

---

# Computer Vision & Image Processing

## CSE 473 / 573

Instructor - Kevin R. Keane, PhD

TAs - Radhakrishna Dasari, Yuhao Du, Niyazi Sorkunlu

Lecture 20

October 16, 2017

Structure from Motion

# Slide credits

---

Svetlana Lazebnik

# Reading

---

F&P chapter 8

# Structure from motion

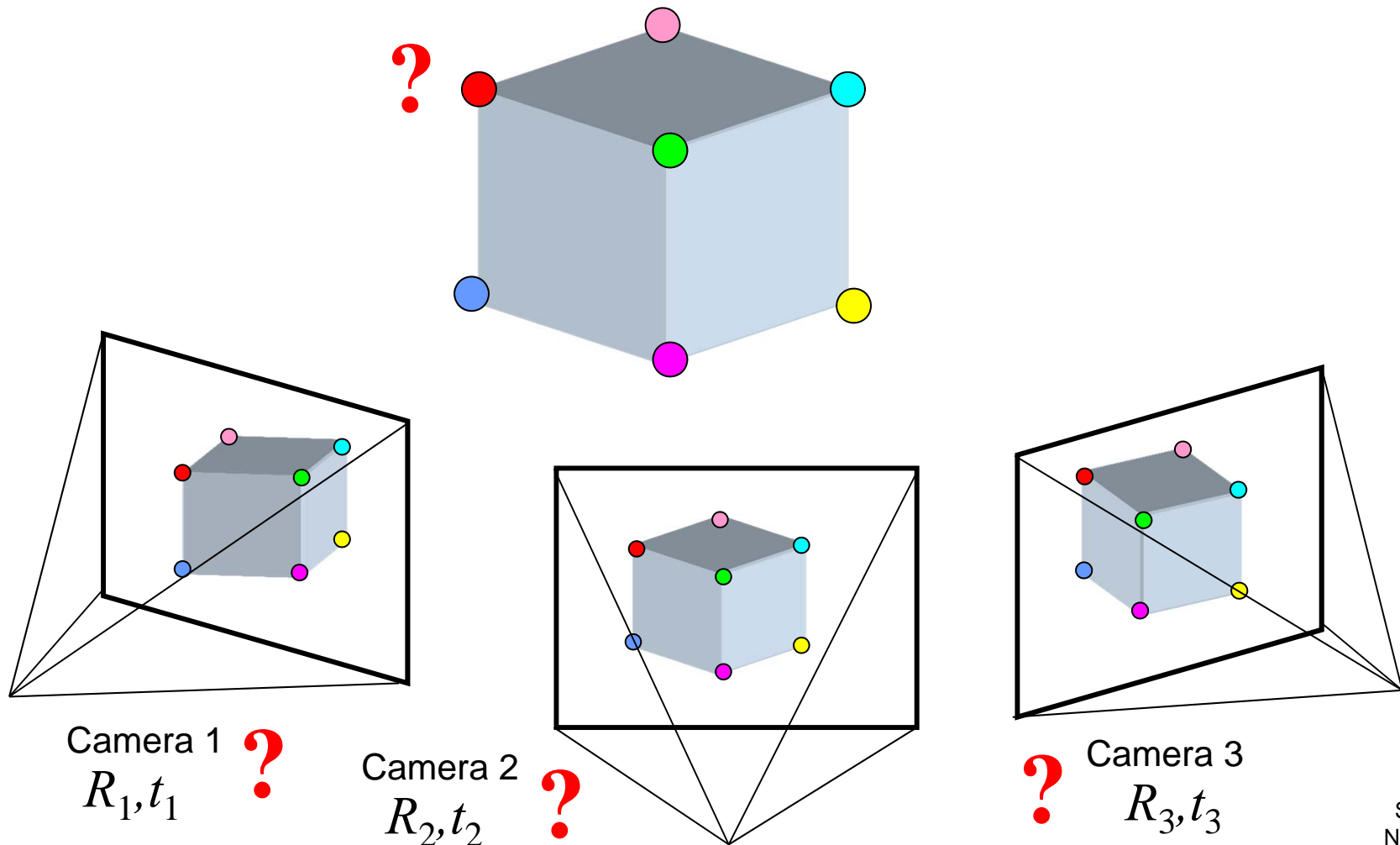
---



Драконъ, видимый подъ различными углами зрѣнія  
По гравюру на мѣди изъ „Oculus artificialis teleiopicus“ Цана. 1702 года.

# Structure from motion

- Given a set of corresponding points in two or more images, compute the camera parameters and the 3D point coordinates



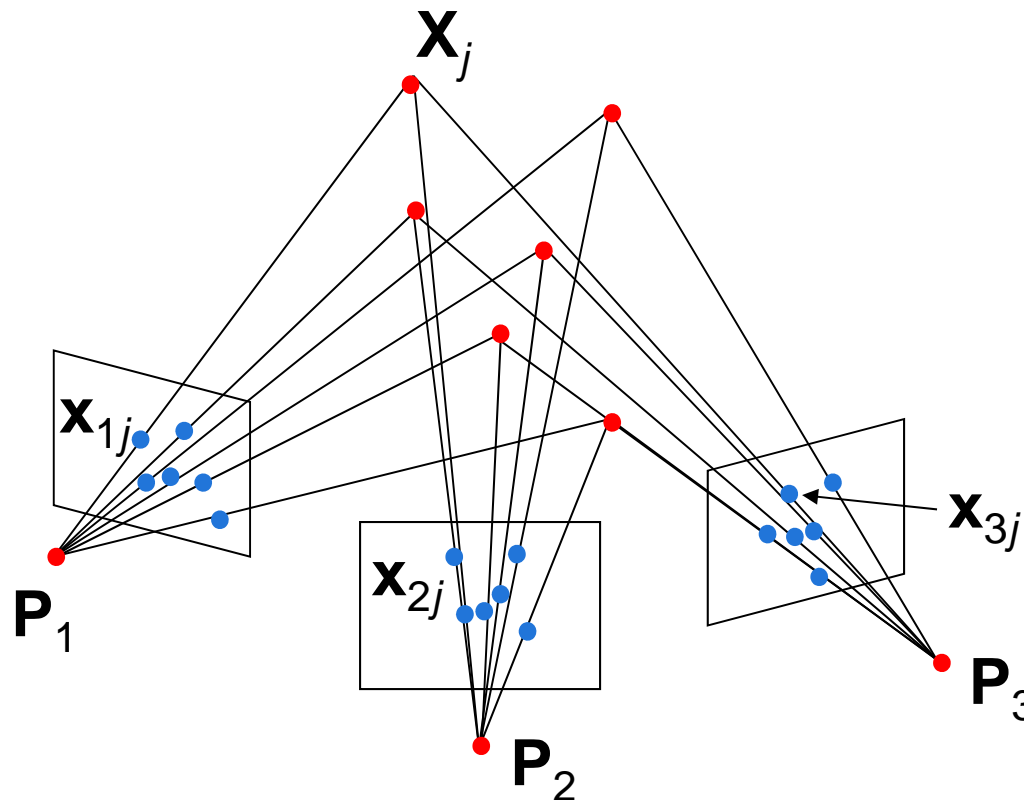
# Structure from motion

---

- Given:  $m$  images of  $n$  fixed 3D points

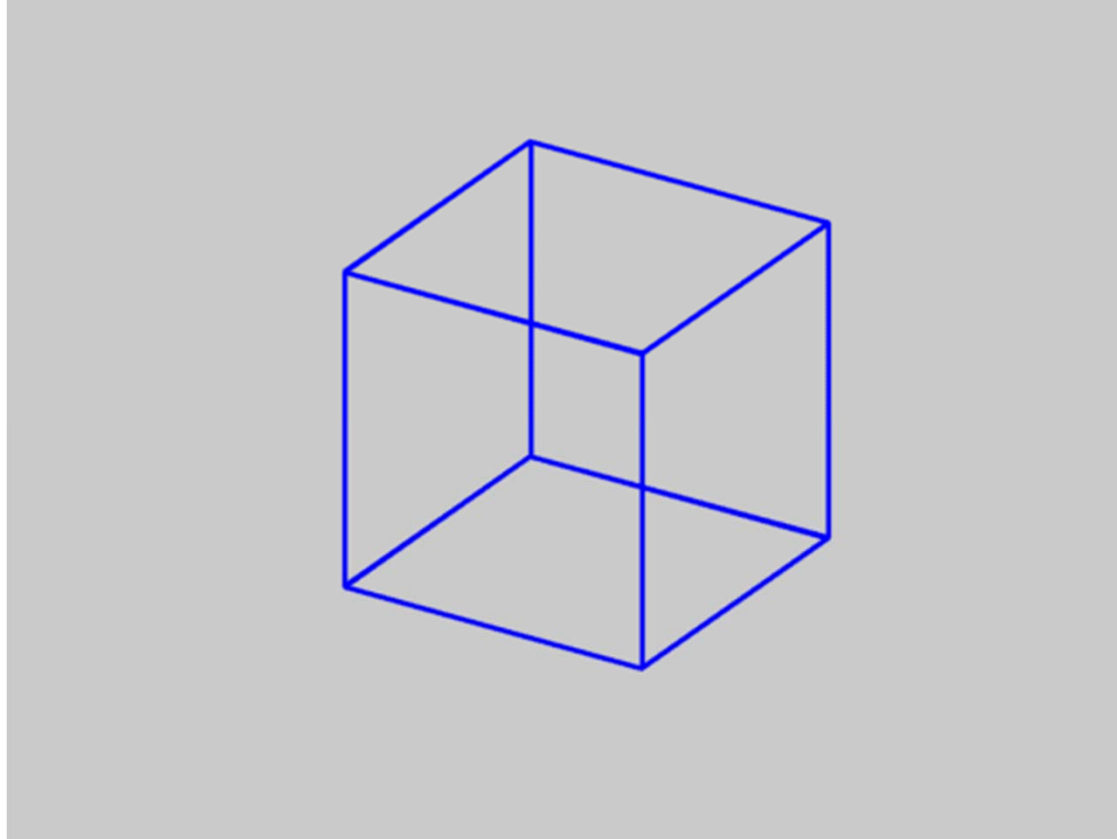
$$\lambda_{ij} \mathbf{x}_{ij} = \mathbf{P}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$

- 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 in the image remain exactly the same:

$$\mathbf{x} = \mathbf{P}\mathbf{X} = \left(\frac{1}{k}\mathbf{P}\right)(k\mathbf{X})$$

It is impossible to recover the absolute scale of the scene!



# 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 in the image remain exactly the same
- More generally, if we transform the scene using a transformation  $\mathbf{Q}$  and apply the inverse transformation to the camera matrices, then the images do not change:

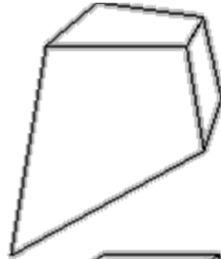
$$\mathbf{x} = \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}^{-1})(\mathbf{Q}\mathbf{X})$$

# Types of ambiguity

---

Projective  
15dof

$$\begin{bmatrix} A & t \\ v^T & v \end{bmatrix}$$



Preserves intersection and tangency

Affine  
12dof

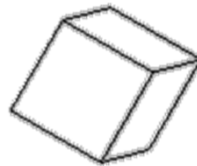
$$\begin{bmatrix} A & t \\ 0^T & 1 \end{bmatrix}$$



Preserves parallelism, volume ratios

Similarity  
7dof

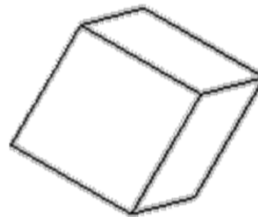
$$\begin{bmatrix} sR & t \\ 0^T & 1 \end{bmatrix}$$



Preserves angles, ratios of length

Euclidean  
6dof

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

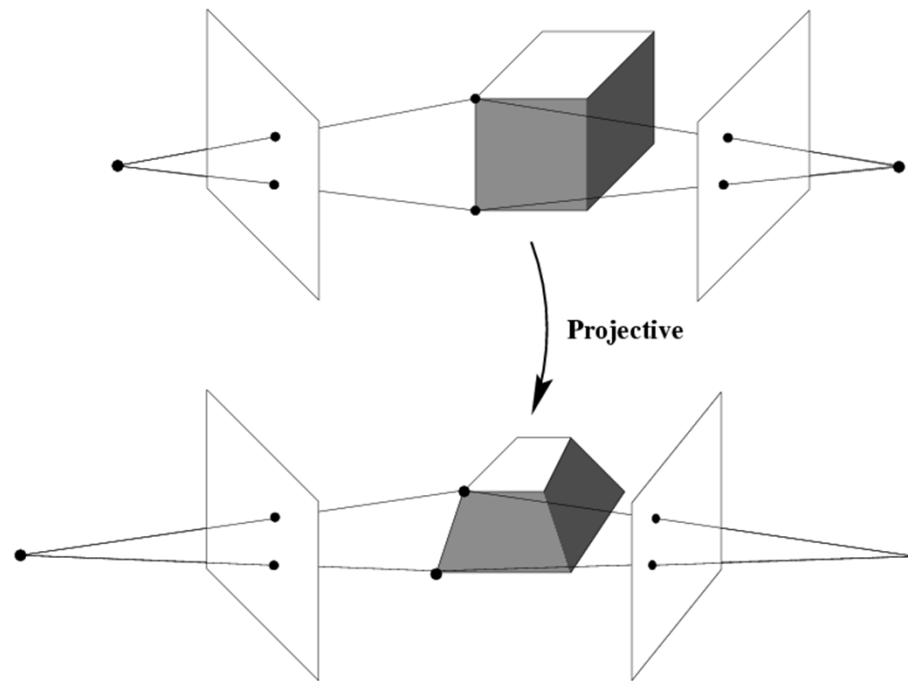


Preserves angles, lengths

- 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

---

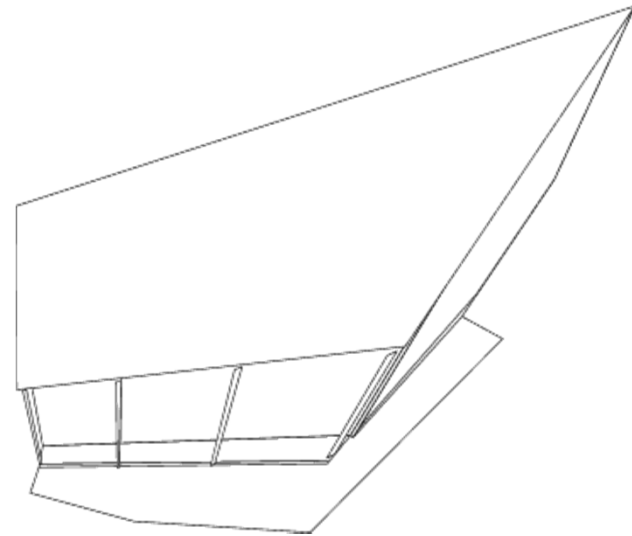
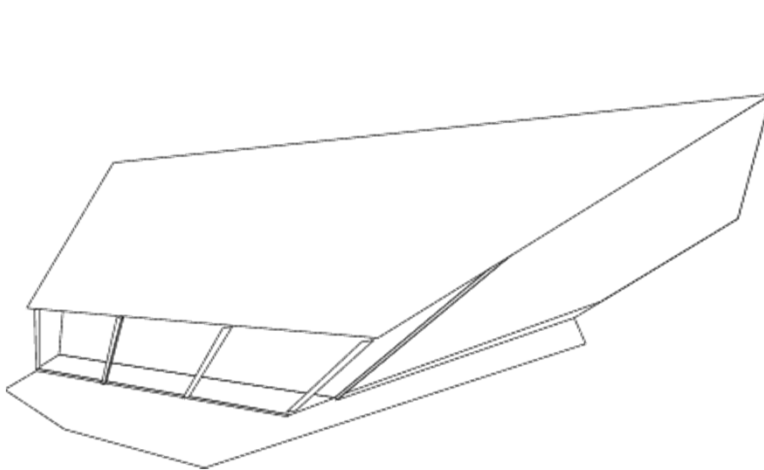


$$\mathbf{Q}_p = \begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{v}^\top & v \end{bmatrix}$$

$$\mathbf{x} = \mathbf{P}\mathbf{X} = \left(\mathbf{P}\mathbf{Q}_p^{-1}\right)(\mathbf{Q}_p \mathbf{X})$$

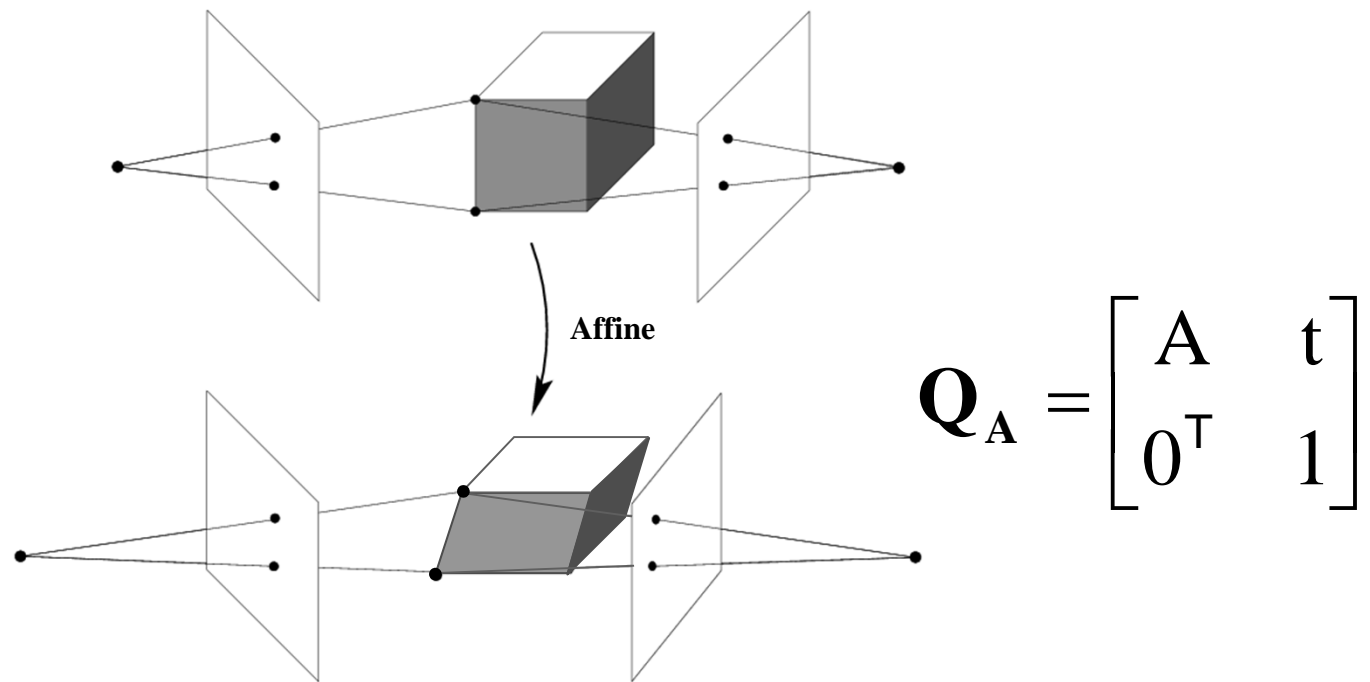
# Projective ambiguity

---



# Affine ambiguity

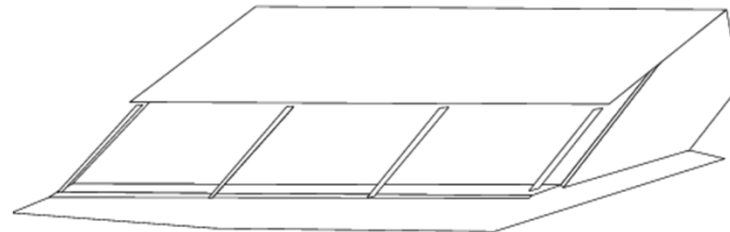
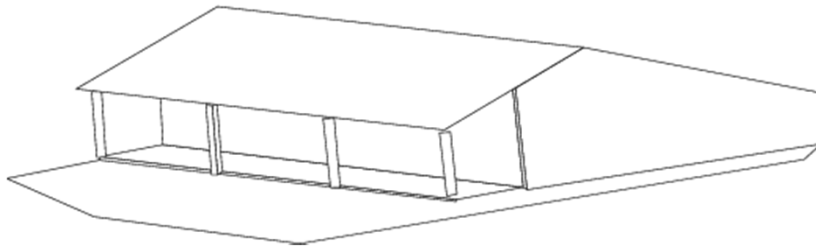
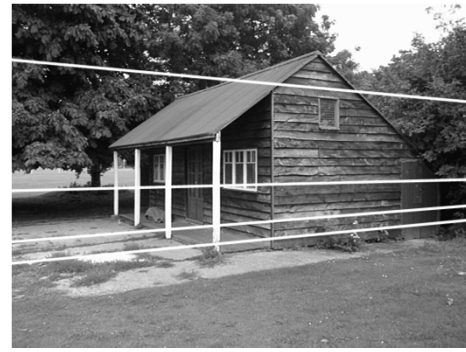
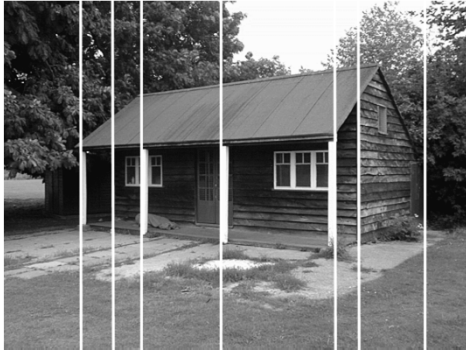
---



$$\mathbf{x} = \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}_A^{-1})(\mathbf{Q}_A \mathbf{X})$$

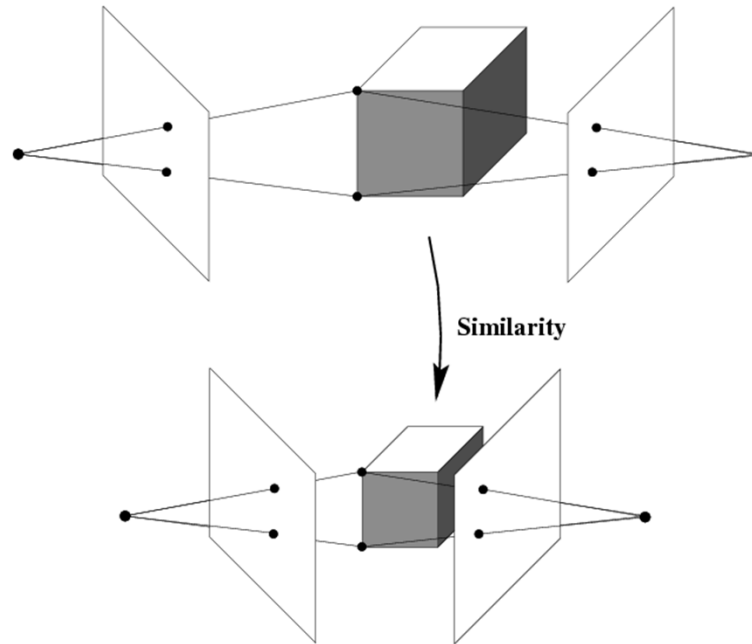
# Affine ambiguity

---



# Similarity ambiguity

---

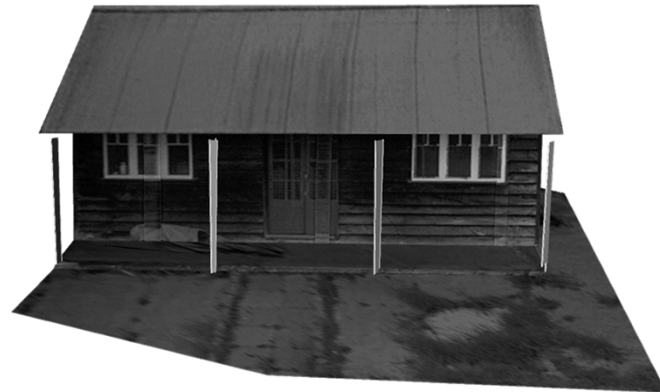
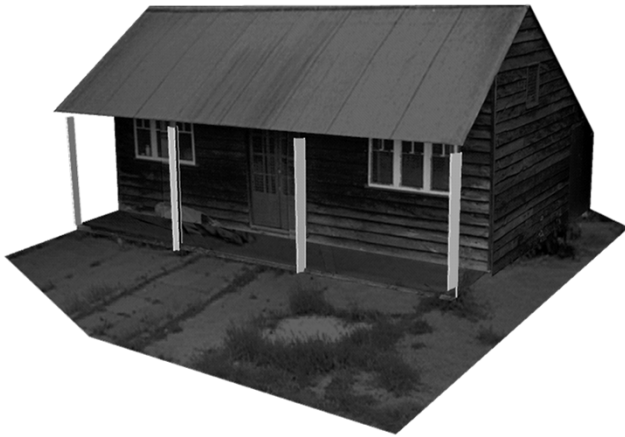
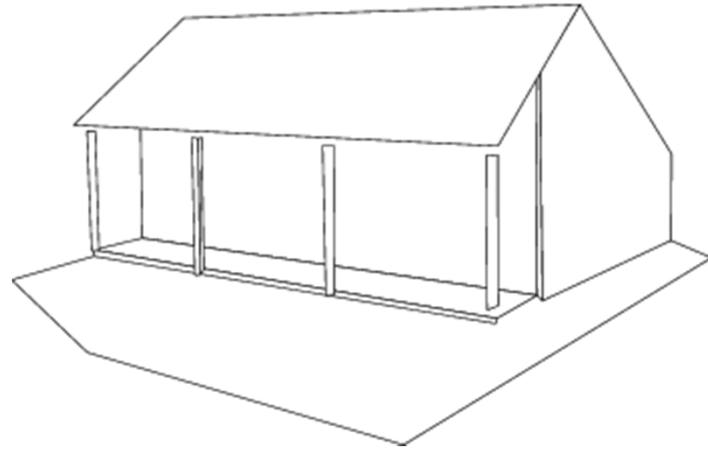
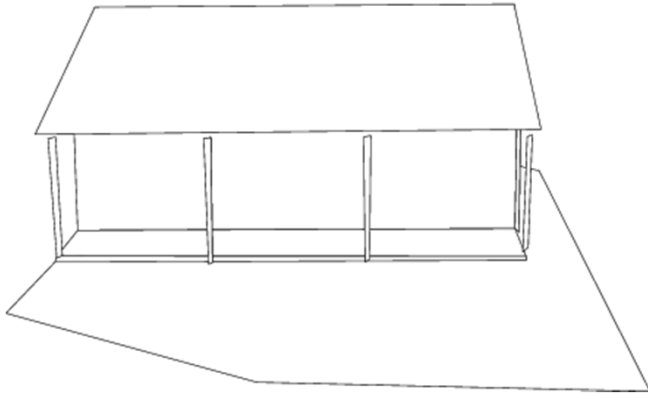


$$\mathbf{Q}_s = \begin{bmatrix} s\mathbf{R} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix}$$

$$\mathbf{x} = \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}_s^{-1})(\mathbf{Q}_s\mathbf{X})$$

# Similarity ambiguity

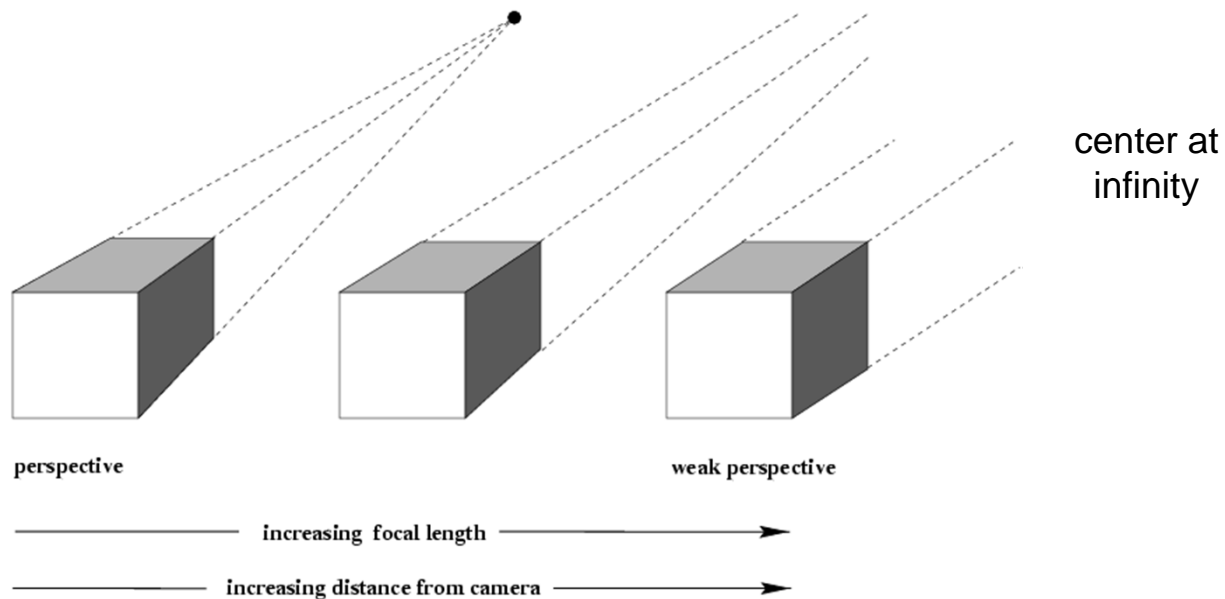
---





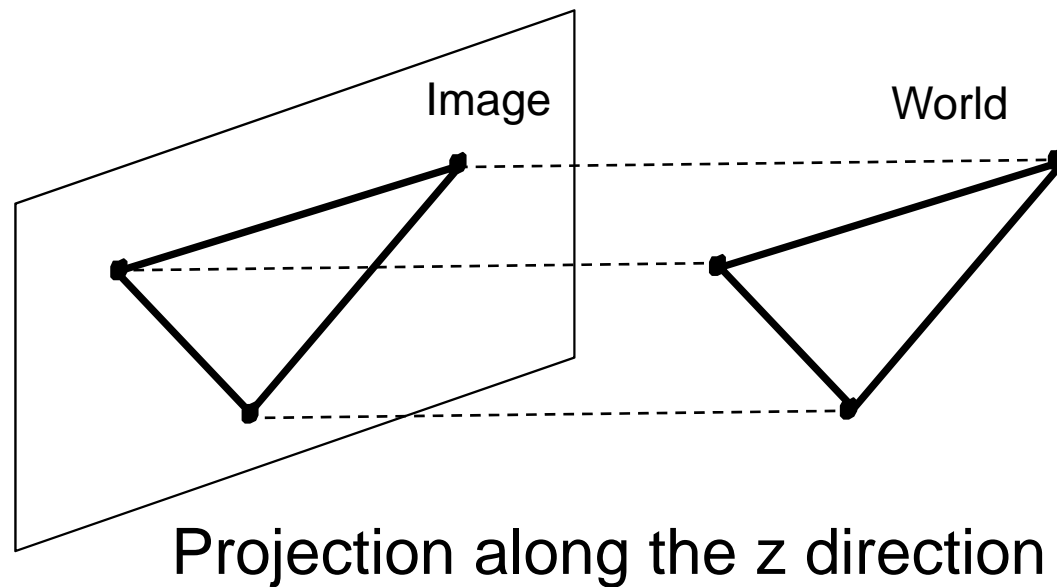
# Structure from motion

- Let's start with *affine* or *weak perspective* cameras (the math is easier)



# Recall: Orthographic Projection

---

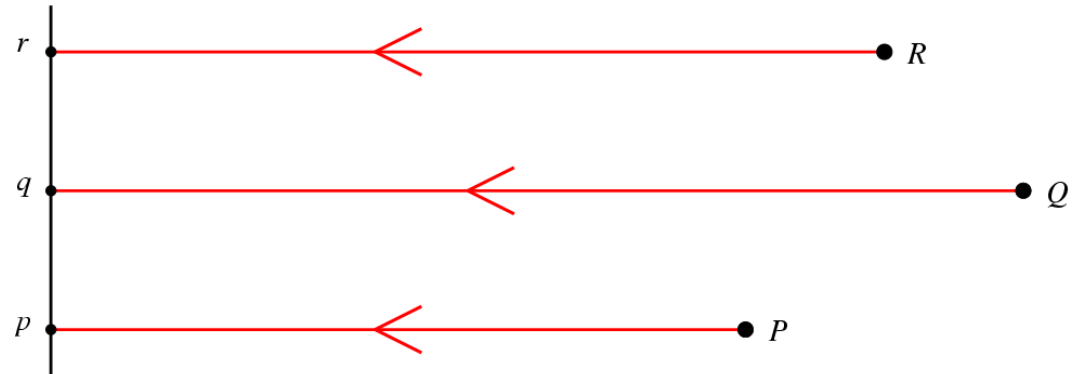


$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \Rightarrow (x, y)$$

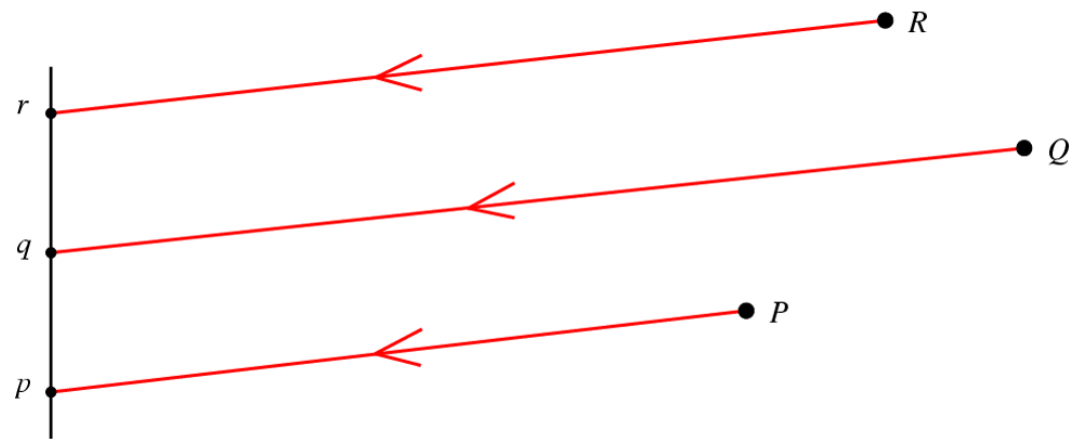
# Affine cameras

---

Orthographic Projection



Parallel Projection



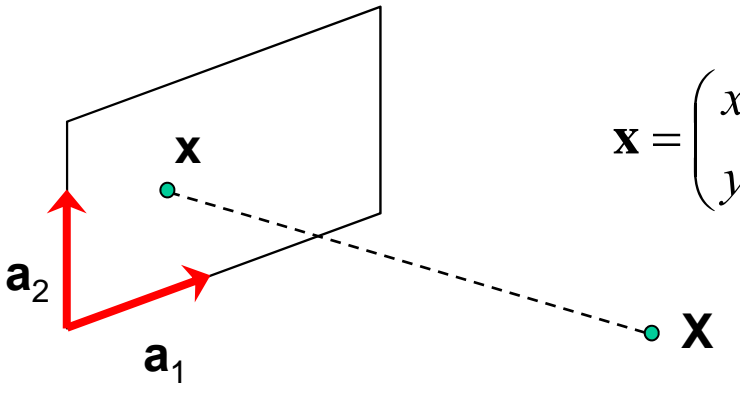
# Affine cameras

---

- A general affine camera combines the effects of an affine transformation of the 3D space, orthographic projection, and an affine transformation of the image:

$$\mathbf{P} = [3 \times 3 \text{ affine}] \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} [4 \times 4 \text{ affine}] = \begin{bmatrix} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & 1 \end{bmatrix}$$

- Affine projection is a linear mapping + translation in non-homogeneous coordinates



$$\mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = \mathbf{A}\mathbf{X} + \mathbf{b}$$

Projection of world origin

# Affine structure from motion

---

- Given:  $m$  images of  $n$  fixed 3D points:

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

- Problem: use the  $mn$  correspondences  $\mathbf{x}_{ij}$  to estimate  $m$  projection matrices  $\mathbf{A}_i$  and translation vectors  $\mathbf{b}_i$ , and  $n$  points  $\mathbf{X}_j$
- The reconstruction is defined up to an arbitrary *affine* transformation  $\mathbf{Q}$  (12 degrees of freedom):

$$\begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & 1 \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & 1 \end{bmatrix} \mathbf{Q}^{-1}, \quad \begin{pmatrix} \mathbf{X} \\ 1 \end{pmatrix} \rightarrow \mathbf{Q} \begin{pmatrix} \mathbf{X} \\ 1 \end{pmatrix}$$

- We have  $2mn$  knowns and  $8m + 3n$  unknowns (minus 12 dof for affine ambiguity)
- Thus, we must have  $2mn \geq 8m + 3n - 12$
- For two views, we need four point correspondences

# Affine structure from motion

---

- Centering: subtract the centroid of the image points in each view

$$\begin{aligned}\hat{\mathbf{x}}_{ij} &= \mathbf{x}_{ij} - \frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ik} = \mathbf{A}_i \mathbf{X}_j + \mathbf{b}_i - \frac{1}{n} \sum_{k=1}^n (\mathbf{A}_i \mathbf{X}_k + \mathbf{b}_i) \\ &= \mathbf{A}_i \left( \mathbf{X}_j - \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k \right) = \mathbf{A}_i \hat{\mathbf{X}}_j\end{aligned}$$

- For simplicity, set the origin of the world coordinate system to the centroid of the 3D points
- After centering, each normalized 2D point is related to the 3D point  $\mathbf{X}_j$  by

$$\hat{\mathbf{x}}_{ij} = \mathbf{A}_i \mathbf{X}_j$$

# Affine structure from motion

---

- Let's create a  $2m \times n$  data (measurement) matrix:

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \cdots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \cdots & \hat{\mathbf{x}}_{2n} \\ & & \ddots & \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \cdots & \hat{\mathbf{x}}_{mn} \end{bmatrix}$$

↓ cameras ( $2m$ )

→ points ( $n$ )

C. Tomasi and T. Kanade. [Shape and motion from image streams under orthography: A factorization method.](#) *IJCV*, 9(2):137-154, November 1992.

# Affine structure from motion

---

- Let's create a  $2m \times n$  data (measurement) matrix:

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \cdots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \cdots & \hat{\mathbf{x}}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \cdots & \hat{\mathbf{x}}_{mn} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \vdots \\ \mathbf{A}_m \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_n \end{bmatrix}$$

points ( $3 \times n$ )

cameras  
( $2m \times 3$ )

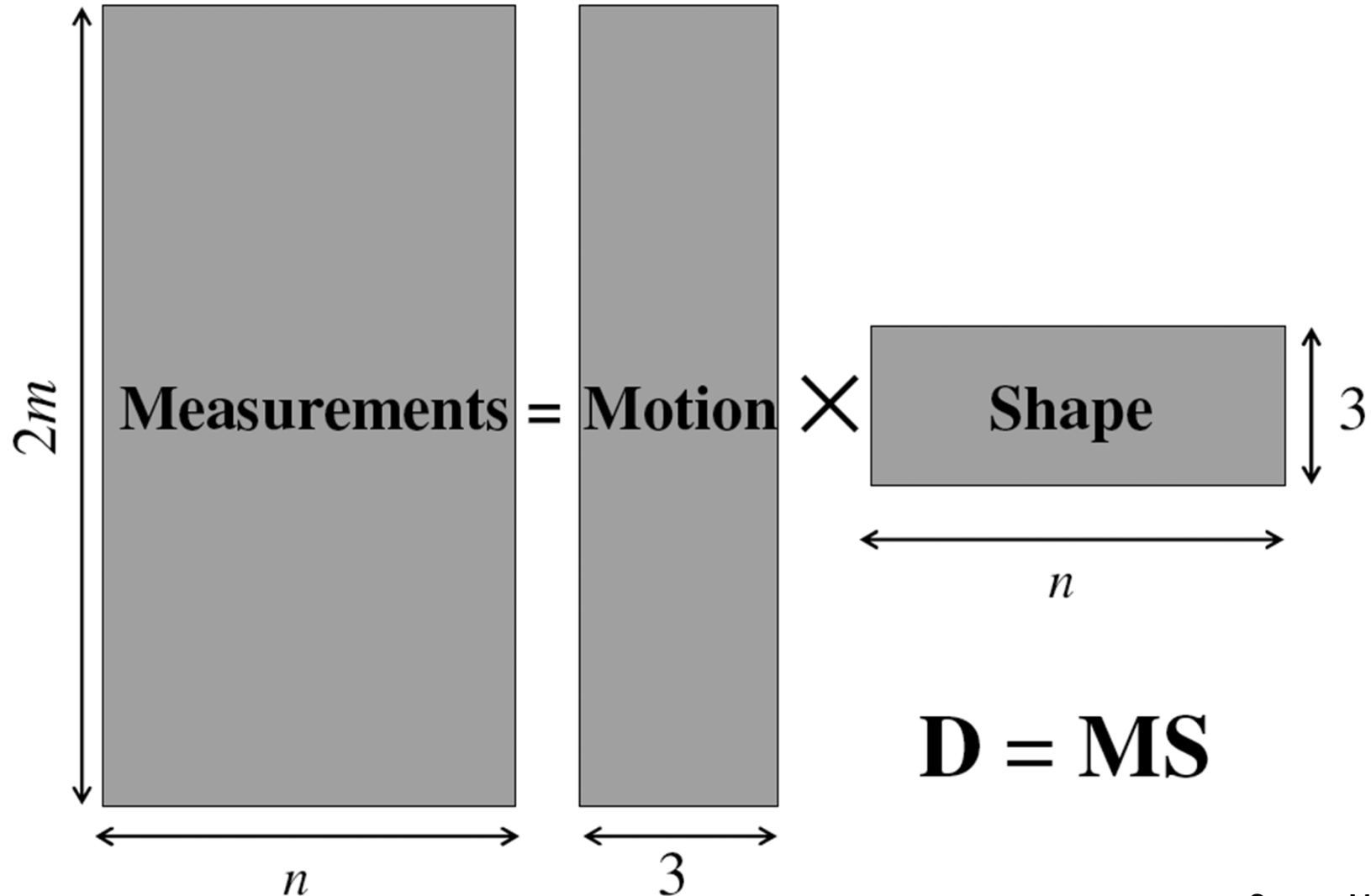
The measurement matrix  $\mathbf{D} = \mathbf{MS}$  must have rank 3!

C. Tomasi and T. Kanade. [Shape and motion from image streams under orthography: A factorization method.](#) *IJCV*, 9(2):137-154, November 1992.



# Factorizing the measurement matrix

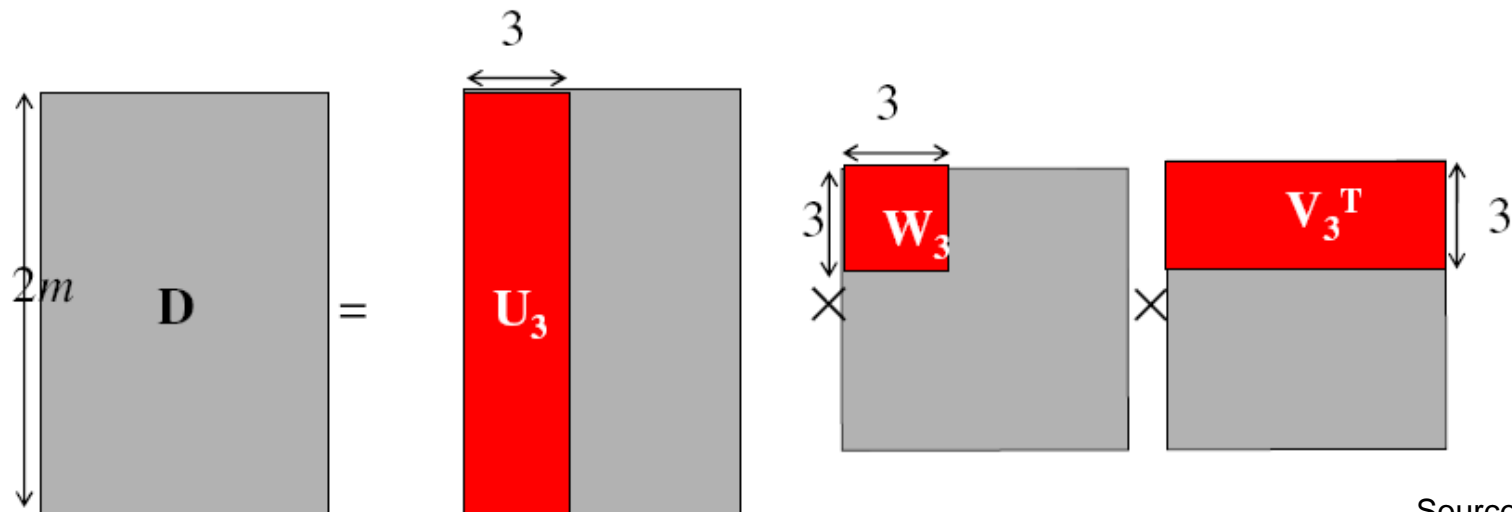
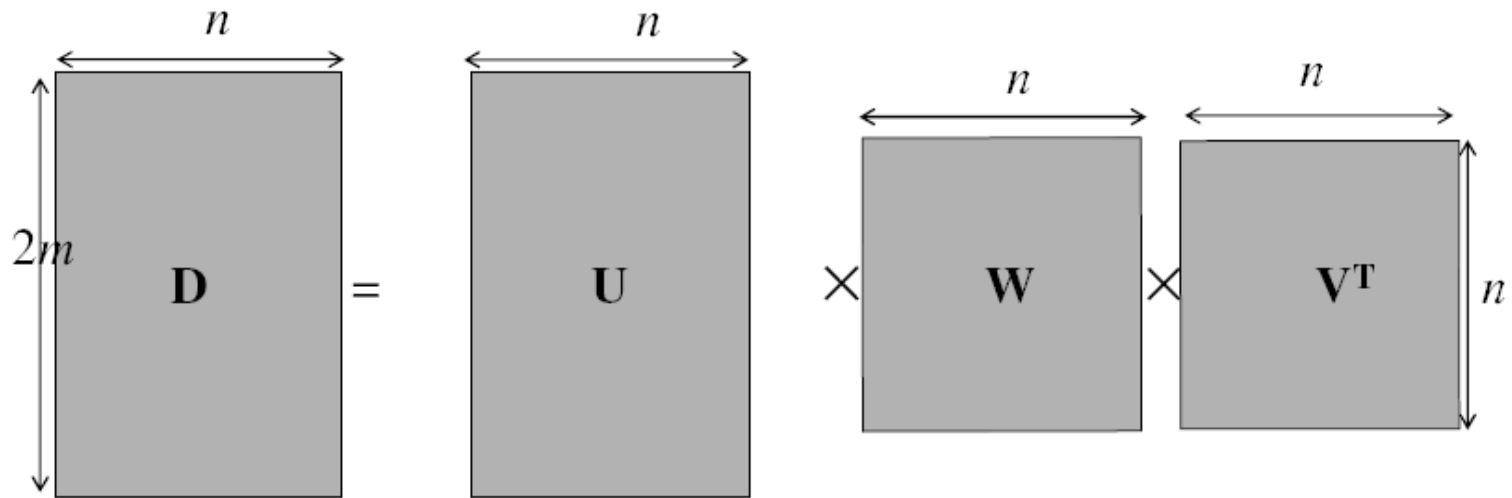
---



# Factorizing the measurement matrix

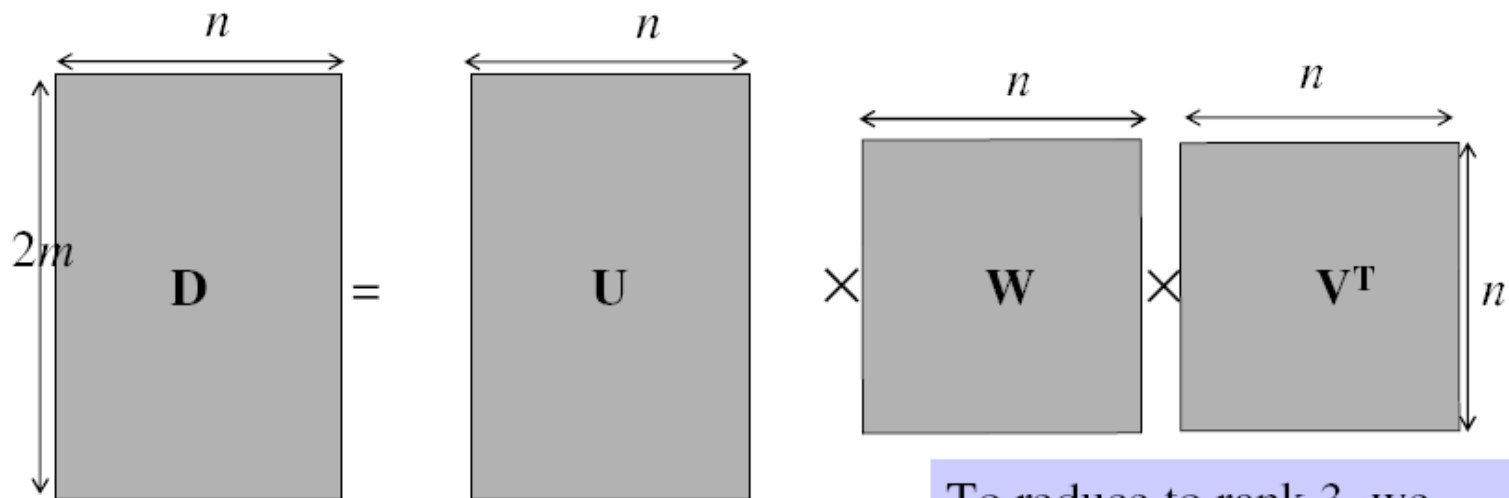
---

- Singular value decomposition of  $D$ :

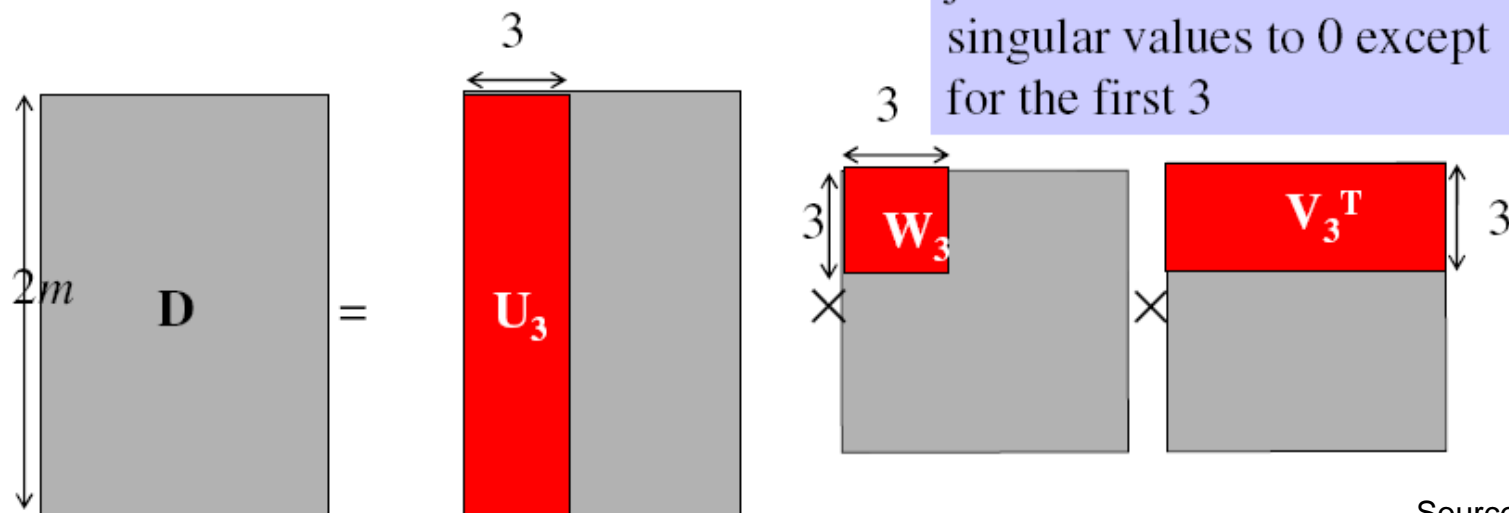


# Factorizing the measurement matrix

- Singular value decomposition of  $D$ :



To reduce to rank 3, we just need to set all the singular values to 0 except for the first 3



# Factorizing the measurement matrix

---

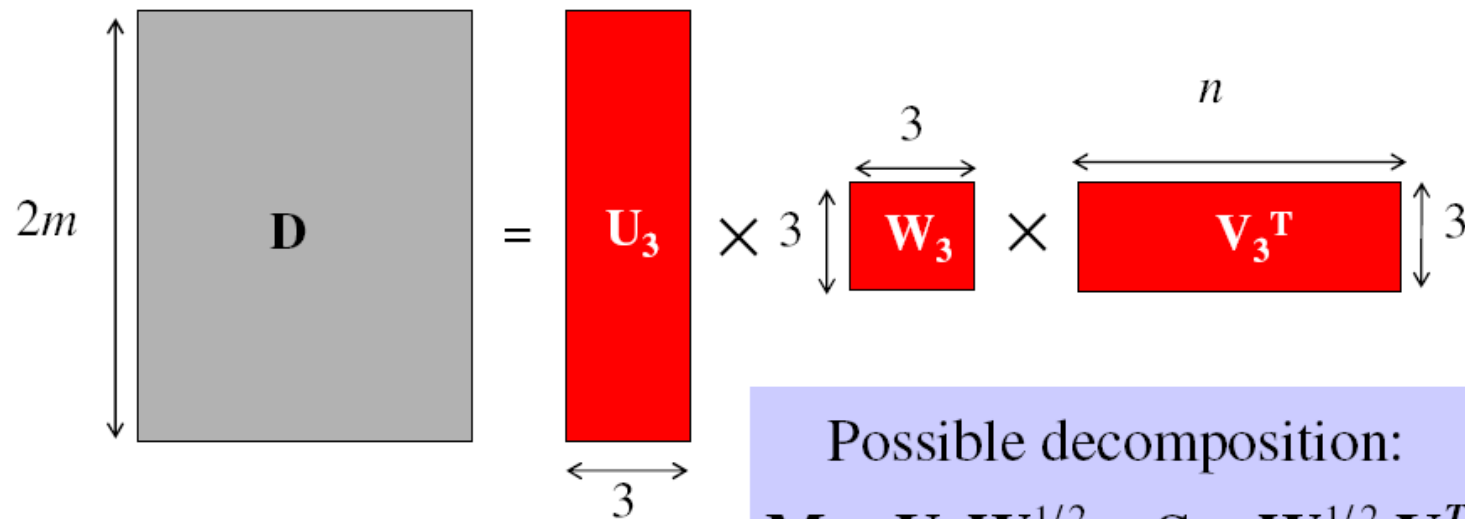
- Obtaining a factorization from SVD:

$$\begin{array}{c} 2m \\ \updownarrow \\ \text{D} \end{array} = \begin{array}{c} \text{U}_3 \\ \leftarrow 3 \end{array} \times \begin{array}{c} 3 \\ \leftarrow \\ \text{W}_3 \\ \updownarrow 3 \end{array} \times \begin{array}{c} n \\ \leftarrow \\ \text{V}_3^T \\ \updownarrow 3 \end{array}$$

# Factorizing the measurement matrix

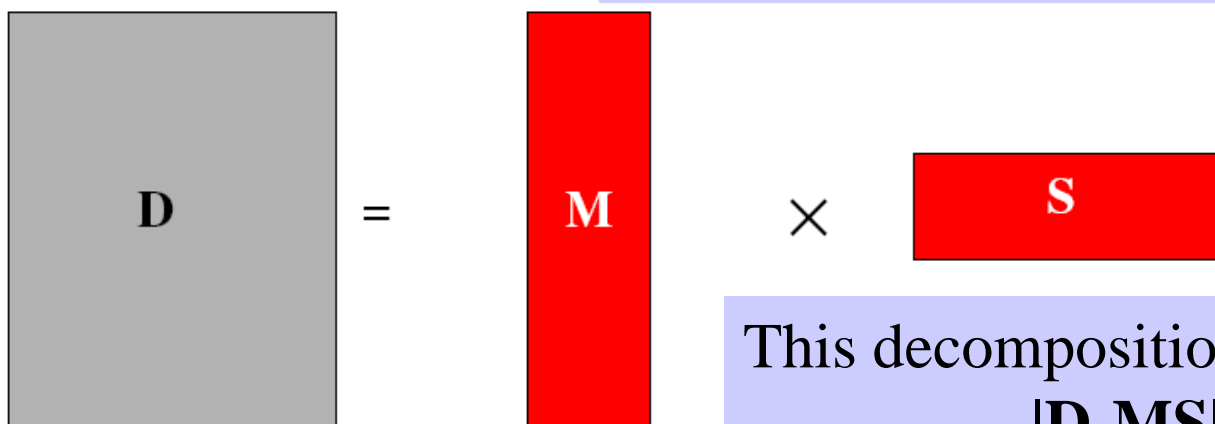
---

- Obtaining a factorization from SVD:



Possible decomposition:

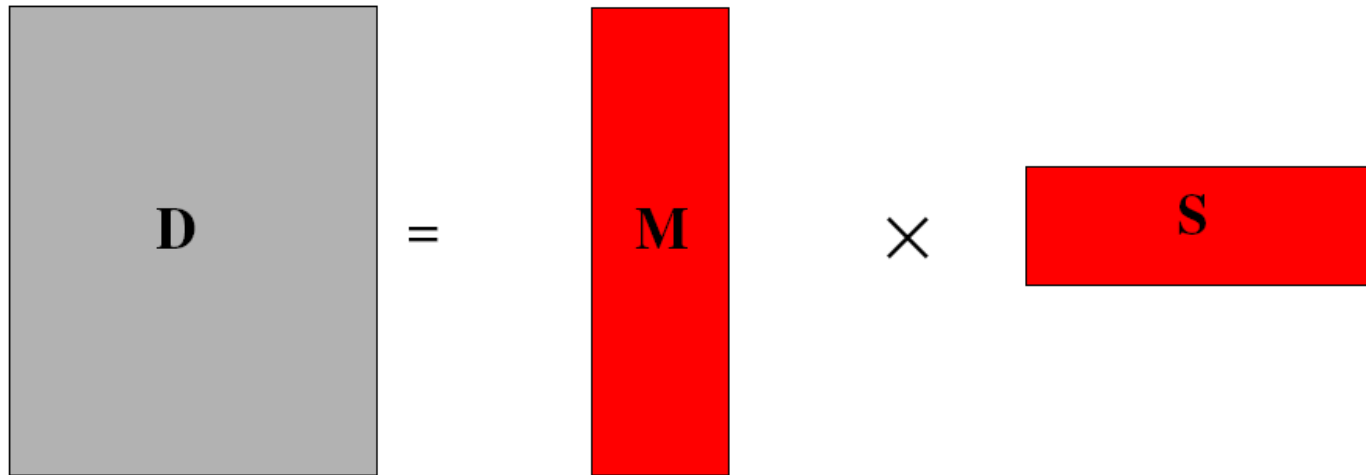
$$M = U_3 W_3^{1/2} \quad S = W_3^{1/2} V_3^T$$



This decomposition minimizes  $|D - MS|^2$

# Affine ambiguity

---

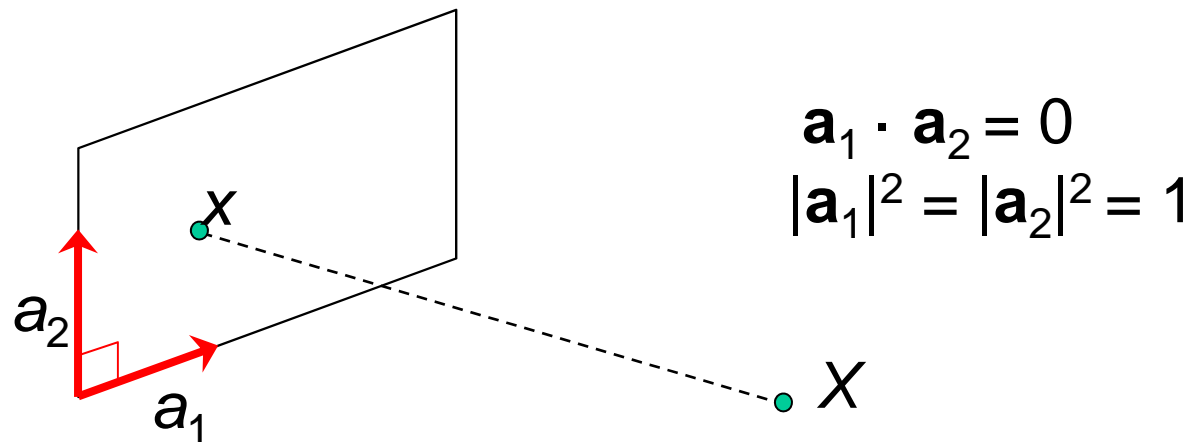

$$\mathbf{D} = \mathbf{M} \times \mathbf{S}$$

- The decomposition is not unique. We get the same  $\mathbf{D}$  by using any  $3 \times 3$  matrix  $\mathbf{C}$  and applying the transformations  $\mathbf{M} \rightarrow \mathbf{MC}$ ,  $\mathbf{S} \rightarrow \mathbf{C}^{-1}\mathbf{S}$
- That is because we have only an affine transformation and we have not enforced any Euclidean constraints (like forcing the image axes to be perpendicular, for example)

# Eliminating the affine ambiguity

---

- Transform each projection matrix  $A$  to another matrix  $AC$  to get orthographic projection
  - Image axes are perpendicular and scale is 1



- This translates into  $3m$  equations:  
$$(\mathbf{A}_i \mathbf{C})(\mathbf{A}_i \mathbf{C})^T = \mathbf{A}_i (\mathbf{C} \mathbf{C}^T) \mathbf{A}_i^T = \mathbf{I}_d, \quad i = 1, \dots, m$$
  - Solve for  $\mathbf{L} = \mathbf{C} \mathbf{C}^T$
  - Recover  $\mathbf{C}$  from  $\mathbf{L}$  by Cholesky decomposition:  $\mathbf{L} = \mathbf{C} \mathbf{C}^T$
  - Update  $\mathbf{M}$  and  $\mathbf{S}$ :  $\mathbf{M} = \mathbf{M} \mathbf{C}$ ,  $\mathbf{S} = \mathbf{C}^{-1} \mathbf{S}$

# Reconstruction results

---



1



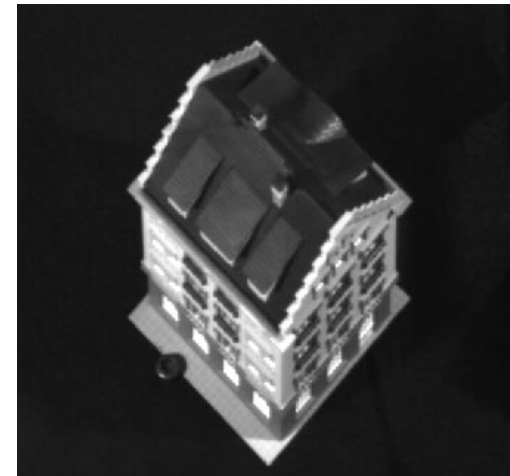
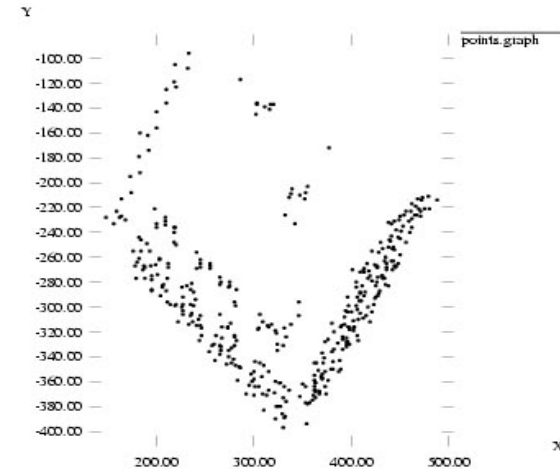
60



120



150



C. Tomasi and T. Kanade. [Shape and motion from image streams under orthography: A factorization method](#). *IJCV*, 9(2):137-154, November 1992.



# Algorithm summary

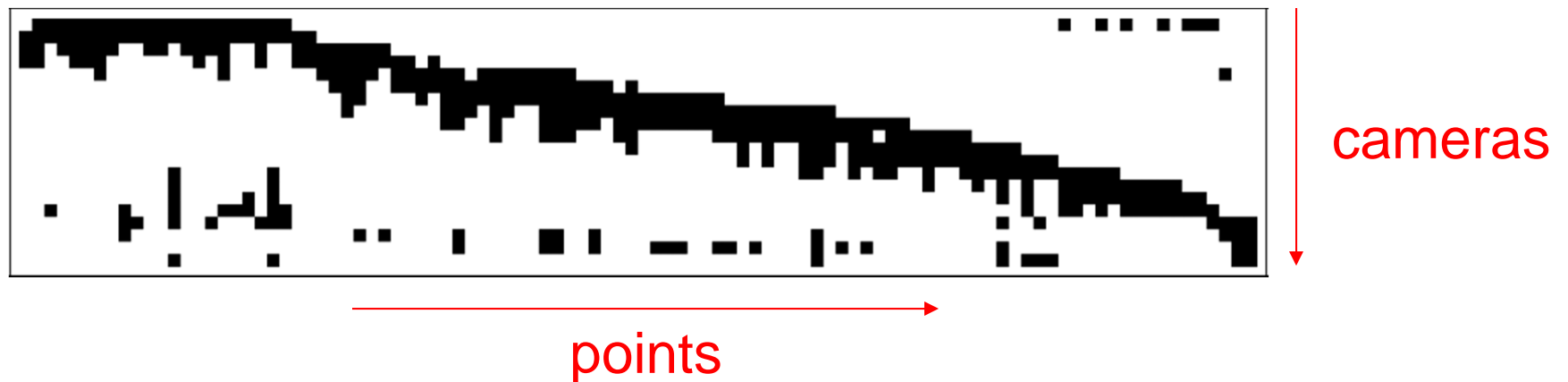
---

- Given:  $m$  images and  $n$  features  $\mathbf{x}_{ij}$
- For each image  $i$ , center the feature coordinates
- Construct a  $2m \times n$  measurement matrix  $\mathbf{D}$ :
  - Column  $j$  contains the projection of point  $j$  in all views
  - Row  $i$  contains one coordinate of the projections of all the  $n$  points in image  $i$
- Factorize  $\mathbf{D}$ :
  - Compute SVD:  $\mathbf{D} = \mathbf{U} \mathbf{W} \mathbf{V}^T$
  - Create  $\mathbf{U}_3$  by taking the first 3 columns of  $\mathbf{U}$
  - Create  $\mathbf{V}_3$  by taking the first 3 columns of  $\mathbf{V}$
  - Create  $\mathbf{W}_3$  by taking the upper left  $3 \times 3$  block of  $\mathbf{W}$
- Create the motion and shape matrices:
  - $\mathbf{M} = \mathbf{U}_3 \mathbf{W}_3^{1/2}$  and  $\mathbf{S} = \mathbf{W}_3^{1/2} \mathbf{V}_3^T$  (or  $\mathbf{M} = \mathbf{U}_3$  and  $\mathbf{S} = \mathbf{W}_3 \mathbf{V}_3^T$ )
- Eliminate affine ambiguity

# Dealing with missing data

---

- So far, we have assumed that all points are visible in all views
- In reality, the measurement matrix typically looks something like this:

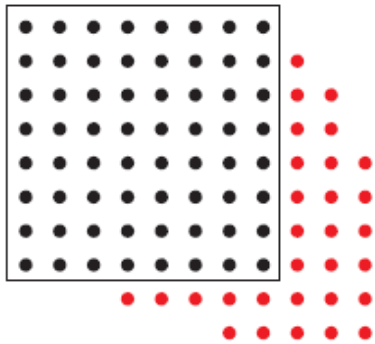


- Possible solution: decompose matrix into dense sub-blocks, factorize each sub-block, and fuse the results
  - Finding dense maximal sub-blocks of the matrix is NP-complete (equivalent to finding maximal cliques in a graph)

# Dealing with missing data

---

- Incremental bilinear refinement



(1) Perform factorization on a dense sub-block

(2) Solve for a new 3D point visible by at least two known cameras (linear least squares)

(3) Solve for a new camera that sees at least three known 3D points (linear least squares)

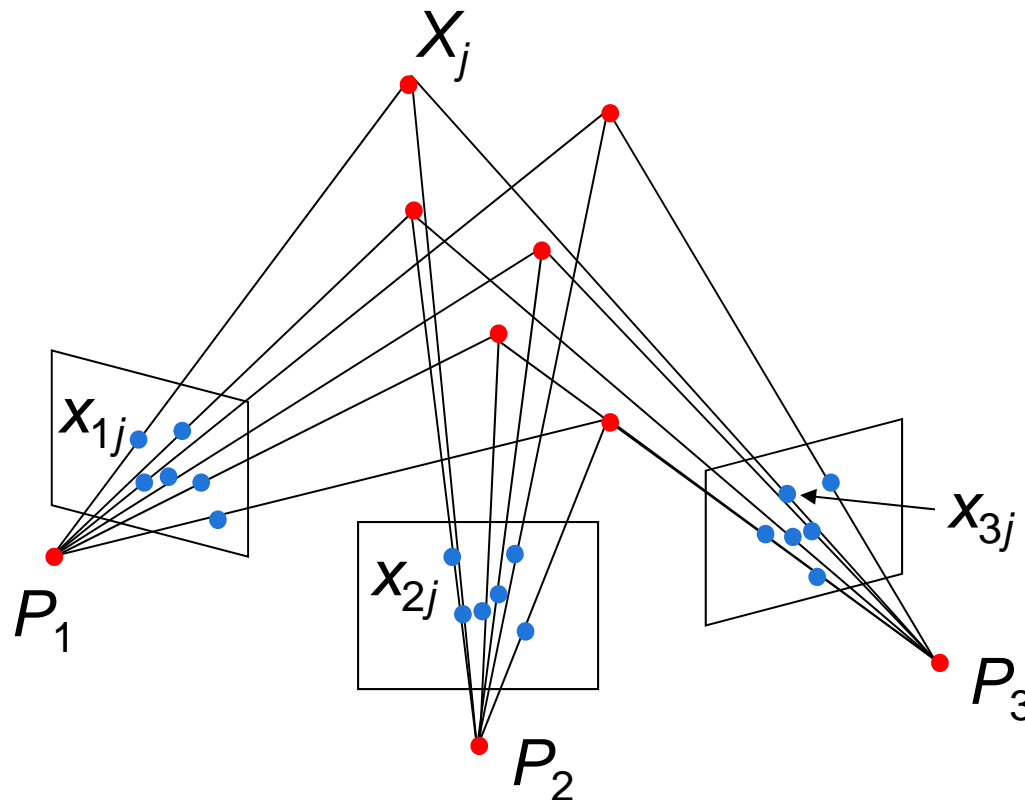
# Projective structure from motion

---

- Given:  $m$  images of  $n$  fixed 3D points

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

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



# Projective structure from motion

---

- Given:  $m$  images of  $n$  fixed 3D points

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

- 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 4x4 projective transformation  $\mathbf{Q}$ :

$$\mathbf{X} \rightarrow \mathbf{QX}, \mathbf{P} \rightarrow \mathbf{PQ}^{-1}$$

- We can solve for structure and motion when

$$2mn \geq 11m + 3n - 15$$

- For two cameras, at least 7 points are needed

# Projective SFM: Two-camera case

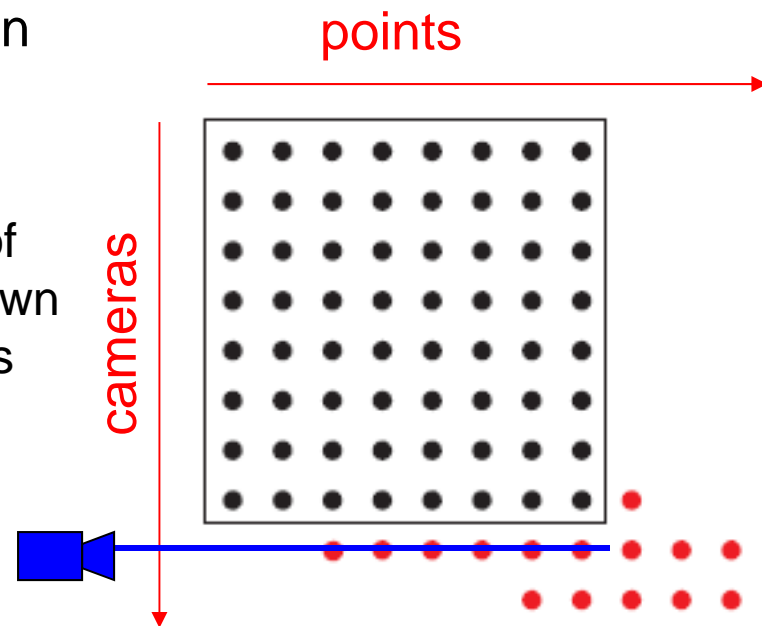
---

- Compute fundamental matrix  $\mathbf{F}$  between the two views
- First camera matrix:  $[\mathbf{I} \mid \mathbf{0}]$
- Second camera matrix:  $[\mathbf{A} \mid \mathbf{b}]$
- Then  $\mathbf{b}$  is the epipole ( $\mathbf{F}^T \mathbf{b} = 0$ ),  $\mathbf{A} = -[\mathbf{b}_x] \mathbf{F}$

# Sequential structure from motion

---

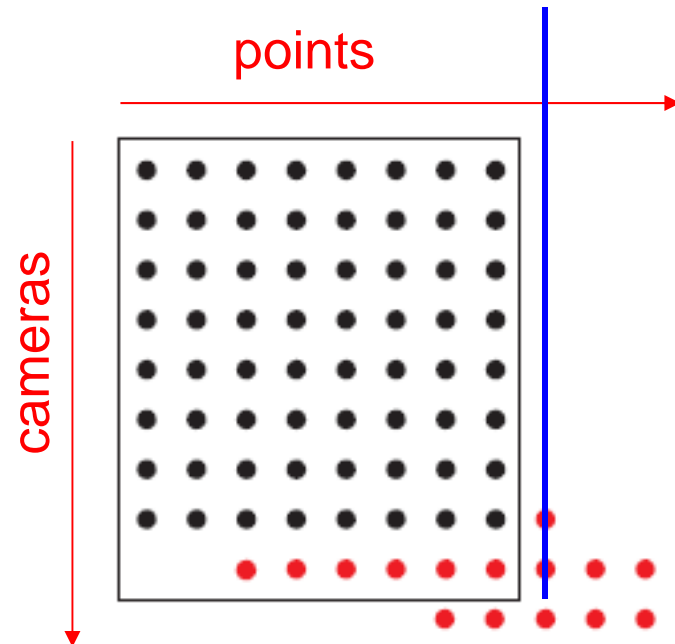
- Initialize motion from two images using fundamental matrix
- Initialize structure by triangulation
- For each additional view:
  - Determine projection matrix of new camera using all the known 3D points that are visible in its image – *calibration*



# Sequential structure from motion

---

- Initialize motion from two images using 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 new 3D points, re-optimize existing points that are also seen by this camera – *triangulation*

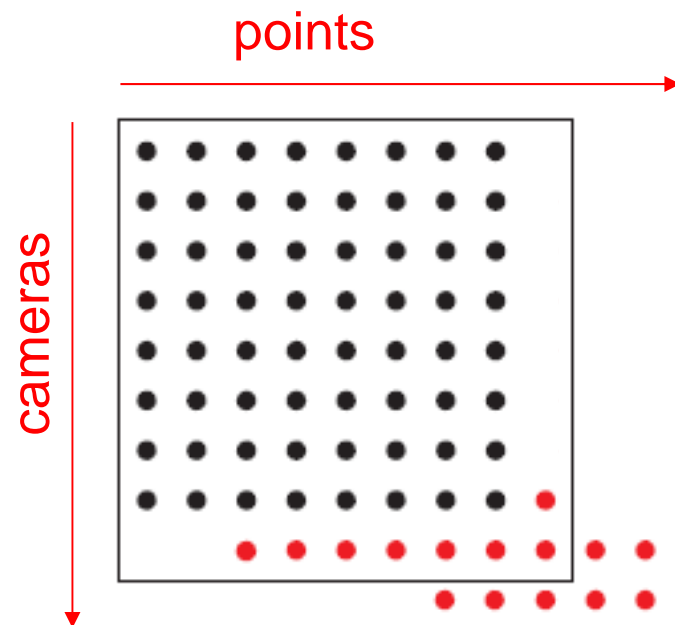




# Sequential structure from motion

---

- Initialize motion from two images using 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 new 3D points, re-optimize existing points that are also seen by this camera – *triangulation*
- Refine structure and motion: bundle adjustment



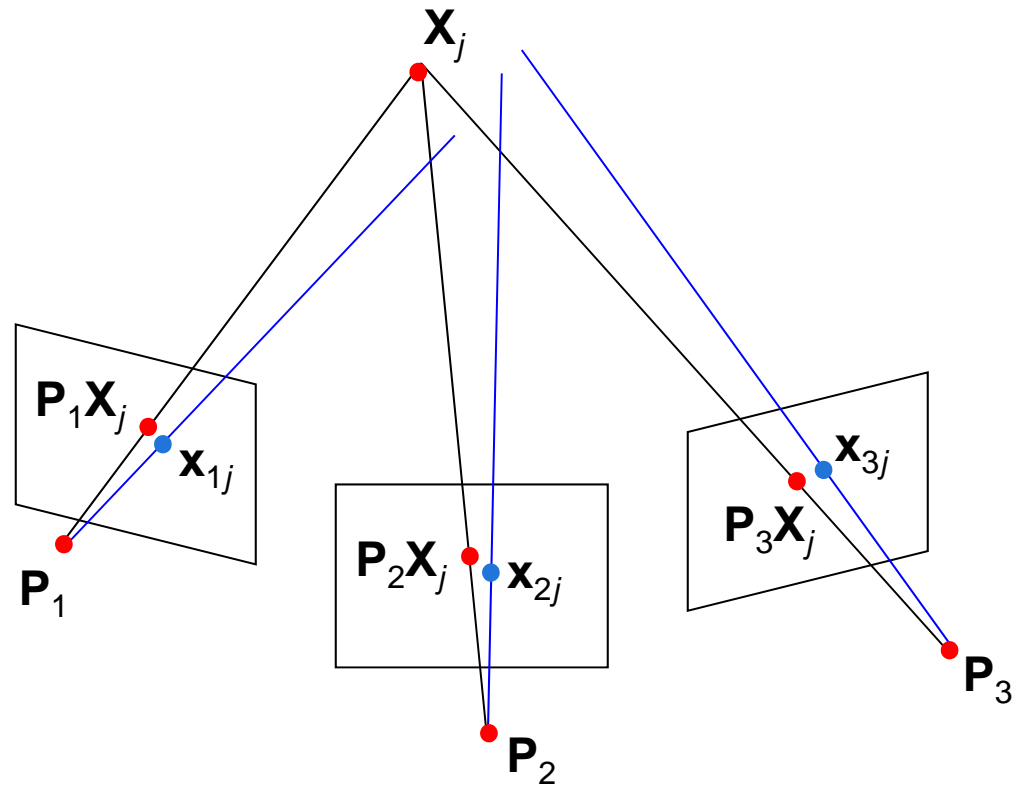
# Bundle adjustment

---

- Non-linear method for refining structure and motion
- Minimize reprojection error

$$\sum_{i=1}^m \sum_{j=1}^n w_{ij} \left\| \mathbf{x}_{ij} - \frac{1}{\lambda_{ij}} \mathbf{P}_i \mathbf{X}_j \right\|^2$$

visibility flag:  
is point  $j$   
visible in  
view  $i$ ?



# Self-calibration

---

- Self-calibration (auto-calibration) is the process of determining intrinsic camera parameters directly from uncalibrated images
- For example, when the images are acquired by a single moving camera, we can use the constraint that the intrinsic parameter matrix remains fixed for all the images
  - Compute initial projective reconstruction and find 3D projective transformation matrix  $\mathbf{Q}$  such that all camera matrices are in the form  $\mathbf{P}_i = \mathbf{K} [\mathbf{R}_i | \mathbf{t}_i]$
- Can use constraints on the form of the calibration matrix: zero skew
- Can use vanishing points

# Modern SFM pipeline

---

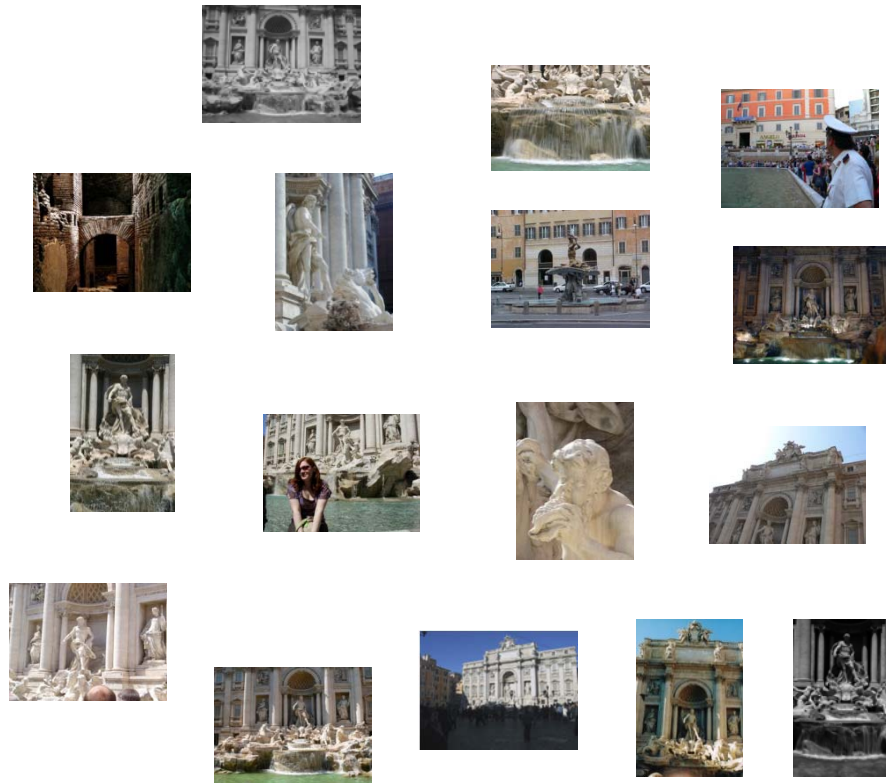


N. Snavely, S. Seitz, and R. Szeliski, ["Photo tourism: Exploring photo collections in 3D,"](#)  
SIGGRAPH 2006.

# Feature detection

---

Detect features using SIFT [Lowe, IJCV 2004]

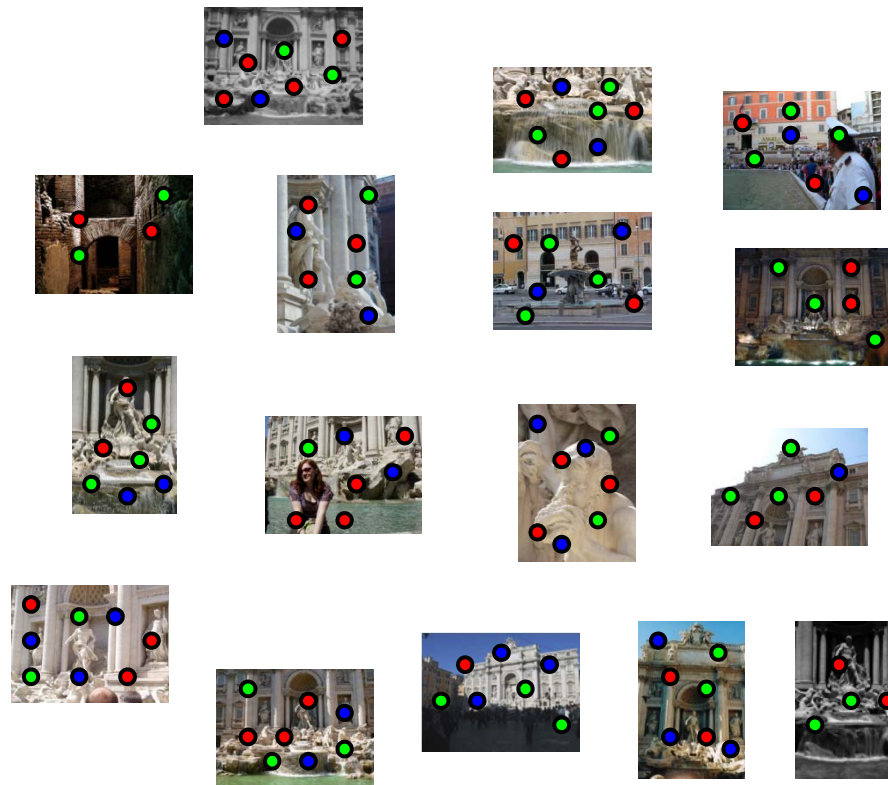


Source: N. Snavely

# Feature detection

---

Detect features using SIFT

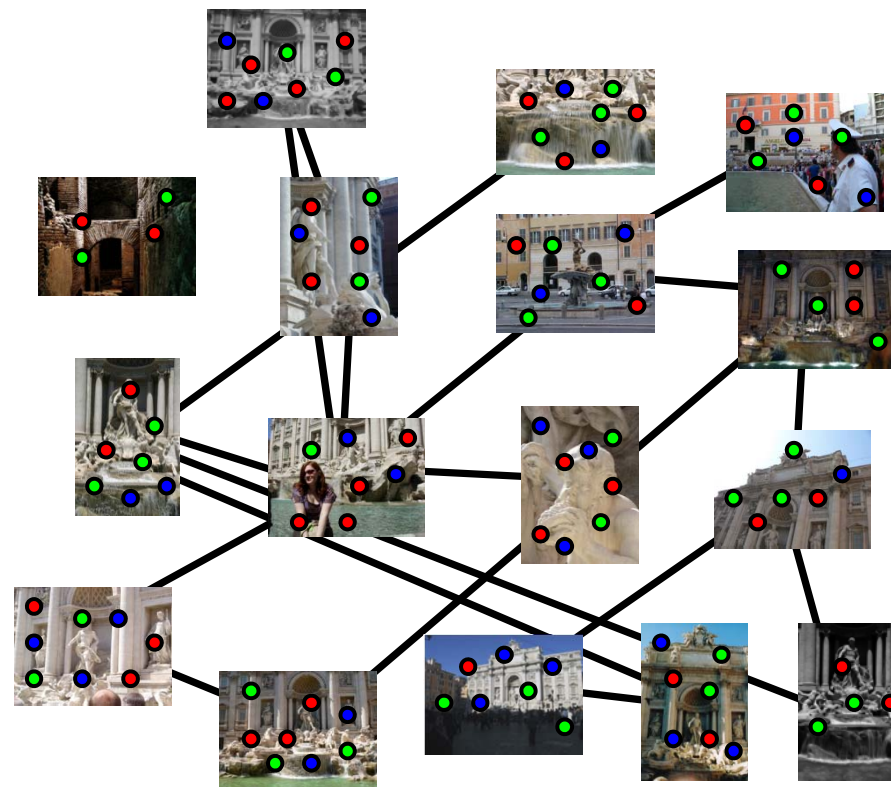


Source: N. Snavely

# Feature matching

---

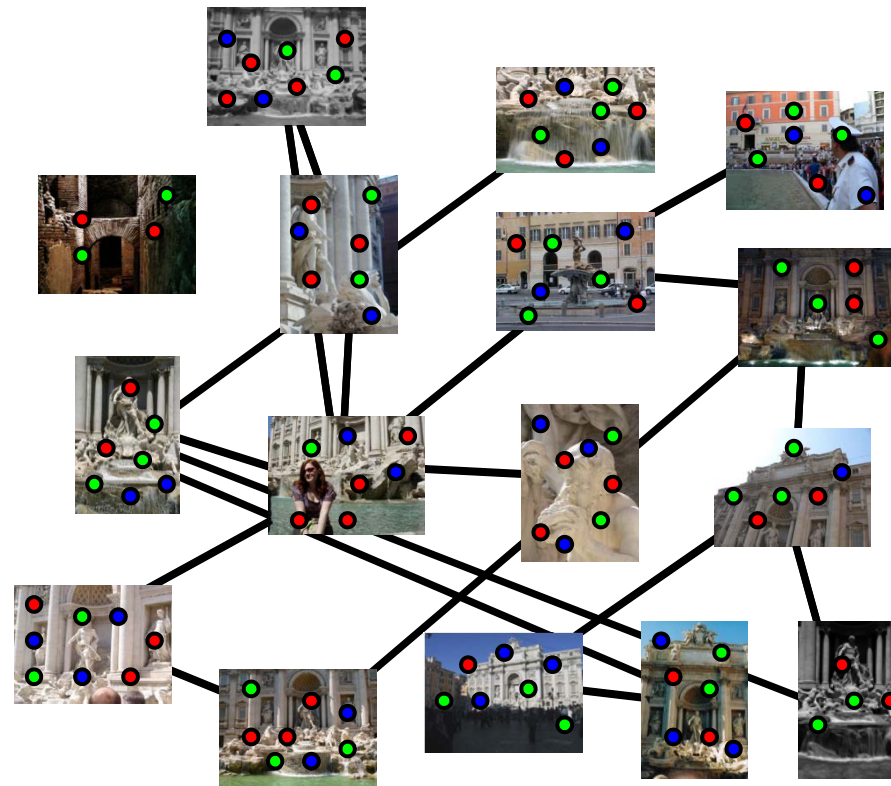
Match features between each pair of images



# Feature matching

---

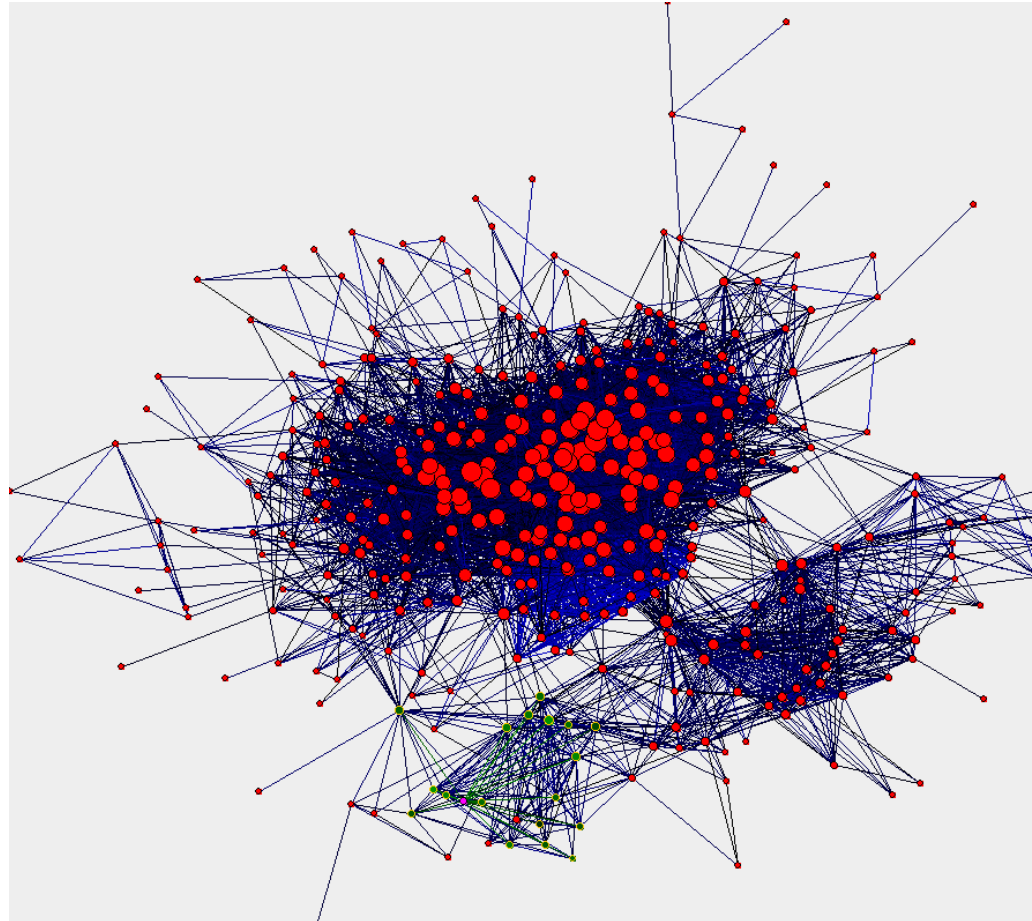
Use RANSAC to estimate fundamental matrix between each pair





# Image connectivity graph

---



(graph layout produced using the Graphviz toolkit: <http://www.graphviz.org/>)

Source: N. Snavely

# Incremental SFM

---

- Pick a pair of images with lots of inliers (and preferably, good EXIF data)
  - Initialize intrinsic parameters (focal length, principal point) from EXIF
  - Estimate extrinsic parameters ( $\mathbf{R}$  and  $\mathbf{t}$ )
    - [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

# The devil is in the details

---

- Handling degenerate configurations (e.g., homographies)
- Eliminating outliers
- Dealing with repetitions and symmetries
- Handling multiple connected components
- Closing loops
- ....

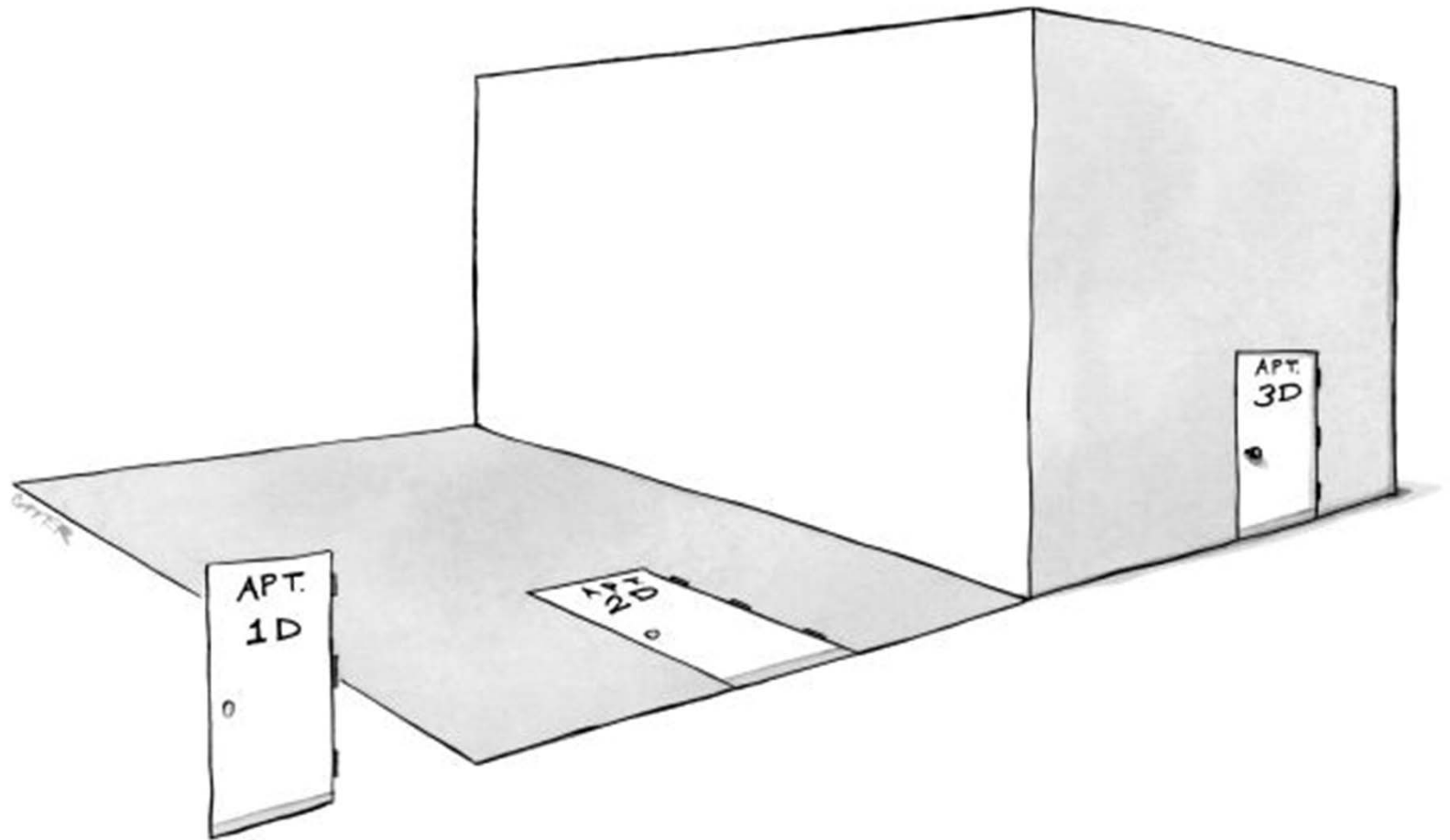
# Review: Structure from motion

---

- Ambiguity
- Affine structure from motion
  - Factorization
- Dealing with missing data
  - Incremental structure from motion
- Projective structure from motion
  - Bundle adjustment
  - Modern structure from motion pipeline

# Summary: 3D geometric vision

---



# Summary: 3D geometric vision

---

- Single-view geometry
  - The pinhole camera model
    - Variation: orthographic projection
  - The perspective projection matrix
  - Intrinsic and extrinsic parameters
  - Calibration
  - Single-view metrology, calibration using vanishing points
- Multiple-view geometry
  - Triangulation
  - The epipolar constraint
    - Essential matrix and fundamental matrix
  - Stereo
    - Binocular, multi-view
  - Structure from motion
    - Reconstruction ambiguity
    - Affine SFM
    - Projective SFM

---

**QUESTIONS?**