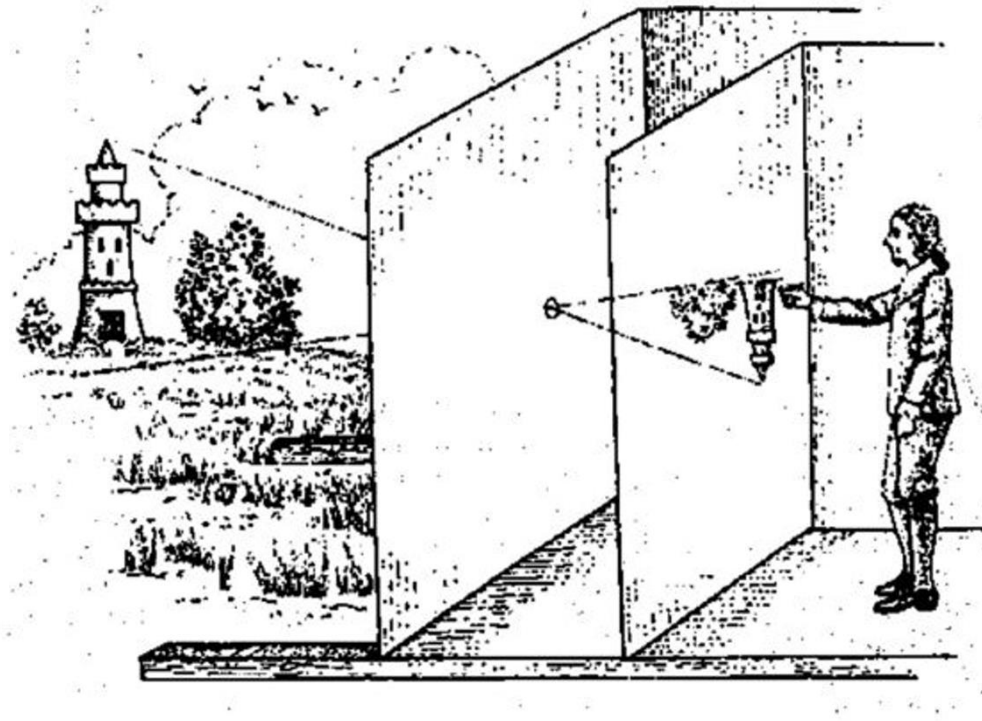# It's a 3D World, After All



## 3D Computer Vision

Elliott Wu

# 3D or Not 3D?

- Provocative blogpost by Vincent Sitzmann:
  https://www.vincentsitzmann.com/blog/bitter_lesson_of_cv/

- Ongoing discussions:
  https://x.com/vincesitzmann/status/2023420051182022774?s=20

VINCENT SITZMANN

Feb 1, 2026

## The flavor of the bitter lesson for computer vision

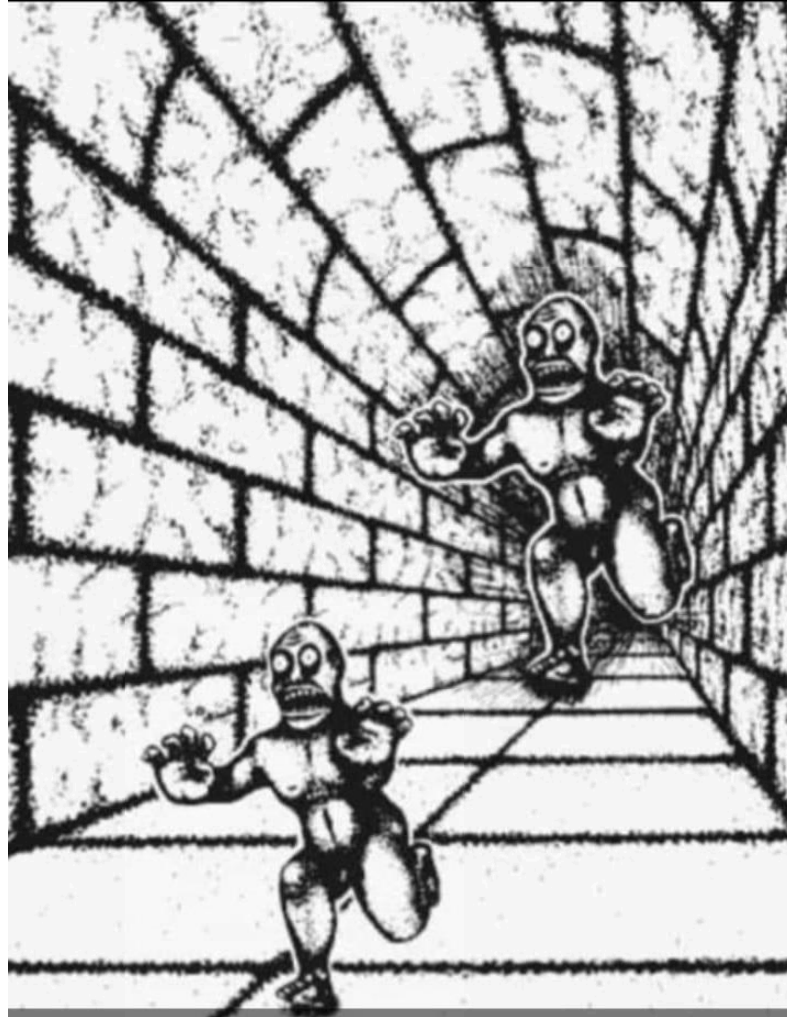I believe that computer vision as we know it is about to go away.

Historically, we have treated vision as a mapping from images to intermediate representations—classes, segmentation masks, or 3D reconstructions. But in the era of the Bitter Lesson, these distinct tasks are becoming qualitatively no different than edge detection: historical artifacts of scoping "solvable intermediate problems" rather than solving intelligence.

While the "LLM moment" in NLP clarified that language modeling is the ultimate objective, the vision community is still debating the flavor of its own revolution. We continue to fine-tune models for specific tasks like point tracking, segmentation, or 3D reconstruction—even as world models emerge, skirt all conventional intermediate representations, and directly solve a problem dramatically more general than everything our community has tackled in the past.
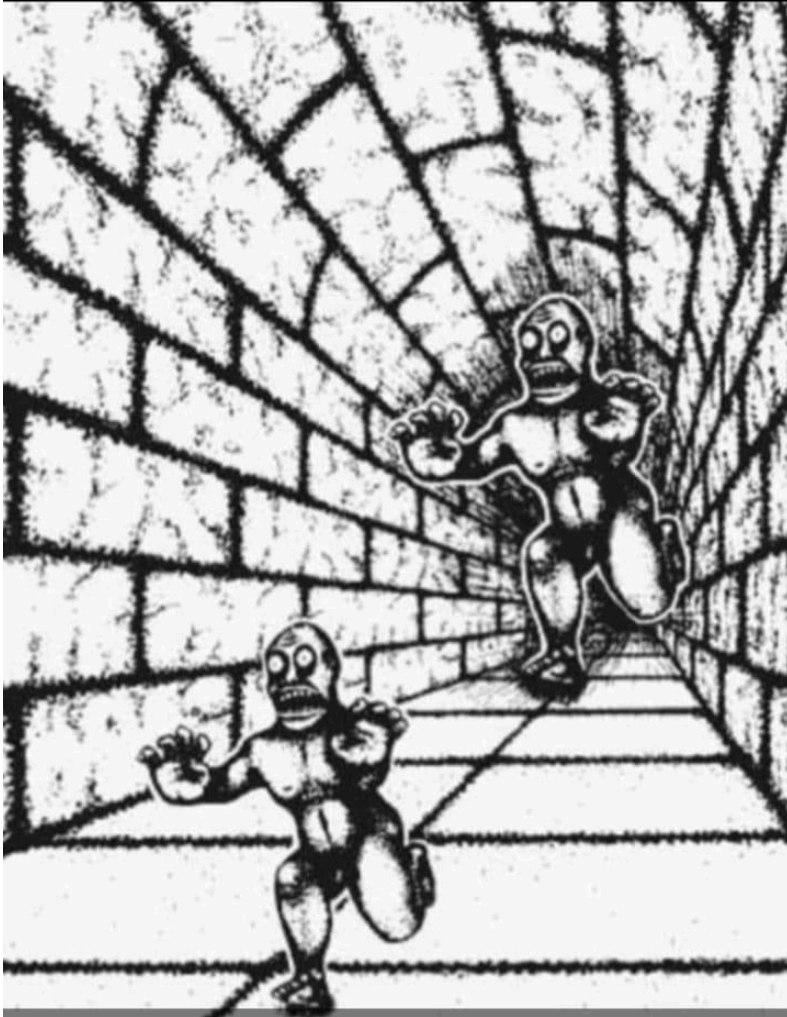
In this post, I argue that the future of computer vision is as part of end-to-end perception-action loops. The historical boundaries between computer vision, robot learning, and control will dissolve. Frontier research will no longer draw a boundary between "seeing" and "learning to act."

As a special case, I will discuss the waning importance of 3D representations: I predict that just as we no longer hand-craft features for detection, we will soon stop using 3D as part of embodied intelligence.

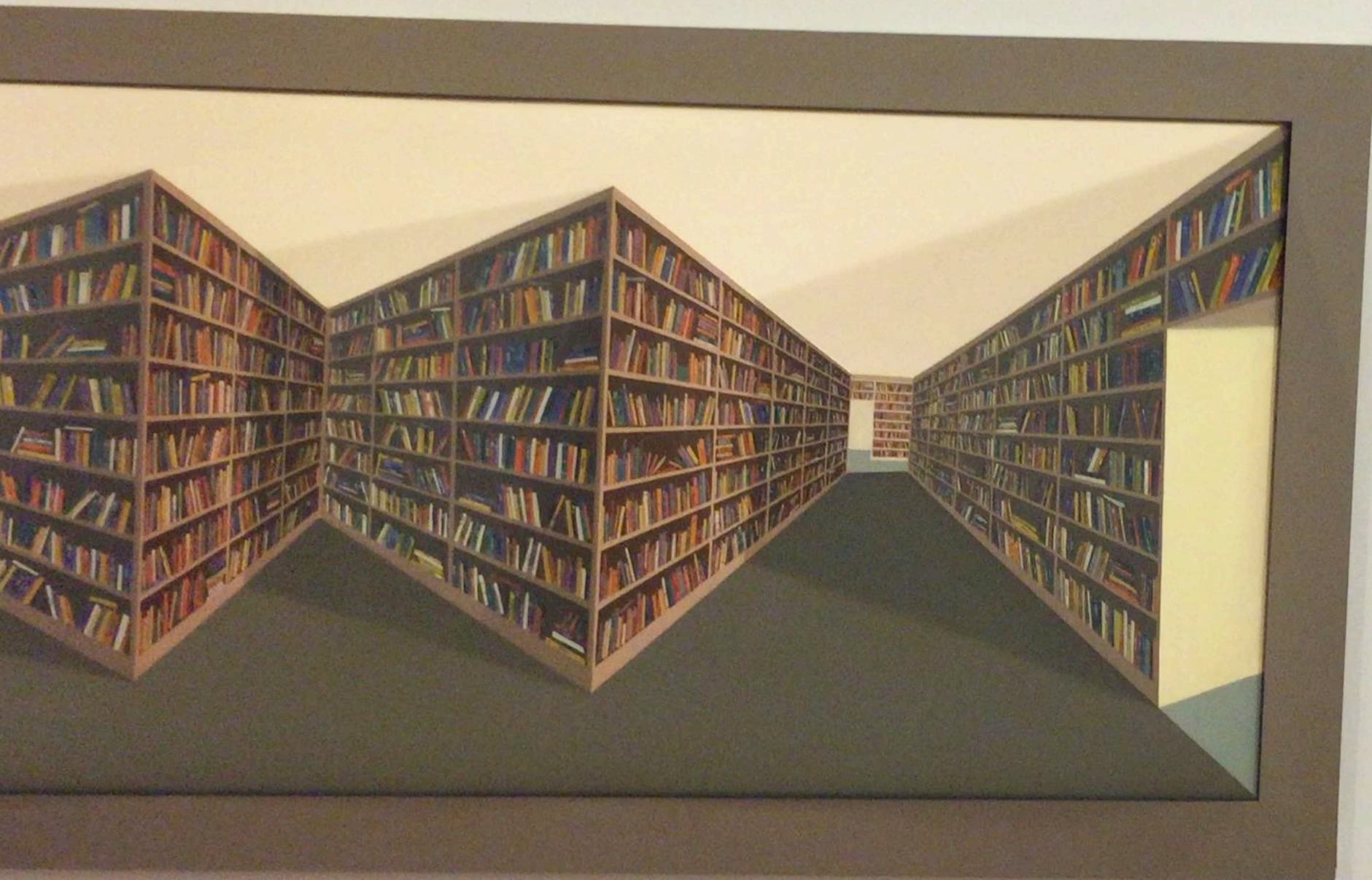# We're good at seeing 3D!

# We're good at seeing 3D!
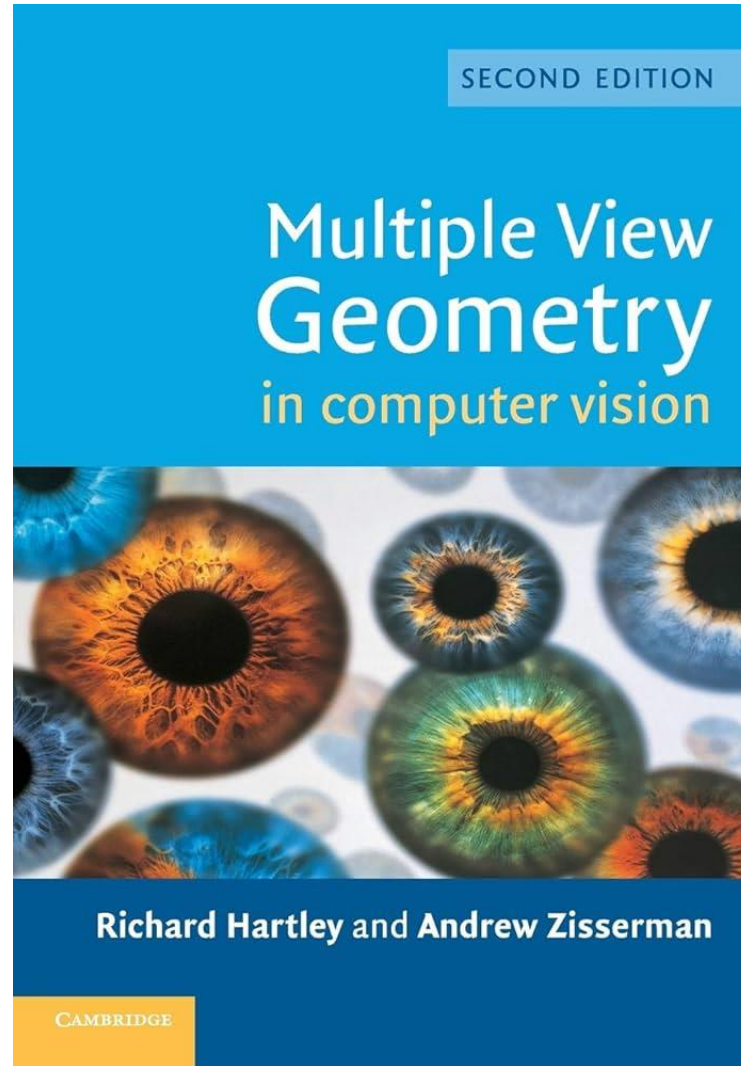
# And sometimes too good..

# And sometimes too good..

These 'Reverspective' paintings are mind-boggling. https://youtu.be/fzZy1FtYsVM

# 3D Computer Vision – Outline

- Part 1 – Camera Model & Projection

- Part 2 – Multi-view Geometry

- Part 3 – 3D Representations & Rendering
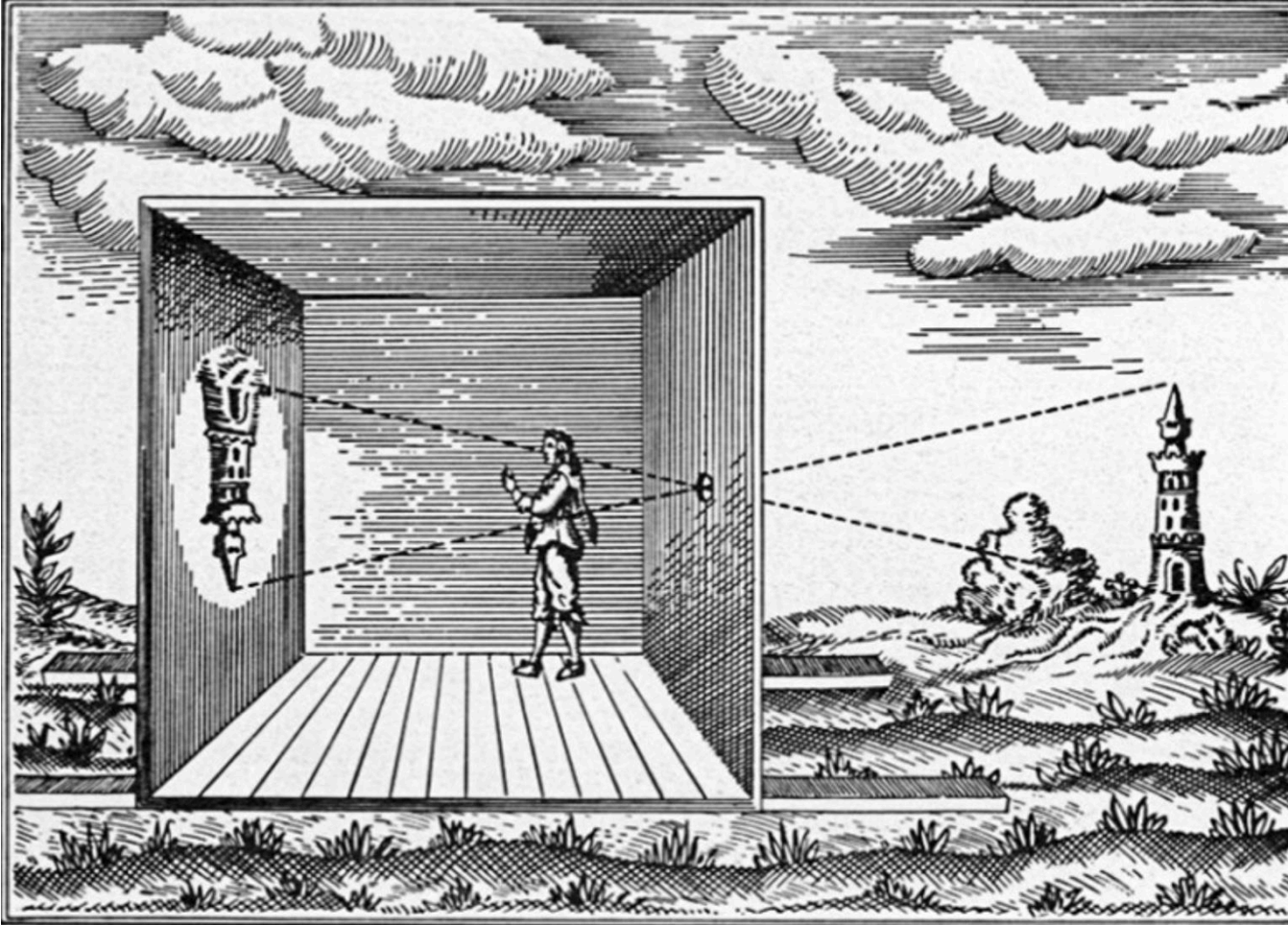
- Part 4 – Learning-based 3D Modeling

# Multi-view Geometry "Bible"

# Part 1 – Camera Model & Projection
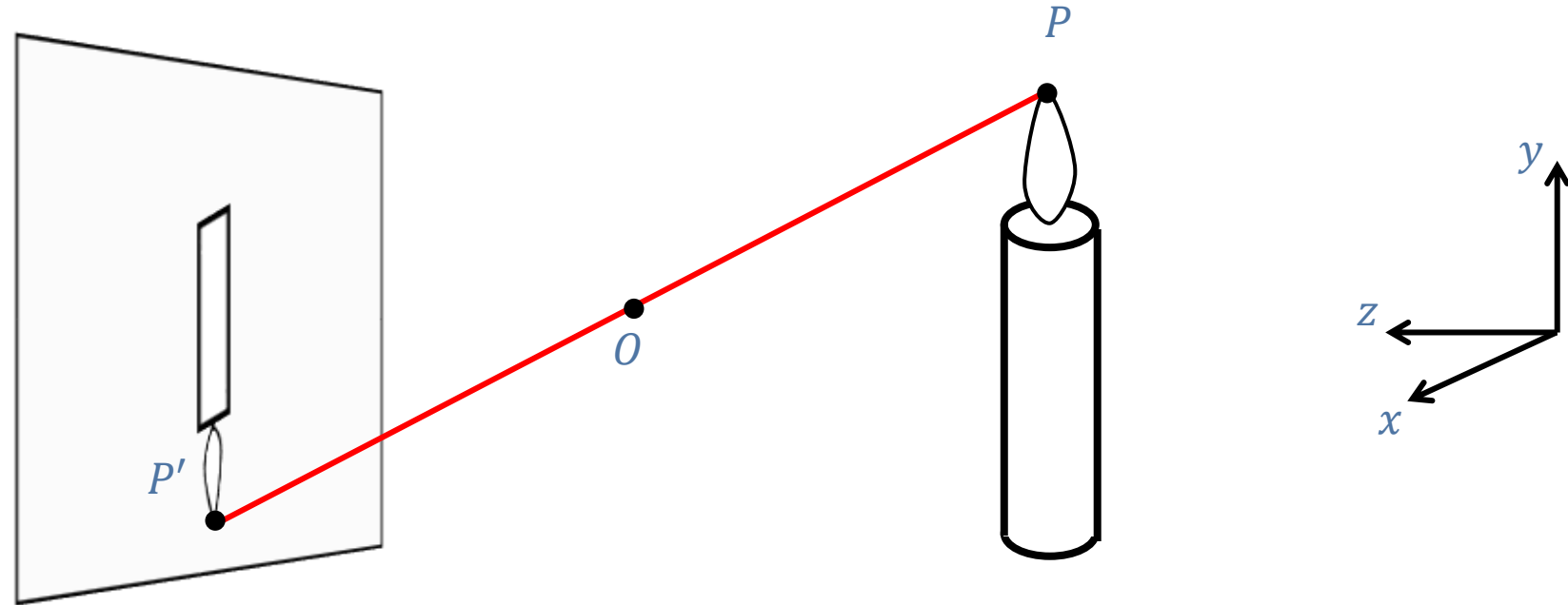
# Image Formation – Camera Obscura

- Basic principle known to Mozi (470-390 BCE), Aristotle (384-322 BCE)

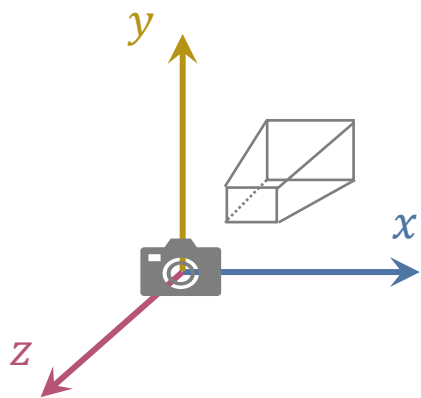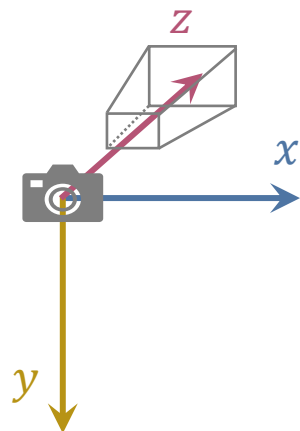- Drawing aid for artists such as Leonardo da Vinci (1452-1519)

# Pinhole Camera



## Canonical coordinate system

- The optical center ($O$) is at the origin
- The z axis is the optical axis perpendicular to the image plane
- The xy plane is parallel to the image plane, x and y axes are horizontal and vertical directions of the image plane

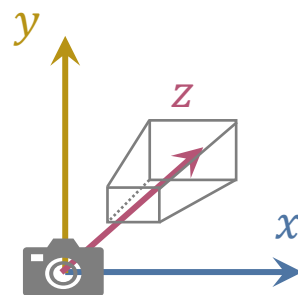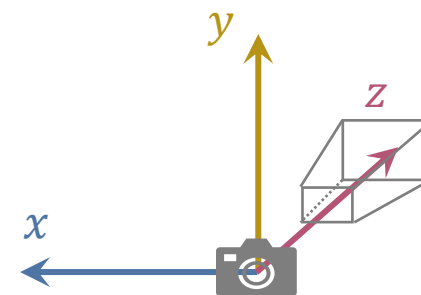# Everybody Agrees… Except on the Axes



- OpenGL
- Blender
- ARKit
- Three.js
- Nerfstudio
- …

- OpenCV
- Open3D
- COLMAP
- gsplat
- …

- Unity
- …

- PyTorch3D
- ?

# Perspective Projection



Focal length

$f$

$y = f \dfrac{Y}{Z}$

$P'$

Image plane

$Y$

$Z$

$O$

$P$

$y$

$z$

$x$

$$\left( f\frac{X}{Z}, f\frac{Y}{Z} \right) \leftarrow (X, Y, Z)$$

2D point        3D point

# Focal Length of Lenses

- In practice, most cameras use lenses, and the focal length is determined by the physical properties of the lens, such as refraction index, thickness, curvature etc.

# Focal Length of Lenses

- In practice, most cameras use lenses, and the focal length is determined by the physical properties of the lens, such as refraction index, thickness, curvature etc.

# Field of View (FOV)

- The field of view is the angular extent of the world observed by the camera.
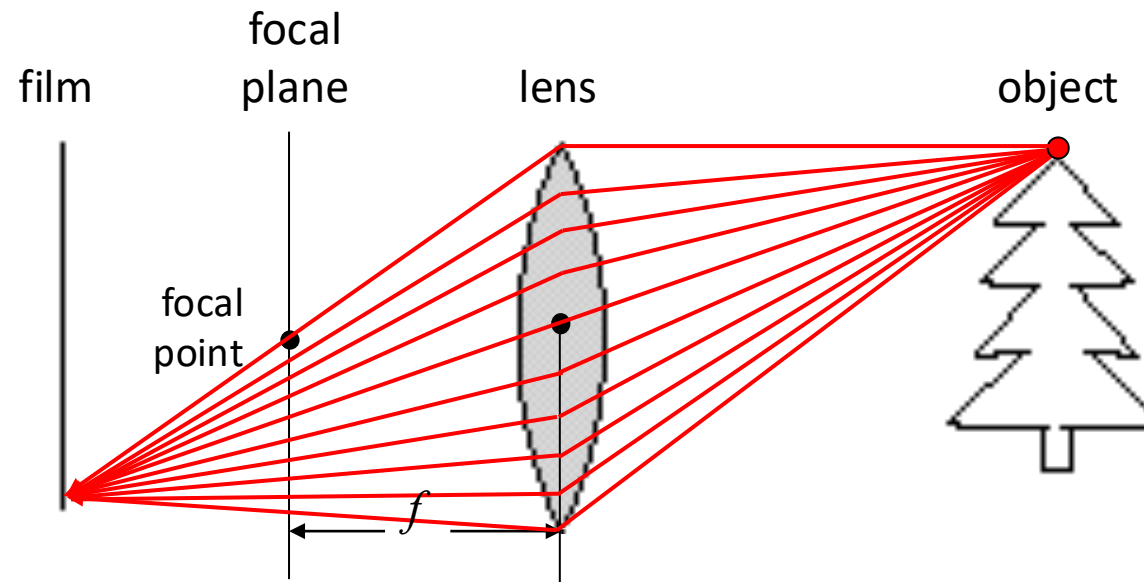
- Focal length ($f$), length of the sensor ($d$):

$$\varphi = \tan^{-1}\frac{d}{2f}$$

- Larger focal length = smaller FOV (assuming a constant sensor size)

# Field of View / Focal Length Ambiguity



Large FOV, small f
Camera close to car

Small FOV, large f
Camera far from the car

# Field of View / Focal Length Ambiguity



wide-angle

standard

telephoto

- What would happen when reconstructing 3D shapes from a single image with *unknown* focal length/FOV?

# Field of View / Focal Length Ambiguity



wide-angle                    standard                    telephoto

(reconstructed by Hunyuan3D V3.1)

- What would happen when reconstructing 3D shapes from a single image with *unknown* focal length/FOV?

# Dolly Zoom Effect

- Continuously adjusting the focal length, while dollying toward or away from the subject.



https://en.wikipedia.org/wiki/Dolly_zoom

# Perspective Projection

- Instead of dealing with an image that is upside down, most of the time we will pretend that the image plane is in front of the camera center.

# Perspective Projection



3D point $(X, Y, Z)$

2D point $(x, y)$

$f\,Y\,/\,Z$

$(X, Y, Z) \rightarrow (x, y)$

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

- Perspective projection is not a linear transformation!
- But most other transformations we will be dealing with are linear.
- Projective geometry provides a handy mathematical tool to unify them.

# Homogeneous Coordinates

- To form homogeneous coordinates from Euclidean coordinates, append 1 as the last entry:

$$(x, y) \Longrightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

homogeneous coordinates
of a 2D point $(x, y)$

$$(X, Y, Z) \Longrightarrow \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

homogeneous coordinates
of a 3D point $(X, Y, Z)$

- To convert homogeneous coordinates to Euclidean coordinates, divide by the last entry:

$$\begin{bmatrix} x \\ y \\ w \end{bmatrix} \Longrightarrow (x/w, y/w)$$

$$\begin{bmatrix} X \\ Y \\ Z \\ W \end{bmatrix} \Longrightarrow (X/W, Y/W, Z/W)$$

- All scalar multiples represent the same point! $\begin{bmatrix} x \\ y \\ w \end{bmatrix} \sim \lambda \begin{bmatrix} x \\ y \\ w \end{bmatrix}, \quad \lambda \neq 0$

# Homogeneous Coordinates

| | **2D** | **3D** | |
|---|---|---|---|
| 2D point $(x, y)$ | $\mathbf{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$ | 3D point $(x, y, z)$ | $\mathbf{X} = \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$ |
| 2D line $\mathbf{l}$ | $\mathbf{l} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$ $\mathbf{l}^\top \mathbf{x} = ax + by + c = 0$ | 3D plane $\boldsymbol{\pi}$ | $\boldsymbol{\pi} = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$ $\boldsymbol{\pi}^\top \mathbf{X} = aX + bY + cZ + d = 0$ |
| 2 points $\mathbf{x}_1, \mathbf{x}_2$ form a line $\mathbf{l}$ | $\mathbf{l} = \mathbf{x}_1 \times \mathbf{x}_2$ | 3 points $\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3$ form a plane $\boldsymbol{\pi}$ | $\boldsymbol{\pi} = \mathbf{x}_1 \times \mathbf{x}_2 \times \mathbf{x}_3$ |
| 2 lines $\mathbf{l}_1, \mathbf{l}_2$ intersect at point $\mathbf{x}$ | $\mathbf{x} = \mathbf{l}_1 \times \mathbf{l}_2$ | 3 planes $\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \boldsymbol{\pi}_3$ intersect at point $\mathbf{X}$ | $\mathbf{X} = \boldsymbol{\pi}_1 \times \boldsymbol{\pi}_2 \times \boldsymbol{\pi}_3$ |

# Perspective Projection Matrix

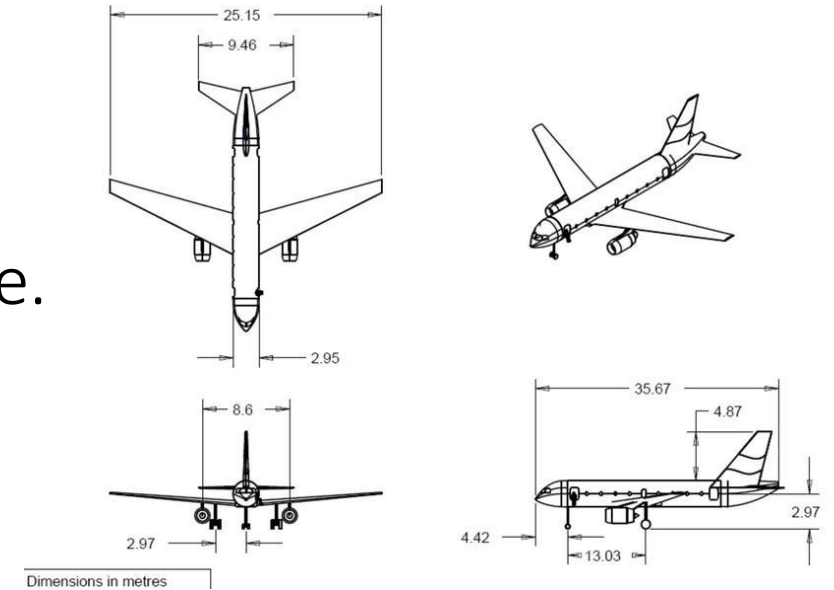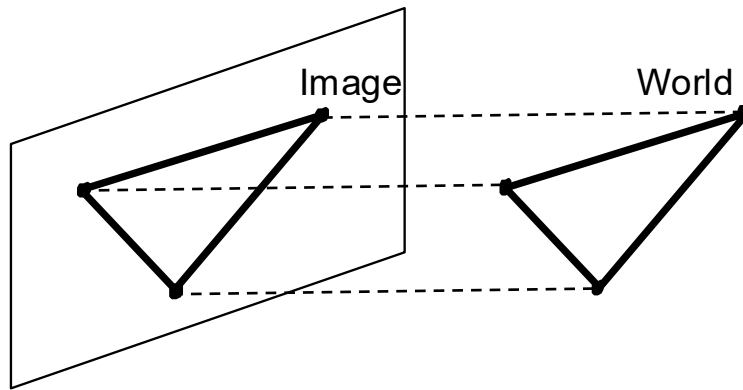- Projection is a matrix multiplication using homogeneous coordinates:

$$\begin{bmatrix} & & & \\ & & & \\ & & & \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = \begin{bmatrix} & \\ & \\ & \end{bmatrix} \Longrightarrow \left( f\frac{X}{Z}, f\frac{Y}{Z} \right)$$

divide by the third coordinate

- What will happen if $f$ is very large?
  - Weak perspective projection $x \approx f\frac{X}{Z_0}$, where $Z_0$ is a constant global depth

# Orthographic Projection

- A form of *parallel projection* where the projection axis is orthogonal to image plane.
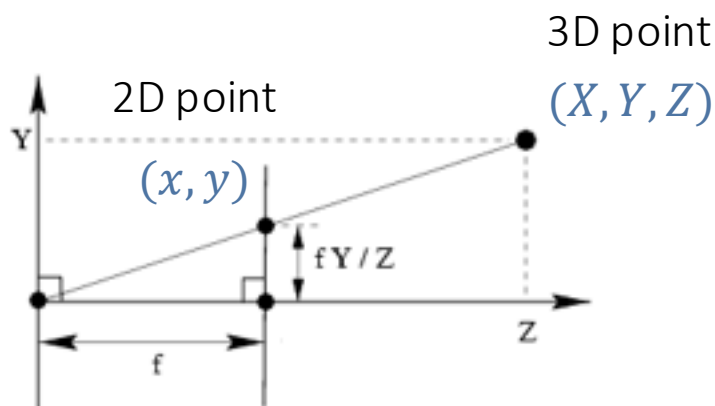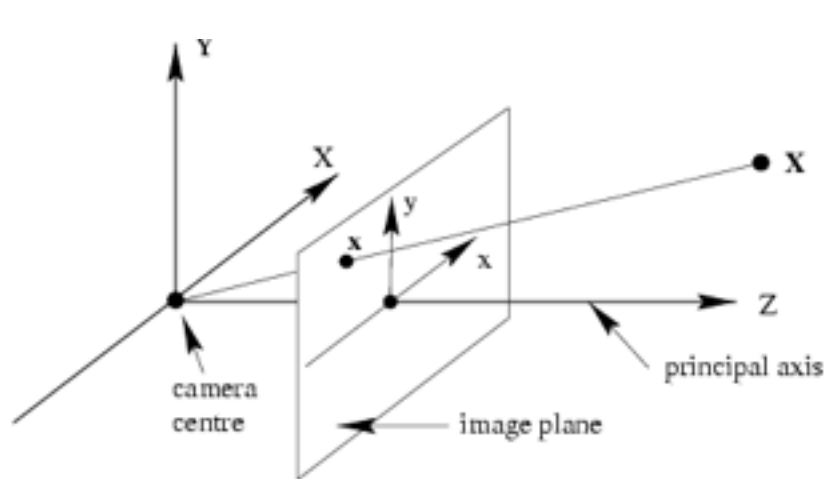


- Assuming projection along the $z$ axis, what's the matrix?

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



Three-view drawing

# Perspective Projection *in Normalized Coordinates*

(assuming camera looking toward z axis and no scaling)



3D point
$(X, Y, Z)$

2D point

$(x, y)$

$f\,Y/Z$

$(X, Y, Z) \to (x, y)$

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

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

$\underbrace{\phantom{x}}_{x}$  $\underbrace{\phantom{P}}_{P}$  $\underbrace{\phantom{X}}_{X}$

$x$

Homogeneous
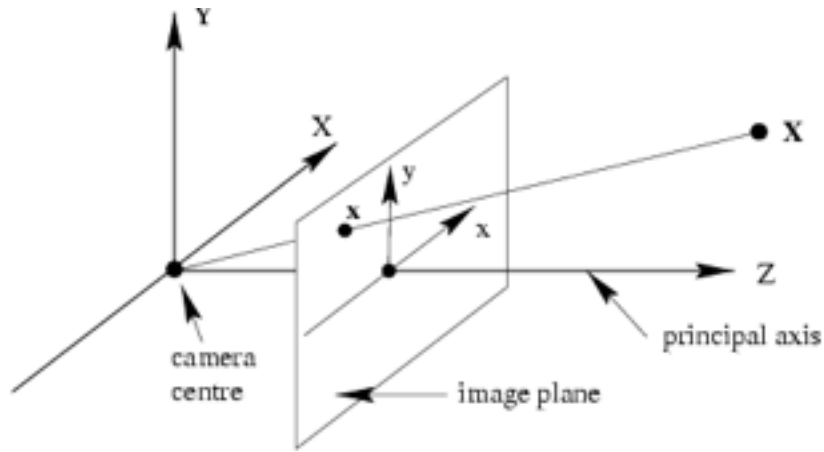coordinates of
image point

$P$

Camera
projection
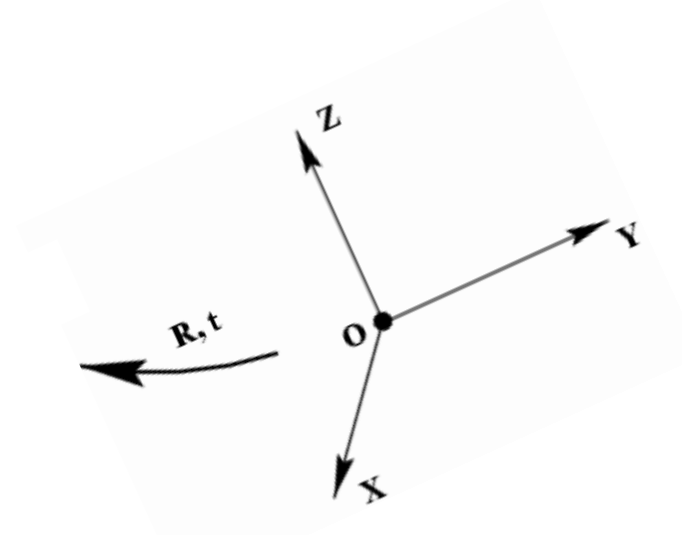matrix

$X$

Homogeneous
coordinates of
3D point

$$\boldsymbol{x} \cong \boldsymbol{PX}$$

Equality up to scale

# Camera Calibration



camera coordinate system

world coordinate system

- **Camera calibration**: figuring out transformation from *world* coordinate system to *image* coordinate system

$$\begin{pmatrix} 2D \\ point \\ \boldsymbol{x} \\ (3x1) \end{pmatrix} \cong \begin{pmatrix} Camera\ to \\ pixel\ coord. \\ trans.\ matrix \\ \boldsymbol{K}\ (3x3) \end{pmatrix} \begin{pmatrix} Canonical \\ projection\ matrix \\ [\boldsymbol{I} \mid \boldsymbol{0}]\ (3x4) \end{pmatrix} \begin{pmatrix} World\ to \\ camera\ coord. \\ trans.\ matrix \\ \begin{bmatrix} \boldsymbol{R} & \boldsymbol{t} \\ \boldsymbol{0}^T & 1 \end{bmatrix}\ (4x4) \end{pmatrix} \begin{pmatrix} 3D \\ point \\ \boldsymbol{X} \\ (4x1) \end{pmatrix}$$

*Intrinsic* **camera parameters**: principal point, scaling factors

*Extrinsic* **camera parameters**: rotation, translation

# Camera Calibration



camera coordinate system                    world coordinate system

- **Camera calibration**: figuring out transformation from *world* coordinate system to *image* coordinate system
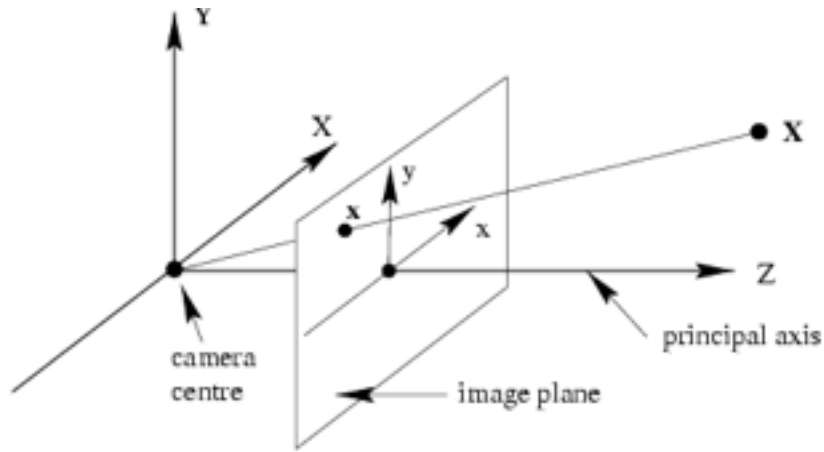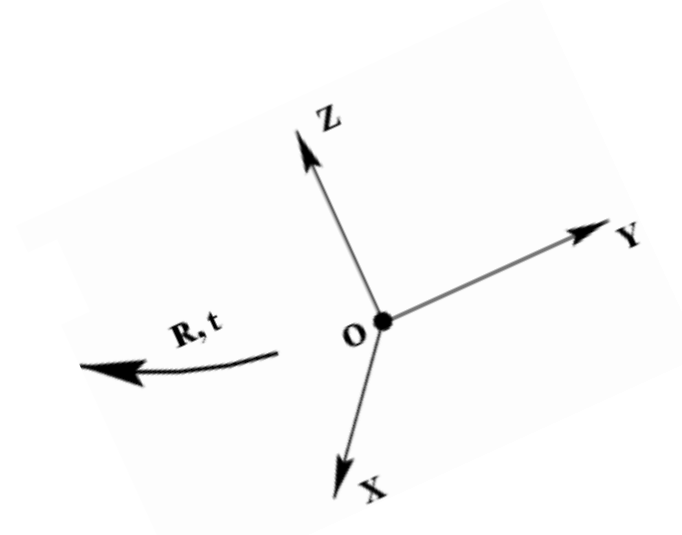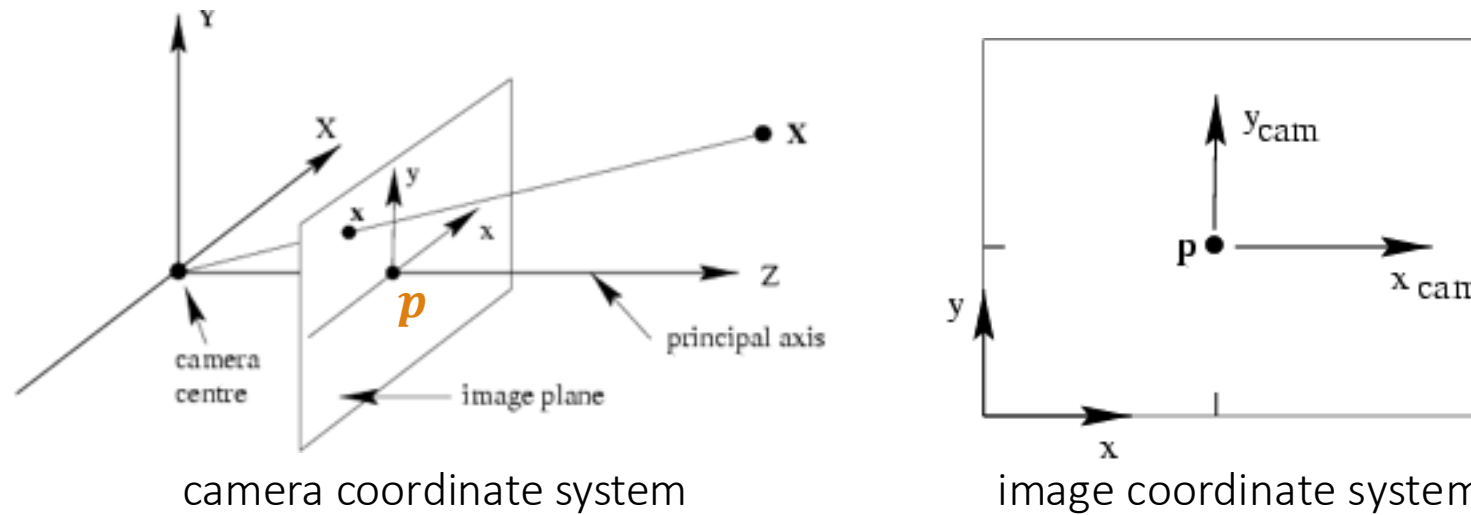
$$
\begin{bmatrix} \text{2D} \\ \text{point} \\ \boldsymbol{x} \\ \text{(3x1)} \end{bmatrix}
\cong
\begin{bmatrix} \text{Camera to} \\ \text{pixel coord.} \\ \text{trans. matrix} \\ \boldsymbol{K} \text{ (3x3)} \end{bmatrix}
\begin{bmatrix} \textit{Canonical} \\ \textit{projection matrix} \\ [\boldsymbol{I} \mid \boldsymbol{0}] \text{ (3x4)} \end{bmatrix}
\begin{bmatrix} \text{World to} \\ \text{camera coord.} \\ \text{trans. matrix} \\ \begin{bmatrix} \boldsymbol{R} & \boldsymbol{t} \\ \boldsymbol{0}^T & 1 \end{bmatrix} \text{ (4x4)} \end{bmatrix}
\begin{bmatrix} \text{3D} \\ \text{point} \\ \boldsymbol{X} \\ \text{(4x1)} \end{bmatrix}
$$

$$\boldsymbol{P} = \boldsymbol{K}[\boldsymbol{R}|\boldsymbol{t}]$$

*General* camera projection matrix

# Intrinsic Parameters – Principal Point



camera coordinate system                    image coordinate system
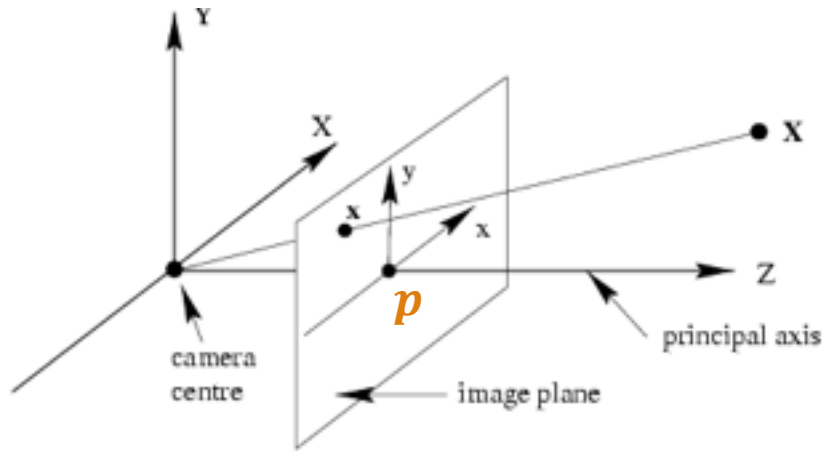
**Principal point ($p$):** point where principal axis intersects the image plane

- In the *normalized* coordinate system, the origin is at the principal point.
- In the *image* coordinate system, the origin is in the corner.

# Intrinsic Parameters – Principal Point
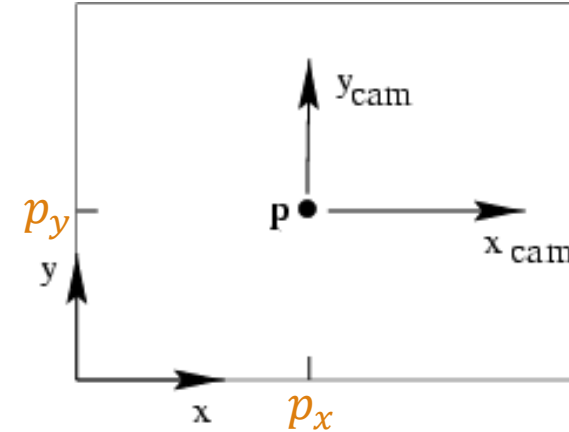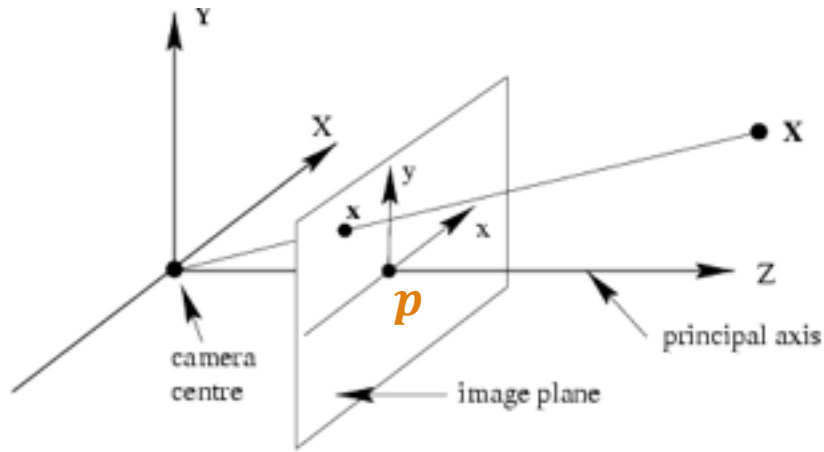


camera coordinate system



image coordinate system

$$x = f\frac{X}{Z} + p_x, \qquad y = f\frac{Y}{Z} + p_y$$

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \cong \begin{pmatrix} fX + Zp_x \\ fY + Zp_y \\ Z \end{pmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}\begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$$

calibration matrix $K$

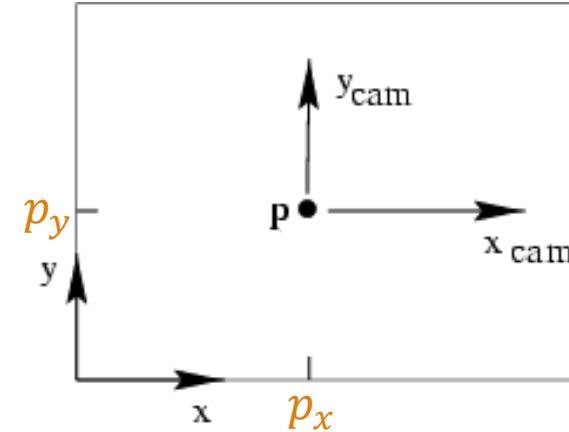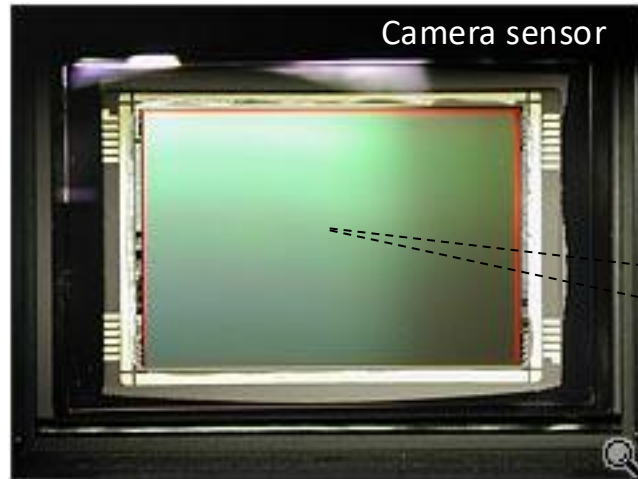# Intrinsic Parameters – Scaling



camera coordinate system



image coordinate system

$$x = f\frac{X}{Z} + p_x, \qquad y = f\frac{Y}{Z} + p_y$$

$$\begin{pmatrix} u \\ v \\ 1 \end{pmatrix} \xleftarrow{\alpha} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \cong \begin{pmatrix} fX + Zp_x \\ fY + Zp_y \\ Z \end{pmatrix} = \begin{bmatrix} f & 0 & p_x & 0 \\ 0 & f & p_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$$

pixels/m

pixels        m        calibration matrix $\boldsymbol{K}$

# Intrinsic Parameters – Scaling



Camera sensor

$m_x$ pixels/m in horizontal direction
$m_y$ pixels/m in vertical direction

Pixel size (m): $\dfrac{1}{m_x} \times \dfrac{1}{m_y}$

Scaling factors

$$\begin{bmatrix} m_x & 0 & 0 \\ 0 & m_y & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

pixels/m

Calibration matrix
$K$ in metric units

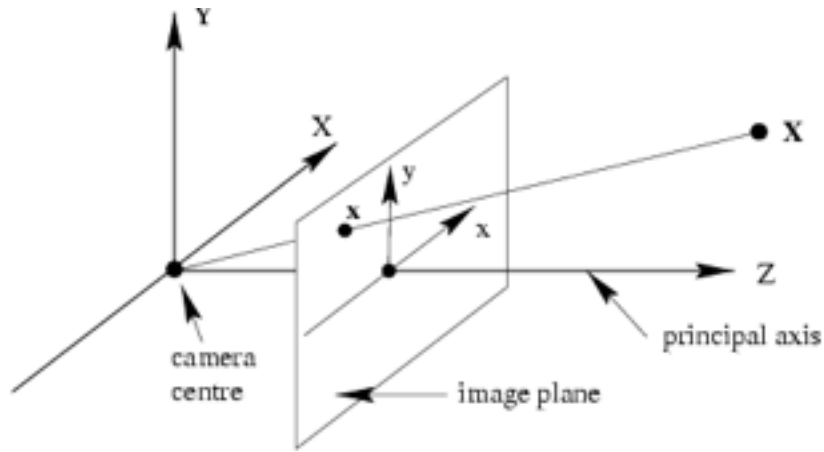$$\begin{bmatrix} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{bmatrix}$$
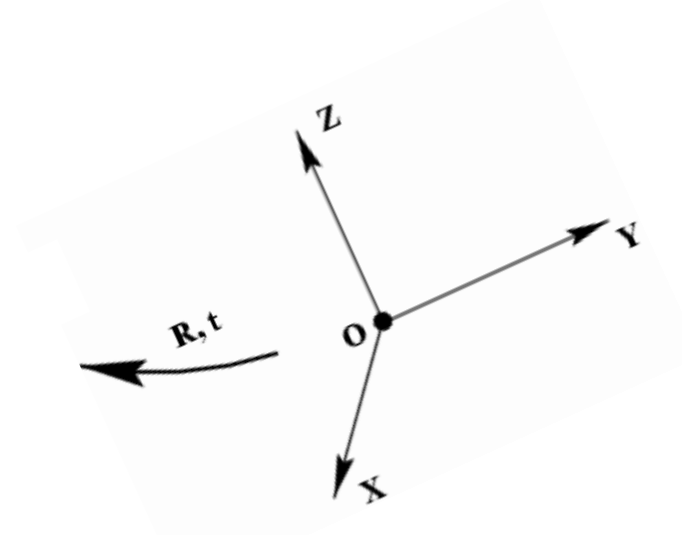
m

Calibration matrix
$K$ in pixel units

$$= \begin{bmatrix} \alpha_x & 0 & \beta_x \\ 0 & \alpha_y & \beta_y \\ 0 & 0 & 1 \end{bmatrix}$$

pixels

# Camera Calibration



camera coordinate system

world coordinate system

- **Camera calibration**: figuring out transformation from *world* coordinate system to *image* coordinate system
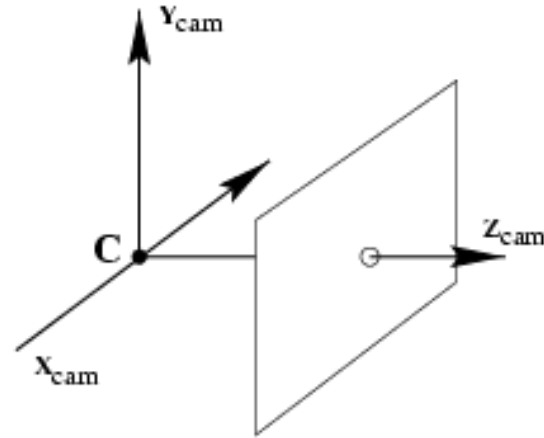
$$\begin{bmatrix} 2D \\ point \\ \boldsymbol{x} \\ (3x1) \end{bmatrix} \cong \begin{bmatrix} Camera\ to \\ pixel\ coord. \\ trans.\ matrix \\ \boldsymbol{K}\ (3x3) \end{bmatrix} \begin{bmatrix} Canonical \\ projection\ matrix \\ [\boldsymbol{I}\ |\ \boldsymbol{0}]\ (3x4) \end{bmatrix} \begin{bmatrix} World\ to \\ camera\ coord. \\ trans.\ matrix \\ \begin{bmatrix} \boldsymbol{R} & \boldsymbol{t} \\ \boldsymbol{0}^T & 1 \end{bmatrix}\ (4x4) \end{bmatrix} \begin{bmatrix} 3D \\ point \\ \boldsymbol{X} \\ (4x1) \end{bmatrix}$$

$$\boldsymbol{P} = \boldsymbol{K}[\boldsymbol{R}|\boldsymbol{t}]$$
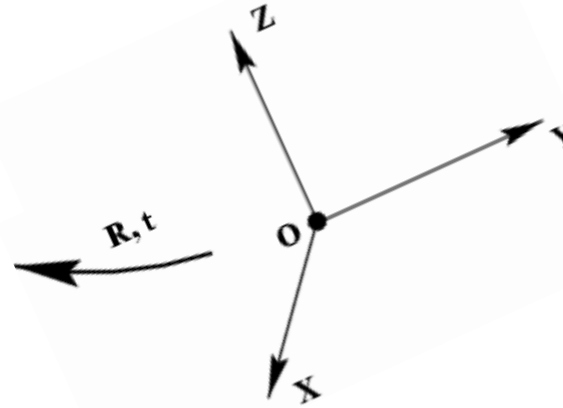
*General* camera projection matrix

# Extrinsic Parameters



camera coordinate system          world coordinate system

- In *non-homogeneous* coordinates, the transformation from world to normalized camera coordinate system is given by:

$$\widetilde{X}_{\text{cam}} = R(\widetilde{X} - \widetilde{C}) = R\widetilde{X} + t$$

coords. of point in normalized camera frame

3x3 rotation matrix

coords. of a point in world frame

coords. of camera center in world frame

W2C: $\widetilde{X}_{\text{cam}} = R\widetilde{X}_{\text{world}} + t$

C2W: $\widetilde{X}_{\text{world}} = R^T \widetilde{X}_{\text{cam}} - R^T t$

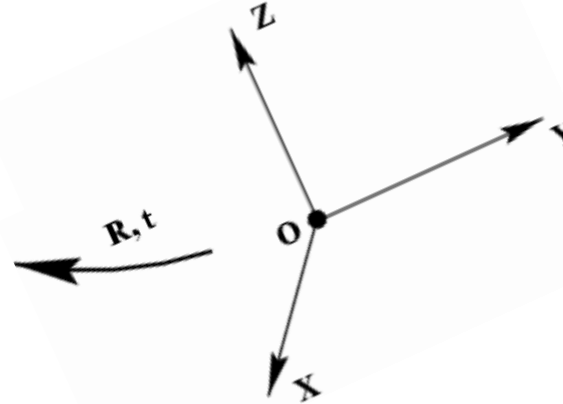# Extrinsic Parameters



camera coordinate system



world coordinate system

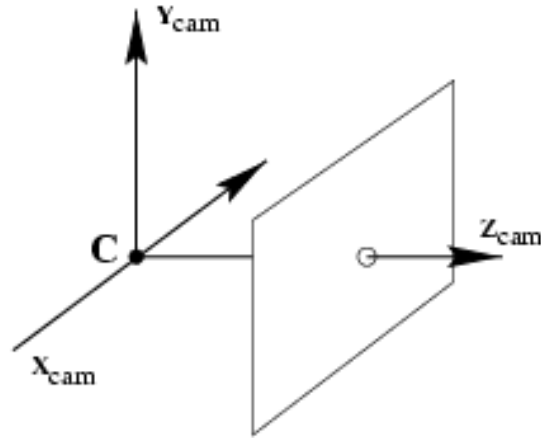In non-homogeneous coordinates:

$$\widetilde{X}_{\text{cam}} = R\widetilde{X} + t$$

In homogeneous coordinates:

$$\begin{pmatrix} \widetilde{X}_{\text{cam}} \\ 1 \end{pmatrix} = \begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix} \begin{pmatrix} \widetilde{X} \\ 1 \end{pmatrix}$$

3D transformation matrix (4 x 4)

# Extrinsic Parameters



camera coordinate system

world coordinate system

In non-homogeneous coordinates:

$$\widetilde{X}_{\text{cam}} = R\widetilde{X} + t$$
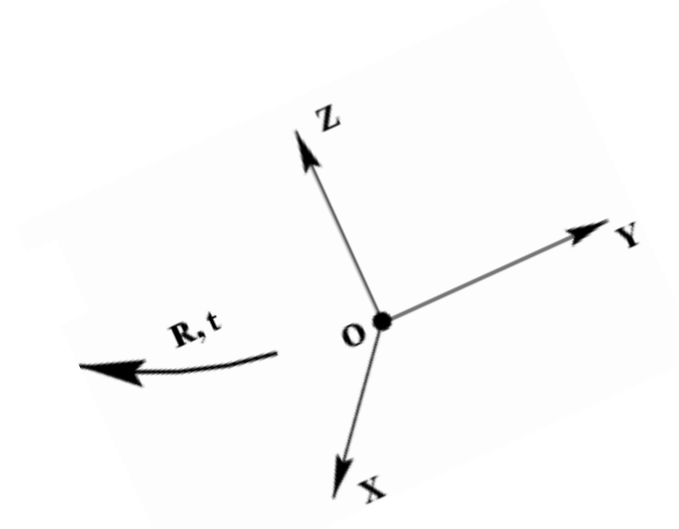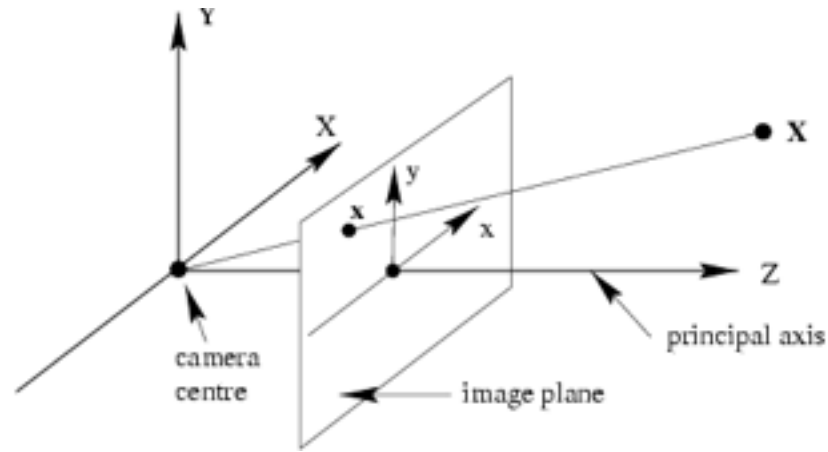
In homogeneous coordinates:

$$X_{\text{cam}} = \begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix} X$$
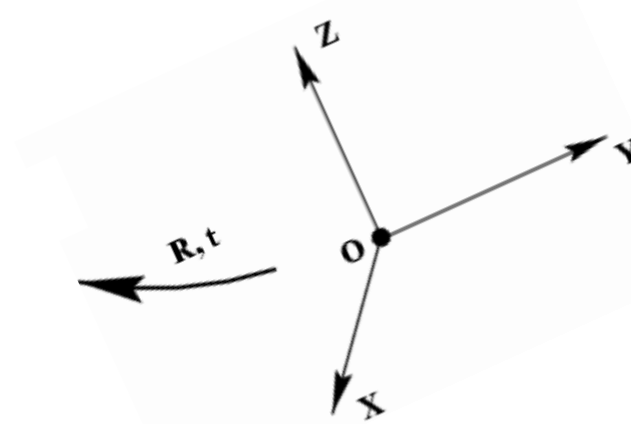
3D transformation matrix (4 x 4)

And then to image coordinates: $x \cong K[I|0]X_{\text{cam}}$

# Camera Calibration



camera coordinate system



world coordinate system

$$x \cong K[R|t]X$$

- **Camera calibration**: figuring out transformation from *world* coordinate system to *image* coordinate system

$$
\begin{bmatrix} \text{2D} \\ \text{point} \\ x \\ \text{(3x1)} \end{bmatrix}
\cong
\begin{bmatrix} \text{Camera to} \\ \text{pixel coord.} \\ \text{trans. matrix} \\ K \text{ (3x3)} \end{bmatrix}
\begin{bmatrix} \textit{Canonical} \\ \text{projection matrix} \\ [I \mid 0] \text{ (3x4)} \end{bmatrix}
\begin{bmatrix} \text{World to} \\ \text{camera coord.} \\ \text{trans. matrix} \\ \begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix} \text{ (4x4)} \end{bmatrix}
\begin{bmatrix} \text{3D} \\ \text{point} \\ X \\ \text{(4x1)} \end{bmatrix}
$$

$$P = K[R|t]$$

*General* camera projection matrix

# Camera Calibration
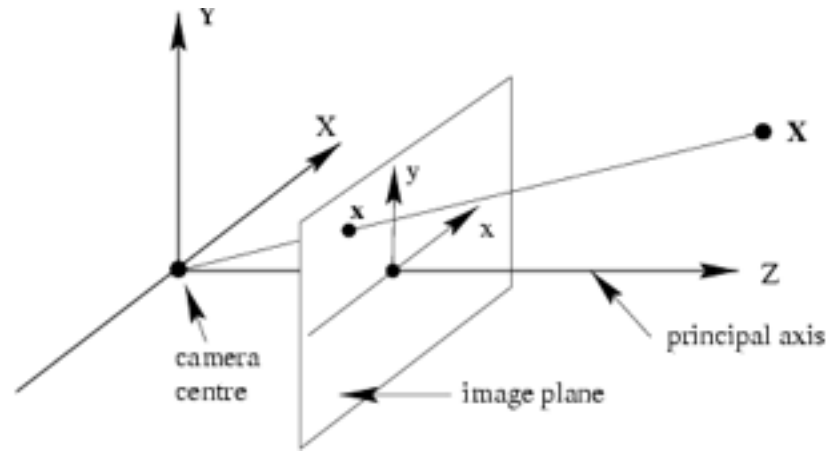


camera coordinate system



world coordinate system

$$x \cong K[R|t]X$$

- **Camera calibration**: figuring out transformation from *world* coordinate system to *image* coordinate system
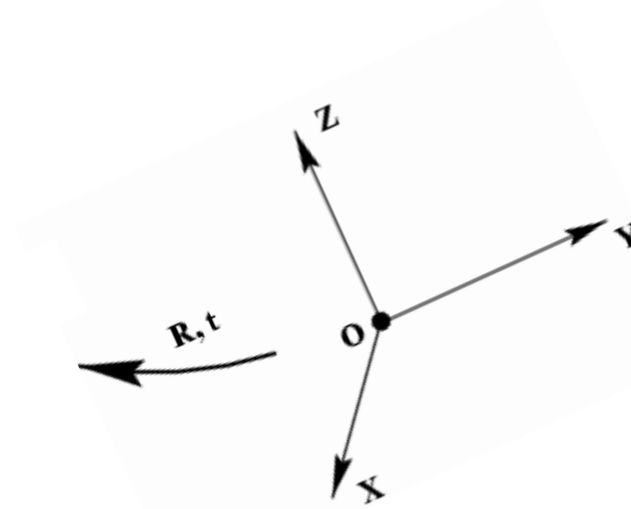
$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \cong \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$$

We could solve this as a linear system (DLT), then perform QR decomposition to get $K$, $R$, $t$. (not robust in practice; sensitive to noise)

# Perspective-n-Point (PnP)

$$x \cong K[R|t]X$$



Camera
**R,t** **?**

**Known:**
- $n$ 3D points $X$
- Corresponding 2D pixel coordinates $x$
- Camera intrinsics $K$

**Solve for:**
- Camera pose (extrinsics) $R$, $t$

**Q:** How many points at least do we need?

**A:** n ≥ 3 (in general). We have 6 DoF unknowns $R$, $t$. Each point provides 2 constraints (2 equations for $x$ and $y$).

# Part 2 – Multi-view Geometry

# Multiple View Geometry



Camera 1
$\mathbf{R_1, t_1}$

Camera 2
$\mathbf{R_2, t_2}$

Camera 3
$\mathbf{R_3, t_3}$

# Multiple View Geometry



Incorrect 3D point

Camera 1
$\mathbf{R_1, t_1}$

Camera 2
$\mathbf{R_2, t_2}$

Camera 3
$\mathbf{R_3, t_3}$

# Multiple View Geometry



Incorrect camera poses

Camera 1
$\mathbf{R_1, t_1}$

Camera 2
$\mathbf{R_2, t_2}$

Camera 3
$\mathbf{R_3, t_3}$

# Multiple View Geometry



**Multi-view constraints** allow us to determine both camera parameters and 3D (up to scale/projective ambiguities)

Camera 1
$\mathbf{R_1, t_1}$ **?**

Camera 2
$\mathbf{R_2, t_2}$ **?**

Camera 3
$\mathbf{R_3, t_3}$ **?**

# Estimating Correspondences



Feature detection

Feature matching

# Estimating Correspondences

# Estimating Correspondences



- **Detect feature points**
  - SIFT descriptors [Lowe, 2004]
  - Learned features – SuperPoint [DeTone, 2018])

# Estimating Correspondences



- **Detect feature points**
  - SIFT descriptors [Lowe, 2004]
  - Learned features – SuperPoint [DeTone, 2018])

- **Match features** between pair of images
  - Use RANSAC to filter outliers
  - Learned matching – SuperGlue [Sarlin, 2020]

# Let's Start with Two Views – Epipolar Geometry

"epi-pole" ≈ upon the pole/pivot



**Epipolar constraints**

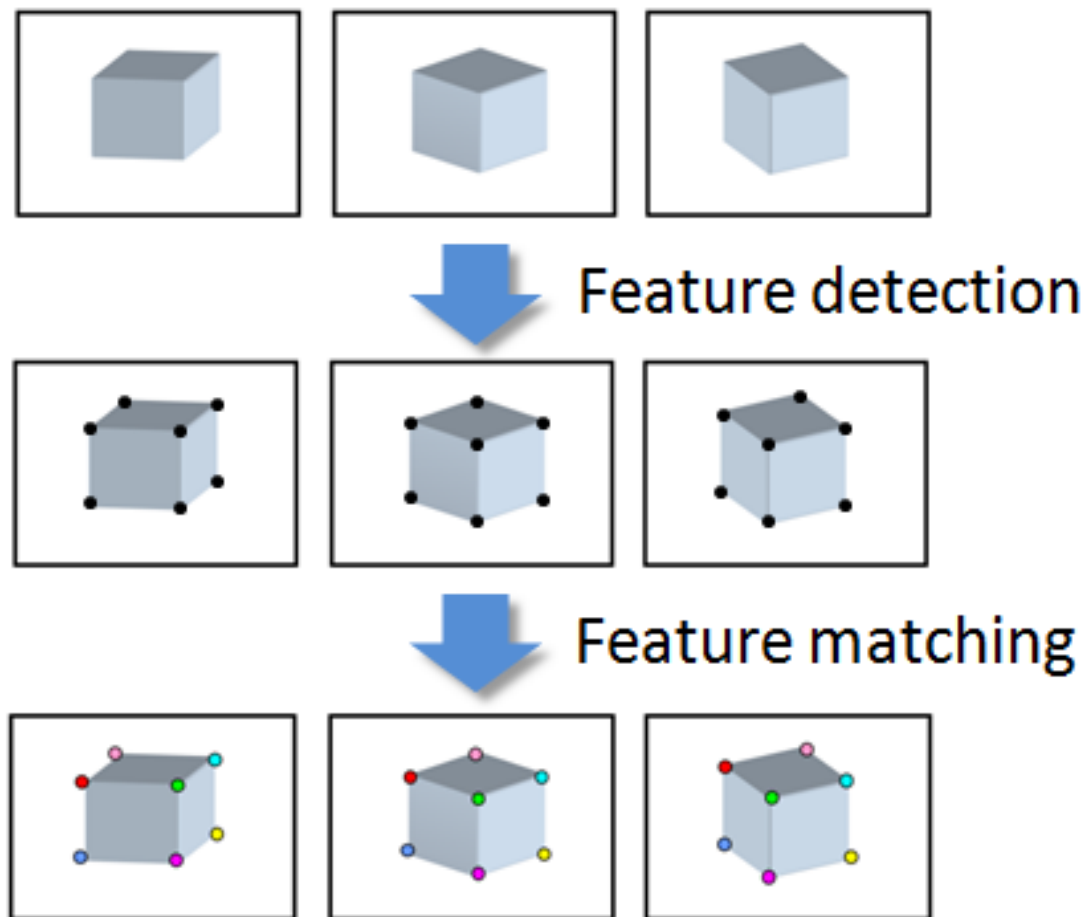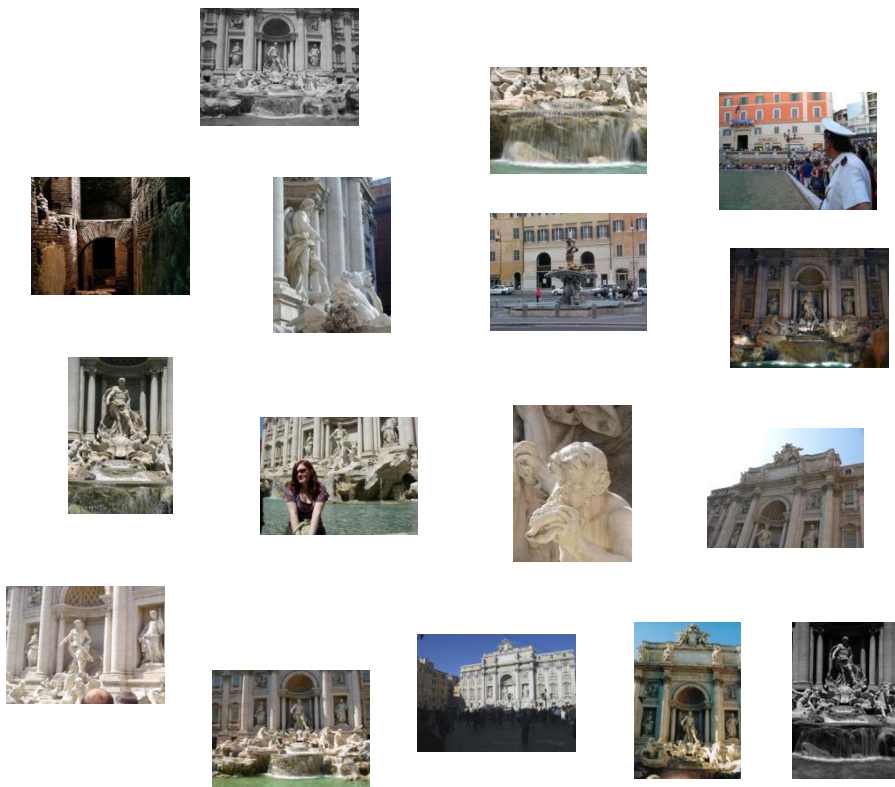~~**Multi-view constraints**~~ allow us to determine both camera parameters and 3D (up to scale/projective ambiguities)
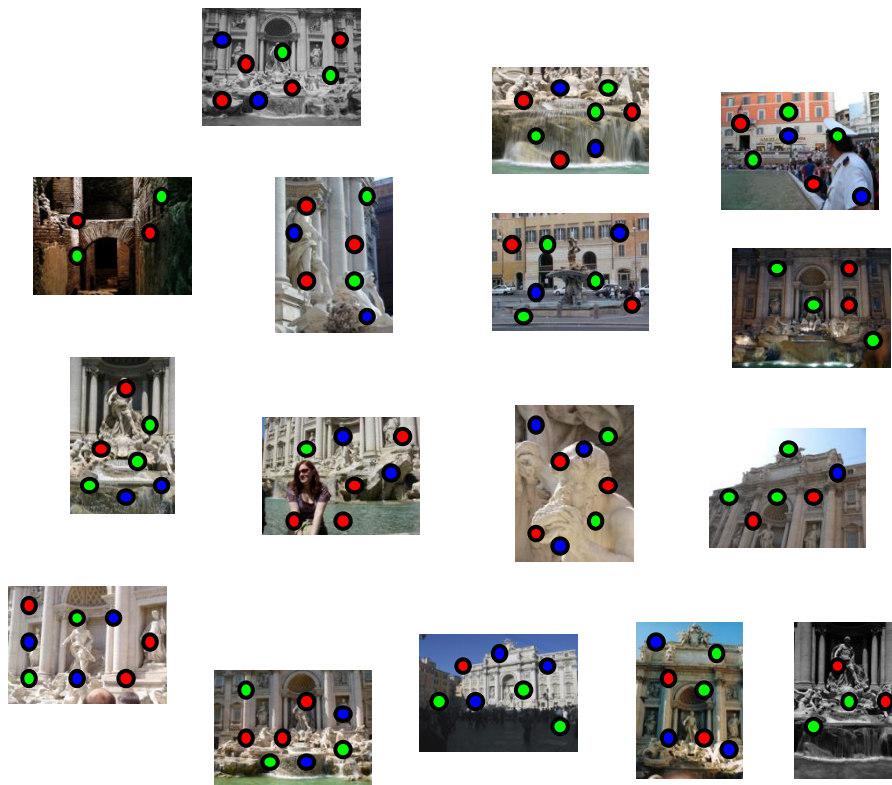
Camera 1
$R_1, t_1$ **?**

Camera 3
$R_3, t_3$ **?**

# Epipolar Constraint



- Projecting points on a ray to another camera forms a **line**, i.e., the corresponding point on the second image *must* lie on this line

# Epipolar Constraint



- Projecting points on a ray to another camera forms a **line**, i.e., the corresponding point on the second image *must* lie on this line

- Hence, if the *matched* 2D point is found along this line, we can determine its 3D location

# Epipolar Geometry



- Suppose we have two cameras with centers $O$, $O'$
- The **baseline** is the line connecting the origins

# Epipolar Geometry



- **Epipoles** $e$, $e'$ are where the baseline intersects the image planes, or projections of the other camera in each view

# Epipolar Geometry



**Baseline**

**Epipoles**

- Consider a **point** $X$, which projects to $x$ and $x'$

# Epipolar Geometry



- The plane formed by $X$, $O$, and $O'$ is called an **epipolar plane**
- There is a family of planes passing through $O$ and $O'$

# Epipolar Geometry



- **Epipolar lines** connect the epipoles to the projections of $X$
- Equivalently, they are intersections of the epipolar plane with the image planes – thus, they come in matching pairs

# Epipolar Geometry



- All **epipolar lines** pass through the **epipoles**
- **Epipoles** can lie outside of the image

# Example of Epipolar Lines

# Parallel to Image Plane



- Where are the epipoles and what do the epipolar lines look like?
- Epipoles **infinitely** far away, epipolar lines parallel

# Perpendicular to Image Plane



- Where are the epipoles and what do the epipolar lines look like?

- Epipole is "focus of expansion" and coincides with the principal point of the camera
- Epipolar lines go out from principal point

# Epipolar Constraint: Calibrated Case

$$x \cong K[R|t]X$$



- Assume the intrinsic and extrinsic parameters of the cameras are known (*calibrated*), and world coordinate system is set to that of the first camera

- Then the projection matrices are given by $K[I \mid 0]$ and $K'[R \mid t]$

- We can pre-multiply the projection matrices (and the image points) by the inverse calibration matrices to get *normalized* image coordinates:

$$x_{\text{norm}} = K^{-1}x_{\text{pixel}} \cong [I \mid 0]X, \qquad x'_{\text{norm}} = K'^{-1}x'_{\text{pixel}} \cong [R \mid t]X$$

# Epipolar Constraint: Calibrated Case $\boxed{x \cong K[R|t]X}$

$$x_{\text{norm}} \cong [I \mid 0]X \qquad\qquad X = (\rho x, 1)^T \qquad\qquad x'_{\text{norm}} \cong [R \mid t]X$$

$$[I \mid 0] \binom{\rho x}{1} \qquad x \qquad x' \qquad [R \mid t] \binom{\rho x}{1} = \rho Rx + t$$

$$t \qquad R$$

$$x' \cong \rho Rx + t \quad \Longrightarrow \quad x' \cdot [t \times (Rx)] = 0$$

- This means the three vectors $x'$, $Rx$, and $t$ are linearly dependent, i.e., lying on the same plane
- This constraint can be written using the *triple product*

# Epipolar Constraint: Calibrated Case

$$x \cong K[R|t]X$$

$$x_{\text{norm}} \cong [I \mid 0]X \qquad\qquad X = (\rho x, 1)^T \qquad\qquad x'_{\text{norm}} \cong [R \mid t]X$$

$$[I \mid 0]\binom{\rho x}{1} \qquad x \qquad\qquad x' \qquad [R \mid t]\binom{\rho x}{1} = \rho R x + t$$

$$t \qquad R$$

$$x' \cong \rho R x + t \implies x' \cdot [t \times (Rx)] = 0 \implies x'^T[t_\times]Rx = 0$$

$$\text{Recall: } a \times b = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix}\begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} = [a_\times]b \implies \text{skew-symmetric matrix: } A^T = -A$$

# Epipolar Constraint: Calibrated Case

$$x \cong K[R|t]X$$

$$x_{\text{norm}} \cong [I \mid 0]X$$

$$x'_{\text{norm}} \cong [R \mid t]X$$

$$X = (\rho x, 1)^T$$



$$[I \mid 0]\begin{pmatrix} \rho x \\ 1 \end{pmatrix}$$

$$x$$

$$x'$$

$$[R \mid t]\begin{pmatrix} \rho x \\ 1 \end{pmatrix} = \rho R x + t$$

$$t$$

$$R$$

$$x' \cong \rho R x + t \implies x' \cdot [t \times (Rx)] = 0 \implies x'^T[t_\times]Rx = 0 \implies x'^T E x = 0$$

**Essential Matrix**

# The Essential Matrix



$$l' = E^T x$$

$$x'^T \boxed{E x} = 0$$

$$(x', y', 1) \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$

- $Ex$ is the **epipolar line** associated with $x$ ($l' = Ex$)

- Recall: a line is given by $ax + by + c = 0$ or $l^T x = 0$ in homogeneous coordinates, where $l = (a, b, c)^T$ and $x = (x, y, 1)^T$

- $x'^T E x = x'^T l' = 0$ means $x'$ lies on the epipolar line $l'$

# The Essential Matrix



$$l = E^T x'$$

$$x'^T E x = 0$$

$$(x', y', 1) \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$

- $Ex$ is the **epipolar line** associated with $x$ ($l' = Ex$)
- $x'^T E x = x'^T l' = 0$ means $x'$ lies on the epipolar line $l'$
- Equivalently, $E^T x'$ is the **epipolar line** associated with $x'$ ($l = E^T x'$)
- $x'^T E x = l x = 0$ means $x$ lies on the epipolar line $l$

# The Essential Matrix
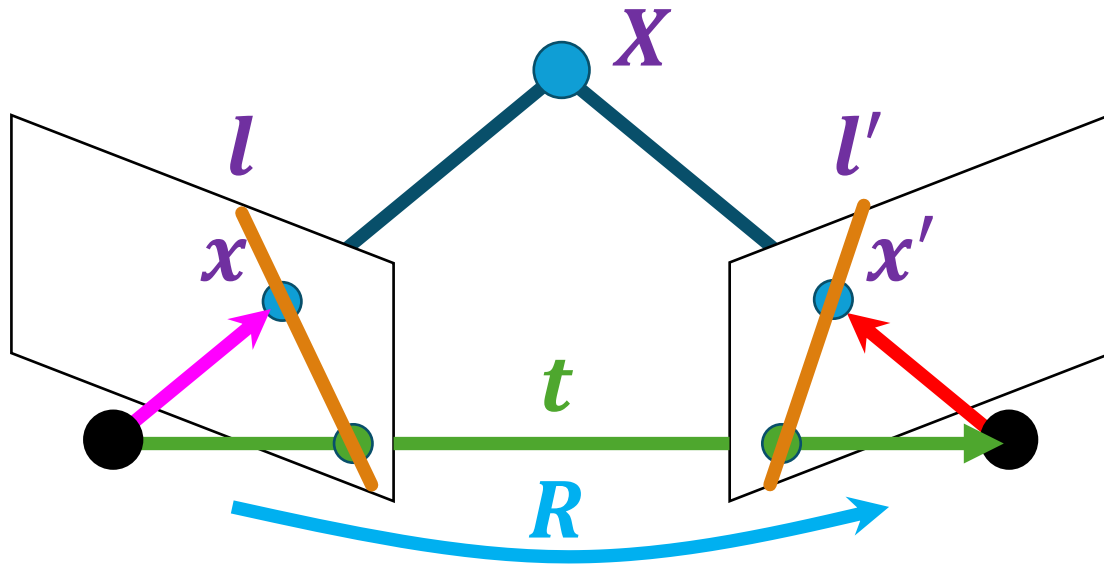


$$x'^T E x = 0$$

$$(x', y', 1) \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$

- $E x$ is the **epipolar line** associated with $x$ ($l' = E x$)
- $E^T x'$ is the **epipolar line** associated with $x'$ ($l = E^T x'$)
- $E e = 0$  and  $E^T e' = 0$ , where $e, e'$ are the epipoles
- $E$ is singular (rank two) and has five degrees of freedom => why?
  - $E = [t_\times] R$! $[t_\times]$ is skew-symmetric; has rank 2
  - $R$: 3 DoF, $t$: 3 DoF, but we lost 1 DoF due to scale; move along $t$ doesn't change $l$

# Epipolar Constraint: Uncalibrated

$$x \cong K[R|t]X$$
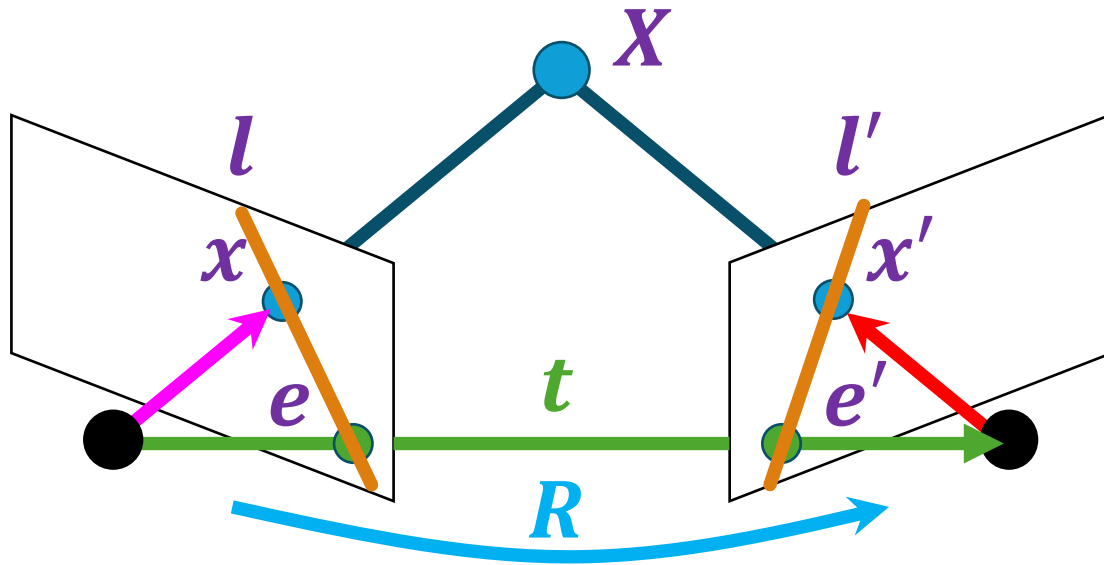


$$x'^{T} E x = 0$$

$$(x', y', 1)\begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix}\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$

- What if camera intrinsics $K$, $K'$ are unknown?
- We can write the epipolar constraint:

$$x'^{T}_{\text{norm}} E x_{\text{norm}} = 0$$

where $x_{\text{norm}} = K^{-1}x$, $x'_{\text{norm}} = K'^{-1}x'$

# Epipolar Constraint: Uncalibrated

$$x \cong K[R|t]X$$
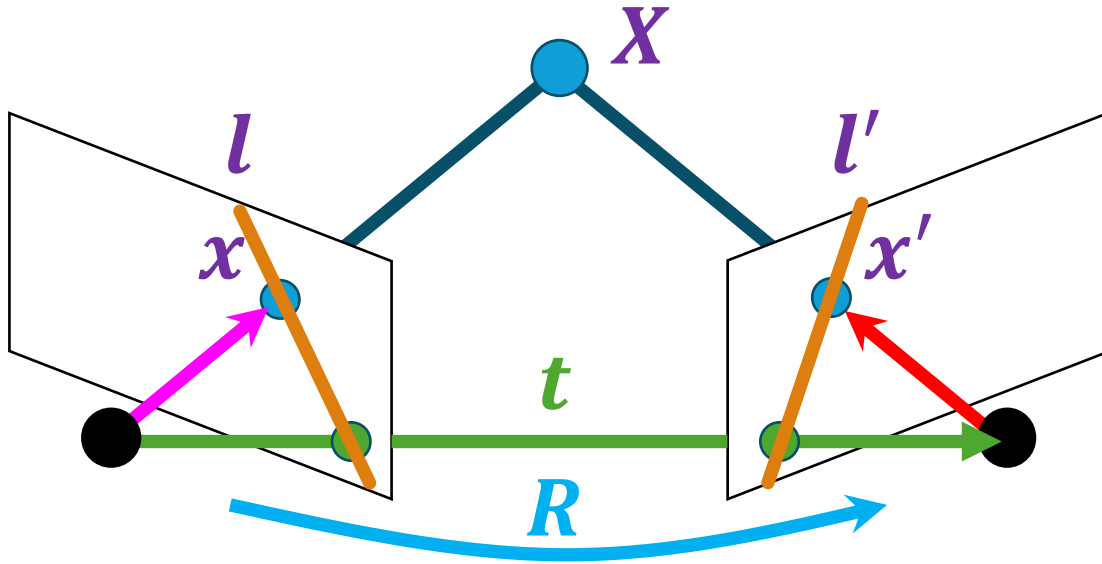


$$x'^T E x = 0$$

$$(x', y', 1) \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$

- What if camera intrinsics $K$, $K'$ are unknown?
- We can write the epipolar constraint:

**Fundamental Matrix**

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

$$x'^T_{\text{norm}} E x_{\text{norm}} = x'^T \boxed{K'^{-T} E K^{-1}} x = 0$$

$$\text{where } x_{\text{norm}} = K^{-1} x, \; x'_{\text{norm}} = K'^{-1} x'$$

# The Fundamental Matrix

$$x \cong K[R|t]X$$



$$x'^T F x = 0$$

$$(x', y', 1) \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0$$
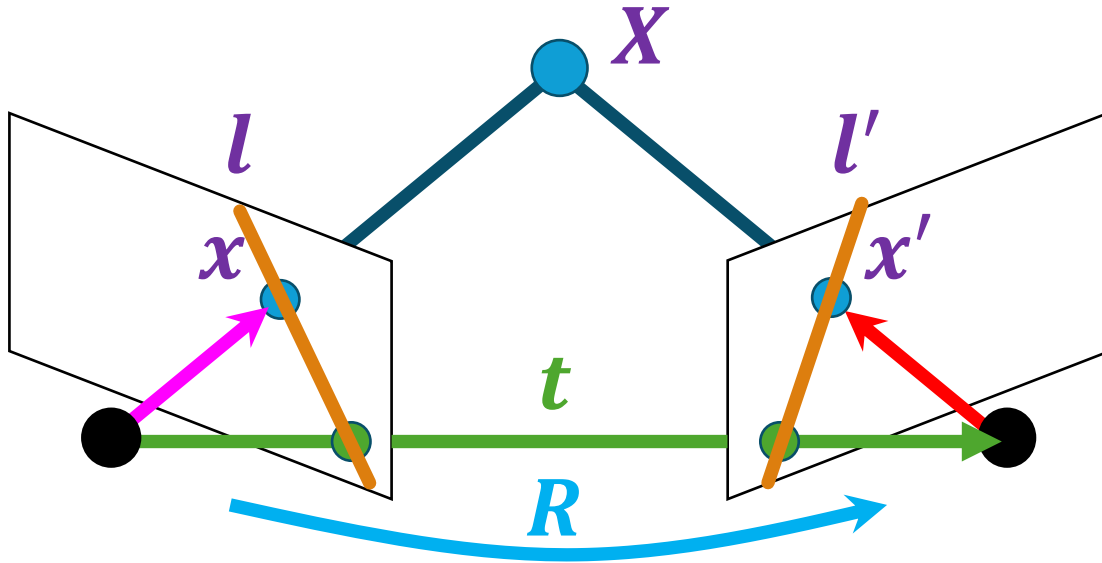
- $Fx$ is the **epipolar line** associated with $x$ ($l' = Fx$)

- $F^T x'$ is the **epipolar line** associated with $x'$ ($l = F^T x'$)

- $Fe = 0$ and $F^T e' = 0$, where $e, e'$ are the epipoles

- $F$ is singular (rank two) and has seven degrees of freedom => why?
  - $F = K'^{-T} [t_\times] R K^{-1}$! $[t_\times]$ is skew-symmetric; has rank 2
  - 9 elements, but we lost 1 DoF to scale, and 1 DoF to rank constraint

# How Can We Use the Epipolar Constraint?

- **Given:** $F$, $x$, $x'$

- **Q:** does there exist a 3D point that projects to $x$ and $x'$?

- **A:** Yes, if $\text{residual}(x'^T F x)$ is sufficiently low

- Note: the interpretation of $\text{residual}(x'^T F x)$ is the distance (geometric or algebraic) between $x$ and $l = F^T x'$, or $x'$ and $l' = F x$



$x$

$x'$

# How Can We Use the Epipolar Constraint?

- **Given:** $F$

- **Q:** how do we find $R$, $t$?

- **A:**

  - **Step 0:** estimate $K$, $K'$ if not known (self-calibration)

  - **Step 1:** compute $E = K'^T F K$

  - **Step 2:** since $E = [t_\times]R$, perform SVD on $E = U\Sigma V^T$. We have 4 solutions $(R_1, \pm t)$, $(R_2, \pm t)$, where $t = U[:, 3]$, $R_1 = UWV^T$, $R_2 = UW^T V^T$, and $W$ is the matrix that rotates 90° about z-axis.

  - **Step 3:** pick the one that gives 3D points in front of both cameras (cheirality check).



$x$        $x'$

# How to Estimate the Fundamental Matrix?

- **Given:** correspondences $\boldsymbol{x_i} = (x_i, y_i, 1)^T$ and $\boldsymbol{x_i'} = (x_i', y_i', 1)^T$
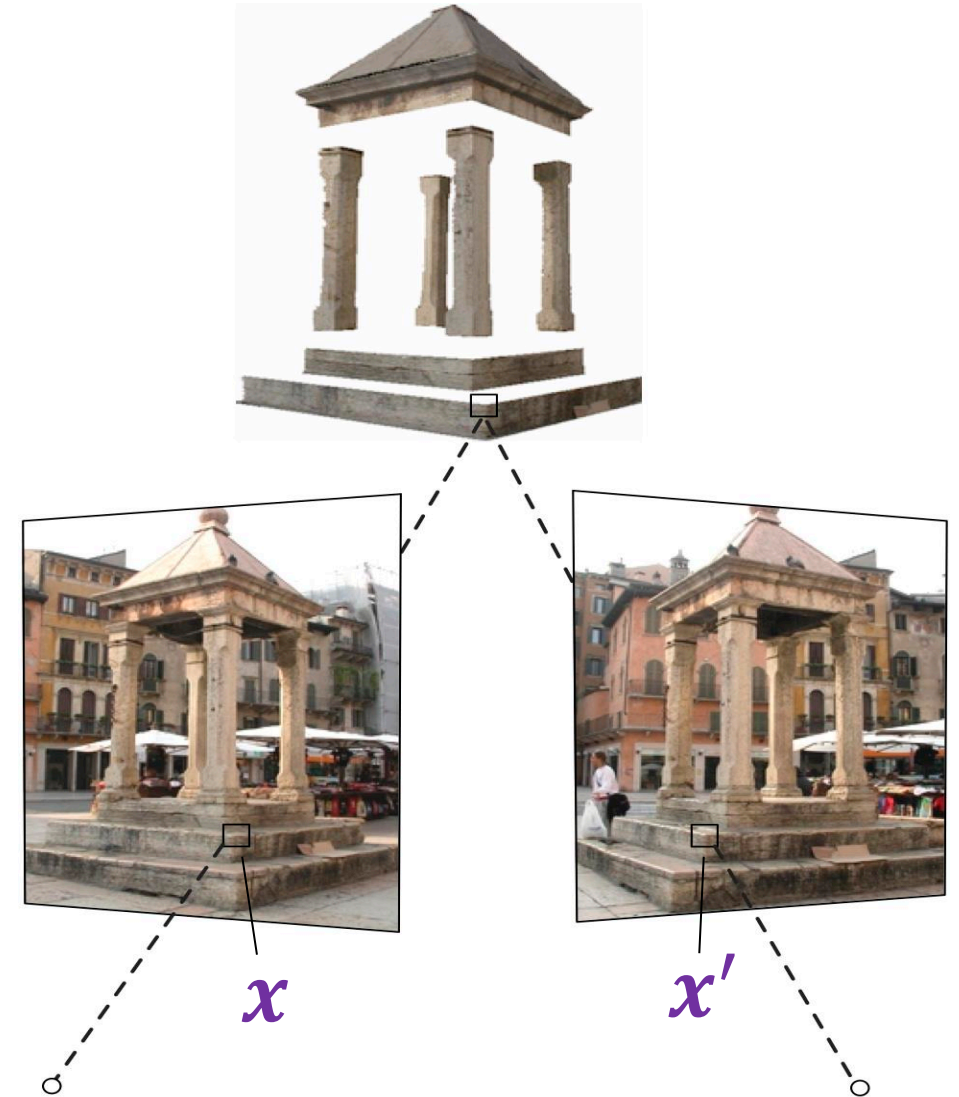- **Constraint:** $\boldsymbol{x'}^T \boldsymbol{F} \boldsymbol{x} = 0$

$$(x', y', 1) \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0 \quad \Longrightarrow \quad (x'x, x'y, x', y'x, y'y, y', x, y, 1) \begin{pmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{pmatrix} = 0$$

# How to Estimate the Fundamental Matrix?

- **Given:** correspondences $x_i = (x_i, y_i, 1)^T$ and $x'_i = (x_i', y_i', 1)^T$
- **Constraint:** $x'^T F x = 0$

$$\overbrace{\begin{bmatrix} x_1'x_1 & x_1'y_1 & x_1' & y_1'x_1 & y_1'y_1 & y_1' & x_1 & y_1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_n'x_n & x_n'y_n & x_n' & y_n'x_n & y_n'y_n & y_n' & x_n & y_n & 1 \end{bmatrix}}^{A} \begin{pmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{pmatrix} = 0$$

Homogeneous least squares to find $f$:

$$\arg\min_{\|f\|=1} \|Af\|_2^2 \quad \Longrightarrow \quad \text{Least eigenvector of } A^T A$$

This is known as the "*eight-point algorithm*" introduced by Christopher Longuet-Higgins (1981).

# Small Trick to Enforce Rank-2 Constraint

- We know $F$ must be singular/rank 2. How do we force that?
- **Solution:** take SVD of the initial estimate and throw out the smallest singular value

$$F_{\text{init}} = U\Sigma V^T$$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix} \longrightarrow \Sigma' = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$F = U\Sigma'V^T$$

# The Fundamental Matrix Song



By Daniel Wedge (2008): https://danielwedge.com/fmatrix/

# The Fundamental Matrix Song – Live!
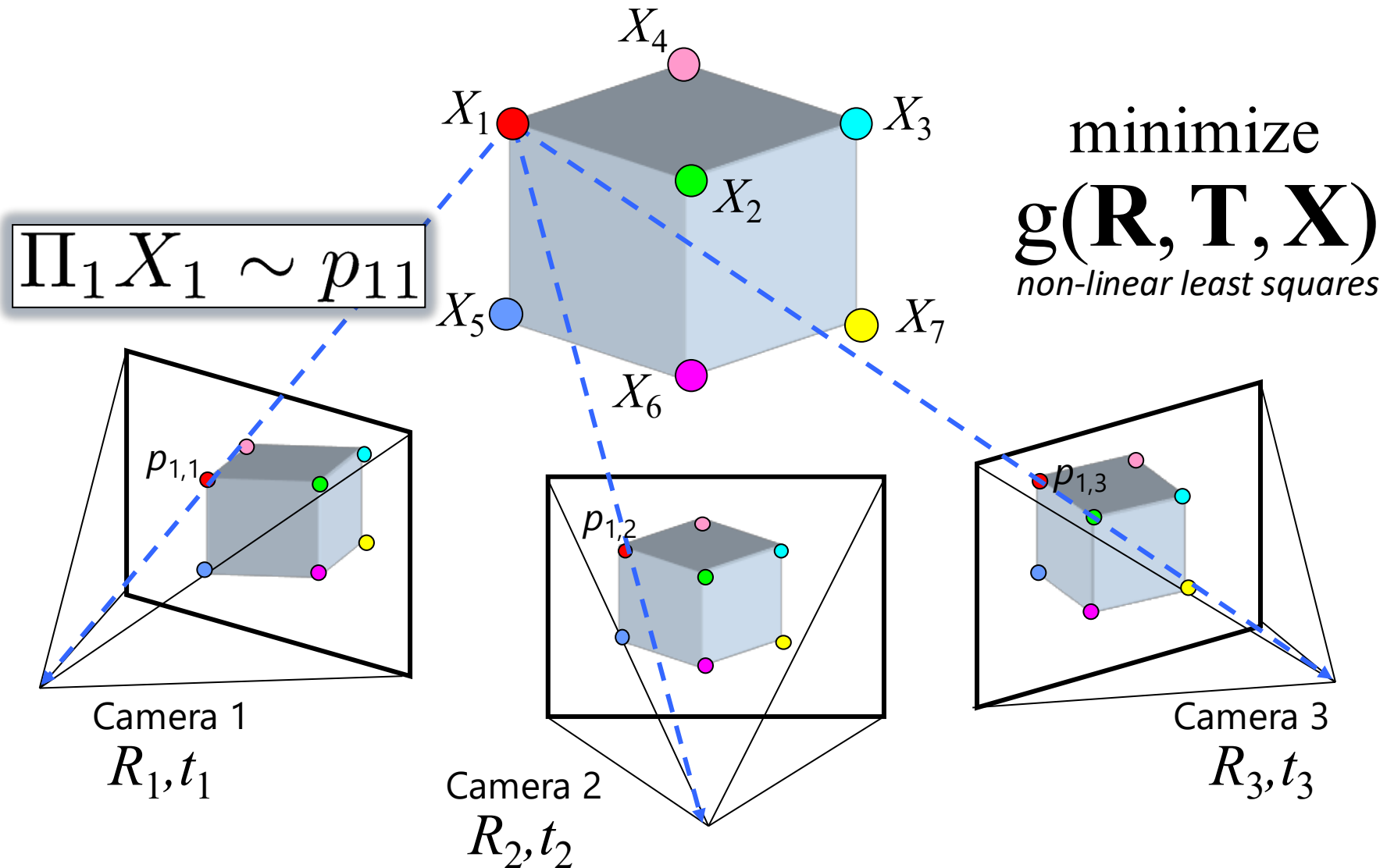


Daniel Wedge and the CVPR house band at CVPR 2023 in Vancouver

# More Than Two Views?

- The geometry of three views is described by a 3 x 3 x 3 tensor called the *trifocal tensor*

- The geometry of four views is described by a 3 x 3 x 3 x 3 tensor called the *quadrifocal tensor*

- After this it starts to get complicated…

- Or we can pose it as an optimization problem

# Structure from Motion (SfM)



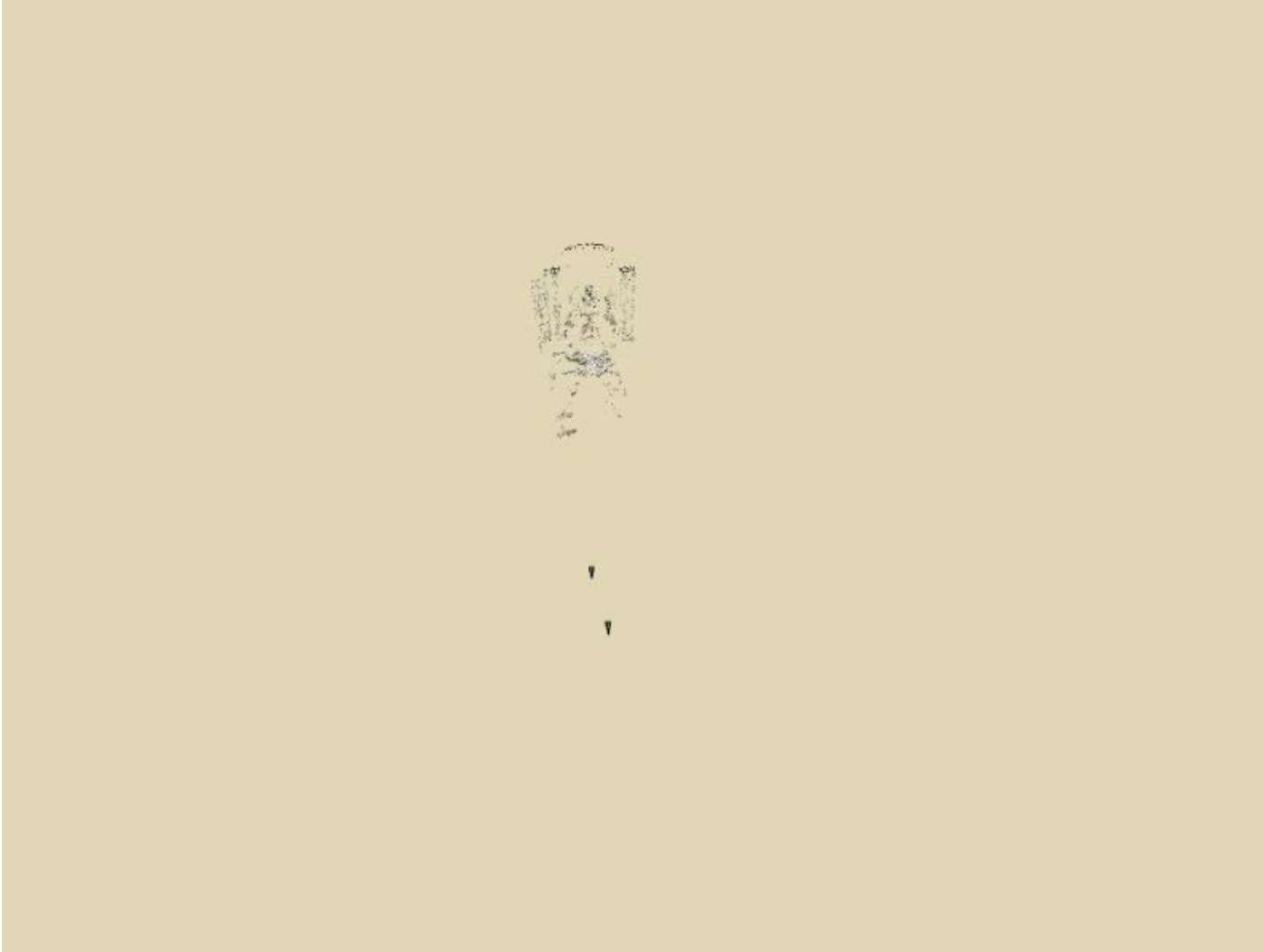$X_4$

$X_1$  $X_3$

$X_2$

minimize
$$g(\mathbf{R}, \mathbf{T}, \mathbf{X})$$
*non-linear least squares*

$$\Pi_1 X_1 \sim p_{11}$$

$X_5$  $X_7$

$X_6$

$p_{1,1}$

$p_{1,2}$

$p_{1,3}$

Camera 1
$R_1, t_1$

Camera 2
$R_2, t_2$

Camera 3
$R_3, t_3$

# Bundle Adjustment

- Minimize reprojection errors:

$$g(\mathbf{X}, \mathbf{R}, \mathbf{T}) = \sum_{i=1}^{m} \sum_{j=1}^{n} w_{ij} \cdot \left\| \mathbf{P}(\mathbf{x}_i, \mathbf{R}_j, \mathbf{t}_j) - \begin{bmatrix} u_{i,j} \\ v_{i,j} \end{bmatrix} \right\|^2$$

*indicator variable*:
is point *i* visible in image *j* ?

*predicted*
image location

*observed*
image location

- Optimized using non-linear least squares, e.g., Levenberg-Marquardt algorithm
- Susceptible to local minima; requires careful initialization

# Incremental Structure from Motion



- Photo Tourism (Snavely et al., SIGGAPH'06)

- Trevis Fountain, Rome
  - 466 Internet photos
  - > 100,000 3D points
  - Very large optimization problem

# Practical SfM Tools

- COLMAP (by Schönberger et al.): https://colmap.github.io/
- nerfstudio COLMAP Python wrapper: https://docs.nerf.studio/quickstart/custom_dataset.html

- Still, with thousands of images, it can take many hours or even days!
- There's a family of methods optimized for efficiency and continuous streams, often referred to as *Simultaneous Localization and Mapping* (**SLAM**), particularly useful for robot navigation for instance
  - One classic, widely-used system is *ORB-SLAM* (Mur-Artal et al., 2015)

- We will discuss learning-based approaches later

# Part 1&2 Summary – Multi-view Geometry

- Camera model and projection
  - Homogeneous coordinates
  - Pinhole camera, perspective projection $x \cong K[R|t]X$
  - Intrinsics $K$, extrinsics $[R|t]$
  - Camera calibration

- Epipolar geometry
  - Epipolar plane, epipolar lines
  - Essential matrix $E$ => $x'^T_{\text{norm}} E x_{\text{norm}} = 0$ (calibrated)
  - Fundamental matrix $F$ => $x'^T E x = 0$ (uncalibrated)

- Structure from Motion (SfM)
  - Bundle adjustment