

CS131: Computer Vision: Foundations and Applications

Multi-view Geometry

Juan Carlos Niebles and Adrien Gaidon



Stanford University

Reference: [Szeliski](#) 11.3 & 12.1, [H&Z ch. 9](#)

Most slides adapted from S. Lazebnik, J. Johnson, D. Fouhey

The story so far

Model the 3D-to-2D camera projection: $\tilde{\mathbf{x}}^I \sim \mathbf{P}\tilde{\mathbf{X}}^W$ with $\mathbf{P} = \mathbf{K}[\mathbf{R}|\mathbf{t}]$

Calibrate cameras (get $\mathbf{K}[\mathbf{R}|\mathbf{t}]$) from N 2D-3D correspondences $(\tilde{\mathbf{x}}_i^I, \tilde{\mathbf{X}}_i^W)$

- cast constraints (2D-3D correspondences) as a linear system $\mathbf{A}\mathbf{p} = \mathbf{0}$
- total least squares ($\operatorname{argmin}_{\mathbf{x}} \|\mathbf{A}\mathbf{p}\|^2$ s. t. $\|\mathbf{p}\|^2 = 1$) gives the best approximation
- closed-form solution via the SVD of \mathbf{A} (its last right singular vector) & Cholesky
- refine by minimizing the 2D reprojection error $\sum_i \|\operatorname{proj}(\mathbf{K}[\mathbf{R}|\mathbf{t}]\mathbf{X}_i; \boldsymbol{\kappa}) - \mathbf{x}_i\|^2$

In general, we don't have **3D** measurements...

... but more than 1 image!

What are the geometric constraints governing

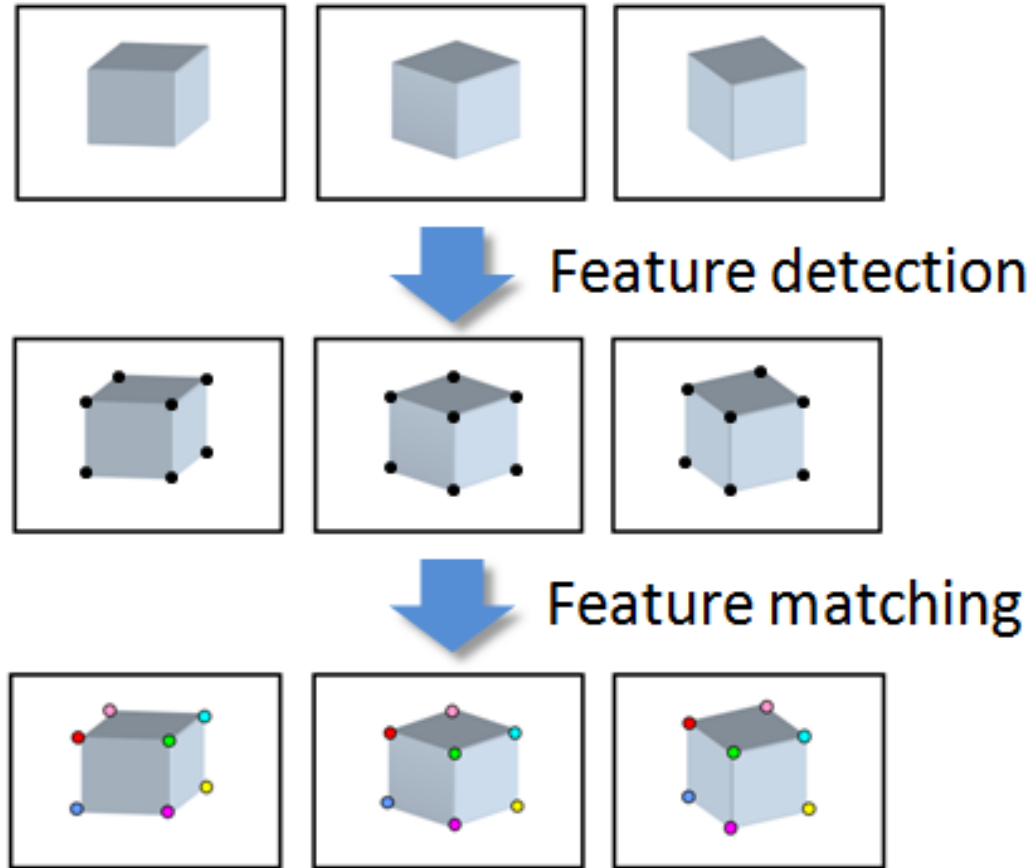
multiple views of the same scene?

→ **2D** correspondences!

Get 3D structure & motion from 2D correspondences

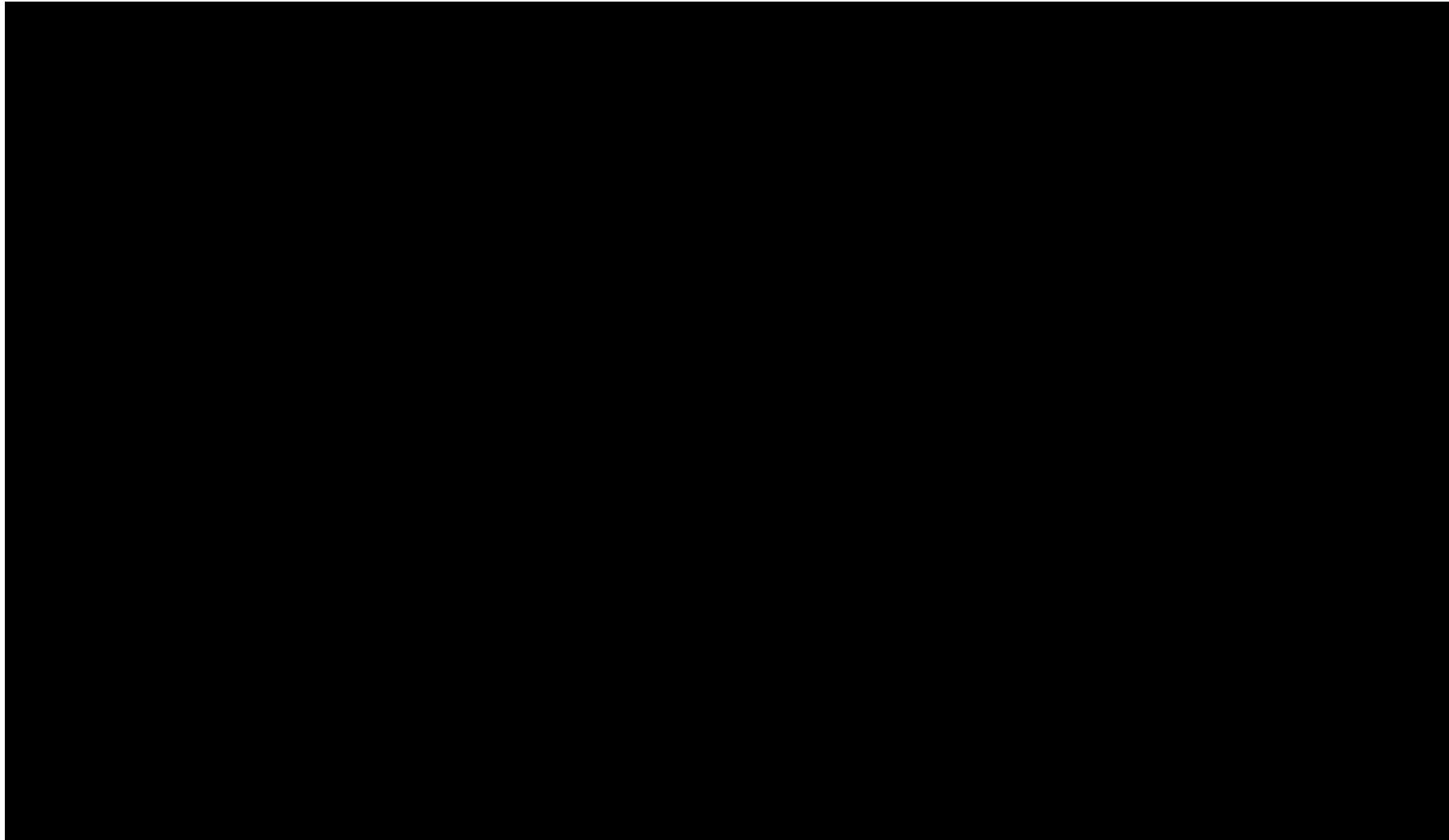


Correspondence estimation



Can use a wide range of features (cf. example in Project 1).
More details in upcoming lectures by Juan Carlos

Get 3D structure & motion from 2D correspondences



Dubrovnik, Croatia. 4,619 images (out of an initial 57,845 downloaded from Flickr). 3.5M points!
Total reconstruction time: 17.5 hours on 352 cores

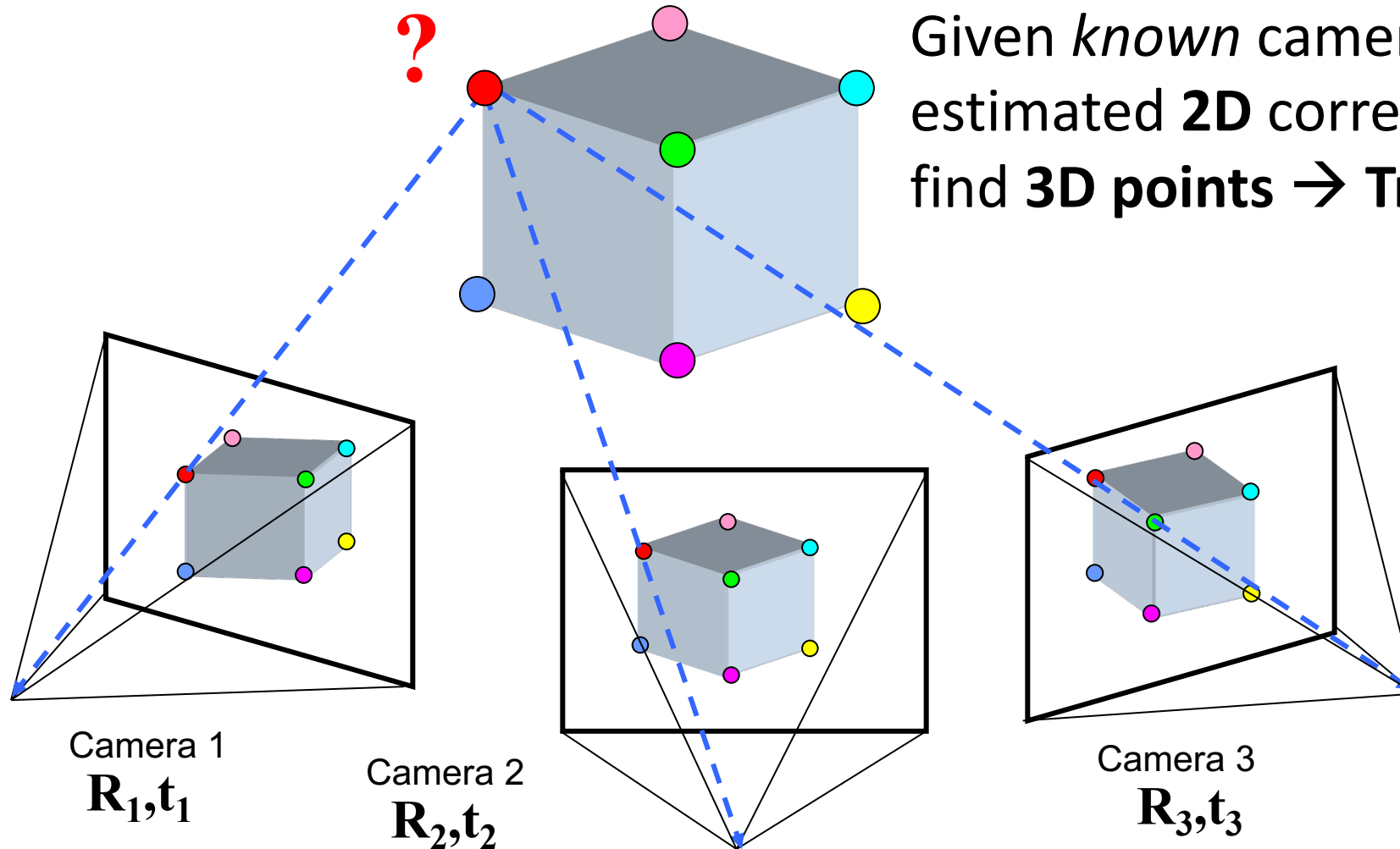
Building Rome in a Day, Agarwal et al, ICCV'09
<http://grail.cs.washington.edu/rome/>

Why do we care about 3D reconstruction?

- Mapping, Localization, Navigation for Robots, [Drones](#), [Cars](#)
(cf. also visual [SLAM](#))
- AR (e.g., HoloLens) and VR (e.g., Oculus)
- Movies ([special FX](#)), Digital Preservation, “[Photo Tourism](#)”, ...
- Software: [COLMAP](#) (SfM), [orb-slam2](#) / [g2o](#) / [gtsam](#) (SLAM)
- Hot topic in industry & academia (top category at CVPR)

Multi-view geometry problems

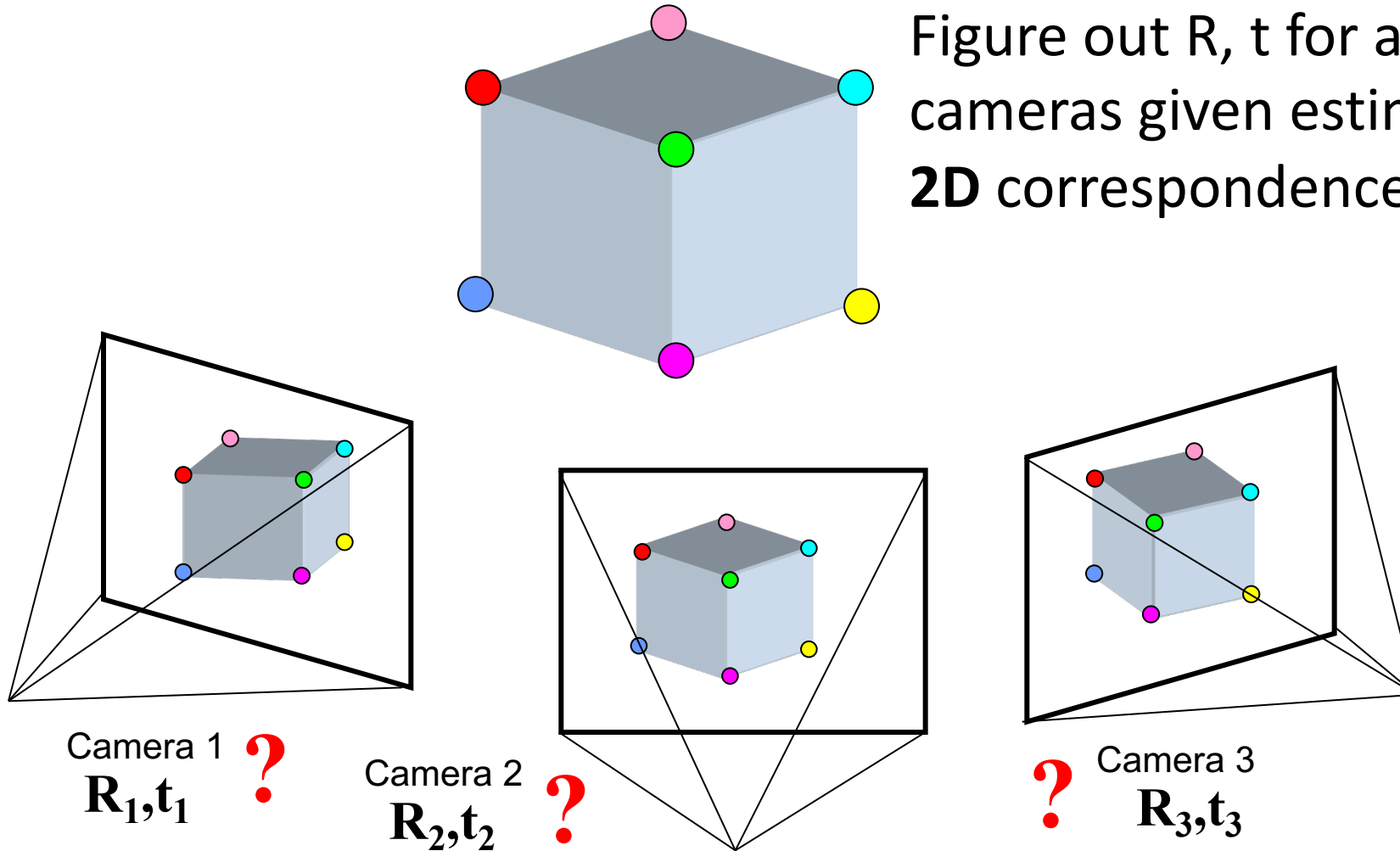
Recovering structure:
Given *known* cameras and
estimated **2D** correspondences,
find **3D points** → **Triangulation**



Multi-view geometry problems

Motion:

Figure out R, t for a set of cameras given estimated **2D** correspondences



What will we learn today?

Triangulation

Epipolar geometry

Stereo

Structure-from-Motion (SfM)

What will we learn today?

Triangulation

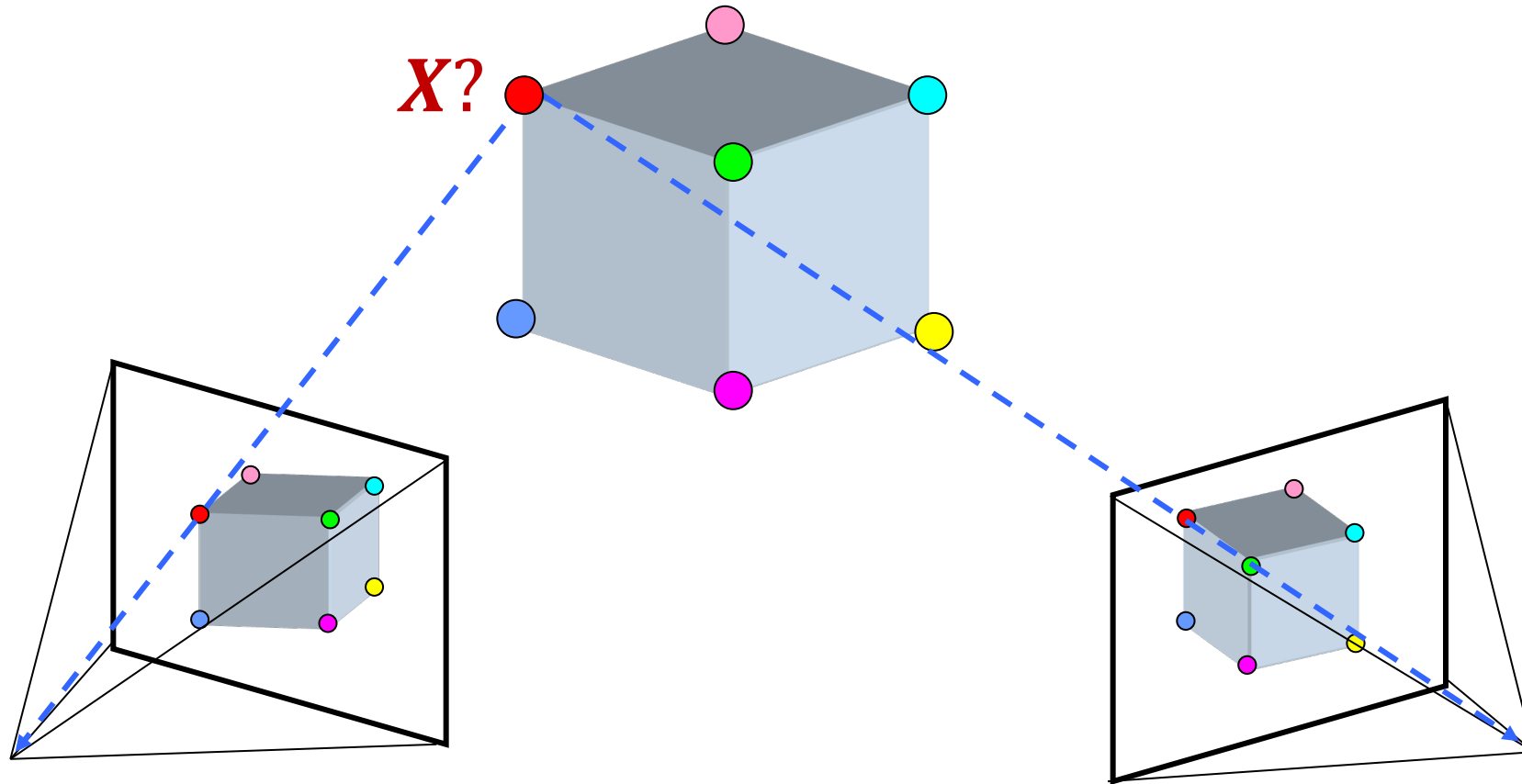
Epipolar geometry

Stereo

Structure-from-Motion (SfM)

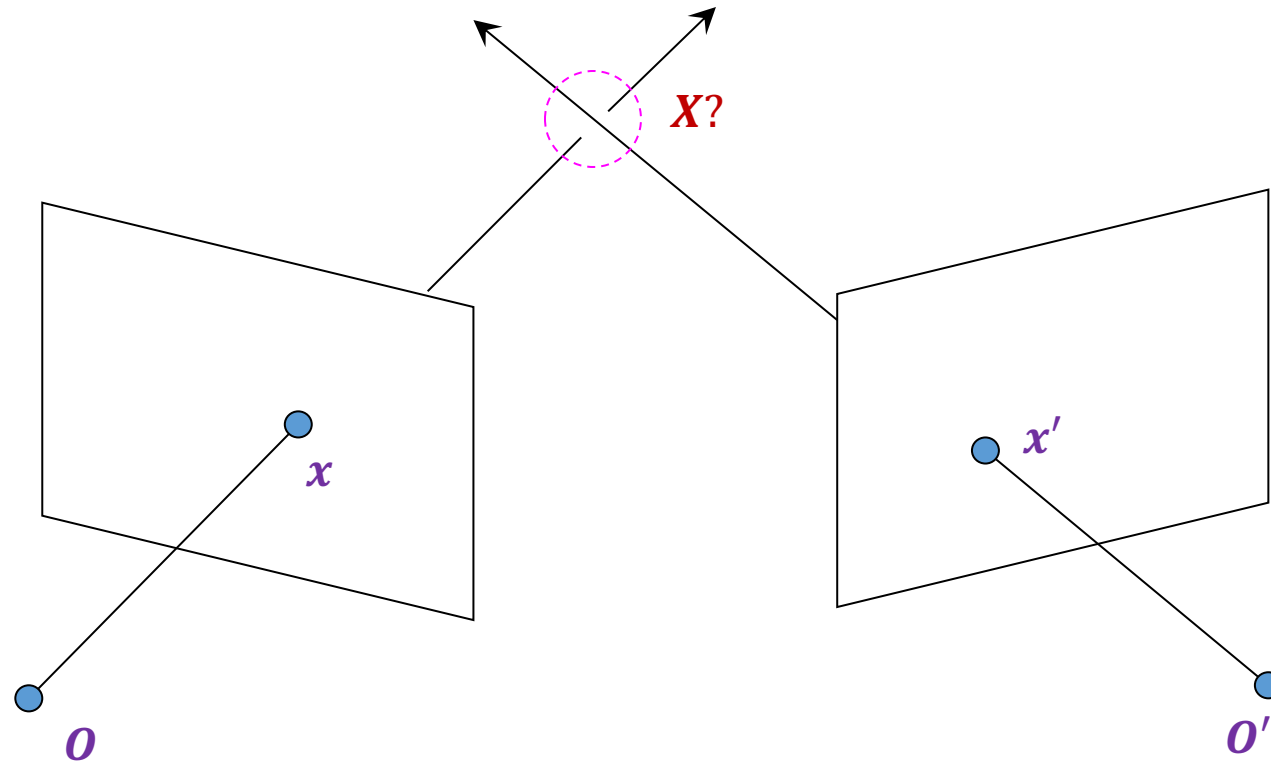
Triangulation

Given projections of a 3D point in two or more images (with known camera matrices), find the coordinates of the point



Triangulation

Given projections of a 3D point in two or more images (with known camera matrices), find the coordinates of the point

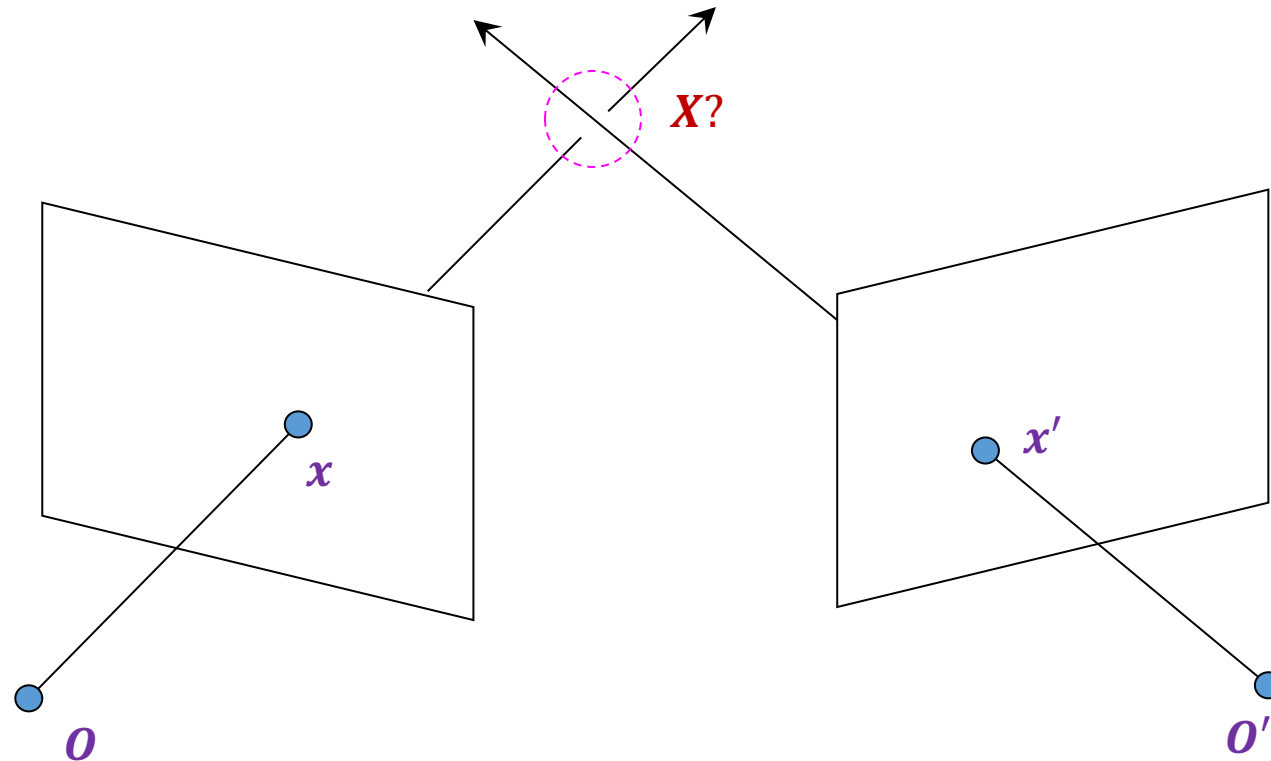


Triangulation

We want to intersect the two visual rays corresponding to x and x'

But do they always intersect *exactly*?

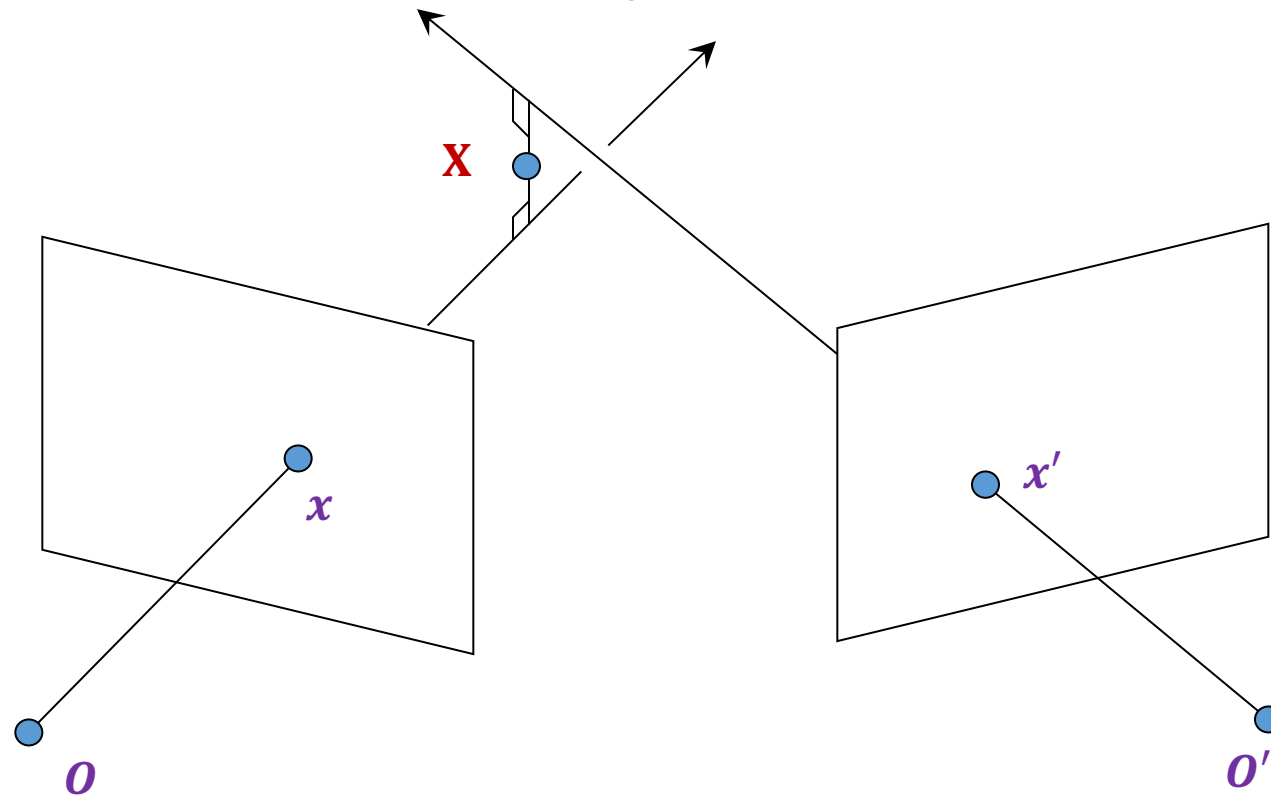
No! Noise in 2D matching or numerical errors



Triangulation: linear approach

Find the shortest segment connecting the two viewing rays

Let \mathbf{X} be the midpoint of that segment: solve for \mathbf{X} !



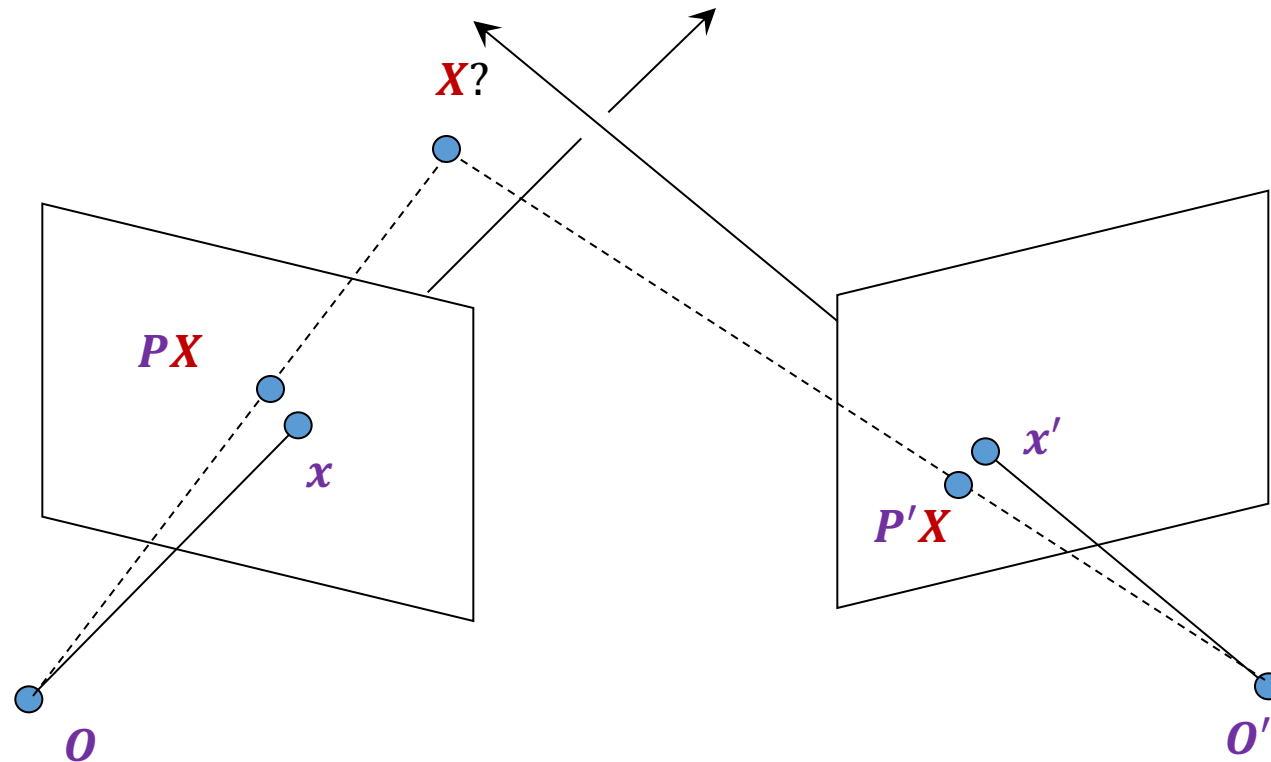
what is \mathbf{A} ?
(answer in appendix)

As for calibration: constraints $(\mathbf{x} \sim \mathbf{P}\mathbf{X}, \mathbf{x}' \sim \mathbf{P}'\mathbf{X}) \rightarrow \mathbf{A}\mathbf{X} = \mathbf{0} \rightarrow$ SVD of \mathbf{A}

Triangulation: *non-linear* approach

Find \mathbf{X} that minimizes the 2D reprojection errors

$$\|\text{proj}(\mathbf{P}\mathbf{X}) - \mathbf{x}\|^2 + \|\text{proj}(\mathbf{P}'\mathbf{X}) - \mathbf{x}'\|^2$$



What will we learn today?

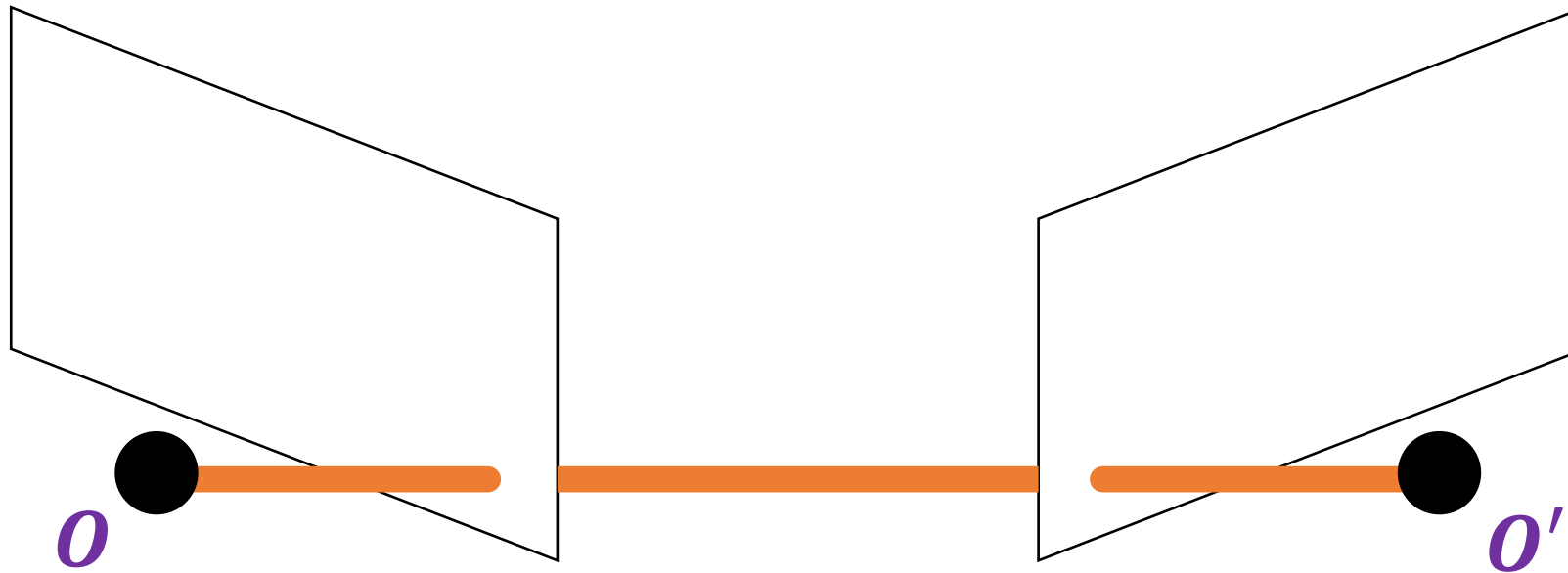
Triangulation

Epipolar geometry

Stereo

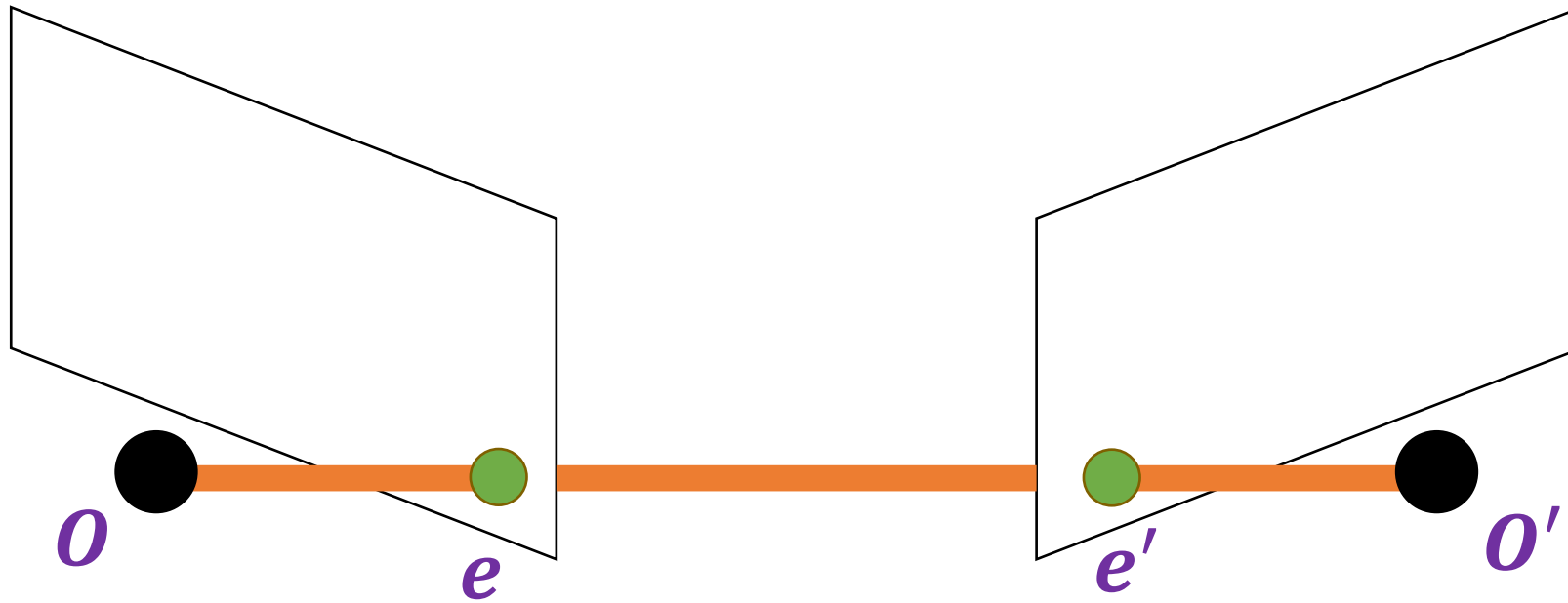
Structure-from-Motion (SfM)

Epipolar geometry setup



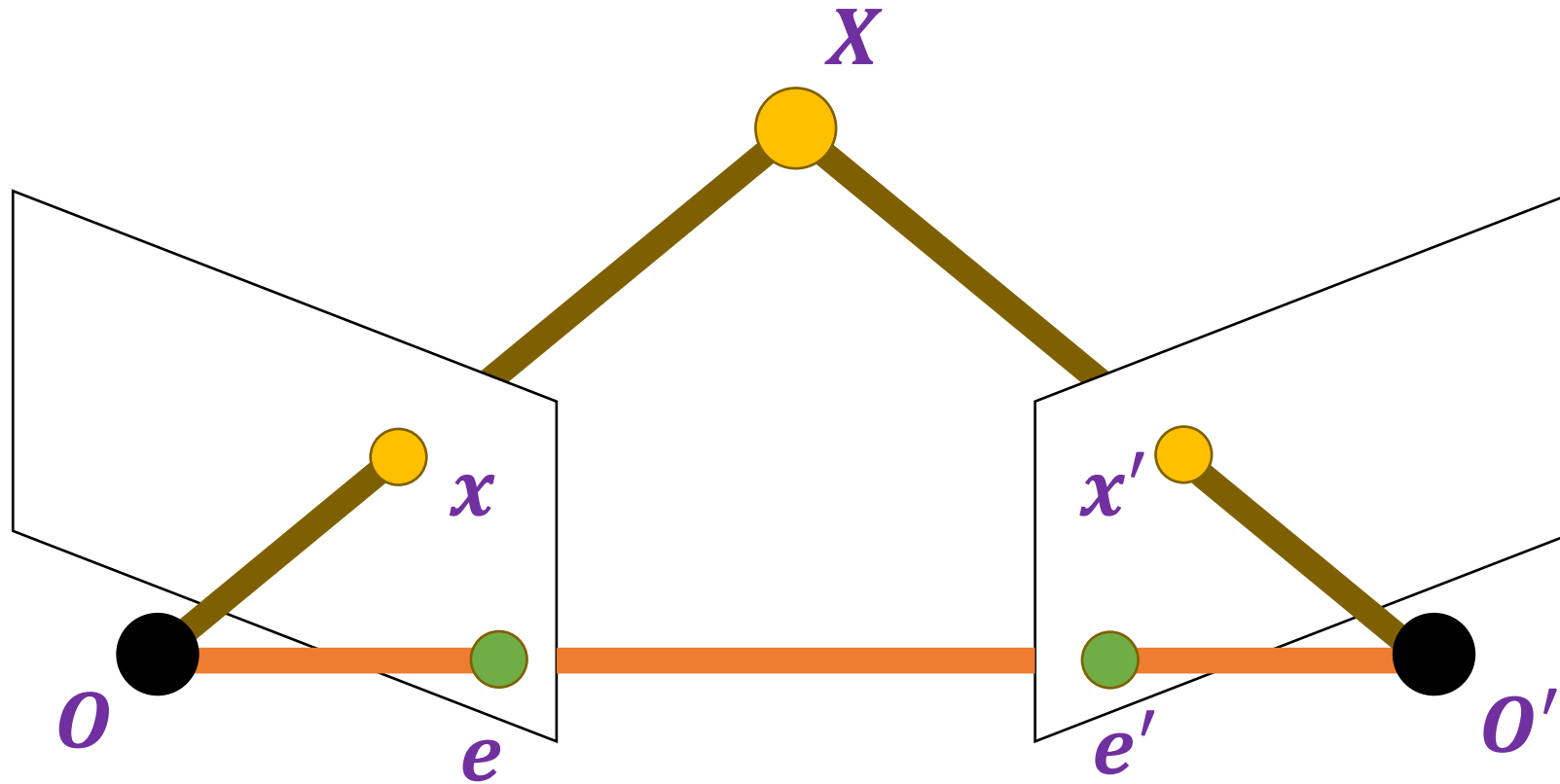
- Suppose we have two cameras with centers O, O'
- The **baseline** is the line connecting the origins

Epipolar geometry setup



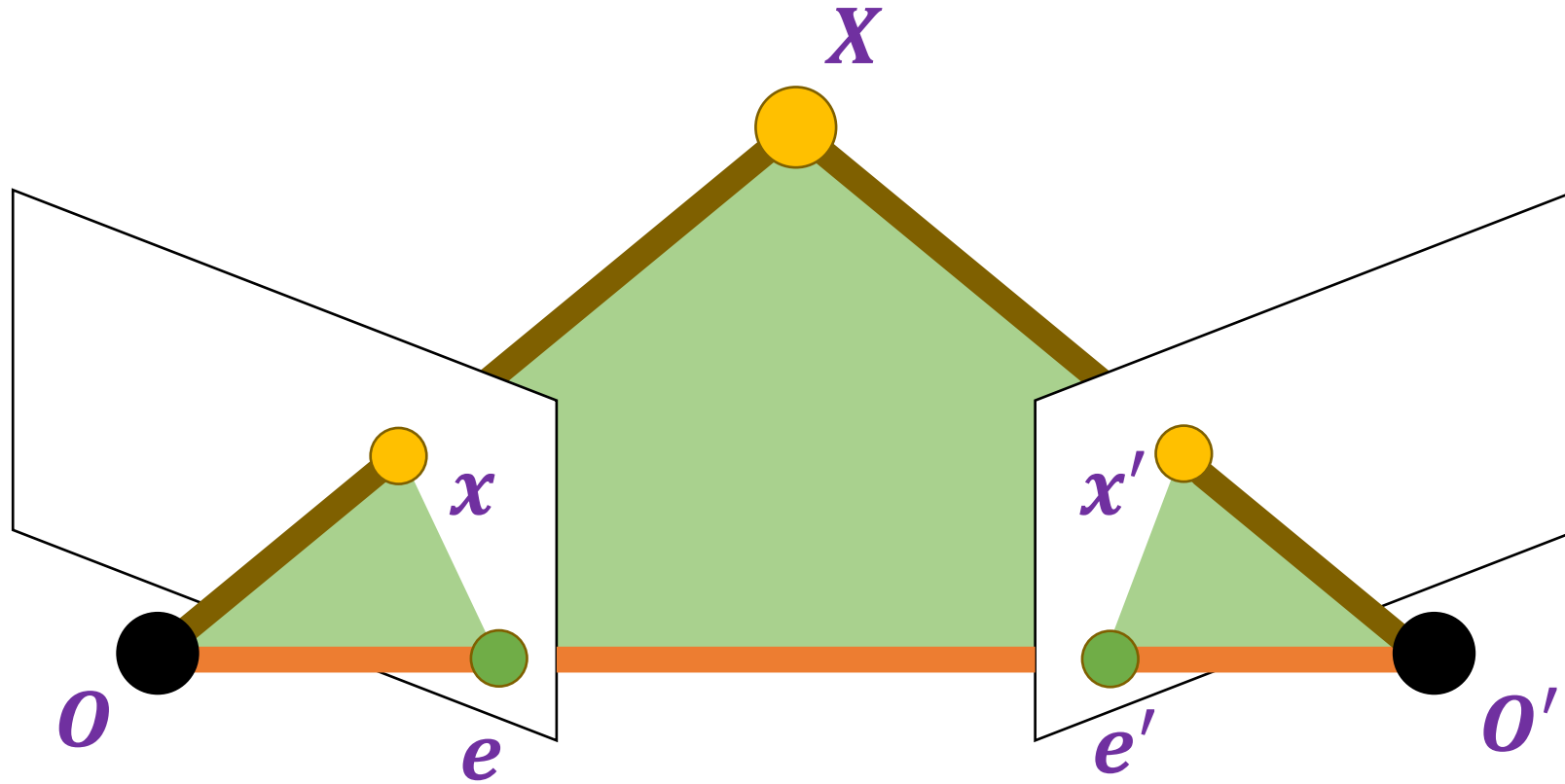
- **Epipoles** e, e' are where the baseline intersects the image planes
- Equivalently: epipoles are projections of the other camera in each view

Epipolar geometry setup



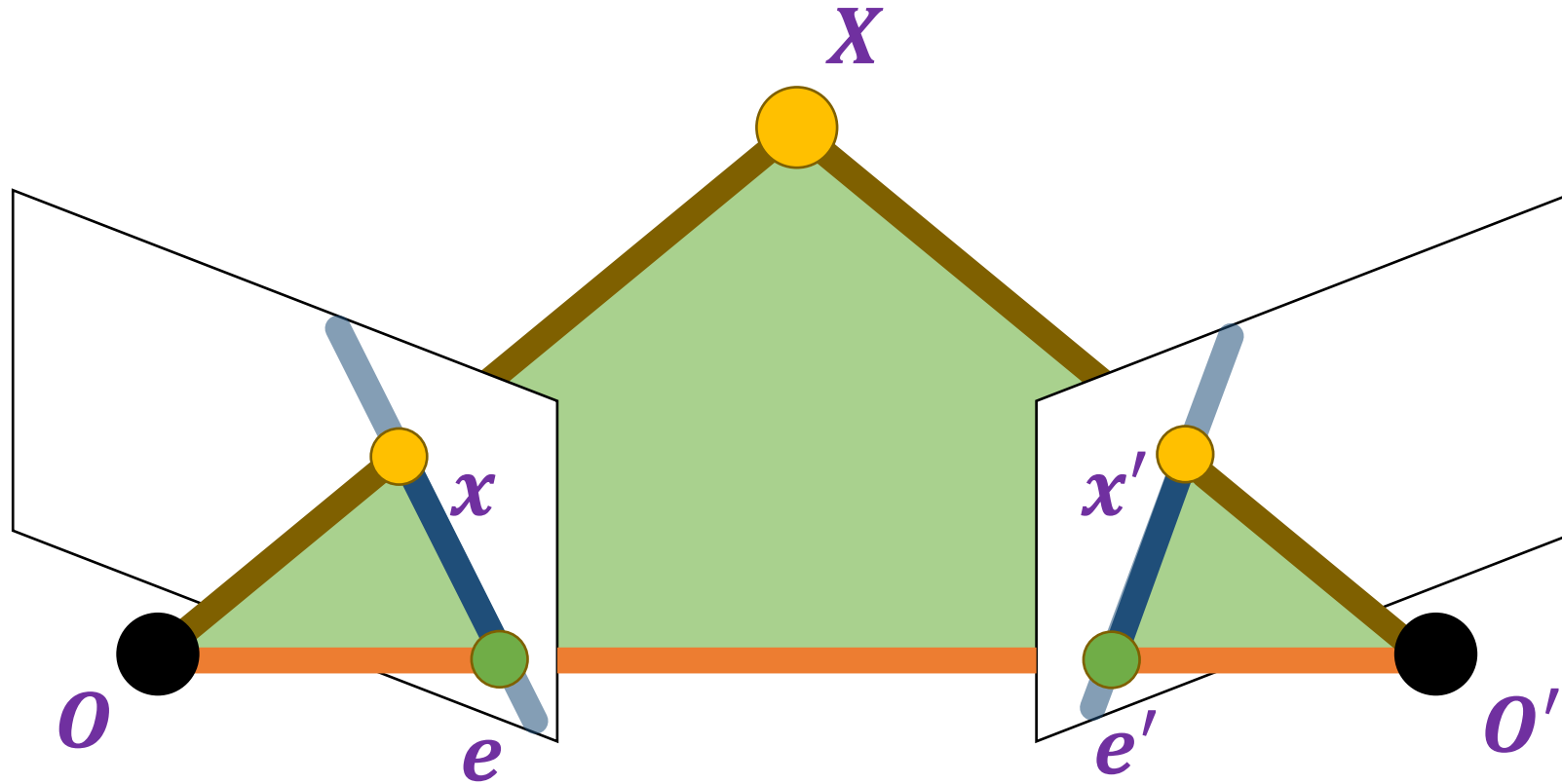
- Consider a **point X** , which projects to x and x'

Epipolar geometry setup



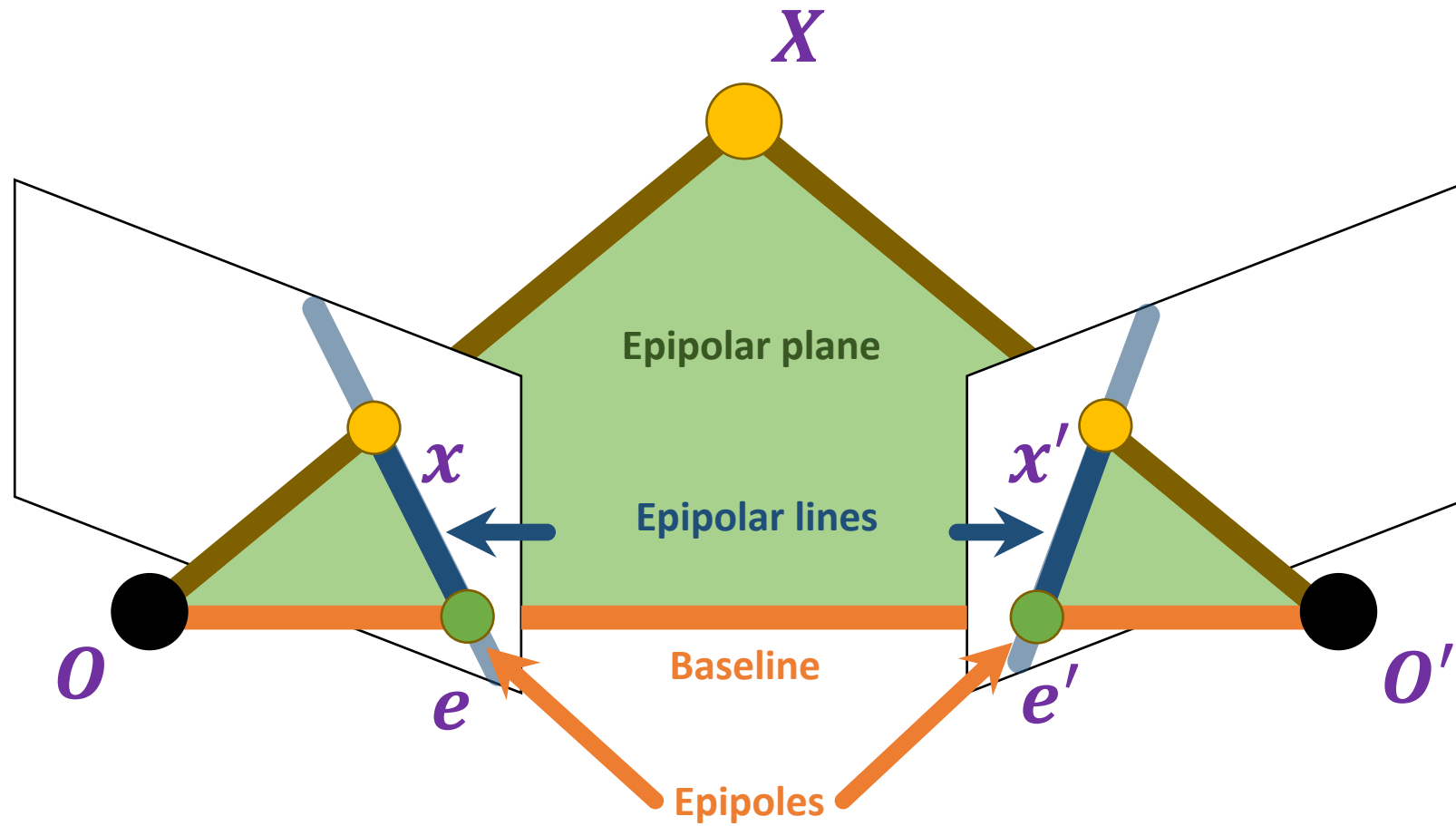
- The plane formed by X , O , and O' is called an **epipolar plane**

Epipolar geometry setup

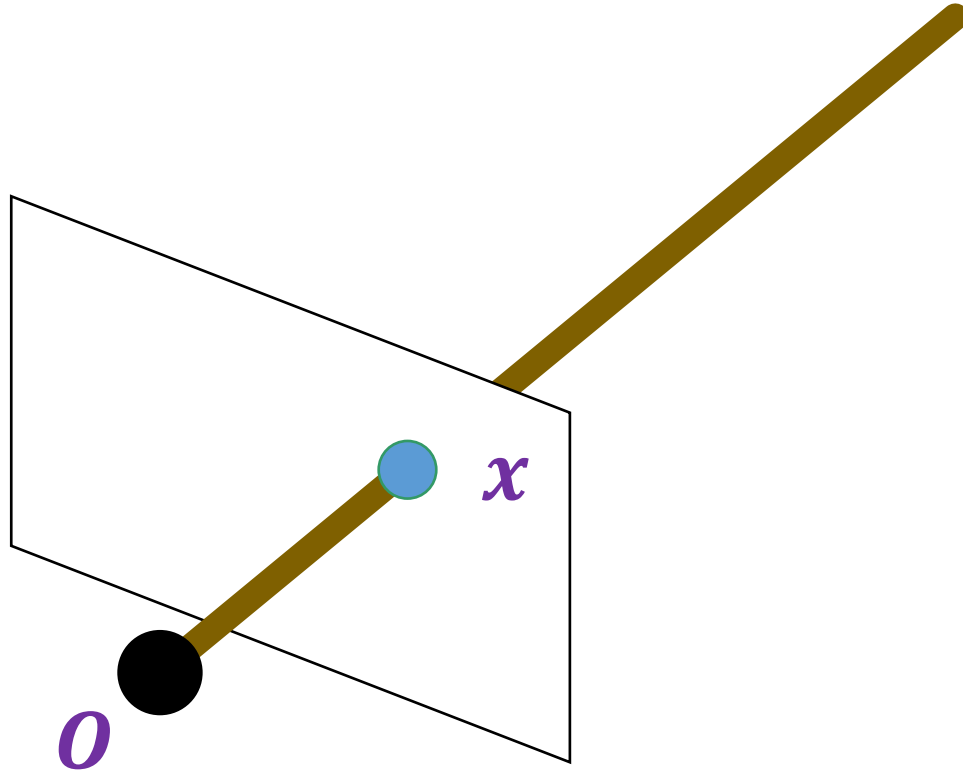


- **Epipolar lines** connect the epipoles to the projections of X
- Equivalently, they are intersections of the epipolar plane with the image planes, come in pairs (for x and x')

Epipolar geometry setup: Summary

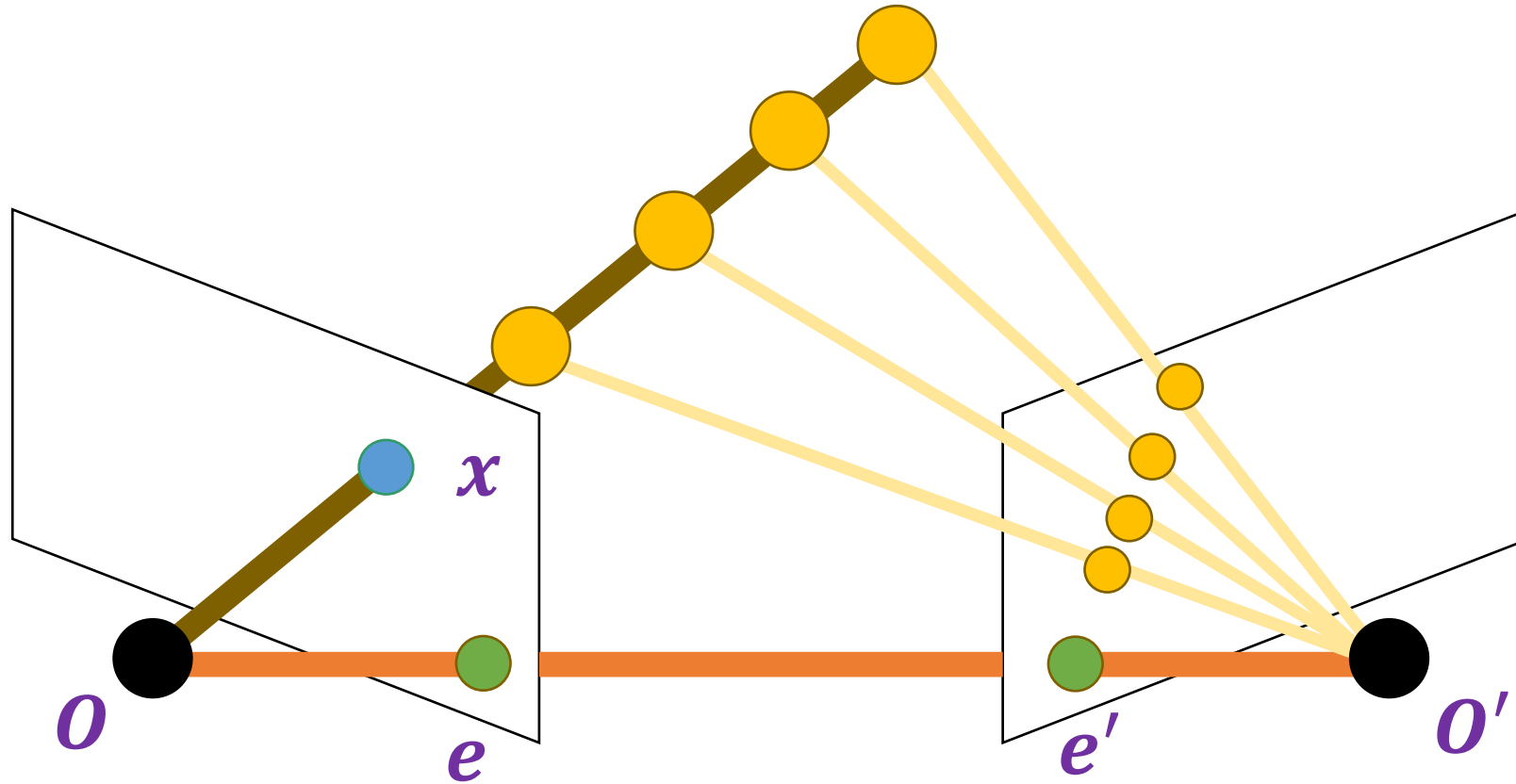


Epipolar constraint



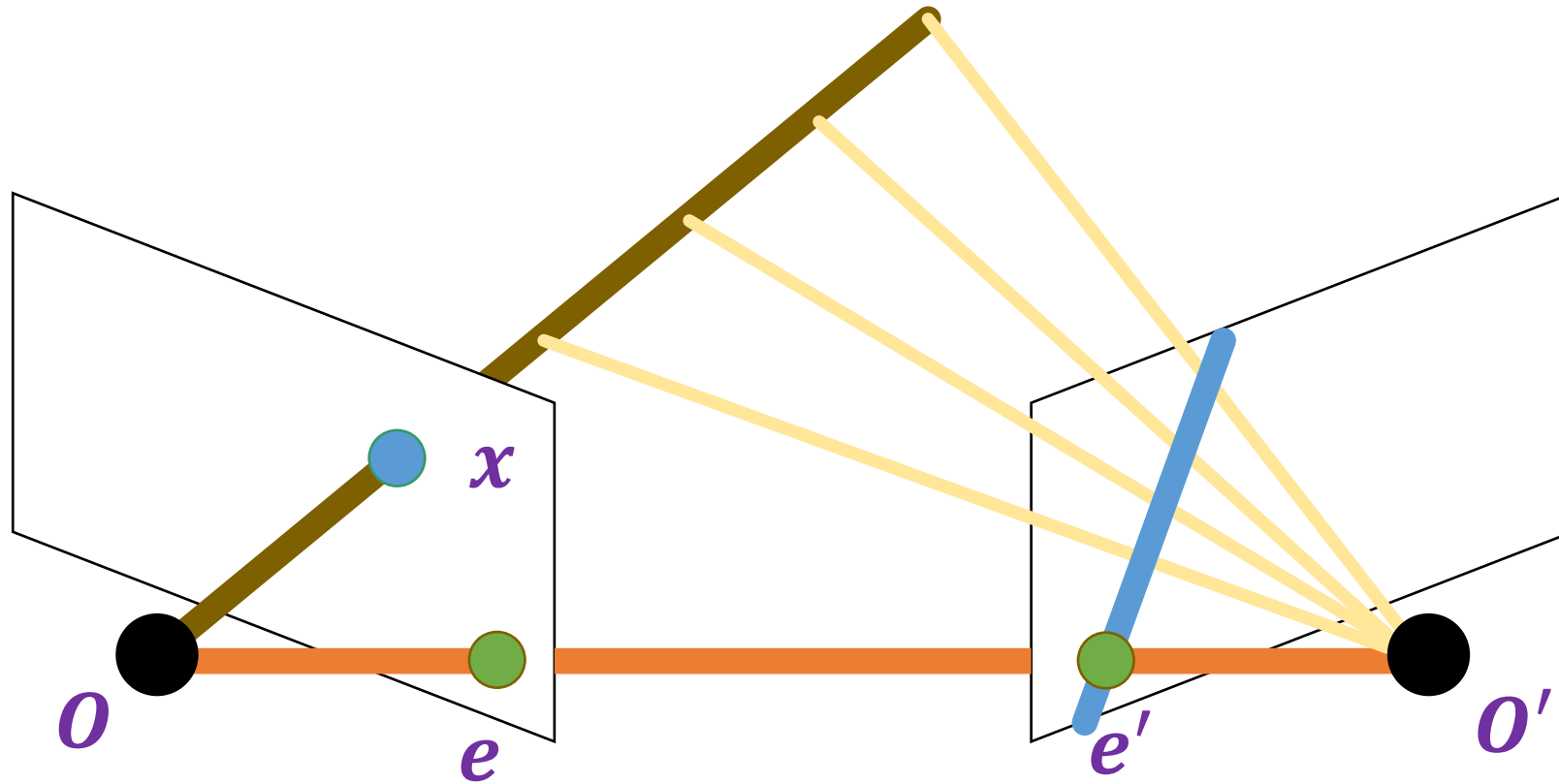
- Suppose we observe a single point x in one image

Epipolar constraint



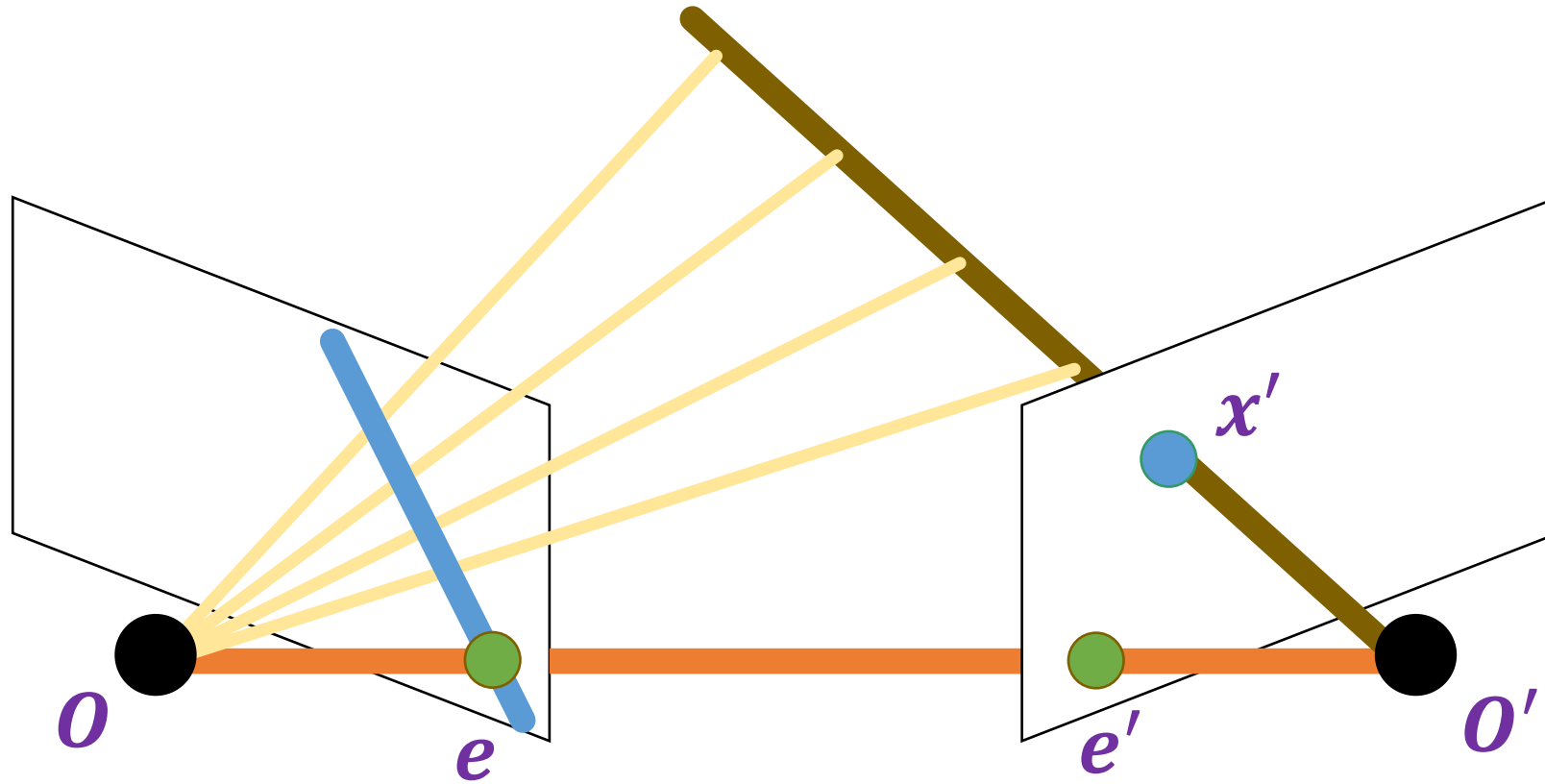
- Where can we find the x' corresponding to x in the other image?

Epipolar constraint



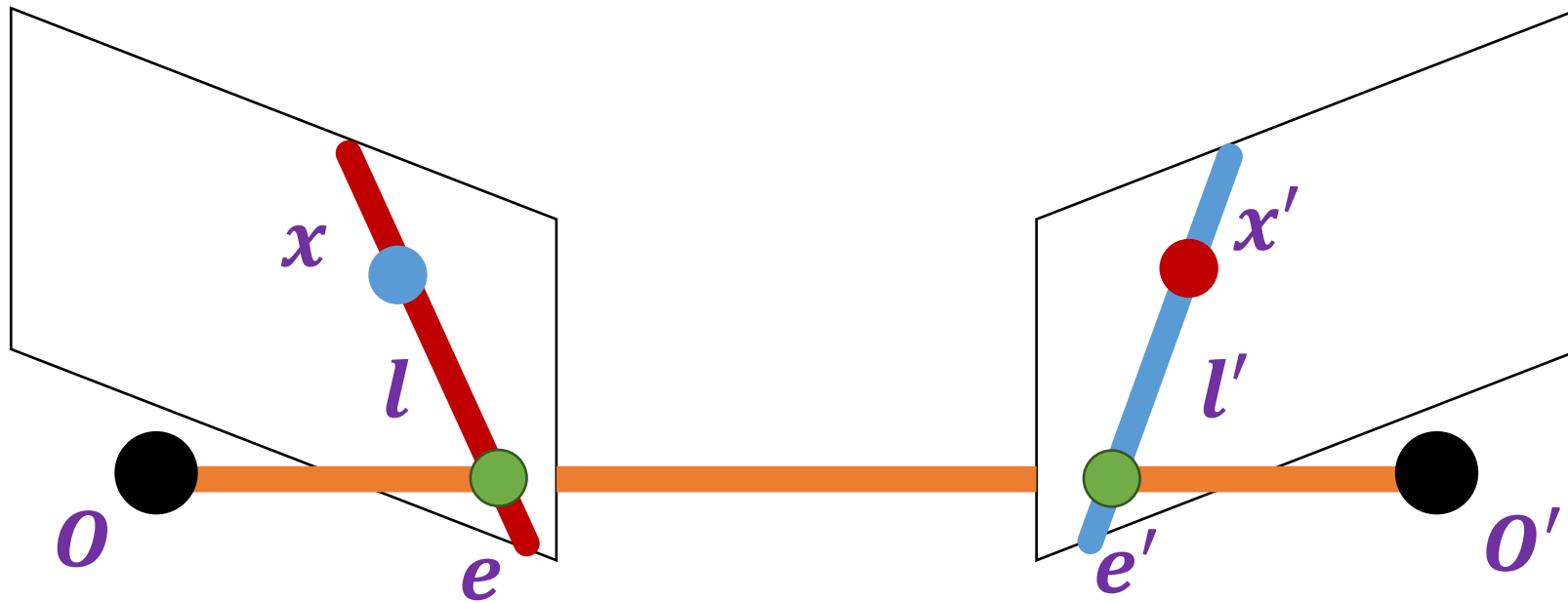
- Where can we find the x' corresponding to x in the other image?
- Along the **epipolar line** corresponding to x (projection of visual ray connecting O with x into the second image plane)

Epipolar constraint



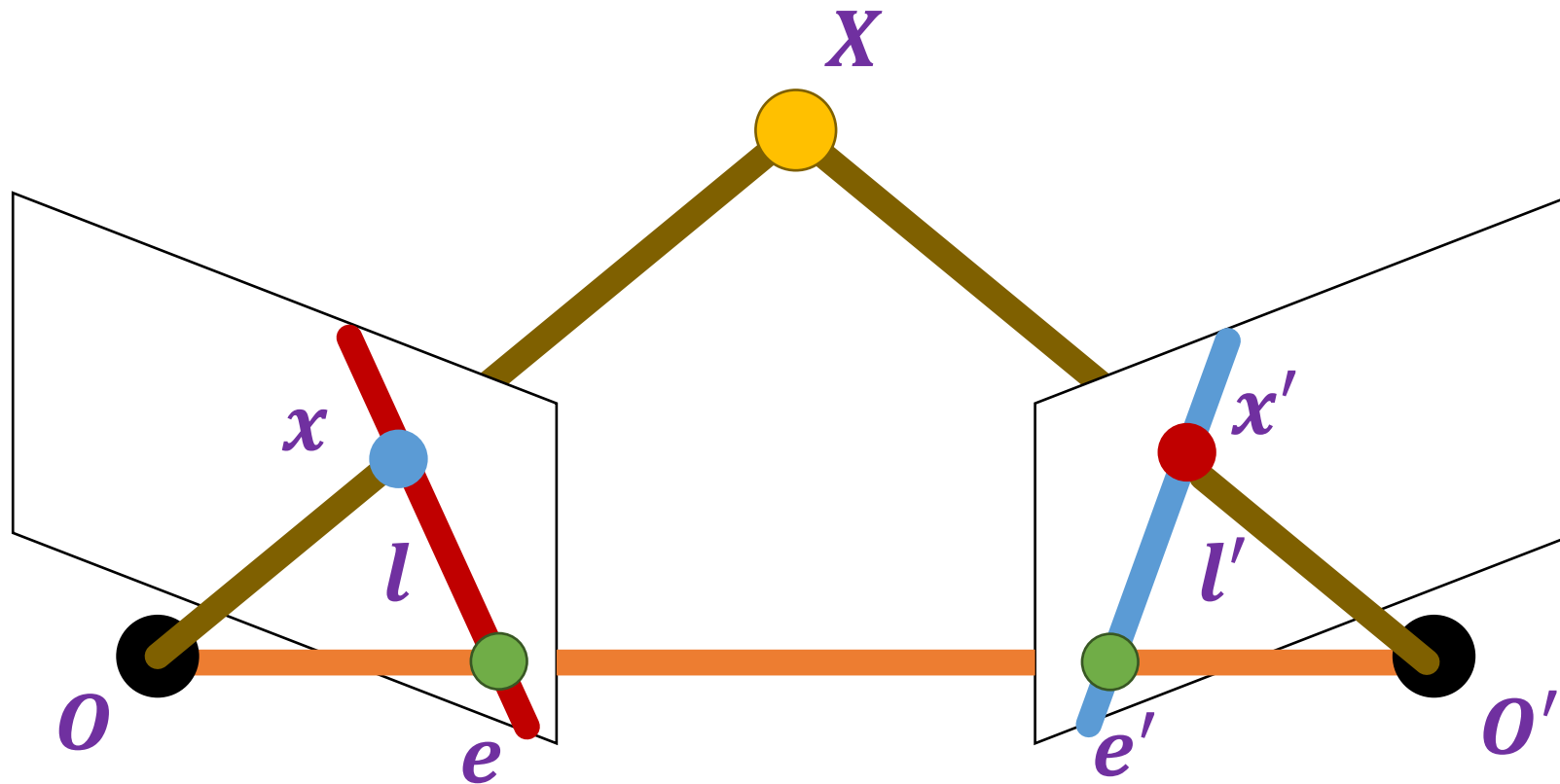
- Similarly, all points in the left image corresponding to x' have to lie along the epipolar line corresponding to x'

Epipolar constraint



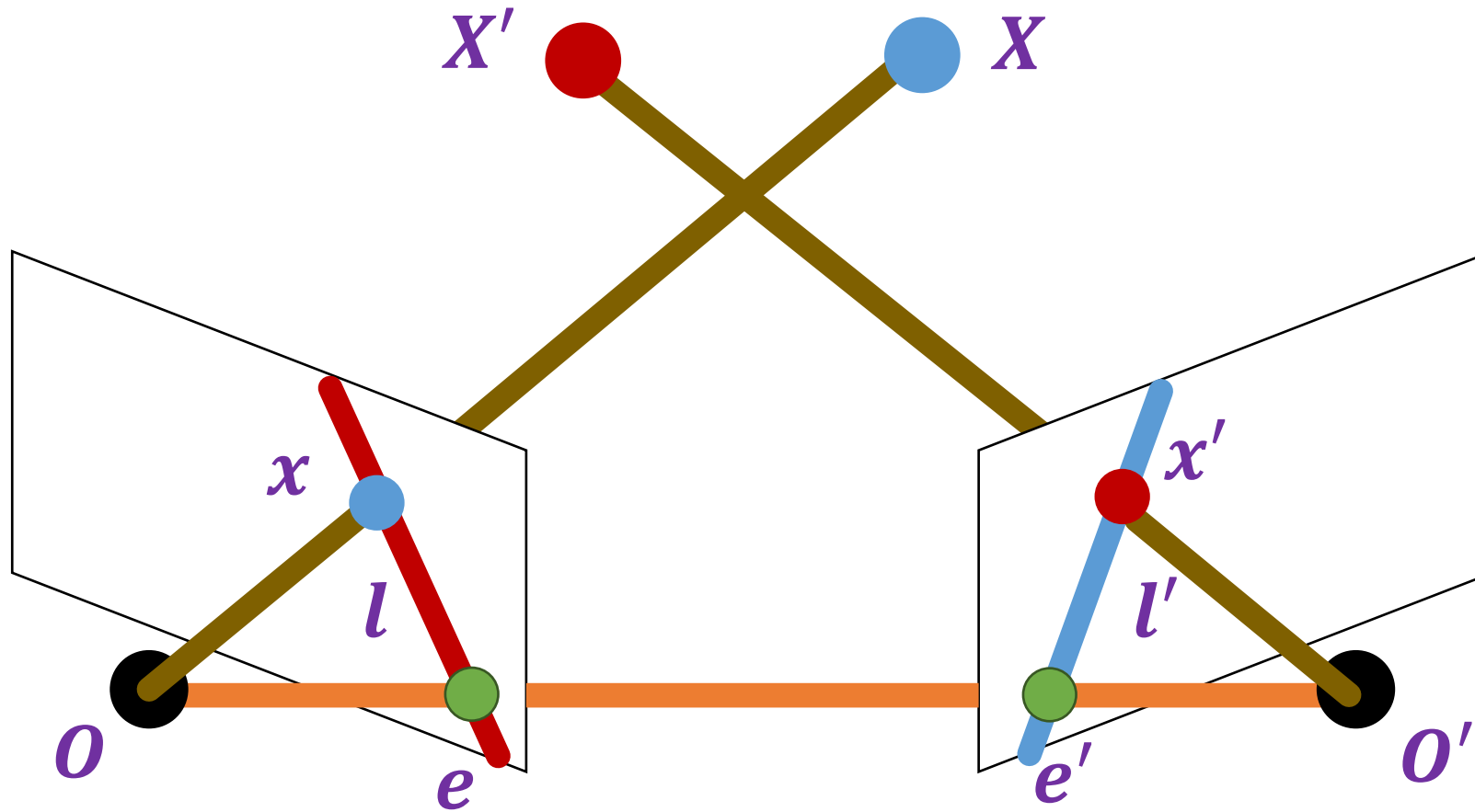
- Potential matches for x have to lie on the matching epipolar line l' and vice-versa \rightarrow need only to search along 1D epipolar line for matching!

Epipolar constraint



- Whenever two points x and x' lie on matching epipolar lines l and l' , the **visual rays** corresponding to them meet in space, i.e., x and x' **could be** projections of the same 3D point X

Epipolar constraint

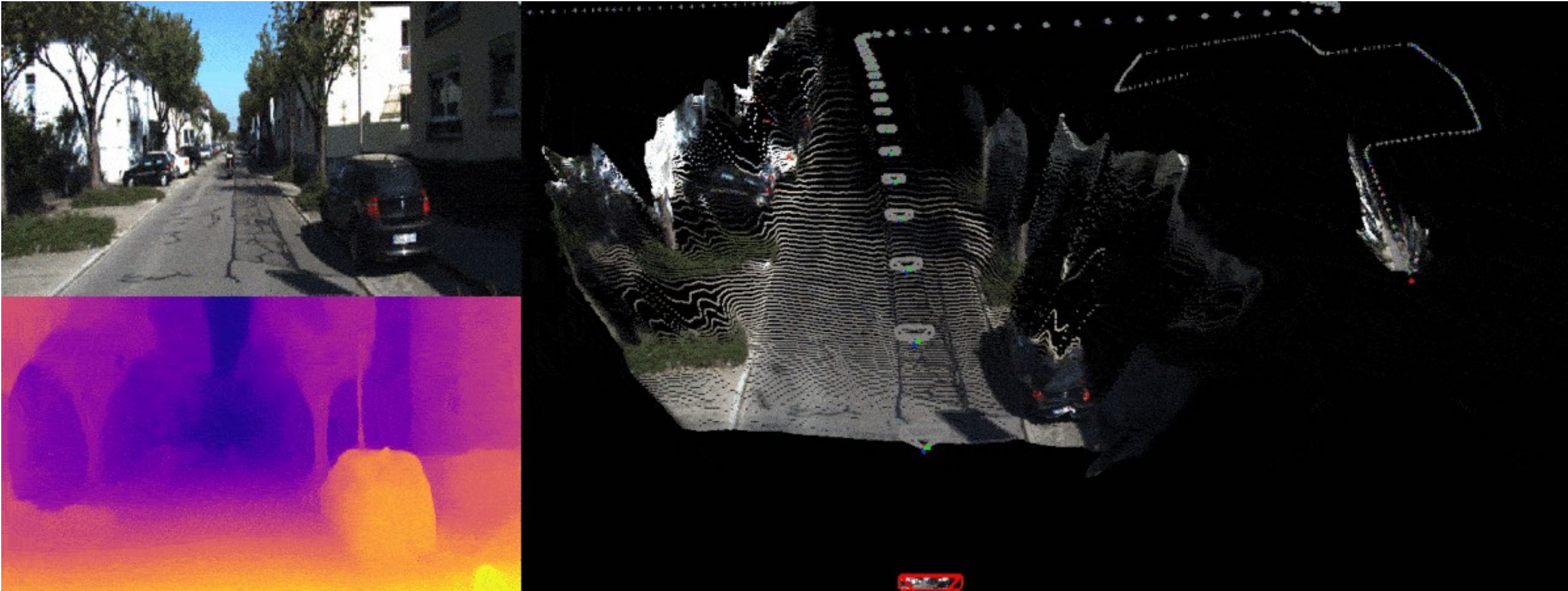


- Caveat: if x and x' satisfy the epipolar constraint, this doesn't mean they *have to be* projections of the same 3D point

Epipolar constraint: Example



Epipolar Geometry & Deep Learning



Multi-Frame Self-Supervised Depth Estimation with Transformers (CVPR 2022)

Vitor Guizilini, Rares Ambrus, Dian Chen, Sergey Zakharov, Adrien Gaidon

Epipolar Geometry & Deep Learning

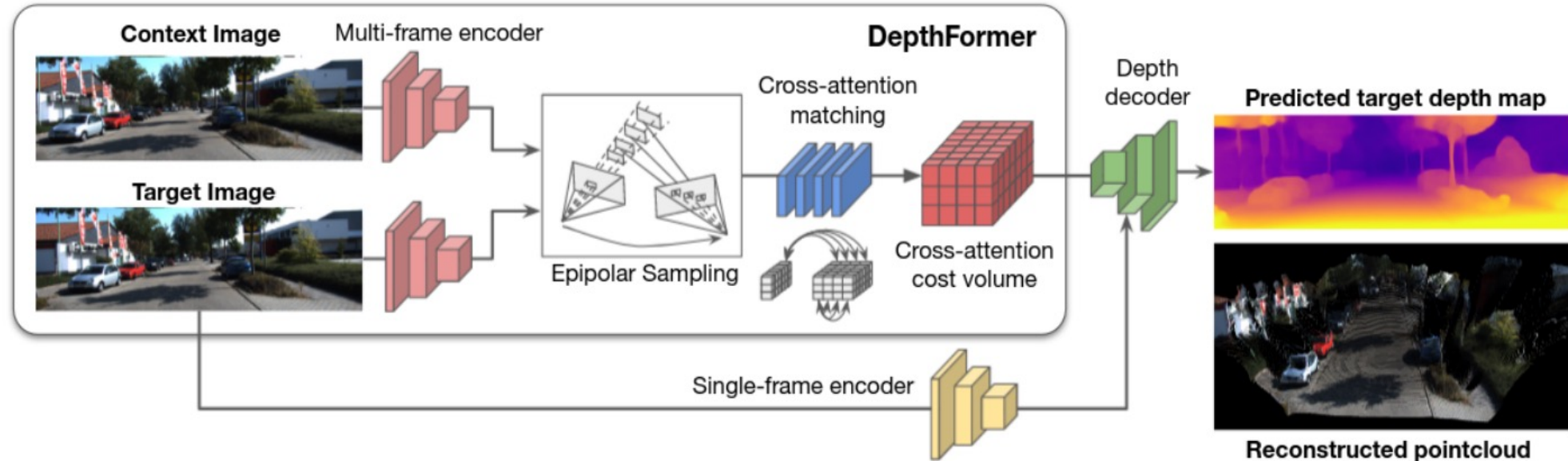
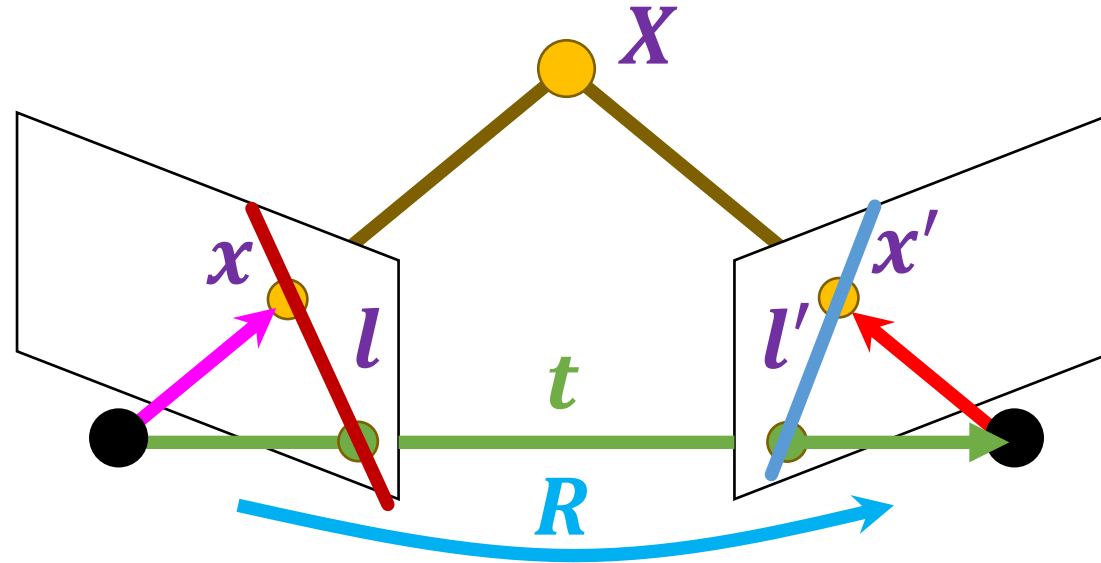


Figure 1. **Our DepthFormer architecture** achieves state-of-the-art multi-frame self-supervised monocular depth estimation by **improving feature matching** across images during cost volume generation.

Multi-Frame Self-Supervised Depth Estimation with Transformers (CVPR 2022)

Vitor Guizilini, Rares Ambrus, Dian Chen, Sergey Zakharov, Adrien Gaidon

The Epipolar Constraint as an Equation



$x'^T F x = 0$ where $F = K'^{-T} E K^{-1}$ is called the **Fundamental Matrix** [1]

and $E = [t]_{\times} R$ is the **Essential Matrix** [2]

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

(sketch of proof in appendix)

[1] [Faugeras et al., \(1992\)](#), [Hartley \(1992\)](#)

[2] H. C. Longuet-Higgins. [A computer algorithm for reconstructing a scene from two projections](#). Nature, 1981

Estimating the fundamental matrix - teaser

- Given: correspondences $\mathbf{x} = (x, y, 1)^T$ and $\mathbf{x}' = (x', y', 1)^T$



Estimating the fundamental matrix - teaser

- Given: 2D correspondences $\mathbf{x} = (x, y, 1)^T$ and $\mathbf{x}' = (x', y', 1)^T$
- Constraints: $\mathbf{x}'^T \mathbf{F} \mathbf{x} = 0$ (1 per correspondence, how many needed?)
- Boils down to another homogeneous linear equation $\mathbf{A} \mathbf{X} = \mathbf{0}$
- Recast once more into total least squares ([Sz.A.2.1](#)) due to noise
- SVD gives the solution as usual + enables enforcing rank 2 constraint by replacing smallest singular value with 0
- This “algebraic” algorithm is called “[normalized] 8-point algorithm” (R. Hartley. [In defense of the eight-point algorithm](#). TPAMI 1997)
- As in calibration and homography fitting: non-linear “geometric” optimization (of reprojected distances) is more precise
- Can be made robust to outliers via the RANSAC algorithm
- See appendix, [H&Z ch. 9](#), [Szeliski](#) 11.3, or take **CS231A** for more!

From epipolar geometry to camera calibration

Estimating the fundamental matrix is known as “weak calibration”

If we know the calibration matrices (K , K') of the two cameras, we can estimate the *essential matrix*: $E = K'^T F K$

The essential matrix gives us the relative rotation and translation between the cameras, or their *extrinsic parameters* ($E = [t]_{\times} R$)

Alternatively, if the calibration matrices are known (or in practice, if good initial guesses of the intrinsics are available), the five-point algorithm can be used to estimate relative camera pose

What will we learn today?

Triangulation

Epipolar geometry

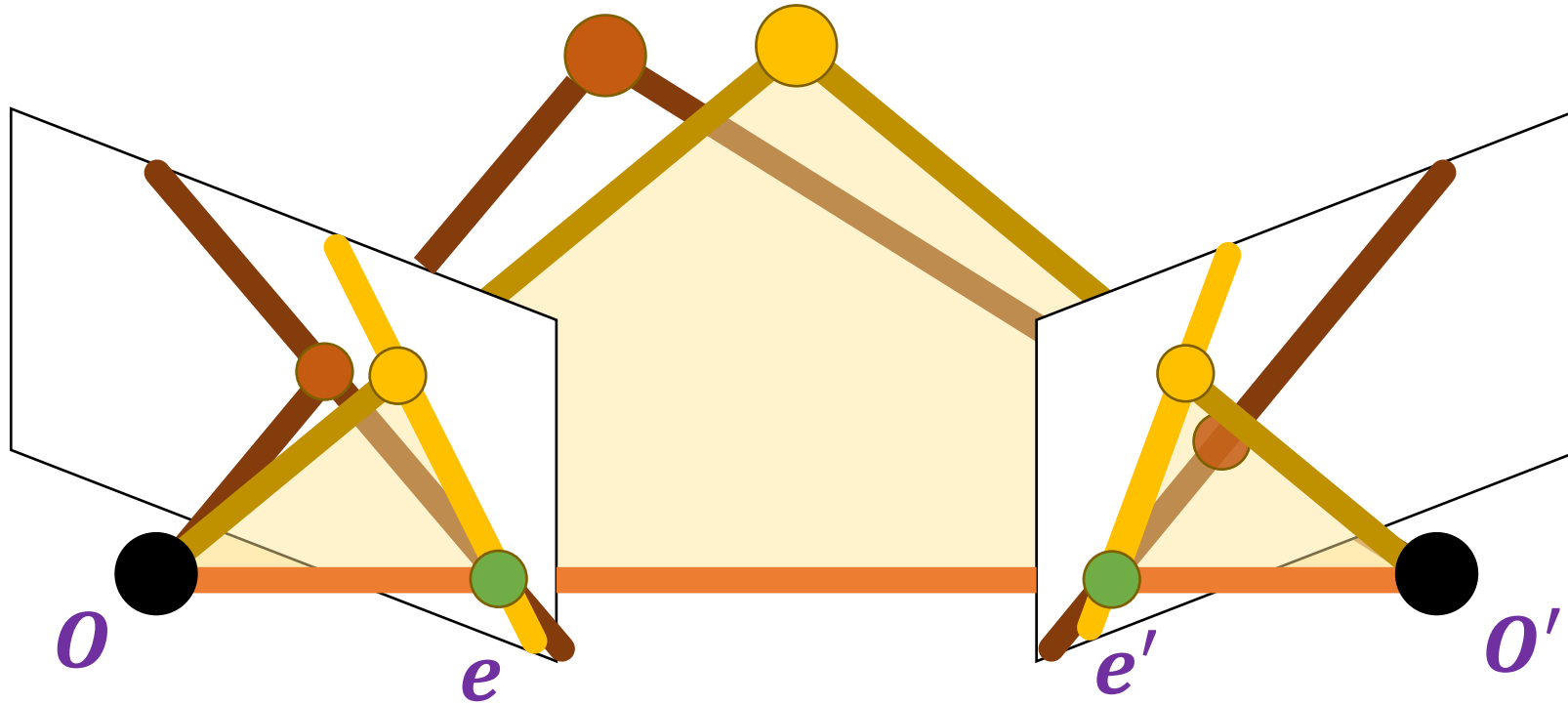
Stereo

Structure-from-Motion (SfM)

Example configuration: Converging cameras

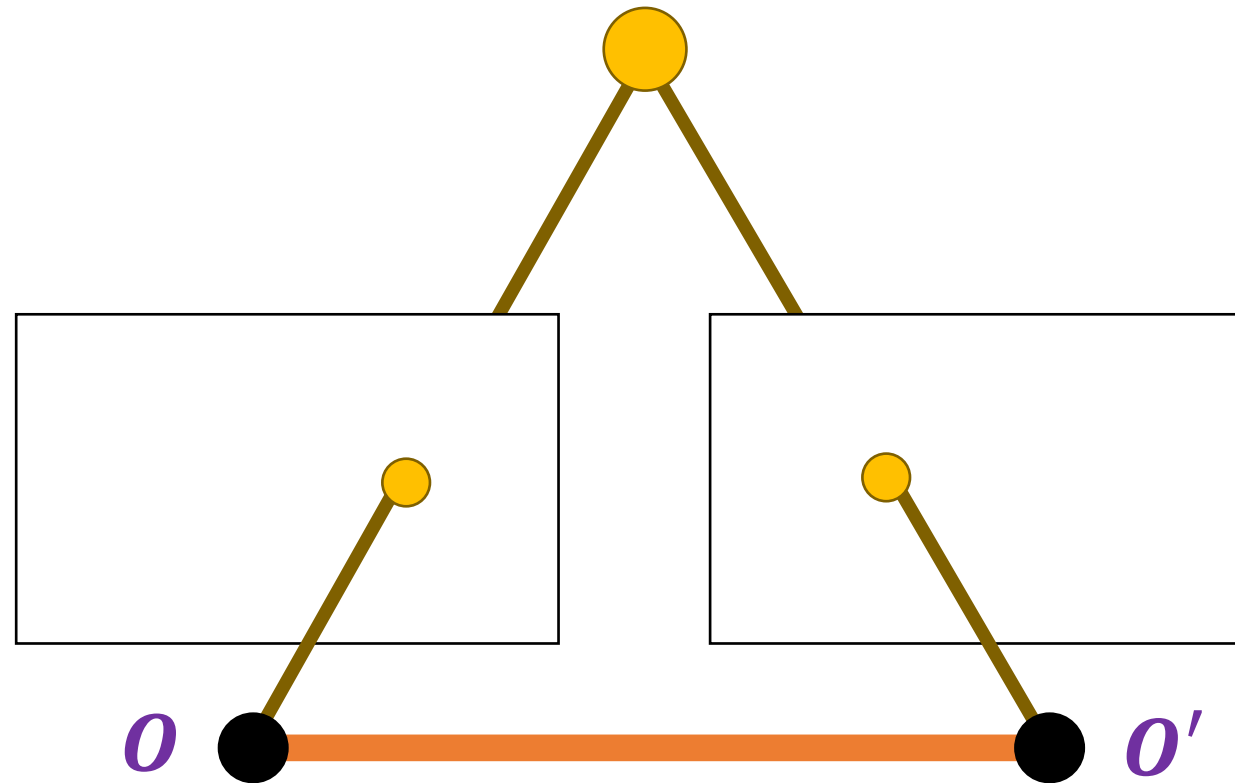


Example configuration: Converging cameras



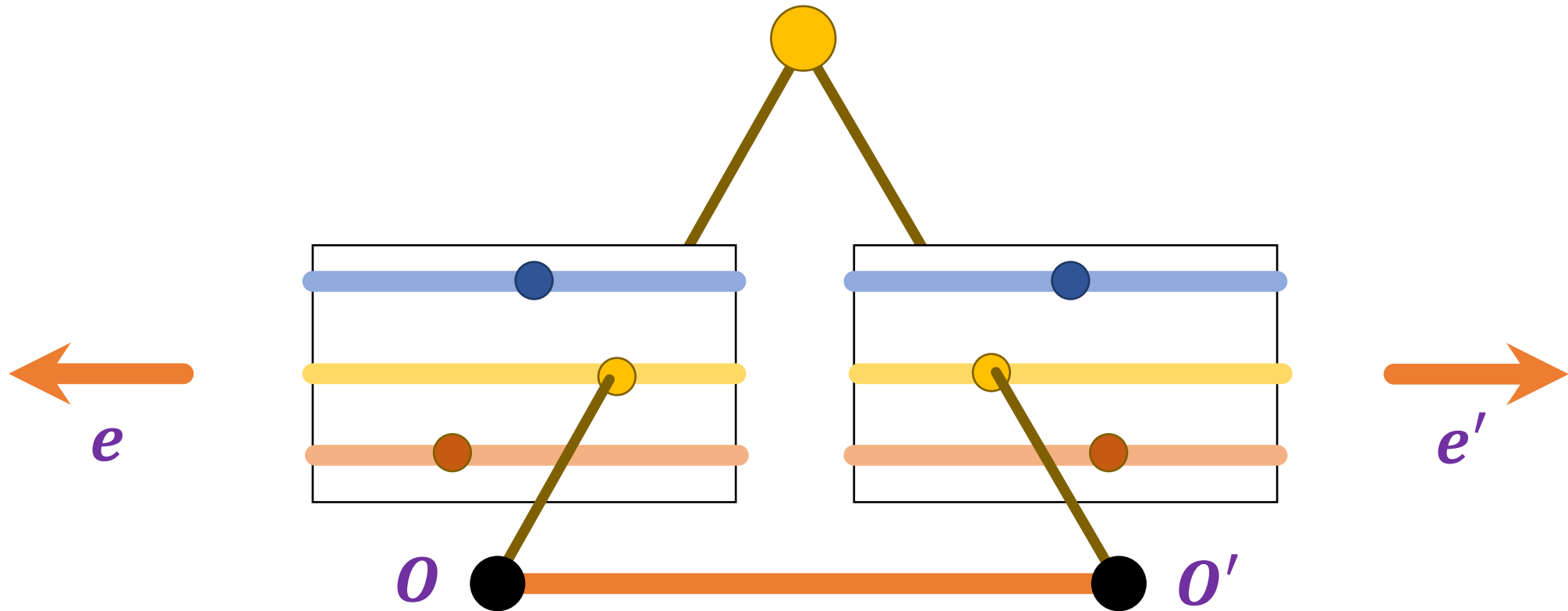
- Epipoles are finite, may be visible in the image

Example configuration: Motion parallel to image plane



Where are the epipoles?
What do the epipolar lines look like?

Example configuration: Motion parallel to image plane



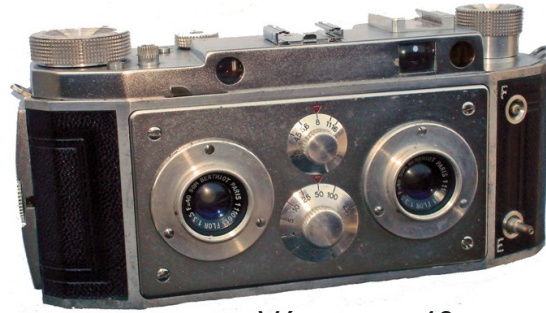
Epipoles *infinitely* far away!

Epipolar lines parallel: "scan lines"

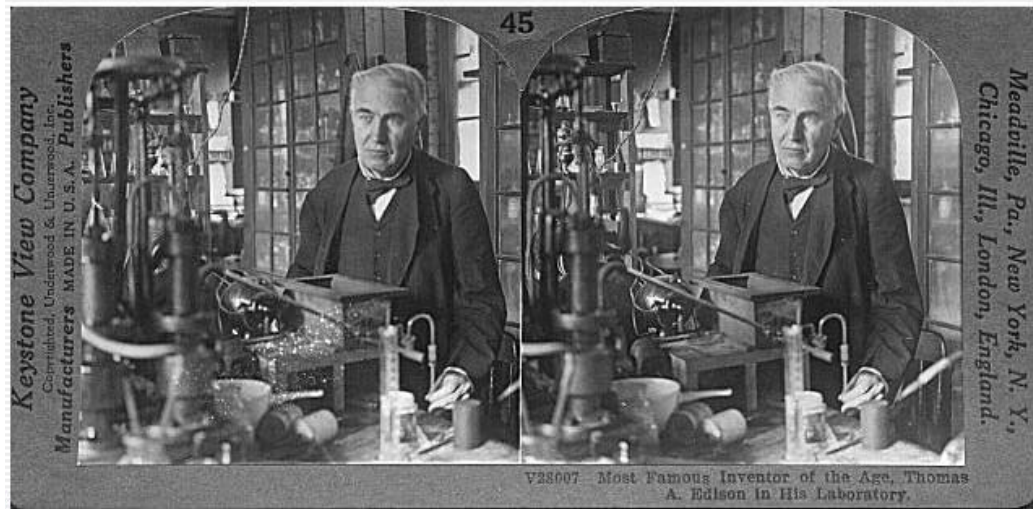
→ Stereo = easier fronto-parallel special case!

History: Stereograms

Humans can fuse pairs of images to get a sensation of depth



Vérascopie 40

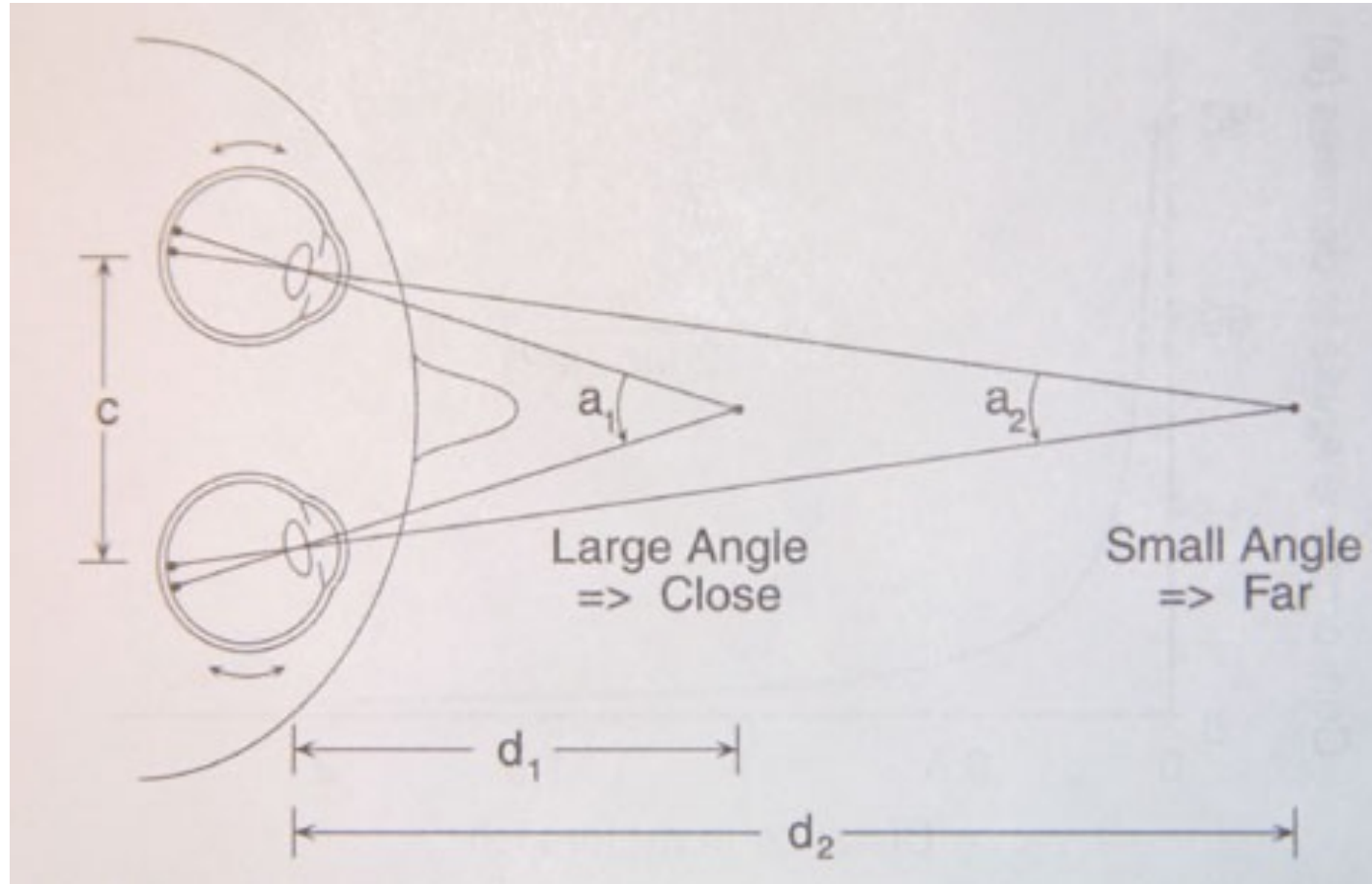


Stereograms: Invented by Sir Charles Wheatstone, 1838



<https://en.wikipedia.org/wiki/Stereoscopy>

Depth from convergence



$$d = \frac{c}{2 \tan(a/2)}$$

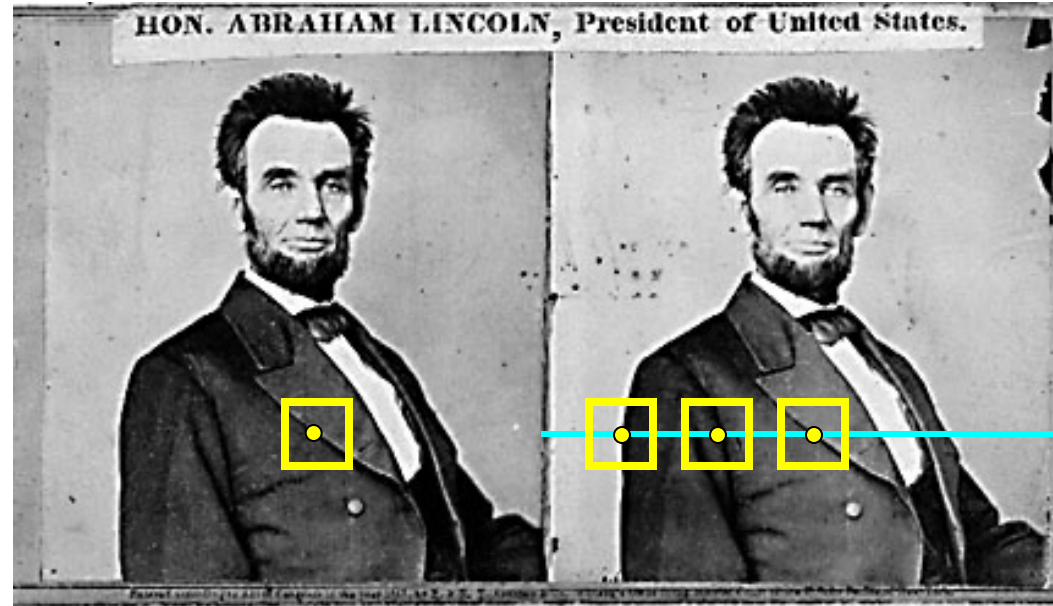
Stereo Matching for Depth Estimation

Given: stereo pair (assumed calibrated)

Wanted: dense depth map



Basic stereo matching algorithm



For each pixel in the first image

Find corresponding epipolar line in the right image: same row!

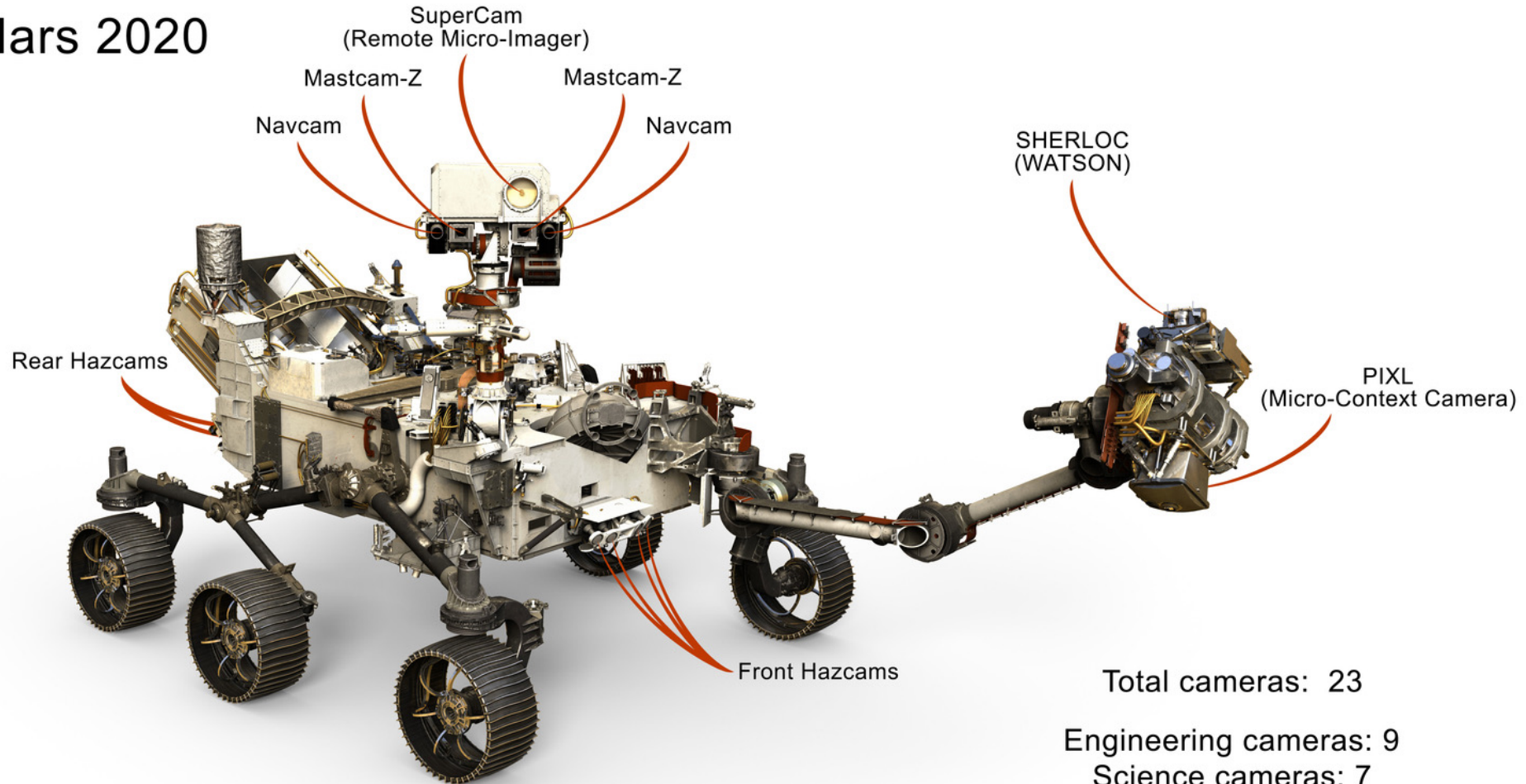
Examine all pixels on the epipolar line and pick the best match

Triangulate the matches to get depth information

More details in appendix: rectification, matching, depth from disparity, etc

Stereo on the Perseverance Mars Rover

Mars 2020



Total cameras: 23

Engineering cameras: 9

Science cameras: 7

Entry, descent and landing cameras: 7

What will we learn today?

Triangulation

Epipolar geometry

Stereo

Structure-from-Motion (SfM)

Reference: [Szeliski](#) 11, [H&Z ch. 9](#)

Most slides adapted from N. Snavely & S. Lazebnik

Structure-from-Motion

Given many images, how can we

- a) figure out where they were all taken from?
- b) build a 3D model of the scene?



N. Snavely, S. Seitz, and R. Szeliski, [Photo tourism: Exploring photo collections in 3D](http://phototour.cs.washington.edu/), SIGGRAPH 2006.
<http://phototour.cs.washington.edu/>

Geometry of more than two views?

2 views: governed by the 3×3 *Fundamental Matrix* (how to go from one point in an image to the epipolar line in the 2nd image)

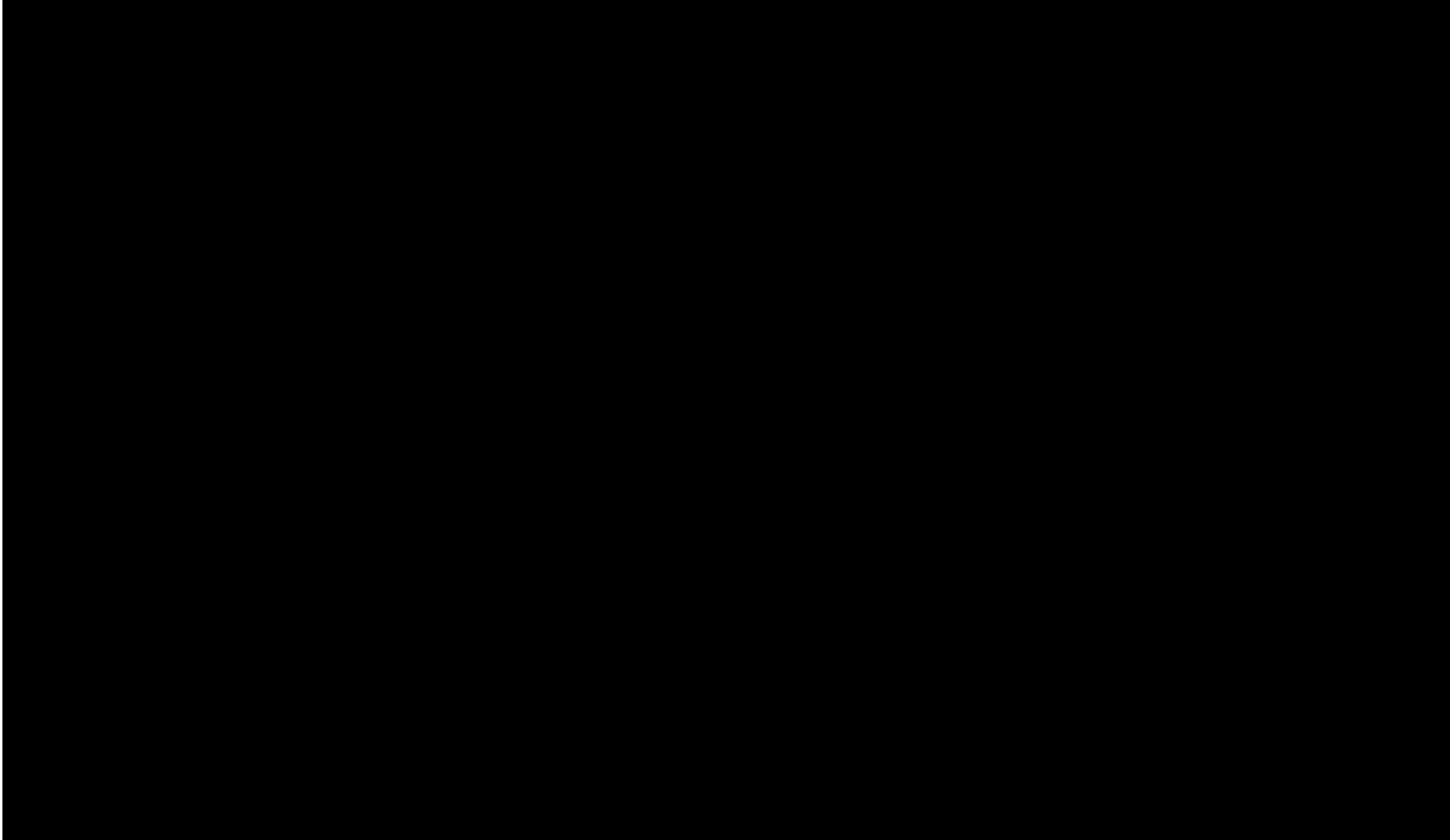
3 views: governed by the $3 \times 3 \times 3$ *Trifocal Tensor*

4 views: governed by the $3 \times 3 \times 3 \times 3$ *Quadrifocal Tensor*

After this it starts to get complicated...

→ *explicitly* solve for camera poses *and* scene geometry

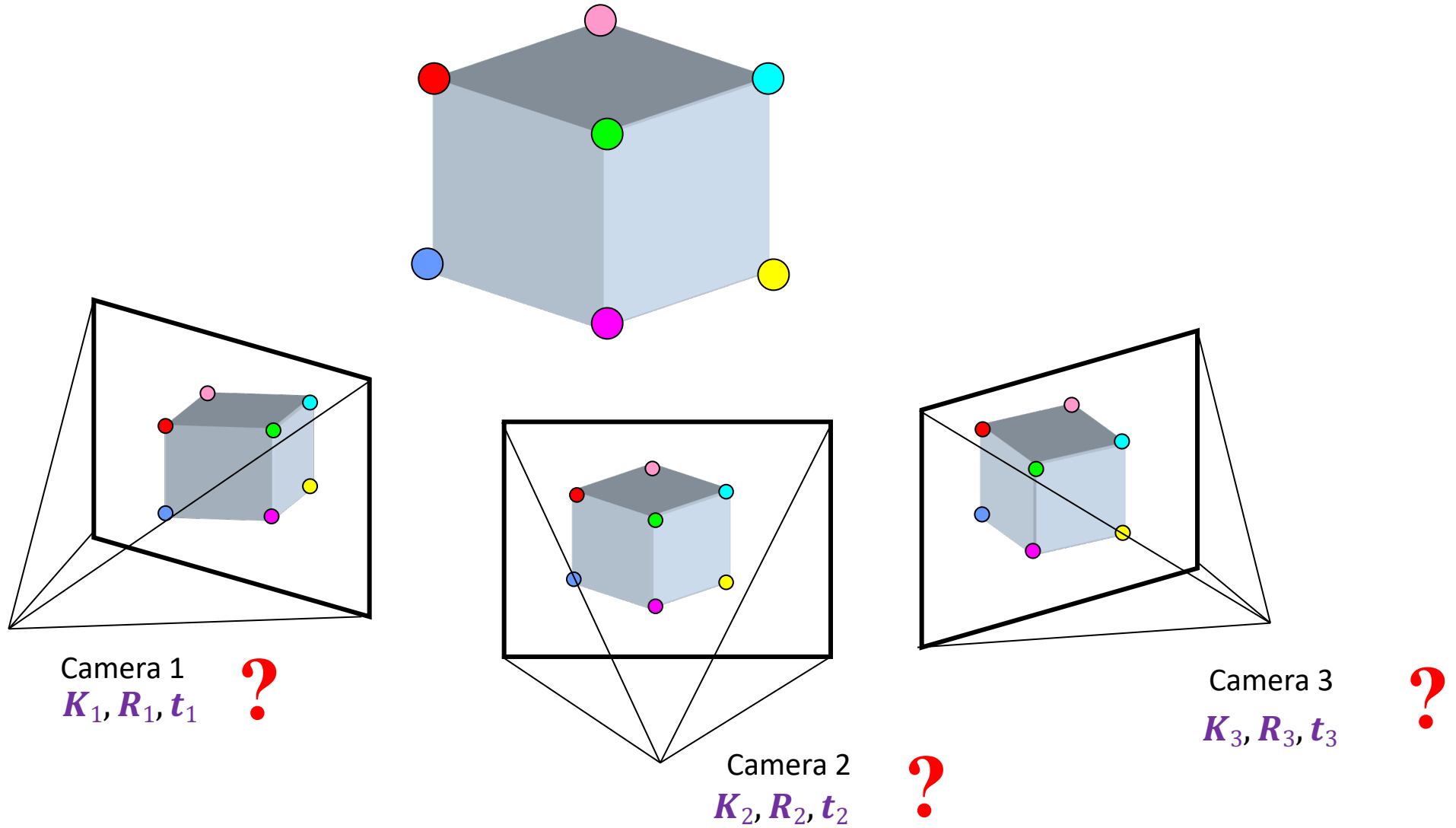
Large-scale structure-from-motion



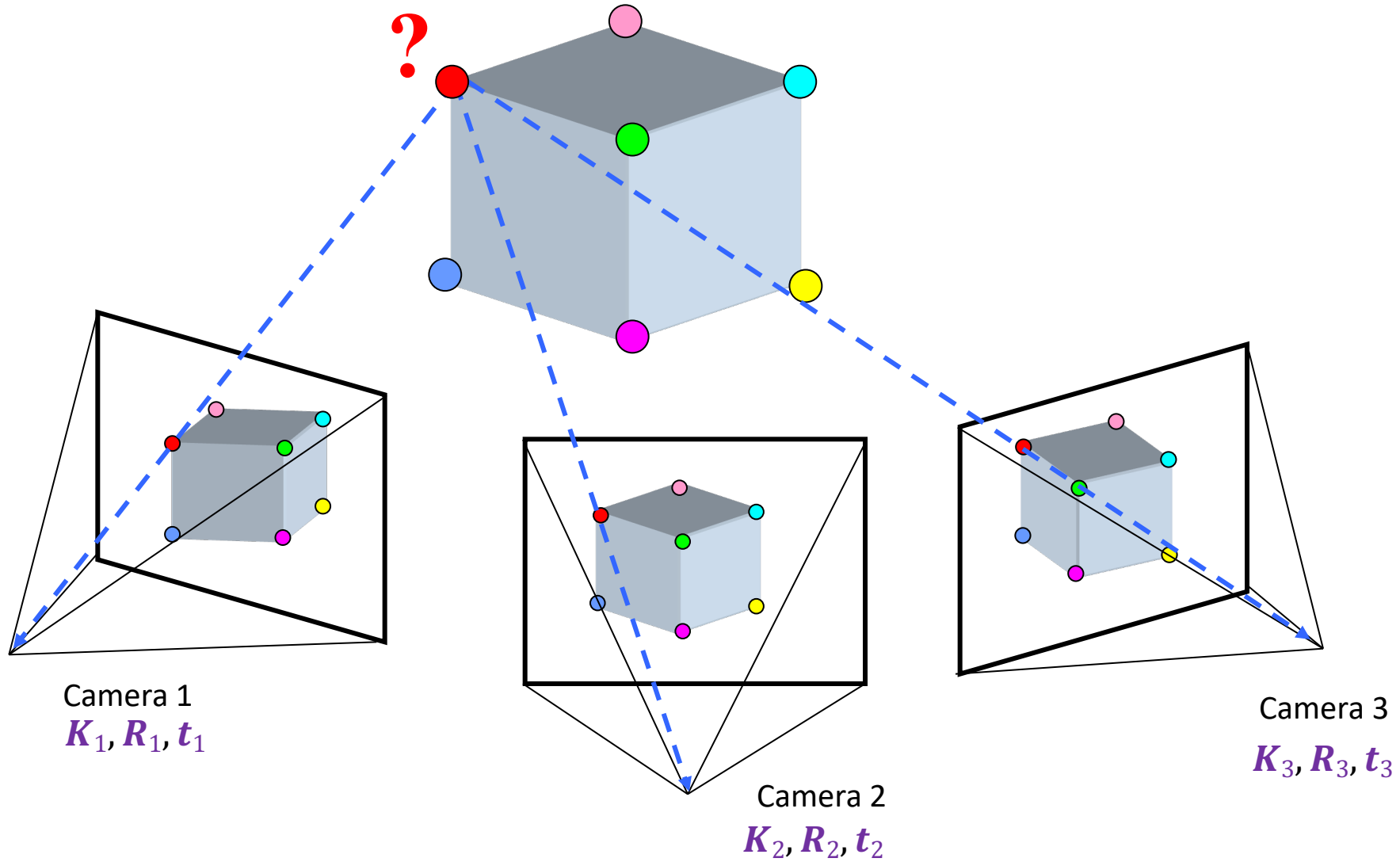
Dubrovnik, Croatia. 4,619 images (out of an initial 57,845 downloaded from Flickr). 3.5M points!
Total reconstruction time: 17.5 hours on 352 cores

Building Rome in a Day, Agarwal *et al*, ICCV'09
<http://grail.cs.washington.edu/rome/>

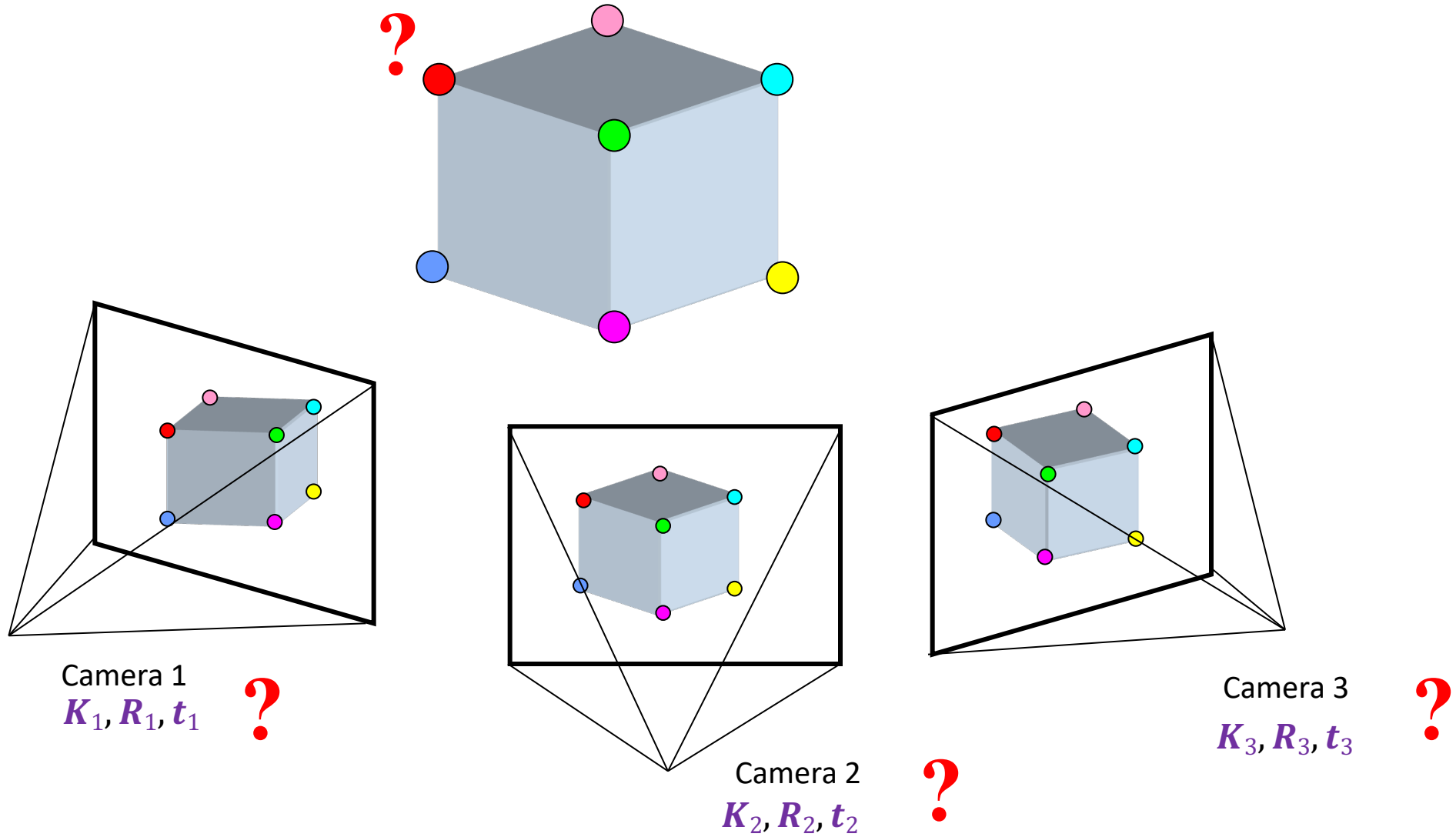
Recall: Calibration



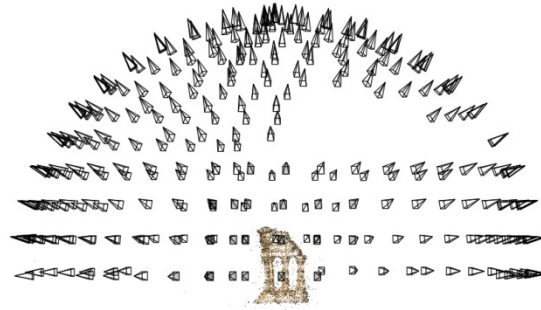
Recall: Triangulation / Multi-view Stereo



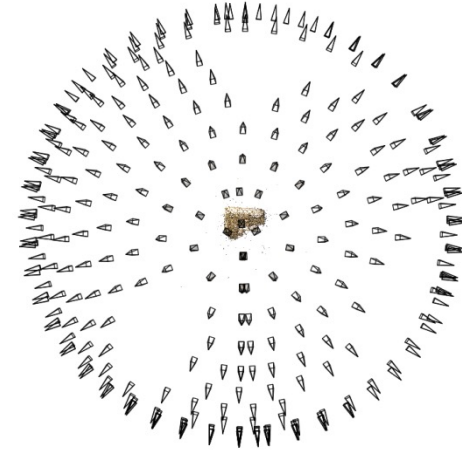
Structure-from-Motion



Structure-from-Motion



Reconstruction (side)

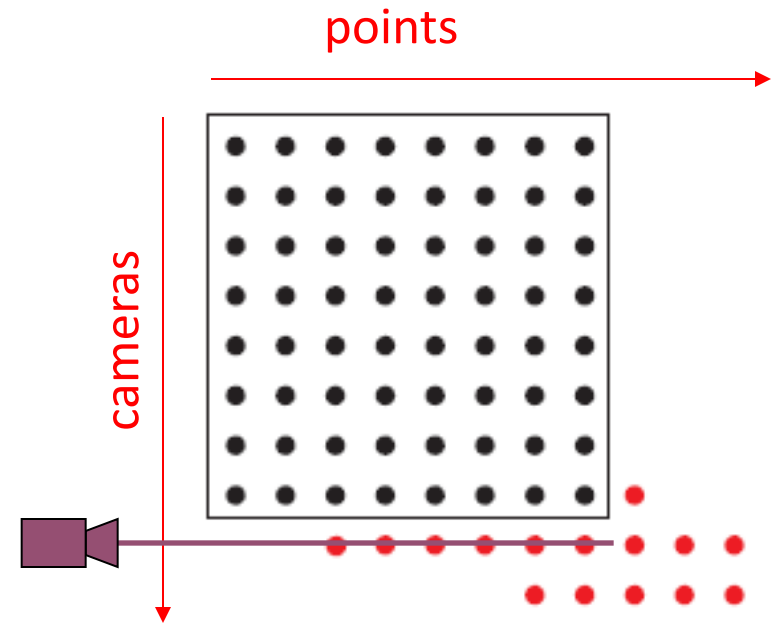


(top)

- Input: images with 2D points \mathbf{x}_{ij} in correspondence
- Output (solved simultaneously now!)
 - structure: 3D location \mathbf{X}_j for each point \mathbf{x}_{ij}
 - motion: camera parameters $\mathbf{R}_i, \mathbf{t}_i$ & possibly \mathbf{K}_i
- Objective function: minimize *reprojection error in 2D*

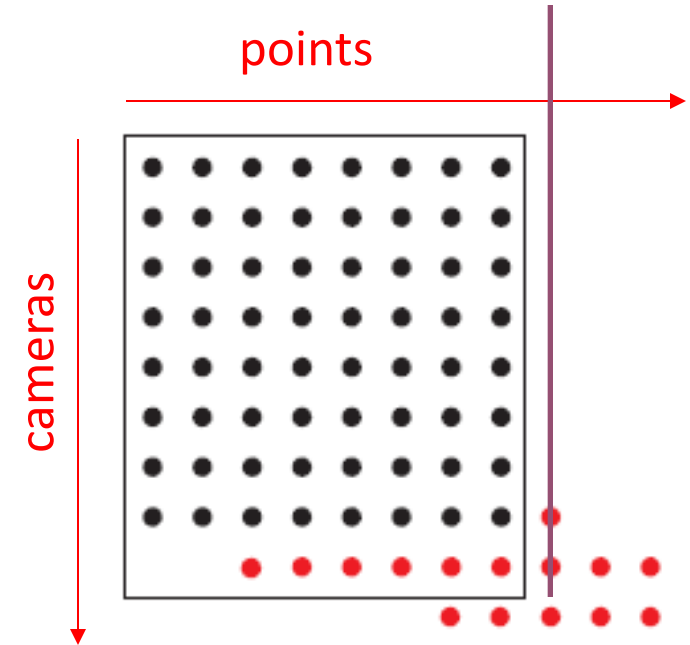
Incremental Structure-from-Motion

- Initialize motion from two images using the fundamental matrix
- Initialize structure by triangulation
- For each additional view:
 - Determine projection matrix of new camera using all the known 3D points that are visible in its image – **calibration**



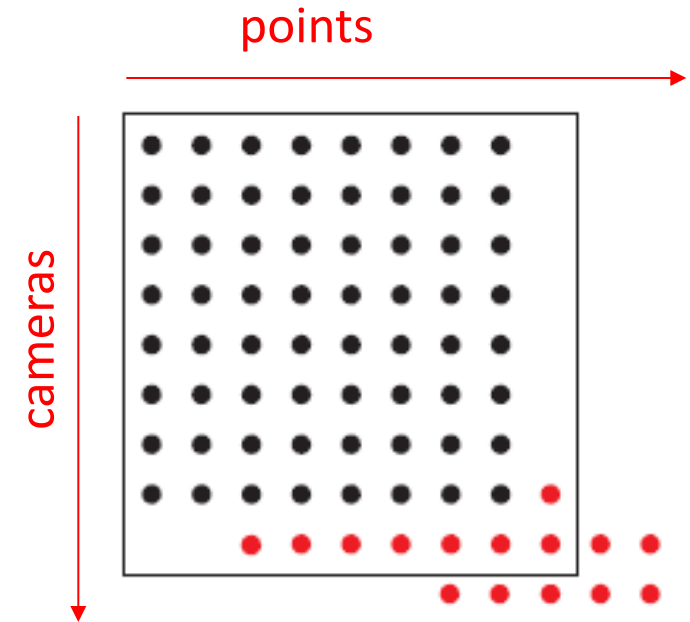
Incremental Structure-from-Motion

- Initialize motion from two images using the fundamental matrix
- Initialize structure by triangulation
- For each additional view:
 - Determine projection matrix of new camera using all the known 3D points that are visible in its image – **calibration**
 - Refine and extend structure: compute newly visible 3D points, re-optimize existing points that are also seen by this camera – **triangulation**



Incremental Structure-from-Motion

- Initialize motion from two images using the fundamental matrix
- Initialize structure by triangulation
- For each additional view:
 - Determine projection matrix of new camera using all the known 3D points that are visible in its image – **calibration**
 - Refine and extend structure: compute newly visible 3D points, re-optimize existing points that are also seen by this camera – **triangulation**



- Refine all cameras & points *jointly*: **bundle adjustment**

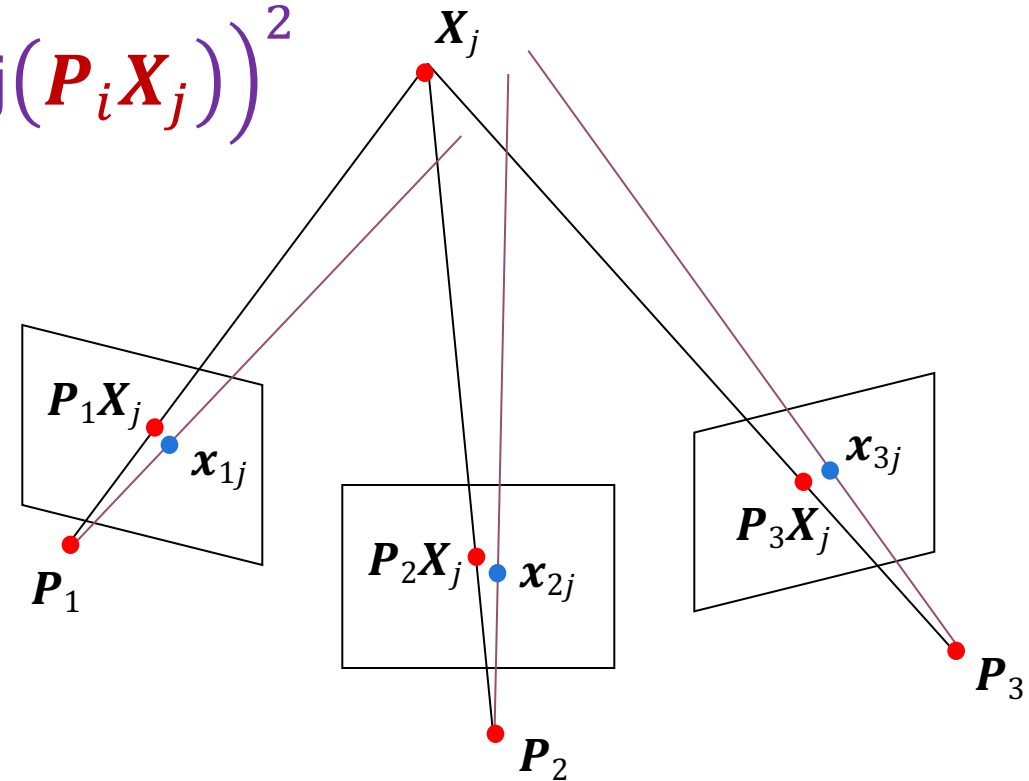
Bundle Adjustment

Non-linear method for refining structure (\mathbf{X}_j) and motion (\mathbf{P}_i)

Minimize reprojection error (with lots of bells and whistles):

$$\sum_{i=1}^m \sum_{j=1}^n w_{ij} d(\mathbf{x}_{ij} - \text{proj}(\mathbf{P}_i \mathbf{X}_j))^2$$

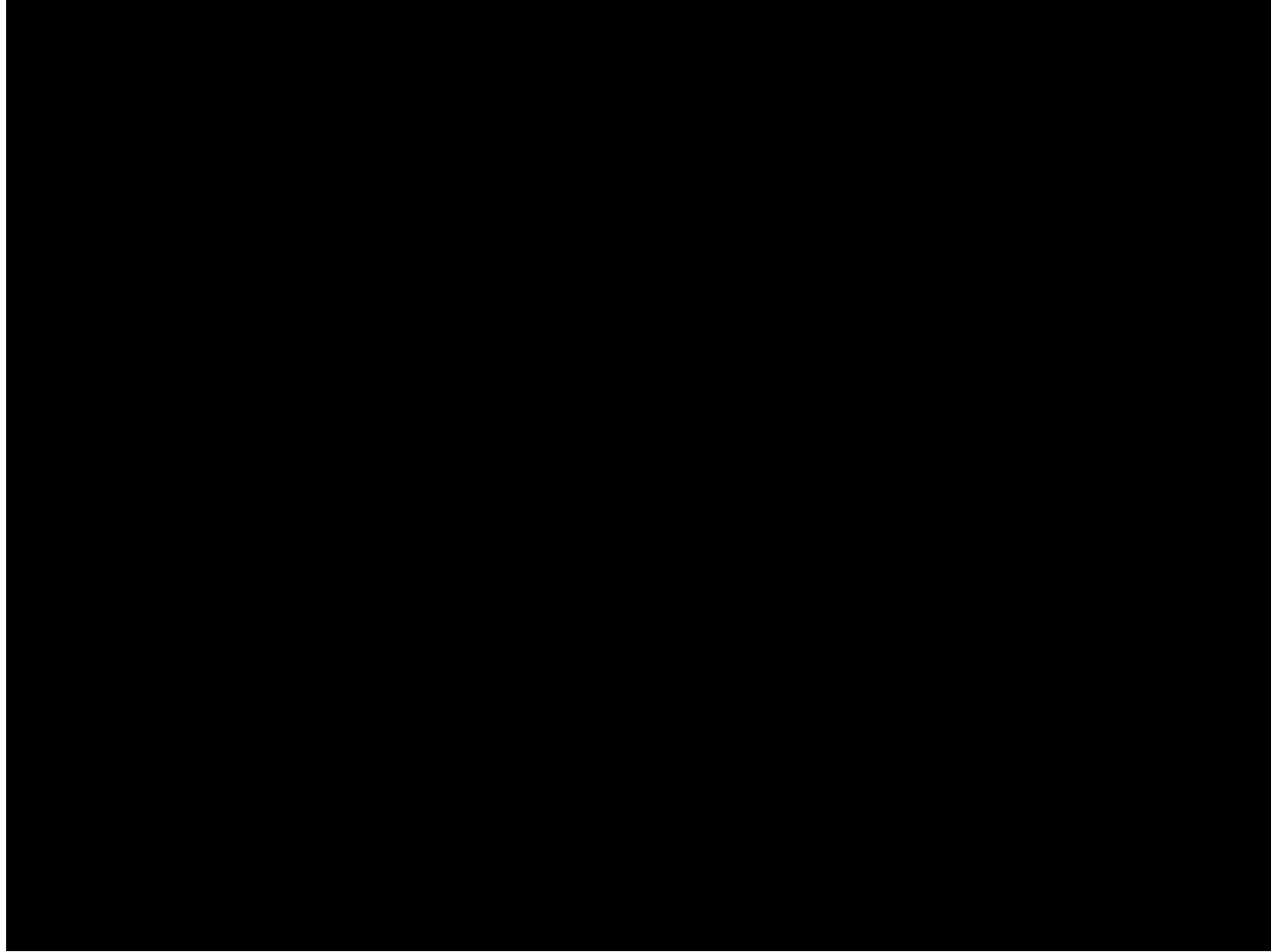
visibility flag: is
point j visible in
view i ?



Incremental SfM in Practice

- Pick a pair of images with lots of inliers (and good EXIF data)
 - Initialize intrinsic parameters (focal length, principal point) from EXIF
 - Estimate extrinsic parameters (R and t) using [five-point algorithm](#)
 - Use triangulation to initialize model points
- While remaining images exist
 - Find an image with many feature matches with images in the model
 - Run RANSAC on feature matches to register new image to model
 - Triangulate new points
 - Perform bundle adjustment to re-optimize everything
 - Optionally, align with GPS from EXIF data or ground control points

Incremental structure from motion



Time-lapse reconstruction of Dubrovnik, Croatia, viewed from above

COLMAP

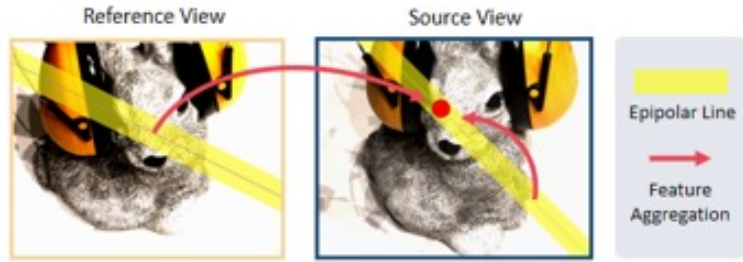


Sparse model of central Rome using 21K photos produced by COLMAP's SfM pipeline.

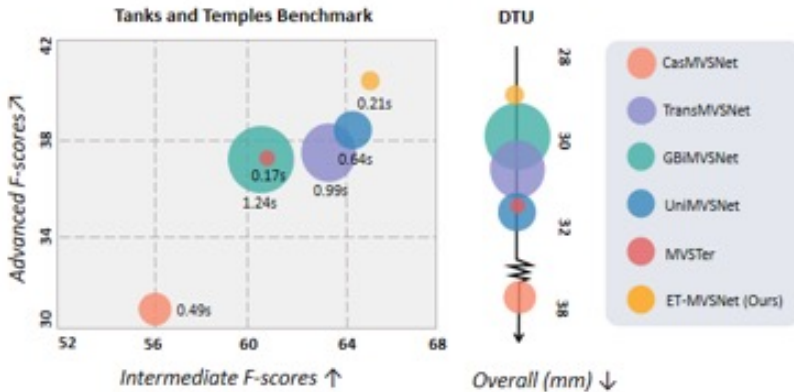


Dense models of several landmarks produced by COLMAP's MVS pipeline.

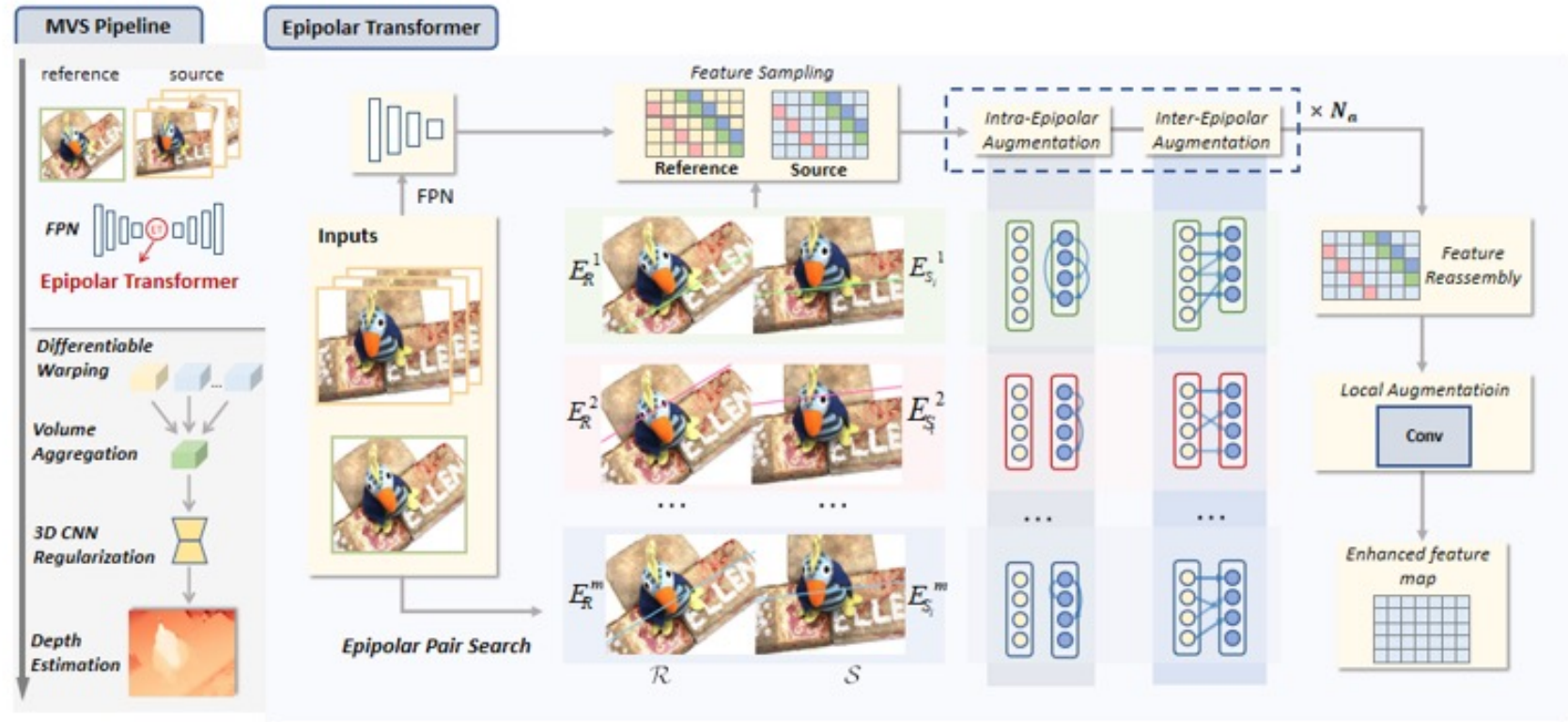
SfM in the age of Deep Learning



(a) Non-Local feature aggregation on epipolar lines



(b) Performance and efficiency comparison with other MVS approaches



[ET-MVSNet: When Epipolar Constraint Meets Non-local Operators in Multi-View Stereo \(ICCV'23\)](#)

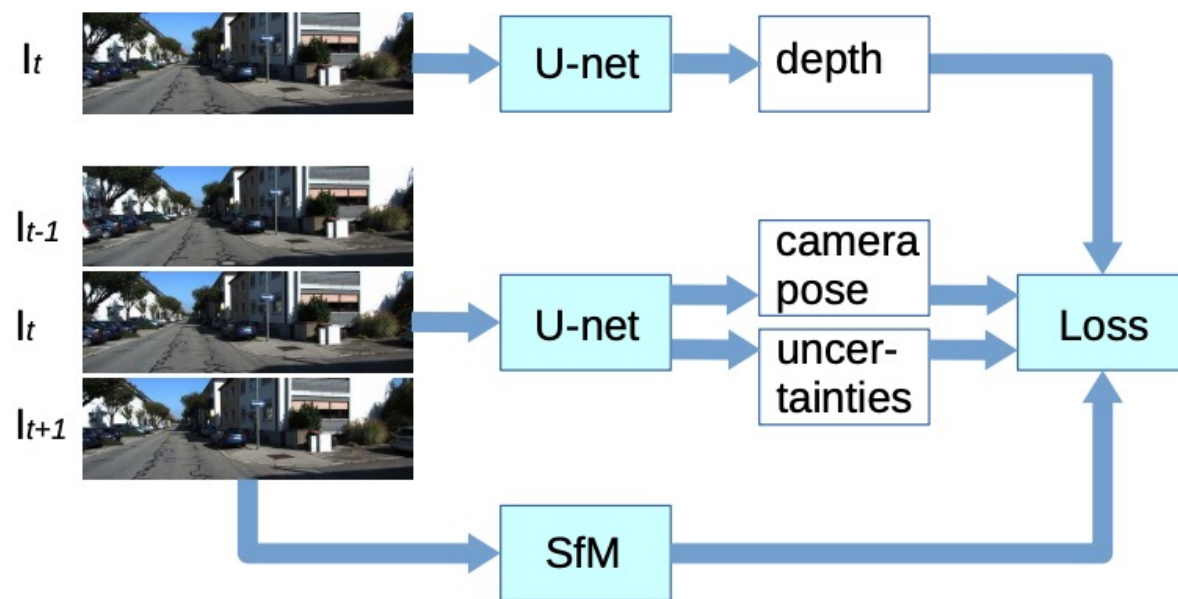
See also [MVSFormer: Multi-View Stereo by Learning Robust Image Features and Temperature-based Depth \(TMLR'23\)](#)

Supervising the new with the old: learning SFM from SFM

Maria Klodt^[0000-0003-3015-9584] and Andrea Vedaldi^[0000-0003-1374-2858]

Visual Geometry Group, University of Oxford
{klodt,vedaldi}@robots.ox.ac.uk

Abstract. Recent work has demonstrated that it is possible to learn deep neural networks for monocular depth and ego-motion estimation from unlabelled video sequences, an interesting theoretical development with numerous advantages in applications. In this paper, we propose a number of improvements to these approaches. First, since such self-supervised approaches are based on the brightness constancy assumption, which is valid only for a subset of pixels, we propose a probabilistic learning formulation where the network predicts distributions over variables rather than specific values. As these distributions are conditioned on the observed image, the network can learn which scene and object types are likely to violate the model assumptions, resulting in more robust learning. We also propose to build on dozens of years of experience in developing handcrafted structure-from-motion (SfM) algorithms. We do so by using an off-the-shelf SfM system to generate a supervisory signal for the deep neural network. While this signal is also noisy, we show that our probabilistic formulation can learn and account for the defects of SfM, helping to integrate different sources of information and boosting the overall performance of the network.



(b) proposed network architecture:
the depth and pose-uncertainty networks
are supervised by traditional SfM.

What did we learn today?

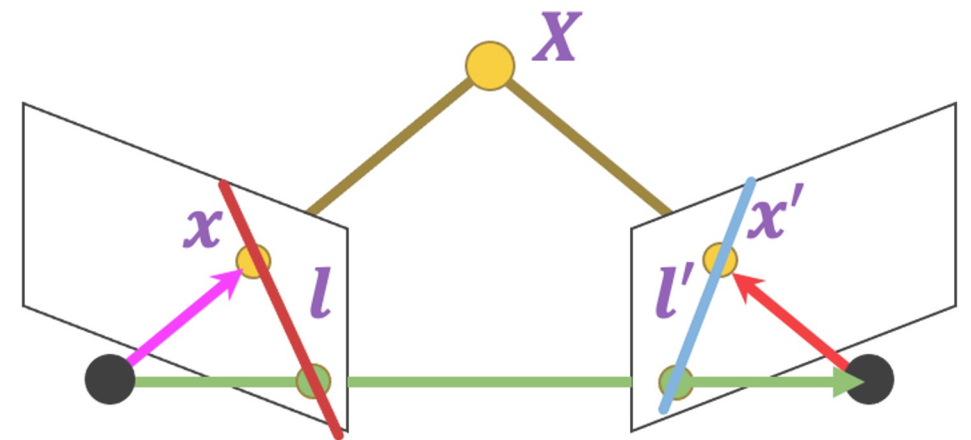
Triangulation

from calibrated cameras and 2D correspondences to 3D points

Epipolar geometry

Epipolar constraint,
Essential & Fundamental matrices

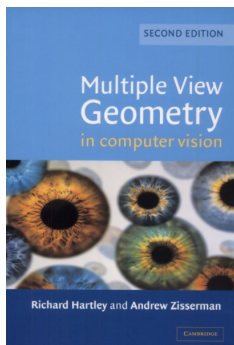
$$x'^T F x = 0$$



Stereo (fronto-parallel special case)

Structure-from-Motion (SfM) (many images / video)

Note: this is just an introduction → for more, take CS231a or get into HZ



Wrapping up Geometric Vision

Homogeneous Coordinates & Projective Space

2D & 3D Transforms as Matrix Multiplication

Pinhole Camera Model $P=K[R|t]$

Calibration from known 3D-to-2D correspondences

Multi-view geom: fundamental matrix, stereo, SfM

Appendix

Triangulation: Linear approach

$$\begin{array}{lll} \lambda \mathbf{x} = \mathbf{P}\mathbf{X} & \mathbf{x} \times \mathbf{P}\mathbf{X} = \mathbf{0} & [\mathbf{x}]_{\times} \mathbf{P}\mathbf{X} = \mathbf{0} \\ \lambda' \mathbf{x}' = \mathbf{P}'\mathbf{X} & \mathbf{x}' \times \mathbf{P}'\mathbf{X} = \mathbf{0} & [\mathbf{x}']_{\times} \mathbf{P}'\mathbf{X} = \mathbf{0} \end{array}$$

Rewrite cross product as matrix multiplication:

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

2 independent equations per 2D point --> 4 equations in matrix form

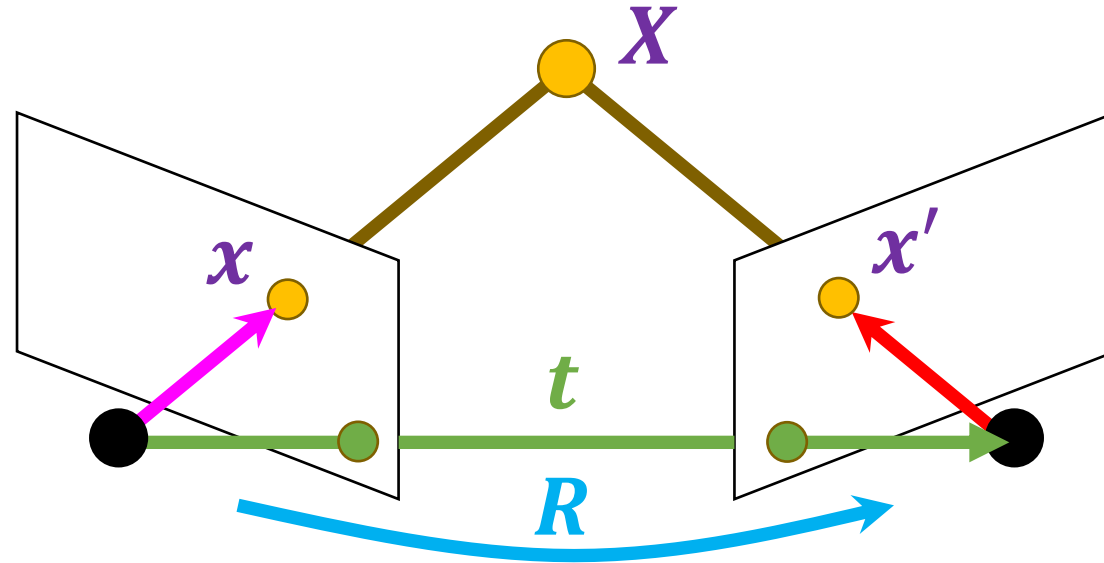
$$\mathbf{P} = \begin{bmatrix} \mathbf{p}_1^{\top} \\ \mathbf{p}_2^{\top} \\ \mathbf{p}_3^{\top} \end{bmatrix}, \mathbf{A}\mathbf{X} = \begin{bmatrix} y\mathbf{p}_3^{\top} - \mathbf{p}_2^{\top} \\ \mathbf{p}_1^{\top} - x\mathbf{p}_3^{\top} \\ y'\mathbf{p}'_3{}^{\top} - \mathbf{p}'_2{}^{\top} \\ \mathbf{p}'_1{}^{\top} - x'\mathbf{p}'_3{}^{\top} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \\ W \end{pmatrix} = \mathbf{0}$$

Total Least Squares:
 $\min_{\mathbf{X}} \|\mathbf{A}\mathbf{X}\|^2 \text{ s.t. } \|\mathbf{X}\|^2 = 1$
 Solved via SVD! (Sz. A.2.1)

Bonus slides - Math of the Epipolar Constraint

Slides credit: S. Lazebnik

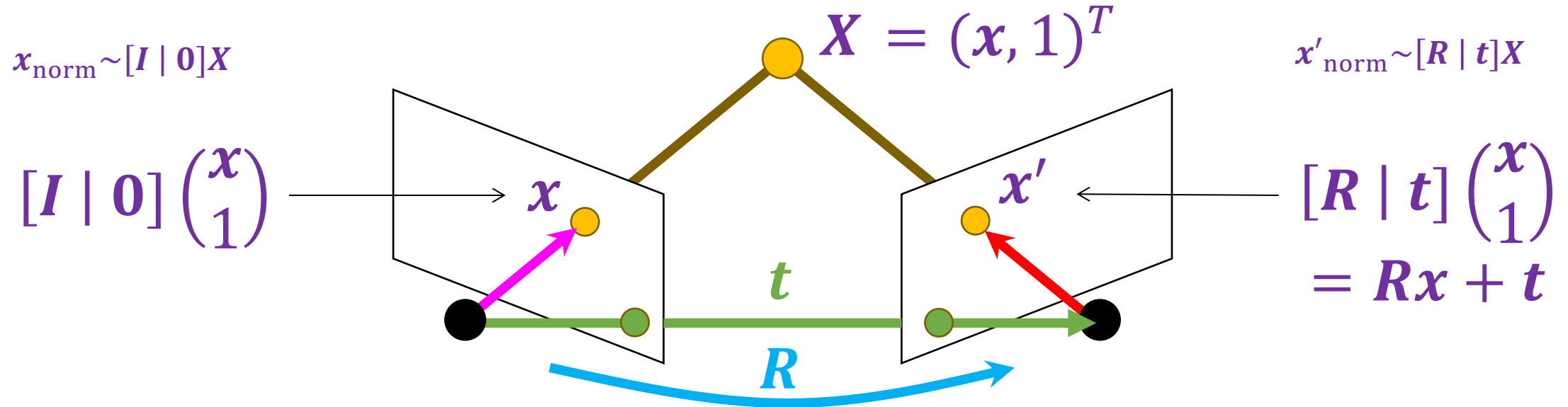
Math of the epipolar constraint: Calibrated case



- Assume the intrinsic and extrinsic parameters of the cameras are known, world coordinate system is set to that of the first camera
- Then the projection matrices are given by $P = K[I \mid \mathbf{0}]$ and $P' = K'[R \mid t]$
- We can pre-multiply the projection matrices (and the image points) by the inverse calibration matrices to get *normalized* image coordinates:

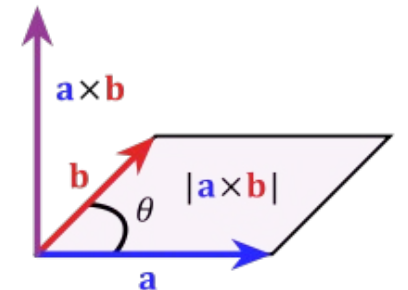
$$\mathbf{x}_{\text{norm}} = K^{-1}\mathbf{x}_{\text{pixel}} \sim [I \mid \mathbf{0}]X, \quad \mathbf{x}'_{\text{norm}} = K'^{-1}\mathbf{x}'_{\text{pixel}} \sim [R \mid t]X$$

Math of the epipolar constraint: Calibrated case

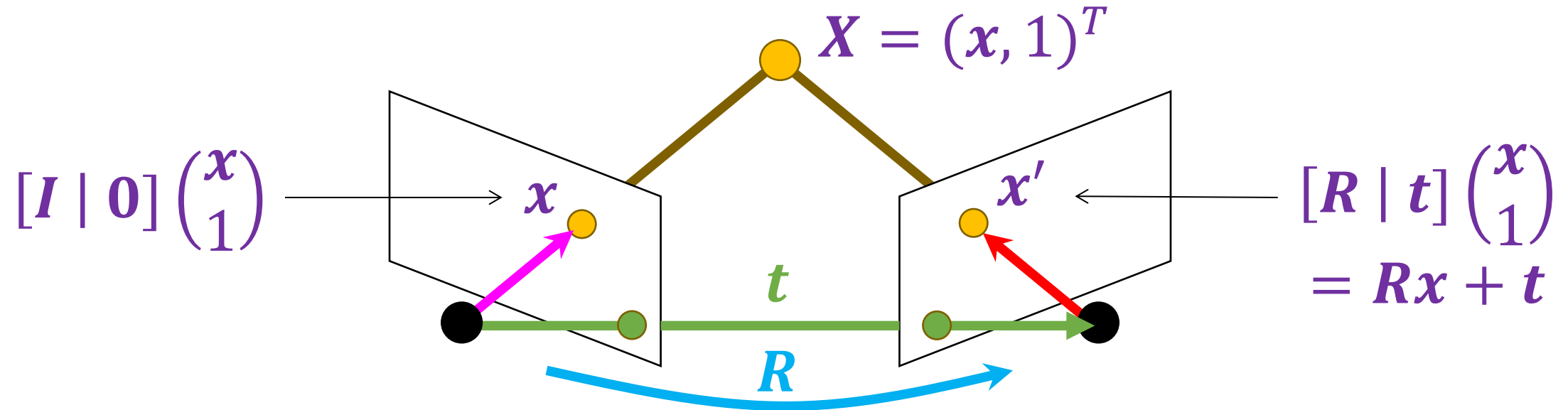


- We have $x' \sim Rx + t$
- This means the three vectors x' , Rx , and t are linearly dependent
- This constraint can be written using the *triple product*

$$x' \cdot [t \times (Rx)] = 0$$



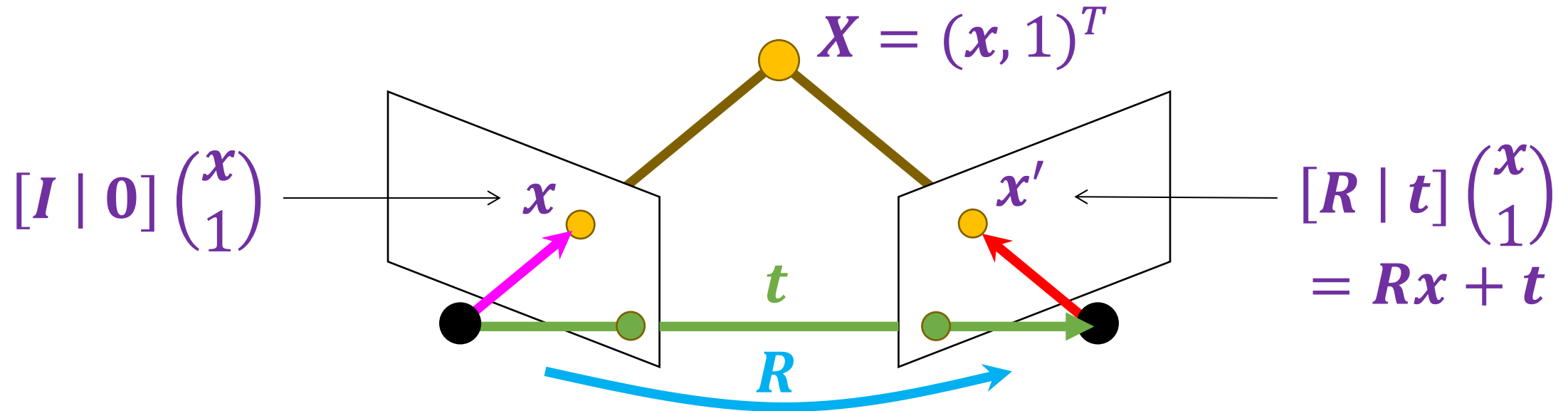
Math of the epipolar constraint: Calibrated case



$$x' \cdot [t \times (Rx)] = 0 \quad \Rightarrow \quad x'^T [t]_{\times} Rx = 0$$

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

Math of the epipolar constraint: Calibrated case



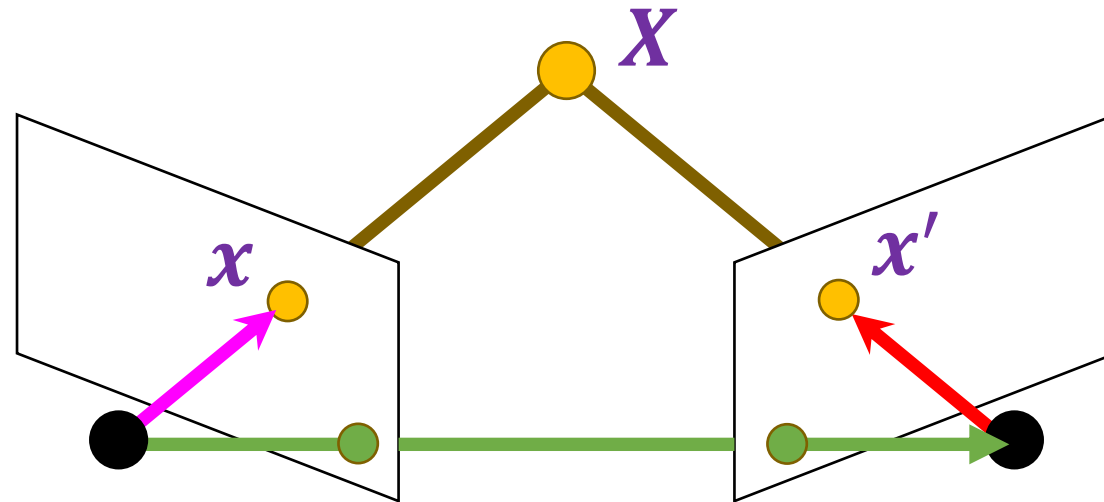
$$[I \mid \mathbf{0}] \begin{pmatrix} x \\ 1 \end{pmatrix}$$

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

$$x' \cdot [t \times (Rx)] = 0 \quad \Rightarrow \quad x'^T [t]_{\times} Rx = 0 \quad \Rightarrow \quad x'^T E x = 0$$

$$E = [t]_{\times} R \quad \text{Essential Matrix}$$

The essential matrix

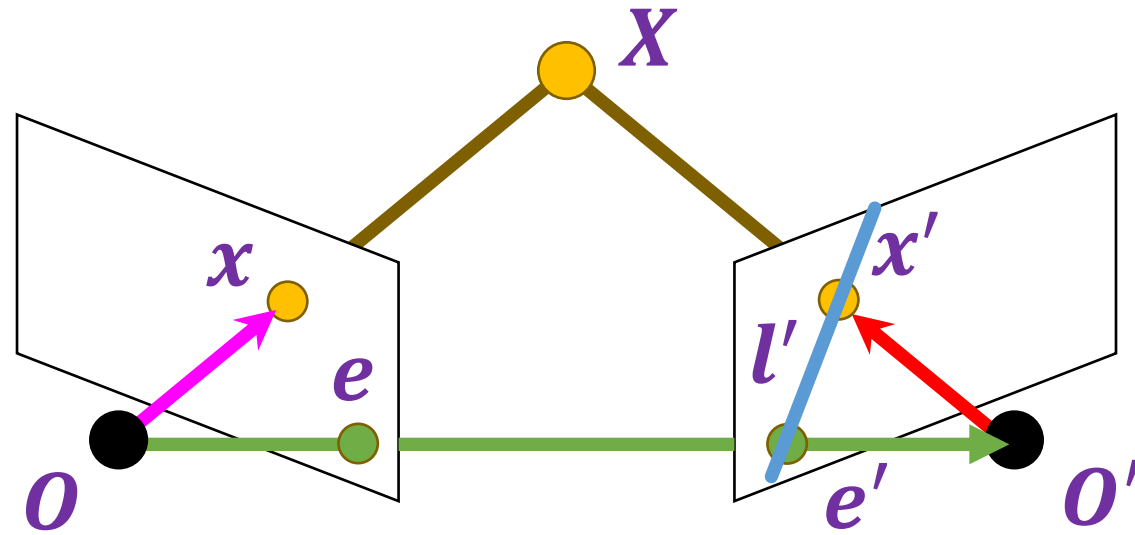


$$x'^T E x = 0$$

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

Essential essential matrix properties

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

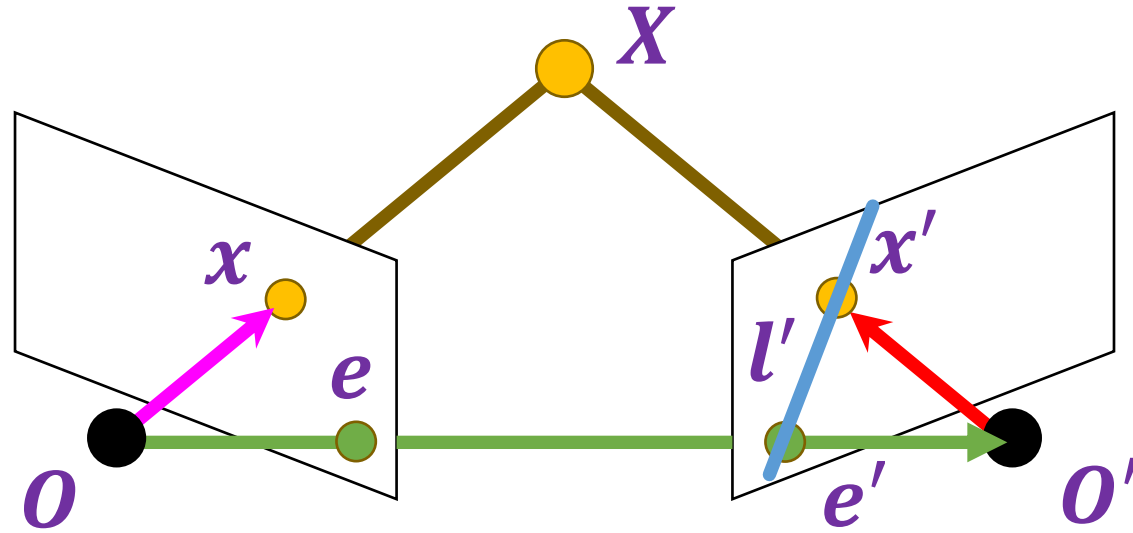


$E x$ is the epipolar line associated with x ($l' = E x$)

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

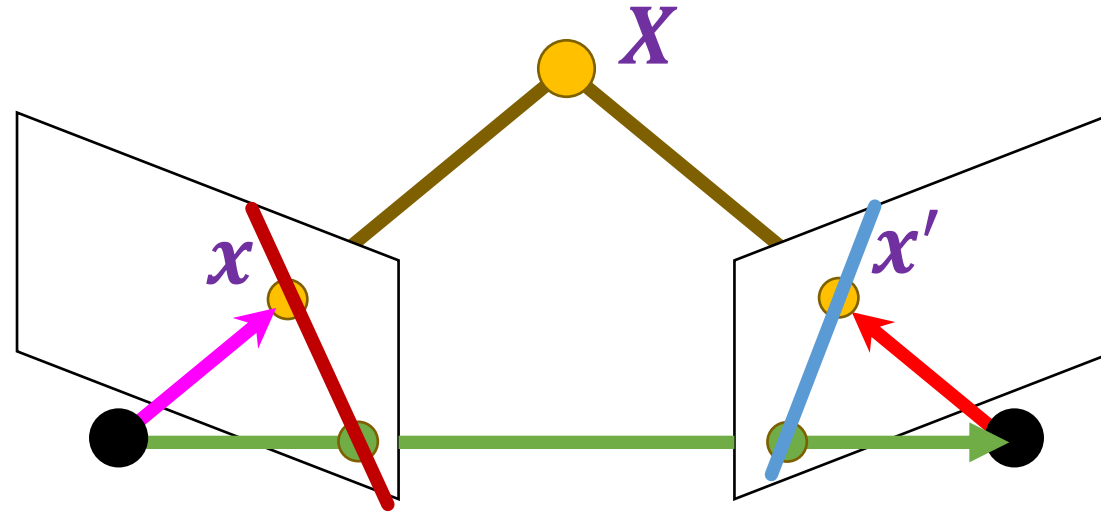
Essential essential matrix properties

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



- $\mathbf{E} \mathbf{x}$ is the epipolar line associated with \mathbf{x} ($\mathbf{l}' = \mathbf{E} \mathbf{x}$)
- $\mathbf{E}^T \mathbf{x}'$ is the epipolar line associated with \mathbf{x}' ($\mathbf{l} = \mathbf{E}^T \mathbf{x}'$)
- Epipoles are the null-spaces of \mathbf{E} : $\mathbf{E} \mathbf{e} = \mathbf{0}$ and $\mathbf{E}^T \mathbf{e}' = \mathbf{0}$
(all epipolar lines pass through epipoles: $\forall \mathbf{x}, \mathbf{l}'^T \mathbf{e}' = \mathbf{x}^T \mathbf{E}^T \mathbf{e}' = 0 \Rightarrow \mathbf{E}^T \mathbf{e}' = \mathbf{0}$)
- \mathbf{E} is singular: rank 2 with 2 non-zero identical singular values
- \mathbf{E} has 5 d.o.f. (3 for R, 3 for t, -1 for scale)

Epipolar constraint: Uncalibrated case

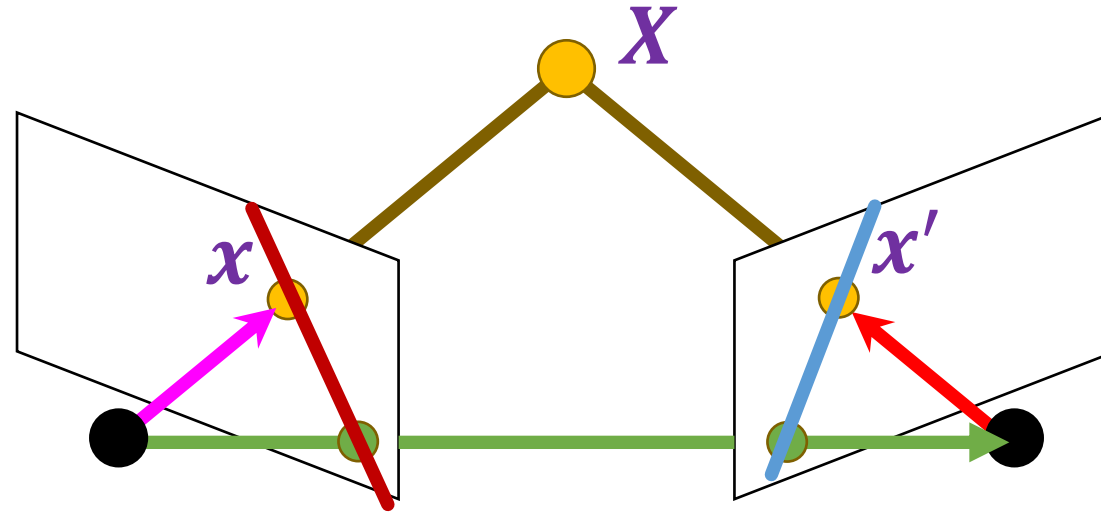


- What if we don't know the camera calibration matrices K and K' ?
- Epipolar constraint in terms of *unknown* normalized coordinates:

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

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

Epipolar constraint: Uncalibrated case



$$\mathbf{x}'_{\text{norm}}{}^T \mathbf{E} \mathbf{x}_{\text{norm}} = 0 \quad \Rightarrow \quad \mathbf{x}'^T \mathbf{F} \mathbf{x} = 0, \text{ where } \mathbf{F} = \mathbf{K}'^{-T} \mathbf{E} \mathbf{K}^{-1}$$

$$\mathbf{x}_{\text{norm}} = \mathbf{K}^{-1} \mathbf{x}$$

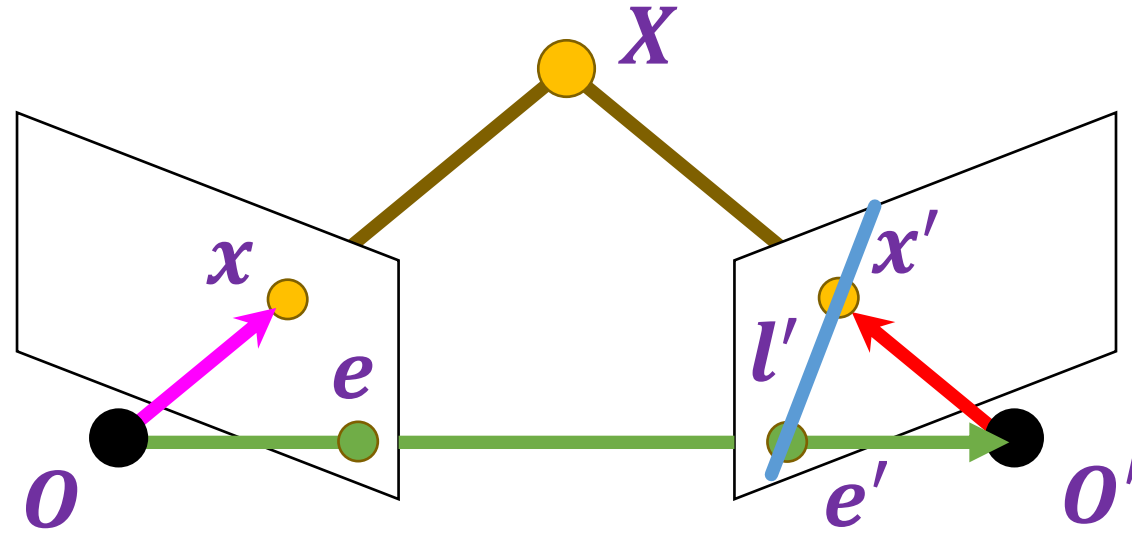
$$\mathbf{x}'_{\text{norm}} = \mathbf{K}'^{-1} \mathbf{x}'$$

Fundamental Matrix

Fundamental matrix properties

Cf. [H&Z ch. 9](#) for proofs

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



- $\mathbf{F} \mathbf{x}$ is the epipolar line associated with \mathbf{x} ($l' = \mathbf{F} \mathbf{x}$)
- $\mathbf{F}^T \mathbf{x}'$ is the epipolar line associated with \mathbf{x}' ($l = \mathbf{F}^T \mathbf{x}'$)
- $\mathbf{F} \mathbf{e} = 0$ and $\mathbf{F}^T \mathbf{e}' = 0$ (can solve for epipoles via... SVD!)
- \mathbf{F} is singular (rank 2)
- \mathbf{F} has 7 d.o.f. (9 parameters, -1 for scale, -1 for singularity)

Estimating the fundamental matrix

- Given: correspondences $\mathbf{x} = (x, y, 1)^T$ and $\mathbf{x}' = (x', y', 1)^T$
- Constraint: $\mathbf{x}'^T \mathbf{F} \mathbf{x} = 0$

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

Enforcing rank-2 constraint

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

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



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



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



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

Enforcing rank-2 constraint

Initial F estimate



Rank-2 estimate

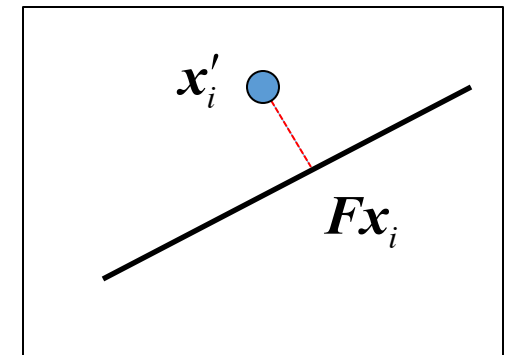
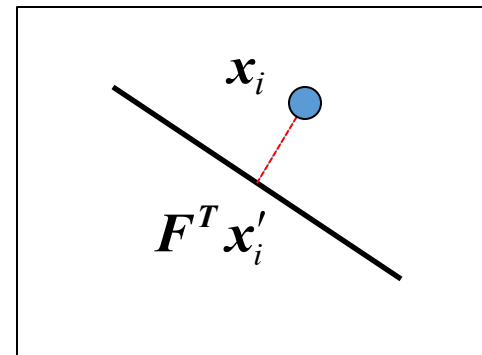


The normalized eight-point algorithm

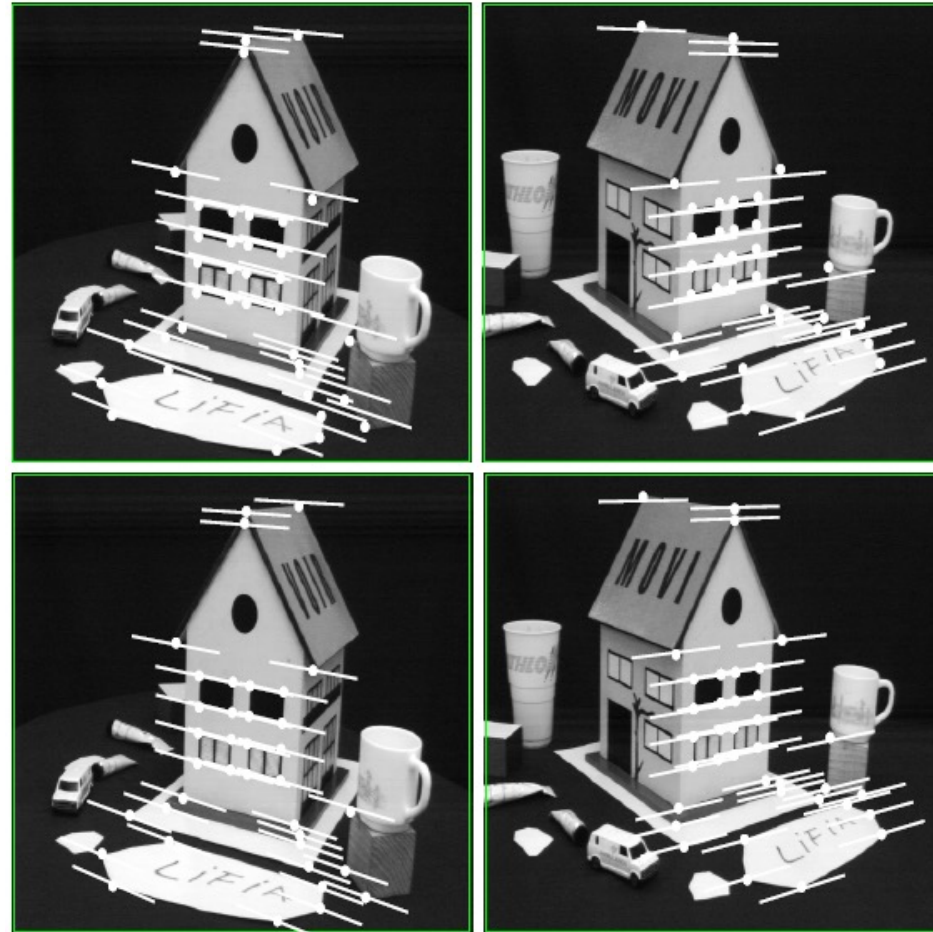
- In each image, center the set of points at the origin, and scale it so the mean squared distance between the origin and the points is 2 pixels
- Use the eight-point algorithm to compute F from the normalized points
- Enforce the rank-2 constraint
- Transform fundamental matrix back to original units: if T and T' are the normalizing transformations in the two images, then the fundamental matrix in original coordinates is $T'^T F T$

Nonlinear estimation

- Linear estimation minimizes the sum of squared *algebraic* distances between points \mathbf{x}'_i and epipolar lines $\mathbf{F}\mathbf{x}_i$ (or points \mathbf{x}_i and epipolar lines $\mathbf{F}^T\mathbf{x}'_i$):
 - $\sum_i (\mathbf{x}'_i{}^T \mathbf{F}\mathbf{x}_i)^2$
- Nonlinear approach: minimize sum of squared *geometric* distances
- $\sum_i [\text{dist}(\mathbf{x}'_i, \mathbf{F}\mathbf{x}_i)^2 + \text{dist}(\mathbf{x}_i, \mathbf{F}^T\mathbf{x}'_i)^2]$



Comparison of estimation algorithms



	8-point	Normalized 8-point	Nonlinear least squares
Av. Dist. 1	2.33 pixels	0.92 pixel	0.86 pixel
Av. Dist. 2	2.18 pixels	0.85 pixel	0.80 pixel

Seven-point algorithm

- Set up least squares system with seven pairs of matches and solve for null space (two vectors) using SVD
- Solve for polynomial equation to get coefficients of linear combination of null space vectors that satisfies $\det(F) = 0$

The Fundamental Matrix Song

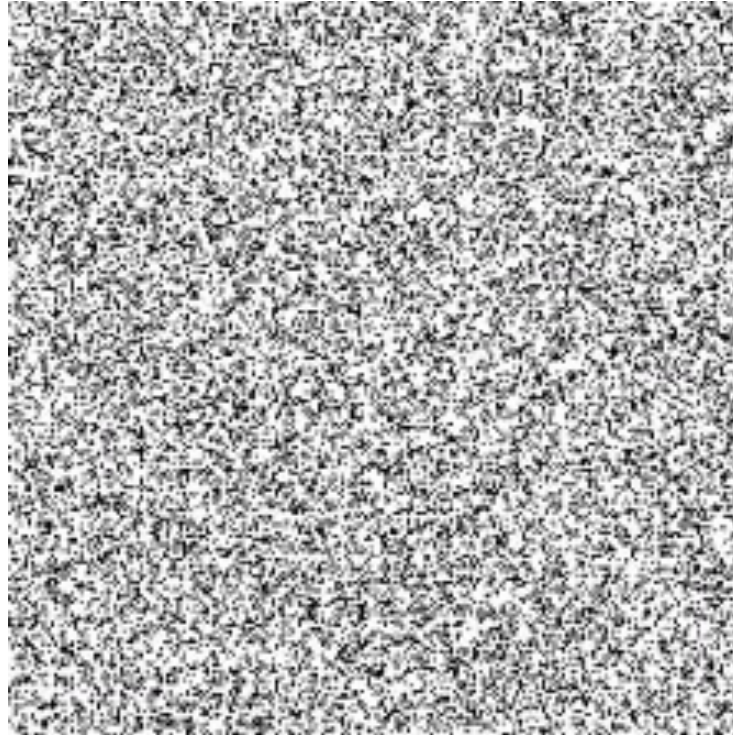


<http://danielwedge.com/fmatrix/>

Bonus slides – More on Stereo

History: Random dot stereograms

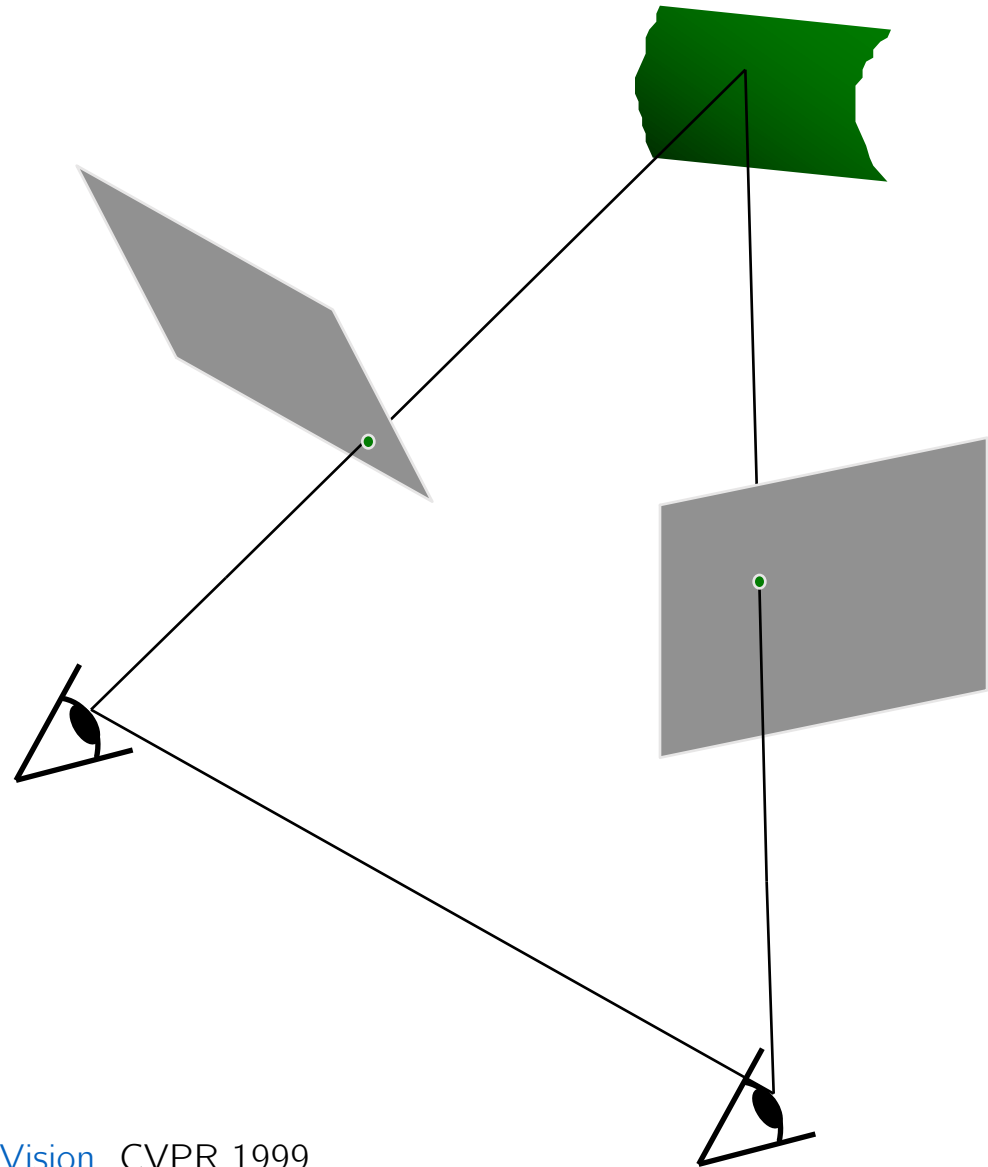
- Invented by [Bela Julesz](#) in the mid-20th century
- Demonstration that stereo perception can happen without any monocular cues



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

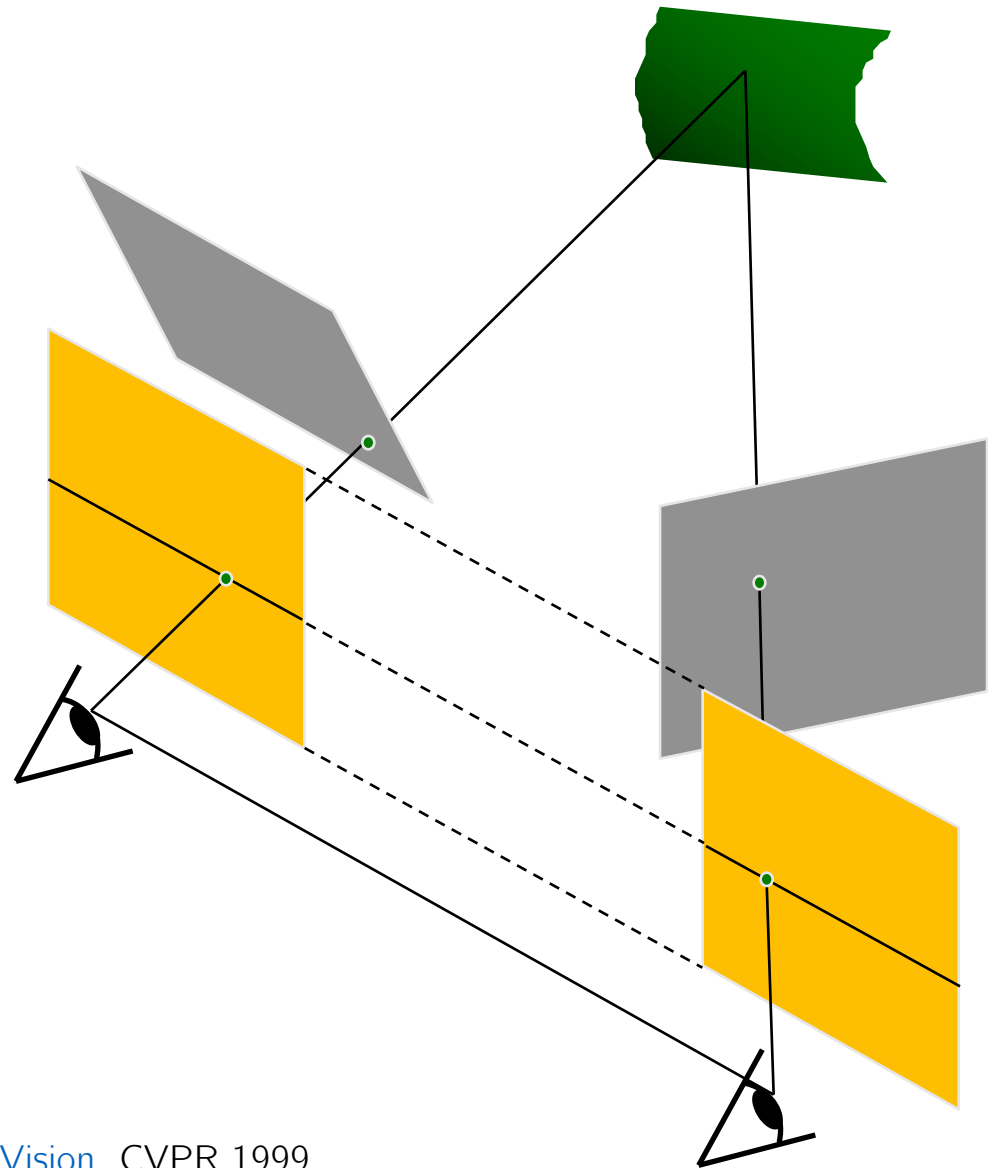
Stereo image rectification

If the image planes are not parallel, we can find homographies to project each view onto a common plane parallel to the baseline



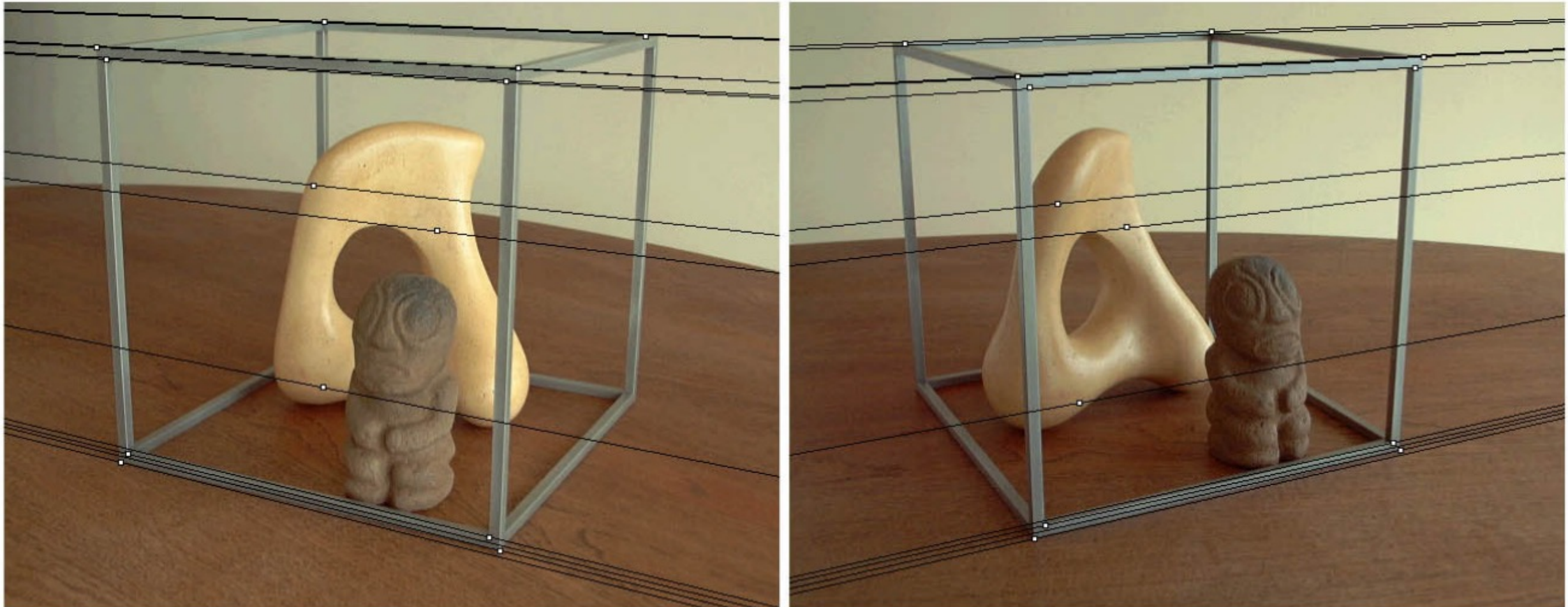
Stereo image rectification

If the image planes are not parallel, we can find homographies to project each view onto a common plane parallel to the baseline



Stereo image rectification

Before rectification:



Stereo image rectification

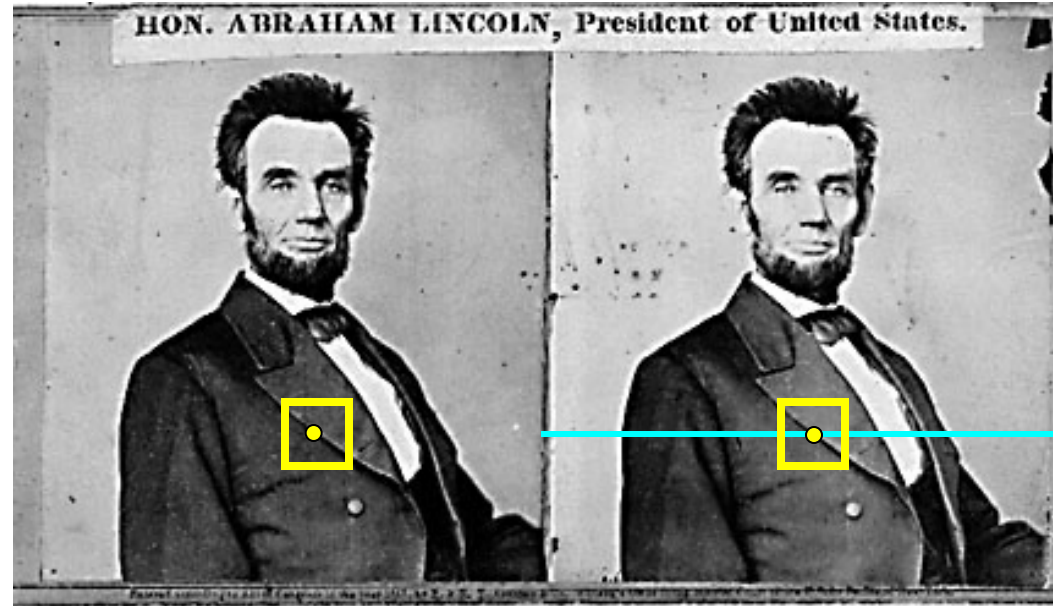
After rectification:



Another rectification example

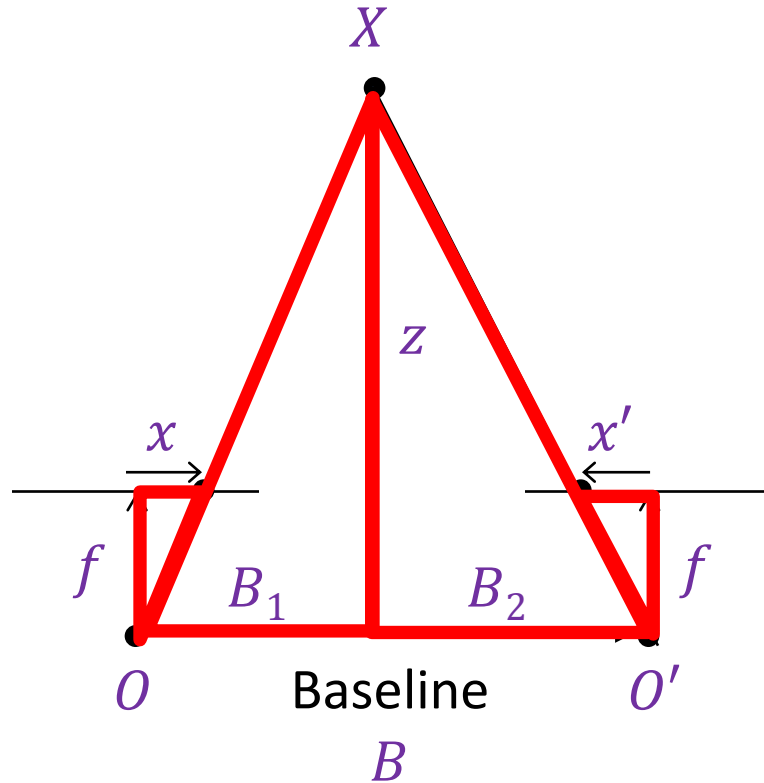


Basic stereo matching algorithm



- If necessary, rectify the two stereo images to transform epipolar lines into horizontal scanlines
- For each pixel x in the first image
 - Find corresponding epipolar scanline in the right image
 - Examine all pixels on the scanline and pick the best match x'
 - Triangulate the matches to get depth information

Depth from disparity



$$\frac{x}{f} = \frac{B_1}{z} \quad \frac{-x'}{f} = \frac{B_2}{z}$$

$$\frac{x - x'}{f} = \frac{B_1 + B_2}{z}$$

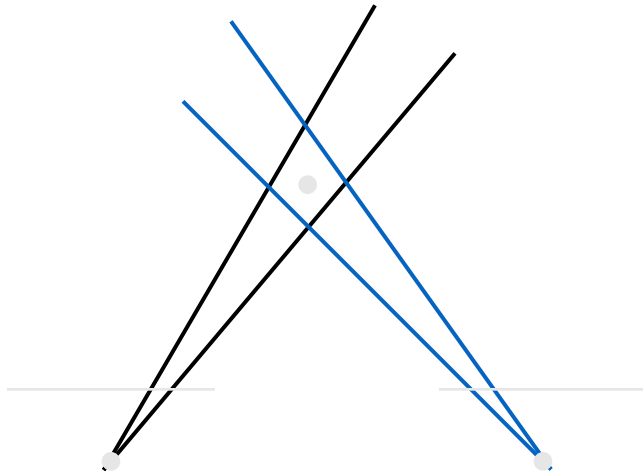
$$\underbrace{x - x'}_{\text{Disparity}} = \frac{fB}{z}$$

Disparity

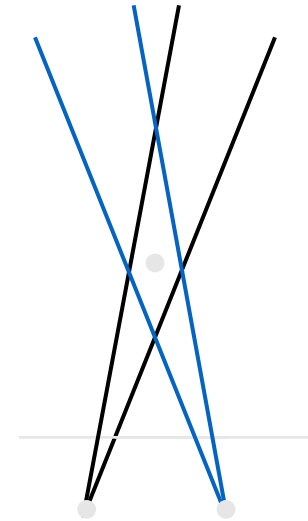
Disparity is inversely proportional to depth!

$$z = \frac{fB}{x - x'}$$

Effect of baseline on stereo results

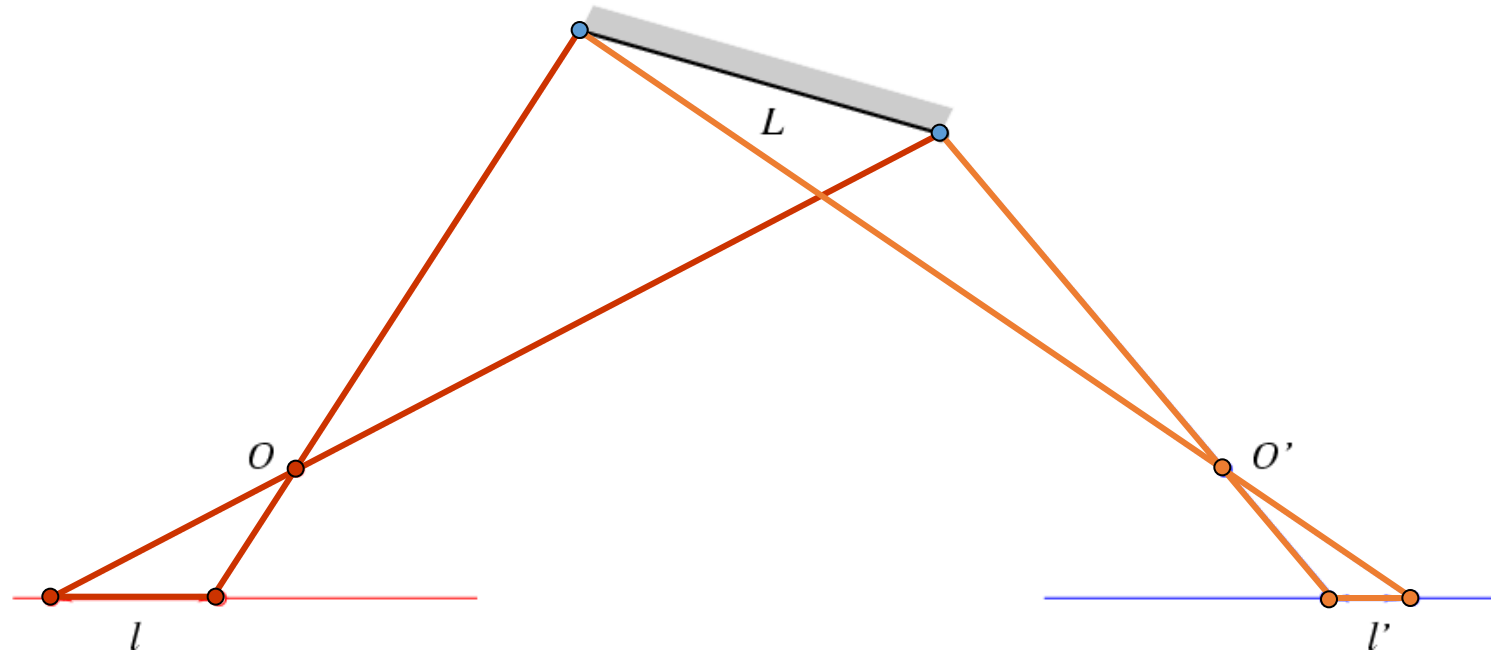


- Larger baseline
 - + Smaller triangulation error
 - Matching is more difficult



- Smaller baseline
 - Higher triangulation error
 - + Matching is easier

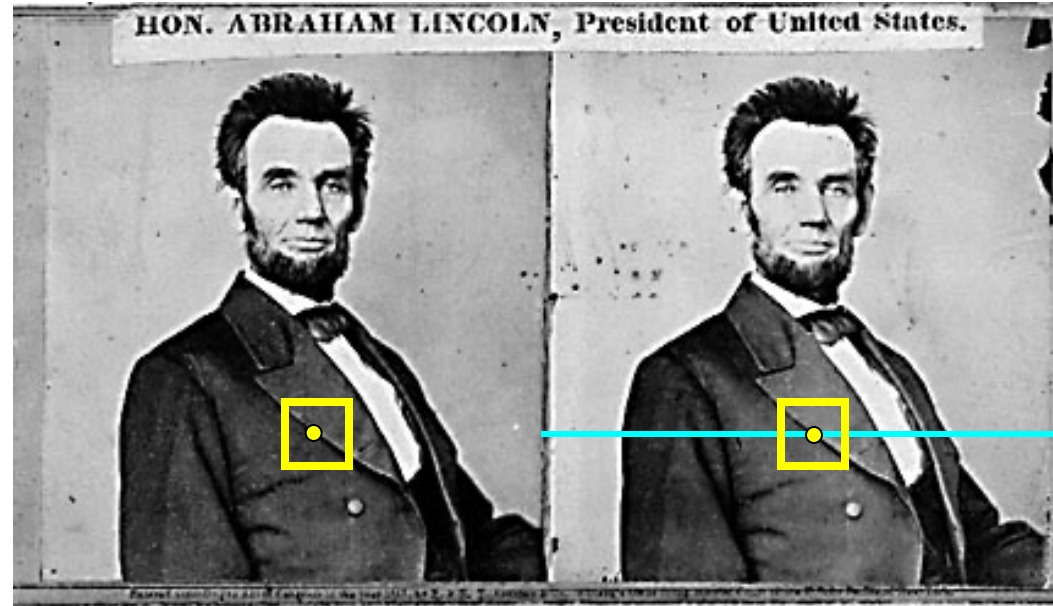
Problem for wide baselines: Foreshortening



- Matching with fixed-size windows will fail!
- Possible solution: adaptively vary window size
- Another solution: *model-based stereo* (CS231a)

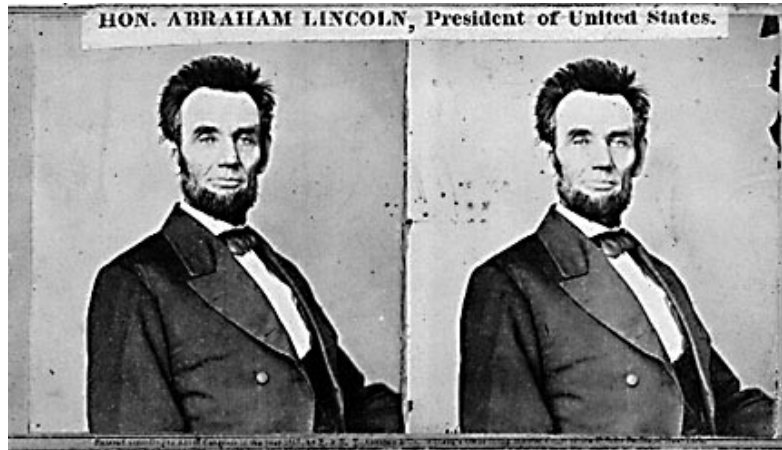
Slide credit: J. Hayes

Basic stereo matching algorithm



- If necessary, rectify the two stereo images to transform epipolar lines into scanlines
- For each pixel x in the first image
 - Find corresponding epipolar scanline in the right image
 - Examine all pixels on the scanline and pick the best match x'
 - Compute disparity $x - x'$ and set $\text{depth}(x) = Bf / (x - x')$

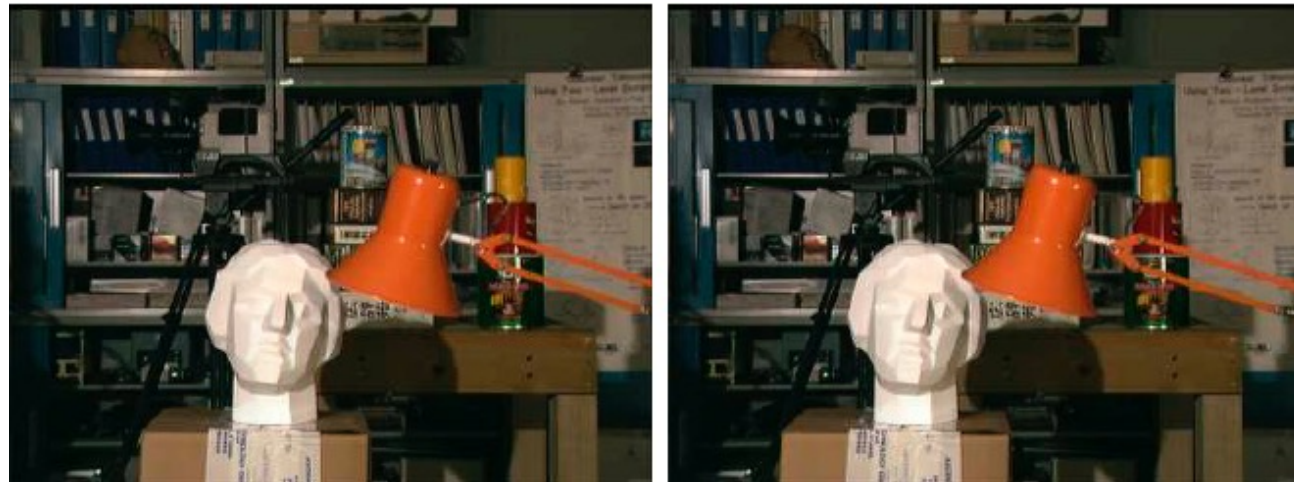
The correspondence problem



Textureless surfaces



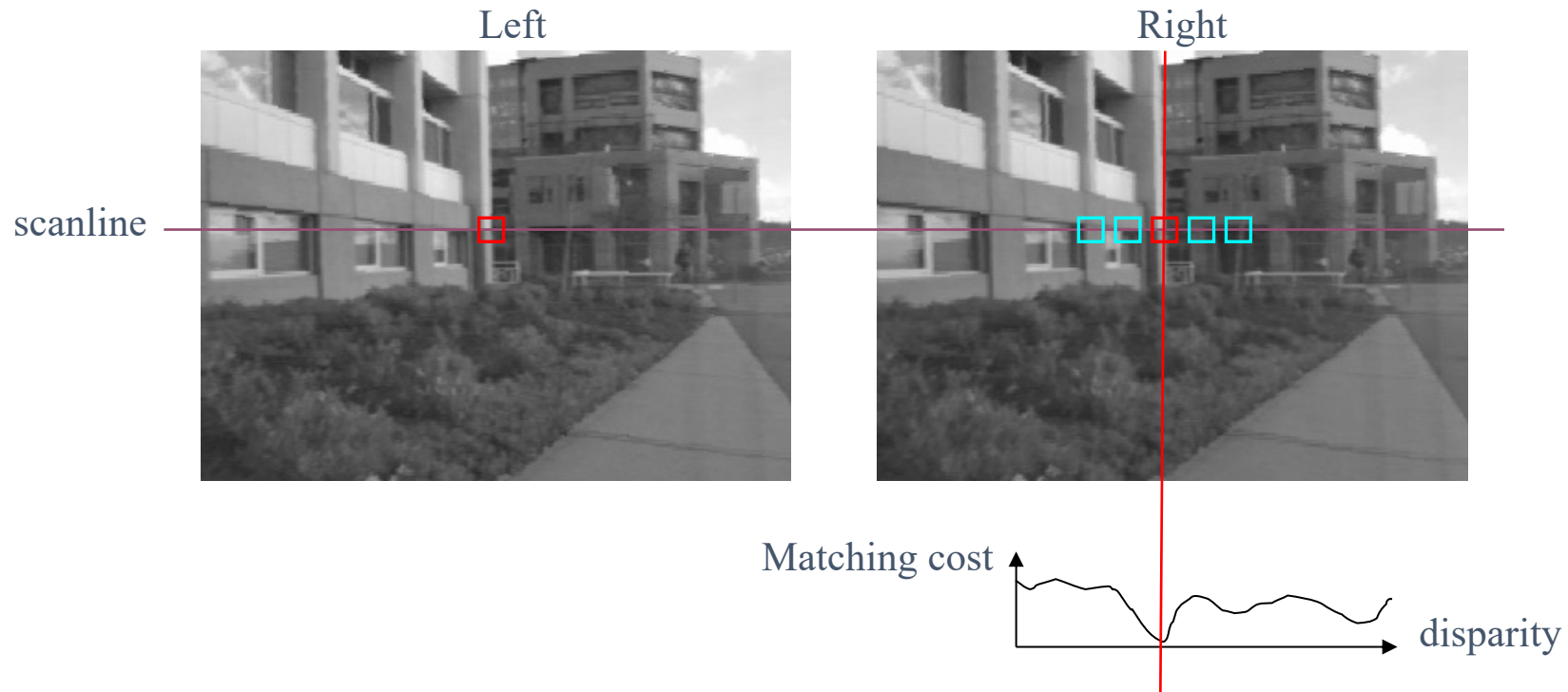
Occlusions, repetition



Non-Lambertian surfaces, specularities

How to find the best match?

- Compare pixels: sliding window + SSD or normalized cross-correlation



How to find the best match?

- Compare pixels: sliding window + SSD or normalized cross-correlation
- Global optimization (Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001)

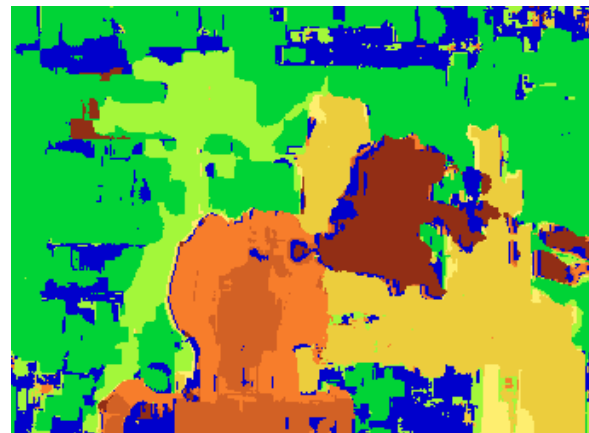
Data



Ground truth



Window-based matching

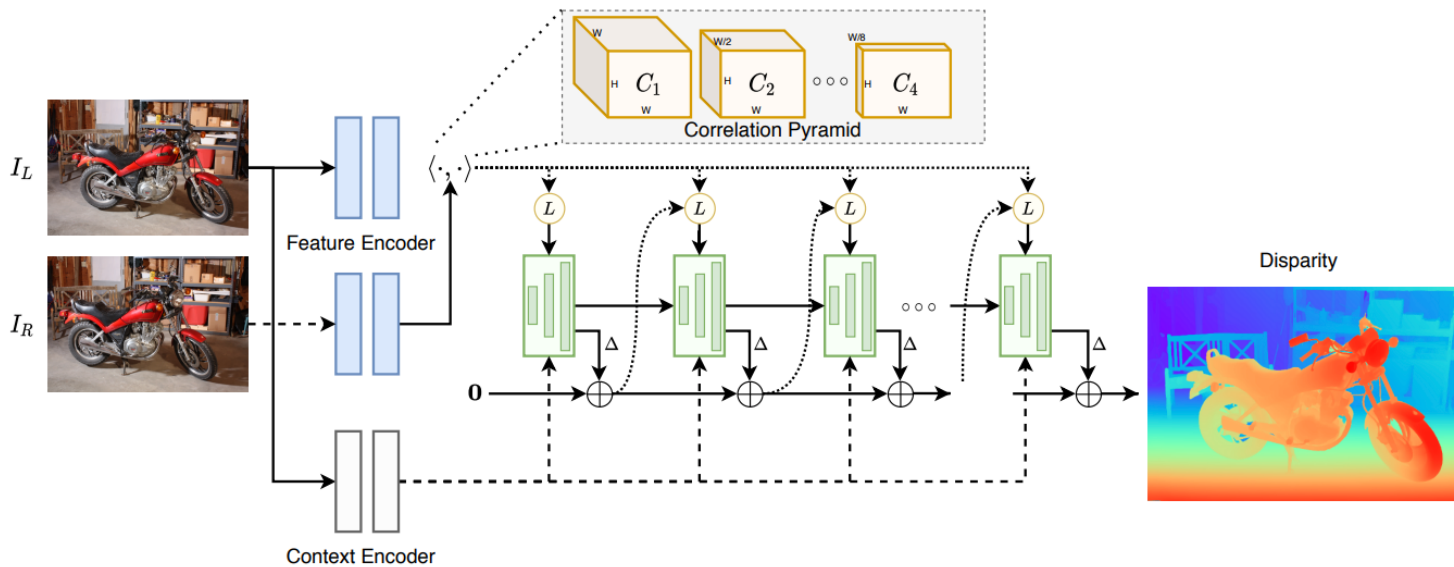


Global optimization method (graph cuts)

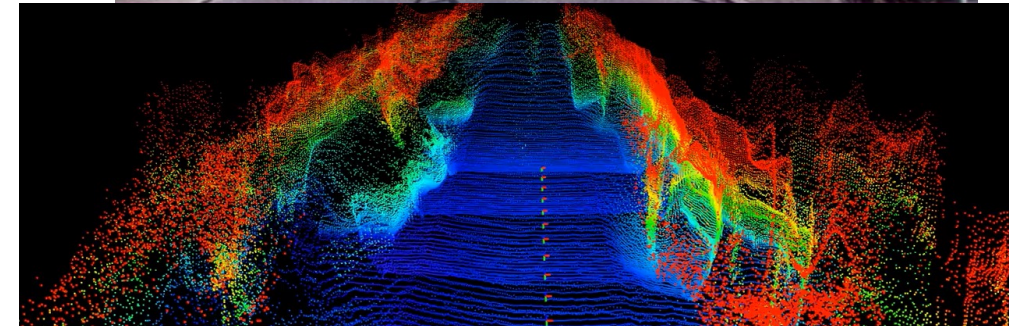


How to find the best match?

- Compare pixels: sliding window + SSD or normalized cross-correlation
- Global optimization (Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001)
- Deep Learning: match learned representations, end-to-end estimation

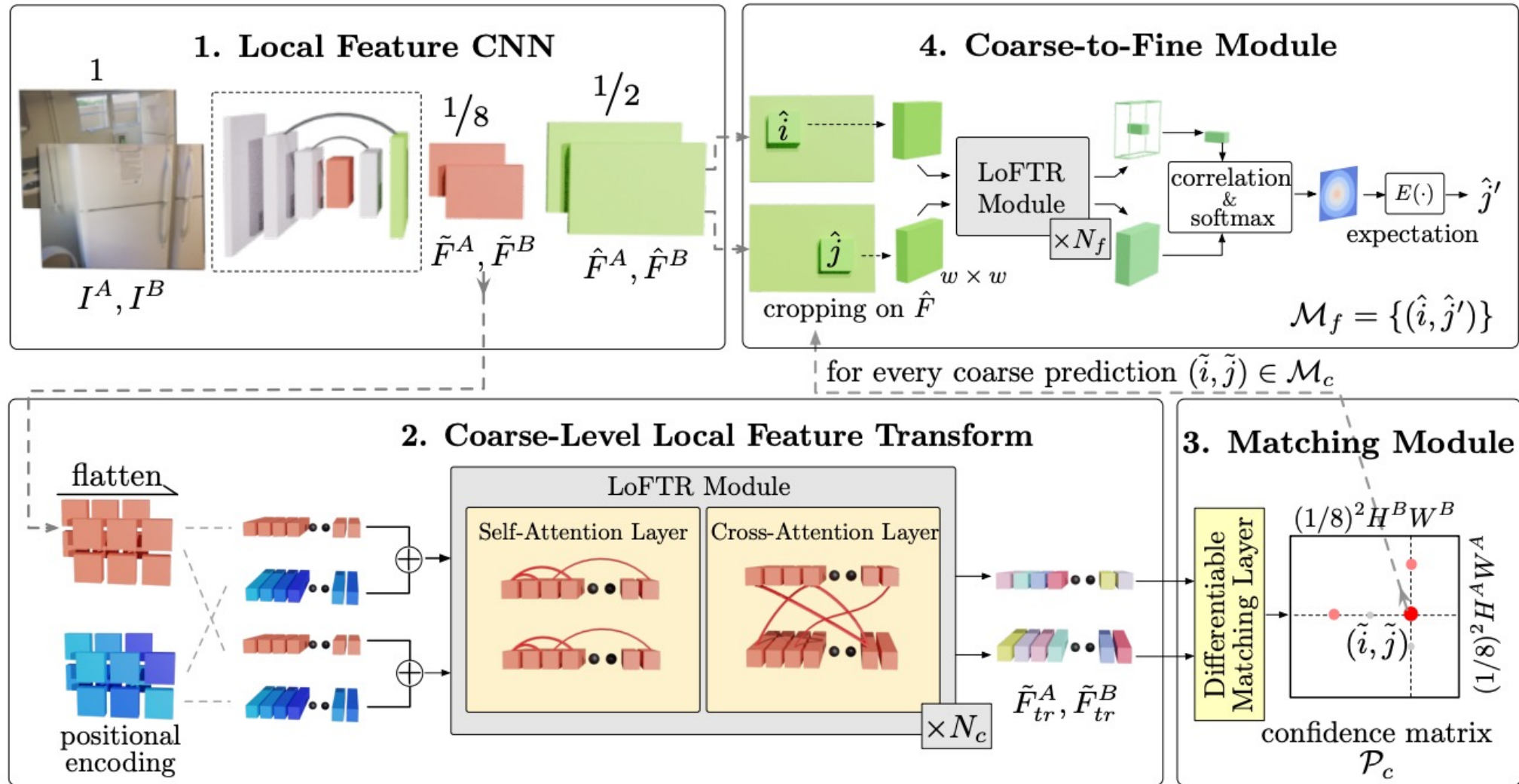


[RAFT-Stereo: Multilevel Recurrent Field Transforms for Stereo Matching](#), L. Lipson et al, arXiv 2021



[SuperDepth: Self-Supervised, Super-Resolved Monocular Depth Estimation](#), Pillai et al, ICRA'19

Correspondences in the age of Deep Learning

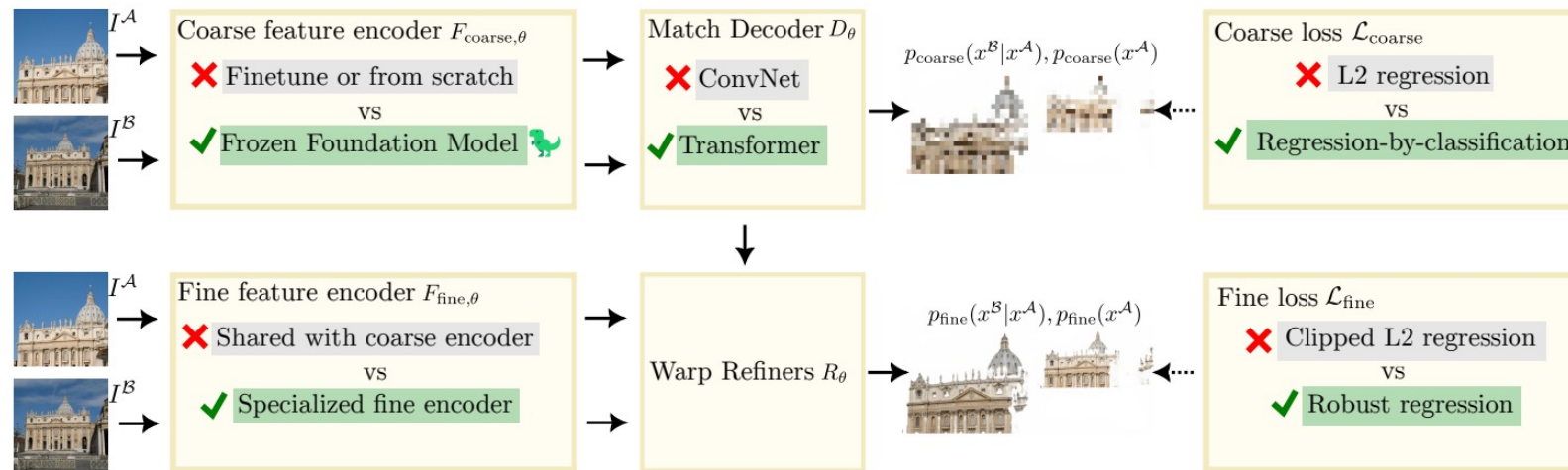


Correspondences in the age of Deep Learning

RoMa: Robust Dense Feature Matching

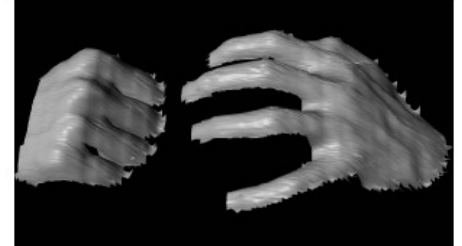
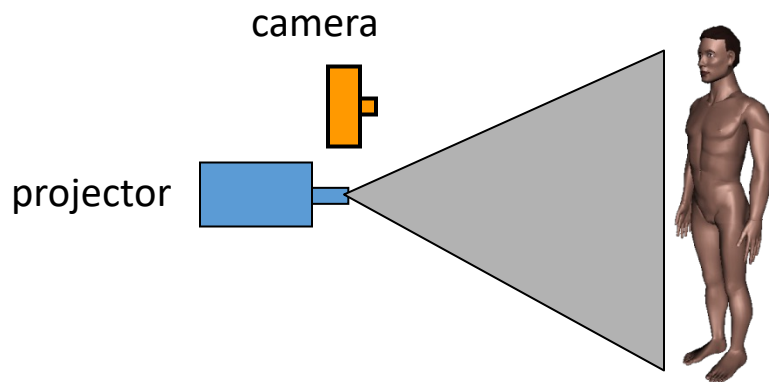
Johan Edstedt¹ Qiyu Sun² Georg Bökman³ Mårten Wadenbäck¹ Michael Felsberg¹
¹Linköping University, ²East China University of Science and Technology, ³Chalmers University of Technology

11 Dec 2023

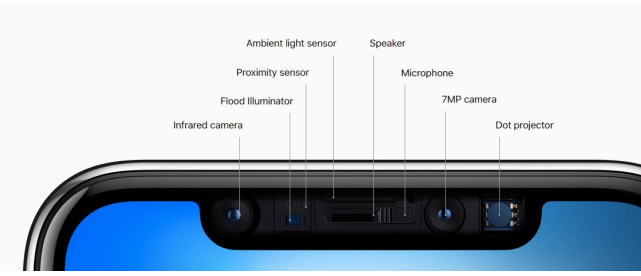


How to find the best match?

- Compare pixels: sliding window + SSD or normalized cross-correlation
- Global optimization (Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001)
- **Deep Learning**: match learned representations, end-to-end estimation
- “Active stereo”: project lasers or structured patterns (e.g., in IR)



L. Zhang, et al. [Rapid Shape Acquisition Using Color Structured Light and Multi-pass Dynamic Programming](#). 3DPVT 2002



<http://bbzipo.wordpress.com/2010/11/28/kinect-in-infrared/>

<https://www.cnet.com/news/apple-face-id-true-depth-how-it-works/>

Bonus slides – SfM ambiguities

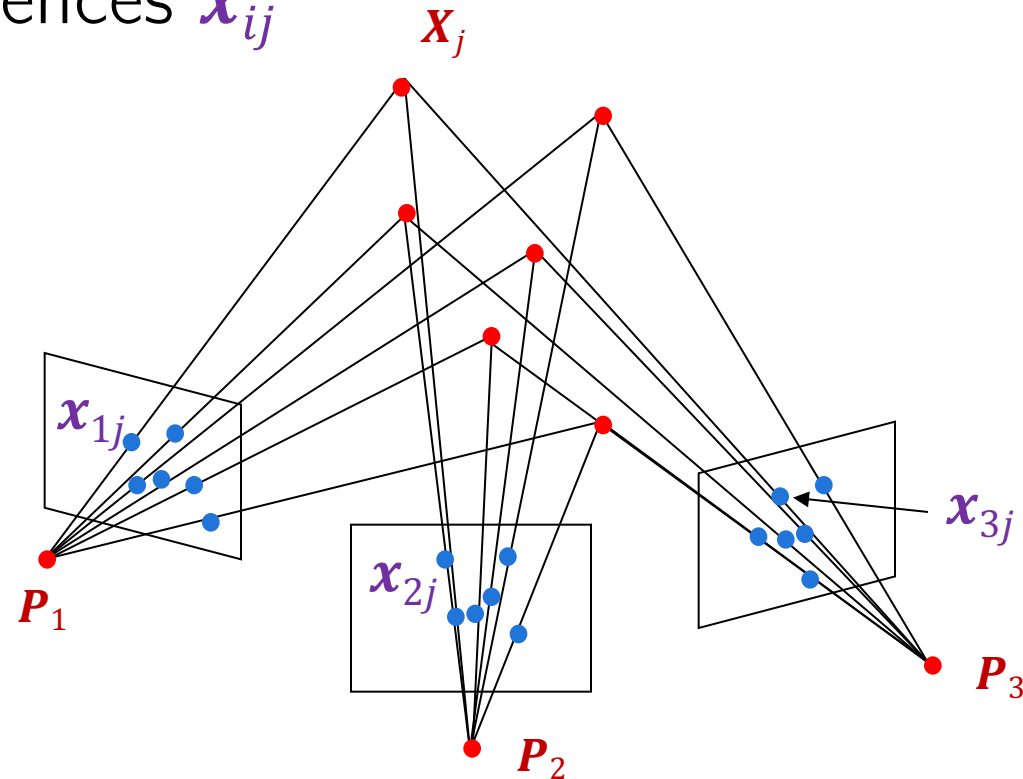
Slides credit: S. Lazebnik

Projective structure from motion

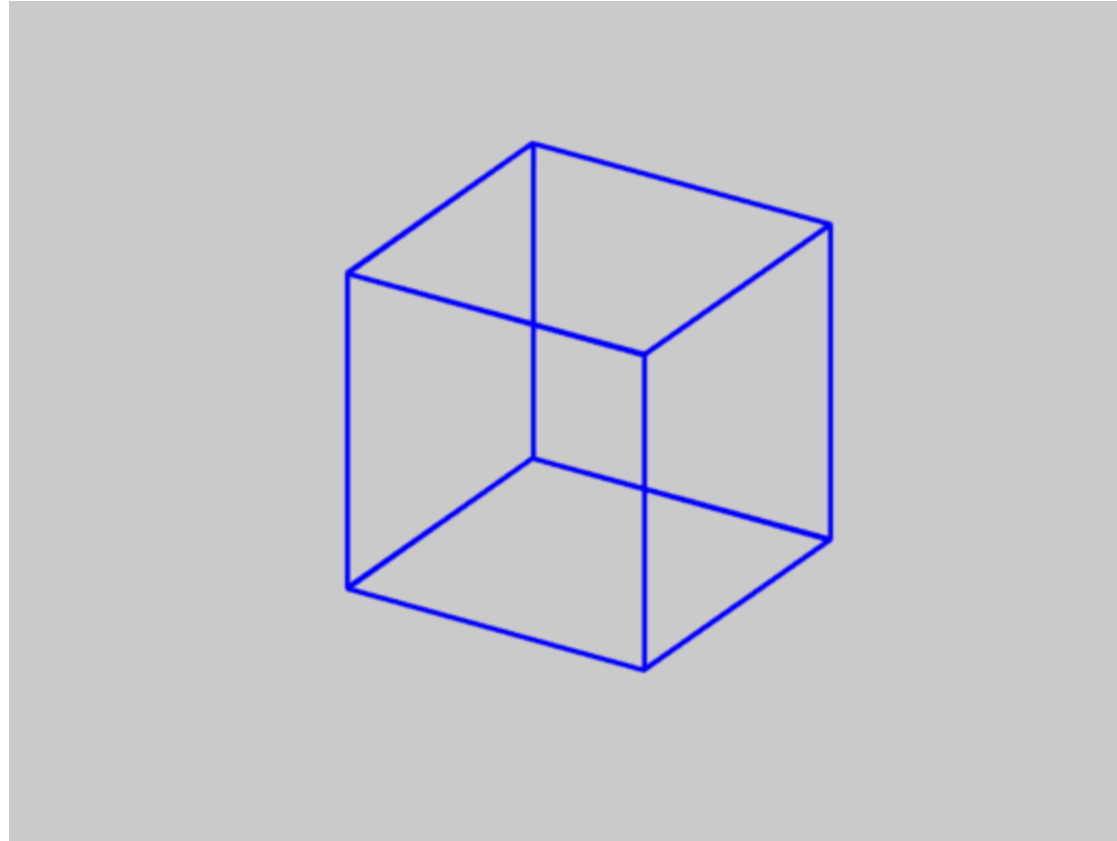
- *Given*: m images of n fixed 3D points such that (ignoring visibility):

$$\mathbf{x}_{ij} \sim \mathbf{P}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$

- *Problem*: estimate m projection matrices \mathbf{P}_i and n 3D points \mathbf{X}_j from the mn correspondences \mathbf{x}_{ij}



Is SFM always uniquely solvable?



Necker cube

Structure from motion ambiguity

If we scale the entire scene by some factor k and, at the same time, scale the camera matrices by the factor of $1/k$, the projections of the scene points remain exactly the same:

$$x \sim PX = \left(\frac{1}{k} P \right) (kX)$$

Without a reference measurement, it is impossible to recover the absolute scale of the scene!

In general, if we transform the scene using a transformation Q and apply the inverse transformation to the camera matrices, then the image observations do not change:

$$x \sim PX = (PQ^{-1})(QX)$$

Projective structure from motion

- *Given*: m images of n fixed 3D points such that (ignoring visibility):

$$\mathbf{x}_{ij} \sim \mathbf{P}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$

- *Problem*: estimate m projection matrices \mathbf{P}_i and n 3D points \mathbf{X}_j from the mn correspondences \mathbf{x}_{ij}
- With no calibration info, cameras and points can only be recovered up to a 4×4 projective transformation \mathbf{Q} :

$$\mathbf{X} \rightarrow \mathbf{QX}, \mathbf{P} \rightarrow \mathbf{PQ}^{-1} \Rightarrow \mathbf{x} \sim \mathbf{PX} = \mathbf{PQ}^{-1}\mathbf{QX}$$

- We can solve for structure and motion when $2mn \geq 11m + 3n - 15$
- Inequality above \rightarrow for $m = 2$ cameras, need at least $n = 7$ points
- Remember: \mathbf{F} has 7 d.o.f., not a coincidence!

Projective SFM: Two-camera case

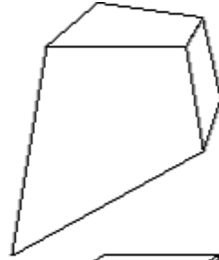
1. Estimate fundamental matrix F between the two views
2. Set first camera matrix to $[I \mid \mathbf{0}]$
3. Then the second camera matrix is given by $[A \mid \mathbf{t}]$ where \mathbf{t} is the epipole ($F^T \mathbf{t} = \mathbf{0}$) and $A = [\mathbf{t}]_{\times} F$

For derivation: cf. CS231A or [H&Z Chap.9](#) result 9.9

Types of ambiguity

Projective
15dof

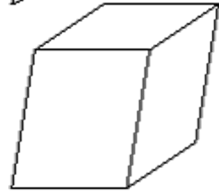
$$\begin{bmatrix} A & t \\ \mathbf{v}^\top & \nu \end{bmatrix}$$



Preserves intersection and tangency

Affine
12dof

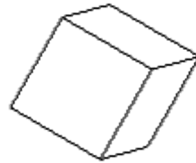
$$\begin{bmatrix} A & t \\ \mathbf{0}^\top & 1 \end{bmatrix}$$



Preserves parallelism, volume ratios

Similarity
7dof

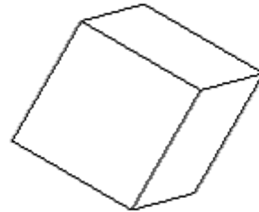
$$\begin{bmatrix} sR & t \\ \mathbf{0}^\top & 1 \end{bmatrix}$$



Preserves angles, ratios of length

Euclidean
6dof

$$\begin{bmatrix} R & t \\ \mathbf{0}^\top & 1 \end{bmatrix}$$



Preserves angles, lengths

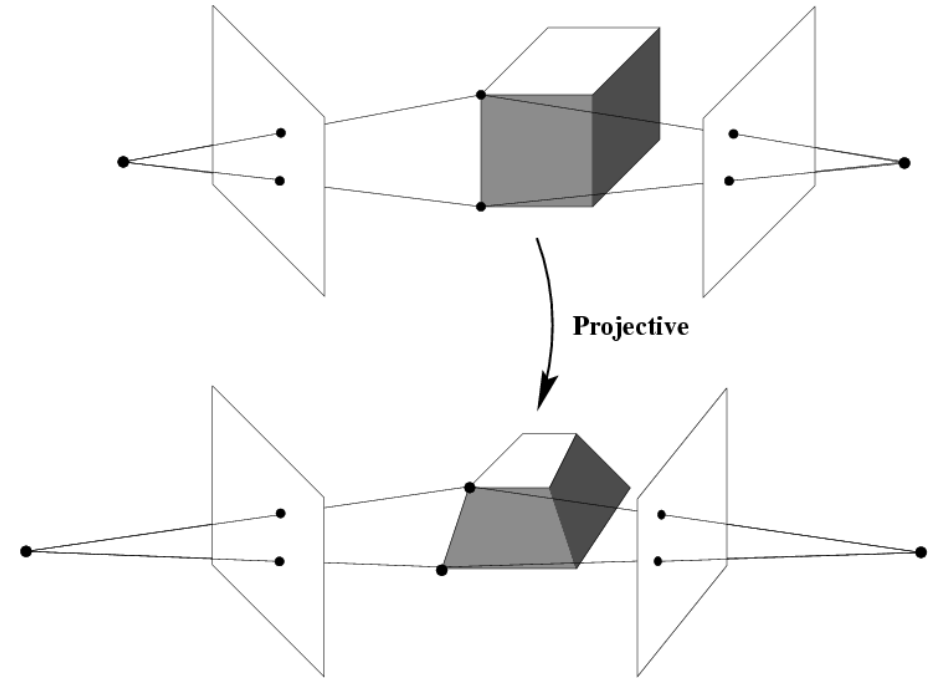
- With no constraints on the camera calibration matrix or on the scene, we get a *projective* reconstruction
- Need additional information to *upgrade* the reconstruction to affine, similarity, or Euclidean

Projective ambiguity

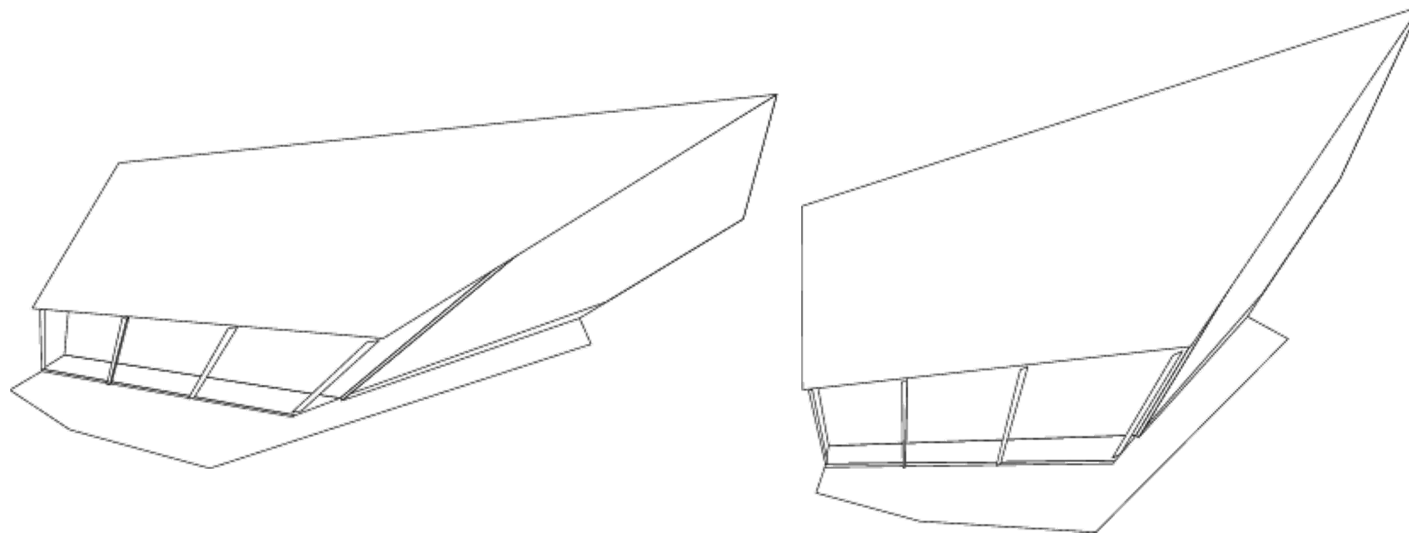
With no constraints on the camera calibration matrices or on the scene, we can reconstruct up to a *projective* ambiguity:

$$x \sim PX = (PQ^{-1})(QX)$$

Q is a general full-rank 4×4 matrix



Projective ambiguity



Affine ambiguity

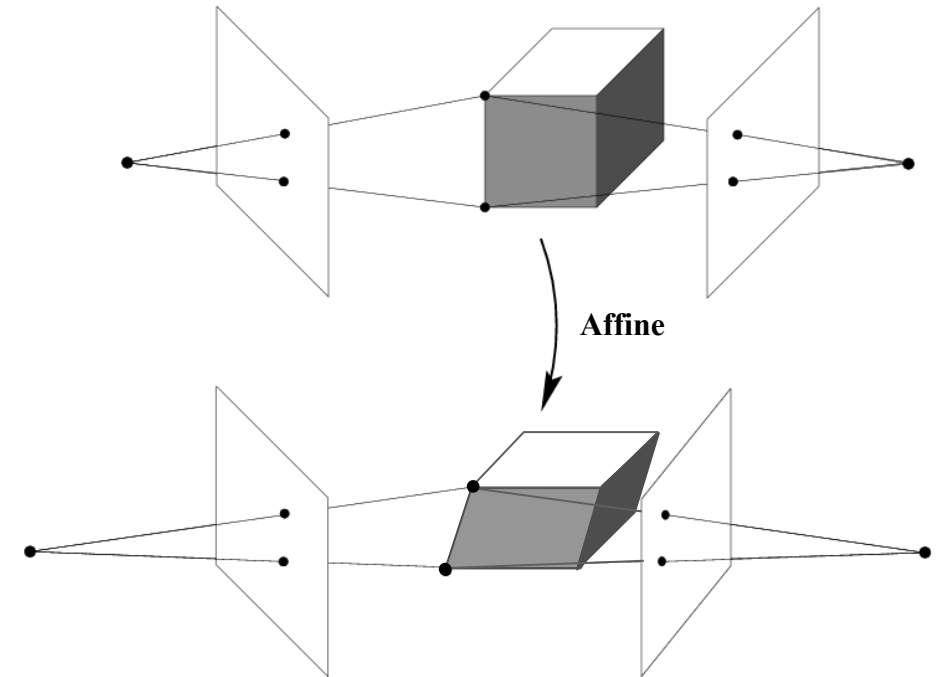
If we impose parallelism constraints, we can get a reconstruction up to an *affine* ambiguity:

$$x \sim PX = (PQ_A^{-1})(Q_A X)$$

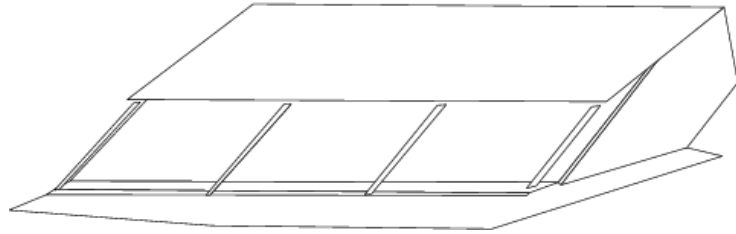
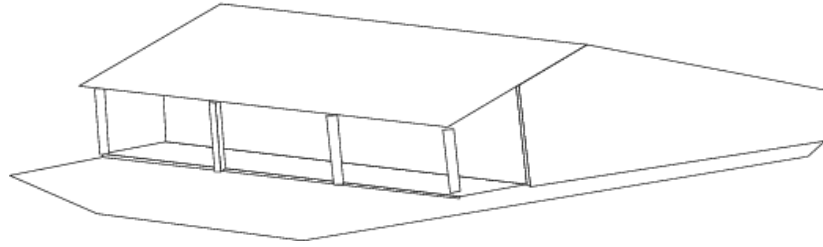
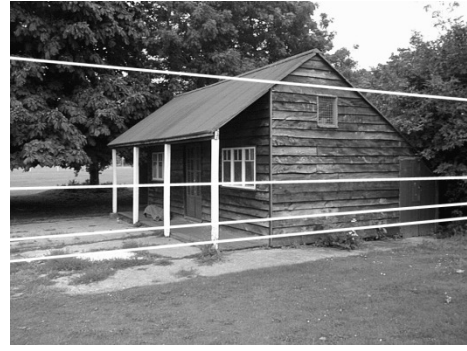
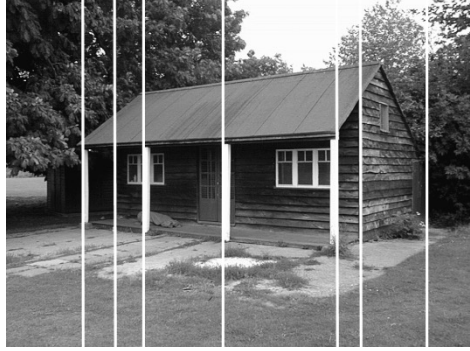
3×3 full-rank matrix

3×1 translation vector

$$Q_A = \begin{bmatrix} A & t \\ \mathbf{0}^T & 1 \end{bmatrix}$$



Affine ambiguity



Similarity ambiguity

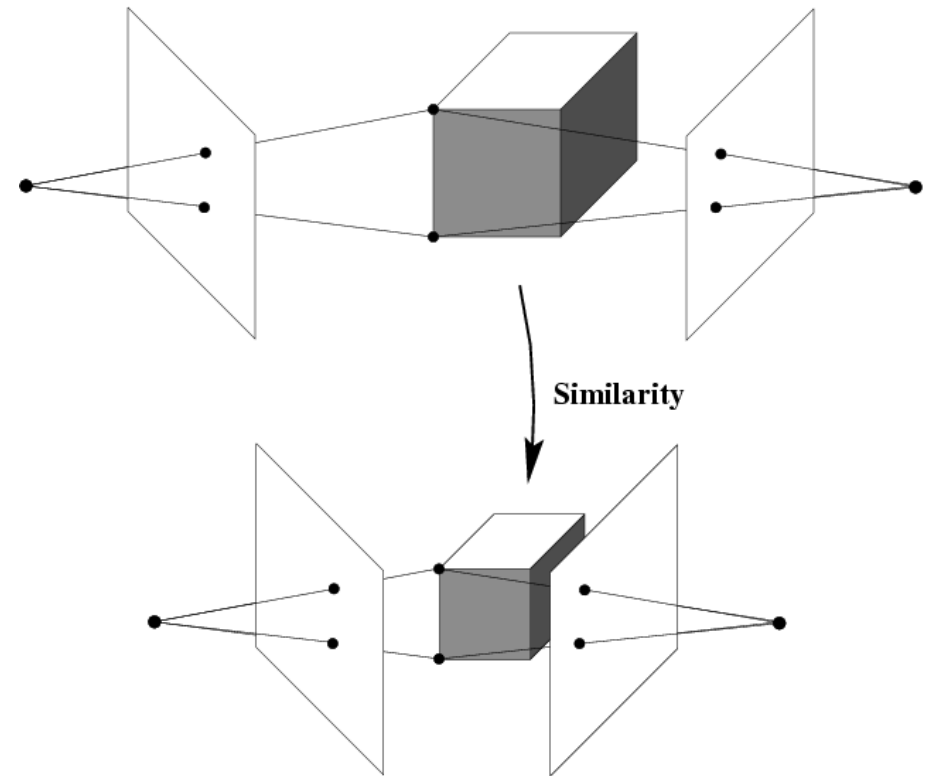
A reconstruction that obeys orthogonality constraints on camera parameters and/or scene

$$x \sim PX = (PQ_S^{-1})(Q_S X)$$

3×3
rotation
matrix

3×1 translation
vector

$$Q_S = \begin{bmatrix} sR & t \\ \mathbf{0}^T & 1 \end{bmatrix}$$



Similarity ambiguity

