# Structure from Motion and SLAM

CS 4391 Introduction to Computer Vision

Professor Yapeng Tian
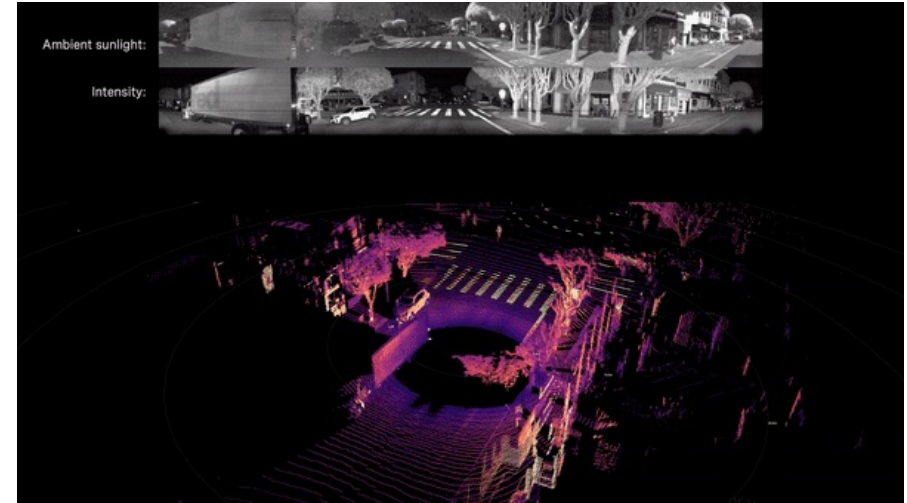
Department of Computer Science

# How to Recover the 3D World from Images?

Structure from Motion (SfM)
- Structure: the geometry of the 3D world
- Motion: camera motion
- Input: a set of images (no need to be videos)
- From computer vision

Simultaneous Localization and Mapping (SLAM)
- Localization: camera pose
- Mapping: build the geometry of the 3D world
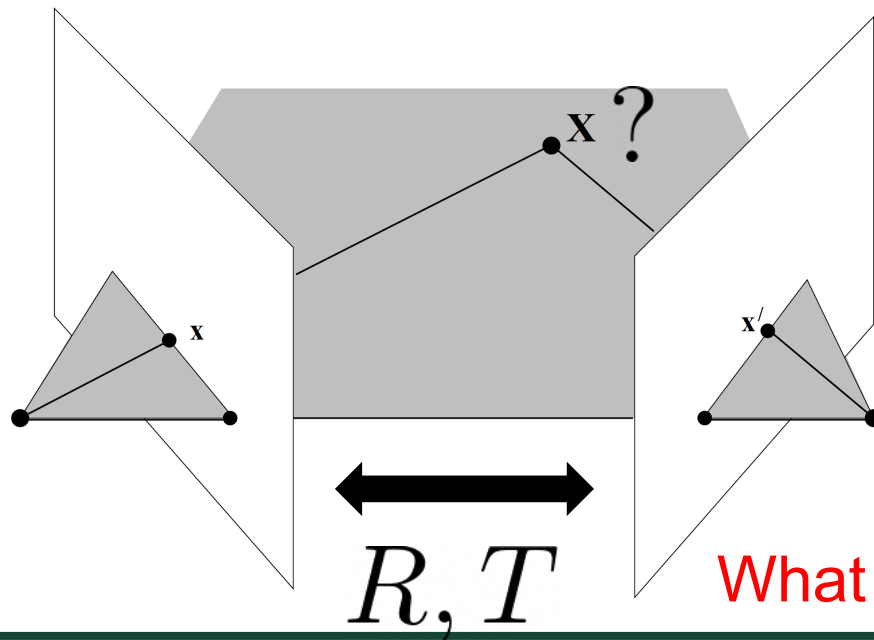- Input: video sequences
- From robotics



*Point cloud captured on an Ouster OS1-128 digital lidar sensor*

# Triangulation

Idea: using images from different views and feature matching

Triangulation from pixel correspondences to compute 3D location

Given $\mathbf{x} \longleftrightarrow \mathbf{x}'$

Intersection of two backprojected lines

$$\mathbf{X} = \mathbf{l} \times \mathbf{l}'$$

What if unknow camera pose?
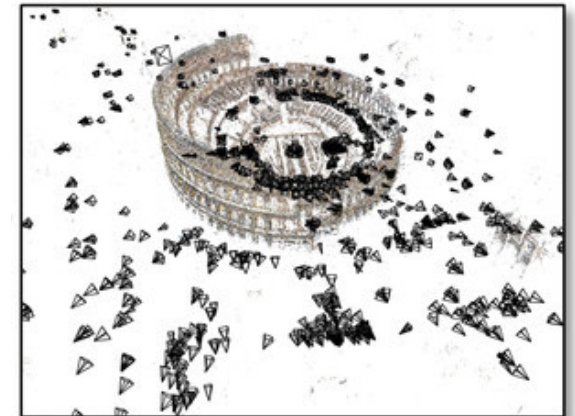
# Colosseum in Rome

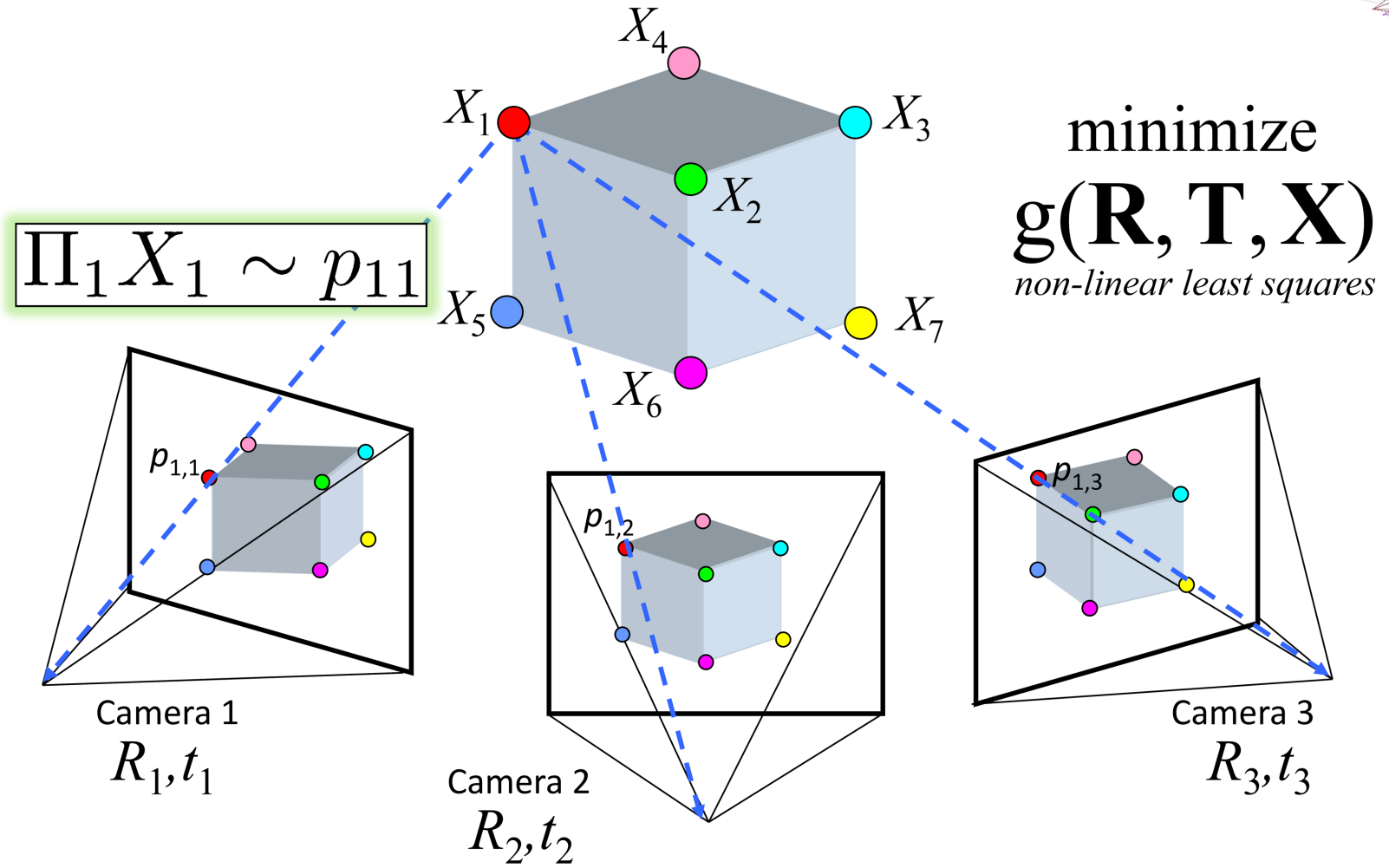# Structure from Motion

Input
- A set of images from different views
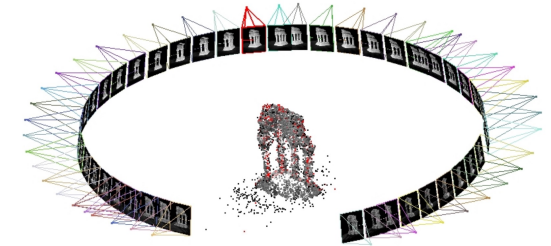
Output
- 3D Locations of all feature points in a world frame
- Camera poses of the images

# Structure from motion



$$\Pi_1 X_1 \sim p_{11}$$

minimize
$$\mathrm{g}(\mathbf{R}, \mathbf{T}, \mathbf{X})$$
*non-linear least squares*

$X_4$

$X_1$    $X_3$

$X_2$

$X_5$    $X_7$

$X_6$

$p_{1,1}$

$p_{1,2}$

$p_{1,3}$

Camera 1
$R_1, t_1$

Camera 2
$R_2, t_2$

Camera 3
$R_3, t_3$
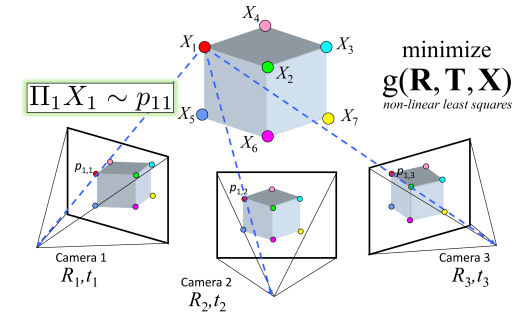
# Structure from Motion

Minimize sum of squared reprojection errors

$$g(\mathbf{X}, \mathbf{R}, \mathbf{T}) = \sum_{i=1}^{m} \sum_{j=1}^{n} w_{ij} \cdot \left\| \mathbf{P}(\mathbf{x}_i, \mathbf{R}_j, \mathbf{t}_j) - \begin{bmatrix} u_{i,j} \\ v_{i,j} \end{bmatrix} \right\|^2$$

m points, n images

predicted image location

observed image location

Indicator variable:
is point i visible in image j?

Projection

$$\begin{bmatrix} x' \\ y' \\ z' \end{bmatrix} = \mathbf{R}\mathbf{x} + \mathbf{t}$$

$$u' = f_x \frac{x'}{z'} + p_x$$

$$v' = f_y \frac{y'}{z'} + p_y$$

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = \mathbf{P}(\mathbf{x}, \mathbf{R}, \mathbf{t})$$

# Structure from Motion

How to minimize

$$g(\mathbf{X}, \mathbf{R}, \mathbf{T}) = \sum_{i=1}^{m} \sum_{j=1}^{n} w_{ij} \cdot \left\| \mathbf{P}(\mathbf{x}_i, \mathbf{R}_j, \mathbf{t}_j) - \begin{bmatrix} u_{i,j} \\ v_{i,j} \end{bmatrix} \right\|^2$$

A non-linear least squares problem (why?)
  • E.g. Levenberg-Marquardt

# The Levenberg-Marquardt Algorithm

Nonlinear least squares $\hat{\beta} \in \text{argmin}_{\beta} \, S\left(\beta\right) \equiv \text{argmin}_{\beta} \sum_{i=1}^{m} \left[y_i - f\left(x_i, \beta\right)\right]^2$

An iterative algorithm
- Start with an initial guess $\beta_0$
- For each iteration $\beta \leftarrow \beta + \delta$

How to get $\delta$ ?
- Linear approximation $f\left(x_i, \beta + \delta\right) \approx f\left(x_i, \beta\right) + \mathbf{J}_i \delta$ $\qquad \mathbf{J}_i = \dfrac{\partial f\left(x_i, \beta\right)}{\partial \beta}$

- Find to $\delta$ minimize the objective $S\left(\beta + \delta\right) \approx \sum_{i=1}^{m} \left[y_i - f\left(x_i, \beta\right) - \mathbf{J}_i \delta\right]^2$

Wikipedia

# The Levenberg-Marquardt Algorithm

Vector notation for

$$S\left(\boldsymbol{\beta}+\boldsymbol{\delta}\right) \approx \sum_{i=1}^{m}\left[y_i - f\left(x_i, \boldsymbol{\beta}\right) - \mathbf{J}_i \boldsymbol{\delta}\right]^2$$

$$
\begin{aligned}
S\left(\boldsymbol{\beta}+\boldsymbol{\delta}\right) &\approx \left\|\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right) - \mathbf{J}\boldsymbol{\delta}\right\|^2 \\
&= \left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right) - \mathbf{J}\boldsymbol{\delta}\right]^{\mathrm{T}}\left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right) - \mathbf{J}\boldsymbol{\delta}\right] \\
&= \left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right]^{\mathrm{T}}\left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right] - \left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right]^{\mathrm{T}}\mathbf{J}\boldsymbol{\delta} - \left(\mathbf{J}\boldsymbol{\delta}\right)^{\mathrm{T}}\left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right] + \boldsymbol{\delta}^{\mathrm{T}}\mathbf{J}^{\mathrm{T}}\mathbf{J}\boldsymbol{\delta} \\
&= \left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right]^{\mathrm{T}}\left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right] - 2\left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right]^{\mathrm{T}}\mathbf{J}\boldsymbol{\delta} + \boldsymbol{\delta}^{\mathrm{T}}\mathbf{J}^{\mathrm{T}}\mathbf{J}\boldsymbol{\delta}.
\end{aligned}
$$

https://www.cs.ubc.ca/~schmidtm/Courses/340-F16/linearQuadraticGradients.pdf

Take derivation with respect to $\delta$ and set to zero $\left(\mathbf{J}^{\mathrm{T}}\mathbf{J}\right)\boldsymbol{\delta} = \mathbf{J}^{\mathrm{T}}\left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right]$

Levenberg's contribution $\left(\mathbf{J}^{\mathrm{T}}\mathbf{J} + \lambda\mathbf{I}\right)\boldsymbol{\delta} = \mathbf{J}^{\mathrm{T}}\left[\mathbf{y} - \mathbf{f}\left(\boldsymbol{\beta}\right)\right]$     damped version

$$\beta \leftarrow \beta + \delta$$

Wikipedia

# Structure from Motion

$$g(\mathbf{X}, \mathbf{R}, \mathbf{T}) = \sum_{i=1}^{m} \sum_{j=1}^{n} w_{ij} \cdot \left\| \mathbf{P}(\mathbf{x}_i, \mathbf{R}_j, \mathbf{t}_j) - \begin{bmatrix} u_{i,j} \\ v_{i,j} \end{bmatrix} \right\|^2$$

*predicted* image location

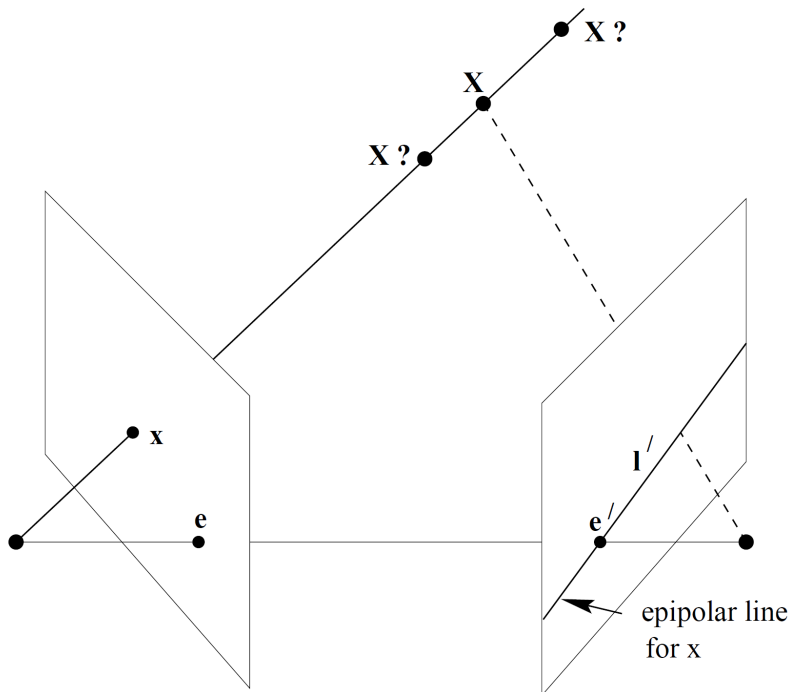*observed* image location

*indicator variable*: is point *i* visible in image *j* ?

$$\beta = (\mathbf{X}, \mathbf{R}, \mathbf{T})$$

How to get the initial estimation $\beta_0$ ?

Random guess is not a good idea.

# Matching Two Views

Fundamental matrix



$\mathbf{x}'$ is on the epiploar line $\mathbf{l}' = F\mathbf{x}$

$$\mathbf{x}'^T F \mathbf{x} = 0$$

$$\begin{bmatrix} x'_i & y'_i & 1 \end{bmatrix} \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} = 0$$

$$x_i x'_i f_{11} + x_i y'_i f_{21} + x_i f_{31} + y_i x'_i f_{12} + y_i y'_i f_{22} + y_i f_{32} + x'_i f_{13} + y'_i f_{23} + f_{33} = 0$$

$$\begin{bmatrix} x_1 x'_1 & x_1 y'_1 & x_1 & y_1 x'_1 & y_1 y'_1 & y_1 & x'_1 & y'_1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_m x'_m & x_m y'_m & x_m & y_m x'_m & y_m y'_m & y_m & x'_m & y'_m & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{21} \\ f_{31} \\ f_{12} \\ f_{22} \\ f_{32} \\ f_{13} \\ f_{23} \\ f_{33} \end{bmatrix} = 0$$

We need 8 points to solve this system.

# Matching Two Views

$$\mathbf{x}'^T F \mathbf{x} = 0$$
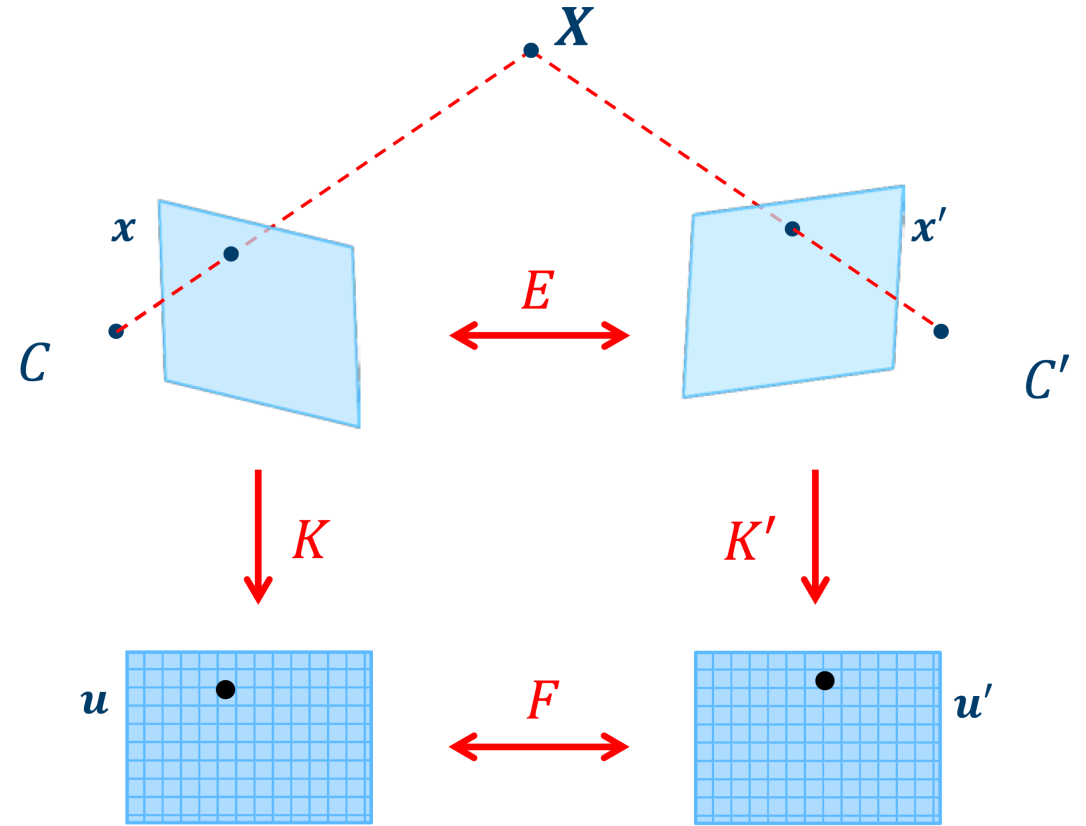
If we know camera intrinsics in SfM

$$(K'^{-1}\mathbf{x}')^T E (K^{-1}\mathbf{x}) = 0$$

**Normalized coordinates**

$$F = K'^{-T} E K^{-1}$$
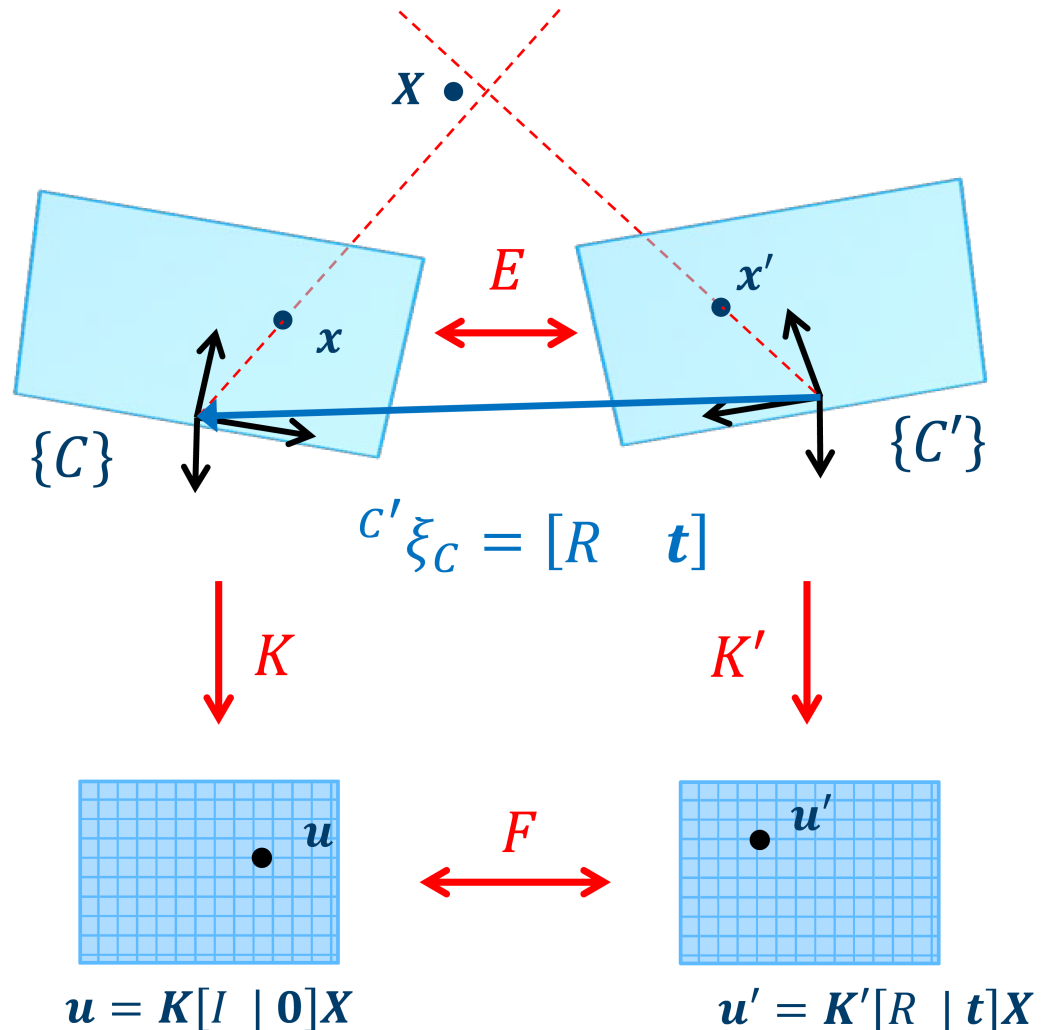
Essential matrix E

$$E = K'^T F K$$

Credit: Thomas Opsahl

# Matching Two Views

Recover the relative pose $R$ and $\boldsymbol{t}$ from the essential matrix E up to the scale of $\boldsymbol{t}$

$$\mathrm{F} = [\mathbf{e}']_\times \mathrm{K}'\mathrm{R}\mathrm{K}^{-1} = \mathrm{K}'^{-\top}[\mathbf{t}]_\times \mathrm{R}\mathrm{K}^{-1}$$

$$E = K'^T F K$$

$$\mathrm{E} = [\mathbf{t}]_\times \mathrm{R}$$



$${}^{C'}\xi_C = [R \quad \boldsymbol{t}]$$

$$\boldsymbol{u} = \boldsymbol{K}[I \mid \mathbf{0}]\boldsymbol{X} \qquad \boldsymbol{u}' = \boldsymbol{K}'[R \mid \boldsymbol{t}]\boldsymbol{X}$$

Credit: Thomas Opsahl

H. C Longuet-Higgins, *A computer algorithm for reconstructing a scene from two projections*, 1981

# Matching Two Views
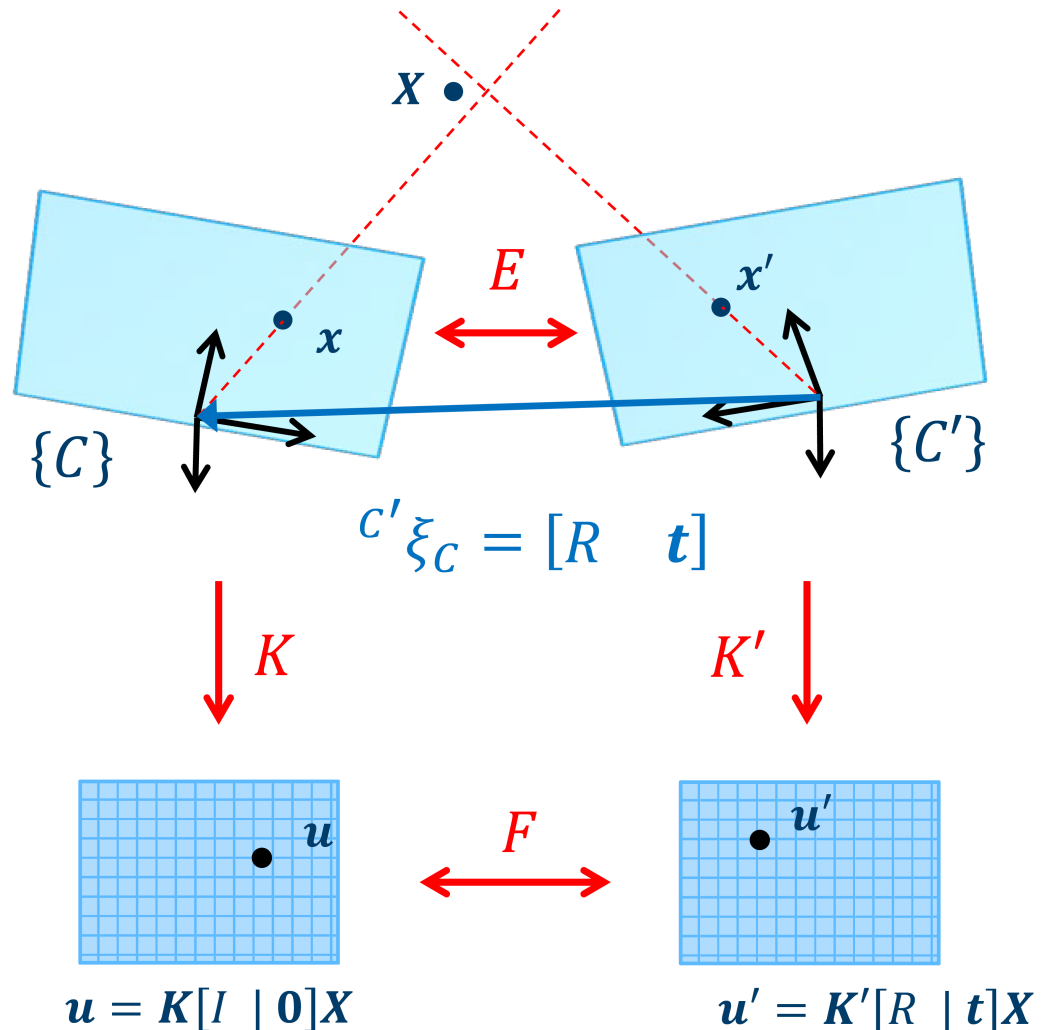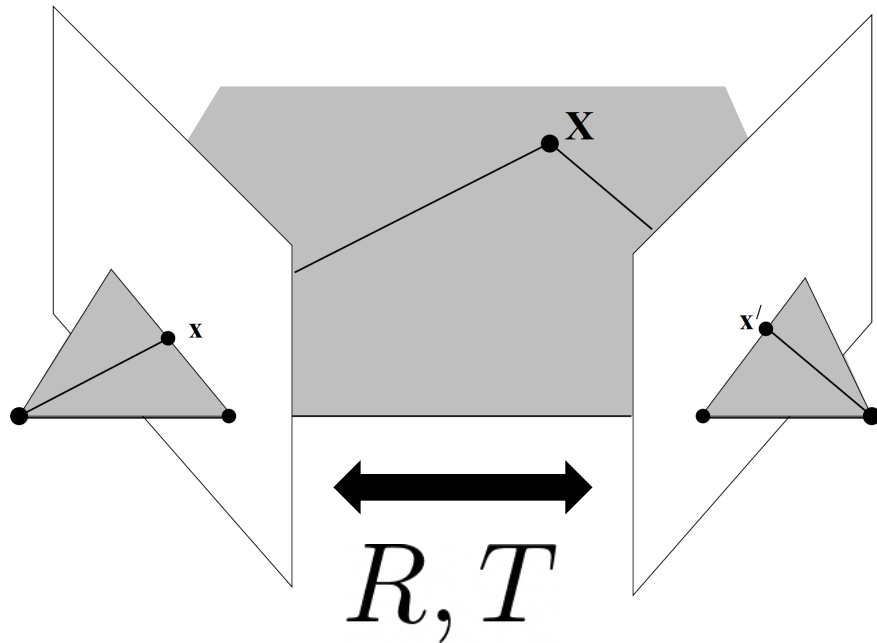
$$\mathrm{E} = [\mathbf{t}]_\times \mathrm{R}$$

$$E \cdot \mathbf{t} = [\mathbf{t}]_\times R \cdot \mathbf{t}$$
$$= (\mathbf{t} \times R) \cdot \mathbf{t} = 0$$

Use SVD to solve for **t**

$$R = -[\mathbf{t}]_\times E$$



$$^{C'}\xi_C = [R \quad t]$$

$$u = K[I \mid 0]X \qquad u' = K'[R \mid t]X$$

Credit: Thomas Opsahl

H. C Longuet-Higgins, *A computer algorithm for reconstructing a scene from two projections*, 1981

# Triangulation



Estimated from essential matrix E

Intersection of two backprojected lines

$$\mathbf{X} = \mathbf{l} \times \mathbf{l}'$$

How to get the initial estimation $\beta_0$ ?

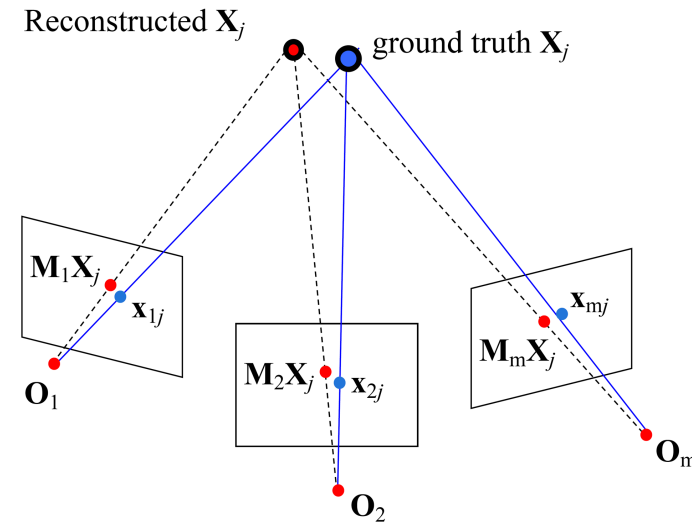$$\beta = (\mathbf{X}, \mathbf{R}, \mathbf{T})$$

# Structure from Motion

Bundle adjustment

- Iteratively refinement of structure (3D points) and motion (camera poses)

- Levenberg-Marquardt algorithm

$$g(\mathbf{X}, \mathbf{R}, \mathbf{T}) = \sum_{i=1}^{m} \sum_{j=1}^{n} w_{ij} \cdot \left\| \mathbf{P}(\mathbf{x}_i, \mathbf{R}_j, \mathbf{t}_j) - \begin{bmatrix} u_{i,j} \\ v_{i,j} \end{bmatrix} \right\|^2$$

*indicator variable*:
is point *i* visible in image *j* ?

*predicted* image location

*observed* image location

Reconstructed $\mathbf{X}_j$     ground truth $\mathbf{X}_j$

$\mathbf{M}_1\mathbf{X}_j$     $\mathbf{x}_{1j}$

$\mathbf{M}_2\mathbf{X}_j$  $\mathbf{x}_{2j}$     $\mathbf{M}_m\mathbf{X}_j$     $\mathbf{x}_{mj}$

$\mathbf{O}_1$

$\mathbf{O}_2$

$\mathbf{O}_m$

Examples: http://vision.soic.indiana.edu/projects/disco/

# Build Rome in One Day



https://grail.cs.washington.edu/rome/

# Simultaneous Localization and Mapping (SLAM)

Localization: camera pose tracking

Mapping: building a 2D or 3D representation of the environment

The goal here is the same as structure from motion but with video input
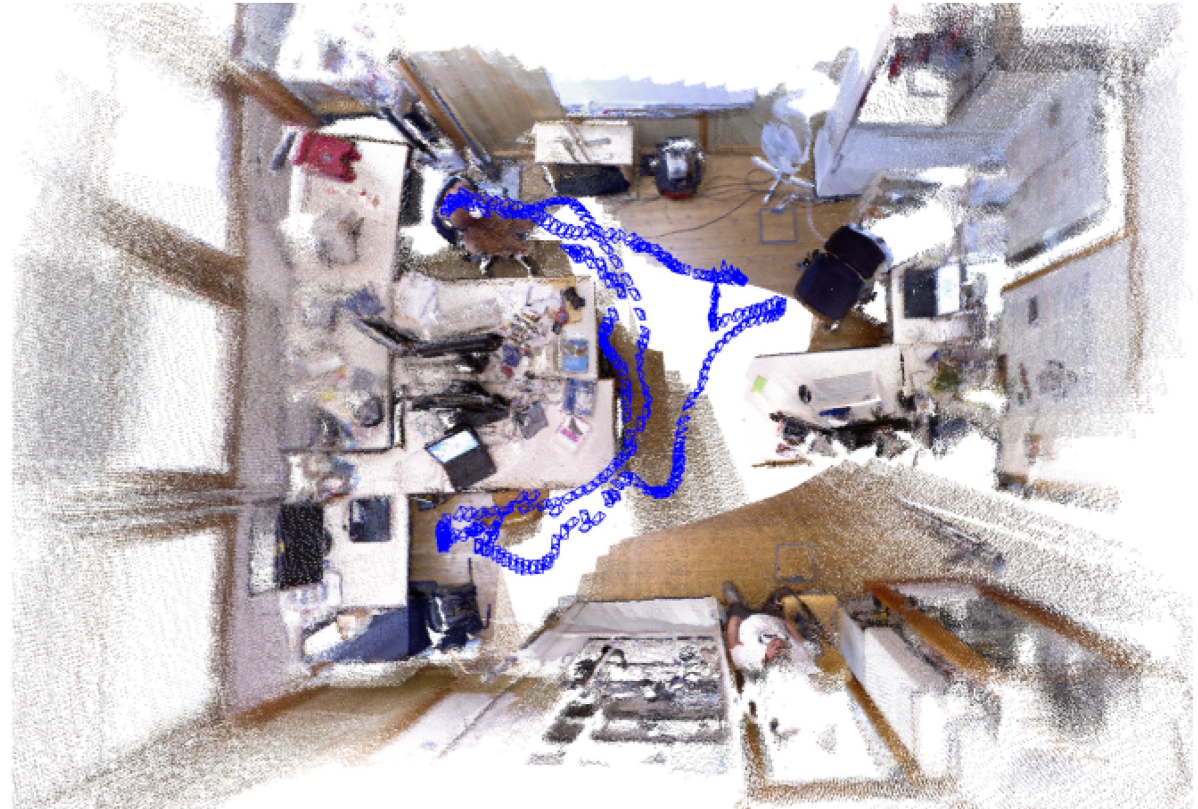


ORB-SLAM2
- Point cloud and camera poses

# Case Study: ORB-SLAM

- Oriented FAST and Rotated BRIEF (ORB)

- Tracking camera poses
  - Motion only Bundle Adjustment (BA)

- Mapping
  - Local BA around camera pose (3D location refinement)

- Loop closing
  - Loop detection



https://webdiis.unizar.es/~raulmur/orbslam/
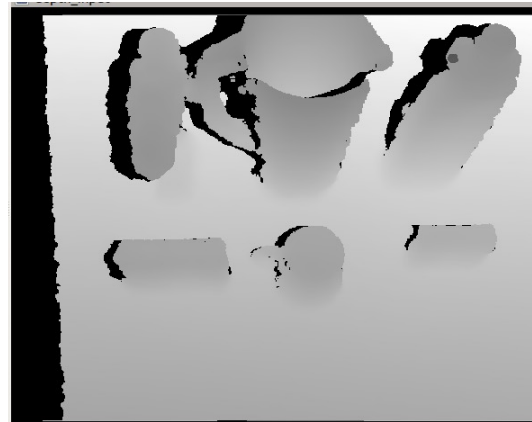
# Case Study: ORB-SLAM

# RGB-D SLAM

RGB-D cameras



Microsoft Kinect



Intel RealSense

Using depth images: 3D points in the camera frame
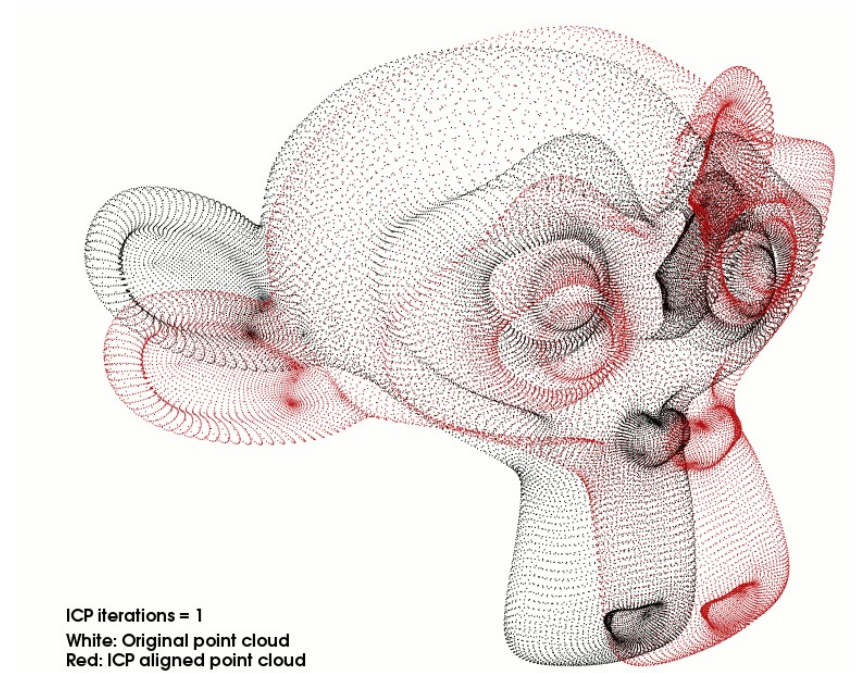


Point Cloud

# RGB-D SLAM

## Camera pose tracking

- Iterative closest point (ICP) algorithm

Input: source point cloud, target point cloud
Output: rigid transformation from source to target

- For i in range(N)
  - For each point in the source, find the closest point in the target (correspondences)
  - Estimation R and T using the correspondences
  - Transform the source points using R and T

ICP iterations = 1
White: Original point cloud
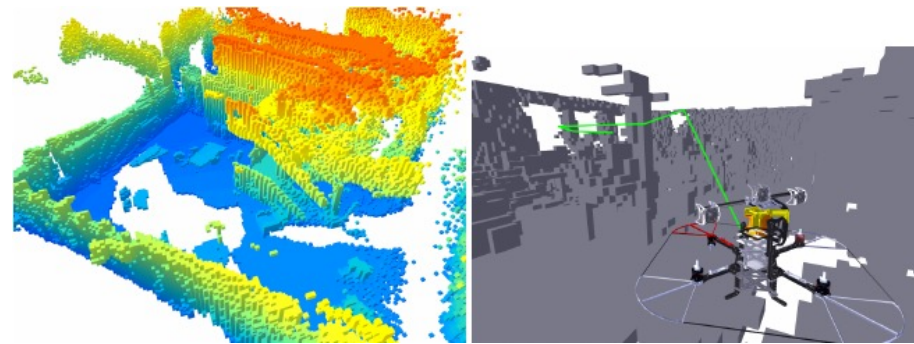Red: ICP aligned point cloud

# RGB-D SLAM

Mapping: fuse point clouds into a global frame
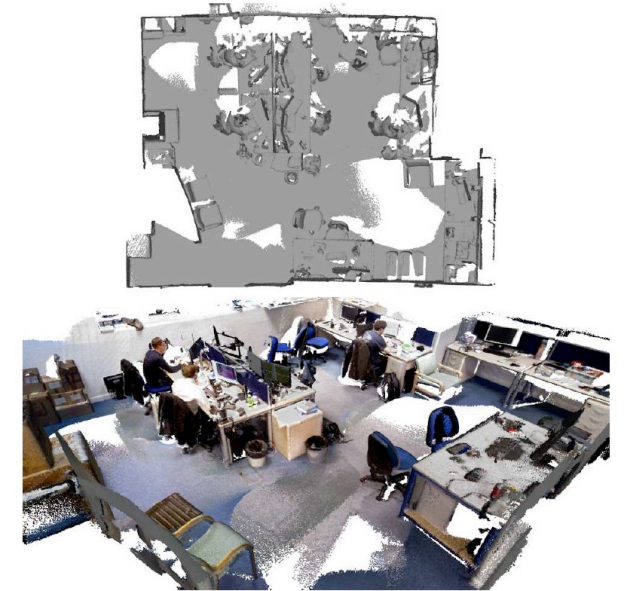
Map representation



Point clouds

ORB-SLAM

Voxels

Visual Odometry and Mapping for Autonomous
Flight Using an RGB-D Camera. Huang, et al. 2011

Surfels (small 3D surface)

ElasticFusion

# KinectFusion



https://youtu.be/of6d7C_ZWwc

# Further Reading

Chapter 11, Computer Vision, Richard Szeliski

KinectFusion: Real-Time Dense Surface Mapping and Tracking. Newcombe et al., ISMAR'11

ORB-SLAM https://webdiis.unizar.es/~raulmur/orbslam/