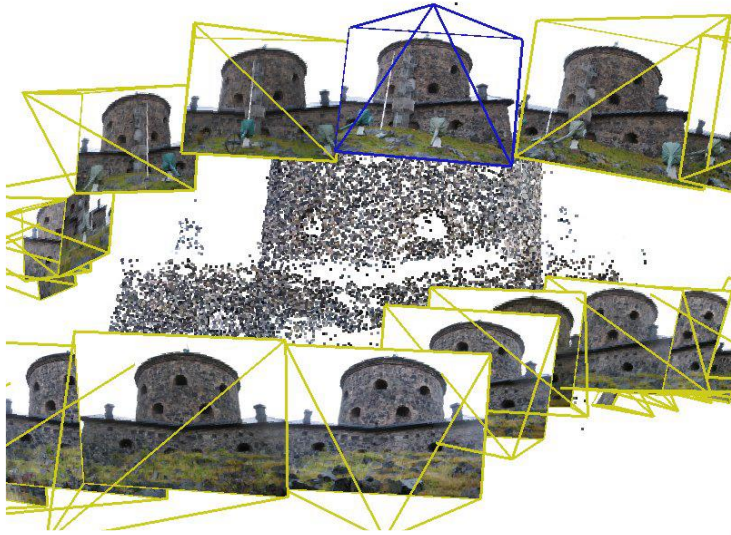
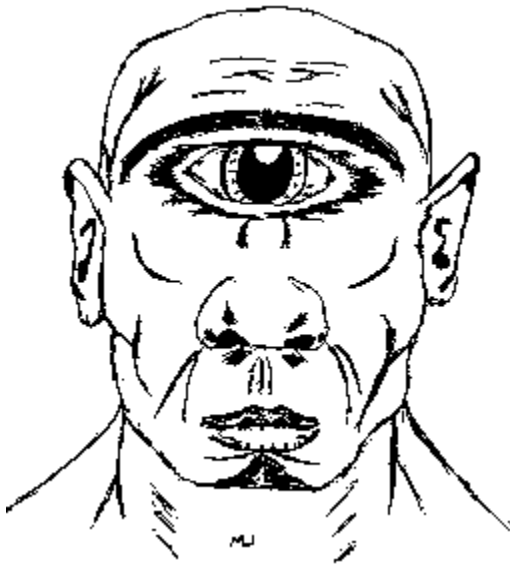


Multi-View Geometry



...with materials from H&Z and Carl Olsson

Motivation: why do we have two eyes?



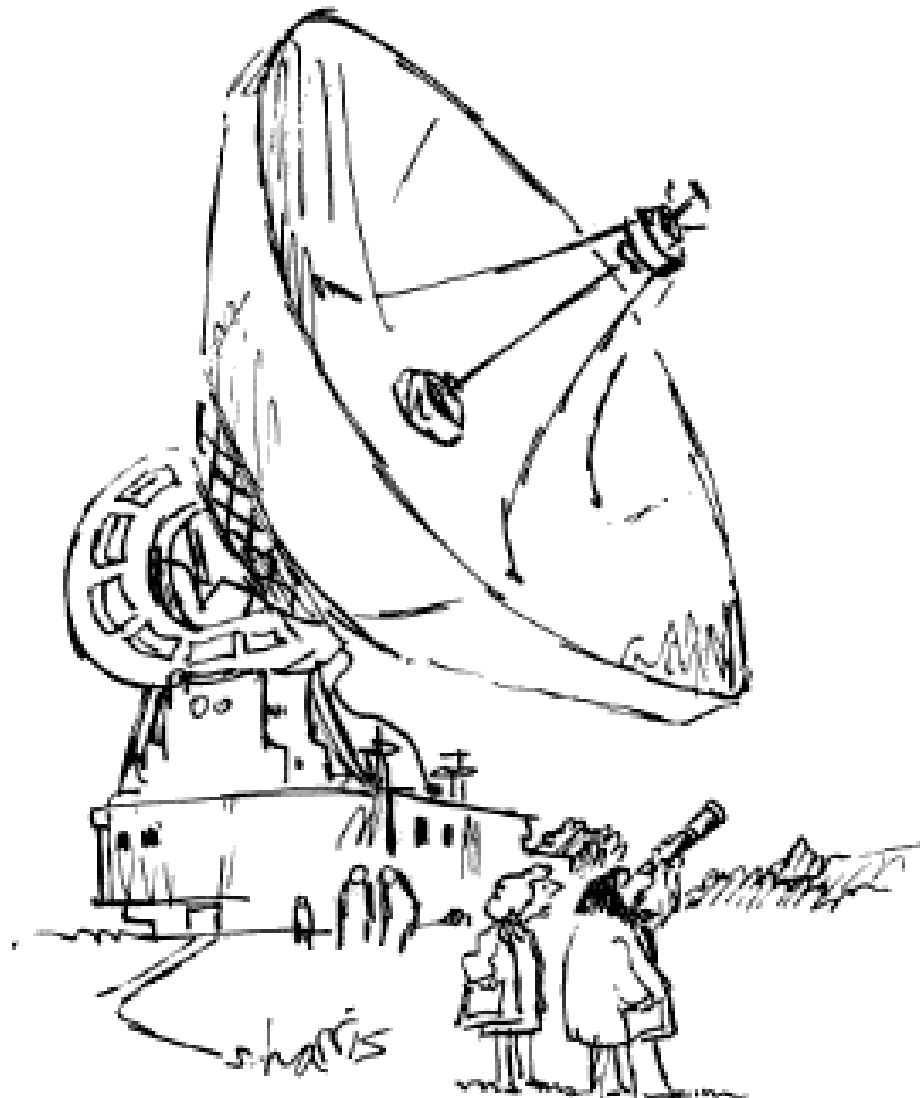
Cyclope

vs.



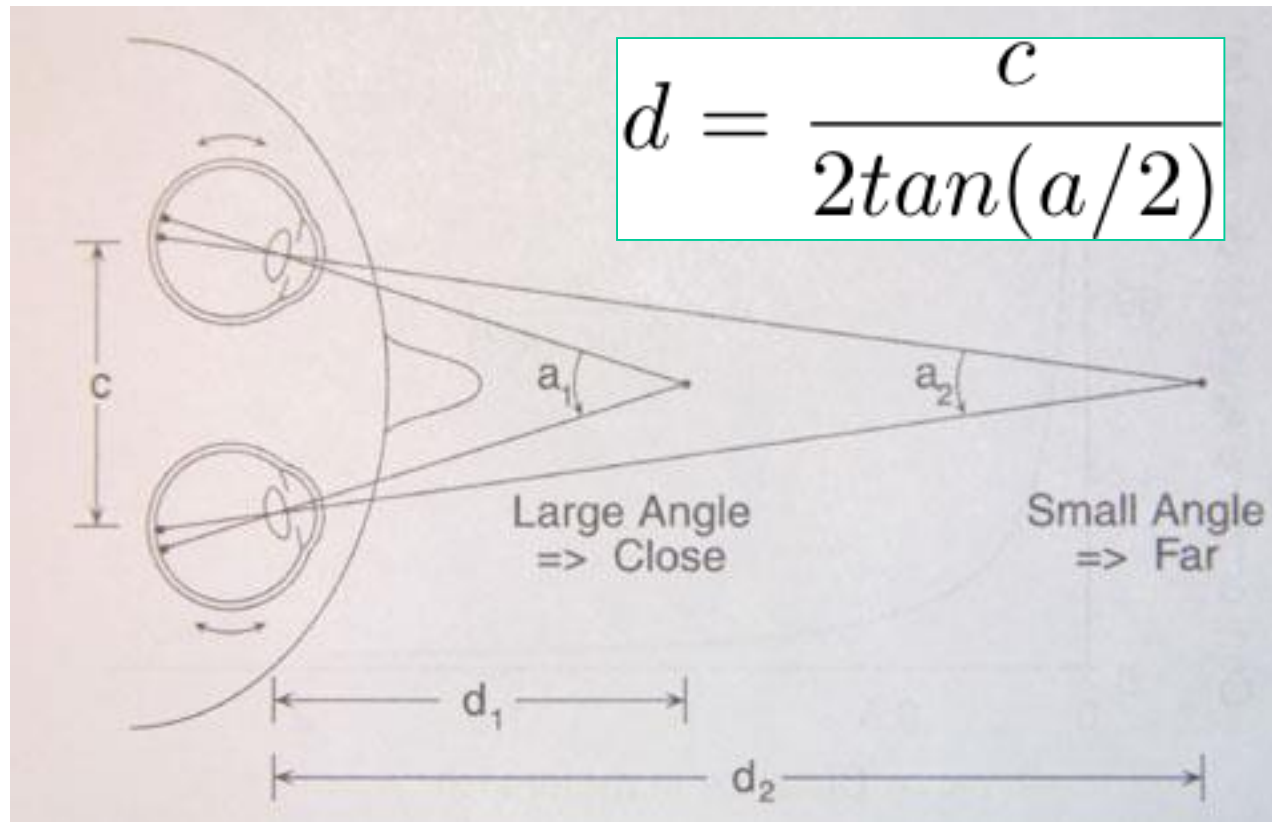
Odysseus

Motivation: two is better than one



"Just checking."

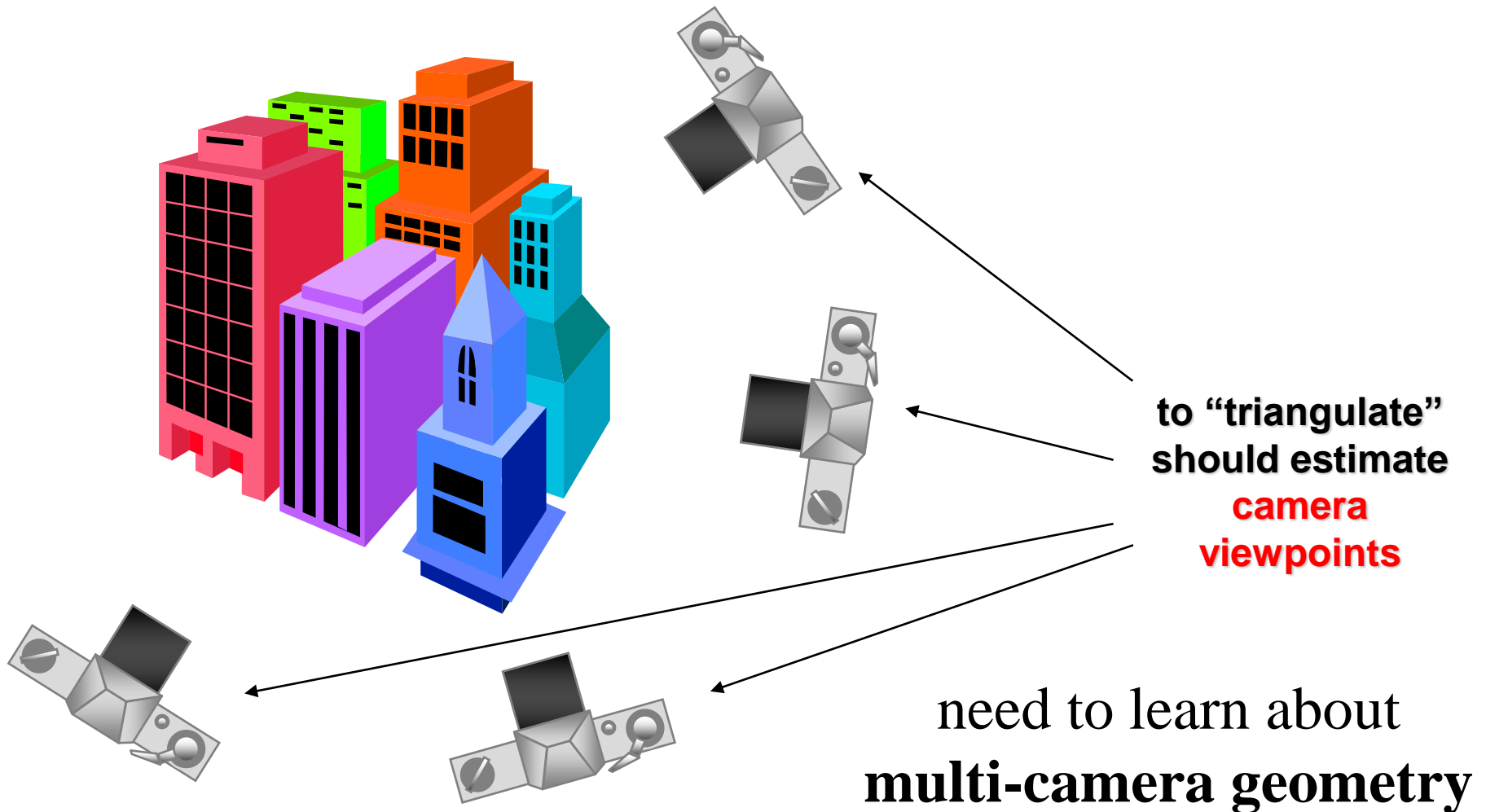
Motivation: **triangulation** gives depth



Human performance: up to 6-8 feet

Motivation: reconstruction problems

Multi-view reconstruction: **shape from two or more images**

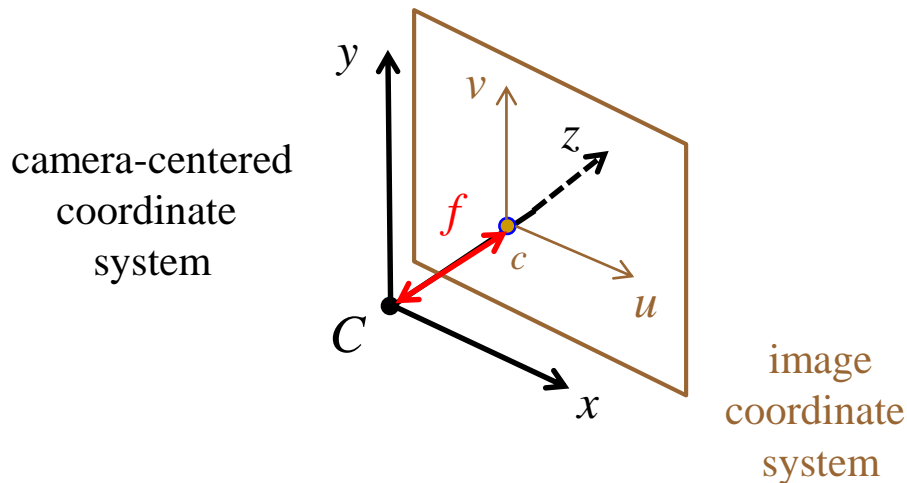


Summary:

- Projective Camera Model
 - intrinsic and extrinsic parameters
 - **projection matrix** (a.k.a. camera matrix)
 - camera calibration (from known 3D points)
 - resection problem
 - estimating intrinsic/extrinsic parameters
- Two cameras (**epipolar geometry**)
 - essential and fundamental matrices: E and F
 - estimating E (from matched features)
 - computing projection matrices from E
- **Structure-from-Motion (SfM)** problem - quick overview
 - at the same time (both are unknown) $\left[\begin{array}{l} \bullet \text{ estimating “motion”: camera positions (projection matrices)} \\ \bullet \text{ estimating “structure”: scene points in 3D space} \end{array} \right.$

Towards projective camera model

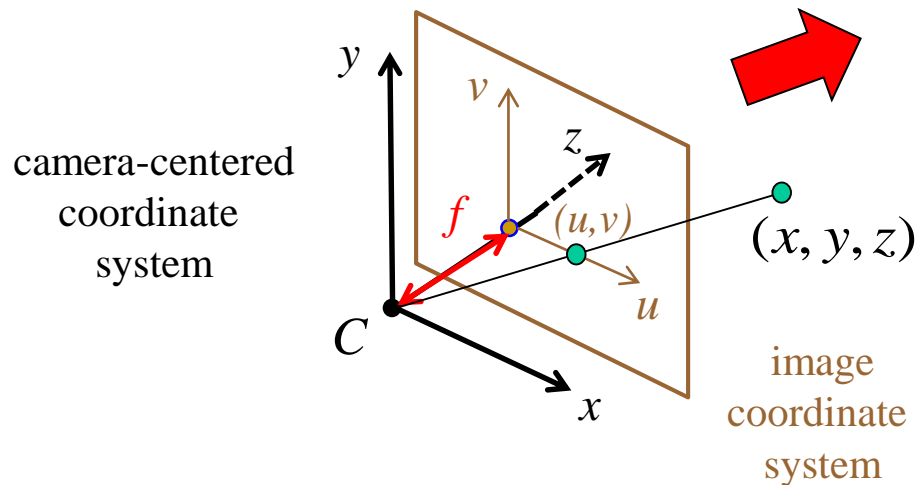
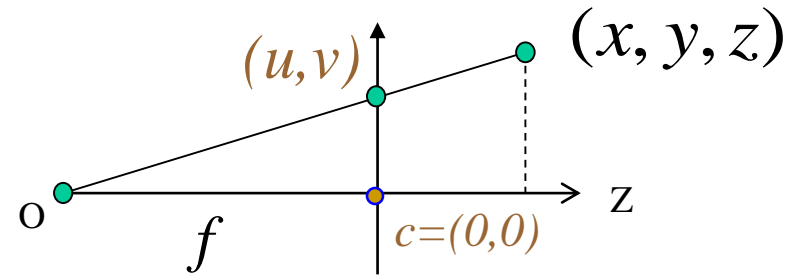
First, if there is only one camera, can use a **camera-centered 3D coordinate system** (x,y,z) :



- optical center is point $(0,0,0)$
- x and y axis are parallel to the image plane
- x and y axis parallel to u and v axis of the image coordinate system
- optical axis (z) intersects image plane at image point $c = (0,0)$

Camera-centered coordinate system

For simplicity,
illustration below assumes
world point (x, y, z)
is inside x-z plane



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

$\underbrace{\hspace{1.5cm}}_u \quad \underbrace{\hspace{1.5cm}}_v$

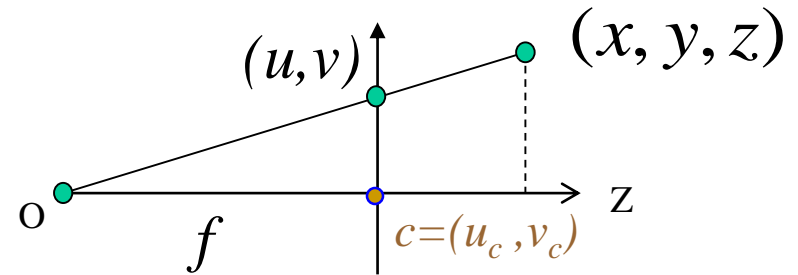
image-based coordinates
of the **projection point**

- optical center is point $(0,0,0)$
- x and y axis are parallel to the image plane
- x and y axis parallel to u and v axis of the image coordinate system
- optical axis (z) intersects image plane at image point $c = (0,0)$

Camera-centered coordinate system

In general, image coordinate center can be anywhere (often in image corner).

Thus, optical axis may intersect image plane at a point with image coordinates $c=(u_c, v_c)$ contributing **additional shift**



camera-centered
coordinate
system

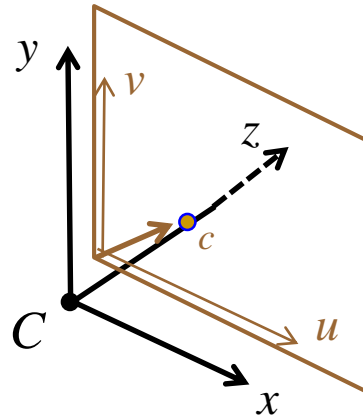


image
coordinate
system

$$(x, y, z) \rightarrow \underbrace{\left(f \frac{x}{z} + u_c\right)}_u, \underbrace{\left(f \frac{y}{z} + v_c\right)}_v$$

image-based coordinates
of the **projection point**

Camera-centered coordinate system

camera projection
can be represented as
matrix multiplication

using **homogeneous representation**
for image points

$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = \begin{bmatrix} f & 0 & u_c \\ 0 & f & v_c \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

K
matrix of intrinsic
camera parameters

$$(x, y, z) \rightarrow \left(\underbrace{f \frac{x}{z} + u_c}_u, \underbrace{f \frac{y}{z} + v_c}_v \right)$$

image-based coordinates
of the **projection point**

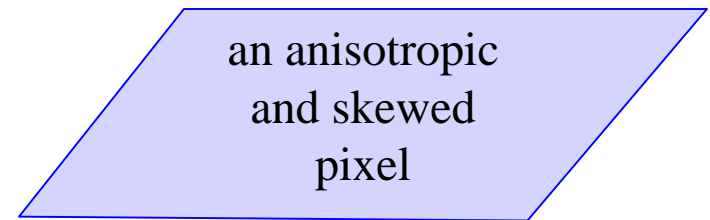
NOTE: $w = z$ (depth)

camera centered coordinates
for 3D world points

Camera-centered coordinate system

Generally, **anisotropic** or **skewed** pixels result in

- different f_x and f_y
- skew coefficient s



using **homogeneous representation**
for image points

$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = \begin{bmatrix} f_x & s & u_c \\ 0 & f_y & v_c \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

K
matrix of intrinsic
camera parameters

s - skew/tilt
 $\frac{f_x}{f_y}$ - aspect ratio

camera centered coordinates
for 3D world points

Camera-centered coordinate system

In general, matrix K of intrinsic camera parameters is **3x3 upper triangular**. It has 5 degrees of freedom.
For square pixels, K has 3 d.o.f.

using **homogeneous representation**
for image points

$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = \underbrace{\begin{bmatrix} f_x & s & u_c \\ 0 & f_y & v_c \\ 0 & 0 & 1 \end{bmatrix}}_K \cdot \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

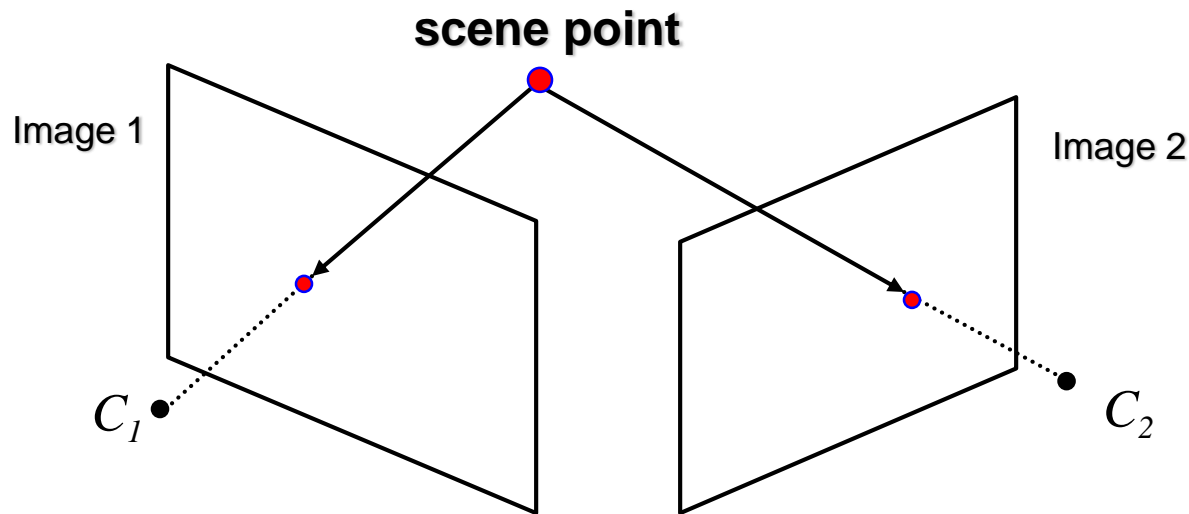
K
matrix of intrinsic
camera parameters

NOTE: here matrix K
maps \mathbb{R}^3 to \mathbb{R}^2 (\mathbb{P}^2)
(**not a homography** $\mathbb{P}^2 \rightarrow \mathbb{P}^2$)

camera centered coordinates
for 3D world points

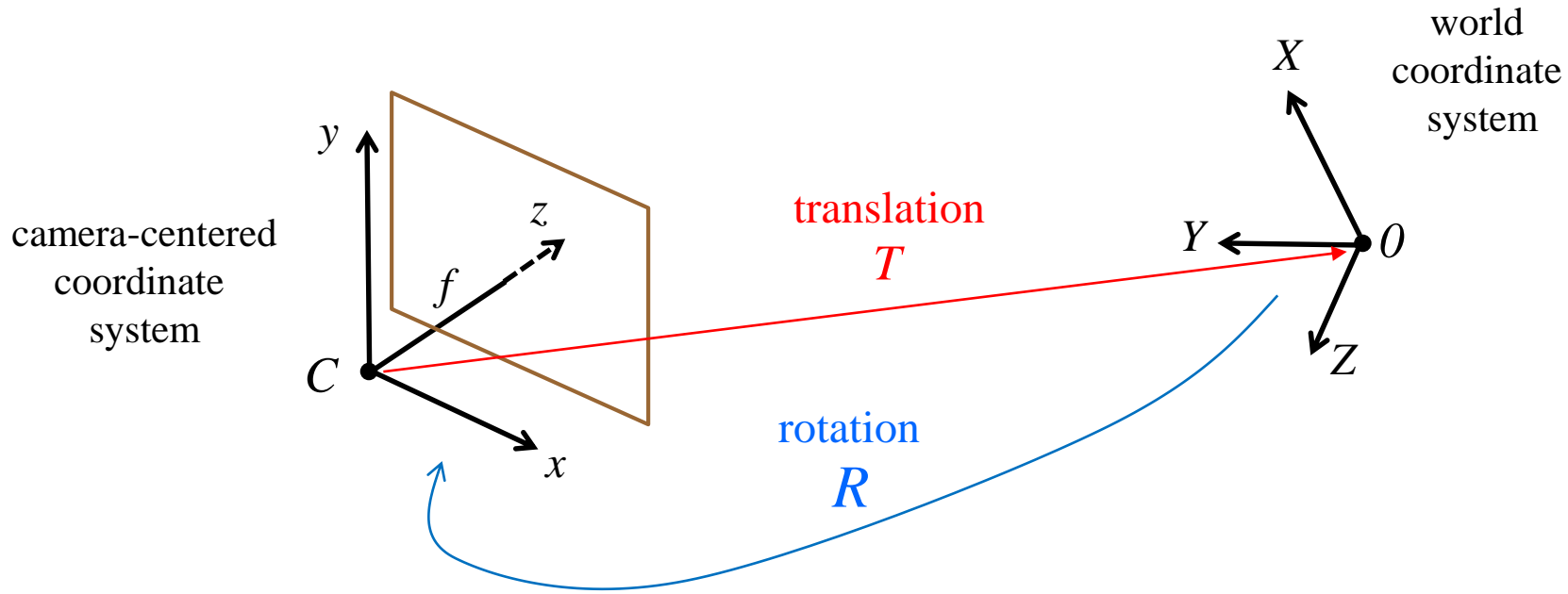
What if there are more than one camera?

Projecting 3D scene onto images with different view-points



Only one camera can serve for world coordinate system.
Other cameras will have their **camera-centered 3D coordinates**
different from the world coordinate system.

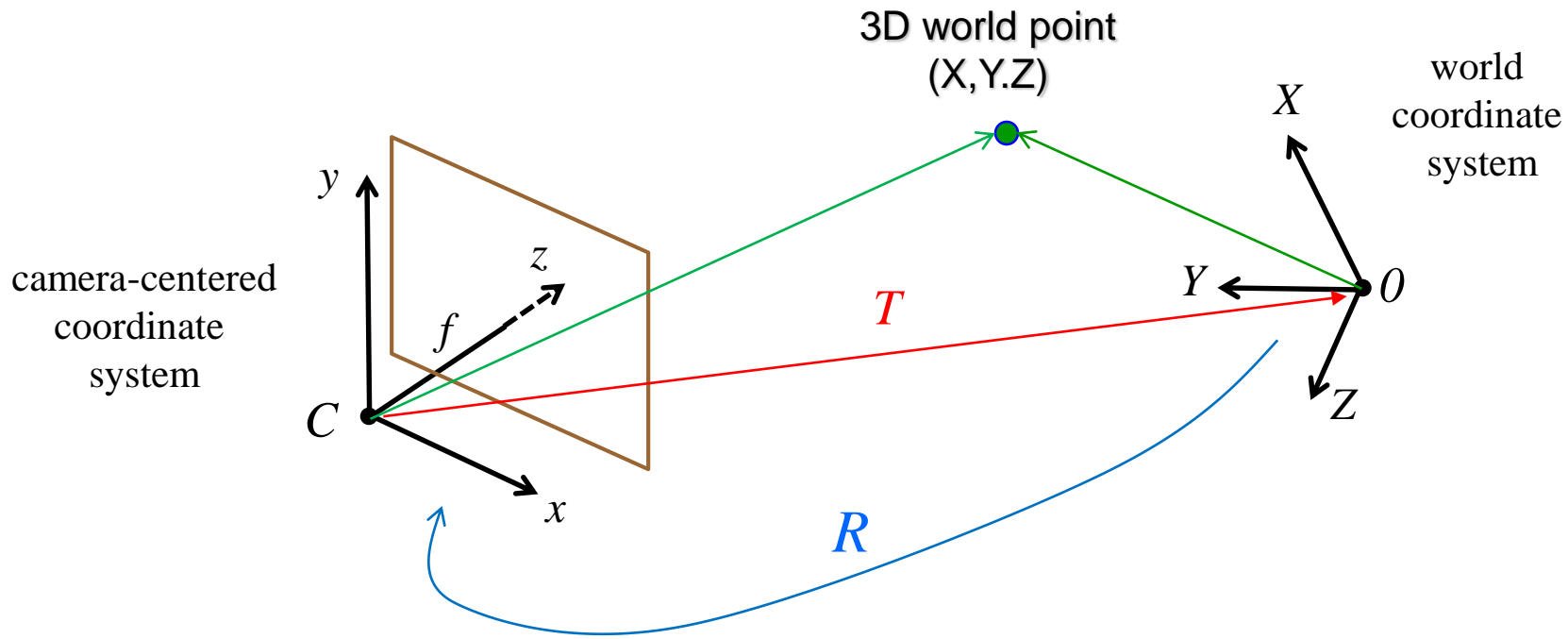
Camera projection matrix



In case of two or more cameras, 3D world coordinate system maybe different from a camera-based coordinate system:

- T is a (**translation**) vector defining relative position of camera's center
- orientation of x, y, z -axis of the camera-based coordinate system can be related to the axis of the world coordinate system via **rotation matrix R**

Camera projection matrix



Converting world coordinates of a point into camera-based 3D coordinate system

using **homogeneous representation** for 3D points in world coordinate system

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = R \cdot \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} + T$$

camera-based 3D coordinates world 3D coordinates

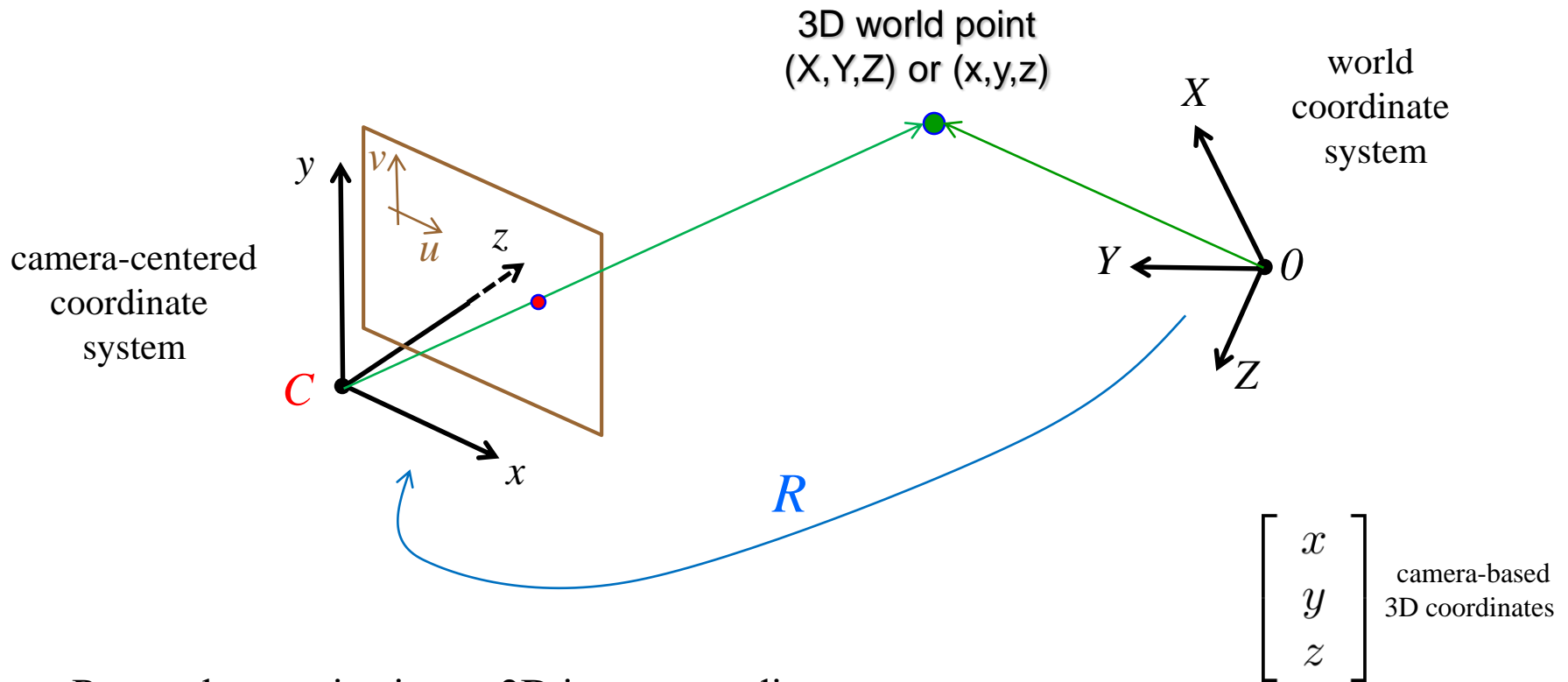
$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \underbrace{\begin{bmatrix} R & T \end{bmatrix}}_{3 \times 4} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}_{4 \times 1}$$

3x1 3x4 4x1

(here vector T is world's center in camera's coordinates)

we get a **linear transformation (matrix multiplication)**

Camera projection matrix



Remember, projecting to 2D image coordinates...

$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = K \cdot \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

homogeneous
image coordinates

camera-based
3D coordinates

\Rightarrow

$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = K \cdot \begin{bmatrix} R & T \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

3x3

3x3

3x4

4x1

project

rotate

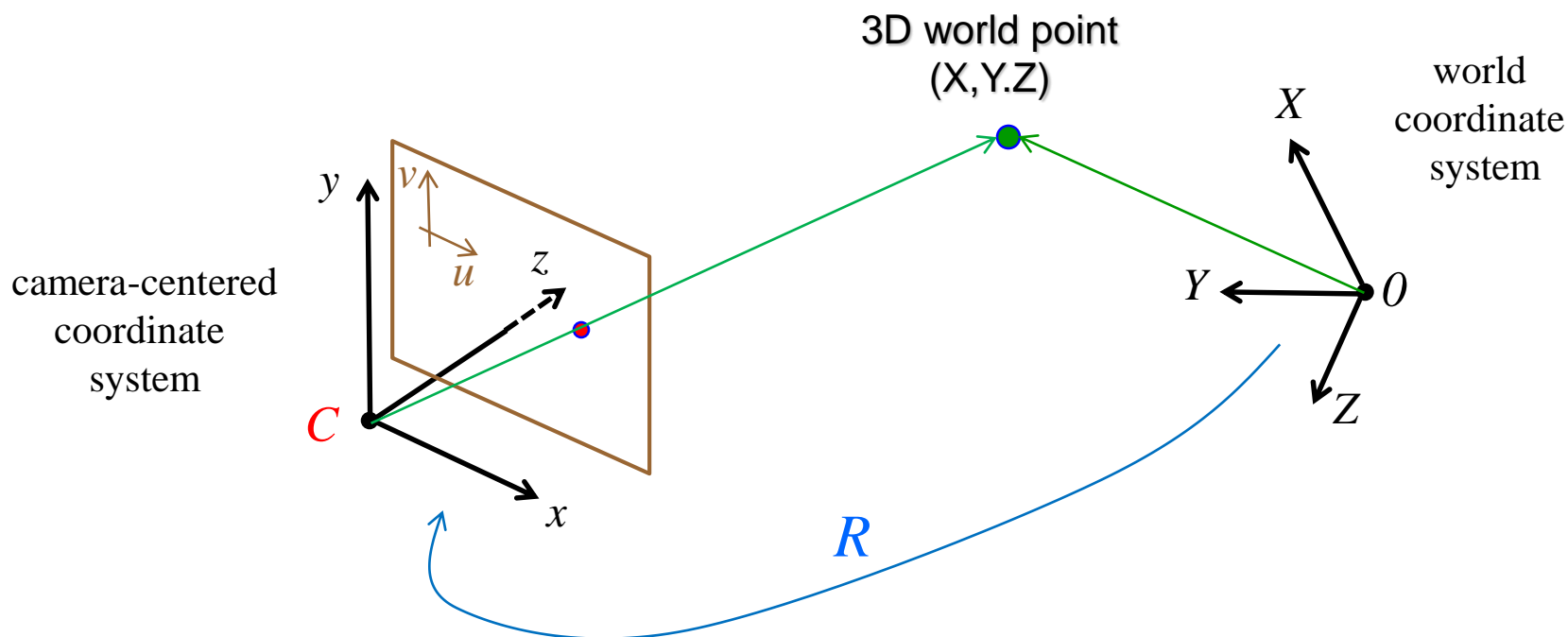
translate

5 d.o.f

3 d.o.f

3 d.o.f

Camera projection matrix



$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = K \cdot \begin{bmatrix} R & T \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \Leftrightarrow \tilde{p} = P \cdot \tilde{X}$$

homogeneous 2D image coordinates

intrinsic camera parameters

extrinsic camera parameters

homogeneous 3D world coordinates

3x1 3x4 4x1

Homogeneous coordinates in 2D and 3D

Trick of adding one more coordinate

- translation becomes matrix multiplication
- 2D points become 3D rays

$$\begin{array}{ccc}
 \text{in } \mathbb{R}^2 & (u, v) \Rightarrow \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \sim \begin{bmatrix} wu \\ wv \\ w \end{bmatrix} & \text{in } \mathbb{P}^2 \\
 & \text{homogeneous 2D image} & \\
 & \text{coordinates} &
 \end{array}
 \qquad
 \begin{array}{ccc}
 \text{in } \mathbb{R}^3 & (X, Y, Z) \Rightarrow \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \sim \begin{bmatrix} wX \\ wY \\ wZ \\ w \end{bmatrix} & \text{in } \mathbb{P}^3 \\
 & \text{homogeneous 3D scene} & \\
 & \text{coordinates} &
 \end{array}$$

Converting *from* homogeneous coordinates

$$\begin{array}{ccc}
 \begin{bmatrix} x \\ y \\ w \end{bmatrix} \Rightarrow (x/w, y/w) & \begin{bmatrix} X \\ Y \\ Z \\ w \end{bmatrix} \Rightarrow (X/w, Y/w, Z/w) \\
 \text{in } \mathbb{P}^2 & \text{in } \mathbb{R}^2 & \text{in } \mathbb{R}^3 \\
 & & \text{in } \mathbb{P}^3
 \end{array}$$

Camera calibration

Goal: estimate intrinsic camera parameters

- focal length f , image center (u_c, v_c) , other elements of **matrix K**
- if needed, corrections for lens distortions (*radial distortion* in fish eye lenses)
not represented by K

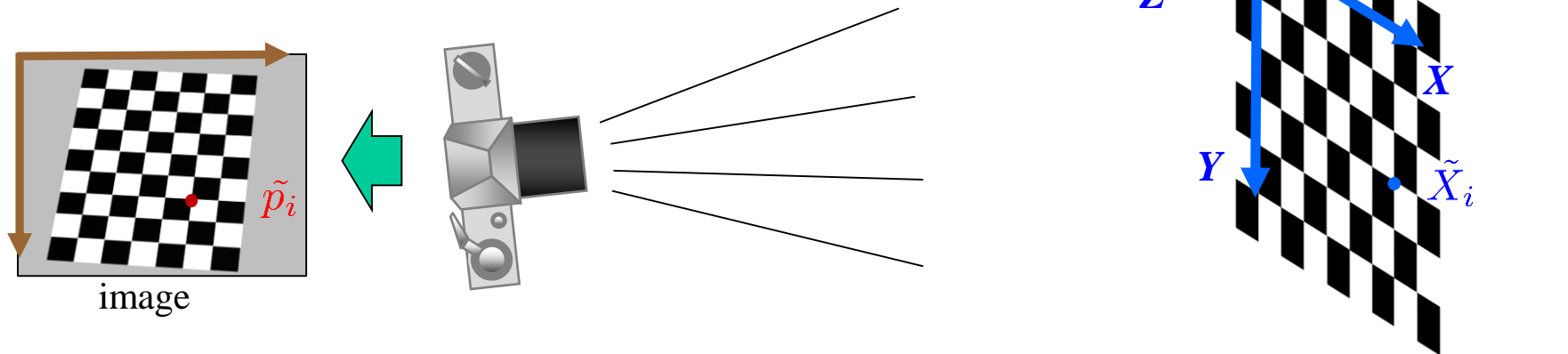
Motivation:

- if K is known, only 6 *d.o.f* remains in projection matrix $P = K \cdot (R/T)$
(3 *d.o.f*. for each rotation R and translation T)
=> it becomes **easier to estimate projection matrices**
corresponding to different viewpoints as camera(s) move around
- using *calibrated* camera(s) is a way to **remove projective ambiguity**
in *structure from motion* 3D reconstruction (more later)

Camera calibration

Basic calibration technique:

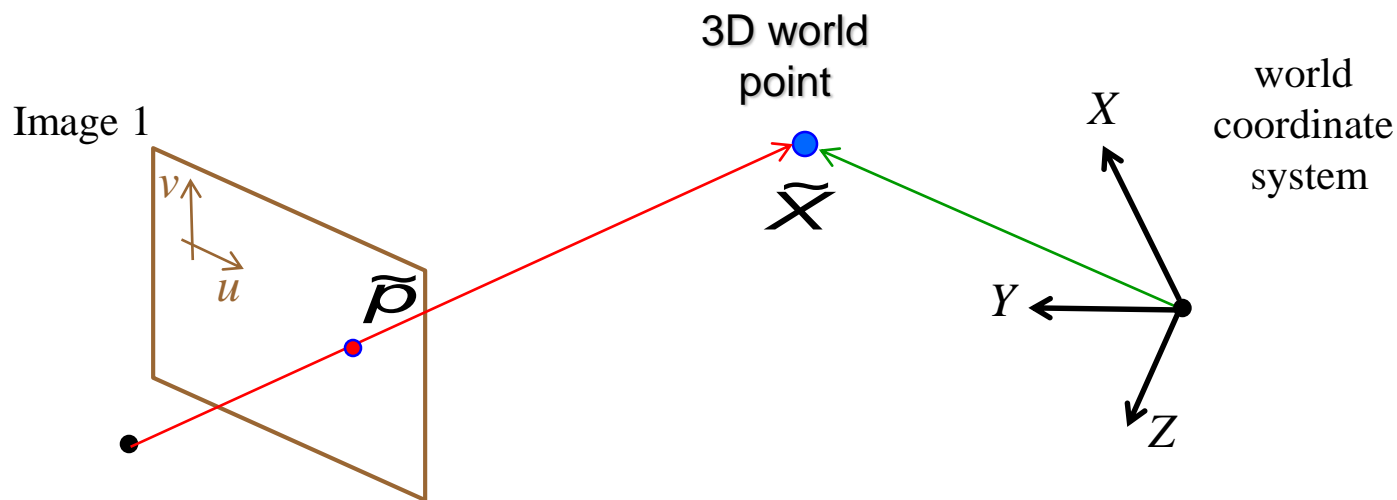
assume a set of 3D points $\{\tilde{X}_i\}$
with known **world coordinates**
and a set of matching **image points** $\{\tilde{p}_i\}$



- find camera matrix P from known matches
(**resection problem**)
- then, find intrinsic and extrinsic parameters
(use **matrix factorization**)

$$\tilde{X}_i \leftrightarrow \tilde{p}_i$$

Camera projection matrix (estimating from $\tilde{X}_i \leftrightarrow \tilde{p}_i$)



P has 12 entries, 11 d.o.f.

Q: How many matched pairs

$\tilde{X}_i \leftrightarrow \tilde{p}_i$
are needed? **A:** 5.5

Q: Solving for a, b, \dots, k, l ?

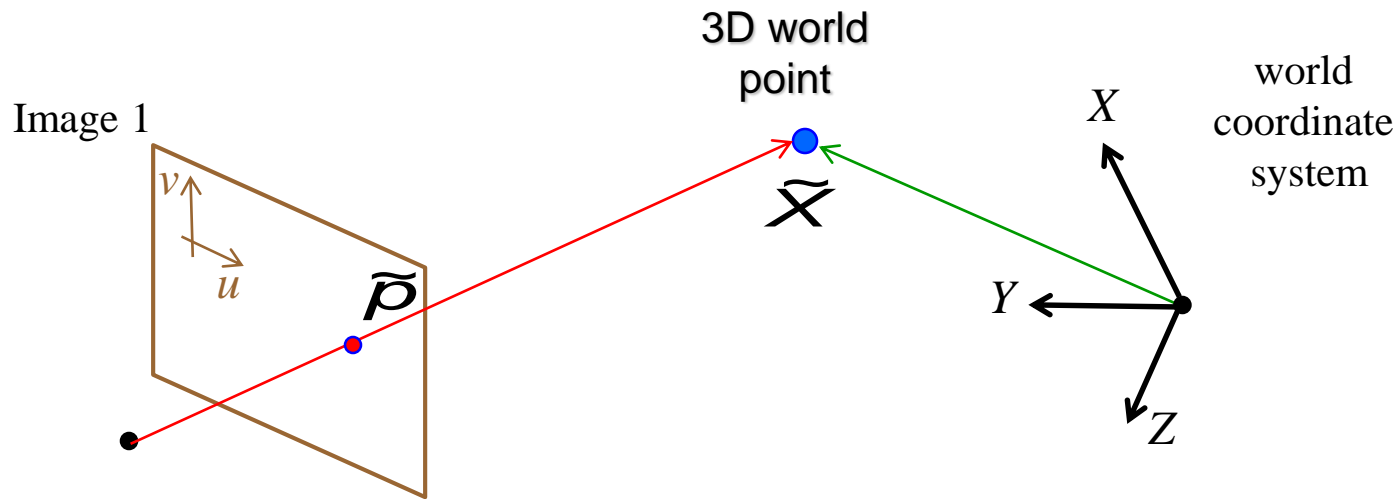
A: similar to estimating
homographies
(see Topic 3, or H&Z p.179)

$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & g & k & l \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

estimate unknown
projection matrix P

(resection problem)

Camera projection matrix (estimating from $\tilde{X}_i \leftrightarrow \tilde{p}_i$)



$$\begin{bmatrix} wu \\ wv \\ w \end{bmatrix} = \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & g & k & l \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

estimate unknown
projection matrix P

(resection problem)

- Use more than 6 matched pairs

$$\tilde{X}_i \leftrightarrow \tilde{p}_i$$

to compensate for errors
(*homogeneous least squares*)

Extracting intrinsic parameters from P

Now, assume that 3×4 projection matrix P is already estimated

$$P = \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & g & k & l \end{bmatrix} = \underbrace{K}_{3 \times 3} \cdot \underbrace{\begin{bmatrix} R & T \end{bmatrix}}_{3 \times 4}$$

known

unknown

How can we get K (as well as R, T) from P ?

Extracting intrinsic parameters from P

$$P = \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & g & k & l \end{bmatrix} \stackrel{?}{=} K \cdot \left[\begin{array}{c|c} R & T \end{array} \right]$$

matrix factorization: H&Z Sec 6.2.4 (p. 163)

Theorem [\mathcal{QR} factorization]: for any $n \times n$ matrix A there is an orthogonal matrix \mathcal{Q} and an upper (or *right*) triangular matrix \mathcal{R} such that $A = \mathcal{R}\mathcal{Q}$.

(If A is invertible and the diagonal elements in \mathcal{R} are chosen positive then the factorization is **unique**.)

$$P = \left[\begin{array}{ccc|c} a & b & c & d \\ e & f & g & h \\ i & g & k & l \end{array} \right] \stackrel{A = \mathcal{R}\mathcal{Q}}{=} \underbrace{\mathcal{R}}_{\textcircled{K}} \cdot \left[\begin{array}{c|c} \underbrace{\mathcal{Q}}_{\textcircled{R}} & \underbrace{\mathcal{R}^{-1}a}_{\textcircled{T}} \end{array} \right]$$

Calibrated Camera (*normalization*)

Once intrinsic parameters K are known

- can “**normalize**” the camera:

switch to a **new image coordinate system** (\tilde{u}, \tilde{v}) defined as

$$\begin{bmatrix} w\tilde{u} \\ w\tilde{v} \\ w \end{bmatrix} = K^{-1} \cdot \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \quad \text{Q: what kind of transform is this for camera's image?}$$

- then, camera's **new projection matrix** \tilde{P} becomes

$$\tilde{P} = K^{-1}P = \cancel{K^{-1}} \cdot K \cdot \left[\begin{array}{c|c} R & T \end{array} \right] = \boxed{\left[\begin{array}{c|c} R & T \end{array} \right]}$$

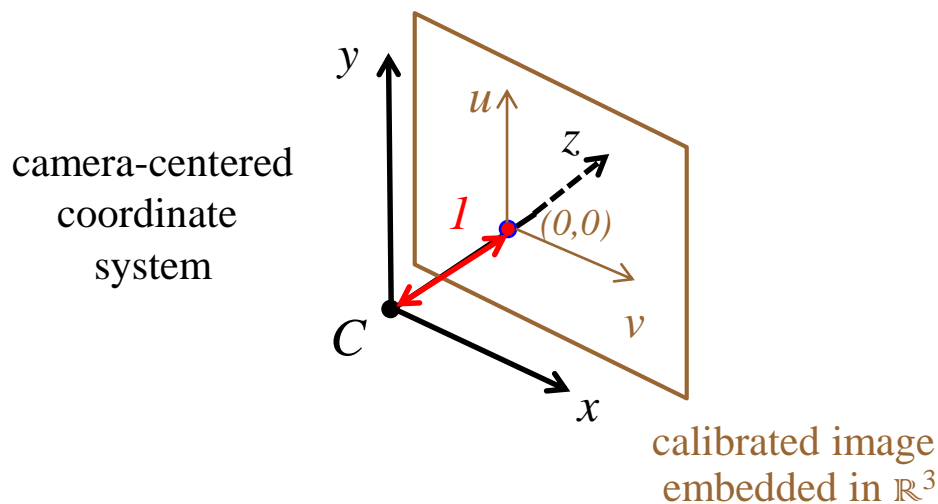
rotation and translation only

Calibrated Camera

After normalization, “effective” intrinsic parameters form an **identity matrix**

$$\left[\begin{array}{c|c} R & T \end{array} \right] = \underbrace{\left[\begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right]}_{\tilde{K} = I} \cdot \left[\begin{array}{c|c} R & T \end{array} \right]$$

extrinsic
parameters



Geometric interpretation:

focal length $f = 1$

point $(0,0)$ = intersection of image plane with optical axis

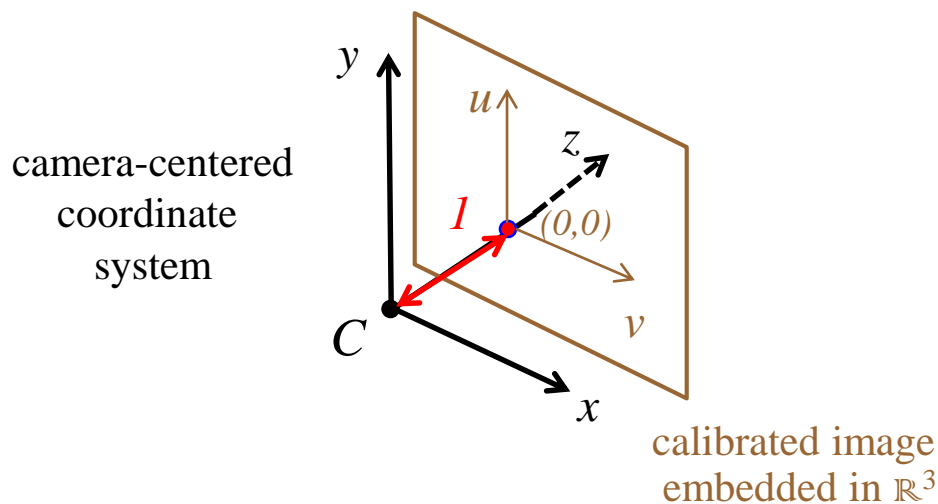
Calibrated Camera

To project onto a **calibrated camera** (a.k.a. normalized camera) one needs only its position (**translation+rotation**) in world coordinates

calibrated/normalized camera's projection matrix

$$P = \left[\begin{array}{c|c} R & T \end{array} \right]$$

still 3x4 matrix
but only 6 d.o.f



Property for normalized camera:
(homogeneous) image coordinates for
any pixel p coincide with this pixel's
camera-centered world coordinates
(treating pixel p as a point in \mathbb{R}^3)

$$p \propto \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} x_p \\ y_p \\ z_p \end{bmatrix}$$

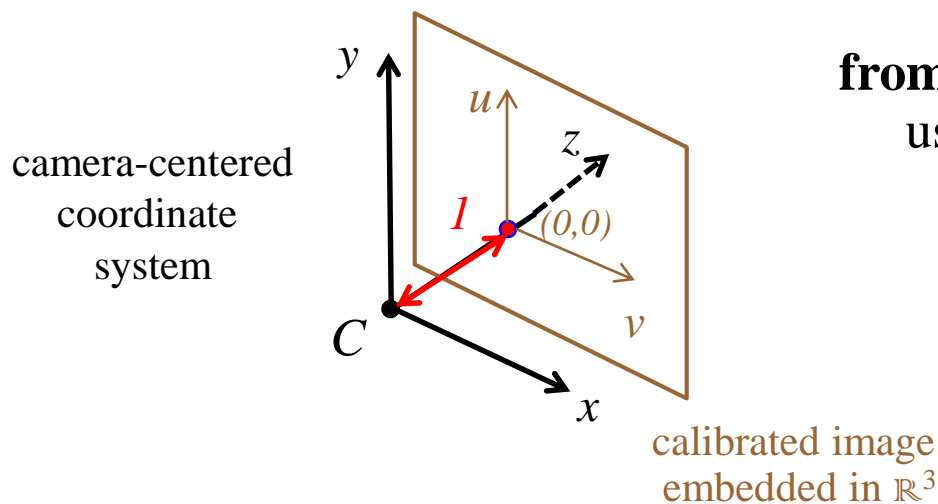
Calibrated Camera

To project onto a **calibrated camera** (a.k.a. normalized camera) one needs only its position (**translation+rotation**) in world coordinates

calibrated/normalized camera's projection matrix

$$P = \left[\begin{array}{c|c} R & T \end{array} \right]$$

still 3x4 matrix
but only 6 d.o.f



from **normalized** back to **original** camera:

use K as a warp $\tilde{p} = Kp$ ($\mathbb{P}^2 \Rightarrow \mathbb{P}^2$)

$\Rightarrow K$ can be interpreted as
a homography mapping calibrated image
embedded in \mathbb{R}^3 to the “digital space”
(i.e. pixels in the original image)

Q: why restrict K to upper triangular ?

hint: $K = QR$

Calibrated Camera

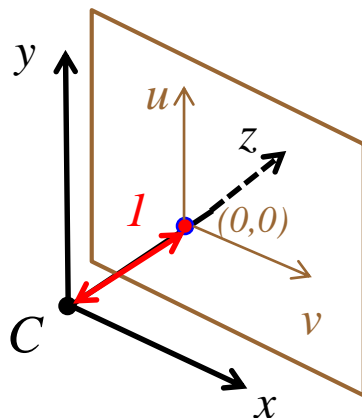
To project onto a **calibrated camera** (a.k.a. normalized camera) one needs only its position (**translation+rotation**) in world coordinates

calibrated/normalized camera's projection matrix

$$P = \left[\begin{array}{c|c} R & T \end{array} \right]$$

still 3x4 matrix
but only 6 d.o.f

camera-centered
coordinate
system



calibrated image
embedded in \mathbb{R}^3

Main point of calibration:
converts any camera to a
“standard” pin hole camera
model shown on the left. After
calibration, images are
independent of how the camera is
made, and depend only on
camera's location/orientation.

Calibrated Camera

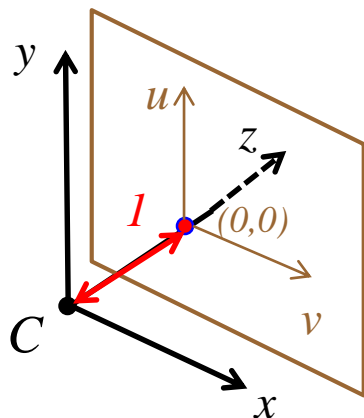
To project onto a **calibrated camera** (a.k.a. normalized camera) one needs only its position (**translation+rotation**) in world coordinates

calibrated/normalized camera's projection matrix

$$P = \left[\begin{array}{c|c} R & T \end{array} \right]$$

still 3x4 matrix
but only 6 d.o.f

camera-centered
coordinate
system



calibrated image
embedded in \mathbb{R}^3

**Estimating multiple viewpoints P_n
is the “motion” part of the
structure-from-motion problem**

NOTE: *camera calibration* uses
known 3D points $\{\tilde{X}_i\}$.
The “**structure**” part of *SfM* problem
estimates unknown 3D scene points.
(later in this topic)

Calibrated Camera

**For simplicity, the rest of this topic assumes
that all images are normalized (calibrated cameras)**

unless explicitly stated otherwise

Two cameras geometry

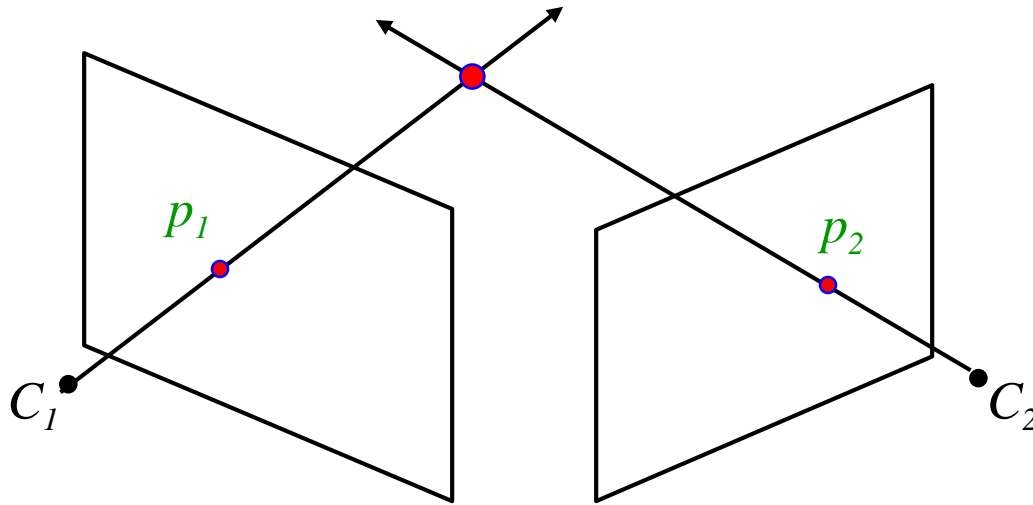
Epipolar geometry

essential & fundamental matrices

Motivation: helps reconstruction

Stereo reconstruction

From 2D images back to 3D scene

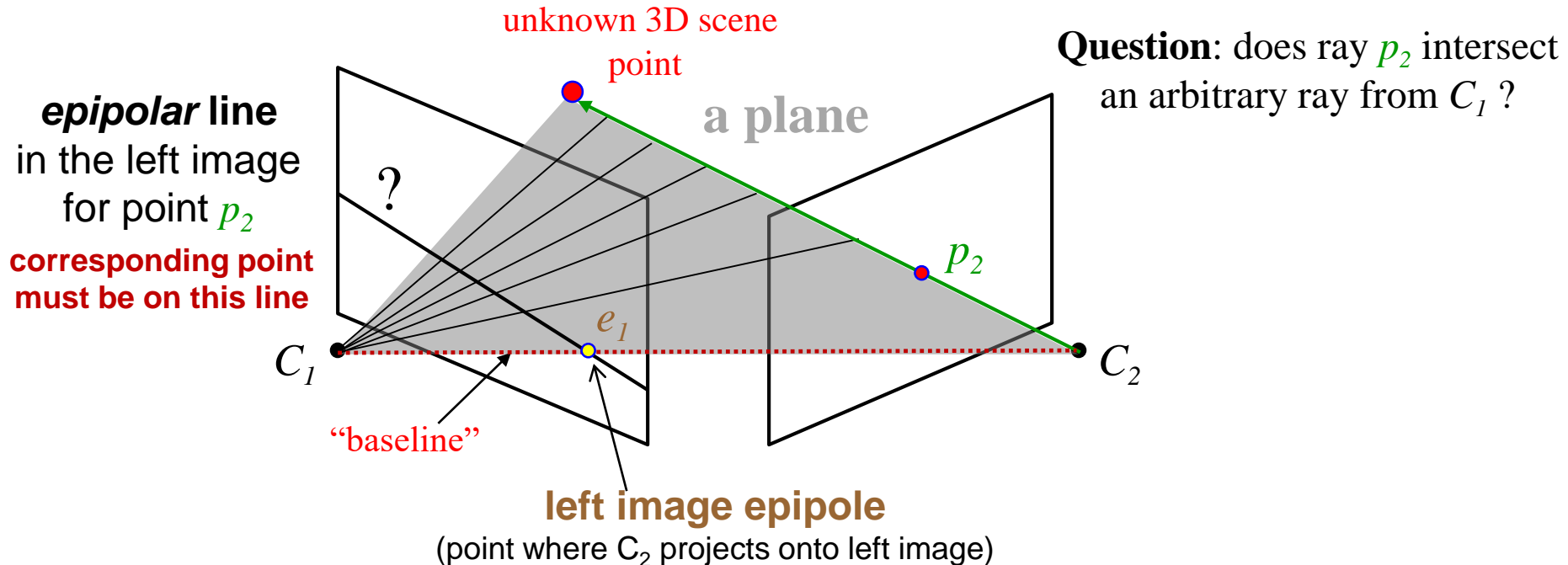


Triangulation: can reconstruct a point as an intersection of two rays, assuming...

- known projection matrix (camera position)
- known **point correspondence**

Epipolar lines

- Find pairs of points that correspond to same scene point
 - not trivial (remember mosaicing)

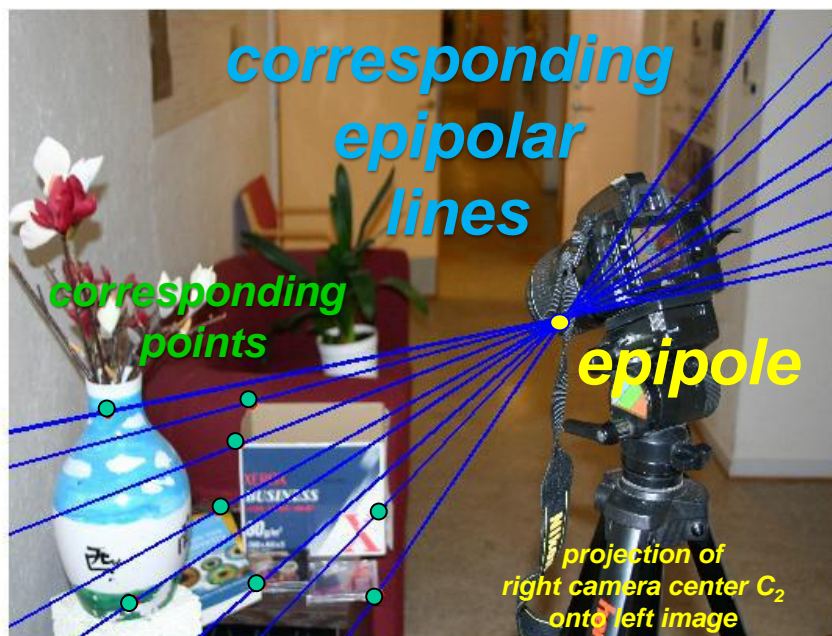


Any right image point p_2 corresponds to some left image **epipolar line**.

It is a projection of **ray** $C_2 \rightarrow p_2$ (ray $C_2 \rightarrow$ **unknown 3D scene point**).

Epipolar lines

Example [from Carl Olsson]
(two stationary cameras)



left camera image
(contains the right camera)

consider some features
in the right image
(projections of some 3D points)



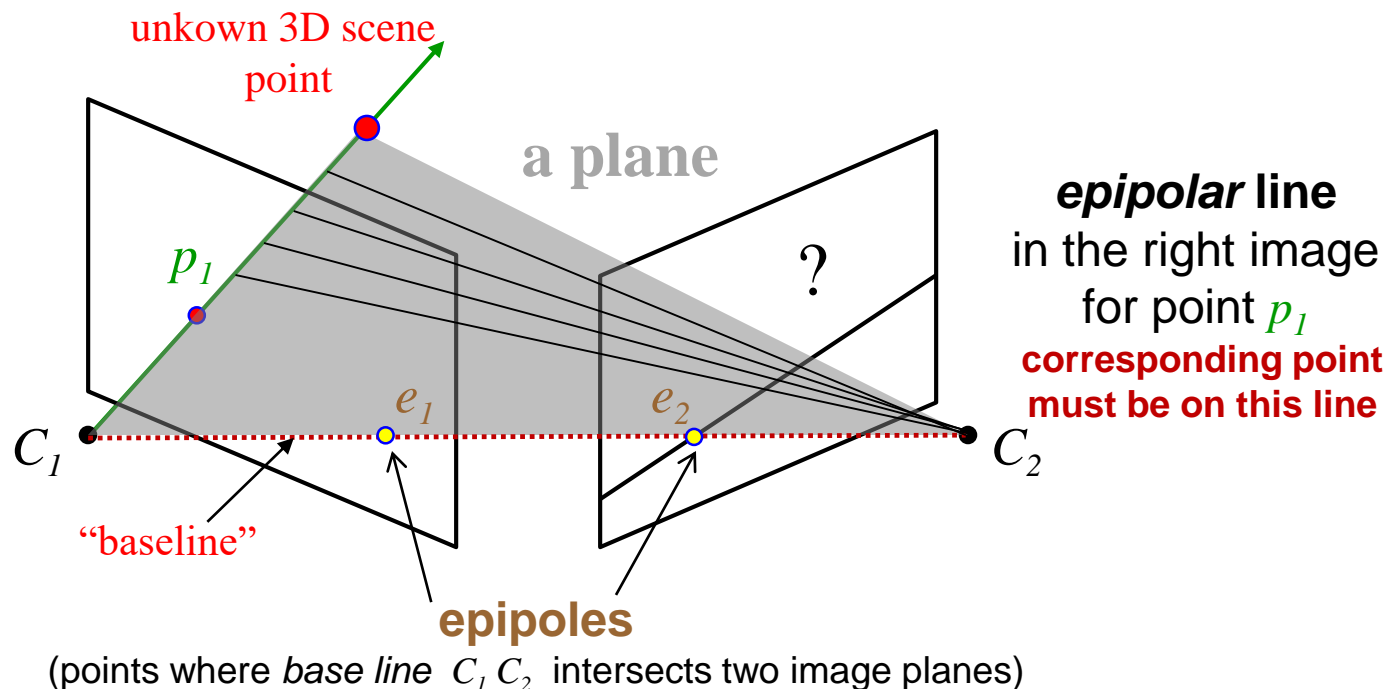
right camera image

Any right image point p_2 corresponds to some left image **epipolar line**.

It is a projection of **ray** $C_2 \rightarrow p_2$ (ray $C_2 \rightarrow$ **unknown 3D scene point**).

Epipolar lines

Similarly, for any given point p_1 in the left image...

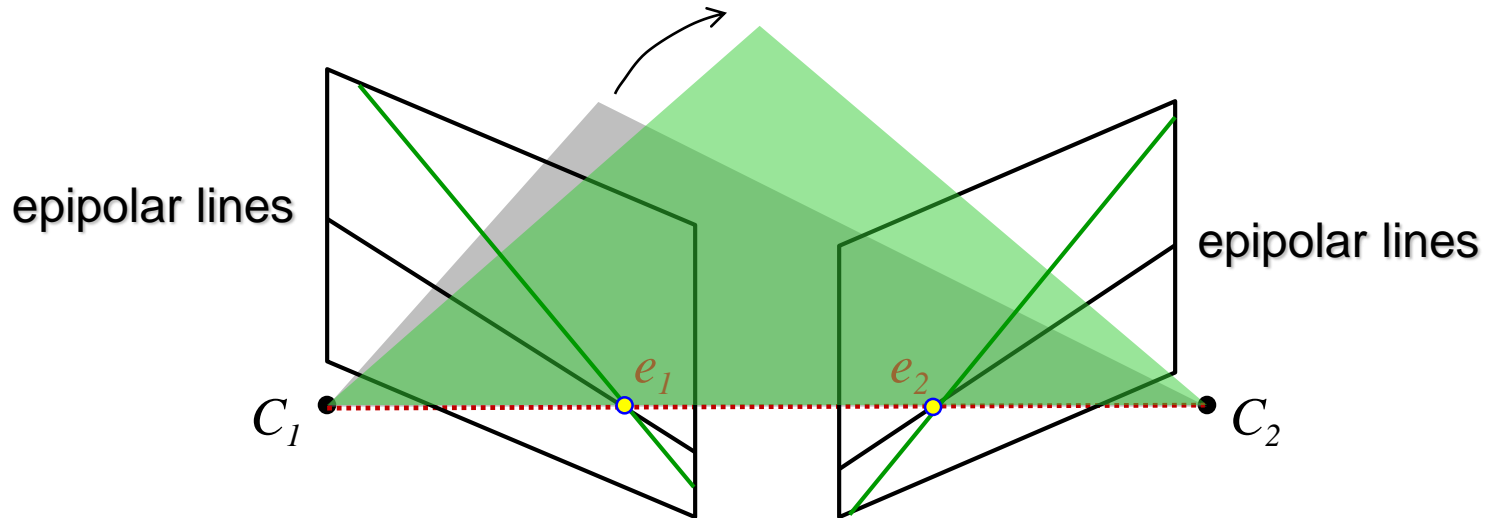


epipolar constraint for the right image: for any point p_1 in the left image, the corresponding point in the right image must be on the line where plane $p_1 C_1 C_2$ intersects the right image (right image *epipolar line*)

- reduces correspondence problem to 1D search along conjugate *epipolar lines*

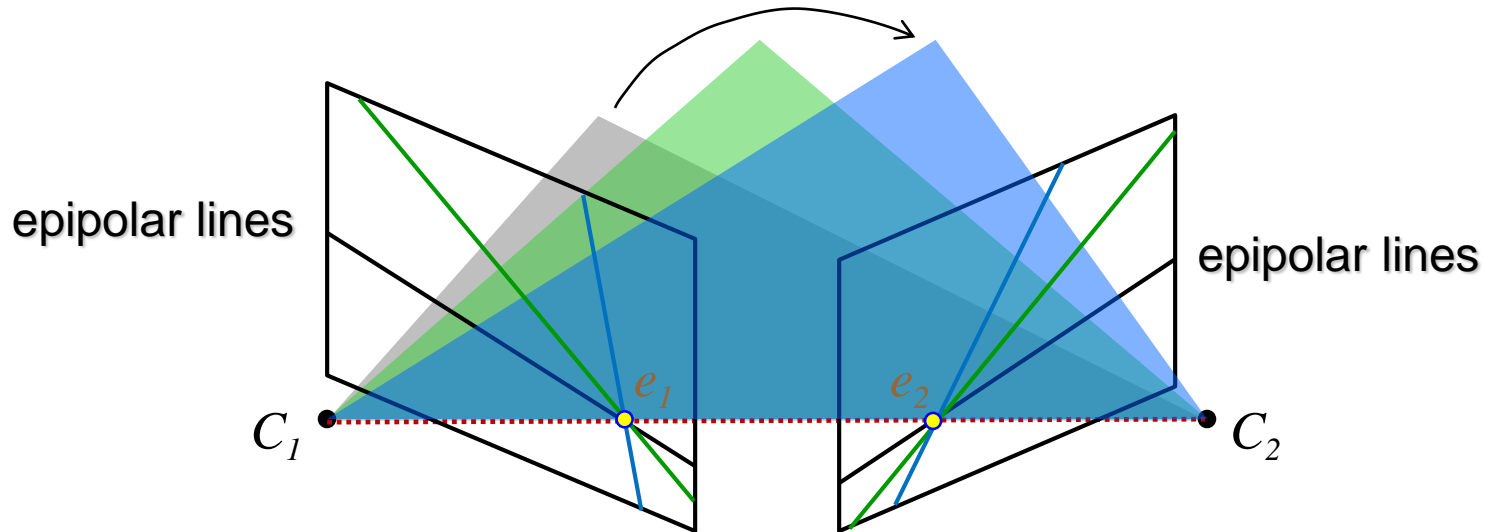
Epipolar lines

System of corresponding epipolar lines depends only on camera set up and it does not depend on 3D scene.



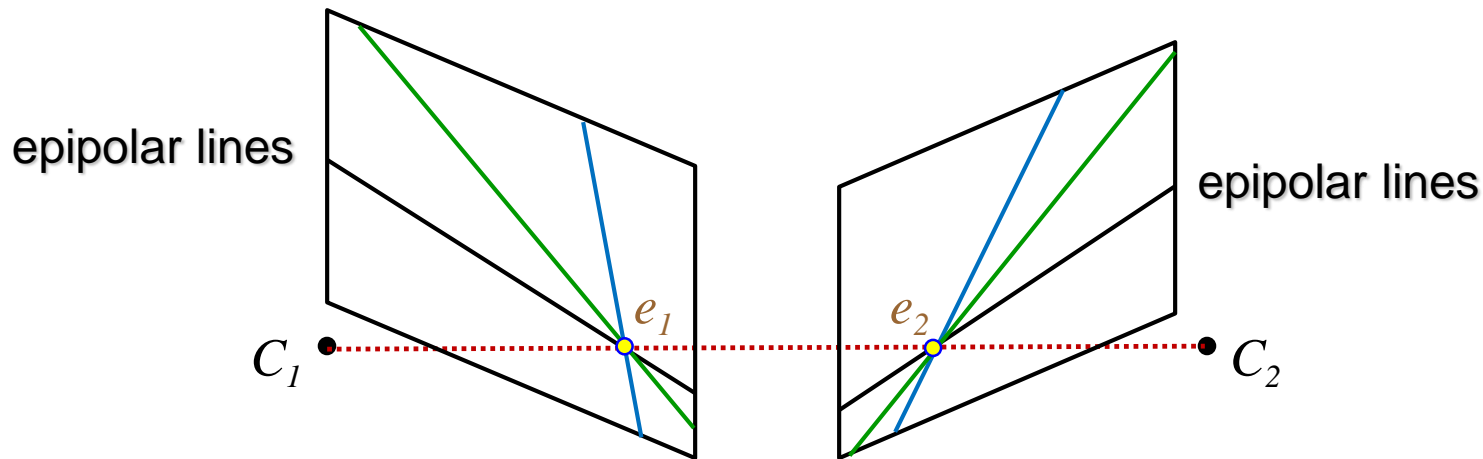
Epipolar lines

System of corresponding epipolar lines depends only on camera set up and it does not depend on 3D scene.



- Intersection of **epipolar planes** (planes containing base line C_1C_2) with image planes define a system of corresponding *epipolar lines*
- Corresponding points can be only on corresponding epipolar lines
 - important to know such lines when searching for corresponding pairs of points

Epipolar lines



- **How can we compute epipolar lines for a given pair of images?**

- if known, camera projection matrices P_1 and P_2 contain all information

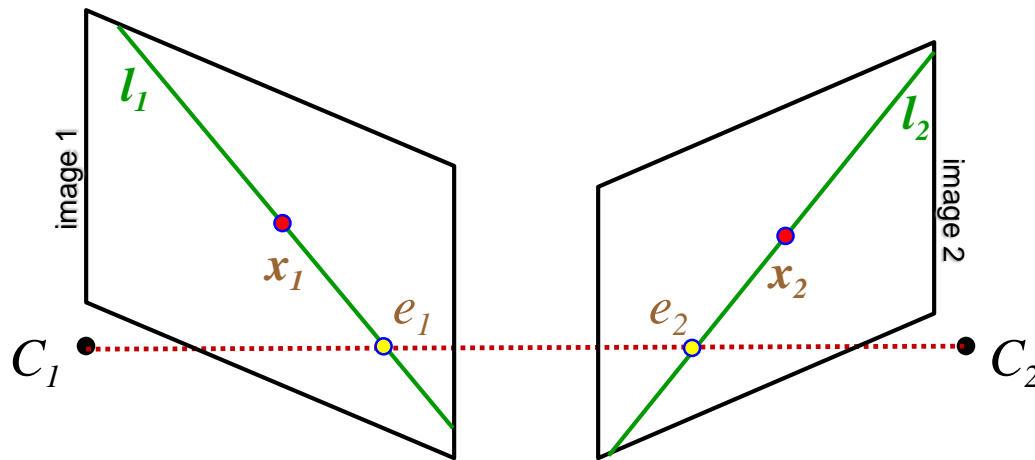
$$e_1 = P_1 C_2 \quad e_2 = P_2 C_1 \quad x_1 = P_1 X \quad x_2 = P_2 X \quad (X - \text{any 3D point})$$

- but only relative position of two cameras really matters:

can estimate a single 3x3 *essential matrix* rather than two 3x4 matrices $P = (R|T) \dots$

Essential matrix E (definition)

The system of corresponding epipolar lines is fully described by a 3x3 matrix E in equation below



3x3 matrix

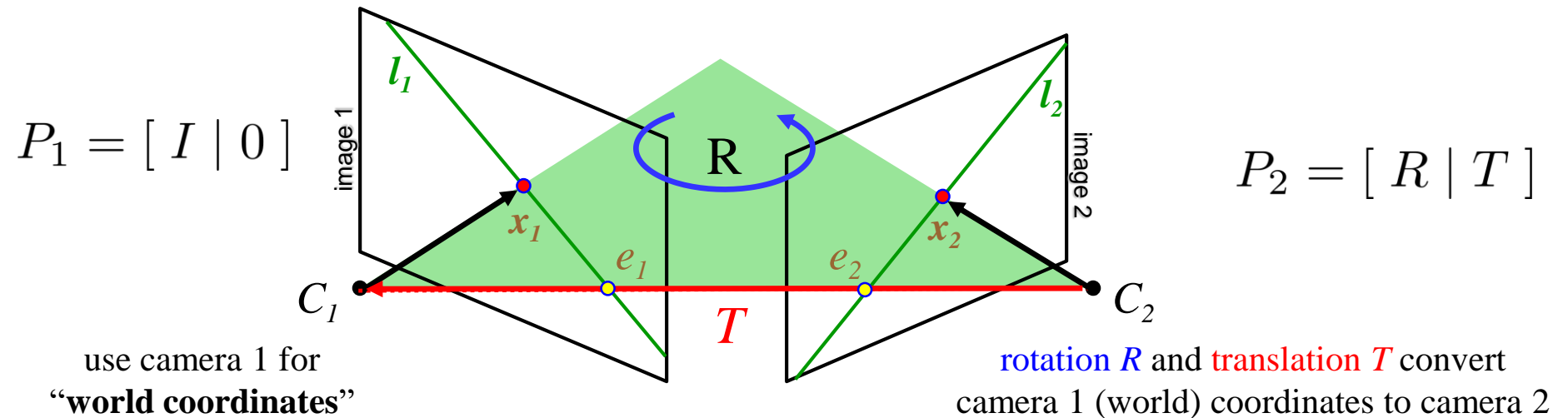
$$\underbrace{x_2^T}_{(l_1)^T} \underbrace{E x_1}_{l_2} = 0$$

for any pair of pixels/points x_1 and x_2
on the corresponding epipolar lines
(assuming calibrated cameras)

NOTE: given x_1 in image 1 vector $l_2 = E x_1$ gives equation $x_2 \cdot l_2 = 0$ (a line in image 2)
given x_2 in image 2 vector $l_1 = E^T x_2$ gives equation $x_1 \cdot l_1 = 0$ (a line in image 1)

Essential matrix E (proof of existence)

Recall: assuming calibrated cameras,
pixels \mathbf{x}_1 and \mathbf{x}_2 in (homogeneous) image coordinates
can be treated as **3D points (vectors)** in the
corresponding camera-centered coordinates of 3D space



dot product cross product

$$\mathbf{x}_2 \cdot [T \times (R\mathbf{x}_1)] = 0$$

for any pair of pixels/points \mathbf{x}_1 and \mathbf{x}_2
on the corresponding epipolar lines
(assuming calibrated cameras)

co-planarity constraint for \mathbf{x}_1 and \mathbf{x}_2
treating \mathbf{x}_1 and \mathbf{x}_2 as vectors in \mathbb{R}^3

NOTE: $R\mathbf{x}_1$ is vector \mathbf{x}_1 in camera 2 coordinates and $T \times R\mathbf{x}_1$ is the green plane's normal (camera 2 coordinates)

Essential matrix E (proof of existence)

NOTE: cross product $a \times b$ can be represented as matrix multiplication

$$a = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \quad \Rightarrow \quad a \times b = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

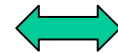
$$a \times b \equiv [a]_{\times} b$$

notation: $[a]_{\times}$

3x3 skew-symmetric matrix
(rank 2)

dot product cross product

$$x_2 \cdot [T]_{\times} (Rx_1) = 0$$



$$x_2^T [T]_{\times} Rx_1 = 0$$

co-planarity constraint for x_1 and x_2
treating x_1 and x_2 as vectors in \mathbb{R}^3

matrix expression

Essential matrix E (proof of existence)

NOTE: due to homogeneous coordinates, scale of E is arbitrary

$$x_2^T E x_1 = 0$$

E is defined by a relative position of two cameras (R and T), as expected

$$E = [T]_{\times} R$$

Q: What is the rank of E ?



essential
matrix

E

$$x_2^T [T]_{\times} R x_1 = 0$$

matrix expression

Essential matrix E

E is defined by a relative position of two cameras (R and T), as expected

$$E = [T]_{\times} R$$

Theorem: 3x3 matrix E is *essential* ($\exists R, T : E = [T]_{\times} R$) if and only if two of its singular values are equal, and the third is zero.

[H&Z. Sec 9.6 p.257]

Then, SVD for essential matrix is
(scale ambiguity allows to use 1 for singular values)

$$E = U \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T$$

Fundamental matrix F

$$x_2^T E x_1 = 0$$

This assumes calibrated camera coordinates

Remember: $\tilde{x} = K^{-1} x$

calibrated
(normalized)
coordinates

original image
coordinates

$$\Rightarrow x_2^T \underbrace{K^{-T} E K^{-1}}_{F} x_1 = 0$$

F - fundamental matrix

$$x_2^T F x_1 = 0$$

defines epipolar lines for
uncalibrated cameras

Essential and Fundamental matrices

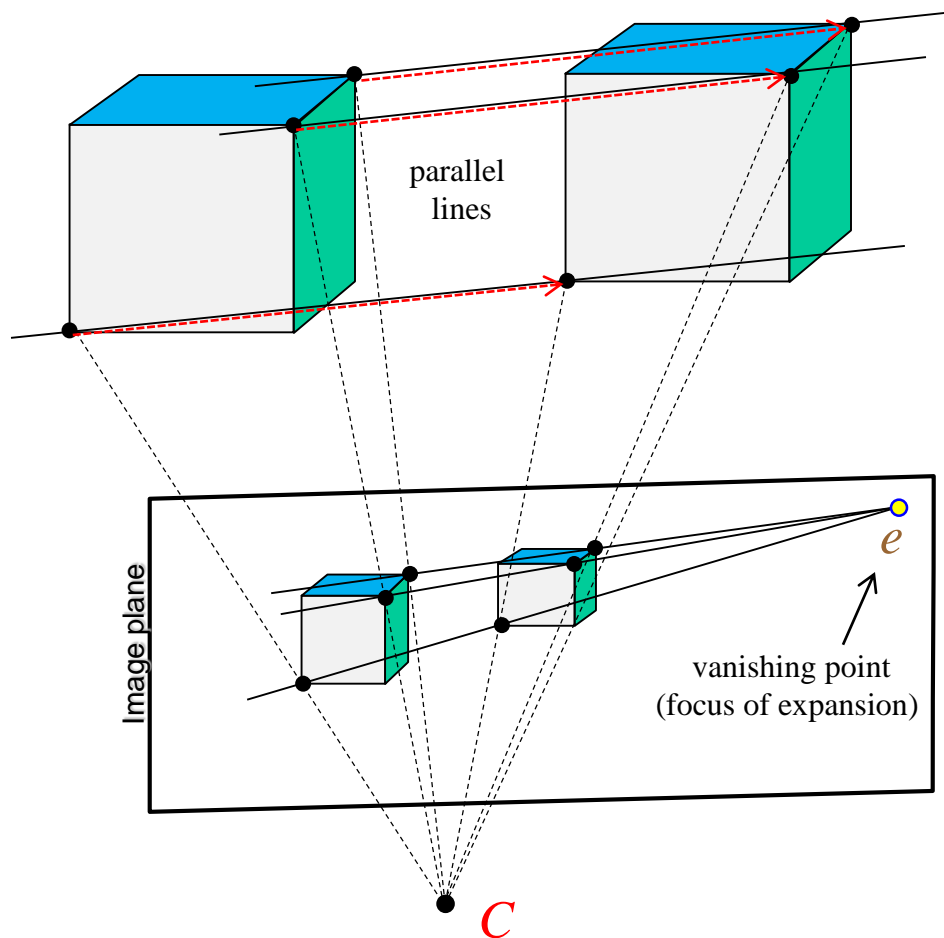
essential matrix E

- epipolar lines $x_2^T E x_1 = 0$
(for two calibrated cameras)
- rank 2 $E = [T]_{\times} R$
- 5 d.o.f (6 d.o.f from R & T , - scale)
- two equal non-zero singular values

fundamental matrix F

- epipolar lines $x_2^T F x_1 = 0$
(for two arbitrary cameras)
- rank 2 $F = K^{-T} E K^{-1}$
- 7 d.o.f (9 par., - scale & $\det F=0$)
- two non-zero singular values

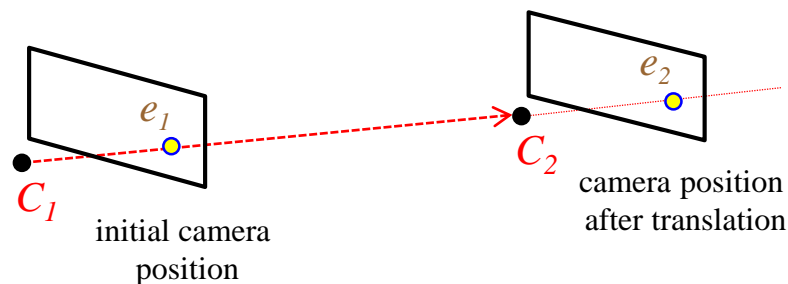
Example: camera translation



Hint: instead of moving the camera, can **equivalently** assume that all 3D scene points translate by vector C_1C_2

objects slide along the epipolar lines

- assume no camera rotation
- vector of camera translation is the same as the base line C_1C_2

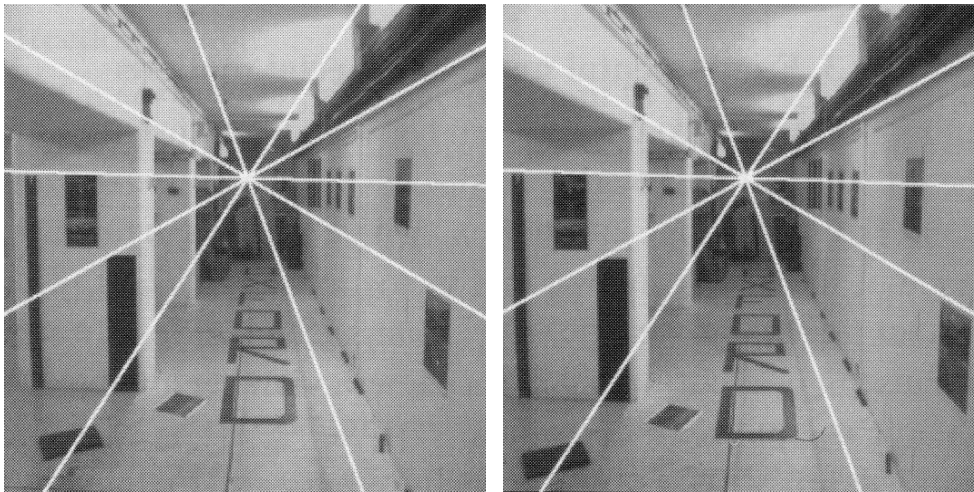


- epipole has the same position in both images $e_1 = e_2 = e$
- the epipole is the projection of the point at infinity for all lines parallel to the base line C_1C_2 (epipole is a *vanishing point* for such lines, also called *focus of expansion* in this case)

- $$e = \begin{bmatrix} a \\ b \\ 1 \end{bmatrix} \Rightarrow E = [e]_{\times} = \begin{bmatrix} 0 & -1 & b \\ 1 & 0 & -a \\ -b & a & 0 \end{bmatrix}$$

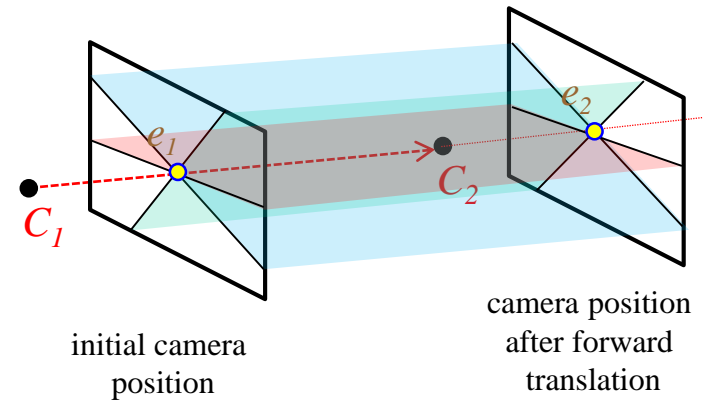
Example: camera translation

Example: forward camera motion



(images from H&Z p.248)

note how objects slide along the epipolar lines



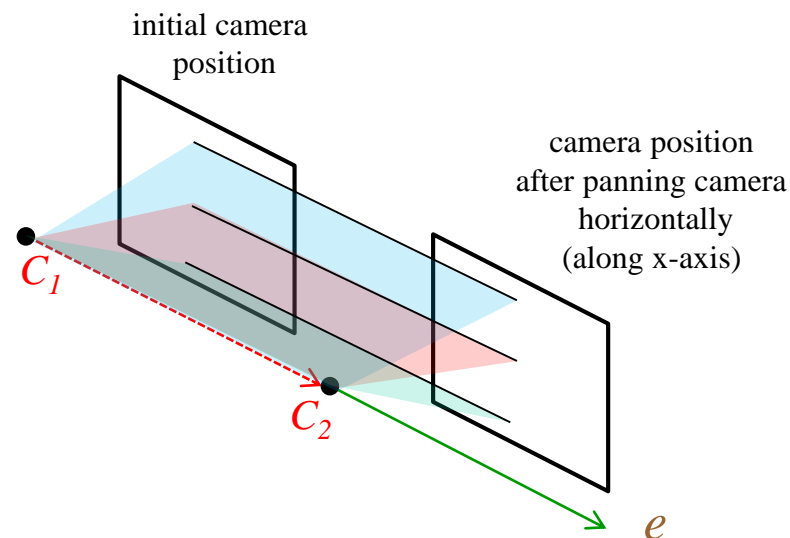
$$e = \begin{bmatrix} a \\ b \\ 1 \end{bmatrix} \Rightarrow E = [e]_{\times} = \begin{bmatrix} 0 & -1 & b \\ 1 & 0 & -a \\ -b & a & 0 \end{bmatrix}$$

Example: camera translation

Example: panning camera motion



note how objects slide along the epipolar lines



- epipole $e_1 = e_2 = e$ is a point at infinity for the image pane
- epipolar lines are parallel lines $y_1 = y_2$

$$e = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \Rightarrow E = [e]_{\times} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix}$$

What's left to cover

- Estimation of E and F
 - simpler **8-point method** (no explicit enforcement of rank or other constraints for E or F)
 - more advanced **5-point method** (see H&Z book, we do not cover this in class)
 - similarly to homography estimation in previous topics, we cover only least squares for *algebraic errors* (*reprojection* errors use more advanced optimization)
- Extraction of cameras (projection matrices) from E
- **Structure from Motion**
 - match, find E , find cameras (estimate pose), **triangulate** (estimate structure)
 - bundle adjustment
 - reconstruction ambiguities

Estimating F or E from $N \geq 8$ matches

8-point method

Assume corresponding points $\mathbf{x}_i \leftrightarrow \bar{\mathbf{x}}_i$ in two images
(matched pair corresponding to a projection of unknown 3D point X_i)

They must lie on the corresponding epipolar lines, thus

$$\bar{\mathbf{x}}_i^T F \mathbf{x}_i = 0 \quad (\text{use } E \text{ for calibrated images})$$

If $\mathbf{x}_i = (x_i, y_i, z_i)$ and $\bar{\mathbf{x}}_i = (\bar{x}_i, \bar{y}_i, \bar{z}_i)$ then

$$\begin{aligned} \bar{\mathbf{x}}_i^T F \mathbf{x}_i = & F_{11}\bar{x}_i x_i + F_{12}\bar{x}_i y_i + F_{13}\bar{x}_i z_i \\ & + F_{21}\bar{y}_i x_i + F_{22}\bar{y}_i y_i + F_{23}\bar{y}_i z_i \\ & + F_{31}\bar{z}_i x_i + F_{32}\bar{z}_i y_i + F_{33}\bar{z}_i z_i = 0 \end{aligned}$$

One matching pair $\mathbf{x}_i \leftrightarrow \bar{\mathbf{x}}_i$ gives **only one linear equation**.

Eight is enough to determine elements of 3x3 matrix F (as scale is arbitrary)

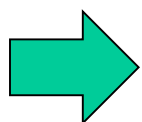
Note: enforcing known properties (e.g. rank=2) allows to use fewer points.

Estimating F or E from $N \geq 8$ matches

In matrix form: one row for each of $N \geq 8$ correspondences

$$\underbrace{\begin{bmatrix} \bar{x}_1 x_1 & \bar{x}_1 y_1 & \bar{x}_1 z_1 & \cdots & \bar{z}_1 z_1 \\ \bar{x}_2 x_2 & \bar{x}_2 y_2 & \bar{x}_2 z_2 & \cdots & \bar{z}_2 z_2 \\ \bar{x}_3 x_3 & \bar{x}_3 y_3 & \bar{x}_3 z_3 & \cdots & \bar{z}_3 z_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \bar{x}_N x_N & \bar{x}_N y_N & \bar{x}_N z_N & \cdots & \bar{z}_N z_N \end{bmatrix}}_{\mathbf{A}} \begin{bmatrix} F_{11} \\ F_{12} \\ F_{13} \\ \vdots \\ F_{33} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

\mathbf{A}
 \mathbf{f}
 $\mathbf{0}$



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

If matched points
have some errors
(not exact locations) ?

Estimating F or E from $N \geq 8$ matches

solve *homogeneous least squares*

$$\min_{\|\mathbf{f}\|=1} \|\mathbf{A}\mathbf{f}\|$$

as in homography estimation,
constraint $\|\mathbf{f}\|=1$ fixes the scale of \mathbf{f} (i.e. F)

$$\begin{bmatrix} E_{11} \\ E_{12} \\ E_{13} \\ \vdots \\ E_{33} \end{bmatrix}$$

for E use \mathbf{e}
instead of \mathbf{f}

Use eigen vector for the smallest eigen value of 9x9 matrix $\mathbf{A}^T \mathbf{A}$

Need $N \geq 8$ to get a unique minimizer \mathbf{f} (up to sign, $-\mathbf{f}$ also works).

If $N = 8$ then perfect (zero) least squares loss is achieved at a unique solution (up to sign).

If $N \leq 7$ then the problem is under-constrained. The (right) null space of \mathbf{A} has dimension ≥ 2 and there are many unit norm solutions \mathbf{f} achieving zero loss.

Estimating F from $N \geq 8$ matches

solve *homogeneous least squares*

$$\min_{\|\mathbf{f}\|=1} \|\mathbf{A}\mathbf{f}\|$$

as in homography estimation,
constraint $\|\mathbf{f}\|=1$ fixes the scale of \mathbf{f} (i.e. F)

Issue: optimal F may not satisfy $\det(F)=0$ and $\text{rank}(F)=2$.

One “solution”: find the “closest” rank 2 matrix \tilde{F} s.t.

$$\min_{\text{rank}(\tilde{F})=2} \|\tilde{F} - F\|$$

where $\|\tilde{F} - F\|$ Frobenius norm $:= \sqrt{\sum_{ij} (\tilde{F}_{ij} - F_{ij})^2}$

Estimating F from $N \geq 8$ matches

Theorem (*low rank approximation*) [Eckart-Young-Mirsky]:

Assuming SVD for $m \times n$ matrix $A = U \operatorname{diag}(s_1, s_2, \dots, s_n) V^T$

$\min_{\operatorname{rank}(\tilde{A})=k} \|\tilde{A} - A\|$ is solved by $\tilde{A} = U \operatorname{diag}(s_1, \dots, s_k, 0, \dots, 0) V^T$
 $\underbrace{\hspace{10em}}$
 k largest singular values of A

the minimizer is unique iff $s_{k+1} \neq s_k$

Issue: optimal F may not satisfy $\det(F)=0$ and $\operatorname{rank}(F)=2$.

One “solution”: find the “closest” rank 2 matrix \tilde{F} s.t.

$$\min_{\operatorname{rank}(\tilde{F})=2} \|\tilde{F} - F\|$$

where $\|\tilde{F} - F\|_{\text{Frobenius norm}} := \sqrt{\sum_{ij} (\tilde{F}_{ij} - F_{ij})^2}$

Estimating F from $N \geq 8$ matches

Theorem (*low rank approximation*) [Eckart-Young-Mirsky]:

Assuming SVD for $m \times n$ matrix $A = U \operatorname{diag}(s_1, s_2, \dots, s_n) V^T$

$$\min_{\operatorname{rank}(\tilde{A})=k} \|\tilde{A} - A\| \text{ is solved by } \tilde{A} = U \operatorname{diag}(s_1, \dots, s_k, 0, \dots, 0) V^T$$

$\underbrace{\hspace{10em}}_{k \text{ largest singular values of } A}$

the minimizer is unique iff $s_{k+1} \neq s_k$

Issue: optimal F may not satisfy $\det(F)=0$ and $\operatorname{rank}(F)=2$.

One “solution”: find the “closest” rank 2 matrix \tilde{F} s.t.

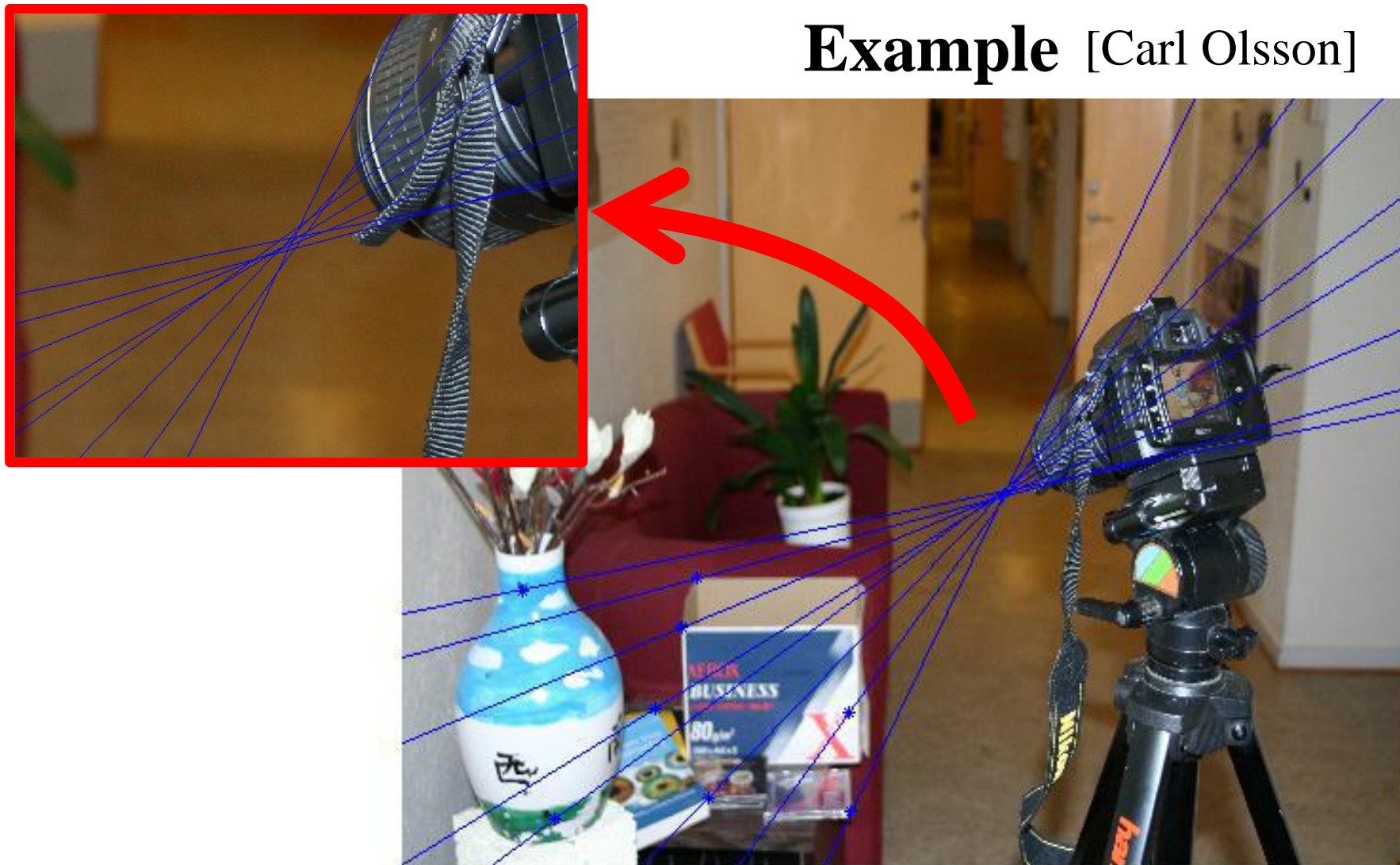
$$\min_{\operatorname{rank}(\tilde{F})=2} \|\tilde{F} - F\|$$



$$\tilde{F} = U \begin{bmatrix} s_1 & 0 & 0 \\ 0 & s_2 & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T$$

Estimating F from $N \geq 8$ matches

Example [Carl Olsson]



Epipole is not well defined if rank is not constrained to 2

Estimating F from $N \geq 8$ matches

YouTube video: search “The Fundamental Matrix Song”



Estimating E from $N \geq 8$ matches

If point matches $\mathbf{x}_i \leftrightarrow \bar{\mathbf{x}}_i$ are in normalized camera images,
solve *homogeneous least squares*

$$\min_{\|\mathbf{e}\|=1} \|\mathbf{A}\mathbf{e}\|$$

$$\mathbf{e} = \begin{bmatrix} E_{11} \\ E_{12} \\ E_{13} \\ \vdots \\ E_{33} \end{bmatrix}$$

as in homography estimation,
constraint $\|\mathbf{e}\|=1$ fixes the scale of \mathbf{e} (i.e. E)

Issue: optimal E may not have SVD $E = U \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T$.

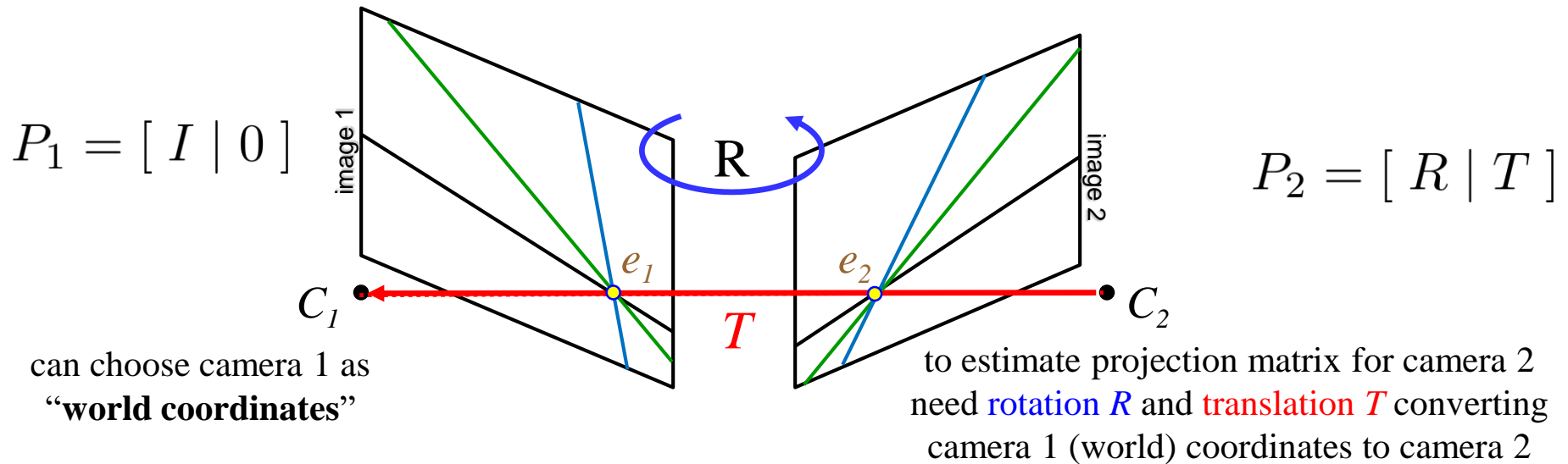
One “solution”: find the “closest” essential matrix \tilde{E}

[H&Z, Sec.11.7.3, p.294]

$$\min_{s_1=s_2, s_3=0} \|\tilde{E} - E\| \rightarrow \tilde{E} = U \begin{bmatrix} \frac{s_1+s_2}{2} & 0 & 0 \\ 0 & \frac{s_1+s_2}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T$$

Extracting cameras from essential matrix E

Now assume essential matrix E is given, need to find P_1 and P_2



Given essential matrix $E = U \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T$

find rotation R and translation T such that $E = [T]_{\times} R$

mathematical formulation of the problem

Extracting cameras from essential matrix E

Four distinct R, T solutions

(up to scale)

Assume SVD decomposition $E = U \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T$

such that $\det(UV^T) = 1$ (if $\det(UV^T) = -1$ switch the sign of the last column in V).

Then, using special matrix $W := \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ we have

$$E = [T]_{\times} R \text{ for any combination of } R = UWV^T \text{ or } UW^T V^T$$

and $T = \pm U_3$ (scale is arbitrary)

see [H&Z:sec 9.6.2, p.258] for proof

↑
the last column of U

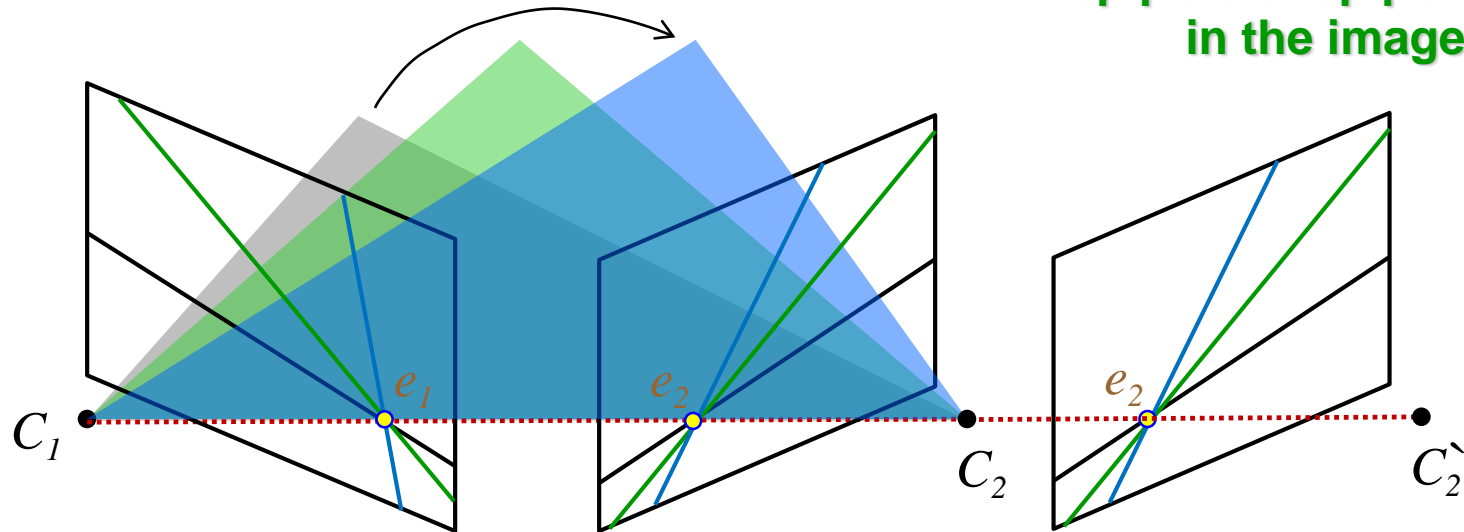
↑
Q: Why?

Extracting cameras from essential matrix E

Four distinct R, T solutions

(up to scale)

baseline length $|T|$
does not change
epipole or epipolar lines
in the images



$$E = [T]_{\times} R \text{ for any combination of } R = UWV^T \text{ or } UW^T V^T$$

and $T = \pm U_3$ (scale is arbitrary)

see [H&Z:sec 9.6.2, p.258] for proof

↑
the last column of U

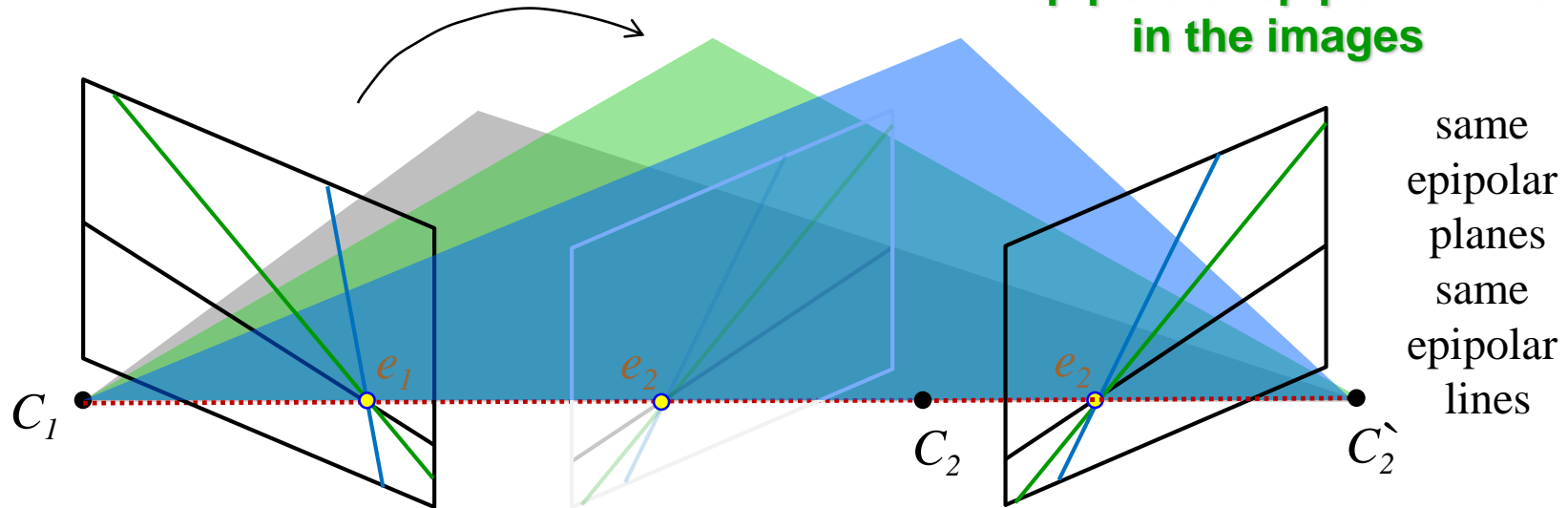
↑
Q: Why?

Extracting cameras from essential matrix E

Four distinct R, T solutions

(up to scale)

baseline length $|T|$
does not change
epipole or epipolar lines
in the images



$$E = [T]_{\times} R \text{ for any combination of } R = UWV^T \text{ or } UW^T V^T$$

and $T = \pm U_3$ (scale is arbitrary)

see [H&Z:sec 9.6.2, p.258] for proof

↑
the last column of U

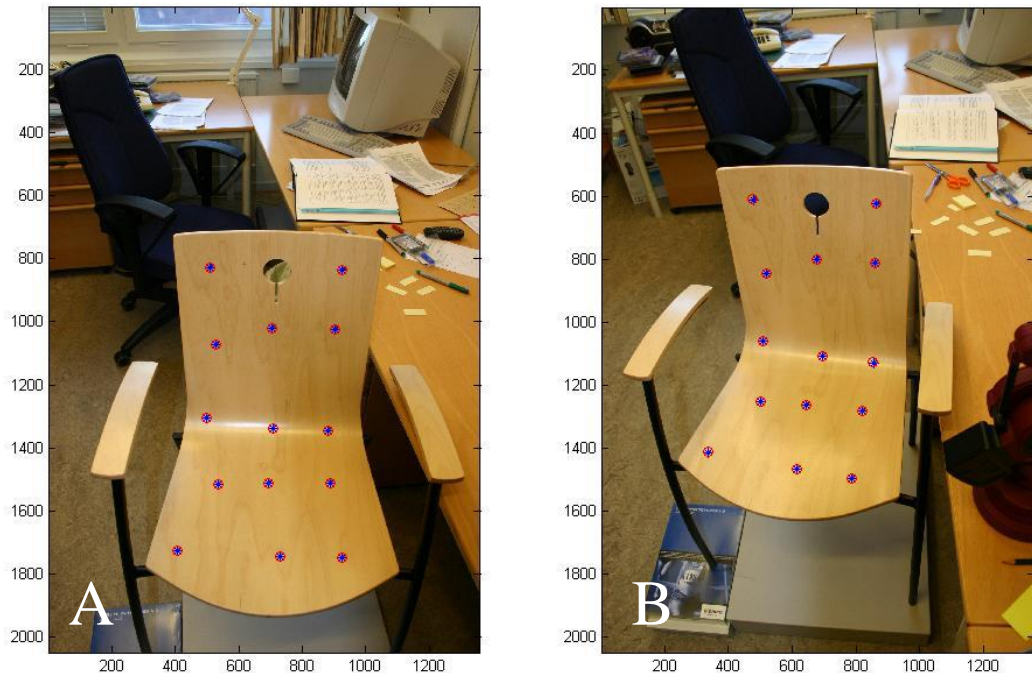
↑
Q: Why?

Extracting cameras from essential matrix E

Four distinct R, T solutions
(up to scale)

Example:
[from Carl Olsson]

Two given views of a chair



14 known correspondences allow to estimate
essential matrix E assuming K is known
(e.g. 8 point method)

Extracting cameras from essential matrix E

Four distinct R, T solutions

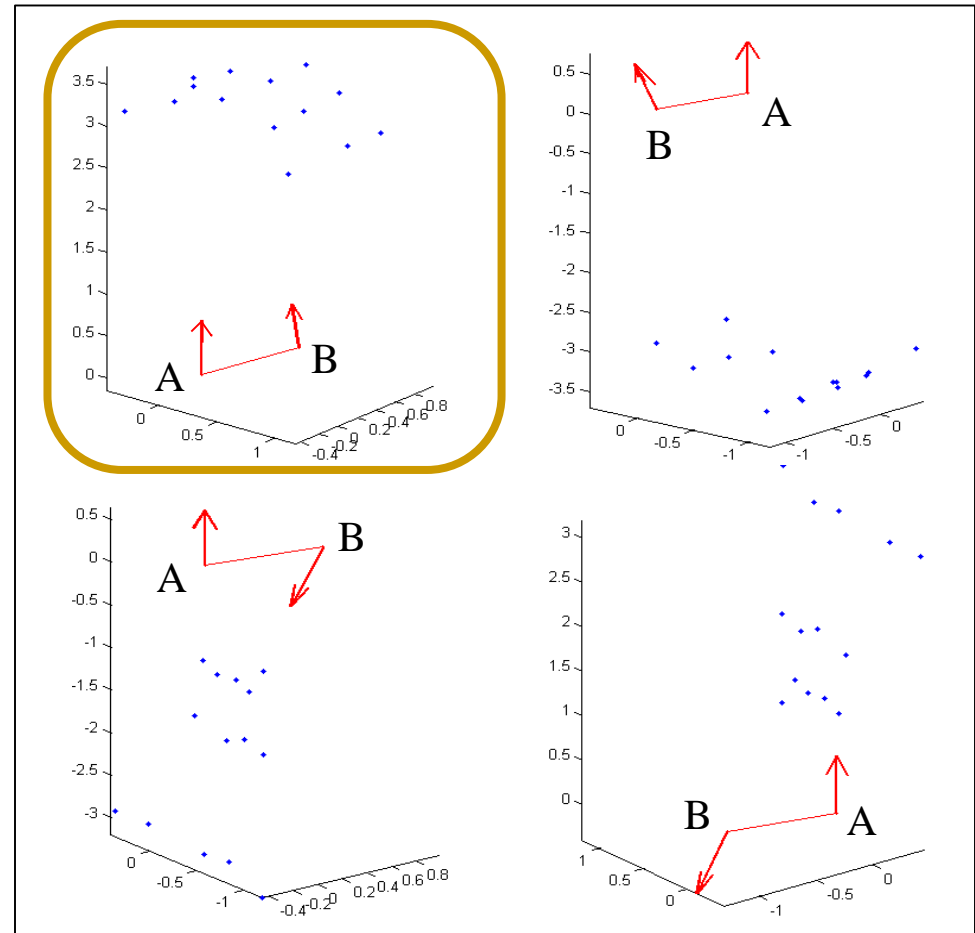
(up to scale)

Example:

[from Carl Olsson]

- four distinct **relative camera positions** (motion R, T) computed from E (up to scale)
- 3D structure $\{X_i\}$ computed from correspondences $\mathbf{x}_i \leftrightarrow \bar{\mathbf{x}}_i$ by *triangulation* (more soon...) up to a *similarity transformation* (i.e. scale+position+orientation)

baseline reversal ($T = \pm U_3$)

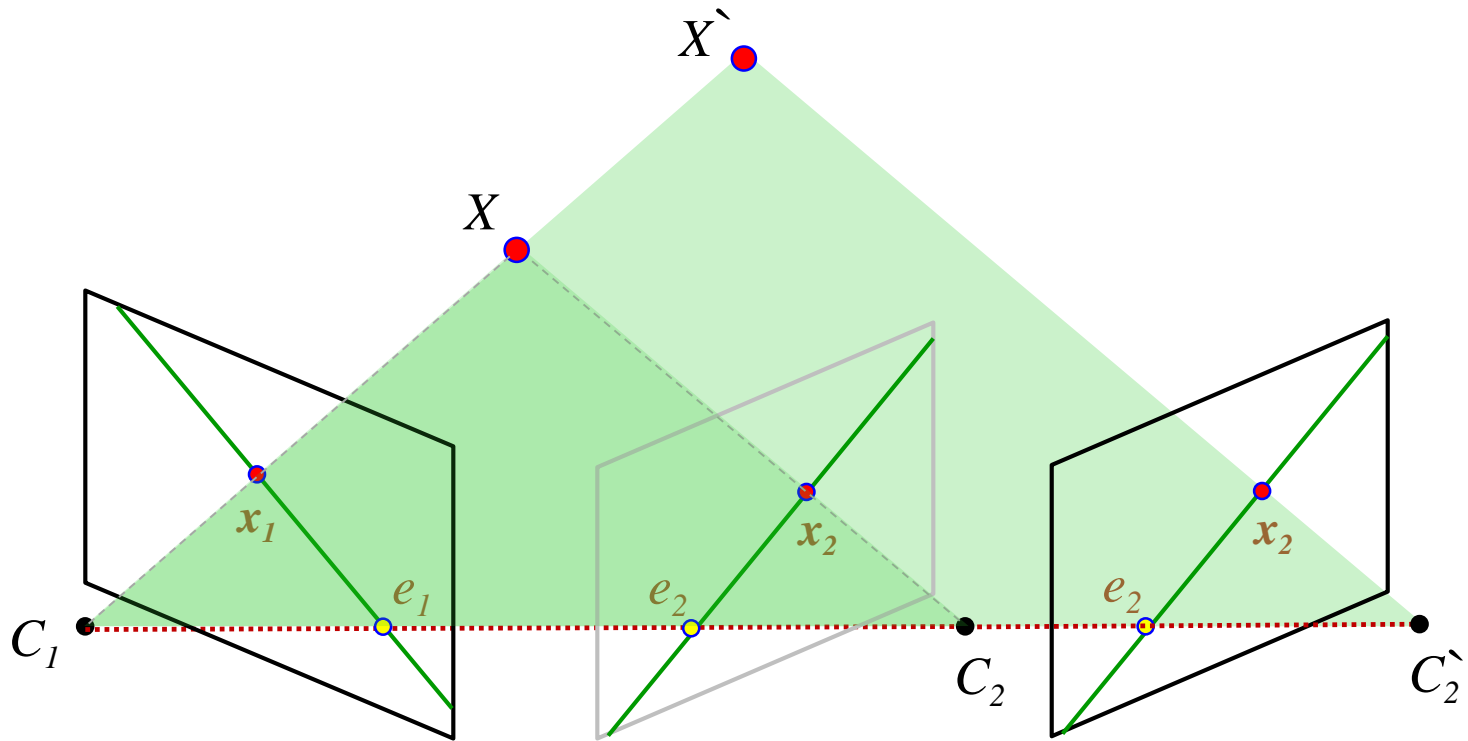


camera B orientation flips ($R = UWV^T$ or UW^TV^T)

Note: only one solution has positive “depths” for both cameras

Causes for 3D reconstruction ambiguity:

- **scale** remember: epipolar geometry can not help to estimate baseline length $|T|$



baseline $T = C_1C_2$

larger baseline
 $T' = C_1C'_2$

Causes for 3D reconstruction ambiguity:

- **scale**
- **position+orientation ?**

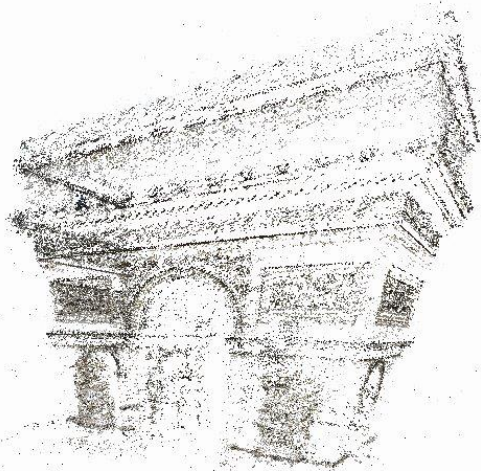
epipolar geometry determines only relative camera positions

Extracting cameras from fundamental matrix F

One can also estimate camera projection matrices from **fundamental matrix**, but there are more ambiguities [see H&Z]

Examples

[from Carl Olsson]



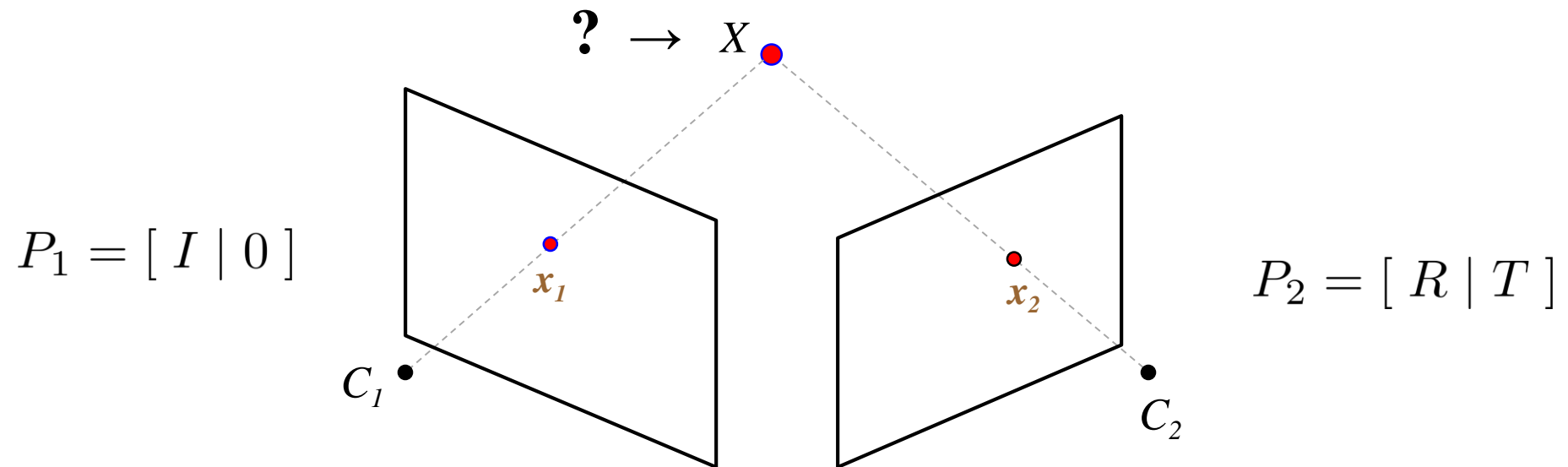
3D reconstruction with
“projective” ambiguity
(cameras estimated from F)



3D reconstruction with
similarity transform ambiguity
(cameras estimated from E)

Triangulation

Now, assume known projection matrices P_1 , P_2 and a match $\mathbf{x}_1 \leftrightarrow \mathbf{x}_2$



projection constraints

$$\begin{bmatrix} w_1 u_1 \\ w_1 v_1 \\ w_1 \end{bmatrix} = P_1 \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad \begin{bmatrix} w_2 u_2 \\ w_2 v_2 \\ w_2 \end{bmatrix} = P_2 \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

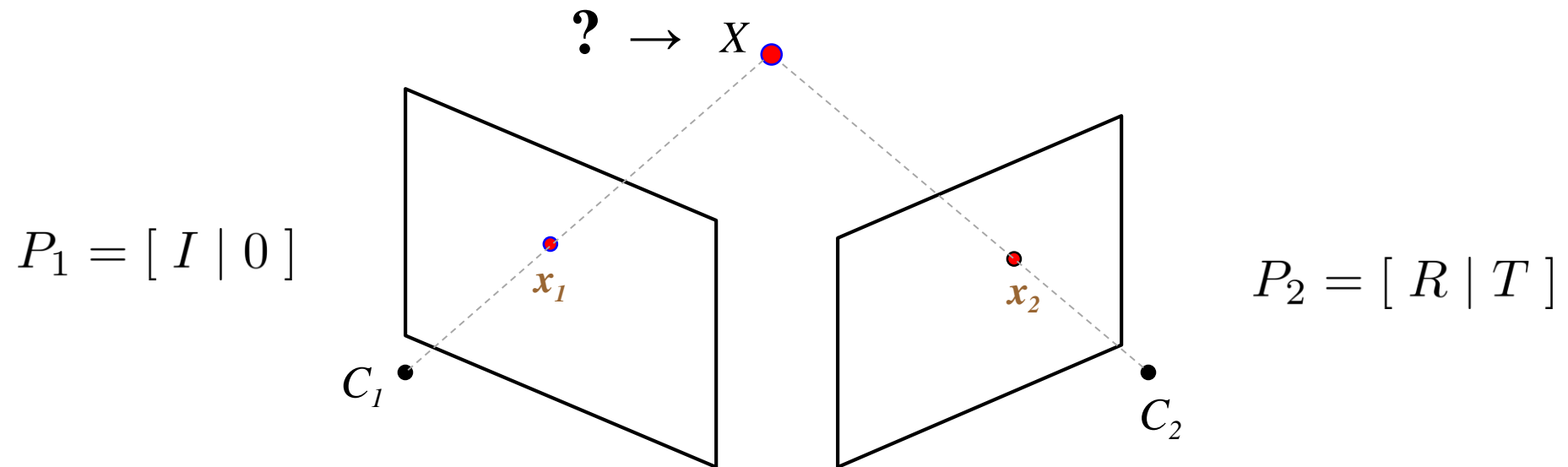
6 equations with 5 unknown (X, Y, Z, w_1, w_2)

But, we do not care about w_1 & w_2 – **eliminate** them (*à la* slide 26 topic 5)

\Rightarrow 4 equations with 3 unknown (X, Y, Z)

Triangulation

Now, assume known projection matrices P_1 , P_2 and a match $\mathbf{x}_1 \leftrightarrow \mathbf{x}_2$



projection constraints

$$\begin{bmatrix} w_1 u_1 \\ w_1 v_1 \\ w_1 \end{bmatrix} = P_1 \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad \begin{bmatrix} w_2 u_2 \\ w_2 v_2 \\ w_2 \end{bmatrix} = P_2 \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

One equation is redundant only if points x_1 , x_2 are exactly on the corresponding epipolar lines (the corresponding rays intersect in 3D).

Due to errors, use least squares.

Structure-from-Motion workflow

Basic sequential reconstruction

- For the first two images, use 8 point algorithm to estimate essential matrix E , cameras, and triangulate some points $\{X_i\}$.
- Each new view should see some previously reconstructed scene points $\{X_i\}$ (“feature matches” with previous cameras). Use such points to estimate new camera position (*resection problem*).
- Add new scene points using triangulation, e.g. for new “matches” with previously non-matched (and non-triangulated) features in earlier views.
- If there are more cameras, iterate previous two steps.
- **Issues**
 - errors can accumulate
 - new views are used only to add new 3D points, but they can help to improve accuracy for previously reconstructed scene

Structure-from-Motion workflow

“Bundle adjustment”

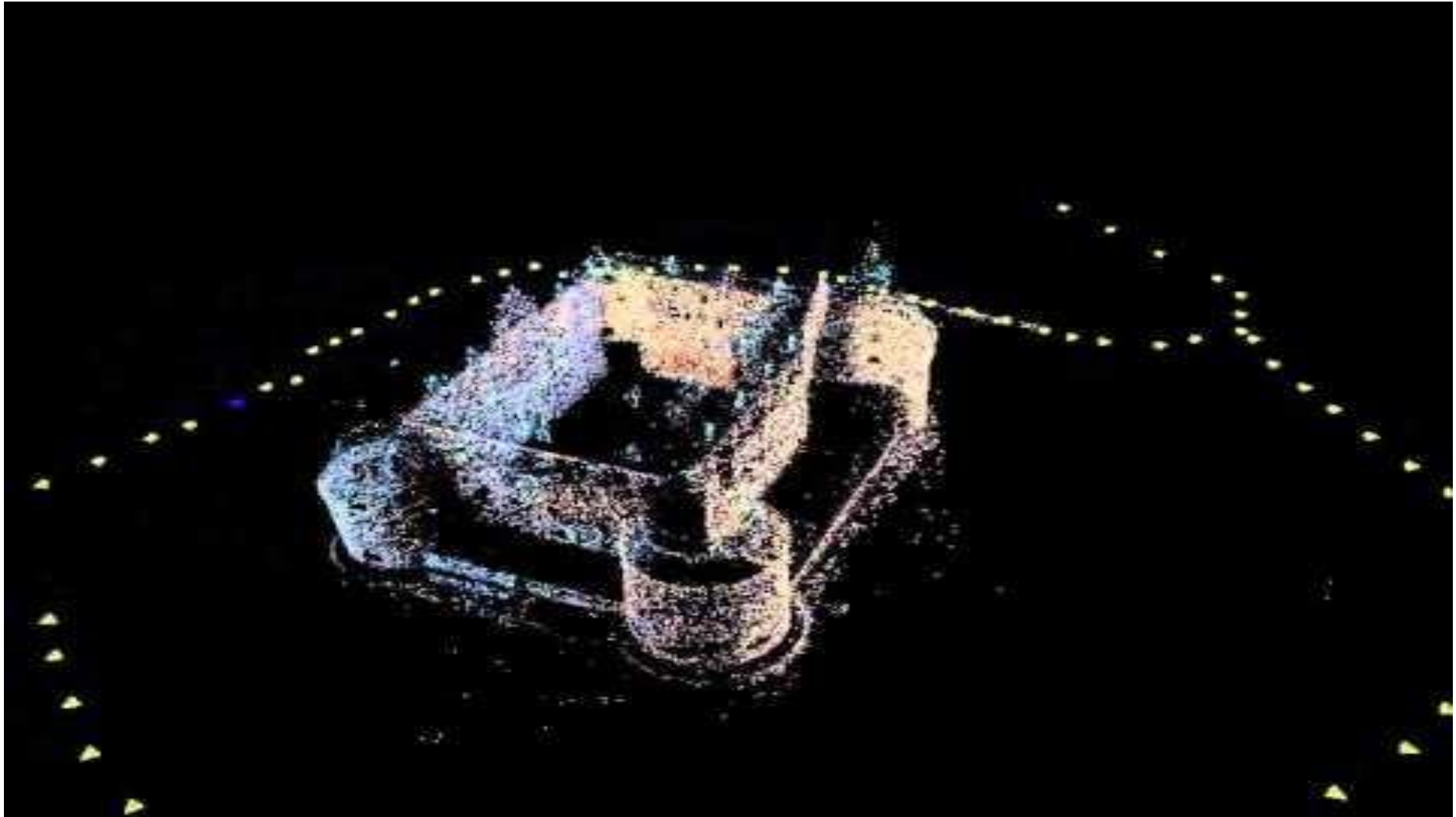
i -th “feature track” $tr_i := \{x_{ik} | k \in V(i)\}$

↑
feature i
location
in image k
↑
set of images
where feature i
is visible

$$\min_{\{P_k\}, \{X_i\}} \sum_i \sum_{k \in V(i)} \|x_{ik} - P_k X_i\|$$

projection error

Structure-from-Motion workflow



from Carl Olsson

Applications of multi-view geometry:

Pose estimation

Rigid motion segmentation

Augmented reality

Special effects in video

Volumetric 3D reconstruction

Depth reconstruction (stereo-next topic)

Examples:

We were fitting a single essential/fundamental matrix to a pair of images corresponding to two different view points

Q: Can matched features in two images support more than one fundamental matrix?

Examples:

We were fitting a single essential/fundamental matrix to a pair of images corresponding to two different view points

Q: Can matched features in two images support **more than one fundamental matrix?**

Hint: we assumed that the scene is stationary and only camera moved or, equivalently, that the camera is stationary but the whole 3D scene moved (R, T).

Examples:

We were fitting a single essential/fundamental matrix to a pair of images corresponding to two different view points

Q: Can matched features in two images support **more than one fundamental matrix**?

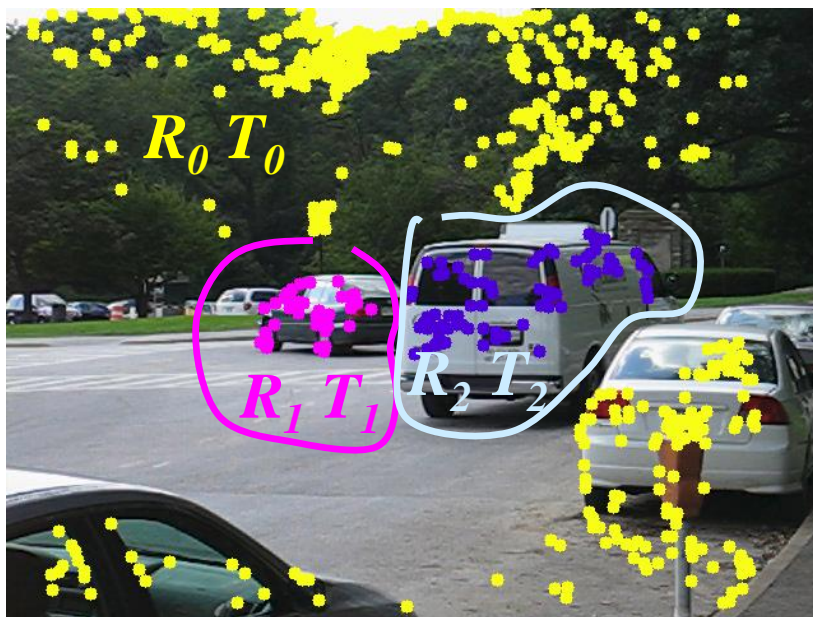


multi-model fitting with
fundamental matrices
(UFL, previous topic)

Examples: Rigid Motion Estimation

We were fitting a single essential/fundamental matrix to a pair of images corresponding to two different view points

Q: Can matched features in two images support **more than one fundamental matrix**?



multi-model fitting with
fundamental matrices
(UFL, previous topic)

Examples: Augmented Reality

- if camera position C and orientation R are known (in addition to K) then can insert “new” objects into the 3D scene
- particularly useful for movies: camera path can be computed
 - can generate correct views of new objects



www.2d3.com
 (“boujou”)

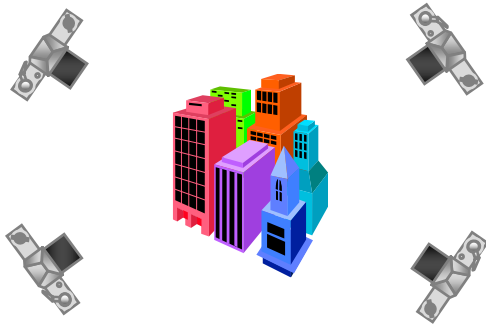
Examples: Dense 3D Reconstruction

Sparse reconstruction (points in 3D) is done by triangulating point correspondences (part of Structure-from-Motion problem)

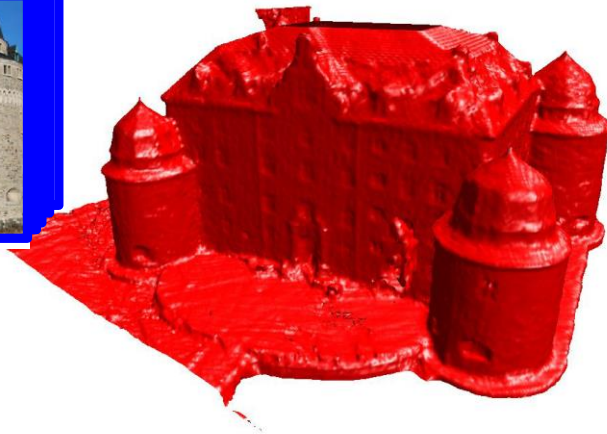
How about **dense** 3D reconstruction from n views?

Examples: Dense 3D Reconstruction

A. computing *Surfaces in 3D Volumes* (volumetric reconstruction)

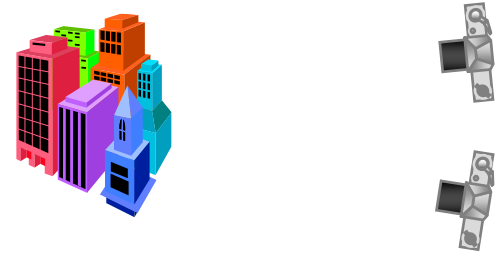


multiple wide baseline views

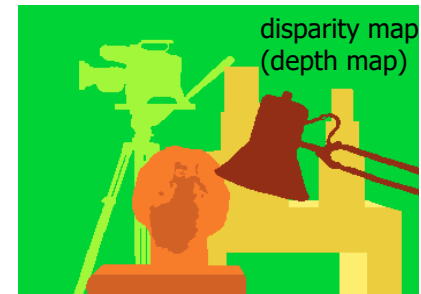


Relates to volumetric **segmentation** (topic 9)

B. computing dense *Depth Maps* (stereo)



two narrow baseline views



will discuss in **topic 8: stereo**

From sparse features to dense reconstructions

- Find & match features in 2 or more images (**sparse points**)
- Estimate epipolar geometry
 - also gives camera position
 - triangulate for **sparse** 3D reconstruction
- Now, we can move towards **denser reconstruction**
 - find many more matches (correspondences) using known epipolar lines: constrained search space significantly reduces ambiguity for feature matching
 - use “**regularization**” to estimate surfaces or depth maps