



Lecture 14 – Structure-from-Motion

Dec 18th, 2018

Danping Zou,
Associate Professor
Institute for Sensing and Navigation



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



Structure from motion



- We have learnt the most of knowledge to build a **Structure from Motion (SfM)** system.
 - Single view :
 - Perspective projection
 - Camera intrinsic, distortion
 - Camera calibration
 - Pose estimation
 - Two view :
 - Epipolar constraint
 - Essential/Fundamental matrix, RANSAC
 - Triangulation
 - Stereo rectify, calibration, and matching
 - Mathematics:
 - Homogeneous linear system (SVD)
 - Nonlinear least square (Gauss-Newton, Levenberg Marquadt)



Outline

- What's structure-from-motion problem
 - Preliminarily analysis
- Factorization approach
 - Orthographic projection
 - SVD factorization get initial decomposition
 - Enforce the orthogonality for rotation matrices
 - Get the final factorization
- Incremental approach
 - Initialization
 - Pose estimation
 - Triangulation
 - Bundle adjustment





What is structure from motion ?

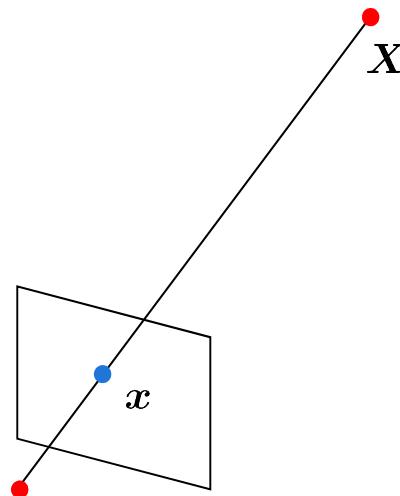


- **Structure from motion** is problem that tries to recover both 3D point clouds and camera poses from multiple views of 2D observations.
 - Inputs:
 - 2D feature points
 - Outputs:
 - Camera pose for each view
 - 3D coordinates for each point



Structure from motion

- Single view imaging model – perspective projection



$$\boldsymbol{x} \sim [\mathbf{R} \mid \boldsymbol{t}] \boldsymbol{X} \quad (\boldsymbol{m} \sim \mathbf{K} \boldsymbol{x})$$



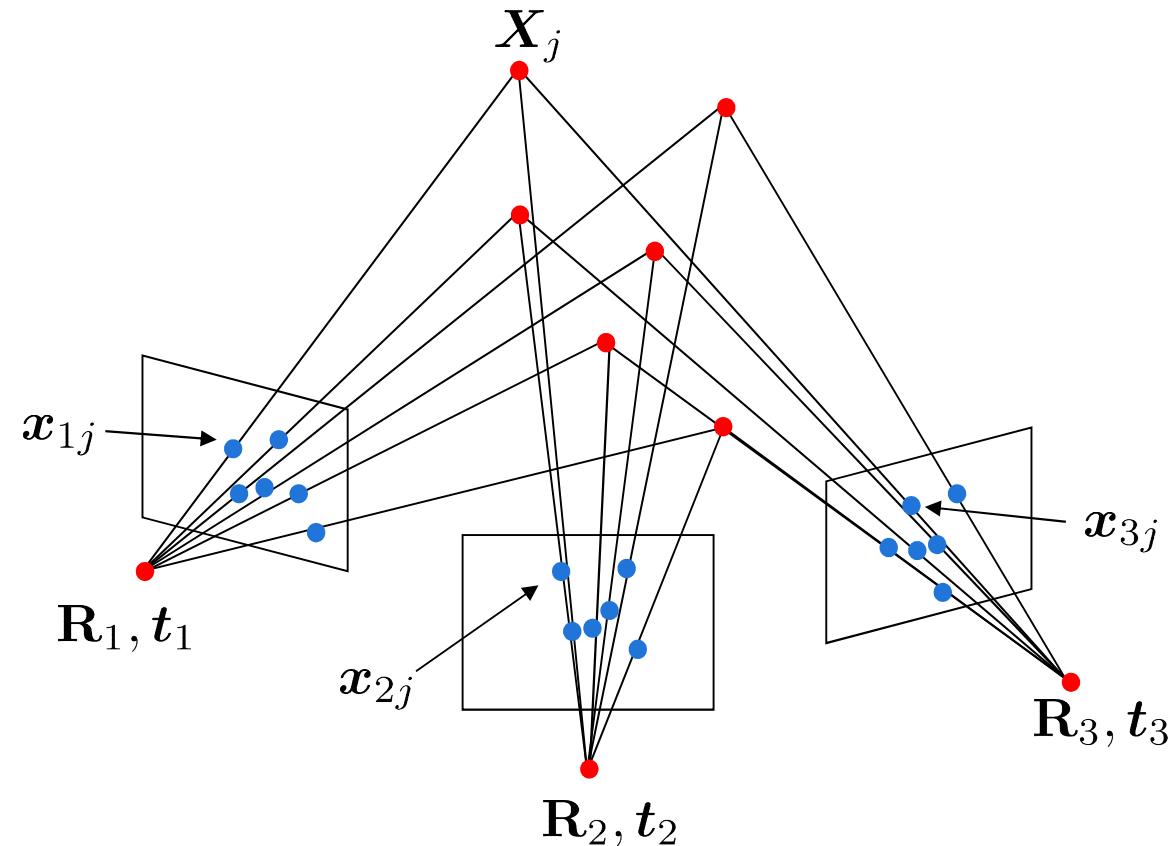
Pixel coordinates-> Image plane

$$\boldsymbol{x} = \begin{bmatrix} \frac{\mathbf{R}_{(1,:)} \boldsymbol{X} + t_1}{\mathbf{R}_{(3,:)} \boldsymbol{X} + t_3} \\ \frac{\mathbf{R}_{(2,:)} \boldsymbol{X} + t_2}{\mathbf{R}_{(3,:)} \boldsymbol{X} + t_3} \end{bmatrix} \text{ where } \boldsymbol{t} = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix}$$



Structure from motion

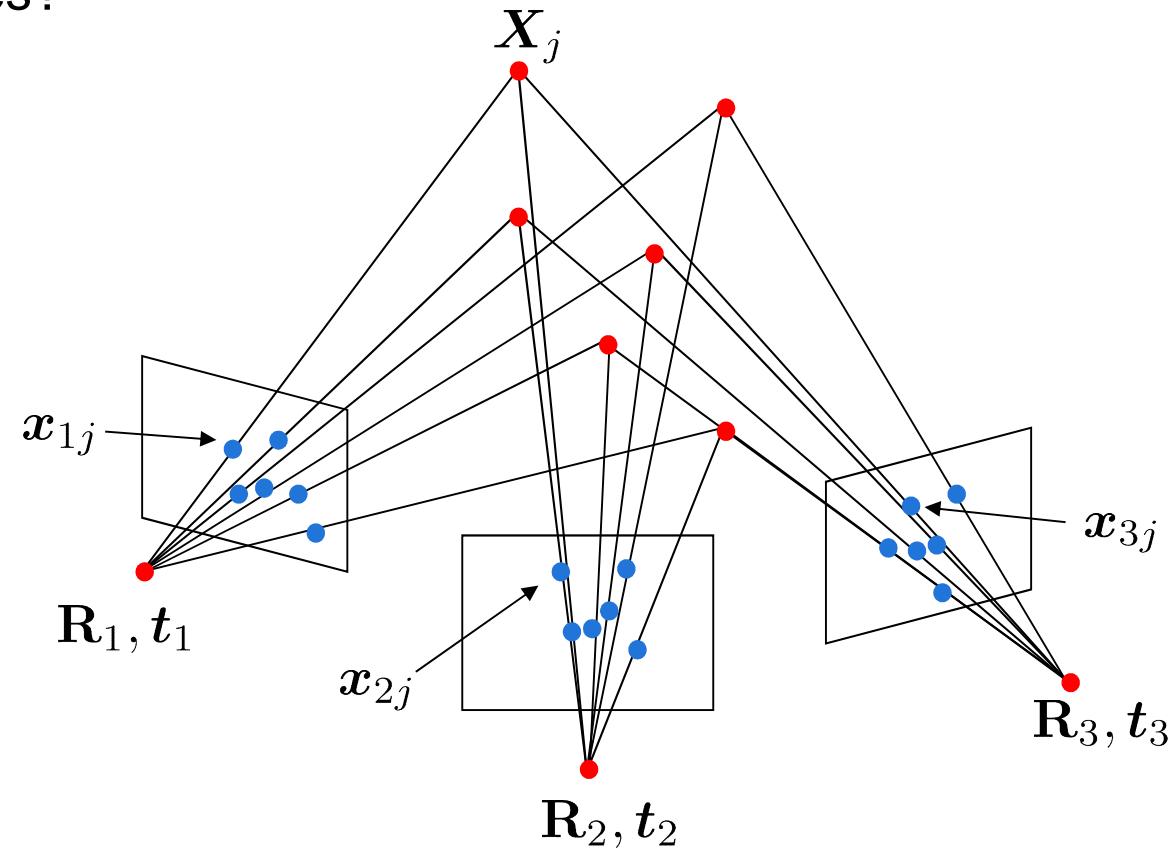
- If we have multiple views of 2D observations,





Structure from motion

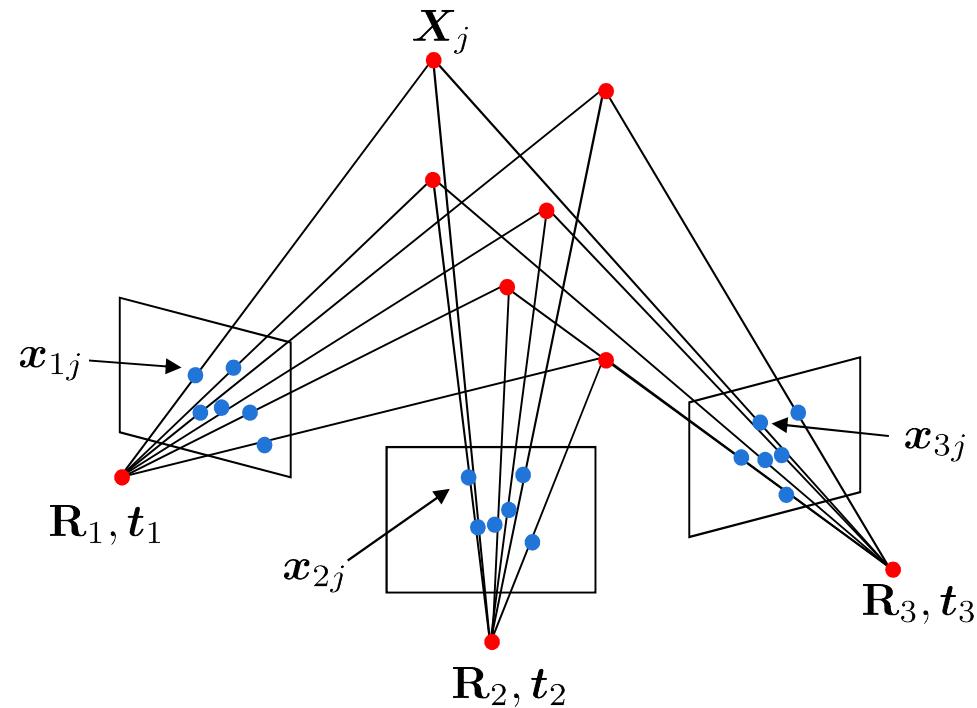
- How do we compute both the 3D point coordinates and the camera poses?





Structure from motion

- Suppose there are N 3D points and M views.
- Each point is observable at every views.
- The number of 2D observations is $N \times M$





Can we solve this problem ?



Structure from motion



- The number of unknowns:
 - Camera pose (rotation : 9 unknowns + translation : 3 unknowns)
 - 3D point (3 unknowns)

$$\#unkown = 12M + 3N$$

- The number of equations:
 - Each 2D feature point gives rise to two equations (x,y)

$$\#equation = 2M \times N$$

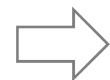


Structure from motion



- If the number of equations is not less than the number of unknowns, the problem can be solved.

$$\begin{aligned}\#\text{equation} &= 2M \times N \\ &\geq \#\text{unkown} = 12M + 3N \\ (2N - 12)M &\geq 3N\end{aligned}$$



$$N > 6, M \geq \frac{3N}{2N-12}$$

$$N = 7, M \geq 11$$

$$N = 10, M \geq 4$$

...

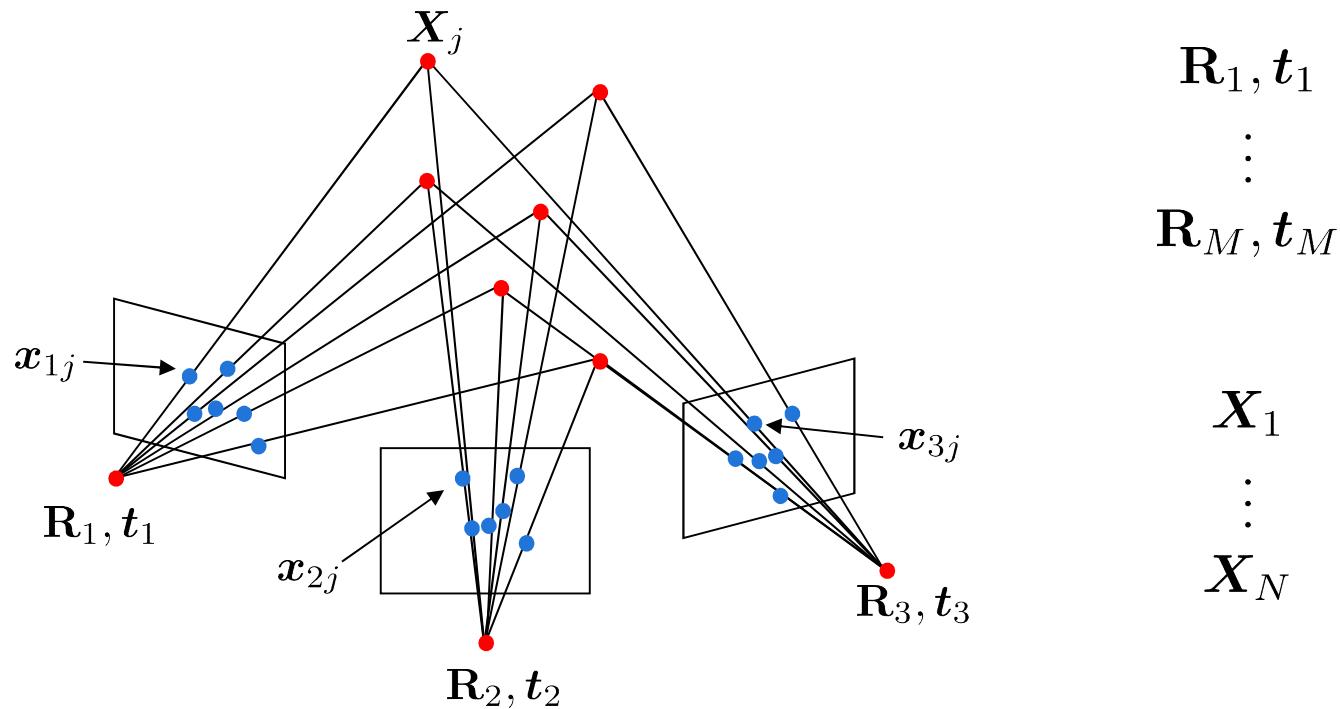
$$N = 30, M \geq 2$$

The more points, the fewer views are required.



Structure from motion

- It is possible to solve this problem, but it is a highly nonlinear problem.
- The key is to get the initial solution.





Structure-from-motion



1990

Factorization/Batch

Tomasi C, Kanade T

Shape and motion from
image streams under
orthography: a factorization
method, IJCV, 1992

2000

Incremental approach

- Photo Tourism (Bundler)
- Rome in a Day





上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



Factorization approach



Structure from motion



▪ Factorization method

- Using orthographic or affine projection instead of perspective projection
- Close-form solution to the 3D structure and the camera poses all together

Tomasi C, Kanade T. Shape and motion from image streams under orthography: a factorization method, IJCV, 1992

《像外行一样思考,像专家一样实践》
-金出武雄

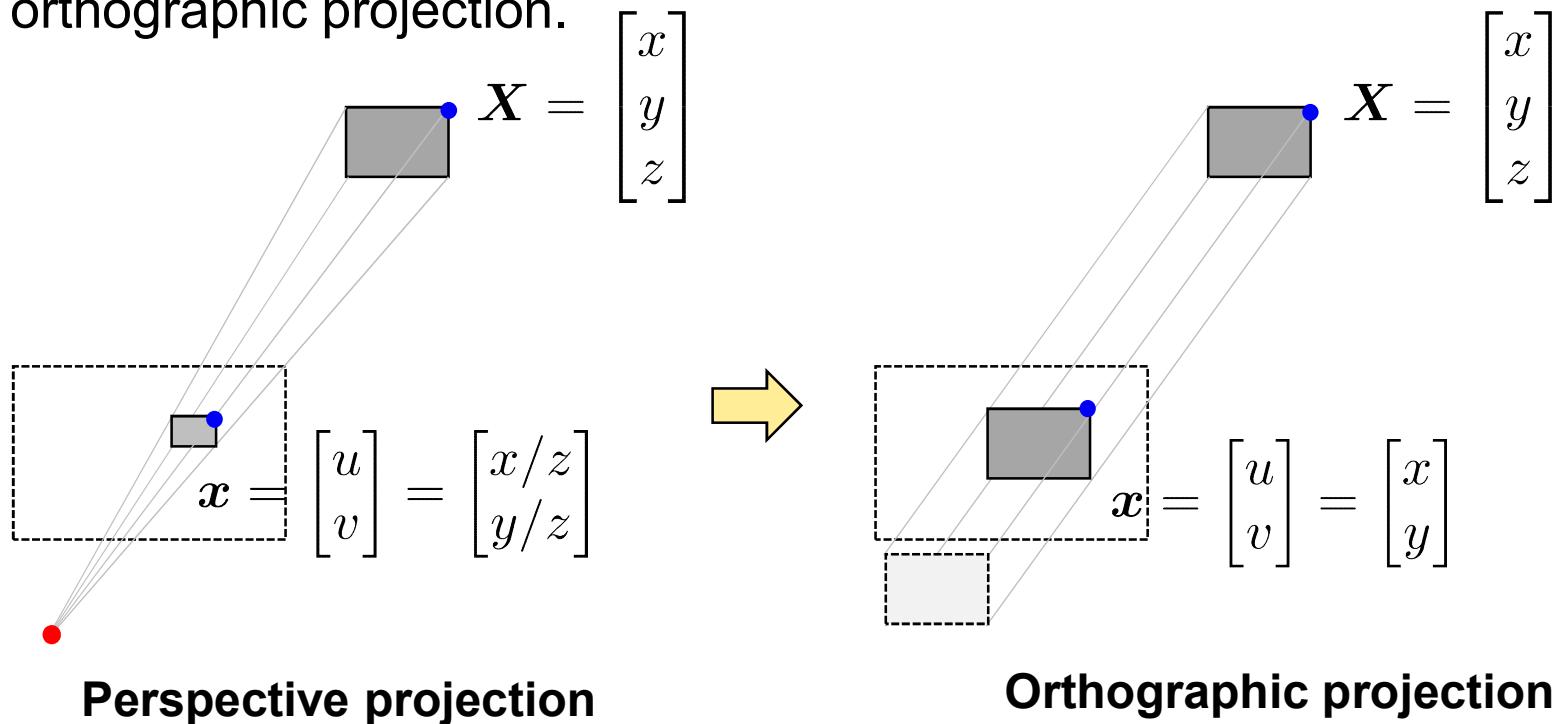


金出武雄
(The Robotics Institute, CMU)



Structure from motion

- The key idea is to approximate the perspective projection by the orthographic projection.

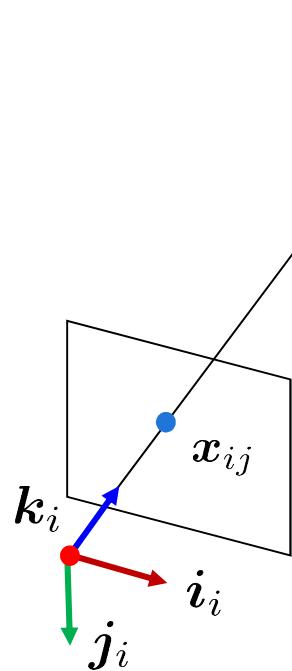


The approximation is better when the object is far away.



Factorization methods

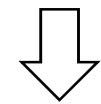
- The camera pose of the i-th view is represented by $(\mathbf{R}_i, \mathbf{t}_i)$,
(where $\mathbf{R}_i = [i, j, k]^T$) which is the transformation from the world frame to the camera frame.



$$\tilde{x}_{ij} \sim \mathbf{R}_i \quad \mathbf{X}_j + \mathbf{t}_i$$

$3 \times 1 \quad 3 \times 3 \quad 3 \times 1$

$$x_{ij} = [\mathbf{i}_i, \mathbf{j}_i]^T \mathbf{X}_j + \mathbf{t}_{i,1:2}$$



$$u_{ij} = \mathbf{i}_i^T \mathbf{X}_j + t_{i,1}$$

$$v_{ij} = \mathbf{j}_i^T \mathbf{X}_j + t_{i,2}$$



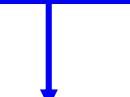
Factorization methods

- Remove the translation :
 - Choose the image origin be the centroid of 2D points
 - Choose the scene origin be the centroid of 3D points
- We can drop the camera translation

$$\boldsymbol{x}_{ij} = [\boldsymbol{i}_i, \boldsymbol{j}_i]^T \boldsymbol{X}_j + \cancel{\boldsymbol{t}_{i,1:2}}$$



$$\boldsymbol{x}_{ij} = [\boldsymbol{i}_i, \boldsymbol{j}_i]^T \boldsymbol{X}_j$$



$$\boldsymbol{x}_{ij} = \boldsymbol{\pi}_i \boldsymbol{X}_j$$



Factorization methods

- Projection of N feature points in the i-th view :

$$[x_{i1} \cdots x_{ij} \cdots x_{iN}] = \pi_i [X_1 \cdots X_j \cdots X_N]$$

- Projection of N feature points in M views :

$$\begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1N} \\ x_{21} & x_{22} & \cdots & x_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ x_{M1} & x_{M2} & \cdots & x_{MN} \end{bmatrix} = \begin{bmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_M \end{bmatrix} [X_1, X_2 \cdots X_N]_{3 \times N}$$

$2M \times N$ $2M \times 3$

$$\mathbf{M} = \boldsymbol{\Pi} \mathbf{S}$$



Factorization methods



- We can write the measurement equations by

$$\text{Known } \mathbf{M}_{2M \times N} = \mathbf{\Pi}_{2M \times 3} \mathbf{S}_{3 \times N} \text{ Unknowns to be solved}$$

- Our goal is to factorize the observation matrix \mathbf{M} to get the rotation matrix and the shape matrix.
- We observed that \mathbf{M} is at most rank 3, since
 - Rank of $\mathbf{\Pi}$ is 3
 - Rank of \mathbf{S} is 3
 - Product of 2 matrices of rank 3 has rank 3.
- But with noise, $\text{rank}(\mathbf{M})$ might be > 3 .



Factorization methods

- SVD factorization:
 - Apply SVD to the observation matrix, we have

$$\begin{array}{cccc} \mathbf{M} & = & \mathbf{U} & \Sigma & \mathbf{V}^T \\ 2M \times N & & 2M \times N & N \times N & N \times N \end{array}$$

- As \mathbf{M} has rank 3, it has three non-zero singular values, we get corresponding columns of \mathbf{U} and \mathbf{V} .

$$\mathbf{M} = [\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3, \dots] \text{diag}(\sigma_1, \sigma_2, \sigma_3, 0, \dots) [\mathbf{V}_1, \mathbf{V}_2, \mathbf{V}_3, \dots]^T$$

$$\begin{array}{ccccc} \mathbf{M} & = & \mathbf{U}_{1:3} & \text{diag}(\sigma_1, \sigma_2, \sigma_3) & \mathbf{V}_{1:3}^T \\ 2M \times N & & 2M \times 3 & 3 \times 3 & 3 \times N \\ & & \downarrow & \downarrow & \\ & & \mathbf{M} = \mathbf{\Pi}' & \mathbf{S}' & \end{array}$$



Factorization methods

- Π' differ from Π by a linear transformation A :

$$\begin{aligned} M &= \Pi S \\ M &= \Pi' S' \\ \Pi' &= \Pi A^{3 \times 3} \end{aligned}$$

- If we can find $A \in \mathbb{R}^{3 \times 3}$, we can solve both Π and S .
- We solve A by enforcing the orthogonality on Π .



Factorization methods

- We revisit the projection matrix Π

$$\Pi = \begin{bmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_M \end{bmatrix} \xrightarrow{\hspace{1cm}} \pi_i = \begin{bmatrix} i_i^T \\ j_i^T \end{bmatrix} \xrightarrow{i_i \perp j_i} \pi_i \pi_i^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

- The projection matrix of an orthographic camera has two orthonormal rows.
- Enforce the orthonormal constraints on each π'_i , such that :

$$\pi'_i \mathbf{A} \mathbf{A}^T \pi'^T_i = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



Factorization methods

- From $\pi'_i \mathbf{A} \mathbf{A}^T \pi'^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, letting $\mathbf{Q} = \mathbf{A} \mathbf{A}^T$, we have

$$\pi'_i = \begin{bmatrix} i'^T \\ j'^T \end{bmatrix}$$

$$\begin{aligned} i'^T \mathbf{Q} i' &= 1 \\ j'^T \mathbf{Q} j' &= 1 \\ i'^T \mathbf{Q} j' &= 0 \end{aligned}$$

$$\mathbf{Q} = \mathbf{A} \mathbf{A}^T$$

$$\mathbf{Q} = \begin{bmatrix} q_1 & q_2 & q_3 \\ q_2 & q_4 & q_5 \\ q_3 & q_5 & q_6 \end{bmatrix}$$

- Finally we can get a linear system of equations:

$$\begin{array}{ccc} \mathbf{H} & \mathbf{q} & = & \mathbf{b} \\ 3M \times 6 & 6 \times 1 & & 3M \times 1 \end{array}$$

$$(\text{where } \mathbf{q} = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_6 \end{bmatrix})$$



Factorization methods



- This is a linear least squared problem, which can be solved by :

$$\mathbf{H}q = \mathbf{b} \iff q = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H} \mathbf{b}$$

- Or using the backslash '\' operator in Matlab to get the solution:

$$q = \mathbf{H} \setminus \mathbf{b}$$

- After that, we can get \mathbf{A} through Cholesky decomposition :

In Matlab, we have

$$\mathbf{Q} = \mathbf{A} \mathbf{A}^T$$

$$\mathbf{A} = \text{chol}(\mathbf{Q})^T$$



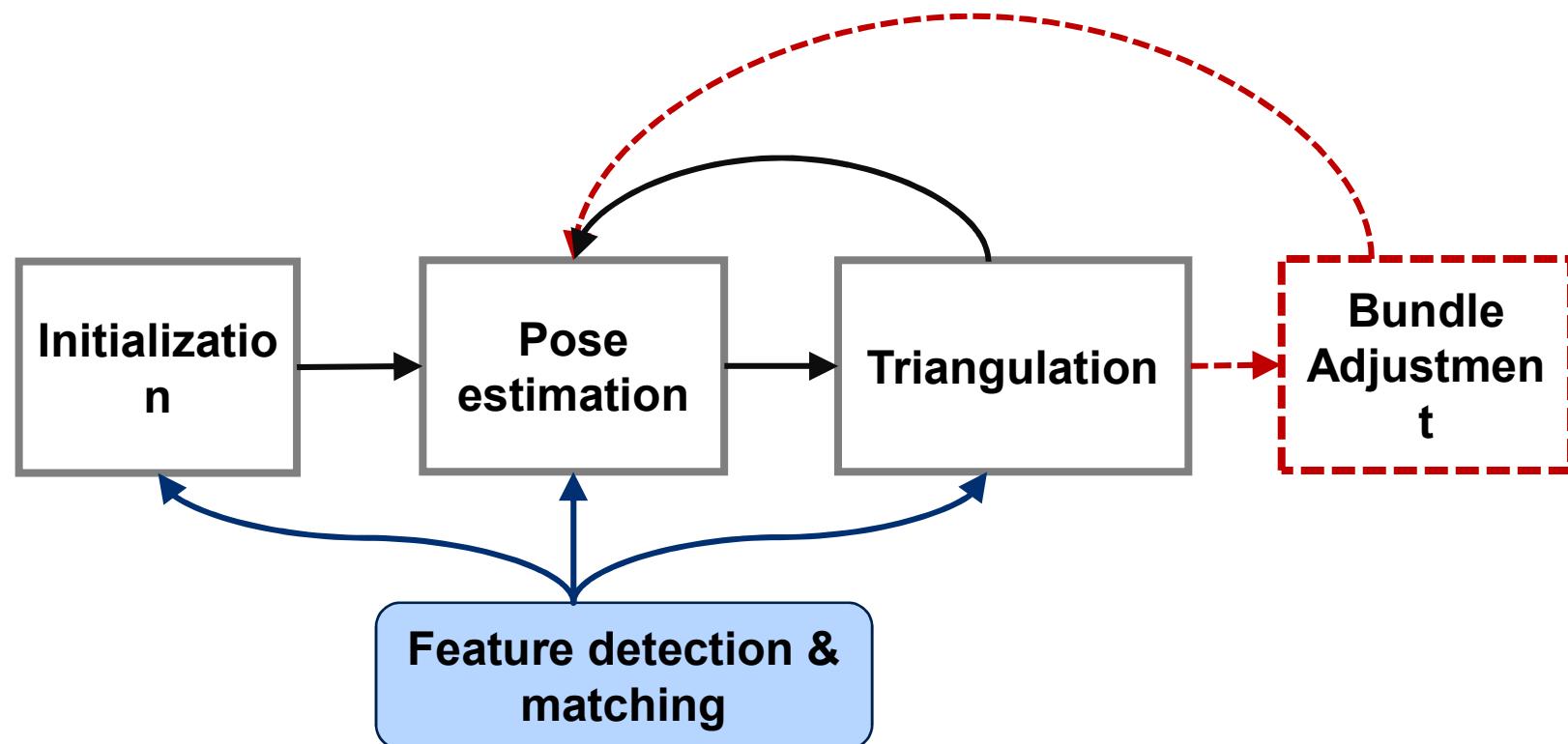
Incremental approach



Incremental SfM



- A typical pipeline of a incremental structure from motion system :

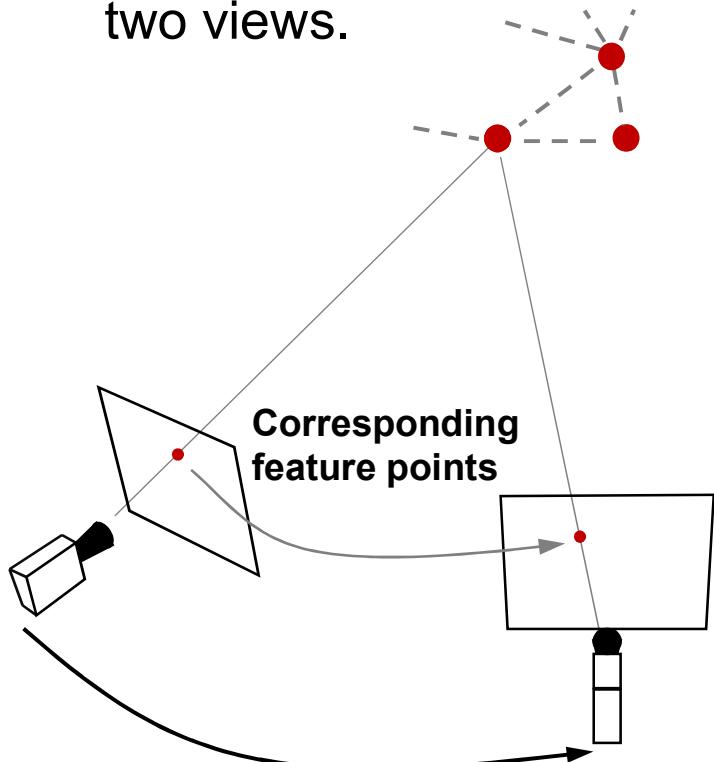




Incremental SfM



- **Initialization** – Generate the 3D points and camera poses of the initial two views.



- Feature detection & matching
 - Find correspondences between two views
 - SIFT, SURF, ORB
- Find the essential matrix
 - RANSAC, 5-point algorithm, 8-point algorithm
- Extract the camera poses of the two views
 - SVD
- Generate the initial 3D points
 - Triangulation

Initialization

Pose estimation

Triangulation

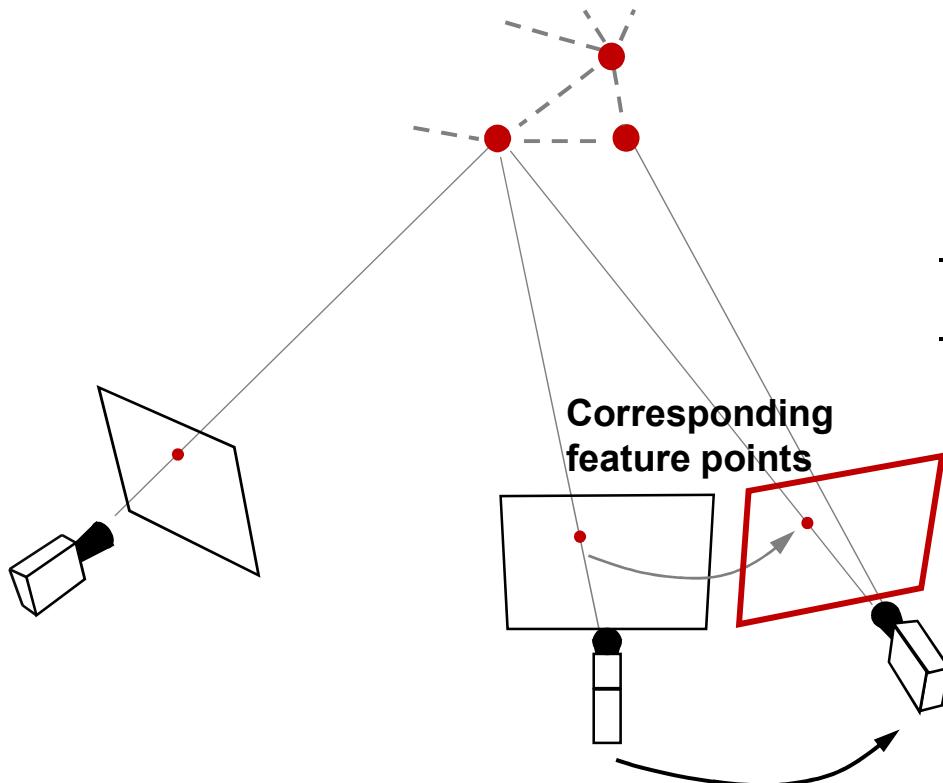
Bundle adjustment



Incremental SfM



- **Pose estimation**



- Feature detection & matching
 - Find correspondences between the current view and previous view (KLT, ORB)
- Solve the PnP problem
- Handle the outliers (RANSAC / M-estimator)

Initialization

Pose estimation

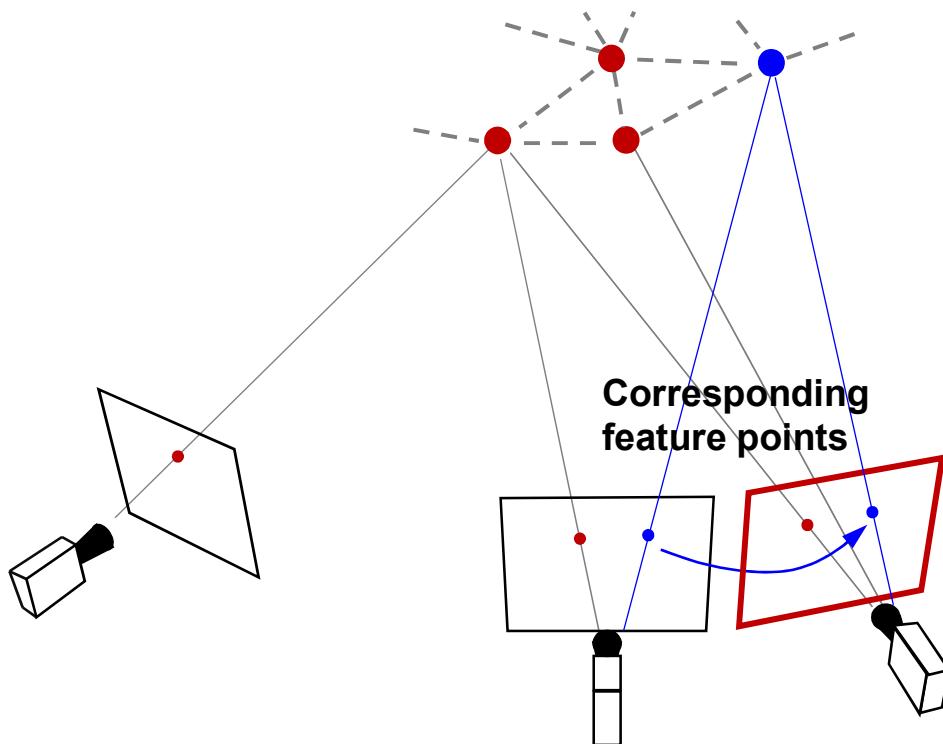
Triangulation

Bundle adjustment



Incremental SfM

- Triangulation



Generate new 3D points
from corresponding 2D
points.

$$\begin{aligned} \mathbf{x} \times \mathbf{P}\mathbf{X} &= \mathbf{0} \\ \mathbf{x}' \times \mathbf{P}'\mathbf{X} &= \mathbf{0} \end{aligned} \quad \Rightarrow \quad \mathbf{A}_{6 \times 4}\mathbf{X} = \mathbf{0}_{6 \times 1}$$

$$\begin{aligned} u &= \frac{\mathbf{P}_1^T \mathbf{X}}{\mathbf{P}_3^T \mathbf{X}} & v &= \frac{\mathbf{P}_2^T \mathbf{X}}{\mathbf{P}_3^T \mathbf{X}} \\ u' &= \frac{\mathbf{P}'_1^T \mathbf{X}}{\mathbf{P}'_3^T \mathbf{X}} & v' &= \frac{\mathbf{P}'_2^T \mathbf{X}}{\mathbf{P}'_3^T \mathbf{X}} \end{aligned} \quad \Rightarrow \quad \mathbf{A}_{4 \times 3}\mathbf{x} = \mathbf{b}_{4 \times 1}$$

Initialization

Pose estimation

Triangulation

Bundle adjustment

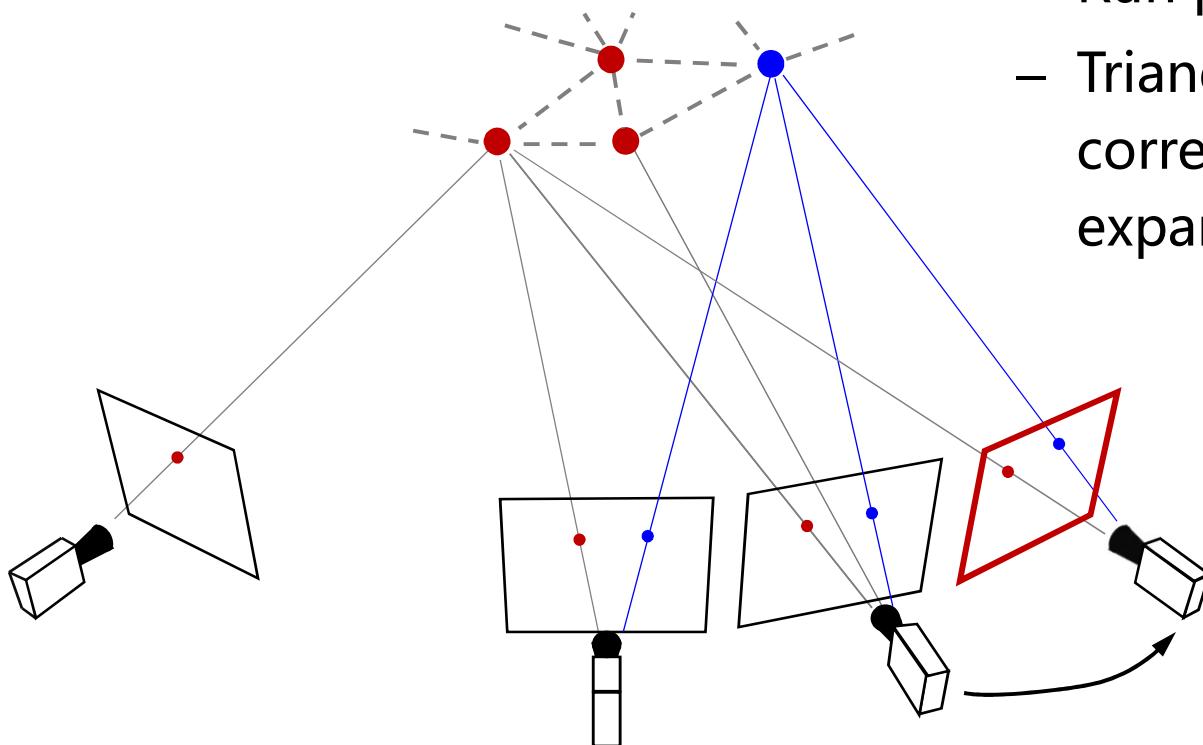


Incremental SfM



- Repeat the above steps

- Run pose estimation again
- Triangulate more feature correspondences to expand the 3D point cloud



Initialization

Pose estimation

Triangulation

Bundle adjustment



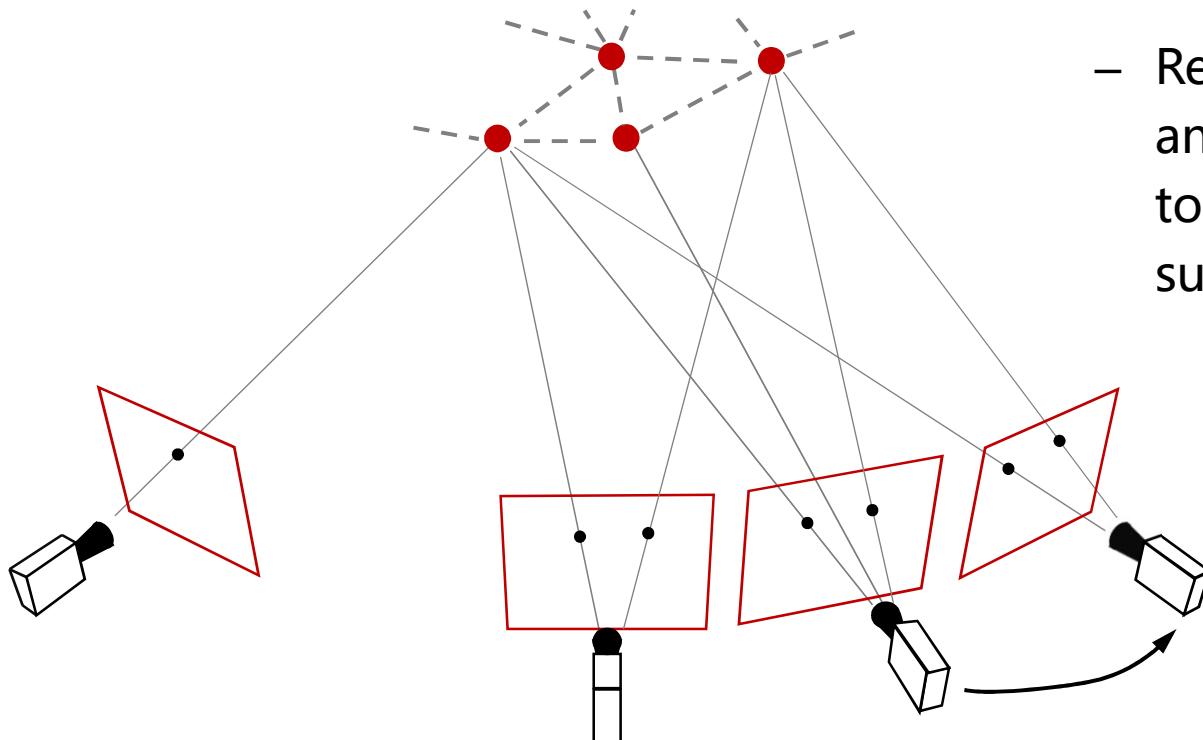
Incremental SfM



- Call bundle adjustment from time to time to refine the result

– Refine the 3D coordinates and the camera poses together by minimizing the sum of re-projection errors:

$$\sum_{ij} \mathbf{z}_{ij}^T \Sigma_{im}^{-1} \mathbf{z}_{ij}$$



Initialization

Pose estimation

Triangulation

Bundle adjustment



Bundle Adjustment



- Bundle adjustment is to refine the 3D coordinates and the camera poses together by minimizing the sum of re-projection errors:

$$\sum_{ij} \mathbf{z}_{ij}^T \Sigma_{im}^{-1} \mathbf{z}_{ij}$$

- \mathbf{z}_{ij} is the re-projection error of j -th point in the i -th view, defined as

$$\mathbf{z}_{ij} = \mathbf{m}_{ij} - \mathbf{x}_{ij} = \mathbf{m}_{ij} - \mathcal{P}(\theta_i, \mathbf{X}_j)$$

- $\Sigma_{im} \in \mathbb{R}^{2 \times 2}$ represents the covariance of observation noises (1~4 pixels)

$$\Sigma_{im} \in \begin{bmatrix} \sigma_{im}^2 & 0 \\ 0 & \sigma_{im}^2 \end{bmatrix}$$



Bundle Adjustment

- This is a nonlinear least squares problem, where the parameter vector is

$$\Theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_M \\ \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_N \end{bmatrix} = \begin{bmatrix} \mathbf{R}_1 \\ \mathbf{t}_1 \\ \vdots \\ \mathbf{R}_M \\ \mathbf{t}_M \\ \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_N \end{bmatrix}$$

$\Theta \leftarrow \Theta \boxplus \Delta\Theta$

$\mathbf{R}_i \leftarrow \mathbf{R}_i \oplus \exp(\Delta\theta^\wedge)$

$\mathbf{t}_i \leftarrow \mathbf{t}_i + \Delta\mathbf{t}_i$

$\mathbf{X}_i \leftarrow \mathbf{X}_i + \Delta\mathbf{X}_i$



Bundle Adjustment

- The normal equation becomes very large to be solved

$$\text{Gauss-Newton: } \mathbf{H}^T \mathbf{H} \Delta \mathbf{x} = \mathbf{H}^T \mathbf{r}$$

$$\text{Levenberg-Marquardt: } (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}) \Delta \mathbf{x} = \mathbf{H}^T \mathbf{r}$$

- For example, $N = 1000, M = 10$. The Jacobian matrix becomes

$$\mathbf{H} \in \mathbb{R}^{2MN \times (6N+3M)}$$

$$\mathbb{R}^{20000 \times 6030}$$



SfM systems

- Bundler (2006,2007 Noah Snavely)



c++ open source

<https://www.cs.cornell.edu/~snavely/bundler/>



SfM systems

- Photo tourism / Modeling 3D world from massive photo collection

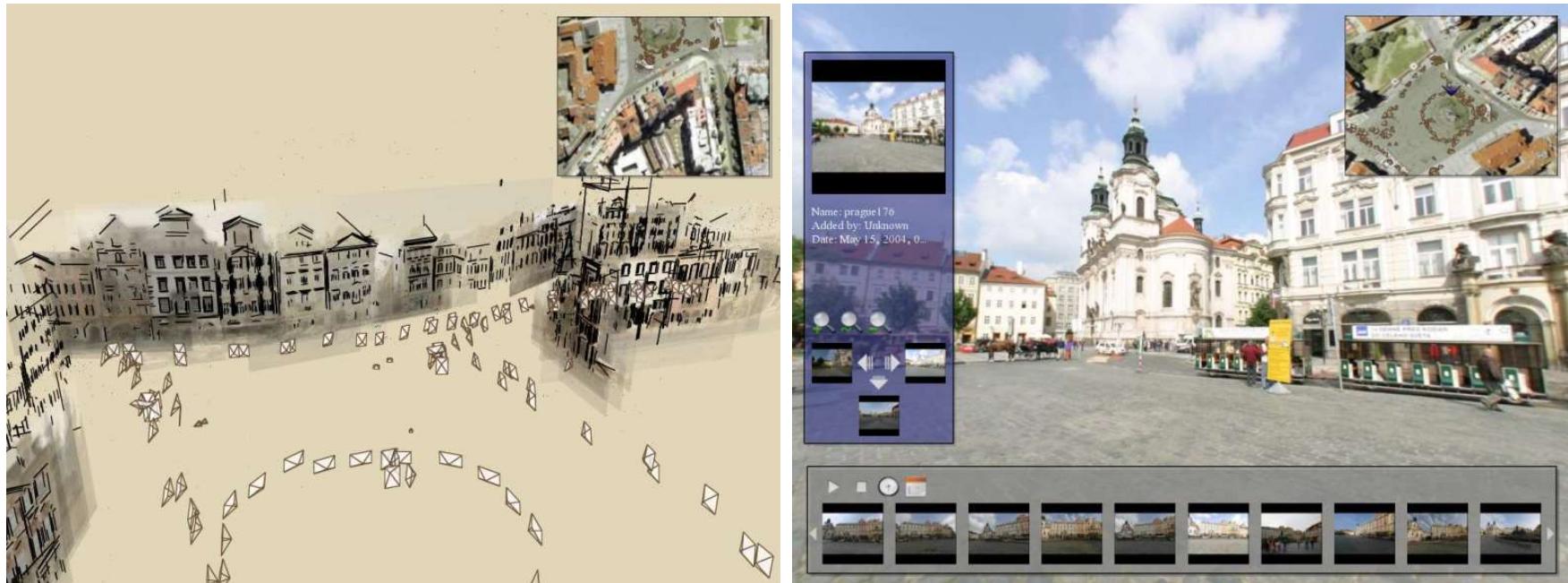


Photo Tourism

Exploring photo collections in 3D

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

SIGGRAPH 2006

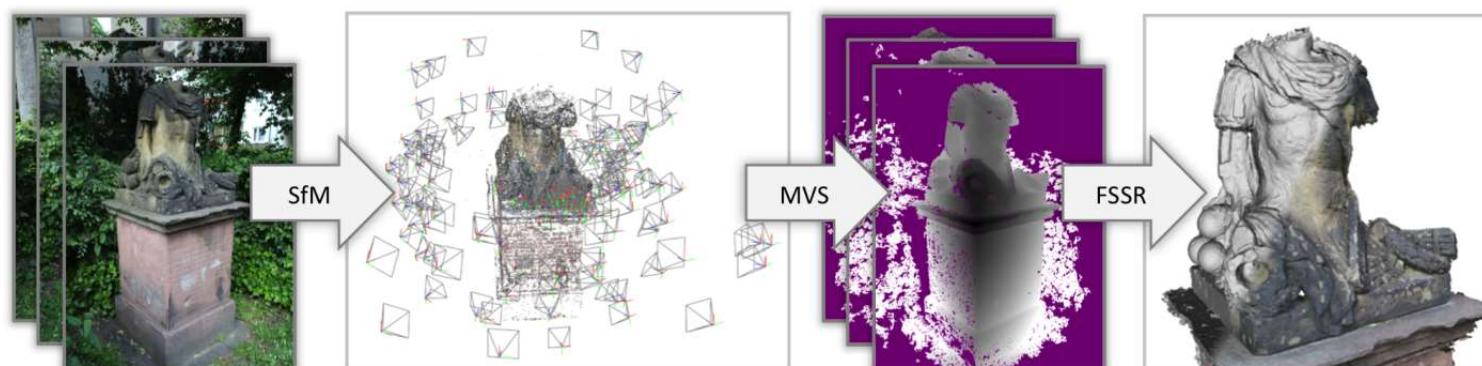


SfM systems

- Other implementations:
 - Visual SfM Changchang Wu(<http://ccwu.me/vsfm/>)



- MVE (Multi-View Reconstruction Environment)
2014, Simon Fuhrmann
(<https://github.com/simonfuhrmann/mve/wiki/MVE-Users-Guide>)





SfM systems

- OpenMVG



Pierre Moulon

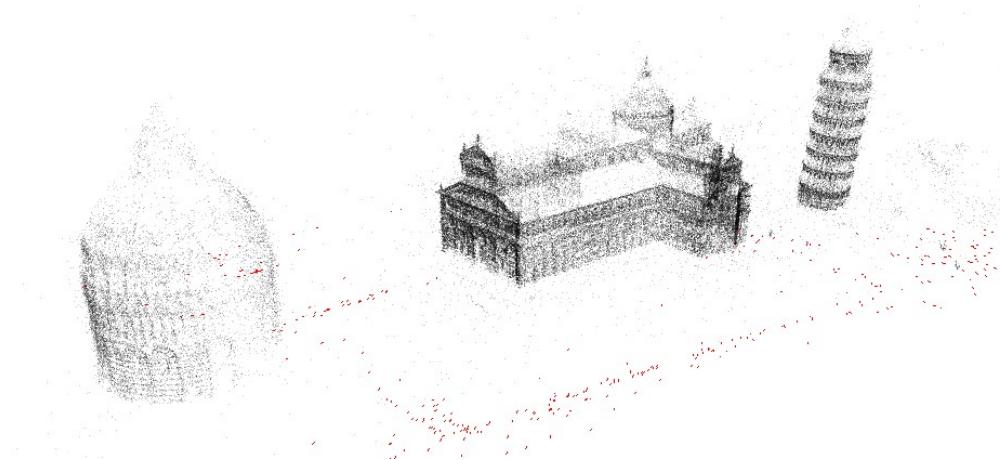
(<https://libraries.io/github/openMVG/openMVG>)

The screenshot shows the documentation page for the openMVG libraries. The top navigation bar includes a logo for 'openMVG latest', a search bar, and a link to 'Search docs'. The main content area is titled 'openMVG libraries' and contains a sidebar with links to various modules: image, numeric, features, cameras, multiview, linear programming, robust_estimation, matching, tracks, geometry, geodesy, and sfm. Below the sidebar, there's a section for 'openMVG samples', 'openMVG softwares & tools', 'patented', 'dependencies', 'third_party', 'FAQ', and 'Bibliography'. A callout box on the right side of the sidebar says 'Reach 7 million devs each month when you advertise with Read the Docs.' At the bottom of the sidebar, there's a button for 'Read the Docs' and a link to 'v:latest'.

- Theia Vision Library



(<http://www.theia-sfm.org/sfm.html>)





SfM systems



- Matlab 2018b – Computer vision system toolbox (jbox.sjtu.edu.cn)
 - Camera Calibration and 3-D Vision

bundleAdjustment	Refine camera poses and 3-D points
cameraMatrix	Camera projection matrix
cameraPoseToExtrinsics	Convert camera pose to extrinsics
epipolarLine	Compute epipolar lines for stereo images
estimateCameraParameters	Calibrate a single or stereo camera
estimateEssentialMatrix	Estimate essential matrix from corresponding points in a pair of images
estimateFundamentalMatrix	Estimate fundamental matrix from corresponding points in stereo images
estimateWorldCameraPose	Estimate camera pose from 3-D to 2-D point correspondences
extrinsics	Compute location of calibrated camera
extrinsicsToCameraPose	Convert extrinsics to camera pose
isEpipoleInImage	Determine whether image contains epipole
lineToBorderPoints	Intersection points of lines in image and image border
relativeCameraPose	Compute relative rotation and translation between camera poses
triangulate	3-D locations of undistorted matching points in stereo images
triangulateMultiview	3-D locations of undistorted points matched across multiple images
undistortImage	Correct image for lens distortion
undistortPoints	Correct point coordinates for lens distortion
cameraParameters	Object for storing camera parameters



Summary



- Structure from motion problem is **solved** for static scenes (around 2006) .
- Lots of SfM software are available nowadays, both in open source and commercial products.
- SfM system is the first step to reconstruct the dense 3D model of the world.
 - Sparse 3D points
 - Camera poses -> Dense 3D points (Multi-view stereo)
- SfM technique can be applied to SLAM problem (Next week)