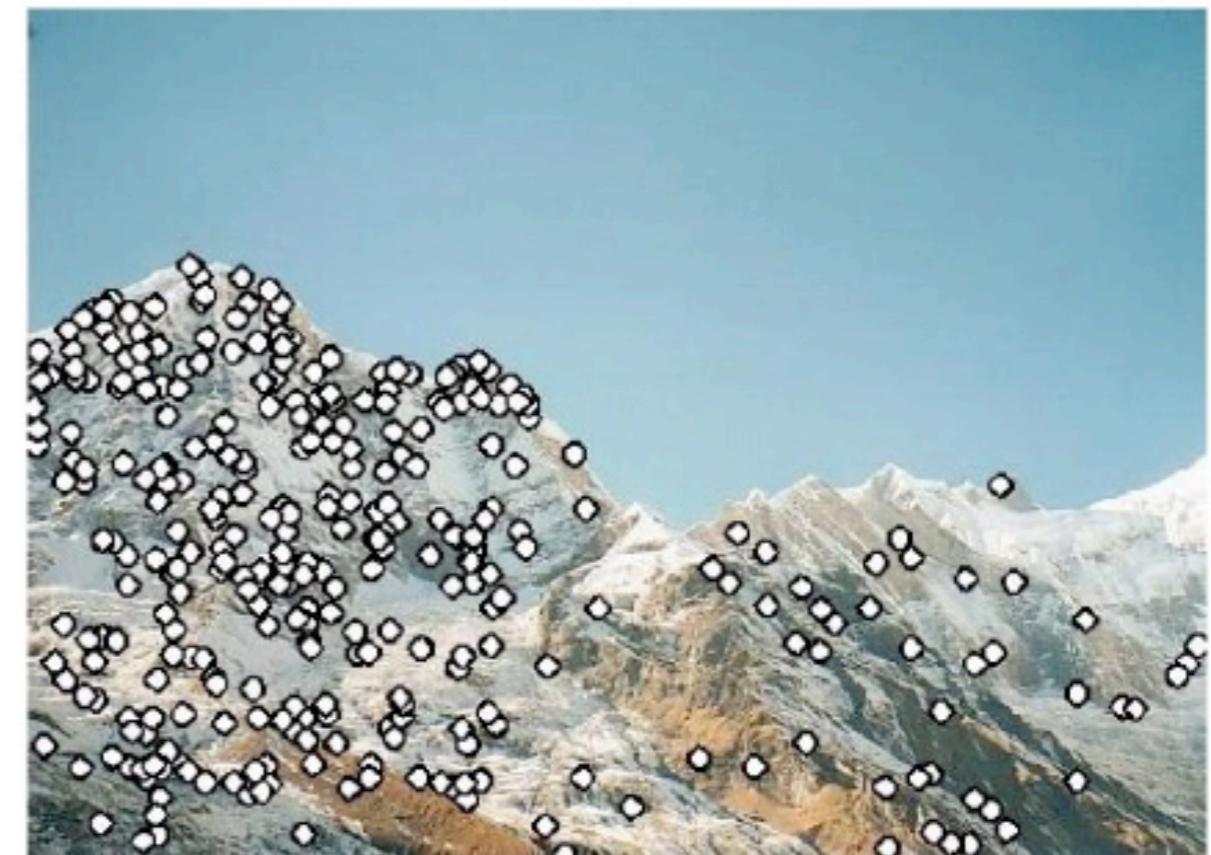
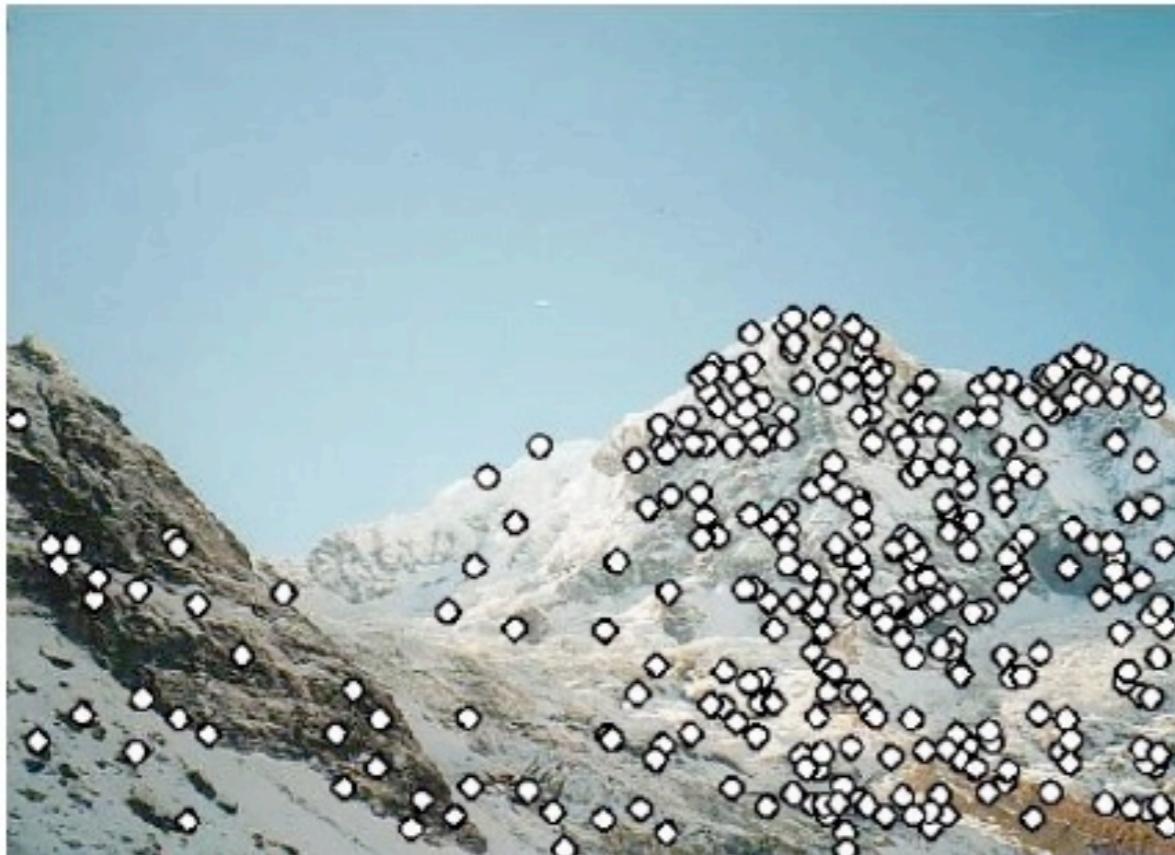
$$x' = K'[R|0]X = K'R\tilde{X}$$

$$x' = \boxed{K'RK^{-1}}x$$

$$x' = Hx$$

Rotating camera

- ◆ Application
 - ◆ Panorama stitching
 - ◆ estimate the homography between two overlapping images, transform one into the image plane of the other



Rotating camera

- ◆ Application
 - ◆ Panorama stitching
 - ◆ estimate the homography between two overlapping images, transform one into the image plane of the other



Estimating the homography

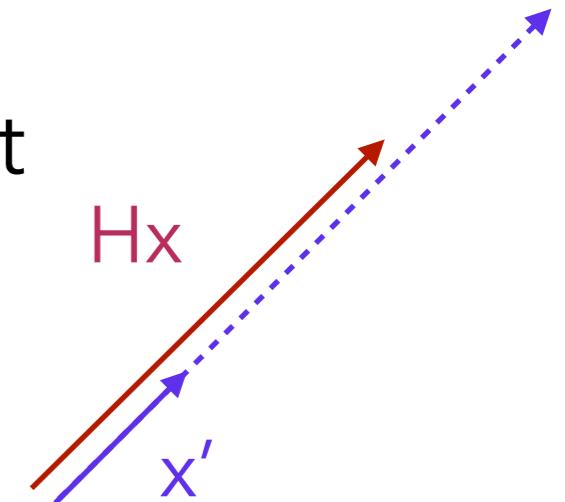
- ◆ Eight unknowns
- ◆ Two linearly independent equations per point
- ◆ Need 4 point correspondences

$$\lambda x' = Hx \quad , \quad x' \propto Hx$$

Estimating the homography

- ◆ Eight unknowns
- ◆ Two linearly independent equations per point
- ◆ Need 4 point correspondences

$$x' \propto Hx \quad \rightarrow \quad Hx \times x' = 0$$



$$\begin{bmatrix} 0 & 0 & 0 & x & y & 1 & -xy' & -yy' & -y' \\ -x & -y & -1 & 0 & 0 & 0 & xx' & yx' & x' \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \\ h_{21} \\ h_{22} \\ h_{23} \\ h_{31} \\ h_{32} \\ h_{33} \end{bmatrix} = 0$$

Estimating the homography

- ◆ Stack equations from at least 4 correspondences
- ◆ Remember
 - ◆ a homogeneous linear equation system

Solution ?

Estimating the homography

- ◆ Stack equations from at least 4 correspondences
- ◆ Remember
 - ◆ a homogeneous linear equation system
 - ◆ solve s.t. $\|h\|=1$, using SVD

$$Ah = 0$$

$$[U, D, V] = \text{svd}(A)$$

$$h = \begin{bmatrix} v_{19} \\ v_{29} \\ \vdots \\ v_{99} \end{bmatrix}$$

Estimating the homography

- ◆ A note on numerical stability
 - ◆ Coefficients of an equation system should be in the same order of magnitude so as not to lose significant digits
 - ◆ in pixels: $xx' \sim 1E6$
 - ◆ Conditioning: scale and shift points to be in [-1..1]
 - ◆ apply inverse transform to estimated homography
 - ◆ a general rule, not only for homography

$$\begin{bmatrix} 0 & 0 & 0 & x & y & 1 & -xy' & -yy' & -y' \\ -x & -y & -1 & 0 & 0 & 0 & xx' & yx' & x' \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \\ h_{21} \\ h_{22} \\ h_{23} \\ h_{31} \\ h_{32} \\ h_{33} \end{bmatrix} = 0$$

Estimating the homography

- ◆ A note on numerical stability
 - ◆ Conditioning: scale and shift points to be in [-1..1]

$$s = \frac{1}{2} \max_i (\|x_i\|)$$

$$t = \text{mean}(x_i)$$

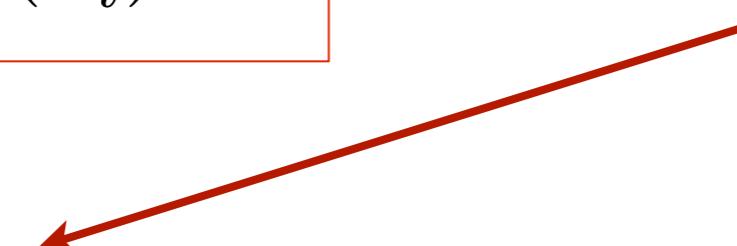
$$s' = \frac{1}{2} \max_i (\|x'_i\|)$$

$$t' = \text{mean}(x'_i)$$



$$T = \begin{bmatrix} \frac{1}{s} & 0 & -\frac{t_x}{s} \\ 0 & \frac{1}{s} & -\frac{t_y}{s} \\ 0 & 0 & 1 \end{bmatrix}$$

$$T' = \begin{bmatrix} \frac{1}{s'} & 0 & -\frac{t'_x}{s'} \\ 0 & \frac{1}{s'} & -\frac{t'_y}{s'} \\ 0 & 0 & 1 \end{bmatrix}$$



$$u = Tx$$

SVD

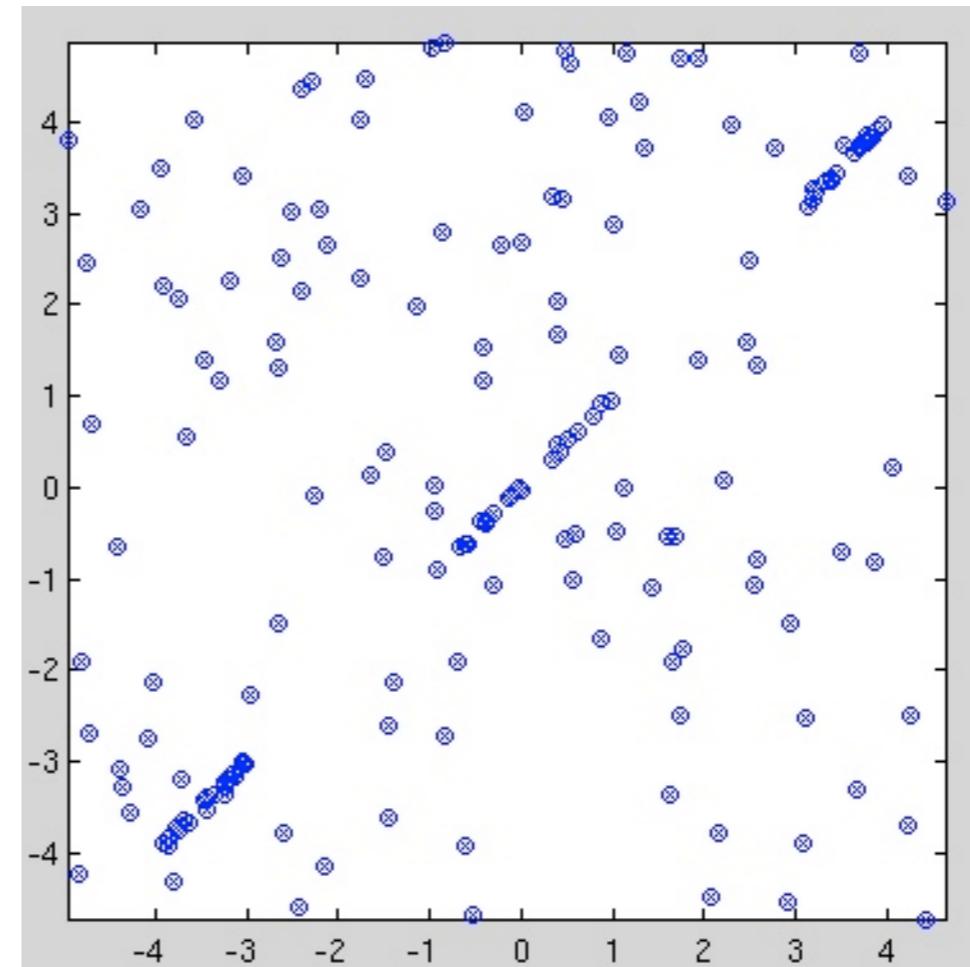
$$u' = T'x'$$

$$u' = \bar{H}u$$

$$\begin{aligned} T'x' &= \bar{H}Tx \\ H &= (T')^{-1}\bar{H}T \end{aligned}$$

Robust Estimation

- ◆ Problem: some correspondences may be wrong
- ◆ Using false correspondences will corrupt the estimate
- ◆ Need a method to get rid of matching errors
- ◆ Again general, not only for homography



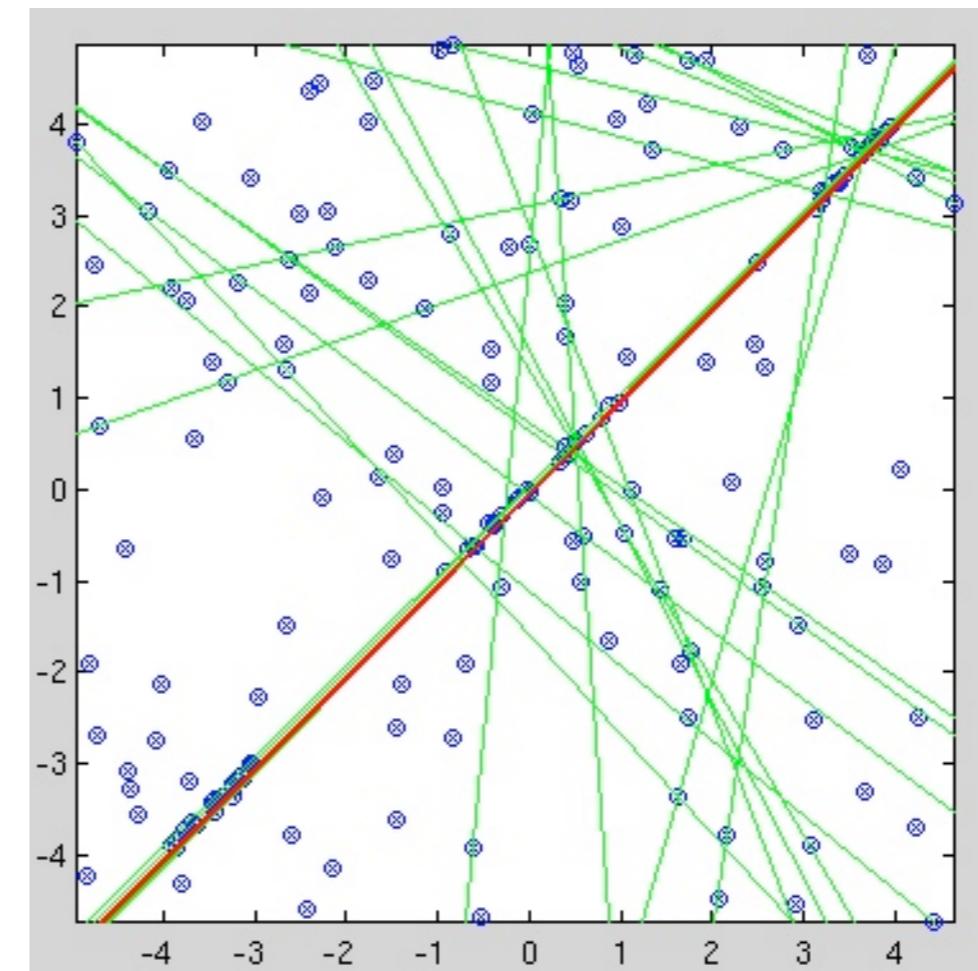
Robust Estimation

◆ RANSAC (random sample consensus)

◆ iterate:

- ◆ randomly pick a minimal set of points
- ◆ construct the model (a line, a homography, ...) from the points
- ◆ measure the support (count number of points with small error)

◆ keep the best hypothesis



Robust Estimation



- ◆ RANSAC (random sample consensus)
 - ◆ we cannot try all possible samples - how many is "enough"?



Robust Estimation

- ◆ RANSAC (random sample consensus)
 - ◆ we cannot try all possible samples - how many is "enough"?
 - ◆ assume we have a conservative guess of the fraction of "inliers" (correct points)

- ◆ say 20% of the correspondences are wrong
- ◆ how many random samples do we have to pick, in order to "almost certainly" find the right homography?



Robust Estimation

◆ RANSAC (random sample consensus)

- ◆ we cannot try all possible samples - how many is "enough"?
- ◆ assume we have a conservative guess of the fraction of "inliers" (correct points)

- probability of picking an uncontaminated sample

$$w^p$$

- probability of seeing no uncontaminated sample in k trials

$$z = (1 - w^p)^k$$

- required k to be $z\%$ (say, 99%) sure there is an uncontaminated sample

$$k = \frac{\log(z)}{\log(1 - w^p)}$$

Robust Estimation

◆ RANSAC (random sample consensus)

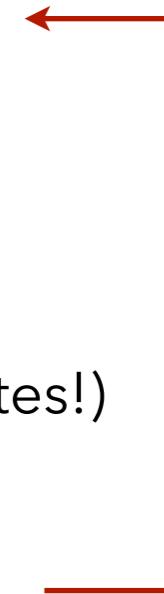
- ◆ an immensely successful brute-force method
- ◆ even the best hypothesis may be a rather bad fit
- ◆ should always refit to all inliers

p	proportion of outliers							
	5%	10%	20%	25%	30%	40%	50%	
2	2	3	5	6	7	11	17	
3	3	4	7	9	11	19	35	
4	3	5	9	13	17	34	72	
5	4	6	12	17	26	57	146	
6	4	7	16	24	37	97	293	
7	4	8	20	33	54	163	588	
8	5	9	26	44	78	272	1177	

Estimating the homography

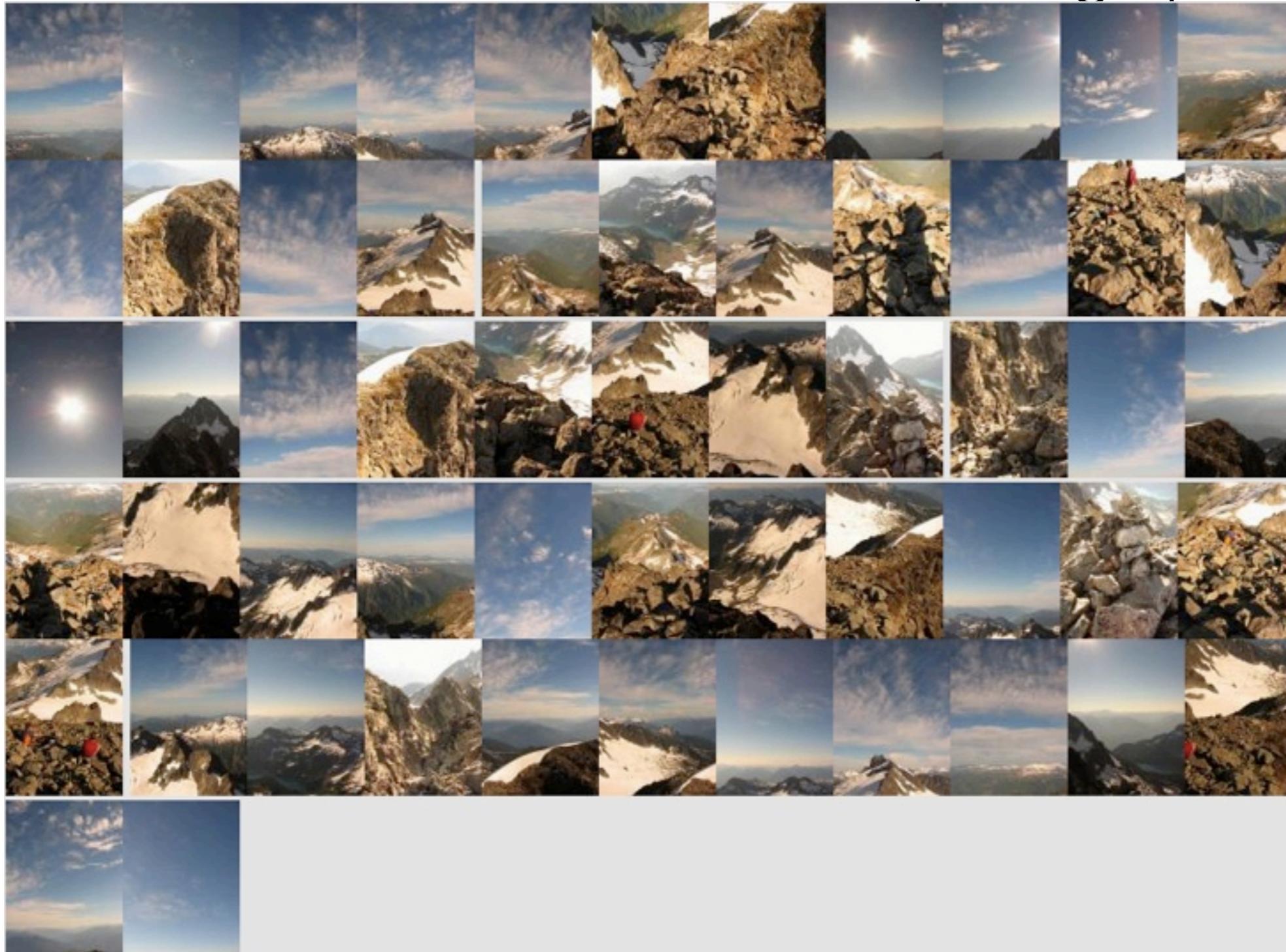
- ◆ The cookbook recipe:

- ◆ condition all image coordinates
- ◆ RANSAC loop: iterate
 - ◆ pick 4 correspondences
 - ◆ build equations, estimate homography
 - ◆ transform all points x
 - ◆ measure distance to all points x' (in non-homogeneous coordinates!)
 - ◆ count inliers (scale down threshold too!)
 - ◆ re-estimate final homography
- ◆ undo conditioning



Application: AutoStitch

- ◆ Given a collection of unordered photographs...



Application: AutoStitch

- ◆ ... automatically create a panorama
 - ◆ find images related by a homography
 - ◆ extract relative rotation
 - ◆ remap to a sphere
 - ◆ blend colors along seams



Application: AutoStitch

- ◆ ... automatically create a panorama
 - ◆ find images related by a homography
 - ◆ extract relative rotation
 - ◆ remap to a sphere
 - ◆ blend colors along seams

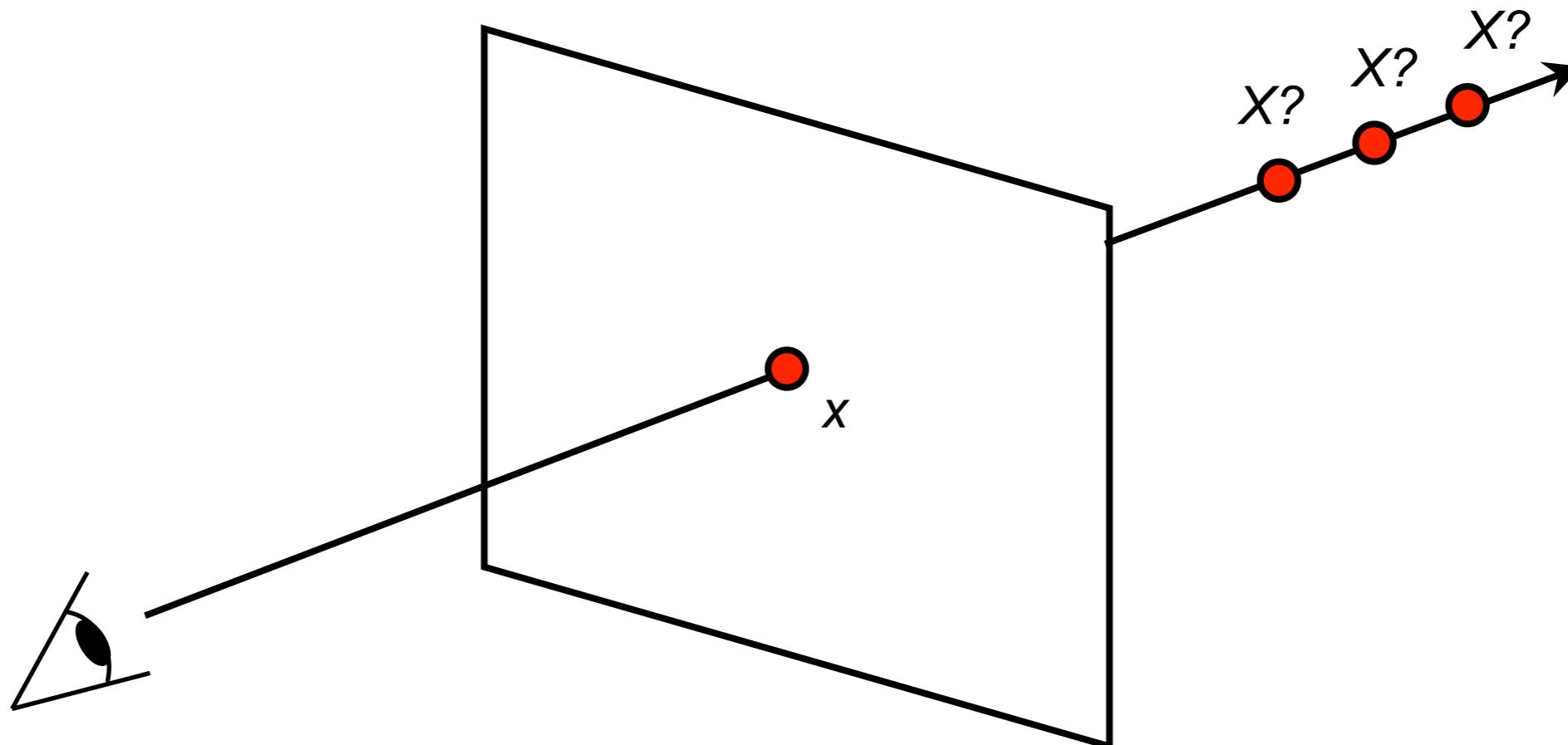


Application: AutoStitch

- ◆ ... automatically create a panorama
 - ◆ find images related by a homography
 - ◆ extract relative rotation
 - ◆ remap to a sphere
 - ◆ blend colors along seams



Our goal: Recovery of 3D structure



Our goal: Recovery of 3D structure

- ◆ Recovery of depth from one image is inherently ambiguous



Our goal: Recovery of 3D structure

- ◆ Recovery of depth from one image is inherently ambiguous



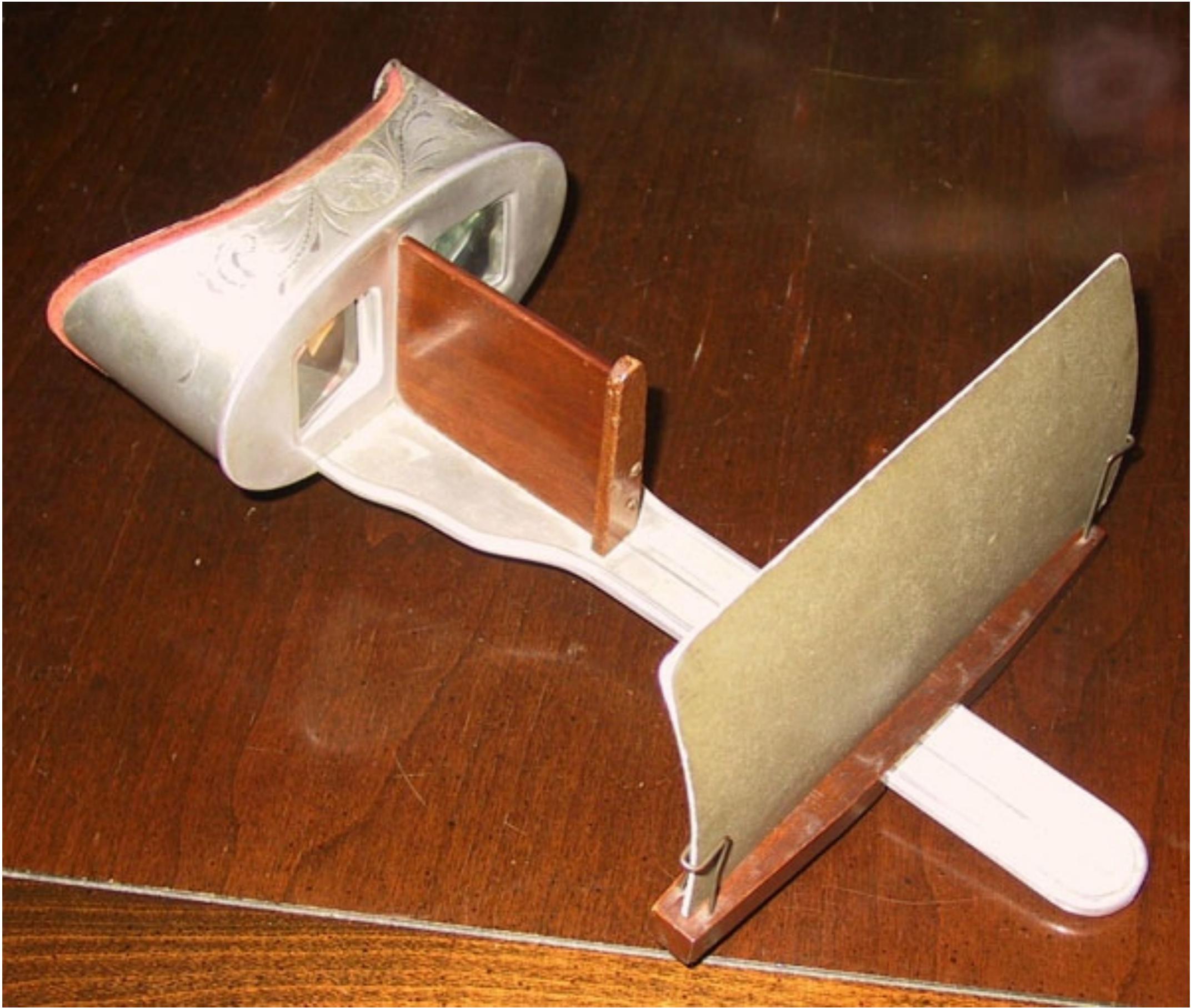
Our goal: Recovery of 3D structure

- ◆ Next: geometric methods, that rely on two views



What is Stereo (Vision)?

- ◆ Stereo vision, stereopsis, or short stereo is the perception or measurement of depth from two projections.
- ◆ Was first described in technical detail by Charles Wheatstone in the 1830s with the invention of the stereoscope.
- ◆ The human visual system heavily relies on stereo vision:
 - ◆ Our eyes give two slightly shifted projections of the scene.
 - ◆ But only relative depth can be judged accurately by humans.



[Black]



[Black]



[Black]

Left image



[Black]

Right image



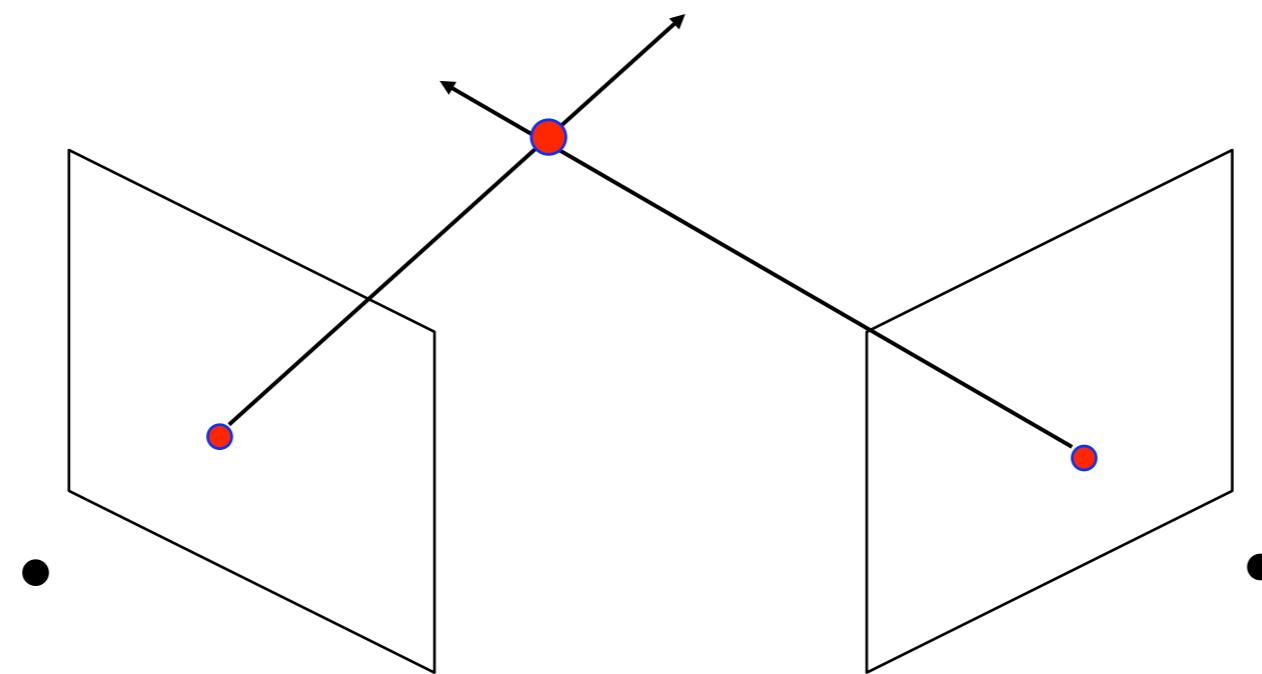
[Black]

Outline

- ◆ Part 2: two-view geometry
- ◆ Triangulation
 - ◆ 3D point reconstruction by ray intersection
 - ◆ simple geometric method
 - ◆ linear method
 - ◆ non-linear geometric method
- ◆ Application: robot SLAM

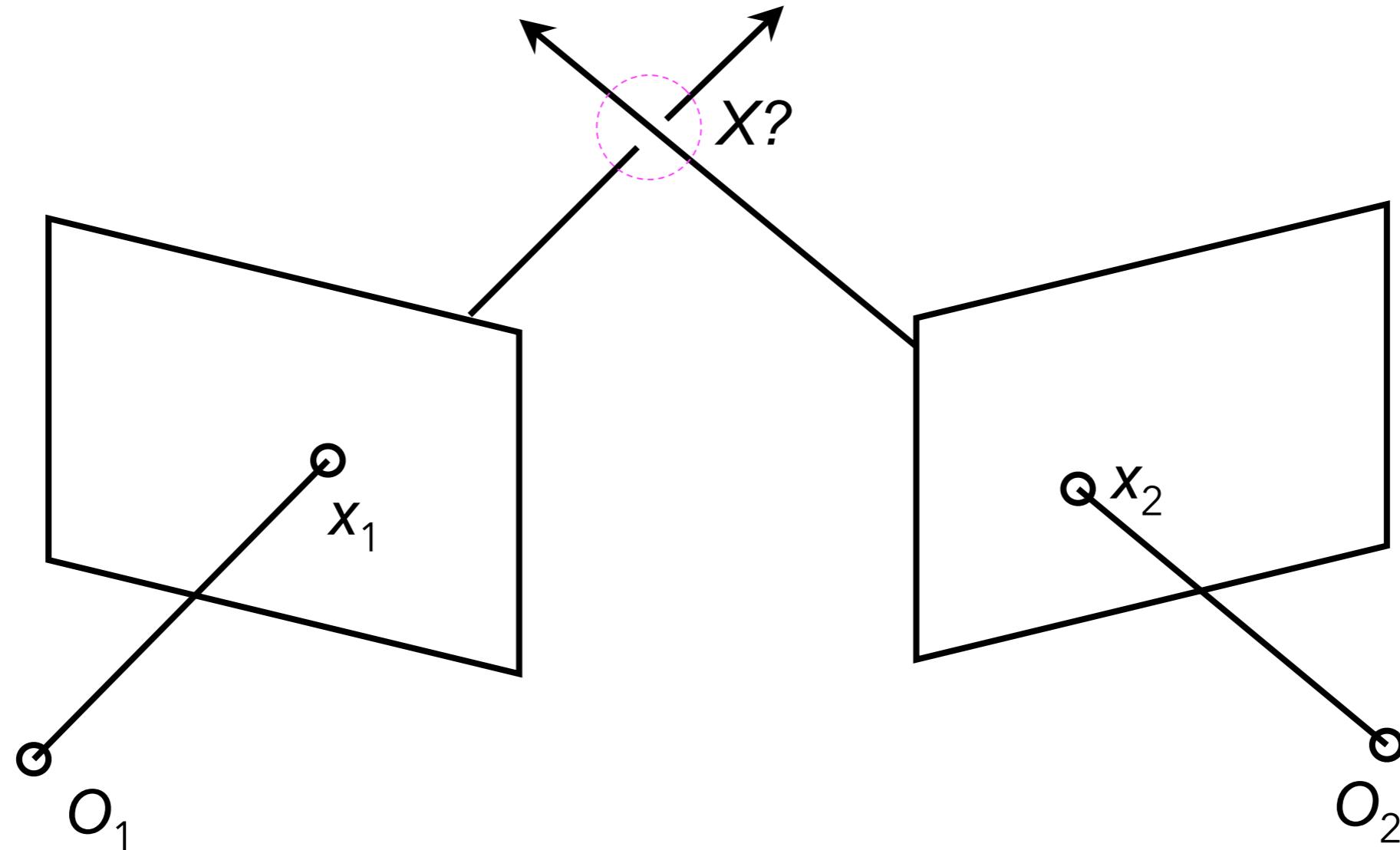
Depth from Stereo

- ◆ Triangulation
 - ◆ intersect two rays originating from the same point in the scene
 - ◆ Requires correspondences (knowledge which pixels are images of "the same point")
 - ◆ requires camera pose (in order to construct the 3D rays)



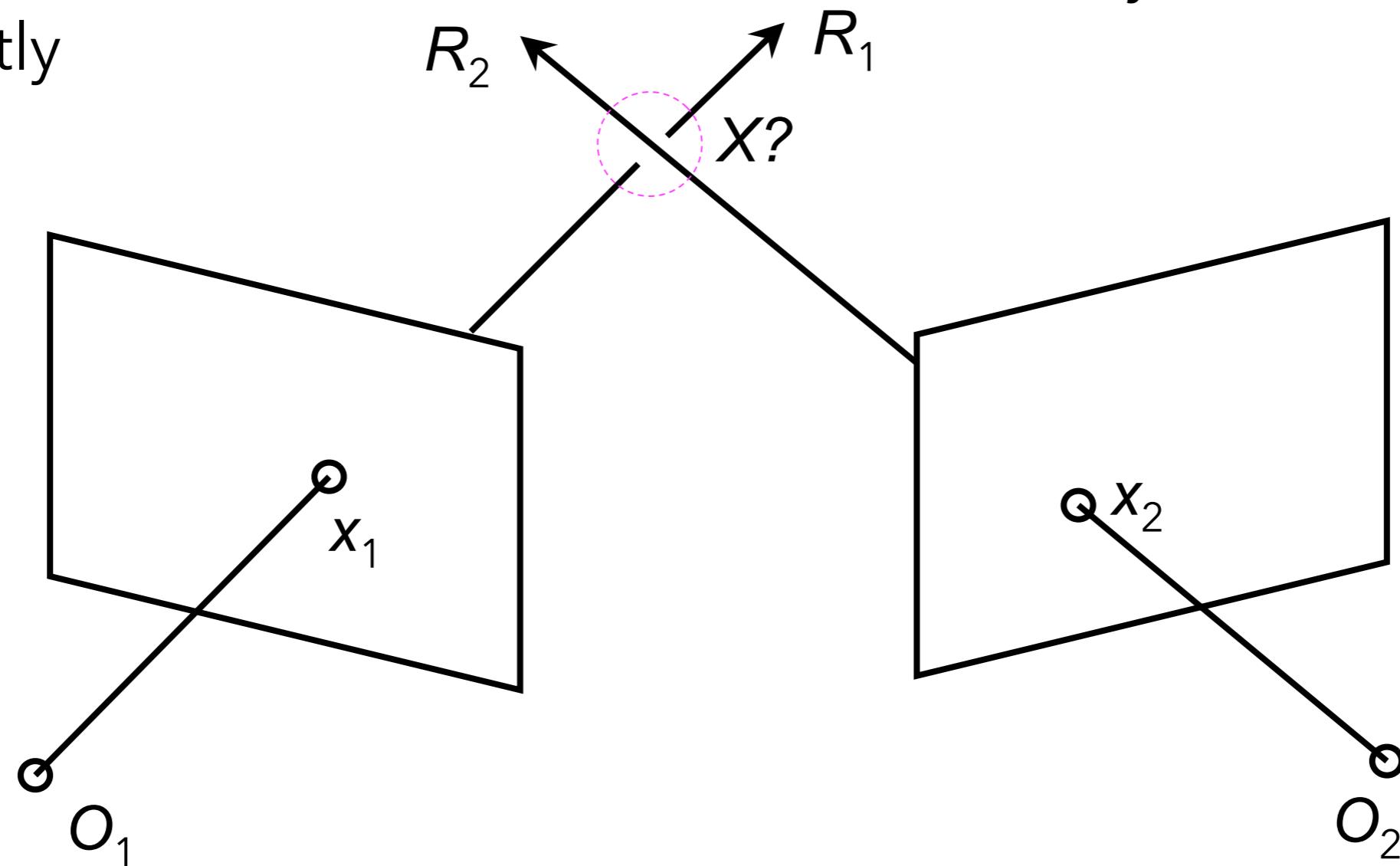
Triangulation

- Given projections of a 3D point in two or more images (with known camera matrices), find the coordinates of the point



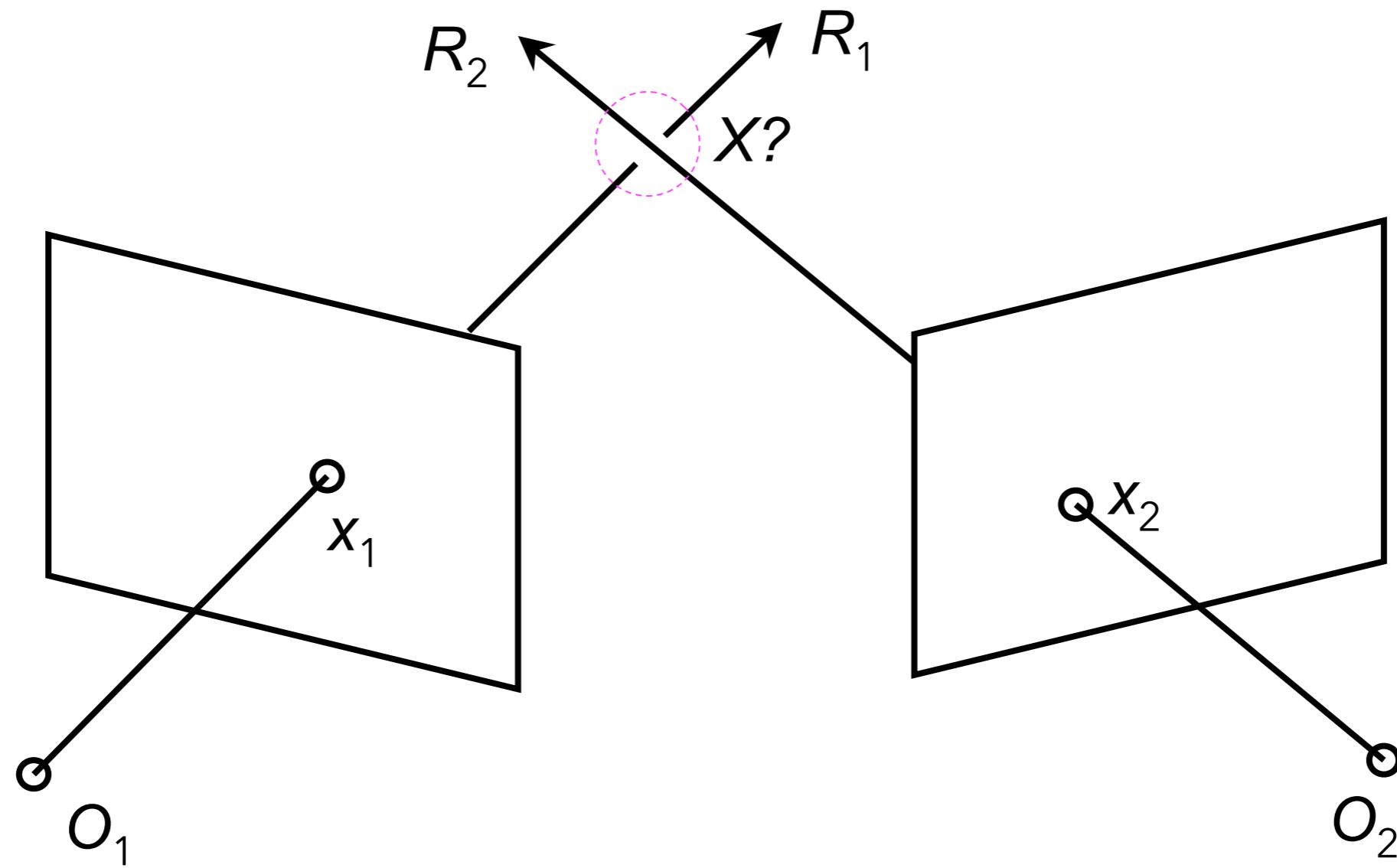
Triangulation

- ◆ We want to intersect the two viewing rays corresponding to x_1 and x_2
- ◆ Because of noise and numerical errors, they do not meet exactly



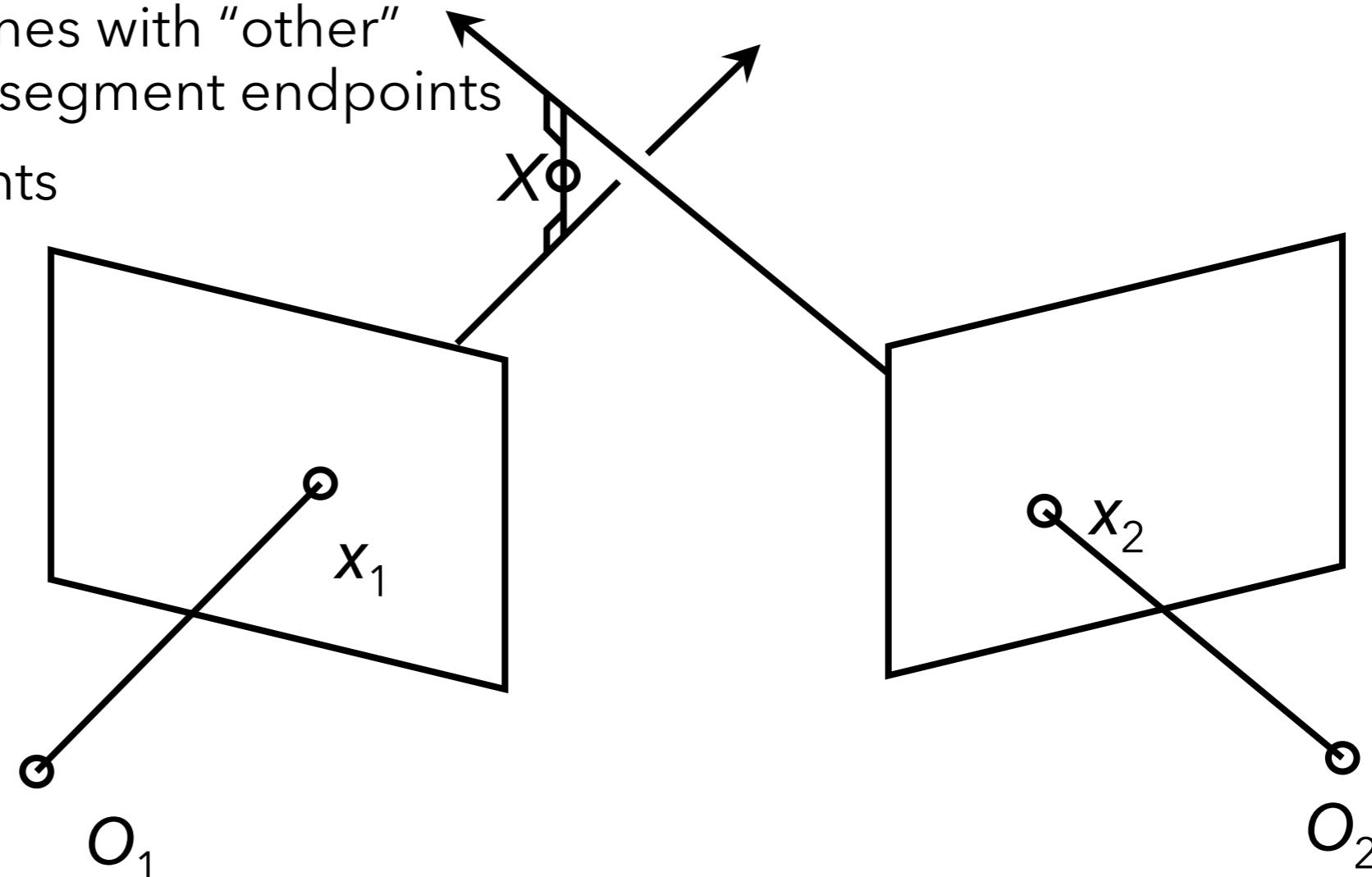
Triangulation

- ◆ Ideas?



Triangulation: geometric midpoint

- ◆ Find shortest segment connecting the two viewing rays and let X be the midpoint of that segment
 - ◆ find the segment direction (normal to both rays)
 - ◆ construct 2 planes each containing the segment and one ray
 - ◆ intersect planes with “other” rays to yield segment endpoints
 - ◆ average points



Triangulation: Linear approach

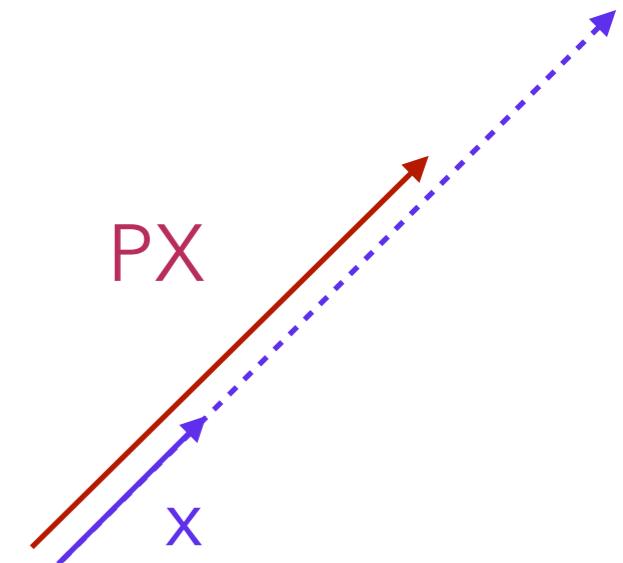
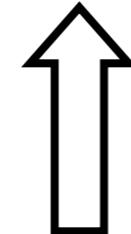


$$\lambda_1 \mathbf{x}_1 = \mathbf{P}_1 \mathbf{X}$$

$$\mathbf{x}_1 \times \mathbf{P}_1 \mathbf{X} = 0$$

$$\lambda_2 \mathbf{x}_2 = \mathbf{P}_2 \mathbf{X}$$

$$\mathbf{x}_2 \times \mathbf{P}_2 \mathbf{X} = 0$$



- ◆ Two independent equations per image point for three unknown entries of \mathbf{X}
- ◆ ... again, a homogeneous linear equation system
- ◆ ... again, solve with SVD
- ◆ note: directly generalizes to >2 views, just stack more equations

Side note

- ◆ An algebraic trick (useful here and later):
- ◆ how to get there without expanding the cross-product and re-ordering: cross-product can be written as a matrix-vector multiplication

$$x \times PX = 0$$

$$a \times b = 0$$

$$[a_{\times}]b = 0$$

$$[a_{\times}] = \begin{bmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{bmatrix}$$

$$[x_{\times}]P X = 0$$

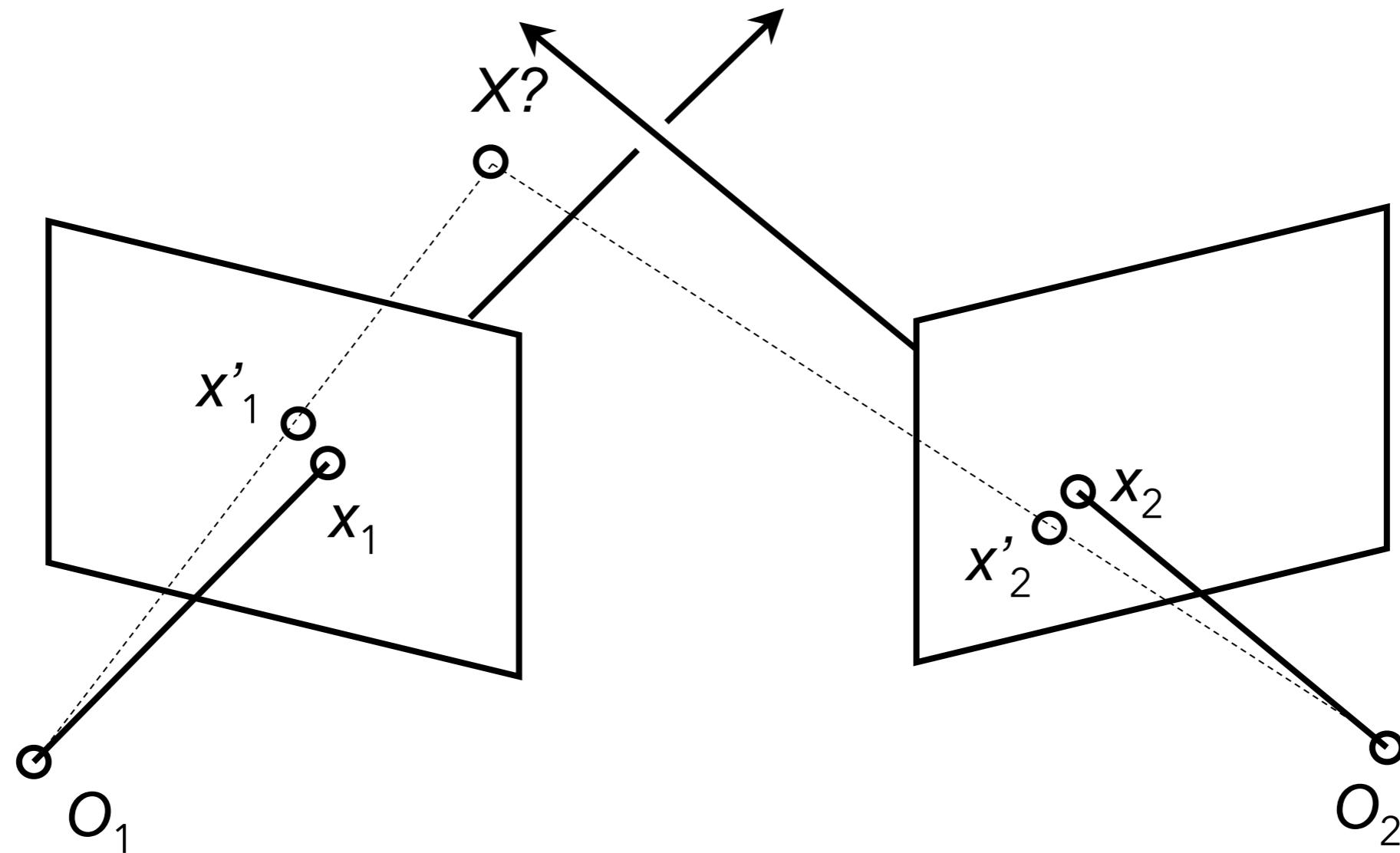


a matrix

Triangulation: non-linear approach

- ◆ Find X that minimizes the (squared) reprojection error:

$$d^2(x_1, P_1 X) + d^2(x_2, P_2 X)$$



Triangulation: non-linear approach

- ◆ Find X that minimizes the (squared) reprojection error:

$$d^2(x_1, P_1 X) + d^2(x_2, P_2 X)$$

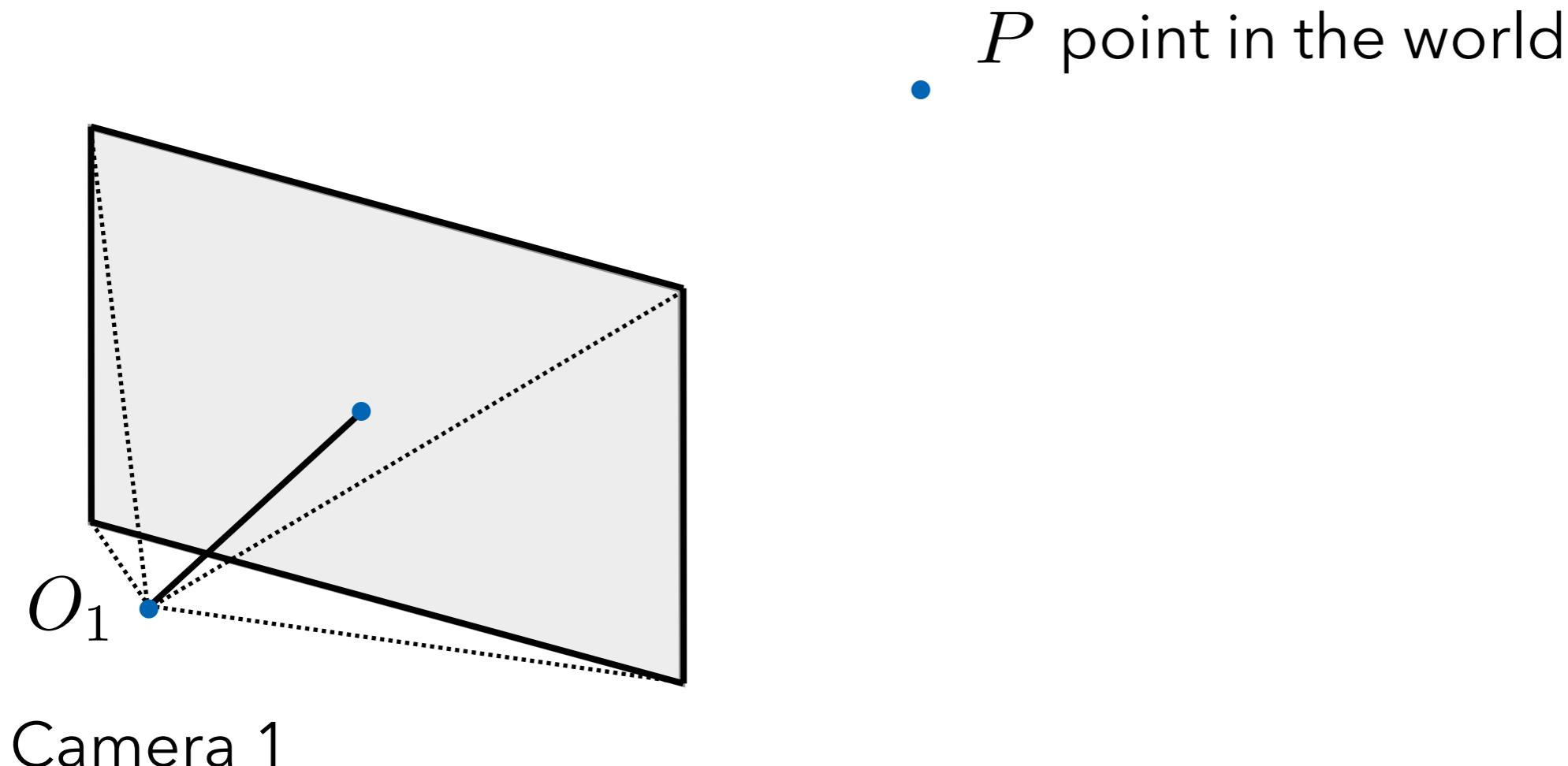
- ◆ The most accurate method, but is more complex than the other two
- ◆ With 2 cameras
 - ◆ find roots of a 6th degree polynomial (Hartley & Zisserman, chapter 12)
- ◆ With >2 cameras
 - ◆ initialize with linear estimate
 - ◆ optimize with iterative methods (Gauss-Newton, Levenberg-Marquardt - HZ, appendix 6)



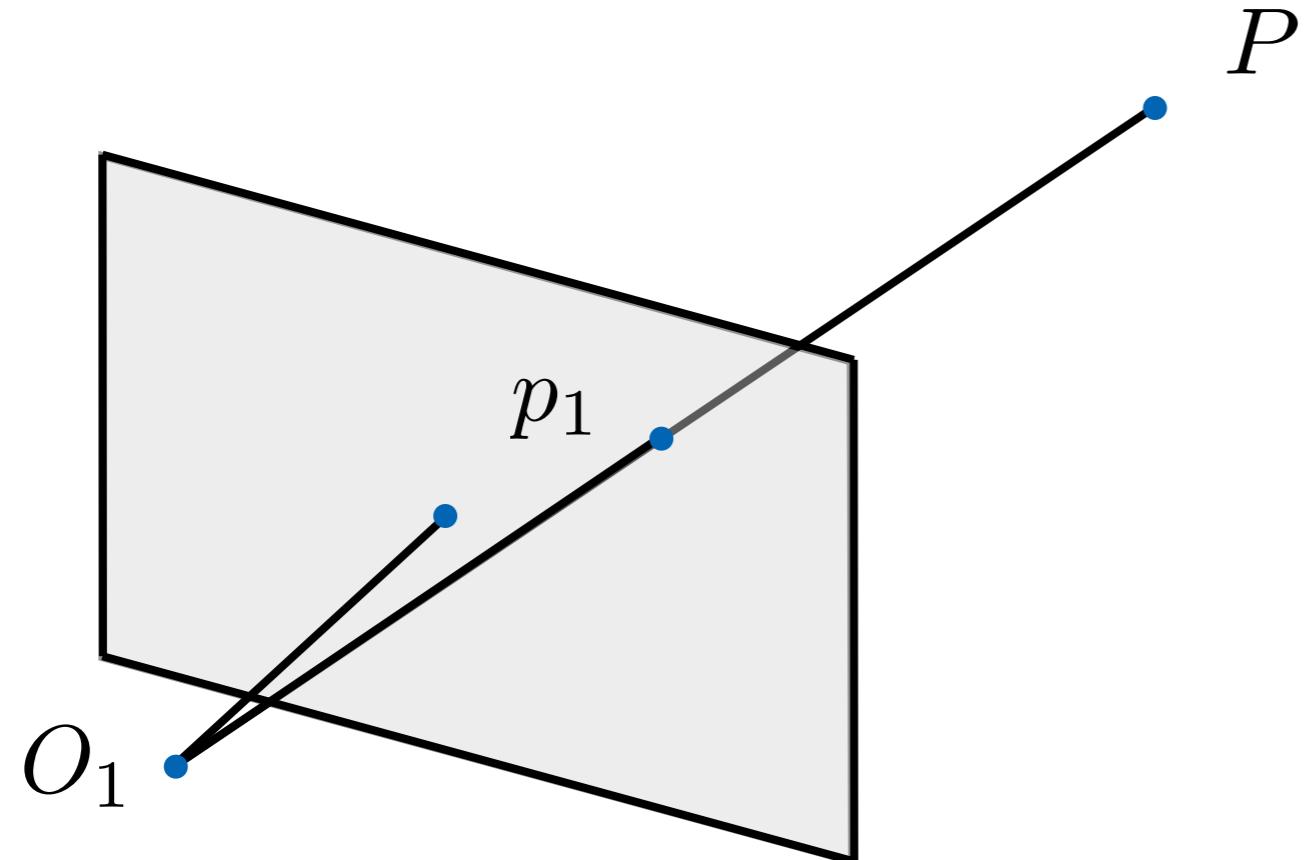
Outline

- ◆ Goal: find the **relation** between two different views of the same scene
 - ◆ recover camera placement
 - ◆ matching and triangulation for multiple views
 - ◆ scene modeling only from images
- ◆ **Geometrical relations**
 - ◆ epipolar constraint
 - ◆ epipolar geometry of 2 views
- ◆ **Algebraic relations**
 - ◆ calibrated cameras: essential matrix
 - ◆ uncalibrated cameras: fundamental matrix
 - ◆ the normalized 8-point algorithm
 - ◆ relation to camera matrices

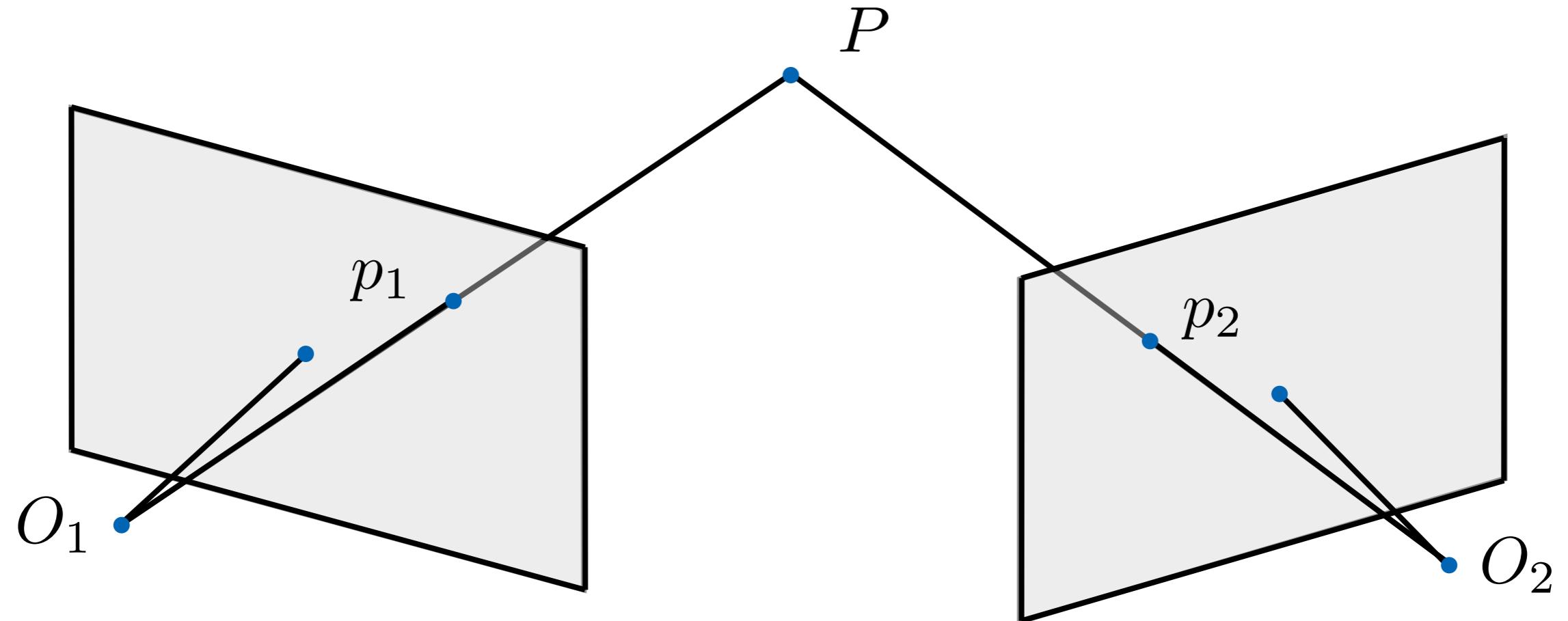
Epipolar Geometry



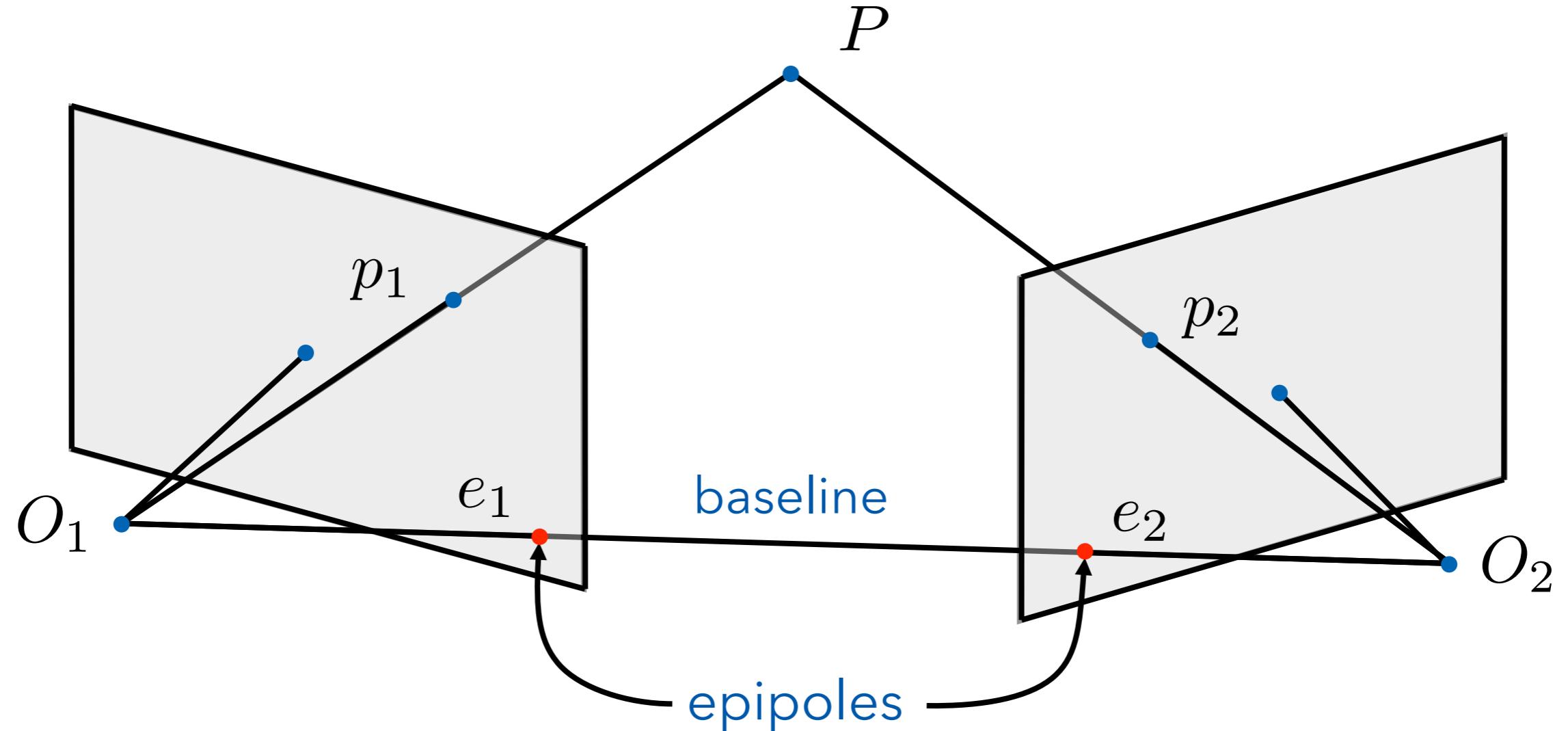
Epipolar Geometry



Epipolar Geometry

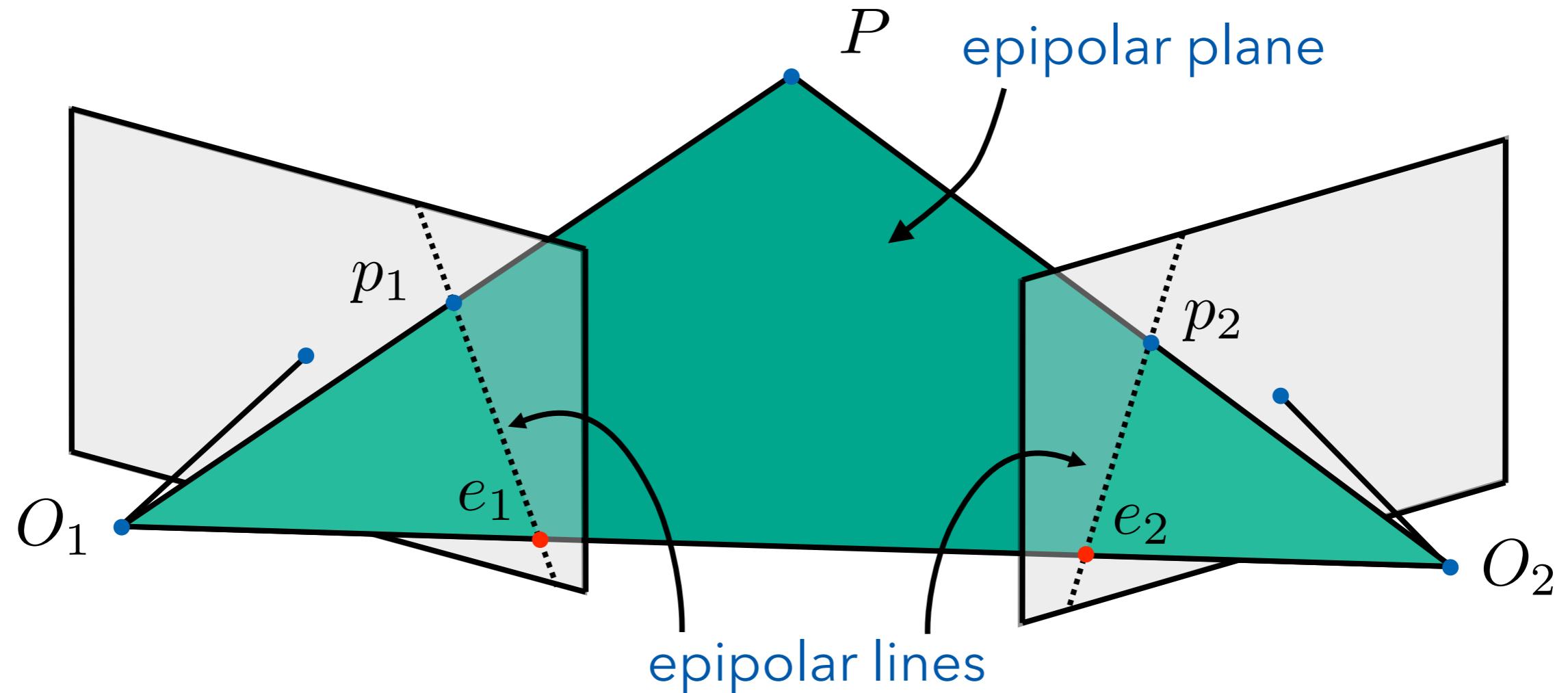


Epipolar Geometry



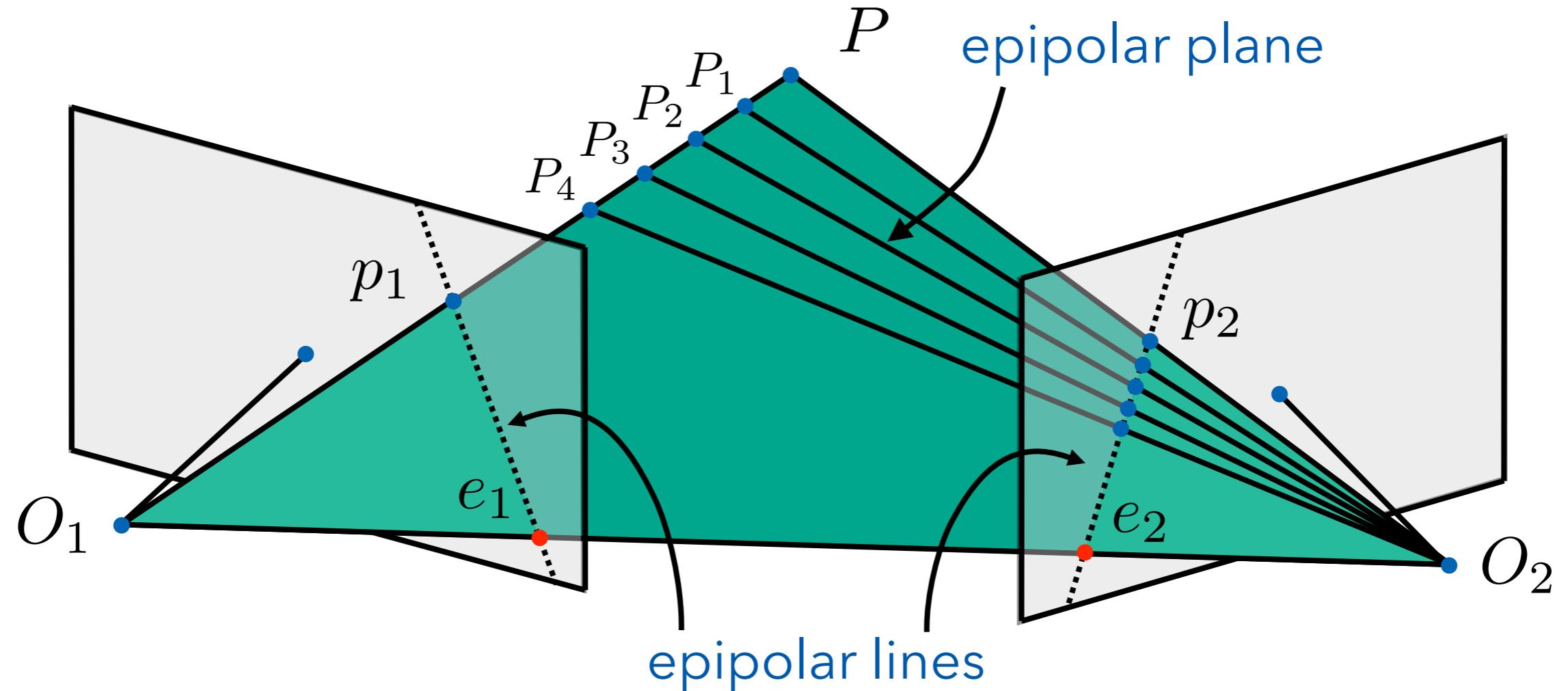
- ◆ Epipole: Image location of the optical center of the other camera.
 - ◆ Can be outside of the visible area.

Epipolar Geometry



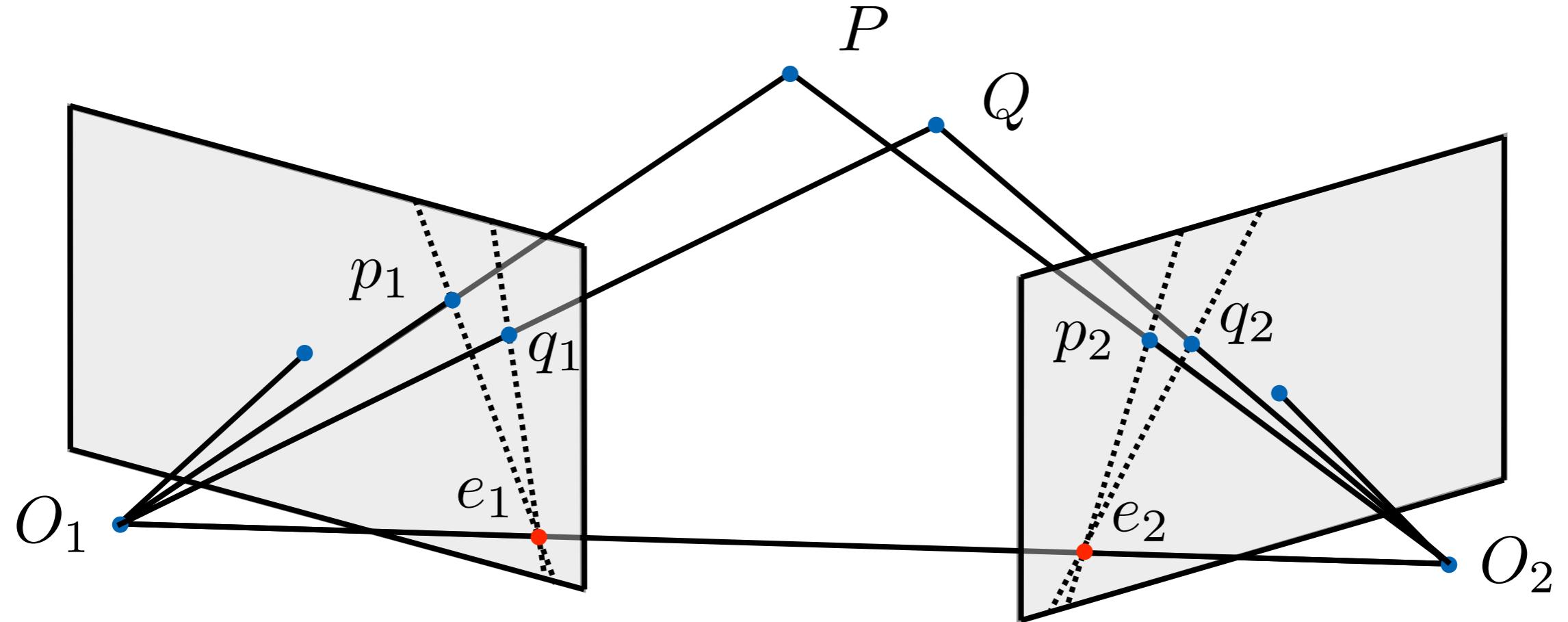
- ◆ Epipolar plane: Plane through both camera centers and world point.

Epipolar Geometry



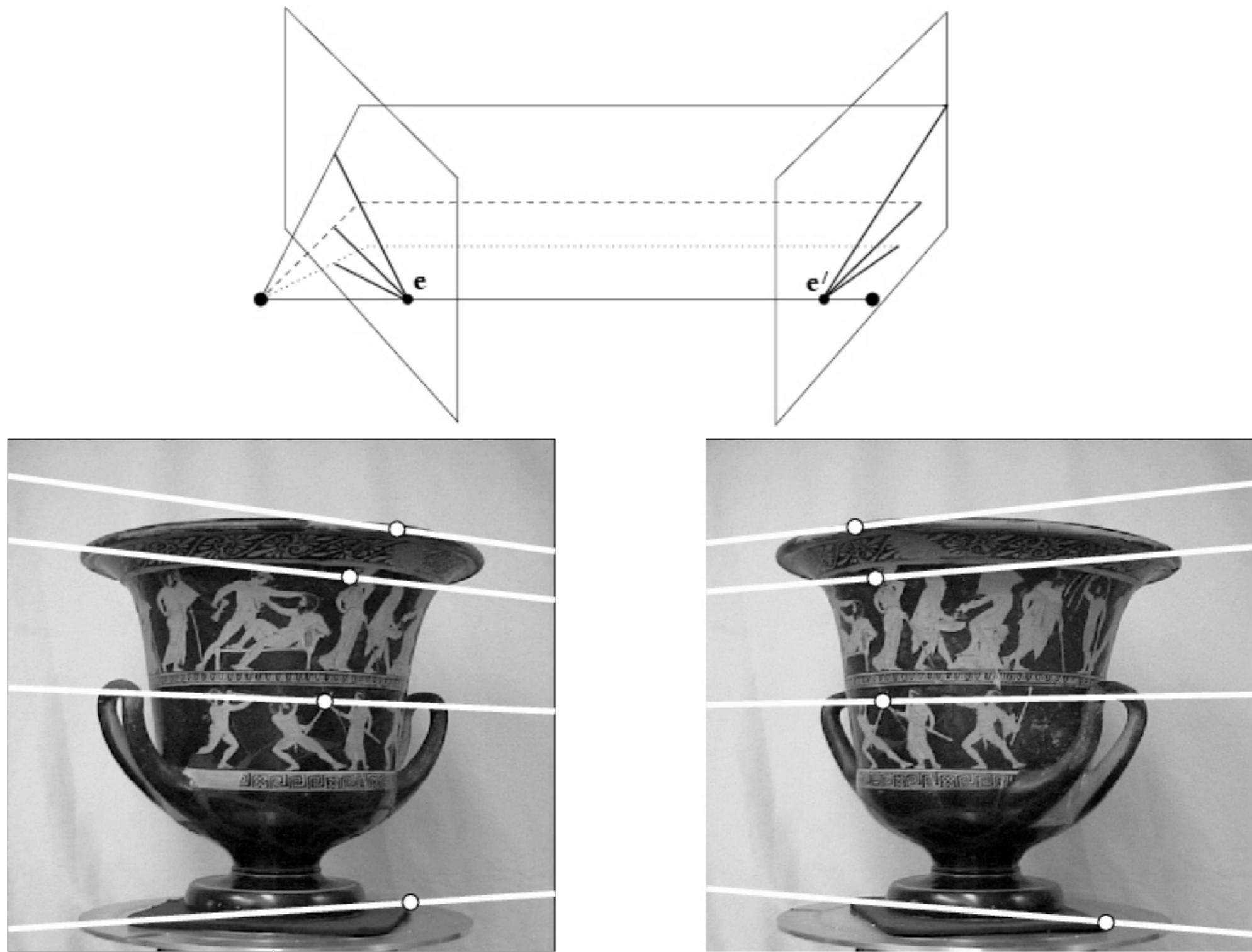
- ◆ Epipolar line: Constrains the location where a particular feature from one view can be found in the other.
- ◆ Here: Possible matches for p_1

Epipolar Geometry

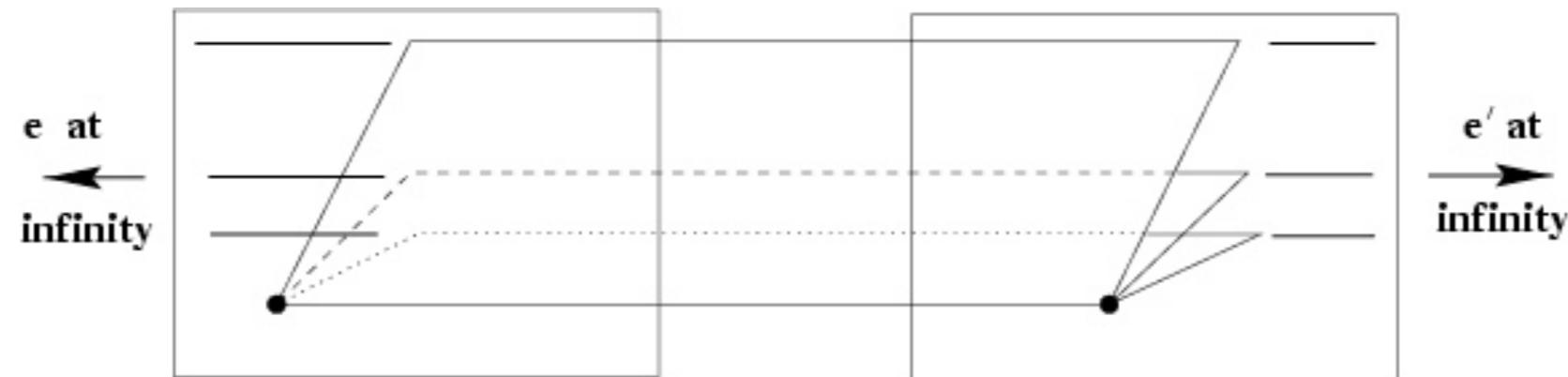


- ◆ Epipolar lines:
 - ◆ Intersect at the epipoles.
 - ◆ In general not parallel.

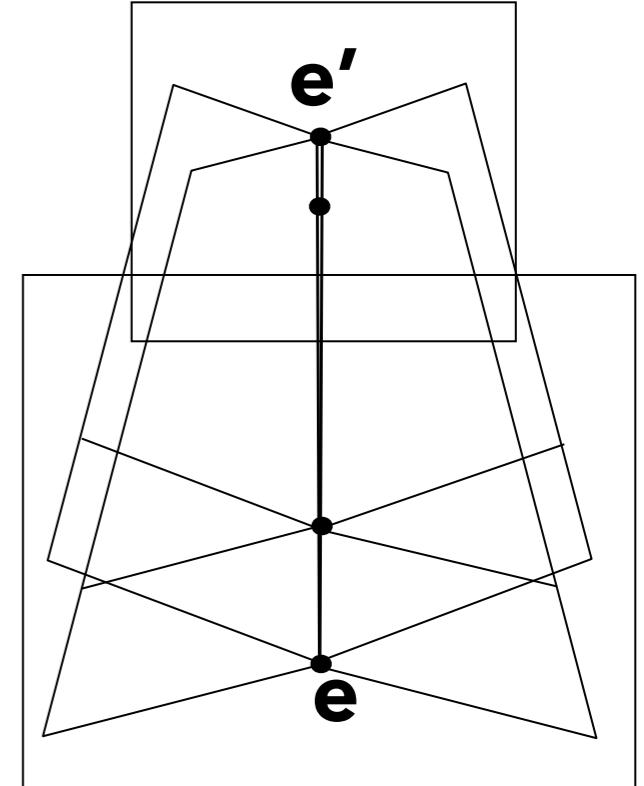
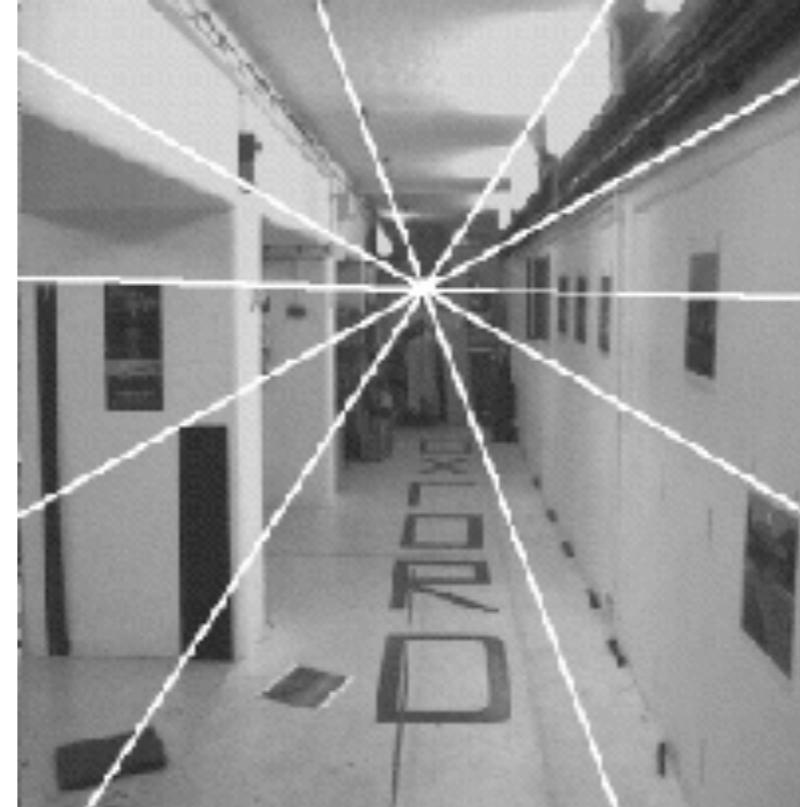
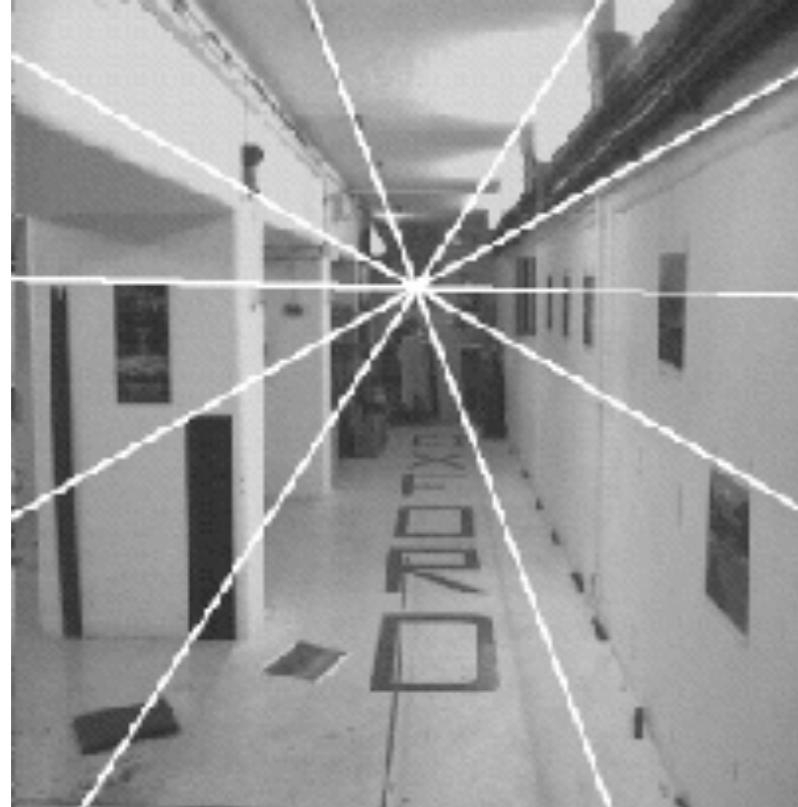
Example: Converging Cameras



Example: Motion Parallel to Image Plane

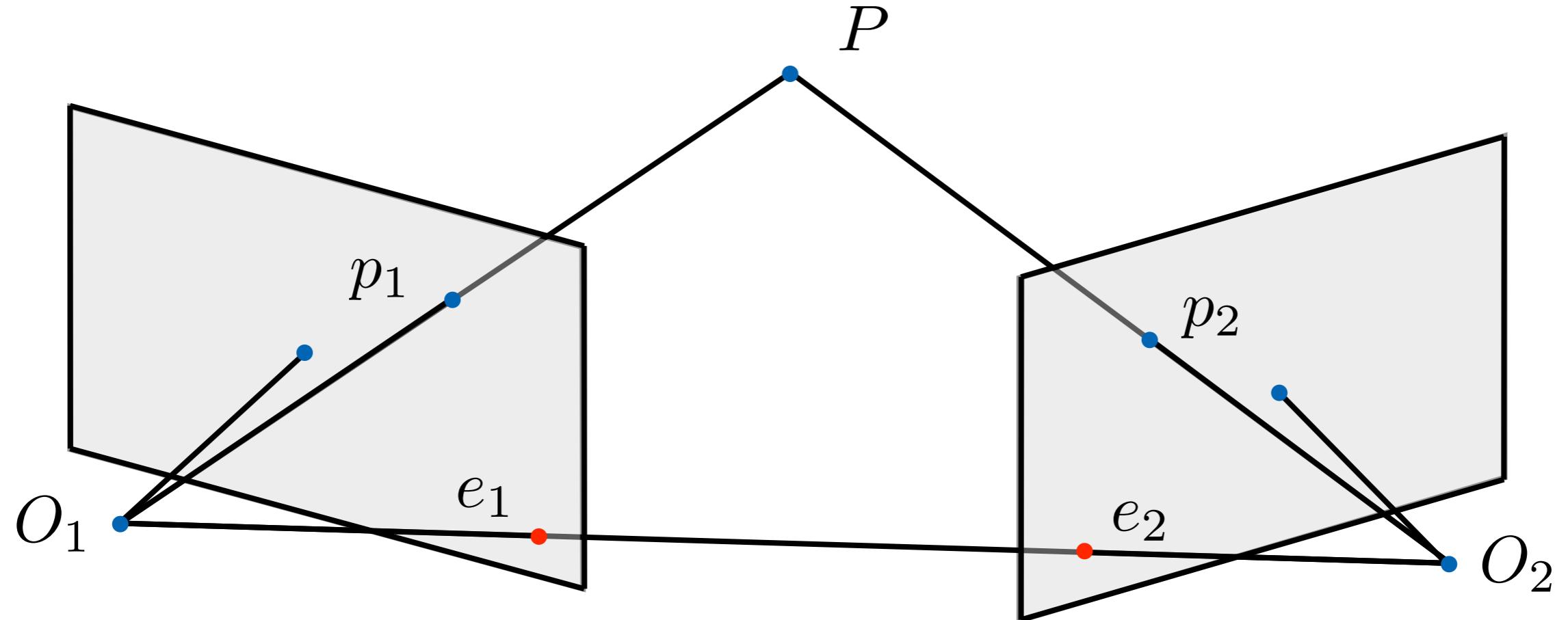


Example: Forward Motion



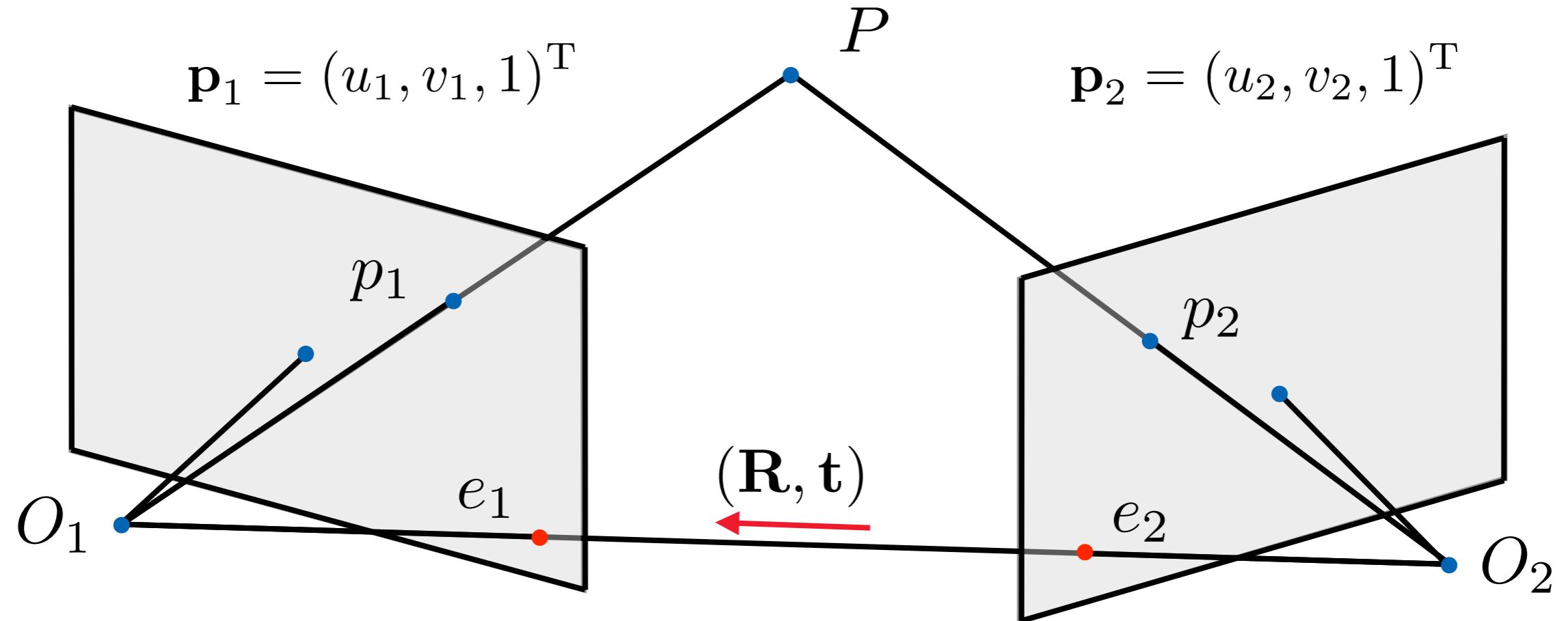
- ◆ Epipole has same coordinate in both images
- ◆ Point move along lines radiating from epipole
 - “focus of expansion”

Epipolar Geometry



- ◆ Mathematically: $\overrightarrow{O_1P}, \overrightarrow{O_2P}, \overrightarrow{O_1O_2}$ are coplanar.
- ◆ Equivalently: $\overrightarrow{O_1p_1}, \overrightarrow{O_2p_2}, \overrightarrow{O_1O_2}$ are coplanar.
- ◆ In other terms: $\overrightarrow{O_1p_1} \cdot [\overrightarrow{O_1O_2} \times \overrightarrow{O_2p_2}] = 0$

Epipolar Geometry



- ◆ Rewrite in camera coordinate system of camera 1:

$$\overrightarrow{O_1 p_1} \cdot [\overrightarrow{O_1 O_2} \times \overrightarrow{O_2 p_2}] = 0$$
$$\Leftrightarrow \mathbf{p}_1^T [\mathbf{t} \times (\mathbf{R} \mathbf{p}_2)] = 0$$

Epipolar Geometry

- ◆ Epipolar constraint: $\mathbf{p}_1^T [\mathbf{t} \times (\mathbf{R}\mathbf{p}_2)] = 0$
- ◆ Express cross-product as matrix: $\mathbf{t} \times \mathbf{v} = [\mathbf{t}]_{\times} \mathbf{v}$
- ◆ **Essential matrix:** $\mathbf{E} = [\mathbf{t}]_{\times} \mathbf{R}$
- ◆ captures the relation of two views [Longuet-Higgins, 1981]
- ◆ Epipolar constraint with essential matrix:

$$\begin{aligned}
 0 &= \mathbf{p}_1^T [\mathbf{t} \times (\mathbf{R}\mathbf{p}_2)] \\
 &= \mathbf{p}_1^T [([\mathbf{t}]_{\times} \mathbf{R})\mathbf{p}_2] \\
 &= \boxed{\mathbf{p}_1^T \mathbf{E} \mathbf{p}_2}
 \end{aligned}$$

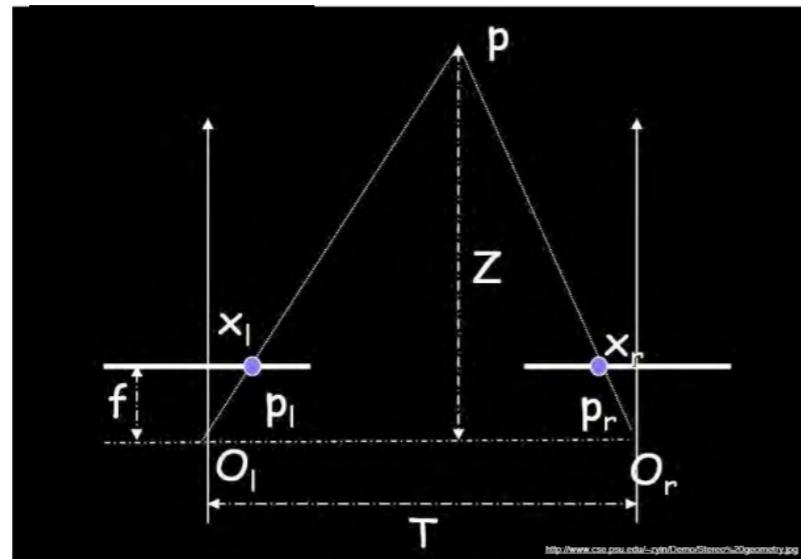
Epipolar Geometry

- ◆ Two views of the same 3D point must satisfy $\mathbf{p}_1^T \mathbf{E} \mathbf{p}_2 = 0$
- ◆ Other important properties...
- ◆ Epipolar lines: $\mathbf{l}_1 = \mathbf{E} \mathbf{p}_2$ and $\mathbf{l}_2 = \mathbf{E}^T \mathbf{p}_1$
- ◆ Epipoles are the (left/right) null-space of the essential matrix:

$$\mathbf{e}_1^T \mathbf{E} = \mathbf{E}^T \mathbf{e}_1 = 0 \quad \mathbf{E} \mathbf{e}_2 = 0$$

- ◆ Essential matrix is singular; has rank 2.
- ◆ The two remaining eigenvalues are equal.
- ◆ 5 degrees of freedom (translation + rotation have 6, but scale is arbitrary)

Binocular Stereo: Parallel cameras



$$\mathbf{R} = \mathbf{I}$$

$$\mathbf{t} = [-b, 0, 0]$$

$$\mathbf{E} = [\mathbf{t}]_{\times} \mathbf{R} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & b \\ 0 & -b & 0 \end{bmatrix}$$

$$0 = \mathbf{p}_1^T \mathbf{E} \mathbf{p}_2$$

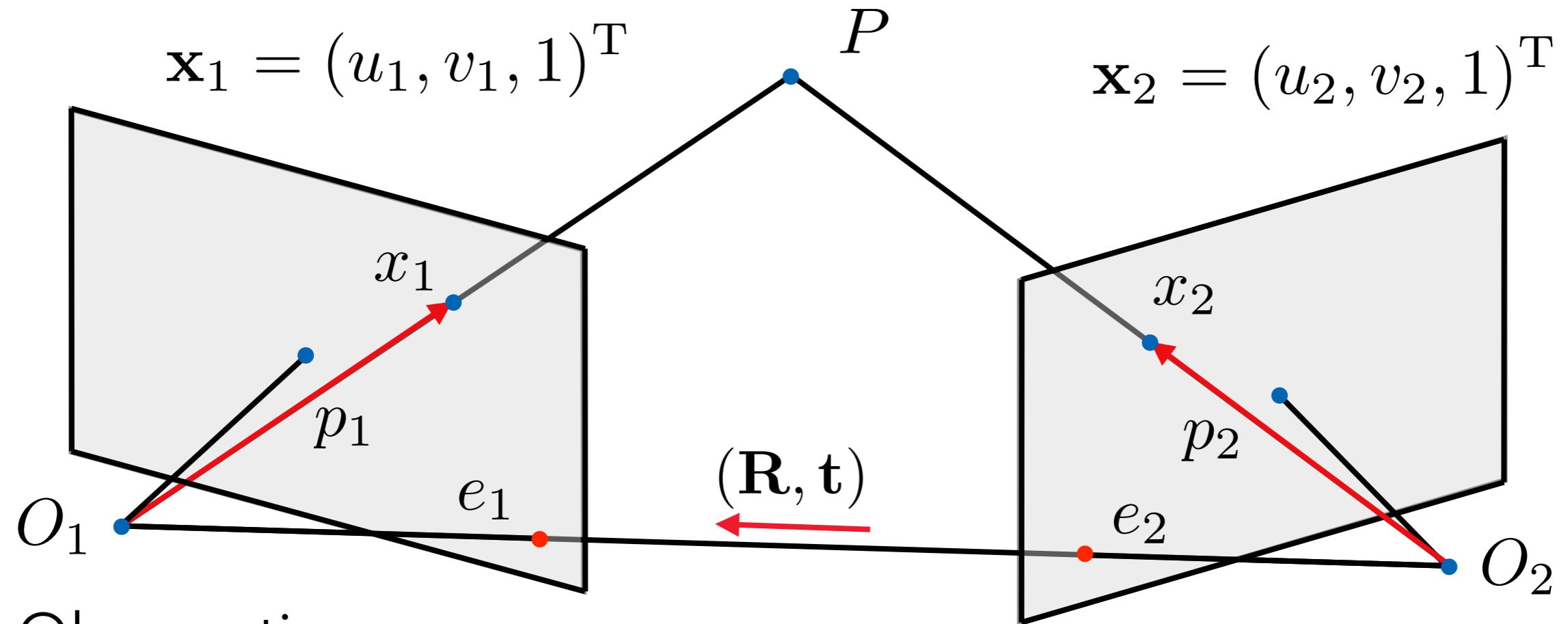
$$= [x_1, y_1, 1] \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & b \\ 0 & -b & 0 \end{bmatrix} \begin{bmatrix} x_2 \\ y_2 \\ 1 \end{bmatrix}$$

$$= [x_1, y_1, 1] \begin{bmatrix} 0 \\ b \\ -by_2 \end{bmatrix} = by_1 - by_2$$



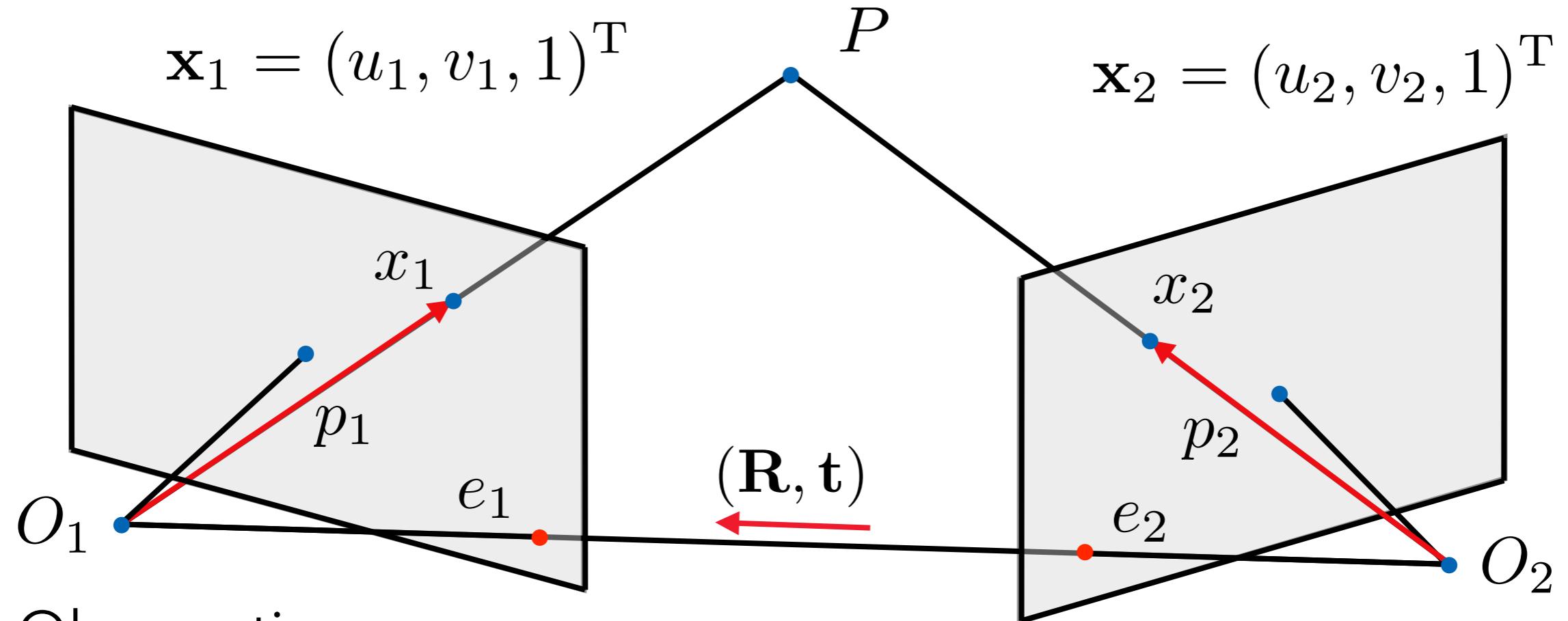
$y_1 = y_2$
 x_2 arbitrary

Uncalibrated Cameras



- ◆ Observation:
 - ◆ the epipolar constraint is derived by intersecting lines and planes
 - ◆ these lines and planes exist for any two pinhole cameras
 - ◆ if we do not know the calibration, the constraint must still hold

Uncalibrated Cameras



- ◆ Observation:

- ◆ The transformation from rays \mathbf{p} to image points \mathbf{x} is the calibration \mathbf{K}
 - ◆ The transformation exists and is invertible, even if it is unknown

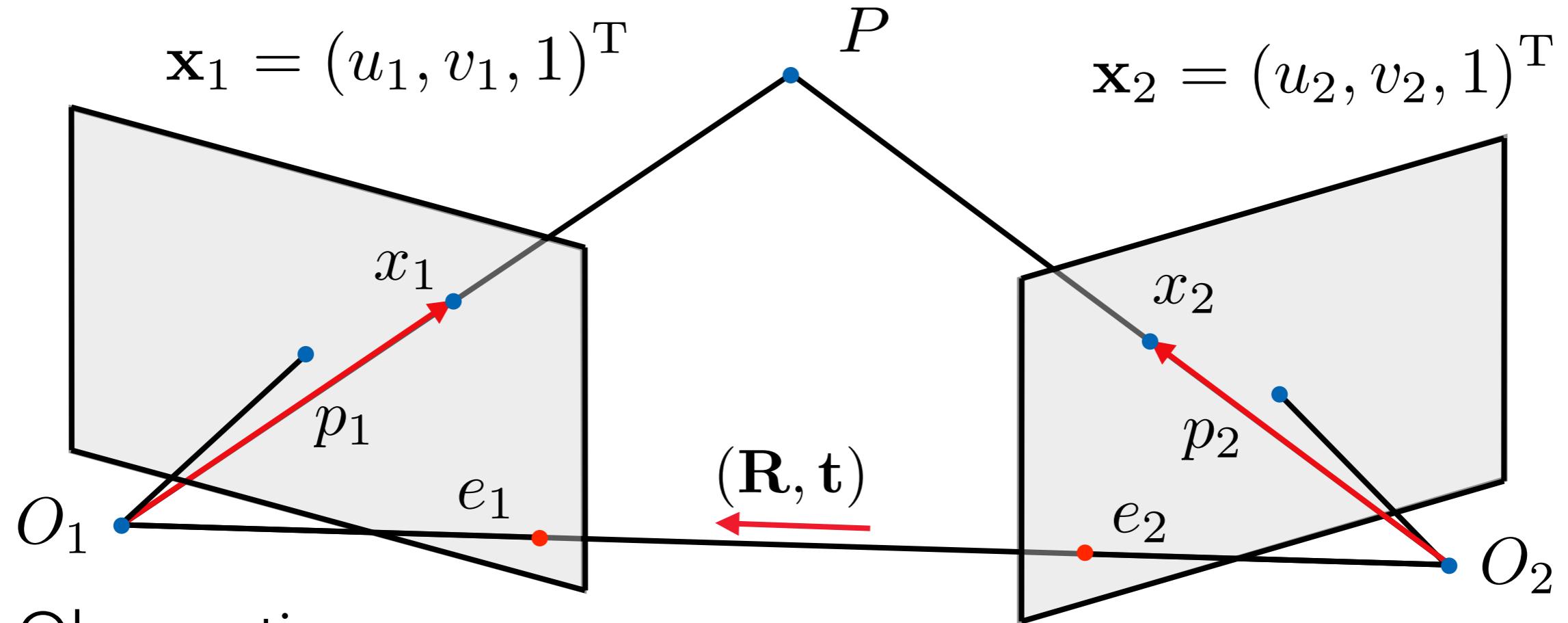
$$\mathbf{x}_1 = \mathbf{K}_1 \mathbf{p}_1$$

$$\mathbf{p}_1 = \mathbf{K}_1^{-1} \mathbf{x}_1$$

$$\mathbf{x}_2 = \mathbf{K}_2 \mathbf{p}_2$$

$$\mathbf{p}_2 = \mathbf{K}_2^{-1} \mathbf{x}_2$$

Uncalibrated Cameras



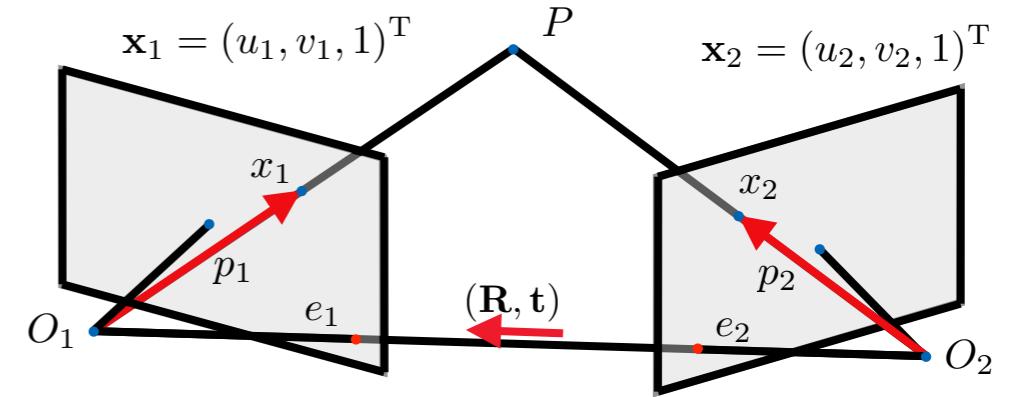
- ◆ Observation:
 - ◆ We silently assumed that everything was done in world coordinates
 - ◆ We can simply plug in the pixel coordinates:

$$0 = \mathbf{p}_1^T \mathbf{E} \mathbf{p}_2 = \mathbf{x}_1^T \mathbf{K}_1^{-T} \mathbf{E} \mathbf{K}_2^{-1} \mathbf{x}_2$$

$$0 = \mathbf{x}_1^T \mathbf{F} \mathbf{x}_2$$

Fundamental matrix

$$\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2 = 0$$



- ◆ holds for calibrated and uncalibrated cameras
- ◆ geometrically: two rays must be coplanar, independent of any specific \mathbf{K} , \mathbf{R} , \mathbf{t}
- ◆ \mathbf{F} is called the fundamental matrix, and encapsulates the relative geometry of the two uncalibrated views
[Faugeras and Luong 1992, Hartley 1992]

Fundamental matrix

- ◆ Two image coordinates of the same 3D point must satisfy

$$\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2 = 0$$

- ◆ Other important properties...
- ◆ Epipolar lines: $\mathbf{l}_1 = \mathbf{F} \mathbf{x}_2$ and $\mathbf{l}_2 = \mathbf{F}^T \mathbf{x}_1$
- ◆ Epipoles are the (left/right) null-space of the fundamental matrix:

$$\mathbf{e}_1^T \mathbf{F} = \mathbf{F}^T \mathbf{e}_1 = 0 \quad \mathbf{F} \mathbf{e}_2 = 0$$

- ◆ Fundamental matrix is singular; has rank 2.
- ◆ 7 degrees of freedom

Estimating the fundamental matrix

- ◆ Seven degrees of freedom
- ◆ One equation per correspondence
- ◆ Need at least 7 pairs of corresponding points, but the solution is non-linear
- ◆ Linear solution with at least 8 point pairs: the normalized eight-point algorithm



Eight-point algorithm

- ◆ Stack equations from at least 8 correspondences
- ◆ Remember
 - ◆ a homogeneous linear equation system
 - ◆ solve s.t. $\| f \| = 1$,
using SVD

$$\begin{bmatrix} x' & y' & 1 \end{bmatrix} \begin{bmatrix} F_{11} & F_{12} & F_{13} \\ F_{21} & F_{22} & F_{23} \\ F_{31} & F_{32} & F_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = 0$$

$$\begin{bmatrix} xx' & yx' & x' & xy' & yy' & y' & x & y & 1 \end{bmatrix} \begin{bmatrix} F_{11} \\ F_{12} \\ F_{13} \\ F_{21} \\ F_{22} \\ F_{23} \\ F_{31} \\ F_{32} \\ F_{33} \end{bmatrix} = 0$$

$$Af = 0$$

Eight-point algorithm

- ◆ Remember: numerical stability
 - ◆ Coefficients of an equation system should be in the same order of magnitude
 - ◆ in pixels: $xx' \sim 1E6$
 - ◆ Conditioning: scale and shift points to be in [-1..1]
 - ◆ apply inverse transform to estimated matrix

$$s = \frac{1}{2} \max_i (\|x_i\|)$$

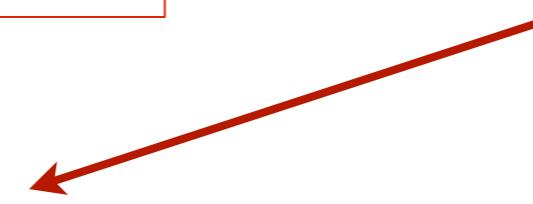
$$t = \text{mean}(x_i)$$



$$T = \begin{bmatrix} \frac{1}{s} & 0 & -\frac{t_x}{s} \\ 0 & \frac{1}{s} & -\frac{t_y}{s} \\ 0 & 0 & 1 \end{bmatrix}$$

$$u = Tx$$

$$u' = T'x'$$



$$u'^\top \bar{F} u = 0$$

Eight-point algorithm

- ◆ Problem: the result should be of rank(2)
- ◆ Find the most similar rank(2)-matrix, i.e. force the smallest eigenvalue to be 0

$$A\bar{f} = 0$$

SVD ↓ to ensure unit scale

$$U_A D_A V_A^\top = A$$

$$\tilde{F} = \begin{bmatrix} v_{19} & v_{29} & v_{39} \\ v_{49} & v_{59} & v_{69} \\ v_{79} & v_{89} & v_{99} \end{bmatrix}$$

rank(3)

SVD

$$U_F \tilde{D}_F V_F^\top = \tilde{F}$$

$$D_F = \tilde{D}_F \quad , \quad D_{F,33} = 0$$

rank(2)

$$\bar{F} = U_F D_F V_F^\top$$

un-conditioning

$$F = T'^\top \bar{F} T$$



Estimating the fundamental matrix

- ◆ Calibration not required, hence can deal with
 - ◆ archival images
 - ◆ photos not taken for reconstruction purposes
 - ◆ varying intrinsics
- ◆ So, where to start?
- ◆ The cookbook recipe
 - ◆ find interest points in both images
 - ◆ match them without epipolar constraint
 - ◆ compute epipolar geometry
 - ◆ ... refine

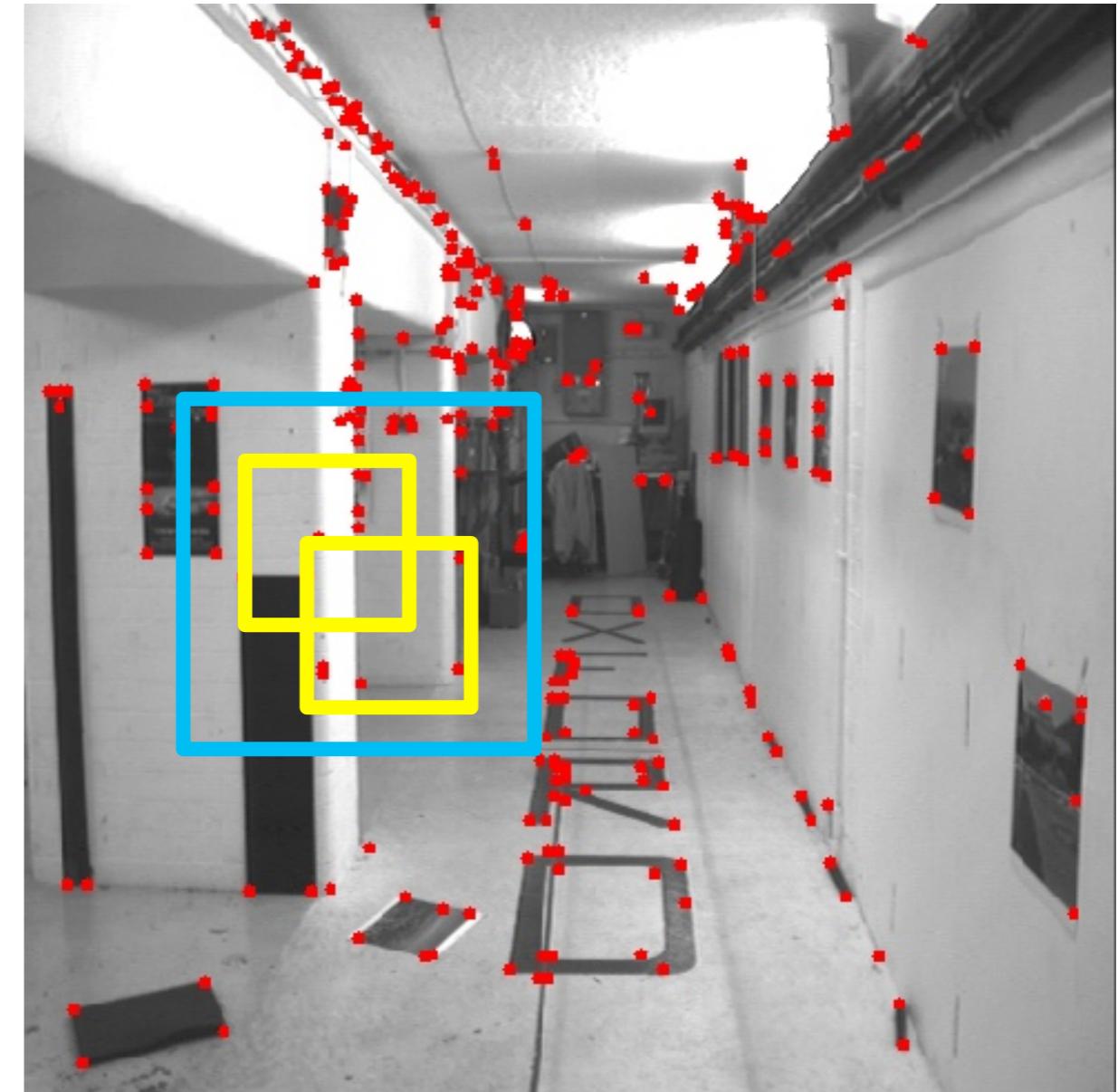
Example

- ◆ Interest point (here: Harris corner detector)



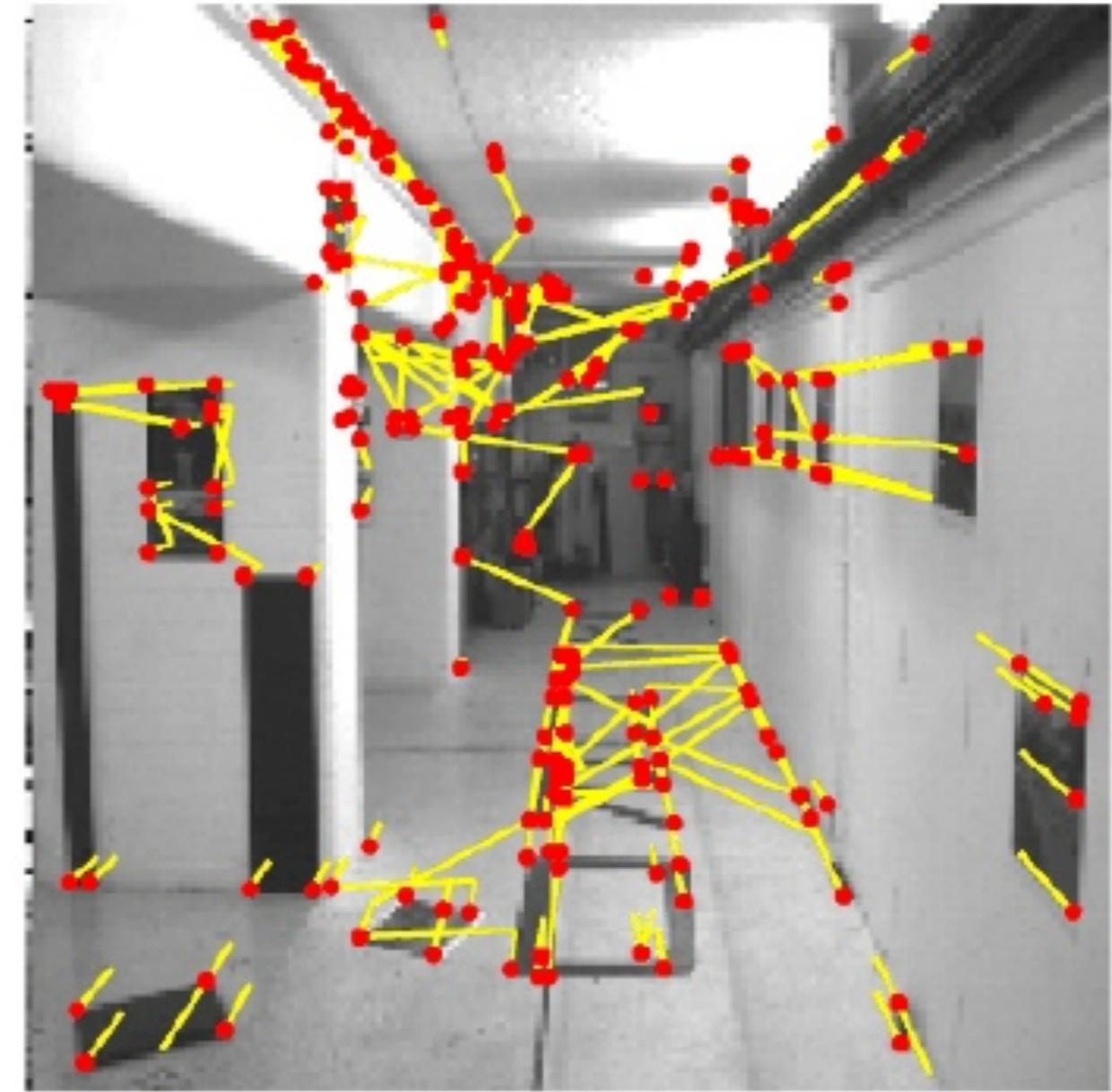
Example

- ◆ Match using only weak constraints (here: proximity)



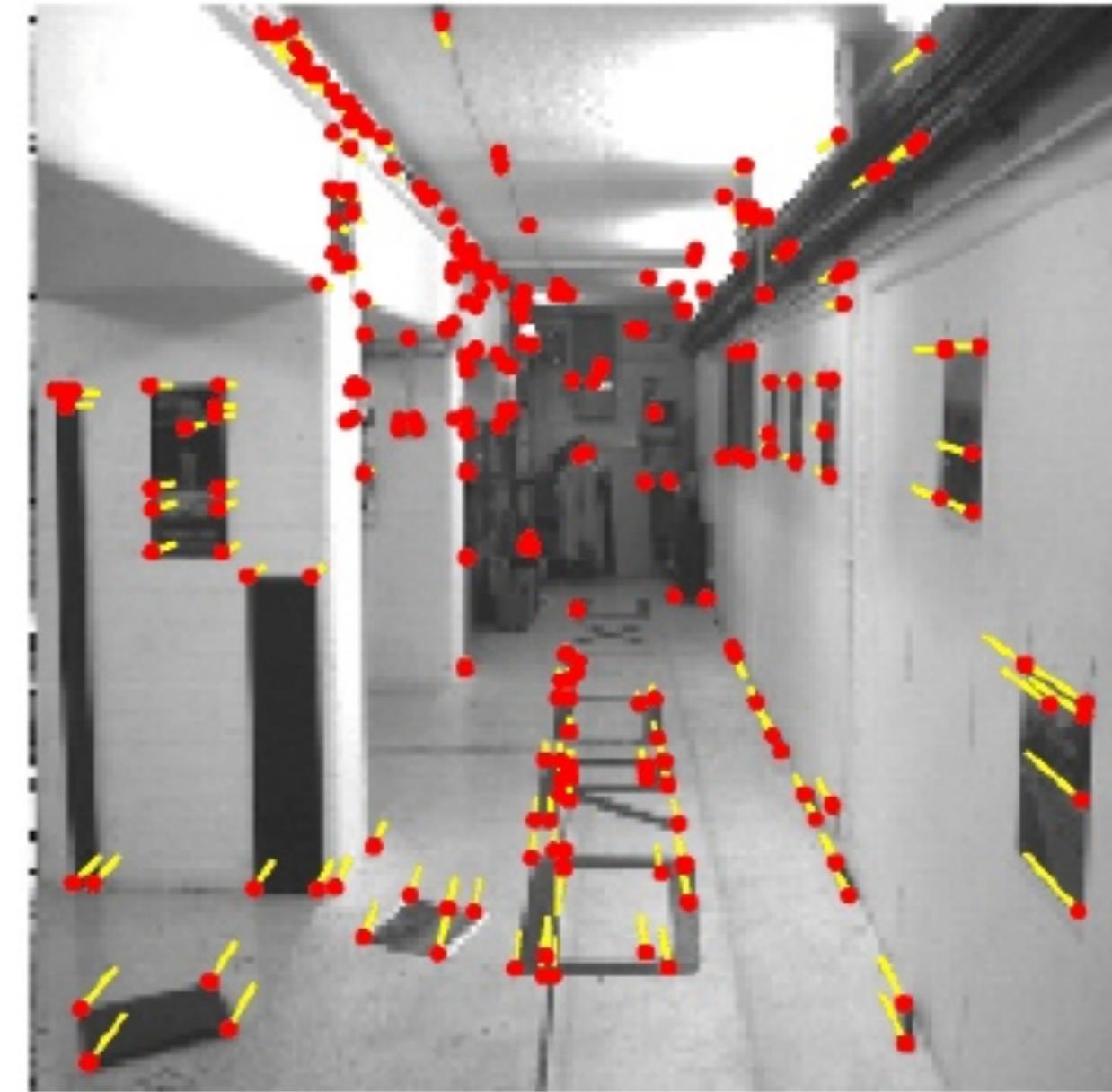
Example

- ◆ many wrong matches, but enough to estimate \mathbf{F} (RANSAC)



Example

- ◆ correspondences consistent with epipolar geometry



Example

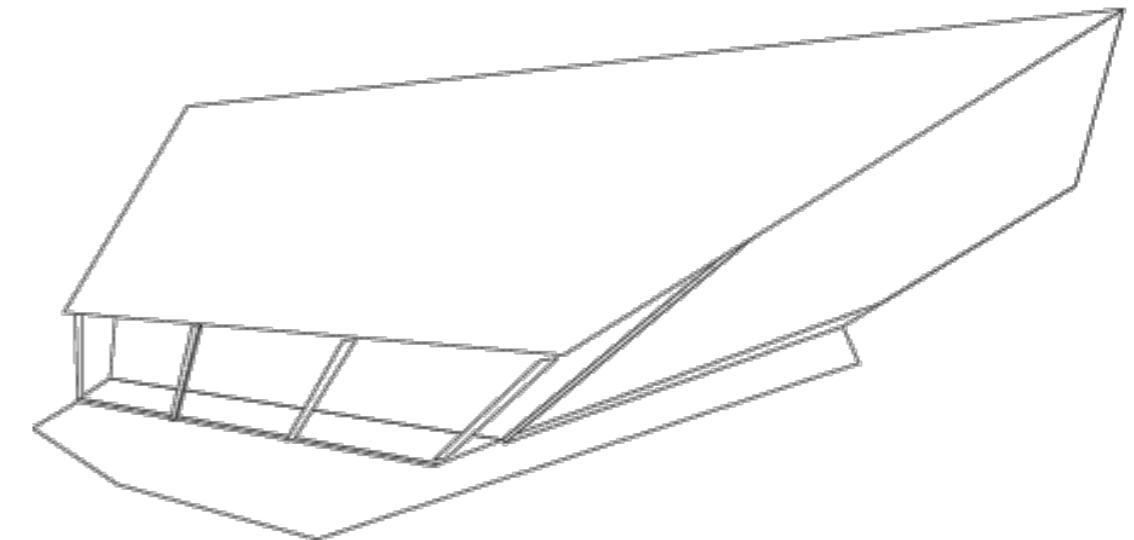
- ◆ Resulting epipolar geometry



What can we do with 2 views?



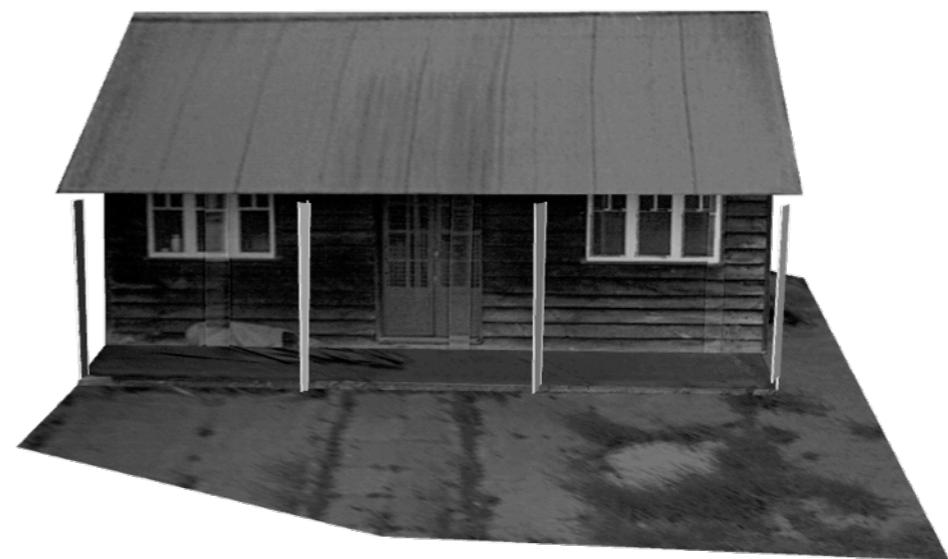
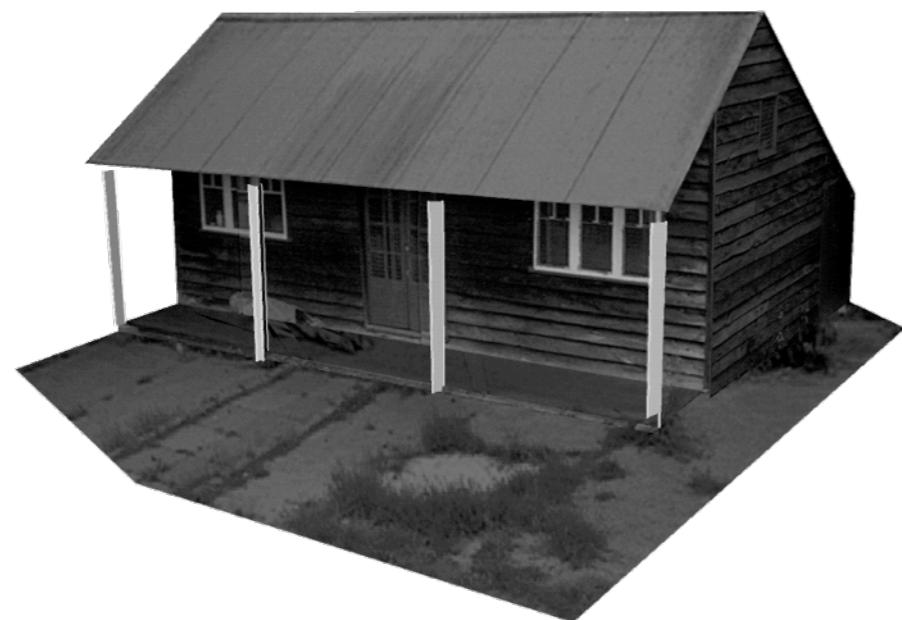
- ◆ Given two views of a static scene, and a small number of correspondences
- ◆ If the views were recorded with any pinhole camera, we can
 - ◆ reconstruct the scene up to a projective ambiguity
 - ◆ determine line intersections, coplanarity,...
 - ◆ need at least a 3rd view to find the calibration



What can we do with 2 views?



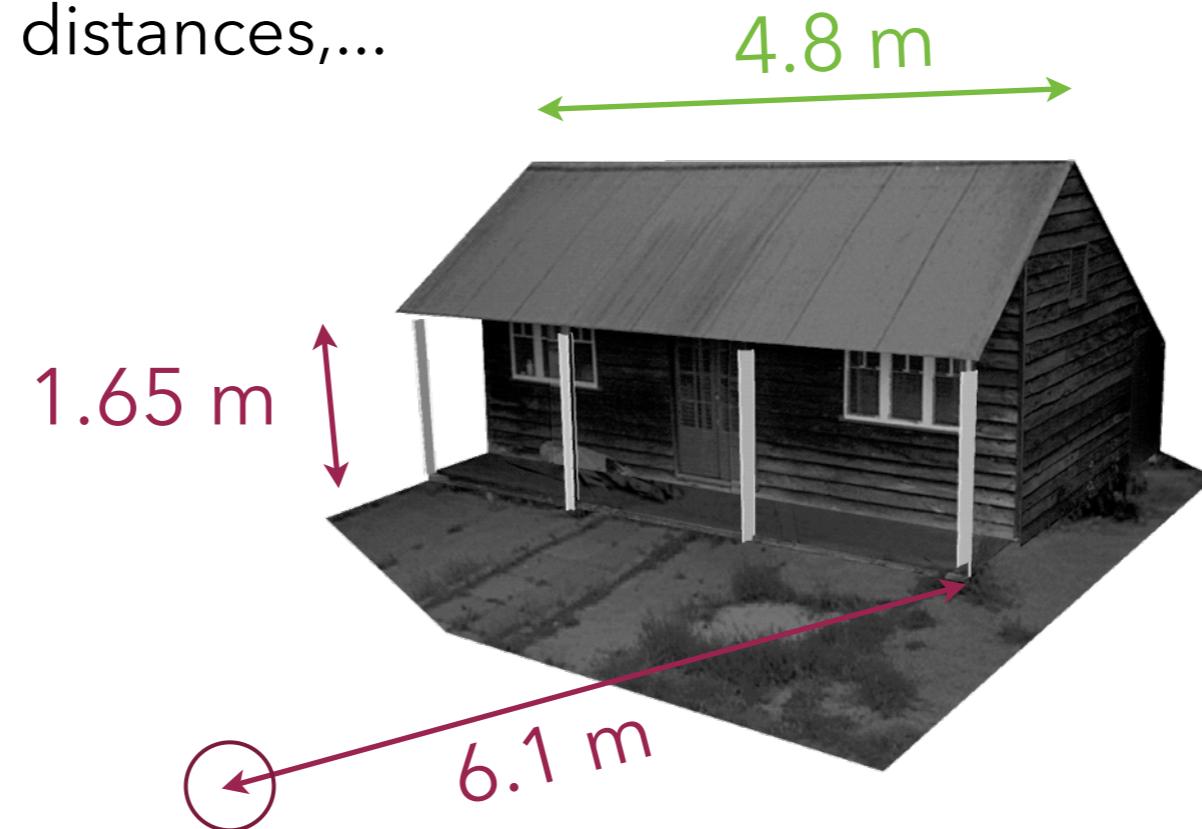
- ◆ Given two views of a static scene, and a small number of correspondences
- ◆ If the internal camera parameters have been calibrated, we can
 - ◆ reconstruct the scene up to scale
 - ◆ determine angles, relative distances,...
 - ◆ allows 3D modeling



What can we do with 2 views?



- ◆ Given two views of a static scene, and a small number of correspondences
- ◆ If the baseline length (or any other distance in the scene) is known, we can
 - ◆ determine the global scale
 - ◆ measure absolute distances,...
 - ◆ allows navigation

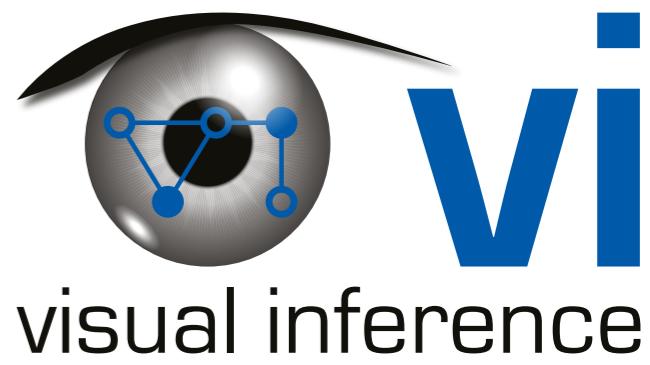


Computer Vision I

Dense Stereo



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Stereo Vision



- ◆ Given two views of the same scene, estimate its 3D shape
- ◆ Densely, not just at interest points...



Stereo Vision - Easier problem



- ◆ Narrower formulation: given two views with similar (but not identical!) viewpoint, fuse them to obtain a depth image

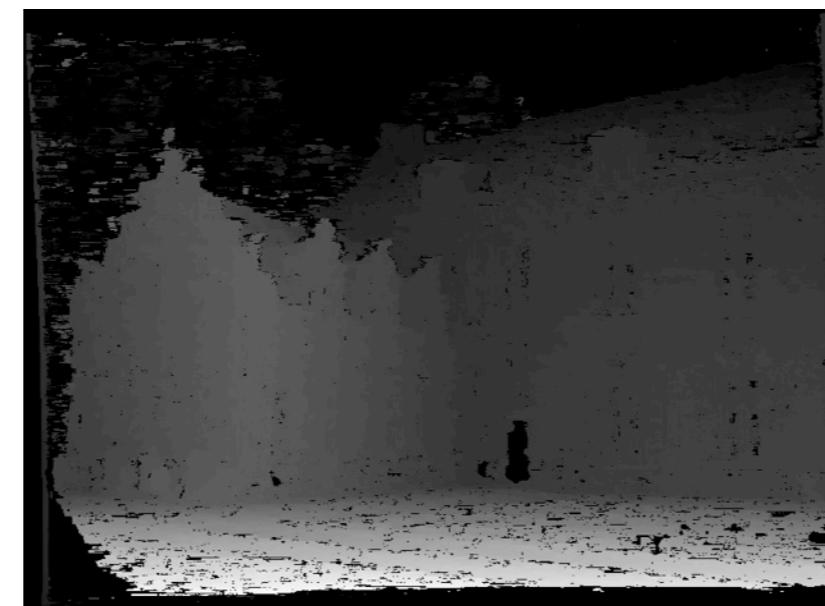
Image 1



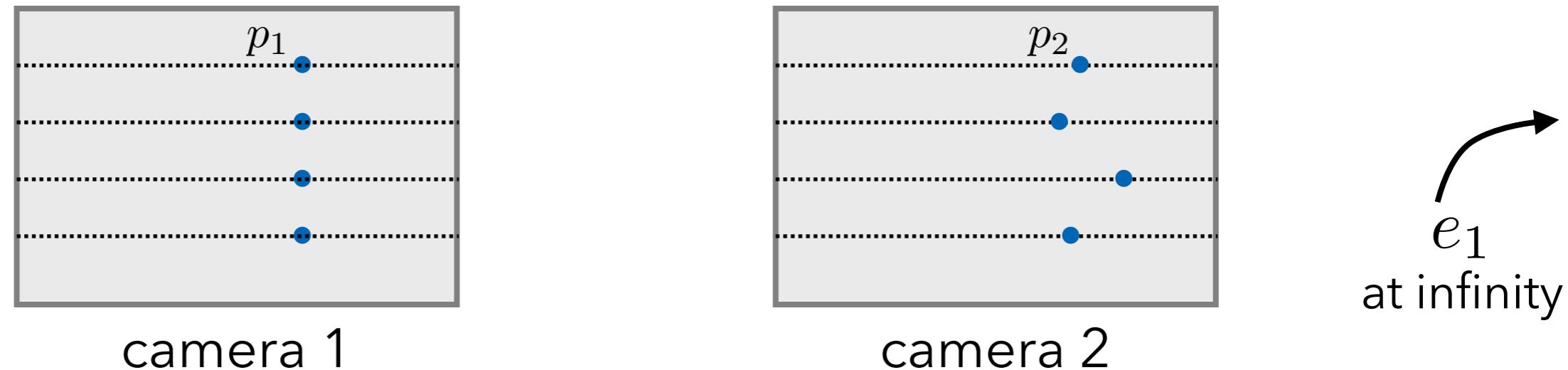
Image 2



Dense depth map



Binocular Stereo

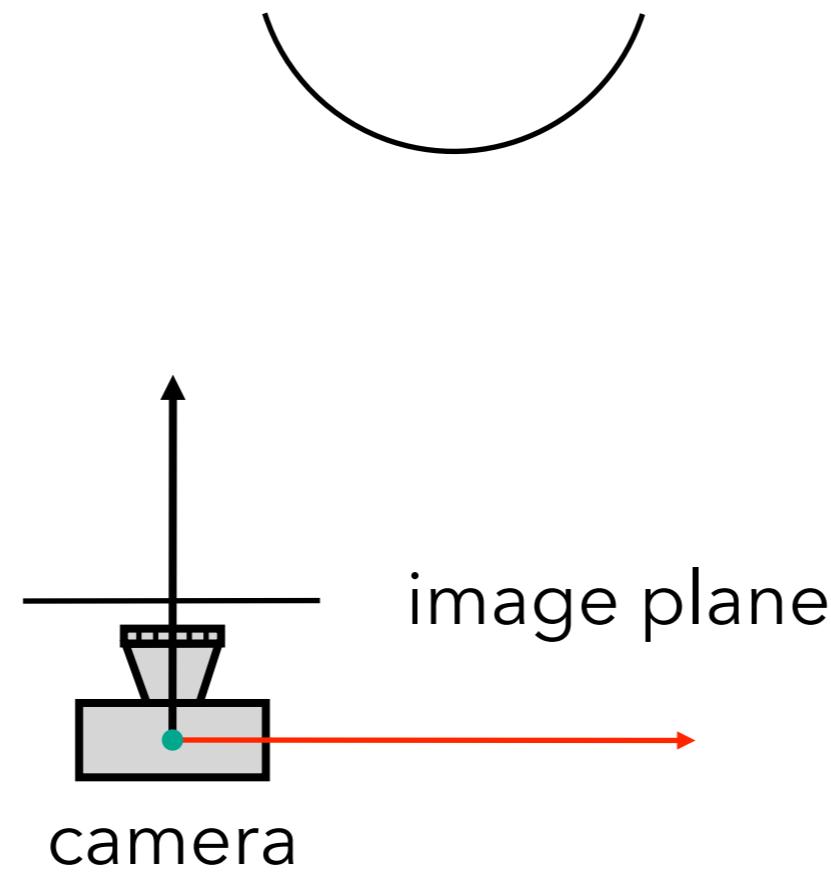


- ◆ Special case of epipolar geometry for stereo cameras with a standard binocular setup:
 - ◆ Epipoles are at infinity.
 - ◆ Epipolar lines are parallel.
 - ◆ Points correspond along the scanlines of the image.

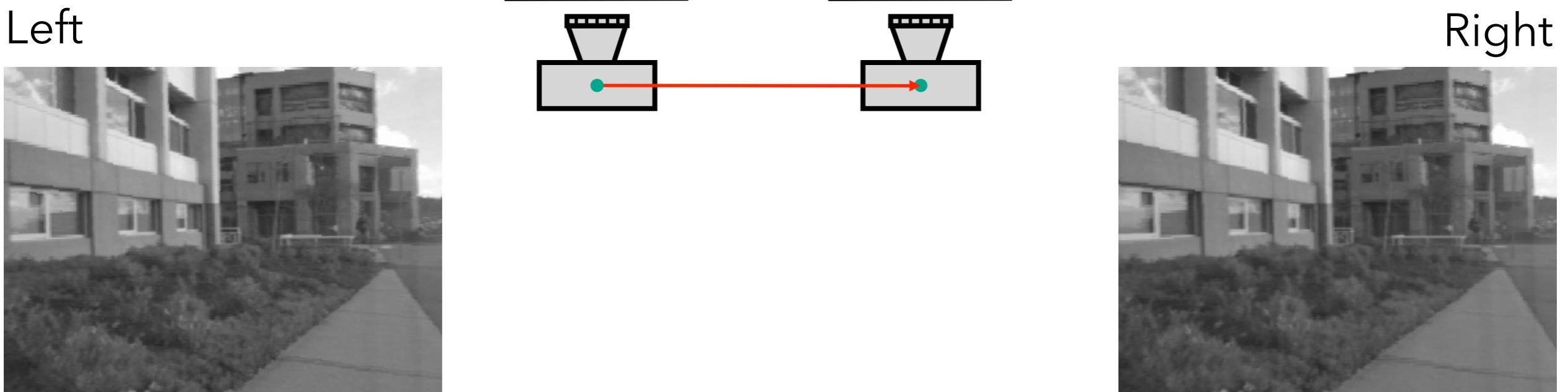
Binocular Stereo



Left

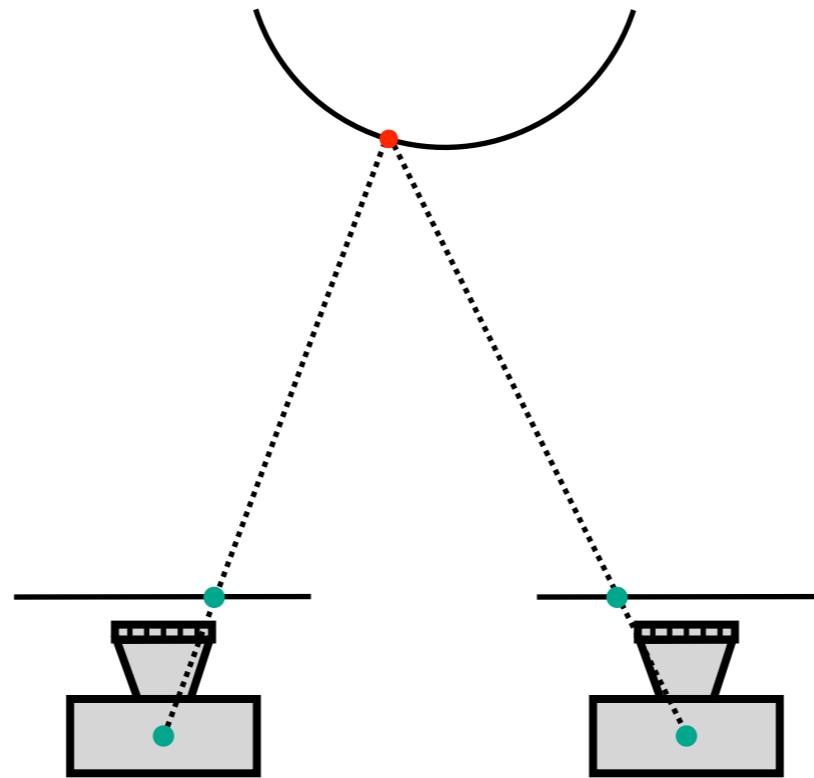


Binocular Stereo



Binocular Stereo

Left

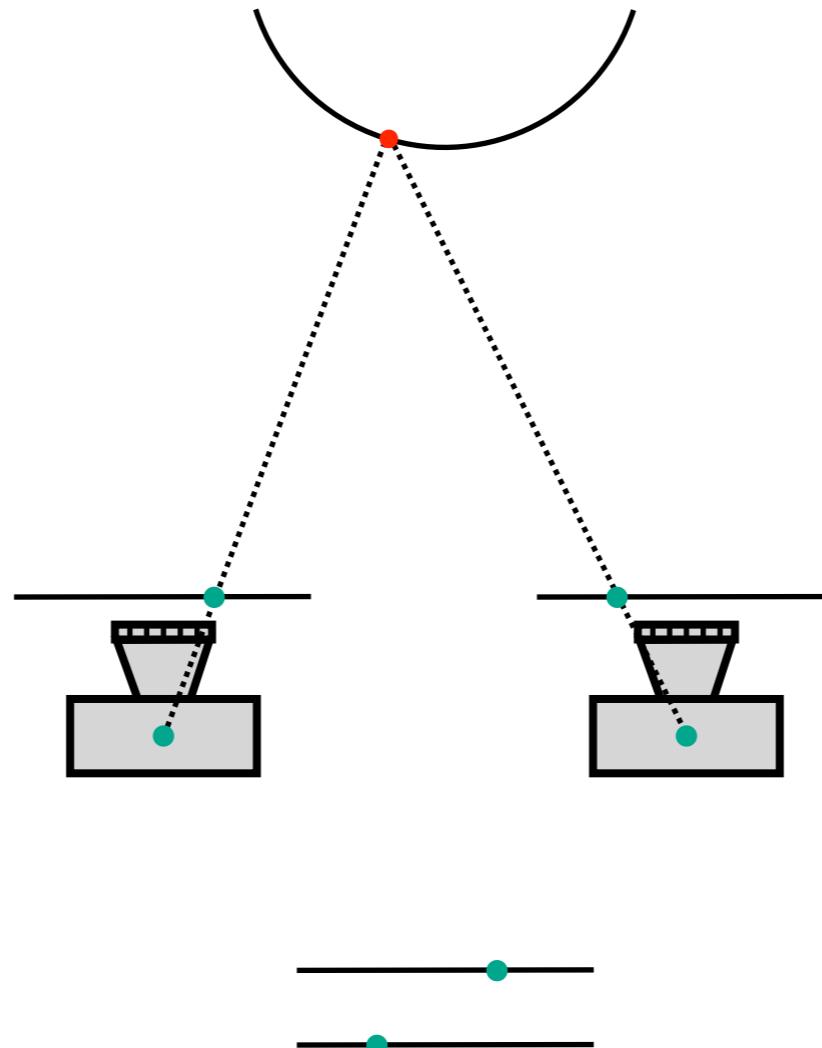


Right



Binocular Stereo

Left

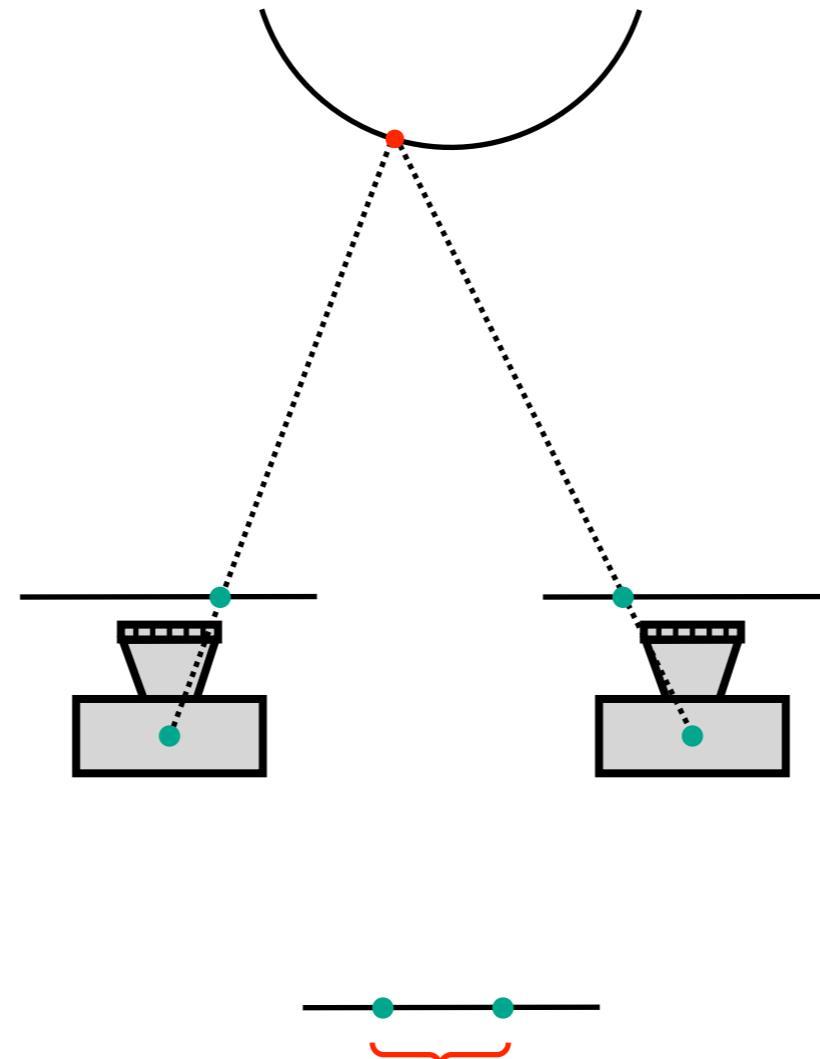


Right



Binocular Stereo

Left



binocular disparity

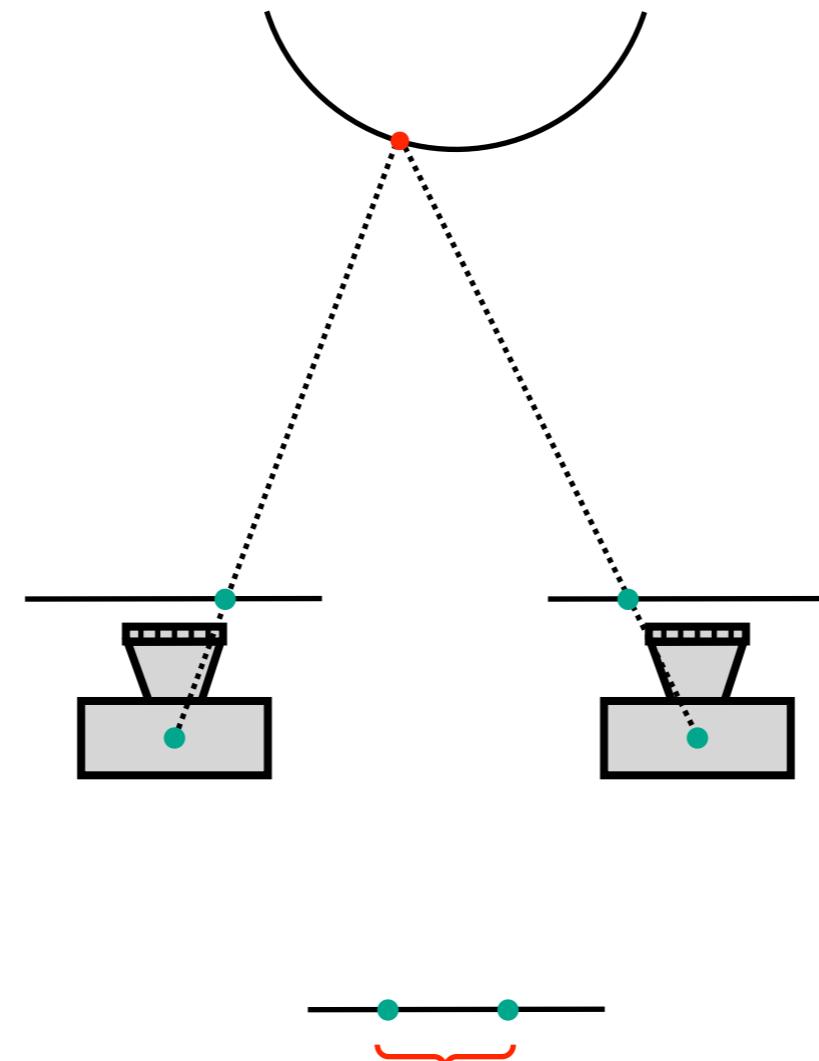
Right



Binocular Stereo



Left



binocular disparity

From known geometry of the cameras and estimated disparity, recover depth in the scene

Right



Triangulation

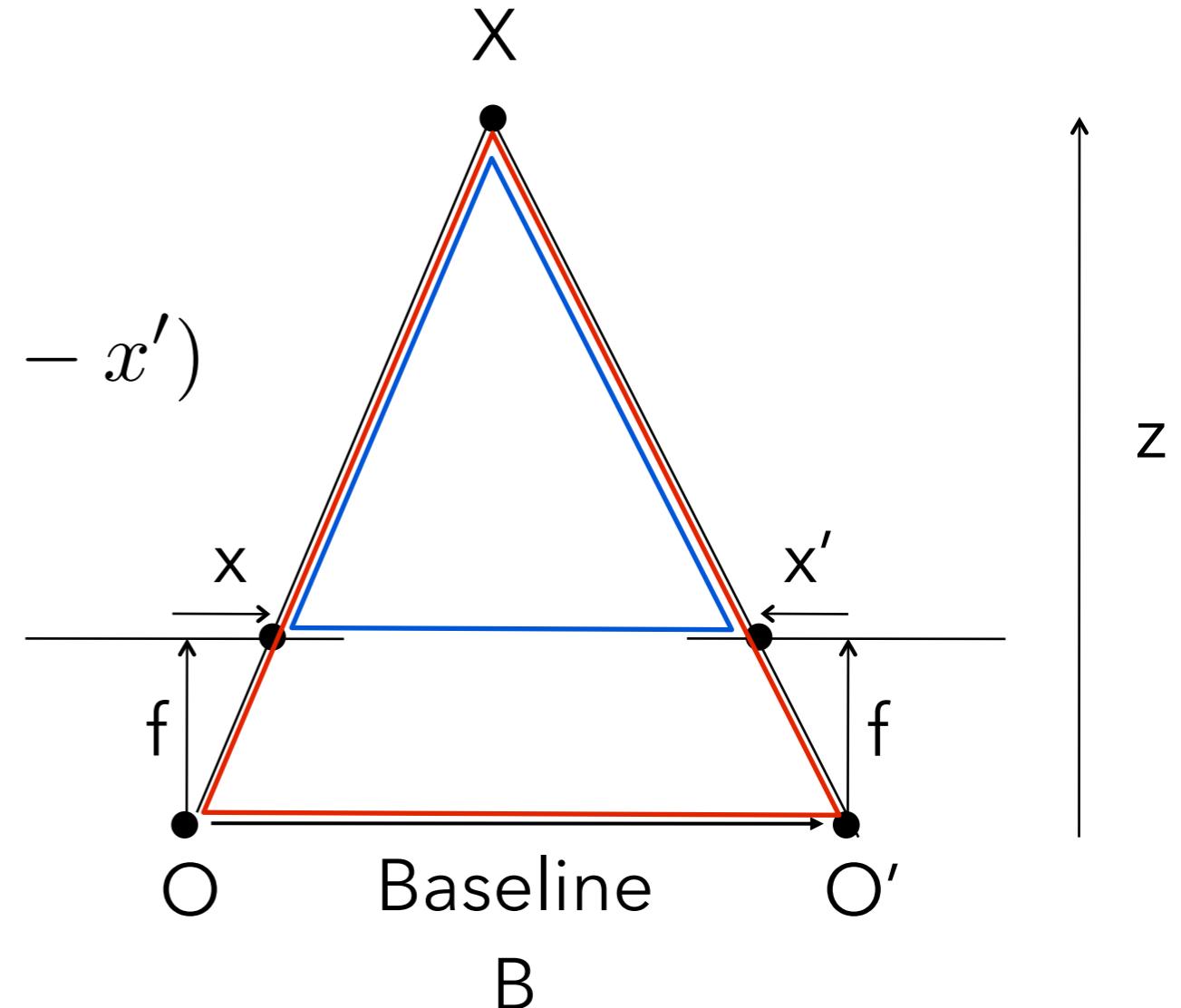
$$\frac{Z - f}{B - (x - x')} = \frac{Z}{B}$$

$$B \cdot Z - B \cdot f = B \cdot Z - Z \cdot (x - x')$$

$$Z = \frac{B \cdot f}{x - x'}$$

$$d = x - x' = \frac{B \cdot f}{Z}$$

disparity



- ◆ Parallel image planes: disparity from equal triangles

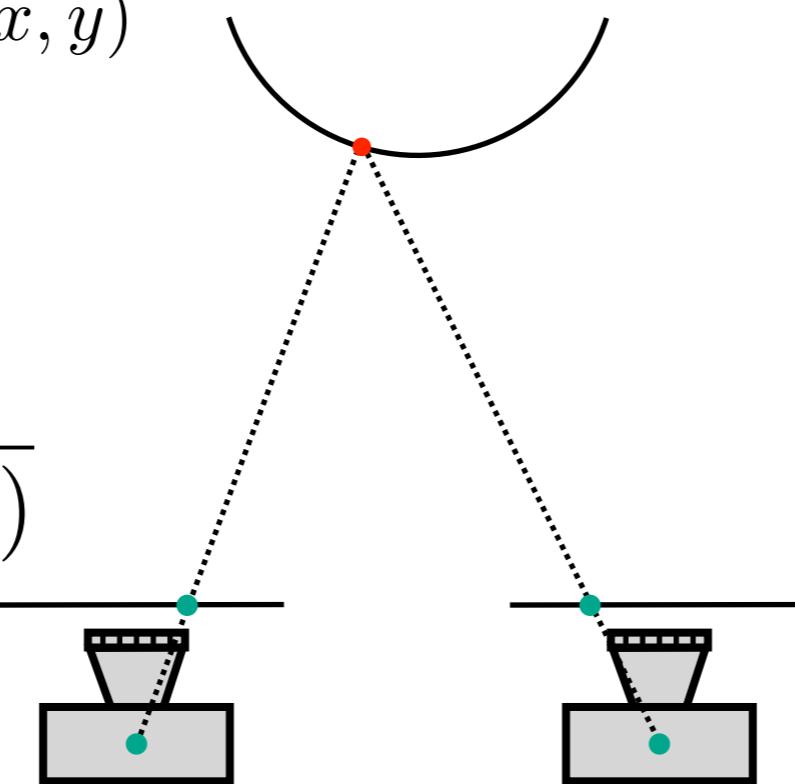
Binocular Disparity

$Z(x, y)$ is depth at pixel (x, y)
 $d(x, y)$ is disparity

Estimate:

$$Z(x, y) = \frac{fB}{d(x, y)}$$

Left



Search for best match

Right



[Black]

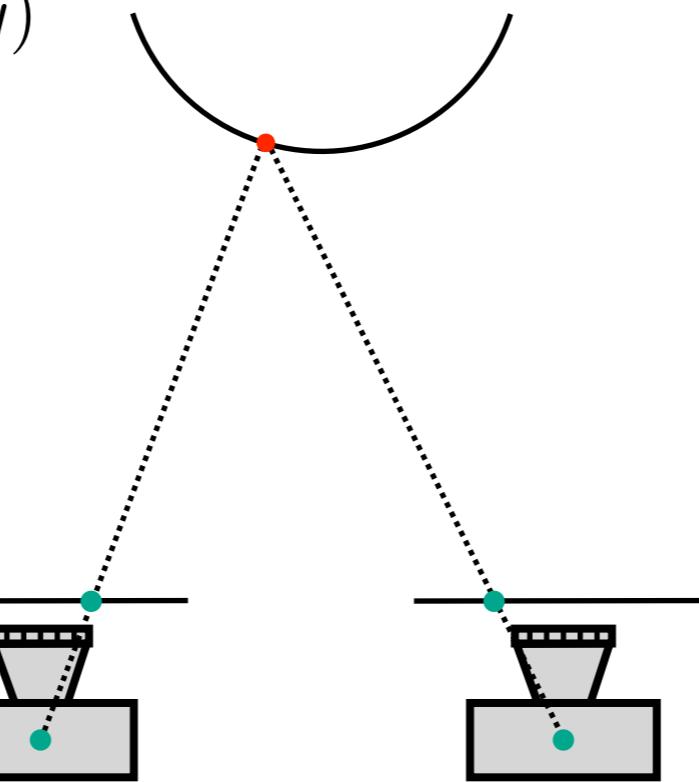
Binocular Disparity

$Z(x, y)$ is depth at pixel (x, y)
 $d(x, y)$ is disparity

Estimate:

$$Z(x, y) = \frac{fB}{d(x, y)}$$

Left



Right

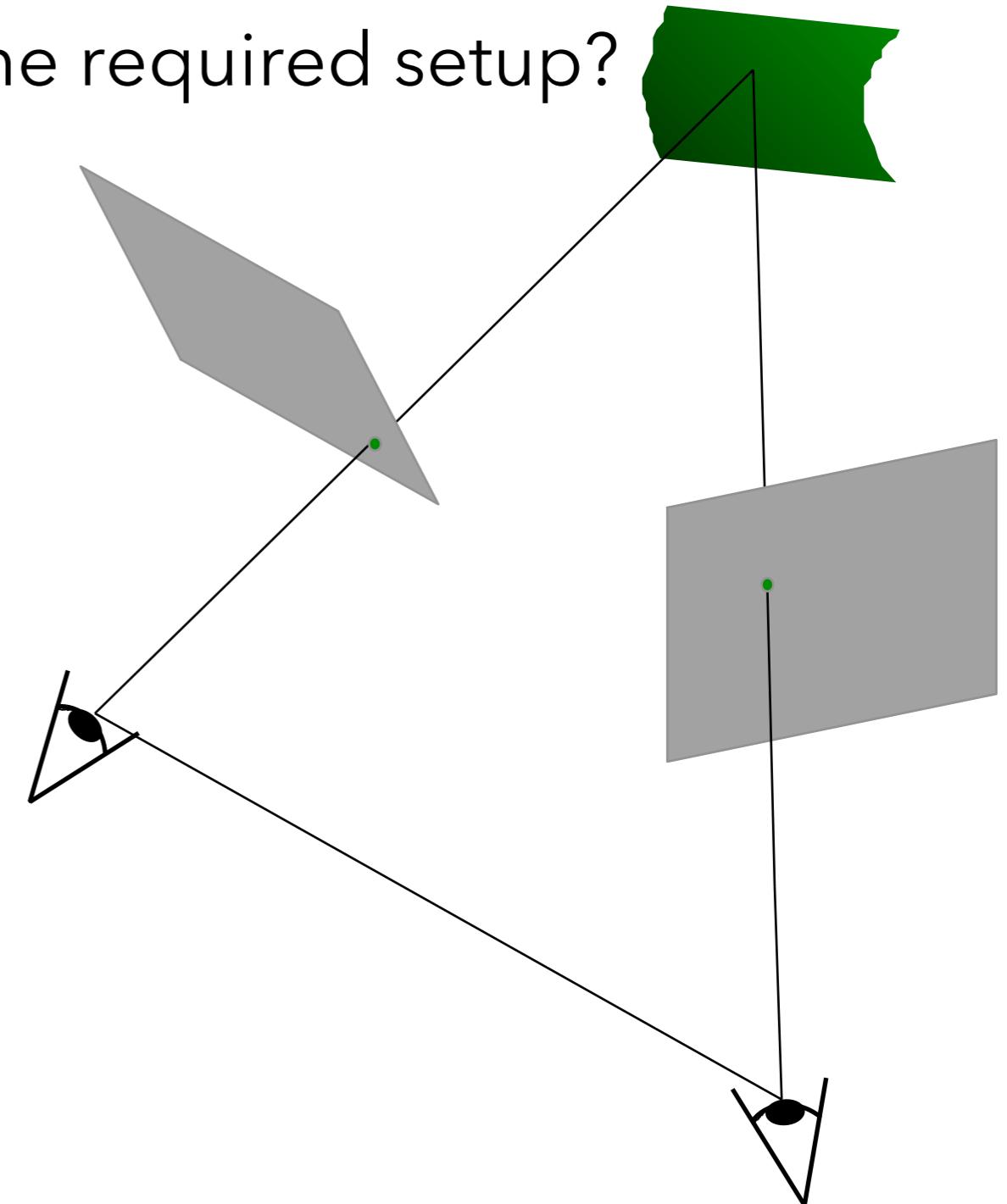
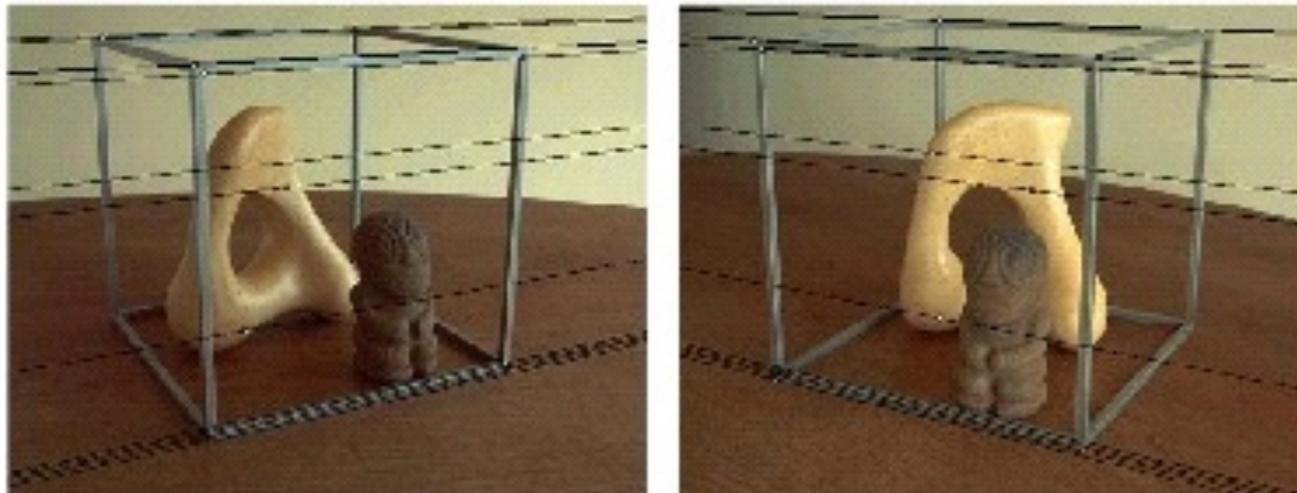


Do I need to consider
this region?

[Black]

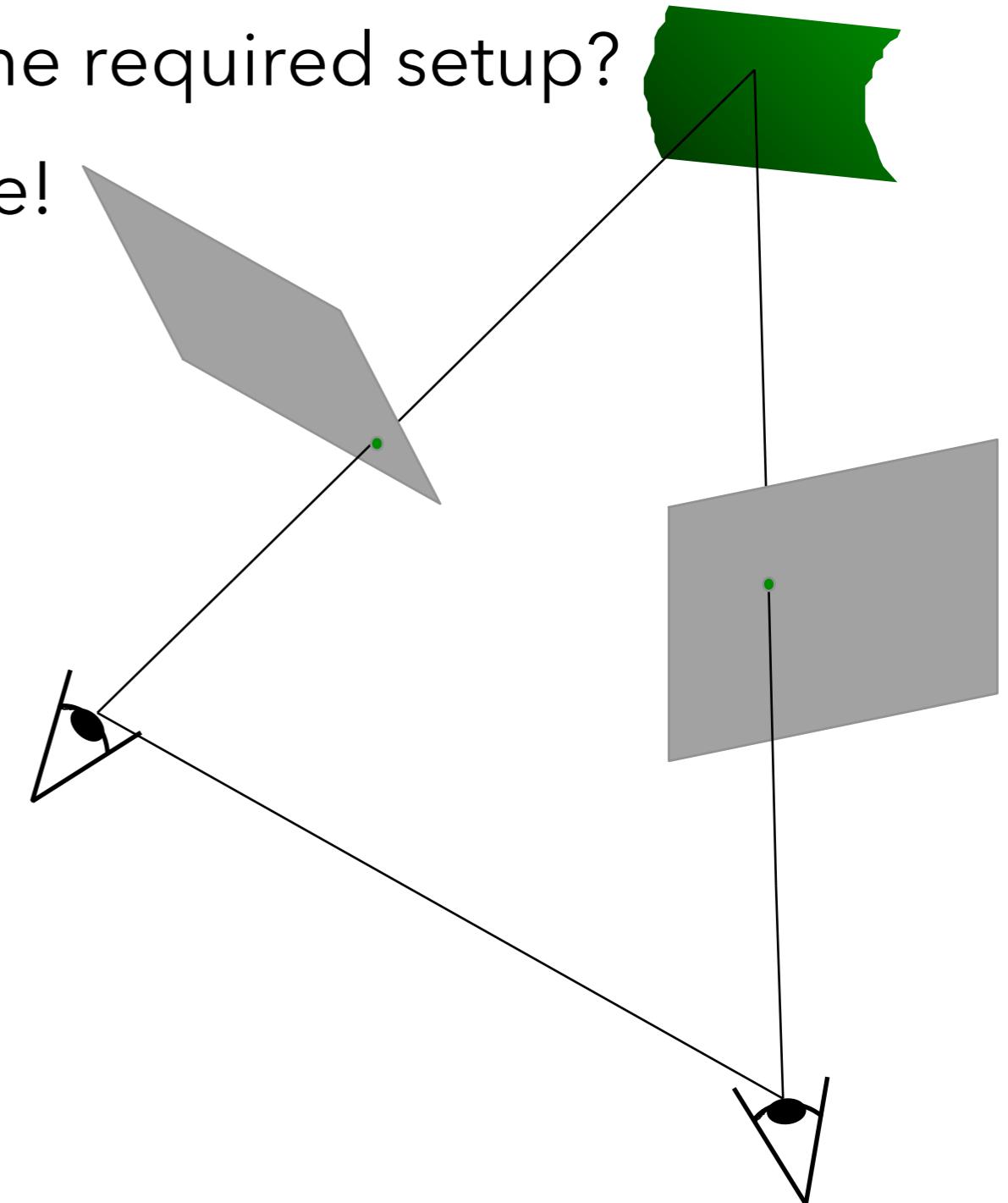
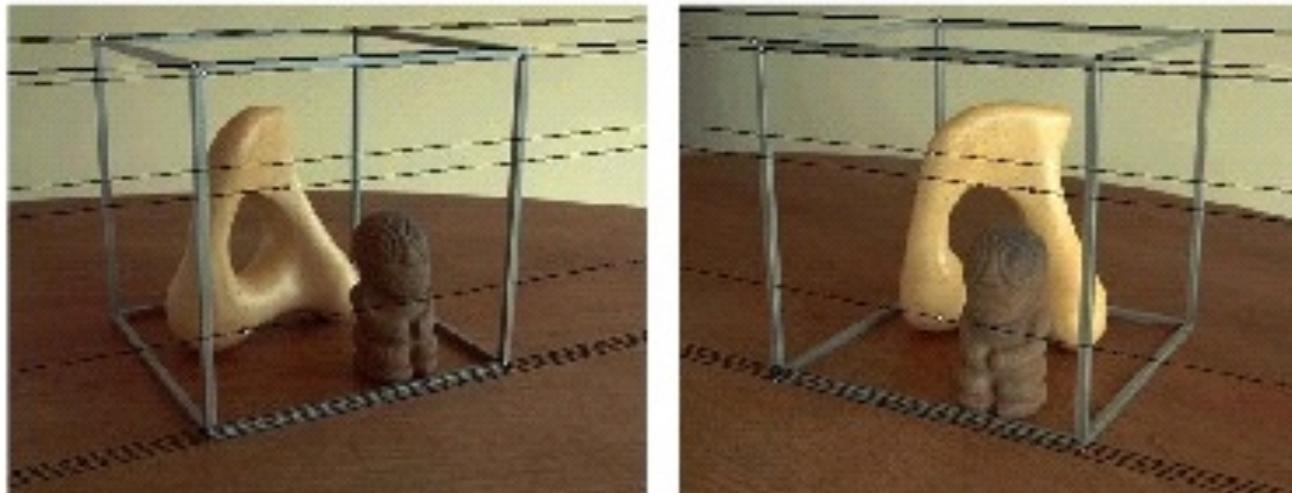
Stereo rectification

- ◆ What if the images are not in the required setup?



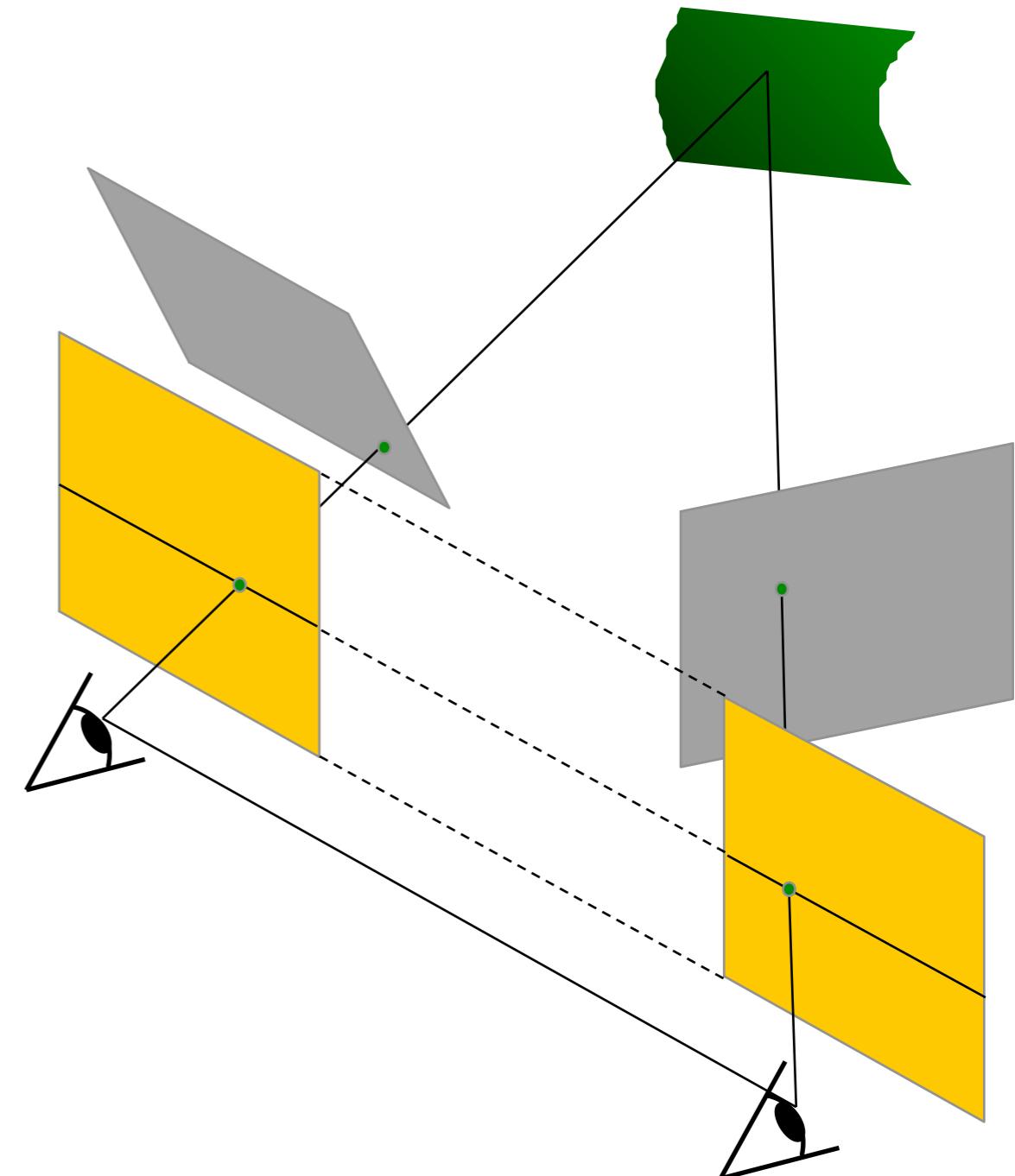
Stereo rectification

- ◆ What if the images are not in the required setup?
- ◆ Rewarp them such that they are!

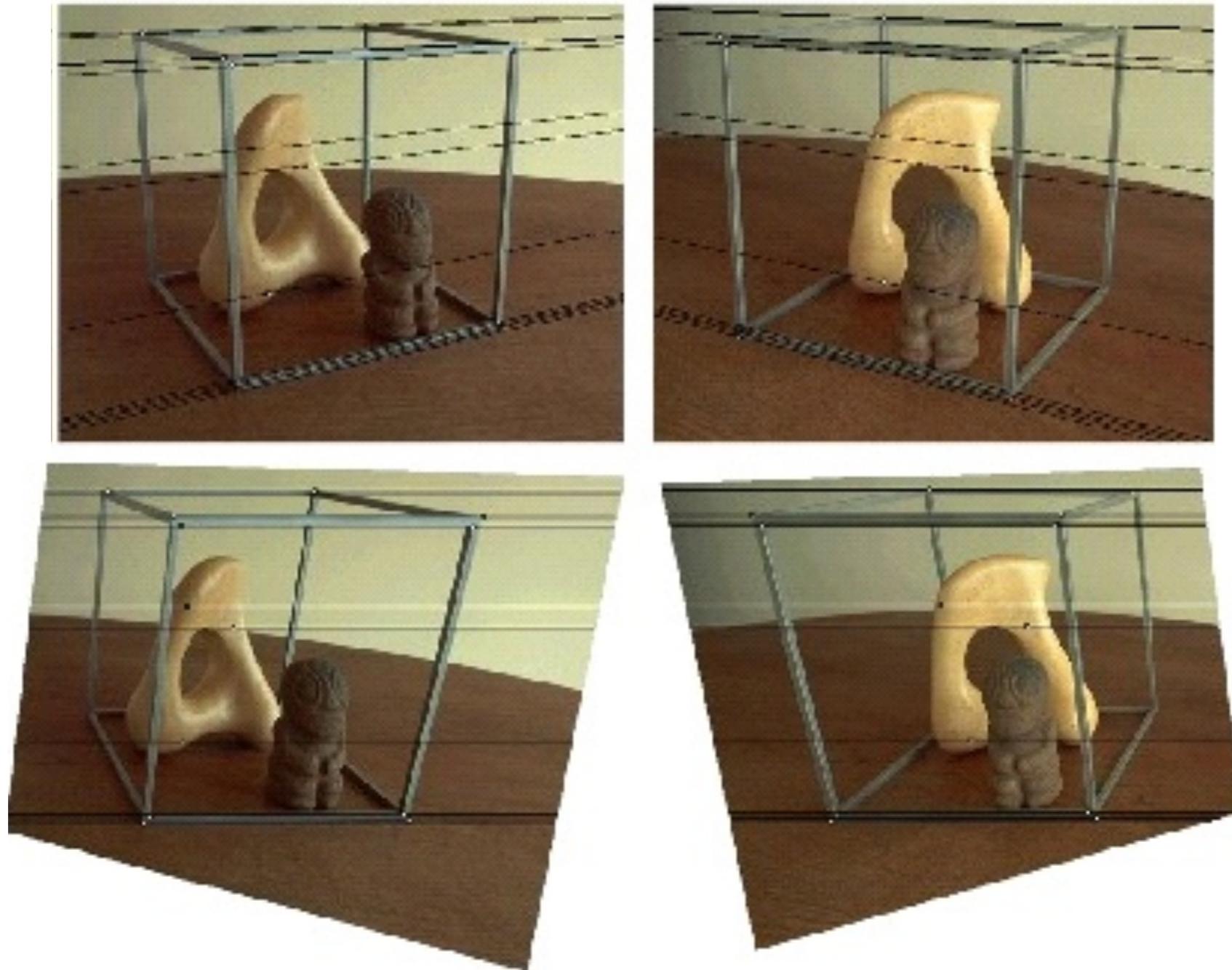


Stereo rectification

- ◆ Solution: map both image planes to a common plane parallel to the baseline
- ◆ requires a homography for each image
- ◆ after the transform, pixel motion is horizontal again

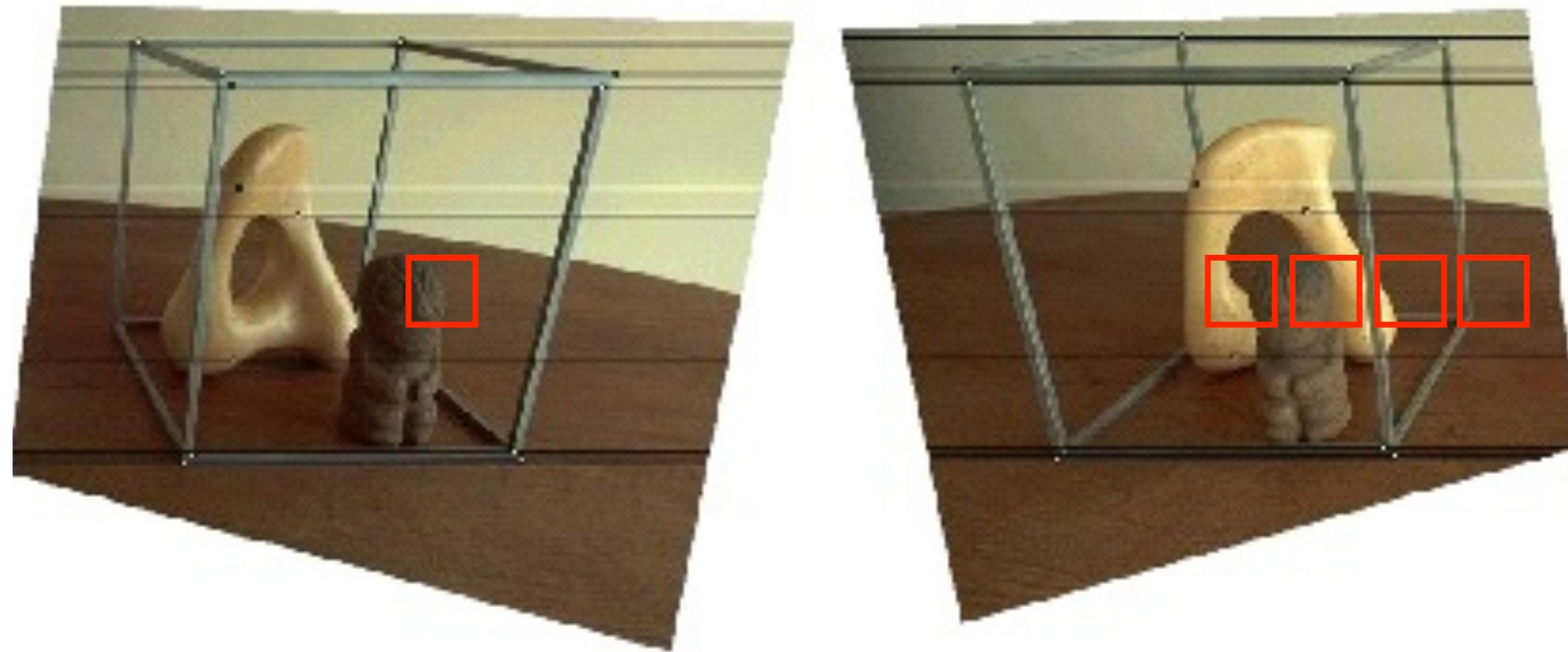


Stereo rectification



[Szeliski and Fleet]

Matching



- ◆ Matching only has to occur along epipolar lines.
- ◆ Now in the simpler binocular case where the cameras are pointing forward.

[Szeliski and Fleet]

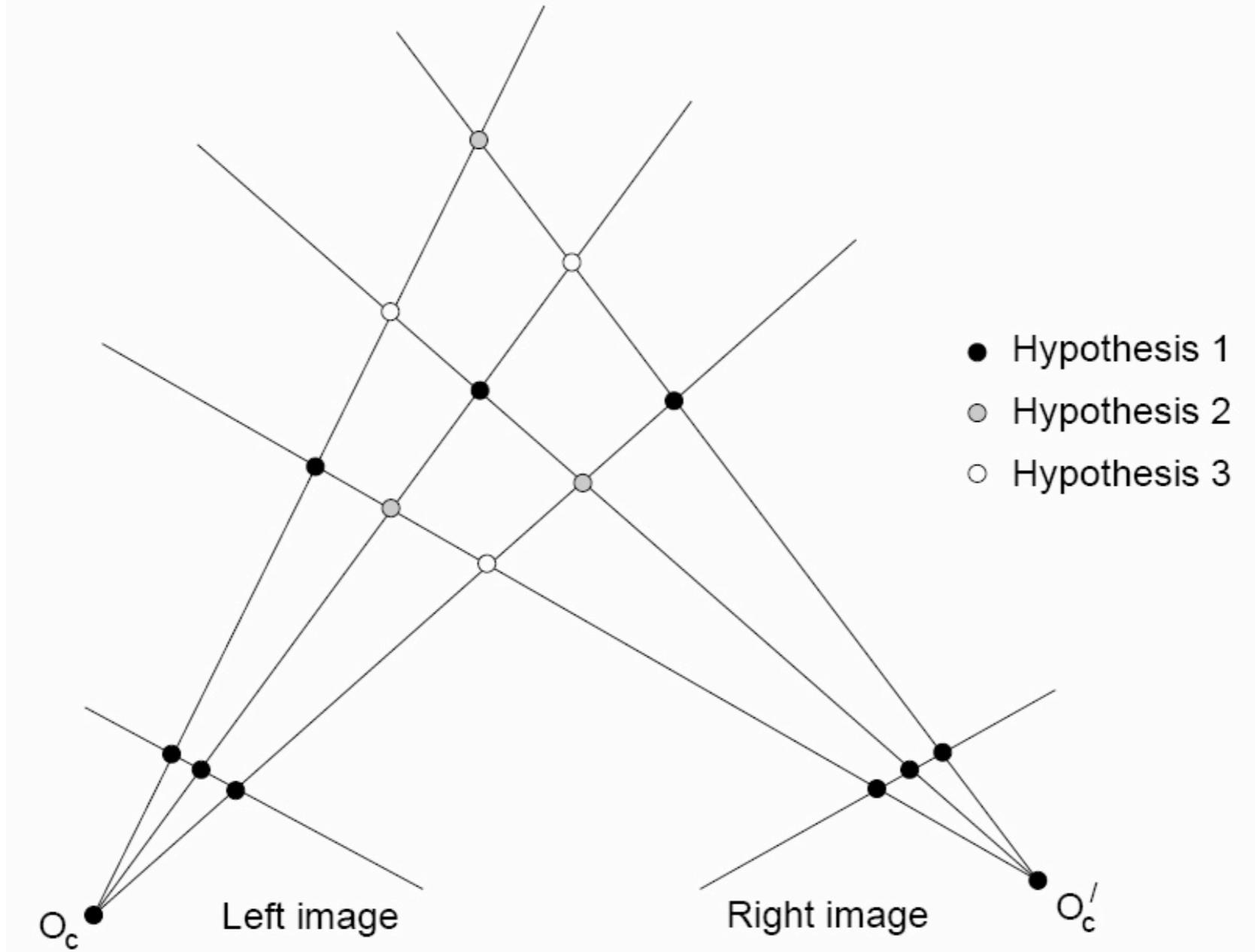
Stereo Correspondence

- ◆ Search over disparity to find correspondences
- ◆ Range of disparities to search over can change dramatically within a single image pair.



[Black]

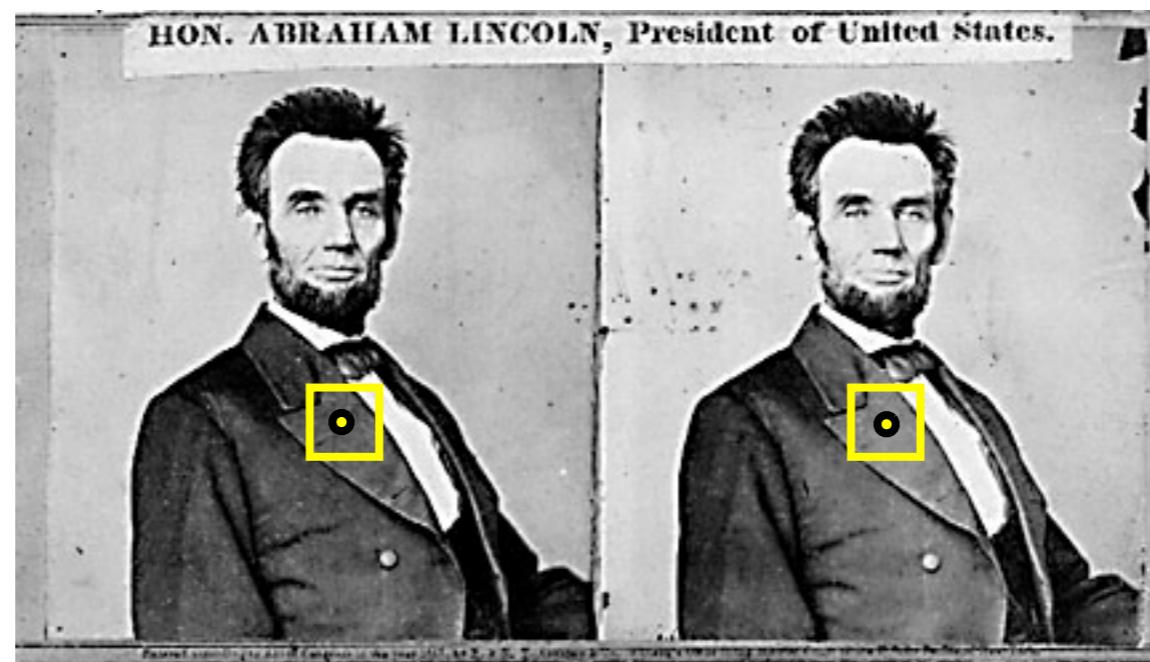
Correspondence problem



Even when constrained to 1D, there are multiple matching hypotheses - which one is correct?

Correspondence problem

- ◆ Let's make some assumptions to simplify the matching problem
 - ◆ The baseline is relatively small (compared to the depth of scene points) - "narrow-baseline stereo"
 - ◆ Matching regions are similar in appearance
 - ◆ Most scene points are visible in both views



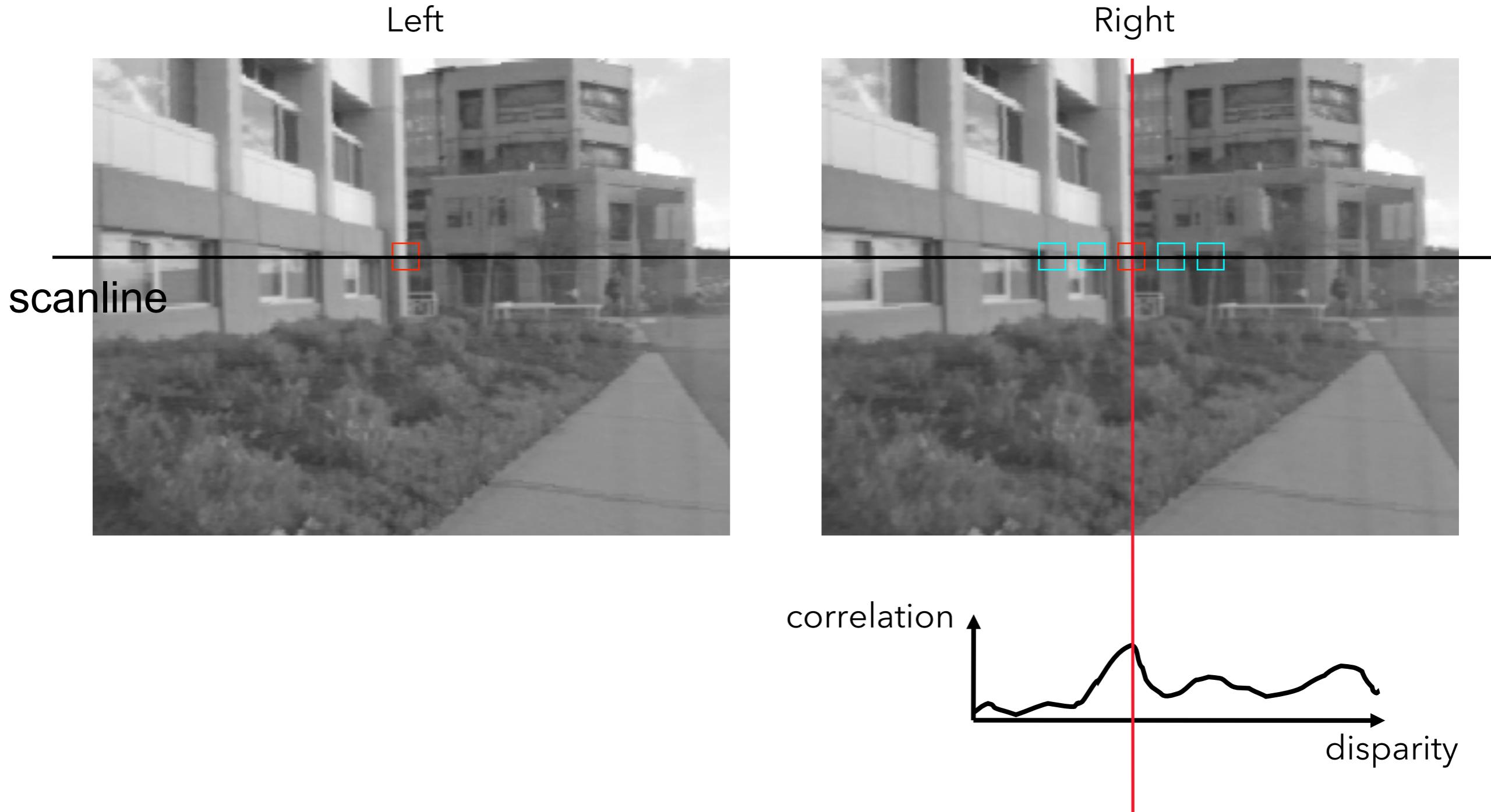
note: the 1st assumption makes the others a lot more likely

Correspondence problem

- ◆ Matching with narrow baseline
 - ◆ What are two “similar pixels”?
 - ◆ An ill-posed question, a single intensity sample does not reveal the local image structure
 - ◆ Measure similarity of the surrounding region
 - ◆ Area-based similarity measures



Correspondence Using Correlation



Normalized Correlation

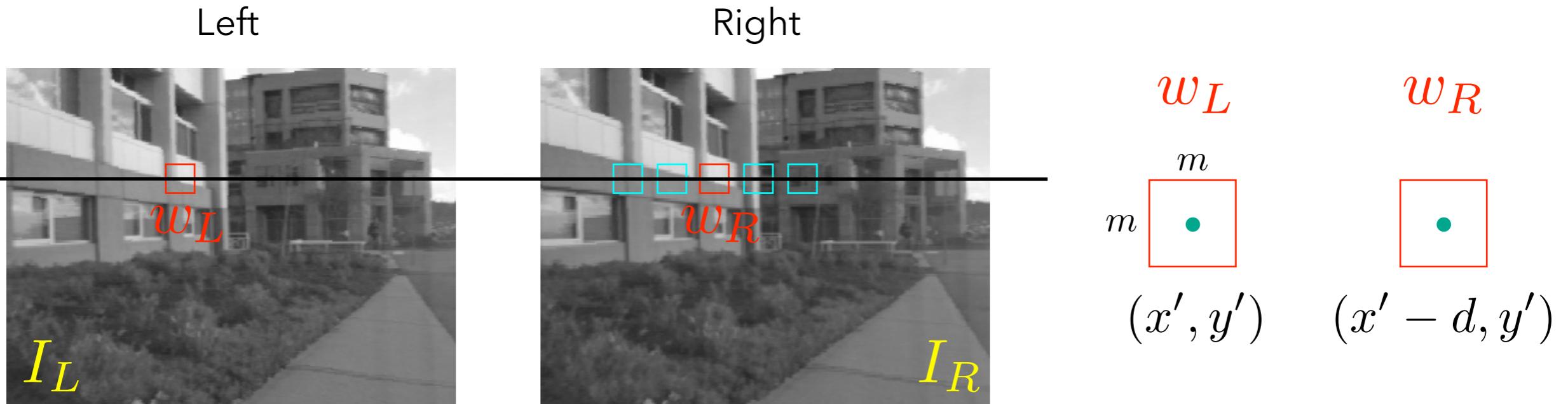


- ◆ w_L and w_R are corresponding $m \times m$ windows of pixels. We also write them as vectors \mathbf{w}_L and \mathbf{w}_R .
- ◆ The normalized correlation computes the cosine of the angle between the patches:

$$\text{NC}(x, y, d) = \frac{(\mathbf{w}_L(x, y) - \bar{\mathbf{w}}_L(x, y))^T (\mathbf{w}_R(x - d, y) - \bar{\mathbf{w}}_R(x - d, y))}{|\mathbf{w}_L(x, y) - \bar{\mathbf{w}}_L(x, y)| |\mathbf{w}_R(x - d, y) - \bar{\mathbf{w}}_R(x - d, y)|}$$

patch mean

Even simpler: Sum of Squared (Pixel) Differences



- ◆ w_L and w_R are corresponding $m \times m$ windows of pixels.
- ◆ The SSD cost measures the intensity difference as a function of disparity:

$$\text{SSD}_R(x, y, d) = \sum_{(x', y') \in w_L(x, y)} (I_L(x', y') - I_R(x' - d, y'))^2$$

[Black]

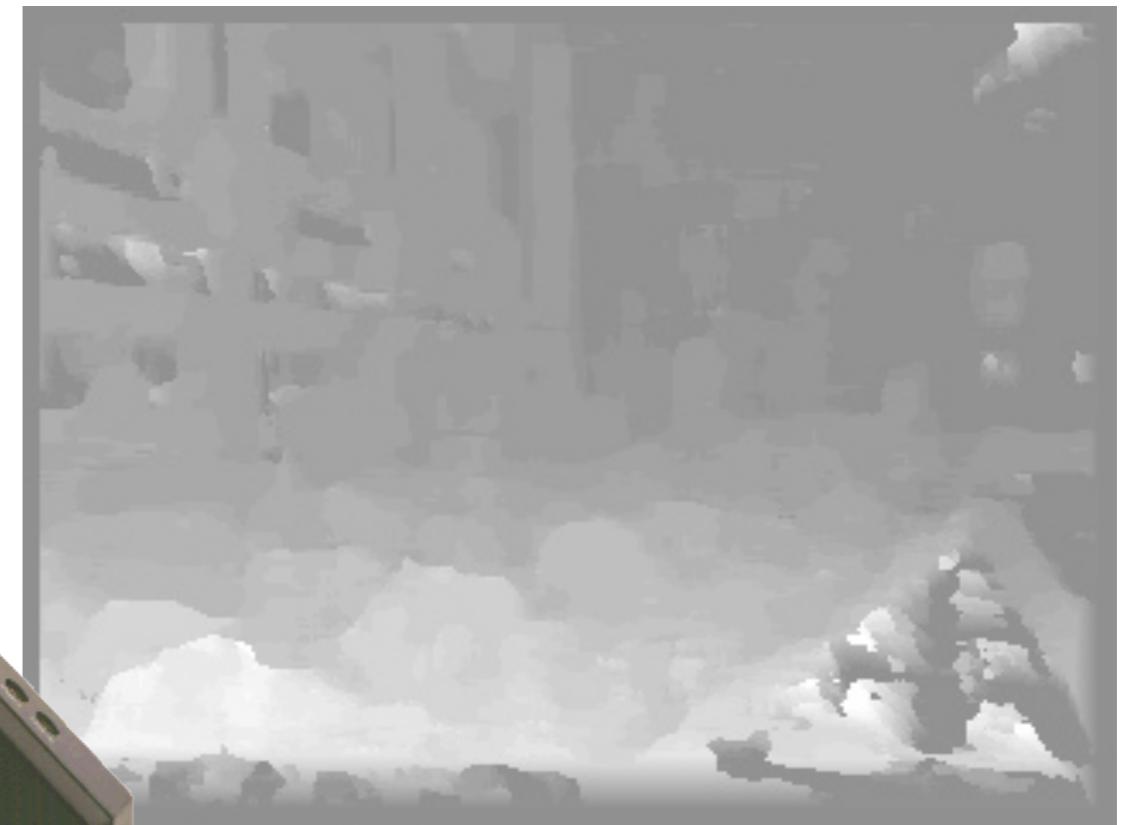
Window-based matching

- ◆ In practice:
 - ◆ Choose some disparity range $[0, d_{\max}]$
 - ◆ For all pixels (x, y) try all disparities and choose the one that maximizes the normalized correlation or minimizes the SSD.

Left



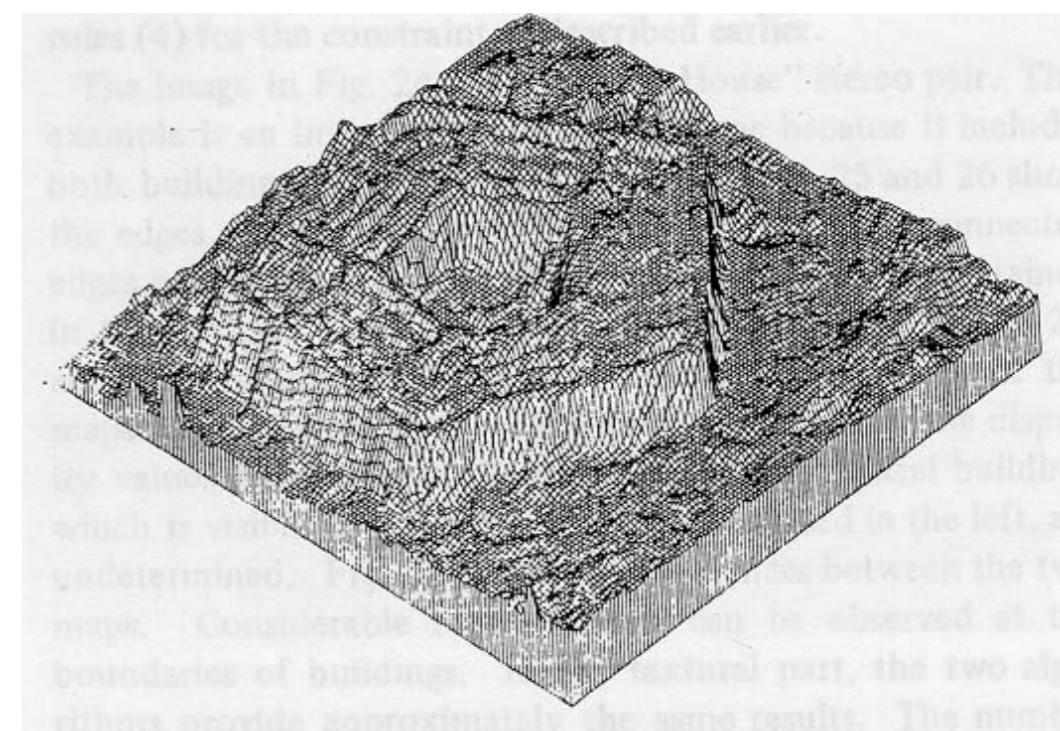
Disparity Map



[Point Grey Research]

Convert disparity to depth

- ◆ Of course, we can convert the disparity back to depth:



From [Ohta & Kanade, 1985]

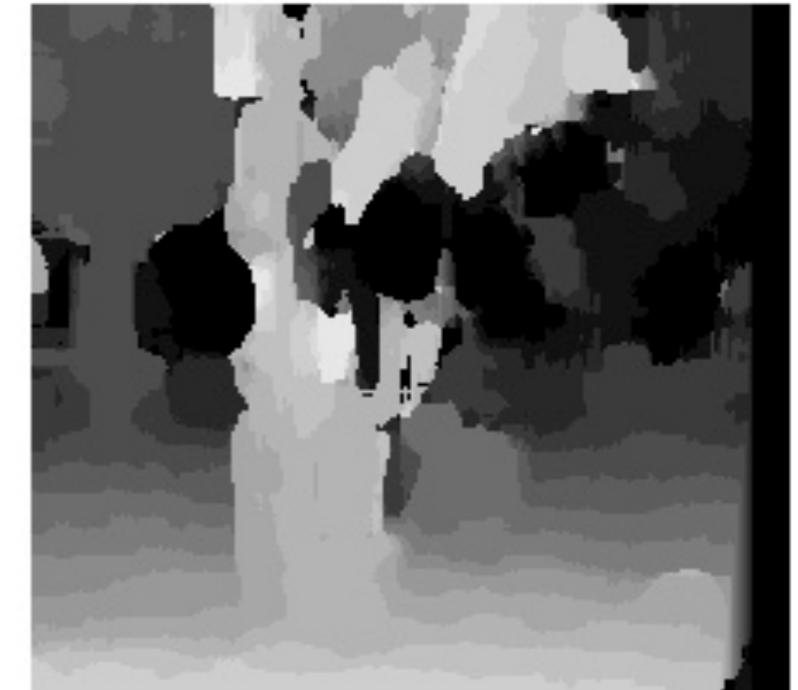
Window-based matching

- ◆ In practice:
 - ◆ Choose some disparity range $[0, d_{\max}]$
 - ◆ For all pixels (x, y) try all disparities and choose the one that maximizes the normalized correlation or minimizes the SSD.
- ◆ Challenges and Problems:
 - ◆ How do we choose the right window size m ?
 - ◆ Mismatches often lead to relatively poor results quality.

Influence of window size



$m = 3$



$m = 20$

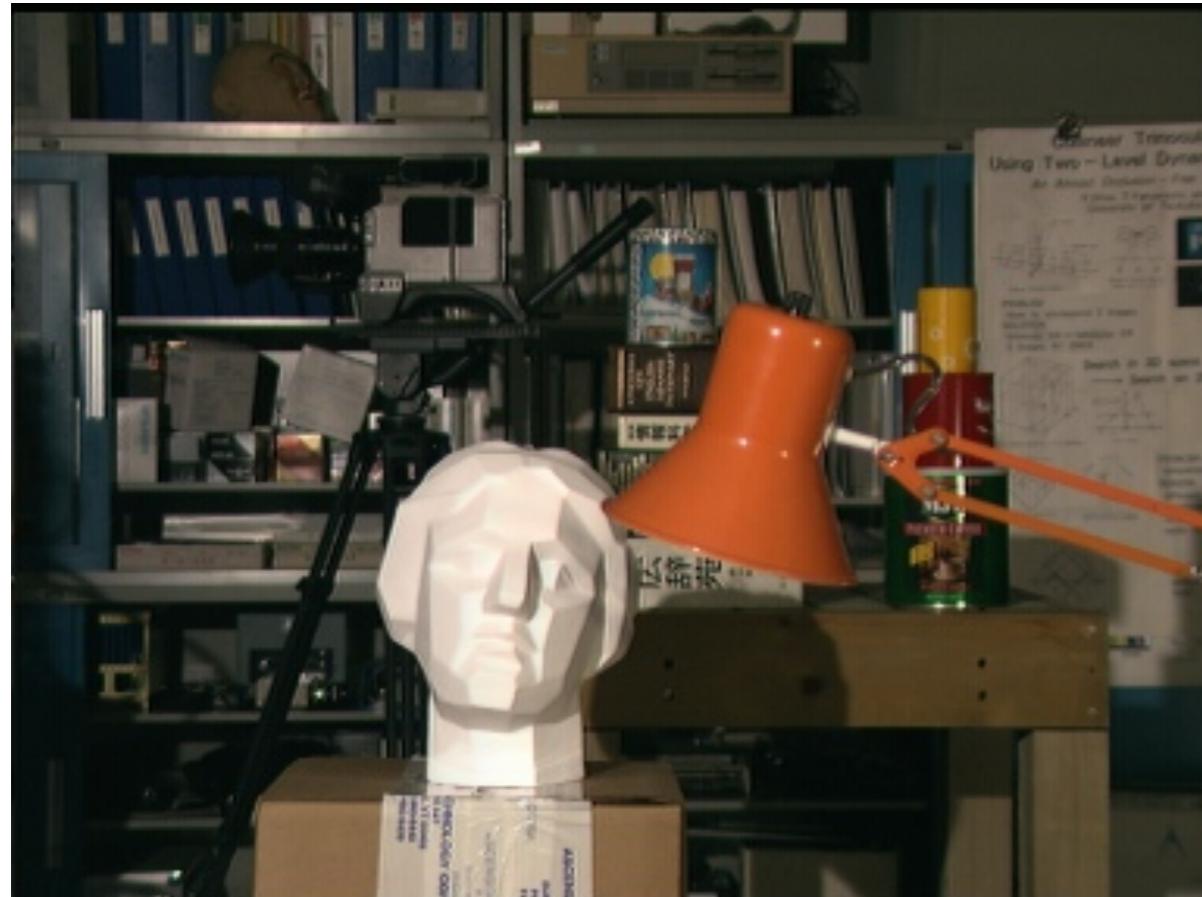
- ◆ Better results with adaptive window:

- ◆ T. Kanade and M. Okutomi: A Stereo Matching Algorithm with an Adaptive Window: Theory and Experiment, Proc. International Conference on Robotics and Automation, 1991.
- ◆ D. Scharstein and R. Szeliski: Stereo matching with nonlinear diffusion. International Journal of Computer Vision, 28(2):155-174, July 1998 .

[Seitz]

Preview: Stereo results

- ◆ Data from University of Tsukuba:



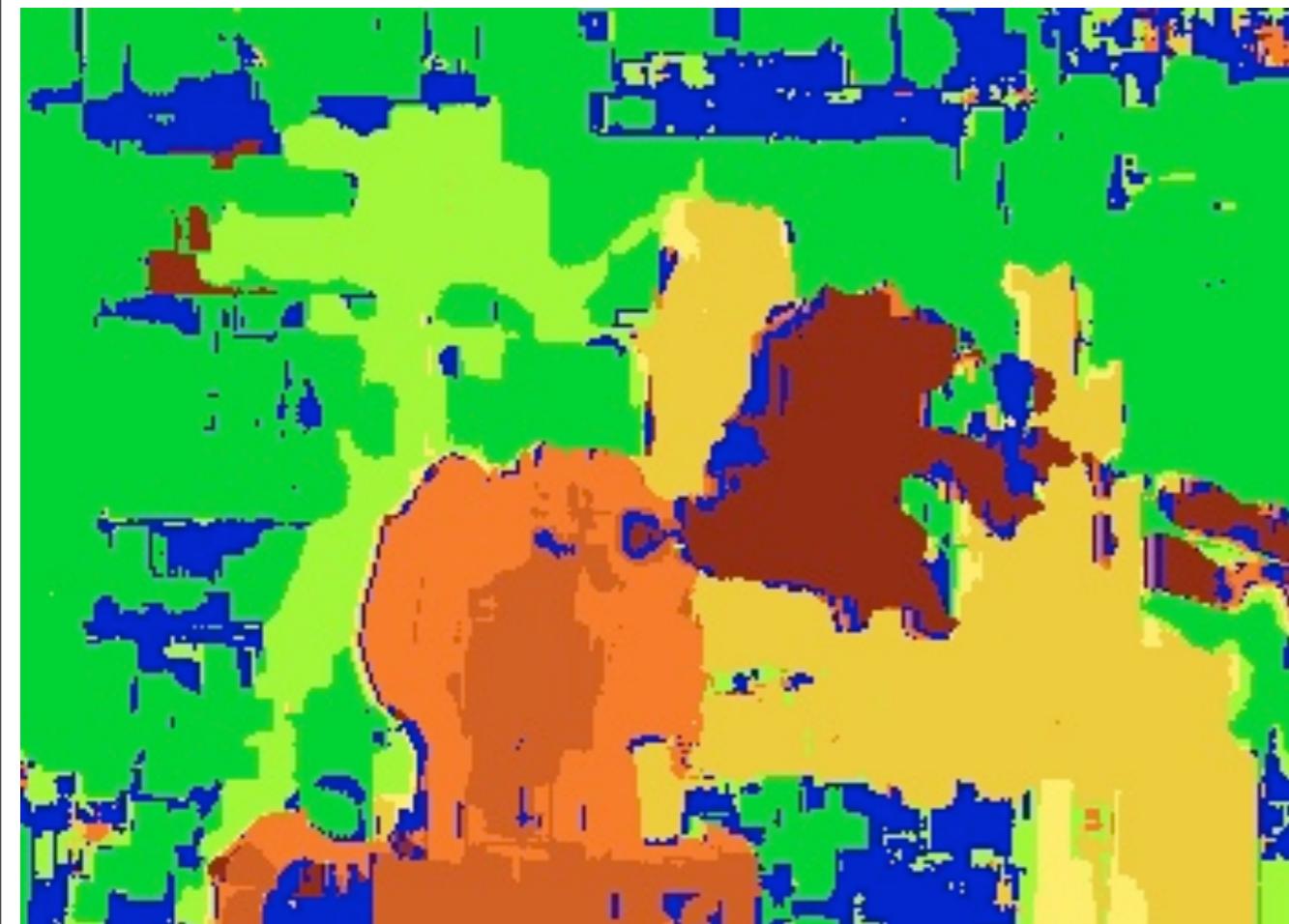
Scene



Ground truth

[Seitz]

Results with window correlation



Window-based matching
(best window size)

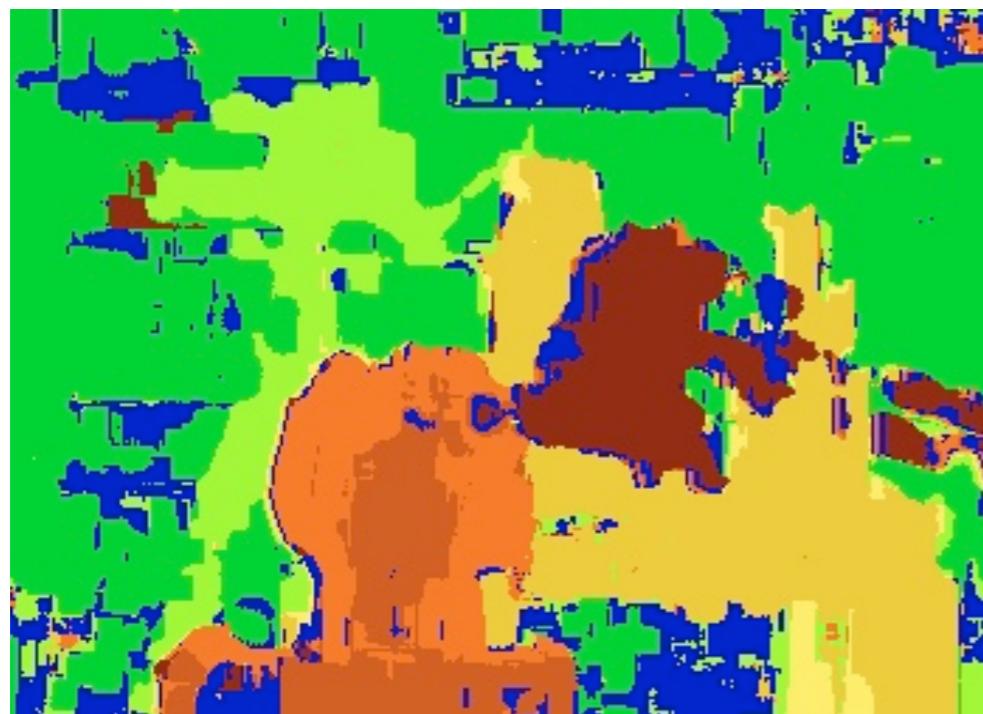


Ground truth

[Seitz]

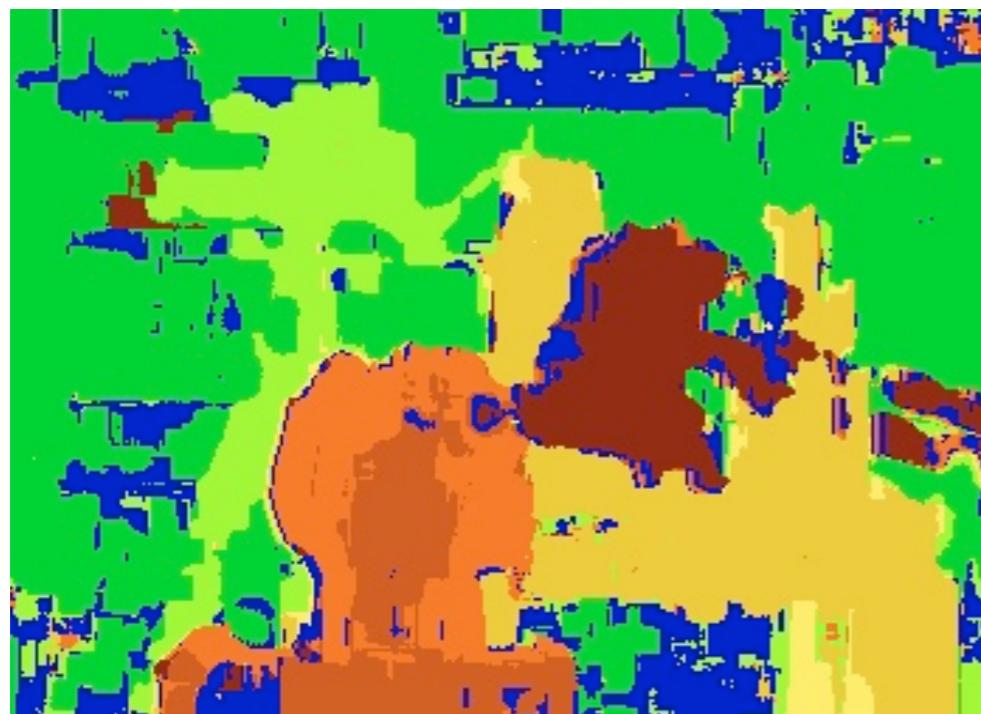
How can we improve window-based matching?

- ◆ The similarity constraint is local
 - ◆ each reference window is matched independently
 - ◆ other points do not influence the result



How can we improve window-based matching?

- ◆ The similarity constraint is local
 - ◆ each reference window is matched independently
 - ◆ other points do not influence the result
- ◆ Need to enforce non-local correspondence constraints



Where do we go from here?



Graph cuts



Ground truth

- ◆ How can be improve the results?
 - ◆ We need to impose spatial regularity to improve results!
 - ◆ Stay tuned for Computer Vision 2