

Stereo

CSE 576

Ali Farhadi

Several slides from Larry Zitnick and Steve Seitz

**WHY DO WE HAVE
TWO EYES?**



Why do we perceive depth?



What do humans use as depth cues?

Motion

Convergence

When watching an object close to us, our eyes point slightly inward. This difference in the direction of the eyes is called convergence. This depth cue is effective only on short distances (less than 10 meters).

Binocular Parallax

As our eyes see the world from slightly different locations, the images sensed by the eyes are slightly different. This difference in the sensed images is called binocular parallax. Human visual system is very sensitive to these differences, and binocular parallax is the most important depth cue for medium viewing distances. The sense of depth can be achieved using binocular parallax even if all other depth cues are removed.

Monocular Movement Parallax

If we close one of our eyes, we can perceive depth by moving our head. This happens because human visual system can extract depth information in two similar images sensed after each other, in the same way it can combine two images from different eyes.

Focus

Accommodation

Accommodation is the tension of the muscle that changes the focal length of the lens of eye. Thus it brings into focus objects at different distances. This depth cue is quite weak, and it is effective only at short viewing distances (less than 2 meters) and with other cues.

What do humans use as depth cues?

Image cues

Retinal Image Size

When the real size of the object is known, our brain compares the sensed size of the object to this real size, and thus acquires information about the distance of the object.

Linear Perspective

When looking down a straight level road we see the parallel sides of the road meet in the horizon. This effect is often visible in photos and it is an important depth cue. It is called linear perspective.

Texture Gradient

The closer we are to an object the more detail we can see of its surface texture. So objects with smooth textures are usually interpreted being farther away. This is especially true if the surface texture spans all the distance from near to far.

Overlapping

When objects block each other out of our sight, we know that the object that blocks the other one is closer to us. The object whose outline pattern looks more continuous is felt to lie closer.

Aerial Haze

The mountains in the horizon look always slightly bluish or hazy. The reason for this are small water and dust particles in the air between the eye and the mountains. The farther the mountains, the hazier they look.

Shades and Shadows

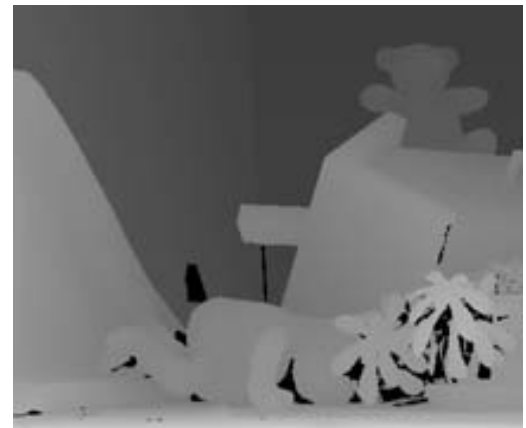
When we know the location of a light source and see objects casting shadows on other objects, we learn that the object shadowing the other is closer to the light source. As most illumination comes downward we tend to resolve ambiguities using this information. The three dimensional looking computer user interfaces are a nice example on this. Also, bright objects seem to be closer to the observer than dark ones.



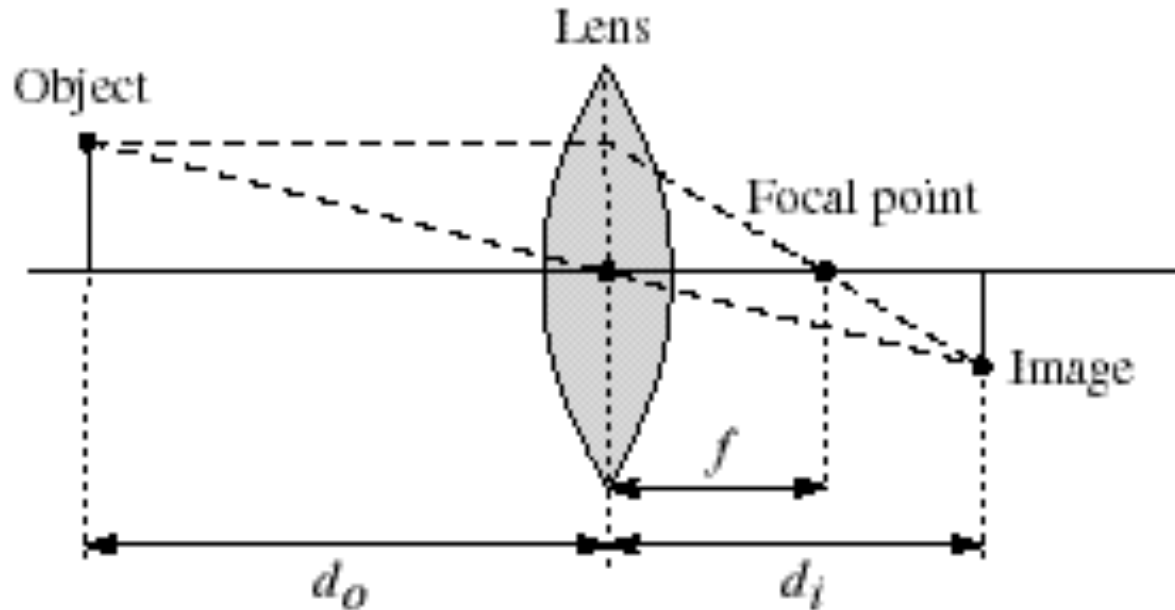


Amount of horizontal movement is ...

...inversely proportional to the distance from the camera



Cameras

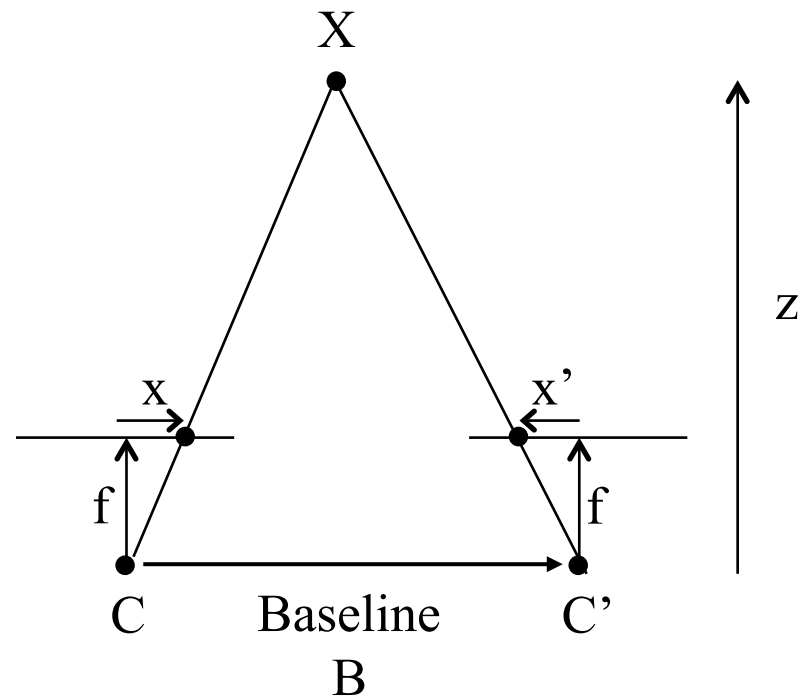
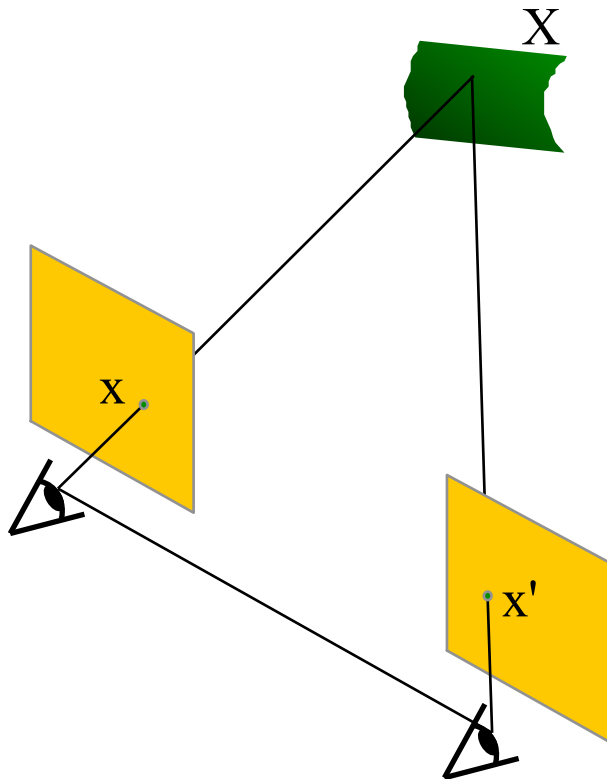


Thin lens equation:
$$\frac{1}{d_o} + \frac{1}{d_i} = \frac{1}{f}$$

- Any object point satisfying this equation is in focus

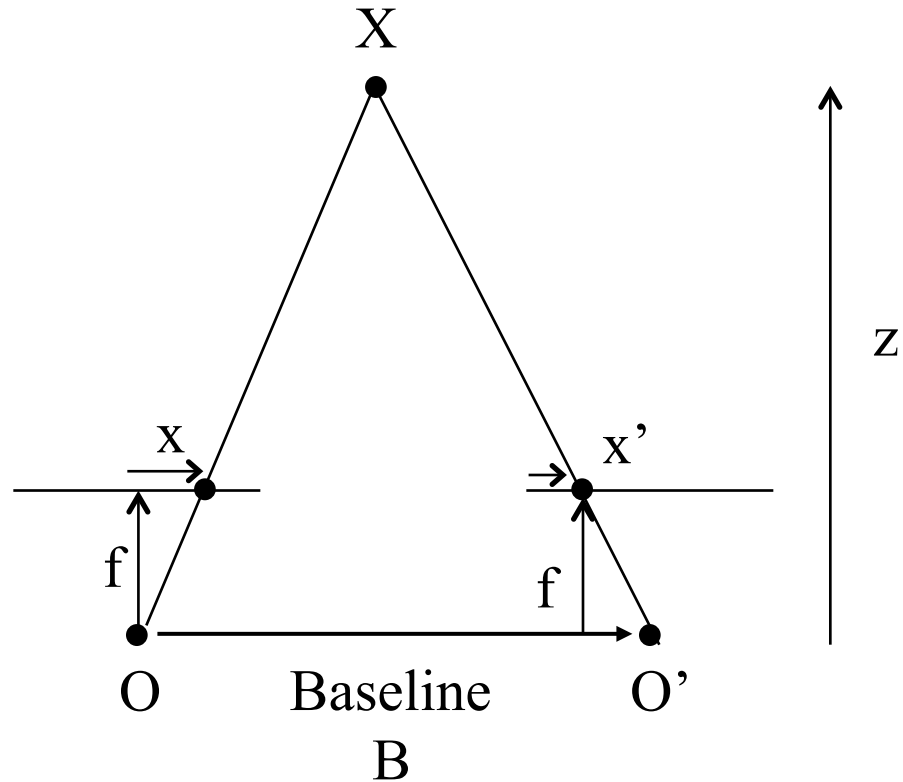
Depth from Stereo

Goal: recover depth by finding image coordinate x' that corresponds to x



Depth from disparity

$$\frac{x - x'}{O - O'} = \frac{f}{z}$$



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

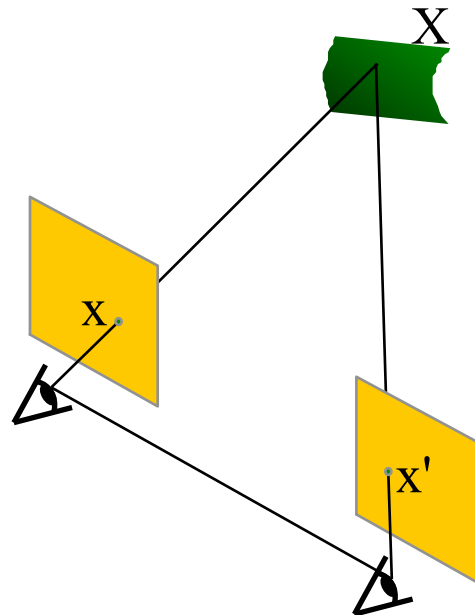
Disparity is inversely proportional to depth.

Depth from Stereo

Goal: recover depth by finding image coordinate x' that corresponds to x

Sub-Problems

1. Calibration: How do we recover the relation of the cameras (if not already known)?
2. Correspondence: How do we search for the matching point x' ?



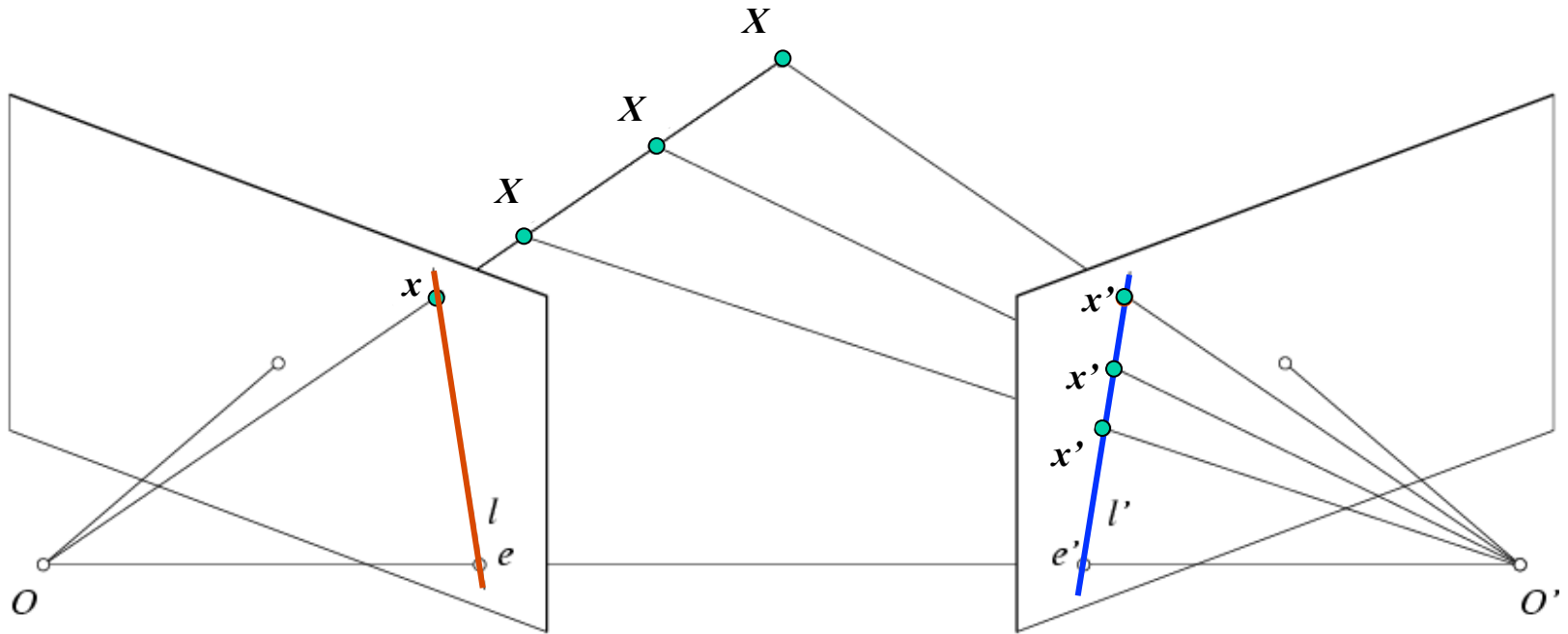
Correspondence Problem



We have two images taken from cameras with different intrinsic and extrinsic parameters

How do we match a point in the first image to a point in the second? How can we constrain our search?

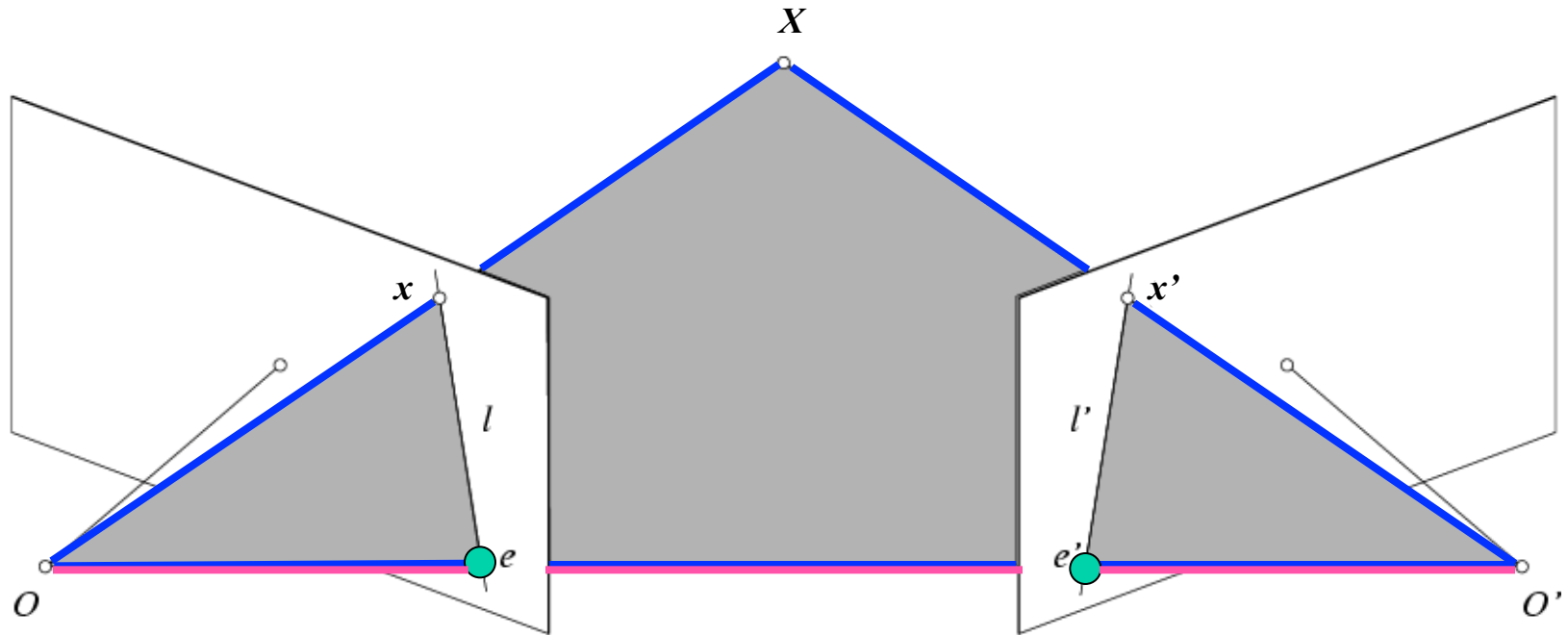
Key idea: Epipolar constraint



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

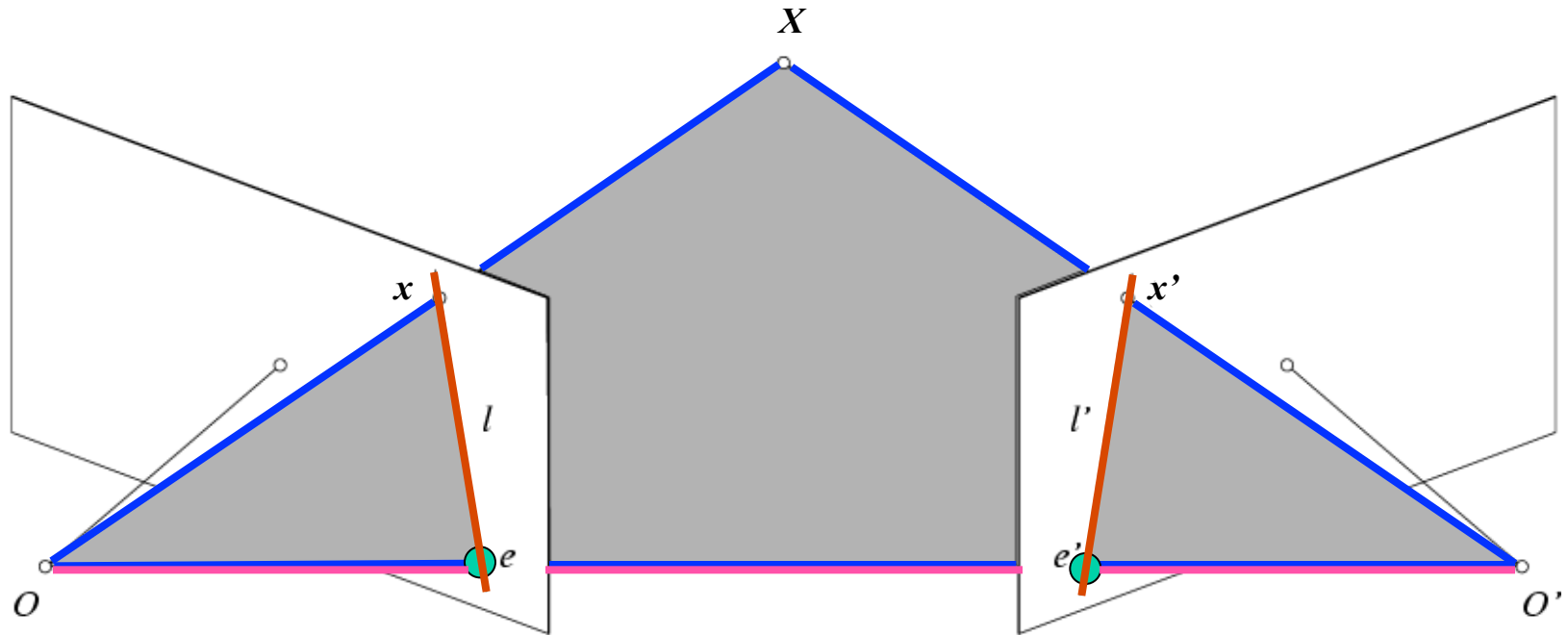
Potential matches for x' have to lie on the corresponding line l .

Epipolar geometry: notation



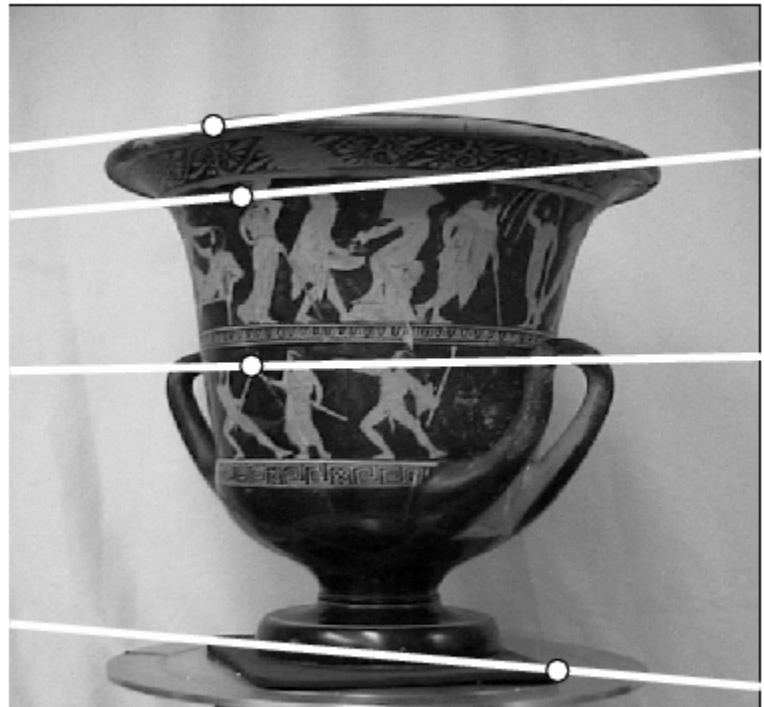
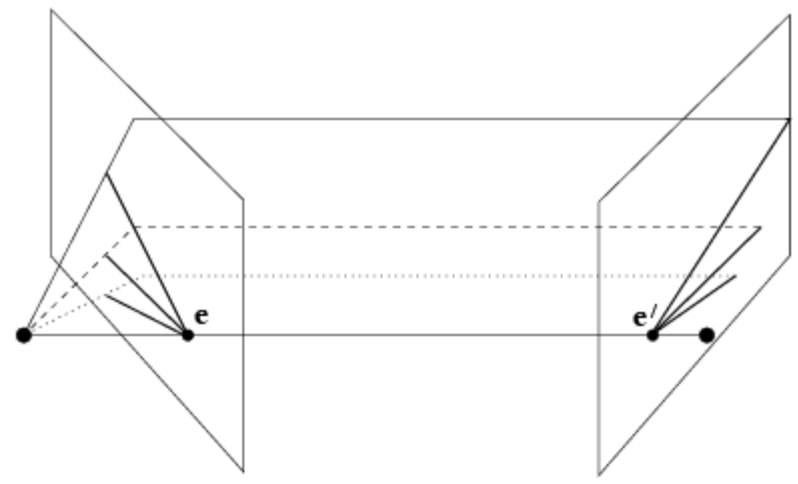
- **Baseline** – line connecting the two camera centers
- **Epipoles**
= intersections of baseline with image planes
= projections of the other camera center
- **Epipolar Plane** – plane containing baseline (1D family)

Epipolar geometry: notation

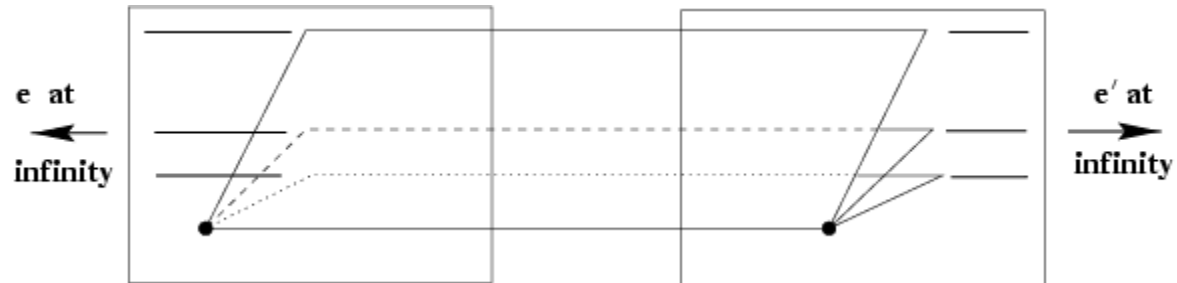


- **Baseline** – line connecting the two camera centers
- **Epipoles**
= intersections of baseline with image planes
= projections of the other camera center
- **Epipolar Plane** – plane containing baseline (1D family)
- **Epipolar Lines** - intersections of epipolar plane with image planes (always come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

Example: Converging cameras



Example: Motion parallel to image plane



Example: Forward motion

What would the epipolar lines look like if the camera moves directly forward?

Example: Motion perpendicular to image plane

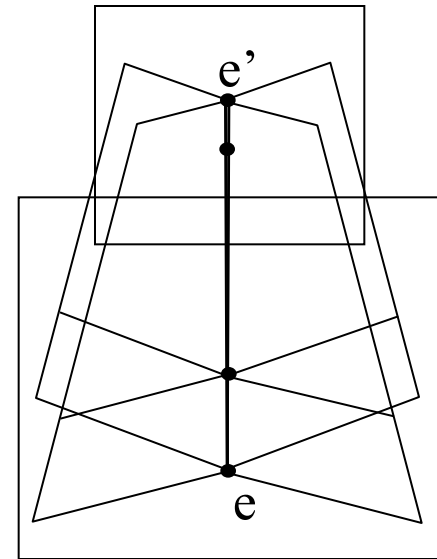
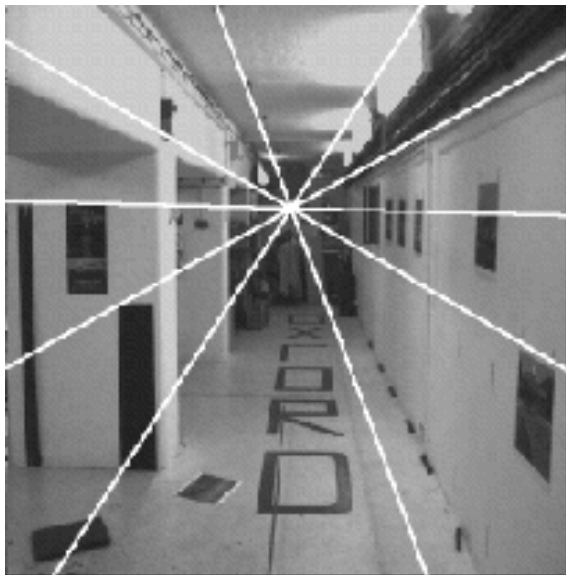
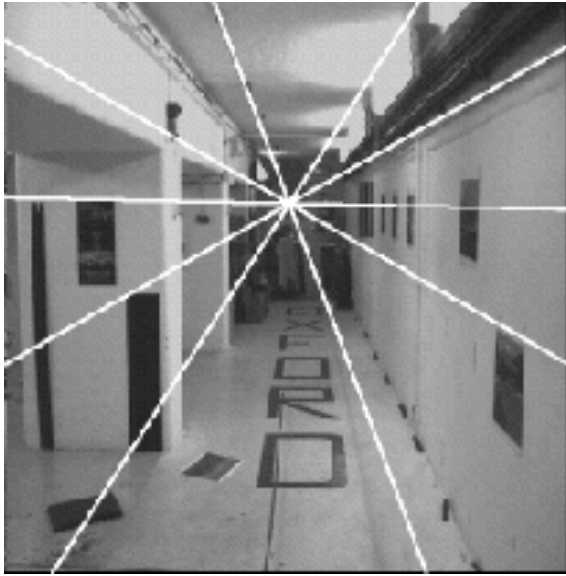


Example: Motion perpendicular to image plane



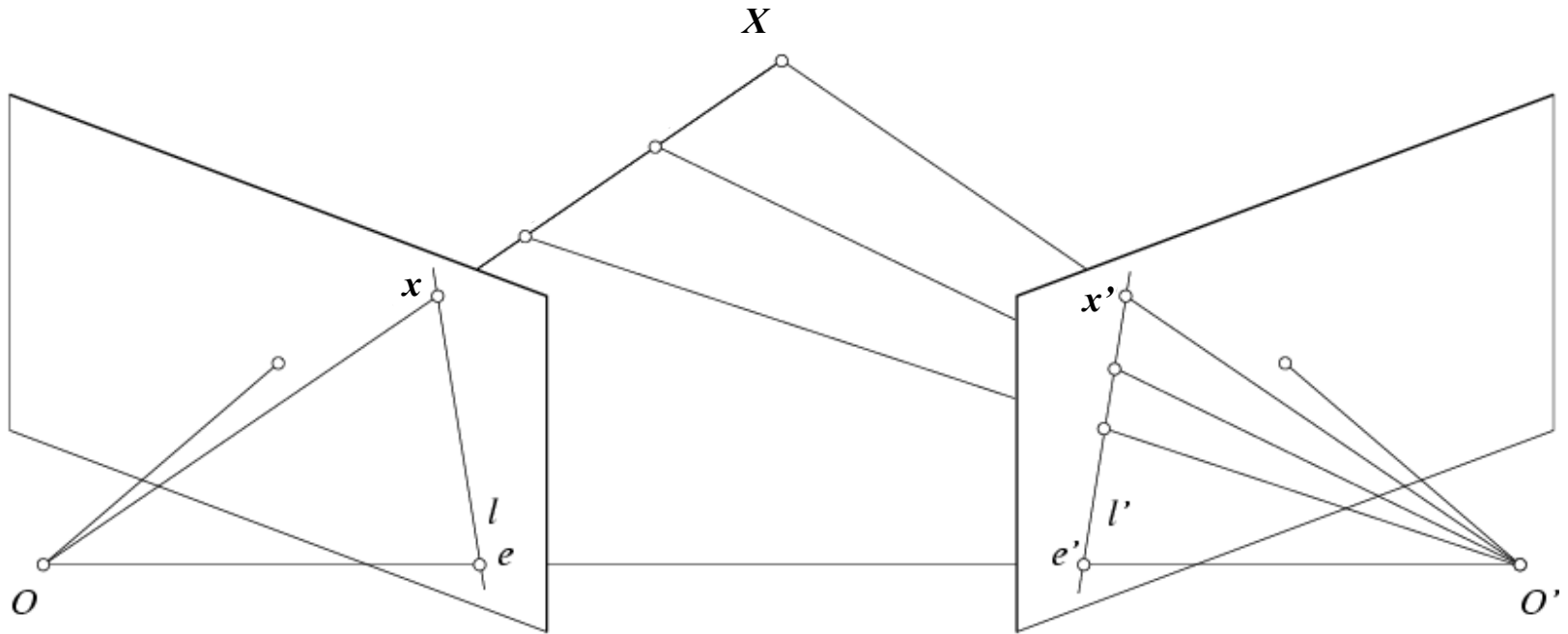
- Points move along lines radiating from the epipole: “focus of expansion”
- Epipole is the principal point

Example: Forward motion



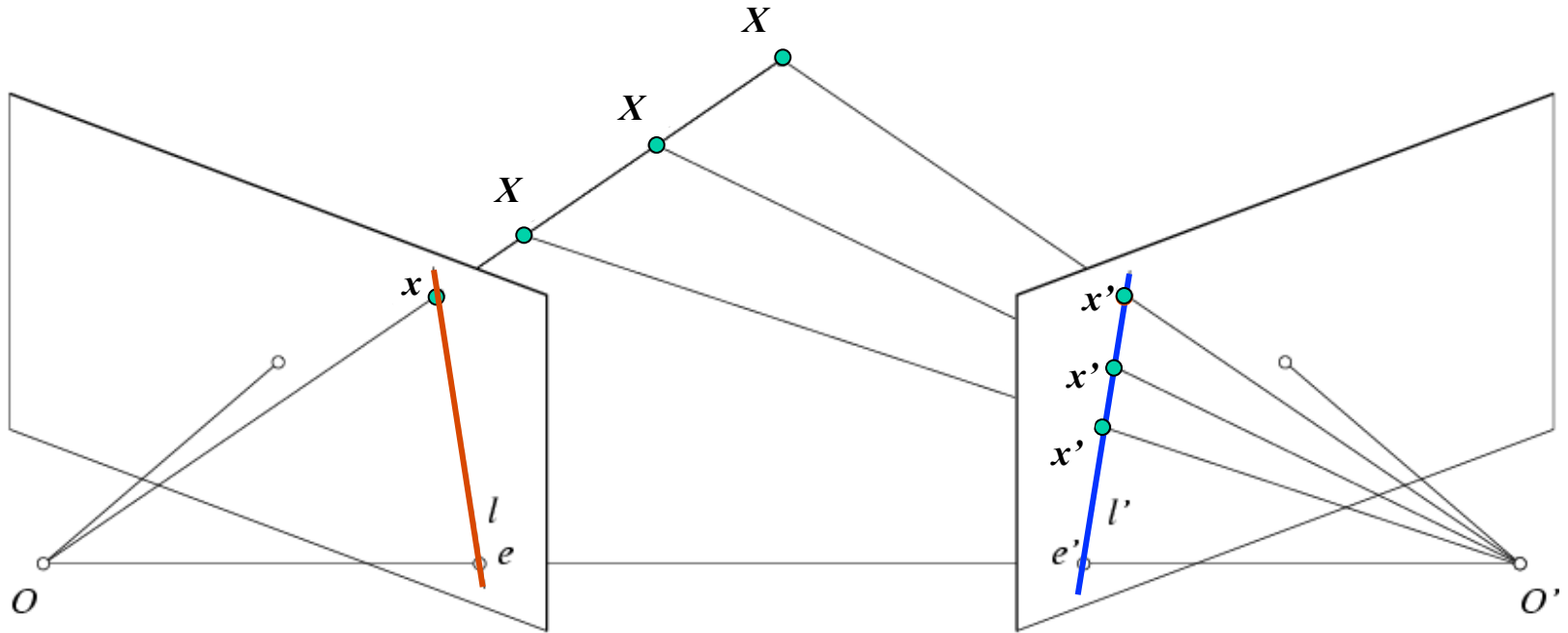
Epipole has same coordinates in both images.
Points move along lines radiating from “Focus of expansion”

Epipolar constraint



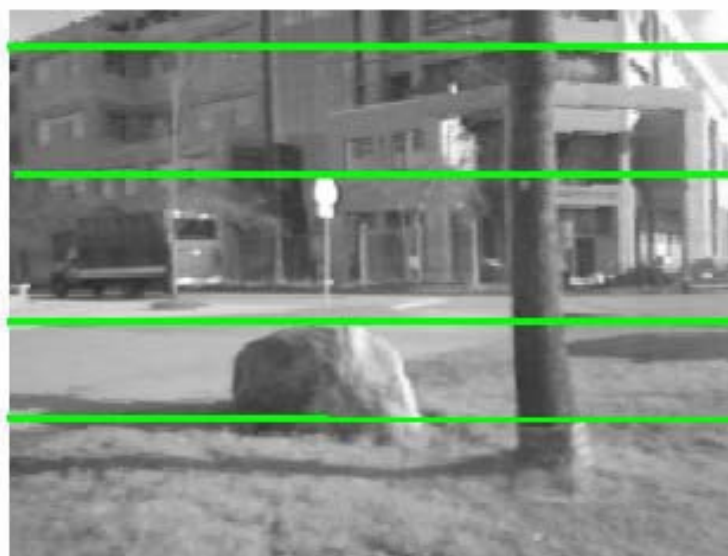
- If we observe a point \mathbf{x} in one image, where can the corresponding point \mathbf{x}' be in the other image?

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

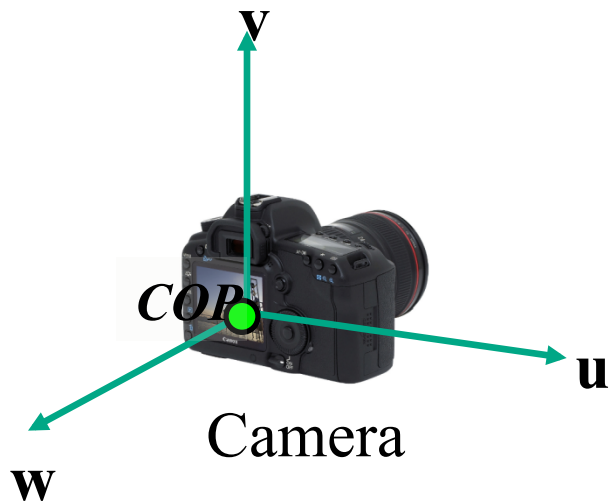


Camera parameters

How many numbers do we need to describe a camera?

- We need to describe its *pose* in the world
- We need to describe its internal *parameters*

A Tale of Two Coordinate Systems



Two important coordinate systems:

1. *World* coordinate system
2. *Camera* coordinate system



“The World”

Camera parameters

- To project a point (x,y,z) in *world* coordinates into a camera
- First transform (x,y,z) into *camera* coordinates
- Need to know
 - Camera position (in world coordinates)
 - Camera orientation (in world coordinates)
- Then project into the image plane
 - Need to know camera *intrinsics*
- These can all be described with matrices

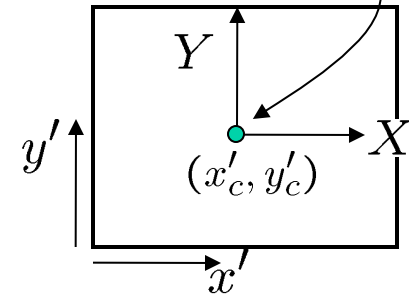
Camera parameters

A camera is described by several parameters

- Translation **T** of the optical center from the origin of world coords
- Rotation **R** of the image plane
- focal length **f**, principle point (x'_c, y'_c) , pixel size (s_x, s_y)
- blue parameters are called “extrinsics,” red are “intrinsics”

Projection equation

$$\mathbf{x} = \begin{bmatrix} wx \\ wy \\ w \end{bmatrix} = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = \mathbf{\Pi} \mathbf{X}$$



- The projection matrix models the cumulative effect of all parameters
- Useful to decompose into a series of operations

$$\mathbf{\Pi} = \begin{bmatrix} -fs_x & 0 & x'_c \\ 0 & -fs_y & y'_c \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{R}_{3 \times 3} & \mathbf{0}_{3 \times 1} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{I}_{3 \times 3} & \mathbf{T}_{3 \times 1} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix}$$

intrinsics projection rotation translation

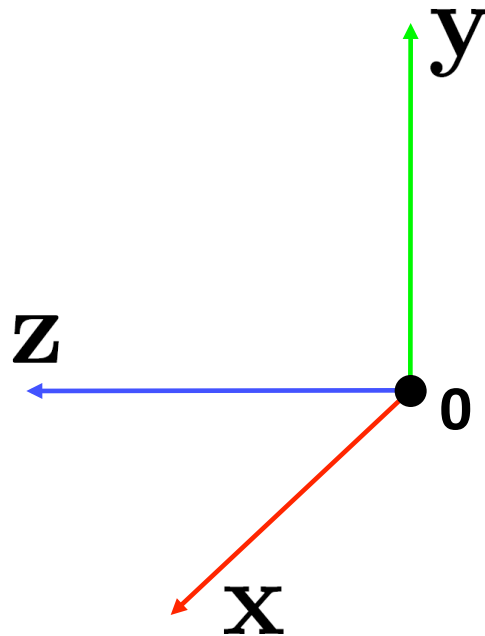
identity matrix

- The definitions of these parameters are not completely standardized
 - especially intrinsics—varies from one book to another

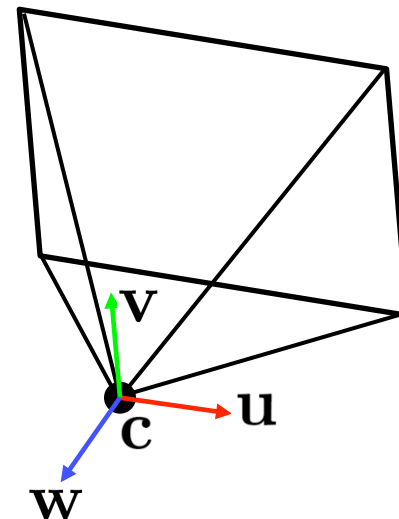
Extrinsics

How do we get the camera to “canonical form”?

- (Center of projection at the origin, x-axis points right, y-axis points up, z-axis points backwards)



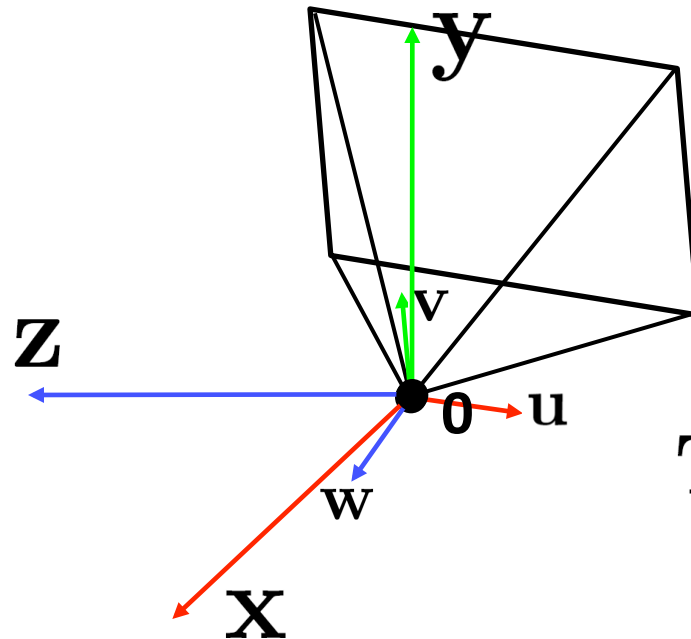
Step 1: Translate by $-c$



Extrinsics

How do we get the camera to “canonical form”?

- (Center of projection at the origin, x-axis points right, y-axis points up, z-axis points backwards)



Step 1: Translate by $-\mathbf{c}$

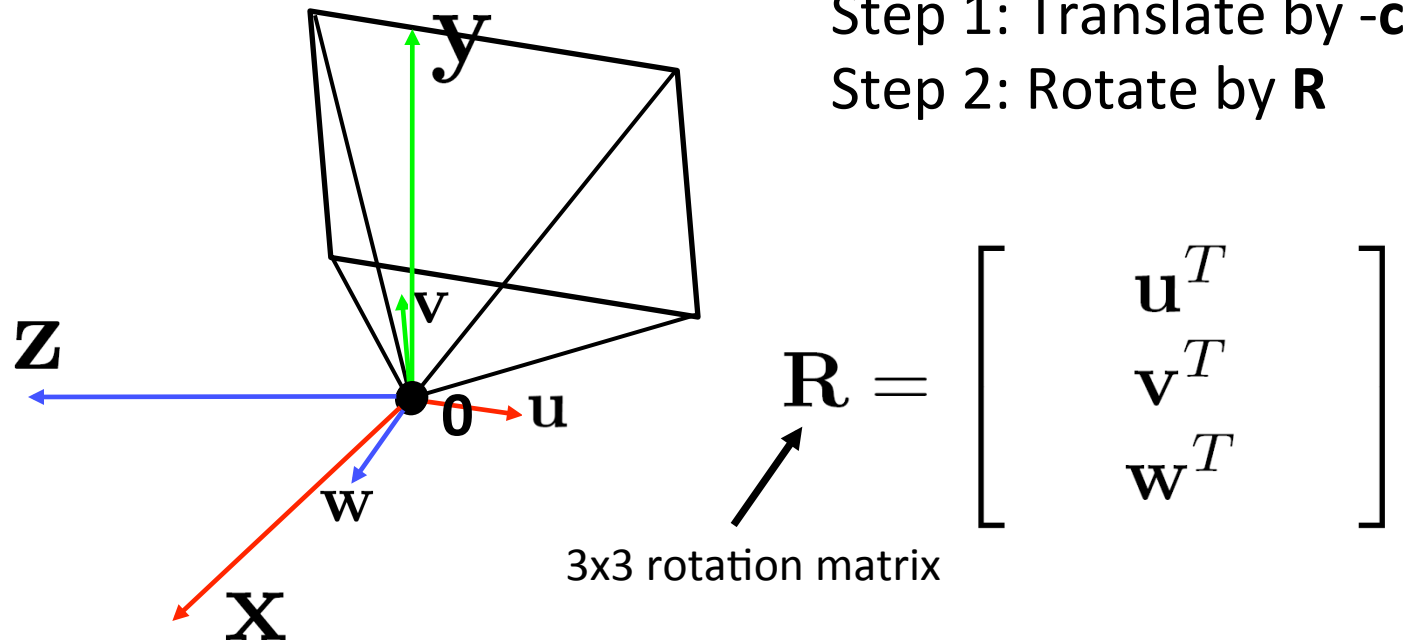
How do we represent translation as a matrix multiplication?

$$\mathbf{T} = \begin{bmatrix} \mathbf{I}_{3 \times 3} & -\mathbf{c} \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Extrinsics

How do we get the camera to “canonical form”?

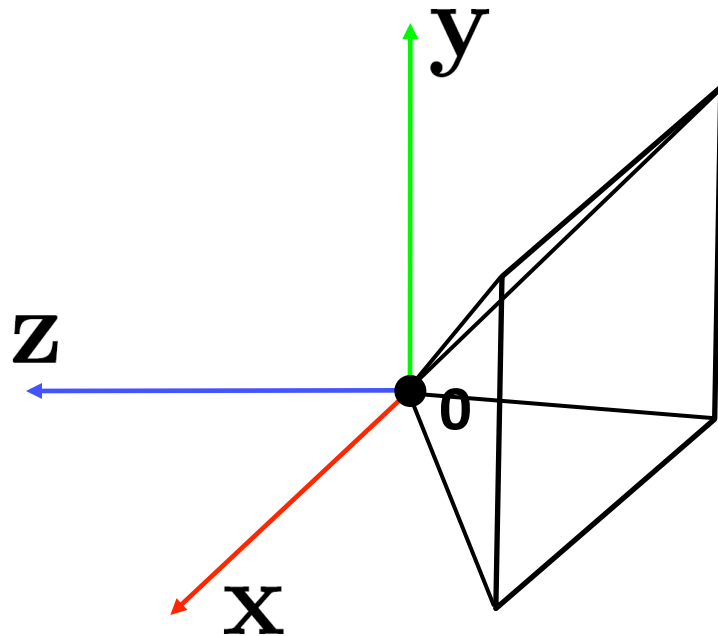
- (Center of projection at the origin, x-axis points right, y-axis points up, z-axis points backwards)



Extrinsics

How do we get the camera to “canonical form”?

- (Center of projection at the origin, x-axis points right, y-axis points up, z-axis points backwards)



Step 1: Translate by $-c$

Step 2: Rotate by \mathbf{R}

$$\mathbf{R} = \begin{bmatrix} \mathbf{u}^T \\ \mathbf{v}^T \\ \mathbf{w}^T \end{bmatrix}$$

Perspective projection

$$\underbrace{\begin{bmatrix} -f & 0 & 0 \\ 0 & -f & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{K}} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

K
(intrinsic)

(converts from 3D rays in camera coordinate system to pixel coordinates)

in general, $\mathbf{K} = \begin{bmatrix} -f & s & c_x \\ 0 & -\alpha f & c_y \\ 0 & 0 & 1 \end{bmatrix}$ (upper triangular matrix)

α : **aspect ratio** (1 unless pixels are not square)

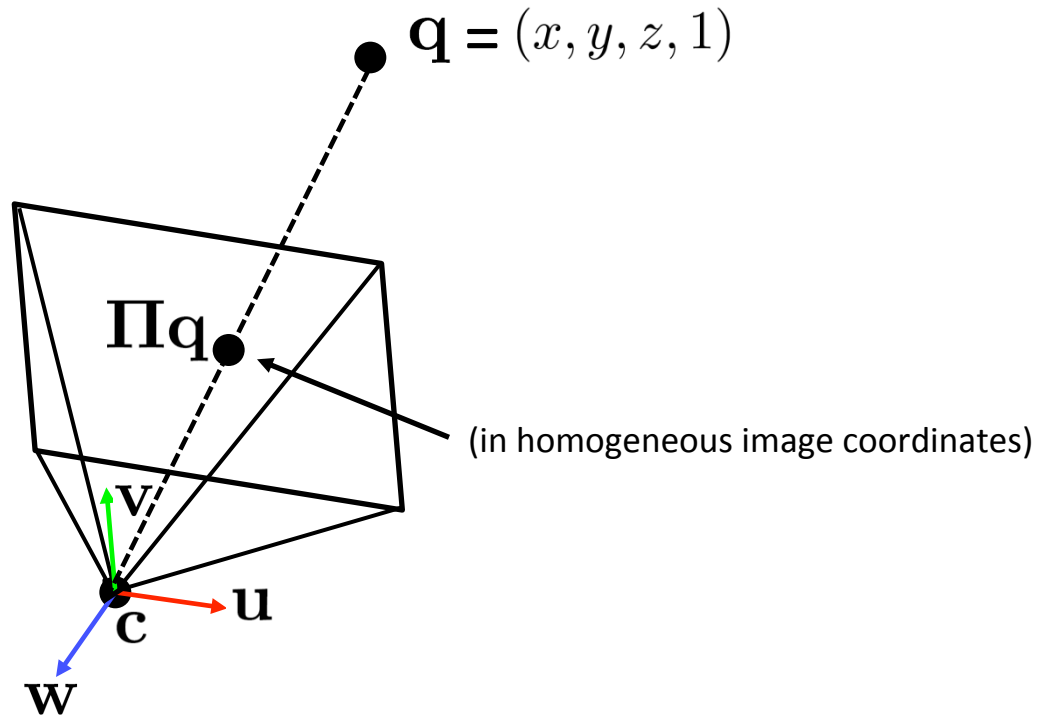
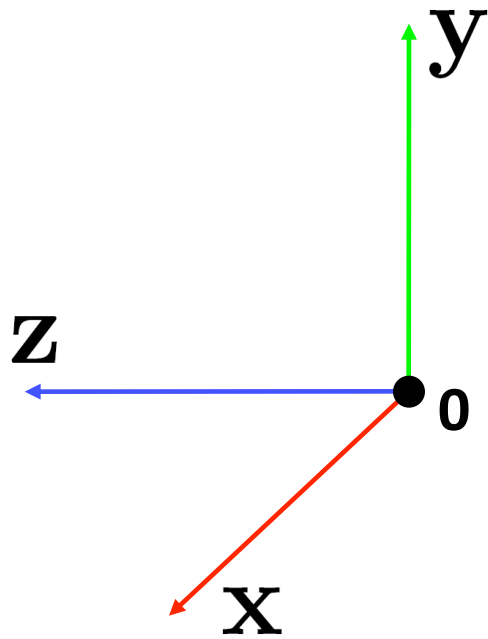
s : **skew** (0 unless pixels are shaped like rhombi/parallelograms)

(c_x, c_y) : **principal point** ((0,0) unless optical axis doesn't intersect projection plane at origin)

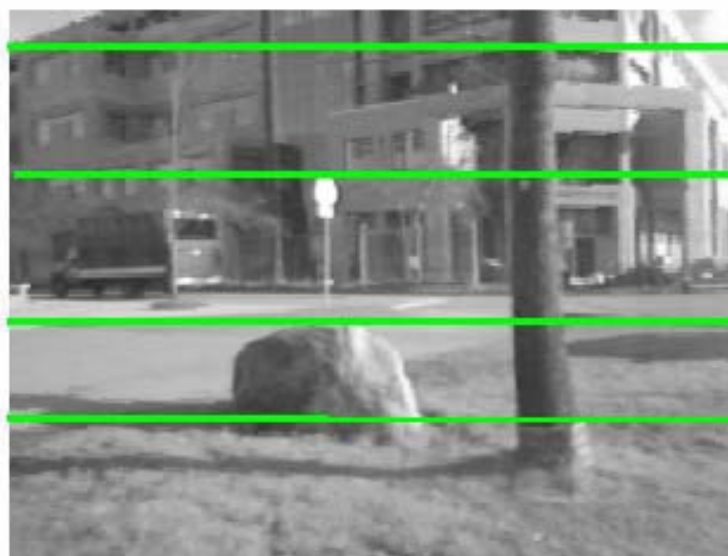
Projection matrix

$$\mathbf{\Pi} = \mathbf{K} \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}}_{\text{projection}} \underbrace{\begin{bmatrix} \mathbf{R} & \begin{matrix} 0 \\ 0 \\ 0 \end{matrix} \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{\text{rotation}} \underbrace{\begin{bmatrix} \mathbf{I}_{3 \times 3} & -\mathbf{c} \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{\text{translation}}$$

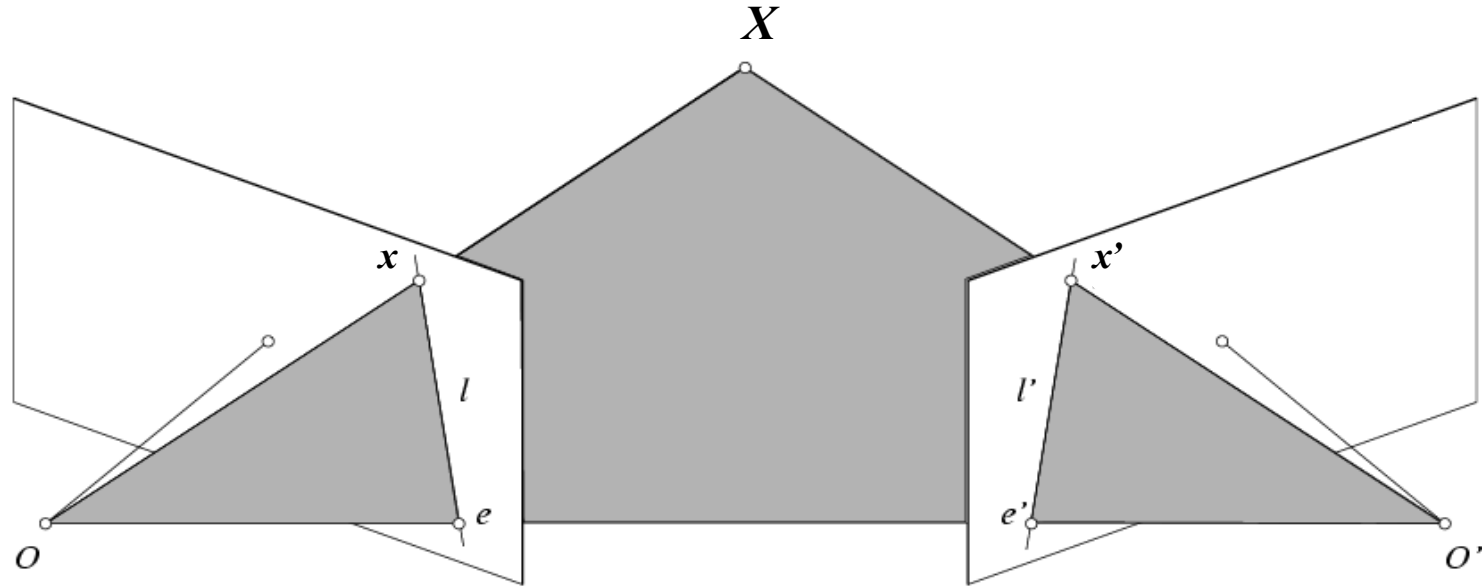
Projection matrix



Epipolar constraint example

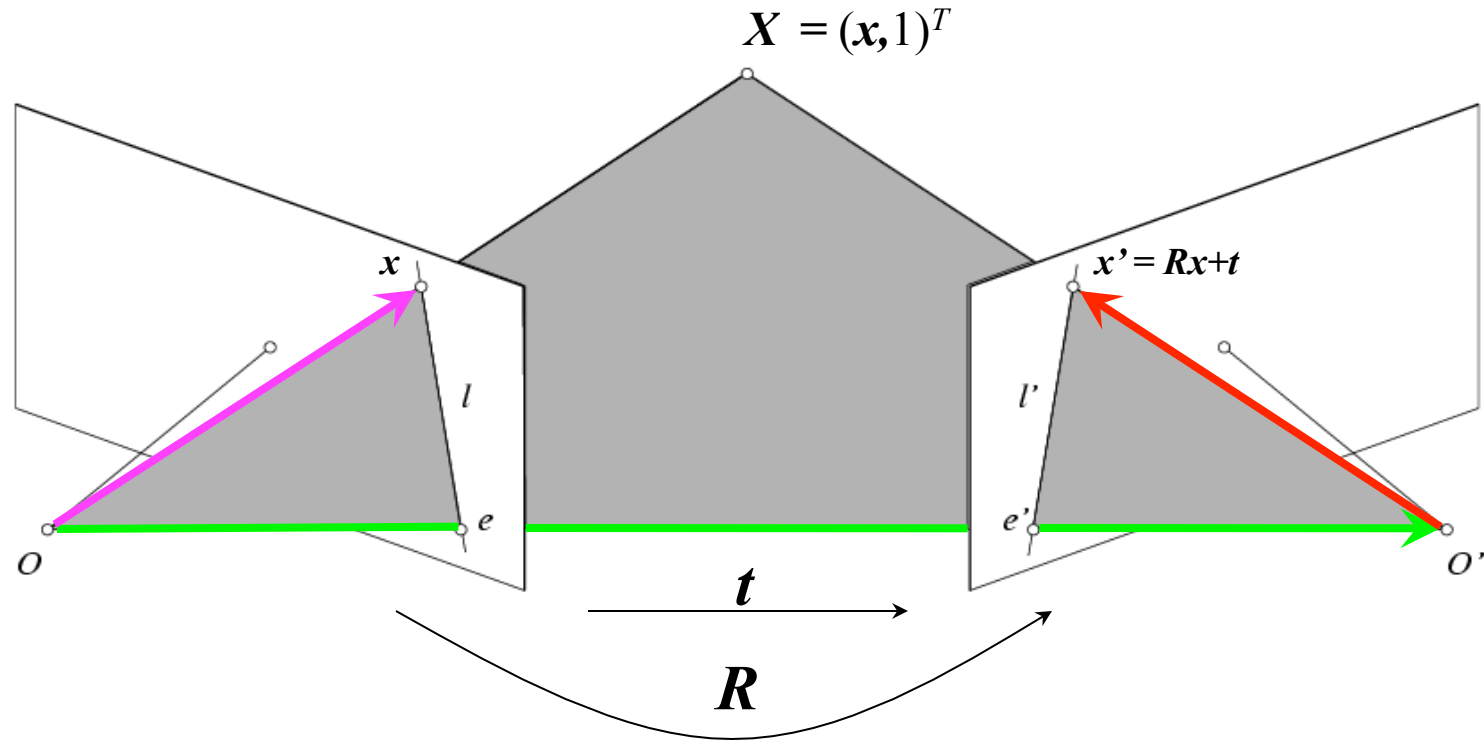


Epipolar constraint: Calibrated case



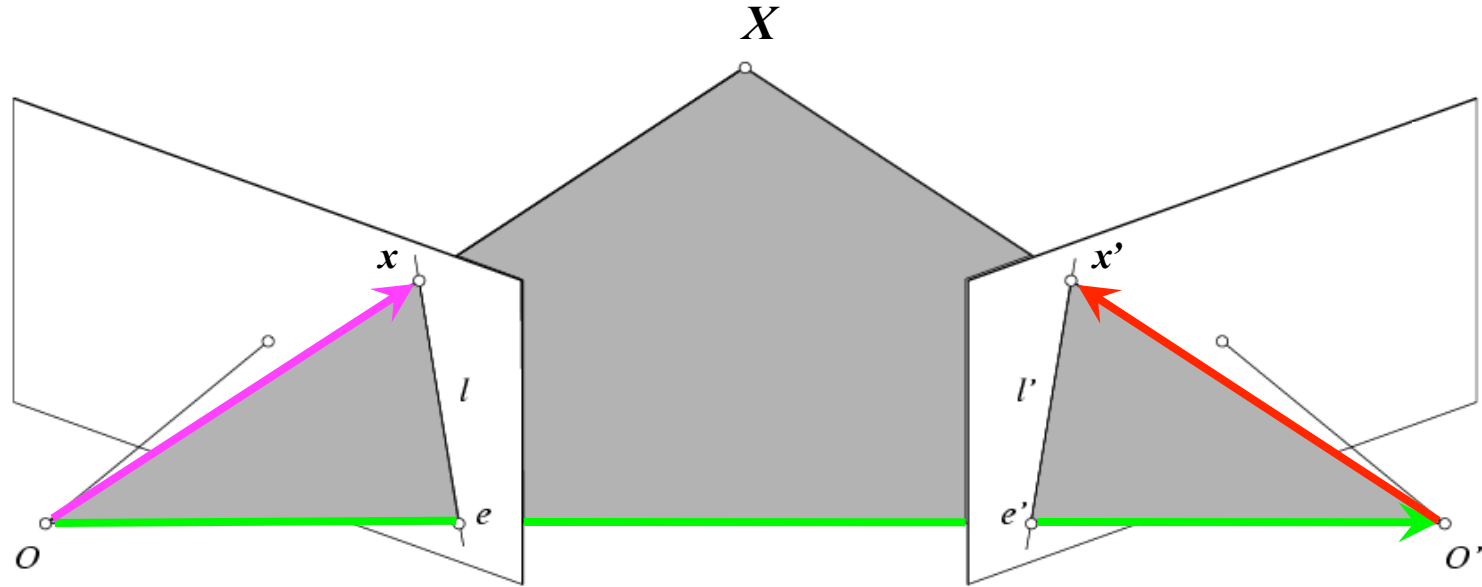
- Assume that the intrinsic and extrinsic parameters of the cameras are known
- We can multiply the projection matrix of each camera (and the image points) by the inverse of the calibration matrix to get *normalized* image coordinates
- We can also set the global coordinate system to the coordinate system of the first camera. Then the projection matrices of the two cameras can be written as $[\mathbf{I} \mid \mathbf{0}]$ and $[\mathbf{R} \mid \mathbf{t}]$

Epipolar constraint: Calibrated case



The vectors Rx , t , and x' are coplanar

Epipolar constraint: Calibrated case

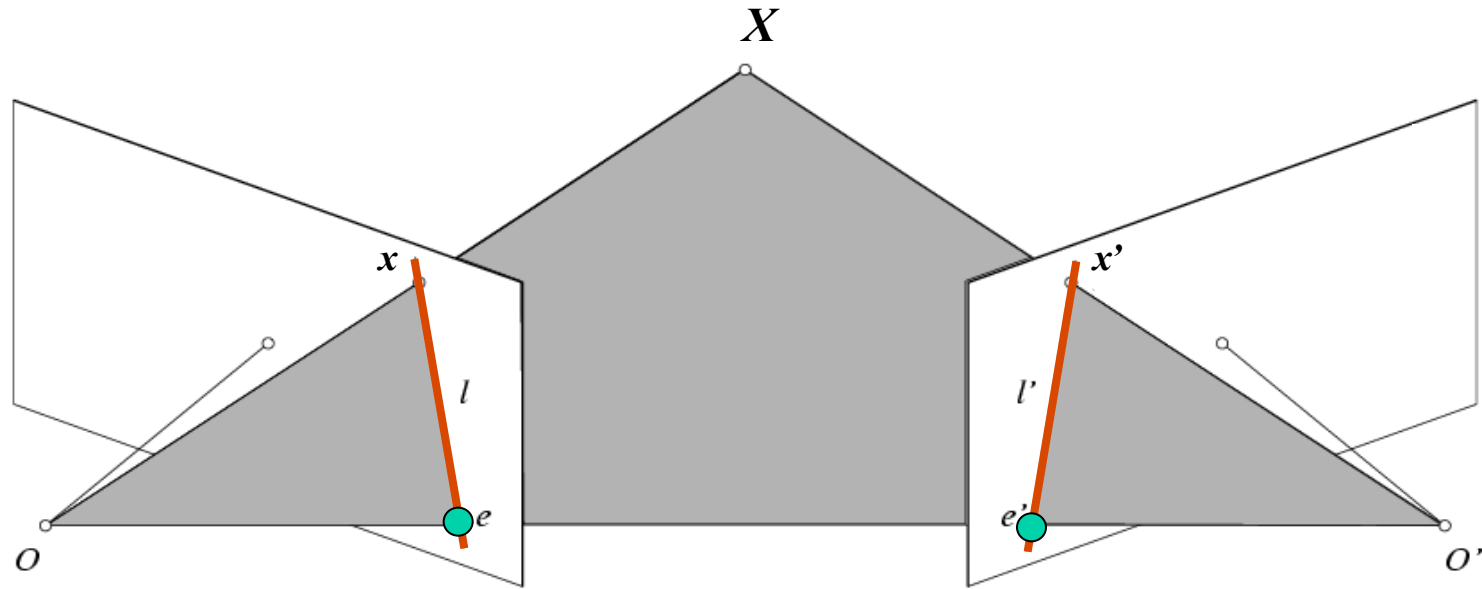


$$\mathbf{x}' \cdot [\mathbf{t} \times (\mathbf{R}\mathbf{x})] = 0 \quad \Rightarrow \quad \mathbf{x}'^T \mathbf{E} \mathbf{x} = 0 \quad \text{with} \quad \mathbf{E} = [\mathbf{t}_\times] \mathbf{R}$$

Essential Matrix
(Longuet-Higgins, 1981)

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

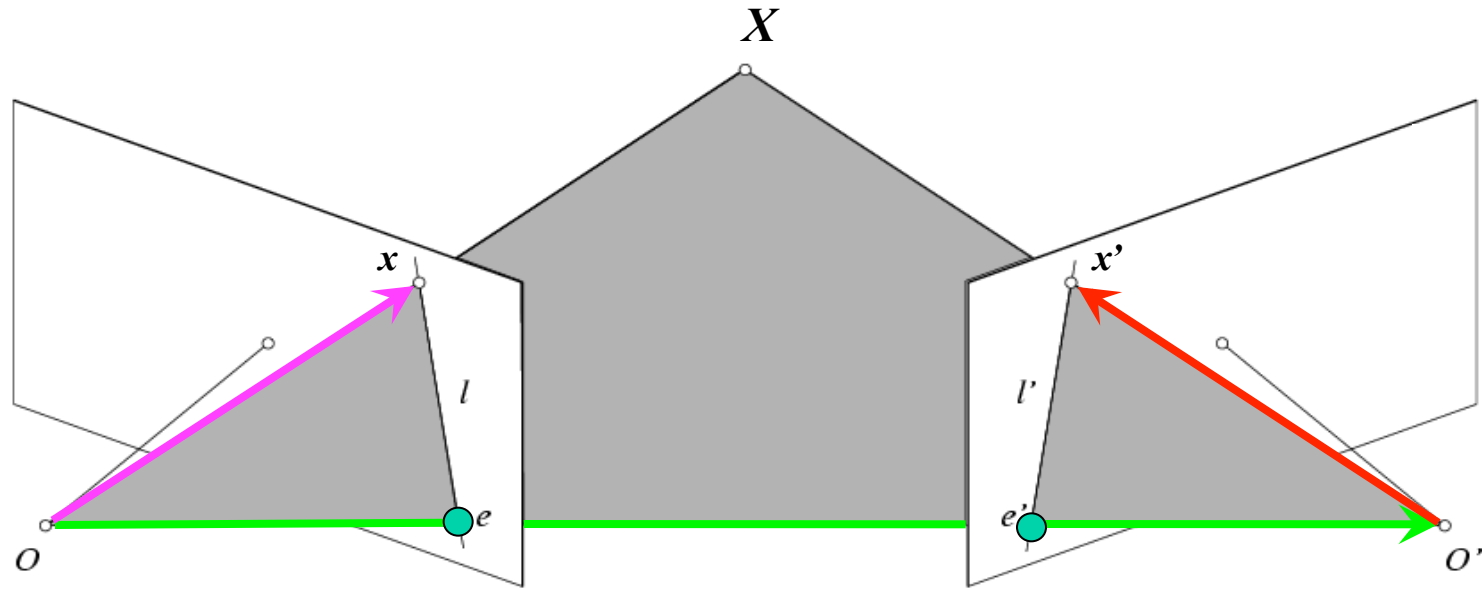
Epipolar constraint: Calibrated case



$$\mathbf{x}' \cdot [\mathbf{t} \times (\mathbf{R}\mathbf{x})] = 0 \quad \Rightarrow \quad \mathbf{x}'^T \mathbf{E} \mathbf{x} = 0 \quad \text{with} \quad \mathbf{E} = [\mathbf{t}_\times] \mathbf{R}$$

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

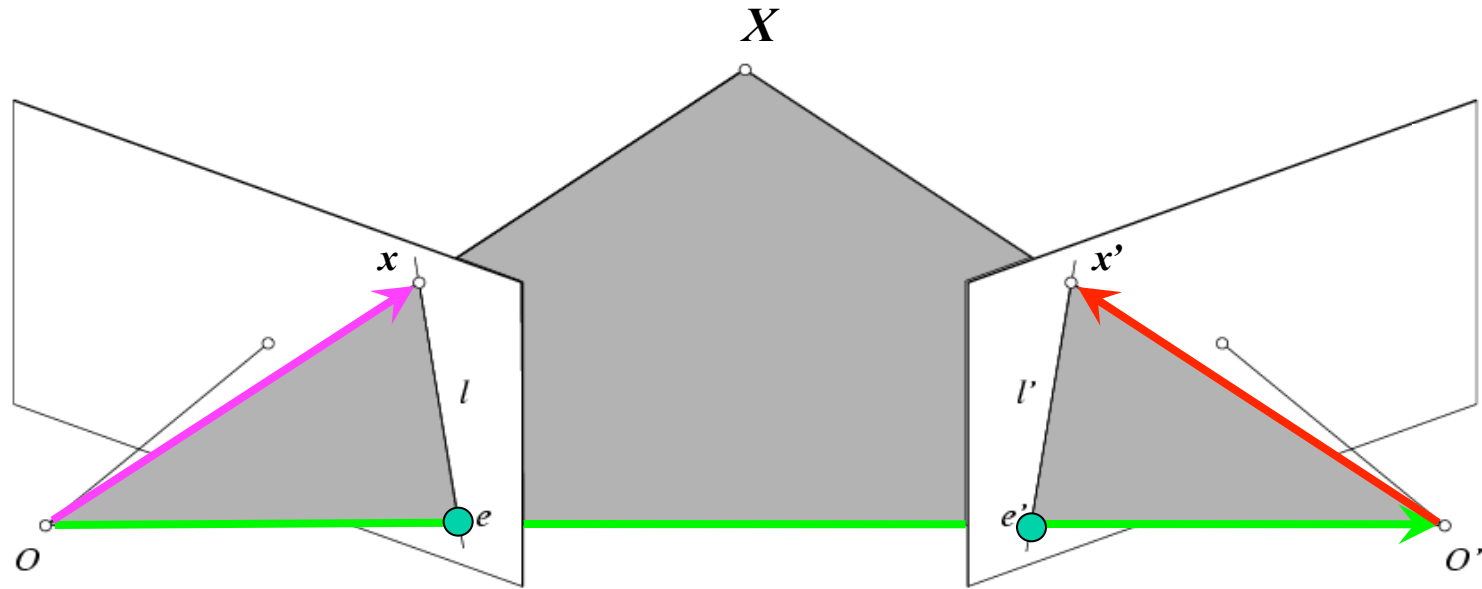
Epipolar constraint: Uncalibrated case



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

$$\hat{\mathbf{x}}'^T \mathbf{E} \hat{\mathbf{x}} = 0 \quad \hat{\mathbf{x}} = \mathbf{K}^{-1} \mathbf{x}, \quad \hat{\mathbf{x}}' = \mathbf{K}'^{-1} \hat{\mathbf{x}}'$$

Epipolar constraint: Uncalibrated case



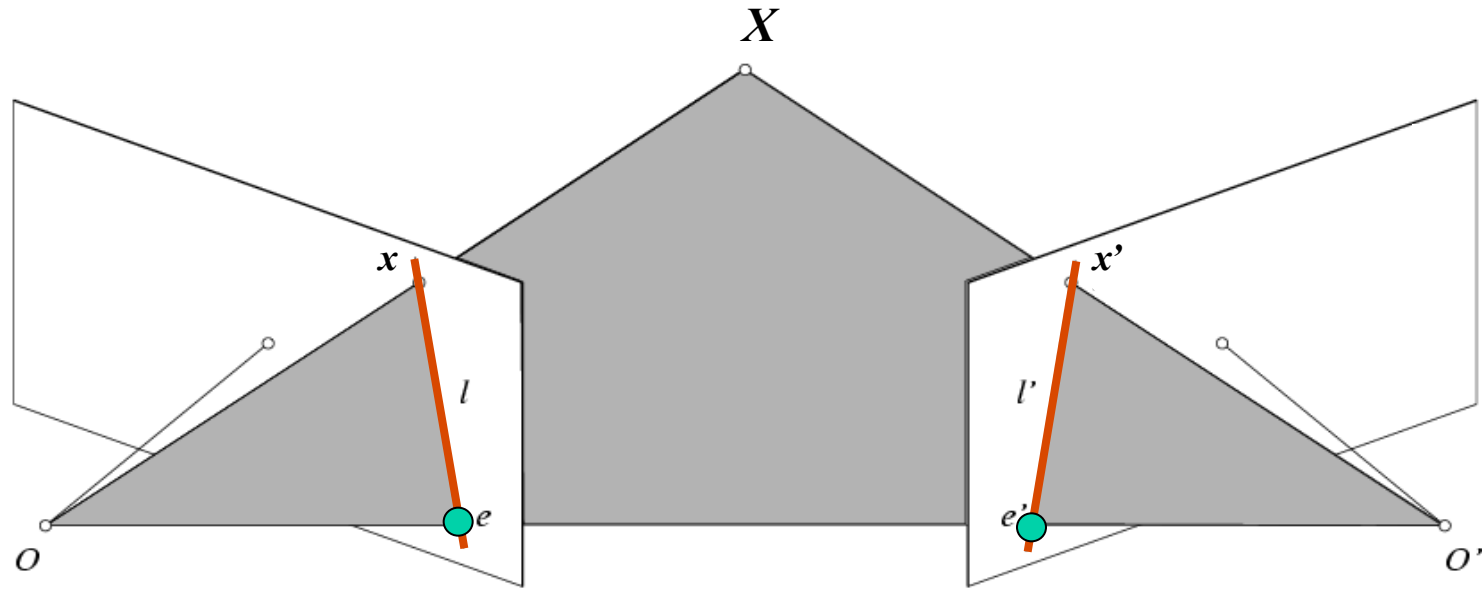
$$\hat{\mathbf{x}}'^T \mathbf{E} \hat{\mathbf{x}} = 0 \quad \Rightarrow \quad \mathbf{x}'^T \mathbf{F} \mathbf{x} = 0 \quad \text{with} \quad \mathbf{F} = \mathbf{K}'^{-T} \mathbf{E} \mathbf{K}^{-1}$$

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

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

Fundamental Matrix
(Faugeras and Luong, 1992)

Epipolar constraint: Uncalibrated case



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

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

The eight-point algorithm

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

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

\mathbf{A}

Minimize:

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

under the constraint

$$\|\mathbf{F}\|^2 = 1$$

Smallest
eigenvalue of
 $\mathbf{A}^T \mathbf{A}$