



Structure from Motion (SfM)

Jiayuan Gu

gujy1@shanghaitech.edu.cn

Outline

- Problem Formulation
- Projective structure from motion
- SfM pipeline

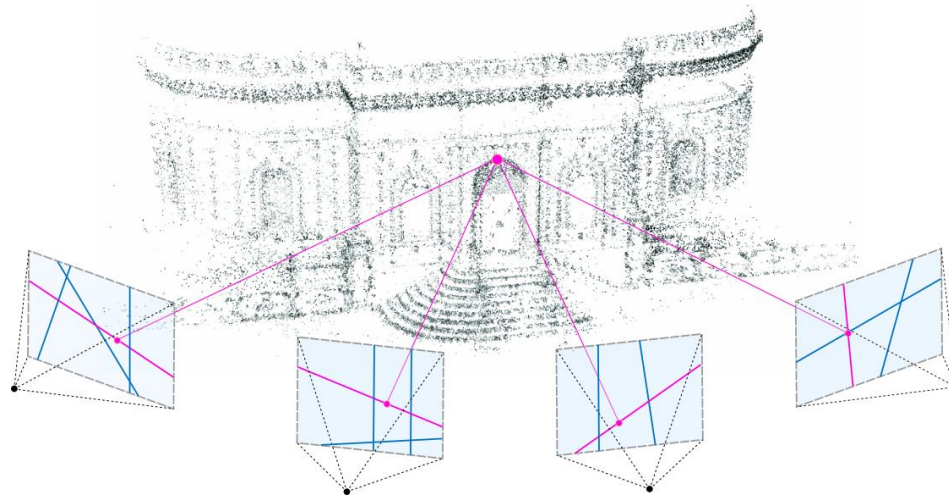
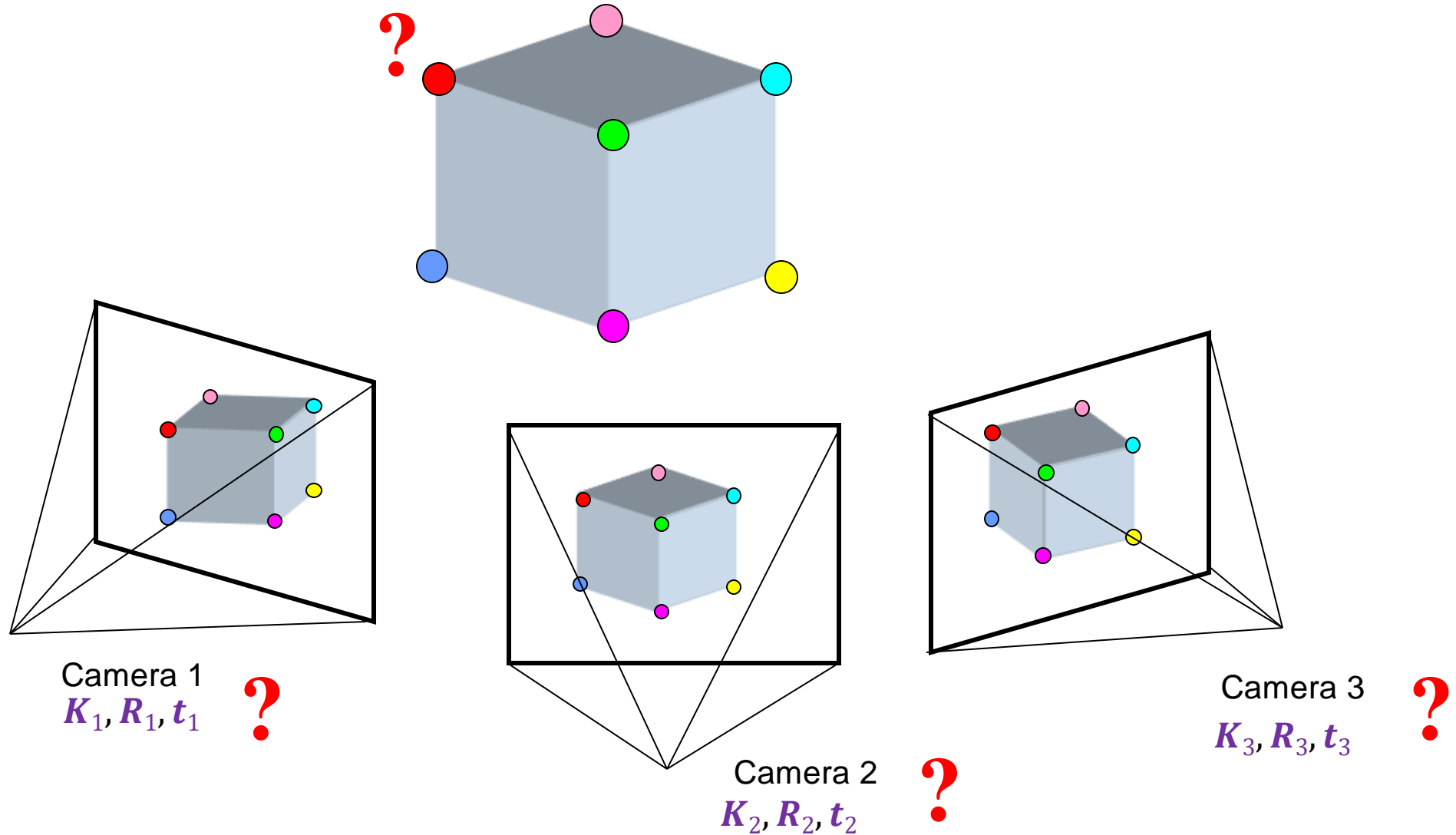
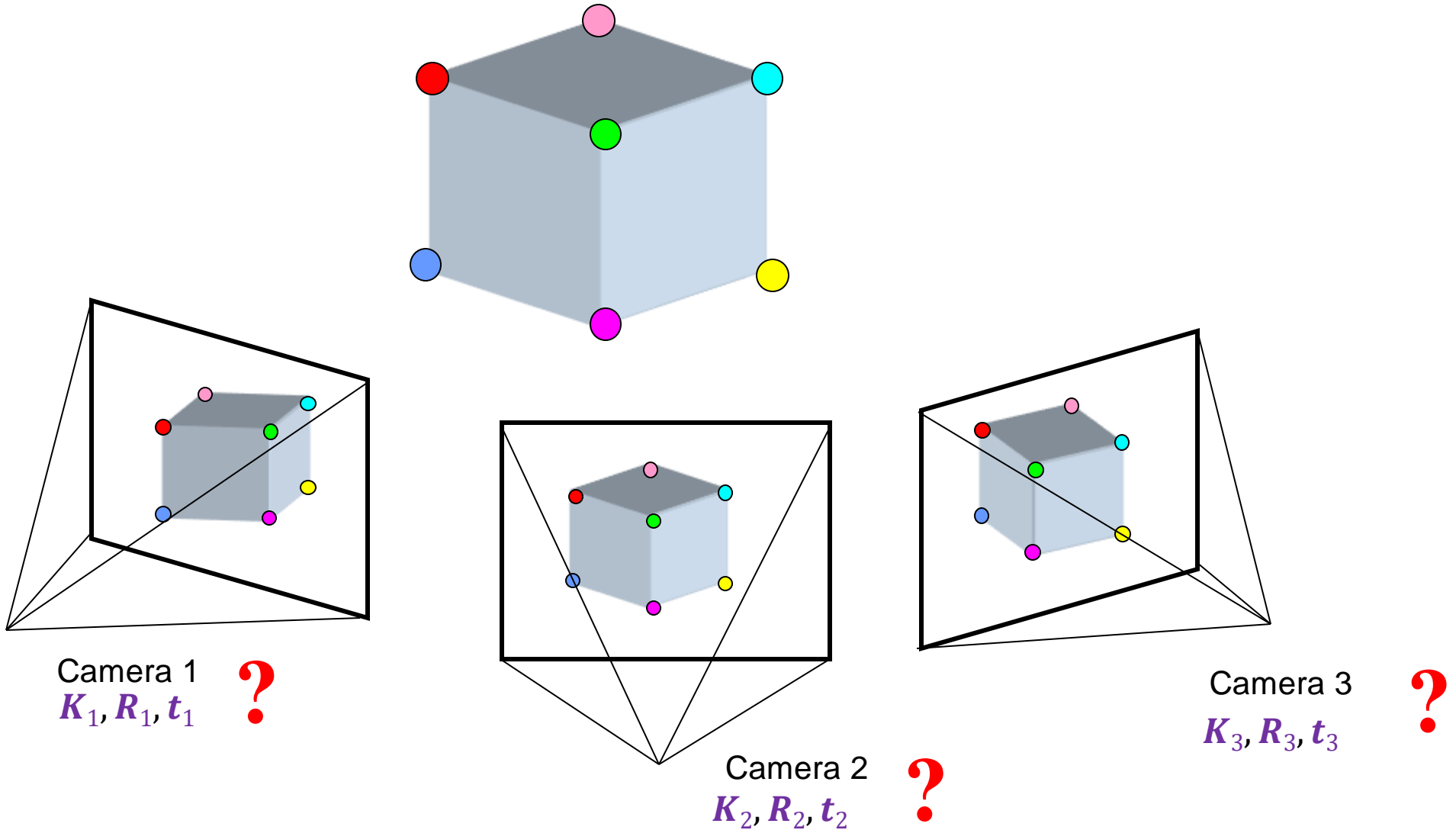


Figure from [blog](#)

Structure from motion

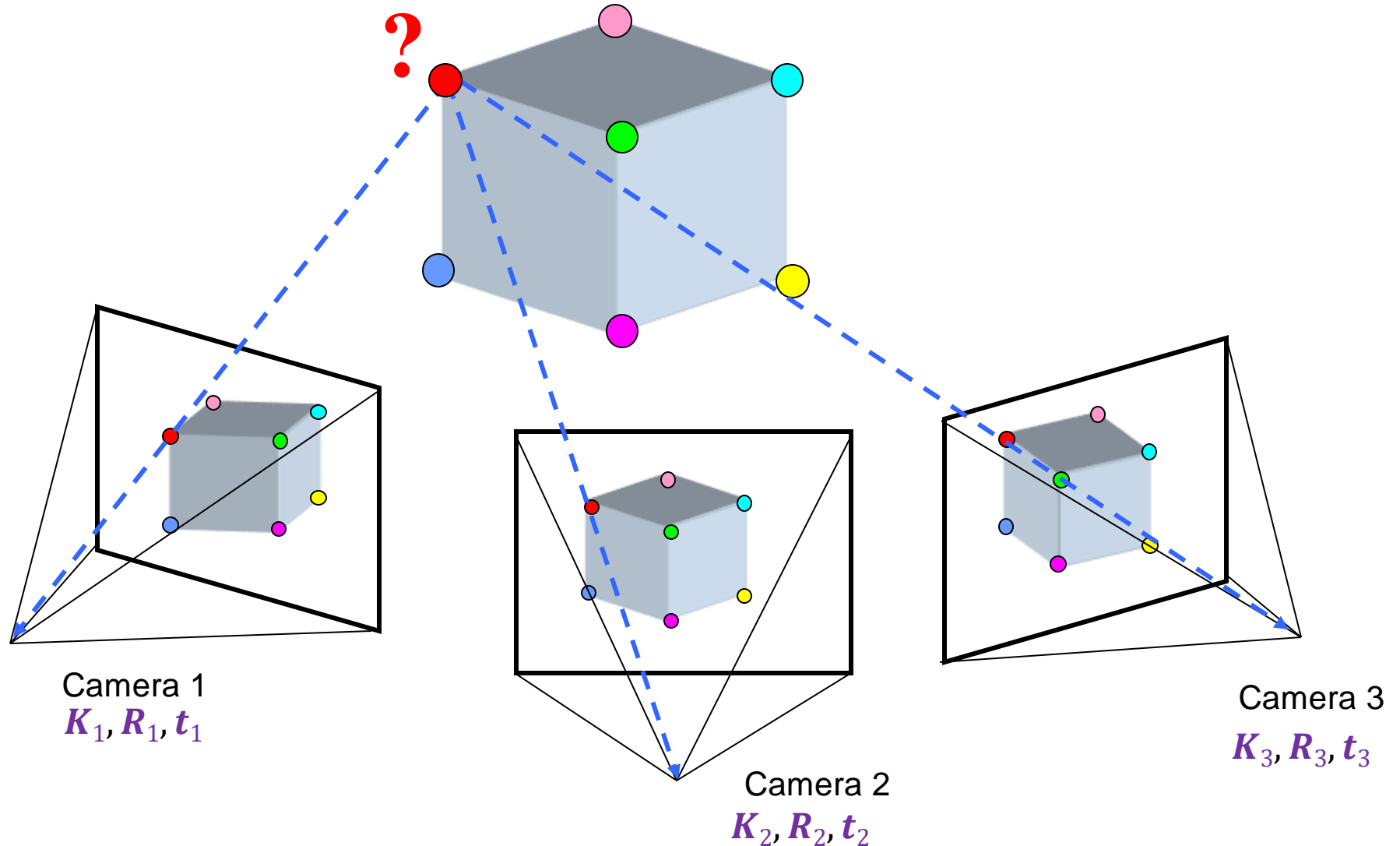


Recall: Calibration



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

Recall: Triangulation



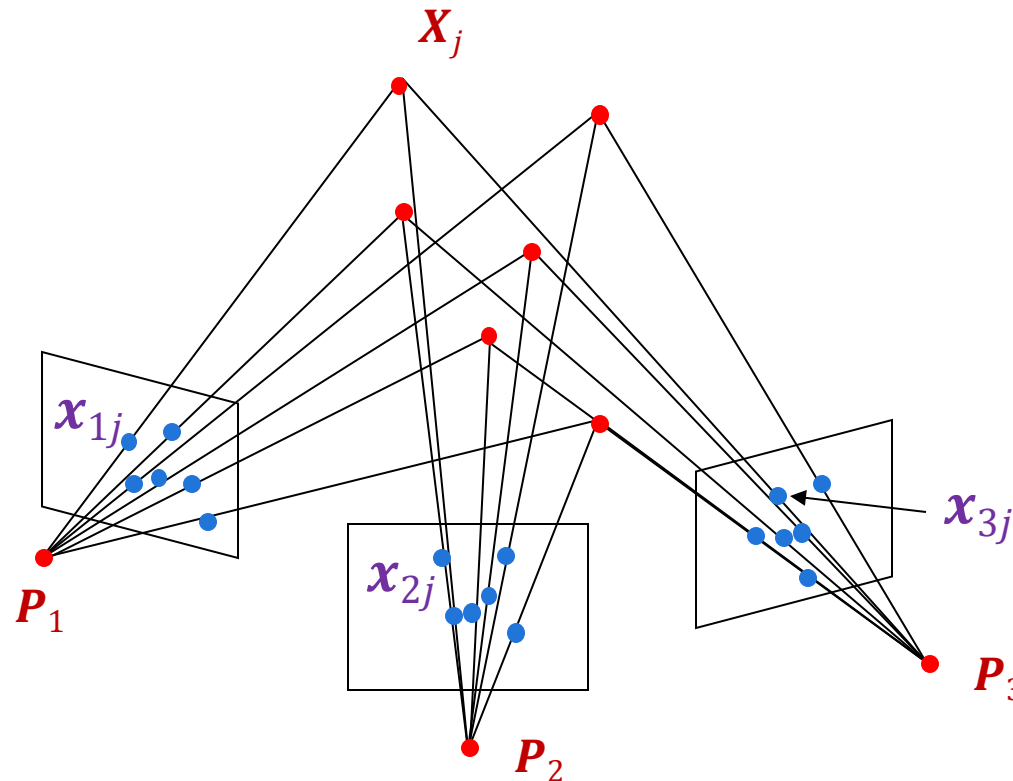
Given *known cameras* and projections of the same 3D point in two or more images, compute the 3D coordinates of that point

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 \mathbf{P}_i and n 3D points \mathbf{X}_j from the mn correspondences \mathbf{x}_{ij}



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:

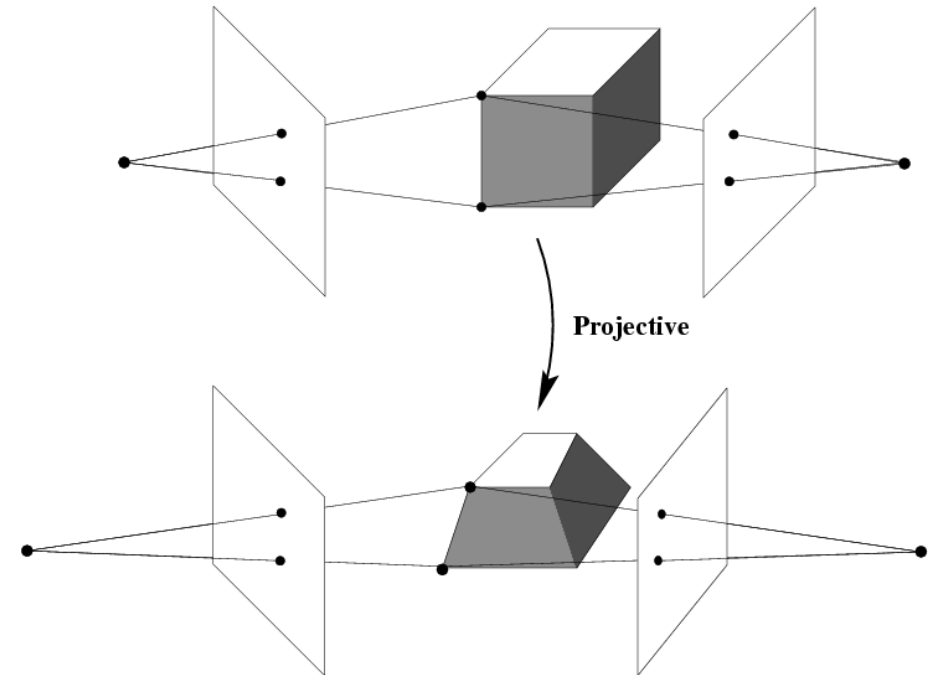
$$x \cong PX = (PQ^{-1})(QX)$$

Projective Ambiguity

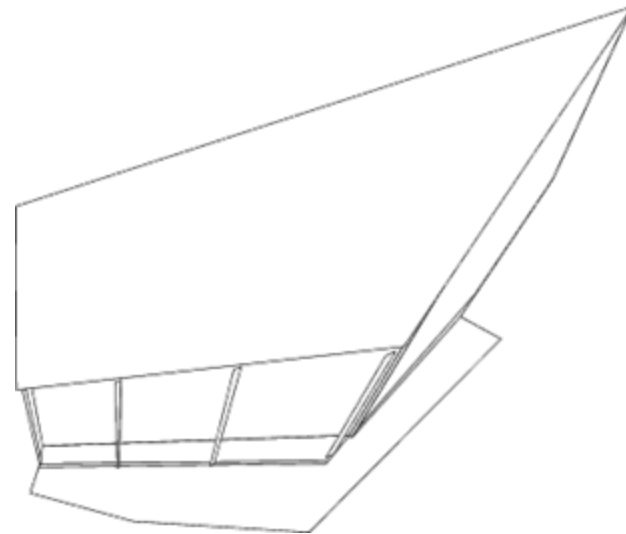
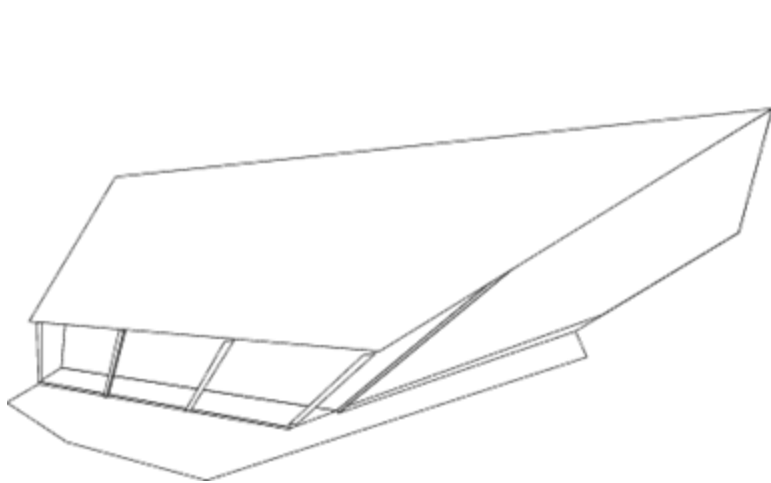
- With no constraints on the camera calibration matrices or on the scene, we can reconstruct up to a *projective* ambiguity:

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

Q is a general full-rank 4×4 matrix



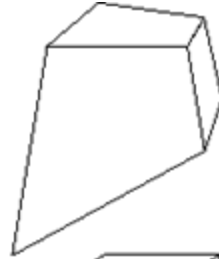
Projective Ambiguity



Types of Transformation

Projective
15dof

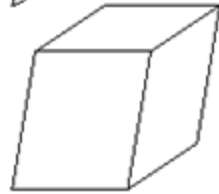
$$\begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{v}^\top & v \end{bmatrix}$$



Preserves intersection and tangency

Affine
12dof

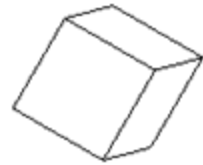
$$\begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix}$$



Preserves parallelism, volume ratios

Similarity
7dof

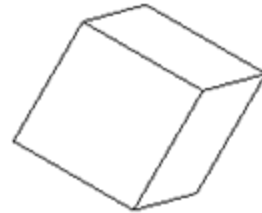
$$\begin{bmatrix} s\mathbf{R} & \mathbf{t} \\ \mathbf{0}^\top & 1 \end{bmatrix}$$



Preserves angles, ratios of length

Euclidean
6dof

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



Preserves angles, lengths

Projective Structure from Motion

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

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

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

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

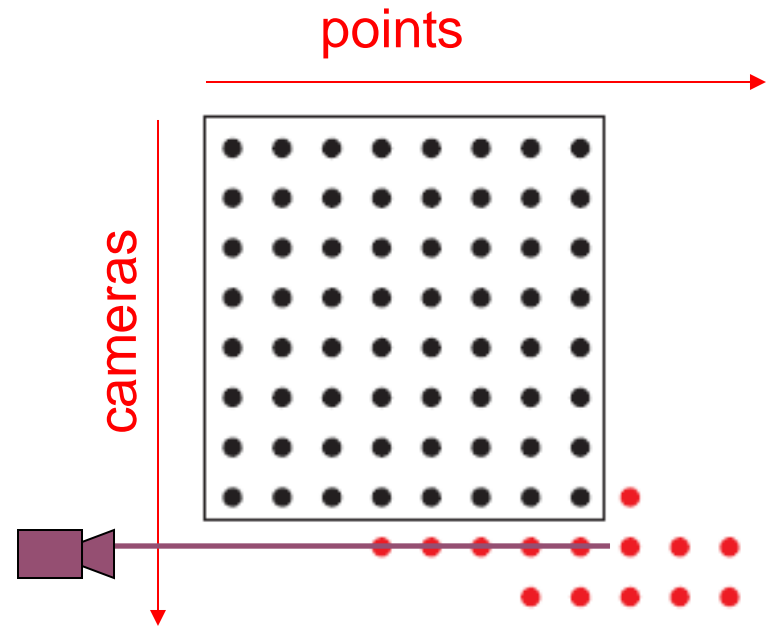
- We can solve for structure and motion when $2mn \geq 11m + 3n - 15$
- For two cameras, at least 7 points are needed

Projective SFM: Two-Camera Case

1. Estimate fundamental matrix F between the two views
 2. Set first camera matrix to $[I \mid 0]$
 3. Then the second camera matrix is given by $[A \mid t]$ where t is the epipole ($F^T t = 0$) and $A = -[t]_{\times} F$
- In practice, SFM pipelines use the guesses of intrinsic parameters and the [five-point algorithm](#) (for essential matrix)
 - Recall 7-point or 8-point algorithm for fundamental matrix

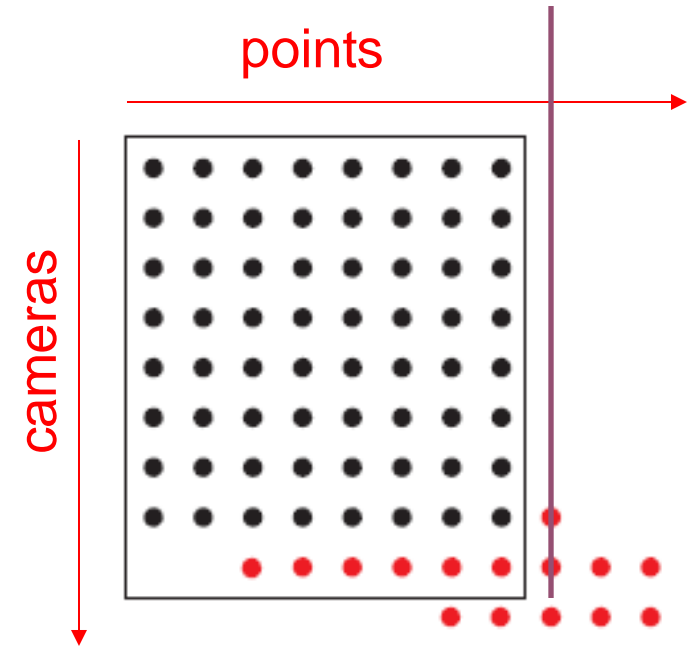
Incremental Structure from Motion

- 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**



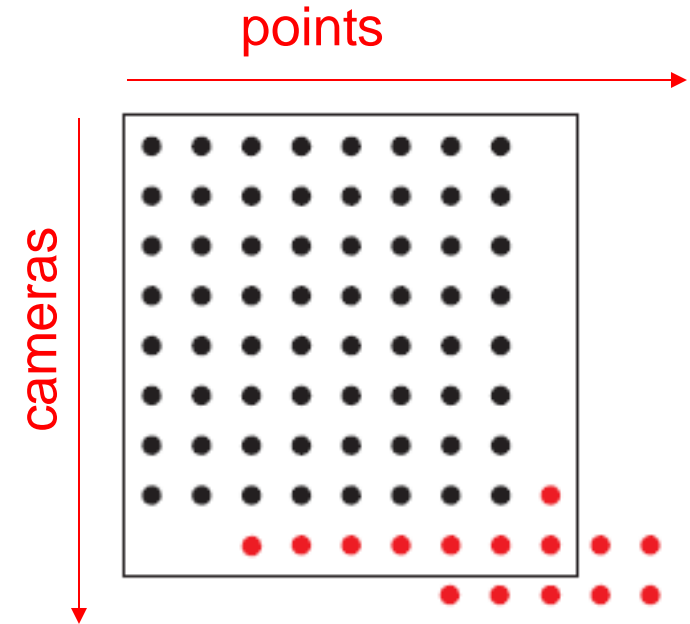
Incremental Structure from Motion

- 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

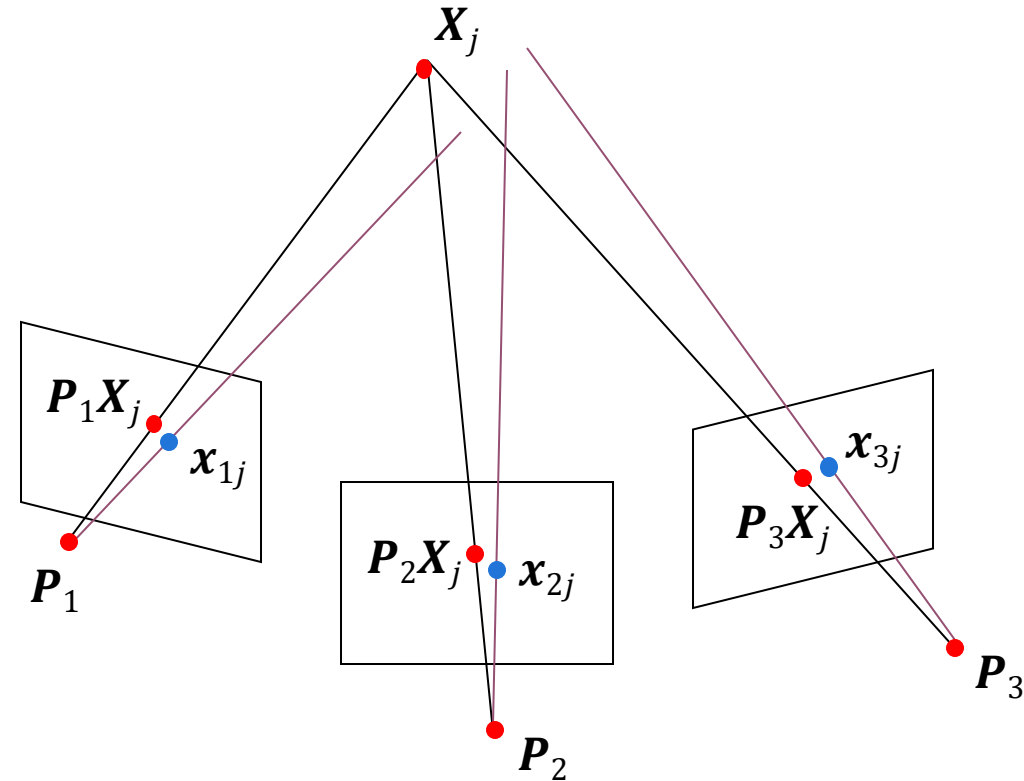


Bundle Adjustment

- Non-linear method for refining structure and motion
- Minimize reprojection error (with lots of bells and whistles):

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

visibility flag: is point j visible in view i ?



Outline

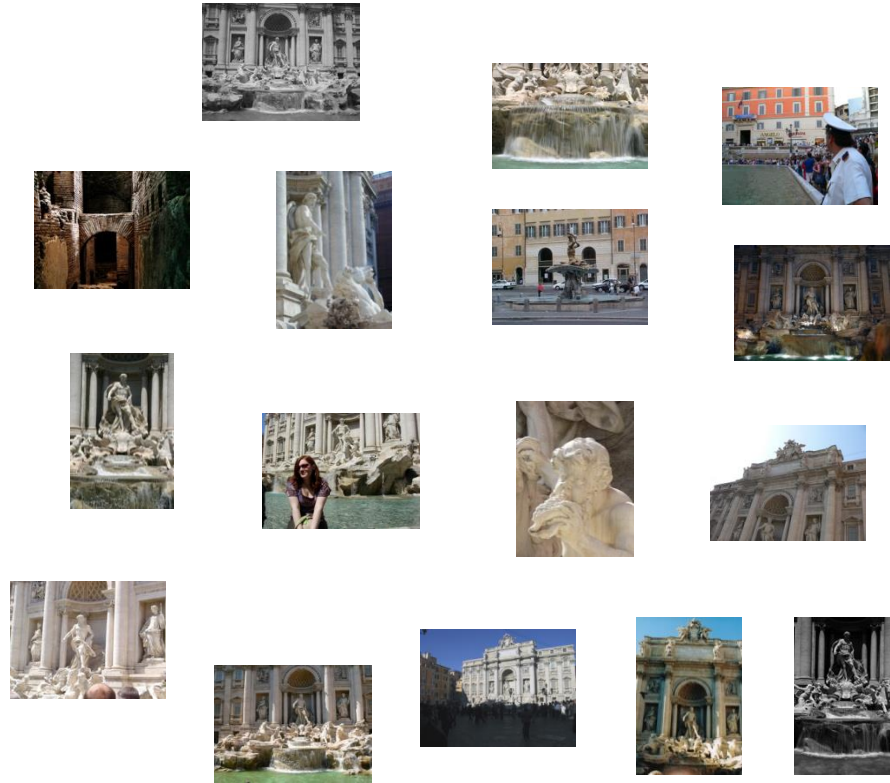
- Problem Formulation
- Projective structure from motion
- SfM pipeline

Representative SFM Pipeline



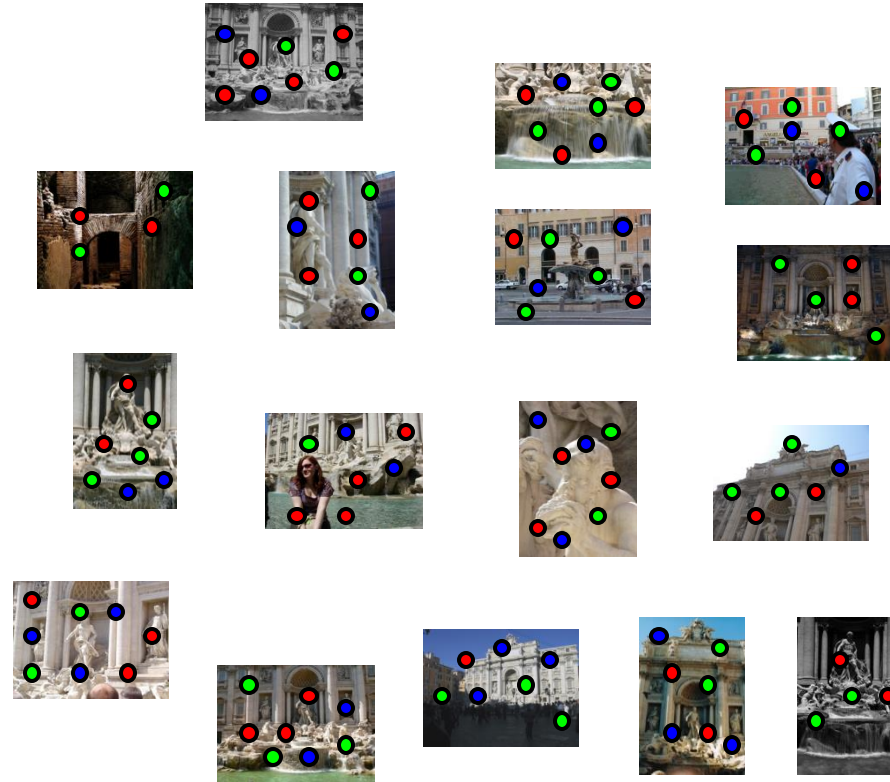
N. Snavely, S. Seitz, and R. Szeliski. [Photo tourism: Exploring photo collections in 3D](http://phototour.cs.washington.edu/). SIGGRAPH 2006
<http://phototour.cs.washington.edu/>

Feature Detection



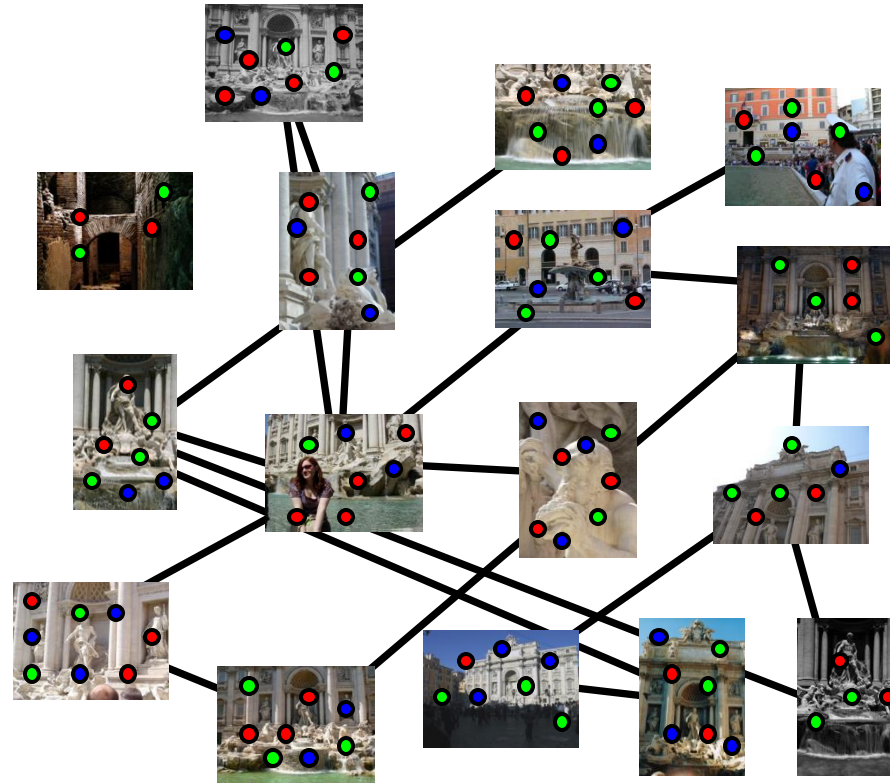
Detect SIFT features

Feature Detection



Other popular feature types: [SURF](#), [ORB](#), [BRISK](#), ...

Feature Matching



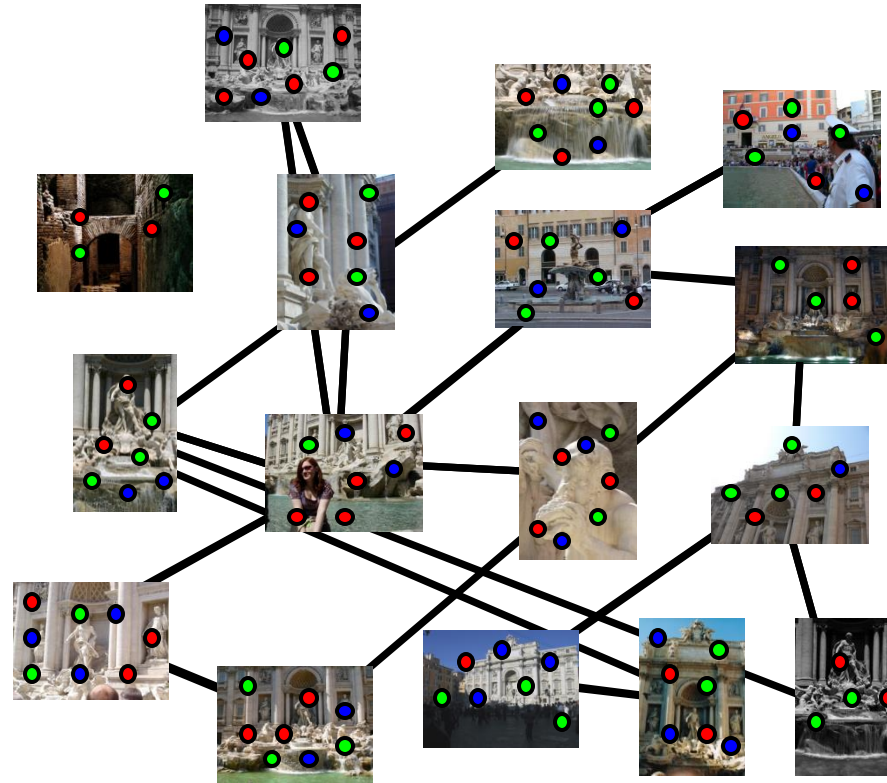
Match features between each pair of images

Feature Matching



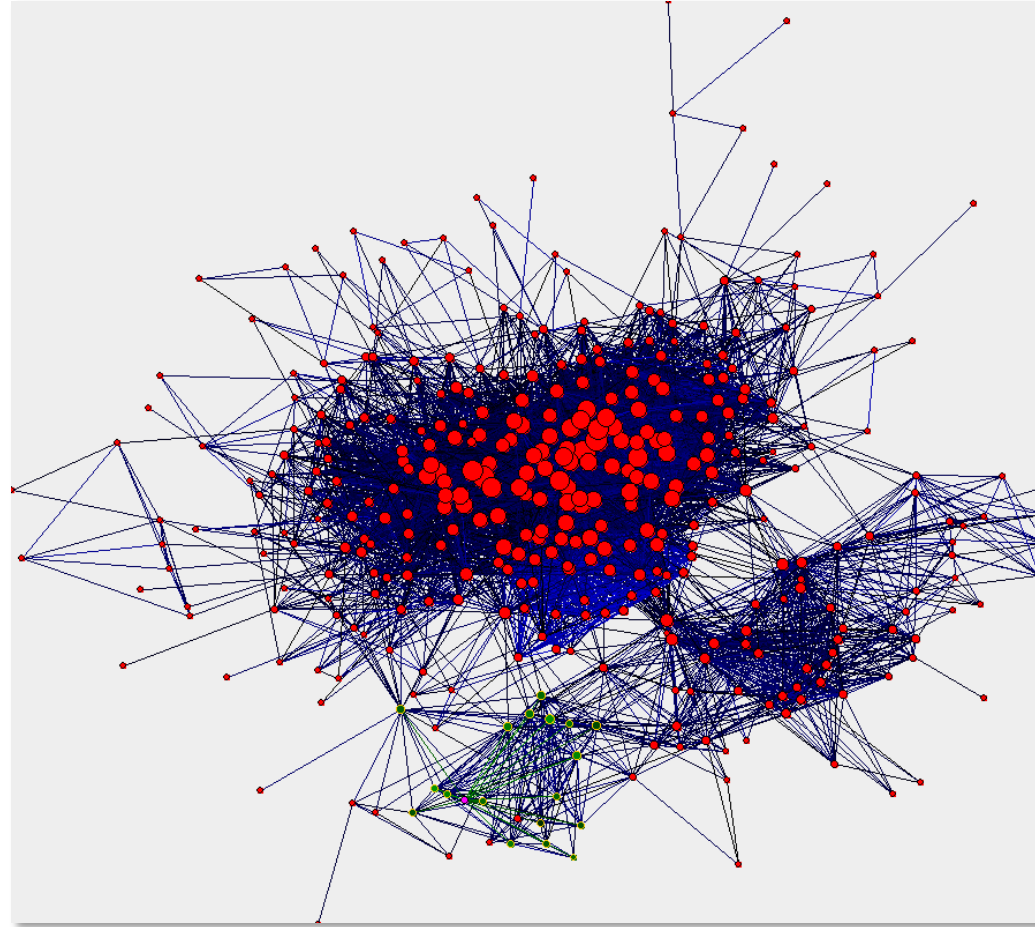
Use RANSAC to estimate fundamental matrix between each pair

Feature Matching



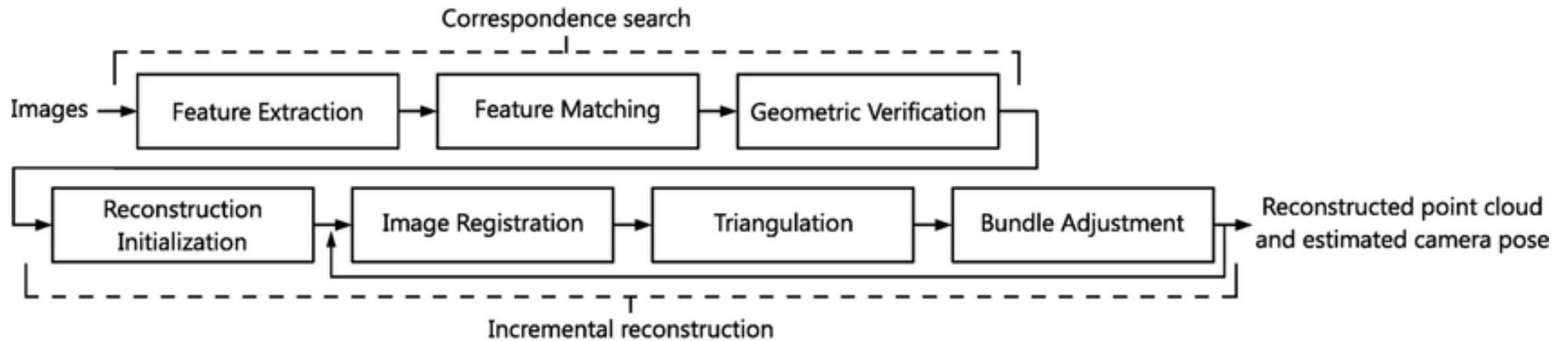
Use RANSAC to estimate fundamental matrix between each pair

Image Connectivity Graph



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

Structure from Motion (SfM) Pipeline



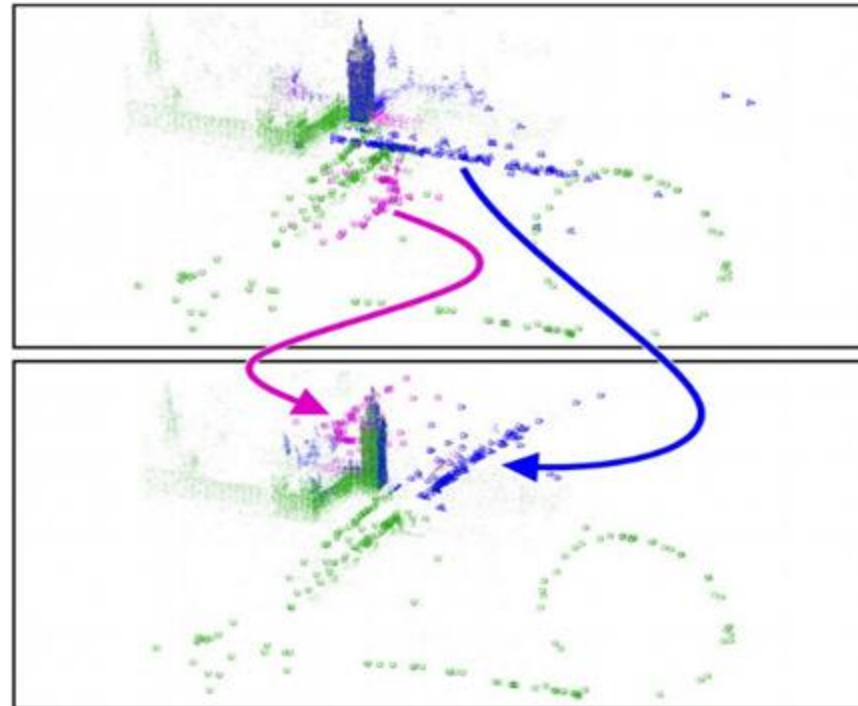
Incremental SFM

- 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 [five-point algorithm](#)
 - Use triangulation to initialize model points
- While remaining images exist
 - Find an image with many feature matches with images in the model
 - Run RANSAC on feature matches to register new image to model
 - Triangulate new points
 - Perform bundle adjustment to re-optimize everything
 - Optionally, align with GPS from EXIF data or ground control points

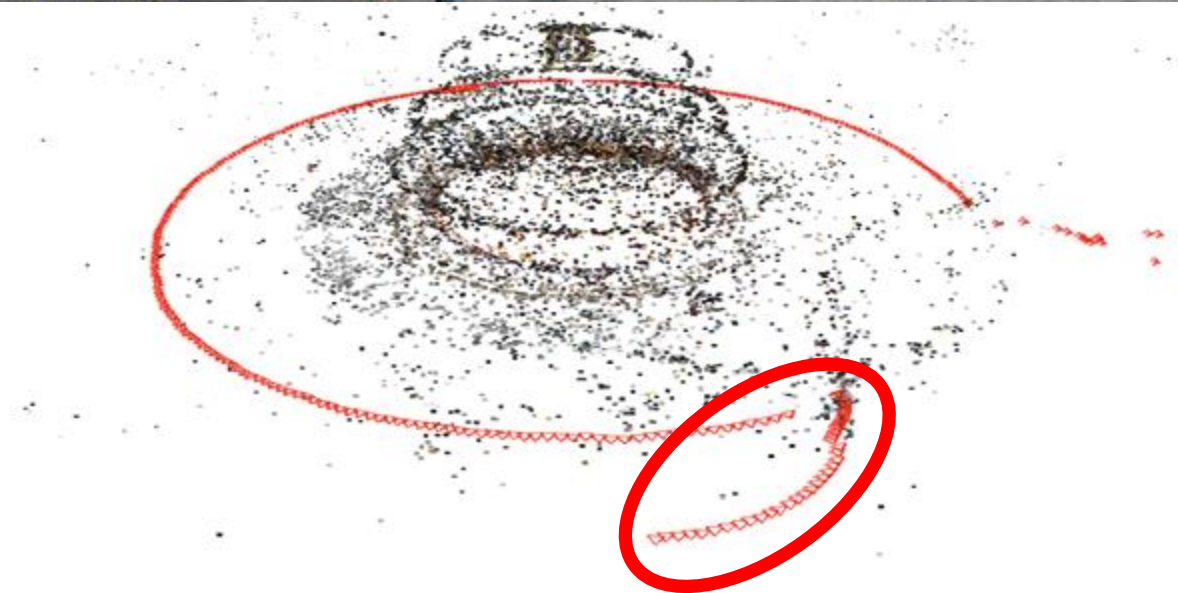
The devil is in the details

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

Repetitive structures cause catastrophic failures



Repetitive structures cause catastrophic failures

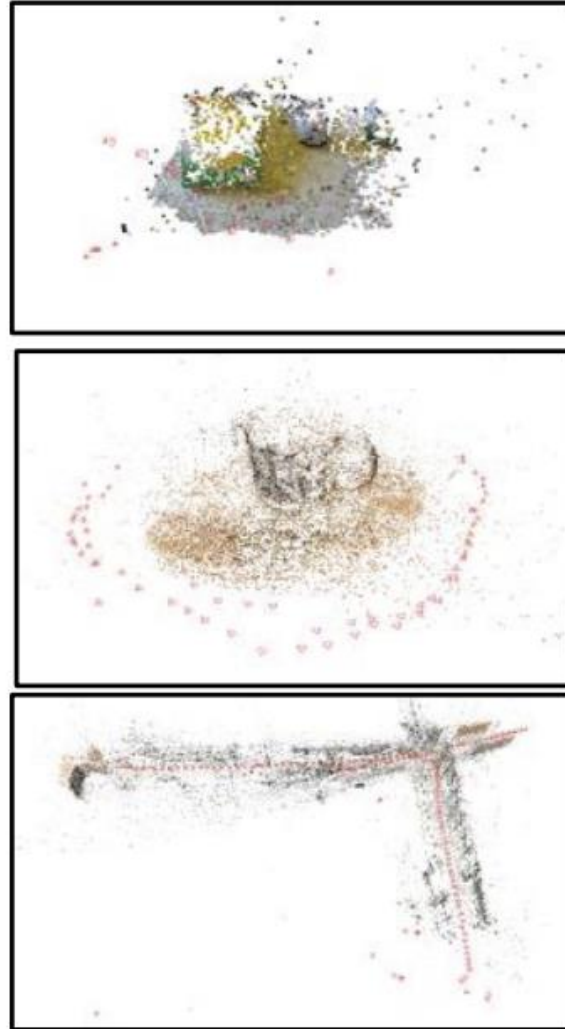


Repetitive structures cause catastrophic failures

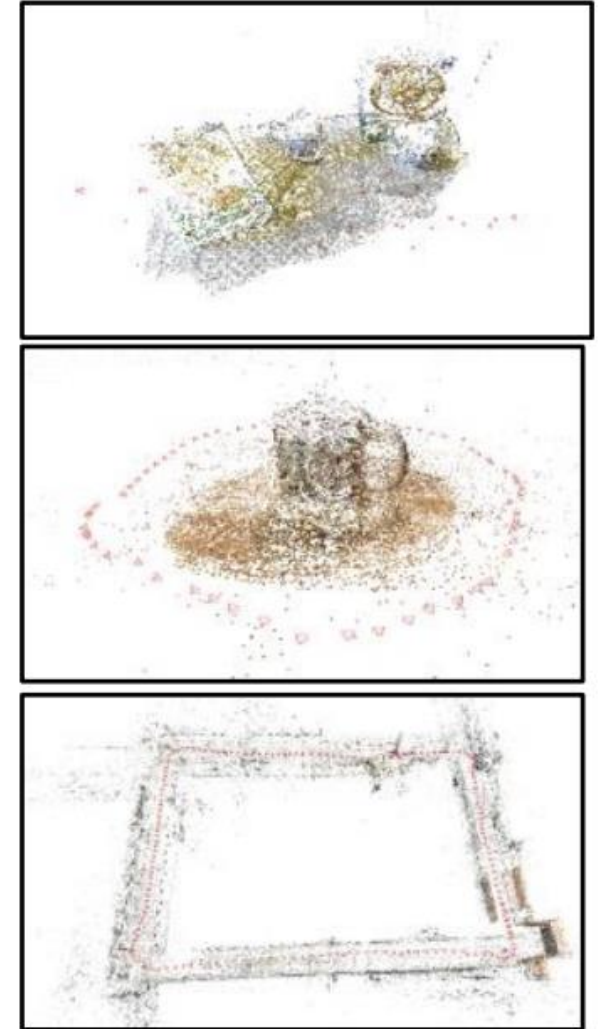
Erroneously Matching Images



Baseline Reconstruction



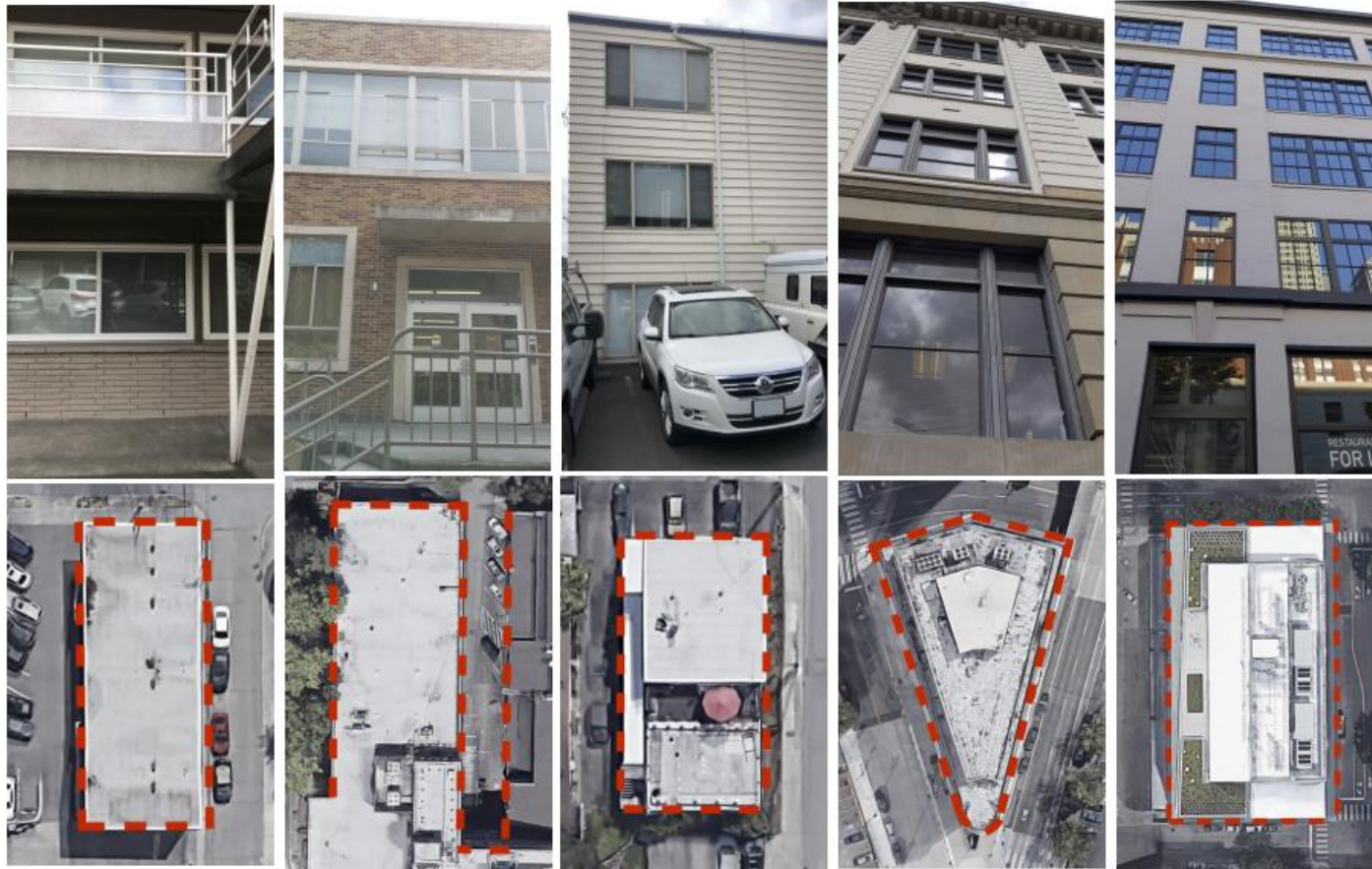
Our Reconstruction



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



seattle1

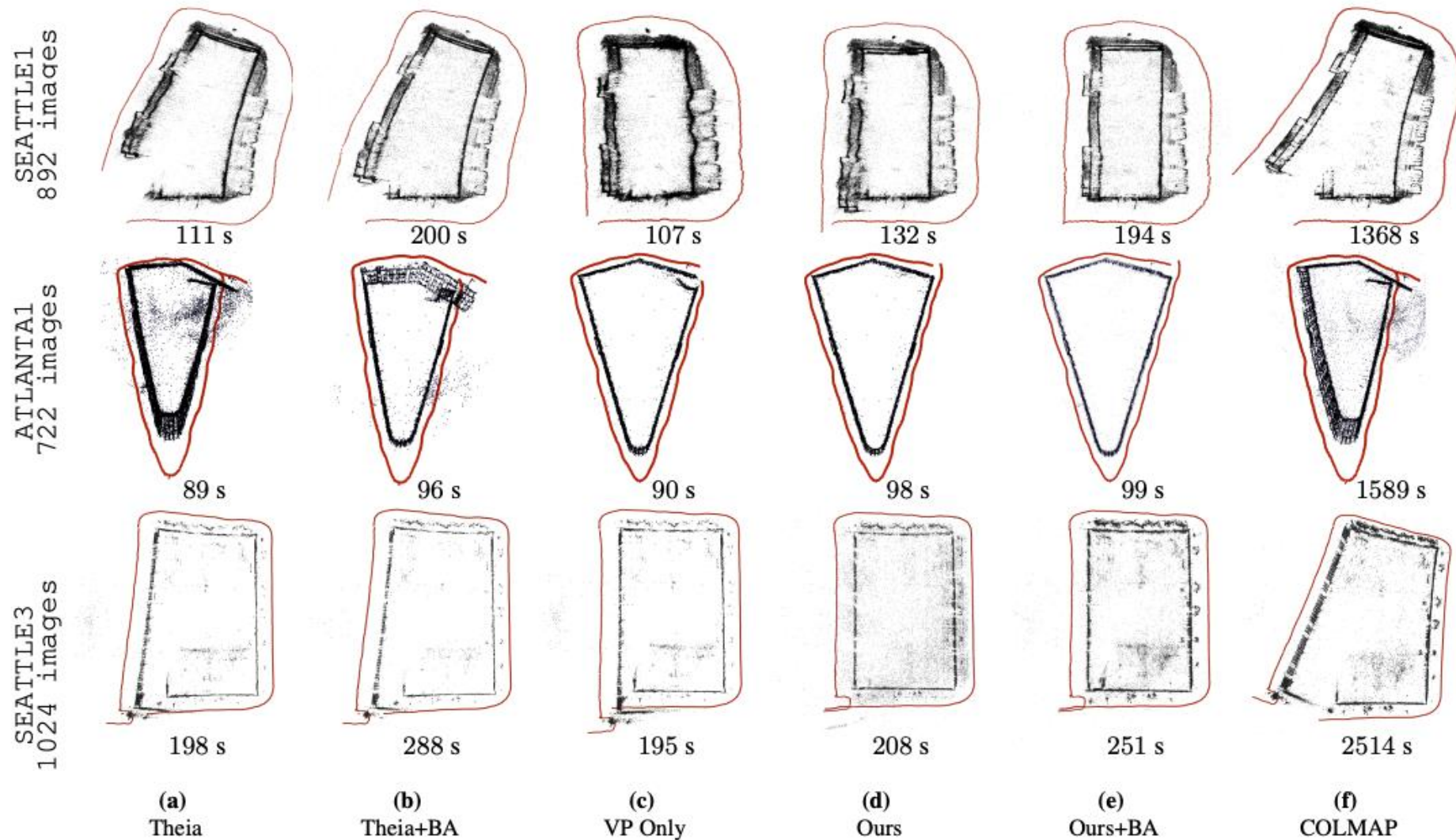
more_half

seattle2

atlanta1

seattle3

Reducing error accumulation and closing loops



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

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