

Advanced Computer Vision

Week 09

Nov. 1, 2022
Seokju Lee

Structure-from-Motion (SfM)

Contents of Structure-from-Motion

Multi-view 3D reconstruction

- Pinhole camera model
- Two-view geometry
- Epipolar geometry
- Essential matrix
- Fundamental matrix
- Bundle adjustment
- ...

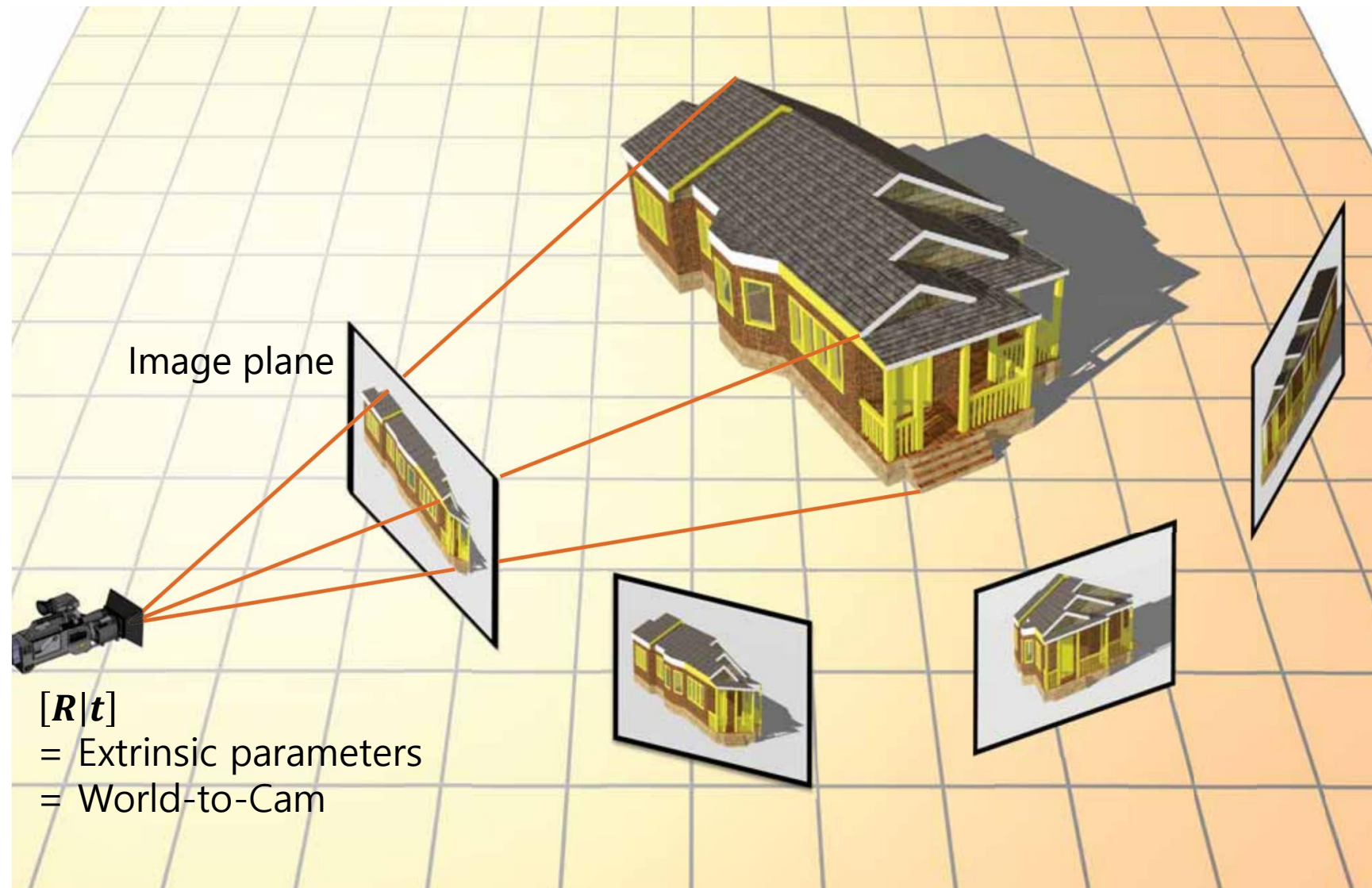
Pinhole Camera Model

$$x = K[R|t]X = PX$$

$$s \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & \text{skew_}cf_x & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$
$$= K[R|t] \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

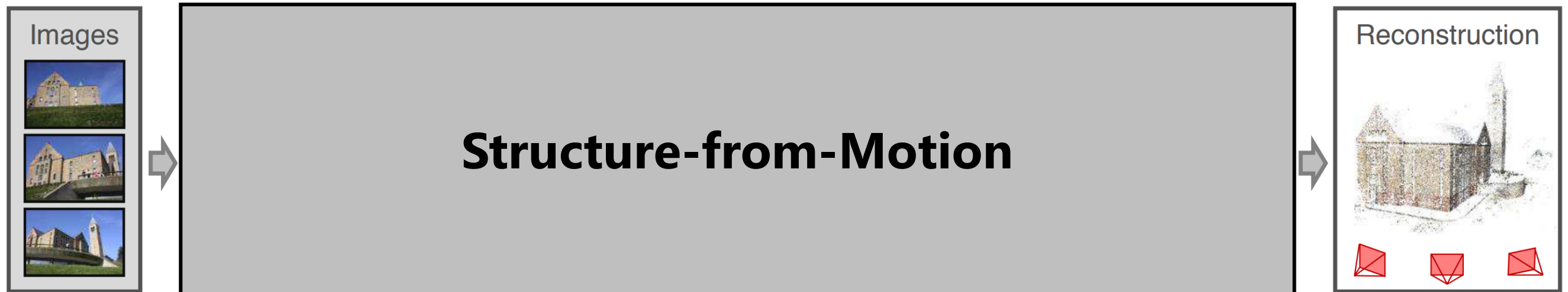
- (X, Y, Z) : 3D point in the world coordinate
- $[R|t]$: extrinsic parameters to convert the world coordinate into the camera coordinate
- K : intrinsic parameters to represent the camera characteristics
- $K[R|t]$: camera projection matrix
- (x, y) : 2D pixel location in the image plane
- s : scale factor

Multi-View 3D Reconstruction



Structure-from-Motion Pipeline

Given only the 2D multi-view images of a scene,
recover the underlying 3D **structure** and the camera **motion**.



Structure-from-Motion Pipeline

Q1. How to find **2D-3D points** to reconstruct?

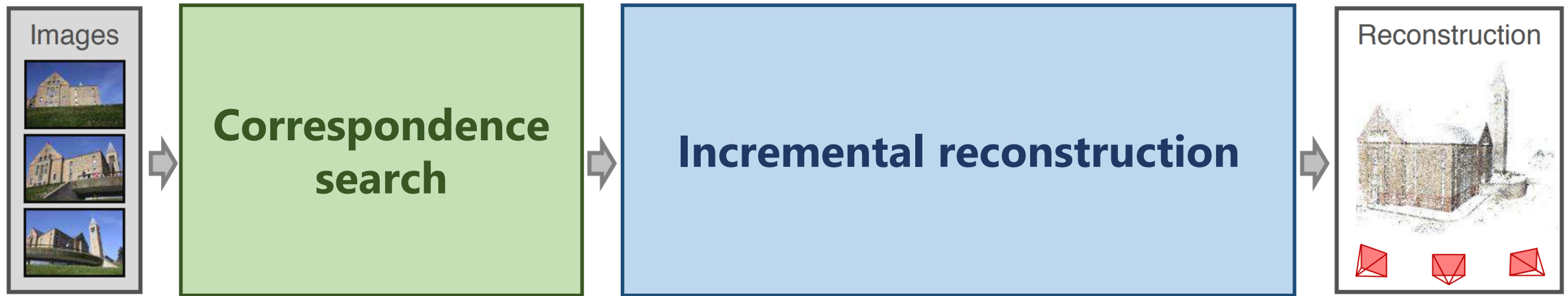
Q2. How to find an **optimal** 3D structure and camera poses for multiple view?



Structure-from-Motion Pipeline

Q1. How to find **2D-3D points** to reconstruct?

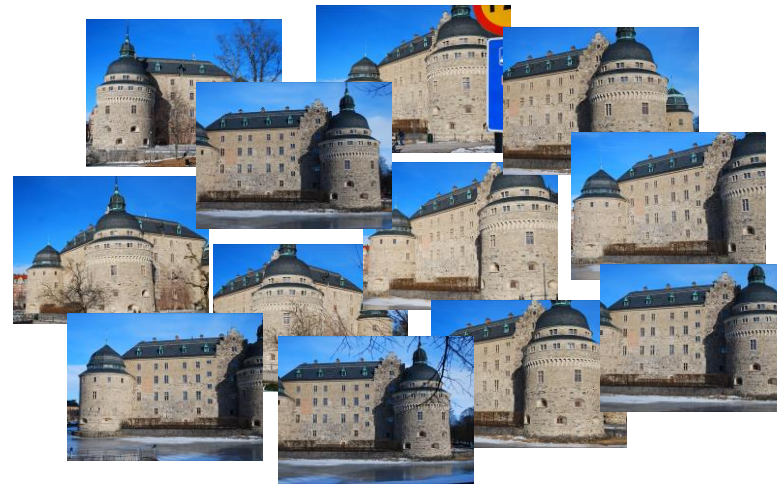
Q2. How to find an **optimal** 3D structure and camera poses for multiple view?



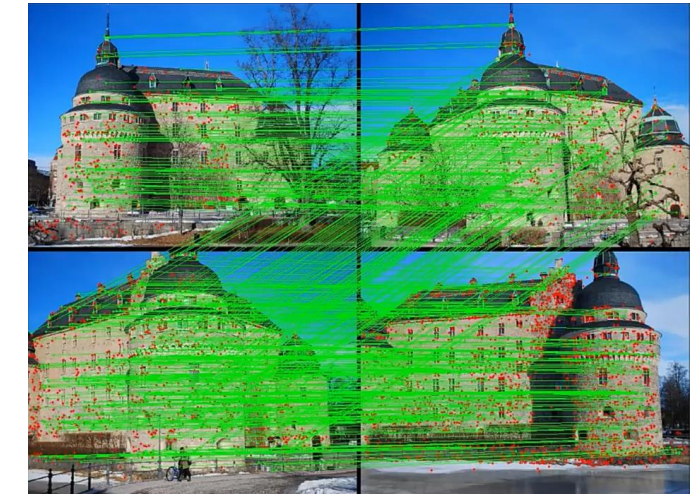
Structure-from-Motion (1): Correspondence Search

Correspondence search

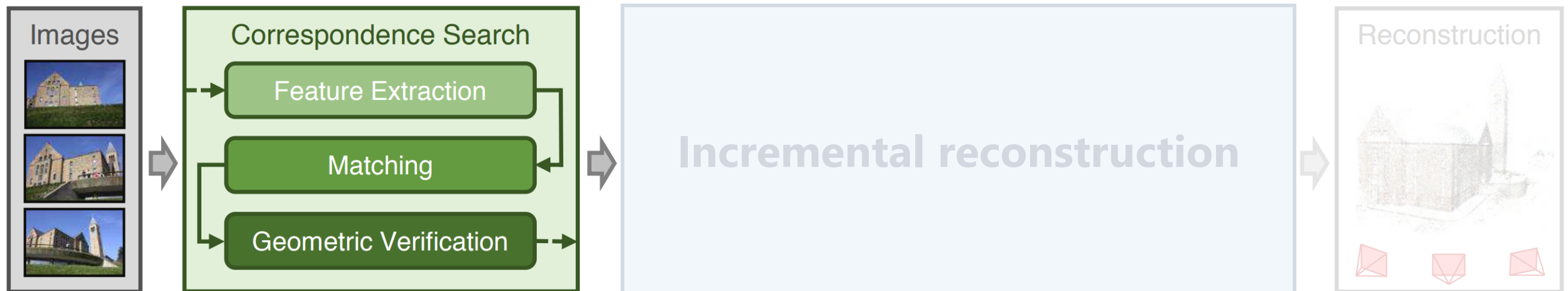
- Feature extraction
- Matching
- Geometric verification



A set of images



Feature extraction & matching (graph)



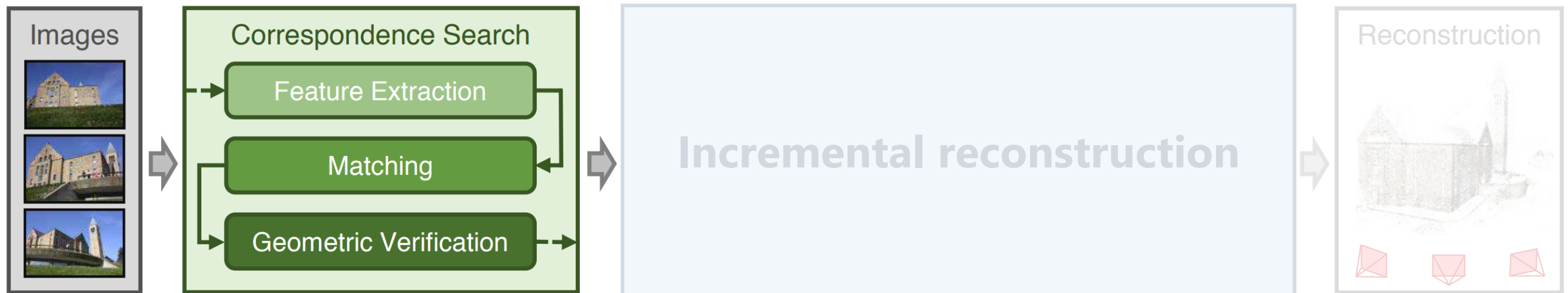
Structure-from-Motion (1): Correspondence Search

c.f.) Homography matrix

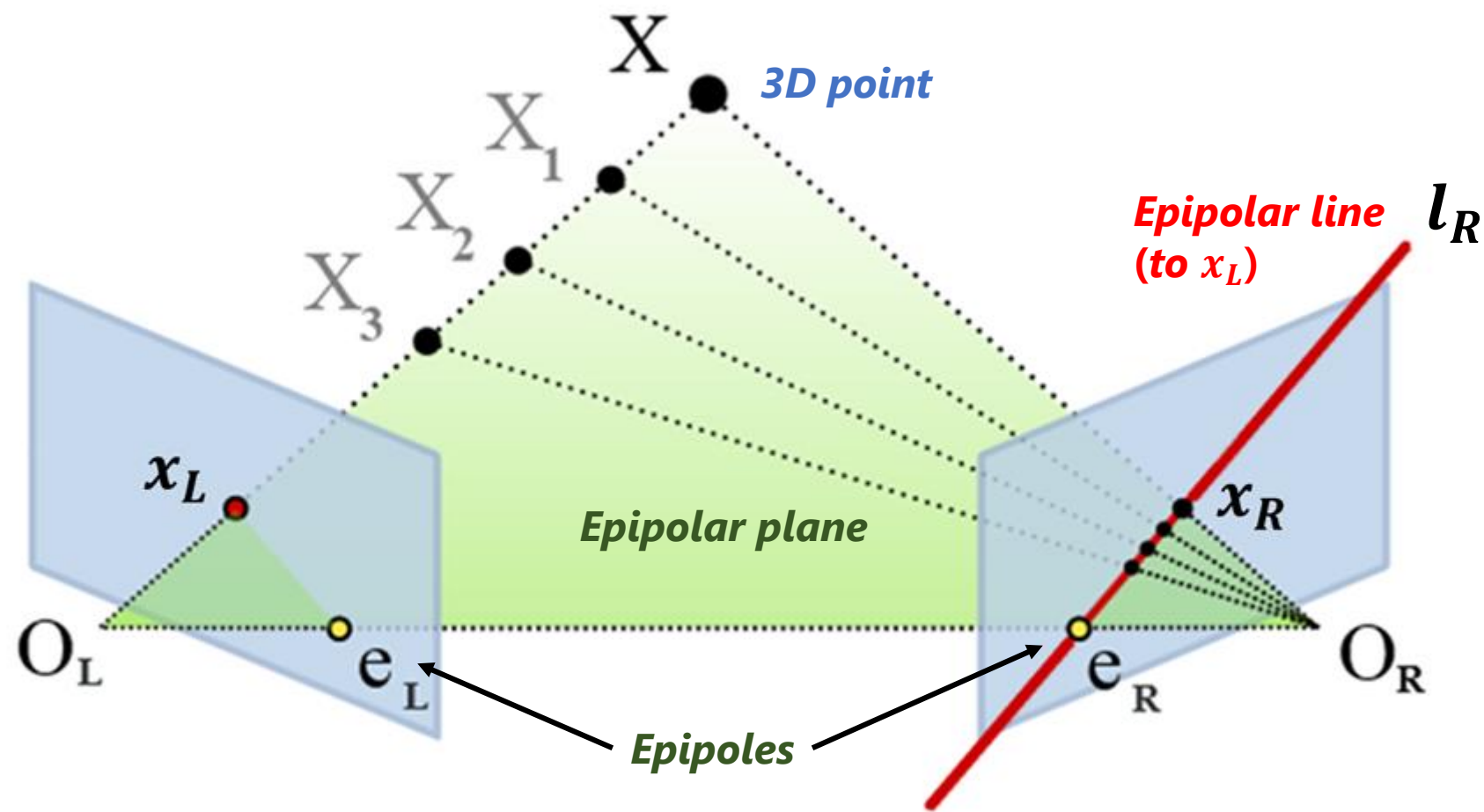
Correspondence search

- Feature extraction
- Matching
- **Geometric verification**

- Epipolar geometry describes the relation of moving cameras **(5 or 8-point algorithm)** through the essential matrix, $\mathbf{E} \in \mathbb{R}^{3 \times 3}$, or the fundamental matrix, $\mathbf{F} \in \mathbb{R}^{3 \times 3}$
- If estimated \mathbf{E} projects a sufficient number of features between the images, it is verified! **(RANSAC)**



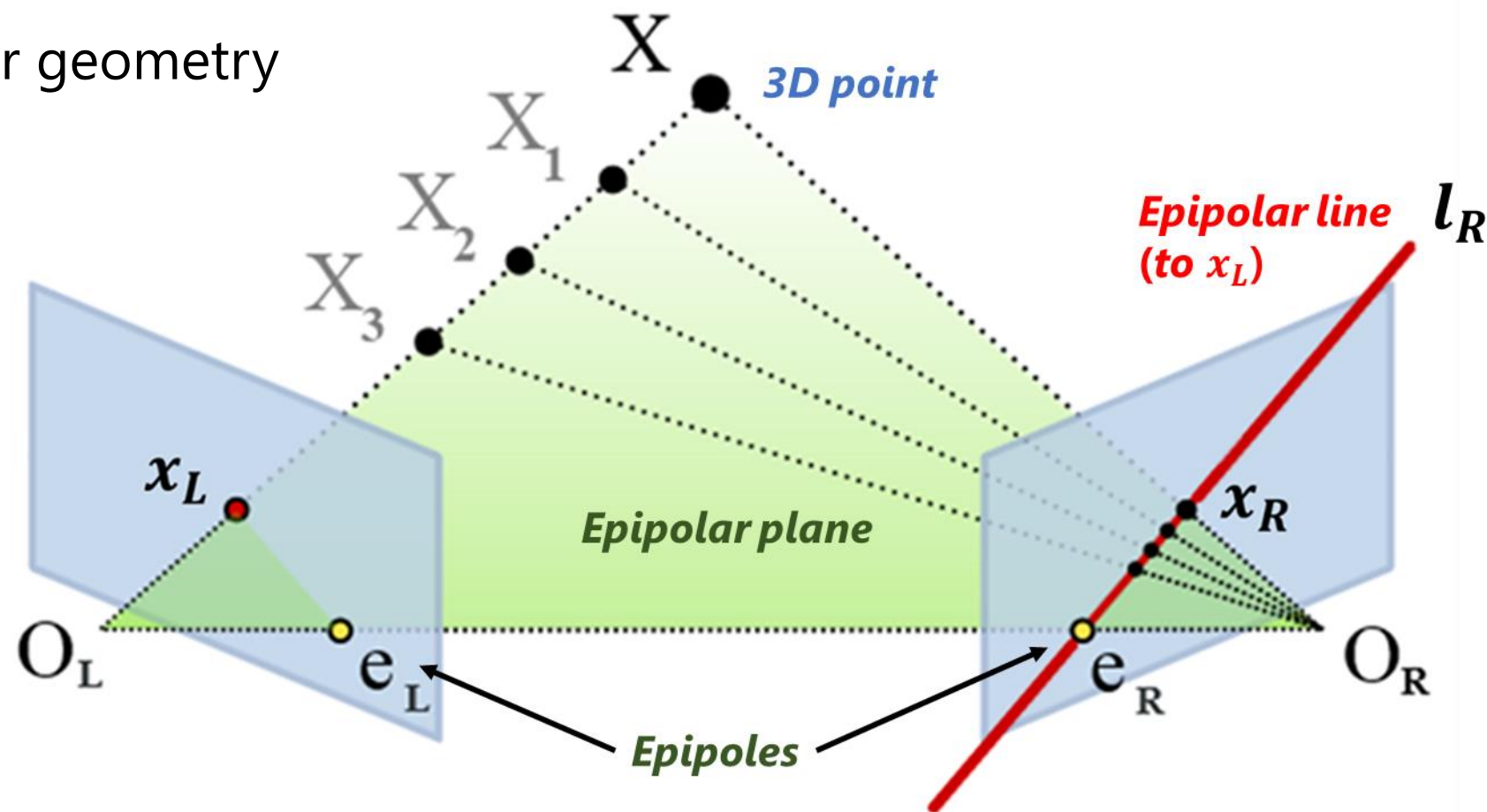
Epipolar Geometry



Potential matches for x_L are on the epipolar line l_R

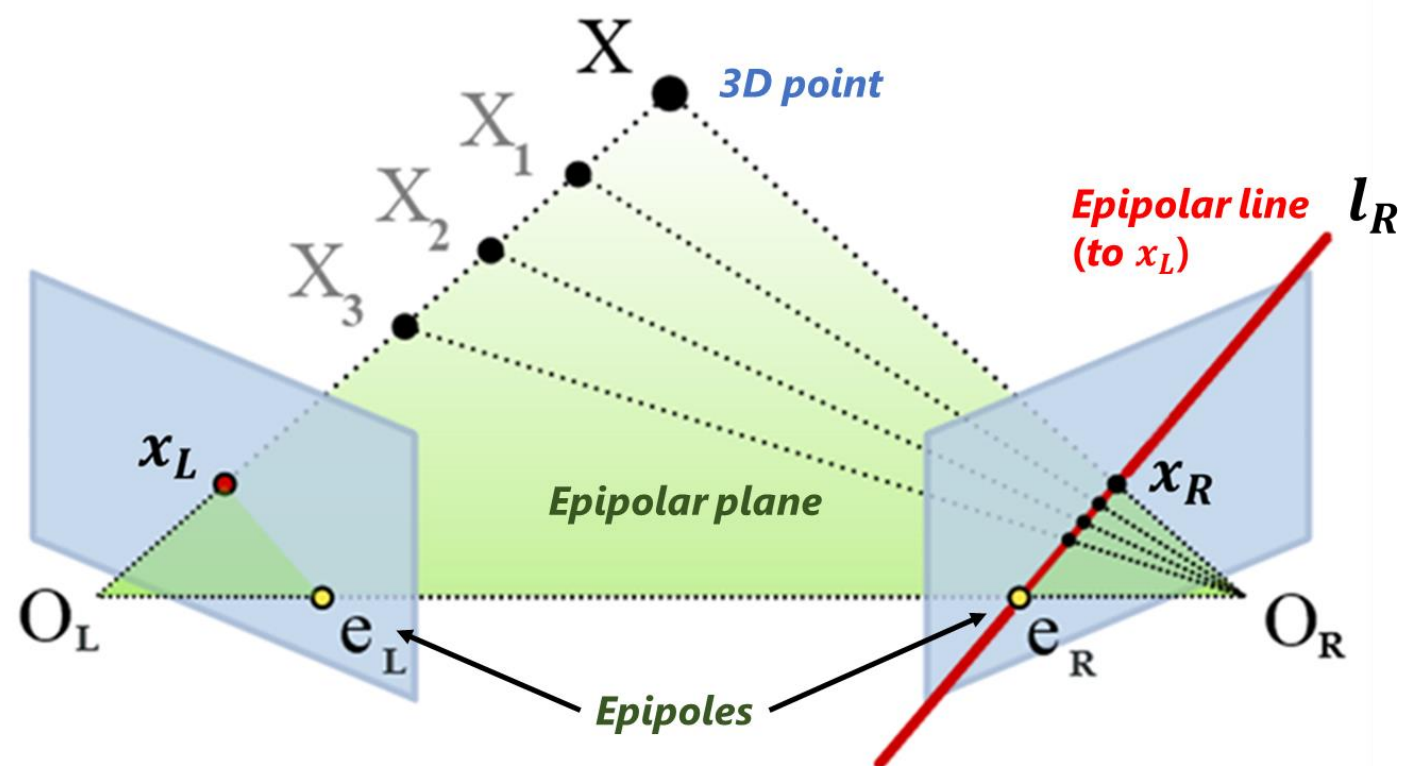
Essential Matrix

: Encodes epipolar geometry



Given a point x_L in one image, multiplying by the **essential matrix** $\mathbf{E} \in \mathbb{R}^{3 \times 3}$ will make the epipolar line in the right view: $\mathbf{E}x_L = l_R$

Epipolar Constraint



For a **linear** equation,

$$ax + by + c = 0 \text{ in vector form } l = \begin{bmatrix} a \\ b \\ c \end{bmatrix}.$$

If the point x_R is on the epipolar line l_R ,

$$x_R^\top l_R = 0 \quad \text{or, } l_R^\top x_R = 0$$

Since $E x_L = l_R$,

$$x_R^\top E x_L = 0$$

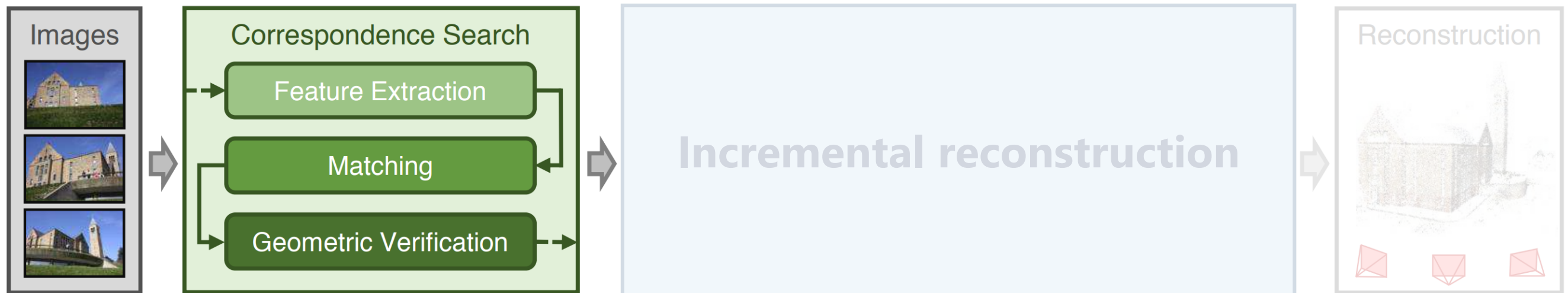
Structure-from-Motion (1): Correspondence Search

Correspondence search: Geometric verification

- **5-point / 8-point** algorithm:

Estimating the **essential matrix** from the feature point correspondences.

→ Direct linear transform: $\mathbf{x}_R^\top \mathbf{E} \mathbf{x}_L = 0 \rightarrow \mathbf{A} \mathbf{p} = 0 \rightarrow$ Use SVD to solve!



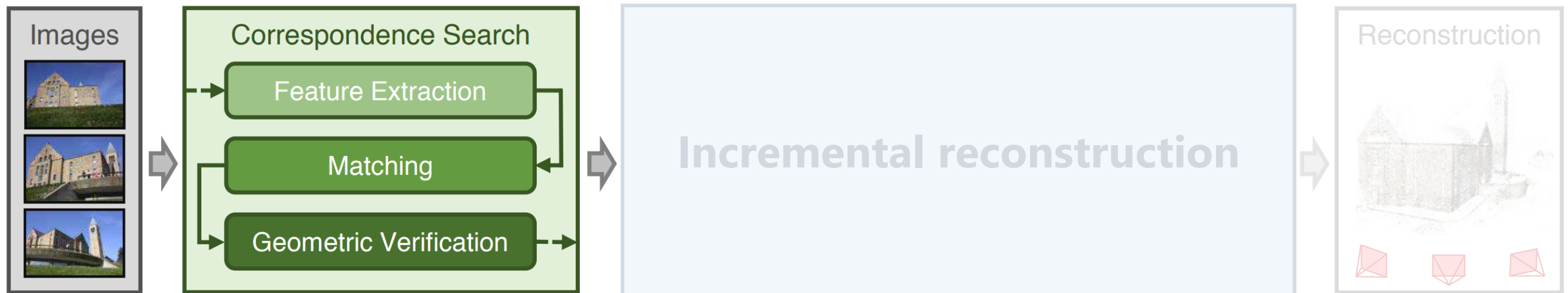
Structure-from-Motion (1): Correspondence Search

Correspondence search: Geometric verification

- **5-point / 8-point** algorithm: $\mathbf{x}_R^\top \mathbf{E} \mathbf{x}_L = 0 \rightarrow \mathbf{A} \mathbf{p} = 0$

Sample at least **5** or **8** points and compute the **essential matrix**.

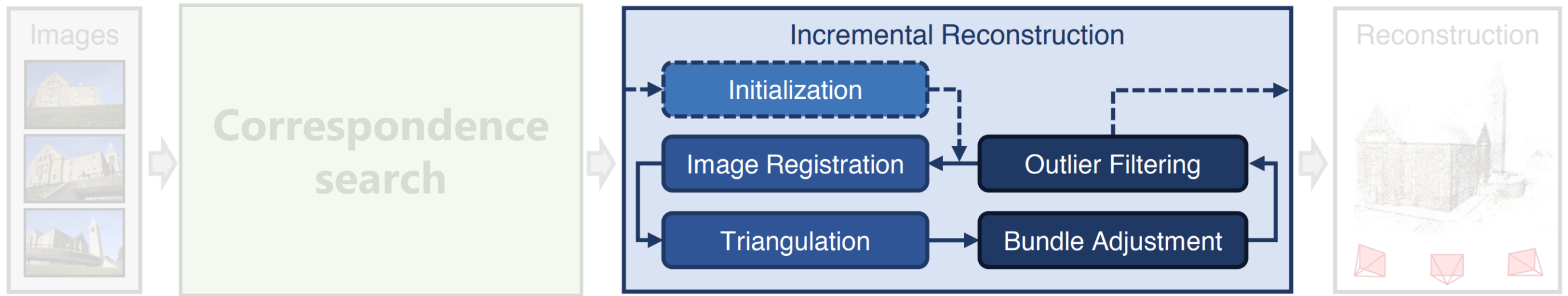
→ **RANSAC** to discriminate **inliers/outliers** and the **best** essential matrix!



Structure-from-Motion (2): Incremental Reconstruction

Incremental reconstruction

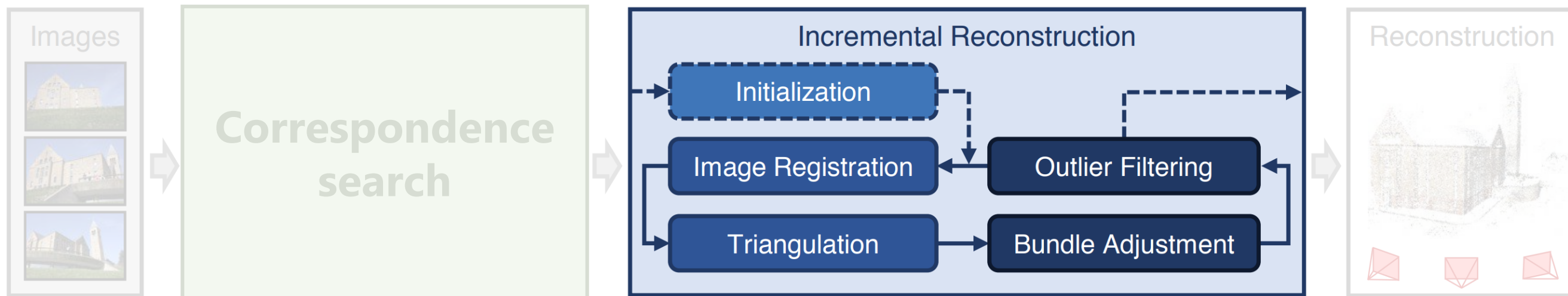
- Camera initialization
- Triangulation
- Bundle adjustment (refinement)



Structure-from-Motion (2): Incremental Reconstruction

Starting from two views, aggregate more views to refine the estimation

- Camera initialization → Estimate camera **projection** matrices from $\mathbf{E} \in \mathbb{R}^{3 \times 3}$
- Triangulation → **Lift** 2D points to 3D spaces using camera projection matrices
- Bundle adjustment (refinement)

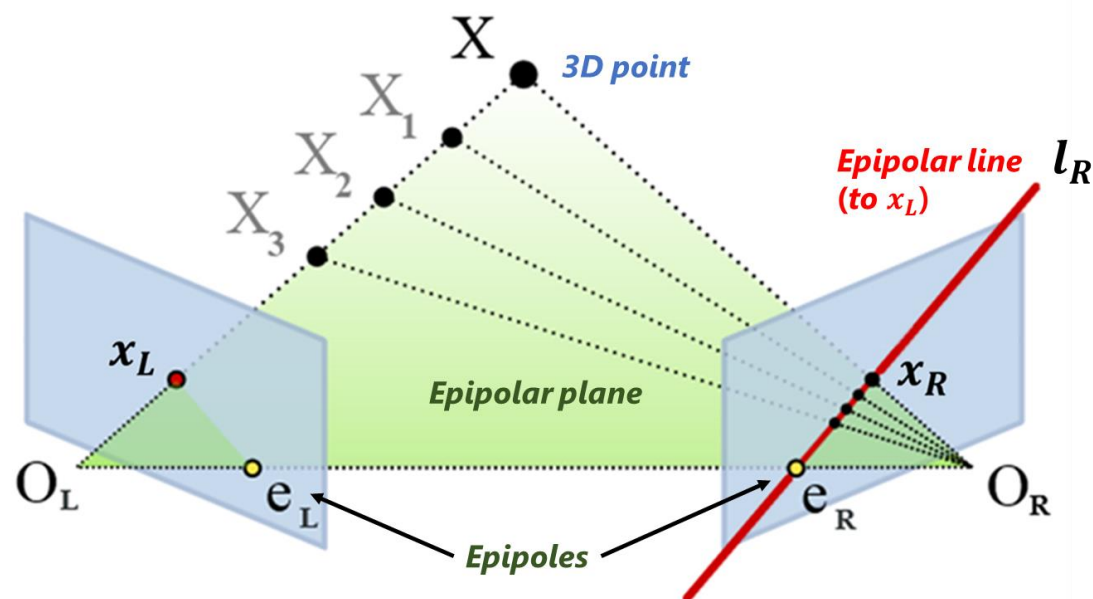


Fundamental Matrix

(focal length = 1)

Essential matrix ($\mathbf{E} \in \mathbb{R}^{3 \times 3}$): Relation between normalized pixel correspondences

Fundamental matrix ($\mathbf{F} \in \mathbb{R}^{3 \times 3}$): Relation between camera pixel correspondences



real pixel normalized

$$\hat{x}_L = K_L x_L$$

$$\hat{x}_R = K_R x_R$$

$$x_L = K_L^{-1} \hat{x}_L$$

$$x_R = K_R^{-1} \hat{x}_R$$

$$x_R^\top \mathbf{E} x_L = (K_R^{-1} \hat{x}_R)^\top \mathbf{E} (K_L^{-1} \hat{x}_L)$$

$$= x_R^\top \left(K_R^\top \right)^{-1} \mathbf{E} (K_L^{-1}) \hat{x}_L = 0$$

$$= \mathbf{F}$$

Structure-from-Motion (2): Incremental Reconstruction

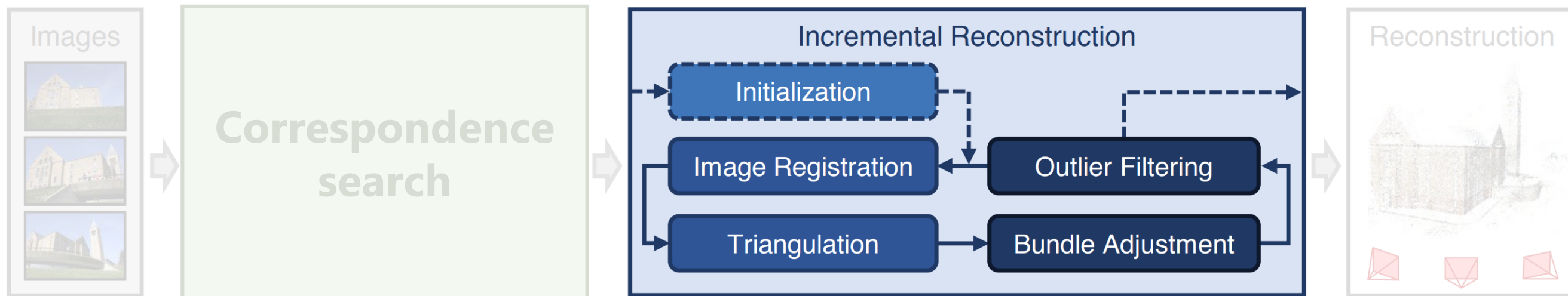
Starting from two views, aggregate more views to refine the estimation

- Camera initialization → Estimate camera **projection** matrices from $\mathbf{E} \in \mathbb{R}^{3 \times 3}$
- If we capture images with the **same** camera,

$$\mathbf{p}_{img}'^T \mathbf{F} \mathbf{p}_{img} = 0$$

$$\mathbf{F} = (\mathbf{K}^T)^{-1} \mathbf{E} \mathbf{K}^{-1}$$

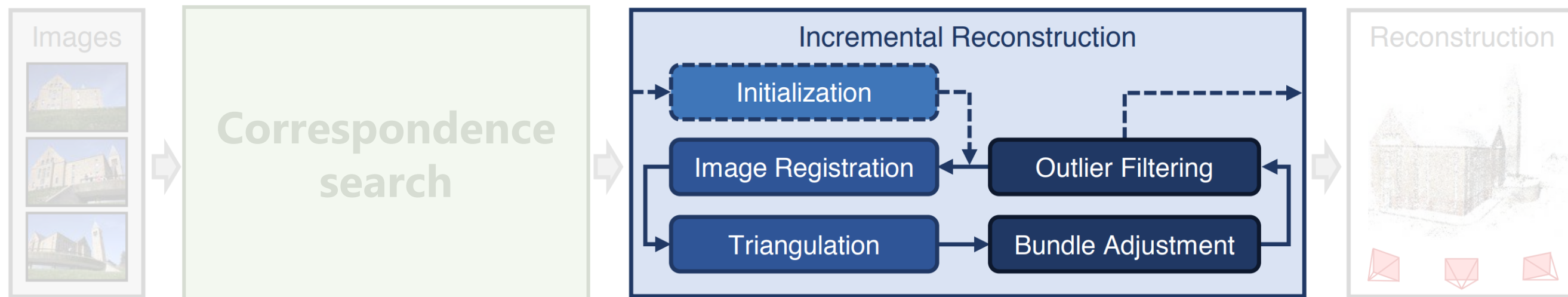
$$\mathbf{E} = \mathbf{K}^T \mathbf{F} \mathbf{K}$$



Structure-from-Motion (2): Incremental Reconstruction

Camera initialization: Estimate camera projection matrices P, P' from $E \in \mathbb{R}^{3 \times 3}$

- Suppose $P = [I|0] \in \mathbb{R}^{3 \times 4}$ (center of the world-coordinate)
- From derivation, there are 4 possible candidates for P'
- Triangulate 2D point correspondences from two views, and physically verify P'

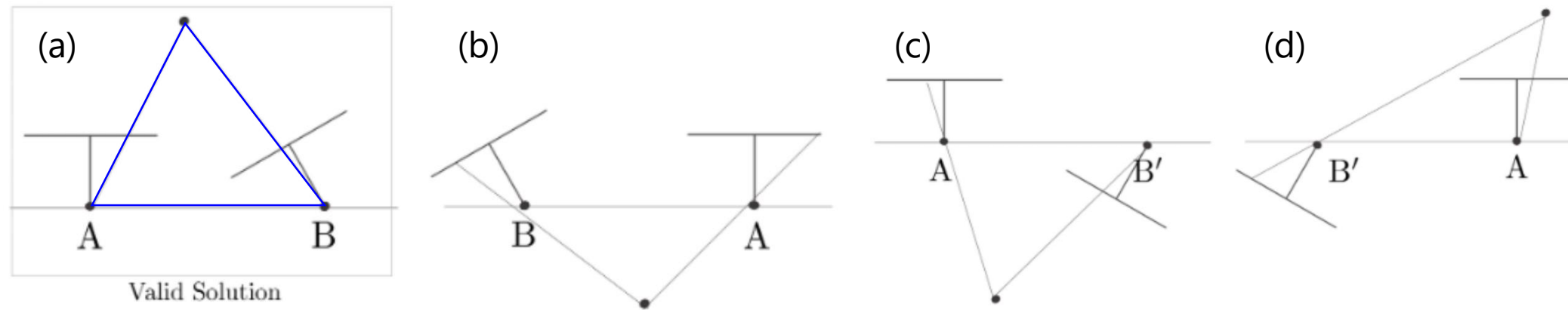


Structure-from-Motion (2): Incremental Reconstruction

Camera initialization: Estimate camera projection matrices P, P' from $E \in \mathbb{R}^{3 \times 3}$

- Suppose $P = [I|0] \in \mathbb{R}^{3 \times 4}$ (center of the world-coordinate)
- From derivation, there are 4 possible candidates for P'
- Triangulate 2D point correspondences from two views, and physically verify P'

Target 3D point is located in front of the two camera views

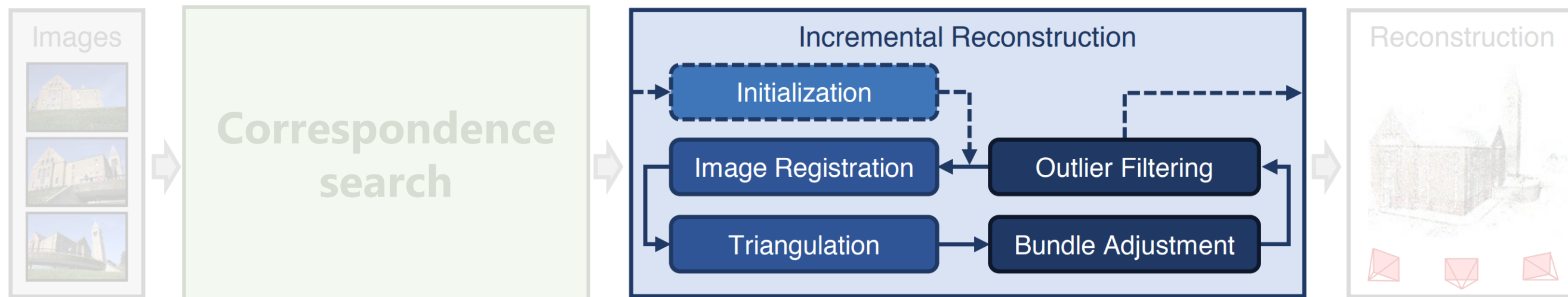


Four solutions of essential matrix

Structure-from-Motion (2): Incremental Reconstruction

Incremental reconstruction

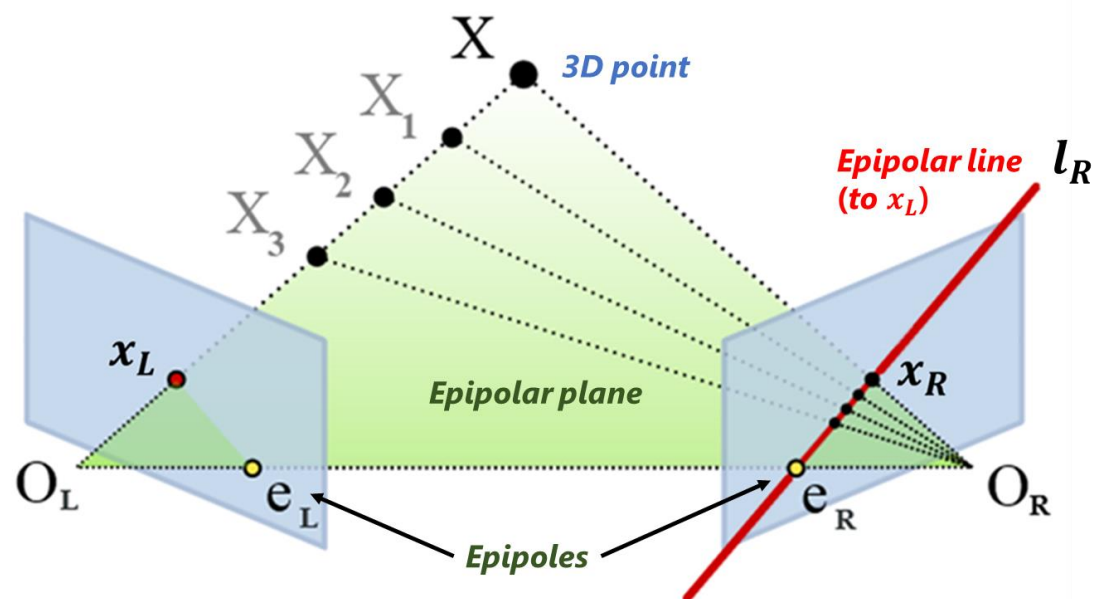
- Camera initialization → Estimate camera **projection** matrices from $\mathbf{E} \in \mathbb{R}^{3 \times 3}$
- Triangulation → **Lift** 2D points to 3D spaces using camera projection matrices
- Bundle adjustment (refinement)



Triangulation

Given a set of corresponding points x_L, x_R and camera matrices \mathbf{P}, \mathbf{P}' , estimate \mathbf{X}

Fundamental matrix ($\mathbf{F} \in \mathbb{R}^{3 \times 3}$): Relation between camera pixel correspondences



$$\mathbf{x} = \mathbf{P}\mathbf{X} = \alpha\mathbf{P}\mathbf{X}$$

homogeneous
coordinate

same direction but
differs by a scale factor

$$\mathbf{x} \times \mathbf{P}\mathbf{X} = \mathbf{0}$$

cross product is zero
(this equality removes scale factor)

$$\begin{aligned} x(\mathbf{p}_{3,row}^\top \mathbf{X}) - (\mathbf{p}_{1,row}^\top \mathbf{X}) &= 0 \\ y(\mathbf{p}_{3,row}^\top \mathbf{X}) - (\mathbf{p}_{2,row}^\top \mathbf{X}) &= 0 \\ x(\mathbf{p}_{2,row}^\top \mathbf{X}) - y(\mathbf{p}_{1,row}^\top \mathbf{X}) &= 0 \end{aligned}$$

One 2D to 3D point correspondence
gives **two** equations

$$\begin{bmatrix} x\mathbf{p}_{3,row}^\top - \mathbf{p}_{1,row}^\top \\ y\mathbf{p}_{3,row}^\top - \mathbf{p}_{2,row}^\top \\ x'\mathbf{p}'_{3T} - \mathbf{p}'_{1T} \\ y'\mathbf{p}'_{3T} - \mathbf{p}'_{2T} \end{bmatrix} \mathbf{X} = \mathbf{0}$$

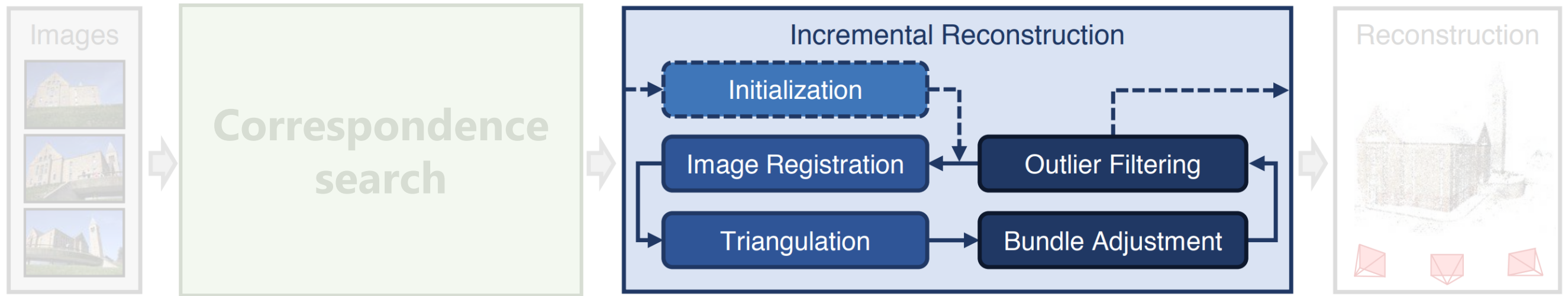
Now we can make system equations.
Use SVD to solve!

→ Two-view geometry

Structure-from-Motion (2): Incremental Reconstruction

Incremental reconstruction

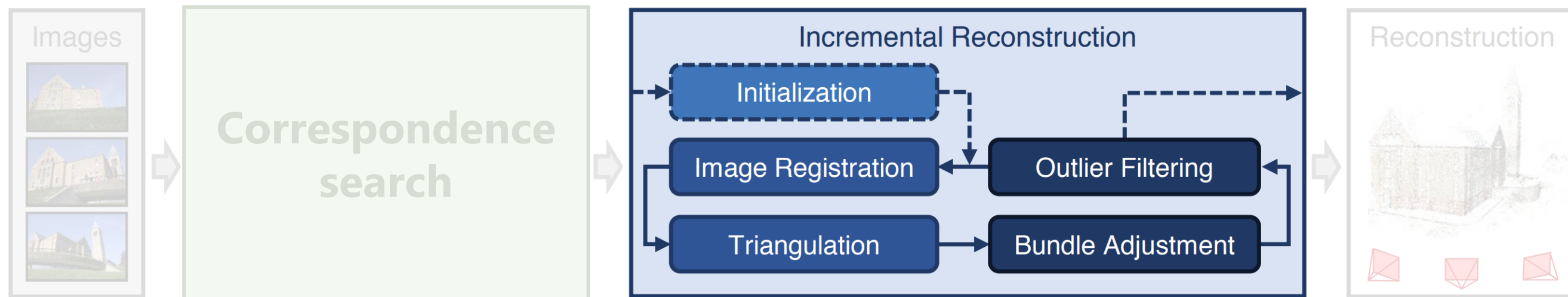
- Camera initialization → Estimate camera **projection** matrices from $\mathbf{E} \in \mathbb{R}^{3 \times 3}$
- Triangulation → **Lift** 2D points to 3D spaces using camera projection matrices
- Bundle adjustment (refinement)



Structure-from-Motion (2): Incremental Reconstruction

Bundle adjustment: Non-linear method for jointly refining SfM

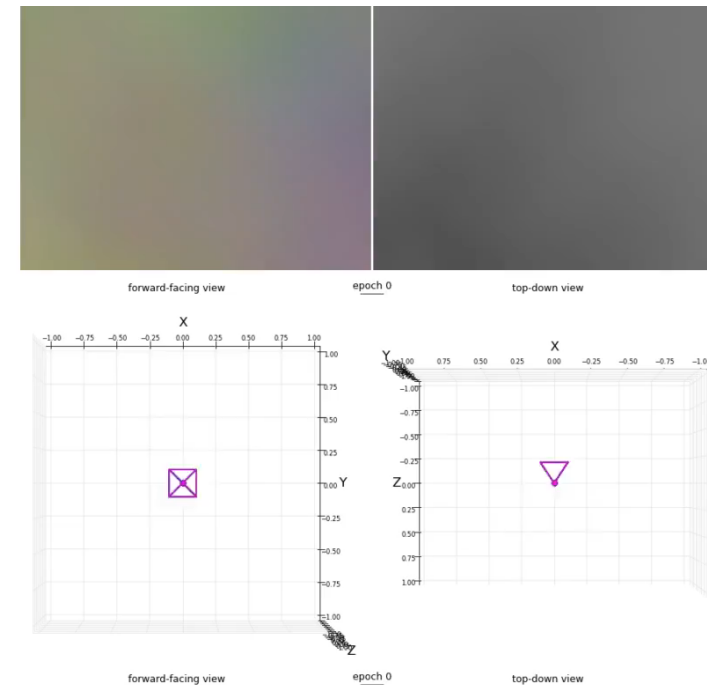
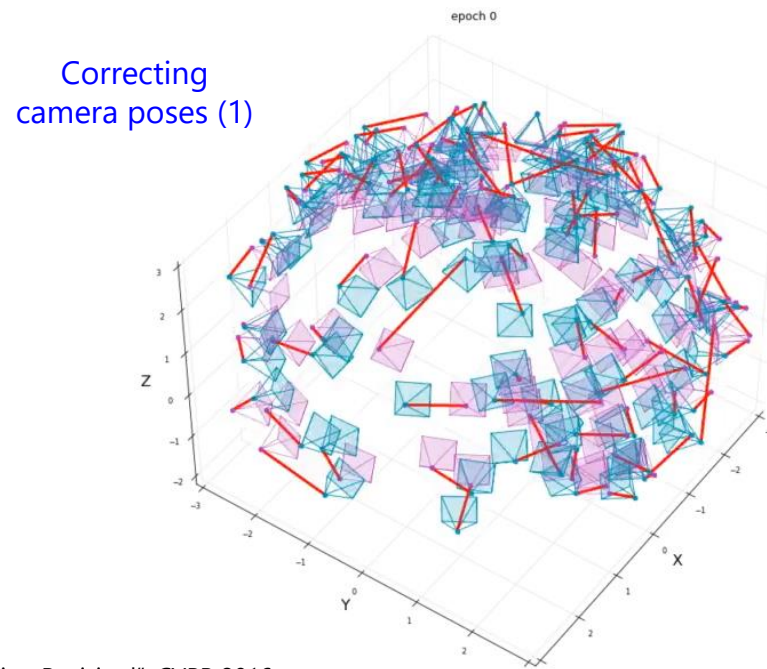
- **Bundle** of light rays leaving each 3D points and **converging** on each camera center
- Refining **3D points** and the camera **poses** by minimizing **reprojection error**
- Usually optimized with Levenberg-Marquardt optimization



Structure-from-Motion (2): Incremental Reconstruction

Bundle adjustment: Non-linear method for jointly refining SfM

- **Bundle** of light rays leaving each 3D points and **converging** on each camera center
- Refining **3D points** and the camera **poses** by minimizing **reprojection error**
- Usually optimized with Levenberg-Marquardt optimization



Correcting camera poses (2)

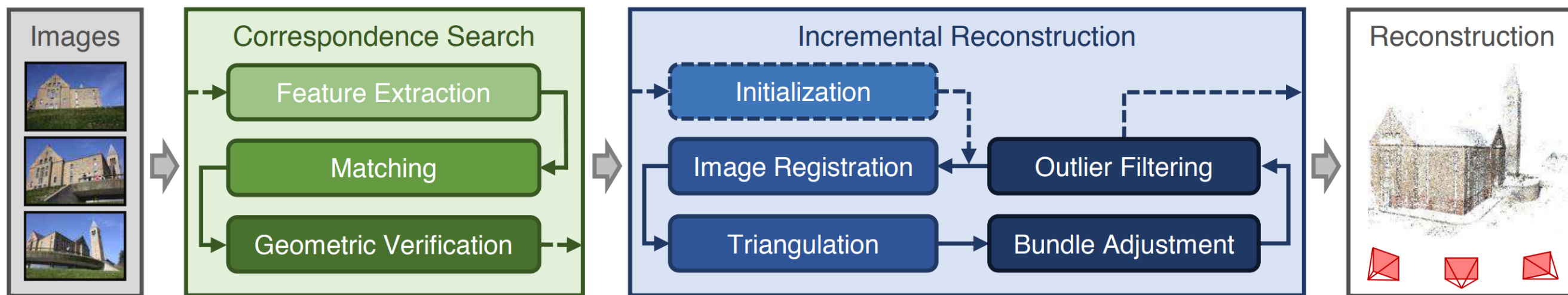
Structure-from-Motion Pipeline

Correspondence search

- Feature extraction
- Matching
- Geometric verification

Incremental reconstruction

- Camera initialization
- Triangulation
- Bundle adjustment (refinement)



Experiment: Build Your Own 3D Structure

<https://view.kentech.ac.kr/f088fa7f-874e-44bc-bd6d-6084b42dfdf7>

