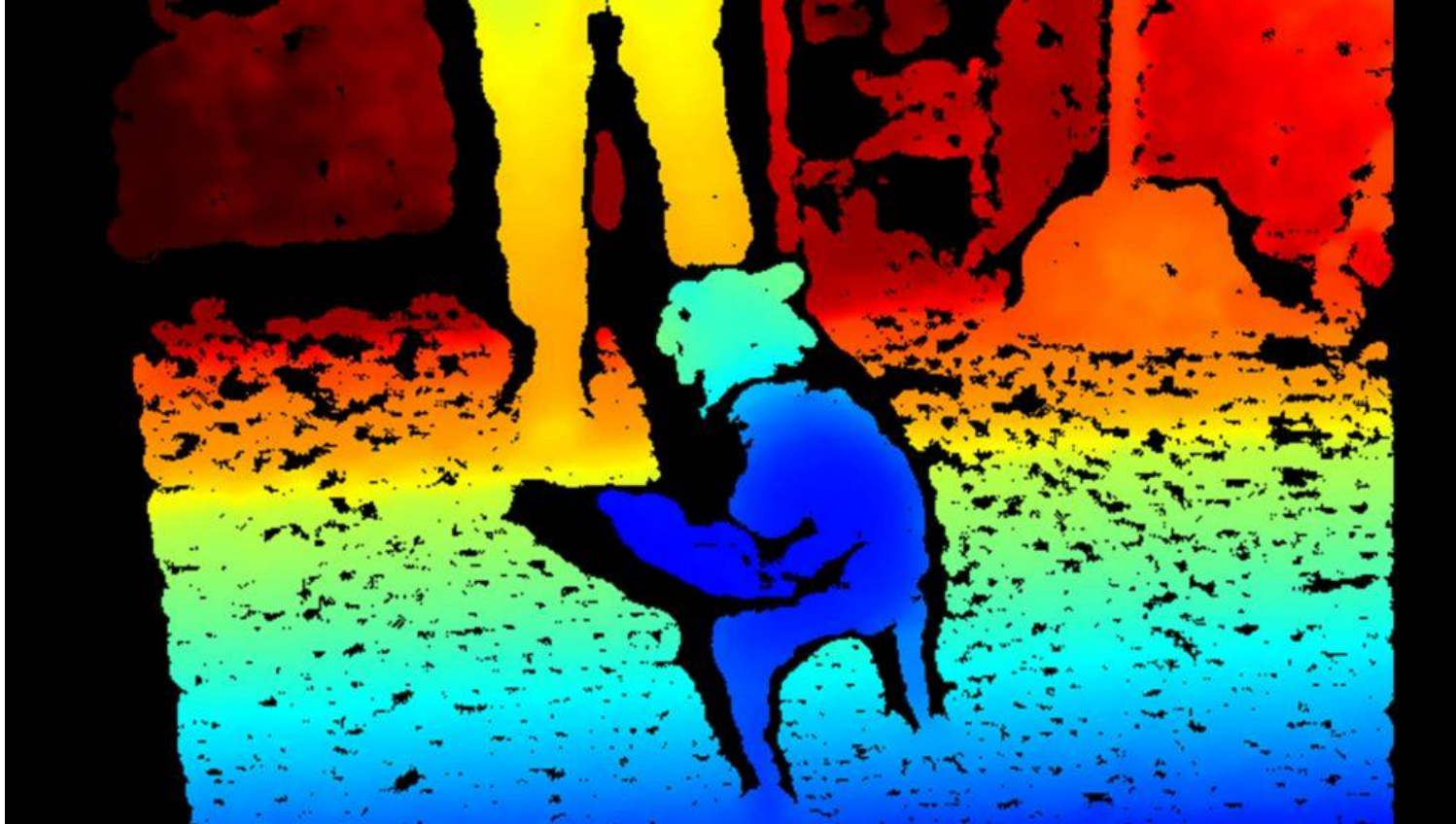
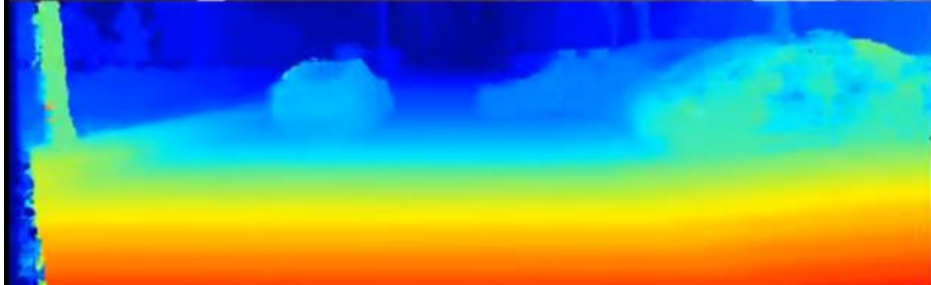


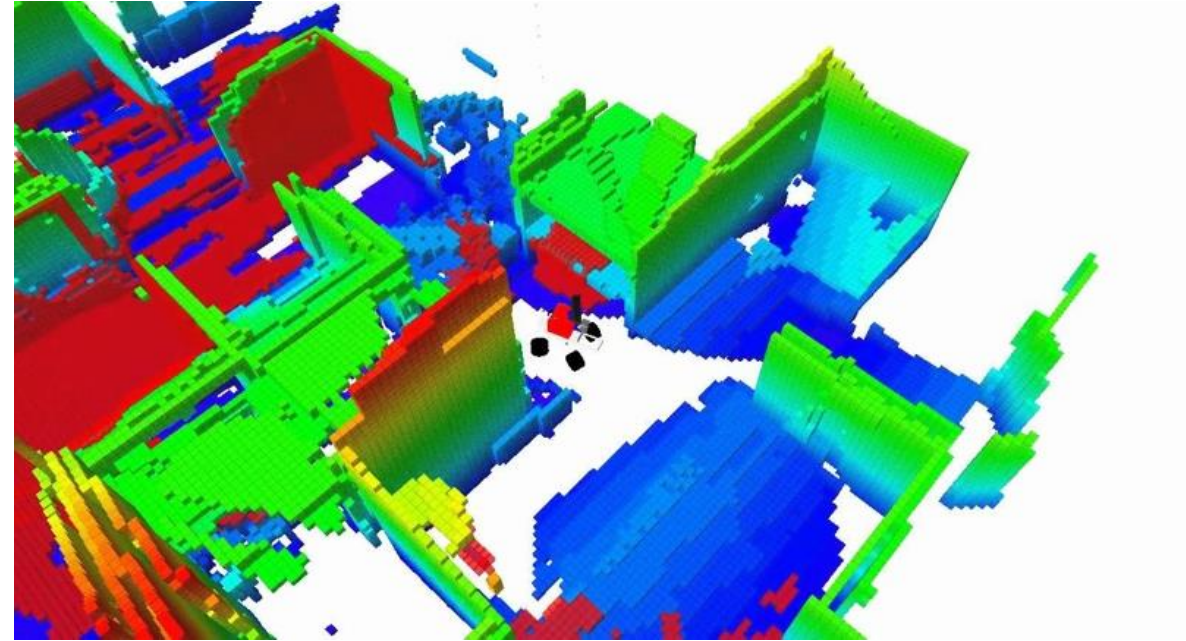
# Stereo



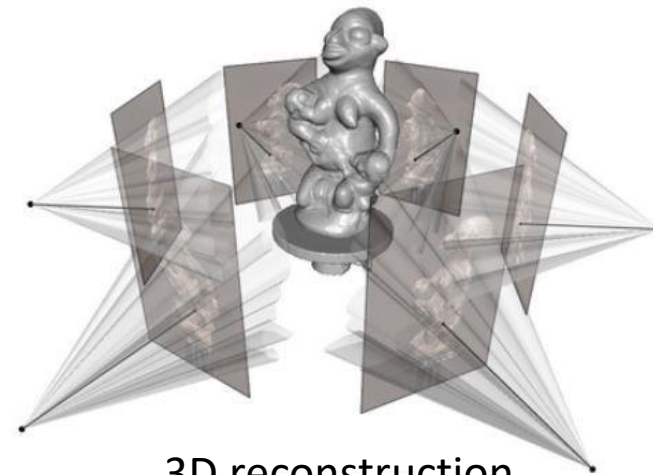
# What can be done with stereo vision?



Autonomous driving



SLAM- robot navigation



3D reconstruction

# References

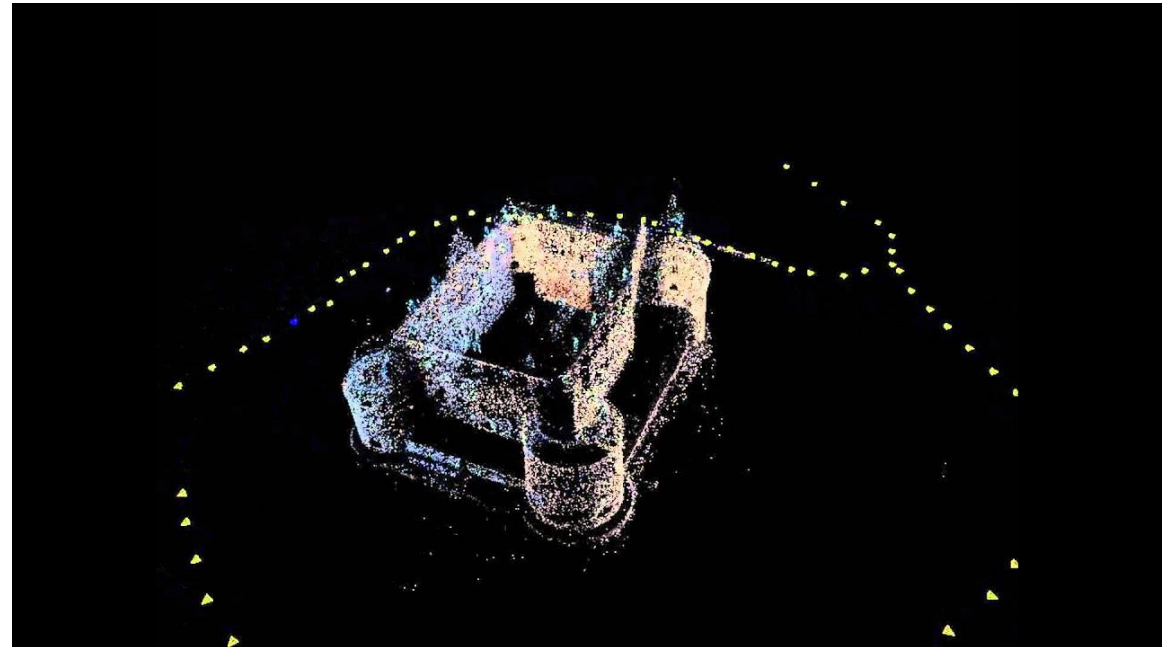
- <http://szeliski.org/Book/>
- <http://www.cs.cornell.edu/courses/cs5670/2019sp/lectures/lectures.html>
- <http://www.cs.cmu.edu/~16385/>

# Contents

- **Structure from motion**
- Triangulation
- Stereo matching
- Camera rectification
- Epipolar geometry
  - Essential matrix
  - Fundamental matrix
  - Estimating the fundamental matrix
- Other 3D sensors

# Structure from motion

- **Structure from motion (SfM)** is the process of estimating the 3-D structure of a scene from a set of 2-D images. SfM is used in many applications, such as 3-D scanning and augmented reality.
  - [Mathworks]
- SfM is also known as **3D reconstruction**.
- **Stereo vision** is a subcategory of SfM in which we are dealing only with 2 images.



# Structure and motion

	Structure (3D model of world)	Motion (6 DOFs of cameras)
Pose Estimation (camera pose estimation)	Known	<b>Estimate</b>
Triangulation	<b>Estimate</b>	Known
3D reconstruction/ SfM/ stereo vision	<b>Estimate</b>	<b>Estimate</b>

# Structure and motion

- So essentially one can say that “structure from motion” is the wrong name...
  - Structure and motion is more precise, but nobody will understand what are you talking about.
- In this class we will learn about 3D reconstruction from two cameras (and triangulation as a subtopic).

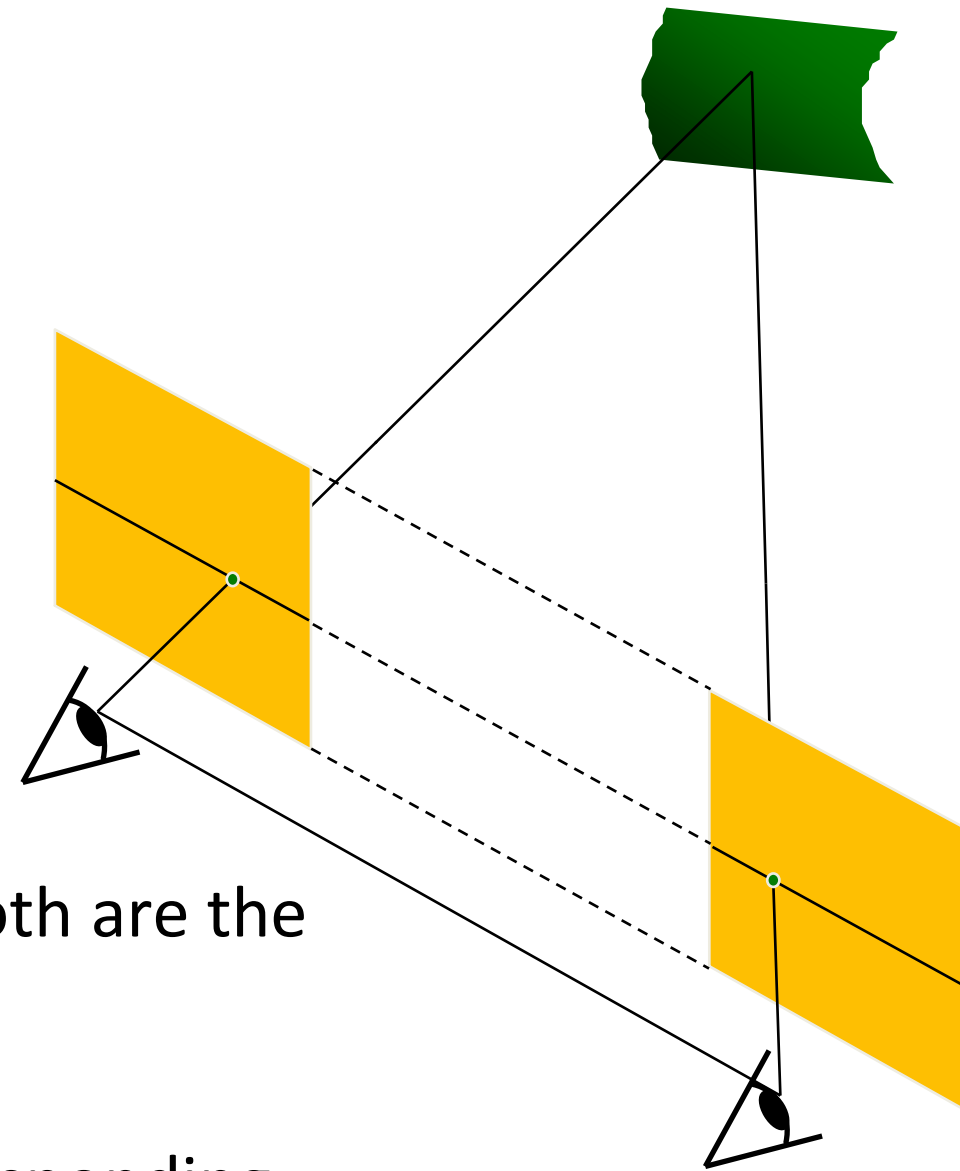


# Contents

- Structure from motion
- **Triangulation**
- Stereo matching
- Camera rectification
- Epipolar geometry
  - Essential matrix
  - Fundamental matrix
  - Estimating the fundamental matrix
- Other 3D sensors



# Triangulation



- Assume both cameras are rectified- 6 DOF of both are the same except the horizontal translation.
- Assume same focal length  $f$  in both cameras
- Assume we know for each pixel in left the corresponding pixel in right.
- From this we want to get a depth image using triangulation.

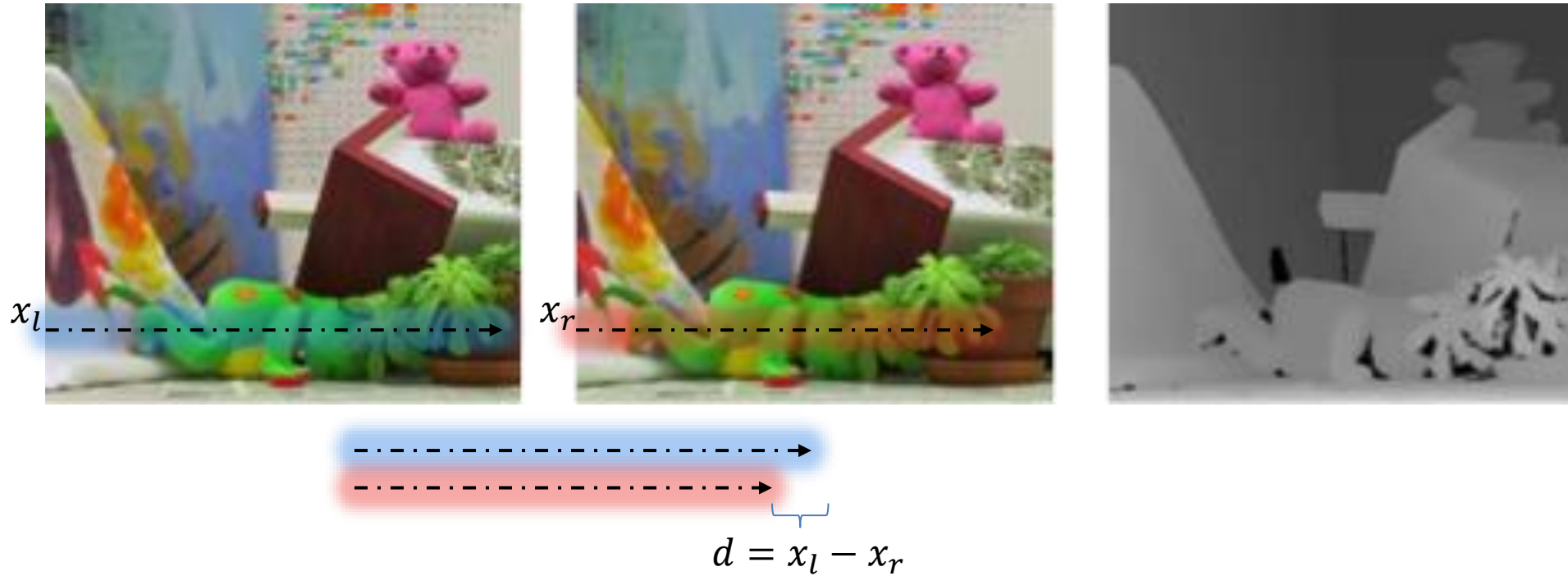


Left



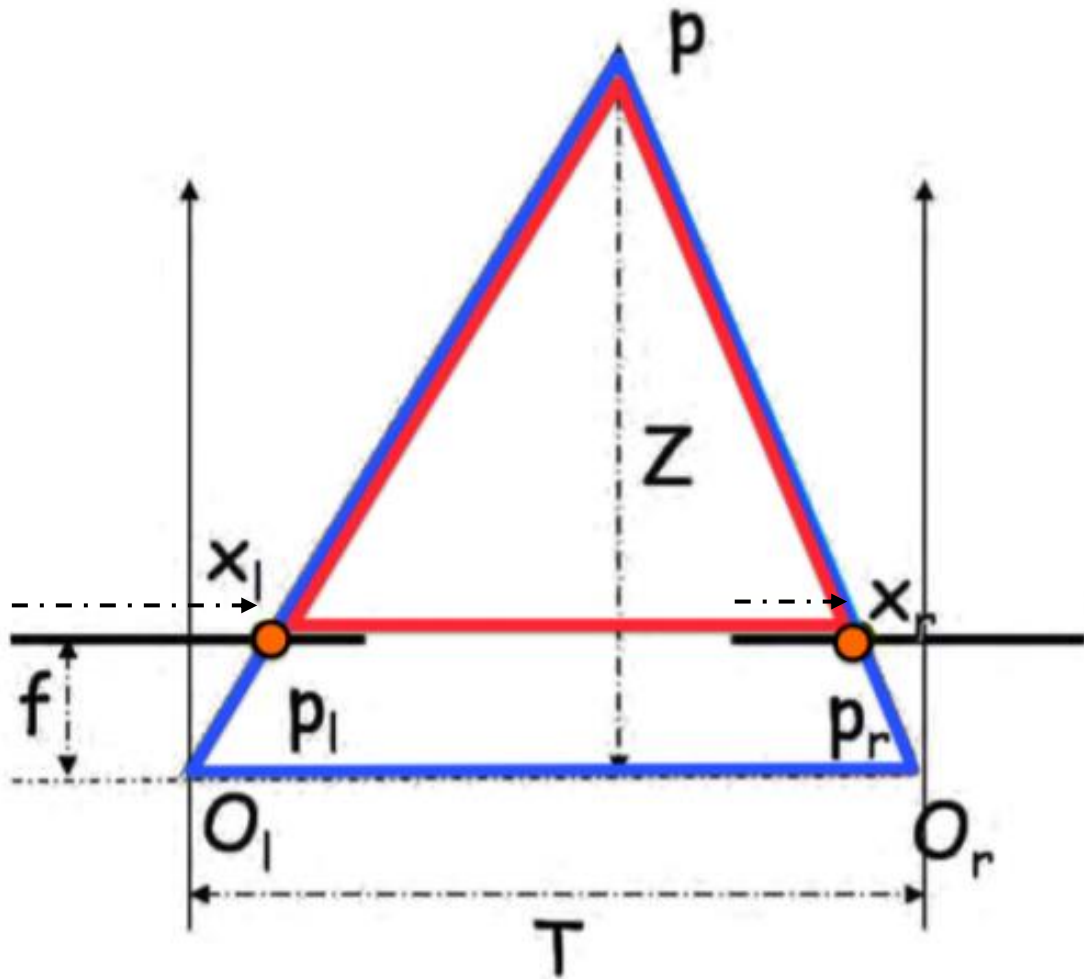
Right

# Triangulation



- The amount of horizontal movement is inversely proportional to the distance from the camera.
- The amount of horizontal movement == disparity ( $d = x_l - x_r$ ).
- Distance from the camera == depth (or Z).
- Note:  $x_l$  &  $x_r$  are in normalized image coordinate system:  $x = K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$

# Triangulation



Similar triangles:

$$\frac{T}{Z} = \frac{\overbrace{T + x_r - x_l}^{T - d}}{Z - f}$$

An orange curved arrow points from the similar triangles equation to the final formula. The formula is enclosed in an orange rounded rectangle:

$$Z = \frac{f \cdot T}{x_l - x_r}$$

Labels with arrows point to the components of the formula:

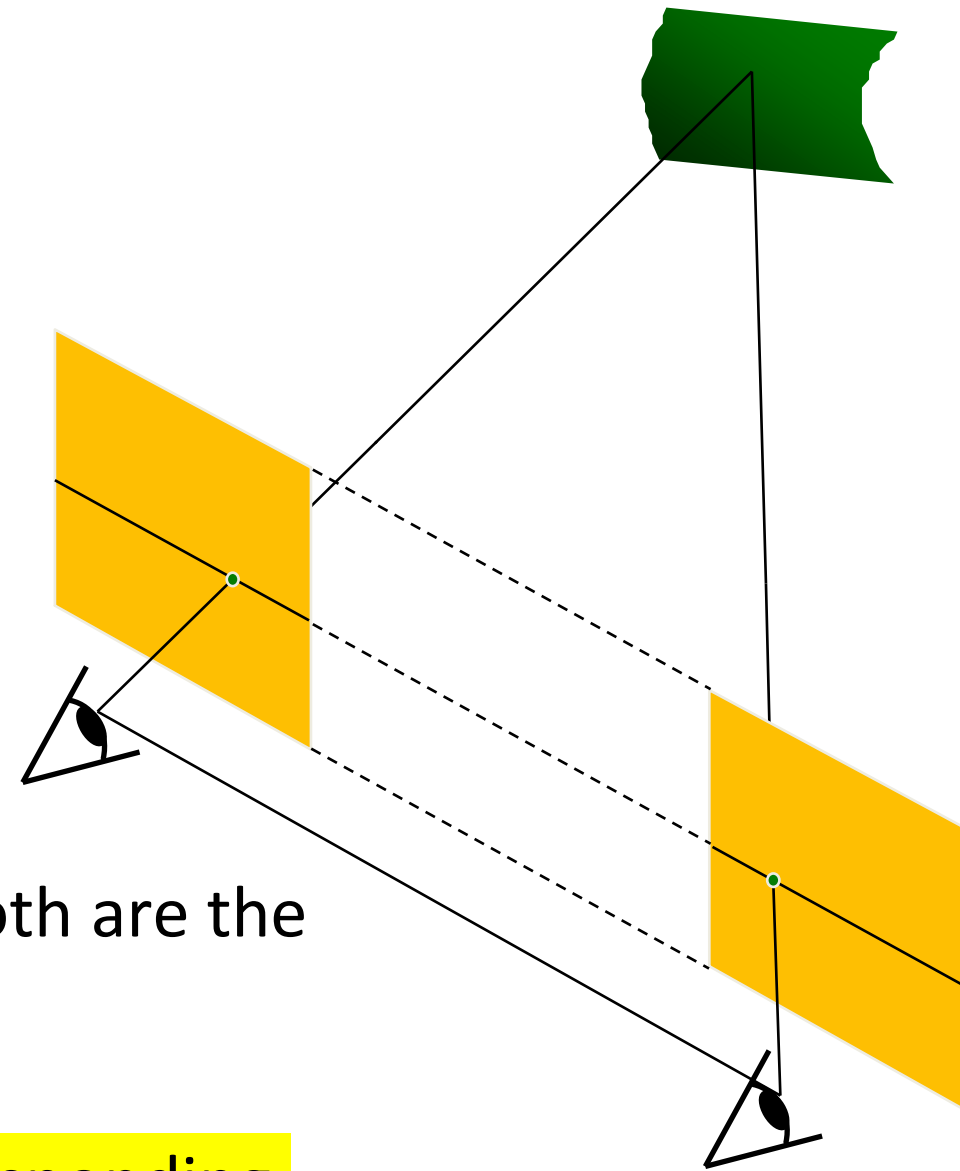
- baseline** points to  $T$  in the numerator.
- focal length** points to  $f$  in the numerator.
- disparity** points to  $x_l - x_r$  in the denominator, which is enclosed in a blue dotted rectangle.

# Contents

- Structure from motion
- Triangulation
- **Stereo matching**
- Camera rectification
- Epipolar geometry
  - Essential matrix
  - Fundamental matrix
  - Estimating the fundamental matrix
- Other 3D sensors



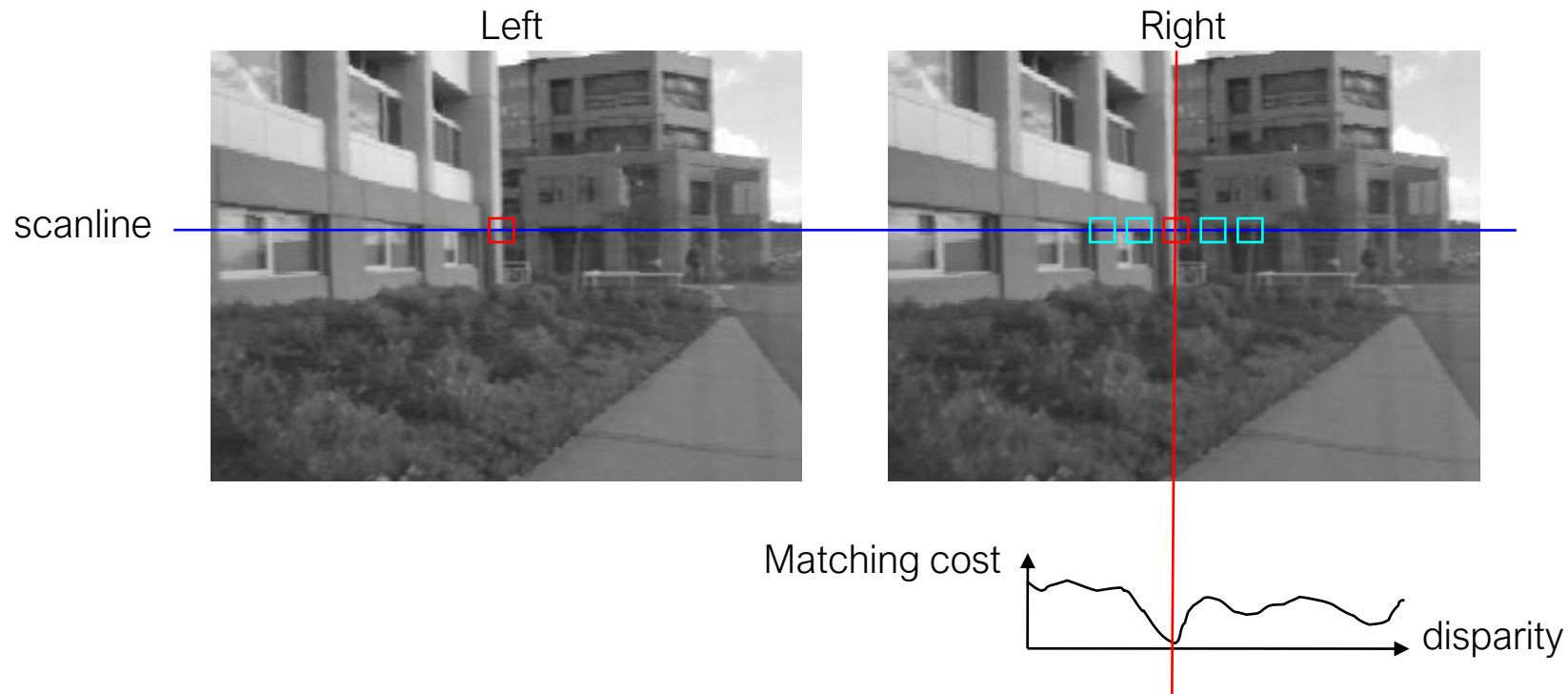
# Triangulation



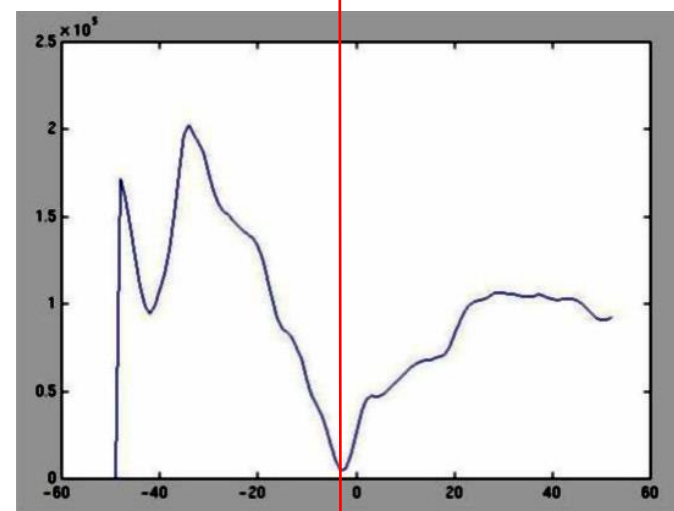
- Assume both cameras are rectified- 6 DOF of both are the same except the horizontal translation.
- Assume same focal length  $f$  in both cameras
- Assume we know for each pixel in left the corresponding pixel in right.
- From this we want to get a depth image using triangulation.



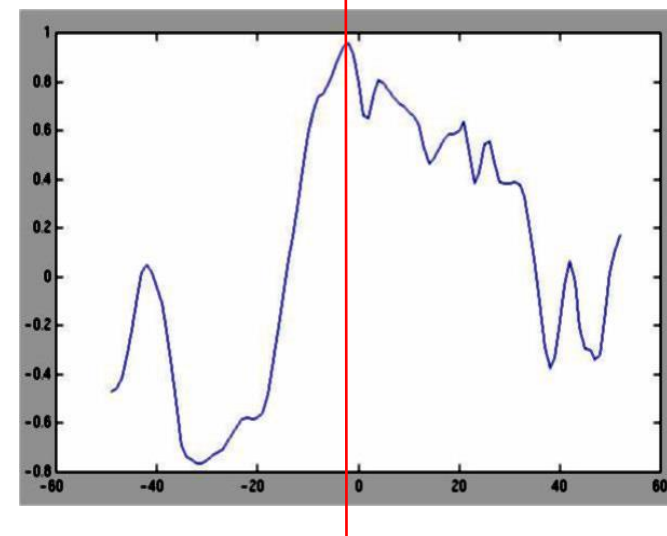
# Stereo Block Matching



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



SSD



Normalized cross-correlation

# Effect of window size



$W = 3$



$W = 20$

# Effect of window size



$W = 3$

## **Smaller window**

- + More detail
- More noise

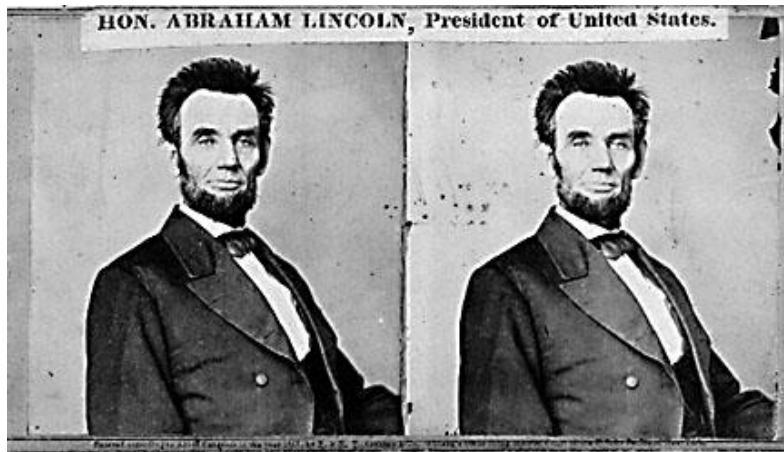


$W = 20$

## **Larger window**

- + Smoother disparity maps
- Less detail
- Fails near boundaries

## *When will stereo block matching fail?*





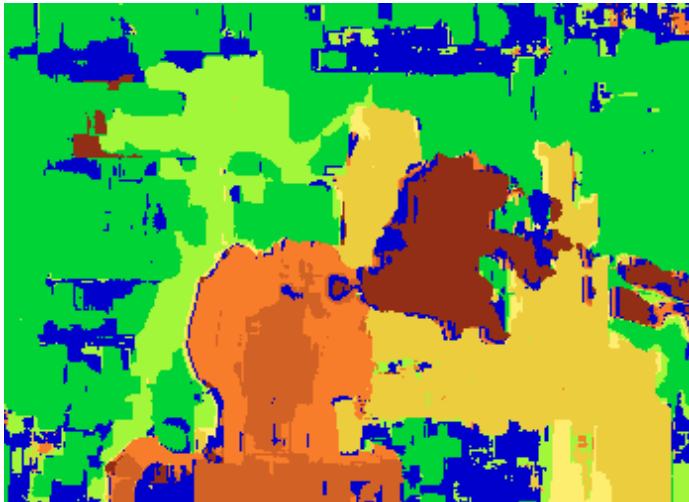
## *When will stereo block matching fail?*







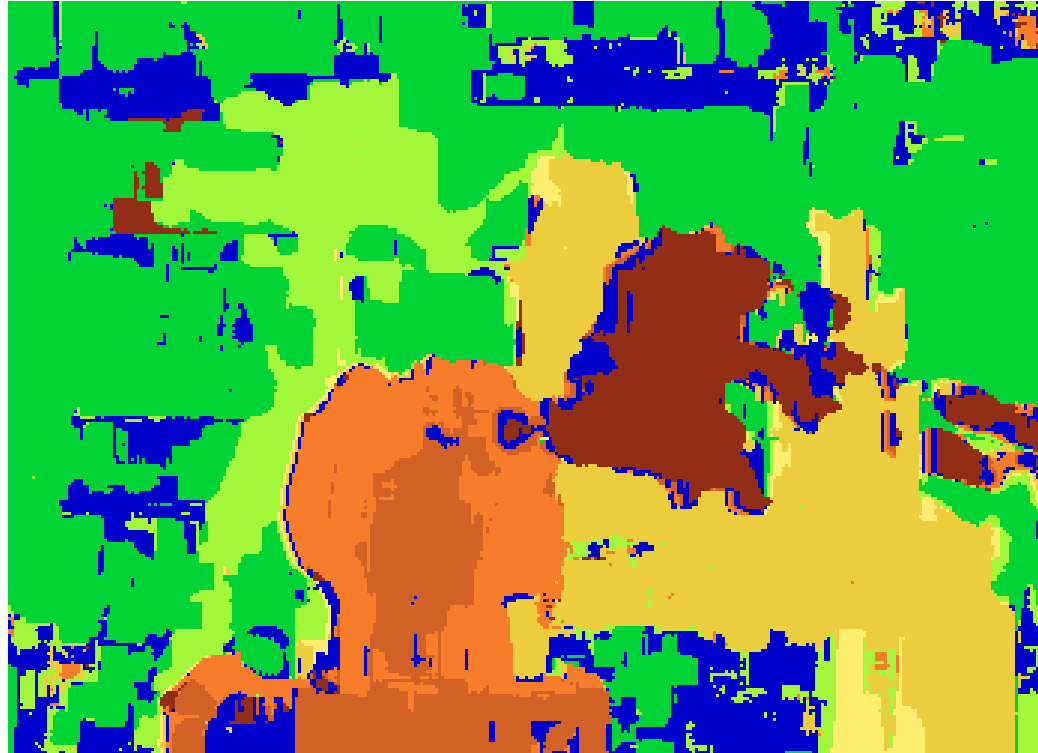
Block matching



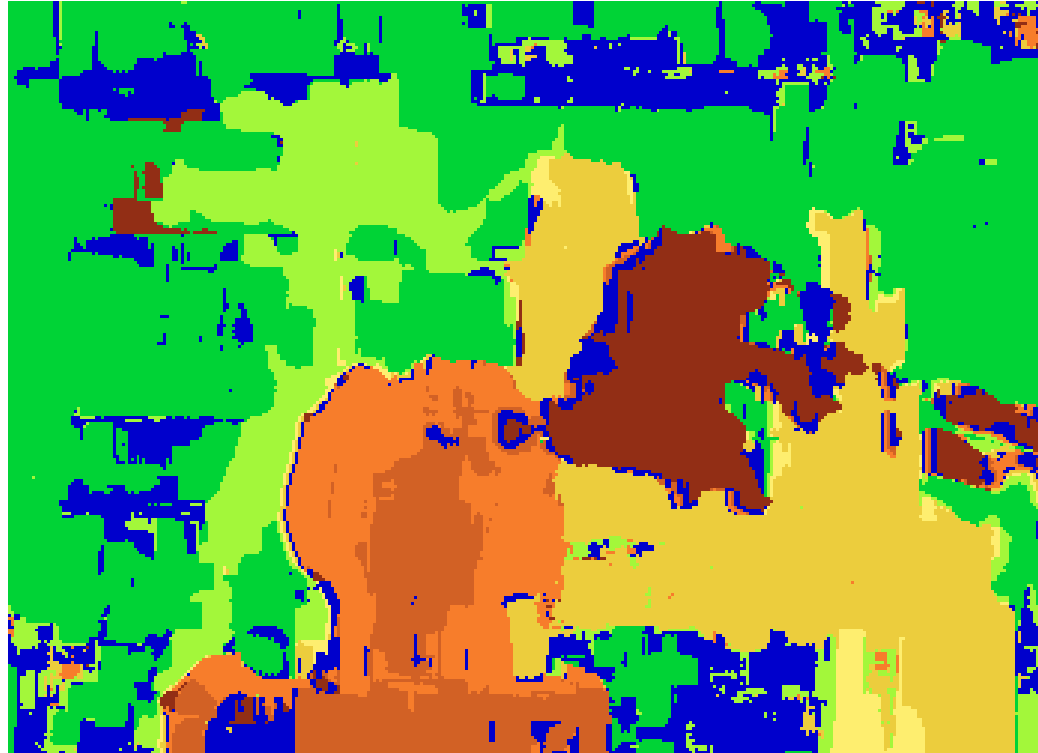
Ground truth



*What are some problems with the result?*



*How can we improve depth estimation?*

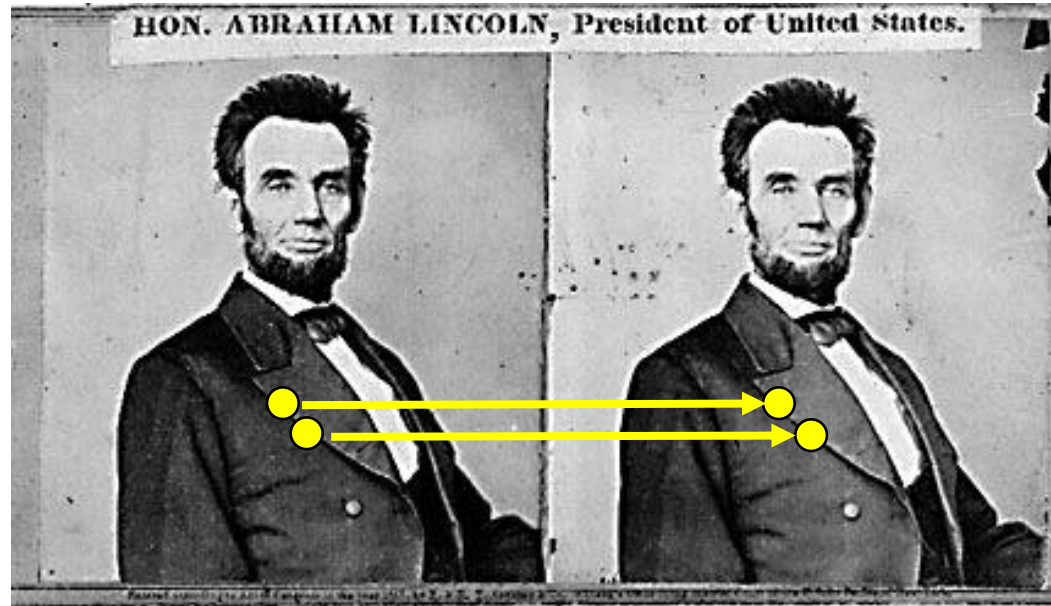


*How can we improve depth estimation?*

Too many discontinuities.  
We expect disparity values to change slowly.

Let's make an assumption:  
**depth should change smoothly**

# Energy Minimization



What defines a good stereo correspondence?

- 1. Match quality**

- Want each pixel to find a good match in the other image

- 2. Smoothness**

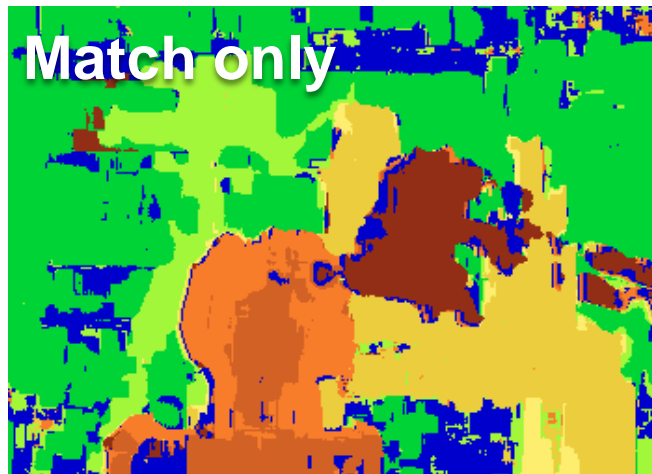
- If two pixels are adjacent, they should (usually) move about the same amount

energy function  
(for one pixel)

$$E(d) = \underbrace{E_d(d)}_{\text{data term}} + \lambda \underbrace{E_s(d)}_{\text{smoothness term}}$$

Want each pixel to find a good match  
in the other image  
(block matching result)

Adjacent pixels should (usually)  
move about the same amount  
(smoothness function)



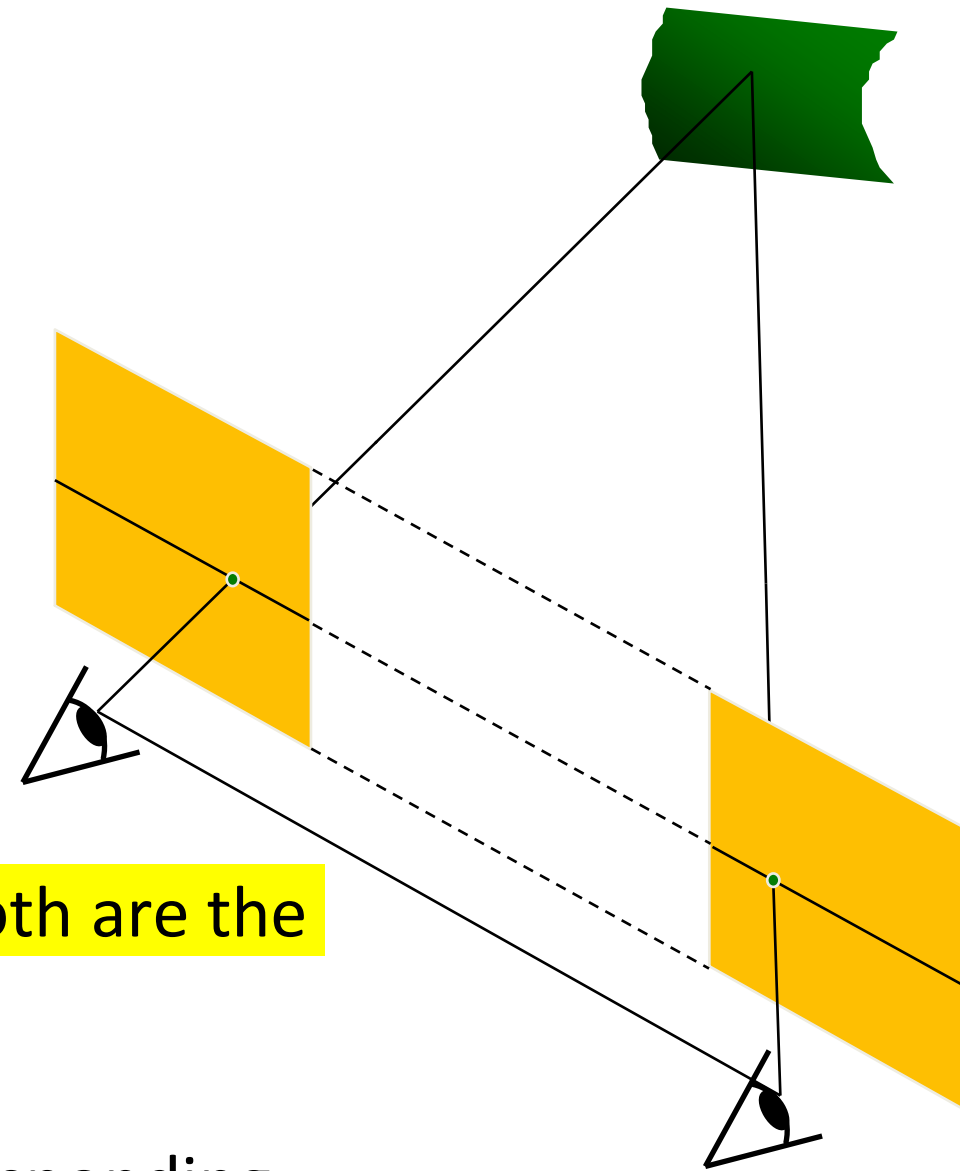
Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001

# Contents

- Structure from motion
- Triangulation
- Stereo matching
- **Camera rectification**
- Epipolar geometry
  - Essential matrix
  - Fundamental matrix
  - Estimating the fundamental matrix
- Other 3D sensors



# Triangulation



- Assume both cameras are rectified- 6 DOF of both are the same except the horizontal translation.
- Assume same focal length  $f$  in both cameras.
- Assume we know for each pixel in left the corresponding pixel in right.
- From this we want to get a depth image using triangulation.



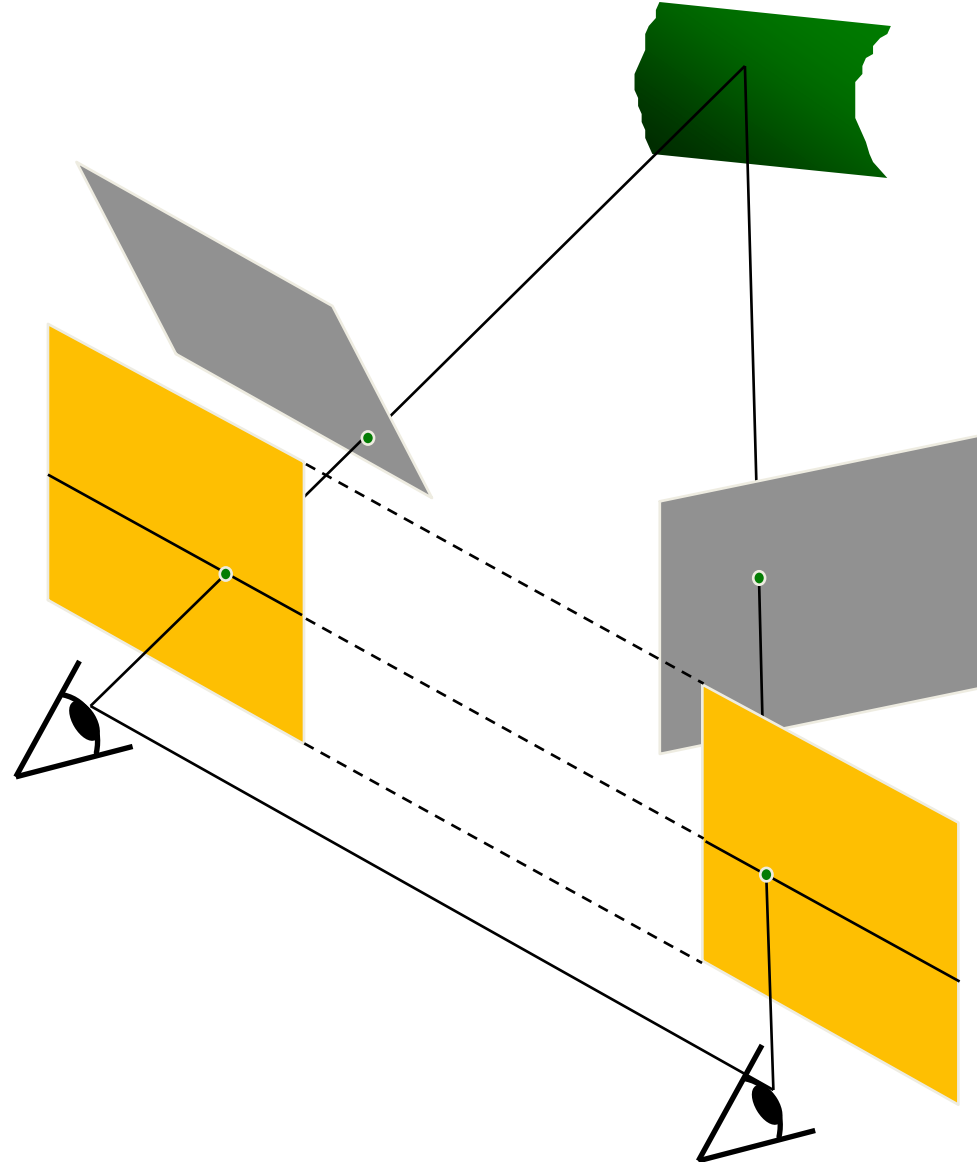
Original stereo pair



After rectification

# Stereo image rectification

- Out of scope...
- Assuming we know  $R, t$  of cam2 relative to cam1 we can use homography on both images to rectify them
  - Rectification proof here:  
[https://www.cs.cmu.edu/~16385/s17/Slides/13.1\\_Stereo\\_Rectification.pdf](https://www.cs.cmu.edu/~16385/s17/Slides/13.1_Stereo_Rectification.pdf)
- We want to find the relative  $R, t$  of the images.



# Contents

- Structure from motion
- Triangulation
- Stereo matching
- Camera rectification
- **Epipolar geometry**
  - Essential matrix
  - Fundamental matrix
  - Estimating the fundamental matrix
- Other 3D sensors

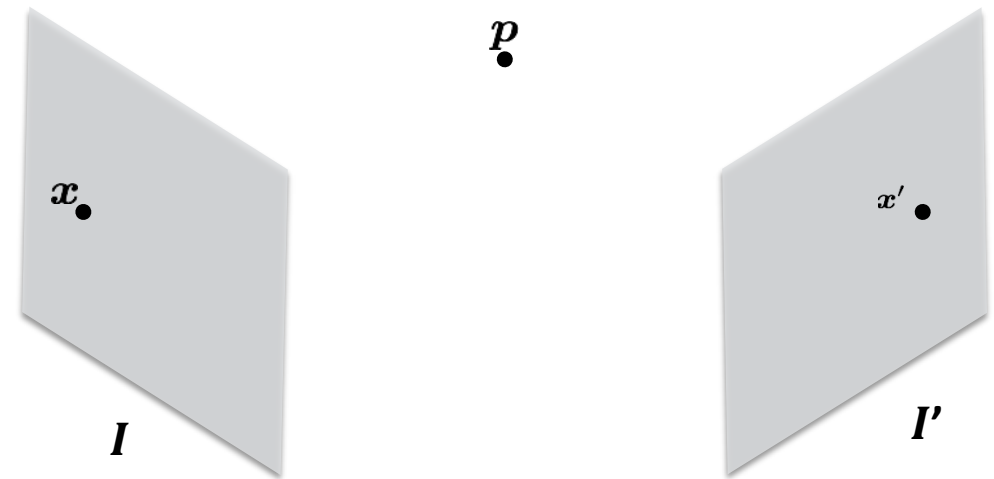
# Epipolar geometry

- **Epipolar geometry** is the geometry of stereo vision. When two cameras view a 3D scene from two distinct positions, there are a number of geometric relations between the 3D points and their projections onto the 2D images that lead to constraints between the image points.
  - [Wikipedia]

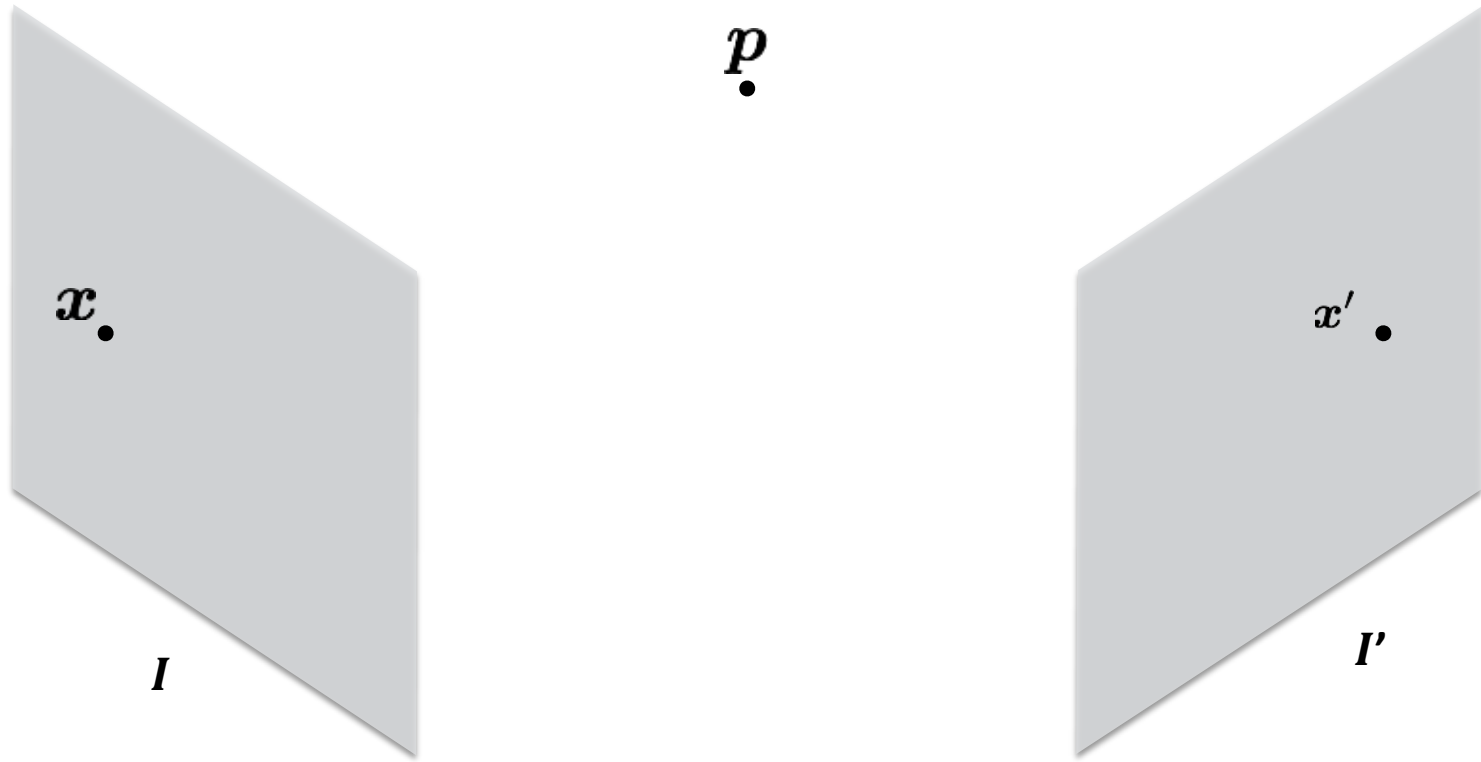
# Epipolar geometry - The triangulation problem

- Given:
  - two 2D points in the **normalized image coordinate system**  $(x, x')$  in two different images  $(I, I')$  that describes the same point  $p$  in 3D space.
  - Rotation and translation between the two cameras.
- Find  $p$ .

- Normalized image coordinate system:  $x = K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$



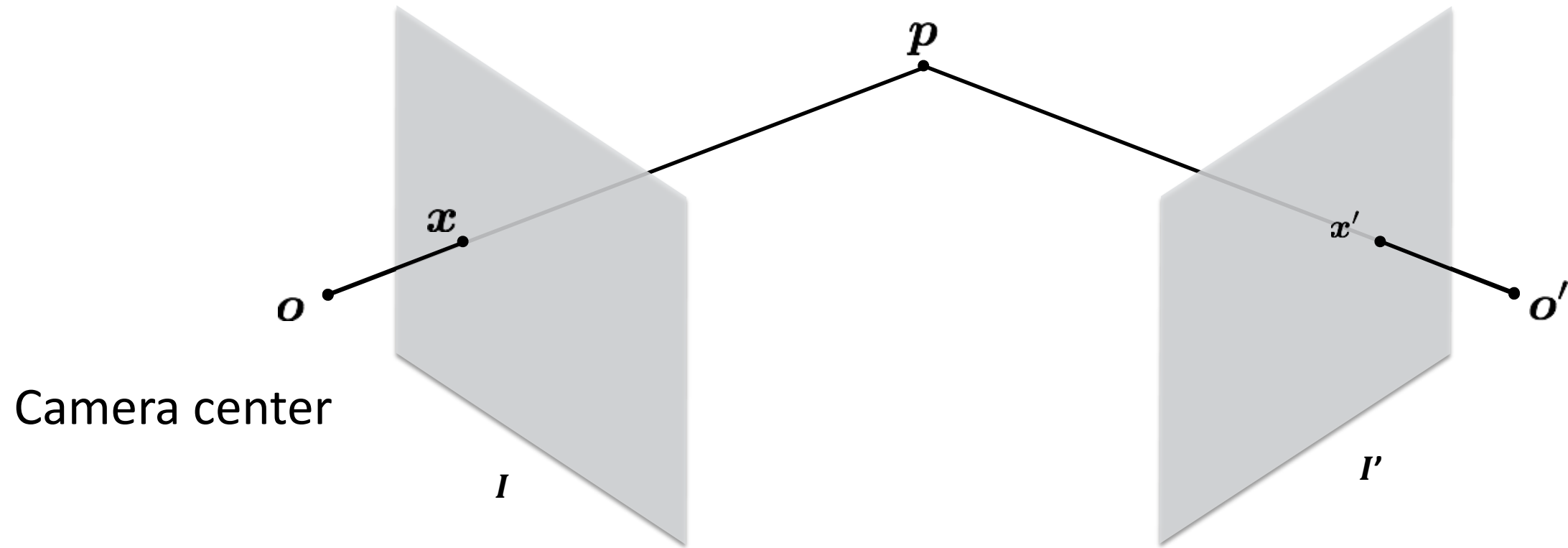
# Epipolar geometry





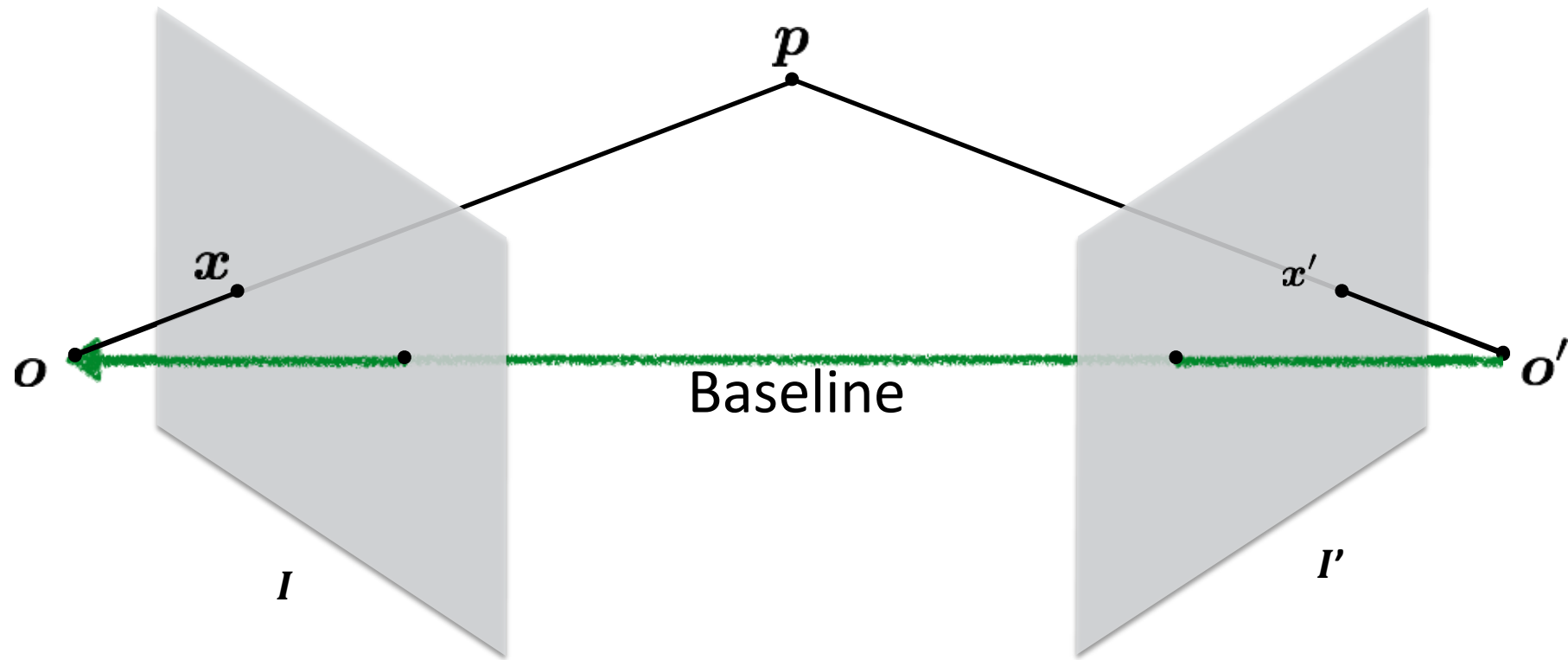
# Epipolar geometry

- We can trace lines from the **camera center** of each image, through the given 2D point to the 3D point  $p$ .



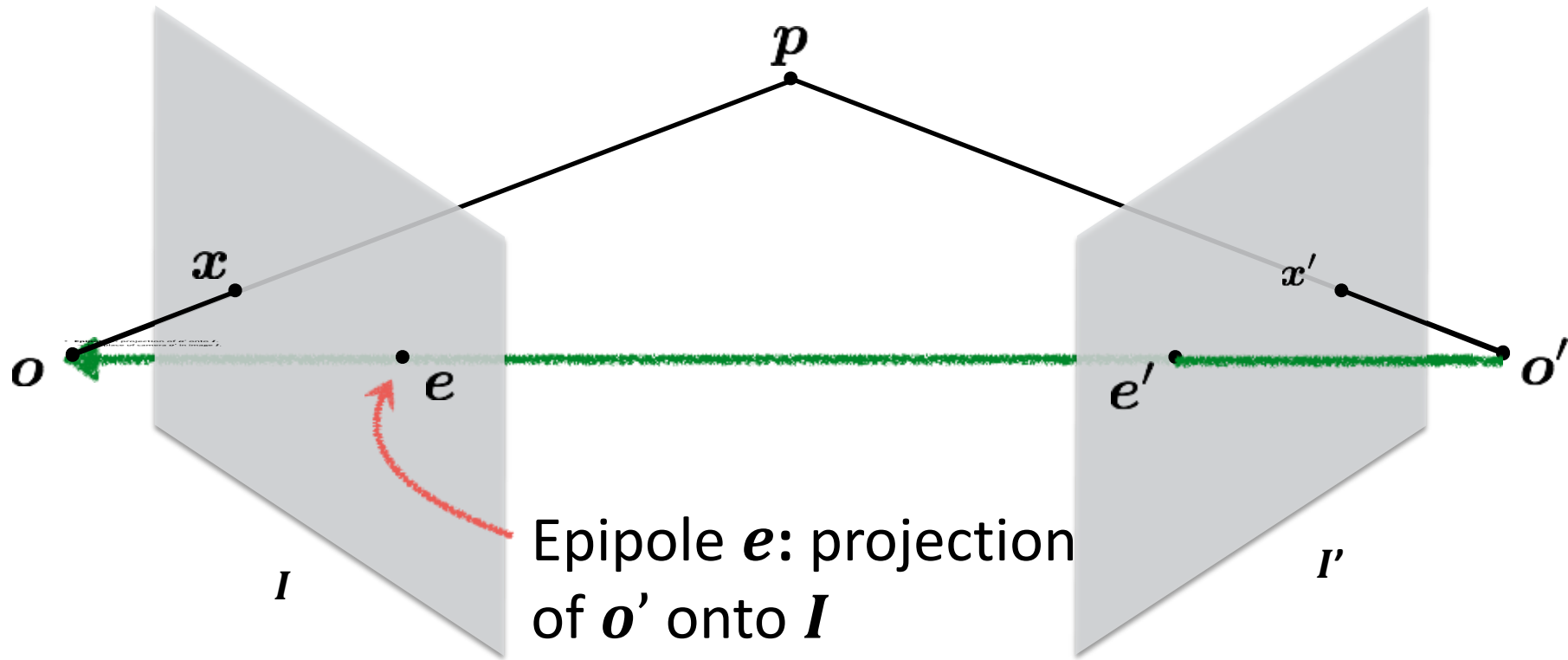
# Epipolar geometry

- **Baseline** is a vector that represent the translation between two cameras



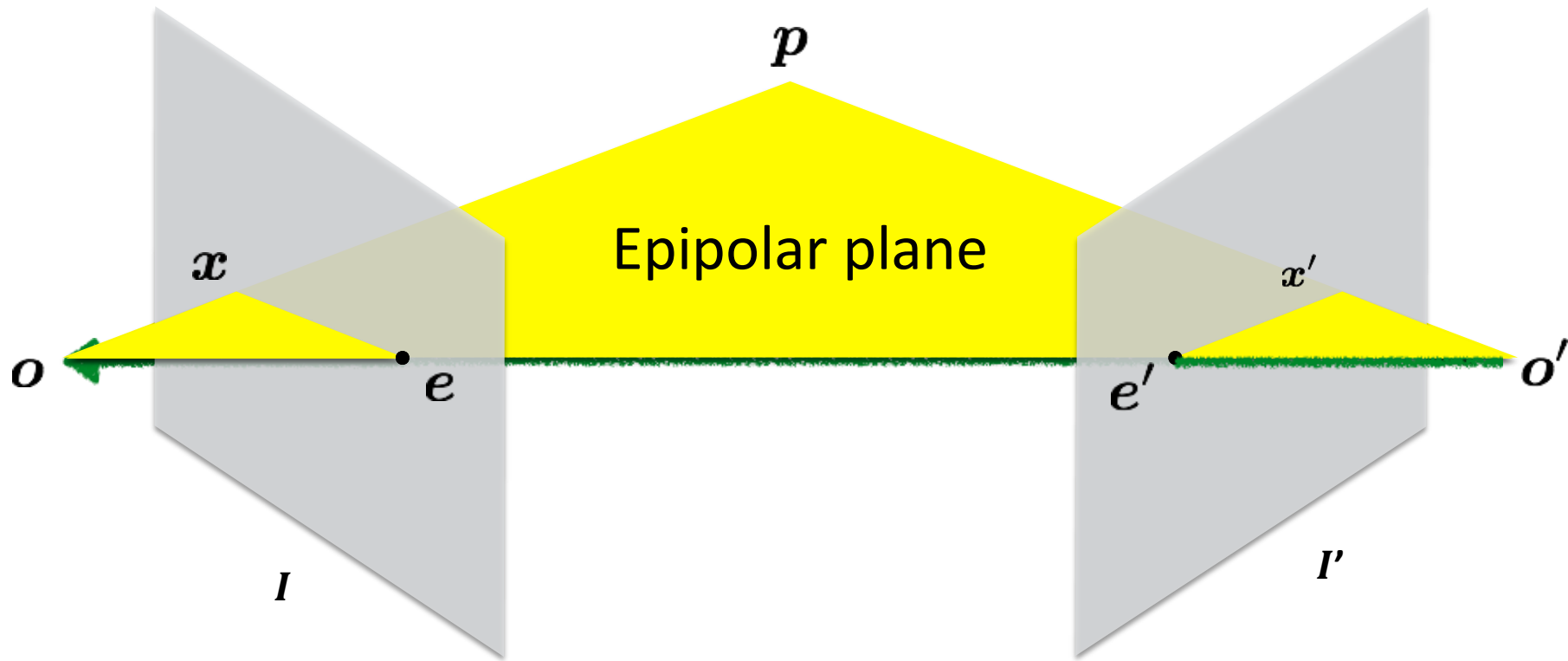
# Epipolar geometry

- **Epipole  $e$ :** projection of  $o'$  onto  $I$ .
  - The place of camera  $o'$  in image  $I$ .



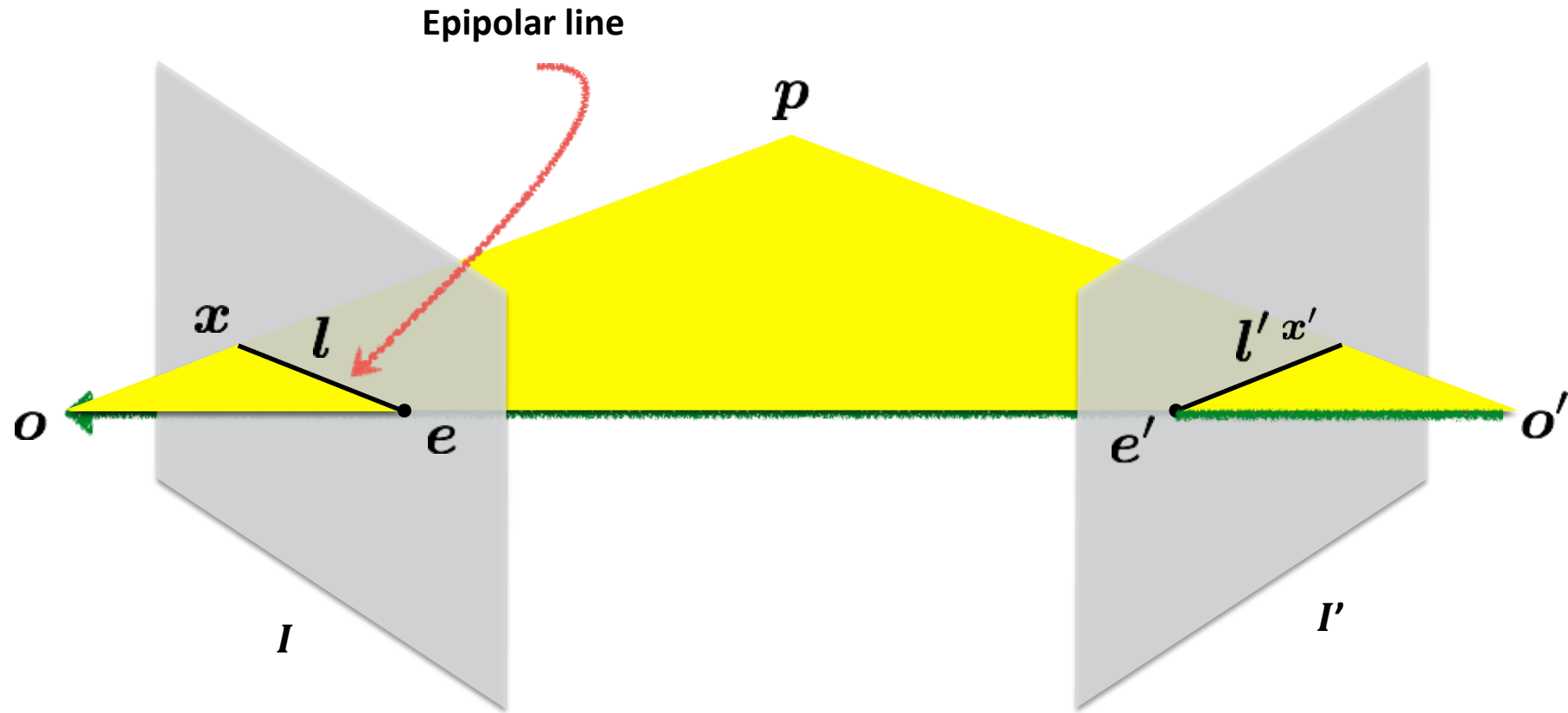
# Epipolar geometry

- **Epipolar plane:** the plane that is constructed from the 3 points  $(p, o, o')$ .



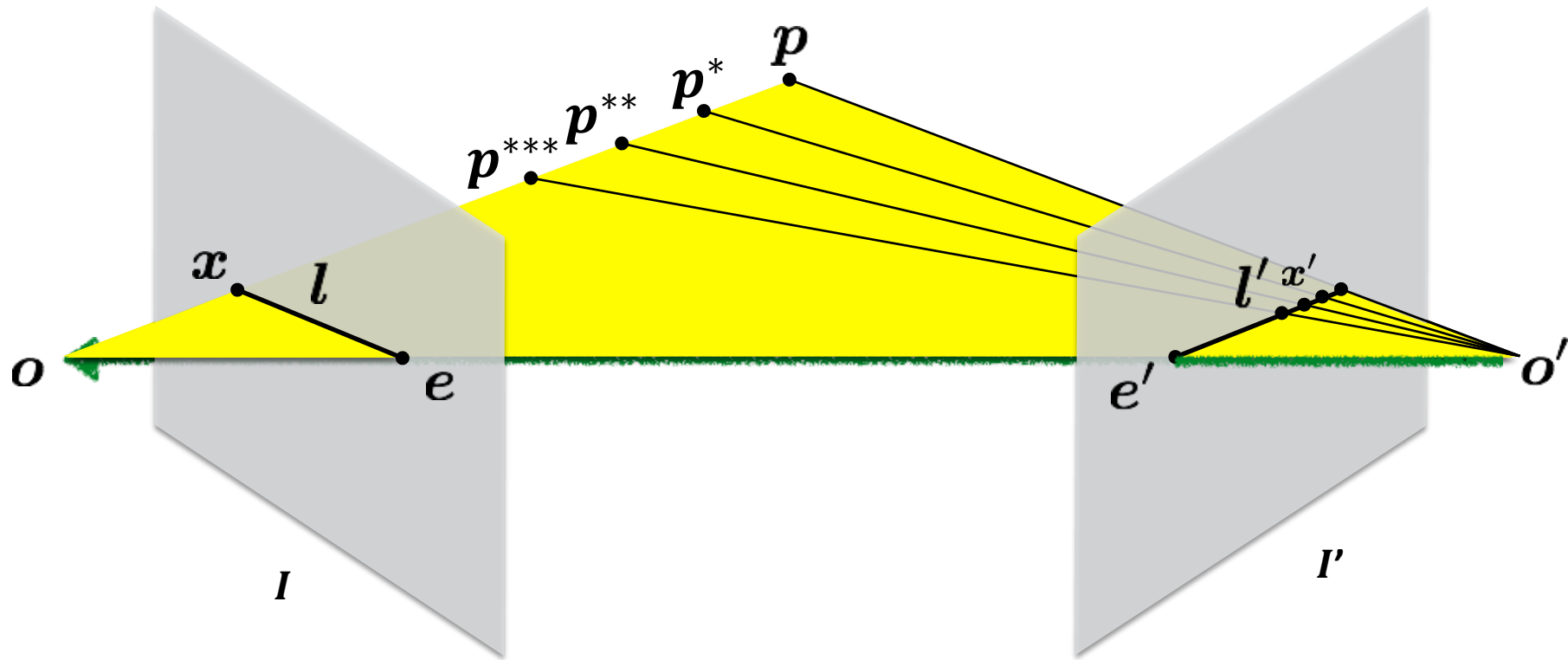
# Epipolar geometry

- **Epipolar line:** intersection of Epipolar plane and image plane.



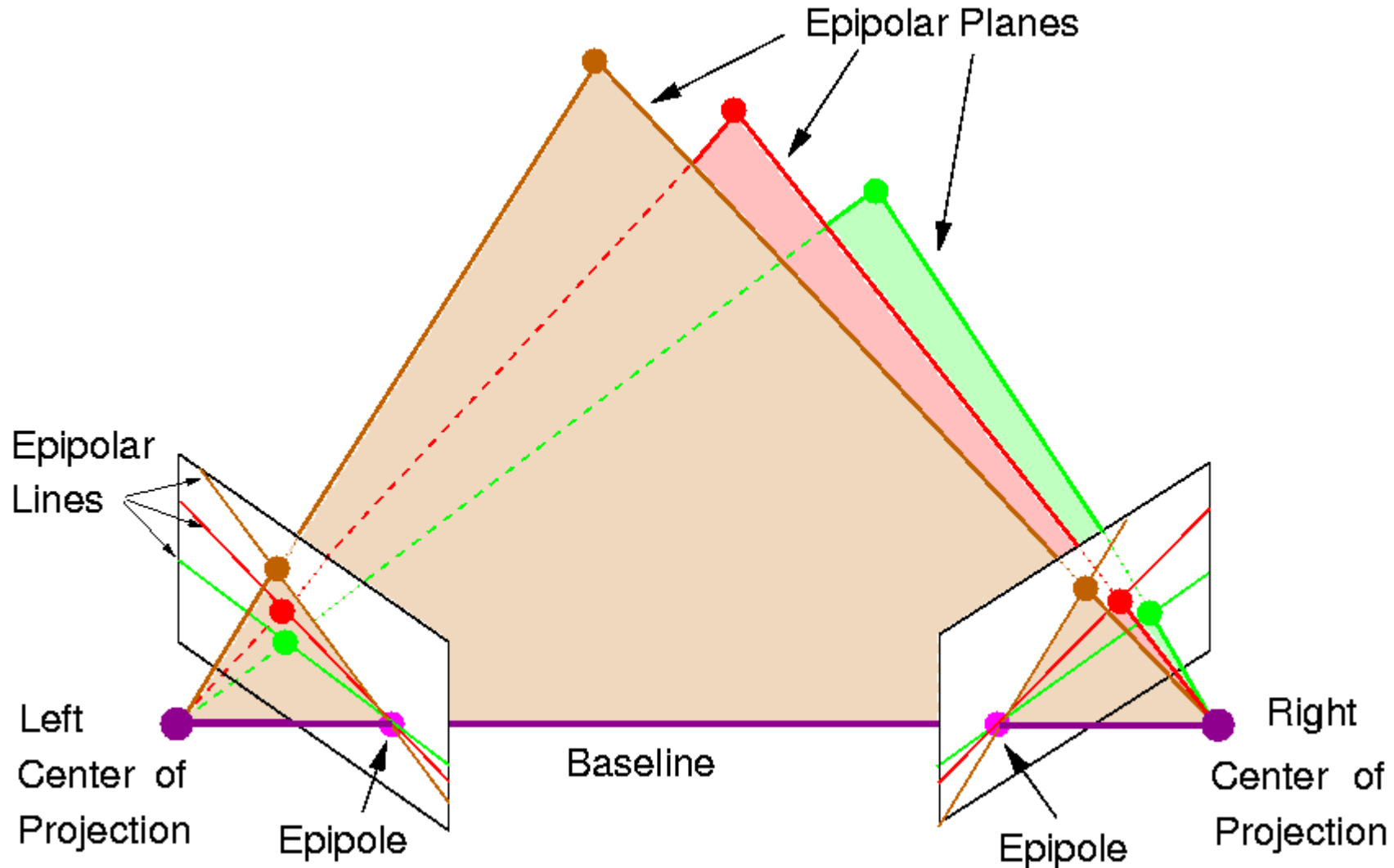
# Epipolar constraint

- **The epipolar constraint:** a point  $x$  in image  $I$  is mapped onto an epipolar line  $l'$  in image  $I'$ .
  - This happens since we don't know  $p$  in advance.



# Epipolar geometry

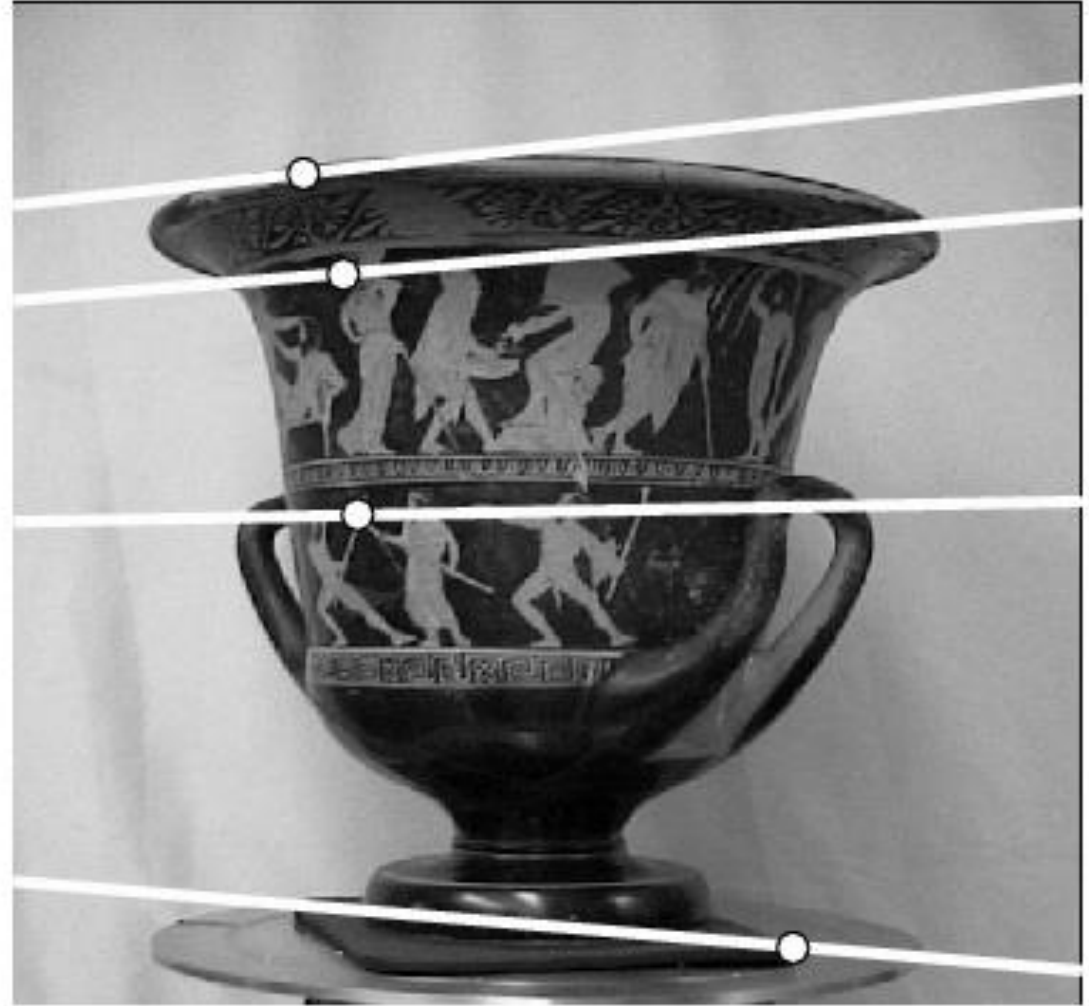
- Note: all epipolar lines pass through the epipole.





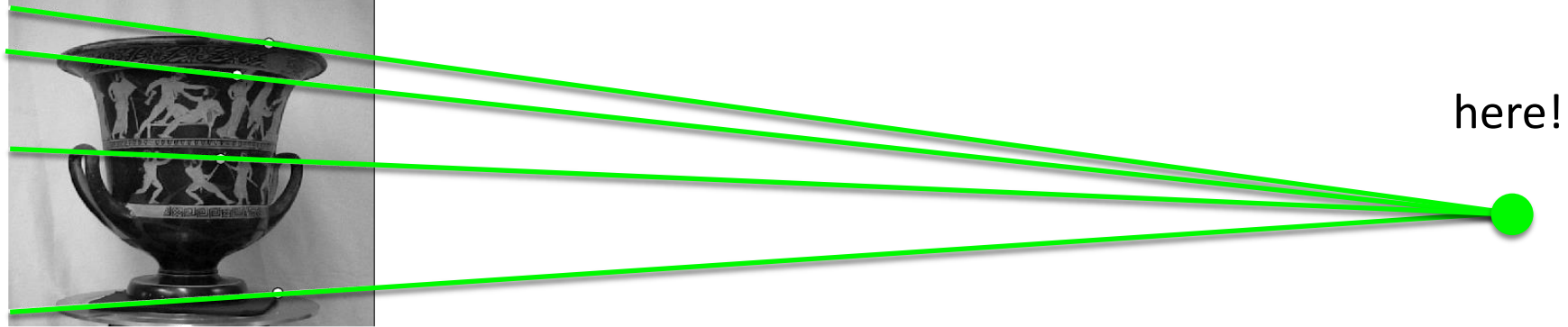
# Epipolar geometry

- Where is the epipole in this images?



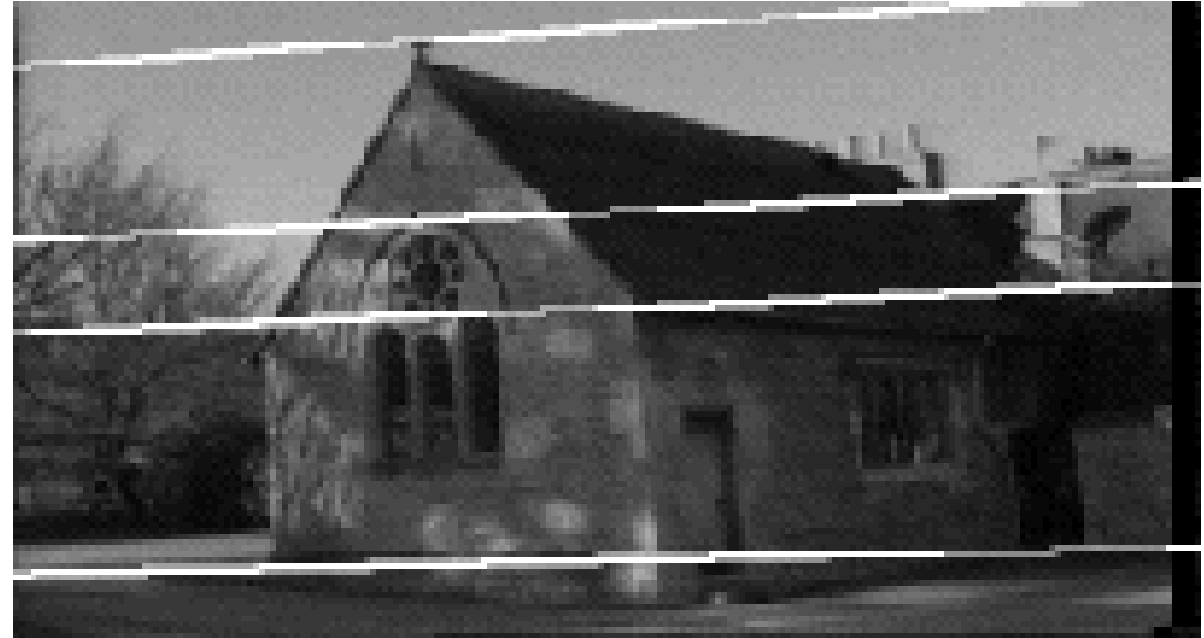
# Epipolar geometry

- Where is the epipole in this images? The epipole doesn't have to be inside the image!



# Epipolar geometry

- Where is the epipole in this image?



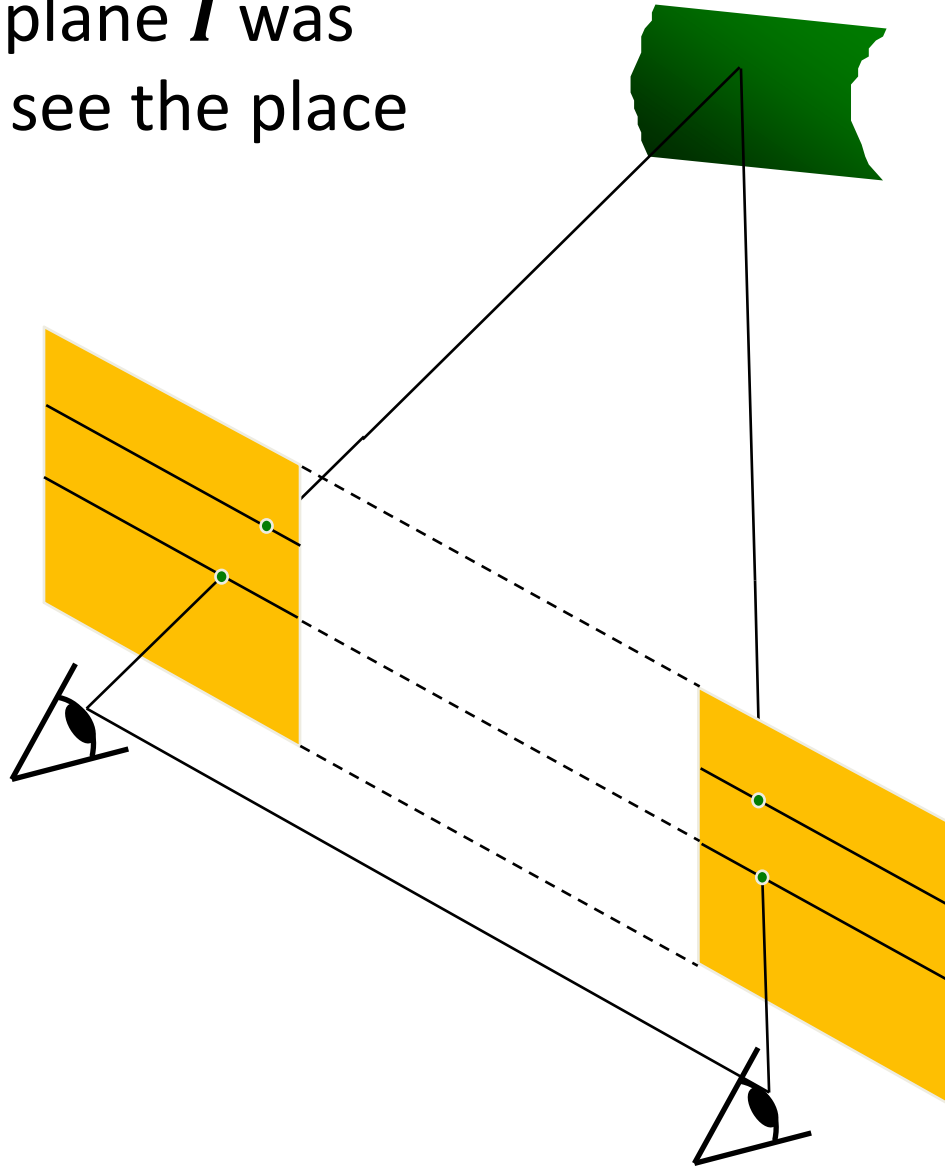
# Epipolar geometry

- Where is the epipole in this image? The epipolar lines doesn't converge since the baseline (translation) is parallel to the image plane!



# Epipolar geometry

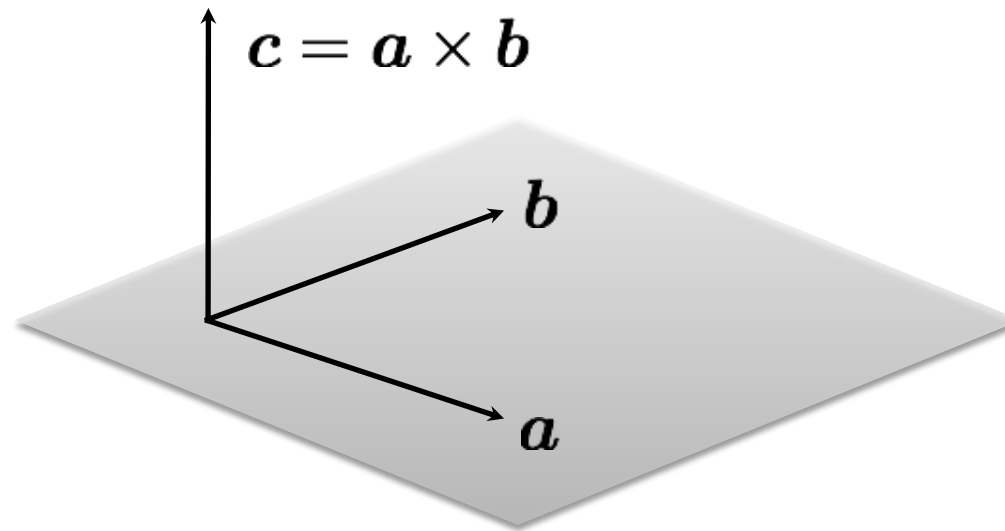
- Even if the image plane  $I$  was infinite, you can't see the place of  $o'$ .



# Contents

- Structure from motion
- Triangulation
- Stereo matching
- Camera rectification
- Epipolar geometry
  - **Essential matrix**
  - Fundamental matrix
  - Estimating the fundamental matrix
- Other 3D sensors

# Recall: Dot Product



$$c \cdot a = 0$$

$$c \cdot b = 0$$

Dot product of two orthogonal vectors is zero.

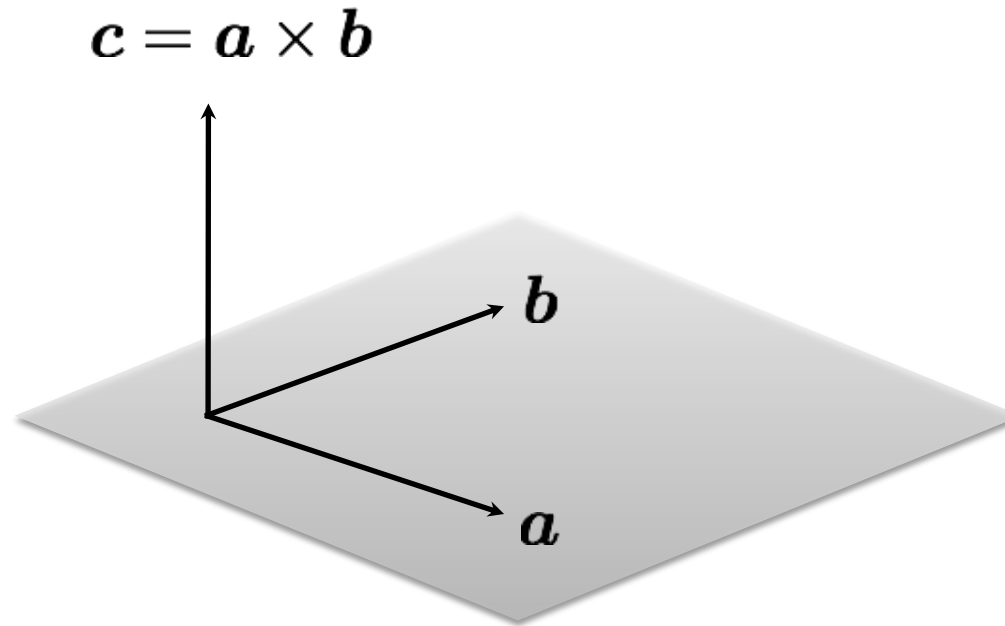
Dot product can also be written as vector multiplication [ $a, b$  are size  $3 \times 1$ ]:

$$a \cdot b = 0 \leftrightarrow a^T b = 0$$

# Recall: Cross Product

## Vector (cross) product

takes two vectors and returns a vector perpendicular to both



$$c \cdot a = 0$$

$$c \cdot b = 0$$



# Recall: Cross Product

$$\mathbf{a} \times \mathbf{b} = \begin{bmatrix} a_2b_3 - a_3b_2 \\ a_3b_1 - a_1b_3 \\ a_1b_2 - a_2b_1 \end{bmatrix}$$

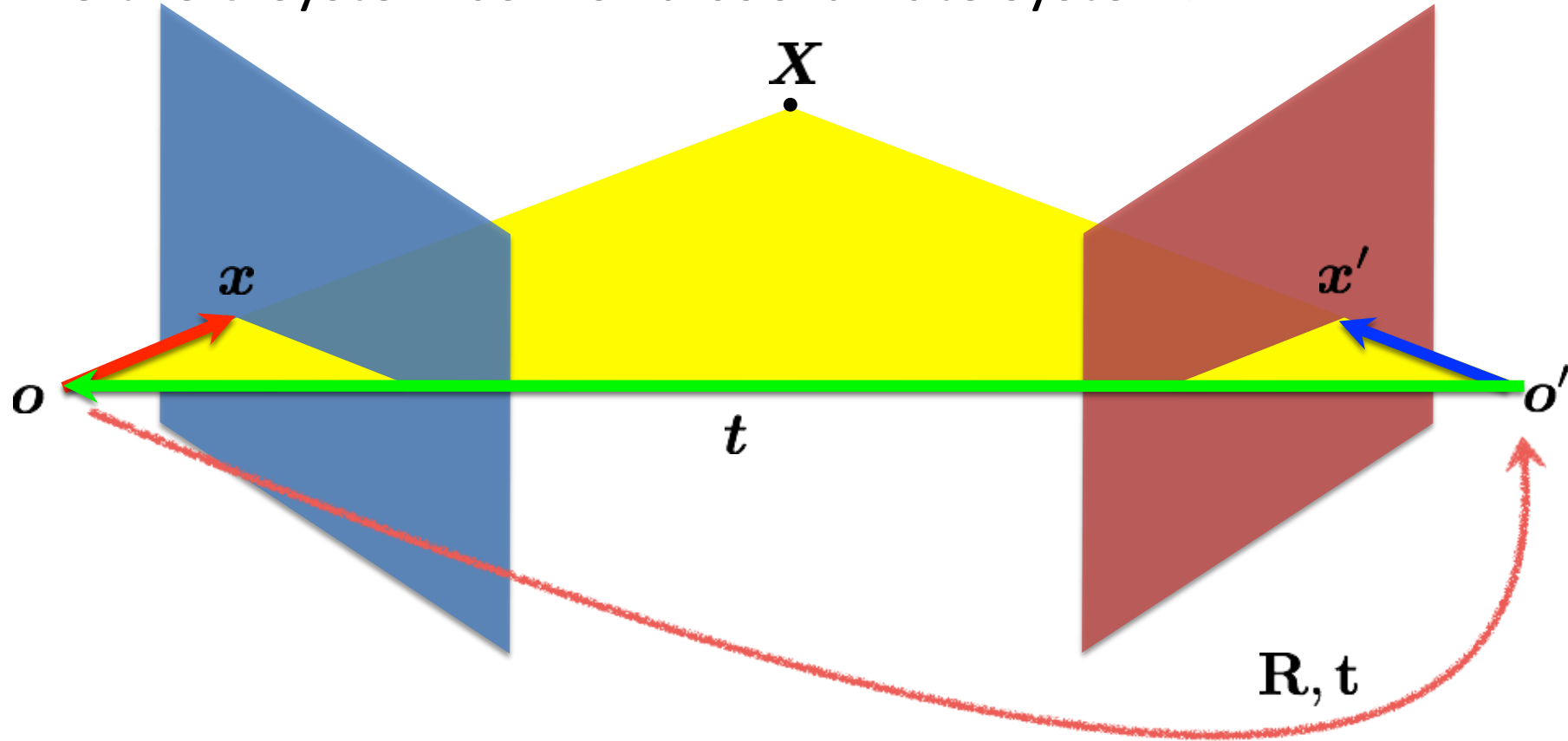
Can also be written as a matrix multiplication

$$\mathbf{a} \times \mathbf{b} = [\mathbf{a}]_{\times} \mathbf{b} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

**Skew symmetric**

# Building the essential matrix

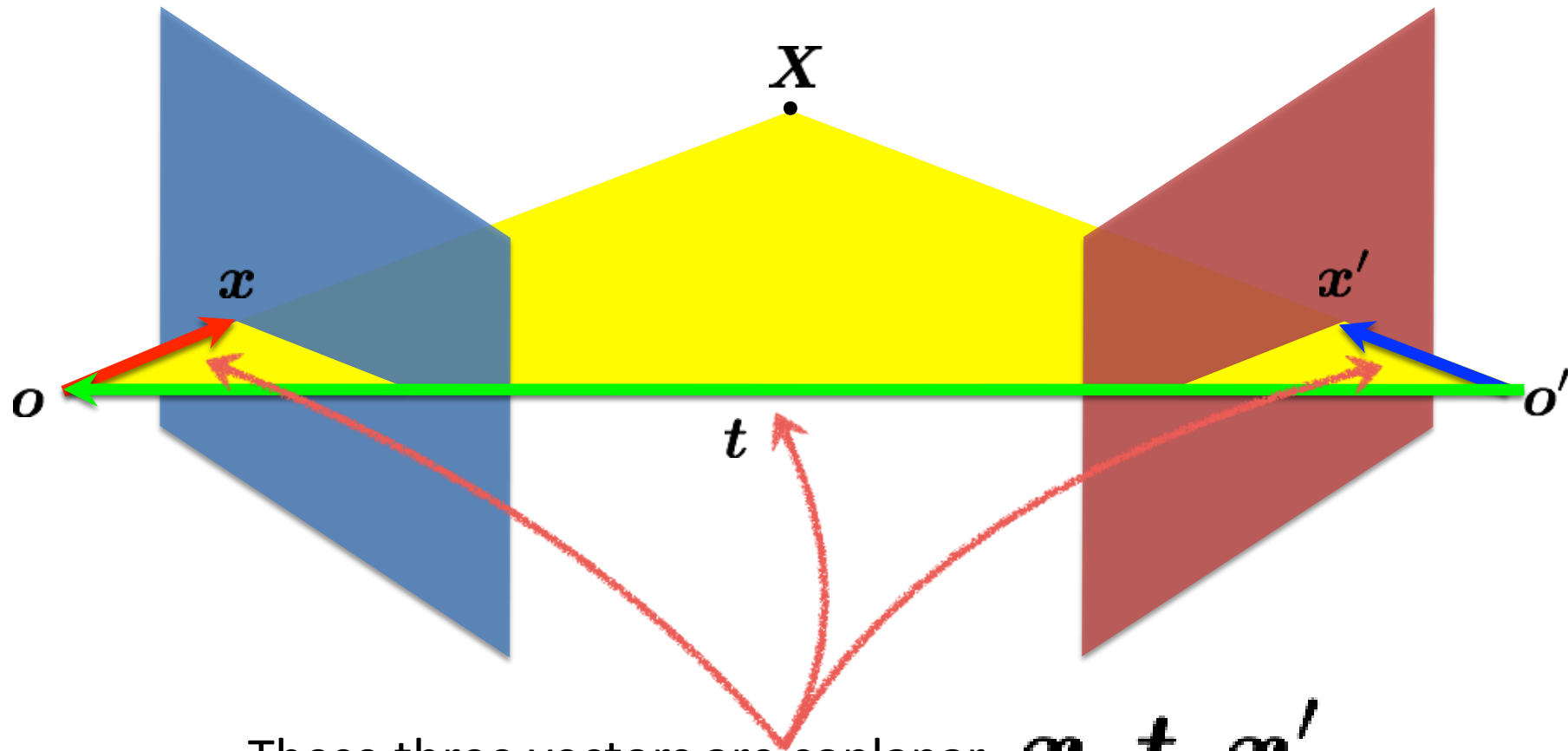
- Let's define the  $\mathbf{o}$  system as world coordinate system.



$$\mathbf{x}' = \mathbf{R}(\mathbf{x} - \mathbf{t})$$

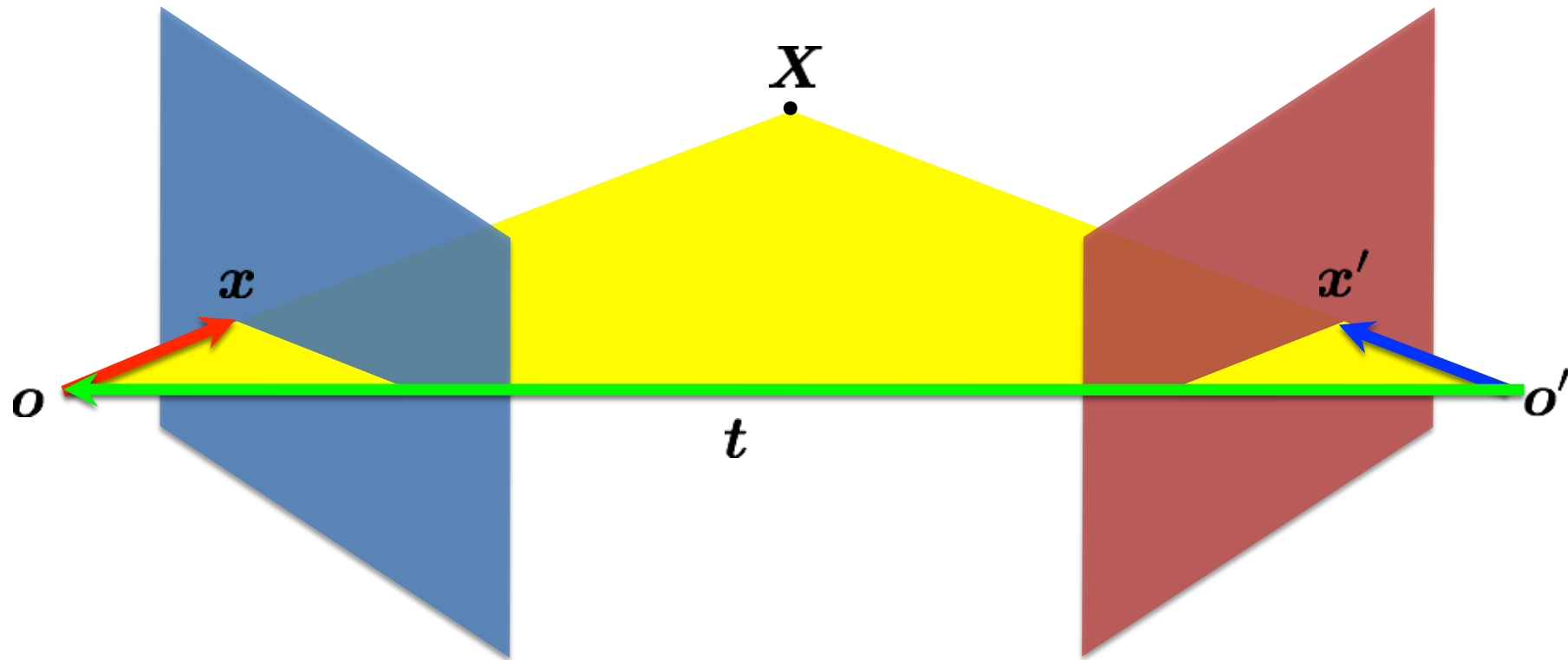
\* $\mathbf{t}$  here is  $\mathbf{c}$  from camera calibration class

# Building the essential matrix



These three vectors are coplanar  $x, t, x'$

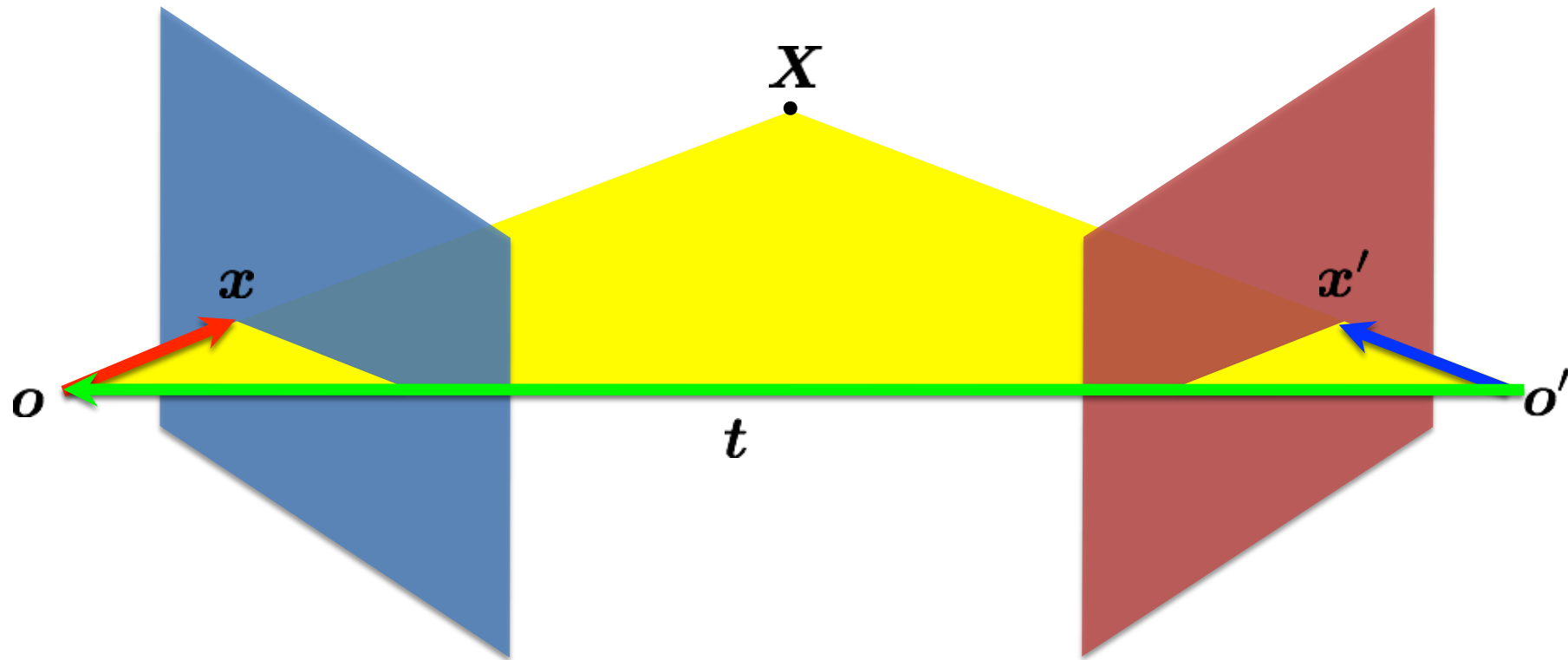
# Building the essential matrix



These three vectors are coplanar  $x, t, x'$

$$x^\top (t \times x) = ?$$

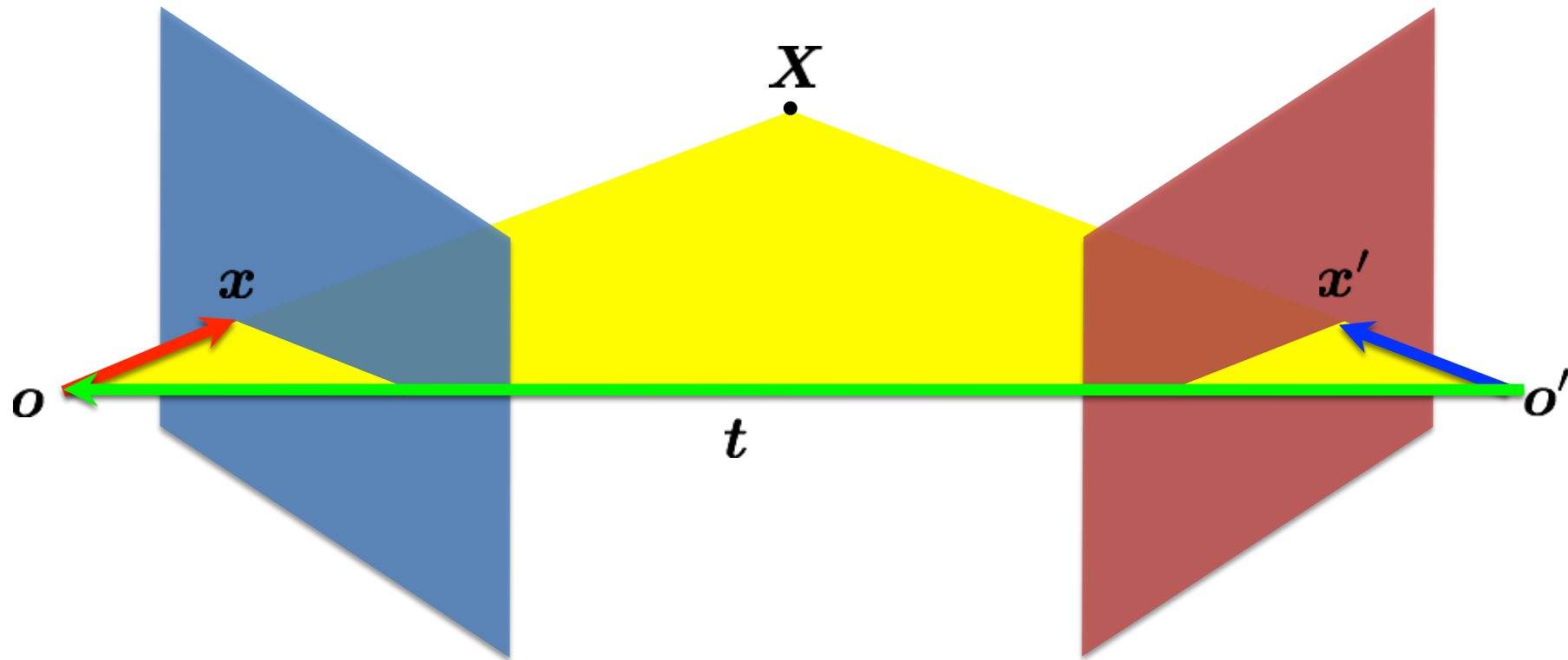
# Building the essential matrix



These three vectors are coplanar  $x, t, x'$

$$x^\top (t \times x) = 0$$

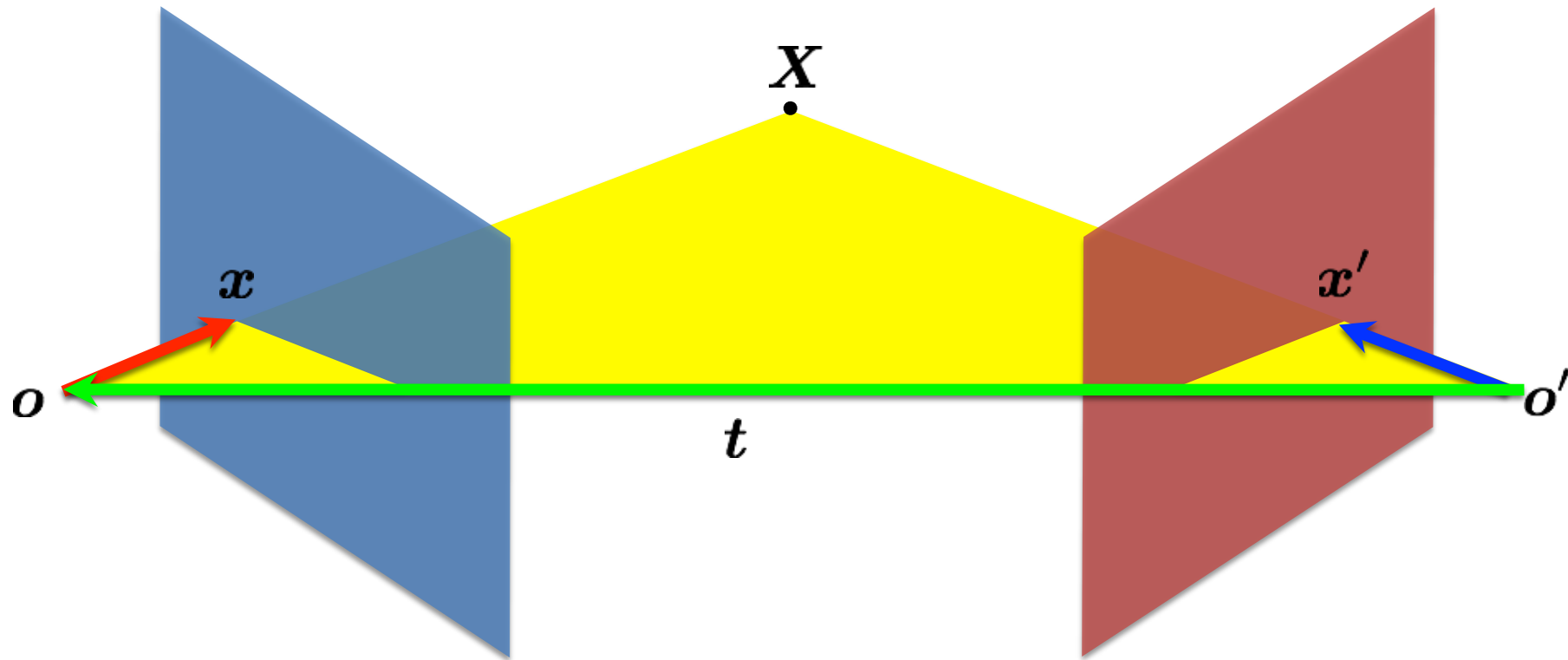
# Building the essential matrix



These three vectors are coplanar  $x, t, x'$

$$(x - t)^\top (t \times x) = ?$$

# Building the essential matrix



These three vectors are coplanar  $x, t, x'$

$$(x - t)^\top (t \times x) = 0$$

# Building the essential matrix

rigid motion

coplanarity

$$\begin{aligned} R^T &= R^{-1} \begin{cases} \mathbf{x}' = \mathbf{R}(\mathbf{x} - \mathbf{t}) \\ R^T \mathbf{x}' = \mathbf{x} - \mathbf{t} \\ \mathbf{x}'^T \mathbf{R} = (\mathbf{x} - \mathbf{t})^T \end{cases} \end{aligned}$$

$$(\mathbf{x} - \mathbf{t})^\top (\mathbf{t} \times \mathbf{x}) = 0$$

$$(\mathbf{x}'^\top \mathbf{R})(\mathbf{t} \times \mathbf{x}) = 0$$



# Building the essential matrix

rigid motion

coplanarity

$$\begin{aligned} R^T &= R^{-1} \begin{cases} \mathbf{x}' = \mathbf{R}(\mathbf{x} - \mathbf{t}) \\ R^T \mathbf{x}' = \mathbf{x} - \mathbf{t} \\ \mathbf{x}'^T R = (\mathbf{x} - \mathbf{t})^T \end{cases} \end{aligned}$$

$$(\mathbf{x} - \mathbf{t})^\top (\mathbf{t} \times \mathbf{x}) = 0$$

$$(\mathbf{x}'^\top \mathbf{R})(\mathbf{t} \times \mathbf{x}) = 0$$

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

# Building the essential matrix

rigid motion

coplanarity

$$\begin{aligned} R^T &= R^{-1} \begin{cases} \mathbf{x}' = \mathbf{R}(\mathbf{x} - \mathbf{t}) \\ R^T \mathbf{x}' = \mathbf{x} - \mathbf{t} \\ \mathbf{x}'^T R = (\mathbf{x} - \mathbf{t})^T \end{cases} \end{aligned}$$

$$(\mathbf{x} - \mathbf{t})^\top (\mathbf{t} \times \mathbf{x}) = 0$$

$$(\mathbf{x}'^\top \mathbf{R})(\mathbf{t} \times \mathbf{x}) = 0$$

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

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

# Building the essential matrix

rigid motion

coplanarity

$$\begin{aligned} R^T &= R^{-1} \begin{cases} \mathbf{x}' = \mathbf{R}(\mathbf{x} - \mathbf{t}) \\ R^T \mathbf{x}' = \mathbf{x} - \mathbf{t} \\ \mathbf{x}'^T R = (\mathbf{x} - \mathbf{t})^T \end{cases} \end{aligned}$$

$$(\mathbf{x} - \mathbf{t})^\top (\mathbf{t} \times \mathbf{x}) = 0$$

$$(\mathbf{x}'^\top \mathbf{R})(\mathbf{t} \times \mathbf{x}) = 0$$

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

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

$$\mathbf{x}'^\top \mathbf{E} \mathbf{x} = 0$$

# Building the essential matrix

rigid motion

coplanarity

$$\begin{aligned} R^T &= R^{-1} \begin{cases} \mathbf{x}' = \mathbf{R}(\mathbf{x} - \mathbf{t}) \\ R^T \mathbf{x}' = \mathbf{x} - \mathbf{t} \\ \mathbf{x}'^T R = (\mathbf{x} - \mathbf{t})^T \end{cases} \end{aligned}$$

$$(\mathbf{x} - \mathbf{t})^\top (\mathbf{t} \times \mathbf{x}) = 0$$

$$(\mathbf{x}'^\top \mathbf{R})(\mathbf{t} \times \mathbf{x}) = 0$$

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

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

$$\boxed{\mathbf{x}'^\top \mathbf{E} \mathbf{x} = 0}$$

**Essential Matrix**

$$E = R[\mathbf{t}]_x$$

# Contents

- Structure from motion
- Triangulation
- Stereo matching
- Camera rectification
- Epipolar geometry
  - Essential matrix
  - **Fundamental matrix**
  - Estimating the fundamental matrix
- Other 3D sensors

# Fundamental matrix

$$\hat{x}'^{\top} \mathbf{E} \hat{x} = 0$$

The Essential matrix operates on image points expressed in

**normalized coordinates**

(points have been aligned (normalized) to camera coordinates)

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

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

camera point                      image point

# Fundamental matrix

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

The Essential matrix operates on image points expressed in **normalized coordinates**  
(points have been aligned (normalized) to camera coordinates)

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

camera point                      image point

Writing out the epipolar constraint in terms of image coordinates

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

$$x'^T (\mathbf{K}'^{-T} \mathbf{E} \mathbf{K}^{-1}) x = 0$$

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

Fundamental  
Matrix

$$\mathbf{F} = \mathbf{K}'^{-T} \mathbf{E} \mathbf{K}^{-1}$$

# Contents

- Structure from motion
- Triangulation
- Stereo matching
- Camera rectification
- Epipolar geometry
  - Essential matrix
  - Fundamental matrix
  - **Estimating the fundamental matrix**
- Other 3D sensors



# Estimating $F$

- Given enough correspondence point between the two images, one can reconstruct the fundamental matrix  $F$ .
- If  $K_1, K_2$  are known, we can find  $E$ .
  - We can then decompose  $E$  to  $R, t$  between the two images (This part is out of scope for this lecture).
  - $t$  is found up to a scale in the estimation but it's easy to get a good measure of it with a ruler.



# Estimating F – 8-point algorithm

- The fundamental matrix F is defined by

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

for any pair of matches  $\mathbf{x}$  and  $\mathbf{x}'$  in two images.

- Let  $\mathbf{x}=(u,v,1)^T$  and  $\mathbf{x}'=(u',v',1)^T$ ,  
each match gives a linear equation

$$\mathbf{F} = \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix}$$

$$uu' f_{11} + vu' f_{12} + u' f_{13} + uv' f_{21} + vv' f_{22} + v' f_{23} + uf_{31} + vf_{32} + f_{33} = 0$$

# 8-point algorithm

- Like with homographies, instead of solving  $\mathbf{A}\mathbf{f} = 0$ , we seek  $\mathbf{f}$  to minimize  $\|\mathbf{A}\mathbf{f}\|$ , least eigenvector of  $\mathbf{A}^T \mathbf{A}$ .

$$\begin{bmatrix} u_1 u_1' & v_1 u_1' & u_1' & u_1 v_1' & v_1 v_1' & v_1' & u_1 & v_1 & 1 \\ u_2 u_2' & v_2 u_2' & u_2' & u_2 v_2' & v_2 v_2' & v_2' & u_2 & v_2 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ u_n u_n' & v_n u_n' & u_n' & u_n v_n' & v_n v_n' & v_n' & u_n & v_n & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = 0$$

## 8-point algorithm – Problem?

- $\mathbf{F}$  should have rank 2
- To enforce that  $\mathbf{F}$  is of rank 2,  $\mathbf{F}$  is replaced by  $\mathbf{F}'$  that minimizes  $\|\mathbf{F} - \mathbf{F}'\|$  subject to the rank constraint.

- This is achieved by SVD. Let  $\mathbf{F} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  where

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix}, \text{ let } \mathbf{\Sigma}' = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

then  $\mathbf{F}' = \mathbf{U}\mathbf{\Sigma}'\mathbf{V}^T$  is the solution.

# 8-point algorithm

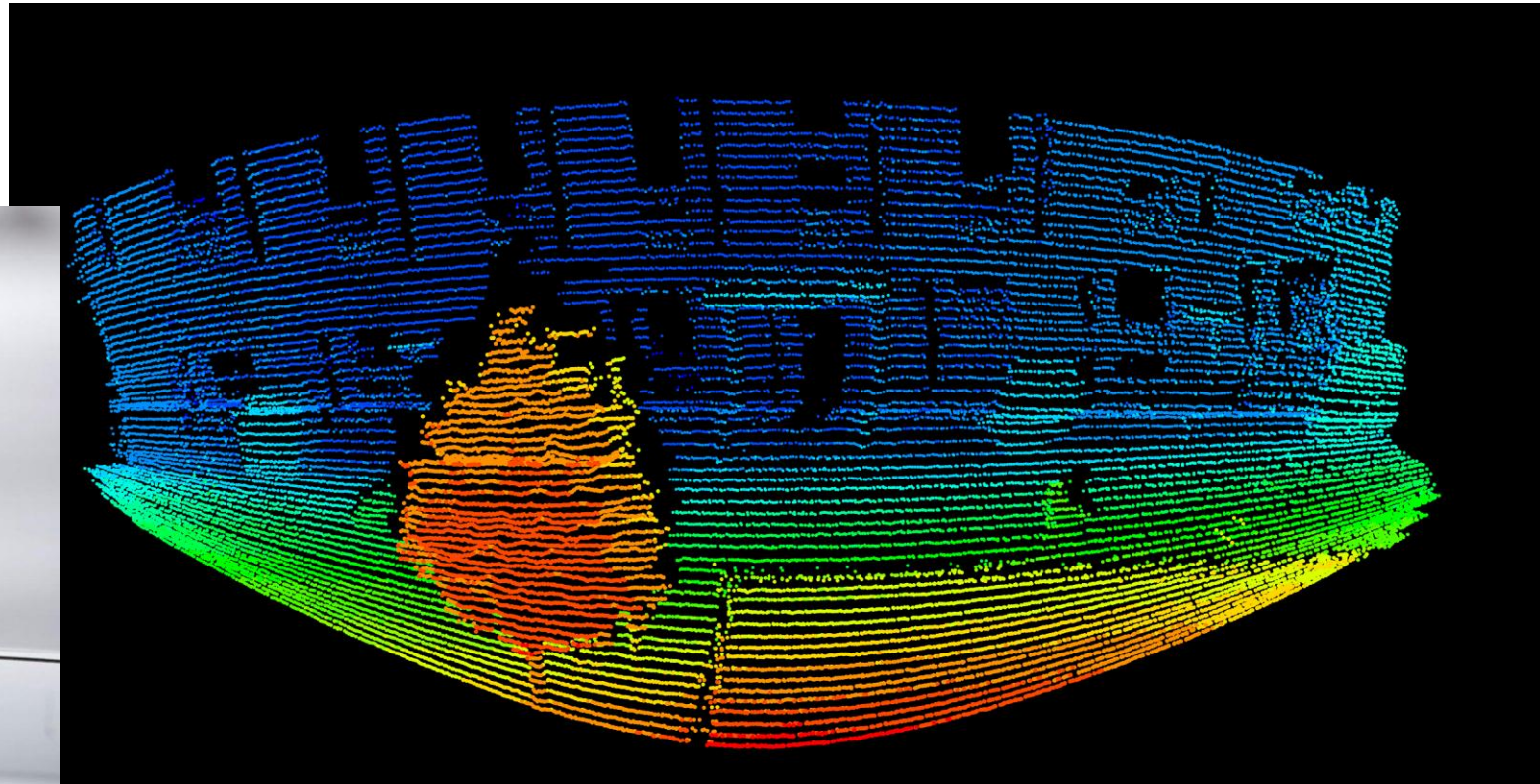
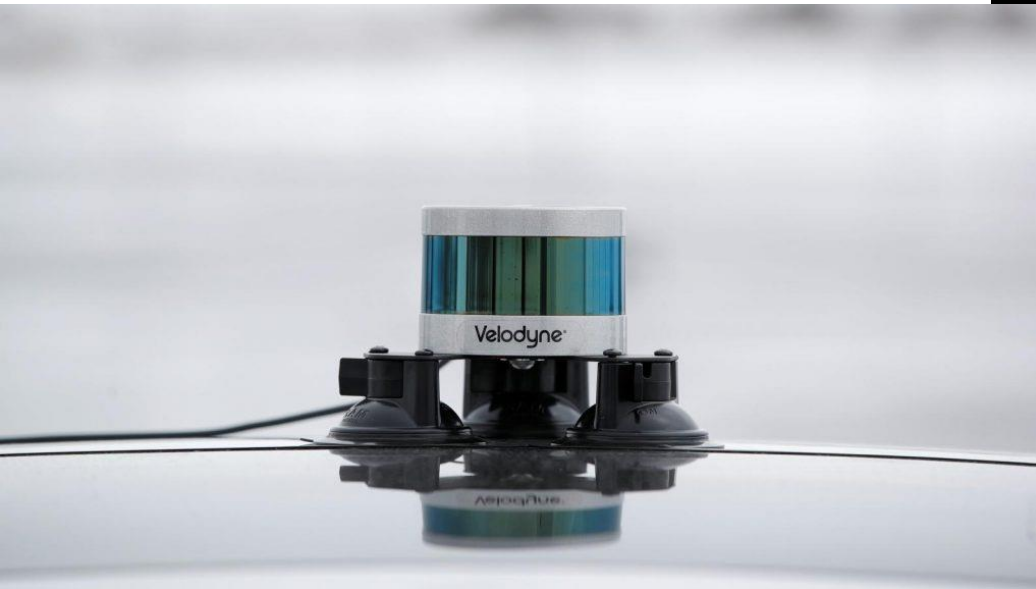
- Pros: it is linear, easy to implement and fast
- Cons: susceptible to noise.
  - Solutions: (all out of scope)
    - normalized 8 points algorithm.
    - 7 points algorithm.
    - Finding  $K, K'$  with single camera intrinsics calibration and then search for  $E$  (only 5 DOFs instead of 8/7).

# Contents

- Structure from motion
- Triangulation
- Stereo matching
- Camera rectification
- Epipolar geometry
  - Essential matrix
  - Fundamental matrix
  - Estimating the fundamental matrix
- **Other 3D sensors**

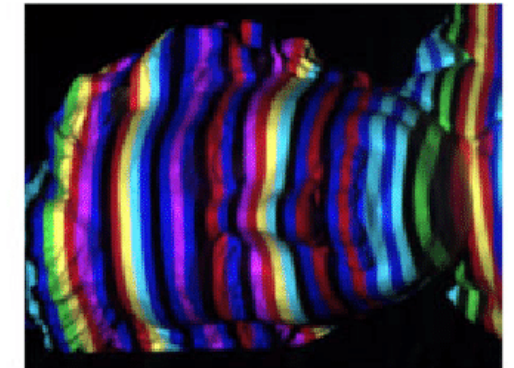
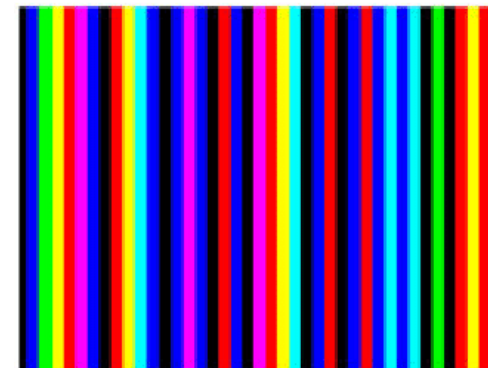
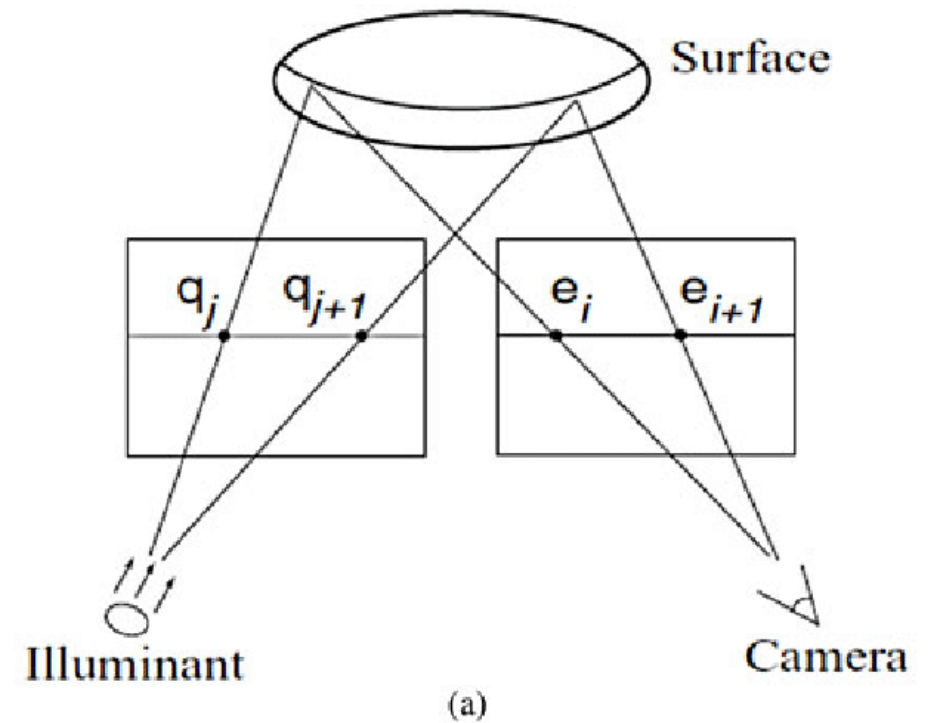
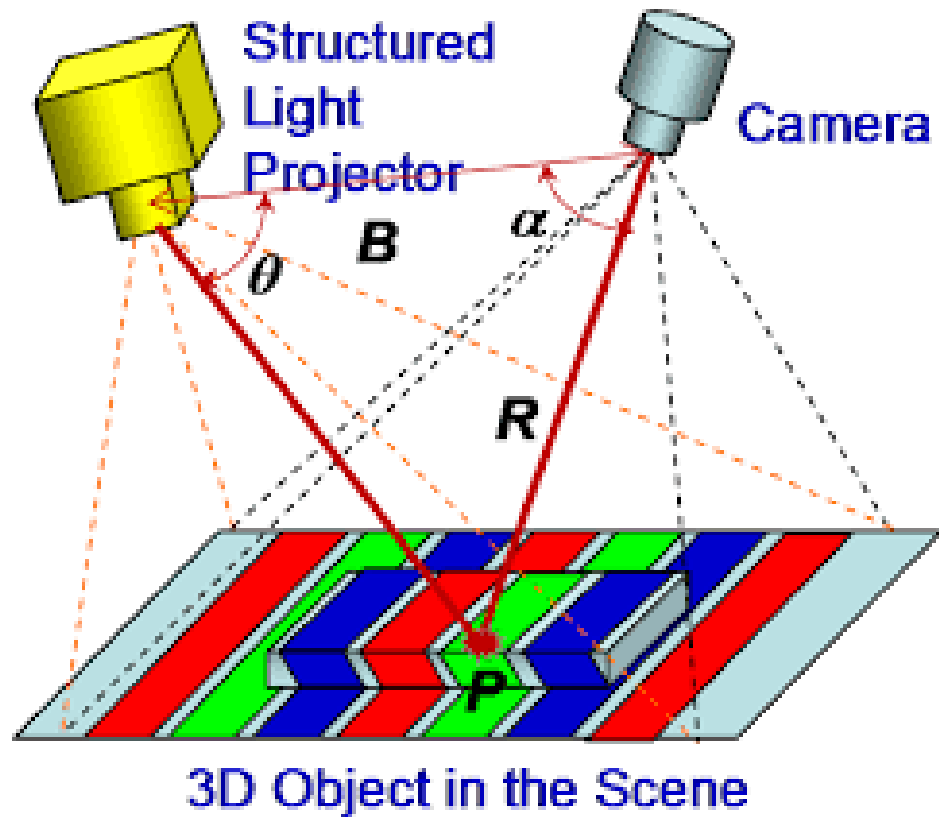
# Other types of 3D sensors

- **LIDAR**, which stands for **L**ight **D**etection and **R**anging (or light radar), is a remote sensing method that uses light in the form of a pulsed laser to measure ranges.
- Most known: velodyne projector.



# Other types of 3D sensors

- Structured light





# Other types of 3D sensors

- Coded light
- Realsense SR305
- <https://www.youtube.com/watch?v=PluL7WTIKrM>



# Other types of 3D sensors

- *Light Coding*
- *Used in Kinect v1- Kinect for xbox 360.*
- *Iphone x front camera*

