

# No. 4

## 動きからの3次元復元

# Structure from Motion and SLAM

担当教員：向川康博・田中賢一郎

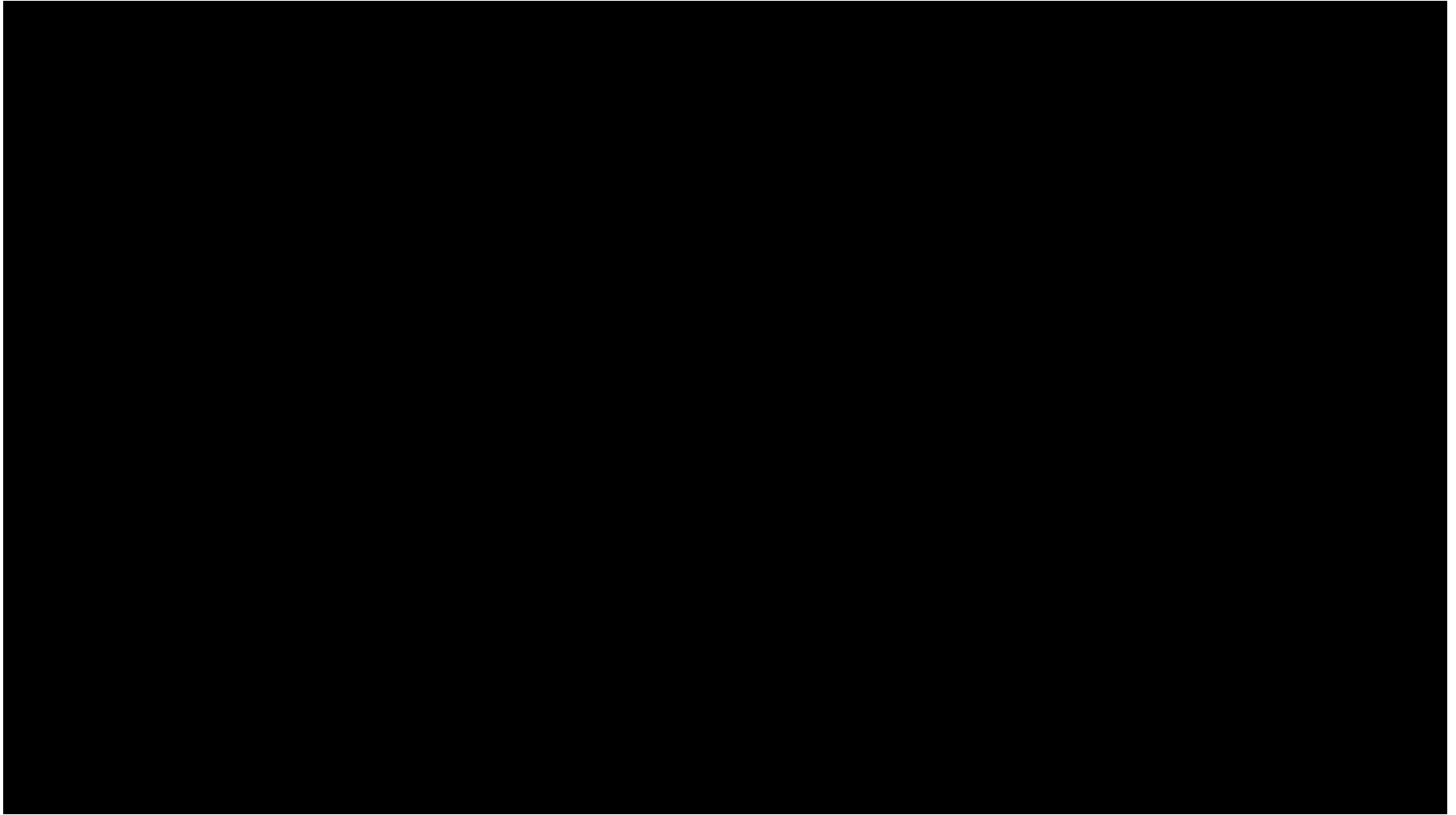
# Slide credits

- Special Thanks: Some slides are adopted from other instructors' slides.
  - Tomokazu Sato, former NAIST CV1 class
  - James Tompkin, Brown CSCI 1430 Fall 2017
  - Ioannis Gkioulekas, CMU 16-385 Spring 2018
  
  - We also thank many other instructors for sharing their slides.

# Today's topics

- Reminder: calibration and stereo
- 2-view Structure from Motion (SfM)
- Multi-view SfM
  - Large-scale SfM
- visual SLAM

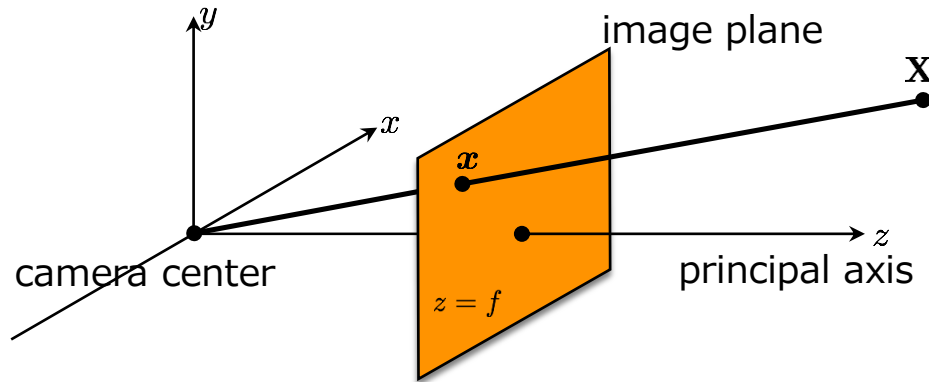
# Direct Sparse Odometry



<https://www.youtube.com/watch?v=C6-xwSOOdqQ>

# Recap: Calibration and Stereo

- Camera calibration

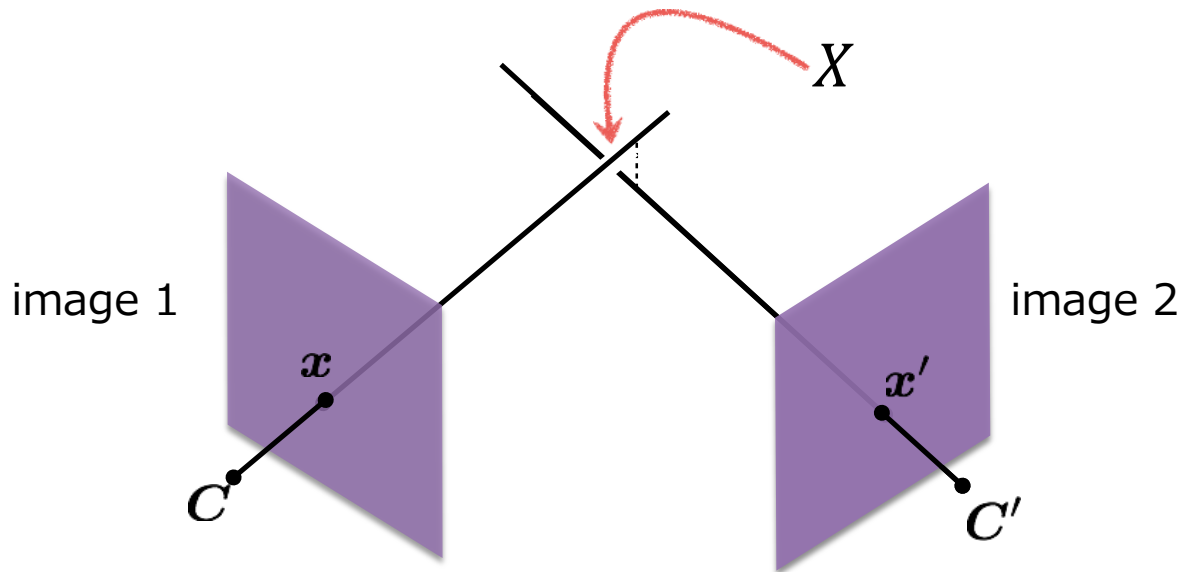


$$\mathbf{x} = \mathbf{M}\mathbf{X}$$

known      estimate      known

Solve Perspective  
n-Points (PnP) problem

- Stereo



$$\mathbf{x} = \mathbf{M}\mathbf{X}$$

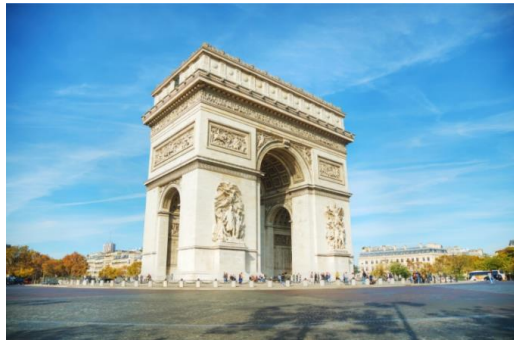
known      known      estimate

Solve Triangulation problem

# Reconstruction (再構成)

	Structure (scene geometry)	Motion (camera geometry)	Measurements
Pose Estimation	known	<b>estimate</b>	3D to 2D correspondences
Triangulation	<b>estimate</b>	known	2D to 2D correspondences
Reconstruction	<b>estimate</b>	<b>estimate</b>	2D to 2D correspondences

# How do you reconstruct?



We don't know both scene geometry and camera geometry

# Reconstruction Problem

- Problems so far

Camera calibration

$$\mathbf{x} = \mathbf{M}\mathbf{X}$$

known estimate known

Triangulation (Stereo)

$$\mathbf{x} = \mathbf{M}\mathbf{X}$$

known known estimate

- Can we jointly estimate  $\mathbf{M}$  and  $\mathbf{X}$ ?

$$\mathbf{x} = \mathbf{M}\mathbf{X}$$

known estimate estimate



# 2-view SfM

Structure from Motion

# Two-view SfM

1. Compute the Fundamental Matrix  $F$  from points correspondences

Mini-report 1: How many point pairs are necessary?



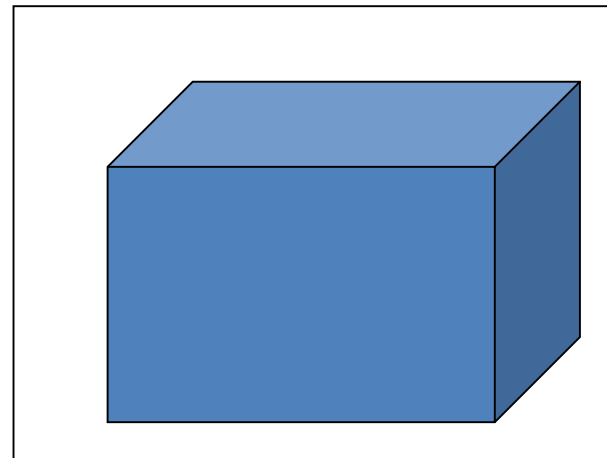
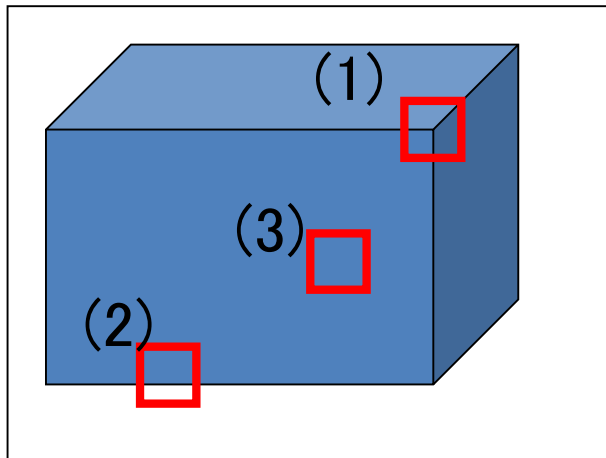
Image feature matching

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

Estimate  $F$  from matched pairs

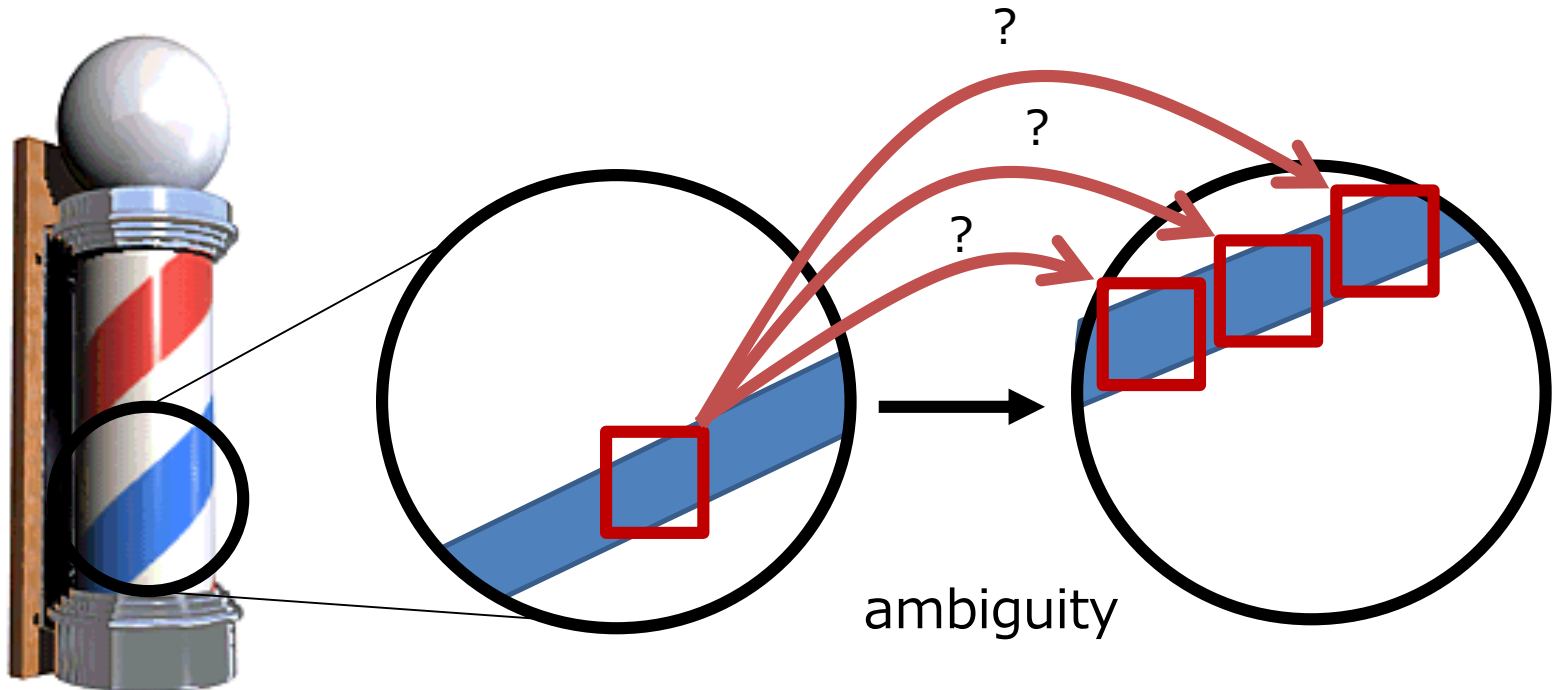
# Good feature point

- Mini-report:2
  - Which is the good feature point to find corresponding point uniquely? What is the reason?
    - (1) corner
    - (2) edge
    - (3) flat



# Aperture problem (窓問題)

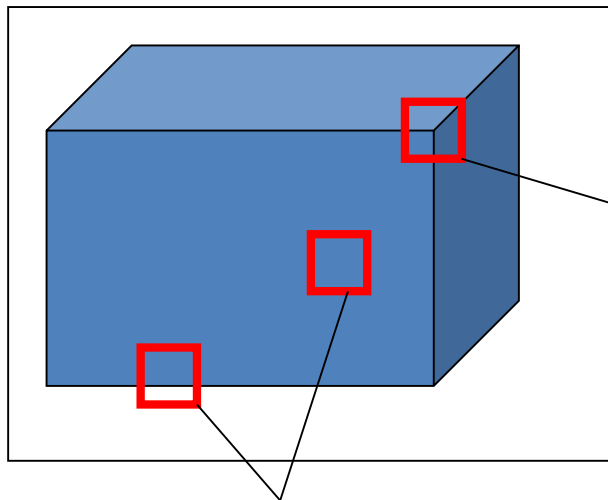
- Human eyes interpret conveniently.
- Corresponding point cannot be uniquely determined.



barber-pole illusion

# Interest operators

- Corner detectors to find small areas which are suitable as feature points
  - Moravec corner detector (1980)
  - Harris corner detector (1988)
  - KLT(Kanade-Locus-Tomasi) corner detector(1991)
- Check whether small area contains edges in multiple directions



Good feature point



Bad feature points: corresponding point cannot be determined uniquely.

# Harris corner detector

- Principal component analysis of gradient in small area
- Classification by Eigenvalues  $\lambda_1$  and  $\lambda_2$

1. X and Y direction differentiation

$$X = \partial I / \partial x \quad Y = \partial I / \partial y$$

2. Calculation of variance and covariance

$$A = X^2 \otimes w \quad B = Y^2 \otimes w \quad C = (XY) \otimes w$$

$$w_{u,v} = \exp - (u^2 + v^2) / 2\sigma^2$$

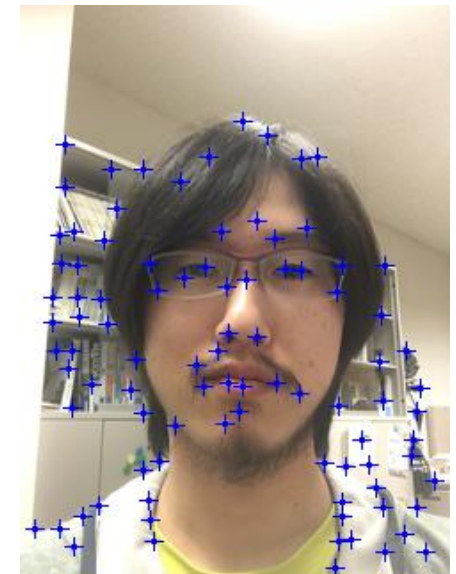
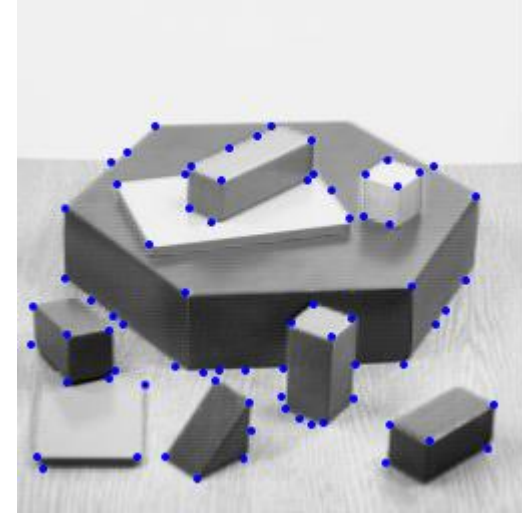
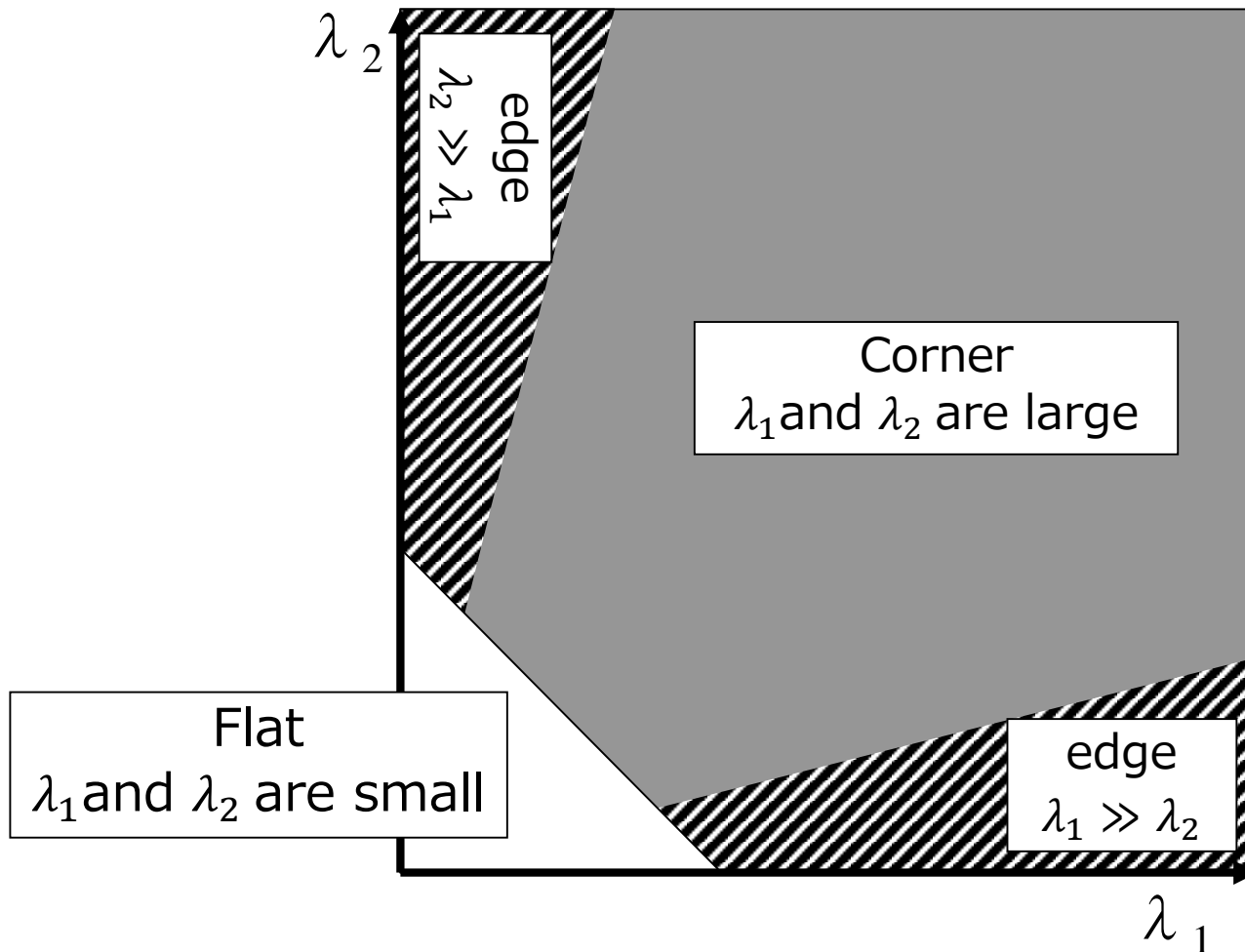
3. Eigenvalues of the variance-covariance matrix

$$M = \begin{bmatrix} A & C \\ C & B \end{bmatrix} \quad \lambda_1, \lambda_2$$

Principal component analysis  
of image gradient

# Harris corner detector

- Eigenvalue-based classification



# Image feature descriptor

## - SIFT -

- Features that are invariant to changes in scale and rotation of the image
- Intensity gradient histogram
  - Rotate coordinate axis in gradient direction (invariant to rotational change)
  - Normalize the vector sum (to reduce the influence of lighting changes)

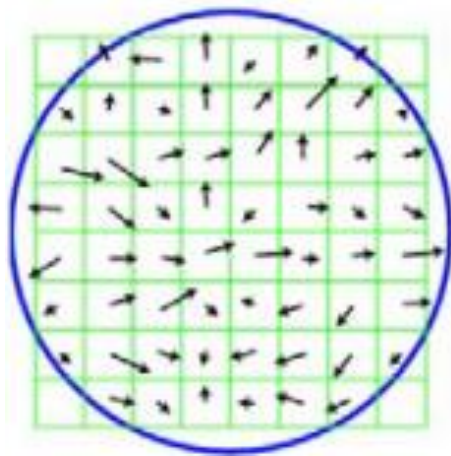
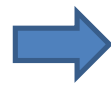
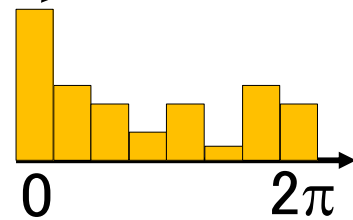


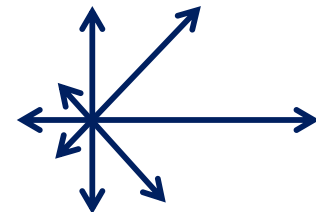
Image gradient



Detect peak position  
(gradient direction)



Angle histogram

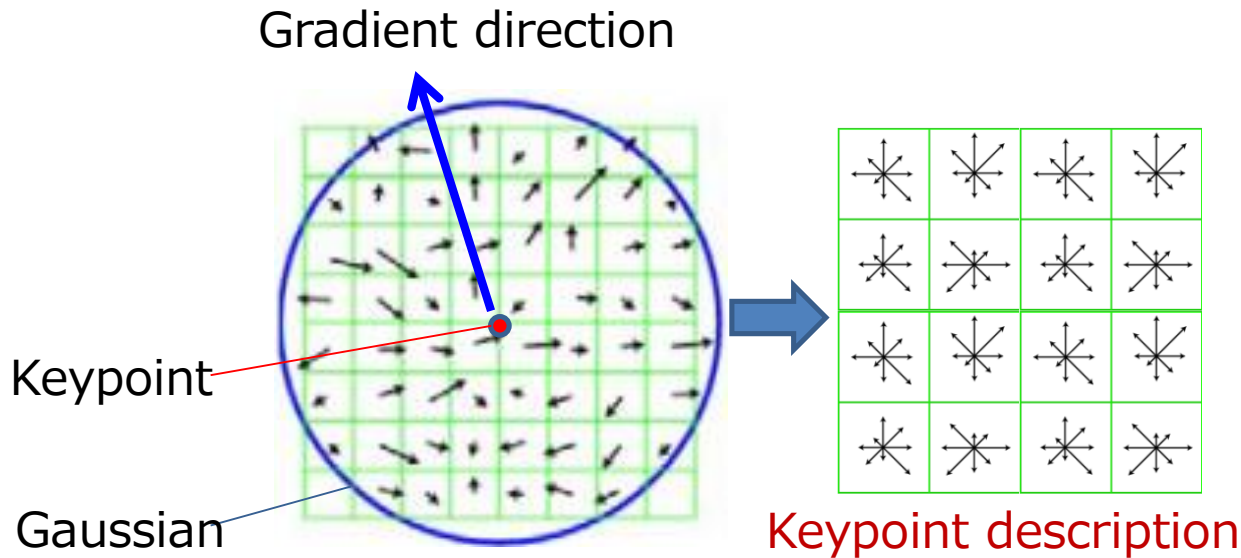




# Image feature descriptor

## - SIFT -

- Accumulate while overlapping gradient information around key points with Gaussian function
- Histogram in 8 directions every  $4 \times 4 = 16$  blocks
  - 128 dimensional feature vector



# Two-view SfM

1. Compute the Fundamental Matrix  $F$  from points correspondences  
**8-point algorithm**



Image feature matching

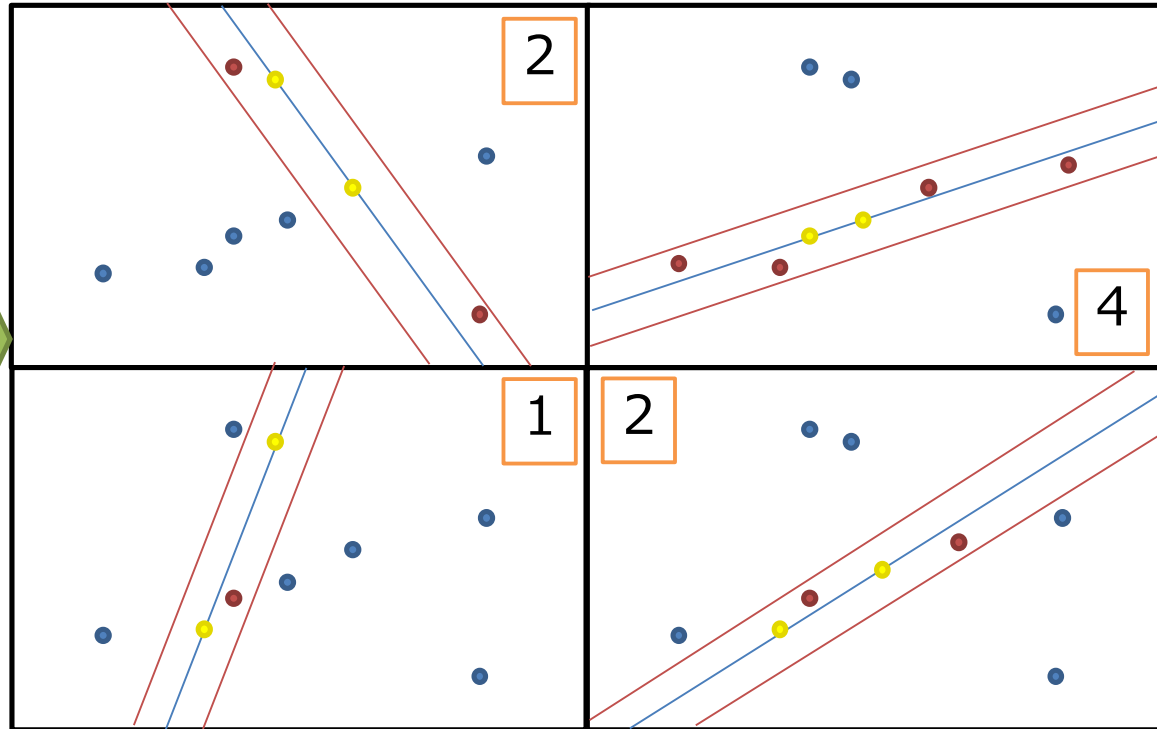
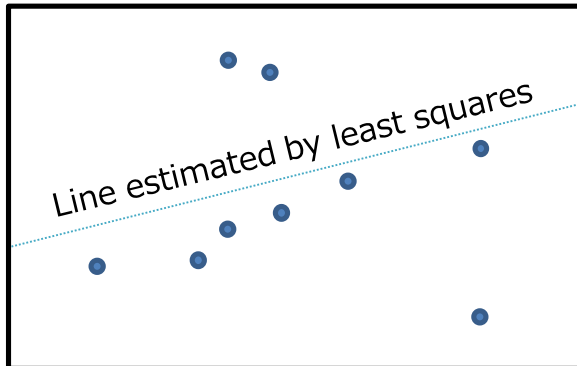
$$\mathbf{p}'^T \mathbf{F} \mathbf{p} = 0$$

Estimate  $F$  from matched pairs

# RANSAC (RANdOm SAMpling Consensus)

Eliminating outliers (外れ値).

line fitting example

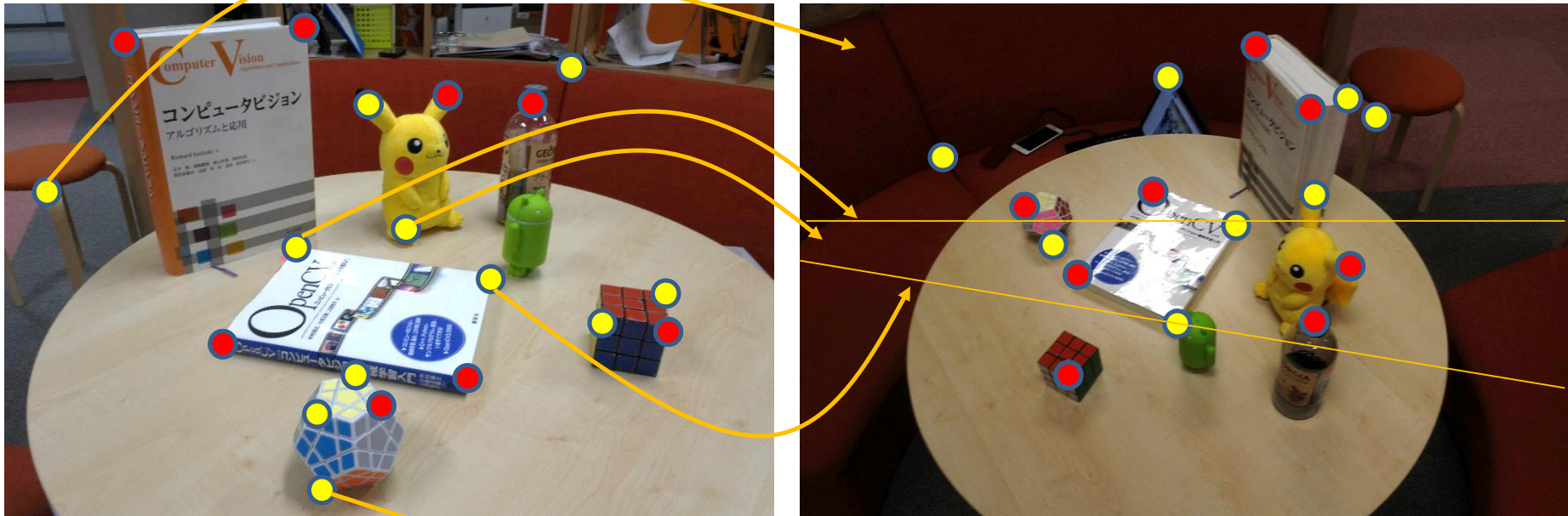


1. Randomly sample two points from given points
3. Repeat 1 and 2 for given number of iterations.

2. Count **inliers** for the line that passes selected two points.
4. Select the line that maximize the number of inliers.

# 8-point algorithm

- Randomly sampled 8 points



Do the remaining points agree?

# Two-view SfM

1. Compute the Fundamental Matrix  $F$  from points correspondences

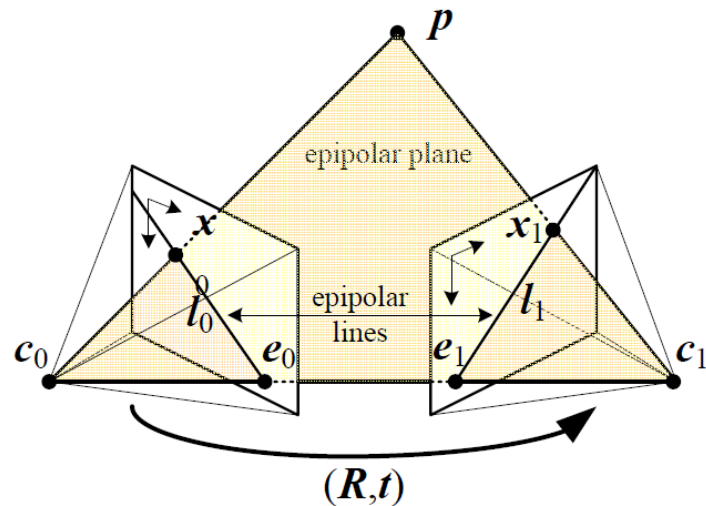
**8-point algorithm**

2. Compute the camera matrices  $M$  from the Fundamental matrix

$$M = [I|0] \text{ and } M' = [[e_x]F|e']$$



# Estimation of camera pose from pairs of feature points



Even if we do not know 3D positions of feature points, we can estimate camera's rotation  $\mathbf{R}$  and translation **direction**  $\mathbf{t}$  by epipolar constraint from multiple 2D positions of corresponding feature pairs.

$$\mathbf{M} = [\mathbf{I} | \mathbf{0}] \text{ and } \mathbf{M}' = [[\mathbf{e}_x] \mathbf{F} | \mathbf{e}']$$

It should be noted that, we **cannot recover the scale of  $\mathbf{t}$**  (actual distance between  $c_0$  and  $c_1$ ).



東武ワールドスクエア

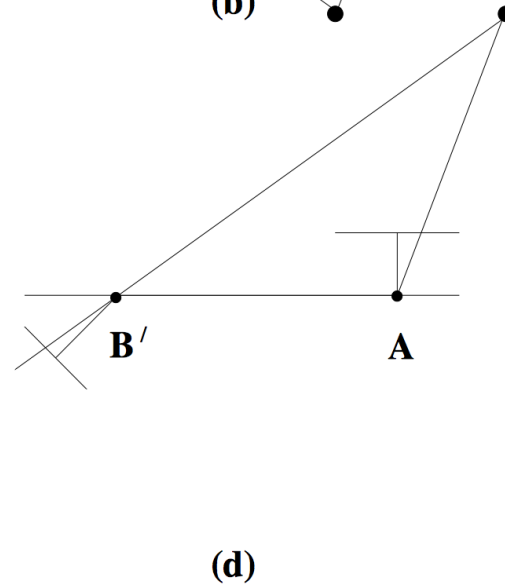
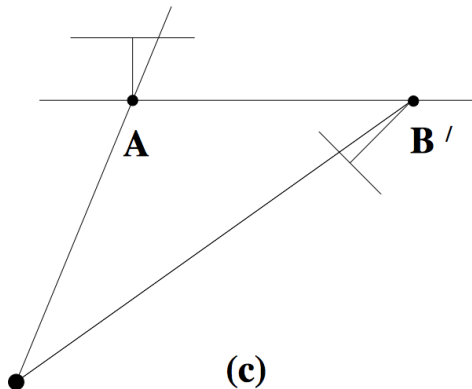
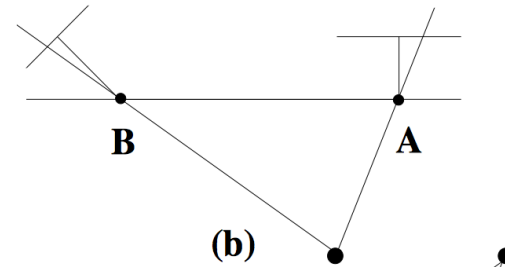
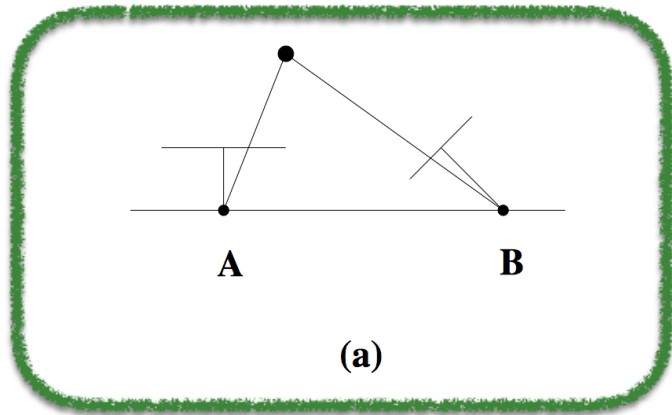
# Unknown real size

- Miniature effect



# Front/back ambiguity

- Find the configuration where the points is in front of both cameras





# Two-view SfM

1. Compute the Fundamental Matrix  $F$  from points correspondences

**8-point algorithm**

2. Compute the camera matrices  $M$  from the Fundamental matrix

$$M = [I|0] \text{ and } M' = [[e_x]F|e']$$

3. For each point correspondence, compute the point  $X$  in 3D space

**Triangulate** with  $x = MX$  and  $x' = M'X$

# 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 in the image remain exactly the same:

$$\mathbf{x} = \mathbf{M}\mathbf{X} = \left( \frac{1}{k} \mathbf{M} \right) (k\mathbf{X})$$

It is impossible to recover the absolute scale of the scene!

# Structure from motion ambiguity

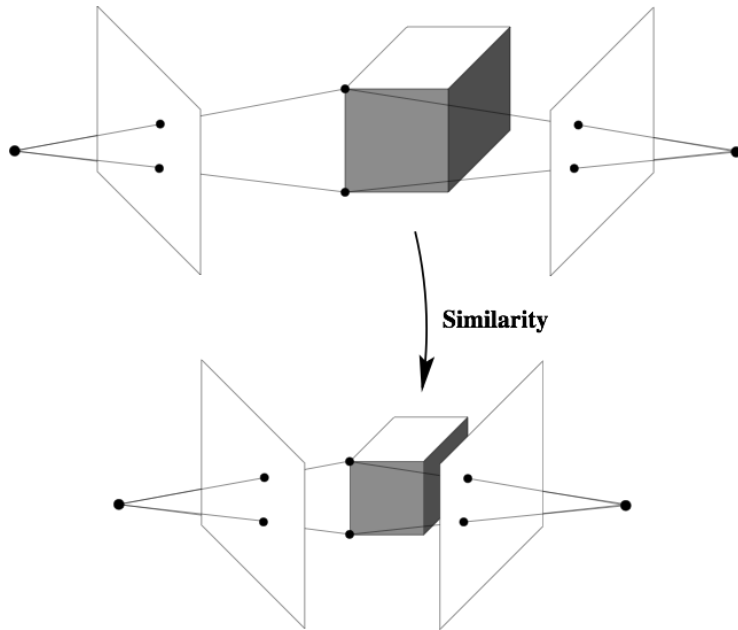
- If we scale the entire scene by some factor  $k$  and, at the same time, scale the camera matrices by the factor of  $1/k$ , the projections of the scene points in the image remain exactly the same
- More generally: if we transform the scene using a transformation  $Q$  and apply the inverse transformation to the camera matrices, then the images do not change

$$\mathbf{x} = \mathbf{M}\mathbf{X} = (\mathbf{M}\mathbf{Q}^{-1})(\mathbf{Q}\mathbf{X})$$

# Reconstruction ambiguity

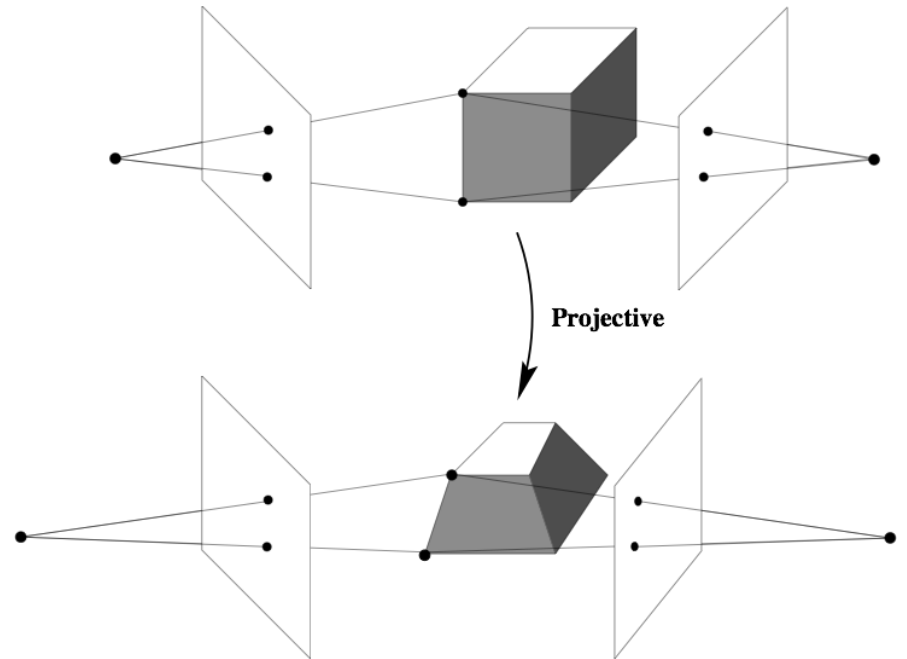
## Calibrated cameras

(similarity projection ambiguity)



## Uncalibrated cameras

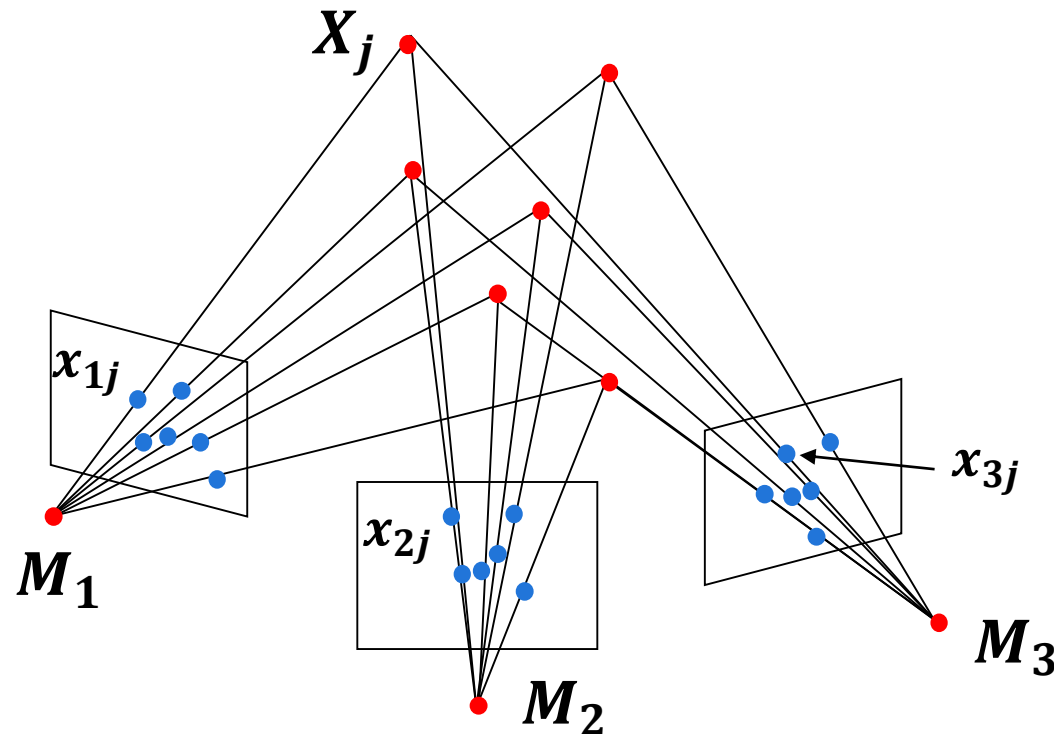
(projective projection ambiguity)



# Multi-view SfM

# Projective structure from motion

- Given:  $m$  images of  $n$  fixed 3D points
  - $x_{ij} = M_i X_j, i = (1, \dots, m,) j = (1, \dots, n)$
- Problem: estimate  $m$  projection matrices  $M_i$  and  $n$  3D points  $X_j$  from the  $mn$  correspondences  $x_{ij}$



# Sequential SfM

Multi-view, ordered images

# Initialization

- Initialize by 2-view SfM

1. Compute the Fundamental Matrix  $F$  from points correspondences

**8-point algorithm**

2. Compute the camera matrices  $M$  from the Fundamental matrix

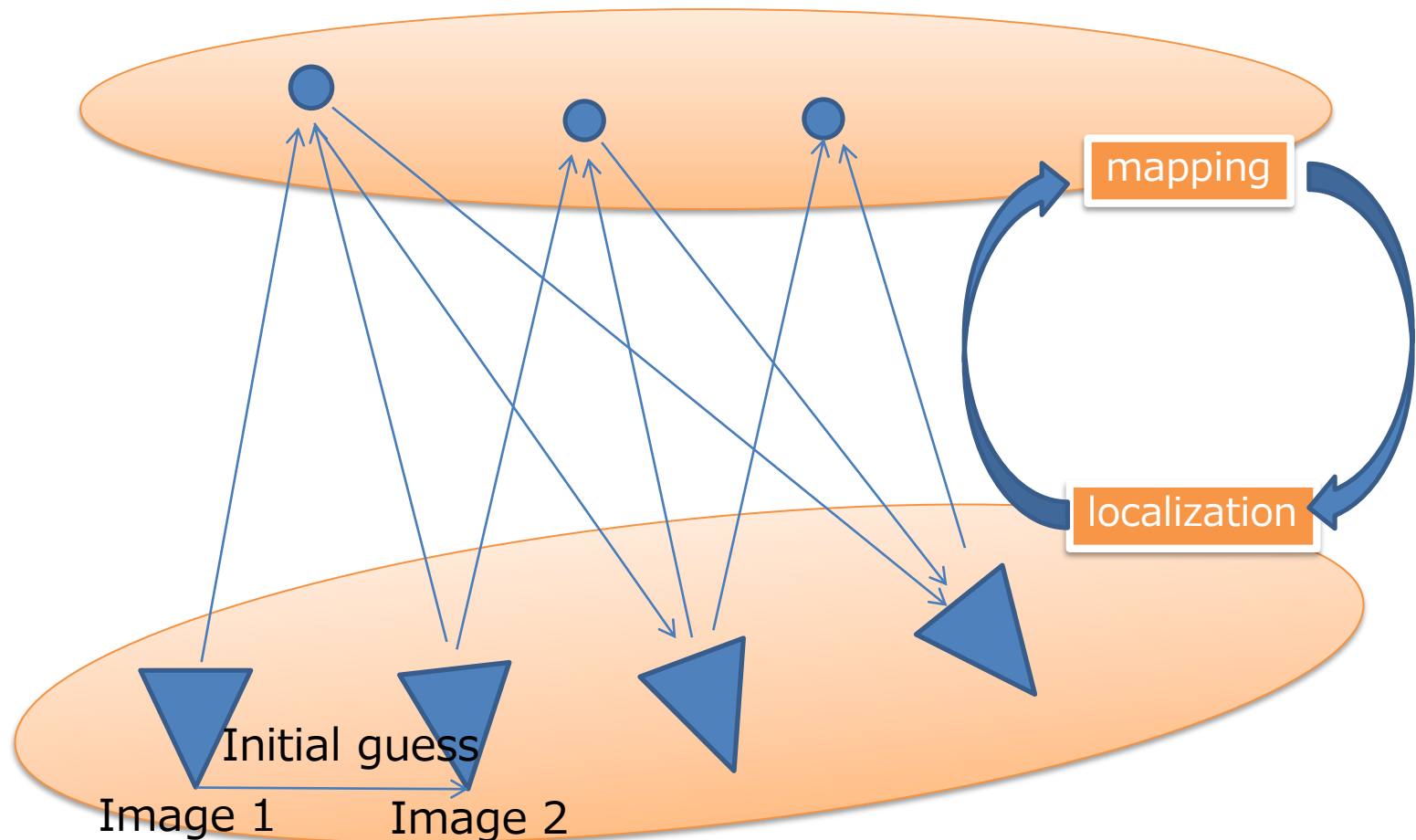
$$\mathbf{M} = [I | \mathbf{0}] \text{ and } \mathbf{M}' = [[\mathbf{e}_x]F | \mathbf{e}']$$

3. For each point correspondence, compute the point  $\mathbf{X}$  in 3D space

**Triangulate** with  $\mathbf{x} = \mathbf{M}\mathbf{X}$  and  $\mathbf{x}' = \mathbf{M}'\mathbf{X}$

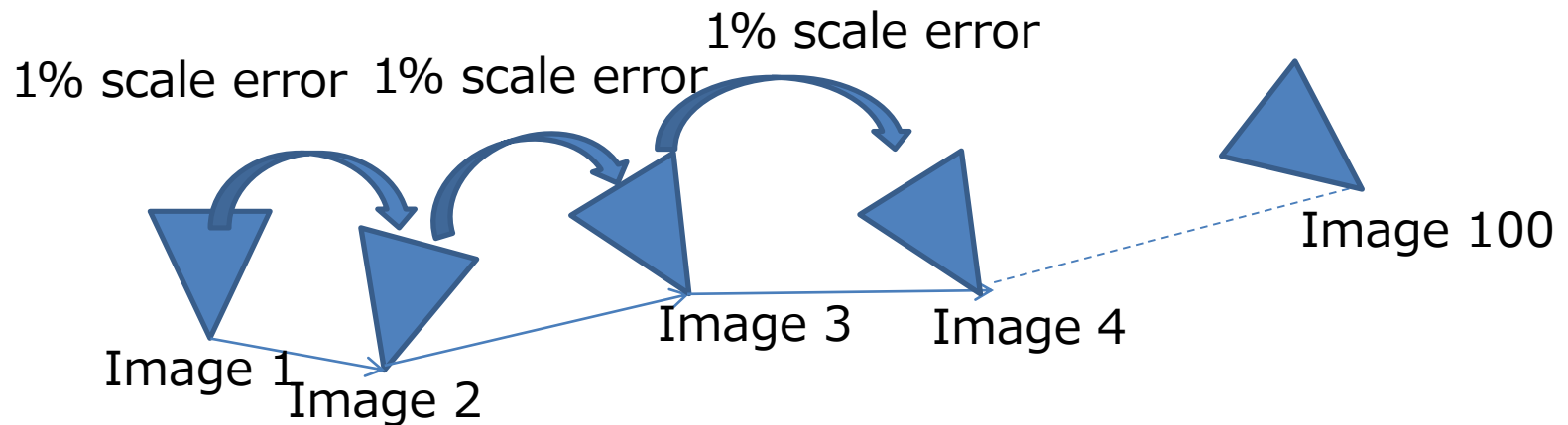


# Idea for sequential SfM



If we know camera poses for a pair of image 1 and 2, we can continue to estimate camera poses and 3-D structure for new input by repeating 'mapping' and 'localization'.

# Problem of accumulative errors by chaining relative poses



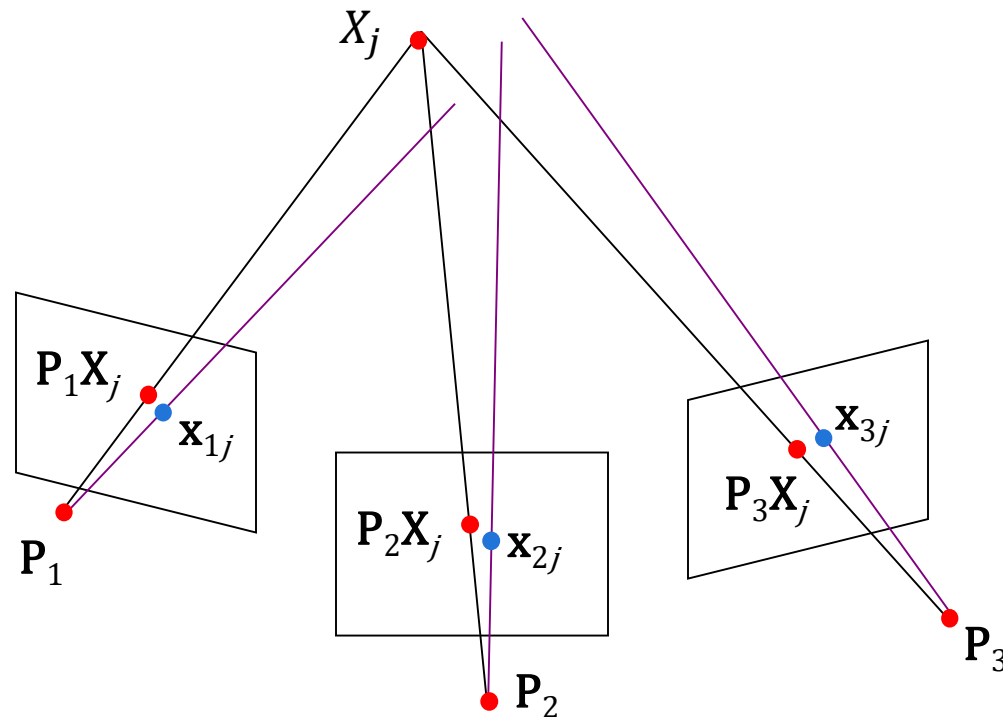
If relative camera poses are estimated with +1% biased scale error for each pair, the scale error at 100 frame will be  $1.01^{100} = 2.70 = 270\%$ .

\*This kind of effect is called as scale drift.

# Bundle adjustment

- Non-linear method for refining structure and motion
- Minimizing reprojection error

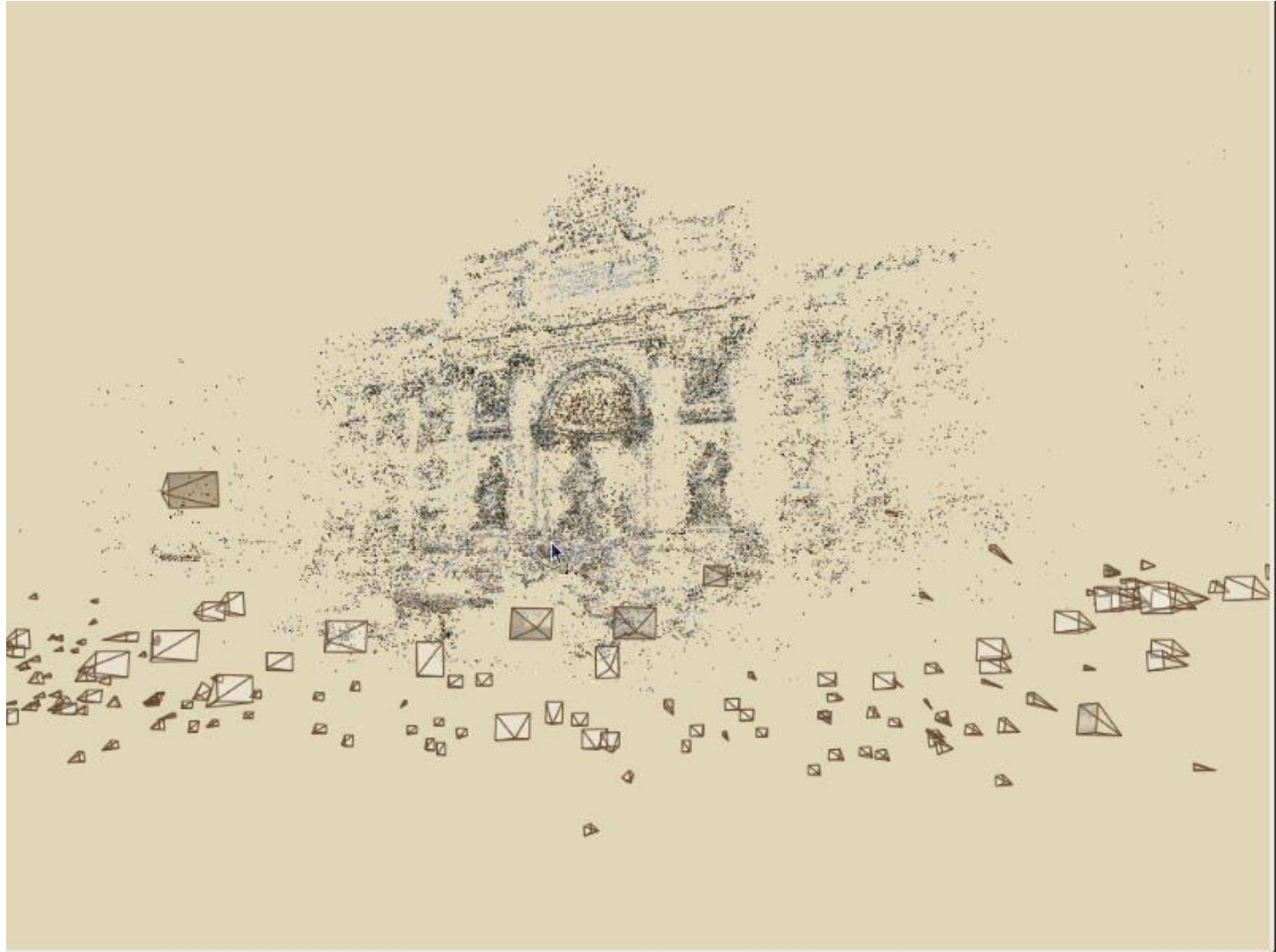
$$E(\mathbf{M}, \mathbf{X}) = \sum_{i=1}^m \sum_{j=1}^n D(\mathbf{x}_{ij}, \mathbf{M}_i \mathbf{X}_j)^2$$



# Large-scale SfM

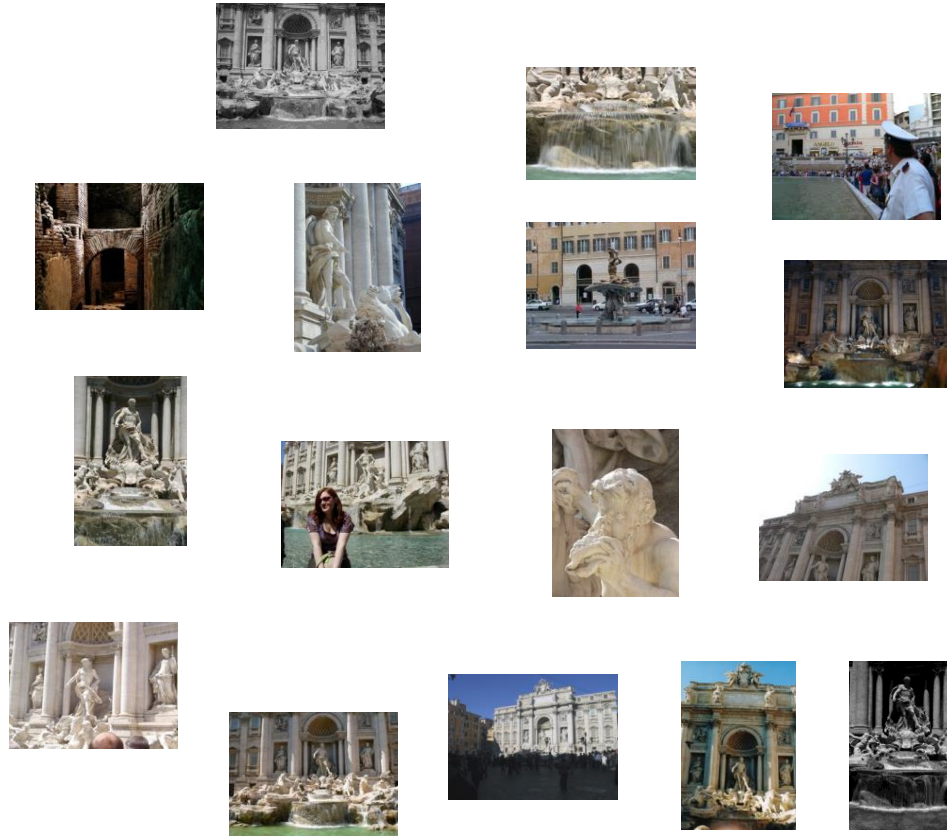
Multi-view, non-ordered images

# Photo Tourism



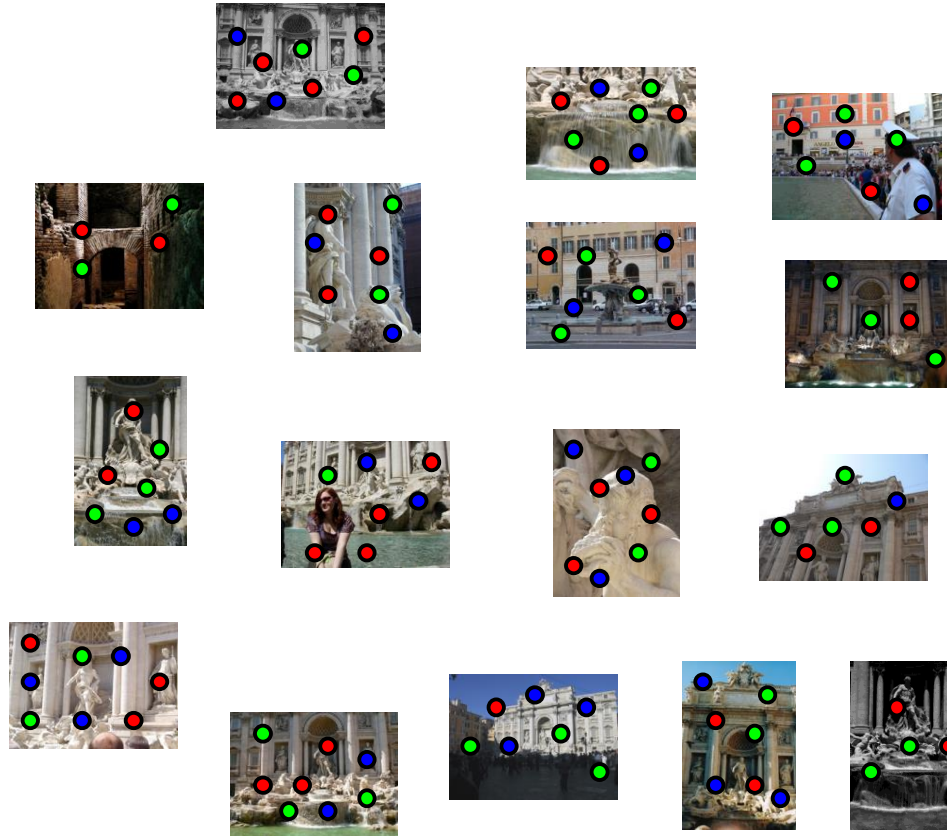
# Feature detection

Detect features using SIFT [Lowe, IJCV 2004]



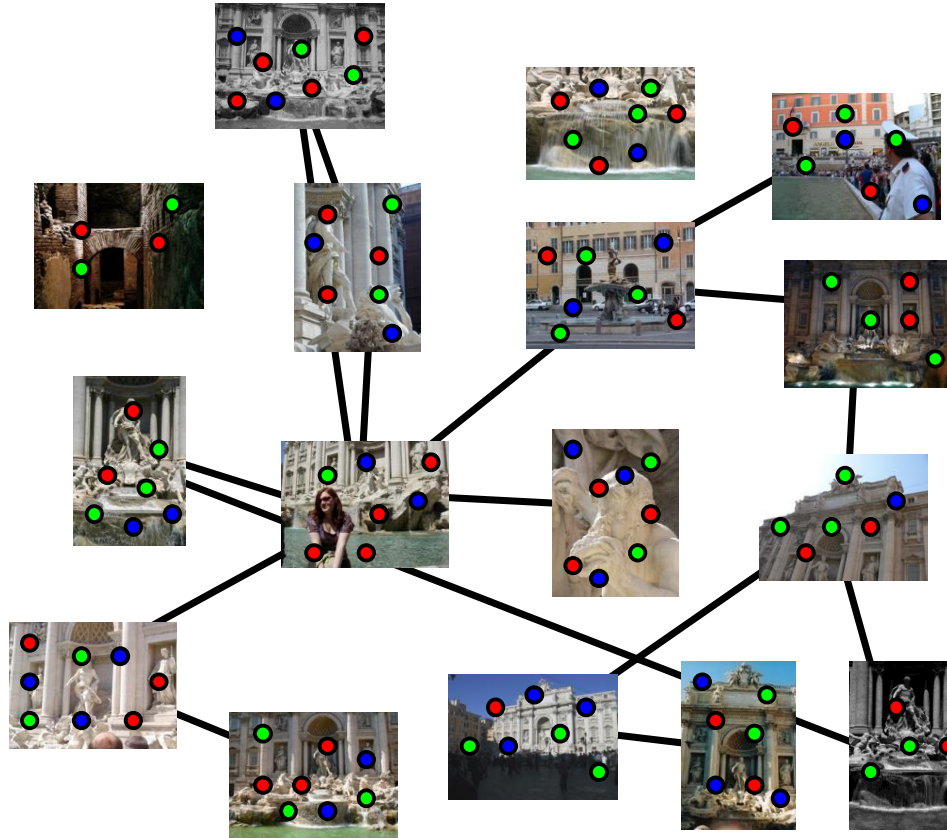
# Feature description

Describe features using SIFT [Lowe, IJCV 2004]



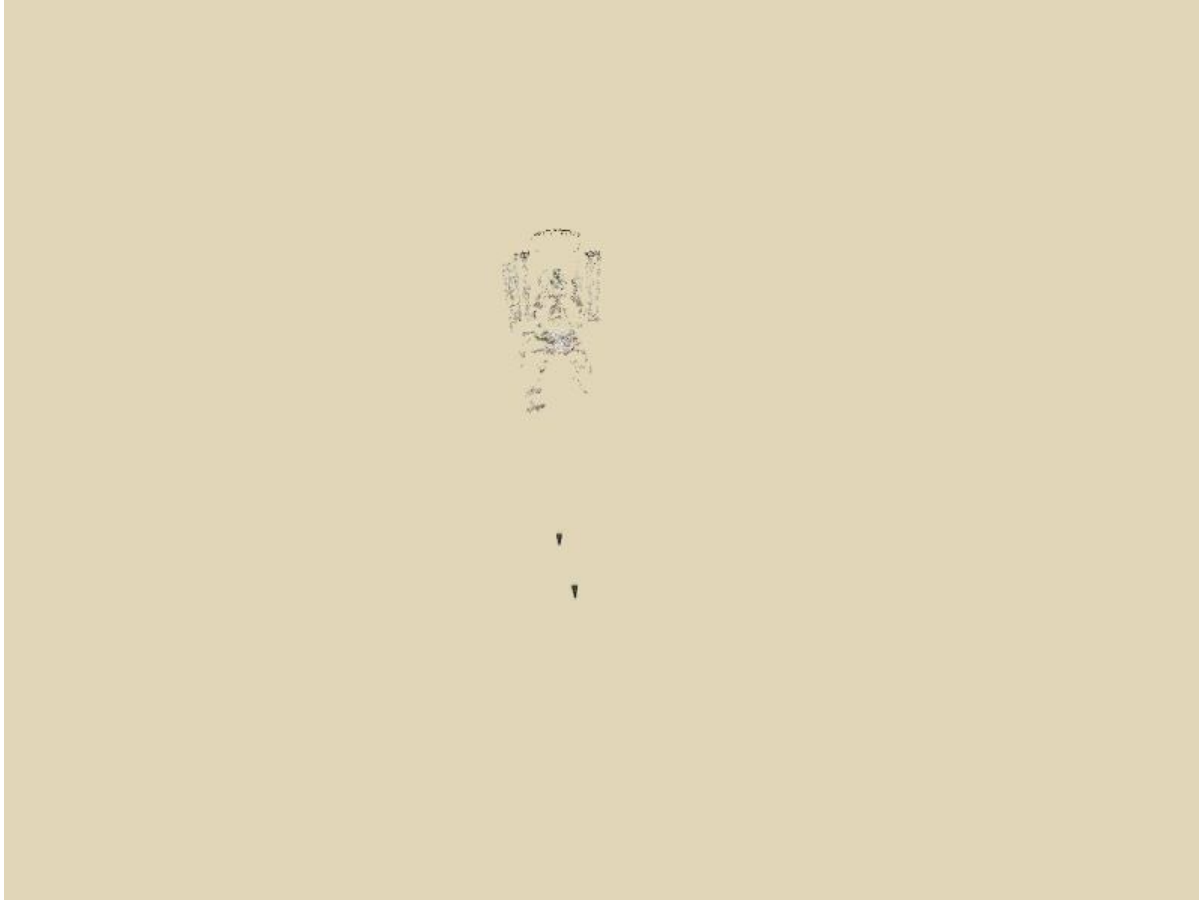
# Feature matching

Refine matching using RANSAC to estimate fundamental matrix between each pair





# Incremental structure from motion



# Final reconstruction

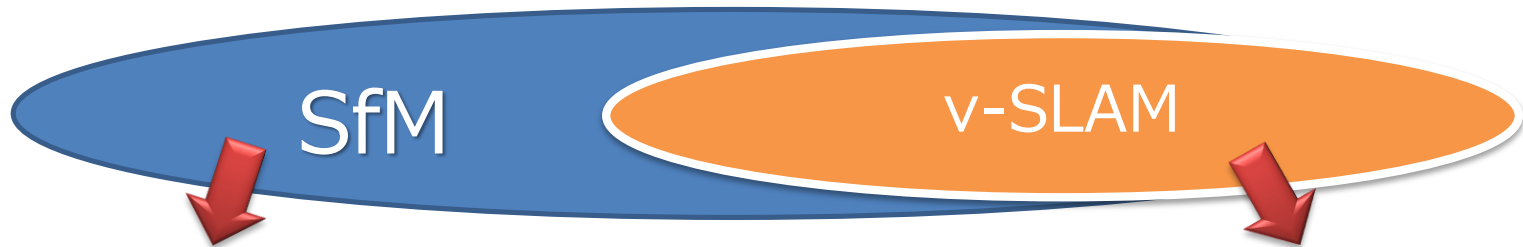




# visual SLAM

Simultaneous Localization and Mapping

# Relationship between SfM and v-SLAM

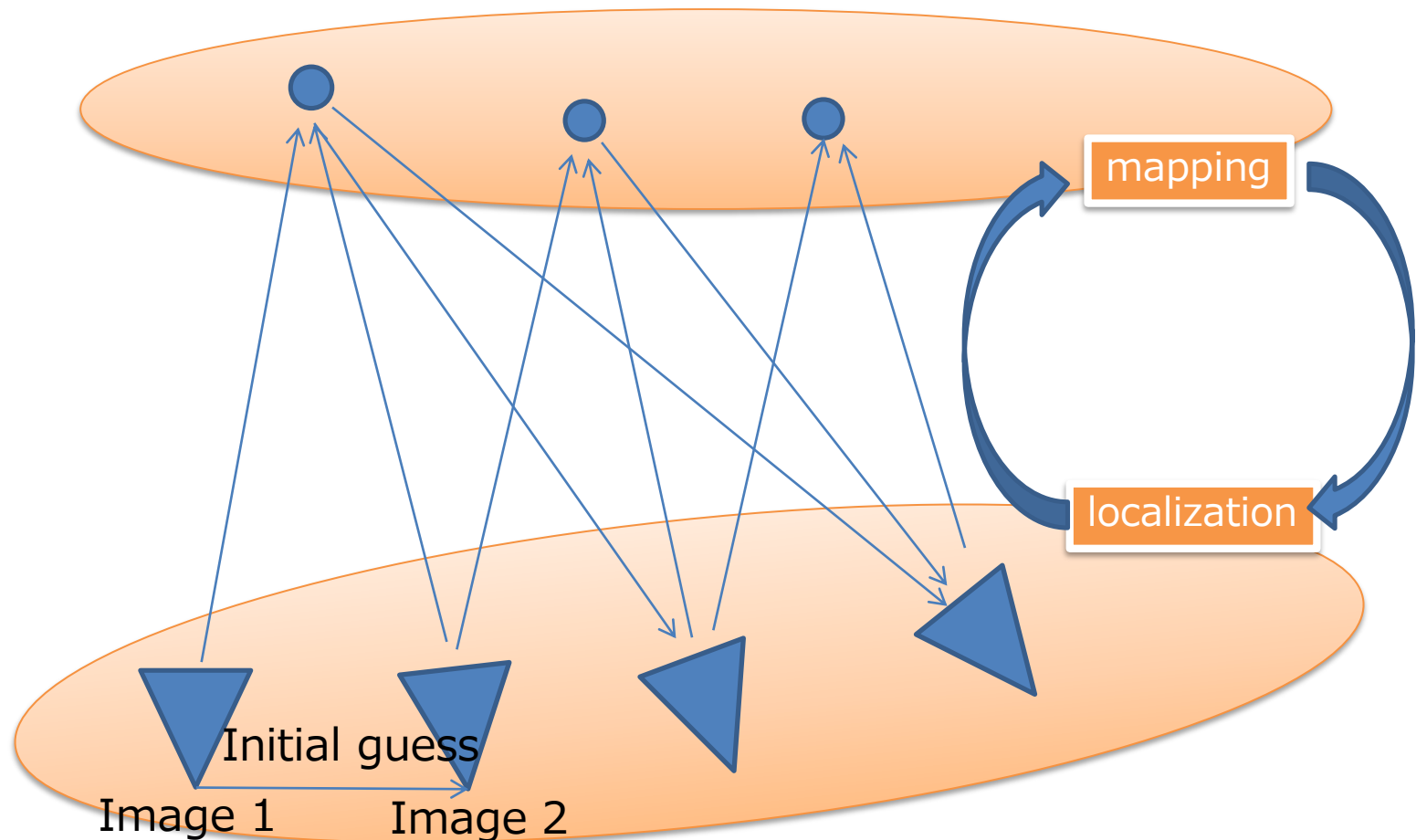


Offline, high accuracy

Real-time, successive output

	SfM	v-SLAM
Input	Images (non-ordered is OK)	Video
Real-time	Not necessary	Required
Output	Batch	Successive
Available information for estimation	All frames	Data acquired before the target frame
Recovery from failure	Easier	Not easy
Feature tracking	Not easy (for non-successive image input)	Easier
Accumulation of errors	Smaller	Bigger / Faster

# Idea for sequential SfM



If we know camera poses for a pair of image 1 and 2, we can continue to estimate camera poses and 3-D structure for new input by repeating 'mapping' and 'localization'.

# Basic pipeline of sequential visual SLAM

Camera pose initialization by Two-view SfM

Iterative process

Feature point tracking

Camera pose estimation

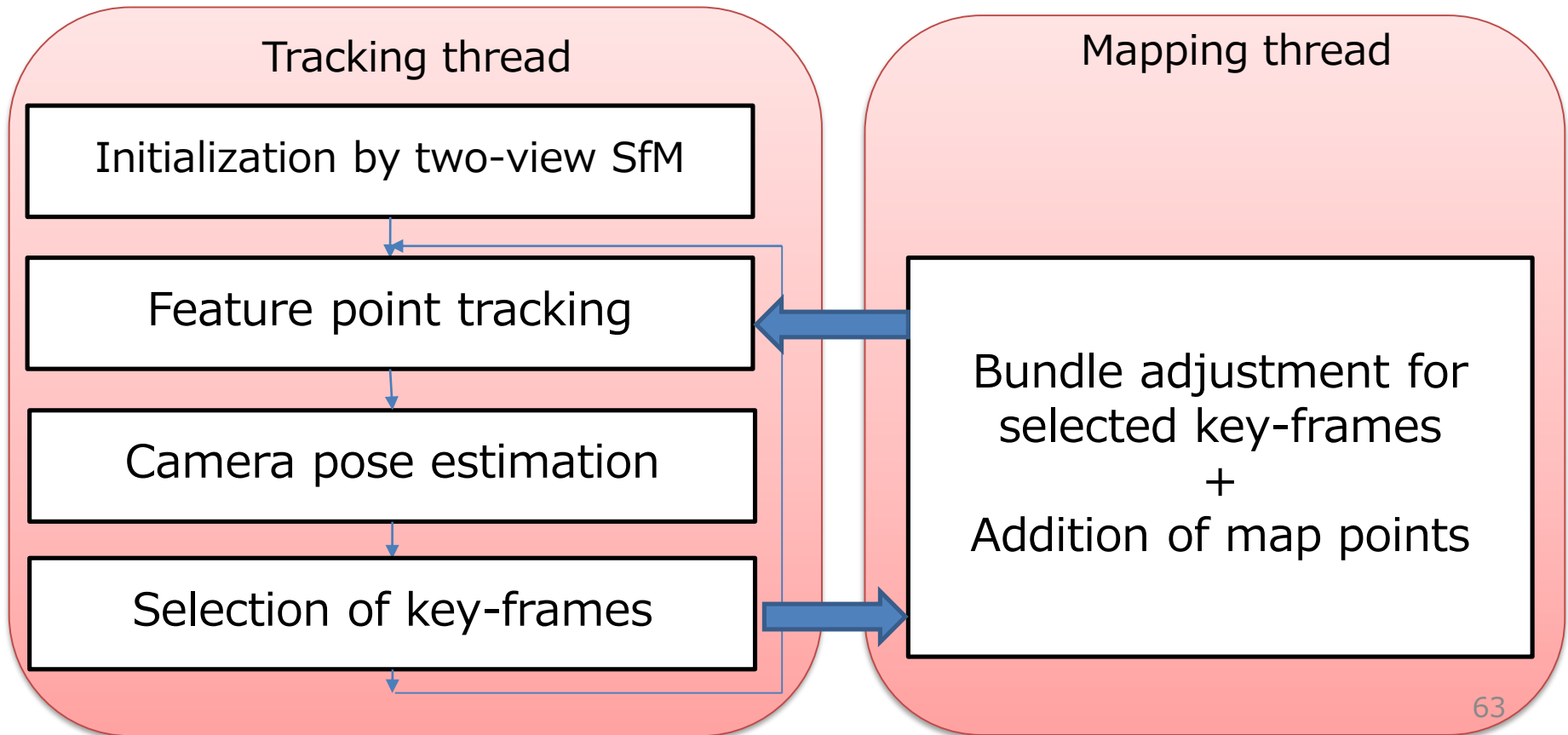
Estimation of 3D points

Addition and deletion of feature points

(option) Local / Global bundle adjustment

# Parallel Tracking and Mapping

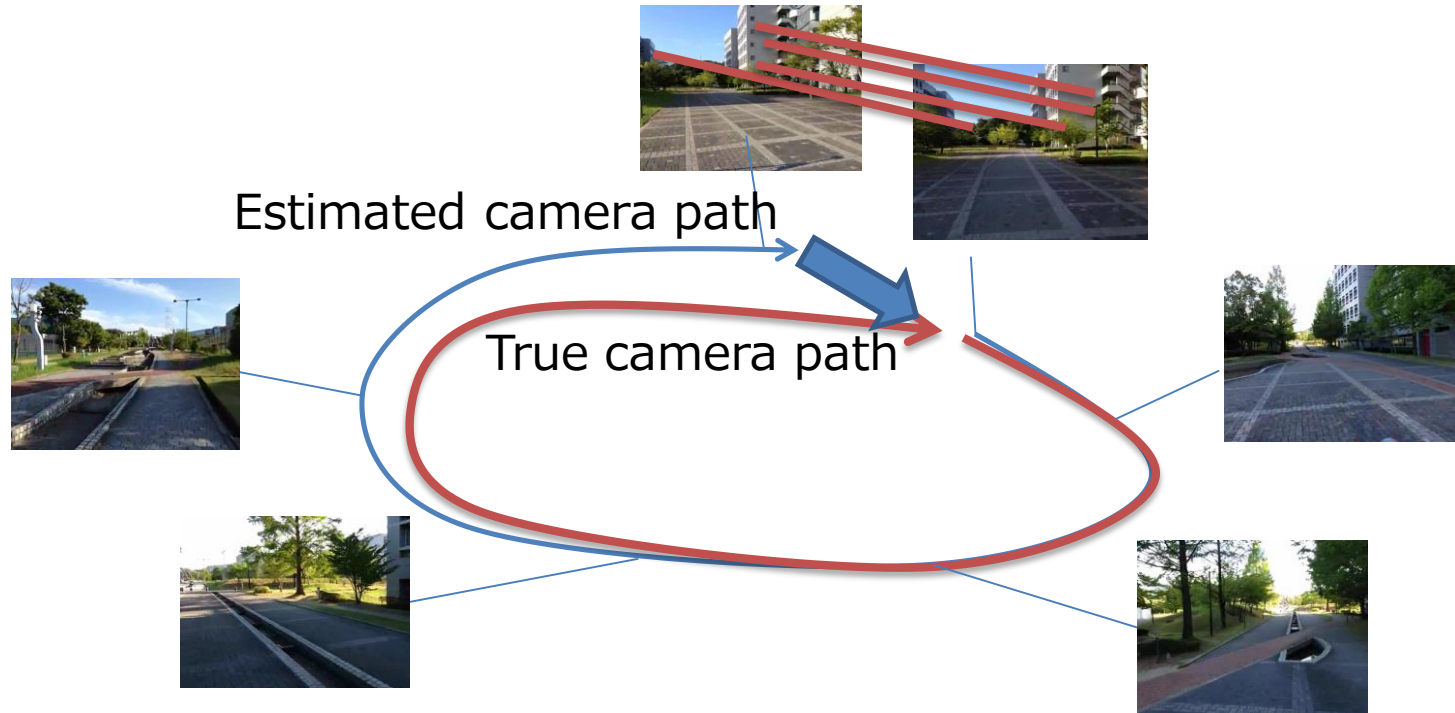
Local bundle adjustment is asynchronously processed for minimizing accumulation of errors using selected key-frames in order not to prevent real-time processing.





# Reduction of accumulation errors / re-start from tracking failure

## Loop closing, re-localization



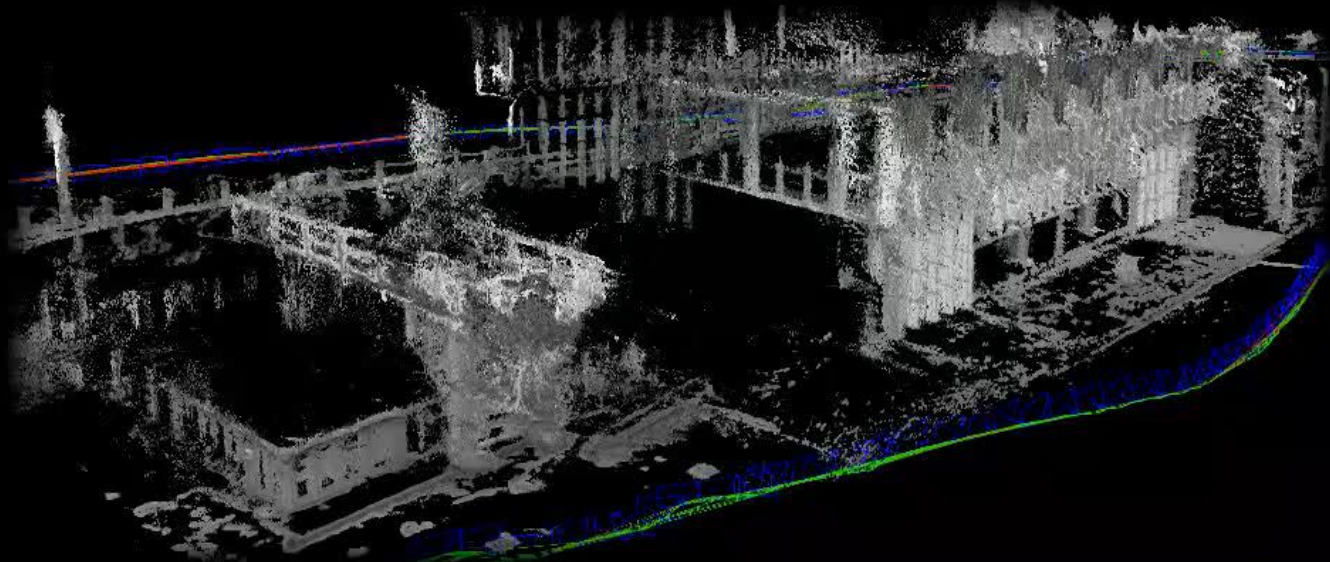
### Basic flow

1. Find similar image of current input from already observed images
2. Find corresponding points between current and selected images.
3. Optimize data using bundle adjustment / Estimate camera pose.

# LSD-SLAM: Large Scale Direct Monocular SLAM

## LSD-SLAM: Large-Scale Direct Monocular SLAM

Jakob Engel, Thomas Schöps, Daniel Cremers  
**ECCV 2014, Zurich**



Computer Vision Group  
Department of Computer Science  
Technical University of Munich



# Final report

- Explain the reason of these strange views

