

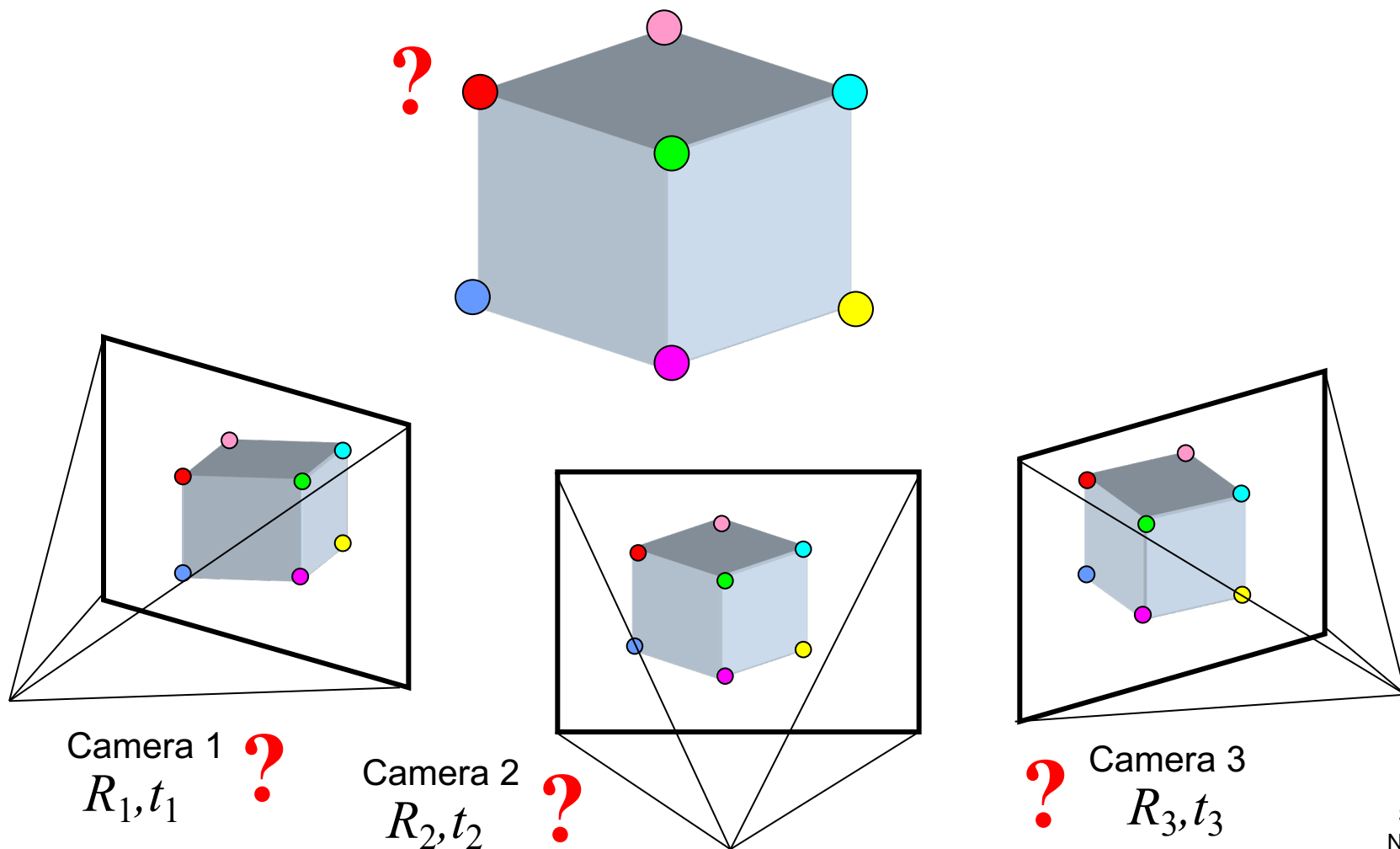
Structure from motion

Outline

- Representative SfM pipeline
 - Incremental SfM
 - Bundle adjustment
- Ambiguities in SfM
- Special Case: Affine structure from motion
 - Factorization
- SfM in practice

Structure from motion

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



Representative SFM pipeline



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

<http://phototour.cs.washington.edu/>

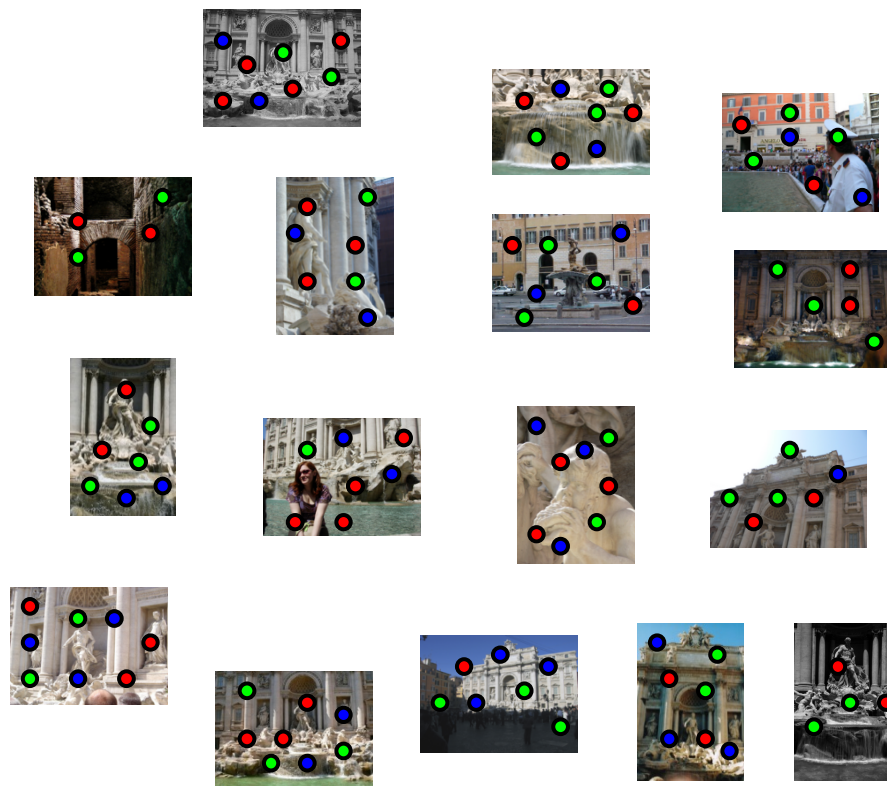
Feature detection

Detect SIFT features



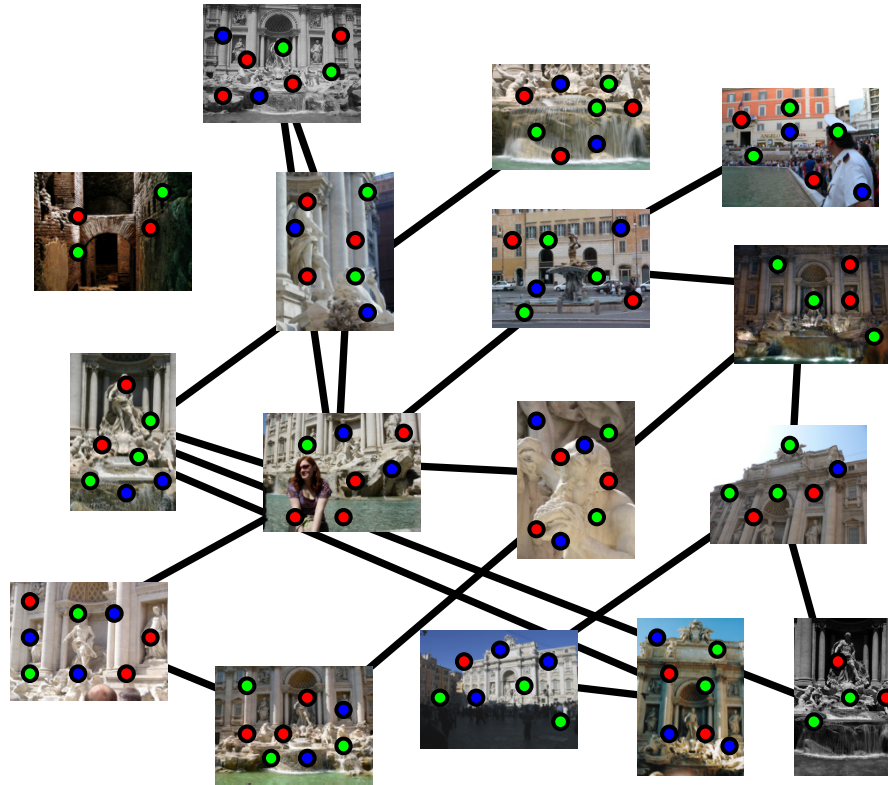
Feature detection

Detect SIFT features



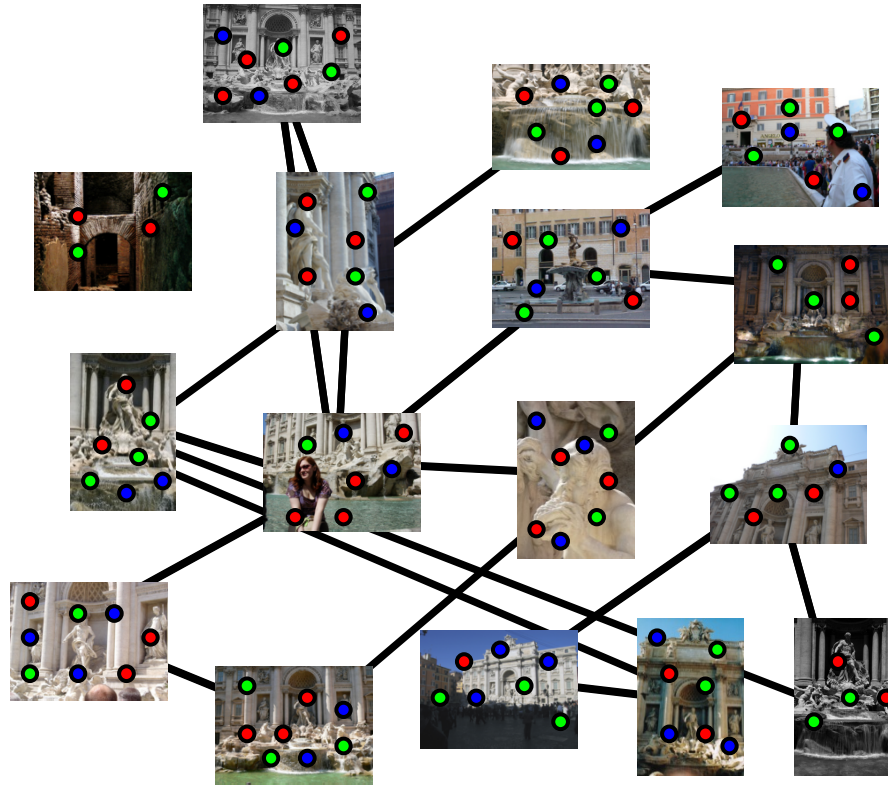
Feature matching

Match features between each pair of images



Feature matching

Use RANSAC to estimate fundamental matrix between each pair



Feature matching

Use RANSAC to estimate fundamental matrix between each pair



Feature matching

Use RANSAC to estimate fundamental matrix between each pair

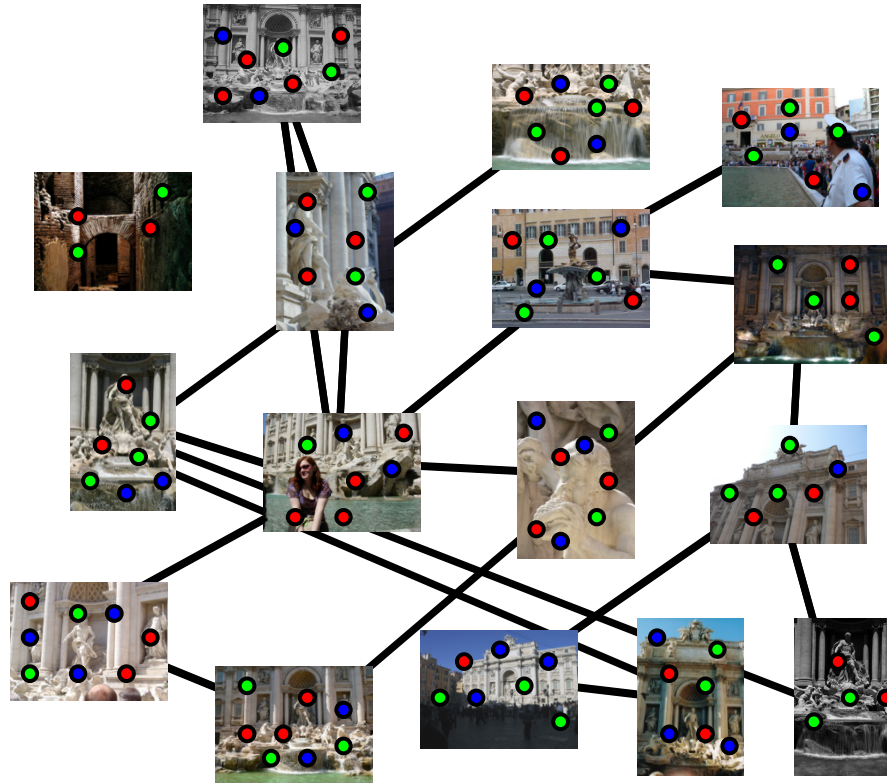
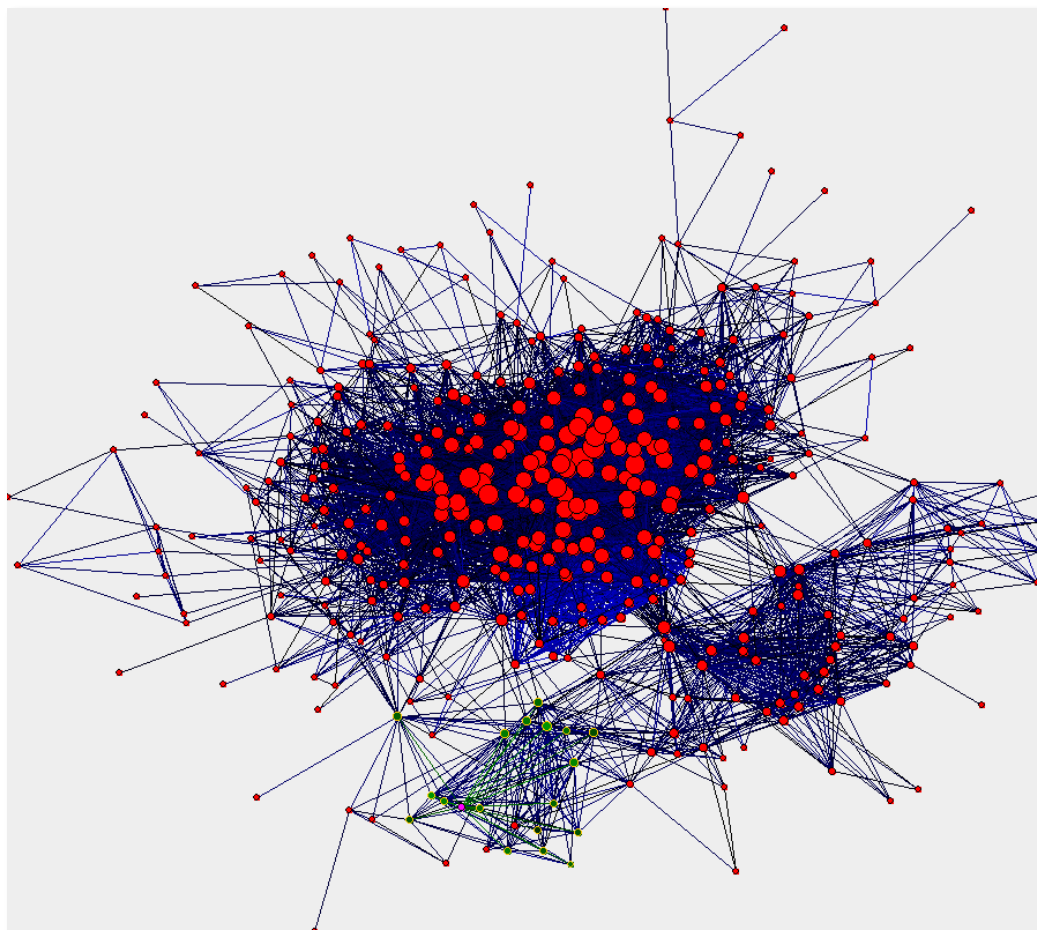


Image connectivity graph



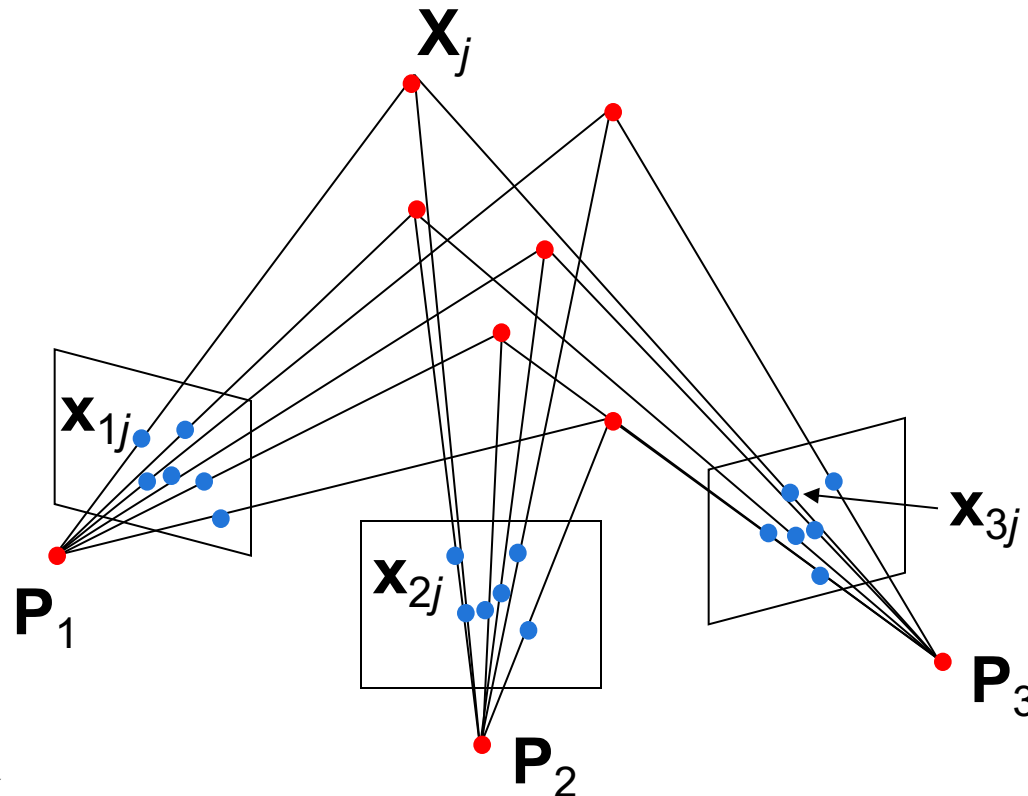
(graph layout produced using the Graphviz toolkit: <http://www.graphviz.org/>)

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}

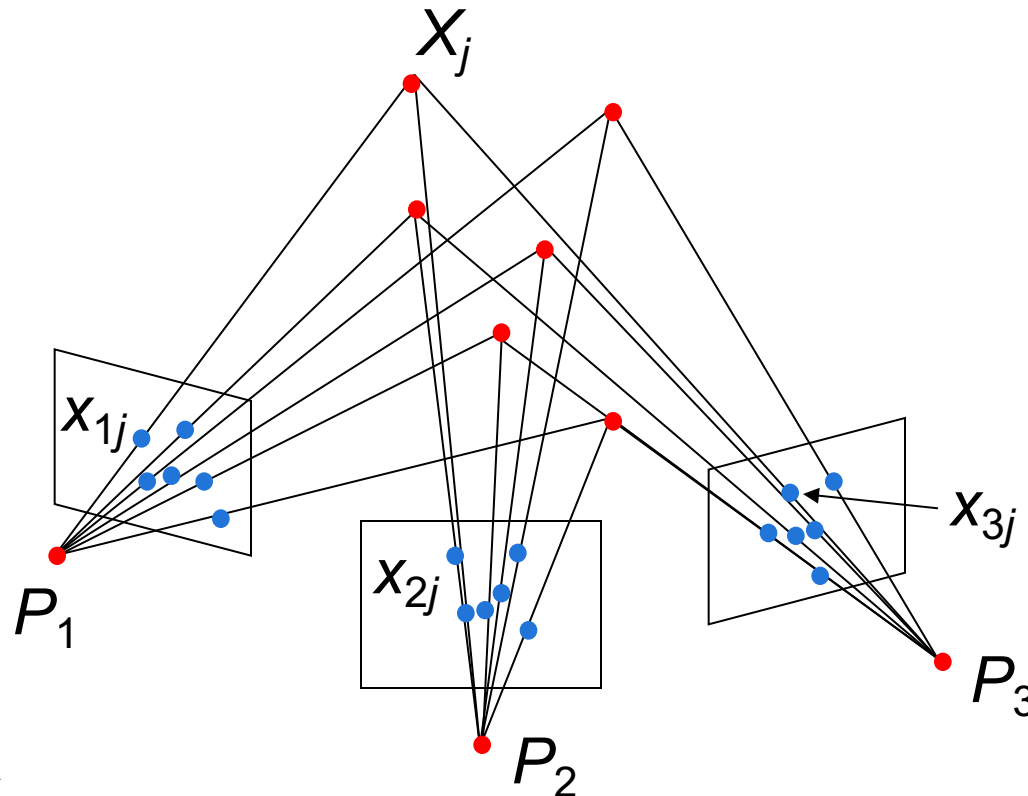


Projective structure from motion

- 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}



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} :

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

- We can solve for structure and motion when

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

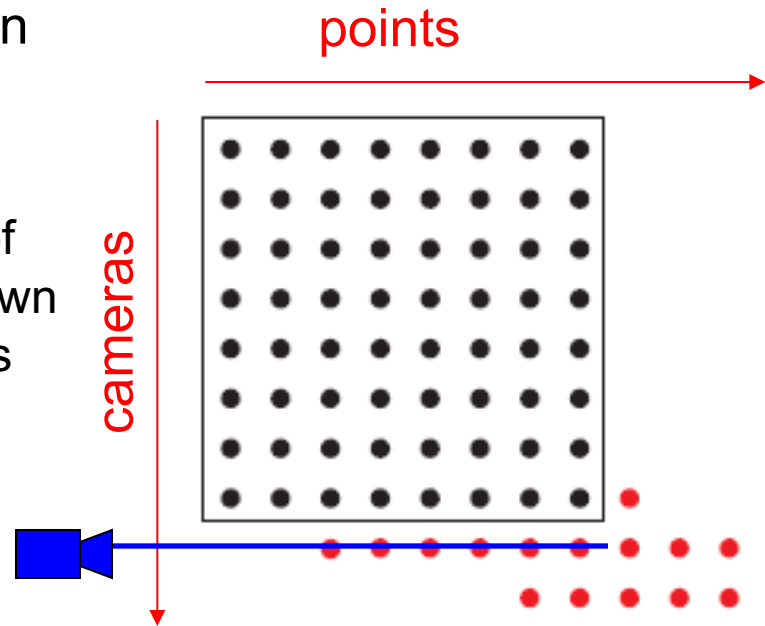
- For two cameras, at least 7 points are needed

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 ($\mathbf{F}^T \mathbf{b} = \mathbf{0}$), $\mathbf{A} = -[\mathbf{b}_\times] \mathbf{F}$

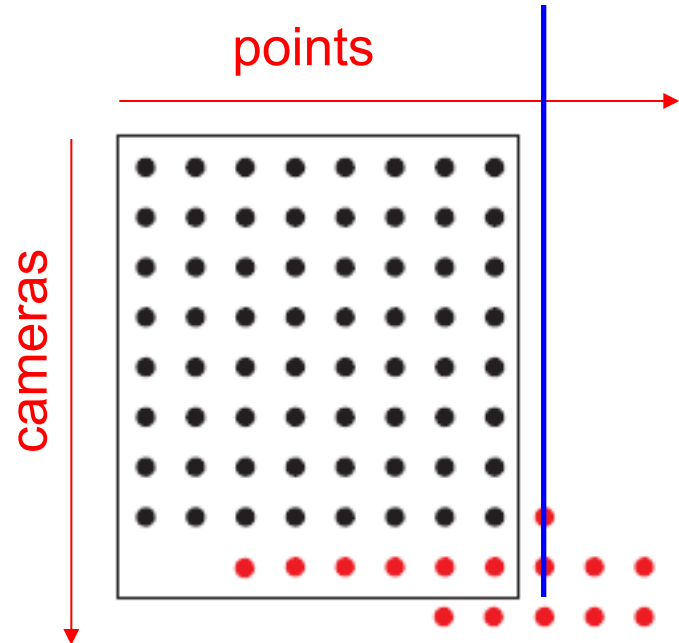
Incremental structure from motion

- 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*



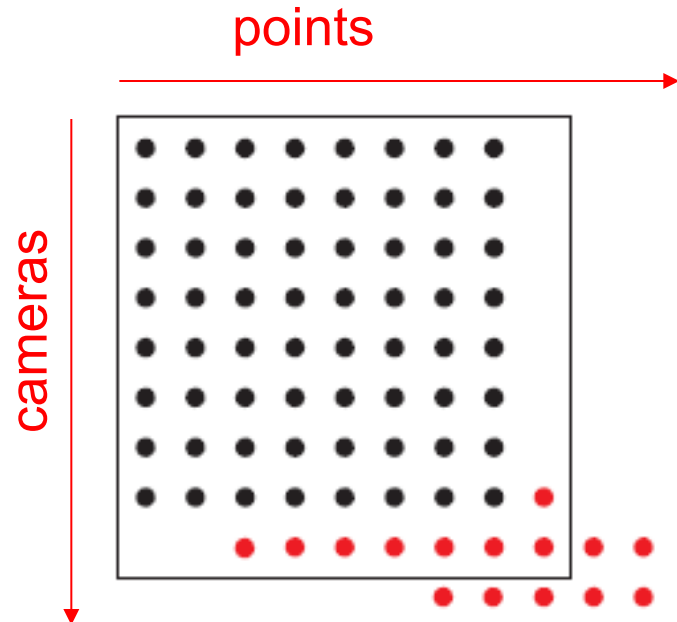
Incremental structure from motion

- 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*



Incremental 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

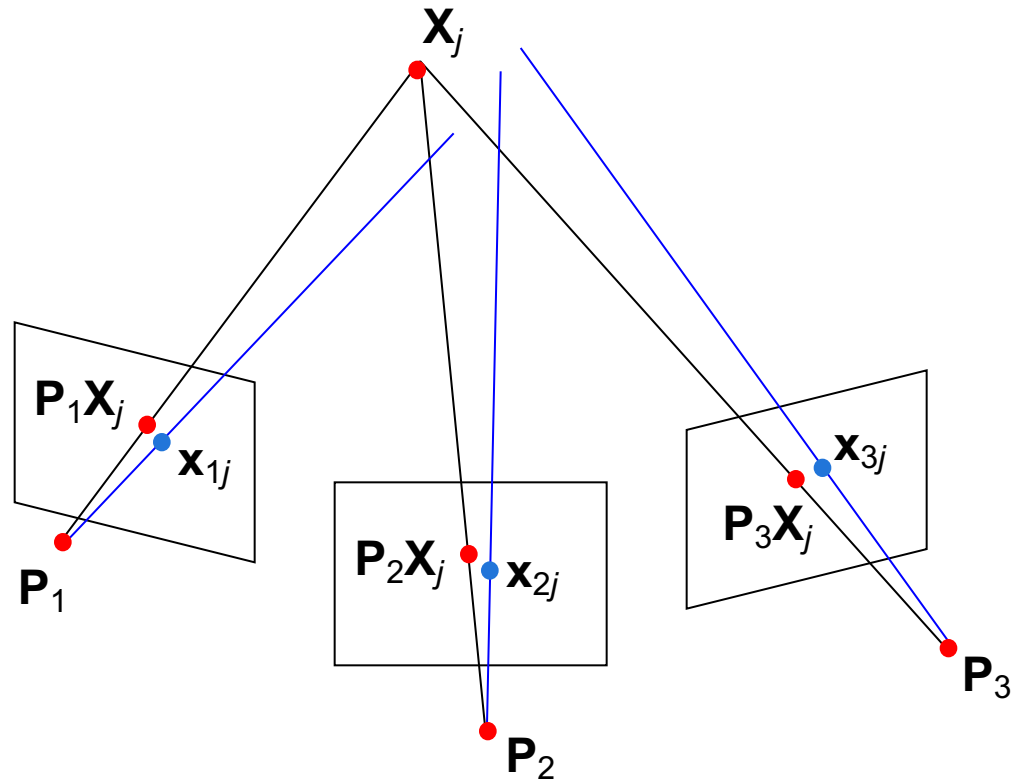


Bundle adjustment

- 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 ?



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}) 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

Photo Tourism

Exploring photo collections in 3D

Noah Snavely Steven M. Seitz Richard Szeliski
University of Washington *Microsoft Research*

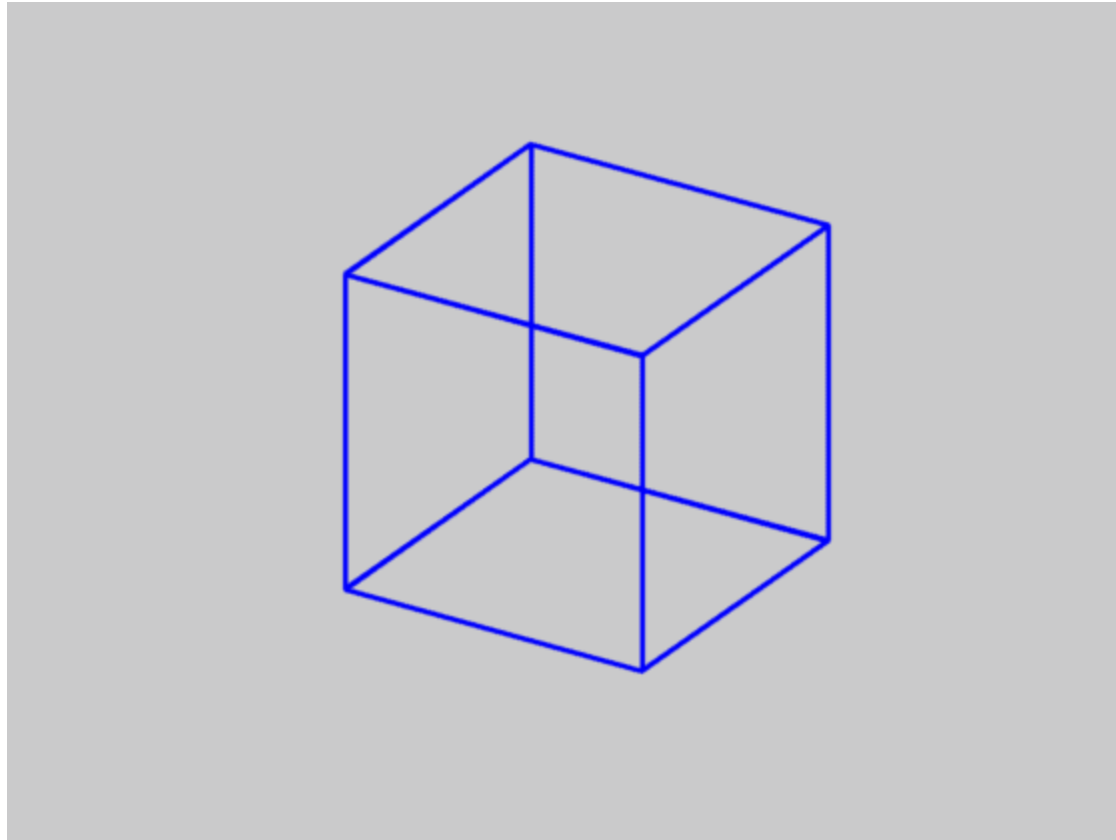
SIGGRAPH 2006

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/>
See also: <http://grail.cs.washington.edu/projects/rome/>

Outline

- Representative SfM pipeline
 - Incremental SfM
 - Bundle adjustment
- Ambiguities in SfM
- Special Case: Affine structure from motion
 - Factorization
- SfM in practice

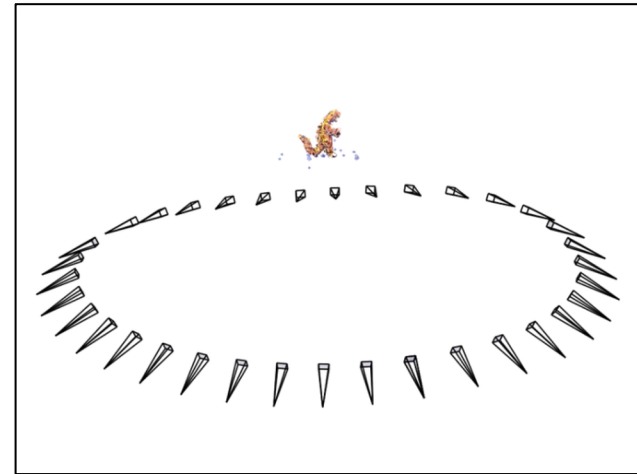
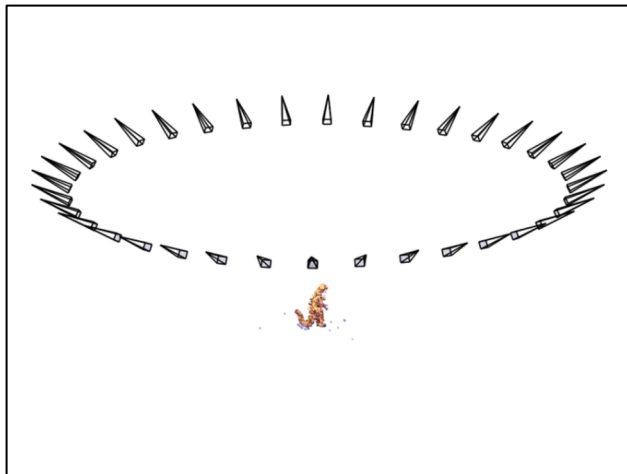
Is SFM always uniquely solvable?



Necker cube

Is SFM always uniquely solvable?

- Necker reversal



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:

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:

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

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:

Structure from motion ambiguity

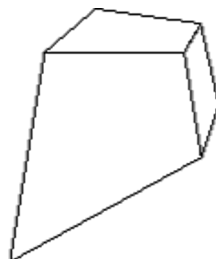
- If we scale the entire scene by some factor k and, at the same time, scale the camera matrices by the factor of $1/k$, the projections of the scene points in the image remain exactly the same
- More generally, if we transform the scene using a transformation \mathbf{Q} and apply the inverse transformation 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

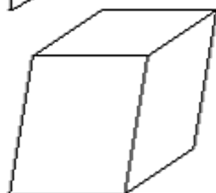
$$\begin{bmatrix} A & t \\ v^T & v \end{bmatrix}$$



Preserves intersection and tangency

Affine
12dof

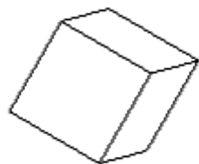
$$\begin{bmatrix} A & t \\ 0^T & 1 \end{bmatrix}$$



Preserves parallelism, volume ratios

Similarity
7dof

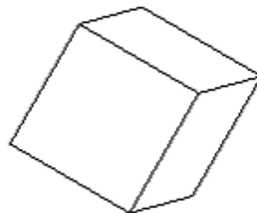
$$\begin{bmatrix} sR & t \\ 0^T & 1 \end{bmatrix}$$



Preserves angles, ratios of length

Euclidean
6dof

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

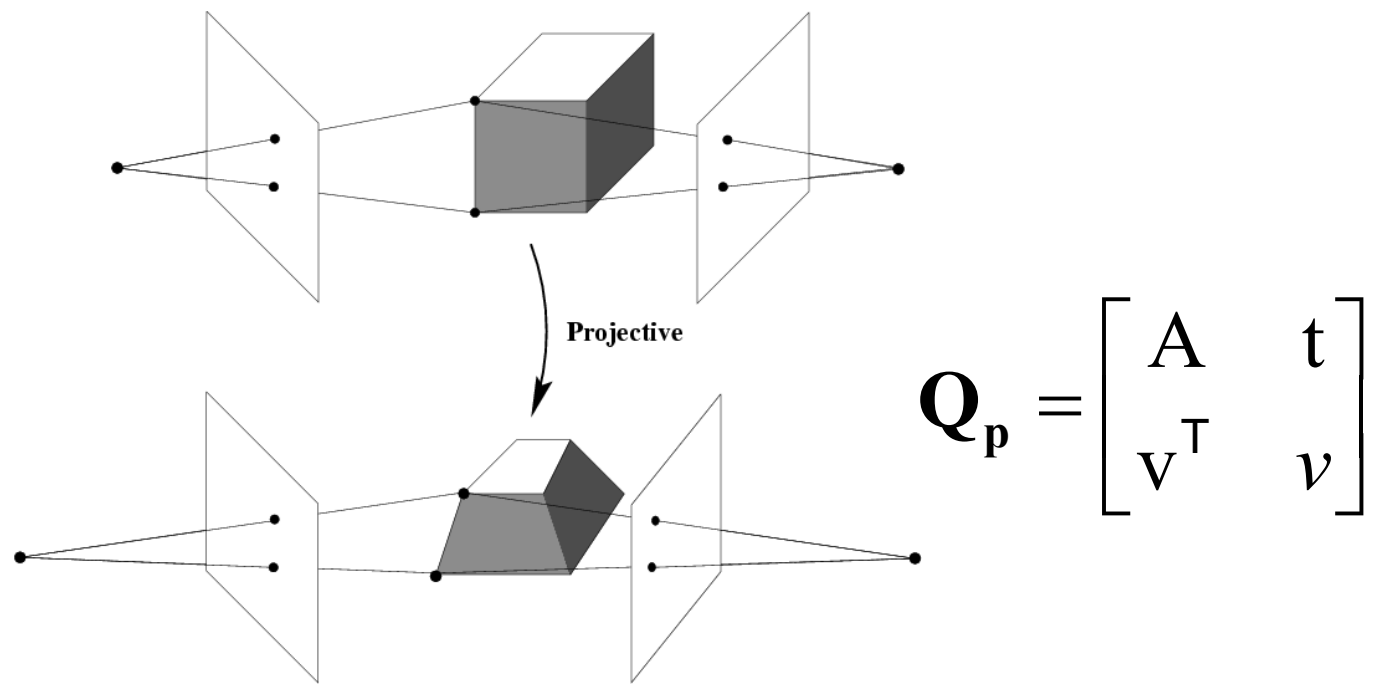


Preserves angles, lengths

- 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

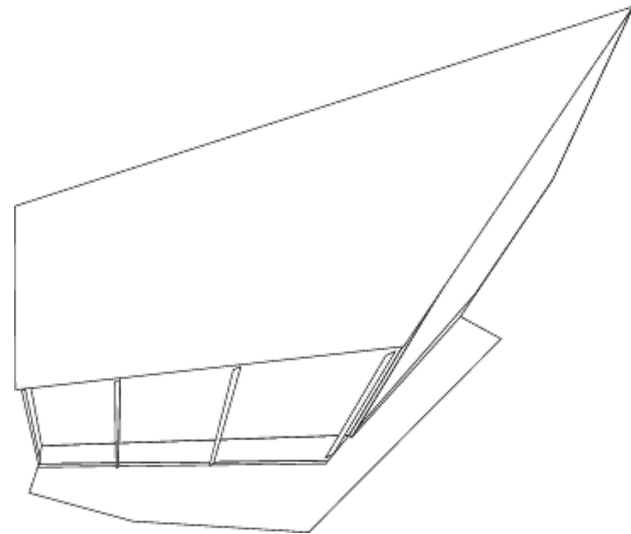
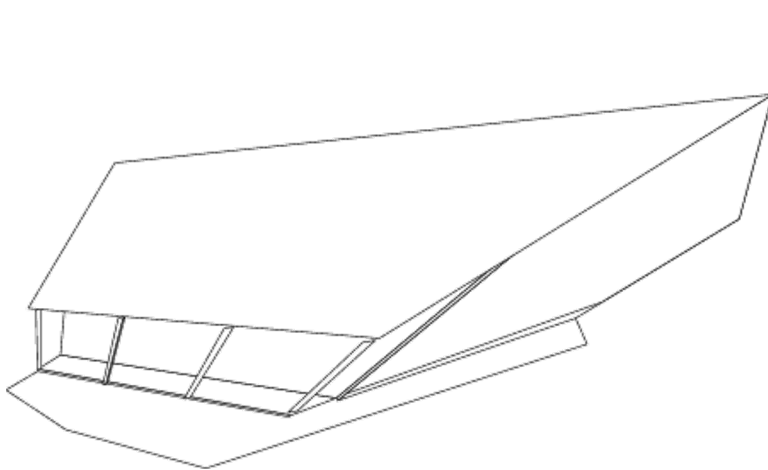
Projective ambiguity

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



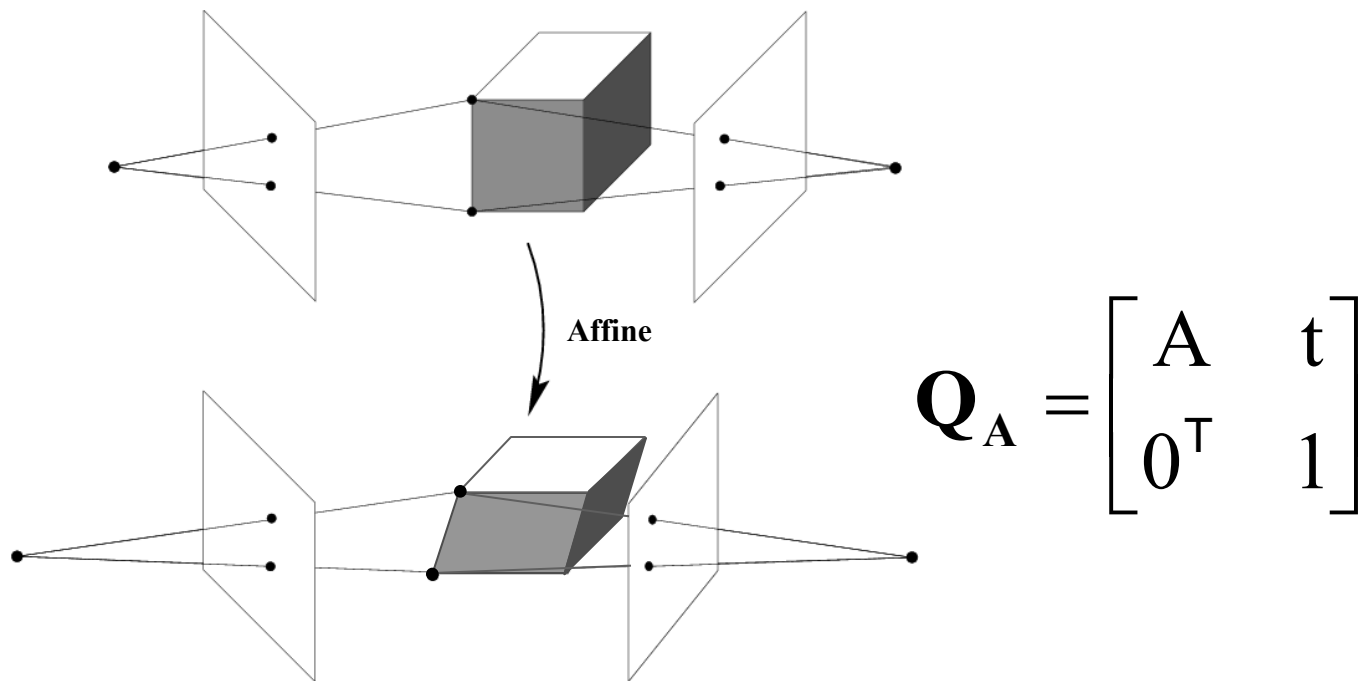
$$\mathbf{x} = \mathbf{P}\mathbf{X} = \left(\mathbf{P}\mathbf{Q}_p^{-1}\right)\left(\mathbf{Q}_p\mathbf{X}\right)$$

Projective ambiguity



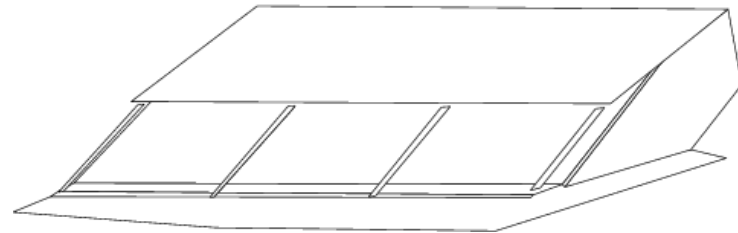
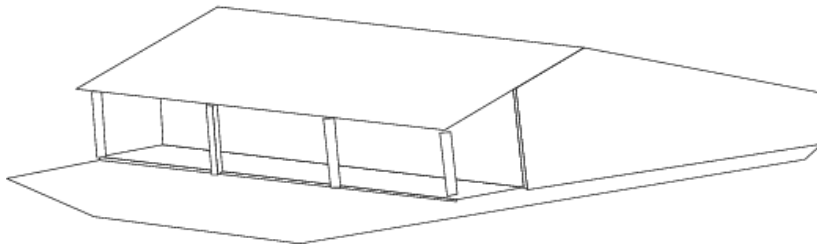
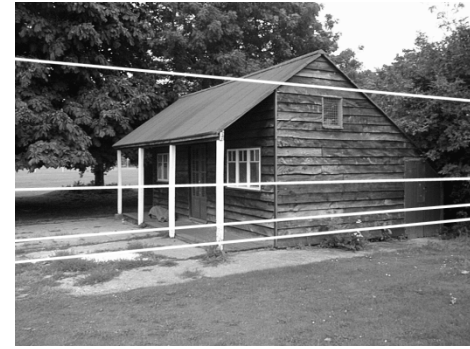
Affine ambiguity

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



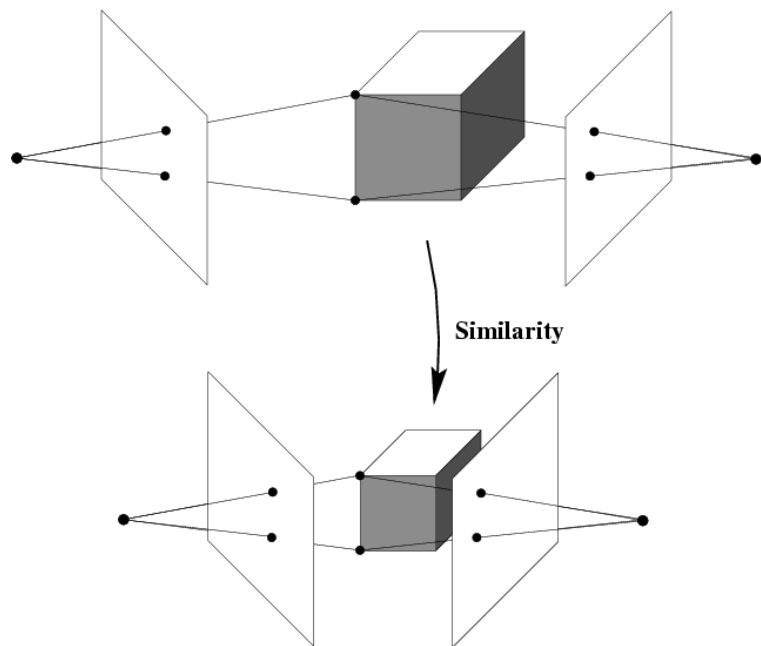
$$\mathbf{x} = \mathbf{P}\mathbf{X} = \left(\mathbf{P}\mathbf{Q}_A^{-1} \right) \left(\mathbf{Q}_A \mathbf{X} \right)$$

Affine ambiguity



Similarity ambiguity

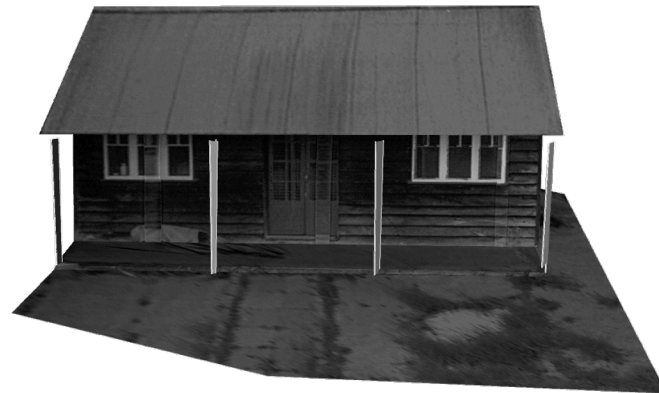
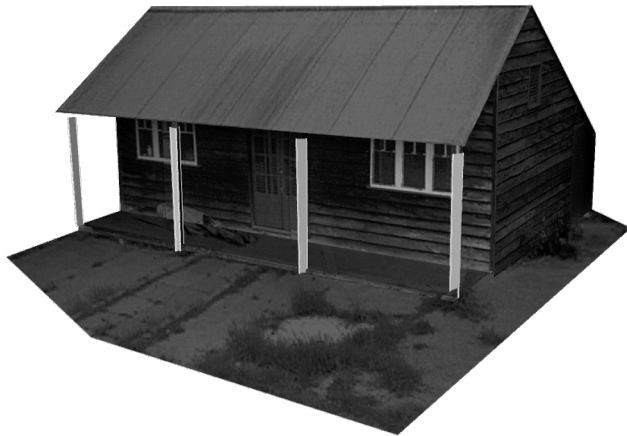
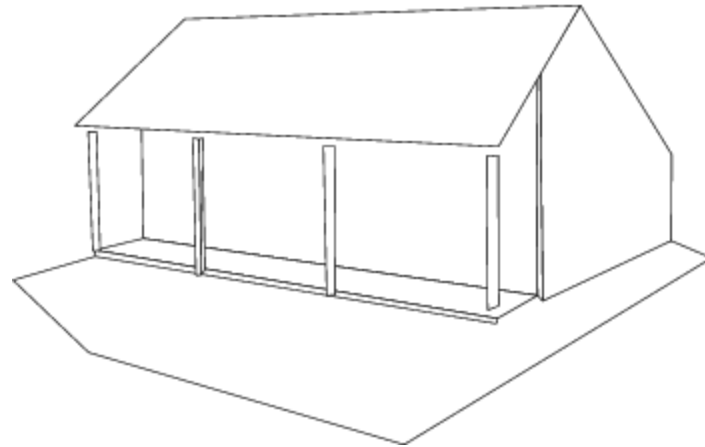
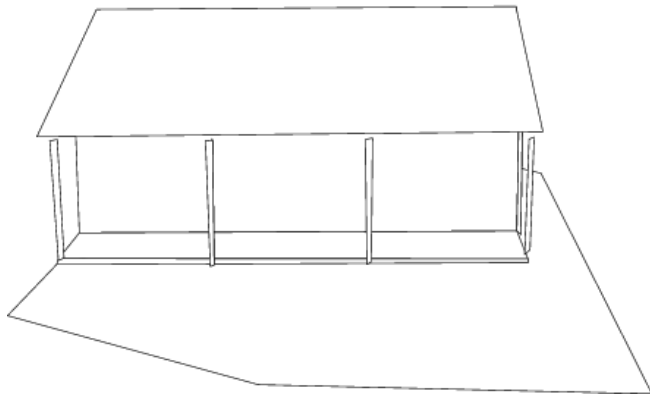
- A reconstruction that obeys orthogonality constraints on camera parameters and/or scene



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

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

Similarity ambiguity

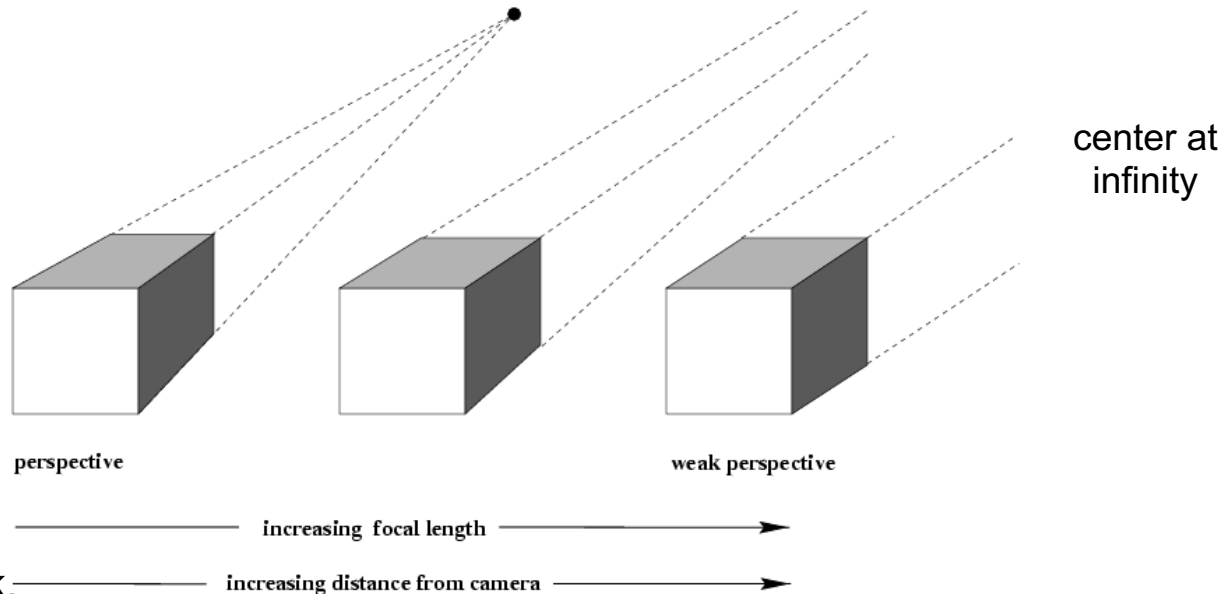


Outline

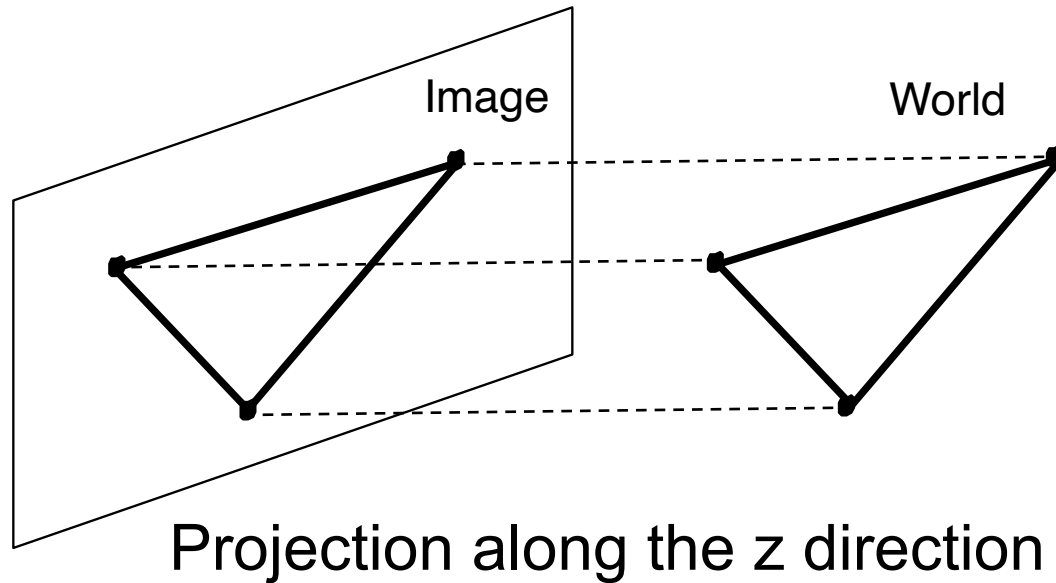
- Representative SfM pipeline
 - Incremental SfM
 - Bundle adjustment
- Ambiguities in SfM
- Special Case: Affine structure from motion
 - Factorization
- SfM in practice

Special Case: Affine structure from motion

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



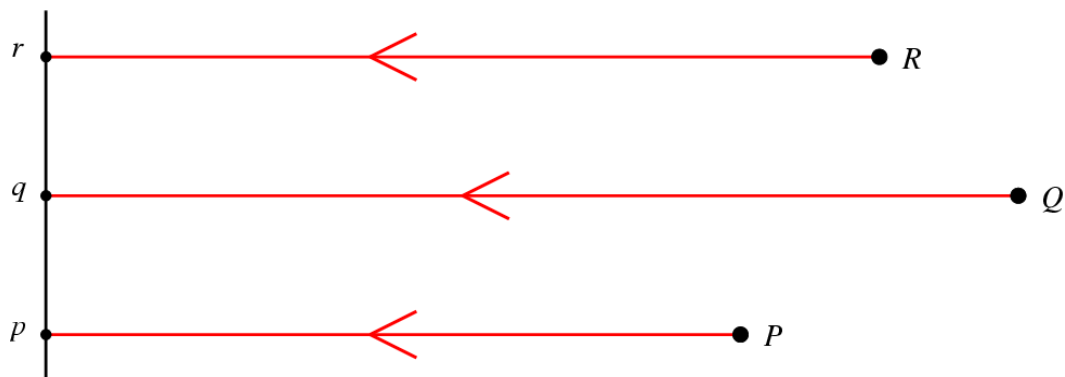
Recall: Orthographic Projection



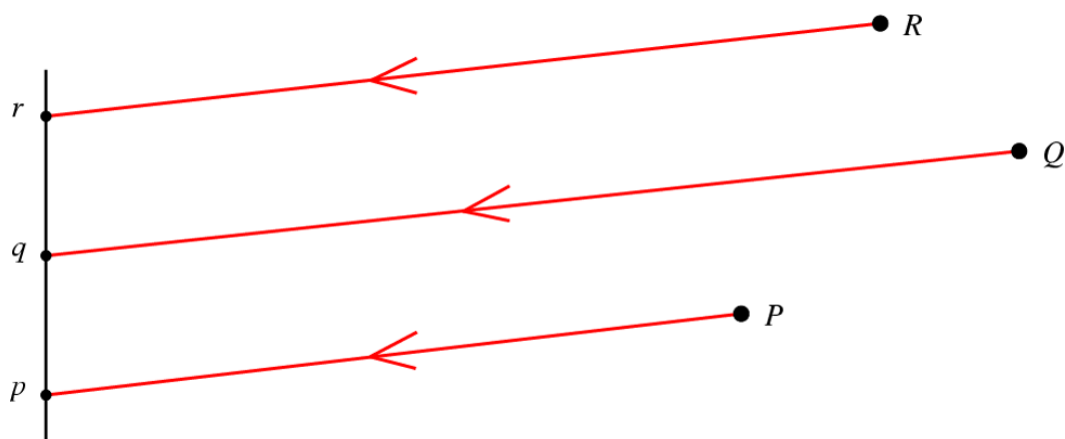
$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \Rightarrow (x, y)$$

Affine cameras

Orthographic Projection



Parallel Projection

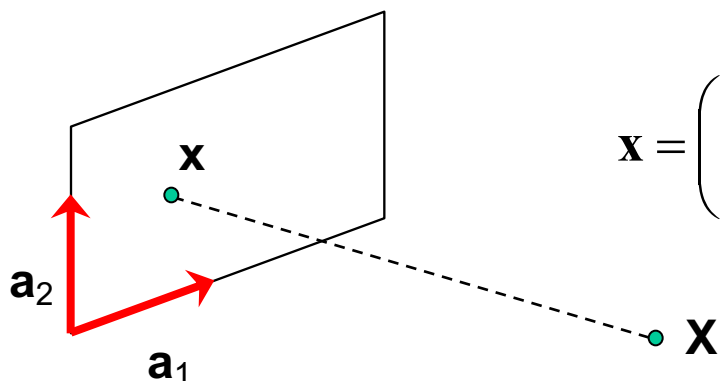


Affine cameras

- 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:

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

↓ cameras ($2m$)

→ points (n)

Affine structure from motion

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

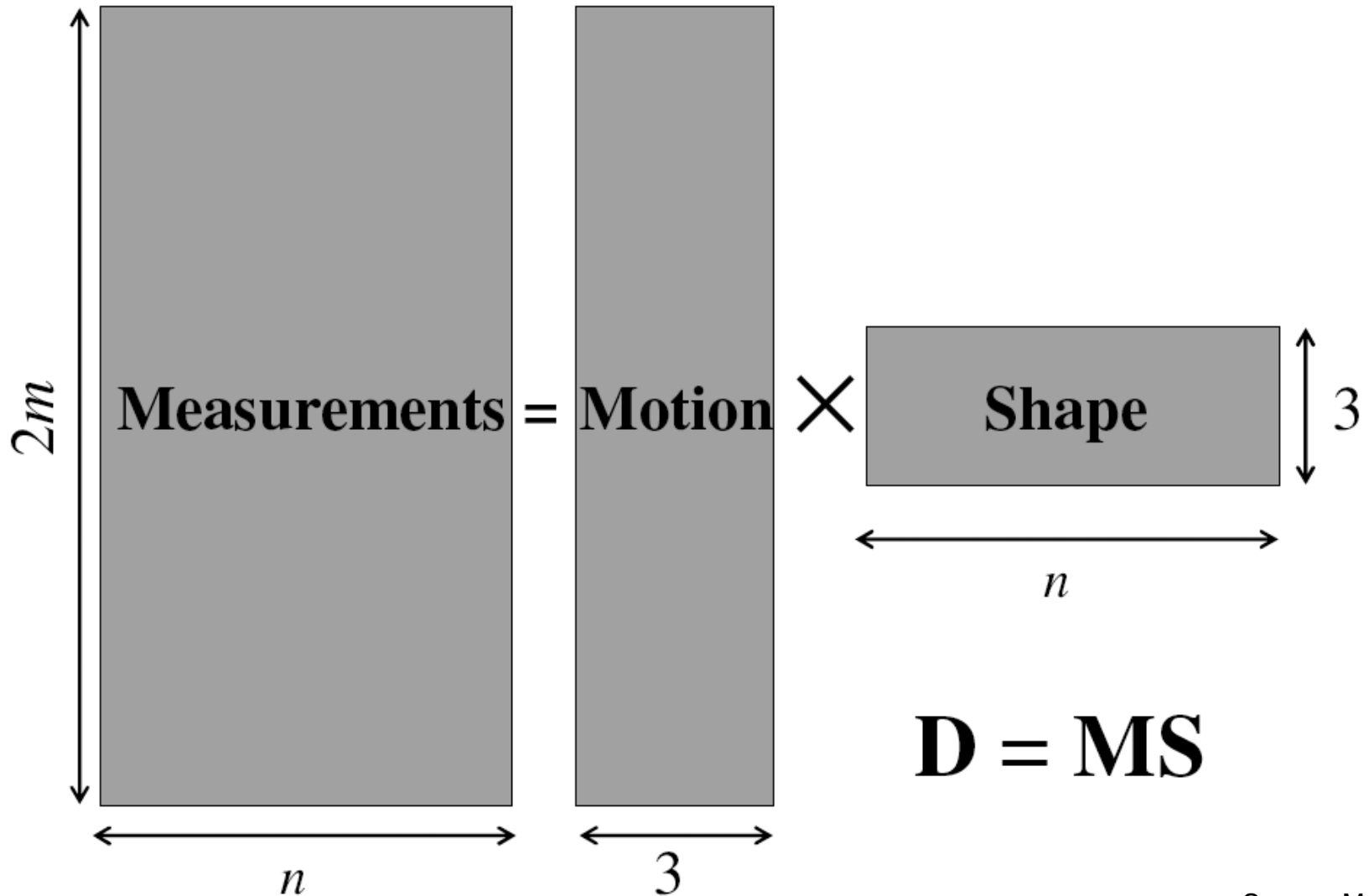
$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \cdots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \cdots & \hat{\mathbf{x}}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \cdots & \hat{\mathbf{x}}_{mn}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

points ($3 \times n$)

cameras
($2m \times 3$)

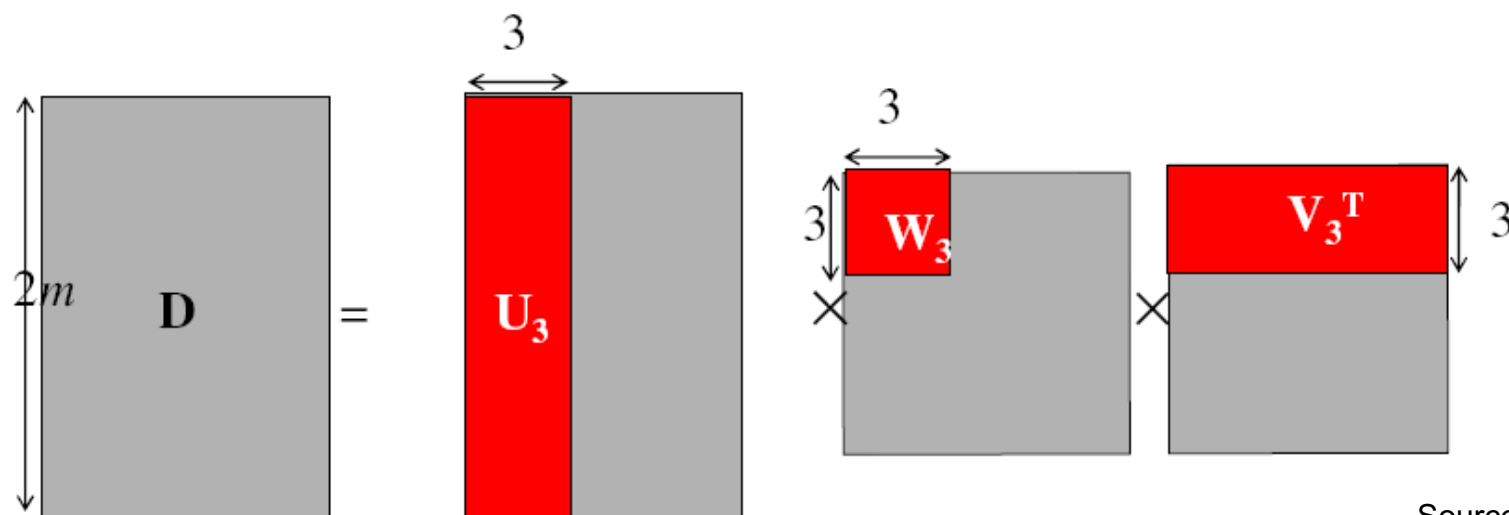
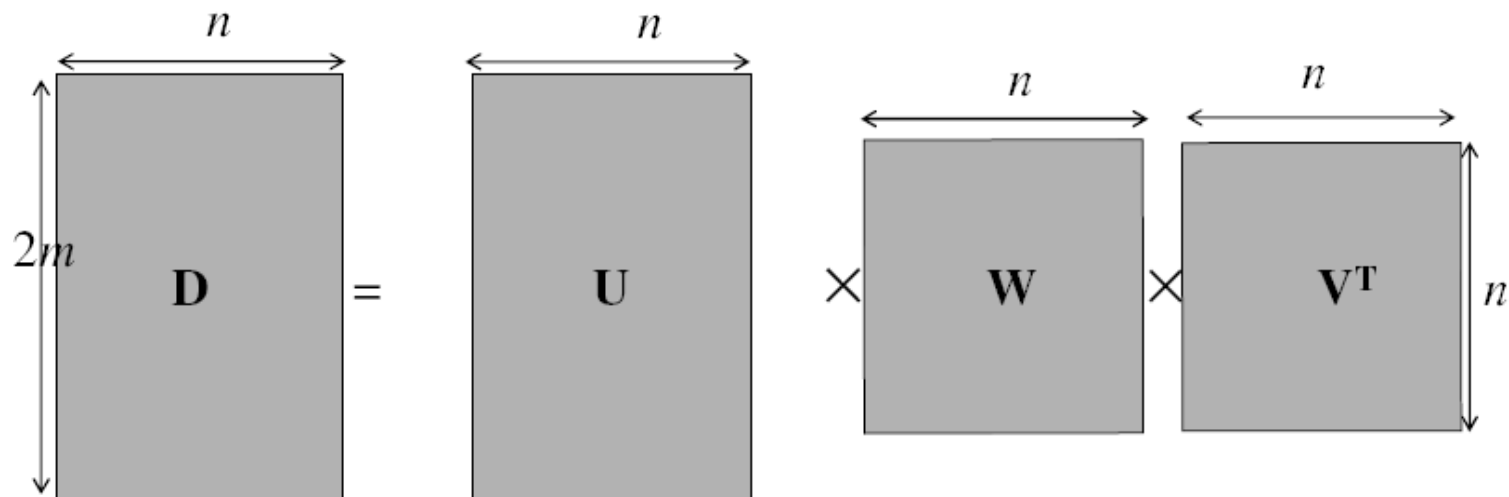
The measurement matrix $\mathbf{D} = \mathbf{MS}$ must have rank 3!

Factorizing the measurement matrix



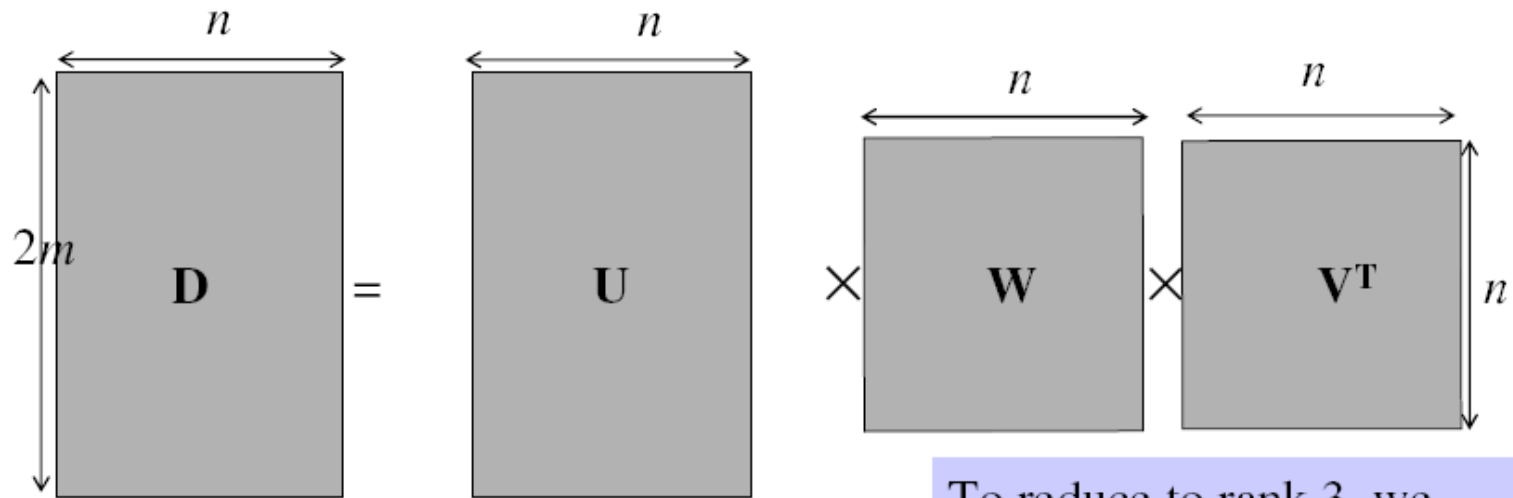
Factorizing the measurement matrix

- Singular value decomposition of D :

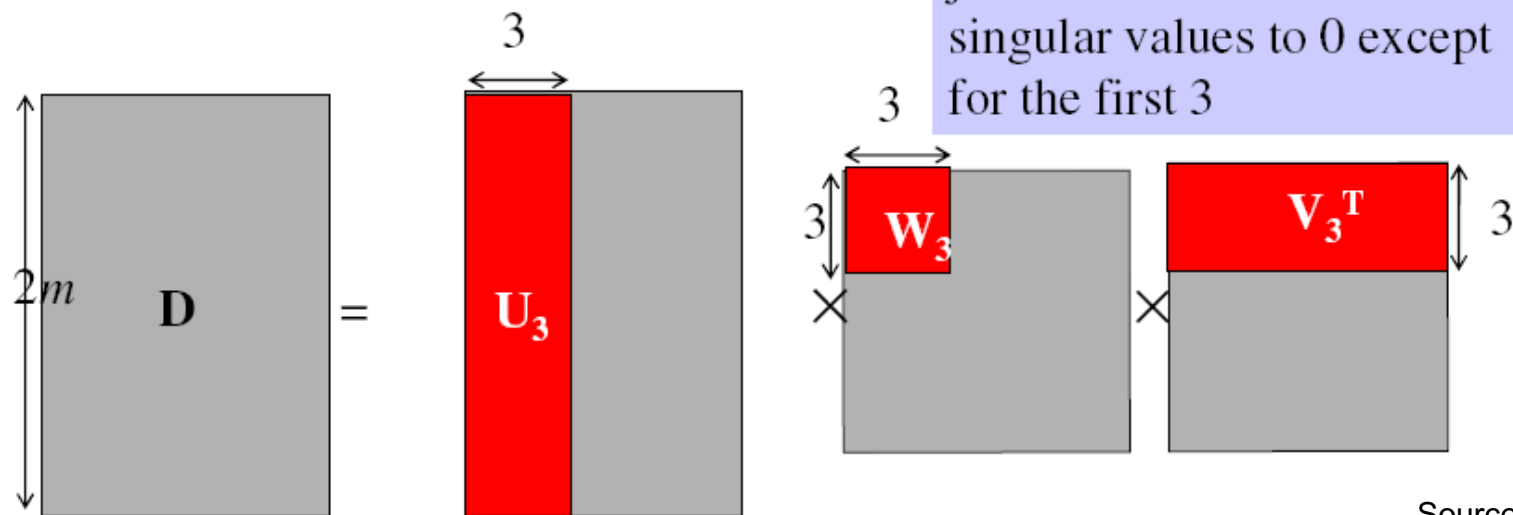


Factorizing the measurement matrix

- Singular value decomposition of D :

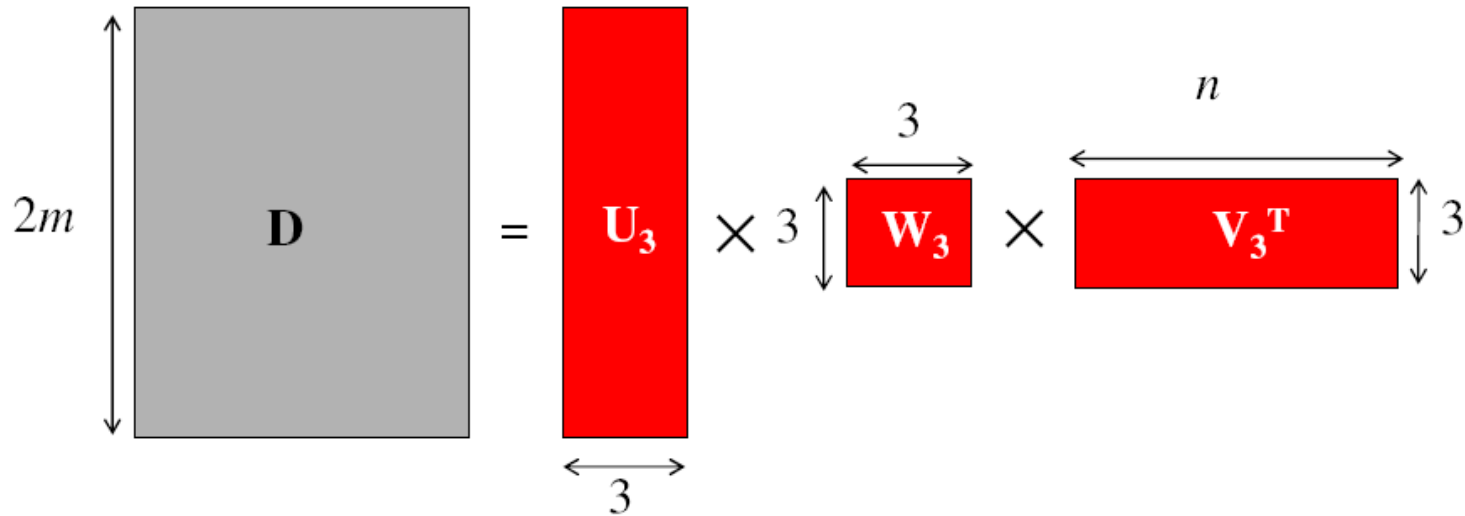


To reduce to rank 3, we just need to set all the singular values to 0 except for the first 3



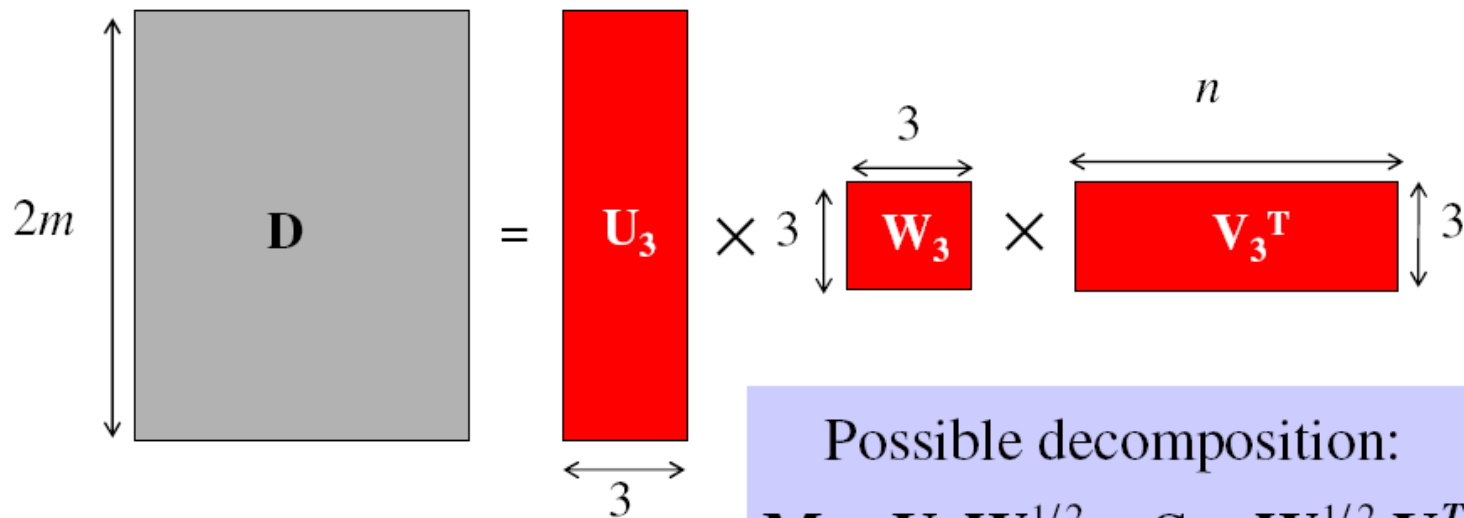
Factorizing the measurement matrix

- Obtaining a factorization from SVD:



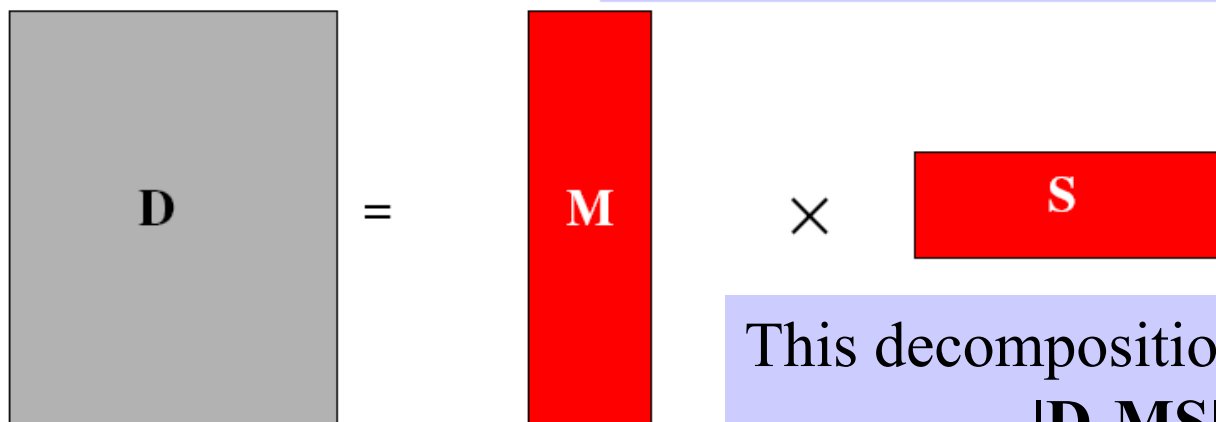
Factorizing the measurement matrix

- Obtaining a factorization from SVD:



Possible decomposition:

$$\mathbf{M} = \mathbf{U}_3 \mathbf{W}_3^{1/2} \quad \mathbf{S} = \mathbf{W}_3^{1/2} \mathbf{V}_3^T$$



This decomposition minimizes $|\mathbf{D} - \mathbf{MS}|^2$

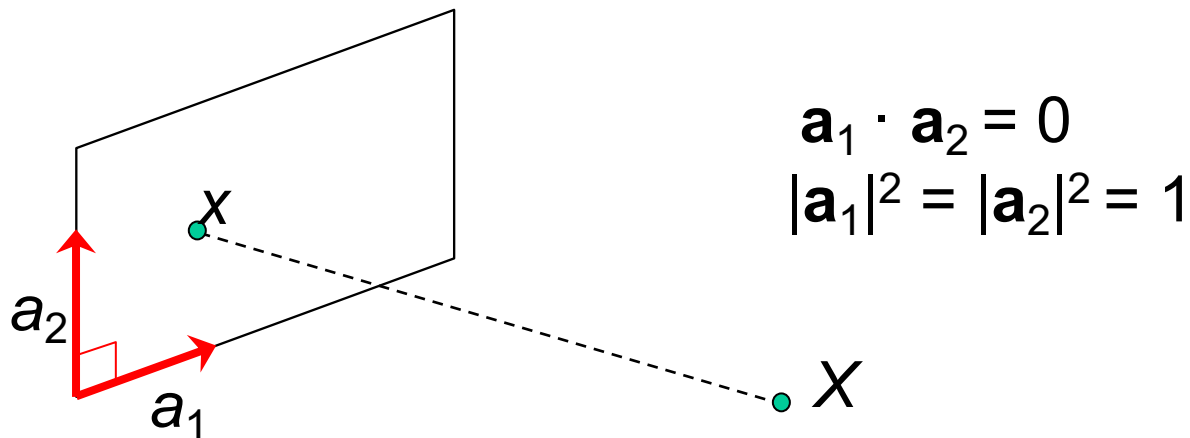
Affine ambiguity

The diagram shows a gray square labeled **D** on the left. To its right is an equals sign. Further right is a red vertical rectangle labeled **M**. To its right is a multiplication symbol \times . To the right of the multiplication symbol is a red horizontal rectangle labeled **S**.

- The decomposition is not unique. We get the same **D** by using any 3×3 matrix **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:
$$(\mathbf{A}_i \mathbf{C})(\mathbf{A}_i \mathbf{C})^T = \mathbf{A}_i (\mathbf{C} \mathbf{C}^T) \mathbf{A}_i = \mathbf{I}_d, \quad i = 1, \dots, m$$
 - 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



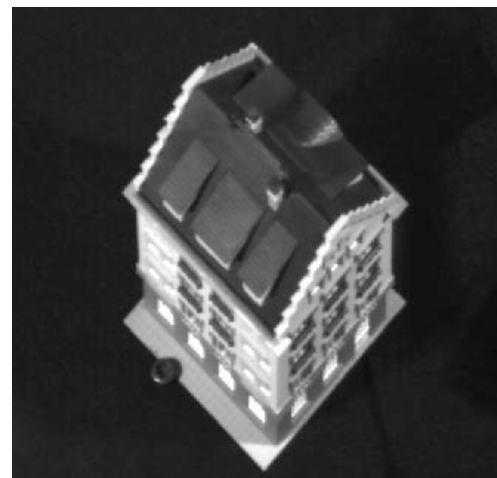
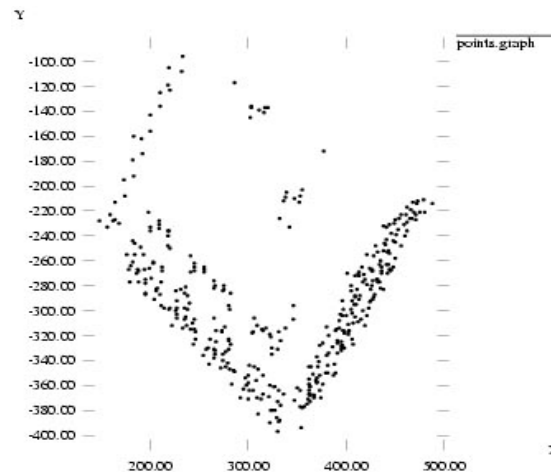
60



120



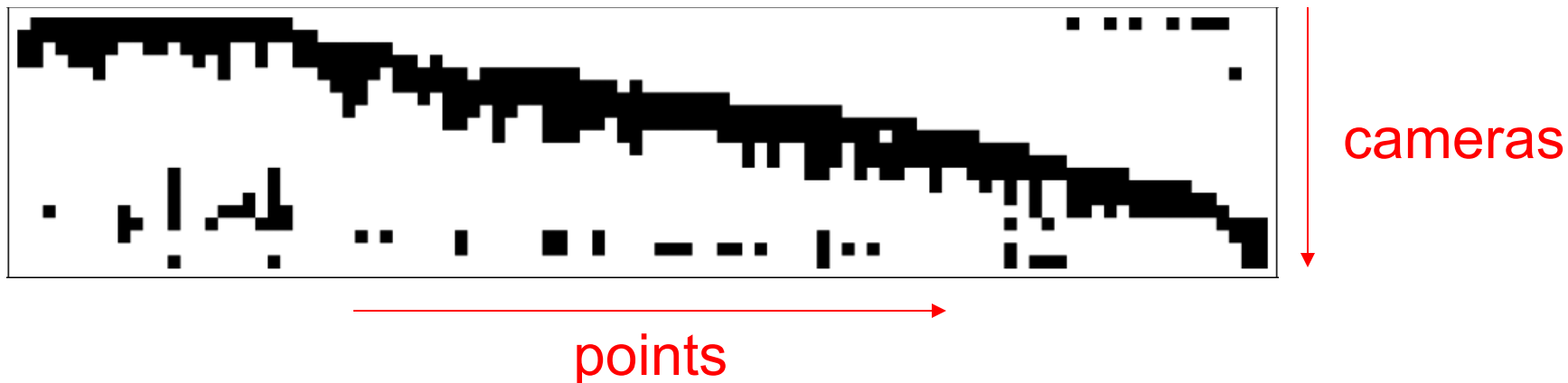
150



C. Tomasi and T. Kanade, [Shape and motion from image streams under orthography: A factorization method](#), IJCV 1992

Dealing with missing data

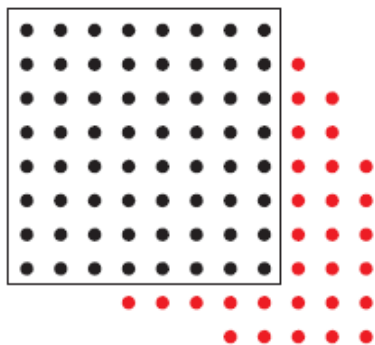
- So far, we have assumed that all points are visible in all views
- In reality, the measurement matrix typically looks something like this:



- 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
(*triangulation*)

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

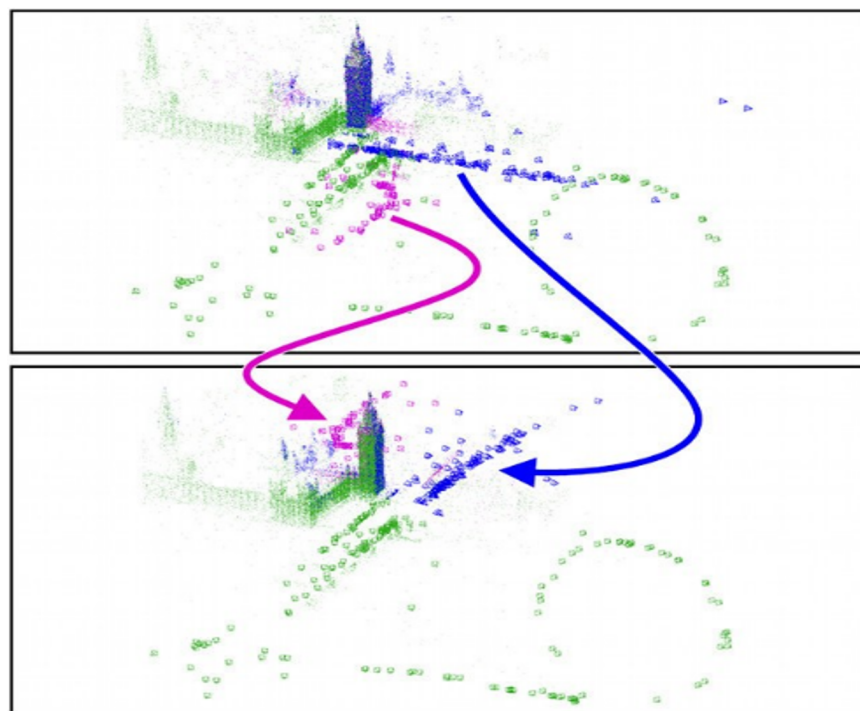
Outline

- Representative SfM pipeline
 - Incremental SfM
 - Bundle adjustment
- Ambiguities in SfM
- Special Case: Affine structure from motion
 - Factorization
- SfM in practice

The devil is in the details

- Handling degenerate configurations (e.g., homographies)
- Eliminating outliers
- Dealing with repetitions and symmetries

Repetitive structures



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
- Making the whole thing efficient!
 - See, e.g., [Towards Linear-Time Incremental Structure from Motion](#)

SfM software

- [Bundler](#)
- [OpenSfM](#)
- [OpenMVG](#)
- [VisualSfM](#)
- See also [Wikipedia's list of toolboxes](#)

Outline

- Representative SfM pipeline
 - Incremental SfM
 - Bundle adjustment
- Ambiguities in SfM
- Special Case: Affine structure from motion
 - Factorization
- SfM in practice