

多视角几何重建

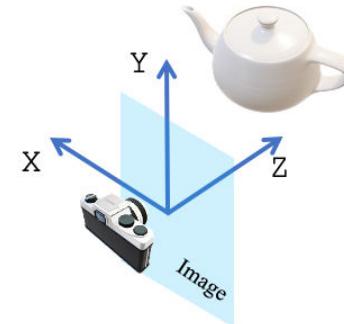
3D from Multiview



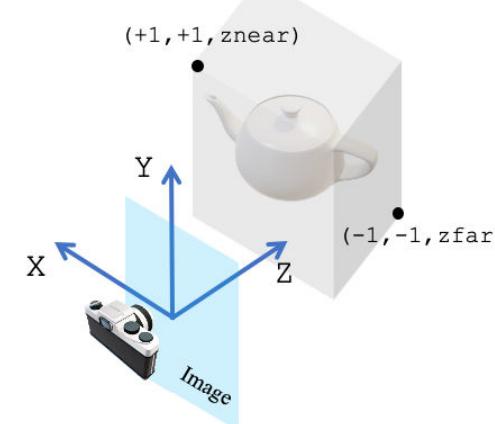
Basic Knowledge — Camera Model



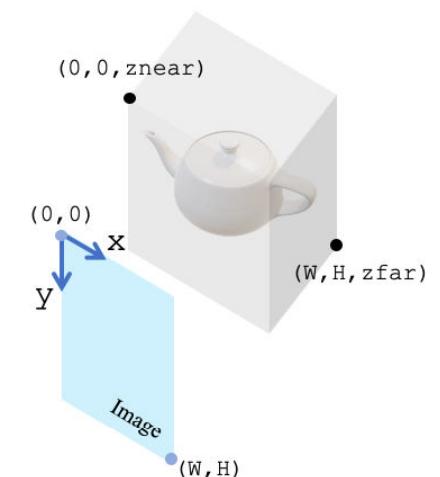
R, T
↻



World
Coordinate System



NDC
Coordinate System



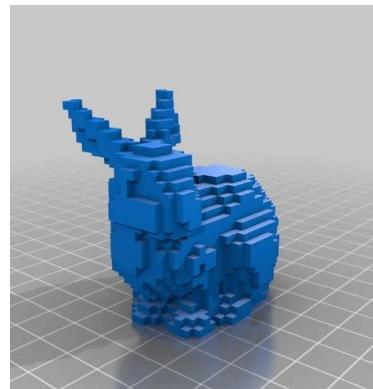
Screen
Coordinate System

All Homogeneous Matrix

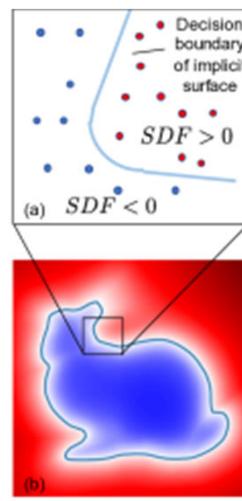
Basic Knowledge — 3D Representation&Rendering



2.5D / Image Based Rendering

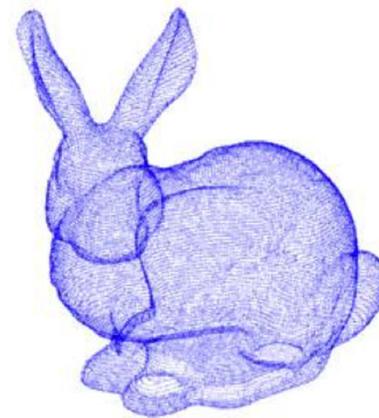


Explicit

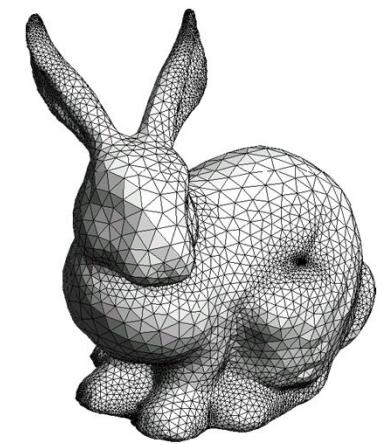


Volumetrics

Implicit

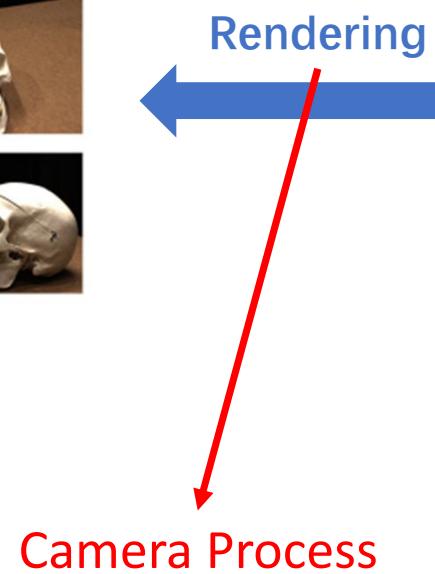


Point clouds



Meshes

What we can do now in 3DV



3D Surface



Light & Reflectance

3D Representation

What we cover today



Recon



3D Surface

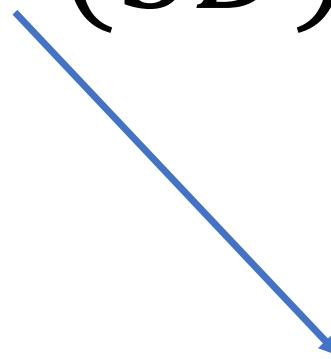


Light & Reflectance

How to obtain good 3D (in certain representation) from Multi-view Images?

Mathematic Formulation

$$F(3D) = 2D$$

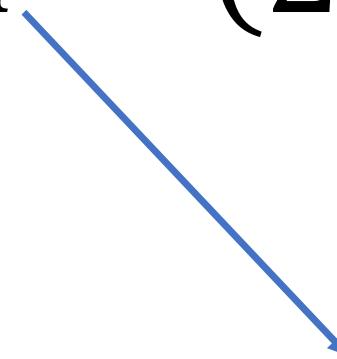


World space to camera space

Camera space projecting to (normalized space to) image space

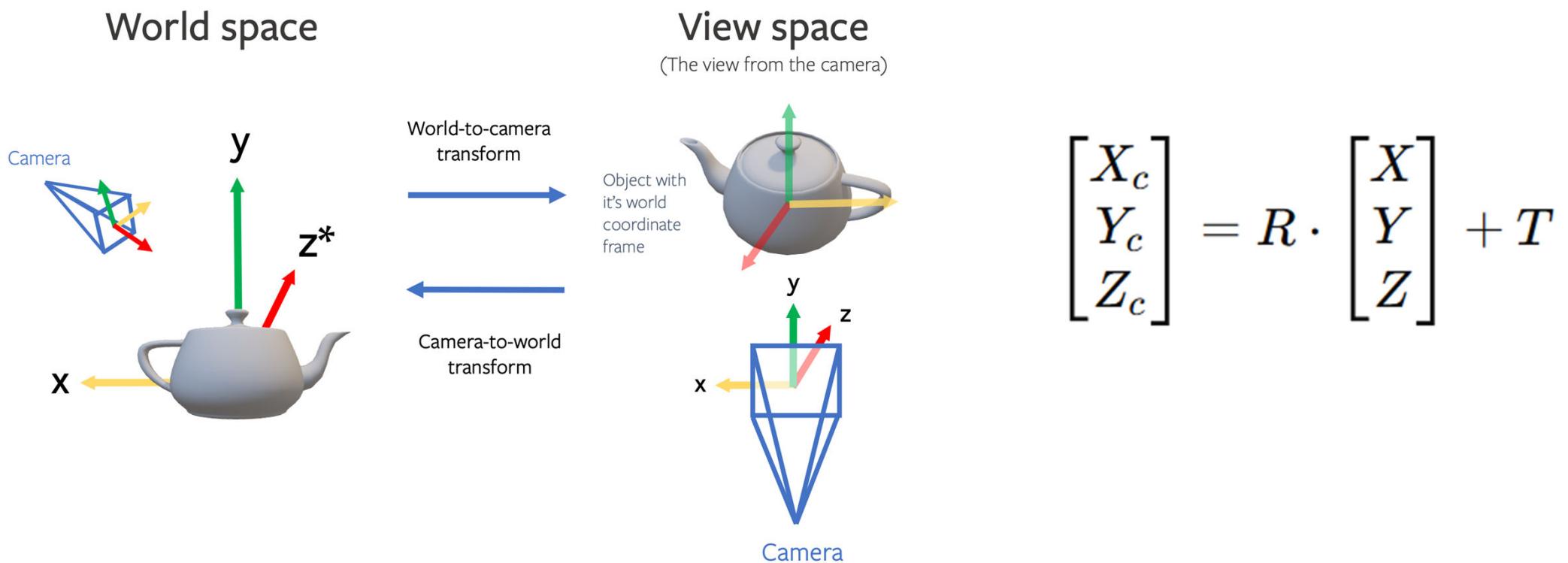
Mathematic Formulation

$$F^{-1}(2D) = 3D$$

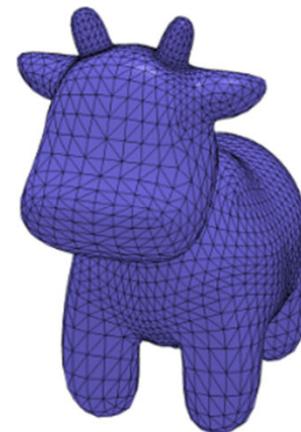
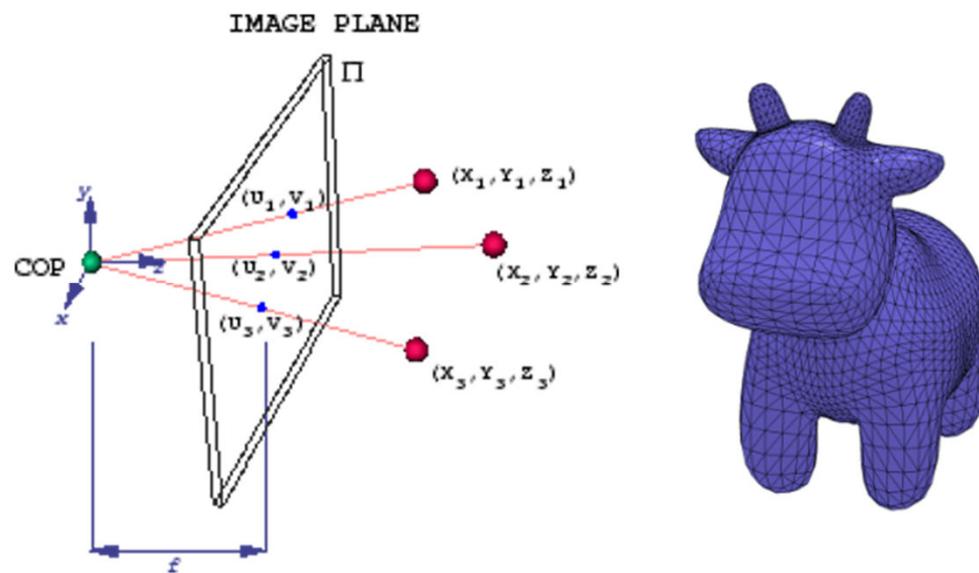


Reconstruction: the inverse function
Solve for 3D from 2D observation

Detailed Formulation of $F(3D) = 2D$



Detailed Formulation of $F(3D) = 2D$



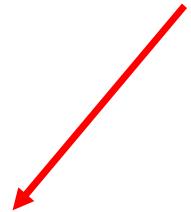
$$u = f \cdot \frac{X_c}{Z_c} + c_x$$
$$v = f \cdot \frac{Y_c}{Z_c} + c_y$$

Detailed Formulation of F(3D) = 2D

$$u = f \cdot \frac{R_{11}X + R_{12}Y + R_{13}Z + T_1}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_x$$
$$v = f \cdot \frac{R_{21}X + R_{22}Y + R_{23}Z + T_2}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_y$$

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} R & T \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad K = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

$$F(3D) = 2D$$



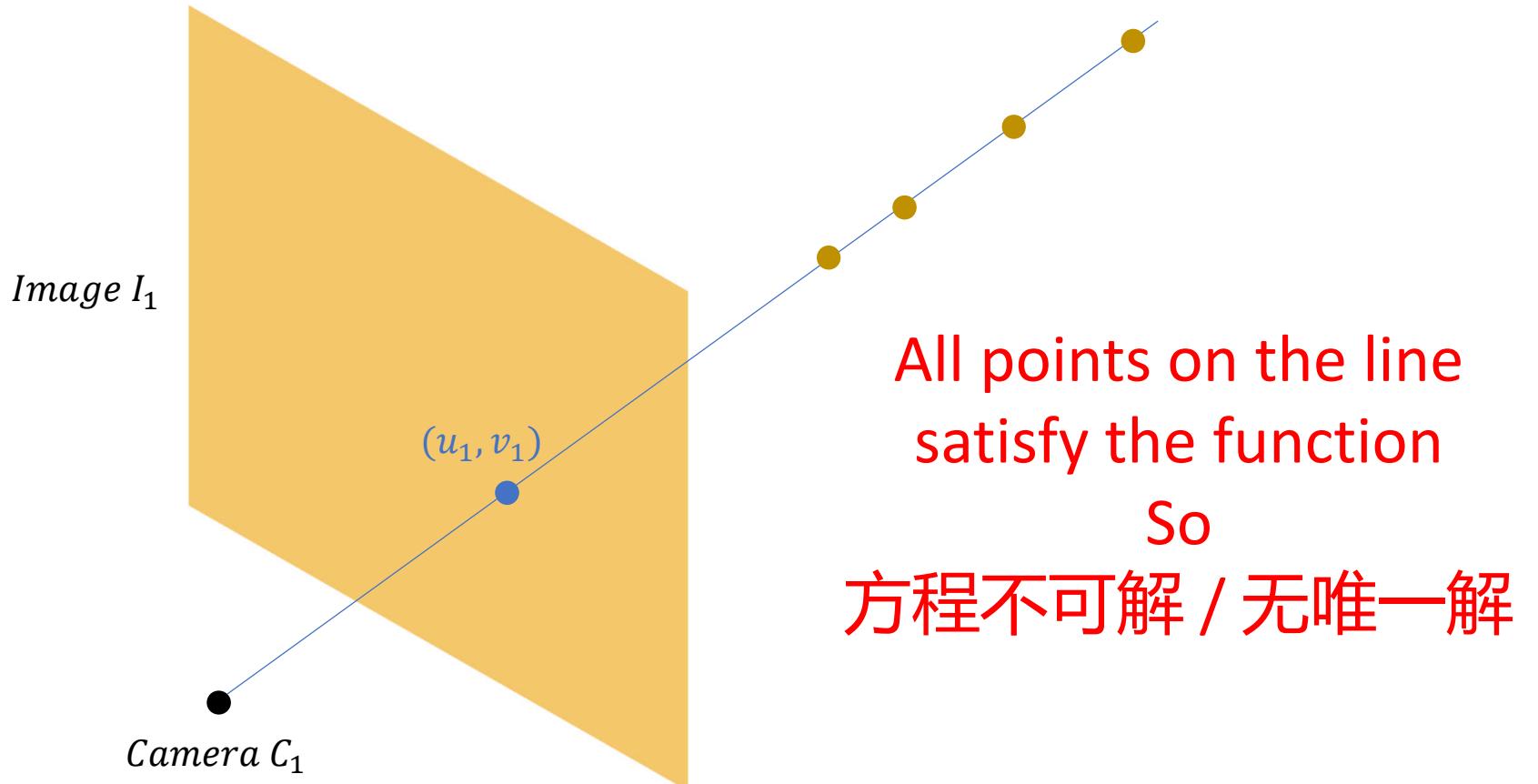
$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} R & T \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad K = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}$$

给定什么样的输入 / 列什么样的等式约束才能让
方程的未知量：3D (XYZ) 是可解的？

Assume for a Point (X,Y,Z)
We have 1 Camera / Image
and we know the Camera Parameters

$$u = f \cdot \frac{R_{11}X + R_{12}Y + R_{13}Z + T_1}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_x$$
$$v = f \cdot \frac{R_{21}X + R_{22}Y + R_{23}Z + T_2}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_y$$

two observations / equations (u,v)
three unknown variables (X, Y, Z)



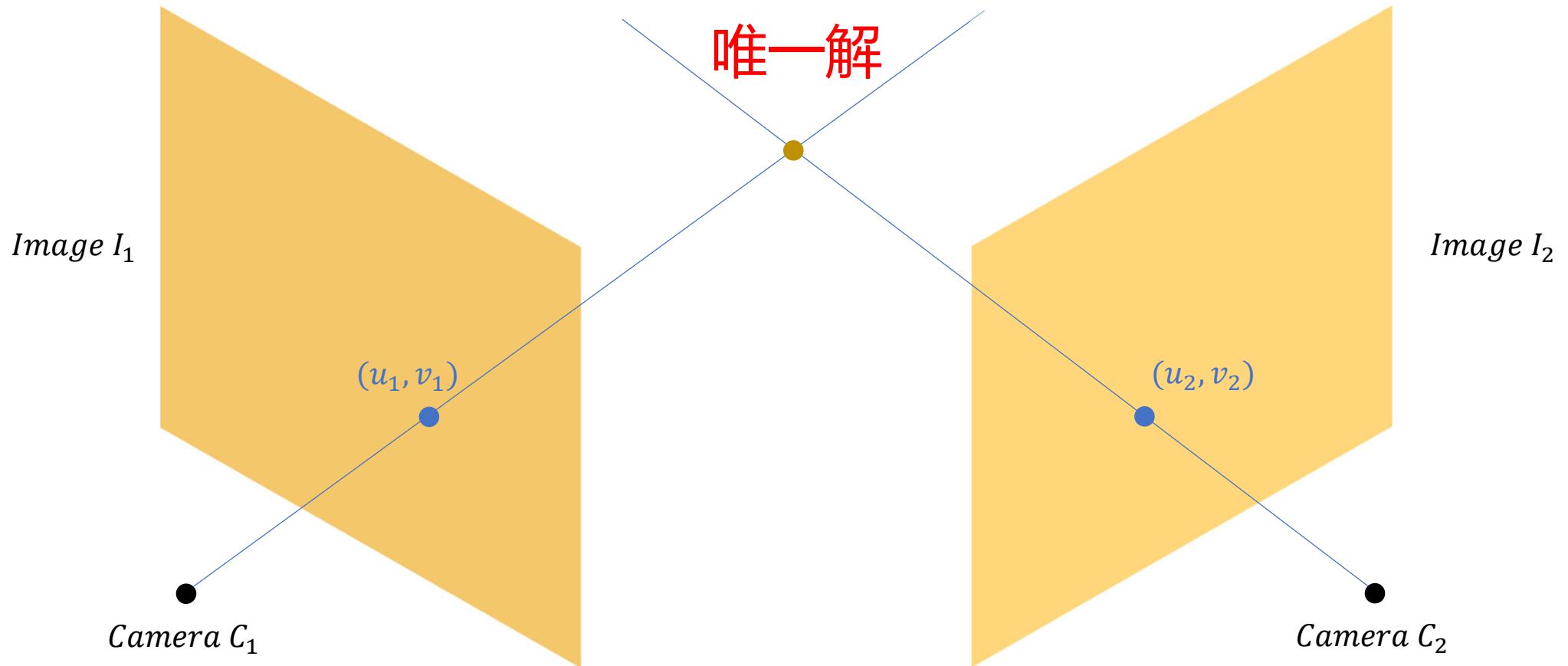
Assume for a Point (X,Y,Z)
We have 2 Camera / Image
and we know Cameras' Parameters

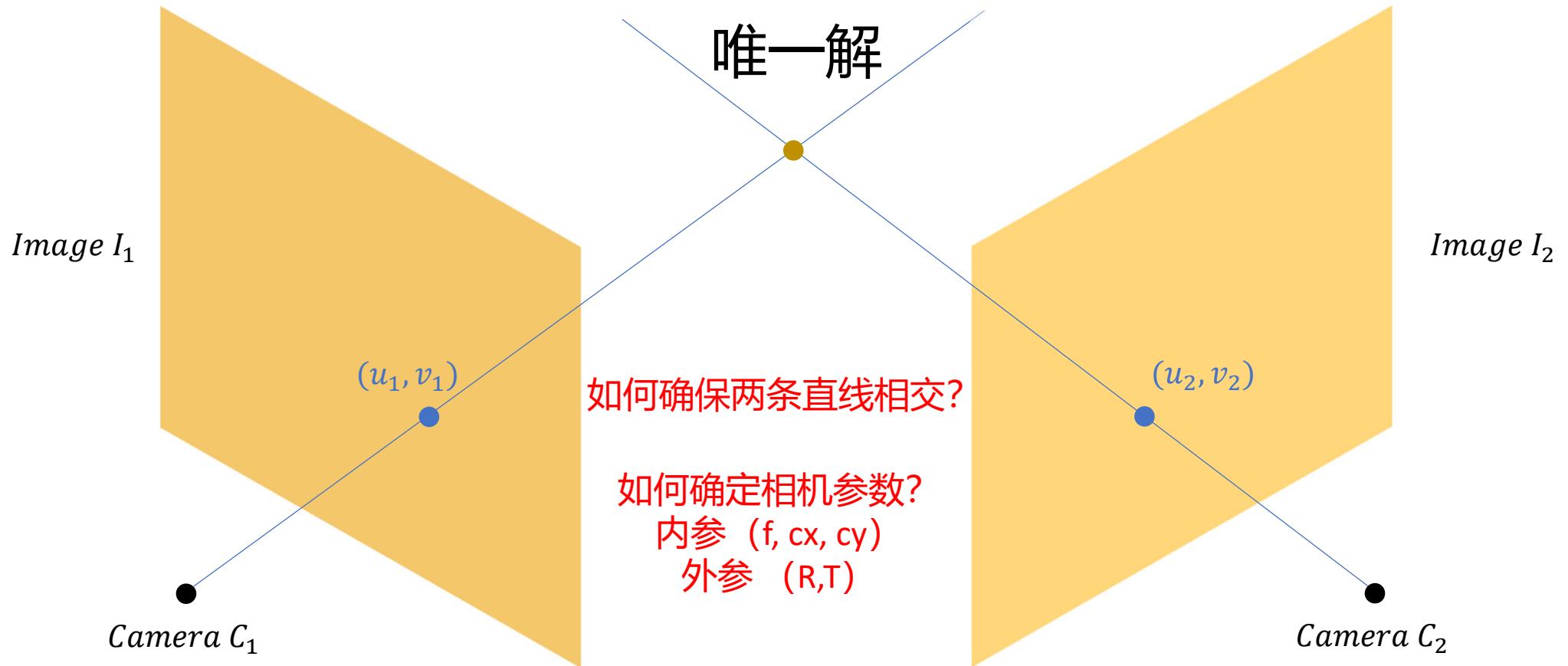
$$u = f \cdot \frac{R_{11}X + R_{12}Y + R_{13}Z + T_1}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_x$$
$$v = f \cdot \frac{R_{21}X + R_{22}Y + R_{23}Z + T_2}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_y$$

4 observations / equations (u_i, v_i)

3 unknown variables (X, Y, Z), so

数学上是可解的





确保相交



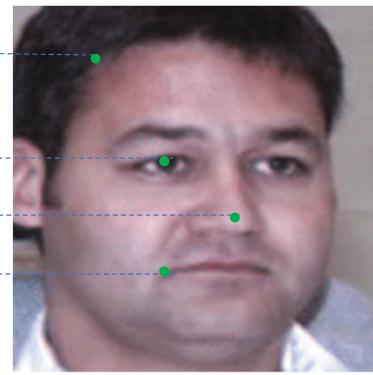
如何确保对应同一个3D点



通过视觉（颜色）
上的相似度

右图中是否存在和左图中红点同一个（3D）点的像素？

光流 (Optical Flow) : 局部区域颜色匹配找对应



$$I'[y + \Delta y, x + \Delta x] = I[y, x]$$

Horn&Schunck Optical Flow

Brightness constancy assumption

$$f(x, y, t) = f(x + dx, y + dy, t + dt)$$

Taylor Series

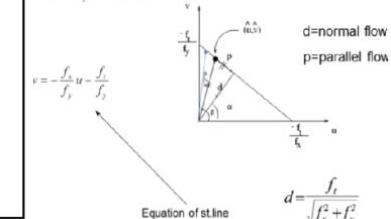
$$f(x, y, t) = f(x, y, t) + \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy + \frac{\partial f}{\partial t} dt$$

$$f_x dx + f_y dy + f_t dt = 0$$

$$f_x u + f_y v + f_t = 0$$

Interpretation of optical flow eq

$$f_x u + f_y v + f_t = 0$$



Lucas & Kanade (Least Squares)

Optical flow eq

$$f_x u + f_y v = -f_t$$

Consider 3 by 3 window

$$f_{x1} u + f_{y1} v = -f_{t1}$$

$$\vdots$$

$$f_{x9} u + f_{y9} v = -f_{t9}$$

$$\boxed{Au = f_t}$$

$$Au = f_t$$

$$A^T Au = A^T f_t$$

$$u = (A^T A)^{-1} A^T f_t$$

Pseudo Inverse

$$\min \sum_i (f_{xi} u + f_{yi} v + f_{ti})^2$$



$$\sum (f_{xi} u + f_{yi} v + f_{ti}) f_{xi} = 0$$

$$\sum (f_{xi} u + f_{yi} v + f_{ti}) f_{yi} = 0$$

Lucas & Kanade

$$\sum (f_{xi} u + f_{yi} v + f_{ti}) f_{xi} = 0$$

$$\sum (f_{xi} u + f_{yi} v + f_{ti}) f_{yi} = 0$$

$$\sum f_{xi}^2 u + \sum f_{xi} f_{yi} v = -\sum f_{xi} f_{ti}$$

$$\sum f_{xi} f_{yi} u + \sum f_{yi}^2 v = -\sum f_{yi} f_{ti}$$

$$\begin{bmatrix} \sum f_{xi}^2 & \sum f_{xi} f_{yi} \\ \sum f_{xi} f_{yi} & \sum f_{yi}^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} -\sum f_{xi} f_{ti} \\ -\sum f_{yi} f_{ti} \end{bmatrix}$$

$$u = \frac{-\sum f_{xi}^2 \sum f_{xi} f_{ti} + \sum f_{xi} f_{yi} \sum f_{yi} f_{ti}}{\sum f_{xi}^2 \sum f_{yi}^2 - (\sum f_{xi} f_{yi})^2}$$

Least Squares Fit

Lucas-Kanade
without pyramids

Fails in areas of large motion

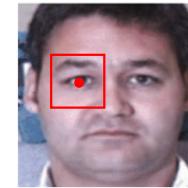
Lucas-Kanade with Pyramids

Taken from the Lecture video from the UCF CRCV course by Prof Mubarak Shah:

<https://www.youtube.com/watch?v=KoMTYnI>NNnc>

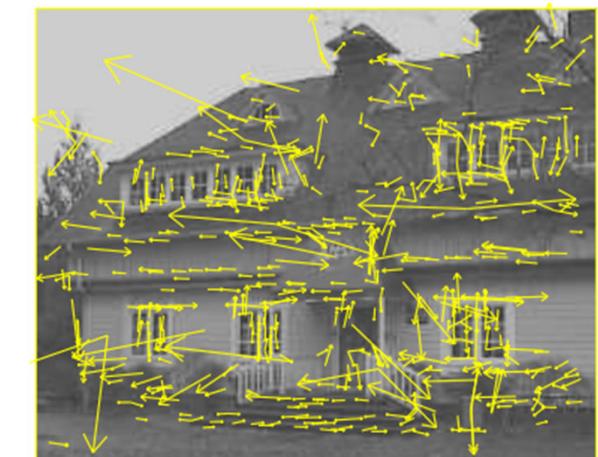
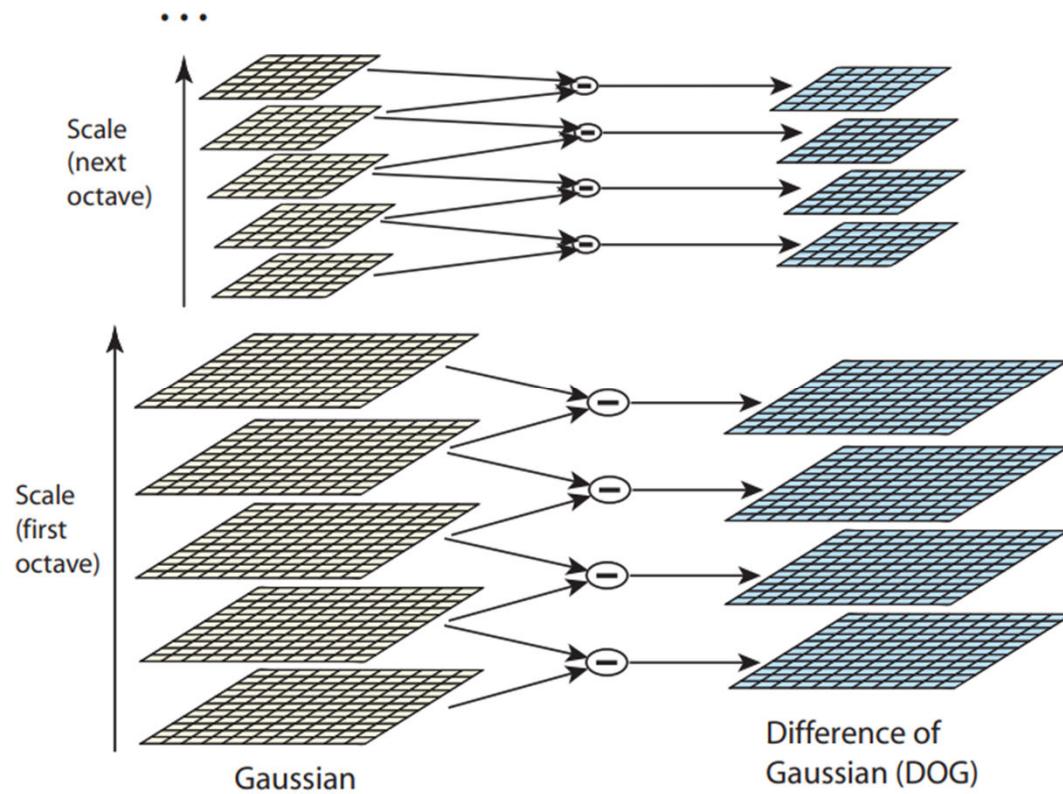
$$v = \frac{\sum f_{xi} f_{ti} \sum f_{xi} f_{yt} - \sum f_{xi}^2 \sum f_{yt} f_{ti}}{\sum f_{xi}^2 \sum f_{yt}^2 - (\sum f_{xi} f_{yt})^2}$$

But 局部颜色对光照、Scale、Rotation等不鲁棒



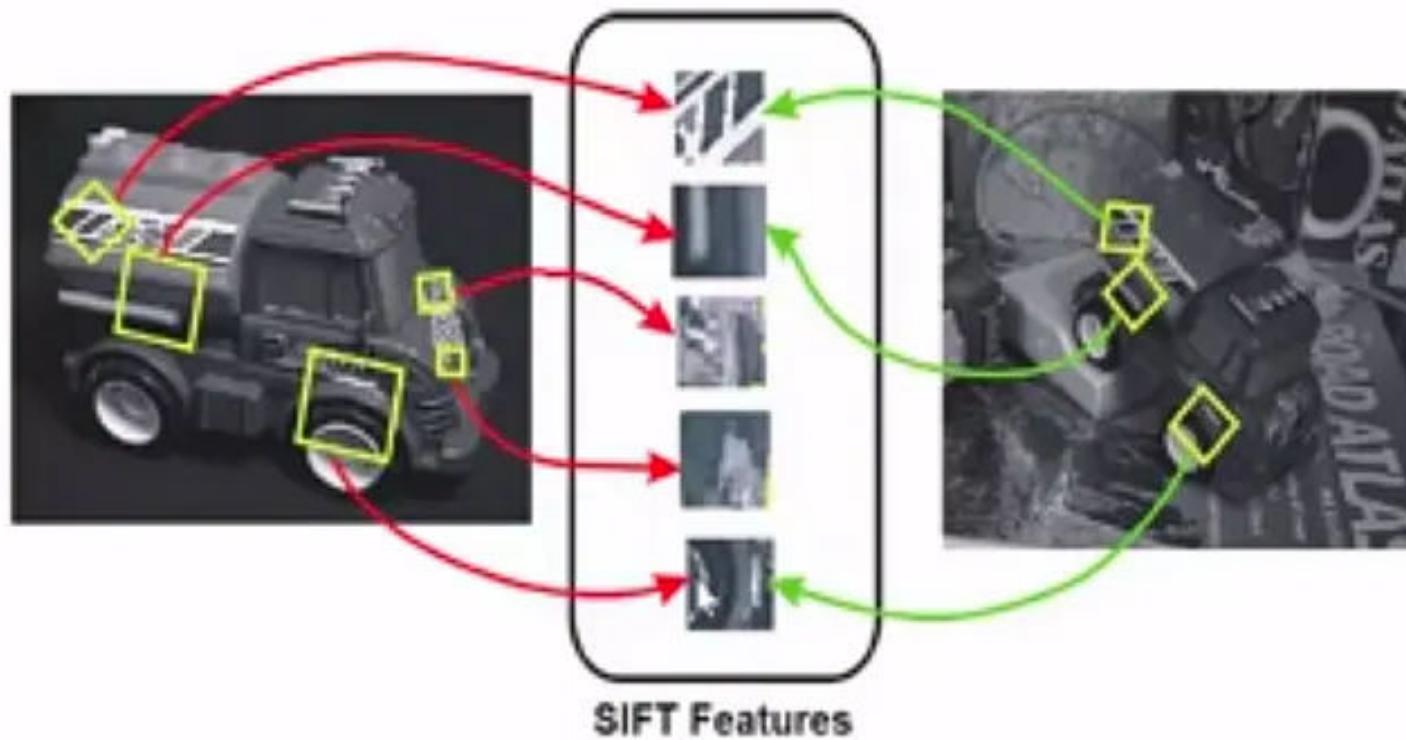
同一个对应点在几种case下的局部颜色度量相差很大

Scale-invariant feature transform (SIFT)



Distinctive Image Features from Scale-Invariant Keypoint. IJCV 2004.

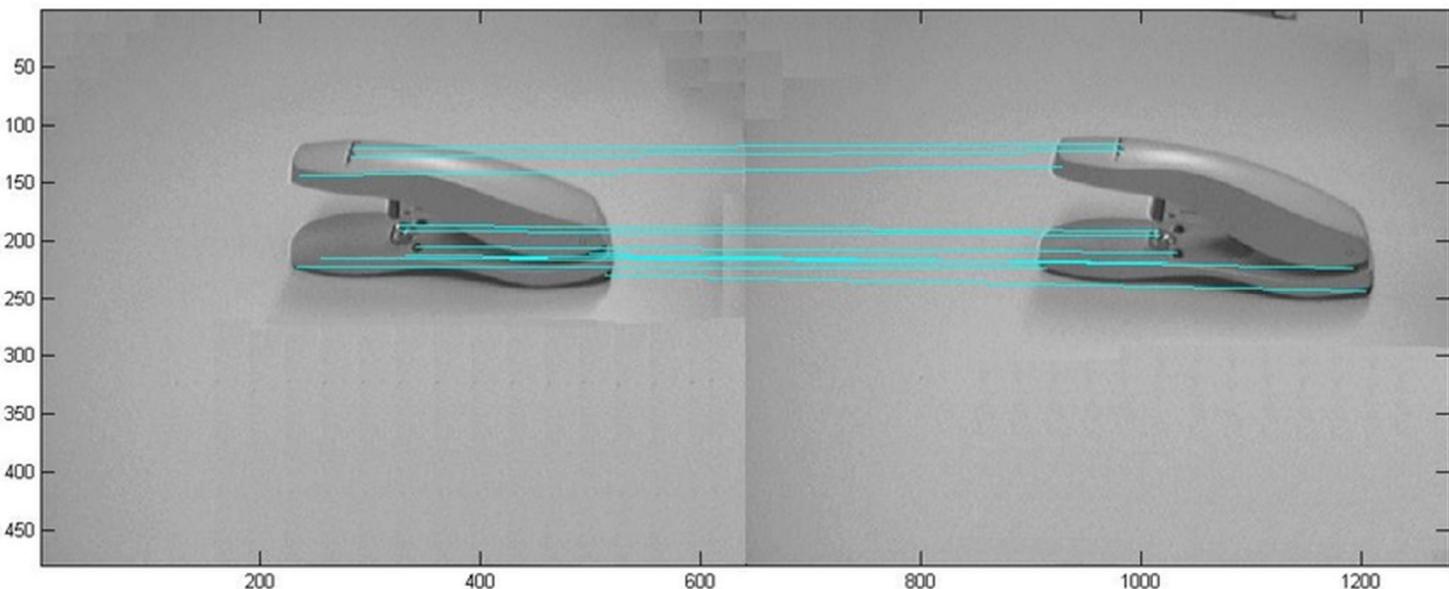
Scale-invariant feature transform (SIFT)



Distinctive Image Features from Scale-Invariant Keypoint. IJCV 2004.

Cases SIFT Fails

Hand-crafted
Small overlap
Large motion
Textureless
Sparse matching



Deep Learning for Correspondence

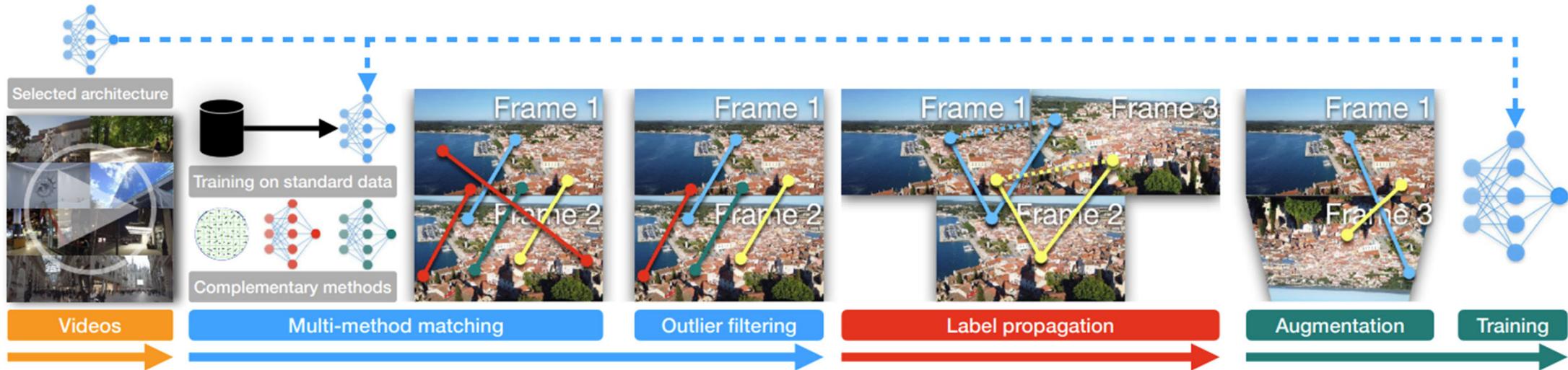
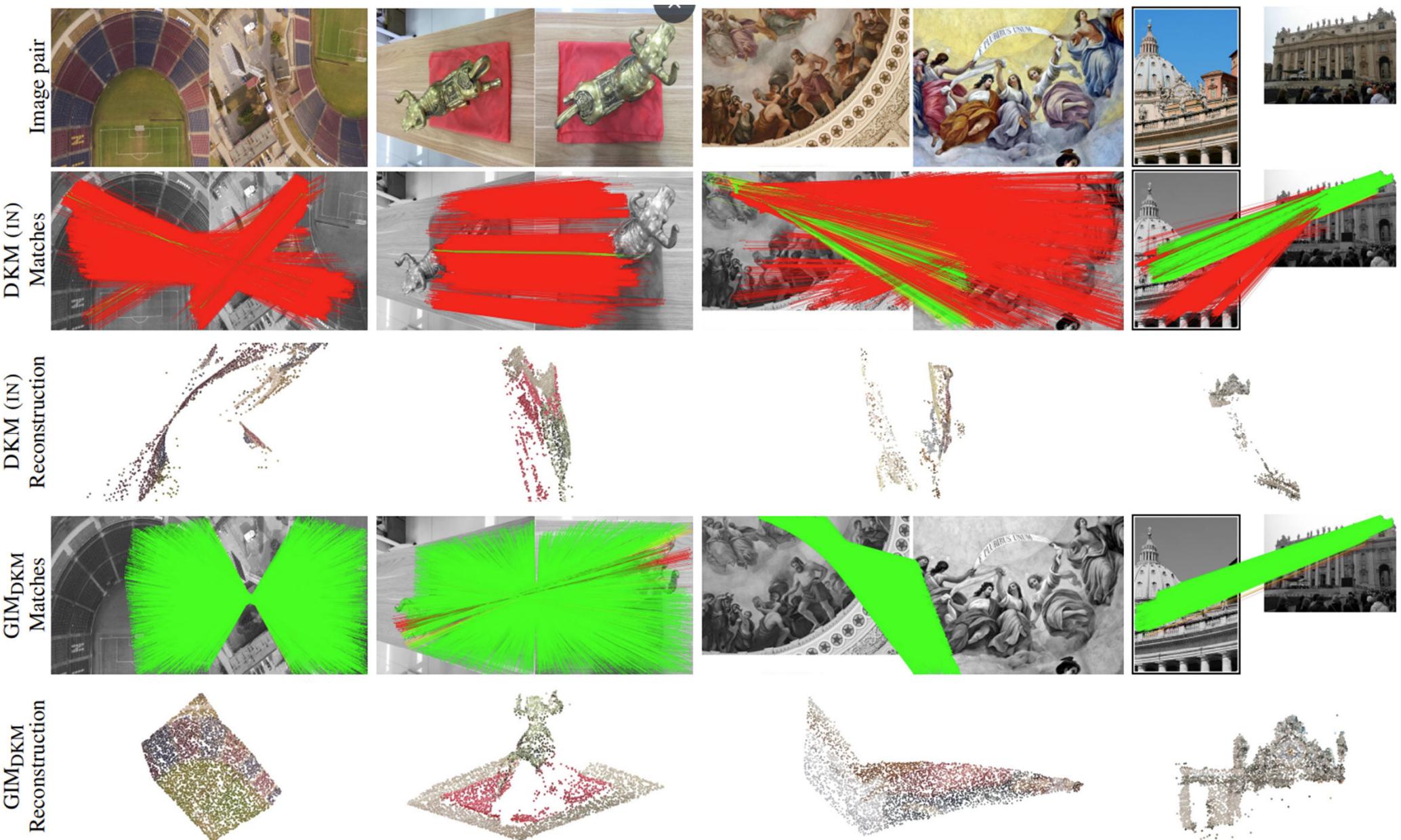
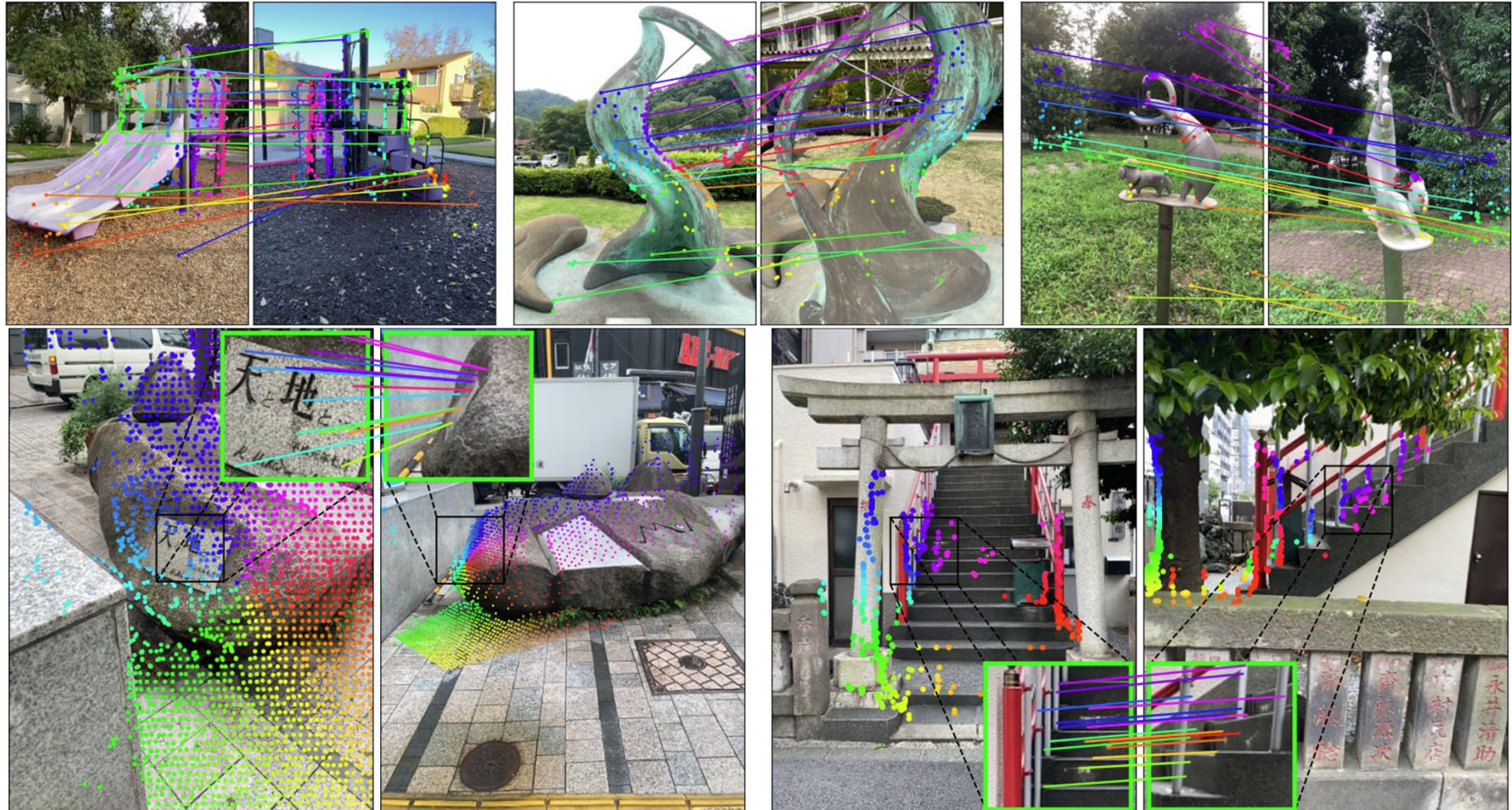


Figure 2: **GIM framework**. We start by downloading a large amount of internet videos. Then, given a selected architecture, we first train it on standard datasets, and generate correspondences between nearby frames by using the trained model with multiple complementary image matching methods. The self-training signal is then enhanced by 1) filtering outlier correspondences with robust fitting, 2) propagating correspondences to distant frames and 3) injecting strong data augmentations.

GIM: LEARNING GENERALIZABLE IMAGE MATCHER FROM INTERNET VIDEOS. ICLR 2024.



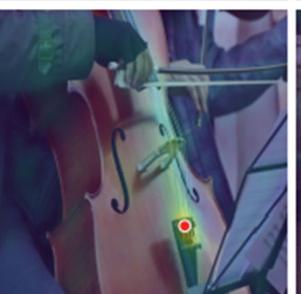
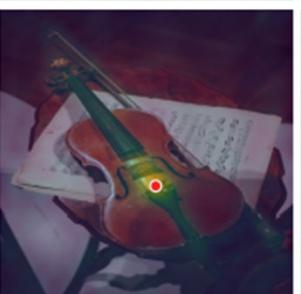
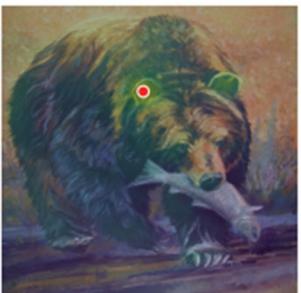
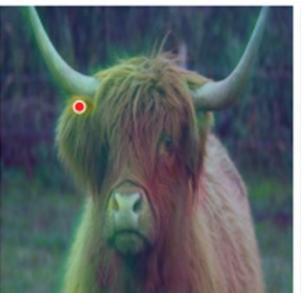
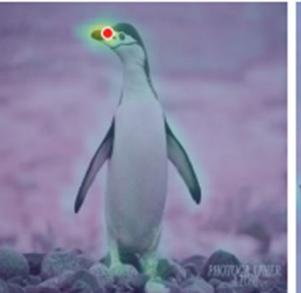
Correspondence for Large Pose Difference



Grounding Image Matching in 3D with MAST3R. ArXiv 2024.

Correspondence between Different Objects

Source Point



cross-instance

cross-category

cross-domain

Emergent Correspondence from Image Diffusion. NeurIPS 2023.

More Applications of Dense Correspondence



CoDeF: Content Deformation Fields for Temporally Consistent Video Processing. CVPR 2024.



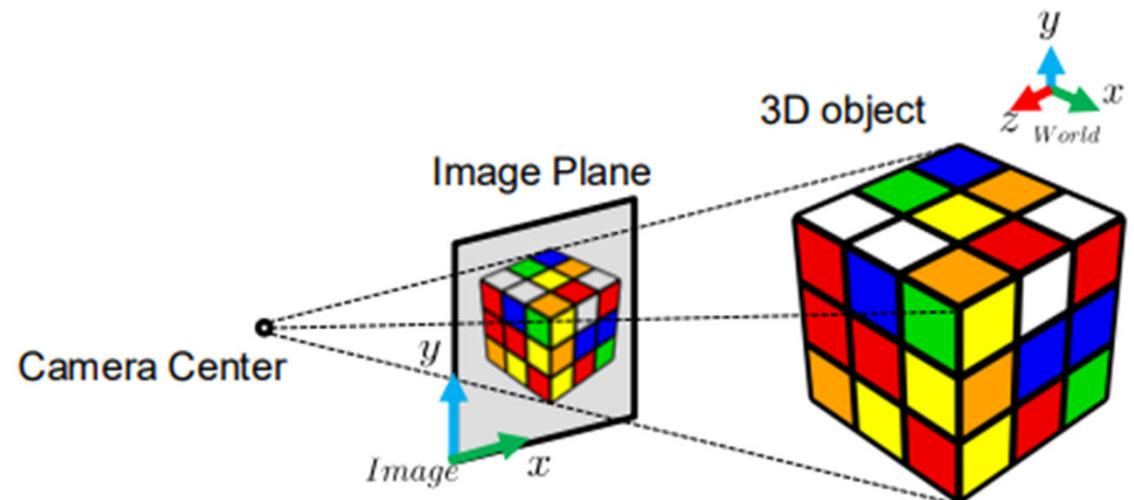
相机标定 (Camera Calibration)

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} R & T \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

$$\mathbf{P} = \mathbf{K}[\mathbf{R}|\mathbf{t}]$$

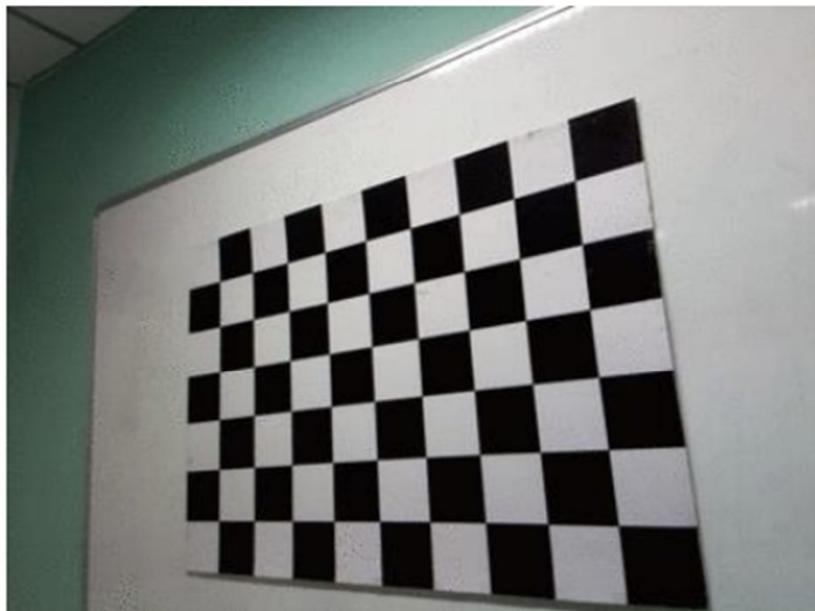
$$\mathbf{P} = \left[\begin{array}{ccc} f & 0 & p_x \\ 0 & f & p_y \\ 0 & 0 & 1 \end{array} \right] \left[\begin{array}{ccc|c} r_1 & r_2 & r_3 & t_1 \\ r_4 & r_5 & r_6 & t_2 \\ r_7 & r_8 & r_9 & t_3 \end{array} \right]$$

intrinsic parameters extrinsic parameters

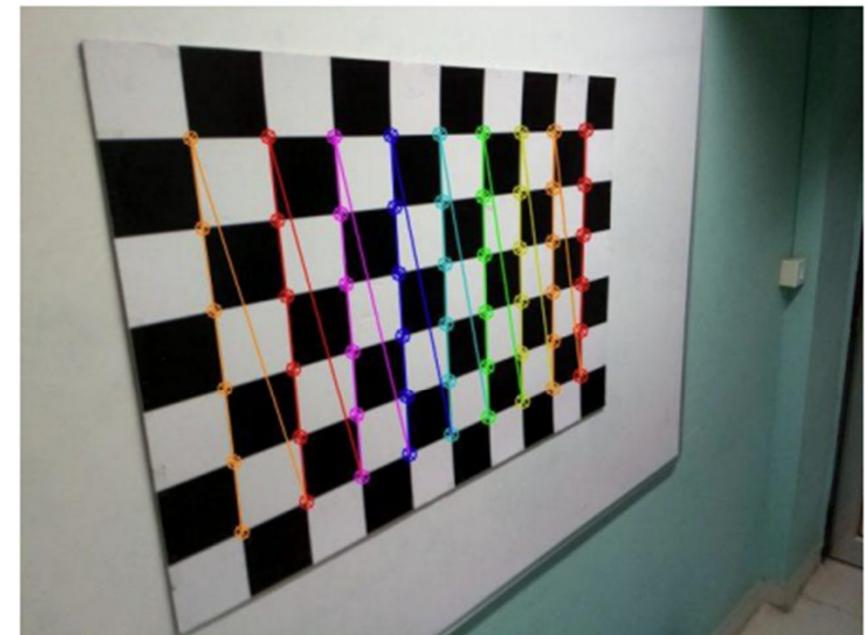


Assume I have a Checkerboard (标定板)

Capture multiple images of the checkerboard from different viewpoints



Find checkerboard corners



Know the 2D & 3D

If I know N 2D & 3D pairs across K Cameras / Images

$$u = f \cdot \frac{R_{11}X + R_{12}Y + R_{13}Z + T_1}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_x$$
$$v = f \cdot \frac{R_{21}X + R_{22}Y + R_{23}Z + T_2}{R_{31}X + R_{32}Y + R_{33}Z + T_3} + c_y$$

2*N*K observations / equations (u_i, v_i)

K*9 unknown variables (3 for (f, cx, cy), 6 for (R, T)), so

If N ≥ 5 , it could be solved

How to Solve?

Optimize with Non-Linear Least Squares

$$\min_{\{\mathbf{R}_k, \mathbf{t}_k\}, \mathbf{K}} \sum_i \sum_k \|\mathbf{x}_{ik} - \pi_{\mathbf{K}}([\mathbf{R}_k | \mathbf{t}_k] \mathbf{X}_i)\|_2^2$$

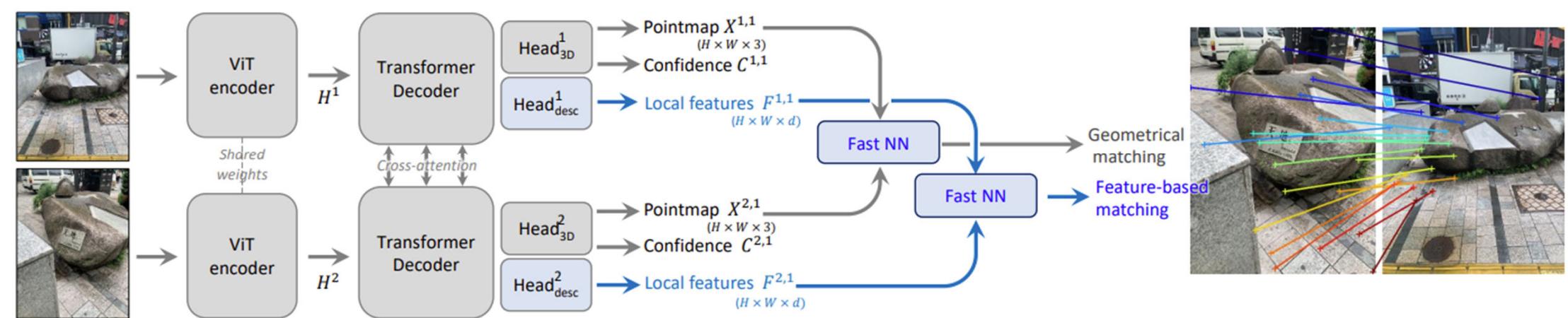
The diagram illustrates the components of the optimization equation:

- Extrinsic**: Points to the term $\{\mathbf{R}_k, \mathbf{t}_k\}$.
- Intrinsic**: Points to the term \mathbf{K} .
- Detection Corner Points**: Points to the term \mathbf{x}_{ik} .
- Perspective Projection**: Points to the term $\pi_{\mathbf{K}}([\mathbf{R}_k | \mathbf{t}_k] \mathbf{X}_i)$.
- Known 3D Location**: Points to the term $[\mathbf{R}_k | \mathbf{t}_k] \mathbf{X}_i$.

$$\boldsymbol{\beta}^{(s+1)} = \boldsymbol{\beta}^{(s)} - (\mathbf{J}_r^T \mathbf{J}_r)^{-1} \mathbf{J}_r^T \mathbf{r}(\boldsymbol{\beta}^{(s)})$$

What if we do not have Checkerboard?

But we could have corresponding points



We know 2D

Assume N pairs in 2 Cameras / Images

$$\min_{\{\mathbf{R}_k, \mathbf{t}_k\}, \mathbf{K}} \sum_i \sum_k \|\mathbf{x}_{ik} - \pi_{\mathbf{K}}([\mathbf{R}_k | \mathbf{t}_k] \mathbf{X}_i)\|_2^2$$

Extrinsic Intrinsic Detection Corner Points Perspective Projection Known 3D Location

2*N*2 observations / equations (u_i, v_i)

Camera: 2*9 unknown variables

3D: 3*N unknown, so

If $N \geq 18$, it could be solved

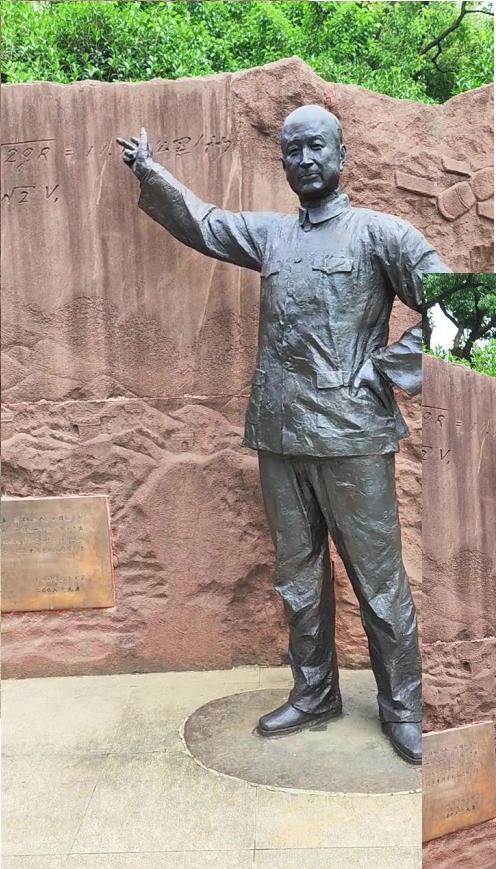
Also could be solved with Gauss-Newton

Extend to Multiple Views — Bundle Adjustment



Step by Step Implementation

Record Images / Videos



Bundle Adjustment Library / Software

colmap



Meshroom



Metashape





Workspace



Workspace (2 chunks, 18 images)

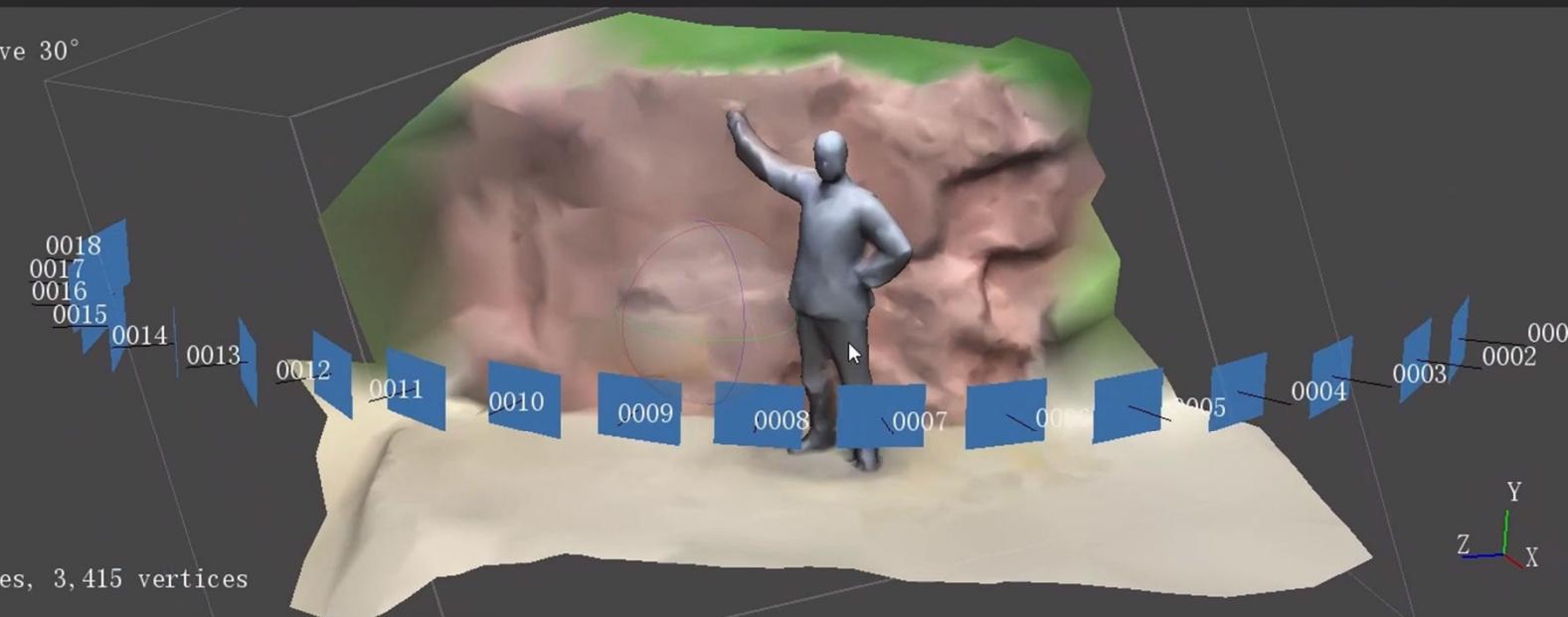
Chunk 1

Chunk 2 (18 images, 6,105 tie points)

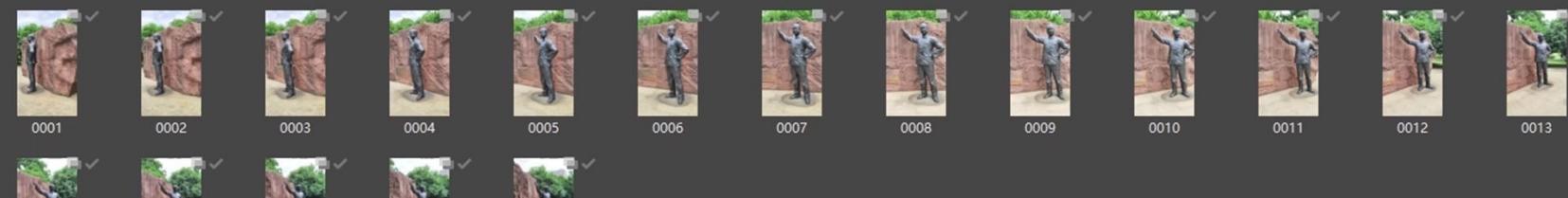
Model

Perspective 30°

6,757 faces, 3,415 vertices



Photos



Photos Console Jobs

Improve Rendering with 3DGS





中国科学技术大学

University of Science and Technology of China

谢谢观看！