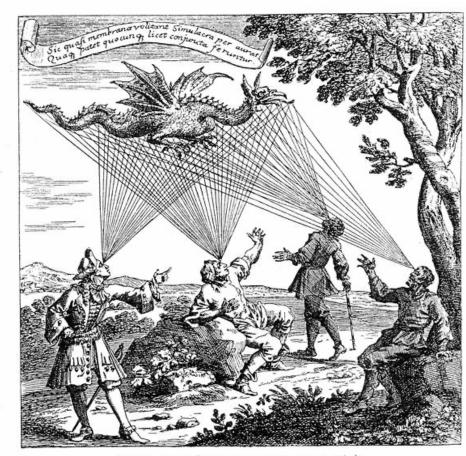
Structure from motion

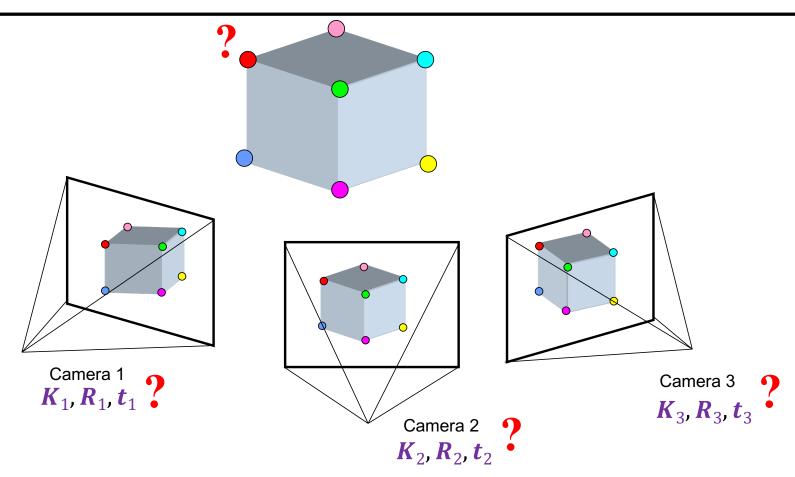


Драконь, видимый подъ различными углами эрінія По граворі на міли изи "Oculus artificialis teledioptricus" Цана. 1702 года.

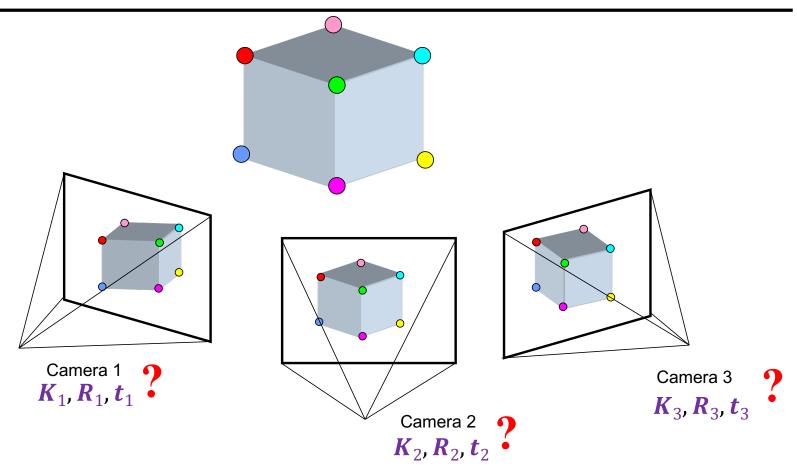
Outline: Structure from motion

- Problem definition and ambiguities
- Affine structure from motion
 - Factorization
- Projective structure from motion
 - Bundle adjustment
- Modern structure from motion pipeline

Structure from motion



Recall: Calibration



• Given a set of known 3D points seen by a camera, compute the camera parameters

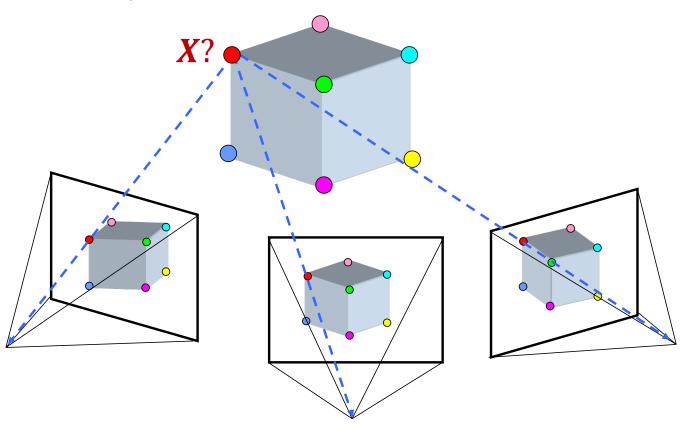
• Given projections of a 3D point in two or more images (with known camera matrices), find the coordinates of the point



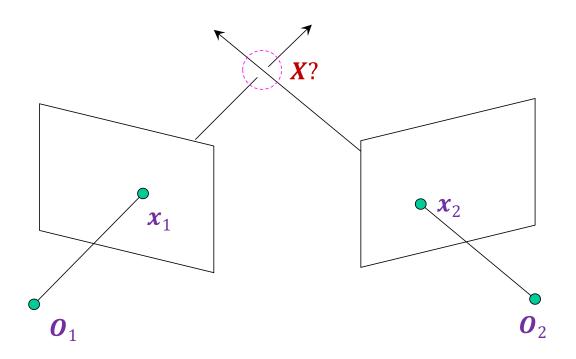




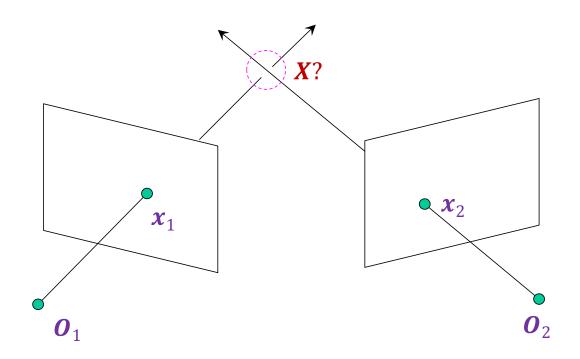
• Given projections of a 3D point in two or more images (with known camera matrices), find the coordinates of the point



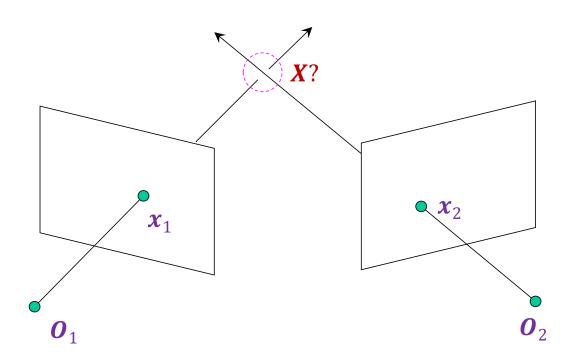
• Given projections of a 3D point in two or more images (with known camera matrices), find the coordinates of the point



• We want to intersect the two visual rays corresponding to x_1 and x_2 , but because of noise and numerical errors, they don't meet exactly



Ignored camera intrinsics cause cameras are known!



$$\mathbf{x}_1 \equiv \left[\mathcal{I}|\mathbf{0}\right] \mathbf{X}$$

$$\mathbf{x}_2 \equiv \left[\mathcal{R}|\mathbf{t}\right]\mathbf{X}$$

$$\mathbf{x}_1 \equiv [\mathcal{I}|\mathbf{0}] \, \mathbf{X}$$

$$\mathbf{x}_2 \equiv [\mathcal{R}|\mathbf{t}] \, \mathbf{X}$$

$$X_1 = x_1 X_3$$

Affine coordinates of image point

$$X_2 = \overset{\downarrow}{y_1} X_3$$

Remember camera is known, And substitute

$$\mathbf{x}_1 \equiv [\mathcal{I}|\mathbf{0}] \, \mathbf{X}$$

$$\mathbf{x}_2 \equiv [\mathcal{R}|\mathbf{t}] \, \mathbf{X}$$

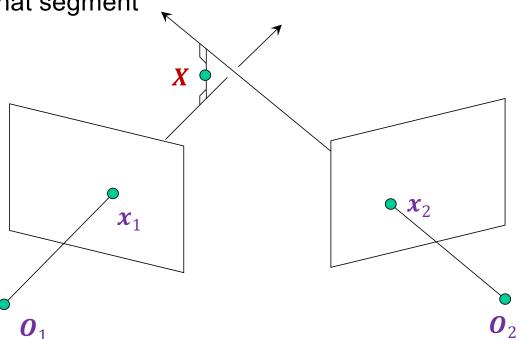
$$x_2 = \frac{r_{11}x_1X_3 + r_{12}y_1X_3 + r_{13}X_3 + t_1}{r_{31}x_1X_3 + r_{32}y_1X_3 + r_{32}X_3 + t_3}$$

$$y_2 = \frac{r_{21}x_1X_3 + r_{22}y_1X_3 + r_{23}X_3 + t_2}{r_{31}x_1X_3 + r_{32}y_1X_3 + r_{32}X_3 + t_3}$$

Triangulation – Straightforward Approaches

- Above gives two possible points, average them
- Choose least squares X_3

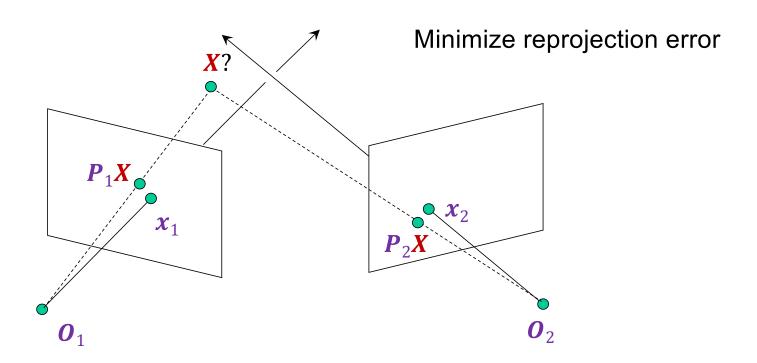
 Find shortest segment connecting the two viewing rays and let X be the midpoint of that segment



Triangulation: Nonlinear approach

Find X that minimizes

$$\|\operatorname{proj}(\boldsymbol{P}_1\boldsymbol{X}) - \boldsymbol{x}_1\|_2^2 + \|\operatorname{proj}(\boldsymbol{P}_2\boldsymbol{X}) - \boldsymbol{x}_2\|_2^2$$

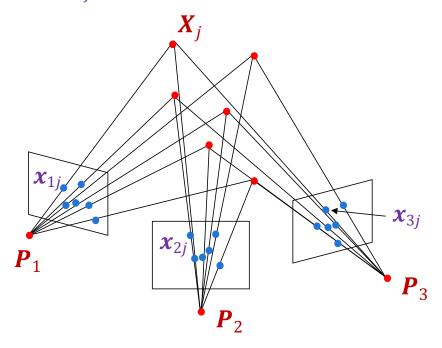


Structure from motion: Problem formulation

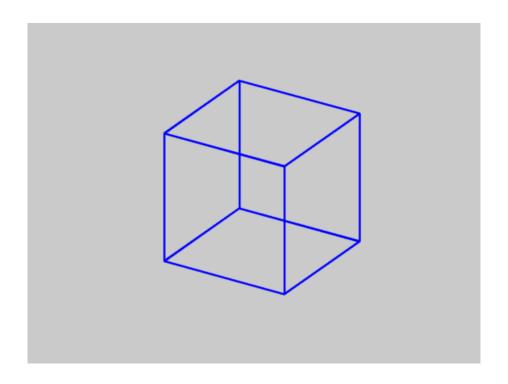
Given: m images of n fixed 3D points such that (ignoring visibility)

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

• Problem: estimate m projection matrices P_i and n 3D points X_j from the mn correspondences x_{ij}



Is SFM always uniquely solvable?



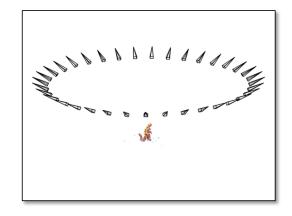
Necker cube

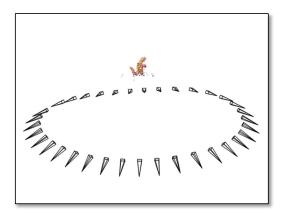
Source: N. Snavely

Is SFM always uniquely solvable?

Could actually happen in affine structure from motion:







Source: N. Snavely

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:

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

- 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:

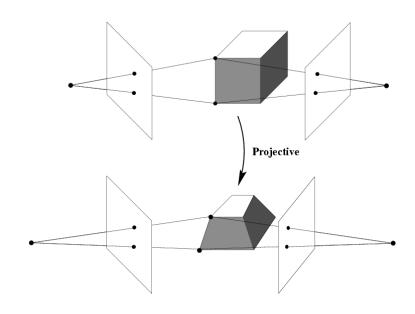
$$x \cong PX = (PQ^{-1})(QX)$$

Projective ambiguity

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

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

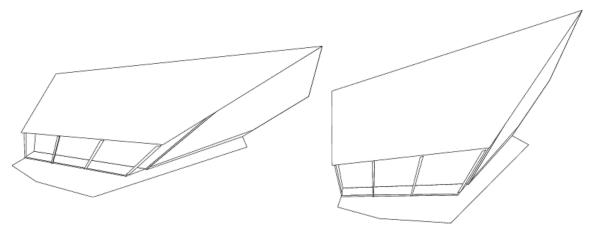
Q is a general full-rank 4×4 matrix



Projective ambiguity



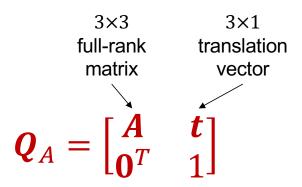


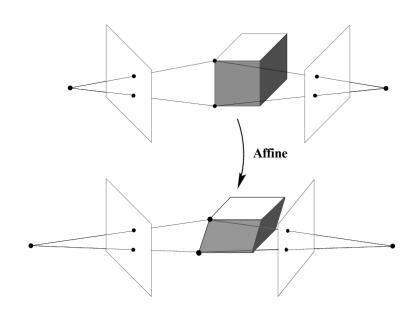


Affine ambiguity

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

$$x \cong PX = (PQ_A^{-1})(Q_AX)$$



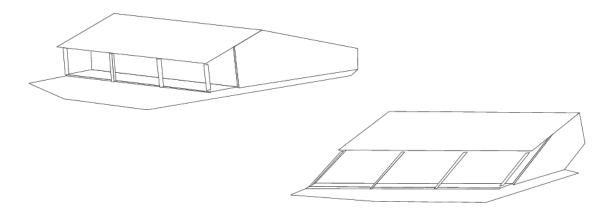


Affine ambiguity







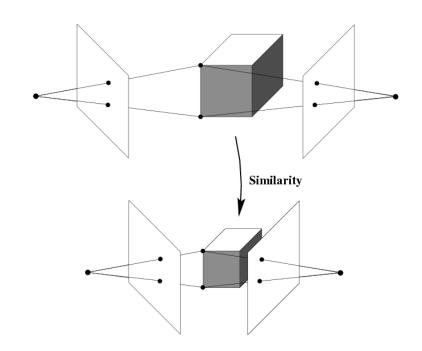


Similarity ambiguity

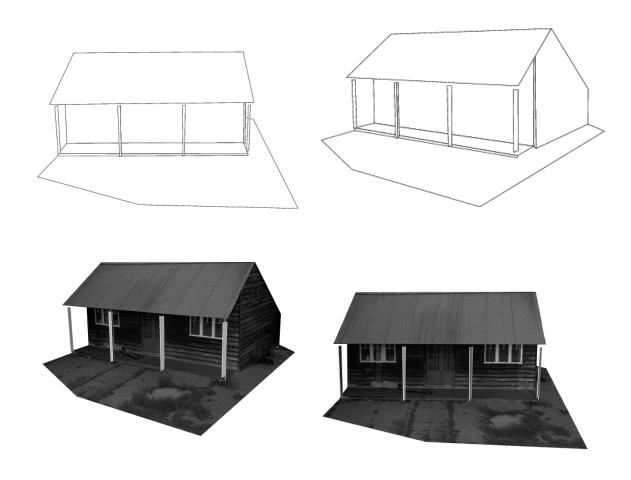
 A reconstruction that obeys orthogonality constraints on camera parameters and/or scene

$$x \cong PX = (PQ_S^{-1})(Q_SX)$$

$$3 \times 3 \qquad 3 \times 1$$
rotation translation vector
$$Q_S = \begin{bmatrix} SR & t \\ OT & 1 \end{bmatrix}$$



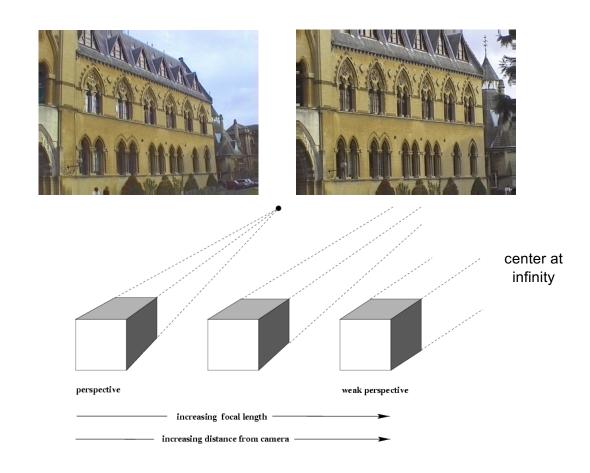
Similarity ambiguity



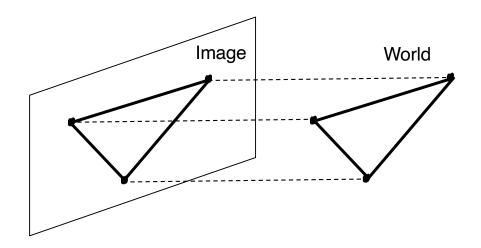
Outline: Structure from motion

- Problem definition and ambiguities
- Affine structure from motion
 - Factorization

• Let's start with affine or weak perspective cameras



Recall: Orthographic projection



Just drop the z coordinate!

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{bmatrix} & & \\ & & \\ 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$

General affine projection

 A general affine projection is a 3D-to-2D linear mapping plus translation:

$$\mathbf{P} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & t_1 \\ a_{21} & a_{22} & a_{23} & t_2 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{0}^T & 1 \end{bmatrix}$$

 a_1 , a_2 : rows of projection matrix

In non-homogeneous coordinates:

• **Given**: *m* images of *n* fixed 3D points such that

$$x_{ij} = A_i X_j + t_i, \quad i = 1, ..., m, j = 1, ..., n$$

- **Problem**: use the mn correspondences x_{ij} to estimate m projection matrices A_i and translation vectors t_i , and n points X_j
- The reconstruction is defined up to an arbitrary affine transformation Q (12 degrees of freedom):

$$\begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{0}^T & 1 \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{0}^T & 1 \end{bmatrix} \mathbf{Q}^{-1}, \qquad \begin{pmatrix} \mathbf{X}_j \\ 1 \end{pmatrix} \rightarrow \mathbf{Q} \begin{pmatrix} \mathbf{X}_j \\ 1 \end{pmatrix}$$

- How many knowns and unknowns for m images and n points?
 - 2mn knowns and 8m + 3n unknowns
 - To be able to solve this problem, we must have $2mn \ge 8m + 3n 12$ (affine ambiguity takes away 12 dof)
 - E.g., for two views, we need four point correspondences

 First, center the data by subtracting the centroid of the image points in each view:

$$\widehat{x}_{ij} = x_{ij} - \frac{1}{n} \sum_{k=1}^{n} x_{ik}$$

$$= A_i X_j + t_i - \frac{1}{n} \sum_{k=1}^{n} (A_i X_k + t_i)$$

$$= A_i \left(X_j - \frac{1}{n} \sum_{k=1}^{n} X_k \right)$$

$$= A_i \widehat{X}_j$$

• After centering, each normalized 2D point \hat{x}_{ij} is related to the 3D point by

$$\widehat{\mathbf{x}}_{ij} = A_i \widehat{\mathbf{X}}_j$$

 We can get rid of the need to center the 3D data (and the translation ambiguity) by defining the origin of the world coordinate system as the centroid of the 3D points

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

$$\boldsymbol{D} = \begin{bmatrix} \widehat{\boldsymbol{x}}_{11} & \widehat{\boldsymbol{x}}_{12} & \cdots & \widehat{\boldsymbol{x}}_{1n} \\ \widehat{\boldsymbol{x}}_{21} & \widehat{\boldsymbol{x}}_{22} & \cdots & \widehat{\boldsymbol{x}}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \widehat{\boldsymbol{x}}_{m1} & \widehat{\boldsymbol{x}}_{m2} & \cdots & \widehat{\boldsymbol{x}}_{mn} \end{bmatrix} \quad \text{cameras} \quad (2m) \quad \widehat{\boldsymbol{x}}_{ij} = \boldsymbol{A}_i \boldsymbol{X}_j$$
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.

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

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

$$\boldsymbol{M}$$

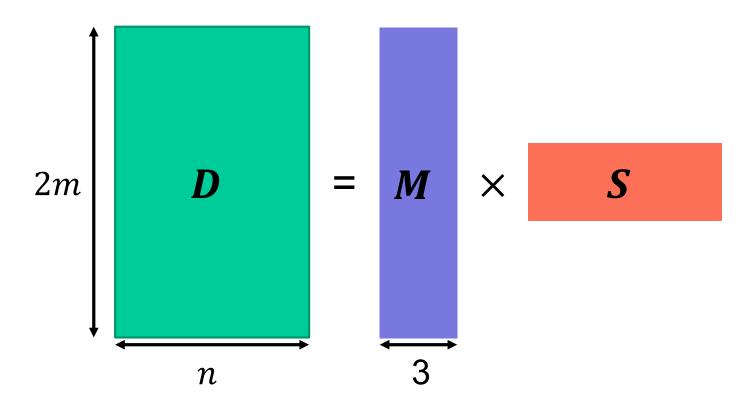
$$\boldsymbol{Cameras}$$

$$(2m \times 3)$$

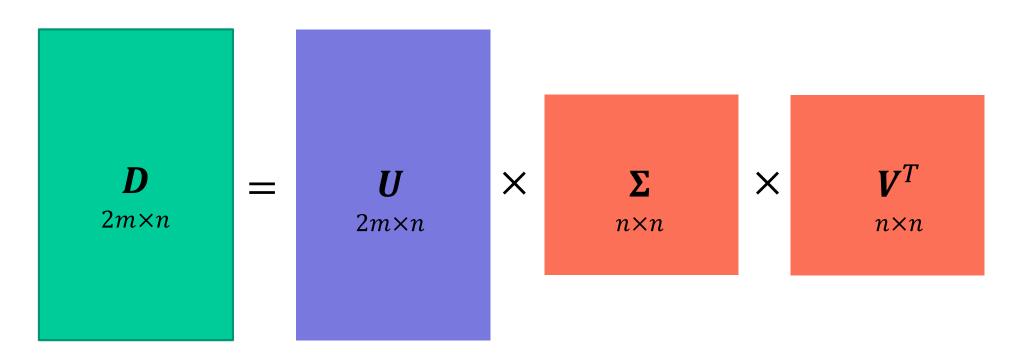
What must be the rank of the measurement matrix D = MS?

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

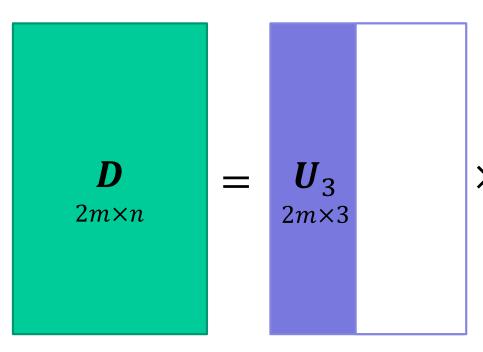
• We want:



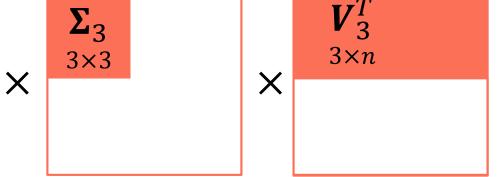
• Perform SVD of **D**:



Keep top 3 singular values:

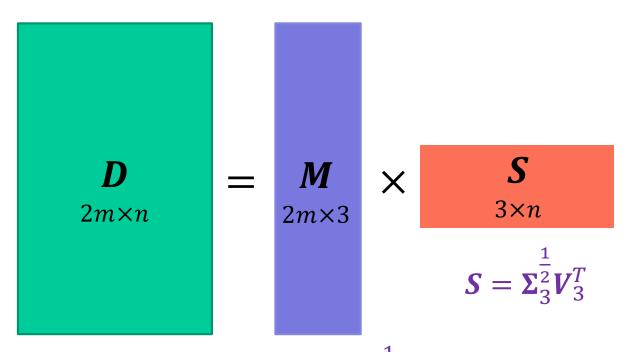


This is the closest approximation of D
with a rank-3 matrix in terms of
Frobenius norm



- What to do about Σ_3 ?
 - One solution: $\mathbf{M} = \mathbf{U}_3 \mathbf{\Sigma}_3^{\frac{1}{2}}, \mathbf{S} = \mathbf{\Sigma}_3^{\frac{1}{2}} \mathbf{V}_3^T$

One possible solution:

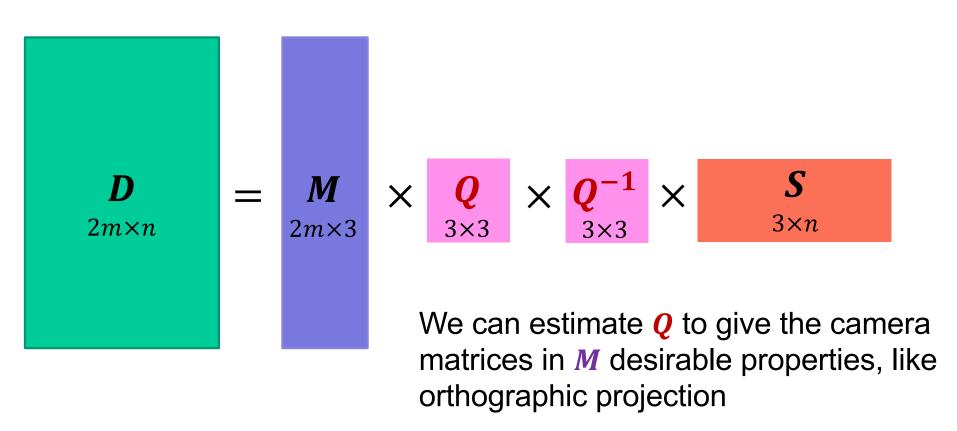


 $\boldsymbol{M} = \boldsymbol{U}_3 \boldsymbol{\Sigma}_3^{\frac{1}{2}}$

Are there other solutions?

Factorizing the measurement matrix

Other possible solutions:



Eliminating the affine ambiguity

So far, we have obtained one solution:

$$\mathbf{D} = \begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_m \end{bmatrix} \begin{bmatrix} X_1 & X_2 & \cdots & X_n \end{bmatrix}$$

$$2m \times 3$$

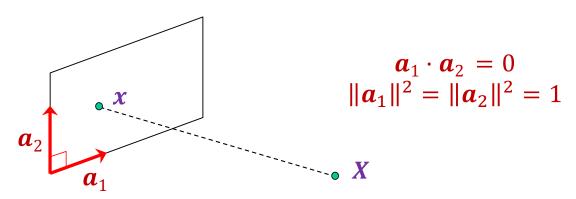
We want:

$$\boldsymbol{D} = \begin{bmatrix} A_1 \boldsymbol{Q} \\ A_2 \boldsymbol{Q} \\ \vdots \\ A_m \boldsymbol{Q} \end{bmatrix} [\boldsymbol{Q}^{-1} \boldsymbol{X}_1 \quad \boldsymbol{Q}^{-1} \boldsymbol{X}_2 \quad \cdots \quad \boldsymbol{Q}^{-1} \boldsymbol{X}_n]$$

such that each camera matrix $A_i Q$ represents orthographic projection, i.e., has orthonormal axes (rows)

Eliminating the affine ambiguity

• Let a_1 and a_2 be the rows of a 2×3 orthographic projection matrix. Then

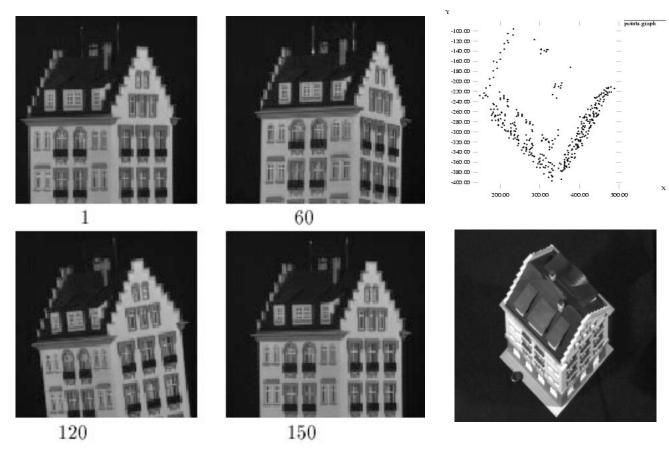


• This translates into 3m constraints on the 9 entries of Q:

$$(\boldsymbol{A}_{i}\boldsymbol{Q})(\boldsymbol{A}_{i}\boldsymbol{Q})^{T} = \boldsymbol{A}_{i}(\boldsymbol{Q}\boldsymbol{Q}^{T})\boldsymbol{A}_{i}^{T} = \boldsymbol{I}_{2\times2}, \qquad i = 1, ..., m$$

- Are the constraints linear?
- First, solve for $L = QQ^T$
- Recover *Q* from *L* by Cholesky decomposition
- Update M to MQ, S to $Q^{-1}S$

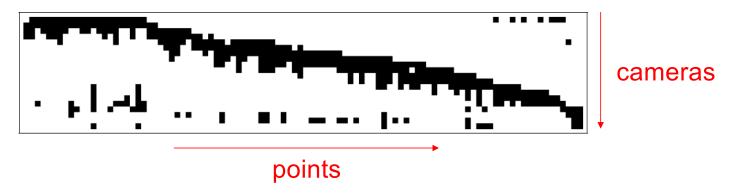
Reconstruction results



C. Tomasi and T. Kanade, <u>Shape and motion from image streams under orthography:</u>
<u>A factorization method</u>, IJCV 1992

Dealing with missing data

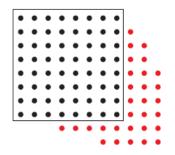
- So far, we have assumed that all points are visible in all views
- In reality, the measurement matrix typically looks something like this:



- Possible solution: decompose matrix into dense sub-blocks, factorize each sub-block, and fuse the results
 - Unfortunately, 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:



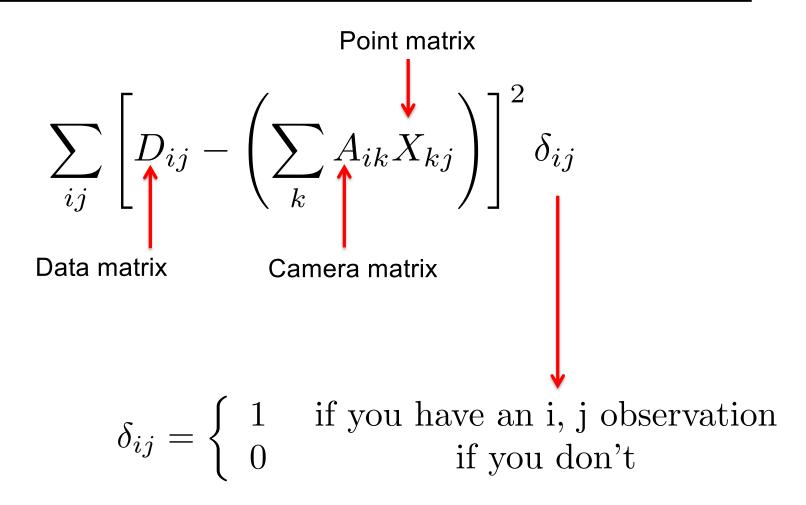
Perform factorization on a dense sub-block

Solve for a new 3D point visible by at least two known cameras – triangulation

Solve for a new camera that sees at least three known 3D points – calibration

F. Rothganger et al. Segmenting, Modeling, and Matching Video Clips Containing Multiple Moving Objects. PAMI 2007.

There is an optimization problem here



Outline: Structure from motion

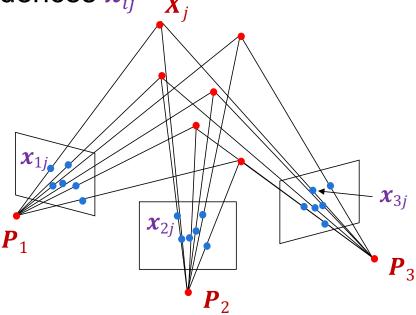
- Problem definition and ambiguities
- Affine structure from motion
 - Factorization
- Projective structure from motion

Projective structure from motion

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

$$\boldsymbol{x}_{ij} \cong \boldsymbol{P}_i \boldsymbol{X}_j$$
, $i = 1, ..., m$, $j = 1, ..., n$

• **Problem**: estimate m projection matrices P_i and n 3D points X_j from the mn correspondences x_{ij}



Projective structure from motion

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

$$\boldsymbol{x}_{ij} \cong \boldsymbol{P}_i \boldsymbol{X}_j, i = 1, ..., m, j = 1, ..., n$$

- **Problem**: estimate m projection matrices P_i and n 3D points X_j from the mn correspondences x_{ij}
- With no calibration info, cameras and points can only be recovered up to a 4×4 projective transformation Q:

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

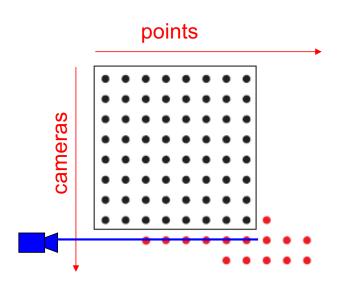
- We can solve for structure and motion when $2mn \ge 11m + 3n 15$
- For two cameras, at least 7 points are needed

Projective SFM: Two-camera case

- 1. Estimate fundamental matrix *F* between the two views
- 2. Set first camera matrix to [I | 0]
- 3. Then the second camera matrix is given by $[A \mid t]$ where t is the epipole $(F^Tt = 0)$ and $A = -[t_{\times}]F$
- In practice, SFM pipelines use guesses of intrinsic parameters and the <u>five-point algorithm</u>

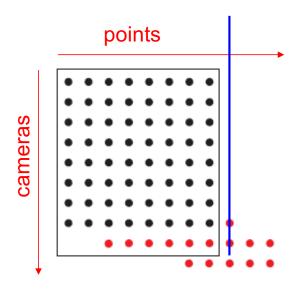
Incremental structure from motion

- 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



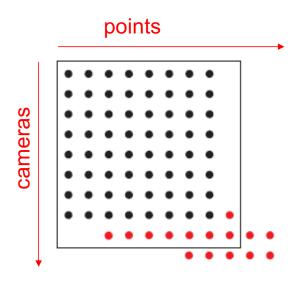
Incremental structure from motion

- 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



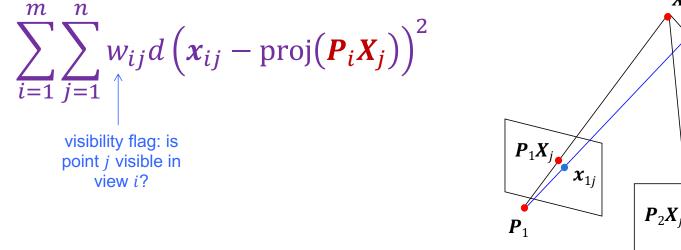
Incremental structure from motion

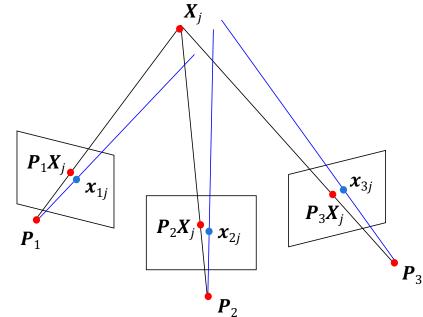
- Initialize motion from two images using fundamental matrix
- Initialize structure by triangulation
- For each additional view:
 - Determine projection matrix of new camera using all the known 3D points that are visible in its image – calibration
 - Refine and extend structure: compute new 3D points, 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 (with lots of bells and whistles):





B. Triggs et al. <u>Bundle adjustment – A modern synthesis</u>. International Workshop on Vision Algorithms, 1999

Outline: Structure from motion

- Problem definition and ambiguities
- Affine structure from motion
 - Factorization
- Projective structure from motion
 - Incremental reconstruction, bundle adjustment
- Modern structure from motion pipeline

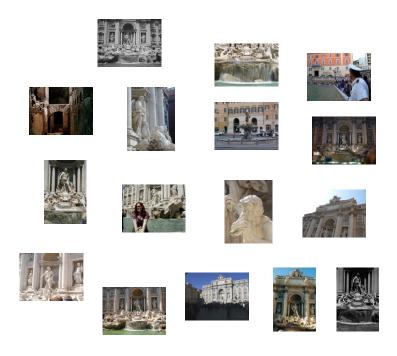
Representative SFM pipeline



N. Snavely, S. Seitz, and R. Szeliski. http://phototour.cs.washington.edu/

Feature detection

Detect SIFT features



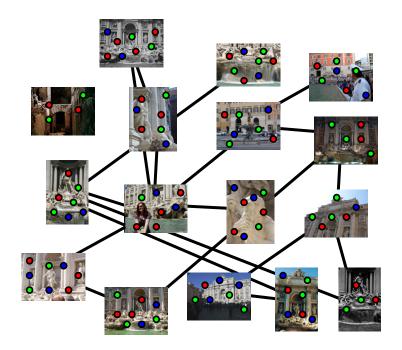
Feature detection

Detect SIFT features

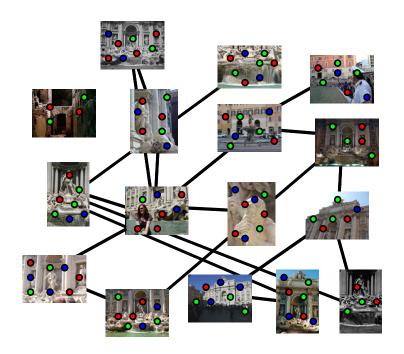
Other popular feature types: <u>SURF</u>, <u>ORB</u>, <u>BRISK</u>, ...



Match features between each pair of images



Use RANSAC to estimate fundamental matrix between each pair



Use RANSAC to estimate fundamental matrix between each pair





Use RANSAC to estimate fundamental matrix between each pair

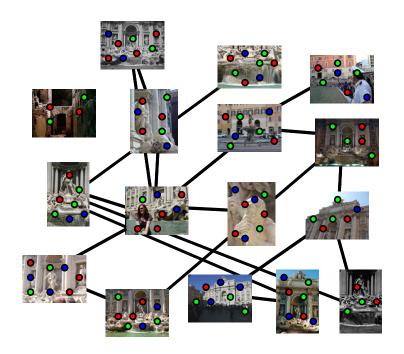
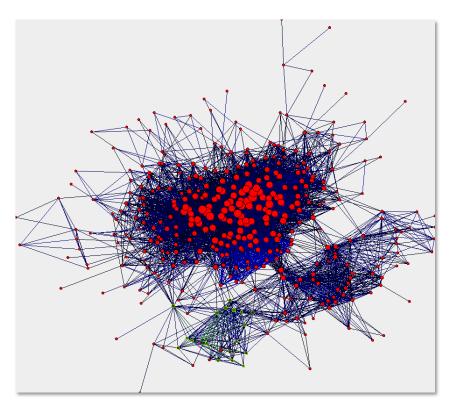


Image connectivity graph



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

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 (R and t) using <u>five-point algorithm</u>
 - 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

The devil is in the details

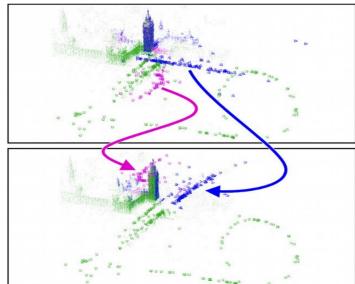
- Handling degenerate configurations (e.g., homographies)
- Filtering out incorrect matches
- Dealing with repetitions and symmetries

Repetitive structures cause catastrophic failures





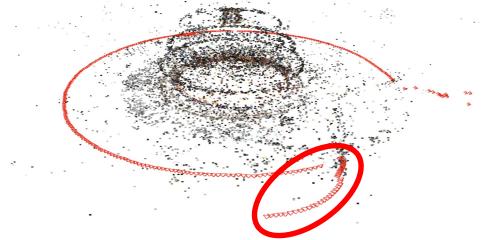




https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf

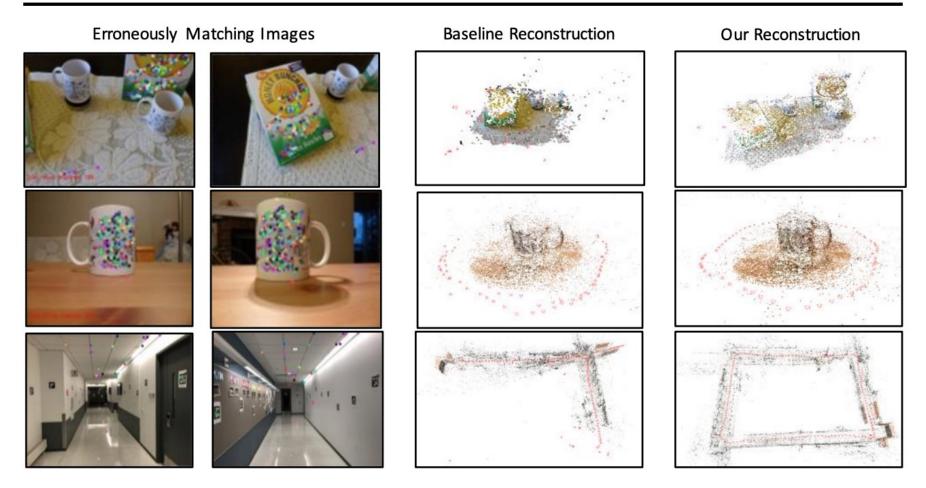
Repetitive structures cause catastrophic failures





R. Kataria et al. <u>Improving Structure from Motion with Reliable Resectioning</u>. 3DV 2020

Repetitive structures cause catastrophic failures



R. Kataria et al. Improving Structure from Motion with Reliable Resectioning. 3DV 2020

The devil is in the details

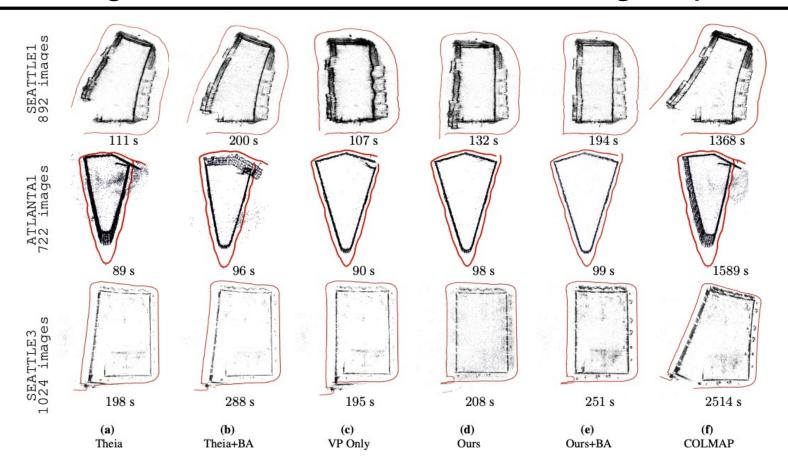
- Handling degenerate configurations (e.g., homographies)
- Filtering out incorrect matches
- Dealing with repetitions and symmetries
- Reducing error accumulation and closing loops

Reducing error accumulation and closing loops



A. Holynski et al. Reducing Drift in Structure From Motion Using Extended Features. arXiv 2020

Reducing error accumulation and closing loops



A. Holynski et al. Reducing Drift in Structure From Motion Using Extended Features. arXiv 2020

The devil is in the details

- Handling degenerate configurations (e.g., homographies)
- Filtering out incorrect matches
- Dealing with repetitions and symmetries
- Reducing error accumulation and closing loops
- Making the whole thing efficient!
 - See, e.g., Towards Linear-Time Incremental Structure from Motion

SFM software

- Bundler
- OpenSfM
- OpenMVG
- VisualSFM
- COLMAP
- See also Wikipedia's list of toolboxes