

An Invitation ~~Introduction~~ to 3D Vision : A Tutorial for Everyone

Sunglok Choi, Senior Researcher
sunglok@etri.re.kr | <http://sites.google.com/site/sunglok>
Electronics and Telecommunication Research Institute (ETRI)

What is Computer Vision?

Computer vision is an interdisciplinary field that deals with how computers can be made to gain high-level understanding from digital images or videos. From the perspective of engineering, it seeks to automate tasks that the human visual system can do.[1][2][3]

"Computer vision is concerned with the automatic extraction, analysis and understanding of useful information from a single image or a sequence of images. It involves the development of a theoretical and algorithmic basis to achieve automatic visual understanding." [9] As a scientific discipline, computer vision is concerned with the theory behind artificial systems that extract information from images. The image data can take many forms, such as video sequences, views from multiple cameras, or multi-dimensional data from a medical scanner.[10] As a technological discipline, computer vision seeks to apply its theories and models for the construction of computer vision systems.

Computer Vision by [Wikipedia](#)



What is Computer Vision?

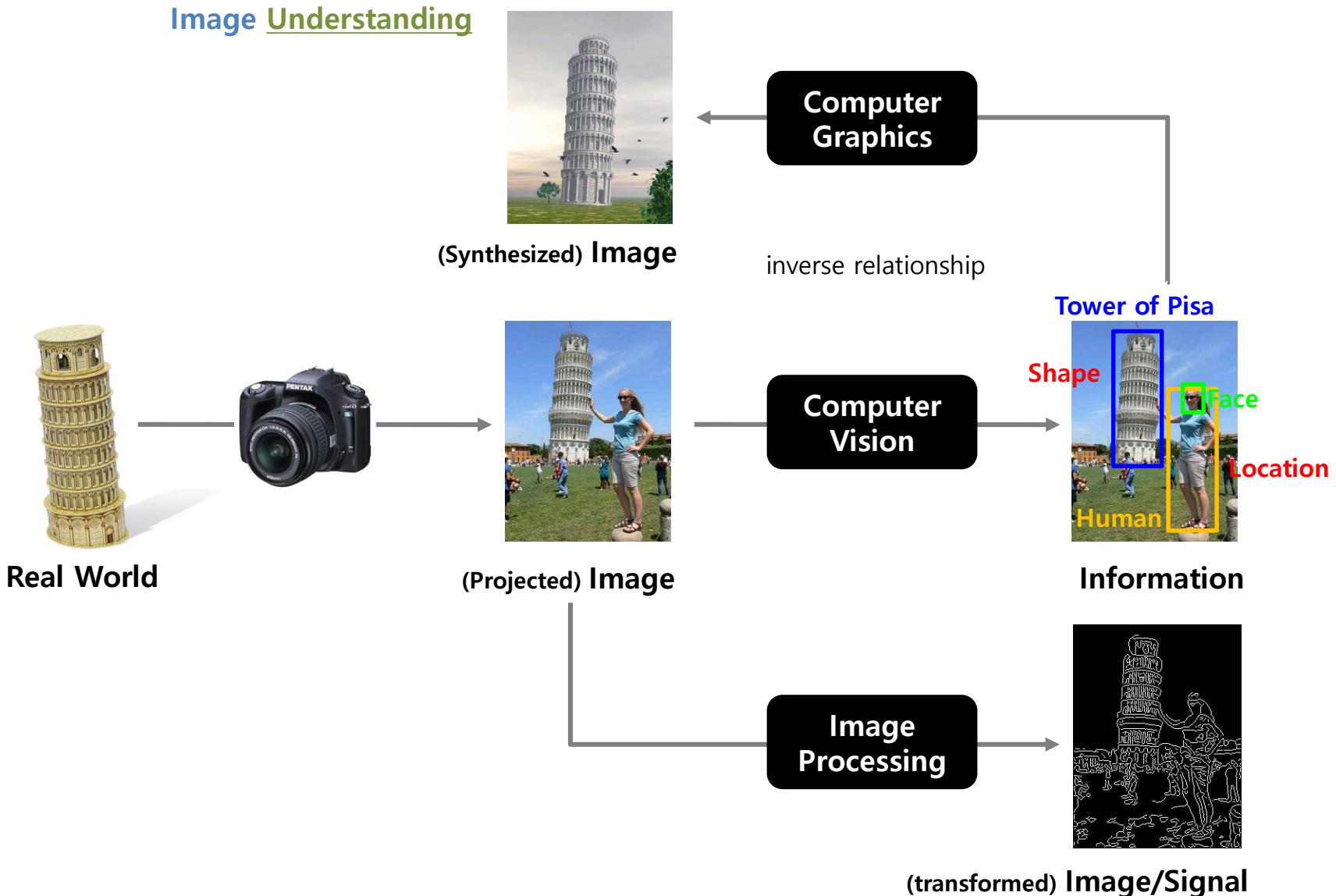
Computer vision is an interdisciplinary field that deals with how computers can be made to gain high-level understanding from digital images or videos. From the perspective of engineering, it seeks to automate tasks that the human visual system can do.[1][2][3]

"Computer vision is concerned with the automatic extraction, analysis and understanding of useful information from a single image or a sequence of images. It involves the development of a theoretical and algorithmic basis to achieve automatic visual understanding."^[9] As a scientific discipline, computer vision is concerned with the theory behind artificial systems that extract information from images. The image data can take many forms, such as video sequences, views from multiple cameras, or multi-dimensional data from a medical scanner.^[10] As a technological discipline, computer vision seeks to apply its theories and models for the construction of computer vision systems.

Computer Vision by [Wikipedia](#)



What is Computer Vision?



What is Computer Vision?

Computer Vision

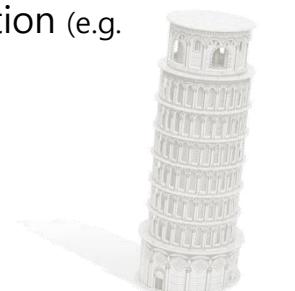
What is it?

- Label (e.g. Tower of Pisa)
- Shape (e.g.)



Where am I?

- Place (e.g. Piazza del Duomo, Pisa, Italy)
- Location (e.g.)



(84, 10, 18) [m]



What is 3D Vision?

Visual Geometry (Multiple View Geometry)

Geometric Vision

Computer Vision

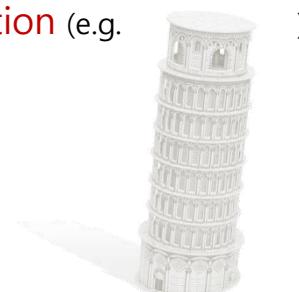
What is it?

- **Label** (e.g. Tower of Pisa)
- **Shape** (e.g.)



Where am I?

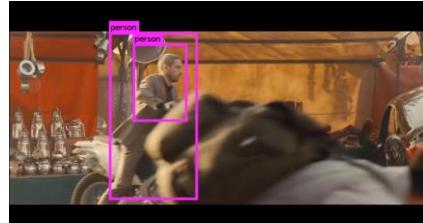
- **Place** (e.g. Piazza del Duomo, Pisa, Italy)
- **Location** (e.g.)



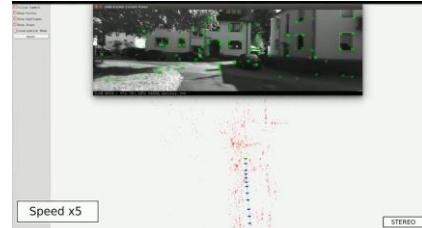
(84, 10, 18) [m]



Recognition Problems v.s. Reconstruction Problems



YOLO v2 (2016)

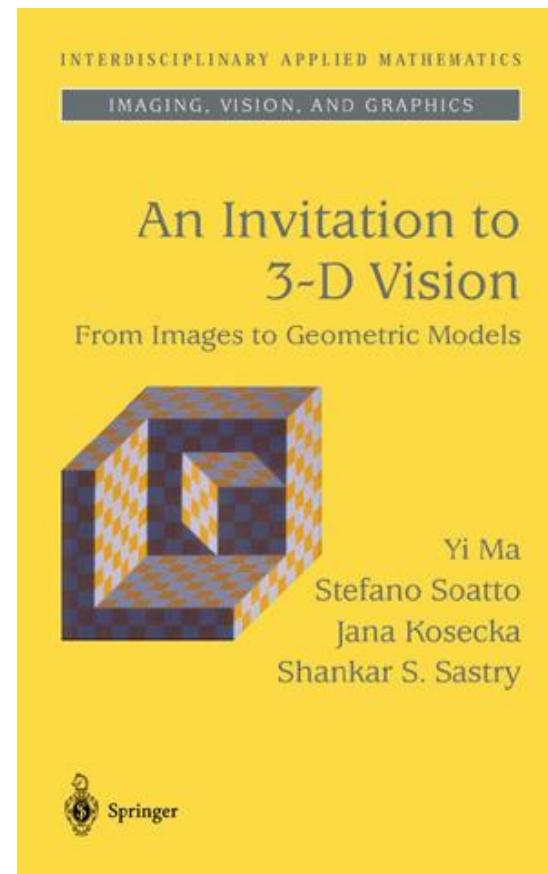
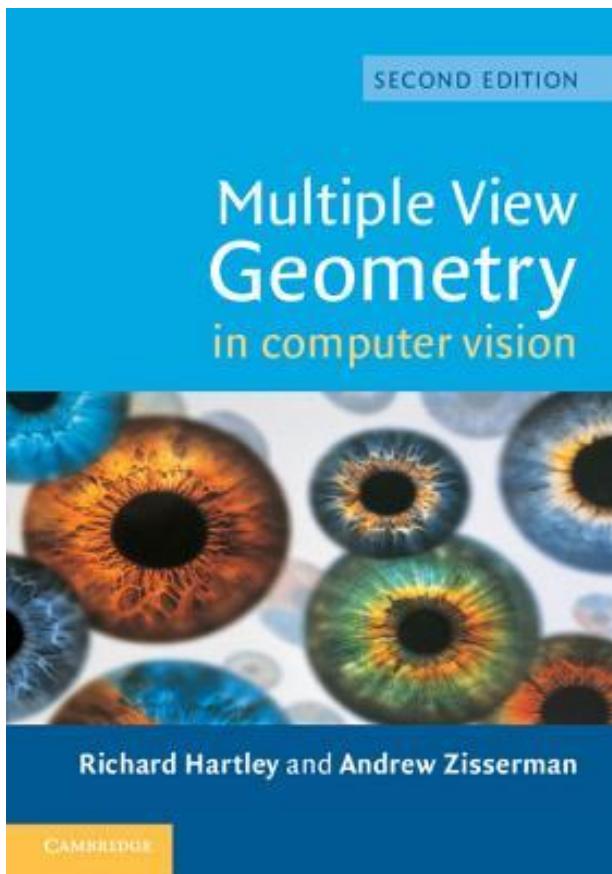


ORB-SLAM2 (2016)



What is 3D Vision?

- Reference Books



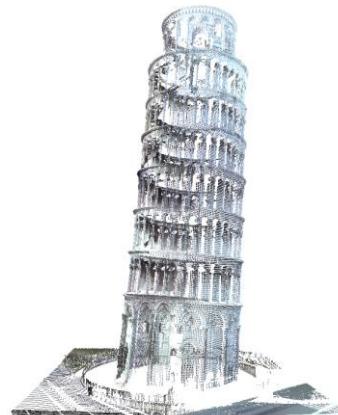
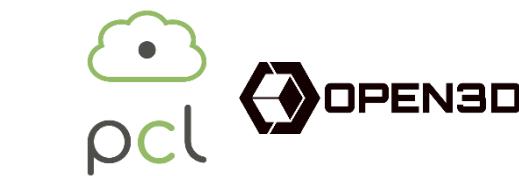


What is 3D Vision?

- OpenCV (Open Source Computer Vision)
 - calib3d Module: Camera Calibration and 3D Reconstruction
 - OpenCV API Reference
 - Introduction
 - core. The Core Functionality
 - imgproc. Image Processing
 - imgcodecs. Image file reading and writing
 - videoio. Media I/O
 - highgui. High-level GUI and Media I/O
 - video. Video Analysis
 - calib3d. Camera Calibration and 3D Reconstruction
 - features2d. 2D Features Framework
 - objdetect. Object Detection
 - ml. Machine Learning
 - flann. Clustering and Search in Multi-Dimensional Spaces
 - photo. Computational Photography
 - stitching. Images stitching
 - cuda. CUDA-accelerated Computer Vision
 - cudaarithm. CUDA-accelerated Operations on Matrices
 - cudabgsegm. CUDA-accelerated Background Segmentation
 - cudacodec. CUDA-accelerated Video Encoding/Decoding
 - cudafeatures2d. CUDA-accelerated Feature Detection and Description
 - cudafilters. CUDA-accelerated Image Filtering
 - cudaimproc. CUDA-accelerated Image Processing
 - cudaoptflow. CUDA-accelerated Optical Flow
 - cudastereo. CUDA-accelerated Stereo Correspondence
 - cudawarping. CUDA-accelerated Image Warping
 - shape. Shape Distance and Matching
 - superres. Super Resolution
 - videotab. Video Stabilization
 - viz. 3D Visualizer

cf. All examples in this tutorial are mostly **less than 100 lines** and based on **recent OpenCV (> 3.0.0)**.

What is 3D Vision?



depth image, range data, point cloud, polygon mesh, ...



Perspective Camera

3D Vision

v.s.

3D Data Processing



RGB-D Camera
(Stereo, Structured Light, ToF, Light Field)



Omni-directional Camera



Range Sensor
(LiDAR, RADAR)

What is 3D Vision?



image



Perspective Camera

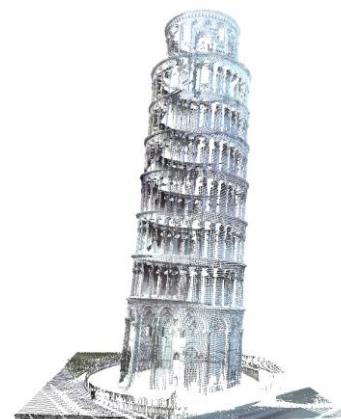


Omni-directional Camera

3D Vision

v.s.

3D Data Processing



depth image, range data, point cloud, polygon mesh, ...



RGB-D Camera
(Stereo, Structured Light, ToF, Light Field)



Range Sensor
(LiDAR, RADAR)

Applications: Photo Browsing

- [Photo Tourism](#) (2006)



Applications: 3D Reconstruction

- [Building Rome in a Day](#) (2009)

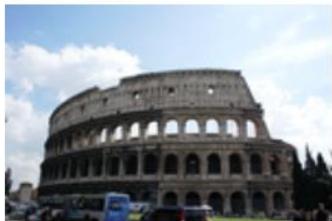


photo by [Ceci And Brandon](#)



2106 images, 819242 points
[click here to see more views](#)



photo by [TORISFERICK](#)



1936 images, 656699 points
[click here to see more views](#)



photo by [act2win](#)



1815 images, 422593 points
[click here to see more views](#)



photo by [angle_in_soul](#)



1381 images, 499044 points
[click here to see more views](#)



Applications: 3D Reconstruction from Cellular Phones

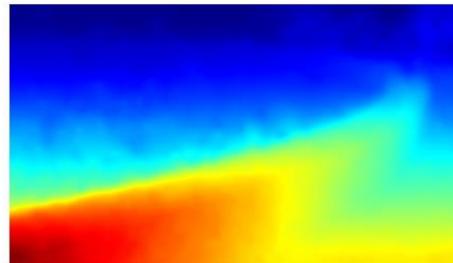
- [Structure from Small Motion](#) (SfSM; 2015)



(a) Reference images



(b) SfSM results



(c) Depth maps

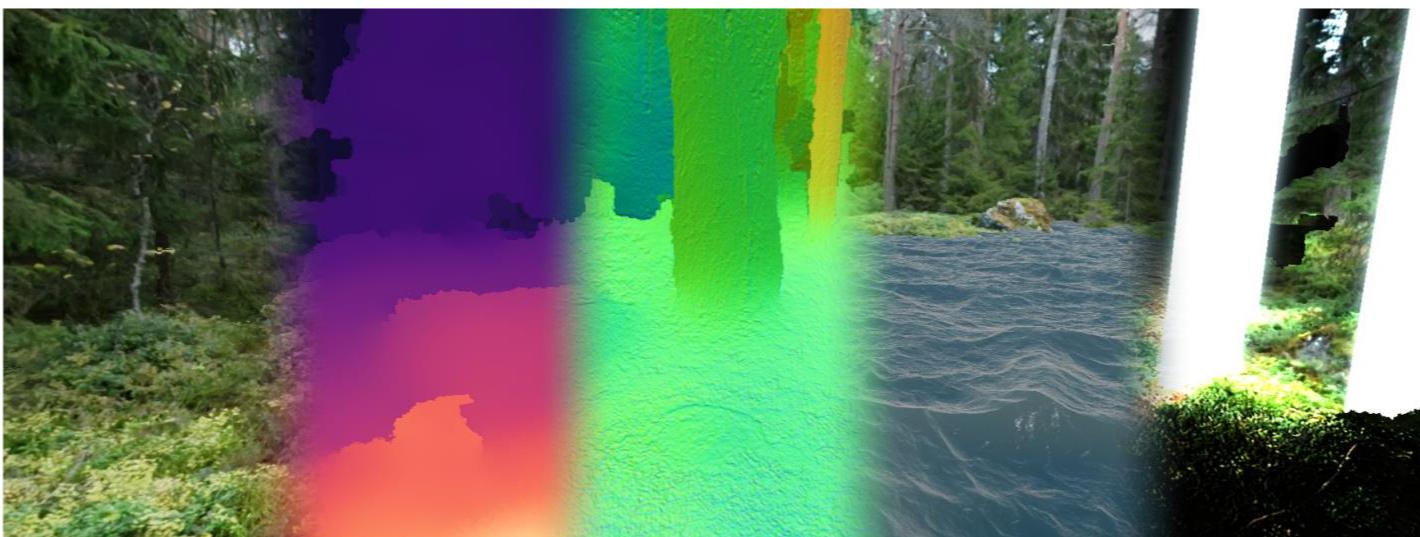


(d) Our 3D meshes

- [Casual 3D Photography](#) (2017)



Casual 3D photo capture



Color

Depth
Reconstruction

Normal map

Geometry-aware
Example Effects

Lighting

[Reference] Im et al., **High Quality Structure from Small Motion for Rolling Shutter Cameras**, ICCV, 2015

[Reference] Hedman et al., **Causal 3D Photography**, SIGGRAPH Asia, 2017

Applications: Real-time Visual SLAM

- [ORB-SLAM](#) (2014)



Instituto Universitario de Investigación
en Ingeniería de Aragón
Universidad Zaragoza

ORB-SLAM2: an Open-Source SLAM System
for Monocular, Stereo and RGB-D Cameras

Raúl Mur-Artal and Juan D. Tardós

raulmur@unizar.es

tardos@unizar.es

Applications: Augmented Reality

- [PTAM: Parallel Tracking and Mapping \(2007\)](#)

Parallel Tracking and Mapping
for Small AR Workspaces

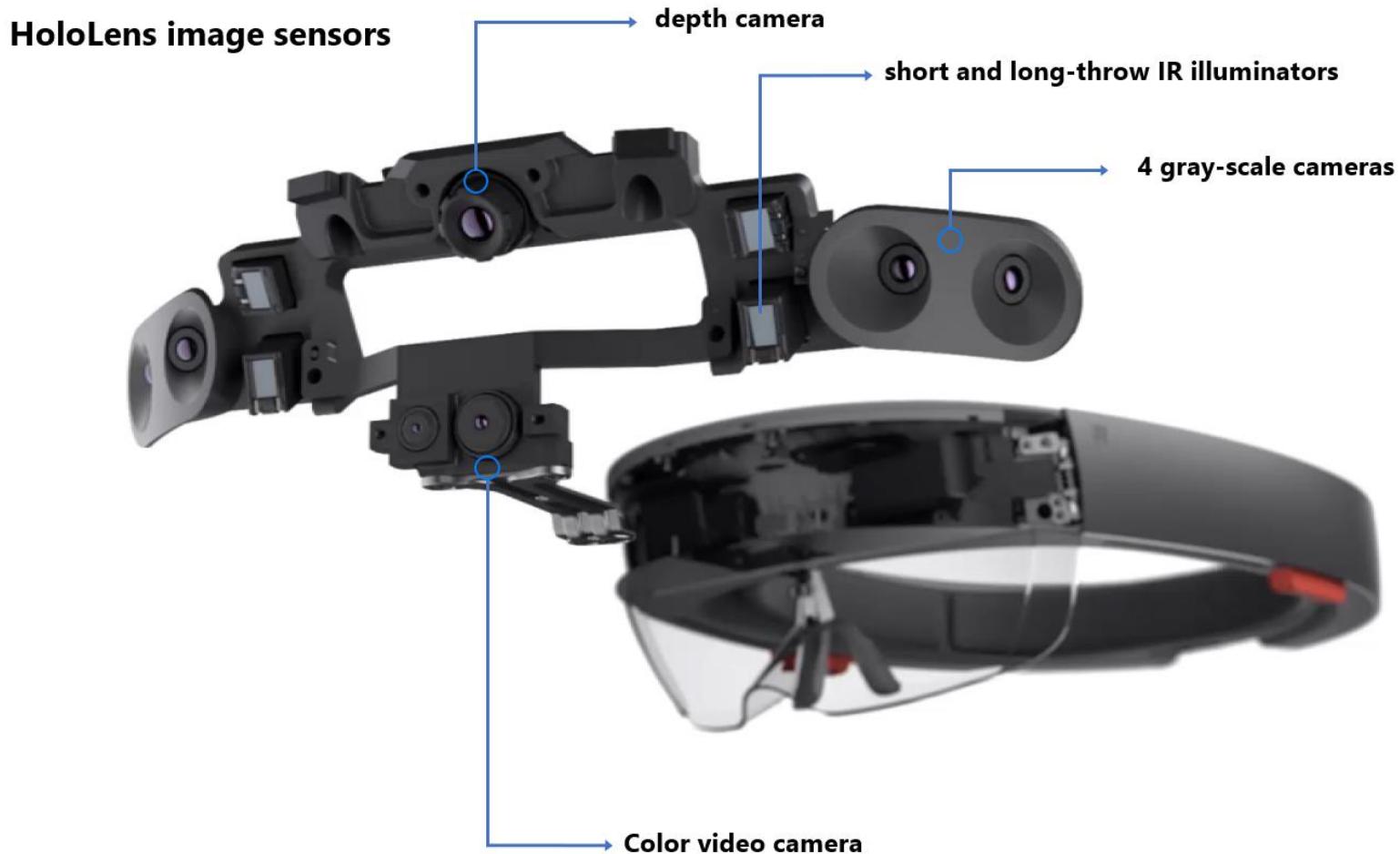
ISMAR 2007 video results

Georg Klein and David Murray
Active Vision Laboratory
University of Oxford



Applications: Mixed Reality

- [Microsoft Hololens 2](#) (2019)
 - Head tracking: 4 x visible light cameras



Applications: Virtual Reality

- [Oculus Quest](#) (2019)

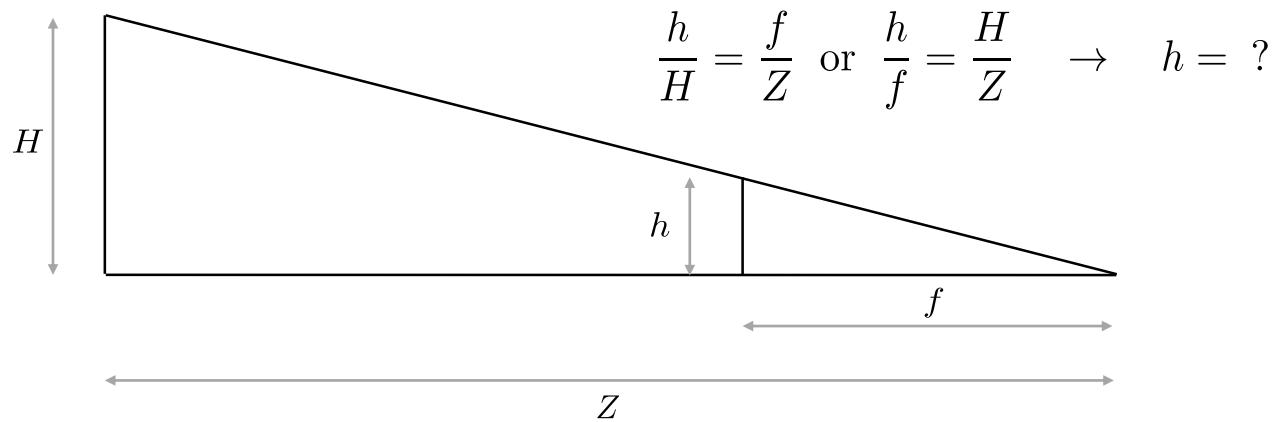


Table of Contents

- What is 3D Vision?
- Single-view Geometry
 - Camera Projection Model
 - General 2D-3D Geometry
- Two-view Geometry
 - Planar 2D-2D Geometry (**Projective Geometry**)
 - General 2D-2D Geometry (**Epipolar Geometry**)
- Correspondence Problem
- Multi-view Geometry
- Summary

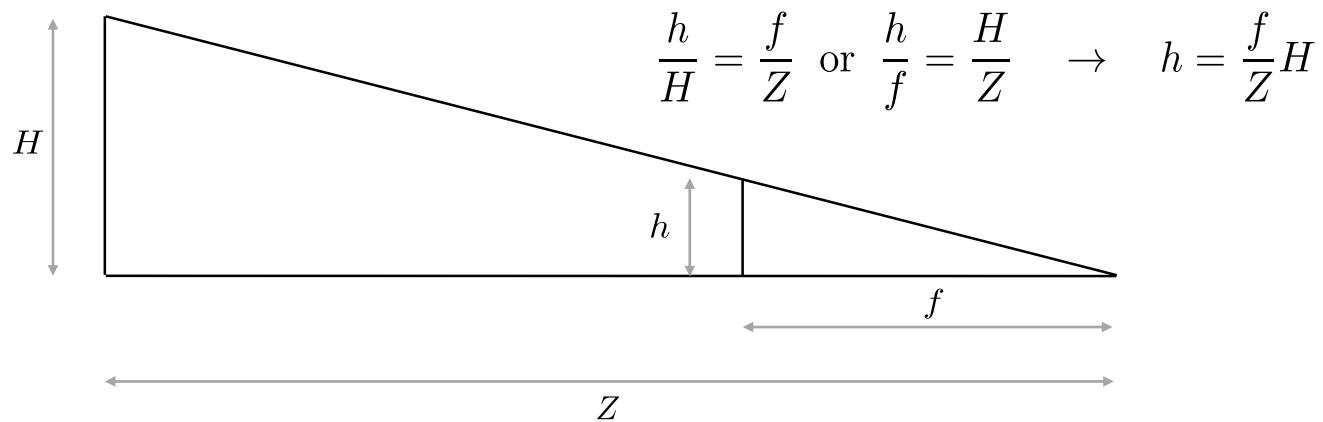
Getting Started with 2D

- Similarity



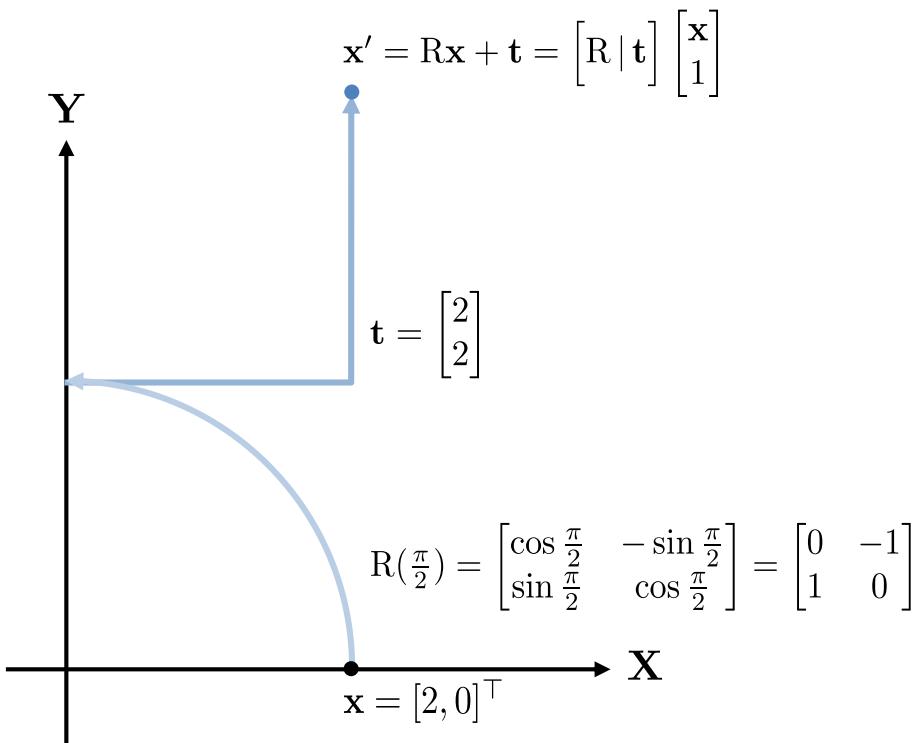
Getting Started with 2D

- Similarity



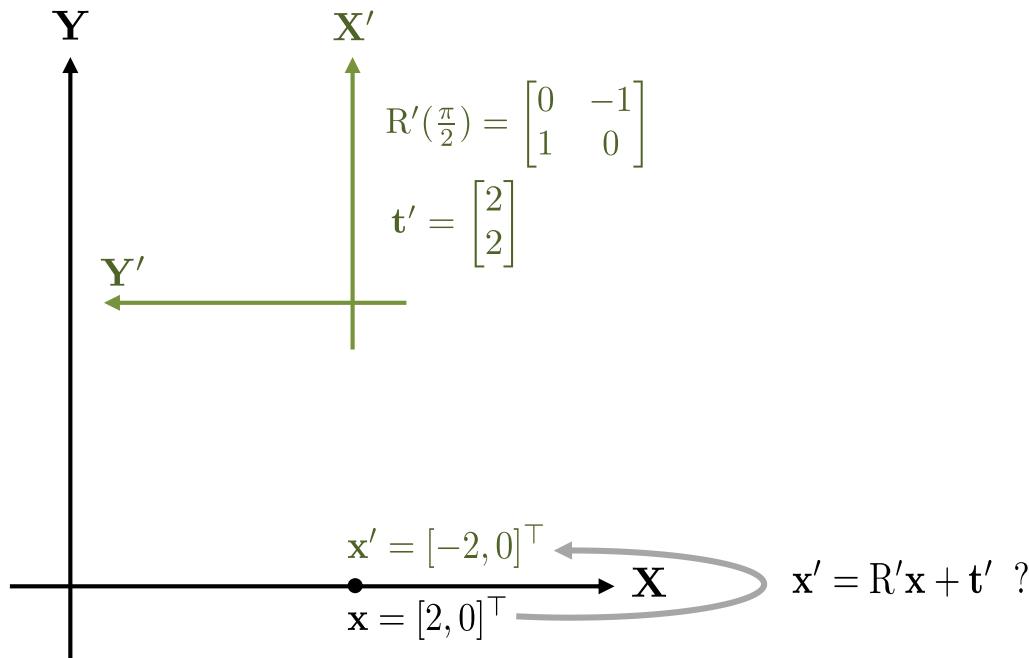
Getting Started with 2D

- Point transformation



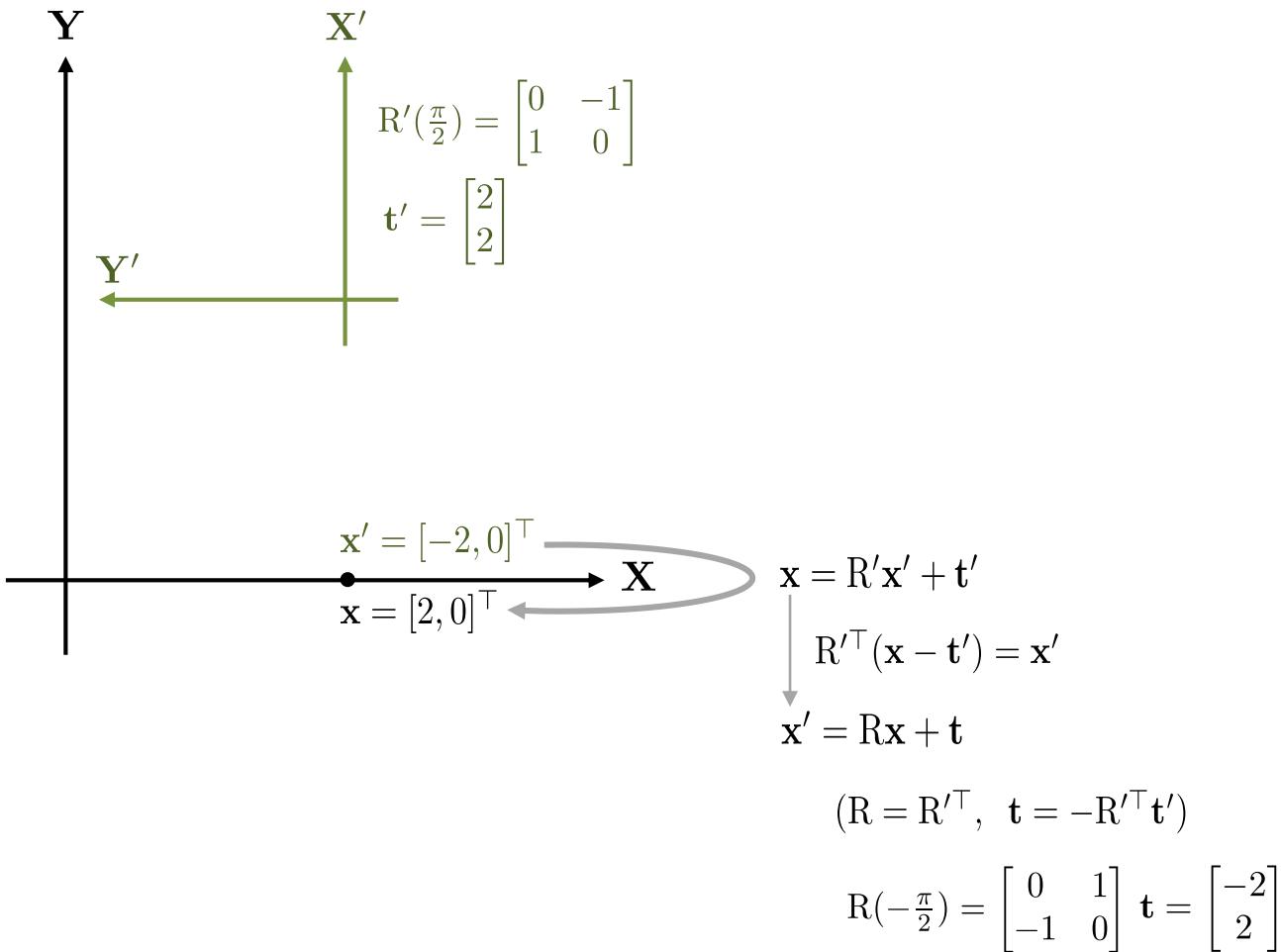
Getting Started with 2D

- Coordinate transformation



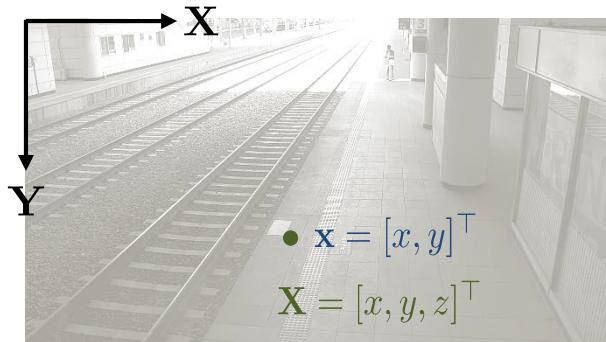
Getting Started with 2D

- Coordinate transformation is the **inverse** of point transformation.

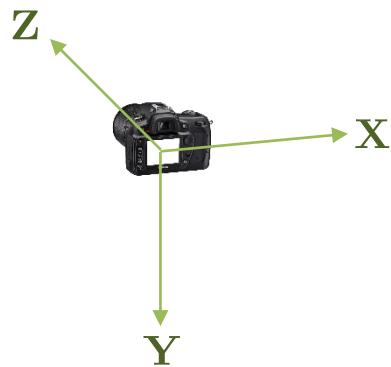


Getting Started with 2D

- Image coordinate (unit: [pixel])



- Camera coordinate (unit: [meter])



Getting Started with 2D

- 2D rotation matrix

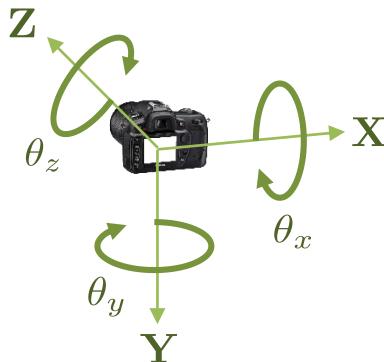


$$R(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

cf. Properties of a rotation matrix

- $R^{-1} = R^T$ (orthogonal matrix)
- $\det(R) = 1$

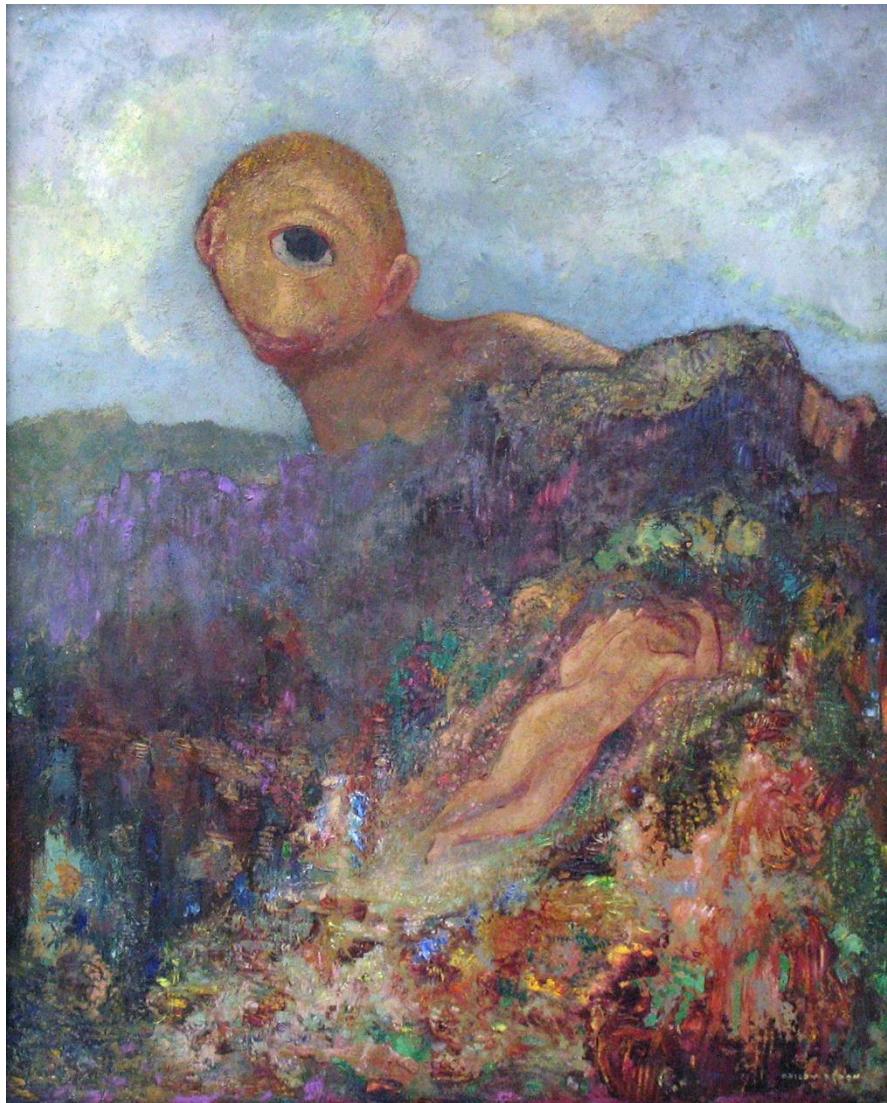
- 3D rotation matrix



	Cameras	Vehicles	Airplanes	Telescopes
θ_x	Tilt	Pitch	Attitude	Elevation
θ_y	Pan	Yaw	Heading	Azimuth
θ_z	Roll	Roll	Bank	Horizon

$$R(\theta_x, \theta_y, \theta_z) = R_z(\theta_z)R_y(\theta_y)R_x(\theta_x) = \begin{bmatrix} \cos \theta_z & -\sin \theta_z & 0 \\ \sin \theta_z & \cos \theta_z & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \theta_y & 0 & \sin \theta_y \\ 0 & 1 & 0 \\ -\sin \theta_y & 0 & \cos \theta_y \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \theta_x & -\sin \theta_x \\ 0 & \sin \theta_x & \cos \theta_x \end{bmatrix}$$

Single-view Geometry



Camera Projection Model

- Pinhole Camera Model
- Geometric Distortion

General 2D-3D Geometry

- Camera Calibration
- Absolute Camera Pose Estimation

The Cyclops, gouache and oil by Odilon Redon

Camera Projection Model

- Pinhole Camera Model



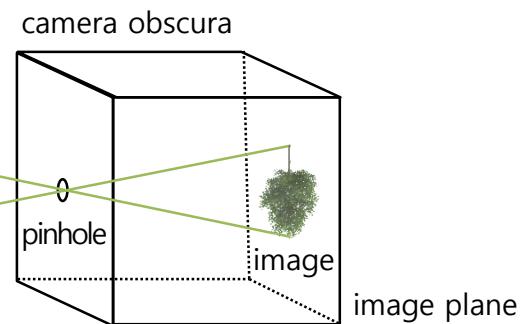
A large-scale camera obscura
at San Francisco, California



A modern-day camera obscura



An Image in camera obscura
at Portslade, England

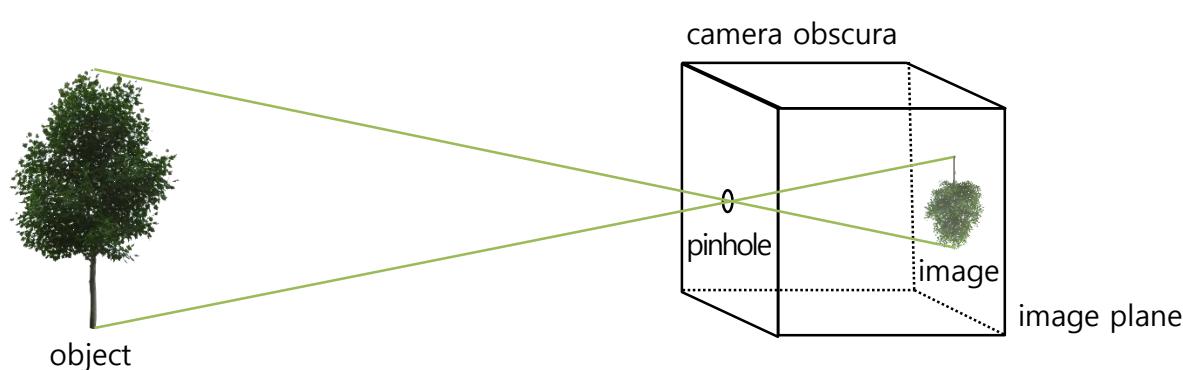
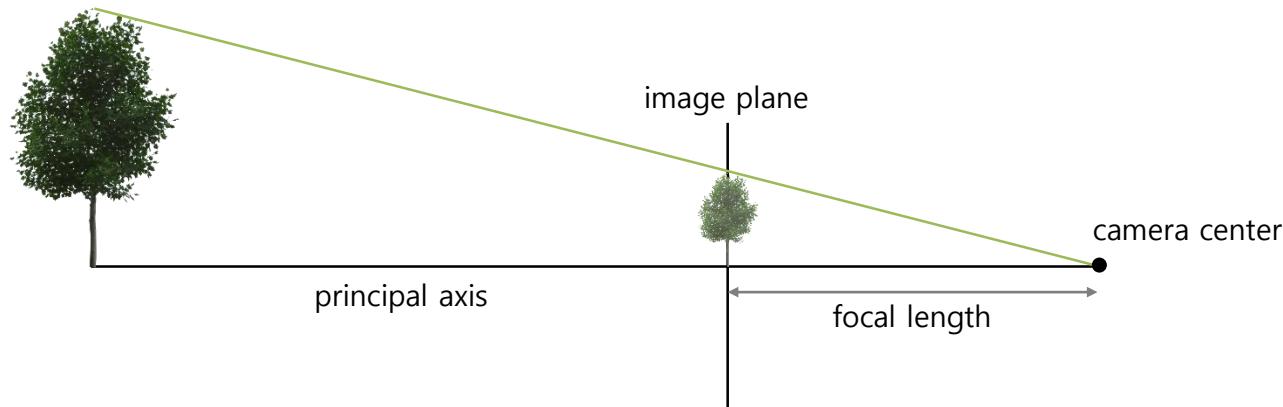


Camera Projection Model

- Pinhole Camera Model

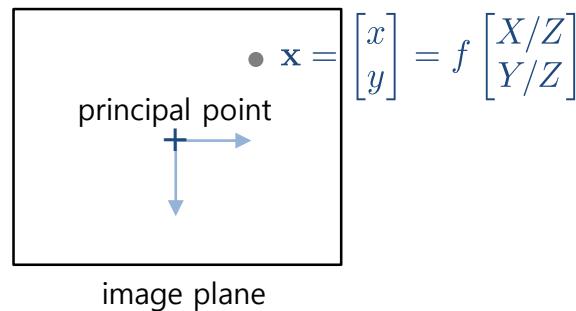
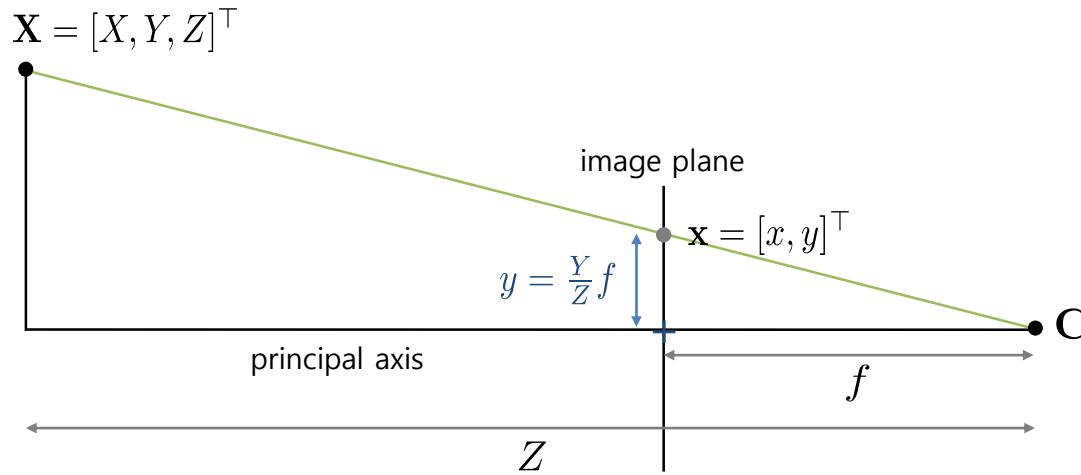
cf. conclusion in advance

$$\mathbf{x} = \mathbf{P}\mathbf{X} \quad (\mathbf{P} = \mathbf{K}[\mathbf{R} | \mathbf{t}])$$



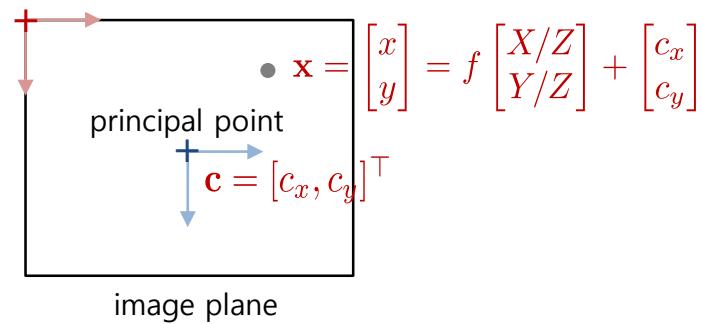
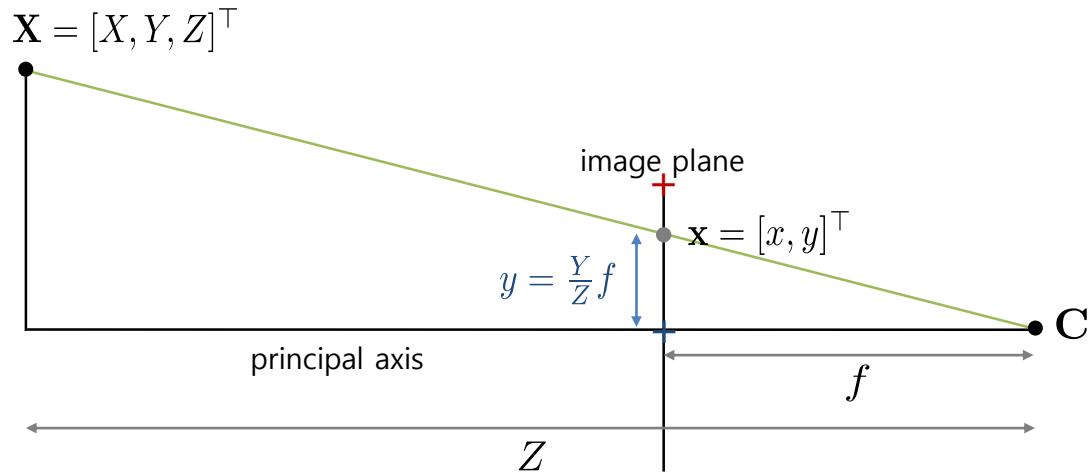
Camera Projection Model

- Pinhole Camera Model



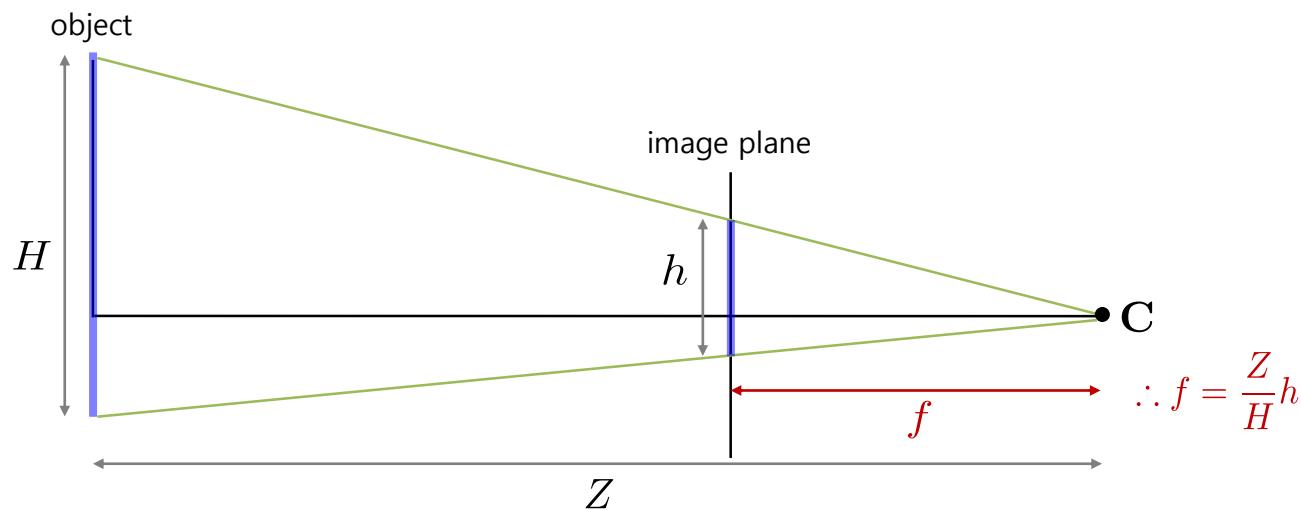
Camera Projection Model

- Pinhole Camera Model



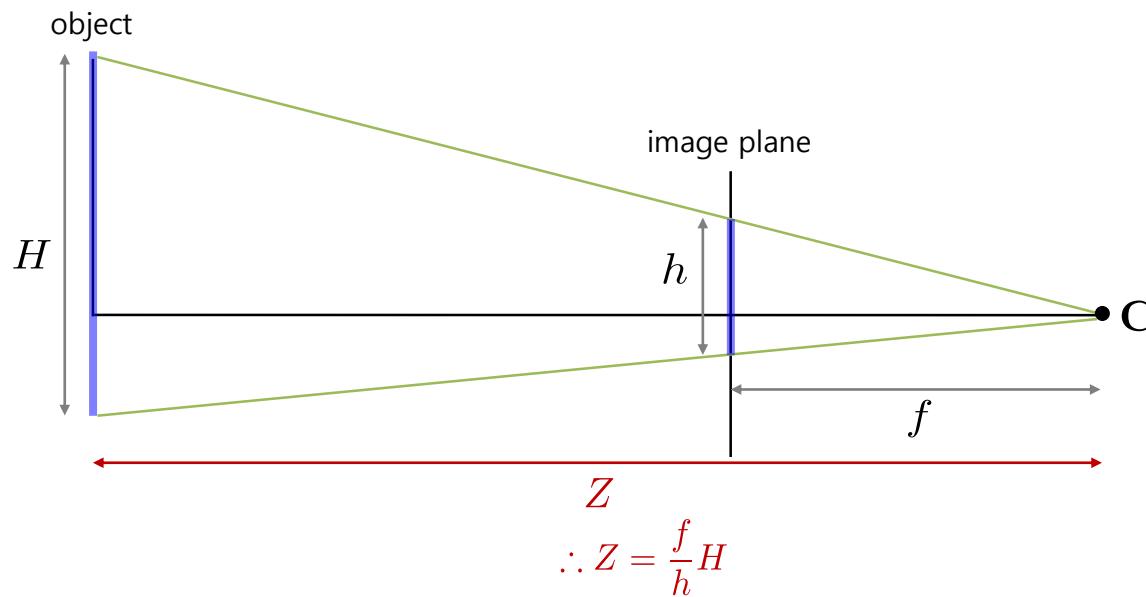
Camera Projection Model

- Example: **Simple Camera Calibration**
 - Unknown: **Focal length** (f) (unit: [pixel])
 - Given: The observed object height (h) on the image plane (unit: [pixel])
 - Assumptions
 - The object height (H) and distance (Z) from the camera are known.
 - The object is aligned with the image plane.



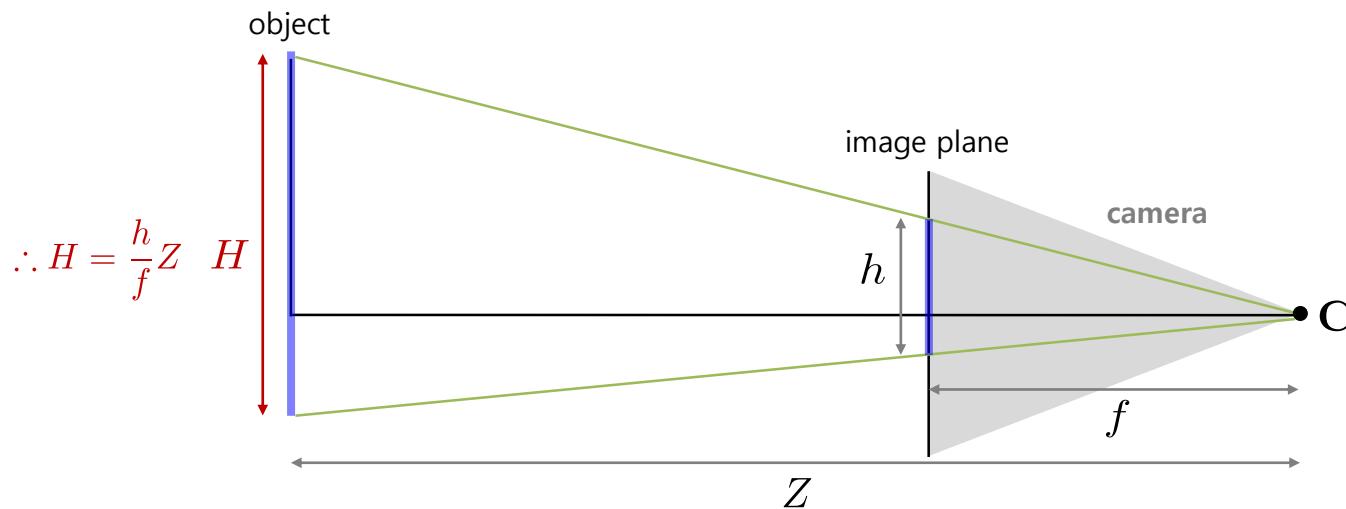
Camera Projection Model

- Example: **Simple Distance Measurement**
 - Unknown: **Object distance** (Z) (unit: [m])
 - Given: The observed object height (h) on the image plane (unit: [pixel])
 - Assumptions
 - The object height (H) and focal length (f) are known.
 - The object is aligned with the image plane.



Camera Projection Model

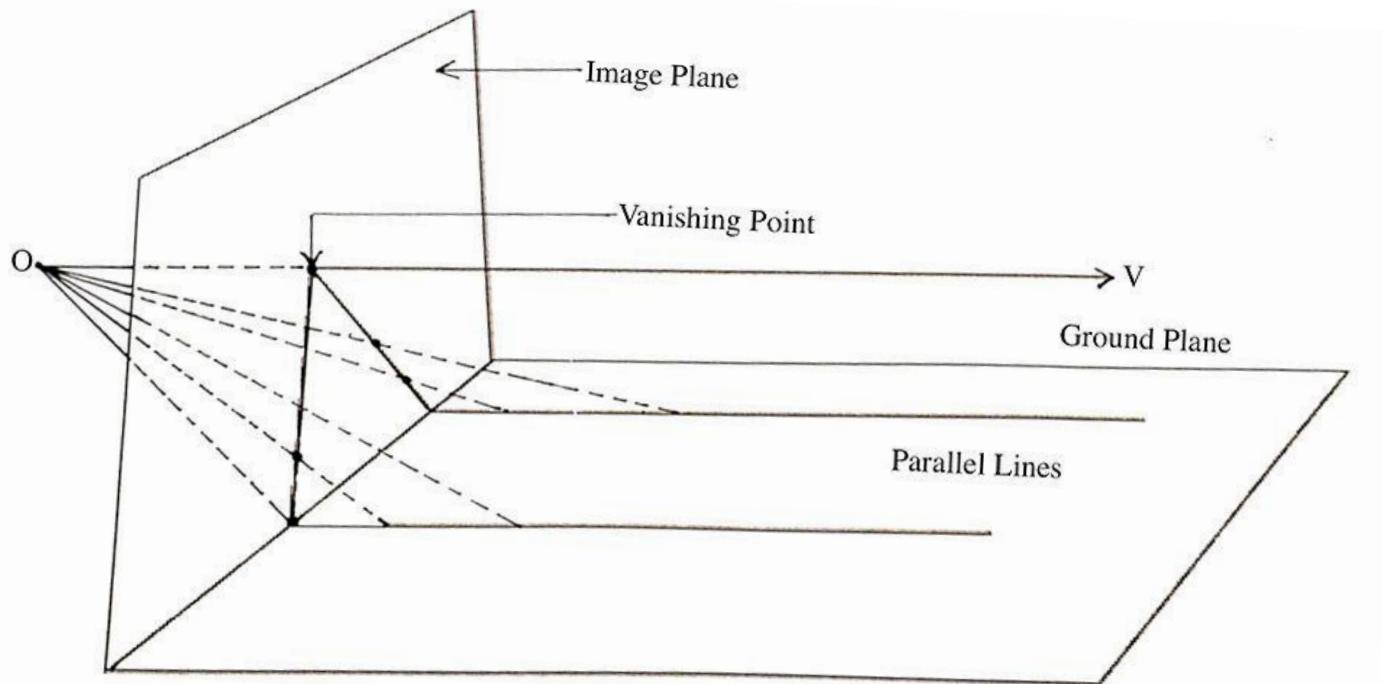
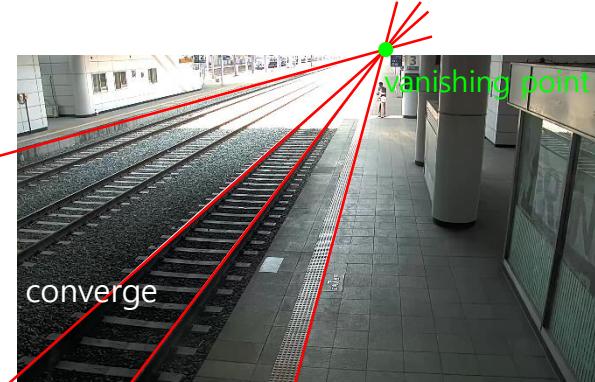
- Example: **Simple Height Measurement**
 - Unknown: **Object height** (H) (unit: [m])
 - Given: The observed object height (h) on the image plane (unit: [pixel])
 - Assumptions
 - The object distance (Z) from the camera and focal length (f) are known.
 - The object is aligned with the image plane.



Camera Projection Model

▪ Vanishing Points

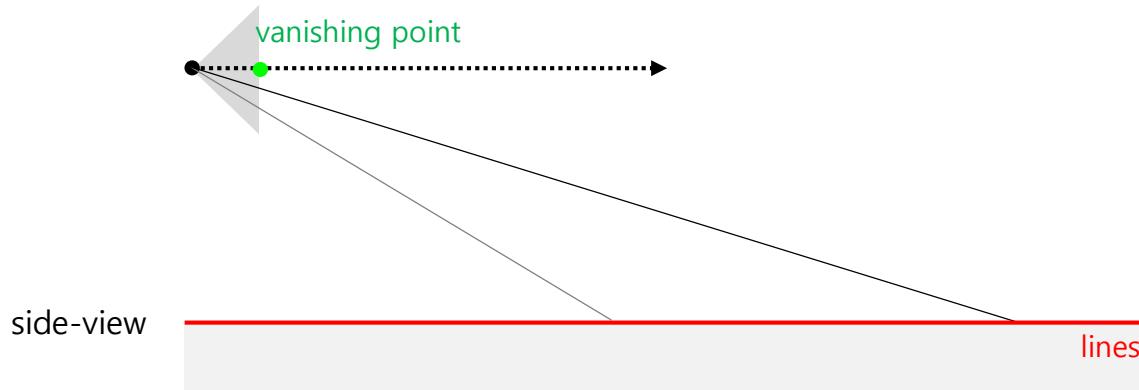
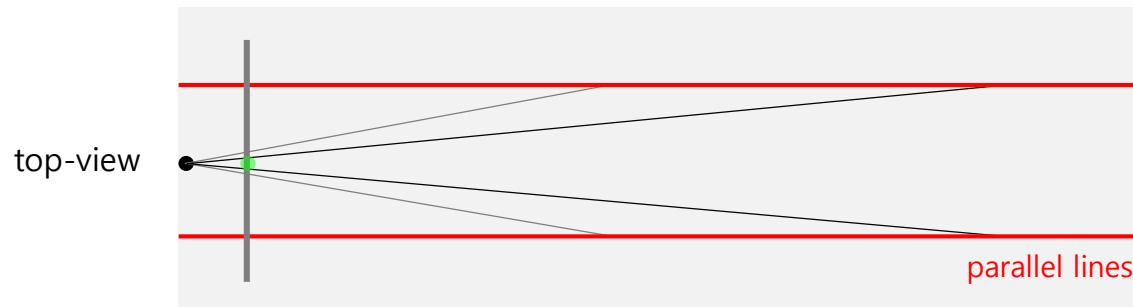
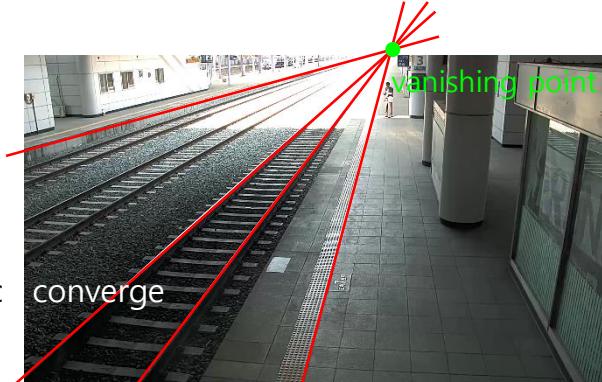
- A point on the image plane where mutually parallel lines in 3D space converge
 - A vector to the vanishing point is parallel to the lines.
 - A vector to the vanishing point is parallel to the reference plane made by the lines.



Camera Projection Model

▪ Vanishing Points

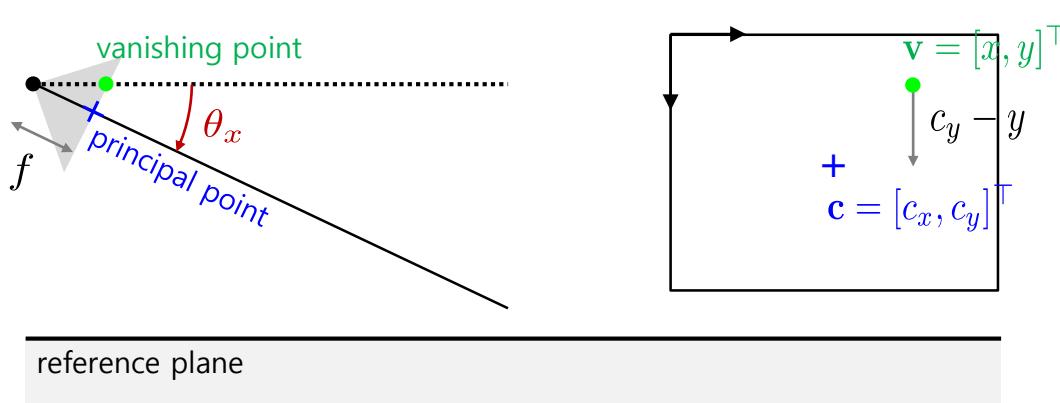
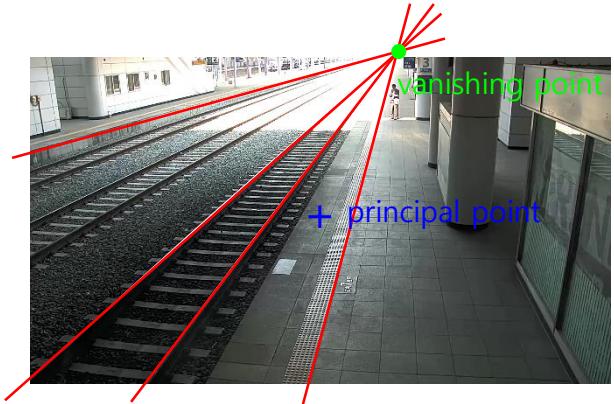
- A point on the image plane where mutually parallel lines in 3D space converge
 - A vector to the vanishing point is parallel to the lines.
 - A vector to the vanishing point is parallel to the reference plane made by the lines.



Camera Projection Model

- Example: Simple Visual Compass
 - Unknown: Tilt angle (θ_x) (unit: [rad])
 - Given: A vanishing point (x, y) from the reference plane
 - Assumptions
 - The focal length (f) is known.
 - The principal point (c_x, c_y) is known or selected as the center of images.
 - The camera has no roll, $\theta_z = 0$.

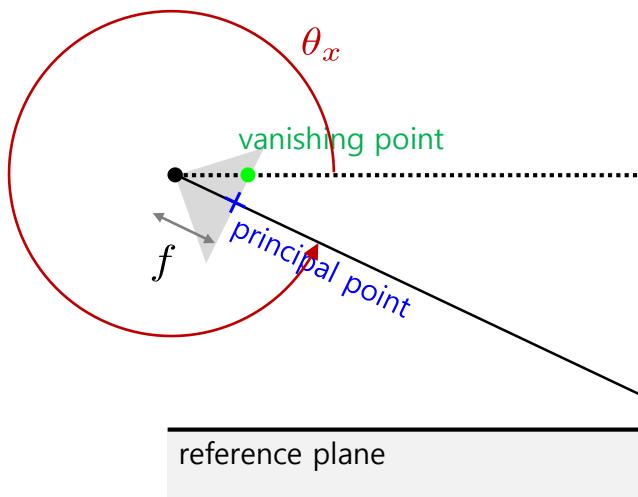
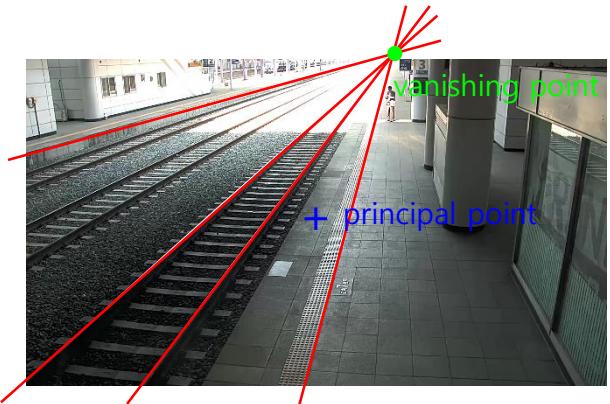
cf. The tilt angle in this page is defined as the opposite direction of the common notation.



Camera Projection Model

- Example: Simple Visual Compass
 - Unknown: Tilt angle (θ_x) (unit: [rad])
 - Given: A vanishing point (x, y) from the reference plane
 - Assumptions
 - The focal length (f) is known.
 - The principal point (c_x, c_y) is known or selected as the center of images.
 - The camera has no roll, $\theta_z = 0$.

cf. Similarly, the pan angle (θ_y) w.r.t. rails also can be calculated using x instead of y .



$$\therefore \theta_x = \tan^{-1} \frac{y - c_y}{f}$$

$\mathbf{v} = [x, y]^\top$
+
 $\mathbf{c} = [c_x, c_y]^\top$

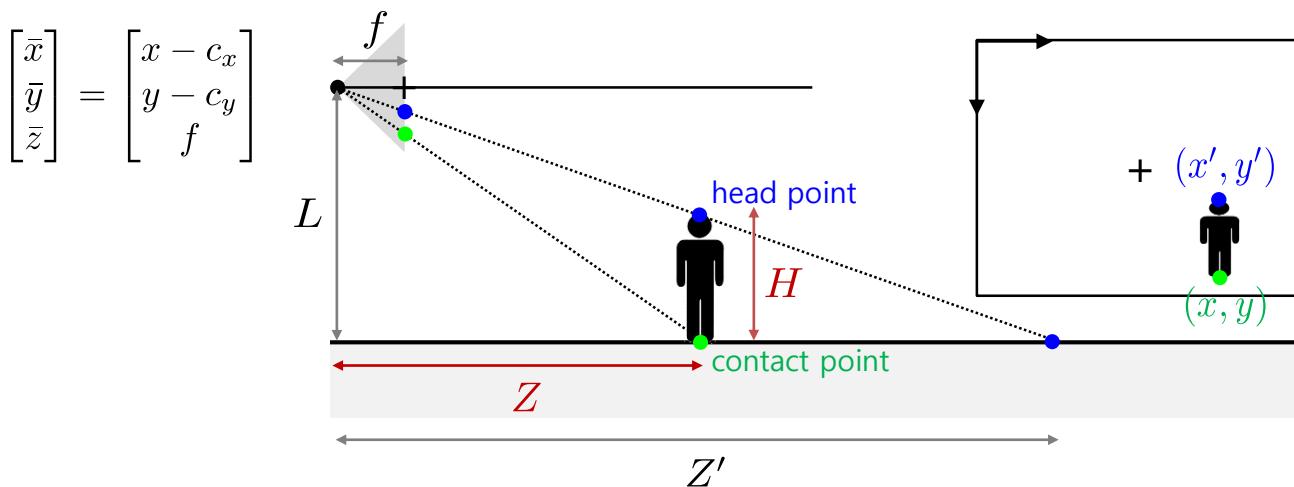
Camera Projection Model

- Example: **Object Localization and Measurement #1**

- Unknown: Object position and height (unit: [m])
- Given: The object's contact and head points on the image (unit: [pixel])
- Assumptions
 - The focal length, principal points, and camera height, are known.
 - The camera is aligned to the reference plane.
 - The object is on the reference plane.



$$\therefore Z = \frac{\bar{z}}{\bar{y}} L \quad X = \frac{\bar{x}}{\bar{y}} L \quad H = \left(\frac{\bar{y}}{\bar{z}} - \frac{\bar{y}'}{\bar{z}'} \right) Z$$



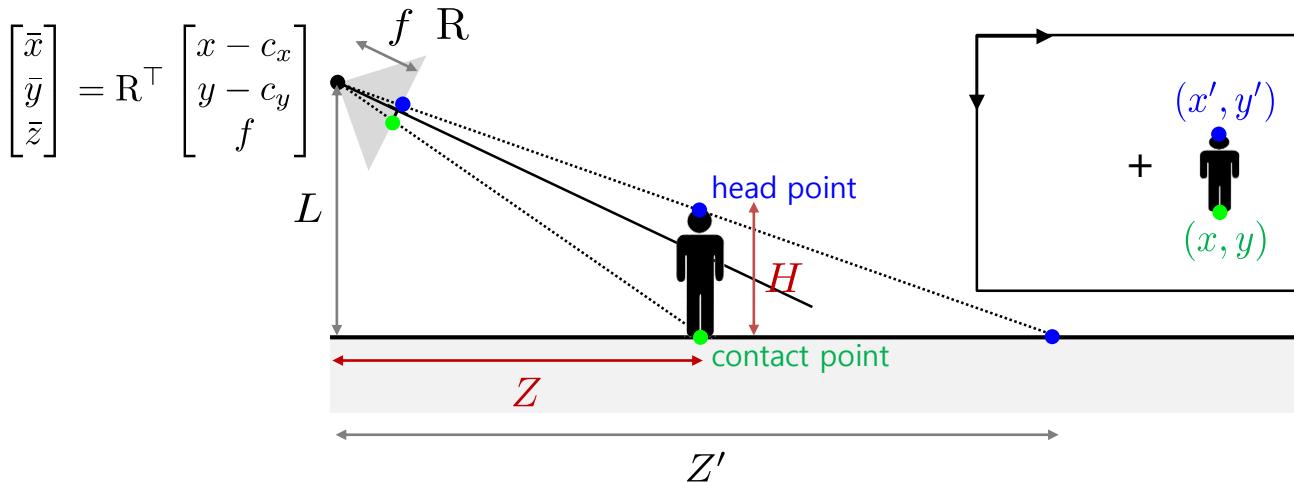
Camera Projection Model

- Example: Object Localization and Measurement #2
 - Unknown: Object position and height (unit: [m])
 - Given: The object's contact and head points on the image (unit: [pixel])
 - Assumptions
 - The focal length, principal points, and camera height, are known.
 - ~~The camera is aligned to the reference plane. The camera orientation is known.~~
 - The object is on the reference plane.



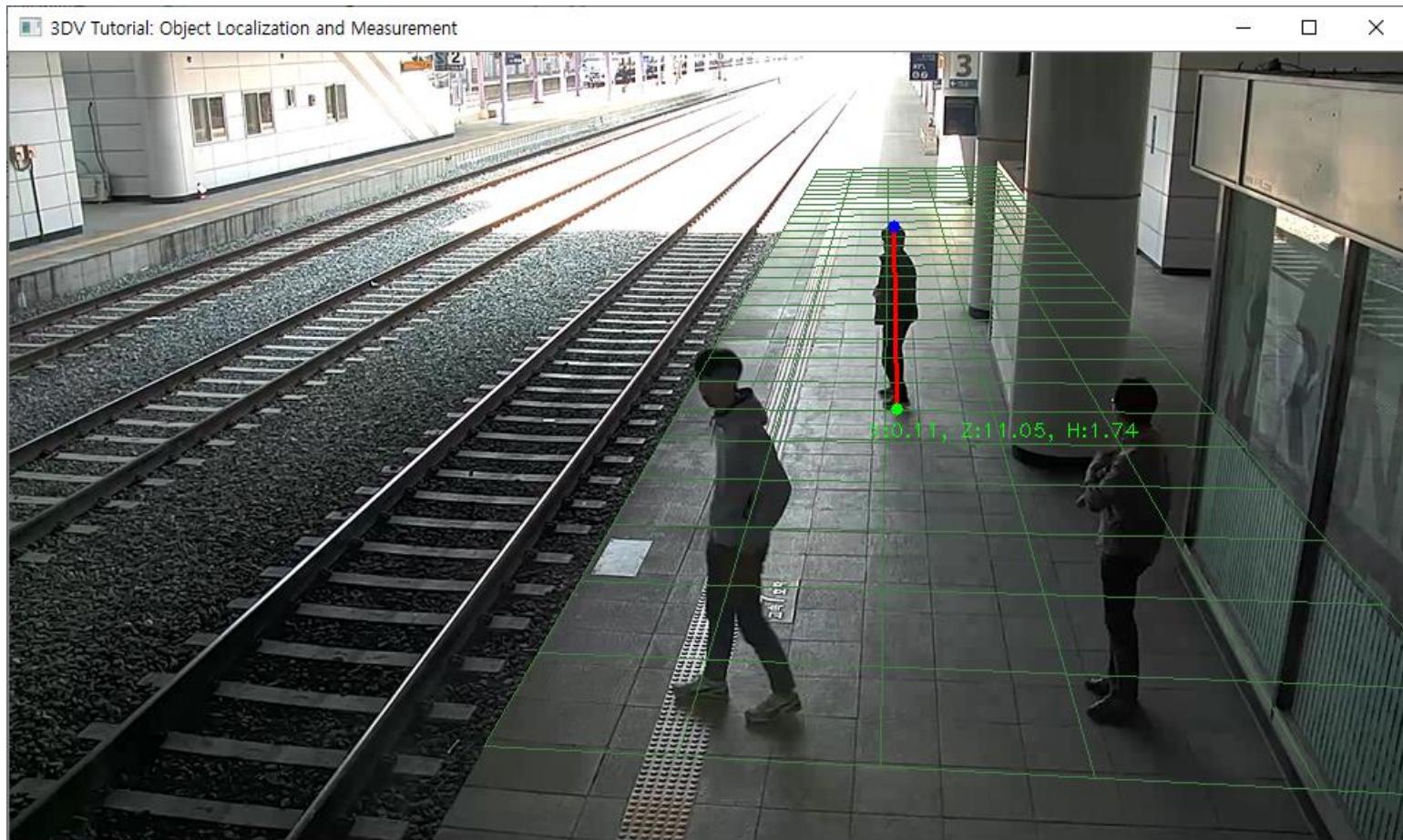
$$\therefore Z = \frac{\bar{z}}{\bar{y}} L \quad X = \frac{\bar{x}}{\bar{y}} L \quad H = \left(\frac{\bar{y}}{\bar{z}} - \frac{\bar{y}'}{\bar{z}'} \right) Z$$

\because same camera center



Camera Projection Model

- Example: **Object Localization and Measurement #2** [object_localization.cpp]





```

1. #include "opencv2/opencv.hpp"
2. ...
3. int main()
4. {
5.     const char* input = "data/daejeon_station.png";
6.     double f = 810.5, cx = 480, cy = 270, L = 3.31;
7.     cv::Point3d cam_ori(DEG2RAD(-18.7), DEG2RAD(-8.2), DEG2RAD(2.0));
8.     cv::Range grid_x(-2, 3), grid_z(5, 35);

9.     ...

10.    while (true)
11.    {
12.        cv::Mat image_copy = image.clone();
13.        if (drag.end.x > 0 && drag.end.y > 0)
14.        {
15.            // Calculate object location and height
16.            cv::Point3d c = R.t() * cv::Point3d(drag.start.x - cx, drag.start.y - cy, f);
17.            if (c.y < DBL_EPSILON) continue; // Skip the degenerate case (beyond the horizon)
18.            cv::Point3d h = R.t() * cv::Point3d(drag.end.x - cx, drag.end.y - cy, f);
19.            double Z = c.z / c.y * L, X = c.x / c.y * L, H = (c.y / c.z - h.y / h.z) * Z;

20.            // Draw head/contact points and location/height
21.            cv::line(image_copy, drag.start, drag.end, cv::Vec3b(0, 0, 255), 2);
22.            cv::circle(image_copy, drag.end, 4, cv::Vec3b(255, 0, 0), -1);
23.            cv::circle(image_copy, drag.start, 4, cv::Vec3b(0, 255, 0), -1);
24.            cv::putText(image_copy, cv::format("X:%.2f, Z:%.2f, H:%.2f", X, Z, H), ...);
25.        }

26.        // Show the image
27.        cv::imshow("3DV Tutorial: Object Localization and Measurement", image_copy);
28.        int key = cv::waitKey(1);
29.        if (key == 27) break; // 'ESC' key: Exit
30.    }
31.    return 0;
32. }

```

$$\begin{bmatrix} \bar{x} \\ \bar{y} \\ \bar{z} \end{bmatrix} = R^T \begin{bmatrix} x - c_x \\ y - c_y \\ f \end{bmatrix}$$

$$Z = \frac{\bar{z}}{\bar{y}} L$$

$$X = \frac{\bar{x}}{\bar{y}} L$$

$$H = \left(\frac{\bar{y}}{\bar{z}} - \frac{\bar{y}'}{\bar{z}'} \right) Z$$

Camera Projection Model

- Camera Matrix K

$\bullet \quad \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix} = f \begin{bmatrix} X/Z \\ Y/Z \end{bmatrix} + \begin{bmatrix} c_x \\ c_y \end{bmatrix}$

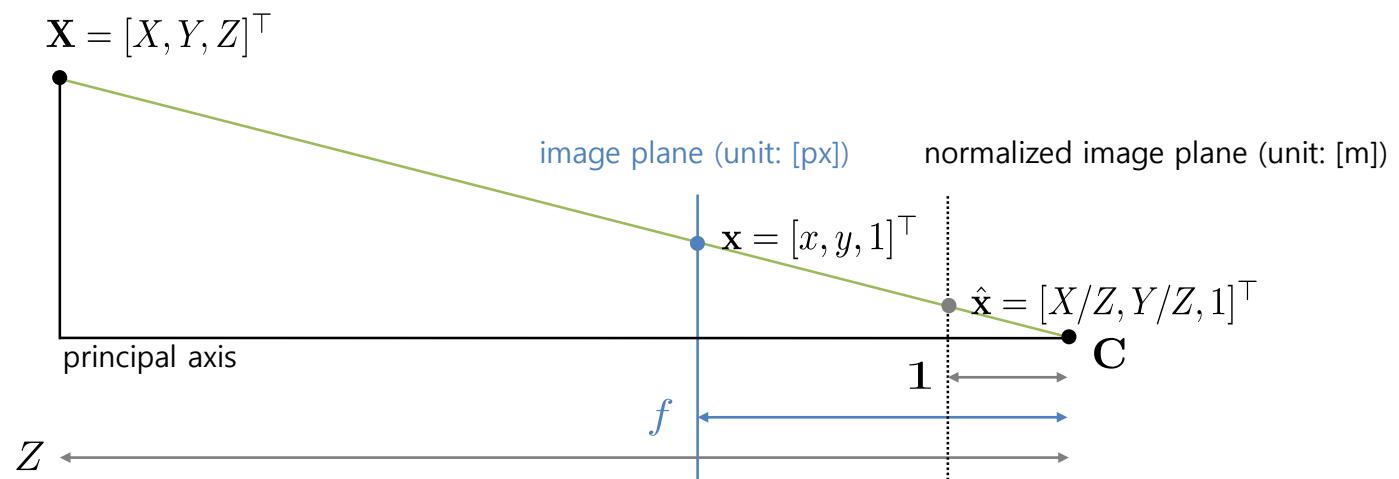
+ $\rightarrow \mathbf{x} = \mathbf{K}\hat{\mathbf{x}}$ where $\mathbf{K} = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}$, $\mathbf{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$, and $\hat{\mathbf{x}} = \begin{bmatrix} X/Z \\ Y/Z \\ 1 \end{bmatrix}$

specialized ↑
generalized ↓

$$\mathbf{K} = \begin{bmatrix} f & 0 & w/2 \\ 0 & f & h/2 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{K} = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$



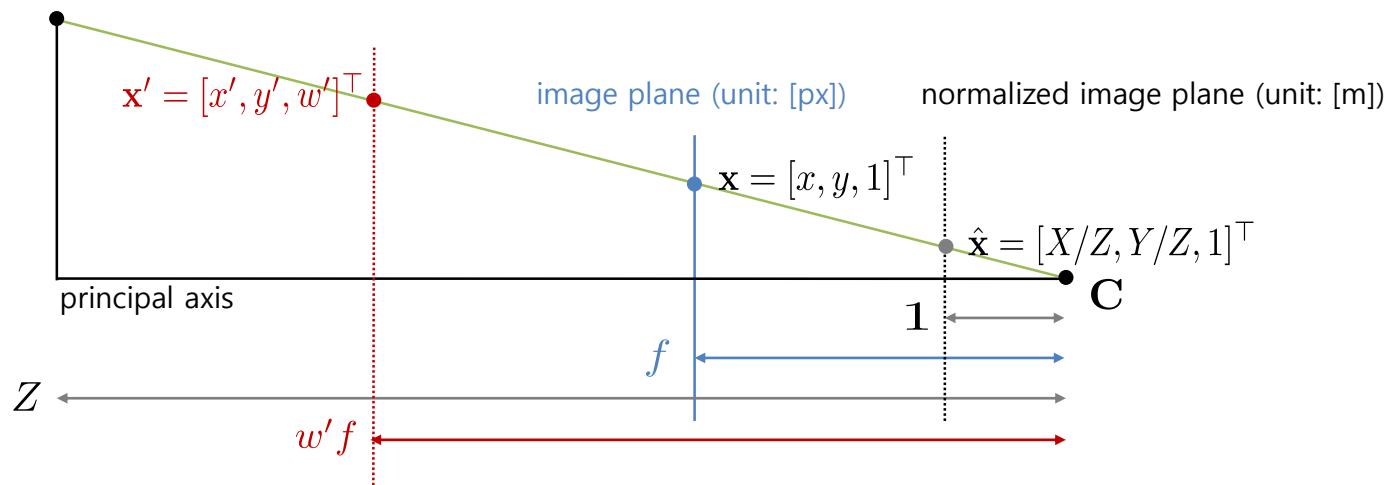
Camera Projection Model

- **Homogeneous Coordinates** (a.k.a. projective coordinates)
 - It describes n -dimensional project space as $n + 1$ -dimensional coordinate system. It holds non-conventional equivalence relationship: $(x_1, x_2, \dots, x_{n+1}) \sim (\lambda x_1, \lambda x_2, \dots, \lambda x_{n+1})$ such that $(0 \neq \lambda \in \mathbb{R})$
 - e.g. $(5, 12)$ is written as $(5, 12, 1)$ which is also equal to $(10, 24, 2)$ or $(30, 36, 3)$ or ...
 - On the previous slide,

$$\mathbf{x} = K\hat{\mathbf{x}} \quad \text{where} \quad K = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}, \quad \text{and} \quad \hat{\mathbf{x}} = \begin{bmatrix} X/Z \\ Y/Z \\ 1 \end{bmatrix}$$

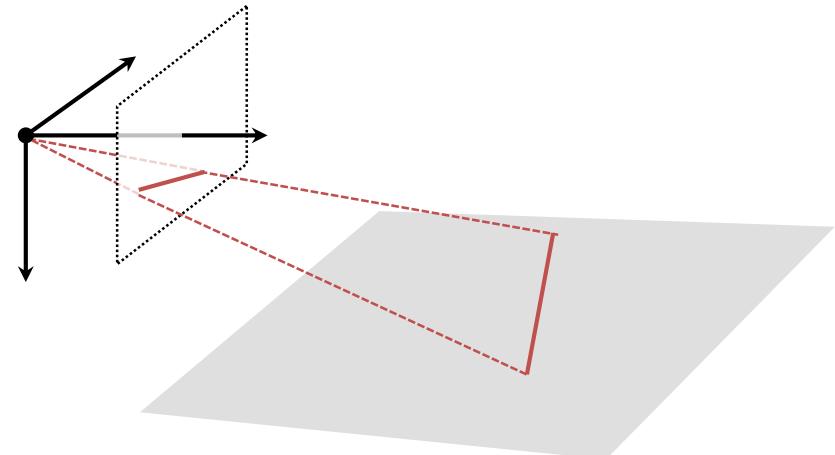
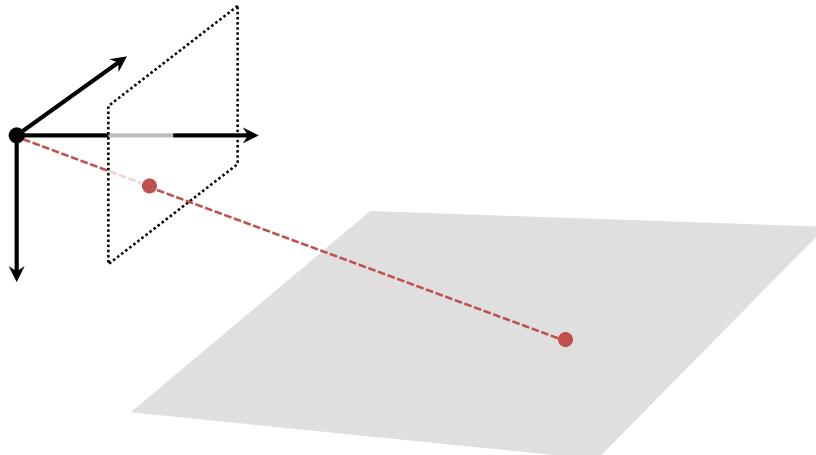
$$\mathbf{x}' = K\mathbf{X} \quad \text{where} \quad K = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{x}' = \begin{bmatrix} x' \\ y' \\ w' \end{bmatrix}, \quad \text{and} \quad \mathbf{X} = \begin{bmatrix} X \\ Y \\ Z \end{bmatrix}$$

$$\mathbf{X} = [X, Y, Z]^\top$$



Camera Projection Model

- Why Homogeneous Coordinates?
 - An affine transformation ($\mathbf{y} = \mathbf{Ax} + \mathbf{b}$) is formulated by a single matrix multiplication.
 - A light ray (line at the camera center) is observed as a point on the image plane.
 - A plane at the camera center is observed as a line on the image plane.
 - A conic whose peak is at the camera center is observed as a conic section on the image plane.
 - A point at infinity (a.k.a. ideal point) is numerically represented by $w = 0$.
 - A point and line ($ax + by + c = 0$) are described beautifully as like $\mathbf{l}^T \mathbf{x} = 0$ or $\mathbf{x}^T \mathbf{l} = 0$ ($\mathbf{l} = [a, b, c]^T$).
 - Intersection of two lines: $\mathbf{x} = \mathbf{l}_1 \times \mathbf{l}_2$
 - A line by two points: $\mathbf{l} = \mathbf{x}_1 \times \mathbf{x}_2$
 - ...



Camera Projection Model

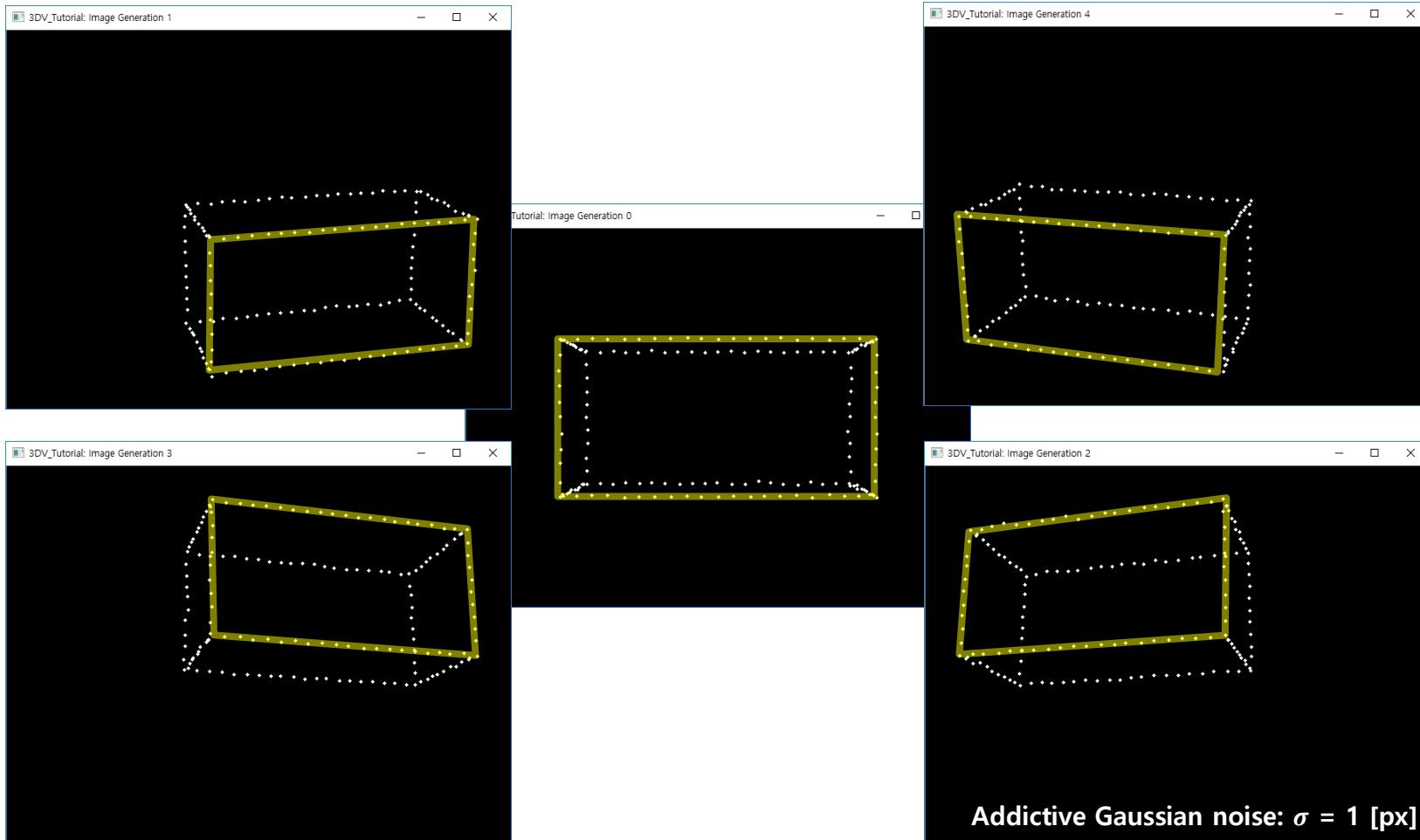
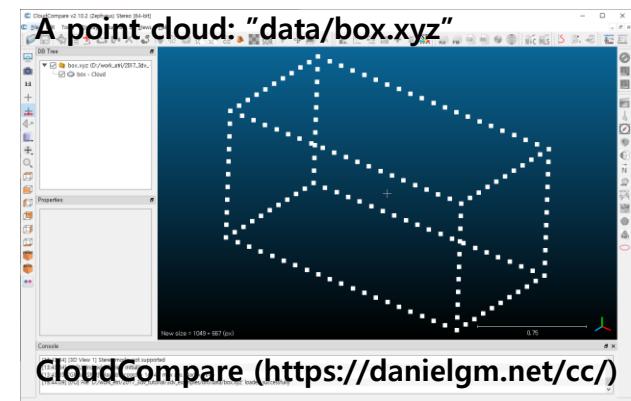
- **Projection Matrix P**
 - If a 3D point is not based on the camera coordinate, the point should be transformed.
-
- $\mathbf{X} = [X, Y, Z]^\top$
 $\mathbf{X}' = \mathbf{R}\mathbf{X} + \mathbf{t} \longrightarrow \mathbf{X}' = [\mathbf{R} \mid \mathbf{t}] \mathbf{X}$ where $\mathbf{X} = \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$
 $\mathbf{x} = \mathbf{K}\mathbf{X}'$
 $\mathbf{x} = \mathbf{K}[\mathbf{R} \mid \mathbf{t}] \mathbf{X}$
 $\mathbf{x} = \mathbf{P}\mathbf{X}$ where $\mathbf{P} = \mathbf{K}[\mathbf{R} \mid \mathbf{t}]$
 3×4 matrix
 3×4 matrix
- The camera pose can be derived from \mathbf{R}^\top and $-\mathbf{R}^\top \mathbf{t}$.

$$\text{c.f. } \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{R}^\top & -\mathbf{R}^\top \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

- **Camera Parameters** ~ projection matrix
 - **Intrinsic parameters** ~ camera matrix
 - e.g. focal length, principle point, skew, distortion coefficient, ...
 - **Extrinsic parameters** ~ point transformation
 - e.g. rotation and translation

Camera Projection Model

- Example: **Image Formation** [image_formation.cpp]



Camera Projection Model



- Example: **Image Formation** [image_formation.cpp]

```
1. #include "opencv2/opencv.hpp"

2. #define Rx(rx)      (cv::Mat<double>(3, 3) << 1, 0, 0, 0, cos(rx), -sin(rx), 0, sin(rx), cos(rx))
3. #define Ry(ry)      (cv::Mat<double>(3, 3) << cos(ry), 0, sin(ry), 0, 1, 0, -sin(ry), 0, cos(ry))
4. #define Rz(rz)      (cv::Mat<double>(3, 3) << cos(rz), -sin(rz), 0, sin(rz), cos(rz), 0, 0, 0, 1)

5. int main()
6. {
7.     // The given camera configuration: focal length, principal point, image resolution, position, and orientation
8.     double f = 1000, cx = 320, cy = 240, noise_std = 1;
9.     cv::Size img_res(640, 480);
10.    std::vector<cv::Point3d> cam_pos = { cv::Point3d(0, 0, 0), cv::Point3d(-2, -2, 0), ... };
11.    std::vector<cv::Point3d> cam_ori = { cv::Point3d(0, 0, 0), cv::Point3d(-CV_PI / 12, CV_PI / 12, 0), ... };

12.    // Load a point cloud in the homogeneous coordinate
13.    FILE* fin = fopen("data/box.xyz", "rt");
14.    if (fin == NULL) return -1;
15.    cv::Mat X;
16.    while (!feof(fin))
17.    {
18.        double x, y, z;
19.        if (fscanf(fin, "%lf %lf %lf", &x, &y, &z) == 3) X.push_back(cv::Vec4d(x, y, z, 1));
20.    }
21.    fclose(fin);
22.    X = X.reshape(1).t(); // Convert to a 4 x N matrix
```

```

24. // Generate images for each camera pose
25. cv::Mat K = (cv::Mat<double>(3, 3) << f, 0, cx, 0, f, cy, 0, 0, 1);    ∴ K = 
$$\begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

26. for (size_t i = 0; i < cam_pos.size(); i++)
27. {
28.     // Derive a projection matrix
29.     cv::Mat Rc = Rz(cam_ori[i].z) * Ry(cam_ori[i].y) * Rx(cam_ori[i].x);
30.     cv::Mat tc(cam_pos[i]);
31.     cv::Mat Rt;
32.     cv::hconcat(Rc.t(), -Rc.t() * tc, Rt);
33.     cv::Mat P = K * Rt;
34. 
35.     // Project the points (cf. OpenCV provides 'cv::projectPoints()' with consideration of distortion.)
36.     cv::Mat x = P * X;
37.     x.row(0) = x.row(0) / x.row(2);
38.     x.row(1) = x.row(1) / x.row(2);
39.     x.row(2) = 1;
40. 
41.     // Add Gaussian noise
42.     cv::Mat noise(2, x.cols, x.type());
43.     cv::randn(noise, cv::Scalar(0), cv::Scalar(noise_std));
44.     x.rowRange(0, 2) = x.rowRange(0, 2) + noise;
45. 
46.     // Show and store the points
47.     cv::Mat image = cv::Mat::zeros(img_res, CV_8UC1);
48.     for (int c = 0; c < x.cols; c++)
49.     {
50.         cv::Point p(x.col(c).rowRange(0, 2));
51.         if (p.x >= 0 && p.x < img_res.width && p.y >= 0 && p.y < img_res.height)
52.             cv::circle(image, p, 2, 255, -1);
53.     }
54.     cv::imshow(cv::format("3DV_Tutorial: Image Formation %d", i), image);
55. 
56.     FILE* fout = fopen(cv::format("imageFormation%d.xyz", i).c_str(), "wt");
57.     if (fout == NULL) return -1;
58.     for (int c = 0; c < x.cols; c++)
59.         fprintf(fout, "%f %f 1\n", x.at<double>(0, c), x.at<double>(1, c));
60.     fclose(fout);
61. }
62. 
63. cv::waitKey(0);
64. return 0;
65. }
```

$$\therefore \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix}^{-1} = \begin{bmatrix} R^\top & -R^\top t \\ 0 & 1 \end{bmatrix}$$

$$\therefore \begin{bmatrix} x \\ y \\ w \end{bmatrix} \sim \begin{bmatrix} x/w \\ y/w \\ 1 \end{bmatrix}$$

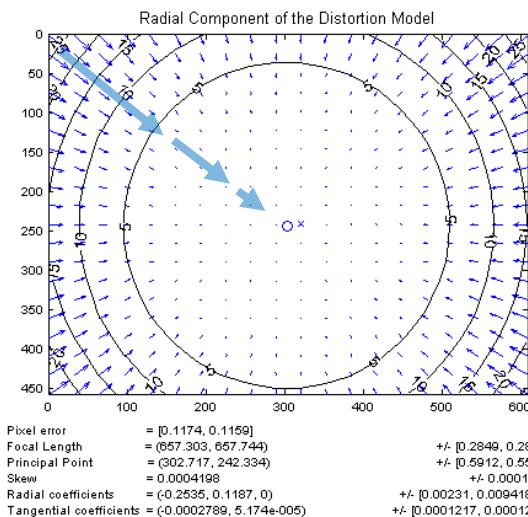


Camera Projection Model

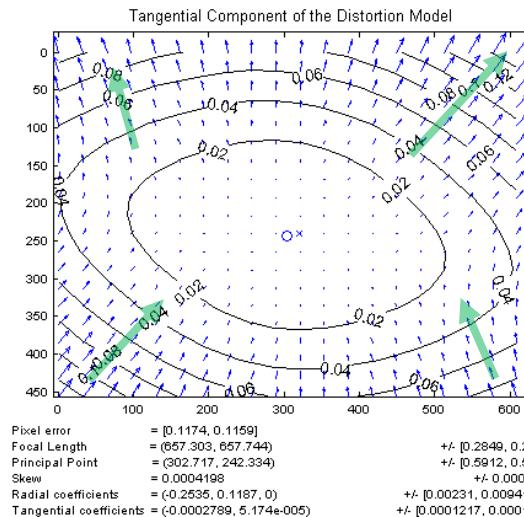
▪ Geometric Distortion

- Geometric distortion models are usually defined on the normalized image plane.
- **Polynomial distortion model** (a.k.a. Brown-Conrady model)
 - **Radial distortion:** k_1, k_2, \dots
 - **Tangential distortion:** p_1, p_2, \dots (usually negligible)

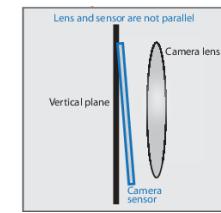
$$\begin{bmatrix} \hat{x}_d \\ \hat{y}_d \end{bmatrix} = (1+k_1r^2+k_2r^4+\dots) \begin{bmatrix} \hat{x} \\ \hat{y} \end{bmatrix} + (1+p_3r^2+p_4r^4+\dots) \begin{bmatrix} 2p_1\hat{x}\hat{y} + p_2(r^2 + 2\hat{x}^2) \\ 2p_2\hat{x}\hat{y} + p_1(r^2 + 2\hat{y}^2) \end{bmatrix} \quad \text{where} \quad r^2 = \hat{x}^2 + \hat{y}^2$$



radial distortion



tangential distortion



Why? lens and sensor are not parallel.

– Field-of-view (FOV) model

- Radial distortion is modeled by a tangent function (good for severe distortion; e.g. fisheye lens).

Camera Projection Model

- **Geometric Distortion Correction**

- Input: The original image
- Output: Its rectified image (w/o geometric distortion)
- Given: Its camera matrix and distortion coefficient
- Solutions for the polynomial distortion model
 - OpenCV `cv::undistort()` and `cv::undistortPoints()` (cf. included in `imgproc` module)
↔ `cv::projectPoints()` (cf. included in `calib3d` module)
 - Camera Distortion Correction: Theory and Practice, <http://darkpgmr.tistory.com/31>



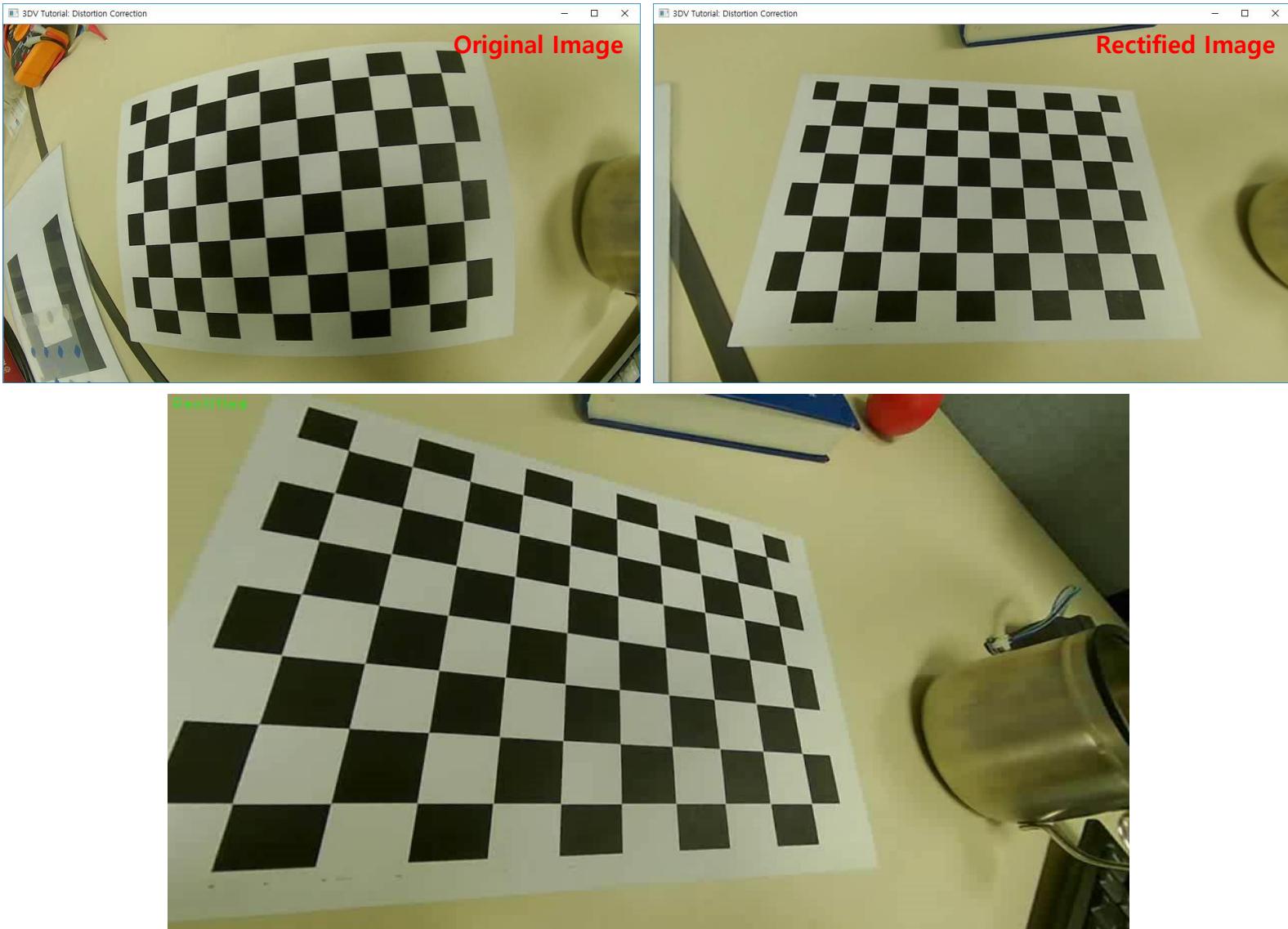
distortion correction →

K1: 1.105763E-01
K2: 1.886214E-02
K3: 1.473832E-02
P1: -8.448460E-03
P2: -7.356744E-03



Camera Projection Model

- Example: Geometric Distortion Correction [distortion_correction.cpp]





```
1. #include <opencv2/opencv.hpp>

2. int main()
3. {
4.     const char* input = "data/chessboard.avi";
5.     cv::Matx3d K(432.7390364738057, 0, 476.0614994349778, 0, 431.2395555913084, 288.7602152621297, 0, 0, 1);
6.     std::vector<double> dist_coeff = { -0.2852754904152874, 0.1016466459919075, ... };

7.     // Open a video
8.     cv::VideoCapture video;
9.     if (!video.open(input)) return -1;

10.    // Run distortion correction
11.    bool show_rectify = true;
12.    cv::Mat map1, map2;
13.    while (true)
14.    {
15.        // Grab an image from the video
16.        cv::Mat image;
17.        video >> image;
18.        if (image.empty()) break;

19.        // Rectify geometric distortion (cf. 'cv::undistort()' can be applied for one-time remapping.)
20.        if (show_rectify)
21.        {
22.            if (map1.empty() || map2.empty())
23.                cv::initUndistortRectifyMap(K, dist_coeff, cv::Mat(), cv::Mat(), image.size(), CV_32FC1, map1, map2);
24.            cv::remap(image, image, map1, map2, cv::InterpolationFlags::INTER_LINEAR);
25.        }

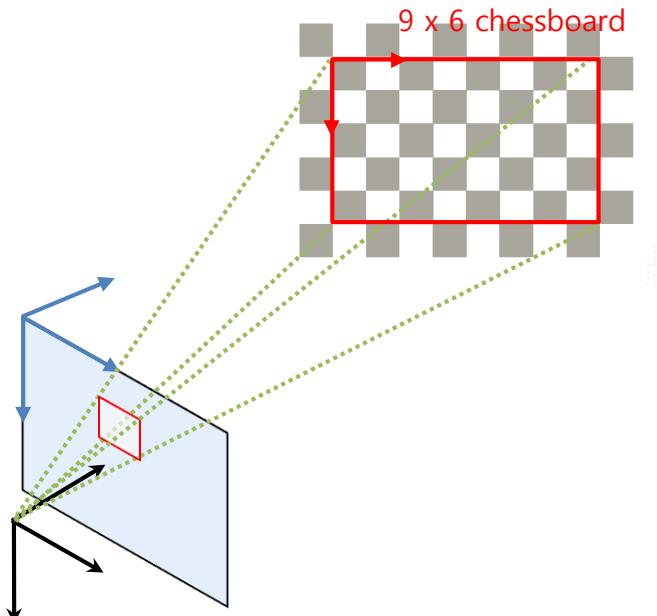
26.        // Show the image
27.        cv::imshow("3DV Tutorial: Distortion Correction", image);
28.        int key = cv::waitKey(1);
29.        if (key == 27) break;                                // 'ESC' key: Exit
30.        else if (key == 9) show_rectify = !show_rectify;    // 'Tab' key: Toggle rectification
31.        ...
32.    }

33.    video.release();
34.    return 0;
35. }
```

General 2D-3D Geometry

▪ Camera Calibration

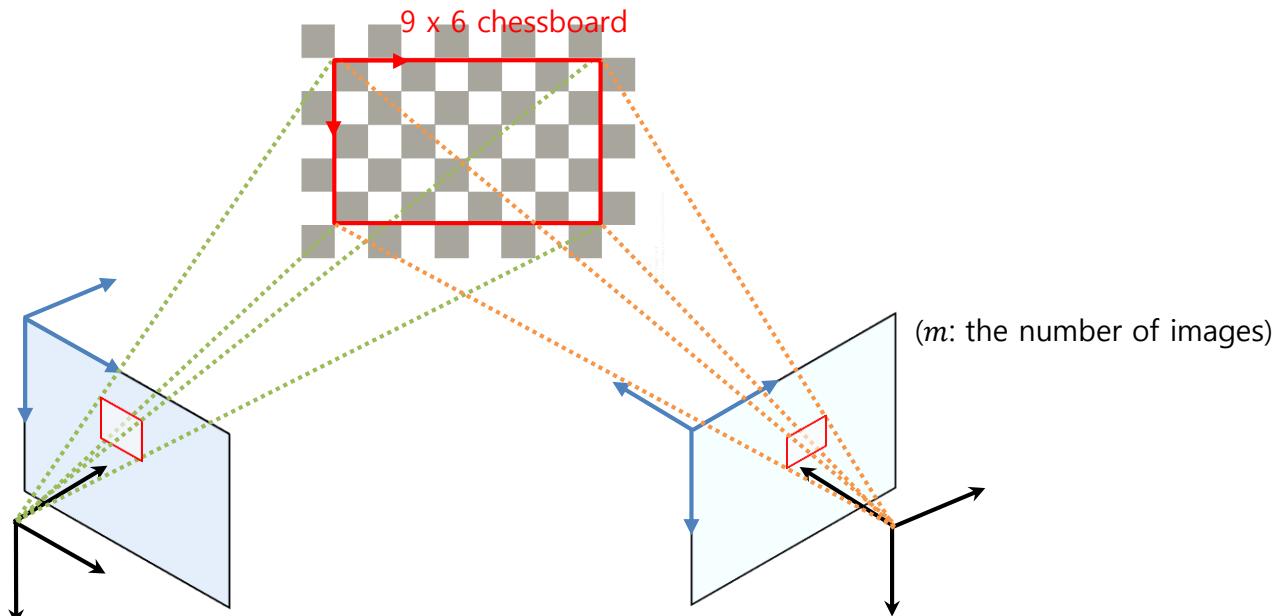
- Unknown: Intrinsic + extrinsic parameters (5^* + 6 DoF)
 - The number of intrinsic parameters* can be varied w.r.t. user preference.
- Given: 3D points $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ and their projected points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$
- Constraints: $n \times$ projection $\mathbf{x}_i = \mathbf{K} [\mathbf{R} | \mathbf{t}] \mathbf{X}_i$



General 2D-3D Geometry

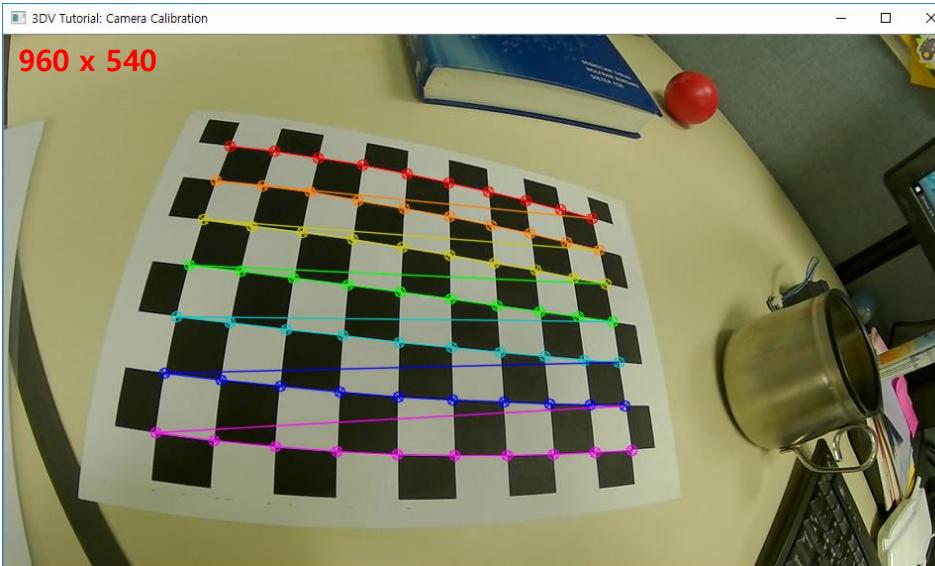
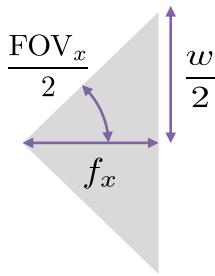
▪ Camera Calibration

- Unknown: Intrinsic + $m \times$ extrinsic parameters ($5^* + m \times 6$ DoF)
- Given: 3D points $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ and their projected points from j th camera \mathbf{x}_i^j
- Constraints: $n \times m \times$ projection $\mathbf{x}_i^j = \mathbf{K} [\mathbf{R}_j | \mathbf{t}_j] \mathbf{X}_i$
- Solutions
 - OpenCV `cv::calibrateCamera()` and `cv::initCameraMatrix2D()`
 - Camera Calibration Toolbox for MATLAB, http://www.vision.caltech.edu/bouguetj/calib_doc/
 - GML C++ Camera Calibration Toolbox, <http://graphics.cs.msu.ru/en/node/909>
 - DarkCamCalibrator, <http://darkpgmr.tistory.com/139>



General 2D-3D Geometry

- Example: Camera Calibration [camera_calibration.cpp]



field-of-view (FOV) = focal length

$$\text{- Horizontal: } 2 \times \tan^{-1} \frac{f_x}{w/2} = 84^\circ$$

$$\text{- Vertical: } 2 \times \tan^{-1} \frac{f_y}{h/2} = 116^\circ$$

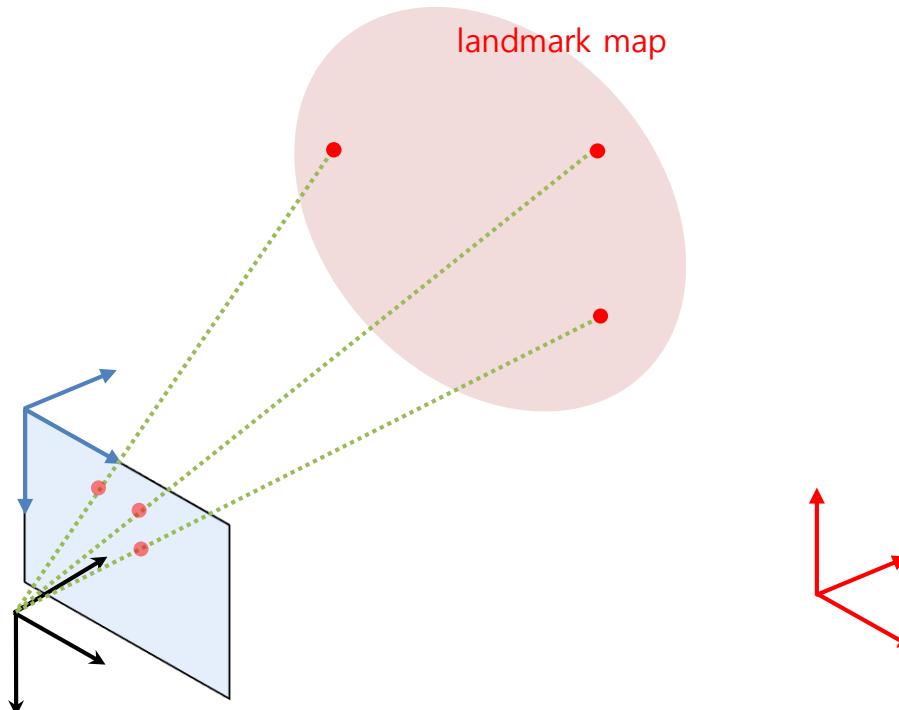
```
## Camera Calibration Results
* The number of applied images = 22
* RMS error = 0.473353
* Camera matrix (K) =
  [432.7390364738057, 0, 476.0614994349778] ~ close to the center of the image, (480, 270)
  [0, 431.2395555913084, 288.7602152621297]
  [0, 0, 1]
* Distortion coefficient (k1, k2, p1, p2, k3, ...) =
  [-0.2852754904152874, 0.1016466459919075,
   -0.0004420196146339175, 0.0001149909868437517,
   -0.01803978785585194] ~ close to zero (usually negligible)
```



```
1. #include "opencv2/opencv.hpp"
2. int main()
3. {
4.     const char* input = "data/chessboard.avi";
5.     cv::Size board_pattern(10, 7);
6.     float board_cellsize = 0.025f;
7.
8.     // Select images
9.     std::vector<cv::Mat> images;
10.    ...
11.
12.    // Find 2D corner points from the given images
13.    std::vector<std::vector<cv::Point2f>> img_points;
14.    for (size_t i = 0; i < images.size(); i++)
15.    {
16.        std::vector<cv::Point2f> pts;
17.        if (cv::findChessboardCorners(images[i], board_pattern, pts))
18.            img_points.push_back(pts);
19.    }
20.    if (img_points.empty()) return -1;
21.
22.    // Prepare 3D points of the chess board
23.    std::vector<std::vector<cv::Point3f>> obj_points(1);
24.    for (int r = 0; r < board_pattern.height; r++)
25.        for (int c = 0; c < board_pattern.width; c++)
26.            obj_points[0].push_back(cv::Point3f(board_cellsize * c, board_cellsize * r, 0));
27.    obj_points.resize(img_points.size(), obj_points[0]); // Copy
28.
29.    // Calibrate the camera
30.    cv::Mat K = cv::Mat::eye(3, 3, CV_64F);
31.    cv::Mat dist_coeff = cv::Mat::zeros(4, 1, CV_64F);
32.    std::vector<cv::Mat> rvecs, tvecs;
33.    double rms = cv::calibrateCamera(obj_points, img_points, images[0].size(), K, dist_coeff, rvecs, tvecs);
34.
35.    // Report calibration results      X1,X2,...      X1,X2,...
36.    ...
37.    return 0;
38. }
```

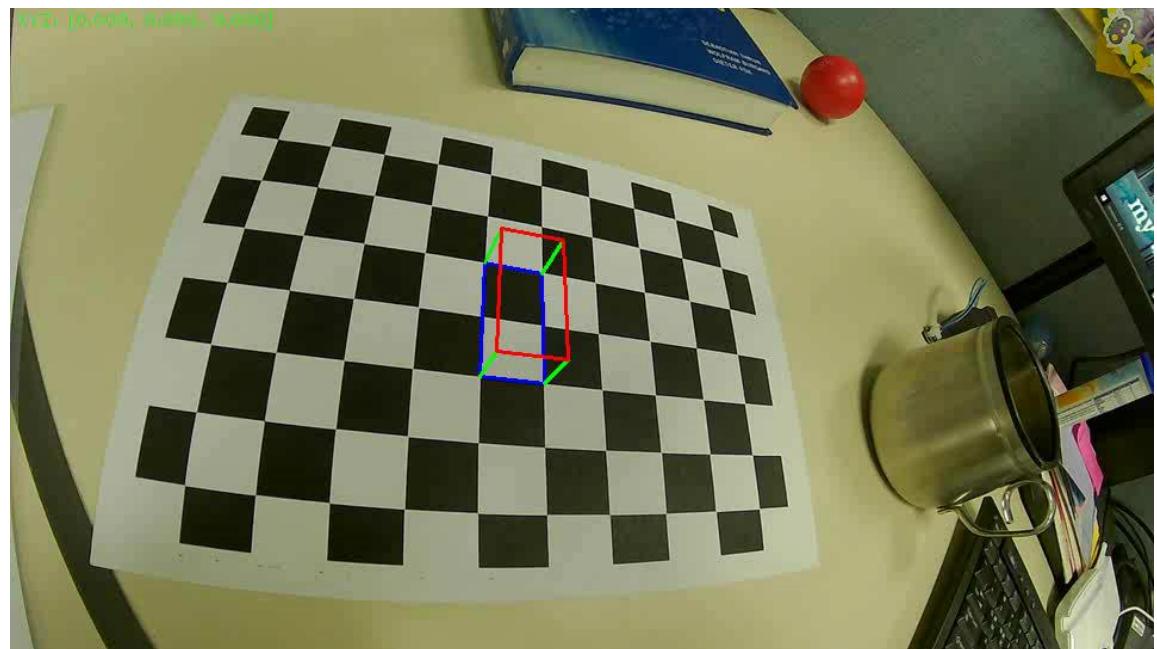
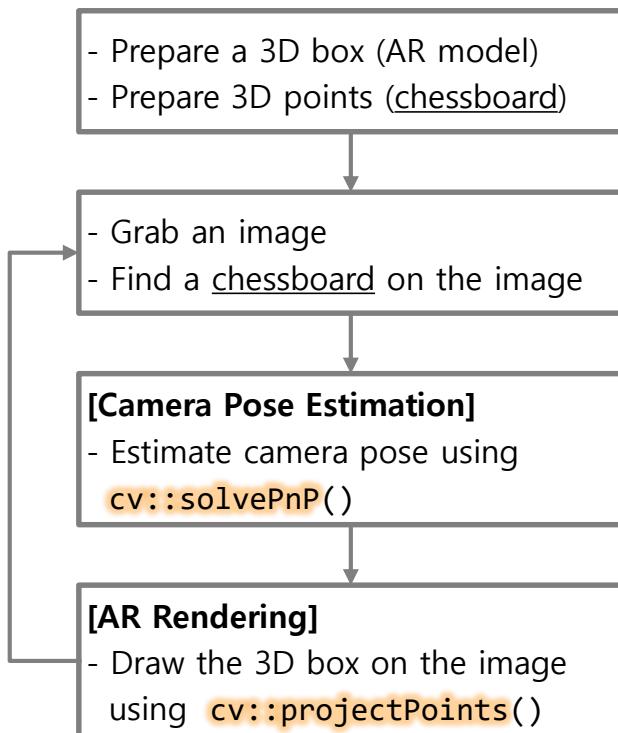
General 2D-3D Geometry

- **Absolute Camera Pose Estimation** (Perspective-n-Point; PnP)
 - Unknown: Camera pose (6 DoF)
 - Given: 3D points $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$, their projected points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, and camera matrix K
 - Constraints: $n \times$ projection $\mathbf{x}_i = \mathbf{K} [\mathbf{R} | \mathbf{t}] \mathbf{X}_i$
 - Solutions ($n \geq 3$) → 3-point algorithm
 - OpenCV `cv::solvePnP()` and `cv::solvePnPRansac()`
 - Efficient PnP (EPnP), <http://cvlab.epfl.ch/EPnP/>



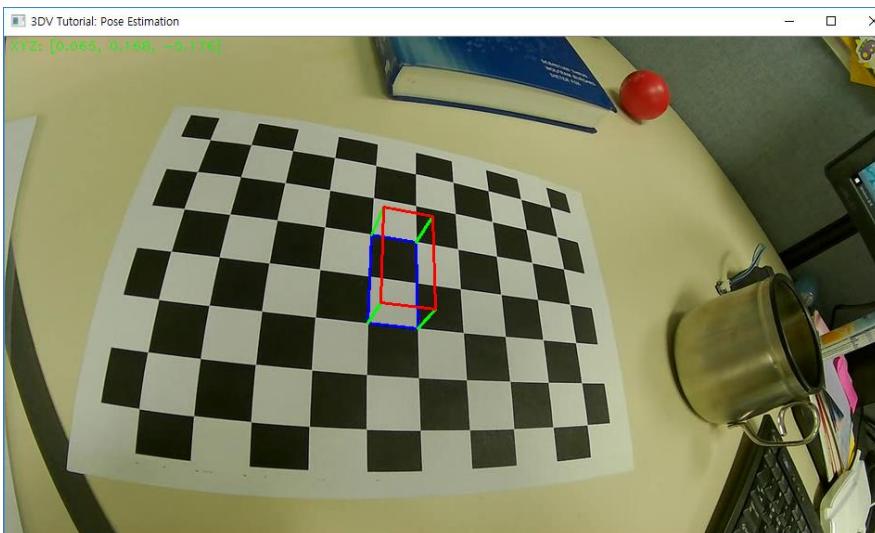
General 2D-3D Geometry

- Example: **Pose Estimation (Chessboard)** [pose_estimation_chessboard.cpp]





```
1. #include "opencv2/opencv.hpp"
2. int main()
3. {
4.     const char* input = "data/chessboard.avi";
5.     cv::Matx33d K(432.7390364738057, 0, 476.0614994349778, 0, 431.2395555913084, 288.7602152621297, 0, 0, 1);
6.     std::vector<double> dist_coeff = { -0.2852754904152874, 0.1016466459919075, ... };
7.     cv::Size board_pattern(10, 7);
8.     double board_cellsize = 0.025;
9.
10.    // Open a video
11.    cv::VideoCapture video;
12.    if (!video.open(input)) return -1;
13.
14.    // Prepare a 3D box for simple AR
15.    std::vector<cv::Point3d> box_lower = { cv::Point3d(4 * board_cellsize, 2 * board_cellsize, 0), ... };
16.    std::vector<cv::Point3d> box_upper = { cv::Point3d(4 * board_cellsize, 2 * board_cellsize, -board_cellsize), ... } ;
17.
18.    // Prepare 3D points on a chessboard
19.    std::vector<cv::Point3d> obj_points;
20.    for (int r = 0; r < board_pattern.height; r++)
21.        for (int c = 0; c < board_pattern.width; c++)
22.            obj_points.push_back(cv::Point3d(board_cellsize * c, board_cellsize * r, 0));
```





```
31. // Run pose estimation
32. while (true)
33. {
34.     // Grab an image from the video
35.     cv::Mat image;
36.     video >> image;
37.     if (image.empty()) break;

38.     // Estimate camera pose
39.     std::vector<cv::Point2d> img_points;
40.     bool success = cv::findChessboardCorners(image, board_pattern, img_points, ...);
41.     if (success)
42.     {
43.         cv::Mat rvec, tvec;
44.         cv::solvePnP(obj_points, img_points, K, dist_coeff, rvec, tvec);
45.             X1, X2, ..., X1, X2, ...
46.         // Draw the box on the image
47.         cv::Mat line_lower, line_upper;
48.         cv::projectPoints(box_lower, rvec, tvec, K, dist_coeff, line_lower);
49.         cv::projectPoints(box_upper, rvec, tvec, K, dist_coeff, line_upper);
50.         ...

51.         // Print camera position
52.         cv::Mat R;
53.         cv::Rodrigues(rvec, R);
54.         cv::Mat p = -R.t() * tvec;
55.         cv::String info = cv::format("XYZ: [% .3f, % .3f, % .3f]", cv::Point3d(p));
56.         cv::putText(image, info, cv::Point(5, 15), cv::FONT_HERSHEY_PLAIN, 1, cv::Vec3b(0, 255, 0));
57.     }

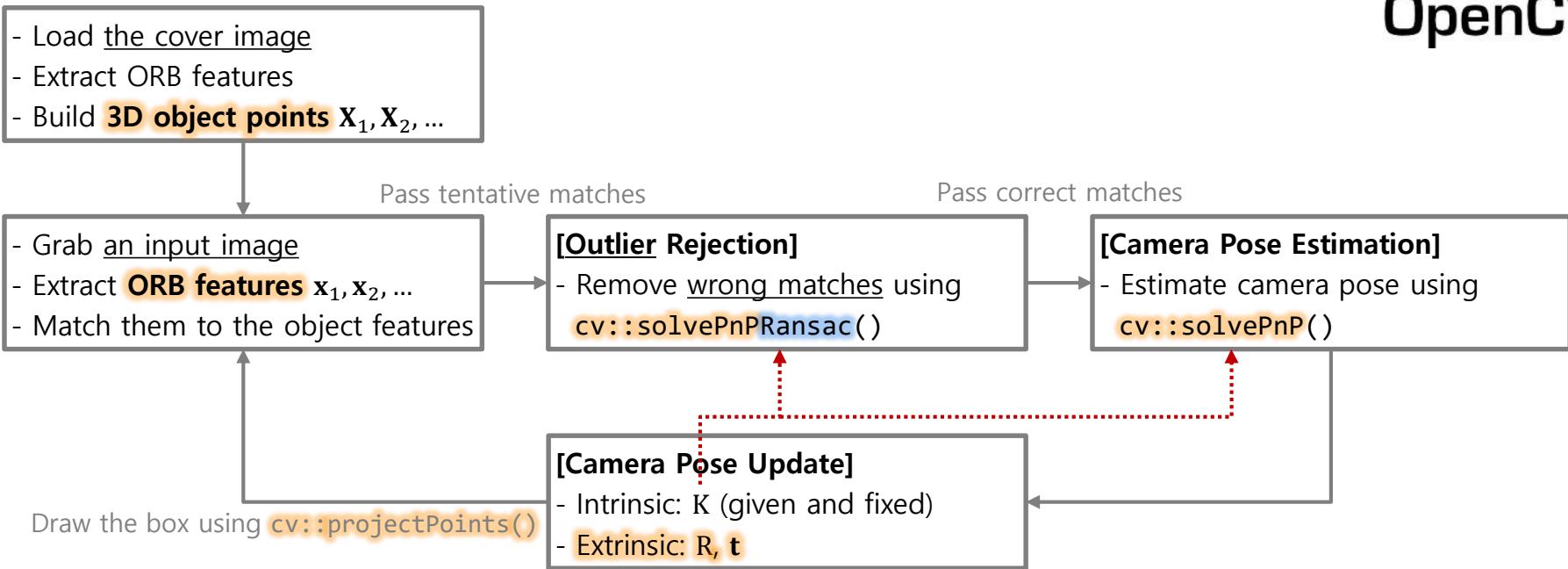
58.     // Show the image
59.     cv::imshow("3DV Tutorial: Pose Estimation", image);
60.     int key = cv::waitKey(1);
61.     ...

62.     video.release();
63.     return 0;
64. }
```

General 2D-3D Geometry



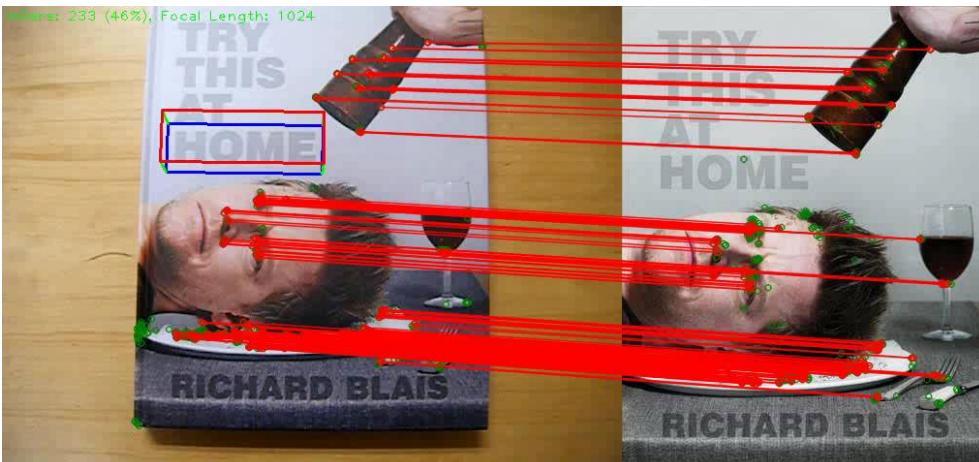
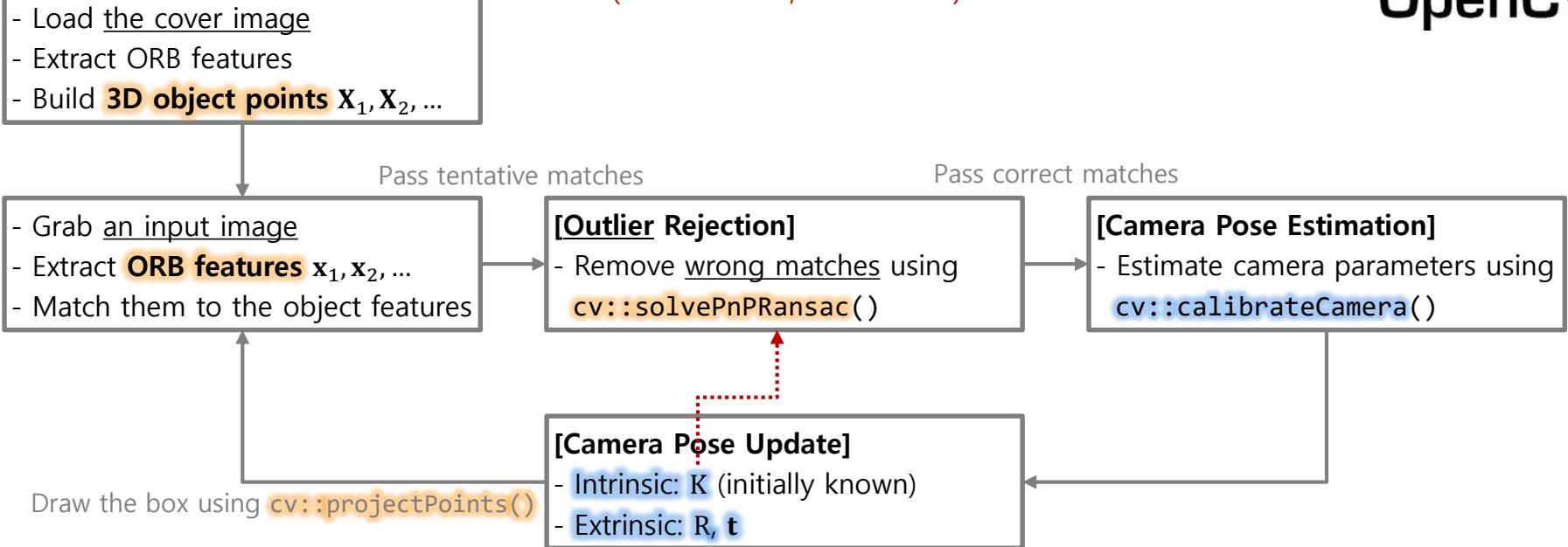
- Example: **Pose Estimation (Book)** [pose_estimation_book1.cpp]



General 2D-3D Geometry



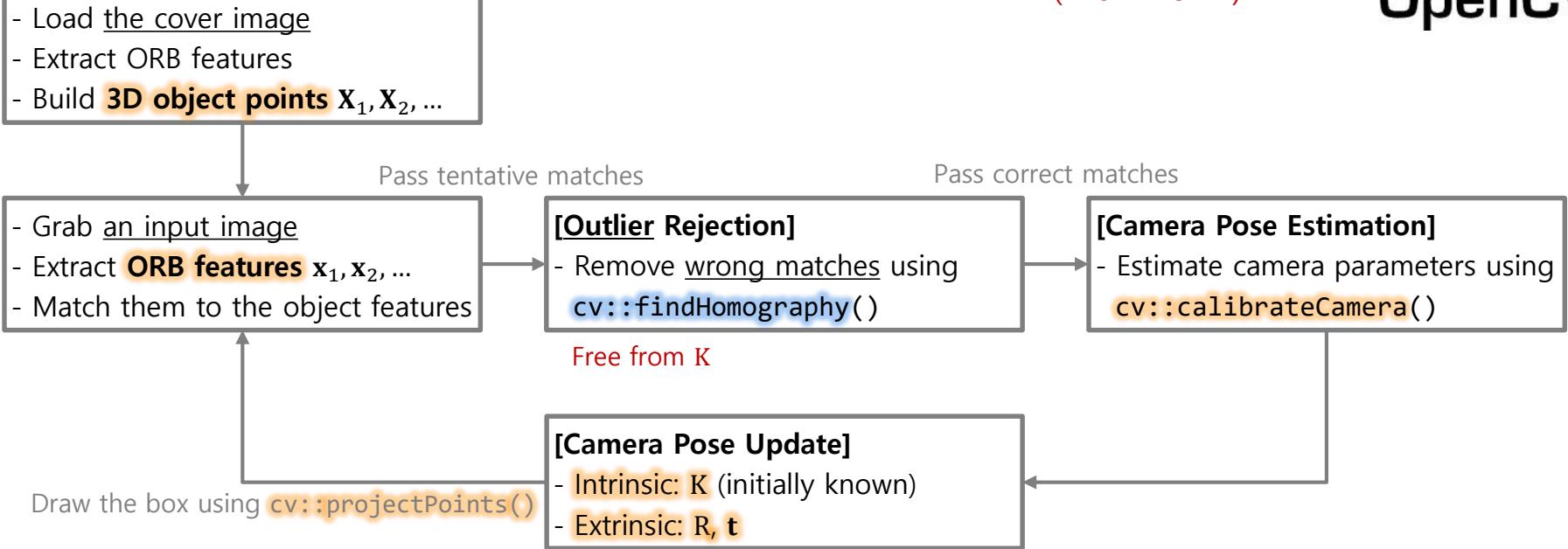
- Example: **Pose Estimation (Book) + Camera Calibration** [pose_estimation_book2.cpp]
(\because unknown, autofocus)



General 2D-3D Geometry

[pose_estimation_book3.cpp]

- Example: **Pose Estimation (Book) + Camera Calibration – Initially Given K**
(\therefore unknown)



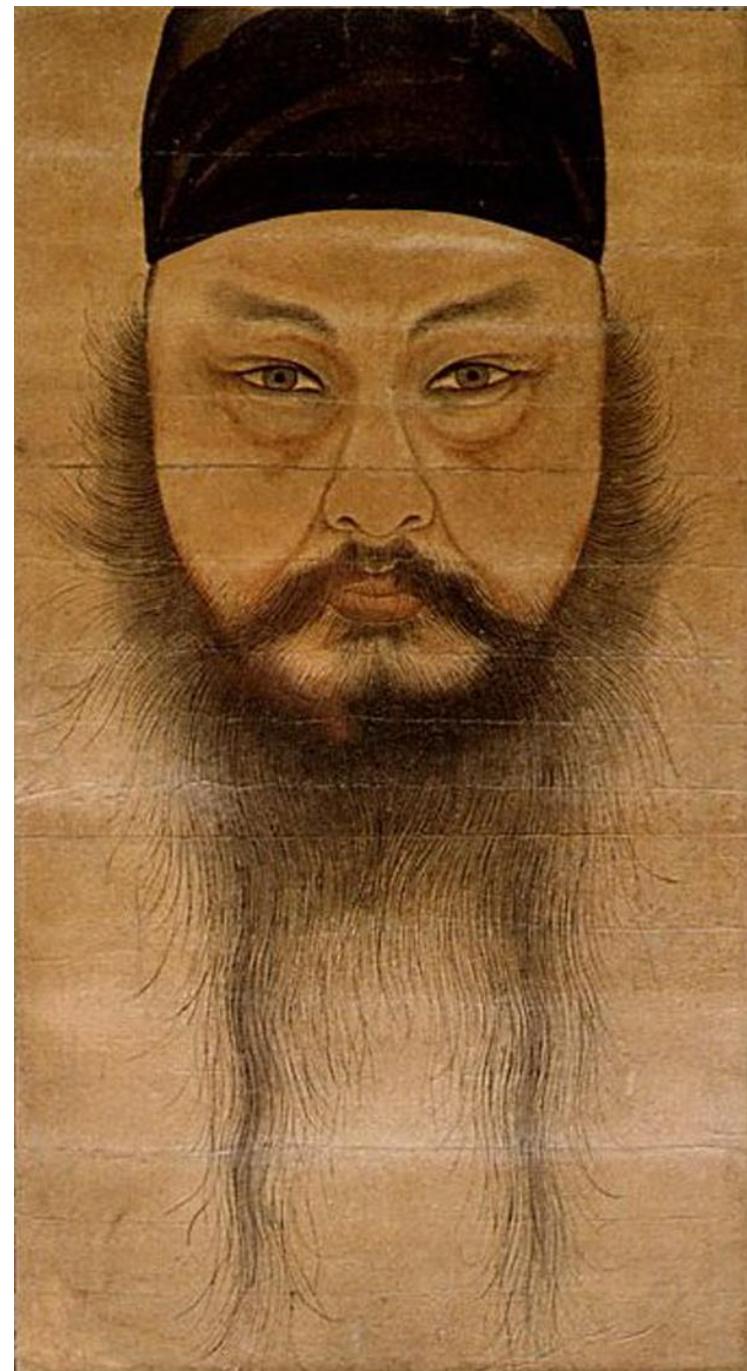
Two-view Geometry

Planar 2D-2D Geometry (Projective Geometry)

- Planar Homography Estimation

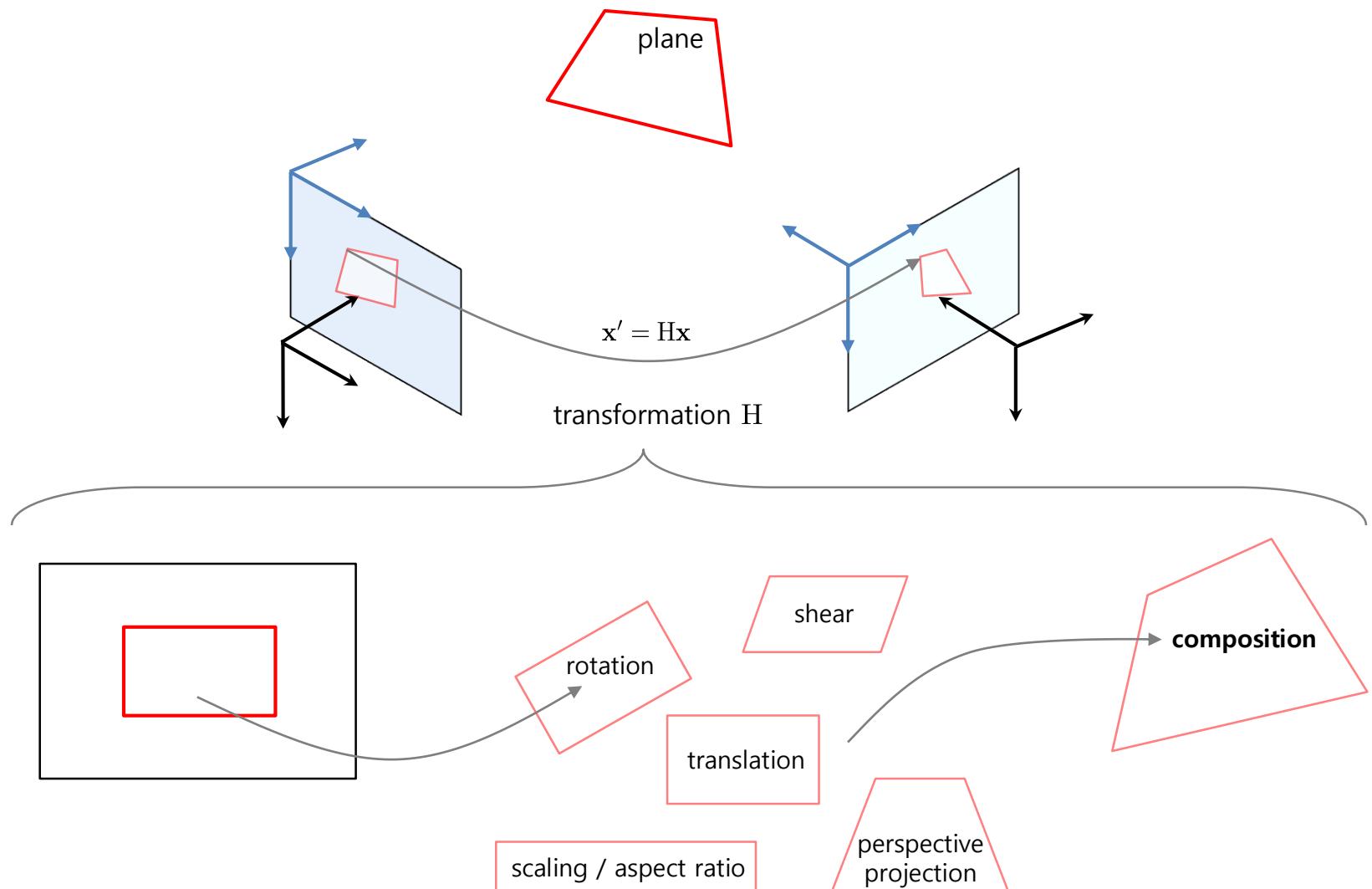
General 2D-2D Geometry (Epipolar Geometry)

- Relative Camera Pose Estimation
- Triangulation



윤두서(1668-1715) 자화상, 국보 제240호
Korean National Treasure No. 240

Planar 2D-2D Geometry (Projective Geometry)



Overview of Projective Geometry

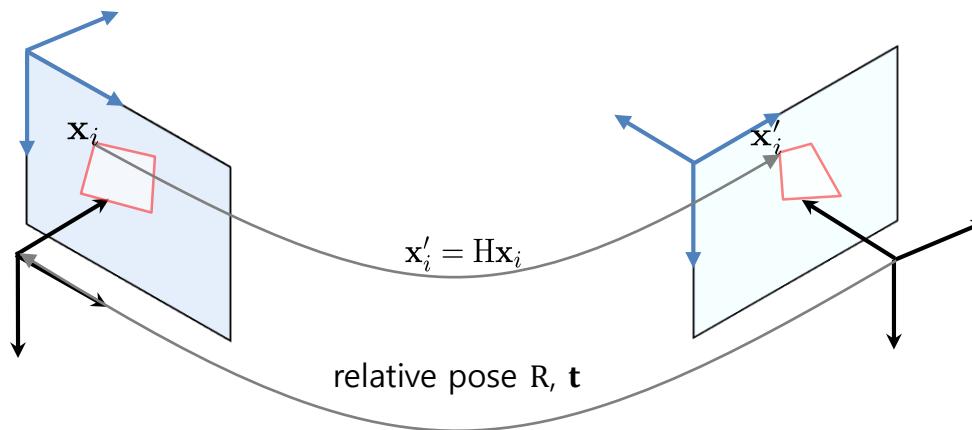
	Euclidean Transform (a.k.a. Rigid Transform)	Similarity Transform	Affine Transform	Projective Transform (a.k.a. Planar Homography)
Matrix Forms H	$\begin{bmatrix} \cos \theta & -\sin \theta & t_x \\ \sin \theta & \cos \theta & t_y \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} s \cos \theta & -s \sin \theta & t_x \\ s \sin \theta & s \cos \theta & t_y \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} a_{11} & a_{12} & t_x \\ a_{21} & a_{22} & t_y \\ 0 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} a_{11} & a_{12} & t_x \\ a_{21} & a_{22} & t_y \\ v_1 & v_2 & 1 \end{bmatrix}$
DoF	3	4	6	8
Transformations	O O X X X X	O O O X X X	O O O O O X	O O O O O O
Invariants	O O O O O O	X O O O O O	X X X O O O	X X X X O O
OpenCV Functions			<code>cv::getAffineTransform()</code> <code>cv::estimateRigidTransform()</code> - <code>cv::warpAffine()</code>	<code>cv:: getPerspectiveTransform()</code> - <code>cv::findHomography()</code> <code>cv::warpPerspective()</code>

cf. Similarly **3D transformations** (3D-3D geometry) are represented as **4x4 matrices**.

Planar 2D-2D Geometry (Projective Geometry)

▪ Planar Homography Estimation

- Unknown: Planar homography (8 DoF)
- Given: Point correspondence $(\mathbf{x}_1, \mathbf{x}'_1), \dots, (\mathbf{x}_n, \mathbf{x}'_n)$
- Constraints: $n \times$ projective transformation $\mathbf{x}'_i = \mathbf{H}\mathbf{x}_i$
- Solutions ($n \geq 4$) → 4-point algorithm
 - OpenCV `cv::getPerspectiveTransform()` and `cv::findHomography()`
 - cf. More simplified transformations need less number of minimal correspondence.
 - Affine ($n \geq 3$), similarity ($n \geq 2$), Euclidean ($n \geq 2$)
- [Note] Planar homography can be decomposed as relative camera pose.
 - OpenCV `cv::decomposeHomographyMat()`
 - cf. However, the decomposition needs to know camera matrices.



Planar 2D-2D Geometry (Projective Geometry)

- Example: Perspective Distortion Correction [perspective_correction.cpp]





```
1. #include <opencv2/opencv.hpp>

2. void MouseEventHandler(int event, int x, int y, int flags, void* param)
3. {
4.     ...
5.     std::vector<cv::Point> *points_src = (std::vector<cv::Point> *)param;
6.     points_src->push_back(cv::Point(x, y));
7.     ...
8. }

9. int main()
10.{
11.    const char* input = "data/sunglok_desk.jpg";
12.    cv::Size card_size(450, 250);

13.    // Prepare the rectified points
14.    std::vector<cv::Point> points_dst; x'_1, x'_2, ...
15.    points_dst.push_back(cv::Point(0, 0));
16.    points_dst.push_back(cv::Point(card_size.width, 0));
17.    points_dst.push_back(cv::Point(0, card_size.height));
18.    points_dst.push_back(cv::Point(card_size.width, card_size.height));

19.    // Load an image
20.    cv::Mat original = cv::imread(input);
21.    if (original.empty()) return -1;

22.    // Get the matched points from a user's mouse
23.    std::vector<cv::Point> points_src; X_1, X_2, ...
24.    cv::namedWindow("3DV Tutorial: Perspective Correction");
25.    cv::setMouseCallback("3DV Tutorial: Perspective Correction", MouseEventHandler, &points_src);
26.    while (points_src.size() < 4)
27.    {
28.        ...
29.    }
30.    if (points_src.size() < 4) return -1;

31.    // Calculate planar homography and rectify perspective distortion
32.    cv::Mat H = cv::findHomography(points_src, points_dst);
33.    cv::Mat rectify;
34.    cv::warpPerspective(original, rectify, H, card_size);

35.    // Show the rectified image
36.    ...
37. }
```

Planar 2D-2D Geometry (Projective Geometry)



- Example: Planar Image Stitching [image_stitching.cpp]

```
1. #include "opencv2/opencv.hpp"  
2. int main()  
3. {  
4.     // Load two images (cf. We assume that two images have the same size and type)  
5.     cv::Mat image1 = cv::imread("data/hill01.jpg");  
6.     cv::Mat image2 = cv::imread("data/hill02.jpg");  
7.     if (image1.empty() || image2.empty()) return -1;  
  
8.     // Retrieve matching points  
9.     cv::Ptr<cv::FeatureDetector> fdetector = cv::ORB::create();  
10.    std::vector<cv::KeyPoint> keypoint1, keypoint2;  
11.    cv::Mat descriptor1, descriptor2;  
12.    fdetector->detectAndCompute(image1, cv::Mat(), keypoint1, descriptor1);  
13.    fdetector->detectAndCompute(image2, cv::Mat(), keypoint2, descriptor2);  
14.    cv::Ptr<cv::DescriptorMatcher> fmatcher = cv::DescriptorMatcher::create("BruteForce-Hamming");  
15.    std::vector<cv::DMatch> match;  
16.    fmatcher->match(descriptor1, descriptor2, match);  
  
17.    // Calculate planar homography and merge them  
18.    std::vector<cv::Point2f> points1, points2;  
19.    for (size_t i = 0; i < match.size(); i++)  
20.    {  
21.        points1.push_back(keypoint1.at(match.at(i).queryIdx).pt);  
22.        points2.push_back(keypoint2.at(match.at(i).trainIdx).pt);  
23.    }  
24.    cv::Mat H = cv::findHomography(points2, points1, cv::RANSAC);  
25.    cv::Mat merge;  
26.    cv::warpPerspective(image2, merge, H, cv::Size(image1.cols * 2, image1.rows));  
27.    merge.colRange(0, image1.cols) = image1 * 1; // Copy  
  
28.    // Show the merged image  
29.    ...  
30. }
```



Planar 2D-2D Geometry (Projective Geometry)



- Example: 2D Video Stabilization [video_stabilization.cpp]

```
1. #include "opencv2/opencv.hpp"  
2. int main()  
3. {  
4.     // Open a video and get the reference image and feature points  
5.     cv::VideoCapture video;  
6.     if (!video.open("data/traffic.avi")) return -1;  
7.  
8.     cv::Mat gray_ref;  
9.     video >> gray_ref;  
10.    if (gray_ref.empty())  
11.    {  
12.        video.release();  
13.        return -1;  
14.    }  
15.    if (gray_ref.channels() > 1) cv::cvtColor(gray_ref, gray_ref, cv::COLOR_RGB2GRAY);  
16.    std::vector<cv::Point2f> point_ref;  
17.    cv::goodFeaturesToTrack(gray_ref, point_ref, 2000, 0.01, 10);  
18.    if (point_ref.size() < 4)  
19.    {  
20.        video.release();  
21.        return -1;  
22.    }
```



A shaking CCTV video: "data/traffic.avi"



```
22.     // Run and show video stabilization
23.     while (true)
24.     {
25.         // Grab an image from the video
26.         cv::Mat image, gray;
27.         video >> image;
28.         if (image.empty()) break;
29.         if (image.channels() > 1) cv::cvtColor(image, gray, cv::COLOR_RGB2GRAY);
30.         else
31.             gray = image.clone();

32.         // Extract optical flow and calculate planar homography
33.         std::vector<cv::Point2f> point;
34.         std::vector<uchar> m_status;
35.         cv::Mat err, inlier_mask;
36.         cv::calcOpticalFlowPyrLK(gray_ref, gray, point_ref, point, m_status, err);
37.         cv::Mat H = cv::findHomography(point, point_ref, inlier_mask, cv::RANSAC);

38.         // Synthesize a stabilized image
39.         cv::Mat warp;
40.         cv::warpPerspective(image, warp, H, cv::Size(image.cols, image.rows));

41.         // Show the original and rectified images together
42.         for (int i = 0; i < point_ref.size(); i++)
43.         {
44.             if (inlier_mask.at<uchar>(i) > 0) cv::line(image, point_ref[i], point[i], cv::Vec3b(0, 0, 255));
45.             else cv::line(image, point_ref[i], point[i], cv::Vec3b(0, 127, 0));
46.         }
47.         cv::hconcat(image, warp, image);
48.         cv::imshow("3DVT Tutorial: Video Stabilization", image);
49.         if (cv::waitKey(1) == 27) break; // 'ESC' key
50.     }

51.     video.release();
52.     return 0;
53. }
```

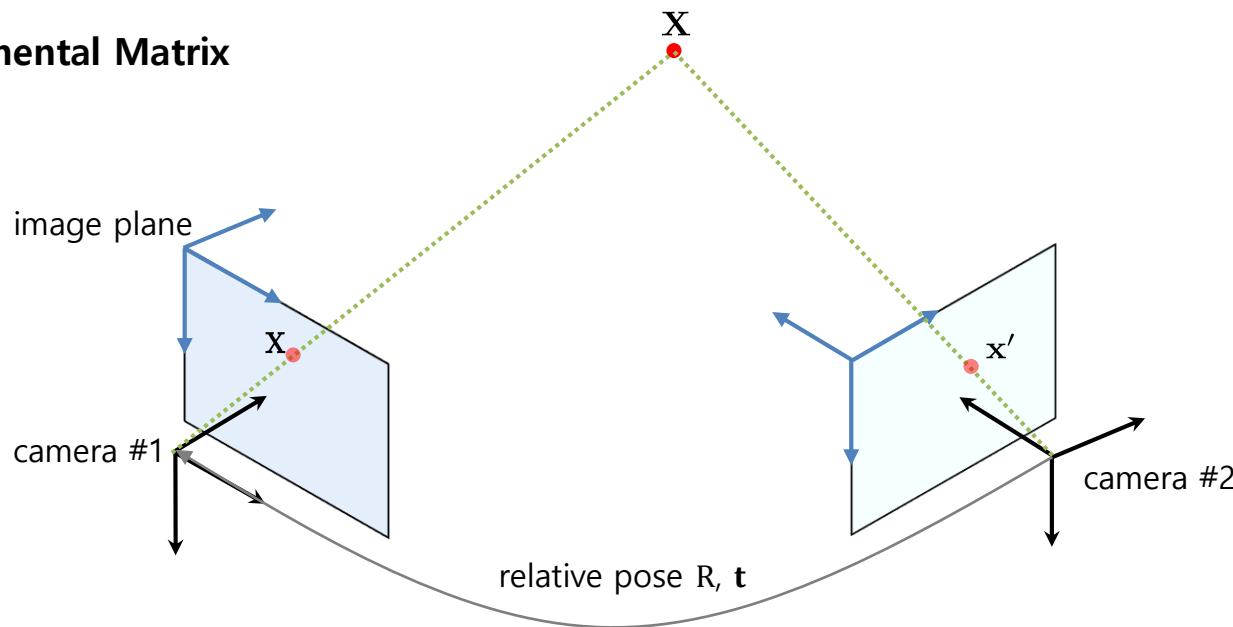
Assumption: **A plane observed by two views**

- Image stitching: Distance \gg depth variation
- Video stabilization: Small motion



General 2D-2D Geometry (Epipolar Geometry)

- **Fundamental Matrix**



Epipolar Constraint:

$$x'^\top F x = 0$$

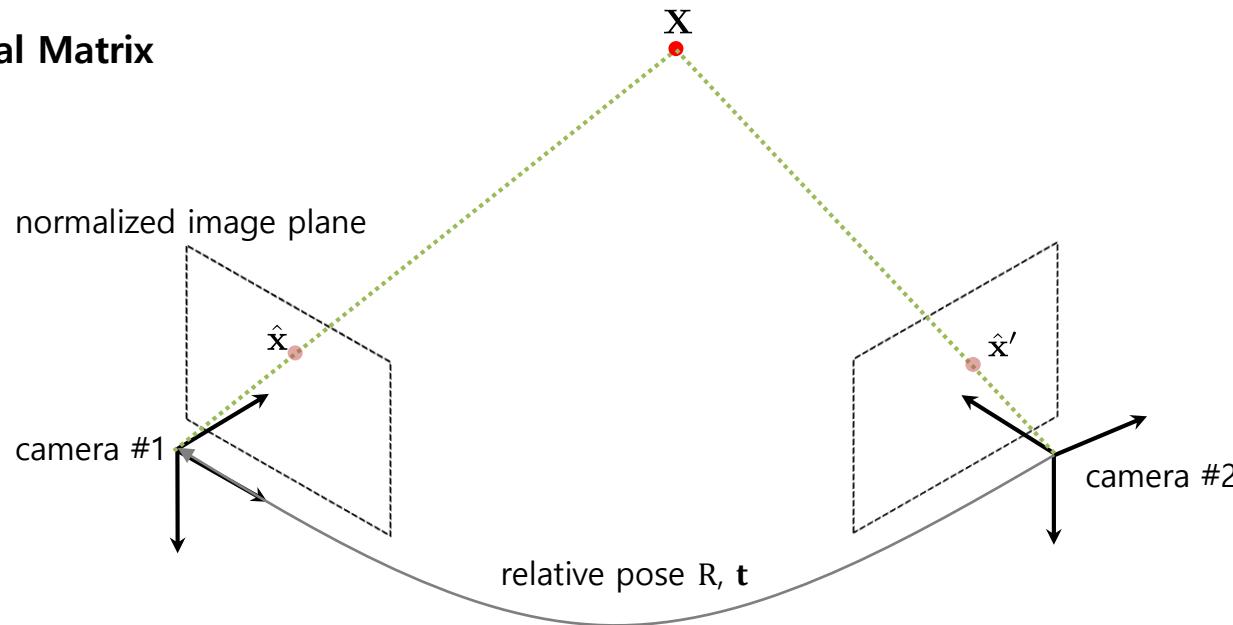
on the image plane

$$\begin{array}{c} \downarrow \\ x = K\hat{x} \\ \downarrow \\ \hat{x}'^\top K'^\top F K \hat{x} = 0 \\ \downarrow \\ \hat{x}'^\top E \hat{x} = 0 \quad (E = K'^\top F K) \end{array}$$

on the normalized image plane

General 2D-2D Geometry (Epipolar Geometry)

- Essential Matrix



Epipolar Constraint:
(Derivation)

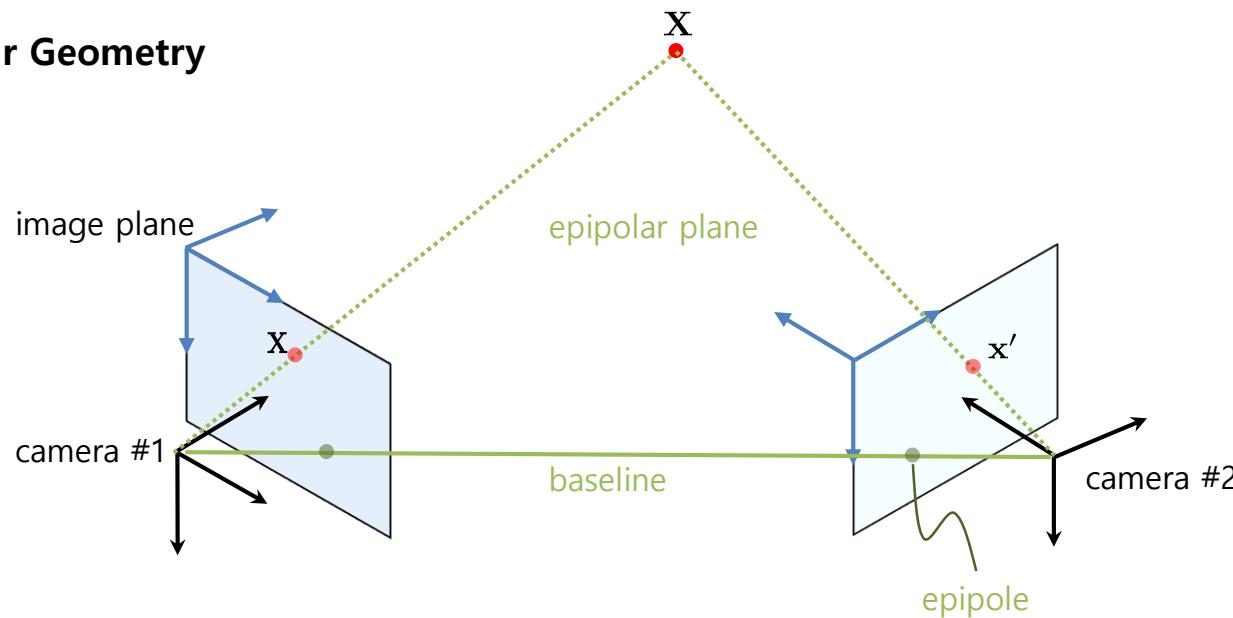
$$\hat{\mathbf{x}}'^\top E \hat{\mathbf{x}} = 0 \quad (E = [\mathbf{t}]_\times R)$$

c.f. $[\mathbf{t}]_\times = \begin{bmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ -t_y & t_x & 0 \end{bmatrix}$

$$\begin{aligned}
 & \hat{\mathbf{x}}' \cdot \mathbf{t} \times R \hat{\mathbf{x}} = 0 \\
 & \hat{\mathbf{x}}' \cdot (\mathbf{t} \times \hat{\mathbf{x}}') = 0 \\
 & \lambda' \hat{\mathbf{x}}' \cdot (\mathbf{t} \times \hat{\mathbf{x}}') = \lambda \hat{\mathbf{x}}' \cdot \mathbf{t} \times R \hat{\mathbf{x}} \\
 & \hat{\mathbf{x}}' \cdot \\
 & \lambda' \mathbf{t} \times \hat{\mathbf{x}}' = \lambda \mathbf{t} \times R \hat{\mathbf{x}} \\
 & \mathbf{t} \times \\
 & \lambda' \hat{\mathbf{x}}' = \lambda R \hat{\mathbf{x}} + \mathbf{t} \quad \xleftarrow{X = \lambda \hat{\mathbf{x}}} \quad \mathbf{X}' = R \mathbf{X} + \mathbf{t}
 \end{aligned}$$

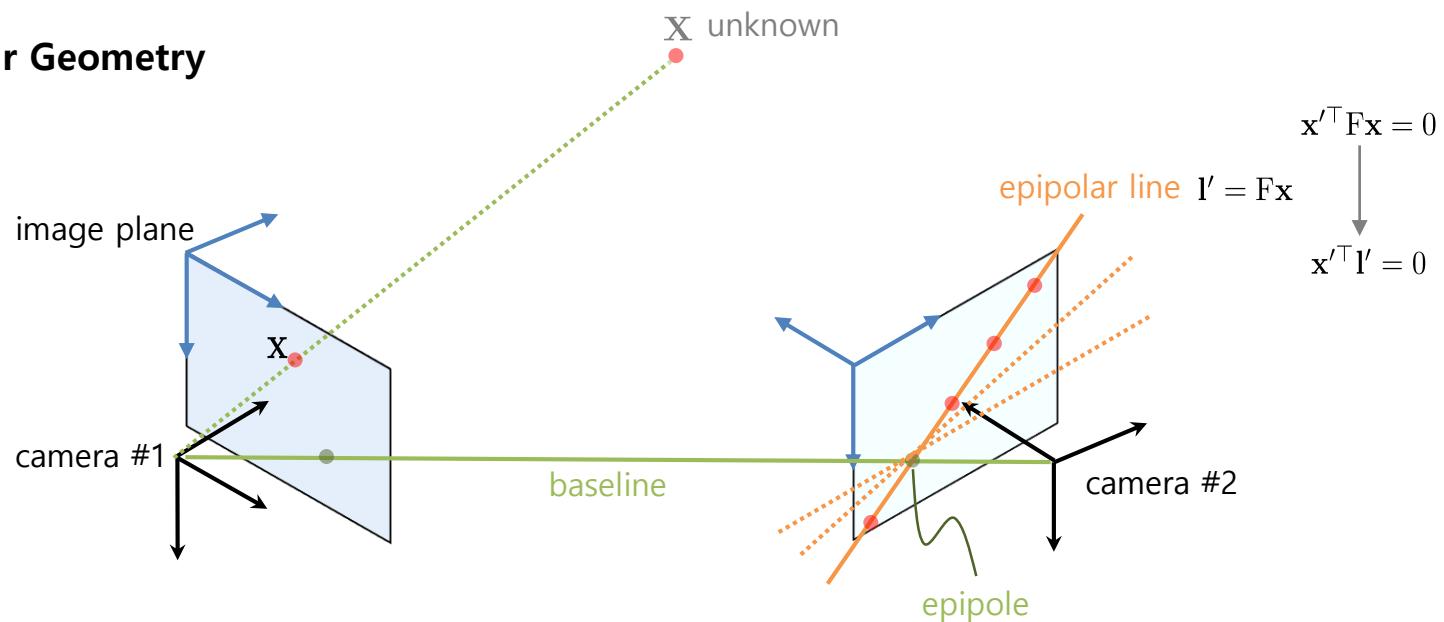
General 2D-2D Geometry (Epipolar Geometry)

- Epipolar Geometry

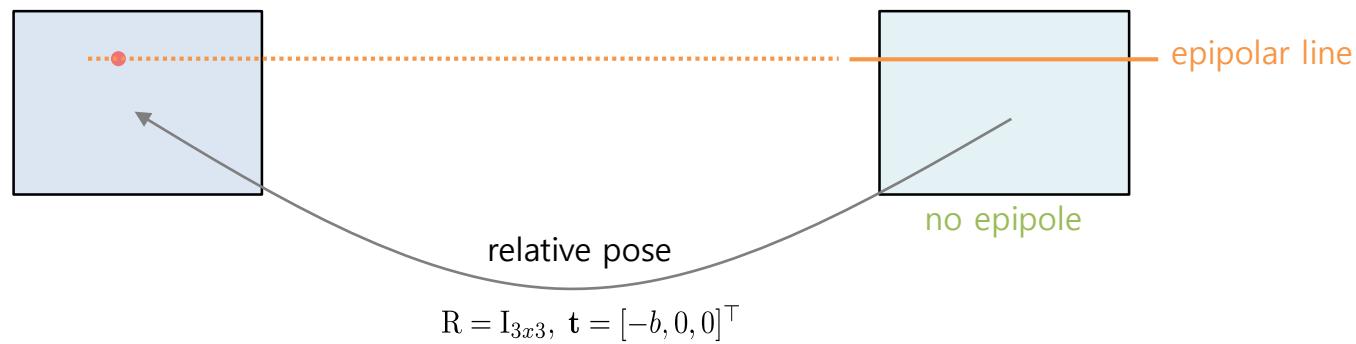


General 2D-2D Geometry (Epipolar Geometry)

- Epipolar Geometry



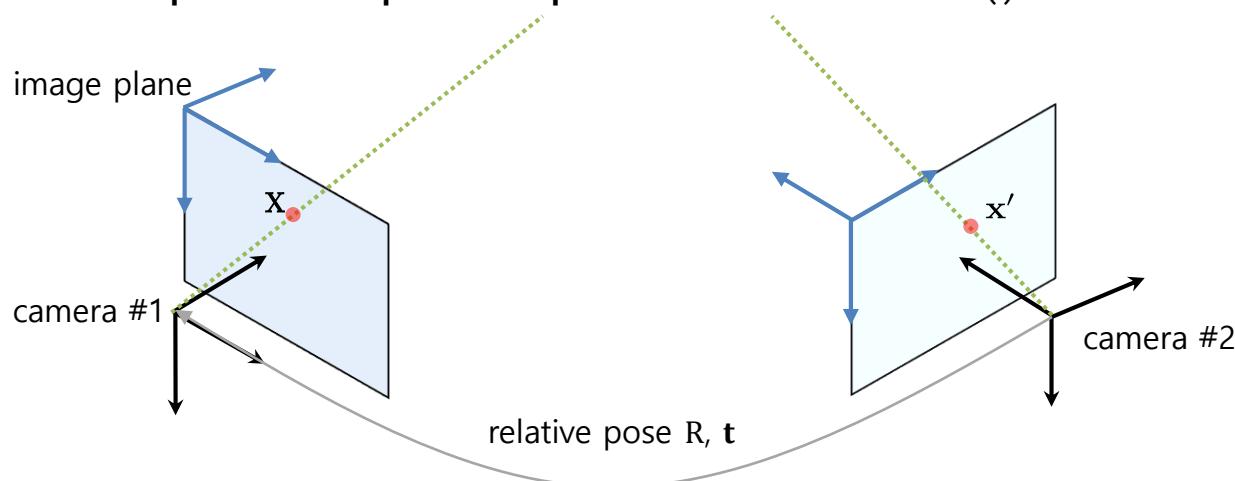
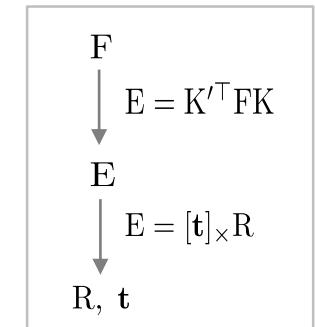
Special case) Stereo Cameras



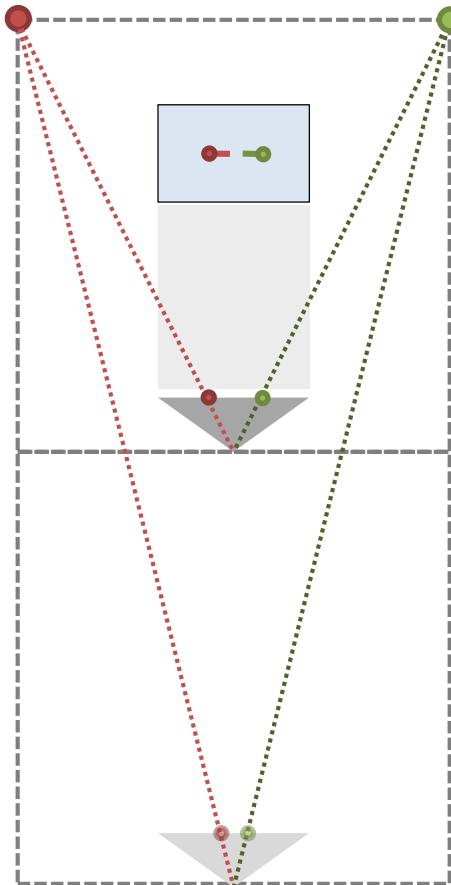
General 2D-2D Geometry (Epipolar Geometry)

- **Relative Camera Pose Estimation (~ Fundamental/Essential Matrix Estimation)**

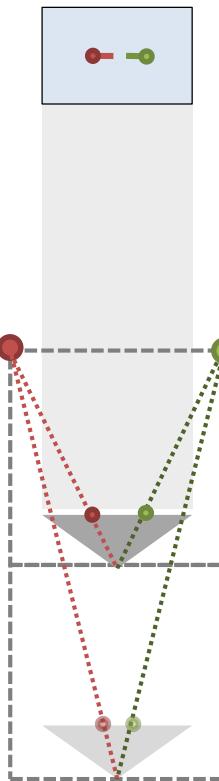
- Unknown: **Rotation and translation R, t** (**5 DoF**; up-to scale “**scale ambiguity**”)
- Given: Point correspondence $(\mathbf{x}_1, \mathbf{x}'_1), \dots, (\mathbf{x}_n, \mathbf{x}'_n)$ and camera matrices K, K'
- Constraints: $n \times$ epipolar constraint $(\mathbf{x}'^\top F \mathbf{x} = 0 \text{ or } \hat{\mathbf{x}}'^\top E \hat{\mathbf{x}} = 0)$
- Solutions (OpenCV)
 - **Fundamental matrix:** 7/8-point algorithm (7 DoF)
 - **Estimation:** `cv::findFundamentalMat()` → 1 solution
 - **Conversion to E:** $E = K'^\top F K$
 - **Essential matrix:** 5-point algorithm (5 DoF)
 - **Estimation:** `cv::findEssentialMat()` → k solutions
 - **Decomposition:** `cv::decomposeEssentialMat()` → 4 solutions “**relative pose ambiguity**”
 - **Decomposition with positive-depth check:** `cv::recoverPose()` → 1 solution



Epipolar Geometry: Scale Ambiguity



(a) 2-meter-wide tunnel



(b) 1-meter-wide tunnel

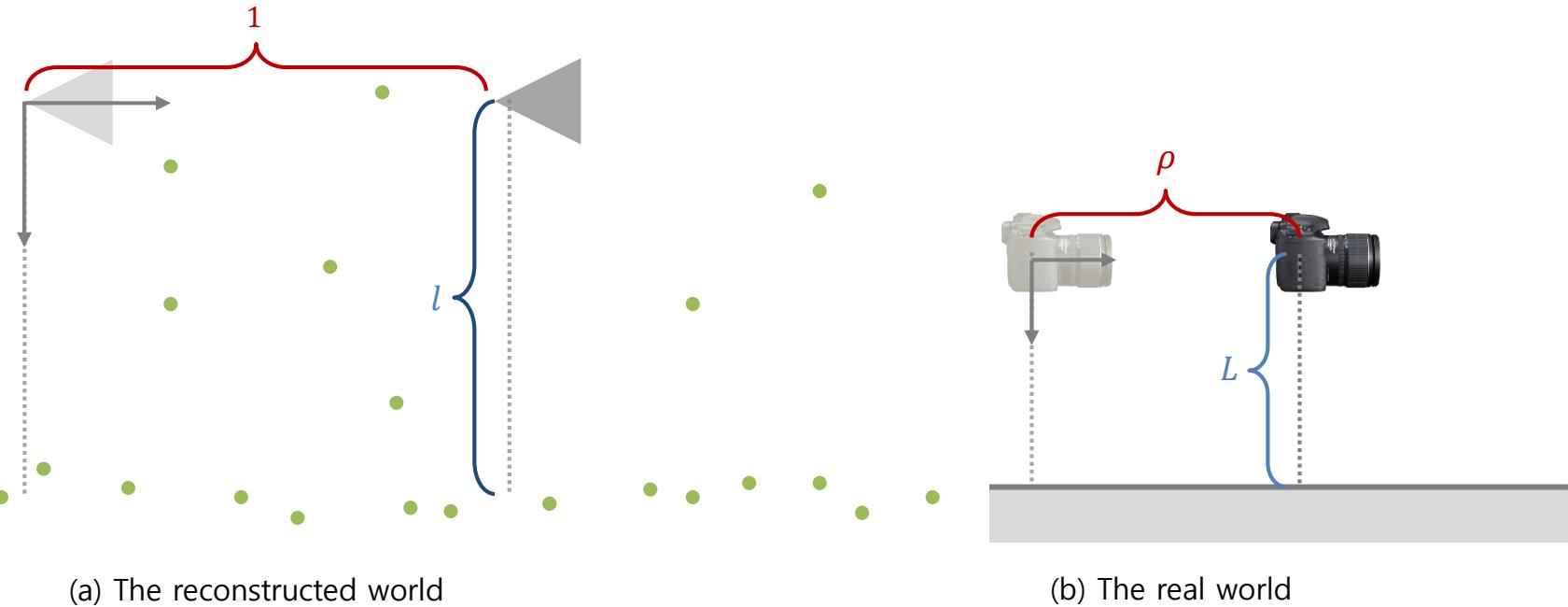
Epipolar Geometry: Scale Ambiguity



The Sandcrawler @ Star Wars IV: A New Hope (1977)



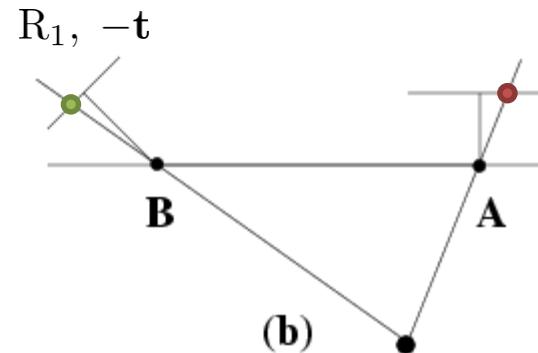
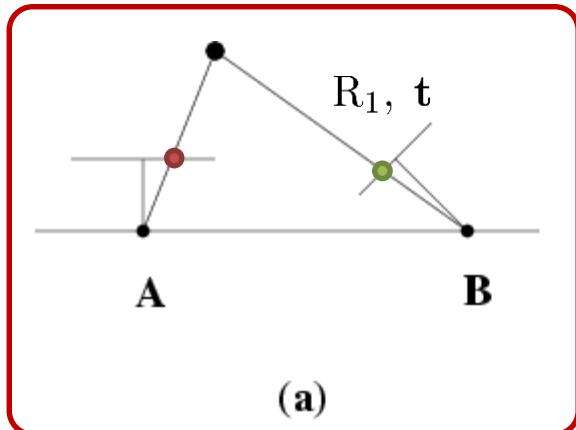
Epipolar Geometry: Resolving Scale Ambiguity



How to resolve scale ambiguity

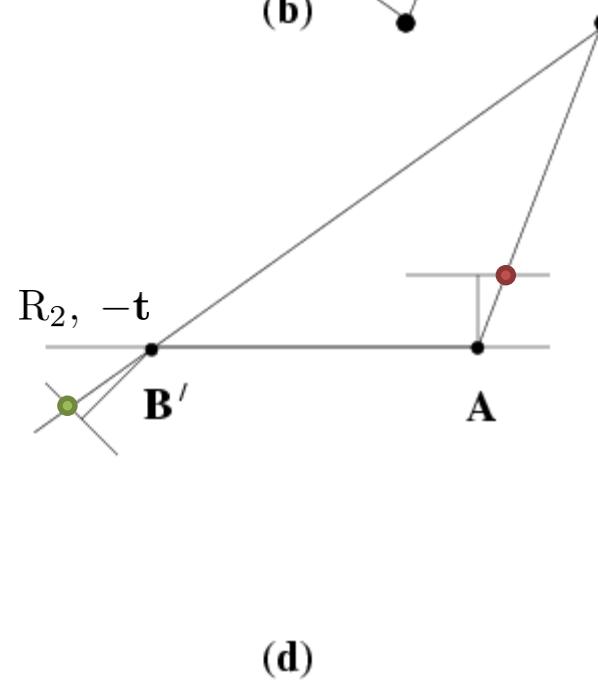
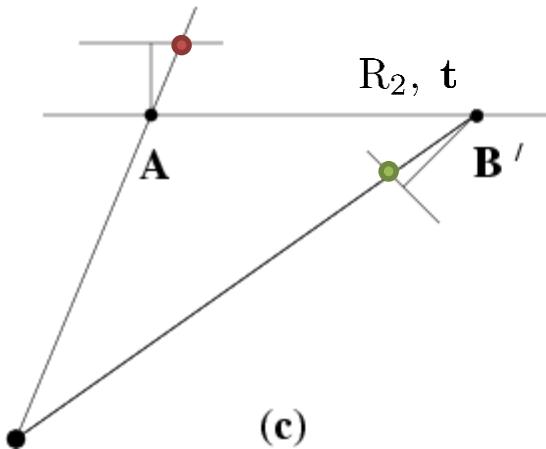
- Additional **sensors**: speedometers, IMUs, GPSs, depth/distance (stereo, RGB-D, LiDAR, ...)
- More **constraints**: known size of objects, **known and fixed height of camera** ($\rho = \frac{L}{l}$)

Epipolar Geometry: Relative Pose Ambiguity



How to resolve pose ambiguity

- Positive depth constraint



General 2D-2D Geometry (Epipolar Geometry)

▪ Relative Camera Pose Estimation

- Unknown: Rotation and translation R, t (**5 DoF**; up-to scale)
- Given: Point correspondence $(\mathbf{x}_1, \mathbf{x}'_1), \dots, (\mathbf{x}_n, \mathbf{x}'_n)$ and camera matrices K, K'
- Constraints: $n \times$ epipolar constraint $(\mathbf{x}'^\top F \mathbf{x} = 0 \text{ or } \hat{\mathbf{x}}'^\top E \hat{\mathbf{x}} = 0)$
- Solutions (OpenCV)

intrinsic & extrinsic
camera parameters

Fundamental matrix: 7/8-point algorithm (**7 DoF**)

$$F = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & 1 \end{bmatrix} \quad \& \quad \text{rank}(F) = 2 \quad \text{det}(F) = 0$$

- **Estimation:** `cv::findFundamentalMat()` $\rightarrow 1$ solution
- **Conversion to E:** $E = K'^\top F K$
- **Degenerate cases:** No translation, correspondence from a single plane

extrinsic
camera parameters

Essential matrix: 5-point algorithm (**5 DoF**)

$$2EE^\top E - \text{tr}(EE^\top)E = 0$$

- **Estimation:** `cv::findEssentialMat()` $\rightarrow k$ solutions
- **Decomposition:** `cv::decomposeEssentialMat()` $\rightarrow 4$ solutions
- **Decomposition with positive-depth check:** `cv::recoverPose()` $\rightarrow 1$ solution
- **Degenerate cases:** No translation ($\because E = [t]_\times R$)

General 2D-2D Geometry (Epipolar Geometry)

▪ Relative Camera Pose Estimation

- Unknown: Rotation and translation R, t (**5 DoF**; up-to scale)
- Given: Point correspondence $(\mathbf{x}_1, \mathbf{x}'_1), \dots, (\mathbf{x}_n, \mathbf{x}'_n)$ and camera matrices K, K'
- Constraints: $n \times$ epipolar constraint $(\mathbf{x}'^\top F \mathbf{x} = 0 \text{ or } \hat{\mathbf{x}}'^\top E \hat{\mathbf{x}} = 0)$
- Solutions (OpenCV)

$$F = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & 1 \end{bmatrix} \quad \& \quad \text{rank}(F) = 2 \quad \text{det}(F) = 0$$

**intrinsic & extrinsic
camera parameters**

• **Fundamental matrix:** 7/8-point algorithm (**7 DoF**)

- **Estimation:** `cv::findFundamentalMat()` → 1 solution
- **Conversion to E:** $E = K'^\top F K$
- **Degenerate cases:** No translation, correspondence from a single plane

**extrinsic
camera parameters**

• **Essential matrix:** 5-point algorithm (**5 DoF**)

$$2EE^\top E - \text{tr}(EE^\top)E = 0$$

- **Estimation:** `cv::findEssentialMat()` → k solutions
- **Decomposition:** `cv::decomposeEssentialMat()` → 4 solutions
- **Decomposition with positive-depth check:** `cv::recoverPose()` → 1 solution
- **Degenerate cases:** No translation ($\because E = [t]_\times R$)

**intrinsic & extrinsic
camera parameters
+ plane normal**

• **Planar homography:** 4-point algorithm (**8 DoF**)

- **Estimation:** `cv::findHomography()` → 1 solutions
- **Conversion to calibrated H:** $\hat{H} = K'^{-1} H K$
- **Decomposition:** `cv::decomposeHomographyMat()` → 4 solutions
- **Degenerate cases:** Correspondence **not** from a single plane

$$H = K' \hat{H} K^{-1}$$

$$\hat{H} = \lambda'' \left(R + \frac{1}{d} t n^\top \right)$$

$$\hat{x}' = \lambda'' \left(R + \frac{1}{d} t n^\top \right) \hat{x} \quad \lambda' \hat{x}' = \lambda R \hat{x} + t$$

$$\frac{1}{d} n^\top \hat{x} = 1 \quad (\because n_x \hat{x} + n_y \hat{y} + n_z \hat{z} - d = 0)$$

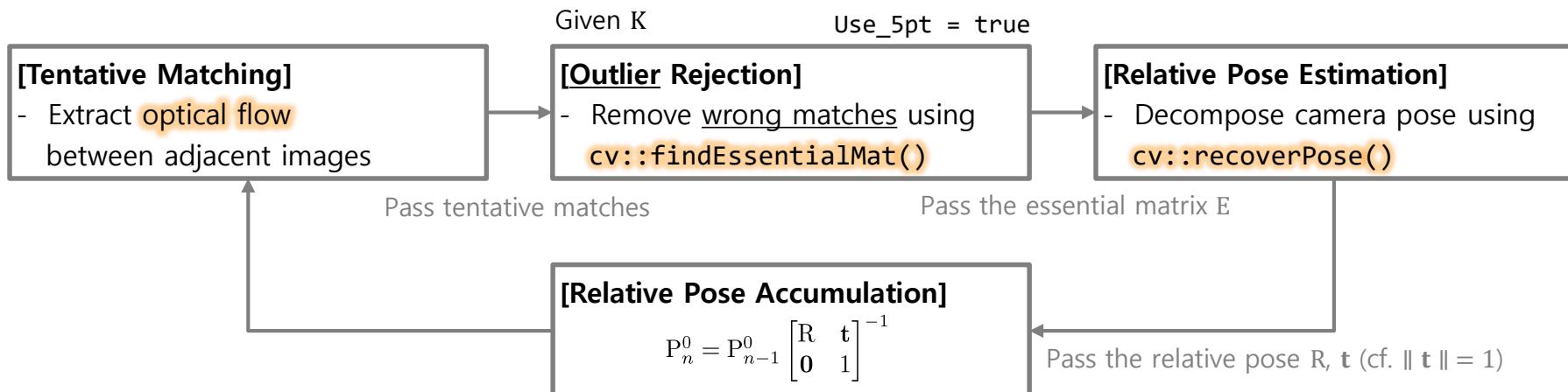
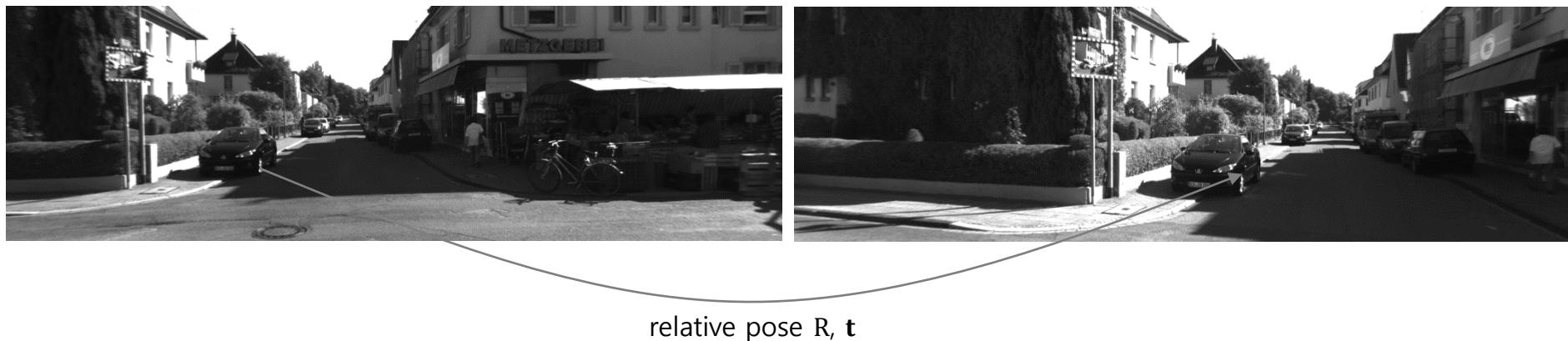
Overview of Epipolar Geometry

	Items	General 2D-2D Geometry	Planar 2D-2D Geometry
On Image Planes	Model	Fundamental Matrix (7 DoF)	Planar Homography (8 DoF)
	Formulation	$F = K'^{-\top} E K^{-1}$	$E = K'^{\top} F K$
	Estimation	<ul style="list-style-type: none"> - 7-point algorithm ($n \geq 7$) $\rightarrow k$ solution - (normalized) 8-point algorithm $\rightarrow 1$ solution - <code>cv::findFundamentalMat()</code> 	<ul style="list-style-type: none"> - 4-point algorithm ($n \geq 4$) $\rightarrow 1$ solution - <code>cv::findHomography()</code>
	Input	- (x_i, x'_i) [px] on the image plane	- (x_i, x'_i) [px] on a plane in the image plane
	Degenerate Cases	<ul style="list-style-type: none"> - no translational motion - correspondence from a single plane 	<ul style="list-style-type: none"> - correspondence not from a single plane
	Decomposition to R and t	- convert to an essential matrix and decompose it	- <code>cv::decomposeHomographyMat()</code>
On Normalized Image Planes	Model	Essentials Matrix (5 DoF)	(Calibrated) Planar Homography (8 DoF)
	Formulation	$E = [\mathbf{t}]_x R$	$\hat{H} = R + \frac{1}{d} \mathbf{t} n^T$
	Estimation	<ul style="list-style-type: none"> - 5-point algorithm ($n \geq 5$) $\rightarrow k$ solution - <code>cv::findEssentialMat()</code> 	<ul style="list-style-type: none"> - 4-point algorithm ($n \geq 4$) $\rightarrow 1$ solution - <code>cv::findHomography()</code>
	Input	- (\hat{x}_i, \hat{x}'_i) [m] on the normalized image plane	- (\hat{x}_i, \hat{x}'_i) [m] on a plane in the normalized image plane
	Degenerate Cases	<ul style="list-style-type: none"> - no translational motion 	<ul style="list-style-type: none"> - correspondence not from a single plane
	Decomposition to R and t	<ul style="list-style-type: none"> - <code>cv::decomposeEssentialMat()</code> - <code>cv::recoverPose()</code> 	<ul style="list-style-type: none"> - <code>cv::decomposeHomographyMat()</code> with $K = I_{3 \times 3}$

General 2D-2D Geometry (Epipolar Geometry)

- Example: Monocular Visual Odometry (Epipolar Version) [vo_epipolar.cpp]

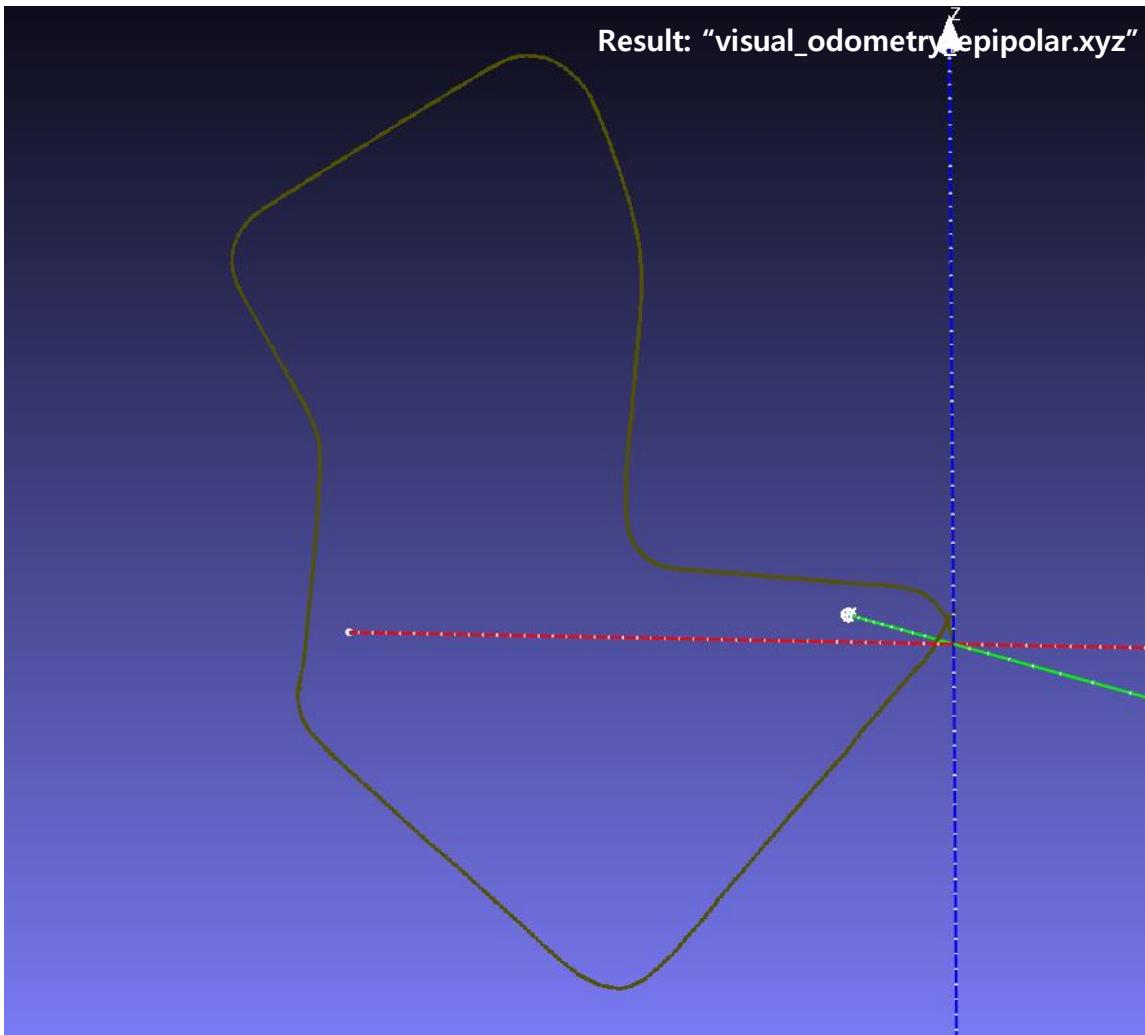
KITTI Odometry Dataset (#: 00, Left): "data/KITTI_07_L/%06d.png"



Inliers: 2 (0%), XYZ: [0.000, 0.000, 0.000]



Result: "visual_odometry.epipolar.xyz"





```
1. #include "opencv2/opencv.hpp"
2. int main()
3. {
4.     const char* input = "data/KITTI_07_L/%06d.png";
5.     double f = 707.0912;
6.     cv::Point2d c(601.8873, 183.1104);
7.     bool use_5pt = true;
8.     int min_inlier_num = 100;
9.     ...
10.    // Open a video and get the initial image
11.    cv::VideoCapture video;
12.    if (!video.open(input)) return -1;
13.
14.    cv::Mat gray_prev;
15.    video >> gray_prev;
16.    ...
17.
18.    // Run and record monocular visual odometry
19.    cv::Mat camera_pose = cv::Mat::eye(4, 4, CV_64F);
20.    while (true)
21.    {
22.        // Grab an image from the video
23.        cv::Mat image, gray;
24.        video >> image;
25.        if (image.empty()) break;
26.        if (image.channels() > 1) cv::cvtColor(image, gray, cv::COLOR_RGB2GRAY);
27.        else
28.            gray = image.clone();
29.
30.        // Extract optical flow
31.        std::vector<cv::Point2f> point_prev, point;
32.        cv::goodFeaturesToTrack(gray_prev, point_prev, 2000, 0.01, 10);
33.        std::vector<uchar> m_status;
34.        cv::Mat err;
35.        cv::calcOpticalFlowPyrLK(gray_prev, gray, point_prev, point, m_status, err);
36.        gray_prev = gray;
```

```

33.     // Calculate relative pose
34.     cv::Mat E, inlier_mask;
35.     if (use_5pt)
36.     {
37.         E = cv::findEssentialMat(point_prev, point, f, c, cv::RANSAC, 0.99, 1, inlier_mask);
38.     }
39.     else
40.     {
41.         cv::Mat F = cv::findFundamentalMat(point_prev, point, cv::FM_RANSAC, 1, 0.99, inlier_mask);
42.         cv::Mat K = (cv::Mat<double>(3, 3) << f, 0, c.x, 0, f, c.y, 0, 0, 1);
43.         E = K.t() * F * K;
44.     }
45.     cv::Mat R, t;
46.     int inlier_num = cv::recoverPose(E, point_prev, point, R, t, f, c, inlier_mask);

47.     // Accumulate relative pose if result is reliable
48.     if (inlier_num > min_inlier_num)
49.     {
50.         cv::Mat T = cv::Mat::eye(4, 4, R.type());
51.         T(cv::Rect(0, 0, 3, 3)) = R * 1.0;
52.         T.col(3).rowRange(0, 3) = t * 1.0;
53.         camera_pose = camera_pose * T.inv();    ::  $P_n^0 = P_{n-1}^0 \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix}^{-1}$ 
54.     }

55.     // Show the image and write camera pose
56.     ...
57. }

58. video.release();
59. fclose(camera_traj);
60. return 0;
61. }

```

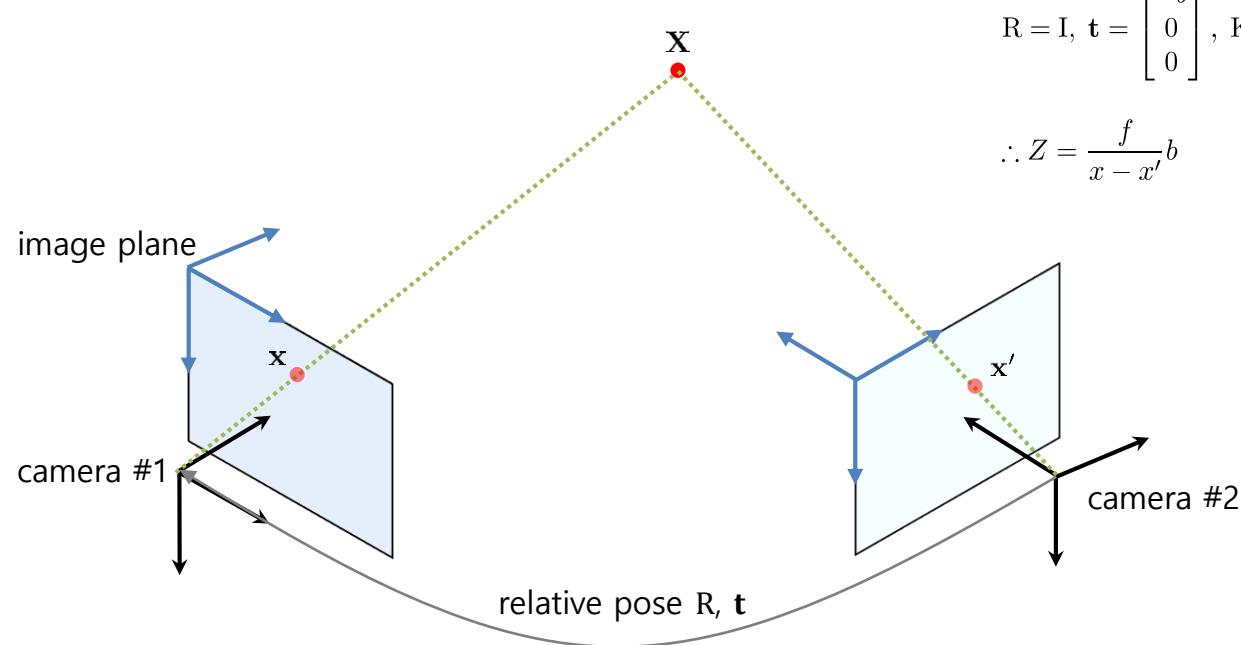


General 2D-2D Geometry (Epipolar Geometry)

- **Triangulation** (Point Localization)

- Unknown: Position of a 3D point \mathbf{X} (3 DoF)
- Given: Point correspondence $(\mathbf{x}, \mathbf{x}')$, camera matrices $(\mathbf{K}, \mathbf{K}')$, and relative pose (\mathbf{R}, \mathbf{t})
- Constraints: $2 \times$ projection $\mathbf{x} = \mathbf{K} [\mathbf{I} | \mathbf{0}] \mathbf{X}$, $\mathbf{x}' = \mathbf{K}' [\mathbf{R} | \mathbf{t}] \mathbf{X}$
- Solution (OpenCV): `cv::triangulatePoints()`

Special case) Stereo cameras



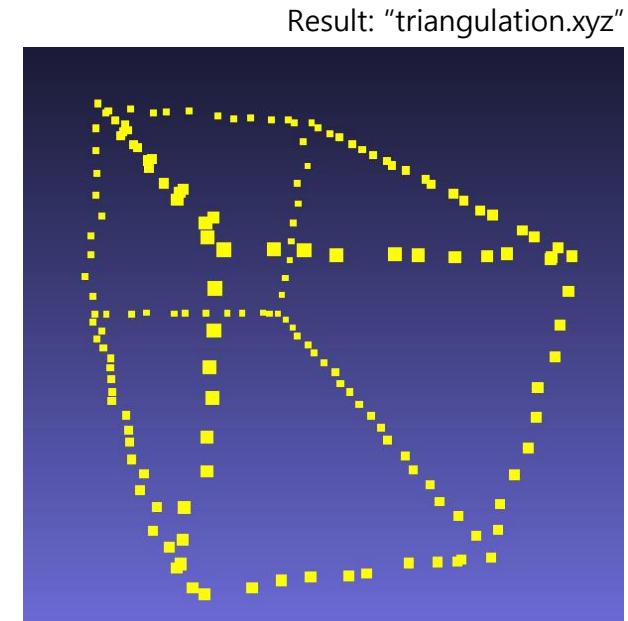
$$\mathbf{R} = \mathbf{I}, \mathbf{t} = \begin{bmatrix} -b \\ 0 \\ 0 \end{bmatrix}, \mathbf{K} = \mathbf{K}' = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

$$\therefore Z = \frac{f}{x - x'} b$$

- Example: **Triangulation (Two-view Reconstruction)** [triangulation.cpp]

```

1. #include "opencv2/opencv.hpp"
2. int main()
3. {
4.     // cf. You need to run 'image_formation.cpp' to generate point observation.
5.     const char *input0 = "image_formation0.xyz", *input1 = "image_formation1.xyz";
6.     double f = 1000, cx = 320, cy = 240;
7.
8.     // Load 2D points observed from two views
9.     std::vector<cv::Point2d> points0, points1;
10.    ...
11.
12.    // Estimate relative pose of two views
13.    cv::Mat F = cv::findFundamentalMat(points0, points1, cv::FM_8POINT);
14.    cv::Mat K = (cv::Mat<double>(3, 3) << f, 0, cx, 0, f, cy, 0, 0, 1);
15.    cv::Mat E = K.t() * F * K;
16.    cv::Mat R, t;
17.    cv::recoverPose(E, points0, points1, K, R, t);
18.
19.    // Reconstruct 3D points (triangulation)
20.    cv::Mat P0 = K * cv::Mat::eye(3, 4, CV_64F);           :: x = K[I|0]X
21.    cv::Mat Rt, X;
22.    cv::hconcat(R, t, Rt);                                :: x' = K'[R|t]X
23.    cv::Mat P1 = K * Rt;
24.
25.    cv::triangulatePoints(P0, P1, points0, points1, X);
26.    X.row(0) = X.row(0) / X.row(3);
27.    X.row(1) = X.row(1) / X.row(3);
28.    X.row(2) = X.row(2) / X.row(3);
29.    X.row(3) = 1;
30.
31.    // Store the 3D points
32.    FILE* fout = fopen("triangulation.xyz", "wt");
33.    if (fout == NULL) return -1;
34.    for (int c = 0; c < X.cols; c++)
35.        fprintf(fout, "%f %f %f\n", X.at<double>(0, c), X.at<double>(1, c), X.at<double>(2, c));
36.    fclose(fout);
37.    return 0;
38. }
```

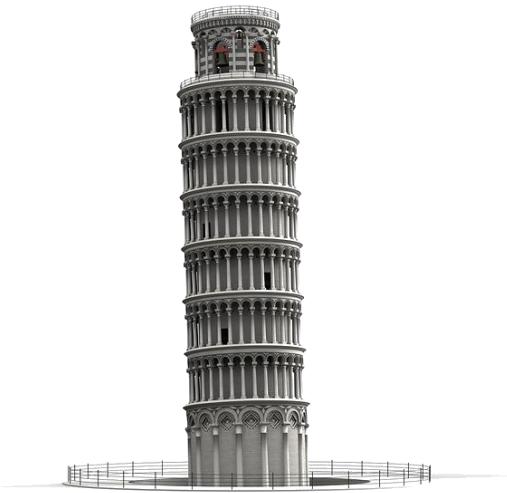


Correspondence Problem

Feature Correspondence

Robust Parameter Estimation

- RANSAC
- M-estimator



Images are from [pixabay](#).

Quiz



600 pixels

What is it?

Quiz

What is it?

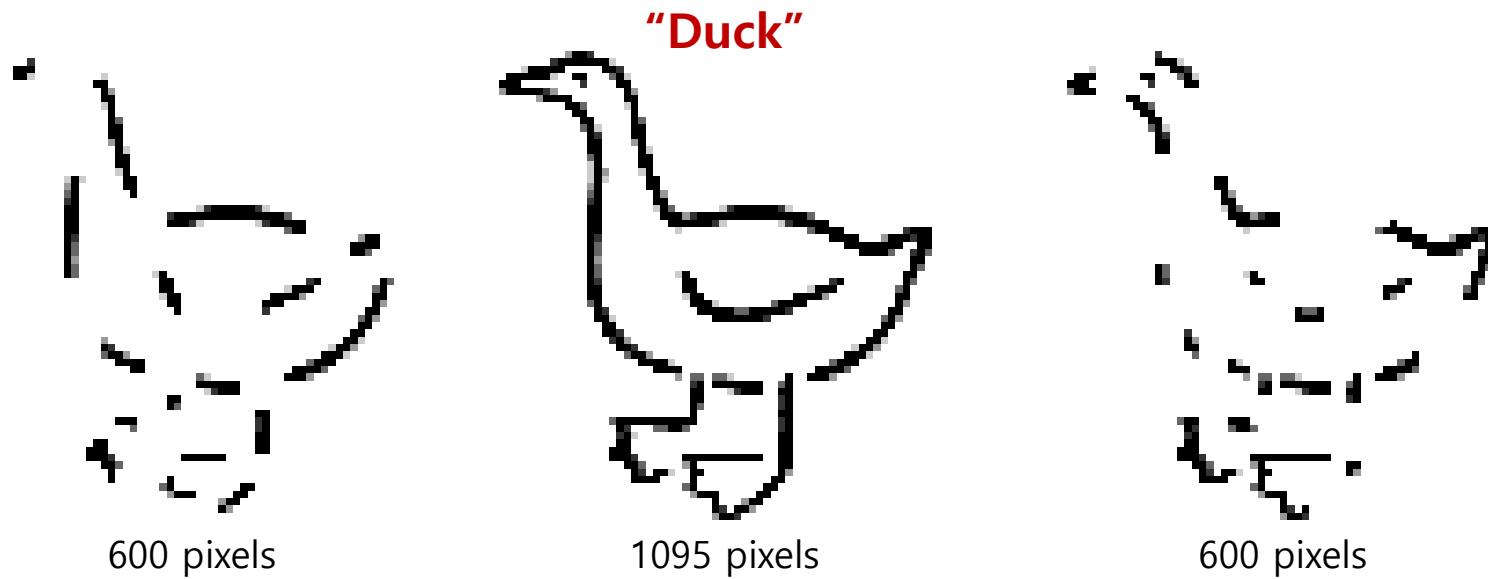


600 pixels

Quiz

- **Why corners (junction)?**

- a.k.a. keypoints, interest points, salient points, and feature points
- cf. \subset local invariant features (e.g. corner, edge, region, ...)

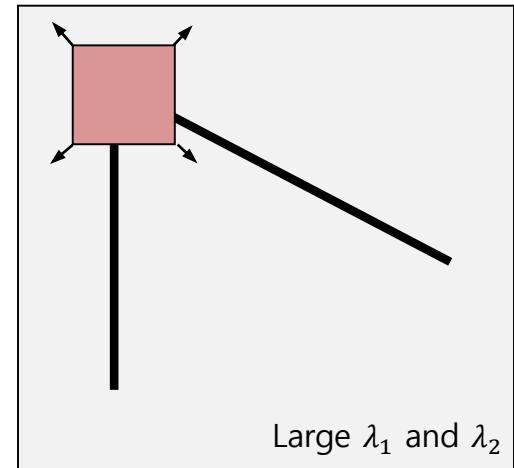
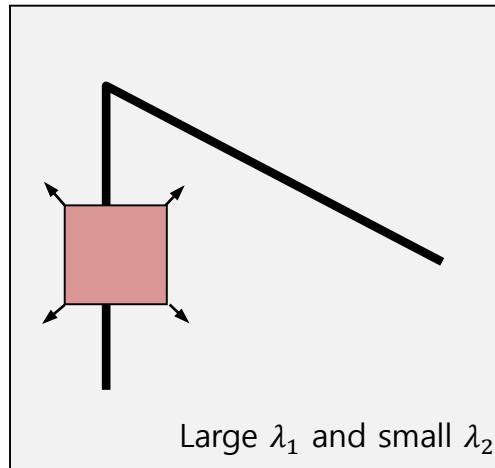
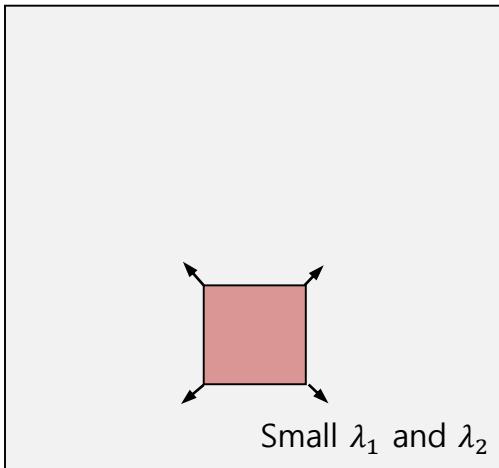


- **Requirements of local invariance features**

- **repeatability** (invariance, robustness), **distinctiveness**, **locality** (due to occlusion)
- quantity, accuracy, efficiency

Harris Corner (1988)

- Key idea: **Sliding window**



“flat” region:
no change
in all directions

“edge”:
no change
along the edge direction

“corner”:
significant change
in all directions

$$C(\Delta_x, \Delta_y) = \sum_{(x,y) \in W} \left(I(x + \Delta_x, y + \Delta_y) - I(x, y) \right)^2$$

$$\approx [\Delta_x \quad \Delta_y] \begin{bmatrix} \sum_W I_x^2 & \sum_W I_x I_y \\ \sum_W I_x I_y & \sum_W I_y^2 \end{bmatrix} \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix}$$

M

c.f. $I(x + \Delta_x, y + \Delta_y) \approx I(x, y) + [I_x(x, y) \quad I_y(x, y)] \begin{bmatrix} \Delta_x \\ \Delta_y \end{bmatrix}$ and $I_x = \frac{\partial I}{\partial x}$

$R = \min(\lambda_1, \lambda_2)$

Harris corner response:

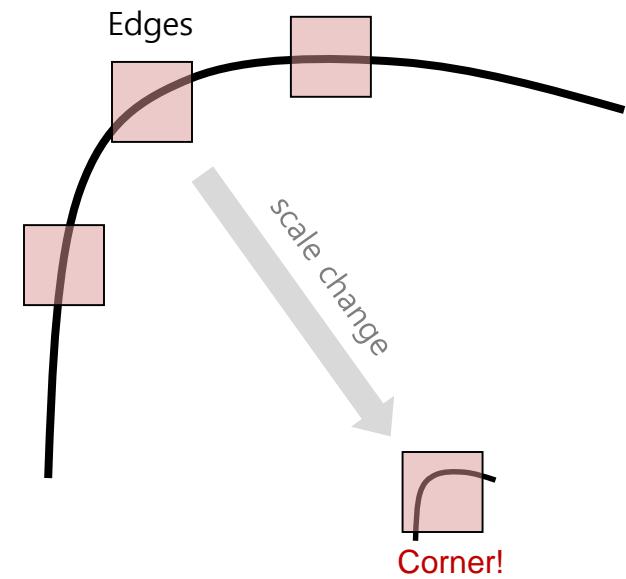
$$R = \det(M) - k \operatorname{trace}(M)^2$$

c.f. $\det(M) = \lambda_1 \lambda_2$, $\operatorname{trace}(M) = \lambda_1 + \lambda_2$, $k \in [0.04, 0.06]$

cf. Good-Feature-to-Track (Shi-Tomasi; 1994):

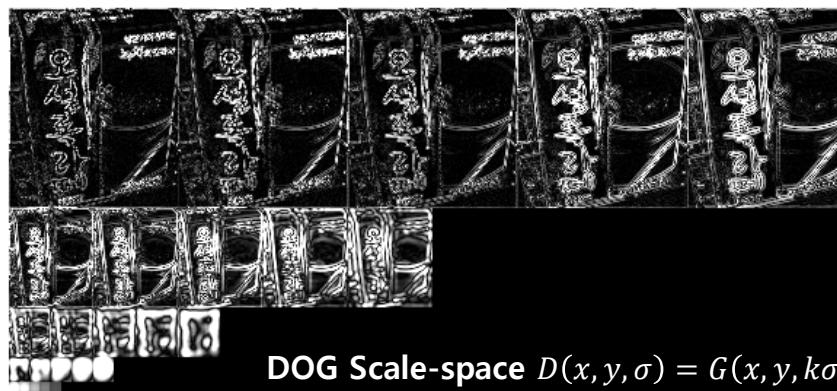
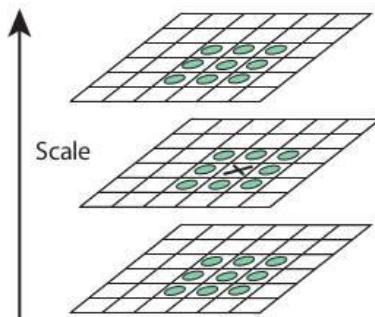
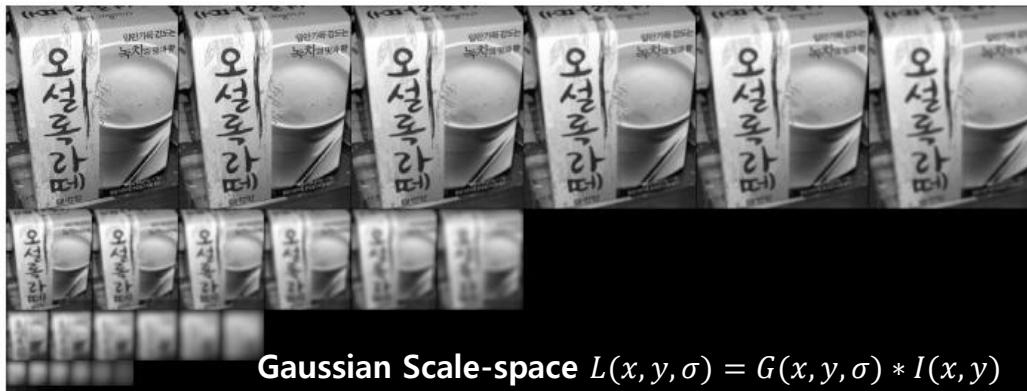
Harris Corner (1988)

- Properties
 - Invariant to translation, rotation, and intensity shift ($I \rightarrow I + b$) ~~intensity scaling ($I \rightarrow aI$)~~
 - But variant to **image scaling**



SIFT (Scale-Invariant Feature Transform; 1999)

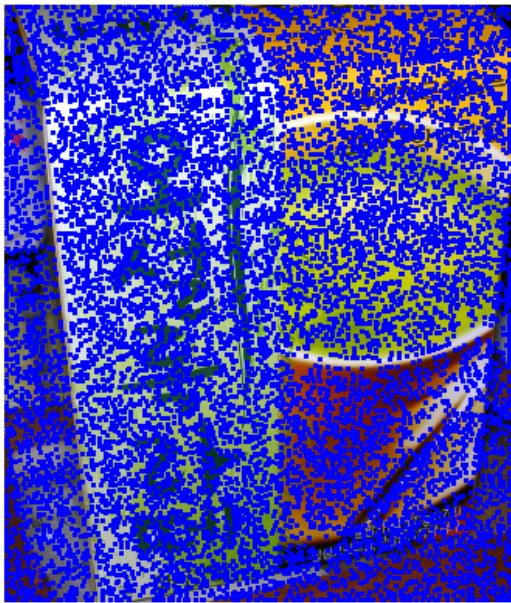
- Key idea: **Scale-space** (~ image pyramid)



- Part #1) **Feature point detection**

- Find **local extrema** (minima and maxima) in DOG scale-space
- Localize its position accurately (sub-pixel level) using 3D quadratic function
- Eliminate **low contrast candidates**, $|D(\mathbf{x})| < \tau$
- Eliminate **candidates on edges**, $\frac{\text{trace}(H)^2}{\det(H)} < \frac{(r+1)^2}{r}$ where $H = \begin{bmatrix} D_{xx} & D_{xy} \\ D_{xy} & D_{yy} \end{bmatrix}$

SIFT (Scale-Invariant Feature Transform; 1999)



local extrema (N: 11479)



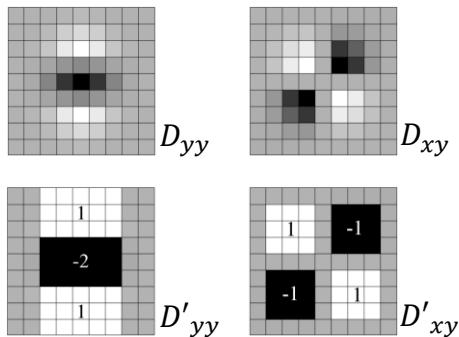
feature points (N: 971)

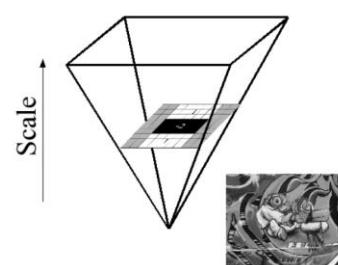
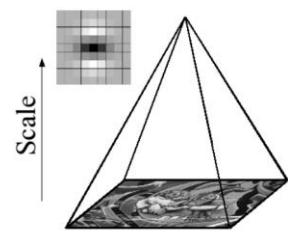


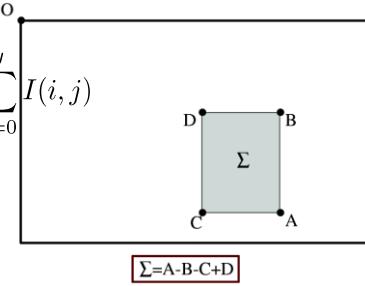
feature scales and orientations

cf. SURF (Speeded Up Robust Features; 2006):

- Key idea: Approximation of SIFT using **integral image** and ...


$$D_{yy} = \begin{bmatrix} & & \\ & -1 & \\ & & \end{bmatrix}$$
$$D_{xy} = \begin{bmatrix} & & \\ & -1 & \\ & & \end{bmatrix}$$
$$D'_{yy} = \begin{bmatrix} 1 & & \\ & -2 & \\ & & 1 \end{bmatrix}$$
$$D'_{xy} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

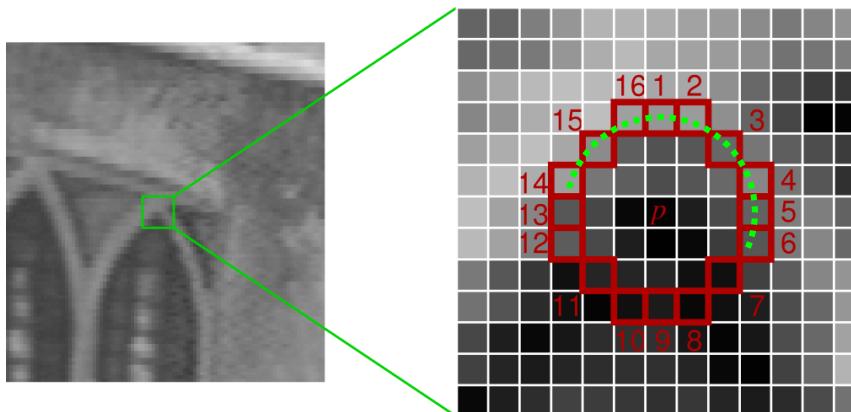


$$S(x, y) = \sum_{i=0}^x \sum_{j=0}^y I(i, j)$$

$$\Sigma = A - B - C + D$$

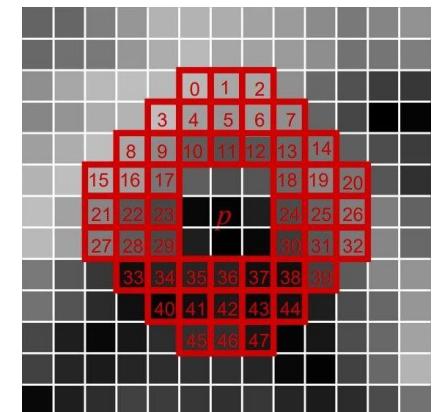
FAST (Features from Accelerated Segment Test; 2006)

- Key idea: **Continuous arc of N or more pixels**

- Is this patch a corner?
 - Is the segment brighter than $p + t$? Is the segment darker than $p - t$?
 - t : The threshold of similar intensity
 - Too many corners! it needs non-maximum suppression.

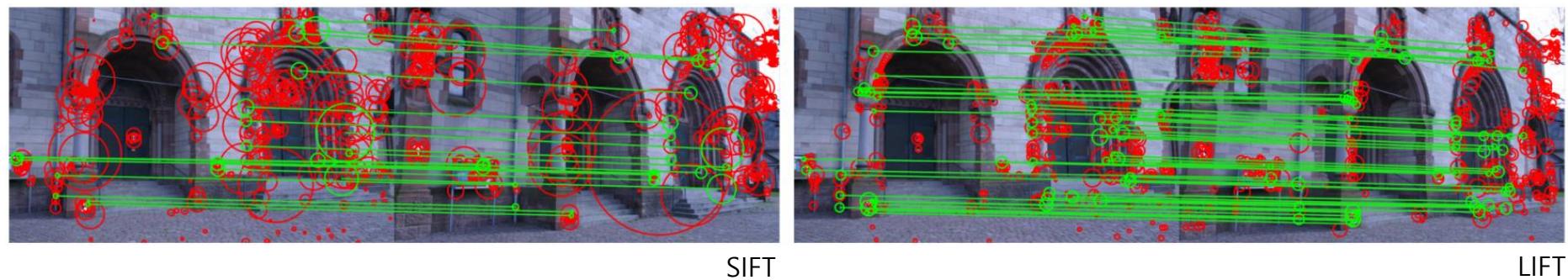
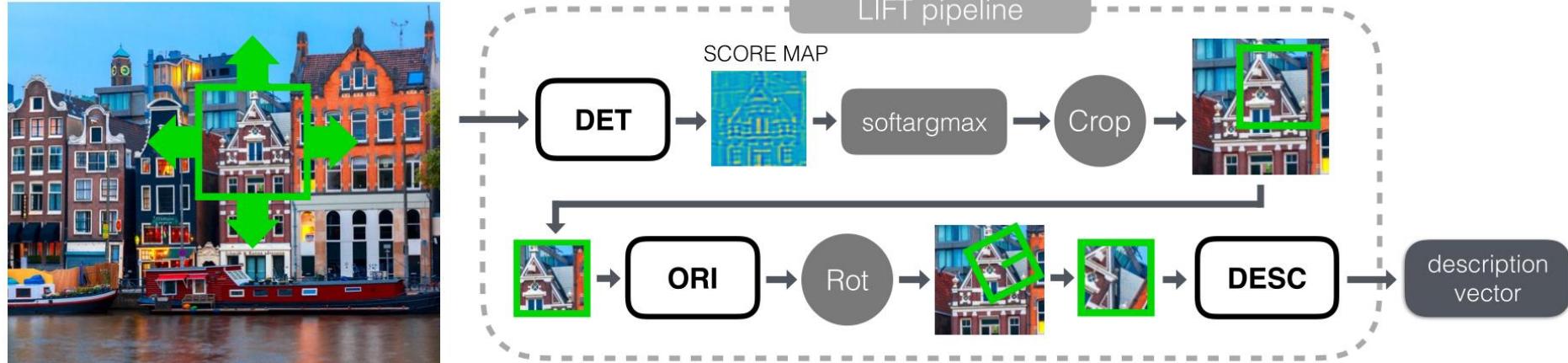


- Versions
 - **FAST-9** ($N: 9$), FAST-12 ($N: 12$), ...
 - FAST-ER: Training a decision tree to enhance repeatability with more pixels



LIFT (Learned Invariant Feature Transform; 2016)

- Key idea: **Deep neural network**
 - **DET** (feature detector) + **ORI** (orientation estimator) + **DESC** (feature descriptor)



Overview of Feature Correspondence

- **Features**
 - **Corners:** Harris corner, GFTT (Shi-Tomasi corner), SIFT, SURF, FAST, LIFT, ...
 - Edges, line segments, regions, ...
- **Feature Descriptors and Matching**
 - **Patch:** Raw intensity
 - Measures: SSD (sum of squared difference), ZNCC (zero normalized cross correlation), ...
 - **Floating-point descriptors:** SIFT, SURF, (DAISY), LIFT, ... → e.g. A 128-dim. vector (a histogram of gradients)
 - Measures: Euclidean distance, cosine distance, (the ratio of first and second bests)
 - Matching: Brute-force matching ($O(N^2)$), ANN (approximated nearest neighborhood) search ($O(\log N)$)
 - Pros (+): **High discrimination power**
 - Cons (-): **Heavy computation**
 - **Binary descriptors:** BRIEF, ORB, (BRISK), (FREAK), ... → e.g. A 128-bit string (a series of intensity comparison)
 - Measures: Hamming distance
 - Matching: Brute-force matching ($O(N^2)$)
 - Pros (+): **Less storage and faster extraction/matching**
 - Cons (-): **Less performance**
- **Feature Tracking (a.k.a. Optical Flow)**
 - **Optical flow:** (Horn-Schunck method), Lukas-Kanade method
 - Measures: SSD (sum of squared difference)
 - Tracking: Finding displacement of a similar patch
 - Pros (+): **No descriptor and matching** (faster and compact)
 - Cons (-): **Not working in wide baseline**

SIFT (Scale-Invariant Feature Transform; 1999)

- Part #2) Orientation assignment

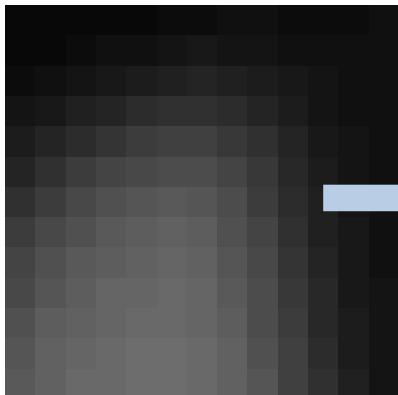
- Derive magnitude and orientation of gradient of each patch

$$m(x, y) = \sqrt{(L(x + 1, y) - L(x - 1, y))^2 + (L(x, y + 1) - L(x, y - 1))^2}$$

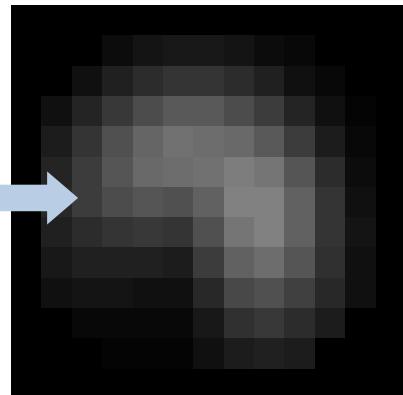
$$\theta(x, y) = \tan^{-1} \frac{L(x, y + 1) - L(x, y - 1)}{L(x + 1, y) - L(x - 1, y)}$$

- Find the strongest orientation

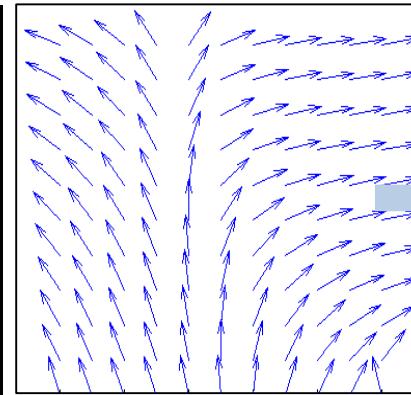
- Histogram voting (36 bins) with Gaussian-weighted magnitude



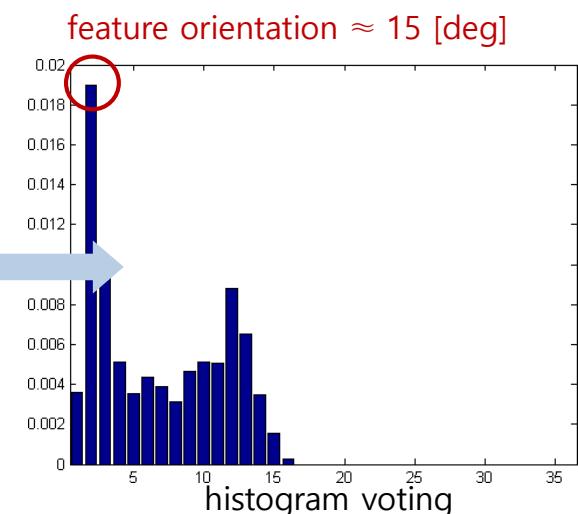
feature patch



gradient magnitude



gradient orientation

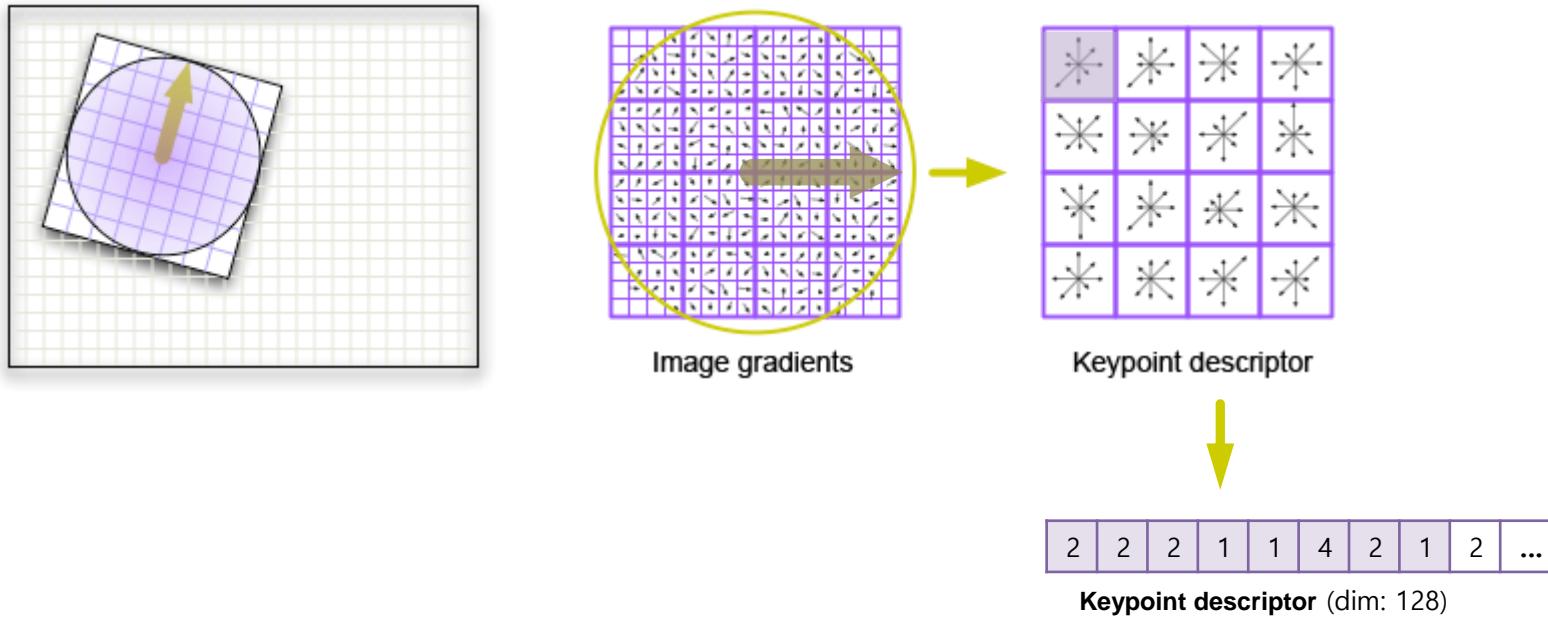


SIFT (Scale-Invariant Feature Transform; 1999)

- Part #3) **Feature descriptor extraction**
 - Build a 4x4 gradient histogram (8 bins) from each patch (16x16 pixels)
 - Use Gaussian-weighted magnitude again
 - Use relative angles w.r.t. the assigned feature orientation
 - Encode the histogram into a 128-dimensional vector

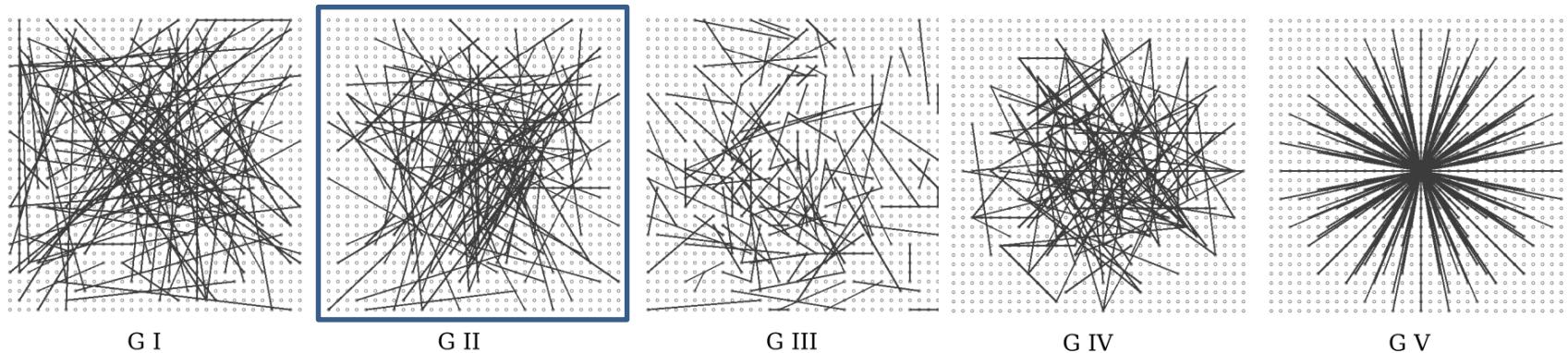


feature scales and orientations



BRIEF (Binary Robust Independent Elementary Features; 2010)

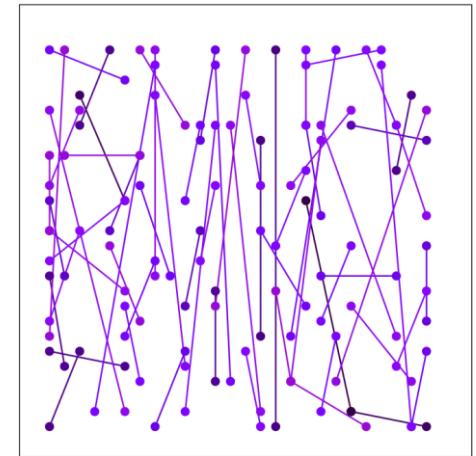
- Key idea: **A sequence of intensity comparison of random pairs**
 - Applying smoothing for stability and repeatability
 - Path size: 31×31 pixels



- Versions: The number of tests
 - BRIEF-32, BRIEF-64, BRIEF-128, BRIEF-256 ...
- Examples of combinations
 - CenSurE detector (a.k.a. Star detector) + BRIEF descriptor
 - SURF detector + BRIEF descriptor

ORB (Oriented FAST and rotated BRIEF, 2011)

- Key idea: **Adding rotation invariance to BRIEF**
 - **Oriented FAST**
 - Generate scale pyramid for scale invariance
 - Detect *FAST-9* points (filtering with Harris corner response)
 - Calculate feature orientation by *intensity centroid*
$$\theta = \tan^{-1} \frac{m_{01}}{m_{10}} \quad \text{where} \quad m_{pq} = \sum_{x,y} x^p y^q I(x, y)$$
 - **Rotation-aware BRIEF**
 - Extract BRIEF descriptors w.r.t. the known orientation
 - Use better comparison pairs trained by greedy search
- Combination: **ORB**
 - FAST-9 detector (with orientation) + BRIEF-256 descriptor (with trained pairs)
- Computing time
 - ORB: **15.3 [msec]** / SURF: 217.3 [msec] / SIFT: 5228.7 [msec] @ 24 images (640x480) in Pascal dataset



Lukas-Kanade Optical Flow

- Key idea: **Finding movement of a patch**

Brightness constancy constraint: $I(x, y, t) = I(x + \Delta_x, y + \Delta_y, t + \Delta_t)$

(if same patch)

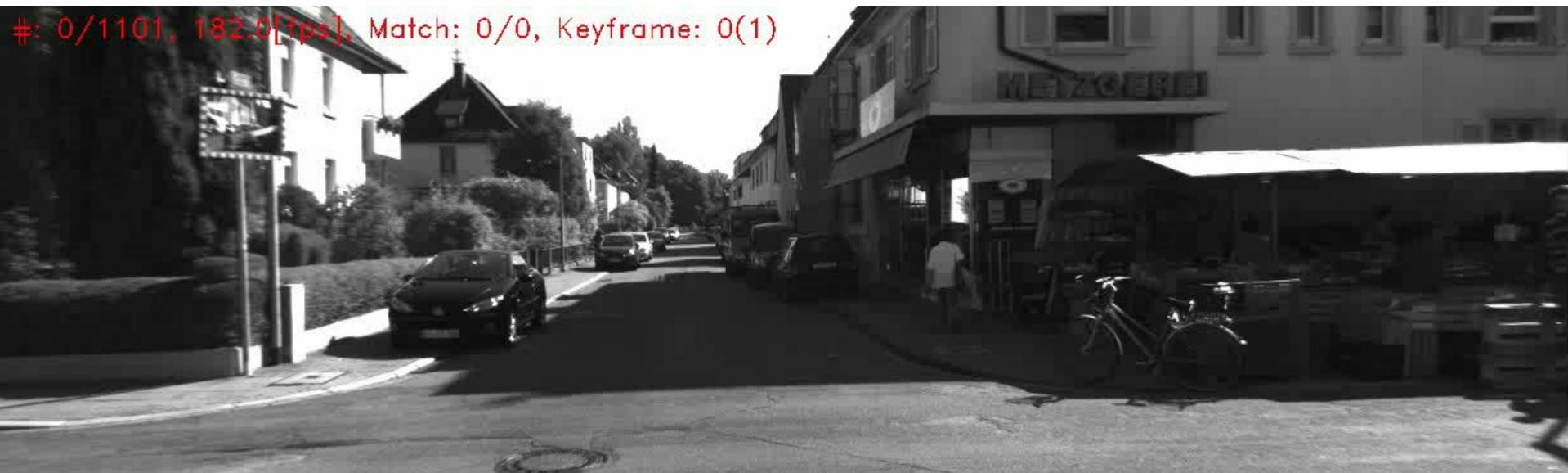
$$I_x \frac{\Delta_x}{\Delta_t} + I_y \frac{\Delta_y}{\Delta_t} + I_t = 0 \quad \text{because} \quad I(x + \Delta_x, y + \Delta_y, t + \Delta_t) \approx I(x, y, t) + I_x \Delta_x + I_y \Delta_y + I_t \Delta_t$$

$$\mathbf{A}\mathbf{v} = \mathbf{b} \quad \text{where} \quad \mathbf{A} = \begin{bmatrix} I_x(p_1) & I_y(p_1) \\ \vdots & \vdots \\ I_x(p_n) & I_y(p_n) \end{bmatrix}, \mathbf{v} = \begin{bmatrix} V_x \\ V_y \end{bmatrix}, \mathbf{b} = \begin{bmatrix} -I_t(p_1) \\ \vdots \\ -I_t(p_n) \end{bmatrix}, \text{and } p_i \in W$$

$$\therefore \mathbf{v} = \mathbf{A}^\dagger \mathbf{b} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}$$

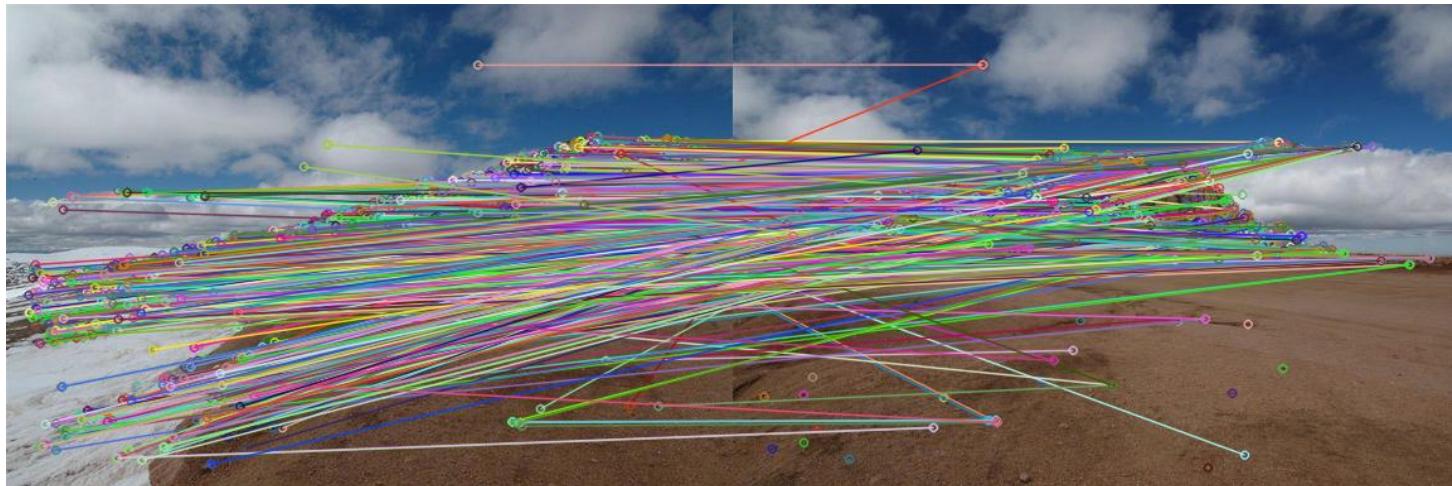
- Combination: **KLT tracker**

- Shi-Tomasi detector (a.k.a. GFTT) + Lukas-Kanade optical flow

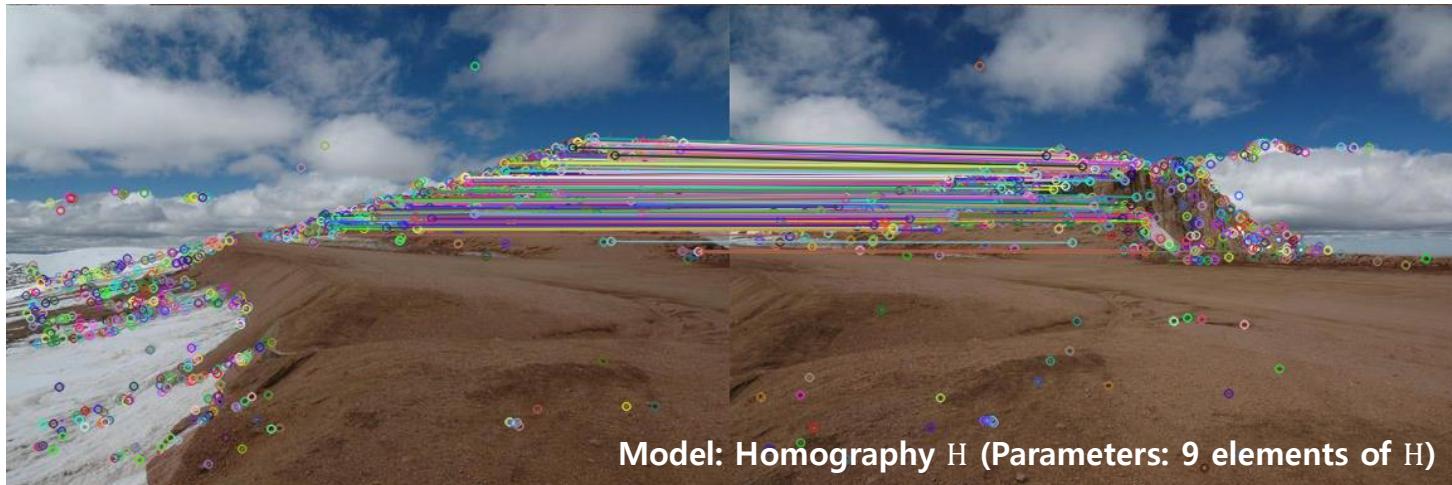


Why Outliers?

Putative matches (inliers + outliers)



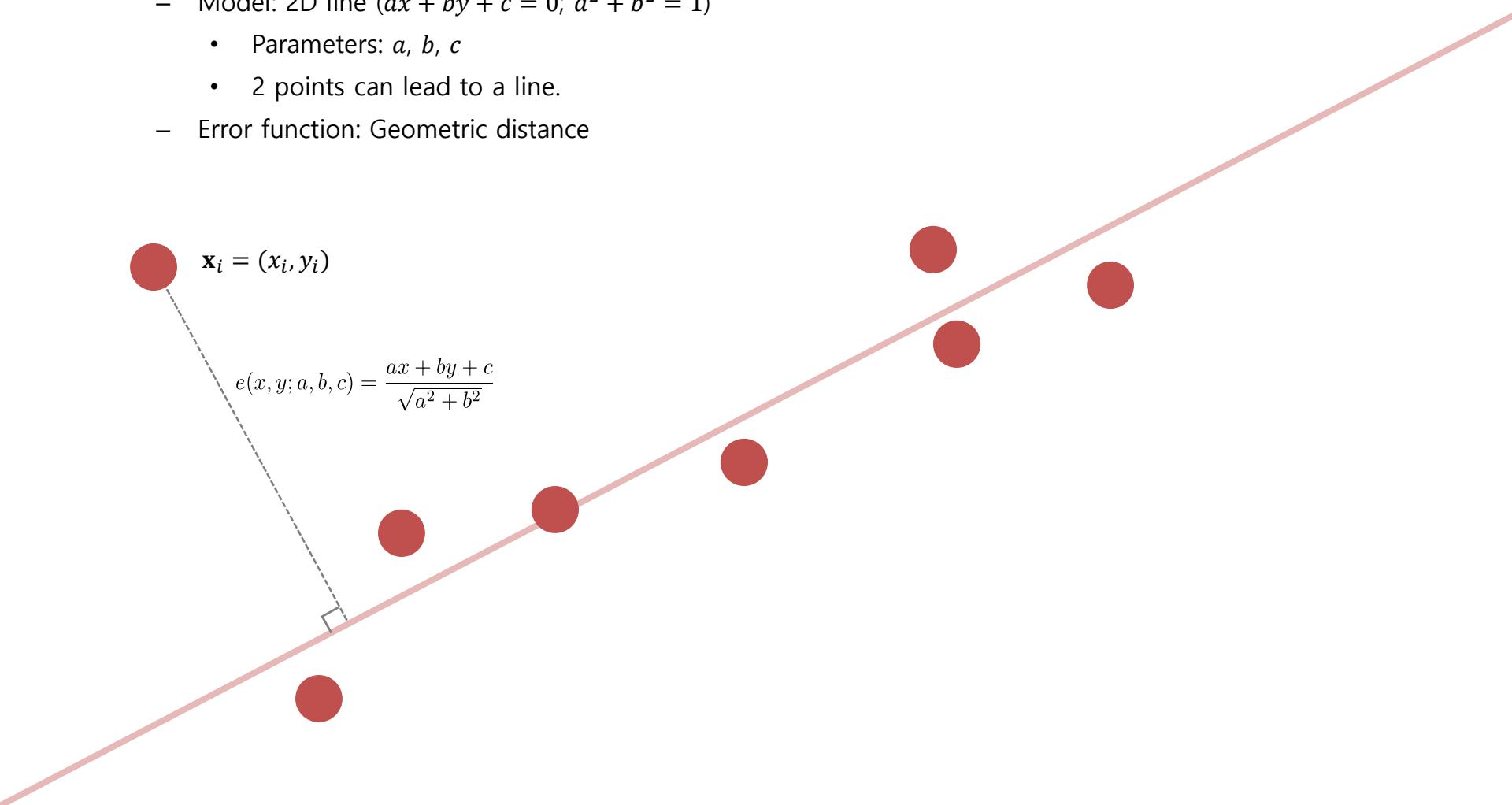
After applying RANSAC (inliers)



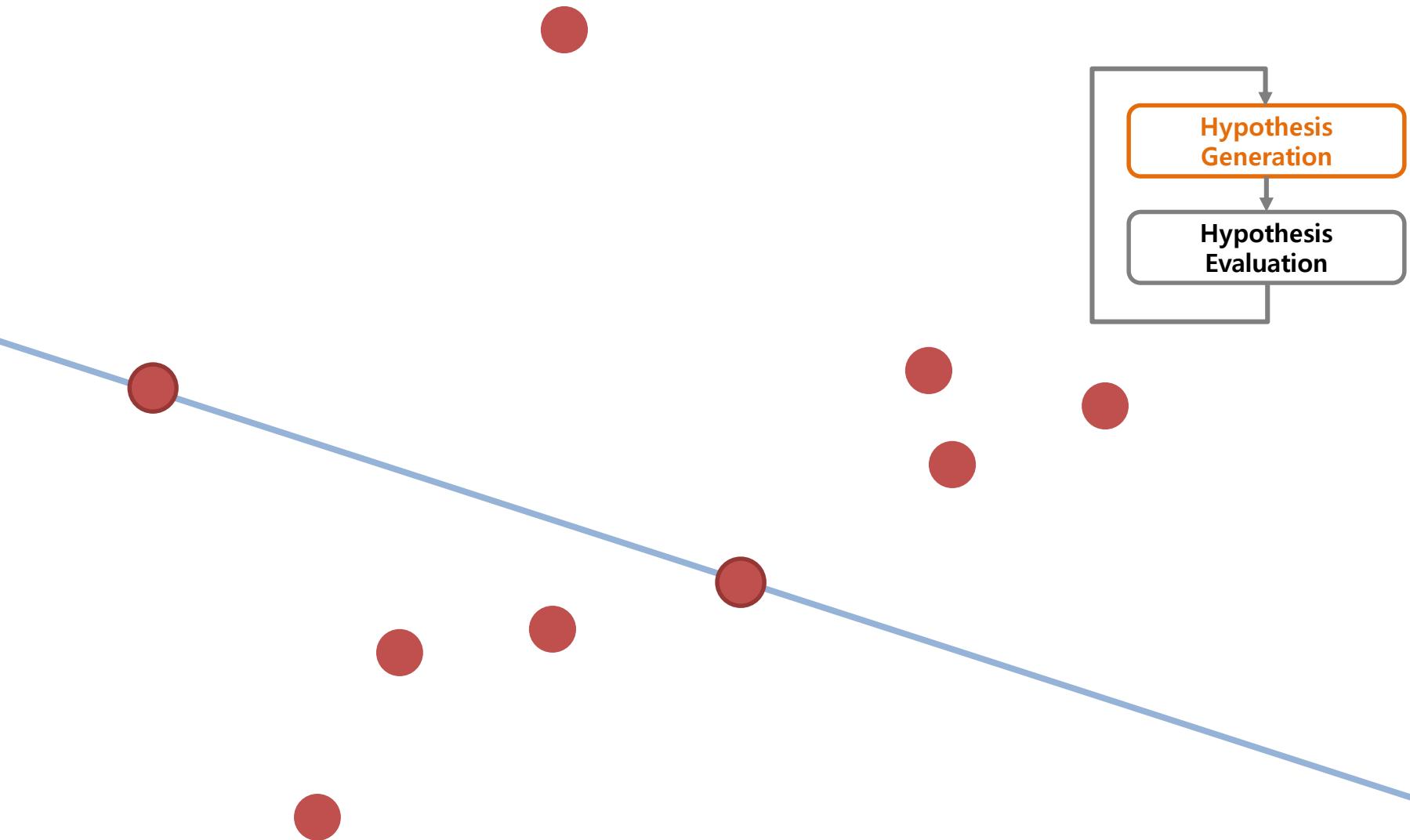
Model: Homography H (Parameters: 9 elements of H)

RANSAC: Random Sample Consensus

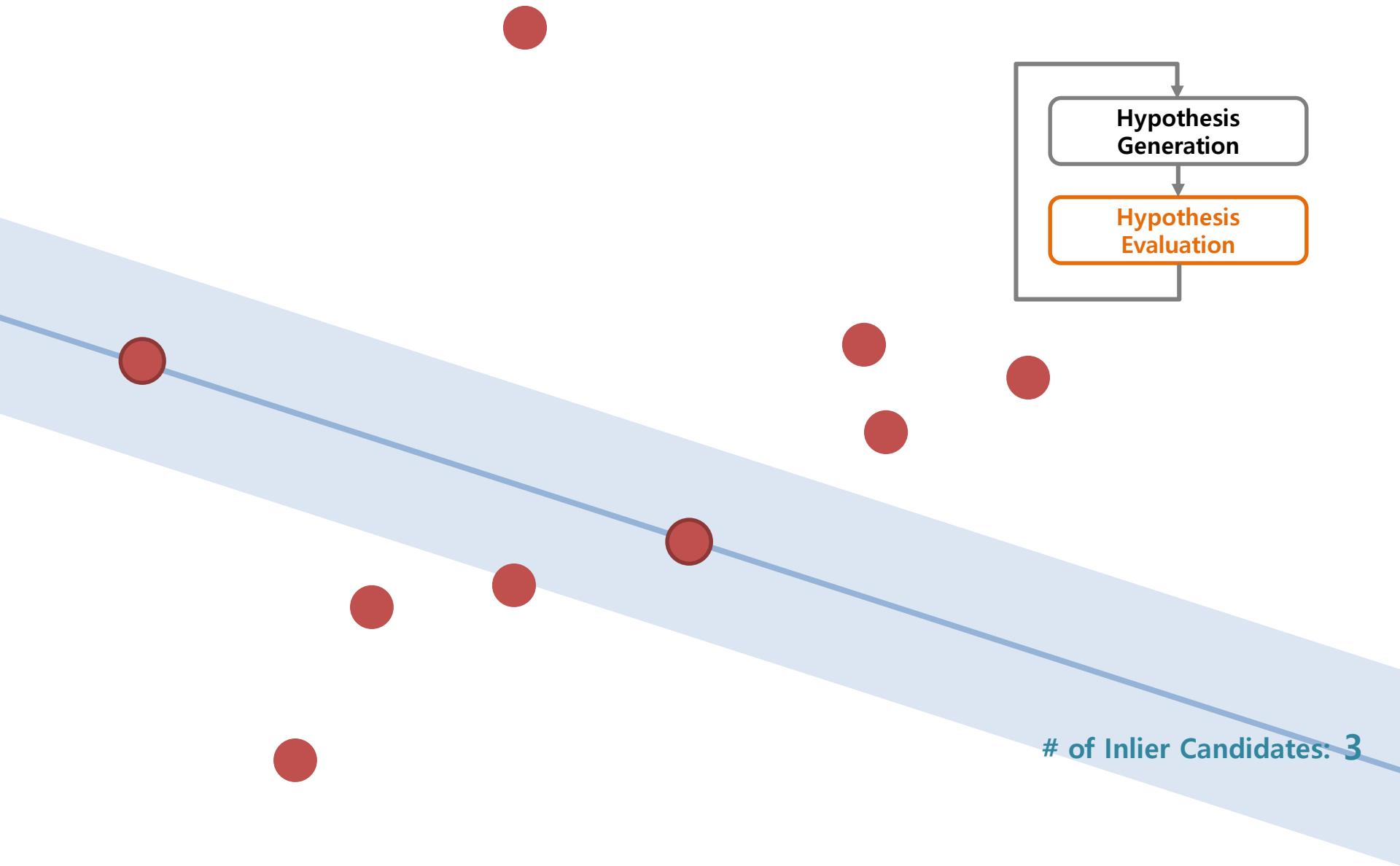
- Example: **Line Fitting with RANSAC**
 - Model: 2D line ($ax + by + c = 0; a^2 + b^2 = 1$)
 - Parameters: a, b, c
 - 2 points can lead to a line.
 - Error function: Geometric distance



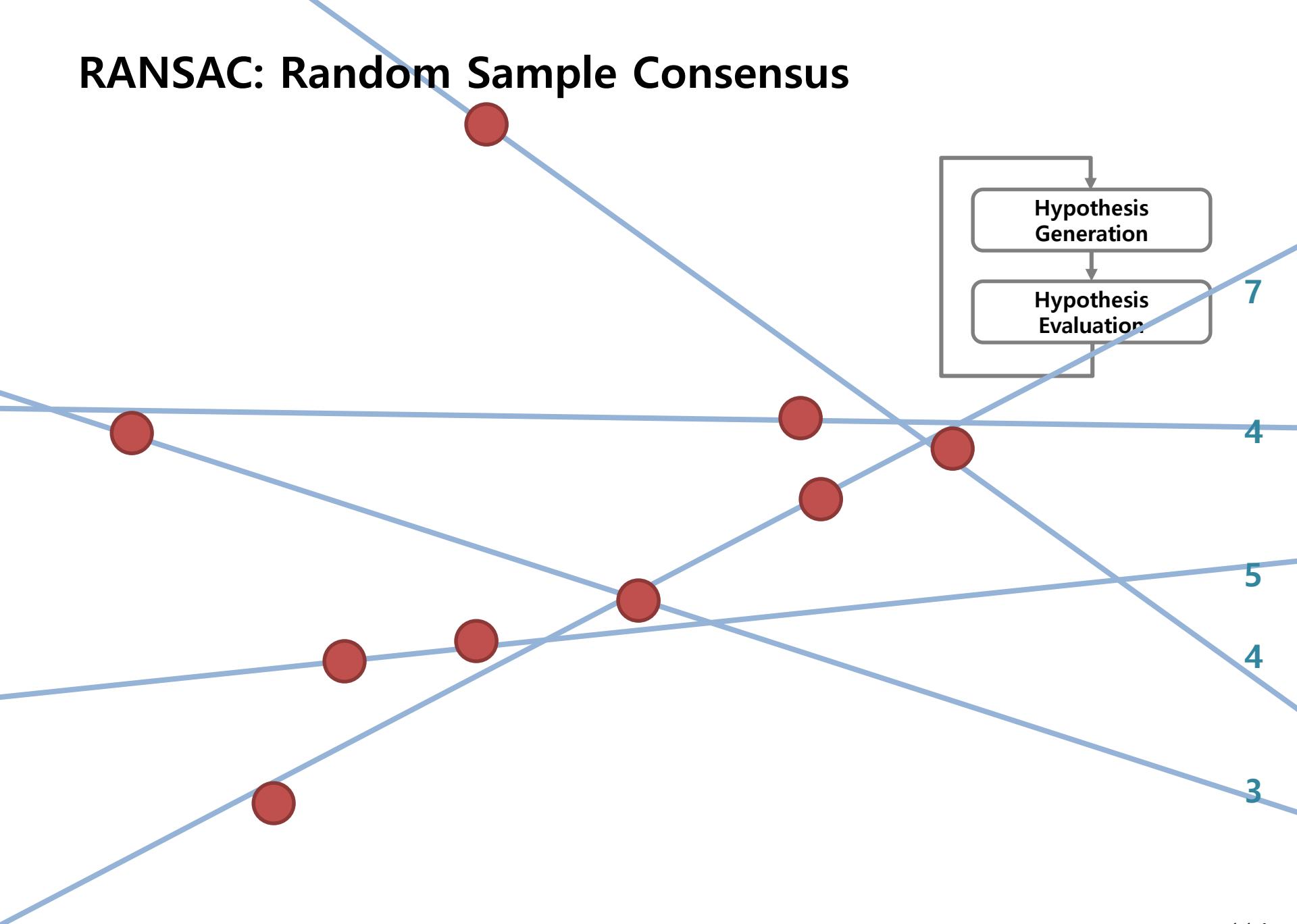
RANSAC: Random Sample Consensus



RANSAC: Random Sample Consensus



RANSAC: Random Sample Consensus



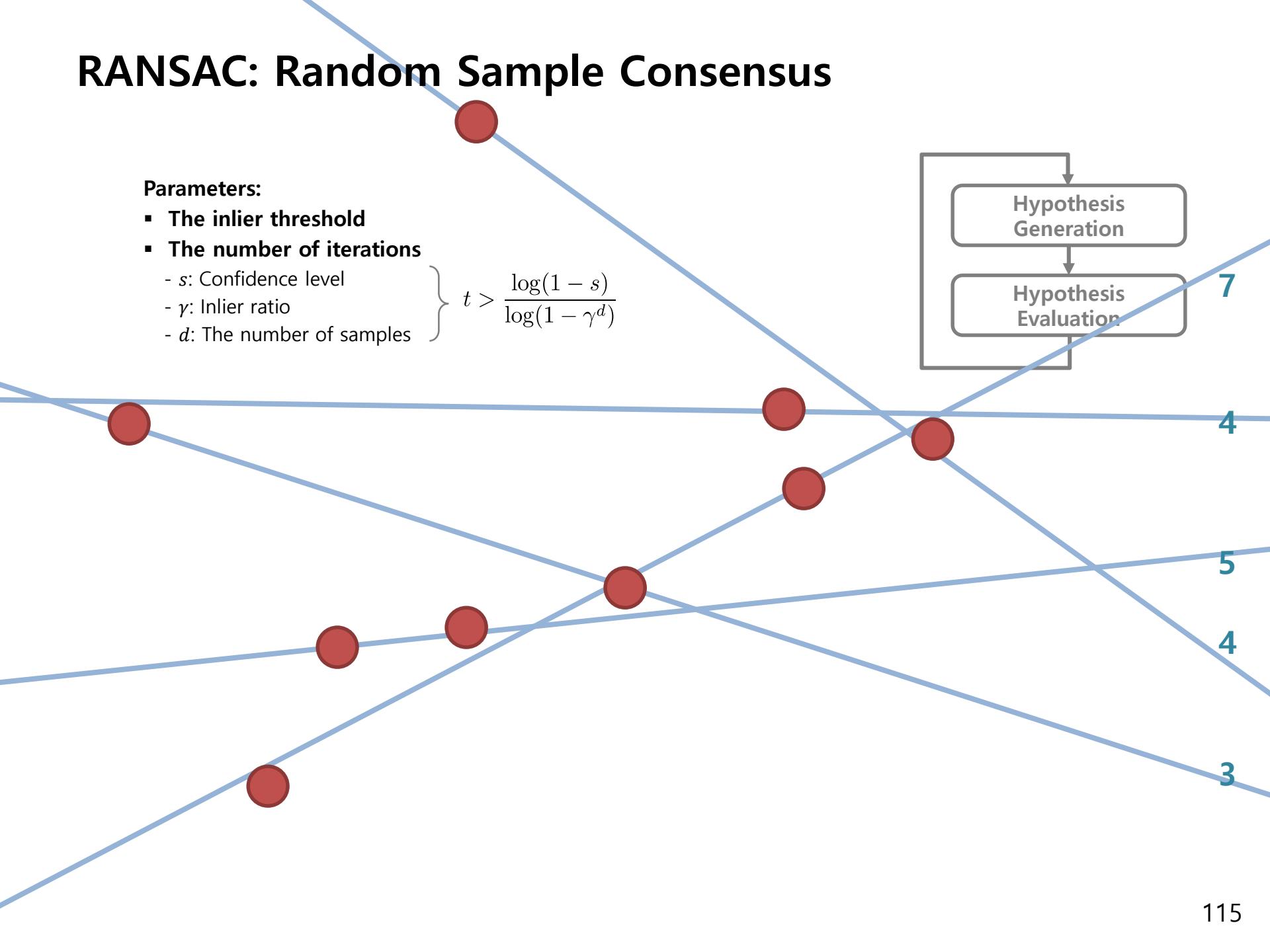
RANSAC: Random Sample Consensus

Parameters:

- The inlier threshold
- The number of iterations

- s : Confidence level
- γ : Inlier ratio
- d : The number of samples

$$t > \frac{\log(1 - s)}{\log(1 - \gamma^d)}$$



- Example: **Line Fitting with RANSAC** [line_fitting_ransac.cpp]

```

1. #include "opencv2/opencv.hpp"

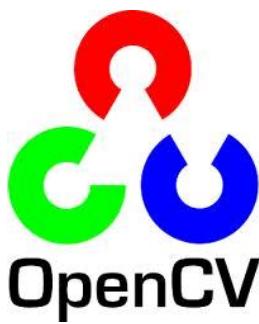
2. // Convert a line format, [n_x, n_y, x_0, y_0] to [a, b, c]
3. // cf. A line model in OpenCV: n_x * (x - x_0) = n_y * (y - y_0)
4. #define CONVERT_LINE(line) (cv::Vec3d(line[0], -line[1], -line[0] * line[2] + line[1] * line[3]))

5. int main()
6. {
7.     cv::Vec3d truth(1.0 / sqrt(2.0), 1.0 / sqrt(2.0), -240.0); // The line model: a*x + b*y + c = 0 (a^2 + b^2 = 1)
8.     int ransac_trial = 50, ransac_n_sample = 2;
9.     double ransac_thresh = 3.0; // 3 x 'data_inlier_noise'
10.    int data_num = 1000;
11.    double data_inlier_ratio = 0.5, data_inlier_noise = 1.0;

12.    // Generate data
13.    std::vector<cv::Point2d> data;
14.    cv::RNG rng;
15.    for (int i = 0; i < data_num; i++)
16.    {
17.        if (rng.uniform(0.0, 1.0) < data_inlier_ratio)
18.        {
19.            double x = rng.uniform(0.0, 480.0);
20.            double y = (truth(0) * x + truth(2)) / -truth(1);
21.            x += rng.gaussian(data_inlier_noise);
22.            y += rng.gaussian(data_inlier_noise);
23.            data.push_back(cv::Point2d(x, y)); // Inlier
24.        }
25.        else data.push_back(cv::Point2d(rng.uniform(0.0, 640.0), rng.uniform(0.0, 480.0))); // Outlier
26.    }

27.    // Estimate a line using RANSAC ...
28.    // Estimate a line using least-squares method (for reference) ...

29.    // Display estimates
30.    printf("The Truth: %.3f, %.3f, %.3f\n", truth[0], truth[1], truth[2]);
31.    printf("Estimate (RANSAC): %.3f, %.3f, %.3f (Score: %d)\n", best_line[0], best_line[1], ..., best_score);
32.    printf("Estimate (LSM): %.3f, %.3f, %.3f\n", lsm_line[0], lsm_line[1], lsm_line[2]);
33.    return 0;
34. }
```



$$\text{c.f. } t > \frac{\log(1 - s)}{\log(1 - \gamma^d)} = \frac{\log(1 - 0.999)}{\log(1 - 0.5^2)} = 24$$

```

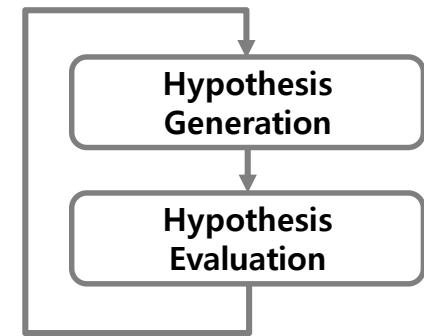
27. // Estimate a line using RANSAC
28. int best_score = -1;
29. cv::Vec3d best_line;
30. for (int i = 0; i < ransac_trial; i++)
31. {
32.     // Step 1: Hypothesis generation
33.     std::vector<cv::Point2d> sample;
34.     for (int j = 1; j < ransac_n_sample; j++)
35.     {
36.         int index = rng.uniform(0, int(data.size()));
37.         sample.push_back(data[index]);
38.     }
39.     cv::Vec4d nnxy;
40.     cv::fitLine(sample, nnxy, CV_DIST_L2, 0, 0.01, 0.01);
41.     cv::Vec3d line = CONVERT_LINE(nnxy);

42.     // Step 2: Hypothesis evaluation
43.     int score = 0;
44.     for (size_t j = 0; j < data.size(); j++)
45.     {
46.         double error = fabs(line(0) * data[j].x + line(1) * data[j].y + line(2));
47.         if (error < ransac_thresh) score++;
48.     }

49.     if (score > best_score)
50.     {
51.         best_score = score;
52.         best_line = line;
53.     }
54. }

55. // Estimate a line using least squares method (for reference)
56. cv::Vec4d nnxy;
57. cv::fitLine(data, nnxy, CV_DIST_L2, 0, 0.01, 0.01);
58. cv::Vec3d lsm_line = CONVERT_LINE(nnxy);

```



Line Fitting Result

```

* The Truth: 0.707, 0.707, -240.000
* Estimate (RANSAC): 0.712, 0.702, -242.170 (Score: 434)
* Estimate (LSM): 0.748, 0.664, -314.997

```

Least Squares Method, RANSAC, and M-estimator

- **Least Squares Method**

- Find a model while minimizing sum of squared errors,
- $$\arg \min_{a,b,c} \sum_i e(\mathbf{x}_i; a, b, c)^2$$

- **RANSAC**

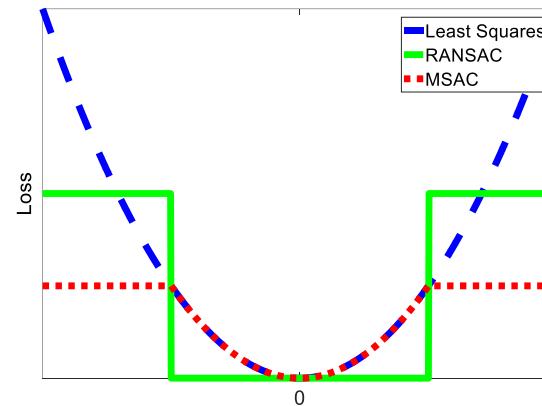
- Find a model while maximizing the number of supports (~ inlier candidates)
 - ~ minimizing the number of outlier candidates

- **Why RANSAC was robust to outliers?**

Problem: $\arg \min_{a,b,c} \sum_i \rho(e(\mathbf{x}_i; a, b, c))$

In the view of **loss functions** ρ ,

- Least squares method: $\rho(x) = x^2$
 - RANSAC:
- $$\rho(x) = \begin{cases} 0 & \text{if } |x| < \tau \\ 1 & \text{o.w.} \end{cases}$$



- **M-estimator (~ weighted least squares) | MSAC (in OpenCV)**

- Find a model while minimizing sum of (squared) errors **with a truncated loss function**



One-page Tutorial for Ceres Solver

▪ Ceres Solver?

- An **open source C++** library for modelling and solving large and complicated **optimization** problems.
 - Since 2010 by Google (BSD license)
- Problem types: 1) **Non-linear least squares** (with bounds), 2) General unconstrained minimization
- Homepage: <http://ceres-solver.org/>

▪ Solving Non-linear Least Squares

$$\arg \min_{\mathbf{m}} \sum_i \rho_i(\|r_i(\mathbf{m})\|^2)$$

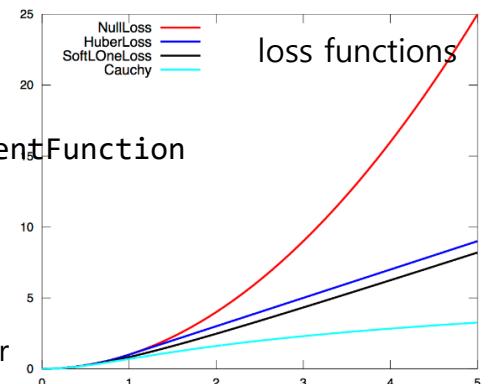
1. Define residual functions (or cost function or error function)
2. Instantiate `ceres::Problem` and add residuals using its member function, `AddResidualBlock()`
 - Instantiate each residual r_i in the form of `ceres::CostFunction` and add it
 - Select how to calculate its derivative (Jacobian)

`(ceres::AutoDiffCostFunction` or `ceres::NumericDiffCostFunction` or `ceres::SizedCostFunction`)
 - cf. Automatic derivation (using the chain rule) is recommended for convenience and performance.
 - Instantiate its `ceres::LossFunction` ρ_i and add it (if the problem needs robustness against outliers)
3. Instantiate `ceres::Solver::Options` (and also `ceres::Solver::Summary`) and configure the option
4. Run `ceres::Solve()`

▪ Solving General Minimization

$$\arg \min_x f(x)$$

- `ceres::CostFunction` → `ceres::FirstOrderFunction`, `ceres::GradientFunction`
- `ceres::Problem` → `ceres::GradientProblem`
- `ceres::Solver` → `ceres::GradientProblemSolver`



■ Example: Line Fitting with M-estimator [line_fitting_m_est.cpp]

```

1. #include "opencv2/opencv.hpp"
2. #include "ceres/ceres.h"
3. ...
4. struct GeometricError
5. {
6.     GeometricError(const cv::Point2d& pt) : datum(pt) { }
7.     template<typename T>
8.     bool operator()(const T* const line, T* residual) const
9.     {
10.         residual[0] = (line[0] * T(datum.x) + line[1] * T(datum.y) + line[2]) / sqrt(line[0] * line[0] + line[1] * line[1]);
11.         return true;
12.     }
13. private:
14.     const cv::Point2d datum;
15. };

16. int main()
17. {
18.     ...
19.     // Estimate a line using M-estimator
20.     cv::Vec3d opt_line(1, 0, 0);
21.     ceres::Problem problem;
22.     for (size_t i = 0; i < data.size(); i++)           2) Instantiate a problem and add a residual for each datum
23.     {
24.         ceres::CostFunction* cost_func = new ceres::AutoDiffCostFunction<GeometricError, 1, 3>(new GeometricError(data[i]));
25.         ceres::LossFunction* loss_func = NULL;
26.         if (loss_width > 0) loss_func = new ceres::CauchyLoss(loss_width);           The dimension of a residual
27.         problem.AddResidualBlock(cost_func, loss_func, opt_line.val);                The dimension of the first model parameter
28.     }
29.     ceres::Solver::Options options;           3) Instantiate options and configure it
30.     options.linear_solver_type = ceres::ITERATIVE_SCHUR;
31.     options.num_threads = 8;
32.     options.minimizer_progress_to_stdout = true;
33.     ceres::Solver::Summary summary;
34.     ceres::Solve(options, &problem, &summary);           4) Solve the minimization problem
35.     std::cout << summary.FullReport() << std::endl;
36.     opt_line /= sqrt(opt_line[0] * opt_line[0] + opt_line[1] * opt_line[1]); // Normalize
37.     ...
38.     return 0;
39. }

```

1) Define a residual as C++ generic function ($T \sim \text{double}$)
 cf. The generic is necessary for automatic differentiation.

c.f. $e(x, y; a, b, c) = \frac{ax + by + c}{\sqrt{a^2 + b^2}}$

2) Instantiate a problem and add a residual for each datum

3) Instantiate options and configure it

4) Solve the minimization problem

Overview of Robust Parameter Estimation

- Bottom-up Approaches (~ Voting) e.g. line fitting, relative pose estimation
 - **Hough transform**
 - A datum votes multiple parameter candidates.
 - cf. The parameter space is maintained as a multi-dimensional histogram (discretization).
 - Score: The number of hits by data
 - Selection: Finding a peak on the histogram after voting
 - **RANSAC family**
 - A sample of data votes a single parameter candidate.
 - Score: The number of inlier candidates (whose error is within threshold)
 - Selection: Keeping the best model during RANSAC's iterations
 - cf. RANSAC involves many iterations of parameter estimation and error calculation.
- Top-down Approaches e.g. graph SLAM, multi-view reconstruction
 - **M-estimator**
 - All data aims to find the best parameter (from its initial guess).
 - Score: A cost function
 - The cost function includes a truncated loss function.
 - Selection: Minimizing the cost function (following its gradient)
 - cf. Nonlinear optimization is computationally heavy and leads to a local minima.

Multi-view Geometry

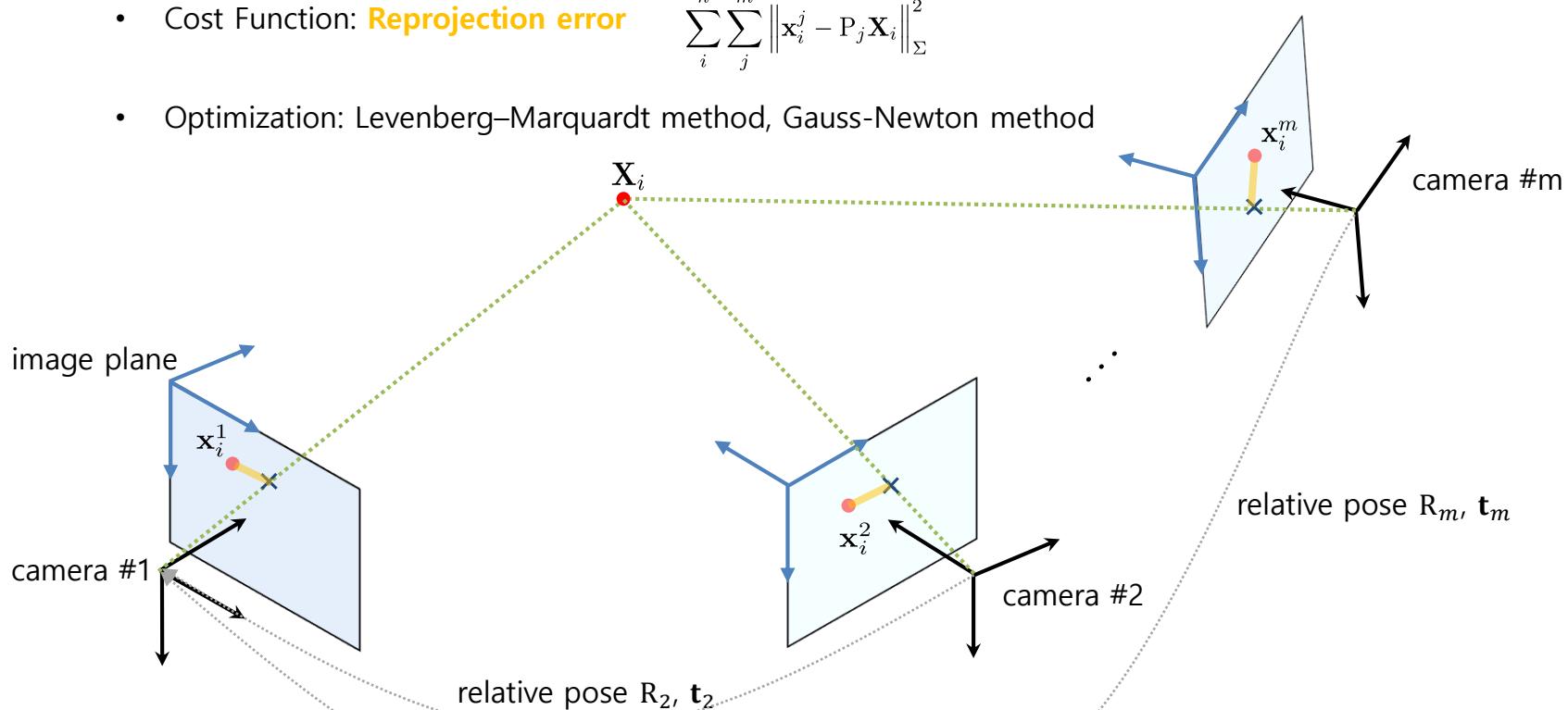


The Stanford Multi-Camera Array,
Stanford Computer Graphics Lab.

Bundle Adjustment

▪ Bundle Adjustment

- Unknown: Position of 3D points and each camera's relative pose ($6n + 3m$ DoF)
- Given: Point correspondence, camera matrices, position of 3D points, and each camera's relative pose
cf. initial values
- Constraints: $n \times m \times$ projection $\mathbf{x}_i^j = \mathbf{P}_j \mathbf{X}_i = \mathbf{K}_j [\mathbf{R}_j \mid \mathbf{t}_j] \mathbf{X}_i$
- Solution: **Non-linear least-square optimization**
 - Cost Function: **Reprojection error** $\sum_i^n \sum_j^m \| \mathbf{x}_i^j - \mathbf{P}_j \mathbf{X}_i \|^2_{\Sigma}$
 - Optimization: Levenberg–Marquardt method, Gauss–Newton method



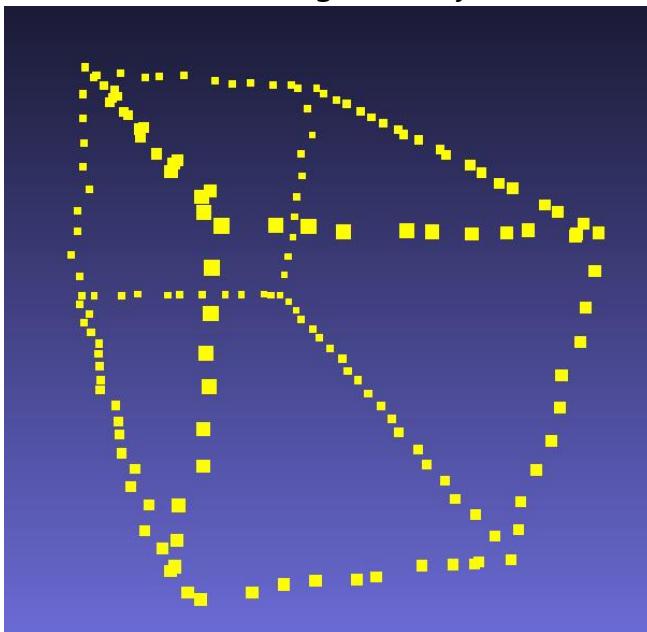
Bundle Adjustment

- **Bundle Adjustment**
 - Unknown: Position of 3D points and each camera's relative pose ($6n + 3m$ DoF)
 - Given: Point correspondence, camera matrices, position of 3D points, and each camera's relative pose
 - Constraints: $n \times m \times$ projection $\mathbf{x}_i^j = \mathbf{P}_j \mathbf{X}_i = \mathbf{K}_j [\mathbf{R}_j \mid \mathbf{t}_j] \mathbf{X}_i$ cf. initial values
 - Solution: **Non-linear least-square optimization**
 - Cost Function: **Reprojection error** $\sum_i^n \sum_j^m v_i^j \left\| \mathbf{x}_i^j - \mathbf{P}_j \mathbf{X}_i \right\|_{\Sigma}^2$ where $v_i^j = \begin{cases} 1 & \text{if } \mathbf{x}_i^j \text{ is visible} \\ 0 & \text{otherwise} \end{cases}$
 - Optimization: Levenberg–Marquardt method, Gauss–Newton method
- **Bundle Adjustment and Optimization Tools**
 - OpenCV: No function (but [cvSBA](#) is available as a [SBA](#) wrapper.)
 - Bundle adjustment: [SBA](#) (Sparse Bundle Adjustment), [SSBA](#) (Simple Sparse Bundle Adjustment), [pba](#) (Multicore Bundle Adjustment)
 - (Graph) optimization: [g2o](#) (General Graph Optimization), [iSAM](#) (Incremental Smoothing and Mapping), [GTSAM](#), [Ceres Solver](#) (Google)

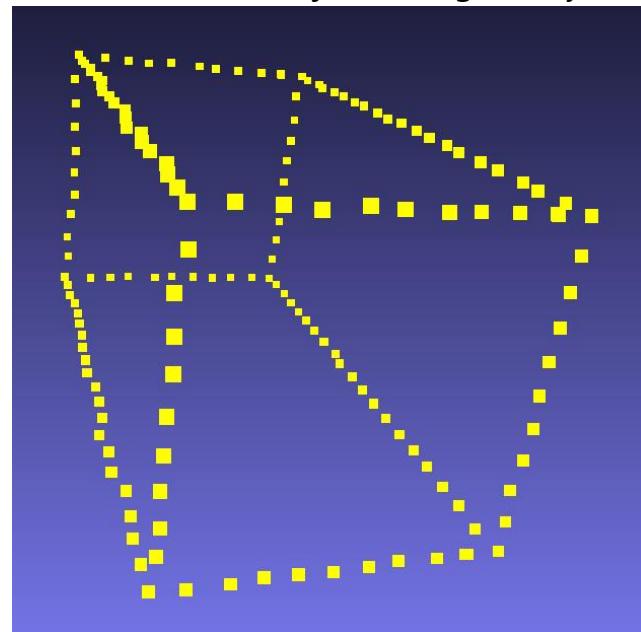
Bundle Adjustment

- Example: **Bundle Adjustment (Global)** [bundle_adjustment_global.cpp, bundle_adjustment.hpp]

Result: "triangulation.xyz"

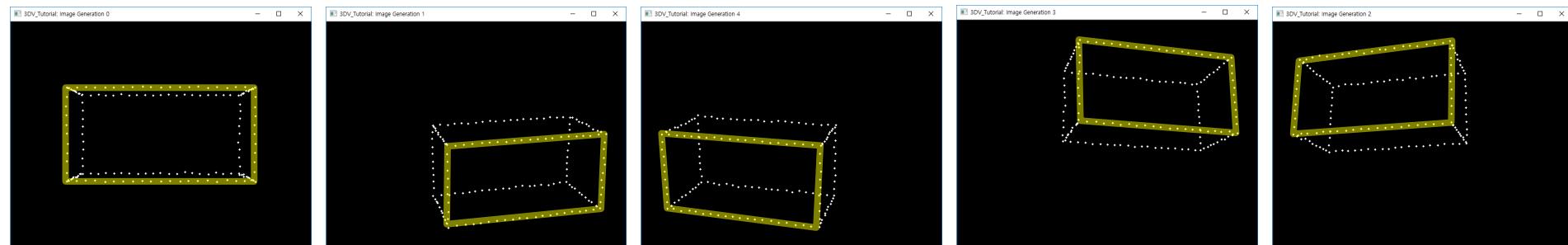


Result: "bundle_adjustment_global.xyz"



[Given]

- **Intrinsic parameters:** camera matrices K_j (all cameras have same and fixed camera matrix.)
- **2D point correspondence \mathbf{x}_i^j** (all points are visible on all camera views.)



[Unknown]

- **Extrinsic parameters:** camera poses

$$\mathbf{R}_0 = \mathbf{I}_{3 \times 3}$$

$$\mathbf{t}_0 = [0, 0, 0]^T$$

$$\mathbf{R}_1 = \mathbf{I}_{3 \times 3}$$

$$\mathbf{t}_1 = [0, 0, 0]^T$$

$$\mathbf{R}_2 = \mathbf{I}_{3 \times 3}$$

$$\mathbf{t}_2 = [0, 0, 0]^T$$

$$\mathbf{R}_3 = \mathbf{I}_{3 \times 3}$$

$$\mathbf{t}_3 = [0, 0, 0]^T$$

$$\mathbf{R}_4 = \mathbf{I}_{3 \times 3}$$

$$\mathbf{t}_4 = [0, 0, 0]^T$$

- **3D points $\mathbf{X}_i = [0, 0, 5.5]^T$** initial values

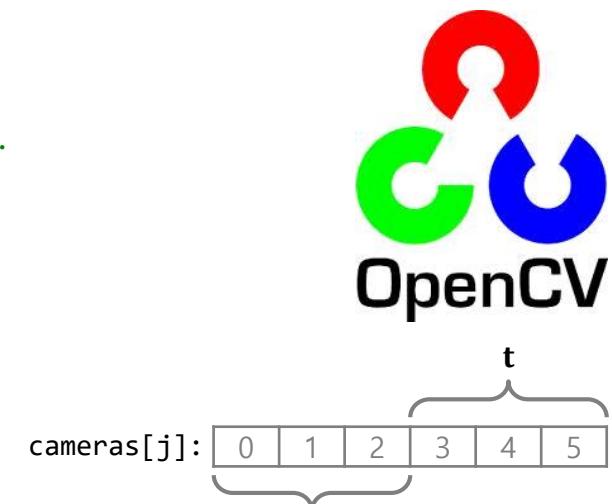
Ceres Solver (non-linear least-square optimization)

$$\sum_i^n \sum_j^m v_i^j \left\| \mathbf{x}_i^j - \mathbf{P}_j \mathbf{X}_i \right\|_{\Sigma}^2 \quad \text{where} \quad v_i^j = 1 \quad (\because \text{all points are visible})$$

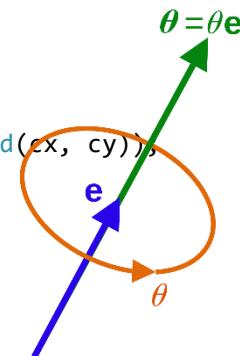
```

1. #include "bundle_adjustment.hpp"
2. int main()
3. {
4.     // cf. You need to run 'image_formation.cpp' to generate point observation.
5.     const char* input = "image_formation%d.xyz";
6.     int input_num = 5;
7.     double f = 1000, cx = 320, cy = 240;
8.
9.     // Load 2D points observed from multiple views
10.    std::vector<std::vector<cv::Point2d>> xs;
11.    ...
12.
13.    // Initialize cameras and 3D points
14.    std::vector<cv::Vec6d> cameras(xs.size());
15.    std::vector<cv::Point3d> Xs(xs.front().size(), cv::Point3d(0, 0, 5.5));
16.
17.    // Run bundle adjustment
18.    ceres::Problem ba;
19.    for (size_t j = 0; j < xs.size(); j++)
20.    {
21.        for (size_t i = 0; i < xs[j].size(); i++)
22.        {
23.            ceres::CostFunction* cost_func = ReprojectionError::create(xs[j][i], f, cv::Point2d(cx, cy));
24.            double* view = (double*)&(cameras[j]);
25.            double* X = (double*)&(Xs[i]);
26.            ba.AddResidualBlock(cost_func, NULL, view, X);
27.        }
28.    }
29.    ceres::Solver::Options options;
30.    options.linear_solver_type = ceres::ITERATIVE_SCHUR;
31.    options.num_threads = 8;
32.    options.minimizer_progress_to_stdout = true;
33.    ceres::Solver::Summary summary;
34.    ceres::Solve(options, &ba, &summary);
35.
36.    // Store the 3D points and camera pose to an XYZ file
37.    ...
38.    cv::Vec3d rvec(cameras[j][0], cameras[j][1], cameras[j][2]), t(cameras[j][3], cameras[j][4], cameras[j][5]);
39.    cv::Matx33d R;
40.    cv::Rodrigues(rvec, R);
41.
42.    ...
43.    return 0;
44. }

```



Axis-angle representation of R
(a.k.a. Rodrigues' rotation formula)





```
1. #include "opencv2/opencv.hpp"
2. #include "ceres/ceres.h"
3. #include "ceres/rotation.h"

4. // Reprojection error for bundle adjustment
5. struct ReprojectionError
6. {
7.     ReprojectionError(const cv::Point2d& _x, double _f, cv::Point2d& _c) : x(_x), f(_f), c(_c) { }

8.     template <typename T>
9.     bool operator()(const T* const camera, const T* const point, T* residuals) const
10.    {
11.        //  $X' = R*X + t$ 
12.        T X[3];
13.        ceres::AngleAxisRotatePoint(camera, point, X);
14.        X[0] += camera[3];
15.        X[1] += camera[4];
16.        X[2] += camera[5];

17.        //  $x' = K*X'$ 
18.        T x_p = f * X[0] / X[2] + cx;
19.        T y_p = f * X[1] / X[2] + cy;

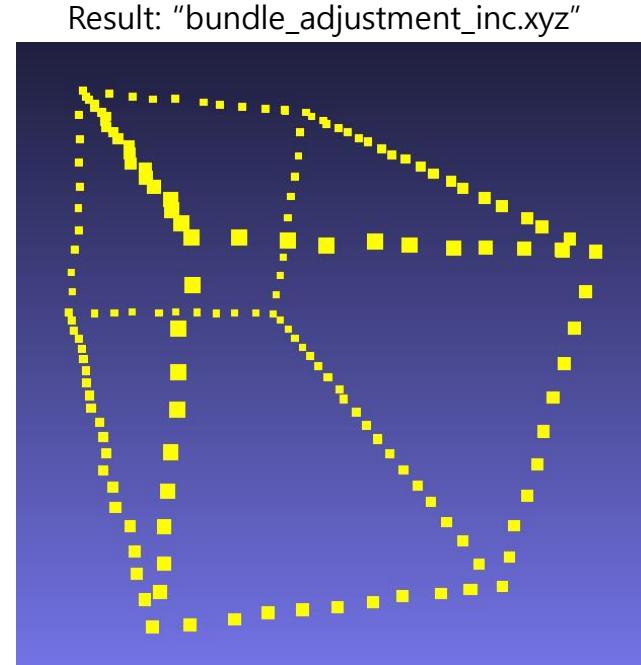
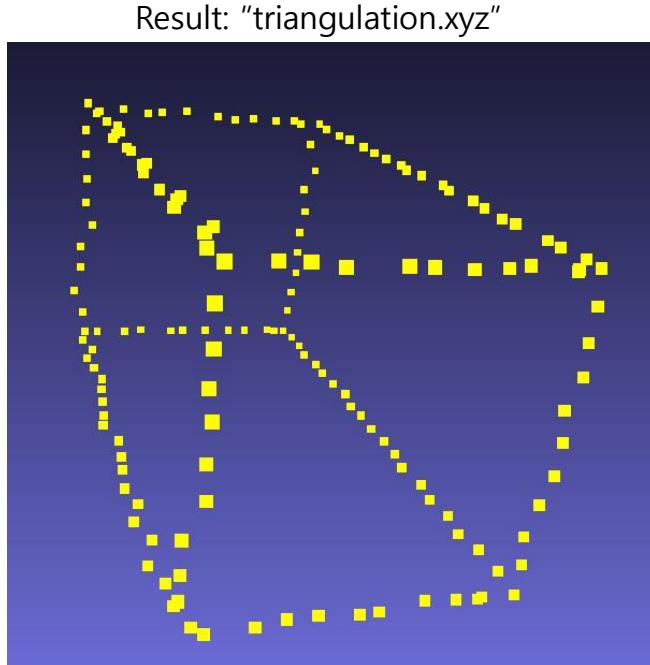
20.        // residual = x - x'
21.        residuals[0] = T(x.x) - x_p;           x - PX
22.        residuals[1] = T(x.y) - y_p;
23.        return true;
24.    }

25.    static ceres::CostFunction* create(const cv::Point2d& _x, double _f, cv::Point2d& _c)
26.    {
27.        return (new ceres::AutoDiffCostFunction<ReprojectionError, 2, 6, 3>(new ReprojectionError(_x, _f, _c)));
28.    }

29. private:
30.     const cv::Point2d x;
31.     const double f;
32.     const cv::Point2d c;
33. };
```

Bundle Adjustment

- Example: **Bundle Adjustment (Incremental)** [bundle_adjustment_inc.cpp, bundle_adjustment.hpp]



[bundle_adjustment_global.xyz]

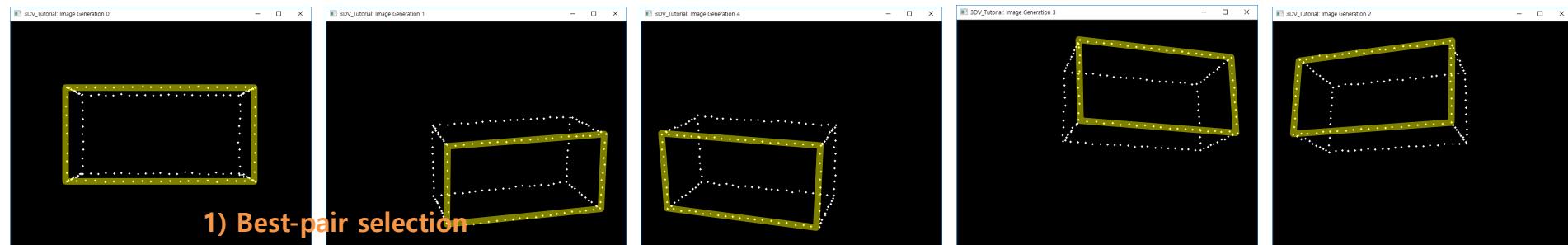
- # of iterations: 20
- Reprojection error: 0.113 ← **2637.626**

[bundle_adjustment_inc.xyz]

- # of iterations: **17 + 4 + 6**
- Reprojection error: 0.113 ← **722.350**

[Given]

- **Intrinsic parameters:** camera matrices K_j (all cameras have same and fixed camera matrix.)
- **2D point correspondence x_i^j** (all points are visible on all camera views.)



[Unknown]

- **Extrinsic parameters:** camera poses

$$R_0 = I_{3 \times 3}$$

$$t_0 = [0, 0, 0]^\top$$

$$R_1, t_1$$

$$R_2, t_2$$

$$R_3, t_3$$

$$R_4, t_4$$

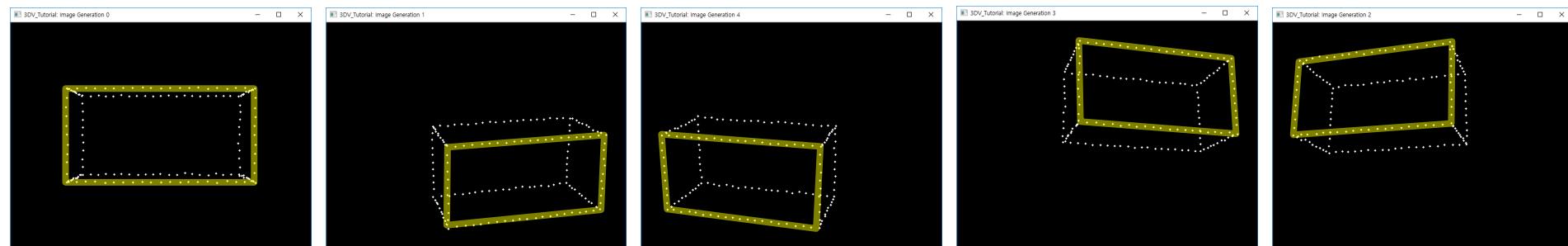
2) Epipolar geometry

3) Triangulation

- **3D points X_i**

[Given]

- **Intrinsic parameters:** camera matrices K_j (all cameras have same and fixed camera matrix.)
- **2D point correspondence x_i^j** (all points are visible on all camera views.)



[Unknown]

- **Extrinsic parameters:** camera poses

$$\begin{aligned} R_0 &= I_{3 \times 3} \\ t_0 &= [0, 0, 0]^\top \end{aligned}$$

$$R_1, t_1$$

$$R_2, t_2$$

$$R_3, t_3$$

$$R_4, t_4$$

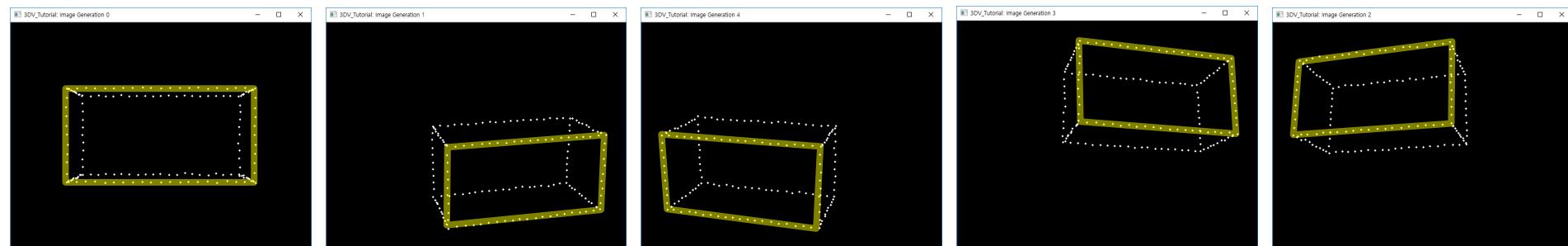
- **3D points X_i**

5) Perspective-n-point (PnP)

6) Triangulation

[Given]

- **Intrinsic parameters:** camera matrices K_j (all cameras have same and fixed camera matrix.)
- **2D point correspondence x_i^j** (all points are visible on all camera views.)



[Unknown]

- **Extrinsic parameters:** camera poses

4) Next-view selection

Repeat 4 – 7 until the last image

$$R_0 = I_{3 \times 3}$$

$$t_0 = [0, 0, 0]^\top$$

$$R_1, t_1$$

$$R_2, t_2$$

$$R_3, t_3$$

$$R_4, t_4$$

- 3D points X_i

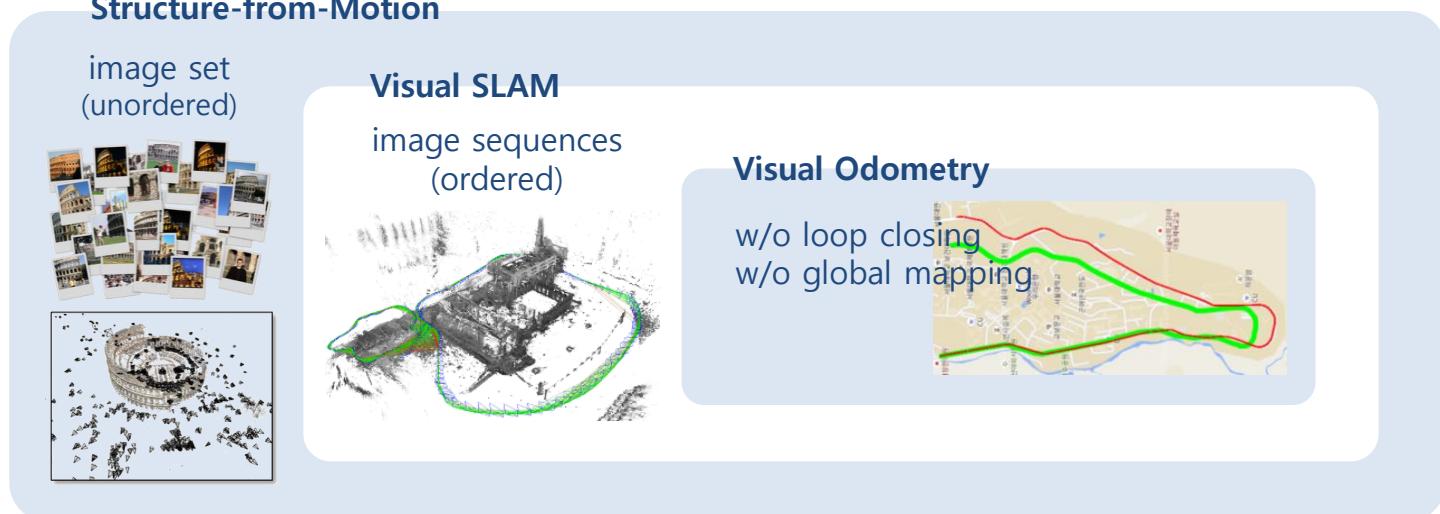
5) Perspective-n-point (PnP)

6) Triangulation

7) Bundle Adjustment

Applications

- **Structure-from-Motion (SfM)** → 3D Reconstruction, Photo Browsing
 - [Bundler](#), [COLMAP](#), [MVE](#), [Theia](#), [openMVG](#), [OpenSfM](#), [Toy SfM](#) / [VisualSfM](#) (GUI, binary only)
- **Visual SLAM** → Augmented Reality, Navigation (Mapping and Localization)
 - [PTAM](#) (Parallel Tracking and Mapping), [DTAM](#) (Dense Tracking and Mapping), [ORB-SLAM2](#), [LSD-SLAM](#)
 - cf. Visual loop closure (a.k.a. visual place recognition): [DBoW2](#), [FBoW](#), [PoseNet](#), [NetVLAD](#)
- **Visual Odometry** → Navigation (Localization)
 - [LIBVISO2](#) (C++ Library for Visual Odometry 2), [SVO](#) (Semi-direct Monocular Visual Odometry), [DVO](#) (Direct Sparse Odometry), [DeepVO](#), [UnDeepVO](#)



Structure-from-Motion

- Example: **Structure-from-Motion (Global)** [sfm_global.cpp]

0) Load images and extract features

0) Match features and find good matches (which has enough inliers)



Structure-from-Motion

- Example: **Structure-from-Motion (Global)** [sfm_global.cpp]

0) Load images and extract features

0) Match features and find good matches (which has enough inliers)



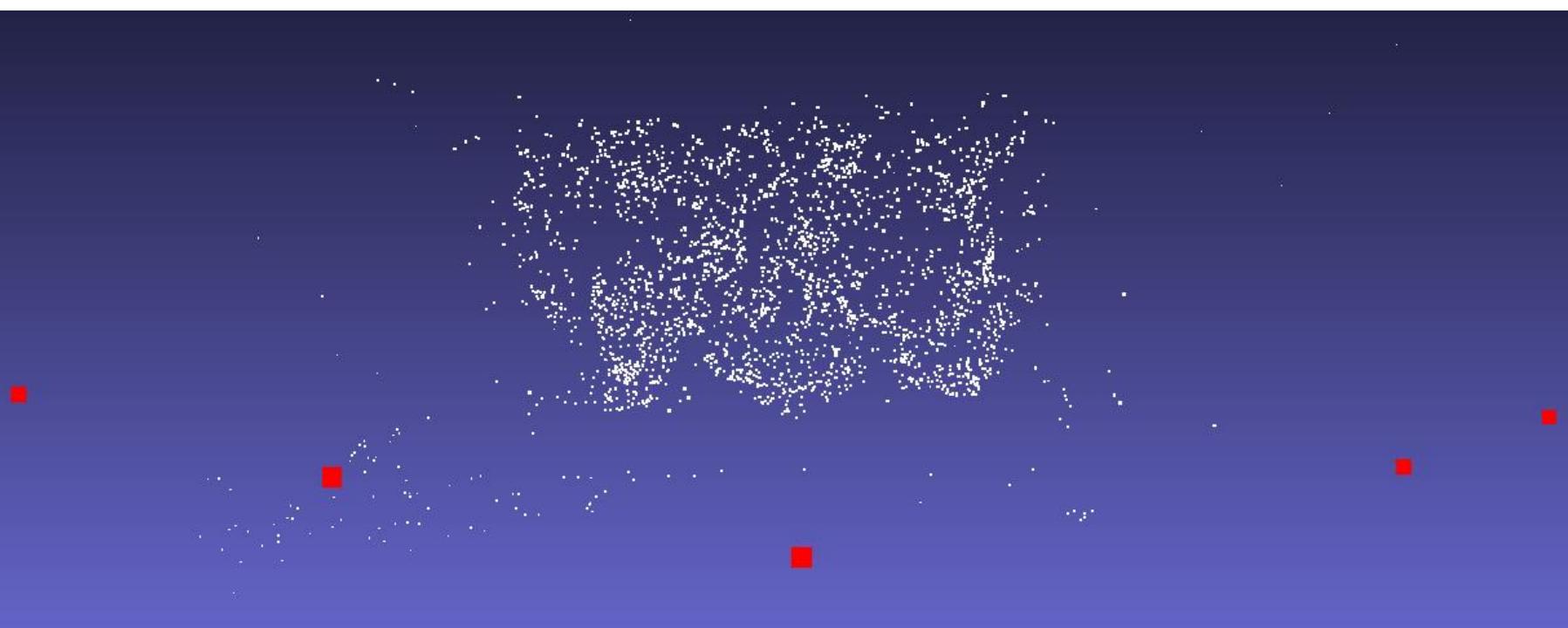
- 1) Initialize cameras (R, t, f, \dots)
- 2) Initialize 3D points and build a visibility graph
- 3) Optimize camera pose and 3D points together (BA)

$$K_i = \begin{bmatrix} f & 0 & w/2 \\ 0 & f & h/2 \\ 0 & 0 & 1 \end{bmatrix}, \quad R_i = I_{3 \times 3}, \quad \text{and} \quad t_i = [0, 0, 0]^\top$$
$$X_i = [0, 0, 2]^\top$$

Ceres Solver (non-linear least-square optimization)

Structure-from-Motion

- Example: **Structure-from-Motion (Global)** [sfm_global.cpp]



Structure-from-Motion

- Example: **Structure-from-Motion (Incremental)** [sfm_inc.cpp]

0) Build a viewing graph (well-matched pairs) while finding inliers

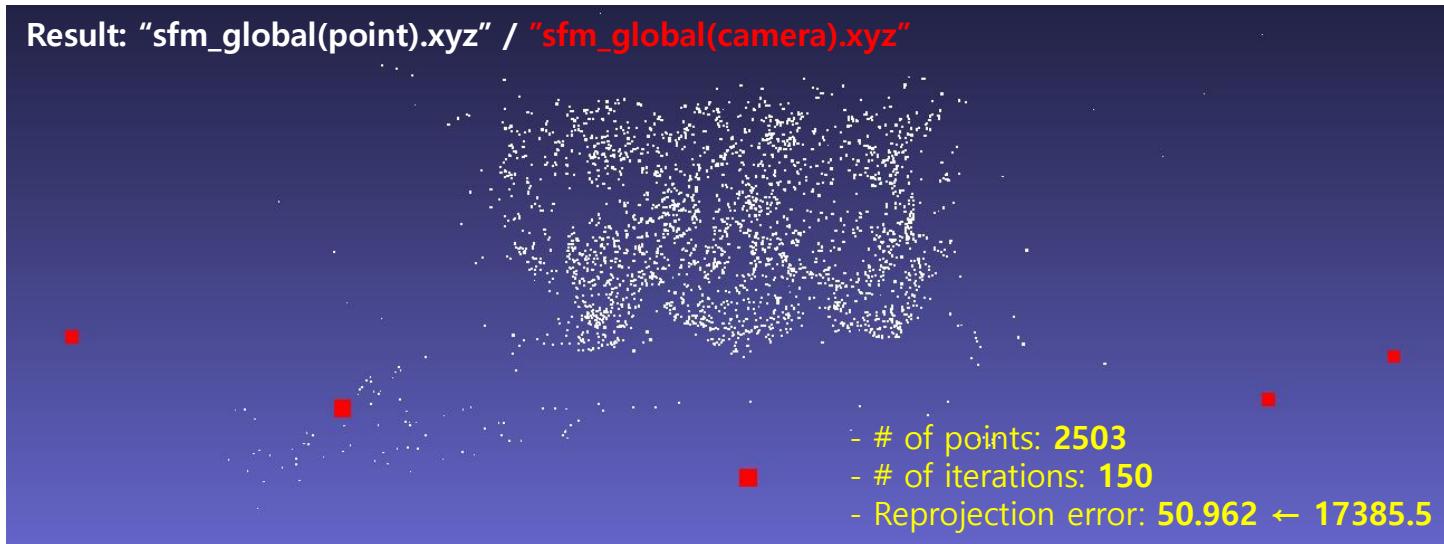


- 1) Select the best pair
- 2) Estimate relative pose from the best two views (epipolar geometry)
- 3) Reconstruct 3D points of the best two views (triangulation)
- 4) Select the next image to add
Separate points into known and unknown for PnP (known) and triangulation (unknown)
- 5) Estimate relative pose of the next view (PnP)
- 6) Reconstruct newly observed 3D points (triangulation)
- 7) Optimize camera pose and 3D points together (BA)

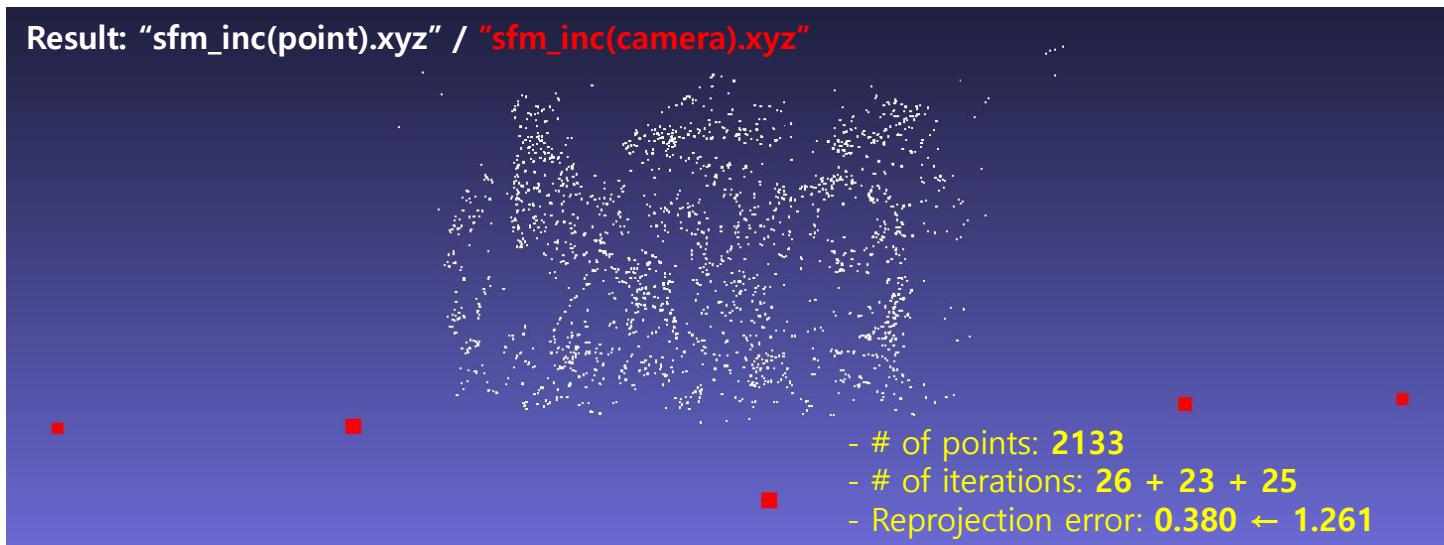
Incrementally add more views

Structure-from-Motion

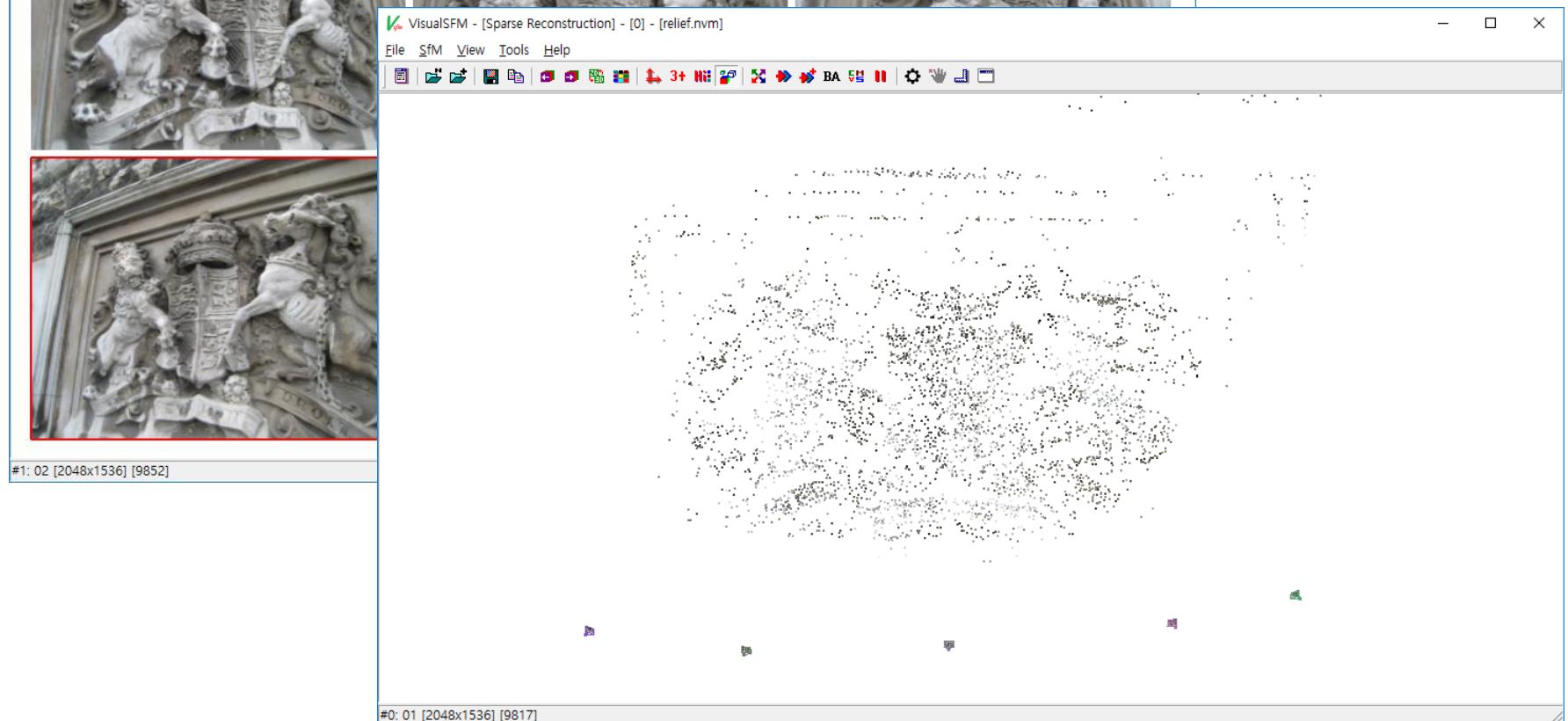
Result: "sfm_global(point).xyz" / "sfm_global(camera).xyz"



Result: "sfm_inc(point).xyz" / "sfm_inc(camera).xyz"

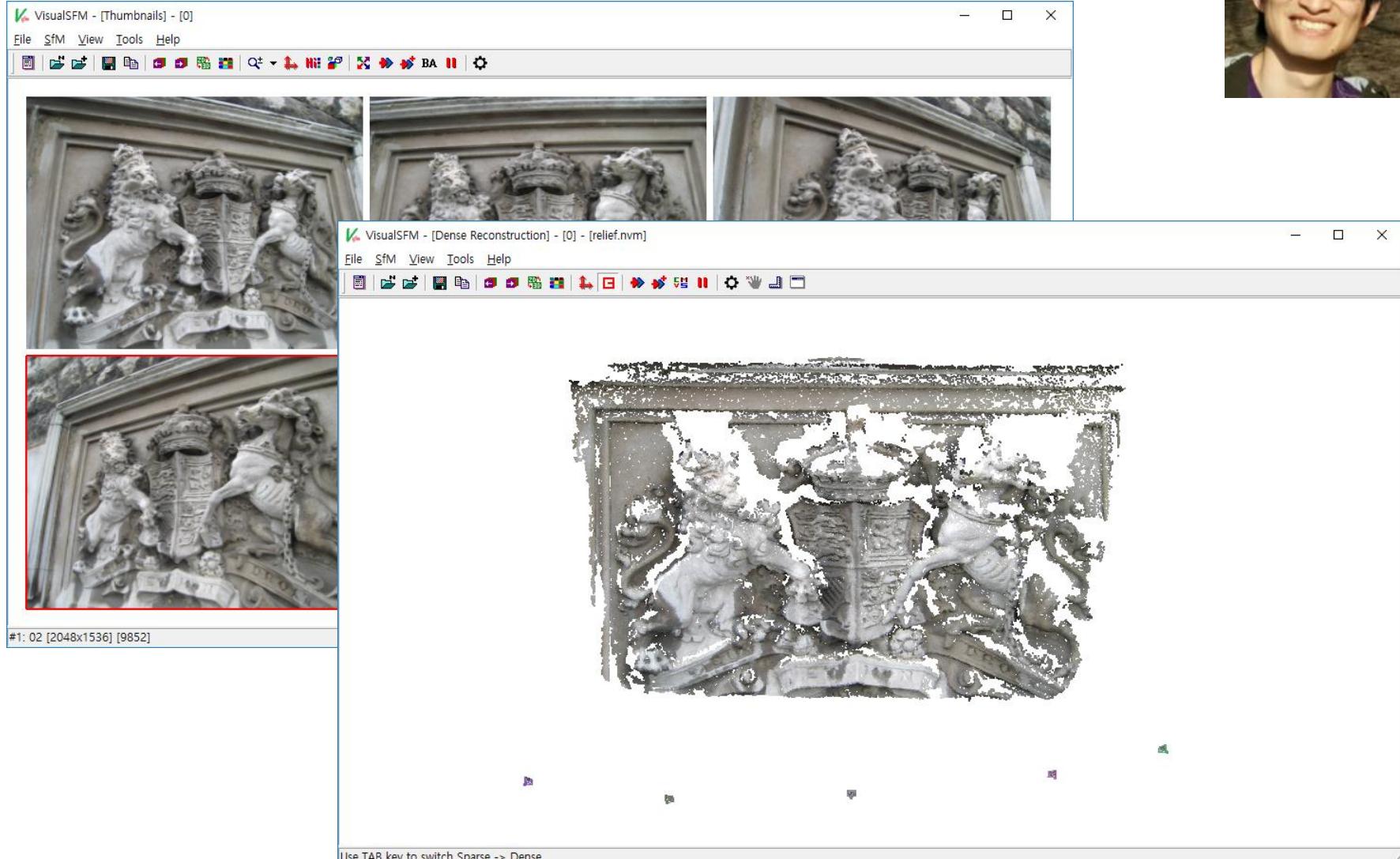


Structure-from-Motion using VisualSfM



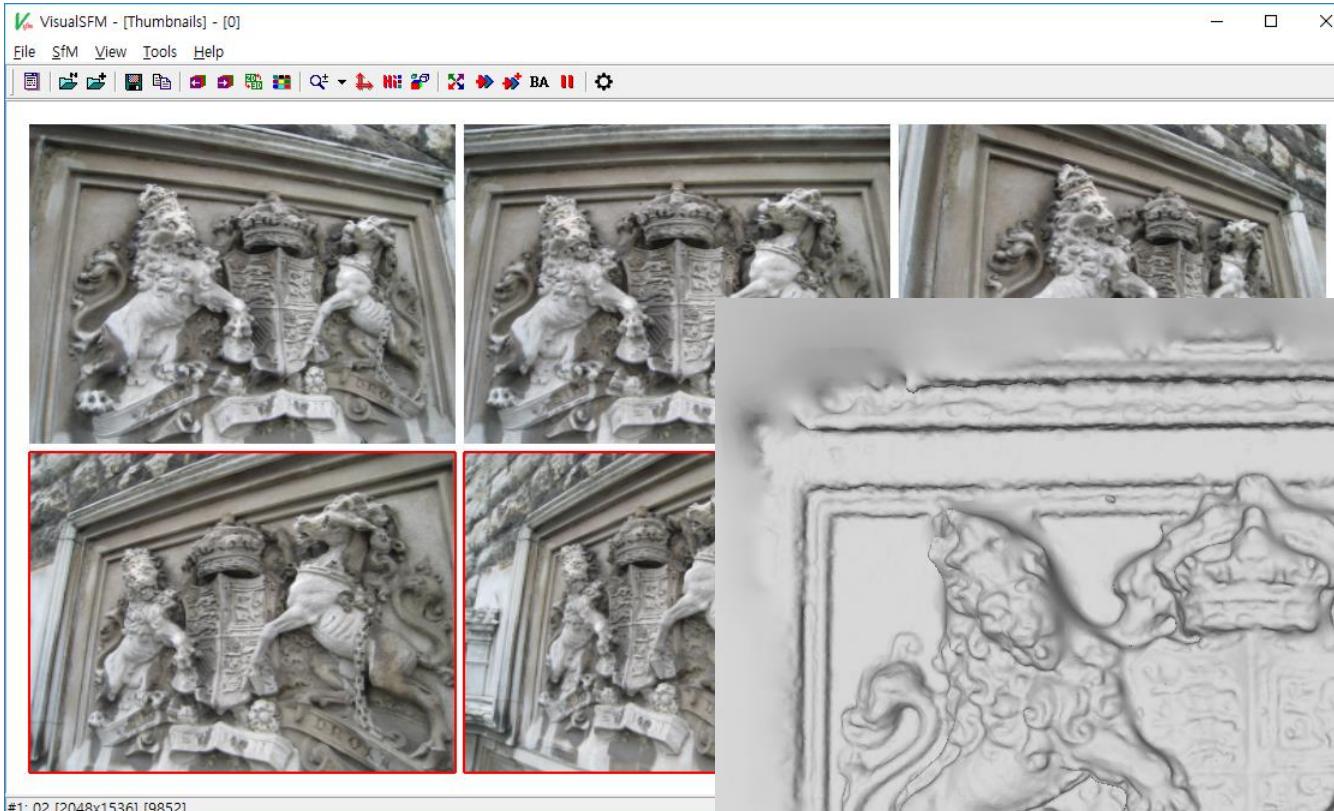
Sparse 3D Reconstruction by VisualSfM (SiftGPU and pba)

Structure-from-Motion using VisualSfM



Dense 3D Reconstruction by VisualSfM (CMVS and PMVS2)

Structure-from-Motion using VisualSfM

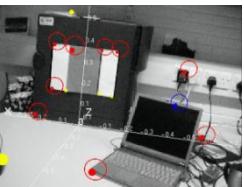


Surface Reconstruction by VisualSfM (Poisson Surface Reconstruction) [\[URL\]](#)



Visual Odometry and SLAM: History

Monocular Camera



MonoSLAM (2003): features, Bayesian filtering
V.S.

SLAM in Robotics



SfM in Computer Vision



PTAM (2007): more features, bundle adjustment (keyframes)

DTAM (2011): dense volumetric

RGB-D Camera



KineticFusion (2011): dense volumetric

ORB-SLAM (2014)

better feature management

LSD-SLAM (2014)

direct, semi-dense

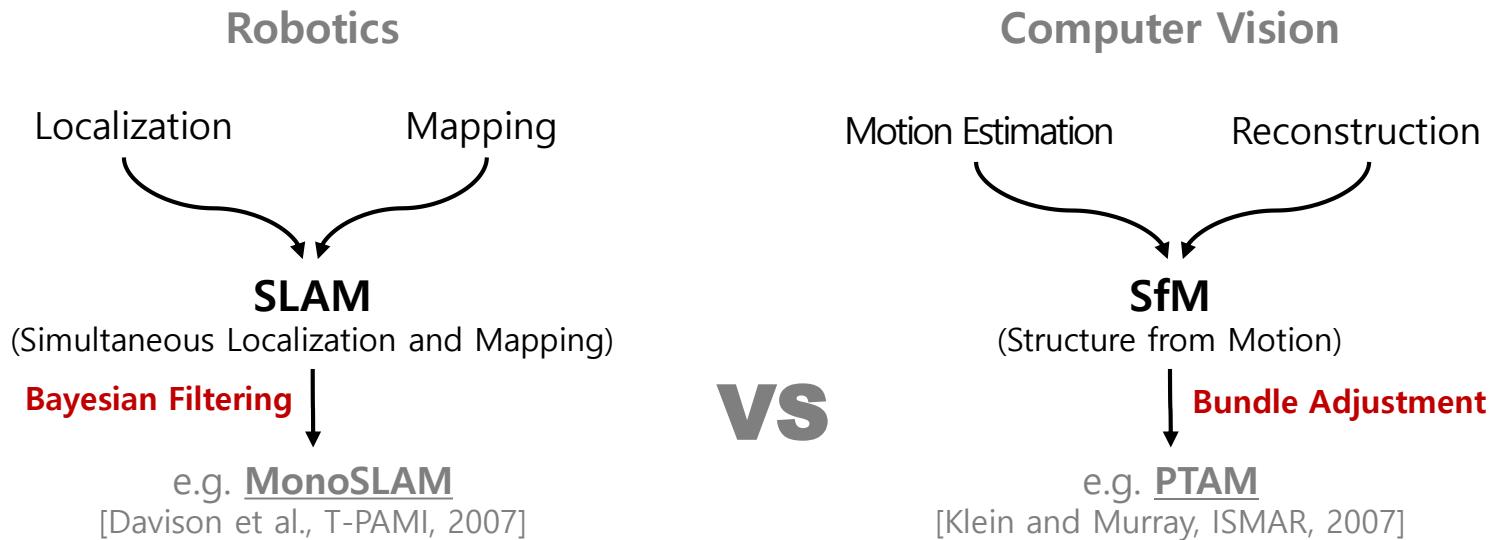


ElasticFusion (2015)
dense surfels



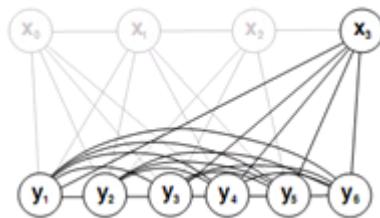
DynamicFusion (2015)
dynamic dense

Paradigm #1: Bayesian Filtering v.s. Bundle Adjustment

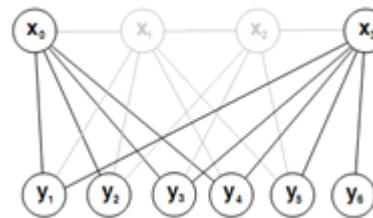


“Real-time Monocular SLAM: Why Filter?”

[Strasdat et al., ICRA, 2010]



1. Global optimization
2. # of features (100 vs 4000)



Paradigm #2: Feature-based Method v.s. Direct Method

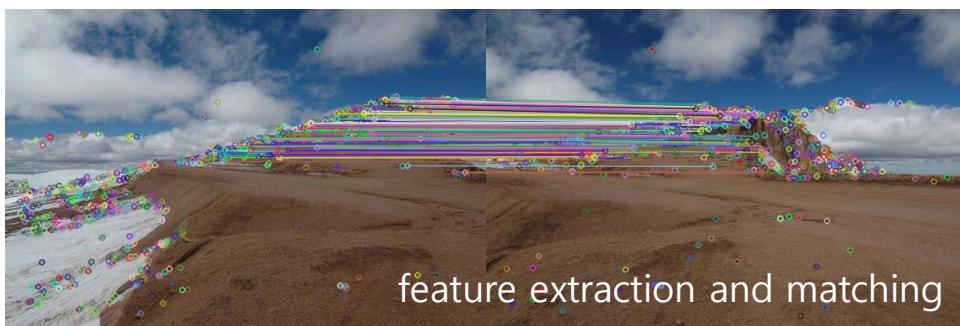
- e.g. Image stitching



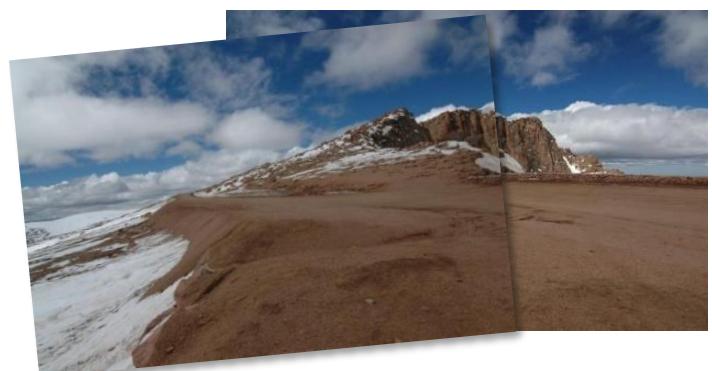
Feature-based Method

VS

Direct Method



$$\arg \min_H \sum_i \left\| H \mathbf{x}_i - \mathbf{x}'_i \right\|^2$$



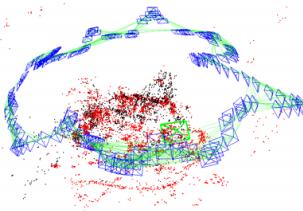
$$\arg \min_H \sum_{u,v} \left\| I \begin{bmatrix} u \\ 1 \end{bmatrix} - I' \left(H \begin{bmatrix} u \\ 1 \end{bmatrix} \right) \right\|^2$$

Paradigm #2: Feature-based Method v.s. Direct Method

Feature-based/Direct SLAM



Feature-Based SLAM

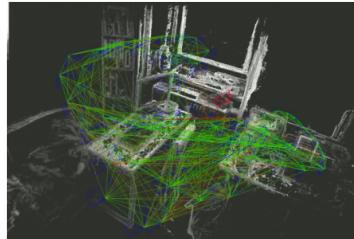


Minimize **Feature Reprojection Error**

Sparse Reconstruction

PTAM, ORB-SLAM

Direct SLAM



Minimize **Photometric Error**

Semi Dense / Dense Reconstruction

DTAM, LSD-SLAM, DPPTAM



The Future of Real-Time SLAM (ICCV'15 Workshop). Raúl Mur Ardal. University of Zaragoza.

Why should we still use features?



Robustness

Reliable two-view monocular initialization

Good invariance to viewpoint and illumination

Less affected by auto-gain and auto-exposure

Less affected by dynamic elements

Accuracy

Bundle adjustment (joint map-trajectory optimization)

Place Recognition (loop detection, relocalization)

Bags of Words

But sparse reconstructions ...

Comparison



Feature-Based

Input Images

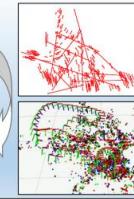


Extract & Match Features (SIFT / SURF / BRIEF / ...)



abstract images to feature observations

Track:
min. reprojection error
(point distances)



Map:
est. feature-parameters
(3D points / normals)

Chiuso '02, Nistér '04, Eade '06, Klein '06,
Davison '07, Strasdat '10, Mur-Artal '14, ...

Jakob Engel

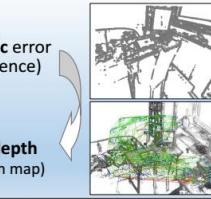
Direct

Input Images



keep full image

Track:
min. photometric error
(intensity difference)



Map:
est. per-pixel depth
(semi-dense depth map)

Matthies '88, Hanna '91, Comport '06,
Newcombe '11, Engel '13, ...

Semi-Dense Direct SLAM

23

Comparison



Feature-Based

can only use & reconstruct corners

faster

flexible: outliers can be removed retroactively.

robust to inconsistencies in the model/system (rolling shutter).

decisions (KP detection) based on less complete information.

no need for good initialization.

~20+ years of intensive research

can use & reconstruct whole image

slower (but good for parallelism)

inflexible: difficult to remove outliers retroactively.

not robust to inconsistencies in the model/system (rolling shutter).

decision (linearization point) based on more complete information.

needs good initialization.

~4 years of research (+5years 25 years ago)

145



The Future of Real-Time SLAM (ICCV'15 Workshop). Raúl Mur Ardal. University of Zaragoza.

Jakob Engel

Semi-Dense Direct SLAM

24

Why Visual Odometry?



Visual odometry

VS



Wheel odometry

- + direct motion measure
- + six degree-of-freedoms
- + easy to install
- heavy computation
- visual disturbance (e.g. moving objects)

- indirect motion measure (e.g. slippage)
- two degree-of-freedoms
- necessary to be on rotor/shaft
- + simple calculation

Why Visual Odometry?



Visual odometry

VS



Wheel odometry



- indirect motion measure (e.g. slippage)
- two degree-of-freedoms
- necessary to be on rotor/shaft
- + simple calculation

no wheels / rough terrains, in-the-sky, under-the-water

Why Visual Odometry?



Visual odometry

VS



Wheel odometry



- indirect motion measure (e.g. slippage)
- two degree-of-freedoms
- necessary to be on rotor/shaft
- + simple calculation

no wheels / rough terrains, in-the-sky, under-the-water



Visual Odometry

VS



Visual SLAM

no assumption on trajectories
→ navigation / large space (outdoor)

closed-loop is preferred for convergence
→ mapping / small space (indoor, campus)

Why Visual Odometry?



Visual odometry

VS



Wheel odometry



- indirect motion measure (e.g. slippage)
- two degree-of-freedoms
- necessary to be on rotor/shaft
- + simple calculation

no wheels / rough terrains, in-the-sky, under-the-water



Visual Odometry

VS



Visual SLAM

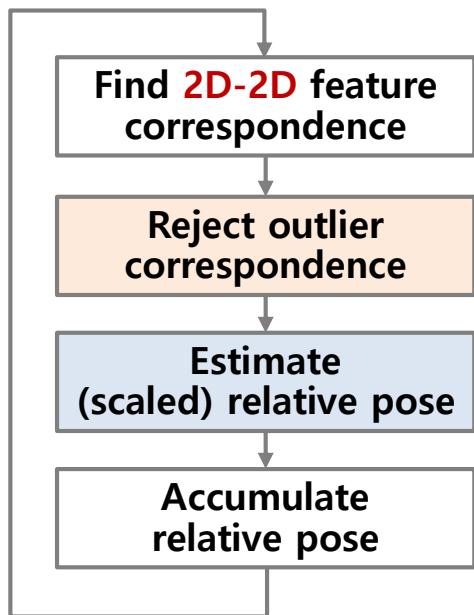


closed-loop is preferred for convergence
→ mapping / small space (indoor, campus)

real navigation situations

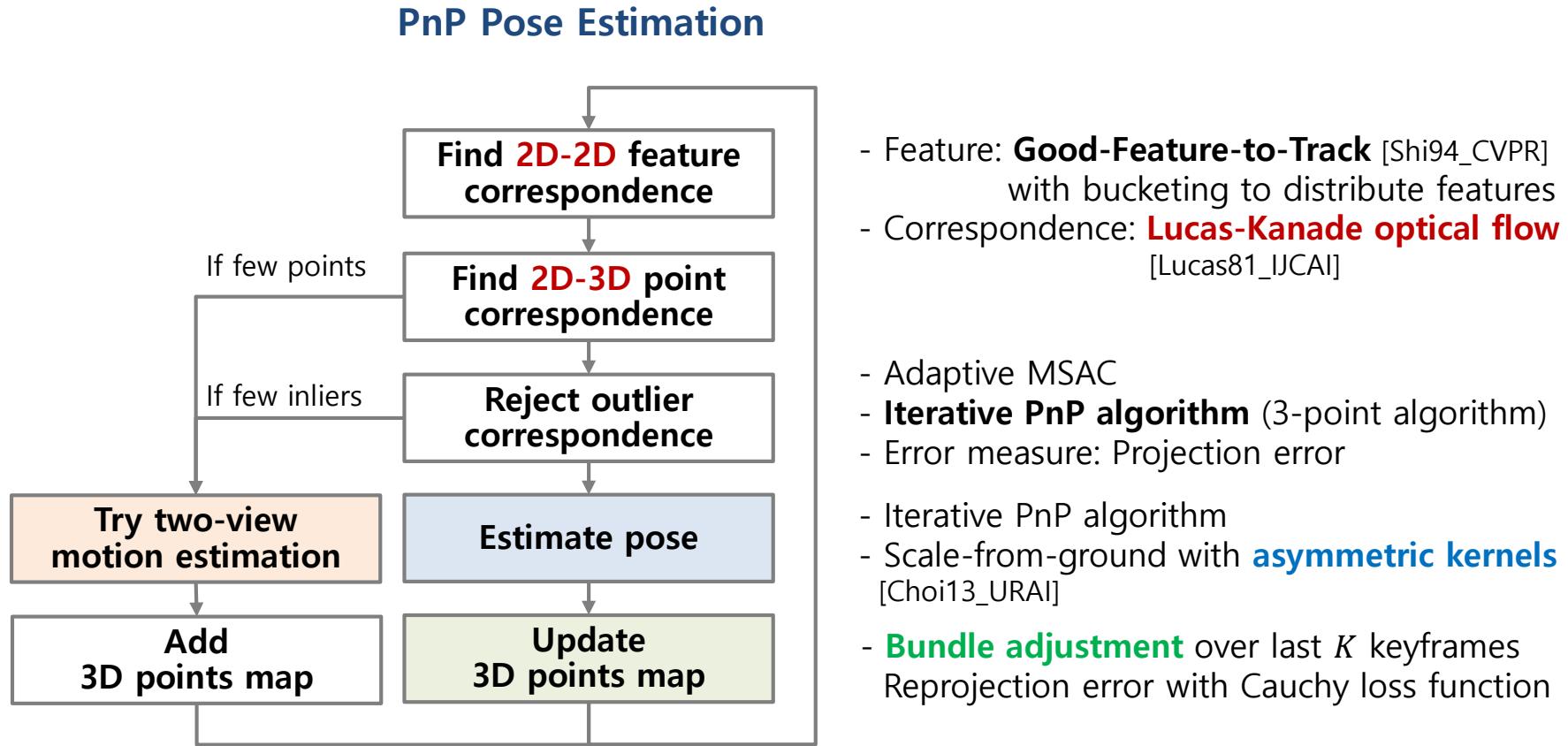
Feature-based Monocular Visual Odometry

Two-view Motion Estimation

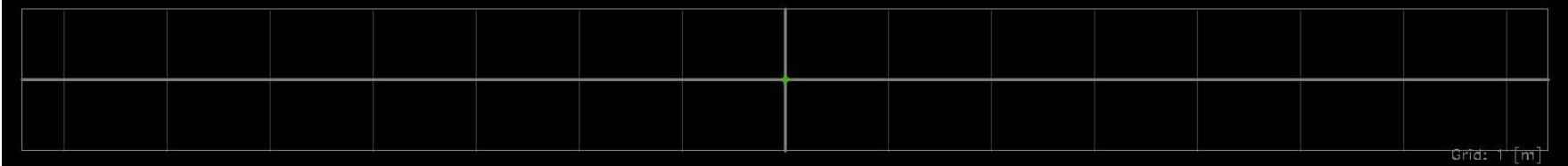
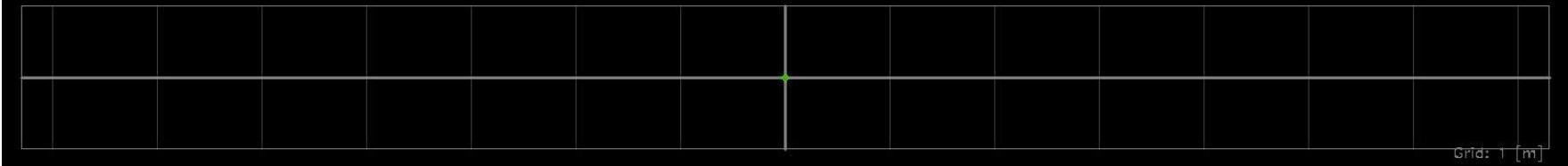


- Feature: **Good-Feature-to-Track** [Shi94_CVPR]
with bucketing to distribute features
- Correspondence: **Lucas-Kanade optical flow** [Lucas81_IJCAI]
- Adaptive MSAC [Choi09_IROS]
- **Iterative 5-point algorithm** [Choi15_IJCAS]
- Error measure: Sampson distance
- Normalized 8-point algorithm
- Scale-from-ground with **asymmetric kernels** [Choi13_URAI]

Feature-based Monocular Visual Odometry



Feature-based Monocular Visual Odometry



Summary

- **What is 3D Vision?**
- **Single-view Geometry**
 - **Camera Projection Model**
 - Pinhole Camera Model
 - Geometric Distortion Models
 - **General 2D-3D Geometry**
 - Camera Calibration
 - Absolute Camera Pose Estimation (PnP Problem)
- **Two-view Geometry**
 - **Planar 2D-2D Geometry (Projective Geometry)**
 - Planar Homography
 - **General 2D-2D Geometry (Epipolar Geometry)**
 - Fundamental/Essential Matrix
 - Relative Camera Pose Estimation
 - Triangulation (Point Localization)
- **Multi-view Geometry**
 - Bundle Adjustment (Non-linear Optimization)
 - Applications: Structure-from-motion, Visual SLAM, and Visual Odometry
- **Correspondence Problem**
 - Feature Correspondence: Feature Matching and Tracking
 - Robust Parameter Estimation: (Hough Transform), RANSAC, M-estimator

$$\begin{aligned} \therefore \mathbf{x} &= \mathbf{P}\mathbf{X} \quad (\mathbf{P} = \mathbf{K}[\mathbf{R} | \mathbf{t}]) \\ \mathbf{x} &= \mathbf{K}\hat{\mathbf{x}} \end{aligned}$$

$$\begin{aligned} \therefore \mathbf{x}' &= \mathbf{H}\mathbf{x} \\ \hat{\mathbf{x}}' &= \hat{\mathbf{H}}\hat{\mathbf{x}} \quad (\hat{\mathbf{H}} = \mathbf{R} + \frac{1}{d}\mathbf{t}\mathbf{n}^\top) \end{aligned}$$

$$\begin{aligned} \therefore \mathbf{x}'^\top \mathbf{F} \mathbf{x} &= 0 \\ \hat{\mathbf{x}}'^\top \mathbf{E} \hat{\mathbf{x}} &= 0 \quad (\mathbf{E} = [\mathbf{t}]_\times \mathbf{R}) \end{aligned}$$

$$\therefore \operatorname{argmin} \sum_i^n \sum_j^m \left\| \mathbf{x}_i^j - \mathbf{P}_j \mathbf{X}_i \right\|_{\Sigma}^2$$

Applications in Deep Learning Era

- **There are still many researches and applications.**

- 3D reconstruction
 - Real-time visual odometry/SLAM
 - Augmented reality (mixed reality), virtual reality



model-based problem solving

(e.g. calculating camera pose; minimizing a cost function)



unknown models and procedures

(e.g. recognizing objects, finding correspondence)

- 3D understanding of **results** from deep learning

- e.g. Bounding boxes from object recognition → metric position and size

- Designing **cost functions** for deep learning

- e.g. Left-right consistency in [MonoDepth](#) [CVPR 2017]
 - e.g. [Eigendecomposition-free training](#) of deep learning [ECCV 2018]

- Building **datasets** for deep learning

- e.g. [LIFT](#) [ECCV 2016]

Appendix: How to Solve Problems

▪ Linear Equations

- Inhomogeneous cases: Multiplying a pseudo-inverse

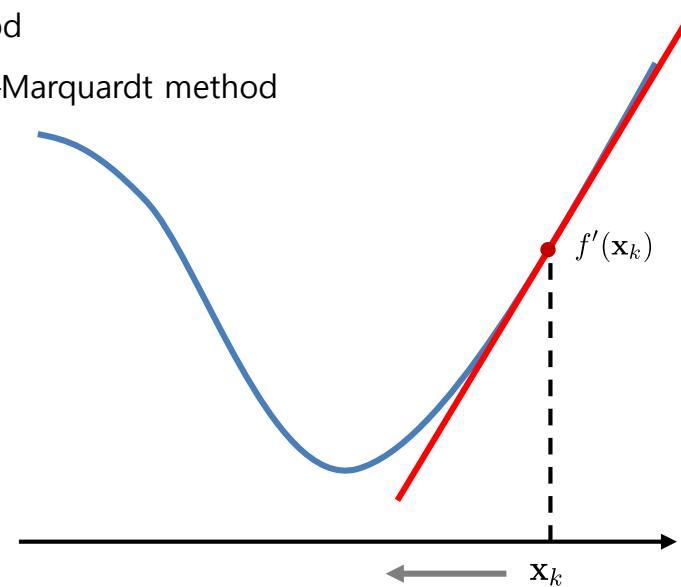
$$Ax = b \quad \longrightarrow \quad x = A^\dagger b$$

- Homogenous cases: Finding a null vector (or a vector which has the smallest singular value)

$$Ax = 0 \quad \longrightarrow \quad x = V [0, 0, \dots, 1]^\top \quad \text{where} \quad A = UDV^\top$$

▪ Non-linear Equations

- Non-linear optimization $\hat{x} = \underset{x}{\operatorname{argmin}} f(x)$
- General cases: Gradient-descent method, Newton method
- Least squares cases: Gauss-Newton method, Levenberg–Marquardt method



Appendix: Further Information

- **Beyond Point Features**
 - Other features: [OPVO](#), [Kimera](#)
 - Direct methods (w/o features): Already mentioned (including deep learning)
- ~~Real-time / Large-scale SfM~~
- **(Spatially / Temporally) Non-static SfM**
 - Deformable (or moving) objects: [Non-rigid SfM](#)
- **Depth, Object, and Semantic Recognition**
- **Sensor Fusion and New Sensors**
 - + Depth (ORB-SLAM2, LSD-SLAM, DVO have their variants with depth.)
 - RGB-D: [ElasticFusion](#), [RGB-D SLAM](#), [BaMVO](#), [RTAB-Map](#)
 - Stereo cameras: [S-PTAM](#), [ProSLAM](#)
 - LiDAR: [LIMO](#) / cf. [Cartographer](#)
 - + IMU: [OKVIS](#), [ROVIO](#), [VINS-Mono](#)
 - Visual-inertia Calibration: [Kalib](#)
 - + GPS
 - Omni-directional cameras: [Multi-FoV datasets](#)
 - Light-field cameras
 - Event camera: [ETAM](#)*
- **Minimal Solvers / Self-calibration**
 - [OpenGV](#), [Minimal Problems in Computer Vision](#)
- **Public Datasets and Evaluations**
 - Datasets: [Awesome Robotics Datasets](#)
 - Evaluations: [KITTI Odometry/SLAM Evaluation 2012](#), [GSLAM](#), [evo](#)