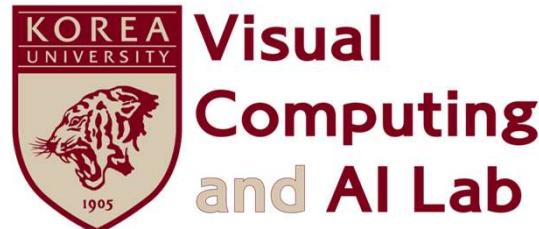


INTRODUCTION TO COMPUTER VISION

Lecture 4 – Camera Calibration

Gyeongsik Moon
Visual Computing and AI Lab
Korea University



Slides Credit: [Prof. Dr. Davide Scaramuzza](#)

1. Camera calibration

2. Camera localization

1. Camera calibration

Camera Calibration

3D → 2D
focal
principal

3D → 3D

Goal:

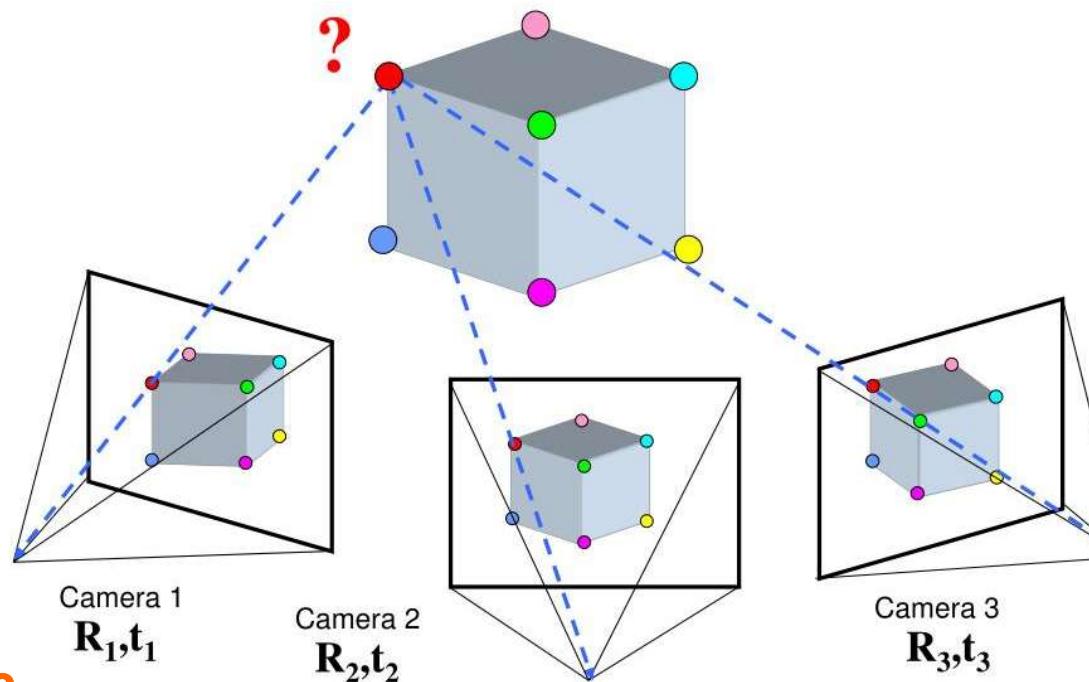
**Getting intrinsic/extrinsic
camera parameters**

The diagram illustrates the goal of camera calibration. At the top left, handwritten blue text reads "3D → 2D" above "focal principal", with a blue arrow pointing from "focal principal" down towards the word "Goal". At the top right, handwritten blue text reads "3D → 3D", with a blue arrow pointing from it up towards the word "Goal". Below these, the word "Goal:" is written in large, bold black font. Underneath "Goal:", the text "Getting intrinsic/extrinsic camera parameters" is also written in large, bold black font.

Multi-View Geometry (MVG) Theory

According to MVG (we'll learn this later), if we know

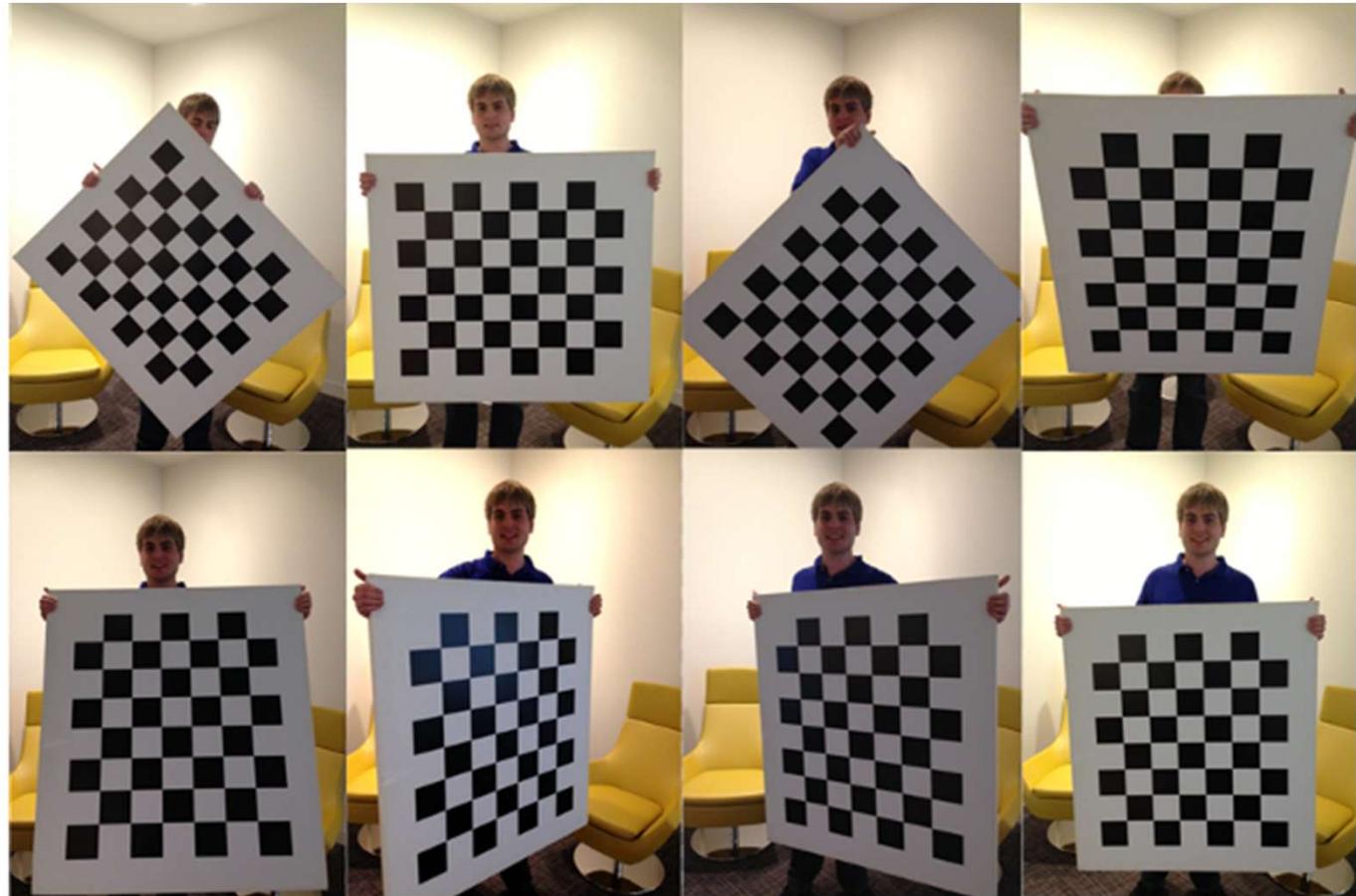
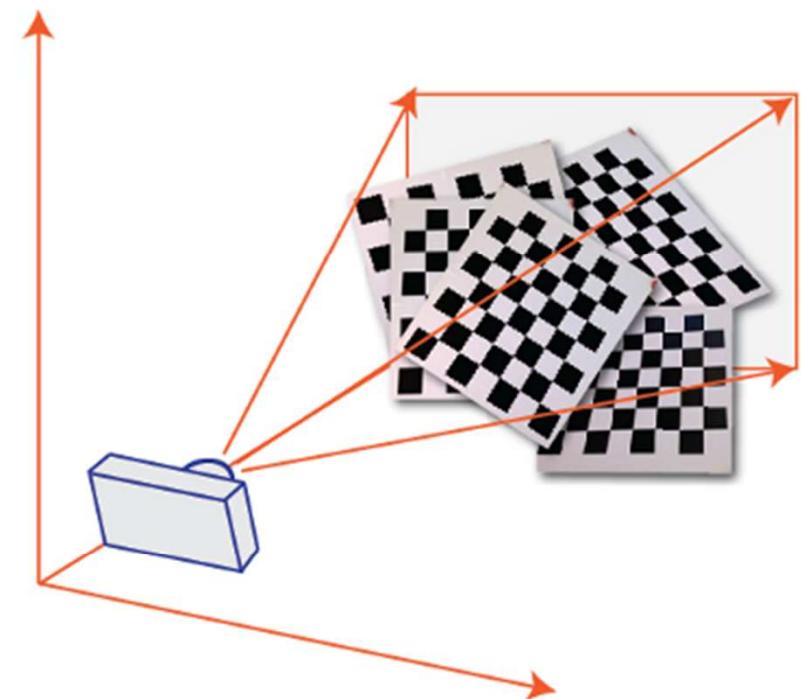
- **Camera intrinsic/extrinsic parameters** (we learnt what are they in prev. classes)
 - *Matched points across multiple viewpoints (same-colored dots in images)*
- , then, we can lift the multi-view observations to the 3D space



Slide from Lecture 3

Slide credit:
Noah Snavely

How to calibrate cameras?



How to calibrate cameras?



Camera Calibration

$$\begin{pmatrix} f_x & 0 & C_x \\ 0 & f_y & C_y \\ 0 & 0 & 1 \end{pmatrix}$$

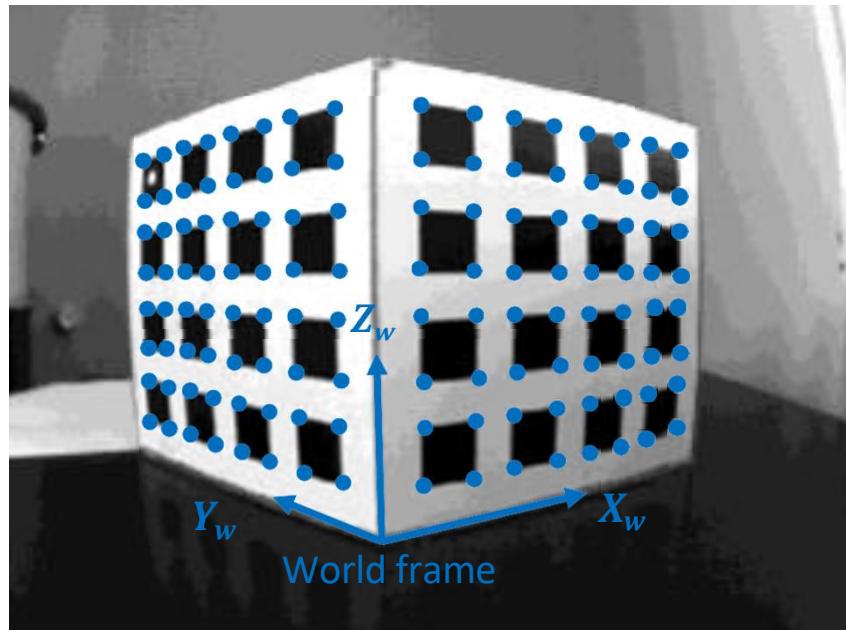
- Calibration is the process to determine the **intrinsic** (K plus lens distortion) **and extrinsic** (R, T) parameters of a camera. For now, we will **neglect the lens distortion** and see later how it can be determined.
- The solution for K, R, T can be found by applying the perspective projection equation:

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R | T] \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

- There are two popular methods:
 - Tsai's method:** uses 3D objects
 - Zhang's method:** uses planar grids

Tsai's Method: Calibration from 3D Objects

- This method was proposed in 1987 by Tsai and consists of measuring the 3D position of $n \geq 6$ control points on a 3D calibration target and the 2D coordinates of their projection in the image.
- Assumption: we know 2D and 3D coordinates of control points
 - 2D: image pre-processing (e.g., corner detectors) *distinctive points, (corner...)*
 - 3D: we know actual size of the 3D object and actual positions of control points as well



Tsai, Roger Y. (1987) "A Versatile Camera Calibration Technique for High-Accuracy 3D Machine Vision Metrology Using Off-the-Shelf TV Cameras and Lenses," *IEEE Journal of Robotics and Automation*, 1987. [PDF](#).

Applying the Direct Linear Transform (DLT) algorithm

The idea of the DLT is to rewrite the perspective projection equation as a **homogeneous linear equation** and solve it by standard methods. Let's write the perspective equation for a generic 3D-2D point correspondence:

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R | T] \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \Rightarrow$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha_u r_{11} + u_0 r_{31} & \alpha_u r_{12} + u_0 r_{32} & \alpha_u r_{13} + u_0 r_{33} & \alpha_u t_1 + u_0 t_3 \\ \alpha_v r_{21} + v_0 r_{31} & \alpha_v r_{22} + v_0 r_{32} & \alpha_v r_{23} + v_0 r_{33} & \alpha_v t_2 + v_0 t_3 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

m

Applying the Direct Linear Transform (DLT) algorithm

The idea of the DLT is to rewrite the perspective projection equation as a **homogeneous linear equation** and solve it by standard methods. Let's write the perspective equation for a generic 3D-2D point correspondence:

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R | T] \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \Rightarrow$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \boxed{\begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix}} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

Applying the Direct Linear Transform (DLT) algorithm

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = M \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \quad \text{goal is getting } M.$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} m_1^T \\ m_2^T \\ m_3^T \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

where m_i^T is the i -th row of M

Applying the Direct Linear Transform (DLT) algorithm

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} m_1^T \\ m_2^T \\ m_3^T \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \xrightarrow{\hspace{1cm}} P$$

Conversion back from homogeneous coordinates to pixel coordinates leads to:

$$\begin{aligned} u &= \frac{\lambda u}{\lambda} = \frac{m_1^T \cdot P}{m_3^T \cdot P} & \Rightarrow & (m_1^T - u_i m_3^T) \cdot P_i = 0 \\ v &= \frac{\lambda v}{\lambda} = \frac{m_2^T \cdot P}{m_3^T \cdot P} & \Rightarrow & (m_2^T - v_i m_3^T) \cdot P_i = 0 \end{aligned}$$

Applying the Direct Linear Transform (DLT) algorithm

- By re-arranging the terms, we obtain

$$\begin{aligned}(m_1^T - u_i m_3^T) \cdot P_i &= 0 \\ (m_2^T - v_i m_3^T) \cdot P_i &= 0\end{aligned} \Rightarrow \begin{pmatrix} P_1^T & 0^T & -u_1 P_1^T \\ 0^T & P_1^T & -v_1 P_1^T \end{pmatrix} \begin{pmatrix} m_1 \\ m_2 \\ m_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

- For n points, we can stack all these equations into a big matrix:

$$\begin{pmatrix} P_1^T & 0^T & -u_1 P_1^T \\ 0^T & P_1^T & -v_1 P_1^T \\ \vdots & \vdots & \vdots \\ P_n^T & 0^T & -u_n P_n^T \\ 0^T & P_n^T & -v_n P_n^T \end{pmatrix} \begin{pmatrix} m_1 \\ m_2 \\ m_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

Applying the Direct Linear Transform (DLT) algorithm

/2

$$2n \left(\begin{array}{cccccccccc} X_w^1 & Y_w^1 & Z_w^1 & 1 & 0 & 0 & 0 & -u_1 X_w^1 & -u_1 Y_w^1 & -u_1 Z_w^1 & -u_1 \\ 0 & 0 & 0 & 0 & X_w^1 & Y_w^1 & Z_w^1 & 1 & -v_1 X_w^1 & -v_1 Y_w^1 & -v_1 Z_w^1 & -v_1 \\ X_w^n & Y_w^n & Z_w^n & 1 & 0 & 0 & 0 & -u_n X_w^n & -u_n Y_w^n & -u_n Z_w^n & -u_n \\ 0 & 0 & 0 & 0 & X_w^n & Y_w^n & Z_w^n & 1 & -v_n X_w^n & -v_n Y_w^n & -v_n Z_w^n & -v_n \end{array} \right) = \begin{pmatrix} m_{11} \\ m_{12} \\ m_{13} \\ m_{14} \\ m_{21} \\ m_{22} \\ m_{23} \\ m_{24} \\ m_{31} \\ m_{32} \\ m_{33} \\ m_{34} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Rightarrow Q \cdot M = 0$$

Q (this matrix is known)

$U, V : 2D$ processing $\frac{52}{50}$.

$P : \frac{45}{50}$.

Applying the Direct Linear Transform (DLT) algorithm

$$Q \cdot M = 0$$

Minimal solution

- $Q_{(2n \times 12)}$ should have rank 11 to have a unique (up to a scale) *non-zero* solution M
- Because each 3D-to-2D point correspondence provides 2 independent equations, then $5 + \frac{1}{2}$ point correspondences are needed (in practice **6 point** correspondences!) *6개를 알면 세울 수 있음*

개를 줄이고.

Over-determined solution

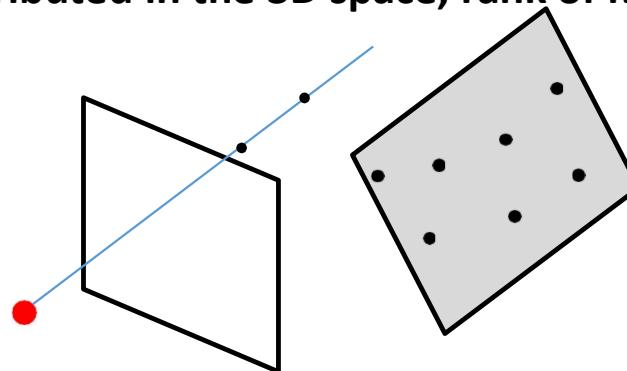
- For $n \geq 6$ points, a solution is the **Least Square solution**, which minimizes the sum of squared residuals, $\|QM\|^2$, subject to the constraint $\|M\|^2 = 1$. It can be solved through Singular Value Decomposition (SVD). The solution is the eigenvector corresponding to the smallest eigenvalue of the matrix $Q^T Q$ (because it is the unit vector x that minimizes $\|Qx\|^2 = x^T Q^T Qx$).
- Matlab instructions:
 - $[U, S, V] = \text{SVD}(Q);$
 - $M = V(:, 12);$*practice 때기 6개보다 많은 points 사용.*

Applying the Direct Linear Transform (DLT) algorithm

Degenerate configurations
(cases where DLT does not work)

$$Q \cdot M = 0$$

1. Points lying on a plane and/or along a single line passing through the center of projection
 - If points are not sufficiently distributed in the 3D space, rank of M is not enough to get solutions



2. Camera and points on a twisted cubic (i.e., smooth curve in 3D space of degree 3)
 - Points on a twisted cubic can be represented with polynomial equation, which can strongly limit rank of M



Applying the Direct Linear Transform (DLT) algorithm

- Once we have determined M , we can recover the intrinsic and extrinsic parameters by remembering that:

$$M = K(R \mid T)$$

$$\begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix}$$

Applying the Direct Linear Transform (DLT) algorithm

- Once we have determined M , we can recover the intrinsic and extrinsic parameters by remembering that:

$$M = K(R \mid T)$$

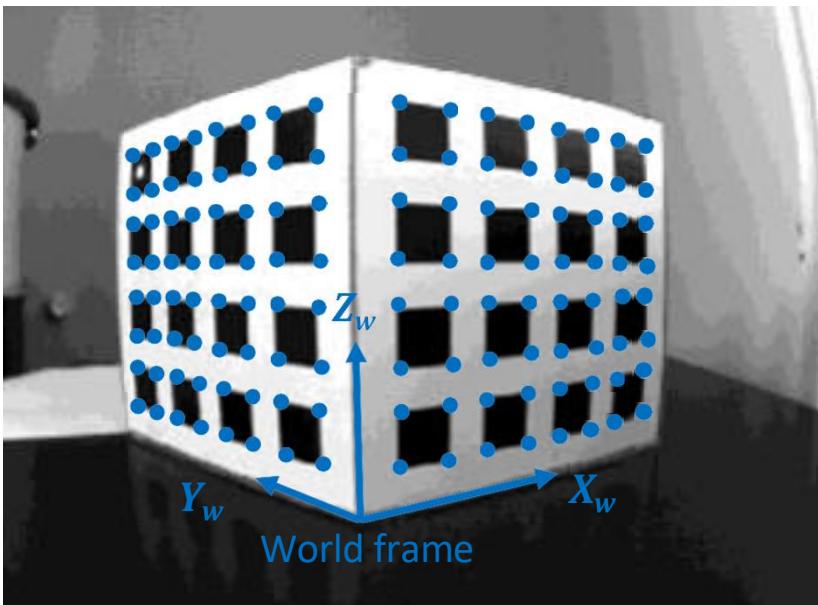
$k \perp (R) = I$
 $R R^T = I$

$$\begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} = \begin{bmatrix} \alpha r_{11} + u_0 r_{31} & \alpha r_{12} + u_0 r_{32} & \alpha r_{13} + u_0 r_{33} & \alpha t_1 + u_0 t_3 \\ \alpha r_{21} + v_0 r_{31} & \alpha r_{22} + v_0 r_{32} & \alpha r_{23} + v_0 r_{33} & \alpha t_2 + v_0 t_3 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix}$$

- However, notice that we are not enforcing the constraint that R is orthogonal, i.e., $R \cdot R^T = I$
- To do this, we can use the so-called **QR factorization of M** , which decomposes M into a R (orthogonal), T , and an upper triangular matrix (i.e., K)

Example of Tsai's Calibration Results

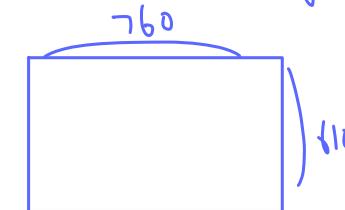
Recommendation: use many **more than 6 points** (ideally more than 20) and **non coplanar**



Corners can be detected with accuracy < 0.1 pixels

α_u	α_u/α_v	K_{12}	u_0	v_0	Average Reprojection error
1673.3	1.0063	1.39	379.96	305.78	0.365

$$\begin{pmatrix} \alpha_u & K_{12} & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{pmatrix}$$

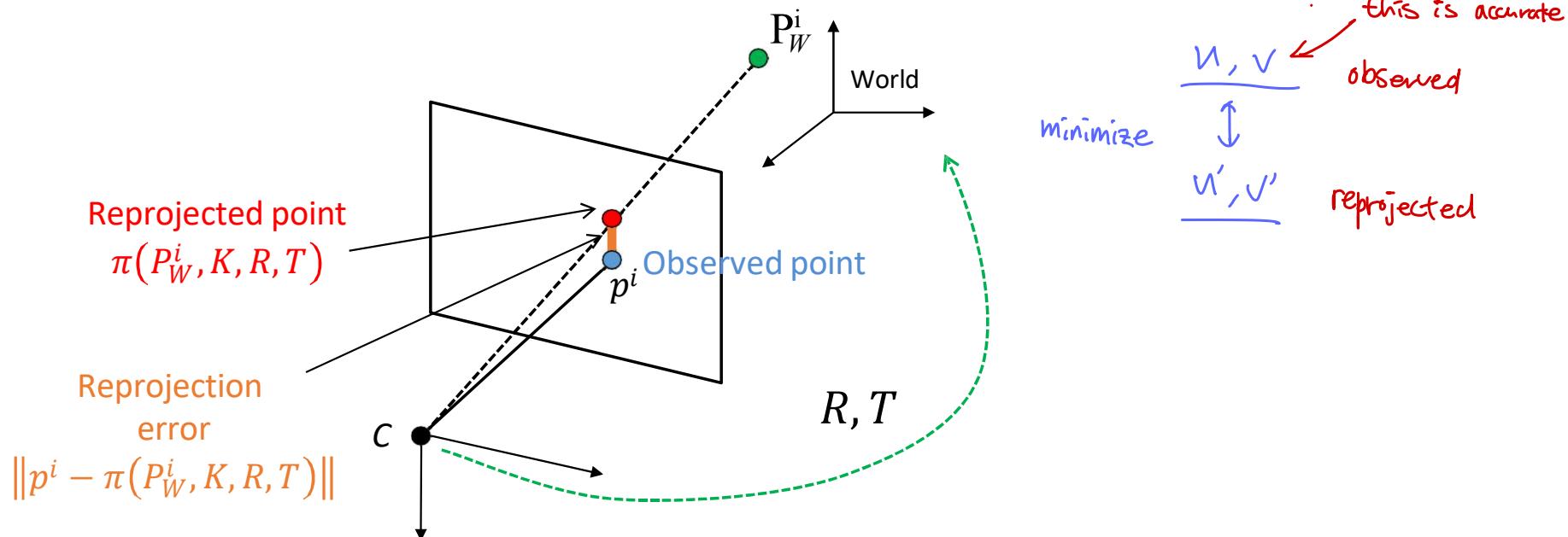


설정은 0에 가까워야 함. 너무 크다.

How can we estimate the lens distortion parameters?
How can we enforce $\alpha_u = \alpha_v$ and $K_{12} = 0$?

Reprojection Error

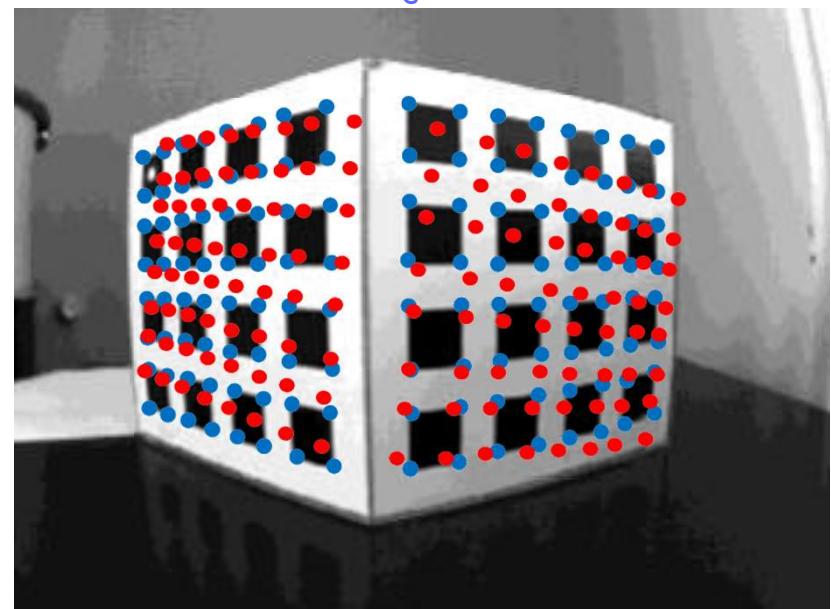
- The reprojection error is the **Euclidean distance** (in pixels) between an **observed image point** and the **corresponding 3D point reprojected** onto the camera frame.
- The reprojection error gives us a **quantitative measure of the accuracy** of the calibration (**ideally it should be zero**).



Reprojection Error

- The reprojection error can be used to assess the quality of the camera calibration
- What are the sources of the reprojection error? -> limitation of DLT noise of observation.

1. Lens distortion
2. Errors in corner detections
3. Least square optimization is prone to noise



- Control points
(observed points)
- Reprojected points
 $\pi(P_W^i, K, R, T)$

Non-Linear Calibration Refinement

Tsai 3 calibration ကြတ်လဲ ရှိသူ minimize,

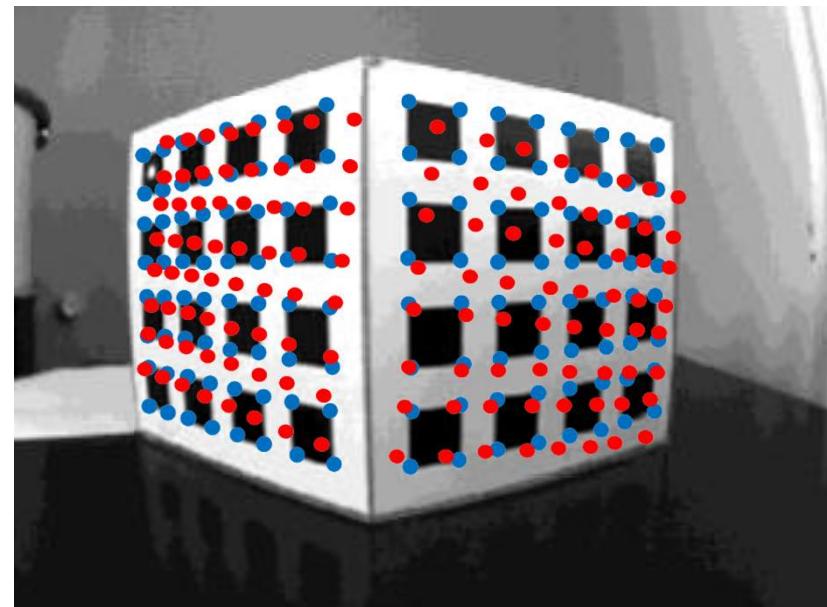
- The calibration parameters K, R, T determined by the DLT can be refined by minimizing the following cost:

$$K, R, T, \text{lens distortion} = \underset{argmin_{K,R,T,\text{lens}}}{\sum_{i=1}^n} \|p^i - \pi(P_W^i, K, R, T)\|^2 \quad \text{minimize}$$

- This time we also include the **lens distortion** (can be set to 0 for initialization)
- Can be minimized using **Levenberg–Marquardt** (more robust than Gauss–Newton to local minima)

Hessian

Jacobian (SGD)



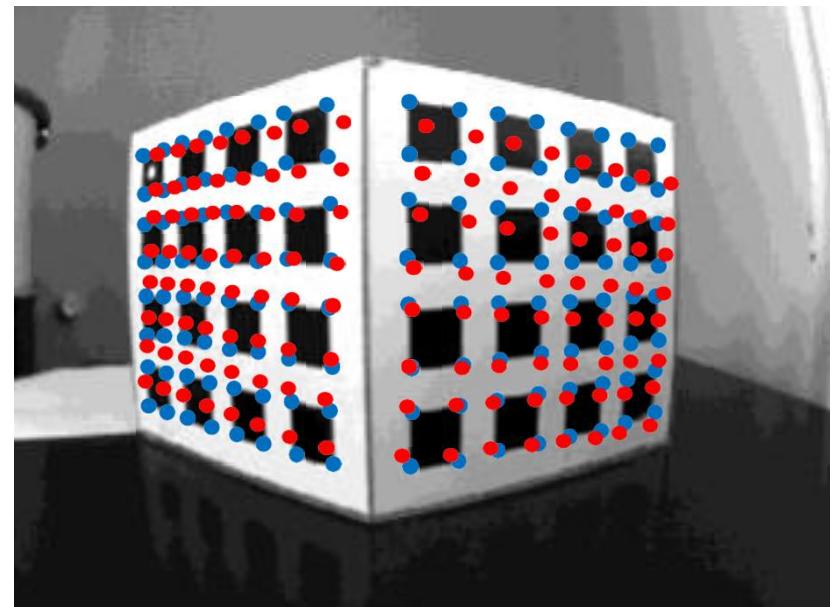
- Control points (observed points)
- Reprojected points $\pi(P_W^i, K, R, T)$

Non-Linear Calibration Refinement

- The calibration parameters K, R, T determined by the DLT can be refined by minimizing the following cost:

$$K, R, T, \text{lens distortion} = \\ \operatorname{argmin}_{K, R, T, \text{lens}} \sum_{i=1}^n \|p^i - \pi(P_W^i, K, R, T)\|^2$$

- This time we also include the **lens distortion** (can be set to 0 for initialization)
- Can be minimized using **Levenberg–Marquardt** (more robust than Gauss-Newton to local minima)



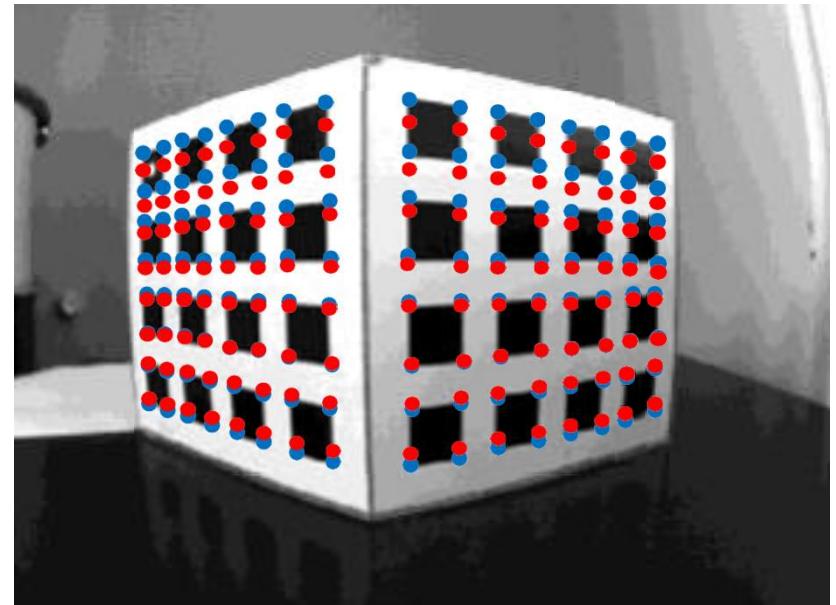
- Control points (observed points)
- Reprojected points $\pi(P_W^i, K, R, T)$

Non-Linear Calibration Refinement

- The calibration parameters K, R, T determined by the DLT can be refined by minimizing the following cost:

$$K, R, T, \text{lens distortion} = \\ \operatorname{argmin}_{K, R, T, \text{lens}} \sum_{i=1}^n \|p^i - \pi(P_W^i, K, R, T)\|^2$$

- This time we also include the **lens distortion** (can be set to 0 for initialization)
- Can be minimized using **Levenberg–Marquardt** (more robust than Gauss-Newton to local minima)



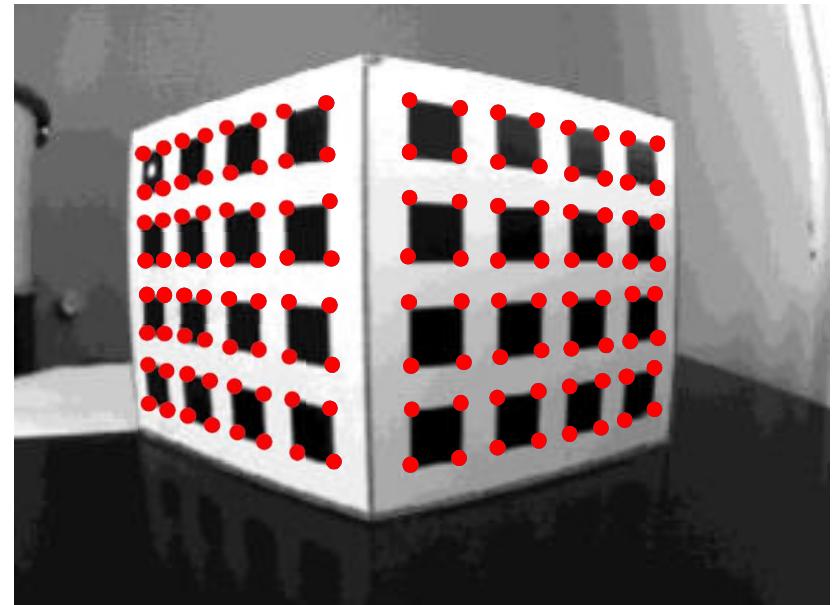
- Control points
(observed points)
- Reprojected points
 $\pi(P_W^i, K, R, T)$

Non-Linear Calibration Refinement

- The calibration parameters K, R, T determined by the DLT can be refined by minimizing the following cost:

$$K, R, T, \text{lens distortion} = \\ \operatorname{argmin}_{K, R, T, \text{lens}} \sum_{i=1}^n \|p^i - \pi(P_W^i, K, R, T)\|^2$$

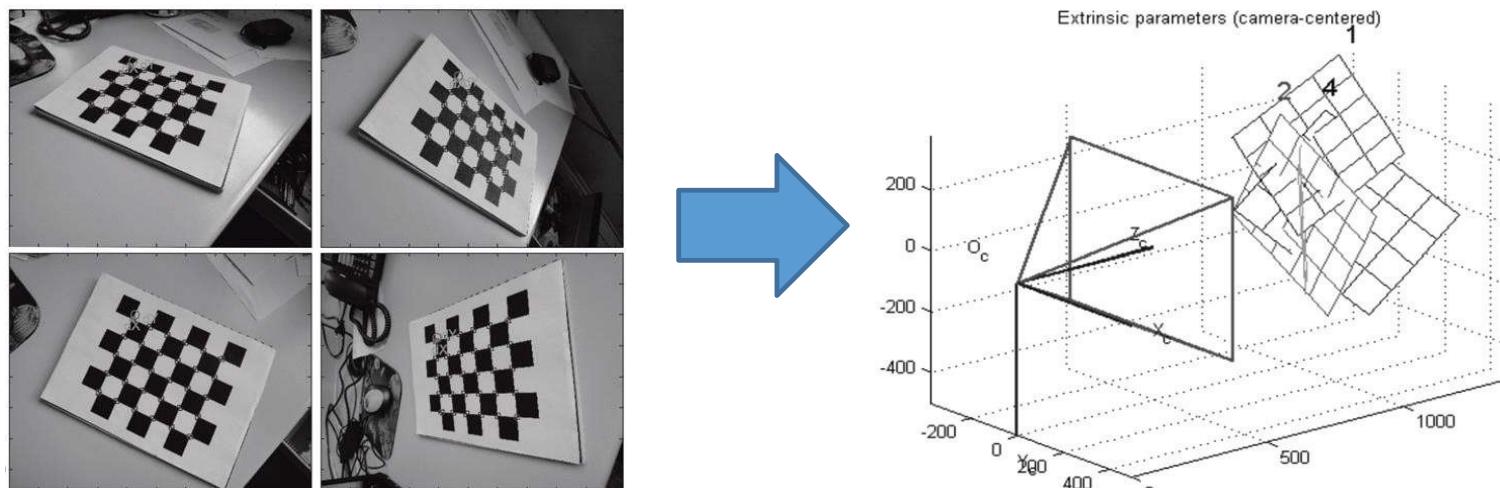
- This time we also include the **lens distortion** (can be set to 0 for initialization)
- Can be minimized using **Levenberg–Marquardt** (more robust than Gauss-Newton to local minima)



- Control points
(observed points)
- Reprojected points
 $\pi(P_W^i, K, R, T)$

Zhang's Algorithm: Calibration from Planar Grids

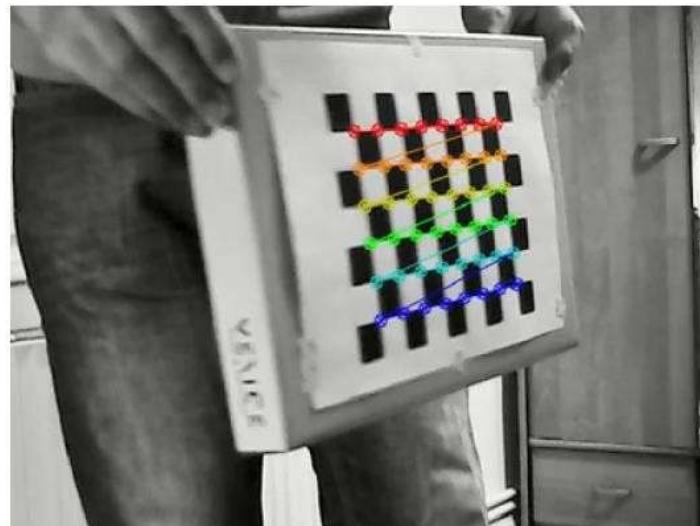
- **Tsai's calibration** requires that the world's 3D points are non-coplanar, which is **not very practical**
- **Today's camera calibration toolboxes** ([Matlab](#), [OpenCV](#)) use **multiple views of a planar grid** (e.g., a checker board)
- They are based on a method developed in 2000 by Zhang (Microsoft Research)



Zhang, A flexible new technique for camera calibration, IEEE Transactions on Pattern Analysis and Machine Intelligence, 2000. [PDF](#).

Zhang's Algorithm: Calibration from Planar Grids

- **Tsai's calibration** requires that the world's 3D points are non-coplanar, which is **not very practical**
- **Today's camera calibration toolboxes** ([Matlab](#), [OpenCV](#)) use **multiple views** of a **planar grid** (e.g., a checker board)
- They are based on a method developed in 2000 by Zhang (Microsoft Research)



Zhang, A flexible new technique for camera calibration, IEEE Transactions on Pattern Analysis and Machine Intelligence, 2000. [PDF](#).

Applying the Direct Linear Transform (DLT) algorithm

As in Tsai's method, we start by writing the perspective projection equation (again, we neglect the radial distortion). However, in **Zhang's method the points are all coplanar**, i.e., $Z_w = 0$, and thus we can write:

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R | T] \cdot \begin{bmatrix} X_w \\ Y_w \\ 0 \\ 1 \end{bmatrix} \Rightarrow$$

in World coord system

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ 0 \\ 1 \end{bmatrix}$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & t_1 \\ r_{21} & r_{22} & t_2 \\ r_{31} & r_{32} & t_3 \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ 1 \end{bmatrix}$$

Applying the Direct Linear Transform (DLT) algorithm

As in Tsai's method, we start by writing the perspective projection equation (again, we neglect the radial distortion). However, in **Zhang's method the points are all coplanar**, i.e., $\mathbf{Z}_w = \mathbf{0}$, and thus we can write:

$$\begin{aligned} \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} &= K[R | T] \cdot \begin{bmatrix} X_w \\ Y_w \\ 0 \\ 1 \end{bmatrix} \Rightarrow \\ \Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} &= \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ 0 \\ 1 \end{bmatrix} \\ \Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} &= \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ 1 \end{bmatrix} \end{aligned}$$

Applying the Direct Linear Transform (DLT) algorithm

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ 1 \end{bmatrix}$$

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = H \cdot \begin{bmatrix} X_w \\ Y_w \\ 1 \end{bmatrix}$$

This matrix is called
Homography

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} h_1^T \\ h_2^T \\ h_3^T \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ 1 \end{bmatrix}$$

where h_i^T is the *i-th* row of H

Applying the Direct Linear Transform (DLT) algorithm

$$\Rightarrow \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} h_1^T \\ h_2^T \\ h_3^T \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ 1 \end{bmatrix} \xrightarrow{P}$$

Conversion back from homogeneous coordinates to pixel coordinates leads to:

$$\begin{aligned} u &= \frac{\lambda u}{\lambda} = \frac{h_1^T \cdot P}{h_3^T \cdot P} & \Rightarrow (h_1^T - u h_3^T) \cdot P_i = 0 \\ v &= \frac{\lambda v}{\lambda} = \frac{h_2^T \cdot P}{h_3^T \cdot P} & \Rightarrow (h_2^T - v h_3^T) \cdot P_i = 0 \end{aligned}$$

Applying the Direct Linear Transform (DLT) algorithm

- By re-arranging the terms, we obtain:

$$\begin{aligned}(h_1^T - u_i h_3^T) \cdot P_i &= 0 &\Rightarrow P_i^T \cdot h_1 + 0 \cdot h_2^T - u_i P_i^T \cdot h_3^T &= 0 &\Rightarrow \begin{pmatrix} P_i^T & 0^T & -u_i P_i^T \\ 0^T & P_i^T & -v_i P_i^T \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\(h_2^T - v_i h_3^T) \cdot P_i &= 0 &\Rightarrow 0 \cdot h_1^T + P_i^T \cdot h_2^T - v_i P_i^T \cdot h_3^T &= 0\end{aligned}$$

- For n points (from a **single view**), we can stack all these equations into a big matrix:

$$\underbrace{\begin{pmatrix} P_1^T & 0^T & -u_1 P_1^T \\ 0^T & P_1^T & -v_1 P_1^T \\ \dots & \dots & \dots \\ P_n^T & 0^T & -u_n P_n^T \\ 0^T & P_n^T & -v_n P_n^T \end{pmatrix}}_{Q} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix} \Rightarrow Q \cdot H = 0$$

Q (this matrix is **known**) H (this matrix is **unknown**)

Applying the Direct Linear Transform (DLT) algorithm

$$Q \cdot H = 0$$

$M \in \mathbb{R}^{2n \times 12}$
 $H \in \mathbb{R}^{2n \times 9}$

Minimal solution

- $Q_{(2n \times 9)}$ should have rank 8 to have a unique (up to a scale) non-trivial solution H
- Each point correspondence provides 2 independent equations x, y .
- Thus, a minimum of **4 non-collinear points** is required $8/2 = 4$

Over-determined solution

- $n \geq 4$ points
- It can be solved through Singular Value Decomposition (SVD) (same considerations as before)

How to recover K, R, T

- H can be decomposed by recalling that:
- Differently from Tsai's, the decomposition of H into K, R, T requires at least two views *objekt* *zwei*
- In practice **the more views the better**, e.g., 20-50 views spanning the entire field of view of the camera for the best calibration results
- Notice that now each view j has a **different homography H^j** (and so a different R^j and T^j). However, **K is the same for all views**:

$$\begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11} & r_{12} & t_1 \\ r_{21} & r_{22} & t_2 \\ r_{31} & r_{32} & t_3 \end{bmatrix}$$

$$\begin{bmatrix} h_{11}^j & h_{12}^j & h_{13}^j \\ h_{21}^j & h_{22}^j & h_{23}^j \\ h_{31}^j & h_{32}^j & h_{33}^j \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{11}^j & r_{12}^j & t_1^j \\ r_{21}^j & r_{22}^j & t_2^j \\ r_{31}^j & r_{32}^j & t_3^j \end{bmatrix}$$

How to recover K, R, T from H and from multiple views?

Details won't be asked
at the exam ☺

1. Estimate the homography H_i for each i -th view using the DLT algorithm.

2. Determine the intrinsics K of the camera from a set of homographies:

1. Each homography $H_i \sim K(\mathbf{r}_1, \mathbf{r}_2, \mathbf{t})$ provides two linear equations in the 6 entries of the matrix $B := K^{-T}K^{-1}$. Letting $\mathbf{w}_1 := Kr_1$, $\mathbf{w}_2 := Kr_2$, the rotation constraints $\mathbf{r}_1^T \mathbf{r}_1 = \mathbf{r}_2^T \mathbf{r}_2 = 1$ and $\mathbf{r}_1^T \mathbf{r}_2 = 0$ become $\mathbf{w}_1^T B \mathbf{w}_1 - \mathbf{w}_2^T B \mathbf{w}_2 = 0$ and $\mathbf{w}_1^T B \mathbf{w}_2 = 0$.

R 은 orthogonal matrix 이다.

2. Stack $2N$ equations from N views, to yield a linear system $A\mathbf{b} = \mathbf{0}$. Solve for \mathbf{b} (i.e., B) using the Singular Value Decomposition (SVD).
3. Use Cholesky decomposition to obtain K from B .

$$r_1 = k^{-1} h_1, r_2 = k^{-1} h_2$$

$$\begin{aligned} r_1^T r_2 &= (k^{-1} h_1)^T (k^{-1} h_2) = h_1^T k^{-T} k^{-1} h_2 = 0 \\ &= h_1^T B h_2 = 0 \end{aligned}$$

3. The extrinsic parameters for each view can be computed using K :

$\mathbf{r}_1 \sim \lambda K^{-1} H_i(:, 1)$, $\mathbf{r}_2 \sim \lambda K^{-1} H_i(:, 2)$, $\mathbf{r}_3 = \mathbf{r}_1 \times \mathbf{r}_2$ and $T_i = \lambda K^{-1} H_i(:, 3)$, with $\lambda = 1/K^{-1} H_i(:, 1)$. Finally, build $R_i = (\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3)$ and enforce rotation matrix constraints.

$$h_i^T B h_j = \sum_{m=1}^3 \sum_{n=1}^3 h_{im} B_{mn} h_{jn} = v_{ij}^T b$$

$$B_{mn} = B_{nm}$$

$$B = k^{-T} k^{-1}$$

$k \in \text{square up}$ $k^{-T} = (k^{-1})^T$

$$B^T = (k^{-T} k^{-1})^T = k^{-T} k^{-1} = B$$

$B^T = B$ 면 대칭행렬이여.

$$\|h_1\|^2 = \|r_2\|^2$$

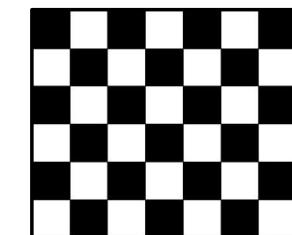
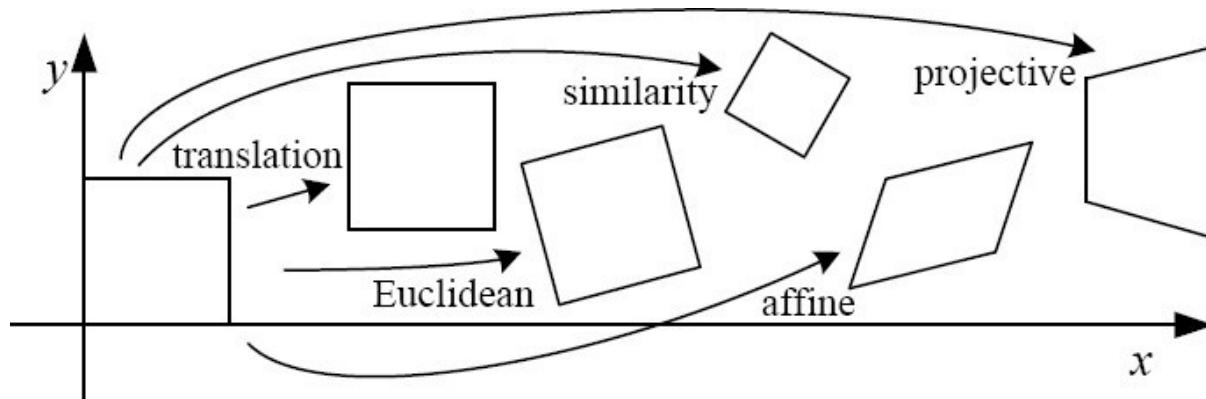
$$\Rightarrow (k^{-1} h_1)^T (k^{-1} h_1) = h_1^T B h_1 = h_2^T B h_2$$

$$\begin{bmatrix} B_{11} & B_{12} & B_{13} \\ B_{21} & B_{22} & B_{23} \\ B_{31} & B_{32} & B_{33} \end{bmatrix}$$

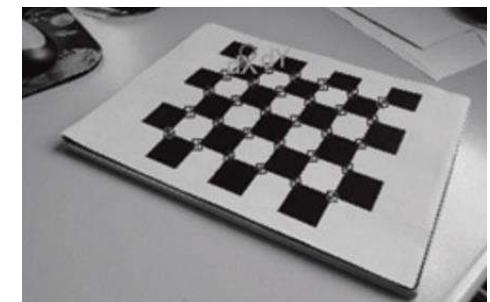
제한 조건

$$b = [B_{11}, B_{12}, B_{13}, B_{21}, B_{22}, B_{23}, B_{31}, B_{32}, B_{33}]^T$$

Types of 2D Transformations (from Lecture 2)



H



Name	Matrix	# D.O.F.	Preserves:	Icon
translation	$[I \mid t]_{2 \times 3}$	2	orientation + ...	
rigid (Euclidean)	$[R \mid t]_{2 \times 3}$	3	lengths + ...	
similarity	$[sR \mid t]_{2 \times 3}$	4	angles + ...	
affine	$[A]_{2 \times 3}$	6	parallelism + ...	
projective	$[\tilde{H}]_{3 \times 3}$	8	straight lines	

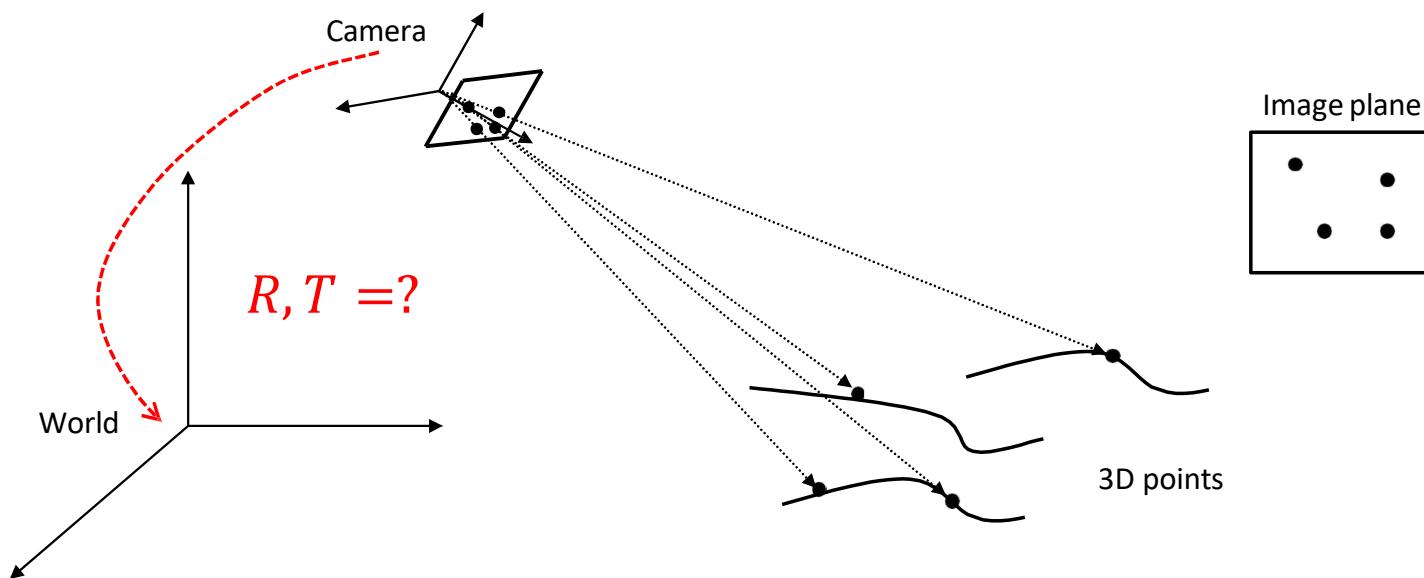
This matrix is called **Homography**

2. Camera localization

Camera Localization (or Perspective from n Points: PnP)

Degree of Freedom

- This is the problem of determining the **6DoF pose of a camera** (position and orientation) with respect to the world frame **from a set of 3D-2D point correspondences**.
- It assumes the **camera** to be **already calibrated**
- **In other words, the goal is getting extrinsics (R and T) while intrinsics are given**
- The **DLT can be used** to solve this problem **but is suboptimal**. We want to study **algebraic solutions** to the problem.



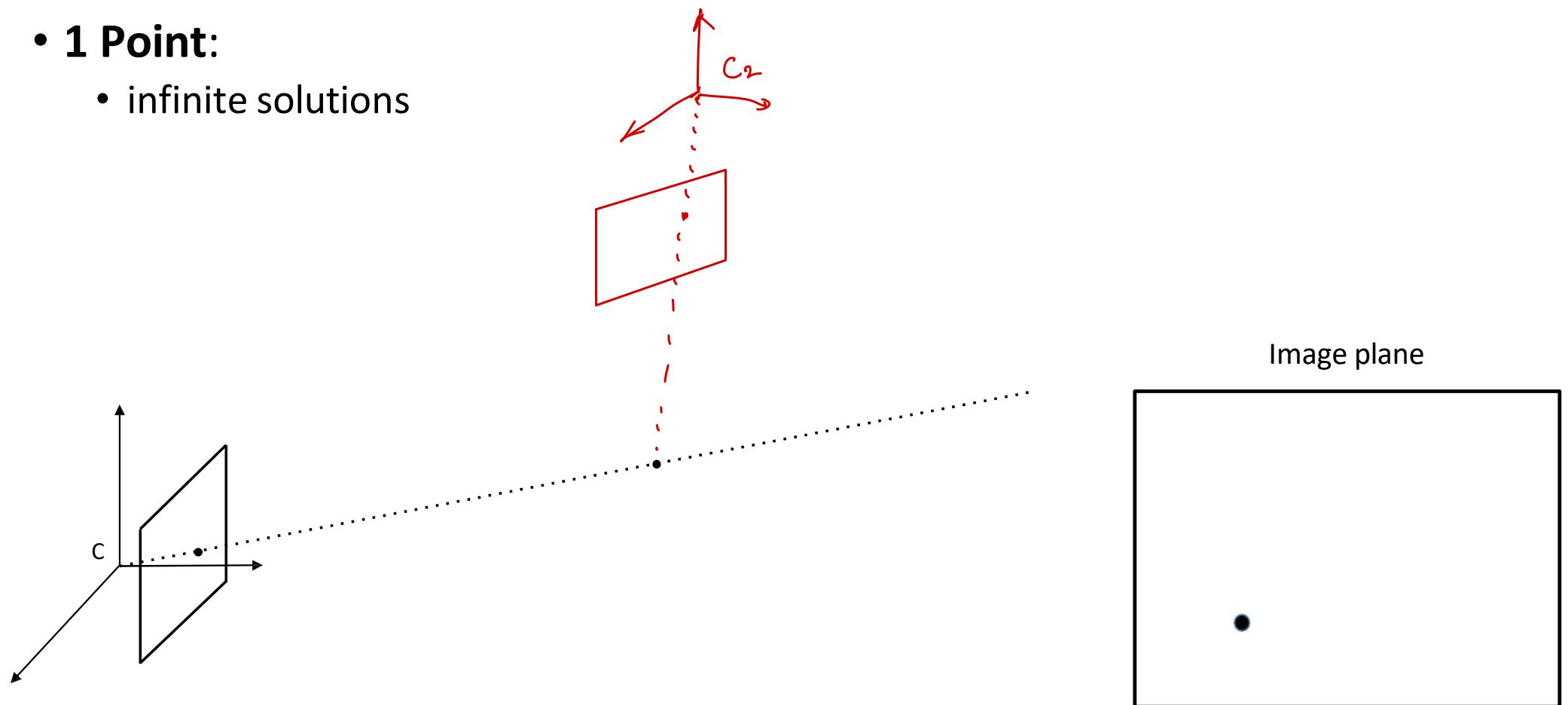
How Many Points are Enough?

- Here, ‘points’ means correspondence between 3D and 2D of points
 - E.g., ‘2 points’ means we know 1) 3D and 2D positions of point A and 2) 3D and 2D positions of point B
-
- **1 Point:**
 - infinite solutions
 - **2 Points:**
 - infinitely many solutions, but bounded
 - **3 Points (non collinear):**
 - up to 4 solution
 - **4 Points:**
 - Unique solution

Zhang

1 Point

- **1 Point:**
 - infinite solutions



2 Points

- **2 Points:**
 - infinite solutions, but bounded

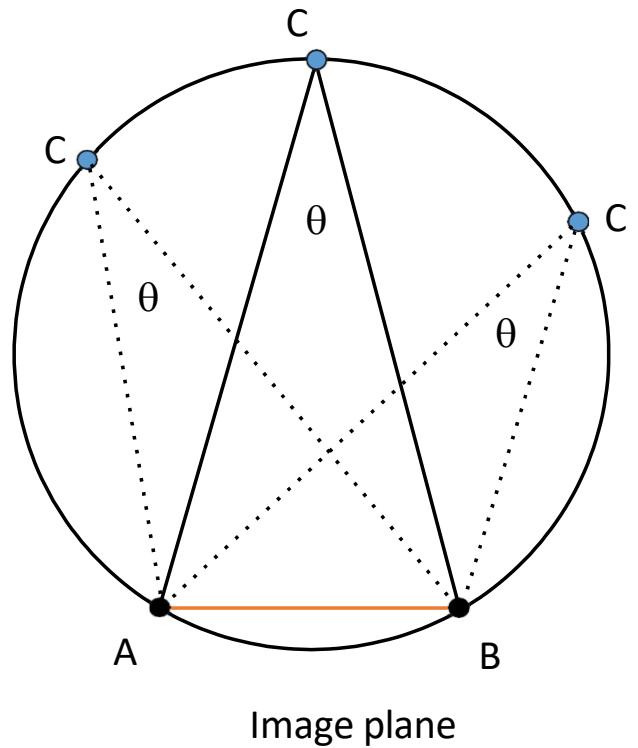
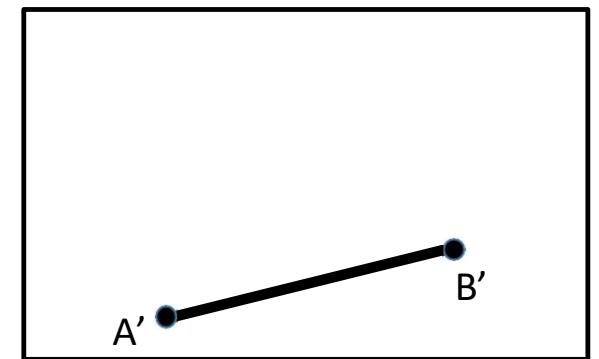
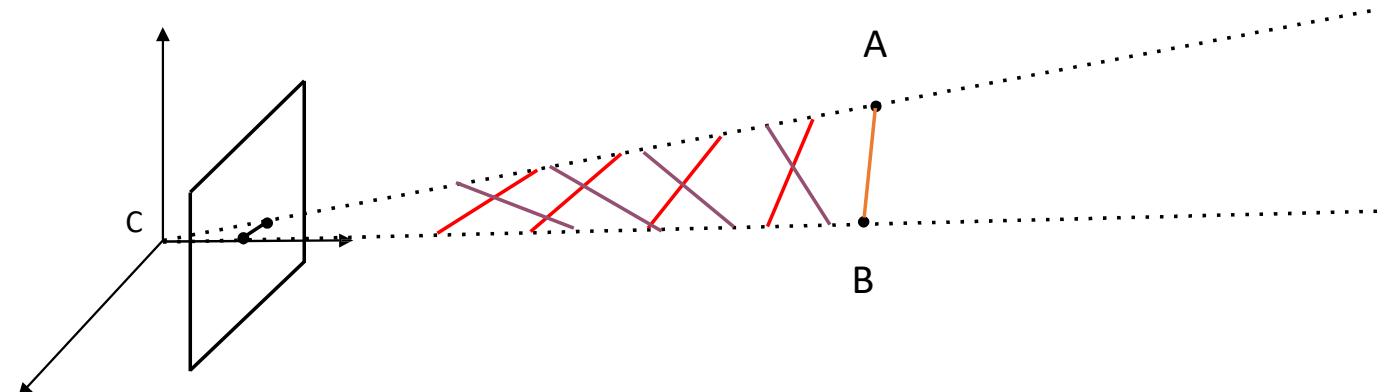


Image plane



3 Points (P3P problem)

From Carnot's Theorem:

- **3 Points (non collinear):**

- up to 4 solution
- L_A, L_B , and L_C are unknown variables

$$|AB|^2 = s_1^2 = L_B^2 + L_A^2 - 2L_B L_A \cos \theta_{AB}$$

$$|AC|^2 = s_2^2 = L_A^2 + L_C^2 - 2L_A L_C \cos \theta_{AC}$$

$$|BC|^2 = s_3^2 = L_B^2 + L_C^2 - 2L_B L_C \cos \theta_{BC}$$

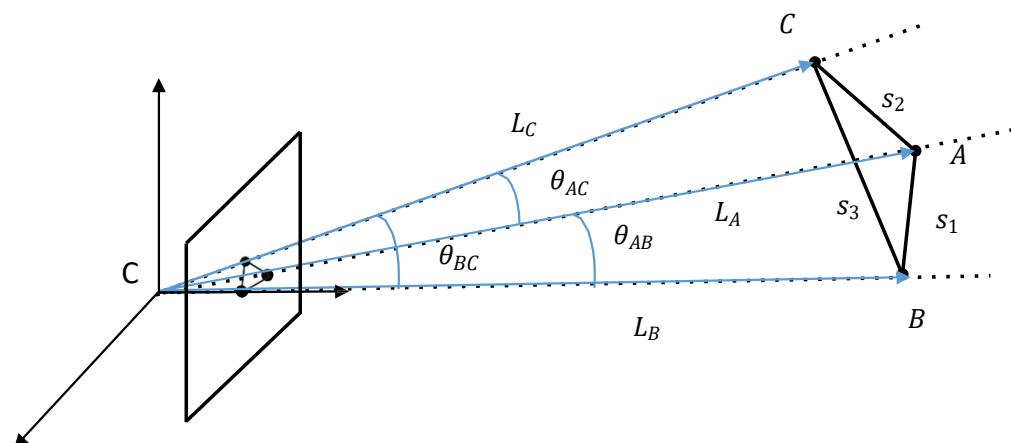
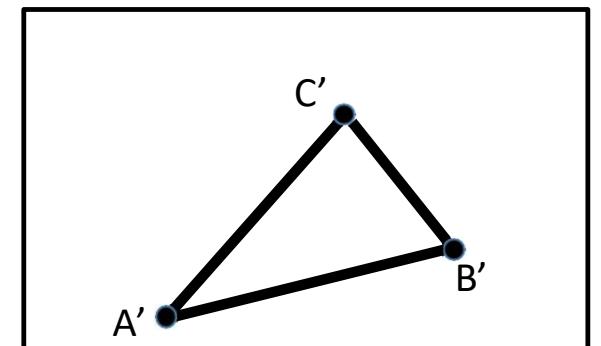


Image plane



Algebraic Approach: reduce to 4th order equation

$$||AB||^2 = s_1^2 = L_B^2 + L_A^2 - 2L_B L_A \cos \theta_{AB}$$

$$||AC||^2 = s_2^2 = L_A^2 + L_C^2 - 2L_A L_C \cos \theta_{AC}$$

$$||BC||^2 = s_3^2 = L_B^2 + L_C^2 - 2L_B L_C \cos \theta_{BC}$$

- It is known that **n independent polynomial equations**, in **n unknowns**, can have no more **solutions** than the **product of their respective degrees?** Thus, the system can have a maximum of **($2 \times 2 \times 2 = 8$)** solutions. However, because every term in the system is either a constant or of **second degree**, for every real positive solution **there is a negative solution.** $L_A \ L_B \ L_C$; positive.
- Thus, with 3 points, there are at most **($2 \times 2 \times 2 / 2 = 4$) valid (positive) solutions.**

Algebraic Approach: reduce to 4th order equation

$$s_1^2 = L_B^2 + L_A^2 - 2L_B L_A \cos \theta_{AB}$$

$$s_2^2 = L_A^2 + L_C^2 - 2L_A L_C \cos \theta_{AC}$$

$$s_3^2 = L_B^2 + L_C^2 - 2L_B L_C \cos \theta_{BC}$$

- By defining $x = L_B/L_A$, it can be shown that the system can be reduced to a 4th order equation:

$$G_0 + G_1 x + G_2 x^2 + G_3 x^3 + G_4 x^4 = 0$$

How can we disambiguate the 4 solutions? How do we determine R and T ?

- A 4th point can be used to disambiguate the solutions. A classification of the four solutions and the determination of R and T from the point distances was given Gao's algorithm, implemented in OpenCV ([solvePnP P3P](#))
- In real worlds, we have more than three points, so use one of them to get only one (R,T) pair out of four pairs.

Gao, Hou, Tang, Cheng. Complete Solution Classification for the Perspective-Three-Point Problem.
IEEE Transactions on Pattern Analysis and Machine Intelligence, 2003. [PDF](#).

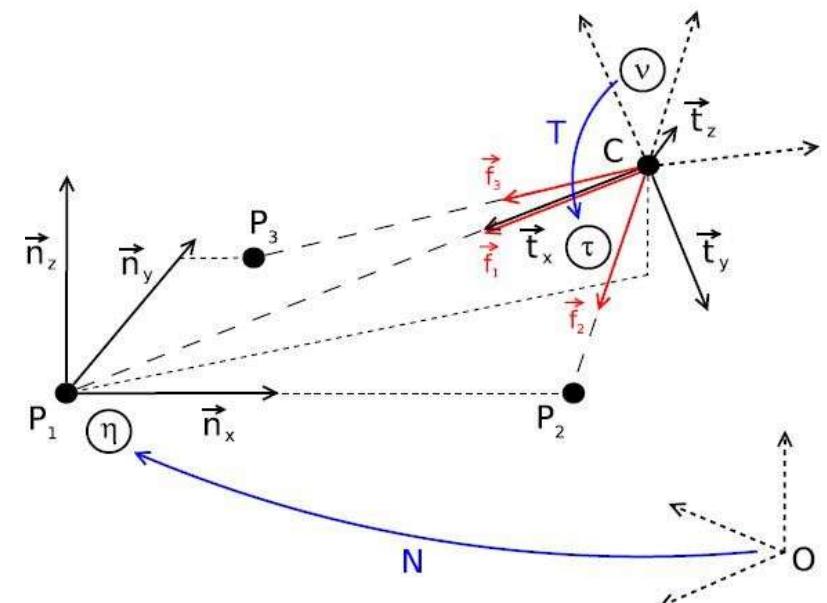
Modern Solution to P3P

Details won't be asked
at the exam ☺

Modern은 다른 strategy
사용된다~.

A more **modern version of P3P** was developed by Kneip in 2011 and **directly solves for the camera's pose** (not distances from the points). This solution inspired the algorithm currently used in OpenCV ([solvePnP AP3P](#)), by Ke'17, which consists of two steps:

1. Eliminate the camera's position and the features' distances to yield a system of 3 equations *in the camera's orientation alone*.
 2. Successively eliminate two of the unknown 3-DOFs (angles) algebraically and arrive at a **quartic polynomial equation**.
- Outperforms previous methods in terms of speed, accuracy, and robustness to close-to-singular cases.



Kneip, Scaramuzza, Siegwart. A Novel Parameterization of the Perspective-Three-Point Problem for a Direct Computation of Absolute Camera Position and Orientation. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2011. [PDF](#).

Ke, Roumeliotis. An Efficient Algebraic Solution to the Perspective-Three-Point Problem. CVPR'17. [PDF](#).

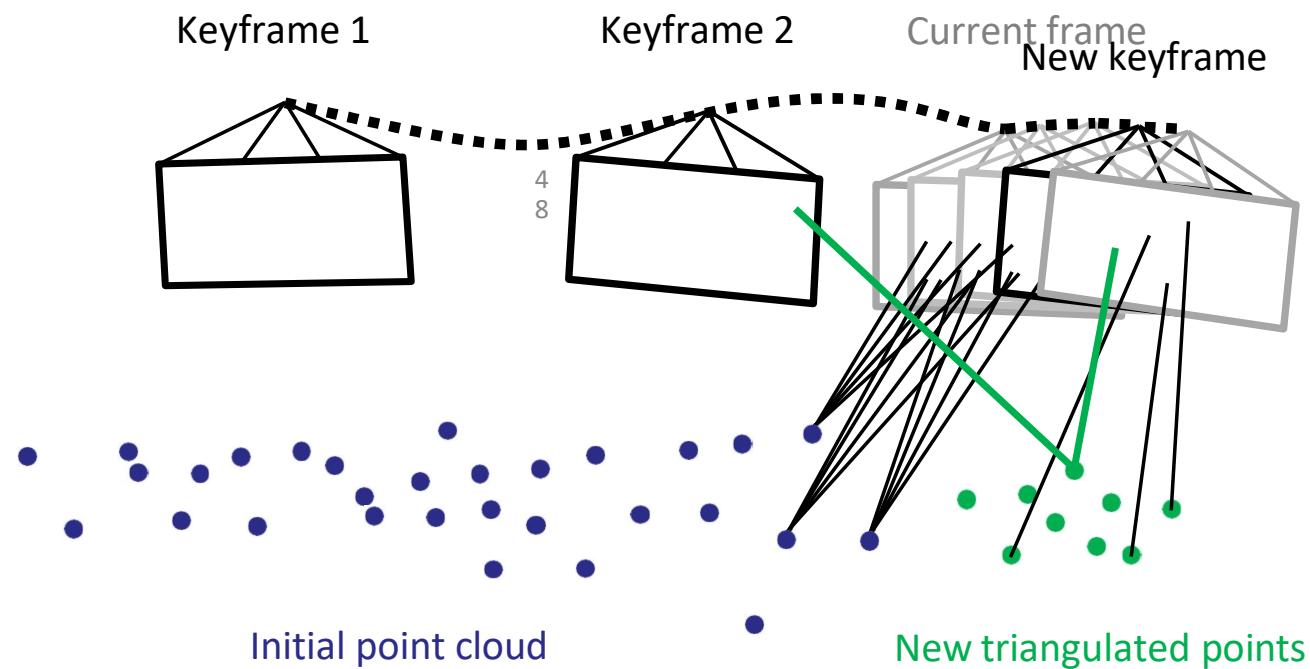
Solution to PnP for $n \geq 4$

Details won't be asked
at the exam ☺

An efficient algebraic solution to the PnP problem for $n \geq 4$ was developed by Lepetit in 2009 and coined **EPnP** (Efficient PnP) and can be found in OpenCV ([solvePnP EPnP](#))

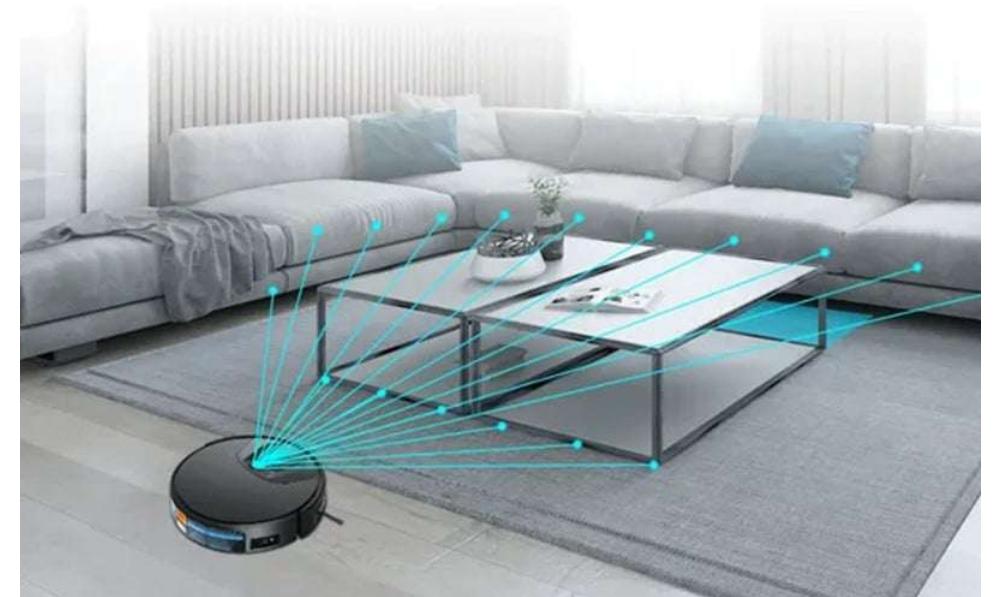
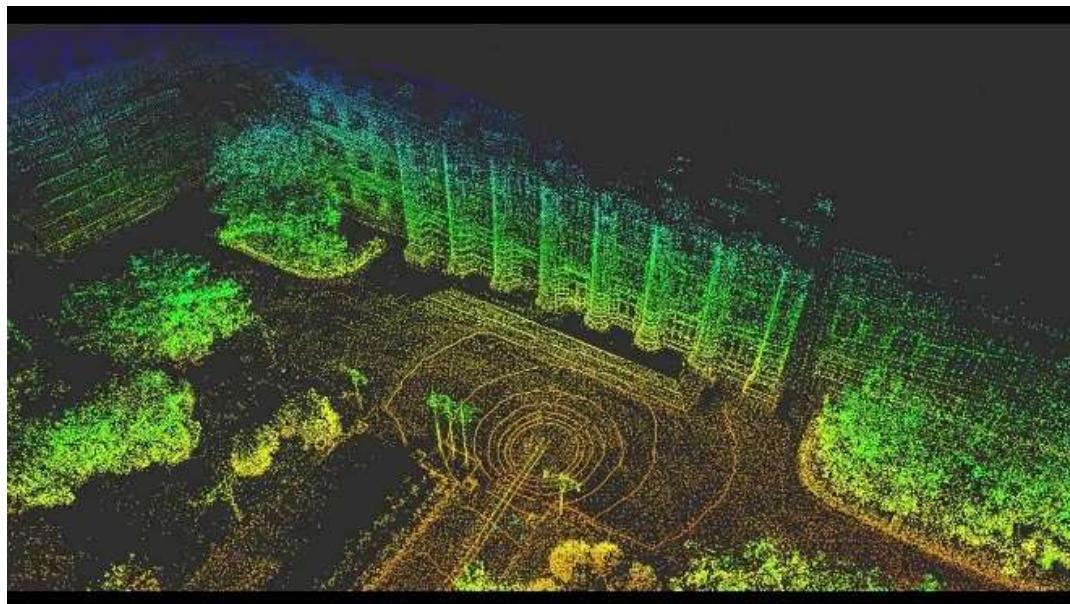
- EPnP expresses the n world's points as a weighted sum of **four virtual control points**
- The coordinates of these virtual control points become the **unknowns of the problem**, which can be solved in $O(n)$ time by solving a **constant number of quartic polynomial equations**
- The final pose of the camera is then solved from the control points

Application to Monocular Visual Odometry



SLAM

(Simultaneous Localization and Mapping)



Robust Estimation in Presence of Outliers

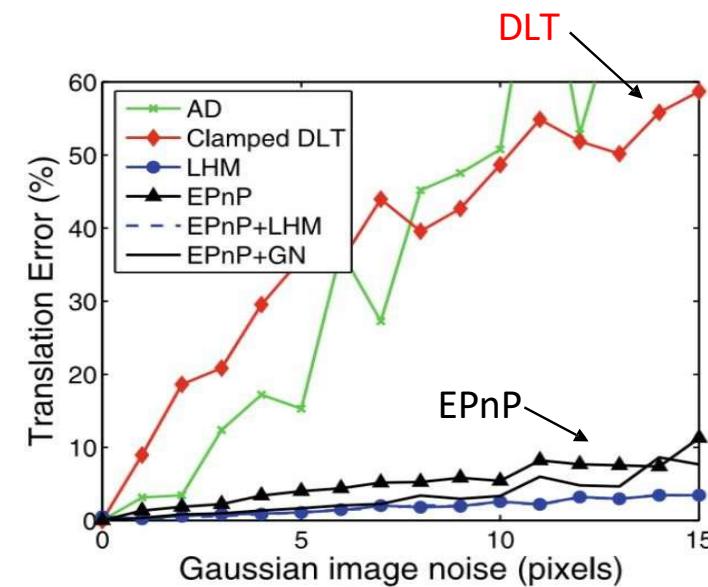
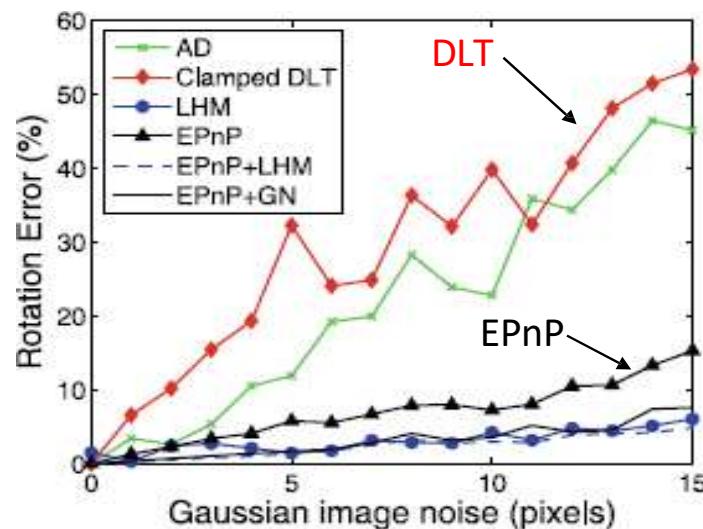
- **All PnP problems** (solved by DLT, EPnP, or P3P algorithms) are prone to errors if there are outliers in the set of 3D-2D point correspondences.
- The **RANSAC** algorithm can be used, in conjunction with the PnP algorithm, to **remove the outliers**.
- PnP with RANSAC can be found in OpenCV's ([*solvePnPRansac*](#))

EPnP vs. DLT

If a camera is calibrated, only R and T need to be determined. In this case, should we use DLT or EPnP?

EPnP vs. DLT: Accuracy vs noise

EPnP is more up to **10 × more accurate and more efficient** than DLT

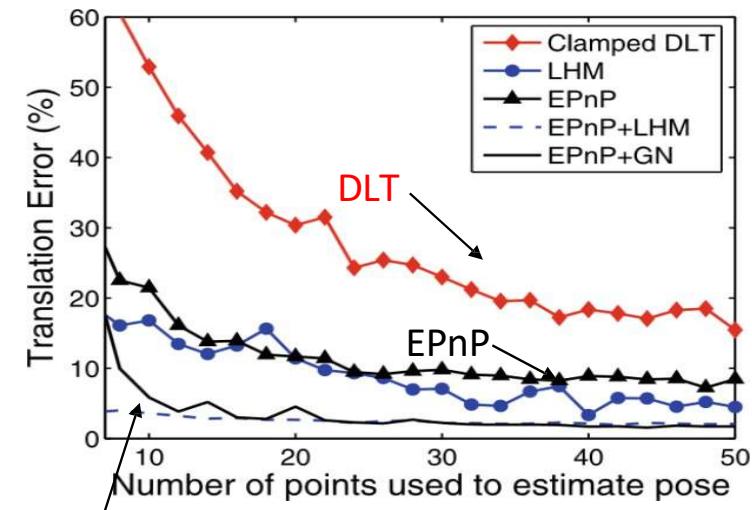
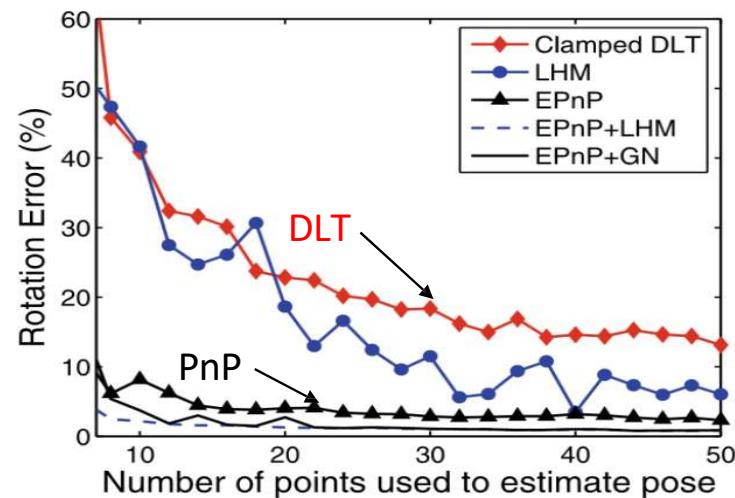


Plots from

Lepetit, Moreno Noguer, Fua, EPnP: An Accurate $O(n)$ Solution to the PnP Problem, International Journal of Computer Vision. [PDF](#).

EPnP vs. DLT: Accuracy vs number of points

EPnP is more up to **10 × more accurate and more efficient** than DLT



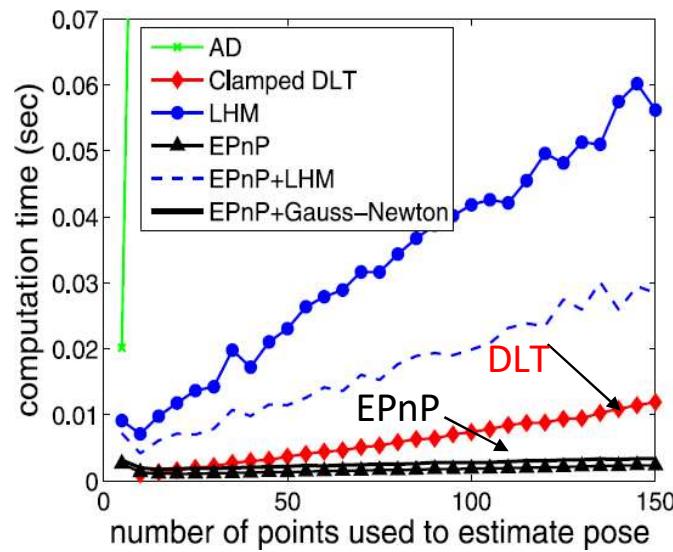
EPnP+ Gauss Newton

Plots from

Lepetit, Moreno Noguer, Fua, EPnP: An Accurate $O(n)$ Solution to the PnP Problem, International Journal of Computer Vision. [PDF](#).

EPnP vs. DLT: Timing

EPnP is more up to **10 × more accurate and more efficient** than DLT



Plots from

Lepetit, Moreno Noguer, Fua, EPnP: An Accurate $O(n)$ Solution to the PnP Problem, International Journal of Computer Vision. [PDF](#).

PnP problem: Recap

Calibrated camera
(i.e., intrinsic parameters are known)

Uncalibrated camera
(i.e., intrinsic parameters unknown)

Either DLT or EPnP can be used

Only DLT can be used

EPnP: minimum number of points: **3 (P3P) +1** for disambiguation

DLT: Minimum number of points: **4 if coplanar, 6 if non-coplanar**

The output of both DLT and EPnP can be refined via **non-linear optimization**
by minimizing the sum of squared reprojection errors

Real world cases

- Small capture setup with few Kinect cameras
- I and my colleagues made this to capture simple data quickly when I was in Meta



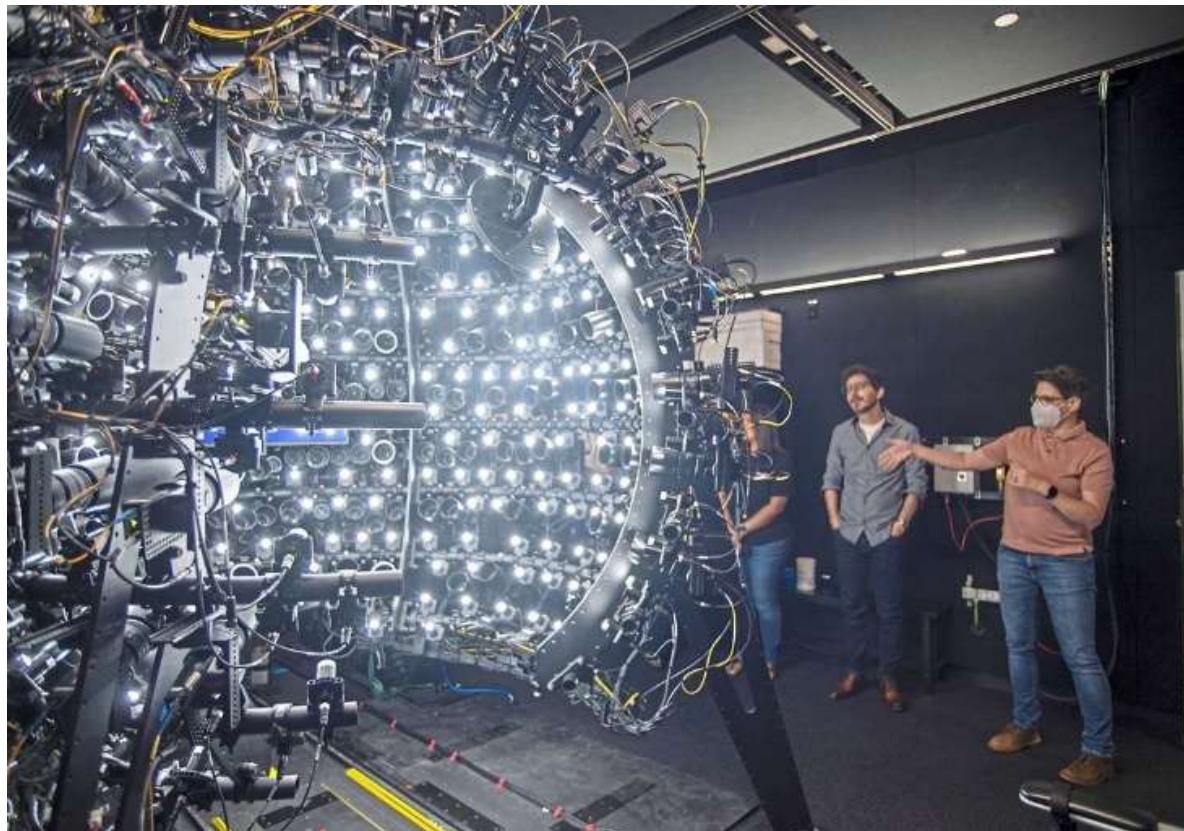
Real world cases

- Captured data with EMG wrist band



Real world cases

- Calibrating over 100 cameras...
- Calibrate cameras for every capture due to temperature/humidity changes



Mugsy studio in
Meta

Real world cases

- We need to synchronize cameras in addition to calibrating them..!
- Tricky way: clap (!)
- Advanced way: use synchronized triggers

