

# Embodied Vision

Image Formation and Multi-view Geometry

Tsung-Wei Ke

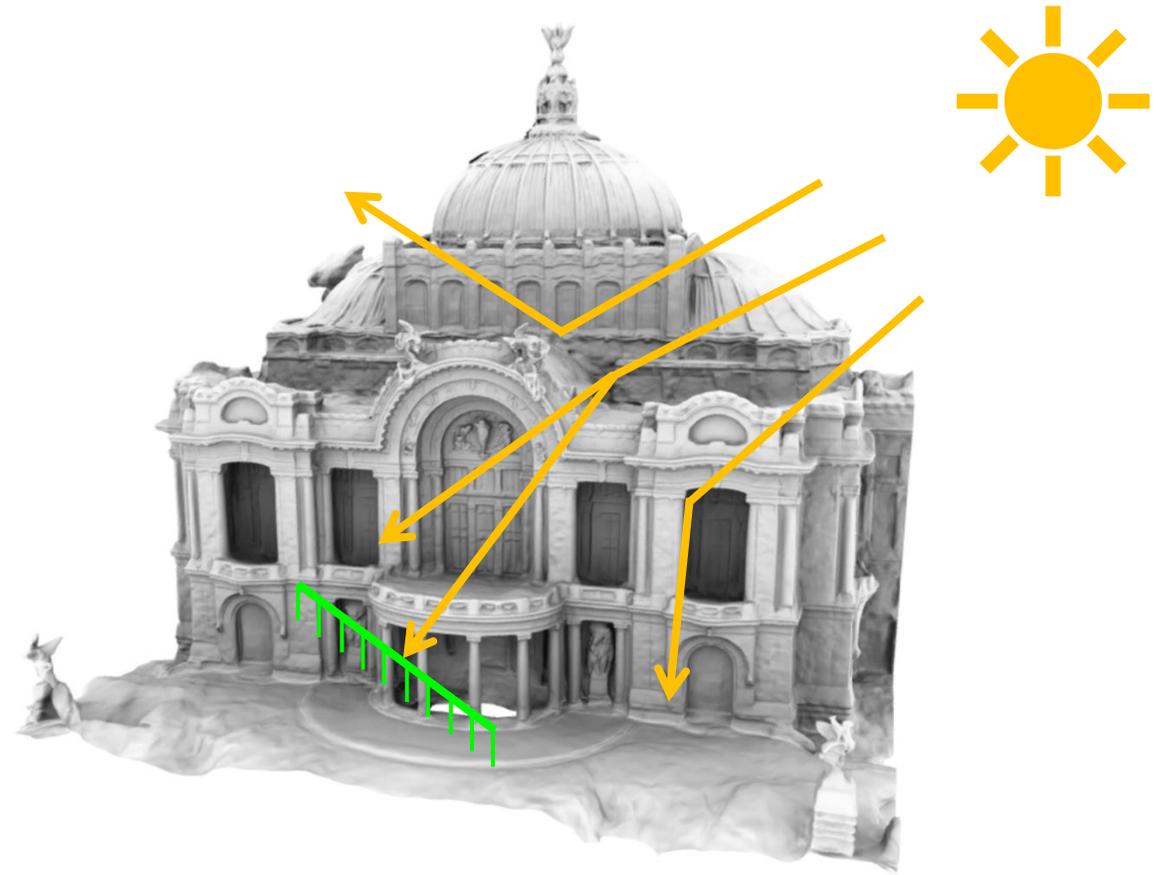
Spring 2026



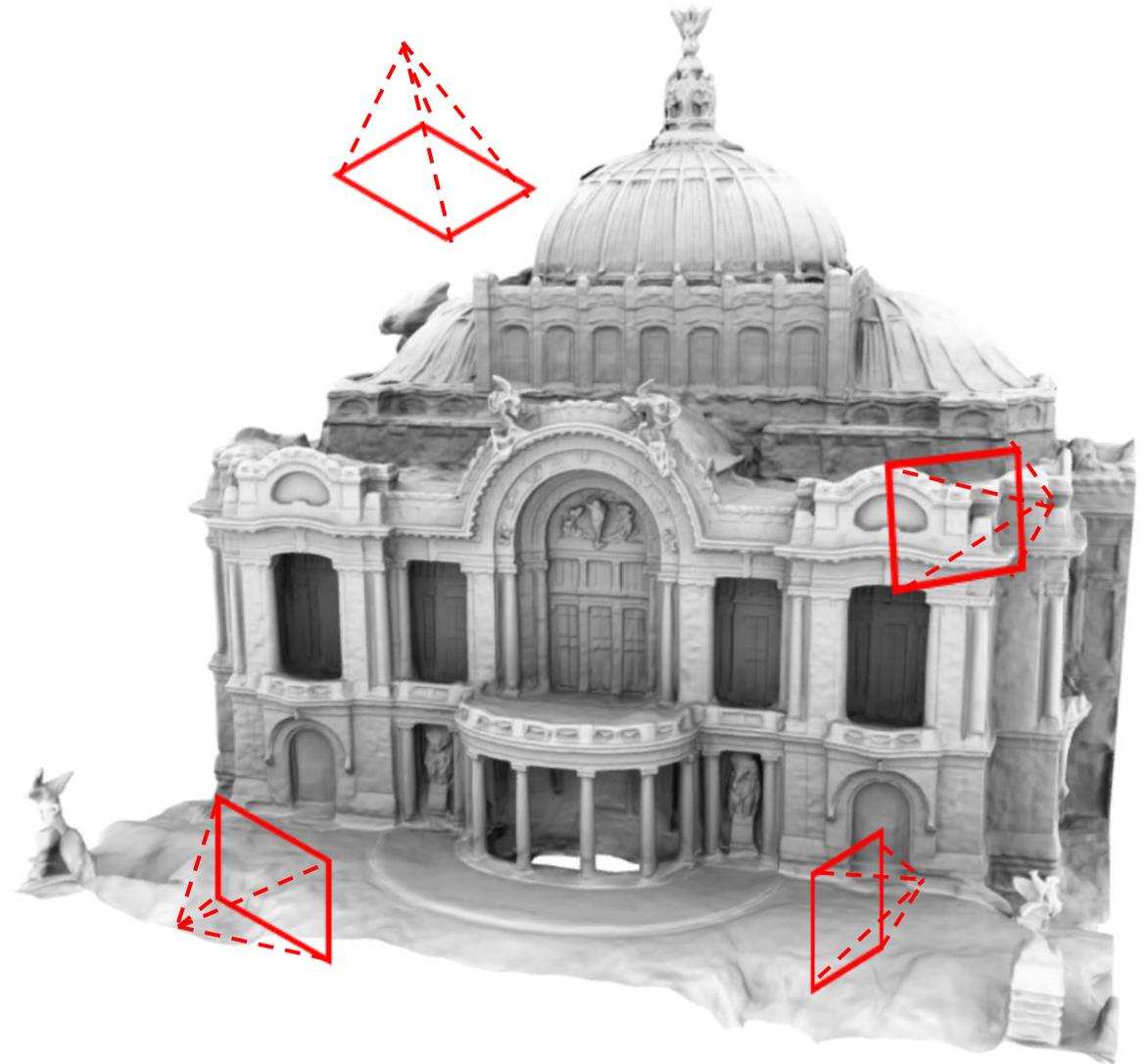
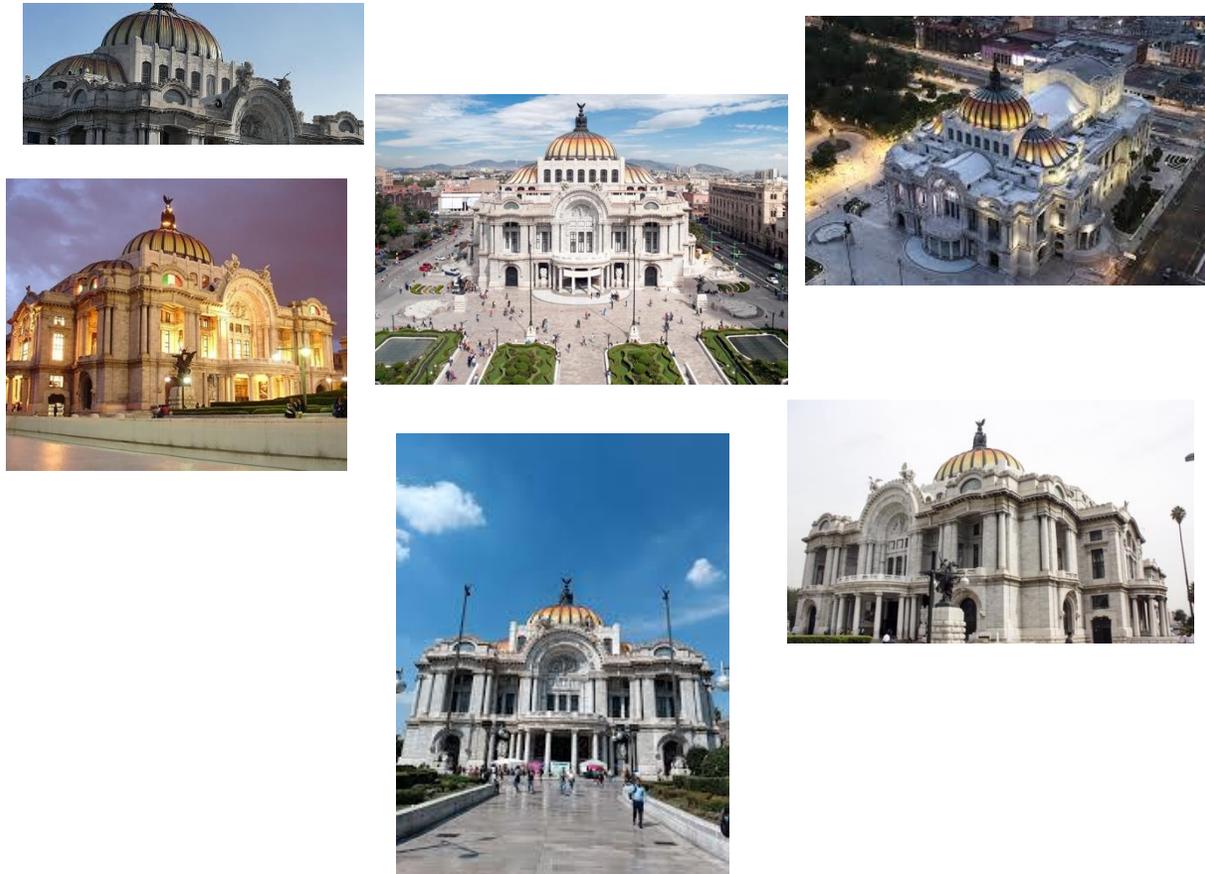
# Disclaimer

- This lecture borrows contents heavily from
  - [Machine Learning for Inverse Graphics](#) by Vincent Sitzmann at MIT
  - [Computer Vision](#) by Andreas Geiger at University Tübingen
  - [Intro to Computer Vision and Computational Photography](#) by Alyosha Efros at UC Berkeley

# Image Formation: How Do We See 2D Images from the 3D World?



# Multi-view Geometry: How Do We Reconstruct the 3D World from 2D Images?



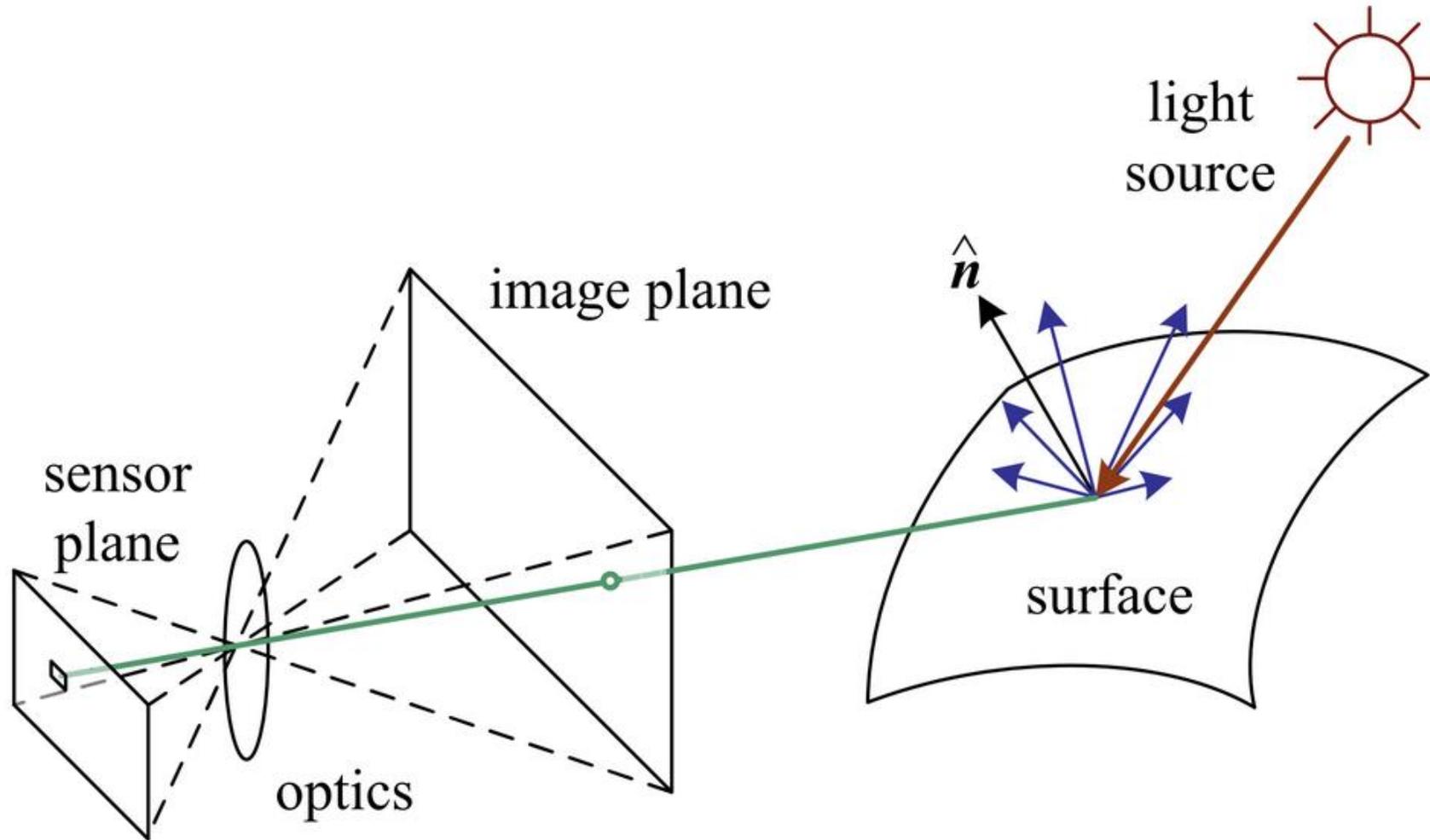
# Content

- Image Formation
  - Pinhole camera model
  - 2D transformation
  - 3D transformation
- Geometry Reconstruction
  - Reconstruction with depths
  - Multi-view geometry reconstruction with camera motions and correspondence
  - Epipolar geometry reconstruction

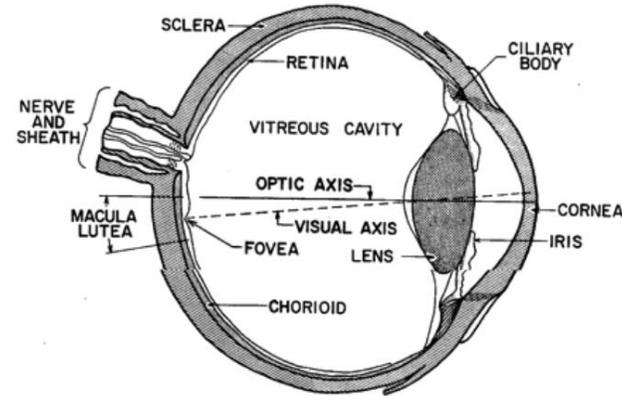
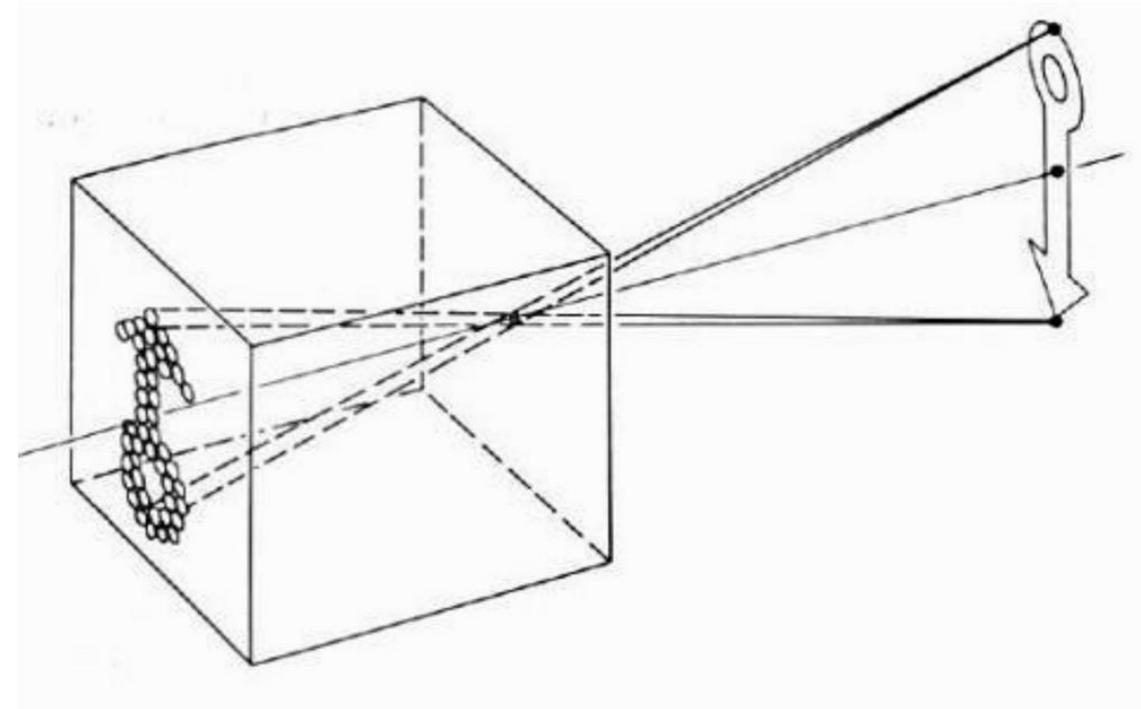
# Content

- Image Formation
  - Pinhole camera model
  - 2D transformation
  - 3D transformation
- Geometry Reconstruction
  - Reconstruction with depths
  - Multi-view geometry reconstruction with camera motions and correspondence
  - Epipolar geometry reconstruction

# 2D Image Formation from the 3D World: A simplified model of photometric image formation



# Perspective Projection



Animal eye: a long time ago



Gemma Frisius, 1558

Pinhole Perspective Projection:  
Brunelleschi, 15th Century



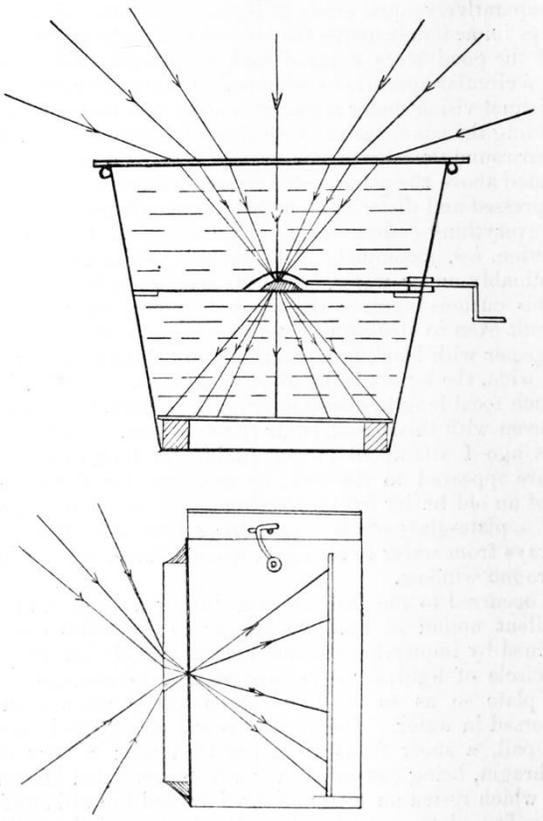
Photographic camera:  
Niepce, 1816

Image source:

- US Navy Manual of Basic Optics and Optical Instruments, prepared by Bureau of Naval Personnel. Reprinted by Dover Publications, Inc., 1969.
- Remote Sensing. A.L. Nowicki, "Stereoscopy." Manual of Photogrammetry, Thompson, Radlinski, and Speert (eds.), third edition, 1966

# Other Types of Camera Projection

Ellipsoidal projection



Parallel projection

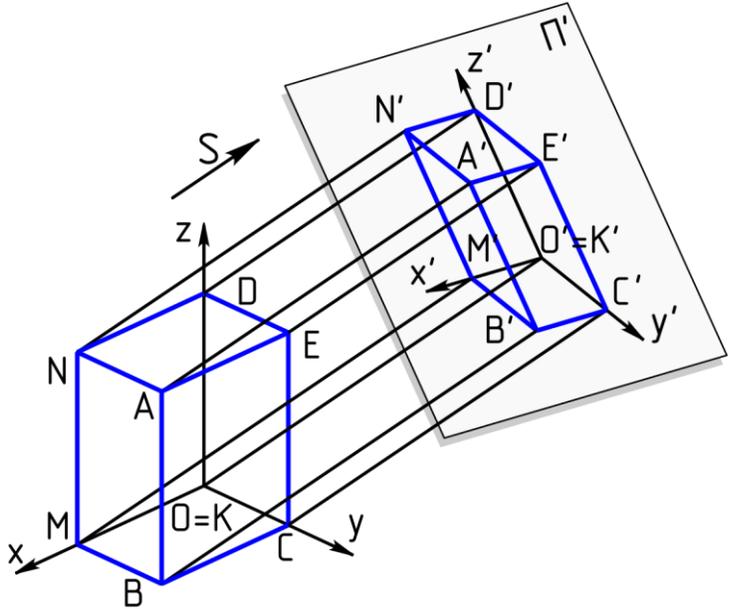
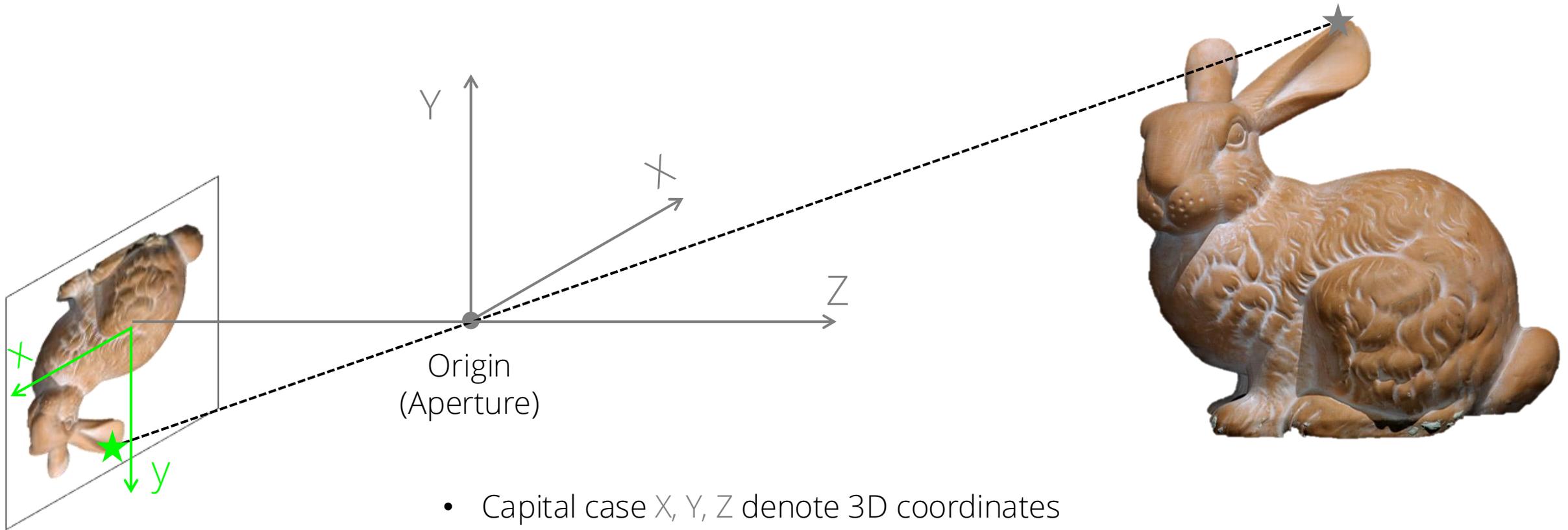


Image source Robert W. Wood, "Fish-Eye Views, and Vision under Water" (August 1906)

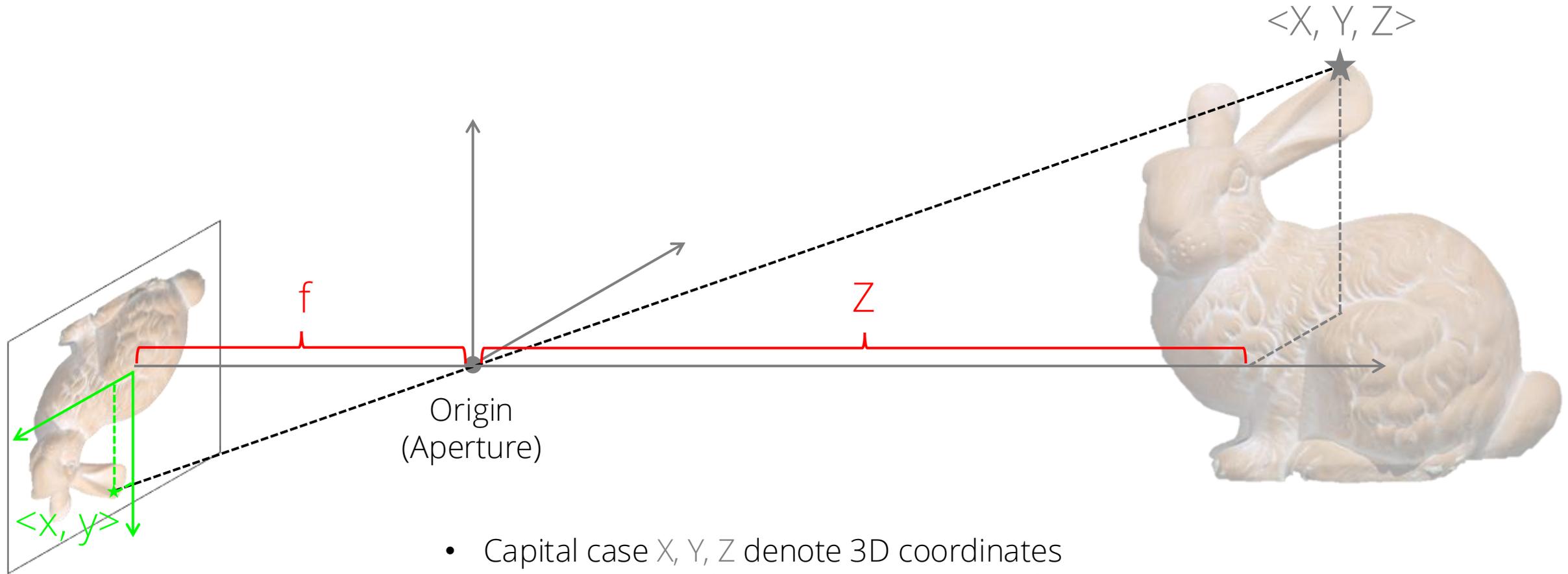
[https://en.wikipedia.org/wiki/Parallel\\_projection#/media/File:Axonometric\\_projection.svg](https://en.wikipedia.org/wiki/Parallel_projection#/media/File:Axonometric_projection.svg)

# Perspective Projection



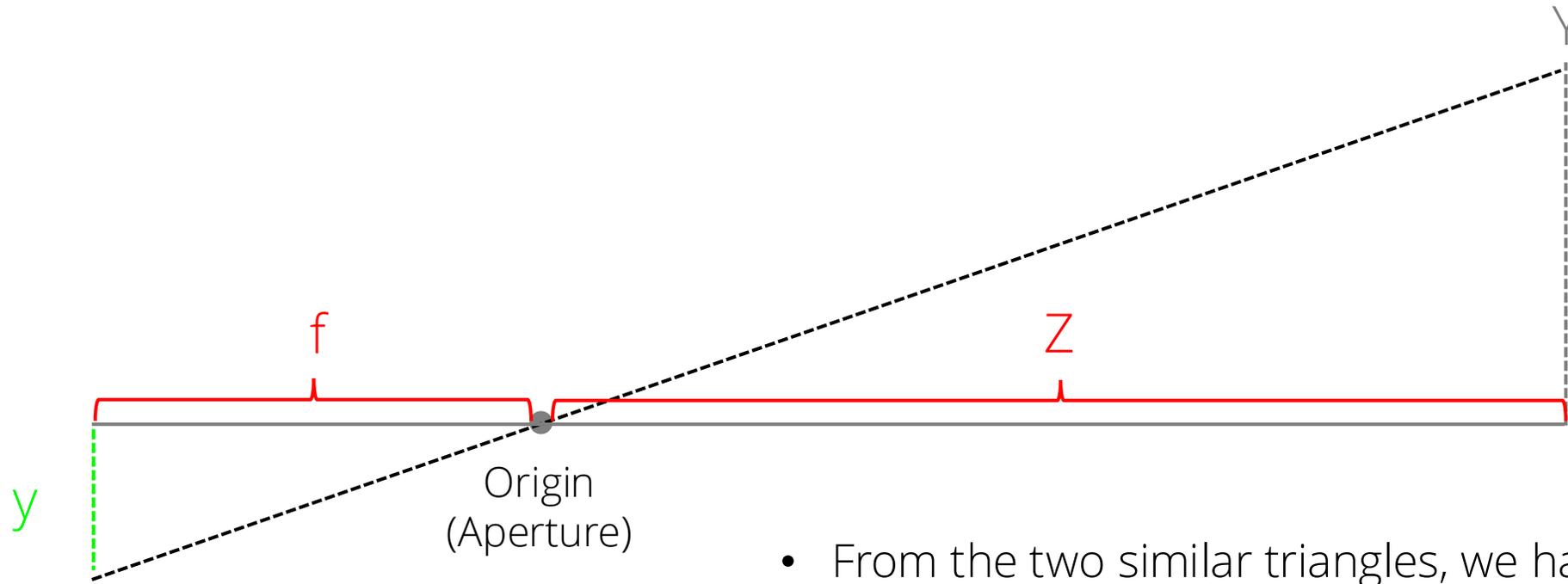
- Capital case  $X, Y, Z$  denote 3D coordinates
- Lower case  $x, y$  denotes 2D Image coordinates
- How to derive the relation between  $\langle x, y \rangle$  and  $\langle X, Y, Z \rangle$  ?

# Perspective Projection



- Capital case  $X, Y, Z$  denote 3D coordinates
- Lower case  $x, y$  denotes 2D Image coordinates
- $f$  denotes the camera focal length

# Perspective Projection



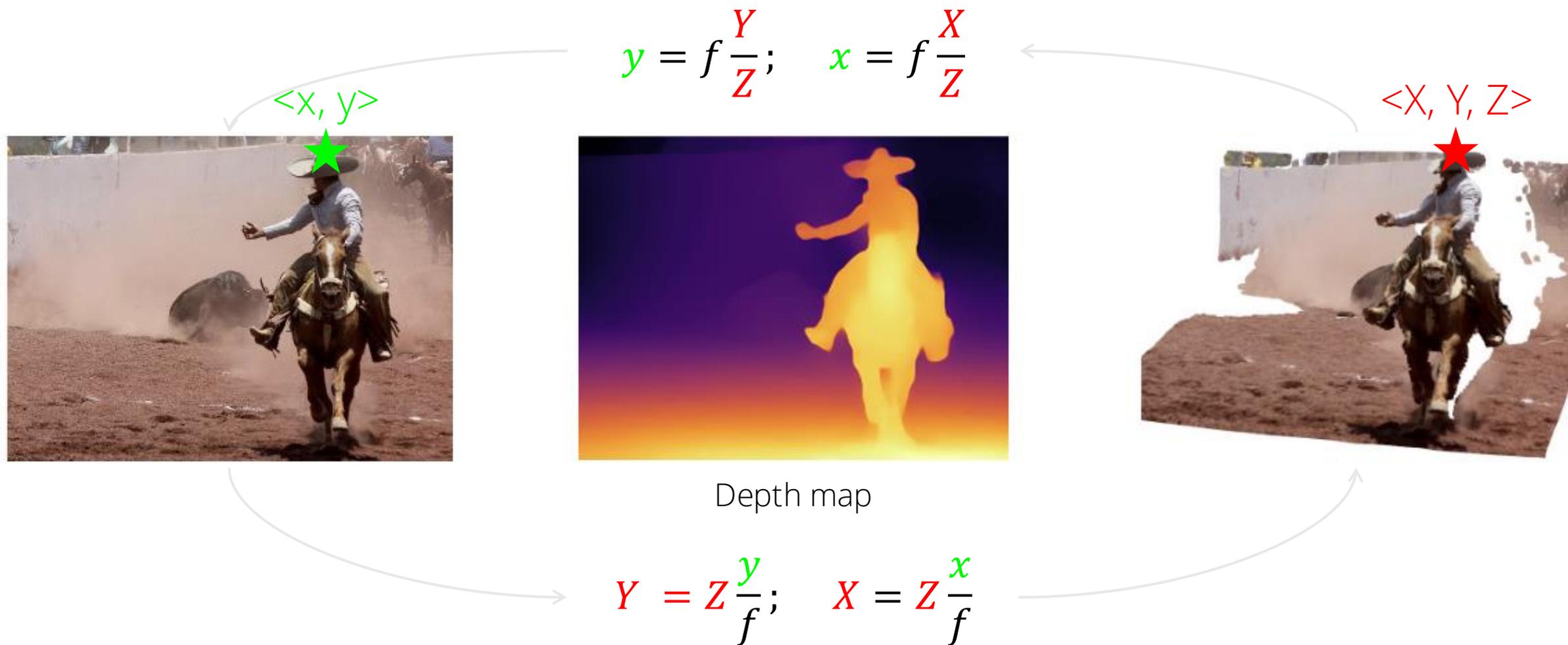
- From the two similar triangles, we have:

$$\frac{y}{f} = \frac{Y}{Z} \Rightarrow y = f \frac{Y}{Z}$$

- Likewise

$$\frac{x}{f} = \frac{X}{Z} \Rightarrow x = f \frac{X}{Z}$$

# We Can Transfer Between 2D Pixel Coordinates and 3D Coordinates, if $f$ and $Z$ are Known



# Information Loss by Perspective Projection



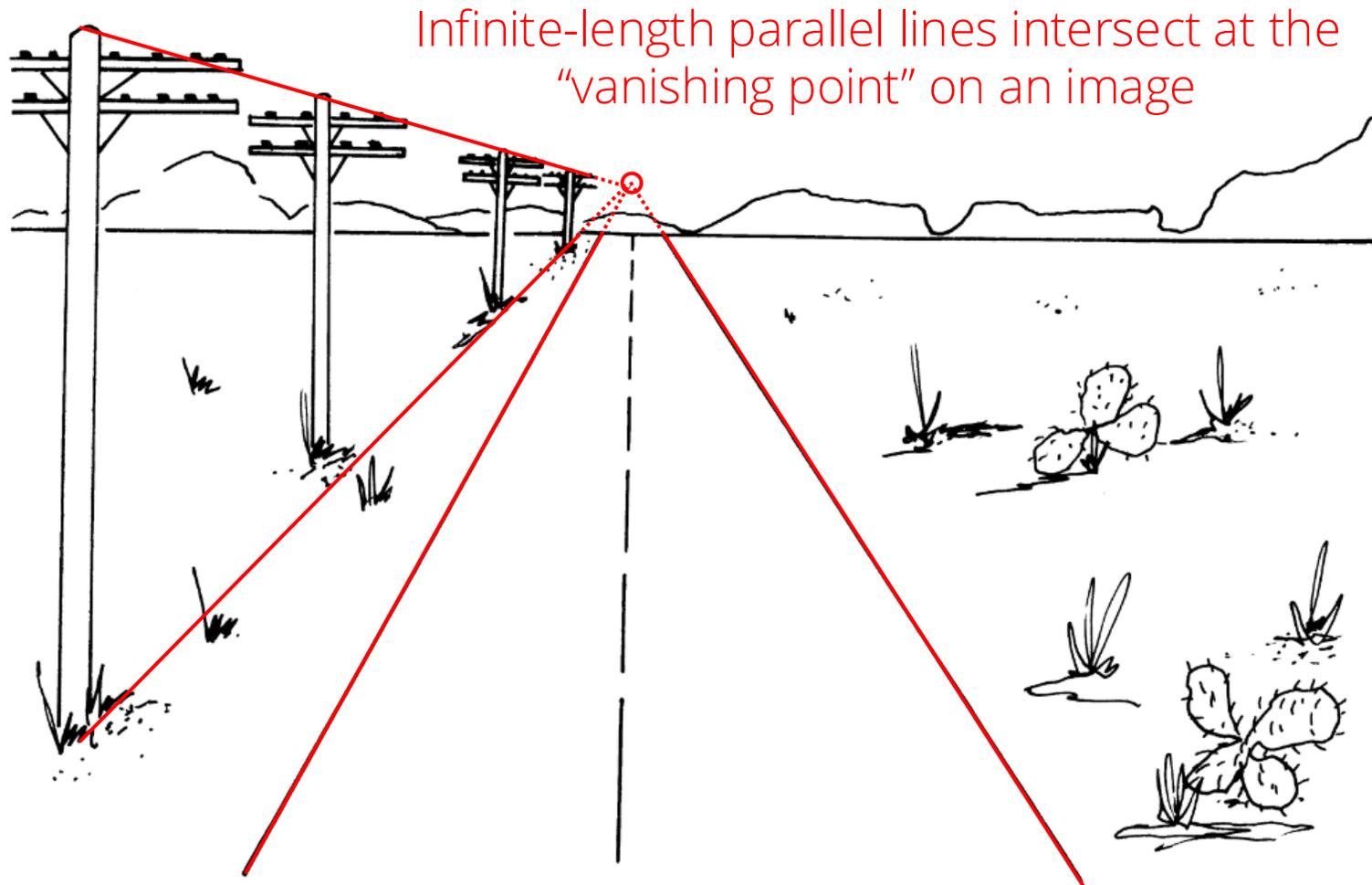
Slide credit A. Efros  
Image source: Fooling the eye Making of 3D sidewalk art:  
<http://www.youtube.com/watch?v=3SNYtd0Ayt0>

# Information Loss by Perspective Projection

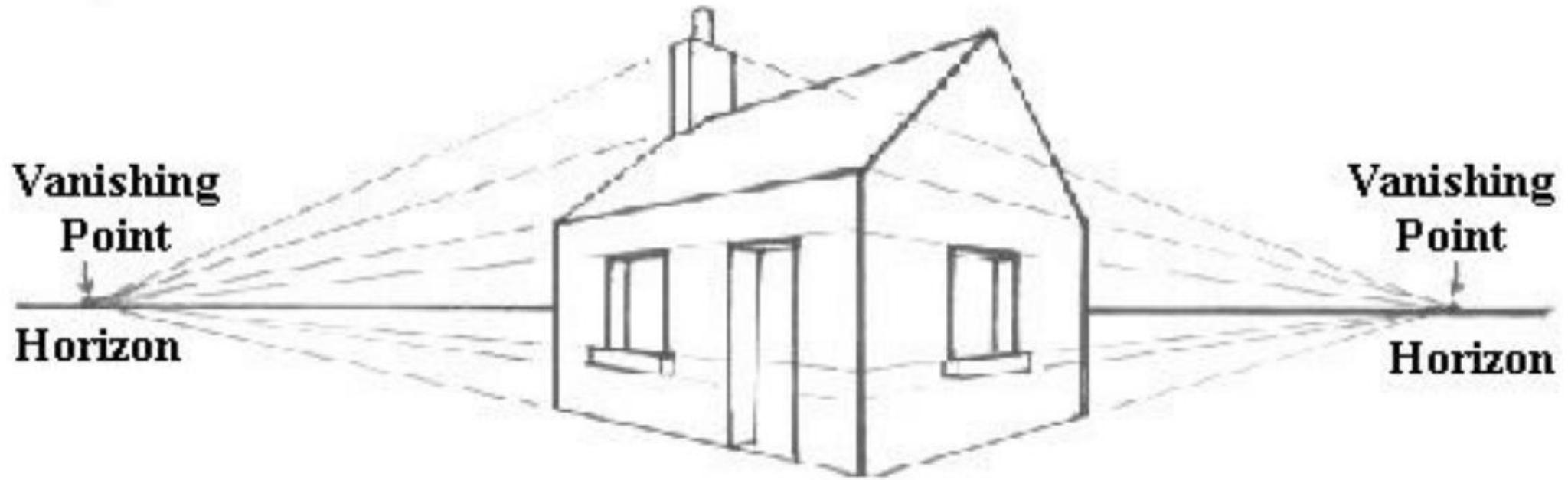


Slide credit A. Efros  
Image source: Fooling the eye Making of 3D sidewalk art:  
<http://www.youtube.com/watch?v=3SNYtd0Ayt0>

# Angles Are Not Preserved by Perspective Projection



# Each Family of Parallel Lines has its Own Vanishing Point



# Proof

- Line in 3D space:

$$X(t) = X_0 + at$$

$$Y(t) = Y_0 + bt$$

$$Z(t) = Z_0 + ct$$

- Perspective projection:

$$y(t) = f \frac{Y(t)}{Z(t)} = f \frac{Y_0 + bt}{Z_0 + ct}$$

$$x(t) = f \frac{X(t)}{Z(t)} = f \frac{X_0 + at}{Z_0 + ct}$$

- Projective lines when  $t \rightarrow \infty$ :

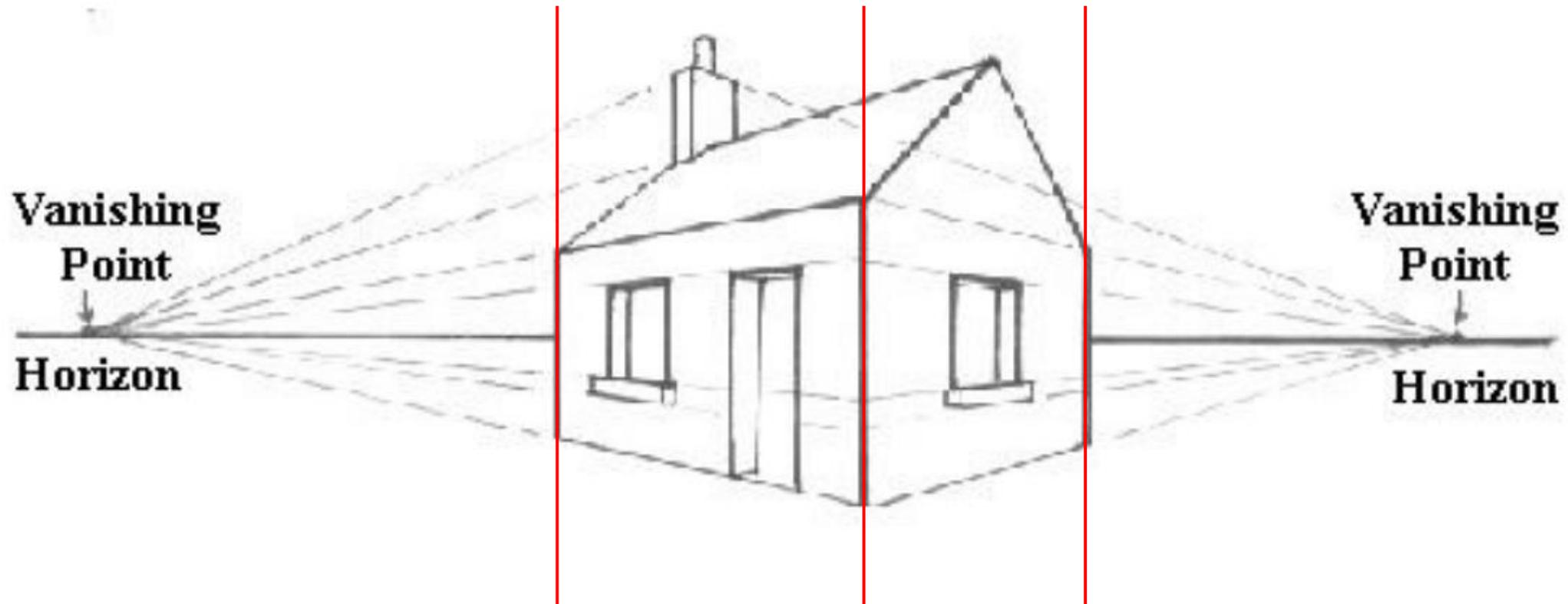
$$\lim_{t \rightarrow \infty} y(t) = \lim_{t \rightarrow \infty} f \frac{Y_0 + bt}{Z_0 + ct} = f \frac{b}{c}$$

$$\lim_{t \rightarrow \infty} x(t) = \lim_{t \rightarrow \infty} f \frac{X_0 + at}{Z_0 + ct} = f \frac{a}{c}$$

This is a point!

# Each Family of Parallel Lines has its Own Vanishing Point

Where is the “vanishing point” of these vertical parallel lines?



# Proof

- Line in 3D space:

$$X(t) = X_0 + 0t = X_0$$

$$Y(t) = Y_0 + bt$$

$$Z(t) = Z_0 + 0t = Z_0$$

- Perspective projection:

$$y(t) = f \frac{Y(t)}{Z(t)} = f \frac{Y_0 + bt}{Z_0}$$

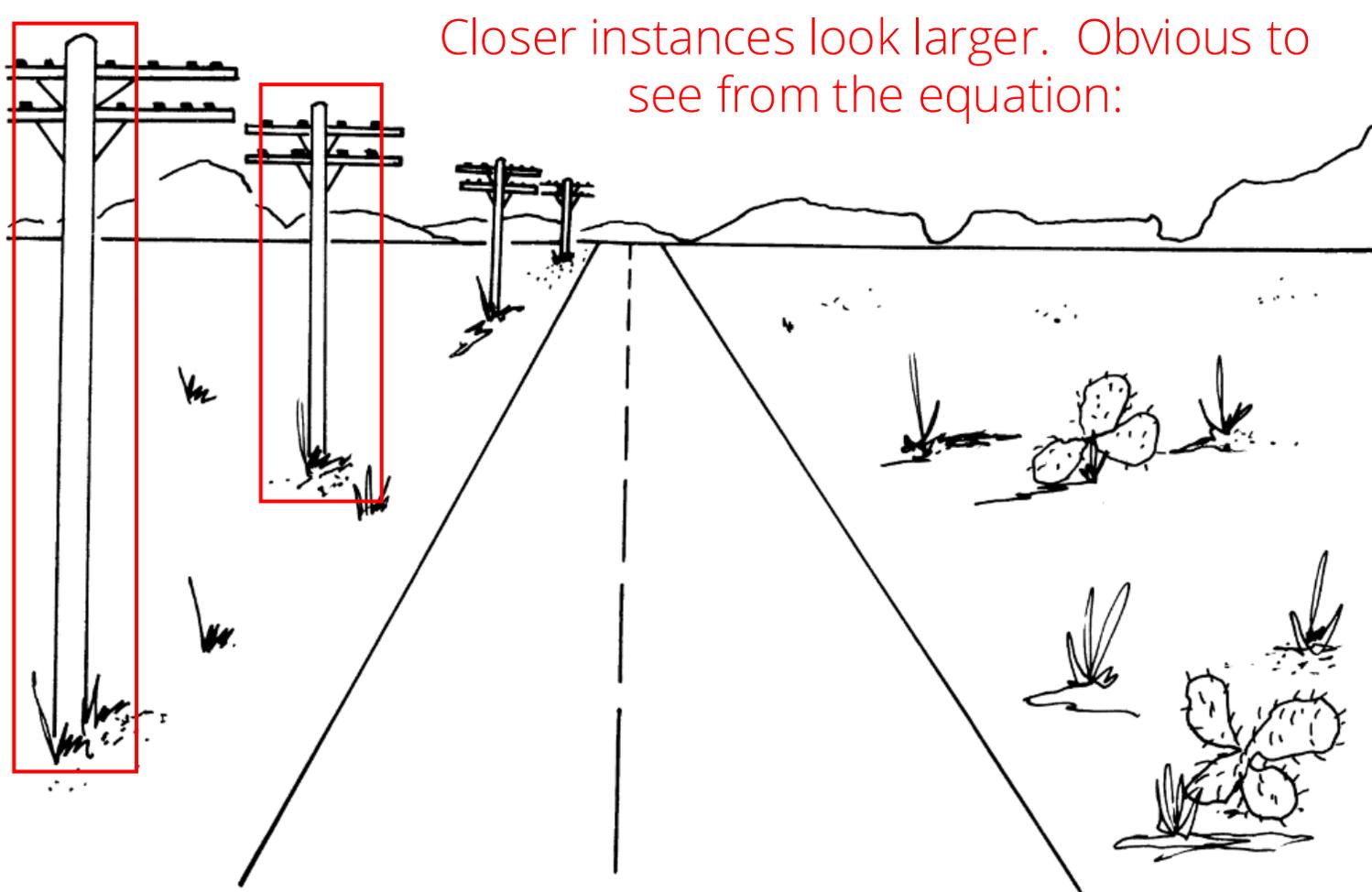
$$x(t) = f \frac{X(t)}{Z(t)} = f \frac{X_0}{Z_0}$$

- Projective lines when  $t \rightarrow \infty$ :

$$\lim_{t \rightarrow \infty} y(t) = \lim_{t \rightarrow \infty} f \frac{Y_0 + bt}{Z_0} = f \frac{b}{0}$$

$$\lim_{t \rightarrow \infty} x(t) = \lim_{t \rightarrow \infty} f \frac{X_0}{Z_0}$$

# Scales Are Not Preserved by Perspective Projection



$$y = f \frac{Y}{Z}$$
$$x = f \frac{X}{Z}$$

# Illusion of Perspective Projection





# Homogeneous Coordinates

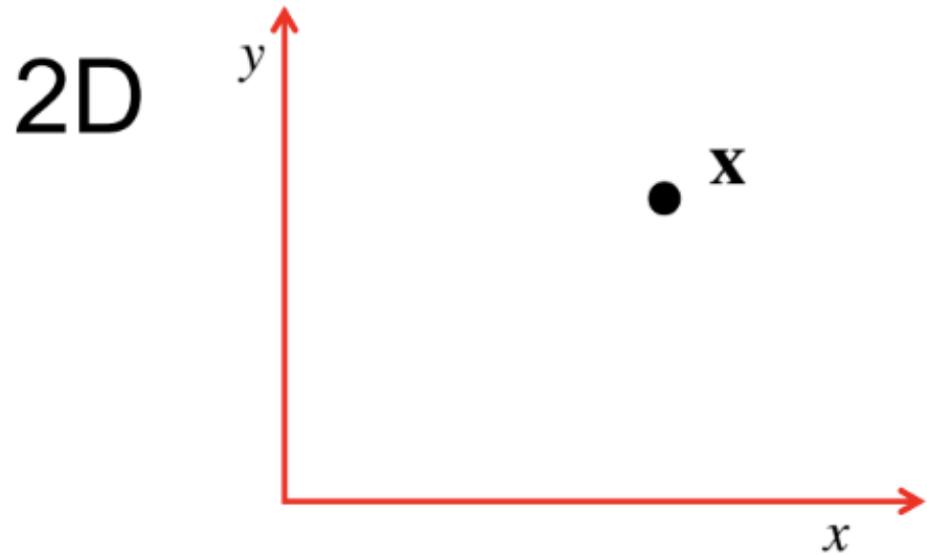
- Representing n-dim coordinates with (n+1)-dim

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} wx_1 \\ wx_2 \\ wx_3 \\ \vdots \\ wx_n \\ w \end{bmatrix}$$

Represent the same point in the  
n-dimensional space

- We can think of this as a line through the origin in (n+1)-dimensional space

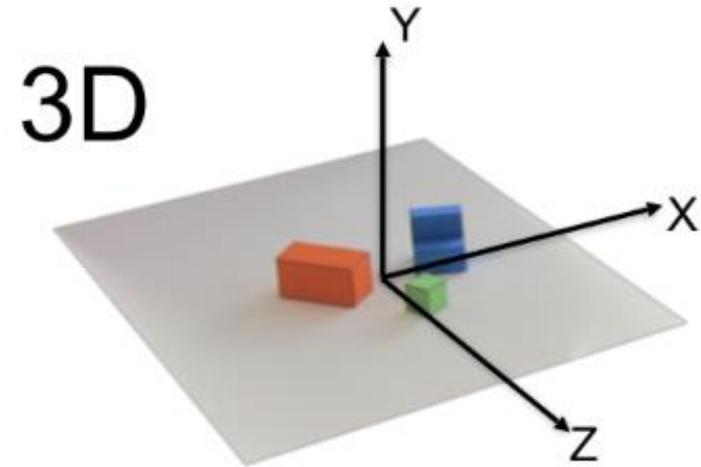
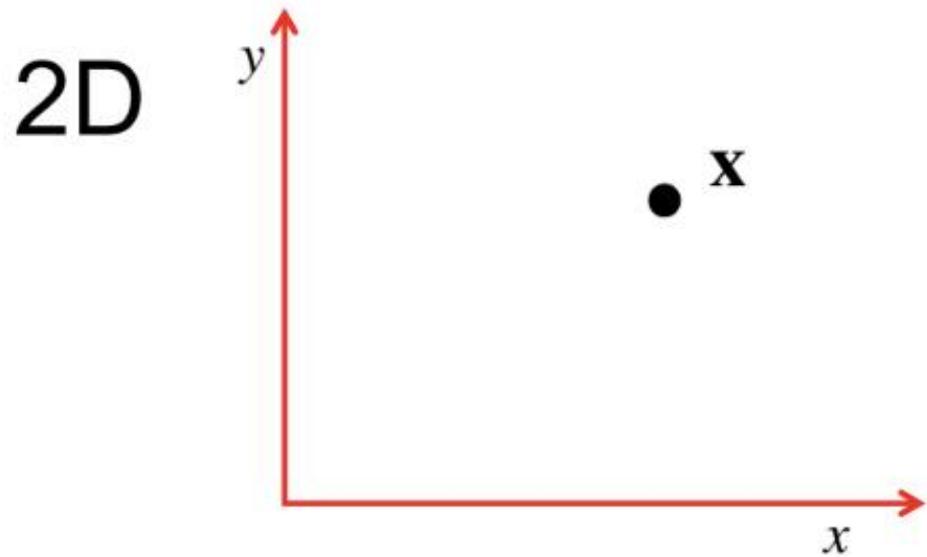
# Homogeneous Coordinates



$$\mathbf{x} = (x, y) \quad \rightarrow \quad \tilde{\mathbf{x}} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

heterogeneous coordinates                      homogeneous coordinates

# Homogeneous Coordinates



$$\mathbf{x} = (x, y) \quad \longrightarrow \quad \tilde{\mathbf{x}} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

heterogeneous coordinates                      homogeneous coordinates

$$\mathbf{X} = (X, Y, Z) \quad \longrightarrow \quad \tilde{\mathbf{X}} = \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

# Homogeneous Coordinates

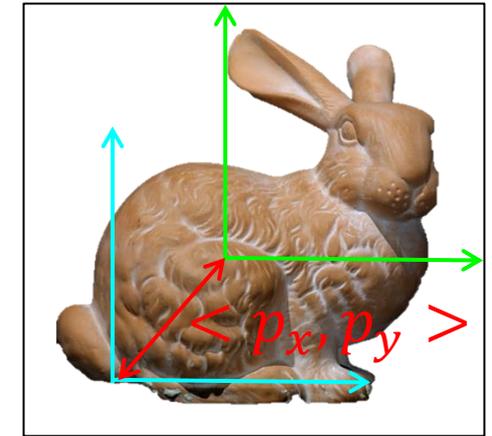
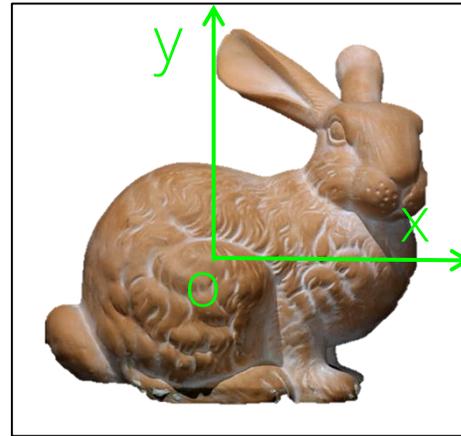
- Homogeneous image coordinates:  $(x, y) \Rightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$
- Homogeneous scene coordinates:  $(X, Y, Z) \Rightarrow \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$
- We can represent the pin-hole camera model with homogeneous coordinates

$$\begin{aligned} x &= f \frac{X}{Z} \\ y &= f \frac{Y}{Z} \end{aligned} \longrightarrow \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = \begin{bmatrix} fX \\ fY \\ Z \end{bmatrix} \Rightarrow \begin{bmatrix} f \frac{X}{Z} \\ f \frac{Y}{Z} \\ 1 \end{bmatrix}$$

# Content

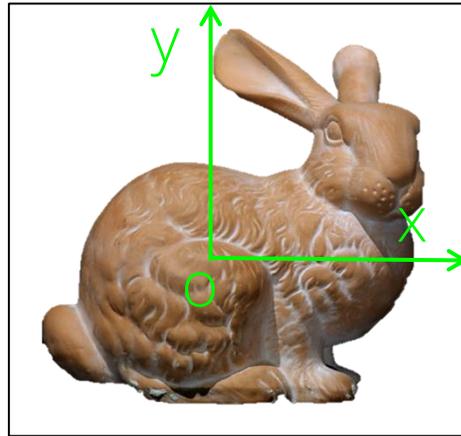
- Image Formation
  - Pinhole camera model
  - 2D transformation
  - 3D transformation
- Geometry Reconstruction
  - Reconstruction with depths
  - Multi-view geometry reconstruction with camera motions and correspondence
  - Epipolar geometry reconstruction

# What is the Matrix Form of Perspective Projection with



$$\begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

# Decompose the Projection Matrix

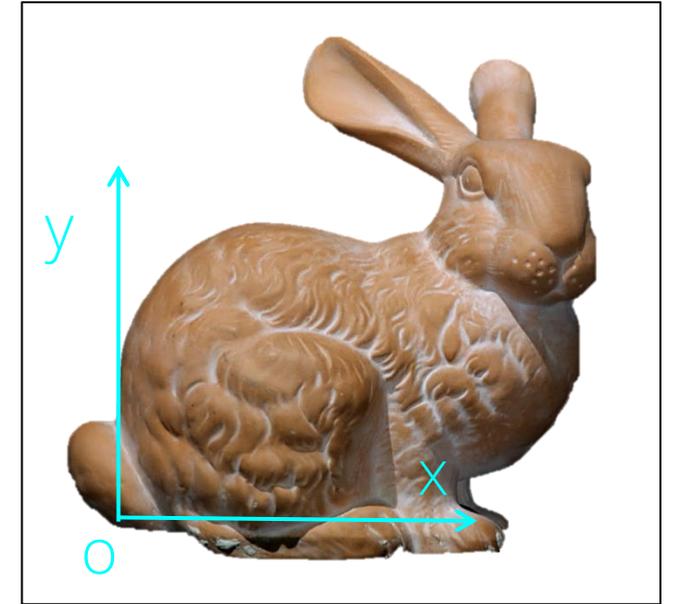
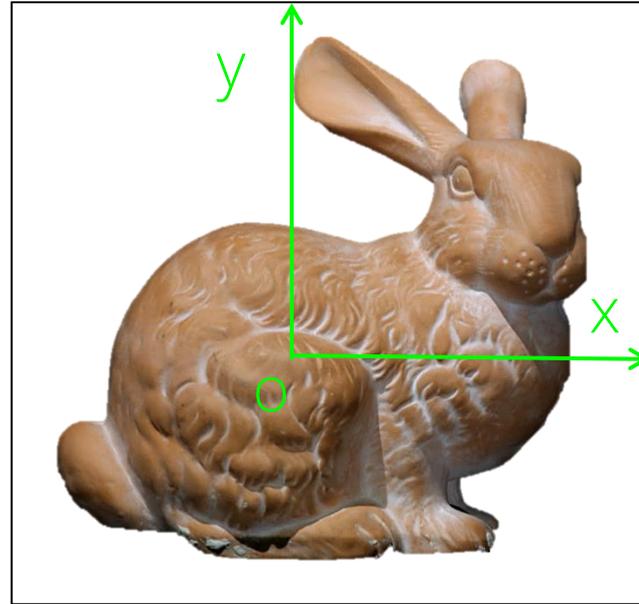
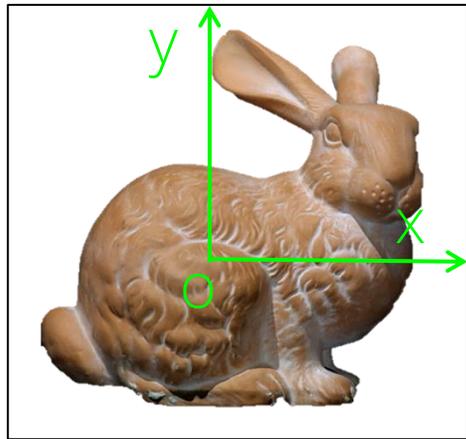


Perspective projection from 3D to 2D, assuming image plane at  $z=1$  and shared camera/image origin

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \Rightarrow \begin{bmatrix} X/Z \\ Y/Z \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

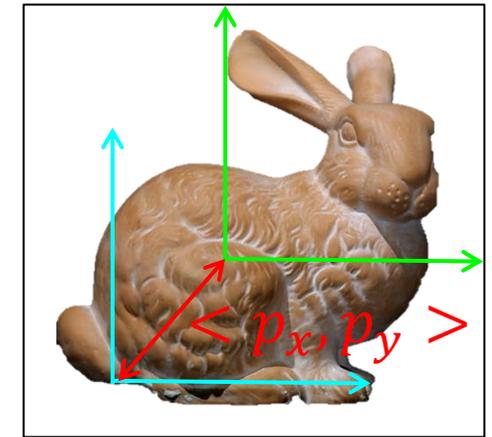
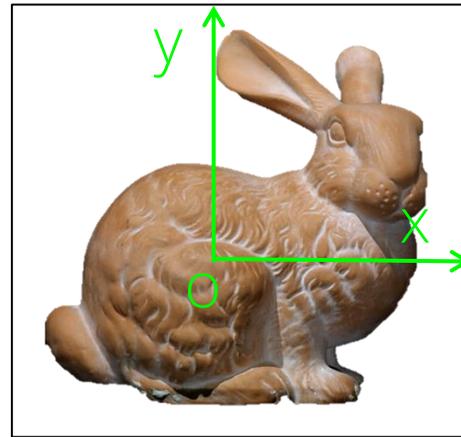
# Decompose the Projection Matrix



$$\begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X/Z \\ Y/Z \\ 1 \end{bmatrix} = \begin{bmatrix} f X/Z + p_x \\ f Y/Z + p_y \\ 1 \end{bmatrix}$$

Transformation from 2D to 2D, accounting for not unit focal length and origin shift

# What is the Matrix Form of Perspective Projection with



$$\begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$K_{img}^{pix}$

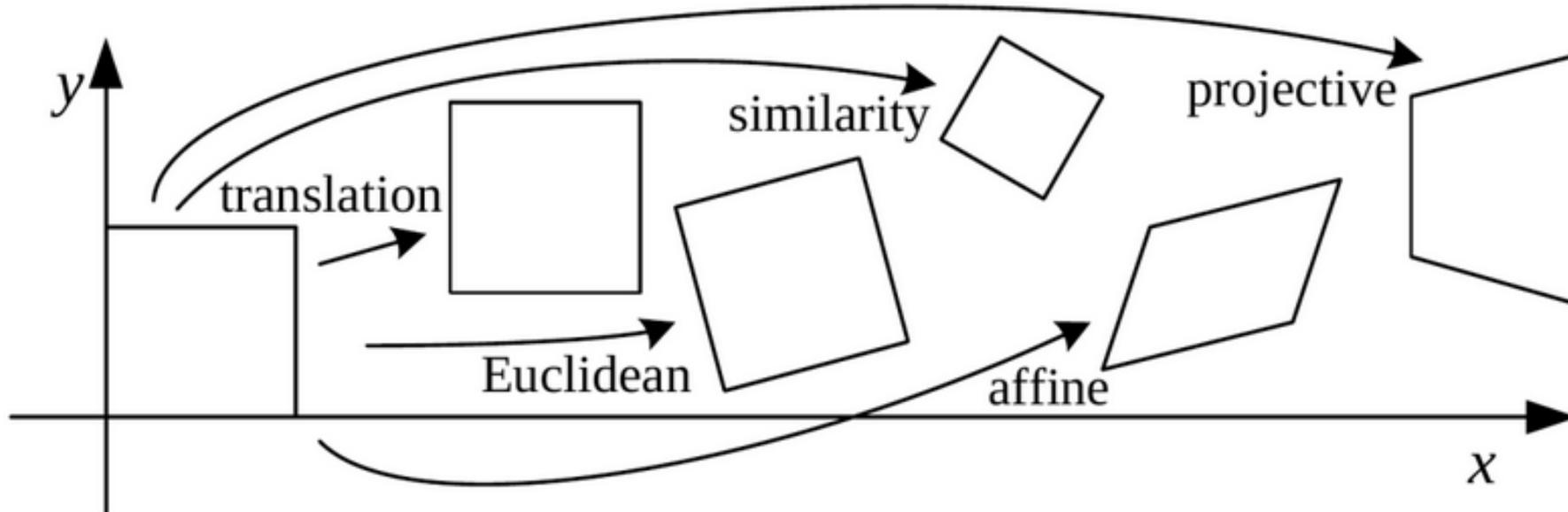
$K_{cam}^{img}$

$K_{img}^{pix} K_{cam}^{img}$  is also called camera intrinsics matrix

# 2D-2D Transformation

$$\begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X/Z \\ Y/Z \\ 1 \end{bmatrix} = \begin{bmatrix} f X/Z + p_x \\ f Y/Z + p_y \\ 1 \end{bmatrix}$$

$K$        $\chi^{old}$        $\chi^{new}$

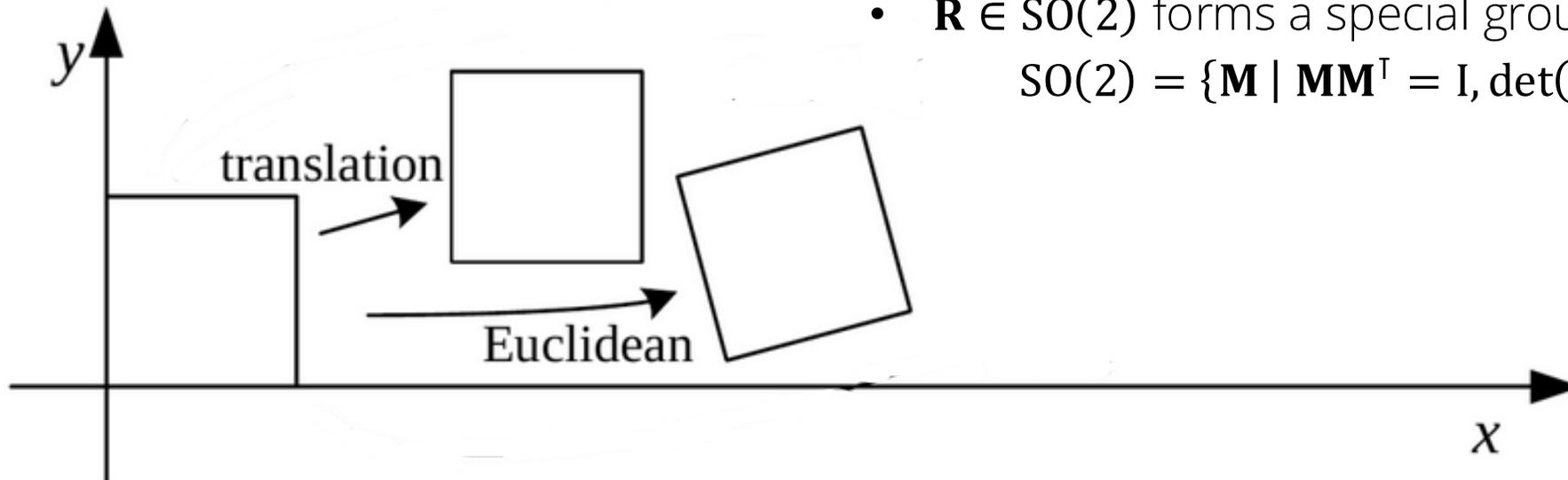


# Euclidean Transformations

$$\mathbf{R}\mathbf{x} + \mathbf{t} = \mathbf{x}' \Leftrightarrow \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix}$$

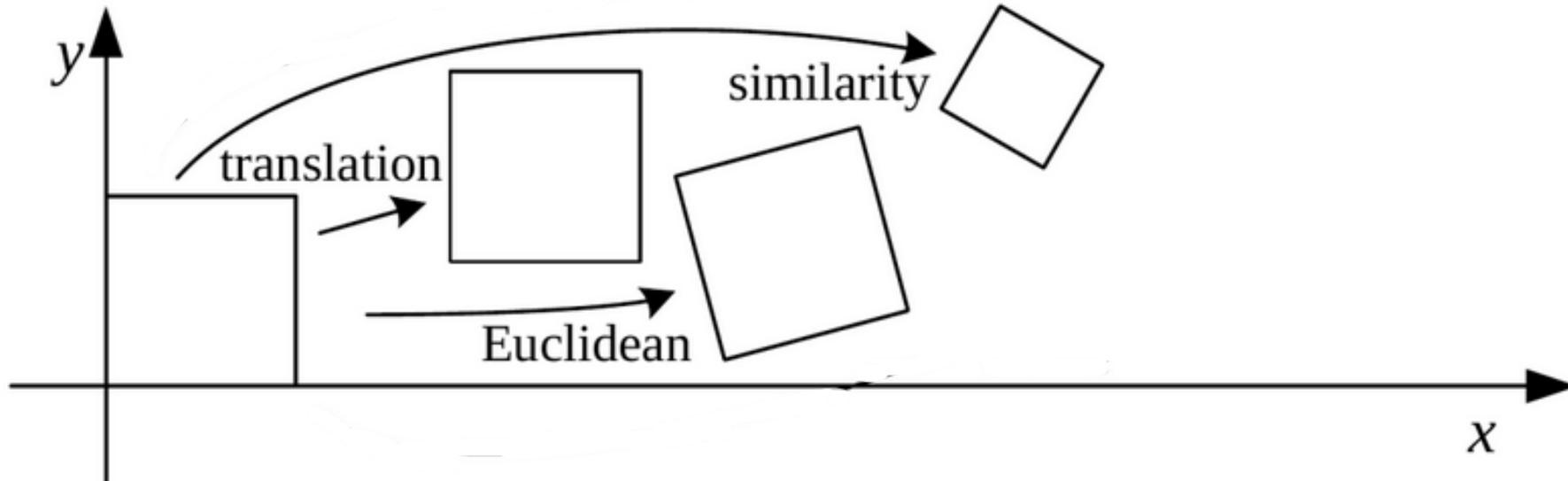
Bold fonts denote vector

- Translation is a special case when  $\mathbf{R} = \mathbf{I}$
- $\mathbf{R} \in SO(2)$  forms a special group  
 $SO(2) = \{\mathbf{M} \mid \mathbf{M}\mathbf{M}^\top = \mathbf{I}, \det(\mathbf{M}) = 1\}$



# Similarity Transformations

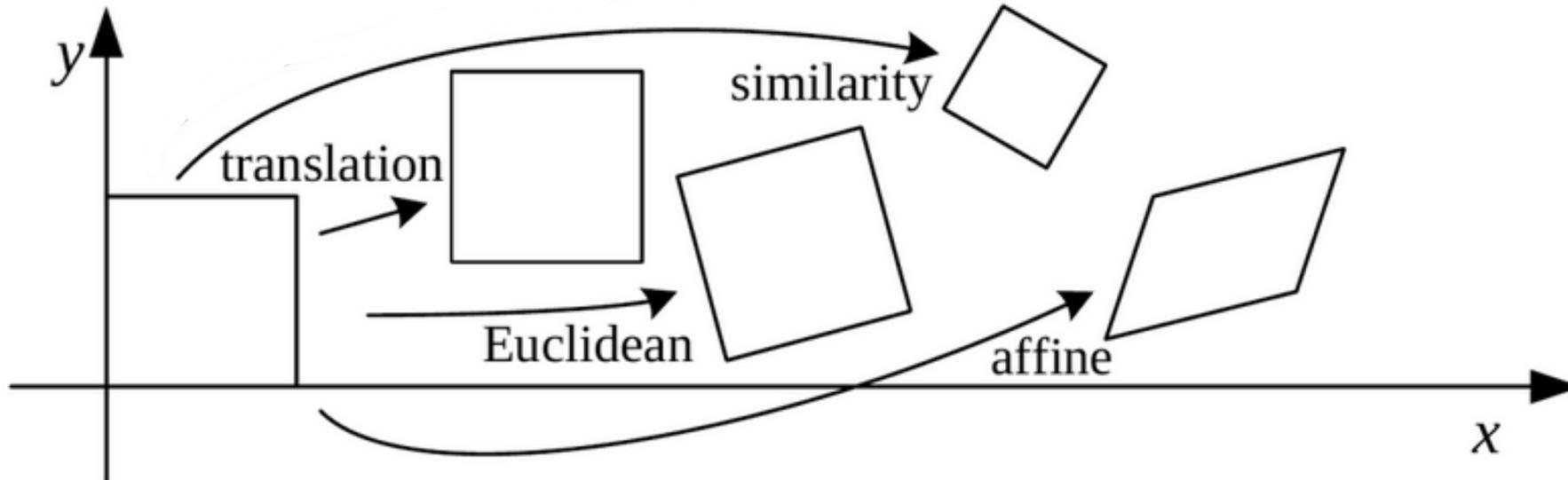
$$s\mathbf{R}\mathbf{x} + \mathbf{t} = \mathbf{x}' \Leftrightarrow \begin{bmatrix} s\mathbf{R} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix}$$



# Affine Transformations

$$\mathbf{Ax} + \mathbf{t} = \mathbf{x}' \Leftrightarrow \begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix}$$

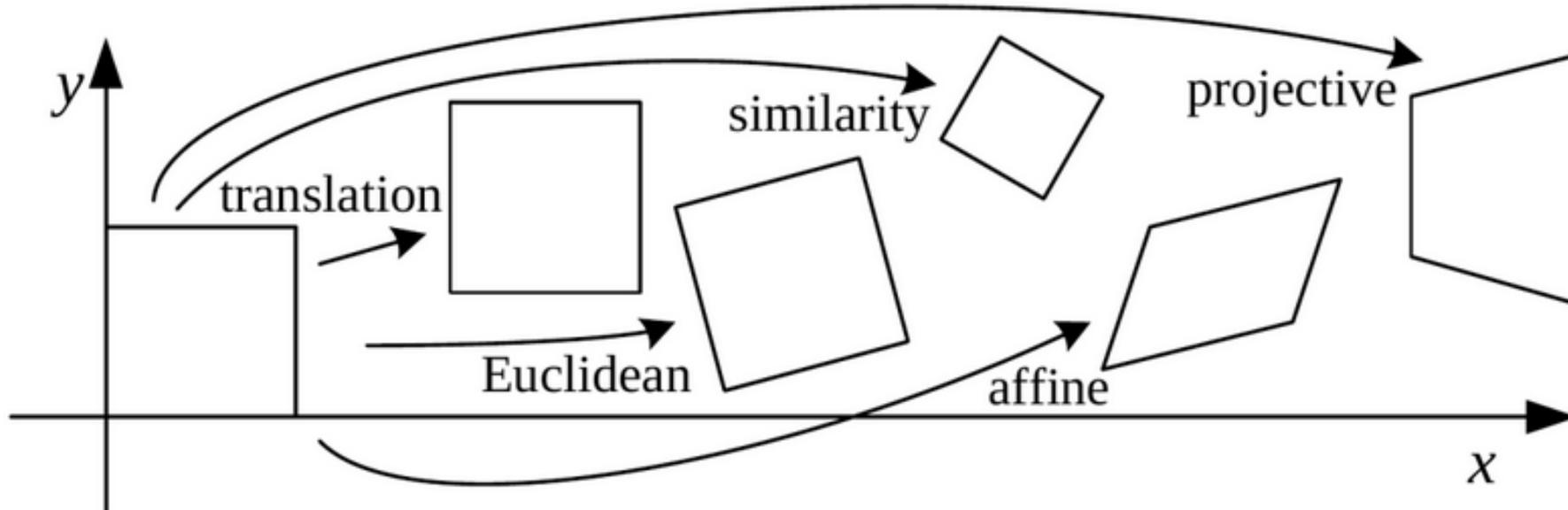
- $\mathbf{A}$  is an arbitrary  $2 \times 2$  matrix
- Parallel lines remain parallel after affine transformations



# Projective Transformations

$$\mathbf{H}\mathbf{x} = \mathbf{x}' \Leftrightarrow \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix}$$

- $\mathbf{H}$  is an arbitrary  $3 \times 3$  matrix, also called homography matrix
- Projective transformations preserve straight lines

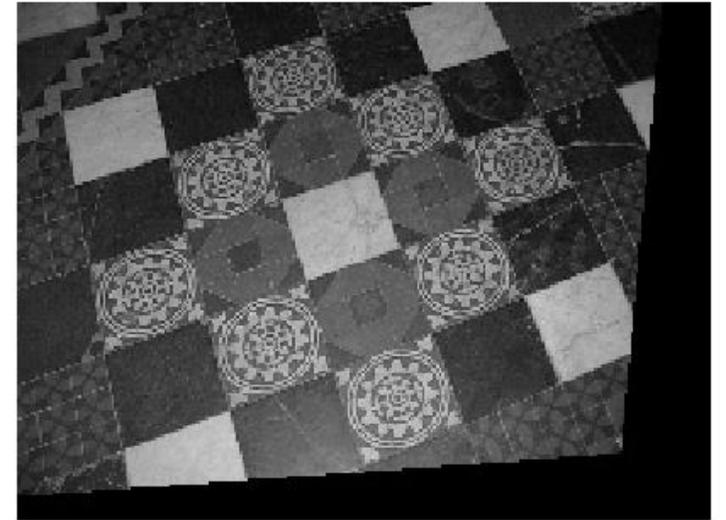


# Application: Image Rectification with Projective Transformations

Image 1



Image 2



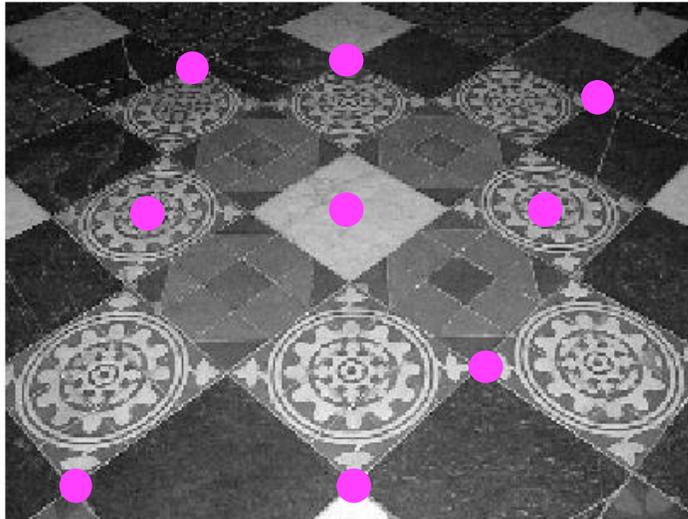
$$Hx = x'$$

# How Does Image Rectification Work?

Image 1



Image 2



- We have a set of points  $\mathbf{x}$  from image 2 and our goal is to rectify image 2 so that it looks like image 1
- Transformation matrix  $\mathbf{H}$  and transformed points are  $\mathbf{x}'$  unknown
- The equation is not solvable with two unknown variables



$$\mathbf{H}\mathbf{x} = \mathbf{x}'$$

# Idea: Obtain $\mathbf{x}'$ from correspondence

Image 1

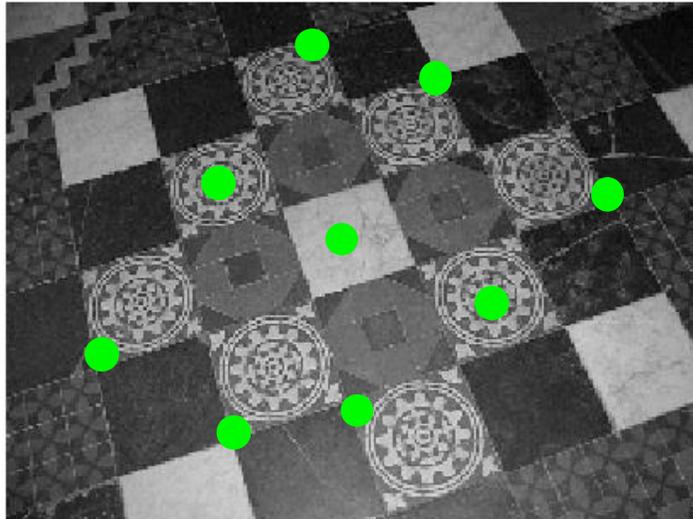
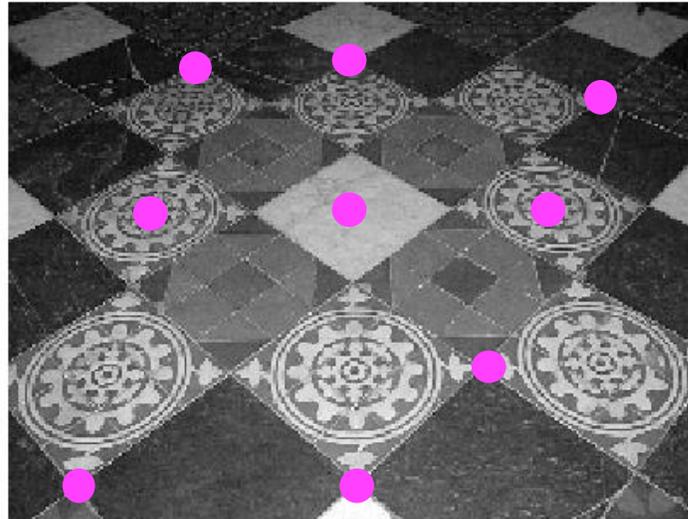


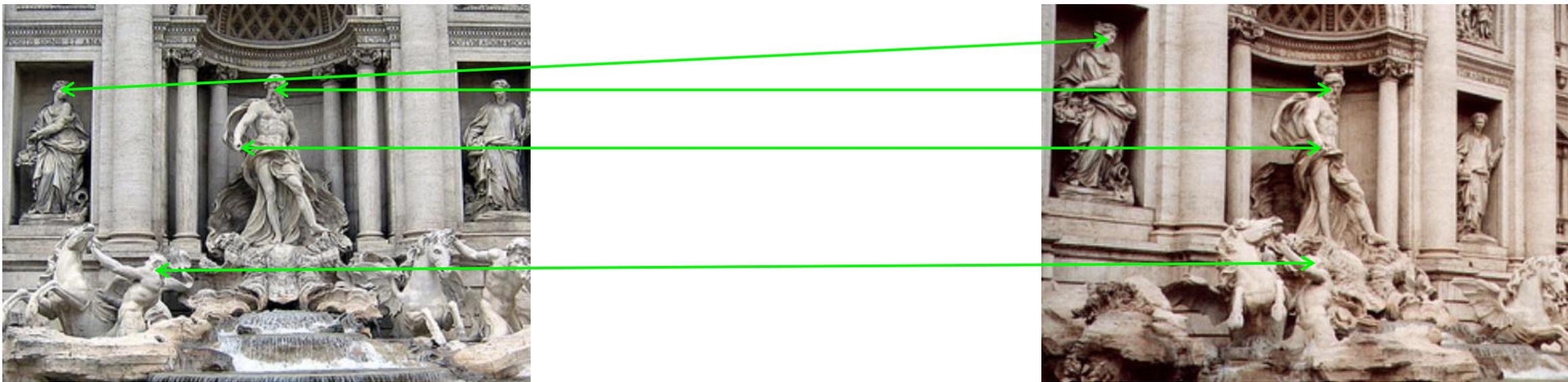
Image 2



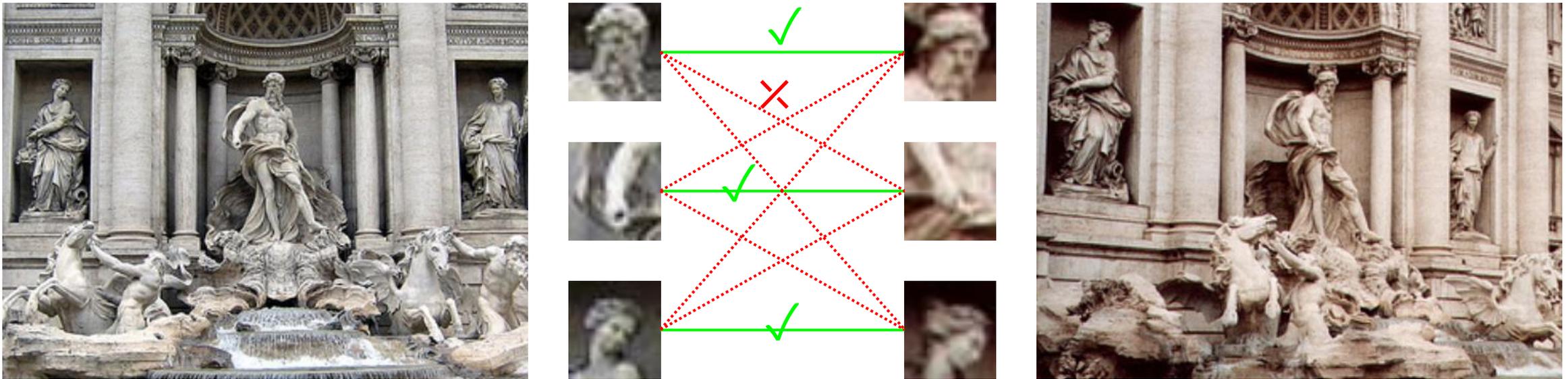
- We have a set of points  $\mathbf{x}$  from image 2 and our goal is to rectify image 2 so that it looks like image 1
- We obtain  $\mathbf{x}'$  through pixel correspondence


$$\mathbf{H}\mathbf{x} = \mathbf{x}'$$

# Idea: Obtain $\mathbf{x}'$ from correspondence



# Idea: Obtain $\mathbf{x}'$ from correspondence



- How to obtain pixel correspondence between two images?
  - Use a good feature extractor (e.g. SIFT / DINO.v2 ...) to encode each pixel into a feature vector
  - Find the best matching between two sets of feature vectors

# Solve $\mathbf{H}$ given $\mathbf{x}$ and $\mathbf{x}'$

Image 1

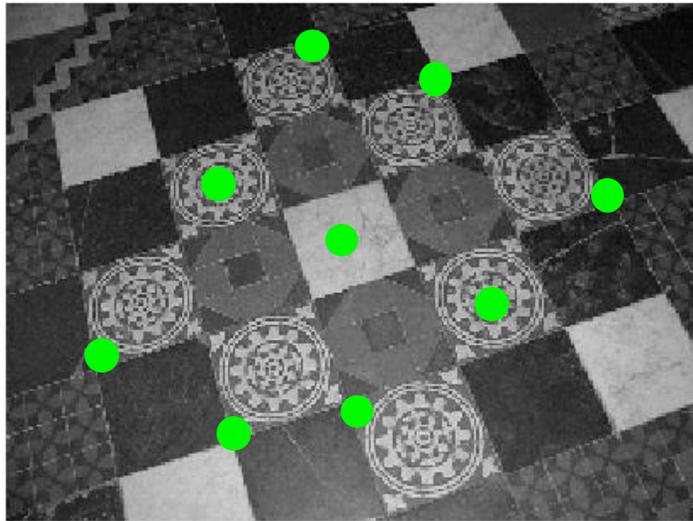
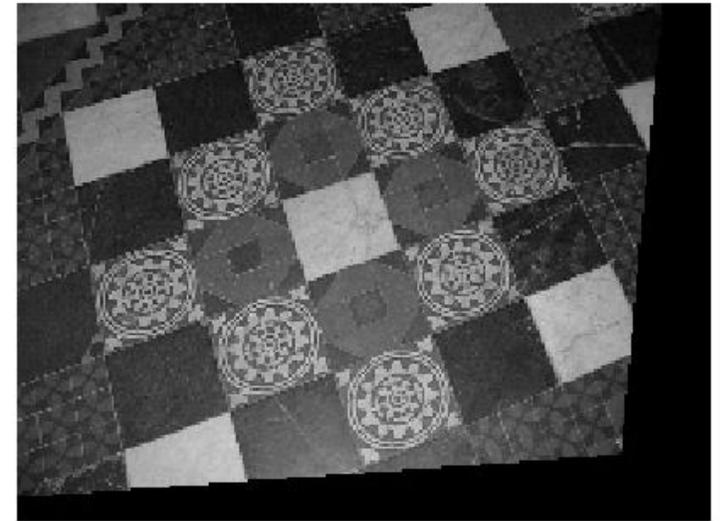
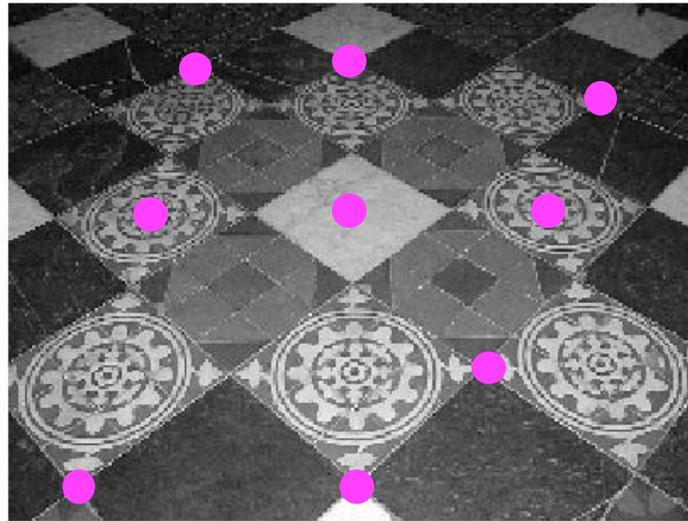


Image 2



- We have a set of points  $\mathbf{x}$  from image 2 and our goal is to rectify image 2 so that it looks like image 1
- We obtain  $\mathbf{x}'$  through pixel correspondence
- Transformation matrix  $\mathbf{H}$
- The problem is to find  $\mathbf{H}$  such that  $\mathbf{H}\mathbf{x} = \mathbf{x}'$



$$\mathbf{H}\mathbf{x} = \mathbf{x}'$$

# The Direct Linear Transformation Algorithm

- Let's start with one point correspondence:  $\mathbf{H}\mathbf{x}_i = \mathbf{x}'_i$
- Note that  $\mathbf{x}_i$  and  $\mathbf{x}'_i$  are homogeneous coordinates, thus  $\mathbf{H}\mathbf{x}_i$  and  $\mathbf{x}'_i$  are in the same direction, might differ in magnitude

$$\begin{bmatrix} \mathbf{h}_1^\top \\ \mathbf{h}_2^\top \\ \mathbf{h}_3^\top \end{bmatrix} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} = \begin{bmatrix} x'_i \\ y'_i \\ w'_i \end{bmatrix} \Rightarrow \mathbf{x}'_i \times \mathbf{H}\mathbf{x}_i = 0 \Rightarrow \begin{bmatrix} x'_i \\ y'_i \\ w'_i \end{bmatrix} \times \begin{bmatrix} \mathbf{h}_1^\top \mathbf{x}_i \\ \mathbf{h}_2^\top \mathbf{x}_i \\ \mathbf{h}_3^\top \mathbf{x}_i \end{bmatrix} = 0$$

$$\Rightarrow \begin{bmatrix} y'_i \mathbf{h}_3^\top \mathbf{x}_i - w'_i \mathbf{h}_2^\top \mathbf{x}_i \\ w'_i \mathbf{h}_1^\top \mathbf{x}_i - x'_i \mathbf{h}_3^\top \mathbf{x}_i \\ x'_i \mathbf{h}_2^\top \mathbf{x}_i - y'_i \mathbf{h}_1^\top \mathbf{x}_i \end{bmatrix} = 0 \Rightarrow \begin{bmatrix} \mathbf{0}^\top & -w'_i \mathbf{x}_i^\top & y'_i \mathbf{x}_i^\top \\ w'_i \mathbf{x}_i^\top & \mathbf{0}^\top & -x'_i \mathbf{x}_i^\top \\ -y'_i \mathbf{x}_i^\top & x'_i \mathbf{x}_i^\top & \mathbf{0}^\top \end{bmatrix} \begin{pmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \mathbf{h}_3 \end{pmatrix} = 0$$

3 × 9 matrix

9 × 1 vector

Remember bold-font variable denotes vector

# The Direct Linear Transformation Algorithm

- Let's start with one point correspondence:  $\mathbf{H}\mathbf{x}_i = \mathbf{x}'_i$
- Note that  $\mathbf{x}_i$  and  $\mathbf{x}'_i$  are homogeneous coordinates, thus  $\mathbf{H}\mathbf{x}_i$  and  $\mathbf{x}'_i$  are in the same direction, might differ in magnitude

$$\begin{aligned} \begin{bmatrix} \mathbf{h}_1^\top \\ \mathbf{h}_2^\top \\ \mathbf{h}_3^\top \end{bmatrix} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} &= \begin{bmatrix} x'_i \\ y'_i \\ w'_i \end{bmatrix} \Rightarrow \mathbf{x}'_i \times \mathbf{H}\mathbf{x}_i = 0 \Rightarrow \begin{bmatrix} x'_i \\ y'_i \\ w'_i \end{bmatrix} \times \begin{bmatrix} \mathbf{h}_1^\top \mathbf{x}_i \\ \mathbf{h}_2^\top \mathbf{x}_i \\ \mathbf{h}_3^\top \mathbf{x}_i \end{bmatrix} = 0 \\ \Rightarrow \begin{bmatrix} y'_i \mathbf{h}_3^\top \mathbf{x}_i - w'_i \mathbf{h}_2^\top \mathbf{x}_i \\ w'_i \mathbf{h}_1^\top \mathbf{x}_i - x'_i \mathbf{h}_3^\top \mathbf{x}_i \\ x'_i \mathbf{h}_2^\top \mathbf{x}_i - y'_i \mathbf{h}_1^\top \mathbf{x}_i \end{bmatrix} = 0 &\Rightarrow \begin{bmatrix} \mathbf{0}^\top & -w'_i \mathbf{x}_i^\top & y'_i \mathbf{x}_i^\top \\ w'_i \mathbf{x}_i^\top & \mathbf{0}^\top & -x'_i \mathbf{x}_i^\top \\ -y'_i \mathbf{x}_i^\top & x'_i \mathbf{x}_i^\top & \mathbf{0}^\top \end{bmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} = 0 \end{aligned}$$

By adding one constraint  $\left\| \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} \right\| = 1$ , we

need 4 points with correspondence

3 × 9 matrix with 2 × 9 9 × 1 vector  
degree of freedom

# The Direct Linear Transformation Algorithm: A Solution with SVD

## Objective

Find the general solution to a set of equations  $\mathbf{Ax} = \mathbf{b}$  where  $\mathbf{A}$  is an  $m \times n$  matrix of rank  $r < n$ .

## Algorithm

- (i) Find the SVD  $\mathbf{A} = \mathbf{UDV}^T$ , where the diagonal entries  $d_i$  of  $\mathbf{D}$  are in descending numerical order.
- (ii) Set  $\mathbf{b}' = \mathbf{U}^T \mathbf{b}$ .
- (iii) Find the vector  $\mathbf{y}$  defined by  $y_i = b'_i/d_i$  for  $i = 1, \dots, r$ , and  $y_i = 0$  otherwise.
- (iv) The solution  $\mathbf{x}$  of minimum norm  $\|\mathbf{x}\|$  is  $\mathbf{Vy}$ .
- (v) The general solution is  $\mathbf{x} = \mathbf{Vy} + \lambda_{r+1} \mathbf{v}_{r+1} + \dots + \lambda_n \mathbf{v}_n$ , where  $\mathbf{v}_{r+1}, \dots, \mathbf{v}_n$  are the last  $n - r$  columns of  $\mathbf{V}$ .

We only need the eigen vector of the least eigen value for  $\mathbf{Ax} = \mathbf{b}$

# The Direct Linear Transformation Algorithm

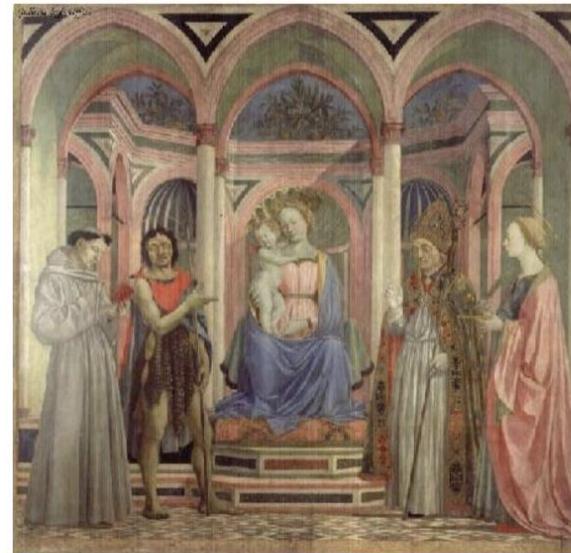
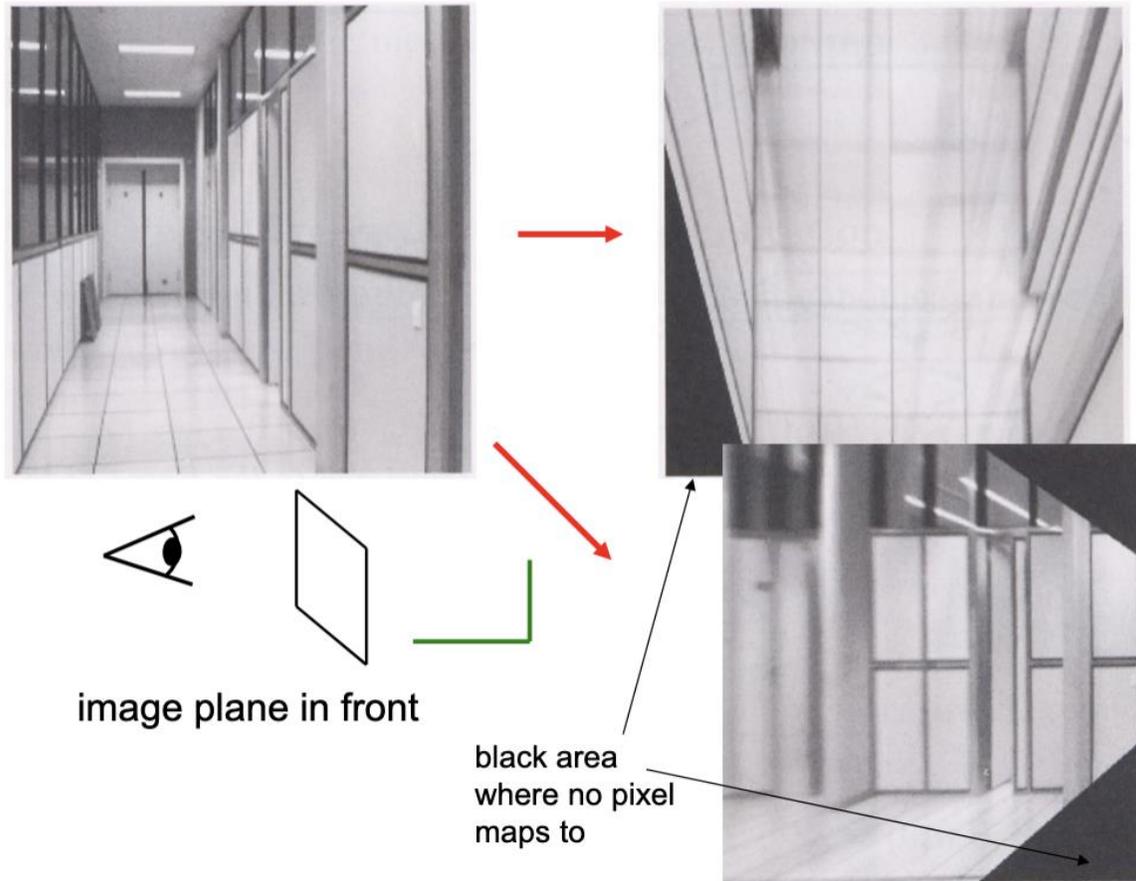
## Objective

Given  $n \geq 4$  2D to 2D point correspondences  $\{\mathbf{x}_i \leftrightarrow \mathbf{x}'_i\}$ , determine the 2D homography matrix  $H$  such that  $\mathbf{x}'_i = H\mathbf{x}_i$ .

## Algorithm

- (i) For each correspondence  $\mathbf{x}_i \leftrightarrow \mathbf{x}'_i$  compute the matrix  $A_i$  from (4.1). Only the first two rows need be used in general.
- (ii) Assemble the  $n$   $2 \times 9$  matrices  $A_i$  into a single  $2n \times 9$  matrix  $A$ .
- (iii) Obtain the SVD of  $A$  (section A4.4(p585)). The unit singular vector corresponding to the smallest singular value is the solution  $\mathbf{h}$ . Specifically, if  $A = UDV^T$  with  $D$  diagonal with positive diagonal entries, arranged in descending order down the diagonal, then  $\mathbf{h}$  is the last column of  $V$ .
- (iv) The matrix  $H$  is determined from  $\mathbf{h}$  as in (4.2).

# Image Rectification: Virtual Camera Rotation with 2D Transformation



*St. Lucy Altarpiece, D. Veneziano*  
Slide from Criminisi

What is the (complicated) shape of the floor pattern?



**Automatically rectified floor**



From Martin Kemp, *The Science of Art (manual reconstruction)*

# Stitching Images with Image Rectification

## Naïve Stitching

---



left on top

right on top

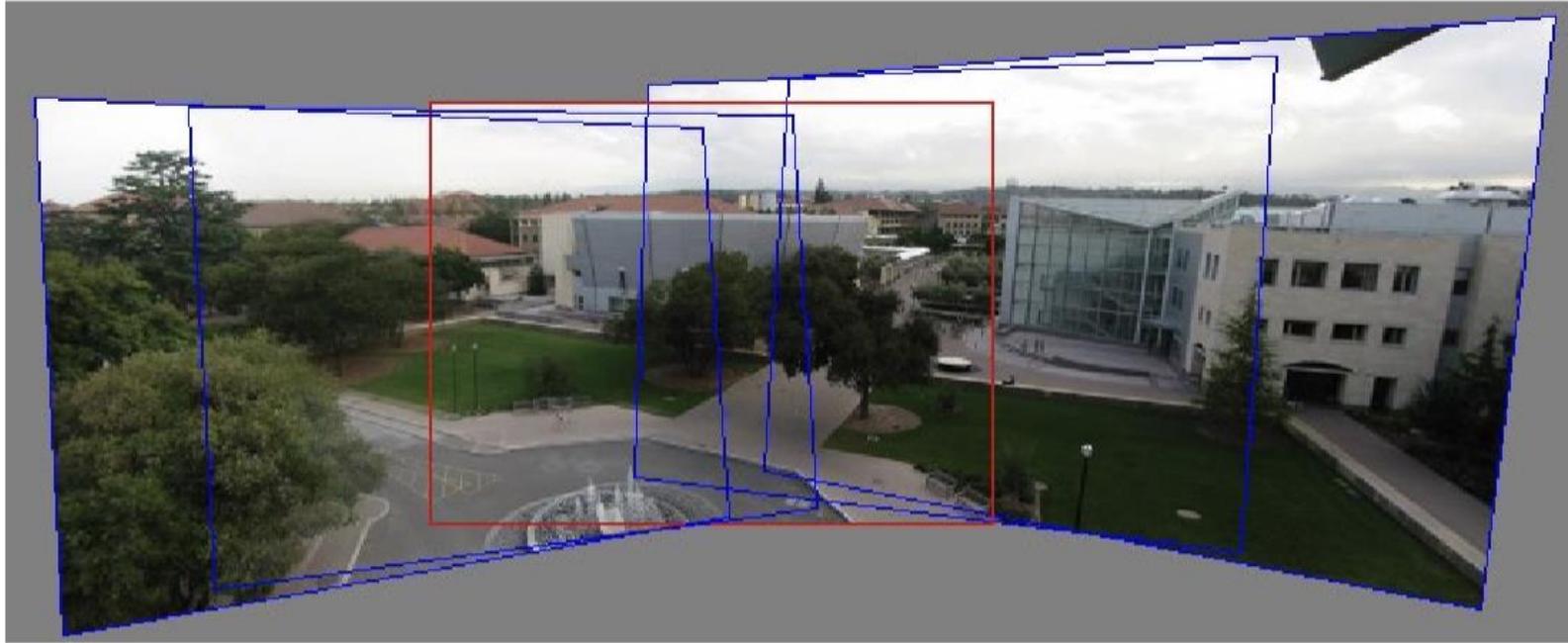


Translations are not enough to align the images



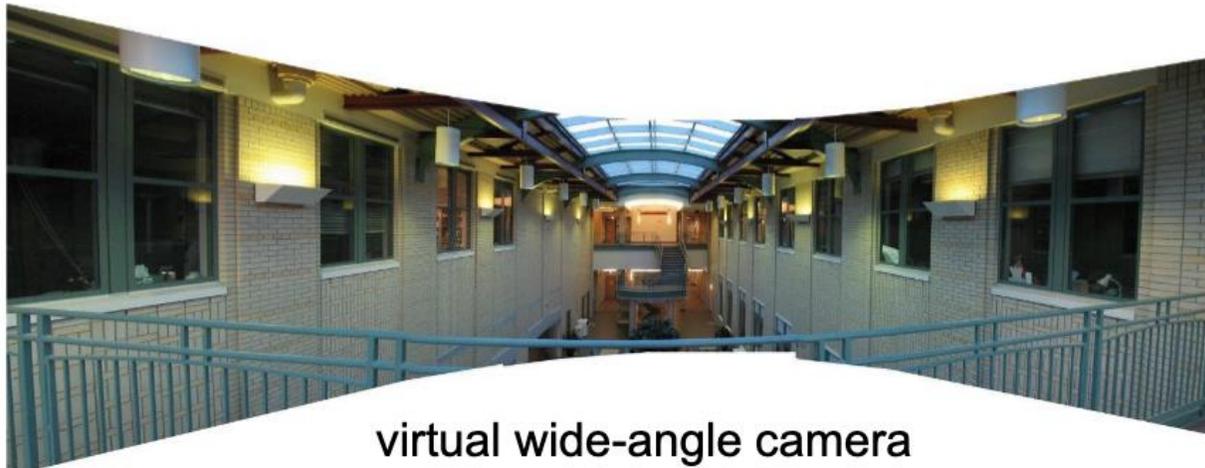
# Stitching Images with Image Rectification

## Panoramas



1. Pick one image (red)
2. Warp the other images towards it (usually, one by one)
3. blend

# Image Panoramas

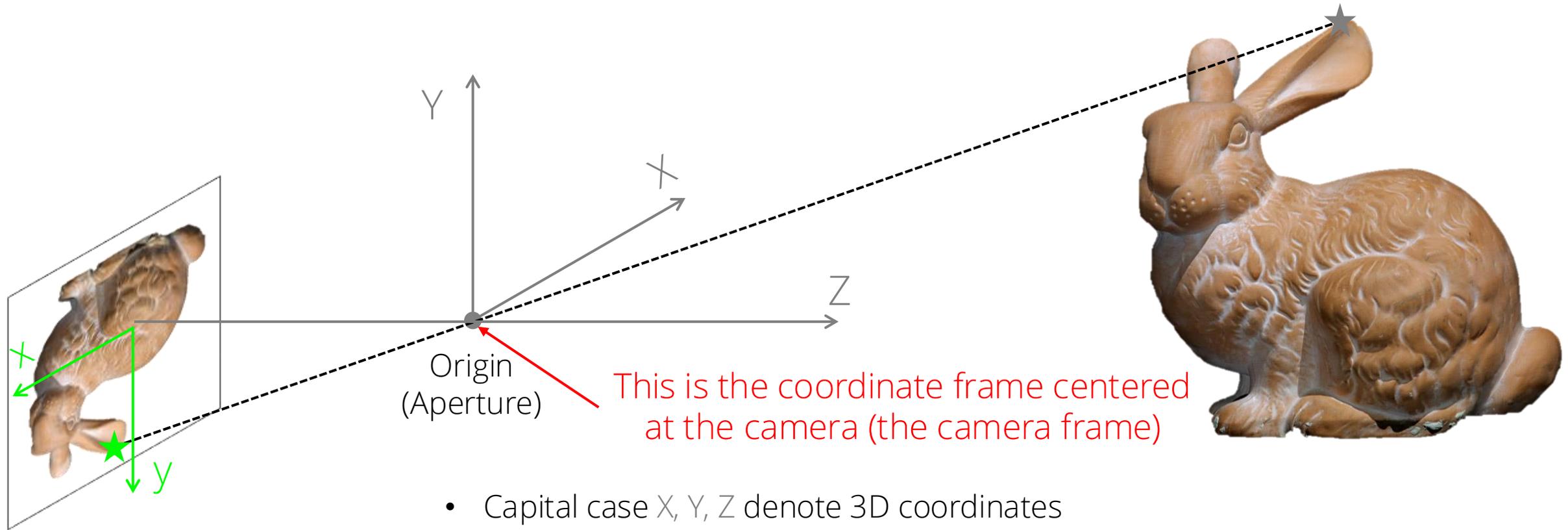


virtual wide-angle camera

# Content

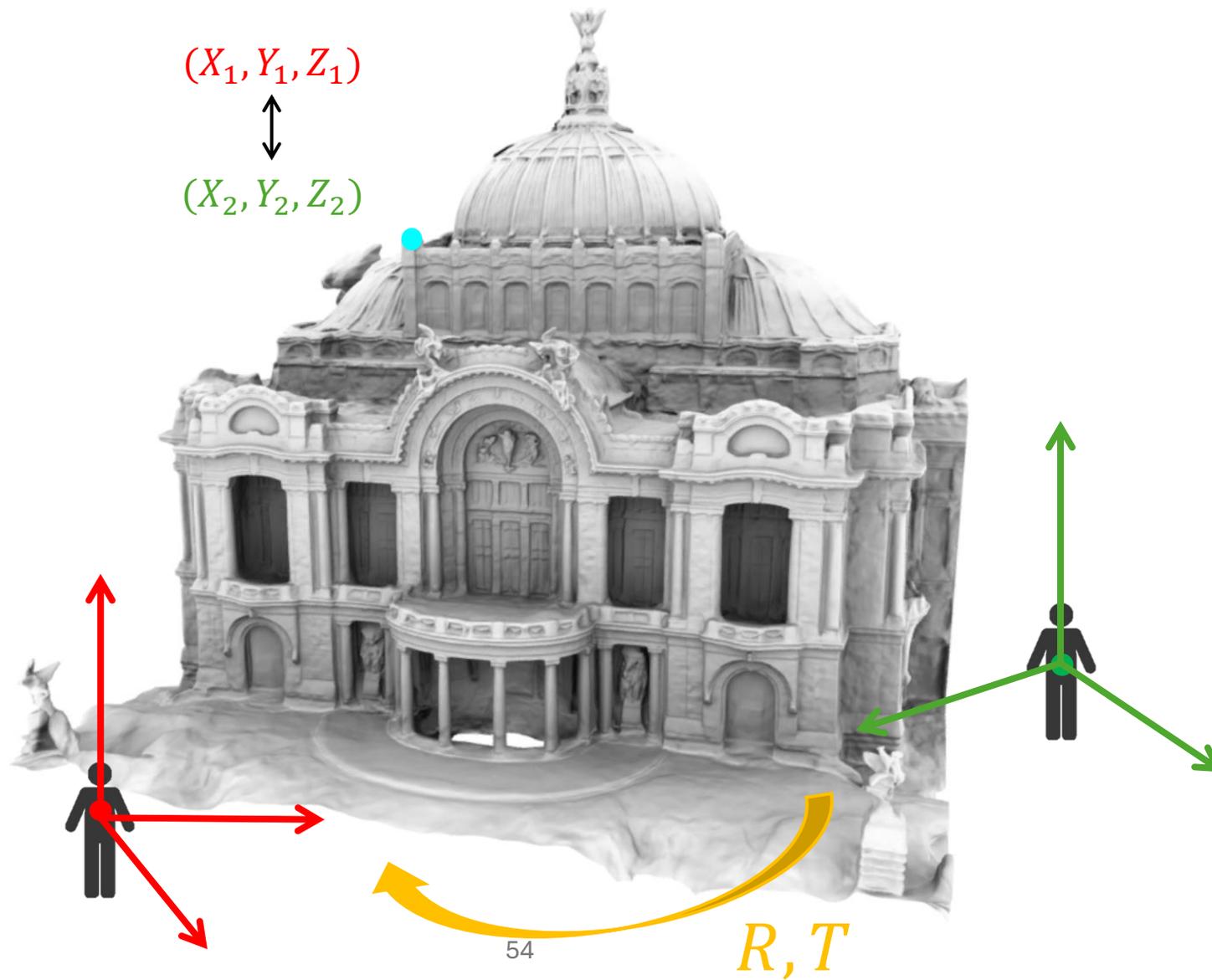
- Image Formation
  - Pinhole camera model
  - 2D transformation
  - 3D transformation
- Geometry Reconstruction
  - Reconstruction with depths
  - Multi-view geometry reconstruction with camera motions and correspondence
  - Epipolar geometry reconstruction

# Let's Get Back to Perspective Projection.// Which Coordinate Frame did We Use?

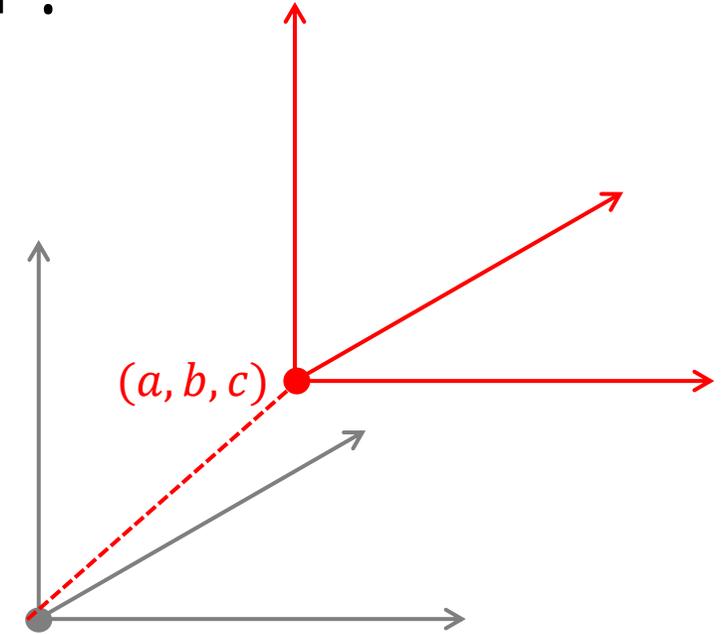
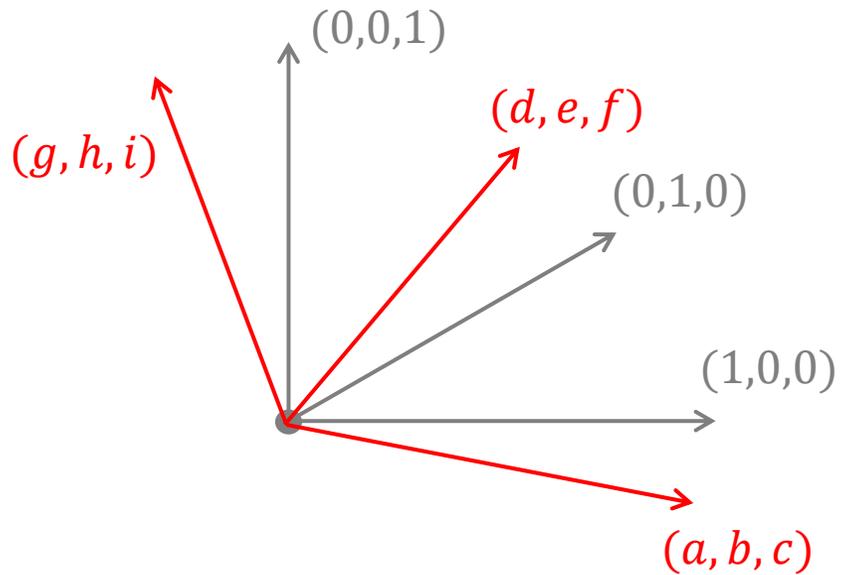


- Capital case  $X, Y, Z$  denote 3D coordinates
- Lower case  $x, y$  denotes 2D Image coordinates
- How to derive the relation between  $\langle x, y \rangle$  and  $\langle X, Y, Z \rangle$  ?

# What if You Have More Than 1 Camera? How to Merge 3D Points from Other Cameras?



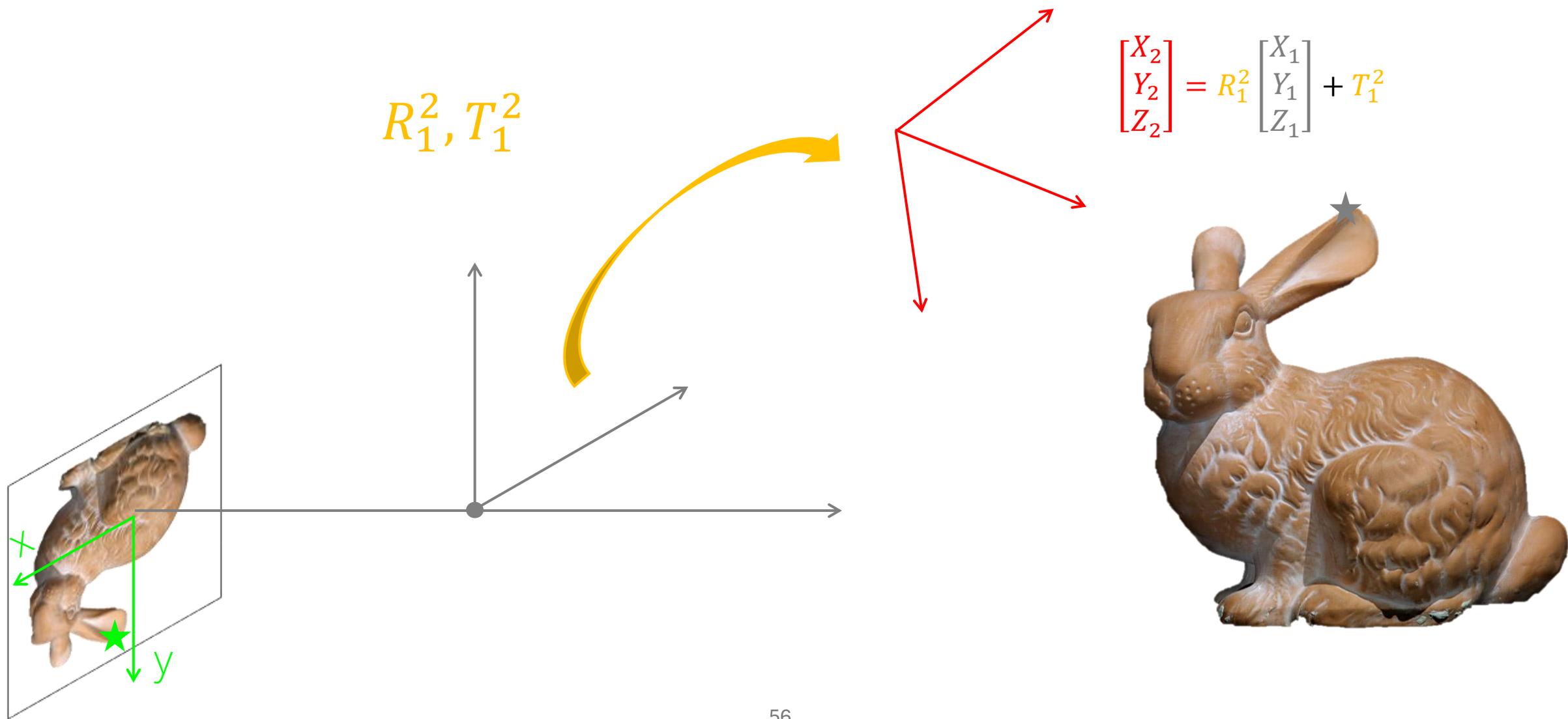
# How to Transform a 3D Point from 1 Coordinate Frame to Another?



$$R = \begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix} \Rightarrow \begin{bmatrix} X_2 \\ Y_2 \\ Z_2 \end{bmatrix} = R \begin{bmatrix} X_1 \\ Y_1 \\ Z_1 \end{bmatrix}$$

$$T = \begin{bmatrix} a \\ b \\ c \end{bmatrix} \Rightarrow \begin{bmatrix} X_2 \\ Y_2 \\ Z_2 \end{bmatrix} = \begin{bmatrix} X_1 \\ Y_1 \\ Z_1 \end{bmatrix} + \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

# How to Transform a 3D Point from 1 Coordinate Frame to Another?



# Coordinate Transformation with Homogeneous Coordinates

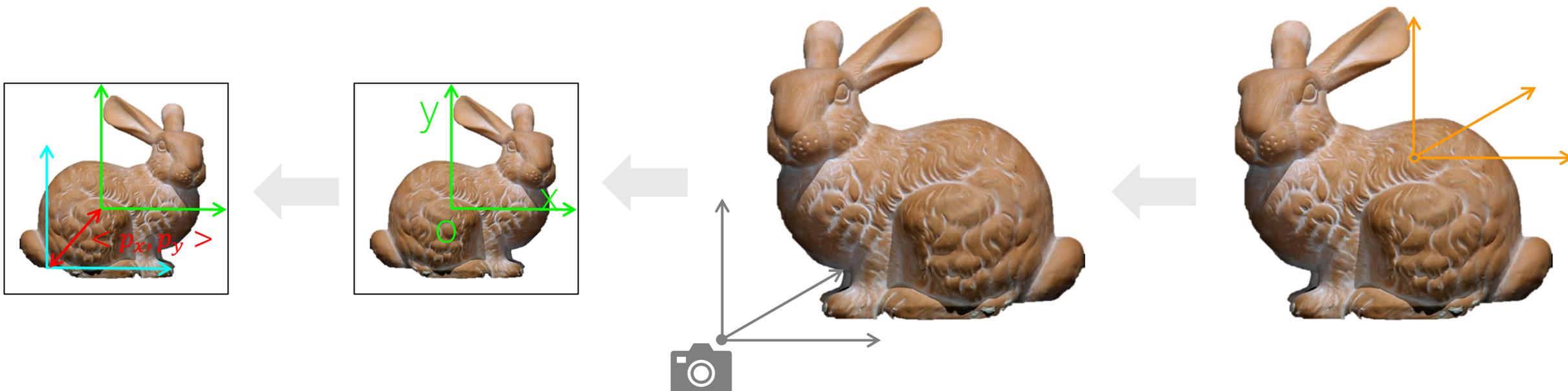
- We can express 3D coordinate transformation in a matrix form using homogeneous coordinates

$$\begin{bmatrix} X_{cam_2} \\ Y_{cam_2} \\ Z_{cam_2} \\ w \end{bmatrix} = \begin{bmatrix} R_{cam_1}^{cam_2} & T_{cam_1}^{cam_2} \\ \mathbf{0} & 1 \end{bmatrix} \begin{bmatrix} X_{cam_1} \\ Y_{cam_1} \\ Z_{cam_1} \\ 1 \end{bmatrix}$$

$M$

- $M$  is called extrinsic matrix
- In fact,  $M$  forms a special group  $SE(3)$

# Let's Put Everything Together



$$\begin{bmatrix} x_{pix} \\ y_{pix} \\ 1 \end{bmatrix} = \begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} R_w^c & T_w^c \\ \mathbf{0} & 1 \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$K_{img}^{pix}$

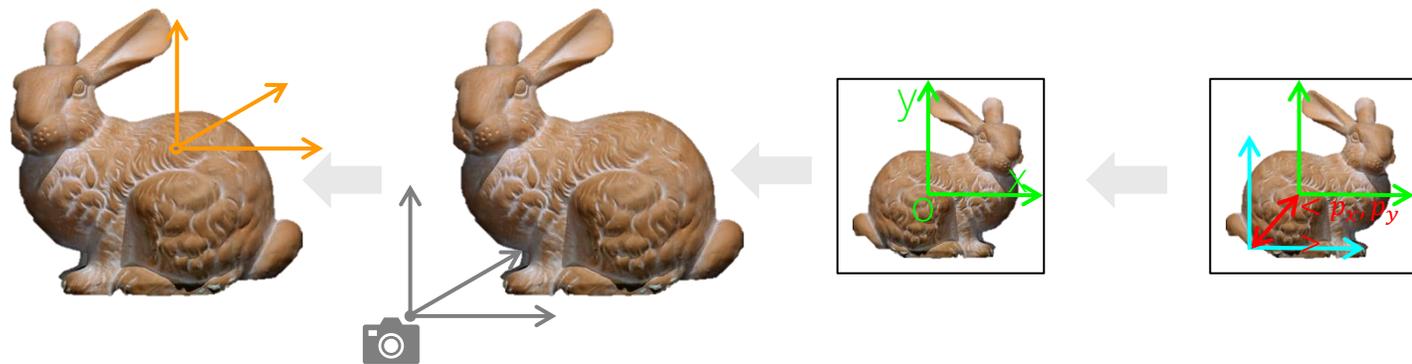
$K_{cam}^{img}$

$M_{world}^{cam}$

# Content

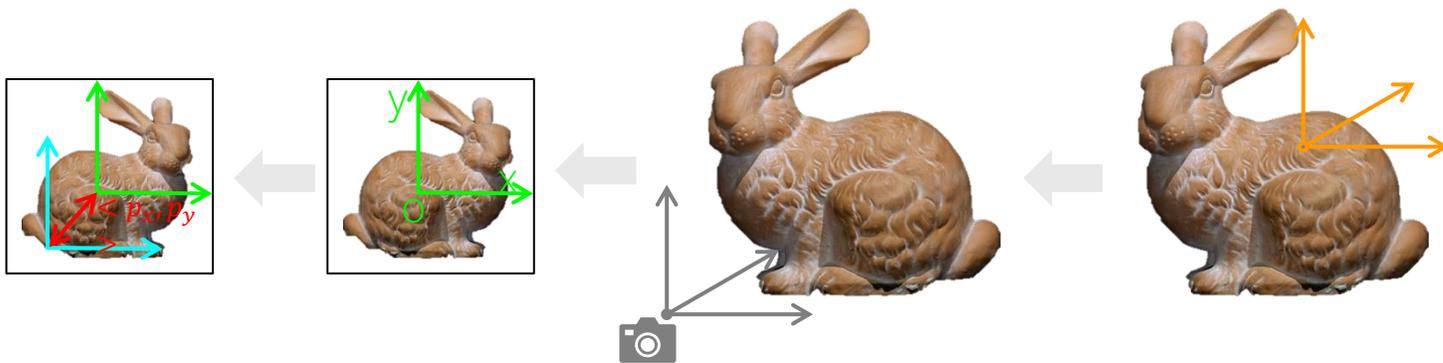
- Image Formation
  - Pinhole camera model
  - 2D transformation
  - 3D transformation
- Geometry Reconstruction
  - Reconstruction with depths
  - Multi-view geometry reconstruction with camera motions and correspondence
  - Epipolar geometry reconstruction

## 2D-to-3D reconstruction



$$\begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{R}_w^c & \mathbf{T}_w^c \\ \mathbf{0} & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}^{-1} \begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} x_{pix} \\ y_{pix} \\ 1 \end{bmatrix}$$

## 3D-to-2D perspective projection



$$\begin{bmatrix} x_{pix} \\ y_{pix} \\ 1 \end{bmatrix} = \begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{R}_w^c & \mathbf{T}_w^c \\ \mathbf{0} & 1 \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

# 3D Model Reconstruction from 2D Images

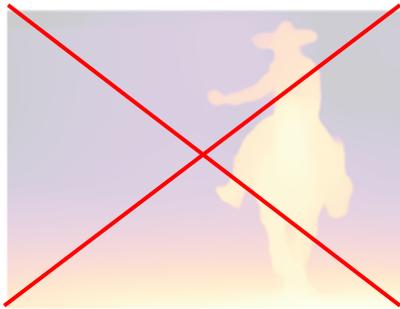


Image source <https://3dvision.princeton.edu/courses/SFMedu/>  
Image credit S. Tulsiani

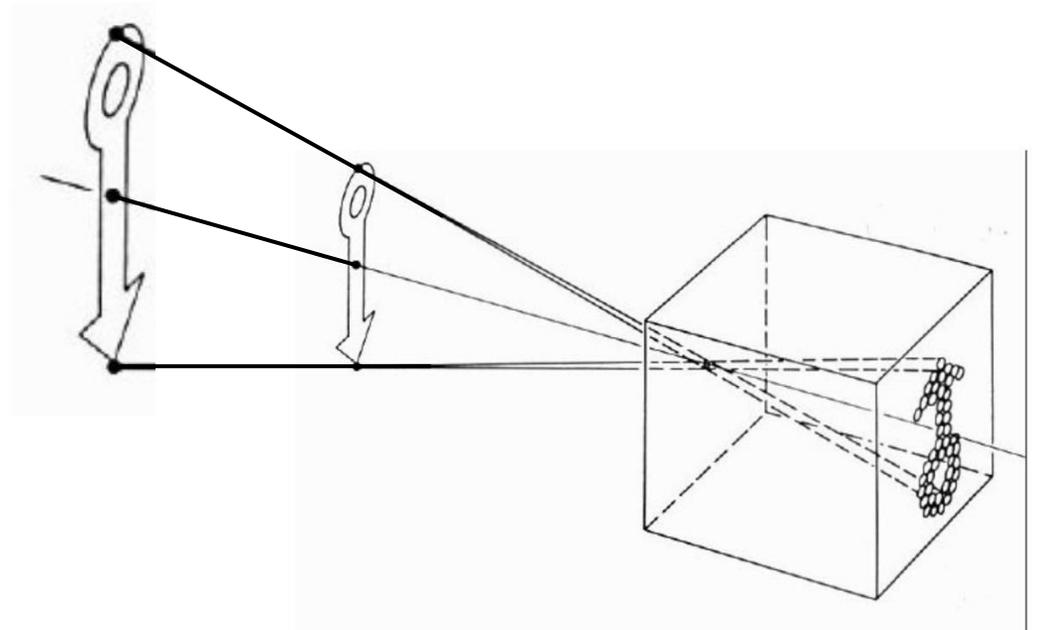
# Content

- Image Formation
  - Pinhole camera model
  - 2D transformation
  - 3D transformation
- Geometry Reconstruction
  - Reconstruction with depths
  - Multi-view geometry reconstruction with camera motions and correspondence
  - Epipolar geometry reconstruction

# We Said to Reconstruct 3D Models from a 2D Image with the Depth Map. What if We Don't have the Depth Map?



Depth map



2D to 3D is inherently ambiguous:

$$x = PX = PHH^{-1}X = \tilde{P}\tilde{X}$$

# Resolve Ambiguous 3D Model Reconstruction with Multiview 2D Images

## Stereoscope

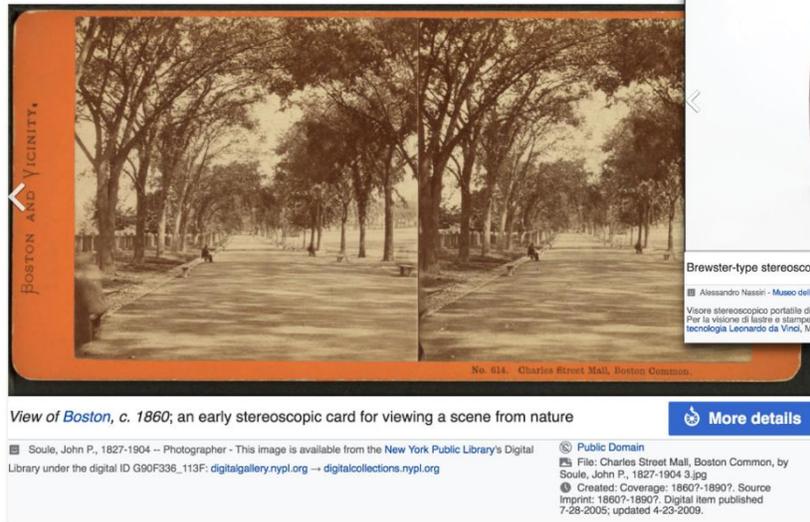


Image source <https://3dvision.princeton.edu/courses/SFMedu/>  
Image credit S. Tulsiani

Image credit A. Efors

# Resolve Ambiguous 3D Model Reconstruction with Multiview 2D Images



(a)

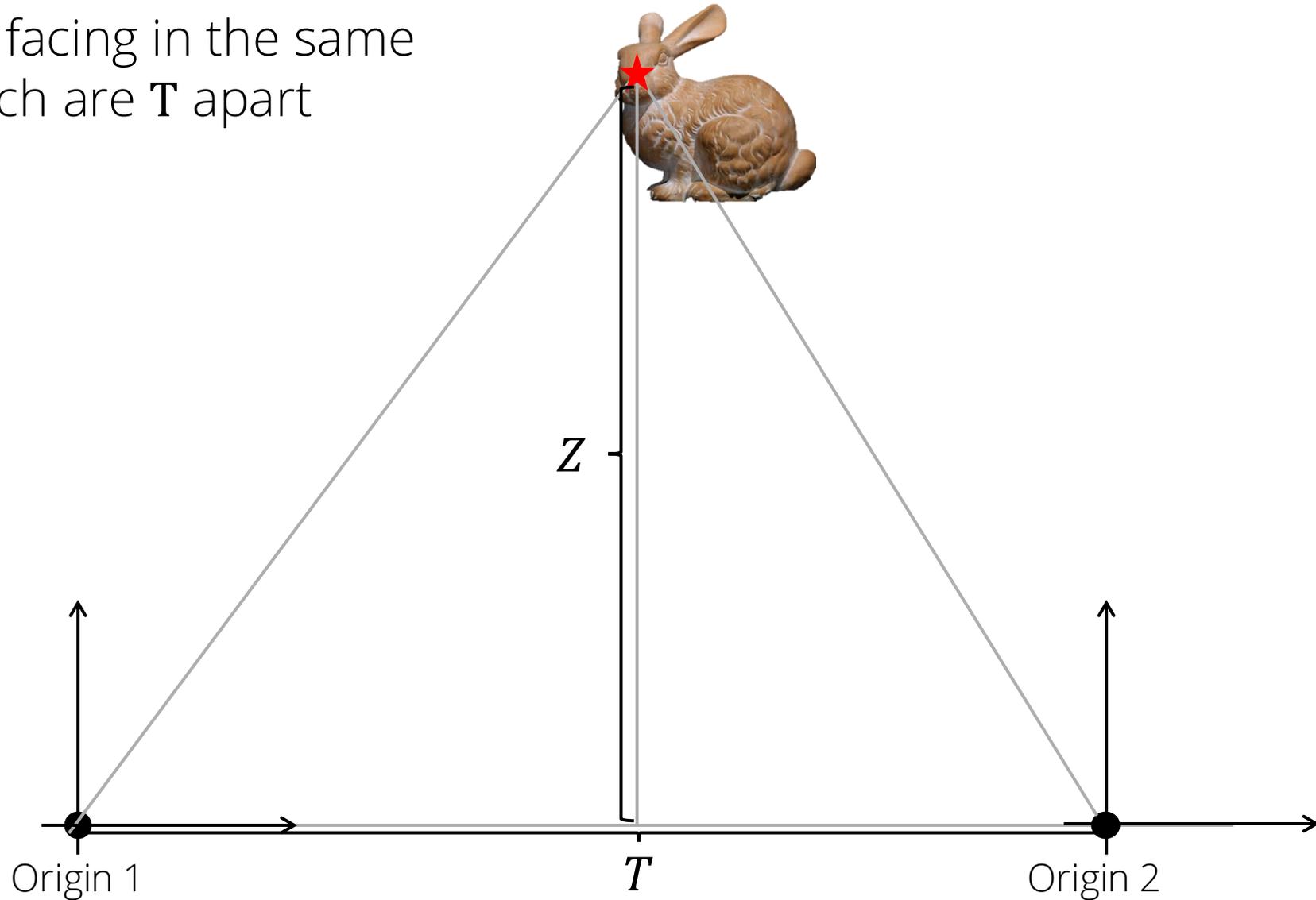


(b)

Figure 1.1: (a) Stereo anaglyph of the ocean liner, the Titanic [McManus2022]. The red image shows the right eye's view, and cyan the left eye's view. When viewed through stereo red/cyan stereo glasses, as in (b), the cyan contrast appears in the left eye image and the red variations appear to the right eye, creating a the perception of 3d.

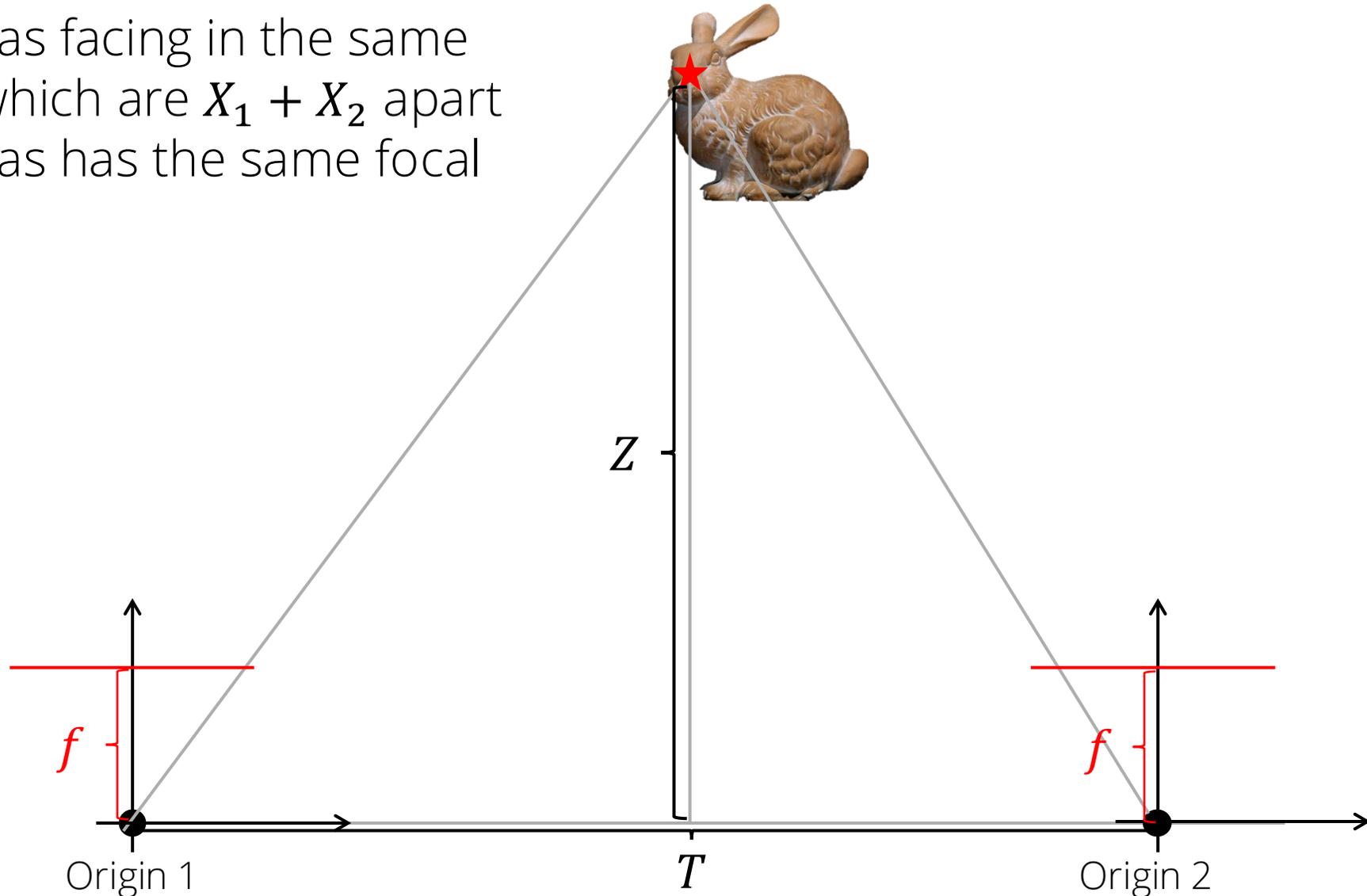
# Geometry for a Simple Stereo System

- Two cameras facing in the same direction, which are  $T$  apart



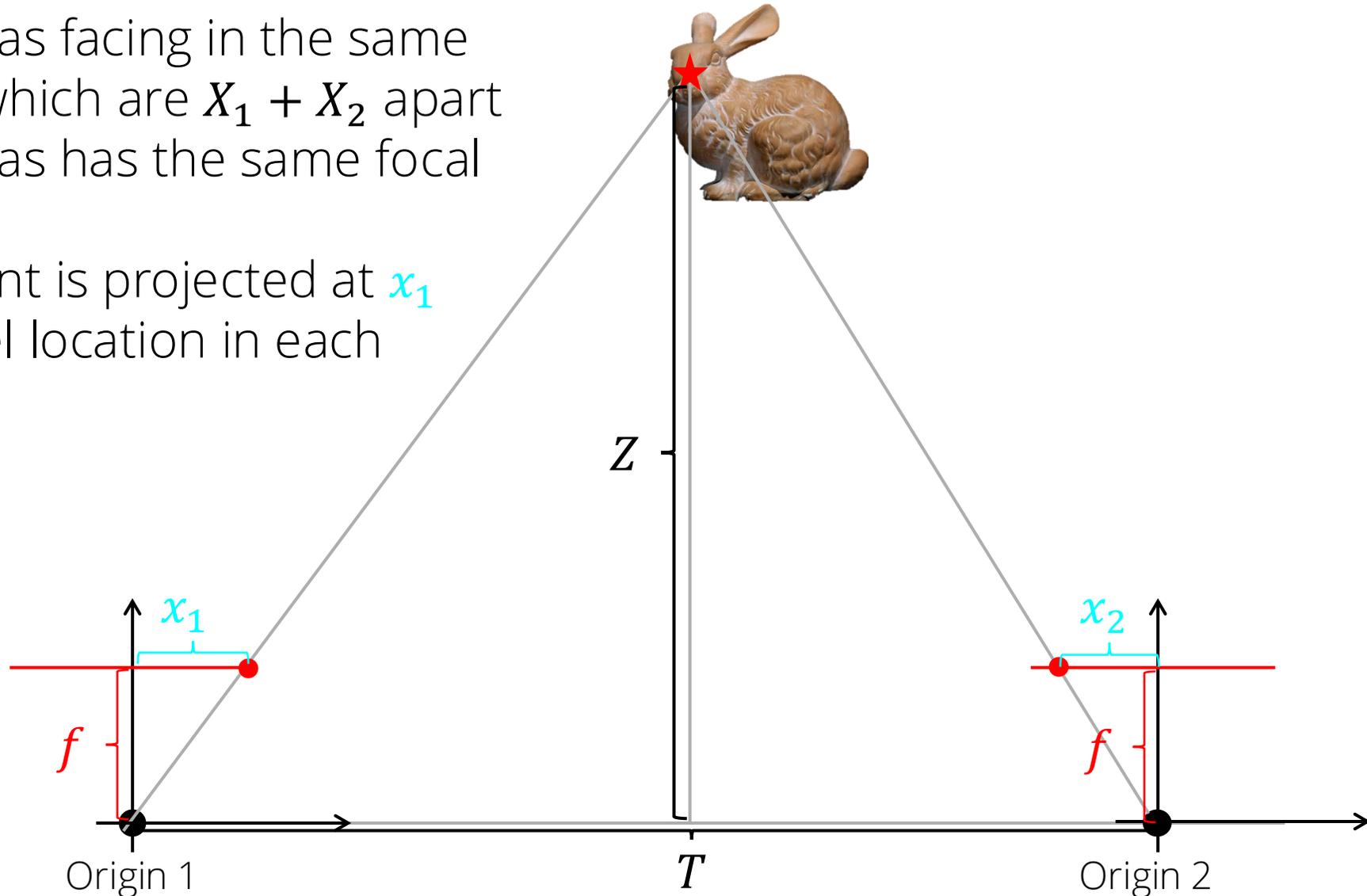
# Geometry for a Simple Stereo System

- Two cameras facing in the same direction, which are  $X_1 + X_2$  apart
- Two cameras has the same focal length



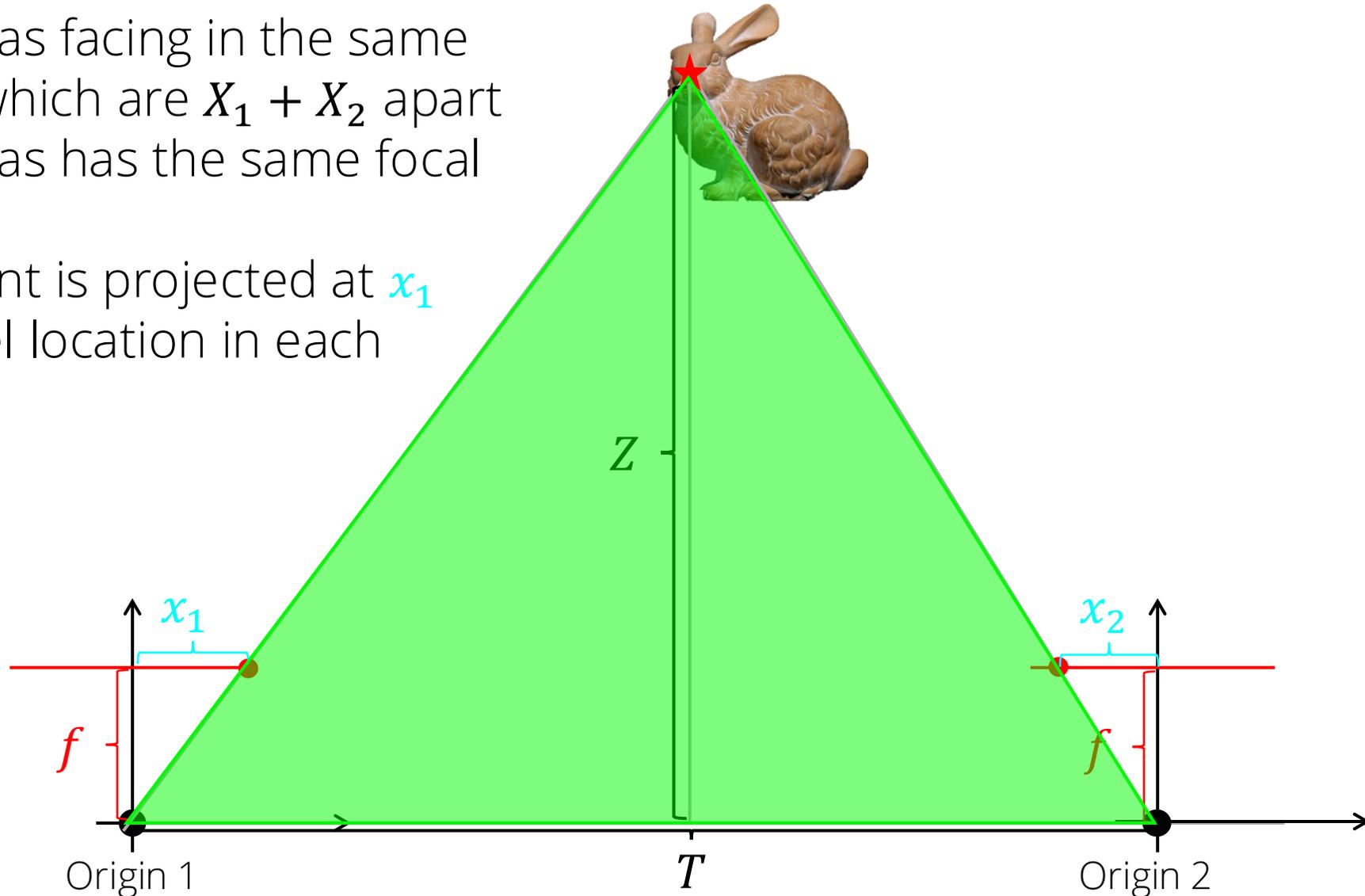
# Geometry for a Simple Stereo System

- Two cameras facing in the same direction, which are  $X_1 + X_2$  apart
- Two cameras has the same focal length
- The 3D point is projected at  $x_1$  and  $x_2$  pixel location in each camera



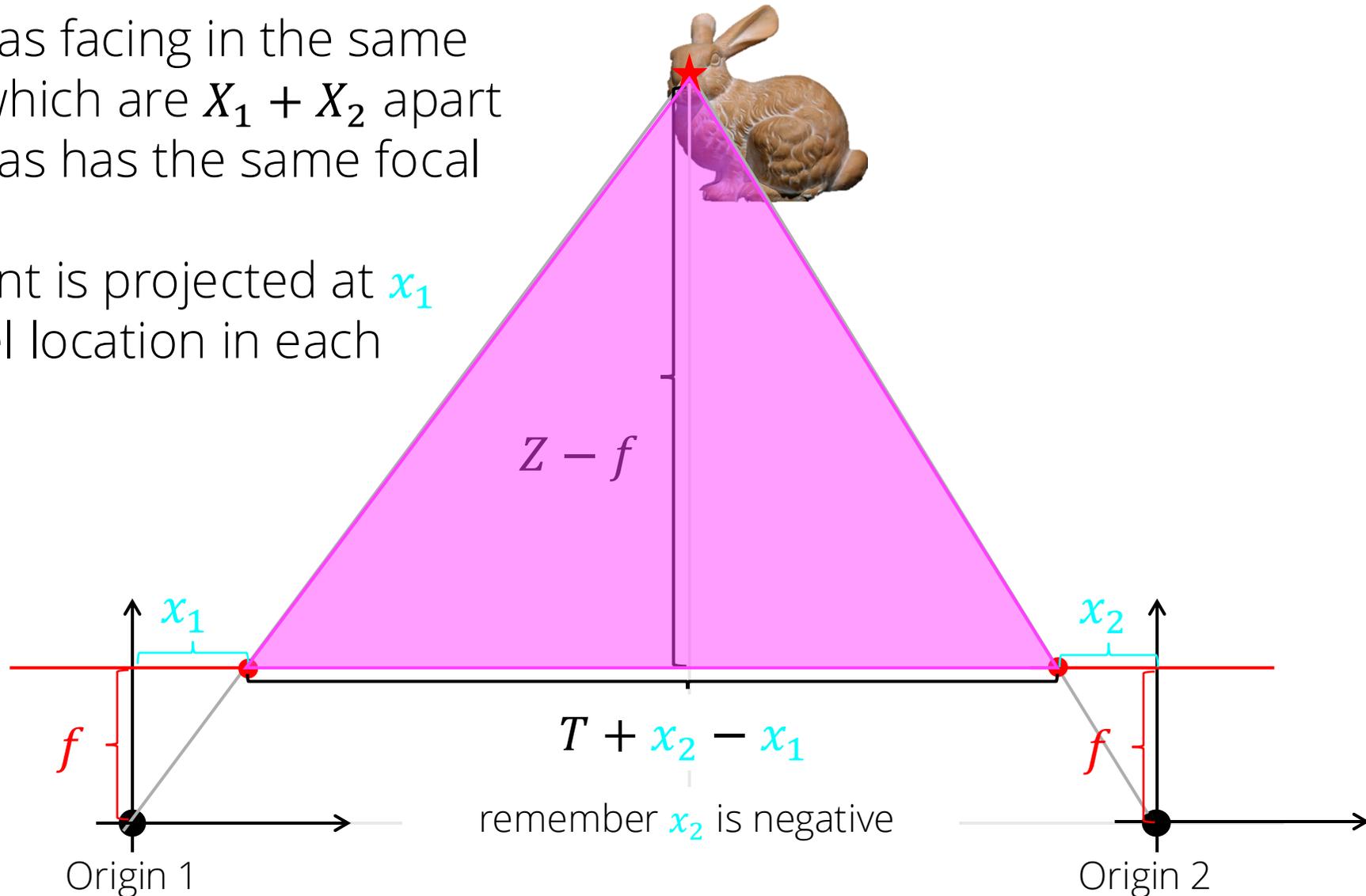
# Geometry for a Simple Stereo System

- Two cameras facing in the same direction, which are  $X_1 + X_2$  apart
- Two cameras has the same focal length
- The 3D point is projected at  $x_1$  and  $x_2$  pixel location in each camera



# Geometry for a Simple Stereo System

- Two cameras facing in the same direction, which are  $X_1 + X_2$  apart
- Two cameras has the same focal length
- The 3D point is projected at  $x_1$  and  $x_2$  pixel location in each camera

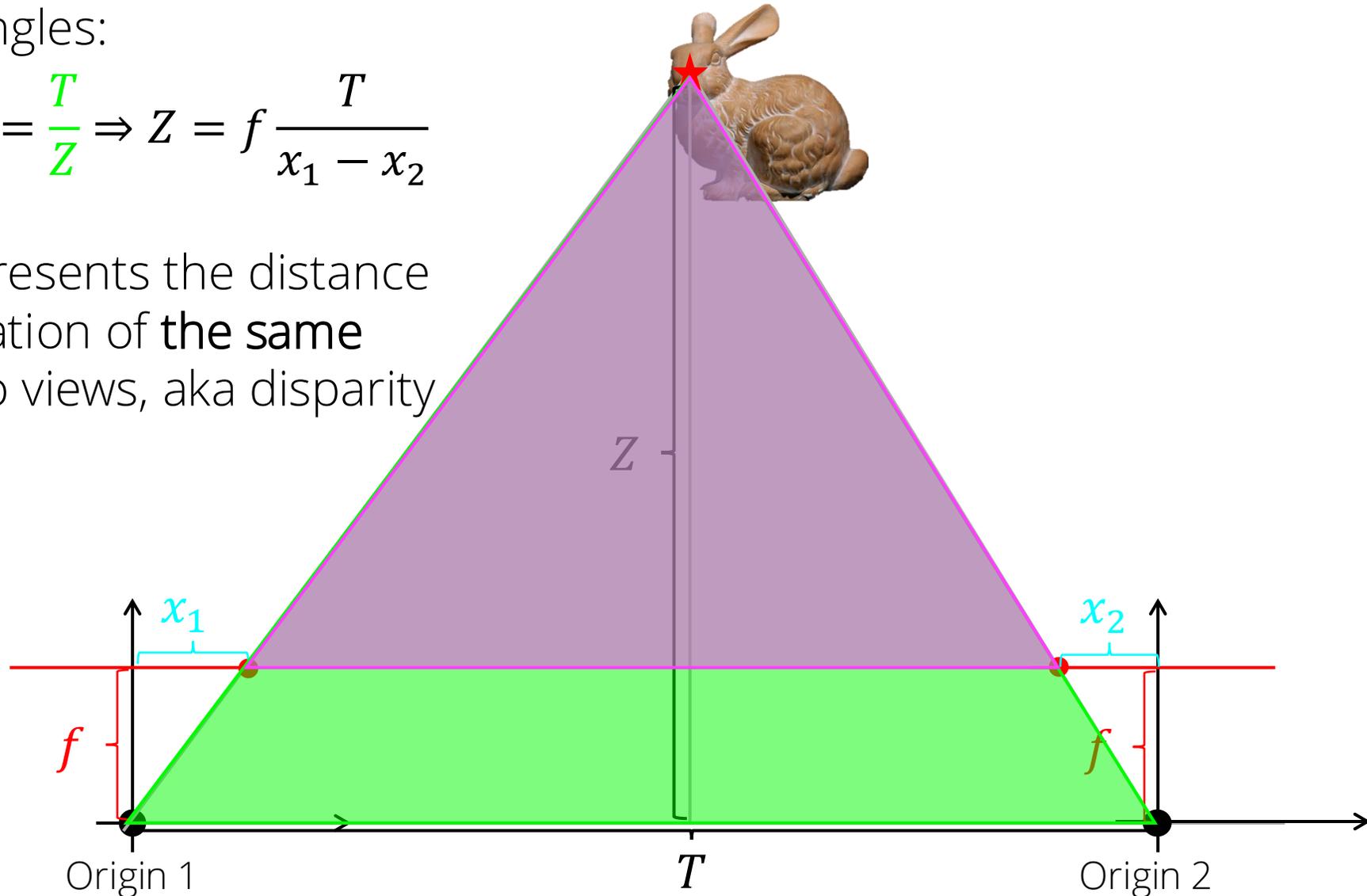


# Geometry for a Simple Stereo System

- Similar triangles:

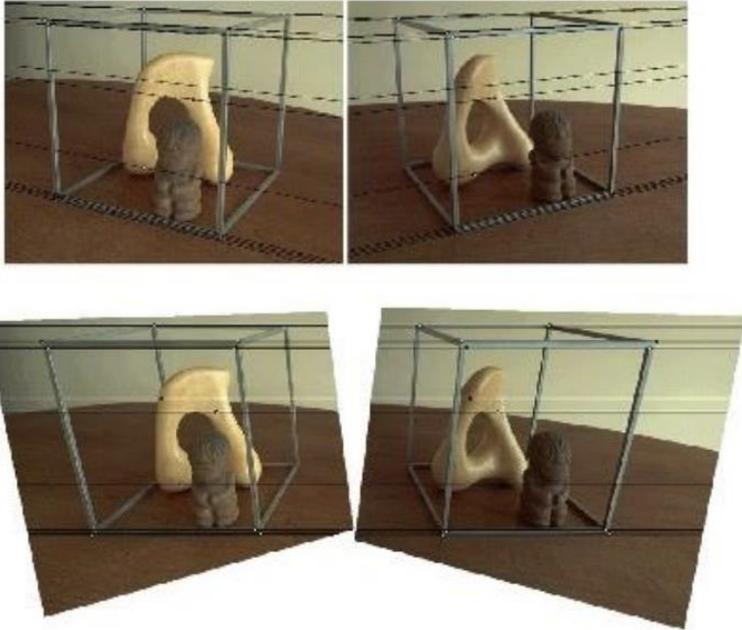
$$\frac{T + x_2 - x_1}{Z - f} = \frac{T}{Z} \Rightarrow Z = f \frac{T}{x_1 - x_2}$$

- $x_1 - x_2$  represents the distance of pixel location of **the same point** in two views, aka disparity



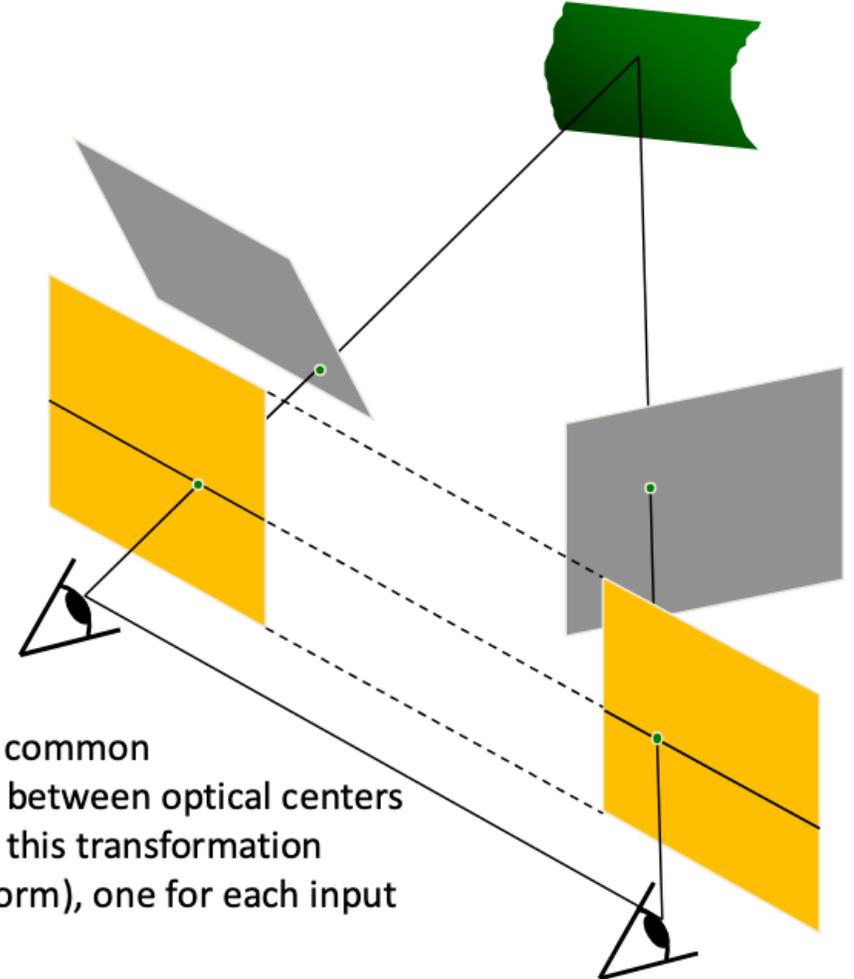
# What if the Two Cameras Are Not Parallel?

## Rectification example



- We'll talk about how to obtain the homography projection that makes 2 cameras parallel next

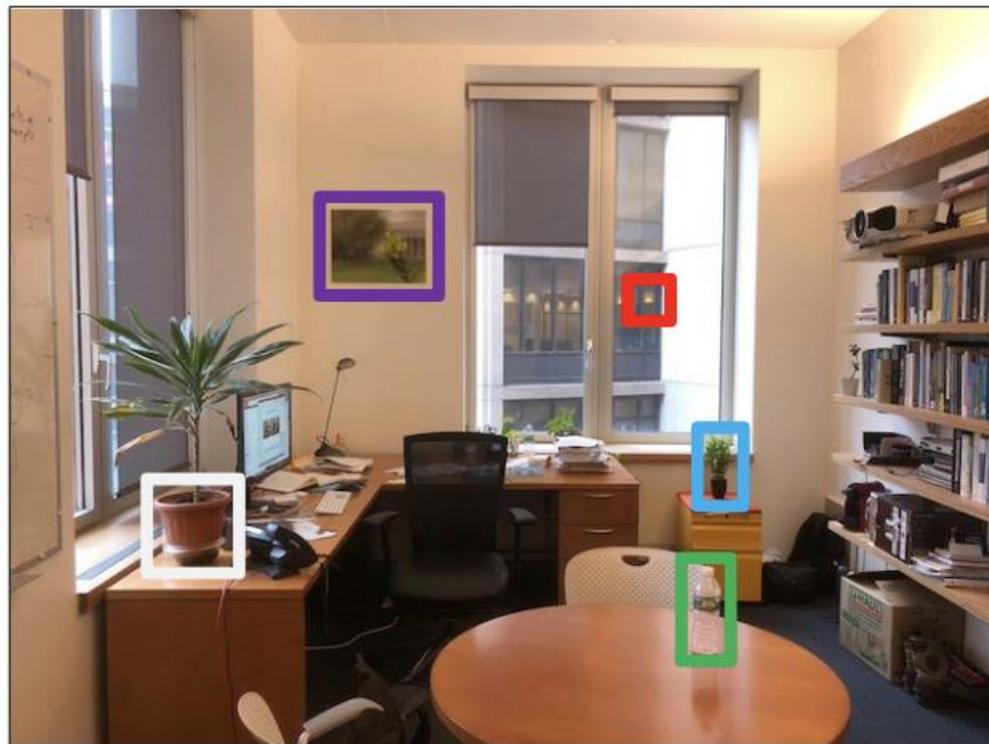
## Stereo image rectification



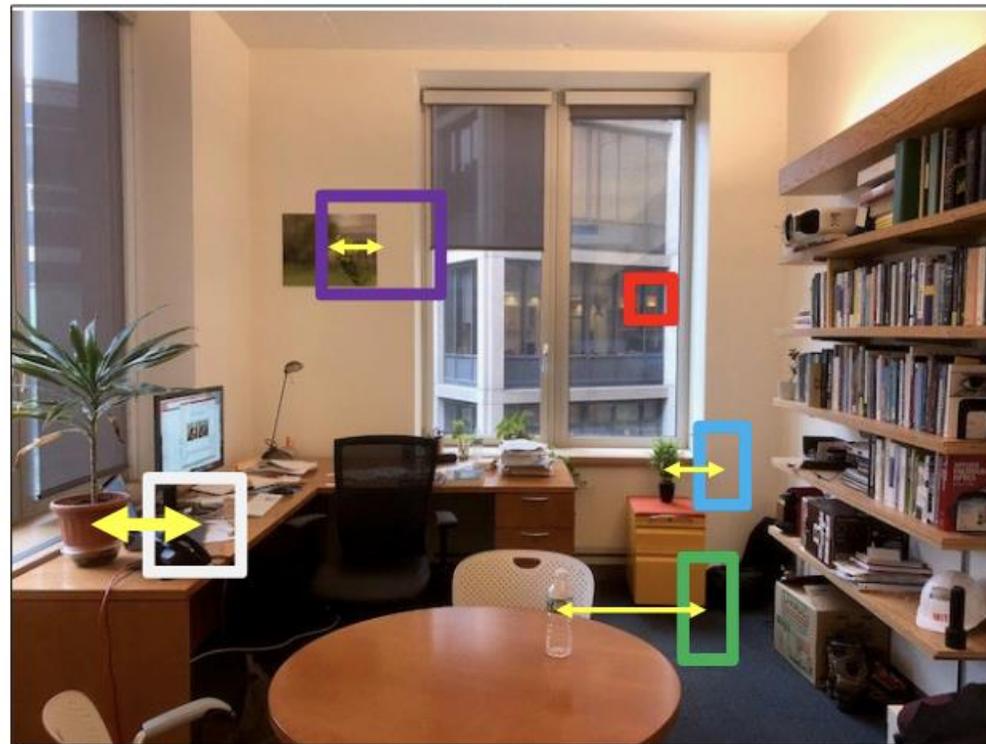
- Reproject image planes onto a common plane parallel to the line between optical centers
- Pixel motion is horizontal after this transformation
- Two homographies (3x3 transform), one for each input image reprojection

•C. Loop and Z. Zhang. [Computing Rectifying Homographies for Stereo Vision](#). IEEE Conf. Computer Vision and Pattern Recognition, 1999.

Disparity Map: the distance of pixel location of the same point in two views



Left image



Right image

# Disparity Map: the distance of pixel location of the same point in two views

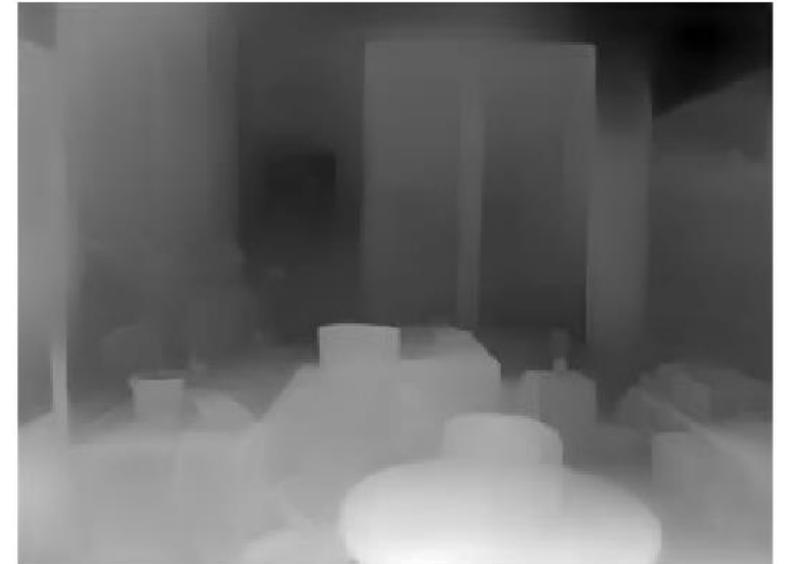
$I(x,y)$



$I'(x,y) = I(x+D(x,y), y)$



$D(x,y)$



- Distance derived from the disparity map

$$Z(x,y) = \frac{a}{D(x,y)}$$

# Disparity Map: the distance of pixel location of the same point in two views

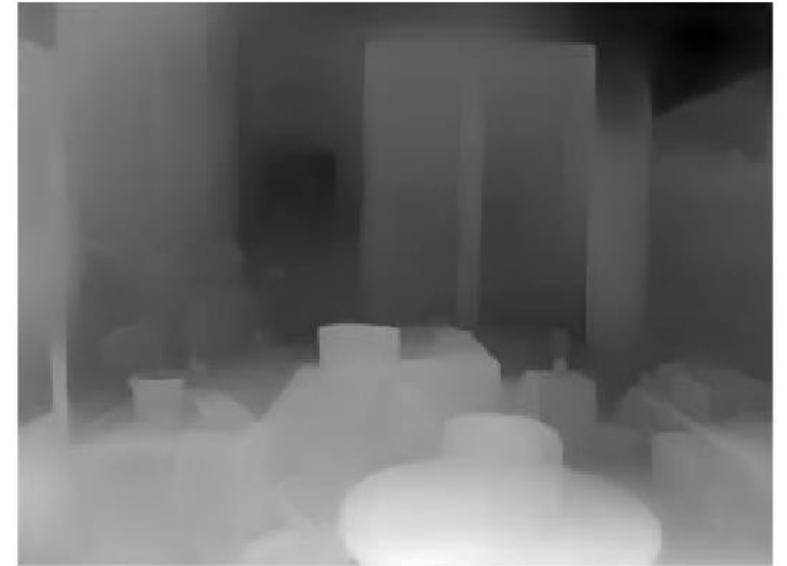
$I(x,y)$



$I'(x,y) = I(x+D(x,y), y)$



$D(x,y)$

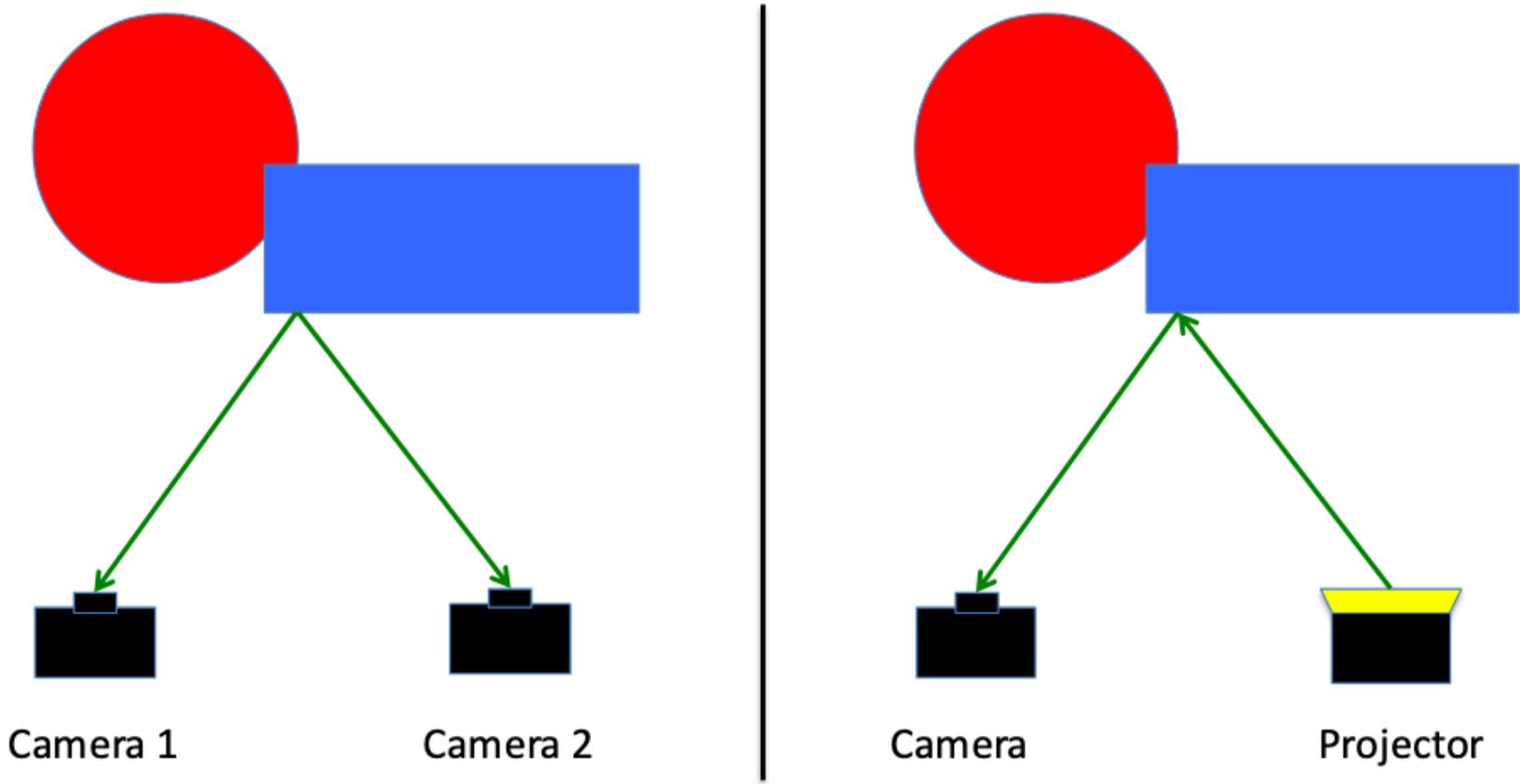


- Disparity map is derived from pixel correspondence
- How to derive pixel correspondence?

- Distance derived from the disparity map

$$Z(x,y) = \frac{a}{D(x,y)}$$

# Depth from Triangulation



Passive Stereopsis

Active Stereopsis

Active sensing simplifies the problem of estimating point correspondences

# Range Sensors



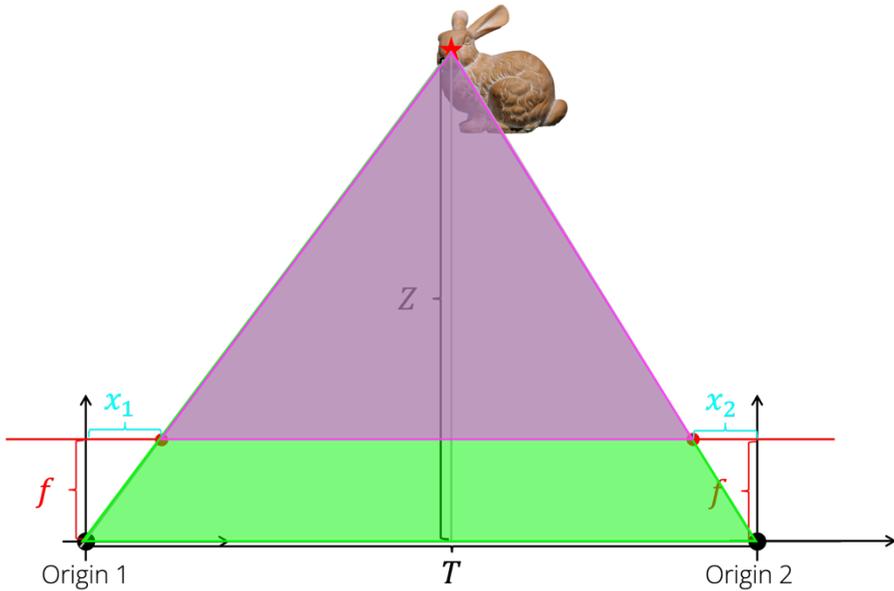
primesense sensor (used in Kinect)



Velodyne LIDAR Sensor

<http://www.primesense.com/>, <http://www.ifixit.com/>,  
[http://mirror.umd.edu/roswiki/kinect\\_calibration\(2f\)technical.html](http://mirror.umd.edu/roswiki/kinect_calibration(2f)technical.html)  
<http://velodynelidar.com/lidar/lidar.aspx>

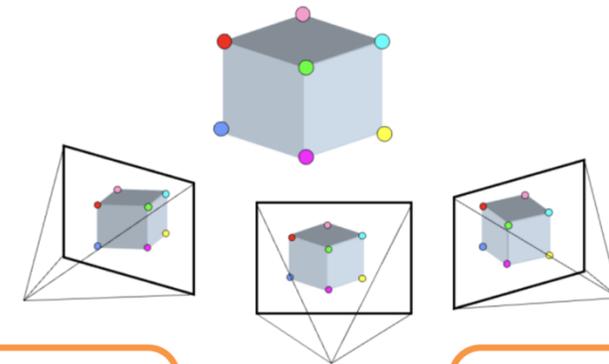
# Let's Take a Short Break



$$Z = f \frac{T}{x_1 - x_2}$$

- $x_1 - x_2$  represents the disparity, which requires pixel correspondence
- $T$  represents the relative camera pose

3D Points  
(Structure)



Correspondences

Camera  
(Motion)

Image credit A. Efros

# Content

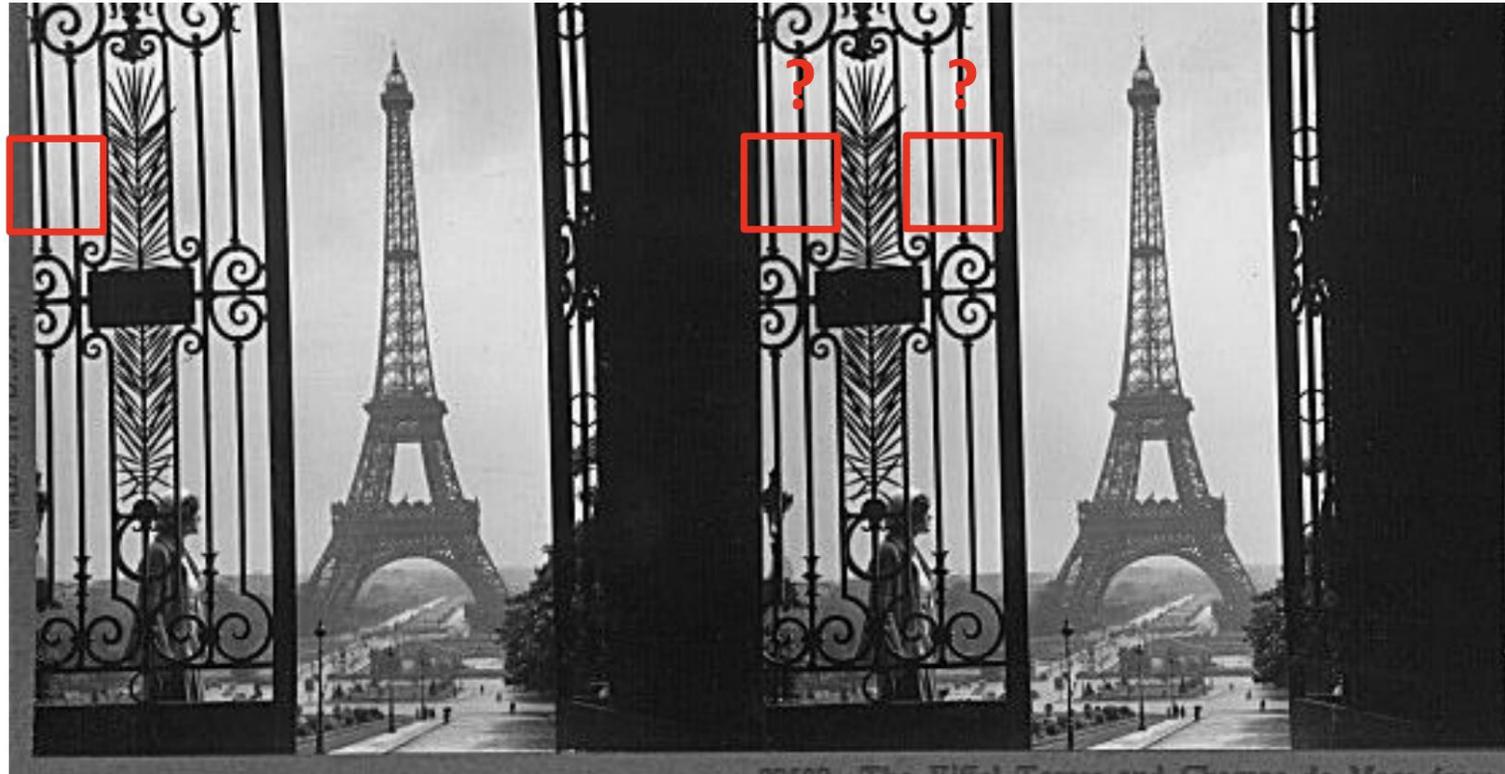
- Image Formation
  - Pinhole camera model
  - 2D transformation
  - 3D transformation
- Geometry Reconstruction
  - Reconstruction with depths
  - Multi-view geometry reconstruction with camera motions and correspondence
  - Epipolar geometry reconstruction

# How Can I Find the Corresponding Pixel in Another View?



- We assume the camera extrinsics are known
- Idea:
  1. Extract features: SIFT / Deep Neural Networks
  2. Do Nearest Neighbor Search on all pixels in another view

# Failures of Correspondence Search



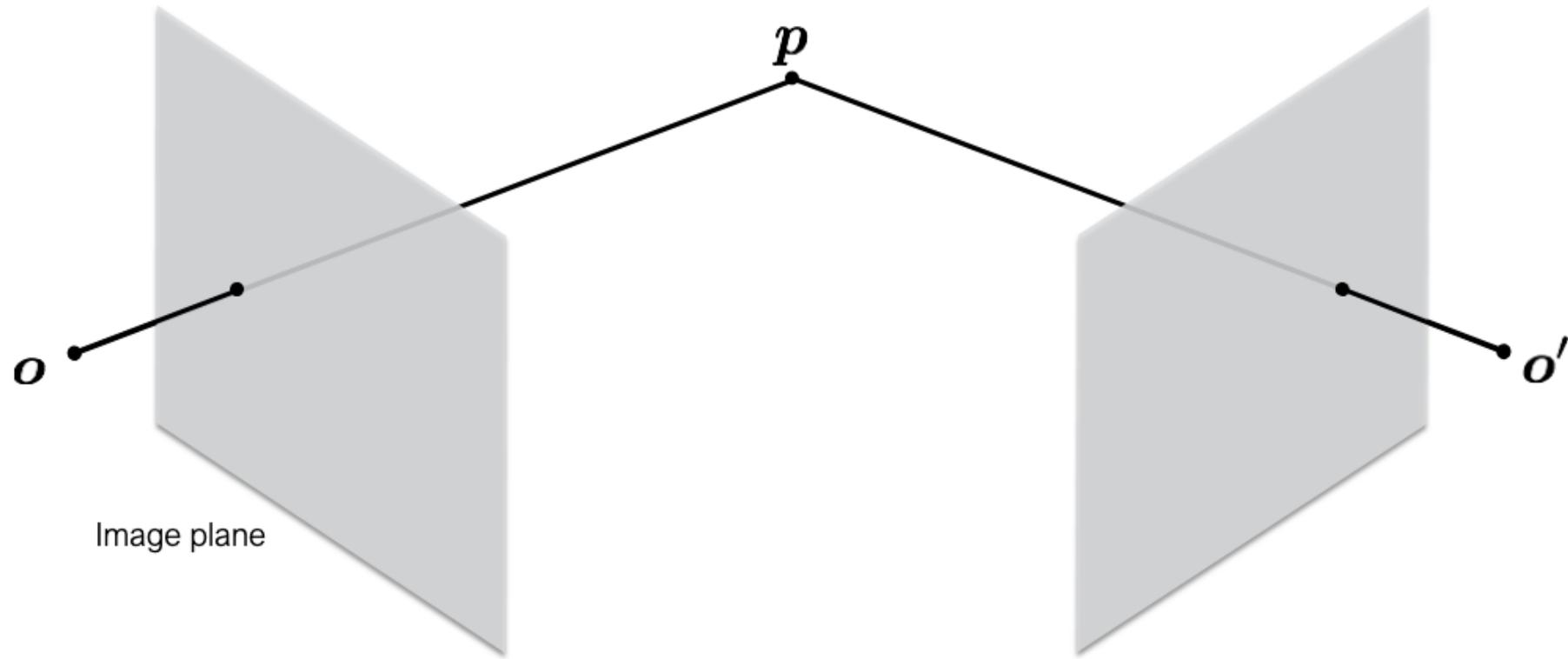
Is there any geometrical structure embedded in multiple views, that could facilitate correspondence search?

# If Camera Poses are Known, the Search Space of Pixel Correspondence is Constraint

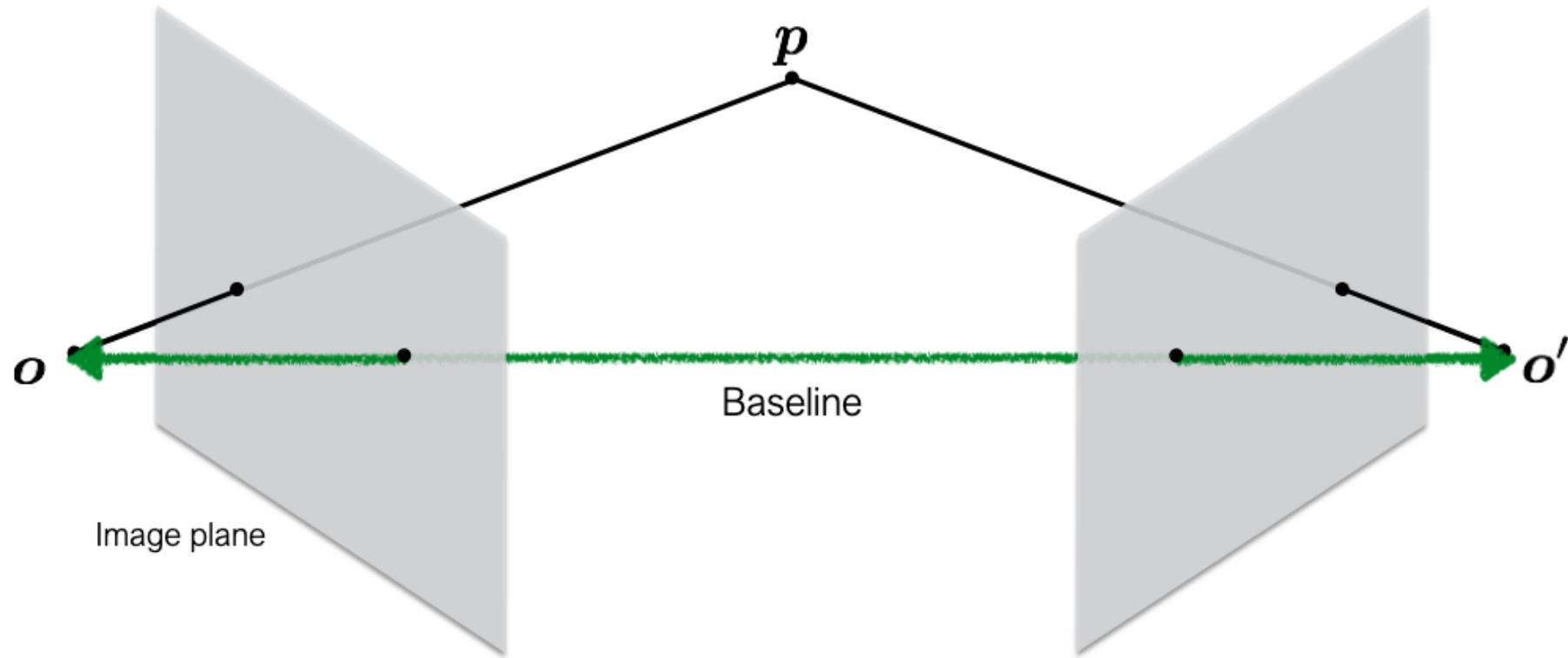


We don't need to search over the whole image, but the patches along the cyan line.

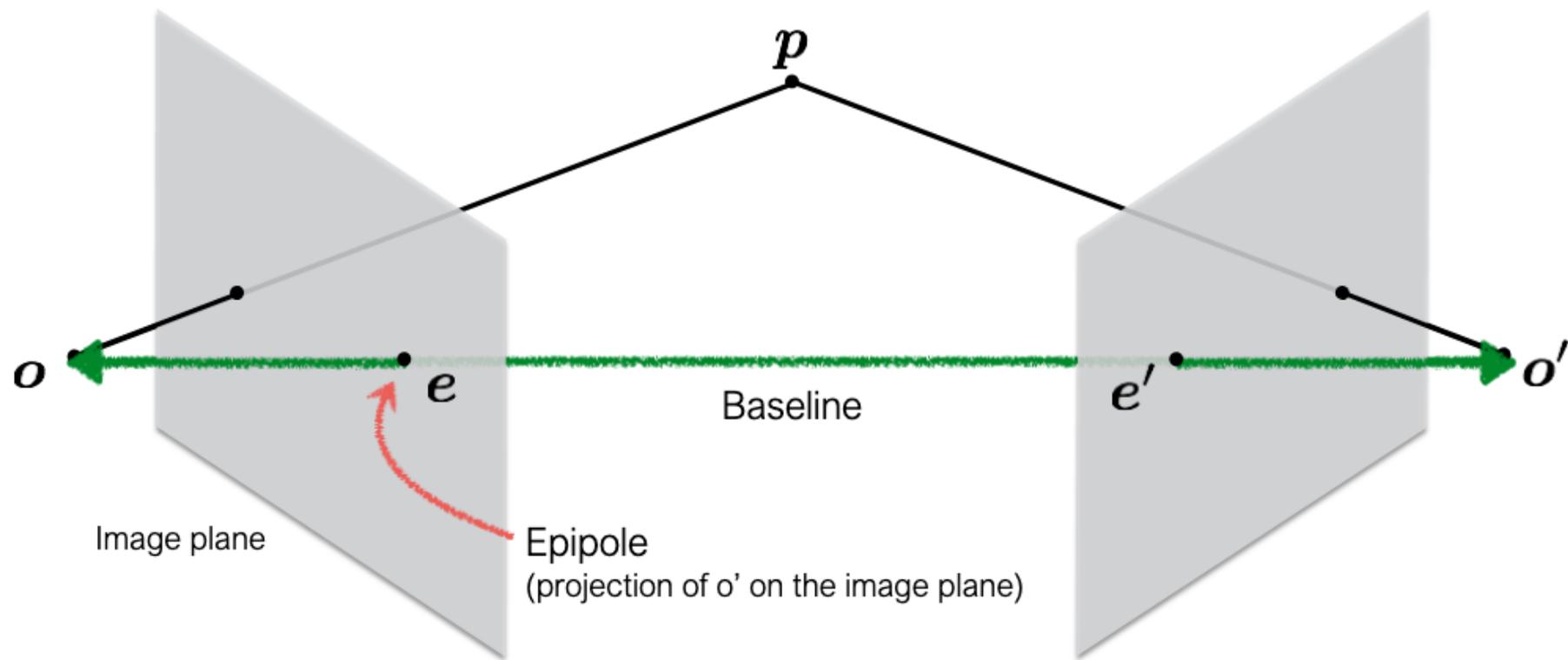
# Epipolar geometry



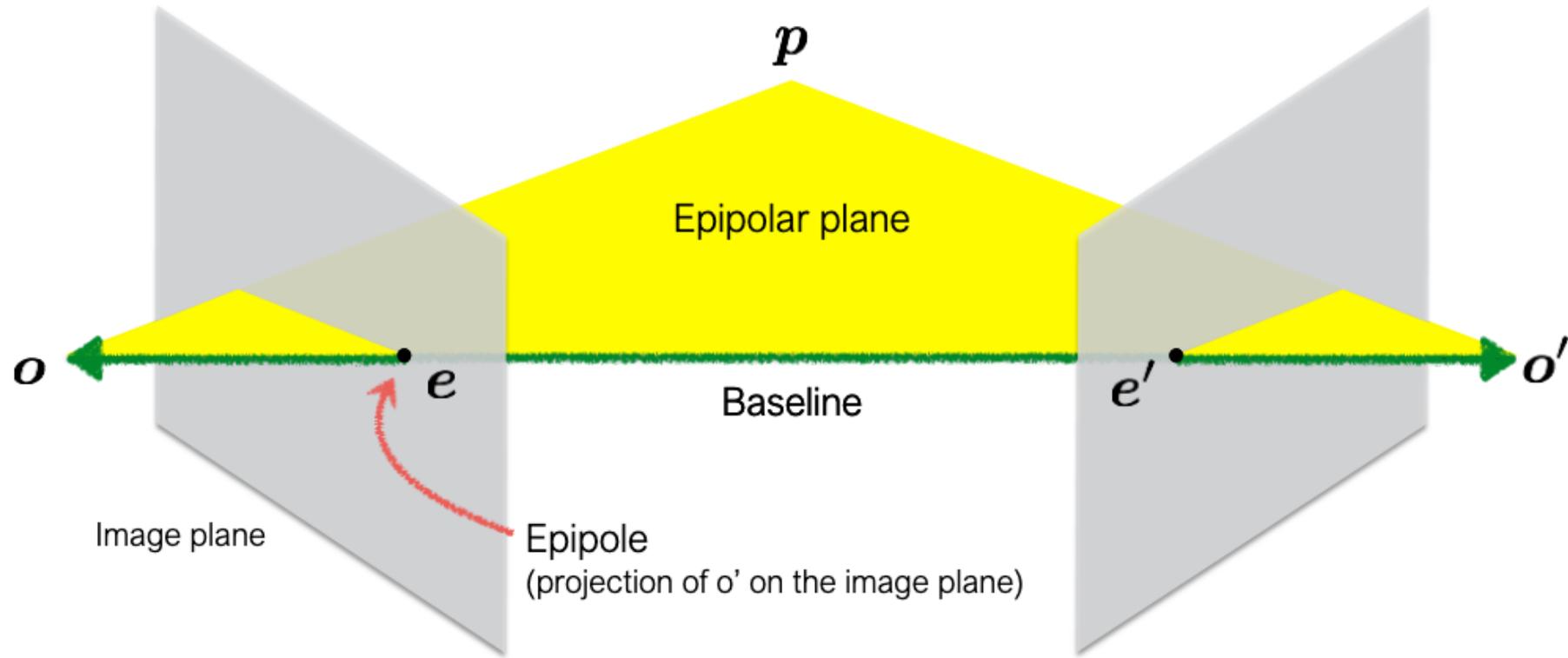
# Epipolar geometry



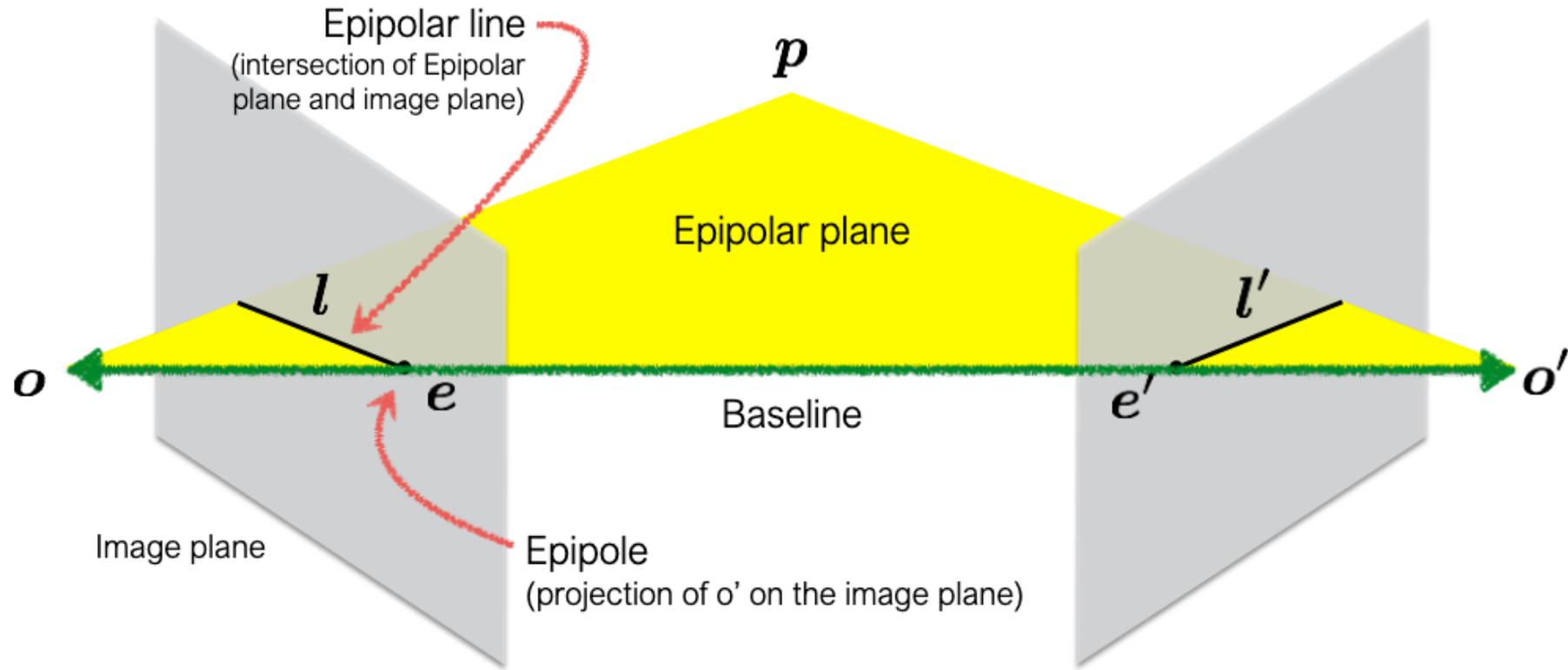
# Epipolar geometry



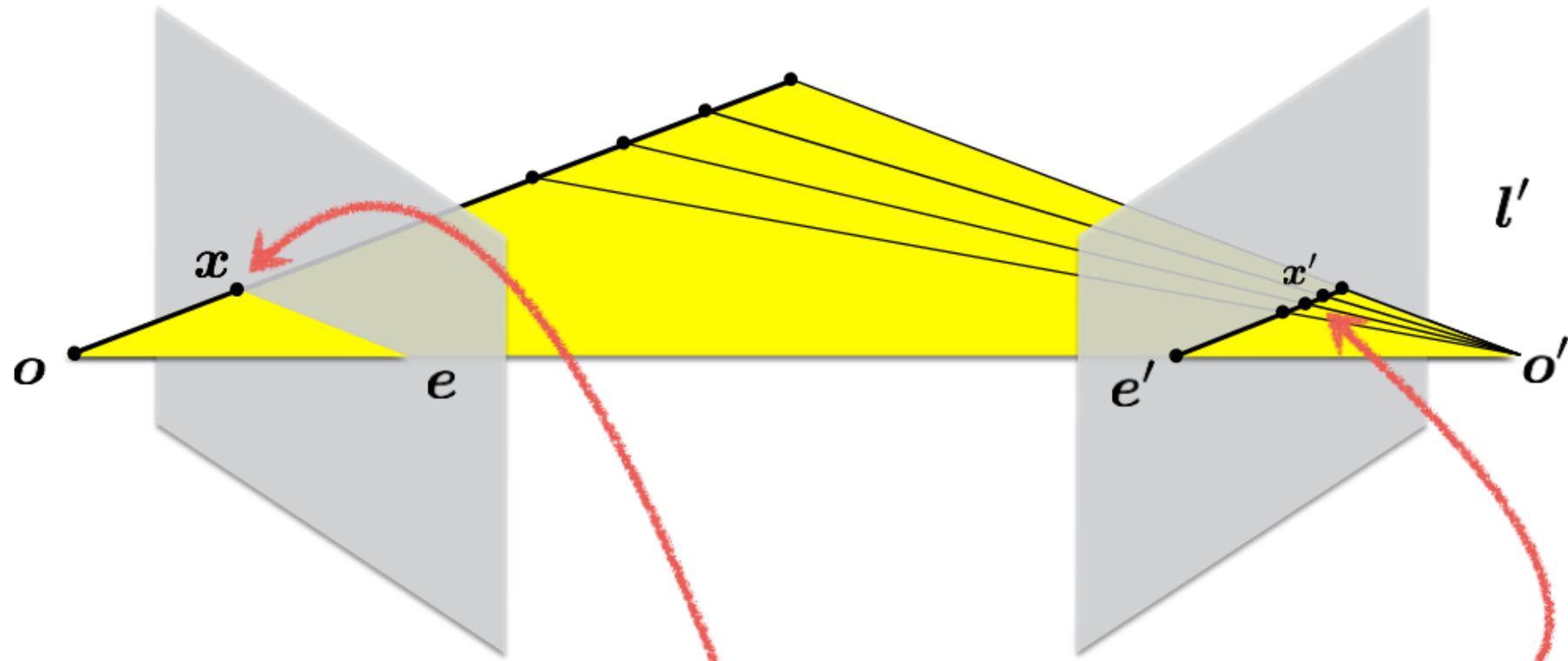
# Epipolar geometry



# Epipolar geometry



# Epipolar constraint



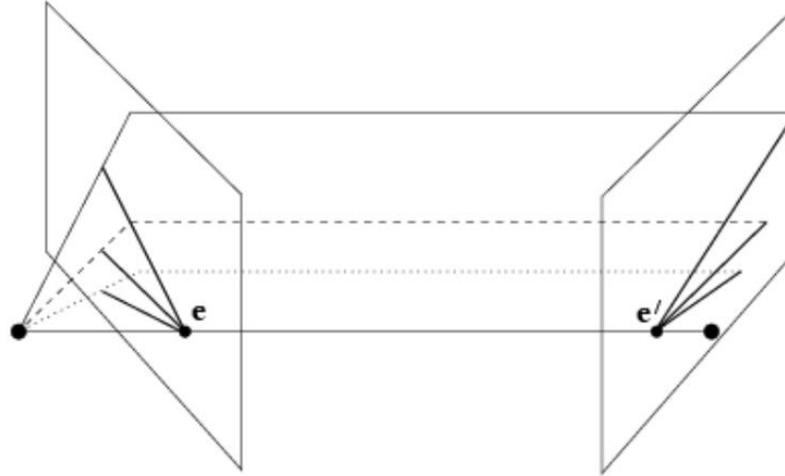
Potential matches for  $x$  lie on the epipolar line  $l'$

# If Camera Poses are Known, the Search Corresponded Pixels Along the Epipolar Line



We don't need to search over the whole image, but the patches along the **epipolar** line.

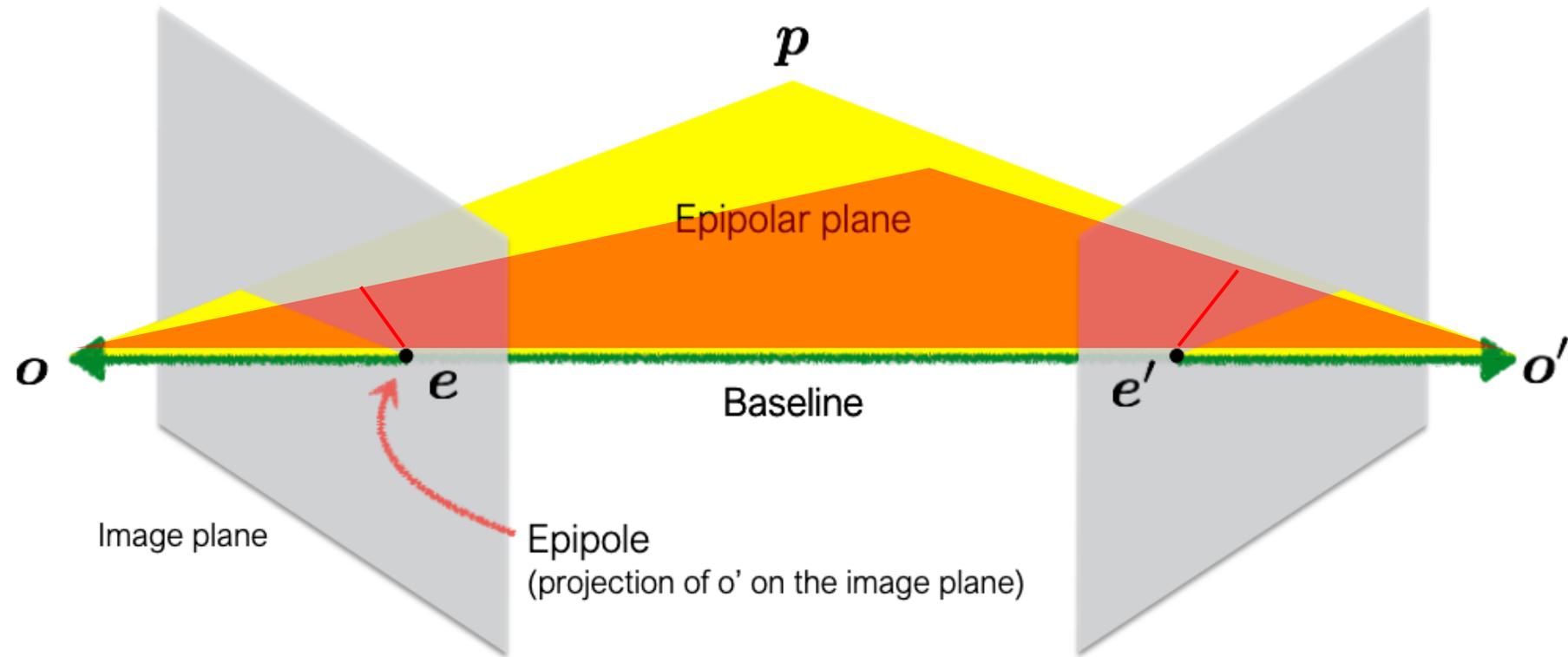
# Example: converging cameras



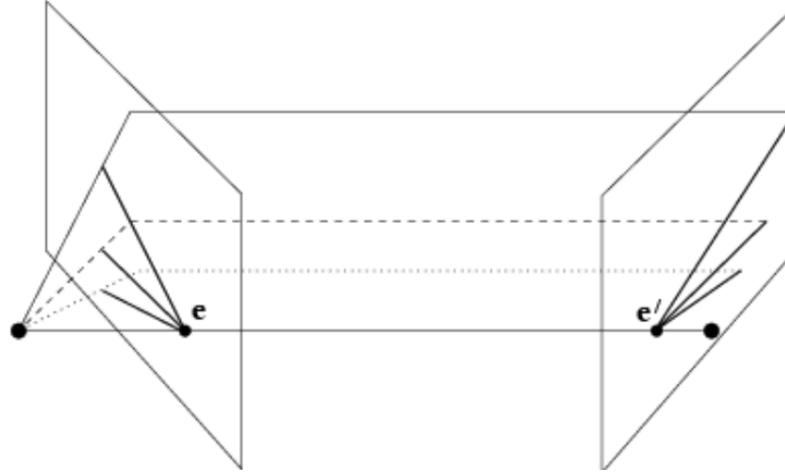
As position of 3d point varies, epipolar lines “rotate” about the baseline



# Example: Converging Camera



# Example: converging cameras

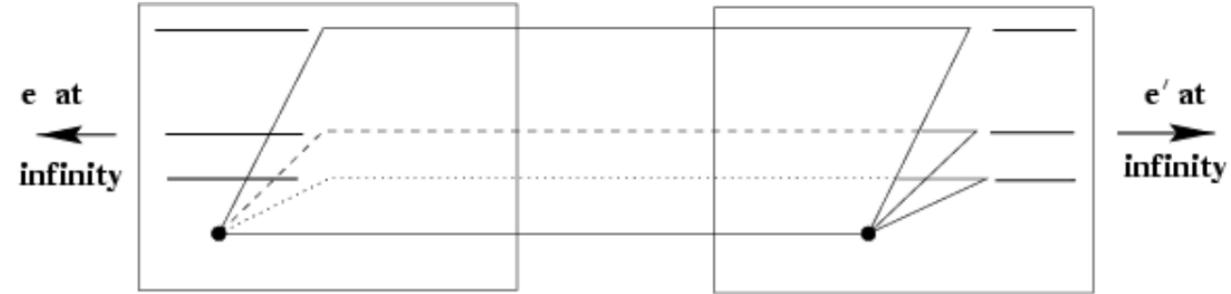


As position of 3d point varies, epipolar lines “rotate” about the baseline



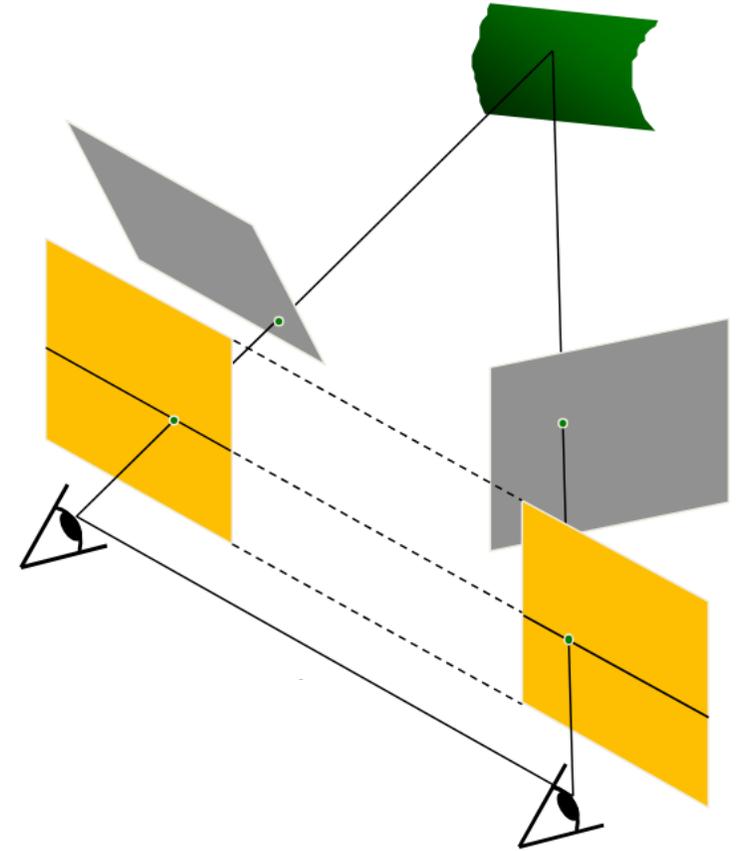
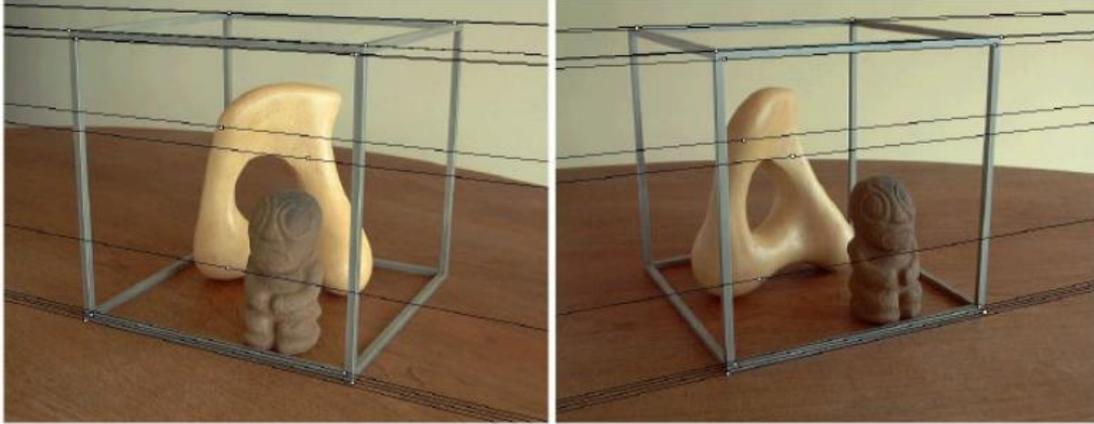
Here is the Epipole!

# Example: motion parallel with image plane



Epipolar lines are horizontally parallel

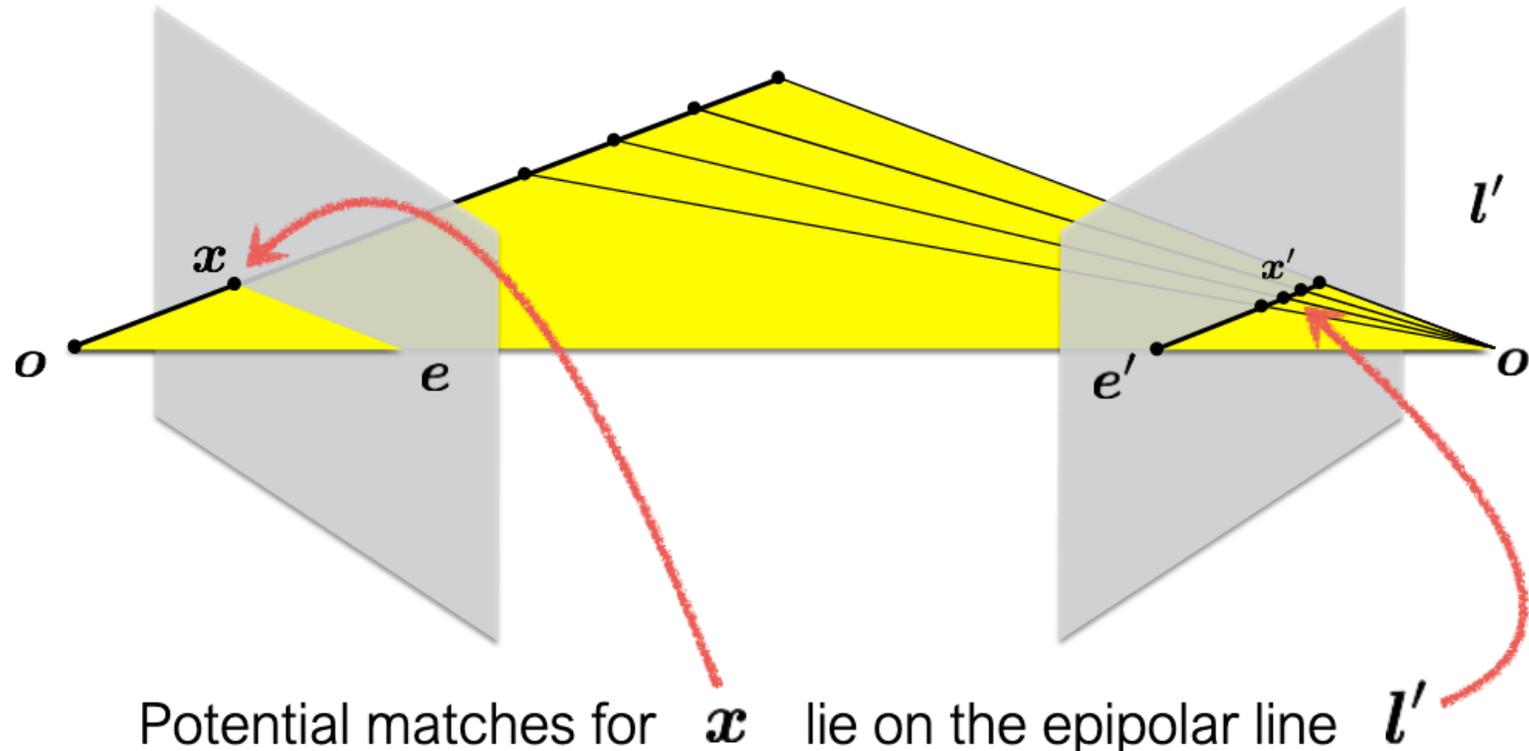
# Stereo Image Rectification



# Can We Represent the Epipolar Constraints Mathematically?

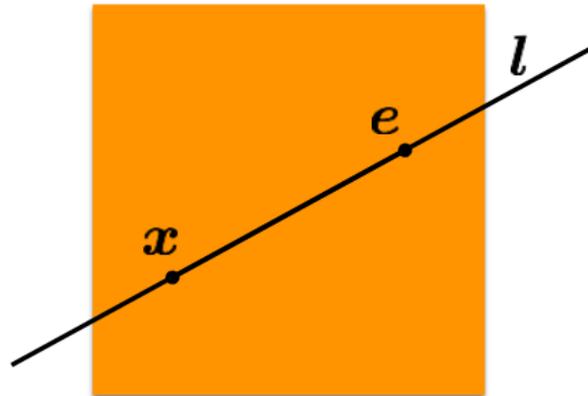
Given a pixel  $x$  in one view, describe the epipolar line  $l'$  in another view

- If the camera poses are known,  $l' = Ex$ , where  $E$  is the essential matrix
- If the camera poses are unknown,  $l' = Fx$ , where  $F$  is the fundamental matrix



# How to Represent a Line?

$$ax + by + c = 0 \quad \text{in vector form} \quad \mathbf{l} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

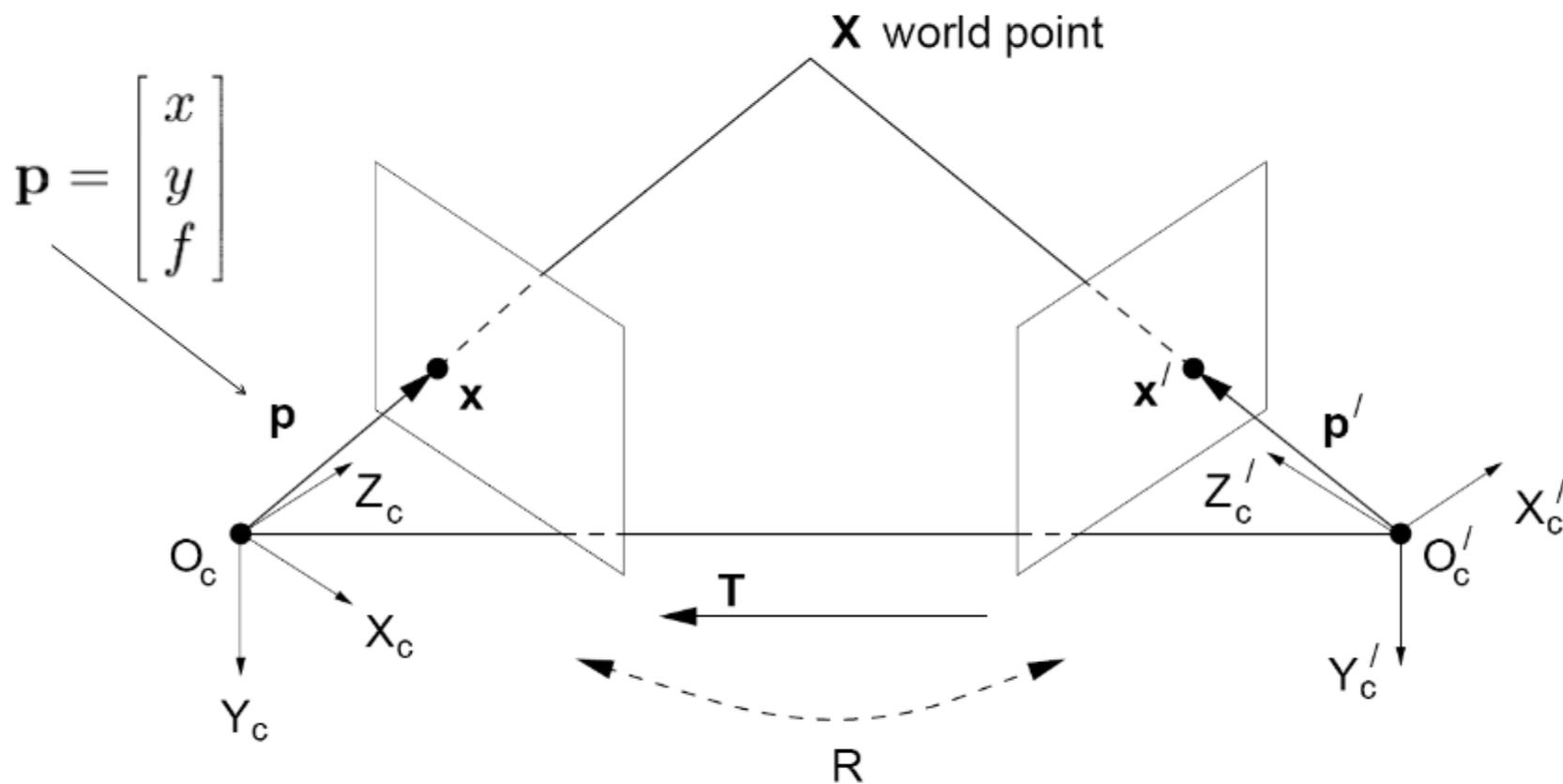


If the point  $\mathbf{x}$  is on the epipolar line  $\mathbf{l}$  then

$$\mathbf{x}^\top \mathbf{l} = 0$$

# Deriving the Essential Matrix:

## Stereo geometry, with calibrated cameras



Camera-centered coordinate systems are related by known rotation  $\mathbf{R}$  and translation  $\mathbf{T}$ :

$$\mathbf{X}' = \mathbf{R}\mathbf{X} + \mathbf{T}$$

# Review: Cross product

$$\vec{a} \times \vec{b} = \vec{c}$$

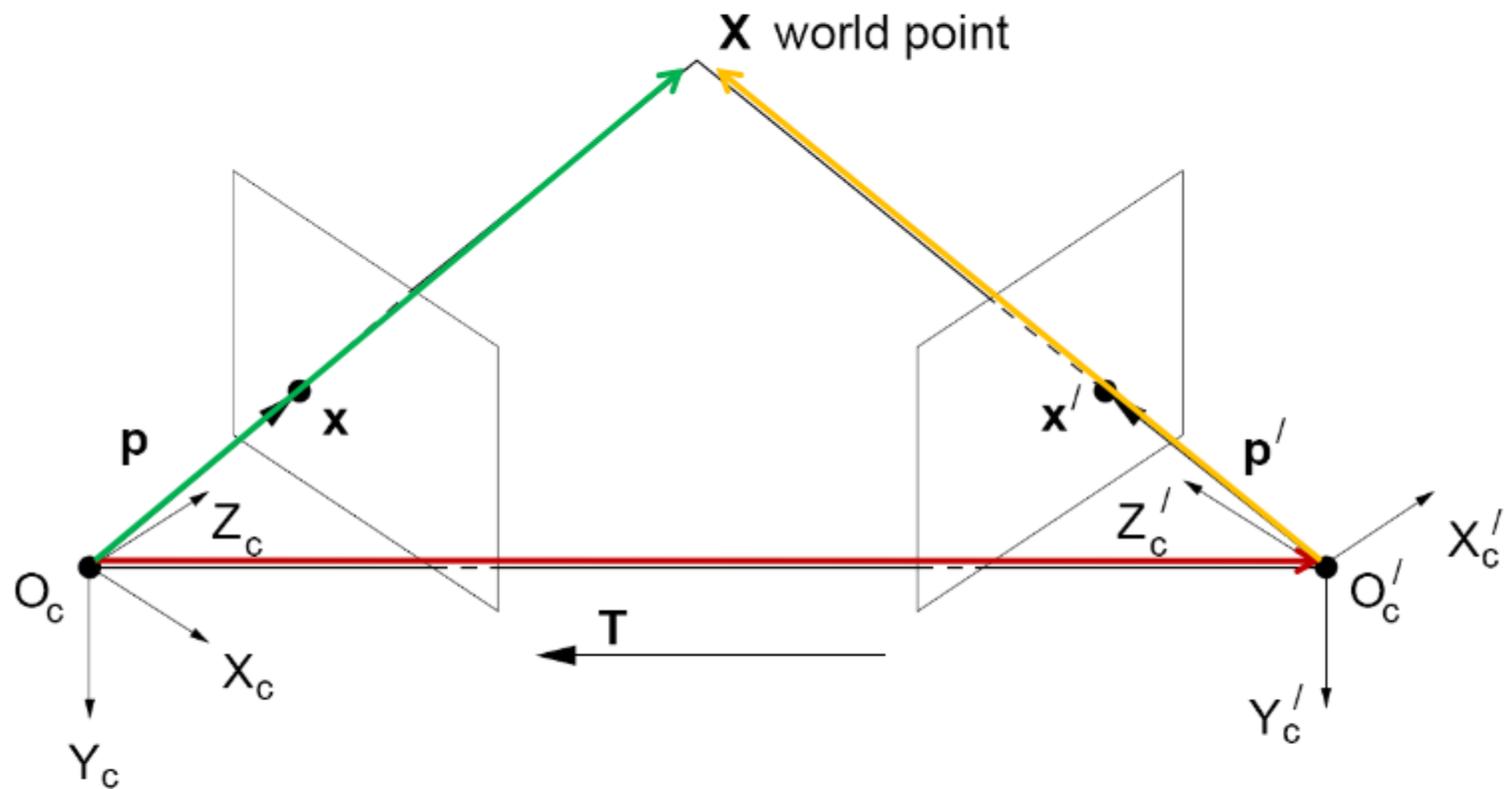
$$\vec{a} \cdot \vec{c} = 0$$

$$\vec{b} \cdot \vec{c} = 0$$

Vector cross product takes two vectors and returns a third vector that's perpendicular to both inputs.

So here,  $c$  is perpendicular to both  $a$  and  $b$ , which means the dot product = 0.

# Deriving the Essential Matrix: From geometry to algebra



$$\mathbf{X}' = \mathbf{R}\mathbf{X} + \mathbf{T}$$

$$\underbrace{\mathbf{T} \times \mathbf{X}'}_{\text{Normal to the plane}} = \mathbf{T} \times \mathbf{R}\mathbf{X}$$

$$\mathbf{X}' \cdot (\mathbf{T} \times \mathbf{X}') = \mathbf{X}' \cdot (\mathbf{T} \times \mathbf{R}\mathbf{X}) = 0$$

# Representing Cross Product in a Matrix Form

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \text{ and } \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

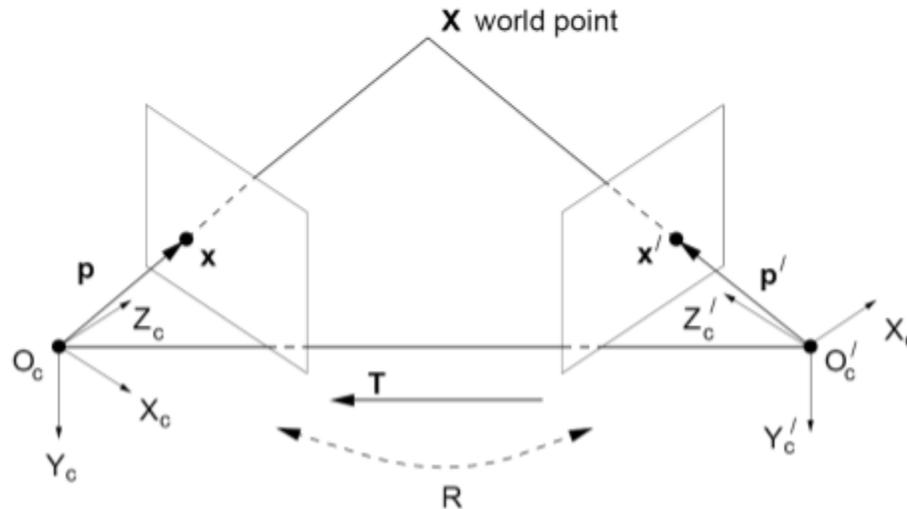
$$\mathbf{a} \times \mathbf{b} = \begin{bmatrix} a_2 b_3 - a_3 b_2 \\ a_3 b_1 - a_1 b_3 \\ a_1 b_2 - a_2 b_1 \end{bmatrix} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$[\mathbf{a}]_{\times}$ : skew matrix

# Essential matrix

$$\mathbf{X}' \cdot (\mathbf{T}_x \mathbf{R}\mathbf{X}) = 0$$

Let  $\mathbf{E} = \mathbf{T}_x \mathbf{R}$



This holds for the rays  $\mathbf{p}$  and  $\mathbf{p}'$  that are parallel to the camera-centered position vectors  $\mathbf{X}$  and  $\mathbf{X}'$ , so we have:

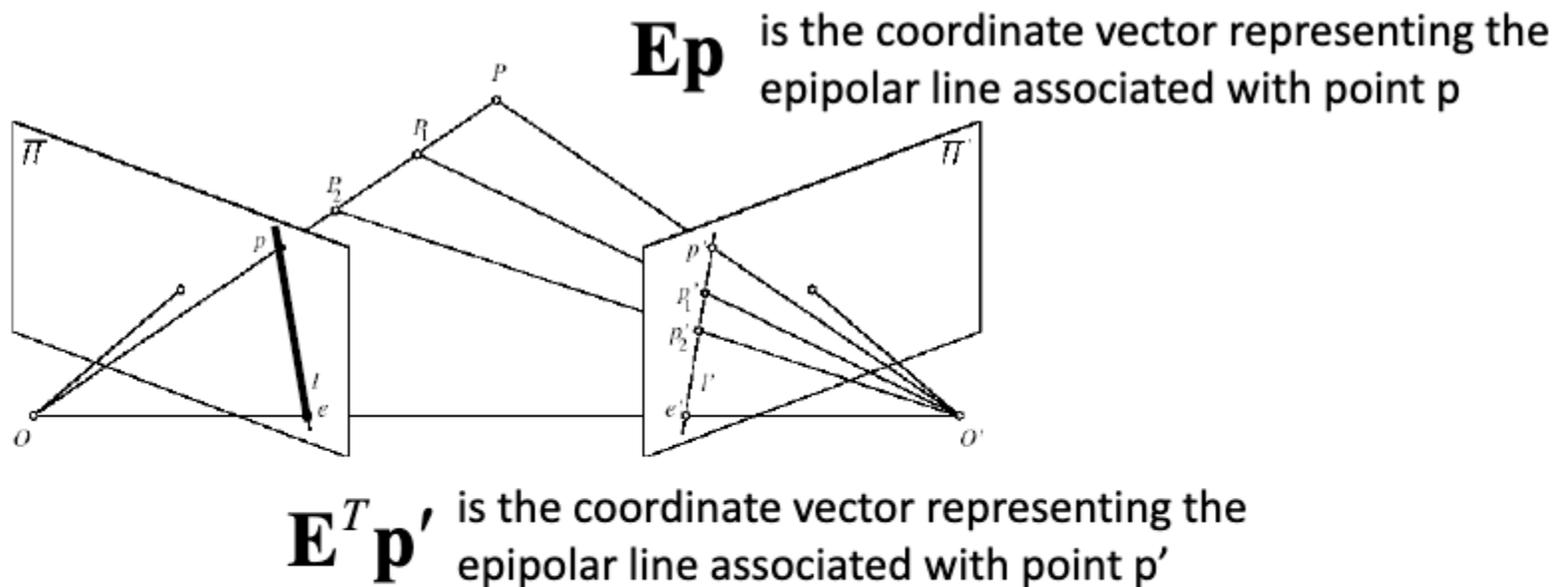
$$\mathbf{p}'^T \mathbf{E} \mathbf{p} = 0$$

$\mathbf{E}$  is called the **essential matrix**, which relates corresponding image points [Longuet-Higgins 1981]

# Essential matrix and epipolar lines

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

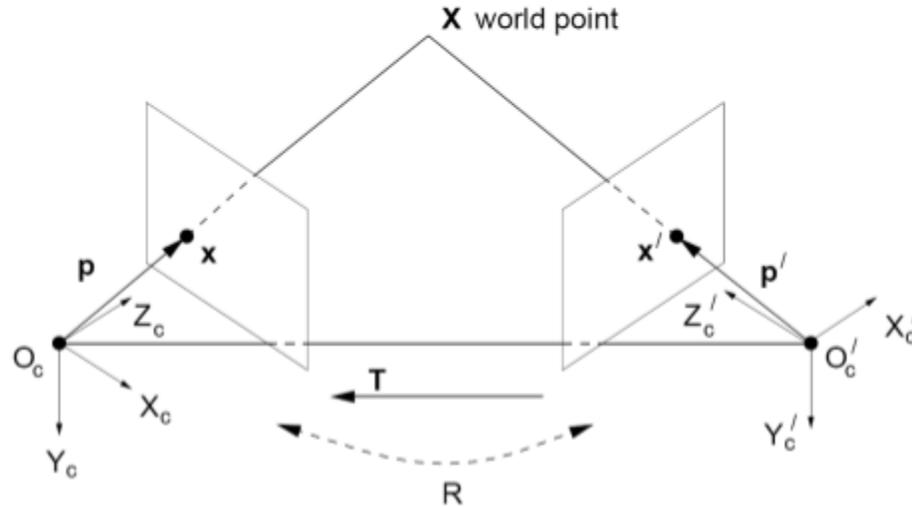
Epipolar constraint: if we observe point  $\mathbf{p}$  in one image, then its position  $\mathbf{p}'$  in second image must satisfy this equation.



# Uncalibrated cameras

$$\mathbf{E} = \mathbf{T}_x \mathbf{R}$$

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



- For an *uncalibrated* stereo rig, can we express the epipolar constraints algebraically via the **Essential Matrix**?
- No, we do not know  $T$  or  $R$
- However we can use the **Fundamental Matrix**
  - Estimate epipolar geometry from a (redundant) set of point correspondences between two uncalibrated cameras

# Uncalibrated case

For a given camera:

$$\bar{\mathbf{p}} = \mathbf{M}_{\text{int}} \mathbf{p}$$

← Camera coordinates

So, for two cameras (left and right):

$$\begin{aligned} \mathbf{p}_{(left)} &= \mathbf{M}_{left,int}^{-1} \bar{\mathbf{p}}_{(left)} \\ \mathbf{p}_{(right)} &= \mathbf{M}_{right,int}^{-1} \bar{\mathbf{p}}_{(right)} \end{aligned}$$

← Camera coordinates

← Image pixel coordinates

Internal calibration matrices, one per camera

$$\mathbf{p}_{(left)} = \mathbf{M}_{left,int}^{-1} \bar{\mathbf{p}}_{(left)}$$

$$\mathbf{p}_{(right)} = \mathbf{M}_{right,int}^{-1} \bar{\mathbf{p}}_{(right)}$$

## Uncalibrated case: **Fundamental matrix**

$$\mathbf{p}_{(right)}^T \mathbf{E} \mathbf{p}_{(left)} = 0$$

From before, the  
**essential matrix E.**

$$\left( \mathbf{M}_{right,int}^{-1} \bar{\mathbf{p}}_{right} \right)^T \mathbf{E} \left( \mathbf{M}_{left,int}^{-1} \bar{\mathbf{p}}_{left} \right) = 0$$

$$\bar{\mathbf{p}}_{right}^T \left( \mathbf{M}_{right,int}^{-T} \mathbf{E} \mathbf{M}_{left,int}^{-1} \right) \bar{\mathbf{p}}_{left} = 0$$

$$\bar{\mathbf{p}}_{right}^T \mathbf{F} \bar{\mathbf{p}}_{left} = 0$$

Fundamental matrix

Grauman

# Fundamental matrix

- Relates pixel coordinates in the two views
- More general form than essential matrix: we remove need to know intrinsic parameters
- If we estimate fundamental matrix from correspondences in pixel coordinates, can reconstruct epipolar geometry without intrinsic or extrinsic parameters

# If Intrinsics are Unknown, How to Get Fundamental Matrix? Correspondence!

Each point  
correspondence  
generates one  
constraint on  $F$

$$\bar{\mathbf{p}}_{right}^T \mathbf{F} \bar{\mathbf{p}}_{left} = 0$$

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

- $(u', v')$  and  $(u, v)$  are pixel coordinates of a pair of corresponding pixel from two views
- $\mathbf{F}$  has 9 variables but we can impose a constraint on the scale of  $\mathbf{F}$ , ending up with 8 degree of freedom.
- **The Eight-point algorithm:** we can solve the system with at least 8 pairs of corresponding pixels.

# If Intrinsic are Unknown, How to Get Fundamental Matrix? Correspondence!

Each point  
correspondence  
generates one  
constraint on F

$$\bar{\mathbf{p}}_{right}^T \mathbf{F} \bar{\mathbf{p}}_{left} = 0$$

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

Collect n of these  
constraints

$$\begin{bmatrix} u'_1 u_1 & u'_1 v_1 & u'_1 & v'_1 u_1 & v'_1 v_1 & v'_1 & u_1 & v_1 & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = \mathbf{0}$$

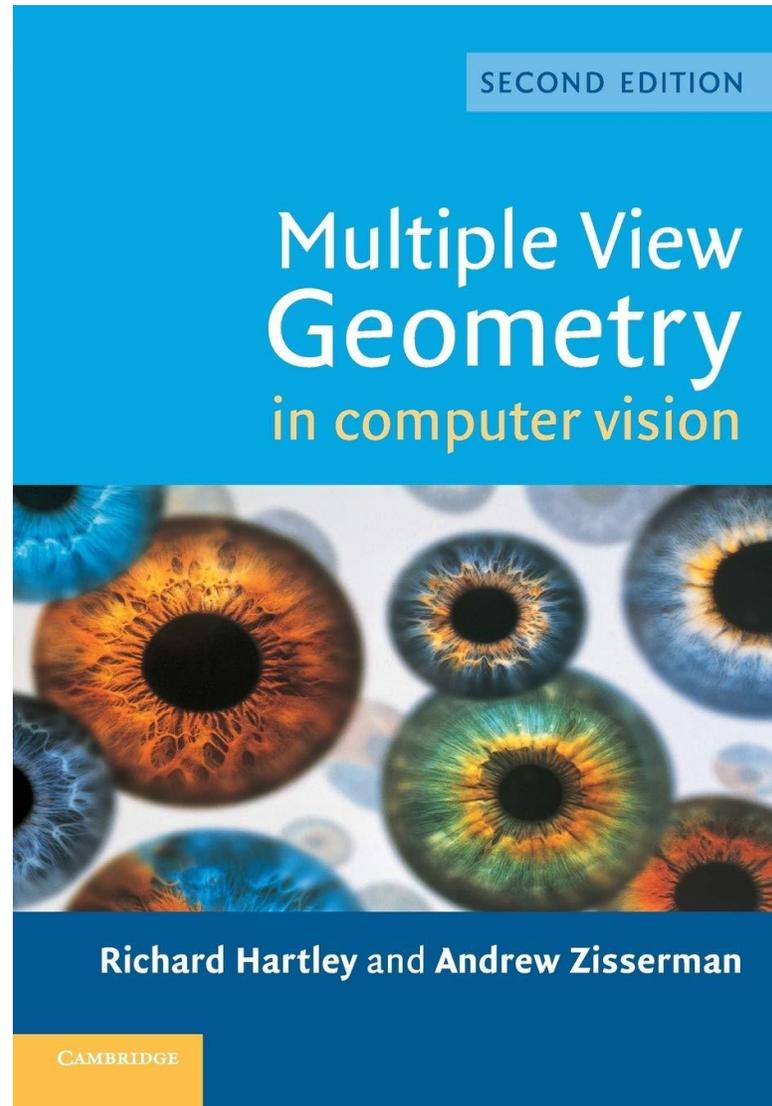
Solve for f , vector of parameters.

Grauman

# Why Do We Need to Know Multi-view Geometry?



# Learn More About Multiview Geometry

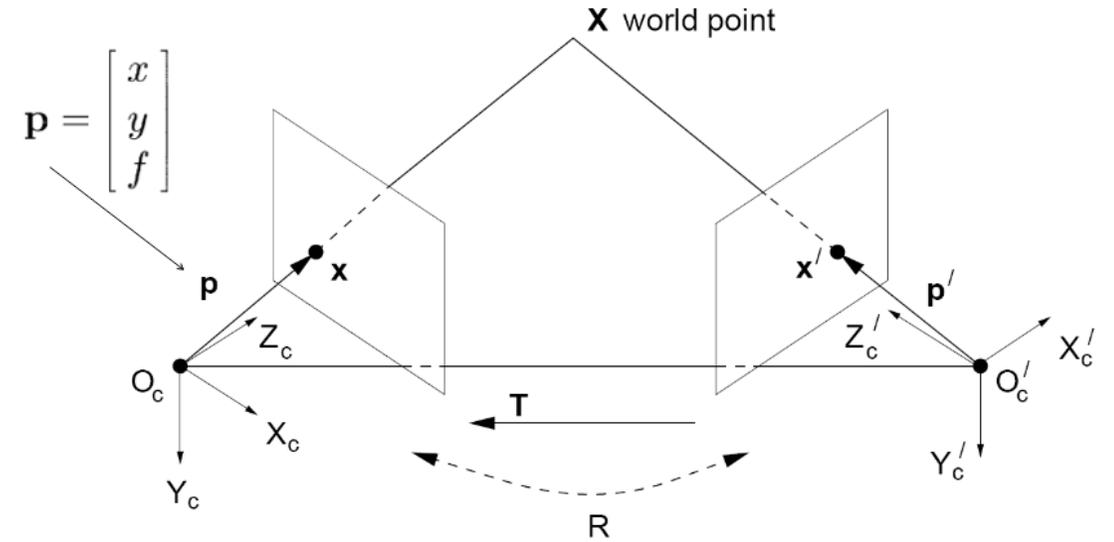
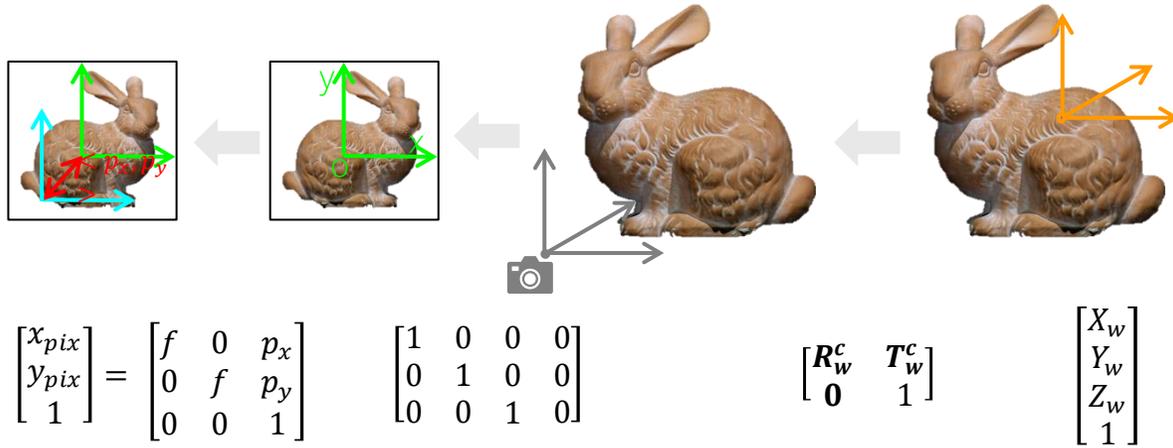


# Learn More About Multiview Geometry

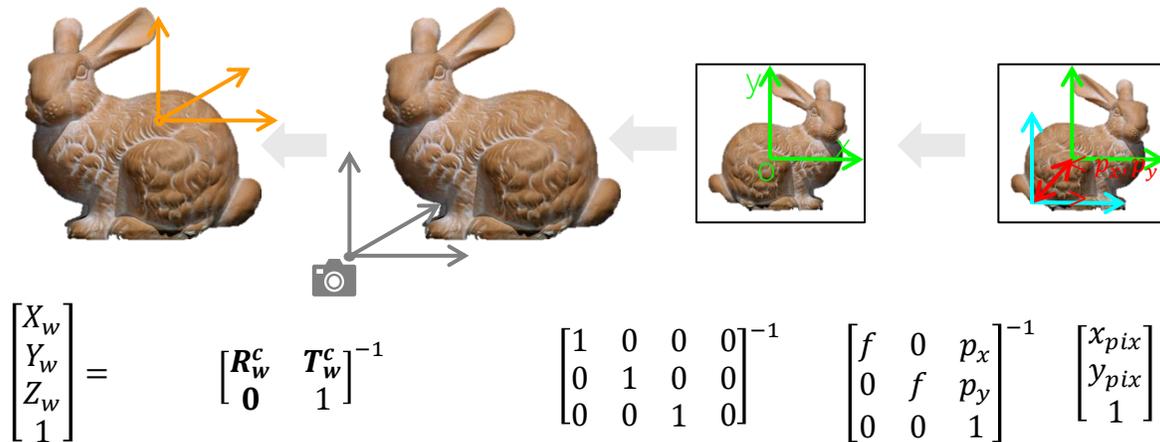
- A more advanced lecture in geometry-based vision
  - [Geometry-based Methods in Vision](#) by Shubham Tulsiani at CMU

# Summary

## 3D-to-2D perspective projection



## 2D-to-3D reconstruction



# What We Will Cover Next Week

- Structure from motion and 3D foundation models
- 3D Scene Representations