

# Computer Vision

CS-E4850, 5 study credits

Lecturer: Juho Kannala

# Lecture 11: Two-view geometry & stereo vision

- **Two-view geometry** (a.k.a. epipolar geometry) describes the geometric constraints between two views
- **Stereo vision** is the principle of using two views to measure depths of scene points

**Acknowledgement:** many slides from Svetlana Lazebnik, Steve Seitz, Yuri Boykov, Noah Snavely, and others (detailed credits on individual slides)

# Reading

- Szeliski's book, Section 7.2 and Chapter 11

and/or

- Hartley & Zisserman book, Chapters 9-12

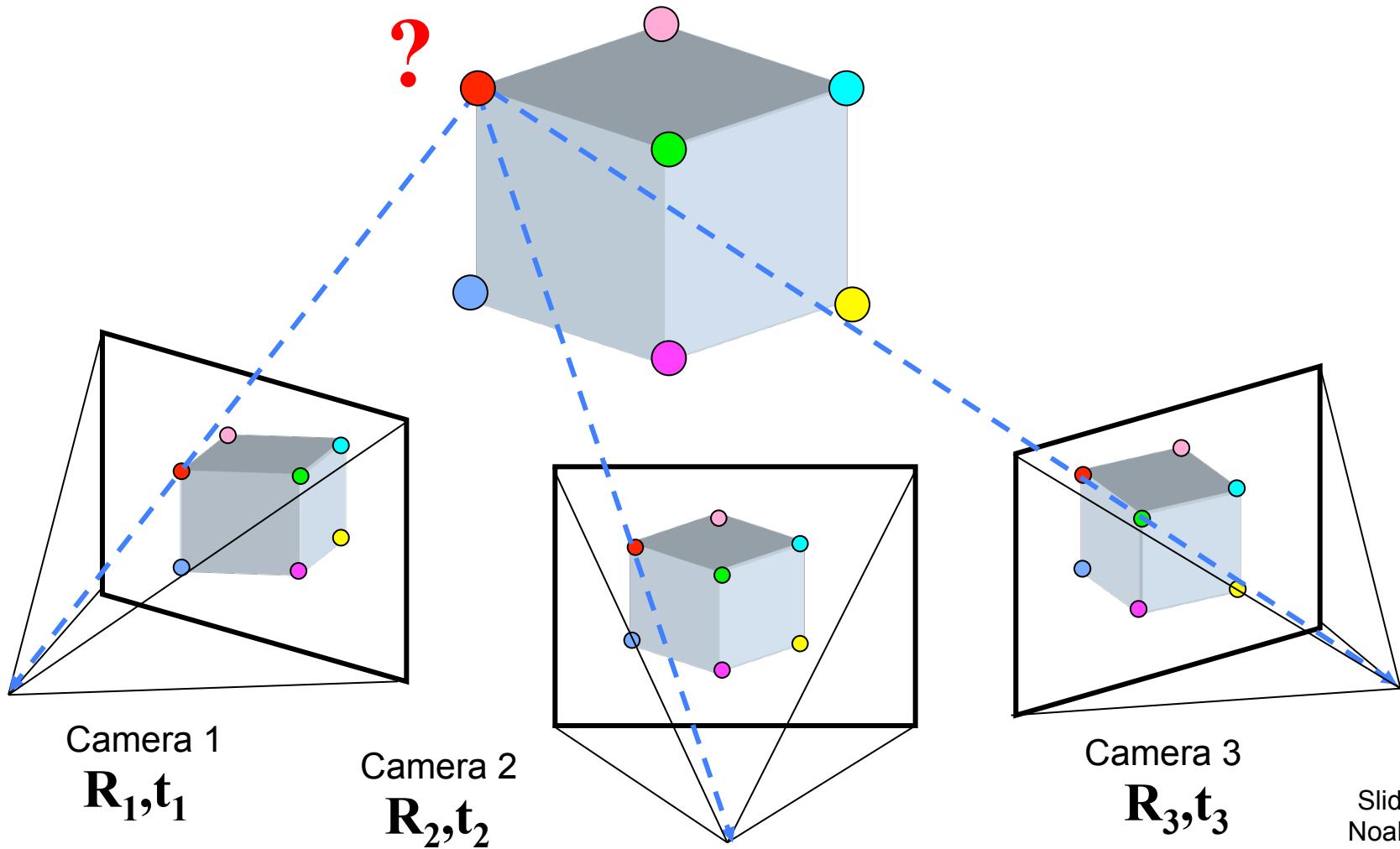
# Multi-view geometry

---



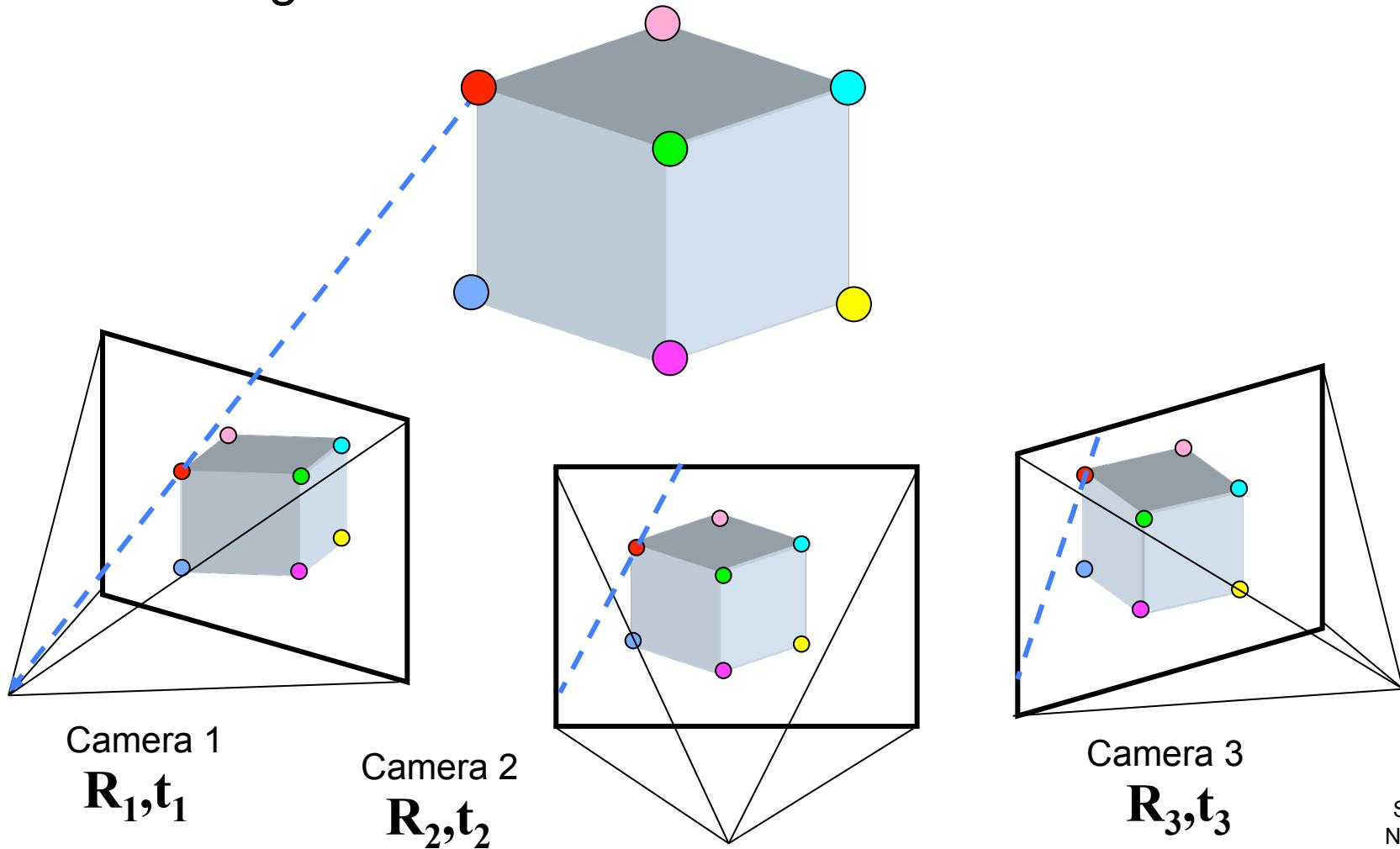
# Multi-view geometry problems

- **Structure:** Given projections of the same 3D point in two or more images, compute the 3D coordinates of that point



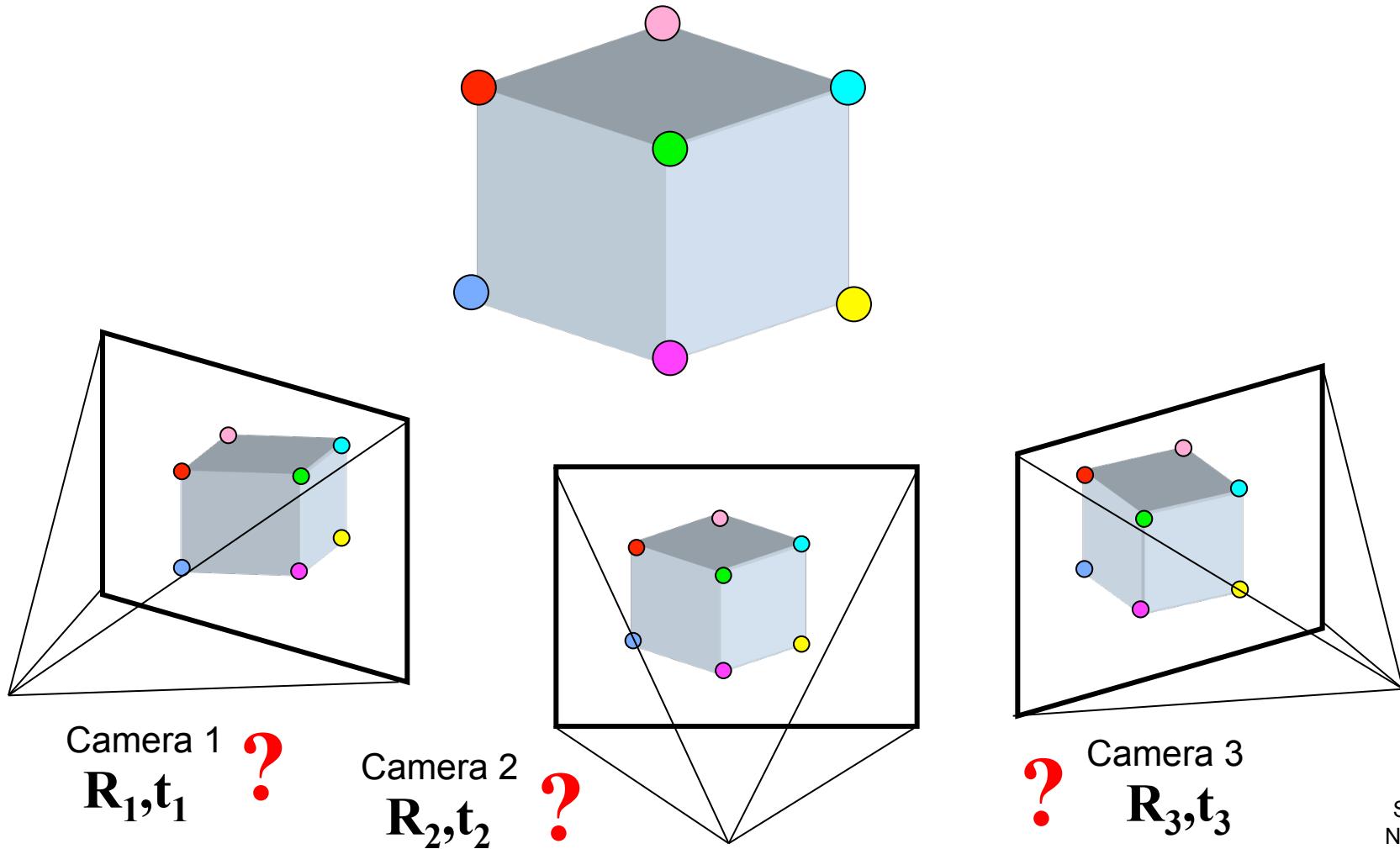
# Multi-view geometry problems

- **Stereo correspondence:** Given a point in one of the images, where could its corresponding points be in the other images?



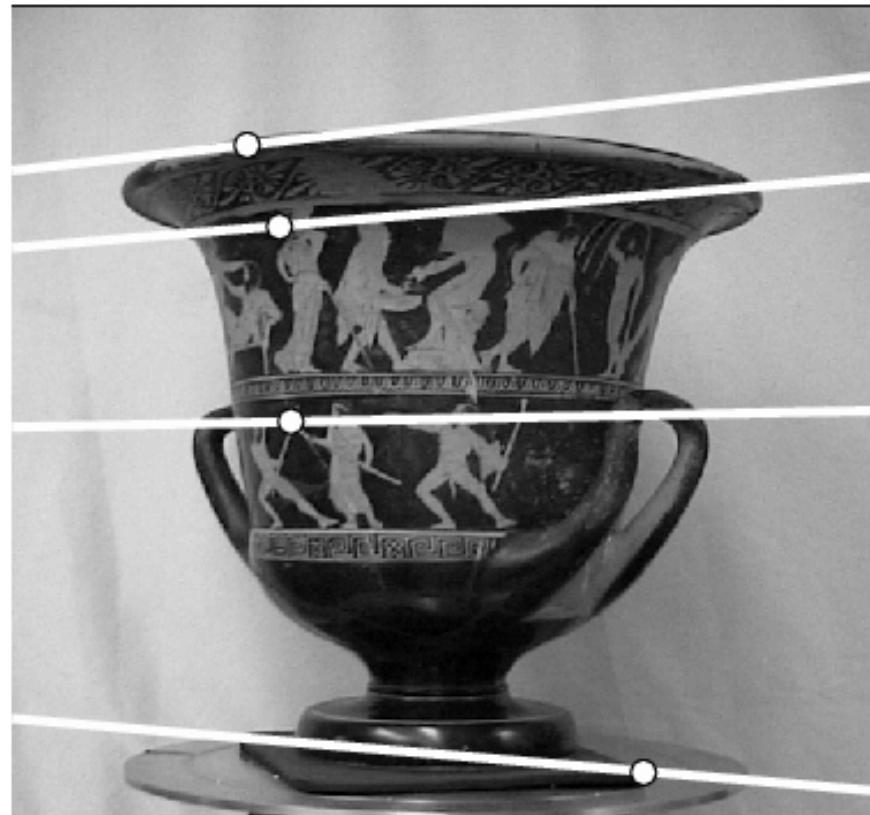
# Multi-view geometry problems

- **Motion:** Given a set of corresponding points in two or more images, compute the camera parameters

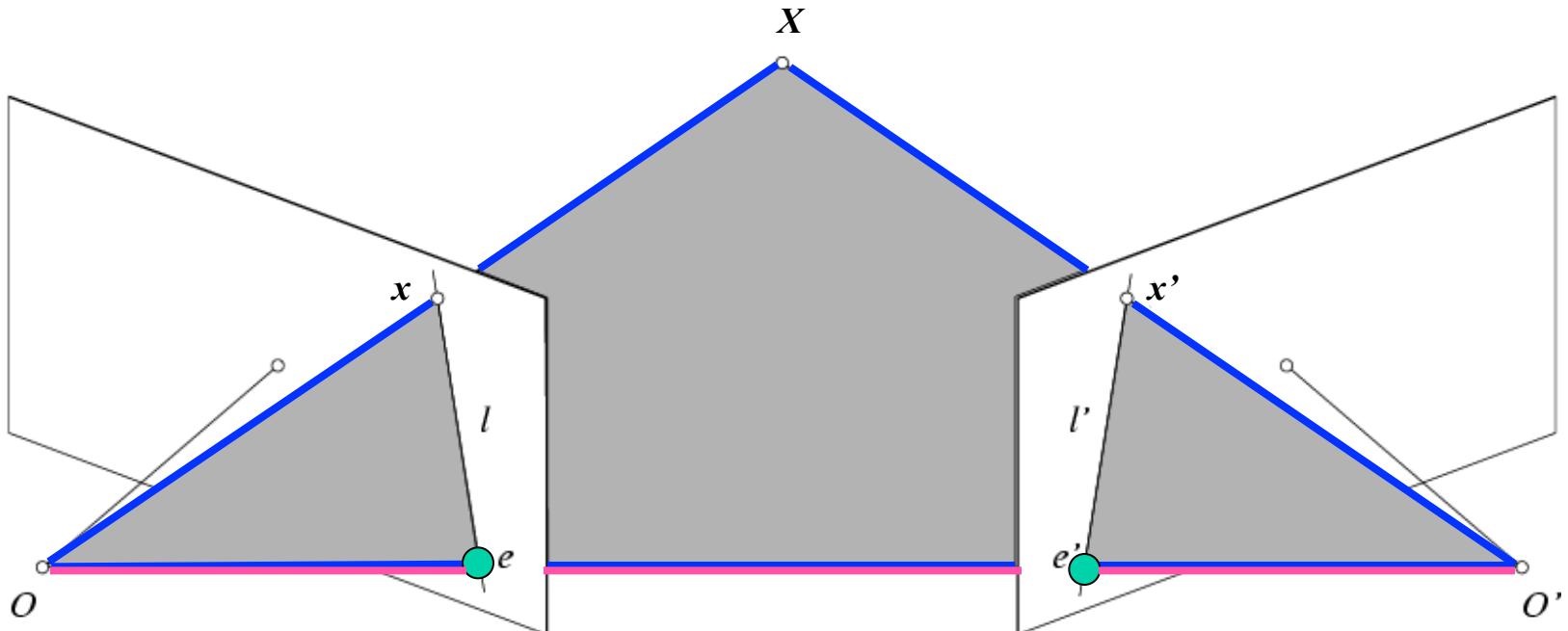


# Two-view geometry

---

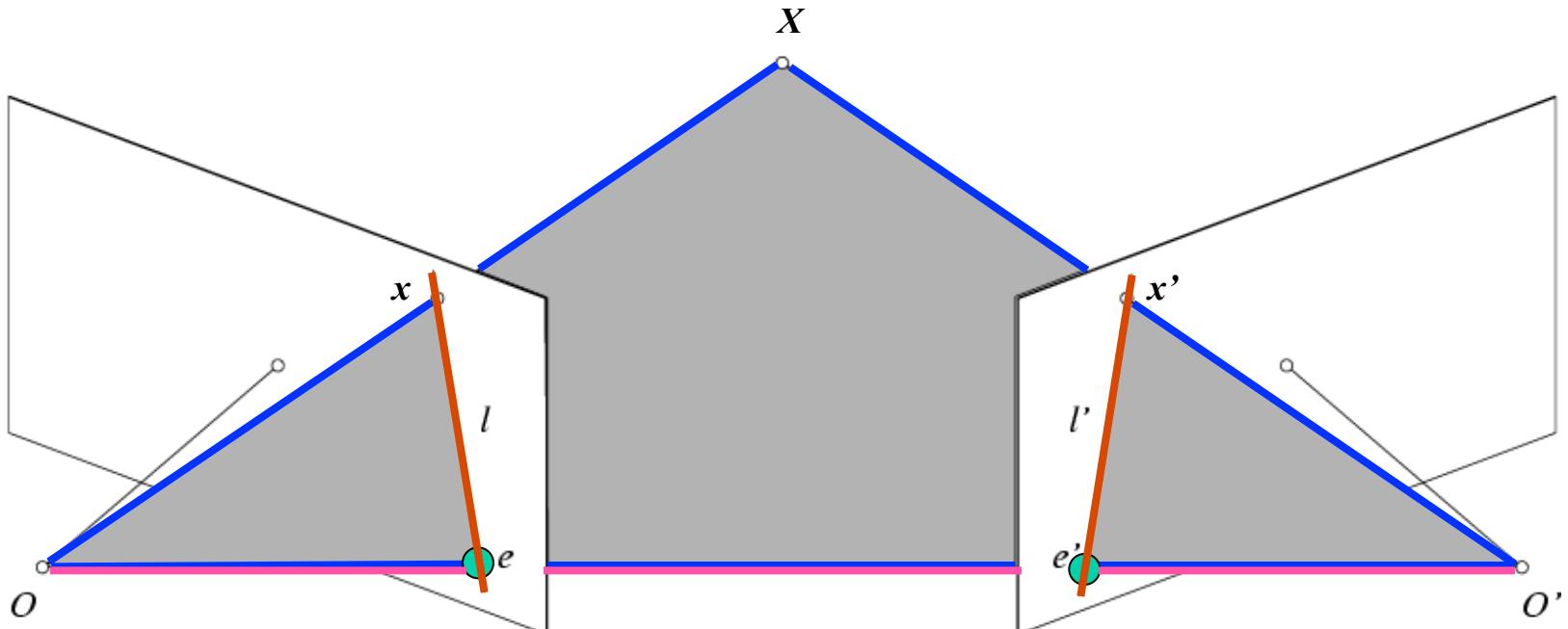


# Epipolar geometry



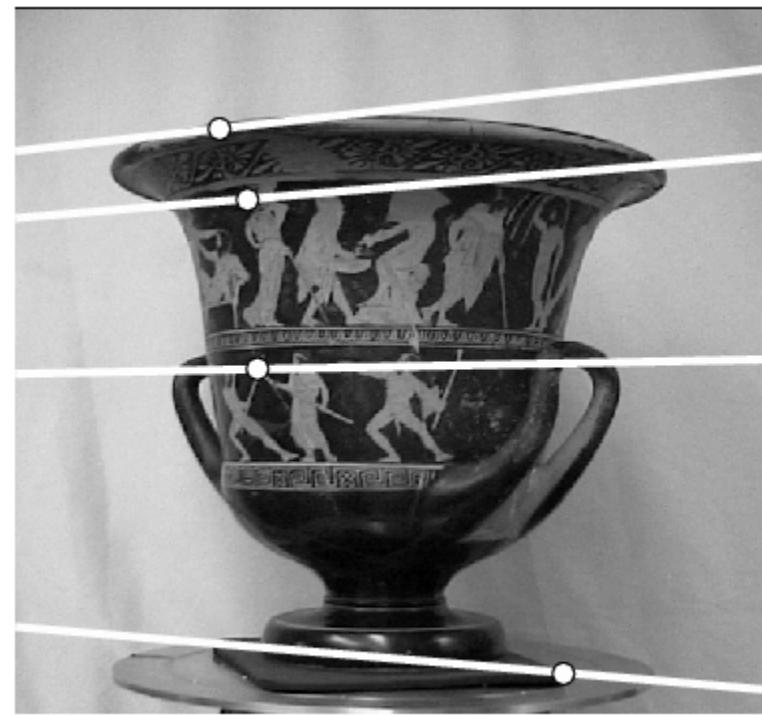
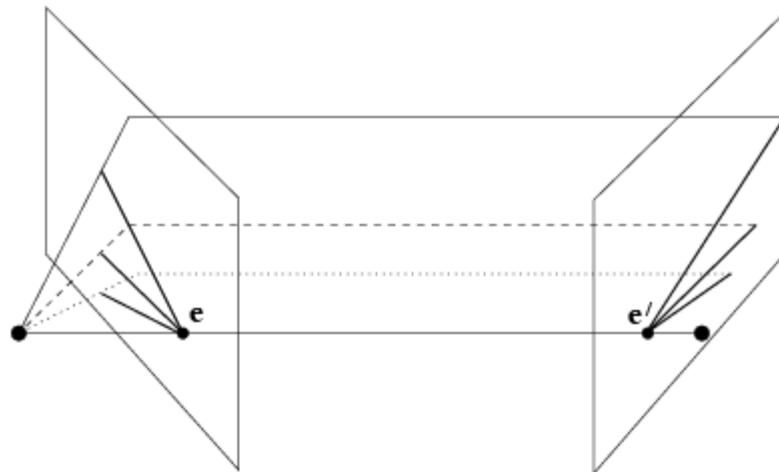
- **Baseline** – line connecting the two camera centers
- **Epipolar Plane** – plane containing baseline (1D family)
- **Epipoles**
  - = intersections of baseline with image planes
  - = projections of the other camera center
  - = vanishing points of the motion direction

# Epipolar geometry



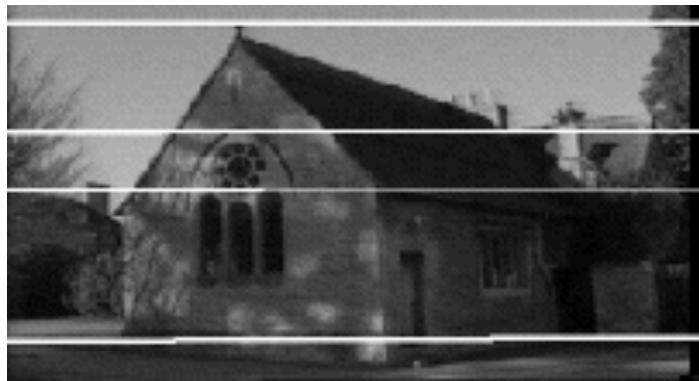
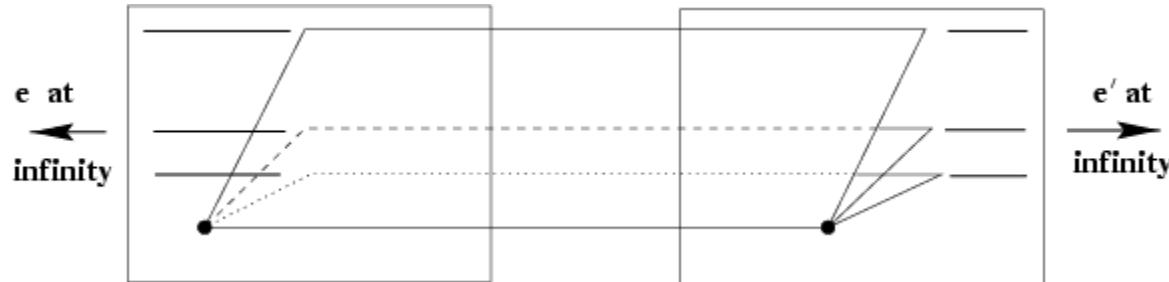
- **Baseline** – line connecting the two camera centers
- **Epipolar Plane** – plane containing baseline (1D family)
- **Epipoles**
  - = intersections of baseline with image planes
  - = projections of the other camera center
  - = vanishing points of the motion direction
- **Epipolar Lines** - intersections of epipolar plane with image planes (always come in corresponding pairs)

# Example: Converging cameras



# Example: Motion parallel to image plane

---



# Example: Motion perpendicular to image plane

---



Source: S. Lazebnik

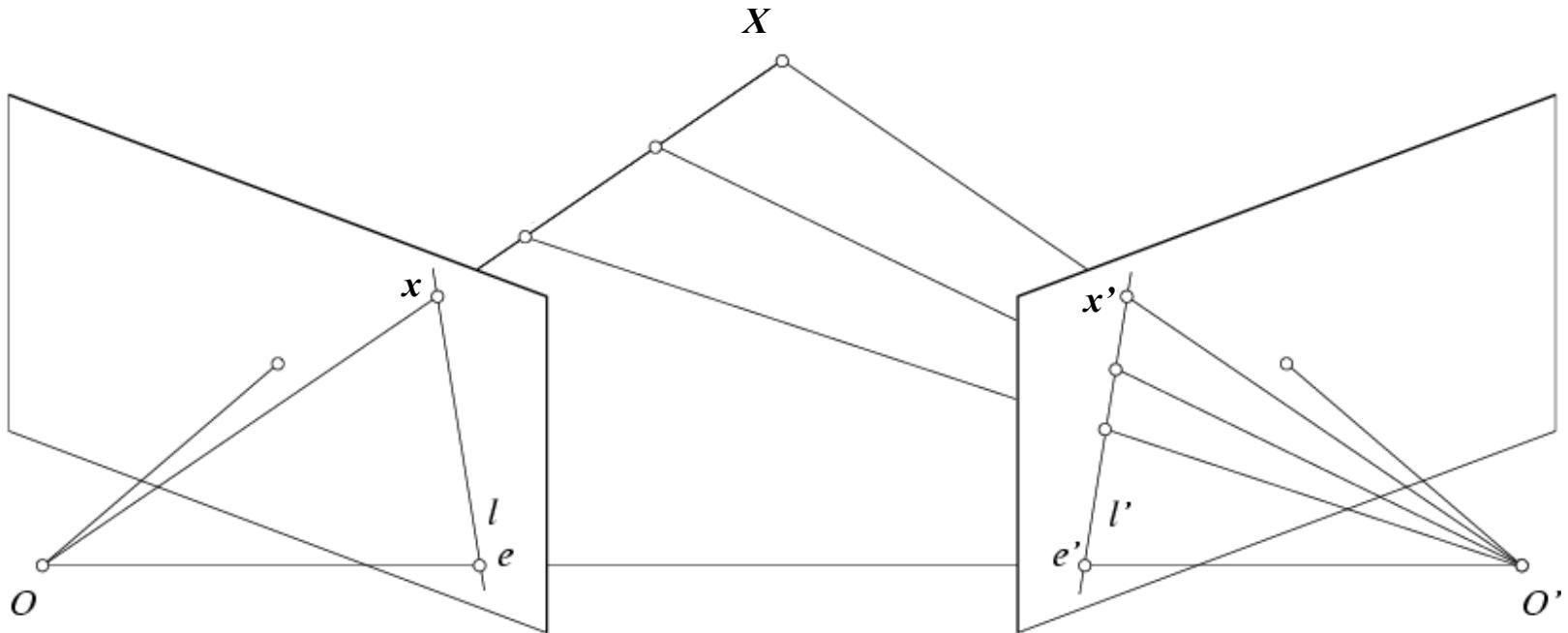
# Example: Motion perpendicular to image plane

---



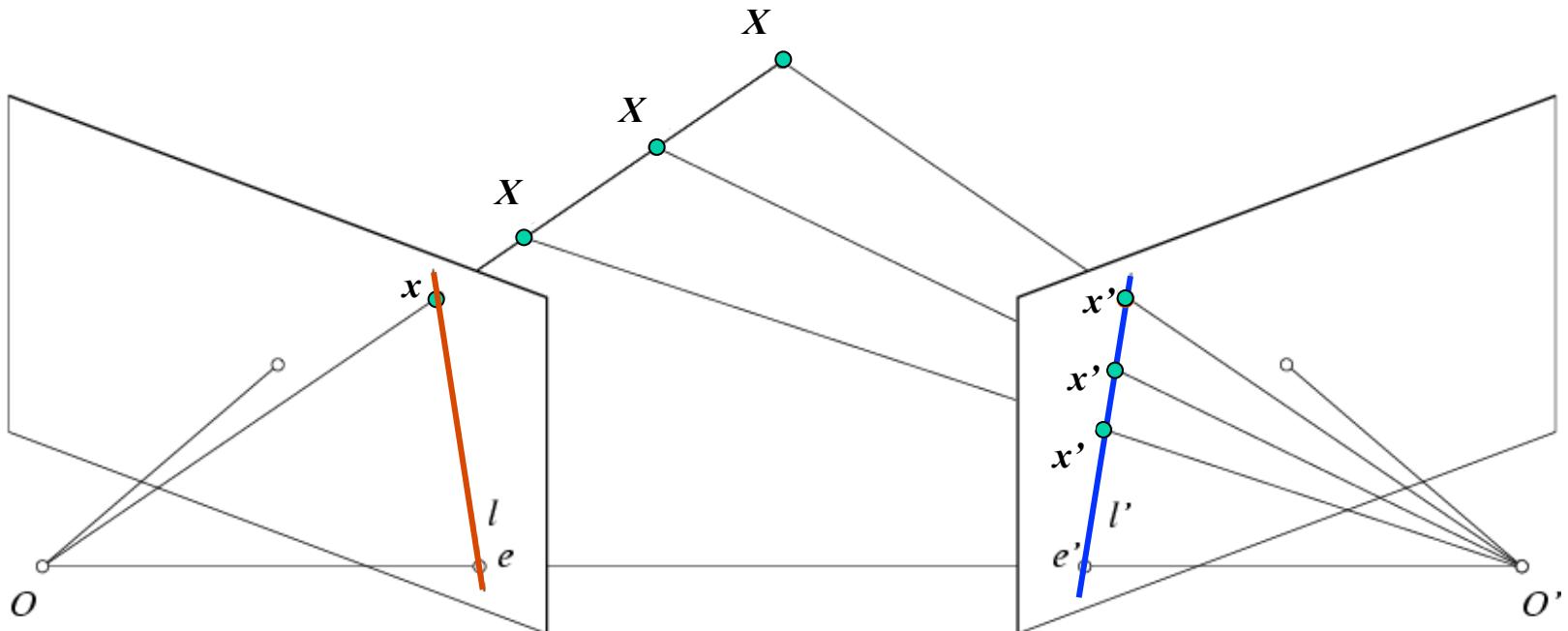
- Points move along lines radiating from the epipole: “focus of expansion”
- Epipole is the principal point

# Epipolar constraint



- If we observe a point  $x$  in one image, where can the corresponding point  $x'$  be in the other image?

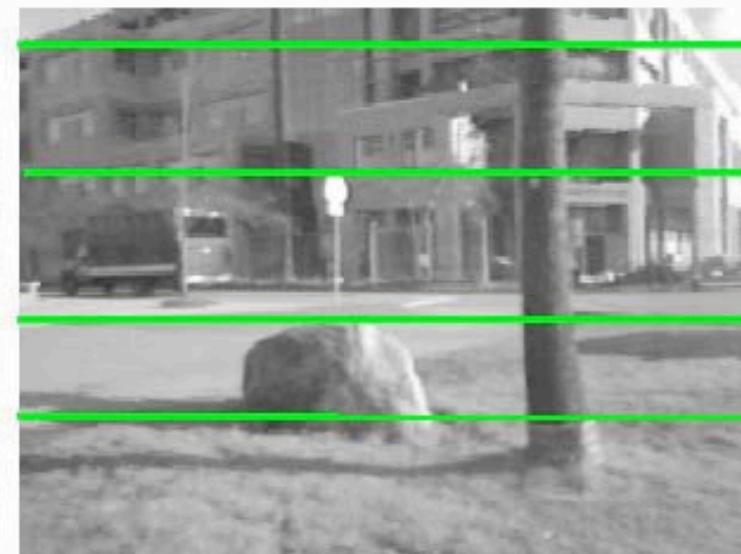
# Epipolar constraint



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

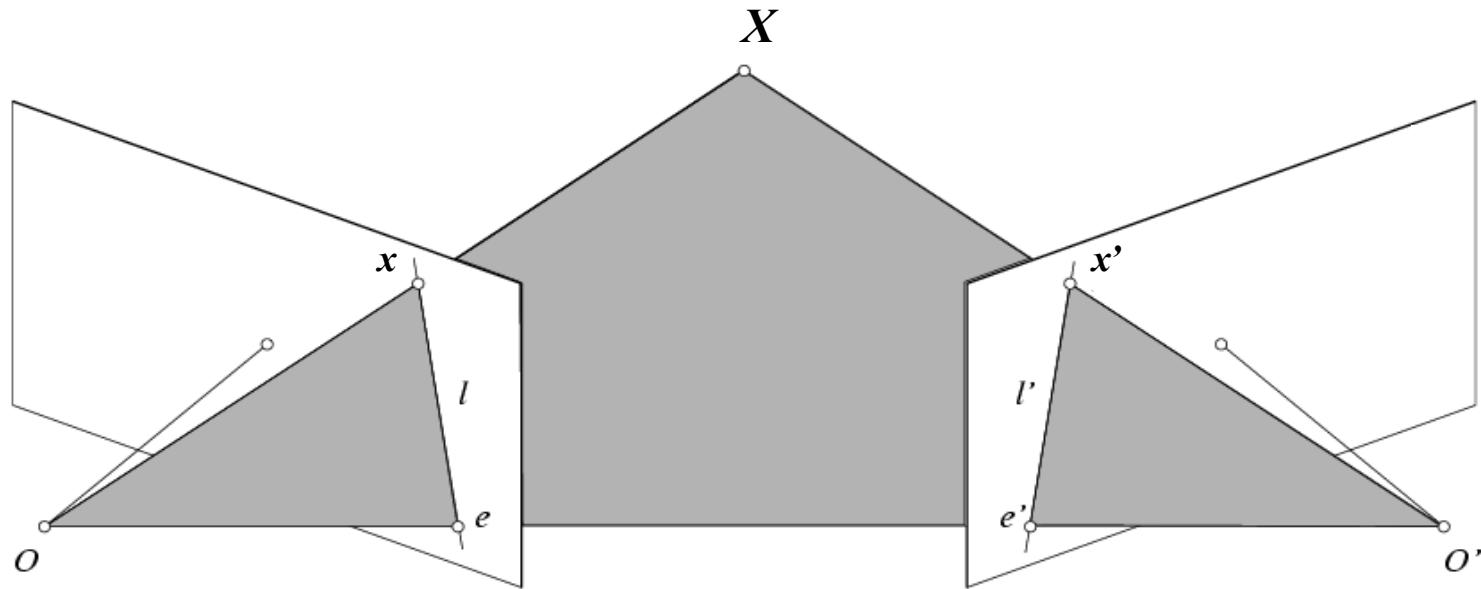
# Epipolar constraint example

---



# Epipolar constraint: Calibrated case

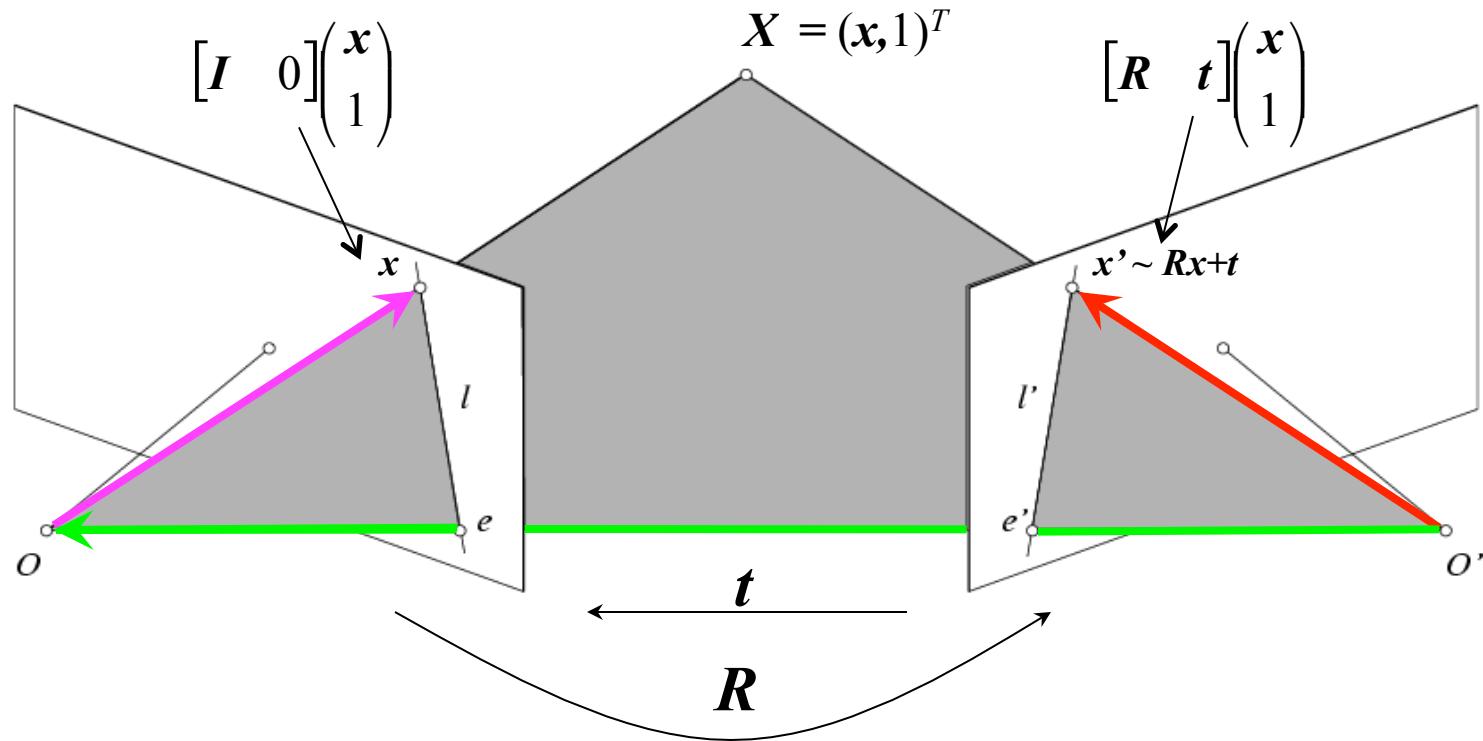
---



- Intrinsic and extrinsic parameters of the cameras are known, world coordinate system is set to that of the first camera
- Then the projection matrices are given by  $K[I \mid \mathbf{0}]$  and  $K'[R \mid t]$
- We can multiply the projection matrices (and the image points) by the inverse of the calibration matrices to get *normalized* image coordinates:

$$\mathbf{x}_{\text{norm}} = K^{-1} \mathbf{x}_{\text{pixel}} = [I \ 0] X, \quad \mathbf{x}'_{\text{norm}} = K'^{-1} \mathbf{x}'_{\text{pixel}} = [R \ t] X$$

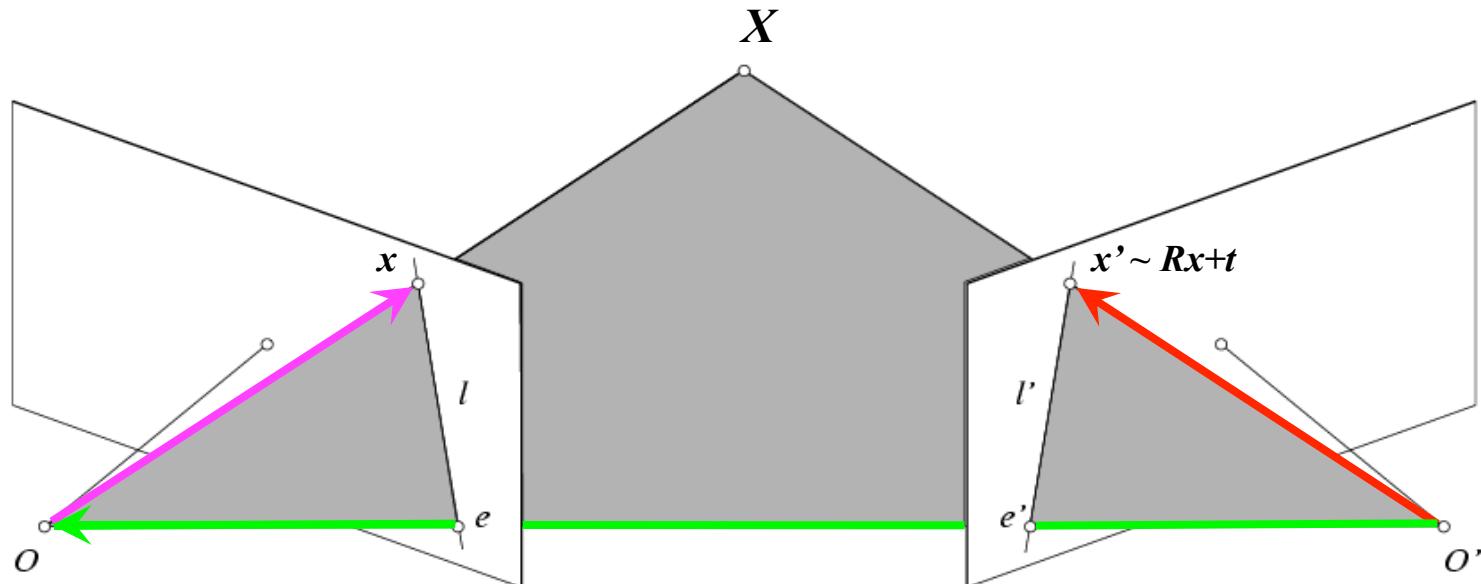
# Epipolar constraint: Calibrated case



The vectors  $Rx$ ,  $t$ , and  $x'$  are coplanar

# Epipolar constraint: Calibrated case

---

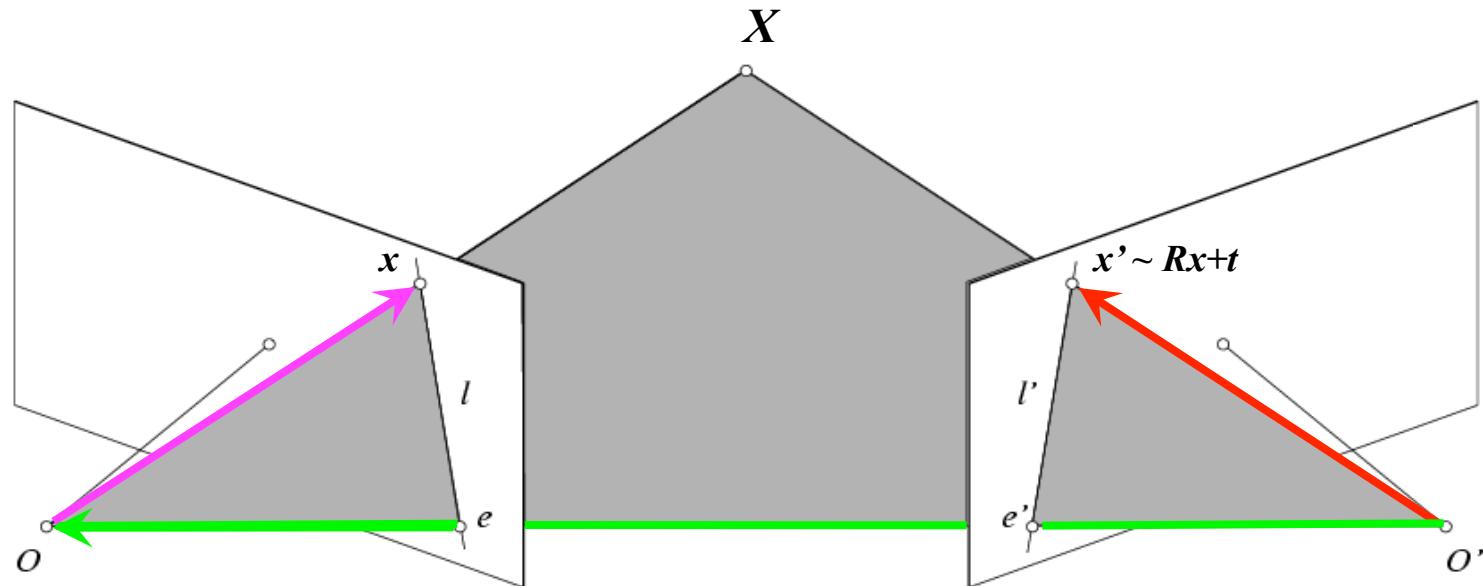


$$\mathbf{x}' \cdot [\mathbf{t} \times (\mathbf{R}\mathbf{x})] = 0 \quad \rightarrow \quad \mathbf{x}'^T [\mathbf{t}_x]^T \mathbf{R} \mathbf{x} = 0$$

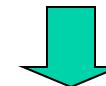
Recall:  $\mathbf{a} \times \mathbf{b} = \begin{bmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{bmatrix} \begin{bmatrix} b_x \\ b_y \\ b_z \end{bmatrix} = [\mathbf{a}_x] \mathbf{b}$

The vectors  $\mathbf{Rx}$ ,  $\mathbf{t}$ , and  $\mathbf{x}'$  are coplanar

# Epipolar constraint: Calibrated case



$$x' \cdot [t \times (Rx)] = 0 \quad \rightarrow \quad x'^T [t_x] R x = 0 \quad \rightarrow \quad x'^T E x = 0$$

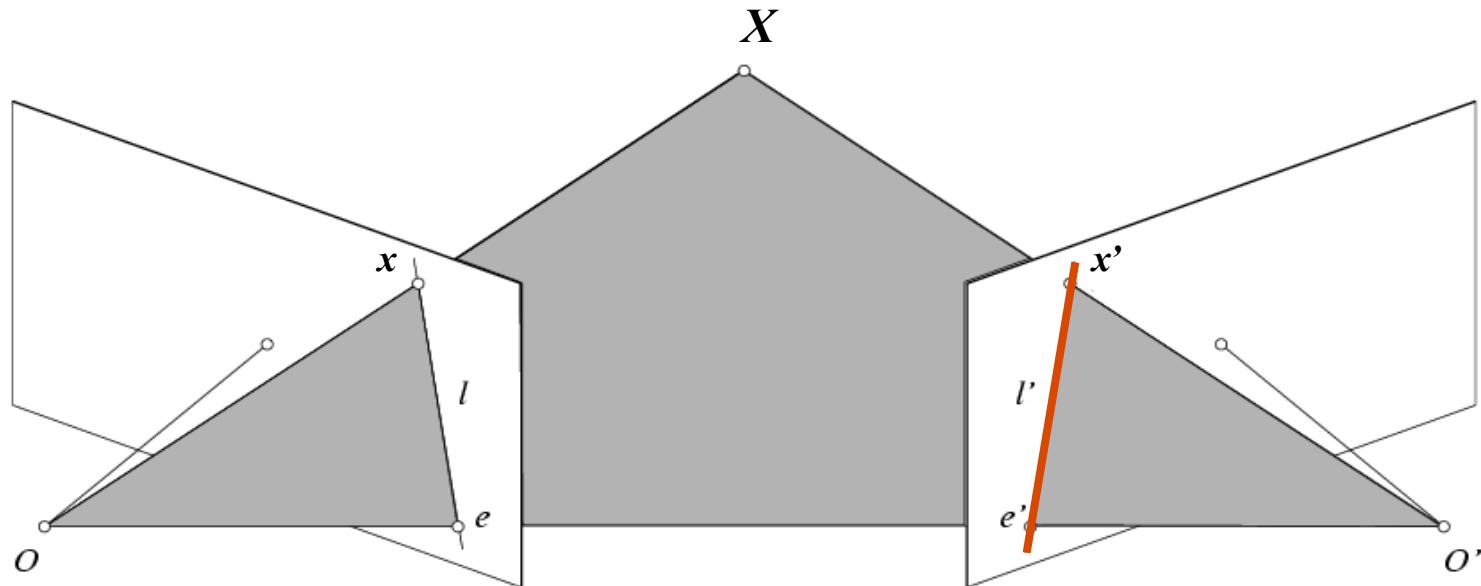


**Essential Matrix**  
(Longuet-Higgins, 1981)

The vectors  $Rx$ ,  $t$ , and  $x'$  are coplanar

# Epipolar constraint: Calibrated case

---



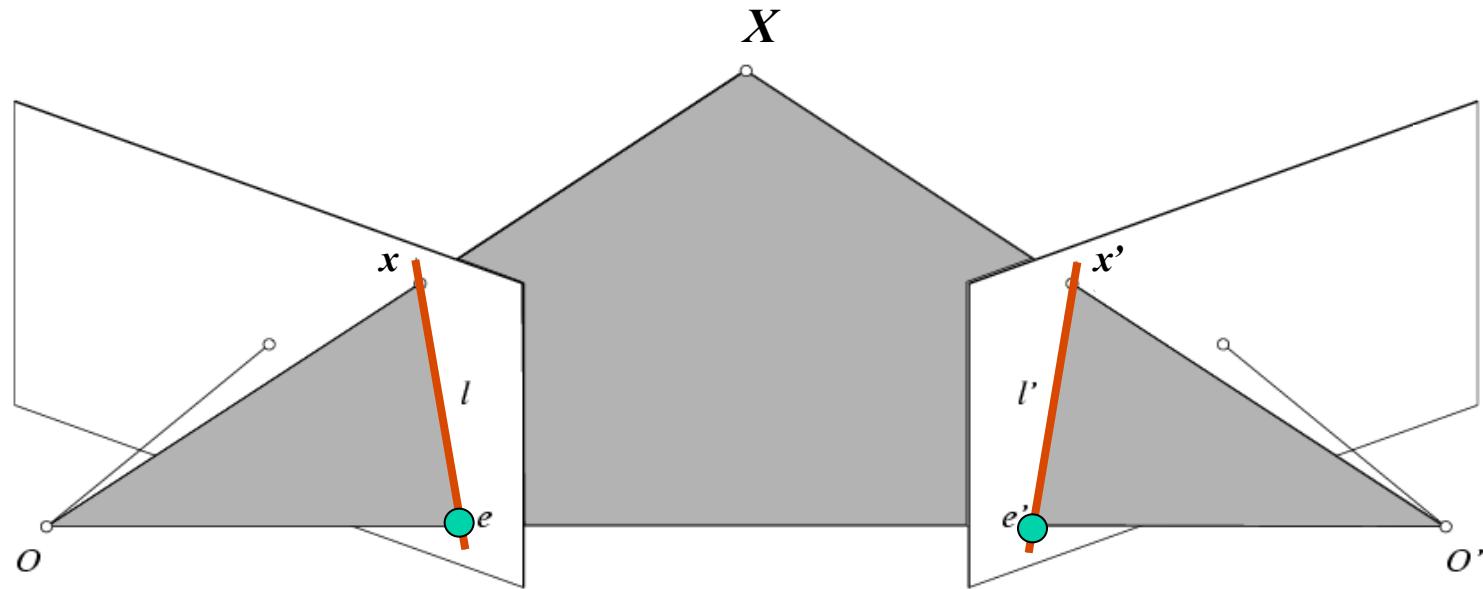
$$\mathbf{x}'^T \mathbf{E} \mathbf{x} = 0$$

- $\mathbf{E} \mathbf{x}$  is the epipolar line associated with  $\mathbf{x}$  ( $\mathbf{l}' = \mathbf{E} \mathbf{x}$ )
  - Recall: a line is given by  $ax + by + c = 0$  or

$$\mathbf{l}^T \mathbf{x} = 0 \quad \text{where} \quad \mathbf{l} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

# Epipolar constraint: Calibrated case

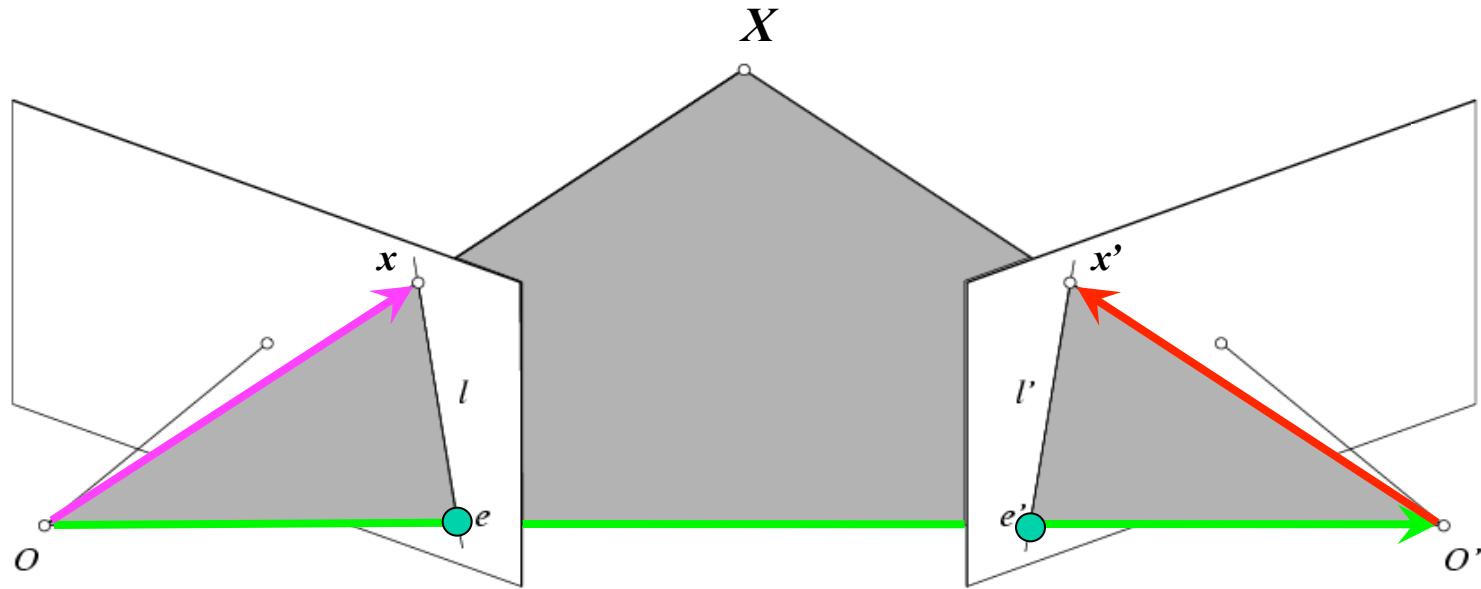
---



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

- $\mathbf{E} \mathbf{x}$  is the epipolar line associated with  $\mathbf{x}$  ( $\mathbf{l}' = \mathbf{E} \mathbf{x}$ )
- $\mathbf{E}^T \mathbf{x}'$  is the epipolar line associated with  $\mathbf{x}'$  ( $\mathbf{l} = \mathbf{E}^T \mathbf{x}'$ )
- $\mathbf{E} \mathbf{e} = 0$  and  $\mathbf{E}^T \mathbf{e}' = 0$
- $\mathbf{E}$  is singular (rank two)
- $\mathbf{E}$  has five degrees of freedom

# Epipolar constraint: Uncalibrated case

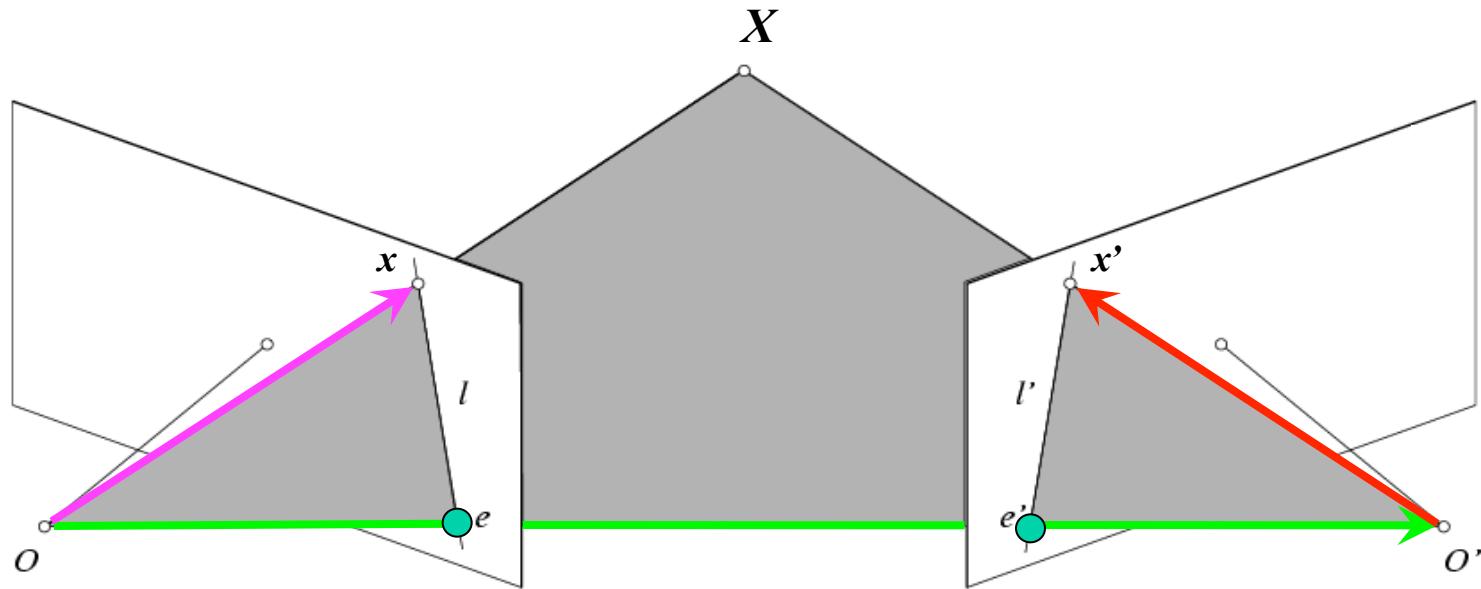


- The calibration matrices  $K$  and  $K'$  of the two cameras are unknown
- We can write the epipolar constraint in terms of *unknown* normalized coordinates:

$$\hat{x}'^T E \hat{x} = 0$$

$$\hat{x} = K^{-1}x, \quad \hat{x}' = K'^{-1}x'$$

# Epipolar constraint: Uncalibrated case



$$\hat{x}'^T E \hat{x} = 0 \quad \xrightarrow{\text{red arrow}} \quad x'^T F x = 0 \quad \text{with} \quad F = K'^{-T} E K^{-1}$$

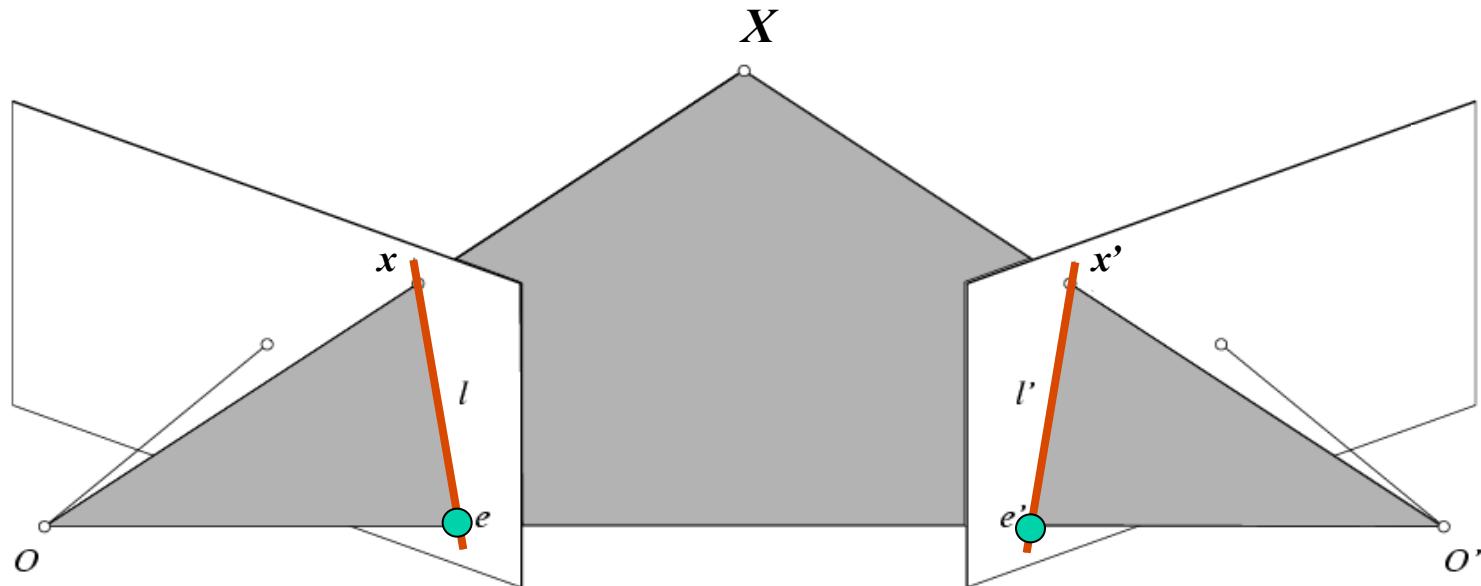
$$\hat{x} = K^{-1} x$$

$$\hat{x}' = K'^{-1} x'$$

**Fundamental Matrix**  
(Faugeras and Luong, 1992)

# Epipolar constraint: Uncalibrated case

---



$$\hat{x}'^T E \hat{x} = 0 \quad \rightarrow \quad x'^T F x = 0 \quad \text{with} \quad F = K'^{-T} E K^{-1}$$

- $Fx$  is the epipolar line associated with  $x$  ( $I' = Fx$ )
- $F^T x'$  is the epipolar line associated with  $x'$  ( $I = F^T x'$ )
- $F\mathbf{e} = 0$  and  $F^T \mathbf{e}' = 0$
- $F$  is singular (rank two)
- $F$  has seven degrees of freedom

# Estimating the fundamental matrix

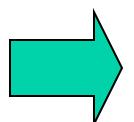
---



# The eight-point algorithm

$$\mathbf{x} = (u, v, 1)^T, \quad \mathbf{x}' = (u', v', 1)$$

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

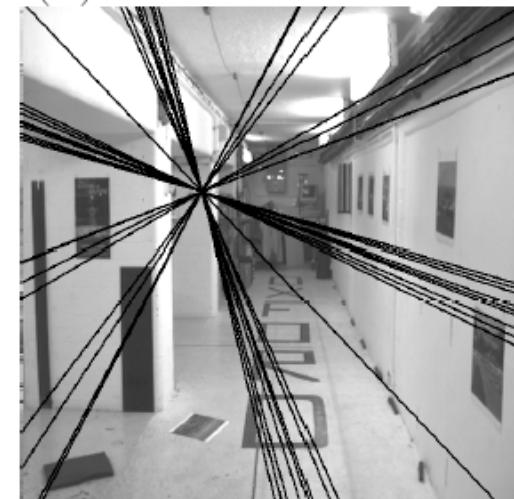
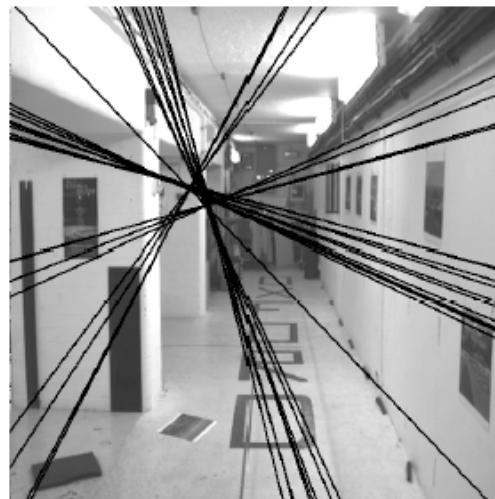


$$\begin{bmatrix} u'u & u'v & u' & v'u & v'v & v' & u & v & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = 0$$

Solve homogeneous  
linear system using  
eight or more matches



Enforce rank-2  
constraint (take SVD  
of  $\mathbf{F}$  and throw out the  
smallest singular value)



# Problem with eight-point algorithm

---

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

# Problem with eight-point algorithm

---

250906.36	183269.57	921.81	200931.10	146766.13	738.21	272.19	198.81
2692.28	131633.03	176.27	6196.73	302975.59	405.71	15.27	746.79
416374.23	871684.30	935.47	408110.89	854384.92	916.90	445.10	931.81
191183.60	171759.40	410.27	416435.62	374125.90	893.65	465.99	418.65
48988.86	30401.76	57.89	298604.57	185309.58	352.87	846.22	525.15
164786.04	546559.67	813.17	1998.37	6628.15	9.86	202.65	672.14
116407.01	2727.75	138.89	169941.27	3982.21	202.77	838.12	19.64
135384.58	75411.13	198.72	411350.03	229127.78	603.79	681.28	379.48

$$\begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \end{bmatrix} = -1$$

Poor numerical conditioning

Can be fixed by rescaling the data

# The normalized eight-point algorithm

---

(Hartley, 1995)

- Center the image data at the origin, and scale it so the mean squared distance between the origin and the data points is 2 pixels
- Use the eight-point algorithm to compute  $\mathbf{F}$  from the normalized points
- Enforce the rank-2 constraint (for example, take SVD of  $\mathbf{F}$  and throw out the smallest singular value)
- Transform fundamental matrix back to original units: if  $\mathbf{T}$  and  $\mathbf{T}'$  are the normalizing transformations in the two images, than the fundamental matrix in original coordinates is  $\mathbf{T}'^T \mathbf{F} \mathbf{T}$

# Nonlinear estimation

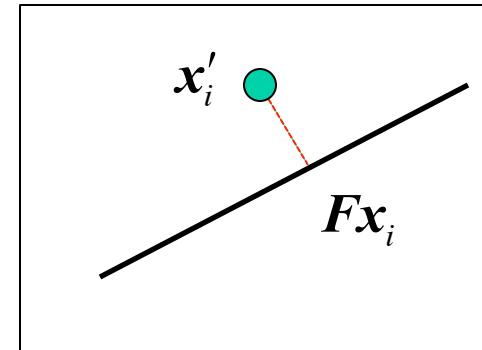
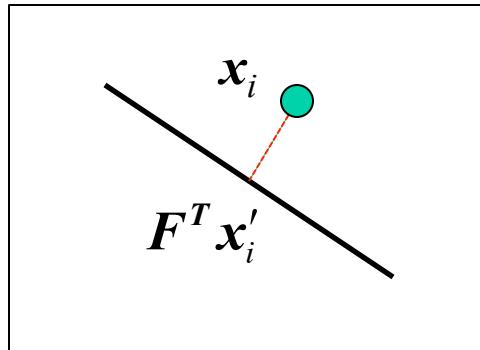
---

- Linear estimation minimizes the sum of squared *algebraic* distances between points  $\mathbf{x}'_i$  and epipolar lines  $\mathbf{F} \mathbf{x}_i$  (or points  $\mathbf{x}_i$  and epipolar lines  $\mathbf{F}^T \mathbf{x}'_i$ ):

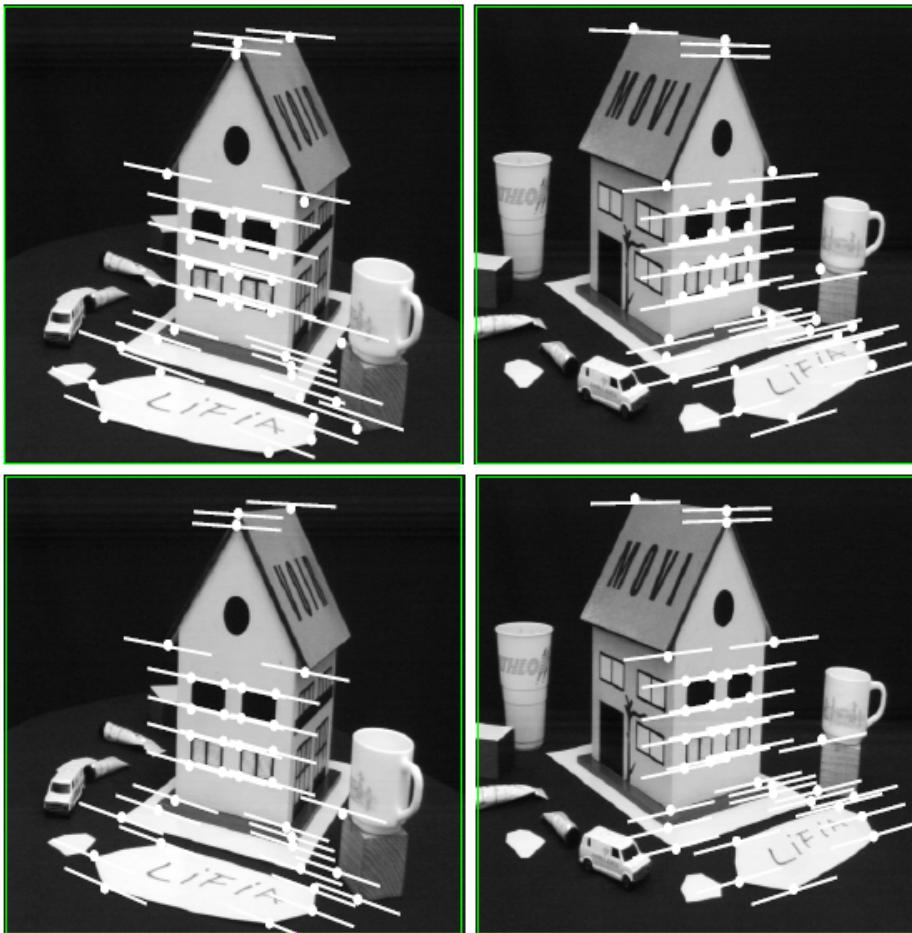
$$\sum_{i=1}^N (\mathbf{x}'_i^T \mathbf{F} \mathbf{x}_i)^2$$

- Nonlinear approach: minimize sum of squared *geometric* distances

$$\sum_{i=1}^N [d^2(\mathbf{x}'_i, \mathbf{F} \mathbf{x}_i) + d^2(\mathbf{x}_i, \mathbf{F}^T \mathbf{x}'_i)]$$



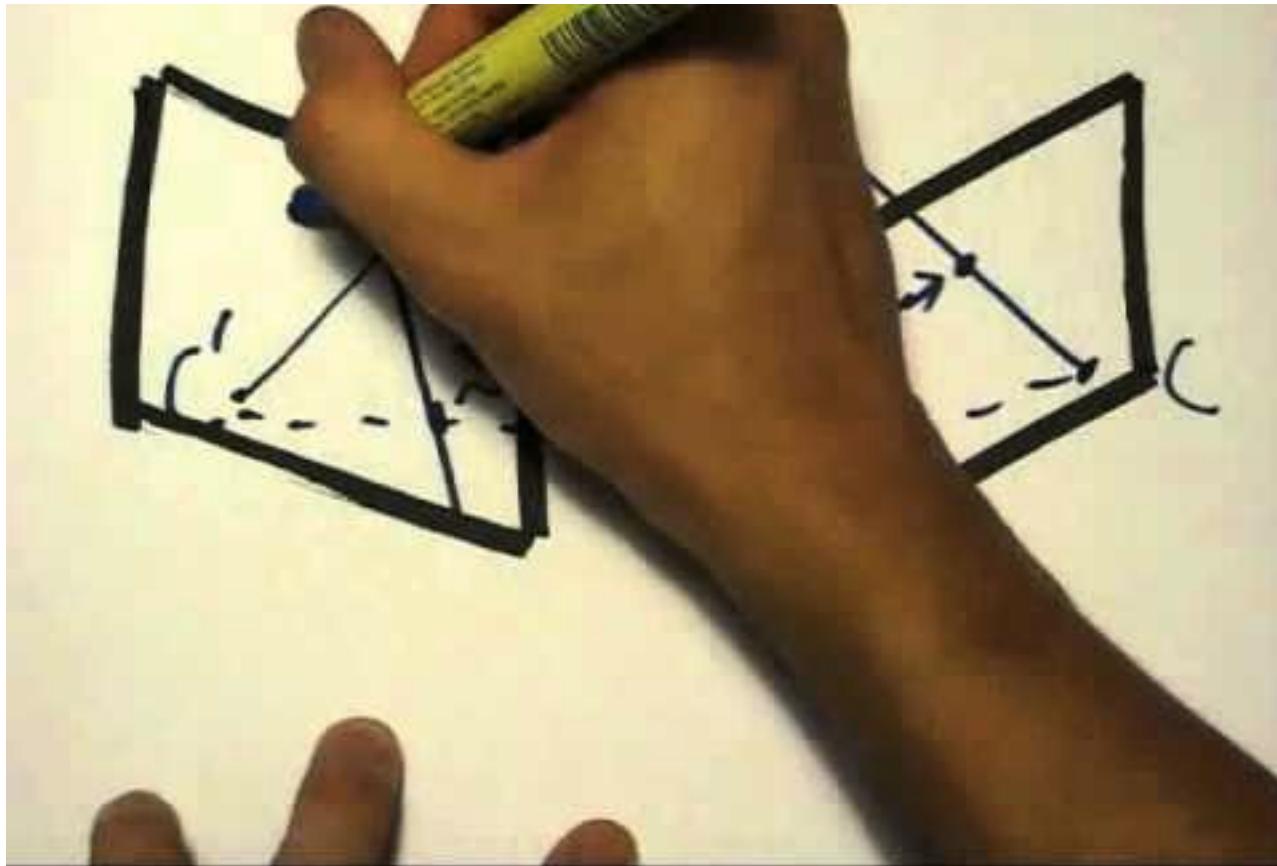
# Comparison of estimation algorithms



	8-point	Normalized 8-point	Nonlinear least squares
Av. Dist. 1	2.33 pixels	0.92 pixel	0.86 pixel
Av. Dist. 2	2.18 pixels	0.85 pixel	0.80 pixel

# The Fundamental Matrix Song

---



<http://danielwedge.com/fmatrix/>

# From epipolar geometry to camera calibration

---

- Estimating the fundamental matrix is known as “weak calibration”
- If we know the calibration matrices of the two cameras, we can estimate the essential matrix:  $E = K'^T F K$
- The essential matrix gives us the relative rotation and translation between the cameras, or their extrinsic parameters

# Stereo

---



Many slides adapted from Steve Seitz

# Binocular stereo

---

- Given a calibrated binocular stereo pair, fuse it to produce a depth image

image 1



image 2



Dense depth map



# Binocular stereo

---

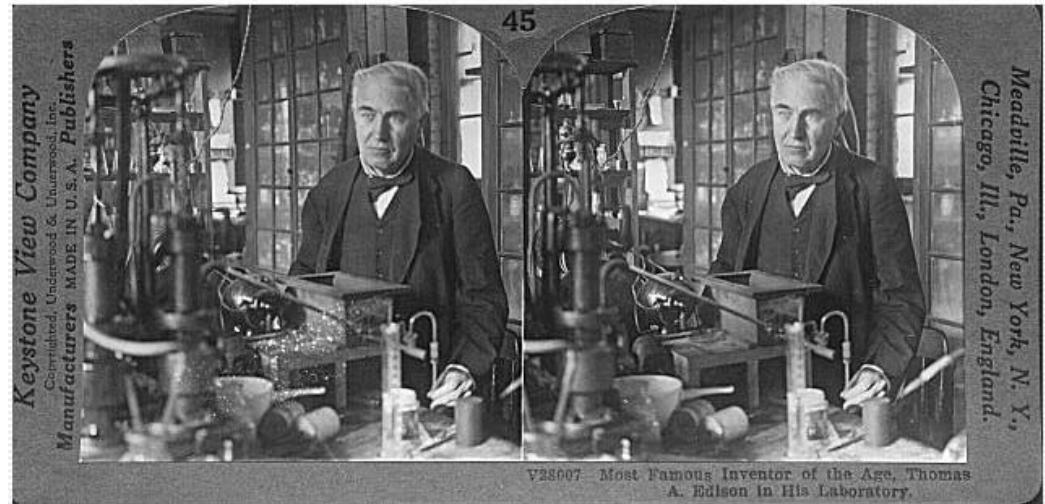
- Given a calibrated binocular stereo pair, fuse it to produce a depth image



Where does the depth information come from?

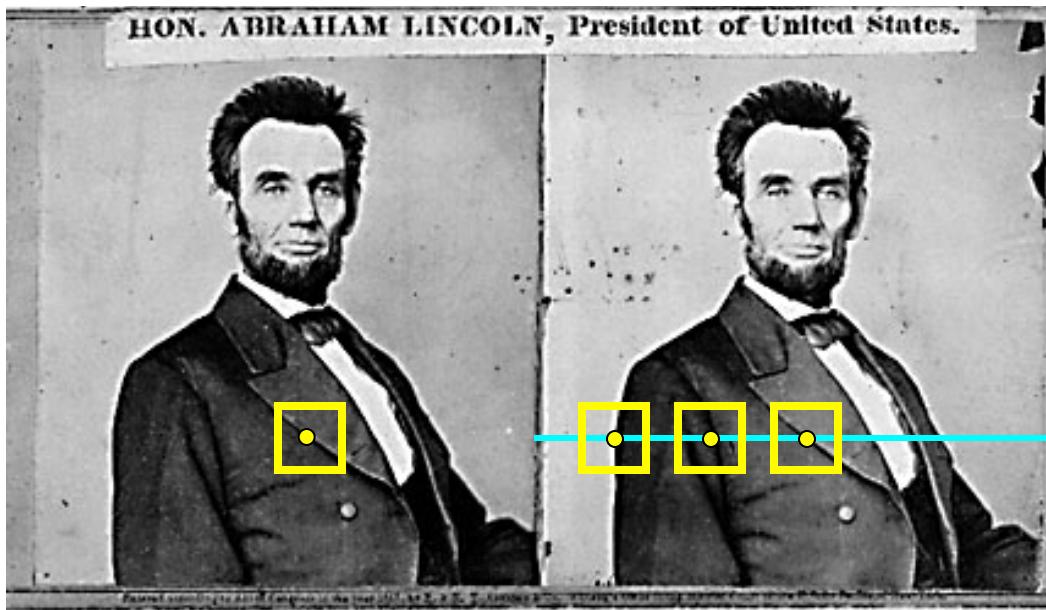
# Binocular stereo

- Given a calibrated binocular stereo pair, fuse it to produce a depth image
  - Humans can do it



Stereograms: Invented by Sir Charles Wheatstone, 1838

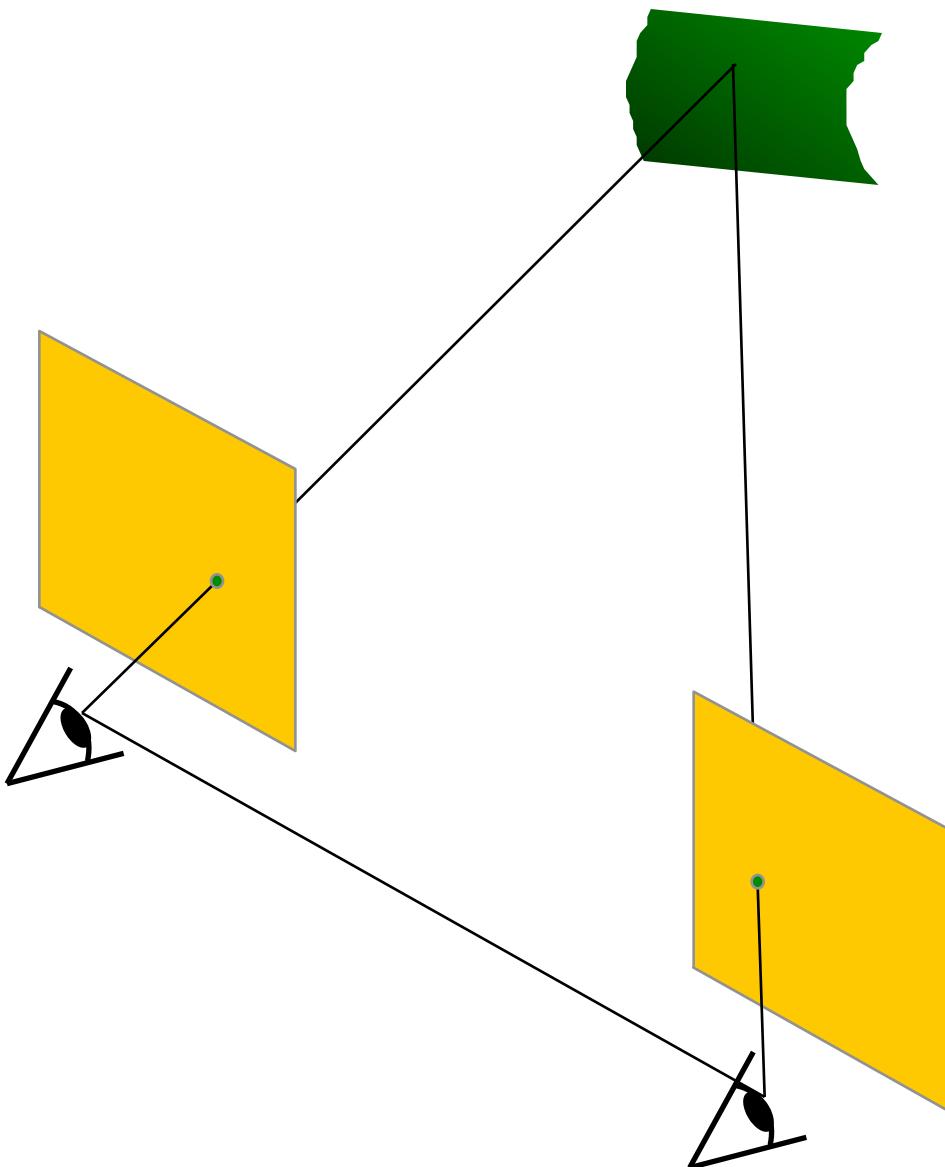
# Basic stereo matching algorithm



- For each pixel in the first image
  - Find corresponding epipolar line in the right image
  - Examine all pixels on the epipolar line and pick the best match
  - Triangulate the matches to get depth information
- Simplest case: epipolar lines are corresponding scanlines
  - When does this happen?

# Simplest Case: Parallel images

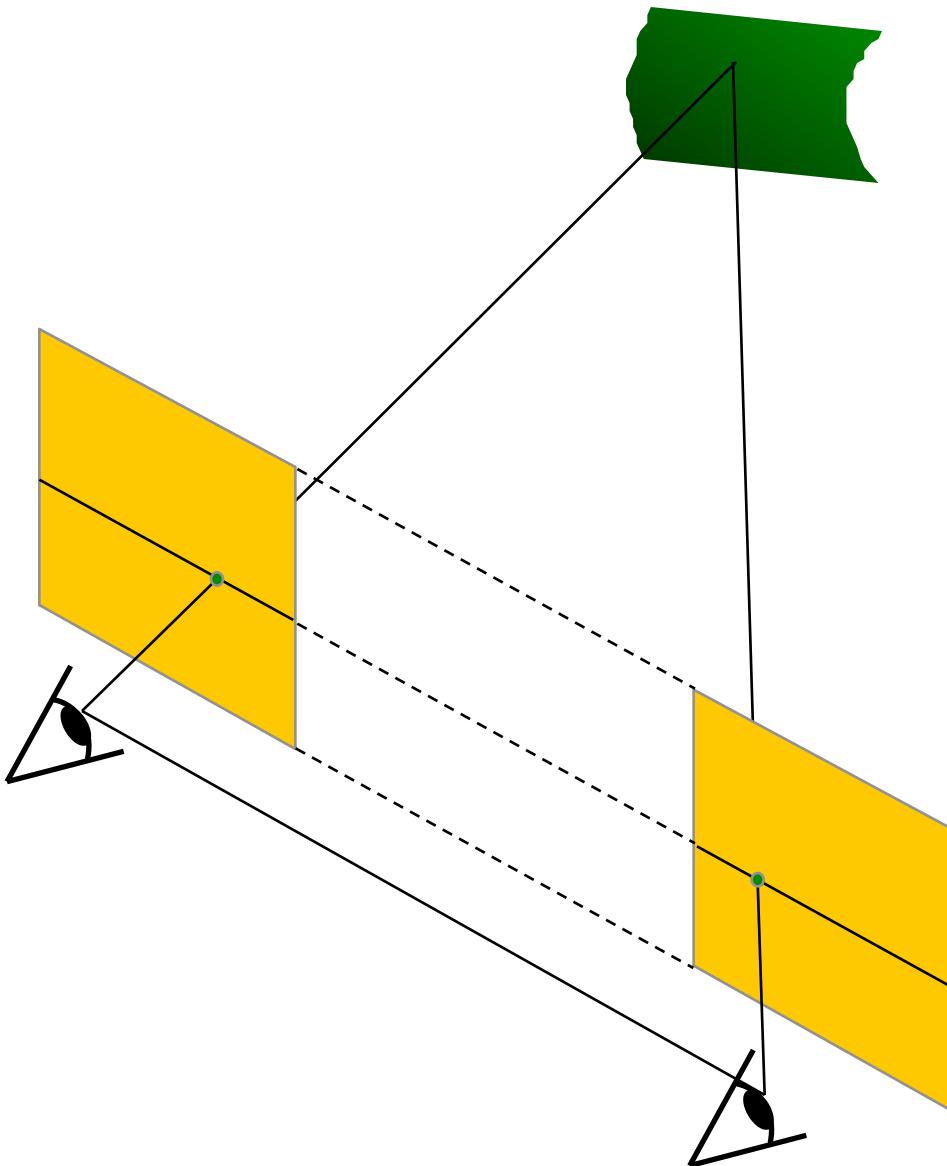
---



- Image planes of cameras are parallel to each other and to the baseline
- Camera centers are at same height
- Focal lengths are the same

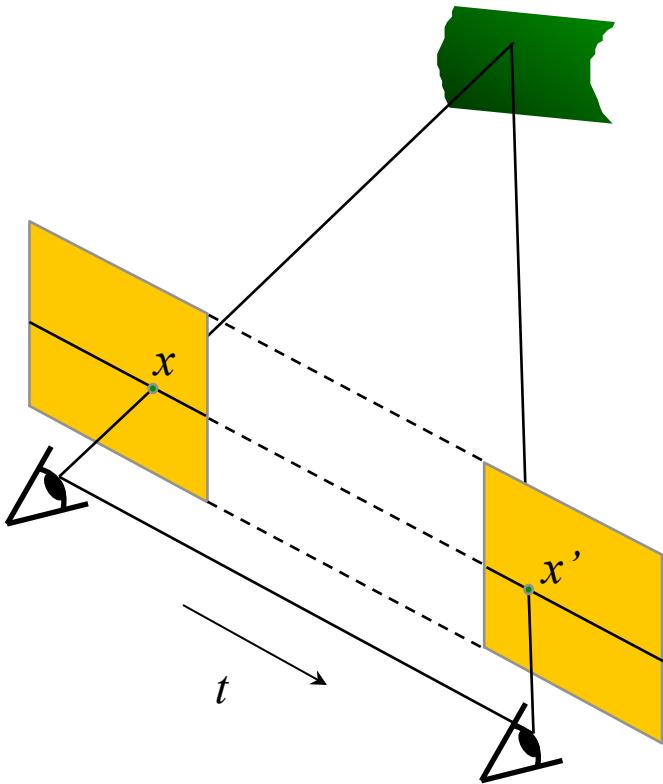
# Simplest Case: Parallel images

---



- Image planes of cameras are parallel to each other and to the baseline
- Camera centers are at same height
- Focal lengths are the same
- Then epipolar lines fall along the horizontal scan lines of the images

# Essential matrix for parallel images



Epipolar constraint:

$$\mathbf{x}'^T \mathbf{E} \mathbf{x} = 0, \quad \mathbf{E} = [\mathbf{t}_x] \mathbf{R}$$

$$\mathbf{R} = \mathbf{I} \quad \mathbf{t} = (T, 0, 0)$$

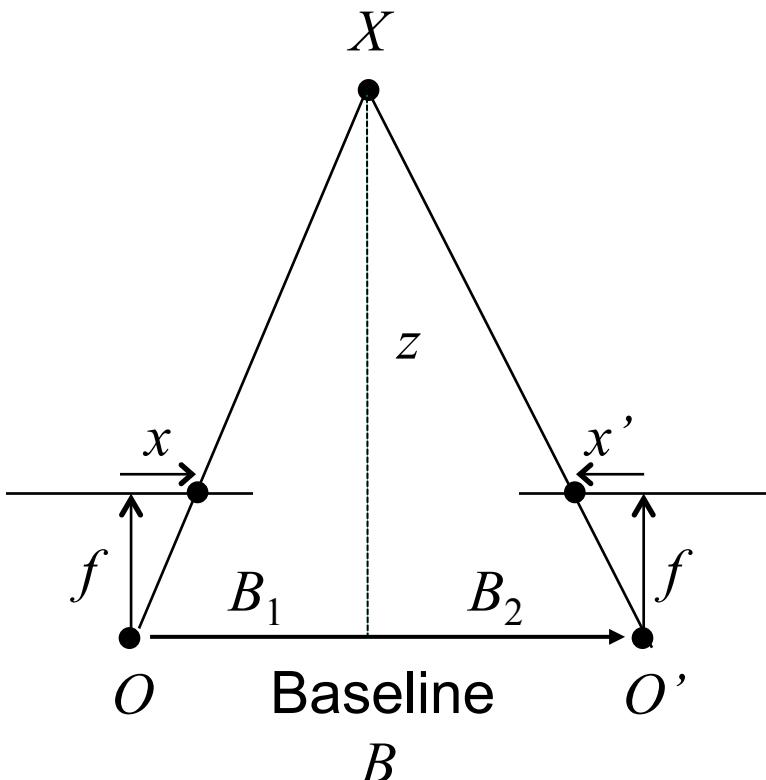
$$\mathbf{E} = [\mathbf{t}_x] \mathbf{R} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -T \\ 0 & T & 0 \end{bmatrix}$$

$$(u' \quad v' \quad 1) \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -T \\ 0 & T & 0 \end{bmatrix} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = 0 \quad (u' \quad v' \quad 1) \begin{pmatrix} 0 \\ -T \\ Tv \end{pmatrix} = 0 \quad Tv' = Tv$$

The y-coordinates of corresponding points are the same!

# Depth from disparity

---



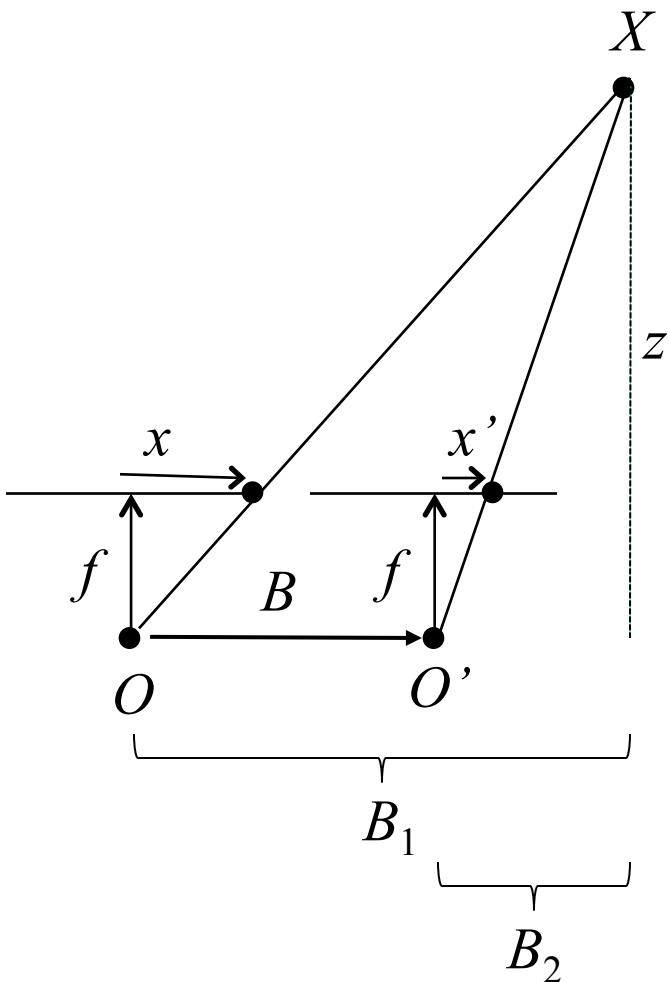
$$\frac{x}{f} = \frac{B_1}{z} \quad \frac{-x'}{f} = \frac{B_2}{z}$$

$$\frac{x - x'}{f} = \frac{B_1 + B_2}{z}$$

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

Disparity is inversely proportional to depth!

# Depth from disparity

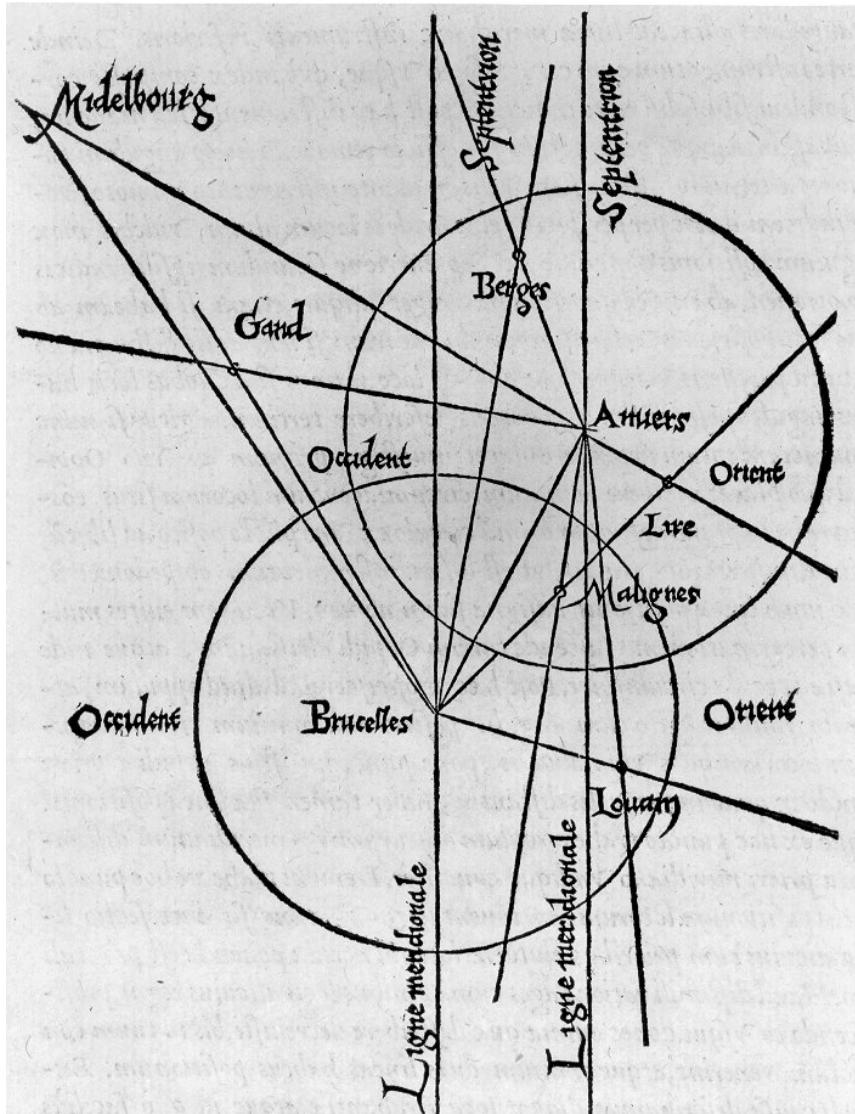


$$\frac{x}{f} = \frac{B_1}{z} \quad \frac{x'}{f} = \frac{B_2}{z}$$

$$\frac{x - x'}{f} = \frac{B_1 - B_2}{z}$$

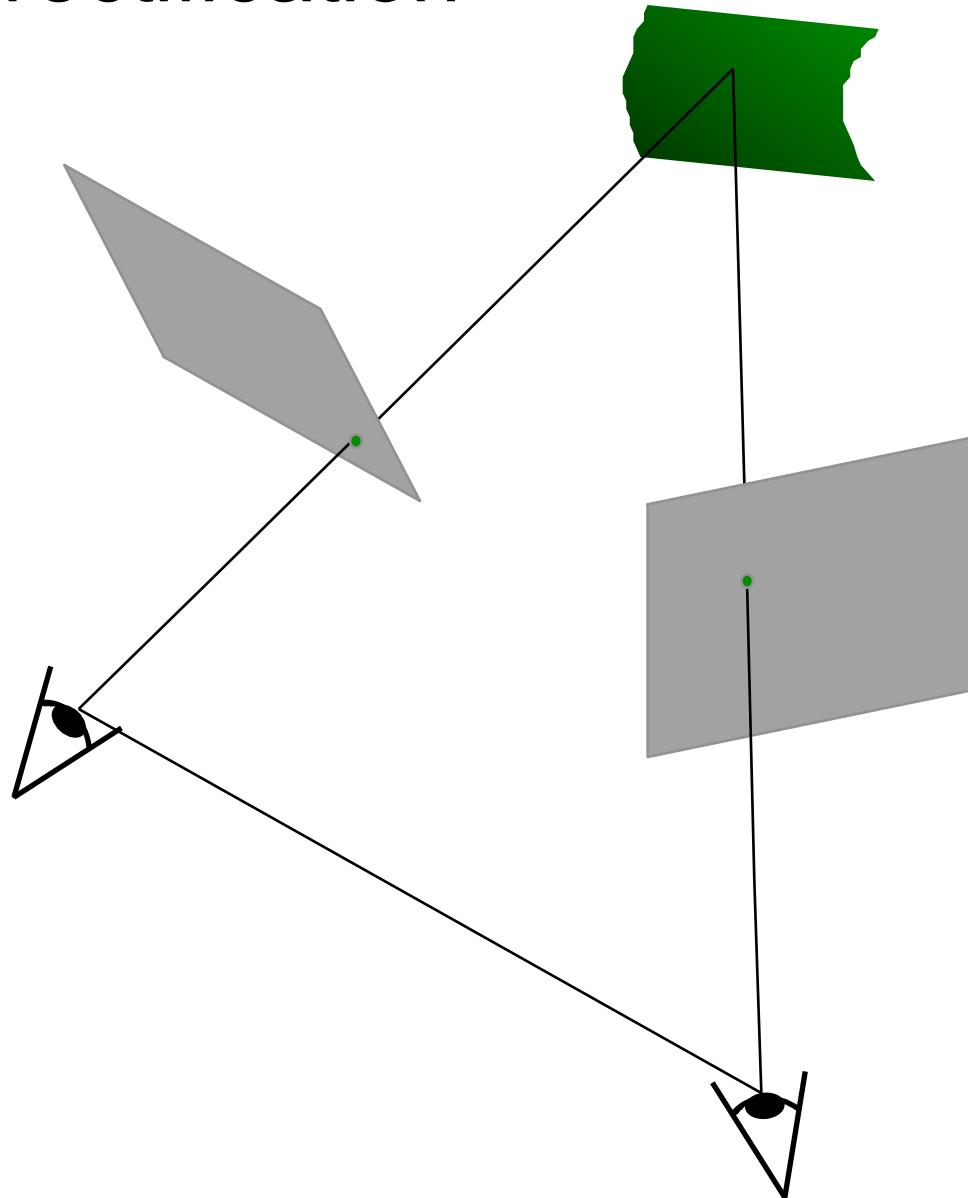
$$\text{disparity} = x - x' = \frac{B \cdot f}{z}$$

# Triangulation: History

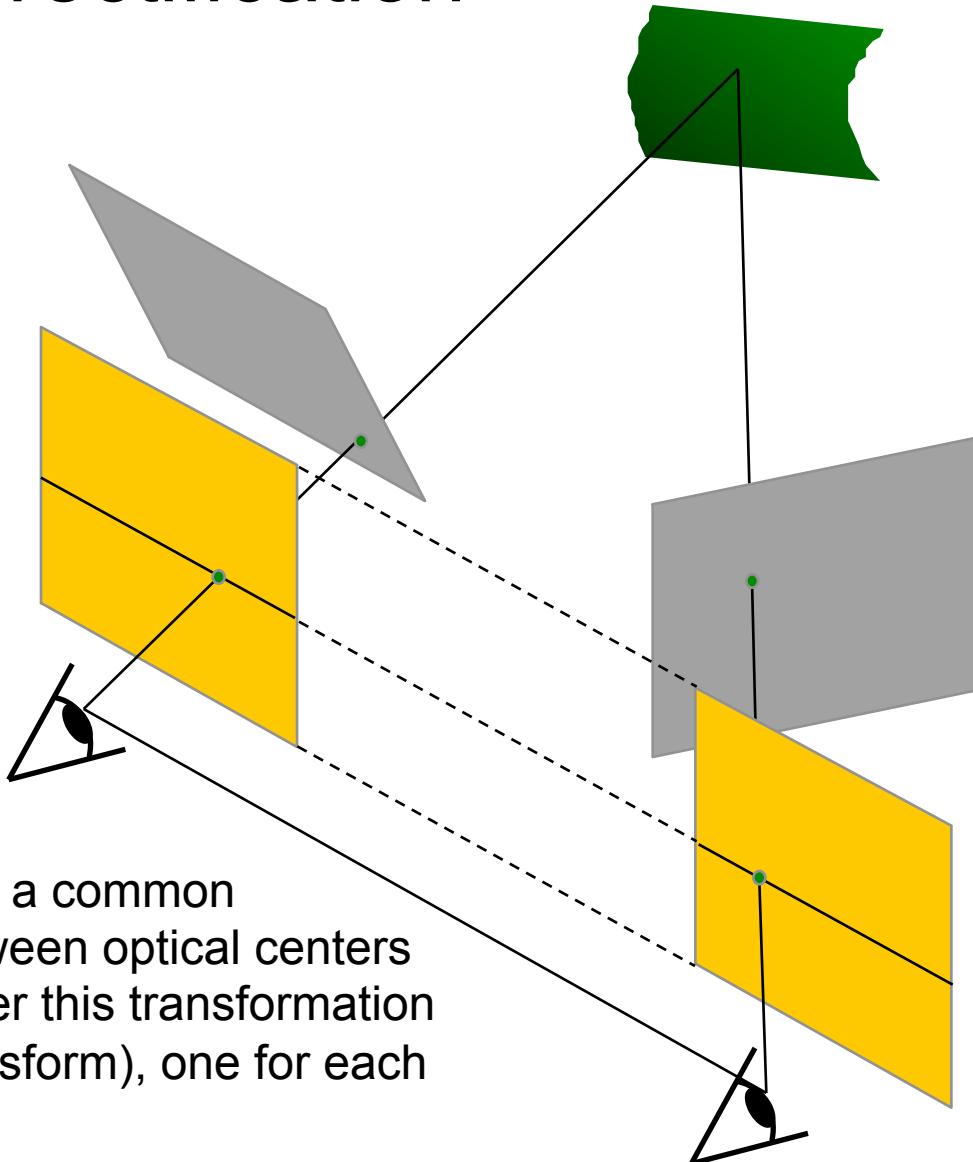


From [Wikipedia](#): Gemma Frisius's 1533 diagram introducing the idea of triangulation into the science of surveying. Having established a baseline, e.g. the cities of Brussels and Antwerp, the location of other cities, e.g. Middelburg, Ghent etc., can be found by taking a compass direction from each end of the baseline, and plotting where the two directions cross. This was only a theoretical presentation of the concept — due to topographical restrictions, it is impossible to see Middelburg from either Brussels or Antwerp. Nevertheless, the figure soon became well known all across Europe.

# Stereo image rectification



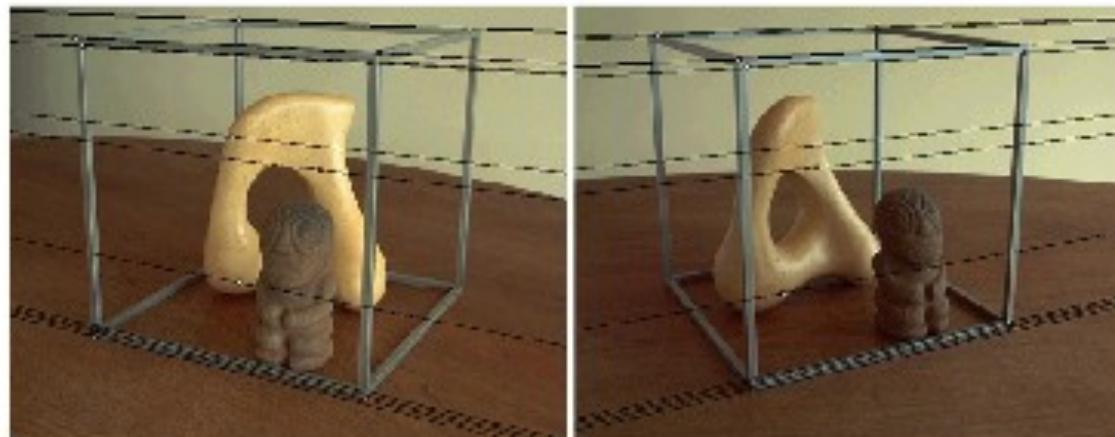
# Stereo image rectification



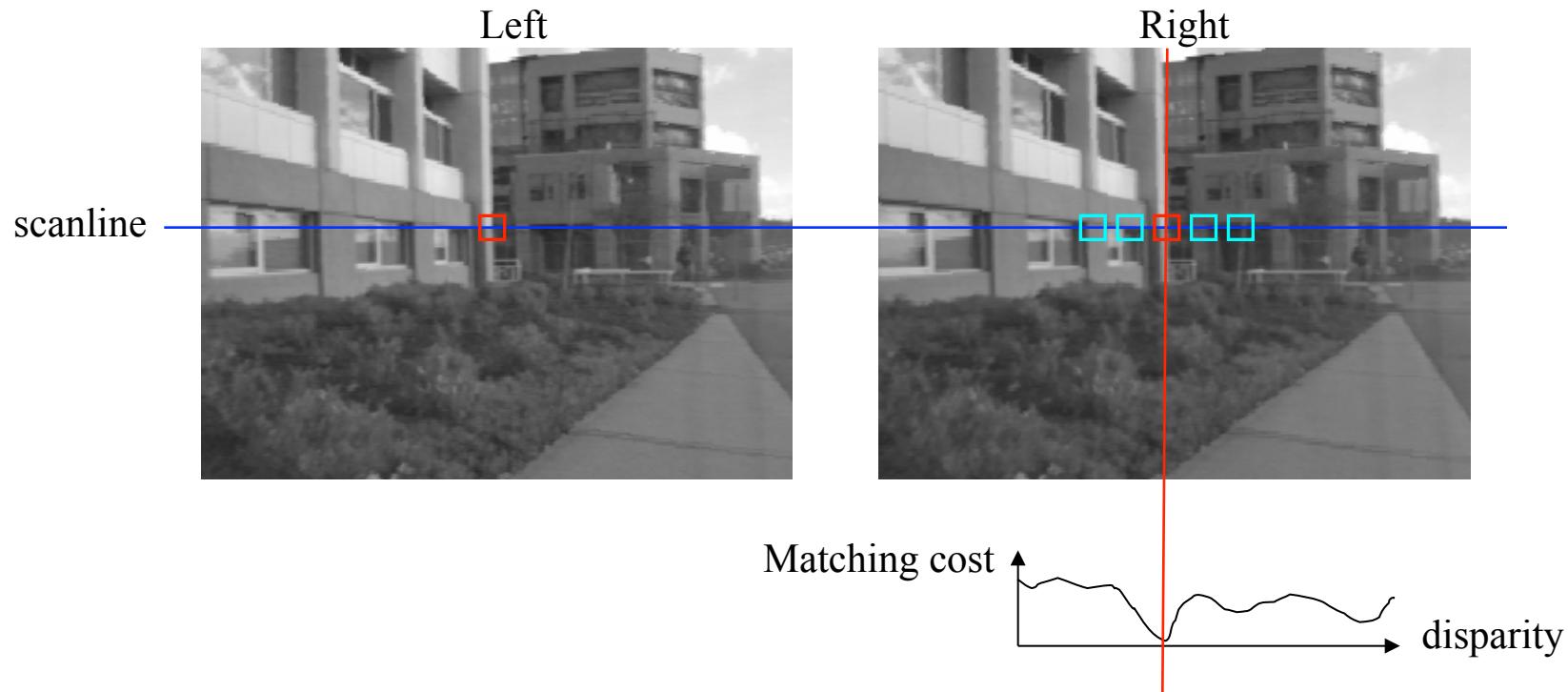
- Reproject image planes onto a common plane parallel to the line between optical centers
- Pixel motion is horizontal after this transformation
- Two homographies (3x3 transform), one for each input image reprojection

# Rectification example

---

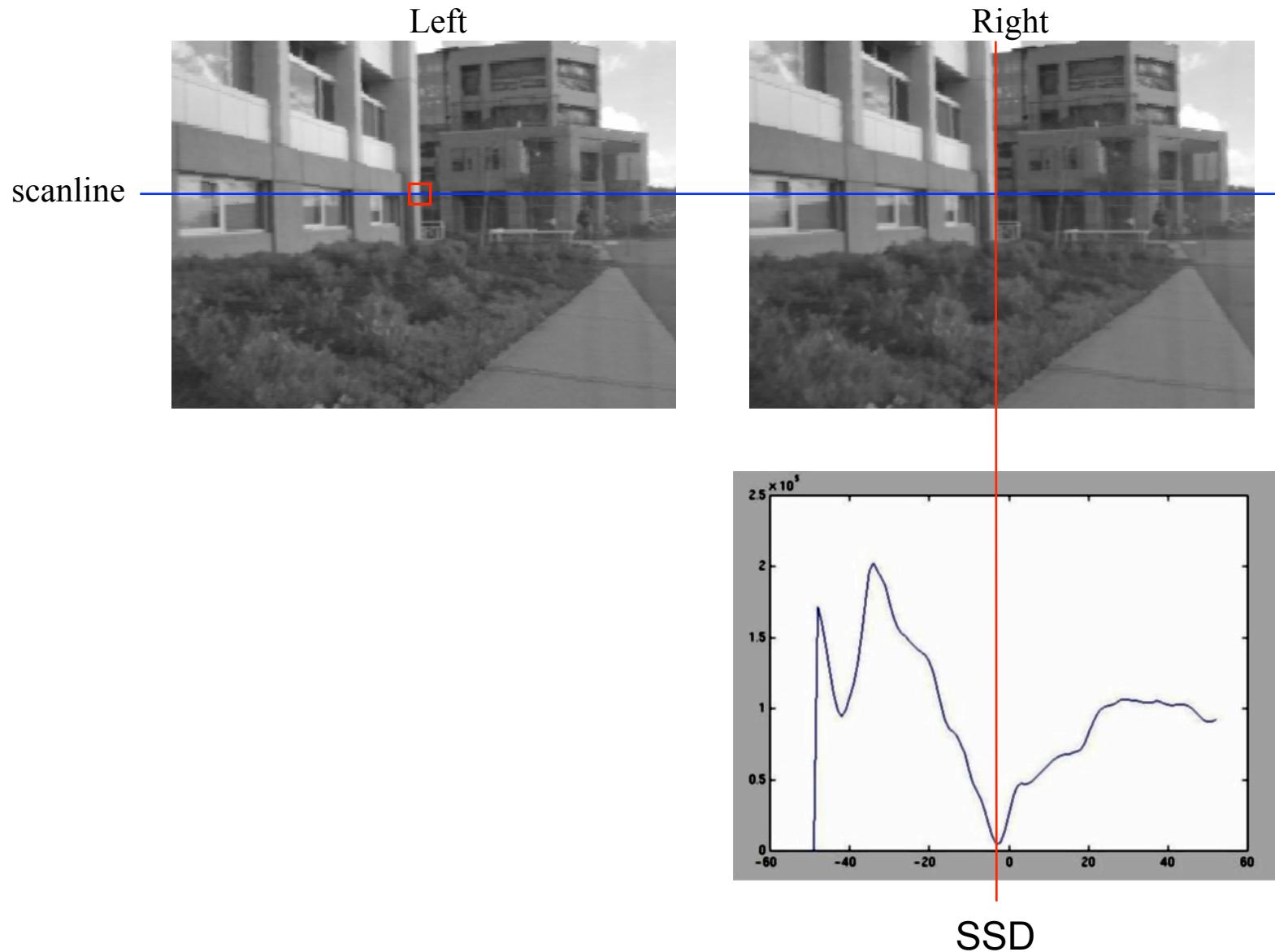


# Correspondence search

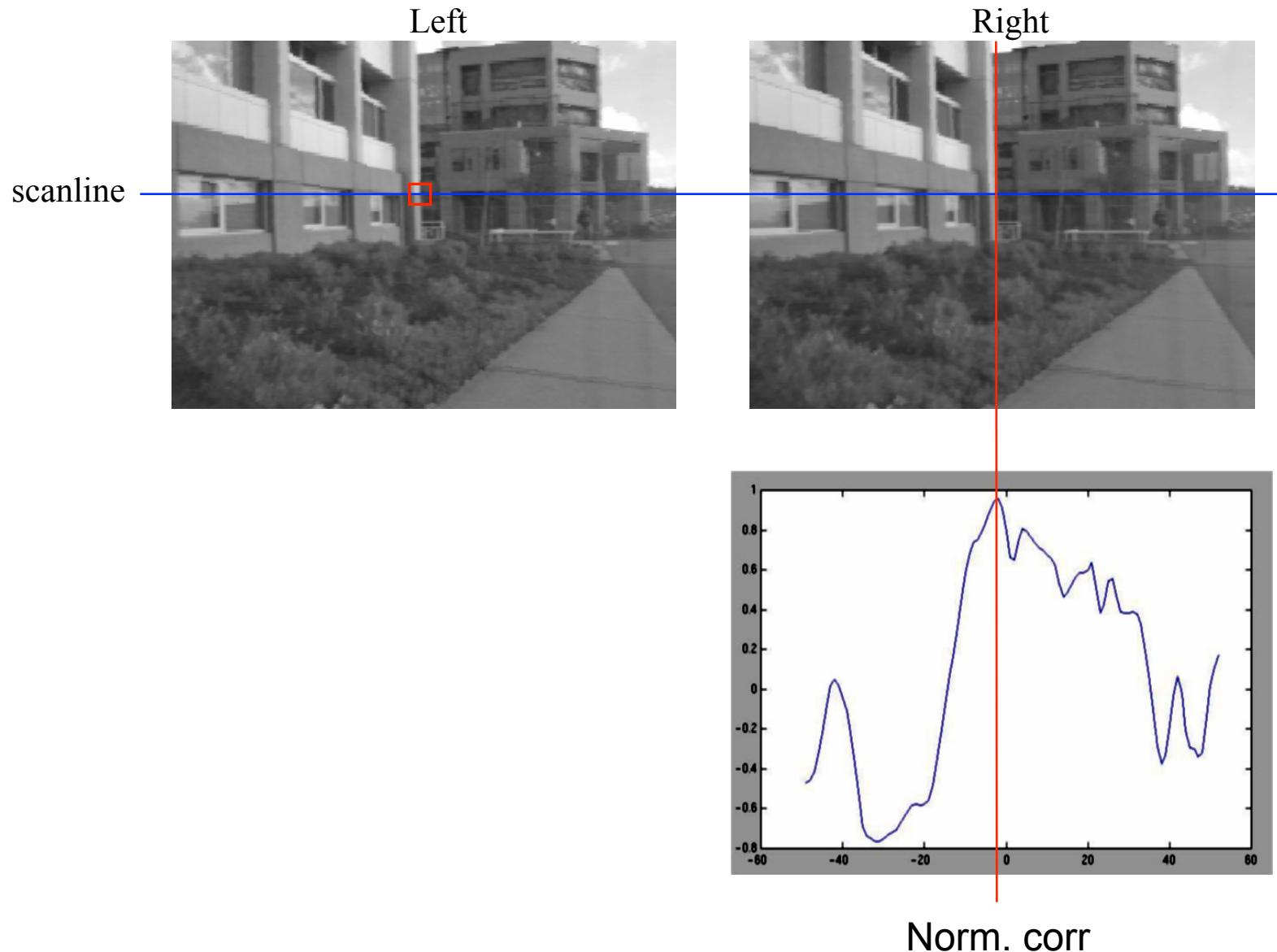


- Slide a window along the right scanline and compare contents of that window with the reference window in the left image
- Matching cost: SSD or normalized correlation

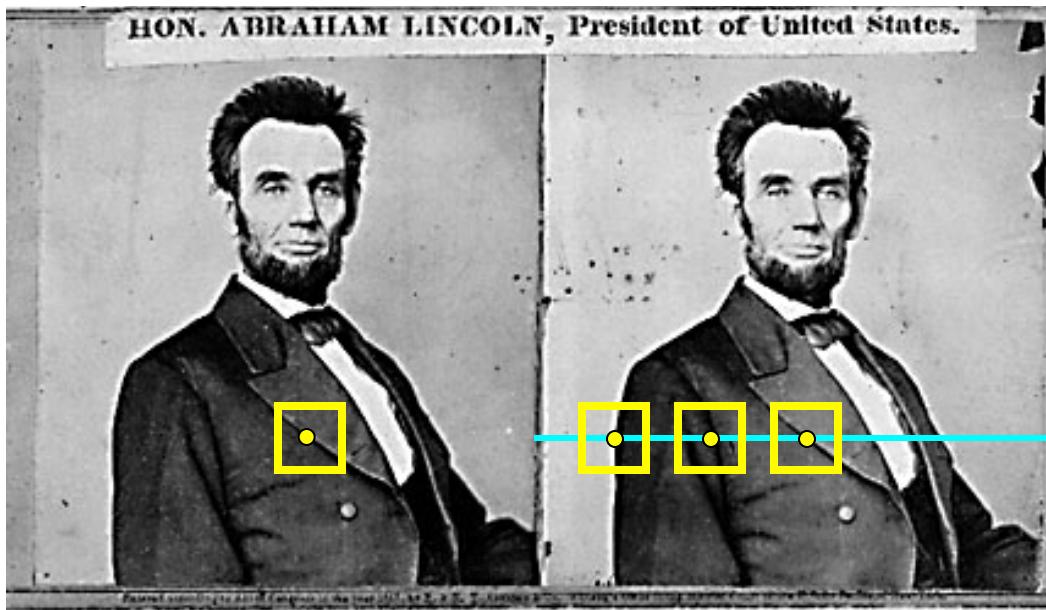
# Correspondence search



# Correspondence search

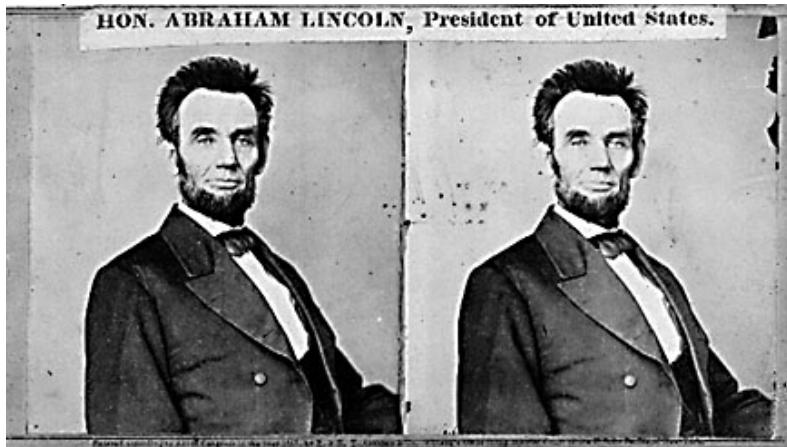


# Basic stereo matching algorithm



- If necessary, rectify the two stereo images to transform epipolar lines into scanlines
- For each pixel  $x$  in the first image
  - Find corresponding epipolar scanline in the right image
  - Examine all pixels on the scanline and pick the best match  $x'$
  - Compute disparity  $x-x'$  and set  $\text{depth}(x) = B*f/(x-x')$

# Failures of correspondence search



Textureless surfaces



Occlusions, repetition



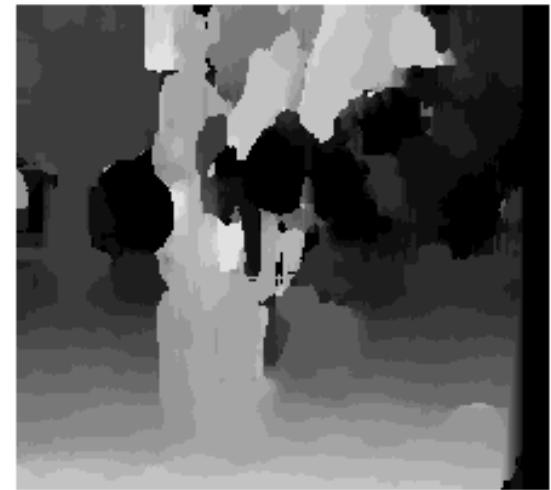
Non-Lambertian surfaces, specularities

# Effect of window size

---



$$W = 3$$



$$W = 20$$

- Smaller window
  - + More detail
  - More noise
  
- Larger window
  - + Smoother disparity maps
  - Less detail

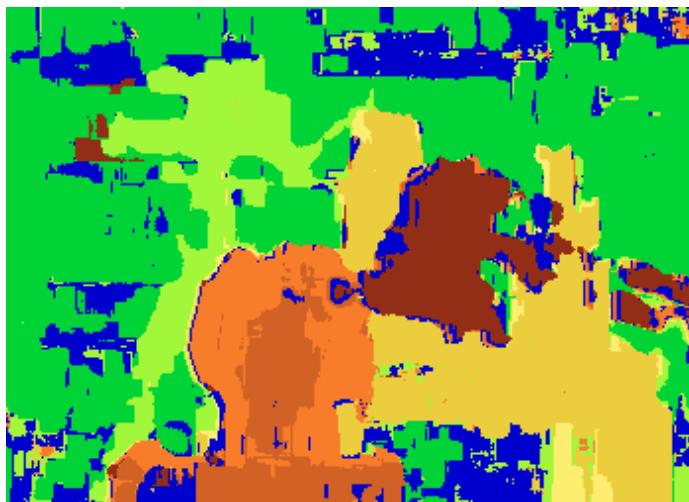
# Results with window search

---

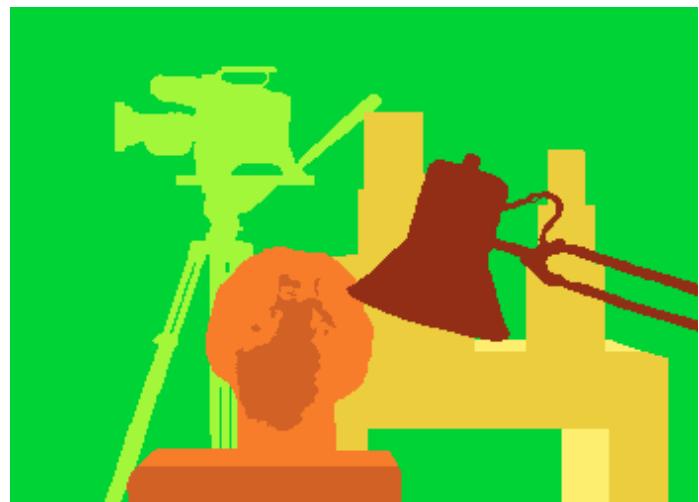
Data



Window-based matching



Ground truth



# Better methods exist...

---



Graph cuts



Ground truth

Y. Boykov, O. Veksler, and R. Zabih,

Fast Approximate Energy Minimization via Graph Cuts, PAMI 2001

For the latest and greatest: <http://www.middlebury.edu/stereo/>

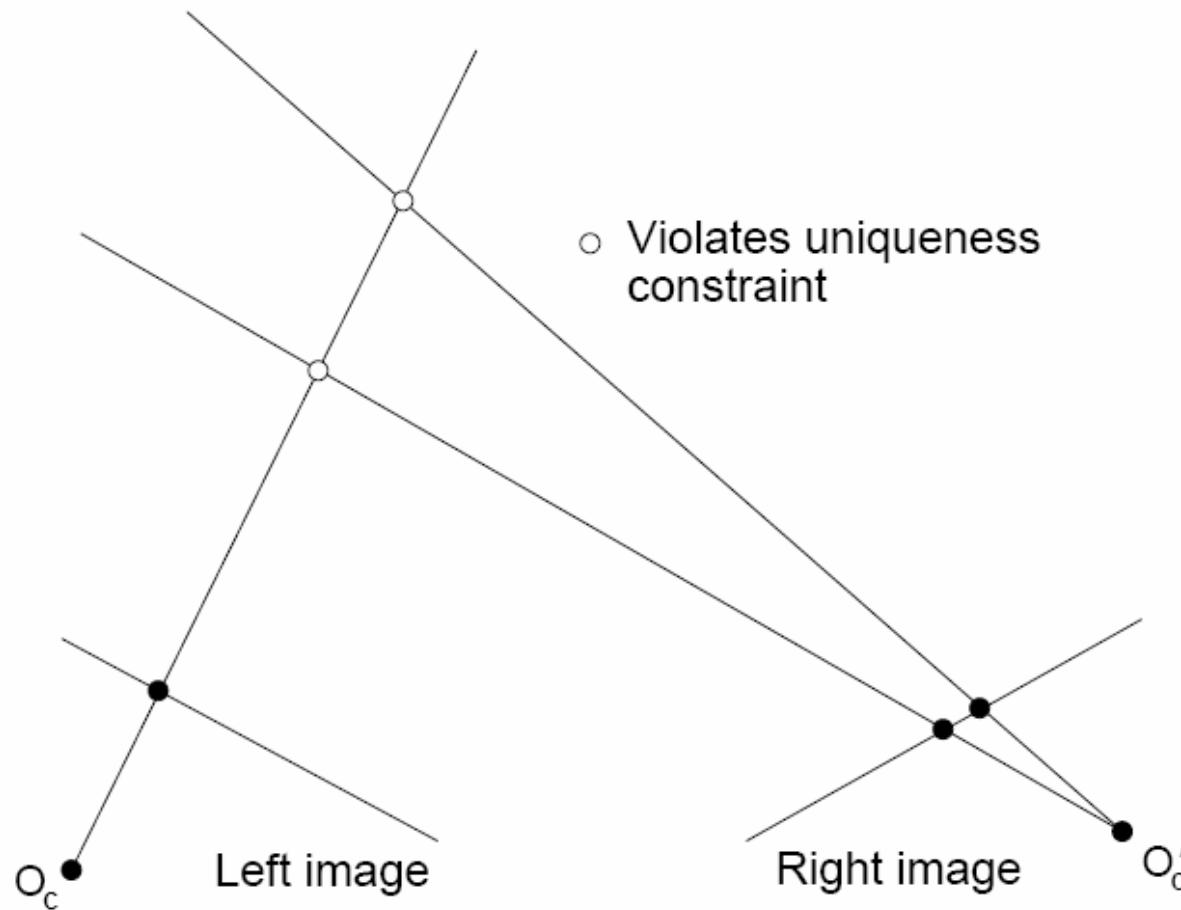
# How can we improve window-based matching?

---

- The similarity constraint is **local** (each reference window is matched independently)
- Need to enforce **non-local** correspondence constraints

# Non-local constraints

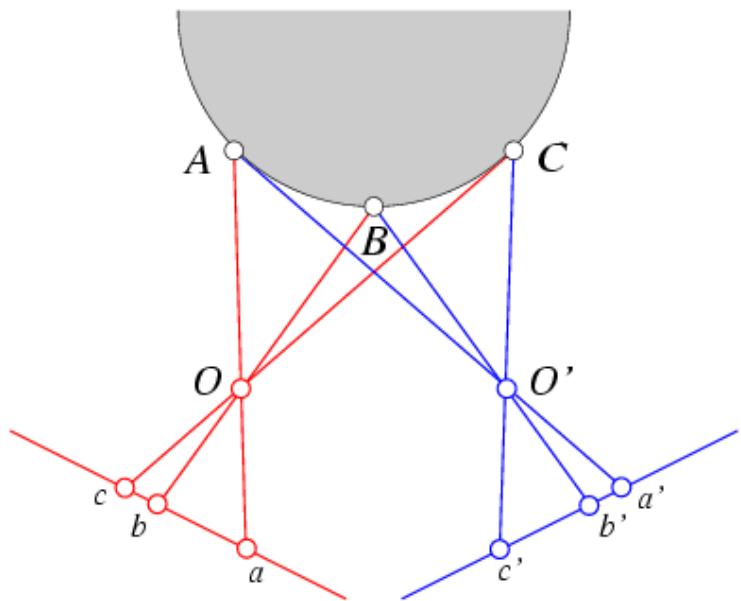
- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image



# Non-local constraints

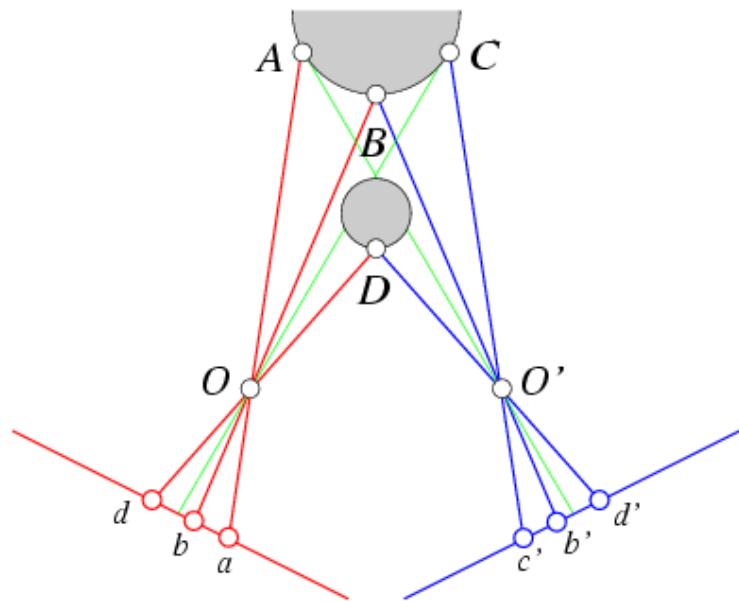
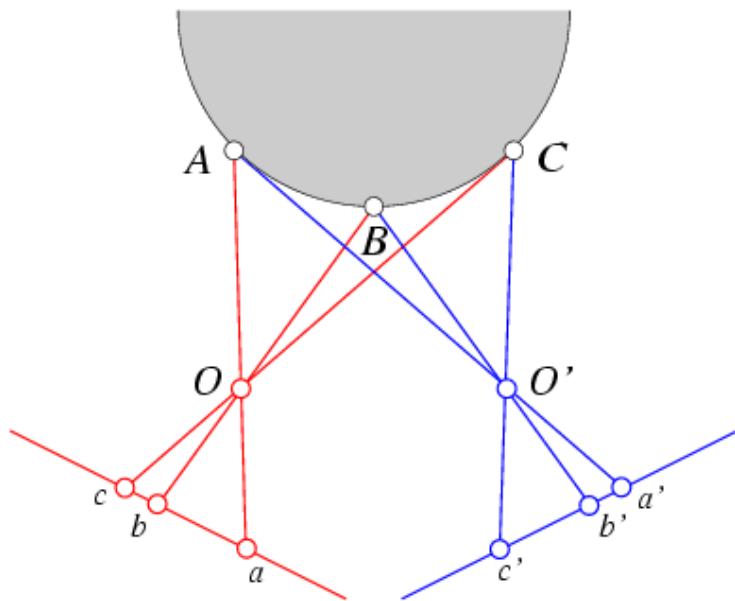
---

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering
  - Corresponding points should be in the same order in both views



# Non-local constraints

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering
  - Corresponding points should be in the same order in both views



Ordering constraint doesn't hold

# Non-local constraints

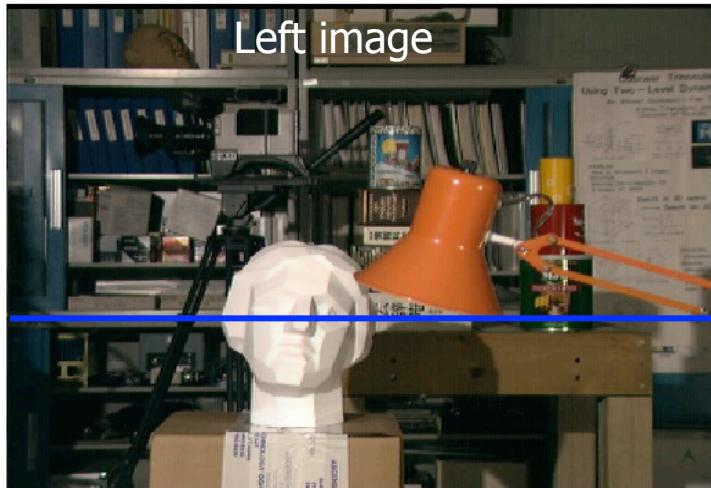
---

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering
  - Corresponding points should be in the same order in both views
- Smoothness
  - We expect disparity values to change slowly (for the most part)

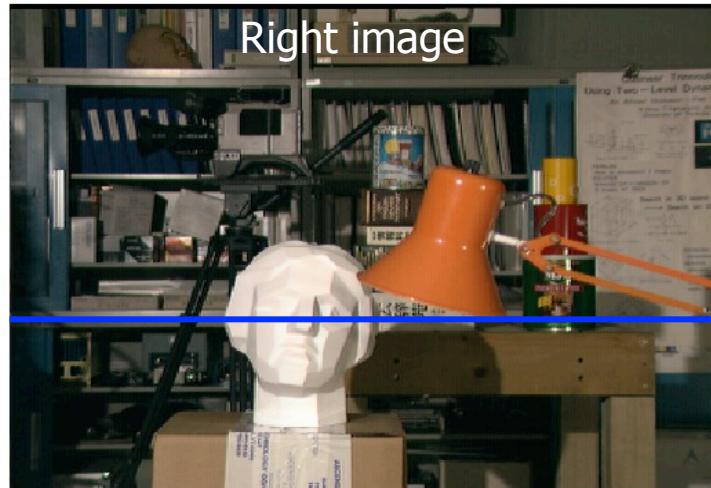
# Scanline stereo

---

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

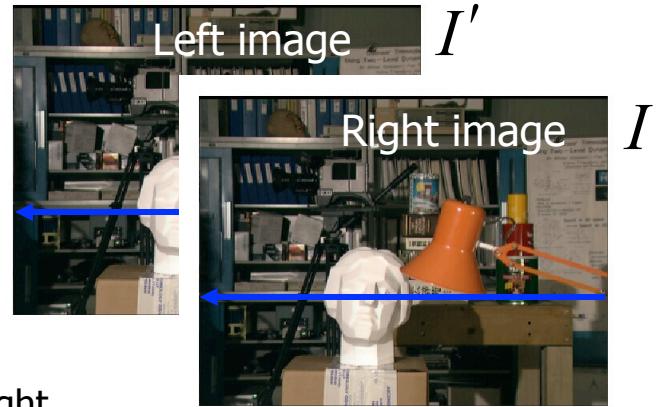
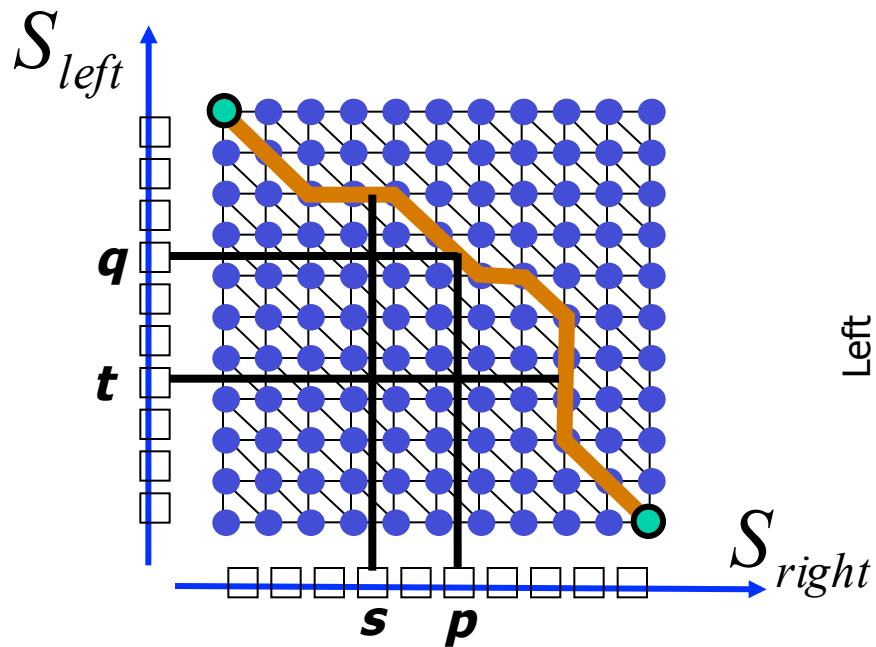


Left image



Right image

# “Shortest paths” for scan-line stereo

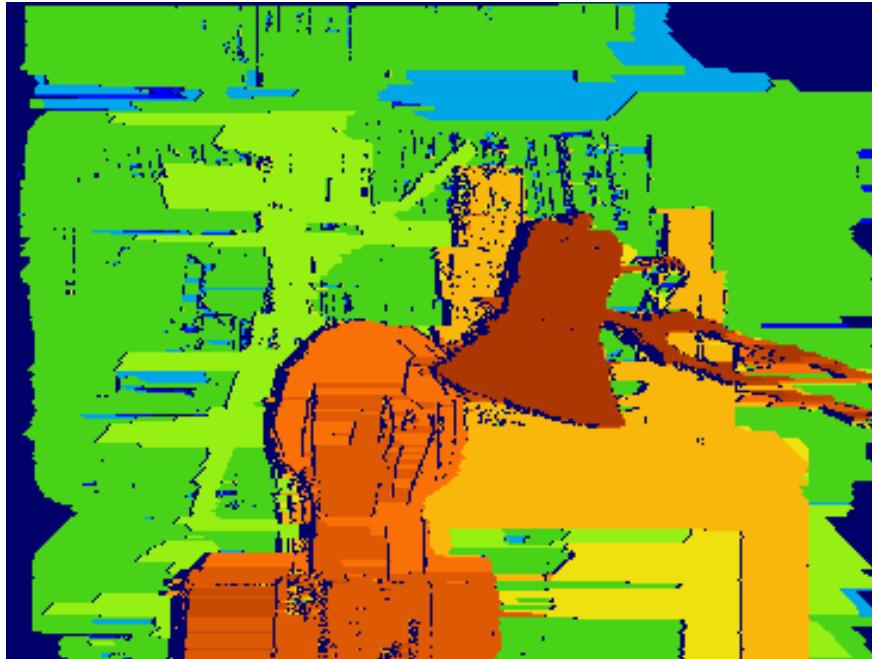


Can be implemented with dynamic programming  
Ohta & Kanade '85, Cox et al. '96

# Coherent stereo on 2D grid

---

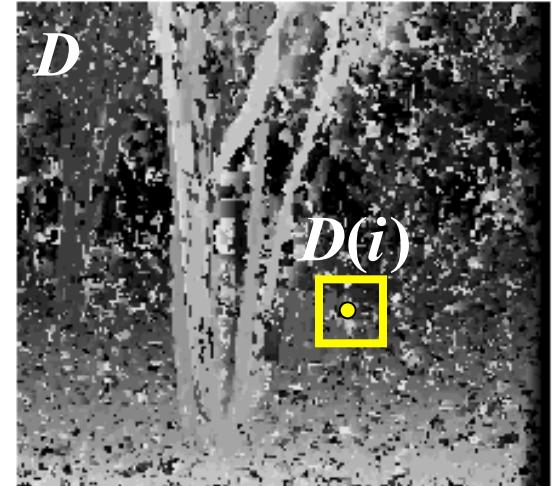
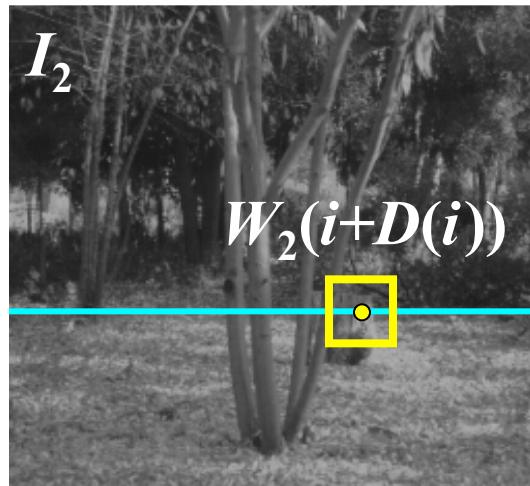
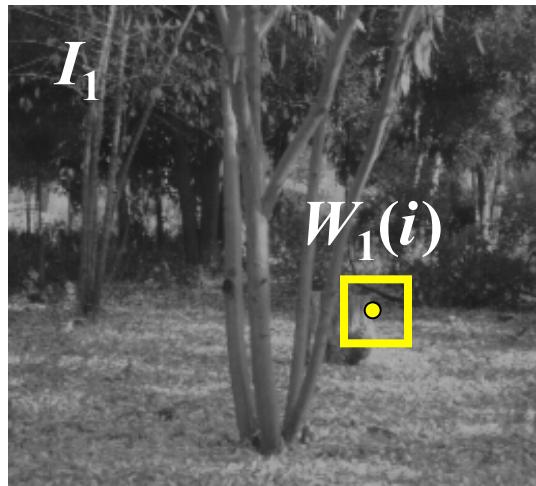
- Scanline stereo generates streaking artifacts



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

# Stereo matching as energy minimization

---



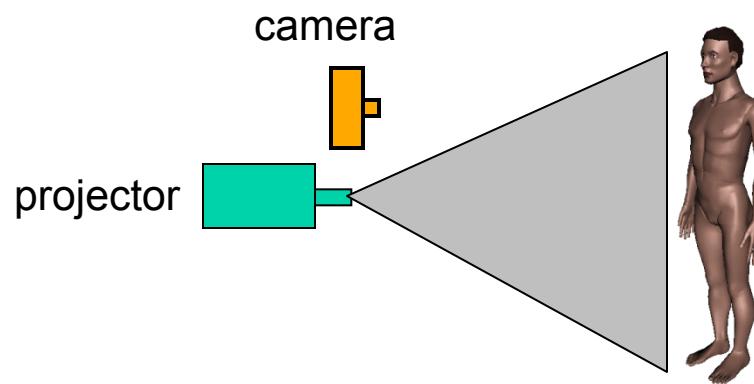
$$E(D) = \underbrace{\sum_i (W_1(i) - W_2(i + D(i)))^2}_{\text{data term}} + \lambda \underbrace{\sum_{\text{neighbors } i,j} \rho(D(i) - D(j))}_{\text{smoothness term}}$$

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

# Active stereo with structured light



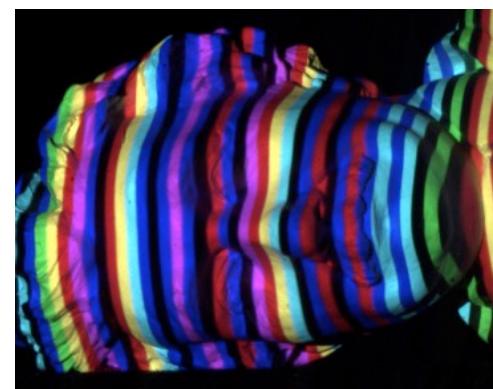
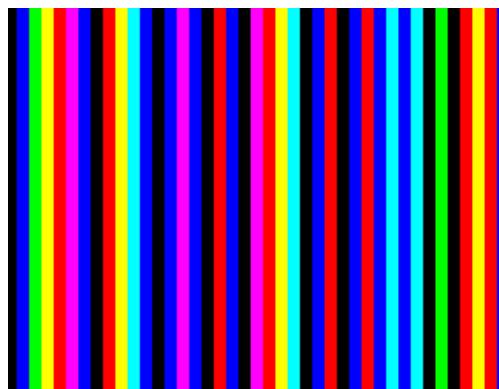
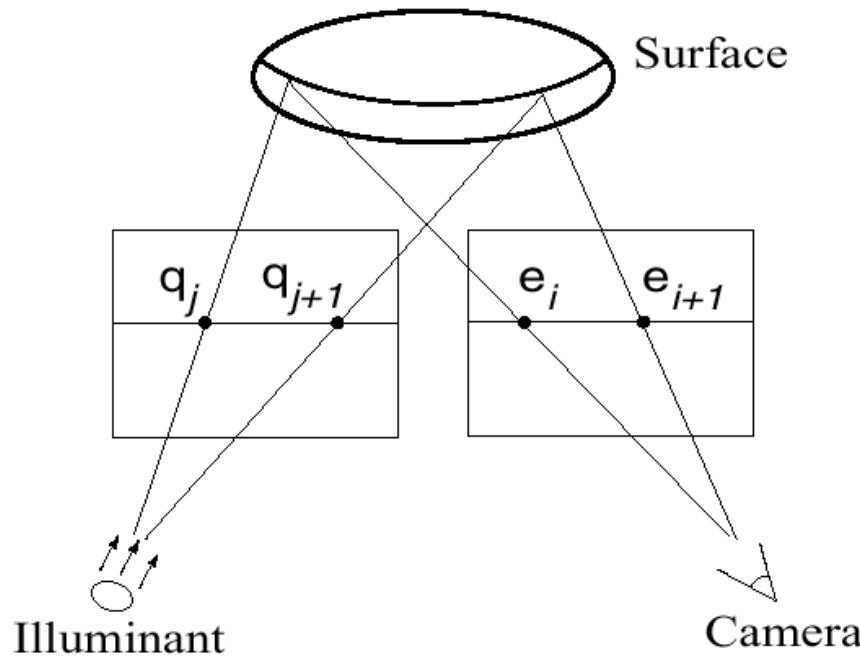
- Project “structured” light patterns onto the object
  - Simplifies the correspondence problem
  - Allows us to use only one camera



L. Zhang, B. Curless, and S. M. Seitz.

[Rapid Shape Acquisition Using Color Structured Light and Multi-pass Dynamic Programming. 3DPVT 2002](#)

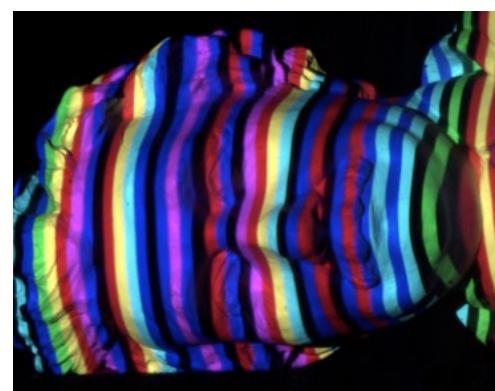
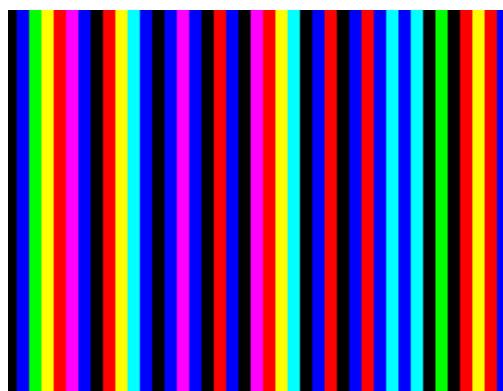
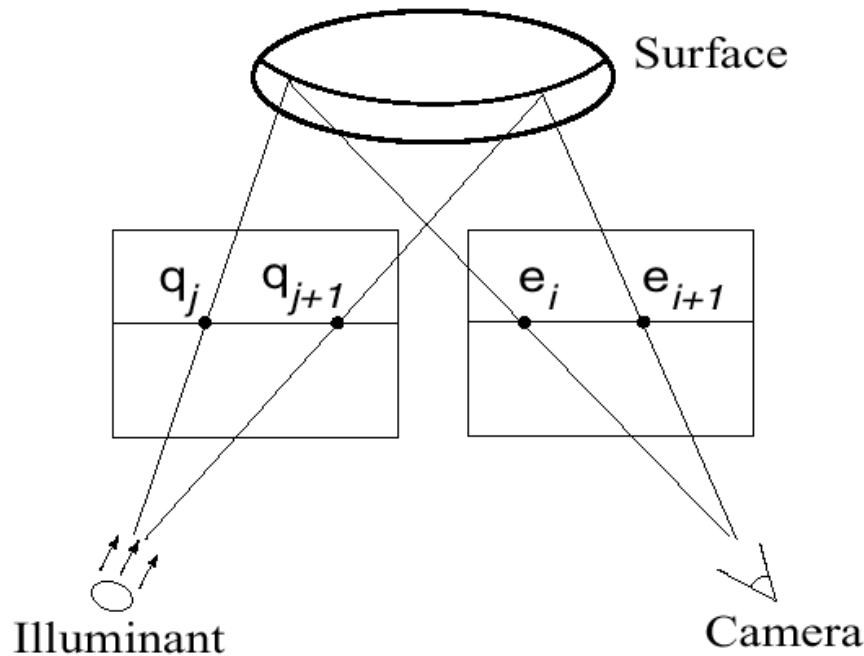
# Active stereo with structured light



L. Zhang, B. Curless, and S. M. Seitz.

[Rapid Shape Acquisition Using Color Structured Light and Multi-pass Dynamic Programming.](#) 3DPVT 2002

# Active stereo with structured light



[http://en.wikipedia.org/wiki/Structured-light\\_3D\\_scanner](http://en.wikipedia.org/wiki/Structured-light_3D_scanner)

# Kinect: Structured infrared light

---



<http://bbzippo.wordpress.com/2010/11/28/kinect-in-infrared/>