



# Machine Vision

## Lecture 8: Multi view geometry

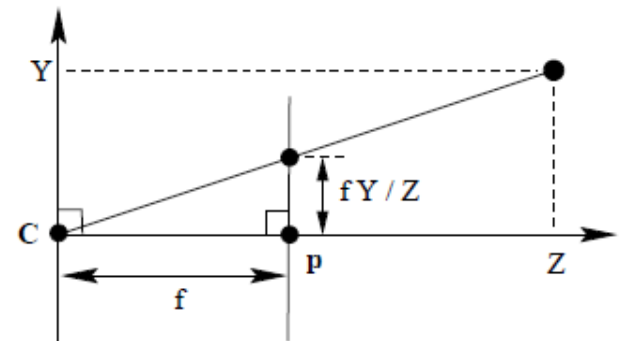
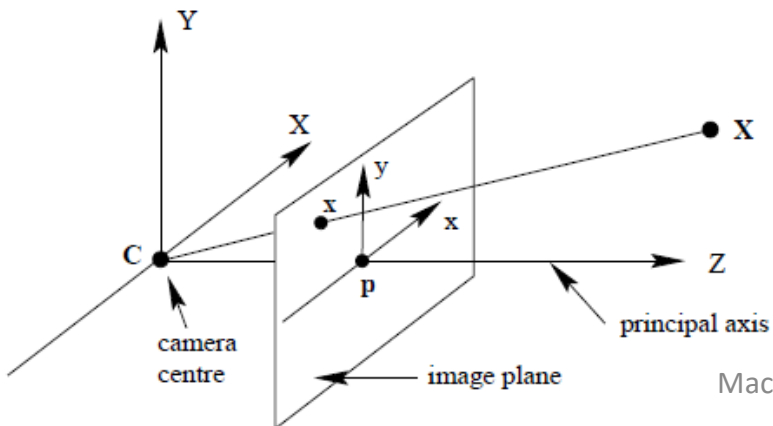
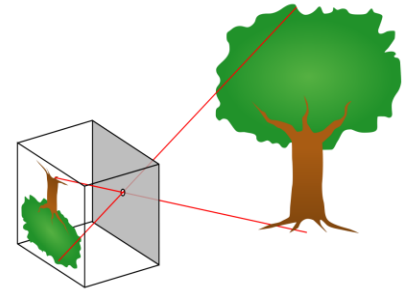
# Central projection model

- We will first look at cameras with finite centre of projection
- A pinhole camera with focal length  $f$  located at the coordinate origin projects a 3d point

$$(X, Y, Z) \rightarrow \left( \frac{fX}{Z}, \frac{fY}{Z} \right)$$

- In homogeneous coordinates this can be expressed as

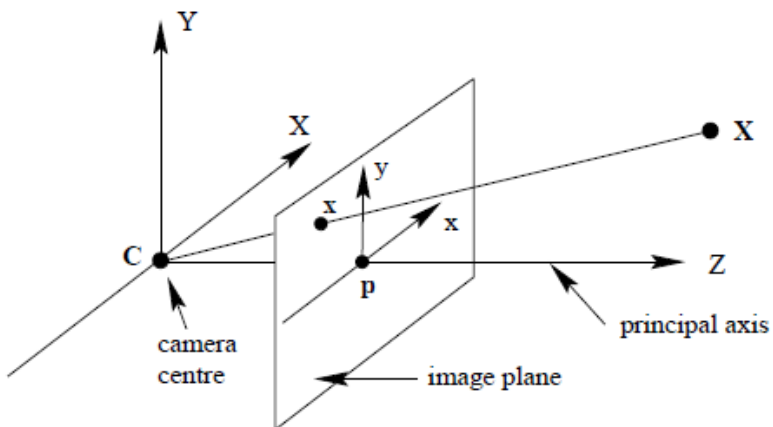
$$\begin{pmatrix} fX \\ fY \\ Z \end{pmatrix} = \begin{pmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$$



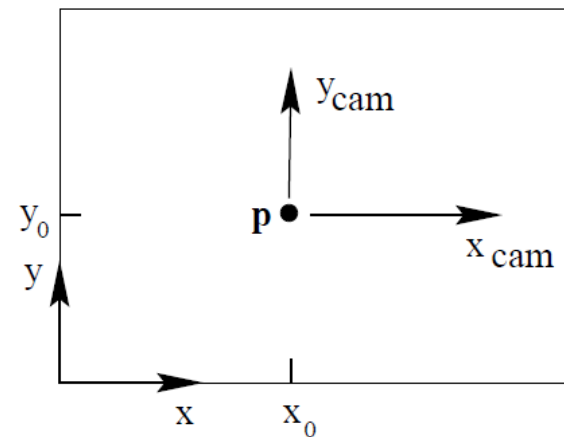
# Camera calibration matrix

- We assumed that the coordinate system of the image plane is the same as the object coordinate system
- This is in general not the case, as our object world is typically measured in metres [ $m^3$ ] and the image world is measured in pixels [ $pel^2$ ]
- Also the image coordinate system origin is not necessarily where the object Z-axis pierces the image plane, which we accommodate with translation by  $(x_0, y_0)$  and scaling by  $m_x$  and  $m_y$  of the image plane

$$\mathbf{x}' = \begin{pmatrix} fm_x & x_0 & 0 \\ fm_y & y_0 & 0 \\ 1 & 0 & 0 \end{pmatrix} \mathbf{X}$$



© Vision



# Camera calibration matrix

- We can summarise this image coordinate system transformation in the  $3 \times 3$  camera calibration matrix

$$K = \begin{pmatrix} fm_x & s & x_0 \\ & fm_y & y_0 \\ & & 1 \end{pmatrix} = \begin{pmatrix} c & s & x_0 \\ & \alpha c & y_0 \\ & & 1 \end{pmatrix}$$

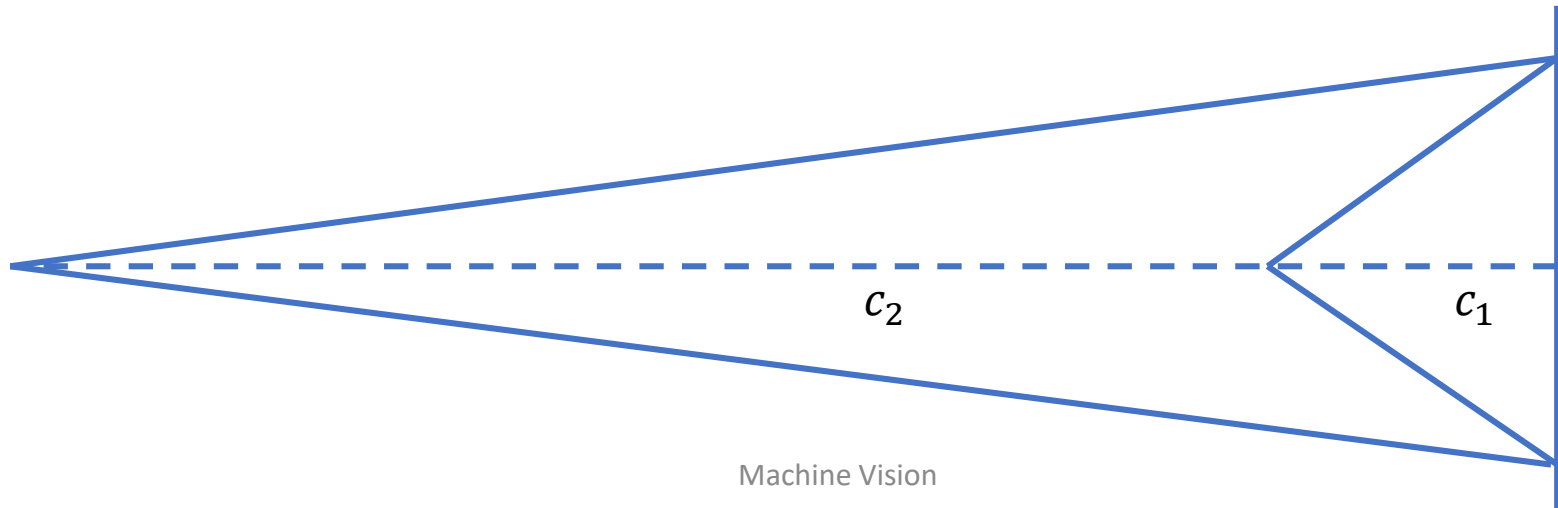
- Which has 5 parameters

- The **principal length**  $c = fm_x$
- The **aspect ratio**  $\alpha = m_y/m_x$
- The **principal point**  $(x_0, y_0)$
- The **image skew**  $s$

(modern CCD cameras are manufactured with a very regular pixel grid, so typically  $s = 0$  for digital images; however,  $s \neq 0$  often happens when processing scanned film)

# Changing the focal length (zoom)

- The principal length  $c$  is the distance between the projection centre and the image plane measured in the unit of the object coordinate system
- The longer the principal length, the narrower is the opening angle and the smaller is the field of view
- **Optical zoom** is changing distance of the image plane to the centre of projection, i.e. affecting the principal length  $c$



# Backward projection of rays

- The calibration matrix not only tells us how 3d points are projected into 2d coordinates
- We can also reverse the equation and obtain a way of calculating the direction in space corresponding to the image point

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

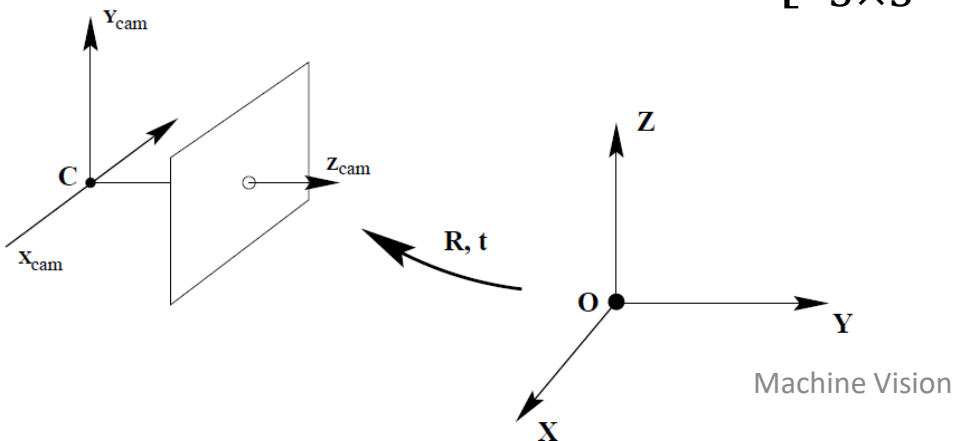
- The distance  $\lambda$  from the camera is unknown, but we know that the 3d scene point is somewhere on the line

$$\mathbf{X} = \lambda \frac{\mathbf{m}}{\sqrt{\mathbf{m}^T \mathbf{m}}}$$

# Object coordinate system

- Thus far we have assumed the camera to be located in the origin of the 3d coordinate system
- Obviously this is not the case, and we need to accommodate this transformation between **camera coordinate system** and **world coordinate system**
- If we apply the translation and rotation to every image point we can express the full transformation from object to image coordinate system in homogeneous coordinates as follows

$$x = KR^T [I_{3 \times 3} \quad -t]X$$



# The projective camera

- The projective camera can be described as a homogeneous  $3 \times 4$  matrix  $\mathbf{P} \in \mathbb{P}^{11}$  transforming homogeneous 3d scene points  $\mathbf{X} \in \mathbb{P}^3$  into homogeneous 2d image coordinates  $\mathbf{x} \in \mathbb{P}^2$  as

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

- The projection matrix can be decomposed into

$$\mathbf{P} = \begin{pmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{pmatrix} = \underbrace{\begin{pmatrix} c & s & x_0 \\ & \alpha c & y_0 \\ & & 1 \end{pmatrix}}_{\mathbf{K}} \underbrace{\begin{pmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{pmatrix}}_{\mathbf{R}}^T \begin{pmatrix} 1 & & & -t_1 \\ & 1 & & -t_2 \\ & & 1 & -t_3 \end{pmatrix}$$

- The 11dof of  $\mathbf{P}$  are distributed across the camera calibration  $\mathbf{K}$  containing 5dof, the rotation of the camera  $\mathbf{R}$  containing 3dof, and the position of the camera  $\mathbf{t}$  containing 3dof



# The projective camera

- To get the camera position  $\mathbf{t}$  from the projection matrix we note that

$$\mathbf{P} \begin{pmatrix} \mathbf{t} \\ 1 \end{pmatrix} = \mathbf{K}\mathbf{R}^T [\mathbf{I}_{3 \times 3} \quad -\mathbf{t}] \begin{pmatrix} \mathbf{t} \\ 1 \end{pmatrix} = \mathbf{K}\mathbf{R}^T \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} = \mathbf{0}$$

- Therefore the centre of projection is the right-null-space of  $\mathbf{P}$
- If we cut

$$\mathbf{P} = [\mathbf{M} \quad \mathbf{p}]$$

into a  $3 \times 3$  submatrix

$$\mathbf{M} = \mathbf{K}\mathbf{R}^T$$

and a 3-vector

$$\mathbf{p} = -\mathbf{K}\mathbf{R}^T \mathbf{t}$$

Then we can easily calculate it using the following insight:

$$-\mathbf{M}^{-1}\mathbf{p} = \mathbf{R}\mathbf{K}^{-1}\mathbf{K}\mathbf{R}^T \mathbf{t} = \mathbf{t}$$

# The projective camera

- To separate

$$\mathbf{M} = \mathbf{K}\mathbf{R}^T$$

- We note that this is a product of an upper diagonal matrix

$$\mathbf{K} = \begin{pmatrix} c & s & x_0 \\ & \alpha c & y_0 \\ & & 1 \end{pmatrix}$$

- and a rotation matrix  $\mathbf{R}^T$ , which we can separate using the RQ decomposition algorithm

```
K,Rt = scipy.linalg.rq(M)
```

# General back-projection

- We already saw that the direction of the of a ray from the origin in the un-rotated coordinate frame is

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

- If we apply the coordinate system transformation, the ray originating from the centre of projection into this direction in the world coordinate frame is

$$\mathbf{X} = \mathbf{t} + \lambda \mathbf{R} \mathbf{K}^{-1} \mathbf{x}'$$

- Putting this all together we obtain a way to calculate the ray of 3d points corresponding to an image coordinate

$$\mathbf{X} = -\mathbf{M}^{-1}\mathbf{p} + \lambda \mathbf{M}^{-1}\mathbf{x}' = \mathbf{M}^{-1}(\lambda \mathbf{x} - \mathbf{p})$$

# Triangulation of points

- Given two corresponding points  $x$  and  $x'$  in two images

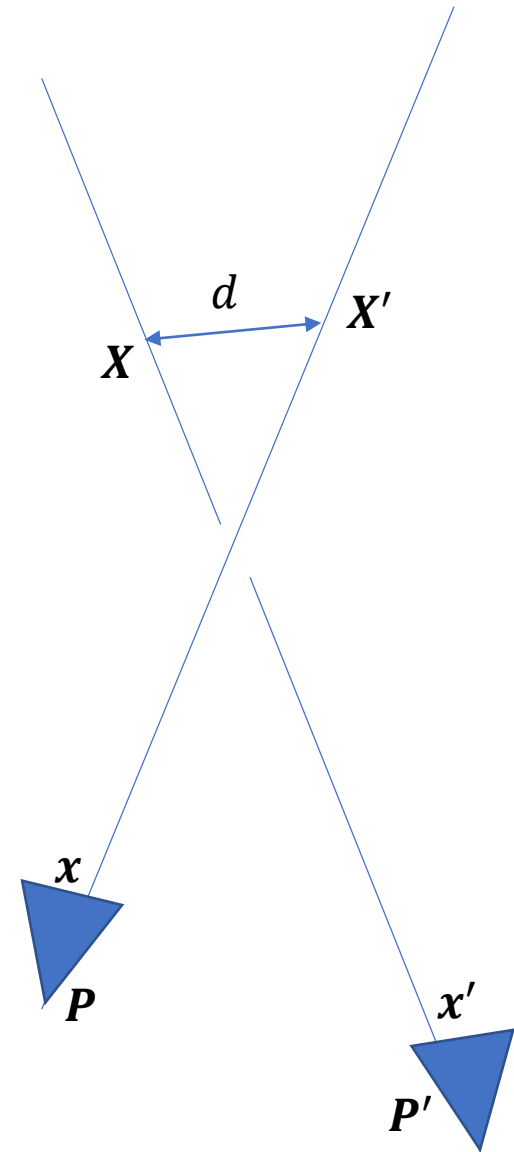
$$\begin{aligned} P &= (M \quad p) \\ P' &= (M' \quad p') \end{aligned}$$

- The 3d point must be on both rays

$$\begin{aligned} X &= M^{-1}(\lambda x - p) \\ X' &= M'^{-1}(\mu x' - p') \end{aligned}$$

- In the presence of noise, this is not exactly the case, therefore we try to find the point where the rays are closest, i.e. minimise the distance

$$d = |M^{-1}(\lambda x - p) - M'^{-1}(\mu x' - p')|^2$$



# Triangulation of points

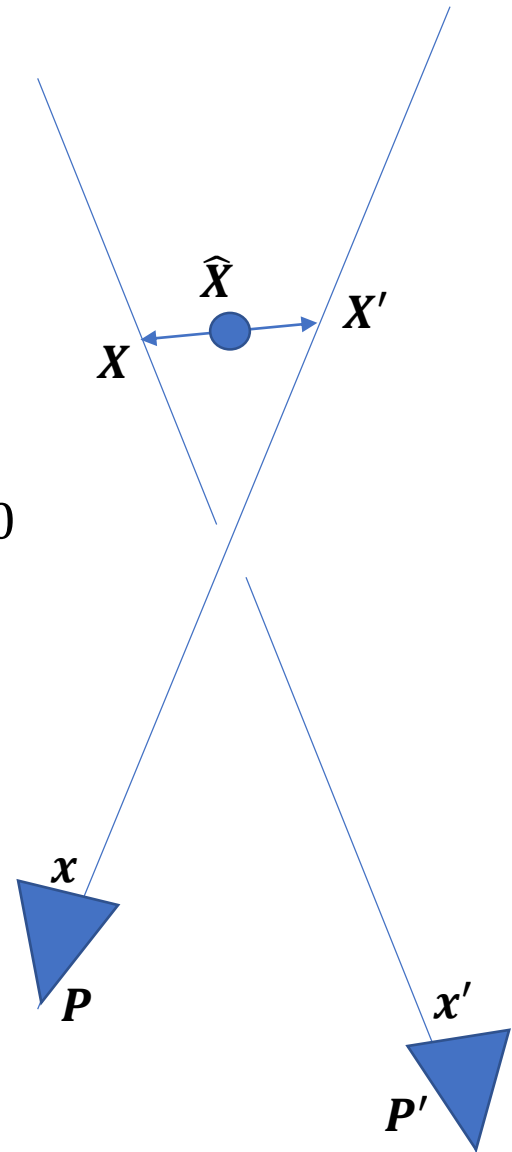
- To minimise the distance we look at the derivatives

$$\frac{\partial d}{\partial \lambda} = 2 \left( \mathbf{M}^{-1}(\lambda \mathbf{x} - \mathbf{p}) - \mathbf{M}'^{-1}(\mu \mathbf{x}' - \mathbf{p}') \right) \mathbf{M}^{-1} \mathbf{x} = 0$$

$$\frac{\partial d}{\partial \mu} = 2 \left( \mathbf{M}^{-1}(\lambda \mathbf{x} - \mathbf{p}) - \mathbf{M}'^{-1}(\mu \mathbf{x}' - \mathbf{p}') \right) \mathbf{M}'^{-1} \mathbf{x}' = 0$$

- This are two linear equations, which we can easily solve for the unknown  $\lambda$  and  $\mu$
- The triangulated point is then determined half-way between the two points closest to each other

$$\hat{\mathbf{X}} = \frac{\mathbf{X} + \mathbf{X}'}{2}$$



# Back projection of lines

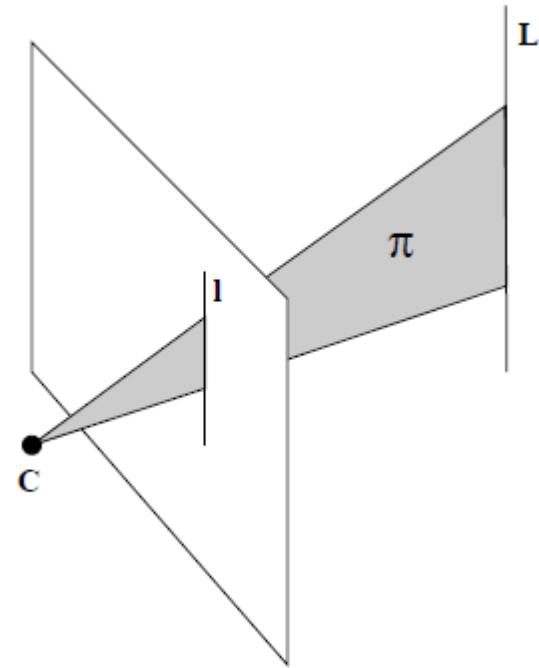
- A point  $\mathbf{x}$  is on a line  $\mathbf{l}$  if  $\mathbf{l}^T \mathbf{x} = 0$
- All 3d points  $\mathbf{X}$  that project somewhere on this line must fulfil

$$\mathbf{l}^T \mathbf{x} = \underbrace{\mathbf{l}^T \mathbf{P}}_{\boldsymbol{\pi}^T} \mathbf{X} = 0$$

- This can be considered as a plane equation of the 3d plane

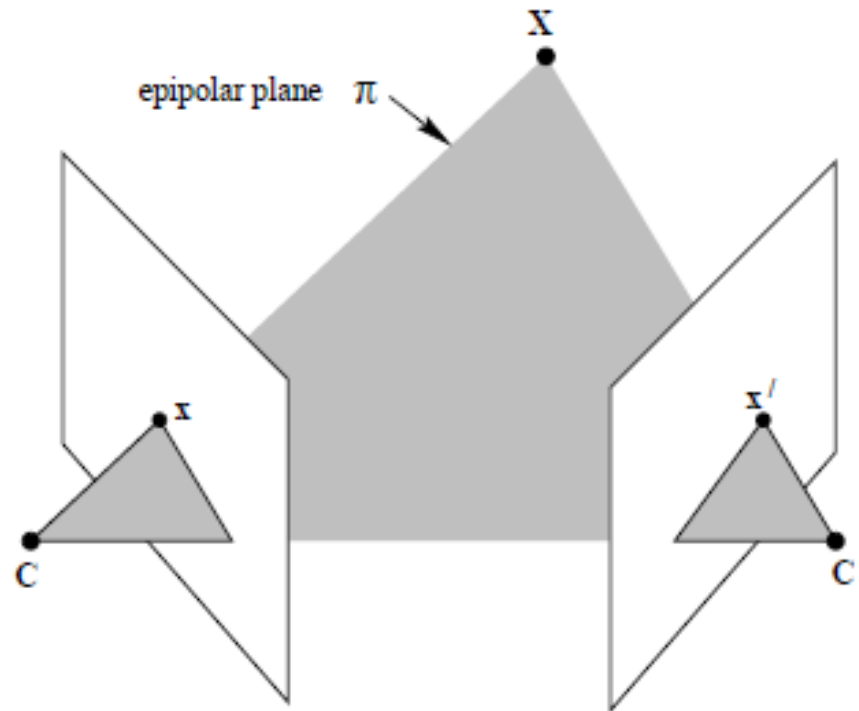
$$\boldsymbol{\pi} = \mathbf{P}^T \mathbf{l}$$

- Which is the back-projection of the image line  $\mathbf{l}$



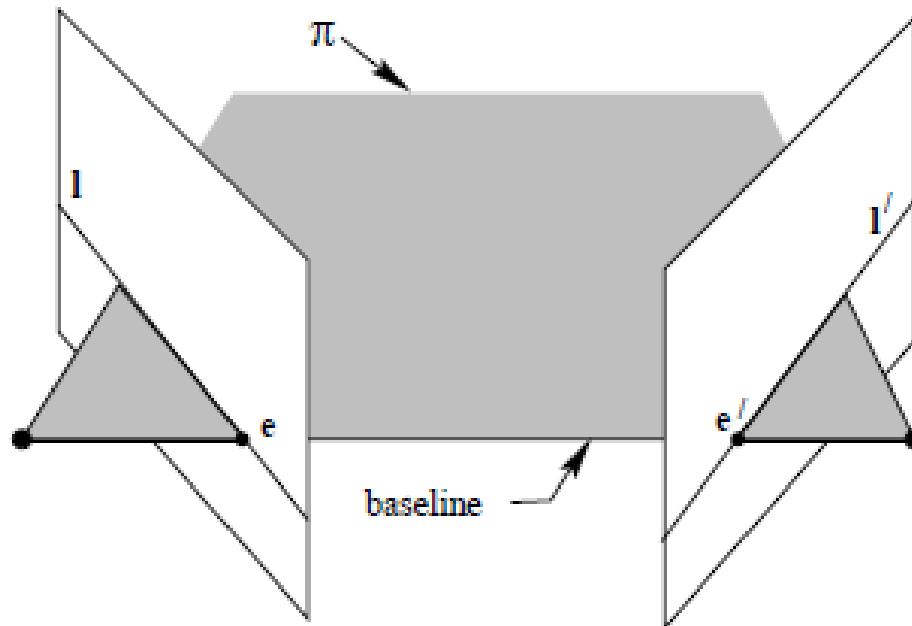
# Epipolar geometry

- A single 3d point that is visible in two images defines an **epipolar plane** in 3d space, which connects the point and the two centres of projection



# Epipolar geometry

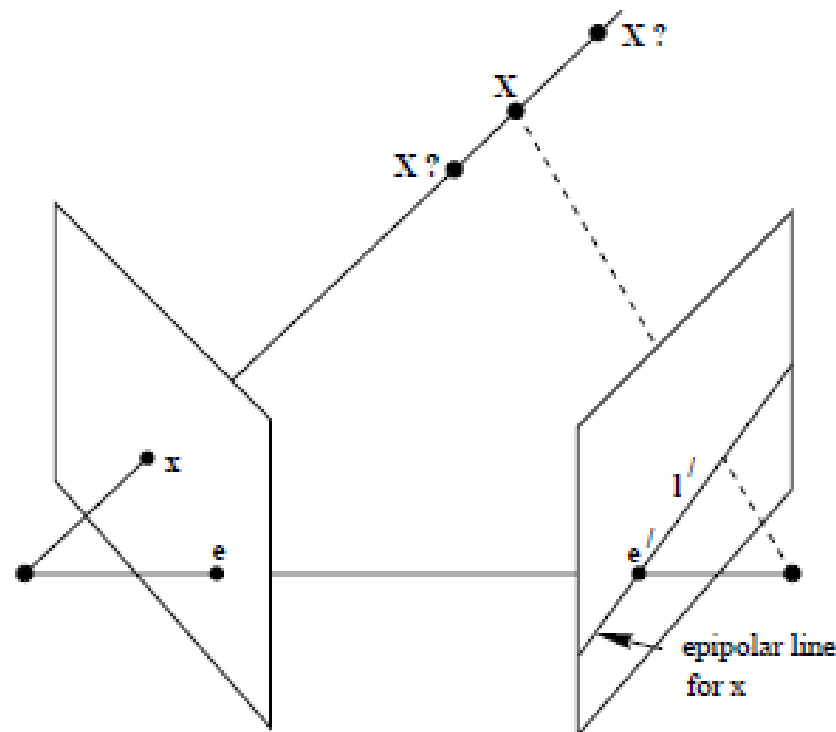
- Because all these epipolar planes go through the centres of projection, and therefore through the **baseline** between the two images, they create corresponding **epipolar lines**
- All epipolar lines intersect in the **epipoles**, which are the intersections of the baseline with the image planes



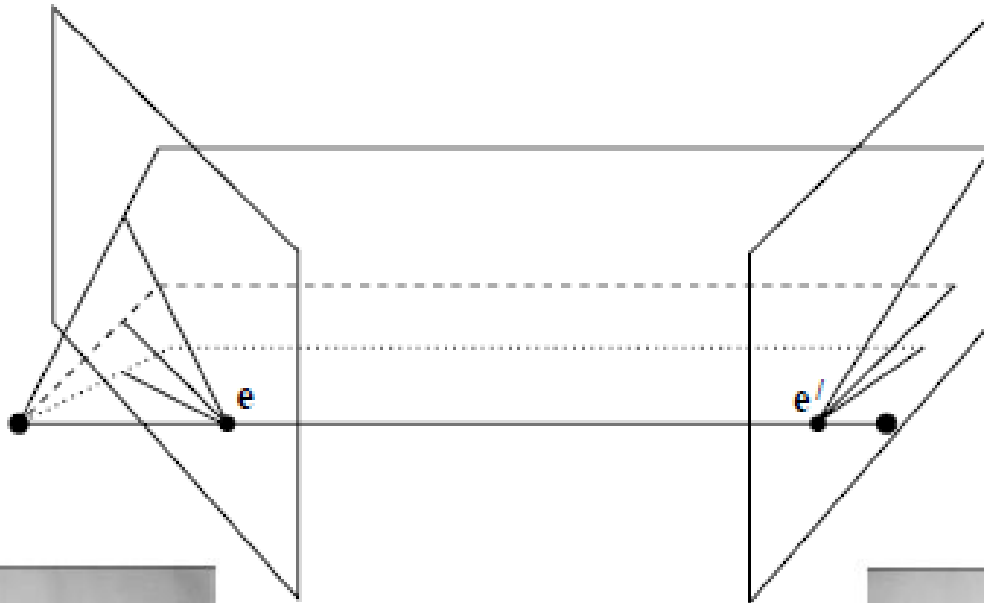


# Epipolar geometry

- When looking for a point correspondence of a point  $x$  in another image, we don't know where it is due to the unknown distance
- However, knowing the epipolar geometry we can restrict the search to the **epipolar line** in the second image



# Epipolar geometry



# Fundamental matrix

- Let the two camera matrices be

$$\begin{aligned} \mathbf{P} &= \mathbf{K}(\mathbf{I} \quad \mathbf{0}) \\ \mathbf{P}' &= \mathbf{K}'\mathbf{R}^T(\mathbf{I} \quad -\mathbf{t}) \end{aligned}$$

- Then the point  $\mathbf{x}$  in the first image back-projects to the line

$$\mathbf{X} = \lambda \mathbf{K}^{-1} \mathbf{x}$$

- which projects into the second image at

$$\mathbf{x}' = \mathbf{K}'\mathbf{R}^T(\mathbf{X} - \mathbf{t}) = \mathbf{K}'\mathbf{R}^T(\lambda \mathbf{K}^{-1} \mathbf{x} - \mathbf{t})$$

- The epipole is the image of the centre of projection ( $\lambda = 0$ )

$$\mathbf{e}' = -\mathbf{K}'\mathbf{R}^T \mathbf{t}$$

- Now the epipolar line through the epipole  $\mathbf{e}'$  and  $\mathbf{x}'$  is

$$\mathbf{l} = \mathbf{e}' \times \mathbf{x}' = \mathbf{e}' \times (\lambda \mathbf{K}'\mathbf{R}^T \mathbf{K}^{-1} \mathbf{x} + \mathbf{e}') = \mathbf{S}[\mathbf{K}'\mathbf{R}^T \mathbf{t}] \mathbf{K}'\mathbf{R}^T \mathbf{K}^{-1} \mathbf{x}$$

# Fundamental matrix

- A point  $\mathbf{x}'$  in the second image is on the epipolar line

$$\mathbf{l} = S[\mathbf{K}'\mathbf{R}^T\mathbf{t}]\mathbf{K}'\mathbf{R}^T\mathbf{K}^{-1}\mathbf{x}$$

- If

$$\mathbf{l}^T\mathbf{x}' = \mathbf{x}'^T\mathbf{F}\mathbf{x} = 0$$

- with the  $3 \times 3$  **fundamental matrix**

$$\mathbf{F} = S[\mathbf{K}'\mathbf{R}^T\mathbf{t}]\mathbf{K}'\mathbf{R}^T\mathbf{K}^{-1}$$

# Fundamental matrix

- In conclusion, two points  $\mathbf{x}$  and  $\mathbf{x}'$  can only refer to the same scene point if they obey the following equation

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

- Obviously, this equation is homogeneous, i.e. the scale of  $\mathbf{F}$  does not alter the result
- Also, because  $\mathcal{S}[\mathbf{e}']$  has rank 2, the fundamental matrix is always singular

$$\det \mathbf{F} = 0$$

- These two conditions mean that the fundamental matrix has 7 degrees of freedom

# Fundamental matrix

- If  $\mathbf{F}$  is the fundamental matrix of the image pair  $(\mathbf{P}, \mathbf{P}')$ , then  $\mathbf{F}^T$  is the fundamental matrix of the image pair  $(\mathbf{P}', \mathbf{P})$
- For a point  $\mathbf{x}$  in the first image, the epipolar line in the second image is

$$\mathbf{l}' = \mathbf{F}\mathbf{x}$$

- For a point  $\mathbf{x}'$  in the second image, the epipolar line in the first image is

$$\mathbf{l} = \mathbf{F}^T\mathbf{x}'$$

# Calculating the epipoles

- The fundamental matrix is singular, and the epipoles are the left and right null-spaces

$$\begin{aligned} \mathbf{F}\mathbf{e} &= \mathbf{0} \\ \mathbf{F}^T\mathbf{e}' &= \mathbf{0} \end{aligned}$$

- To calculate the epipoles we can use the singular value decomposition, with the epipole being the singular vector corresponding to the smallest singular value of  $\mathbf{F}$

```
U, S, V = np.linalg.svd(F)
e = V[2, :]
```

# Estimating the Fundamental matrix

- The fundamental matrix can be calculated from 7 point correspondences  $\mathbf{x}'_i \leftrightarrow \mathbf{x}_i$
- Each point correspondence provides a condition

$$\mathbf{x}'_i{}^T \mathbf{F} \mathbf{x}_i = 0$$

- Or equivalent using the Kronecker product

$$\underbrace{(\mathbf{x}_i^T \otimes \mathbf{x}'_i{}^T)}_{\mathbf{a}_i^T} \text{vec}[\mathbf{F}] = 0$$



# Estimating the Fundamental matrix

- These 7 equations can be stacked into a  $7 \times 9$  matrix

$$A = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_7^T \end{pmatrix}$$

- For which we calculate the two null-vectors  $A\mathbf{f}_1 = \mathbf{0}$  and  $A\mathbf{f}_2 = \mathbf{0}$  using the singular value decomposition

```
U, S, V = np.linalg.svd(A)
F1 = V[8, :].reshape(3, 3).T
F2 = V[7, :].reshape(3, 3).T
```

# Estimating the Fundamental matrix

- The fundamental matrix we are looking for is now

$$\mathbf{F} = \alpha \mathbf{F}_1 + (1 - \alpha) \mathbf{F}_2$$

- To determine the value for alpha we use the singularity constraint

$$\det(\alpha \mathbf{F}_1 + (1 - \alpha) \mathbf{F}_2) = 0$$

- This is a degree 3 polynomial in the unknown  $\alpha$ , which we can easily solve

# Estimating the Fundamental matrix

- In case there are 8 points or more, we can also apply the DLT algorithm we have seen before by stacking all points into

$$A = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_8^T \end{pmatrix}$$

- The fundamental matrix is then found as the singular vector corresponding to the smallest singular value
- In this case the singularity constraint needs to be applied (again using the singular value decomposition)

```
U,S,V = np.linalg.svd(A)
```

```
F = V[8,:].reshape(3,3).T
```

```
U,S,V = np.linalg.svd(F)
```

```
F = np.matmul(U,np.matmul(np.diag([S[0],S[1],0]),V))
```

# Projective invariance

- The fundamental matrix for a pair of cameras  $(\mathbf{P}, \mathbf{P}')$  and a pair of cameras  $(\mathbf{P}\mathbf{H}, \mathbf{P}'\mathbf{H})$  is the same for all 3d homographies  $\mathbf{H}$
- Therefore, the knowing the fundamental matrix determines the 3d scene up to a 3d projective transformation only
- If necessary we can choose the canonical cameras

$$\begin{aligned}\mathbf{P} &= [\mathbf{I} \quad \mathbf{0}] \\ \mathbf{P}' &= [\mathbf{S}[\mathbf{e}']\mathbf{F} \quad \mathbf{e}']\end{aligned}$$

- And determine the necessary homography later from other information (camera calibration)

# Image rectification

- If apply any 2d homographies  $\mathbf{H}$  and  $\mathbf{H}'$  to both images

$$\begin{aligned}\hat{\mathbf{x}} &= \mathbf{H}\mathbf{x} \\ \hat{\mathbf{x}}' &= \mathbf{H}'\mathbf{x}'\end{aligned}$$

- then the fundamental matrix between the two images transforms according to

$$\hat{\mathbf{F}} = \mathbf{H}'^{-T} \mathbf{F} \mathbf{H}^{-1}$$

- This is particularly useful, if we want to transform the images to achieve a given target fundamental matrix  $\hat{\mathbf{F}}$

# Image rectification

- An important special case is related to how our two eyes are arranged:
  - Both eyes are identical:  $\mathbf{K} = \mathbf{K}'$
  - Both eyes look into the same direction:  $\mathbf{R} = \mathbf{I}$
  - The translation is horizontal only:  $\mathbf{t} = (b \quad 0 \quad 0)^T$
- In this case the fundamental matrix is

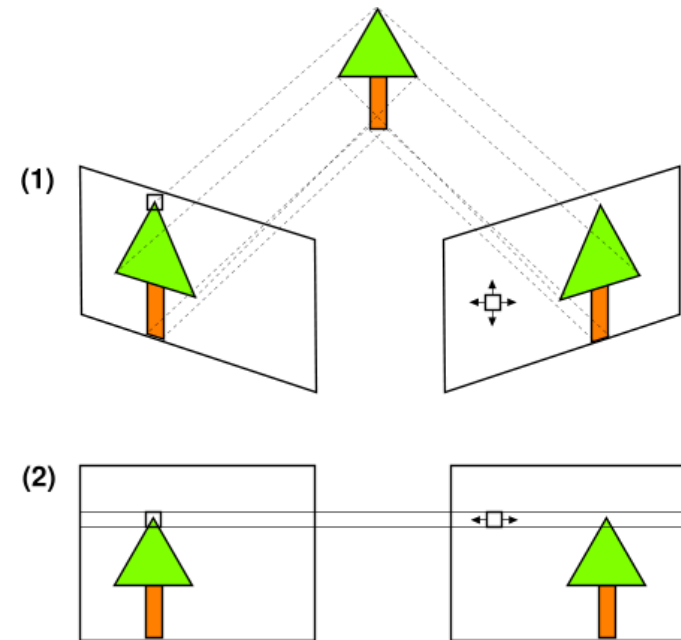
$$\mathbf{F} = \mathbf{S}[\mathbf{K}'\mathbf{R}^T\mathbf{t}]\mathbf{K}'\mathbf{R}^T\mathbf{K}^{-1} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & -cb \\ 0 & cb & 0 \end{pmatrix}$$

# Image rectification

- To achieve this special configuration we therefore need to find homographies so that

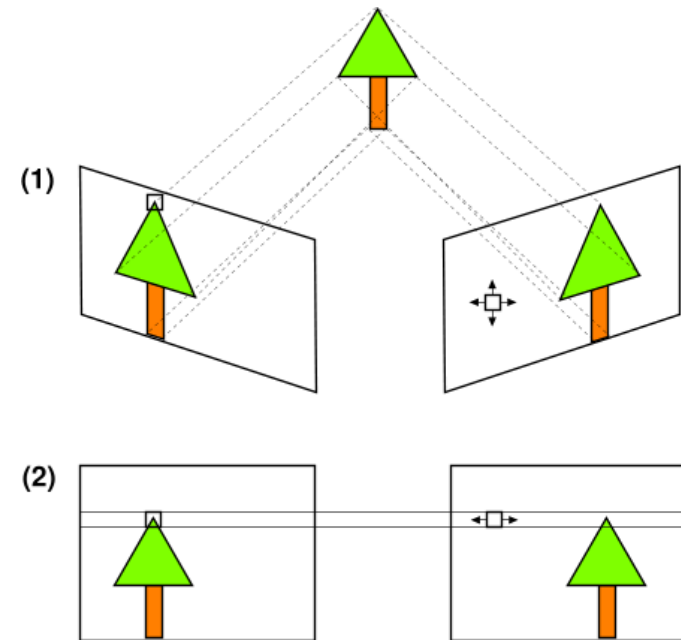
$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{pmatrix} = \mathbf{H}'^{-T} \mathbf{F} \mathbf{H}^{-1}$$

- There are many homographies  $\mathbf{H}$  and  $\mathbf{H}'$  that fulfil these equations
- Typically we will choose these transformations so that they minimally distort the original input images
- We also make sure that corresponding epipolar lines align



# Disparity

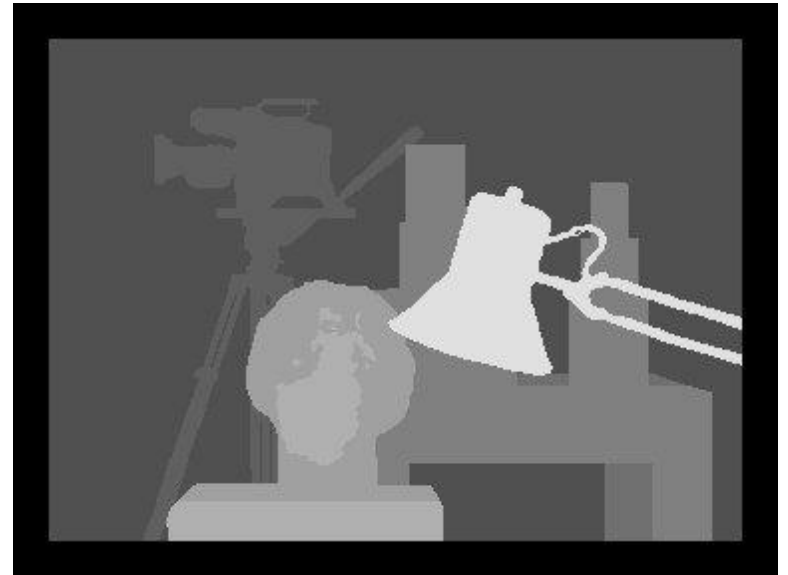
- Corresponding epipolar lines are aligned in a rectified image, therefore the depth of a 3d point only affects the horizontal displacement between the images
- The matching problem for a rectified image pair is therefore 1-dimensional
- This horizontal displacement between the images is called **disparity**





# Dense stereo

- Algorithms that solve this 1d search problem and calculate the disparity, and therefore the depth, for each pixel are called **dense stereo** algorithms



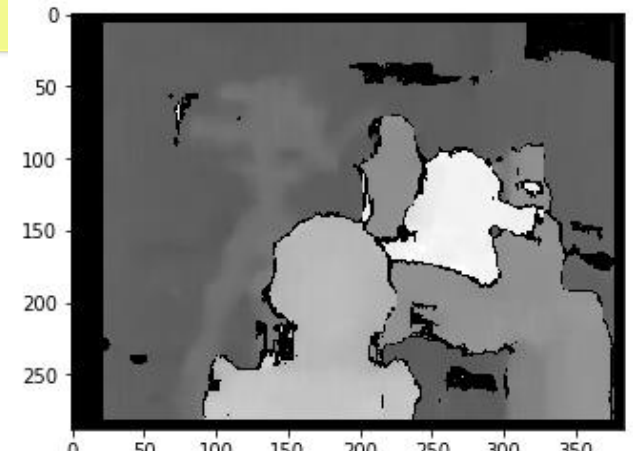
# Dense stereo

- There are several dense stereo algorithms for rectified images, most of them combining smoothness constraints and similarity measures

Block Matching compares patches to calculate similarity metric

```
stereo = cv2.createStereoBM(numDisparities=16, blockSize=15)  
disparity = stereo.compute(imgL, imgR)
```

The disparity is linked to depth, therefore the search range can be limited



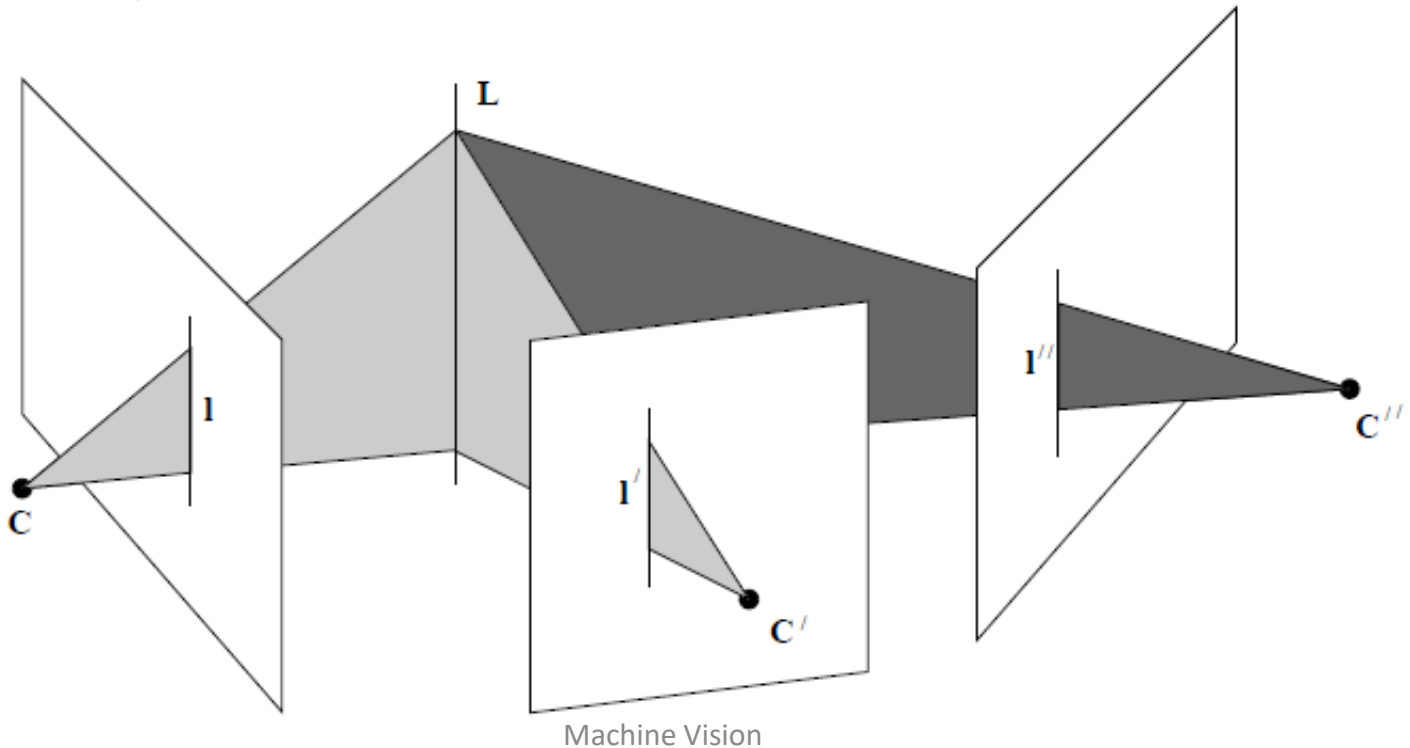
# Stereoscopic images

- Another application of rectified images is to present them to each eye individually
- Because this disparity estimation is how our spatial perception works, humans will perceive the scene in 3d then



# Trifocal tensor

- We now will look very briefly at the geometry of three images
- Three lines  $\mathbf{l}$ ,  $\mathbf{l}'$ ,  $\mathbf{l}''$  must all back-project onto a single line  $\mathbf{L}$  in space
- These back-projections are the three planes  $\pi = \mathbf{P}^T \mathbf{l}$ ,  $\pi' = \mathbf{P}'^T \mathbf{l}'$  and  $\pi'' = \mathbf{P}''^T \mathbf{l}''$



# Trifocal tensor

- We now will look briefly at the geometry of three images
- Three lines  $\mathbf{l}, \mathbf{l}', \mathbf{l}''$  must all back-project onto a single line  $\mathbf{L}$  in space
- These back-projections are the three planes  $\boldsymbol{\pi} = \mathbf{P}^T \mathbf{l}, \boldsymbol{\pi}' = \mathbf{P}'^T \mathbf{l}'$  and  $\boldsymbol{\pi}'' = \mathbf{P}''^T \mathbf{l}''$
- All points  $\mathbf{X}$  on the line  $\mathbf{L}$  must therefore be incident to all three lines, i.e.

$$\begin{pmatrix} \boldsymbol{\pi}^T \\ \boldsymbol{\pi}'^T \\ \boldsymbol{\pi}''^T \end{pmatrix} \mathbf{X} = \mathbf{0}$$

- Because the line  $\mathbf{L}$  is a one-dimensional entity (in addition to the homogeneity of the equation), the null-space of this matrix must be 2-dimensional
- This is called the tri-focal constraint

# Trifocal tensor

- The tri-focal constraint can be expressed as stating that the line

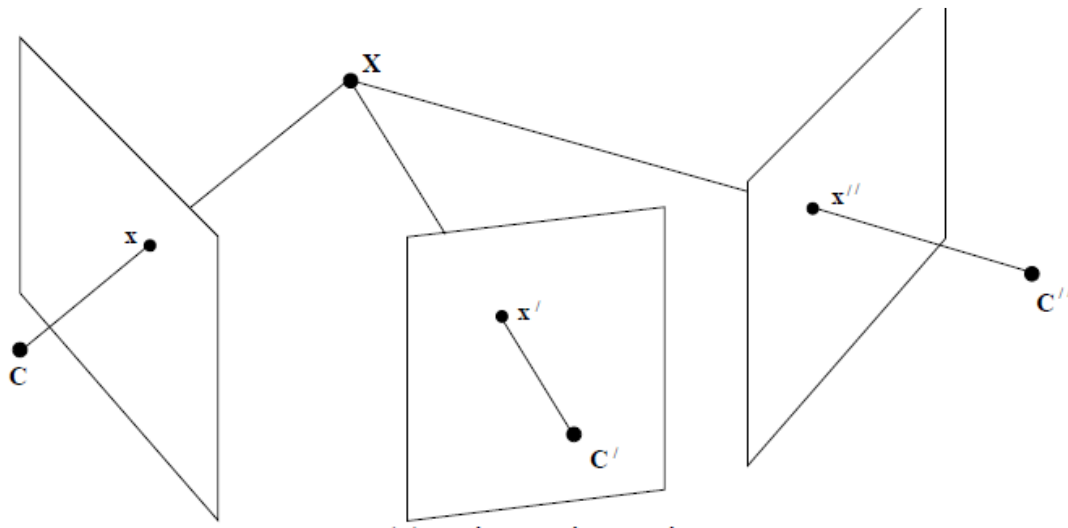
$$l = \begin{pmatrix} l'^T \mathcal{T}_1 l'' \\ l'^T \mathcal{T}_2 l'' \\ l'^T \mathcal{T}_3 l'' \end{pmatrix}$$

- The  $3 \times 3 \times 3$  tensor  $[\mathcal{T}_1, \mathcal{T}_2, \mathcal{T}_3]$  describing this relationship is called the **tri-focal tensor**

# Tri-focal geometry

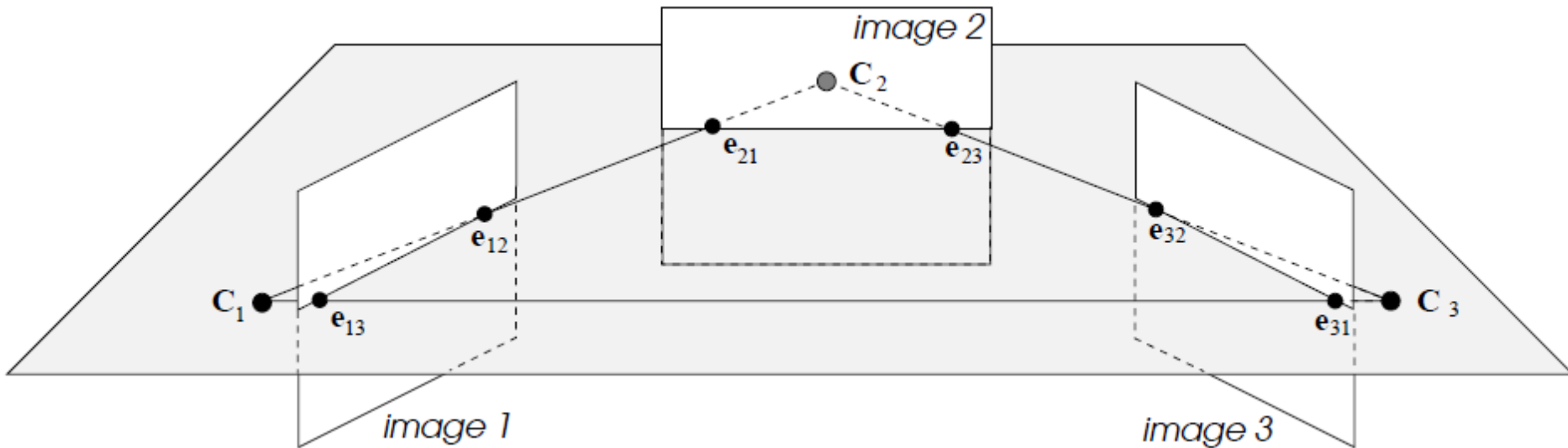
- Why is this important? Why can't we just use the pair-wise relationships provided by the epipolar geometry of mutual pairs of images?
- If we, for example want to transfer a point correspondence  $x \leftrightarrow x'$  from one image pair into a third image, we could simply calculate the intersection of the epipolar lines in that image, i.e.

$$x'' = (F_{31}x) \times (F_{32}x')$$



# Tri-focal geometry

- Why is this important? Why can't we just use the pair-wise relationships provided by the epipolar geometry of mutual pairs of images?
- Unfortunately, this point transfer via epipolar lines does not work in the **tri-focal plane** connecting all three projection centres
- Point transfer via the tri-focal tensor is possible, though





Thank you for your attention!